

Параллельная обработка данных

8.09

Компьютерный мир



→ Среда 1956
300 м²
150 кВт
2000 оп/сек

Совр. ПК ~ 10 Gflops

flops = floating point operation per second

- весовые,
- арифметические,
- логические операции

M - 10⁶, G - 10⁹, T - 10¹², P - 10¹⁵

~ 1-2 GBytes

~ 400 GB (жест. диск)

Top 500.org список самых мощных

CRAY XT4, ~ 36000 проц., AMD Opteron, DC
~ 266 Tflops

IBM BlueGene, 212990 процессоров,
IBM PowerPC440
~ 478 Tflops

SGI Altix, 51000 проц., Intel Xeon, QC
~ 487 Tflops

IBM Roadrunner ~ 6500 AMD Opteron, DC
+ ~ 12200 IBM Cell
~ 1100 Tflops = 1.1 Pflops
ОП = 98 Тбайт

- задачи сложны?
- задачи ватны?

билет
270937

задача: посчитать кол-во
счастливых билетов

```
for (i0=0; i0<9; ++i0)
  for (i1=0; i1<9; ++i1)
    ... // i2
    ... // i3
    ... // i4
    ... // i5
    if (i0+i1+i2 ==
        i3+i4+i5)
      ++count;
```

- 10 цифр ~ 40 сек
- 12 цифр ~ 1.5 часа
- 14 цифр - попутно записывается
- 16 цифр - не решается

10^N , N-цифровое число

Даже супер с нашим алгоритмом
много не решает

А алгоритм сложности N^2

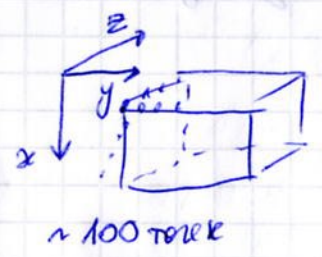
• Задача N ферзей

- 8x8 - бонус
- 25x25 - рекорд
- 14x14 10 ферзей - уже тяжело

Моделирование нефтяного резервуара



Будем считать, это
резервуар - прямоуг. паралл.



5-20 функций в V точке

Решить темн. ур-е для
V ф-ций, т.е. выполнит
200-4000 операций
~ 100-1000 шагов по времени

Всего операций: 10^6 (точек) * 10 (ф-ций) *
 $* 500$ (операций) * 500 (шагов по времени) =
 $= 2.5 \cdot 10^{12}$

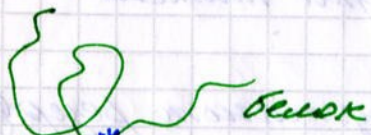
- автомобилестроение

- крас-тесты
- формула 1: BMW в 2006 - кр на
12 Tflops
- Williams в 2007 - суперк на 7 Tflops
- Renault, 2008 36 Tflops

- авиадвигатели

TOP50 - суперк. в СНГ supercomputers.ru

- проектирование лекарств



кратейный угасник

Надо перебрать лекарства

На проверку содержания угасник

от 1 мин до 10 часов

- кинопроизводство

Шрек 3 ~ 12000 кадров

1 кадр ~ 2 часа на 1 процессоре

~ 20 · 10⁶ процессор-часов

Суперк 8 тыс. процессоров

- суперкомпьютеры и наука

- теория

- эксперимент

- вычислит. наука

- whipreol

- Procter & Gamble

... Кинсон

EDSAC, 1949²

15.09

ТАКТ $2 \cdot 10^{-6}$ с

произв-ть 100 см/с

IBM BlueGene 2007-2008

ТАКТ $1,4 \cdot 10^{-9}$ с (700 MHz)

произв-ть $4,78 \cdot 10^{14}$ см/с

~ 1500 раз
 $4,7 \cdot 10^{12}$ раз

Принципы ТОО

- параллельность

- непрерывность



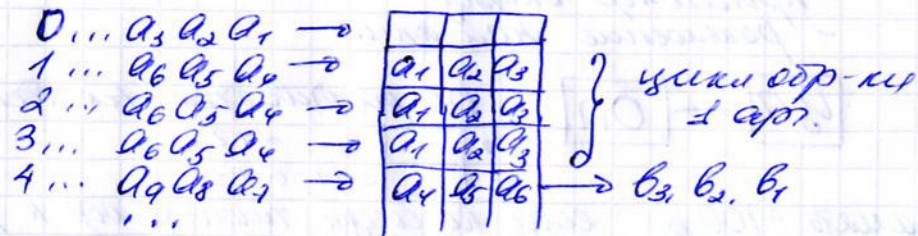
послед. обработка



T = 900 тактов

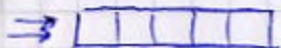
паралл. обработка - N устр-в

в N раз быстрее



T = 300 тактов

Конвейерная обработка



микрооперации (например, при сложении
всущ. чисел: сравнение порядков,
выравнивание порядков, сложение и т.д.)

0... $a_3 a_2 a_1 a_0$ [] [] [] 3 ступени конвейера, каждая 1 такт

1... $a_4 a_3 a_2 a_1 a_0 \rightarrow a_1$ [] []

2... $a_5 a_4 a_3 a_2 a_1 a_0 \rightarrow a_2 a_1$ [] [] []

3... $a_6 a_5 a_4 a_3 a_2 a_1 a_0 \rightarrow a_3 a_2 a_1$ [] [] [] []

4... $a_7 a_6 a_5 a_4 a_3 a_2 a_1 a_0 \rightarrow a_4 a_3 a_2 a_1$ [] [] [] [] []

... $\rightarrow a_5 a_4$

Циклы обработки 1 арт-та

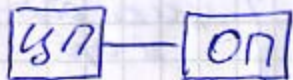
$$T = 3 + 299 = 302 \text{ такта}$$

Экспурс в историю

1) IBM 701 (1953), IBM 704 (1955)
разр-паралл - память
арифметика
ферритовая память

2) IBM 709 (1958)
независимый ввод/вывод

3) IBM STRETCH (1961)
- пропуск вперёд
- равенство памяти

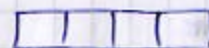
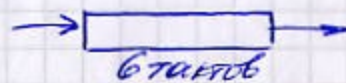


ОП разбивается на банки



память ~100р. если та соседн. тактам сор. к разным
мерилнее банкам, сор-на параллельная
(и по сей день)

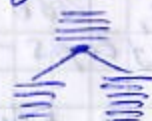
4) ATLAS (1963)
конвейер команд



4 ступени
выбора адреса операции
вычисления и т.д.

(1,5) раза, получили (1,6) р.

вселились:



намагни. проход

5) CDC 6600 (1964)
- 10 незав. Функц. Упр-в

такт 100нс
произв-ть 2-3 млн оп/сек

6) CDC 7600 (1969)
- 8 незав. конвейерных ФУ

такт 27нс
произв-ть 10-15 млн оп/сек

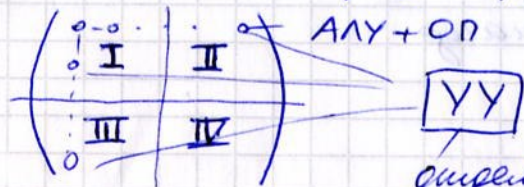
- типовая произв-ть: (недостигали)
число опер-ц
время

- реальная произв-ть (опр. программ)
всего чипов меньше типовых

7) ILLIAC IV (1974)

- матрица процессоров

16x16 - 256 процессоров



в 4 моменты вр. выполняет все команды

Процессор разбито на 4 квадранта

Проект: 40 нс, 256 ПЭ
~1 Gflops - типовая прогн-та

- 1967 - начало проекта
- 1971 - квадрат (8x8)
- 1974 - сдвиг в экпл.

1/4, 80 нс
реальная прогн-та 50 Mflops

Иерархия памяти

- регистры
- кэш-память
- оперативная память
- диск
- лента
- ...



- локальность вычислений

маш. код



- локальность пер-х данных



- 6 регистров в прогн. элементе
- ОП в ПЭ (2048 слов)
- 2 диска по 1 Gbit
- лазерный диск 10¹² bit

Закон Амдала

22.09

ускорение в p раз (# процессоров)

0 ≤ f ≤ 1 - доля операций, кот.

невозможно распараллелить

$$\frac{1-f}{p} = S \leq \frac{1}{f + \frac{1-f}{p}} \quad \text{ускорение } \left(\frac{t_{\text{с прогн}}}{t_{\text{p прогн}}} \right)$$

$$f=0,1, p \rightarrow \infty \Rightarrow S \leq 10$$

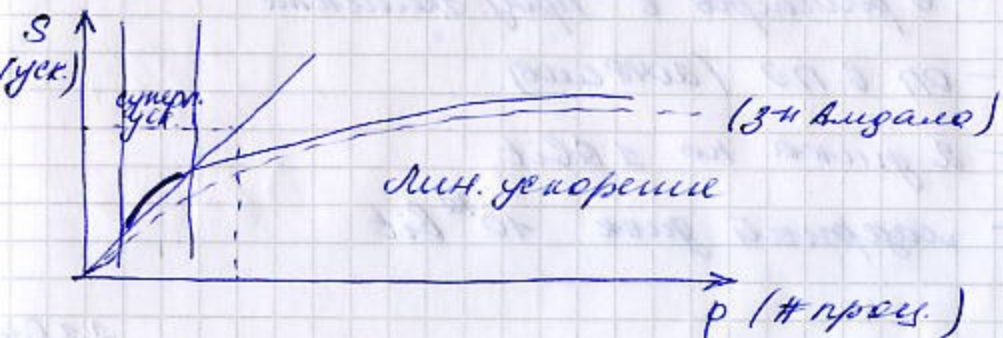
⇒ не так уж трудно увеличить # прогн-в
не получим ускорения > 10 раз

Следствие: для ускор. в q раз
необходимо ускорить в q раз не менее,
чем $(1 - \frac{1}{q})$ -ю часть программы. (5)

$$q = 100 \Rightarrow (1 - \frac{1}{q}) = 99,99\% \text{ пр-во}$$

$$\Rightarrow S = q \text{ задач, тогда } f \leq \frac{1}{q}$$

Суперлинейное ускорение



при эффективном кэшировании

кэш-памяти можно получить больше, чем мин. ускорение

$S = 0;$

```
for (int i=0; i<n; i++)
```

```
    S = S + A[i];
```

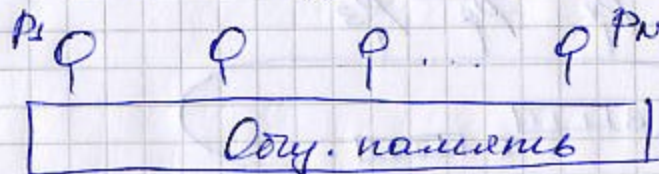
Сравнение (разбивает A на пары)



проблемы: - фрагменты массива
- доступности массива A
- задержки при перемещении
...
...

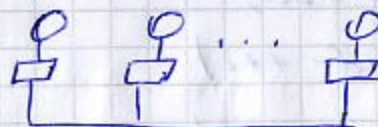
Архитектура компьютеров

1) Комп. с общей памятью



SMP (Shared Memory Processors)
Symmetric Multi Processors

2) Комп. с распредел. памятью



"Утка лесно" - канал связи с ОП
N=64

N=1000

```
for (i=0; i<N; ++i)
```

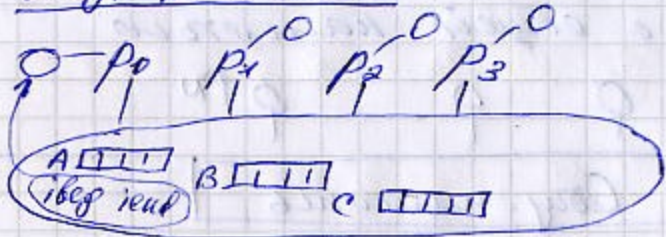
```
    Ai = Bi + Ci * x;
```

```
for (i=0; i<N; ++i)
```

```
    Ai = (Ai + Ci) + BN-i-1;
```



Одн. память



$N = 1000;$

$n-pr = 4;$

$size = N/n-pr;$

$ibeg = proc_id * size;$

$iend = proc_id * size + size;$

for ($i = ibeg; i < iend; ++i$)

for ($i = ibeg; i < iend; ++i$)...

Классы переменных

глобальные: A, B, C, x

локальные: i, ibeg, iend

Распр. память



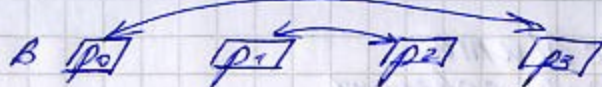
$N = 1000 / 4 = 250$

<обмен 1>

for ($i = 0; i < N; i++$)

...

<обмен 2>



for (...)

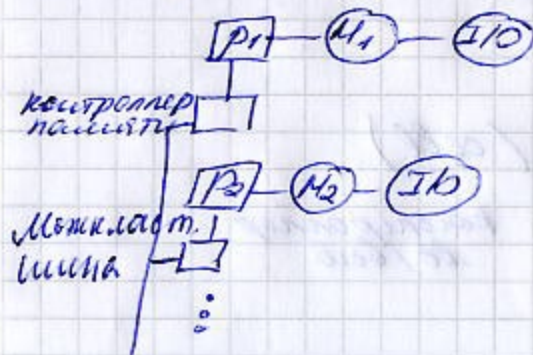
<обмен 3>



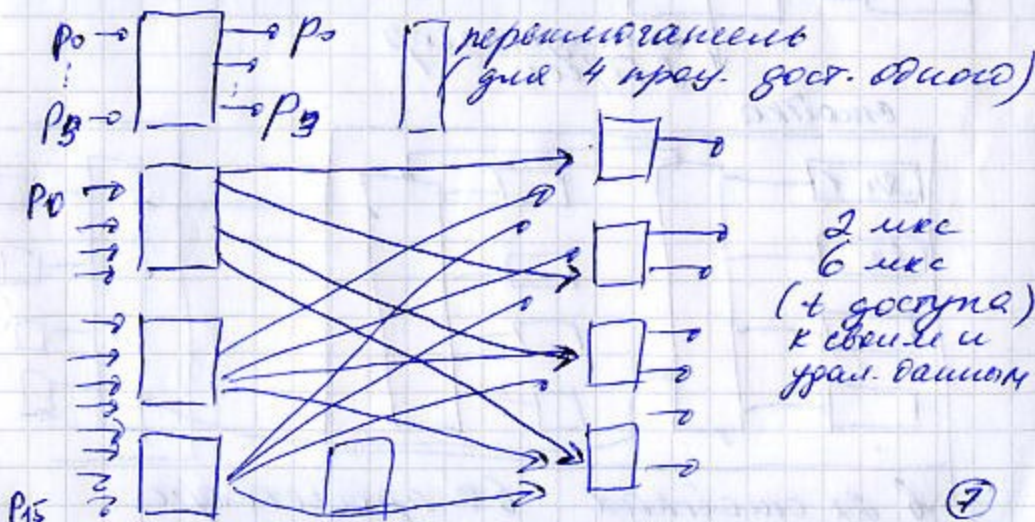
N -pn access UMA (Unified Memory Access)

NUMA (Non-UMA)

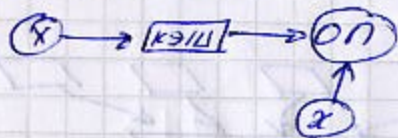
Скелетон - в комм. C_m^*



BBN Butterfly



ccNUMA
cache coherent



HP SuperDome (20002)

до 64 ядр-в

ОП ≤ 256 Гб - 1 Тб

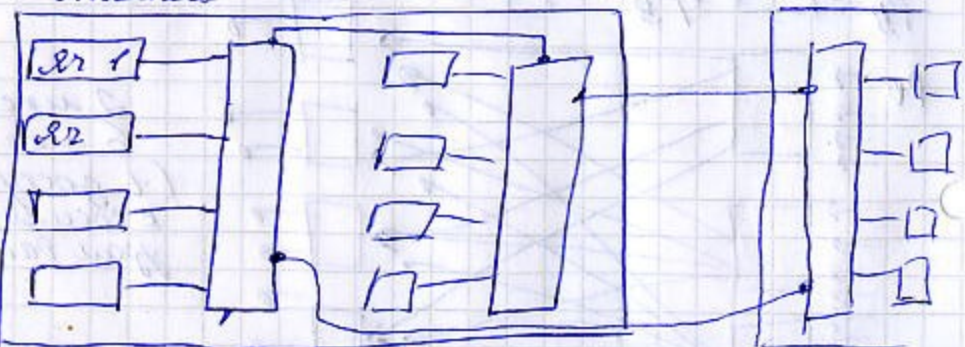
РА 8600 / 8700

Itanium

формат. ячея (cell)



стойка



в 2х стойках 64 процессора

РА 8700

750 МГц

4 яд / макс

мик. 3 Гфрагм

x64

192 Гфрагм

Выделение ячеек произв-ми: 610



Мешают произв-ми:

1) Закон Амгана

$$\text{Ускорение } S < \frac{1}{\alpha + \frac{1-\alpha}{P}}$$



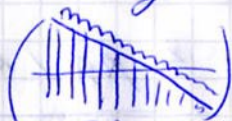
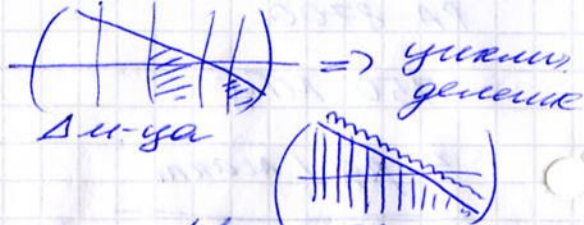
всегда есть последовательная часть

2) арх-ра ccNUMA обращение к памяти

3) ccNUMA кэш

4) Равномерное распределение нагрузки между процессорами

⇒ В Сербанке - 9 машин HP SuperDome (самая большая именная)



5) Производительность процессора

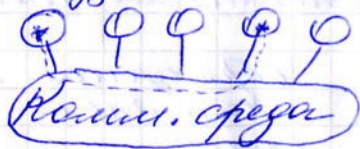
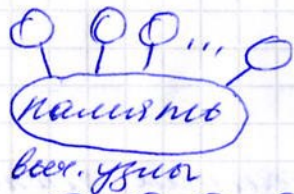
PA-8700
Itanium II

Нормальная произв-ть 15%
ниже произв-ть

→ Эффективность работы 1-3%

6) ...

Компьютеры с распределенной памятью



много времени - на обмен сообщениями

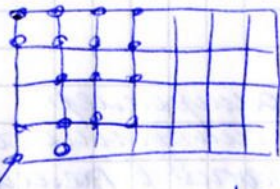
Корано 90%

- процессора (какие?)

- коммуникационная сеть (как упрощена?)

Intel Paragon

Intel Paragon

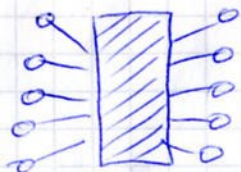


прямая решетка

связь между соседями

процессор i860

IBM SP1/SP2



связь от V к V

комм-р - т упр-я связи увеличено

Массивно-параллельные компьютеры

CRAY 1993-1994

T3D, T3E, ~~T3X~~

не войдут в топ 500

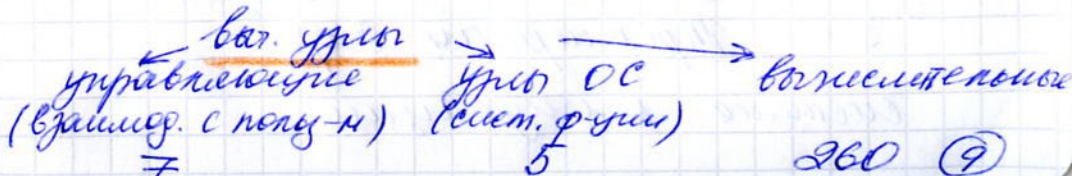
X T3, X T4, X T5 (2 место в топ 500)
в наст. вр.

T3D, T3E - на осн. проц. DEC ALPHA (120)

2 проц. X T3, X T4, X T5 - AMD Opteron (4 ядра)
на узел

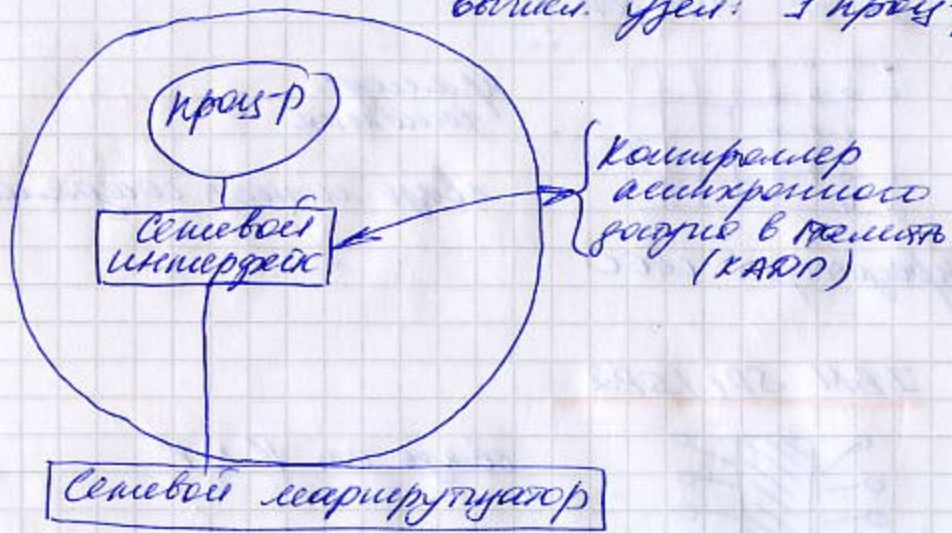
- всег. узел

- коммутир. сеть



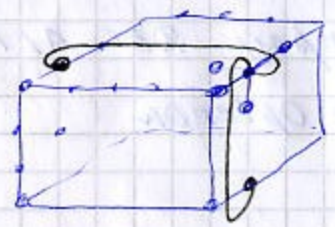
260 (9)

Высш. узел: 1 процессор



С помощью KADP узел получает доступ к памяти др. узла, не нарушая его работу - вытеснение одурин процессоров

Комп. сеть: узлы рашётки



Целочисл. точки
внутри тоже узлы



у 4 узла 6 соседей
"Трёхмерный мор" -
схема коммутации

повреждённые связи можно менять на обходной маршрут
за 3 переключения можно попасть из 1 узла в другой

x_A, y_A, z_A - коорд. А

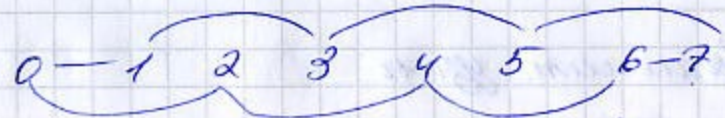
x_B, y_B, z_B - коорд. В

Маршрут из А в В:

$A \xrightarrow{x_A, y_A, z_A} A' \xrightarrow{x_B, y_A, z_A} A'' \xrightarrow{x_B, y_B, z_A} B$

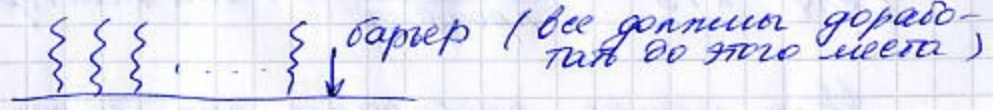
На 4 этапе - крайний нуль

$A \rightarrow B$ и $B \rightarrow A$ отличаются
длина проводника в шест
 $0-1-2-3-4-5-6-7$



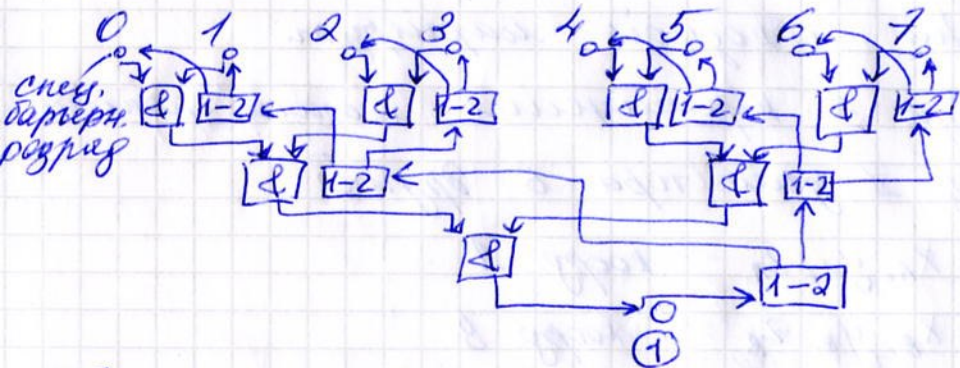
связи более-менее равномерны

Барьерная синхронизация



Коммарное кол-во сообщений -
каждые посыланы всем (для
синхронизации)

Крайнее аппаратное решение:



⊗ можно заменить на ⊕

в зареи поиска в БД, например (до первой находки)

13.10

Вычислительные кластеры

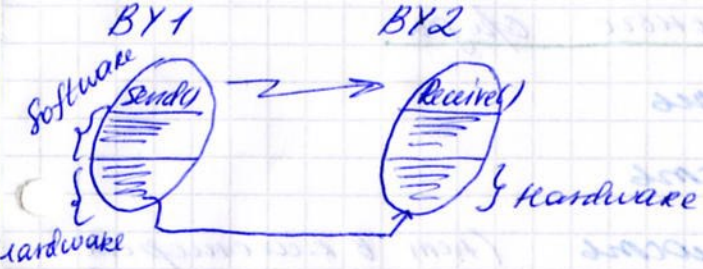
1994

- вычислит. узлы
- сетевая технология
- Linux, GNU, MPI

Cray XT5 $\frac{\text{цена}}{\text{проц}} \ll$ для кластеров

Узлы можно подбирать под задачу

- оп → max
- макс. частота → max
- оп, + HDD
- ускорители (Cell, GPU, ...)



- латентность - время для пересылки сообщения по сетевой линии (задержки)
- скорость передачи по сети
- цена

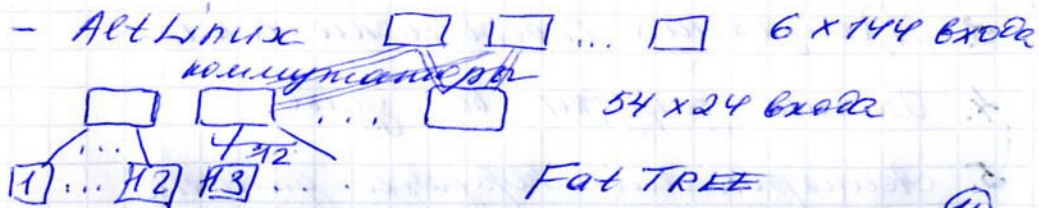
GE (gigabit ethernet) InfiniBand

1 Gb/s - скорость	20-40 Gb/s
80-90 мкс - латентность	1-3 мкс
5 \$ узел - цена	1000 \$ узел

СКУД МГУ "Челюшев"

- вычисл. узлы: Intel Harpertown, 3GHz (625)
 - (8Gb оп, без дисков,
 - 8Gb оп, HDD
 - 16Gb оп, HDD
 - 32Gb оп, HDD)

- InfiniBand DDR
- Gethernet (для поддержки работы DC)
- ServNet (для мониторинга узлов)



Экземплярные среды

- массовость
- удаленность
- динамичность (нет в кластерах и CRAU/X75...)
- гетерогенность
- разная административная принадлежность (разные канал. и др. политики и т.д.)

grid - технологии - все предельно
своей единую среду (метакомпьютер)

cloud computing - облачные вычисления

Идея:



Из сущ.-х ресурсов формируются среда, дающая сервис высокого уровня

Проблемы: (неадекватные расходы)

1. Закон Амдала $S < \frac{1}{\alpha + \frac{1-\alpha}{n}}$

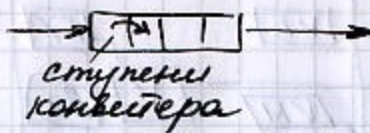
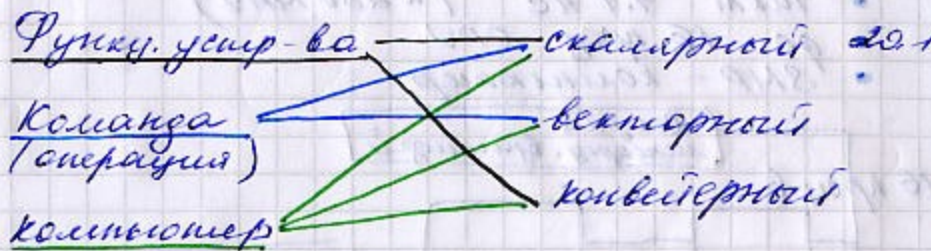
2. Латентность

3. Пропускная способность сети

4. Баланс загрузки выст. узлов

5. Асинхронная передача данных

6. Реальная прощ-ть процессоров



$A: x, y \rightarrow z$ - скаляр

$A: V_1, x \rightarrow V_2$ - вектор

Векторно-конвейерного компьютера

1976 - Cray-1

начало 90х Cray C90 (T90)

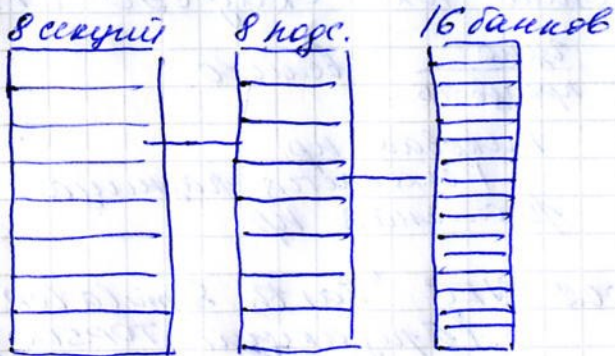
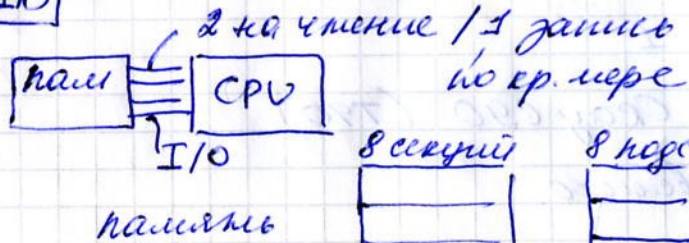
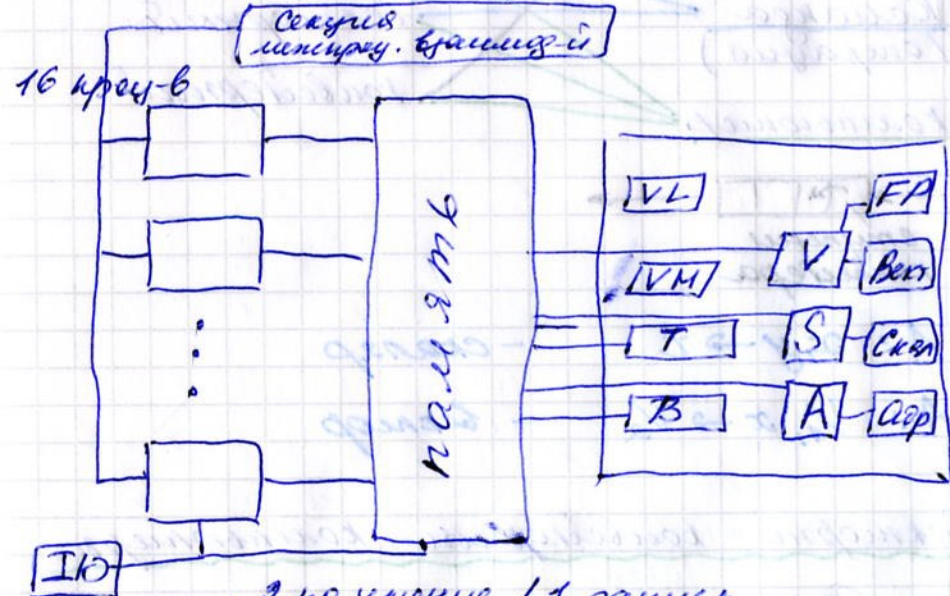
цена
проц-ть велико

микров. пр.
| считается разница
реальная пр.

2002 NEC "Earth Simulator" 40 Tflops
(5 редакций TOP500)

CRAY C90

- такт 4.1 нс (~250 МГц)
- 90 - 16 процес. CPU
- SMP - компьютер



- 0 0e 0n 0b
- 1 1c 0n 0b
- ...
- 7 7c 0n 0b
- 8 0c 1n 0b
- ...
- 63 7c 7n 0b
- 64 0c 0n 1b
- 65 1c 0n 1b
- ...

макс. задержка 7 тактов при сор. к с.ч. по ш. соединенным - с емкостью 64

$$h = A \cdot e^{k/n}$$

$k=6$ - макс. кол-во
 $k=0$ - мин. кол-во

• регистры

- Осн. группа

- Вспомогат. группа

• адресное шир 32 разряд

• скалярное шир 64 разряд

• векторное

вект. рег. - совок. из 128 регистров 64 разряд



Вспомогат. • 64 шир. 32p (B) (для A)

• 64 шир. 64p (T) (для S)

VL - регистр для вектора (Sp)

VM - регистр маски

if (...) $A_i = \dots$
else $A_i = \dots$

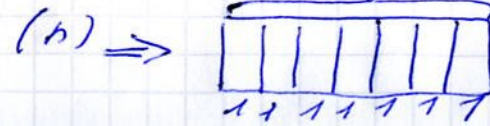


• Функц. устройства

- независимые

- конвейерные

- 4 ступени - 1 такт



(n+1) время в тактах при макс. из n-тов



$\max(d_k)$ - расстояние
полученный результат

1. Агрегированные ПУ

2 шим., 32 разр., свая. операции. A_2 ^{32,0}

2. Скалярные ПУ

4 шим., 64 разр., скал. опер., A_2 ^{64,0}

3. Векторные ПУ

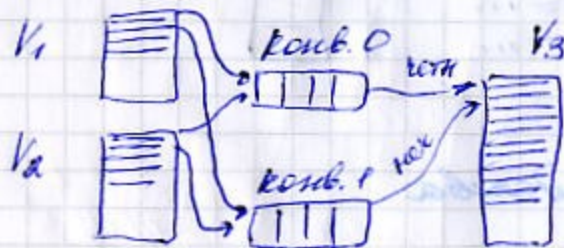
5-7 шим., 64 разр. вект. опер. A_2 ^{64,8}

4. Век. арифметико-FP

3 шим., 64 разр., свая/вект. A_{FP} ^{64,0} A_{FP} ^{64,6}

$\pm, *, \frac{1}{x}$

$A: V_1, V_2 \rightarrow V_3$



"четн. эл-ты на конв. 0
нечетн. - на конв. 1"

$a_i = b_i + \alpha + c_i$



$(e^{*+n-1} + e^{*+n-1})$ тактов на ван.

Зацепление ПУ: z_i сразу идет
на e^* : $(e^{*+e^{*+n-1}})$ - увеличение
в 2 раза на уровне архитектуры

Особенности, помогающие
считать быстро:

- ПУ независимы
- ПУ конвейерные
- векторная обработка

for $(i=0; i < n; ++i)$

$A[i] = B[i] + C[i];$

$L: B \rightarrow V_1$
 $L: C \rightarrow V_2$
 $A: V_1 + V_2 \rightarrow V_3$
 $ST: V_3 \rightarrow A$

} закрепленные
8-10-20 раз

- дублирование конвейеров
- зацепление ПУ
- 16 процессоров

flops - линейная пропуск-ть

$\pm, *, \div$ забудем, т.к. редко используются

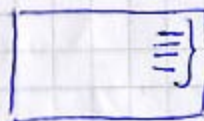
Если же 1 текст: 4 операции ($\pm, *$)

4.1 нс \Rightarrow на 1 проц. 1 Gflops
(10^9 выч. оп / сек)

27.10 2002 - 6 машин NEC "Earth Simulator"
40 Tflops / 2003-2005 #1 в Top 500

640 узлов * 8 процессоров

1 проц. \sim 8 Gflops

 674, 72 вект. регистра * 256 элементов

500 MHz


8 секций

SSE - инструкции - в Intel процессорах, цел. и вещ. векторной обработки

Свой век-конв. машинно производится NEC и CRAY (прототип кодова)

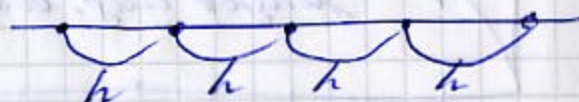


Векторизация программ

 } \rightarrow for (i=0; i<n; ++i) } \rightarrow
 $A[i] = B[i] + C[i]$
 $A^V; B, C \rightarrow A$
векторизация

Условия: - векторы данных
- относительно равномерное
операции

вектор данных:



Один и тот же тип, одно число тип
разм. в памяти



for (i=0; i<n; ++i)
 $A[i] = A[i] + B[i]$

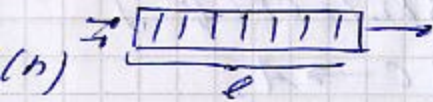
! конвейер работает не лучше, н.г.
• операции равны

Проблемное место:

1. Закон Амгана

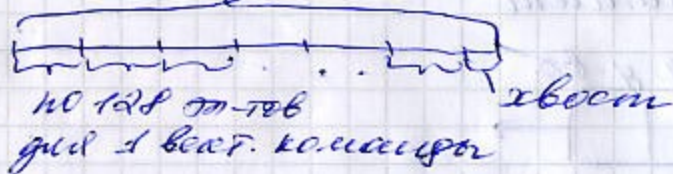
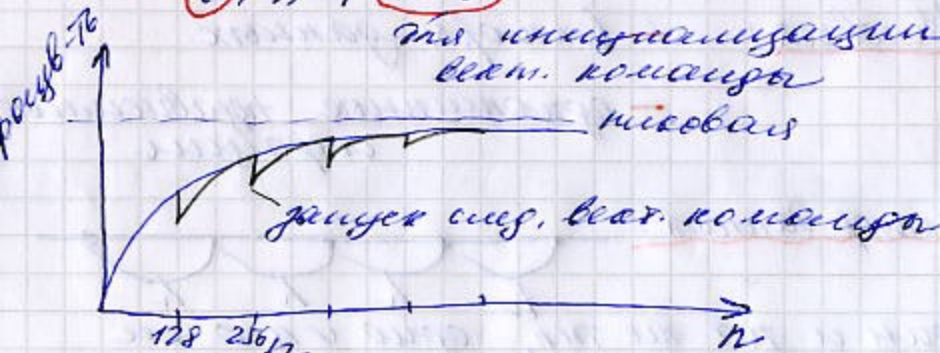
Ускорение $SS \frac{1}{d + \frac{1-d}{P}}$
счит. время
 $\xrightarrow{d_1} \xrightarrow{d_2}$

2. Время разгона конвейера
 3. Симуляция векторных команд



$$T = e + n - 1 + d$$

то инциализация
вект. команды



Операция $a_i = b_i * s + c_i$

n	M-flops
1	7
4	27
64	300
128	430
129	364
256	548
257	490
...	...
8200	800

4. Конфузия в памяти

```
for (i=0; i<n+k; i+=k)
  a[i] = b[i*s] + c[i];
```

n = 1000

k	M-flops
1	405
2	444
4	274
8	142
16	84
32	44
64	22

```
do i=1, n
  do j=1, n
    do k=1, n
      x(i,j,k) = x(i,j,k) + p(k) * y(k,j)
```

$$x(i,j,k) \sim x(i-1,j,k) \quad 1 \quad (x/c_1, c_2, c_3)$$

$$x(i,j,k) \sim x(i,j-1,k) \quad c_1$$

$$x(i,j,k) \sim x(i,j,k-1) \quad c_1 * c_2$$

$x(40, 40, \dots) \Rightarrow$ фазет. 1600 = 64 * 25
 \Rightarrow при векторизации n -то фазетом
 хушим образом

Решение: $x(41, 41, \dots)$
 * фазетно - то будет конфуз. в памяти

Агрегация: $A(I A(I))$
 $I A(I) = I$
 $I A(I) = \text{const} - \text{конфузия}$
 где n то же значение n (не)

5. Ограниченной пропускной способностью канала процессор-памяти

$$a_i = b_i + s_i + c_i$$

\downarrow \downarrow \downarrow \downarrow
 b_{in} s_x b_x c_x - только 2 входа

6. Необходимость цел-х векторных регистров

7. Сбалансированность + / *

8. ...

$A = B + C$ вектора

$A = B + C * I$ произв. возрастает в 2 раз

$$\frac{N_{оп}}{T}$$

Параметры на уровне машинной команды



- сколько-то ПУ
- суперкамарная архитектура
- арх-ра с большими командами, веном (VLIV)

машинные инструкции
 в процессе исполнения происходит перемещение

- путь во всех совр. процессорах

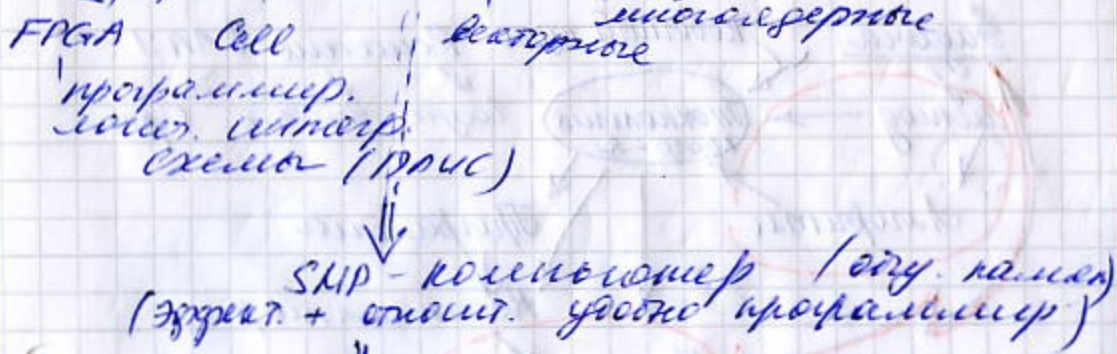
VLIV (Very Large Instruction Word)

Распаралелливается компьютеру



Всё решается в статье - компьютер от цел. текста

GP (AMD) (графические) эффективность GPU
 удобство у разработчиков программ. мир. вычисления



К-р с распр. памятью
 NEC SX-9

- CRAY XT5 (h)
- AMD (6-е ядро)
 - векторные
 - FPGA

распред. ↓ вертикаль. обмен

3.11 Технологии параллельного программирования

Типовая конфигурация 1.2 P900's

Стоимость ~ \$88M

Площадь ~ 1000 м²

Вес 150 т

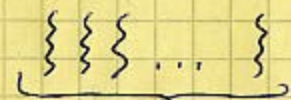
6000+ Intel

5000+ GPU (AMD)



Нужно, чтобы задача обтекло
решалась на компьютере

Модель SPM
(Single Program - Multiple Data)



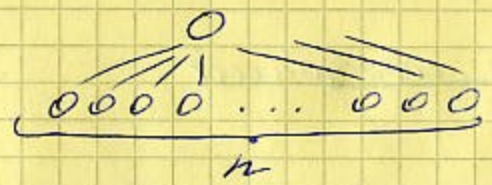
n процессов, у каждого свои данные

Как "подружить" разные данные?

Если спец. функция типа my-id()

```
if (my-id() == 3) { ... }
else { ... }
```

Master/Slaves (Workers)



процессы "общаются" только с
мастером, ком. контролер. ситуация
у мастера своя программа,
у рабов - своя

MIMD

1. Традиц. ЯПроц-я (Fortran, C, ...)
2. Специализированные

Open MP

pragma omp

средство добавить комментарии
к y-ву алгоритма

3. Расширения ЯП

HPF (high performance Fortran)

FORALL

MPC (расширение C)

Ограничение: вопрос к компьютеру:
для каких платформ он спущен?

4. Новые языки

Occam

Sisal (язык ориентированного программирования)

Fortress

5. Сети передачи сообщений

{...}

MPI

PVM

Shmem (shared memory, упрощ. CRAY, активная сторона одна)

6. Параллельные библиотеки

LAPACK (лин. алгебра)

FFTW

MKL (on Intel)

7. Специализированные пакеты

ANSYS, CFX, FlowVision
GAUSSIAN, GAMESS

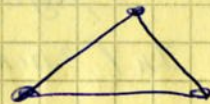
ИПО Сатурн иезелурот иттеперг.
накелот, они очень доро-
ше для коммерческого исп-
тия

Много времени проходит между
новой технологией и внедрением
=> в пакетах устаревшие модели

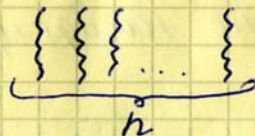
Решение: разработка программ
с открытым кодом

Нужно найти компромисс между

Эффективностью пр-я,
продуктивностью и
переносимостью



Технология прор-я Linda



пространство
кортежей

Нет прямого общения процессов
между собой, только через пр-во
кортежей

кортеж: (17, true, 3.7)

("summa", 28.1, 32.5, 0, 0)

Функции:

1. Out ("x", 12)

кортеж будет занесен в пр-во
кортежей, если такой уже был,
добавится еще один
out никак не блокирует

2 In("y", 28, 13)

забрали у пр-во кортежей
указаний
можно в кон-ве пар-ра указать
шаблон:

In("sum", int i, 13)
какое-то число, если
несколько кортежей поддают,
случайно выпадет один

In("sum", ?i, 13)
если переи. i уже определены в пр-ве
(разные области видимости)

In("sum", i, 13)
вычисл. значение i и
представляется

In("x", ?i, ?j, int k)
Если ты один кортеж не найден,
процесс блокируется

3. Read(— " —)
не запрашивает кортеж у пр-ва,
а только знает значение

4. eval("p", 17, g(x)+12, f(x)*f(x))

уже дано
не запрашивается
заверш. работ
процессов
кортеж
кортеж
процесс
кортеж
кортеж
процесс

a = g(x)+12
b = f(x)*f(x)
→ ("p", 17, a, b) → в пр-во кортежей

out("Next", 0); "Желтая
калочка"

{ In("Next", ?my-id);
Out("Next", my-id+1);

все процессы сдв. переступают
read("Next", ?number)
↑ однее # процессов

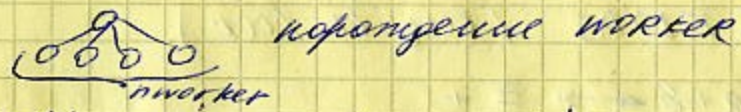
Out("Barrier", n);

In("Barrier", ?b);

b = b - 1; // не ебл-ся ли процесс
последним?

if(b != 0) { out("Barrier", b);
read("Go"); // блокировка,
пока не пойдут
кортеж "Go"

}
out("Go");



for(i=0; i < nworker; ++i)

eval("go", worker(i));

В парам. программы может быть
таблица в терминах Linda

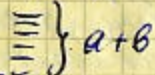
- обращение к пр-во кортежей
- объем памяти (для номеров и
в-м с расп. памятью кода моделир.)
- поиск нужного кортежа

10.11 Оценка производительности параллельных вычисл. систем

- Типовая кривая

- MIPS, маш. команды (миллион оп/сек)

C/C++/F



в перлинах MIPS в прог. характеристика лучше, а фук. времени меньше у сопроцессора

- FLOPS (век. оп. /сек)

$$a+b \leftrightarrow \equiv +$$

1 к 1

- мест Linpack - основе top500
решение СЛАУ с матрицей n-угей (LU-разложение)

1. n-уга 100x100, запрет на изменение текста программы
2. размер 1000x1000 [
 - инициализация
 - решение СЛАУ
 - $\frac{2n^3}{3} + 2n^2$ операций (и.р. пользователя)

3. Размер n-уги $N_{max} + N_{max}$ - которое "влезает" в вычисления, сваленку (разрешено использовать BLAS - базовую библиотеку мин. алгебры)
Если имеет запрет галек от СЛАУ, мест не очень подходит

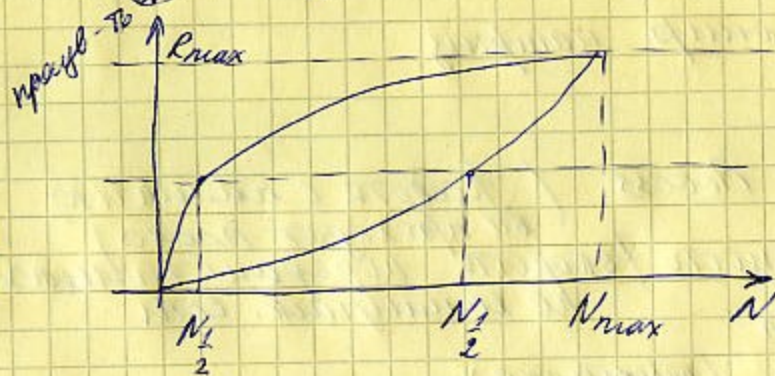
Top 500

R_{peak} - пиковая
 R_{max} - на Linpack

$$\frac{R_{max}}{R_{peak}} \cdot 100\% \sim КПД$$

N_{max}

$$N_1 \sim \frac{R_{max}}{2}$$



- Perfect Club Benchmarks

$\Sigma \sim 50000$ операций

Много приложений из фрагментов облачной
Сложно переиспользовать тест для разных конфигураций - от теста вообще отказываются

- NAS Parallel Benchmarks

- 5 вычислит. ядер (FFT, IS, CG, ...)

- 3 моделих приложения

NPB 1.0 - фикс. задача

NPB 2.0 - можно только
одну компьютеру
допускается внесение
изменений < 5% их текста

A, B, C (D) - классы памяти
(объём)

HPC Challenge

- HPL (Linpack)

- умножение $n \times n$

- транспонир. матриц

- FFT

- Random Access (работа с памятью,
все проходит плохо)

- Латентность / скорость передачи данных
по компьютер. сети

- STREAM

$$\begin{cases} a_i = b_i & \text{(пересылка)} \\ a_i = a_i + b_i \\ a_i = b_i + c_i \\ a_i = b_i * s + c_i \end{cases}$$

Средняя (скорость) $\frac{R_{peak}}{Средняя}$ > 1 для почти всех с-м

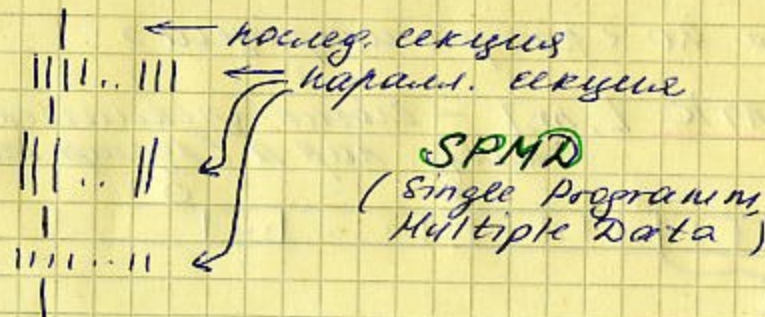
Open MP

C/C++, Fortran (77, 90...)

- спецификациям

- функции

- доступ к переменным среды



pragma omp

!\$OMP L // Fortran

!\$OMP PARALLEL

<парам. среды>

!\$OMP END PARALLEL

OMP_NUM_THREADS - перем. среды

OMP_DYNAMIC - можно ли менять ↑
в процессе программы

OMP_NESTED - разрешение вложенного
параллелизма

<условие> - вход в критич. секцию
при выполнении условия

Распределение работы

1) Ф-ция нижнего уровня

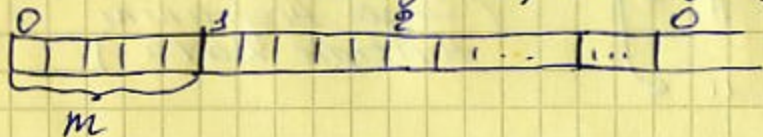
```
IF (get-omp-nit-thread(). EQ. 4) THEN  
  <для типа 4>
```

```
ELSE  
  <для остальных>  
END IF
```

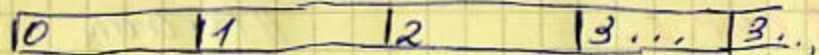
2) Парам. цикла

```
!$OMP DO <размер-с-итерациями>
```

a) STATIC [, m] - делит цикл на m частей
каждой. пар-р (размер блока)



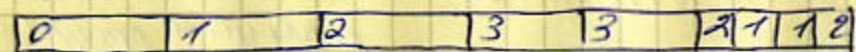
b) DYNAMIC [, m]



Вначале всё так же, а потом
каждую порцию даёт тому, кто быстрее
освободился

В конце - один работает, остальные
простаивают

в) guided [, m]



Блоки всё меньше по размеру

2) runtime

может способ в динамике или
выбор свой

3) Секции

Усл., когда нужно описать конкретную
параллелизм, усл. в программе

```
!$OMP SECTIONS
```

```
!$OMP SECTION  
  <1>
```

```
!$OMP SECTION  
  <2>
```

```
!$OMP SECTION  
  <3>
```

```
!OMP END SECTIONS
```



(NOWAIT) - нет барьеров

4) Single

Фек. параметр когда должен

Парм. исполнен 1 раз (не важно, кем)

```
!$OMP SINGLE  
  <'...>
```

```
!$OMP END SINGLE
```


MPI_Comm_size (MPI_Comm comm, int *size)

Название функции параметризовано:

if (rank == 0) { ... }

double MPI_Wtime (void) - время, прошедшее с тех. момента

Точка - точка, # загов

MPI_Send (void *buf, int count, MPI_Datatype dt, int *dest, int msgtag, MPI_Comm comm)

Для процесса данные принимаются асинхронно коммуникацией

Блокирующая

int -> MPI_INT
double -> MPI_DOUBLE

MPI_Bsend
buffer,

Будет ли работать: нужно задать ему какой-то размер и тип передаваемых данных. Будем контролировать в ядре - нет блокировки

MPI_Ssend
Synchronous - ядро, пока сообщ. не принято

MPI_Rsend
Ready - пока не будет гарантировано, что принимающий процесс готов, и сообщ. без проверки поступает прин. процессу

MPI_Recv (void *buff, int count, MPI_Datatype dt, int source, int msgtag, MPI_Comm comm, MPI_Status *status)

Для MPI_Status - структура, содержит следующие поля: status.MPI_SOURCE / ID отпр, status.MPI_TAG / тэг сообщ, status.MPI_ERROR / код err.

Примеч. процесс не обязан знать идентификатор и тэг - можно указать в аргументах MPI_ANY_SOURCE и MPI_ANY_TAG

Блокирующая: пока не поступит сообщ.

MPI_Get_Count (MPI_Status *status, MPI_Datatype dt, int *count)
Сколько данных было передано

MPI_Send MPI_Recv
MPI_Recv MPI_Send

MPI_Isend (... , MPI_Request *req)
:
IbSend
Isend
Irecv

MPI_Irecv (...)

MPI_Wait (MPI_Request *req, MPI_Status *status)
Онижение завершения фидлоуп. опер.

Неоднопр. - убиваем deadlock

waitall
waitany
waittrue

MPI-Test (req, int *flag, status)
проверяет, завершилась ли операция req

MPI-Testall
Testany
Testsome

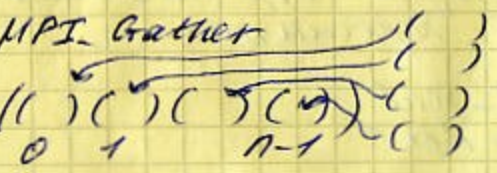
MPI-Sendrecv(...)
deadlock завершается аппаратно

Кеместивное взаимодействие -

грабьют все процесс коммуникатора

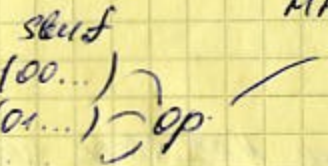
MPI-Beast (void *buf, int count,
MPI-Datatype dt, int source,
MPI-Comm comm)

распознает все процессы в комм



MPI-Scatter - обратное к Gather
(бу посылает одному) - один посыл ^{всем} процессу

MPI-Reduce (void *sbuf, void *rbuf,
int count, MPI-Datatype dt,
MPI-Op op, int root, MPI-Comm comm)



MPI-Barrier (MPI-Comm comm)
барьерная синхронизация

MPI-Comm-split (MPI-Comm comm,
int color, int key, MPI-Comm *newcomm)

Новые коммуникаторы:

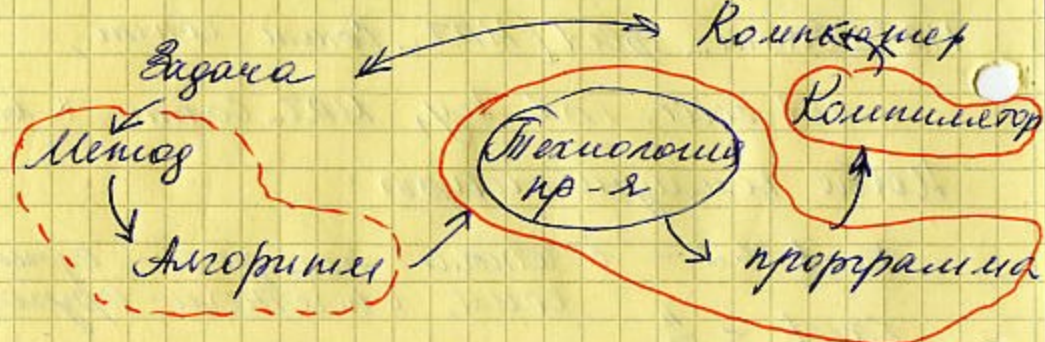
rank % 2 - с четной номерами - один
комм, с нечетными - другой
rank < $\frac{n}{2}$

MPI-Comm-free (MPI-Comm comm)
можно удалить конкрет. коммуник

MPI 2

- 1) гибкий, копирование процессов
- 2) асинхронные коммуникации
send-receive
put/get + асинхронный график
put →
← get
- 3) параметричной блэг/блвбг
- 4) расширение ф-ти конклет. операций
- 5) C++

1.12 Введем в теорию анализа структуриро парал. прогн-е



Пример переменные матриц

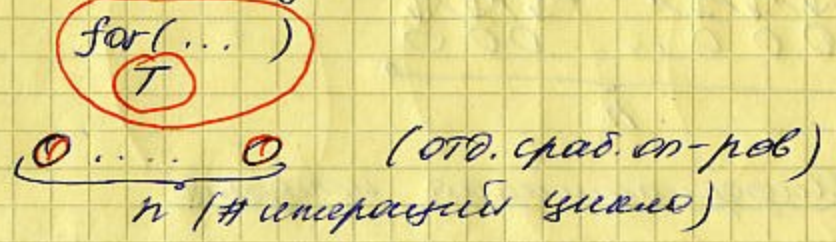
```

for(i=0; i<n; ++i)
  for(j=0; j<n; ++j)
    for(k=0; k<n; ++k)
      A[i][j] = A[i][j] + B[j][k] + C[k][j];
  
```

порядок циклов (i, j, k) ↔ (k, j, i)
 для данной прогн. 4 из 6 порядков допуст.
 преув. можно в ~5-7 раз оптимизир
 в завис. от порядка циклов
 Варианты (k, i, j) и (i, k, j) почти
 лучше всего, т.к. цикл по j один
 внутренний

Графовые модели

Вершины: циклы, операторы, мин.
 участки кода, тело цикла,
 отдельные срабатывания от-б



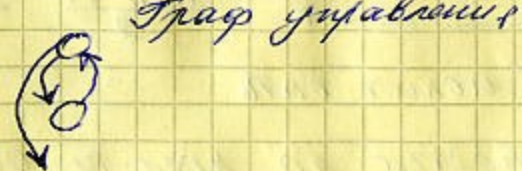
Дуги: - операц. отношения (по пере-
 даче управления)
 $A \rightarrow B \Leftrightarrow B$ н.д. выполнено сразу после A
 - информационное (по передаче
 данных)

4 основные графовые модели:

• Граф управления программы
 (операторы; операц. отношения)

```

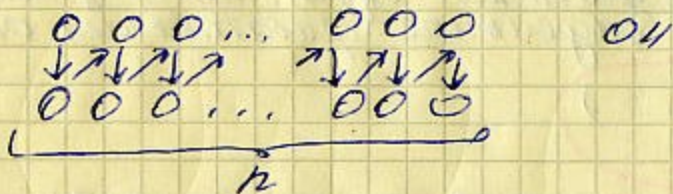
for(i=0; i<n; ++i)
{
  A[i] = A[i-1] + p + 17;
  B[i] = A[i] - t;
}
  
```



• Информационный граф пр-ма
 (операторы; инф. отн-я)

Операционная история

(сработавшие опер-в, операци. отношения)

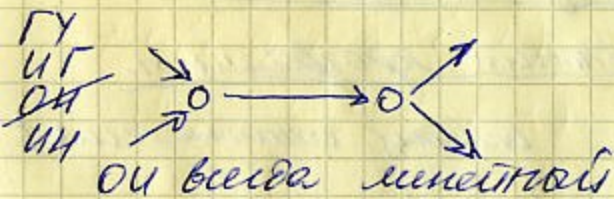


Информационная история

(сработ. опер-в, информаци. отно-я)



ИИ может быть без дуг: $a_i = b_i + c_i$



ОИ всегда линейной

ИИ: завершим, 200 дуг



=> не может быть

ГУ и ИГ сработали по послед-н опер-в

ИИ и ОИ - по таблице

Операци. отношения

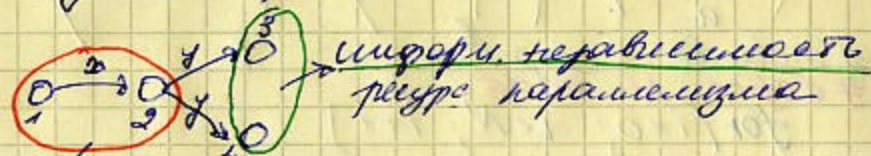


$$x = a + 12i;$$

$$y = x + t - 8;$$

$$t_1 = y + 2j;$$

$$t_2 = y + a - 8;$$



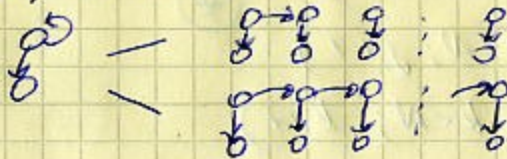
поместить коридор
не можем

Информ. зав-н - критерий отб-в-н предпр.

Критерии:

- компактность (+ ИГ)

- информативность (+ ИИ)



- реализуемость (+ ИГ)

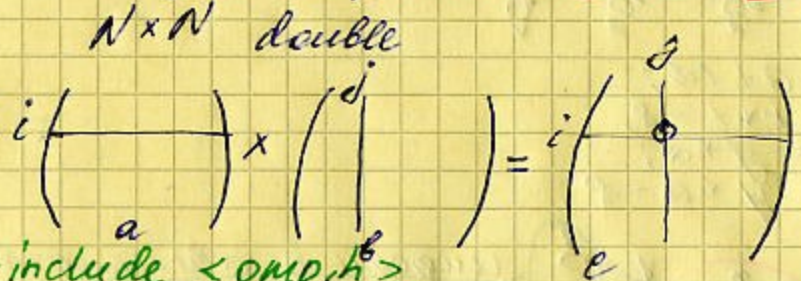
Линейный класс

Граф алгоритма - параметрич. гаран-тая симп. история

+ компактность

+ строгость по тегам пр-т
(реализуемость)

8.12 Задача перемножения матриц



```
#include <omp.h>
#define N 4096
for(i=0; i<N; i++)
{
    for(j=0; j<N; j++)
    {
        a[i][j] = (i+j)*0.0001;
        b[i][j] = (i-j)*0.0001;
        c[i][j] = 0.0;
        t1 = omp_get_wtime();
        #pragma omp parallel shared(a, b, c) private(i, j, k) size(4)
        {
            for(i=0; i<N; i++)
            {
                for(j=0; j<N; j++)
                {
                    for(k=0; k<N; k++)
                    {
                        c[i][j] += a[i][k]*b[k][j];
                    }
                }
            }
        }
    }
}
```

- ресурс параллелизма
- распределение операций
- выбрать технологию
- распредел. данных (MPI)
- передача данных (MPI)

- параллелизм алгоритма
- нагрузка
- эффективность и масштабируемость

```
* size = omp_get_num_threads();
#pragma omp for schedule(static)
// разбивается на ~ равное число
t2 = omp_get_wtime();
for(i=0; s=0.0; i<N; i++)
    for(j=0; j<N; j++)
        s += c[i][j];
printf("N=%d, Nproc=%d, sum=%d, time=%d\n", N, size, s, t2-t1);
```

СКУД МГУ "Чертюков"

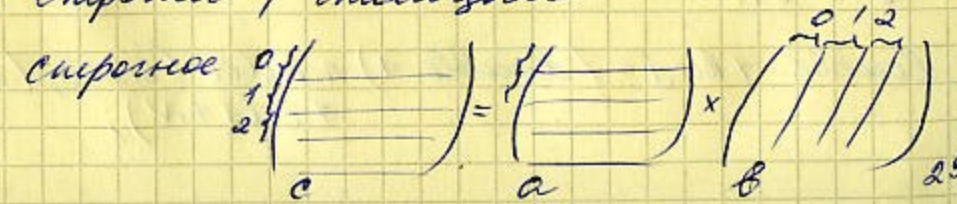
2x4 ядра = 8 ядер

Компьютер Intel 11.1
-03-оригин

# ядер	1	2	4	8	16
время	164	89.8	59.5	36.9	53.4

Распределение данных

- строки / столбцовые



j-й элемент в каждой переменной процессу

всему процессу a_{ij}

Завести матрицу-буфер $n \times s$, $buf()$

Каждый процесс переписать строку в буферный массив

всему процессу: MPI-Beast

```
#include <mpi.h>
```

```
MPI_Init(&argc, &argv);
```

```
MPI_Comm_rank(MPI_COMM_WORLD, &rank);
```

```
MPI_Comm_size(MPI_COMM_WORLD, &size);
```

```
myrows = mycols = rows = cols = (N-1)/size + 1;
```

```
last-rows = last-cols = N - rows * (size - 1);
```

```
if (rank == size - 1) // need process
```

```
myrows = mycols = last-rows;
```

```
double **a = (double **) malloc(sizeof(double) *  
myrows);
```

```
for (i=0; i < myrows; i++)  
a[i] = (double *) malloc(sizeof(double) * N);
```

```
b, c, ...
```

```
double *buf = (double *) malloc(sizeof(double) *  
cols + N);
```

```
/* unimplemented */
```

```
MPI_Barrier(MPI_COMM_WORLD);
```

```
t1 = MPI_Wtime();
```

```
for (m=0; m < size; m++)
```

```
if (rank == m)
```

```
for (j=0; j < mycols; j++)
```

```
for (k=0; k < N; k++)
```

```
buf[k+j*N] = b[k][j];
```

```
if (m == size - 1) l = last-cols + N;
```

```
else l = cols + N;
```

```
MPI_Bcast(buf, l, MPI_Double, m,  
MPI_COMM_WORLD);
```

```
if (m == size - 1)  
jmax = last-cols;  
else  
jmax = cols;
```

```
for (i=0; i < myrows; i++)
```

```
for (j=0; j < jmax; j++)
```

```
for (k=0; k < N; k++)
```

```
c[i][j] += a[i][k] * buf[k+j*N];
```

```
}
```

```
MPI_Barrier(MPI_COMM_WORLD)
```

```
t2 = ...
```



```
for (i=0; sum.=0.0; i<myrows; i++)
```

```
for (j=0; j<N; j++)
```

```
sum += c[i][j];
```

```
MPI_Reduce (&sum, &s, 1, MPI_DOUBLE,
            MPI_SUM, 0, MPI_COMM_WORLD);
```

```
if (rank==0) printf(...
```

#thr	1	2	4	8	16	32	64	128
time	155.3	88.2	79.2	63.2	7.5	2.6	2.34	1.60
	164.1		17.36	10.2 (9#)	8.27	2.5	0.98	0.67
#thr	256	512	1024		(#25)		(#121)	
time	3.77	5.6	12.4					

(заплата на пересылку данных)
 0.99 0.64 0.93 данные для 2 алгоритма
 (#484) (экспер. сум # = n^2)

• выработка стратегии для расщеп-я



Скор. $N^2(2N-1)$

Пересылки size cols * N = $\frac{N^2}{size} + size = N^2$

$N^2(size-1)$ общий объем пересылки для 1 алгоритма

$\left(\frac{N}{\sqrt{size}}\right)^2 = \frac{N^2}{size}$
 (размер пересыл. блока)
 $2(\sqrt{size}-1) * \frac{N^2}{size} + size = 2(\sqrt{size}-1) * N^2$
 для 2 алгоритма

пересылок:

$size + (size-1)$ в 1 алгоритме

$2(\sqrt{size}-1) * size$ во 2 алгоритме

Объем буфера в 1 алг. $\frac{N^2}{size}$

во 2 алг. $2 + \sqrt{size} * \frac{N^2}{size}$

в $2 + \sqrt{size}$ раз больше объем памяти