

# Введение в параллельные вычисления

к.т.н. Алексей Ивахненко



- Что такое параллельные вычисления, насколько они распространены, востребованы и выгодны?
- Как устроены суперкомпьютеры?
- MPI, OpenMP, OpenACC





- 📎 Как работать с большими объемами данных и вычислений?
- Масштабируемость. Закон Амдаля.
- Производительность (FLOPS), пропускная способность памяти.
- Приложения, ограниченные вычислительными возможностями и пропускной способностью памяти.
- 🔊 Частота процессора и частота памяти. Значение кэшей.
- 🔊 Векторизация.
- Наборы векторных инструкций.
- Частота процессора и TDP.
- Примеры многоядерных систем:
  - Intel CPU
  - Cell Broadband Engine
  - Cavium ThunderX

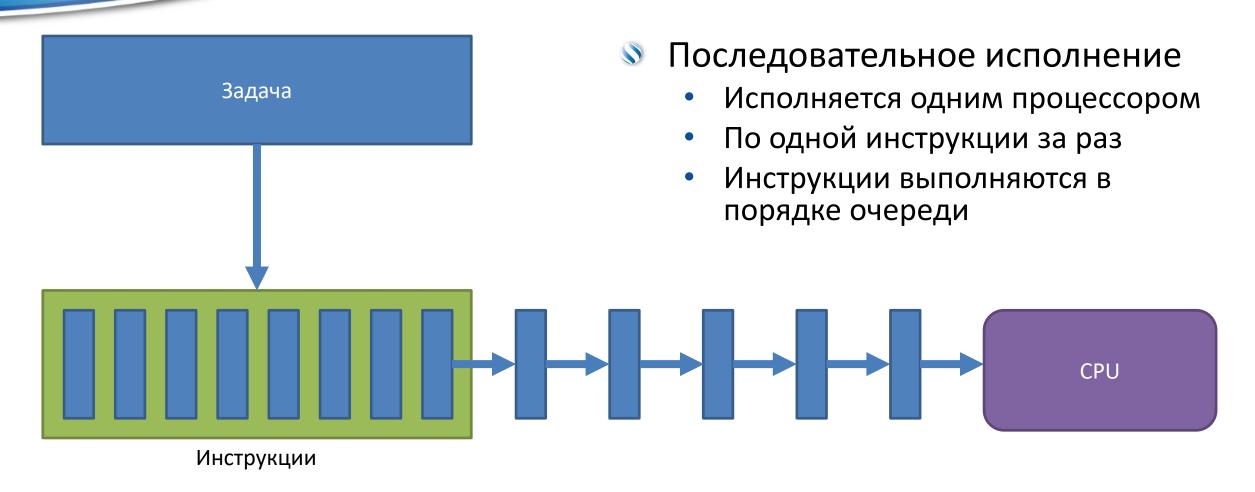




- Roofline модель оценки производительности.
- Распределенная архитектура памяти.
- Важность размещения данных в памяти.
- Схема современного гибридного кластера с множеством сокетов, ядер и ускорителей.
- Классификация Флинна.
- Общая память и распределенная память.
- Ошибки, специфичные для параллельных вычислений: deadlocks, memory races.
- Потоки и процессы.
- Синхронное и Асинхронное исполнение.
- Критические секции, мьютексы и барьеры.
- Сигналы, сообщения, события для коммуникации между процессами.
- Наборы утилит: cpuinfo, Ispci. Управление задачами на сервере с помощью SLURM.



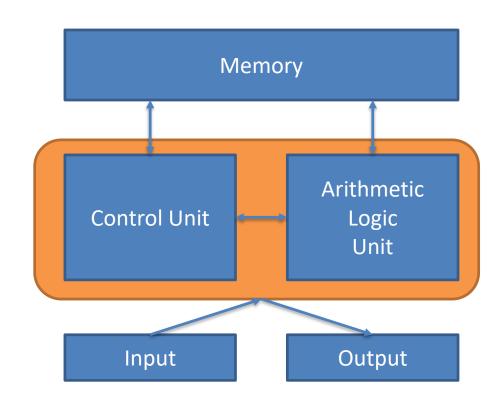
#### Последовательное исполнение





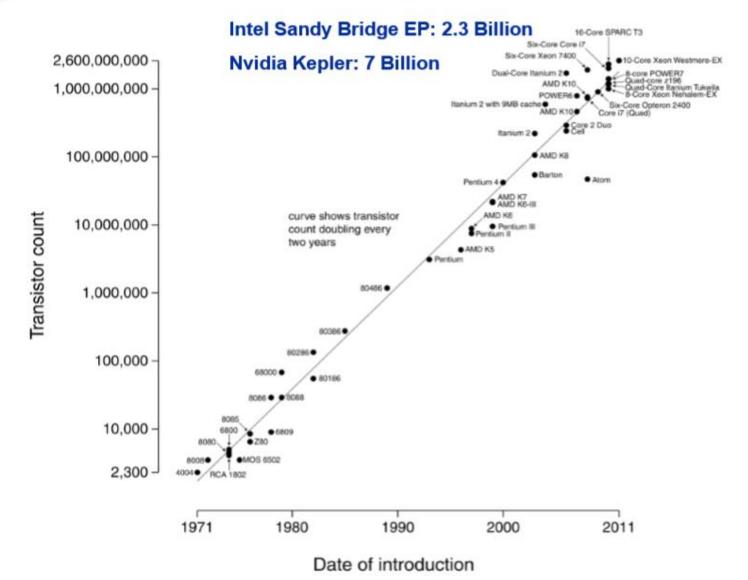
# Архитектура фон Неймана

- Данные и инструкции хранятся в памяти
- Control unit загружает инструкции и данные из памяти, декодирует инструкции и последовательно выполняет операции
- АЛУ выполняет базовые арифметические операции





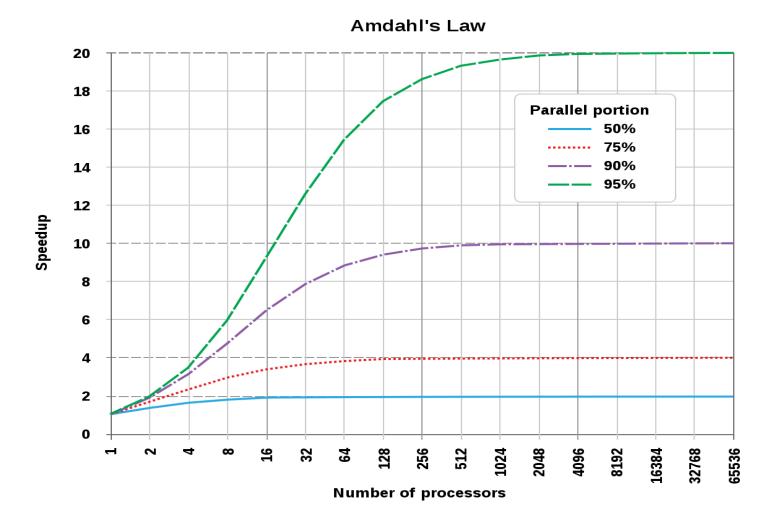
## Закон Мура





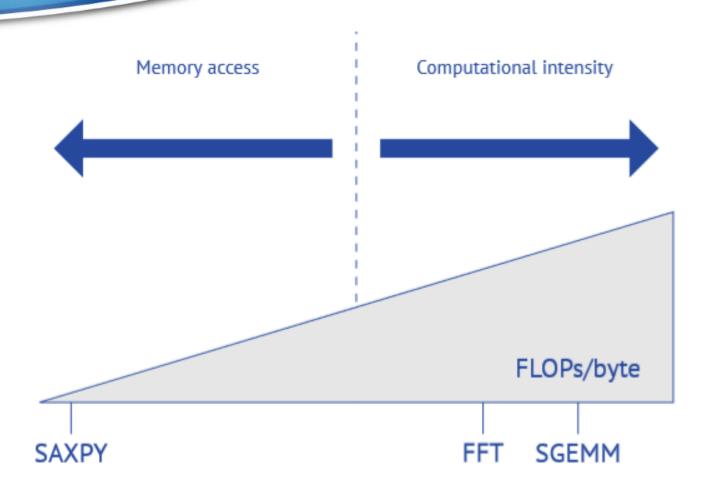
# Закон Амдаля

S(p)≤1/(f + (1 -f)/p) p – количество процессоров f – доля программы, которая должна быть выполнена последовательно S - ускорение





## Ограничения



Optimal computational intensity for different GPUs (FLOPS : byte): Tesla C1060 9.1:1 Tesla C2050 7.1:1 16.9:1 Tesla K20 17.5:1 Tesla K40 23.7:1 Tesla M40 Tesla P100 12.9:1 NVIDIA Titan X 22.8:1



# **Top Supercomputers**

System	Cores	Rmax (TFlop/s)		Power (kW)
Sunway TaihuLight - Sunway MPP, Sunway SW26010 260C 1.45GHz, Sunway,		93,014.60		•
NRCPC	10,010,000	33,011.00	123, 133.30	13,371
National Supercomputing Center in Wuxi				
China				
Tianhe-2 (MilkyWay-2) - TH-IVB-FEP Cluster, Intel Xeon E5-2692 12C 2.200GHz, TH	3,120,000	33,862.70	54,902.40	17,808
Express-2, Intel Xeon Phi 31S1P, NUDT				
National Super Computer Center in Guangzhou				
China				
Piz Daint - Cray XC50, Xeon E5-2690v3 12C 2.6GHz, Aries interconnect, NVIDIA	361,760	19,590.00	25,326.30	2,272
Tesla P100, Cray Inc.				
Swiss National Supercomputing Centre (CSCS)				
Switzerland				
Gyoukou - ZettaScaler-2.2 HPC system, Xeon D-1571 16C 1.3GHz, Infiniband EDR,	19,860,000	19,135.80	28,192.00	1,350
PEZY-SC2 700Mhz , ExaScaler				
Japan Agency for Marine-Earth Science and Technology				
Japan				
Titan - Cray XK7, Opteron 6274 16C 2.200GHz, Cray Gemini interconnect, NVIDIA	560,640	17,590.00	27,112.50	8,209
K20x, Cray Inc.				)
DOE/SC/Oak Ridge National Laboratory				
United States				



### Методика тестирования

- Бенчмарк HPL (High-Performance Linpack): выбирается наибольший размер задачи, помещающийся в память.
- Предоставляет следующую информацию
  - Rpeak: теоретический максимум производительности (FLOPS) системы (FLOPS\*Количество тактов CPU\*количество CPU)
  - Rmax: Максимум, достигнутый для заданного размера задачи.
  - Nmax: размер матрицы
  - N1/2: Размер задачи при котором достигается 50% Rmax. Низкое значение N1/2 показывает способность системы работать со множеством различных размеров задачи (меньше-лучше).

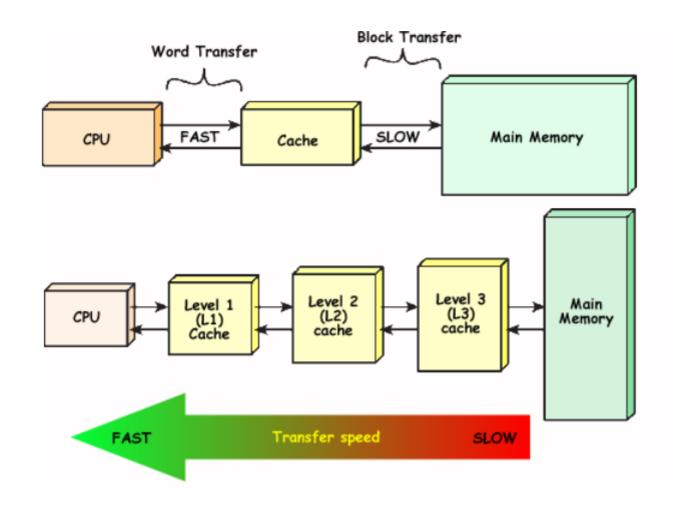


#### Основные ограничители

- Вычислительные возможности
  - FLOPS
- Пропускная способность памяти
  - Чтение/запись
  - Обмен сообщениями и сигналами
  - Кэши
- Латентность

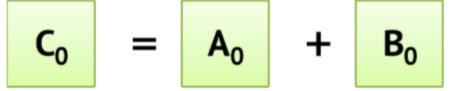


#### Кэш



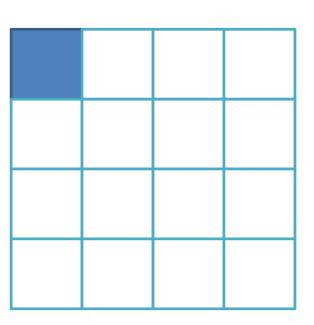


### Последовательное исполнение



$$C_1 = A_1 + B_1$$

$$C_{n-1} = A_{n-1} + B_{n-1}$$

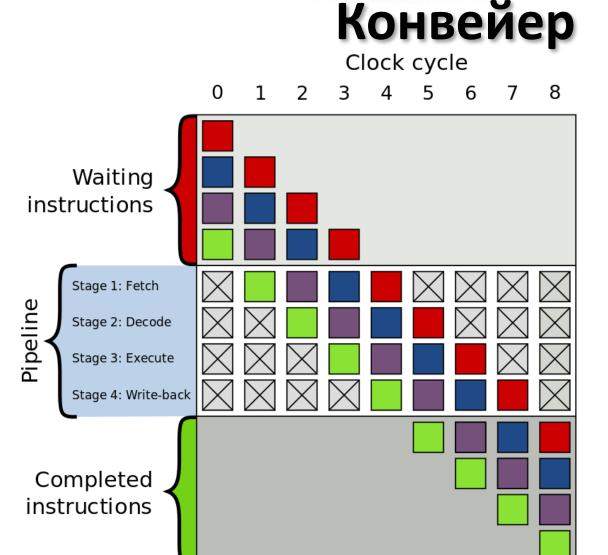


#### Последовательное исполнение

```
multiply_and_add(const float* a, const float* b, const float* c, float* d)
{
    for(int i=0; i<8; i++)
    {
        d[i] = a[i] * b[i];
        d[i] = d[i] + c[i];
    }
}</pre>
```



- Общий конвейер с четырьмя стадиями работы:
  - Получение
  - Раскодирование
  - Выполнение
  - Запись результата

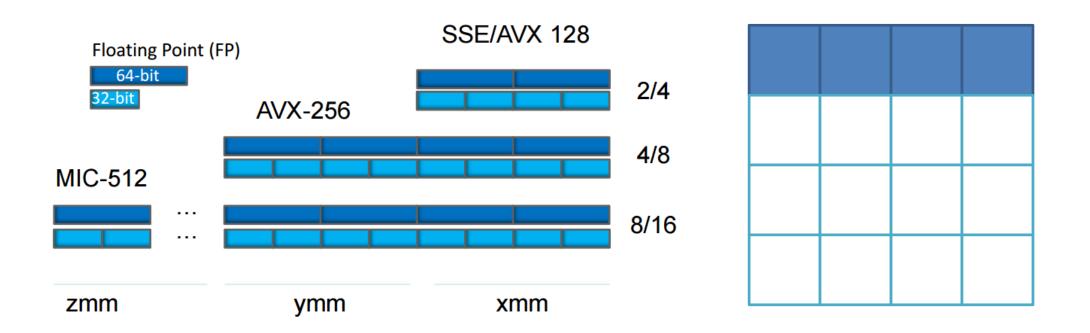




```
multiply_and_add(const float* a, const float* b, const float* c, float* d)
#pragma unroll 4
for(int i=0; i<8; i++)
  d[i] = a[i] * b[i];
  d[i] = d[i] + c[i];
 return _mm256_fmadd_ps(a, b, c);
```

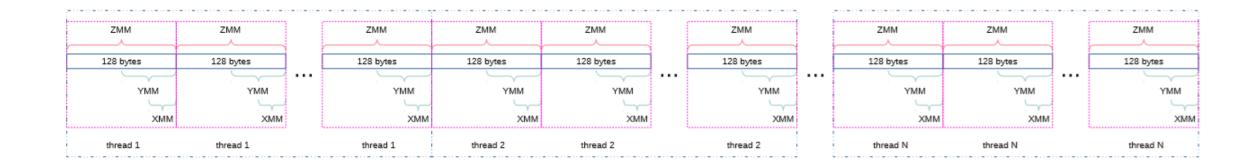


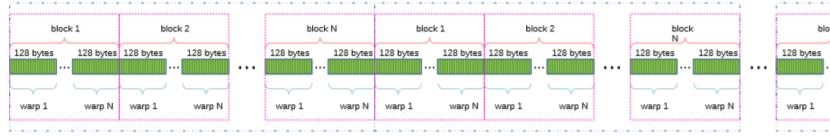
### Векторные инструкции

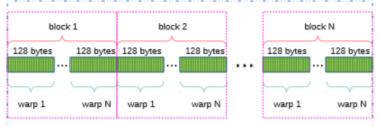




#### Векторные инструкции







I - thread (lightweight GPU-thread)

#### Векторизация

```
multiply_and_add(const float* a, const float* b, const float* c, float* d)
for(int i=0; i<8; i++)</pre>
  d[i] = a[i] * b[i];
  d[i] = d[i] + c[i];
 return _mm256_fmadd_ps(a, b, c);
```



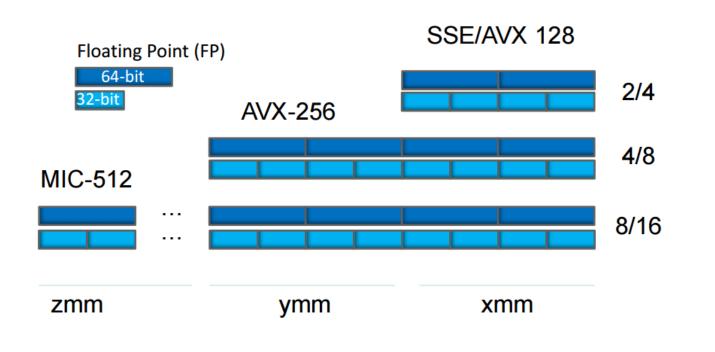
# AVX/AVX2

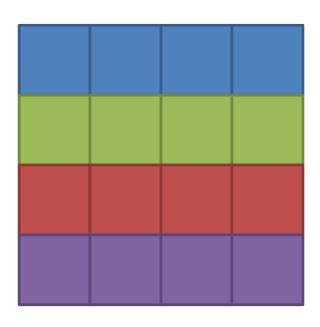
- В общем виде инструкция AVX/AVX2 называется следующим образом:
  - mm<bit width> <name> <data type>
- - ps vectors contain floats (ps stands for packed singleprecision)
  - pd vectors contain doubles (pd stands for packed doubleprecision)
  - epi8/epi16/epi32/epi64 vectors contain 8-bit/16-bit/32bit/64-bit signed integers
  - epu8/epu16/epu32/epu64 vectors contain 8-bit/16-bit/32bit/64-bit unsigned integers
  - si128/si256 unspecified 128-bit vector or 256-bit vector
  - m128/m128i/m128d/m256/m256i/m256d identifies input vector types when they're different than the type of the returned vector

Data Type	Description
m128	128-bit vector
	containing 4 floats
m128d	128-bit vector
	containing
	2 doubles
m128i	128-bit vector
	containing integers
m256	256-bit vector
	containing 8 floats
m256d	256-bit vector
	containing
	4 doubles
m256i	256-bit vector
	containing integers



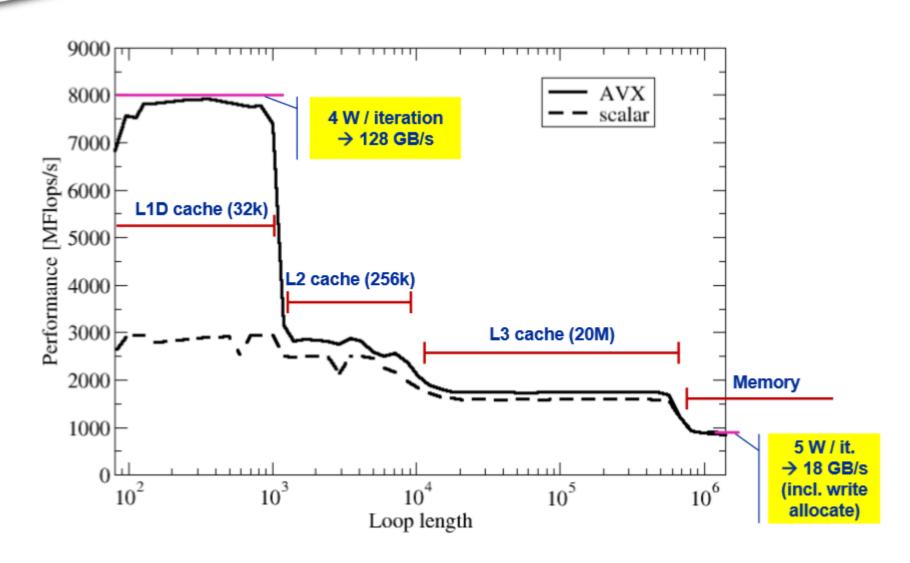
# Векторные инструкции и параллельное исполнение





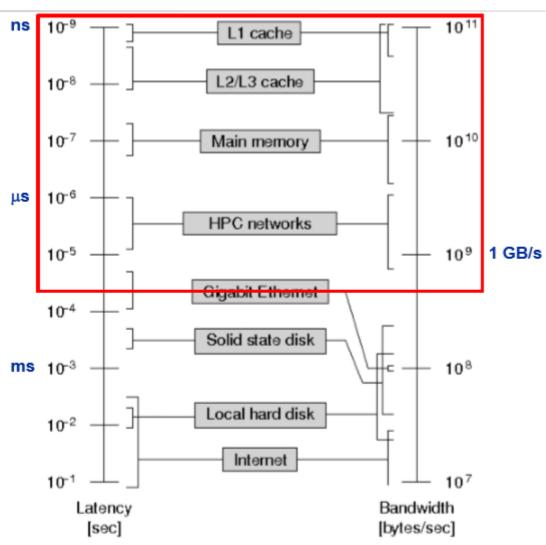


#### Влияние кэшей





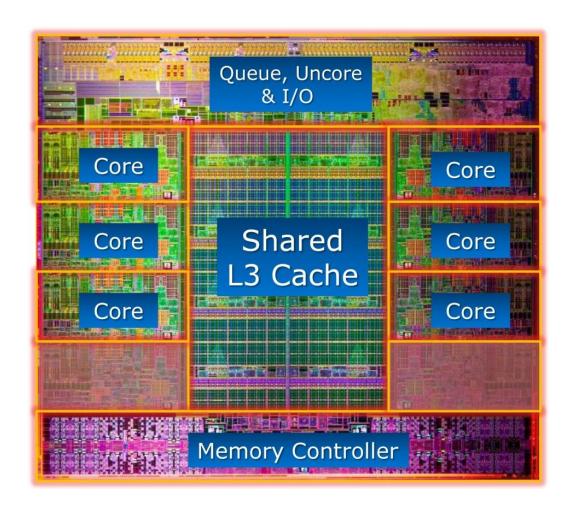
#### Латентность





#### Core i7

- Несколько мощных ядер
  - 2,4,6,8 ядер, 2,66—3,6ГГц каждое
  - Каждое физическое ядро представляется в системе как 2 логических и может исполнять 2 потока конкурентно (Hyper-Threading)
- 3 уровня кэша, большой объем L3
  - На ядро: L1=32KB (data)
     + 32KB (Instructions), L2=256KB
  - Shared L3, до 20 MB
- Запросы к памяти управляются раздельно для каждого потока/процесса

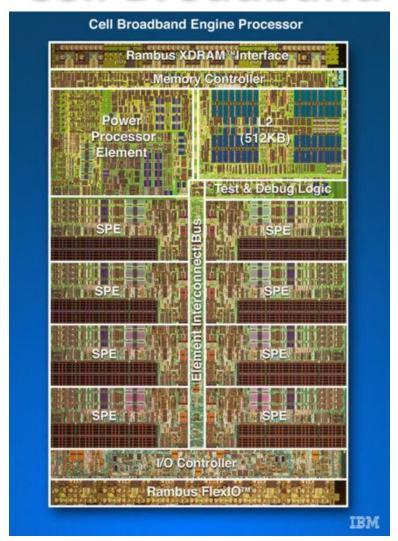




# Cell (Cell Broadband Engine Architecture)

- микропроцессорная архитектура, совместно разработанная Sony, Toshiba и IBM.
- Разработка архитектуры и первые прототипы были созданы в STI Design Center за четырёхлетний период с начала марта 2001 года.
- Cell совмещает ядро общего назначения архитектуры POWER с сопроцессорами, которые значительно ускоряют обработку мультимедиа и векторных вычислений.

#### **Cell Broadband**





#### **Cavium ThunderX**





#### **Nvidia Volta V100**

#### Volta SMM Specifications:

- 84 SMs
- 4096-bit HBM2 memory interface
- 6144 KB L2 cache
- NVLink support
- 12nm FFN
- 21.1 billion transistors
- 7,5 TFLOPS of double precision
- 15 TFLOPS of single precision
- 30 TFLOPS of halfprecision





#### **Nvidia Volta V100**

# Volta SMM Specifications:

- Up to 5376 CUDA cores
- 14336 KB register file size/GPU
- Up to 128 KB of shared memory



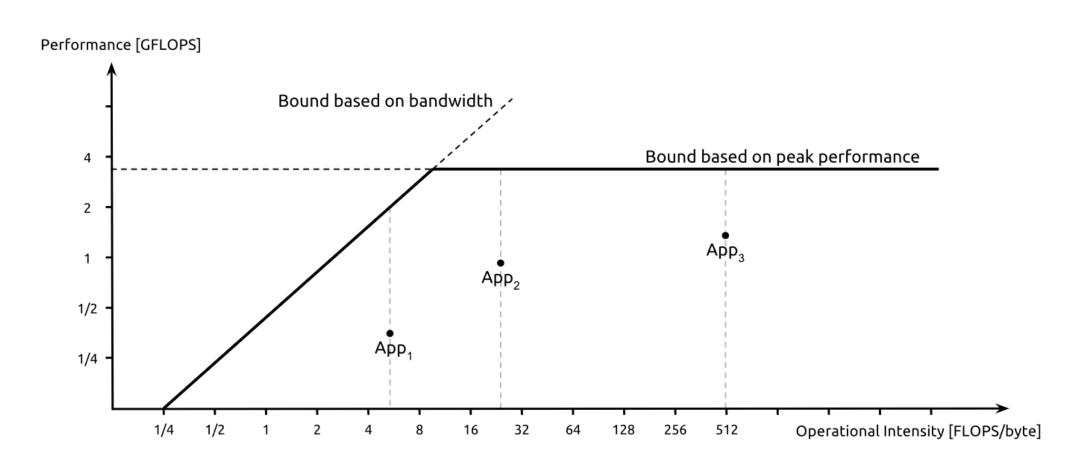




- Noofline модель оценки производительности.
- Распределенная архитектура памяти.
- Важность размещения данных в памяти.
- Схема современного гибридного кластера с множеством сокетов, ядер и ускорителей.
- 📎 Классификация Флинна.
- 🔊 Общая память и распределенная память.
- 🔊 Ошибки, специфичные для параллельных вычислений: deadlocks, memory races.
- 🔊 Потоки и процессы.
- Синхронное и Асинхронное исполнение.
- Критические секции, мьютексы и барьеры.
- 🔊 Сигналы, сообщения, события для коммуникации между процессами.
- 🔊 Наборы утилит: cpuinfo, lspci. Управление задачами на сервере с помощью SLURM.

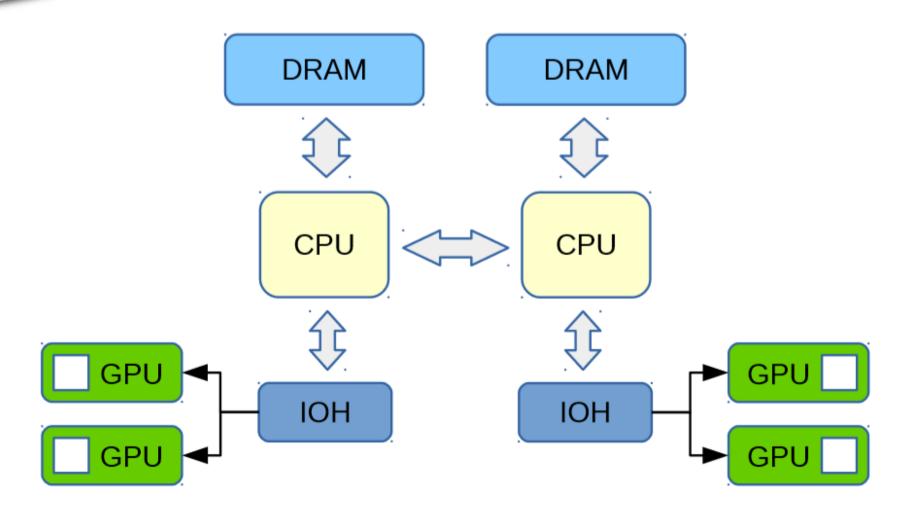


# Roofline модель



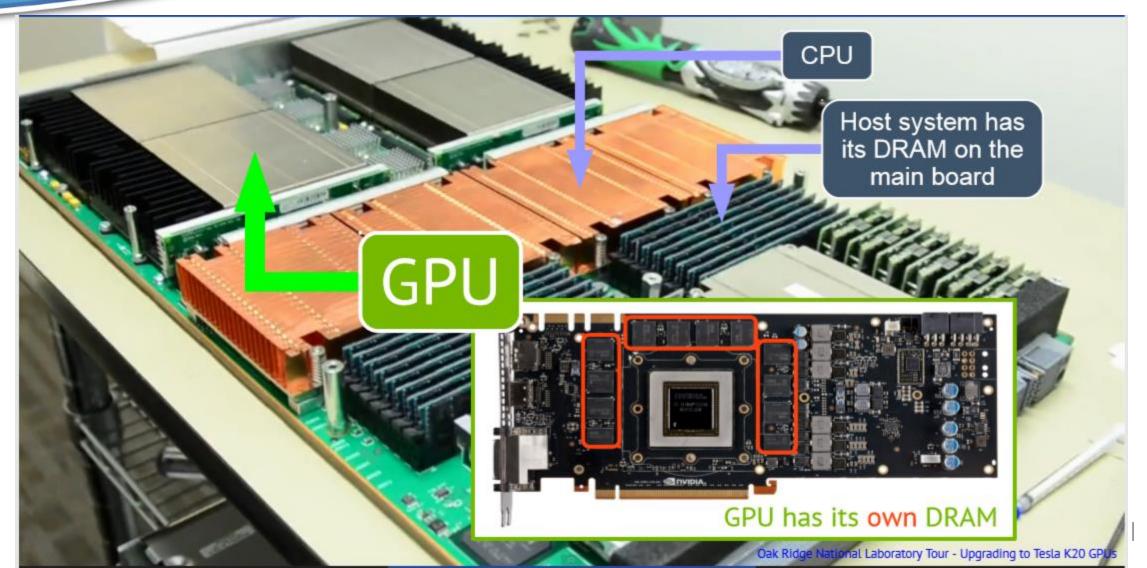


# Современный кластер



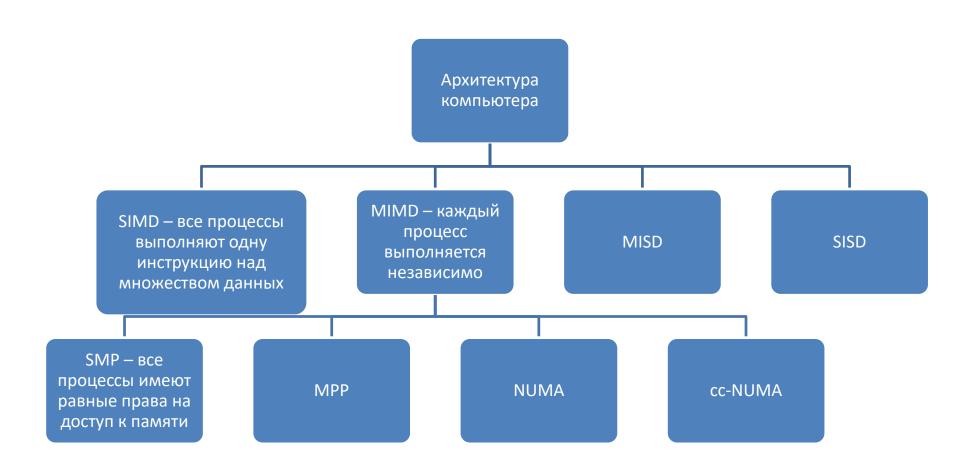


# Современный кластер





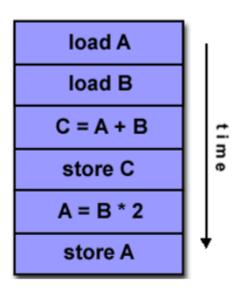
# Классификация Флинна







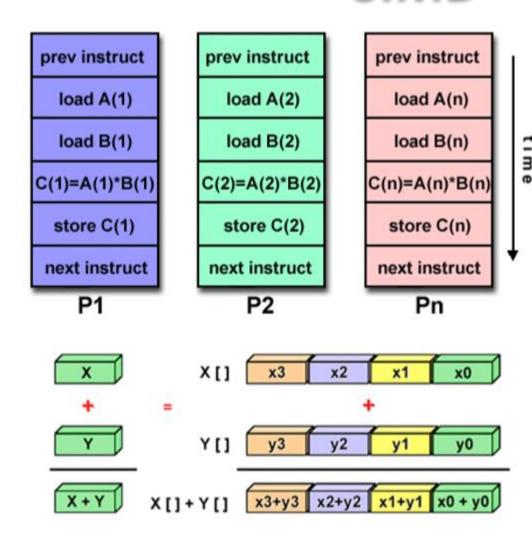
- Последовательный вычислитель
- Детерминированное исполнение





SIMD

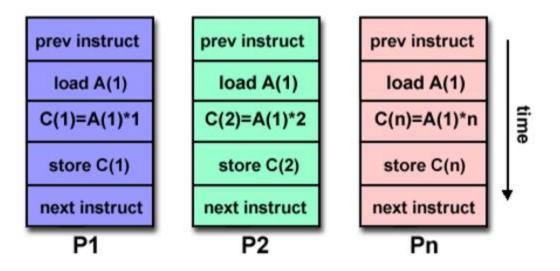
- Все вычислители исполняют одну и ту же инструкцию синхронно по тактам
- Каждый вычислитель может оперировать отдельным данными
- 👂 2 варианта:
  - Набор процессоров
  - Векторное исполнение
- Большинство современных компьютеров относятся к SIMD





## MISD

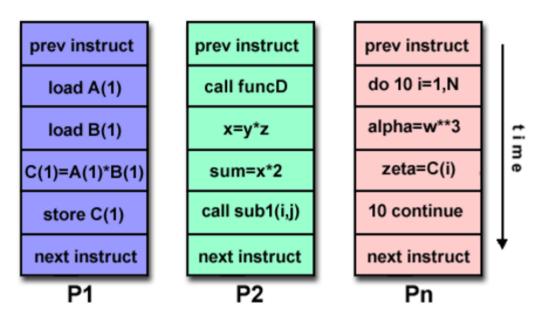
- Поток данных обрабатывается множеством процессоров
- Каждый юнит обрабатывает данные независимо с помощью собственного потока инструкций.
- Примеров не так много: Carnegie-Mellon C.mmp computer (1971).







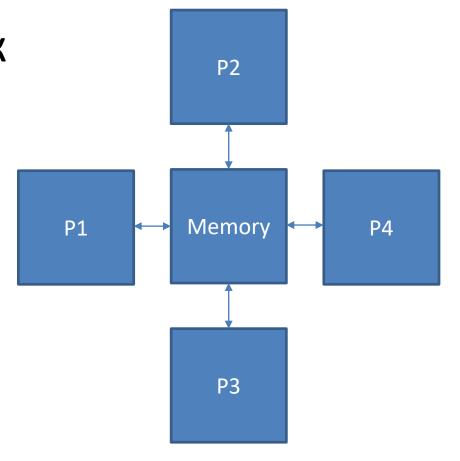
- Самый распространенный тип параллельного компьютера
- Каждый процессор может исполнять свой собственный поток инструкций
- Исполнение может быть синхронным и асинхронным, детерминированным и недетерминированным
- 🔰 Примеры:
  - Почти все современные суперкомпьютеры, многопроцессорные системы, многоядерные процессоры
- Многие MIMD архитектуры состоят из SIMD модулей





# Общая память

- Все процессоры имеют доступ к памяти, как глобальному адресному пространству
- Процессоры работают независимо, но делят ресурсы памяти





# Общая память

## Преимущества

- Глобальное адресное пространство упрощает программирование
- Позволяет распараллеливать код по частям
- Унифицирует и ускоряет обмен данными между CPU

## Недостатки

- Память не масштабируется вместе с количеством процессоров
- Накладные расходы на организацию общего доступа к памяти
- Программист сам отвечает за синхронизацию доступов к памяти
- Soaring expense of internal network.



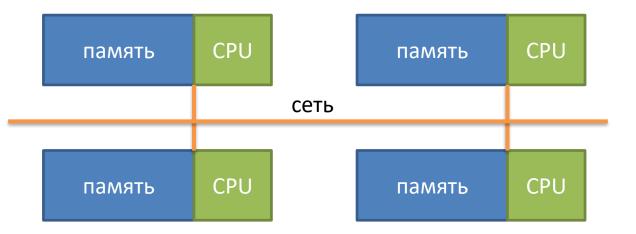
# Виды систем с общей памятью

- Uniform memory access (UMA)
  - Так же называется Symmetric Multi-Processors (SMP)
  - Процессоры равны между собой
  - Время доступа из всех Pn к памяти одинаково
  - Cache Coherent:
    - ✓ Если один процессор обновляет значения в памяти, они обновляются для всех процессоров.
- Non-Uniform memory access (NUMA)
  - Достигается с помощью соединения нескольких SMP
  - Один SMP имеет прямой доступ к памяти другого
  - Не все процессоры имеют одинаковое время доступа
  - Внутри одного SMP доступ к памяти быстрее
  - Медленный доступ к памяти по сети
  - Дополнительная работа для обеспечения когерентности (CC-NUMA)



## Распределенная память

- Каждый процессор обладает своей отдельной памятью
  - Нет глобального адресного пространства
  - Для коммуникации между процессорами используется сеть





## Распределенная память

## Преимущества

- Размер памяти увеличивается с количеством процессоров
- Быстрый доступ к памяти (кроме сети)
- Cost effective (commodity components)

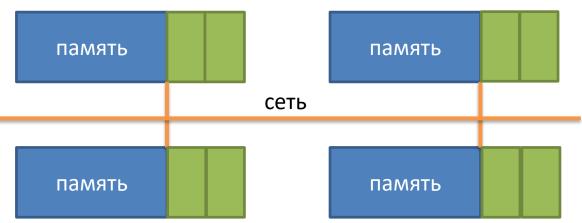
#### Недостатки

- Программист отвечает за коммуникацию
- Преобразование структур данных, основанных на глобальной памяти в распределенную потребует дополнительных усилий
- Неравномерное время доступа к памяти. Зависит от латентности, пропускной способности и загрузки сети
- Сложность поэтапного распараллеливания



# Гибридная модель

- Наиболеераспространенная модель
- Компоненты общей памяти - СС-UMA SMP.
- Глобальные адресные пространства внутри каждого SMP





# Модель общей памяти

- Параллельные задачи используют общее глобальное адресное пространство
- 👀 Операции чтения и записи могут происходить асинхронно
- 🔊 Locks и semaphors для управления доступом к памяти
  - Следует избегать одновременного чтения и записи в одну и ту же ячейку памяти
  - Следует избегать нескольких одновременных операций записи в одну и ту же ячейку памяти
- Компилятор транслирует переменные в адреса глобальной памяти
- Nonbsoвaтель определяет private и shared переменные
- Можно распараллеливать по частям



# Синхронизация

## 🔊 Барьер

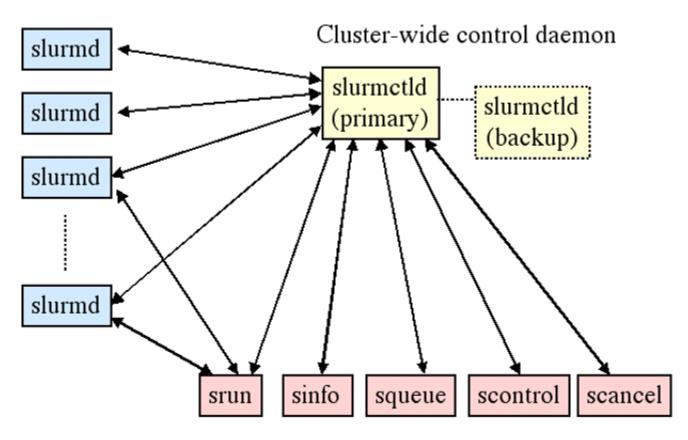
- Каждый поток выполняет свою работу пока не встретит барьер.
   После этого он останавливается (или «блокируется»)
- Когда все потоки достигли барьера, они синхронизированы и продолжают работы с этой точки

## Lock / semaphore

- Первый поток, занявший lock «отмечает» его. Этот поток может безопасно обращаться к защищенной памяти
- Остальные потоки, пытающиеся отметить lock ожидают, когда lock будет освобожден
- Может быть блокирующим или неблокирующим



#### One daemon per node



User and administrator tools



# **SLURM**

#### 2 daemons

- slurmctld controller, optional backup
- slurmd computer node daemon

## 🔊 Команды

- scontrol administration tool, get/set configuration
- sinfo reports general system information
- squeue reports job and job step information
- srun submit/initiate job or job step
- scancel signal or cancel a job or job step