

ОГЛАВЛЕНИЕ

Предисловие	6
Глава 1. Введение	7
1.1. Место компилятора в программном обеспечении	7
1.2. Структура компилятора	8
1.3. Реализация множеств и отображений в Java	11
Глава 2. Языки и их представление	14
2.1. Алфавиты, цепочки и языки	14
2.2. Представление языков	16
2.3. Грамматики	17
2.3.1. Формальное определение грамматики (17). 2.3.2. Типы грамматик и их свойства (18).	
2.4. Машины Тьюринга	20
2.4.1. Неразрешимость проблемы останова (21). 2.4.2. Класс рекурсивных множеств (22).	
2.5. Связь машин Тьюринга и грамматик типа 0	24
2.6. Линейно ограниченные автоматы и их связь с контекстно-зависимыми грамматиками	27
Глава 3. Лексический анализ	33
3.1. Регулярные множества и выражения	35
3.2. Конечные автоматы	37
3.3. Интерпретатор НКА на Java	40
3.4. Алгоритмы построения конечных автоматов	41
3.4.1. Построение недетерминированного конечного автомата по регулярному выражению (41). 3.4.2. Построение детерминированного конечного автомата по недетерминированному (43). 3.4.3. Реализация на Java (44). 3.4.4. Обоснование (45). 3.4.5. Построение детерминированного конечного автомата по регулярному выражению (46). 3.4.6. Построение ДКА по РВ на Java (49). 3.4.7. Обоснование (49).	
3.5. Связь регулярных множеств, конечных автоматов и регулярных грамматик	51
3.5.1. Построение детерминированного конечного автомата с минимальным числом состояний (54). 3.5.2. Проверка эквивалентности регулярных языков (58). 3.5.3. Реализация на Java (61).	
3.6. Программирование лексического анализа	62
3.7. Конструктор лексических анализаторов LEX	66
Глава 4. Синтаксический анализ	70
4.1. Контекстно-свободные грамматики и автоматы с магазинной памятью	70
4.2. Преобразования КС-грамматик	78
4.3. Алгоритм Кока–Янгера–Касами	80
4.4. Разбор сверху-вниз (предсказывающий разбор)	80
4.4.1. Алгоритм разбора сверху-вниз (80). 4.4.2. Функции <i>FIRST</i> и <i>FOLLOW</i> (83). 4.4.3. Конструирование таблицы предсказывающего анализатора (87). 4.4.4. LL(<i>k</i>)-грамматики (87). 4.4.5. Удаление левой рекурсии (90). 4.4.6. Левая факторизация (91). 4.4.7. Рекурсивный спуск (92). 4.4.8. Конструктор LL(1)-анализаторов на Java (93). 4.4.9. Восстановление процесса анализа после синтаксических ошибок (93).	
4.5. Разбор снизу-вверх типа сдвиг-свертка	93
4.5.1. Основа (93). 4.5.2. LR(1)-анализаторы (95). 4.5.3. Конструирование LR(1)-таблицы (98). 4.5.4. Конструктор LR(1)-анализаторов	

на Java (103). 4.5.5. Корректность построения (103). 4.5.6. LR(k)-грамматики (106). 4.5.7. Восстановление процесса анализа после синтаксических ошибок (109). 4.5.8. Варианты LR-анализаторов (109).	
Глава 5. Элементы теории перевода	111
5.1. Преобразователи с магазинной памятью	111
5.2. Синтаксически управляемый перевод	112
5.2.1. Схемы синтаксически управляемого перевода (112). 5.2.2. Обобщенные схемы синтаксически управляемого перевода (115).	
5.3. Атрибутные грамматики	117
5.3.1. Определение атрибутных грамматик (117). 5.3.2. Классы атрибутных грамматик и их реализация (120). 5.3.3. Язык описания атрибутных грамматик (123).	
Глава 6. Проверка контекстных условий	127
6.1. Описание областей видимости и блочной структуры	127
6.2. Занесение в среду и поиск объектов	128
Глава 7. Организация таблиц символов	137
7.1. Таблицы идентификаторов	137
7.2. Таблицы расстановки	139
7.3. Таблицы расстановки со списками	141
7.4. Функции расстановки	143
7.5. Таблицы на деревьях	144
7.6. Реализация блочной структуры	148
7.7. Сравнение методов реализации таблиц	148
Глава 8. Промежуточное представление программы	149
8.1. Представление в виде ориентированного графа	149
8.2. Трехадресный код	150
8.3. Линеаризованные представления	154
8.4. Виртуальная машина Java	156
8.4.1. Организация памяти (156). 8.4.2. Набор команд виртуальной машины (157).	
8.5. Организация информации в генераторе кода	159
8.6. Уровень промежуточного представления	160
Глава 9. Генерация кода	161
9.1. Модель машины	161
9.2. Динамическая организация памяти	164
9.2.1. Организация магазина со статической цепочкой (165). 9.2.2. Организация магазина с дисплеем (169).	
9.3. Назначение адресов	170
9.4. Трансляция переменных	171
9.5. Трансляция целых выражений	174
9.6. Трансляция арифметических выражений	175
9.7. Трансляция логических выражений	183
9.8. Выделение общих подвыражений	190
9.9. Трансляция объектно-ориентированных свойств языков программирования	194
9.9.1. Виртуальные базовые классы (194). 9.9.2. Множественное наследование (195). 9.9.3. Единичное наследование и виртуальные функции (196). 9.9.4. Множественное наследование и виртуальные функции (196). 9.9.5. Виртуальные базовые классы с виртуальными функциями (198).	

9.10. Генерация оптимального кода методами сопоставления образцов	200
9.10.1. Сопоставление образцов (200). 9.10.2. Построение покрытия (203). 9.10.3. Выбор дерева вывода наименьшей стоимости (208). 9.10.4. Синтаксический анализ для Т-грамматик (209). 9.10.5. Атрибутная схема для алгоритма сопоставления образцов (210).	
Глава 10. Системы автоматизации построения трансляторов	215
10.1. Система СУПЕР	215
10.2. Система YACC	217
Приложение. Классы атрибутивных грамматик	220
А.1. Атрибутированное дерево разбора	220
А.2. Незацикленные атрибутивные грамматики	220
А.3. Вычислительные последовательности и корректность. Определение визита.	222
А.4. Чистые многовизитные грамматики	223
А.5. Абсолютно незацикленные атрибутивные грамматики	224
А.6. Простые многовизитные атрибутивные грамматики	227
А.7. Одновизитные атрибутивные грамматики	228
А.8. Многопроходные грамматики.	229
Список литературы	234

ПРЕДИСЛОВИЕ

Предлагаемая вниманию читателя книга основана на курсе лекций, прочитанных на факультете вычислительной математики и кибернетики МГУ им. М.В. Ломоносова и на факультете управления и прикладной математики Московского физико-технического института. Диапазон вопросов, затронутых в книге, достаточно широк — от свойств формальных языков и до техники трансляции. Круг этих вопросов представляет собой базис знаний, необходимых квалифицированному специалисту в области методов трансляции. В книге представлены «классические» разделы теории разработки компиляторов: лексический и синтаксический анализ, организация памяти компилятора (таблицы символов) и периода исполнения (магазина), генерация кода. Рассматриваются такие средства автоматизации процесса разработки трансляторов, как LEX, YACC, СУПЕР, методы генерации оптимального кода. На протяжении всего изложения проводится единая «атрибутная» точка зрения на процесс разработки компилятора. Большая часть алгоритмов, изложенных в книге, реализована в виде программ на языке Java. Эти программы составляют программное приложение к книге, размещенное на сайте tryp.ru.

В тексте книги важную роль играют формальные конструкции. В предложениях, их содержащих, обычные правила пунктуации иногда не соблюдаются.

Книга рассчитана как на студентов и аспирантов программистских специальностей, так и на профессионалов в области программирования.

Автор благодарит всех, кто способствовал выходу этой книги в свет, и с признательностью примет все конструктивные замечания по ее содержанию и оформлению.

Глава 1

ВВЕДЕНИЕ

1.1. Место компилятора в программном обеспечении

Компиляторы составляют существенную часть программного обеспечения ЭВМ. Это связано с тем, что языки высокого уровня стали основным средством разработки программ. Сегодня только очень малая часть программного обеспечения, требующая особой эффективности, разрабатывается с помощью ассемблеров. В настоящее время имеют применение довольно много языков программирования. Наряду с традиционными языками, такими, например, как Фортран, широкое распространение получили так называемые «универсальные» языки (Паскаль, Си, Java и др.), а также некоторые специализированные (например, язык обработки списочных структур Лисп). Кроме того, большое распространение получили языки, связанные с узкими предметными областями, такие, как входные языки пакетов прикладных программ.

Для ряда названных языков имеется довольно много реализаций. Так, на рынке программного обеспечения представлены десятки реализаций языков Паскаля, Java или Си. С другой стороны, постоянно растущая потребность в новых компиляторах связана с бурным развитием архитектур ЭВМ. Это развитие идет по различным направлениям. Наряду с возникновением новых архитектур совершенствуются старые архитектуры как в концептуальном отношении, так и по отдельным, конкретным параметрам. Это можно проиллюстрировать на примере микропроцессора Intel. Последовательные версии этого микропроцессора отличаются не только техническими характеристиками, но и, что более важно, новыми возможностями и, значит, изменением (расширением) системы команд. Естественно, это требует новых компиляторов (или модификации старых). То же можно сказать о микропроцессорах Motorola.

В рамках традиционных последовательных машин развивается большое число различных направлений архитектур. Примерами могут служить архитектуры CISC, RISC. Такие ведущие фирмы, как Intel, Motorola, Sun, начинают переходить на выпуск машин с RISC-архитектурами. Естественно, для каждой новой системы команд требуется полный набор новых компиляторов с распространенных языков.

Наконец, бурно развиваются различные параллельные архитектуры. Среди них отметим векторные, многопроцессорные, с широким командным словом архитектуры (вариантом которых являются суперскалярные ЭВМ). На рынке уже имеются десятки типов ЭВМ с параллельной архитектурой, начиная от супер-ЭВМ (Cray, CDC и др.), через рабочие станции (например, IBM RS/6000) и кончая персональными компьютерами. Естественно, каждая из новых машин снабжается новыми компиляторами для многих языков программирования. Здесь необходимо также отметить, что новые архитектуры требуют разработки совершенно новых подходов к созданию компиляторов, так что наряду с собственно разработкой компиляторов ведется и большая научная работа по созданию новых методов трансляции.

1.2. Структура компилятора

Обобщенная структура компилятора и основные фазы компиляции показаны на рис. 1.1.

На начальной фазе лексического анализа входная программа, представляющая собой поток литер, разбивается на лексемы — слова в соответствии с определениями языка. Основными формализмами, лежащими в основе реализации лексических анализаторов, являются конечные автоматы и регулярные выражения. Лексический анализатор может работать в двух основных режимах: либо как подпрограмма, вызываемая синтаксическим анализатором для получения очередной лексемы, либо как полный проход, результатом которого является файл лексем.

В процессе выделения лексем лексический анализатор может как самостоятельно строить таблицы объектов (чисел, строк, идентификаторов и т. п.), так и выдавать значения для каждой лексемы при очередном к нему обращении. В этом случае таблицы объектов строятся в последующих фазах (например, в процессе синтаксического анализа).

На этапе лексического анализа обнаруживаются некоторые (простейшие) ошибки: недопустимые символы, неправильная запись чисел, идентификаторов и др.

Центральная задача синтаксического анализа — разбор структуры программы. Как правило, под структурой понимается дерево, соответствующее разбору в контекстно-свободной грамматике языка. В настоящее время чаще всего используется либо LL(1)-анализ (и его вариант — рекурсивный спуск), либо LR(1)-анализ и его варианты (LR(0), SLR(1), LALR(1) и др.). Рекурсивный спуск чаще используется при ручном программировании синтаксического анализатора, LR(1) — при использовании систем автоматического построения синтаксических анализаторов.



Рис. 1.1

Результатом синтаксического анализа является синтаксическое дерево со ссылками на таблицы объектов. Ошибки, связанные со структурой программы, также обнаруживаются в процессе синтаксического анализа.

На этапе контекстного анализа выявляются зависимости между частями программы, которые не могут быть описаны контекстно-свободным синтаксисом. Это в основном — связи «описание-использование», в частности,

анализ типов объектов, анализ областей видимости, соответствие параметров, метки и др. В процессе контекстного анализа таблицы объектов пополняются информацией об описаниях (свойствах) объектов.

Основным формализмом, используемым при контекстном анализе, является аппарат атрибутивных грамматик. Результатом контекстного анализа является атрибутивное дерево программы. Информация об объектах может быть как рассредоточена в самом дереве, так и сосредоточена в отдельных таблицах объектов. В процессе контекстного анализа также могут быть обнаружены ошибки, связанные с неправильным использованием объектов.

Затем программа может быть переведена во внутреннее представление. Это делается для целей оптимизации и/или удобства генерации кода. Еще одной целью преобразования программы во внутреннее представление является желание иметь переносимый компилятор. Тогда только последняя фаза (генерация кода) является машинно-зависимой. В качестве внутреннего представления могут использоваться префиксная или постфиксная запись, ориентированный граф, тройки, четверки и другие способы.

Фаз оптимизации может быть несколько. Оптимизации обычно делят на машинно-зависимые и машинно-независимые, локальные и глобальные. Определенная часть машинно-зависимой оптимизации выполняется на фазе генерации кода. Глобальная оптимизация пытается принять во внимание структуру всей программы, локальная — только небольших ее фрагментов. Глобальная оптимизация основывается на глобальном потоковом анализе, который выполняется на графе программы и представляет по существу преобразование этого графа. При этом могут учитываться такие свойства программы, как межпроцедурный анализ, межмодульный анализ, анализ областей жизни переменных и т. п.

Наконец, генерация кода — последняя фаза трансляции. Ее результатом является либо ассемблерный модуль, либо объектный (или загрузочный) модуль. В процессе генерации кода могут выполняться некоторые локальные оптимизации, такие, как распределение регистров, выбор длинных или коротких переходов, учет стоимости команд при выборе конкретной последовательности команд. Для генерации кода разработаны различные методы, такие, как таблицы решений, сопоставление образцов, включающее динамическое программирование, различные синтаксические методы.

Конечно, те или иные фазы транслятора могут либо отсутствовать совсем, либо объединяться. В простейшем случае однопроходного транслятора нет явной фазы генерации промежуточного представления и оптимизации, остальные фазы объединены в одну, причем нет и явно построенного синтаксического дерева.

1.3. Реализация множеств и отображений в Java

В качестве программных приложений к книге на сайте `тряп.рф` приведены программы на языке Java, реализующие изложенные алгоритмы. Во многом реализация этих алгоритмов основана на пакете *Java.util*, в котором определен ряд интерфейсов и классов, позволяющих определять множества и отображения. Основу этой структуры составляет интерфейс *Collection*, который определяет набор методов по работе с конечными наборами данных. Рассмотрим кратко те из них, которые используются в приложениях к этой книге. В нем определены следующие методы.

boolean add(Object obj) — добавляет *obj* к вызывающей коллекции. Возвращает *true*, если *obj* был добавлен к коллекции, и *false*, если *obj* уже является элементом коллекции или если коллекция не допускает дубликатов.

boolean addAll(Collection c) — добавляет все элементы коллекции *c* к вызывающей коллекции. Возвращает *true*, если все элементы были добавлены к коллекции, и *false* в противном случае.

boolean contains(Object obj) — возвращает *true*, если вызывающая коллекция содержит элемент *obj*, и *false* в противном случае.

boolean isEmpty() — возвращает *true*, если вызывающая коллекция пуста, и *false* в противном случае.

Iterator iterator() — Возвращает итератор для вызывающей коллекции (подробнее об итераторах сказано ниже).

int size() — возвращает число элементов, содержащихся в вызывающей коллекции.

Интерфейс *List* расширяет *Collection*, позволяя хранить последовательности (списки) элементов. Элементы могут быть внесены или извлечены в соответствии с их позицией, отсчитываемой от нуля. Список может содержать дублированные элементы. Интерфейс *List* добавляет следующие методы:

void add(int index, Object obj) — вставляет *obj* в вызывающий список в позицию с индексом *index*. Элементы списка, начиная с позиции *index*, сдвигаются (т. е. их позиции увеличиваются на 1).

Object get(int index) — возвращает объект, хранящийся в позиции *index* вызывающей коллекции.

Интерфейс *Set* расширяет интерфейс *Collection*, не допуская дублирования элементов, что сказывается на поведении методов *add* и *addAll*. Собственных методов интерфейс *Set* не добавляет.

Перечисленные интерфейсы реализуются следующими классами:

Класс *AbstractCollection* — реализует большую часть интерфейса *Collection*.

Класс *AbstractList* — расширяет *AbstractCollection* и реализует большую часть интерфейса *List*.

Класс *AbstractSequentialList* — расширяет класс *AbstractList* для реализации последовательного, а не произвольного доступа к элементам.

Класс *LinkedList* — реализует связный список, расширяя класс *AbstractSequentialList*. В классе *LinkedList* определены, в частности, следующие методы:

void addFirst(Object obj) — добавить элемент в начало списка;

void addLast(Object obj) — добавить элемент в конец списка;

Object removeFirst() — удалить первый элемент списка;

Object removeLast() — удалить последний элемент списка;

Object getFirst() — получить первый элемент списка;

Object getLast() — получить последний элемент списка;

Класс *ArrayList* — реализует динамический (расширяющийся) массив, расширяя класс *AbstractList*. Для добавления элемента в позицию *index* используется метод *void add(int index, Object obj)*. Удалить элемент можно либо указав ссылку на него, либо указав его позицию:

Object remove(Object obj);

Object remove(int index).

Класс *AbstractSet* — расширяет класс *AbstractCollection* и реализует большую часть интерфейса *Set*.

Класс *HashSet* — расширяет класс *AbstractSet* для реализации множества в виде хэш-таблицы, *TreeSet* расширяет класс *AbstractSet* для реализации множества в виде дерева.

Для перебора элементов коллекции используются итераторы, определенные в интерфейсе *Iterator*. В этом интерфейсе объявлены следующие методы.

boolean hasNext() — возвращает *true*, если в коллекции имеется следующий элемент, иначе возвращает *false*.

Object next() — возвращает следующий элемент. Выбрасывает исключение *NoSuchElementException*, если следующего элемента нет.

void remove() — удаляет текущий элемент. Выбрасывает исключение *IllegalStateException*, если сделана попытка вызвать метод *remove()*, которому не предшествует метод *next()*.

Для реализаций (конечных) отображений используется интерфейс *Map*. В этом интерфейсе определены методы, позволяющие поместить в коллекцию пару (ключ, значение) и извлечь значение по ключу:

Object put(Object k, Object v) — помещает в вызывающую коллекцию (отображение) объект *v*, доступный по ключу *k*. При этом, если с ключом *k* уже было связано некоторое значение, оно возвращается методом *put*, иначе возвращается *null*.

Object get(Object obj) — возвращает значение, связанное с ключом *k*.

boolean containsKey(Object k) — возвращает *true*, если вызывающая коллекция содержит объект с ключом *k*.

Set *keySet()* — возвращает *Set*-объект, состоящий из ключей вызывающей коллекции.

Set *entrySet()* — возвращает *Set*-объект, который содержит элементы отображения, т. е. пары (ключ, значение), имеющие тип *Map.Entry*. В интерфейсе *Map.Entry* определены следующие методы.

Object *getKey()* — возвращает ключ текущего входа вызывающей коллекции (типа *Map.Entry*).

Object *getValue()* — возвращает значение текущего входа вызывающей коллекции.

Отметим, что сами коллекции являются наследниками класса *Object*, поэтому могут помещаться как элементы коллекций. При этом надо иметь в виду, что в каждом конкретном контексте использования элемента коллекции транслятору нужно указать тип, которому принадлежит используемый элемент коллекции. Это делается с помощью явного приведения типа. Схема использования выглядит следующим образом:

```
MyClass element;
String key;
HashMap collection = new HashMap();
collection.put(key, element);
...
HashSet collectionMap = new HashSet(collection.entrySet());
Iteratot iter = collectionMap.iterator();
while (iter.hasNext())
{Map.Entry me = (Map.Entry)iter.next();
String nextKey = (String) me.getKey();
MyClass nextElement = (MyClass) me.getValue();
}
```

Глава 2

ЯЗЫКИ И ИХ ПРЕДСТАВЛЕНИЕ

2.1. Алфавиты, цепочки и языки

Алфавит, или *словарь* — это конечное множество *символов*. Для обозначения символов мы будем пользоваться цифрами, латинскими буквами и специальными литерами типа: #, \$.

Пусть V — алфавит. *Цепочка в алфавите V* — это любая строка конечной длины, составленная из символов алфавита V . Синонимом цепочки являются *предложение*, *строка* и *слово*. Пустая цепочка (обозначается e) — это цепочка, в которую не входит ни один символ.

Конкатенацией цепочек x и y называется цепочка xy . Заметим, что $xe = ex = x$ для любой цепочки x .

Пусть x, y, z — произвольные цепочки в некотором алфавите. Цепочка y называется *подцепочкой* цепочки xyz . Цепочки x и y называются соответственно *префиксом* и *суффиксом* цепочки xy . Заметим, что любой префикс или суффикс цепочки является подцепочкой этой цепочки. Кроме того, пустая цепочка является префиксом, суффиксом и подцепочкой для любой цепочки.

Пример 2.1. Для цепочки $abbba$ префиксом является любая цепочка из множества $L_1 = \{e, a, ab, abb, abbb, abbba\}$, суффиксом является любая цепочка из множества $L_2 = \{e, a, ba, bba, bbba, abbba\}$, подцепочкой является любая цепочка из множества $L_1 \cup L_2$.

Длиной цепочки w (обозначается $|w|$) называется число символов в ней. Например, $|abababa| = 7$, а $|e| = 0$.

Язык в алфавите V — это некоторое множество цепочек в алфавите V .

Пример 2.2. Пусть дан алфавит $V = \{a, b\}$. Вот некоторые языки в алфавите V :

- а) $L_1 = \emptyset$ — пустой язык;
- б) $L_2 = \{e\}$ — язык, содержащий только пустую цепочку (заметим, что L_1 и L_2 — различные языки);
- в) $L_3 = \{e, a, b, aa, ab, ba, bb\}$ — язык, содержащий цепочки из a и b , длина которых не превосходит 2;
- г) L_4 — язык, включающий всевозможные цепочки из a и b , содержащие четное число a и четное число b ;

д) $L_5 = \{a^{n^2} | n > 0\}$ — язык цепочек из a , длины которых представляют собой квадраты натуральных чисел.

Два последних языка содержат бесконечное число цепочек.

Введем обозначение V^* для множества всех цепочек в алфавите V , включая пустую цепочку. Каждый язык в алфавите V является подмножеством V^* . Для множества всех цепочек в алфавите V , кроме пустой цепочки, будем использовать обозначение V^+ .

Пример 2.3. Пусть $V = \{0, 1\}$. Тогда $V^* = \{e, 0, 1, 00, 01, 10, 11, 000, \dots\}$, $V^+ = \{0, 1, 00, 01, 10, 11, 000, \dots\}$.

Введем некоторые операции над языками.

Пусть L_1 и L_2 — языки в алфавите V . *Конкатенацией* языков L_1 и L_2 называется язык $L_1L_2 = \{xy | x \in L_1, y \in L_2\}$.

Пусть L — язык в алфавите V . *Итерацией* языка L называется язык L^* , определяемый следующим образом.

1. $L^0 = \{e\}$.
2. $L^n = LL^{n-1}$, $n \geq 1$.
3. $L^* = \bigcup_{n=0}^{\infty} L^n$.

Пример 2.4. Пусть $L_1 = \{aa, bb\}$, $L_2 = \{e, a, bb\}$. Тогда $L_1L_2 = \{aa, bb, aaa, bba, aabb, bbbb\}$, $L_1^* = \{e, aa, bb, aaaa, aabb, bbaa, bbbb, aaaaaa, \dots\}$.

Большинство языков, представляющих интерес, содержат бесконечное число цепочек. При этом возникают три важных вопроса.

Во-первых, как *представить* язык (т. е. специфицировать входящие в него цепочки)? Если язык содержит только конечное множество цепочек, то ответ прост: можно просто перечислить его цепочки. Если язык бесконечен, то для него необходимо найти конечное представление. Это конечное представление, в свою очередь, будет строкой символов над некоторым алфавитом вместе с некоторой интерпретацией, связывающей это представление с языком.

Во-вторых, для любого ли языка существует конечное представление? Можно предположить, что ответ отрицателен. Мы увидим, что множество всех цепочек над алфавитом счетно. Язык — это любое подмножество цепочек. Из теории множеств известно, что множество всех подмножеств счетного множества несчетно. Хотя мы и не дали строгого определения того, что является конечным представлением, интуитивно ясно, что любое разумное определение этого понятия ведет только к счетному множеству конечных представлений, поскольку нужно иметь возможность записать такое конечное представление в виде строки символов конечной длины. Поэтому языков значительно больше, чем конечных представлений.

В-третьих, можно спросить: какова *структура* тех классов языков, для которых существует конечное представление?

2.2. Представление языков

Процедура — это конечная последовательность инструкций, которые могут быть механически выполнены. Примером может служить машинная программа. Процедура, которая всегда заканчивается, называется *алгоритмом*.

Один из способов представления языка — дать алгоритм, определяющий, принадлежит ли цепочка языку. Более общий способ состоит в том, чтобы дать процедуру, которая останавливается с ответом «да» для цепочек, принадлежащих языку, и либо останавливается с ответом «нет», либо вообще не останавливается для цепочек, не принадлежащих языку. Говорят, что такая процедура или алгоритм *распознает* язык.

Такой метод представляет язык с точки зрения распознавания. Язык можно также представить методом порождения. А именно, можно дать процедуру, которая систематически *порождает* в определенном порядке цепочки языка.

Если мы можем распознать цепочки языка над алфавитом V либо с помощью процедуры, либо с помощью алгоритма, то мы можем и генерировать язык, поскольку мы можем систематически генерировать все цепочки из V^* , проверять каждую цепочку на принадлежность языку и выдавать список только цепочек языка. Но если процедура не всегда заканчивается при проверке цепочки, то мы не сдвинемся дальше первой цепочки, на которой процедура не заканчивается. Эту проблему можно обойти, организовав проверку таким образом, чтобы процедура никогда не продолжала проверять одну цепочку бесконечно. Для этого введем следующую конструкцию.

Предположим, что V имеет p символов. Мы можем рассматривать цепочки из V^* как числа, представленные в базисе p , плюс пустая цепочка ϵ . Можно занумеровать цепочки в порядке возрастания длины и в «числовом» порядке для цепочек одинаковой длины.

Пусть P — процедура для проверки принадлежности цепочки языку L . Предположим, что P может быть представлена дискретными шагами, так что имеет смысл говорить об i -м шаге процедуры для любой данной цепочки. Прежде чем дать процедуру перечисления цепочек языка L , дадим процедуру нумерации пар положительных чисел.

Все упорядоченные пары положительных чисел можно взаимнооднозначно отобразить на множество положительных чисел и тем самым упорядочить. Теперь можно дать процедуру перечисления цепочек L . Упорядочиваем и нумеруем пары целых положительных чисел: $(1,1)$, $(2,1)$, $(1,2)$, $(3,1)$, $(2,2)$, При нумерации пары (i, j) генерируем i -ю цепочку из V^* и применяем к цепочке первые j шагов процедуры P . Как только мы определили, что сгенерированная цепочка принадлежит L , добавляем цепочку к списку элементов L . Если цепочка i принадлежит L , то это будет определено P за j шагов для некоторого конечного j . При перечислении (i, j)

будет сгенерирована цепочка с номером i . Легко видеть, что эта процедура перечисляет все цепочки L .

Если мы имеем процедуру генерации цепочек языка, то мы всегда можем построить процедуру распознавания предложений языка, но не всегда алгоритм. Для определения того, принадлежит ли x языку L , просто нумеруем предложения L и сравниваем x с каждым предложением. Если сгенерировано x , то процедура останавливается, распознав, что x принадлежит L . Конечно, если x не принадлежит L , то процедура никогда не закончится.

Язык, предложения которого могут быть сгенерированы процедурой, называется *рекурсивно перечислимым*. Язык рекурсивно перечислим, если имеется процедура, распознающая предложения языка. Говорят, что язык *рекурсивен*, если существует алгоритм для распознавания языка. Класс рекурсивных языков является собственным подмножеством класса рекурсивно перечислимых языков. Мало того, существуют языки, не являющиеся даже рекурсивно перечислимыми.

2.3. Грамматики

2.3.1. Формальное определение грамматики. Для нас наибольший интерес представляет одна из систем генерации языков — *грамматики*. Понятие грамматики изначально было формализовано лингвистами для естественных языков. Предполагалось, что это может помочь при машинном переводе. Однако наилучшие результаты в этом направлении достигнуты при описании не естественных языков, а языков программирования. Примером может служить способ описания синтаксиса языков программирования при помощи БНФ — формы Бэкуса–Наура.

Определение 2.1. *Грамматика* — это четверка $G = (N, T, P, S)$, где:

- 1) N — алфавит *нетерминальных символов*;
- 2) T — алфавит *терминальных символов*, $N \cap T = \emptyset$;
- 3) P — конечное множество *правил* вида $\alpha \rightarrow \beta$, где $\alpha \in (N \cup T)^* N (N \cup T)^*$, $\beta \in (N \cup T)^*$;
- 4) $S \in N$ — *начальный знак* (или *аксиома*) грамматики.

Мы будем использовать большие латинские буквы для обозначения нетерминальных символов, малые латинские буквы из начала алфавита для обозначения терминальных символов, малые латинские буквы из конца алфавита для обозначения цепочек из T^* и, наконец, малые греческие буквы для обозначения цепочек из $(N \cup T)^*$.

Будем использовать также сокращенную запись $A \rightarrow \alpha_1 | \alpha_2 | \dots | \alpha_n$ для обозначения группы правил $A \rightarrow \alpha_1, A \rightarrow \alpha_2, \dots, A \rightarrow \alpha_n$.

Определим на множестве $(N \cup T)^*$ бинарное отношение *выводимости* \Rightarrow следующим образом: если $\delta \rightarrow \gamma \in P$, то $\alpha\delta\beta \Rightarrow \alpha\gamma\beta$ для всех $\alpha, \beta \in (N \cup T)^*$. Если $\alpha_1 \Rightarrow \alpha_2$, то говорят, что цепочка α_2 *непосредственно выводима* из α_1 .

Мы будем использовать также рефлексивно-транзитивное и транзитивное замыкания отношения \Rightarrow , а также его степень $k \geq 0$ (обозначаемые соответственно \Rightarrow^* , \Rightarrow^+ и \Rightarrow^k). Если $\alpha_1 \Rightarrow^* \alpha_2$ ($\alpha_1 \Rightarrow^+ \alpha_2$, $\alpha_1 \Rightarrow^k \alpha_2$), то говорят, что цепочка α_2 *выводима* (*нетривиально выводима*, *выводима за k шагов*) из α_1 .

Если $\alpha \Rightarrow^k \beta$ ($k \geq 0$), то существует последовательность шагов

$$\gamma_0 \Rightarrow \gamma_1 \Rightarrow \gamma_2 \Rightarrow \dots \Rightarrow \gamma_{k-1} \Rightarrow \gamma_k,$$

где $\alpha = \gamma_0$ и $\beta = \gamma_k$. Последовательность цепочек $\gamma_0, \gamma_1, \gamma_2, \dots, \gamma_k$ в этом случае называют *выводом* β из α .

Сентенциальной формой грамматики G называется цепочка, выводимая из ее начального символа.

Языком, порождаемым грамматикой G (обозначается $L(G)$), называется множество всех ее терминальных сентенциальных форм, т. е.

$$L(G) = \{w \mid w \in T^*, S \Rightarrow^+ w\}.$$

Грамматики G_1 и G_2 называются *эквивалентными*, если они порождают один и тот же язык, т. е. $L(G_1) = L(G_2)$.

Пример 2.5. Грамматика $G = (\{S, B, C\}, \{a, b, c\}, P, S)$, где $P = \{S \rightarrow aSBC, S \rightarrow aBC, CB \rightarrow BC, aB \rightarrow ab, bB \rightarrow bb, bC \rightarrow bc, cC \rightarrow cc\}$, порождает язык $L(G) = \{a^n b^n c^n \mid n > 0\}$.

Действительно, применяем $n - 1$ раз правило 1 и получаем $a^{n-1}S(BC)^{n-1}$, затем один раз правило 2 и получаем $a^n(BC)^n$, затем $n(n - 1)/2$ раз правило 3 и получаем $a^n B^n C^n$.

Затем используем правило 4 и получаем $a^n b B^{n-1} C^n$. Далее применяем $n - 1$ раз правило 5 и получаем $a^n b^n C^n$. Затем, применив правило 6 и $n - 1$ раз правило 7, получаем $a^n b^n c^n$. Можно показать, что язык $L(G)$ состоит из цепочек только такого вида.

Пример 2.6. Рассмотрим грамматику $G = (\{S\}, \{0, 1\}, \{S \rightarrow 0S1, S \rightarrow 01\}, S)$. Легко видеть, что цепочка $000111 \in L(G)$, так как существует вывод

$$S \Rightarrow 0S1 \Rightarrow 00S11 \Rightarrow 000111.$$

Нетрудно показать, что грамматика порождает язык $L(G) = \{0^n 1^n \mid n > 0\}$.

Пример 2.7. Рассмотрим грамматику $G = (\{S, A\}, \{0, 1\}, \{S \rightarrow 0S, S \rightarrow 0A, A \rightarrow 1A, A \rightarrow 1\}, S)$. Нетрудно показать, что грамматика порождает язык $L(G) = \{0^n 1^m \mid n, m > 0\}$.

2.3.2. Типы грамматик и их свойства. Рассмотрим классификацию грамматик (предложенную Н. Хомским), основанную на виде их правил.

Определение 2.2. Пусть дана грамматика $G = (N, T, P, S)$. Тогда:

- 1) если на правила грамматики не наложены никакие ограничения, то ее называют грамматикой типа 0, или грамматикой *без ограничений*.
- 2) если
 - а) каждое правило грамматики, кроме $S \rightarrow e$, имеет вид $\alpha \rightarrow \beta$, где $|\alpha| \leq |\beta|$, и
 - б) в том случае, когда $S \rightarrow e \in P$, символ S не встречается в правых частях правил,

то грамматiku называют грамматикой типа 1, или *неукорачивающей*, или *контекстно-зависимой* (КЗ-грамматикой, КЗГ), или *контекстно-чувствительной* (КЧ-грамматикой);

- 3) если каждое правило грамматики имеет вид $A \rightarrow \beta$, где $A \in N$, $\beta \in (N \cup T)^*$, то ее называют грамматикой типа 2, или *контекстно-свободной* (КС-грамматикой);
- 4) если каждое правило грамматики имеет вид либо $A \rightarrow xB$, либо $A \rightarrow x$, где $A, B \in N$, $x \in T^*$, то ее называют грамматикой типа 3, или *праволинейной*.

Легко видеть, что грамматика в примере 2.5 — неукорачивающая, в примере 2.6 — контекстно-свободная, в примере 2.7 — праволинейная.

Язык, порождаемый грамматикой типа i , называют языком типа i . Язык типа 0 называют также *языком без ограничений*, язык типа 1 — *контекстно-зависимым* (КЗ), язык типа 2 — *контекстно-свободным* (КС), язык типа 3 — *праволинейным*.

Теорема 2.1. *Каждый контекстно-свободный язык может быть порожден неукорачивающей контекстно-свободной грамматикой.*

Доказательство. Пусть L — контекстно-свободный язык. Тогда существует контекстно-свободная грамматика $G = (N, T, P, S)$, порождающая L .

Построим новую грамматiku $G' = (N', T, P', S')$ следующим образом.

1. Если в P есть правило вида $A \rightarrow \alpha_0 B_1 \alpha_1 \dots B_k \alpha_k$, где $k \geq 0$, $B_i \Rightarrow^+ e$ для $1 \leq i \leq k$, и ни из одной цепочки α_j ($0 \leq j \leq k$) не выводится e , то включить в P' все правила (кроме $A \rightarrow e$) вида

$$A \rightarrow \alpha_0 X_1 \alpha_1 \dots X_k \alpha_k,$$

где X_i — это либо B_i , либо e .

2. Если $S \Rightarrow^+ e$, то включить в P' правила $S' \rightarrow S$, $S' \rightarrow e$ и положить $N' = N \cup \{S'\}$; в противном случае положить $N' = N$ и $S' = S$.

Порождает ли грамматика пустую цепочку, можно установить следующим простым алгоритмом.

Шаг 1. Строим множество $N_0 = \{N \mid N \rightarrow e\}$.

Шаг 2. Строим множество $N_i = N_{i-1} \cup \{N \mid N \rightarrow \alpha, \alpha \in \{N_{i-1}\}^*\}$.

Шаг 3. Если $N_i = N_{i-1}$, то перейти к шагу 4, иначе к шагу 2.

Шаг 4. Если $S \in N_i$, то $S \rightarrow *e$.

Легко видеть, что G' — неукорачивающая грамматика. Можно показать по индукции, что $L(G') = L(G)$. ■

Пусть K_i — класс всех языков типа i . Доказано, что справедливо следующее (строгое) включение: $K_3 \subset K_2 \subset K_1 \subset K_0$.

Заметим, что если язык порождается некоторой грамматикой, то это не означает, что он не может быть порожден грамматикой с более сильными ограничениями на правила. Приводимый ниже пример иллюстрирует этот факт.

Пример 2.8. Рассмотрим грамматику $G = (\{S, A, B\}, \{0, 1\}, \{S \rightarrow AB, A \rightarrow 0A, A \rightarrow 0, B \rightarrow 1B, B \rightarrow 1\}, S)$. Эта грамматика является контекстно-свободной. Легко показать, что $L(G) = \{0^n 1^m \mid n, m > 0\}$. Однако в примере 2.7 приведена праволинейная грамматика, порождающая тот же язык.

2.4. Машины Тьюринга

Формально *машина Тьюринга* — это совокупность $M = (Q, \Gamma, \Sigma, D, q_0, F)$, где:

Q — конечное множество состояний;

$F \subseteq Q$ — множество заключительных состояний;

Γ — множество допустимых ленточных символов; один из них, обычно обозначаемый B , — *пустой* символ;

Σ — множество входных символов: подмножество Γ , не включающее B ;

D — функция переходов, отображение из $(Q - F) \times \Gamma \rightarrow Q \times \Gamma \times \{L, R\}$; для некоторых аргументов функция D может быть не определена.

q_0 — начальное состояние.

Схематически машина Тьюринга изображена на рис. 2.1.

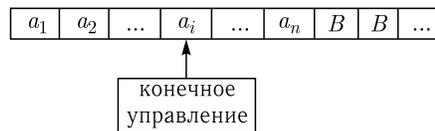


Рис. 2.1 Машина Тьюринга

Определенная таким образом машина Тьюринга называется *детерминированной*. *Недетерминированная* машина Тьюринга для каждой пары $(Q - F) \times \Gamma$ может иметь несколько возможных переходов. В начале n ячеек ленты содержат вход $w \in (\Gamma - \{B\})^*$, остальная часть ленты содержит пустые символы. Обозначим конфигурацию машины Тьюринга как (q, w, i) , где

$q \in Q$ — текущее состояние, i — выделенный элемент строки, «положение головки», w — текущее содержимое занятого участка ленты. Если головка сдвигается с ячейки, то машина должна записать в нее символ, так что лента всегда состоит из участка, состоящего из конечного числа непустых символов и бесконечного количества пустых символов.

Определим на конфигурациях отношение \vdash (шаг) следующим образом.

Пусть $(q, A_1, A_2, \dots, A_n, i)$ — конфигурация M , где $1 \leq i \leq n + 1$.

Если $1 \leq i \leq n$ и $D(q, A_i) = (p, A, R)$ (R от англ. right — правый), то $A \neq B$ и $(q, A_1 A_2 \dots A_n, i) \vdash_M (p, A_1 A_2 \dots A_{i-1} A A_{i+1} \dots A_n, i + 1)$, т. е. M печатает символ A и передвигается вправо.

Если $2 \leq i \leq n$ и $D(q, A_i) = (p, A, L)$ (L от англ. left — левый), то при $i = n$ допустимо $A = B$ и $(q, A_1 A_2 \dots A_n, i) \vdash_M (p, A_1 A_2 \dots A_{i-1} A A_{i+1} \dots A_n, i - 1)$, M печатает A и передвигается влево, но не за конец ленты.

Если $i = n + 1$, головка просматривает пустой символ B .

Если $D(q, B) = (p, A, R)$, то $A \neq B$ и $(q, A_1 A_2 \dots A_n, n + 1) \vdash_M (p, A_1 A_2 \dots A_n A, n + 2)$.

Если $D(q, B) = (p, A, L)$, то допустимо $A = B$ и $(q, A_1 A_2 \dots A_n, n + 1) \vdash_M (p, A_1 A_2 \dots A_n A, n)$.

Если две конфигурации связаны отношением \vdash_M , то мы говорим, что вторая получается из первой за один шаг. Если вторая получается из первой за конечное, включая 0, число шагов, то такое отношение будем обозначать \vdash_M^* .

Язык, допускаемый M , — это множество таких слов из T^* , которые, будучи расположены в левом конце ленты, переводят M из начального состояния q_0 с начальным положением головки в самом левом конце ленты в конечное состояние. Формально, язык, допускаемый M , это $L = \{w \mid w \in \Sigma^* \text{ и } (q_0, w, 1) \vdash_M^* (q, u, i) \text{ для некоторых } q \in F, u \in \Gamma^* \text{ и целого } i\}$.

Если M распознаёт L , то M останавливается, т. е. не имеет переходов после того, как слово допущено. Однако если слово не допущено, то возможно, что M не останавливается.

Язык, допускаемый некоторой M , называется *рекурсивно перечислимым*. Если M останавливается на всех входах, то говорят, что M задает алгоритм, а язык называется *рекурсивным*.

Существует машина Тьюринга, которая по некоторому описанию произвольной M и кодированию слова x моделирует поведение M со входом x . Такая машина Тьюринга называется *универсальной машиной Тьюринга*.

2.4.1. Неразрешимость проблемы останова. Проблема останова для машины Тьюринга формулируется следующим образом: можно ли по данной машине Тьюринга в произвольной конфигурации со строкой конечной длины непустых символов на ленте определить, остановится ли она? Говорят, что

эта проблема *рекурсивно неразрешима*: это означает, что не существует алгоритма, который для любой M в произвольной конфигурации определял бы, остановится ли в конце концов M .

Перенумеруем все машины Тьюринга и все возможные входы над алфавитом Σ . Рассмотрим язык

$$L_1 = \{x_i \mid x_i \text{ не допускается } T_i\}.$$

Ясно, что L_1 не допускается никакой M . Допустим, что это не так. Пусть L_1 допускается T_j . Следовательно $x_j \in L_1$ тогда и только тогда, когда x_j не допускается T_j . Но поскольку T_j допускает L_1 , $x_j \in L_1$ тогда и только тогда, когда допускается T_j , — противоречие. Поэтому L_1 не является рекурсивно перечислимым множеством.

Предположим, что мы имеем алгоритм (т.е. машину Тьюринга, которая всегда останавливается), позволяющий определить, остановится ли машина Тьюринга в данной конфигурации. Тогда машину Тьюринга T , допускающую L_1 , можно построить следующим образом.

1. Если дано слово x , то T перечисляет слова x_1, x_2, \dots , пока не будет $x_i = x$.

2. T генерирует кодировку машины Тьюринга T_i .

3. Управление передается гипотетической машине, которая определяет, останавливается ли T_i на входе x_i .

4. Если выясняется, что T_i не останавливается на входе x_i , то T останавливается и допускает x_i .

5. Если выясняется, что T_i останавливается на входе x_i , то управление передается универсальной машине Тьюринга, которая моделирует T_i на входе x_i . Поскольку T_i останавливается, универсальная машина Тьюринга тоже останавливается и определяет, допускает ли T_i слово x_i .

В любом случае T останавливается, допуская x_i в том случае, когда T_i отвергает x_i , и отвергая x_i , когда T_i допускает x_i .

Таким образом, из нашего предположения, что существует машина Тьюринга, которая определяет, останавливается ли произвольная машина Тьюринга, следует, что L_1 допускается некоторой машиной Тьюринга, а это противоречит доказанному выше. Это, в свою очередь, дает следующую теорему.

Теорема 2.2. *Не существует алгоритма для определения того, остановится ли произвольная машина Тьюринга в произвольной конфигурации.*

2.4.2. Класс рекурсивных множеств. Теперь можно показать, что класс рекурсивных множеств является собственным подклассом класса рекурсивно перечислимых множеств. Это значит, что существует множество, слова которого могут быть распознаны машиной Тьюринга, которая не останавливается на некоторых словах, не принадлежащих множеству, но оно

не может быть распознано никакой машиной Тьюринга, которая останавливается на всех словах. Примером такого множества является дополнение к L_1 .

Лемма 2.1. *Если множество рекурсивно, то и его дополнение рекурсивно.*

Доказательство. Если L — рекурсивное множество, $L \subseteq T^*$, то существует M , допускающая L и гарантированно останавливающаяся. Можно считать, что после допуска M больше не делает шагов. Построим M_1 по M , добавив новое состояние q как единственное допускающее состояние M_1 . Правила M_1 включают все правила M , так что M_1 симулирует M . Кроме того, для каждой пары, составленной из недопускающего состояния и ленточного символа M , для которых у M переход не определен, M_1 переходит в состояние q и затем останавливается.

Таким образом M_1 симулирует M вплоть до остановки. Если M останавливается в одном из допускающих состояний, то M_1 останавливается без допуска. Если M останавливается в одном из недопускающих состояний, то, значит, M_1 не допускает входное слово. Тогда M_1 делает еще один переход в состояние q и допускает. Ясно, что M_1 допускает $T^* \setminus L$. ■

Лемма 2.2. *Пусть x_1, x_2, \dots — нумерация всех слов некоторого конечного алфавита Σ и T_1, T_2, \dots — нумерация всех машин Тьюринга с ленточными символами, выбранными из некоторого конечного алфавита, включающего Σ . Пусть $L_2 = \{x_i \mid x_i \text{ допускается } T_i\}$. Тогда L_2 — рекурсивно перечислимое множество, дополнение которого не рекурсивно перечислимо.*

Доказательство. Слова L_2 допускаются некоторой M , работающей следующим образом (отметим, что M не обязательно останавливается на словах не из L_2).

1. Если дано x , то M перечисляет предложения x_1, x_2, \dots пока не найдет $x_i = x$, определяя тем самым, что x — это i -е слово в перечислении.
2. M генерирует M_i и передает управление универсальной машине Тьюринга, которая симулирует M_i со входом x .
3. Если M_i останавливается со входом x_i и допускает, то M останавливается и допускает; если M_i останавливается и отвергает x_i , то M останавливается и отвергает x_i . Наконец, если M_i не останавливается, то M не останавливается.
4. Таким образом, L_2 рекурсивно перечислимо, поскольку L_2 — это множество, допускаемое M . Но дополнение $\sim L_2$ к L_2 не может быть рекурсивно перечислимо, поскольку если T_j — машина Тьюринга, допускающая $\sim L_2$, то $x_j \in \sim L_2$ тогда и только тогда, когда x_j не допускается M_j . Это противоречит утверждению, что $\sim L_2$ — это язык, допускаемый T_j . ■

Теорема 2.3. *Существует рекурсивно перечислимое множество, не являющееся рекурсивным.*

Доказательство. По лемме 2.2 L_2 — рекурсивно перечислимое множество, дополнение которого не рекурсивно перечислимо. Если бы L_2 было рекурсивно, то по лемме 2.1 $\sim L_2$ было бы рекурсивно и, следовательно, рекурсивно перечислимо, что противоречит утверждению леммы 2.2. ■

2.5. Связь машин Тьюринга и грамматик типа 0

Докажем, что язык распознаётся машиной Тьюринга, если и только если он генерируется грамматикой типа 0. Для доказательства части «если» мы построим недетерминированную машину Тьюринга, которая будет недетерминированно выбирать выводы в грамматике и смотреть, является ли вывод входом. Если да, то машина допускает вход.

Для доказательства части «только если» мы построим грамматику, которая недетерминированно генерирует представления терминальной строки и затем симулирует машину Тьюринга на этой строке. Если строка допускается некоторой M , то строка конвертируется в терминальные символы, которые она представляет.

Теорема 2.4. *Если L генерируется грамматикой типа 0, то L распознаётся машиной Тьюринга.*

Доказательство. Пусть $G = (N, \Sigma, P, S)$ — грамматика типа 0 и $L = L(G)$. Опишем неформально недетерминированную машину Тьюринга M , допускающую L :

$$M = (Q, \Sigma, \Gamma, D, q_0, F)$$

где $\Gamma = N \cup \Sigma \cup \{B, \#, X\}$.

Предполагается, что последние три символа не входят в $N \cup \Sigma$.

Вначале M содержит на входной ленте $w \in \Sigma^*$. Затем M вставляет $\#$ перед w , сдвигая все символы w на одну ячейку вправо, и $\#S\#$ после w , так что содержимым ленты становится $\#w\#S\#$.

Теперь M недетерминированно симулирует вывод в G , начиная с S . Каждая сентенциальная форма вывода появляется по порядку между последними двумя $\#$. Если некоторый выбор переходов ведет к терминальной строке, то она сравнивается с w . Если они совпадают, то M допускает.

Формально, пусть M имеет на ленте $\#w\#A_1A_2 \dots A_k\#$. M передвигает недетерминированно головку по $A_1A_2 \dots A_k$, выбирая позицию i и константу r между 1 и максимальной длиной левой части любого правила вывода в P . Затем M проверяет подстроки $A_iA_{i+1} \dots A_{i+r-1}$. Если $A_iA_{i+1} \dots A_{i+r-1}$ — левая часть некоторого правила вывода из P , то она может быть заменена на правую часть. Тогда M может сдвинуть $A_{i+r}A_{i+r+1} \dots A_k\#$ либо влево,

либо вправо, освобождая или заполняя место, если правая часть имеет длину, отличную от r .

Из этой простой симуляции выводов в G видно, что M печатает на ленте строку вида $\#w\#y\#$, $y \in V^*$, в точности, если $S \Rightarrow_{G^*} y$. Если $y = w$, то M допускает L . ■

Теорема 2.5. *Если L распознаётся некоторой (детерминированной) машиной Тьюринга, то L генерируется грамматикой типа 0.*

Доказательство. Пусть $M = (Q, \Sigma, \Gamma, D, q_0, F)$ допускает L . Построим грамматику G , которая недетерминированно генерирует две копии представления некоторого слова из Σ^* и затем симулирует поведение M на одной из копий. Если M допускает слово, то G трансформирует вторую копию в терминальную строку. Если M не допускает L , то вывод никогда не приводит к терминальной строке.

Формально, пусть

$$G = (N, \Sigma, P, A_1), \quad \text{где } N = ([\Sigma \cup \{e\}] \times \Gamma) \cup Q \cup \{A_1, A_2, A_3\},$$

и P состоит из следующих правил.

1. $A_1 \rightarrow q_0 A_2$.
2. $A_2 \rightarrow [a, a] A_2$ для каждого $a \in \Sigma$.
3. $A_2 \rightarrow A_3$.
4. $A_3 \rightarrow [e, B] A_3$.
5. $A_3 \rightarrow e$.
6. $q[a, C] \rightarrow [a, E] p$ для каждого $a \in \Sigma \cup \{e\}$ и каждого $q \in Q$ и такого $C \in \Gamma$, что $D(q, C) = (p, E, R)$.
7. $[b, I] q[a, C] \rightarrow p[b, I][a, J]$ для каждого C, J, I из Γ , a и b из $\Sigma \cup \{e\}$ и таких q из Q , что $D(q, C) = (p, J, L)$.
8. $[a, C] q \rightarrow q a q$, $q[a, C] \rightarrow q a q$, $q \rightarrow e$ для каждого $a \in \Sigma \cup \{e\}$, $C \in \Gamma$, $q \in F$.

Используя правила 1 и 2, получаем

$$A_1 \Rightarrow^* q_0 [a_1, a_1][a_2, a_2] \dots [a_n, a_n] A_2,$$

где $a_i \in \Sigma$ для некоторого i .

Предположим, что M допускает строку $a_1 a_2 \dots a_n$. Тогда найдется m , для которого M использует не более, чем m ячеек справа от входа. Используя правило 3, затем m раз правило 4 и, наконец, правило 5, имеем

$$A_1 \Rightarrow^* q_0 [a_1, a_1][a_2, a_2] \dots [a_n, a_n][e, B]^m.$$

Начиная с этого момента могут быть использованы только правила 6 и 7, пока не сгенерируется допускающее состояние. Отметим, что первые компоненты ленточных символов в $(\Sigma \cup \{e\}) \times \Gamma$ никогда не меняются. Индукцией по числу шагов M можно показать, что если

$$(q_0, a_1 a_2 \dots a_n, 1) \vdash_M^* (q, X_1 X_2 \dots X_s, r),$$

то

$$q_0[a_1, a_1][a_2, a_2] \dots [a_n, a_n][e, B]^m \Rightarrow_G^* \\ \Rightarrow_G^*[a_1, X_1][a_2, X_2] \dots [a_{r-1}, X_{r-1}]q[a_r, X_r] \dots [a_{n+m}, X_{n+m}],$$

где a_1, a_2, \dots, a_n принадлежат Σ , $a_{n+1} = a_{n+2} = \dots = a_{n+m} = e$, X_1, X_2, \dots, X_{n+m} принадлежат Γ и $X_{s+1} = X_{s+2} = \dots = X_{n+m} = B$.

Предположение индукции тривиально для 0 шагов. Предположим, что оно справедливо для $k - 1$ шагов. Пусть

$$(q_0, a_1 a_2 \dots a_n, 1) \vdash_M^* (q_1, X_1 X_2 \dots X_r, j_1) \vdash_M (q_2, Y_1 Y_2 \dots Y_s, j_2)$$

за k шагов. По предположению индукции

$$q_0[a_1, a_1][a_2, a_2] \dots [a_n, a_n][e, B]^m \Rightarrow_G^* \\ \Rightarrow_G^*[a_1, X_1][a_2, X_2] \dots [a_{r-1}, X_{r-1}]q_1[a_{j_1}, X_{j_1}] \dots [a_{n+m}, X_{n+m}].$$

Пусть $E = L$, если $j_2 = j_1 - 1$, и $E = R$, если $j_2 = j_1 + 1$. В этом случае $D(q_1, X_{j_1}) = (q_2, Y_{j_1}, E)$.

По правилу 6 или 7

$$q_1[a_{j_1}, X_{j_1}] \rightarrow [a_{j_1}, Y_{j_1}]q_2$$

или

$$[a_{j_1-1}, X_{j_1-1}]q_1[a_{j_1}, X_{j_1}] \rightarrow q_2[a_{j_1-1}, X_{j_1-1}][a_{j_1}, Y_{j_1}],$$

в зависимости от того, равно ли E значению R или L . Теперь $X_i = Y_i$ для всех $i \neq j_1$.

Таким образом,

$$q_0[a_1, a_1][a_2, a_2] \dots [a_n, a_n][e, B]^m \Rightarrow_G^*[a_1, Y_1]q_2[a_{j_2}, Y_{j_2}] \dots [a_{n+m}, Y_{n+m}],$$

что доказывается предположением индукции.

По правилу 8, если $q \in F$, легко показать, что

$$[a_1, X_1] \dots q[a_j, X_j] \dots [a_{n+m}, X_{n+m}] \Rightarrow_G^* a_1 a_2 \dots a_n.$$

Таким образом, G может генерировать $a_1 a_2 \dots a_n$, если $a_1 a_2 \dots a_n$ допускается M . Таким образом, $L(G)$ включает все слова, допускаемые M . Для завершения доказательства необходимо показать, что все слова из $L(G)$ допускаются M . Индукцией доказывается, что $A_1 \Rightarrow_G^* w$, только если w допускается M .

Тем самым мы доказали, что если слово допускается M , то оно порождается G . Если слово не допускается, то M не попадает в заключительное состояние и слово не порождается. ■

2.6. Линейно ограниченные автоматы и их связь с контекстно-зависимыми грамматиками

Каждый КЗ-язык является рекурсивным, но обратное не верно. Покажем, что существует алгоритм, позволяющий для произвольного КЗ-языка L в алфавите T и произвольной цепочки $w \in T^*$ определить, принадлежит ли w языку L .

Теорема 2.6. *Каждый контекстно-зависимый язык является рекурсивным языком.*

Доказательство. Пусть L — контекстно-зависимый язык. Тогда существует некоторая неукорачивающая грамматика $G = (N, T, P, S)$, порождающая L .

Пусть $w \in T^*$ и $|w| = n$. Если $n = 0$, т. е. $w = e$, то принадлежность $w \in L$ проверяется тривиальным образом. Так что будем предполагать, что $n > 0$.

Определим множество T_m как множество строк $u \in (N \cup T)^+$ длины не более n , таких, что вывод $S \Rightarrow^* u$ имеет не более m шагов. Ясно, что $T_0 = \{S\}$.

Легко показать, что T_m можно получить из T_{m-1} , просматривая, какие строки с длиной, меньшей или равной n , можно вывести из строк из T_{m-1} применением одного правила, т. е.

$$T_m = T_{m-1} \cup \{u \mid v \Rightarrow u \text{ для некоторого } v \in T_{m-1}, \text{ где } |u| \leq n\}.$$

Если $S \Rightarrow^* u$ и $|u| \leq n$, то $u \in T_m$ для некоторого m . Если из S не выводится u или $|u| > n$, то u не принадлежит T_m ни для какого m .

Очевидно, что $T_m \supseteq T_{m-1}$ для всех $m \geq 1$. Поскольку T_m зависит только от T_{m-1} , то из $T_m = T_{m-1}$ следует $T_m = T_{m+1} = T_{m+2} = \dots$. Процедура будет вычислять T_1, T_2, T_3, \dots пока для некоторого m не окажется $T_m = T_{m-1}$. Если w не принадлежит T_m , то w не принадлежит и $L(G)$, поскольку для $j > m$ выполняется $T_j = T_m$. Если $w \in T_m$, то $S \Rightarrow^* w$.

Покажем, что существует такое m , что $T_m = T_{m-1}$. Поскольку для каждого $i \geq 1$ справедливо $T_i \supseteq T_{i-1}$, то из $T_i \neq T_{i-1}$ следует, что число элементов в T_i по крайней мере на 1 больше, чем в T_{i-1} . Пусть $|N \cup T| = k$. Тогда число строк в $(N \cup T)^+$ длиной, меньшей или равной n , равно $k + k^2 + \dots + k^n \leq nk^n$. Только эти строки могут быть в любом T_i . Значит, $T_m = T_{m-1}$ для некоторого $m \leq nk^n$. Таким образом, процедура, вычисляющая T_i для всех $i \geq 1$ до тех пор, пока не будут найдены два равных множества, гарантированно заканчивается; значит, это алгоритм. ■

Линейно ограниченный автомат (ЛОА) — это недетерминированная машина Тьюринга с одной лентой, которая никогда не выходит за пределы $|w|$ ячеек, где w — вход. Формально линейно ограниченный автомат обозначается

$M = (Q, \Sigma, \Gamma, D, q_0, F)$. Обозначения имеют тот же смысл, что и для машин Тьюринга: Q — это множество состояний, $F \subseteq Q$ — множество заключительных состояний, Γ — множество ленточных символов, $\Sigma \subseteq \Gamma$ — множество входных символов, $q_0 \in Q$ — начальное состояние, D — отображение из $Q \times \Gamma$ в подмножество $Q \times \Gamma \times \{L, R\}$.

Множество Σ содержит два специальных символа, обычно обозначаемые \textcircled{C} и \textcircled{R} , — левый и правый концевые маркеры соответственно. Эти символы располагаются сначала по концам входа, и их функция — предотвратить переход головки за пределы области, в которой расположен вход.

Конфигурация M и отношение \vdash_M , связывающее две конфигурации, если вторая может быть получена из первой применением D , определяются так же, как и для машин Тьюринга. Конфигурация M обозначается как $(q, A_1, A_2, \dots, A_n, i)$, где $q \in Q$, $A_1, A_2, \dots, A_n \in \Gamma$, i — целое от 1 до n . Предположим, что $(p, A, L) \in D(q, A_i)$ и $i > 1$. Будем говорить, что

$$(q, A_1, A_2 \dots A_n, i) \vdash_M (p, A_1, A_2 \dots A_{i-1} A A_{i+1} \dots A_n, i-1).$$

Если $(p, A, R) \in D(q, A_i)$ и $i < n$, то будем говорить, что

$$(q, A_1, A_2, \dots, A_n, i) \vdash_M (p, A_1, A_2 \dots A_{i-1} A A_{i+1} \dots A_n, i+1),$$

т. е. M печатает A поверх A_i , меняет состояние на p и передвигает головку влево или вправо, но не за пределы области, в которой символы располагались исходно. Как обычно, определим отношение \vdash_M^* как

$$(q, \alpha, i) \vdash_M^* (q, \alpha, i);$$

если

$$(q_1, \alpha_1, i_1) \vdash_M^* (q_2, \alpha_2, i_2)$$

и

$$(q_2, \alpha_2, i_2) \vdash^* (q_3, \alpha_3, i_3),$$

то

$$(q_1, \alpha_1, i_1) \vdash_M^* (q_3, \alpha_3, i_3).$$

Язык, допускаемый M , — это $\{w \mid w \in (\Sigma \setminus \{\textcircled{C}, \textcircled{R}\})^*$ и $(q_0, \textcircled{C}w\textcircled{R}, 1) \vdash_M^* (q, \alpha, i)$ для некоторого $q \in F$, $\alpha \in \Gamma^*$ и целого $i\}$.

Будем называть M *детерминированным*, если $D(q, A)$ содержит не более одного элемента для любых $q \in Q$, $A \in \Gamma$. Неизвестно, совпадают ли классы множеств, допускаемых детерминированными и недетерминированными ЛОА. Ясно, что любое множество, допускаемое недетерминированным ЛОА, допускается некоторой детерминированной МТ (машиной Тьюринга). Однако число ячеек ленты, требуемой этой МТ, может экспоненциально зависеть от длины входа.

Класс множеств, допускаемых ЛОА, в точности совпадает с классом контекстно-зависимых языков.

Теорема 2.7. *Если L — контекстно-зависимый язык, то L допускается ЛОА.*

Доказательство. Пусть $G = (V_N, V_T, P, S)$ — контекстно-зависимая грамматика. Построим ЛОА M , допускающий язык $L(G)$. Не вдаваясь в детали построения M , поскольку он довольно сложен, рассмотрим схему его работы. В качестве ленточных символов будем рассматривать пары (s_i^1, s_i^2) , где $s_i^1 \in \Sigma$, $\Sigma = V_T \cup \{ @, \$ \}$, $s_i^2 \in \Gamma$, $\Gamma = V_T \cup V_N \cup \{ B \}$. В начальной конфигурации лента содержит $(@, B), (a_1, B), \dots (a_n, B), (\$, B)$, где $a_1, \dots, a_n = w$ — входная цепочка, $n = |w|$. Цепочку символов s_1^1, \dots, s_n^1 будем называть первым треком, s_1^2, \dots, s_n^2 — вторым треком. Первый трек будет содержать входную строку x с концевыми маркерами. Второй трек будет использоваться для вычислений. На первом шаге M помещает символ S в самую левую ячейку второго трека. Затем M выполняет процедуру генерации в соответствии со следующими шагами.

1. Процедура выбирает подстроку символов α из второго трека, такую, что $\alpha \rightarrow \beta \in P$.
2. Процедура заменяет подстроку α на β , перемещая, если необходимо, вправо символы справа от α . Если эта операция могла бы привести к перемещению символа за правый концевой маркер, то ЛОА останавливается.
3. Процедура недетерминированно выбирает, перейти на шаг 1 или завершиться.

На выходе из процедуры первый трек все еще содержит строку x , а второй трек содержит строку γ , такую, что $S \Rightarrow_G^* \gamma$. ЛОА сравнивает символы первого трека с соответствующими символами второго трека. Если сравнение неуспешно, то строки символов первого и второго треков неодинаковы и ЛОА останавливается без допуска. Если строки одинаковы, то ЛОА останавливается и допускает.

Если $x \in L(G)$, то существует некоторая последовательность шагов, на которой ЛОА строит x на втором треке и допускает вход. Аналогично, для того, чтобы ЛОА допустил x , должна существовать последовательность шагов, такая, что x может быть построен на втором треке. Таким образом, должен быть вывод x из S в G . ■

Отметим схожесть этих рассуждений и рассуждений в случае произвольной грамматики. Тогда промежуточные сентенциальные формы могли иметь длину, произвольно большую по сравнению с длиной входа. Как следствие, требовалась вся мощь машин Тьюринга. В случае контекстно-зависимых

грамматик промежуточные сентенциальные формы не могут быть длиннее входа.

Теорема 2.8. *Если L допускается ЛОА, то L — контекстно-зависимый язык.*

Доказательство. Конструкция КЗГ по ЛОА аналогична конструкции грамматики типа 0, моделирующей машину Тьюринга. Различие заключается в том, что нетерминалы КЗГ должны указывать не только текущее и исходное содержимое ячеек ленты ЛОА, но и то, является ли ячейка соседней справа или слева с концевым маркером. Кроме того, состояние ЛОА должно комбинироваться с символом под головкой, поскольку КЗГ не может иметь отдельные символы для концевых маркеров и состояния ЛОА, так как эти символы должны были бы быть заменены на ϵ , когда строка превращается в терминальную.

1-я группа правил — моделирование начальной конфигурации вида

$(q_0, @w\#, 1), |w| > 1:$

$A_1 \rightarrow [q_0, @, a.a]A_2;$

$A_2 \rightarrow [a.a]A_2;$

$A_2 \rightarrow [a.a, \#].$

2-я группа правил — моделирование движения на левом конце цепочки при $q \in F:$

$[q, @, X, a] \rightarrow [@, p, X, a],$ если $(p, @, R) \in D(q, @);$

$[@, q, X, a] \rightarrow [p, @, Y, a],$ если $(p, Y, L) \in D(q, X);$

$[@, q, X, a][Z, b] \rightarrow [@, Y, a][p, Z, b],$ если $(p, Y, R) \in D(q, X).$

3-я группа правил — моделирование движения в середине цепочки при $q \in F:$

$[q, X, a][Z, b] \rightarrow [Y, a][p, Z, b],$ если $(p, Y, R) \in D(q, X);$

$[Z, b][q, X, a] \rightarrow [p, Z, b][Y, a],$ если $(p, Y, L) \in D(q, X).$

$[q, X, a][Z, b, \#] \rightarrow [Y, a][p, Z, b, \#],$ если $(p, Y, R) \in D(q, X).$

4-я группа правил — моделирование движения на правом конце цепочки при $q \in F:$

$[q, X, a, \#] \rightarrow [Y, a, p, \#],$ если $(p, Y, R) \in D(q, X);$

$[X, a, q, \#] \rightarrow [p, X, a, \#],$ если $(p, \#, L) \in D(q, \#);$

$[Z, b][q, X, a, \#] \rightarrow [p, Z, b][Y, a, \#],$ если $(p, Y, L) \in D(q, X).$

5-я группа правил — восстановление входного символа в состоянии $q \in F:$

на левом конце цепочки

$[q, \#, a, X] \rightarrow a, [@, q, X, a] \rightarrow a;$

в середине цепочки

$[q, X, a] \rightarrow a;$

на правом конце цепочки

$[q, X, a, \#] \rightarrow a, [X, a, q, \#] \rightarrow a.$

6-я группа правил — восстановление входной цепочки слева-направо:

$$\begin{aligned} a[x, b] &\rightarrow ab, \\ a[x, b, \#] &\rightarrow ab. \end{aligned}$$

7-я группа правил — восстановление входной цепочки справа-налево:

$$\begin{aligned} [x, a]b &\rightarrow ab, \\ [a, x, a]b &\rightarrow ab. \end{aligned}$$

Отдельно для односимвольных цепочек — моделирование начальной конфигурации вида $(q_0, @, a, \#, 1)$:

$$A_1 \rightarrow [q_0, @, a, a, \#].$$

Если $q \in F$:

$$[q, @, X, a, \#] \rightarrow [a, p, X, a, \#], \text{ если } (p, @, R) \in D(q, @);$$

$$[@, q, X, a, \#] \rightarrow [p, @, Y, a, \#], \text{ если } (p, Y, L) \in D(q, X);$$

$$[@, q, X, a, \$] \rightarrow [a, Y, a, p, \#], \text{ если } (p, Y, R) \in D(q, X);$$

$$[@, X, a, q, \#] \rightarrow [a, p, X, a, \#], \text{ если } (p, \$, L) \in D(q, \$).$$

Если $q \in F$:

$$[q, @, X, a, \#] \rightarrow a;$$

$$[@, q, X, a, \#] \rightarrow a;$$

$$[@, X, a, q, \#] \rightarrow a. \quad \blacksquare$$

Теорема 2.9. *Существуют рекурсивные множества, не являющиеся контекстно-зависимыми.*

Доказательство. Все строки в $\{0, 1\}^*$ можно занумеровать. Пусть x_i — i -е слово. Мы можем занумеровать все грамматики типа 0, терминальными символами которых являются 0 и 1. Поскольку имена переменных не важны и каждая грамматика имеет конечное их число, можно предположить, что существует счетное число переменных.

Представим переменные в двоичной кодировке как 01, 011, 0111, 01111 и т.д. Предположим, что 01 всегда является стартовым символом. Кроме того, в этой кодировке терминал 0 будет представляться как 00, а терминал 1 как 001. Символ « \rightarrow » представляется как 0011, а запятая как 00111. Любая грамматика с терминалами 0 и 1 может быть представлена строкой правил, использующей стрелку (0011) для разделения левой и правой частей, и запятой (00111) для разделения правил. Строки, представляющие символы, используемые в правилах, — это 00, 001 и 01^i для $i = 1, 2, \dots$. Множество используемых переменных определяется неявно правилами.

Отметим, что не все строки из 0 и 1 представляют грамматики, причем не обязательно КЗГ. Однако по данной строке легко можно сказать, представляет ли она КЗГ. i -ю грамматику можно найти, генерируя двоичные строки в описанном порядке, пока не сгенерируется i -я строка, являющаяся КЗГ. Поскольку имеется бесконечное число КЗГ, их можно занумеровать в некотором порядке G_1, G_2, \dots

Определим $L = \{x_i \mid x_i \notin L(G_i)\}$. Ясно, что L рекурсивно. По строке x_i легко можно определить i и затем определить G_i . По теореме 2.6 имеется алгоритм, определяющий для x_i , принадлежит ли она $L(G_i)$, поскольку G_i — КЗГ. Таким образом, имеется алгоритм, определяющий для любой x , принадлежит ли она G .

Покажем теперь, что L не генерируется никакой КЗ-грамматикой. Допустим, что L генерируется КЗ-грамматикой G_i . Сначала предположим, что $x_i \in L$. Поскольку $L(G_i) = L$, $x_i \in L(G_i)$. Но тогда по определению $x_i \notin L(G_i)$ — противоречие. Таким образом, предположим, что $x_i \notin L$. Поскольку $L(G_i) = L$, $x_i \notin L(G_i)$. Но тогда по определению $x_i \in L(G_i)$ — снова противоречие. Из этого можно заключить, что L не генерируется G_i . Поскольку приведенный выше аргумент справедлив для каждой КЗ-грамматики G_i в перечислении и поскольку перечисление содержит все КЗ-грамматики, можно заключить, что L — не КЗ-язык. Поэтому L — рекурсивное множество, не являющееся контекстно-зависимым. ■

Глава 3

ЛЕКСИЧЕСКИЙ АНАЛИЗ

Основная задача лексического анализа — разбить входной текст, состоящий из последовательности одиночных символов, на последовательность слов, или лексем, т. е. выделить эти слова из непрерывной последовательности символов. Все символы входной последовательности с этой точки зрения разделяются на символы, принадлежащие каким-либо лексемам, и символы, разделяющие лексемы (разделители). В некоторых случаях между лексемами может и не быть разделителей. С другой стороны, в некоторых языках лексемы могут содержать незначащие символы (например, символ пробела в Фортране). В Си разделительное значение символов-разделителей может блокироваться («\» в конце строки внутри «...»).

Обычно все лексемы делятся на классы. Примерами таких классов являются числа (целые, восьмеричные, шестнадцатеричные, действительные и т. п.), идентификаторы, строки. Отдельно выделяются ключевые слова и символы пунктуации (иногда их называют символами-ограничителями). Как правило, ключевые слова — это некоторое конечное подмножество идентификаторов. В некоторых языках (например, ПЛ/1) смысл лексемы может зависеть от ее контекста и невозможно провести лексический анализ в отрыве от синтаксического.

Для осуществления двух дальнейших фаз анализа лексический анализатор (ЛА) выдает информацию двух типов: для синтаксического анализатора, работающего вслед за лексическим, существенна информация о последовательности классов лексем, ограничителей и ключевых слов, а для контекстного анализатора, работающего вслед за синтаксическим, существенна информация о конкретных значениях отдельных лексем (идентификаторов, чисел и т. п.).

Таким образом, общая схема работы лексического анализатора такова. Сначала выделяется отдельная лексема (при этом, возможно, используются символы-разделители). Ключевые слова распознаются явным выделением непосредственно из текста, либо сначала выделяется идентификатор, а затем делается проверка на принадлежность его множеству ключевых слов.

Если выделенная лексема является ограничителем, то этот ограничитель (точнее, некоторый его признак) выдается как результат лексического

анализа. Если выделенная лексема является ключевым словом, то выдается признак соответствующего ключевого слова. Если выделенная лексема является идентификатором — выдается признак идентификатора, а сам идентификатор сохраняется отдельно. Наконец, если выделенная лексема принадлежит какому-либо из других классов лексем (например, лексема представляет собой число, строку и т. п.), то выдается признак соответствующего класса, а значение лексемы сохраняется отдельно.

Лексический анализатор может быть как самостоятельной фазой трансляции, так и подпрограммой, работающей по принципу «дай лексему». В первом случае (рис. 3.1, а) выходом анализатора является файл лексем, во втором — (рис. 3.1, б) лексема выдается при каждом обращении к анализатору (при этом, как правило, признак класса лексемы возвращается как результат функции «лексический анализатор», а значение лексемы передается через глобальную переменную). С точки зрения обработки значений лексем, анализатор может либо просто выдавать значение каждой лексемы, при этом построение таблиц объектов (идентификаторов, строк, чисел и т. п.) переносится на более поздние фазы, либо он может самостоятельно строить таблицы объектов. В этом случае в качестве значения лексемы выдается указатель на вход в соответствующую таблицу.



Рис. 3.1

Работа лексического анализатора задается некоторым конечным автоматом. Однако непосредственное описание конечного автомата неудобно с практической точки зрения. Поэтому для задания лексического анализатора, как правило, используется либо регулярное выражение, либо праволинейная грамматика. Все три формализма (конечных автоматов, регулярных выражений и праволинейных грамматик) имеют одинаковую выразительную мощность. В частности, по регулярному выражению или праволинейной грамматике можно сконструировать конечный автомат, распознающий тот же язык.

3.1. Регулярные множества и выражения

Введем понятие регулярного множества, играющее важную роль в теории формальных языков.

Регулярное множество в алфавите T определяется рекурсивно следующим образом.

1. \emptyset (пустое множество) — регулярное множество в алфавите T .
2. $\{e\}$ — регулярное множество в алфавите T (e — пустая цепочка).
3. $\{a\}$ — регулярное множество в алфавите T для каждого $a \in T$.
4. Если P и Q — регулярные множества в алфавите T , то регулярными являются и множества:

а) $P \cup Q$ (объединение),

б) PQ (конкатенация, т. е. множество $\{pq | p \in P, q \in Q\}$),

в) P^* (итерация: $P^* = \bigcup_{n=0}^{\infty} P^n$).

5. Ничто другое не является регулярным множеством в алфавите T .

Итак, множество в алфавите T регулярно тогда и только тогда, когда оно либо \emptyset , либо $\{e\}$, либо $\{a\}$ для некоторого $a \in T$, либо его можно получить из этих множеств применением конечного числа операций объединения, конкатенации и итерации.

Приведенное выше определение регулярного множества позволяет ввести следующую удобную форму его записи, называемую регулярным выражением.

Регулярное выражение в алфавите T и обозначаемое им регулярное множество в алфавите T определяются рекурсивно следующим образом.

1. \emptyset — регулярное выражение, обозначающее регулярное множество \emptyset .
2. e — регулярное выражение, обозначающее регулярное множество $\{e\}$.
3. a — регулярное выражение, обозначающее регулярное множество $\{a\}$.
4. Если p и q — регулярные выражения, обозначающие регулярные множества P и Q соответственно, то:

а) $(p|q)$ — регулярное выражение, обозначающее регулярное множество $P \cup Q$,

б) (pq) — регулярное выражение, обозначающее регулярное множество PQ ,

в) (p^*) — регулярное выражение, обозначающее регулярное множество P^* .

5. ничто другое не является регулярным выражением в алфавите T .

Мы будем опускать лишние скобки в регулярных выражениях, договорившись о том, что операция итерации имеет наивысший приоритет, затем идет операция конкатенации, наконец, операция объединения имеет наименьший приоритет.

Кроме того, мы будем пользоваться записью p^+ для обозначения pp^* . Таким образом, запись $(a|((ba)(a^*)))$ эквивалентна $a|ba^+$.

Мы также будем использовать запись $L(r)$ для регулярного множества, обозначаемого регулярным выражением r .

Пример 3.1. Несколько регулярных выражений и обозначаемых ими регулярных множеств:

- а) $a(e|a)|b$ обозначает множество $\{a, b, aa\}$;
- б) $a(a|b)^*$ обозначает множество всевозможных цепочек, состоящих из a и b , начинающихся с a ;
- в) $(a|b)^*(a|b)(a|b)^*$ обозначает множество всех непустых цепочек, состоящих из a и b , т. е. множество $\{a, b\}^+$;
- г) $((0|1)(0|1)(0|1))^*$ обозначает множество всех цепочек, состоящих из нулей и единиц, длины которых делятся на 3.

Ясно, что для каждого регулярного множества можно найти регулярное выражение, обозначающее это множество, и наоборот. Более того, для каждого регулярного множества существует бесконечно много обозначающих его регулярных выражений.

Будем говорить, что регулярные выражения *равны*, или *эквивалентны* ($=$), если они обозначают одно и то же регулярное множество.

Существуют алгебраические законы, позволяющие осуществлять эквивалентные преобразования регулярных выражений.

Лемма 3.1. Пусть p , q и r — регулярные выражения. Тогда справедливы следующие соотношения:

- 1) $p|q = q|p$; 5) $p(q|r) = pq|pr$; 9) $p^* = p|p^*$;
- 2) $\emptyset^* = \epsilon$; 6) $(p|q)r = pr|qr$; 10) $(p^*)^* = p^*$;
- 3) $p|(q|r) = (p|q)|r$; 7) $pe = ep = p$; 11) $p|p = p$;
- 4) $p(qr) = (pq)r$; 8) $\emptyset p = p\emptyset = \emptyset$; 12) $p|\emptyset = p$.

Следствие. Для любого регулярного выражения существует эквивалентное регулярное выражение, которое либо есть \emptyset , либо не содержит в своей записи \emptyset .

В дальнейшем будем рассматривать только регулярные выражения, не содержащие в своей записи \emptyset .

При практическом описании лексических структур бывает полезно сопоставлять регулярным выражениям некоторые имена и ссылаться на них

по этим именам. Для определения таких имен мы будем использовать запись вида

$$\begin{aligned}d_1 &= r_1, \\d_2 &= r_2, \\&\dots, \\d_n &= r_n,\end{aligned}$$

где d_i — различные имена, а каждое r_i — регулярное выражение над символами $T \cup \{d_1, d_2, \dots, d_{i-1}\}$, т.е. символами основного алфавита и ранее определенными символами (именами). Таким образом, для любого r_i можно построить регулярное выражение над T , повторно заменяя имена регулярных выражений на обозначаемые ими регулярные выражения.

Пример 3.2. *Использование имен для обозначения регулярных выражений.*

а) Регулярные выражения для множества идентификаторов:

$$\begin{aligned}Letter &= a|b|c| \dots |x|y|z, \\Digit &= 0|1| \dots |9, \\Identifier &= Letter(Letter|Digit)^*.\end{aligned}$$

б) Регулярные выражения для множества чисел в десятичной записи:

$$\begin{aligned}Digit &= 0|1| \dots |9, \\Integer &= Digit^+, \\Fraction &= .Integer|e, \\Exponent &= (E(+|-|e)Integer)|e, \\Number &= Integer Fraction Exponent.\end{aligned}$$

3.2. Конечные автоматы

Регулярные выражения, введенные ранее, служат для *описания* регулярных множеств. Для *распознавания* регулярных множеств служат *конечные автоматы*.

Недетерминированный конечный автомат (НКА) — это пятерка $M = (Q, T, D, q_0, F)$, где:

- 1) Q — конечное множество *состояний*;
- 2) T — конечное множество допустимых *входных символов* (входной алфавит);
- 3) D — *функция переходов* (отображающая множество $Q \times (T \cup \{e\})$ во множество подмножеств Q), определяющая поведение управляющего устройства;
- 4) $q_0 \in Q$ — *начальное состояние* управляющего устройства;
- 5) $F \subseteq Q$ — множество *заключительных состояний*.

Работа конечного автомата представляет собой некоторую последовательность шагов, или тактов. *Такт* определяется текущим состоянием управляющего устройства и входным символом, обозреваемым в данный момент входной головкой. Сам шаг состоит из изменения состояния и, возможно, сдвига входной головки на одну ячейку вправо (рис. 3.2).

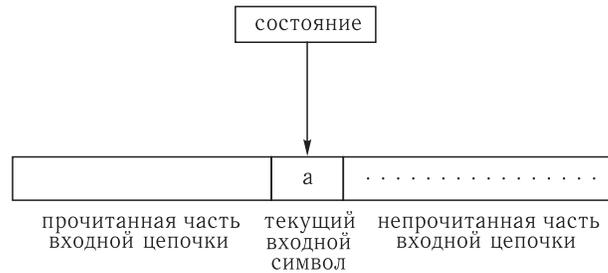


Рис. 3.2

Недетерминизм автомата заключается в том, что, во-первых, находясь в некотором состоянии и обозревая текущий символ, автомат может перейти в одно из, вообще говоря, *нескольких* возможных состояний и, во-вторых, автомат может делать переходы по ϵ .

Пусть $M = (Q, T, D, q_0, F)$ — НКА. *Конфигурацией* автомата M называется пара $(q, w) \in Q \times T^*$, где q — текущее состояние управляющего устройства, а w — цепочка символов на входной ленте, состоящая из символа под головкой и всех символов справа от него. Конфигурация (q_0, w) называется *начальной*, а конфигурация (q, ϵ) , где $q \in F$, — *заключительной* (или *допускающей*). *Тактом* автомата M называется бинарное отношение \vdash , определенное на конфигурациях M следующим образом: если $p \in D(q, a)$, где $a \in T \cup \{\epsilon\}$, то $(q, aw) \vdash (p, w)$ для всех $w \in T^*$.

Будем обозначать символом \vdash^+ (\vdash^*) транзитивное (рефлексивно-транзитивное) замыкание отношения \vdash .

Будем говорить, что автомат M *допускает* цепочку w , если $(q_0, w) \vdash^* (q, \epsilon)$ для некоторого $q \in F$. Языком, *допускаемым* (*распознаваемым, определяемым*) автоматом M (обозначается $L(M)$), называется множество входных цепочек, допускаемых автоматом M :

$$L(M) = \{w | w \in T^* \text{ и } (q_0, w) \vdash^* (q, \epsilon) \text{ для некоторого } q \in F\}.$$

Важным частным случаем недетерминированного конечного автомата является детерминированный конечный автомат, который на каждом такте работы имеет возможность перейти не более чем в одно состояние и не может делать переходы по ϵ .

Пусть $M = (Q, T, D, q_0, F)$ — НКА. Будем называть M *детерминированным конечным автоматом* (ДКА), если выполняются следующие два условия:

- 1) $D(q, e) = \emptyset$ для любого $q \in Q$;
- 2) $D(q, a)$ содержит не более одного элемента для любых $q \in Q$ и $a \in T$.

Так как функция переходов ДКА содержит не более одного элемента для любой пары аргументов, для ДКА мы будем пользоваться записью $D(q, a) = p$ вместо $D(q, a) = \{p\}$.

Конечный автомат может быть изображен графически в виде *диаграммы*, представляющей собой ориентированный граф, в котором каждому состоянию соответствует вершина, а дуга, помеченная символом $a \in T \cup \{e\}$, соединяет две вершины p и q , если $p \in D(q, a)$. На диаграмме выделяются начальное и заключительные состояния (в примерах ниже — соответственно входящей стрелкой и двойным контуром).

Пример 3.3. Пусть $L = L(r)$, где $r = (a|b)^* a(a|b)(a|b)$.

- а) Недетерминированный конечный автомат M , допускающий язык L :

$$M = \{\{1, 2, 3, 4\}, \{a, b\}, D, 1, \{4\}\},$$

где функция переходов D определяется так:

$$\begin{aligned} D(1, a) &= \{1, 2\}, & D(3, a) &= \{4\}, \\ D(1, b) &= \{1\}, & D(2, b) &= \{3\}, \\ D(2, a) &= \{3\}, & D(3, b) &= \{4\}. \end{aligned}$$

Диаграмма автомата приведена на рис. 3.3, а.

- б) Детерминированный конечный автомат M , допускающий язык L :

$$M = \{\{1, 2, 3, 4, 5, 6, 7, 8\}, \{a, b\}, D, 1, \{3, 5, 6, 8\}\},$$

где функция переходов D определяется так:

$$\begin{aligned} D(1, a) &= 2, & D(5, a) &= 8, \\ D(1, b) &= 1, & D(5, b) &= 6, \\ D(2, a) &= 4, & D(6, a) &= 2, \\ D(2, b) &= 7, & D(6, b) &= 1, \\ D(3, a) &= 3, & D(7, a) &= 8, \\ D(3, b) &= 5, & D(7, b) &= 6, \\ D(4, a) &= 3, & D(8, a) &= 4, \\ D(4, b) &= 5, & D(8, b) &= 7. \end{aligned}$$

Диаграмма автомата приведена на рис. 3.3, б.

Пример 3.4. Диаграмма автомата, допускающего множество чисел в десятичной записи, приведена на рис. 3.4.

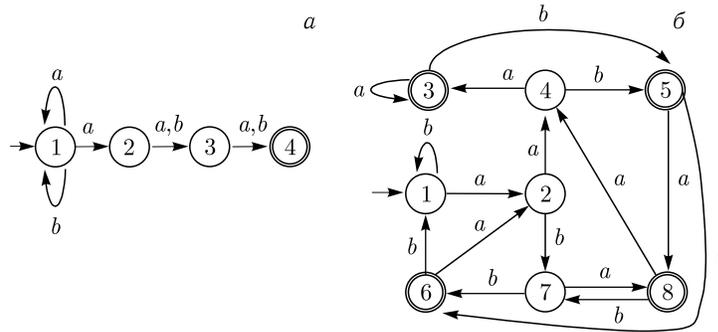


Рис. 3.3

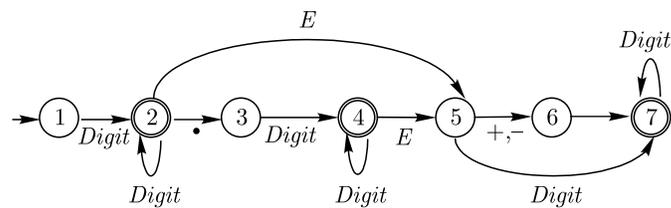


Рис. 3.4

Пример 3.5. Анализ цепочек.

- а) При анализе цепочки $w = ababa$ автомат из примера 3.3, *a* может сделать следующую последовательность тактов:

$$(1, ababa) \vdash (1, baba) \vdash (1, aba) \vdash (2, ba) \vdash (3, a) \vdash (4, e).$$

Состояние 4 является заключительным, следовательно, цепочка w допускается этим автоматом.

- б) При анализе цепочки $w = ababab$ автомат из примера 3.3, *b* должен сделать следующую последовательность тактов:

$$(1, ababab) \vdash (2, babab) \vdash (7, abab) \vdash (8, bab) \vdash (7, ab) \vdash (8, b) \vdash (7, e).$$

Так как состояние 7 не является заключительным, цепочка w не допускается этим автоматом.

3.3. Интерпретатор НКА на Java

В пакете Automata рассматривается реализация интерпретатора конечных автоматов. Если интерпретация детерминированного конечного автомата достаточно очевидна, то интерпретация недетерминированного конечного автомата требует некоторых пояснений. Недетерминизм в определении НКА означает следующее: из данной конфигурации НКА могут быть порождены несколько новых. Таким образом, мы имеем дерево возможных конфигураций

НКА на данном входе. Поскольку в общем случае НКА может иметь ϵ -переходы, которые могут образовывать циклы на графе автомата, то дерево конфигураций на данном входе, вообще говоря, бесконечно. Автомат допускает, если среди этих конфигураций есть допускающая. Таким образом, возникает вопрос о сведении потенциально бесконечного множества конфигураций к некоторому конечному, которое нужно просмотреть, чтобы убедиться, допускает ли НКА строку или нет. С другой стороны, ясно, что количество различных конфигураций НКА на данном входе длины n не превосходит $n * |Q|$, где Q — множество состояний.

Основные структуры данных в программе — это две функции. Одна из них — `HashMap TransitionFunction` — отображает состояния в отображения `HashMap SymbolToState` — для каждого символа в данном состоянии, для которого определен переход; значением этой функции является множество состояний, в которые возможен переход по данному символу (включая ϵ). Основной цикл составляет просмотр списка конфигураций `LinkedList configurations`. Вначале этот список состоит из одной начальной конфигурации. Затем выбираем первый элемент этого списка, удаляем его из списка, определяем конфигурации, в которые можно попасть из данной по текущему символу, и включаем их в список. Для каждой конфигурации определяется число — длина непросмотренного входа, который надо прочитать из данного состояния. Если мы попадаем повторно в состояние и длина непросмотренного входа не изменилась, то конфигурацию не надо включать в список непросмотренных конфигураций. Если попадаем в заключительную конфигурацию, то останавливаемся.

3.4. Алгоритмы построения конечных автоматов

3.4.1. Построение недетерминированного конечного автомата по регулярному выражению. Рассмотрим алгоритм построения по регулярному выражению недетерминированного конечного автомата, допускающего тот же язык.

Алгоритм 3.1. Построение недетерминированного конечного автомата по регулярному выражению.

Вход. Регулярное выражение r в алфавите T .

Выход. НКА M , такой, что $L(M) = L(r)$.

Метод. Автомат для выражения строится композицией из автоматов, соответствующих подвыражениям. На каждом шаге построения строящийся автомат имеет в точности одно заключительное состояние, в начальное состояние нет переходов из других состояний, и нет переходов из заключительного состояния в другие.

1. Для выражения \emptyset строится автомат согласно рис. 3.5.



Рис. 3.5

2. Для выражения e строится автомат согласно рис. 3.6.

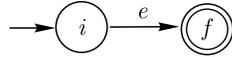


Рис. 3.6

3. Для выражения a ($a \in T$) строится автомат согласно рис. 3.7.

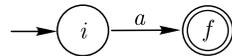


Рис. 3.7

4. Пусть $M(s)$ и $M(t)$ — НКА для регулярных выражений s и t соответственно.

- а) Для выражения $s|t$ автомат $M(s|t)$ строится, как показано на рис. 3.8. Здесь i — новое начальное состояние и f — новое заключительное состояние. Заметим, что имеет место переход по e из i в начальные состояния $M(s)$ и $M(t)$ и переход по e из заключительных состояний $M(s)$ и $M(t)$ в f . Начальное и заключительное состояния автоматов $M(s)$ и $M(t)$ не являются таковыми для автомата $M(s|t)$.

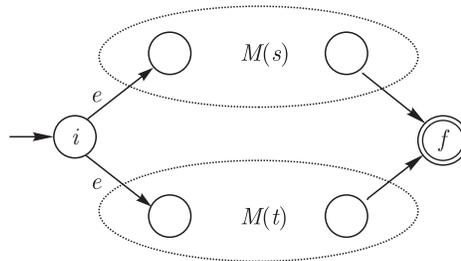


Рис. 3.8

- б) Для выражения st автомат $M(st)$ строится следующим образом (рис. 3.9).

Начальное состояние автомата $M(s)$ становится начальным для нового автомата, а заключительное состояние $M(t)$ становится заключительным для нового автомата. Начальное состояние $M(t)$ и заключительное состояние $M(s)$ сливаются, т.е. все переходы

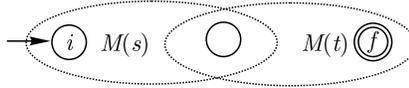


Рис. 3.9

из начального состояния $M(t)$ становятся переходами из заключительного состояния $M(s)$. В новом автомате это объединенное состояние не является ни начальным, ни заключительным.

- в) Для выражения s^* автомат $M(s^*)$ строится следующим образом (рис. 3.10):

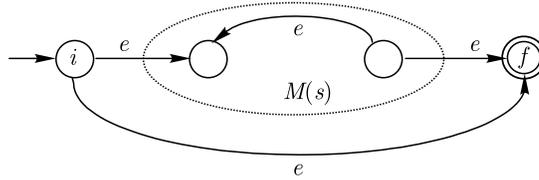


Рис. 3.10

Здесь i — новое начальное состояние, а f — новое заключительное состояние.

3.4.2. Построение детерминированного конечного автомата по недетерминированному. Рассмотрим алгоритм построения по недетерминированному конечному автомату детерминированного конечного автомата, допускающего тот же язык.

В алгоритме будут использоваться следующие функции:

$e\text{-closure}(R)$ ($R \subseteq Q$) — множество состояний НКА, достижимых из состояний, входящих в R , посредством только переходов по e , т. е. множество

$$S = \bigcup_{q \in R} \{p \mid (q, e) \vdash^* (p, e)\};$$

$move(R, a)$ ($R \subseteq Q$) — множество состояний НКА, в которые есть переход на входе a для состояний из R , т. е. множество

$$S = \bigcup_{q \in R} \{p \mid p \in D(q, a)\}.$$

Функция $e\text{-closure}(R)$ вычисляется следующим простым алгоритмом.

1. Внести все состояния из R в список *list*.
2. while (*list* не пуст) {
 - выбрать первый элемент списка r и удалить его из списка;
 - for (каждого состояния u с дугой от r в u , помеченной e)
 - do if ($u \notin e\text{-closure}(R)$) {
 - добавить u к $e\text{-closure}(R)$;
 - поместить u в список *list*; }

Алгоритм 3.2. Построение детерминированного конечного автомата по недетерминированному.

Вход. НКА $M = (Q, T, D, q_0, F)$.

Выход. ДКА $M' = (Q', T, D', q'_0, F')$, такой, что $L(M) = L(M')$.

Метод. Каждое состояние результирующего ДКА — это некоторое множество состояний исходного НКА.

Вначале Q' и D' пусты. Выполнить шаги 1–4:

1. Определить $q'_0 = e\text{-closure}(\{q_0\})$.
2. Добавить q'_0 в Q' как непомеченное состояние.
3. Выполнить следующую процедуру:


```
while (в  $Q'$  есть непомеченное состояние  $R$ ) {
  пометить  $R$ ;
  for (каждого входного символа  $a \in T$ ) {
     $S = e\text{-closure}(\text{move}(R, a))$ ;
    if ( $S \neq \emptyset$ ) {
      if ( $S \notin Q'$ )
        добавить  $S$  в  $Q'$  как непомеченное состояние;
      определить  $D'(R, a) = S$ ;
    }
  }
}
```
4. Определить $F' = \{S \mid S \in Q', S \cap F \neq \emptyset\}$.

Пример 3.6. Результат применения алгоритма 3.2 приведен на рис. 3.11.

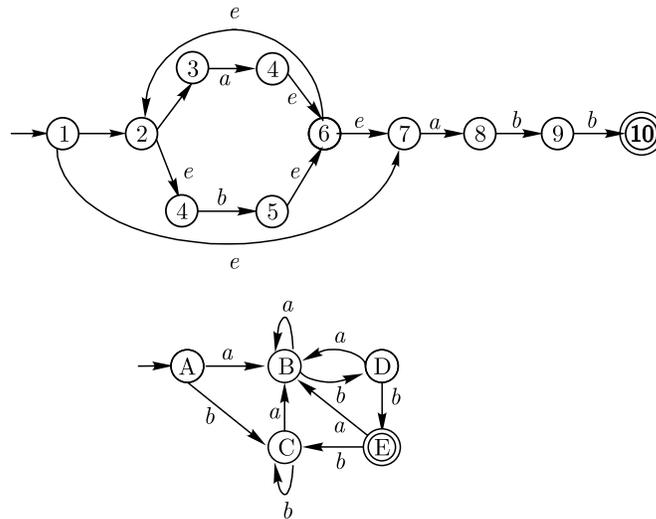


Рис. 3.11

3.4.3. Реализация на Java. Алгоритм преобразования недетерминированного конечного автомата в детерминированный реализован в пакете Automata функцией NFAtoDFA(). Две основные структуры данных — это функция, отображающая имя состояния (номер) в подмножество состояний исходного НКА, образующих данное состояние ДКА, и список непомяченных состояний. Последнее состояние выбирается из списка непомяченных состояний. Когда в результате работы функции move формируется новое подмножество, оно проверяется на совпадение с каким-либо из имеющихся. Если такое подмножество уже есть, переход делается в него. Если нет — формируется новое состояние, оно попадает в список состояний и в список непомяченных состояний.

3.4.4. Обоснование. Определим *расширенную функцию переходов* для ДКА рекурсивно следующим образом:

для ДКА

$$D_{DFA}^e(q, a) = D(q, a), \quad a \in T,$$

$$D_{DFA}^e(q, wa) = D_{DFA}(D_{DFA}^e(q, w), a), \quad w \in T^+;$$

для НКА

$$D_{NFA}^e(q, w), \quad w \in T^*,$$

— это множество состояний, в которых может оказаться автомат по прочтении цепочки w ;

$$D_{NFA}^e(q, e) = e\text{-closure}(\{q\}).$$

Если $w = ua$, $D_{NFA}^e(q, u) = \{p_1, p_2, \dots, p_k\}$, то

$$D_{NFA}^e(q, w) = e\text{-closure}\left(\bigcup_i D_{NFA}(p_i, a)\right) = e\text{-closure}(\text{move}(p_1, p_2, \dots, p_k, a)),$$

$p_i \in D_{NFA}^e(q, u)$.

Теорема 3.1. $D_{DFA}^e(q_0^e, w) = D_{NFA}^e(q_0, w)$, $q_0^e = e\text{-closure}(\{q_0\})$.

Доказательство. Индукция по длине слова.

Базис. Если $|w| = 0$, то $w = e$, $D_{NFA}^e(q_0, e) = e\text{-closure}(\{q_0\}) = q_0^e = D_{DFA}^e(q_0^e, e)$.

Шаг индукции. Пусть $w = ua$. По предположению индукции

$$D_{DFA}^e(q_0^e, u) = D_{NFA}^e(q_0, u) = \{p_1, p_2, \dots, p_k\} = r.$$

Тогда по построению ДКА и определению D_{DFA}^e имеем

$$\begin{aligned} D_{DFA}^e(q_0^e, ua) &= D_{DFA}(D_{DFA}^e(q, u), a) = D_{DFA}(r, a) = D_{DFA}(\{p_1, p_2, \dots, p_k\}, a) = \\ &= e\text{-closure}(\text{move}(\{p_1, p_2, \dots, p_k\}, a)) = \\ &= e\text{-closure}\left(\bigcup_i D_{NFA}(p_i, a)\right), p_i \in \{p_1, p_2, \dots, p_k\}. \end{aligned}$$

Для НКА имеем по определению D_{NFA}^e

$$D_{NFA}^e(q, ua) = e\text{-closure}\left(\bigcup_i D_{NFA}(p_i, a)\right), p_i \in \{p_1, p_2, \dots, p_k\},$$

т. е. имеет место равенство. Таким образом, прочитав w , НКА окажется в p_i .

Если $p_i \in F_{NFA}$, $p_i \in p^e = \{p_1, p_2, \dots, p_k\}$ -состояние ДКА, то $p^e \in F_{DFA}$ по построению. ■

3.4.5. Построение детерминированного конечного автомата по регулярному выражению. Приведем теперь алгоритм построения по регулярному выражению детерминированного конечного автомата, допускающего тот же язык [3].

Пусть дано регулярное выражение r в алфавите T . Для однозначности заключим это выражение в скобки и добавим маркер конца: $(r)\#$. Такое регулярное выражение будем называть *пополненным*. В процессе своей работы алгоритм будет использовать пополненное регулярное выражение.

Алгоритм будет оперировать с синтаксическим деревом для пополненного регулярного выражения $(r)\#$, каждый лист которого помечен символом $a \in T \cup \{e, \#\}$, а каждая внутренняя вершина помечена знаком одной из операций: \cdot (конкатенация), $|$ (объединение), $*$ (итерация).

Каждому листу дерева (кроме e -листьев) присвоим уникальный номер, называемый *позицией*, и будем использовать его, с одной стороны, для ссылки на лист в дереве и, с другой стороны, для ссылки на символ, соответствующий этому листу. Заметим, что если некоторый символ используется в регулярном выражении несколько раз, то он имеет несколько позиций.

Обойдем дерево T снизу-вверх слева-направо и вычислим четыре функции: *nullable*, *firstpos*, *lastpos* и *followpos*. Три первые функции — *nullable*, *firstpos* и *lastpos* — определены на узлах дерева, а *followpos* — на множестве позиций. Значением всех функций, кроме *nullable*, является множество позиций. Функция *followpos* вычисляется через остальные три функции.

Функция *firstpos*(n) для каждого узла n синтаксического дерева регулярного выражения дает множество позиций, которые соответствуют первым символам в подцепочках, генерируемых подвыражением с вершиной в n . Аналогично, *lastpos*(n) дает множество позиций, которым соответствуют последние символы в подцепочках, генерируемых подвыражениями с вершиной n . Для узла n , поддеревья которого (т. е. деревья, у которых узел n является корнем) могут породить пустое слово, определим *nullable*(n) = *true*, а для остальных узлов *nullable*(n) = *false*.

Выражения для вычисления функций *nullable*, *firstpos* и *lastpos* приведены в табл. 3.1.

Таблица 3.1

узел n	$nullable(n)$	$firstpos(n)$	$lastpos(n)$
лист e	<i>true</i>	\emptyset	\emptyset
лист i (не e)	<i>false</i>	$\{i\}$	$\{i\}$
$ $	$nullable(u)$	$firstpos(u)$	$lastpos(u)$
\wedge	or	\cup	\cup
$u v$	$nullable(v)$	$firstpos(v)$	$lastpos(v)$
\cdot	$nullable(u)$	if $nullable(u)$ then	if $nullable(v)$ then
\wedge	and	$firstpos(u)$	$lastpos(u)$
$u v$	$nullable(v)$	\cup	\cup
		$firstpos(v)$	$lastpos(v)$
		else $firstpos(u)$	else $lastpos(v)$
$*$	<i>true</i>	$firstpos(v)$	$lastpos(v)$
$ $			
v			

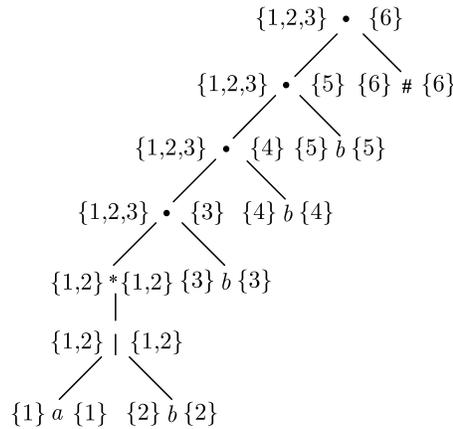


Рис. 3.12

Пример 3.7. На рис. 3.12 приведено синтаксическое дерево для пополненного регулярного выражения $(a|b)^*abb\#$ с результатом вычисления функций $firstpos$ и $lastpos$. Слева от каждого узла расположено значение $firstpos$, справа от узла — значение $lastpos$. Заметим, что эти функции могут быть вычислены за один обход дерева.

Если i — позиция, то $followpos(i)$ есть множество позиций j , таких, что существует некоторая строка $\dots cd\dots$, входящая в язык, описываемый

регулярным выражением, причем такая, что позиция i соответствует этому вхождению c , а позиция j — вхождению d .

Функция *followpos* может быть вычислена также за один обход дерева снизу-вверх следующим двум правилам.

1. Пусть n — внутренний узел с операцией \cdot (конкатенация), u и v — его потомки. Тогда для каждой позиции i , входящей в *lastpos*(u), добавляем к множеству значений *followpos*(i) множество *firstpos*(v).
2. Пусть n — внутренний узел с операцией $*$ (итерация), u — его потомок. Тогда для каждой позиции i , входящей в *lastpos*(u), добавляем к множеству значений *followpos*(i) множество *firstpos*(u).

Пример 3.8. Результат вычисления функции *followpos* для регулярного выражения из предыдущего примера приведен в табл. 3.2.

Таблица 3.2

позиция	<i>followpos</i>
1	{1, 2, 3}
2	{1, 2, 3}
3	{4}
4	{5}
5	{6}
6	\emptyset

Алгоритм 3.3. Прямое построение ДКА по регулярному выражению.

Вход. Регулярное выражение r в алфавите T .

Выход. ДКА $M = (Q, T, D, q_0, F)$, такой, что $L(M) = L(r)$.

Метод. Состояния ДКА соответствуют множествам позиций.

Вначале Q и D пусты. Выполнить шаги 1–6:

1. Построить синтаксическое дерево для пополненного регулярного выражения $(r)\#$.
2. Обходя синтаксическое дерево, вычислить значения функций *nullable*, *firstpos*, *lastpos* и *followpos*.
3. Определить $q_0 = \text{firstpos}(\text{root})$, где *root* — корень синтаксического дерева.
4. Добавить q_0 в Q как непомеченное состояние.
5. Выполнить следующую процедуру:


```
while (в  $Q$  есть непомеченное состояние  $R$ ) {
  пометить  $R$ ;
  for (каждого входного символа  $a \in T$ ,
    такого, что в  $R$  имеется позиция,
    которой соответствует  $a$ ) {
```

```

    пусть символ  $a$  в  $R$  соответствует позициям
     $p_1, \dots, p_n$ , и пусть  $S = \bigcup_{1 \leq i \leq n} followpos(p_i)$ ;
    if ( $S \neq \emptyset$ ) {
        if ( $S \notin Q$ )
            добавить  $S$  в  $Q$  как непомеченное
            состояние;
        определить  $D(R, a) = S$ ;
    }
}
}

```

6. Определить F как множество всех состояний из Q , содержащих позиции, связанные с символом $\#$.

Пример 3.9. Результат применения алгоритма 3.3 к регулярному выражению $(a|b)^*abb$ приведен на рис. 3.13.

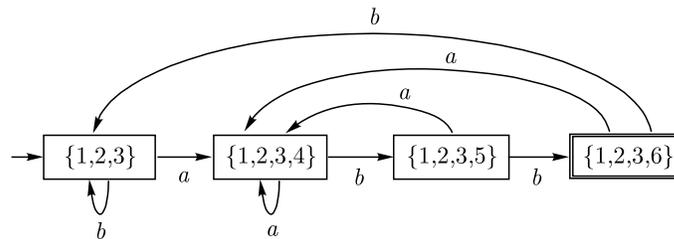


Рис. 3.13

Лемма 3.2. Построенный автомат не зависит от порядка расстановки скобок ассоциативных операций.

Доказательство. Нужно рассмотреть непосредственные построения для двух ассоциативных операций — конкатенации и объединения. ■

3.4.6. Построение ДКА по РВ на Java. Регулярное выражение записывается в соответствии с синтаксисом, задаваемым следующими правилами:

```

Expr ::= Union '#'
Union ::= Factor ('|' Factor)
Factor ::= Term Term*
Term ::= Atom ['*']
Atom ::= String | '(' Expr ')' | '$'

```

В соответствии с этой грамматикой в пакете *regularExpression* написан синтаксический анализатор методом рекурсивного спуска (см. раздел 4.3). Одновременно с синтаксическим анализом осуществляется вычисление атрибутов *nullable*, *firstPos*, *lastPos* и функции *followPos*. Затем по функции *followPos* практически так же, как и при построении ДКА по НКА, строится ДКА.

3.4.7. Обоснование. Пусть задано регулярное выражение PB со своим деревом, и пусть ДКА — детерминированный конечный автомат, построенный по PB в соответствии с алгоритмом 3.3. Введем обозначение $p \approx a$, если позиции p соответствует символ a . Каждой позиции соответствует единственный символ, но обратное несправедливо: один и тот же символ может соответствовать разным позициям.

Лемма 3.3. Если $p_1, p_2, \dots, p_n, p_{n+1}$ — последовательность позиций, такая, что:

- 1) $p_1 \in \text{firstpos}(\text{root})$;
- 2) $p_{i+1} \in \text{followpos}(p_i)$, $1 \leq i \leq n$;
- 3) $p_{n+1} \approx \#$,

то $a_1 a_2 \dots a_n \in L(PB)$, $p_i \approx a_i$.

И наоборот, если $a_1 a_2 \dots a_n \in L(PB)$, то существует $p_1, p_2, \dots, p_n, p_{n+1}$ — последовательность позиций, такая, что выполняются перечисленные выше условия.

Доказательство. Достаточность. Условие $p_1 \in \text{firstpos}(\text{root})$ означает, что по определению строка может начинаться с символа a_1 , $p_1 \approx a_1$. Если $p_{i+1} \in \text{followpos}(p_i)$, $1 \leq i \leq n$, то это означает, что пара символов $a_i a_{i+1}$ может встретиться в какой-либо строке, принадлежащей языку. Запись $p_{n+1} \approx \#$ означает, что символ a_n может быть заключительным в строке.

Необходимость. Обратно, пусть $a_1 a_2 \dots a_n \in L(PB)$. Тогда по построению $\exists p_1 \approx a_1$, $p_1 \in \text{firstpos}(\text{root})$. Поскольку в строке, принадлежащей языку, a_{i+1} следует за a_i , то это значит по определению, что среди позиций, соответствующих каждому символу, мы можем выбрать такие, что $p_{i+1} \approx a_{i+1}$, $p_{i+1} \in \text{followpos}(p_i)$. Поскольку $\#$ добавили для строки, принадлежащей языку, $p_{n+1} \in \text{followpos}(p_n)$.

Пусть $\{p_{i+1}^j\}$ — множество таких p_{i+1}^j , что $p_{i+1}^j \approx a_{i+1}$, а $\{p_i^k\}$ — множество таких p_i^k , что $p_i^k \approx a_i$. Поскольку символ a_{i+1} следует в строке за символом a_i , а строка принадлежит языку, каждый элемент $p_{i+1}^j \approx a_{i+1}$ есть образ некоторого $p_i^k \in \{p_i^k\}$, т.е. $p_{i+1}^j \in \text{followpos}(p_i^k)$. Идя таким образом от последнего символа $\#$ к первому a_1 , мы построим последовательность позиций, удовлетворяющих условию. ■

Теорема 3.2. $w \in L(PB) \text{ тогда и только тогда, когда } w \in L(\text{ДКА})$. (Здесь и далее *туткк* — сокращение словосочетания «тогда и только тогда, когда».)

Доказательство. Достаточность. Пусть $w \in L(PB)$, $w = a_1 a_2 \dots a_n$. Существует последовательность позиций $p_1, p_2, \dots, p_n, p_{n+1}$, порождающих w (в смысле $p_i \approx a_i$), такая, что $p_1 \in \text{firstpos}(\text{root})$, $p_{i+1} \in \text{followpos}(p_i)$, $1 \leq i \leq n$, $p_{n+1} \approx \#$. Покажем, что тогда существует последовательность состояний ДКА $\bar{q}_0, \bar{q}_1, \dots, \bar{q}_n$, такая, что $p_i \in \bar{q}_{i-1}$, $\bar{q}_i = D(\bar{q}_{i-1}, a_i)$, $1 \leq i \leq n$. Докажем это индукцией по i .

Базис. Возьмем $\overline{q_0} = \text{firstpos}(\text{root}) = q_0$. Тогда, поскольку $w \in L(\text{PB})$, $p_1 \in \text{firstpos}(\text{root}) = \overline{q_0}$.

Шаг. Пусть после прочтения $a_1 a_2 \dots a_i$ автомат оказался в состоянии $\overline{q_i}$. По предположению индукции $p_i \in \overline{q_{i-1}}$.

По построению ДКА

$$D(\overline{q_{i-1}}, a_i) = \bigcup_{p_i \in \overline{q_{i-1}}, p_i \approx a_i} \text{followpos}(p_i) = \overline{q_i} \neq \emptyset,$$

поскольку $p_i \approx a_i$, $p_{i+1} \in \text{followpos}(p_i)$, $\overline{q_i}$ — состояние ДКА по построению и $p_{i+1} \in \overline{q_i}$.

Если $i = n$, то $p_{n+1} \in \text{followpos}(p_n)$ и по построению $p_{n+1} \in \overline{q_n} \in F$, поскольку $p_{n+1} \approx \#$, так что цепочка допускается ДКА.

Необходимость. Пусть $w \in L(\text{ДКА})$. Пусть $\overline{q_0}, \overline{q_1}, \dots, \overline{q_n}$ — последовательность состояний ДКА, которые проходит автомат, допуская w . Имеем: $\overline{q_i} = D(\overline{q_{i-1}}, a_i)$, $1 \leq i \leq n$. $\overline{q_i} = \{p_i^j\}$, $\overline{q_{i-1}} = \{p_{i-1}^l\}$. Поскольку переход из $\overline{q_{i-1}}$ в $\overline{q_i}$ определен, то $\forall p_i^j \exists p_{i-1}^l$, такое, что $p_i^j \in \text{followpos}(p_{i-1}^l)$ и $p_i^j \approx a_i$. Значит, можно выбрать последовательность позиций, удовлетворяющих лемме. Следовательно, $a_1 a_2 \dots a_n \in L(\text{PB})$, $p_{n+1} \in \overline{q_n}$, поскольку $\overline{q_n} \in F$; $p_{n+1} \approx \#$. ■

3.5. Связь регулярных множеств, конечных автоматов и регулярных грамматик

В подразделе 3.3.3 приведен алгоритм построения детерминированного конечного автомата по регулярному выражению. Рассмотрим теперь, как по описанию конечного автомата построить регулярное множество, совпадающее с языком, допускаемым конечным автоматом.

Теорема 3.3. *Язык, допускаемый детерминированным конечным автоматом, является регулярным множеством.*

Доказательство. Пусть L — язык, допускаемый детерминированным конечным автоматом

$$M = (\{q_1, \dots, q_n\}, T, D, q_1, F).$$

Обозначим посредством R_{ij}^k множество всех слов x , таких, что $D^e(q_i, x) = q_j$ и если $D^e(q_i, y) = q_s$ для любой цепочки y — префикса x , отличного от x и e , то $s \leq k$.

Иными словами, R_{ij}^k есть множество всех слов, которые переводят конечный автомат из состояния q_i в состояние q_j , не проходя ни через какое состояние q_s для $s > k$. Однако i и j могут быть больше k .

R_{ij}^k может быть определено рекурсивно следующим образом:

$$R_{ij}^0 = \{a \mid a \in T, D(q_i, a) = q_j\},$$

$$R_{ij}^k = R_{ij}^{k-1} \cup R_{ik}^{k-1} (R_{kk}^{k-1})^* R_{kj}^{k-1}, \text{ где } 1 \leq k \leq n.$$

Таким образом, определение R_{ij}^k означает, что для входной цепочки w , переводящей M из q_i в q_j без перехода через состояния с номерами, большими k , справедливо ровно одно из следующих двух утверждений.

- 1) Цепочка w принадлежит R_{ij}^{k-1} , т.е. при анализе цепочки w автомат никогда не достигает состояний с номерами, большими или равными k .
- 2) Цепочка w может быть представлена как $w = w_1 w_2 w_3$, где $w_1 \in R_{ik}^{k-1}$ (подцепочка w_1 переводит M сначала в q_k), $w_2 \in (R_{kk}^{k-1})^*$ (подцепочка w_2 переводит M из q_k обратно в q_k , не проходя через состояния с номерами, большими или равными k), $w_3 \in R_{kj}^{k-1}$ (подцепочка w_3 переводит M из состояния q_k в q_j) — рис. 3.14.

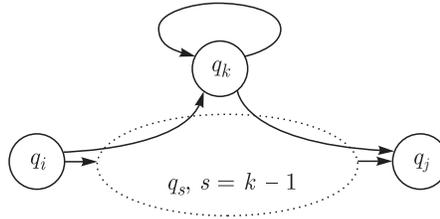


Рис. 3.14

Тогда $L = \bigcup_{q_j \in F} R_{1j}^n$. Индукцией по k можно показать, что это множество регулярно. ■

Таким образом, для всякого регулярного множества имеется конечный автомат, допускающий в точности это регулярное множество, и наоборот — язык, допускаемый конечным автоматом, есть регулярное множество.

Рассмотрим теперь соотношение между теми языками, которые порождаются праволинейными грамматиками, и теми, которые допускаются конечными автоматами.

Праволинейная грамматика $G = (N, T, P, S)$ называется *регулярной*, если:

- 1) каждое ее правило, кроме $S \rightarrow e$, имеет вид либо $A \rightarrow aB$, либо $A \rightarrow a$, где $A, B \in N$, $a \in T$;
- 2) в том случае, когда $S \rightarrow e \in P$, начальный символ S не встречается в правых частях правил.

Лемма 3.4. Пусть G — праволинейная грамматика. Тогда существует регулярная грамматика G' такая, что $L(G) = L(G')$.

Доказательство предоставляется читателю в качестве упражнения.

Теорема 3.4. Пусть $G = (N, T, P, S)$ — праволинейная грамматика. Тогда существует недетерминированный конечный автомат $M = (Q, T, D, q_0, F)$ для которого $L(M) = L(G)$.

Доказательство. На основании леммы 3.4 без ограничения общности можно считать, что G — регулярная грамматика.

Построим НКА M следующим образом:

- 1) состояниями M будут нетерминалы G плюс новое состояние R , не принадлежащее N , так что $Q = N \cup \{R\}$;
- 2) в качестве начального состояния M примем S , т. е. $q_0 = S$;
- 3) если P содержит правило $S \rightarrow e$, то $F = \{S, R\}$, иначе $F = \{R\}$ (напомним, что S не встречается в правых частях правил, если $S \rightarrow e \in P$);
- 4) состояние $R \in D(A, a)$, если $A \rightarrow a \in P$; кроме того, $D(A, a)$ содержит все B , такие, что $A \rightarrow aB \in P$, причем $D(R, a) = \emptyset$ для каждого $a \in T$.

Читая вход w , M моделирует вывод w в грамматике G . Покажем, что $L(M) = L(G)$. Пусть $w = a_1a_2 \dots a_n \in L(G)$, $n \geq 1$. Тогда $S \Rightarrow a_1A_1 \Rightarrow \dots \Rightarrow a_1a_2 \dots a_{n-1}A_{n-1} \Rightarrow a_1a_2 \dots a_{n-1}a_n$ для некоторой последовательности нетерминалов A_1, A_2, \dots, A_{n-1} . По определению, $D(S, a_1)$ содержит A_1 , $D(A_1, a_2)$ содержит A_2 , и т. д., $D(A_{n-1}, a_n)$ содержит R . Значит, $w \in L(M)$, поскольку $D^e(S, w)$ содержит R , а $R \in F$. Если $e \in L(G)$, то $S \in F$, так что $e \in L(M)$.

Аналогично, если $w = a_1a_2 \dots a_n \in L(M)$, $n \geq 1$, то существует последовательность состояний $S, A_1, A_2, \dots, A_{n-1}, R$, такая, что $D(S, a_1)$ содержит A_1 , $D(A_1, a_2)$ содержит A_2 , и т. д. Поэтому $S \Rightarrow a_1A_1 \Rightarrow a_1a_2A_2 \Rightarrow \dots \Rightarrow a_1a_2 \dots a_{n-1}A_{n-1} \Rightarrow a_1a_2 \dots a_{n-1}a_n$ — вывод в G и $w \in L(G)$. Если $e \in L(M)$, то $S \in F$, так что $S \rightarrow e \in P$ и $e \in L(G)$. ■

Теорема 3.5. Для каждого недетерминированного конечного автомата $M = (Q, T, D, q_0, F)$ существует праволинейная грамматика $G = (N, T, P, S)$, такая, что $L(G) = L(M)$.

Доказательство. Без потери общности можно считать, что автомат M — детерминированный. Определим грамматику G следующим образом:

- 1) нетерминалами грамматики G будут состояния автомата M , так что $N = Q$;
- 2) в качестве начального символа грамматики G примем q_0 , т. е. $S = q_0$;
- 3) $A \rightarrow aB \in P$, если $D(A, a) = B$;
- 4) $A \rightarrow a \in P$, если $D(A, a) = B$ и $B \in F$;
- 5) $S \rightarrow e \in P$, если $q_0 \in F$.

Доказательство того, что $S \Rightarrow^* w$ тогда и только тогда, когда $D^e(q_0, w) \in F$, аналогично доказательству теоремы 3.2. ■

В некоторых случаях для определения того, является ли язык регулярным, может быть полезным необходимое условие, которое называется леммой Огдена о разрастании.

Теорема 3.6. (Лемма о разрастании для регулярных множеств.) *Пусть L — регулярное множество. Тогда существует такая константа k , что если $w \in L$ и $|w| \geq k$, то цепочку w можно представить в виде xyz , где $0 < |y|$, $|xy| \leq k$ и $xy^i z \in L$ для всех $i \geq 0$.*

Доказательство. Пусть $M = (Q, \Sigma, D, q_0, F)$ — детерминированный конечный автомат, допускающий L , т. е. $L(M) = L$ и $k = |Q|$. Пусть $w \in L$ и $|w| \geq k$. Рассмотрим последовательность конфигураций, которые проходит автомат M , допуская цепочку w . Так как в ней по крайней мере $k + 1$ конфигураций, то среди первых $k + 1$ конфигураций найдутся две различные с одинаковыми состояниями. Таким образом, получаем существование такой последовательности тактов, что

$$(q_0, xyz) \vdash^* (q_1, yz) \vdash^r (q_1, z) \vdash^* (q_2, e)$$

для некоторых $q_1 \in Q$, $q_2 \in F$ и $0 < r \leq k$. Поскольку число конфигураций от начальной до (q_1, z) включительно не превосходит $k + 1$, а автомат детерминированный, $|xy| \leq k$. Для любого $i > 0$ автомат может проделать последовательность тактов $(q_0, xy^i z) \vdash^* (q_1, y^i z) \vdash^+ (q_1, y^{i-1} z) \dots \vdash^+ (q_1, yz) \vdash^+ (q_1, z) \vdash^* (q_2, e)$.

Таким образом, $xy^i z \in L$ для всех $i \geq 1$. Случай $i = 0$, т. е. $xy \in L$, также очевиден. ■

С помощью леммы о разрастании можно показать, что не является регулярным множеством язык $L = \{0^n 1^n \mid n \geq 1\}$.

Допустим, что L регулярен. Тогда для достаточно большого n слово $0^n 1^n$ можно представить в виде xyz , причем $y \neq e$ и $xy^i z \in L$ для всех $i \geq 0$. Если $y \in 0^+$ или $y \in 1^+$, то $xz = xy^0 z \notin L$. Если $y \in 0^+ 1^+$, то $xuyz \notin L$. Получили противоречие. Следовательно, L не может быть регулярным множеством.

3.5.1. Построение детерминированного конечного автомата с минимальным числом состояний. Рассмотрим теперь алгоритм построения ДКА с минимальным числом состояний, эквивалентного данному ДКА [2].

Пусть $M = (Q, T, D, q_0, F)$ — ДКА. Будем называть M *всюду определенным*, если $D(q, a) \neq \emptyset$ для всех $q \in Q$ и $a \in T$.

Лемма 3.5. *Пусть $M = (Q, T, D, q_0, F)$ — ДКА, не являющийся всюду определенным. Тогда существует всюду определенный ДКА M' , такой, что $L(M) = L(M')$.*

Доказательство. Рассмотрим автомат

$$M' = (Q \cup \{q'\}, T, D', q_0, F),$$

где $q' \notin Q$ — некоторое новое состояние, а функция D' определяется следующим образом.

1. Для всех $q \in Q$ и таких $a \in T$, что $D(q, a) \neq \emptyset$, определить $D'(q, a) = D(q, a)$.
2. Для всех $q \in Q$ и таких $a \in T$, что $D(q, a) = \emptyset$, определить $D'(q, a) = q'$.
3. Для всех $a \in T$ определить $D'(q', a) = q'$.

Легко показать, что автомат M' допускает тот же язык, что и M . ■

Приведенный ниже алгоритм получает на входе всюду определенный автомат. Если автомат не является всюду определенным, то его можно сделать таковым на основании только леммы 3.5.

Алгоритм 3.4. Построение ДКА с минимальным числом состояний.

Вход. Всяду определенный ДКА $M = (Q, T, D, q_0, F)$.

Выход. ДКА $M' = (Q', T, D', q'_0, F')$, такой, что $L(M) = L(M')$ и M' содержит наименьшее возможное число состояний.

Метод. Выполнить шаги 1-5:

1. Построить начальное разбиение Π множества состояний из двух групп: заключительные состояния Q и остальные $Q - F$, т. е. $\Pi = \{F, Q - F\}$.
2. Применить к Π следующую процедуру и получить новое разбиение Π_{new} :

```
for (каждой группы  $G$  в  $\Pi$ ) {
  разбить  $G$  на подгруппы так, чтобы
  состояния  $s$  и  $t$  из  $G$  оказались
  в одной подгруппе тогда и только тогда,
  когда для каждого входного символа  $a$ 
  состояния  $s$  и  $t$  имеют переходы по  $a$ 
  в состояния из одной и той же группы в  $\Pi$ ;
```

```
  заменить  $G$  в  $\Pi_{\text{new}}$  на множество всех
  полученных подгрупп;
```

```
}
```

3. Если $\Pi_{\text{new}} = \Pi$, то полагаем $\Pi_{\text{res}} = \Pi$ и переходим к шагу 4, иначе повторяем шаг 2 с $\Pi := \Pi_{\text{new}}$.

4. Пусть $\Pi_{\text{res}} = \{G_1, \dots, G_n\}$. Определим:

$Q' = \{G_1, \dots, G_n\}$;

$q'_0 = G$, где группа $G \in Q'$ такова, что $q_0 \in G$;

$F' = \{G \mid G \in Q' \text{ и } G \cap F \neq \emptyset\}$;

$D'(p', a) = q'$, если $D(p, a) = q$, где $p \in p'$ и $q \in q'$.

Таким образом, каждая группа в Π_{res} становится состоянием нового автомата M' . Если группа содержит начальное состояние автомата M , то эта группа становится начальным состоянием автомата M' . Если группа содержит заключительное состояние M , то она становится

заключительным состоянием M' . Отметим, что каждая группа Π_{res} либо состоит только из состояний из F , либо не имеет состояний из F . Переходы определяются очевидным образом.

- 5) Если M' имеет «мертвое» состояние, т. е. состояние, которое не является допускающим и из которого нет путей в допускающие, то удалить его и связанные с ним переходы из M' . Удалить из M' также все состояния, не достижимые из начального.

Пример 3.10. Результат применения алгоритма 3.4 приведен на рис. 3.15.

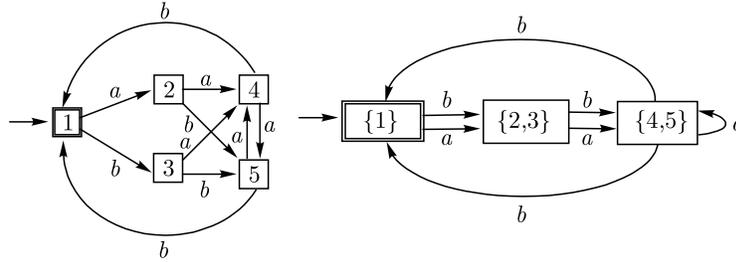


Рис. 3.15

Пусть $D^e(p, w)$ — расширенная функция переходов, которая определяется рекурсивно:

$$D^e(p, a) = D(p, a), \quad D^e(p, wa) = D(D^e(p, w), a).$$

Будем говорить, что состояния p и q эквивалентны, если для всех входных цепочек w состояние $D^e(p, w)$ является допускающим тогда и только тогда, когда состояние $D^e(q, w)$ — допускающее. Состояния $D^e(p, w)$ и $D^e(q, w)$ могут и не совпадать — лишь бы оба они были либо допускающими, либо недопускающими.

Если два состояния p и q не эквивалентны друг другу, то будем говорить, что они различимы, или неэквивалентны, т. е. существует хотя бы одна цепочка w , для которой одно из состояний $D^e(p, w)$ и $D^e(q, w)$ является допускающим, а другое — нет.

Для того чтобы найти эквивалентные состояния, нужно выявить все пары различных состояний. Все пары различных состояний можно найти представленным ниже алгоритмом. Те пары состояний, которые найти не удастся, будут эквивалентными. Алгоритм, который называется *алгоритмом заполнения таблицы*, состоит в рекурсивном обнаружении пар различных состояний ДКА (Q, Σ, D, q_0, F) .

Базис. Если состояние p — допускающее, а q — не допускающее, то пара состояний (p, q) различима.

Индукция. Пусть p и q — состояния, для которых существует входной символ a , приводящий их в различные состояния $r = D(p, a)$ и $s = D(q, a)$.

Тогда (p, q) — пара различных состояний. Это очевидно, потому что если r и s различимы, то должна существовать цепочка w , отличающая r от s . Тогда цепочка w отличает p от q , так как $D^e(r, aw)$ и $D^e(s, aw)$ — это та же пара состояний, что и $D^e(r, w)$ и $D^e(s, w)$.

Пример 3.11. Рассмотрим в качестве примера автомат, изображенный на рис. 3.16.

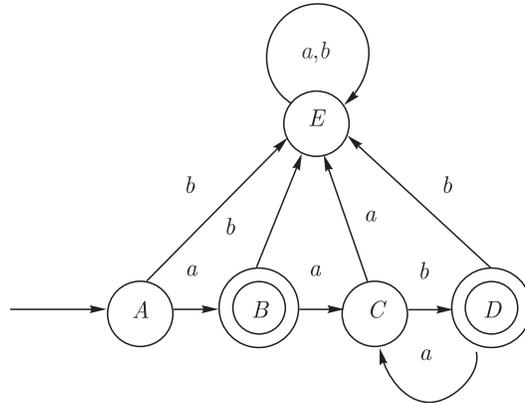


Рис. 3.16

Построим матрицу эквивалентности состояний.

Базис. Все пары незаключительных состояний не эквивалентны паре (B, D) . В пару (D, E) по символу b ведут, соответственно, множества $\{C\}$ и $\{A, B, D, E\}$, что делает неэквивалентными пары (C, E) и (C, A) . В пару (B, E) по символу a ведут, соответственно, множества $\{A\}$ и $\{E, C\}$, что делает неэквивалентными пары (A, E) и (C, A) . Результат представлен в табл. 3.3.

Таблица 3.3

D	x			
C	x	x		
B	x	0	x	
A	x	x	x	x
	E	D	C	B

Таким образом, состояния D и B эквивалентны.

Теорема 3.7. Если два состояния не различаются с помощью алгоритма заполнения таблицы, то они эквивалентны.

Доказательство. Снова рассмотрим ДКА $A = (Q, \Sigma, D, q_0, F)$. Предположим, что утверждение теоремы неверно, т.е. существует хотя бы одна пара состояний (p, q) , для которой выполняются следующие условия:

- 1) состояния p и q различимы, т.е. существует некоторая цепочка w , для которой только одно из состояний $D^e(p, w)$ и $D^e(q, w)$ является допускающим;
- 2) алгоритм заполнения таблицы не может обнаружить, что состояния r и s различимы.

Назовем такую пару состояний *плохой парой*.

Если существуют плохие пары, то среди них должны быть такие, которые различимы с помощью кратчайших из всех цепочек, различающих плохие пары. Пусть плохая пара (p, q) такова, что для нее $w = a_1 a_2 \dots a_n$ — кратчайшая из всех цепочек, различающих плохие пары. Тогда только одно из состояний $D^e(p, w)$ и $D^e(q, w)$ является допускающим.

Заметим, что цепочка w не может быть e , так как если некоторая пара состояний различается с помощью e , то ее можно обнаружить, выполнив базисную часть алгоритма заполнения таблицы. Следовательно, $n \geq 1$.

Рассмотрим состояния $r = D(p, a_1)$ и $s = D(q, a_1)$. Эти состояния можно различить с помощью цепочки $a_2 a_3 \dots a_n$, поскольку она переводит r и s в состояния $D^e(p, w)$ и $D^e(q, w)$. Однако тогда цепочка, отличающая r от s , короче любой цепочки, различающей плохую пару. Следовательно, (r, s) не может быть плохой парой, и алгоритм заполнения таблицы должен был обнаружить, что эти состояния различимы.

Но индуктивная часть алгоритма заполнения таблицы не остановится, пока не придет к выводу, что состояния p и q также различимы, поскольку уже обнаружено, что состояние $D(p, a_1) = r$ отличается от $D(q, a_1) = s$. Получено противоречие с предположением о том, что существуют плохие пары состояний. Но если плохих пар нет, то любую пару различимых состояний можно обнаружить с помощью алгоритма заполнения таблицы, и теорема доказана. ■

3.5.2. Проверка эквивалентности регулярных языков. Эквивалентность регулярных языков легко проверяется с помощью алгоритма заполнения таблицы. Предположим, что языки L и M представлены, например, соответственно регулярным выражением и некоторым НКА. Преобразуем каждое из этих представлений в ДКА. Теперь представим себе ДКА, множество состояний которого равно объединению множеств состояний автоматов для языков L и M . Технически этот ДКА содержит два начальных состояния, но фактически при проверке эквивалентности начальное состояние не играет никакой роли, поэтому любое из этих двух состояний можно принять за единственное начальное.

Далее проверяем эквивалентность начальных состояний двух заданных ДКА, используя алгоритм заполнения таблицы. Если они эквивалентны, то $L = M$, а если нет, то $L \neq M$.

Пример 3.12. Рассмотрим регулярное выражение $a(ab)^*$. Построим по нему детерминированный конечный автомат двумя способами:

- 1) сначала построив промежуточный недетерминированный конечный автомат;
 - 2) построив алгоритмом прямого построения ДКА по регулярному выражению.
- В первом случае получим автомат, изображенный на рис. 3.17, во втором — изображенный на рис. 3.16 (в обоих случаях автоматы пополнены мертвым состоянием).

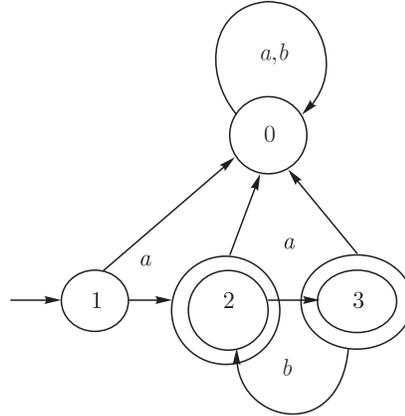


Рис. 3.17

Можно считать, что на рис. 3.16 и 3.17 изображен один ДКА, содержащий девять состояний от A до E и от 0 до 3 . Если применить алгоритм заполнения таблицы к этому автомату, то в результате получим таблицу различимости (табл. 3.4).

Таблица 3.4

3	x	x	0	x	x
2	x	0	x	0	x
1	0	x	x	x	x
0	x	x	x	x	0
	A	B	C	D	E

Поскольку в результате проверки установлено, что состояния A и 1 эквивалентны и являются начальными у двух заданных автоматов, делаем вывод, что эти ДКА действительно допускают один и тот же язык.

Время заполнения таблицы, а значит, и время проверки эквивалентности двух состояний, полиномиально относительно числа состояний. Если число состояний равно n , то количество пар состояний равно C_n^2 , или $n(n - 1)/2$. За один цикл рассматриваются все пары состояний, чтобы определить, является ли одна из пар состояний-преемников различимой. Значит, один цикл занимает время не больше $O(n^2)$. Кроме того, если в некотором цикле

не обнаружены новые пары различных состояний, то алгоритм заканчивается. Следовательно, количество циклов не превышает $O(n^2)$, а верхняя граница времени заполнения таблицы равна $O(n^4)$.

Однако с помощью более аккуратно построенного алгоритма можно заполнить таблицу за время $O(n^2)$. С этой целью для каждой пары состояний (r, s) необходимо составить список пар состояний (p, q) , «зависящих» от (r, s) , т. е. если пара (r, s) различима, то (p, q) также различима. Вначале такие списки создаются путем рассмотрения каждой пары состояний (p, q) , и для каждого входного символа a (а их число фиксировано) пара (p, q) вносится в список для пары состояний-преемников $(D(p, a)$ и $D(q, a)$).

Если обнаруживается, что пара (r, s) различима, то в списке этой пары каждая пока неразличимая пара отмечается как различимая и помещается в очередь пар, списки которых нужно проверить аналогичным образом.

Общее время работы этого алгоритма пропорционально сумме длин списков, так как каждый раз проверяется наличие некоторой пары в списке (когда проходим по списку пары, признанной различимой). Так как размер входного алфавита считается постоянным, то каждая пара состояний попадает в $O(1)$ списков. Поскольку всего пар $O(n^2)$, суммарное время также $O(n^2)$.

Еще одним важным следствием проверки эквивалентности состояний является возможность минимизации ДКА. Это значит, что для каждого ДКА можно найти эквивалентный ему ДКА с наименьшим числом состояний. Более того, для данного языка существует единственный минимальный (с точностью до обозначения состояний). Основная идея минимизации ДКА состоит в том, что понятие эквивалентности состояний позволяет объединять состояния в блоки следующим образом.

1. Все эквивалентные состояния и только они образуют блок.
2. Затем строим ДКА с состояниями — блоками.
3. Переходы определяем естественным образом.
4. Начальным состоянием является блок, содержащий начальное состояние исходного автомата.
5. Заключительным состоянием является блок, содержащий заключительные состояния исходного автомата.

Доказательство основано на том, что эквивалентность состояний *транзитивна*, а кроме того, по определению симметрична и рефлексивна; значит, блоки состояний образуют *разбиение*.

Теорема 3.8. *ДКА, полученный в результате применения алгоритма минимизации, имеет наименьшее возможное число состояний из всех ДКА, эквивалентных данному.*

Доказательство. Пусть N — ДКА, эквивалентный исходному M , но имеющий меньшее число состояний. Применим алгоритм определения

эквивалентных неразличимых состояний к объединению M и N . Их начальные состояния эквивалентны (неразличимы). Для любых двух неразличимых состояний p и q все их (непосредственные) преемники по любому символу неразличимы, поскольку если бы преемники были различимы, то p и q можно было бы различить. Таким образом, каждое состояние p автомата M неотлично хотя бы от одного состояния q автомата N . В самом деле, существует цепочка a_1, \dots, a_k , переводящая M из начального состояния в p . Эта же цепочка переводит N в некоторое q . Это следует из того, что начальные состояния неразличимы, и из предыдущего замечания. Если $|M| > |N|$, то в M имеются два состояния, неразличимые с некоторым q из M . Но это противоречит построению M , в котором все состояния различимы. ■

Следствие. *Между состояниями любых двух минимальных для данного ДКА автоматов можно установить взаимно однозначное соответствие, т. е. эти автоматы эквивалентны с точностью до именованя состояний.*

Доказательство. Достаточно применить предыдущее рассуждение в обе стороны. ■

Таким образом, получаем следующий алгоритм минимизации автомата.

1. Создать все пары (заключительное состояние, незаклучительное состояние) и отметить их как непомяенные.
2. Пока есть непомяенная пара:
 - а) помяенить ее;
 - б) для каждой пары состояний, ведущей в данную пару по каждому (одному и тому же) символу: если эта пара не обозначена как различимая, сделать ее различимой и непомяенной;
 - в) построить множества эквивалентных состояний как транзитивное замыкание, каждое такое множество объявить состоянием;
 - г) множества, состоящие из заключительных состояний, сделать заключительными.

Если применить этот алгоритм к автомату, изображенному на рис. 3.16, то получим автомат, изображенный на рис. 3.17, поскольку состояния B и D эквивалентны. Если применить этот алгоритм одновременно к автоматам, изображенным на рис. 3.16 и 3.17, то получим, что эквивалентными являются множества состояний $\{1, A\}$, $\{3, C\}$, $\{2, B, D\}$ и $\{0, E\}$.

3.5.3. Реализация на Java. Программа на Java минимизации ДКА с проверкой эквивалентности двух ДКА приведена в пакете Finite (программное приложение). Функция `equivalence()` делает следующее.

1. Сделать автомат всюду определенным.
2. Для каждого состояния определить, из какого состояния по каждому символу есть переход.
3. Для каждой пары состояний (p, q) и каждого входного символа a сформировать множество таких пар состояний (r, s) , что $D(r, a) = p$, $D(s, a) = q$.
4. Сформировать список пар различных состояний:
 - а) если есть пара $(p, q) \in (Q \setminus F \times F)$, то занести ее в список непомеченных пар;
 - б) пока список непомеченных пар непуст, выбрать пару (p, q) , пометить ее как различимую и удалить из списка непомеченных пар; если пара (r, s) такова, что $D(r, a) = p$, $D(s, a) = q$, ее нет в списке различных пар и нет в списке непомеченных пар, то внести ее в список непомеченных пар.
5. Построить множество $L = \{(Q \times Q) \setminus \text{список различных пар}\}$:
6. Построить разбиение множества L на классы эквивалентности:
 - а) выбрать пару (p, q) из этого множества;
 - б) построить множество $S_i = \{p, q\}$;
 - в) пока в S_i есть такой элемент s , что во множестве L есть пара (r, s) или (s, r) , добавить r к этому множеству и исключить пару (r, s) (или (s, r)) из L .
7. Каждое множество S_i становится новым состоянием; все состояния из $Q \setminus \{S_i\}$ остаются одиночными состояниями.
8. Множества S_i , состоящие из заключительных состояний, становятся заключительными состояниями; все заключительные состояния, не вошедшие в $\{S_i\}$, остаются заключительными состояниями.

3.6. Программирование лексического анализа

Как уже отмечалось ранее, лексический анализатор (ЛА) может быть оформлен как подпрограмма. При обращении к ЛА вырабатываются как минимум два результата: тип выбранной лексемы и ее значение (если оно есть).

Если ЛА сам формирует таблицы объектов, то он выдает тип лексемы и указатель на соответствующий вход в таблице объектов. Если же ЛА не работает с таблицами объектов, то он выдает тип лексемы, а ее значение передается, например, через некоторую глобальную переменную. Помимо значения лексемы, эта глобальная переменная может содержать некоторую

дополнительную информацию: номер текущей строки, номер символа в строке и т. п. Эта информация может использоваться в различных целях, например, для диагностики.

В основе ЛА лежит диаграмма переходов соответствующего конечного автомата. Отдельная проблема здесь — анализ ключевых слов. Как правило, ключевые слова — это выделенные идентификаторы. Поэтому возможны два основных способа распознавания ключевых слов: либо очередная лексема сначала диагностируется на совпадение с каким-либо ключевым словом и в случае неуспеха делается попытка выделить лексему из какого-либо класса, либо, наоборот, после выборки лексемы идентификатора происходит обращение к таблице ключевых слов на предмет сравнения. Подробнее о механизмах поиска в таблицах будет сказано ниже (гл. 7), здесь отметим только, что поиск ключевых слов может вестись либо в основной таблице имен (и в этом случае в нее до начала работы ЛА загружаются ключевые слова), либо в отдельной таблице. При первом способе все ключевые слова непосредственно встраиваются в конечный автомат ЛА, при втором — конечный автомат содержит только разбор идентификаторов.

В некоторых языках (например, ПЛ/1 или Фортран) ключевые слова могут использоваться в качестве обычных идентификаторов. В этом случае работа ЛА не может идти независимо от работы синтаксического анализатора. В Фортране возможны, например, следующие строки:

```
DO 10 I=1,25
DO 10 I=1.25
```

В первом случае строка — это заголовок цикла DO, во втором — оператор присваивания. Поэтому, прежде чем можно будет выделить лексему, ЛА должен заглянуть довольно далеко.

Еще сложнее дело в ПЛ/1. Здесь возможны такие операторы:

```
IF ELSE THEN ELSE = THEN; ELSE THEN = ELSE;
```

или

```
DECLARE (A1, A2, ... , AN)
```

и только в зависимости от того, что стоит после «)», можно определить, является ли DECLARE ключевым словом или идентификатором. Длина такой строки может быть сколь угодно большой, и уже невозможно отделить фазу синтаксического анализа от фазы лексического анализа.

Рассмотрим несколько подробнее вопросы программирования ЛА. Основная операция ЛА, на которую уходит большая часть времени его работы — это взятие очередного символа и проверка на принадлежность его некоторому диапазону. Например, основной цикл при выборке числа в простейшем случае может выглядеть следующим образом:

```
while (Insym<='9' && Insym>='0')
{ ... }
```

Программу можно значительно улучшить следующим образом [5]. Пусть LETTER, DIGIT, BLANK — элементы перечислимого типа. Введем массив map, входами которого будут символы, значениями — типы символов. Инициализируем массив map следующим образом:

```
map['a']=LETTER;
.....
map['z']=LETTER;
map['0']=DIGIT;
.....
map['9']=DIGIT;
map[' ']=BLANK;
.....
```

Тогда приведенный цикл примет следующую форму:

```
while (map[Insym]==DIGIT)
{ ... }
```

Выделение ключевых слов может осуществляться после выделения идентификаторов. ЛА работает быстрее, если ключевые слова выделяются непосредственно.

Для этого строится конечный автомат, описывающий множество ключевых слов. На рис. 3.18 приведен фрагмент такого автомата.

Рассмотрим пример программирования этого конечного автомата на языке Си, приведенный в [22]:

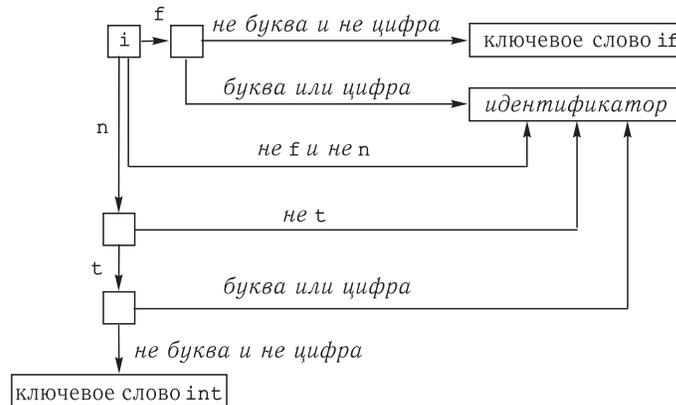


Рис. 3.18

```
.....
case 'i':
  if (cp[0]=='f' &&! (map[cp[2]] & (DIGIT | LETTER)))
  {cp++; return IF;}
  if (cp[0]=='n' && cp[1]=='t'
      &&! (map[cp] & (DIGIT | LETTER)))
```

```
{cp+=2; return INT;}
{ обработка идентификатора }
.....
```

Здесь *cp* — указатель текущего символа. В массиве *map* классы символов кодируются битами.

Поскольку ЛА анализирует каждый символ входного потока, его скорость существенно зависит от скорости выборки очередного символа входного потока. В свою очередь, эта скорость во многом определяется схемой буферизации. Рассмотрим возможные эффективные схемы буферизации.

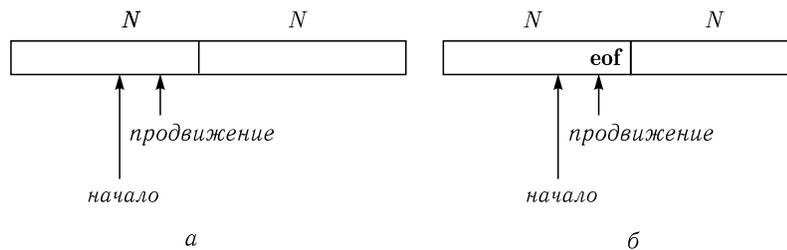


Рис. 3.19

Будем использовать буфер, состоящий из двух одинаковых частей длиной N (рис. 3.19, *а*), где N — размер блока обмена (например, 1024, 2048 и т. п.). Чтобы не читать каждый символ отдельно, в каждую из половин буфера поочередно одной командой чтения считывается N символов. Если на входе осталось меньше N символов, то в буфер помещается специальный символ *eof*. Буфер снабжен двумя указателями: *продвижение* и *начало*. Между указателями размещается текущая лексема. Вначале они оба указывают на первый символ выделяемой лексемы. Один из них, *продвижение*, продвигается вперед, пока не будет выделена лексема, и устанавливается на ее конец. После обработки лексемы оба указателя устанавливаются на символ, следующий за лексемой. Если указатель *продвижение* переходит середину буфера, то правая половина заполняется новыми N символами. Если указатель *продвижение* переходит правую границу буфера, то левая половина заполняется N символами и указатель *продвижение* устанавливается на начало буфера.

При каждом продвижении указателя необходимо проверять, не достигли ли мы границы одной из половин буфера. Для всех символов, кроме лежащих в концах половин буфера, требуются две проверки. Число проверок можно свести к одной, если в конце каждой половины поместить дополнительный «сторожевой» символ, в качестве которого логично взять *eof* (рис. 3.19, *б*).

В этом случае почти для всех символов делается единственная проверка на совпадение с *eof*, и только при совпадении нужно дополнительно проверить, достигли ли мы середины или правого конца.

3.7. Конструктор лексических анализаторов LEX

Для автоматизации разработки ЛА было создано довольно много средств. Как правило, входными языками для них служат или праволинейные грамматики, или языки регулярных выражений. Одной из наиболее распространенных систем является LEX, работающий с расширенными регулярными выражениями. LEX-программа состоит из трех частей:

```
Объявления
%%
Правила трансляции
%%
Вспомогательные подпрограммы
```

Секция объявлений включает объявления переменных, констант и определения регулярных выражений. При определении регулярных выражений могут использоваться следующие конструкции:

[abc]	— либо a , либо b , либо c ;
[a-z]	— диапазон символов;
R^*	— 0 или более повторений регулярного выражения R ;
R^+	— 1 или более повторений регулярного выражения R ;
R_1/R_2	— R_1 , если за ним следует R_2 ;
$R_1 R_2$	— либо R_1 , либо R_2 ;
$R?$	— если есть R , выбрать его;
$R\$$	— выбрать R , если оно последнее в строке;
R	— выбрать R , если оно первое в строке;
[R]	— дополнение к R ;
$R\{n,m\}$	— повторение R от n до m раз;
{имя}	— именованное регулярное выражение;
(R)	— группировка.

Правила трансляции LEX-программ имеют следующий вид:

```
p_1 { действие_1 }
p_2 { действие_2 }
.....
p_n { действие_n }
```

где p_i — регулярное выражение, а действие $_i$ — фрагмент программы, описывающий, какое действие должен сделать ЛА, когда образец p_i сопоставляется лексеме. В LEX действия записываются на Си.

Третья секция содержит вспомогательные процедуры, необходимые для действий. Эти процедуры могут транслироваться отдельно и загружаться с ЛА.

ЛА, сгенерированный LEX, взаимодействует с синтаксическим анализатором следующим образом. При вызове ЛА синтаксическим анализатором он посимвольно читает остаток входа, пока не находит самый длинный префикс, который может быть сопоставлен одному из регулярных выражений `p_i`. Затем он выполняет `действие_i`. Как правило, `действие_i` возвращает управление синтаксическому анализатору. Если это не так, т.е. в соответствующем действии нет возврата, то ЛА продолжает поиск лексем до тех пор, пока действие не вернет управление синтаксическому анализатору. Повторный поиск лексем вплоть до явной передачи управления позволяет ЛА правильно обрабатывать пробелы и комментарии. Синтаксическому анализатору ЛА возвращает единственное значение — тип лексемы. Для передачи информации о типе лексемы используется глобальная переменная `yylval`. Текстовое представление выделенной лексемы хранится в переменной `ytext`, а ее длина в переменной `ylen`.

Пример 3.13. LEX-программа для ЛА, обрабатывающего идентификаторы, числа, ключевые слова `if`, `then`, `else` и знаки логических операций:

```
%{ /*определения констант LT, LE, EQ, NE, GT,
    GE, IF, THEN, ELSE, ID, NUMBER, RELOP, например,
    через DEFINE или скалярный тип*/ %}
/*регулярные определения*/
delim  [ \t\n]
ws     {delim}+
letter [A-Za-z]
digit  [0-9]
id     {letter}({letter}|{digit})*
number {digit}+(\.{digit}+)?(E[+\-]?{digit}+)?
%%
{ws}   /* действий и возврата нет */
if     {return(IF);}
then   {return(THEN);}
else   {return(ELSE);}
{id}   {yylval=install_id(); return(ID);}
{number} {yylval=install_num(); return(NUMBER);}
"<"   {yylval=LT; return(RELOP);}
"<="  {yylval=LE; return(RELOP);}
"="    {yylval=EQ; return(RELOP);}
"<>"  {yylval=NE; return(RELOP);}
">"   {yylval=GT; return(RELOP);}
">="  {yylval=GE; return(RELOP);}
%%
install_id()
/*подпрограмма, которая помещает лексему,
   на первый символ которой указывает ytext,
```

```

        длина которой равна yulen, в таблицу
        символов и возвращает указатель на нее*/
    }
    install_num()
    /*аналогичная подпрограмма для размещения
        лексемы числа*/
    }

```

В разделе объявлений, заключенном в скобки `%{` и `%}`, перечислены константы, используемые правилами трансляции. Все, что заключено в эти скобки, непосредственно копируется в программу ЛА `lex.yy.c` и не рассматривается как часть регулярных определений или правил трансляции. То же касается и вспомогательных подпрограмм третьей секции. В данном примере это подпрограммы `install_id` и `install_num`.

В секцию определений входят также некоторые регулярные определения. Каждое такое определение состоит из имени и регулярного выражения, обозначаемого этим именем. Например, первое определенное имя — это `delim`. Оно обозначает класс символов `{ \t\n\}`, т.е. любой из трех символов: пробел, табуляция или новая строка. Второе определение — разделитель, обозначаемый именем `ws`. Разделитель — это любая последовательность одного или более символов-разделителей. Слово `delim` должно быть заключено в скобки, чтобы отличить его от образца, состоящего из пяти символов `delim`.

В определении `letter` используется класс символов. Сокращение `[A-Za-z]` обозначает любую из прописных букв от `A` до `Z` или строчных букв от `a` до `z`. В пятом определении `id` для группировки используются скобки, являющиеся метасимволами LEX. Аналогично, вертикальная черта — метасимвол LEX, обозначающий объединение.

В последнем регулярном определении `number` символ `«+»` используется как метасимвол «одно или более вхождений», символ `«?»` как метасимвол «ноль или одно вхождение». Обратная черта используется для того, чтобы придать обычный смысл символу, используемому в LEX как метасимвол. В частности, десятичная точка в определении `number` обозначается как `«.»`, поскольку точка сама по себе представляет класс, состоящий из всех символов, за исключением символа новой строки. В классе символов `[+\]` обратная черта перед минусом стоит потому, что знак минус используется как символ диапазона, как в `[A-Z]`.

Если символ имеет смысл метасимвола, то придать ему обычный смысл можно и по-другому, заключив его в кавычки. Так, в секции правил трансляции шесть операций отношения заключены в кавычки.

Рассмотрим правила трансляции, следующие за первым `%%`. Согласно первому правилу, если обнаружено `ws`, т.е. максимальная последовательность пробелов, табуляций и новых строк, никаких действий не производится. В частности, не осуществляется возврат в синтаксический анализатор.

Согласно второму правилу, если обнаружена последовательность букв `if`, нужно вернуть значение `IF`, которое определено как целая константа, понимаемая синтаксическим анализатором как лексема `if`. Аналогично обрабатываются ключевые слова `then` и `else` в двух следующих правилах.

В действии, связанном с правилом для `id`, два оператора. Переменной `yylval` присваивается значение, возвращаемое процедурой `install_id`. Переменная `yylval` определена в программе `lex.yy.c`, выходе LEX и доступна синтаксическому анализатору. Она хранит возвращаемое лексическое значение, поскольку второй оператор в действии, `return(ID)`, может только вернуть код класса лексем. Функция `install_id` заносит идентификаторы в таблицу символов.

Аналогично обрабатываются числа в следующем правиле. В последних шести правилах `yylval` используется для возврата кода операции отношения, возвращаемое же функцией значение — это код лексемы `relop`.

Если, например, в текущий момент ЛА обрабатывает лексему `if`, то этой лексеме соответствуют два образца: `if` и `{id}`, причем более длинной строки, соответствующей образцу, нет. Поскольку образец `if` предшествует образцу для идентификатора, конфликт разрешается в пользу ключевого слова. Такая стратегия разрешения конфликтов позволяет легко резервировать ключевые слова.

Если на входе встречается «`<=`», то первому символу соответствует образец «`<`», но это не самый длинный образец, который соответствует префиксу входа. Стратегия выбора самого длинного префикса легко разрешает такого рода конфликты.

СИНТАКСИЧЕСКИЙ АНАЛИЗ

4.1. Контекстно-свободные грамматики и автоматы с магазинной памятью

Пусть $G = (N, T, P, S)$ — КС-грамматика. Введем несколько важных понятий и определений.

Вывод, в котором в любой сентенциальной форме на каждом шаге делается подстановка самого левого нетерминала, называется *левосторонним*. Если $S \Rightarrow^* u$ в процессе левостороннего вывода, то u — *левая сентенциальная форма*. Аналогично определим *правосторонний* вывод. Обозначим шаги левого (правого) вывода \Rightarrow_l (\Rightarrow_r).

Упорядоченным графом называется пара (V, E) , где V есть множество вершин, а E — множество линейно упорядоченных списков дуг, каждый элемент которого имеет вид $((v, v_1), (v, v_2), \dots, (v, v_n))$. Этот элемент указывает, что из вершины v выходят n дуг, причем первой из них считается дуга, входящая в вершину v_1 , второй — дуга, входящая в вершину v_2 , и т. д.

Упорядоченным помеченным деревом называется упорядоченный граф (V, E) , основой которого является дерево и для которого определена функция $f: V \rightarrow F$ (функция разметки), где F — некоторое множество.

Упорядоченное помеченное дерево D называется *деревом вывода* (или *деревом разбора*) цепочки w в КС-грамматике $G = (N, T, P, S)$, если выполняются следующие условия:

- 1) корень дерева D помечен S ;
- 2) каждый лист помечен либо $a \in T$, либо ϵ ;
- 3) каждая внутренняя вершина помечена нетерминалом $A \in N$;
- 4) если X — нетерминал, которым помечена внутренняя вершина и X_1, \dots, X_n — метки ее прямых потомков в указанном порядке, то $X \rightarrow X_1 \dots X_n$ — правило из множества P ;
- 5) цепочка, составленная из выписанных слева направо меток листьев, равна w .

Процесс определения принадлежности данной строки языку, порожаемо-му данной грамматикой, и, в случае указанной принадлежности, построения дерева разбора для этой строки называется *синтаксическим анализом*. Можно говорить о восстановлении дерева вывода (в частности, правостороннего или левостороннего) для строки, принадлежащей языку. По восстановленному выводу можно строить дерево разбора.

Грамматика G называется *неоднозначной*, если существует цепочка w , для которой имеется два или более различных деревьев вывода в G .

Грамматика G называется *леворекурсивной*, если в ней имеется нетерминал A , такой, что для некоторой цепочки α существует вывод $A \Rightarrow {}^+A\alpha$.

Автомат с магазинной памятью (МП-автомат) состоит из *управляющего устройства*, *входной ленты (входа)* и так называемого *магазина*. Один конец магазина называется *верхушкой*, другой — *дном*. Формально это семерка $M = (Q, T, \Gamma, D, q_0, Z_0, F)$, где:

- 1) Q — конечное множество *состояний* управляющего устройства;
- 2) T — конечный *входной* алфавит;
- 3) Γ — конечный алфавит *магазинных символов*;
- 4) D — отображение множества $Q \times (T \cup \{e\}) \times \Gamma$ в множество конечных подмножеств $Q \times \Gamma^*$, называемое *функцией переходов*;
- 5) $q_0 \in Q$ — *начальное состояние* управляющего устройства;
- 6) $Z_0 \in \Gamma$ — символ, находящийся в магазине в начальный момент (*начальный символ магазина*);
- 7) $F \subseteq Q$ — множество *заключительных состояний*.

Конфигурация МП-автомата — это тройка (q, w, u) , где

- 1) $q \in Q$ — текущее состояние управляющего устройства;
- 2) $w \in T^*$ — непрочитанная часть входной цепочки; первый символ цепочки w находится под входной головкой; если $w = e$, то считается, что вся входная лента прочитана;
- 3) $u \in \Gamma^*$ — содержимое магазина; самый левый символ цепочки u считается верхним символом магазина; если $u = e$, то магазин считается пустым.

Такт работы МП-автомата M будем представлять в виде бинарного отношения \vdash , определенного на конфигурациях. Будем писать

$$(q, aw, Zu) \vdash (p, w, vu),$$

если множество $D(q, a, Z)$ содержит (p, v) , где $q, p \in Q$, $a \in T \cup \{e\}$, $w \in T^*$, $Z \in \Gamma$ и $u, v \in \Gamma^*$ (верхушка магазина слева).

Начальной конфигурацией МП-автомата M называется конфигурация вида (q_0, w, Z_0) , где $w \in \Gamma^*$, т. е. управляющее устройство находится в начальном состоянии, входная лента содержит цепочку, которую нужно проанализировать, а в магазине имеется только начальный символ Z_0 .

Заключительной конфигурацией называется конфигурация вида (q, e, u) , где $q \in F$, $u \in \Gamma^*$, т. е. управляющее устройство находится в одном из заключительных состояний, а входная цепочка целиком прочитана.

Введем транзитивное и рефлексивно-транзитивное замыкания отношения \vdash , а также его степень $k \geq 0$ (обозначаемые \vdash^+ , \vdash^* и \vdash^k соответственно).

Говорят, что цепочка w *допускается* МП-автоматом M , если $(q_0, w, Z_0) \vdash^* (q, e, u)$ для некоторых $q \in F$ и $u \in \Gamma^*$.

Множество всех цепочек, допускаемых автоматом M , называется языком, *допускаемым (распознаваемым, определяемым)* автоматом M (обозначается $L(M)$).

Пример 4.1. Рассмотрим МП-автомат

$$M = (\{q_0, q_1, q_2\}, \{a, b\}, \{Z, a, b\}, D, q_0, Z, \{q_2\}),$$

у которого функция переходов D содержит элементы:

$$\begin{aligned} D(q_0, a, Z) &= \{(q_0, aZ)\}, \\ D(q_0, b, Z) &= \{(q_0, bZ)\}, \\ D(q_0, a, a) &= \{(q_0, aa), (q_1, e)\}, \\ D(q_0, a, b) &= \{(q_0, ab)\}, \\ D(q_0, b, a) &= \{(q_0, ba)\}, \\ D(q_0, b, b) &= \{(q_0, bb), (q_1, e)\}, \\ D(q_1, a, a) &= \{(q_1, e)\}, \\ D(q_1, b, b) &= \{(q_1, e)\}, \\ D(q_1, e, Z) &= \{(q_2, e)\}. \end{aligned}$$

Нетрудно показать, что $L(M) = \{ww^R \mid w \in \{a, b\}^+\}$, где w^R обозначает обращение («переворачивание») цепочки w .

Иногда допустимость определяют несколько иначе: цепочка w допускается МП-автоматом M , если $(q_0, w, Z_0) \vdash^* (q, e, e)$ для некоторого $q \in Q$. В таком случае говорят, что автомат допускает цепочку *опустошением магазина*. Эти определения эквивалентны, ибо справедлива

Теорема 4.1. *Язык допускается МП-автоматом тогда и только тогда, когда он (некоторым другим автоматом) допускается опустошением магазина.*

Доказательство. Пусть $L = L(M)$ для некоторого МП-автомата $M = (Q, T, \Gamma, D, q_0, Z_0, F)$. Построим новый МП-автомат M' , допускающий тот же язык опустошением магазина.

Пусть $M' = (Q \cup \{q'_0, q_e\}, T, \Gamma \cup \{Z'_0\}, D', q'_0, Z'_0, \emptyset)$, где функция переходов D' определена следующим образом:

- 1) если $(r, u) \in D(q, a, Z)$, то $(r, u) \in D'(q, a, Z)$ для всех $q \in Q$, $a \in T \cup \{e\}$ и $Z \in \Gamma$ (моделирование M);
- 2) $D'(q'_0, e, Z'_0) = \{(q_0, Z_0 Z'_0)\}$ (начало работы);
- 3) для всех $q \in F$ и $Z \in \Gamma \cup \{Z'_0\}$ множество $D'(q, e, Z)$ содержит (q_e, e) (переход в состояние сокращения магазина без продвижения);
- 4) $D'(q_e, e, Z) = \{(q_e, e)\}$ для всех $Z \in \Gamma \cup \{Z'_0\}$ (сокращение магазина).

Автомат сначала переходит в конфигурацию $(q_0, w, Z_0 Z'_0)$ соответственно определению D' в п. 2), затем в $(q, e, Y_1 \dots Y_k Z'_0)$, $q \in F$ соответственно п. 1), затем в $(q_e, e, Y_1 \dots Y_k Z'_0)$, $q \in F$ соответственно п. 3), затем в (q_e, e, e) соответственно п. 4). Нетрудно показать по индукции, что $(q_0, w, Z_0) \vdash^+ (q, e, u)$ (где $q \in F$) выполняется для автомата M тогда и только тогда, когда $(q'_0, w, Z'_0) \vdash^+ (q_e, e, e)$ выполняется для автомата M' . Поэтому $L(M) = L'$, где L' — язык, допускаемый автоматом M' опустошением магазина.

Обратно, пусть $M = (Q, T, \Gamma, D, q_0, Z_0, \emptyset)$ — МП-автомат, допускающий опустошением магазина язык L . Построим автомат M' , допускающий тот же язык по заключительному состоянию.

Пусть $M' = (Q \cup \{q'_0, q_f\}, T, \Gamma \cup \{Z'_0\}, D', q'_0, Z'_0, \{q_f\})$, где D' определяется следующим образом:

- 1) $D'(q'_0, e, Z'_0) = \{(q_0, Z_0 Z'_0)\}$ — переход в «режим M »;
- 2) для каждого $q \in Q$, $a \in T \cup \{e\}$, и $Z \in \Gamma$ определим $D'(q, a, Z) = D(q, a, Z)$ — работа в «режиме M »;
- 3) для всех $q \in Q$, $(q_f, e) \in D'(q, e, Z'_0)$ — переход в заключительное состояние.

Нетрудно показать по индукции, что $L = L(M')$. ■

Одним из важнейших результатов теории контекстно-свободных языков является доказательство эквивалентности МП-автоматов и КС-грамматик.

Теорема 4.2. *Язык является контекстно-свободным тогда и только тогда, когда он допускается МП-автоматом.*

Доказательство. Пусть $G = (N, T, P, S)$ — КС-грамматика. Построим МП-автомат, допускающий язык $L(G)$ опустошением магазина.

Пусть $M = (\{q\}, T, N \cup T, D, q, S, \emptyset)$, где D определяется следующим образом:

- 1) если $A \rightarrow u \in P$, то $(q, u) \in D(q, e, A)$;
- 2) $D(q, a, a) = \{(q, e)\}$ для всех $a \in T$.

Фактически этот МП-автомат в точности моделирует все возможные выводы в грамматике G . Нетрудно показать по индукции, что для любой цепочки $w \in T^*$ вывод $S \Rightarrow^+ w$ в грамматике G существует тогда и только тогда, когда существует последовательность тактов $(q, w, S) \vdash^+ (q, e, e)$ автомата M .

Наоборот, пусть дан $M = (Q, T, \Gamma, D, q_0, Z_0, \emptyset)$ — МП-автомат, допускающий опустошением магазина язык L . Построим грамматику G , порождающую язык L .

Пусть $G = (\{ [qZr] \mid q, r \in Q, Z \in \Gamma \} \cup \{S\}, T, P, S)$, где P состоит из правил следующего вида.

1. $S \rightarrow [q_0Z_0q] \in P$ для всех $q \in Q$.
2. Если $(r, e) \in D(q, a, Z)$, то $[qZr] \rightarrow a \in P, a \in T \cup \{e\}$.
3. Если $(r, X_1 \dots X_k) \in D(q, a, Z), k \geq 1$, то

$$[qZs_k] \rightarrow a[rX_1s_1][s_1X_2s_2] \dots [s_{k-1}X_k s_k] \in P$$

для любого набора s_1, s_2, \dots, s_k состояний из Q .

Нетерминалы и правила вывода грамматики определены так, что работе автомата M при обработке цепочки w соответствует левосторонний вывод w в грамматике G .

Лемма 4.1. Если $(q, x, \alpha) \vdash^* (p, y, \beta)$, то $\forall w \in T^*, \gamma \in \Gamma^*(q, xw, \alpha\gamma) \vdash^* (p, yw, \beta\gamma)$.

Доказательство основано на том, что магазинный автомат читает магазин строго сверху-вниз посимвольно и читает вход также строго слева-направо посимвольно.

Интерпретация определенных выше нетерминалов такова.

Теорема 4.3. $[qZp] \Rightarrow^* w$ *туттк* $(q, w, Z) \vdash^* (p, e, e)$.

Доказательство. Необходимость. Пусть $(q, w, Z) \vdash^* (p, e, e)$. Доказываем индукцией по числу переходов.

Базис. $(q, w, Z) \vdash (p, e, e)$, т. е. $(p, e) \in D(q, w, Z)$. По построению (правило 2) $[qZp] \rightarrow w, w \in T \cup \{e\}$.

Шаг. Пусть $(q, w, Z) \vdash^* (p, e, e)$ состоит из $n > 1$ шагов. Рассмотрим первый шаг: $(q, w, Z) \vdash (s_0, u, X_1X_2 \dots X_k) \vdash^* (p, e, e)$, $w = au, (s_0, X_1X_2 \dots X_k) \in D(q, a, Z), a \in T \cup \{e\}$. По построению $[qZs_k] \rightarrow a[s_0Xs_1][s_1Xs_2] \dots [s_{k-1}Xs_k]$ для всех $s_1, \dots, s_k \in Q$. Поскольку автомат читает цепочку u с опустошением магазина, ее можно разбить так, что $u = w_1w_2 \dots w_k$ и $\exists s_1, \dots, s_{k-1}$, такие, что $(s_{i-1}, w_i, X_i) \vdash^* (s_i, e, e)$. При этом используется менее n шагов. По предположению индукции $[s_{i-1}X_i s_i] \Rightarrow^* w_i$. Тогда $[qZs_k] \rightarrow a[s_0Xs_1][s_1Xs_2] \dots [s_{k-1}Xs_k] \Rightarrow^* aw_1[s_1Xs_2] \dots [s_{k-1}Xs_k] \Rightarrow^* aw_1 \dots w_k = w$.

Достаточность. Пусть $[qZp] \Rightarrow^* w$. Индукция по длине вывода.

Базис. Правило $[qZp] \rightarrow w \in P$, значит по правилу 2 $(q, w, Z) \vdash (p, e, e)$, $w \in T \cup e$.

Шаг. Пусть $[qZp] \Rightarrow^* w$ за $n > 1$ шагов. На первом шаге $[qZs_k] \Rightarrow a[s_0X_1s_1][s_1X_2s_2] \dots [s_{k-1}X_k s_k] \Rightarrow^* w (s_k = p)$. Разобьем цепочку $aw_1 \dots w_k$ так, что $[s_{i-1}X_i s_i] \Rightarrow^* w_i, i = 1, \dots, k$. По предположению индукции поскольку все эти выводы короче n , $(s_{i-1}w_i X_i) \vdash^* (s_i, e, e)$. Значит, по лемме $(s_{i-1}, w_i w_{i+1} \dots w_k, X_i X_{i+1} \dots X_k) \vdash^* (s_i, w_{i+1} \dots w_k, X_{i+1} \dots X_k)$. Правило $[qZs_k] \rightarrow a[s_0X_1s_1][s_1X_2s_2] \dots [s_{k-1}X_k s_k]$ соответствует переходу $(s_0, X_1, X_2, \dots, X_k) \in D(q, a, Z)$. Следовательно, $(q, aw_1 \dots w_k) \vdash (s_0, w_1 \dots w_k, X_1 \dots X_k) \vdash^* (s_1, w_2 \dots w_k, X_2 \dots X_k) \vdash^* (s_k, e, e)$.

Следствие. $S \Rightarrow^* w$ *туттк* $[q_0 Z p_0] \Rightarrow^* w$ *туттк* $(q_0, w, Z_0) \vdash^* (p, e, e)$. ■

Пример 4.2. Построение грамматики по МП-автомату.

Пусть задан автомат

$$D(q_0, 0, Z_0) = (q_0, XZ_0)$$

$$D(q_0, 0, X) = (q_0, XX)$$

$$D(q_0, 1, X) = (q_1, e)$$

$$D(q_1, 1, X) = (q_1, e)$$

$$D(q_1, e, X) = (q_1, e)$$

$$D(q_1, e, Z_0) = (q_1, e)$$

Грамматика:

$$S \rightarrow [q_0, Z_0, q_0]$$

$$S \rightarrow [q_0, Z_0, q_1]$$

Из $D(q_0, 0, Z_0) = (q_0, XZ_0)$ получаются:

$$[q_0, Z_0, q_0] \rightarrow 0[q_0, X, q_0][q_0, Z_0, q_0]$$

$$[q_0, Z_0, q_0] \rightarrow 0[q_0, X, q_1][q_1, Z_0, q_0]$$

$$[q_0, Z_0, q_1] \rightarrow 0[q_0, X, q_0][q_0, Z_0, q_1]$$

$$[q_0, Z_0, q_1] \rightarrow 0[q_0, X, q_1][q_0, Z_0, q_1]$$

Из $D(q_0, 0, X) = (q_0, XX)$ получаются:

$$[q_0, X, q_0] \rightarrow 0[q_0, X, q_0][q_0, X, q_0]$$

$$[q_0, X, q_0] \rightarrow 0[q_0, X, q_1][q_1, X, q_0]$$

$$[q_0, X, q_1] \rightarrow 0[q_0, X, q_0][q_0, X, q_1]$$

$$[q_0, X, q_1] \rightarrow 0[q_0, X, q_1][q_1, X, q_1]$$

Из $D(q_0, 1, X) = (q_1, e)$ получается $[q_0, X, q_1] \rightarrow 1$.

Из $D(q_1, e, Z_0) = (q_1, e)$ получается $[q_1, Z_0, q_1] \rightarrow e$.

Из $D(q_1, e, X) = (q_1, e)$ получается $[q_1, X, q_1] \rightarrow e$.

Из $D(q_1, 1, X) = (q_1, e)$ получается $[q_1, X, q_1] \rightarrow 1$.

Нет правил для $[q_1, X, q_0], [q_1, Z_0, q_0]$.

Нет терминальных выводов для $[q_0, Z_0, q_0], [q_0, X, q_0]$.

Остаются:

$$S \rightarrow [q_0, Z_0, q_1]$$

$$[q_0, Z_0, q_1] \rightarrow 0[q_0, X, q_1][q_1, X, q_0]$$

$$[q_0, X, q_1] \rightarrow 0[q_0, X, q_1][q_1, X, q_1]$$

$$[q_1, Z_0, q_1] \rightarrow e$$

$$[q_0, X, q_1] \rightarrow 1$$

$$[q_1, X, q_1] \rightarrow e$$

$$[q_1, X, q_1] \rightarrow 1$$

Или в другой записи:

$$S \rightarrow A$$

$$A \rightarrow 0BC$$

$$B \rightarrow 0BD$$

$$C \rightarrow e$$

$$B \rightarrow 1$$

$$D \rightarrow e$$

$$D \rightarrow 1$$

МП-автомат $M = (Q, T, \Gamma, D, q_0, Z_0, F)$ называется *детерминированным* (ДМП-автоматом), если выполняются следующие условия:

- 1) множество $D(q, a, Z)$ содержит не более одного элемента для любых $q \in Q, a \in T \cup \{e\}, Z \in \Gamma$;
- 2) если $D(q, e, Z) \neq \emptyset$, то $D(q, a, Z) = \emptyset$ для всех $a \in T$.

Допускаемый ДМП-автоматом язык называется *детерминированным* КС-языком.

Так как функция переходов ДМП-автомата содержит не более одного элемента для любой тройки аргументов, мы будем пользоваться записью $D(q, a, Z) = (p, u)$ для обозначения $D(q, a, Z) = \{(p, u)\}$.

Пример 4.3. Рассмотрим ДМП-автомат

$$M = (\{q_0, q_1, q_2\}, \{a, b, c\}, \{Z, a, b\}, D, q_0, Z, \{q_2\}),$$

функция переходов которого определяется следующим образом:

$$D(q_0, X, Y) = (q_0, XY), X \in \{a, b\}, Y \in \{Z, a, b\},$$

$$D(q_0, c, Y) = (q_1, Y), Y \in \{a, b\},$$

$$D(q_1, X, X) = (q_1, e), X \in \{a, b\},$$

$$D(q_1, e, Z) = (q_2, e).$$

Нетрудно показать, что этот детерминированный МП-автомат допускает язык $L = \{wscw^R | w \in \{a, b\}^+\}$.

К сожалению, ДМП-автоматы имеют меньшую распознавательную способность, чем МП-автоматы. Доказано, в частности, что существуют КС-языки, не являющиеся детерминированными КС-языками (таковым, например, является язык из примера 4.1).

Рассмотрим еще один важный вид МП-автомата.

Расширенным автоматом с магазинной памятью назовем семерку $M = (Q, T, \Gamma, D, q_0, Z_0, F)$, где смысл всех символов тот же, что и для обычного МП-автомата, кроме D , представляющего собой отображение конечного подмножества множества $Q \times (T \cup \{e\}) \times \Gamma^*$ во множество конечных подмножеств множества $Q \times \Gamma^*$. Все остальные определения (конфигурации, такта, допустимости) для расширенного МП-автомата остаются такими же, как для обычного.

Теорема 4.4. Пусть $M = (Q, T, \Gamma, D, q_0, Z_0, F)$ — расширенный МП-автомат. Тогда существует МП-автомат M' , такой, что $L(M') = L(M)$.

Расширенный МП-автомат $M = (Q, T, \Gamma, D, q_0, Z_0, F)$ называется *детерминированным*, если выполняются следующие условия:

- 1) множество $D(q, a, u)$ содержит не более одного элемента для любых $q \in Q$, $a \in T \cup \{e\}$, $u \in \Gamma^*$;
- 2) если $D(q, a, u) \neq \emptyset$, $D(q, a, v) \neq \emptyset$ и $u \neq v$, то не существует цепочки x , такой, что $u = vx$ или $v = ux$;
- 3) если $D(q, a, u) \neq \emptyset$, $D(q, e, v) \neq \emptyset$, то не существует цепочки x , такой, что $u = vx$ или $v = ux$.

Теорема 4.5. Пусть $M = (Q, T, \Gamma, D, q_0, Z_0, F)$ — расширенный ДМП-автомат. Тогда существует ДМП-автомат M' , такой, что $L(M') = L(M)$.

ДМП-автомат и расширенный ДМП-автомат лежат в основе рассматриваемых далее в этой главе LL- и LR-анализаторов.

Определение 4.1. Говорят, что КС-грамматика находится в *нормальной форме Хомского*, если каждое правило имеет один из следующих видов:

- 1) $A \rightarrow BC$, A, B, C — нетерминалы;
- 2) $A \rightarrow a$, a — терминал;
- 3) $S \rightarrow e$, и в этом случае S не встречается в правых частях правил.

Утверждение 4.1. Любую КС-грамматику можно преобразовать в эквивалентную ей в нормальной форме Хомского.

Определение 4.2. Назовем *высотой дерева* максимальную длину пути (число внутренних вершин) от корня до листа.

Утверждение 4.2. Если КС-грамматика находится в нормальной форме Хомского, то для любой цепочки α , если $\alpha \in L(G)$ и h — высота дерева вывода с кроной α , $|\alpha| \leq 2^{h-1}$.

Обратно, если $|\alpha| \geq 2^{h-1}$, то высота дерева больше или равна h .

Теорема 4.6. (Лемма о разрастании для контекстно-свободных языков.) Для любого КС-языка L существуют такая константа k , что любая цепочка $\alpha \in L$, $|\alpha| \geq k$, представима в виде $\alpha = uvwxu$, где:

- 1) $|vwx| \leq k$;
- 2) $vx \neq e$;
- 3) $uv^iwx^iy \in L$ для любого $i \geq 0$.

Доказательство. Пусть $L = L(G)$, где $G = (N, \Sigma, P, S)$ — контекстно-свободная грамматика в нормальной форме Хомского. Обозначим через n число нетерминалов, т. е. $n = |N|$, и рассмотрим $k = 2^n$.

Пусть $|\alpha| \geq k = 2^n$. Тогда высота дерева с кроной α больше или равна $n + 1$ и есть путь по дереву (от корня до некоторого листа), который включает не менее, чем $n + 1$ внутренних вершин (нетерминалов). Таким образом, существует хотя бы один нетерминал, который помечает не менее двух вершин этого пути. Среди всех таких нетерминалов на этом пути пусть A — такой, что его второе вхождение, считая от листа, не содержит других нетерминалов, обладающих этим свойством (если бы это было не так, то выбрали бы этот другой). Пусть q — вхождение A , ближайшее к листу, p — ближайшее, расположенное выше. Представим крону α в виде $uvwxy$, где w — крона поддерева D_1 с корнем q и vwx — крона поддерева D_2 с корнем p . Тогда высота поддерева D_2 не более $(n - 1) + 2 = n + 1$ (не более $n - 1$ нетерминалов, отличных от A , плюс два вхождения A), так что $|vwx| \leq 2^n$.

Также очевидно, что $vx \neq \epsilon$, поскольку в силу определения нормальной формы Хомского p имеет двух сыновей, помеченных нетерминалами, из которых не выводится пустая цепочка.

Кроме того, $S \Rightarrow^* uAy \Rightarrow^* uvAxy \Rightarrow^* uvwxy$, а также $A \Rightarrow^* vAx \Rightarrow^* vwx$. Отсюда получаем $A \Rightarrow^* v^iwx^i$ для всех $i \geq 0$ и $S \Rightarrow^* uv^iwx^iy$ для всех $i \geq 0$. ■

Пример 4.4. Покажем, что язык $L = \{a^n b^n c^n \mid n \geq 1\}$ не является контекстно-свободным языком.

Если бы он был КС-языком, то мы взяли бы константу k , которая определяется в лемме о разрастании. Пусть $z = a^k b^k c^k$. Тогда $z = uvwxy$. Так как $|vwx| \leq k$, то в цепочке vwx не могут быть вхождения каждого из символов a, b и c . Таким образом, цепочка uvw , которая по лемме о разрастании принадлежит L , содержит либо k символов a , либо k символов c . Но она не может иметь k вхождений каждого из символов a, b и c , потому что $|uvw| < 3k$. Значит, вхождений какого-то из этих символов в uvw больше, чем другого, и, следовательно, $uvw \notin L$. Полученное противоречие позволяет заключить, что L — не КС-язык.

4.2. Преобразования КС-грамматик

Рассмотрим ряд преобразований, позволяющих «улучшить» свойства контекстно-свободной грамматики без изменения порождаемого ею языка.

Назовем символ $X \in (N \cup T)$ *недостижимым* в КС-грамматике $G = (N, T, P, S)$, если X не появляется ни в одной выводимой цепочке этой грамматики. Иными словами, символ X недостижим, если в G не существует вывода $S \Rightarrow^* \alpha X \beta$, $\alpha, \beta \in (N \cup T)^*$.

Назовем символ $X \in (N \cup T)$ *несводимым (бесплодным)*, если в грамматике не существует вывода вида $X \Rightarrow^* w$, где $w \in T^*$.

Назовем символ *бесполезным*, если он является недостижимым или несводимым.

Бесполезные символы не могут участвовать в порождении терминальных строк языка, поэтому рассмотрим алгоритм построения грамматики, эквивалентной данной, но не содержащей бесполезных символов.

Алгоритм 4.1. Устранение несводимых символов.

Вход. КС-грамматика $G = (N, T, P, S)$.

Выход. КС-грамматика $G' = (N', T', P', S)$ без несводимых символов, такая, что $L(G') = L(G)$.

Метод. Выполнить шаги 1–4:

1. Положить $N_0 = T$ и $i = 1$;
2. Положить $N_i = \{A \mid A \rightarrow \alpha \in P \text{ и } \alpha \in (N_{i-1})^*\} \cup N_{i-1}$;
3. Если $N_i \neq N_{i-1}$, то положить $i = i + 1$ и перейти к шагу 2, в противном случае положить $N_e = N_i$ и перейти к шагу 4;
4. Положить $G_1 = ((N \cap N_e) \cup \{S\}, T, P_1, S)$, где P_1 состоит из правил множества P , содержащих только символы из $N_e \cup T$;

Алгоритм 4.2. Устранение недостижимых символов.

Вход. КС-грамматика $G = (N, T, P, S)$.

Выход. КС-грамматика $G' = (N', T', P', S)$ без недостижимых символов, такая, что $L(G') = L(G)$.

Метод. Выполнить шаги 1–4:

1. Положить $V_0 = \{S\}$ и $i = 1$;
2. Положить $V_i = \{X \mid \text{в } P \text{ есть } A \rightarrow \alpha X \beta \text{ и } A \in V_{i-1}\} \cup V_{i-1}$;
3. Если $V_i \neq V_{i-1}$, положить $i = i + 1$ и перейти к шагу 2, в противном случае перейти к шагу 4;
4. Положить $N' = V_i \cap N$, $T' = V_i \cap T$. Включить в P' все правила из P , содержащие только символы из V_i .

Чтобы устранить все бесполезные символы, необходимо применить к исходной грамматике сначала алгоритм 4.1, а затем алгоритм 4.2.

Пример 4.5. Все символы следующей грамматики

$S \rightarrow AS \mid b$;

$A \rightarrow AB$;

$B \rightarrow a$.

Поскольку все символы грамматики достижимы, применение сначала алгоритма 4.2 не меняет грамматику. Применение алгоритма 4.1 приводит к появлению недостижимых символов.

КС-грамматика без бесполезных символов называется *приведенной*. Легко видеть, что для любой КС-грамматики существует эквивалентная приведенная. В дальнейшем будем предполагать, что все рассматриваемые грамматики — приведенные.

4.3. Алгоритм Кока–Янгера–Касами

Приведем алгоритм синтаксического анализа, применимый для любой грамматики в нормальной форме Хомского.

Алгоритм 4.3 (Кока–Янгера–Касами).

Вход. КС-грамматика $G = (N, T, P, S)$ в нормальной форме Хомского и входная цепочка $w = a_1 a_2 \dots a_n \in T^+$.

Выход. Таблица разбора Tab для w , такая, что $A \in t_{ij}$ тогда и только тогда, когда $A \Rightarrow^+ a_i a_{i+1} \dots a_{i+j-1}$.

Метод. Выполнить шаги 1–3:

1. Положить $t_{i1} = \{A \mid A \rightarrow a_i \in P\}$ для каждого i . Таким образом, если $A \in t_{i1}$, то $A \Rightarrow^+ a_i$.
2. Пусть $t_{ij'}$ вычислено для $1 \leq i \leq n$ и $1 \leq j' < j$. Положим $t_{ij} = \{A \mid \text{для некоторого } 1 \leq k < j \text{ правило } A \rightarrow BC \in P, B \in t_{ik}, C \in t_{i+k, j-k}\}$. Так как $1 \leq k < j$, то $k < j$ и $j - k < j$. Поэтому t_{ik} и $t_{i+k, j-k}$ вычисляются раньше, чем t_{ij} . Если $A \in t_{ij}$, то $A \Rightarrow BC \Rightarrow^+ a_i \dots a_{i+k-1} C \Rightarrow^+ a_i \dots a_{i+k-1} a_{i+k} \dots a_{i+j-1}$.
3. Повторять шаг 2 до тех пор, пока не станут известны t_{ij} для всех $1 \leq i \leq n$ и $1 \leq j \leq n - i + 1$.

Алгоритм 4.4 (нахождения левого разбора по таблице разбора Tab).

Вход. КС-грамматика $G = (N, T, P, S)$ в нормальной форме Хомского с правилами, занумерованными от 1 до p , входная цепочка $w = a_1 a_2 \dots a_n \in T^+$ и таблица разбора Tab .

Выход. Левый разбор цепочки w или сигнал «ошибка».

Метод. Процедура $\text{gen}(i, j, A)$.

1. Если $j = 1$ и $A \rightarrow a_i = p_m$, то выдать m .
2. Пусть $j > 1$ и k — наименьшее из чисел от 1 до $j - 1$, для которых существуют $B \in t_{ik}$, $C \in t_{i+k, j-k}$ и правило $p_m = A \rightarrow BC$. Выдать m и выполнить $\text{gen}(i, k, B)$, затем $\text{gen}(i + k, j - k, C)$. Выполнить $\text{gen}(1, n, S)$, если $S \in t_{1, n}$; иначе «ошибка».

4.4. Разбор сверху-вниз (предсказывающий разбор)

4.4.1. Алгоритм разбора сверху-вниз. Пусть дана КС-грамматика $G = (N, T, P, S)$. Рассмотрим *разбор сверху-вниз (предсказывающий разбор)* для грамматики G .

Главная задача предсказывающего разбора — определение правила вывода, которое нужно применить к нетерминалу. Процесс предсказывающего разбора с точки зрения построения дерева разбора проиллюстрирован на рис. 4.1.

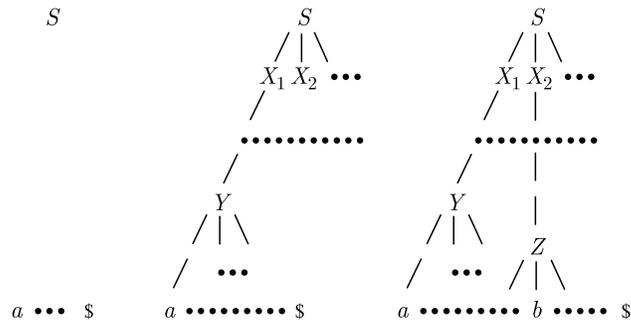


Рис. 4.1

Фрагменты недостроенного дерева соответствуют синтаксическим формам. Вначале дерево состоит только из одной вершины, соответствующей аксиоме S . В этот момент по первому символу входной цепочки предсказывающий анализатор должен определить правило $S \rightarrow X_1 X_2 \dots$, которое должно быть применено к S . Затем необходимо определить правило, которое должно быть применено к X_1 , и т. д. до тех пор, пока в процессе такого построения синтаксической формы, соответствующей левому выводу, не будет применено правило $Y \rightarrow a \dots$. Этот процесс затем применяется для следующего самого левого нетерминального символа синтаксической формы.

На рис. 4.2 условно показана структура предсказывающего анализатора, который определяет очередное правило с помощью таблицы. Такую таблицу можно построить и непосредственно по грамматике. Таблично-управляемый предсказывающий анализатор имеет входную ленту, управляющее устройство (программу), таблицу анализа, магазин и выходную ленту. Входная лента содержит анализируемую строку, заканчивающуюся символом $\$$ — маркером конца строки. Выходная лента содержит последовательность примененных правил вывода.

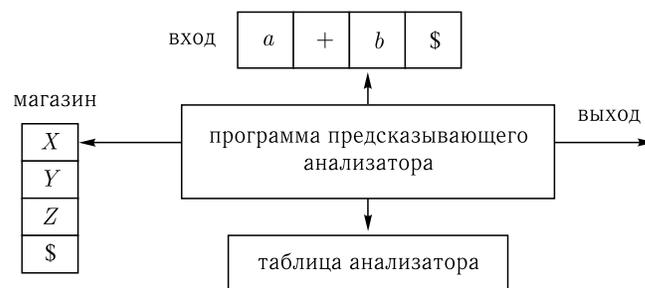


Рис. 4.2

Таблица анализатора — это двумерный массив $M[A, a]$, где A — нетерминал, a — терминал или символ $\$$. Значением $M[A, a]$ может быть некоторое правило грамматики или элемент «ошибка».

Магазин может содержать последовательность символов грамматики с $\$$ на дне. В начальный момент магазин содержит только начальный символ грамматики на вершущке и $\$$ на дне.

Анализатор работает следующим образом. Вначале он находится в конфигурации, в которой магазин содержит $S\$$, на входной ленте $w\$$ (w — анализируемая цепочка), выходная лента пуста. На каждом такте анализатор рассматривает X — символ на вершущке магазина и a — текущий входной символ. Эти два символа определяют действия анализатора. Имеются следующие возможности.

1. Если $X=a=\$,$ то анализатор останавливается, сообщает об успешном окончании разбора и выдает содержимое выходной ленты.
2. Если $X=a \neq \$,$ то анализатор удаляет X из магазина и продвигает указатель входа на следующий входной символ.
3. Если X — терминал, причем $X \neq a,$ то анализатор останавливается и сообщает о том, что входная цепочка не принадлежит языку.
4. Если X — нетерминал, то анализатор обращается к таблице $M[X, a]$. Возможны два случая.
 - а) Значением $M[X, a]$ является правило для X . В этом случае анализатор заменяет X на вершущке магазина на правую часть данного правила, а само правило помещает на выходную ленту. Указатель входа не передвигается.
 - б) Значением $M[X, a]$ является «ошибка». В этом случае анализатор останавливается и сообщает о том, что входная цепочка не принадлежит языку.

Работа анализатора может быть задана следующей программой:

```

Поместить  $\$,$  затем  $S$  в магазин;
do { $X$  = верхний символ магазина;
  if ( $X$  - терминал)
    { if ( $X==InSym$ )
      {удалить  $X$  из магазина;
        $InSym$  = очередной символ;
      }
    else {error(); break;}
    else if ( $X$  - нетерминал)
      if ( $M[X, InSym]==X \rightarrow Y_1 Y_2 \dots Y_k$ )
        {удалить  $X$  из магазина;
         поместить  $Y_k, Y_{k-1}, \dots, Y_1$  в магазин ( $Y_1$  на вершущку);
         вывести правило  $X \rightarrow Y_1 Y_2 \dots Y_k$ ;
        }
      else {error(); break;} /*вход таблицы  $M$  пуст*/

```

```

}
while (X!= $); /*магазин пуст*/
if (InSym != $) error(); /*Не вся строка прочитана*/

```

Пример 4.6. Рассмотрим грамматику арифметических выражений $G = (\{E, E', T, T', F\}, \{id, +, *, (,)\}, P, E)$ с правилами:

$$\begin{aligned}
 E &\rightarrow TE'; & T' &\rightarrow *FT'; \\
 E' &\rightarrow +TE'; & T' &\rightarrow e; \\
 E' &\rightarrow e; & F &\rightarrow (E); \\
 T &\rightarrow FT'; & F &\rightarrow id.
 \end{aligned}$$

Ниже приведена таблица предсказывающего анализатора для этой грамматики (табл. 4.1). Пустые клетки таблицы соответствуют элементу «ошибка».

Таблица 4.1

Нетерминал	Входной символ					
	id	+	*	()	\$
E	$E \rightarrow TE'$			$E \rightarrow TE'$		
E'		$E' \rightarrow +TE'$			$E' \rightarrow e$	$E' \rightarrow e$
T	$T \rightarrow FT'$			$T \rightarrow FT'$		
T'		$T' \rightarrow e$	$T' \rightarrow *FT'$		$T' \rightarrow e$	$T' \rightarrow e$
F	$F \rightarrow id$			$F \rightarrow (E)$		

При разборе входной цепочки $id + id * id\$$ анализатор совершает последовательность шагов, описанную в табл. 4.2. Заметим, что анализатор осуществляет в точности левый вывод. Если за уже просмотренными входными символами поместить символы грамматики в магазине, то можно получить в точности левые сентенциальные формы вывода. Дерево разбора для этой цепочки приведено на рис. 4.3.

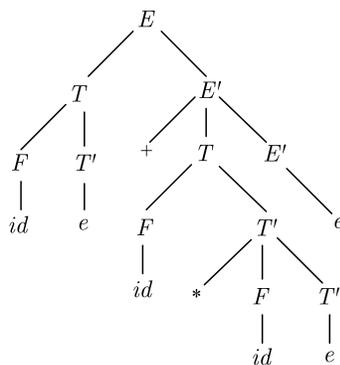


Рис. 4.3

Таблица 4.2

Магазин	Вход	Выход
$E\$,$	$id + id * id\$,$	
$TE'\$,$	$id + id * id\$,$	$E \rightarrow TE'$
$FT'E'\$,$	$id + id * id\$,$	$T \rightarrow FT'$
$id T'E'\$,$	$id + id * id\$,$	$F \rightarrow id$
$T'E'\$,$	$+id * id\$,$	
$E'\$,$	$+id * id\$,$	$T' \rightarrow e$
$+TE'\$,$	$+id * id\$,$	$E' \rightarrow +TE$
$TE'\$,$	$id * id\$,$	
$FT'E'\$,$	$id * id\$,$	$T \rightarrow FT'$
$id T'E'\$,$	$id * id\$,$	$F \rightarrow id$
$T'E'\$,$	$*id\$,$	
$*F'T'E'\$,$	$*id\$,$	$T' \rightarrow *FT'$
$FT'E'\$,$	$id\$,$	
$id T'E'\$,$	$id\$,$	$F \rightarrow id$
$T'E'\$,$	$\$,$	
$E'\$,$	$\$,$	$T' \rightarrow e$
$\$,$	$\$,$	$E' \rightarrow e$

4.4.2. Функции $FIRST$ и $FOLLOW$. При построении таблицы предсказывающего анализатора нам потребуются две функции — $FIRST$ и $FOLLOW$.

Пусть $G = (N, T, P, S)$ — КС-грамматика. Для α — произвольной цепочки, состоящей из символов грамматики, определим $FIRST(\alpha)$ как множество терминалов, с которых начинаются строки, выводимые из α . Если $\alpha \Rightarrow^* e$, то e также принадлежит $FIRST(\alpha)$.

Определим $FOLLOW(A)$ для нетерминала A как множество терминалов a , которые могут появиться непосредственно справа от A в некоторой сентенциальной форме грамматики, т. е. множество терминалов a , таких, что существует вывод вида $S \Rightarrow^* \alpha A a \beta$ для некоторых $\alpha, \beta \in (N \cup T)^*$. Заметим, что между A и a в процессе вывода могут находиться нетерминальные символы, из которых выводится e . Если A может быть самым правым символом некоторой сентенциальной формы, то $\$$ также принадлежит $FOLLOW(A)$.

Рассмотрим алгоритмы вычисления функции $FIRST$.

Алгоритм 4.5. Вычисление $FIRST$ для символов КС-грамматики.

Вход. КС-грамматика $G = (N, T, P, S)$.

Выход. Множество $FIRST(X)$ для каждого символа $X \in (N \cup T)$.

Метод. Выполнить шаги 1–3:

1. Если X — терминал, то положить $FIRST(X) = \{X\}$; если X — нетерминал, то положить $FIRST(X) = \emptyset$.
2. Если в P имеется правило $X \rightarrow e$, то добавить e к $FIRST(X)$.
3. Пока существуют множества $FIRST(X)$, к которым можно добавлять новые элементы, выполнять:

do {*continue* = false;

 Для каждого нетерминала X

 Для каждого правила $X \rightarrow Y_1Y_2\dots Y_k$

$\{i=1; nonstop = true;$

 while ($i \leq k \ \&\& \ nonstop$)

 {добавить $FIRST(Y_i) \setminus \{e\}$ к $FIRST(X)$;

 if (Были добавлены новые элементы)

continue = true;

 if ($e \notin FIRST(Y_i)$) $nonstop = false$;

 else $i+ = 1$;

 }

 if ($nonstop$) {добавить e к $FIRST(X)$;

continue = true;

 } } }

while (*continue*)

Алгоритм 4.6. Вычисление $FIRST$ для цепочки.

Вход. КС-грамматика $G = (N, T, P, S)$.

Выход. Множество $FIRST(X_1X_2 \dots X_n)$, $X_i \in (N \cup T)$.

Метод. Выполнить шаги 1–3:

1. При помощи алгоритма 4.5 вычислить $FIRST(X)$ для каждого $X \in (N \cup T)$.
2. Положить $FIRST(X_1X_2 \dots X_n) = \emptyset$.
3. $\{i = 1; nonstop = true;$
 while ($i \leq n \ \&\& \ nonstop$)
 {добавить $FIRST(X_i) \setminus \{e\}$ к $FIRST(u)$;
 if ($e \notin FIRST(X_i)$) $nonstop = false$;

```

    else  $i+ = 1$ ;
  }
  if (nonstop) {добавить  $e$  к  $FIRST(u)$ ;
  }

```

Рассмотрим алгоритм вычисления функции $FOLLOW$.

Алгоритм 4.7. Вычисление $FOLLOW$ для нетерминалов грамматики.

Вход. КС-грамматика $G = (N, T, P, S)$.

Выход. Множество $FOLLOW(X)$ для каждого символа $X \in N$.

Метод. Выполнить шаги 1–4:

1. Положить $FOLLOW(X) = \emptyset$ для каждого символа $X \in N$.
2. Добавить $\$$ к $FOLLOW(S)$.
3. Если в P есть правило вывода $A \rightarrow \alpha B \beta$, где $\alpha, \beta \in (N \cup T)^*$, то все элементы из $FIRST(\beta)$, за исключением e , добавить к $FOLLOW(B)$.
4. Пока существуют множества $FOLLOW(X)$, к которым можно добавлять новые элементы, выполнять:

если в P есть правило $A \rightarrow \alpha B$ или $A \rightarrow \alpha B \beta$, $\alpha, \beta \in (N \cup T)^*$, где $FIRST(\beta)$ содержит e ($\beta \Rightarrow^* e$), то все элементы из $FOLLOW(A)$ добавить к $FOLLOW(B)$.

Пример 4.7. Рассмотрим грамматику из примера 4.3. Для нее:

$FIRST(E) = FIRST(T) = FIRST(F) = \{ (, id \}$;

$FIRST(E') = \{ +, e \}$;

$FIRST(T') = \{ *, e \}$;

$FOLLOW(E) = FOLLOW(E') = \{), \$ \}$;

$FOLLOW(T) = FOLLOW(T') = \{ +,), \$ \}$;

$FOLLOW(F) = \{ +, *,), \$ \}$.

Например, id и левая скобка добавляются к $FIRST(F)$ на шаге 3 при $i = 1$, поскольку $FIRST(id) = \{ id \}$ и $FIRST(() = \{ (\}$ в соответствии с шагом 1. На шаге 3 при $i = 1$, в соответствии с правилом вывода $T \rightarrow FT'$, к $FIRST(T)$ добавляются также id и левая скобка. На шаге 2 в $FIRST(E')$ включается e .

При вычислении множеств $FOLLOW$ на шаге 2 в $FOLLOW(E)$ включается $\$$. На шаге 3, на основании правила $F \rightarrow (E)$, к $FOLLOW(E)$ добавляется также правая скобка. На шаге 4, примененном к правилу $E \rightarrow TE'$, в $FOLLOW(E')$ включаются $\$$ и правая скобка. Поскольку $E' \Rightarrow^* e$, они также попадают и во множество $FOLLOW(T)$. В соответствии с правилом вывода $E \rightarrow TE'$ на шаге 3 в $FOLLOW(T)$ включаются и все элементы из $FIRST(E')$, отличные от e .

Мы определили $FIRST$ как множество цепочек длины не более 1. Точно так же можно определить функцию $FIRST_k(\alpha)$, где k — натуральное число

и $\alpha \in (N \cup \Sigma)^*$: $FIRST_k(\alpha) = \{w \in \Sigma^* \mid \text{либо } |w| < k \text{ и } \alpha \Rightarrow_G w, \text{ либо } |w| = k \text{ и } \alpha \Rightarrow_G wx \text{ для некоторого } x \in \Sigma^*\}$.

Если $\alpha \in \Sigma^*$, то $FIRST_k(\alpha) = \{w\}$, где w — это первые k символов цепочки α при $|\alpha| \geq k$ и $w = \alpha$ при $|\alpha| < k$.

4.4.3. Конструирование таблицы предсказывающего анализатора.

Для конструирования таблицы предсказывающего анализатора по грамматике G может быть использован алгоритм, основанный на следующей идее. Предположим, что $A \rightarrow \alpha$ — правило вывода грамматики и $a \in FIRST(\alpha)$. Тогда анализатор делает развертку A по α , если входным символом является a . Трудность возникает, когда $\alpha = \epsilon$ или $\alpha \Rightarrow^* \epsilon$. В этом случае нужно развернуть A в α , если текущий входной символ принадлежит $FOLLOW(A)$ или если достигнут $\$$ и $\$ \in FOLLOW(A)$.

Алгоритм 4.8. Построение таблицы предсказывающего анализатора.

Вход. КС-грамматика $G = (N, T, P, S)$.

Выход. Таблица $M[A, a]$ предсказывающего анализатора, $A \in N$, $a \in T \cup \{\$\}$.

Метод. Для каждого правила вывода $A \rightarrow \alpha$ грамматики выполнить шаги 1 и 2. После этого выполнить шаг 3.

1. Для каждого терминала a из $FIRST(\alpha)$ добавить $A \rightarrow \alpha$ к $M[A, a]$.
2. Если $\epsilon \in FIRST(\alpha)$, то добавить $A \rightarrow \alpha$ к $M[A, b]$ для каждого терминала b из $FOLLOW(A)$. Кроме того, если $\epsilon \in FIRST(\alpha)$ и $\$ \in FOLLOW(A)$, то добавить $A \rightarrow \alpha$ к $M[A, \$]$.
3. Положить все неопределенные входы равными «ошибка».

Пример 4.8. Применим алгоритм 4.7 к грамматике из примера 4.3. Поскольку $FIRST(TE') = FIRST(T) = \{(\text{, id})\}$, в соответствии с правилом вывода $E \rightarrow TE'$ входы $M[E, (\text{, id})]$ и $M[E, id]$ становятся равными $E \rightarrow TE'$.

В соответствии с правилом вывода $E' \rightarrow +TE'$ значение $M[E', +]$ равно $E' \rightarrow +TE'$. В соответствии с правилом вывода $E' \rightarrow \epsilon$ значения $M[E', (\text{, id})]$ и $M[E', \$]$ равны $E' \rightarrow \epsilon$, поскольку $FOLLOW(E') = \{(\text{, id}), \$\}$.

Таблица анализа, построенная по алгоритму 4.8 для этой грамматики, идентична табл. 4.1.

4.4.4. LL(k)-грамматики. Алгоритм 4.8 построения таблицы предсказывающего анализатора может быть применен к любой КС-грамматике. Однако для некоторых грамматик построенная таблица может иметь неоднозначно определенные входы. Например, нетрудно доказать, что если грамматика леворекурсивна или неоднозначна, то таблица будет иметь по крайней мере один неоднозначно определенный вход.

Грамматики, для которых таблица предсказывающего анализатора не имеет неоднозначно определенных входов, называются LL(1)-грамматиками.

Предсказывающий анализатор, построенный для LL(1)-грамматики, называется *LL(1)-анализатором*. Первая буква L в названии связана с тем, что входная цепочка читается слева направо, вторая L обозначает, что строится левый вывод входной цепочки, 1 — что на каждом шаге для принятия решения используется один символ из непрочитанной части входной цепочки.

Алгоритм 4.8 для каждой из LL(1)-грамматик G строит таблицу предсказывающего анализатора, распознающего все цепочки из $L(G)$ и только эти цепочки. Нетрудно доказать также, что если G — LL(1)-грамматика, то $L(G)$ — детерминированный КС-язык.

Справедлив также следующий критерий LL(1)-грамматики. Грамматика $G = (N, T, P, S)$ является LL(1)-грамматикой тогда и только тогда, когда для каждой пары правил $A \rightarrow \alpha$, $A \rightarrow \beta$ из P (т.е. правил с одинаковой левой частью) выполняются следующие два условия:

- 1) $FIRST(\alpha) \cap FIRST(\beta) = \emptyset$;
- 2) Если $e \in FIRST(\alpha)$, то $FIRST(\beta) \cap FOLLOW(A) = \emptyset$.

Пример 4.9. Неоднозначная грамматика не является LL(1). Примером может служить грамматика $G = (\{S, E\}, \{if, then, else, a, b\}, P, S)$ со следующими правилами:

$S \rightarrow if\ E\ then\ S \mid if\ E\ then\ S\ else\ S \mid a$;
 $E \rightarrow b$.

Эта грамматика неоднозначна, что иллюстрируется на рис. 4.4.

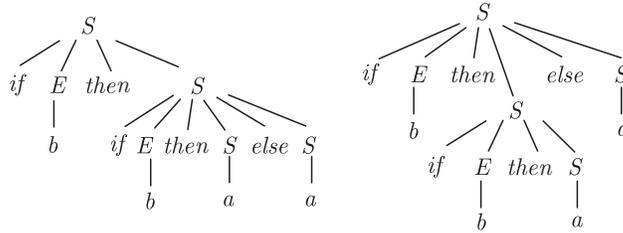


Рис. 4.4

Определение 4.3. КС-грамматика $G = (N, \Sigma, P, S)$ называется *LL(k)-грамматикой* для некоторого фиксированного k , если из

$$1) S \Rightarrow_i^* \omega A \alpha \Rightarrow \iota \omega \beta \alpha \Rightarrow^* \omega x$$

и

$$2) S \Rightarrow_i^* \omega A \alpha \Rightarrow \iota \omega \gamma \alpha \Rightarrow^* \omega y,$$

для которых $FIRST_k(x) = FIRST_k(y)$, вытекает, что $\beta = \gamma$.

Говоря менее формально, G будет *LL(k)-грамматикой*, если для данной цепочки $\omega A \alpha \in (N \cup \Sigma)^*$ и первых k символов (если они существуют), выводящихся из $A \alpha$, существует не более одного правила, которое можно

применить к A , чтобы получить вывод какой-нибудь терминальной цепочки, начинающейся с ω и продолжающейся упомянутыми k терминалами.

Грамматика называется $LL(k)$ -грамматикой, если она $LL(k)$ -грамматика для некоторого k . Доказано, что проблема определения, порождает ли грамматика LL -язык, является алгоритмически неразрешимой.

Теорема 4.7. *КС-грамматика $G = (N, \Sigma, P, S)$ является $LL(k)$ -грамматикой тогда и только тогда, когда для двух различных правил $A \rightarrow \beta$ и $A \rightarrow \gamma$ из P пересечение $FIRST_k(\beta\alpha) \cap FIRST_k(\gamma\alpha)$ пусто при всех таких $\omega A\alpha$, что $S \Rightarrow_i^* \omega A\alpha$.*

Доказательство. Необходимость. Допустим, что ω, A, α, β и γ удовлетворяют условиям теоремы, а $FIRST_k(\beta\alpha) \cap FIRST_k(\gamma\alpha)$ содержит x . Тогда по определению $FIRST$ для некоторых y и z найдутся выводы

$$S \Rightarrow_i^* \omega A\alpha \Rightarrow \omega\beta\alpha \Rightarrow_i^* \omega xy$$

и

$$S \Rightarrow_i^* \omega A\alpha \Rightarrow \omega\gamma\alpha \Rightarrow_i^* \omega xz.$$

(Заметим, что здесь мы использовали тот факт, что N не содержит бесполезных нетерминалов, как это предполагается для всех рассматриваемых грамматик.) Если $|x| < k$, то $y = z = e$. Так как $\beta \neq \gamma$, то G не $LL(k)$ -грамматика.

Достаточность. Допустим, что G не $LL(k)$ -грамматика. Тогда найдутся такие два вывода

$$S \Rightarrow_i^* \omega A\alpha \Rightarrow \omega\beta\alpha \Rightarrow_i^* \omega x$$

и

$$S \Rightarrow_i^* \omega A\alpha \Rightarrow \omega\gamma\alpha \Rightarrow_i^* \omega y,$$

что цепочки x и y совпадают в первых k позициях, но $\beta \neq \gamma$. Поэтому $A \rightarrow \beta$ и $A \rightarrow \gamma$ — различные правила из P и каждое из множеств $FIRST_k(\beta\alpha)$ и $FIRST_k(\gamma\alpha)$ содержит цепочку $FIRST_k(x)$, совпадающую с цепочкой $FIRST_k(y)$. ■

Пример 4.10. Грамматика G , состоящая из двух правил $S \rightarrow aS \mid a$, не будет $LL(1)$ -грамматикой, так как

$$FIRST_1(aS) = FIRST_1(a) = a.$$

Интуитивно это можно объяснить так: видя при разборе цепочки, начинающейся символом a , только этот первый символ, мы не знаем, какое из правил $S \rightarrow aS$ или $S \rightarrow a$ надо применить к S . С другой стороны, G — это $LL(2)$ -грамматика. В самом деле, в обозначениях теоремы 4.7 если $S \Rightarrow_i^* \omega A\alpha$, то $A = S$ и $\alpha = e$. Так как для S даны только два указанных правила, то $\beta = aS$ и $\gamma = a$. Поскольку $FIRST_2(aS) = aS$ и $FIRST_2(a) = a$, то по последней теореме G будет $LL(2)$ -грамматикой.

4.4.5. Удаление левой рекурсии. Основная трудность при использовании предсказывающего анализа — это нахождение такой грамматики для входного языка, по которой можно построить таблицу анализа с однозначно определенными входами. Иногда с помощью некоторых простых преобразований грамматику, не являющуюся LL(1), можно привести к эквивалентной LL(1)-грамматике. Среди этих преобразований наиболее эффективными являются левая факторизация и удаление левой рекурсии. Здесь необходимо сделать два замечания. Во-первых, не всякая грамматика после этих преобразований становится LL(1), и, во-вторых, после таких преобразований получающаяся грамматика может стать менее понятной.

Непосредственную левую рекурсию, т. е. рекурсию вида $A \rightarrow A\alpha$, можно удалить следующим способом. Сначала группируем A -правила:

$$A \rightarrow A\alpha_1 \mid A\alpha_2 \mid \dots \mid A\alpha_m \mid \beta_1 \mid \beta_2 \mid \dots \mid \beta_n,$$

где никакая из строк β_i не начинается с A . Затем заменяем этот набор правил на

$$\begin{aligned} A &\rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_n A'; \\ A' &\rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \alpha_m A' \mid e. \end{aligned}$$

где A' — новый нетерминал. Из нетерминала A можно вывести те же цепочки, что и раньше, но теперь нет левой рекурсии. С помощью этой процедуры удаляются все непосредственные левые рекурсии, но не удаляется левая рекурсия, включающая два или более шагов. Нижеследующий алгоритм позволяет удалить все левые рекурсии из грамматики.

Алгоритм 4.9. Удаление левой рекурсии.

Вход. КС-грамматика G без e -правил (вида $A \rightarrow e$).

Выход. КС-грамматика G' без левой рекурсии, эквивалентная G .

Метод. Выполнить шаги 1 и 2.

1. Упорядочить нетерминалы грамматики G в произвольном порядке.
2. Выполнить следующую процедуру:

```

for (i=1; i<=n; i++) {
  for (j=1; j<=i-1; j++) {
    пусть  $A_j \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_k$  - все текущие правила
    для  $A_j$ ;
    заменить все правила вида  $A_i \rightarrow A_j \alpha$ 
    на правила  $A_i \rightarrow \beta_1 \alpha \mid \beta_2 \alpha \mid \dots \mid \beta_k \alpha$ ;
  }
  удалить правила вида  $A_i \rightarrow A_i$ ;
  удалить непосредственную левую рекурсию в
  правилах для  $A_i$ ;
}

```

После $(i - 1)$ -й итерации внешнего цикла на шаге 2 для любого правила вида $A_k \rightarrow A_s \alpha$, где $k < i$, выполняется $s > k$. В результате на следующей итерации (по i) внутренний цикл (по j) последовательно увеличивает нижнюю границу по t в любом правиле $A_i \rightarrow A_t \alpha$, пока не будет $t \geq i$. Затем, после удаления непосредственной левой рекурсии для A_i -правил, t становится больше i .

Алгоритм 4.9 применим, если грамматика не имеет ϵ -правил (правил вида $A \rightarrow \epsilon$), поскольку при наличии таких правил может нарушиться инвариант цикла при замене правил вида $A_i \rightarrow A_j \alpha$ на правила $A_i \rightarrow \beta_1 \alpha | \beta_2 \alpha | \dots | \beta_k \alpha$. Имеющиеся в грамматике ϵ -правила могут быть удалены предварительно. Получающаяся грамматика без левой рекурсии может иметь ϵ -правила.

4.4.6. Левая факторизация. Основная идея левой факторизации в том, что в том случае, когда неясно, какую из двух альтернатив надо использовать для развертки нетерминала A , нужно изменить A -правила так, чтобы отложить решение до тех пор, пока не будет достаточно информации для принятия правильного решения.

Если $A \rightarrow \alpha \beta_1 | \alpha \beta_2$ — два A -правила и входная цепочка начинается с непустой строки, выводимой из α , то мы не знаем, разворачивать ли по первому правилу или по второму. Можно отложить решение, развернув $A \rightarrow \alpha A'$. Тогда после анализа того, что выводимо из α , можно развернуть по $A' \rightarrow \beta_1$ или по $A' \rightarrow \beta_2$. Левофакторизованные правила принимают вид

$$\begin{aligned} A &\rightarrow \alpha A'; \\ A' &\rightarrow \beta_1 | \beta_2. \end{aligned}$$

Алгоритм 4.10. Левая факторизация грамматики.

Вход. КС-грамматика G .

Выход. Левофакторизованная КС-грамматика G' , эквивалентная G .

Метод. Для каждого нетерминала A найти самый длинный префикс α , общий для двух или более его альтернатив. Если $\alpha \neq \epsilon$, т.е. существует нетривиальный общий префикс, то заменить все A -правила

$$A \rightarrow \alpha \beta_1 | \alpha \beta_2 | \dots | \alpha \beta_n | z,$$

где z обозначает все альтернативы, не начинающиеся с α , на

$$\begin{aligned} A &\rightarrow \alpha A' | z; \\ A' &\rightarrow \beta_1 | \beta_2 | \dots | \beta_n, \end{aligned}$$

где A' — новый нетерминал. Применять это преобразование, пока существует пара альтернатив, имеющих общий префикс.

Пример 4.11. Рассмотрим вновь грамматику условных операторов из примера 4.9:

$$S \rightarrow \text{if } E \text{ then } S \mid \text{if } E \text{ then } S \text{ else } S \mid a;$$

$$E \rightarrow b.$$

После левой факторизации грамматика принимает вид

$$S \rightarrow \text{if } E \text{ then } SS' \mid a;$$

$$S' \rightarrow \text{else } S \mid e;$$

$$E \rightarrow b.$$

К сожалению, грамматика остается неоднозначной, а значит, и не LL(1)-грамматикой.

4.4.7. Рекурсивный спуск. Выше был рассмотрен один из вариантов таблично-управляемого предсказывающего анализа, когда магазин явно использовался в процессе работы анализатора. Возможен иной вариант предсказывающего анализа, в котором каждому нетерминалу сопоставляется процедура (вообще говоря, рекурсивная), а магазин образуется неявно при вызовах таких процедур. Процедуры рекурсивного спуска могут быть записаны, как показано ниже.

В процедуре A для случая, когда имеется правило $A \rightarrow u_i$, такое, что $u_i \Rightarrow^* e$ (напомним, что не может быть больше одного правила, из которого выводится e), приведены два варианта — 1.1 и 1.2. В варианте 1.1 делается проверка, принадлежит ли следующий входной символ $FOLLOW(A)$. Если нет — выдается сообщение об ошибке. В варианте 1.2 этого не делается, так что анализ ошибки перекладывается на процедуру, вызвавшую A .

```
void A(){ //  $A \rightarrow u_1 \mid u_2 \mid \dots \mid u_k$ 
  if ( $InSym \in FIRST(u_i)$ ) // только одному!
    if (parse( $u_i$ ))
      result("A  $\rightarrow u_i$ ");
    else error();
  else
    //Вариант 1:
    if (имеется правило  $A \rightarrow u_i$  такое, что  $u_i \Rightarrow^* e$ )
      //Вариант 1.1
      if ( $InSym \in FOLLOW(A)$ )
        result("A  $\rightarrow u_i$ ");
      else error();
      //Конец варианта 1.1
      //Вариант 1.2:
      result("A  $\rightarrow u_i$ ");
      //Конец варианта 1.2
    //Конец варианта 1
    //Вариант 2:
    if (нет правила  $A \rightarrow u_i$  такого, что  $u_i \Rightarrow^* e$ )
      error();
    //Конец варианта 2
}
```

```

boolean parse(u){ // из u не выводится e!
    v = u;
    while (v ≠ e){ // v = Xz
        if (X - терминал a)
            if (InSym ≠ a)
                return(false);
            else InSym = getInsym();
        else // X - нетерминал B
            B();
            v = z;
    }
    return(true);
}

```

4.4.8. Конструктор LL(1)-анализаторов на Java. В программном приложении приведен пакет LL1, содержащий LL(1)-конструктор и анализатор.

4.4.9. Восстановление процесса анализа после синтаксических ошибок. В приведенных программах рекурсивного спуска была использована процедура реакции на синтаксические ошибки `error()`. В простейшем случае эта процедура выдает диагностику и завершает работу анализатора. Но можно попытаться некоторым разумным образом продолжить работу. Для разбора сверху вниз можно предложить следующий простой алгоритм.

Если в момент обнаружения ошибки на верхушке магазина оказался нетерминальный символ A и для него нет правила, соответствующего входному символу, то сканируем вход до тех пор, пока не встретим символ либо из $FIRST(A)$, либо из $FOLLOW(A)$. В первом случае разворачиваем A по соответствующему правилу, во втором — удаляем A из магазина.

Если на верхушке магазина терминальный символ, то можно удалить все терминальные символы с верхушки магазина вплоть до первого (сверху) нетерминального символа и продолжать так, как это было описано выше.

4.5. Разбор снизу-вверх типа сдвиг-свертка

4.5.1. Основа. В процессе разбора снизу-вверх типа сдвиг-свертка строится дерево разбора входной цепочки, начиная с листьев (снизу) к корню (вверх). Этот процесс можно рассматривать как «свертку» цепочки w к начальному символу грамматики. На каждом шаге свертки подцепочка, которую можно сопоставить правой части некоторого правила вывода, заменяется символом левой части этого правила вывода, и если на каждом шаге выбирается правильная подцепочка, то в обратном порядке прослеживается правосторонний вывод (рис. 4.5). Здесь ко входной цепочке, так же как и при анализе LL(1)-грамматик, приписан концевой маркер $\$$.

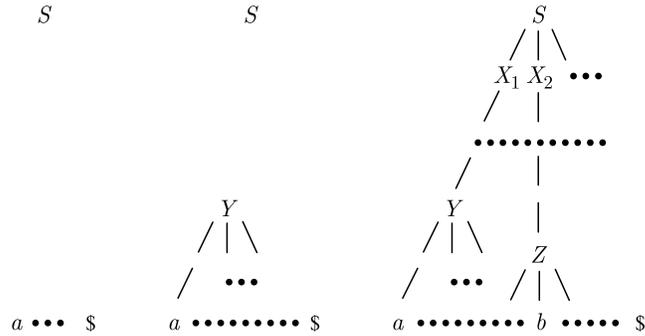


Рис. 4.5

Основой цепочки называется подцепочка сентенциальной формы, которая может быть сопоставлена правой части некоторого правила вывода, свертка по которому к левой части правила соответствует одному шагу в обращении правостороннего вывода. Самая левая подцепочка, которая сопоставляется правой части некоторого правила вывода $A \rightarrow \gamma$, не обязательно является основой, поскольку свертка по правилу $A \rightarrow \gamma$ может дать цепочку, не сводимую к аксиоме.

Формально, основа правой сентенциальной формы z — это правило вывода $A \rightarrow \gamma$ и позиция в z , в которой может быть найдена цепочка γ , такие, что в результате замены γ на A получается предыдущая сентенциальная форма в правостороннем выводе z . Так, если $S \Rightarrow_r^* \alpha A \beta \Rightarrow_r \alpha \gamma \beta$, то $A \rightarrow \gamma$ в позиции, следующей за α , — это основа цепочки $\alpha \gamma \beta$. Подцепочка β справа от основы содержит только терминальные символы.

Вообще говоря, грамматика может быть неоднозначной, поэтому правосторонний вывод $\alpha \gamma \beta$ и основа могут не быть единственными. Если грамматика однозначна, то каждая правая сентенциальная форма грамматики имеет в точности одну основу. Замена основы в сентенциальной форме на нетерминал левой части называется *отсечением* основы. Обращение правостороннего вывода может быть получено с помощью повторного применения отсечения основы, начиная с исходной цепочки w . Если w — слово в рассматриваемой грамматике, то $w = \alpha_n$, где α_n — n -я правая сентенциальная форма еще неизвестного правого вывода $S = \alpha_0 \Rightarrow_r \alpha_1 \Rightarrow_r \dots \Rightarrow_r \alpha_{n-1} \Rightarrow_r \alpha_n = w$.

Чтобы восстановить этот вывод в обратном порядке, выделяем основу γ_n в α_n и заменяем γ_n на левую часть некоторого правила вывода $A_n \rightarrow \gamma_n$, получая $(n - 1)$ -ю правую сентенциальную форму α_{n-1} . Затем повторяем этот процесс, т. е. выделяем основу γ_{n-1} в α_{n-1} и сворачиваем эту основу, получая правую сентенциальную форму α_{n-2} . Если, повторяя этот процесс, мы получаем правую сентенциальную форму, состоящую только из начального символа S , то останавливаемся и сообщаем об успешном завершении

разбора. Обращение последовательности правил, использованных в свертках, есть правый вывод входной строки.

Таким образом, главная задача анализатора типа сдвиг-свертка — это выделение и отсечение основы.

4.5.2. LR(1)-анализаторы. В названии LR(1) символ L указывает на то, что входная цепочка читается слева-направо; R — на то, что строится правый вывод; наконец, 1 указывает на то, что анализатор видит один символ непрочитанной части входной цепочки.

LR(1)-анализ привлекателен по нескольким причинам:

- LR(1)-анализ — наиболее мощный метод анализа без возвратов типа сдвиг-свертка;
- LR(1)-анализ может быть реализован довольно эффективно;
- LR(1)-анализаторы могут быть построены для практически всех конструкций языков программирования;
- класс грамматик, которые могут быть проанализированы LR(1)-методом, строго включает класс грамматик, которые могут быть проанализированы предсказывающими анализаторами (сверху-вниз типа LL(1)).

Схематически структура LR(1)-анализатора изображена на рис. 4.6. Анализатор состоит из входной ленты, выходной ленты, магазина, управляющей программы и таблицы анализа (LR(1)-таблицы), которая имеет две части — *функцию действий (Action)* и *функцию переходов (Goto)*. Управляющая программа одна и та же для всех LR(1)-анализаторов, разные анализаторы отличаются друг от друга только таблицами анализа.

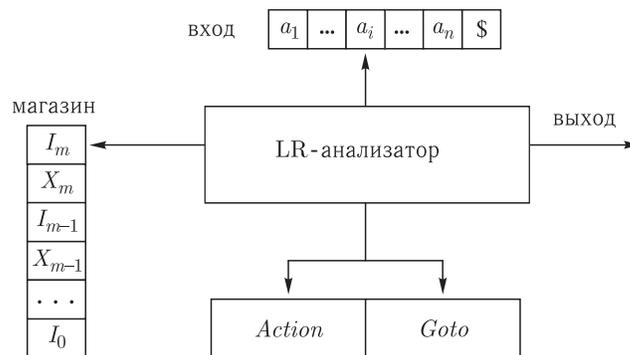


Рис. 4.6

Анализатор читает символы на входной ленте по одному за шаг. В процессе анализа используется магазин, в котором хранятся строки вида $I_0X_1I_1X_2I_2 \dots X_mI_m$ (I_m — верхушка магазина). Каждый X_i — символ грамматики (терминальный или нетерминальный), а I_i — символ *состояния*.

Элемент функции действий $Action[I_m, a_i]$ для символа состояния I_m и входа $a_i \in T \cup \{\$, \}$, может иметь одно из четырех значений:

- 1) `shift I` (сдвиг), где I — символ состояния;
- 2) `reduce $A \rightarrow \gamma$` (свертка по правилу грамматики $A \rightarrow \gamma$);
- 3) `accept` (допуск);
- 4) `error` (ошибка).

Элемент функции переходов $Goto[I_m, A]$ для символа состояния I_m и входа $A \in N$ может иметь одно из двух значений:

- 1) I , где I — символ состояния;
- 2) `error` (ошибка).

Конфигурацией LR(1)-анализатора называется пара, первая компонента которой — содержимое магазина, а вторая — непросмотренный вход:

$$(I_0 X_1 I_1 X_2 I_2 \dots X_m I_m, a_i a_{i+1} \dots a_n \$).$$

Эта конфигурация соответствует правой сентенциальной форме

$$X_1 X_2 \dots X_m a_i a_{i+1} \dots a_n.$$

Если $S \Rightarrow_r^* \gamma A w \Rightarrow_r \gamma \beta w$, то любой префикс δ цепочки $\gamma \beta$ называется *активным* префиксом. Основа сентенциальной формы всегда располагается на верхушке магазина. Таким образом, активный префикс — это такой префикс правой сентенциальной формы, который не переходит правую границу основы этой формы.

Когда анализатор начинает работу, в магазине находится только символ начального состояния I_0 , на входной ленте — анализируемая цепочка с маркером конца.

Каждый очередной шаг анализатора определяется текущим входным символом a_i и символом состояния на верхушке магазина I_m нижеследующим образом.

Пусть LR(1)-анализатор находится в конфигурации

$$(I_0 X_1 I_1 X_2 I_2 \dots X_m I_m, a_i a_{i+1} \dots a_n \$).$$

Анализатор может проделать один из следующих шагов.

1. Если $Action[I_m, a_i] = \text{shift } I$, то анализатор выполняет сдвиг, переходя в конфигурацию

$$(I_0 X_1 I_1 X_2 I_2 \dots X_m I_m a_i I, a_{i+1} \dots a_n \$).$$

Таким образом, в магазин помещаются входной символ a_i и символ состояния I , определяемый $Action[I_m, a_i]$. Текущим входным символом становится a_{i+1} .

2. Если $Action[I_m, a_i] = \text{reduce } A \rightarrow \gamma$, то анализатор выполняет свертку, переходя в конфигурацию

$$(I_0 X_1 I_1 X_2 I_2 \dots X_{m-r} I_{m-r} A I, a_i a_{i+1} \dots a_n \$),$$

где $I = Goto[I_{m-r}, A]$ и r — длина γ , правой части правила вывода.

Анализатор сначала удаляет из магазина $2r$ символов (r символов состояния и r символов грамматики), так что на верхушке оказывается состояние I_{m-r} . Затем анализатор помещает в магазин A — левую часть правила вывода, и I — символ состояния, определяемый $Goto[I_{m-r}, A]$. На шаге свертки текущий входной символ не меняется. Для LR(1)-анализаторов последовательность символов грамматики $X_{m-r+1} \dots X_m$, удаляемых из магазина, всегда соответствует γ — правой части правила вывода, по которому делается свертка.

После осуществления шага свертки генерируется выход LR(1)-анализатора, т. е. исполняются семантические действия, связанные с правилом, по которому делается свертка, например, печатаются номера таких правил.

Заметим, что функция $Goto$ таблицы анализа, построенная по грамматике G , фактически представляет собой функцию переходов детерминированного конечного автомата, распознающего активные префиксы G .

3. Если $Action[I_m, a_i] = \text{accept}$, то разбор успешно завершен.
4. Если $Action[I_m, a_i] = \text{error}$, то анализатор обнаружил ошибку, и выполняются действия по диагностике и восстановлению.

Пример 4.12. Рассмотрим грамматику арифметических выражений $G = (\{E, T, F\}, \{id, +, *\}, P, E)$ с правилами:

- 1) $E \rightarrow E + T$;
- 2) $E \rightarrow T$;
- 3) $T \rightarrow T * F$;
- 4) $T \rightarrow F$;
- 5) $F \rightarrow id$.

В табл. 4.3 представлены функции *Action* и *Goto*, образующие LR(1)-таблицу для этой грамматики. Элемент S_i функции *Action* означает сдвиг и помещение в магазин состояния с номером i , R_j — свертку по правилу номер j , *acc* — допуск, пустая клетка — ошибку. Для функции *Goto* символ i означает помещение в магазин состояния с номером i , пустая клетка — ошибку.

Таблица 4.3

Состояния	<i>Action</i>	<i>Goto</i>
	id + * \$	E T F
0	S6	1 2 3
1	S4 acc	
2	R2 S7 R2	
3	R4 R4 R4	
4	S6	5 3
5	R1 S7 R1	
6	R5 R5 R5	
7	S6	8
8	R3 R3 R3	

В табл. 4.4 описана последовательность состояний магазина и входной ленты на входе $id + id * id$. Например, в первой строке LR-анализатор находится в нулевом состоянии и «видит» первый входной символ id . Действие S_6 в нулевой строке и столбце id в поле *Action* (табл. 4.3) означает сдвиг и помещение символа состояния 6 на верхушку магазина. Это и изображено во второй строке: первый символ id и символ состояния 6 помещаются в магазин, а id удаляется со входной ленты.

Текущим входным символом становится $+$, а действием на вход $+$ в состоянии 6 является свертка по $F \rightarrow id$. Из магазина удаляются два символа (один символ состояния и один символ грамматики). Затем анализируется нулевое состояние. Поскольку *Goto* в нулевом состоянии по символу F — это 3, F и 3 помещаются в магазин. Теперь имеем конфигурацию, соответствующую третьей строке. Остальные шаги определяются аналогично.

Заметим, что в магазине не обязательно должны размещаться одновременно и символы грамматики, и символы состояний. Возможно размещение только символов состояний или только символов грамматики. В этом случае текущее состояние может быть определено путем последовательного чтения символов грамматики в магазине от дна к вершине конечным автоматом, строящимся в следующем подразделе. Одновременное использование и символов грамматики, и символов состояний облегчает понимание поведения LR-анализатора.

Таблица 4.4

Активный префикс	Магазин	Вход	Действие
	0	$id + id * id\$$	сдвиг
id	0 id 6	$+id * id\$$	$F \rightarrow id$
F	0 F 3	$+id * id\$$	$T \rightarrow F$
T	0 T 2	$+id * id\$$	$E \rightarrow T$
E	0 E 1	$+id * id\$$	сдвиг
$E +$	0 E 1 + 4	$id * id\$$	сдвиг
$E + id$	0 E 1 + 4 id 6	$*id\$$	$F \rightarrow id$
$E + F$	0 E 1 + 4 F 3	$*id\$$	$T \rightarrow F$
$E + T$	0 E 1 + 4 T 5	$id\$$	сдвиг
$E + T*$	0 E 1 + 4 T 5 * 7	$id\$$	сдвиг
$E + T * id$	0 E 1 + 4 T 5 * 7 id 6	$\$$	$F \rightarrow id$
$E + T * F$	0 E 1 + 4 T 5 * 7 F 8	$\$$	$T \rightarrow T * F$
$E + T$	0 E 1 + 4 T 5	$\$$	$E \rightarrow E + T$
E	0 E 1		допуск

4.5.3. Конструирование LR(1)-таблицы. Рассмотрим алгоритм конструирования таблицы, управляющей LR(1)-анализатором.

Пусть $G = (N, T, P, S)$ — КС-грамматика. *Пополненной* грамматикой для данной грамматики G называется КС-грамматика

$$G' = (N \cup \{S'\}, T, P \cup \{S' \rightarrow S\}, S'),$$

т. е. эквивалентная грамматика, в которой введены новый начальный символ S' и новое правило вывода $S' \rightarrow S$.

Это дополнительное правило вводится для того, чтобы определить, когда анализатор должен остановить разбор и зафиксировать допуск входа. Таким образом, допуск имеет место тогда и только тогда, когда анализатор готов осуществить свертку по правилу $S' \rightarrow S$.

LR(1)-*ситуацией* называется пара $[A \rightarrow \alpha.\beta, a]$, где $A \rightarrow \alpha\beta$ — правило грамматики, a — терминал или правый концевой маркер $\$$. Вторая компонента ситуации называется *аванцепочкой*.

Будем говорить, что LR(1)-ситуация $[A \rightarrow \alpha.\beta, a]$ *допустима для активного префикса* δ , если существует вывод $S \Rightarrow_r^* \gamma A w \Rightarrow_r \gamma \alpha \beta w$, где $\delta = \gamma \alpha$ и либо a — первый символ w , либо $w = \epsilon$ и $a = \$$.

Будем говорить, что ситуация *допустима*, если она допустима для какого-либо активного префикса.

Пример 4.13. Дана грамматика $G = (\{S, B\}, \{a, b\}, P, S)$ с правилами:

$$S \rightarrow BB;$$

$$B \rightarrow aB \mid b.$$

Существует правосторонний вывод $S \Rightarrow_r^* aaBab \Rightarrow_r^* aaaBab$. Легко видеть, что ситуация $[B \rightarrow a.B, a]$ допустима для активного префикса $\delta = aaa$, если в определении выше положить $\gamma = aa$, $A = B$, $w = ab$, $\alpha = a$, $\beta = B$. Существует также правосторонний вывод $S \Rightarrow_r^* BaB \Rightarrow_r^* BaaB$. Поэтому для активного префикса Vaa допустима ситуация $[B \rightarrow a.B, \$]$.

Центральная идея метода заключается в том, что по грамматике строится детерминированный конечный автомат, распознающий активные префиксы. Для этого ситуации группируются во множества, которые и образуют состояния автомата. Ситуации можно рассматривать как состояния недетерминированного конечного автомата, распознающего активные префиксы, а их группировка на самом деле есть процесс построения детерминированного конечного автомата из недетерминированного.

Анализатор, работающий слева-направо по типу сдвиг-свертка, должен уметь распознавать основы на верхушке магазина. Состояние автомата после прочтения содержимого магазина и текущий входной символ определяют очередное действие автомата. Функцией переходов этого конечного автомата является функция переходов LR-анализатора. Чтобы не просматривать магазин на каждом шаге анализа, на верхушке магазина всегда хранится то состояние, в котором должен оказаться этот конечный автомат после того, как он прочитал символы грамматики в магазине от дна к верхушке.

Рассмотрим ситуацию вида $[A \rightarrow \alpha.B\beta, a]$ из множества ситуаций, допустимых для некоторого активного префикса δ . Тогда существует правосторонний вывод $S\$ \Rightarrow_r^* \gamma Aax \Rightarrow_r^* \gamma \alpha B\beta ax$, где $\delta = \gamma\alpha$, а это значит, что ситуация $[A \rightarrow \alpha B.\beta, a]$ допустима для δB , $B \in N \cup T$. Если $B \in N$ и $B \rightarrow \xi \in R$, то $S\$ \Rightarrow_r^* \gamma \alpha B\beta ax \Rightarrow_r^* \delta Bbw \Rightarrow_r^* \delta \xi bw$. Таким образом, $[B \rightarrow .\xi, b]$ также допустима для δ . Здесь b может быть первым терминалом, выводимым из β , либо из β выводится e в выводе $\beta ax \Rightarrow_r^* bw$ и тогда b равно a . Значит, либо $b \in FIRST(\beta ax)$, либо $\beta \Rightarrow_r^* e$ и $a = \$$.

Рассмотрим теперь недетерминированный конечный автомат, состояниями которого являются LR(1)-ситуации, а переходы определены в соответствии со сказанным выше: если автомат находится в состоянии $[A \rightarrow \alpha.B\beta, a]$, то определен переход в состояние $[A \rightarrow \alpha B.\beta, a]$ по символу B . Если при этом $B \in N$, то определены e -переходы в состояния $[B \rightarrow .\xi, b]$ для всех правил с участием B . Начальным состоянием будем считать $[S' \rightarrow .S, \$]$, все состояния — заключительные.

Теорема 4.8. *Определенный таким образом недетерминированный конечный автомат допускает в точности активные префиксы грамматики.*

Доказательство проводится индукцией по длине активного префикса.

По недетерминированному конечному автомату можно построить эквивалентный детерминированный в соответствии с алгоритмом, изложенным в подразделе 3.4.2. Это осуществляется приведенными ниже функциями `closure`, `goto` и `items`.

Алгоритм 4.11. Конструирование канонической системы множеств допустимых LR(1)-ситуаций.

Вход. КС-грамматика $G = (N, T, P, S)$.

Выход. Каноническая система C множеств допустимых LR(1)-ситуаций для грамматики G .

Метод. Выполнить для пополненной грамматики G' процедуру `items`, которая использует функции `closure` и `goto`.

```

SetOfItemsclosure(SetOfItems I){
  do {SetOfItems J = I;
    for (каждой ситуации  $[A \rightarrow \alpha.B\beta, a]$  из  $J$ ,
        каждого правила вывода  $B \rightarrow \gamma$  из  $G'$ ,
        каждого терминала  $b$  из  $FIRST(\beta a)$ ,
        такого, что  $[B \rightarrow \cdot\gamma, b]$  нет в  $J$ )
      добавить  $[B \rightarrow \cdot\gamma, b]$  к  $J$ ;
    }
  while (к  $J$  можно добавить новую ситуацию);
  return J;
}

SetOfItems goto(SetOfItems I, GrammarSymbol X){
  /* X - символ грамматики */
  Пусть  $J = \{[A \rightarrow \alpha X \beta, a] \mid [A \rightarrow \alpha X \beta, a] \in I\}$ ;
  return closure(J);
}

items(Gramma G') { /* G' - пополненная грамматика */
   $I_0 = \text{closure}(\{[S' \rightarrow \cdot S, \$]\})$ ;
   $C = \{I_0\}$ ;
  do{
    for (каждого множества ситуаций  $I$  из
        системы  $C$ , каждого символа грамматики  $X$ )
      {SetOfItems J = goto(I, X);
        Если  $(J \neq \emptyset) \& (J \notin C)$ 
          добавить J к системе C;
      }
  }
}

```

```

    while (к  $C$  можно добавить новое множество
           ситуаций);
}

```

Если I — множество ситуаций, допустимых для некоторого активного префикса δ , то $\text{goto}(I, X)$ — множество ситуаций, допустимых для активного префикса δX .

Система множеств допустимых LR(1)-ситуаций для всевозможных активных префиксов пополненной грамматики называется *канонической системой* множеств допустимых LR(1)-ситуаций.

Работа алгоритма построения системы C множеств допустимых LR(1)-ситуаций начинается с того, что в C помещается *начальное* множество ситуаций $I_0 = \text{closure}(\{[S' \rightarrow .S, \$]\})$. Затем с помощью функции goto вычисляются новые множества ситуаций и включаются в C . По-существу, $\text{goto}(I, X)$ — переход конечного автомата из состояния I по символу X .

Рассмотрим теперь, как по системе множеств LR(1)-ситуаций строится LR(1)-таблица, т. е. функции действий и переходов LR(1)-анализатора.

Алгоритм 4.12. Построение LR(1)-таблицы.

Вход. Каноническая система $C = \{I_0, I_1, \dots, I_n\}$ множеств допустимых LR(1)-ситуаций для грамматики G .

Выход. Функции $Action$ и $Goto$, составляющие LR(1)-таблицу для грамматики G .

Метод. Для каждого состояния i функции $Action[i, a]$ и $Goto[i, X]$ строятся по множеству ситуаций I_i :

1. Значения функции действия ($Action$) для состояния i определяются следующим образом:
 - а) если $[A \rightarrow \alpha.a\beta, b] \in I_i$ (a — терминал) и $\text{goto}(I_i, a) = I_j$, то полагаем $Action[i, a] = \text{shift } j$;
 - б) если $[A \rightarrow \alpha., a] \in I_i$, причем $A \neq S'$, то полагаем $Action[i, a] = \text{reduce } A \rightarrow \alpha$;
 - в) если $[S' \rightarrow S., \$] \in I_i$, то полагаем $Action[i, \$] = \text{accept}$.
2. Значения функции переходов для состояния i определяются следующим образом: если $\text{goto}(I_i, A) = I_j$, то $Goto[i, A] = j$ (здесь A — нетерминал).
3. Все входы в $Action$ и $Goto$, не определенные шагами 2 и 3, полагаем равными **error**.
4. Начальное состояние анализатора строится из множества, содержащего ситуацию $[S' \rightarrow .S, \$]$.

LR(1)-таблица на основе функций *Action* и *Goto*, полученных в результате работы алгоритма 4.12, называется *канонической*. Работающий с ней LR(1)-анализатор называется *каноническим*.

Пример 4.14. Рассмотрим следующую грамматику, являющуюся пополненной для грамматики из примера 4.8:

- 0) $E' \rightarrow E$;
- 1) $E \rightarrow E + T$;
- 2) $E \rightarrow T$;
- 3) $T \rightarrow T * F$;
- 4) $T \rightarrow F$;
- 5) $F \rightarrow id$.

Множества ситуаций и переходы по *goto* для этой грамматики приведены на рис. 4.7. LR(1)-таблица для этой грамматики приведена в табл. 4.3.

4.5.4. Конструктор LR(1)-анализаторов на Java. В программном приложении приведен пакет LR1, содержащий LR(1)-конструктор и анализатор. Используются следующие структуры данных:

```
HashMap GrammarR = new HashMap();
/* Набор отображений номер правила-> Левая часть, правая часть */
HashMap GrammarN = new HashMap();
/* Для каждого нетерминала Набор отображений номер правила ->
/* правая часть */
HashSet Nonterms = new HashSet(), Terminals = new HashSet();

LinkedList statesList = new LinkedList();
/*
* Список необработанных состояний; содержит номера состояний из
* canonicalSystem
*/
HashMap canonicalSystem = new HashMap();
/* Это отображение Номер состояния -> множество LR1 ситуаций */
HashMap actionTable = new HashMap();
/* Это отображение Номер состояния ->входной символ -> действие */
HashMap gotoTable = new HashMap();

/* Это отображение Номер состояния ->нетерминал -> состояние */
```

4.5.5. Корректность построения.

Теорема 4.9. Если $G = (N, T, P, S')$ — LR(1)-грамматика, то вывод $S \Rightarrow_r \pi w$ существует тогда и только тогда, когда существует цепочка шагов LR(1)-анализатора такова, что $(I_0, w\$, e) \vdash^* (I_0 S I, \$, \pi)$, $Action(I, \$) = accept$.

Здесь π — цепочка номеров примененных правил правостороннего вывода и эта цепочка добавлена в порядке применения свертки (последняя слева) в конфигурацию анализатора.

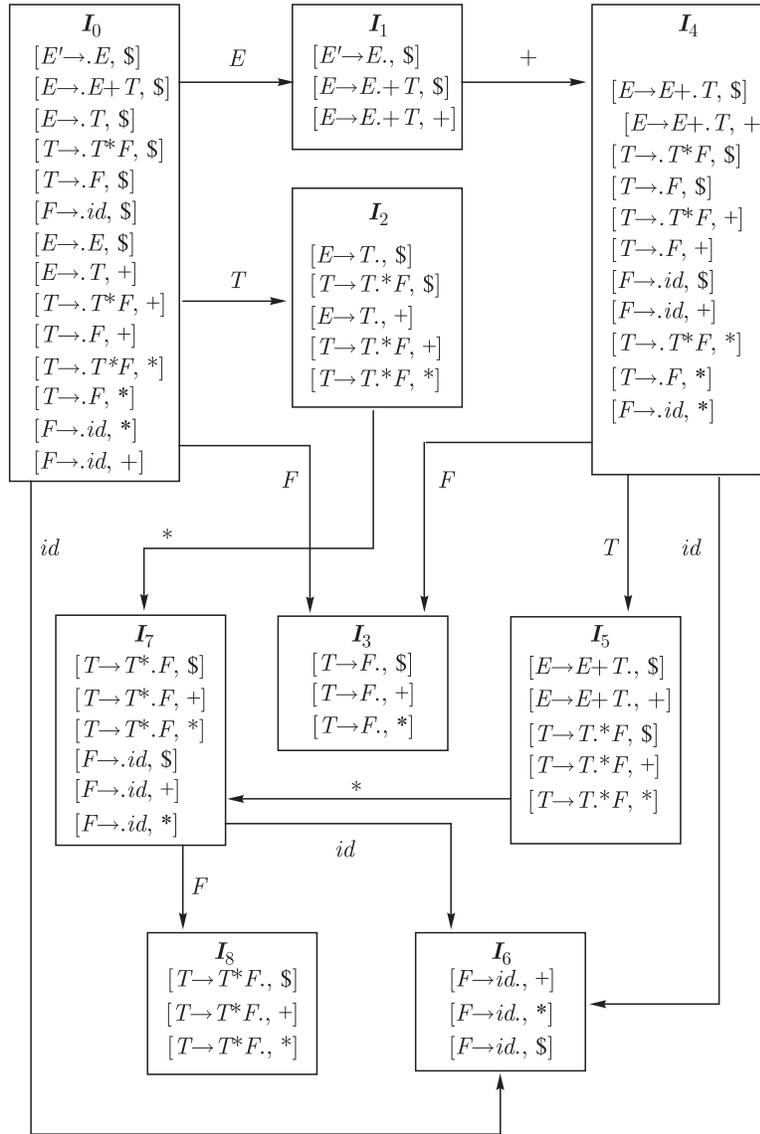


Рис. 4.7

Доказательство. Достаточность. Индукцией по $l = |\pi|$ покажем, что если $S \Rightarrow_r^{\pi} \alpha x$, $\alpha \in (N \cup T)^*$, $x \in T^*$, то $(I_0 X_1 I_1 \dots X_m I_m, x \$, e) \vdash^* (I_0 S I, \$, \pi)$, где $X_1 X_2 \dots X_m = \alpha$, I_0 — начальное состояние $DFAG$, I_i — состояние $DFAG$ после прочтения $X_1 X_2 \dots X_i$ (а не состояние с номером i).

Базис. Пусть $l = 1$, т.е. $\pi = i$. Тогда $S \Rightarrow_r^i \alpha x$, т.е. существует правило $(i)S \rightarrow \alpha x \in P$. Пусть $x = a_1 a_2 \dots a_p$. Согласно определению активных

префиксов и допустимых для них ситуаций:

$$[S \rightarrow \alpha.a_1 \dots a_j \dots a_p, \$] \in I_m = DFA^G(\alpha) = DFA^G(X_1 X_2 \dots X_m);$$

$$[S \rightarrow \alpha a_1.a_2 \dots a_p, \$] \in I_m = DFA^G(\alpha a_1);$$

$$[S \rightarrow \alpha a_1 \dots a_{p-1}.a_p, \$] \in I_m = DFA^G(\alpha a_1 \dots a_{p-1});$$

$$[S \rightarrow \alpha a_1 \dots a_p., \$] \in I_m = DFA^G(\alpha a_1 \dots a_p) = DFA^G(\alpha x).$$

Тогда по построению $Action(I_{m+j-1}, a_j) = shift\ I_{m+j}$, $Action(I_{m+j}, \$) = reduce\ i$. Поэтому $(I_0 X_1 I_1 \dots X_m I_m, x \$, e) = (I_0 X_1 I_1 \dots X_m I_m, a_1 \dots a_p \$, e) \vdash^* (I_0 X_1 I_1 \dots X_m I_m I_{m+1} \dots a_p I_{m+p}, \$, e) \vdash (I_0 ST, \$, i)$, ч. т. д.

Шаг индукции. Предположим, что утверждение выполняется для $l \leq n$ ($n \geq 1$). Покажем, что оно выполняется для $l = n + 1$.

Имеем $S \Rightarrow \overset{\pi'}{r} \alpha' A w \Rightarrow \overset{i}{r} \alpha' \beta w = \alpha x$. Очевидно, что x заканчивается цепочкой w , т. е. $x = zw$. Тогда: $\alpha, \alpha' \in (N \cup T)^*$, $w, x, z \in T^*$, $\pi = \pi' i$; $A \rightarrow \beta \in P$, $\beta = \beta' z$, $\alpha = \alpha' \beta'$.

Имеющийся вывод можно представить как $S \Rightarrow \overset{\pi'}{r} \alpha' A w \Rightarrow \overset{(i)}{r} = \alpha' \beta w = \alpha' \beta' z w$. Имеем $A \rightarrow \beta = \beta' z \in P$, тогда $\alpha' \beta'$ — активный префикс $\alpha' \beta' z w$, причем $DFA^G(\alpha' \beta') = DFA^G(\alpha) = I_m$; поэтому $[A \rightarrow \beta'.z, u] \in I_m$, $u = FIRST(w)$ (либо $u = \$$, если $w = e$). Пусть $z = a_1 \dots a_p$. По построению $Action(I_m, a_1) = shift\ I_{m+1}$, $\dots Action(I_{m+p-1}, a_p) = shift\ I_{m+p}$, $I_{m+p} = DFA^G(\alpha' \beta)$, $Action(I_{m+p}, u) = reduce\ I$, $u = FIRST(w)$ (либо $u = \$$).

Анализатор осуществляет следующие такты:

$$\begin{aligned} & (I_0 X_1 I_1 \dots X_j I_j X_{j+1} I_{j+1} \dots X_m I_m, x \$, e) = \\ & (I_0 X_1 I_1 \dots X_j I_j X_{j+1} I_{j+1} \dots X_m I_m, zw \$, e) = \\ & (I_0 X_1 I_1 \dots X_j I_j X_{j+1} I_{j+1} \dots X_m I_m, a_1 \dots a_p w \$, e) \vdash \\ & (I_0 X_1 I_1 \dots X_j I_j X_{j+1} I_{j+1} \dots X_m I_m a_1 I_{m+1}, a_2 \dots a_p w \$, e) \vdash \\ & (I_0 X_1 I_1 \dots X_j I_j X_{j+1} I_{j+1} \dots X_m I_m a_1 I_{m+1} \dots a_p I_{m+p}, w \$, e) = \\ & (I_0 X_1 I_1 \dots X_j I_j A I'_{j+1}, w \$, i) \vdash^* (I_0 ST, \$, \pi' i) = (I_0 SI, \$, \pi). \end{aligned}$$

Последний переход обоснован предположением индукции с учетом того, что $I_{j+1} = Goto(I_j, A)$, $I'_{j+1} = DFA^G(\alpha' A)$, $\alpha' = X_1 X_2 \dots X_j$. В частности, если $\alpha = e$, т. е. $S \Rightarrow \overset{\pi}{r} x$, то $(I_0, x \$, e) \vdash^* (I_0 SI, \$, \pi)$.

Необходимость. Индукцией по числу l тактов анализатора покажем, что если $(I_0, w \$, e) \vdash^l (I_0 X_1 I_1 \dots X_j I_j X_{j+1} I_{j+1} \dots X_m I_m, x \$, \pi)$, то $\alpha x \Rightarrow \overset{\pi}{r} w$, $\alpha = X_1 \dots X_j \dots X_m$, $x, w \in T^*$, т. е. содержимое магазина в конкатенации с непросмотренной входной строкой представляет правосентенциальную форму, выводимую с помощью последовательности правил π .

Базис. Пусть $l = 1$. Если это $shift$, то пусть $(I_0, w \$, e) = (I_0, X_1 x \$, e) \vdash (I_0 X_1 I, x \$, e)$, $w = X_1 x$. Тогда $X_1 x \Rightarrow w$ за 0 шагов. Если это шаг $reduce\ i$, то, поскольку в начальной конфигурации $(I_0, w \$, e)$ магазин пуст (состояние I_0 не учитываем), основа также пуста. Значит, $[A \rightarrow e., u] \in I_0$, $(I_0, w \$, e) \vdash (I_0 AI, w \$, e)$, $u = FIRST(w)$, (либо $u = \$$), так что $A w \Rightarrow \overset{(i)}{r} w$.

Шаг. Пусть анализатор осуществляет $l = n + 1$ шагов. Первые n шагов дают $(I_0, w\$, e) \vdash^n (I_0 X_1 I_1 \dots X_j I_j X_{j+1} I_{j+1} \dots X_m I_m, x'\$, \pi')$.

По предположению индукции $X_1 \dots X_j \dots X_m x' \Rightarrow_r^{\pi'} w$, затем осуществляется последний шаг. Пусть это *shift*. Тогда:

$$\begin{aligned} & (I_0 X_1 I_1 \dots X_j I_j X_{j+1} I_{j+1} \dots X_m I_m, x'\$, \pi') = \\ & (I_0 X_1 I_1 \dots X_j I_j X_{j+1} I_{j+1} \dots X_m I_m, X_{m+1} x'\$, \pi') \vdash \\ & (I_0 X_1 I_1 \dots X_j I_j X_{j+1} I_{j+1} \dots X_m I_m X_{m+1}, x'\$, \pi'). \\ & X_1 \dots X_j \dots X_m X_{m+1} x' \Rightarrow_r^{\pi'} w, \end{aligned}$$

поскольку

$$x' = X_{m+1} x, X_1 \dots X_j \dots X_m X_{m+1} x = X_1 X_j \dots X_m X_{m+1} x' \Rightarrow_r^{\pi'} w$$

по предположению индукции.

Пусть это *reduce*. За первые n шагов достигнута конфигурация $(I_0 X_1 I_1 \dots X_j I_j X_{j+1} I_{j+1} \dots X_m I_m, x'\$, \pi')$. Затем совершается последний такт — свертка по i -му правилу. По построению $I_m = DFA^G(X_1 \dots X_j \dots X_m)$ и существует LR(1)-ситуация

$$\begin{aligned} & [A \rightarrow Y_1 Y_2 \dots Y_p, u'] \in I_m, u' \in FIRST(x') \text{ (либо \$)}, \\ & A \rightarrow Y_1 Y_2 \dots Y_p \in P, Y_1 Y_2 \dots Y_p = X_{m-p+1} \dots X_m. \end{aligned}$$

Тогда

$$\begin{aligned} & (I_0 X_1 I_1 \dots X_j I_j X_{j+1} I_{j+1} \dots X_m I_m, X_{m+1}, x'\$, \pi') = \\ & = (I_0 X_1 I_1 \dots X_{m-p} I_{m-p} Y_1 I_{m-p+1} Y_2 \dots Y_p I_m, x'\$, \pi') \vdash \\ & \vdash (I_0 X_1 I_1 \dots X_{m-p} I_{m-p} A I'_{m-p+1}, x'\$, i\pi'), \end{aligned}$$

где $I'_{m-p+1} = Goto(I_{m-p}, A)$.

Кроме того, $X_1 X_2 \dots X_{m-p} A x' \Rightarrow_r^{i\pi'} w$, поскольку $X_1 X_2 \dots X_{m-p} A x' \Rightarrow_r^i X_1 X_2 \dots X_{m-p} Y_1 Y_2 \dots Y_p x' = X_1 X_2 \dots X_{m-p} X_{m-p+1} \dots X_{m-1} X_m x' \Rightarrow_r^{\pi'} w$ ввиду предположения индукции. В частности, при $\alpha = S$, $x = e$ получаем, что если $(I_0, w\$, e) \vdash^* (I_0 S I, \$, \pi)$, то $S \Rightarrow_r^{\pi} w$. ■

4.5.6. LR(k)-грамматики. Если для КС-грамматики G функция *Action*, полученная в результате работы алгоритма 4.12, не содержит неоднозначно определенных входов, то грамматика называется *LR(1)-грамматикой*.

Язык L называется *LR(1)-языком*, если он может быть порожден некоторой LR(1)-грамматикой.

Иногда используется другое определение LR(1)-грамматики. Грамматика называется *LR(1)*, если из условий

- 1) $S' \Rightarrow_r^* u A w \Rightarrow_r u v w$;
- 2) $S' \Rightarrow_r^* z B x \Rightarrow_r u v y$;
- 3) $FIRST(w) = FIRST(y)$

следует, что $u A y = z B x$ (т. е. $u = z$, $A = B$ и $x = y$).

В случае произвольного $k \geq 0$ определение таково.

Определение 4.4. Пусть $G = (N, \Sigma, P, S)$ — КС-грамматика и $G' = (N', \Sigma, P', S')$ — полученная из нее пополненная грамматика. Грамматiku G будем называть $LR(k)$ -грамматикой для $k \geq 0$, если из условий

- 1) $S' \Rightarrow_{G'}^* \alpha A \omega \Rightarrow_{G'}^* \alpha \beta \omega$;
- 2) $S' \Rightarrow_{G'}^* \gamma B x \Rightarrow_{G'}^* \alpha \beta y$;
- 3) $FIRST_k(\omega) = FIRST_k(y)$

следует, что $\alpha A y = \gamma B x$ (т. е. $\alpha = \gamma$, $A = B$ и $x = y$).

Грамматика G называется LR -грамматикой, если она $LR(k)$ -грамматика для некоторого $k \geq 0$.

Интуитивно это определение говорит о том, что если $\alpha\beta\omega$ и $\alpha\beta y$ — правывыводимые цепочки пополненной грамматикой, у которых $FIRST_k(\omega) = FIRST_k(y)$, и $A \rightarrow \beta$ — последнее правило, использованное в правом выводе цепочки $\alpha\beta\omega$, то правило $A \rightarrow \beta$ должно использоваться также в правом разборе при свертке $\alpha\beta y$ к $\alpha A y$. Так как A дает β независимо от ω , то $LR(k)$ -условие говорит о том, что в $FIRST_k(\omega)$ содержится информация, достаточная для определения того, что $\alpha\beta$ за один шаг выводится из αA . Поэтому никогда не может возникнуть сомнений относительно того, как свернуть очередную правывыводимую цепочку пополненной грамматикой. Кроме того, работая с $LR(k)$ -грамматикой, мы всегда знаем, допустить ли данную входную цепочку или продолжать разбор.

Пример 4.15. Рассмотрим грамматiku G с правилами

$$S \rightarrow Sa \mid a.$$

Согласно определению, G не $LR(0)$ -грамматика, так как из трех условий:

- 1) $S' \Rightarrow_{G'}^0 S' \Rightarrow_{G'} S$;
- 2) $S' \Rightarrow_{G'} S \Rightarrow_{G'} Sa$;
- 3) $FIRST_0(e) = FIRST_0(a) = e$

не следует, что $S'a = S$. Применяя определение к этой ситуации, имеем $\alpha = e$, $\beta = S$, $\omega = e$, $\gamma = e$, $A = S'$, $B = S$, $x = e$ и $y = a$. Проблема здесь заключается в том, что нельзя установить, является ли S основой правывыводимой цепочки Sa , не видя символа, стоящего после S (т. е. наблюдая «нулевое» количество символов). Интуитивно представляется, что G не должна быть $LR(0)$ -грамматикой, и она не будет ею, если пользоваться первым определением. Это определение мы и будем использовать далее.

Пример 4.16. Пусть G — левосторонняя грамматика с правилами:

$$S \rightarrow Ab \mid Bc;$$

$$A \rightarrow Aa \mid e;$$

$$B \rightarrow Ba \mid e.$$

Заметим, что G не является LR(k)-грамматикой ни для какого k .

Допустим, что G — LR(k)-грамматика. Рассмотрим два правых вывода в пополненной грамматике G' :

$$S' \Rightarrow_r S \Rightarrow_r^* Aa^k b \Rightarrow_r a^k b$$

и

$$S' \Rightarrow_r S \Rightarrow_r^* Ba^k c \Rightarrow_r a^k c.$$

Эти два вывода удовлетворяют условию из определения LR(k)-грамматики при $\alpha = e$, $\beta = e$, $\omega = a^k b$, $\gamma = e$ и $y = a^k c$. Но так как заключение неверно, т. е. $A \neq B$, то G — не LR(k)-грамматика. Более того, так как LR(k)-условие нарушается для всех k , то G — не LR-грамматика.

Если грамматика не является LR(1), то анализатор типа сдвиг–свертка при анализе некоторой цепочки может достигнуть конфигурации, в которой он, зная содержимое магазина и следующий входной символ, не может решить, делать ли сдвиг или свертку (конфликт сдвиг–свертка), или не может решить, какую из нескольких сверток применить (конфликт свертка–свертка).

В частности, неоднозначная грамматика не может быть LR(1). Для доказательства рассмотрим два различных правых вывода:

$$1) S \Rightarrow_r u_1 \Rightarrow_r \dots \Rightarrow_r u_n \Rightarrow_r w, \quad 2) S \Rightarrow_r v_1 \Rightarrow_r \dots \Rightarrow_r v_m \Rightarrow_r w.$$

Нетрудно заметить, что LR(1)-условие (согласно второму определению LR(1)-грамматики) нарушается для наименьшего из чисел i , для которых $u_{n-i} \neq v_{m-i}$.

Пример 4.17. Рассмотрим еще раз грамматику условных операторов:

$$\begin{aligned} S &\rightarrow \text{if } E \text{ then } S \mid \text{if } E \text{ then } S \text{ else } S \mid a; \\ E &\rightarrow b. \end{aligned}$$

Если анализатор типа сдвиг–свертка находится в такой конфигурации, что необработанная часть входной цепочки имеет вид *else ... \$*, а в магазине находится *... if E then S*, то нельзя определить, является ли *if E then S* основой, вне зависимости от того, что лежит в магазине ниже. Это конфликт сдвиг–свертка. В зависимости от того, что следует на входе за *else*, правильной может быть свертка по $S \rightarrow \text{if } E \text{ then } S$ или сдвиг *else*, а затем разбор другого S и завершение основы $\text{if } E \text{ then } S \text{ else } S$. Таким образом, нельзя сказать, нужно ли в этом случае делать сдвиг или свертку, так что грамматика не является LR(1).

Эта грамматика может быть преобразована к LR(1)-виду следующим образом:

$$\begin{aligned} S &\rightarrow M \mid U; \\ M &\rightarrow \text{if } E \text{ then } M \text{ else } M \mid a; \\ U &\rightarrow \text{if } E \text{ then } S \mid \text{if } E \text{ then } M \text{ else } U; \\ E &\rightarrow b. \end{aligned}$$

Основная разница между LL(1)- и LR(1)-грамматиками заключается в следующем. Чтобы грамматика была LR(1)-грамматикой, необходимо распознавать вхождение правой части правила вывода, просмотрев всё, что выведено из этой правой части, и текущий символ входной цепочки. Это требование существенно менее строгое, чем требование для LL(1)-грамматики, когда

необходимо определить применимое правило, видя только первый символ, выводимый из его правой части. Таким образом, класс LL(1)-грамматик есть собственный подкласс класса LR(1)-грамматик.

Справедливы также следующие утверждения [2].

Теорема 4.10. *Каждый LR(1)-язык является детерминированным КС-языком.*

Теорема 4.11. *Если L — детерминированный КС-язык, то существует LR(1)-грамматика, порождающая L .*

Теорема 4.12. *Для любой LR(k)-грамматики при $k > 1$ существует эквивалентная ей LR($k - 1$)-грамматика.*

Доказано, что проблема определения, порождает ли грамматика LR-язык, является алгоритмически неразрешимой.

4.5.7. Восстановление процесса анализа после синтаксических ошибок. Одним из простейших методов восстановления после ошибки при LR(1)-анализе является следующий. При синтаксической ошибке просматриваем магазин от верхушки, пока не найдем состояние s с переходом на выделенный нетерминал A . Затем сканируются входные символы, пока не будет найден такой, который допустим после A . В этом случае на верхушку магазина помещается состояние $Goto[s, A]$ и разбор продолжается. Для нетерминала A возможны несколько таких вариантов. Обычно A — это нетерминал, представляющий одну из основных конструкций языка, например, оператор.

При более детальной проработке реакции на ошибки можно в каждой пустой клетке анализатора поставить обращение к своей подпрограмме. Такая подпрограмма может вставлять или удалять входные символы или символы магазина, менять порядок входных символов.

4.5.8. Варианты LR-анализаторов. Построенные таблицы для LR(1)-анализатора зачастую получаются довольно большими. Поэтому при практической реализации используются различные методы их сжатия. С другой стороны, часто оказывается, что при построении синтаксического анализатора типа «сдвиг-свертка» достаточно более простых методов. Некоторые из этих методов базируются на основе LR(1)-анализаторов.

Одним из способов такого упрощения является LR(0)-анализ — частный случай LR-анализа, когда ни при построении таблиц, ни при анализе не учитывается аванцепочка.

Еще одним вариантом LR-анализа являются так называемые SLR(1)-анализаторы (Simple LR(1)). Они строятся по следующей схеме. Пусть $C = \{I_0, I_1, \dots, I_n\}$ — набор множеств допустимых LR(0)-ситуаций. Состояния анализатора соответствуют I_i . Функции действий и переходов анализатора определяются следующим образом.

1. Если $[A \rightarrow u.av] \in I_i$ и $\text{goto}(I_i, a) = I_j$, то определим $\text{Action}[i, a] = \text{shift } j$.
2. Если $[A \rightarrow u.] \in I_i$, то, для всех $a \in \text{FOLLOW}(A)$, $A \neq S'$, определим $\text{Action}[i, a] = \text{reduce } A \rightarrow u$.
3. Если $[S' \rightarrow S.] \in I_i$, то определим $\text{Action}[i, \$] = \text{акцепт}$.
4. Если $\text{goto}(I_i, A) = I_j$, где $A \in N$, то определим $\text{Goto}[i, A] = j$.
5. Остальные входы для функций Action и Goto определим как **error**.
6. Начальное состояние соответствует множеству ситуаций, содержащему ситуацию $[S' \rightarrow .S]$.

Распространенным вариантом LR(1)-анализа является также LALR(1)-анализ. Он основан на объединении (слиянии) некоторых таблиц. Назовем *ядром* множества LR(1)-ситуаций множество их первых компонент (т.е. во множестве ситуаций не учитываются аванцепочки). Объединим все множества ситуаций с одинаковыми ядрами, а в качестве аванцепочек возьмем объединение аванцепочек. Функции Action и Goto строятся очевидным образом. Если функция Action таким образом построенного анализатора не имеет конфликтов, то он называется *LALR(1)-анализатором* (LookAhead LR(1)). Если грамматика является LR(1), то в таблицах LALR(1)-анализатора могут появиться конфликты типа свертка/свертка (если одно из объединяемых состояний имело ситуации $[A \rightarrow \alpha, a]$ и $[B \rightarrow \beta, b]$, а другое — $[A \rightarrow \alpha, b]$ и $[B \rightarrow \beta, a]$, то в LALR(1) появятся ситуации $[A \rightarrow \alpha, \{a, b\}]$ и $[B \rightarrow \beta, \{b, a\}]$). Конфликты типа сдвиг/свертка появиться не могут, поскольку аванцепочка для сдвига во внимание не принимается.

Глава 5

ЭЛЕМЕНТЫ ТЕОРИИ ПЕРЕВОДА

До сих пор мы рассматривали процесс синтаксического анализа только как процесс анализа допустимости входной цепочки. Однако в компиляторе синтаксический анализ служит основой еще одного важного шага — построения дерева синтаксического анализа. В примерах 4.3 и 4.8 в процессе синтаксического анализа в качестве выхода выдавалась последовательность примененных правил, на основе которой и может быть построено дерево. Построение дерева синтаксического анализа является простейшим частным случаем *перевода* — процесса преобразования некоторой входной цепочки в некоторую выходную.

Определение 5.1. Пусть T — входной алфавит, а Π — выходной алфавит. *Переводом* (или *трансляцией*) с языка $L_1 \subseteq T^*$ на язык $L_2 \subseteq \Pi^*$ называется отображение $\tau : L_1 \rightarrow L_2$. Если $y = \tau(x)$, то цепочка y называется *выходом* для цепочки x .

Мы рассмотрим несколько формализмов для определения переводов: преобразователи с магазинной памятью, схемы синтаксически управляемого перевода и атрибутные грамматики.

5.1. Преобразователи с магазинной памятью

Рассмотрим важный класс абстрактных устройств, называемых преобразователями с магазинной памятью. Эти преобразователи получаются из автоматов с магазинной памятью, если к ним добавить выход и позволить на каждом шаге выдавать выходную цепочку.

Преобразователем с магазинной памятью (МП-преобразователем) называется восьмерка $P = (Q, T, \Gamma, \Pi, D, q_0, Z_0, F)$, где все символы имеют тот же смысл, что и в определении МП-автомата, за исключением того, что Π — конечный выходной алфавит, а D — отображение множества $Q \times (T \cup \{e\}) \times \Gamma$ во множество конечных подмножеств множества $Q \times \Gamma^* \times \Pi^*$.

Определим *конфигурацию* преобразователя P как четверку (q, x, u, y) , где $q \in Q$ — состояние, $x \in T^*$ — цепочка на входной ленте, $u \in \Gamma^*$ —

содержимое магазина, $y \in \Pi^*$ — цепочка на выходной ленте, выданная вплоть до настоящего момента.

Если множество $D(q, a, Z)$ содержит элемент (r, u, z) , то будем писать $(q, ax, Zw, y) \vdash (r, x, uw, yz)$ для любых $x \in T^*$, $w \in \Gamma^*$ и $y \in \Pi^*$. Рефлексивно-транзитивное замыкание отношения \vdash будем обозначать \vdash^* .

Цепочку y назовем *выходом* для x , если $(q_0, x, Z_0, e) \vdash^* (q, e, u, y)$ для некоторых $q \in F$ и $u \in \Gamma^*$. *Переводом* (или *трансляцией*), определяемым МП-преобразователем P (обозначается $\tau(P)$), назовем множество

$$\{(x, y) \mid (q_0, x, Z_0, e) \vdash^* (q, e, u, y) \text{ для некоторых } q \in F \text{ и } u \in \Gamma^*\}.$$

Будем говорить, что МП-преобразователь P является *детерминированным* (ДМП-преобразователем), если выполняются следующие условия:

- 1) для всех $q \in Q$, $a \in T \cup \{e\}$ и $Z \in \Gamma$ множество $D(q, a, Z)$ содержит не более одного элемента;
- 2) если $D(q, e, Z) \neq \emptyset$, то $D(q, a, Z) = \emptyset$ для всех $a \in T$.

Пример 5.1. Рассмотрим перевод τ , отображающий каждую цепочку $x \in \{a, b\}^*\$$, в которой число вхождений символа a равно числу вхождений символа b , в цепочку $y = (ab)^n$, где n — число вхождений a или b в цепочку x . Например, $\tau(abbaab\$) = ababab$.

Этот перевод может быть реализован ДМП-преобразователем

$$P = (\{q_0, q_f\}, \{a, b, \$\}, \{Z, a, b\}, \{a, b\}, D, q_0, Z, \{q_f\})$$

с функцией переходов:

$$\begin{aligned} D(q_0, X, Z) &= \{(q_0, XZ, e)\}, X \in \{a, b\}, \\ D(q_0, \$, Z) &= \{(q_f, Z, e)\}, \\ D(q_0, X, X) &= \{(q_0, XX, e)\}, X \in \{a, b\}, \\ D(q_0, X, Y) &= \{(q_0, e, ab)\}, X \in \{a, b\}, Y \in \{a, b\}, X \neq Y. \end{aligned}$$

5.2. Синтаксически управляемый перевод

Другим формализмом, используемым для определения переводов, является схема синтаксически управляемого перевода. Фактически такая схема представляет собой КС-грамматику, в которой к каждому правилу добавлен элемент перевода. Всякий раз, когда правило участвует в выводе входной цепочки, с помощью элемента перевода вычисляется часть выходной цепочки, соответствующая части входной цепочки, порожденной этим правилом.

5.2.1. Схемы синтаксически управляемого перевода.

Определение 5.2. *Схемой синтаксически управляемого перевода* (или *трансляции*, сокращенно: *СУ-схемой*) называется пятерка $Tr = (N, T, \Pi, R, S)$, где:

- 1) N — конечное множество нетерминальных символов;
- 2) T — конечный входной алфавит;
- 3) Π — конечный выходной алфавит;
- 4) R — конечное множество правил перевода вида

$$A \rightarrow u, v$$

где $u \in (N \cup T)^*$, $v \in (N \cup \Pi)^*$ и вхождения нетерминалов в цепочку v образуют перестановку вхождений нетерминалов в цепочку u , так что каждому вхождению нетерминала B в цепочку u соответствует некоторое вхождение этого же нетерминала в цепочку v ; если нетерминал B встречается более одного раза, то для указания соответствия используются верхние целочисленные индексы;

- 5) S — начальный символ, выделенный нетерминал из N .

Определим *выводимую пару* в схеме Tr следующим образом:

- 1) (S, S) — выводимая пара, в которой символы S соответствуют друг другу;
- 2) если $(xAy, x'Ay')$ — выводимая пара, в цепочках которой вхождения A соответствуют друг другу, и $A \rightarrow u, v$ — правило из R , то $(xuy, x'vy')$ — выводимая пара. Для обозначения такого вывода одной пары из другой будем пользоваться обозначением \Rightarrow : $(xAy, x'Ay') \Rightarrow (xuy, x'vy')$. Рефлексивно-транзитивное замыкание отношение \Rightarrow обозначим \Rightarrow^* .

Переводом $\tau(Tr)$, определяемым СУ-схемой Tr , назовем множество пар

$$\{(x, y) \mid (S, S) \Rightarrow^*(x, y), x \in T^*, y \in \Pi^*\}.$$

Если через P обозначить множество входных правил вывода всех правил перевода, то $G = (N, T, P, S)$ будет *входной грамматикой* для Tr .

СУ-схема $Tr = (N, T, \Pi, R, S)$ называется *простой*, если для каждого правила $A \rightarrow u, v$ из R соответствующие друг другу вхождения нетерминалов встречаются в u и v в одном и том же порядке.

Перевод, определяемый простой СУ-схемой, называется *простым синтаксически управляемым переводом* (простым СУ-переводом).

Пример 5.2. Перевод арифметических выражений в ПОЛИЗ (польскую инверсную запись) можно осуществить простой СУ-схемой с правилами

$$\begin{aligned} E &\rightarrow E + T, & ET+; \\ E &\rightarrow T, & T; \\ T &\rightarrow T * F, & TF*; \\ T &\rightarrow F, & F; \\ F &\rightarrow id, & id; \\ F &\rightarrow (E), & E. \end{aligned}$$

Найдем выход схемы для входа $id * (id + id)$. Нетрудно видеть, что существует последовательность шагов вывода

$$\begin{aligned}
 (E, E) &\Rightarrow (T, T) \Rightarrow (T * F, TF*) \Rightarrow \\
 &\Rightarrow (F * F, FF*) \Rightarrow (id * F, id F*) \Rightarrow (id * (E), id E*) \Rightarrow \\
 &\Rightarrow (id * (E + T), id ET + *) \Rightarrow (id * (T + T), id TT + *) \Rightarrow \\
 &\Rightarrow (id * (F + T), id FT + *) \Rightarrow (id * (id + T), id idT + *) \Rightarrow \\
 &\Rightarrow (id * (id + F), id idF + *) \Rightarrow (id * (id + id), id idid + *),
 \end{aligned}$$

переводящая эту цепочку в цепочку $id id id + *$.

Рассмотрим связь между переводами, определяемыми СУ-схемами и осуществляемыми МП-преобразователями [2].

Теорема 5.1. Пусть P — МП-преобразователь. Тогда существует такая простая СУ-схема Tr , что $\tau(Tr) = \tau(P)$.

Теорема 5.2. Пусть Tr — простая СУ-схема. Тогда существует такой МП-преобразователь P , что $\tau(P) = \tau(Tr)$.

Таким образом, класс переводов, определяемых магазинными преобразователями, совпадает с классом простых СУ-переводов.

Рассмотрим теперь связь между СУ-переводами и детерминированными МП-преобразователями, выполняющими нисходящий или восходящий разбор [2].

Теорема 5.3. Пусть $Tr = (N, T, \Pi, R, S)$ — простая СУ-схема, входной грамматикой которой служит $LL(1)$ -грамматика. Тогда перевод $\{(x, y) | (x, y) \in \tau(Tr)\}$ можно осуществить детерминированным МП-преобразователем.

Существуют простые СУ-схемы, имеющие в качестве входных грамматик $LR(1)$ -грамматики и не реализуемые ни на каком ДМП-преобразователе.

Пример 5.3. Рассмотрим простую СУ-схему с правилами:

$$\begin{aligned}
 S &\rightarrow Sa, \quad aSa; \\
 S &\rightarrow Sb, \quad bSb; \\
 S &\rightarrow e, \quad e.
 \end{aligned}$$

Входная грамматика является $LR(1)$ -грамматикой, но не существует ДМП-преобразователя, определяющего перевод $\{(x, y) | (x, y) \in \tau(Tr)\}$.

Назовем СУ-схему $Tr = (N, T, \Pi, R, S)$ *постфиксной*, если каждое правило из R имеет вид $A \rightarrow u, v$, где $v \in N^*\Pi^*$. Иными словами, каждый элемент перевода представляет собой цепочку из нетерминалов, за которыми следует цепочка выходных символов.

Теорема 5.4. Пусть Tr — простая постфиксная СУ-схема, входная грамматика для которой является $LR(1)$. Тогда перевод $\{(x, y) | (x, y) \in \tau(Tr)\}$ можно осуществить детерминированным МП-преобразователем.

5.2.2. Обобщенные схемы синтаксически управляемого перевода.

Расширим определение СУ-схемы, с тем чтобы выполнять более широкий класс переводов. Во-первых, позволим иметь в каждой вершине дерева разбора несколько переводов. Как и в обычной СУ-схеме, каждый перевод зависит от прямых потомков соответствующей вершины дерева. Во-вторых, позволим элементам перевода быть произвольными цепочками выходных символов и символов, представляющих переводы в потомках. Таким образом, символы перевода могут повторяться или вообще отсутствовать.

Определение 5.3. *Обобщенной схемой синтаксически управляемого перевода* (или трансляции, сокращенно: ОСУ-схемой) называется шестерка $Tr = (N, T, \Pi, \Gamma, R, S)$, где все символы имеют тот же смысл, что и для СУ-схемы, за исключением того, что

- 1) Γ — конечное множество *символов перевода* вида A_i , где $A \in N$ и i — целое число;
- 2) R — конечное множество *правил перевода* вида $A \rightarrow u, A_1 = v_1, \dots, A_m = v_m$, удовлетворяющих следующим условиям:
 - а) $A_j \in \Gamma$ для $1 \leq j \leq m$,
 - б) каждый символ, входящий в v_1, \dots, v_m , либо принадлежит Π , либо является $B_k \in \Gamma$, где B входит в u ,
 - в) если u имеет более одного вхождения символа B , то каждый символ B_k во всех v соотнесен (верхним индексом) с конкретным вхождением B .

Правило $A \rightarrow u$ называют *входным правилом* вывода, A_i — *переводом* нетерминала A , $A_i = v_i$ — *элементом* перевода, связанным с этим правилом перевода. Если в ОСУ-схеме нет двух правил перевода с одинаковым входным правилом вывода, то ее называют *семантически однозначной*.

Выход ОСУ-схемы определим снизу вверх. С каждой внутренней вершиной n дерева разбора (во входной грамматике), помеченной A , свяжем одну цепочку для каждого A_i . Эта цепочка называется *значением* (или переводом) символа A_i в вершине n . Каждое значение вычисляется подстановкой значений символов перевода данного элемента перевода $A_i = v_i$, определенных в прямых потомках вершины n .

Переводом $\tau(Tr)$, *определяемым ОСУ-схемой* Tr , назовем множество $\{(x, y) \mid x \text{ имеет дерево разбора во входной грамматике для } Tr \text{ и } y \text{ — значение выделенного символа перевода } S_k \text{ в корне этого дерева}\}$.

Пример 5.4. Рассмотрим формальное дифференцирование выражений, включающих константы 0 и 1, переменную x , функции \sin и \cos , а также операции $*$ и $+$. Такие выражения порождает грамматика

$$\begin{aligned}
 E &\rightarrow E + T \mid T; \\
 T &\rightarrow T * F \mid F; \\
 F &\rightarrow (E) \mid \sin(E) \mid \cos(E) \mid x \mid 0 \mid 1.
 \end{aligned}$$

Свяжем с каждым из E , T и F два перевода, обозначенных индексами 1 и 2. Индекс 1 указывает на то, что выражение не дифференцировано, 2 — что выражение продифференцировано. Формальная производная — это E_2 . Законы дифференцирования таковы:

$$\begin{aligned}
 d(f(x) + g(x)) &= df(x) + dg(x); \\
 d(f(x) * g(x)) &= f(x) * dg(x) + g(x) * df(x); \\
 d \sin(f(x)) &= \cos(f(x)) * df(x); \\
 d \cos(f(x)) &= -\sin(f(x))df(x); \\
 dx &= 1; \\
 d0 &= 0; \\
 d1 &= 0.
 \end{aligned}$$

Эти законы можно реализовать следующей ОСУ-схемой:

$$\begin{aligned}
 E &\rightarrow E + T, & E_1 &= E_1 + T_1; \\
 & & E_2 &= E_2 + T_2;
 \end{aligned}$$

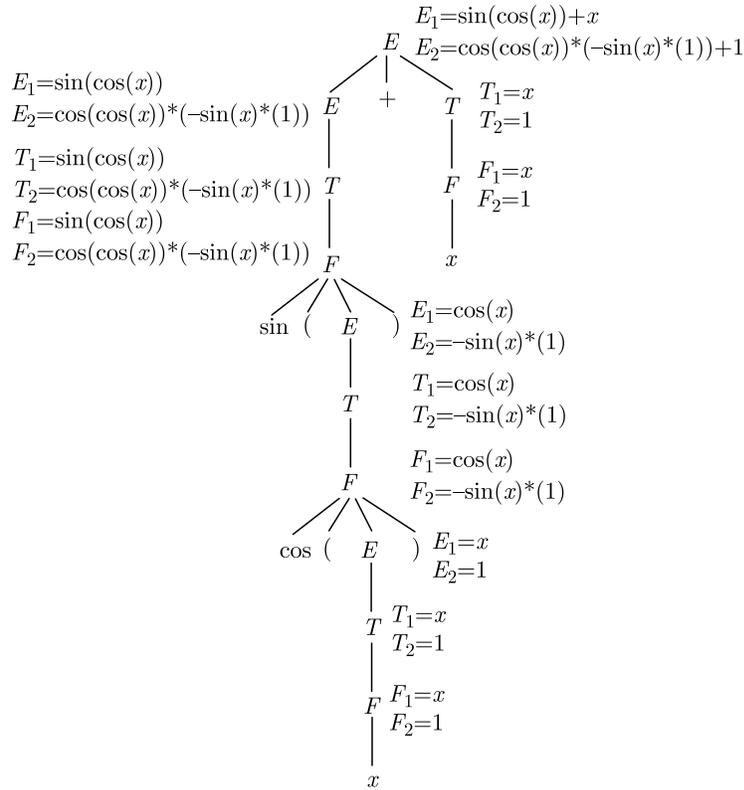


Рис. 5.1

$$\begin{aligned}
E &\rightarrow T, & E_1 &= T_1; \\
&& E_2 &= T_2; \\
T &\rightarrow T * F, & T_1 &= T_1 * F_1; \\
&& T_2 &= T_1 * F_2 + T_2 * F_1; \\
T &\rightarrow F, & T_1 &= F_1; \\
&& T_2 &= F_2; \\
F &\rightarrow (E), & F_1 &= (E_1); \\
&& F_2 &= (E_2); \\
F &\rightarrow \sin(E), & F_1 &= \sin(E_1); \\
&& F_2 &= \cos(E_1) * (E_2); \\
F &\rightarrow \cos(E), & F_1 &= \cos(E_1); \\
&& F_2 &= -\sin(E_1) * (E_2); \\
F &\rightarrow x, & F_1 &= x; \\
&& F_2 &= 1; \\
F &\rightarrow 0, & F_1 &= 0; \\
&& F_2 &= 0; \\
F &\rightarrow 1, & F_1 &= 1; \\
&& F_2 &= 0;
\end{aligned}$$

Дерево вывода для $\sin(\cos(x)) + x$ приведено на рис. 5.1.

5.3. Атрибутные грамматики

Среди всех формальных методов описания языков программирования атрибутные грамматики (введенные Кнудом [7]), по-видимому, наиболее известны и распространены. Причиной этого является то, что формализм атрибутных грамматик основывается на дереве разбора программы в КС-грамматике, что сближает его с хорошо разработанной теорией и практикой построения трансляторов.

5.3.1. Определение атрибутных грамматик. *Атрибутивной грамматикой* называется четверка $AG = (G, A_S, A_I, R)$, где:

- 1) $G = (N, T, P, S)$ — приведенная КС-грамматика;
- 2) A_S — конечное множество *синтезируемых атрибутов*;
- 3) A_I — конечное множество *наследуемых атрибутов*, $A_S \cap A_I = \emptyset$;
- 4) R — конечное множество *семантических правил*.

Атрибутивная грамматика AG сопоставляет каждому символу X из $N \cup T$ множество $A_S(X)$ синтезируемых атрибутов и множество $A_I(X)$ наследуемых атрибутов. Множество всех синтезируемых атрибутов всех символов из $N \cup T$ обозначается A_S , наследуемых — A_I . Атрибуты разных символов являются различными атрибутами. Будем обозначать атрибут a символа X как $a(X)$. Значения атрибутов могут быть произвольных типов, например, представлять собой числа, строки, адреса памяти и т. п.

Пусть правило p из P имеет вид $X_0 \rightarrow X_1 X_2 \dots X_n$. Атрибутная грамматика AG сопоставляет каждому правилу p из P конечное множество $R(p)$ семантических правил вида

$$a(X_i) = f(b(X_j), c(X_k), \dots, d(X_m)),$$

где $0 \leq j, k, \dots, m \leq n$, причем $1 \leq i \leq n$, если $a(X_i) \in A_I(X_i)$ (т. е. $a(X_i)$ — наследуемый атрибут), и $i = 0$, если $a(X_i) \in A_S(X_i)$ (т. е. $a(X_i)$ — синтезируемый атрибут).

Таким образом, семантическое правило определяет значение атрибута a символа X_i на основе значений атрибутов b, c, \dots, d символов X_j, X_k, \dots, X_m соответственно.

В частном случае длина n правой части правила может быть равна нулю; тогда будем говорить, что атрибут a символа X_i получает в качестве значения константу.

В дальнейшем будем считать, что атрибутная грамматика не содержит семантических правил для вычисления атрибутов терминальных символов. Предполагается, что атрибуты терминальных символов — либо предопределенные константы, либо доступны как результат работы лексического анализатора.

Пример 5.5. Рассмотрим атрибутную грамматику, позволяющую вычислить значение вещественного числа, представленного в десятичной записи. Здесь $N = \{Num, Int, Frac\}$, $T = \{digit, .\}$, $S = Num$, а правила вывода и семантические правила определяются следующим образом (верхние индексы используются для ссылки на разные вхождения одного и того же нетерминала):

$$Num \rightarrow Int . Frac \quad v(Num) = v(Int) + v(Frac); \\ p(Frac) = 1;$$

$$Int \rightarrow e \quad v(Int) = 0; \\ p(Int) = 0;$$

$$Int^1 \rightarrow digit Int^2 \quad v(Int^1) = v(digit) * 10^{p(Int^2)} + v(Int^2); \\ p(Int^1) = p(Int^2) + 1;$$

$$Frac \rightarrow e \quad v(Frac) = 0;$$

$$Frac^1 \rightarrow digit Frac^2 \quad v(Frac^1) = v(digit) * 10^{-p(Frac^1)} + v(Frac^2); \\ p(Frac^2) = p(Frac^1) + 1;$$

Для этой грамматики

$$A_S(Num) = \{v\}, \quad A_I(Num) = \emptyset; \\ A_S(Int) = \{v, p\}, \quad A_I(Int) = \emptyset; \\ A_S(Frac) = \{v\}, \quad A_I(Frac) = \{p\}.$$

Пусть даны атрибутная грамматика AG и цепочка, принадлежащая языку, определяемому соответствующей $G = (N, T, P, S)$. Сопоставим этой цепочке «значение» следующим образом. Построим дерево разбора T этой цепочки в грамматике G . Каждый внутренний узел этого дерева помечается нетерминалом X_0 , соответствующим применению p -го правила грамматики; таким образом, у этого узла будет n непосредственных потомков (рис. 5.2).

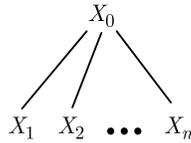


Рис. 5.2

Пусть теперь X — метка некоторого узла дерева и пусть a — атрибут символа X . Если a — синтезируемый атрибут, то $X = X_0$ для некоторого $p \in P$; если же a — наследуемый атрибут, то $X = X_j$ для некоторых $p \in P$ и $1 \leq j \leq n$. В обоих случаях дерево вблизи этого узла имеет вид, приведенный на рис. 5.2. По определению, атрибут a имеет в этом узле значение v , если в соответствующем семантическом правиле

$$a(X_i) = f(b(X_j), c(X_k), \dots, d(X_m))$$

все атрибуты b, c, \dots, d уже определены и имеют в узлах с метками X_j, X_k, \dots, X_m значения v_j, v_k, \dots, v_m соответственно, а $v = f(v_1, v_2, \dots, v_m)$. Процесс вычисления атрибутов на дереве продолжается до тех пор, пока можно вычислить хотя бы один атрибут. Вычисленные в результате атрибуты корня дерева представляют собой «значение», соответствующее данному дереву вывода.

Заметим, что значение синтезируемого атрибута символа в узле синтаксического дерева вычисляется по атрибутам символов в потомках этого узла; значение наследуемого атрибута вычисляется по атрибутам «родителя» и «соседей».

Атрибуты, сопоставленные вхождением символов в дерево разбора, будем называть *вхождениями* атрибутов в дерево разбора, а дерево с сопоставленными каждой вершине атрибутами — *атрибутированным деревом разбора*.

Пример 5.6. Атрибутированное дерево для грамматики из предыдущего примера и цепочки $w = 12,34$ показано на рис. 5.3.

Будем говорить, что семантические правила *заданы корректно*, если они позволяют вычислить все атрибуты произвольного узла в любом дереве вывода.

Между вхождениями атрибутов в дерево разбора существуют зависимости, определяемые семантическими правилами, соответствующими

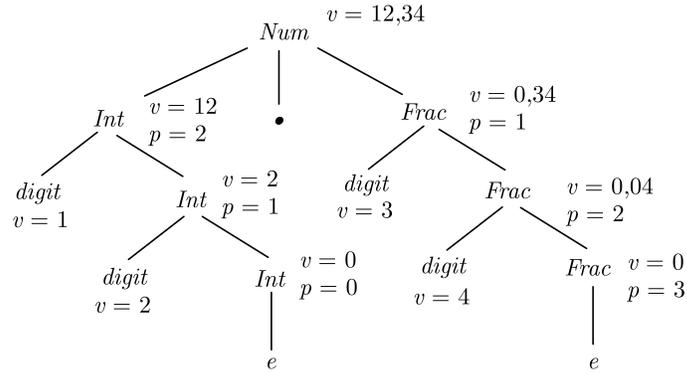


Рис. 5.3

примененным синтаксическим правилам. Эти зависимости могут быть представлены в виде ориентированного графа следующим образом.

Пусть T — дерево разбора. Сопоставим этому дереву ориентированный граф $D(T)$, узлами которого являются пары (n, a) , где n — узел дерева T , a — атрибут символа, служащего меткой узла n . Граф содержит дугу из (n_1, a_1) в (n_2, a_2) тогда и только тогда, когда семантическое правило, вычисляющее атрибут a_2 , непосредственно использует значение атрибута a_1 . Таким образом, узлами графа $D(T)$ являются атрибуты, которые нужно вычислить, а дуги определяют зависимости, подразумевающие, какие атрибуты вычисляются раньше, а какие позже.

Пример 5.7. Граф зависимостей атрибутов для дерева разбора из предыдущего примера показан на рис. 5.4.

Можно показать, что семантические правила корректны тогда и только тогда, когда для любого дерева вывода T соответствующий граф $D(T)$ не содержит циклов (т. е. является ориентированным ациклическим графом).

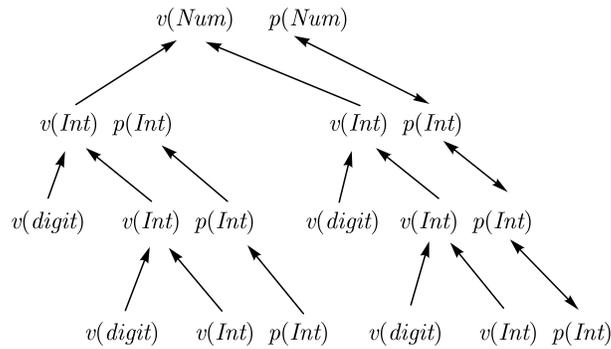


Рис. 5.4

5.3.2. Классы атрибутных грамматик и их реализация. В общем виде реализация вычислителей для атрибутных грамматик вызывает значительные трудности. Это связано с тем, что множество значений атрибутов, связанных с данным деревом, приходится вычислять в соответствии с зависимостями атрибутов, которые образуют ориентированный ациклический граф. На практике стараются осуществлять процесс вычисления атрибутов, привязывая его к тому или иному способу обхода дерева. Рассматривают многовизитные, многопроходные и другие атрибутные вычислители. Это, как правило, ведет к ограничению допустимых зависимостей между атрибутами, поддерживаемых вычислителем.

Простейшими подклассами атрибутных грамматик, для которых вычисление всех атрибутов может быть осуществлено одновременно с синтаксическим анализом, являются S-атрибутные и L-атрибутные грамматики.

Определение 5.4. Атрибутная грамматика называется *S-атрибутной*, если она содержит только синтезируемые атрибуты.

Нетрудно видеть, что для S-атрибутной грамматики на любом дереве разбора все атрибуты могут быть вычислены за один обход дерева снизу-вверх. Таким образом, можно вычислять атрибуты параллельно с восходящим синтаксическим анализом, например, LR(1)-анализом.

Пример 5.8. Рассмотрим S-атрибутную грамматику для перевода арифметических выражений в ПОЛИЗ. Здесь атрибут v имеет строковый тип, \parallel обозначает операцию конкатенации. Правила вывода и семантические правила определяются следующим образом:

$$E^1 \rightarrow E^2 + T, \quad v(E^1) = v(E^2) \parallel v(T) \parallel ' + ';$$

$$E \rightarrow T, \quad v(E) = v(T);$$

$$T \rightarrow T * F, \quad v(T^1) = v(T^2) \parallel v(F) \parallel '*';$$

$$T \rightarrow F, \quad v(T) = v(F);$$

$$F \rightarrow id, \quad v(F) = v(id);$$

$$F \rightarrow (E), \quad v(F) = v(E).$$

Определение 5.5. Атрибутная грамматика называется *L-атрибутной*, если любой наследуемый атрибут любого символа X_j из правой части каждого правила $X_0 \rightarrow X_1 X_2 \dots X_n$ грамматики зависит только от:

- 1) атрибутов символов X_1, X_2, \dots, X_{j-1} , находящихся в правиле слева от X_j ;
- 2) наследуемых атрибутов символа X_0 .

Заметим, что каждая S-атрибутивная грамматика является L-атрибутивной. Все атрибуты на любом дереве для L-атрибутивной грамматики могут быть вычислены за один обход дерева сверху-вниз слева-направо. Таким образом, вычисление атрибутов можно осуществлять параллельно с нисходящим синтаксическим анализом, например, LL(1)-анализом или рекурсивным спуском.

В случае рекурсивного спуска в каждой функции, соответствующей нетерминалу, надо определить формальные параметры, передаваемые по значению, для наследуемых атрибутов, и формальные параметры, передаваемые по ссылке, для синтезируемых атрибутов. В качестве примера рассмотрим реализацию атрибутивной грамматики из примера 5.5 (нетрудно видеть, что грамматика является L-атрибутивной). Классы FloatRef и IntRef введены для реализации передачи параметра по ссылке. Предполагается, что в переменной InSym хранится десятичное значение выбранной цифры.

```
class FloatRef {float val;}
class IntRef {int ref;}

void int_part(FloatRef V0, IntRef P0)
{if (Map[InSym]==Digit)
  { int I=InSym;
    FloatRef V2;
    IntRef P2;
    InSym=getInSym();
    int_part(V2,P2);
    V0.val=I*exp(P2.val*ln(10))+V2.val;
    P0.val=P2.val+1;
  }
  else {V0.val=0;
        P0.val=0;
      }
}
void fract_part(FloatRef V0, IntRef P0)
{if (Map[InSym]==Digit)
  { int I=InSym;
    FloatRef V2;
    int P2=P0+1;
    InSym=getInSym();
    fract_part(V2,P2);
    V0.val=I*exp(-P0.val*ln(10))+V2.val;
  }
  else {V0.val=0;
        }
}
void number()
{ FloatRef V1,V3;
```

```
float V0;
IntRef P;
int_part(V1, P);
if (InSym!='.') error();
fract_part(V3, 1);
V0=V1.val+V3.val;
}
```

5.3.3. Язык описания атрибутных грамматик. Формализм атрибутных грамматик оказался очень удобным средством для описания семантики языков программирования. Вместе с тем выяснилось, что реализация вычислителей для атрибутных грамматик общего вида сталкивается с большими трудностями. В связи с этим было сделано множество попыток рассматривать те или иные классы атрибутных грамматик, обладающих «хорошими» свойствами. К числу таких свойств относятся прежде всего простота алгоритма проверки атрибутной грамматики на зацикленность и простота алгоритма вычисления атрибутов для атрибутных грамматик данного класса.

Атрибутные грамматики использовались для описания семантики языков программирования, и было создано несколько систем автоматизации разработки трансляторов, основанных на формализме атрибутных грамматик. Опыт их использования показал, что «чистый» атрибутный формализм может быть успешно применен для описания семантики языка, но его использование вызывает трудности при создании транслятора. Эти трудности связаны как с самим формализмом, так и с некоторыми технологическими проблемами. К трудностям первого рода можно отнести несоответствие между чисто функциональной природой атрибутного вычислителя и связанной с ней неупорядоченностью процесса вычисления атрибутов (что в значительной степени является преимуществом этого формализма), с одной стороны, и упорядоченностью элементов программы — с другой. Это несоответствие ведет к тому, что приходится идти на искусственные приемы для их сочетания. Технологические трудности связаны с эффективностью трансляторов, полученных с помощью атрибутных систем. Как правило, качество таких трансляторов довольно низко из-за больших расходов памяти и неэффективности вышеупомянутых искусственных приемов.

Учитывая это, мы будем вести дальнейшее изложение на языке, сочетающем особенности атрибутного формализма и обычного языка программирования, в котором предполагается наличие операторов, а значит, и возможность управления порядком их исполнения. Этот порядок может быть привязан к обходу атрибутированного дерева разбора сверху-вниз слева-направо. Что касается грамматики входного языка, то мы не будем предполагать принадлежность ее определенному классу (например, LL(1) или LR(1)). Будем считать, что дерево разбора входной программы уже построено как результат

синтаксического анализа и что атрибутные вычисления осуществляются в результате обхода этого дерева. Таким образом, входная грамматика атрибутного вычислителя может быть даже неоднозначной, что не влияет на процесс атрибутных вычислений.

При записи синтаксиса мы будем использовать расширенную БНФ. Элемент правой части синтаксического правила, заключенный в скобки [], может отсутствовать. Элемент правой части синтаксического правила, заключенный в скобки (), означает возможность повторения один или более раз. Элемент правой части синтаксического правила, заключенный в скобки [()], означает возможность повторения ноль или более раз. В скобках [] или [()] может указываться разделитель конструкций.

Ниже дан синтаксис языка описания атрибутных грамматик. Приведен только синтаксис конструкций, собственно описывающих атрибутные вычисления. Синтаксис обычных выражений и операторов не приводится — он основывается на Си.

```

Атрибутная грамматика ::= 'АЛФАВЕТ'
    ( ОписаниеНетерминала ) ( Правило )
ОписаниеНетерминала ::= ИмяНетерминала
    ':' ':' [ ( ОписаниеАтрибутов / ';' ) ] '.'
ОписаниеАтрибутов ::= Тип ( ИмяАтрибута / ',' )
Правило ::= 'RULE' Синтаксис 'SEMANTICS' Семантика '.'
Синтаксис ::= ИмяНетерминала ':' ':' ПраваяЧасть
ПраваяЧасть ::= [ ( ЭлементПравойЧасти ) ]
ЭлементПравойЧасти ::= ИмяНетерминала
    | Терминал
    | '(' Нетерминал [ '/' Терминал ] ')'
    | '[' Нетерминал '['
    | '[' ( Нетерминал [ '/' Терминал ] ')' ]
Семантика ::= [ ( ЛокальноеОбъявление / ';' ) ]
    [ ( СемантическоеДействие / ';' ) ]
СемантическоеДействие ::= Присваивание
    | [ Метка ] Оператор
Присваивание ::= Переменная ':' '=' Выражение
Переменная ::= ЛокальнаяПеременная
    | Атрибут
Атрибут ::= ЛокальныйАтрибут
    | ГлобальныйАтрибут
ЛокальныйАтрибут ::= ИмяАтрибута '<' Номер '>'
ГлобальныйАтрибут ::= ИмяАтрибута '<' Нетерминал '>'
Метка ::= Целое ':' ':'
    | Целое 'E' ':' ':'
    | Целое 'A' ':' ':'

```

Описание атрибутной грамматики состоит из раздела описания атрибутов и раздела правил. Раздел описания атрибутов определяет состав атрибутов для каждого символа грамматики и тип каждого атрибута. Правила состоят из синтаксической и семантической частей. В синтаксической части используется расширенная БНФ. Семантическая часть правила состоит из локальных объявлений и семантических действий. В качестве семантических действий допускаются как атрибутные присваивания, так и составные операторы.

Метка в семантической части правила привязывает выполнение оператора к обходу дерева разбора сверху-вниз слева-направо. Конструкция `i : оператор` означает, что оператор должен быть выполнен сразу после обхода i -й компоненты правой части. Конструкция `i E : оператор` означает, что оператор должен быть выполнен, только если порождение i -й компоненты правой части пусто. Конструкция `i A : оператор` означает, что оператор должен быть выполнен после разбора каждого повторения i -й компоненты правой части (имеется в виду конструкция повторения).

Каждое правило может иметь локальные определения (типов и переменных). В формулах используются как атрибуты символов данного правила (*локальные атрибуты*), и в этом случае соответствующие символы указываются номерами в правиле (0 для символа левой части, 1 для первого символа правой части, 2 для второго символа правой части, и т. д.), так и атрибуты символов предков левой части правила (*глобальные атрибуты*). В этом случае соответствующий символ указывается именем нетерминала. Таким образом, на дереве образуются *области видимости* атрибутов: атрибут символа имеет область видимости, состоящую из правила, в которое символ входит в правую часть, плюс всё поддереву, корнем которого является символ, за исключением поддеревьев — потомков того же символа в этом поддереве.

Значение терминального символа доступно через атрибут VAL соответствующего типа.

Пример 5.9. Атрибутная грамматика из примера 5.5 записывается следующим образом:

```
ALPHABET Num      :: float V. Int      :: float V;
           int P.
Frac    :: float V;
           int P.
digit   :: int VAL.
```

RULE

```
Num ::= Int '.' Frac
```

SEMANTICS

```
V<0>=V<1>+V<3>; P<3>=1.
```

RULE

Int ::= e
SEMANTICS
 $V\langle 0 \rangle = 0; P\langle 0 \rangle = 0.$

RULE
Int ::= digit Int
SEMANTICS
 $V\langle 0 \rangle = VAL\langle 1 \rangle * 10^{**}P\langle 2 \rangle + V\langle 2 \rangle; P\langle 0 \rangle = P\langle 2 \rangle + 1.$

RULE
Frac ::= e
SEMANTICS
 $V\langle 0 \rangle = 0.$

RULE
Frac ::= digit Frac
SEMANTICS
 $V\langle 0 \rangle = VAL\langle 1 \rangle * 10^{**}(-P\langle 0 \rangle) + V\langle 2 \rangle; P\langle 2 \rangle = P\langle 0 \rangle + 1.$

Глава 6

ПРОВЕРКА КОНТЕКСТНЫХ УСЛОВИЙ

6.1. Описание областей видимости и блочной структуры

Задачей контекстного анализа является установление свойств объектов и их использования. Наиболее часто решаемой задачей является определение существования объекта и соответствия его использования контексту, что осуществляется с помощью анализа типа объекта. Под контекстом здесь понимается вся совокупность свойств текущей точки программы, например множество доступных объектов, тип выражения и т. д.

Таким образом, необходимо хранить объекты и их типы, уметь находить эти объекты и определять их типы, определять характеристики контекста. Совокупность доступных в данной точке объектов будем называть *средой*. Обычно среда программы состоит из частично упорядоченного набора компонент

$$E = \{DS_1, DS_2, \dots, DS_n\}$$

Каждая компонента — это множество объявлений, представляющих собой пары (имя, тип):

$$DS_i = \{(\text{имя}_j, \text{тип}_j) \mid 1 \leq j \leq k_i\}$$

где под типом будем подразумевать полное описание свойств объекта (объектом, в частности, может быть само описание типа).

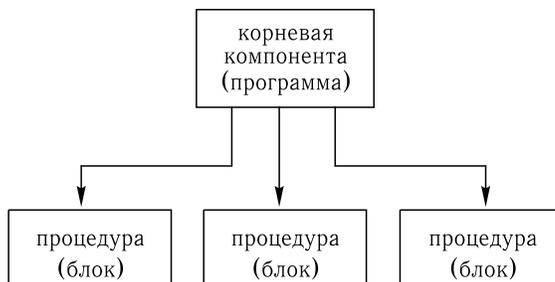


Рис. 6.1

Компоненты образуют дерево, соответствующее этому частичному порядку. Частичный порядок между компонентами обычно определяется статической вложенностью компонент в программе. Эта вложенность может соответствовать блокам, процедурам или классам программы (рис. 6.1). Компоненты среды могут быть именованы. Поиск в среде обычно ведется с учетом упорядоченности компонент. Среда может включать в себя как компоненты, полученные при трансляции «текущего» текста программы, так и «внешние» (например, раздельно компилированные) компоненты.

Для обозначения участков программы, в которых доступны те или иные описания, используются понятия *области действия* и *области видимости*. Областью действия описания является процедура (блок), содержащая описание, со всеми входящими в нее (подчиненными по дереву) процедурами (блоками). Областью видимости описания называется часть области действия, из которой исключены те подобласти, в которых по тем или иным причинам описание недоступно, например, оно перекрыто другим описанием. В разных языках понятия области действия и области видимости уточняются по-разному.

Обычными операциями при работе со средой являются:

- включить объект в компоненту среды;
- найти объект в среде и получить доступ к его описанию;
- образовать в среде новую компоненту, определенным образом связанную с остальными;
- удалить компоненту из среды.

6.2. Занесение в среду и поиск объектов

Рассмотрим схему реализации простой блочной структуры, аналогичной процедурам в Паскале или блокам в Си. Каждый блок может иметь свой набор описаний. Программа состоит из основного именованного блока, в котором имеются описания и операторы. Описания состоят из описаний типов и объявлений переменных. В качестве типа может использоваться целочисленный тип и тип массива. Два типа T1 и T2 считаются эквивалентными, если имеется описание T1=T2 (или T2=T1). Операторами служат операторы присваивания вида *Переменная1=Переменная2* и блоки. Переменная — это либо просто идентификатор, либо выборка из массива. Оператор присваивания считается правильным, если типы переменных левой и правой части эквивалентны. Примером правильной программы может служить

```
program Example;  
  type T1=array 100 of array 200 of integer;  
       T2=T1;
```

```

var   V1:T1;
      V2:T2;
begin
  V1=V2;
  V2[1]=V1[2];

      type  T3=array 300 of T1;
      var   V3:T3;
      begin
        V3[50]=V1;
      end
end

```

Рассматриваемое подмножество языка может быть порождено следующей грамматикой (запись в расширенной БНФ):

```

Prog ::= 'program' Ident ';' Block 'END'
Block ::= [ ( Declaration ) ] 'begin' [ ( Statement ) ] 'end'
Declaration ::= 'type' ( Type_Decl )
Declaration ::= 'var' ( Var_Decl )
Type_Decl ::= ( Ident '=' Type_Defin ';' )
Type_Defin ::= 'ARRAY' Integer 'OF' Type_Defin
Type_Defin ::= Type_Use
Var_Decl ::= ( Ident_List ':' Type_Use ';' )
Type_Use ::= Ident
Ident_List ::= ( Ident / ',' )
Statement ::= Block ';'
Statement ::= Variable '=' Expression ';'
Variable ::= Ident Access
Access ::= '[' Expression ']' Access
Access ::=
Expression ::= ( Term / '+' )
Term ::= ( Factor / '*' )
Factor ::= Variable
Factor ::= '(' Expression ')'
Factor ::= Integer

```

Наш транслятор будет переводить приведенную выше программу в следующую программу на Java:

```

class Example

{ public static void main(String args[])

{ int XXXX_0, XXXX_1, XXXX_2, XXXX_3, XXXX_4, XXXX_6;

    int V1_1[][]=new int[100][200];

    int V2_1[][]=new int[100][200];

```

```

for (XXXX_0=0;XXXX_0<100;XXXX_0++)
    for (XXXX_1=0;XXXX_1<200;XXXX_1++)
        V1_1[XXXX_0][XXXX_1]=V2_1[XXXX_0][XXXX_1];

for (XXXX_1=0;XXXX_1<200;XXXX_1++)
    V2_1[1][XXXX_1]=V1_1[2][XXXX_1];

int V3_2[][][]=new int[300][100][200];

for (XXXX_0=0;XXXX_0<100;XXXX_0++)
    for (XXXX_1=0;XXXX_1<200;XXXX_1++)
        V3_2[50][XXXX_0][XXXX_1]=V1_1[XXXX_0][XXXX_1];

}}

```

Среда состоит из отдельных объектов, реализуемых как записи (в дальнейшем описании мы будем использовать имя `TElement` для имени типа этой записи). Состав полей записи, вообще говоря, зависит от описываемого объекта (тип, переменная и т. п.). В трансляторе, приведенном ниже, элемент среды реализуется классом `TElement`:

```

class TElement {
    int object;           // Переменная или тип
    String objectName;   // Только для отладки
    TElement typeRef = null; // Ссылка на описатель базового типа
    ArrayList IndexList = null; // Список индексов
    int level;           // Block level
}

```

Для реализации синтезируемых атрибутов будем использовать ссылку на этот тип:

```

class TElementRef {
    TElement TER;
}

```

Для реализации некоторых атрибутов (в частности среды, списка идентификаторов и т. п.) в качестве типов данных мы будем использовать различные множества. Множество может быть упорядоченным или неупорядоченным, ключевым или простым. Элементом ключевого множества может быть запись, одним из полей которой является ключ. В Java эти множества реализуются типами данных `LinkedList`, `ArrayList`, `HashMap`.

`HashSet` — простое неупорядоченное множество объектов;

`HashMap` — ключевое неупорядоченное множество объектов ;

`LinkedList` — упорядоченное множество объектов ;

`ArrayList` — список (упорядоченное множество объектов), к элементам которого имеется доступ по индексу.

Над объектами типа множества определены следующие операции:

`T S = new T()` — создать переменную `S` типа `T`;

`S.add(Value)` — включить объект `Value` во множество `S`; если множество упорядоченное, то включение осуществляется в качестве последнего элемента;

`S.put(Name, Value)` — включить объект `Value` в ключевое множество `S`;

`(T)S.get(Key)` — выдать указатель на объект типа `T` с ключом `Key` во множестве `S` и `NULL`, если объект с таким ключом не найден.

Цикл по элементам множества реализуется с помощью итератора:

```
Iterator itr = S.iterator();

while (itr.hasNext()) {T V = (T) itr.next();}
```

Переменная `V` пробегает все значения множества. Если множество упорядочено, то элементы пробегаются в этом порядке, если нет — в произвольном порядке.

Среда представляет собой ключевое множество с ключом — именем объекта. Для реализации среды каждый нетерминал `Block` имеет атрибут `Environ`. Для обеспечения возможности просматривать компоненты среды в соответствии с вложенностью блоков поддерживается стек (список) сред блоков `EnvironStack` (рис. 6.2). Кроме того, среда блока корня дерева (нетерминал `Program`) содержит все предопределенные описания.

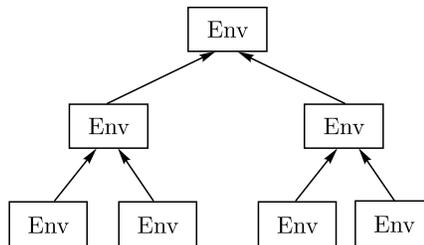


Рис. 6.2

Ниже приведена атрибутивная грамматика языка `ToyLang`, в которой определяются контекстные условия языка. Для каждого объявления переменной и объявления типа создается запись типа `TElement`. В ней указывается имя созданного объекта, его вид (переменная или тип), список индексов и указатель на базовый тип.

ALPHABET

Prog ::= LinkedList EnvironStack.

Ident ::= String Val.

Block ::= HashMap Environ. Ident_List ::= LinkedList IdentSet.

Access ::= TelementRef TypeRefRef; int IndexNo.

Type_Defin, Type_Use, Variable, Expression, Term, Factor ::= TelementRef TypeRefRef.

Declaration, type_Decl, Var_Decl, Statement::.

Rule

Prog ::= 'program' Ident ';' Block 'END'.

Rule

Block ::= [(Declaration)] 'begin' [(Statement)] 'end'

Semantics

```
0: { Environ<0> = new HashMap(); //Новая компонента среды
    if ( EnvironStack<Program>.isEmpty() ) //Блок программы
        Environ.put(intName, intRef); //Предопределенный тип integer
        EnvironStack<Program>.addFirst(Environ);
        //Создать среду нового блока
    }
{ EnvironStack.removeFirst();
  //Удалить среду блока на выходе }.
```

Rule

Declaration ::= 'type' (Type_Decl).

Rule

Declaration ::= 'var' (Var_Decl).

Rule

Type_Decl ::= Ident '=' Type_Defin ';'.

Semantics

```
1: {
if (!((HashMap)EnvironStack<Program>.getFirst()).
    containsKey(Val<1>)) {
    //Если не было объявления с таким идентификатором
    TypeRef<3> = new TELEMENT(ТИПЕОБЪЕКТ);
    //Новый ТИПЕОБЪЕКТ
    TypeRef<3>.objectName = Val<1>;
} else Error("Повторное объявление " + Val<1>); }
3: { ((HashMap)EnvironStack<Program>.getFirst()).put(Val<1>,
    typeRef<3>);
    // typeRef<3> наследуемый атрибут - ссылка на ТИПЕОБЪЕКТ
```

```

        //Поместить в среду текущего блока
    }.

Rule
Type_Defin ::= 'ARRAY' Integer 'OF' Type_Defin
Semantics
2:{if (TypeRef<0>.IndexList == null)
//Первый элемент списка индексов
    TypeRef<0>.IndexList = new ArrayList();
    TypeRef<0>.IndexList.add(new Integer(Val<2>));
    //Добавить новый индекс к списку индексов
    TypeRef<4> = TypeRef<0>;}.

Rule
Type_Defin ::= Type_Use
Semantics
0:{TypeRefRef<1>=newTElementRef();}
    //Синтезируемый атрибут
    { TypeRef<0> = TypeRefRef<1>.TER;}.

Rule
Type_Use ::= Ident
Semantics
{ HashMap Environ = null;
Iterator itr = EnvironStack<Program>.iterator();
    while (itr.hasNext()) { //Пройти по списку сред транслируемых
        //блоков от внутреннего к внешнему
        Environ = (HashMap) itr.next();
        if (Environ.containsKey(Val<1>)) {
            PV = (TElement) Environ.get(Val<1>);
            if (PV.object != TYPEОБЪЕКТ) {
                Error(Val<1>+ " Не идентификатор типа");
                PV = null;
                break;
            }
        }
    }
    TypeRefRef<0>.TER = PV;
    if (PV == null) Error(Val<1>+ " Не объявлен");
}.

Rule
Var_Decl ::= Ident_List ':' Type_Use ';'
Semantics
1:{TypeRefRef<3> = new TElementRef();
    //Синтезируемый атрибут
    IdentSet<1> = new LinkedList();
    //Список идентификаторов }
3:{TElement PV = TypeRefRef<3>.TER;
    Iterator itr = IdentSet<1>.iterator();

```

```

HashMap Env = (HashMap) EnvironStack.getFirst();
while (itr.hasNext()) { //Для каждого идентификатора из списка
    String Ident = (String) itr.next();
    if (!Env.containsKey(Ident)) { //Если идентификатор не объявлен
        TElement newVar = new TElement(VAROBJECT);
        //Создать новый VAROBJECT
        newVar.objectName = Ident;
        newVar.typeRef = PV;
        Env.put(Ident, newVar);
    } else Error("Повторное объявление" + Ident);
}
}

Rule
Ident_List ::= ( Ident / ',' )
Semantics
1A:{IdentSet<0>.addFirst(new String(Val<1>));}.

Rule
Statement ::= Block ';''.

Rule
Statement ::= Variable '=' Expression';'
Semantics
{
TElementRef TypeRefRef<1 >= new TElementRef(); TElementRef
TypeRefRef<3> = new TElementRef();
if (TypeRefRef<1> .TER != TypeRefRef<3> .TER)
    Error("Типы значений различны"); }.

Rule
Variable ::= Ident Access
Semantics
1:{ Iterator
itr=EnvironStack<Program>.iterator();
while (itr.hasNext()) {
    //Пройти по блокам от внутреннего к внешнему
    HashMap Environ = (HashMap) itr.next();
    if (Environ.containsKey(Val<1>)) {
        TElement PV = (TElement) Environ.get(Val<1>);
        if (PV.object == VAROBJECT)
            /* Получаем ссылку на TElement тип переменной */
            typeRefRef<0>.TER = PV.typeRef;
        else {
            Error(Val<1> + " Не идентификатор переменной");
            PV = null;
            break;
        }
    }
}
}
}

```

```

    if (PV == null) Error(Val<1> + " Не объявлен");
IndexNo<2> = 0; //Подсчет числа индексов
TypeRefRef<2> = TypeRefRef<0>; }.

Rule
Access ::= '[' Expression ']' Access
Semantics
TypeRefRef <4> = TypeRefRef <0>;
{if (TypeRefRef <2> .TER != intRef)
    Error("Тип переменной в выражении должен быть integer");
    IndexNo<4> = IndexNo<0>+1; }.

Rule
Access ::=
Semantics
{ Telement eqType = typeRefRef<0>.TER;
while ((eqType.IndexList == null)
    /* Список эквивалентных типов */
    & (eqType.typeRef != null))
    /* Пройти по списку эквивалентных типов но не далее intRef */
    eqType = eqType.typeRef;
    TypeRefRef<0>.TER = eqType;
    if (indexNo<0> != 0) {
        typeRefRef<0>.TER = eqType.typeRef;
        // Получаем тип элемента массива
        if (indexNo<0> != eqType.IndexList.size())
            Error("Выход за пределы размерности массива");
    }
}.

Rule
Expression ::= (Term / '+')
Semantics
0:{ TypeRefRef<1>=TypeRefRef<0>;
    boolean first = true;//Признак первый ли Term
}
1A:{if (!first)
//Если Term не единственный, то все они должны иметь тип integer
if ((TypeRef != intRef) | (TypeRefRef<1>.TER != TypeRef))
// TypeRef - тип предыдущего Term
// TypeRef<1>.TER - тип текущего Term
    Error("Тип переменной в выражении должен быть integer");
    Telement TypeRef = typeRefRef<1>.TER;
    //Тип очередного Term
first = false; }.

Rule
Term ::= (Factor/'*')
Semantics
0:{ typeRefRef<1> =typeRefRef<0>;
    boolean first = true;

```

```
}
1A:{if (!first)
if ((TypeRef != intRef) | (typeRefRef<1>.TER != typeRef))
    Error("Тип переменной в выражении должен быть integer");
    Telement typeRef = typeRefRef<1>.TER;
    first = false; }.
```

```
Rule
Factor ::= Variable
Semantics
{ TypeRefRef<0> =TypeRefRef<1>; }.
```

```
Rule
Factor ::= '(' Expression ')'
Semantics
{ TypeRefRef<0> = TypeRefRef<2>;
  if (TypeRefRef<2>.TER != intRef)
    Error("Тип переменной в выражении должен быть integer");
}.
```

```
Rule
Factor ::= Integer
Semantics
{TypeRefRef<0>.TER = intRef;}
```

Глава 7

ОРГАНИЗАЦИЯ ТАБЛИЦ СИМВОЛОВ

В процессе работы компилятор хранит информацию об объектах программы в специальных таблицах символов. Как правило, информация о каждом объекте состоит из двух основных элементов: имени объекта и описания объекта. Информация об объектах программы должна быть организована таким образом, чтобы поиск ее был по возможности быстрее, а требуемая память по возможности меньше.

Кроме того, со стороны языка программирования могут быть дополнительные требования к организации информации. Имена могут иметь определенную область видимости. Например, поле записи должно быть уникально в пределах структуры (или уровня структуры), но может совпадать с именем объекта вне записи (или другого уровня записи). В то же время имя поля может открываться оператором присоединения, и тогда может возникнуть конфликт имен (или неоднозначность в трактовке имени). Если язык имеет блочную структуру, то необходимо обеспечить такой способ хранения информации, чтобы, во-первых, поддерживать блочный механизм видимости, а во-вторых, эффективно освобождать память при выходе из блока. Некоторые языки (например, Ада) допускают одновременную (в одном блоке) видимость нескольких объектов с одним именем, в других такая ситуация недопустима.

Мы рассмотрим некоторые основные способы организации таблиц символов в компиляторе: таблицы идентификаторов, таблицы расстановки, двоичные деревья и реализацию блочной структуры.

7.1. Таблицы идентификаторов

Как уже было сказано, информацию об объекте обычно можно разделить на две части: имя (идентификатор) и описание. Если длина идентификатора ограничена (или имя идентифицируется по ограниченному числу первых символов идентификатора), то таблица символов может быть организована в виде простого массива строк фиксированной длины, как это изображено на рис. 7.1. Некоторые входы могут быть заняты, некоторые — свободны.

Ясно, что, во-первых, размер массива должен быть не меньше числа идентификаторов, которые могут реально появиться в программе (в противном

Имя объекта					Описание объекта
s	o	r	t		
a					
r	e	a	d		
i					

Рис. 7.1

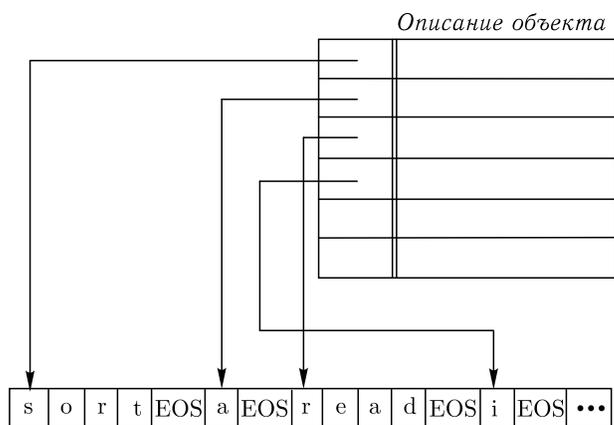


Рис. 7.2

случае возникает переполнение таблицы); во-вторых, как правило, потенциальное число различных идентификаторов существенно больше размера таблицы.

Заметим, что в большинстве языков программирования символьное представление идентификатора может иметь произвольную длину. Кроме того, различные объекты в одной или в разных областях видимости могут иметь одинаковые имена, и нет большого смысла занимать память для повторного хранения идентификатора. Таким образом, удобно имя объекта и его описание хранить по отдельности.

В этом случае идентификаторы хранятся в отдельной таблице — *таблице идентификаторов*. В таблице символов же хранится указатель на соответствующий вход в таблицу идентификаторов. Таблицу идентификаторов можно организовать, например, в виде сплошного массива. Идентификатор в массиве заканчивается каким-либо специальным символом EOS (рис. 7.2). Второй возможный вариант — в качестве первого символа идентификатора в массив заносится его длина.

7.2. Таблицы расстановки

Одним из эффективных способов организации таблицы символов является *таблица расстановки* (или *хеш-таблица*). Поиск в такой таблице может быть организован методом повторной расстановки. Суть его заключается в следующем.

Таблица символов представляет собой массив фиксированного размера N . Идентификаторы могут храниться как в самой таблице символов, так и в отдельной таблице идентификаторов.

Определим некоторую функцию h_1 (*первичную функцию расстановки*), определенную на множестве идентификаторов и принимающую значения от 0 до $N - 1$ (т.е. $0 \leq h_1(id) \leq N - 1$, где id — символьное представление идентификатора). Таким образом, функция расстановки сопоставляет идентификатору некоторый адрес в таблице символов.

Пусть мы хотим найти в таблице идентификатор id . Если элемент таблицы с номером $h_1(id)$ не заполнен, то это означает, что идентификатора в таблице нет. Если же он занят, то это еще не означает, что идентификатор id в таблицу занесен, поскольку (вообще говоря) много идентификаторов могут иметь одно и то же значение функции расстановки. Для того чтобы определить, нашли ли мы нужный идентификатор, сравниваем id с элементом таблицы $h_1(id)$. Если они равны — идентификатор найден, если нет — надо продолжать поиск дальше.

Для этого вычисляется *вторичная функция расстановки* $h_2(h_1)$ (значением которой опять-таки является некоторый адрес в таблице символов). Возможны четыре варианта:

- элемент таблицы не заполнен (т.е. идентификатора в таблице нет),
- идентификатор элемента таблицы совпадает с искомым (т.е. идентификатор найден),
- адрес элемента совпадает с уже просмотренным (т.е. таблица вся просмотрена, но идентификатора нет),
- предыдущие варианты не выполняются, так что необходимо продолжать поиск.

Для продолжения поиска применяется следующая функция расстановки: $h_3(h_2)$, $h_4(h_3)$ и т.д. Как правило, $h_i = h_2$ для $i \geq 2$. Аргументом функции h_2 является целое в диапазоне $[0, N - 1]$, и она может быть устроена по-разному. Приведем три варианта.

- 1) $h_2(i) = (i + 1) \bmod N$.

Берется следующий (циклически) элемент массива. Этот вариант плох тем, что занятые элементы «группируются», образуют последовательные

занятые участки и в пределах этого участка поиск становится по существу линейным.

2) $h_2(i) = (i + k) \bmod N$, где k и N взаимно просты.

По существу это предыдущий вариант, но элементы не накапливаются в последовательных элементах, а «разносятся».

3) $h_2(i) = (a * i + c) \bmod N$ — «псевдослучайная последовательность».

Здесь c и N должны быть взаимно просты, $b = a - 1$ кратно p для любого простого p , являющегося делителем N , b кратно 4, если N кратно 4 [6].

Поиск в таблице расстановки можно описать следующей функцией (здесь NULL некоторое выделенное значение типа int, например, -1):

```
class BoolRef {boolean val;}
class IntRef {int val;}

void Search(String Id, BoolRef Yes, IntRef Point)
{int H0=h1(Id), H=H0;
  while (1)
    {if (Empty(H)==NULL)
      {Yes.val=false;
       Point.val=H;
       return;
      }
      else if (IdComp(H, Id)==0)
        {Yes.val=true;
         Point.val=H;
         return;
        }
      else H=h2(H);
      if (H==H0)
        {Yes.val=false;
         Point.val=NULL;
         return;
        }
    }
}
```

Функция `IdComp(H, Id)` сравнивает элемент таблицы на входе H с идентификатором и вырабатывает 0, если они равны. Функция `Empty(H)` вырабатывает NULL, если вход H пуст. Функция `Search` присваивает параметрам `Yes` и `Pointer` соответственно следующие значения :

`true`, P — если найден требуемый идентификатор, где P — указатель на вход в таблице, соответствующий этому идентификатору;

`false`, NULL — если искомый идентификатор не найден, причем в таблице нет свободного места;

`false`, `P` — если искомый идентификатор не найден, но в таблице есть свободный вход `P`.

Занесение элемента в таблицу можно осуществить следующей функцией:

```
int Insert(String Id)
{BoolRef Yes;
  IntRef Point;
  Search(Id, Yes, Point);
  if (!Yes.val && (Point.val!=NULL)) InsertId(Point, Id);
  return(Point.val);
}
```

Здесь функция `InsertId(Point, Id)` заносит идентификатор `Id` для входа `Point` таблицы.

7.3. Таблицы расстановки со списками

Только что описанная схема страдает одним недостатком — возможностью переполнения таблицы. Рассмотрим ее модификацию, когда все элементы, имеющие одинаковые значения (первичной) функции расстановки, связываются в список (при этом отпадает необходимость использования функций h_i для $i \geq 2$). Таблица расстановки со списками — это массив указателей на списки элементов (рис. 7.3).

Вначале таблица расстановки пуста (все элементы имеют значение `NULL`). При поиске идентификатора `Id` вычисляется функция расстановки $h(\text{Id})$

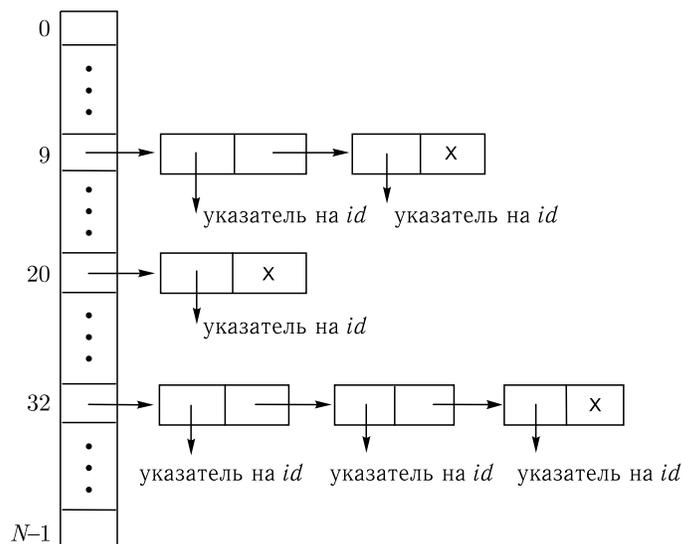


Рис. 7.3

и просматривается соответствующий линейный список. Поиск в таблице может быть описан следующей функцией:

```
class Element
  {String IdentP;
   Element Next;
  };
ArrayList T = new ArrayList(N);
Element Search(String Id)
{ Element P;
 P=(Element)T.get(h(Id));
 while (1)
 {if (P==null) return(null);
  else if (IdComp(P.IdentP, Id)==0) return(P);
  else P=P.Next;
 }
}
```

Занесение элемента в таблицу можно осуществить следующей функцией:

```
Element Insert(String Id)
{ Element P,H;
 P=Search(Id);
 if (P!=null) return(P);
 else {H=H(Id);
       P=new Element();
       P.Next=(Element)T.get(h(Id));
       P.IdentP=Include(Id);
       T.add(H,P);
     }
 return(P);
}
```

Функция Include заносит идентификатор в таблицу идентификаторов. Алгоритм иллюстрируется рис. 7.4.

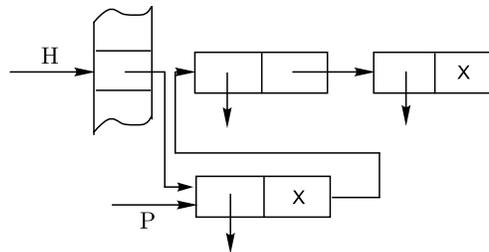


Рис. 7.4

7.4. Функции расстановки

Много внимания исследователи уделили тому, какой должна быть (первичная) функция расстановки. Основные требования к ней очевидны: она должна легко вычисляться и распределять равномерно. Один из возможных подходов здесь заключается в следующем.

1. По символам строки s определяем положительное целое H . Преобразовать одиночные символы в целые обычно можно средствами языка реализации. В Паскале для этого служит функция `ord`, в Си при выполнении арифметических операций символьные значения трактуются как целые.

2. Преобразуем H , вычисленное выше, в номер элемента, т. е. целое между 0 и $N - 1$, где N — размер таблицы расстановки, например, взятием остатка при делении H на N .

Функции расстановки, учитывающие все символы строки, распределяют лучше, чем функции, учитывающие только несколько символов, например, в конце или в середине строки. Но такие функции требуют больше вычислений.

Простейший способ вычисления H — сложение кодов символов. Перед сложением с очередным символом можно умножить старое значение H на константу q , т. е. полагаем $H_0 = 0$, $H_i = q * H_{i-1} + c_i$ для $1 \leq i \leq k$, k — длина строки. При $q = 1$ получаем простое сложение символов. Вместо сложения можно выполнять сложение c_i и $q * H_{j-1}$ по модулю 2. Переполнение при выполнении арифметических операций можно игнорировать.

Функция `Hashpjlw`, приведенная ниже [3], вычисляется, начиная с $H = 0$ (предполагается, что используются 32-битовые целые). Для каждого символа сдвигаем биты H на 4 позиции влево и добавляем очередной символ. Если какой-нибудь из четырех старших битов H равен 1, то сдвигаем эти 4 бита на 24 разряда вправо, затем складываем по модулю 2 с H и устанавливаем в 0 каждый из четырех старших битов, равных 1:

```
final int PRIME = 211;
final int EOS = '\0';
int Hashpjlw(String s)
{char p;
  int H=0, g, i;
  int sLen = s.length();
  for (i=0; i<sLen; i++)
    {H=(H<<4)+s.charAt(i);
     if ((g=H&0xf0000000)!=0)
       {H=H^(g>>24);
        H=H^g;
       }
    }
  return H%PRIME; }
```

7.5. Таблицы на деревьях

Рассмотрим еще один способ организации таблиц символов — с использованием *двоичных деревьев*.

Ориентированное дерево называется *двоичным*, если у него в каждую вершину, кроме одной (*корня*), входит одна дуга и из каждой вершины выходит не более двух дуг. *Ветвью* дерева называется поддерево, состоящее из некоторой дуги данного дерева, ее начальной и конечной вершин, а также всех вершин и дуг, лежащих на всех путях, выходящих из конечной вершины этой дуги. *Высотой* дерева называется максимальная длина пути в этом дереве от корня до листа.

Пусть на множестве идентификаторов задан некоторый линейный (например, лексикографический) порядок \prec , т. е. некоторое транзитивное, антисимметричное и антирефлексивное отношение. Таким образом, для произвольной пары идентификаторов id_1 и id_2 либо $id_1 \prec id_2$, либо $id_2 \prec id_1$, либо id_1 совпадает с id_2 .

Каждой вершине двоичного дерева, представляющего таблицу символов, сопоставим идентификатор. При этом если вершина (которой сопоставлен id) имеет левого потомка (которому сопоставлен id_L), то $id_L \prec id$; если она имеет правого потомка (id_R), то $id \prec id_R$. Элемент таблицы изображен на рис. 7.5.

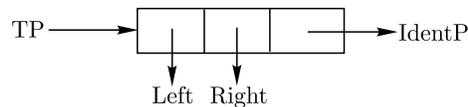


Рис. 7.5

Узел дерева описывается как класс `TreeElement`:

```
class TreeElement
{String IdentP;
  TreeElement Left, Right;
};
```

Поиск в такой таблице может быть описан следующей функцией:

```
TreeElement SearchTree(String Id,
                       TreeElement TP)
{int comp;
  if (TP==null) return null;
  comp=IdComp(Id, TP.IdentP);
  if (comp<0) return(SearchTree(Id, TP.Left));
  if (comp>0) return(SearchTree(Id, TP.Right));
  return TP;
}
```

Занесение в таблицу осуществляется функцией

```

void InsertTree(String Id,
                TreeElement TP)
{int comp=IdComp(Id,TP.IdentP);
  if (comp<0)
    if (TP.Left != null)
      InsertTree(Id,n.Leftt);
    else {n.Left = new TreeElement();
          n.Left.InentP = Id;
          n.Left.Left = null;
          n.Left.Right = null;}
    }
  else if (comp>0)if (TP.Right != null)
    InsertTree(Id,n.Right);
    else {n.Right = new TreeElement();
          n.Right.InentP = Id;
          n.Right.Left = null;
          n.Right.Right = null;}
    }
}
struct TreeElement * Fill(String Id,
                          struct TreeElement * P,
                          struct TreeElement ** FP)
{ if (P==NULL)
  {P=alloc(sizeof(struct TreeElement));
   P->IdentP=Include(Id);
   P->Left=NULL;
   P->Right=NULL;
   *FP=P;
   return(P);
  }
  else return(InsertTree(Id,P));
}

```

Как показано в [10], среднее время поиска в таблице размера n , организованной в виде двоичного дерева, при равных вероятностях появления каждого объекта равно $(2 \ln 2) \log_2 n + O(1)$. Однако на практике случай равных вероятностей появления объектов встречается довольно редко. Поэтому в дереве появляются более длинные и более короткие ветви, и среднее время поиска увеличивается.

Чтобы уменьшить среднее время поиска в двоичном дереве, можно в процессе построения дерева следить за тем, чтобы оно все время оставалось сбалансированным. А именно, назовем дерево *сбалансированным*, если ни для какой вершины высота выходящей из нее правой ветви не отличается от высоты левой более чем на 1. Для того чтобы достичь сбалансированности,

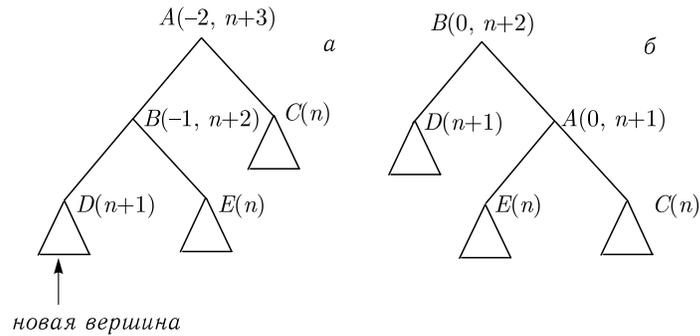


Рис. 7.6

в процессе добавления новых вершин дерево можно слегка перестраивать следующим образом [1].

Определим для каждой вершины дерева характеристику, равную разности высот выходящих из нее правой и левой ветвей. В сбалансированном дереве характеристика вершины может быть равной -1 , 0 или 1 , для листьев она равна 0 .

Пусть мы определили место новой вершины в дереве. Ее характеристика равна 0 . Назовем путь, ведущий от корня к новой вершине, *выделенным*. При добавлении новой вершины могут измениться характеристики только тех вершин, которые лежат на выделенном пути. Рассмотрим *заключительный* отрезок выделенного пути — такой, что до добавления вершины характеристики всех вершин на нем были равны 0 . Если верхним концом этого отрезка является сам корень, то дерево перестраивать не надо, достаточно лишь изменить характеристики вершин на этом пути на 1 или -1 , в зависимости от того, влево или вправо построена новая вершина.

Пусть верхний конец заключительного отрезка — не корень. Рассмотрим вершину A — «родителя» верхнего конца заключительного отрезка. Перед добавлением новой вершины характеристика A была равна ± 1 . Если A имела характеристику 1 (-1) и новая вершина добавляется в левую (правую) ветвь, то характеристика вершины A становится равной 0 , а высота поддерева с корнем в A не меняется. Так что и в этом случае дерево перестраивать не надо.

Пусть теперь характеристика A до перестроения была равна -1 и новая вершина добавлена к левой ветви A (аналогично — для случая 1 и добавления к правой ветви). Рассмотрим вершину B — левого потомка A . Возможны следующие варианты.

Если характеристика B после добавления новой вершины в D стала равна -1 , то дерево имеет структуру, изображенную на рис. 7.6, *а*. Перестроив дерево так, как это изображено на рис. 7.6, *б*, мы добьемся сбалансированности

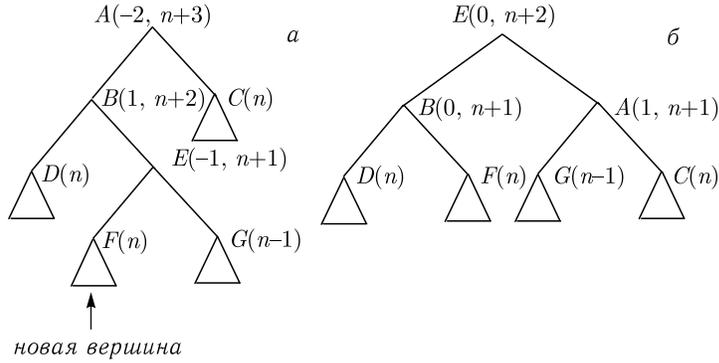


Рис. 7.7

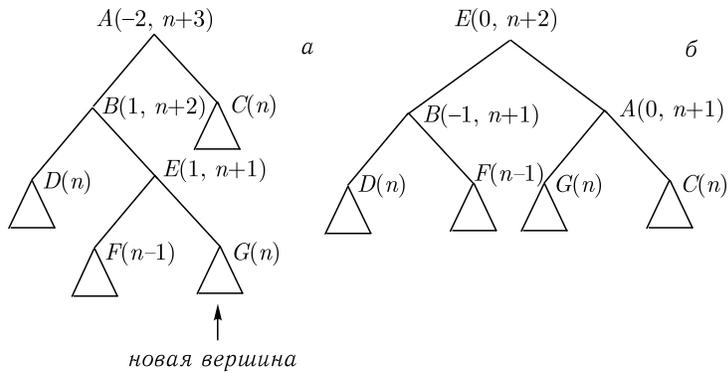


Рис. 7.8

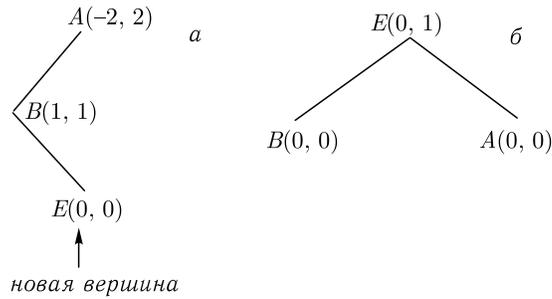


Рис. 7.9

(в скобках указаны характеристики вершин, где это существенно, и соотношения высот после добавления).

Если характеристика вершины B после добавления новой вершины в E стала равна 1, то надо отдельно рассмотреть случаи, когда характеристика вершины E , следующей за B на выделенном пути, стала равна -1 , 1 и 0 (в последнем случае вершина E — новая). Вид дерева для этих случаев показан соответственно на рис. 7.7, 7.8 и 7.9 (a — до, b — после перестроения).

7.6. Реализация блочной структуры

С точки зрения структуры программы блоки (и/или процедуры) образуют дерево. Каждой вершине дерева этого представления, соответствующей блоку, можно сопоставить свою таблицу символов (и, возможно, одну общую таблицу идентификаторов). Работу с таблицами блоков можно организовать в магазинном режиме: при входе в блок создавать таблицу символов, при выходе — уничтожать. При этом сами таблицы должны быть связаны в упорядоченный список, чтобы можно было просматривать их в порядке вложенности. Если таблицы организованы с помощью функций расстановки, то это означает, что для каждой таблицы должна быть создана своя таблица расстановки.

7.7. Сравнение методов реализации таблиц

Рассмотрим преимущества и недостатки рассмотренных методов реализации таблиц с точки зрения техники использования памяти.

Использование динамической памяти, как правило, — довольно дорогая операция, поскольку механизмы поддержания работы с динамической памятью могут быть достаточно сложны. Необходимо поддерживать списки свободной и занятой памяти, выбирать наиболее подходящий кусок памяти при запросе, включать освободившийся кусок в список свободной памяти и, возможно, склеивать куски свободной памяти в списке.

С другой стороны, использование массива требует предварительного отведения довольно большой памяти, а это означает, что значительная память вообще не будет использоваться. Кроме того, часто приходится заполнять не все элементы массива (например, в таблице идентификаторов или в тех случаях, когда в массиве фактически хранятся записи переменной длины, например, если в таблице символов записи для различных объектов имеют различный состав полей). Обращение к элементам массива может означать использование операции умножения при вычислении индексов, что может замедлить исполнение.

Наилучшим, по-видимому, является механизм доступа по указателям с использованием факта магазинной организации памяти в компиляторе. Для этого процедура выделения памяти выдает необходимый кусок из подряд идущей памяти, а при выходе из процедуры вся память, связанная с этой процедурой, освобождается простой перестановкой указателя свободной памяти в состояние перед началом обработки процедуры. В чистом виде это не всегда, однако, возможно. Например, локальный модуль в Модуле-2 может экспортировать некоторые объекты наружу. При этом схему реализации приходится «подгонять» под механизм распределения памяти. В данном случае, например, необходимо экспортированные объекты вынести в среду охватывающего блока и свернуть блок локального модуля.

Глава 8

ПРОМЕЖУТОЧНОЕ ПРЕДСТАВЛЕНИЕ ПРОГРАММЫ

В процессе трансляции компилятор часто используют промежуточное представление (ПП) исходной программы, предназначенное прежде всего для удобства генерации кода и/или проведения различных оптимизаций. Сама форма ПП зависит от целей его использования.

Наиболее часто используемыми формами ПП является ориентированный граф (в частности, абстрактное синтаксическое дерево, в том числе атрибутированное), трехадресный код (в виде троек или четверок), префиксная и постфиксная записи.

8.1. Представление в виде ориентированного графа

Простейшей формой промежуточного представления является синтаксическое дерево программы. Ту же самую информацию о входной программе, но в более компактной форме дает ориентированный ациклический граф (ОАГ), в котором в одну вершину объединены вершины синтаксического дерева, представляющие общие подвыражения. Синтаксическое дерево и ОАГ для оператора присваивания

$$a := b * -c + b * -c$$

приведены на рис. 8.1.

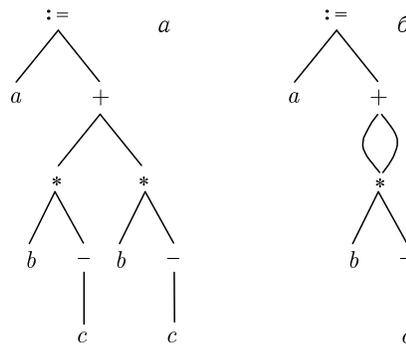


Рис. 8.1

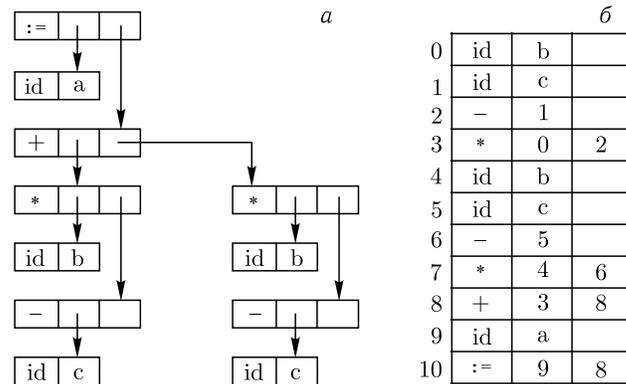


Рис. 8.2

На рис. 8.2 приведены два представления в памяти синтаксического дерева на рис. 8.1, а. Каждая вершина кодируется записью с полем для операции и полями для указателей на потомков. На рис. 8.2, б вершины размещены в массиве записей и индекс (или вход) вершины служит указателем на нее.

8.2. Трехадресный код

Трехадресный код — это последовательность операторов вида $x := y \text{ op } z$, где x , y и z — имена, константы или сгенерированные компилятором временные объекты. Здесь op — двуместная операция, например, операция плавающей или фиксированной арифметики, логическая или побитовая. В правую часть может входить только один знак операции.

Составные выражения должны быть разбиты на подвыражения, при этом могут появиться временные имена (переменные). Смысл термина «трехадресный код» в том, что каждый оператор обычно имеет три адреса: два для операндов и один для результата. Трехадресный код — это линейризованное представление синтаксического дерева или ОАГ, в котором временные имена соответствуют внутренним вершинам дерева или графа. Например, выражение $x + y * z$ может быть протранслировано в последовательность операторов

```
t1 := y * z
t2 := x + t1
```

где $t1$ и $t2$ — имена, сгенерированные компилятором.

В виде трехадресного кода представляются не только двуместные операции, входящие в выражения. В таком же виде представляются операторы управления программой и одноместные операции. В этом случае некоторые из компонент трехадресного кода могут не использоваться. Например, условный оператор

```
if A > B then S1 else S2
```

может быть представлен следующим кодом:

```
t := A - B
JGT t, S2
...
```

Здесь JGT — двуместная операция условного перехода, не вырабатывающая результата.

Разбиение арифметических выражений и операторов управления делает трехадресный код удобным при генерации машинного кода и оптимизации. Использование имен промежуточных значений, вычисляемых в программе, позволяет легко переупорядочивать трехадресный код.

Представления синтаксического дерева и графа рис. 8.1 в виде трехадресного кода имеют следующий вид:

Синтаксическое дерево	Граф
t1 := -c	t1 := -c
t2 := b * t1	t2 := b * t1
t3 := -c	t5 := t2 + t2
t4 := b * t3	a := t5
t5 := t2 + t4	
a := t5	

Трехадресный код — это абстрактная форма промежуточного кода. В реализации трехадресный код может быть представлен записями с полями для операции и операндов. Рассмотрим три способа реализации трехадресного кода: четверки, тройки и косвенные тройки.

Четверка — это запись с четырьмя полями, которые будем называть *op*, *arg1*, *arg2* и *result*. Поле *op* содержит код операции. В операторах с унарными операциями типа $x := -y$ или $x := y$ поле *arg2* не используется. В некоторых операциях (типа «передать параметр») не используются ни *arg2*, ни *result*. Условные и безусловные переходы помещают в *result* метку перехода. В табл. 8.1 приведены четверки для оператора присваивания $a := b * -c + b * -c$. Они получены из трехадресного кода, приведенного выше.

Обычно содержимое полей *arg1*, *arg2* и *result* — это указатели на входы таблицы символов для имен, представляемых этими полями. Временные имена вносятся в таблицу символов по мере их генерации.

Чтобы избежать внесения новых имен в таблицу символов, на временное значение можно ссылаться, используя позицию вычисляющего его оператора. В этом случае трехадресные операторы могут быть представлены записями только с тремя полями: *op*, *arg1* и *arg2*, как это показано в табл. 8.2. Поля *arg1* и *arg2* — это либо указатели на таблицу символов (для имен, определенных программистом, или констант), либо указатели на тройки (для временных значений). Такой способ представления трехадресного кода

Таблица 8.1

	<i>op</i>	<i>arg1</i>	<i>arg2</i>	<i>result</i>
(0)	-	c		t1
(1)	*	b	t1	t2
(2)	-	c		t3
(3)	*	b	t3	t4
(4)	+	t2	t4	t5
(5)	:=	t5		a

называют *тройками*. Тройки соответствуют представлению синтаксического дерева или ОАГ с помощью массива вершин.

Таблица 8.2

	<i>op</i>	<i>arg1</i>	<i>arg2</i>
(0)	-	c	
(1)	*	b	(0)
(2)	-	c	
(3)	*	b	(2)
(4)	+	(1)	(3)
(5)	:=	a	(4)

Числа в скобках — это указатели на тройки, а имена — это указатели на таблицу символов. На практике информация, необходимая для интерпретации различного типа входов в поля *arg1* и *arg2*, кодируется в поле *op* или дополнительных полях. Тройки табл. 8.2 соответствуют четверкам табл. 8.1.

Для представления тройками трехместной операции типа $x[i] := y$ требуется два входа, как это показано в табл. 8.3, представление $x := y[i]$ двумя операциями показано в табл. 8.4.

Таблица 8.3

	<i>op</i>	<i>arg1</i>	<i>arg2</i>
(0)	[]=	x	i
(1)	:=	(0)	y

Трехадресный код может быть представлен не списком троек, а списком указателей на них. Такая реализация обычно называется *косвенными тройками*. Например, тройки табл. 8.2 могут быть реализованы так, как это представлено в табл. 8.5.

Таблица 8.4

	<i>op</i>	<i>arg1</i>	<i>arg2</i>
(0)	= []	<i>y</i>	<i>i</i>
(1)	:=	<i>x</i>	(0)

Таблица 8.5

	оператор	<i>op</i>	<i>arg1</i>	<i>arg2</i>
(0)	(14)	-	<i>c</i>	
(1)	(15)	*	<i>b</i>	(14)
(2)	(16)	-	<i>c</i>	
(3)	(17)	*	<i>b</i>	(16)
(4)	(18)	+	(15)	(17)
(5)	(19)	:=	<i>a</i>	(18)

При генерации объектного кода каждой переменной, как временной, так и определенной в исходной программе, назначается память периода исполнения, адрес которой обычно хранится в таблице генератора кода. При использовании четверок этот адрес легко получить через эту таблицу.

Более существенно преимущество четверок проявляется в оптимизирующих компиляторах, когда может возникнуть необходимость перемещать операторы. Если перемещается оператор, вычисляющий *x*, то не требуется изменений в операторе, использующем *x*. В записи же тройками перемещение оператора, определяющего временное значение, требует изменения всех ссылок на этот оператор в массивах *arg1* и *arg2*. Из-за этого тройки трудно использовать в оптимизирующих компиляторах.

В случае применения косвенных троек оператор может быть перемещен переупорядочиванием списка операторов. При этом не надо менять указатели на *op*, *arg1* и *arg2*. Этим косвенные тройки похожи на четверки. Кроме того, эти два способа требуют примерно одинаковой памяти. Как и в случае простых троек, при использовании косвенных троек выделение памяти для временных значений может быть отложено на этап генерации кода. По сравнению с четверками при использовании косвенных троек можно сэкономить память, если одно и то же временное значение используется более одного раза. Например, в табл. 8.5 можно объединить строки (14) и (16), после чего можно объединить строки (15) и (17).

8.3. Линеаризованные представления

В качестве промежуточных представлений весьма распространены линеаризованные представления деревьев. Линеаризованное представление позволяет сравнительно легко хранить промежуточное представление во внешней памяти и обрабатывать его. Наиболее распространенной формой линеаризованного представления является *польская запись* — префиксная (прямая) или постфиксная (обратная).

Постфиксная запись — это список вершин дерева, в котором каждая вершина следует (при обходе снизу-вверх слева-направо) непосредственно за своими потомками. Дерево на рис. 8.1, а в постфиксной записи может быть представлено следующим образом:

$$a \ b \ c \ - \ * \ b \ c \ - \ * \ + \ :=$$

В постфиксной записи вершины синтаксического дерева не присутствуют явно. Они могут быть восстановлены из порядка, в котором следуют вершины, и из числа операндов соответствующих операций. Восстановление вершин аналогично вычислению выражения в постфиксной записи с использованием магазина.

В префиксной записи сначала указывается операция, а затем ее операнды. Например, для приведенного выше выражения имеем

$$:= \ a \ + \ * \ b \ - \ c \ * \ b \ - \ c$$

Рассмотрим детальнее одну из реализаций префиксного представления — Лидер [12]. Лидер — это аббревиатура от «ЛИнеаризованное ДЕРЕво». Это машинно-независимая префиксная запись. В Лидере сохраняются все объявления и каждому из них присваивается свой уникальный номер, который используется для ссылки на объявление. Рассмотрим пример:

```

module M;
var X, Y, Z: integer;
procedure DIF(A, B: integer): integer;
  var R: integer;
  begin R := A - B;
        return(R);
  end DIF;
begin Z := DIF(X, Y);
end M.
```

Этот фрагмент имеет следующий образ в Лидере:

```

program 'M'
var int
var int
var int
```

```

procbody proc int int end int
  var int
  begin assign var 1 7 end
    int int mi par 1 5 end par 1 6 end
    result 0 int var 1 7 end
    return
  end
begin assign var 0 3 end int
  icall 0 4 int var 0 1 end
  int var 0 2 end end
end

```

Рассмотрим его более детально:

program 'M'	Имя модуля нужно для редактора связей.
var int	Это образ переменных X, Y, Z;
var int	переменным X, Y, Z присваиваются
var int	номера 1, 2, 3 на уровне 0.
procbody proc	Объявление процедуры с двумя
int int end	целыми параметрами, возвращающей целое.
int	Процедура получает номер 4 на уровне 0, и параметры имеют номера 5, 6 на уровне 1.
var int	Переменная R имеет номер 7 на уровне 1.
begin	Начало тела процедуры.
assign	Оператор присваивания.
var 1 7 end	Левая часть присваивания (R).
int	Тип присваиваемого значения.
int mi	Целое вычитание.
par 1 5 end	Уменьшаемое (A).
par 1 6 end	Вычитаемое (B).
result 0	Результат процедуры уровня 0.
int	Результат имеет целый тип.
var 1 7 end	Результат — переменная R.
return	Оператор возврата.

<code>end</code>	Конец тела процедуры.
<code>begin</code>	Начало тела модуля.
<code>assign</code>	Оператор присваивания.
<code>var 0 3 end</code>	Левая часть — переменная <code>Z</code> .
<code>int</code>	Тип присваиваемого значения.
<code>icall 0 4</code>	Вызов локальной процедуры <code>DIF</code> .
<code>int var 0 1 end</code>	Фактические параметры <code>X</code>
<code>int var 0 2 end</code>	и <code>Y</code> .
<code>end</code>	Конец вызова.
<code>end</code>	Конец тела модуля.

8.4. Виртуальная машина Java

Программы на языке Java транслируются в специальное промежуточное представление, которое затем интерпретируется так называемой «виртуальной машиной Java». Виртуальная машина Java представляет собой магазинную машину: она не имеет памяти прямого доступа, все операции выполняются над операндами, расположенными на верхушке магазина. Чтобы, например, выполнить операцию с участием константы (или переменной), ее необходимо предварительно загрузить на верхушку магазина. Код операции — всегда один байт. Если операция имеет операнды, то они располагаются в следующих байтах.

К элементарным типам данных, с которыми работает машина, относятся `short`, `integer`, `long`, `float`, `double` (все знаковые).

8.4.1. Организация памяти. Машина имеет следующие регистры:

`pc` — счетчик команд;

`optop` — указатель вершины магазина операций;

`frame` — указатель на область исполняемого метода в магазине;

`vars` — указатель на 0-ю переменную исполняемого метода.

Все регистры 32-разрядные. Магазиновый фрейм имеет три компонента: локальные переменные, среду исполнения, магазин операндов. Локальные переменные отсчитываются от адреса в регистре `vars`. Среда исполнения служит для поддержания самого магазина. Она включает указатель на предыдущий фрейм, указатель на собственные локальные переменные, на базу стека операций и на верхушку магазина. Кроме того, здесь же хранится некоторая дополнительная информация, например, для отладчика.

Куча сбора мусора содержит экземпляры объектов, которые создаются и уничтожаются автоматически. Область методов содержит коды, таблицы символов и т. п.

С каждым классом связана область констант. Она содержит имена полей, методов и другую подобную информацию, которая используется методами.

8.4.2. Набор команд виртуальной машины. Виртуальная Java-машина имеет следующие команды:

помещение констант в магазин,
помещение локальных переменных в магазин,
запоминание значений из магазина в локальных переменных,
обработка массивов,
управление магазином,
арифметические команды,
логические команды,
преобразования типов,
передача управления,
возврат из функции,
табличный переход,
обработка полей объектов,
вызов метода,
обработка исключительных ситуаций,
прочие операции над объектами,
мониторы,
отладка.

Рассмотрим некоторые команды подробнее.

8.4.2.1. Помещение локальных переменных в магазин. Команда `iload`— загрузить целое из локальной переменной.

Операндом является смещение переменной в области локальных переменных. Указываемое значение копируется на верхушку магазина операций. Имеются аналогичные команды для помещения плавающих, двойных целых, двойных плавающих и т. п.

Команда `istore`— сохранить целое в локальной переменной.

Операндом операции является смещение переменной в области локальных переменных. Значение с верхушки магазина операций копируется в указываемую область локальных переменных. Имеются аналогичные команды для помещения плавающих, двойных целых, двойных плавающих и т. п.

8.4.2.2. Вызов метода. Команда `invokevirtual`.

При трансляции объектно-ориентированных языков программирования из-за возможности перекрытия виртуальных методов, вообще говоря, нельзя статически протранслировать вызов метода объекта. Это связано с тем, что если метод перекрыт в производном классе и вызывается метод объекта-переменной, то статически неизвестно, объект какого класса (базового или производного) хранится в переменной. Поэтому с каждым объектом связывается таблица всех его виртуальных методов: для каждого метода там

помещается указатель на его реализацию в соответствии с принадлежностью самого объекта классу в иерархии классов.

В языке Java различные классы могут реализовывать один и тот же интерфейс. Если объявлена переменная или параметр типа интерфейс, то динамически нельзя определить, объект какого класса присвоен переменной:

```
interface I;
class C1 implements I;
class C2 implements I;
I O;
C1 O1;
C2 O2;
...
O=O1;
...
O=O2;
...
```

В этой точке программы, вообще говоря, нельзя сказать, какого типа значение хранится в переменной O. Кроме того, при работе программы на языке Java имеется возможность использования методов из других пакетов. Для реализации этого механизма в Java-машине используется динамическое связывание.

Предполагается, что магазин операндов содержит `handle` (указатель на описание) объекта или массива и некоторое количество аргументов. Операнд операции используется для конструирования индекса в области констант текущего класса. Элемент по этому индексу в области констант содержит полную сигнатуру метода. Сигнатура метода описывает типы параметров и возвращаемого значения. Из `handle` объекта извлекается указатель на таблицу методов объекта. Просматривается сигнатура метода в таблице методов. Результатом этого просмотра является индекс в таблицу методов именованного класса, для которого найден указатель на блок метода. Блок метода указывает тип метода (`native`, `synchronized` и т.п.) и число аргументов, ожидаемых в магазине операндов.

Если метод помечен как `synchronized`, то запускается монитор, связанный с `handle`. Базис массива локальных переменных для нового магазинного фрейма устанавливается так, что он указывает на `handle` в магазине. Определяется общее число локальных переменных, используемых методом, и, после того как отведено необходимое место для локальных переменных, окружение исполнения нового фрейма помещается в магазин. База магазина операндов для этого вызова метода устанавливается на первое слово после окружения исполнения. Затем исполнение продолжается с первой инструкции вызванного метода.

8.4.2.3. *Обработка исключительных ситуаций.* Команда `athrow` — возбуждает исключительную ситуацию.

С каждым методом связан список операторов `catch`. Каждый оператор `catch` описывает диапазон инструкций, для которых он активен, и тип исключения, который он обрабатывает. Кроме того, с оператором связан набор реализующих его инструкций. При возникновении исключительной ситуации просматривается список операторов `catch`, чтобы установить соответствие. Исключительная ситуация соответствует оператору `catch`, если инструкция, вызвавшая исключительную ситуацию, находится в соответствующем диапазоне, а сама исключительная ситуация принадлежит подтипу типа ситуаций, которые обрабатывает оператор `catch`. Если соответствующий оператор `catch` найден, то управление передается обработчику. Если нет — текущий магазинный фрейм удаляется, и исключительная ситуация возбуждается вновь.

Порядок операторов `catch` в списке важен. Интерпретатор передает управление первому подходящему оператору `catch`.

8.5. Организация информации в генераторе кода

Синтаксическое дерево в чистом виде несет только информацию о структуре программы. На самом деле в процессе генерации кода требуется также информация о переменных (например, их адреса), процедурах (также адреса, уровни), метках и т. п. Для представления этой информации возможны различные решения. Наиболее распространены два:

- информация хранится в таблицах генератора кода;
- информация хранится в соответствующих вершинах дерева.

Рассмотрим, например, структуру таблиц, которые могут быть использованы в сочетании с Лидер-представлением. Поскольку Лидер-представление не содержит информации об адресах переменных, то эту информацию нужно формировать в процессе обработки объявлений и хранить в таблицах. Это касается и описаний массивов, записей и т. п. Кроме того, в таблицах также должна содержаться информация о процедурах (адреса, уровни, модули, в которых процедуры описаны, и т. п.).

При входе в процедуру в таблице уровней процедур заводится новый вход — указатель на таблицу описаний. При выходе указатель восстанавливается на старое значение. Если промежуточное представление — дерево, то информация может храниться в вершинах самого дерева.

8.6. Уровень промежуточного представления

Как видно из приведенных примеров, промежуточное представление программы может быть в различной степени близким либо к исходной программе, либо к машине. Например, промежуточное представление может содержать адреса переменных, и тогда оно уже не может быть перенесено на другую машину. С другой стороны, промежуточное представление может содержать раздел описаний программы, и тогда информацию об адресах можно извлечь из обработки описаний. В то же время ясно, что первое более эффективно, чем второе. Операторы управления в промежуточном представлении могут быть представлены в исходном виде (в виде операторов языка `if`, `for`, `while` и т. п.), а могут содержаться в виде переходов. В первом случае некоторая информация может быть извлечена из самой структуры (например, для оператора `for` — информация о переменной цикла, которую, может быть, разумно хранить на регистре, для оператора `case` — информация о таблице меток и т. п.). Во втором случае представление проще и унифицированней.

Некоторые формы промежуточного представления удобны для различного рода оптимизаций, некоторые — нет (например, косвенные тройки, в отличие от префиксной записи, позволяют эффективное перемещение кода).

Глава 9

ГЕНЕРАЦИЯ КОДА

Задача генератора кода — построение для программы на входном языке эквивалентной машинной программы. Обычно в качестве входа для генератора кода служит некоторое промежуточное представление программы.

Генерация кода включает ряд специфических независимых подзадач: распределение памяти (в частности, распределение регистров), выбор команд, генерацию объектного (или загрузочного) модуля. Конечно, независимость этих подзадач относительна: например, при выборе команд нельзя не учитывать схему распределения памяти, и, наоборот, схема распределения памяти (в частности, регистров) ведет к генерации той или иной последовательности команд. Однако удобно и практично эти задачи все же разделять, обращая при этом внимание на их взаимодействие.

В какой-то мере схема генератора кода зависит от формы промежуточного представления. Ясно, что генерация кода из дерева отличается от генерации кода из троек, а генерация кода из префиксной записи отличается от генерации кода из ориентированного графа. В то же время все генераторы кода имеют много общего, и основные применяемые алгоритмы различаются, как правило, только в деталях, связанных с используемым промежуточным представлением.

В дальнейшем в качестве промежуточного представления мы будем использовать префиксную нотацию. А именно, алгоритмы генерации кода будем излагать в виде атрибутных схем со входным языком Лидер.

9.1. Модель машины

При изложении алгоритмов генерации кода мы будем следовать некоторой модели машины, в основу которой положена система команд микропроцессора Motorola MC68020. В микропроцессоре имеется регистр — счетчик команд PC, 8 регистров данных и 8 адресных регистров.

В системе команд используются следующие способы адресации:

ABS — абсолютная: исполнительным адресом является значение адресного выражения.

IMM — непосредственный операнд: операндом команды является константа, заданная в адресном выражении.

D — прямая адресация через регистр данных, записывается как Xn , операнд находится в регистре Xn .

A — прямая адресация через адресный регистр, записывается как An , операнд находится в регистре An .

INDIRECT — записывается как (An) , адрес операнда находится в адресном регистре An .

POST — пост-инкрементная адресация, записывается как $(An)+$, исполнительный адрес есть значение адресного регистра An , и после исполнения команды значение этого регистра увеличивается на длину операнда.

PRE — пре-инкрементная адресация, записывается как $-(An)$: перед исполнением операции содержимое адресного регистра An уменьшается на длину операнда, исполнительный адрес равен новому содержимому адресного регистра.

INDISP — косвенная адресация со смещением, записывается как (bd, An) , исполнительный адрес вычисляется как $(An)+d$ — содержимое An плюс d .

INDEX — косвенная адресация с индексом, записывается как $(bd, An, Xn*sc)$, исполнительный адрес вычисляется как $(An)+bd+(Xn)*sc$ — содержимое адресного регистра + адресное смещение + содержимое индексного регистра, умноженное на sc .

INDIRPC — косвенная через PC (счетчик команд), записывается как (bd, PC) , исполнительный адрес определяется выражением $(PC)+bd$.

INDEXPC — косвенная через PC со смещением, записывается как $(bd, PC, Xn*sc)$, исполнительный адрес определяется выражением $(PC)+bd+(Xn)*sc$.

INDPRE — пре-косвенная через память, записывается как $([bd, An, sc*Xn], od)$ (схема вычисления адресов для этого и трех последующих способов адресации приведена ниже).

INDPOST — пост-косвенная через память: $([bd, An], sc*Xn, od)$.

INDPREPC — пре-косвенная через PC: $([bd, PC, sc*Xn], od)$.

INDPOSTPC — пост-косвенная через PC: $([bd, PC], Xn, od)$.

Здесь bd — это 16- или 32-битная константа, называемая *смещением*, od — 16- или 32-битная литеральная константа, называемая *внешним смещением*. Эти способы адресации могут использоваться в упрощенных формах без смещений bd и/или od и без регистров An или Xn . Следующие примеры иллюстрируют косвенную постиндексную адресацию:

```
MOVE D0, ([A0])
MOVE D0, ([4, A0])
MOVE D0, ([A0], 6)
MOVE D0, ([A0], D3)
MOVE D0, ([A0], D4, 12)
MOVE D0, ([$12345678, A0], D4, $FF000000)
```

Индексный регистр Xn может масштабироваться (умножаться) на 2, 4, 8, что записывается как $sc*Xn$. Например, в исполнительном адресе

($[24, A0, 4 * D0]$) содержимое квадратных скобок вычисляется как $[A0] + 4 * [D0] + 24$.

Эти способы адресации работают следующим образом. Каждый исполнительный адрес содержит пару квадратных скобок [...] внутри пары круглых скобок, т.е. ($[...], ...$). Сначала вычисляется содержимое квадратных скобок, в результате чего получается 32-битный указатель. Например, если используется пост-индексная форма $[20, A2]$, то исполнительный адрес — это $20 + [A2]$. Аналогично, для преиндексной формы $[12, A4, D5]$ исполнительный адрес — это $12 + [A4] + [D5]$.

Указатель, сформированный содержимым квадратных скобок, используется для доступа в память, чтобы получить новый указатель (отсюда термин «косвенная адресация через память»). К этому новому указателю добавляется содержимое внешних круглых скобок, и таким образом формируется исполнительный адрес операнда.

В дальнейшем изложении будут использованы следующие команды (в частности, рассматриваются только арифметические команды с целыми операндами, но не с плавающими):

MOVEA IA, A — загрузить содержимое по исполнительному адресу IA на адресный регистр A.

MOVE IA1, IA2 — содержимое по исполнительному адресу IA1 переписать по исполнительному адресу IA2.

MOVEM список_регистров, IA — сохранить указанные регистры в памяти, начиная с адреса IA (регистры указываются маской в самой команде).

MOVEM IA, список_регистров — восстановить указанные регистры из памяти, начиная с адреса IA (регистры указываются маской в самой команде).

LEA IA, A — загрузить исполнительный адрес IA на адресный регистр A.

MUL IA, D — умножить содержимое по исполнительному адресу IA на содержимое регистра данных D и результат разместить в D (на самом деле в системе команд имеются две различные команды **MULS** и **MULU** для чисел со знаком и чисел без знака соответственно; для упрощения мы не будем принимать во внимание это различие).

DIV IA, D — разделить содержимое регистра данных D на содержимое по исполнительному адресу IA и результат разместить в D.

ADD IA, D — сложить содержимое по исполнительному адресу IA с содержимым регистра данных D и результат разместить в D.

SUB IA, D — вычесть содержимое по исполнительному адресу IA из содержимого регистра данных D и результат разместить в D.

Команды **SMR** и **TST** формируют разряды регистра состояний. Всего имеется 4 разряда: *Z* — признак нулевого результата, *N* — признак отрицательного результата, *V* — признак переполнения, *C* — признак переноса.

СМР ИА, D — из содержимого регистра данных D вычитается содержимое по исполнительному адресу ИА, при этом формируются все разряды регистра состояний, но содержимое регистра D не меняется.

TST ИА — выработать разряд Z регистра состояний по значению, находящемуся по исполнительному адресу ИА.

VNE ИА — условный переход по признаку $Z = 1$ (не равно) на исполнительный адрес ИА.

VEQ ИА — условный переход по признаку $Z = 0$ (равно) на исполнительный адрес ИА.

VLE ИА — условный переход по признаку *N or Z* (меньше или равно) на исполнительный адрес ИА.

VGT ИА — условный переход по признаку *not N* (больше) на исполнительный адрес ИА.

VLT ИА — условный переход по признаку *N* (меньше) на исполнительный адрес ИА.

VRA ИА — безусловный переход по адресу ИА.

JMP ИА — безусловный переход по исполнительному адресу.

RTD размер_локальных — возврат из подпрограммы с указанием размера локальных.

LINK A, размер_локальных — в магазине сохраняется значение регистра A, в регистр A заносится указатель на это место в магазине, и указатель стека продвигается на размер локальных.

UNLK A — стек сокращается на размер локальных, и регистр A восстанавливается из магазина.

9.2. Динамическая организация памяти

Динамическая организация памяти — это организация памяти периода исполнения программы. Оперативная память программы обычно состоит из нескольких основных разделов: магазин, куча, область статических данных (инициализированных и неинициализированных). Наиболее сложной является работа с магазином. Вообще говоря, магазин периода исполнения необходим для программ не на всех языках программирования. Например, в ранних версиях Фортрана нет рекурсии, так что программа может исполняться без магазина. С другой стороны, исполнение программы с рекурсией может быть реализовано и без магазина (того же эффекта можно достичь, например, и с помощью списковых структур). Однако, для эффективной реализации пользуются магазином, который, как правило, поддерживается на уровне машинных команд.

Рассмотрим схему организации магазина периода выполнения для простейшего случая (как, например, в языке Паскаль), когда все переменные

в магазине (фактические параметры и локальные переменные) имеют смещения, известные при трансляции. Магазин служит для хранения локальных переменных (и параметров) и обращения к ним в языках, допускающих рекурсивные вызовы процедур. Еще одной задачей, которую необходимо решать при трансляции языков с блочной структурой, является обеспечение реализации механизмов статической вложенности. Пусть имеется следующий фрагмент программы на Паскале:

```
procedure P1;
  var V1;
  procedure P2;
    var V2;
  begin
    ...
    P2;
    V1:=...
    V2:=...
    ...
  end;
begin
  ...
  P2;
  ...
end;
```

В процессе выполнения этой программы, находясь в процедуре P2, мы должны иметь доступ к последнему набору значений переменных процедуры P2 и к набору значений переменных процедуры P1, из которой была вызвана P2. Кроме того, необходимо обеспечить восстановление состояния программы при завершении выполнения процедуры.

Мы рассмотрим две возможные схемы динамической организации памяти: схему со статической цепочкой и с дисплеем в памяти. В первом случае все статические контексты связаны в список, который называется *статической цепочкой*; в каждой записи для процедуры в магазине хранится указатель на запись статически охватывающей процедуры (помимо, конечно, указателя *динамической цепочки* — указателя на «базу» динамически предыдущей процедуры). Во втором случае для хранения ссылок на статические контексты используется массив, называемый *дисплеем*. Использование той или иной схемы определяется, помимо прочих условий, прежде всего числом адресных регистров.

9.2.1. Организация магазина со статической цепочкой. Итак, в случае статической цепочки магазин организован так, как это изображено на рис. 9.1.



Рис. 9.1

Таким образом, на запись текущей процедуры в магазине указывает регистр BP (Base Pointer), с которого начинается динамическая цепочка. На статическую цепочку указывает регистр LP (Link Pointer). В качестве регистров BP и LP в различных системах команд могут использоваться универсальные, адресные или специальные регистры. Локальные переменные отсчитываются от регистра BP вверх, фактические параметры — вниз с учетом памяти, занятой точкой возврата и самим сохраненным регистром BP.

Вызов подпрограмм различного статического уровня производится несколько по-разному. При вызове подпрограммы того же статического уровня, что и вызывающая подпрограмма (например, рекурсивный вызов той же самой подпрограммы), выполняются следующие команды:

Занесение фактических параметров в магазин
JSR A

Команда JSR A продвигает указатель SP, заносит PC на верхушку магазина и осуществляет переход по адресу A. После выполнения этих команд состояние магазина становится таким, как это изображено на рис. 9.2. Занесение BP, отведение локальных, сохранение регистров делает вызываемая подпрограмма.

При вызове локальной подпрограммы необходимо установить указатель статического уровня на текущую подпрограмму, а при выходе — восстановить

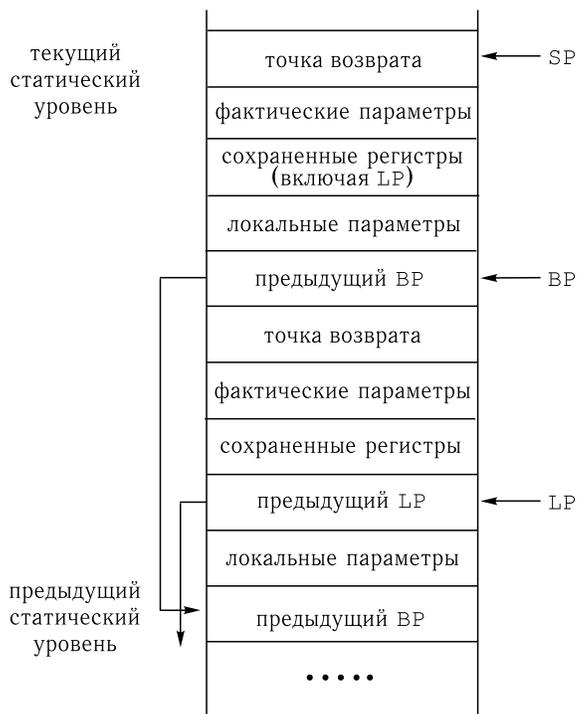


Рис. 9.2

его на старое значение (подпрограммы, охватывающей текущую). Для этого исполняются следующие команды:

```

Занесение фактических параметров в магазин
MOVE BP, LP
SUB Delta, LP
JSR A

```

Здесь *Delta* — размер локальных переменных вызывающей подпрограммы плюс двойная длина слова. Магазин после этого принимает состояние, изображенное на рис. 9.3. Предполагается, что регистр LP уже сохранен среди сохраняемых регистров, причем самым первым (сразу после локальных переменных).

После выхода из подпрограммы в вызывающей подпрограмме выполняется команда

```
MOVE (LP), LP
```

которая восстанавливает старое значение статической цепочки. Если выход осуществлялся из подпрограммы 1-го уровня, эту команду выполнять не надо, поскольку для 1-го уровня нет статической цепочки.

При вызове подпрограммы меньшего, чем вызывающая, уровня выполняются следующие команды:

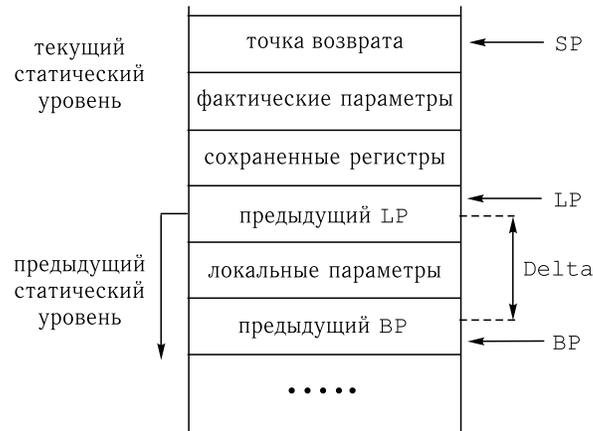


Рис. 9.3

```

Занесение фактических параметров в магазин
MOVE (LP), LP /* количество раз равно разности
                уровней вызывающей и вызываемой ПП */
JSR A

```

Тем самым устанавливается статический уровень вызываемой подпрограммы. После выхода из подпрограммы выполняется команда

```
MOVE -Delta(BP), LP
```

восстанавливающая статический уровень вызывающей подпрограммы.

Тело подпрограммы начинается со следующих команд:

```
LINK BP, -размер_локальных
MOVEM -(SP)
```

Команда `LINK BP, размер_локальных` эквивалентна трем командам:

```
MOVE BP, -(SP)
MOVE SP, BP
ADD -размер_локальных, SP
```

Команда `MOVEM` сохраняет в магазине регистры.

В результате выполнения этих команд магазин приобретает вид, изображенный на рис. 9.1.

Выход из подпрограммы осуществляется следующей последовательностью команд:

```
MOVEM (SP)+
UNLK BP
RTD размер_фактических
```

Команда `MOVEM` восстанавливает регистры из магазина. Команда `UNLK BP` эквивалентна такой последовательности команд:

```
MOVE BP, SP
MOVE (SP), BP
```

```
ADD #4, SP /* 4 - размер слова */
```

Команда `RTD размер_фактических`, в свою очередь, эквивалентна последовательности

```
ADD размер_фактических+4, SP  
JMP -размер_фактических-4(SP)
```

После ее выполнения магазин восстанавливается до состояния, которое было до вызова.

В зависимости от наличия локальных переменных, фактических параметров и необходимости сохранения регистров каждая из этих команд может отсутствовать.

9.2.2. Организация магазина с дисплеем. Рассмотрим теперь организацию магазина с дисплеем. Дисплей — это массив (`DISPLAY`), i -й элемент которого представляет собой указатель на область активации последней вызванной подпрограммы i -го статического уровня. Доступ к переменным самой внутренней подпрограммы осуществляется через регистр `BP`. Дисплей может быть реализован либо через регистры (если их достаточно), либо через массив в памяти.

Для вызова процедуры следующего (по отношению к вызывающей программе) уровня в дисплее отводится очередной элемент. Если вызывающая процедура имеет статический уровень i , то при вызове процедуры уровня $j \leq i$ элементы дисплея j, \dots, i должны быть скопированы (обычно в магазин вызывающей процедуры), текущим уровнем становится j и в `DISPLAY[j]` заносится указатель на область активации вызываемой процедуры. По окончании работы вызываемой процедуры содержимое дисплея восстанавливается из магазина.

Иногда используется комбинированная схема — дисплей в магазине. Дисплей хранится в области активации каждой процедуры. Формирование дисплея для процедуры осуществляется в соответствии с правилами, описанными выше.

Отдельного рассмотрения требует вопрос о технике передачи фактических параметров. Конечно, в случае простых параметров (например, чисел) проблем не возникает. Однако передача массивов по значению — операция довольно дорогая, поэтому с точки зрения экономии памяти целесообразнее сначала в подпрограмму передать адрес массива, а затем уже из подпрограммы по адресу передать в магазин сам массив. В связи с передачей параметров следует упомянуть еще одно обстоятельство.

Рассмотренная схема организации магазина допустима только для языков со статически известными размерами фактических параметров. Однако, например, в языке Модула-2 может быть передан по значению гибкий массив, и в этом случае нельзя статически распределить память для параметров.

Обычно в таких случаях заводят так называемый паспорт массива, в котором хранится вся необходимая информация, а сам массив размещается в магазине — в рабочей области регистров, сохраненных выше.

9.3. Назначение адресов

Назначение адресов переменным, параметрам и полям записей происходит при обработке соответствующих объявлений. В однопроходном трансляторе это может производиться вместе с построением основной таблицы символов, а соответствующие адреса (или смещения) могут храниться в этой же таблице. В промежуточном Лидер-представлении объявления сохранены, что делает это промежуточное представление машинно-независимым. Напомним, что в Лидер-представлении каждому описанию соответствует некоторый номер. В процессе работы генератора кодов поддерживается таблица `Table`, в которой по этому номеру (входу) содержится следующая информация:

- для типа: его размер;
- для переменной: смещение в области процедуры (или глобальной области);
- для поля записи: смещение внутри записи;
- для процедуры: размер локальных параметров;
- для массива: размер массива, размер элемента, значение левой и правой границ.

Функция `IncTab` вырабатывает указатель (вход) на новый элемент таблицы, проверяя при этом наличие свободной памяти.

Для вычисления адресов определим для каждого объявления два синтезируемых атрибута: `DISP` будет обозначать смещение внутри области процедуры (или единицы компиляции), а `SIZE` — размер. Тогда семантика правила для списка объявлений принимает вид

```
RULE
DeclPart ::= ( Decl )
SEMANTICS
  Disp<1>=0;
1A: Disp<1>=Disp<1>+Size<1>;
  Size<0>=Disp<1>.
```

Все объявления, кроме объявлений переменных, имеют нулевой размер. Размер объявления переменной определяется следующим правилом:

```
RULE
Decl ::= 'VAR' TypeDes
SEMANTICS
```

```

Tablentry Entry;
0: Entry=IncTab;
  Size<0>=((Table[VAL<2>]+1) / 2)*2;
  // Выравнивание на границу слова
  Table[Entry]=Disp<0>+Size<0>.

```

В качестве примера трансляции определения типа рассмотрим обработку описания записи:

```

RULE
TypeDes ::= 'REC' ( TypeDes ) 'END'
SEMANTICS
int Disp;
Tablentry Temp;
0: Entry<0>=IncTab;
  Disp=0;
2A: {Temp=IncTab;
     Table[Temp]=Disp;
     Disp=Disp+Table[Entry<2>]+1) / 2)*2;
     // Выравнивание на границу слова
   }
Table[Entry<0>]=Disp.

```

9.4. Трансляция переменных

Переменные отражают всё многообразие механизмов доступа в языке. Переменная имеет синтезированный атрибут ADDRESS — это запись, описывающая адрес в команде MC68020. Этот атрибут сопоставляется всем нетерминалам, представляющим значения. В системе команд MC68020 много способов адресации, и они отражены в структуре значения атрибута ADDRESS, имеющего следующий тип:

```

enum Register
{D0, D1, D2, D3, D4, D5, D6, D7,
 A0, A1, A2, A3, A4, A5, A6, SP, NO};

enum AddrMode
{D, A, Post, Pre, Indirect, IndPre, IndPost, IndirPC,
 IndPrePC, IndPostPC, InDisp, Index, IndexPC, Abs, Imm};

struct AddrType{
  Register AddrReg, IndexReg;
  int IndexDisp, AddrDisp;
  short Scale;
};

```

Значение регистра NO означает, что соответствующий регистр в адресации не используется.

Доступ к переменным осуществляется в зависимости от их уровня: глобальные переменные адресуются с помощью абсолютной адресации; переменные в процедуре текущего уровня адресуются через регистр базы А6.

Если магазин организован с помощью статической цепочки, то переменные предыдущего статического уровня адресуются через регистр статической цепочки А5; переменные остальных уровней адресуются «пробеганием» по статической цепочке с использованием вспомогательного регистра. Адрес переменной формируется при обработке структуры переменной слева направо и передается сначала сверху вниз как наследуемый атрибут нетерминала VarTail, а затем передается снизу-вверх как глобальный атрибут нетерминала Variable. Таким образом, правило для обращения к переменной имеет вид (первое вхождение Number в правую часть — это уровень переменной, второе — ее Лидер-номер):

```

RULE
Variable ::= VarMode Number Number VarTail
SEMANTICS
int Temp;
struct AddrType AddrTmp1, AddrTmp2;
3: if (Val<2>==0) // Глобальная переменная
    {Address<4>.AddrMode=Abs;
      Address<4>.AddrDisp=0;
    }
else // Локальная переменная
    {Address<4>.AddrMode=Index;
      if (Val<2>==Level<Block>) // Переменная
                                // текущего уровня
        Address<4>.AddrReg=A6;
      else if (Val<2>==Level<Block>-1)
        // Переменная предыдущего уровня
        Address<4>.AddrReg=A5;
      else
        {Address<4>.AddrReg=
          GetFree(RegSet<Block>);
          AddrTmp1.AddrMode=Indirect;
          AddrTmp1.AddrReg=A5;
          Emit2(MOVEA, AddrTmp1,
              Address<4>.AddrReg);
          AddrTmp1.AddrReg=Address<4>.AddrReg;
          AddrTmp2.AddrMode=A;
          AddrTmp2.AddrReg=Address<4>.AddrReg;
          for (Temp=Level<Block>-Val<2>;
              Temp>=2; Temp--)
            Emit2(MOVEA, AddrTmp1, AddrTmp2);
        }
    }

```

```

if (Val<2>==Level<Block>)
  Address<4>.AddrDisp=Table[Val<3>];
else
  Address<4>.AddrDisp=
    Table[Val<3>]+Table[LevelTab[Val<2>]];

}.

```

Функция `GetFree` выбирает очередной свободный регистр (либо регистр данных, либо адресный регистр) и отмечает его как использованный в атрибуте `RegSet` нетерминала `Block`. Процедура `Emit2` генерирует двухадресную команду. Первый параметр этой процедуры — код команды, второй и третий параметры имеют тип `AddrType` и служат операндами команды. Смещение переменной текущего уровня отсчитывается от базы (`A6`), а других уровней — от указателя статической цепочки, поэтому оно определяется как алгебраическая сумма размера локальных параметров и величины смещения переменной. Таблица `LevelTab` — это таблица уровней процедур, содержащая указатели на последовательно вложенные процедуры.

Если магазин организован с помощью дисплея, то трансляция для доступа к переменным может быть осуществлена следующим образом:

```

RULE
Variable ::= VarMode Number Number VarTail
SEMANTICS
int Temp;
3: if (Val<2>==0) // Глобальная переменная
  {Address<4>.AddrMode=Abs;
   Address<4>.AddrDisp=0;
  }
else // Локальная переменная
  {Address<4>.AddrMode=Index;
   if (Val<2>=Level<Block>) // Переменная
                               // текущего уровня
     {Address<4>.AddrReg=A6;
      Address<4>.AddrDisp=Table[Val<3>];
     }
   else
     {Address<4>.AddrMode=IndPost;
      Address<4>.AddrReg=NO;
      Address<4>.IndexReg=NO;
      Address<4>.AddrDisp=Display[Val<2>];
      Address<4>.IndexDisp=Table[Val<3>];
     }
  }
}.

```

Рассмотрим трансляцию доступа к полям записи. Она описывается следующим правилом (`Number` — это Лидер-номер описания поля):

```

RULE
VarTail ::= 'FIL' Number VarTail
SEMANTICS
if (Address<0>.AddrMode==Abs)
  {Address<3>.AddrMode=Abs;
   Address<3>.AddrDisp=
     Address<0>.AddrDisp+Table[Val<2>];
  }
else
  {Address<3>=Address<0>;
   if (Address<0>.AddrMode==Index)
     Address<3>.AddrDisp=
       Address<0>.AddrDisp+Table[Val<2>];
   else
     Address<3>.IndexDisp=
       Address<0>.IndexDisp+Table[Val<2>];
  }.

```

9.5. Трансляция целых выражений

Трансляция выражений различных типов управляется синтаксически благодаря наличию указателя типа перед каждой операцией. Мы рассмотрим некоторые наиболее характерные проблемы генерации кода для выражений.

Система команд МС68020 обладает двумя особенностями, сказывающимися на генерации кода для арифметических выражений (то же можно сказать и о генерации кода для выражений типа «множества»):

- 1) один из операндов выражения (правый) должен при выполнении операции находиться на регистре, поэтому если оба операнда не на регистрах, то перед выполнением операции один из них надо загрузить на регистр;
- 2) система команд довольно «симметрична», т. е. нет специальных требований к регистрам при выполнении операций (таких, например, как пары регистров или требования четности и т. п.).

Поэтому выбор команд при генерации арифметических выражений определяется довольно простыми таблицами решений. Например, для целочисленного сложения такая таблица приведена в табл. 9.1.

Здесь имеется в виду, что R — операнд на регистре, V — переменная или константа. Такая таблица решений должна также учитывать коммутативность операций.

```

RULE
IntExpr ::= 'PLUS' IntExpr IntExpr
SEMANTICS
if (Address<2>.AddrMode!=D) &&

```

```

        (Address<3>.AddrMode!=D)
    {Address<0>.AddrMode=D;
    Address<0>.Areg=GetFree(RegSet<Block>);
    Emit2(MOVE,Address<2>,Address<0>);
    Emit2(ADD,Address<2>,Address<0>);
    }
else
    if (Address<2>.AddrMode==D)
        {Emit2(ADD,Address<3>,Address<2>);
        Address<0>:=Address<2>);
        }
    else {Emit2(ADD,Address<2>,Address<3>);
        Address<0>:=Address<3>);
        }.

```

Таблица 9.1

		Правый операнд A2	
		R	V
Левый операнд A1	R	ADD A1, A2	ADD A2, A1
	V	ADD A1, A2	MOVE A1, R ADD A2, R

9.6. Трансляция арифметических выражений

Одной из важнейших задач при генерации кода является распределение регистров. Рассмотрим хорошо известную технику распределения регистров при трансляции арифметических выражений, называемую алгоритмом Сети-Ульмана.

Замечание. Для большей наглядности в данном разделе мы немного отступаем от семантики арифметических команд MC68020 и предполагаем, что команда

Op Arg1, Arg2

выполняет действие Arg2:=Arg1 Op Arg2.

Пусть система команд машины имеет неограниченное число универсальных регистров, в которых выполняются арифметические команды. Рассмотрим, как можно сгенерировать код, используя для данного арифметического выражения минимальное число регистров.

Пусть имеется синтаксическое дерево выражения. Предположим сначала, что распределение регистров осуществляется по простейшей схеме

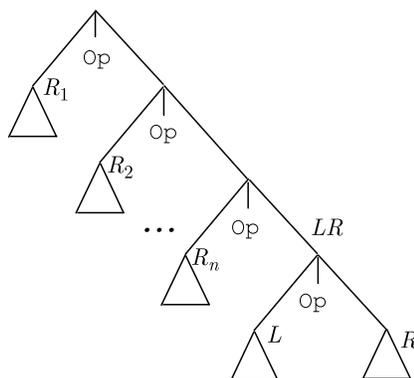


Рис. 9.4

сверху-вниз слева-направо, как изображено на рис. 9.4. Тогда к моменту генерации кода для поддерева LR занято n регистров. Пусть поддерево L требует n_l регистров, а поддерево R — n_r регистров. Если $n_l = n_r$, то при вычислении L будет использовано n_l регистров и под результат будет занят $(n + 1)$ -й регистр. Еще $n_r (= n_l)$ регистров будет использовано при вычислении R . Таким образом, общее число использованных регистров будет равно $n + n_l + 1$.

Если $n_l > n_r$, то при вычислении L будет использовано n_l регистров. При вычислении R будет использовано $n_r < n_l$ регистров, и всего будет использовано не более чем $n + n_l$ регистров. Если $n_l < n_r$, то после вычисления L под результат будет занят один регистр (предположим, $(n + 1)$ -й), а n_r регистров будет использовано для вычисления R . Всего будет использовано $n + n_r + 1$ регистров.

Видно, что для деревьев, совпадающих с точностью до порядка потомков каждой вершины, минимальное число регистров при распределении их слева-направо достигается на дереве, у которого в каждой вершине слева расположено более «сложное» поддерево, требующее большего числа регистров. Таким образом, если в каждой внутренней вершине дерева правое поддерево требует меньшего числа регистров, чем левое, то, обходя дерево слева направо, можно оптимально распределить регистры. Без перестроения дерева это означает, что если в некоторой вершине дерева справа расположено более сложное поддерево, то сначала сгенерируем код для него, а затем уже для левого поддерева.

Алгоритм работает следующим образом. Сначала осуществляется разметка синтаксического дерева. Разметка позволяет определить для каждого поддерева, сколько регистров потребуется для его вычисления. Разметка осуществляется по следующим правилам.

Правила разметки:

1. Если вершина — правый лист или дерево состоит из единственной вершины, то помечаем эту вершину числом 1 (рис. 9.5, б); если вершина — левый лист, то помечаем ее нулем (рис. 9.5, а).

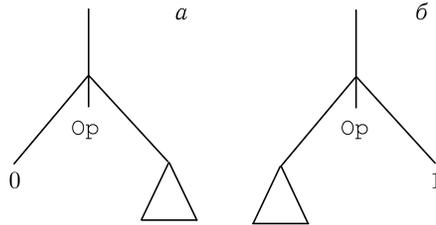


Рис. 9.5

2. Если вершина имеет прямых потомков с метками l_1 и l_2 , то в качестве метки этой вершины выбираем наибольшее из чисел l_1 , l_2 либо число $l_1 + 1$, если $l_1 = l_2$ (рис. 9.6).

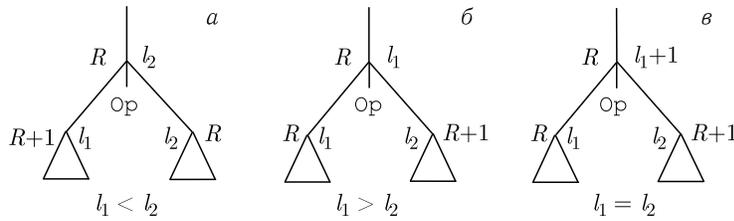


Рис. 9.6

Далее осуществляется распределение регистров для результатов операций. Распределение регистров основывается на следующих рассуждениях. Пусть корню поддерева назначен регистр с номером n . Индуктивное утверждение заключается в том, что к моменту вычисления поддерева все регистры с номерами $1, \dots, n$ заняты под промежуточные результаты уже вычисленных поддеревьев, и, возможно, при вычислении данного поддерева используются регистры с номерами, большими n . После завершения вычислений в данном поддереве значение хранится в регистре с номером n . Назначение регистров осуществляется по следующим правилам:

1) Корню назначается первый регистр.

2) Если метка левого потомка меньше или равна метке правого, то левому потомку назначается регистр на единицу больший, чем предку, а правому — с тем же номером (сначала вычисляется правое поддерево и его результат помещается в регистр R — рис. 9.6, а, в), так что регистры занимают последовательно. Если же метка левого потомка больше метки правого потомка, то наоборот, правому потомку назначается регистр на единицу больший, чем

предку, а левому — с тем же номером (сначала вычисляется левое поддерево и его результат помещается в регистр R — рис. 9.6, б).

После этого формируется код по следующим правилам.

1. Если вершина — правый лист с меткой 1, то ей соответствует код

MOVE X, R

где R — регистр, назначенный этой вершине, а X — адрес переменной, связанной с вершиной (рис. 9.7, б).

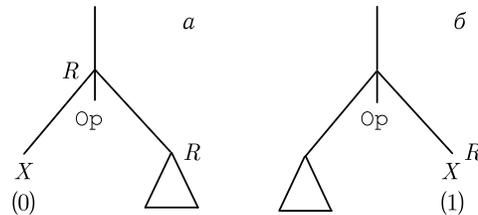


Рис. 9.7

2. Если вершина внутренняя и ее левый потомок — лист с меткой 0, то ей соответствует код

Код правого поддерева

Op X, R

где R — регистр, назначенный этой вершине, X — адрес переменной, связанной с вершиной, а Op — операция, примененная в вершине (рис. 9.7, а).

3. Если непосредственные потомки вершины — не листья и метка правой вершины больше или равна метке левой, то вершине соответствует код

Код правого поддерева

Код левого поддерева

Op $R+1, R$

где R — регистр, назначенный внутренней вершине, а операция Op, вообще говоря, — не коммутативная (рис. 9.8, б);

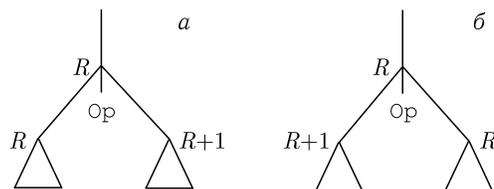


Рис. 9.8

4. Если непосредственные потомки вершины — не листья и метка правой вершины меньше метки левой вершины, то вершине соответствует код

Код левого поддерева

Код правого поддерева

Op $R, R+1$

MOVE $R+1, R$

Последняя команда генерируется для того, чтобы получить результат в нужном регистре (в случае коммутативной операции ее операнды можно поменять местами и избежать дополнительной пересылки — рис. 9.8, а).

Рассмотрим атрибутную схему, реализующую эти правила генерации кода (для большей наглядности входная грамматика соответствует обычной инфиксной записи, а не Лидер-представлению). В этой схеме генерация кода не происходит непосредственно в процессе обхода дерева, как раньше, а из-за необходимости переставлять поддеревья код строится в виде текста с помощью операции конкатенации. Практически, конечно, это нецелесообразно: разумнее управлять обходом дерева непосредственно, однако для простоты мы будем пользоваться конкатенацией.

RULE

Expr ::= IntExpr

SEMANTICS

Reg<1>=1; Code<0>=Code<1>; Left<1>=true.

RULE

IntExpr ::= IntExpr Op IntExpr

SEMANTICS

Left<1>=true; Left<3>=false;

Label<0>=(Label<1>==Label<3>)

? Label<1>+1

: Max(Label<1>,Label<3>);

Reg<1>=(Label<1> <= Label<3>)

? Reg<0>+1

: Reg<0>;

Reg<3>=(Label<1> <= Label<3>)

? Reg<0>

: Reg<0>+1;

Code<0>=(Label<1>==0)

? Code<3> + Code<2>

+ Code<1> + "," + Reg<0>

: (Label<1> < Label<3>)

? Code<3> + Code<1> + Code<2> +

(Reg<0>+1) + "," + Reg<0>

: Code<1> + Code<3> + Code<2> +

Reg<0> + "," + (Reg<0>+1)

+ "MOVE" + (Reg<0>+1) + "," + Reg<0>.

RULE

IntExpr ::= Ident

SEMANTICS

Label<0>=(Left<0>) ? 0 : 1;

Code<0>=(!Left<0>)

? "MOVE" + Reg<0> + "," + Val<1>

```
    : Val<1>.
```

```
RULE
IntExpr ::= '(' IntExpr ')'
SEMANTICS
Label<0>=Label<2>;
Reg<2>=Reg<0>;
Code<0>=Code<2>.
```

```
RULE
Op ::= '+'
SEMANTICS
Code<0>="ADD".
```

```
RULE
Op ::= '-'
SEMANTICS
Code<0>="SUB".
```

```
RULE
Op ::= '*'
SEMANTICS
Code<0>="MUL".
```

```
RULE
Op ::= '/'
SEMANTICS
Code<0>="DIV".
```

Атрибутированное дерево для выражения $A*B+C*(D+E)$ приведено на рис. 9.9. При этом будет сгенерирован следующий код:

```
MOVE B, R1
MUL  A, R1
MOVE E, R2
ADD  D, R2
MUL  C, R2
ADD  R1, R2
MOVE R2, R1
```

Приведенная атрибутивная схема требует двух проходов по дереву выражения. Рассмотрим теперь другую атрибутивную схему, в которой генерация программы для выражений с оптимальным распределением регистров требует лишь одного обхода [9].

Пусть мы произвели разметку дерева разбора так же, как и в предыдущем алгоритме. В качестве регистра результата вычисления поддерева будем назначать регистр с номером, равным метке вершины.

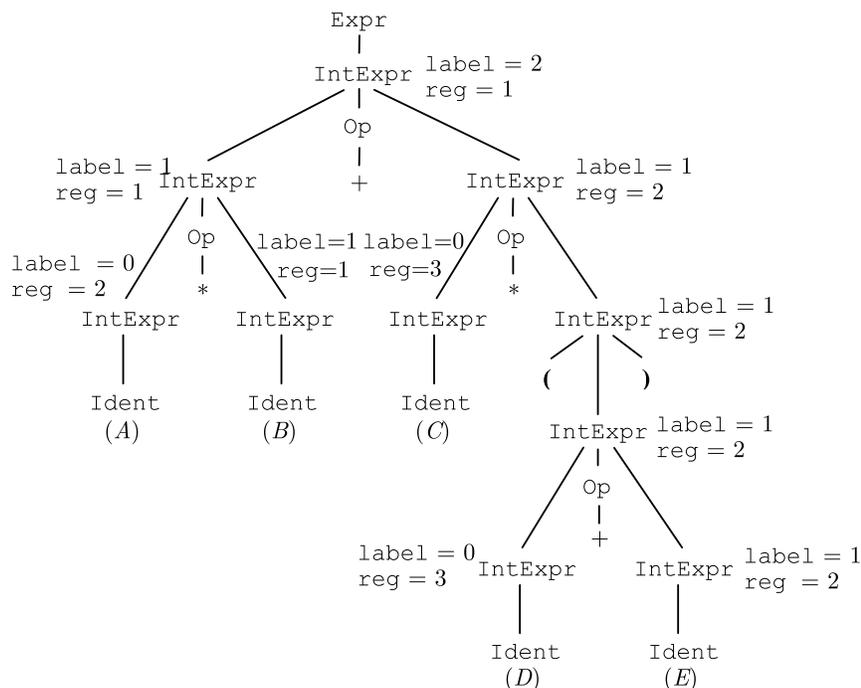


Рис. 9.9

Левому потомку всегда назначается регистр, равный его метке, а правому — его метке, если она не равна метке его левого брата, и метке, увеличенной на единицу, если метки равны. Поскольку более сложное поддерево всегда вычисляется раньше более простого, его регистр результата имеет больший номер, чем любой регистр, используемый при вычислении более простого поддерева, что гарантирует правильность использования регистров.

Приведенные соображения реализуются следующей атрибутивной схемой:

RULE

Expr ::= IntExpr

SEMANTICS

Code<0>=Code<1>; Left<1>=true.

RULE

IntExpr ::= IntExpr Op IntExpr

SEMANTICS

Left<1>=true; Left<3>=false;

Label<0>=(Label<1>==Label<3>)

? Label<1>+1

: Max(Label<1>,Label<3>);

Code<0>=(Label<3> > Label<1>)

? (Label<1>==0)

? Code<3> + Code<2> + Code<1>

```

      + ",R" + Label<3>
    : Code<3> + Code<1> + Code<2> + "R"+
      Label<1> + ",R" + Label<3>
  : (Label<3> < Label<1>)
    ? Code<1> + Code<3> + Code<2> + "R"+
      Label<1> + ", " + Label<3> +
      "MOVE" + Label<3> + ",R" + Label<1>
  : // метки равны
    Code<3> + "MOVE R" + Label<3> +
    ",R" + (Label<3>+1) + Code<1> +
    Code<2> + "R" + Label<1> + ",R" +
    (Label<1>+1).

```

RULE

IntExpr ::= Ident

SEMANTICS

Label<0>=(Left<0>) ? 0 : 1;

Code<0>=(Left<0>) ? Val<1>
 : "MOVE" + Val<1> + "R1".

RULE

IntExpr ::= '(' IntExpr ')'

SEMANTICS

Label<0>=Label<2>;

Code<0>=Code<2>.

RULE

Op ::= '+'

SEMANTICS

Code<0>="ADD".

RULE

Op ::= '-'

SEMANTICS

Code<0>="SUB".

RULE

Op ::= '*'

SEMANTICS

Code<0>="MUL".

RULE

Op ::= '/'

SEMANTICS

Code<0>="DIV".

Команды пересылки требуются для согласования номеров регистров, в которых выполняется операция, с регистрами, в которых должен быть выдан результат. Это имеет смысл, когда эти регистры разные. Получиться это может из-за того, что по приведенной схеме результат выполнения операции всегда находится в регистре с номером метки, а метки левого и правого поддеревьев могут совпадать.

Для выражения $A*B+C*(D+E)$ будет сгенерирован следующий код:

```
MOVE E, R1

ADD D, R1
MUL C, R1
MOVE R1, R2
MOVE B, R1
MUL A, R1
ADD R1, R2
```

В приведенных атрибутных схемах предполагалось, что регистров достаточно для трансляции любого выражения. Если это не так, то приходится усложнять схему трансляции и при необходимости сбрасывать содержимое регистров в память (или магазин).

Реализация этой атрибутной схемы приведена в программном приложении в пакете ToyLang.

9.7. Трансляция логических выражений

Логические выражения, включающие логическое умножение, логическое сложение и отрицание, можно вычислять как непосредственно, используя таблицы истинности, так и с помощью условных выражений, основанных на следующих простых правилах:

A AND B эквивалентно if A then B else False,

A OR B эквивалентно if A then True else B.

Если в выражение могут входить функции с побочным эффектом, то, вообще говоря, результат вычисления зависит от способа вычисления. В некоторых языках программирования не оговаривается, каким способом должны вычисляться логические выражения (например, в Паскале); в некоторых требуется, чтобы вычисления производились конкретным способом (например, в Модуле-2 требуется, чтобы выражения вычислялись по приведенным формулам); в некоторых языках есть возможность явно задать способ вычисления (Си, Ада). Вычисление логических выражений непосредственно по таблицам истинности аналогично вычислению арифметических выражений, поэтому мы не будем их рассматривать отдельно. Рассмотрим подробнее способ вычисления с помощью приведенных выше формул (будем называть его *вычислением*

с условными переходами). Иногда такой способ рассматривают как оптимизацию вычисления логических выражений.

Рассмотрим следующую атрибутивную грамматику со входным языком логических выражений:

```
RULE
Expr ::= BoolExpr
SEMANTICS
FalseLab<1>=False; TrueLab<1>=True;
Code<0>=Code<1>.
```

```
RULE
BoolExpr ::= BoolExpr 'AND' BoolExpr
SEMANTICS
FalseLab<1>=FalseLab<0>; TrueLab<1>=NodeLab<3>;
FalseLab<3>=FalseLab<0>; TrueLab<3>=TrueLab<0>;
Code<0>=NodeLab<0> + ":" + Code<1> + Code<3>.
```

```
RULE
BoolExpr ::= BoolExpr 'OR' BoolExpr
SEMANTICS
FalseLab<1>=NodeLab<3>; TrueLab<1>=TrueLab<0>;
FalseLab<3>=FalseLab<0>; TrueLab<3>=TrueLab<0>;
Code<0>=NodeLab<0> + ":" + Code<1> + Code<3>.
```

```
RULE
BoolExpr ::= F
SEMANTICS
Code<0>=NodeLab<0> + ":" + "GOTO" + FalseLab<0>.
```

```
RULE
BoolExpr ::= T
SEMANTICS
Code<0>=NodeLab<0> + ":" + "GOTO" + TrueLab<0>.
```

Здесь предполагается, что все вершины дерева занумерованы и номер вершины дает атрибут NodeLab. Метки вершин передаются, как это изображено на рис. 9.10.

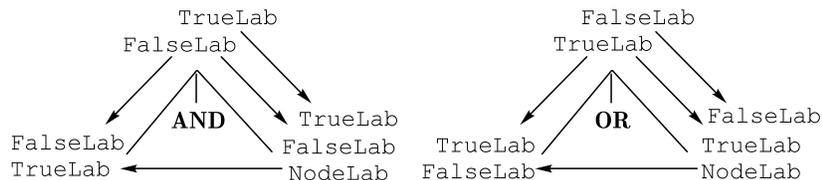


Рис. 9.10

Таким образом, каждому атрибутированному дереву в этой атрибутивной грамматике сопоставляется код, полученный в результате обхода дерева сверху-вниз слева-направо следующим образом. При входе в вершину `BoolExpr` генерируется ее номер, в вершине `F` генерируется текст `GOTO значение атрибута FalseLab<0>`, в вершине `T` — `GOTO значение атрибута TrueLab<0>`. Например, для выражения

`F OR (F AND T AND T) OR T`

получим атрибутированное дерево, изображенное на рис. 9.11, и код

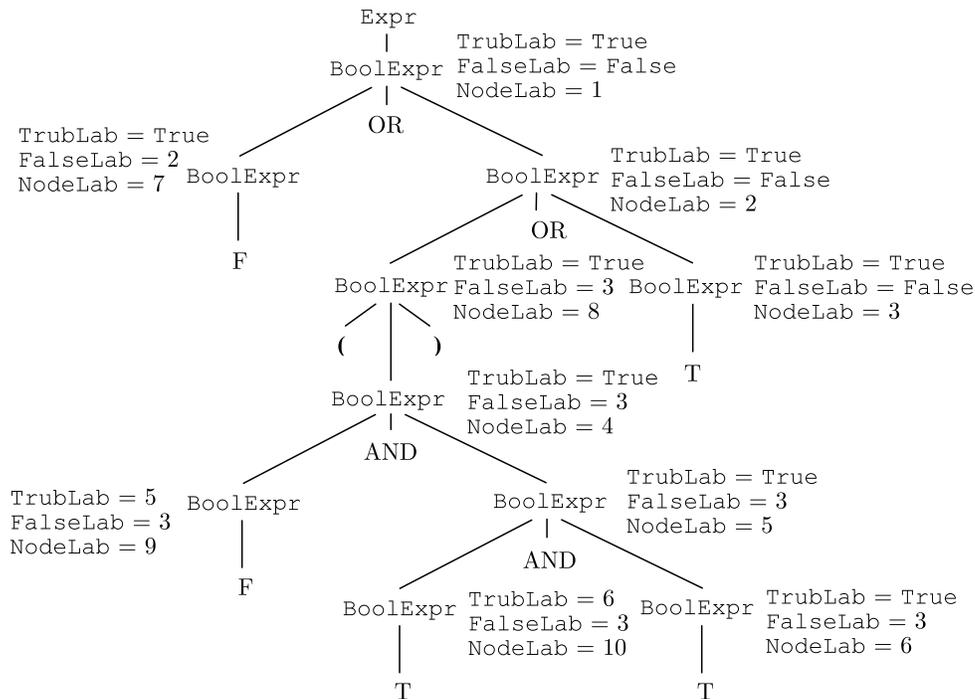


Рис. 9.11

```

1:7:      GOTO 2
2:8:4:9:  GOTO 3
5:10:     GOTO 6
6:        GOTO True
3:        GOTO True
True:     ...
False:    ...

```

Эту линейризованную запись можно трактовать как программу вычисления логического значения: каждая строка может быть помечена номером вершины и содержать либо переход на другую строку, либо переход на `True` или `False`, что соответствует значению выражения `true` или `false`. Будем

говорить, что полученная программа *вычисляет* (или *интерпретирует*) значение выражения, если в результате ее выполнения (от первой строки) мы приходим к строке, содержащей GOTO True или GOTO False.

Утверждение 9.1. В результате интерпретации поддерева с некоторыми значениями атрибутов FalseLab и TrueLab в его корне выполняется команда GOTO TrueLab, если значение выражения истинно, или команда GOTO FalseLab, если значение выражения ложно.

Доказательство. Применим индукцию по высоте дерева. Для деревьев высоты 1, соответствующих правилам

BoolExpr ::= F, BoolExpr ::= T,

справедливость утверждения следует из соответствующих атрибутивных правил. Пусть дерево имеет высоту $n > 1$. Зависимость атрибутов для дизъюнкции и конъюнкции приведена на рис. 9.12.

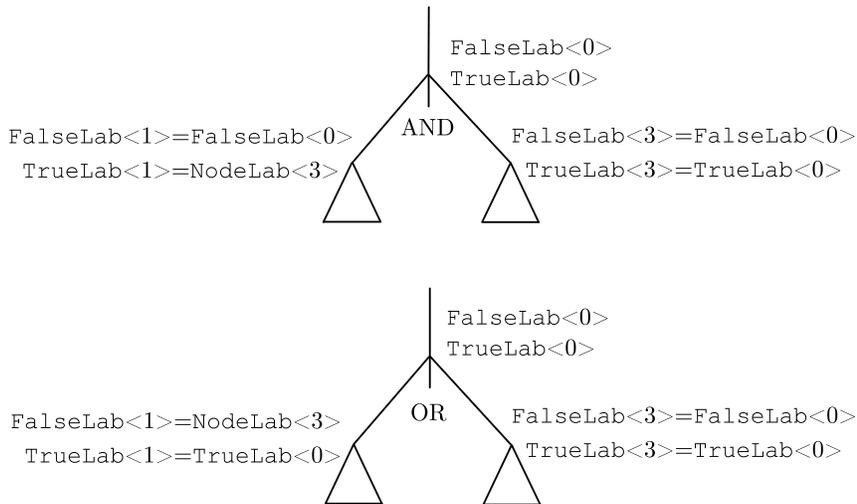


Рис. 9.12

Если для конъюнкции значение левого поддерева ложно и по индукции вычисление левого поддерева завершается командой GOTO FalseLab<1>, то получаем, что вычисление всего дерева завершается командой перехода GOTO FalseLab<0> (= FalseLab<3>). Если же значение левого поддерева истинно, то его вычисление завершается командой перехода GOTO TrueLab<1> (= NodeLab<3>). Если значение правого поддерева ложно, то вычисление всего дерева завершается командой GOTO FalseLab<0> (= FalseLab<3>). Если же оно истинно, вычисление всего дерева завершается командой перехода GOTO TrueLab<0> (= TrueLab<3>). Аналогично — для дизъюнкции. ■

Утверждение 9.2. Для любого логического выражения, состоящего из констант, программа, полученная в результате обхода дерева этого

выражения, завершается со значением логического выражения в обычной интерпретации, т. е. осуществляется переход на True для значения, равного true, и переход на метку False для значения false.

Доказательство. Это утверждение является частным случаем предыдущего. Его справедливость следует из того, что метки корня дерева равны соответственно TrueLab = True и FalseLab = False. ■

Добавим теперь новое правило в предыдущую грамматику:

```

RULE
BoolExpr ::= Ident
SEMANTICS
Code<0>=NodeLab<0> + ":" + "if (" + Val<1> + ")GOTO"
+ TrueLab<0> + "else GOTO" + FalseLab<0>.

```

Тогда, например, для выражения A OR (B AND C AND D) OR E получим следующую программу:

```

1:7:      if (A) GOTO True else GOTO 2
2:8:4:9:  if (B) GOTO 5 else GOTO 3
5:10:     if (C) GOTO 6 else GOTO 3

6:        if (D) GOTO True else GOTO 3
3:        if (E) GOTO True else GOTO False
True:    ...
False:   ...

```

При каждом конкретном наборе данных эта программа превращается в программу вычисления логического значения.

Утверждение 9.3. В каждой строке программы, сформированной предыдущей атрибутивной схемой, одна из меток внутри условного оператора совпадает с меткой следующей строки.

Доказательство. Действительно, по правилам наследования атрибутов TrueLab и FalseLab, в правилах для дизъюнкции и конъюнкции либо атрибут FalseLab, либо атрибут TrueLab принимает значение метки следующего поддерева. Кроме того, как значение FalseLab, так и значение TrueLab, передаются в правое поддерево от предка. Таким образом, самый правый потомок всегда имеет одну из меток TrueLab или FalseLab, равную метке правого брата соответствующего поддерева. Учитывая порядок генерации команд, получаем справедливость утверждения. ■

Дополним теперь атрибутивную грамматику следующим образом:

```

RULE
Expr ::= BoolExpr
SEMANTICS
FalseLab<1>=False; TrueLab<1>=True;
Sign<1>=false;

```

```
Code<0>=Code<1>.
```

```
RULE
```

```
BoolExpr ::= BoolExpr 'AND' BoolExpr
```

```
SEMANTICS
```

```
FalseLab<1>=FalseLab<0>; TrueLab<1>=NodeLab<3>;
```

```
FalseLab<3>=FalseLab<0>; TrueLab<3>=TrueLab<0>;
```

```
Sign<1>=false; Sign<3>=Sign<0>;
```

```
Code<0>=NodeLab<0> + ":" + Code<1> + Code<3>.
```

```
RULE
```

```
BoolExpr ::= BoolExpr 'OR' BoolExpr
```

```
SEMANTICS
```

```
FalseLab<1>=NodeLab<3>; TrueLab<1>=TrueLab<0>;
```

```
FalseLab<3>=FalseLab<0>; TrueLab<3>=TrueLab<0>;
```

```
Sign<1>=true; Sign<3>=Sign<0>;
```

```
Code<0>=NodeLab<0> + ":" + Code<1> + Code<3>.
```

```
RULE
```

```
BoolExpr ::= 'NOT' BoolExpr
```

```
SEMANTICS
```

```
FalseLab<2>=TrueLab<0>; TrueLab<2>=FalseLab<0>;
```

```
Sign<2>=! Sign<0>;
```

```
Code<0>=Code<2>.
```

```
RULE
```

```
BoolExpr ::= F
```

```
SEMANTICS
```

```
Code<0>=NodeLab<0> + ":" + "GOTO" + FalseLab<0>.
```

```
RULE
```

```
BoolExpr ::= T
```

```
SEMANTICS
```

```
Code<0>=NodeLab<0> + ":" + "GOTO" + TrueLab<0>.
```

```
RULE
```

```
BoolExpr ::= Ident
```

```
SEMANTICS
```

```
Code<0>=NodeLab<0> + ":" + "if (" + Val<1> + ") GOTO"
```

```
+ TrueLab<0> + "else GOTO" + FalseLab<0>.
```

Правила наследования атрибута `Sign` приведены на рис. 9.13.

Программу желательно сформировать таким образом, чтобы `else`-метка была как раз меткой следующей вершины. Это можно сделать исходя из следующего утверждения.

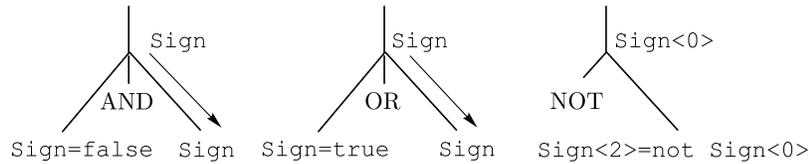


Рис. 9.13

Утверждение 9.4. В каждой терминальной вершине метка ближайшего правого для нее поддерева равна значению атрибута `FalseLab` этой вершины тогда и только тогда, когда значение атрибута `Sign` этой вершины равно `true`, а метка ближайшего правого для нее поддерева равна значению атрибута `TrueLab` этой вершины тогда и только тогда, когда значение атрибута `Sign` равно `false`.

Доказательство. Действительно, если ближайшей общей вершиной является `AND`, то в левого потомка передается `NodeLab` правого потомка в качестве `TrueLab`, причем `Sign` правого потомка равен `true`. Если же ближайшей общей вершиной является `OR`, то в левого потомка передается `NodeLab` правого потомка в качестве `FalseLab`, тогда `Sign` правого потомка равен `false`. Во все же правые потомки значения `TrueLab`, `FalseLab` и `Sign` передаются из предка (за исключением правила для `NOT`, в котором `TrueLab` и `FalseLab` меняются местами, но одновременно меняется на противоположное и значение `Sign`). ■

Эти два утверждения (3 и 4) позволяют заменить последнее правило атрибутной грамматики следующим образом:

```

RULE
BoolExpr ::= Ident
SEMANTICS
Code<0>=NodeLab<0> + ":" +
  (Sign<0>
  ? "if (" + Val<1> + ") GOTO" + TrueLab<0>
  : "if (" + Val<1> + ") GOTO" + FalseLab<0>).

```

В свою очередь, при генерации машинных команд это правило можно заменить на следующее:

```

RULE
BoolExpr ::= Ident
SEMANTICS
Code<0>="TST" + Val<1> +
  (Sign<0>
  ? "BNE" + TrueLab<0>
  : "BEQ" + FalseLab<0>).

```

Таким образом, для выражения `A OR (B AND C AND D) OR E` получим следующий код на командах перехода:

```
1:7:      TST A
          BNE True
2:8:4:9:  TST B
          BEQ 3
5:10:     TST C
          BEQ 3
6:        TST D
          BNE True
3:        TST E
          BEQ False
True:    ...
False:   ...
```

Если элементом логического выражения является сравнение, то генерируется команда, соответствующая знаку сравнения (BEQ для =, BNE для <>, BGE для >= и т. д.), если атрибут *Sign* соответствующей вершины имеет значение *true*, или команда, соответствующая противоположному знаку (BNE для =, BEQ для <>, BLT для >= и т. д.), если атрибут *Sign* имеет значение *false*.

Реализация этой атрибутной схемы приведена в программном приложении в пакете Bool.

9.8. Выделение общих подвыражений

Выделение общих подвыражений относится к области оптимизации программ. В общем случае трудно (или даже невозможно) провести границу между оптимизацией и «качественной трансляцией». Оптимизация — это и есть качественная трансляция. Обычно термин «оптимизация» употребляют, когда для повышения качества программы используют ее глубокие преобразования, например, такие, как перевод в графовую форму для изучения нетривиальных свойств программы.

В этом смысле выделение общих подвыражений — одна из простейших оптимизаций. Для ее осуществления требуется некоторое преобразование программы, а именно, построение ориентированного ациклического графа, о котором говорилось в главе 8.

Линейный участок — это последовательность операторов, в которую управление входит в начале и выходит в конце без остановки и перехода изнутри.

Рассмотрим дерево линейного участка, в котором вершинами служат операции, а потомками — операнды. Будем говорить, что две вершины образуют *общее подвыражение*, если поддеревья для них совпадают, т. е. имеют одинаковую структуру и, соответственно, одинаковые операции во внутренних вершинах и одинаковые операнды в листьях. Выделение общих подвыражений

позволяет генерировать для них код один раз, хотя может привести и к некоторым трудностям, о чем вкратце будет сказано ниже.

Выделение общих подвыражений проводится на линейном участке и основывается на двух положениях.

1. Поскольку на линейном участке переменная может получать несколько присваиваний, то при выделении общих подвыражений необходимо различать вхождения переменных до и после присваивания. Для этого каждая переменная снабжается счетчиком. Вначале счетчики всех переменных устанавливаются равными 0. При каждом присваивании переменной ее счетчик увеличивается на 1.

2. Выделение общих подвыражений осуществляется при обходе дерева выражения снизу-вверх слева-направо. При достижении очередной вершины (пусть операция, примененная в этой вершине, есть бинарная *op*; в случае унарной операции рассуждения те же) просматриваем общие подвыражения, связанные с *op*. Если имеется выражение, связанное с *op* и такое, что его левый операнд есть общее подвыражение с левым операндом нового выражения, а правый операнд - общее подвыражение с правым операндом нового выражения, то объявляем новое выражение общим с найденным и в новом выражении запоминаем указатель на найденное общее выражение. Базисом построения служит переменная: если операндами обоих выражений являются одинаковые переменные с одинаковыми счетчиками, то они являются общими подвыражениями. Если выражение не выделено как общее, то оно заносится в список операций, связанных с *op*.

Рассмотрим теперь реализацию алгоритма выделения общих подвыражений. Поддерживаются следующие глобальные переменные:

Table — таблица переменных; для каждой переменной хранится ее счетчик (**Count**) и указатель на вершину дерева выражений, в которой переменная встретилась в последний раз в правой части (**Last**);

OpTable — таблица списков (типа **ListType**) общих подвыражений, связанных с каждой операцией. Каждый элемент списка хранит указатель на вершину дерева (поле **Addr**) и продолжение списка (поле **List**).

С каждой вершиной дерева выражения связана запись типа **NodeType** со следующими полями:

Left — левый потомок вершины,

Right — правый потомок вершины,

Comm — указатель на предыдущее общее подвыражение,

Flag — признак, указывающий, является ли поддерево общим подвыражением,

Varbl — признак, указывающий, является ли вершина переменной,

VarCount — счетчик переменной.

Выделение общих подвыражений и построение дерева осуществляются приведенными ниже правилами. Атрибут `Entry` нетерминала `Variable` дает указатель на переменную в таблице `Table`. Атрибут `Val` символа `Op` дает код операции. Атрибут `Node` символов `IntExpr` и `Assignment` дает указатель на запись типа `NodeType` соответствующего нетерминала.

```

RULE
Assignment ::= Variable IntExpr
SEMANTICS
Table[Entry<1>].Count=Table[Entry<1>].Count+1.
// Увеличить счетчик присваиваний переменной

RULE
IntExpr ::= Variable
SEMANTICS
if ((Table[Entry<1>].Last!=NULL)
    // Переменная уже была использована
    && (Table[Entry<1>].Last.VarCount
        == Table[Entry<1>].Count ))
    // С тех пор переменной не было присваивания
    {Node<0>.Flag=true;
    // Переменная - общее подвыражение
    Node<0>.Comm= Table[Entry<1>].Last;
    // Указатель на общее подвыражение
    }
else Node<0>.Flag=false;

Table[Entry<1>].Last=Node<0>;
// Указатель на последнее использование переменной
Node<0>.VarCount= Table[Entry<1>].Count;
// Номер использования переменной
Node<0>.Varbl=true.
// Выражение - переменная

RULE
IntExpr ::= Op IntExpr IntExpr
SEMANTICS
ListType L; //Тип списков операции
if ((Node<2>.Flag) && (Node<3>.Flag))
    // И справа, и слева - общие подвыражения
    {L=OpTable[Val<1>];
    // Начало списка общих подвыражений для операции
    while (L!=NULL)

        if ((Node<2>==L.Left)
            && (Node<3>==L.Right))
            // Левое и правое поддеревья совпадают

```

```

        break;
    else L=L.List; // Следующий элемент списка
    }
else L=NULL; // Не общее подвыражение

Node<0>.Varbl=false; // Не переменная
Node<0>.Comm=L;
// Указатель на предыдущее общее подвыражение
// или NULL

if (L!=NULL)
{Node<0>.Flag=true; //Общее подвыражение
Node<0>.Left=Node<2>;
// Указатель на левое поддерево
Node<0>.Right=Node<3>;
// Указатель на правое поддерево
}
else {Node<0>.Flag=false;
// Данное выражение не может рассматриваться
// как общее. Если общего подвыражения с
// данным не было, включить данное в список
// для операции
L=new LisType();
L.Addr=Node<0>;
L.List=OpTable[Val<1>];
OpTable[Val<1>]=L;
}.

```

Рассмотрим теперь некоторые простые правила распределения регистров при наличии общих подвыражений. Если число регистров ограничено, можно выбрать один из следующих двух вариантов.

1. При обнаружении подвыражения, общего с подвыражением в уже просмотренной части дерева (и, значит, с уже распределенными регистрами) проверяем, расположено ли его значение на регистре. Если да и если регистр после этого не менялся, то заменяем вычисление поддерева на значение в регистре. Если регистр менялся, то вычисляем подвыражение заново.

2. Вводим еще один проход. На первом проходе распределяем регистры. Если в некоторой вершине обнаруживается, что ее поддерево — общее с уже вычисленным ранее, но значение регистра утрачено, то в такой вершине на втором проходе необходимо сгенерировать команду сброса регистра в рабочую память. Выигрыш в коде будет, если стоимость команды сброса регистра + доступ к памяти в повторном использовании этой памяти не превосходит стоимости заменяемого поддерева. Поскольку стоимость команды MOVE известна, можно сравнить стоимости и принять оптимальное решение: пометить предыдущую вершину для сброса либо вычислять поддерево полностью.

9.9. Трансляция объектно-ориентированных свойств языков программирования

В этом разделе будут рассмотрены механизмы трансляции базовых конструкций объектно-ориентированных языков программирования, а именно, наследования и виртуальных функций на примере языка C++.

9.9.1. Виртуальные базовые классы. К описателю базового класса можно добавить ключевое слово *virtual*. В этом случае единственный под-объект виртуального базового класса разделяется каждым базовым классом, в котором тот (исходный) базовый класс определен как виртуальный.

Пусть мы имеем следующую иерархию наследования:

```
class L {...}
class A : public virtual L {...}
class B : public virtual L {...}
class C : public A, public B {...}
```

Это можно изобразить диаграммой классов, изображенной на рис. 9.14.

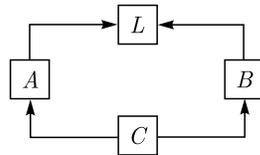


Рис. 9.14 Диаграмма классов

Каждый из объектов *A* и *B* будет содержать *L*, но в объекте *C* будет существовать лишь один объект класса *L*. Ясно, что представление объекта виртуального базового класса *L* не может быть в одной и той же позиции относительно и *A*, и *B* для всех объектов. Следовательно, во всех объектах классов, которые включают класс *L* как виртуальный базовый класс, должен храниться указатель на *L*. Реализация объектов *A*, *B* и *C* могла бы выглядеть так, как показано на рис. 9.15.

```

A* aptr = new A;  ┌───▶ Часть A
                  │
                  └───▶ Часть L

B* bptr = new B;  ┌───▶ Часть B
                  │
                  └───▶ Часть L

C* cptr = new C;  ┌───▶ Часть A
                  │
                  ├───▶ Часть B
                  │
                  ├───▶ Часть C
                  └───▶ Часть L
  
```

Рис. 9.15 Реализация объектов *A*, *B* и *C*

9.9.2. Множественное наследование. Имея два класса

```
class A {...af (int);},
class B {...bf (int);},
```

можно объявить третий класс с этими двумя в качестве базовых:

```
class C : public A, public B {...}.
```

Объект класса *C* может быть размещен как непрерывный объект (рис. 9.16).



Рис. 9.16

Как и в случае с единичным наследованием, здесь не гарантируется порядок выделения памяти для базовых классов, поэтому объект класса *C* может выглядеть и так, как показано на рис. 9.17.



Рис. 9.17

Доступ к члену класса *A*, *B* или *C* реализуется в точности так же, как и для единичного наследования: компилятор знает положение каждого члена в объекте и порождает соответствующий код.

Если объект размещен в памяти согласно первой диаграмме: сначала часть *A* объекта, а затем части *B* и *C*, то вызов функции-члена класса *A* или *C* будет таким же, как вызов функции-члена при единичном наследовании. Вызов функции-члена класса *B* для объекта, заданного указателем на *C*, реализуется несколько сложнее. Рассмотрим следующий пример:

```
C* pc = new C;
pc → bf(2);
```

Функция *B :: bf()* естественным образом предполагает, что ее параметр *this* является указателем на *B*. Чтобы получить указатель на часть *B* объекта *C*, следует добавить к указателю *pc* смещение *B* относительно *C* - константу времени компиляции, которую мы будем называть *delta(B)*. Соотношение указателя *pc* и указателя *this*, передаваемого в *B::bf*, показано на рис. 9.18.

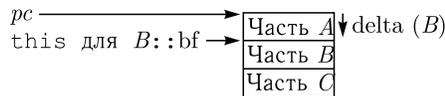


Рис. 9.18

9.9.3. Единичное наследование и виртуальные функции. Если класс *base* содержит виртуальную функцию *vf*, а класс *derived*, порожденный по классу *base*, также содержит функцию *vf* того же типа, то обращение к *vf* для объекта класса *derived* вызывает *derived::vf* даже при доступе через указатель или ссылку на *base*. В таком случае говорят, что функция производного класса подменяет (*override*) функцию базового класса. Если, однако, типы этих функций различны, то функции считаются различными и механизм виртуальности не включается.

Виртуальные функции можно реализовать при помощи таблицы указателей на виртуальные функции *vtbl*. В случае единичного наследования таблица виртуальных функций класса будет содержать ссылки на соответствующие функции, а каждый объект данного класса будет содержать указатель на таблицу *vtbl*.

```
class A {
public:
    int a;
    virtual void f(int);
    virtual void g(int);
    virtual void h(int);
};
class B : public A {
public:
    int b;
    void g(int);
};
class C : public B {
public:
    int c;
    void h(int);
};
```

Объект класса *C* будет выглядеть примерно так, как показано на рис. 9.19.

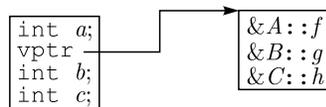


Рис. 9.19

9.9.4. Множественное наследование и виртуальные функции. При множественном наследовании виртуальные функции реализуются несколько сложнее. Рассмотрим следующие объявления:

```

class A {
    public:
        virtual void f(int);
};
class B : {
    public:
        virtual void f(int);
        virtual void g(int);
};
class C : public A, public B {
    public:
        void f();
};

```

Поскольку класс A порожден по классам A и B , каждый из следующих вызовов будет обращаться к $C :: f()$ (при условии, что каждый из трех указателей смотрит на объект класса C):

```

pa → f()
pb → f()
pc → f()

```

Рассмотрим для примера вызов $pb \rightarrow f()$. При входе в $C :: f$ указатель $this$ должен указывать на начало объекта C , а не на часть B в нем. Во время компиляции, вообще говоря, не известно, указывает ли pb на часть B в C , — например, из-за того, что идентификатору pb может быть присвоен просто указатель на объект B . Так что величина $delta(B)$, упомянутая выше, может быть различной для разных объектов в зависимости от структуры классов, порождаемых из B , и должна где-то храниться во время выполнения.

Поскольку это смещение нужно только для виртуального вызова функции, логично хранить его в таблице виртуальных функций.

Указатель $this$, передаваемый виртуальной функции, может быть вычислен путем вычитания смещения объекта, для которого была определена виртуальная функция, из смещения объекта, для которого она вызвана, а затем вычитания этой разности из указателя, используемого при вызове. Здесь значение $delta(B)$ будет необходимо для поиска начала объекта (в нашем случае C), содержащего B , по указателю на B . Сгенерированный код вычтет значение $delta(B)$ из значения указателя, так что хранится смещение со знаком минус, $-\delta(B)$. Объект класса C будет выглядеть так, как показано на рис. 9.20.

Таблица виртуальных функций $vtbl$ для B в C отличается от $vtbl$ для отдельно размещенного B . Каждая комбинация базового и производного классов имеет свою таблицу $vtbl$. В общем случае объекту производного

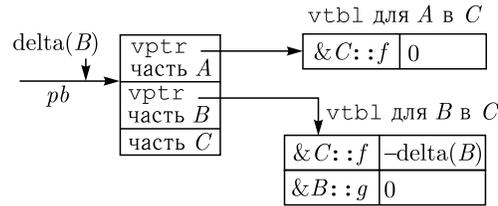


Рис. 9.20

классу требуется таблица *vtbl* для каждого базового класса плюс таблица для производного класса, не считая того, что производный класс может разделять таблицу *vtbl* со своим первым базовым классом. Таким образом, для объекта типа *C* в этом примере требуются две таблицы *vtbl* (таблица для *A* в *C* объединена с таблицей для объекта *C*, и еще одна таблица нужна для объекта *B* в *C*).

9.9.5. Виртуальные базовые классы с виртуальными функциями.

При наличии виртуальных базовых классов построение таблиц для вызовов виртуальных функций становится более сложным. Рассмотрим следующие объявления:

```
class W {
public:
    virtual void f();
    virtual void g();
    virtual void h();
    virtual void k();
};
class MW : public virtual W {
public:
    void g();
};
class BW : public virtual W {
public:
    void f();
};
class BMW : public BW, public MW, public virtual W {
public:
    void h();
};
```

Отношение наследования для этого примера может быть изображено в виде ациклического графа (рис. 9.21).

Функции-члены класса *BMW* могут использоваться, например, так:

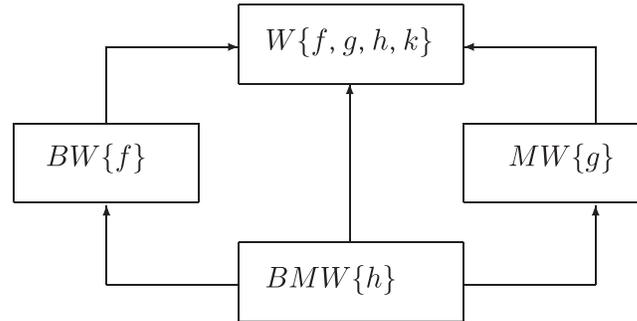


Рис. 9.21

```

void g(BMW*pbmw)
{pbmw → f(); // вызывает BW :: f()
 pbmw → g(); // вызывает MW :: g()
 pbmw → h(); // вызывает BMW :: h()
}

```

Рассмотрим теперь следующий вызов виртуальной функции $f()$:

```

void h(BMW*pbmw)
{MW*pmw = pbmw;
 pmw → f(); // вызывает BW :: f(), потому что
 // pmw указывает на BMW, для которого f берется из BW!
}

```

Виртуальный вызов функции по одному пути в структуре наследования может породить обращение к функции, переопределенной на другом пути.

Структуры объектов класса BMW и его таблиц виртуальных функций $vtbl$ могут выглядеть так, как показано на рис. 9.22.

Виртуальной функции должен быть передан указатель $this$ на объект класса, в котором эта функция описана. Поэтому следует хранить смещение для каждого указателя функции из $vtbl$. Когда объект размещен в памяти так, как это изображено на рис. 9.22, смещение, хранимое с указателем виртуальной функции, исчисляется вычитанием смещения класса, для которого эта таблица $vtbl$ создана, из смещения класса, поставляющего эту функцию. Рассмотрим пример:

```

void callvirt(w*pw)
{ pw → f();
}
main ()
{ callvirt(new BMW);
}

```

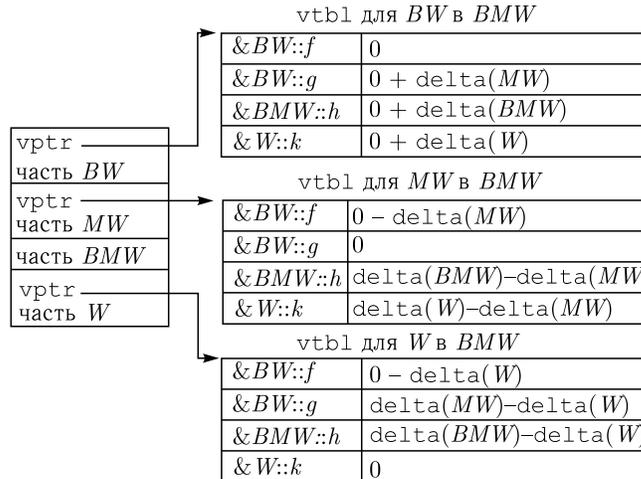


Рис. 9.22

В функции *main* вызов *callvirt* с указателем на *BMW* требует приведения к указателю на *W*, поскольку функция *callvirt* ожидает параметр типа *W**. Так как функция *callvirt* вызывает *f()* (через указатель на *BMW*, преобразованный к указателю на *W*), будет использована таблица *vtbl* класса *W* (в *BMW*), где указано, что экземпляром виртуальной функции *f()*, которую нужно вызвать, является *BW :: f()*. Чтобы передать функции *BW :: f()* указатель *this* на *BW*, указатель *pw* должен быть вновь приведен к указателю на *BMW* (вычитанием смещения для *W*), а затем к указателю на *BW* (добавлением смещения *BW* в объекте *BMW*). Значение смещения *BW* в объекте *BMW* минус смещение *W* в объекте *BMW* и есть смещение, хранимое в строке таблицы *vtbl* для *w* в *BMW* для функции *BW :: f()*.

9.10. Генерация оптимального кода методами сопоставления образцов

9.10.1. Сопоставление образцов. Техника генерации кода, рассмотренная выше, основывалась на однозначном соответствии структуры промежуточного представления с описывающей это представление грамматикой. Недостатком такого «жесткого» подхода является то, что, как правило, одну и ту же программу на промежуточном языке можно реализовать многими различными способами в системе команд машины. Эти разные реализации могут иметь различную длину, время выполнения и другие характеристики. Для генерации более качественного кода применим подход, изложенный в настоящей главе.

Этот подход основан на понятии «сопоставления образцов»: командам машины сопоставляются некоторые «образцы», вхождения которых ищутся

в промежуточном представлении программы, и делается попытка «покрыть» промежуточную программу такими образцами. Если это удастся, то по образцам восстанавливается программа уже в кодах. Каждое такое покрытие соответствует некоторой программе, реализующей одно и то же промежуточное представление.

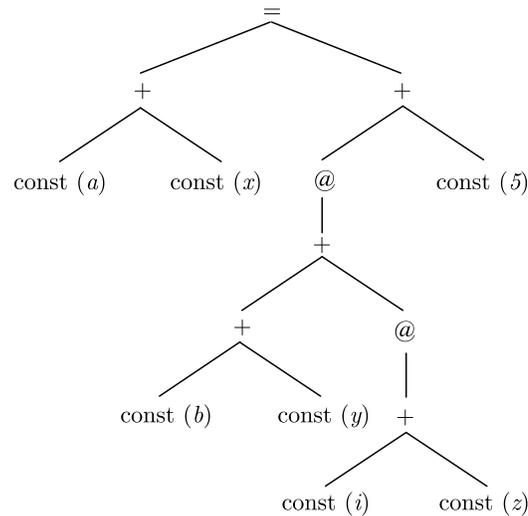


Рис. 9.23

На рис. 9.23 показано промежуточное дерево для оператора $a = b[i] + 5$, где a , b , i — локальные переменные, хранимые со смещениями x , y , z соответственно в областях данных с одноименными адресами.

Элемент массива b занимает память в одну машинную единицу. 0-местная операция `const` возвращает значение атрибута соответствующей вершины промежуточного дерева, указанного на рисунке в скобках после оператора. Одноместная операция `@` означает косвенную адресацию и возвращает содержимое регистра или ячейки памяти, имеющей адрес, задаваемый аргументом операции.

В табл. 9.2 показан пример сопоставления образцов машинным командам. Приведены два варианта задания образца: в виде дерева и в виде правила контекстно-свободной грамматики. Для каждого образца указана машинная команда, реализующая этот образец, и стоимость этой команды.

В каждом дереве-образце корень или лист может быть помечен терминальным и/или нетерминальным символом. Внутренние вершины помечены терминальными символами — знаками операций. При наложении образца на дерево выражения, во-первых, терминальный символ образца должен соответствовать терминальному символу дерева и, во-вторых, образцы должны «склеиваться» по типу нетерминального символа, т.е. тип корня образца

должен совпадать с типом вершины, в которую образец подставляется корнем. Допускается использование «цепных» образцов, т. е. образцов, корню которых не соответствует терминальный символ и которые имеют единственный элемент в правой части. Цепные правила служат для приведения вершин к одному типу. Например, в рассматриваемой системе команд одни и те же регистры используются как для целей адресации, так и для вычислений. Если бы в системе команд для этих целей использовались разные группы регистров, то в грамматике команд могли бы использоваться разные нетерминалы, а для пересылки из адресного регистра в регистр данных могли бы использоваться соответствующая команда и образец.

Таблица 9.2

№	Образец	Правило грамматики	Команда/стоимость	
1	$\begin{array}{c} \text{Reg} \\ \\ \text{const} \end{array}$	$\text{Reg} \rightarrow \text{const}$	MOVE #const, R	2
2	$\begin{array}{c} = \text{Stat} \\ / \quad \backslash \\ + \quad \text{Reg}(j) \\ / \quad \backslash \\ \text{Reg}(i) \quad \text{const} \end{array}$	$\text{Stat} \rightarrow '=' '+' \text{Reg const Reg}$	MOVE Rj, const (Ri)	4
3	$\begin{array}{c} @ \quad \text{Reg}(j) \\ \\ + \\ / \quad \backslash \\ \text{Reg}(i) \quad \text{const} \end{array}$	$\text{Reg} \rightarrow '@' '+' \text{Reg const}$	MOVE const (Ri), Rj	4
4	$\begin{array}{c} + \quad \text{Reg} \\ / \quad \backslash \\ \text{Reg} \quad \text{const} \end{array}$	$\text{Reg} \rightarrow '+' \text{Reg const}$	ADD #const, R	3
5	$\begin{array}{c} + \quad \text{Reg}(i) \\ / \quad \backslash \\ \text{Reg}(i) \quad \text{Reg}(j) \end{array}$	$\text{Reg} \rightarrow '+' \text{Reg Reg}$	ADD Rj, Ri	2
6	$\begin{array}{c} + \quad \text{Reg}(i) \\ / \quad \backslash \\ \text{Reg}(i) \quad @ \\ \\ + \\ / \quad \backslash \\ \text{Reg}(j) \quad \text{const} \end{array}$	$\text{Reg} \rightarrow '+' \text{Reg '@' '+' \text{Reg const}$	ADD const (Rj), Ri	4
7	$\begin{array}{c} @ \quad \text{Reg}(i) \\ \\ \text{Reg}(j) \end{array}$	$\text{Reg} \rightarrow '@' \text{Reg}$	MOVE (Rj), Ri	2

Нетерминалы *Reg* на образцах могут быть помечены индексом (*i* или *j*), что (неформально) соответствует номеру регистра и служит лишь для пояснения смысла использования регистров. Отметим, что при генерации кода рассматриваемым методом не осуществляется распределение регистров. Это является отдельной задачей.

Стоимость может определяться различными способами, например, числом обращений к памяти при выборке и исполнении команды. Здесь мы не рассматриваем этого вопроса. На рис. 9.24 приведен пример покрытия промежуточного дерева рис. 9.23 образцами из табл. 9.2. В рамки заключены фрагменты дерева, сопоставленные образцу правила, номер которого указывается в левом верхнем углу рамки. В квадратных скобках указаны результирующие вершины.

Приведенное покрытие дает такую последовательность команд:

```
1  MOVE #a, Ra
1  MOVE #b, Rb
4  ADD  #y, Rb
1  MOVE #i, Ri
6  ADD  #z (Ri), Rb
7  MOVE (Rb), Rb
4  ADD  #5, Rb
2  MOVE Rb, #x (Ra)
```

Основная идея подхода заключается в том, что каждая команда машины описывается в виде такого образца. Различные покрытия дерева промежуточного представления соответствуют различным последовательностям машинных команд. Задача выбора команд состоит в том, чтобы выбрать наилучший способ реализации того или иного действия или последовательности действий, т. е. выбрать оптимальное (в некотором смысле) покрытие.

Для выбора оптимального покрытия было предложено несколько интересных алгоритмов, в частности, использующих динамическое программирование [14, 18]. Мы здесь рассмотрим алгоритм, предложенный в [16].

9.10.2. Построение покрытия. В качестве промежуточного представления будем использовать абстрактное дерево программы. Абстрактное дерево в качестве внутренних вершин содержит операции. В нашем случае при построении абстрактного дерева будем исходить из следующих предположений:

- 1) каждая переменная программы адресуется суммой двух констант — адреса области памяти и смещения в этой области;
- 2) для выборки значения по адресу используется операция @ — косвенная адресация;
- 3) все массивы одномерные, причем размер элемента массива — одна адресуемая единица.

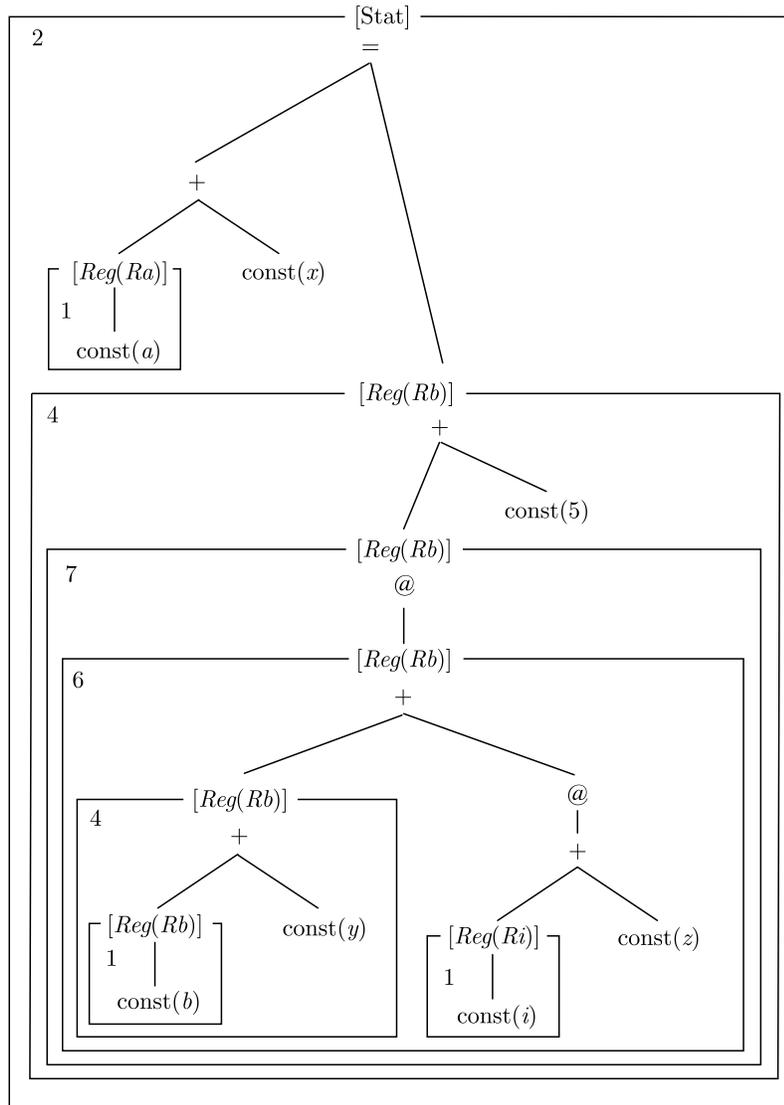


Рис. 9.24

Корень каждого образца помечен символом, который мы будем называть *нетерминальным* (или *нетерминалом*). Смысл использования этого термина будет объяснен ниже при использовании синтаксического анализа для построения покрытия. Некоторые листья образцов также могут быть помечены нетерминалами. Содержательно можно считать, что нетерминал задает тип результата, вычисляемого поддеревом. В нашем случае нетерминалов всего два: *Reg* говорит о том, что результат должен размещаться в регистре, *Stat* — что результат не вырабатывается.

Задача построения покрытия заключается в таком наложении образцов на абстрактное дерево программы, что:

- 1) всё дерево покрыто набором образцов;
- 2) накрывающие образцы не пересекаются;
- 3) во внутренних вершинах образцы «склеиваются» по нетерминалам, т. е. тип вершины-корня образца, примененного в вершине n абстрактного дерева, совпадает с типом вершины образца, для которого вершина n является листом.

Для согласования типов вершин могут использоваться «цепные» образцы, т. е. образцы типа Нетерминал \rightarrow Нетерминал. Цепные образцы не зависят от операций, следовательно, их необходимо проверять отдельно. Применение одного цепного образца может зависеть от применения другого цепного образца. Следовательно, применение цепных образцов необходимо проверять до тех пор, пока можно применить хотя бы один из цепных образцов. Мы предполагаем, что в грамматике нет циклов в применении цепных образцов.

Построение всех вариантов покрытия дано ниже в алгоритме 9.1. В этом алгоритме дерево выражения обходится сверху-вниз, и в нем ищутся поддеревья, сопоставимые с образцами. Обход дерева осуществляется процедурой `PARSE`. После обхода поддерева данной вершины в ней применяется процедура `MATCHED`, которая пытается найти все образцы, сопоставимые поддереву данной вершины. С этой целью для каждого образца осуществляется рекурсивный обход поддерева. В результате обхода дерева снизу-вверх для каждой вершины n (как внутренней, так и листа) для каждого нетерминала находятся все покрытия поддерева с корнем n , помеченные этим нетерминалом (это множество может быть и пустым). Точнее, для вершины n находится множество образцов, применимых в этой вершине с учетом найденных образцов в вершинах дерева — потомках n , соответствующих листьям-нетерминалам образца, и цепных правил, примененных, в частности, к терминальным листьям образца (что соответствует листьям дерева программы). Для последующего отбора образцов при поиске оптимального покрытия и генерации оптимального кода используется принцип динамического программирования: из всех возможных способов покрытия поддерева, соответствующих данному нетерминалу, выбирается тот, который дает минимальную суммарную стоимость. В данном случае это возможно, поскольку стоимость полного покрытия определяется как сумма стоимостей его составляющих.

Для реализации этого принципа каждой вершине ставится в соответствие множество `HashMap поTerm`, хранящее для каждого нетерминала минимальную стоимость покрытия поддерева для этой вершины, обеспечивающего соответствие нетерминала корню образца, и образец, дающий эту

минимальную стоимость (opt OptChoice). Список sons содержит потомков данной вершины, HashMap rules содержит множество проавил для каждого нетерминала. HashMap invert дает для каждого нетерминала множество левых частей для цепных правил с данной правой частью.

Алгоритм 9.1

```

class Tnode {
    Integer op;
    ArrayList sons;
    HashMap nonTerm; // Для каждого нетерминала OptChoice
}

class Pattern {
    Integer op;
    ArrayList sons; // Список Pattern
}

class Rule {
    int cost;
    Pattern pat;
}

class OptChoice {
    int cost;
    Rule rule;
}

HashMap invert; // Отображение нетерминала во множество
                // цепных правил с данной правой частью

HashSet NTerms; // Нетерминалы

HashMap rules; // Отображение нетерминала во множество правил

HashSet operations;

/** ***** */

void Parse(Tnode n) {
    int i;
    Integer nextNT;
    OptChoice opt;
    Rule nextRule;
    Iterator NTiter = NTerms.iterator();
    while (NTiter.hasNext()) { // Для каждого нетерминала
        nextNT = (Integer) NTiter.next();
    }
}

```

```

    opt = (OptChoice) n.nonTerm.get(nextNT);
    opt.cost = maxCost;
    opt.rule = undef;
    HashSet currentRules = (HashSet) rules.get(nextNT);
    // Множество правил для нетерминала
    Iterator ruleIter = currentRules.iterator();
    while (ruleIter.hasNext()) {
        nextRule = (Rule) ruleIter.next();
        if (nextRule.pat.op == n.op) {
            int cost = Matched(n, nextRule);
            if (cost != maxCost)
                if (opt.cost > cost) {
                    opt.cost = cost;
                    opt.rule = nextRule;
                    n.nonTerm.put(nextNT, opt);
                }
        }
    }
    boolean stop = false;
    while (!stop) {
        stop = true;
        NTiter = NTerms.iterator();
        while (NTiter.hasNext()) {
            nextNT = (Integer) NTiter.next();
            opt = (OptChoice) n.nonTerm.get(nextNT);
            if (opt.rule != undef) {
                HashSet inv = (HashSet) invert.get(nextNT);
                Iterator invIter = inv.iterator();
                while (invIter.hasNext()) {
                    nextRule = (Rule) invIter.next();
                    opt = (OptChoice) n.nonTerm.get(nextNT);
                    if (nextRule.cost < opt.cost) {
                        opt.cost = nextRule.cost;
                        opt.rule = nextRule;
                        n.nonTerm.put(nextNT, opt);
                    }
                }
            }
        }
    }
}

/** ***** */

int Matched(Tnode n, Rule r) {

```

```

int m, i;
if (operations.contains(r.pat.op))
    // r.pat.op может содержать операции и нетерминалы
    if (r.pat.op == n.op)
        if ((m = Arity(r.pat.op)) == 0)
            // нуль-местная операция, например, const
            return 0;
        else { // внутренняя вершина
            int cost = 0, costSum = r.cost;
            for (i = 0; i < m; i++) {
cost = Matched((Tnode) n.sons.get(i), (Rule) r.pat.sons.get(i));
                if (cost != maxCost)
                    costSum += cost;
                else {
                    costSum = maxCost;
                    break;
                }
            }
            return costSum;
        }
    else
        return maxCost;
else
    // Нетерминал
    return ((OptChoice) n.nonTerm.get(r.pat.op)).cost;
}

```

Объем вычислений пропорционален объему дерева, умноженному на число нетерминалов и числу правил для каждого нетерминала, т.е. имеет временную и емкостную сложности порядка $O(n)$.

Рассмотрим вновь пример рис. 9.23.

В табл. 9.3 приведен результат работы алгоритма.

9.10.3. Выбор дерева вывода наименьшей стоимости. *Дерево наименьшей стоимости* определяется как дерево, соответствующее минимальной стоимости покрытия. Дереву наименьшей стоимости соответствует набор команд, соответствующих образцам, реализующих это дерево. Выбор оптимального покрытия и соответствующих команд осуществляется в процессе обхода дерева: при движении сверху-вниз выбирается оптимальное покрытие, при движении снизу-вверх осуществляется генерация команд. Обход дерева реализуется двумя взаимно рекурсивными функциями: функция CodeGen(Tnode n, NonTerm NT) осуществляет обход поддерева с вершиной n и (оптимальным) образцом, примененным в этой вершине для нетерминала NT. Функция TreePass(Tnode n, Rule r) осуществляет обход поддерева с вершиной n согласно образцу r.

Таблица 9.3

Операция	Длина	Правила (стоимость)
=	15	2 (22)
+	3	4 (5) 5 (6)
a	1	1 (2)
x	1	1 (2)
+	11	4 (16) 5 (17)
@	9	7 (11)
+	8	5 (13) 6 (11)
+	3	4 (5) 5 (6)
b	1	1 (2)
y	1	1 (2)
@	4	7 (7) 3 (6)
+	3	4 (5) 5 (6)
i	1	1 (2)
z	1	1 (2)
5	1	1 (2)

```

void CodeGen(Tnode n, Integer NT) {
    Rule rule = ((OptChoice) n.nonTerm.get(NT)).rule;
    // Выбираем оптимальное правило для NT
    TreePass(n, rule);
    GenInstr(rule);
}

/** ***** */

void TreePass(Tnode n, Rule r) {
    int m, i;
    if (operations.contains(r.pat.op))
        // Проверка: операция или нетерминал
        if ((m = Arity(r.pat.op)) != 0)
            for (i = 0; i < m; i++)
                TreePass((Tnode) n.sons.get(i), (Rule) r.pat.sons.get(i));
    else
        CodeGen(n, r.pat.op);
}

```

9.10.4. Синтаксический анализ для T-грамматик. Обычно код генерируется из некоторого промежуточного языка с довольно жесткой структурой. В частности, для каждой операции известна ее размерность, т.е. число операндов, большее или равное 0. Операции задаются терминальными символами, и наоборот — будем считать все терминальные символы знаками операций. Назовем грамматики, удовлетворяющие этим ограничениям, *T-грамматиками*. Правая часть каждой продукции в T-грамматике есть правильное префиксное выражение, которое может быть задано следующим определением.

- 1) Операция размерности 0 является правильным префиксным выражением.
- 2) Нетерминал является правильным префиксным выражением.
- 3) Префиксное выражение, начинающееся со знака операции размерности $n > 0$, является правильным, если после знака операции следует n правильных префиксных выражений.
- 4) Ничто другое не является правильным префиксным выражением.

Длины всех выражений из входной цепочки $a_1 \dots a_n$ можно предварительно вычислить (под длиной выражения имеется в виду длина подстроки, начинающейся с символа кода операции и заканчивающейся последним символом, входящим в выражение для этой операции). Поэтому можно проверить, сопоставимо ли некоторое правило с подцепочкой $a_i \dots a_k$ входной цепочки $a_1 \dots a_n$, проходя слева-направо по $a_i \dots a_k$. В процессе прохода по цепочке предварительно вычисленные длины префиксных выражений используются для того, чтобы перейти от одного терминала к следующему терминалу, пропуская подцепочки, соответствующие нетерминалам правой части правила.

Рассмотрим вновь пример рис. 9.23. В префиксной записи приведенный фрагмент программы записывается следующим образом:

```
= + a x + @ + + b y @ + i z 5
```

Образцы, соответствующие машинным командам, задаются правилами грамматики (вообще говоря, неоднозначной). Генератор кода анализирует входное префиксное выражение и строит одновременно все возможные деревья разбора, например, с помощью алгоритма Кока–Янгера–Касами. После окончания разбора выбирается дерево с наименьшей стоимостью. Затем по этому единственному оптимальному дереву генерируется код.

9.10.5. Атрибутная схема для алгоритма сопоставления образцов. Алгоритмы 9.1 и 9.2 являются «универсальными» в том смысле, что конкретные грамматики выражений и образцов являются, по существу, параметрами этих алгоритмов. В то же время для каждой конкретной грамматики можно написать свой алгоритм поиска образцов. Например, в случае нашей

грамматики выражений и образцов, приведенных в табл. 9.2, алгоритм 9.2 может быть представлен атрибутивной грамматикой, приведенной ниже.

Наследуемый атрибут `Match` содержит упорядоченный список (вектор) образцов для сопоставления в поддереве данной вершины. Каждый из образцов либо имеет вид $\langle op\ op\text{-}list \rangle$ (op — операция в данной вершине, а $op\text{-}list$ — список ее операндов), либо представляет собой нетерминал N . В первом случае $op\text{-}list$ «распределяется» по потомкам вершины для дальнейшего сопоставления. Во втором случае сопоставление считается успешным, если есть правило $N \rightarrow op\ \{Pat_i\}$, где w состоит из образцов, успешно сопоставленных потомкам данной вершины. В этом случае по потомкам в качестве образцов распределяются элементы правой части правила. Эти два множества образцов могут пересекаться. Синтезируемый атрибут `Pattern` — вектор логических значений — дает результат сопоставления по вектору-образцу `Match`.

Таким образом, при сопоставлении образцов могут встретиться два случая.

1. Вектор образцов содержит образец $\langle op\ \{Pat_i\} \rangle$, где op — операция, примененная в данной вершине. Тогда распределяем образцы Pat_i по потомкам и считаем сопоставление по данному образцу успешным (истинным), если успешны сопоставления элементов этого образца по всем потомкам.
2. Образцом является нетерминал N . Тогда рассматриваем все правила вида $N \rightarrow op\ \{Pat_i\}$. Вновь распределяем образцы Pat_i по потомкам и считаем сопоставление успешным (истинным), если успешны сопоставления по всем потомкам. В общем случае успешным может быть сопоставление по нескольким образцам.

Отметим, что в общем случае по потомкам одновременно распределяются несколько образцов для сопоставления.

В приведенной ниже атрибутивной схеме не рассматриваются правила выбора покрытия наименьшей стоимости (см. подраздел 9.10.3). Выбор оптимального покрытия может быть сделан еще одним проходом по дереву аналогично тому, как это было сделано выше. Например, в правиле с '+' имеются несколько образцов для `Reg`, но реального выбора одного из них не осуществляется. Кроме того, не уточнены некоторые детали реализации, в частности, конкретный способ формирования векторов `Match` и `Pattern`. В тексте употребляется термин «добавить», что означает добавление очередного элемента к вектору образцов. Векторы образцов записаны в угловых скобках.

RULE

Stat ::= '=' Reg Reg

SEMANTICS

```
Match<2>=<'+' Reg Const>;
Match<3>=<Reg>;
Pattern<0>[1]=Pattern<2>[1]&Pattern<3>[1].
```

Этому правилу соответствует один образец 2. Поэтому в качестве образцов потомков через их атрибуты Match передаются, соответственно,

<'+' Reg Const> и <Reg>.

RULE

```
Reg ::= '+' Reg Reg
```

SEMANTICS

```
if (Match<0> содержит Reg в позиции i)
  {Match<2>=<Reg, Reg, Reg>;
   Match<3>=<Const, Reg, '@' '+' Reg Const>>;
  }
if (Match<0> содержит образец <'+' Reg Const>
    в позиции j)
  {добавить Reg к Match<2> в некоторой позиции k;
   добавить Const к Match<3> в некоторой позиции k;
  }
if (Match<0> содержит образец <'+' Reg Const>
    в позиции j)
  Pattern<0>[j]=Pattern<2>[k]&Pattern<3>[k];
if (Match[0] содержит Reg в i-й позиции)
  Pattern<0>[i]=(Pattern<2>[1]&Pattern<3>[1])
               |(Pattern<2>[2]&Pattern<3>[2])
               |(Pattern<2>[3]&Pattern<3>[3]).
```

Образцы, соответствующие этому правилу, следующие:

- (4) $Reg \rightarrow '+' Reg Const$,
- (5) $Reg \rightarrow '+' Reg Reg$,
- (6) $Reg \rightarrow '+' Reg '@' '+' Reg Const$.

Атрибутам Match второго и третьего символов в качестве образцов при сопоставлении могут быть переданы векторы <Reg, Reg, Reg> и <Const, Reg, '@' '+' Reg Const>> соответственно. Из анализа других правил можно заключить, что при сопоставлении образцов предков левой части данного правила атрибуту Match символа левой части может быть передан образец <'+' Reg Const> (из образцов 2, 3, 6) или образец Reg.

RULE

```
Reg ::= '@' Reg
```

SEMANTICS

```
if (Match<0> содержит Reg в i-й позиции)
  Match<2>=<<'+' Reg Const>, Reg>;
if (Match<0> содержит '@' '+' Reg Const)
```



```
RULE
Reg ::= Const
SEMANTICS
if (Pattern<0> содержит Const в j-й позиции)
Pattern<0>[j]=true;
if (Pattern<0> содержит Reg в i-й позиции)
Pattern<0>[i]=true.
```

Для дерева рис. 9.23 получим значения атрибутов, приведенные на рис. 9.25. Здесь М обозначает Match, Р — Pattern, С — Const, R — Reg.

Глава 10

СИСТЕМЫ АВТОМАТИЗАЦИИ ПОСТРОЕНИЯ ТРАНСЛЯТОРОВ

Системы автоматизации построения трансляторов (САПТ) предназначены для автоматизации процесса разработки трансляторов. Очевидно, что для того, чтобы описать транслятор, необходимо иметь формализм для описания. Этот формализм затем реализуется в виде входного языка САПТ. Как правило, формализмы основаны на атрибутивных грамматиках. Ниже описаны две САПТ, получившие распространение: СУПЕР [4] и Yacc. В основу первой системы положены LL(1)-грамматики и L-атрибутивные вычислители, в основу второй — LALR(1)-грамматики и S-атрибутивные вычислители.

10.1. Система СУПЕР

Программа на входном языке СУПЕР («метапрограмма») состоит из следующих разделов:

- заголовок;
- раздел констант;
- раздел типов;
- алфавит;
- раздел файлов;
- раздел библиотеки;
- атрибутивная схема.

Заголовок определяет имя атрибутивной грамматики, первые три буквы имени задают расширение имени входного файла для реализуемого транслятора.

Раздел констант содержит описание констант, раздел типов — описание типов.

Алфавит содержит перечисление нетерминальных символов и классов лексем, а также атрибутов (и их типов), сопоставленных этим символам. Классы лексем являются терминальными символами с точки зрения синтаксического анализа, но могут иметь атрибуты, вычисляемые в процессе лексического анализа. Определение класса лексем состоит в задании имени класса, имен атрибутов для этого класса и типов этих атрибутов.

Определение нетерминальных символов содержит перечисление этих символов с указанием приписанных им атрибутов и их типов. Аксиома грамматики указывается первым символом в списке нетерминалов.

Раздел библиотеки содержит заголовки процедур и функций, используемых в формулах атрибутивной грамматики.

Раздел файлов содержит описание файловых переменных, используемых в формулах атрибутивной грамматики. Файловые переменные можно рассматривать как атрибуты аксиомы.

Атрибутная схема состоит из списка синтаксических правил и сопоставленных им семантических правил. Для описания синтаксиса языка используется расширенная форма Бэкуса–Наура. Терминальные символы в правой части заключаются в кавычки, классы лексем и нетерминалы задаются их именами. Для задания необязательных символов в правой части используются скобки [], для задания повторяющихся конструкций используются скобки (). В этом случае может быть указан разделитель символов (после /). Например,

$$A ::= B [C] (D) (E / ' , ')$$

Первым правилом в атрибутивной схеме должно быть правило для аксиомы.

Каждому синтаксическому правилу могут быть сопоставлены семантические действия. Каждое такое действие — это оператор, который может использовать атрибуты как символов данного правила (локальные атрибуты), так и символов, могущих быть (динамически) предками символа левой части данного правила в дереве разбора (глобальные атрибуты). Для ссылки на локальные атрибуты символы данного правила (как терминальные, так и нетерминальные) нумеруются от 0 (для символа левой части). При ссылке на глобальные атрибуты надо иметь в виду, что атрибуты имеют области видимости на дереве разбора. Областью видимости атрибута вершины, помеченной N, является все поддерево N, за исключением его поддеревьев, также помеченных N.

Исполнение операторов семантической части правила привязывается к обходу дерева разбора сверху-вниз слева-направо. Для этого каждый оператор может быть снабжен меткой, состоящей из номера ветви правила, к выполнению которой должен быть привязан оператор, и, возможно, одного из суффиксов A, B, E, M.

Суффикс A задает выполнение оператора перед каждым вхождением синтаксической конструкции, заключенной в скобки повторений (). Суффикс B задает выполнение оператора после каждого вхождения синтаксической конструкции, заключенной в скобки повторений (). Суффикс M задает выполнение оператора между вхождениями синтаксической конструкции, заключенной в скобки повторений (). Суффикс E задает выполнение

оператора в том случае, когда конструкция, заключенная в скобки [], отсутствует.

Пример использования меток атрибутивных формул:

```
D ::= 'd' =>
    $0.y:=$0.x+1.
A ::= B (C) [D] =>
    $2.x:=1;
2M: $2.x:=$2.x+1;
    $3.x:=$2.x;
3E: $3.y:=$3.x;
3:  writeln($3.y).
```

Процедура `writeln` напечатает число вхождений символа `C` в `C`-список, если `D` опущено. В противном случае напечатанное число будет на единицу больше.

10.2. Система YACC

В основу системы YACC положен синтаксический анализатор типа LALR(1), генерируемый по входной (мета)программе. Эта программа состоит из трех частей:

```
%{
Си-текст
}%
%token Список имен лексем
%%
Список правил трансляции
%%
Служебные Си-подпрограммы
```

Си-текст (который вместе с окружающими скобками `%{` и `%}` может отсутствовать) обычно содержит Си-объявления (включая `#include` и `#define`), используемые в тексте ниже. Этот Си-текст может содержать и объявления (или предобъявления) функций.

Список имен лексем содержит имена, которые преобразуются YACC-препроцессором в объявления констант (`#define`). Как правило, эти имена используются как имена классов лексем и служат для определения интерфейса с лексическим анализатором.

Каждое правило трансляции имеет вид

```
Левая_часть : альтернатива_1
              {семантические_действия_1}
            | альтернатива_2 {семантические_действия_2}
            | ...
            | альтернатива_n {семантические_действия_n}
            ;
```

Каждое семантическое действие — это последовательность операторов Си. При этом каждому нетерминалу может быть сопоставлен один синтезируемый атрибут. На атрибут нетерминала левой части ссылка осуществляется посредством значка \$\$, на атрибуты символов правой части — посредством значков \$1, \$2, ..., \$n, причем номер соответствует порядку элементов правой части, включая семантические действия. Каждое семантическое действие может вырабатывать значение в результате выполнения присваивания \$\$=Выражение. Исполнение такого оператора в последнем семантическом действии определяет значение атрибута символа левой части.

В некоторых случаях допускается использование грамматик, имеющих конфликты. При этом синтаксический анализатор разрешает конфликты следующим образом:

- конфликты типа свертка/свертка разрешаются выбором правила, предшествующего во входной метапрограмме;

- конфликты типа сдвиг/свертка разрешаются предпочтением сдвига.

Поскольку этих правил не всегда достаточно для правильного определения анализатора, допускается определение старшинства и ассоциативности терминалов.

Например, объявление

```
%left '+' '-'
```

определяет + и – имеющими одинаковый приоритет и левую ассоциативность. Операцию можно определить как правоассоциативную в результате объявления:

```
%right '^'
```

Бинарную операцию можно определить как неассоциативную (т. е. не допускающую объединения двух подряд идущих знаков операции):

```
%nonassoc '<'
```

Символы, перечисленные в одном объявлении, имеют одинаковое старшинство. Старшинство выше для каждого последующего объявления. Конфликты разрешаются путем присваивания старшинства и ассоциативности каждому правилу грамматики и каждому терминалу, участвующим в конфликте. Если необходимо выбрать между сдвигом входного символа s и сверткой по правилу $A \rightarrow w$, то свертка делается, если старшинство правила больше старшинства s или если старшинства одинаковы, а правило левоассоциативно. В противном случае делается сдвиг.

Обычно за старшинство правила принимается старшинство самого правого терминала правила. В тех случаях, когда самый правый терминал не дает нужного приоритета, этот приоритет можно назначить следующим объявлением:

```
%prec терминал
```

Старшинство и ассоциативность правила в этом случае будут такими же, как у указанного терминала.

YACC не сообщает о конфликтах, разрешаемых с помощью ассоциативности и приоритетов. Восстановление после ошибок управляется пользователем с помощью введения в грамматику «правил ошибки» вида

$$A \rightarrow \text{error } w.$$

Здесь `error` — ключевое слово YACC. Когда встречается синтаксическая ошибка, анализатор трактует состояние, набор ситуаций которого содержит правило для `error`, некоторым специальным образом: символы из магазина выталкиваются до тех пор, пока на его верхушке не будет обнаружено состояние, для которого набор ситуаций содержит ситуацию вида $[A \rightarrow \text{.error } w]$, после чего в магазин фиктивно помещается символ `error`, как если бы он встретился во входной строке.

Если w пусто, то делается свертка. После этого анализатор пропускает входные символы, пока не найдет такой, с которым можно продолжить нормальный разбор.

Если w не пусто, просматривается входная строка и делается попытка свернуть w . Если w состоит только из терминалов, то эта строка ищется во входном потоке.