



CS Space
Клуб технологий и науки

Видеокарты: что они могут?
Могут ли они хоть что-то?
Давайте выясним!



Vulkan

OpenCL™

NVIDIA
CUDA



@UnicornGlade



@PolarNick239



polarnick239@gmail.com



Николай Полярный

Activate Ubuntu
Go to Settings to activate Ubuntu.

План лекции

1) **Вычисления:**

shared instruction pointer, code divergence

2) **Работа с памятью:**

hyper-threading, latency hiding, occupancy, registers pressure/spilling, coalesced memory access pattern, cache lines, local/shared memory

3) **Общая картина ЭВМ-архитектуры: CPU - RAM - PCI-E - VRAM - GPU**

4) **Модель вычислений массового параллелизма**

5) **Профилирование и оптимизация GPGPU-алгоритма**

6) **Примеры: A+B, максимум по массиву, merge-sort**

7) **Матрицы: транспонирование, умножение, tensor cores, DeepSeek**

8) **Ray Tracing: real-time BVH, ray tracing cores**

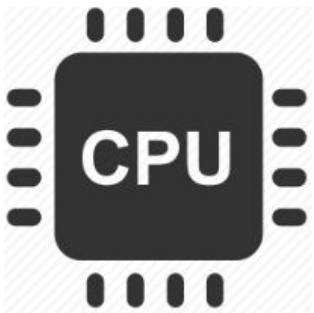
9) **Выводы: какие алгоритмы ускоряются на GPU? OpenCL, CUDA или Vulkan?**

Глава 1: Вычисления

shared instruction pointer, code divergence

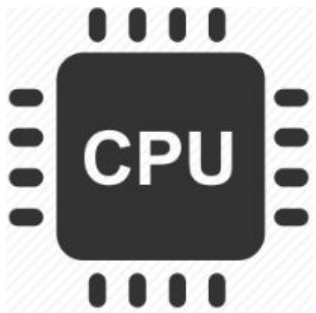
Архитектура

$100 \cdot 10^9$ FLOPS (**F**loating-point operations **p**er **s**econd)

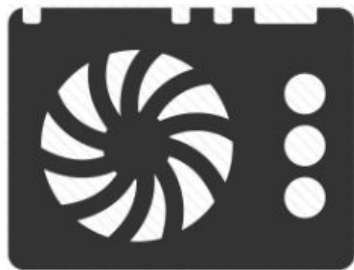


Архитектура

$100 \cdot 10^9$ FLOPS (**F**loating-point operations **p**er **s**econd)



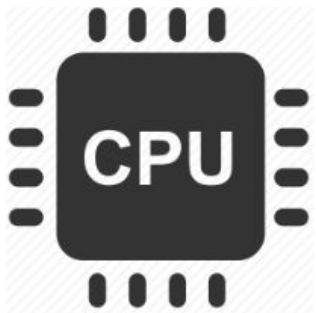
$100 \cdot 10^{12}$ FLOPS (**х1000 раз больше**)



GPU

Архитектура

$100 \cdot 10^9$ FLOPS



40 GB/s memory bandwidth

$100 \cdot 10^{12}$ FLOPS (**х1000 раз больше**)

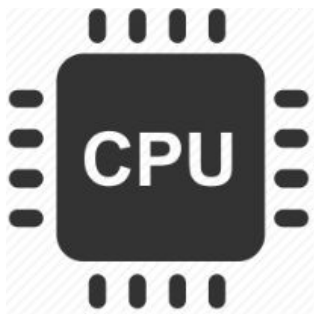


1000 GB/s (**х25 раз больше**)

GPU

Архитектура

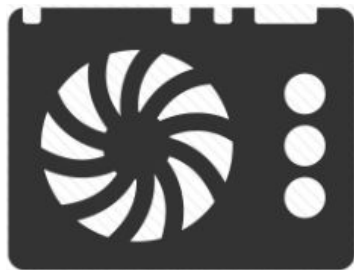
$100 \cdot 10^9$ FLOPS



Мало ядер, но они **МОЩНЫЕ**



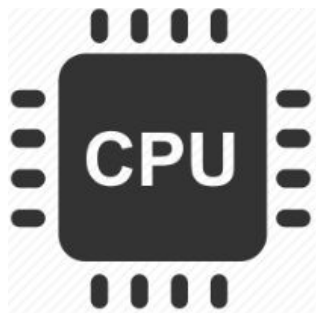
$100 \cdot 10^{12}$ FLOPS



GPU

Архитектура

$100 \cdot 10^9$ FLOPS



$100 \cdot 10^{12}$ FLOPS



GPU

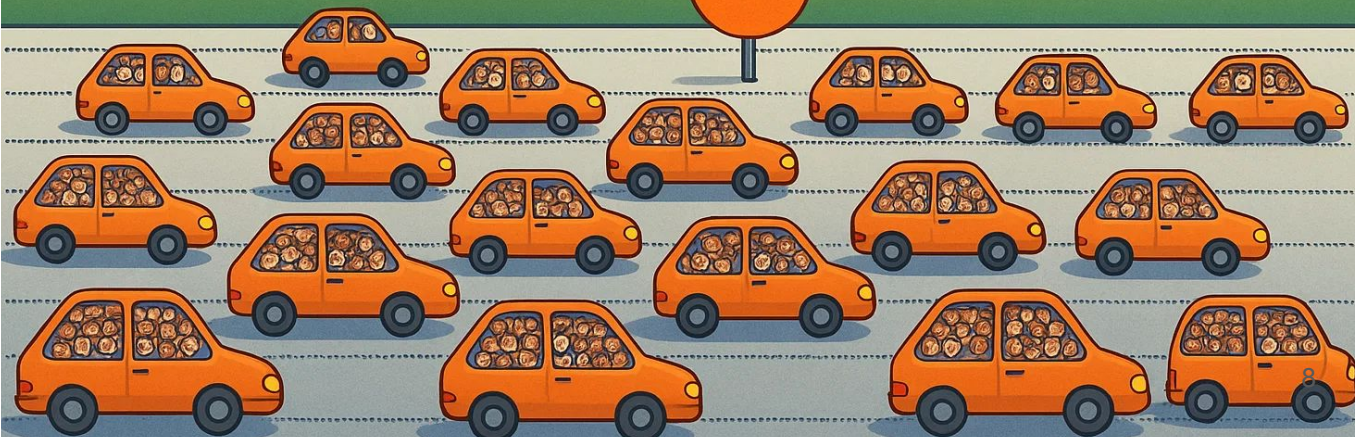
Мало ядер, но они **МОЩНЫЕ**

6 GHz



ТЫСЯЧИ слабых ядер

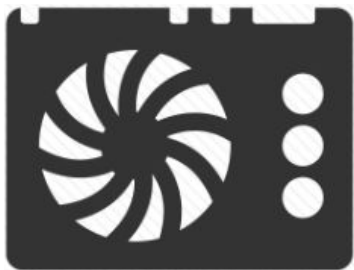
1 GHz



Архитектура

Как уместить ТЫСЯЧИ ядер?
На CPU ведь это почему-то невозможно!

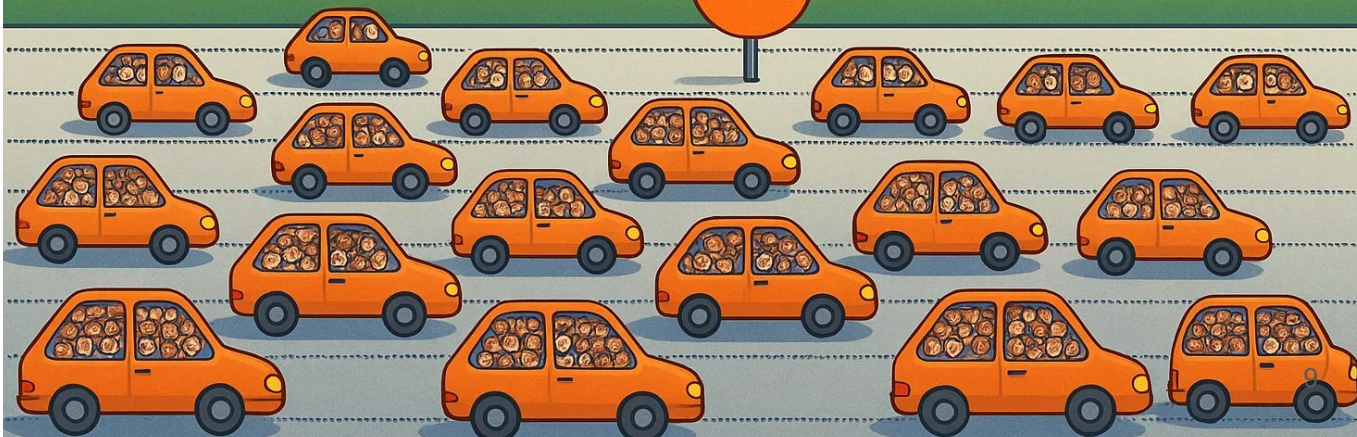
$100 \cdot 10^{12}$ FLOPS



GPU

ТЫСЯЧИ слабых ядер

1 GHz



Архитектура

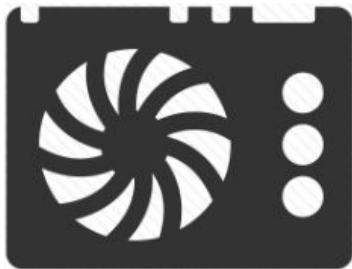


warp

32 слабых медленных
CUDA ядер

Как уместить **ТЫСЯЧИ** ядер?
На CPU ведь это почему-то невозможно!

$100 \cdot 10^{12}$ FLOPS



GPU



ТЫСЯЧИ слабых ядер

1 GHz

SM - Streaming Multiprocessor



Архитектура

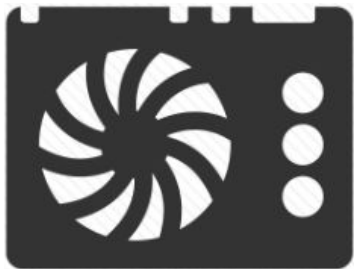


warp

32 слабых медленных
CUDA ядер

Как уместить ТЫСЯЧИ ядер?
На CPU ведь это почему-то невозможно!

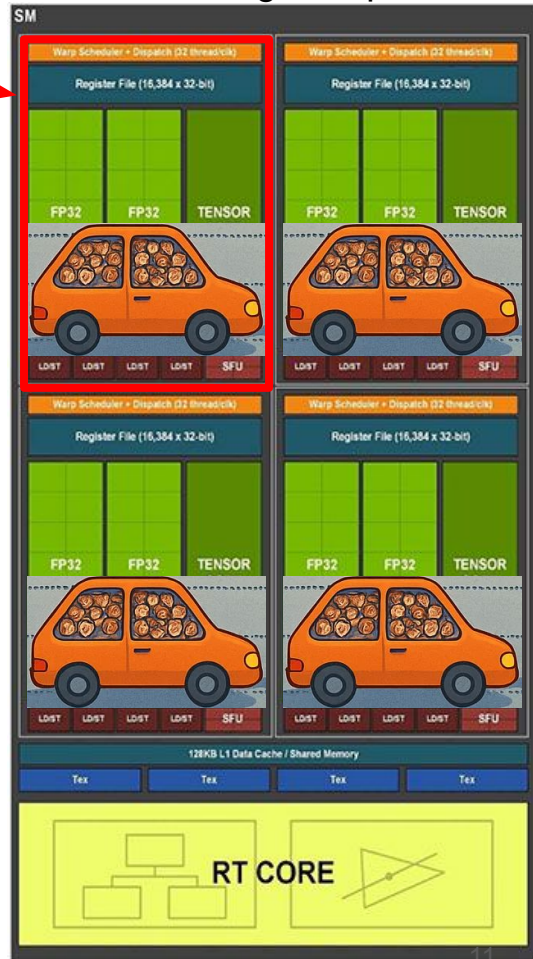
$100 \cdot 10^{12}$ FLOPS



GPU



SM - Streaming Multiprocessor



RTX 3090: **10496 CUDA cores** = 82 SM · 4 warps · 32 ALUs



SM - Streaming Multiprocessor

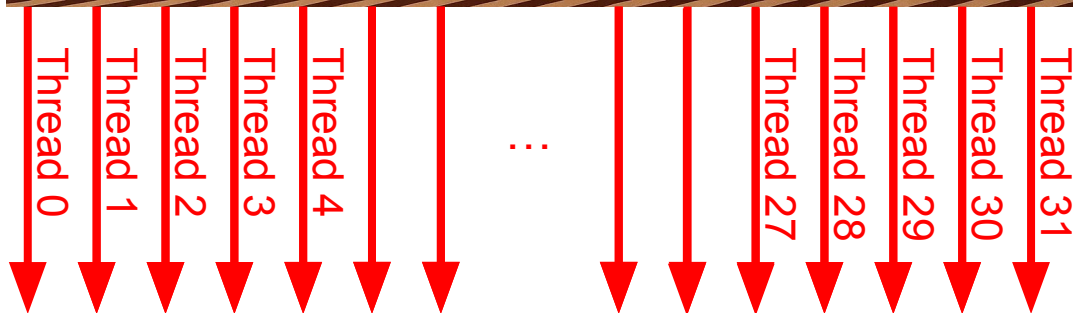


<https://www.techpowerup.com/review/nvidia-geforce-ampere-architecture-board-design-gaming-tech-software/3.html>

CPU ядро



GPU warp



Как уместить ТЫСЯЧИ ядер?
На CPU ведь это почему-то невозможно!

Архитектура

GPU warp



Thread 0

Thread 1

Thread 2

Thread 3

Thread 4

⋮

Thread 27

Thread 28

Thread 29

Thread 30

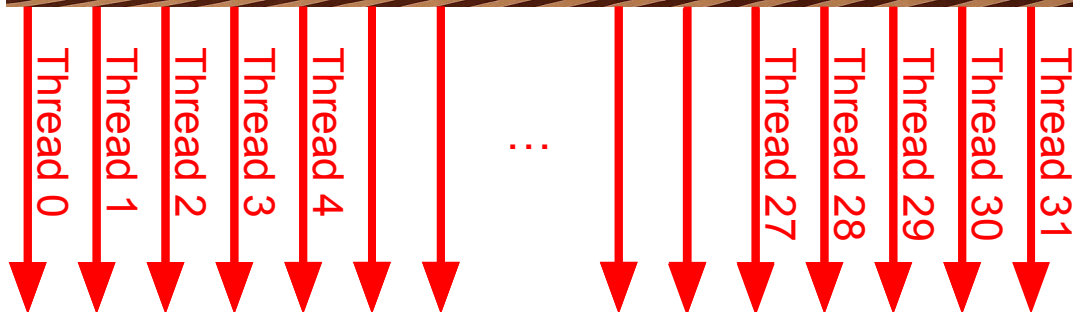
Thread 31

```
239 → int warpThreadID = threadIdx.x % 32;  
240     int result = 0;  
241     if (warpThreadID < 16) {  
242         result = dataA[warpThreadID];  
243     } else {  
244         result = dataB[warpThreadID];  
245     }
```

Один Instruction pointer
на все потоки warp-a!

Архитектура

GPU warp

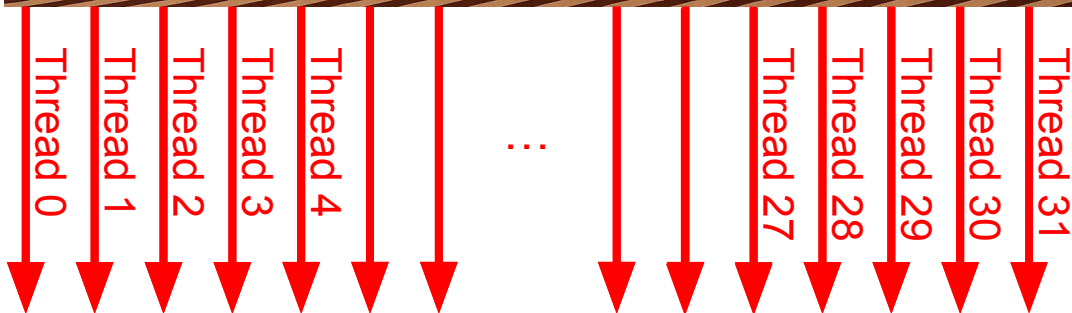


Один Instruction pointer
на все потоки warp-a!

```
239 int warpThreadID = threadIdx.x % 32;  
240 → int result = 0;  
241 if (warpThreadID < 16) {  
242     result = dataA[warpThreadID];  
243 } else {  
244     result = dataB[warpThreadID];  
245 }
```

Архитектура

GPU warp

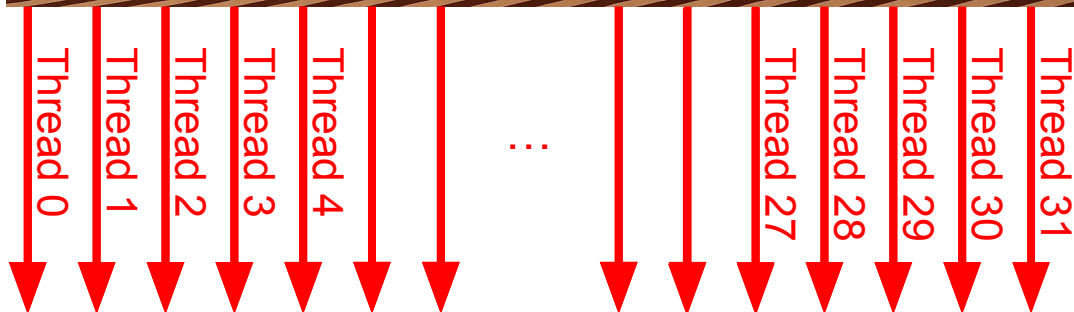


Один Instruction pointer
на все потоки warp-a!

```
239 int warpThreadID = threadIdx.x % 32;  
240 int result = 0;  
241 → if (warpThreadID < 16) {  
242     result = dataA[warpThreadID];  
243 } else {  
244     result = dataB[warpThreadID];  
245 }
```


Архитектура

GPU warp

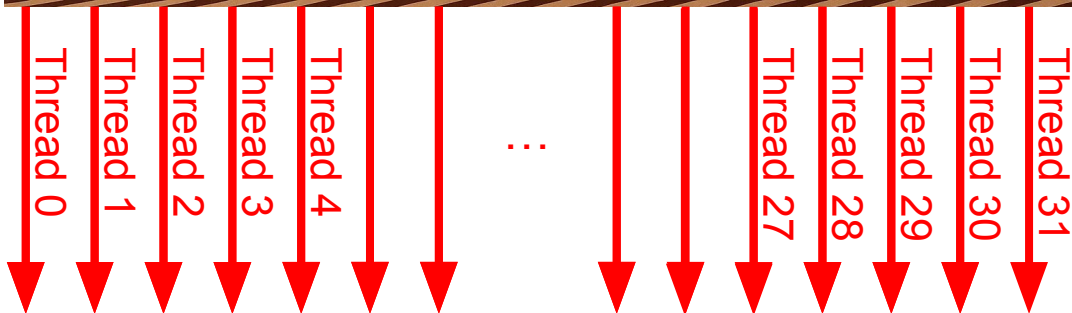


Один Instruction pointer
на все потоки warp-a!

```
239 int warpThreadID = threadIdx.x % 32;  
240 int result = 0;  
241 if (warpThreadID < 16) {  
242     result = dataA[warpThreadID];  
243 } else {  
244     result = dataB[warpThreadID];  
245 }
```

Архитектура

GPU warp

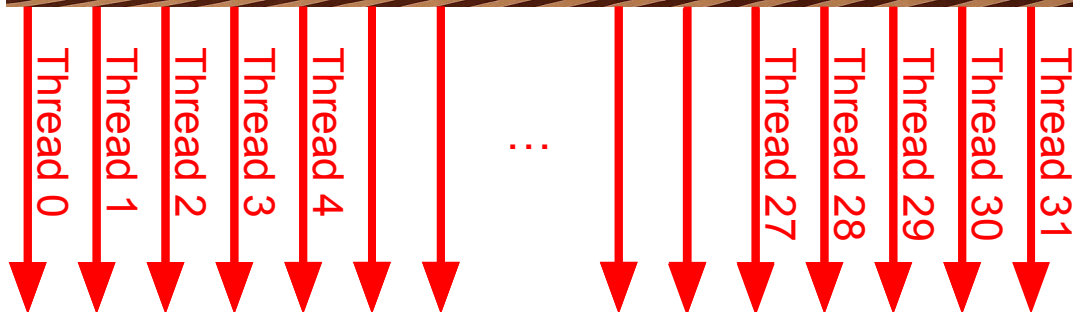


Один Instruction pointer
на все потоки warp-a!

```
239 int warpThreadID = threadIdx.x % 32;  
240 int result = 0;  
241 if (warpThreadID < 16) {  
242     result = dataA[warpThreadID];  
243 } else {  
244     result = dataB[warpThreadID];  
245 }
```

Архитектура

GPU warp

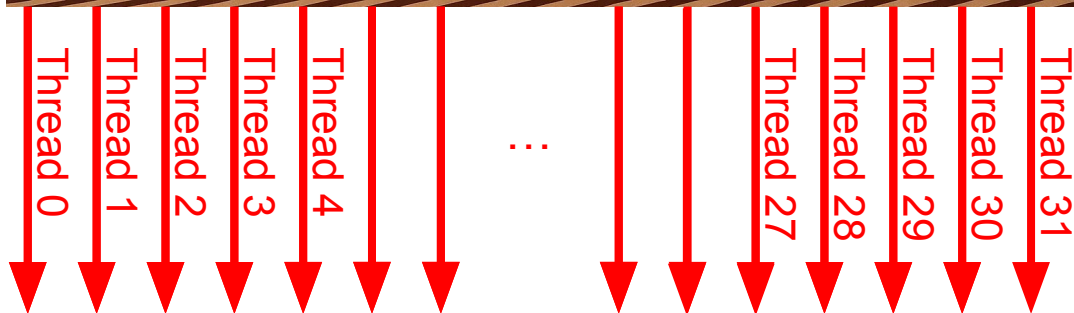


Один Instruction pointer
на все потоки warp-a!

```
239 int warpThreadID = threadIdx.x % 32;  
240 int result = 0;  
241 if (warpThreadID < 16) {  
242     result = dataA[warpThreadID];  
243 } else {  
244 → result = dataB[warpThreadID];  
245 }
```

Архитектура

GPU warp



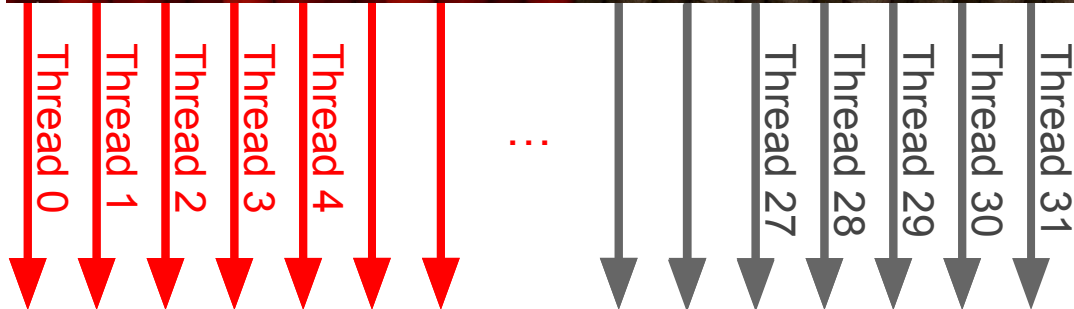
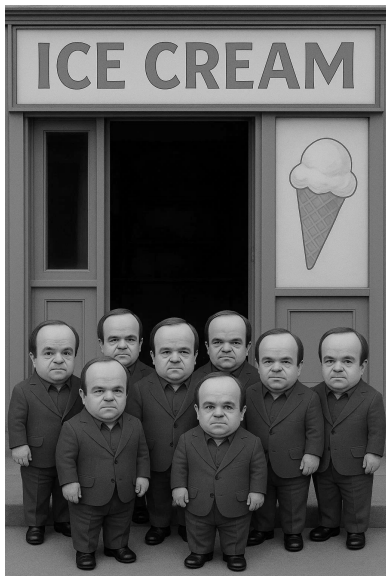
Один Instruction pointer
на все потоки warp-a!

```
239 int warpThreadID = threadIdx.x % 32;  
240 int result = 0;  
241 if (warpThreadID < 16) {  
242     result = dataA[warpThreadID];  
243 } else {  
244     result = dataB[warpThreadID];  
245 }
```

Выходит зайдём в обе ветки if-а?

Архитектура

GPU warp

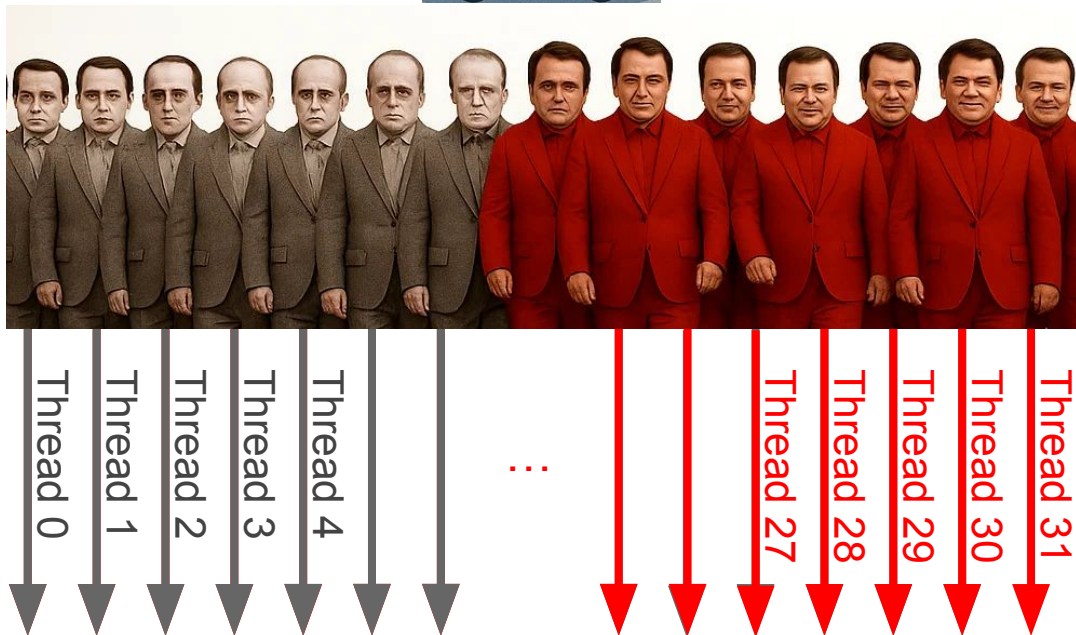


```
239 int warpThreadID = threadIdx.x % 32;
240 int result = 0;
241 if (warpThreadID < 16) {
242     result = dataA[warpThreadID]; // exec mask: [++++ ++++++ ++++++ ++++++ ---- - - - -]
243 } else {
244     result = dataB[warpThreadID]; // exec mask: [- - - - - - - - - - ++++++ ++++++ ++++++ ++++++]
245 }
```

Один Instruction pointer
на все потоки warp-a!

Архитектура

GPU warp

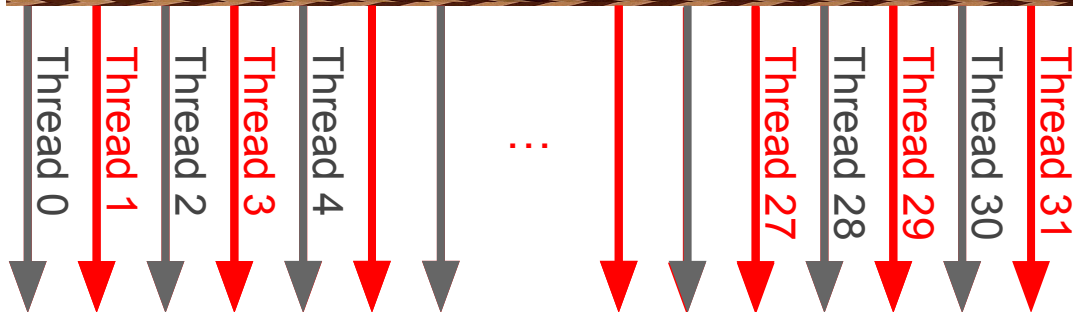


Один Instruction pointer
на все потоки warp-a!

```
239 int warpThreadID = threadIdx.x % 32;  
240 int result = 0;  
241 if (warpThreadID < 16) {  
242     result = dataA[warpThreadID]; // exec mask: [++++ +++++ +++++ +++++ ---- ---- ---- ----]  
243 } else {  
244     result = dataB[warpThreadID]; // exec mask: [---- ---- ---- ---- +++++ +++++ +++++ +++++]  
245 }
```

Архитектура

GPU warp

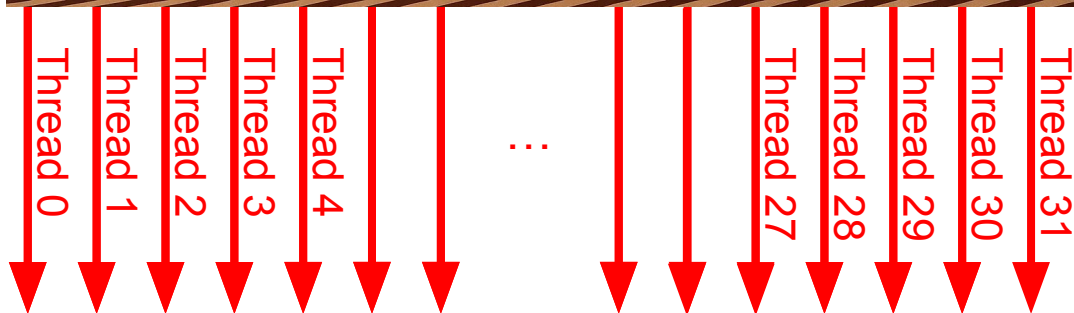


Один Instruction pointer
на все потоки warp-a!

```
239 int warpThreadID = threadIdx.x % 32;  
240 int result = 0;  
241 if (warpThreadID % 2 == 0) {  
242     result = dataA[warpThreadID]; // exec mask: [+--+ +--+ +--+ +--+ +--+ +--+ +--+]  
243 } else {  
244     result = dataB[warpThreadID]; // exec mask: [-+++ -+++ -+++ -+++ -+++ -+++ -+++ -+++]  
245 }
```


Архитектура

GPU warp



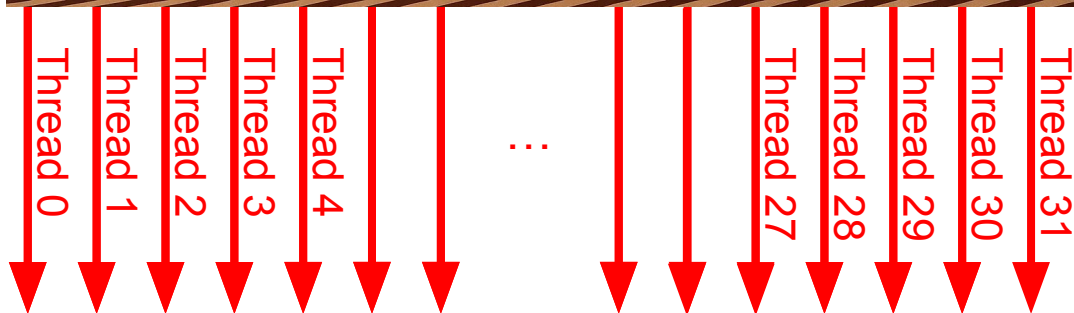
Один Instruction pointer
на все потоки warp-a!

```
241 if (predicate1) {  
242     if (predicate2) {  
243         if (predicate3) {  
244             // A1  
245         } else {  
246             // A2  
247         }  
248     } elif (predicate4) {  
249         // A3  
250     }  
251 } else {  
252     // A4  
253 }
```

Сколько “времени” будет работать warp?

Архитектура

GPU warp



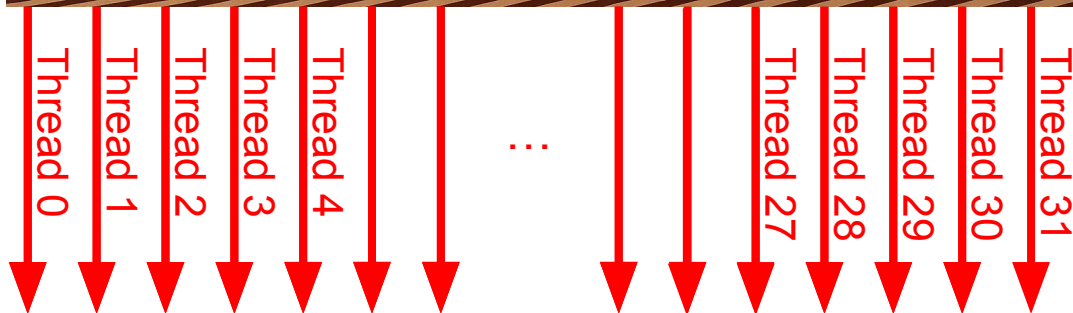
Один Instruction pointer
на все потоки warp-a!

```
241 if (predicate1) {  
242     if (predicate2) {  
243         if (predicate3) {  
244             // A1  
245         } else {  
246             // A2  
247         }  
248     } elif (predicate4) {  
249         // A3  
250     }  
251 } else {  
252     // A4  
253 }
```

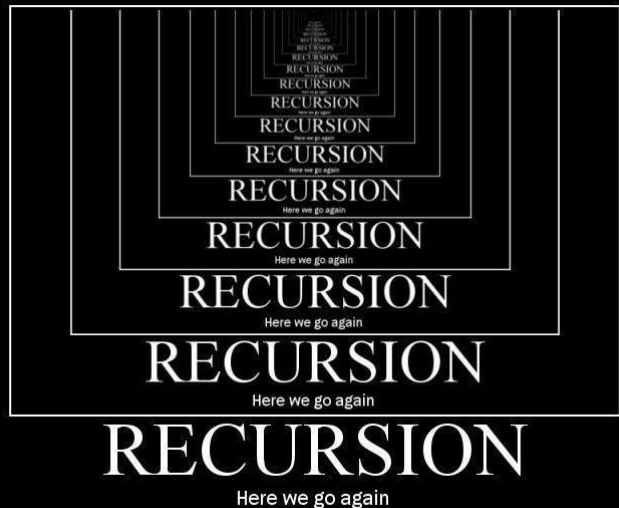
Сколько “времени” будет работать warp?
A1 + A2 + A3 + A4 при **code divergence**

Архитектура

GPU warp



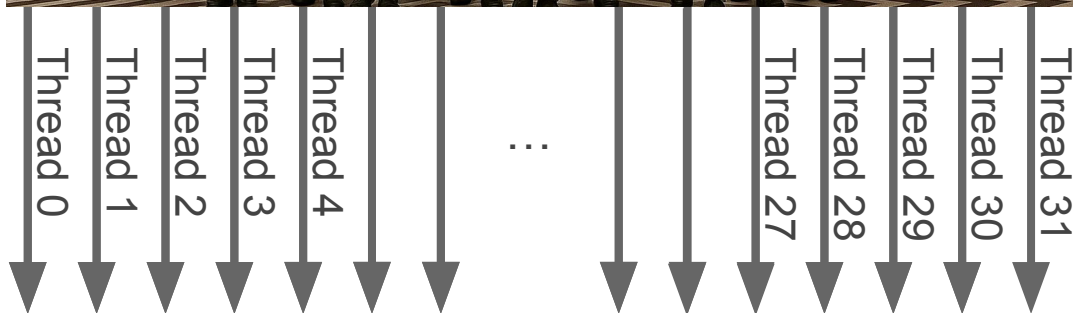
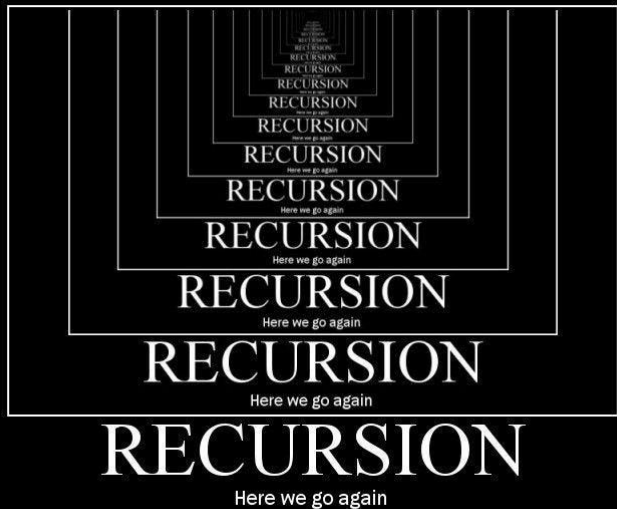
Один Instruction pointer
на все потоки warp-a!



Сколько “времени” будет работать достаточно глубокая рекурсия?

Архитектура

GPU warp

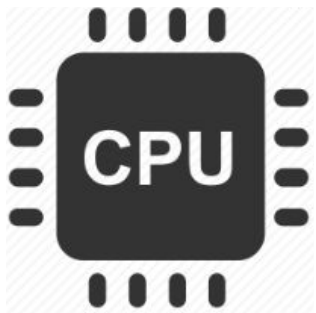


Один Instruction pointer
на все потоки warp-a!

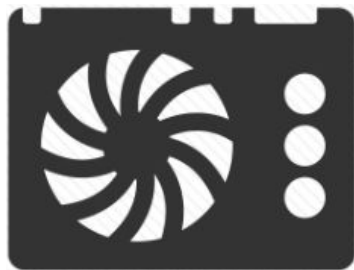
Глава 2: Работа с памятью

hyper-threading, latency hiding, occupancy, registers pressure/spilling,
coalesced memory access pattern, cache lines, local/shared memory

Архитектура



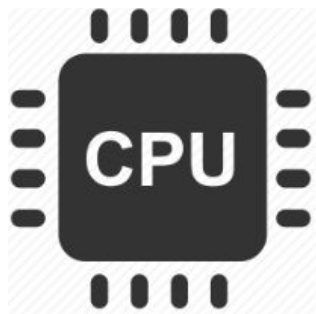
40
GB/s



1000
GB/s

GPU

Архитектура



40
GB/s

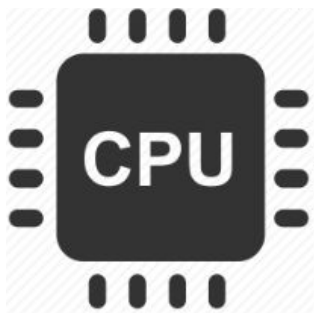
Память - малая пропускная способность
НИЗКАЯ LATENCY



1000
GB/s

GPU

Архитектура



40
GB/s



GPU

1000
GB/s

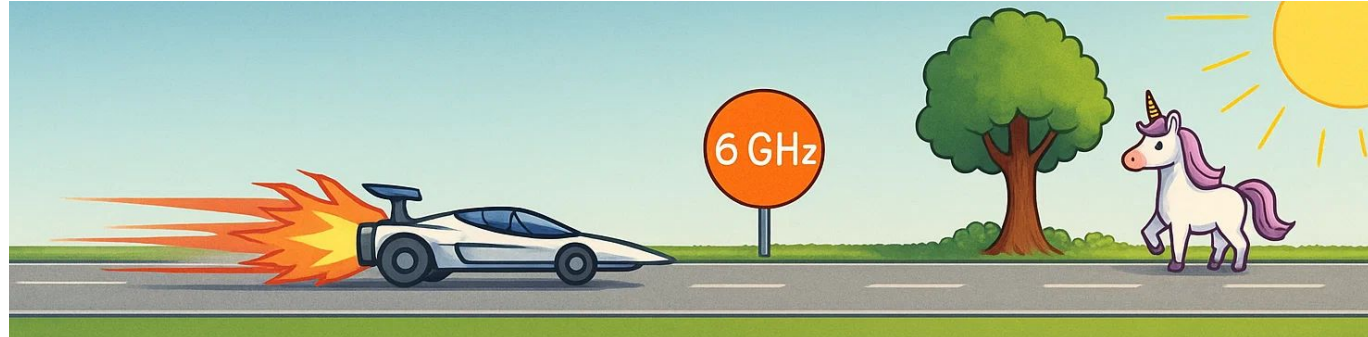
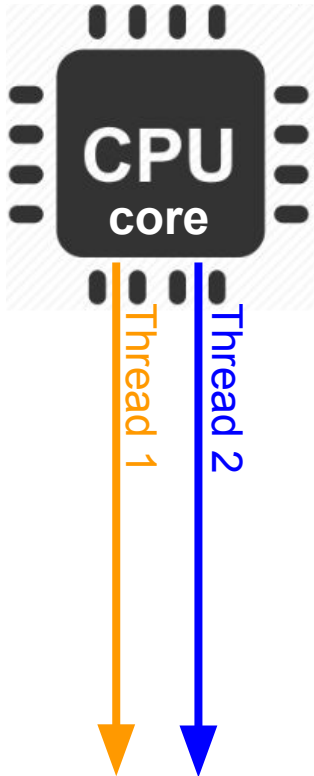
Память - малая пропускная способность
НИЗКАЯ LATENCY



Память - **ОГРОМНАЯ** пропускная способность
но большая latency

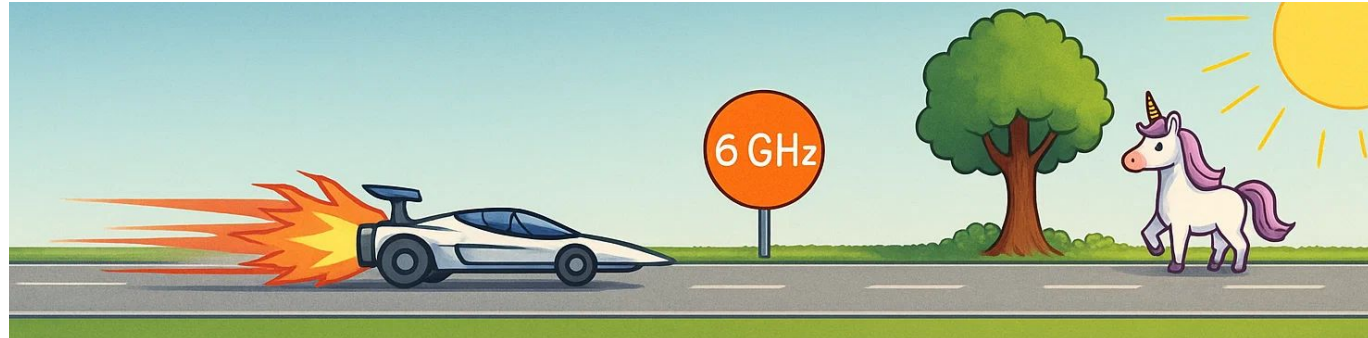
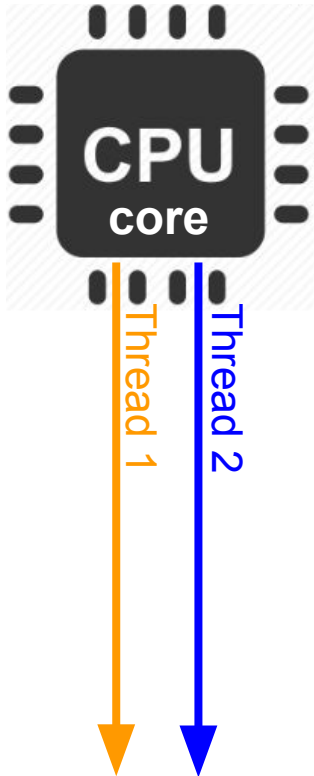


Архитектура CPU: многопоточность



Многопоточность благодаря переключению контекста!

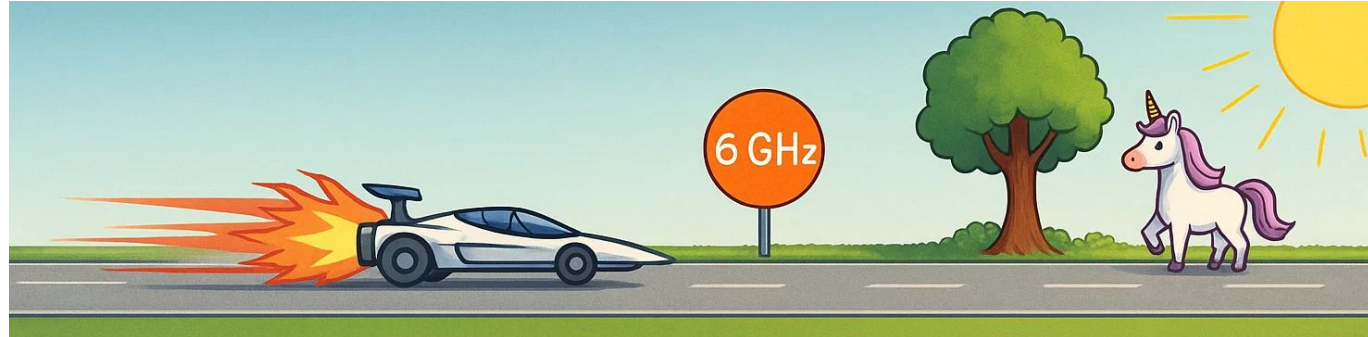
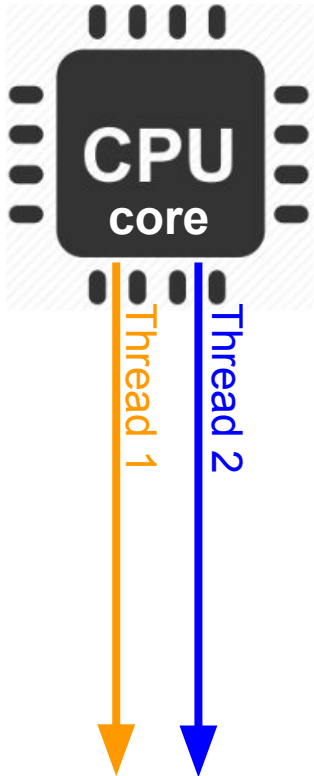
Архитектура CPU: многопоточность



Многопоточность благодаря переключению контекста!

Что нужно перещелкнуть в состоянии ЦПУ ядра для другого потока?

Архитектура CPU: многопоточность

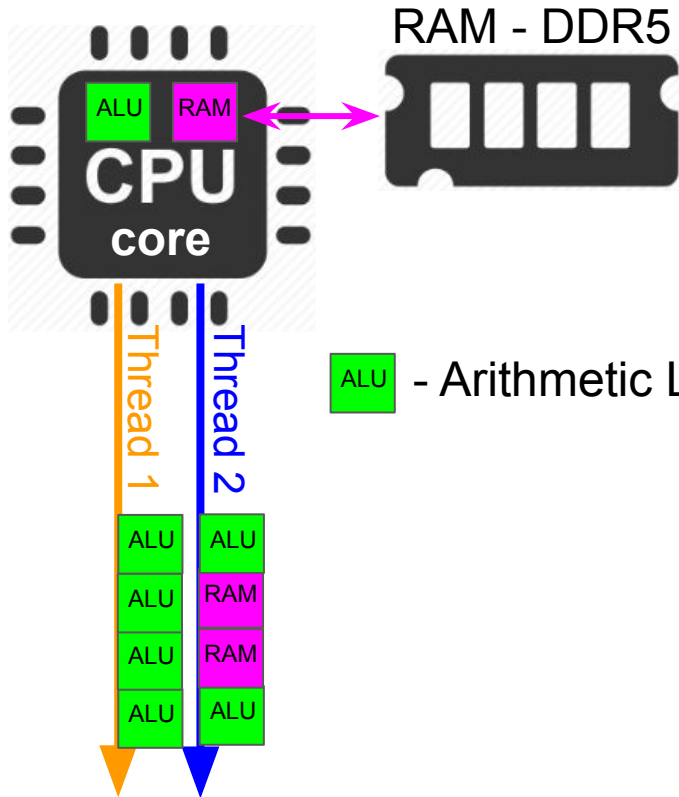


Многопоточность благодаря переключению контекста!

Context switch:

- Обновляет Instruction pointer (указатель на строку кода)
- Подгружает значения регистров процессора (откуда?)

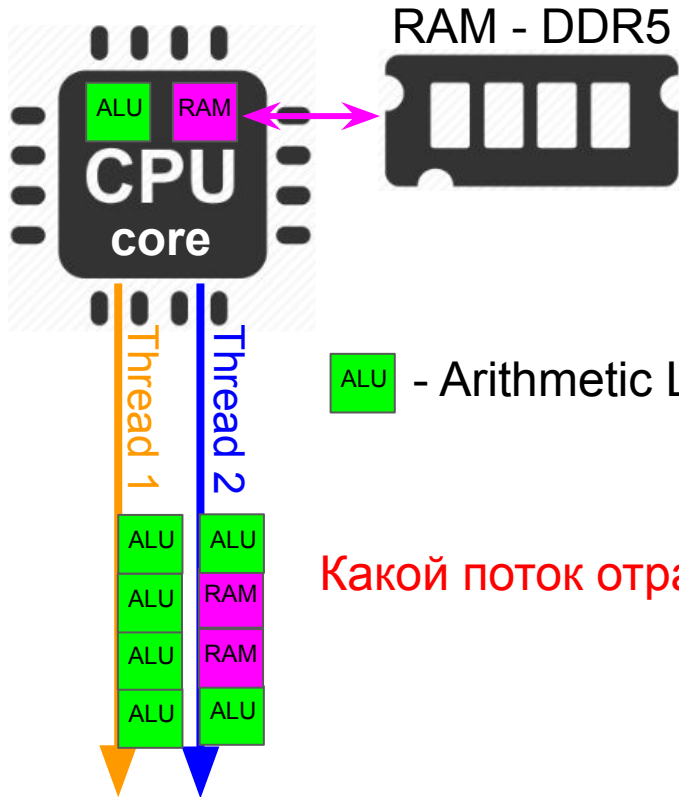
Архитектура CPU: Hyper-Threading, SMT



 - Arithmetic Logical Unit



Архитектура CPU: Hyper-Threading, SMT

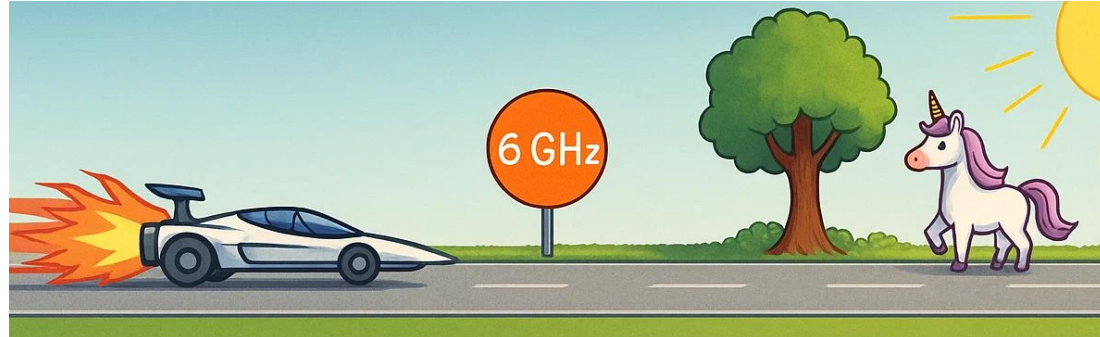
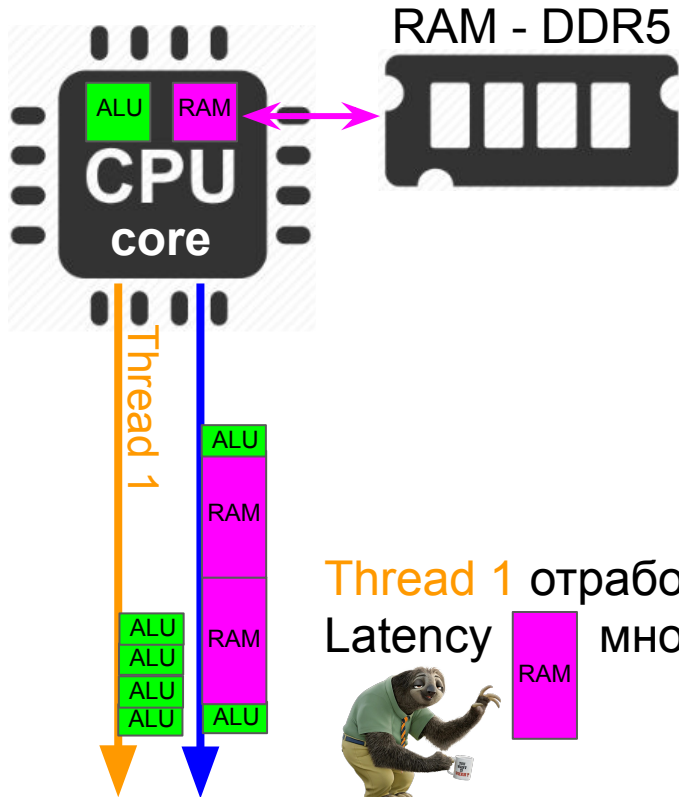


 - Arithmetic Logical Unit

Какой поток отработает быстрее?



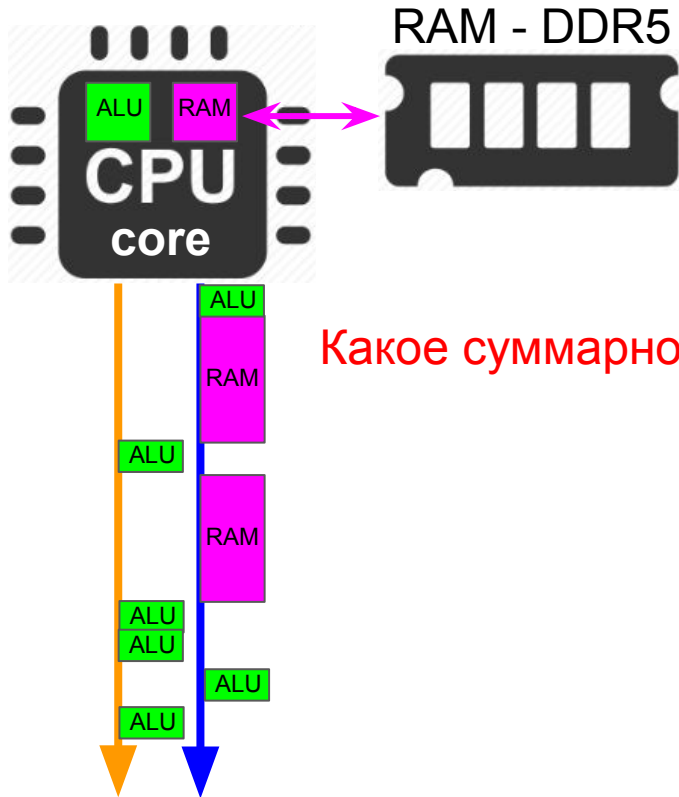
Архитектура CPU: Hyper-Threading, SMT



Thread 1 отработает быстрее!
Latency много больше чем у

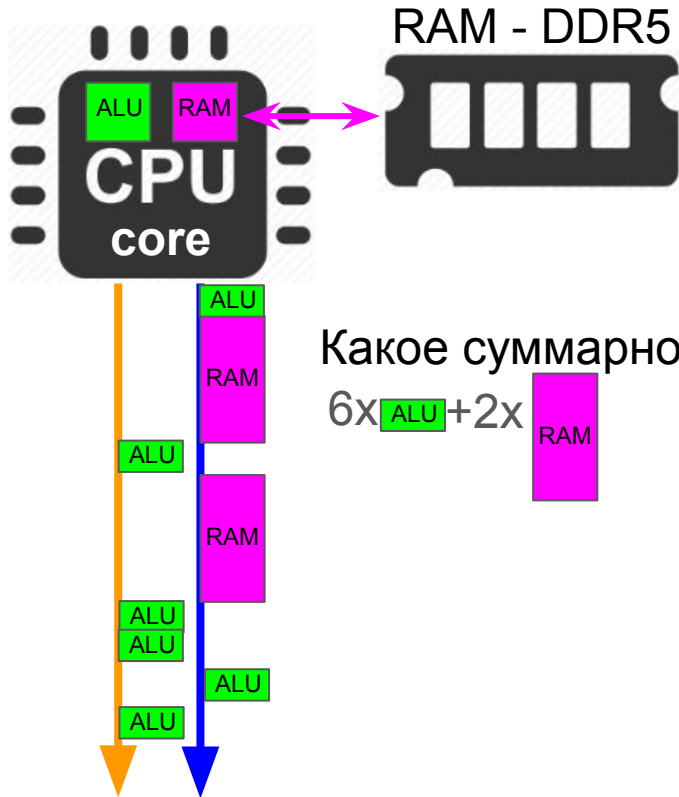


Архитектура CPU: Hyper-Threading, SMT



Какое суммарное время работы двух таких потоков?

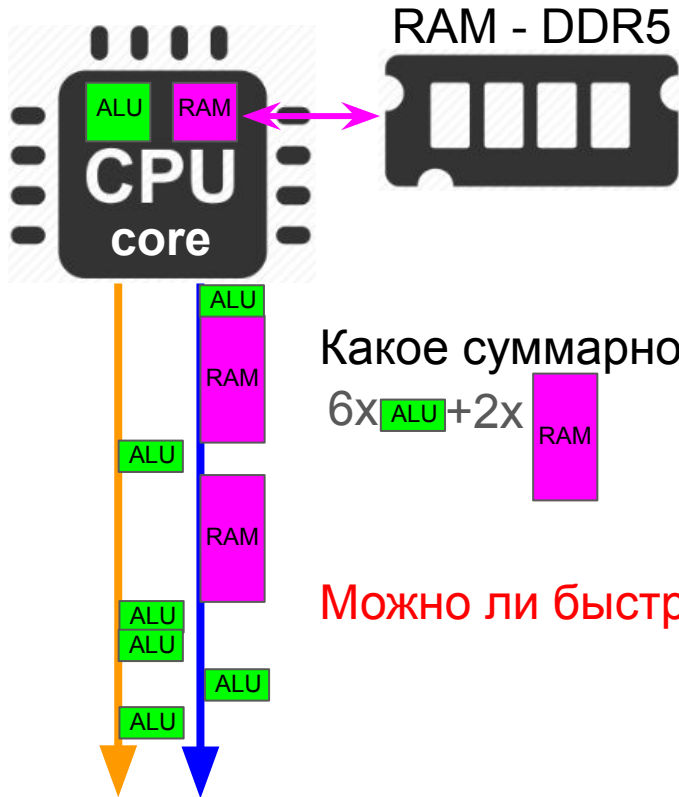
Архитектура CPU: Hyper-Threading, SMT



Какое суммарное время работы двух таких потоков?

$$6 \times \text{ALU} + 2 \times \text{RAM}$$

Архитектура CPU: Hyper-Threading, SMT

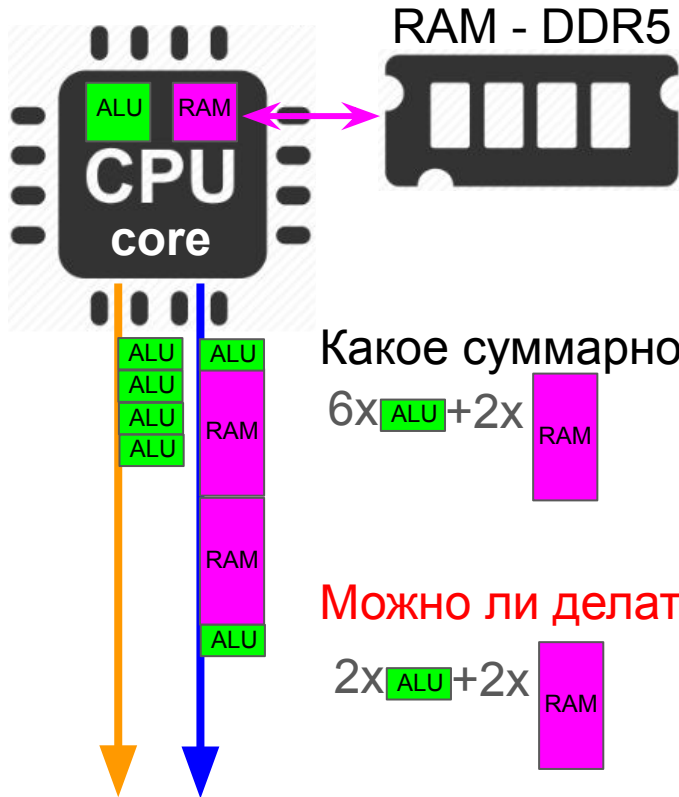


Какое суммарное время работы двух таких потоков?

$$6 \times \text{ALU} + 2 \times \text{RAM}$$

Можно ли быстрее?

Архитектура CPU: Hyper-Threading, SMT



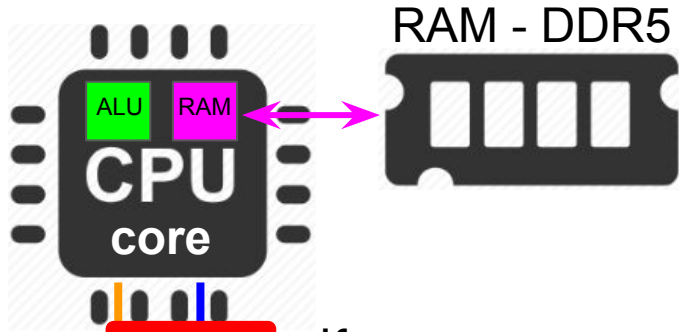
Какое суммарное время работы двух таких потоков?

$$6 \times \text{ALU} + 2 \times \text{RAM}$$

Можно ли делать **ALU** пока ждем **RAM** ?

$$2 \times \text{ALU} + 2 \times \text{RAM}$$

Архитектура CPU: Hyper-Threading, SMT

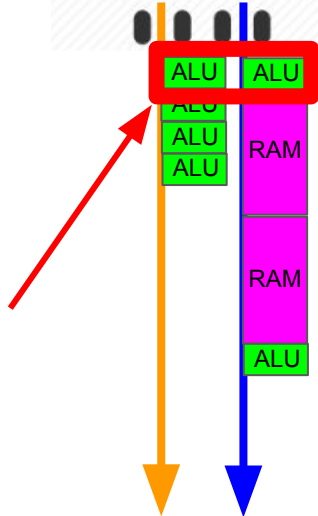


Какое суммарное время работы двух таких потоков?

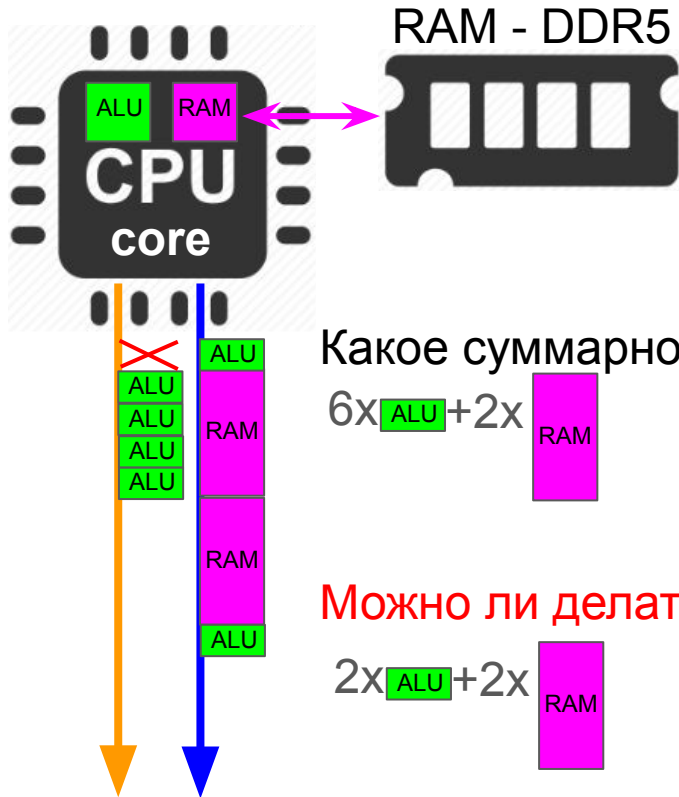
$$6 \times \text{ALU} + 2 \times \text{RAM}$$

Можно ли делать **ALU** пока ждем **RAM** ?

$$2 \times \text{ALU} + 2 \times \text{RAM}$$



Архитектура CPU: Hyper-Threading, SMT



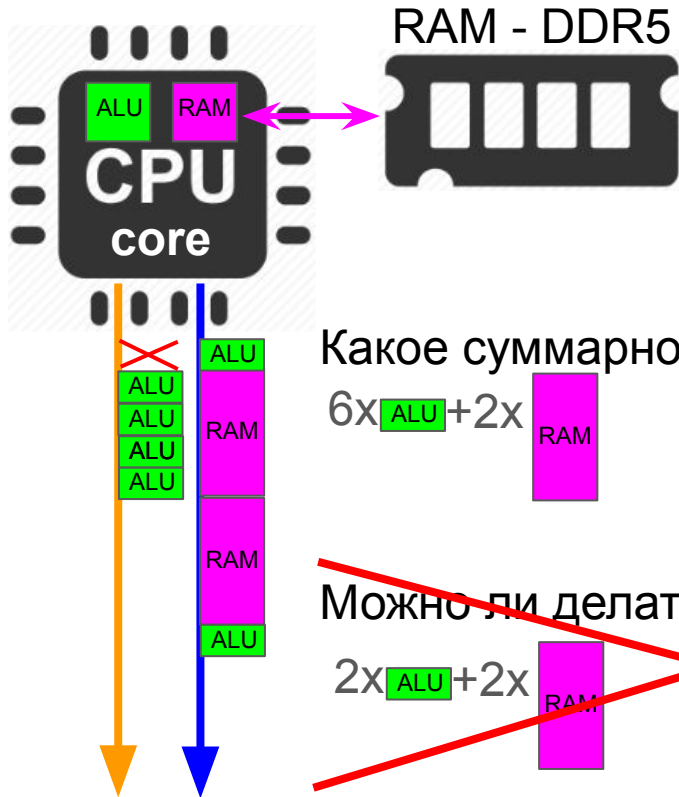
Какое суммарное время работы двух таких потоков?

$$6 \times \text{ALU} + 2 \times \text{RAM}$$

Можно ли делать **ALU** пока ждем **RAM** ?

$$2 \times \text{ALU} + 2 \times \text{RAM}$$

Архитектура CPU: Hyper-Threading, SMT



Какое суммарное время работы двух таких потоков?

$$6 \times \text{ALU} + 2 \times \text{RAM}$$

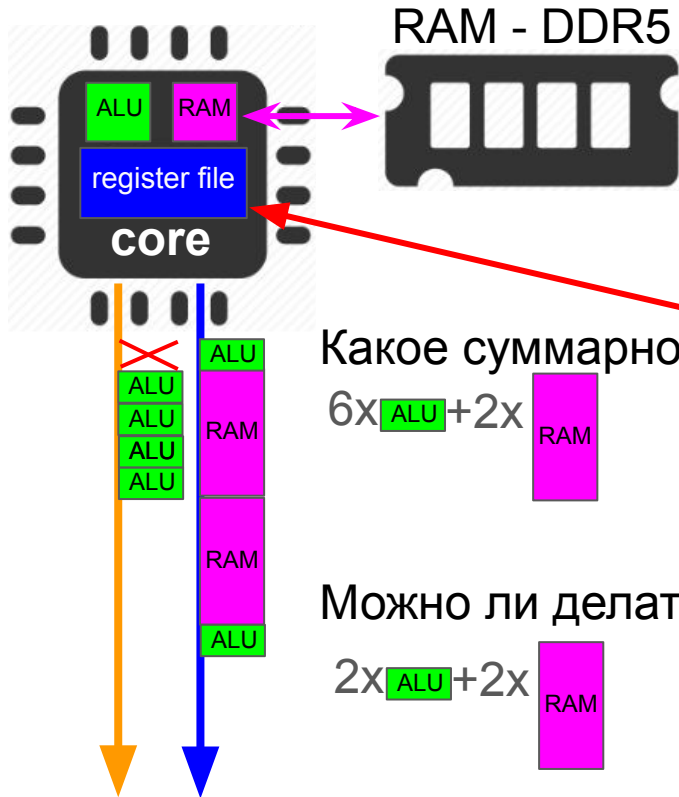
Можно ли делать **ALU** пока ждем **RAM** ?

$$2 \times \text{ALU} + 2 \times \text{RAM}$$

Context switch - дорого!
Долго подгружает регистры!



Архитектура CPU: Hyper-Threading, SMT



Какое суммарное время работы двух таких потоков?

$$6 \times \text{ALU} + 2 \times \text{RAM}$$

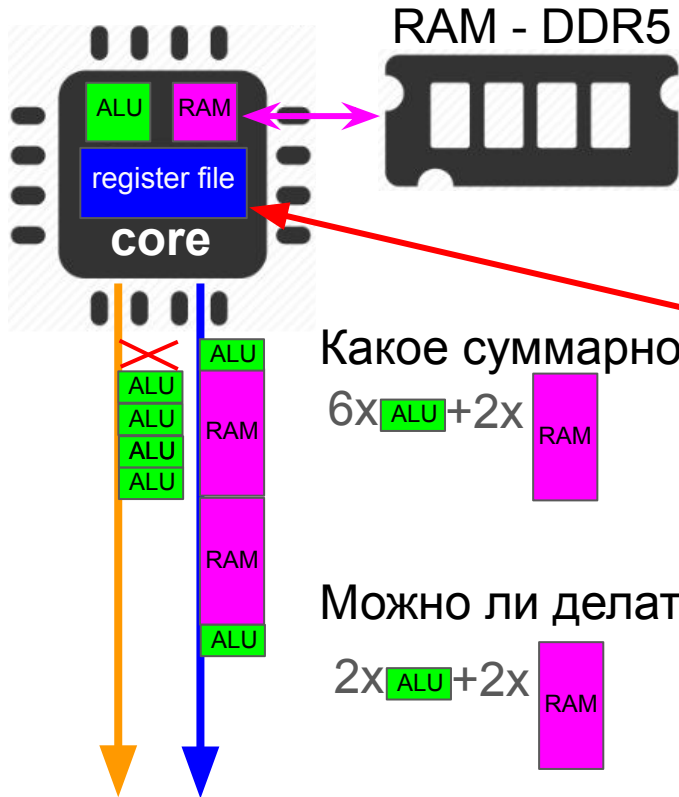
Можно ли делать **ALU** пока ждем **RAM** ?

$$2 \times \text{ALU} + 2 \times \text{RAM}$$

Давайте держать два множества регистров!



Архитектура CPU: Hyper-Threading, SMT



Какое суммарное время работы двух таких потоков?

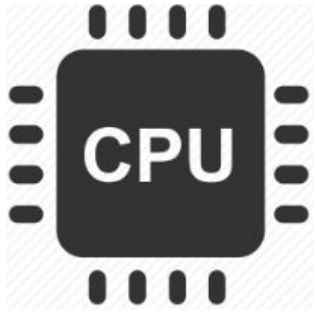
$$6 \times \text{ALU} + 2 \times \text{RAM}$$

Давайте держать два множества регистров!

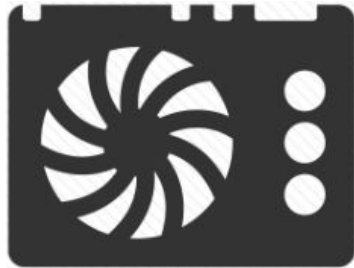
Можно ли делать **ALU** пока ждем **RAM** ? Это и есть **SMT** и **HT**!

$$2 \times \text{ALU} + 2 \times \text{RAM}$$

Архитектура GPU



40
GB/s



1000
GB/s

Память - малая пропускная способность
НИЗКАЯ LATENCY



Память - **ОГРОМНАЯ** пропускная способность
но большая latency



Архитектура GPU



warp

32 слабых медленных CUDA ядер

SM - Streaming Multiprocessor



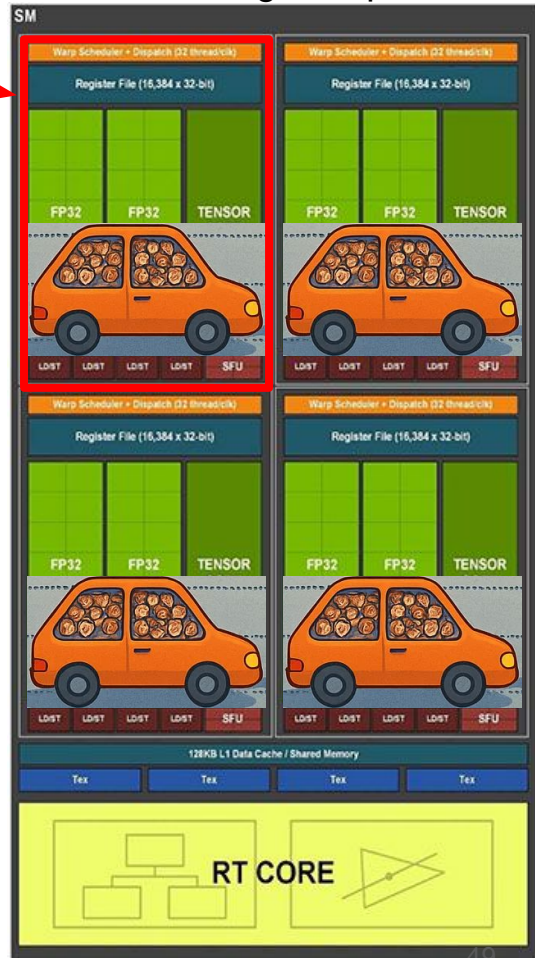
Архитектура GPU



warp

32 слабых медленных
CUDA ядер

SM - Streaming Multiprocessor



Архитектура GPU



warp

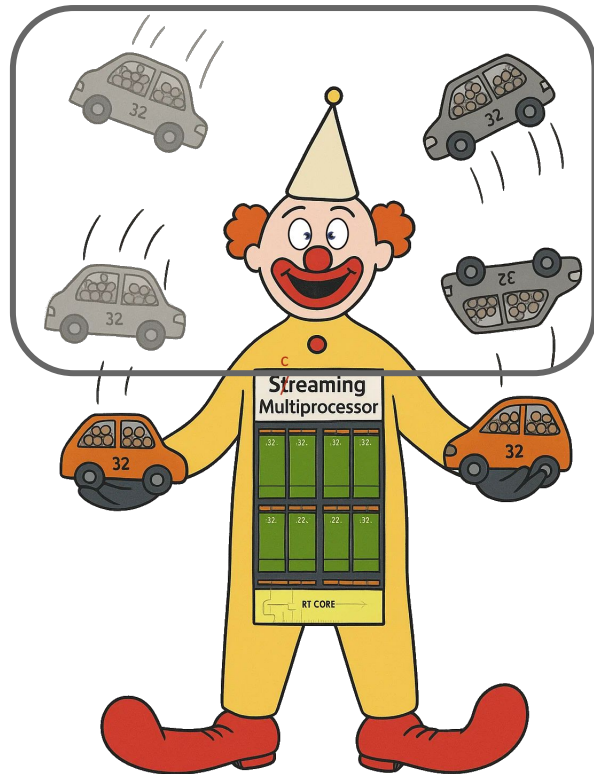
32 слабых медленных
CUDA ядер

Register File для быстрого
жонглирования warp-ами!
(быстрый "context switch")

SM - Streaming Multiprocessor



Архитектура GPU



warp

32 слабых медленных
CUDA ядер

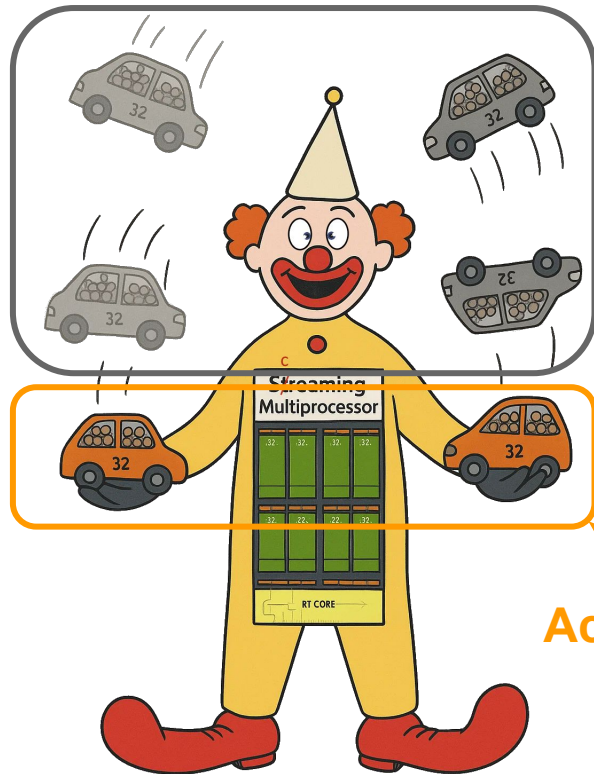
Register File для быстрого
жонглирования warp-ами!
(быстрый "context switch")

Inactive warps: ждут когда
придут данные из VRAM

SM - Streaming Multiprocessor



Архитектура GPU



warp

32 слабых медленных
CUDA ядер

Register File для быстрого
жонглирования warp-ами!
(быстрый "context switch")

Inactive warps: ждут когда
придут данные из VRAM

Active warps: данные в регистрах
Работаем-работаем! 🧐

SM - Streaming Multiprocessor



Архитектура GPU



warp

32 слабых медленных
CUDA ядер

Register File для быстрого
жонглирования warp-ами!
(быстрый "context switch")

Осцирансу = насколько
много доступно warp-ов для
жонглирования.

SM - Streaming Multiprocessor



Архитектура GPU



warp

32 слабых медленных
CUDA ядер

Register File для быстрого
жонглирования warp-ами!
(быстрый "context switch")

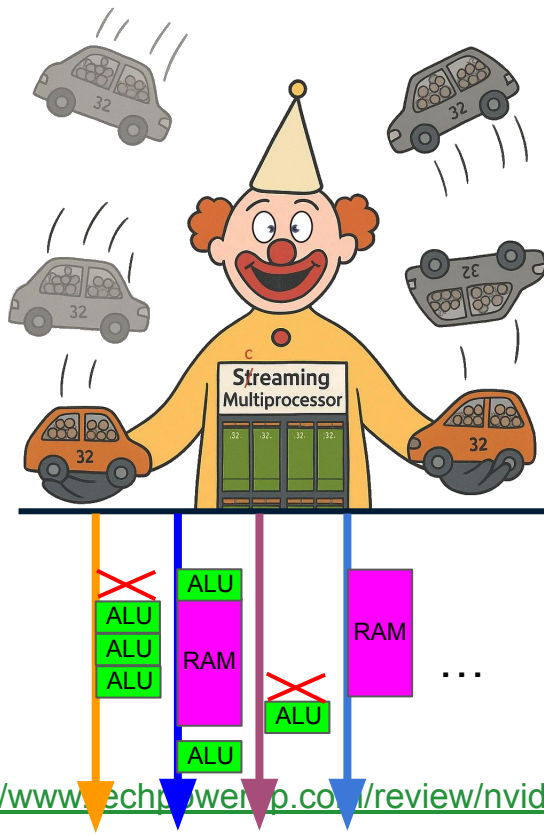
Осцирансу = насколько
много доступно warp-ов для
жонглирования.

**Если осцирансу высокая,
то что?**

SM - Streaming Multiprocessor



Архитектура GPU



warp

32 слабых медленных
CUDA ядер

Register File для быстрого
жонглирования warp-ами!
(быстрый "context switch")

Осцирансу = насколько
много доступно warp-ов для
жонглирования.

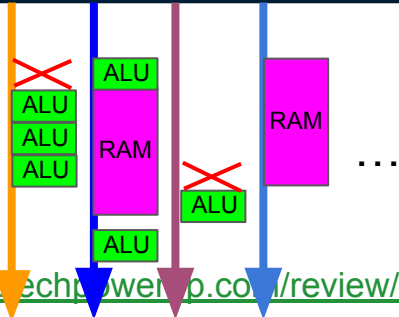
Если осцирансу высокая, то
хорошо скрывается latency
доступа к памяти т.к. всегда
найдется готовый warp!

SM - Streaming Multiprocessor



Архитектура GPU

А почему у SM может быть разное число warp-ов?



warp

32 слабых медленных CUDA ядер

Register File для быстрого жонглирования warp-ами! (быстрый "context switch")

Осцирансу = насколько много доступно warp-ов для жонглирования.

Если осцирансу высокая, то хорошо скрывается latency доступа к памяти т.к. всегда найдется готовый warp!

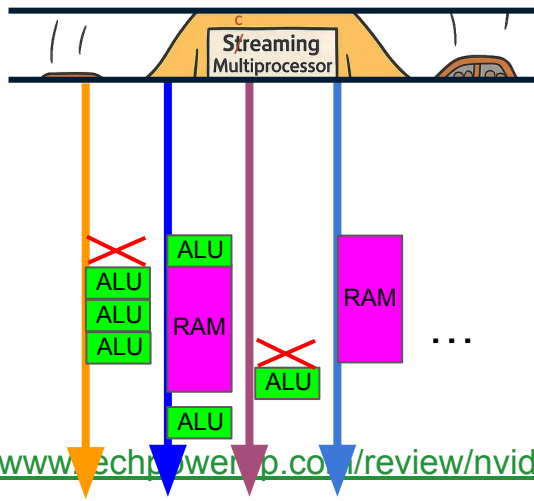
SM - Streaming Multiprocessor



Архитектура GPU

А почему у SM может быть
разное число warp-ов?

Registers Pressure!



32 слабых медленных
CUDA ядер

Register File для быстрого
жонглирования warp-ами!
(быстрый “context switch”)

Осцирансу = насколько
много доступно warp-ов для
жонглирования.

Если осцирансу высокая, то
хорошо скрывается latency
доступа к памяти т.к. всегда
найдется готовый warp!



warp

SM - Streaming Multiprocessor

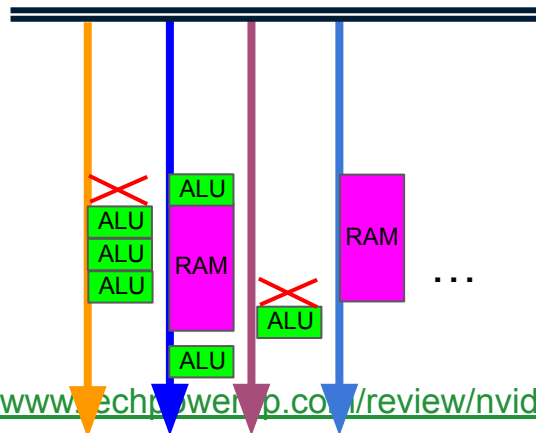


Архитектура GPU

А почему у SM может быть
разное число warp-ов?

Registers Pressure!

Вплоть до **Registers Spilling!**



Register File для быстрого
жонглирования warp-ами!
(быстрый “context switch”)

Осцирансу = насколько
много доступно warp-ов для
жонглирования.

Если осцирансу высокая, то
хорошо скрывается latency
доступа к памяти т.к. всегда
найдется готовый warp!



warp

32 слабых медленных
CUDA ядер

SM - Streaming Multiprocessor



RTX 3090: 10496 CUDA cores = 82 SM · 4 warps · 32 ALUs



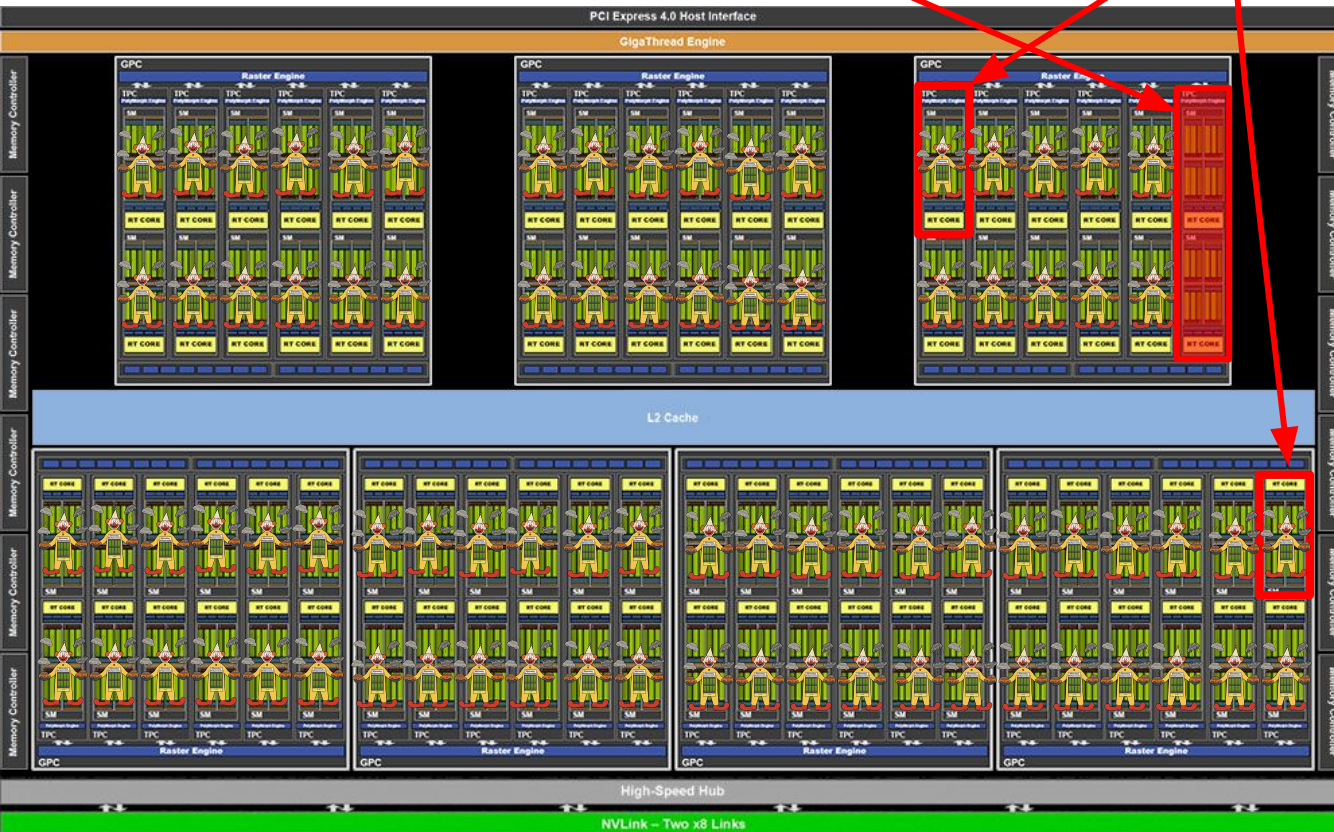
SM - Streaming Multiprocessor



<https://www.techpowerup.com/review/nvidia-geforce-ampere-architecture-board-design-gaming-tech-software/3.html>

RTX 3090: 10496 CUDA cores = 82 SM · 4 warps · 32 ALUs

Куда сбежали два клоуна?



SM - Streaming Multiprocessor



Архитектура VRAM (**coalesced** memory access pattern)



```
float value = data[index];
```

Архитектура VRAM (**coalesced** memory access pattern)

128 bytes = **Cache line** = 32 x float



```
float value = data[index];
```

Архитектура VRAM (**coalesced** memory access pattern)

128 bytes = **Cache line** = 32 x float

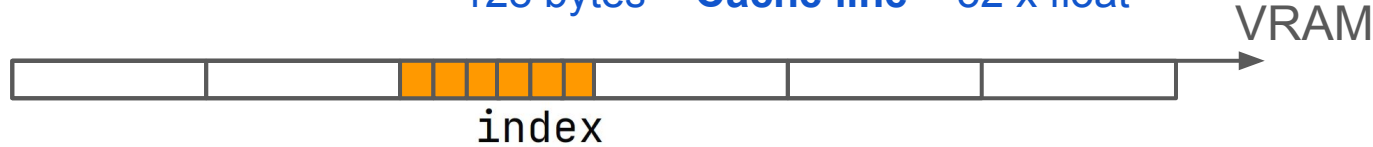


```
float value = data[index];
```



Архитектура VRAM (**coalesced** memory access pattern)

128 bytes = **Cache line** = 32 x float



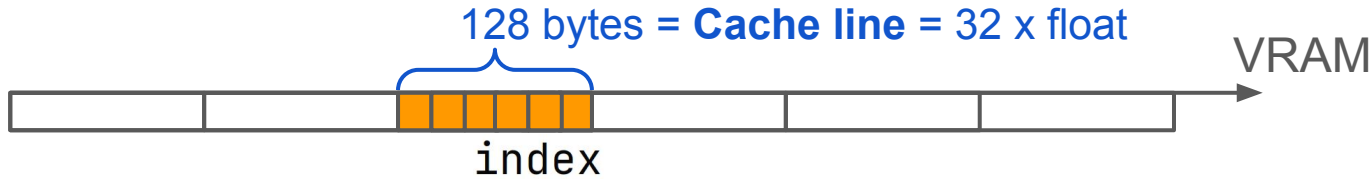
`float value = data[index];` index | cache lines count

[1024+0; 1024+32) ???

Сколько **cache line**-ов чтобы покрыть заказ?
Сколько потребуется **транзакций** из VRAM?



Архитектура VRAM (**coalesced** memory access pattern)



<code>float value = data[index];</code>	index	cache lines count
---	-------	-------------------

<code>[1024+0; 1024+32)</code>	1 транзакция (coalesced)
--------------------------------	--------------------------------------



Архитектура VRAM (**coalesced** memory access pattern)

128 bytes = **Cache line** = 32 x float



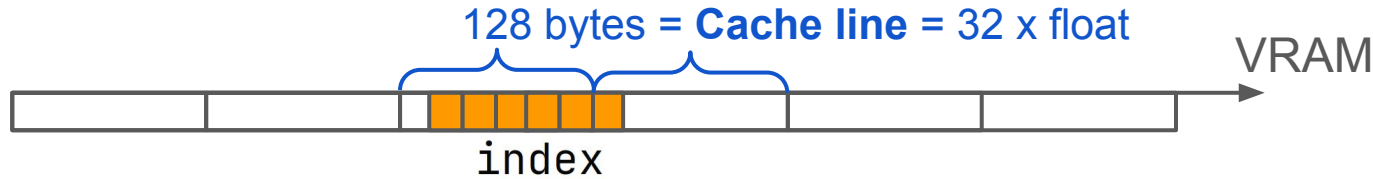
float value = data[index]; index | cache lines count

[1024+0; 1024+32)	1 транзакция (coalesced)
[1024+1; 1024+33)	???

Сколько **cache line**-ов чтобы покрыть заказ?
Сколько потребуется **транзакций** из VRAM?



Архитектура VRAM (**coalesced** memory access pattern)



`float value = data[index];` index | cache lines count

[1024+0; 1024+32)

1 транзакция
(**coalesced**)

[1024+1; 1024+33)

2 транзакции
(**coalesced**)

Насколько просела достигнутая
полезная пропускная способность VRAM?

Архитектура VRAM (**coalesced** memory access pattern)

128 bytes = **Cache line** = 32 x float



float value = data[index]; index | cache lines count



[1024+0; 1024+32)

1 транзакция
(**coalesced**)

[1024+1; 1024+33)

2 транзакции
(**coalesced**)

{1024 + i*32}

???

Сколько **cache line**-ов чтобы покрыть заказ?
Сколько потребуется **транзакций** из VRAM?

Архитектура VRAM (**coalesced** memory access pattern)

128 bytes = **Cache line** = 32 x float



float value = data[index]; index | cache lines count

[1024+0; 1024+32)

1 транзакция
(**coalesced**)

[1024+1; 1024+33)

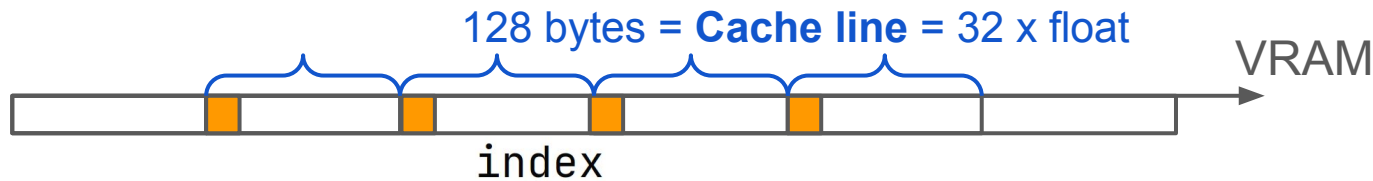
2 транзакции
(**coalesced**)

{1024 + i*32}

???

Сколько **cache line**-ов чтобы покрыть заказ?
Сколько потребуется **транзакций** из VRAM?

Архитектура VRAM (**coalesced** memory access pattern)



`float value = data[index];` index cache lines count

[1024+0; 1024+32)

1 транзакция
(**coalesced**)

[1024+1; 1024+33)

2 транзакции
(**coalesced**)

{1024 + i*32}

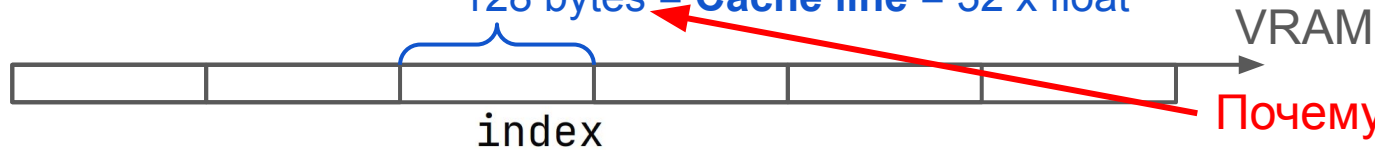
32 транзакции
(**uncoalesced**)



Насколько просела достигнутая
полезная пропускная способность VRAM?⁷⁰

Архитектура VRAM (**coalesced** memory access pattern)

128 bytes = **Cache line** = 32 x float



`float value = data[index];` index cache lines count



[1024+0; 1024+32)

1 транзакция
(**coalesced**)

[1024+1; 1024+33)

2 транзакции
(**coalesced**)

{1024 + i*32}

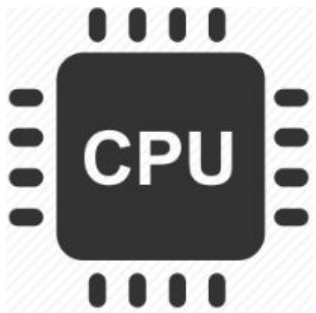
32 транзакции
(**uncoalesced**)

Глава 3: Общая картина ЭВМ-архитектуры

CPU - RAM - PCI-E - VRAM - GPU

Архитектура

$100 \cdot 10^9$ FLOPS

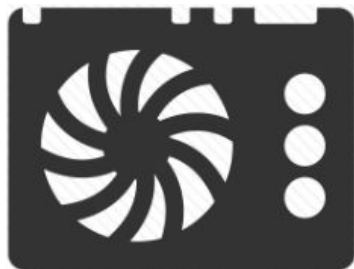


40 GB/s
low latency

RAM - DDR5

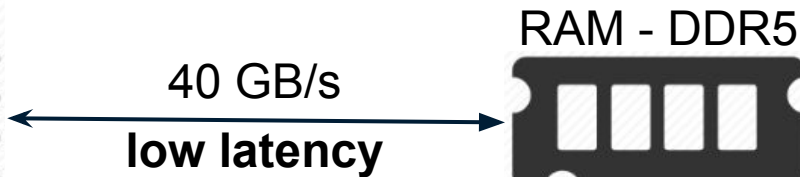
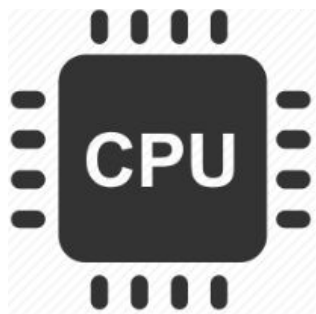


$100 \cdot 10^{12}$ FLOPS



Архитектура

$100 \cdot 10^9$ FLOPS



RAM - DDR5



$100 \cdot 10^{12}$ FLOPS

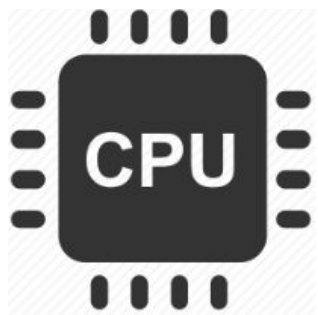


VRAM - GDDR6 или HBM (high bandwidth)

Архитектура

Где лучше исполнять OS?
(например реагировать на клики пользователя)
Какие есть метрики качества?

$100 \cdot 10^9$ FLOPS



40 GB/s
low latency

RAM - DDR5



$100 \cdot 10^{12}$ FLOPS



1000 GB/s
(x25 раз больше)

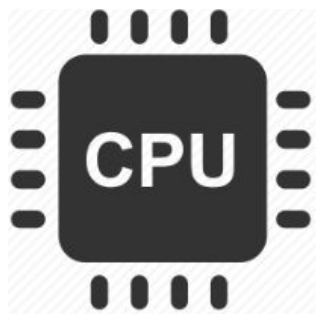


VRAM - GDDR6 или HBM (high bandwidth)

Архитектура

Где быстрее сложить два массива чисел?

$100 \cdot 10^9$ FLOPS



40 GB/s

RAM - DDR5



$100 \cdot 10^{12}$ FLOPS



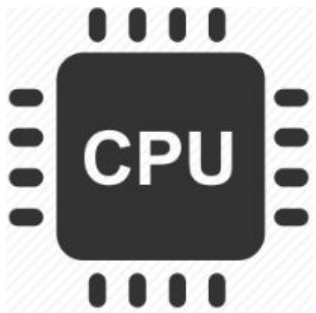
1000 GB/s



VRAM - GDDR6 или HBM (**high bandwidth**)

Архитектура

$100 \cdot 10^9$ FLOPS



Где быстрее сложить два массива чисел?
И во что мы упираемся - в память или в ALUs?
(ALUs = Arithmetic Logic Units)

40 GB/s

RAM - DDR5



$100 \cdot 10^{12}$ FLOPS



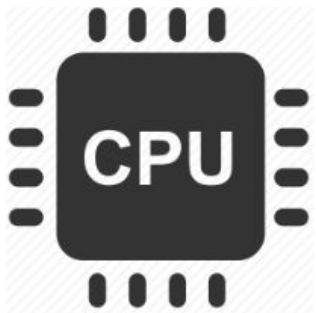
1000 GB/s



VRAM - GDDR6 или HBM (**high bandwidth**)

Архитектура

$100 \cdot 10^9$ FLOPS



40 GB/s

RAM - DDR5



$100 \cdot 10^{12}$ FLOPS



1000 GB/s



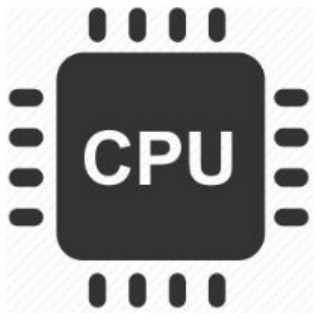
VRAM - GDDR6

PCI-E 5.0 x16
16 GB/s

Где быстрее сложить два массива чисел?
И во что мы упираемся - в память или в ALUs?
А что если данные находятся в **RAM**?

Архитектура

$100 \cdot 10^9$ FLOPS



Где быстрее сложить два массива чисел?
И во что мы упираемся - в память или в ALUs?
А что если данные находятся в **RAM**?

40 GB/s

RAM - DDR5



$100 \cdot 10^{12}$ FLOPS



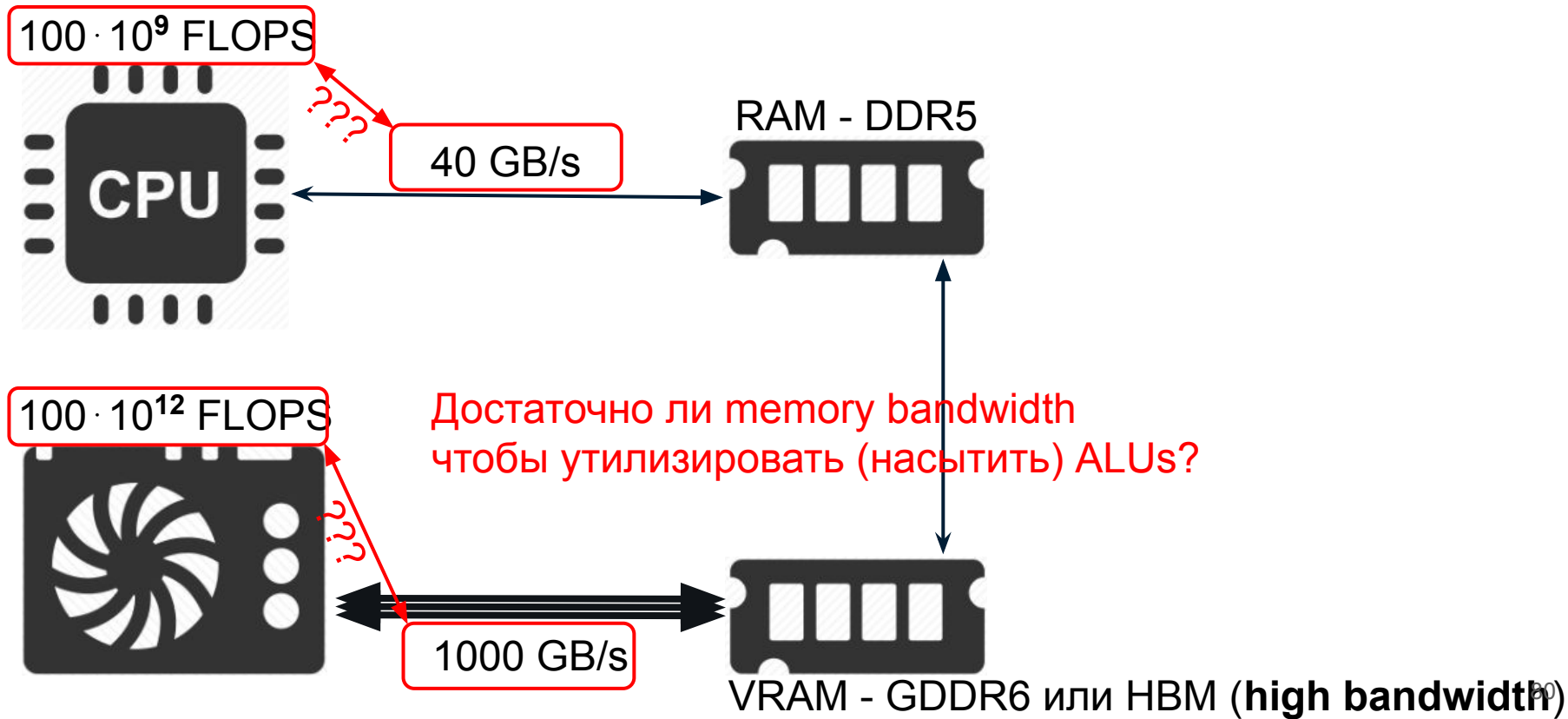
PCI-E 5.0 x16
16 GB/s

1000 GB/s

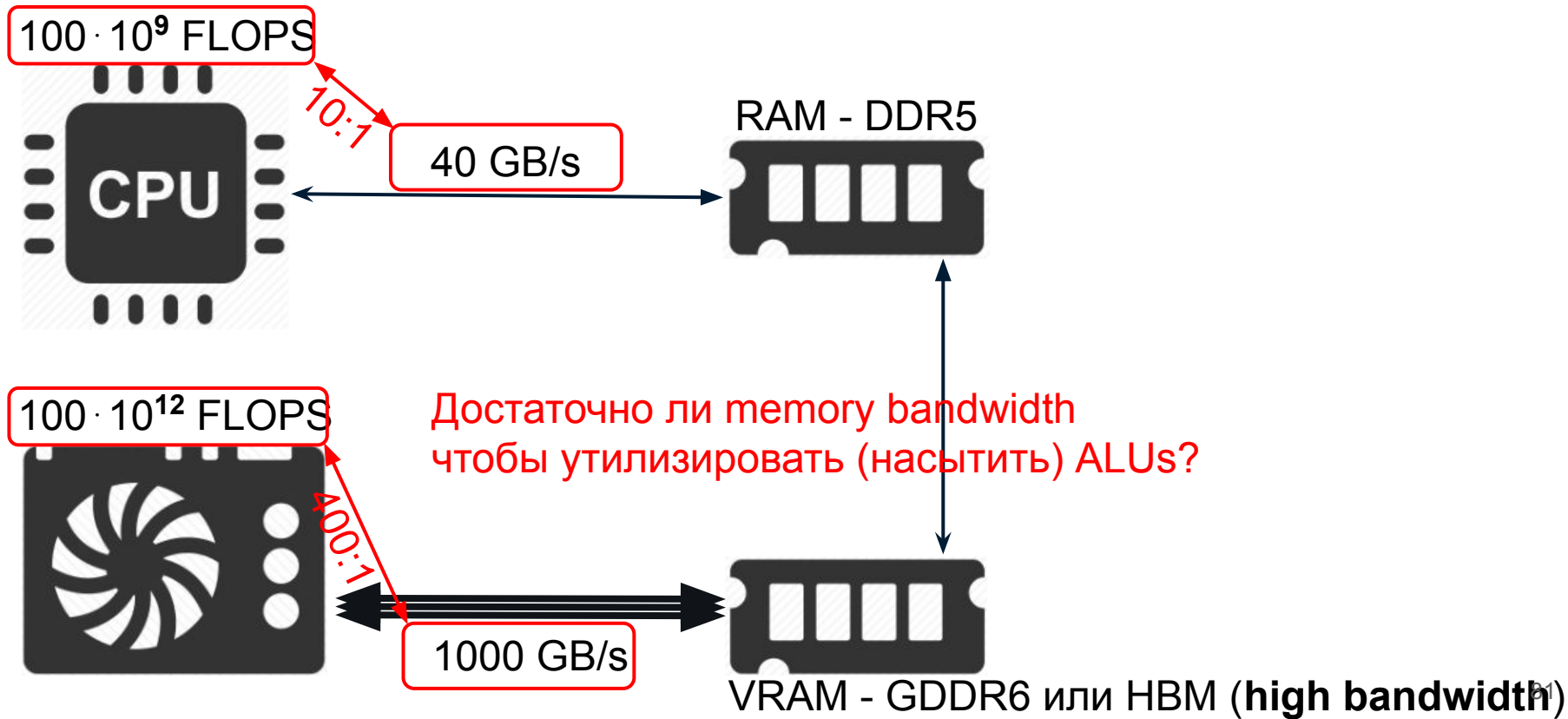


VRAM - GDDR6

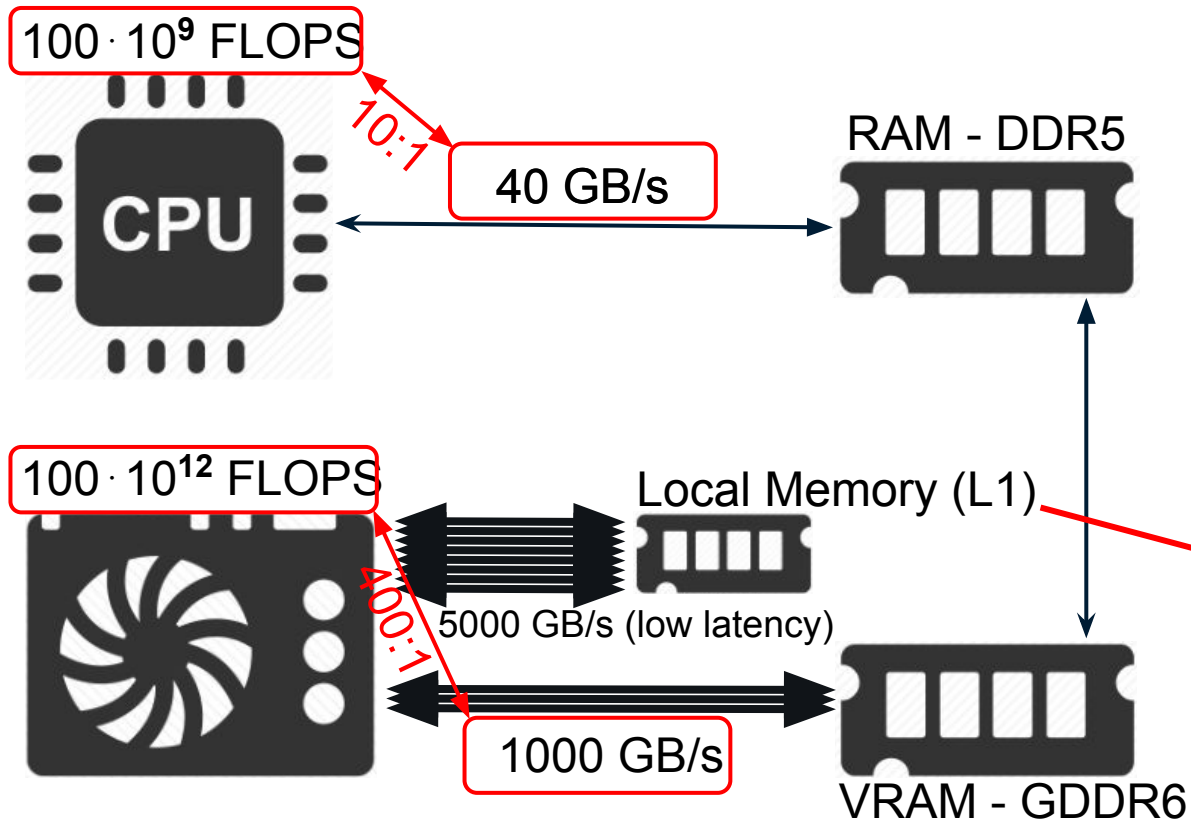
Архитектура



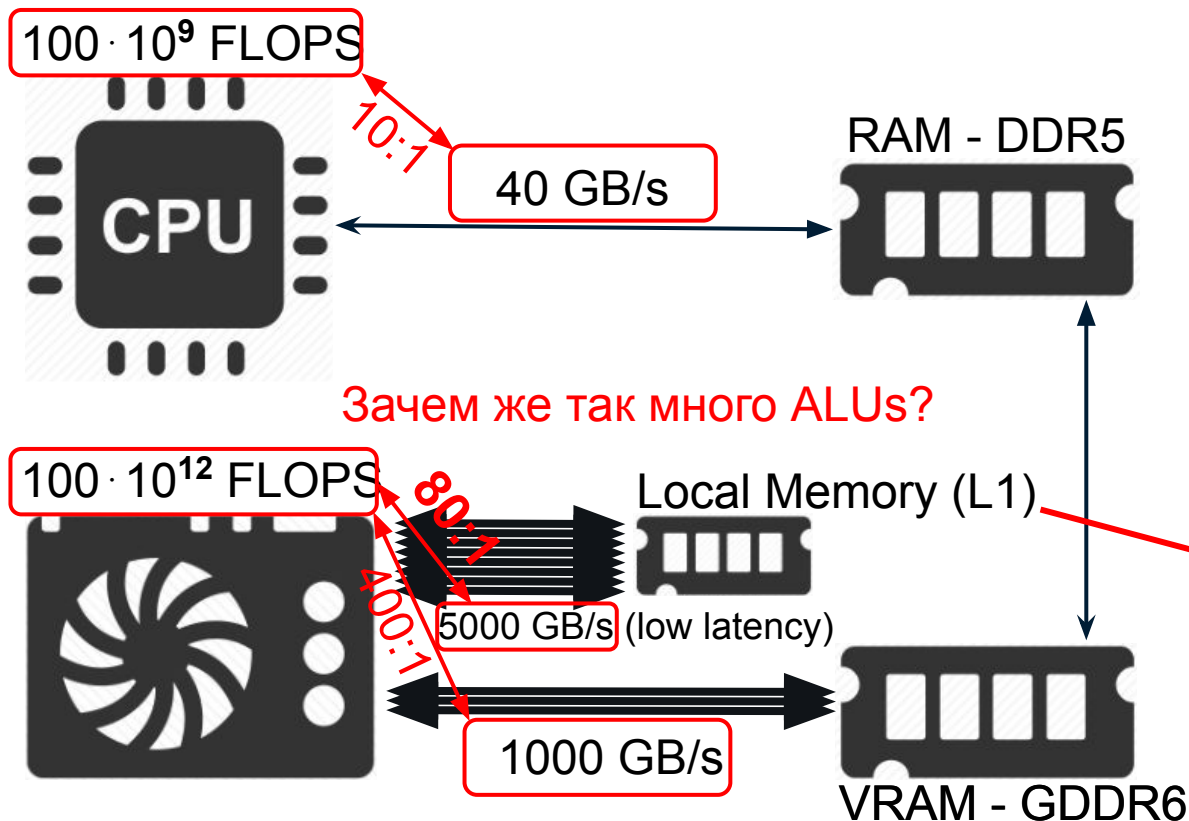
Архитектура



Архитектура

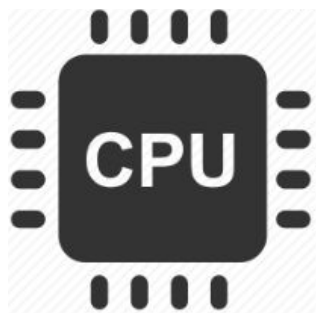


Архитектура



Архитектура

$100 \cdot 10^9$ FLOPS



40 GB/s



RAM - DDR5

Можно ли ее использовать
так же как VRAM?

$100 \cdot 10^{12}$ FLOPS



5000 GB/s (low latency)

Local Memory (L1)



1000 GB/s



VRAM - GDDR6

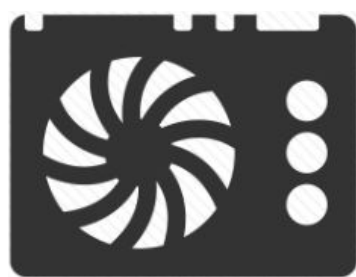




SM - Streaming Multiprocessor



100 · 10¹² FLOPS



Local Memory (L1)



5000 GB/s (low latency)

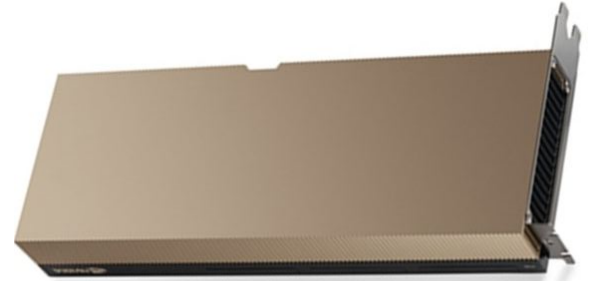
1000 GB/s



VRAM - GDDR6

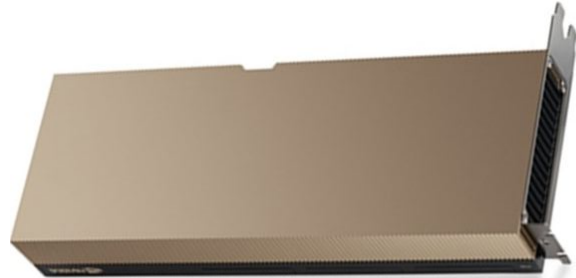
Архитектура GPU

- 1) Десятки тысяч ядер (много **FLOPs**):
 - сгруппированы по 32 в **warp**-ы (по 64 в **wavefront**-ы у AMD)



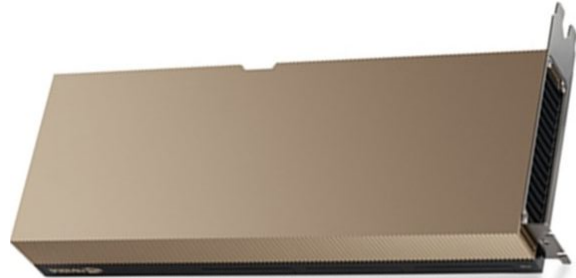
Архитектура GPU

- 1) Десятки тысяч ядер (много **FLOPs**):
 - сгруппированы по 32 в **warp**-ы (по 64 в **wavefront**-ы у AMD)
 - но слабые ядра

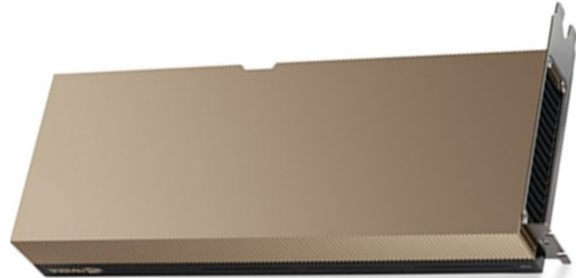


Архитектура GPU

- 1) Десятки тысяч ядер (много **FLOPs**):
 - сгруппированы по 32 в **warp**-ы (по 64 в **wavefront**-ы у AMD)
 - но слабые ядра
 - но у warp-а общий **instruction pointer** (опасность **code divergence**)

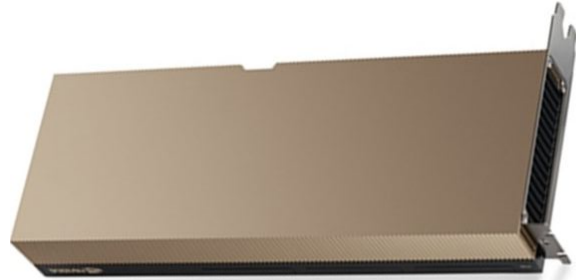


Архитектура GPU



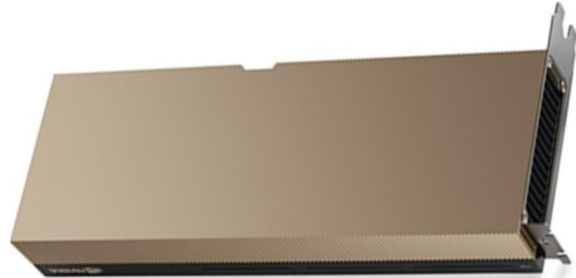
- 1) Десятки тысяч ядер (много **FLOPs**):
 - сгруппированы по 32 в **warp**-ы (по 64 в **wavefront**-ы у AMD)
 - но слабые ядра
 - но у warp-а общий **instruction pointer** (опасность **code divergence**)
 - активен **warp** который не ждет данных из **VRAM** (при высокой **occupancy**)

Архитектура GPU



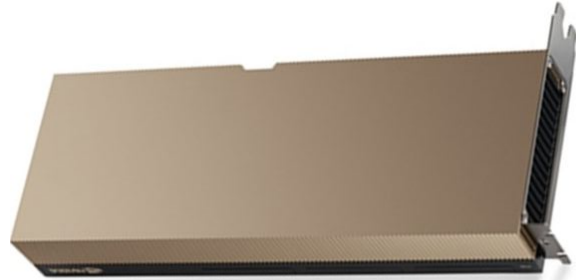
- 1) Десятки тысяч ядер (много **FLOPs**):
 - сгруппированы по 32 в **warp**-ы (по 64 в **wavefront**-ы у AMD)
 - но слабые ядра
 - но у warp-а общий **instruction pointer** (опасность **code divergence**)
 - активен **warp** который не ждет данных из **VRAM** (при высокой **occupancy**)
- 2) Огромная пропускная способность **VRAM**:

Архитектура GPU



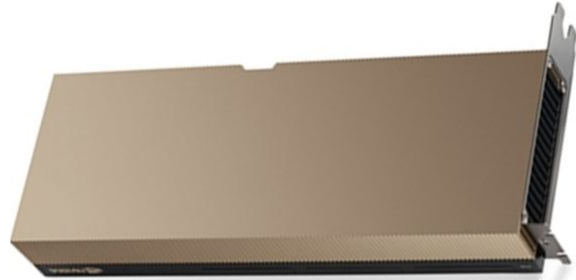
- 1) Десятки тысяч ядер (много **FLOPs**):
 - сгруппированы по 32 в **warp**-ы (по 64 в **wavefront**-ы у AMD)
 - но слабые ядра
 - но у warp-а общий **instruction pointer** (опасность **code divergence**)
 - активен **warp** который не ждет данных из **VRAM** (при высокой **occupancy**)
- 2) Огромная пропускная способность **VRAM**:
 - но большая задержка (**latency**)

Архитектура GPU



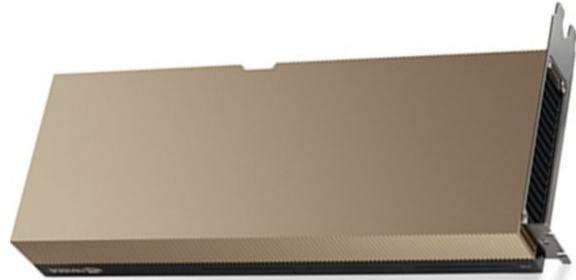
- 1) Десятки тысяч ядер (много **FLOPs**):
 - сгруппированы по 32 в **warp**-ы (по 64 в **wavefront**-ы у AMD)
 - но слабые ядра
 - но у warp-а общий **instruction pointer** (опасность **code divergence**)
 - активен **warp** который не ждет данных из **VRAM** (при высокой **occupancy**)
- 2) Огромная пропускная способность **VRAM**:
 - но большая задержка (**latency**)
 - скрывается за счет **сверх-SMT**/Hyper-Threading (при высокой **occupancy**)

Архитектура GPU



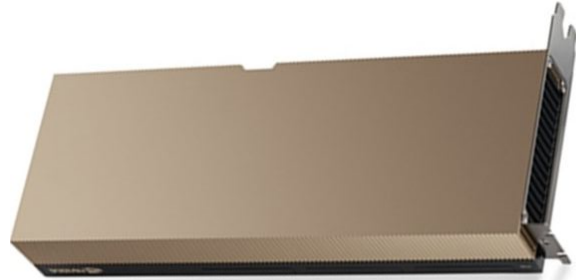
- 1) Десятки тысяч ядер (много **FLOPs**):
 - сгруппированы по 32 в **warp**-ы (по 64 в **wavefront**-ы у AMD)
 - но слабые ядра
 - но у warp-а общий **instruction pointer** (опасность **code divergence**)
 - активен **warp** который не ждет данных из **VRAM** (при высокой **occupancy**)
- 2) Огромная пропускная способность **VRAM**:
 - но большая задержка (**latency**)
 - скрывается за счет **сверх-SMT**/Hyper-Threading (при высокой **occupancy**)
 - транзакции разбиты на 128-байтные **cache line**

Архитектура GPU



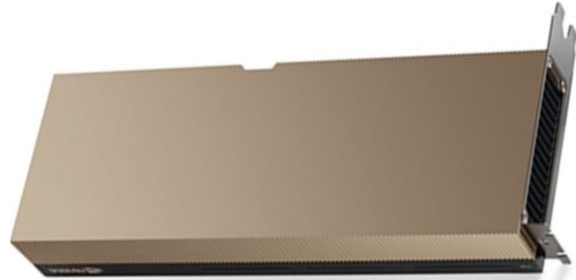
- 1) Десятки тысяч ядер (много **FLOPs**):
 - сгруппированы по 32 в **warp**-ы (по 64 в **wavefront**-ы у AMD)
 - но слабые ядра
 - но у warp-а общий **instruction pointer** (опасность **code divergence**)
 - активен **warp** который не ждет данных из **VRAM** (при высокой **occupancy**)
- 2) Огромная пропускная способность **VRAM**:
 - но большая задержка (**latency**)
 - скрывается за счет **сверх-SMT**/Hyper-Threading (при высокой **occupancy**)
 - транзакции разбиты на 128-байтные **cache line**
 - но при **un-coalesced** memory access pattern - малая проп. способность

Архитектура GPU



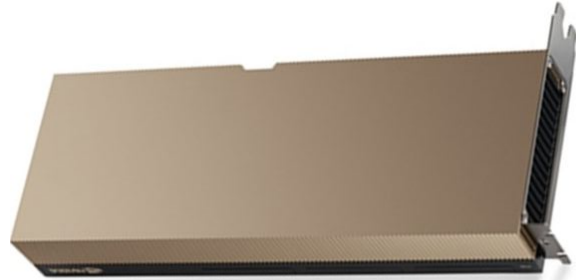
- 1) Десятки тысяч ядер (много **FLOPs**):
 - сгруппированы по 32 в **warp**-ы (по 64 в **wavefront**-ы у AMD)
 - но слабые ядра
 - но у warp-а общий **instruction pointer** (опасность **code divergence**)
 - активен **warp** который не ждет данных из **VRAM** (при высокой **occupancy**)
- 2) Огромная пропускная способность **VRAM**:
 - но большая задержка (**latency**)
 - скрывается за счет **сверх-SMT**/Hyper-Threading (при высокой **occupancy**)
 - транзакции разбиты на 128-байтные **cache line**
 - но при **un-coalesced** memory access pattern - малая проп. способность
- 3) Локальная память (**local/shared memory**):

Архитектура GPU



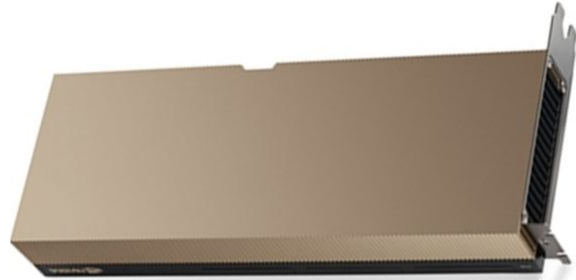
- 1) Десятки тысяч ядер (много **FLOPs**):
 - сгруппированы по 32 в **warp**-ы (по 64 в **wavefront**-ы у AMD)
 - но слабые ядра
 - но у warp-а общий **instruction pointer** (опасность **code divergence**)
 - активен **warp** который не ждет данных из **VRAM** (при высокой **occupancy**)
- 2) Огромная пропускная способность **VRAM**:
 - но большая задержка (**latency**)
 - скрывается за счет **сверх-SMT**/Hyper-Threading (при высокой **occupancy**)
 - транзакции разбиты на 128-байтные **cache line**
 - но при **un-coalesced** memory access pattern - малая проп. способность
- 3) Локальная память (**local/shared memory**):
 - **титаническая** пропускная способность + низкая задержка

Архитектура GPU



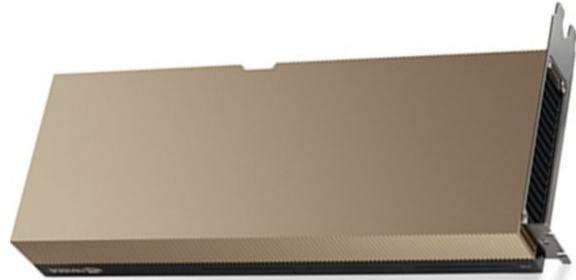
- 1) Десятки тысяч ядер (много **FLOPs**):
 - сгруппированы по 32 в **warp**-ы (по 64 в **wavefront**-ы у AMD)
 - но слабые ядра
 - но у warp-а общий **instruction pointer** (опасность **code divergence**)
 - активен **warp** который не ждет данных из **VRAM** (при высокой **occupancy**)
- 2) Огромная пропускная способность **VRAM**:
 - но большая задержка (**latency**)
 - скрывается за счет **сверх-SMT**/Hyper-Threading (при высокой **occupancy**)
 - транзакции разбиты на 128-байтные **cache line**
 - но при **un-coalesced** memory access pattern - малая проп. способность
- 3) Локальная память (**local/shared memory**):
 - **титаническая** пропускная способность + низкая задержка
 - нет проблемы с coalesced-паттерном

Архитектура GPU



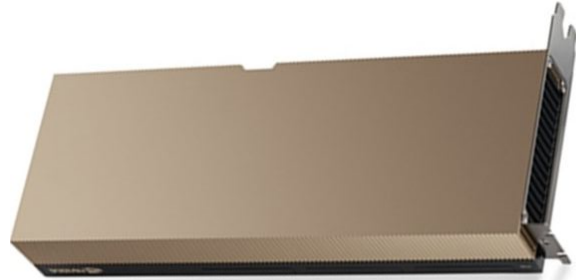
- 1) Десятки тысяч ядер (много **FLOPs**):
 - сгруппированы по 32 в **warp**-ы (по 64 в **wavefront**-ы у AMD)
 - но слабые ядра
 - но у warp-а общий **instruction pointer** (опасность **code divergence**)
 - активен **warp** который не ждет данных из **VRAM** (при высокой **occupancy**)
- 2) Огромная пропускная способность **VRAM**:
 - но большая задержка (**latency**)
 - скрывается за счет **сверх-SMT**/Hyper-Threading (при высокой **occupancy**)
 - транзакции разбиты на 128-байтные **cache line**
 - но при **un-coalesced** memory access pattern - малая проп. способность
- 3) Локальная память (**local/shared memory**):
 - **титаническая** пропускная способность + низкая задержка
 - нет проблемы с coalesced-паттерном
 - **но?**

Архитектура GPU

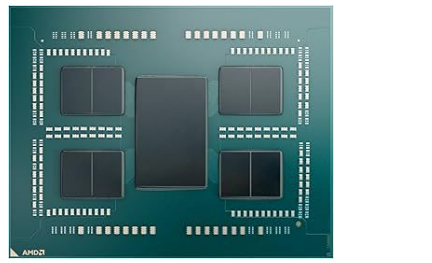


- 1) Десятки тысяч ядер (много **FLOPs**):
 - сгруппированы по 32 в **warp**-ы (по 64 в **wavefront**-ы у AMD)
 - но слабые ядра
 - но у warp-а общий **instruction pointer** (опасность **code divergence**)
 - активен **warp** который не ждет данных из **VRAM** (при высокой **occupancy**)
- 2) Огромная пропускная способность **VRAM**:
 - но большая задержка (**latency**)
 - скрывается за счет **сверх-SMT**/Hyper-Threading (при высокой **occupancy**)
 - транзакции разбиты на 128-байтные **cache line**
 - но при **un-coalesced** memory access pattern - малая проп. способность
- 3) Локальная память (**local/shared memory**):
 - **титаническая** пропускная способность + низкая задержка
 - нет проблемы с coalesced-паттерном
 - но **небольшая** + **локальная** для каждого **warp WorkGroup=Block**

Архитектура GPU



- 1) Десятки тысяч ядер (много **FLOPs**):
 - сгруппированы по 32 в **warp**-ы (по 64 в **wavefront**-ы у AMD)
 - но слабые ядра
 - но у warp-а общий **instruction pointer** (опасность **code divergence**)
 - активен **warp** который не ждет данных из **VRAM** (при высокой **occupancy**)
- 2) Огромная пропускная способность **VRAM**:
 - но большая задержка (**latency**)
 - скрывается за счет **сверх-SMT**/Hyper-Threading (при высокой **occupancy**)
 - транзакции разбиты на 128-байтные **cache line**
 - но при **un-coalesced** memory access pattern - малая проп. способность
- 3) Локальная память (**local/shared memory**):
 - **титаническая** пропускная способность + низкая задержка
 - нет проблемы с coalesced-паттерном
 - но **небольшая** + **локальная** для каждого **warp WorkGroup=Block**
И в этом секрет успеха! Рецепт к масштабируемости!



AMD Ryzen Threadripper 7980X

AMD Ryzen Threadripper 7980X Core-to-Core Latency

CPU0	17.2	18.0	18.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.5	17.
------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	-----

64 ядра = 8 x Чипетов по 8 ядер

AMD Ryzen Threadripper 7980X Core-to-Core Latency

		0-100										100-200										200-300										300-400										400-500										500-600										600-700										700-800										800-900										900-1000										1000-1100										1100-1200										1200-1300										1300-1400										1400-1500										1500-1600										1600-1700										1700-1800										1800-1900										1900-2000										2000-2100										2100-2200										2200-2300										2300-2400										2400-2500										2500-2600										2600-2700										2700-2800										2800-2900										2900-3000										3000-3100										3100-3200										3200-3300										3300-3400										3400-3500										3500-3600										3600-3700										3700-3800										3800-3900										3900-4000										4000-4100										4100-4200										4200-4300										4300-4400										4400-4500										4500-4600										4600-4700										4700-4800										4800-4900										4900-5000										5000-5100										5100-5200										5200-5300										5300-5400										5400-5500										5500-5600										5600-5700										5700-5800										5800-5900										5900-6000										6000-6100										6100-6200										6200-6300										6300-6400										6400-6500										6500-6600										6600-6700										6700-6800										6800-6900										6900-7000										7000-7100										7100-7200										7200-7300										7300-7400										7400-7500										7500-7600										7600-7700										7700-7800										7800-7900										7900-8000										8000-8100										8100-8200										8200-8300										8300-8400										8400-8500										8500-8600										8600-8700										8700-8800										8800-8900										8900-9000										9000-9100										9100-9200										9200-9300										9300-9400										9400-9500										9500-9600										9600-9700										9700-9800										9800-9900										9900-10000									
0-100	100-200	200-300	300-400	400-500	500-600	600-700	700-800	800-900	900-1000	1000-1100	1100-1200	1200-1300	1300-1400	1400-1500	1500-1600	1600-1700	1700-1800	1800-1900	1900-2000	2000-2100	2100-2200	2200-2300	2300-2400	2400-2500	2500-2600	2600-2700	2700-2800	2800-2900	2900-3000	3000-3100	3100-3200	3200-3300	3300-3400	3400-3500	3500-3600	3600-3700	3700-3800	3800-3900	3900-4000	4000-4100	4100-4200	4200-4300	4300-4400	4400-4500	4500-4600	4600-4700	4700-4800	4800-4900	4900-5000	5000-5100	5100-5200	5200-5300	5300-5400	5400-5500	5500-5600	5600-5700	5700-5800	5800-5900	5900-6000	6000-6100	6100-6200	6200-6300	6300-6400	6400-6500	6500-6600	6600-6700	6700-6800	6800-6900	6900-7000	7000-7100	7100-7200	7200-7300	7300-7400	7400-7500	7500-7600	7600-7700	7700-7800	7800-7900	7900-8000	8000-8100	8100-8200	8200-8300	8300-8400	8400-8500	8500-8600	8600-8700	8700-8800	8800-8900	8900-9000	9000-9100	9100-9200	9200-9300	9300-9400	9400-9500	9500-9600	9600-9700	9700-9800	9800-9900	9900-10000																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																						
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00																																																																																																																																																																																																																																																																																																																																																																																																																																													

64 ядра = 8 x Чиплетов (по 8 ядер в чиплете)

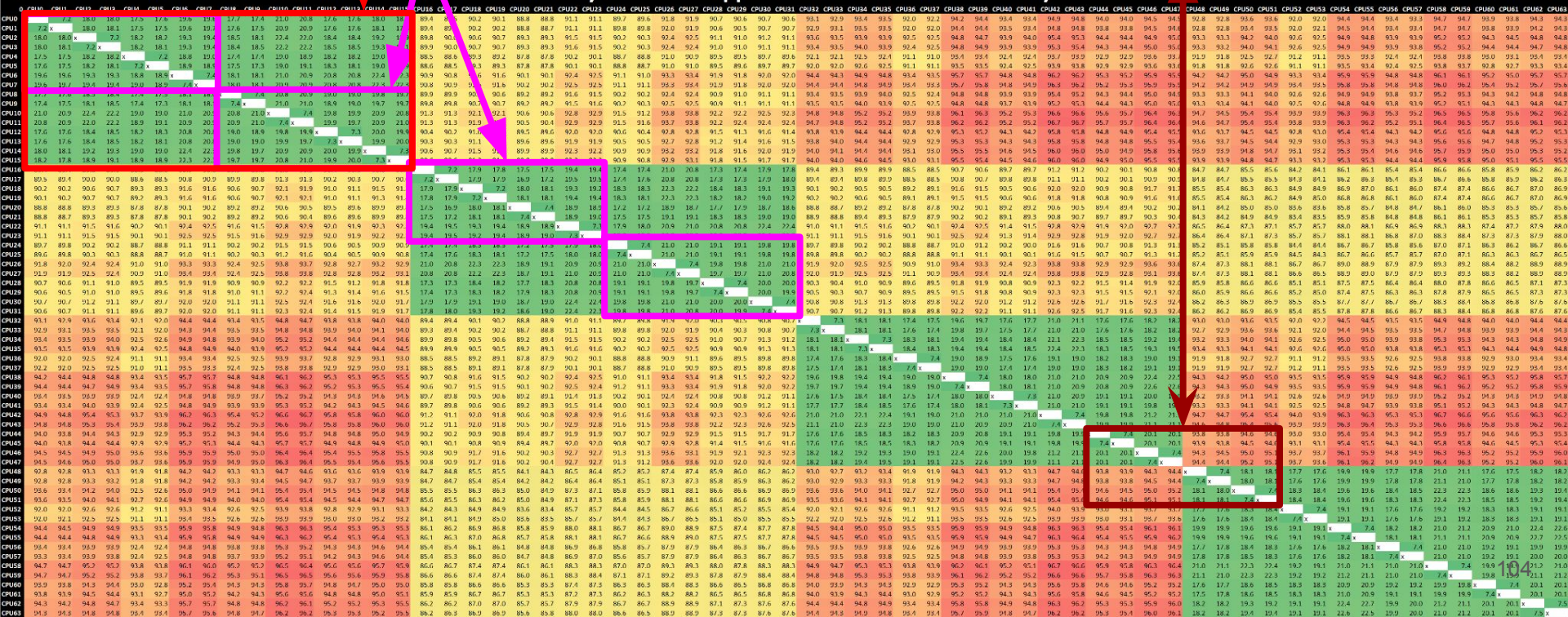
AMD Ryzen Threadripper 7980X Core-to-Core Latency

CPU0	180	180	175	175	175	177	174	170	168	176	176	160	18	CPU8	902	902	901	888	888	911	911	897	896	918	919	907	906	907	906	913	929	934	935	920	922	942	944	934	934	939	948	940	940	945	945	928	928	936	936	920	920	944	944	944	944	933	947	947	939	938	943	943		
CPU1	180	180	175	175	175	177	174	170	168	176	176	160	18	CPU9	902	902	901	888	888	911	911	897	896	918	919	907	906	907	906	913	929	934	935	920	922	942	944	934	934	939	948	940	940	945	945	928	928	936	936	920	920	944	944	944	944	933	947	947	939	938	943	943		
CPU2	180	181	182	182	182	181	183	183	184	185	182	182	184	182	1	CPU10	902	902	901	888	888	911	911	897	896	918	919	907	906	907	906	913	929	934	935	920	922	942	944	934	934	939	948	940	940	945	945	928	928	936	936	920	920	944	944	944	944	933	947	947	939	938	943	943
CPU3	180	181	182	182	182	181	183	183	184	185	182	182	184	182	1	CPU11	902	902	901	888	888	911	911	897	896	918	919	907	906	907	906	913	929	934	935	920	922	942	944	934	934	939	948	940	940	945	945	928	928	936	936	920	920	944	944	944	944	933	947	947	939	938	943	943
CPU4	175	175	182	182	182	182	183	183	184	185	182	182	184	182	1	CPU12	902	902	901	888	888	911	911	897	896	918	919	907	906	907	906	913	929	934	935	920	922	942	944	934	934	939	948	940	940	945	945	928	928	936	936	920	920	944	944	944	944	933	947	947	939	938	943	943
CPU5	176	176	183	183	183	183	184	184	185	182	182	184	182	1	CPU13	902	902	901	888	888	911	911	897	896	918	919	907	906	907	906	913	929	934	935	920	922	942	944	934	934	939	948	940	940	945	945	928	928	936	936	920	920	944	944	944	944	933	947	947	939	938	943	943	
CPU6	176	176	183	183	183	183	184	184	185	182	182	184	182	1	CPU14	902	902	901	888	888	911	911	897	896	918	919	907	906	907	906	913	929	934	935	920	922	942	944	934	934	939	948	940	940	945	945	928	928	936	936	920	920	944	944	944	944	933	947	947	939	938	943	943	
CPU7	176	176	183	183	183	183	184	184	185	182	182	184	182	1	CPU15	902	902	901	888	888	911	911	897	896	918	919	907	906	907	906	913	929	934	935	920	922	942	944	934	934	939	948	940	940	945	945	928	928	936	936	920	920	944	944	944	944	933	947	947	939	938	943	943	
CPU8	176	176	183	183	183	183	184	184	185	182	182	184	182	1	CPU16	902	902	901	888	888	911	911	897	896	918	919	907	906	907	906	913	929	934	935	920	922	942	944	934	934	939	948	940	940	945	945	928	928	936	936	920	920	944	944	944	944	933	947	947	939	938	943	943	
CPU9	176	176	183	183	183	183	184	184	185	182	182	184	182	1	CPU17	902	902	901	888	888	911	911	897	896	918	919	907	906	907	906	913	929	934	935	920	922	942	944	934	934	939	948	940	940	945	945	928	928	936	936	920	920	944	944	944	944	933	947	947	939	938	943	943	
CPU10	210	209	224	222	220	191	200	210	208	218	7	198	219	7	CPU18	902	902	901	888	888	911	911	897	896	918	919	907	906	907	906	913	929	934	935	920	922	942	944	934	934	939	948	940	940	945	945	928	928	936	936	920	920	944	944	944	944	933	947	947	939	938	943	943	
CPU11	208	209	224	222	220	191	200	210	208	218	7	198	219	7	CPU19	902	902	901	888	888	911	911	897	896	918	919	907	906	907	906	913	929	934	935	920	922	942	944	934	934	939	948	940	940	945	945	928	928	936	936	920	920	944	944	944	944	933	947	947	939	938	943	943	
CPU12	208	209	224	222	220	191	200	210	208	218	7	198	219	7	CPU20	902	902	901	888	888	911	911	897	896	918	919	907	906	907	906	913	929	934	935	920	922	942	944	934	934	939	948	940	940	945	945	928	928	936	936	920	920	944	944	944	944	933	947	947	939	938	943	943	
CPU13	176	176	183	183	183	183	184	184	185	182	182	184	182	1	CPU21	902	902	901	888	888	911	911	897	896	918	919	907	906	907	906	913	929	934	935	920	922	942	944	934	934	939	948	940	940	945	945	928	928	936	936	920	920	944	944	944	944	933	947	947	939	938	943	943	
CPU14	180	181	182	183	183	183	184	184	185	182	182	184	182	1	CPU22	902	902	901	888	888	911	911	897	896	918	919	907	906	907	906	913	929	934	935	920	922	942	944	934	934	939	948	940	940	945	945	928	928	936	936	920	920	944	944	944	944	933	947	947	939	938	943	943	
CPU15	180	181	182	183	183	183	184	184	185	182	182	184	182	1	CPU23	902	902	901	888	888	911	911	897	896	918	919	907	906	907	906	913	929	934	935	920	922	942	944	934	934	939	948	940	940	945	945	928	928	936	936	920	920	944	944	944	944	933	947	947	939	938	943	943	
CPU16	180	181	182	183	183	183	184	184	185	182	182	184	182	1	CPU24	902	902	901	888	888	911	911	897	896	918	919	907	906	907	906	913	929	934	935	920	922	942	944	934	934	939	948	940	940	945	945	928	928	936	936	920	920	944	944	944	944	933	947	947	939	938	943	943	
CPU17	180	181	182	183	183	183	184	184	185	182	182	184	182	1	CPU25	902	902	901	888	888	911	911	897	896	918	919	907	906	907	906	913	929	934	935	920	922	942	944	934	934	939	948	940	940	945	945	928	928	936	936	920	920	944	944	944	944	933	947	947	939	938	943	943	
CPU18	180	181	182	183	183	183	184	184	185	182	182	184	182	1	CPU26	902	902	901	888	888	911	911	897	896	918	919	907	906	907	906	913	929	934	935	920	922	942	944	934	934	939	948	940	940	945	945	928	928	936	936	920	920	944	944	944	944	933	947	947	939	938	943	943	
CPU19	180	181	182	183	183	183	184	184	185	182	182	184	182	1	CPU27	902	902	901	888	888	911	911	897	896	918	919	907	906	907	906	913	929	934	935	920	922	942	944	934	934	939	948	940	940	945	945	928	928	936	936	920	920	944	944	944	944	933	947	947	939	938	943	943	
CPU20	180	181	182	183	183	183	184	184	185	182	182	184	182	1	CPU28	902	902	901	888	888	911	911	897	896	918	919	907	906	907	906	913	929	934	935	920	922	942	944	934	934	939	948	940	940	945	945	928	928	936	936	920	920	944	944	944	944	933	947	947	939	938	943	943	
CPU21	180	181	182	183	183	183	184	184	185	182	182	184	182	1	CPU29	902	902	901	888	888	911	911	897	896	918	919	907	906	907	906	913	929	934	935	920	922	942	944	934	934	939	948	940	940	945	945	928	928	936	936	920	920	944	944	944	944	933	947	947	939	938	943	943	
CPU22	180	181	182	183	183	183	184	184	185	182	182	184	182	1	CPU30	902	902	901	888	888	911	911	897	896	918	919	907	906	907	906	913	929	934	935	920	922	942	944	934	934	939	948	940	940	945	945	928	928	936	936	920	920	944	944	944	944	933	947	947	939	938	943	943	
CPU23	180	181	182	183	183	183	184	184	185	182	182	184	182	1	CPU31	902	902	901	888	888	911	911	897	896	918	919	907	906	907	906	913	929	934	935	920	922	942	944	934	934	939	948	940	940	945	945	928	928	936	936	920	920	944	944	944	944	933	947	947	939	938	943	943	
CPU24	180	181	182	183	183	18																																																										

Но почему задержка маленькими группами по 16 ядер?

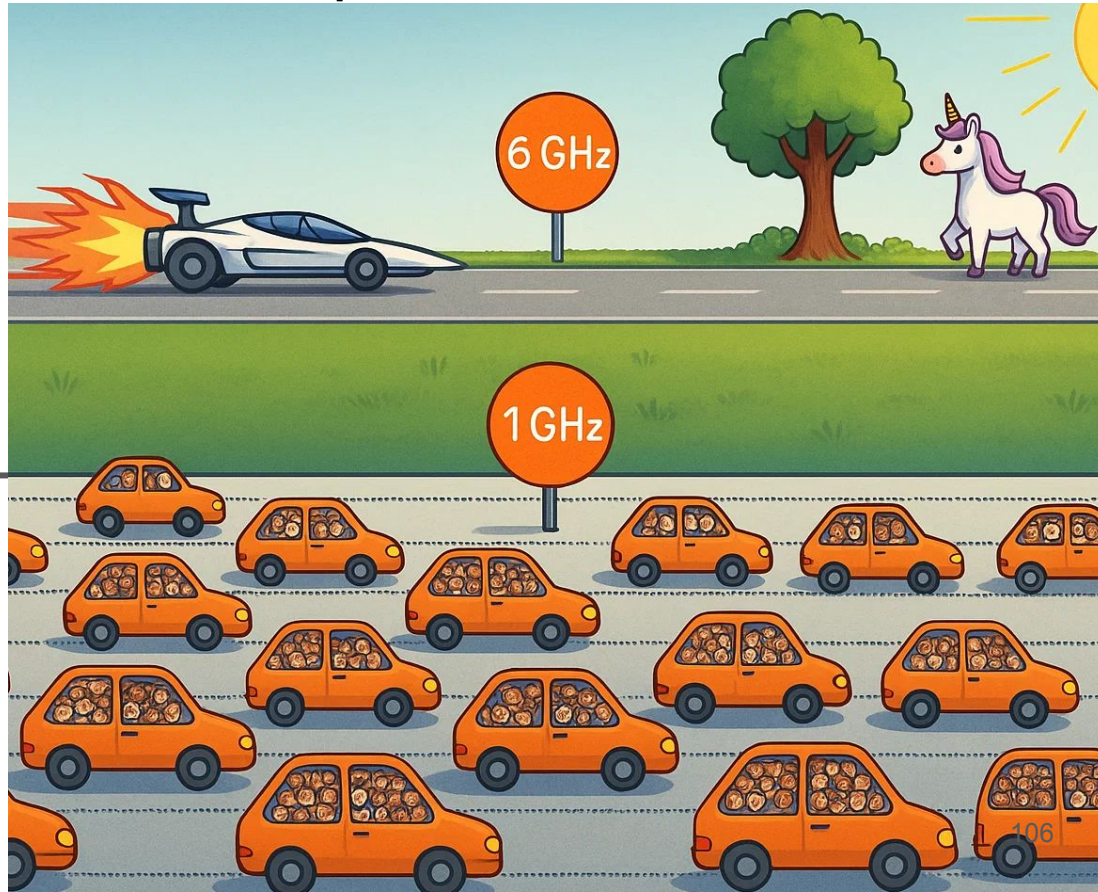
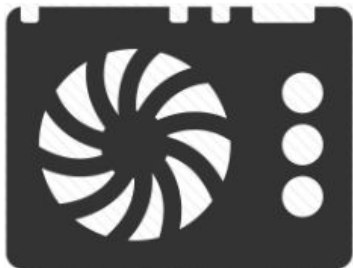
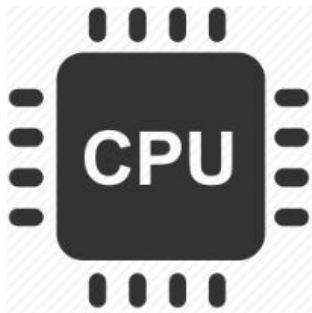
64 ядра = 8 x чипов по 8 ядер

AMD Ryzen Threadripper 7980X Core-to-Core Latency



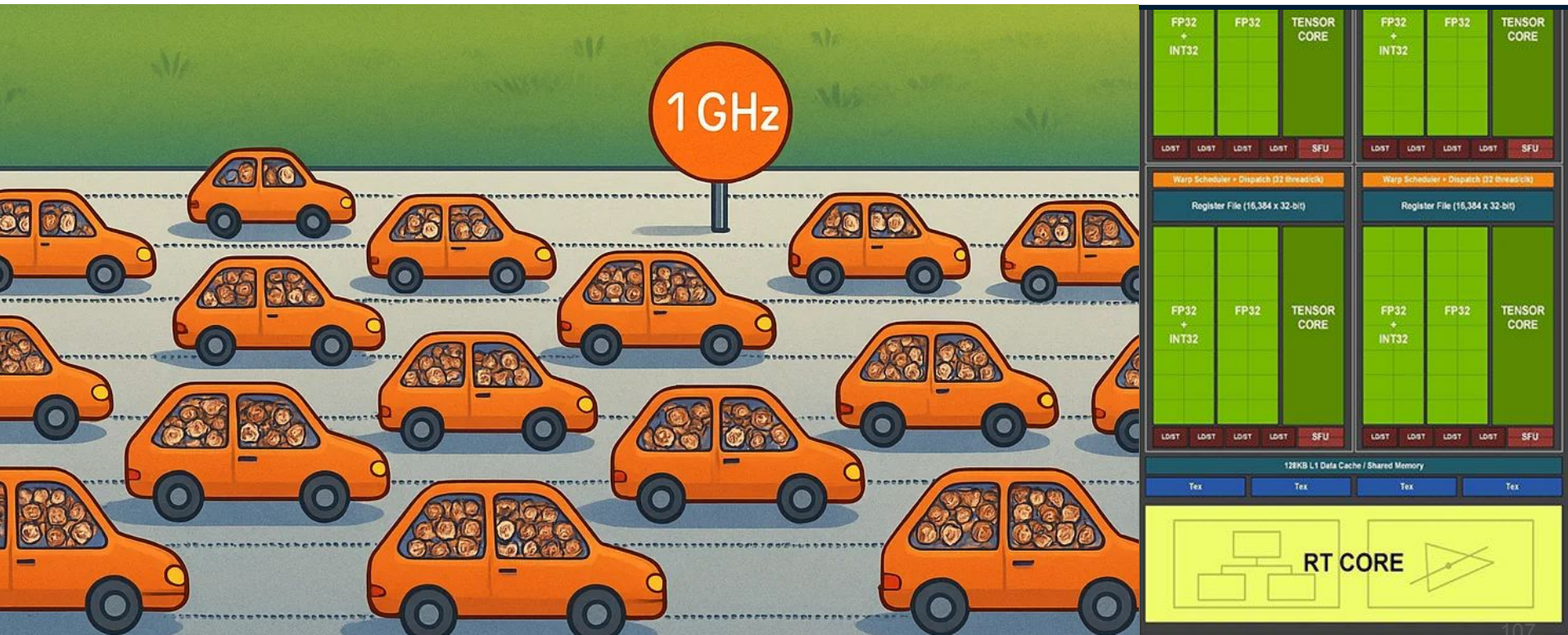
Глава 4: Модель вычислений массового параллелизма

Модель вычислений массового параллелизма



Модель вычислений массового параллелизма

SM - Streaming Multiprocessor



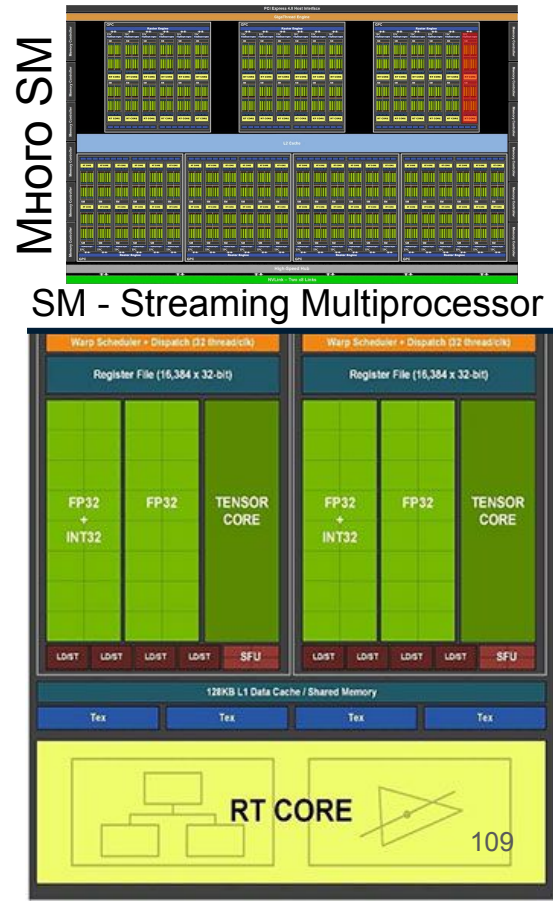
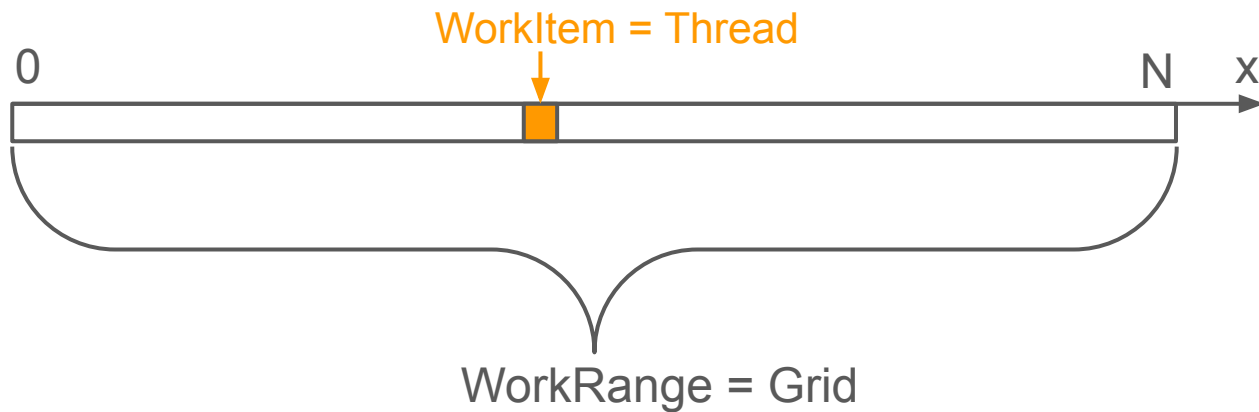
Модель вычислений массового параллелизма



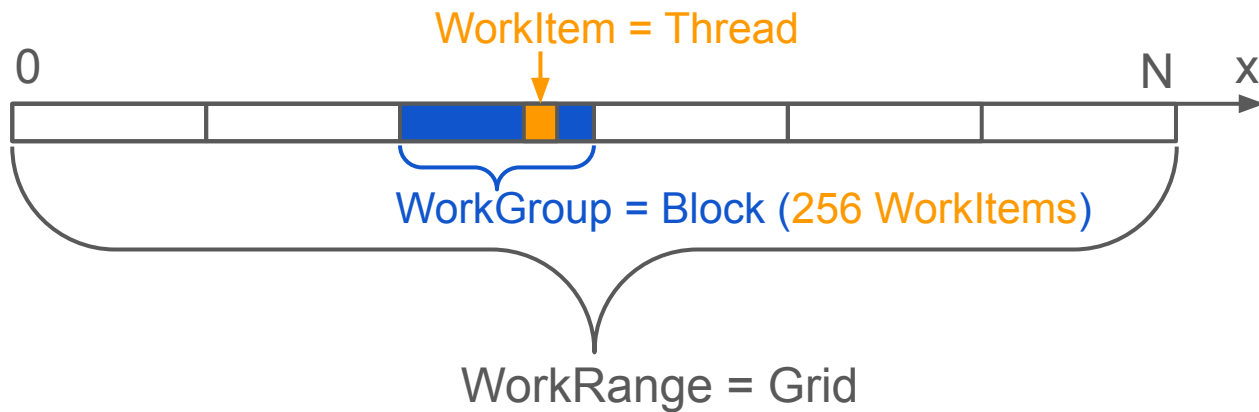
SM - Streaming Multiprocessor



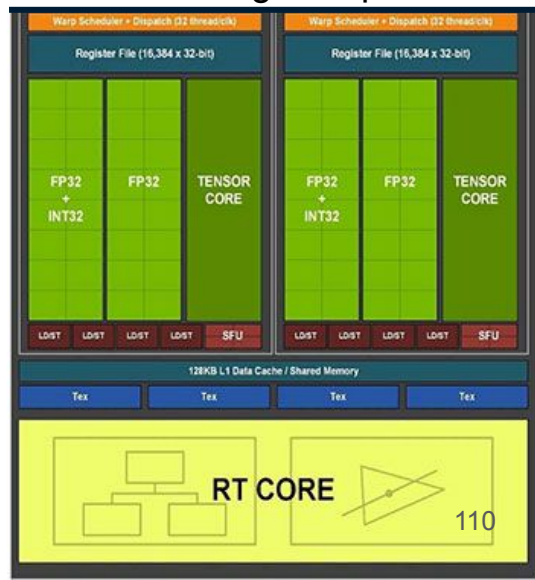
Модель вычислений массового параллелизма



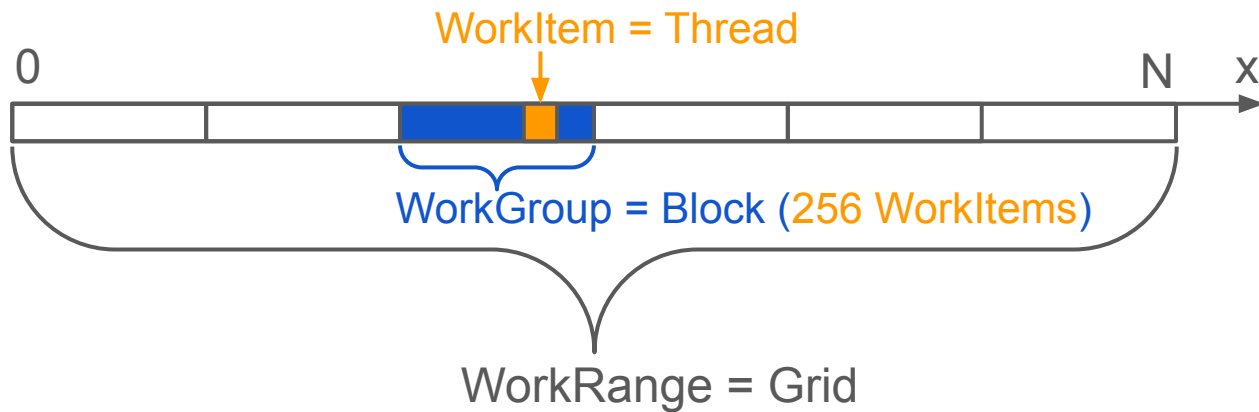
Модель вычислений массового параллелизма



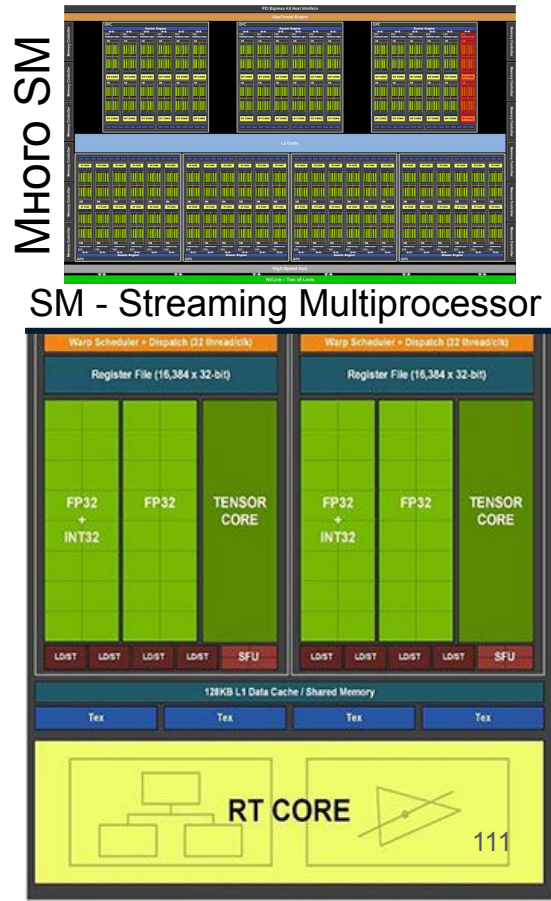
SM - Streaming Multiprocessor



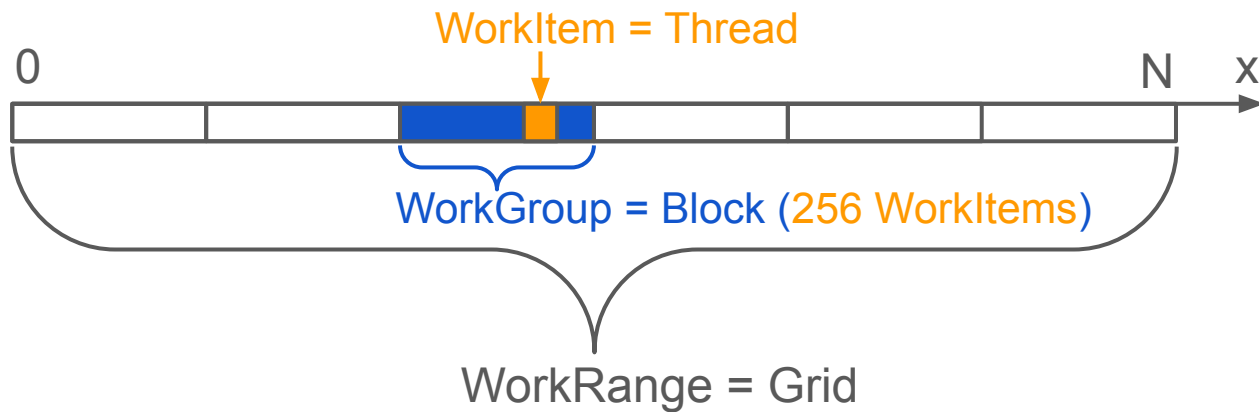
Модель вычислений массового параллелизма



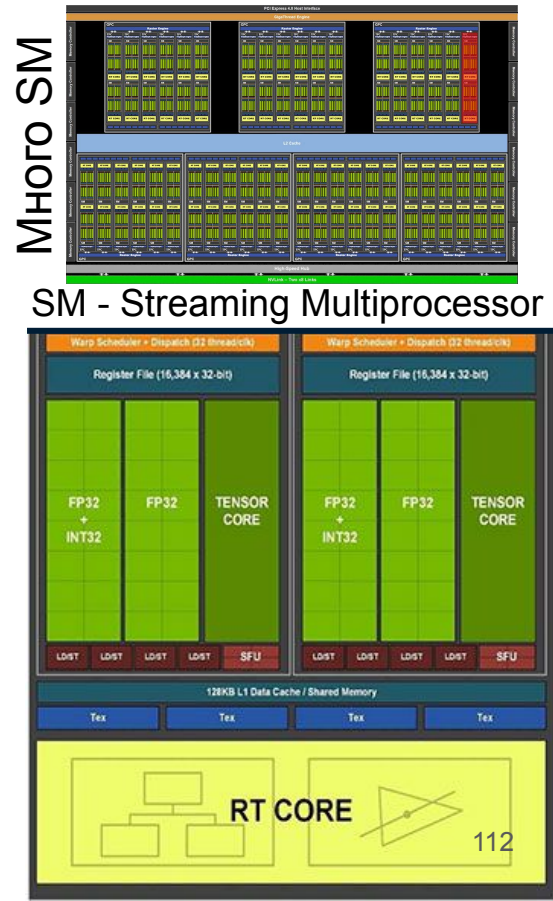
Но ведь в одном warp 32 потока!
(у **AMD**: 64 потока в wavefront)
Как так?



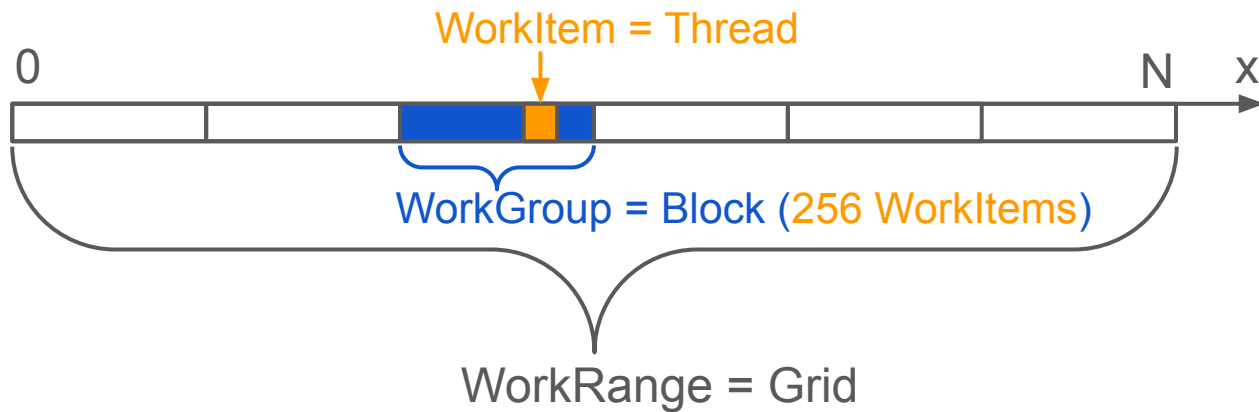
Модель вычислений массового параллелизма



WorkGroup = Block (256 WorkItems) = 8 x Warps (32 threads)

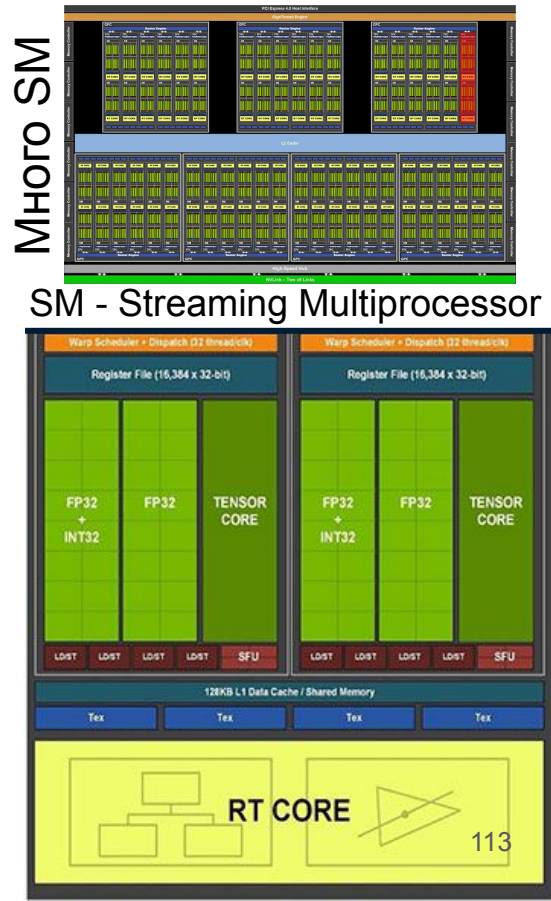


Модель вычислений массового параллелизма

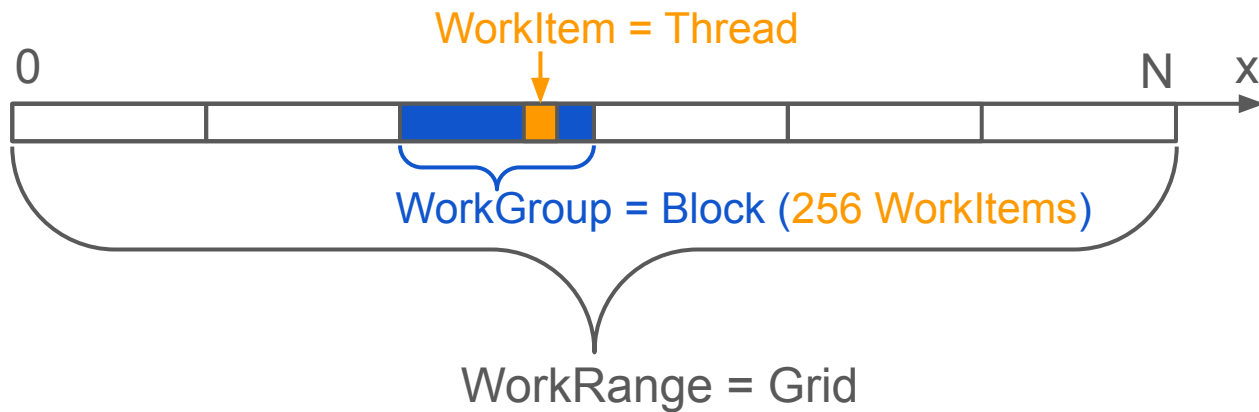


WorkGroup = Block (256 WorkItems) = 8 x Warps (32 threads)

Зачем нам контроль над **WorkGroups**?
Зачем на уровне API знать про **warps**?

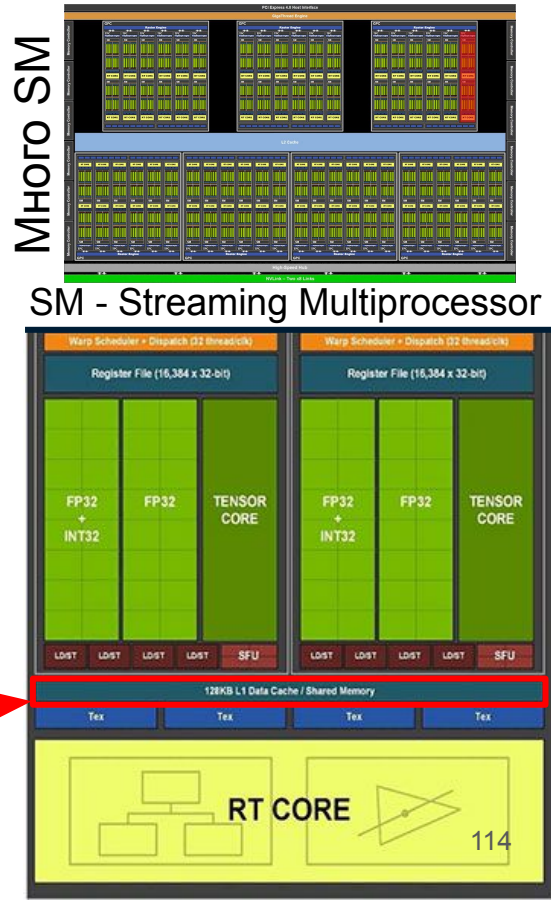


Модель вычислений массового параллелизма

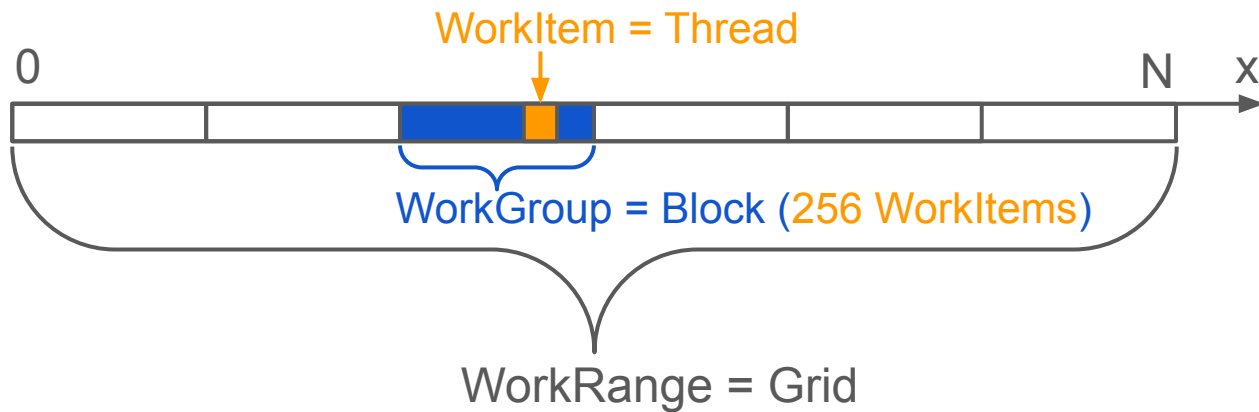


WorkGroup = Block (256 WorkItems) = 8 x Warps (32 threads)

- WorkItems коммуницируют через VRAM
- WorkItems в рамках одной WorkGroup общаются эффективнее - через **Shared/Local Memory (L1)**



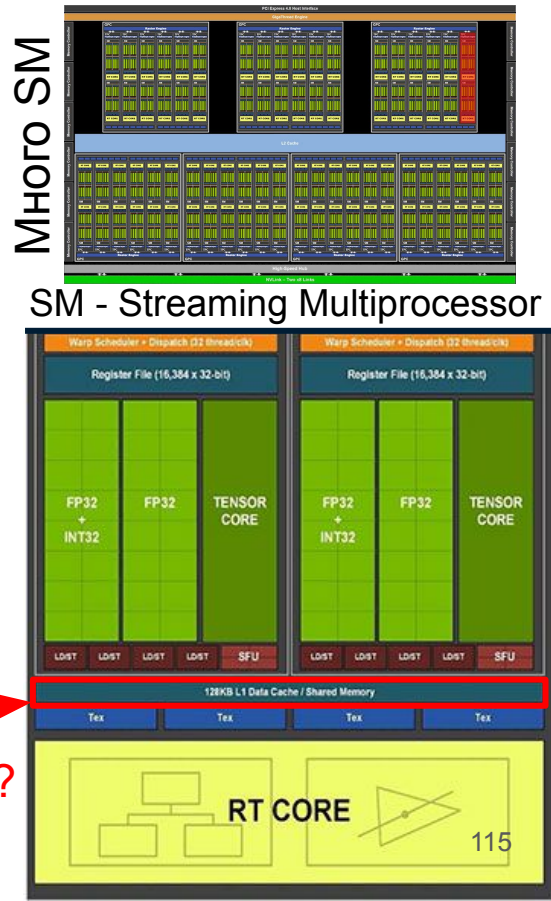
Модель вычислений массового параллелизма



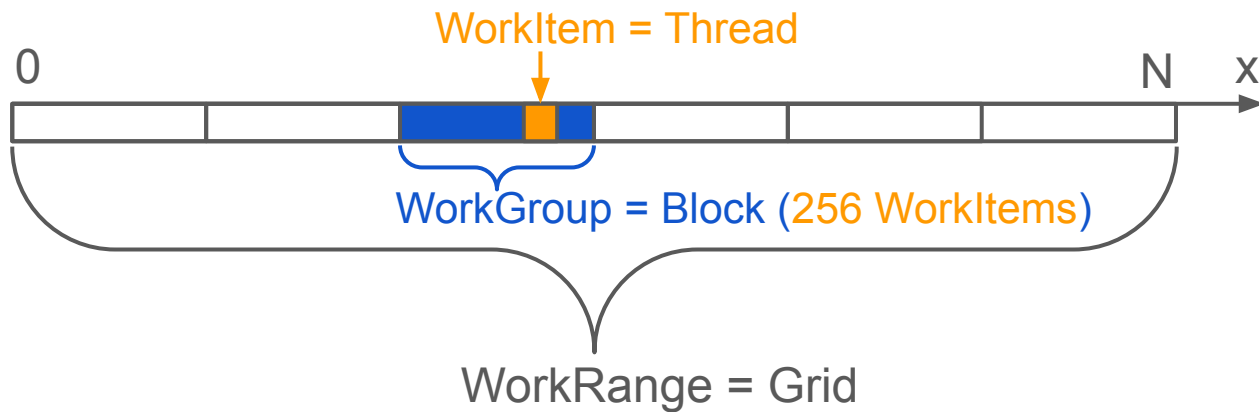
WorkGroup = Block (256 WorkItems) = 8 x Warps (32 threads)

- WorkItems коммуницируют через VRAM
- WorkItems в рамках одной WorkGroup общаются эффективнее - через **Shared/Local Memory (L1)**

Могут ли **warp**-ы одной **WorkGroup** исполняться на разных SM?



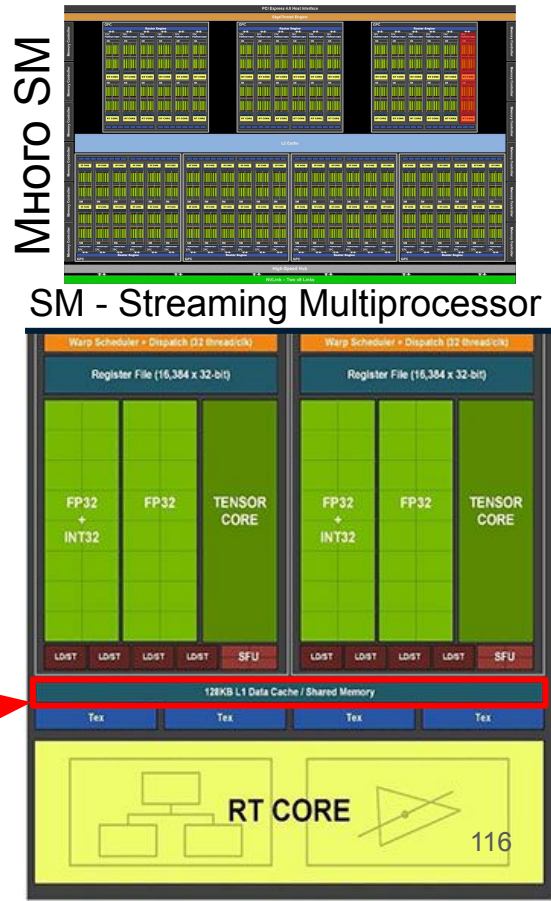
Модель вычислений массового параллелизма



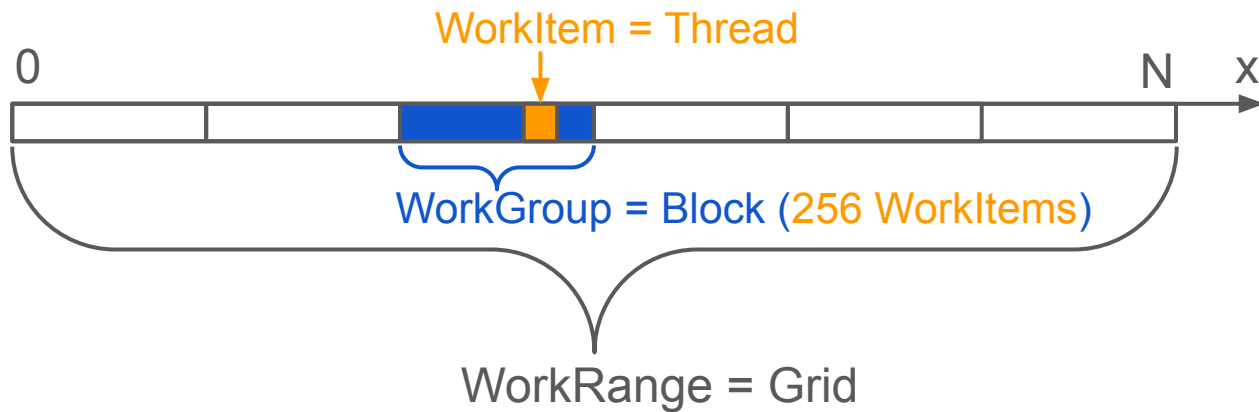
WorkGroup = Block (256 WorkItems) = 8 x Warps (32 threads)

- **WorkItems** коммуницируют через VRAM
- **WorkItems** в рамках одной **WorkGroup** общаются эффективнее - через **Shared/Local Memory (L1)**

А если один такой поток записал в **Local Memory** число - увидит ли другой поток этой **WorkGroup**?



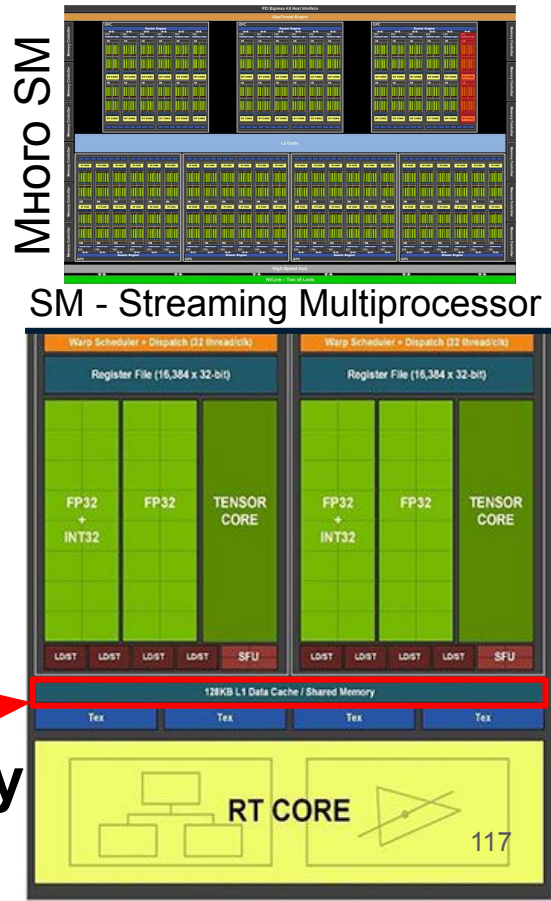
Модель вычислений массового параллелизма



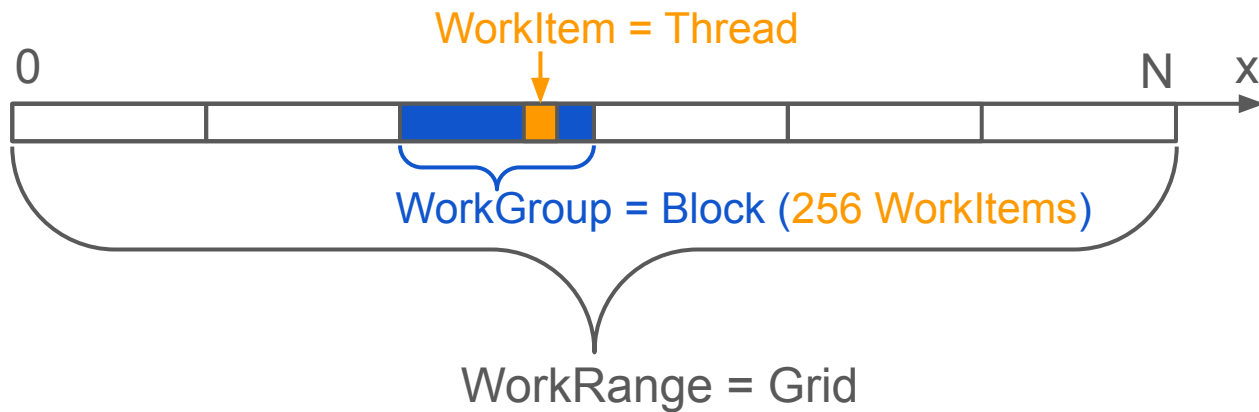
WorkGroup = Block (256 WorkItems) = 8 x Warps (32 threads)

- WorkItems коммуницируют через VRAM
- WorkItems в рамках одной WorkGroup общаются эффективнее - через **Shared/Local Memory (L1)**

А если один такой поток записал в **Local Memory** число - увидит ли другой поток этой **WorkGroup**? **barrier!**

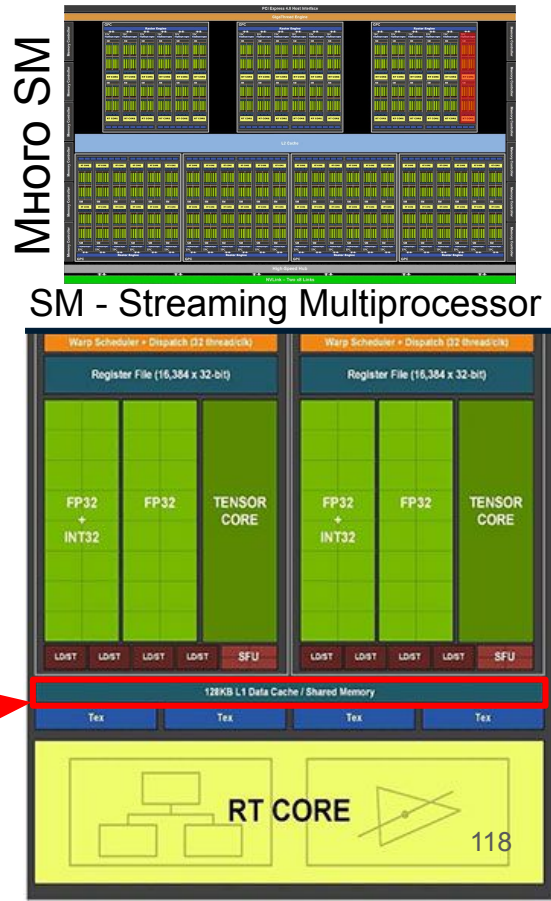


Модель вычислений массового параллелизма



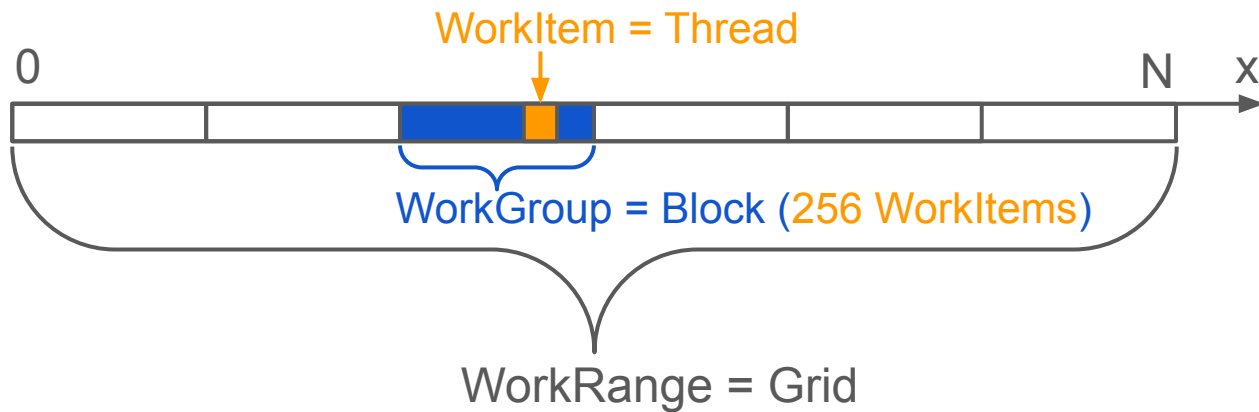
WorkGroup = Block (256 WorkItems) = 8 x Warps (32 threads)

- WorkItems коммуницируют через VRAM
- WorkItems в рамках одной WorkGroup общаются эффективнее - через **Shared/Local Memory (L1)**



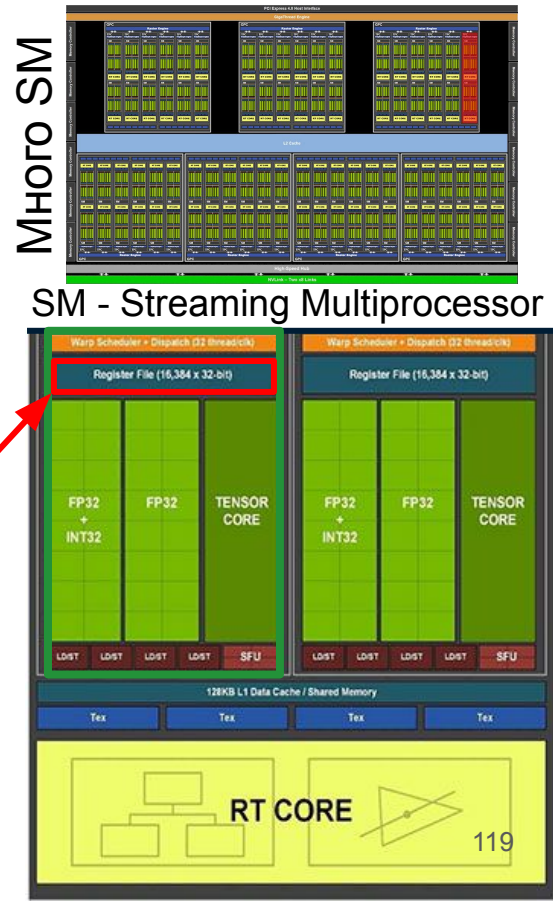
Могут ли потоки одного warp-а общаться еще эффективнее?

Модель вычислений массового параллелизма

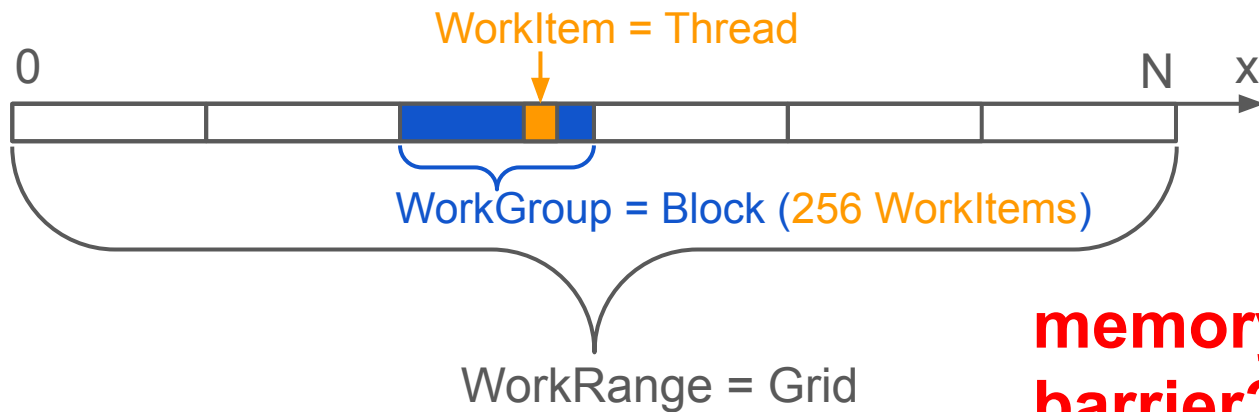


WorkGroup = Block (256 WorkItems) = 8 x Warps (32 threads)

- WorkItems коммуницируют через VRAM
- WorkItems в рамках одной WorkGroup общаются эффективнее - через **Shared/Local Memory (L1)**
- WorkItems в рамках одного warp-а могут подглядывать друг другу в **регистры** (shuffle instr.)



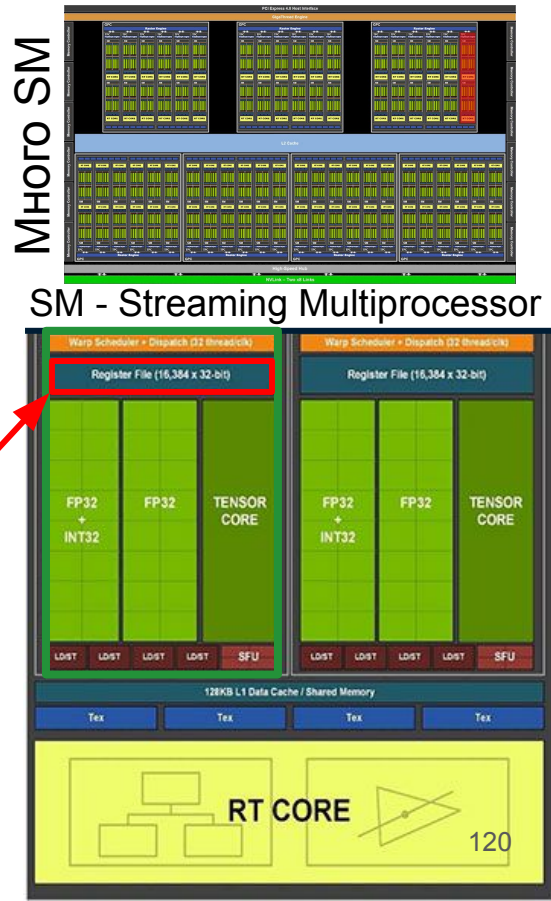
Модель вычислений массового параллелизма



WorkGroup = Block (256 WorkItems) = 8 x Warps (32 threads)

- WorkItems коммуницируют через VRAM
- WorkItems в рамках одной WorkGroup общаются эффективнее - через **Shared/Local Memory (L1)**
- WorkItems в рамках одного warp-а могут подглядывать друг другу в **регистры** (shuffle instr.)

memory barrier?

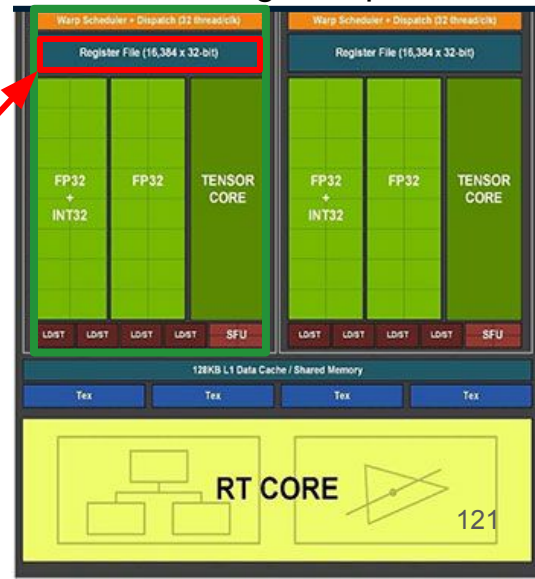


GLSL (Graphics Library Shading Language)

dFdx(...), dFdy(...) — return the partial derivative of an argument
with respect to x or y

- **WorkItems** коммуницируют через VRAM
- **WorkItems** в рамках одной **WorkGroup** общаются эффективнее - через **Shared/Local Memory (L1)**
- **WorkItems** в рамках одного warp-а могут подглядывать друг другу в **регистры** (**shuffle instr.**)

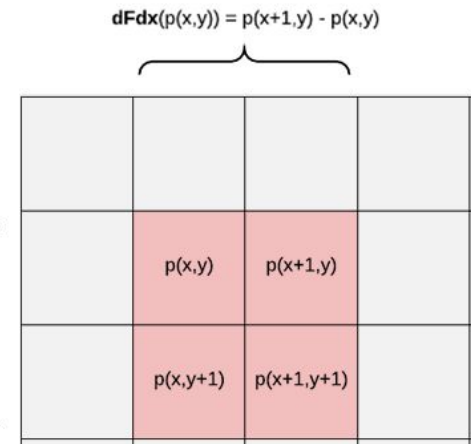
SM - Streaming Multiprocessor



GLSL (Graphics Library Shading Language)

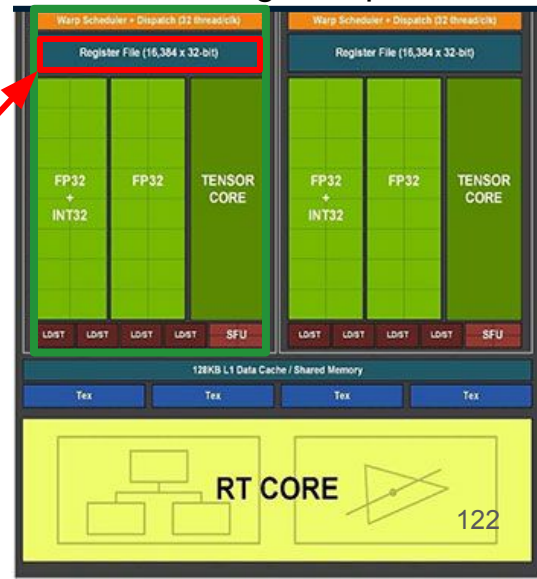
dFdx(...), dFdy(...) — return the partial derivative of an argument with respect to x or y

$$dFdy(p(x,y)) = p(x,y+1) - p(x,y)$$



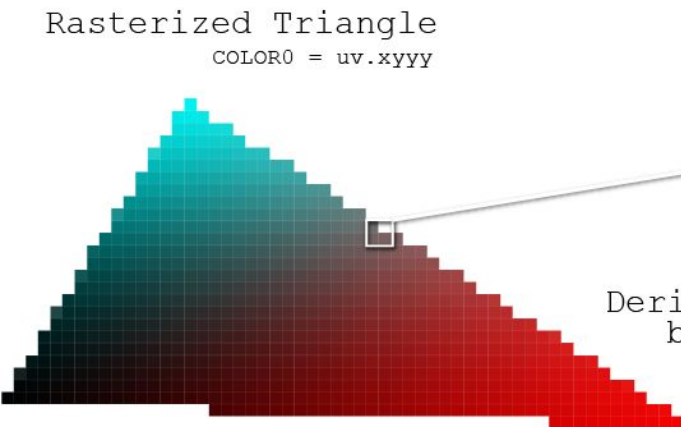
SM - Streaming Multiprocessor

- **WorkItems** коммуницируют через VRAM
- **WorkItems** в рамках одной **WorkGroup** общаются эффективнее - через **Shared/Local Memory (L1)**
- **WorkItems** в рамках одного warp-а могут подглядывать друг другу в **регистры** (**shuffle instr.**)



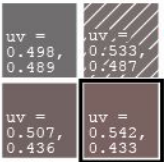
GLSL (Graphics Library Shading Language)

dFdx(...), dFdy(...) — return the partial derivative of an argument with respect to x or y



Rasterized Triangle
COLOR0 = uv.xyyy

Shading Quad
(2x2 screen pixels)

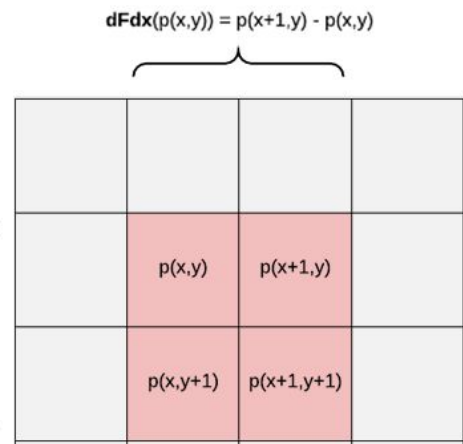


$dFdy(p(x,y)) = p(x,y+1) - p(x,y)$

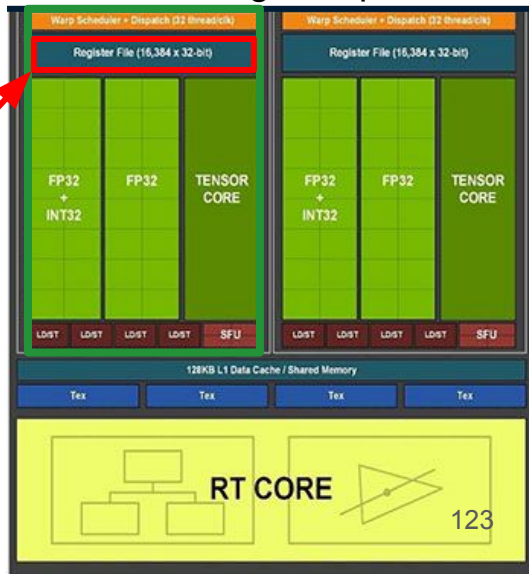
$ddy(uv.x) = 0.533 - 0.542 = -0.009$
 $ddy(uv.y) = 0.487 - 0.433 = 0.054$

Derivatives for
bottom-right
fragment
(high-precision)

$ddx(uv.x) = 0.542 - 0.507 = 0.035$
 $ddx(uv.y) = 0.433 - 0.436 = -0.003$

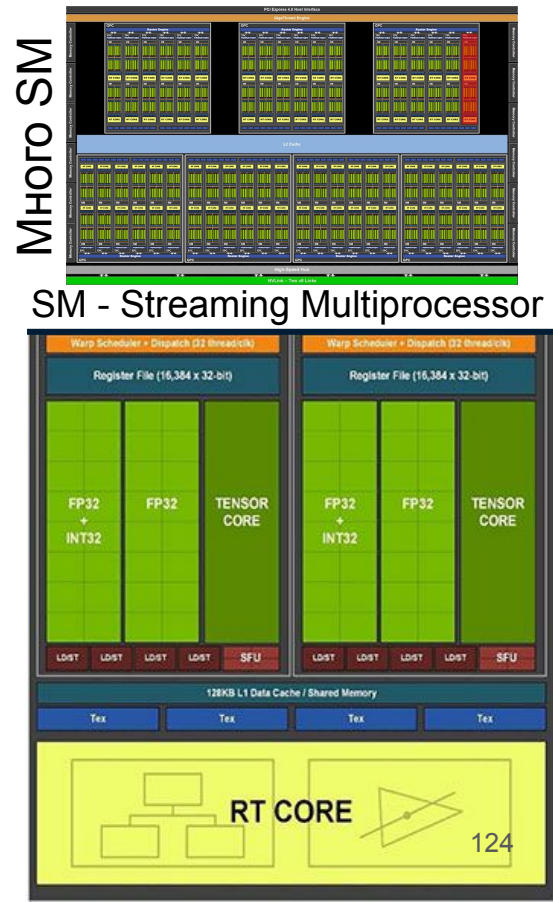
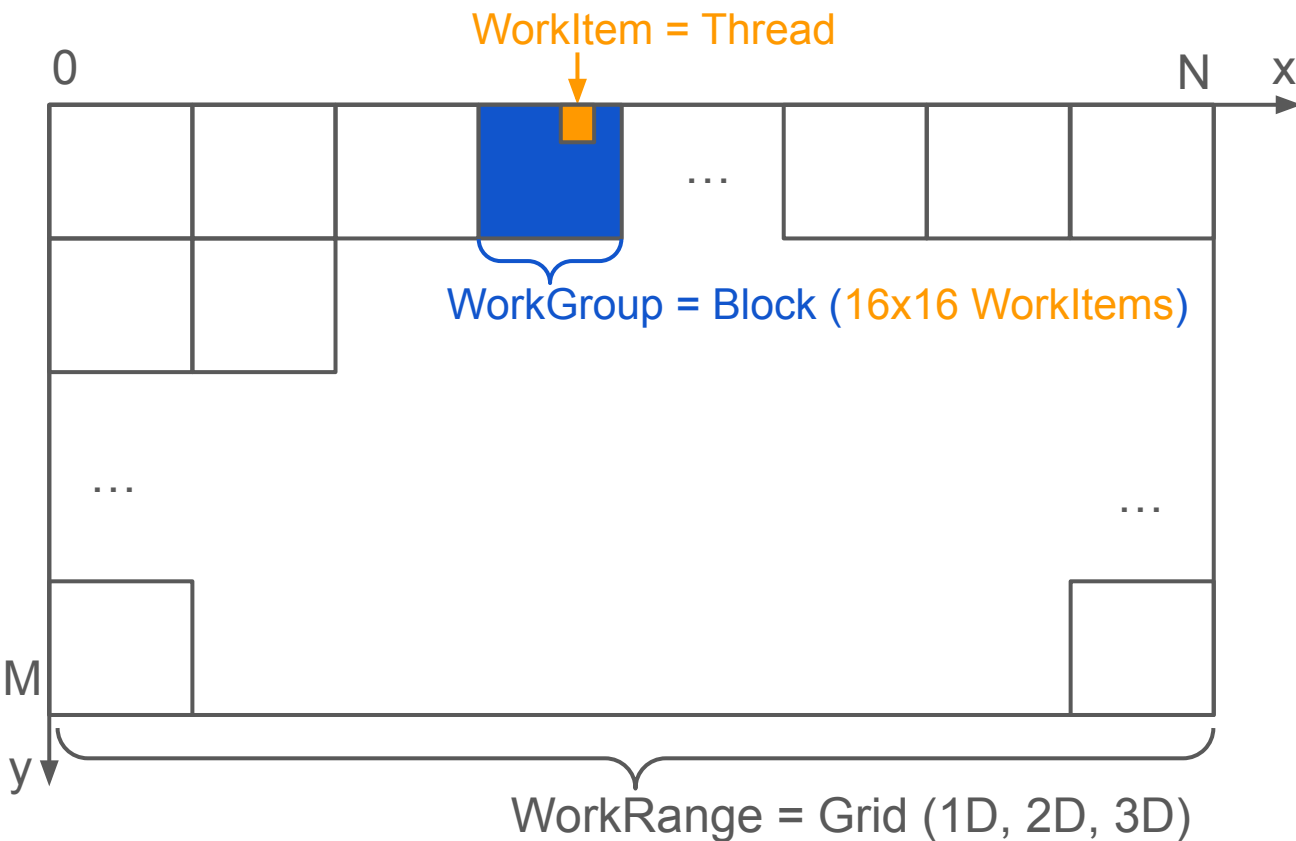


SM - Streaming Multiprocessor

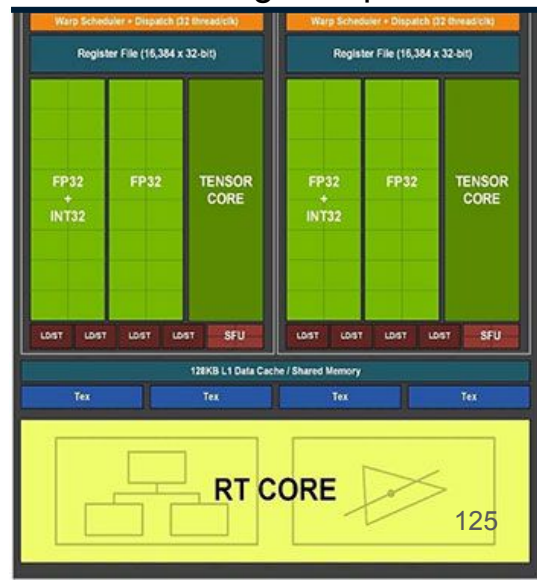
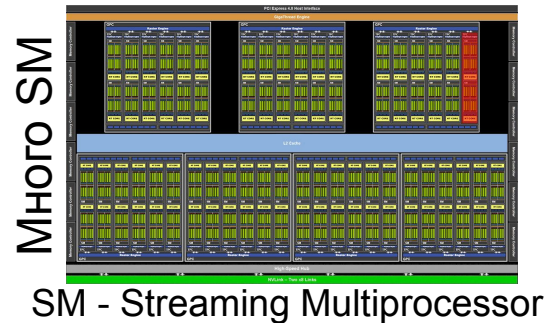
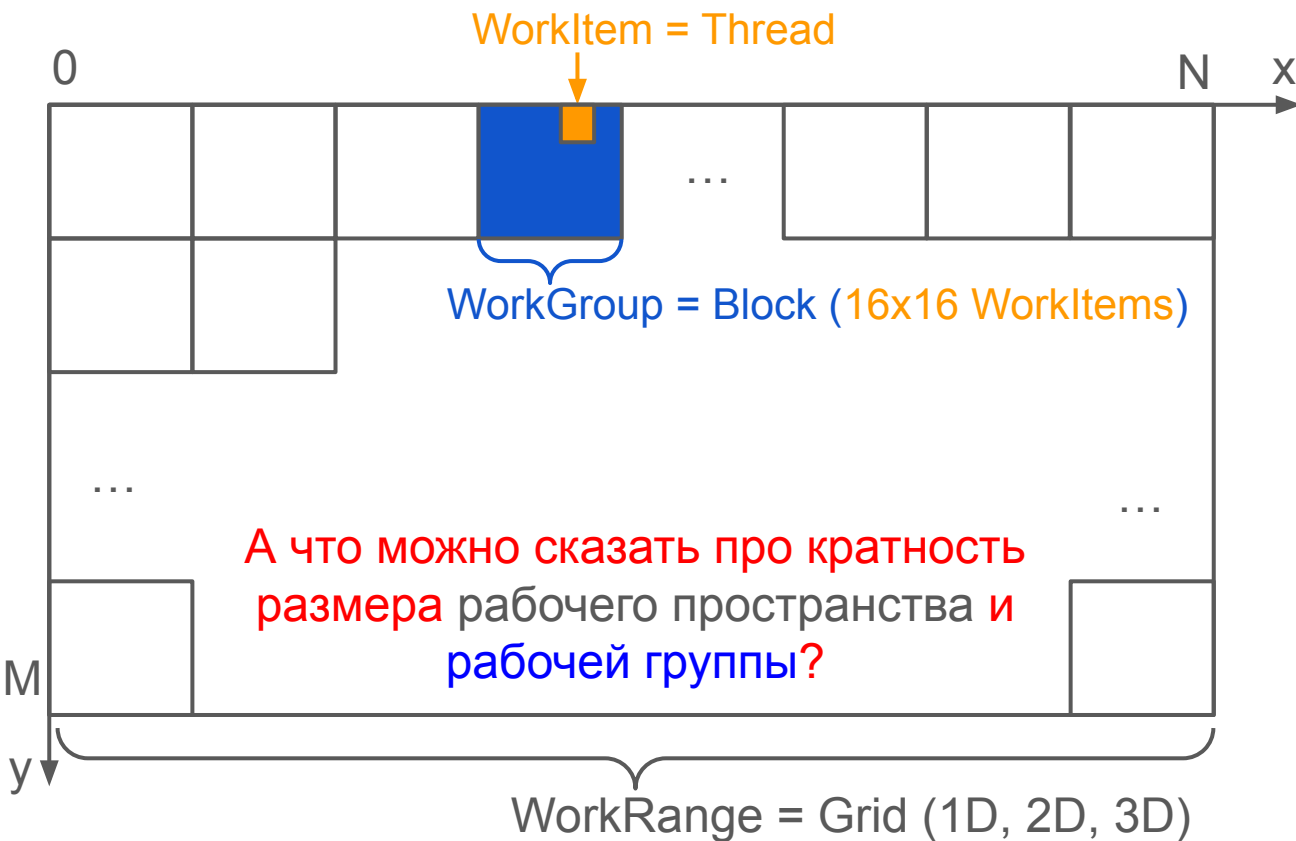


- **WorkItems** коммуницируют через VRAM
- **WorkItems** в рамках одной **WorkGroup** общаются эффективнее - через **Shared/Local Memory (L1)**
- **WorkItems** в рамках одного warp-а могут подглядывать друг другу в **регистры** (**shuffle instr.**)

Модель вычислений массового параллелизма



Модель вычислений массового параллелизма

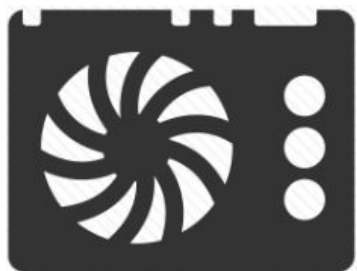


Глава 5: Профилирование и оптимизация

МНОГО вычислений? (**ALU**)

МНОГО читать+писать? (**VRAM**)

$100 \cdot 10^{12}$ FLOPS



1000 GB/s



VRAM - GDDR6

Чего больше?
Как сравнить яблоки и
апельсины?

МНОГО вычислений? (**ALU**)

МНОГО читать+писать? (**VRAM**)

$100 \cdot 10^{12}$ FLOPS



1000 GB/s



VRAM - GDDR6

Чего больше?

МНОГО вычислений? (**ALU**)

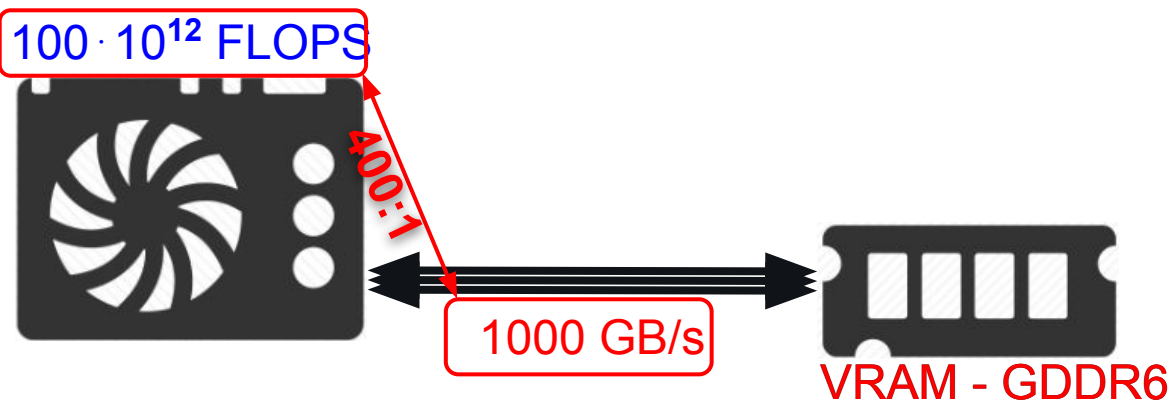
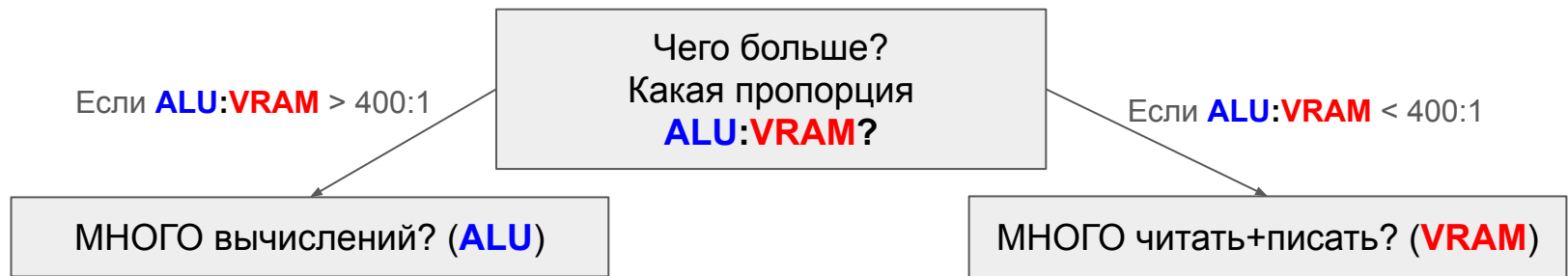
МНОГО читать+писать? (**VRAM**)

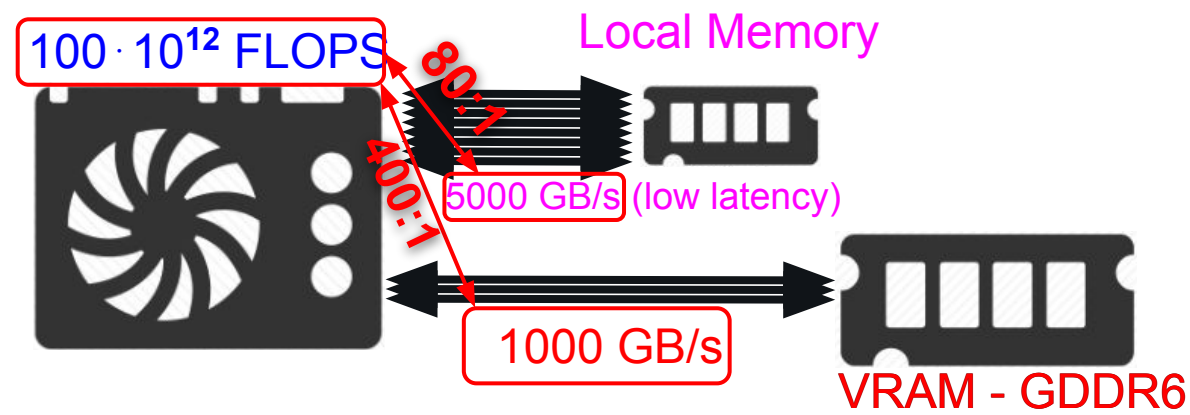
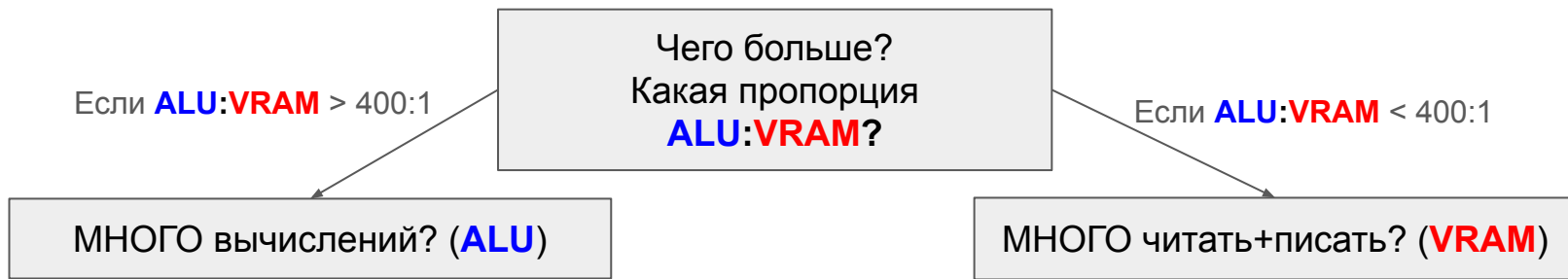
$100 \cdot 10^{12}$ FLOPS

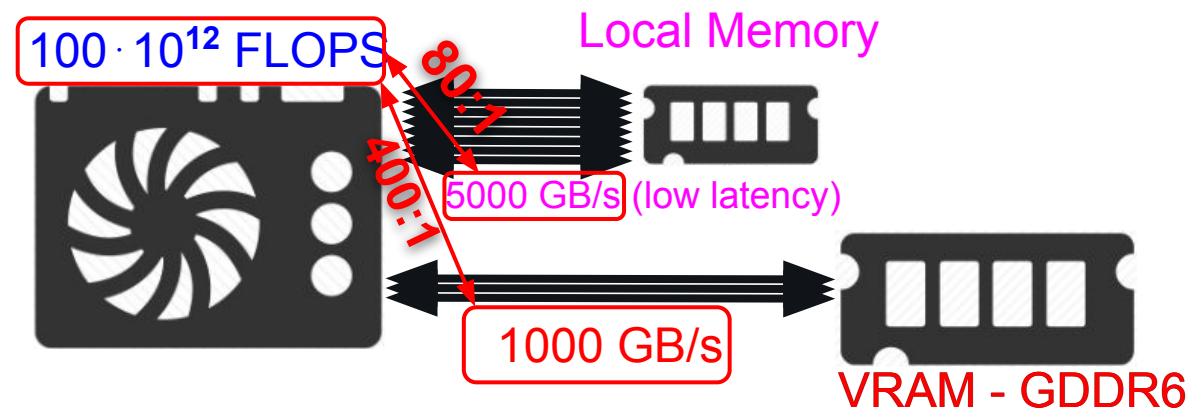
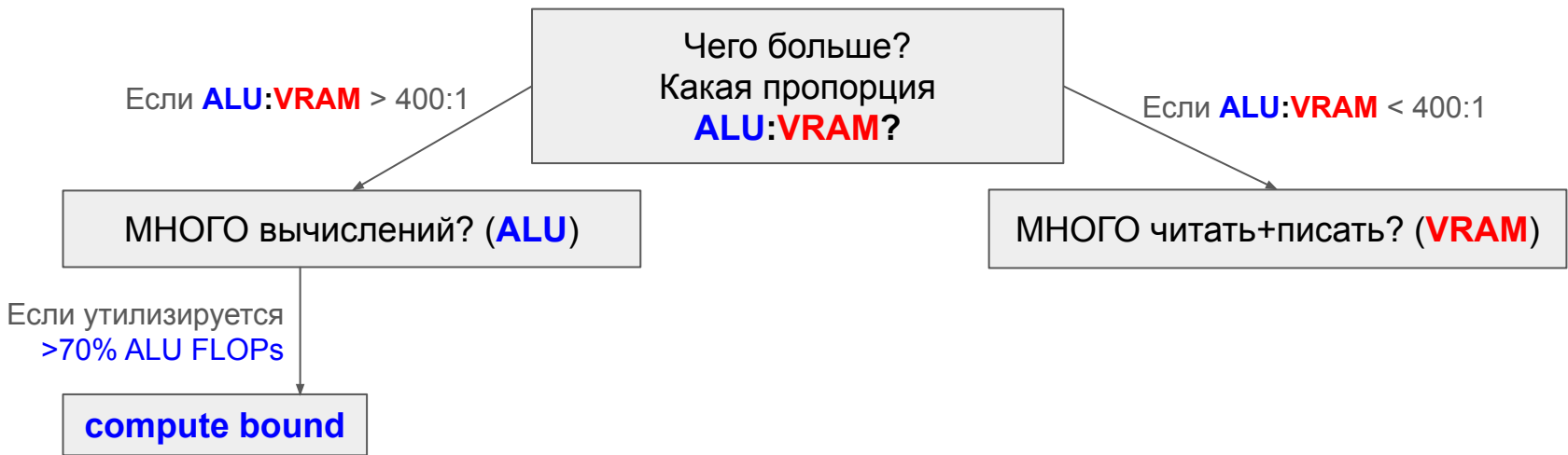
400:1

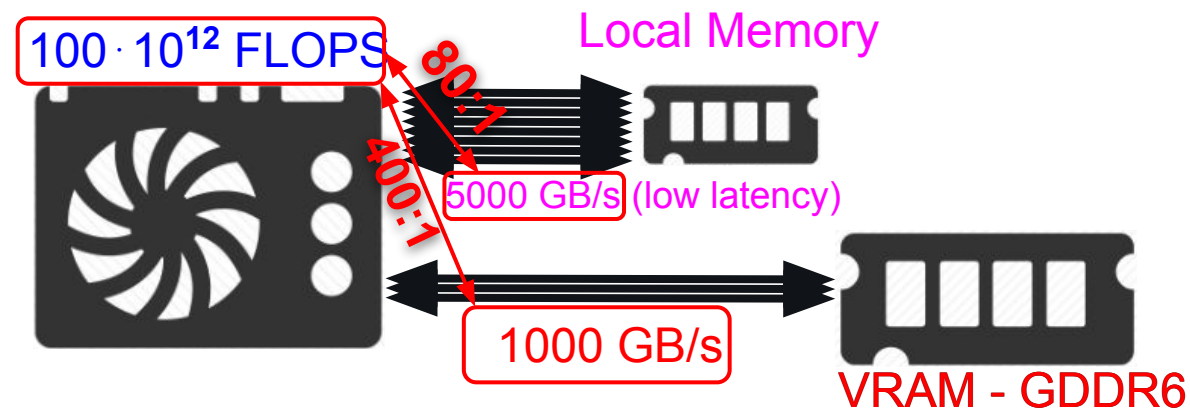
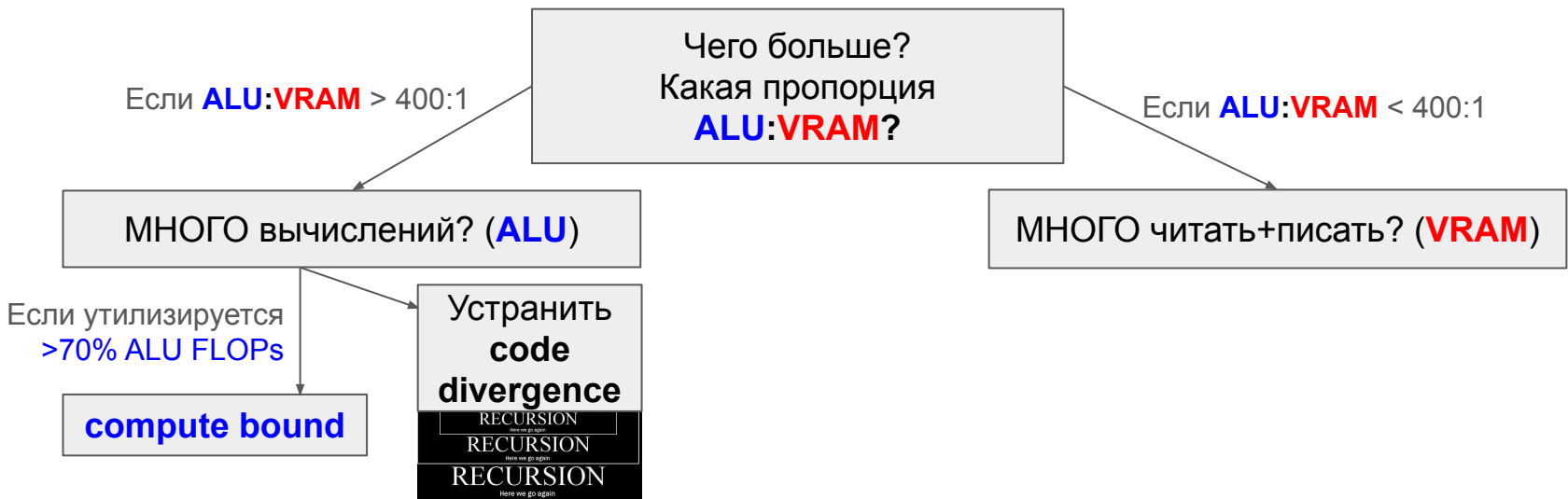
1000 GB/s

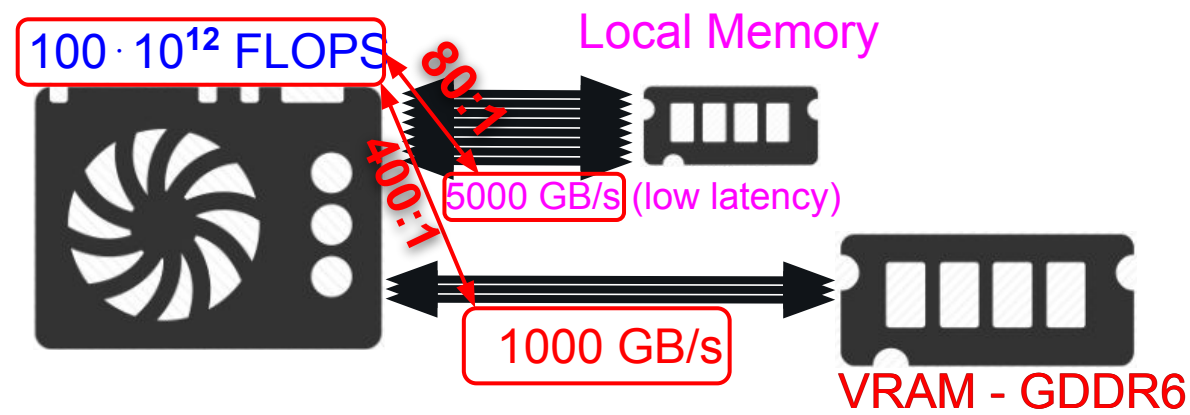
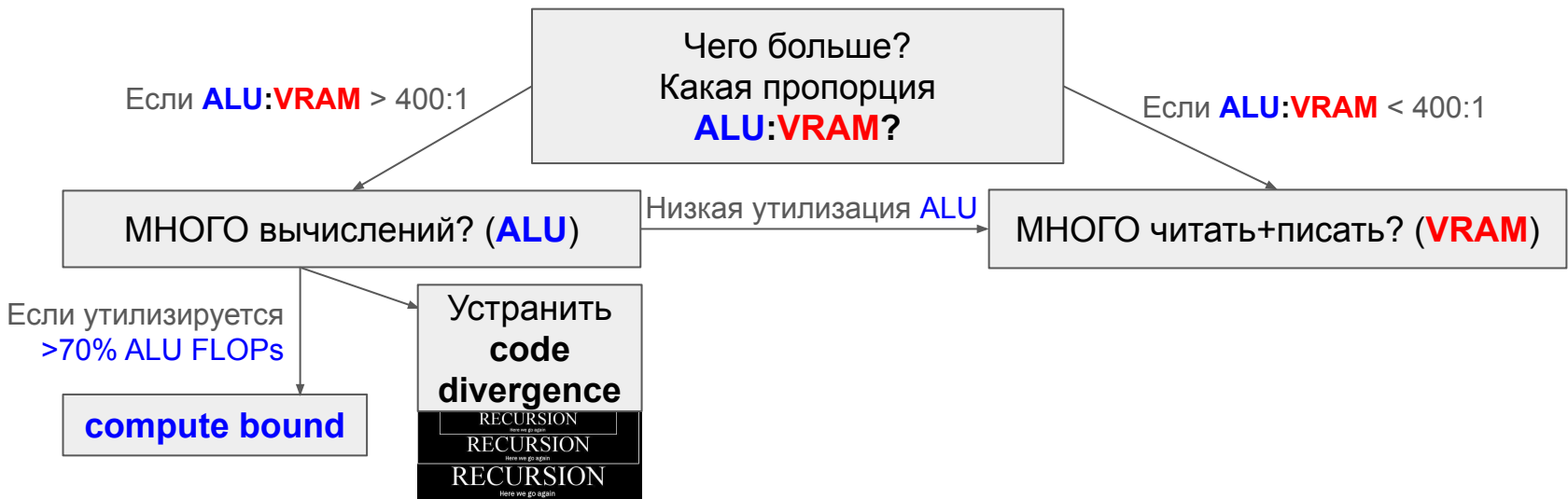
VRAM - GDDR6

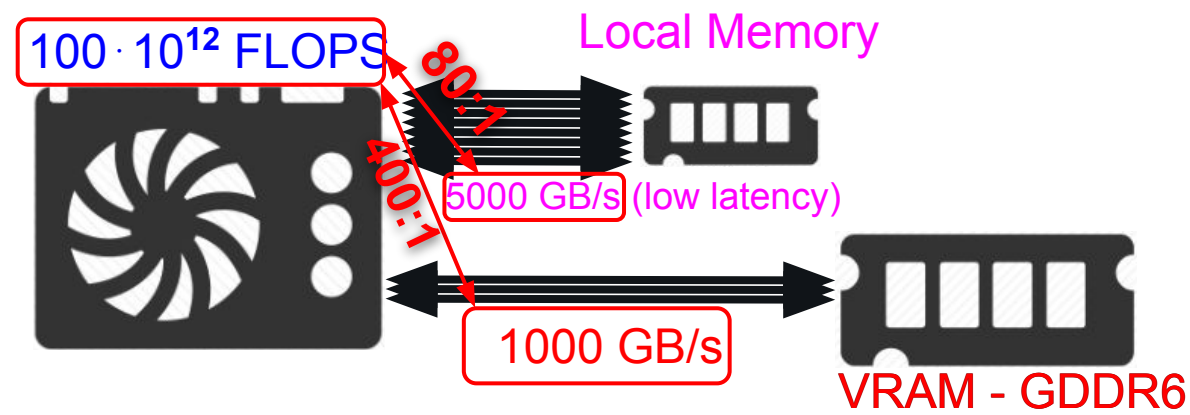


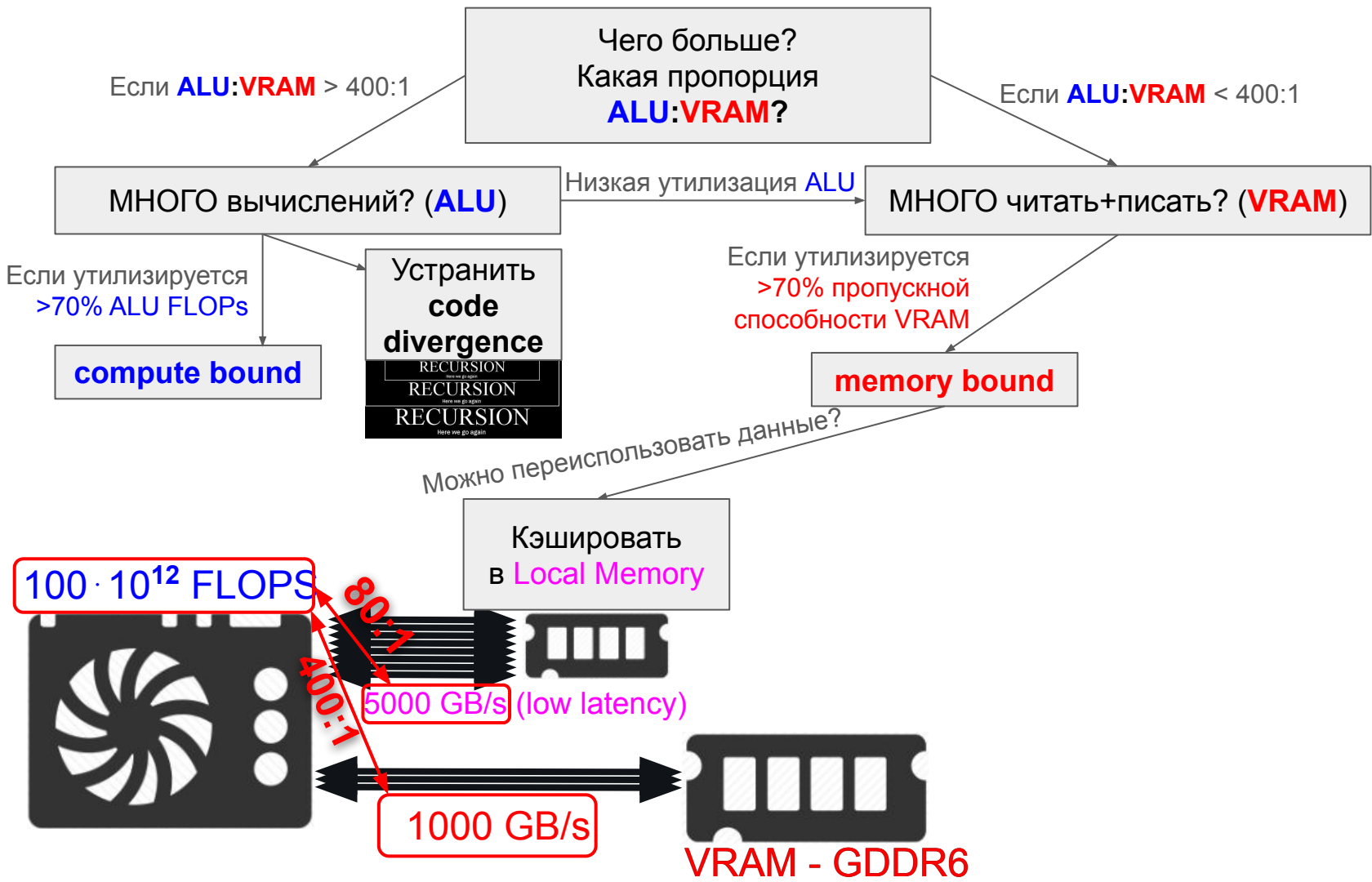


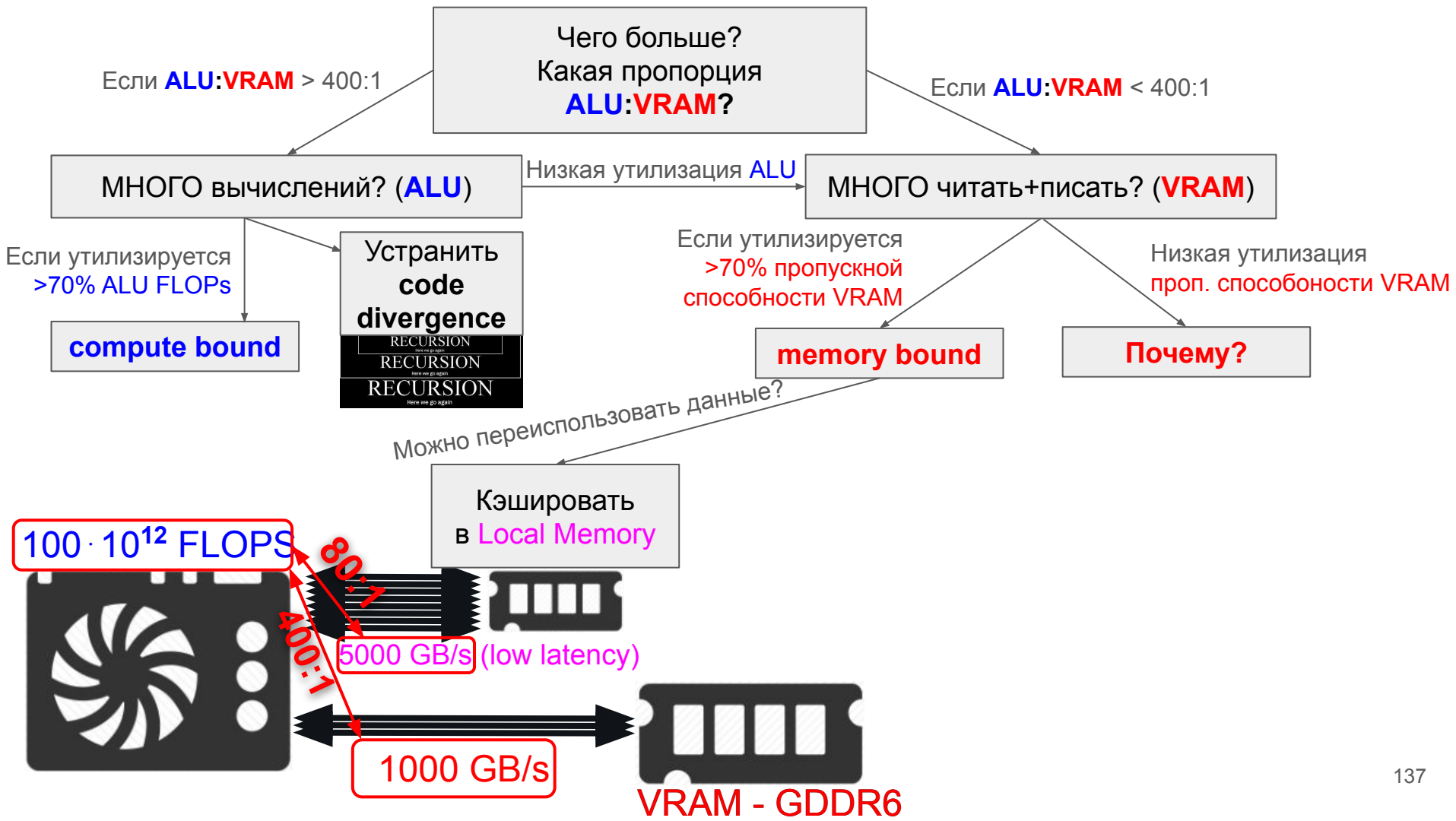


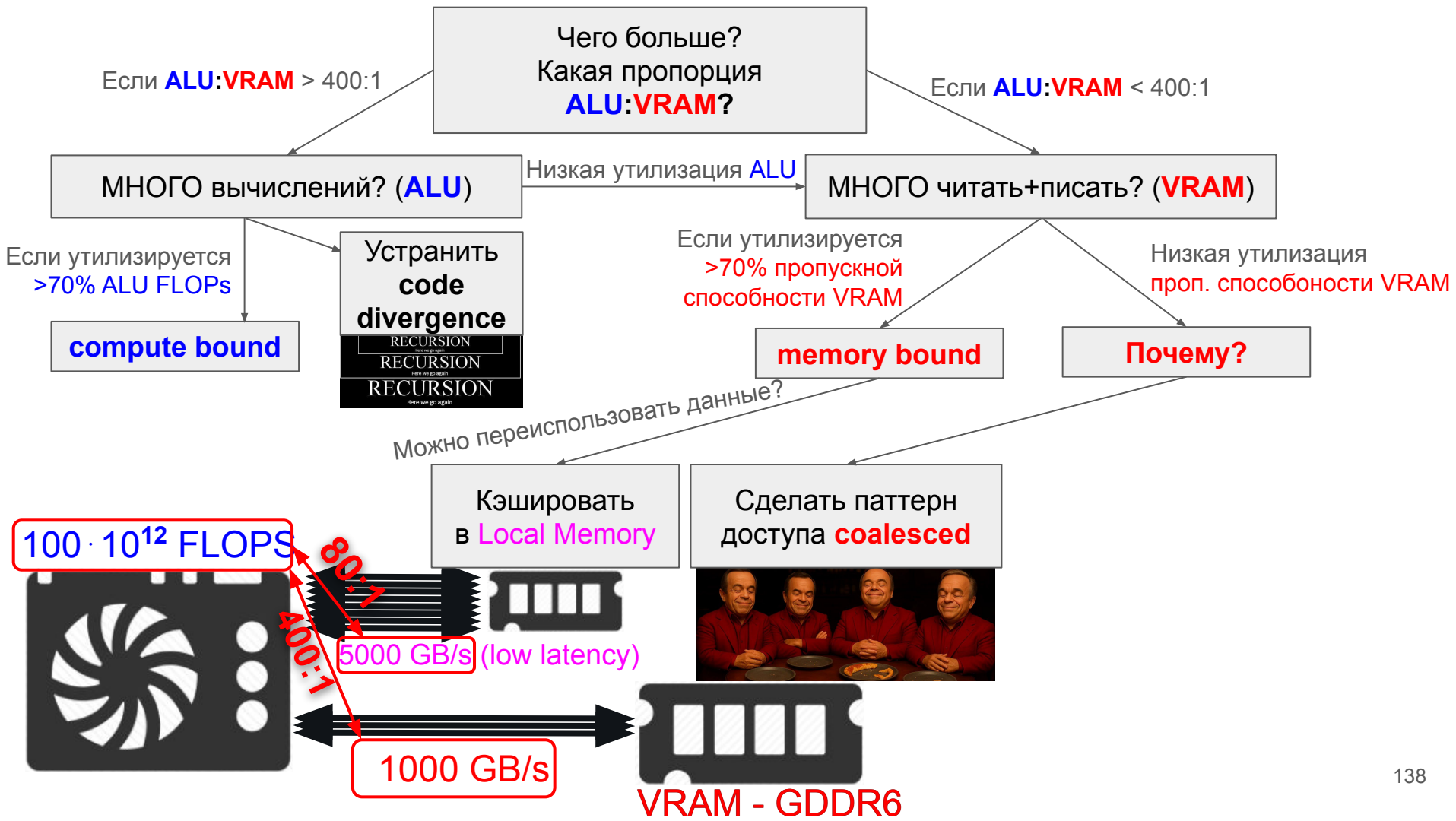


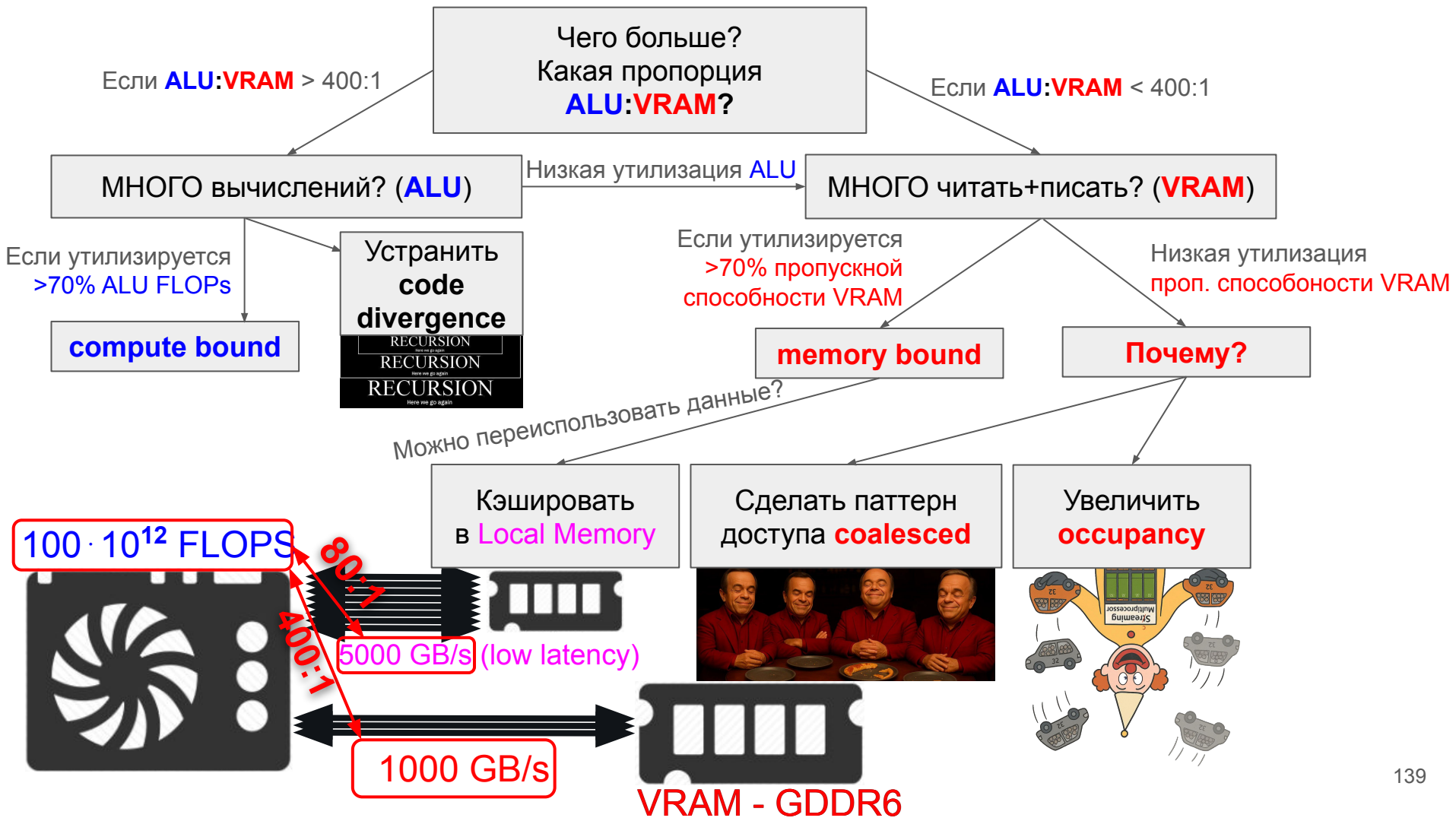


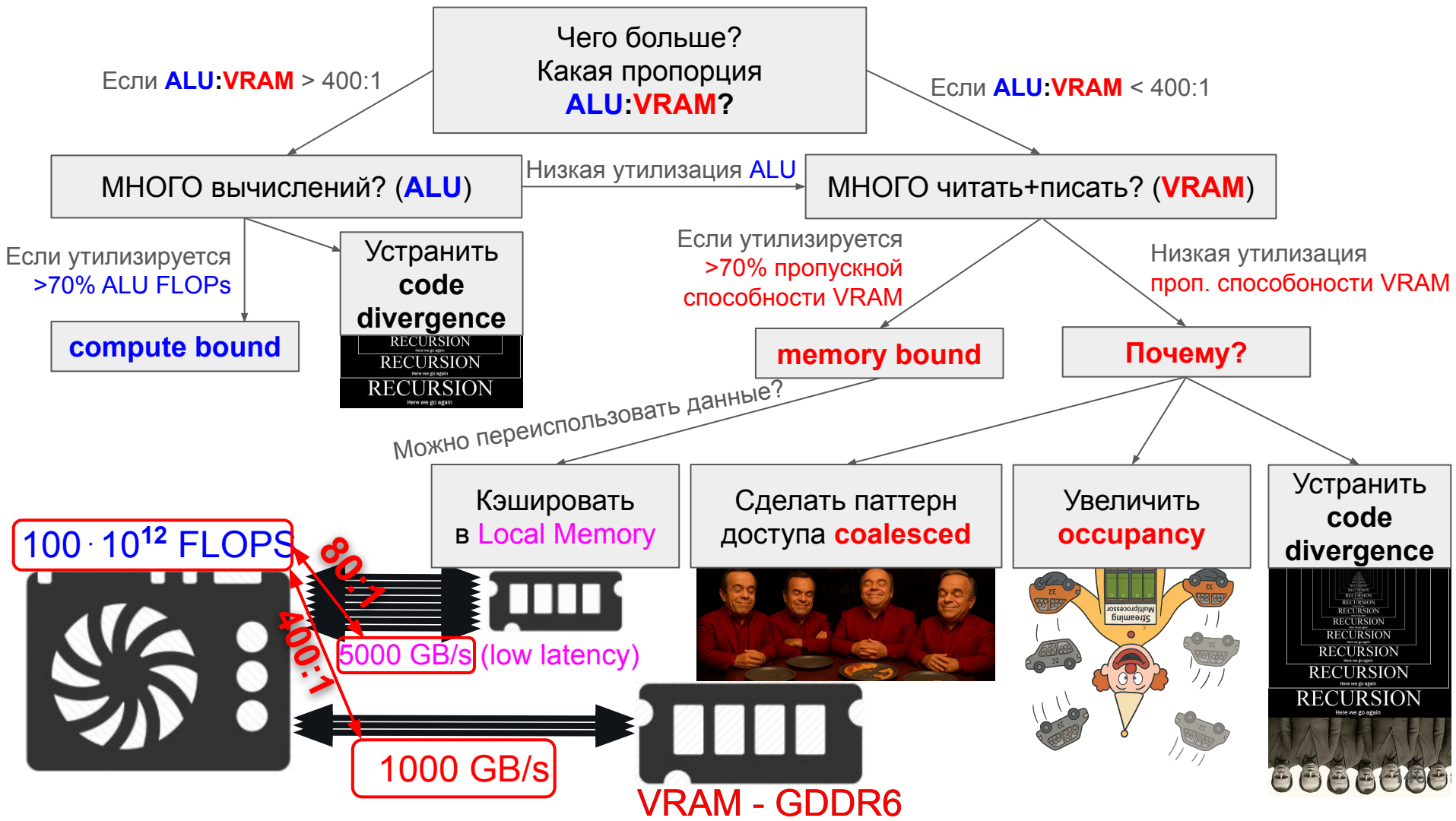


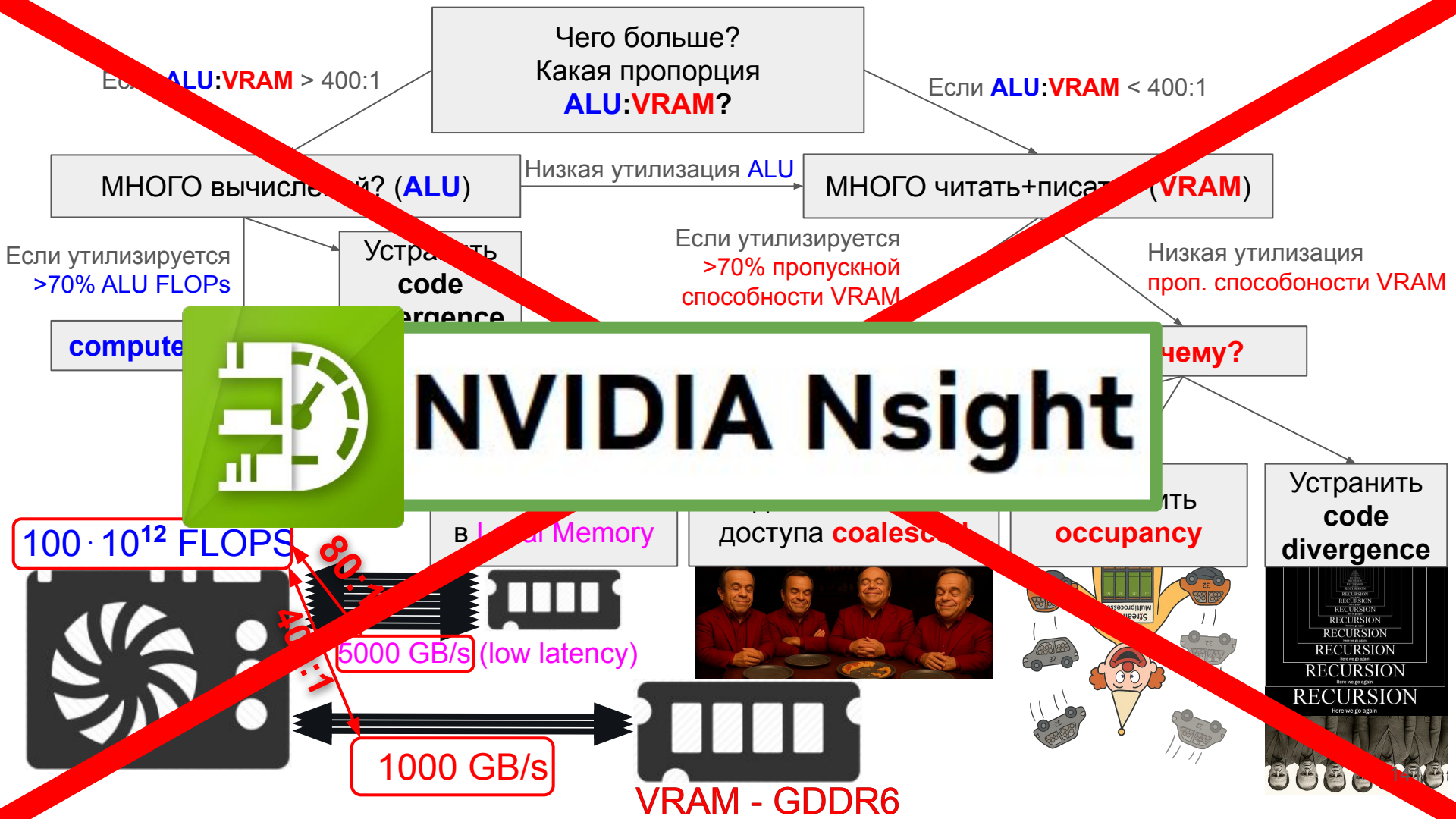






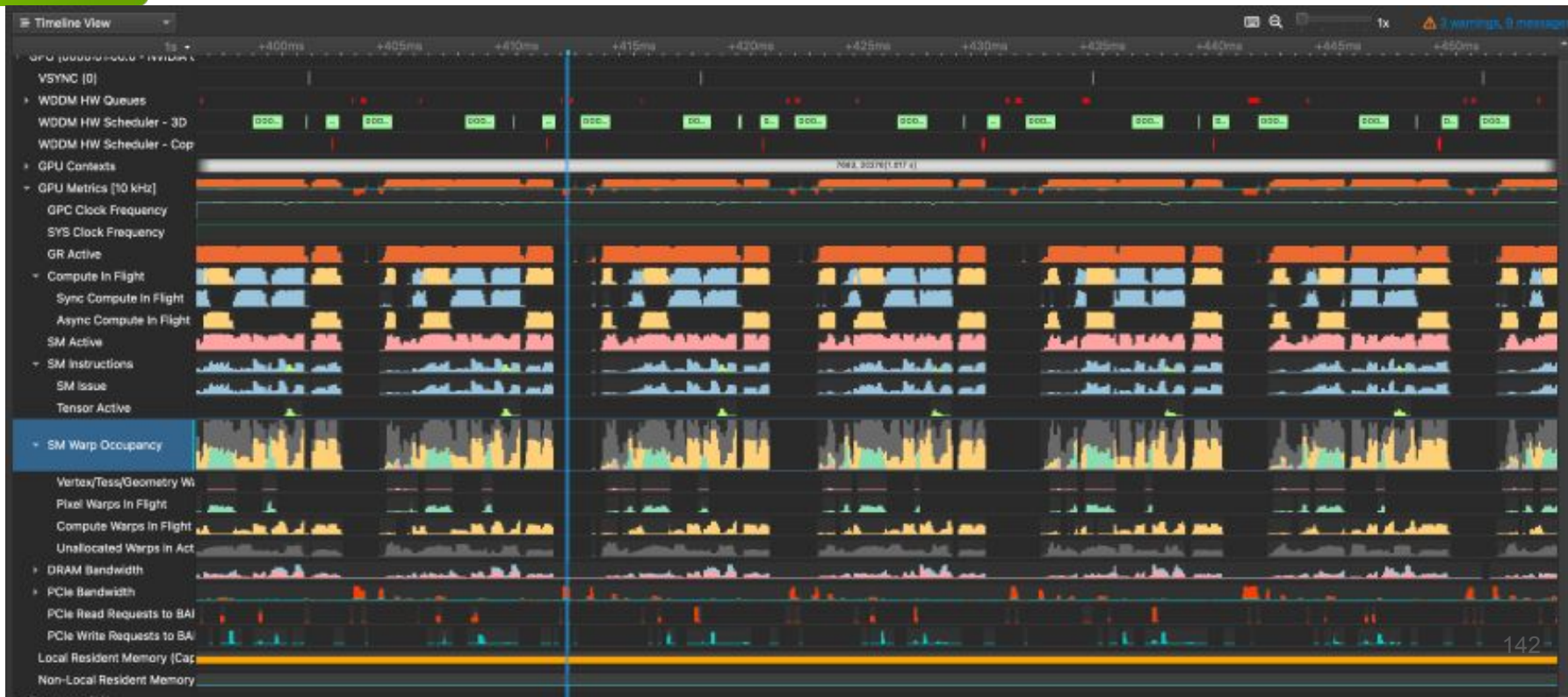


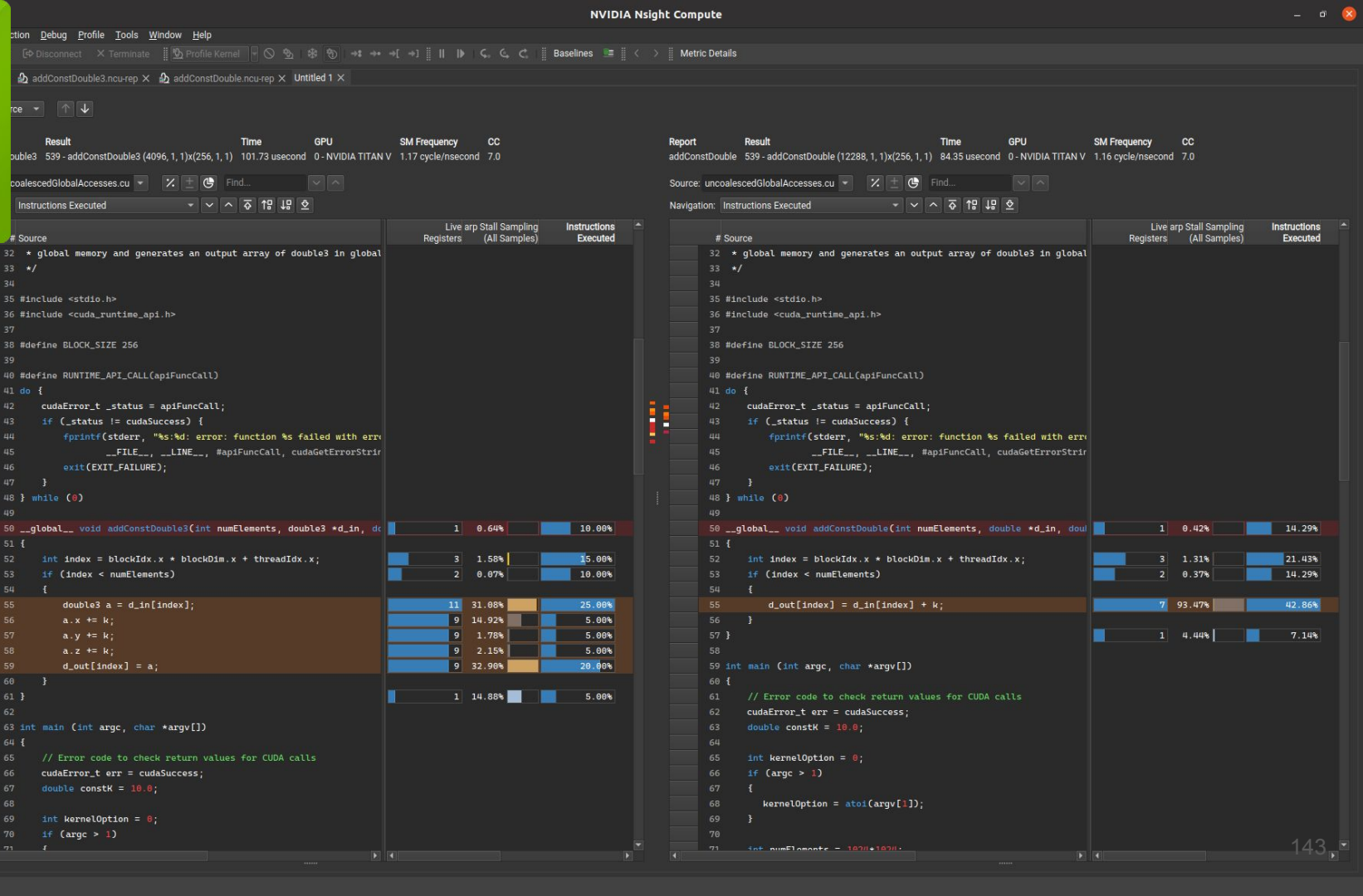






NVIDIA Nsight

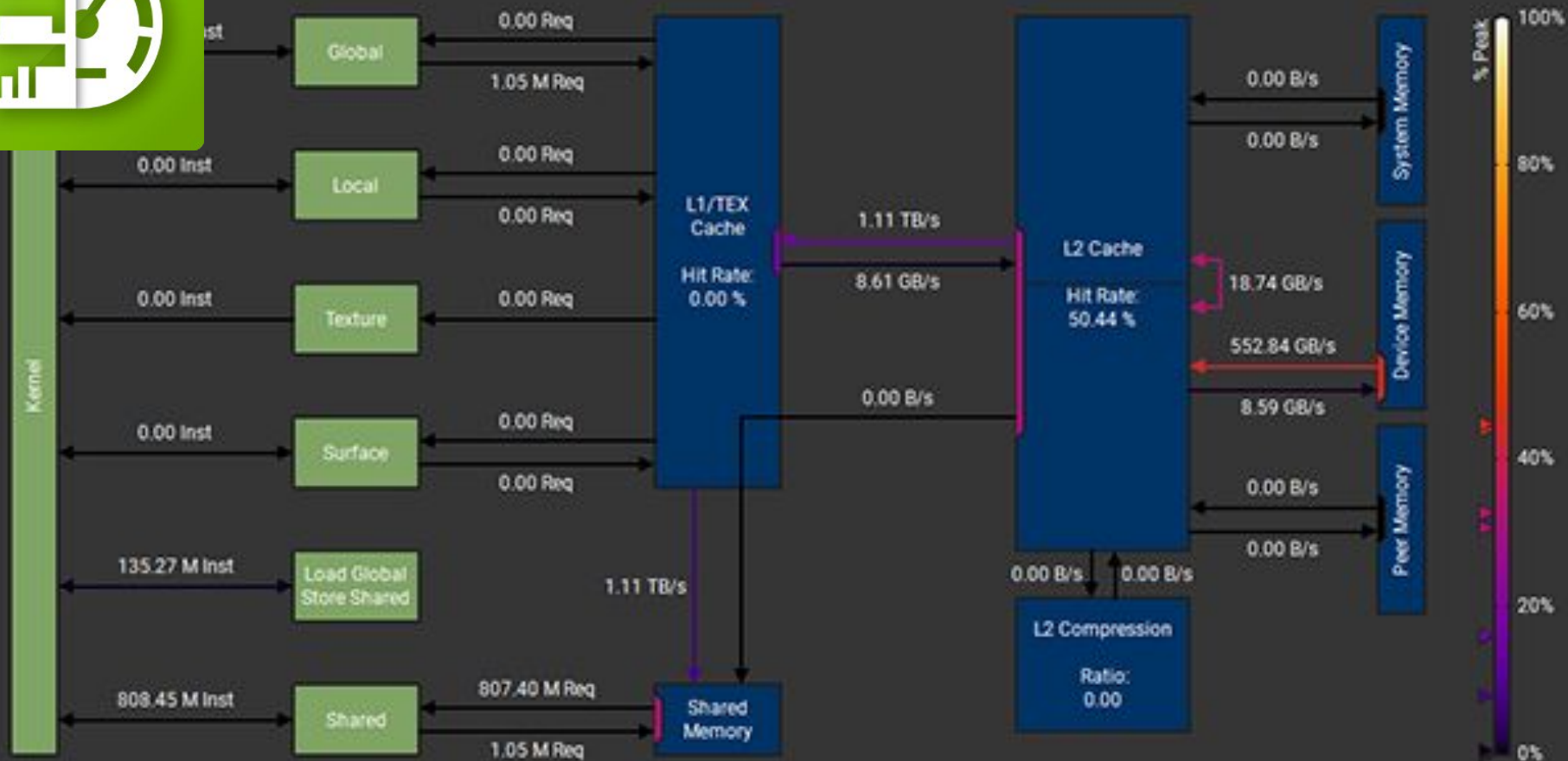






Memory Chart

Show As: **Throughput**





2 - 566 - mergeRanksAndIndicesKernel



Clear Baselines

Apply Rules

Occupancy Calculator

Source Comparison

Copy as Image

	Time	Cycles	GPU	SM Frequency	Process	Attributes
mergeRanksAndIndicesKernel (64, 1, 1)x(256, 1, 1)	4.19 us	4,326	0 - NVIDIA RTX A4500	1.03 cycle/ns	[762349] CuMergeSort	
mergeSortSharedKernel (4096, 1, 1)x(512, 1, 1)	555.10 us	5,79,537	0 - NVIDIA RTX A4500	1.04 cycle/ns	[762349] CuMergeSort	
generateSampleRanksKernel (64, 1, 1)x(256, 1, 1)	26.72 us	27,984	0 - NVIDIA RTX A4500	1.05 cycle/ns	[762349] CuMergeSort	

GPU Speed Of Light Throughput

GPU Throughput Chart



High-level overview of the throughput for compute and memory resources of the GPU. For each unit, the throughput reports the achieved percentage of utilization with respect to the theoretical maximum. Breakdowns show the throughput for each individual sub-metric of Compute and Memory to clearly identify the highest contributor. High-level overview of the utilization for compute and memory resources of the GPU presented as a roofline chart.

Compute (SM) Throughput [%]	10.84 (-75.36%, z=-0.61)	Duration [us]	4.19 (-98.56%, z=-0.75)
Memory Throughput [%]	7.88 (-87.90%, z=-1.23)	Elapsed Cycles [cycle]	4,326 (-98.58%, z=-0.75)
L1/TEX Cache Throughput [%]	10.84 (-79.01%, z=-0.82)	SM Active Cycles [cycle]	2,192.52 (-99.26%, z=-0.74)
L2 Cache Throughput [%]	7.88 (-51.10%, z=-0.75)	SM Frequency [cycle/ns]	1.03 (-1.32%, z=-1.38)
DRAM Throughput [%]	5.53 (-83.26%, z=-1.08)	DRAM Frequency [cycle/ns]	7.55 (-0.53%, z=-1.39)

GPU Throughput



PM Sampling



Timeline view of PM metrics sampled periodically over the workload duration. Data is collected across multiple passes. Use this section to understand how workload behavior changes over its runtime.

Baselines

Difference Bars Green/Red

	Name	Launch	Report	Time	Cycles	Registers	GPU	SM Frequency	CC	Process
✓	mergeSort	560 - mergeSortSharedKernel (4096, ...	/home/prabhanshuc/Documents/me...	555.10 us	5,79,537	17	0 - NVIDIA RTX ...	1.04 cycle/ns	8.6	[762349] CuMergeSort
✓	generateSample	563 - generateSampleRanksKernel (6...	/home/prabhanshuc/Documents/me...	26.72 us	27,984	18	0 - NVIDIA RTX ...	1.05 cycle/ns	8.6	[762349] CuMergeSort
✓	mergeRanks	566 - mergeRanksAndIndicesKernel (...	/home/prabhanshuc/Documents/me...	4.19 us	4,326	16	0 - NVIDIA RTX ...	1.03 cycle/ns	8.6	[762349] CuMergeSort

Глава 6: Примеры

A+B, максимум по массиву, merge-sort

Пример: A + B

A

10	34	12	34	54	113	...	1
----	----	----	----	----	-----	-----	---

+

B

32	12	57	12	14	126	...	5
----	----	----	----	----	-----	-----	---

||

C

42	46	69	46	68	239	...	6
----	----	----	----	----	-----	-----	---

Пример: A + B

A

10	34	12	34	54	113	...	1
----	----	----	----	----	-----	-----	---

+

B

32	12	57	12	14	126	...	5
----	----	----	----	----	-----	-----	---

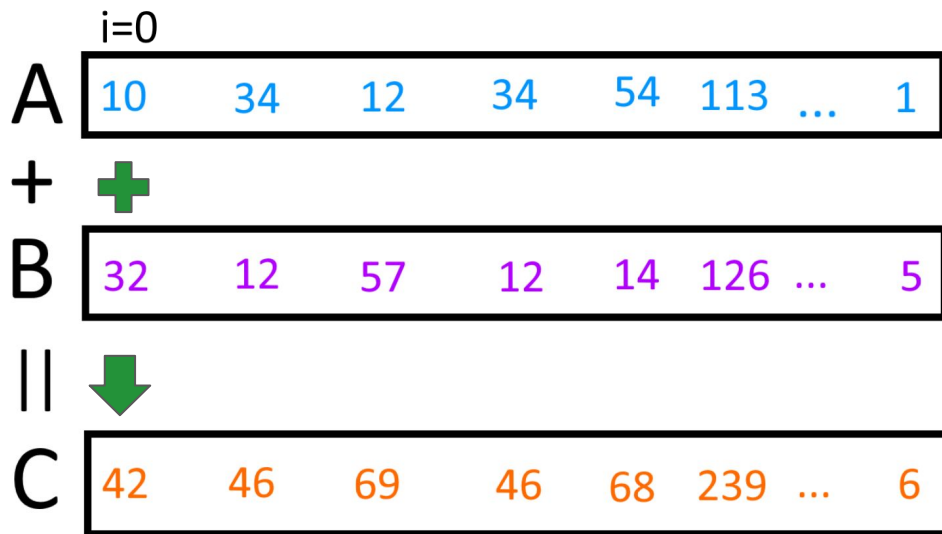
||

C

42	46	69	46	68	239	...	6
----	----	----	----	----	-----	-----	---

```
void solveCPU(int[] a, int[] b, int c[], int n) {  
    for (int i = 0; i < n; ++i) {  
        int sum = a[i] + b[i];  
        c[i] = sum;  
    }  
}
```

Пример: A + B



```
void solveCPU(int[] a, int[] b, int c[], int n) {  
    for (int i = 0; i < n; ++i) {  
        int sum = a[i] + b[i];  
        c[i] = sum;  
    }  
}
```

Пример: A + B

	i=0	i=1					
A	10	34	12	34	54	113 ...	1
+	+	+					
B	32	12	57	12	14	126 ...	5
	↓	↓					
C	42	46	69	46	68	239 ...	6

```
void solveCPU(int[] a, int[] b, int c[], int n) {  
    for (int i = 0; i < n; ++i) {  
        int sum = a[i] + b[i];  
        c[i] = sum;  
    }  
}
```

Пример: A + B

	i=0	i=1	i=2				
A	10	34	12	34	54	113 ...	1
+	+	+	+				
B	32	12	57	12	14	126 ...	5
	↓	↓	↓				
C	42	46	69	46	68	239 ...	6

```
void solveCPU(int[] a, int[] b, int c[], int n) {  
    for (int i = 0; i < n; ++i) {  
        int sum = a[i] + b[i];  
        c[i] = sum;  
    }  
}
```

Пример: A + B

	i=0	i=1	i=2	i=3	i=4	i=5	i=N-1
A	10	34	12	34	54	113 ...	1
+	+	+	+	+	+	+	+
B	32	12	57	12	14	126 ...	5
	↓	↓	↓	↓	↓	↓	↓
C	42	46	69	46	68	239 ...	6

```
void solveCPU(int[] a, int[] b, int c[], int n) {  
    for (int i = 0; i < n; ++i) {  
        int sum = a[i] + b[i];  
        c[i] = sum;  
    }  
}
```

**Какая
асимптотика?**

Пример: A + B

	i=0	i=1	i=2	i=3	i=4	i=5	i=N-1
A	10	34	12	34	54	113 ...	1
+	+	+	+	+	+	+	+
B	32	12	57	12	14	126 ...	5
	↓	↓	↓	↓	↓	↓	↓
C	42	46	69	46	68	239 ...	6

```
void solveCPU(int[] a, int[] b, int c[], int n) {  
    for (int i = 0; i < n; ++i) {  
        int sum = a[i] + b[i];  
        c[i] = sum;  
    }  
}
```

O(N)

Пример: A + B

A	10	34	12	34	54	113	...	1
+								
B	32	12	57	12	14	126	...	5
C	42	46	69	46	68	239	...	6



NVIDIA RTX 4090 - 16 тысяч ядер!



Пример: A + B

A	10	34	12	34	54	113	...	1
+								
B	32	12	57	12	14	126	...	5
C	42	46	69	46	68	239	...	6



```
void solveCPU(int[] a, int[] b, int c[], int n) {  
    for (int i = 0; i < n; ++i) {  
        int sum = a[i] + b[i];  
        c[i] = sum;  
    }  
}
```

O(N)

```
void solveGPU(int[] a, int[] b, int c[], int n) {  
    const int i = get_global_id(0);  
    c[i] = b[i] + a[i];  
}
```

Супер многопоточно!

NVIDIA RTX 4090 - 16 тысяч ядер!



Пример: A + B

A	10	34	12	34	54	113	...	1
+								
B	32	12	57	12	14	126	...	5
C	42	46	69	46	68	239	...	6



```
void solveCPU(int[] a, int[] b, int c[], int n) {  
    for (int i = 0; i < n; ++i) {  
        int sum = a[i] + b[i];  
        c[i] = sum;  
    }  
}
```

O(N)

```
void solveGPU(int[] a, int[] b, int c[], int n) {  
    const int i = get_global_id(0);  
    c[i] = b[i] + a[i];  
}
```

Какая асимптотика?

NVIDIA RTX 4090 - 16 тысяч ядер!



Пример: A + B

	i=0	i=1	i=2	i=3	i=4	i=5	i=N-1
A	10	34	12	34	54	113 ...	1
+	+	+	+	+	+	+	+
B	32	12	57	12	14	126 ...	5
	↓	↓	↓	↓	↓	↓	↓
C	42	46	69	46	68	239 ...	6

```
void solveCPU(int[] a, int[] b, int c[], int n) {  
    for (int i = 0; i < n; ++i) {  
        int sum = a[i] + b[i];  
        c[i] = sum;  
    }  
}
```

~~O(N)~~

```
void solveGPU(int[] a, int[] b, int c[], int n) {  
    const int i = get_global_id(0);  
    c[i] = b[i] + a[i];  
}
```

O(N / 16384)

NVIDIA RTX 4090 - 16 тысяч ядер!



OpenCL™

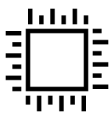
NVIDIA
CUDA®

Vulkan

Пример: A + B (N=100.000.000)

```
for (size_t i = 0; i < n; ++i) {  
    cs[i] = as[i] + bs[i];  
}
```

Пример: A + B (N=100.000.000)



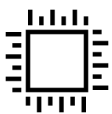
CPU - Intel 13700K

```
for (size_t i = 0; i < n; ++i) {  
    cs[i] = as[i] + bs[i];  
}
```

a + b median time: 0.068 sec (+-0.00481664)

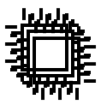
a + b median RAM bandwidth: 16.4351 GB/s

Пример: $A + B$ ($N=100.000.000$)



CPU - Intel 13700K

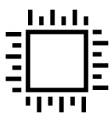
```
for (size_t i = 0; i < n; ++i) {  
    cs[i] = as[i] + bs[i];  
}  
  
a + b median time: 0.068 sec (+-0.00481664)  
a + b median RAM bandwidth: 16.4351 GB/s
```



CPU - Intel 13700K

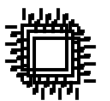
```
#pragma omp parallel for  
for (ptrdiff_t i = 0; i < n; ++i) {  
    cs[i] = as[i] + bs[i];  
}  
  
a + b median time: 0.047 sec (+-0.00153623)  
a + b median RAM bandwidth: 23.7784 GB/s
```

Пример: $A + B$ ($N=100.000.000$)



CPU - Intel 13700K

```
for (size_t i = 0; i < n; ++i) {  
    cs[i] = as[i] + bs[i];  
}  
  
a + b median time: 0.068 sec (+-0.00481664)  
a + b median RAM bandwidth: 16.4351 GB/s
```



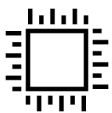
CPU - Intel 13700K

```
#pragma omp parallel for  
for (ptrdiff_t i = 0; i < n; ++i) {  
    cs[i] = as[i] + bs[i];  
}  
  
a + b median time: 0.047 sec (+-0.00153623)  
a + b median RAM bandwidth: 23.7784 GB/s
```

**Почему разница такая
незначительная?**

**Как проверить гипотезу?
Как спровоцировать
значительную разницу?**

Пример: **100*cos**(A + B) (N=100.000.000)

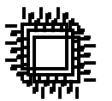


CPU - Intel 13700K

```
for (size_t i = 0; i < n; ++i) {  
    cs[i] = 100.0*cos( Left: as[i] + bs[i]);  
}
```

median time: 1.178 sec (+-0.0235009)

median RAM bandwidth: 0.948716 GB/s



CPU - Intel 13700K

```
#pragma omp parallel for
```

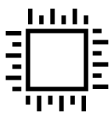
```
for (ptrdiff_t i = 0; i < n; ++i) {  
    cs[i] = 100*cos( Left: as[i] + bs[i]);  
}
```

median time: 0.109 sec (+-0.0121988)

median RAM bandwidth: 10.2531 GB/s

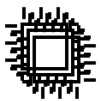
x10.9

Пример: A + B (N=100.000.000)



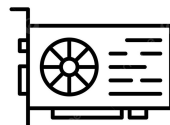
CPU - Intel 13700K

```
for (size_t i = 0; i < n; ++i) {  
    cs[i] = as[i] + bs[i];  
}  
a + b median time: 0.068 sec (+-0.00481664)  
a + b median RAM bandwidth: 16.4351 GB/s
```



CPU - Intel 13700K

```
#pragma omp parallel for  
for (ptrdiff_t i = 0; i < n; ++i) {  
    cs[i] = as[i] + bs[i];  
}  
a + b median time: 0.047 sec (+-0.00153623)  
a + b median RAM bandwidth: 23.7784 GB/s
```

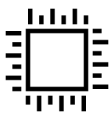


GPU - NVIDIA RTX 4090

```
__kernel void aplusb(__global const float* a,  
                    __global const float* b,  
                    __global float* c,  
                    unsigned int n)  
{  
    const unsigned int index = get_global_id( dimindx: 0 );  
    if (index >= n)  
        return;  
  
    c[index] = a[index] + b[index];  
}
```

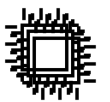
a + b median time: 0.002 sec (+-0.000632456)
a + b median VRAM bandwidth: 558.794 GB/s

Пример: A + B (N=100.000.000)



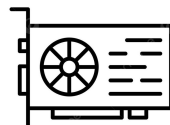
CPU - Intel 13700K

```
for (size_t i = 0; i < n; ++i) {  
    cs[i] = as[i] + bs[i];  
}  
  
a + b median time: 0.068 sec (+-0.00481664)  
a + b median RAM bandwidth: 16.4351 GB/s
```



CPU - Intel 13700K

```
#pragma omp parallel for  
for (ptrdiff_t i = 0; i < n; ++i) {  
    cs[i] = as[i] + bs[i];  
}  
  
a + b median time: 0.047 sec (+-0.00153623)  
a + b median RAM bandwidth: 23.7784 GB/s
```



GPU - NVIDIA RTX 4090

```
__kernel void aplusb(__global const float* a,  
                    __global const float* b,  
                    __global float* c,  
                    unsigned int n)  
{  
  
    const unsigned int index = get_global_id( dimindx: 0 );  
    if (index >= n)  
        return;  
  
    c[index] = a[index] + b[index];  
}  
  
a + b median time: 0.002 sec (+-0.000632456)  
a + b median VRAM bandwidth: 558.794 GB/s
```

x23.5

Пример: A + B (N=100.000.000)

```
1  #ifdef __CLION_IDE__
2  #include <libgpu/opencl/cl/clion_defines.cl>
3  #endif
4
5  #line 5
6
7
8
9
10 __attribute__((reqd_work_group_size(256, 1, 1)))
11
12 __kernel void aplusb(__global const float* a,
13                     __global const float* b,
14                     __global float* c,
15                     unsigned int n)
16 {
17     const unsigned int index = get_global_id(dimindx: 0);
18     if (index >= n)
19         return;
20
21     c[index] = a[index] + b[index];
22 }
```



Пример: A + B (N=100.000.000)

```
1 #ifdef __CLION_IDE__
2 #include <libgpu/opencl/cl/clion_defines.cl>
3 #endif
```

```
5 #line 5
```



```
10 __attribute__((reqd_work_group_size(256, 1, 1)))
11
12 __kernel void aplusb(__global const float* a,
13                     __global const float* b,
14                     __global float* c,
15                     unsigned int n)
16 {
17     const unsigned int index = get_global_id(dimindx: 0);
18     if (index >= n)
19         return;
20
21     c[index] = a[index] + b[index];
22 }
```



```
blockIdx.x * blockDim.x + threadIdx.x;
```

```
1 __global__ void aplusb(const float* a,
2                       const float* b,
3                       float* c,
4                       unsigned int n)
5 {
6     const unsigned int index = blockIdx.x * blockDim.x + threadIdx.x;
7     if (index >= n)
8         return;
9
10    c[index] = a[index] + b[index];
11 }
```

Пример: A + B (N=100.000.000)



```
1  __global__ void aplusb(const float* a,  
2                          const float* b,  
3                          float* c,  
4                          unsigned int n)  
5  {  
6      const unsigned int index = blockIdx.x * blockDim.x + threadIdx.x;  
7      if (index >= n)  
8          return;  
9  
10     c[index] = a[index] + b[index];  
11 }
```


Пример: A + B (N=100.000.000)



```
12 void cuda_aplusb(const gpu::WorkSize &workSize,  
13                 const float* a, const float* b, float* c, unsigned int n,  
14                 cudaStream_t stream)  
15 {  
16     const unsigned int blockSize = 256;  
17     const unsigned int gridSize = (n + blockSize - 1) / blockSize;  
18     aplusb<<<gridSize, blockSize, 0, stream>>>(a, b, c, n);  
19     CUDA_CHECK_KERNEL(stream);  
20 }  
21
```

```
1 __global__ void aplusb(const float* a,  
2                       const float* b,  
3                       float* c,  
4                       unsigned int n)  
5 {  
6     const unsigned int index = blockIdx.x * blockDim.x + threadIdx.x;  
7     if (index >= n)  
8         return;  
9  
10    c[index] = a[index] + b[index];  
11 }
```

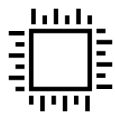
Пример: A + B (N=100.000.000)

GLSL (Graphics Library Shading Language)



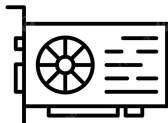
```
1 #version 450
2
3 #include <libgpu/vulkan/vk/common.vk>
4
5 layout (std430, binding = 0) readonly buffer AsIn { uint as[]; };
6 layout (std430, binding = 1) readonly buffer BsIn { uint bs[]; };
7 layout (std430, binding = 2) writeonly buffer CsOut { uint cs[]; };
8
9 layout (push_constant) uniform PushConstants {
10     uint n;
11 } params;
12
13 layout (local_size_x = 256) in;
14
15 void main()
16 {
17     const uint index = gl_GlobalInvocationID.x;
18     if (index >= params.n)
19         return;
20
21     cs[index] = as[index] + bs[index];
22 }
```

Пример: максимум по массиву



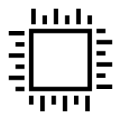
CPU Intel 13700K

```
max = -FLT_MAX;  
for (size_t i = 0; i < n; ++i) {  
    max = std::max(max, as[i]);  
}
```



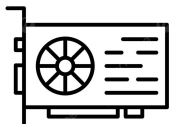
GPU NVIDIA RTX 4090

Пример: максимум по массиву



CPU Intel 13700K

```
max = -FLT_MAX;  
for (size_t i = 0; i < n; ++i) {  
    max = std::max(max, as[i]);  
}
```



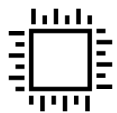
GPU NVIDIA RTX 4090

Как это реализовать в модели массового параллелизма?

Какой размер рабочего пространства?

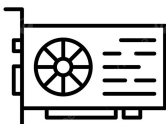
Что делает каждый поток (work item)?

Пример: максимум по массиву

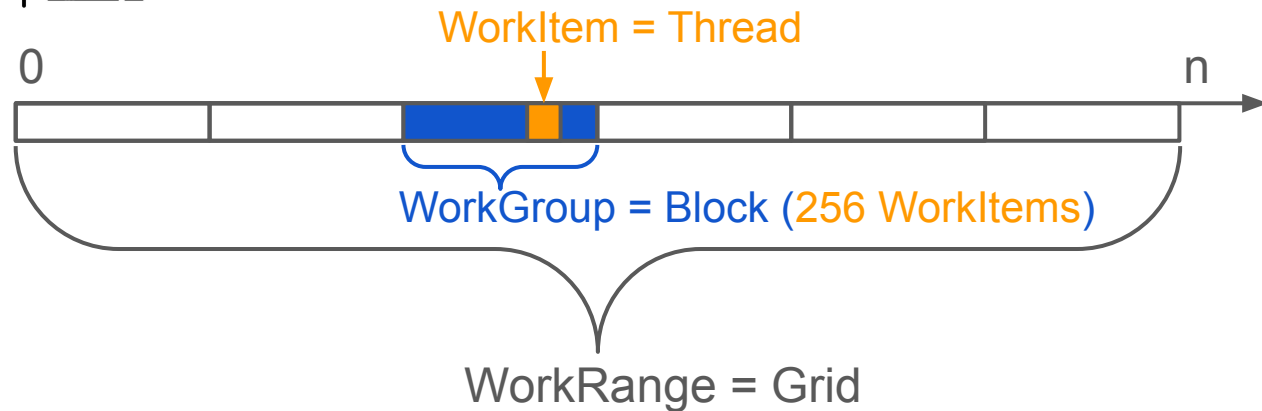


CPU Intel 13700K

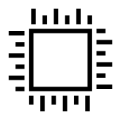
```
max = -FLT_MAX;  
for (size_t i = 0; i < n; ++i) {  
    max = std::max(max, as[i]);  
}
```



GPU NVIDIA RTX 4090

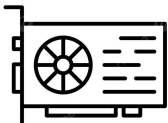


Пример: максимум по массиву

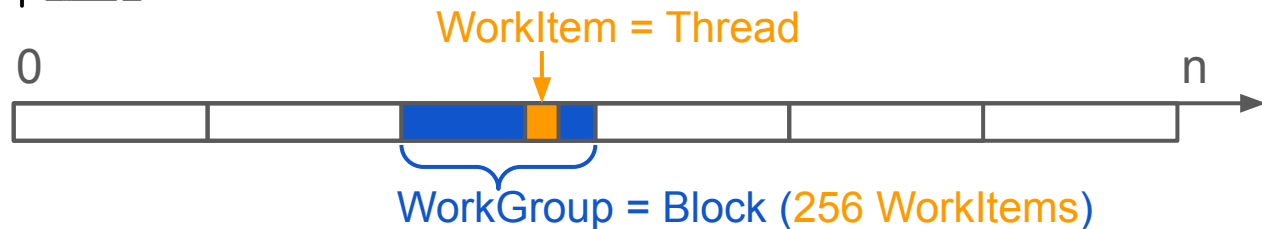


CPU Intel 13700K

```
max = -FLT_MAX;  
for (size_t i = 0; i < n; ++i) {  
    max = std::max(max, as[i]);  
}
```

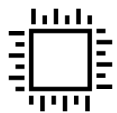


GPU NVIDIA RTX 4090



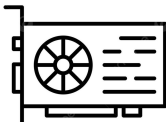
WorkGroup обречен считать как минимум 32
элемента подряд. **Почему?**

Пример: максимум по массиву

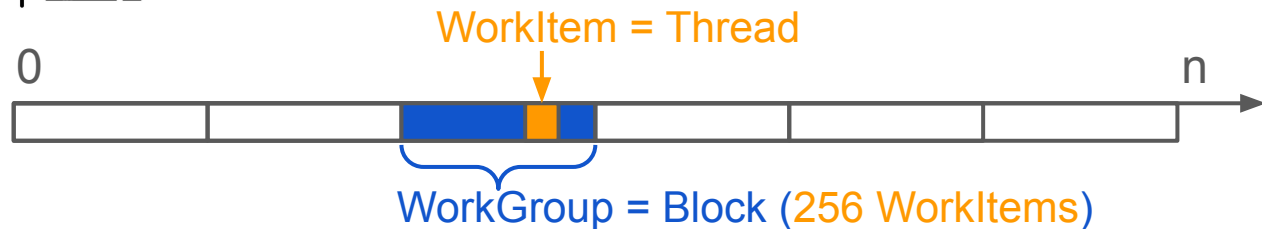


CPU Intel 13700K

```
max = -FLT_MAX;
for (size_t i = 0; i < n; ++i) {
    max = std::max(max, as[i]);
}
```



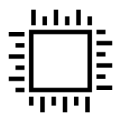
GPU NVIDIA RTX 4090



WorkGroup обречен считывать как минимум 32 элемента подряд ради **coalesced** access pattern.

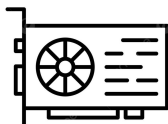
Что будем делать в **WorkGroup**?

Пример: максимум по массиву

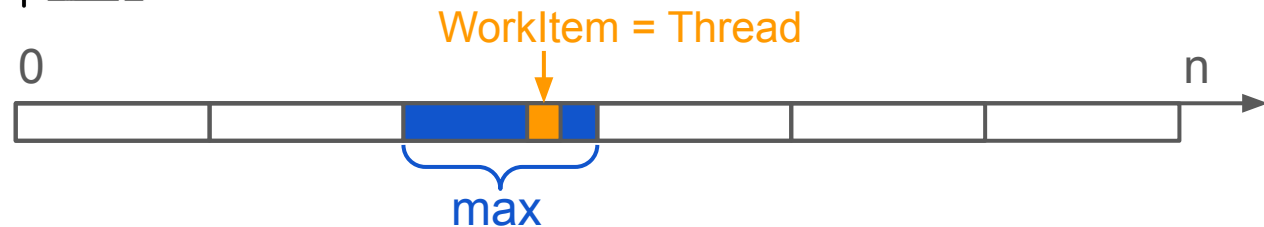


CPU Intel 13700K

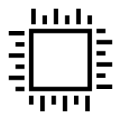
```
max = -FLT_MAX;  
for (size_t i = 0; i < n; ++i) {  
    max = std::max(max, as[i]);  
}
```



GPU NVIDIA RTX 4090

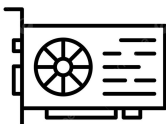


Пример: максимум по массиву

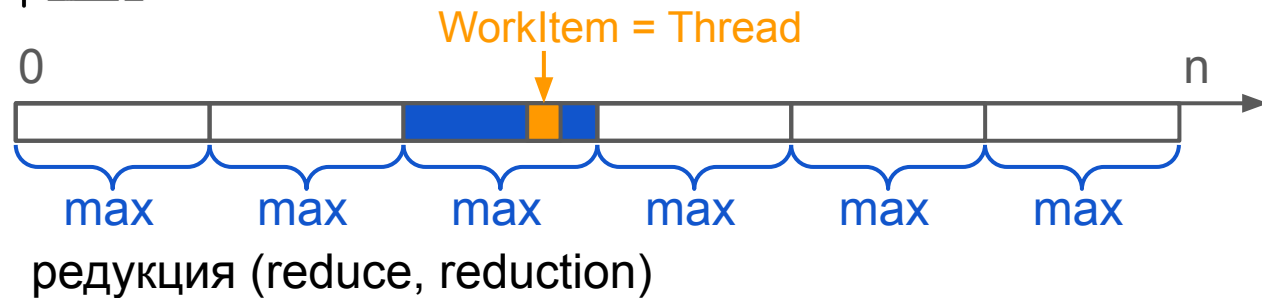


CPU Intel 13700K

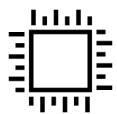
```
max = -FLT_MAX;  
for (size_t i = 0; i < n; ++i) {  
    max = std::max(max, as[i]);  
}
```



GPU NVIDIA RTX 4090

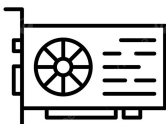


Пример: максимум по массиву

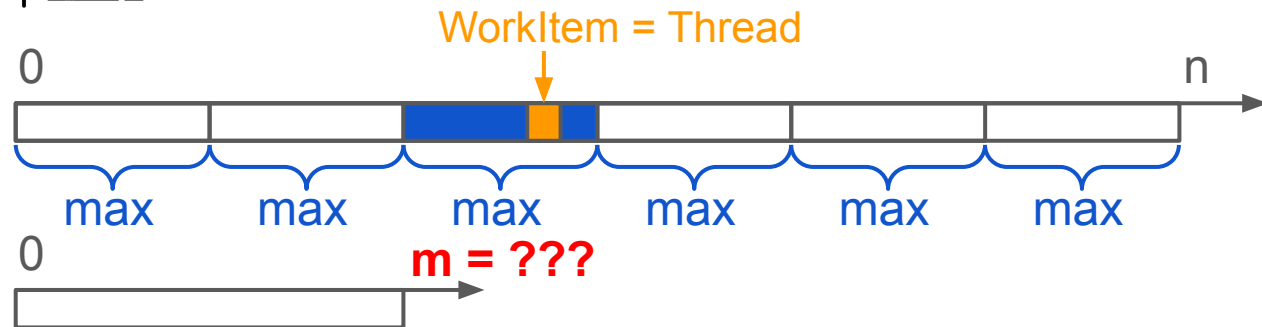


CPU Intel 13700K

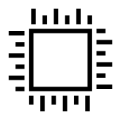
```
max = -FLT_MAX;  
for (size_t i = 0; i < n; ++i) {  
    max = std::max(max, as[i]);  
}
```



GPU NVIDIA RTX 4090

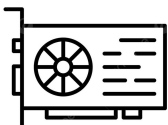


Пример: максимум по массиву

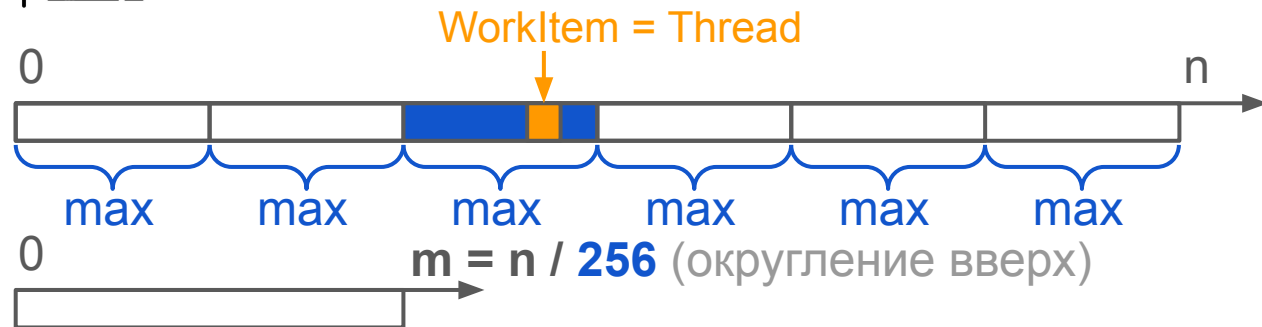


CPU Intel 13700K

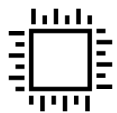
```
max = -FLT_MAX;  
for (size_t i = 0; i < n; ++i) {  
    max = std::max(max, as[i]);  
}
```



GPU NVIDIA RTX 4090

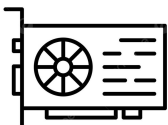


Пример: максимум по массиву

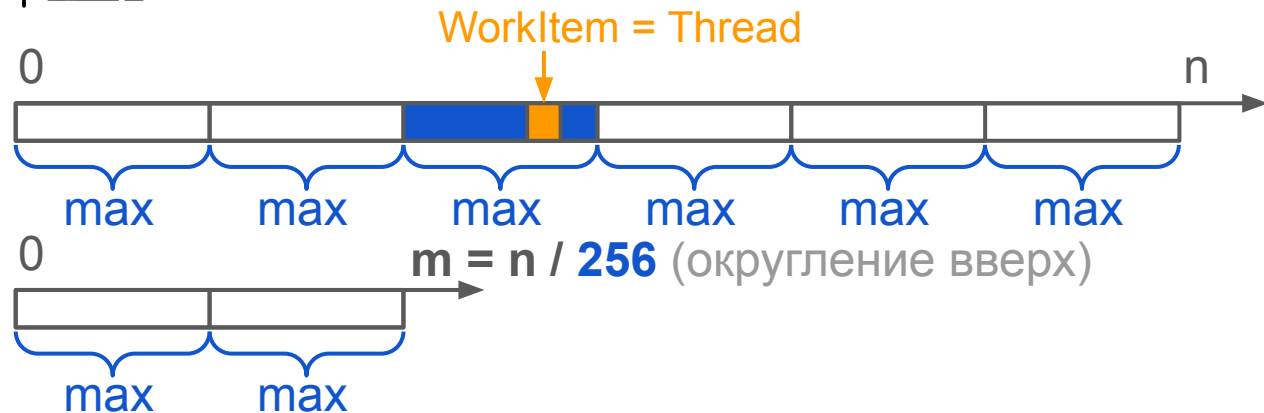


CPU Intel 13700K

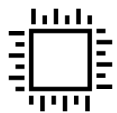
```
max = -FLT_MAX;  
for (size_t i = 0; i < n; ++i) {  
    max = std::max(max, as[i]);  
}
```



GPU NVIDIA RTX 4090

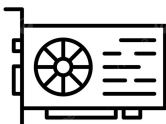


Пример: максимум по массиву

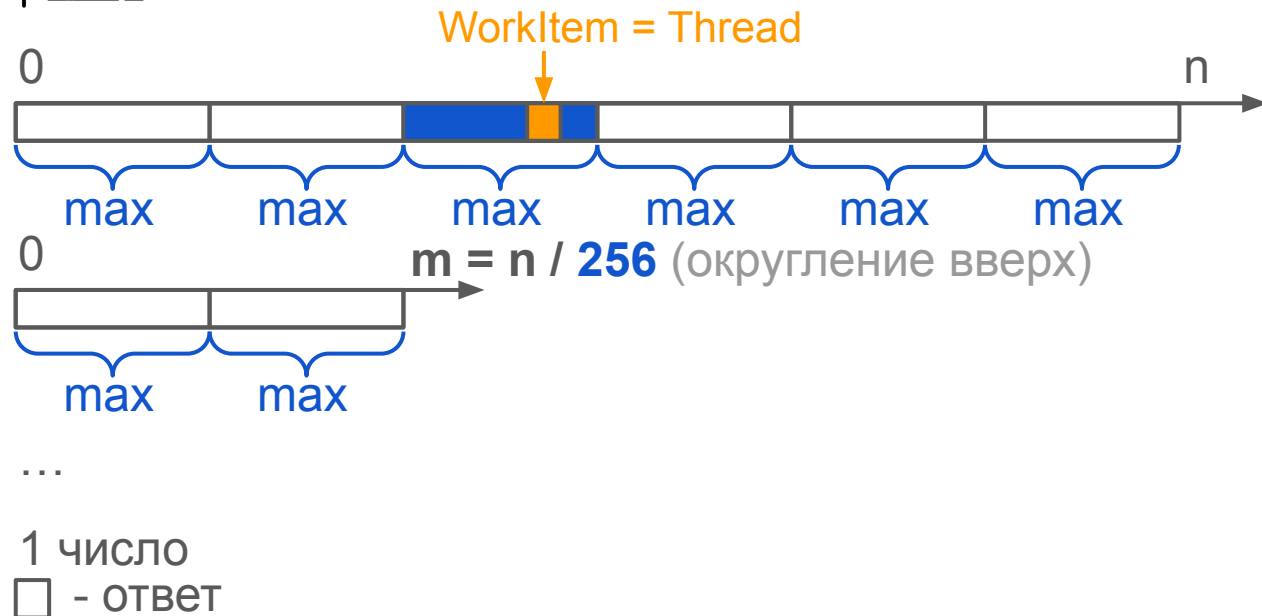


CPU Intel 13700K

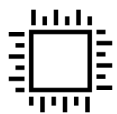
```
max = -FLT_MAX;
for (size_t i = 0; i < n; ++i) {
    max = std::max(max, as[i]);
}
```



GPU NVIDIA RTX 4090

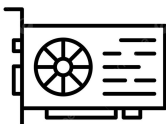


Пример: максимум по массиву

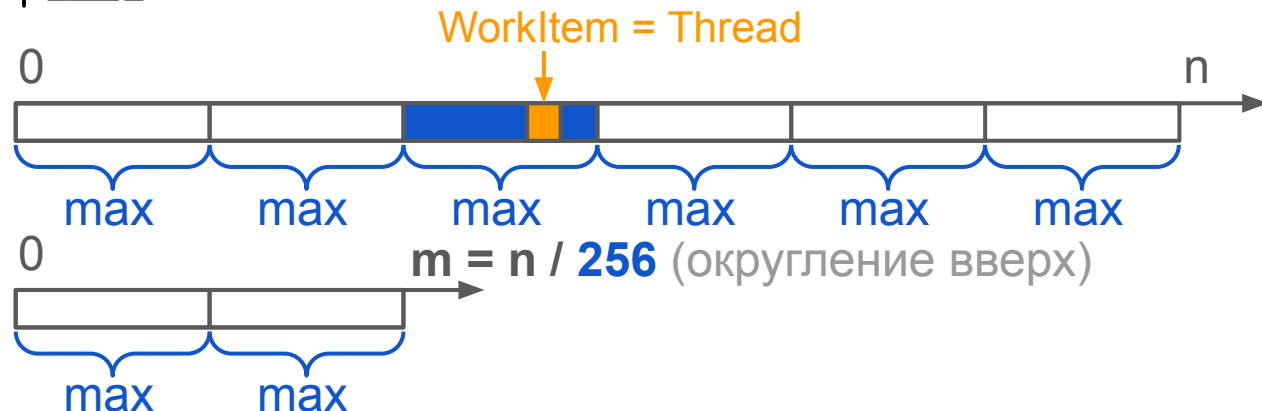


CPU Intel 13700K

```
max = -FLT_MAX;
for (size_t i = 0; i < n; ++i) {
    max = std::max(max, as[i]);
}
```



GPU NVIDIA RTX 4090



...

1 число

☐ - ответ

Все ли особенности мы обсудили?

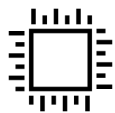
max median time: 0.195 sec (+-0.00200998)

max median RAM bandwidth: **1.91041 GB/s**

max median time: 0.003 sec (+-0.00067082)

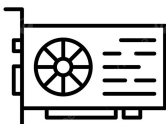
max median VRAM bandwidth: **124.176 GB/s**

Пример: максимум по массиву

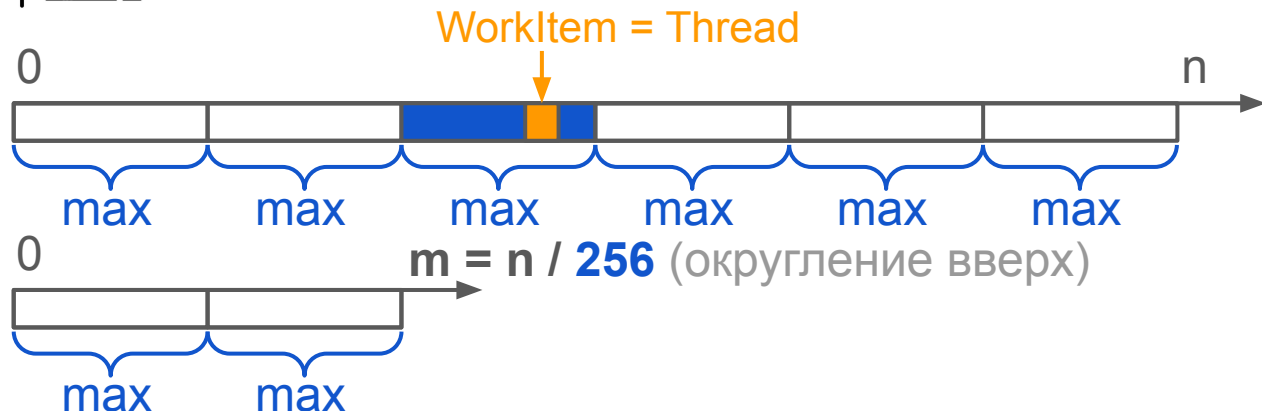


CPU Intel 13700K

```
max = -FLT_MAX;
for (size_t i = 0; i < n; ++i) {
    max = std::max(max, as[i]);
}
```



GPU NVIDIA RTX 4090



...

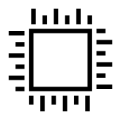
1 число
□ - ответ

- 1) Как реализовать редукцию?
- 2) Как найти локальный max?
- 3) Как синхронизировать этапы?

max median time: 0.195 sec (+-0.00200998)
max median RAM bandwidth: 1.91041 GB/s

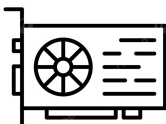
max median time: 0.003 sec (+-0.00067082)
max median VRAM bandwidth: 124.176 GB/s

Пример: максимум по массиву

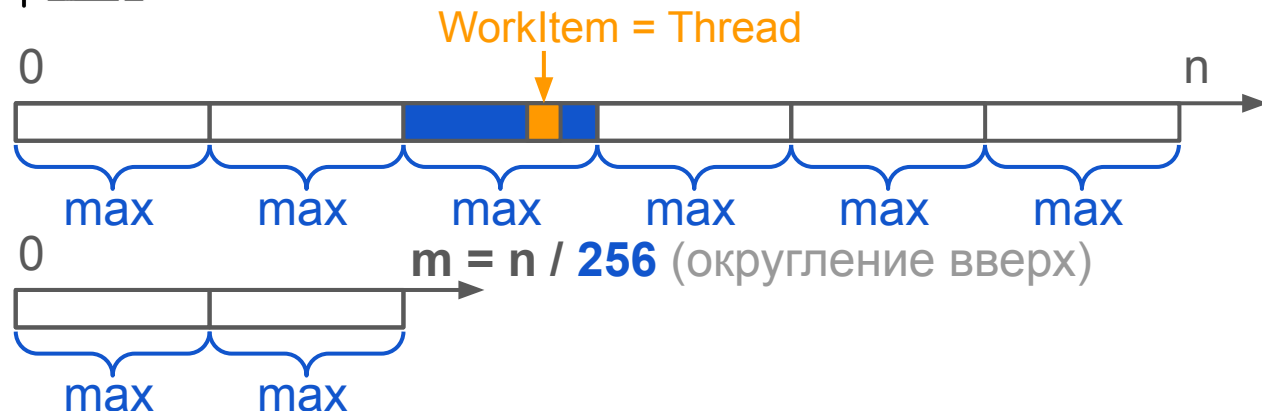


CPU Intel 13700K

```
max = -FLT_MAX;
for (size_t i = 0; i < n; ++i) {
    max = std::max(max, as[i]);
}
```



GPU NVIDIA RTX 4090



1) Грузим 256 элемент в local memory

1 число

□ - ОТВЕТ

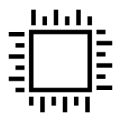
max median time: 0.195 sec (+-0.00200998)

max median RAM bandwidth: 1.91041 GB/s

max median time: 0.003 sec (+-0.00067082)

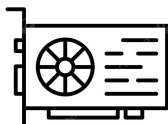
max median VRAM bandwidth: 124.176 GB/s

Пример: максимум по массиву

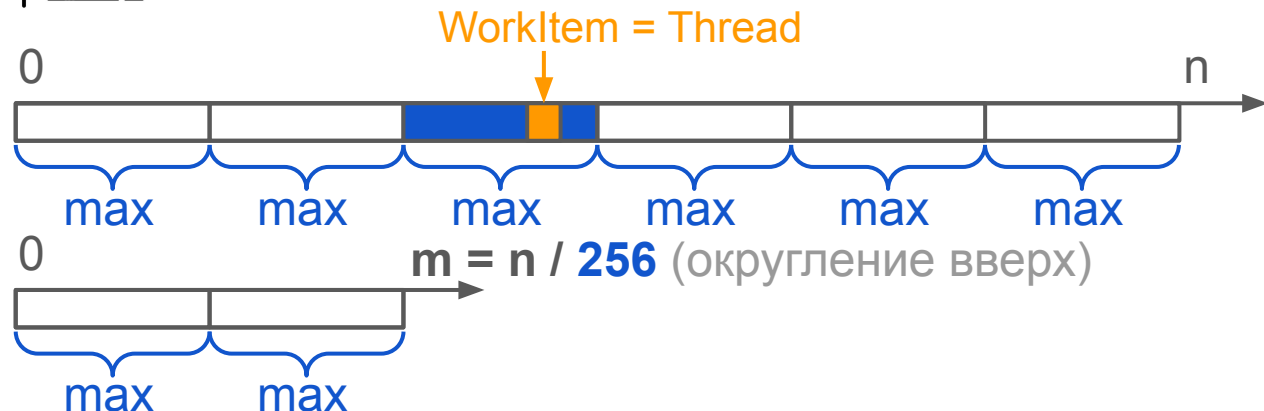


CPU Intel 13700K

```
max = -FLT_MAX;
for (size_t i = 0; i < n; ++i) {
    max = std::max(max, as[i]);
}
```



GPU NVIDIA RTX 4090



- 1) Грузим 256 элемент в **local memory**
- 2) Мастер-поток группы ищет **max**

1 число

□ - ответ

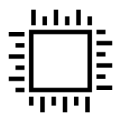
max median time: 0.195 sec (+-0.00200998)

max median RAM bandwidth: **1.91041 GB/s**

max median time: 0.003 sec (+-0.00067082)

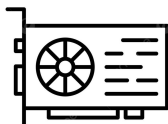
max median VRAM bandwidth: **124.176 GB/s**

Пример: максимум по массиву

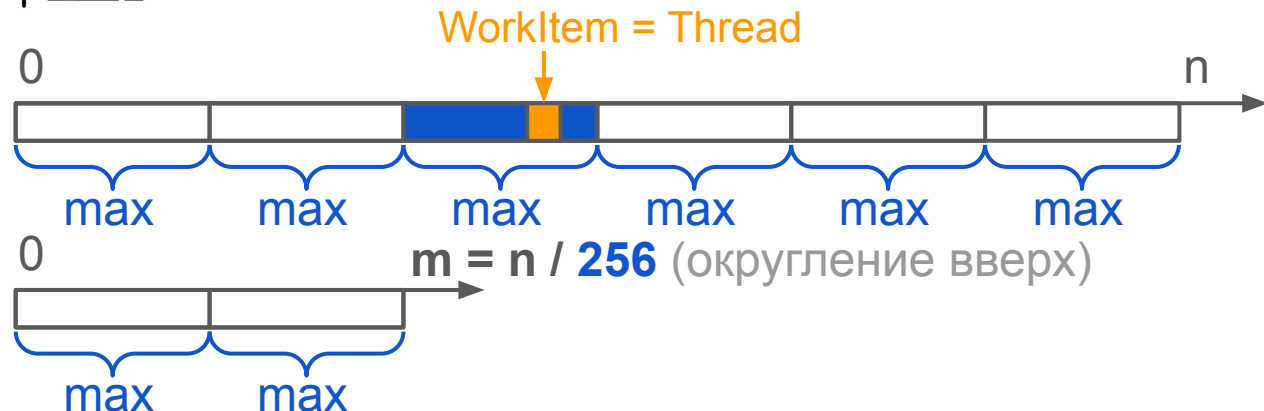


CPU Intel 13700K

```
max = -FLT_MAX;
for (size_t i = 0; i < n; ++i) {
    max = std::max(max, as[i]);
}
```



GPU NVIDIA RTX 4090



...

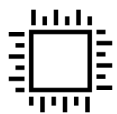
1 число
□ - ответ

- 1) Грузим 256 элемент в **local memory**
- 2) Мастер-поток группы ищет **max**
- 3) Мастер-поток группы пишет **max**

max median time: 0.195 sec (+-0.00200998)
max median RAM bandwidth: **1.91041 GB/s**

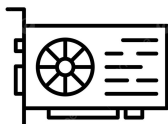
max median time: 0.003 sec (+-0.00067082)
max median VRAM bandwidth: **124.176 GB/s**

Пример: максимум по массиву

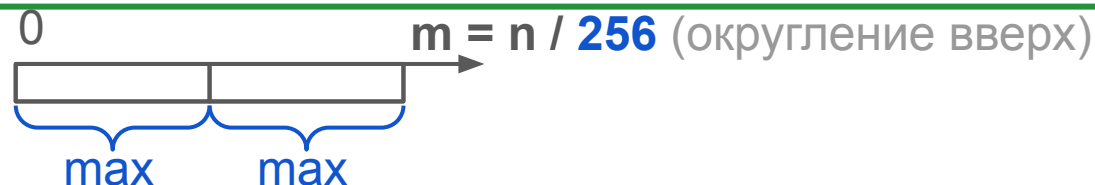
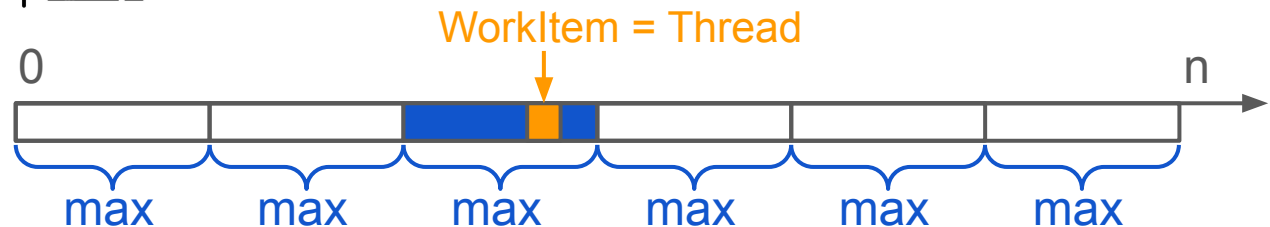


CPU Intel 13700K

```
max = -FLT_MAX;
for (size_t i = 0; i < n; ++i) {
    max = std::max(max, as[i]);
}
```



GPU NVIDIA RTX 4090



...

1 число
□ - ответ

- 1) Грузим 256 элемент в **local memory**
- 2) Мастер-поток группы ищет **max**
- 3) Мастер-поток группы пишет **max**
- 4) **Синхронизация** запусками kernel-ов

max median time: 0.195 sec (+-0.00200998)
max median RAM bandwidth: **1.91041 GB/s**

max median time: 0.003 sec (+-0.00067082)
max median VRAM bandwidth: **124.176 GB/s**

Пример: максимум по массиву



```
29 // Чтение данных
30 if (global_id < params.n) {
31     local_data[local_id] = a[global_id];
32 } else {
33     // Если размер входного массива не кратен рабочей группе
34     // то прощсе потокам за пределами массива подсунуть
35     // нейтральный элемент (с точки зрения поиска максимума)
36     local_data[local_id] = -FLT_MAX;
37 }
```

1) Грузим 256 элемент в **local memory**

Пример: максимум по массиву



```
29 // Чтение данных
30 if (global_id < params.n) {
31     local_data[local_id] = a[global_id];
32 } else {
33     // Если размер входного массива не кратен рабочей группе
34     // то проще потокам за пределами массива подсунуть
35     // нейтральный элемент (с точки зрения поиска максимума)
36     local_data[local_id] = -FLT_MAX;
37 }
```

```
42 // Мастер поток ищет максимум и записывает в выход
43 if (local_id == 0) {
44     float max_val = local_data[0];
45     for (uint i = 1; i < group_size; ++i) {
46         if (local_data[i] > max_val) {
47             max_val = local_data[i];
48         }
49     }
```

- 1) Грузим 256 элемент в **local memory**
- 2) Мастер-поток группы ищет **max**

Пример: максимум по массиву



```
29 // Чтение данных
30 if (global_id < params.n) {
31     local_data[local_id] = a[global_id];
32 } else {
33     // Если размер входного массива не кратен рабочей группе
34     // то проще потокам за пределами массива подсунуть
35     // нейтральный элемент (с точки зрения поиска максимума)
36     local_data[local_id] = -FLT_MAX;
37 }
```

```
42 // Мастер поток ищет максимум и записывает в выход
43 if (local_id == 0) {
44     float max_val = local_data[0];
45     for (uint i = 1; i < group_size; ++i) {
46         if (local_data[i] > max_val) {
47             max_val = local_data[i];
48         }
49     }
50     a_reduced[group_id] = max_val;
51 }
```

- 1) Грузим 256 элемент в **local memory**
- 2) Мастер-поток группы ищет **max**
- 3) Мастер-поток группы пишет **max**

Пример: максимум по массиву



```
29 // Чтение данных
30 if (global_id < params.n) {
31     local_data[local_id] = a[global_id];
32 } else {
33     // Если размер входного массива не кратен рабочей группе
34     // то проще потокам за пределами массива подсунуть
35     // нейтральный элемент (с точки зрения поиска максимума)
36     local_data[local_id] = -FLT_MAX;
37 }
```

```
42 // Мастер поток ищет максимум и записывает в выход
43 if (local_id == 0) {
44     float max_val = local_data[0];
45     for (uint i = 1; i < group_size; ++i) {
46         if (local_data[i] > max_val) {
47             max_val = local_data[i];
48         }
49     }
50     a_reduced[group_id] = max_val;
51 }
```

- 1) Грузим 256 элемент в **local memory**
- 2) Мастер-поток группы ищет **max**
- 3) Мастер-поток группы пишет **max**

Нет ли баги?



Пример: максимум по массиву



```
29 // Чтение данных
30 if (global_id < params.n) {
31     local_data[local_id] = a[global_id];
32 } else {
33     // Если размер входного массива не кратен рабочей группе
34     // то проще потокам за пределами массива подсунуть
35     // нейтральный элемент (с точки зрения поиска максимума)
36     local_data[local_id] = -FLT_MAX;
37 }
```

```
38
39 // Барьер для синхронизации всей рабочей группы
40 barrier();
```

```
41
42 // Мастер поток ищет максимум и записывает в выход
```

```
43 if (local_id == 0) {
44     float max_val = local_data[0];
45     for (uint i = 1; i < group_size; ++i) {
46         if (local_data[i] > max_val) {
47             max_val = local_data[i];
48         }
49     }
50     a_reduced[group_id] = max_val;
51 }
```

- 1) Грузим 256 элемент в **local memory**
- 2) Мастер-поток группы ищет **max**
- 3) Мастер-поток группы пишет **max**

Пример: максимум по массиву



```
29 // Чтение данных
30 if (global_id < params.n) {
31     local_data[local_id] = a[global_id];
32 } else {
33     // Если размер входного массива не кратен рабочей группе
34     // то проще потокам за пределами массива подсунуть
35     // нейтральный элемент (с точки зрения поиска максимума)
36     local_data[local_id] = -FLT_MAX;
37 }
```

```
38
39 // Барьер для синхронизации всей рабочей группы
40 barrier();
41
```

```
42 // Мастер поток ищет максимум и записывает в выход
43
```

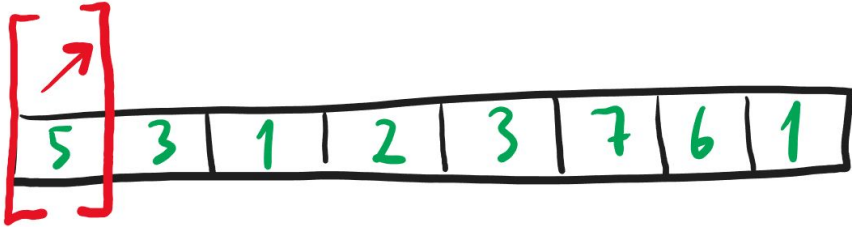
```
44 if (local_id == 0) {
45     float max_val = local_data[0];
46     for (uint i = 1; i < group_size; ++i) {
47         if (local_data[i] > max_val) {
48             max_val = local_data[i];
49         }
50     }
51     a_reduced[group_id] = max_val;
52 }
```

- 1) Грузим 256 элемент в **local memory**
- 2) Мастер-поток группы ищет **max**
- 3) Мастер-поток группы пишет **max**
- 4) **Синхронизация** запусками kernel-ов

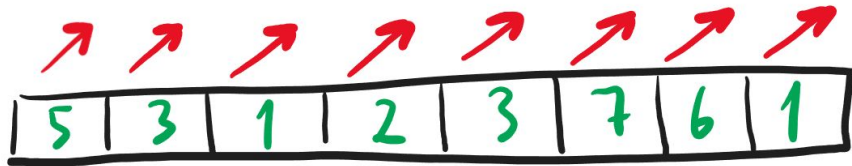
Merge-sort

5	3	1	2	3	7	6	1
---	---	---	---	---	---	---	---

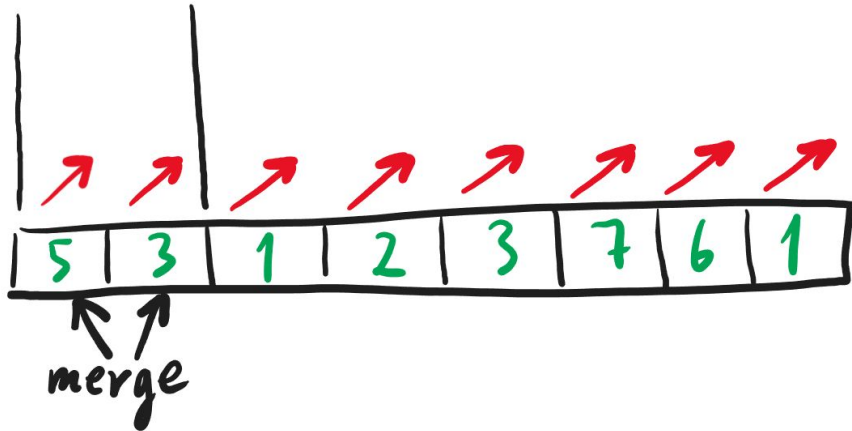
Merge-sort



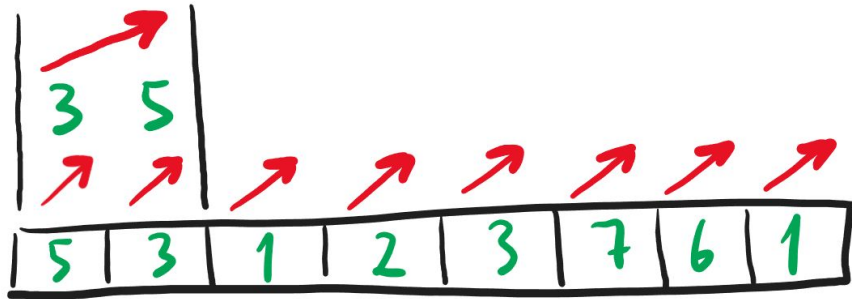
Merge-sort



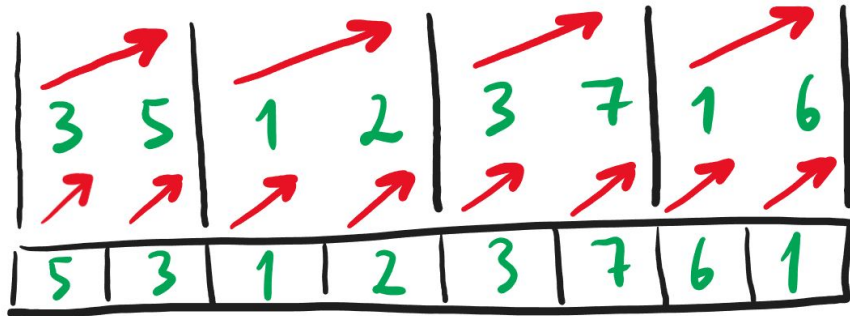
Merge-sort



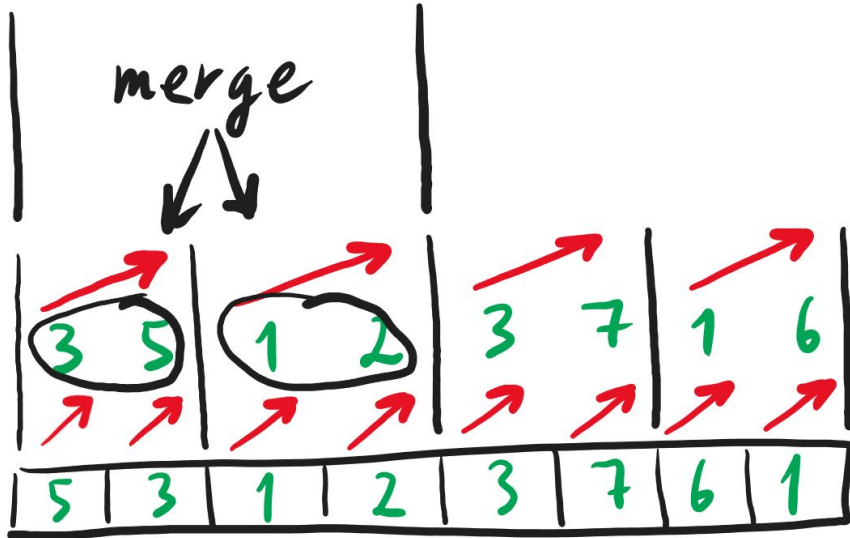
Merge-sort



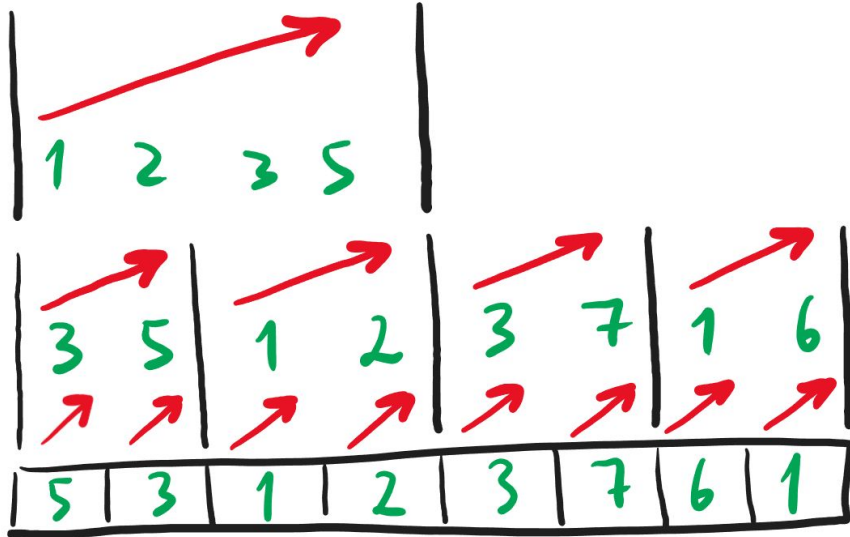
Merge-sort



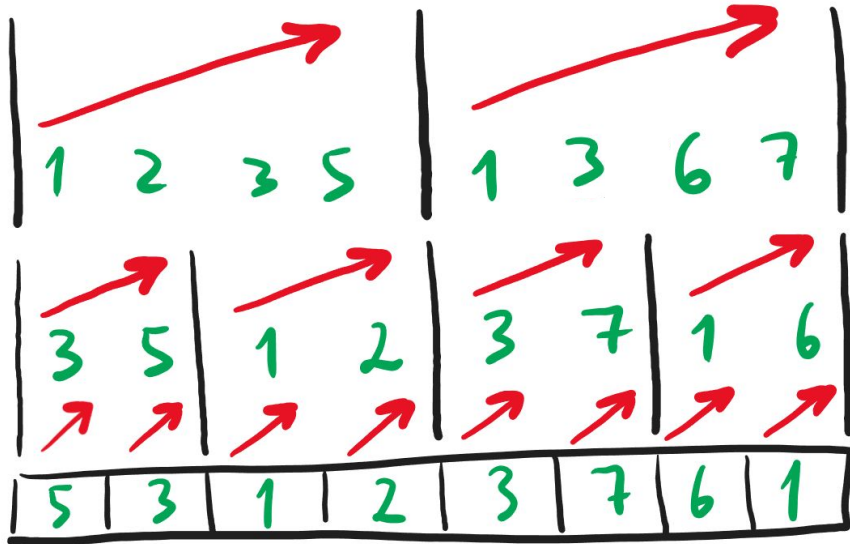
Merge-sort



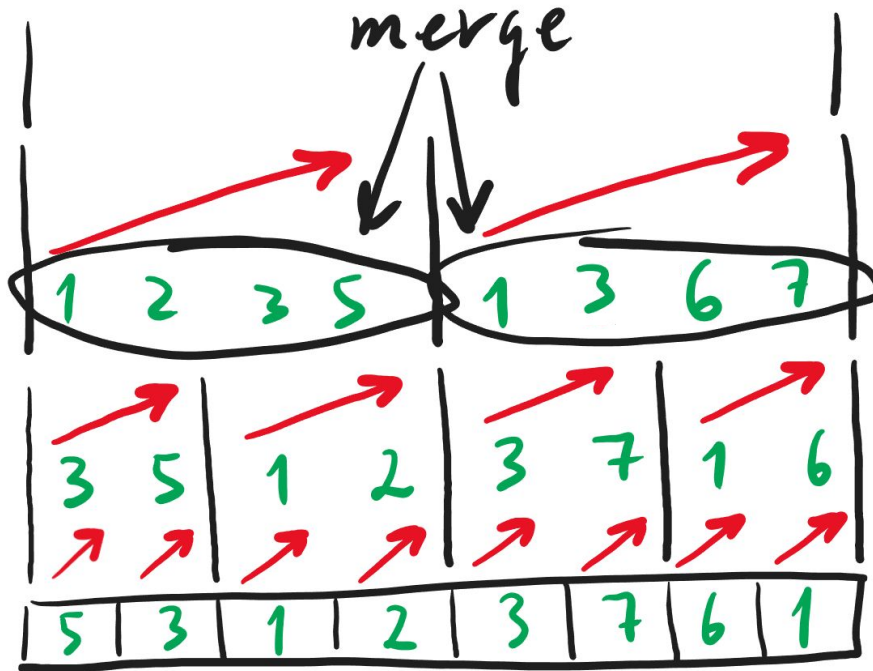
Merge-sort



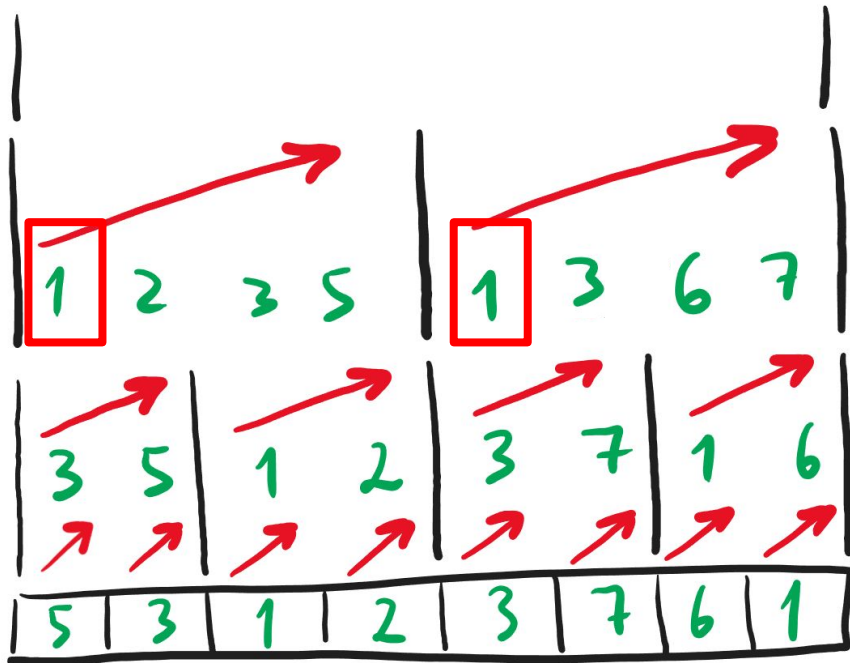
Merge-sort



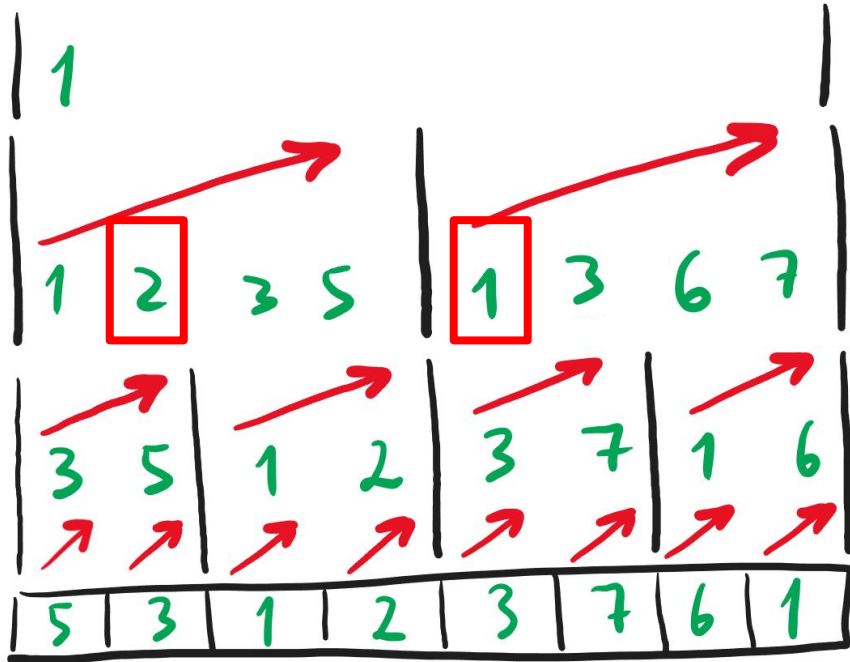
Merge-sort



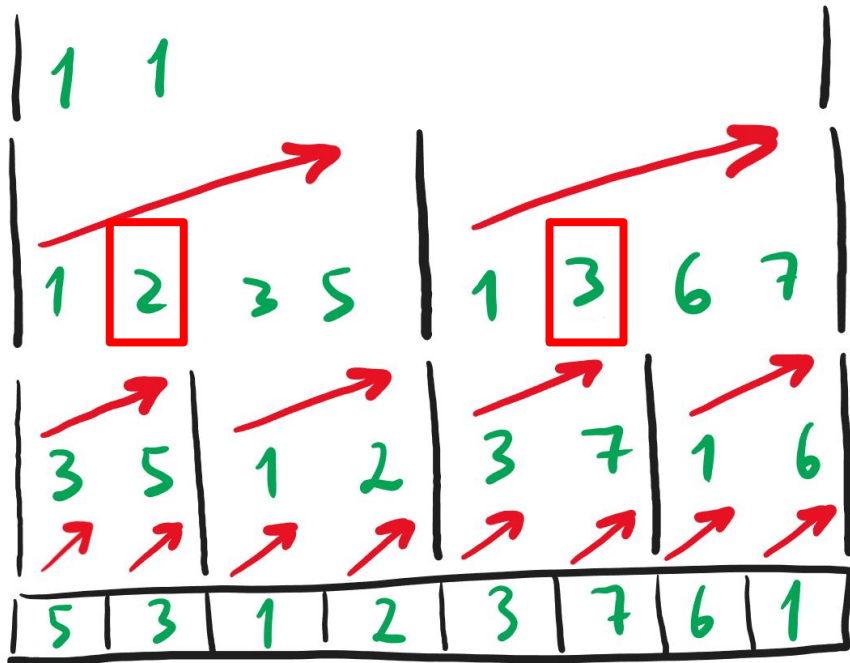
Merge-sort



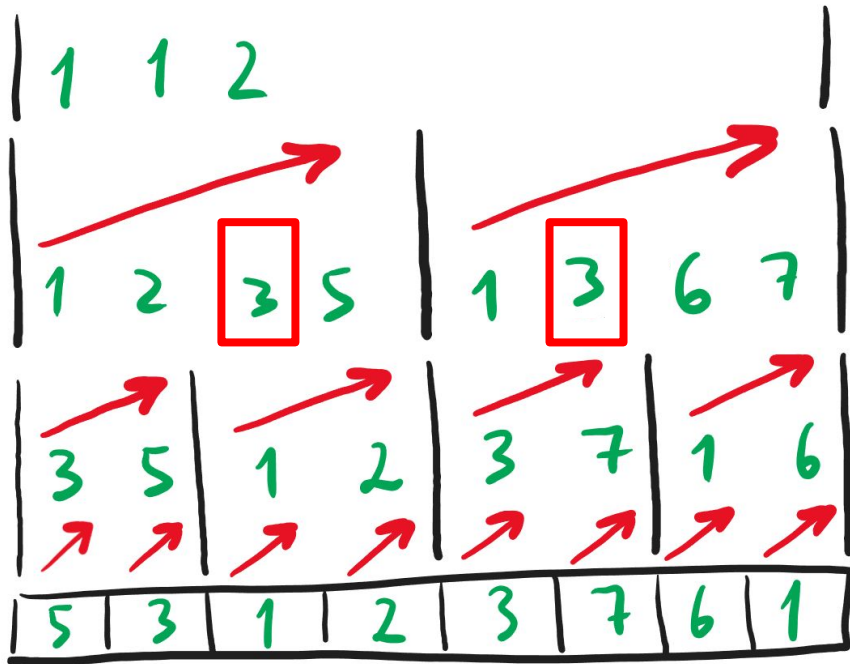
Merge-sort



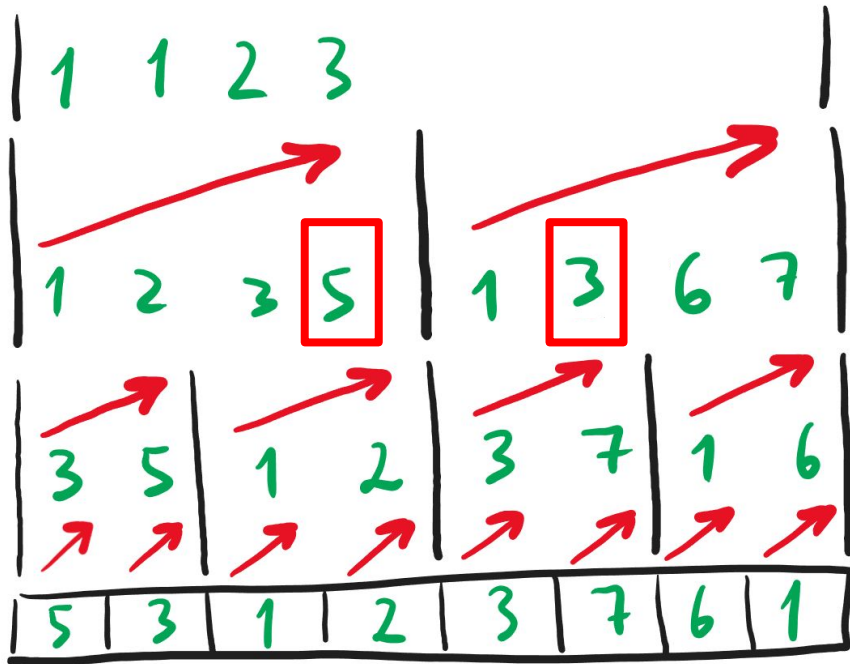
Merge-sort



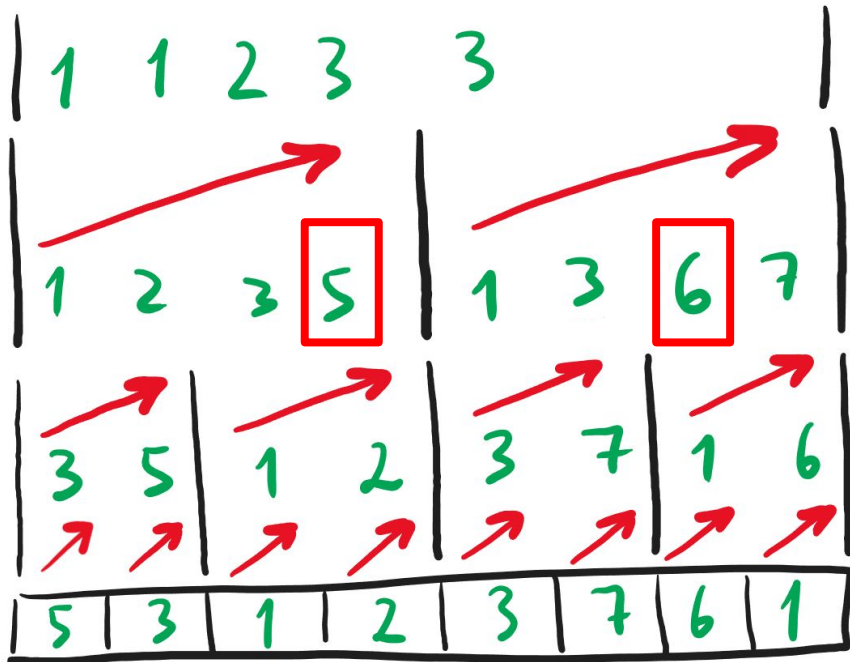
Merge-sort



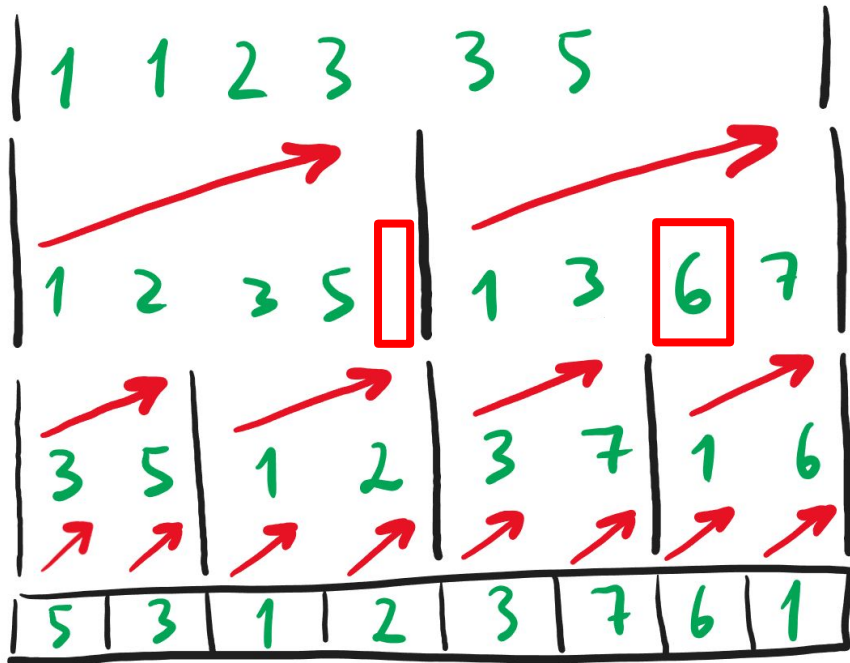
Merge-sort



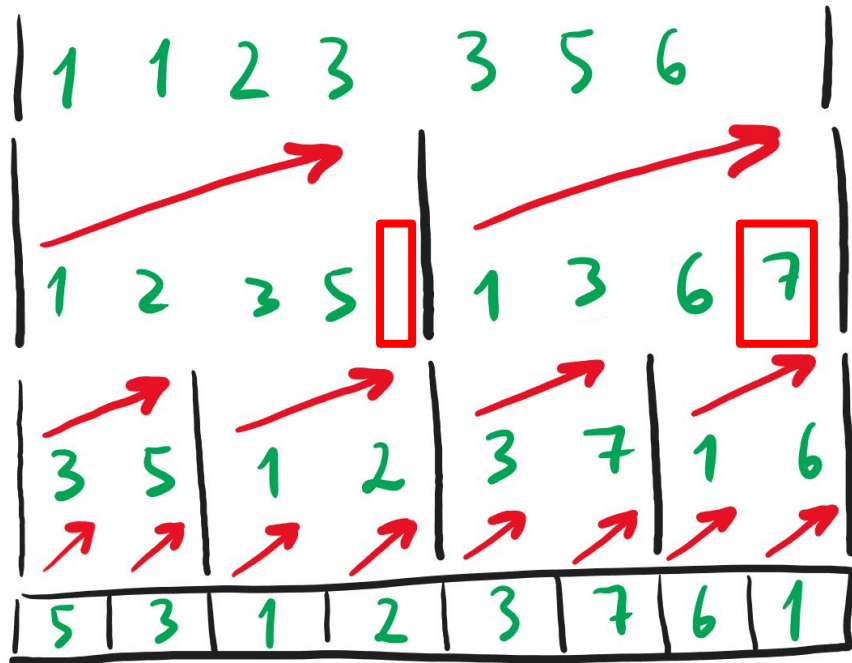
Merge-sort



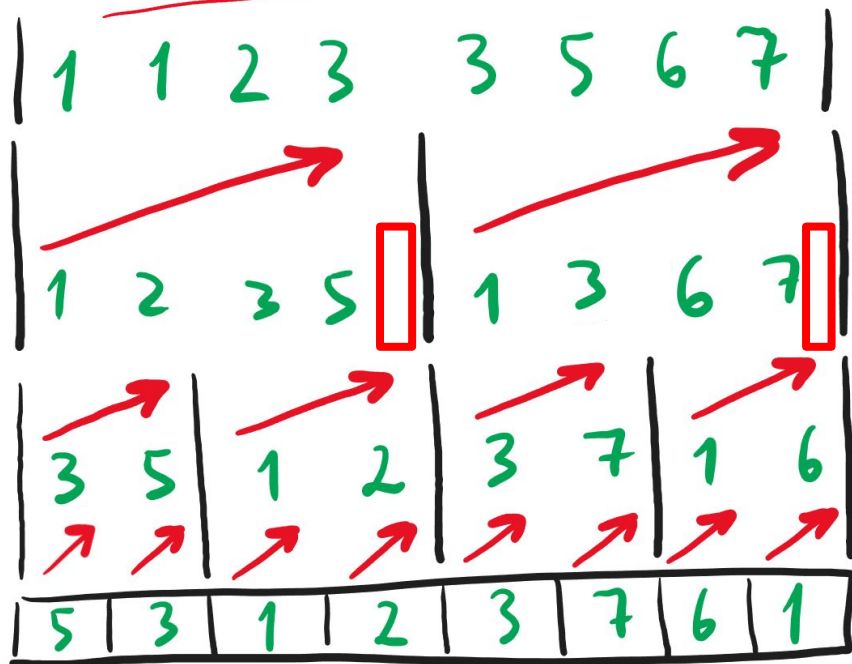
Merge-sort



Merge-sort

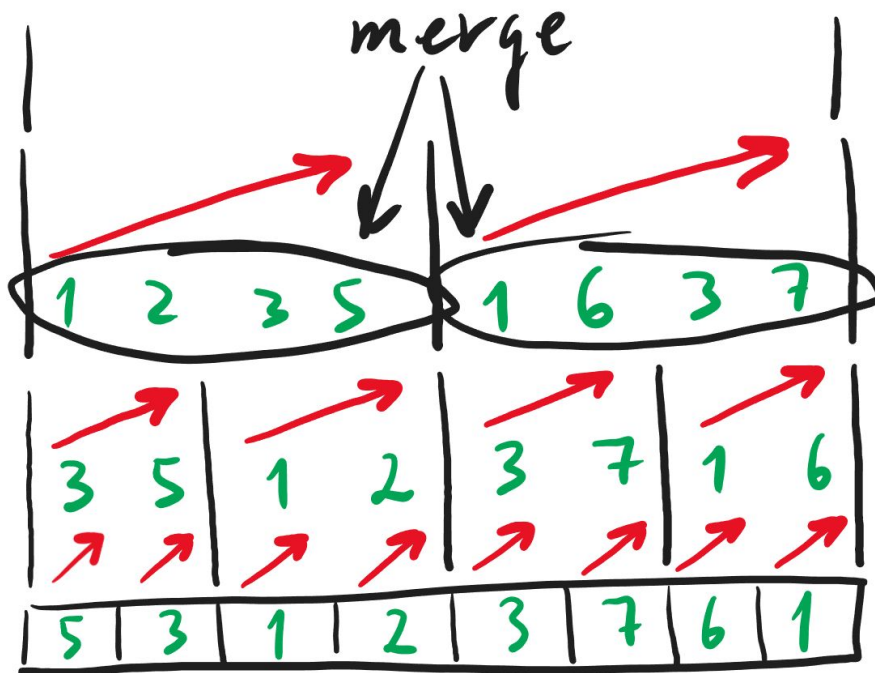


Merge-sort



Как отсортировать на видеокарте?
Т.е. как отсортировать супер
многопоточно?

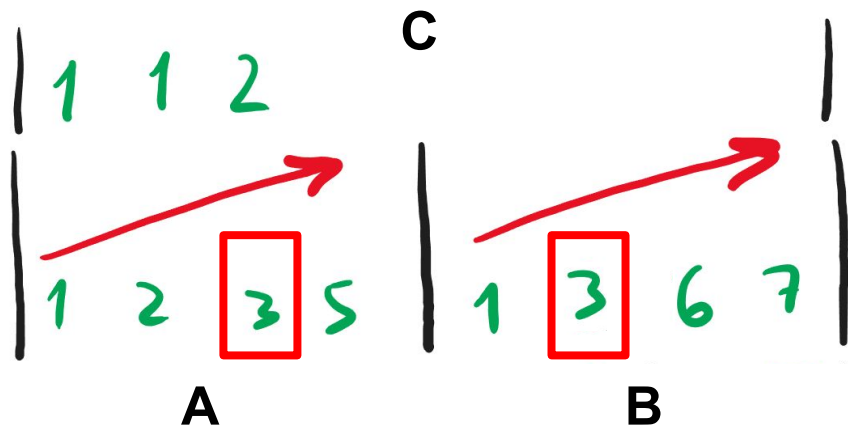
Merge-sort



Как отсортировать на видеокарте?
Т.е. как отсортировать супер
многопоточно?

Достаточно научиться выполнять
операцию merge супер многопоточно!

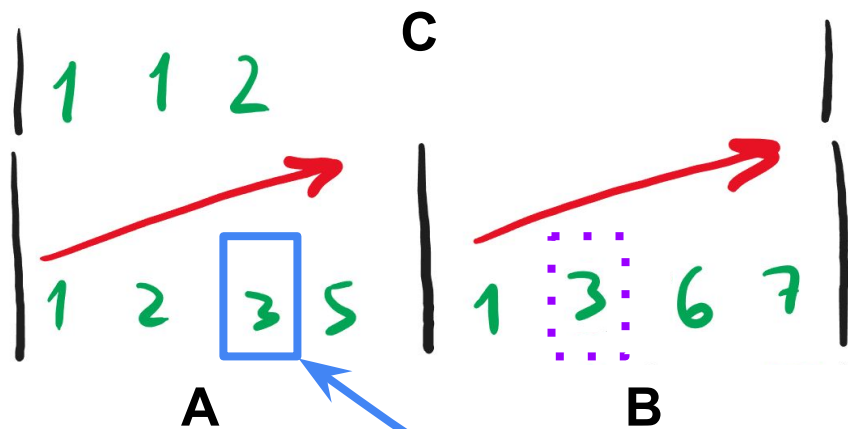
Merge-sort



Как отсортировать на видеокарте?
Т.е. как отсортировать супер
многопоточно?

Достаточно научиться выполнять
операцию merge супер многопоточно!

Merge-sort

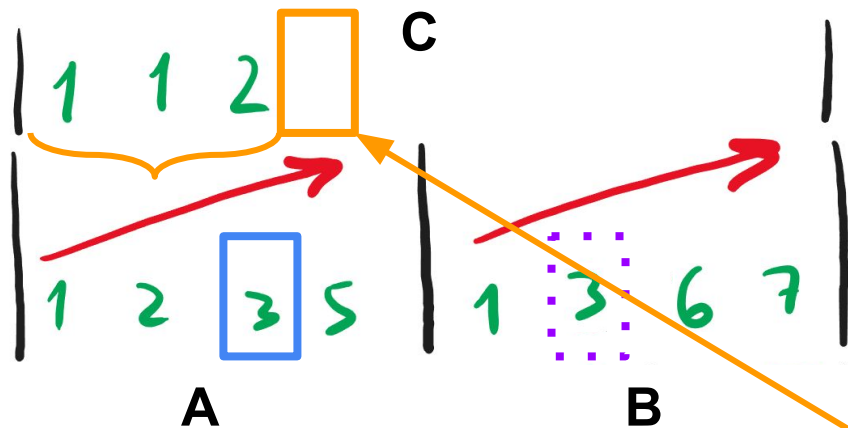


Как отсортировать на видеокарте?
Т.е. как отсортировать супер
многопоточно?

Достаточно научиться выполнять
операцию merge супер многопоточно!

Например один поток - **одно число из A**

Merge-sort



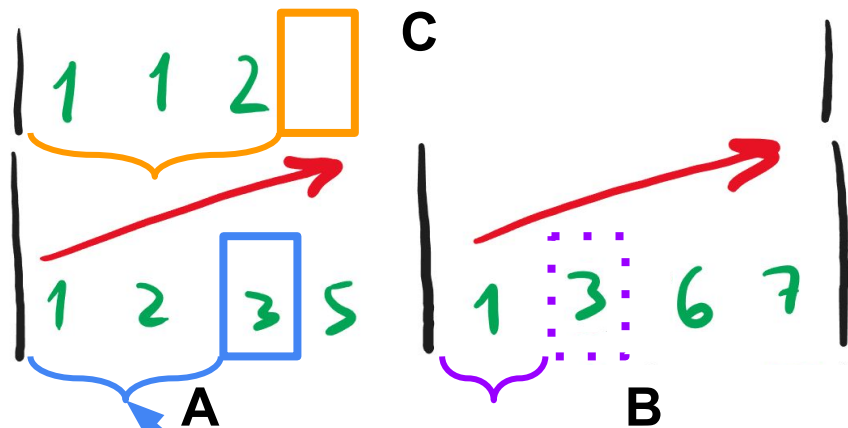
Как отсортировать на видеокарте?
Т.е. как отсортировать супер
многопоточно?

Достаточно научиться выполнять
операцию merge супер многопоточно!

Например один поток - **одно число из A**

На какую **позицию в C** записать его?

Merge-sort



Знаем ли мы сколько чисел до нас в A?

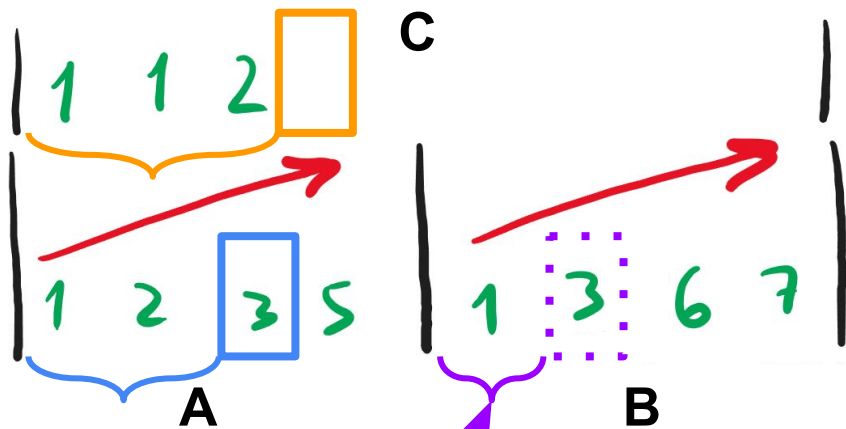
Как отсортировать на видеокарте?
Т.е. как отсортировать супер
многопоточно?

Достаточно научиться выполнять
операцию merge супер многопоточно!

Например один поток - одно число из A

На какую позицию в C записать его?

Merge-sort



Знаем ли мы **сколько** чисел до нас в A?

А как найти **сколько** элементов в B
меньше чем **наше число из A**?

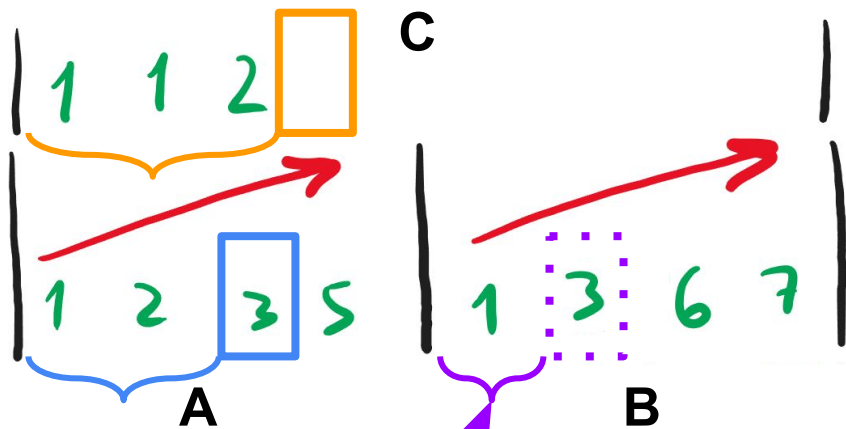
Как отсортировать на видеокарте?
Т.е. как отсортировать супер
многопоточно?

Достаточно научиться выполнять
операцию merge супер многопоточно!

Например один поток - **одно число из A**

На какую **позицию в C** записать его?

Merge-sort



Знаем ли мы **сколько** чисел до нас в A?

А как найти **сколько** элементов в B
меньше чем **наше число из A**?

Бинарный поиск! $O(\log N)$

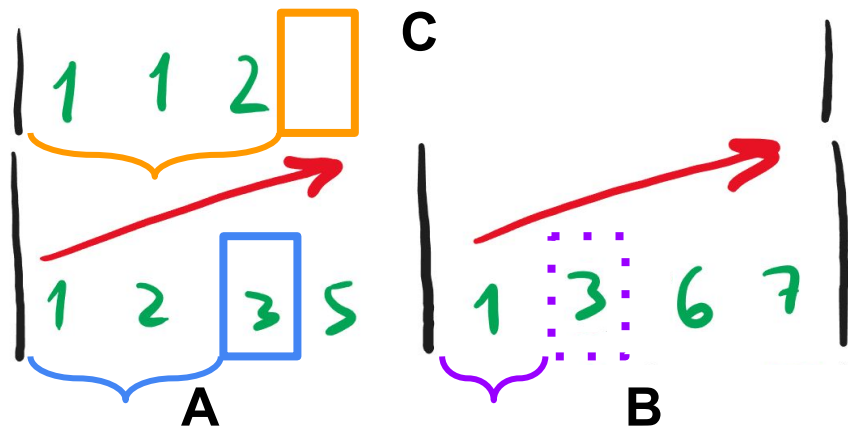
**Как отсортировать на видеокарте?
Т.е. как отсортировать супер
многопоточно?**

Достаточно научиться выполнять
операцию merge супер многопоточно!

Например один поток - **одно число из A**

На какую **позицию в C** записать его?

Merge-sort



Пусть **размер** **A** и **B** - **N** элементов.
Пусть у нас **K** **ядер** в видеокарте.

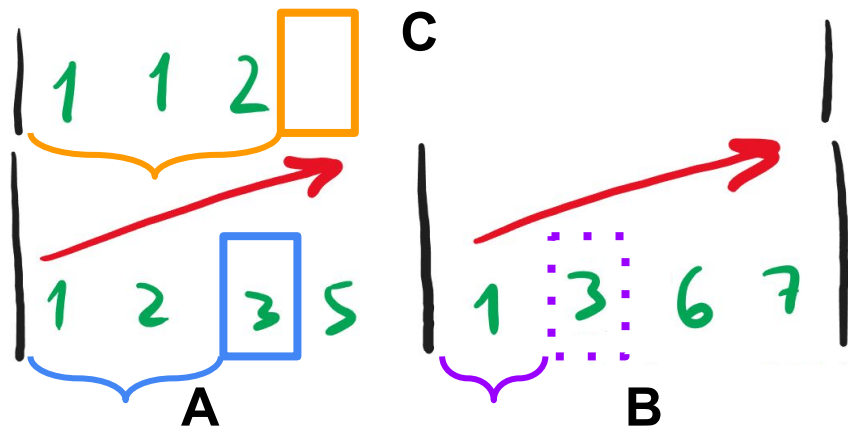
Какая асимптотика **merge** операции?

Знаем ли мы **сколько** чисел до нас в **A**?

А как найти **сколько** элементов в **B**
меньше чем **наше** **число** из **A**?

Бинарный поиск! $O(\log N)$

Merge-sort



Пусть **размер** **A** и **B** - **N** элементов.
Пусть у нас **K** **ядер** в видеокарте.

Какая асимптотика **merge** операции?

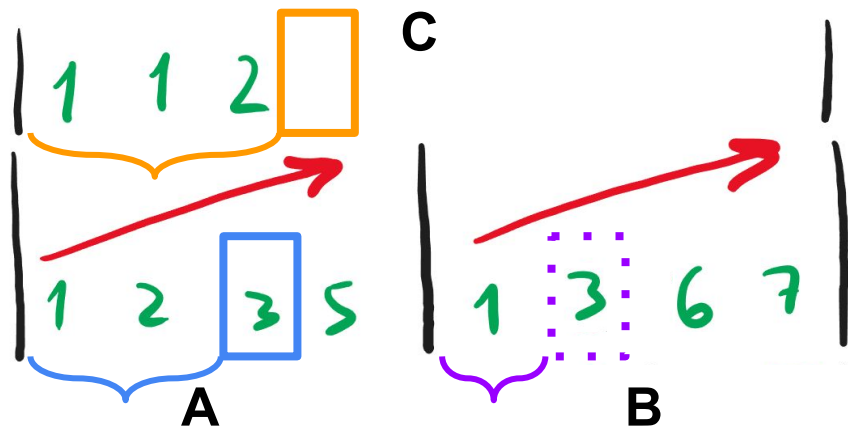
$$O(?) * \boxed{O(\log N)} / ?$$

Знаем ли мы **сколько** чисел до нас в **A**?

А как найти **сколько** элементов в **B**
меньше чем **наше** число из **A**?

Бинарный поиск! $\boxed{O(\log N)}$

Merge-sort



Пусть **размер A и B** - N элементов.
Пусть у нас **K ядер** в видеокарте.

Какая асимптотика **merge** операции?

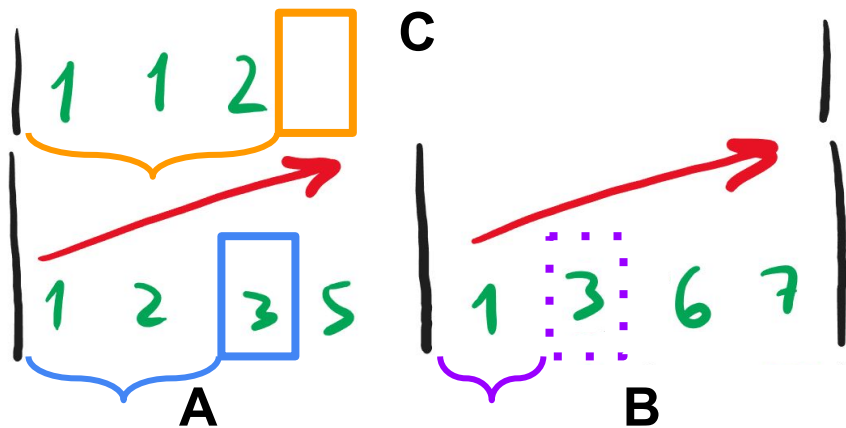
$O(N * O(\log N) / ?)$

Знаем ли мы **сколько чисел до нас в A**?

А как найти **сколько элементов в B**
меньше чем **наше число из A**?

Бинарный поиск! $O(\log N)$

Merge-sort



Пусть **размер A и B** - N элементов.
Пусть у нас **K ядер** в видеокарте.

Какая асимптотика **merge** операции?

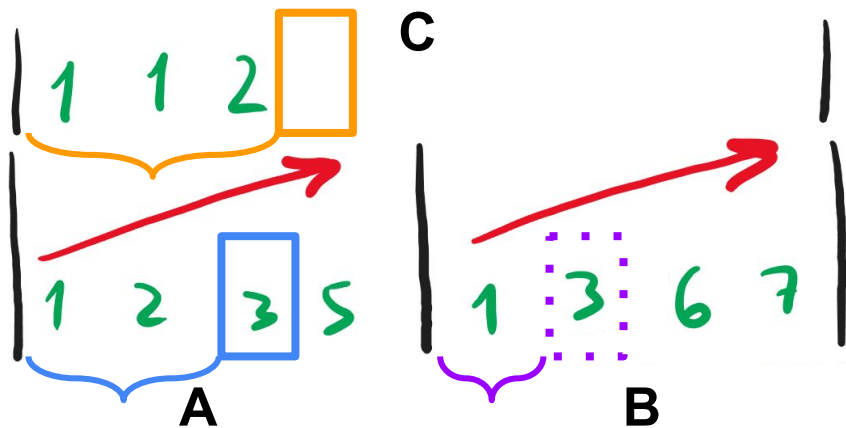
$$O(N * O(\log N) / K)$$

Знаем ли мы **сколько чисел до нас в A**?

А как найти **сколько элементов в B**
меньше чем **наше число из A**?

Бинарный поиск! $O(\log N)$

Merge-sort



Пусть **размер A и B** - N элементов.
Пусть у нас **K ядер** в видеокарте.

Какая асимптотика **merge** операции?

$O(N * O(\log N) / K)$

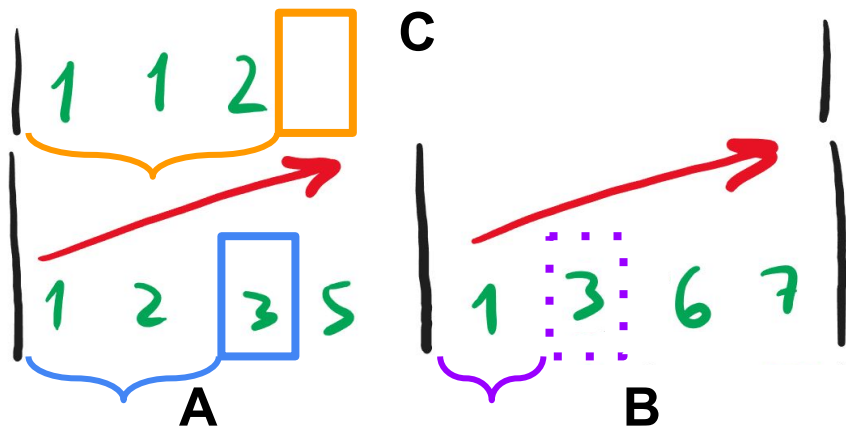
Всего у нас $O(\log N)$ операций слияния

Знаем ли мы **сколько чисел до нас в A**?

А как найти **сколько элементов в B**
меньше чем **наше число из A**?

Бинарный поиск! $O(\log N)$

Merge-sort



Пусть **размер A и B** - **N** элементов.
Пусть у нас **K** ядер в видеокарте.

Какая асимптотика **merge** операции?

$O(N * O(\log N) / K)$

Всего у нас $O(\log N)$ операций слияния

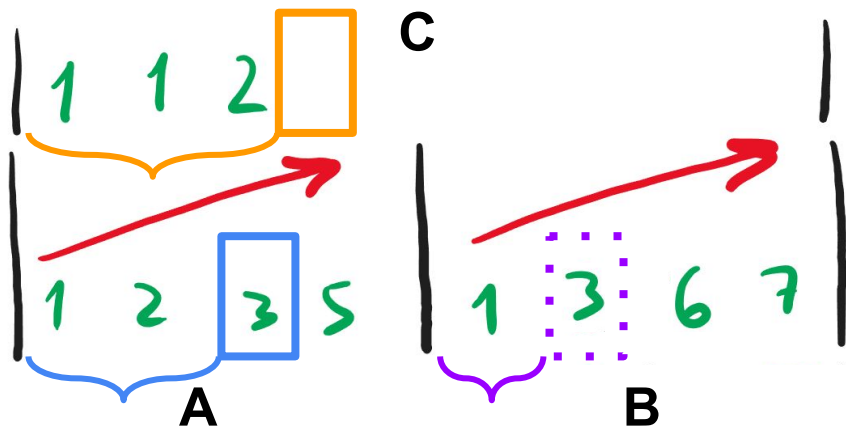
Знаем ли мы **сколько чисел до нас в A**?

А как найти **сколько элементов в B**
меньше чем **наше число из A**?

Бинарный поиск! $O(\log N)$

Какая асимптотика на CPU?

Merge-sort



Пусть **размер A и B** - N элементов.
Пусть у нас **K ядер** в видеокарте.

Какая асимптотика **merge** операции?

$O(N * O(\log N) / K)$

Всего у нас $O(\log N)$ операций слияния

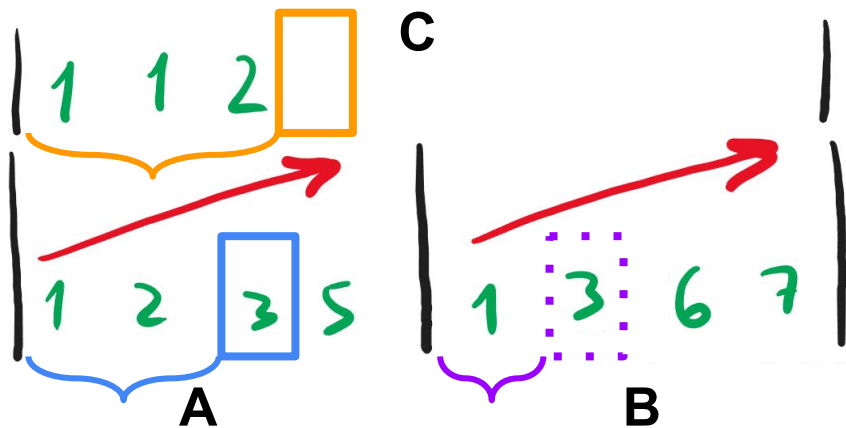
Какая асимптотика на CPU? $O(\log N * N)$

Знаем ли мы **сколько чисел до нас в A**?

А как найти **сколько элементов в B**
меньше чем **наше число из A**?

Бинарный поиск! $O(\log N)$

Merge-sort



Знаем ли мы сколько чисел до нас в A?

А как найти сколько элементов в B меньше чем наше число из A?

Бинарный поиск! $O(\log N)$

Пусть размер A и B - N элементов.
Пусть у нас K ядер в видеокарте.

Какая асимптотика merge операции?

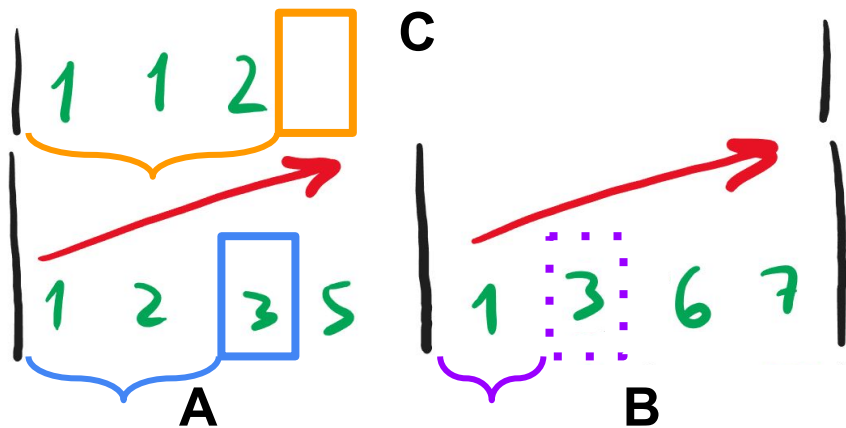
$O(N * O(\log N) / K)$

Всего у нас $O(\log N)$ операций слияния

Какая асимптотика на CPU? $O(\log N * N)$

А на GPU? $O(\log N * ?)$

Merge-sort



Знаем ли мы сколько чисел до нас в A?

А как найти сколько элементов в B меньше чем наше число из A?

Бинарный поиск! $O(\log N)$

Пусть размер A и B - N элементов.
Пусть у нас K ядер в видеокарте.

Какая асимптотика merge операции?

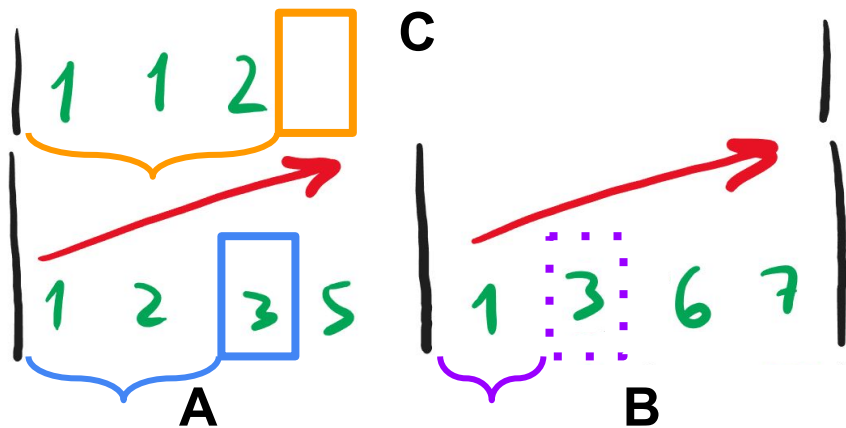
$O(N * O(\log N) / K)$

Всего у нас $O(\log N)$ операций слияния

Какая асимптотика на CPU? $O(\log N * N)$

А на GPU? $O(\log N * N * \log N / K)$

Merge-sort



Знаем ли мы сколько чисел до нас в A?

А как найти сколько элементов в B меньше чем наше число из A?

Бинарный поиск! $O(\log N)$

Пусть размер A и B - N элементов.
Пусть у нас K ядер в видеокарте.

Какая асимптотика merge операции?

$O(N * O(\log N) / K)$

Всего у нас $O(\log N)$ операций слияния

Какая асимптотика на CPU? $O(\log N * N)$

А на GPU? $O(\log N * N * \log N / K)$

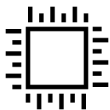
Отличается в $(\log N / K)$ раз, K - большое!

Merge-sort

$N = 33.554.432 \sim 3 \cdot 10^7$

CPU - Intel 13700K

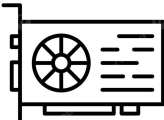
$O(\log N * N)$



6.462 секунд

GPU - NVIDIA RTX 4090

$O(\log N * N * \log N / K)$



0.011 секунд

Быстрее в 582 раза!

Отличается в $(\log N / K)$ раз,
где $K = 16$ тысяч ядер



Massively Parallel Multiview Stereopsis by Surface Normal Diffusion

Silvano Galliani

Katrin Lasinger

Konrad Schindler

Photogrammetry and Remote Sensing, ETH Zurich

Abstract

We present a new, massively parallel method for high-quality multiview matching. Our work builds on the Patchmatch idea: starting from randomly generated 3D planes in scene space, the best-fitting planes are iteratively propagated and refined to obtain a 3D depth and normal field per view, such that a robust photo-consistency measure over all images is maximized. Our main novelties are on the one hand to formulate Patchmatch in scene space, which makes it possible to aggregate image similarity across multiple views and obtain more accurate depth maps. And on the other hand a modified, diffusion-like propagation scheme that can be massively parallelized and delivers dense multiview correspondence over ten 1.9-Megapixel images in 3 seconds, on a consumer-grade GPU. Our method uses a slanted support window and thus has no fronto-parallel bias; it is completely local and parallel, such that computation time scales linearly with image size, and inversely proportional to the number of parallel threads. Further-

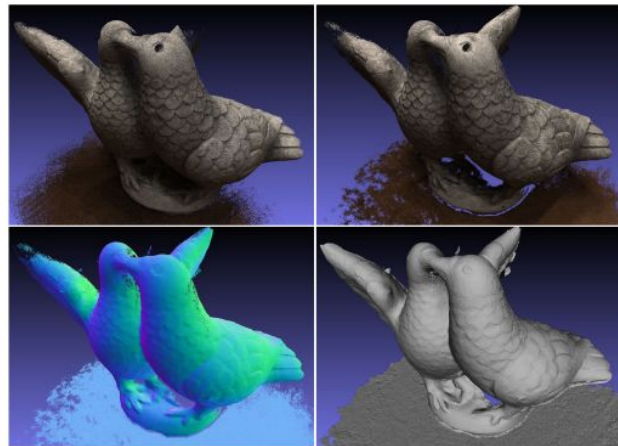
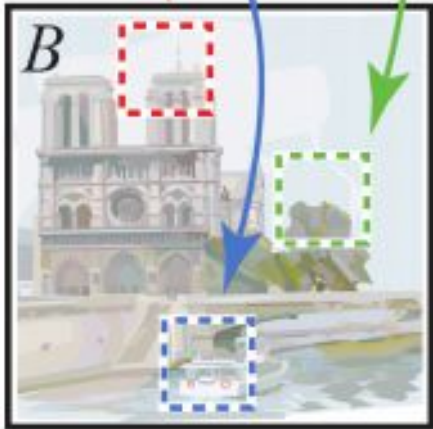
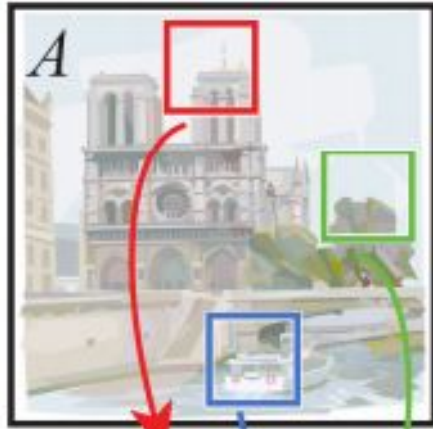
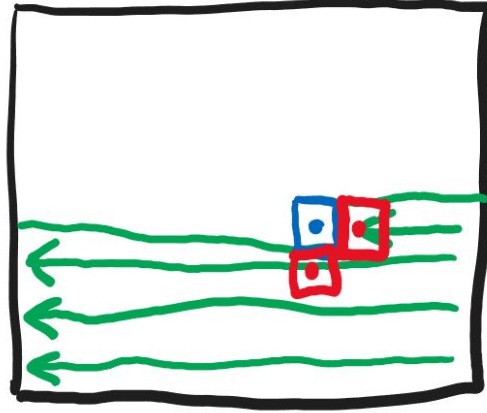
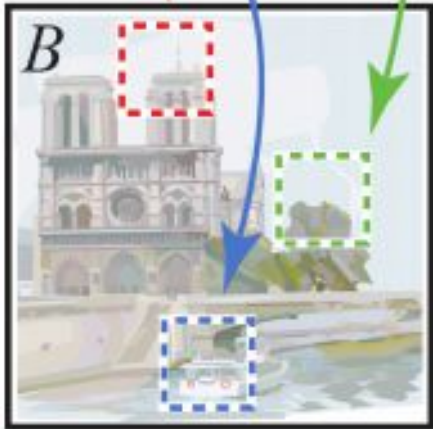
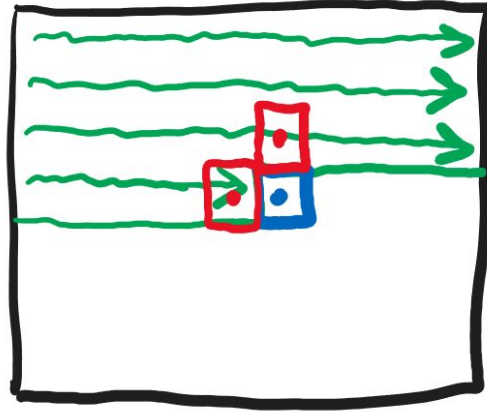
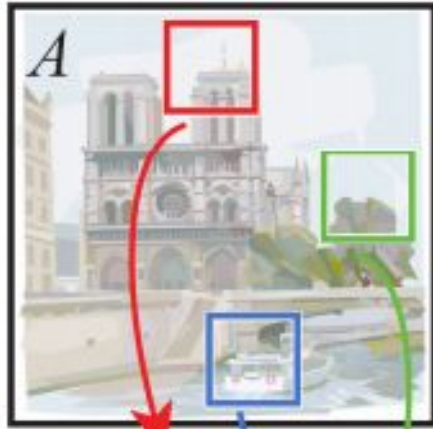


Figure 1: Results on one of the 80 evaluated objects on the DTU benchmark [22]. *Top left*: Ground truth point cloud; *top right*: reconstructed point cloud with texture; *bottom left*: color-coded surface normals; *bottom right*: reconstructed surface.

Gipuma

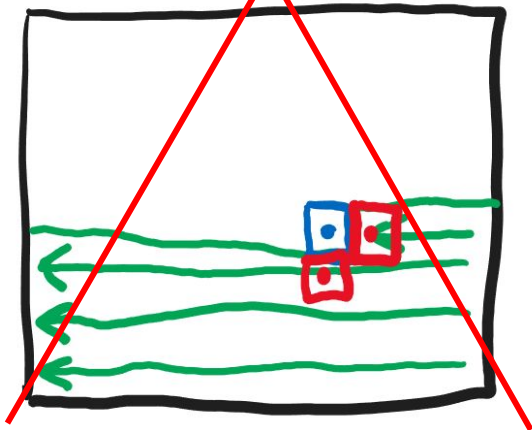
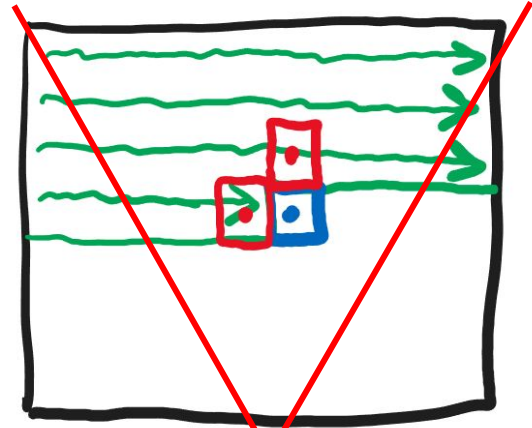
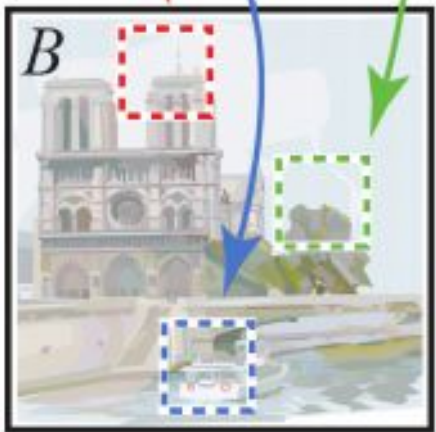
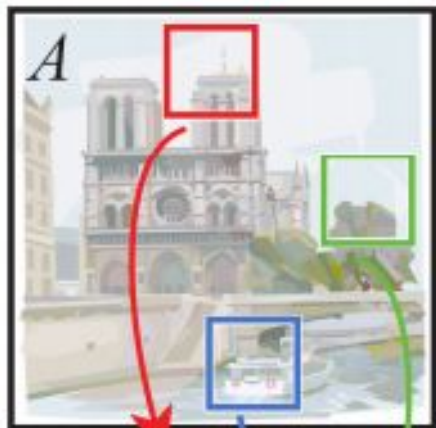


Gipuma



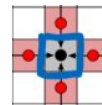
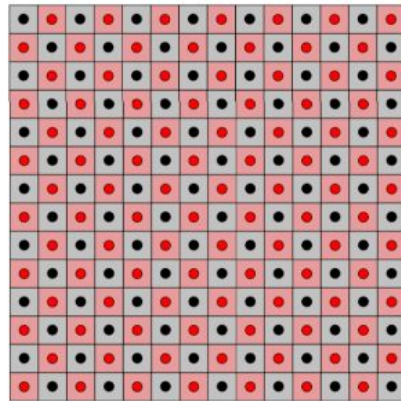
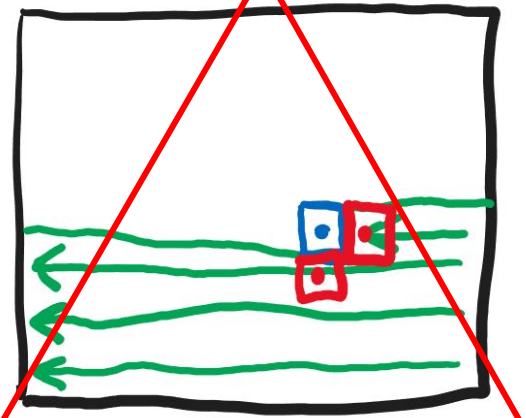
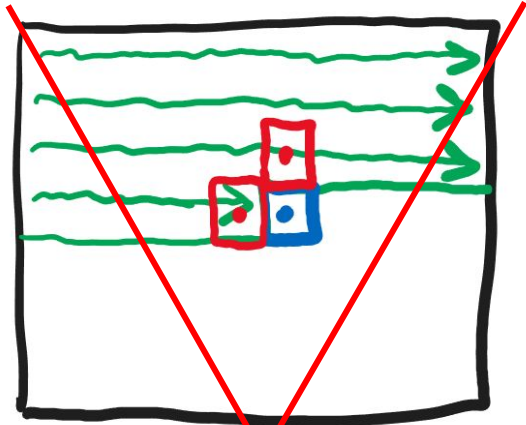
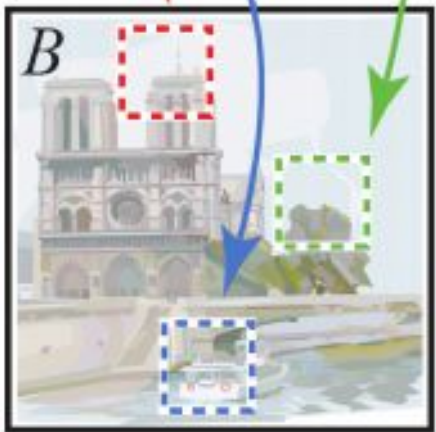
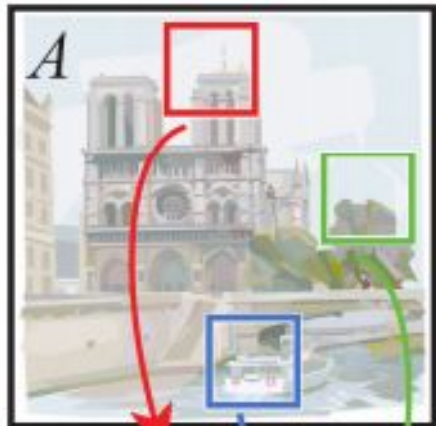
Giruma

Как адаптировать для
массового параллелизма?



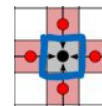
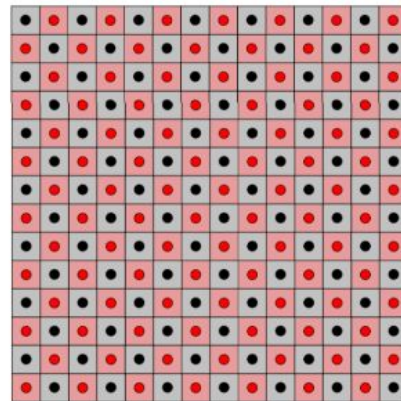
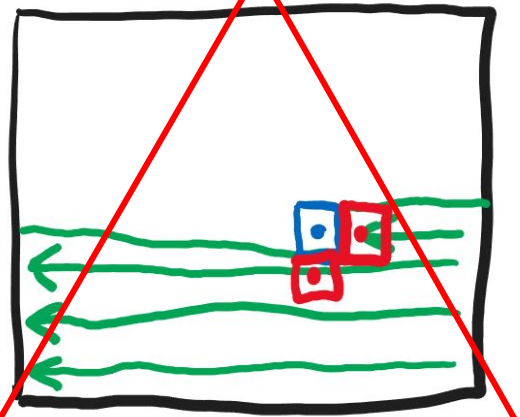
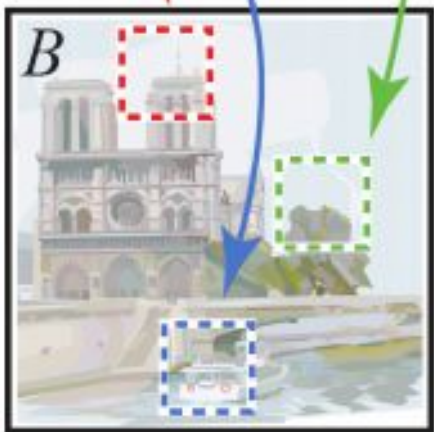
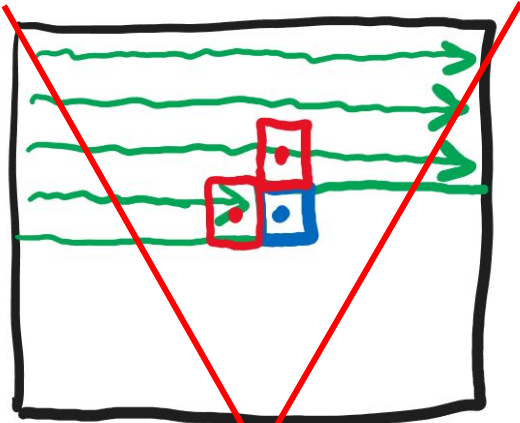
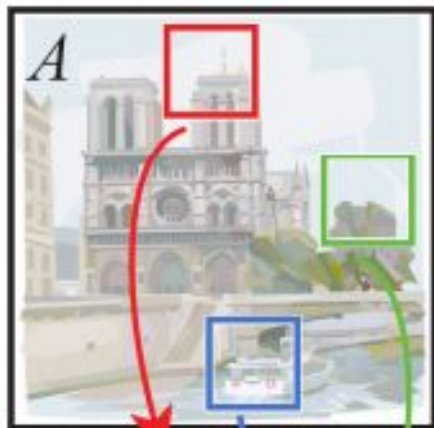
Gipuma

Как адаптировать для
массового параллелизма?



Giruma

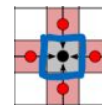
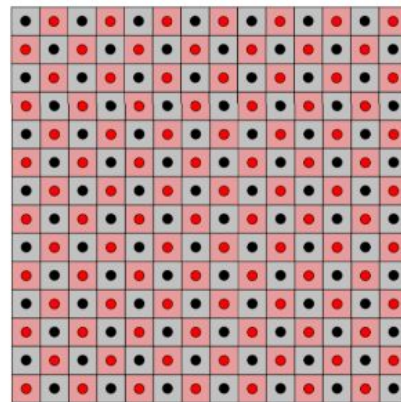
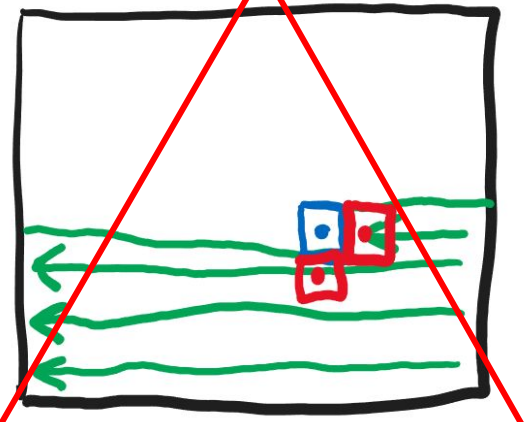
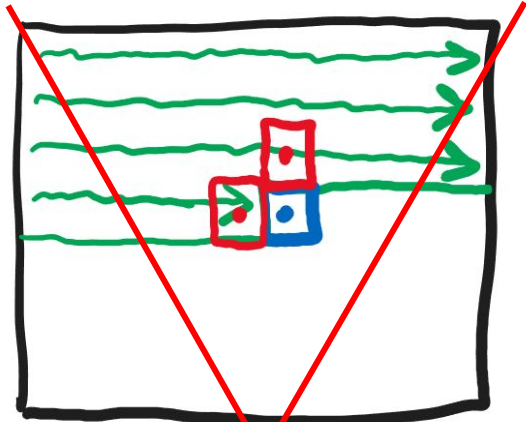
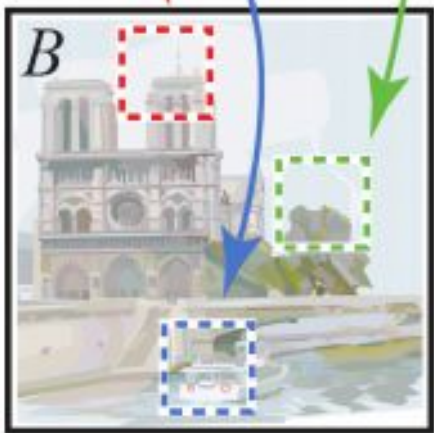
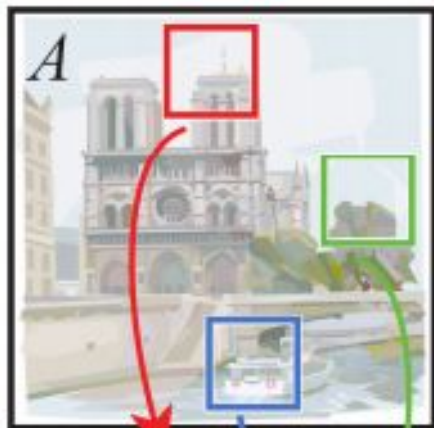
Как адаптировать для
массового параллелизма?



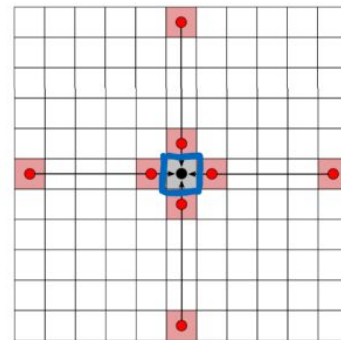
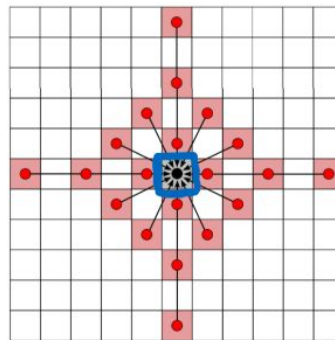
Как ускорить распространение
информации? (правильного ответа)

Gipuma

Как адаптировать для
массового параллелизма?



Как ускорить распространение
информации? (правильного ответа)




Глава 7: Матрицы

транспонирование, умножение, tensor cores, DeepSeek

Транспонирование матрицы

2	4	-1
-10	5	11
18	-7	6



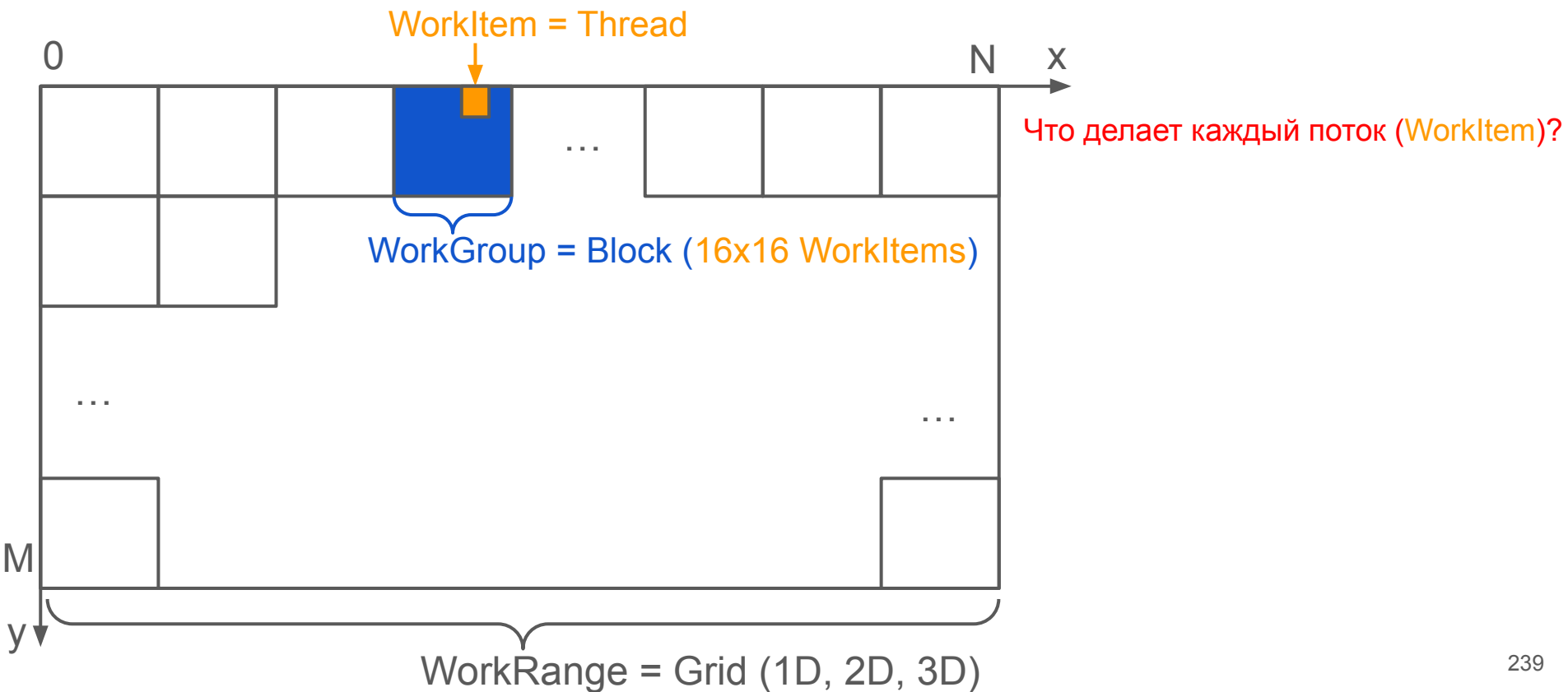
2	-10	18
4	5	-7
-1	11	6

Как это реализовать в модели массового параллелизма?

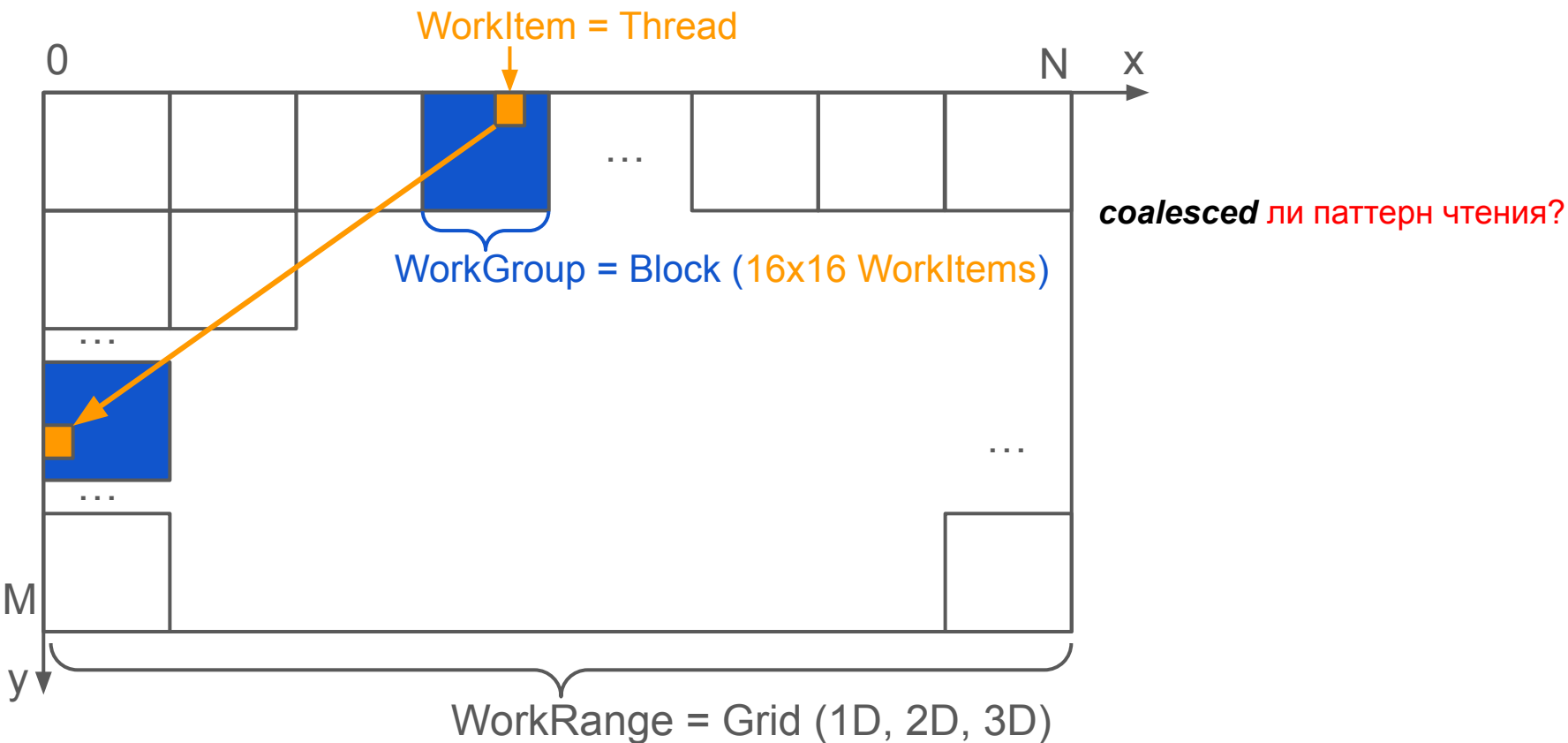
Какой размер рабочего пространства?

Что делает каждый поток (work item)?

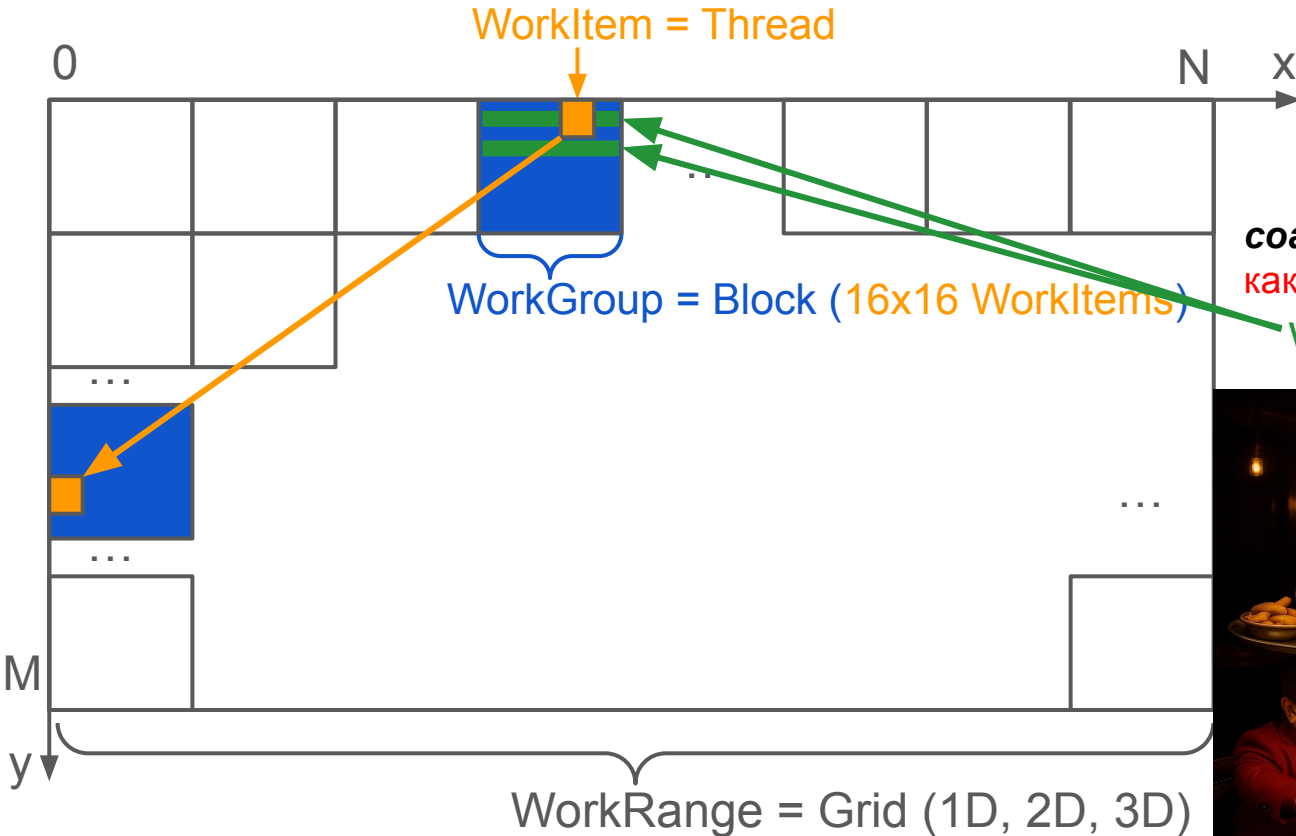
Транспонирование матрицы



Транспонирование матрицы (*coalesced memory access*)



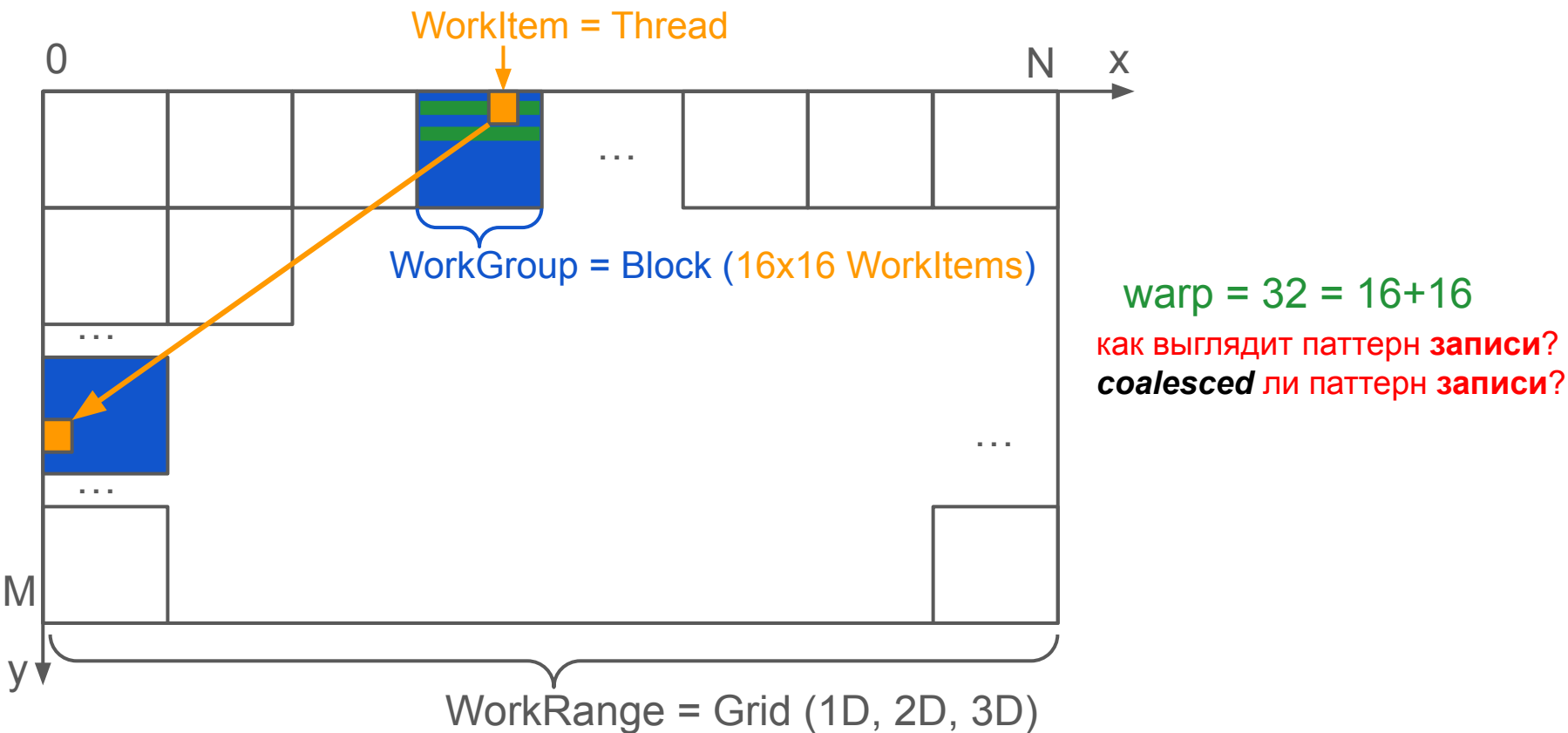
Транспонирование матрицы (*coalesced memory access*)



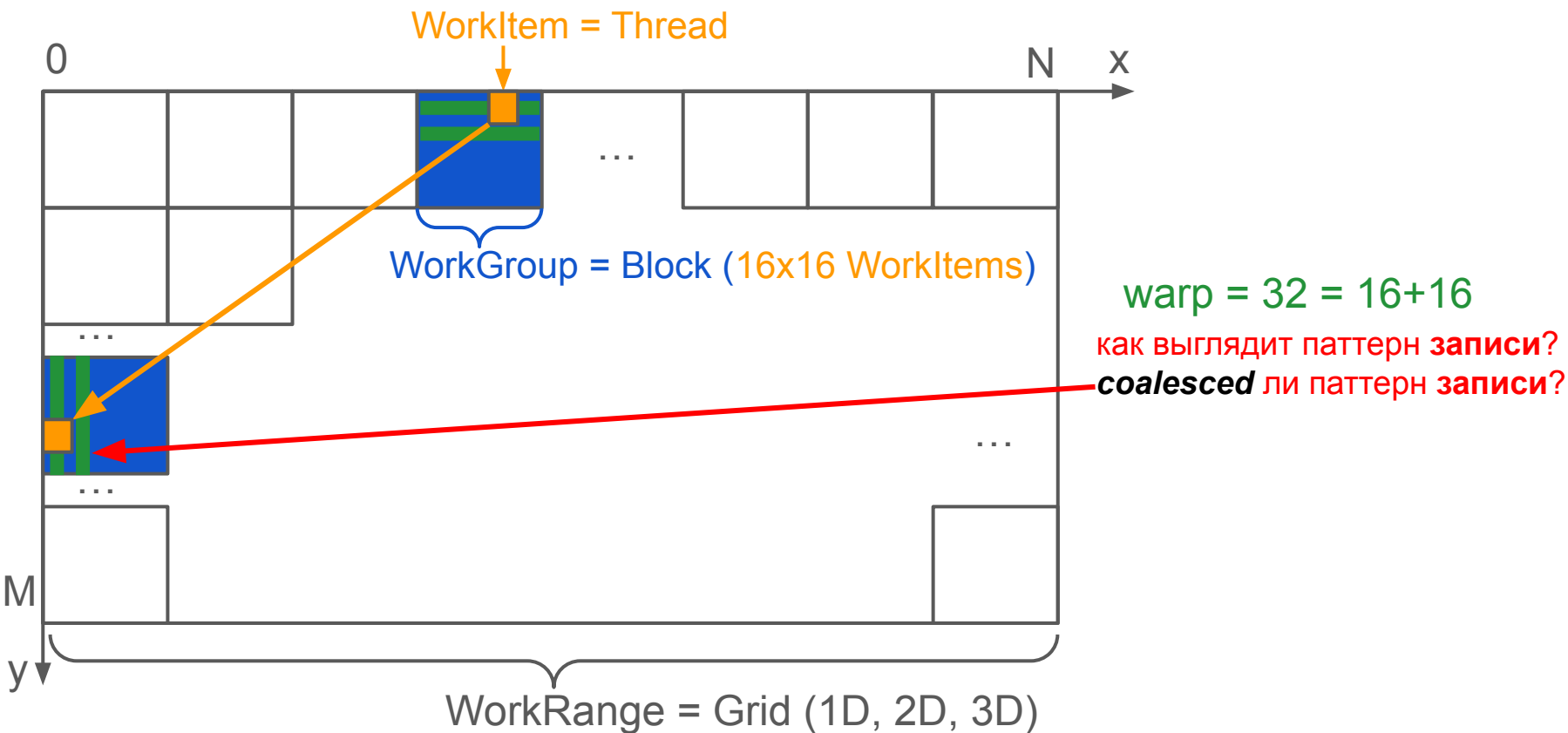
coalesced ли паттерн чтения?
как выглядит выкладка **warp**?
warp = 32 = 16+16



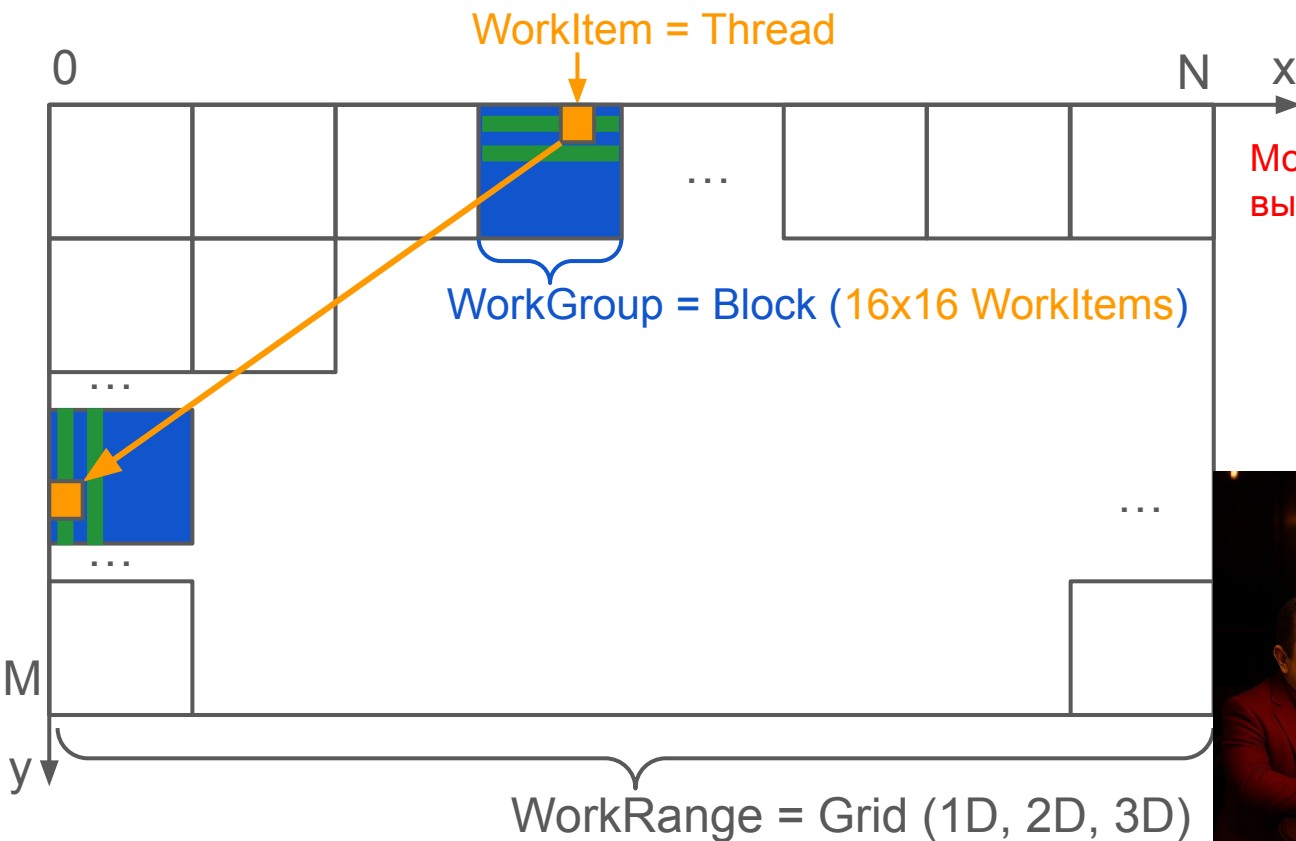
Транспонирование матрицы (*coalesced memory access*)



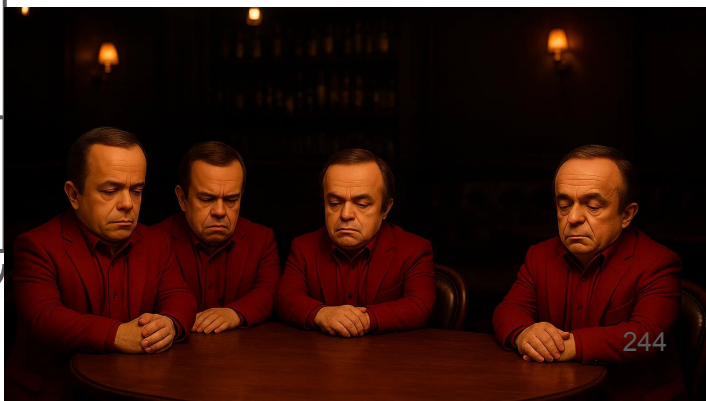
Транспонирование матрицы (*coalesced memory access*)



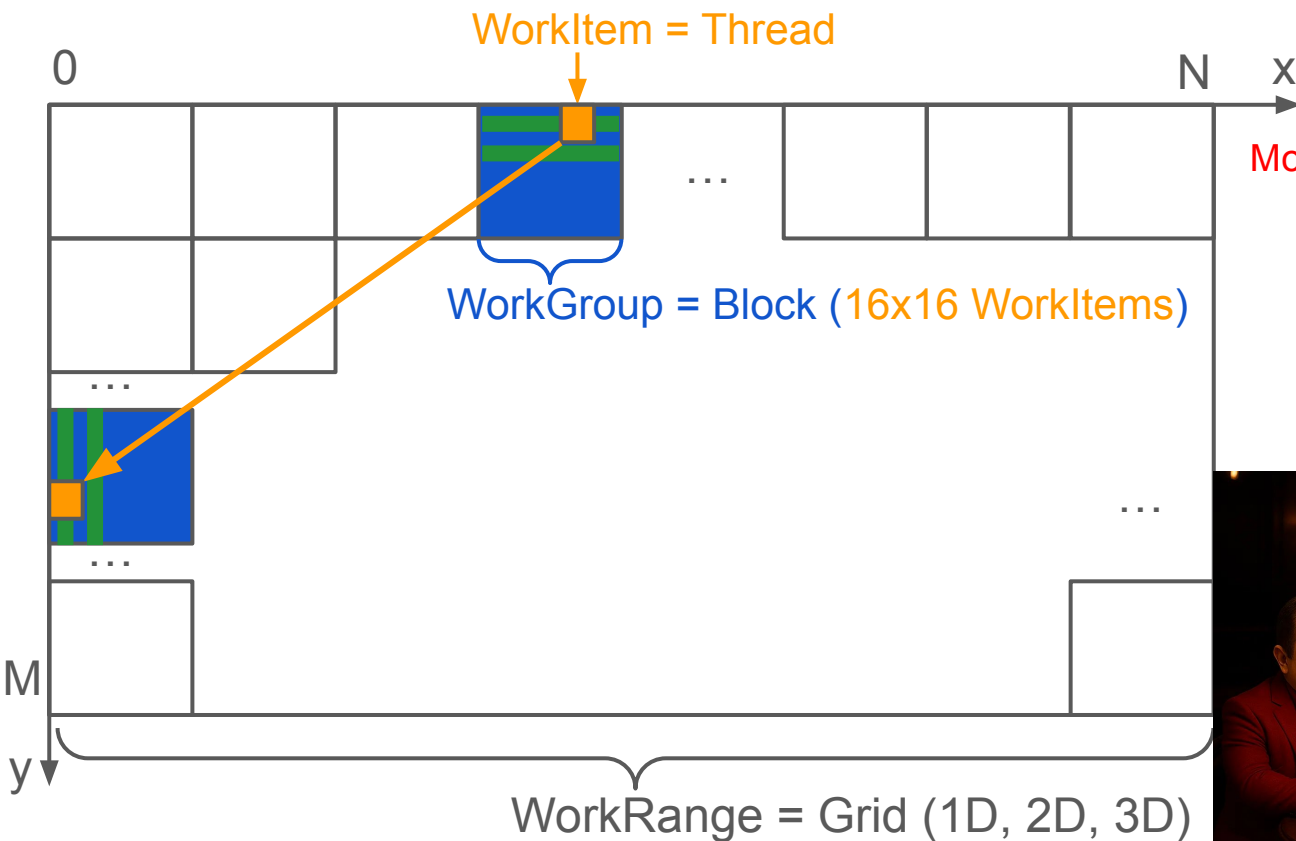
Транспонирование матрицы (*coalesced memory access*)



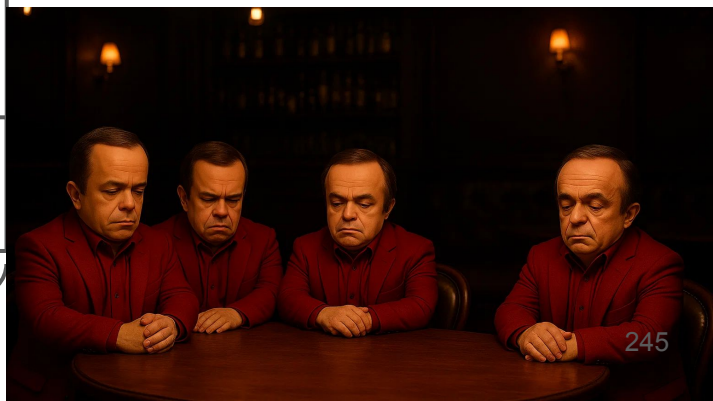
Можем ли мы исправить это поменяв выкладку **warp**-а?



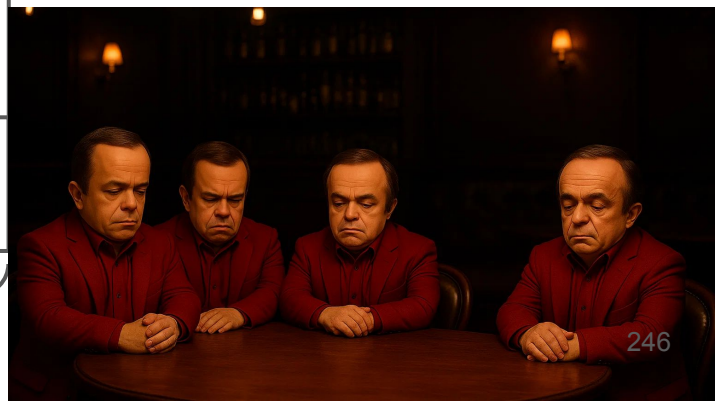
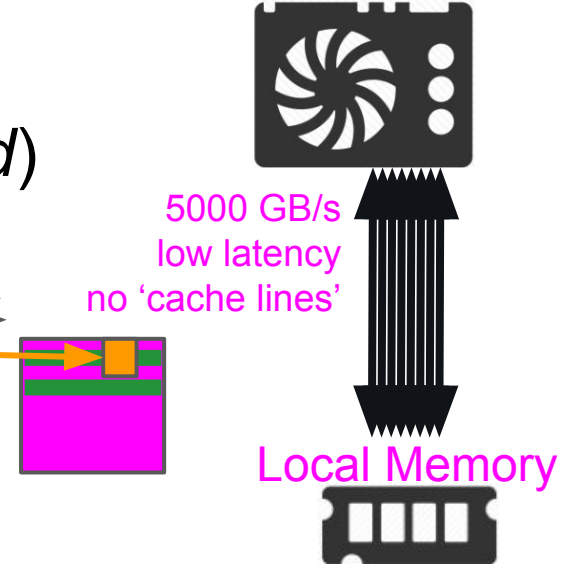
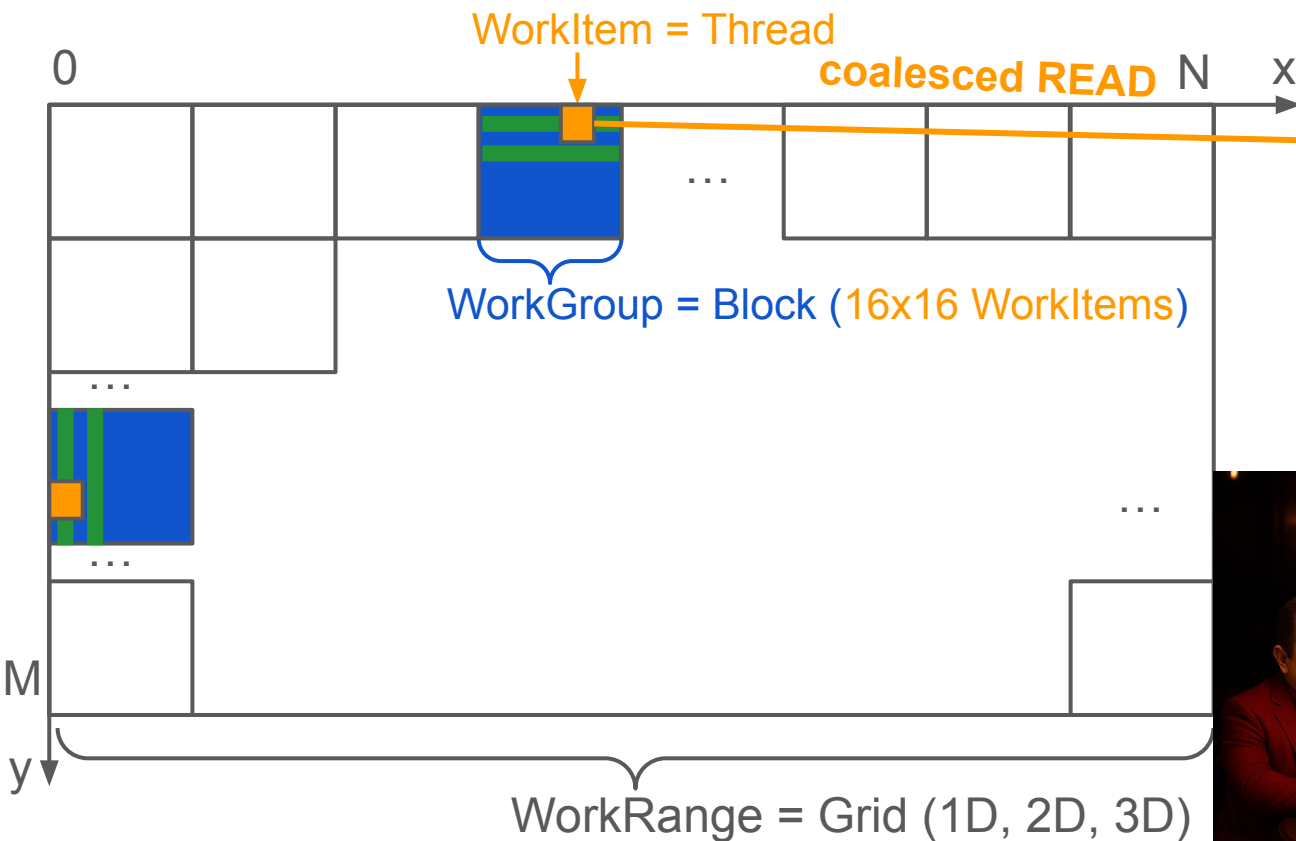
Транспонирование матрицы (*coalesced memory access*)



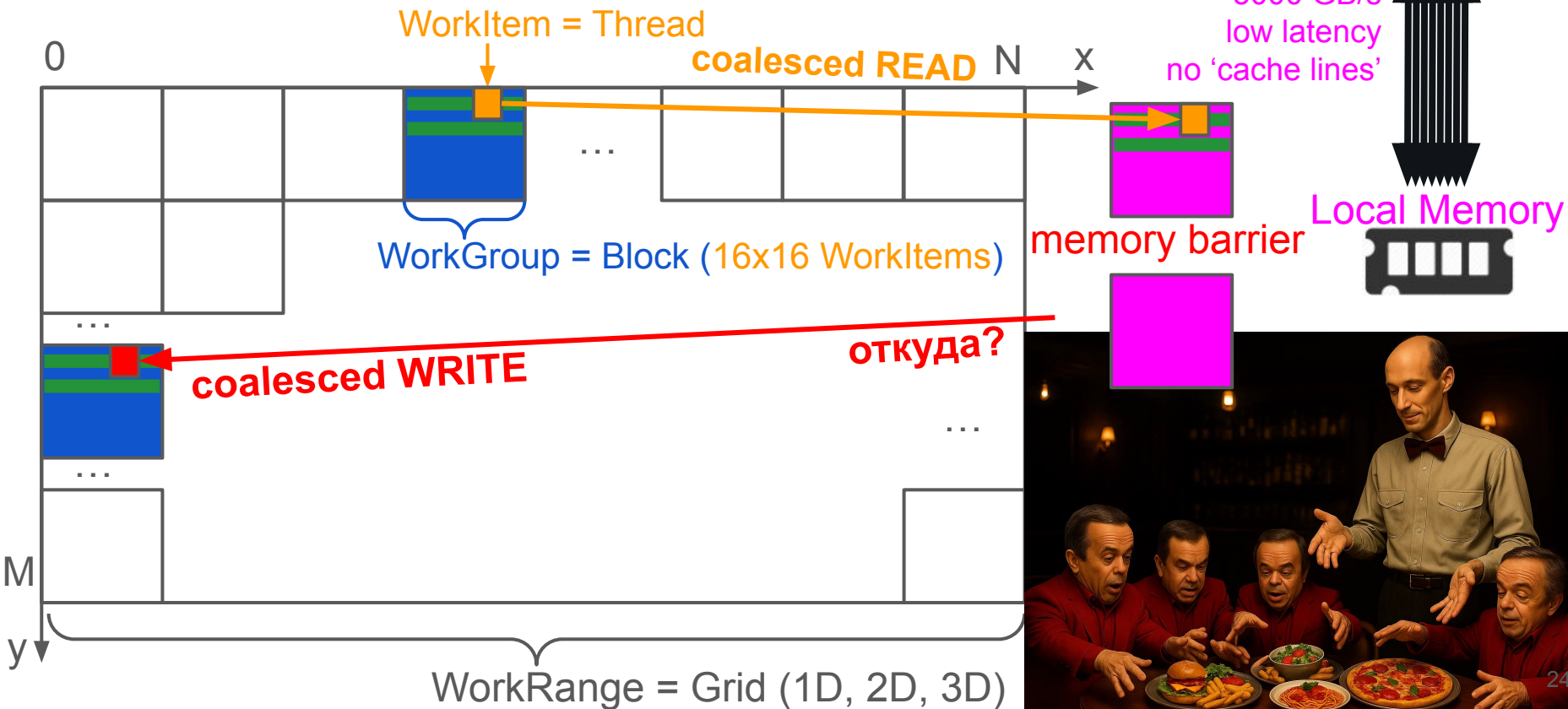
Можем ли мы исправить это?



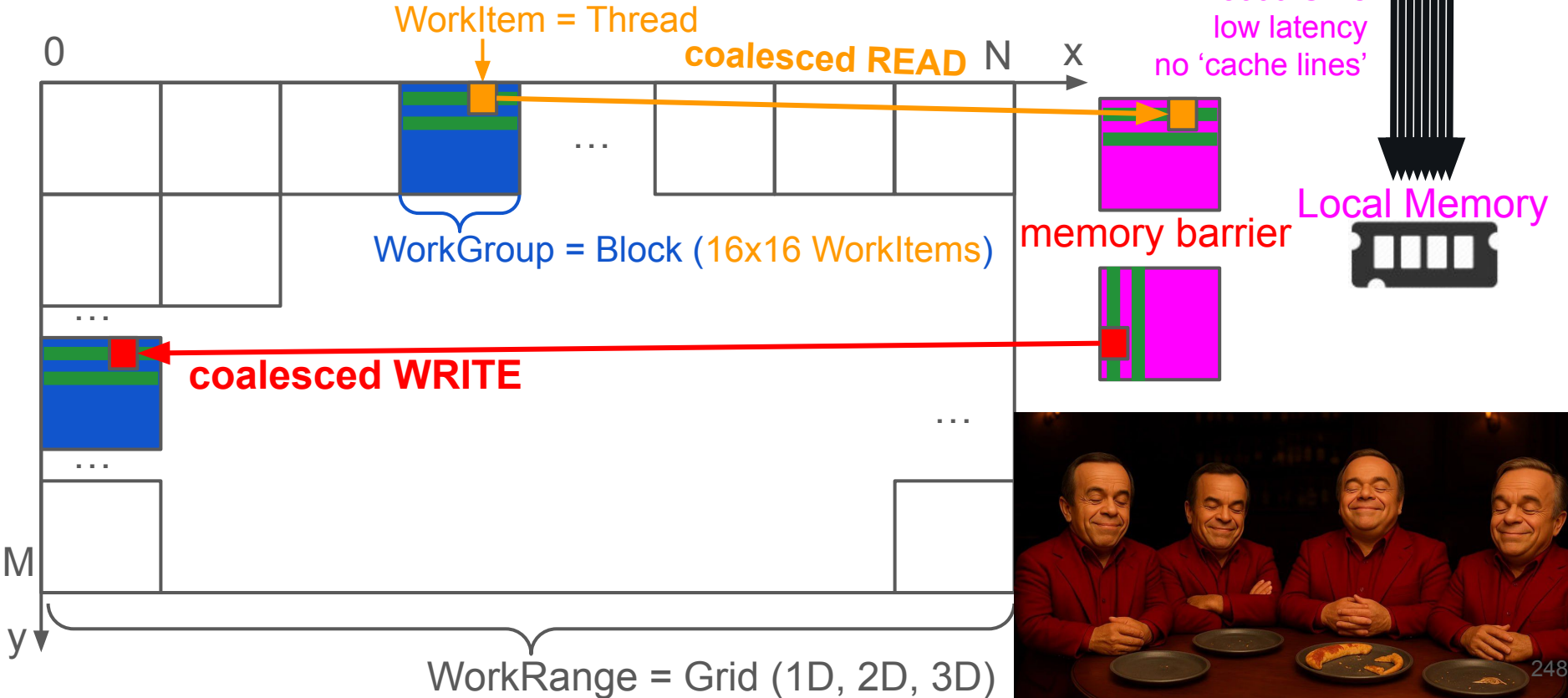
Транспонирование матрицы (coalesced)



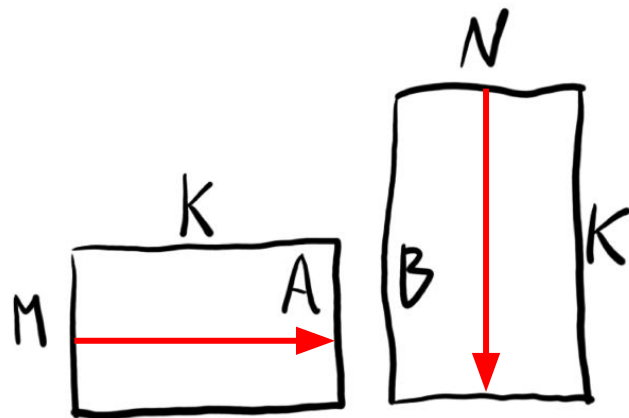
Транспонирование матрицы (coalesced)



Транспонирование матрицы (coalesced)



Умножение матриц



$$C = A \times B$$

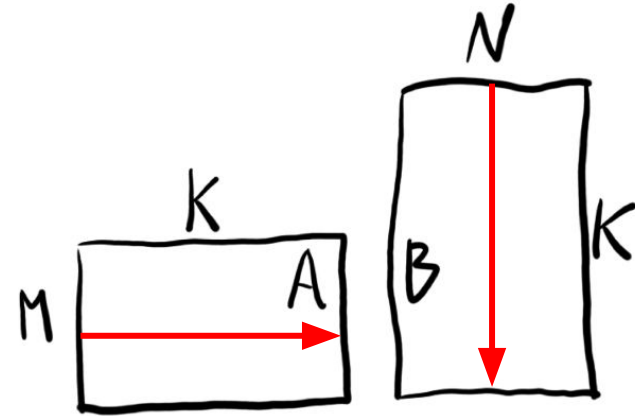
Below the equation, a square represents the resulting matrix C. The top side is labeled 'N' and the left side is labeled 'M'. Inside the square, there is a red circle with a dot in the center, and the letter 'C' is written in the bottom-left corner.

Умножение матриц

Сколько у нас вычислений?

Сколько у нас чтений/записей данных?

Какая пропорция?



$$C = A \times B$$

Diagram illustrating the resulting matrix C. Matrix C is labeled with dimensions M (rows) and N (columns). A red circle with a dot inside is shown within the matrix C, representing a single element of the result.

Умножение матриц

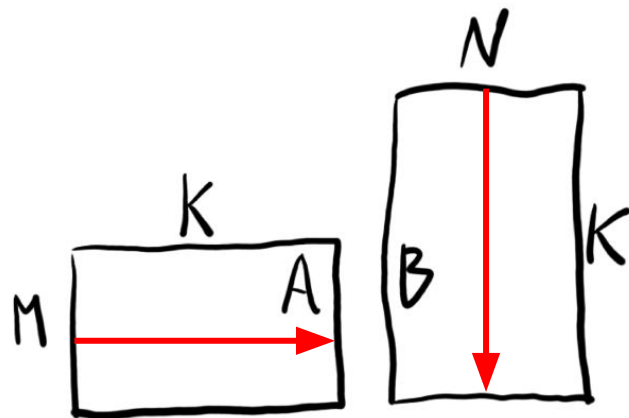
Сколько у нас вычислений?

$O(N \times M \times K)$

Сколько у нас чтений/записей данных?

$O(N \times M \times K)$

Какая пропорция?



$$C = A \times B$$

Diagram illustrating the dimensions of the resulting matrix C. Matrix C is a square with height M and width N. A red circle with a dot inside is shown, representing the result of the multiplication.

Умножение матриц

Сколько у нас вычислений?

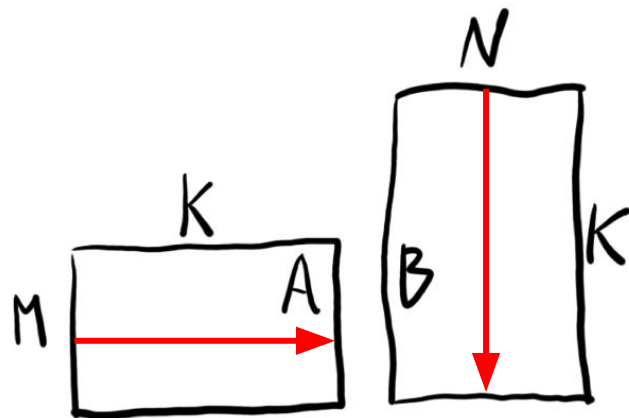
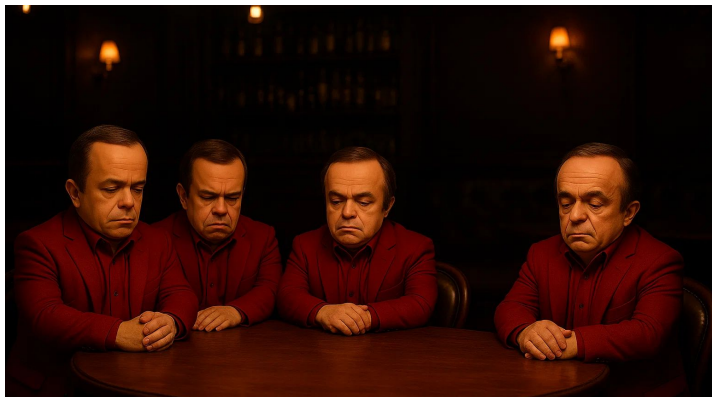
$$O(N \times M \times K)$$

Сколько у нас чтений/записей данных?

$$O(N \times M \times K)$$

Какая пропорция?

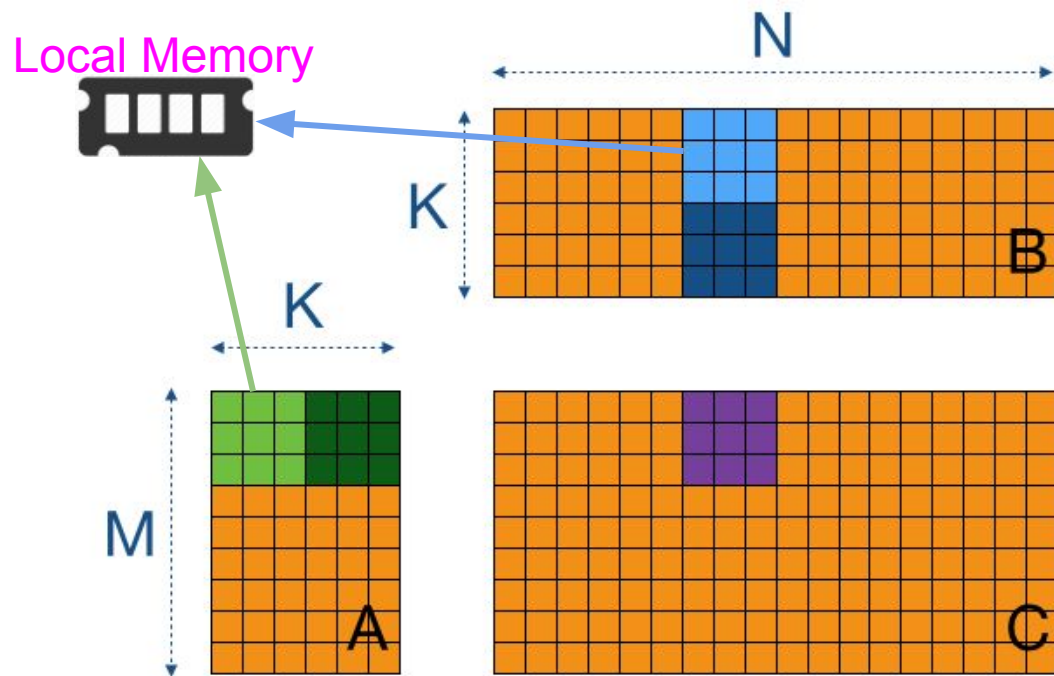
1:1



$$C = A \times B$$
A hand-drawn diagram of the resulting matrix 'C'. It is a square with dimensions 'M' (height) and 'N' (width). Inside the square, there is a red circle with a dot in the center, and the letter 'C' is written in the bottom-left corner.

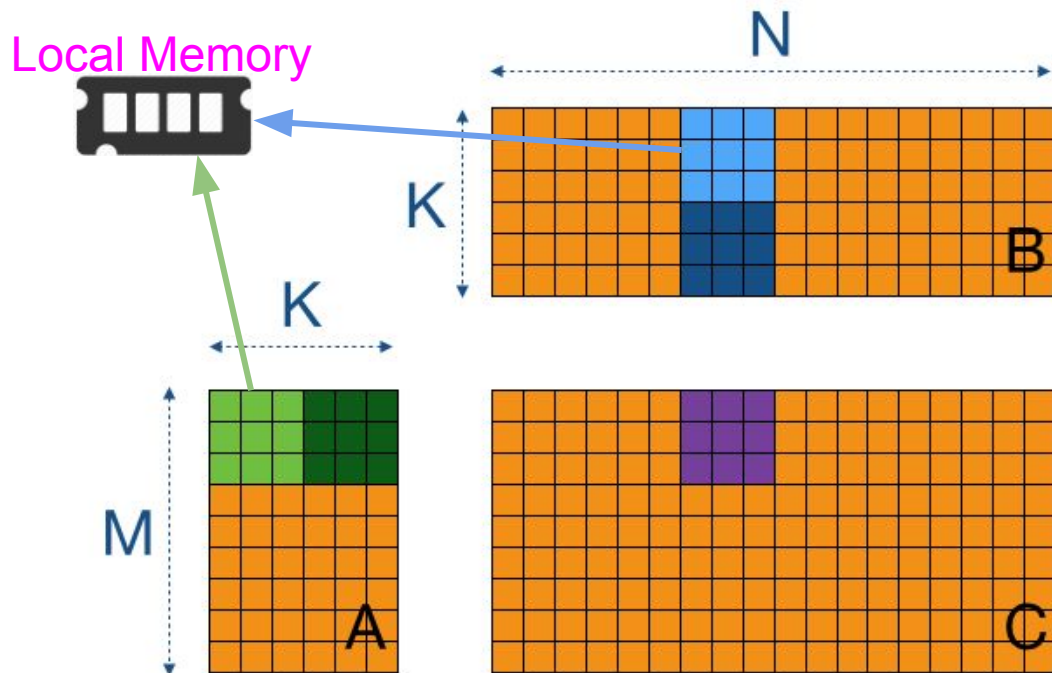
Как увеличить объем вычислений на считанный байт?

Умножение матриц

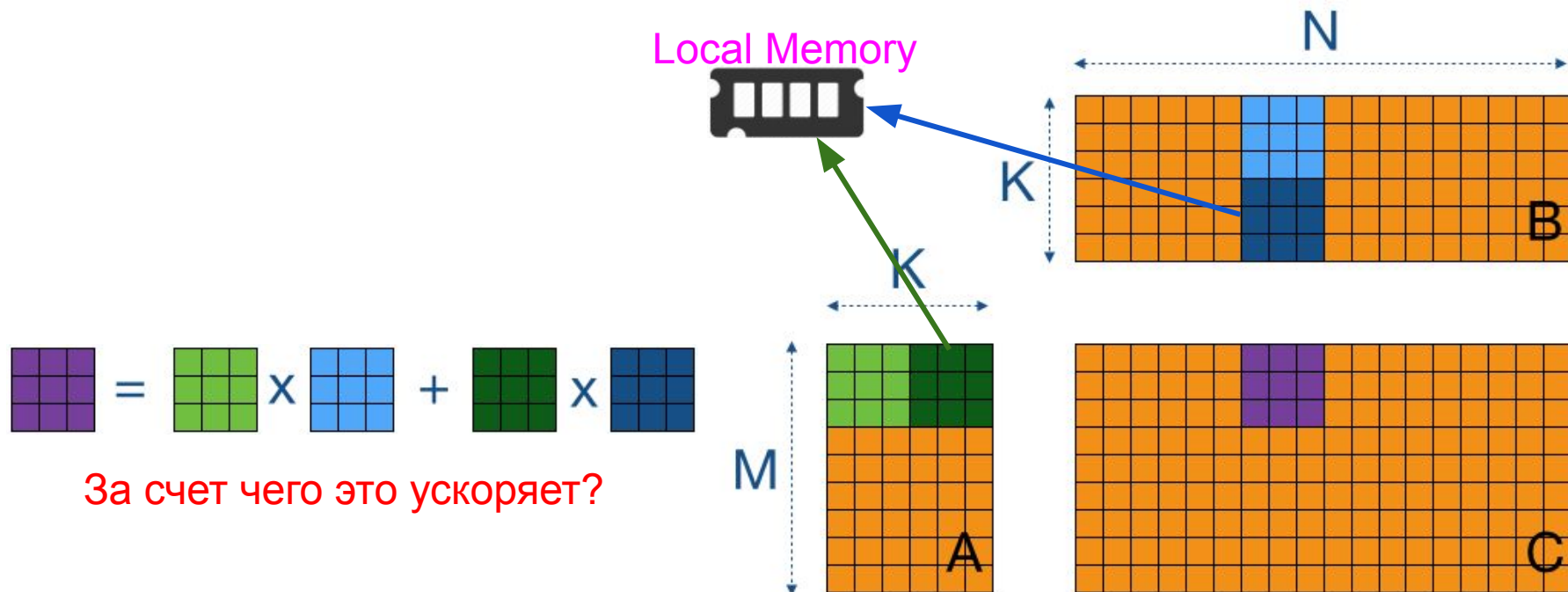


Умножение матриц

$$\begin{bmatrix} \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \end{bmatrix} = \begin{bmatrix} \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \end{bmatrix} \times \begin{bmatrix} \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \end{bmatrix} + \dots$$



Умножение матриц



Умножение матриц

Сколько у нас вычислений?

Сколько у нас чтений/записей?

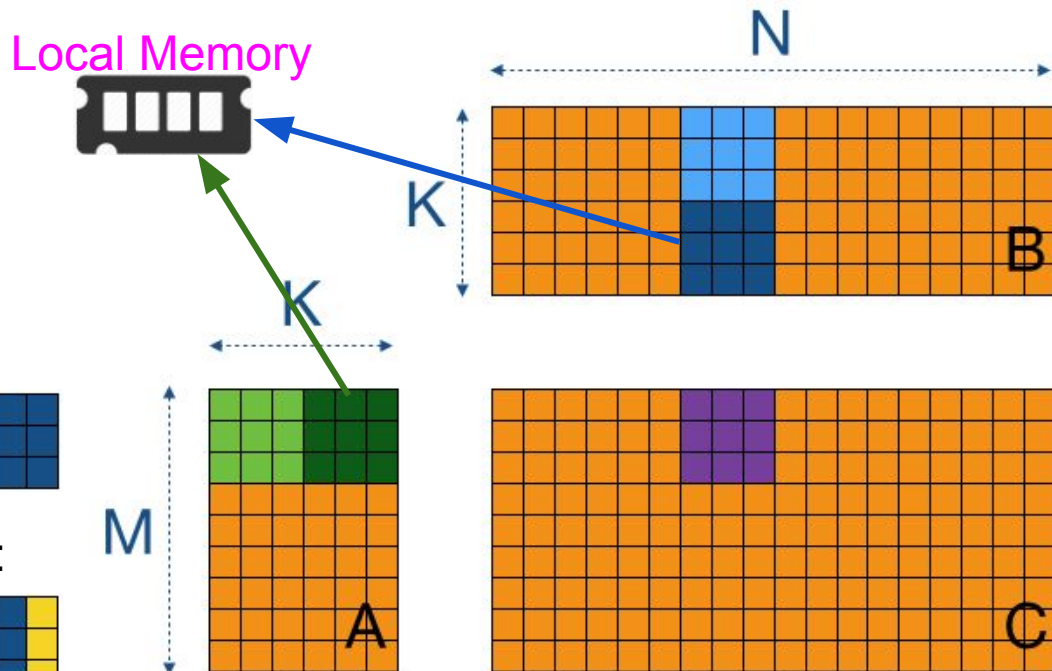
Какая пропорция?

$$\begin{bmatrix} \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \end{bmatrix} = \begin{bmatrix} \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \end{bmatrix} \times \begin{bmatrix} \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \end{bmatrix} + \begin{bmatrix} \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \end{bmatrix} \times \begin{bmatrix} \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \end{bmatrix}$$

Переиспользование данных:

$$\left\{ \begin{bmatrix} \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \end{bmatrix} \right\}_{32} = \begin{bmatrix} \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \end{bmatrix} \times \begin{bmatrix} \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \end{bmatrix} + \begin{bmatrix} \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \end{bmatrix} \times \begin{bmatrix} \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \end{bmatrix}$$

32



Умножение матриц

Сколько у нас вычислений?

$$O(N \times M \times K)$$

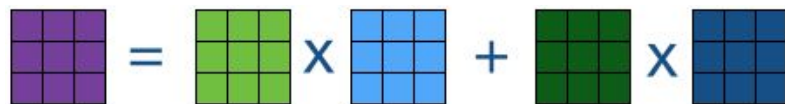
Сколько у нас чтений/записей?

$$O(N \times M \times K /$$

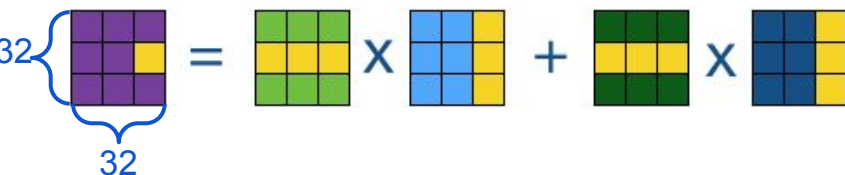
32)

Какая пропорция?

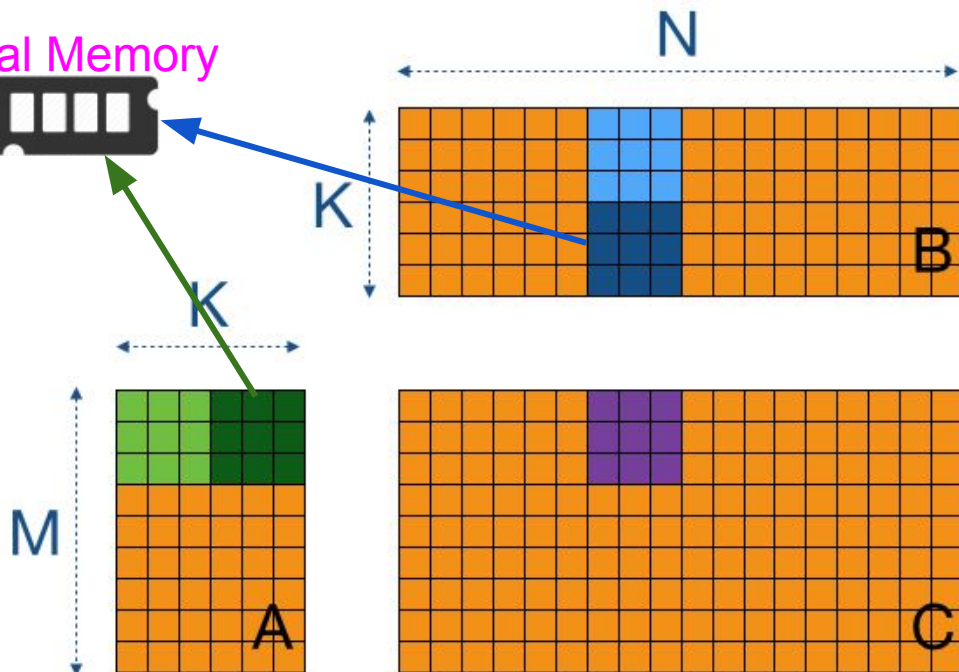
32:1


$$\begin{bmatrix} \text{green} \end{bmatrix} \times \begin{bmatrix} \text{blue} \end{bmatrix} + \begin{bmatrix} \text{green} \end{bmatrix} \times \begin{bmatrix} \text{blue} \end{bmatrix} = \begin{bmatrix} \text{purple} \end{bmatrix}$$

Переиспользование данных:


$$\begin{bmatrix} \text{purple} \end{bmatrix} = \begin{bmatrix} \text{green} \end{bmatrix} \times \begin{bmatrix} \text{blue} \end{bmatrix} + \begin{bmatrix} \text{green} \end{bmatrix} \times \begin{bmatrix} \text{blue} \end{bmatrix}$$

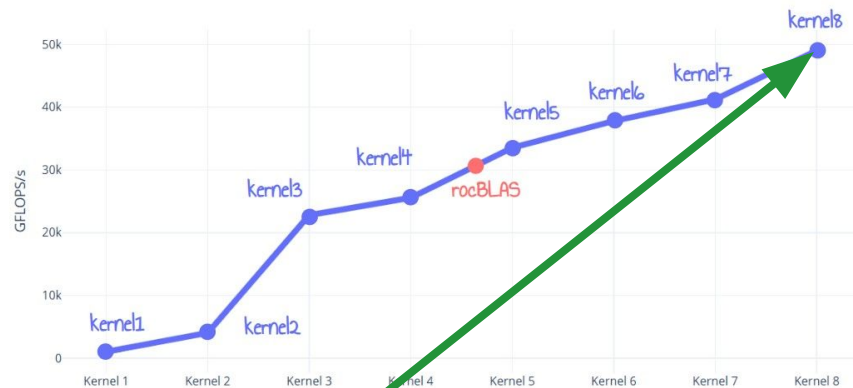
Local Memory



Умножение матриц

Неравная битва за гигафлопсы при умножении матриц
(хорошо описанная аналитика, профилирование, оптимизация):

- AMD RDNA3 - <https://seb-v.github.io/optimization/update/2025/01/20/Fast-GPU-Matrix-multiplication.html>
- NVIDIA Kepler - <https://cnugteren.github.io/tutorial/pages/page15.html>
- <https://siboehm.com/articles/22/CUDA-MMM>



Умножение матриц

Tensor Cores

	VRAM <i>TB/s</i>	FP 32 <i>TFlops</i>	FP 16 <i>TFlops</i>	FP 16 (tensor) <i>TFlops</i>
RTX 3090	0.93	29	29	142
RTX 4090	1.00	73	73	330
RTX 5090	1.79	105	105	419
Tesla H100	3.35	67	268 (4:1)	990 (16:1)



Умножение матриц

Tensor Cores

$$D = \begin{pmatrix} A_{0,0} & A_{0,1} & A_{0,2} & A_{0,3} \\ A_{1,0} & A_{1,1} & A_{1,2} & A_{1,3} \\ A_{2,0} & A_{2,1} & A_{2,2} & A_{2,3} \\ A_{3,0} & A_{3,1} & A_{3,2} & A_{3,3} \end{pmatrix} \begin{pmatrix} B_{0,0} & B_{0,1} & B_{0,2} & B_{0,3} \\ B_{1,0} & B_{1,1} & B_{1,2} & B_{1,3} \\ B_{2,0} & B_{2,1} & B_{2,2} & B_{2,3} \\ B_{3,0} & B_{3,1} & B_{3,2} & B_{3,3} \end{pmatrix} + \begin{pmatrix} C_{0,0} & C_{0,1} & C_{0,2} & C_{0,3} \\ C_{1,0} & C_{1,1} & C_{1,2} & C_{1,3} \\ C_{2,0} & C_{2,1} & C_{2,2} & C_{2,3} \\ C_{3,0} & C_{3,1} & C_{3,2} & C_{3,3} \end{pmatrix}$$

FP16 or FP32 FP16 FP16 or FP32

4x4



Умножение матриц

Tensor Cores

$$D = \begin{pmatrix} A_{0,0} & A_{0,1} & A_{0,2} & A_{0,3} \\ A_{1,0} & A_{1,1} & A_{1,2} & A_{1,3} \\ A_{2,0} & A_{2,1} & A_{2,2} & A_{2,3} \\ A_{3,0} & A_{3,1} & A_{3,2} & A_{3,3} \end{pmatrix} \begin{pmatrix} B_{0,0} & B_{0,1} & B_{0,2} & B_{0,3} \\ B_{1,0} & B_{1,1} & B_{1,2} & B_{1,3} \\ B_{2,0} & B_{2,1} & B_{2,2} & B_{2,3} \\ B_{3,0} & B_{3,1} & B_{3,2} & B_{3,3} \end{pmatrix} + \begin{pmatrix} C_{0,0} & C_{0,1} & C_{0,2} & C_{0,3} \\ C_{1,0} & C_{1,1} & C_{1,2} & C_{1,3} \\ C_{2,0} & C_{2,1} & C_{2,2} & C_{2,3} \\ C_{3,0} & C_{3,1} & C_{3,2} & C_{3,3} \end{pmatrix}$$

FP16 or FP32 FP16 FP16 4x4 FP16 or FP32

Tensor Cores (CUDA kernel)

$$D = \begin{pmatrix} A_{0,0} & A_{0,1} & A_{0,\dots} & A_{0,15} \\ A_{1,0} & A_{1,1} & A_{1,\dots} & A_{1,15} \\ A_{\dots,0} & A_{\dots,1} & A_{\dots,\dots} & A_{\dots,15} \\ A_{15,0} & A_{15,1} & A_{15,\dots} & A_{15,15} \end{pmatrix} \begin{pmatrix} B_{0,0} & B_{0,1} & B_{0,\dots} & B_{0,15} \\ B_{1,0} & B_{1,1} & B_{1,\dots} & B_{1,15} \\ B_{\dots,0} & B_{\dots,1} & B_{\dots,\dots} & B_{\dots,15} \\ B_{15,0} & B_{15,1} & B_{15,\dots} & B_{15,15} \end{pmatrix} + \begin{pmatrix} C_{0,0} & C_{0,1} & C_{0,\dots} & C_{0,15} \\ C_{1,0} & C_{1,1} & C_{1,\dots} & C_{1,15} \\ C_{\dots,0} & C_{\dots,1} & C_{\dots,\dots} & C_{\dots,15} \\ C_{15,0} & C_{15,1} & C_{15,\dots} & C_{15,15} \end{pmatrix}$$

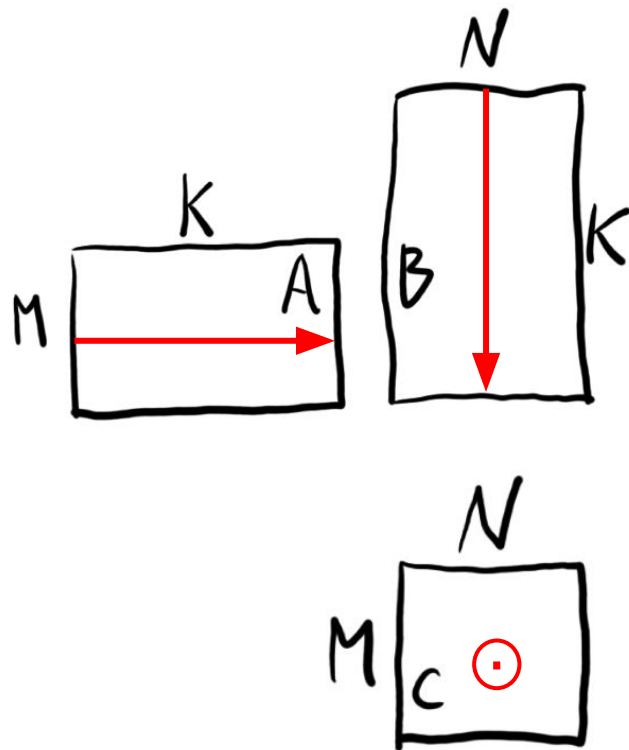
FP16 or FP32 FP16 FP16 16x16 FP16 or FP32



69 // Performs an MxNxK GEMM ($C = \alpha * A * B + \beta * C$)

75 `__global__ void wmma_example(half *a, half *b, float *c, int M, int N, int K, float alpha, float beta) {`

$$C = \alpha A * B + \beta C$$



```
69 // Performs an MxNxK GEMM (C=alpha*A*B + beta*C)
```

```
75 __global__ void wmma_example(half *a, half *b, float *c, int M, int N, int K, float alpha, float beta) {
```

```
85 // Declare the fragments
```

```
86 wmma::fragment<wmma::matrix_a, WMMA_M, WMMA_N, WMMA_K, half, wmma::col_major> a_frag;
```

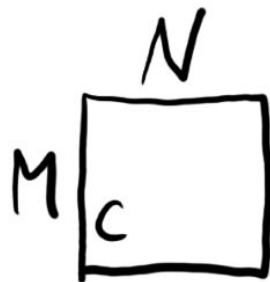
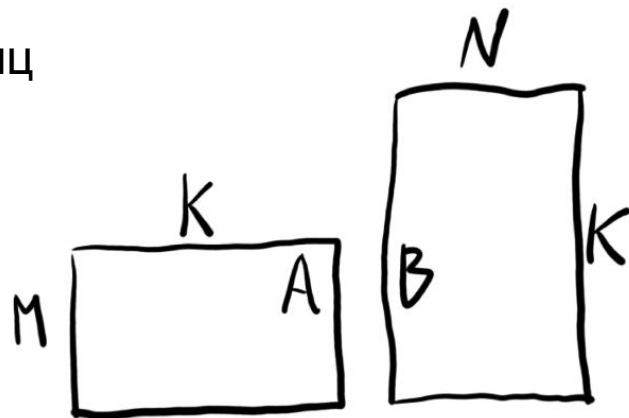
```
87 wmma::fragment<wmma::matrix_b, WMMA_M, WMMA_N, WMMA_K, half, wmma::col_major> b_frag;
```

```
88 wmma::fragment<wmma::accumulator, WMMA_M, WMMA_N, WMMA_K, float> acc_frag;
```

```
89 wmma::fragment<wmma::accumulator, WMMA_M, WMMA_N, WMMA_K, float> c_frag;
```

Общие на весь warp 16x16 фрагменты матриц
WMMA = Warp Matrix Multiply-Accumulate

$$C = \alpha A * B + \beta C$$



69 // Performs an MxNxK GEMM (C=alpha*A*B + beta*C)

75 `__global__ void wmma_example(half *a, half *b, float *c, int M, int N, int K, float alpha, float beta) {`

85 `// Declare the fragments`

86 `wmma::fragment<wmma::matrix_a, WMMA_M, WMMA_N, WMMA_K, half, wmma::col_major> a_frag;`

87 `wmma::fragment<wmma::matrix_b, WMMA_M, WMMA_N, WMMA_K, half, wmma::col_major> b_frag;`

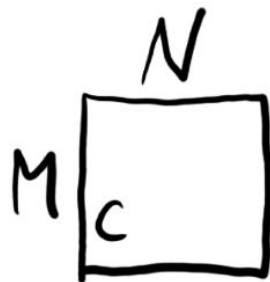
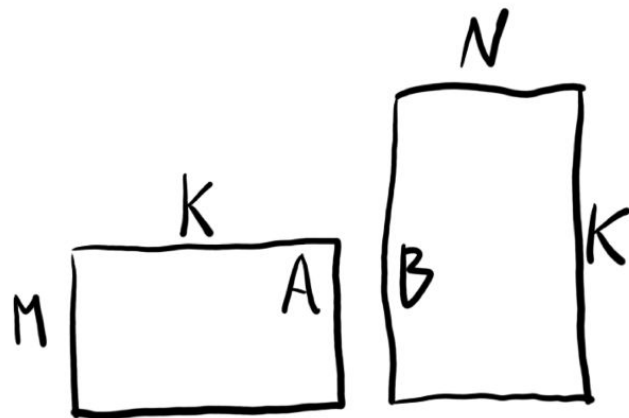
88 `wmma::fragment<wmma::accumulator, WMMA_M, WMMA_N, WMMA_K, float> acc_frag;`

89 `wmma::fragment<wmma::accumulator, WMMA_M, WMMA_N, WMMA_K, float> c_frag;`

91 `wmma::fill_fragment(acc_frag, 0.0f);`

 **16x16**
acc_frag

$$C = \alpha A * B + \beta C$$

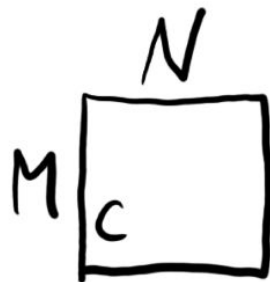
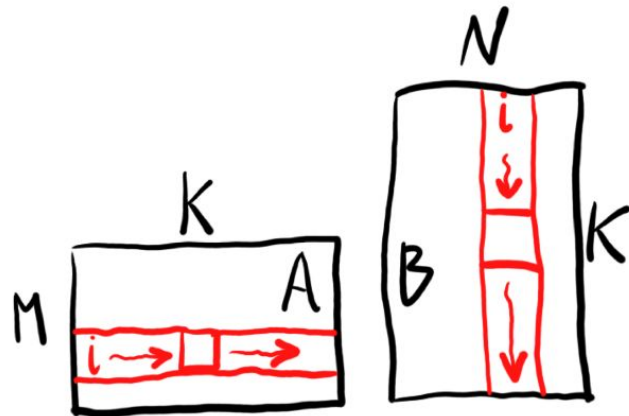


69 // Performs an MxNxK GEMM (C=alpha*A*B + beta*C)

```
75 __global__ void wmma_example(half *a, half *b, float *c, int M, int N, int K, float alpha, float beta) {  
85     // Declare the fragments  
86     wmma::fragment<wmma::matrix_a, WMMA_M, WMMA_N, WMMA_K, half, wmma::col_major> a_frag;  
87     wmma::fragment<wmma::matrix_b, WMMA_M, WMMA_N, WMMA_K, half, wmma::col_major> b_frag;  
88     wmma::fragment<wmma::accumulator, WMMA_M, WMMA_N, WMMA_K, float> acc_frag;  
89     wmma::fragment<wmma::accumulator, WMMA_M, WMMA_N, WMMA_K, float> c_frag;  
  
91     wmma::fill_fragment(acc_frag, 0.0f);  
94     for (int i = 0; i < K; i += WMMA_K) {  
95         int aRow = warpM * WMMA_M;  
96         int aCol = i;  
  
97  
98         int bRow = i;  
99         int bCol = warpN * WMMA_N;  
  
100  
101         // Bounds checking  
102         if (aRow < M && aCol < K && bRow < K && bCol < N) {  
103             // Load the inputs  
104             wmma::load_matrix_sync(a_frag, a + aRow + aCol * lda, lda);  
105             wmma::load_matrix_sync(b_frag, b + bRow + bCol * ldb, ldb);
```

16x16
acc_frag


$$C = \alpha A * B + \beta C$$

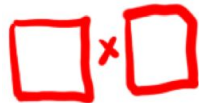


```

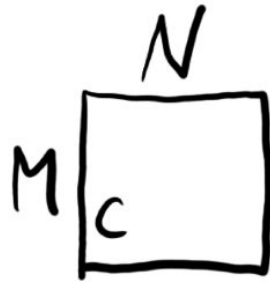
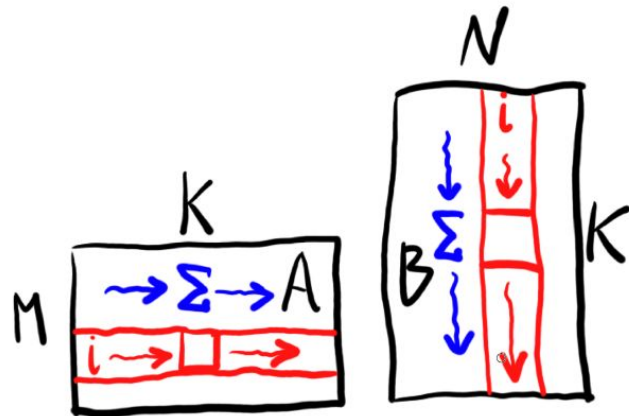
69 // Performs an MxNxK GEMM (C=alpha*A*B + beta*C)
71
72 __global__ void wmma_example(half *a, half *b, float *c, int M, int N, int K, float alpha, float beta) {
73
74     // Declare the fragments
75     wmma::fragment<wmma::matrix_a, WMMA_M, WMMA_N, WMMA_K, half, wmma::col_major> a_frag;
76     wmma::fragment<wmma::matrix_b, WMMA_M, WMMA_N, WMMA_K, half, wmma::col_major> b_frag;
77     wmma::fragment<wmma::accumulator, WMMA_M, WMMA_N, WMMA_K, float> acc_frag;
78     wmma::fragment<wmma::accumulator, WMMA_M, WMMA_N, WMMA_K, float> c_frag;
79
80     wmma::fill_fragment(acc_frag, 0.0f);
81
82     for (int i = 0; i < K; i += WMMA_K) {
83         int aRow = warpM * WMMA_M;
84         int aCol = i;
85
86         int bRow = i;
87         int bCol = warpN * WMMA_N;
88
89         // Bounds checking
90         if (aRow < M && aCol < K && bRow < K && bCol < N) {
91             // Load the inputs
92             wmma::load_matrix_sync(a_frag, a + aRow + aCol * lda, lda);
93             wmma::load_matrix_sync(b_frag, b + bRow + bCol * ldb, ldb);
94
95             // Perform the matrix multiplication
96             wmma::mma_sync(acc_frag, a_frag, b_frag, acc_frag);
97         }
98     }
99 }

```

 16x16
acc_frag

16x16
 → acc_frag

$$C = \alpha A * B + \beta C$$



69 // Performs an MxNxK GEMM (C=alpha*A*B + beta*C)

75 `__global__ void wmma_example(half *a, half *b, float *c, int M, int N, int K, float alpha, float beta) {`

85 `// Declare the fragments`

86 `wmma::fragment<wmma::matrix_a, WMMA_M, WMMA_N, WMMA_K, half, wmma::col_major> a_frag;`

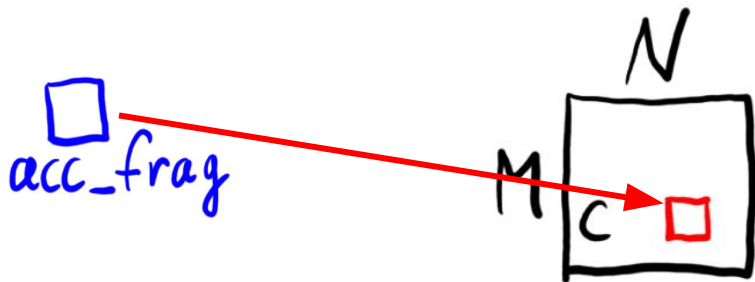
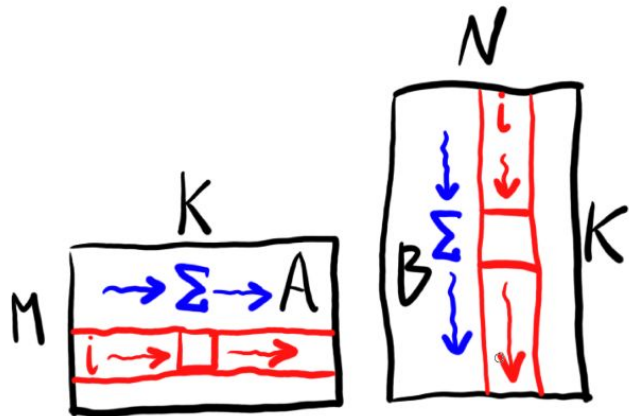
87 `wmma::fragment<wmma::matrix_b, WMMA_M, WMMA_N, WMMA_K, half, wmma::col_major> b_frag;`

88 `wmma::fragment<wmma::accumulator, WMMA_M, WMMA_N, WMMA_K, float> acc_frag;`

89 `wmma::fragment<wmma::accumulator, WMMA_M, WMMA_N, WMMA_K, float> c_frag;`

.....

$$C = \alpha A * B + \beta C$$

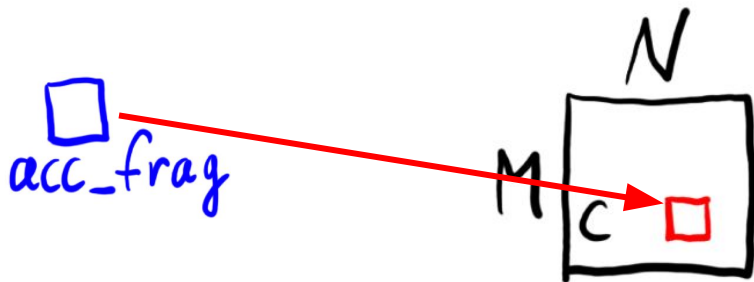
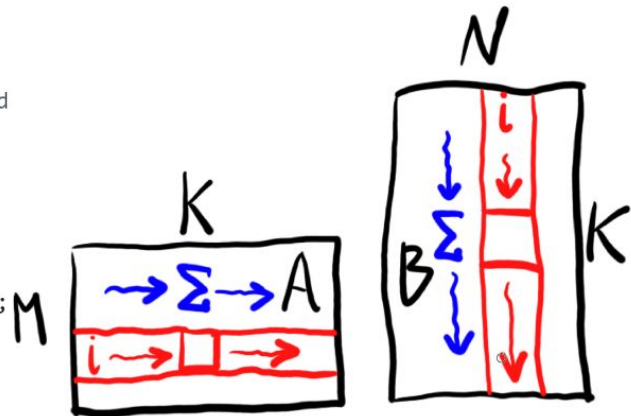


```

69 // Performs an MxNxK GEMM (C=alpha*A*B + beta*C)
71 __global__ void wmma_example(half *a, half *b, float *c, int M, int N, int K, float alpha, float beta) {
73     // Declare the fragments
74     wmma::fragment<wmma::matrix_a, WMMA_M, WMMA_N, WMMA_K, half, wmma::col_major> a_frag;
75     wmma::fragment<wmma::matrix_b, WMMA_M, WMMA_N, WMMA_K, half, wmma::col_major> b_frag;
76     wmma::fragment<wmma::accumulator, WMMA_M, WMMA_N, WMMA_K, float> acc_frag;
77     wmma::fragment<wmma::accumulator, WMMA_M, WMMA_N, WMMA_K, float> c_frag;
78
79     .....
81
82     // Load in the current value of c, scale it by beta, and add this our result scaled
83     int cRow = warpM * WMMA_M;
84     int cCol = warpN * WMMA_N;
85
86     if (cRow < M && cCol < N) {
87         wmma::load_matrix_sync(c_frag, c + cRow + cCol * ldc, ldc, wmma::mem_col_major);
88
89 #pragma unroll
90         for(int i=0; i < c_frag.num_elements; i++) {
91             c_frag.x[i] = alpha * acc_frag.x[i] + beta * c_frag.x[i];
92         }
93
94         // Store the output
95         wmma::store_matrix_sync(c + cRow + cCol * ldc, c_frag, ldc, wmma::mem_col_major);
96     }
97 }

```

$$C = \alpha A * B + \beta C$$



Умножение матриц

Tensor Cores

<https://developer.nvidia.com/blog/programming-tensor-cores-cuda-9/>

<https://github.com/NVIDIA-developer-blog/code-samples/blob/master/posts/tensor-cores/simpleTensorCoreGEMM.cu>

https://github.com/NVIDIA/cutlass/blob/main/examples/00_basic_gemm/basic_gemm.cu

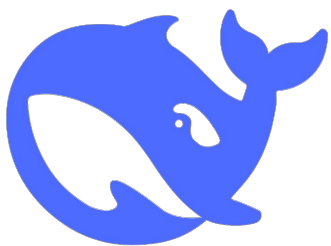
<https://developer.nvidia.com/blog/cutlass-linear-algebra-cuda>

<https://docs.nvidia.com/cuda/cuda-c-programming-guide/>

Неравная битва за гигафлопсы при умножении матриц
(хорошо описанная аналитика, профилирование, оптимизация):

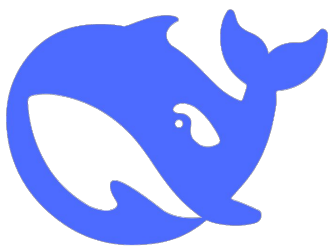
- AMD RDNA3 - <https://seb-v.github.io/optimization/update/2025/01/20/Fast-GPU-Matrix-multiplication.html>
- NVIDIA Kepler - <https://cnugeteren.github.io/tutorial/pages/page15.html>
- <https://siboehm.com/articles/22/CUDA-MMM>





DeerSeek: x2 ускорение обучения (fp16 → fp8)

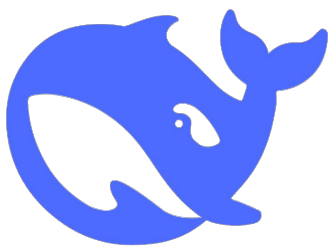
	sign	exponent						mantissa									
FP16	0	0	1	1	0	1	1	0	0	1	0	1	0	0	1	1	= 0.395264
BF16	0	0	1	1	1	1	1	0	1	1	0	0	1	0	1	0	= 0.394531
FP8 E4M3	0	0	1	0	1	1	0	1									= 0.40625
FP8 E5M2	0	0	1	1	0	1	1	0									= 0.375



DeerSeek: x2 ускорение обучения (fp16 → fp8)

NVIDIA H100 (почти то же что и H800*):

- Memory bandwidth: **3.35 TB/sec**
- FP 32: **67 TFlops**
- FP 16: **268 TFlops**
- FP 32 (tensor): **495 TFlops**
- FP 16 (tensor): **990 TFlops**
- FP 8 (tensor): **1979 TFlops**

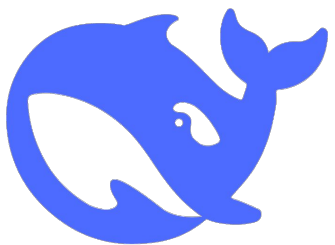


DeerSeek: x2 ускорение обучения (fp16 → fp8)

NVIDIA H100 (почти то же что и H800*):

- Memory bandwidth: **3.35 TB/sec**
- FP 32: **67 TFlops**
- FP 16: **268 TFlops**
- FP 32 (tensor): **495 TFlops**
- FP 16 (tensor): **990 TFlops**
- FP 8 (tensor): **1979 TFlops**

x2



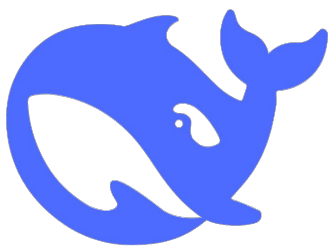
DeerSeek: x2 ускорение обучения (fp16 → fp8)

NVIDIA H100 (почти то же что и H800*):

- Memory bandwidth: **3.35 TB/sec**
- FP 32: **67 TFlops**
- FP 16: **268 TFlops**
- FP 32 (tensor): **495 TFlops**
- FP 16 (tensor): **990 TFlops**
- FP 8 (tensor): **1979 TFlops**

Хватит ли пропускной способности памяти чтобы насытить ALU (tensor cores)?

x2



DeerSeek: x2 ускорение обучения (fp16 → fp8)

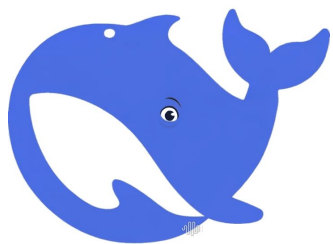
NVIDIA H100 (почти то же что и H800*):

- Memory bandwidth: **3.35 TB/sec**
- FP 32: **67 TFlops**
- FP 16: **268 TFlops**
- FP 32 (tensor): **495 TFlops**
- FP 16 (tensor): **990 TFlops**
- FP 8 (tensor): **1979 TFlops**

x2

Какие риски?

Почему так не делают все?



DeepSeek: fine-grained fp8 quantization

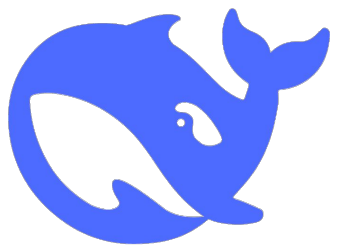


NVIDIA H100 (почти то же что и H800*):

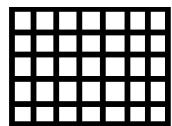
- Memory bandwidth: **3.35 TB/sec**
- FP 32: **67 TFlops**
- FP 16: **268 TFlops**
- FP 32 (tensor): **495 TFlops**
- FP 16 (tensor): **990 TFlops**
- FP 8 (tensor): **1979 TFlops**
- DeepGEMM достиг **1550 TFlops** на H800!

DeepSeek-V3 Technical Report - <https://arxiv.org/abs/2412.19437>

Репозиторий - <https://github.com/deepseek-ai/DeepGEMM> (GEMM - General Matrix Multiplications)



DeepSeek: fine-grained fp8 quantization

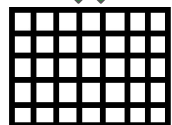


fp32

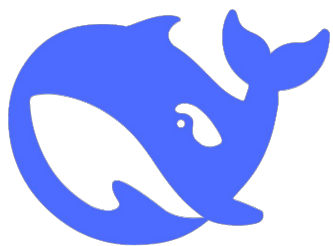


Scaling Factor

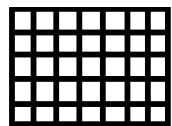
fp32



fp8



DeepSeek: fine-grained fp8 quantization



fp32

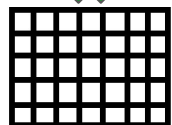
Input Values



Scaling Factor

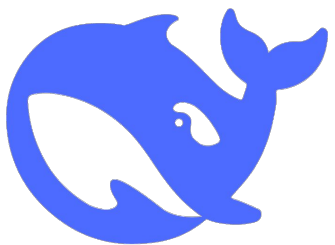
fp32

$\text{absmax}(\text{Input Values}) / 448$

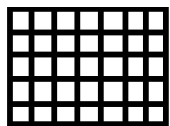


fp8

Input Values / Scaling Factor



DeepSeek: fine-grained fp8 quantization



fp32

Input Values

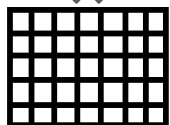


Scaling Factor

fp32

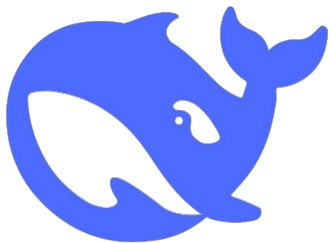
$\text{absmax}(\text{Input Values}) / 448$

Что это за волшебное число? 🦄

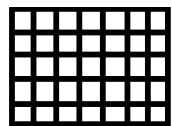


fp8

Input Values / Scaling Factor



DeepSeek: fine-grained fp8 quantization



fp32

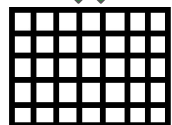
Input Values



Scaling Factor

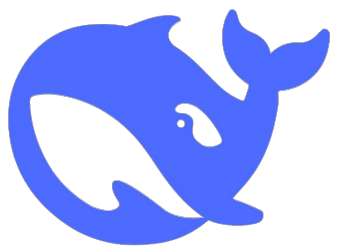
fp32

$\text{absmax}(\text{Input Values}) / 448 = \text{FP8_MAX}$



fp8

$\text{Input Values} / \text{Scaling Factor} \in [-448; +448]$

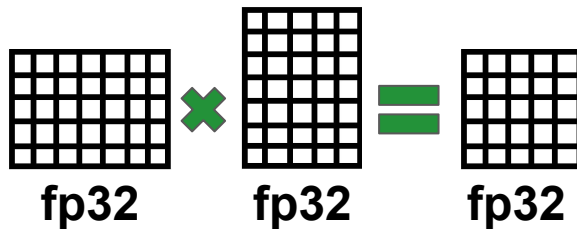


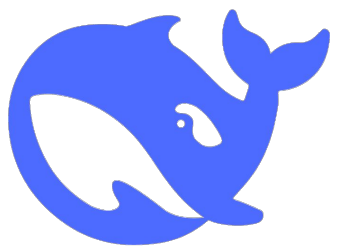
DeepSeek: fine-grained fp8 quantization



Scaling Factor

fp32





DeepSeek: fine-grained fp8 quantization

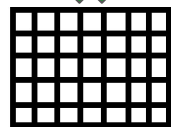


fp32

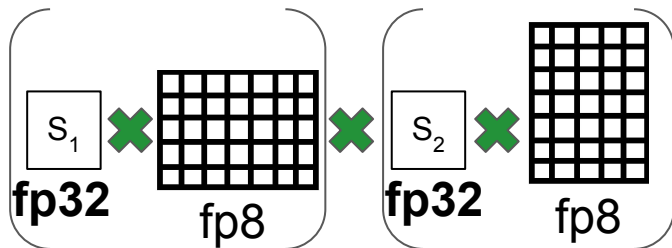
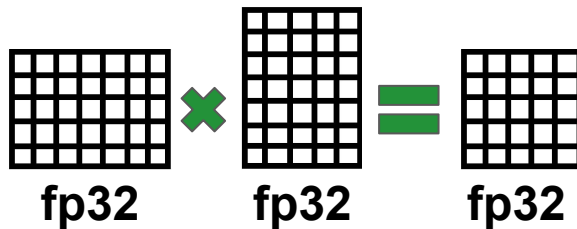


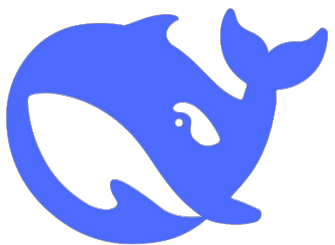
Scaling Factor

fp32

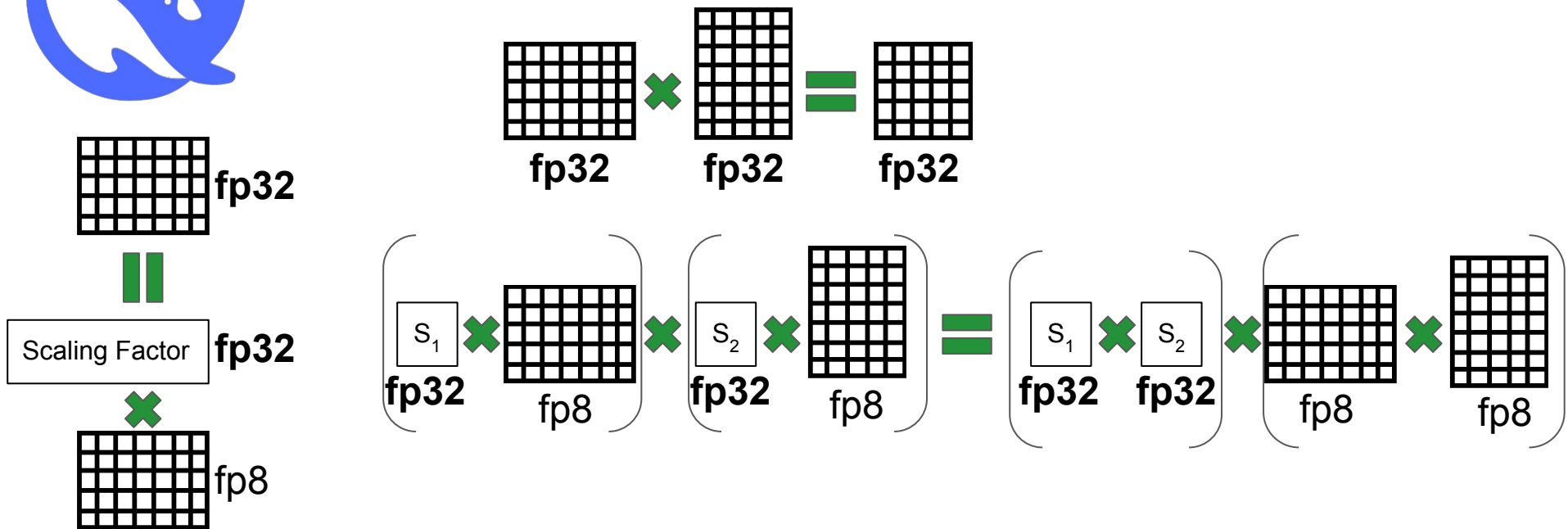


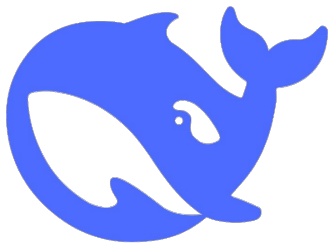
fp8





DeepSeek: fine-grained fp8 quantization





DeepSeek: fine-grained fp8 quantization

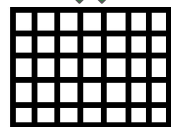


fp32

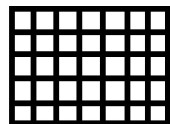


Scaling Factor

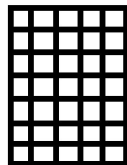
fp32



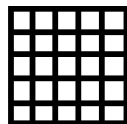
fp8



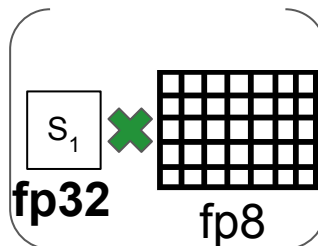
fp32



fp32

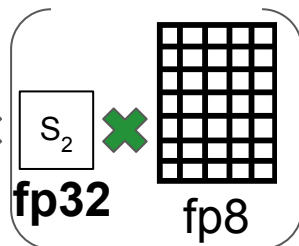


fp32



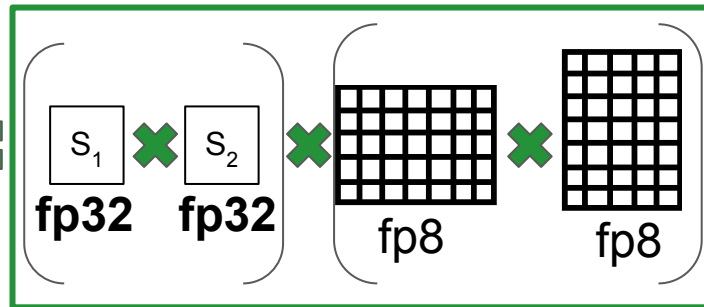
fp32

fp8

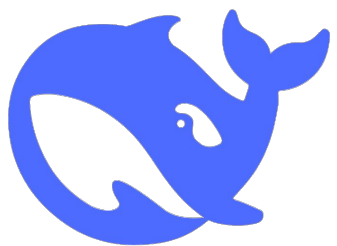


fp32

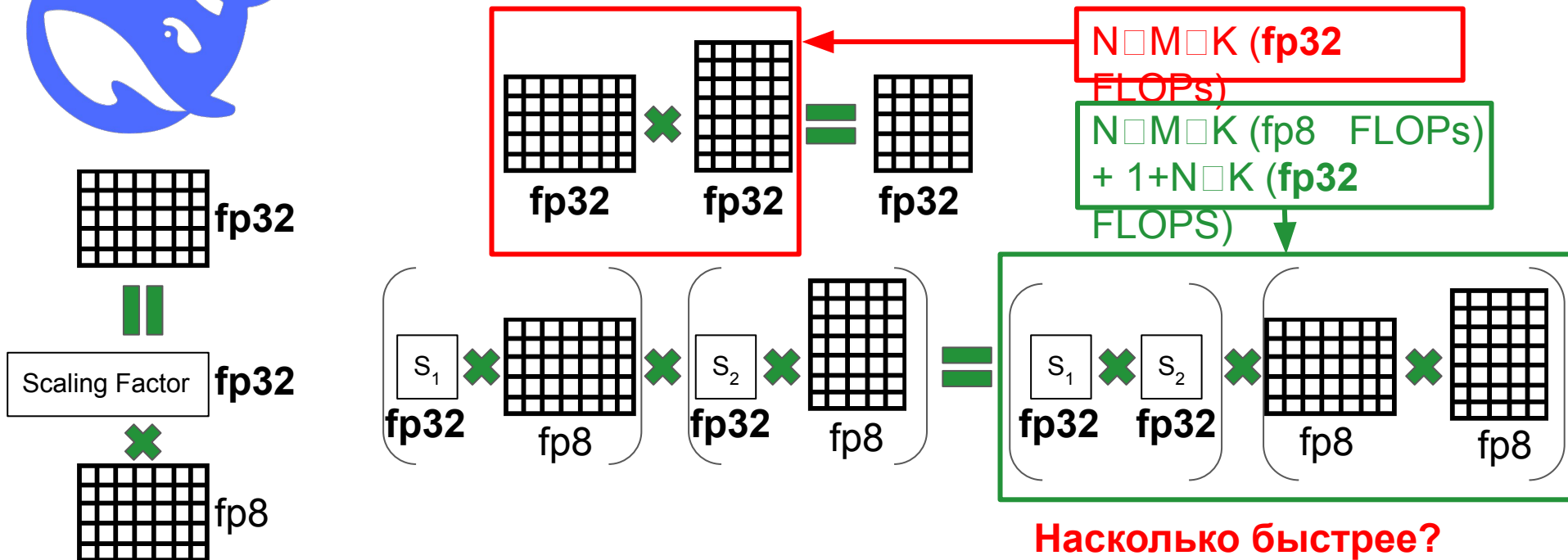
fp8

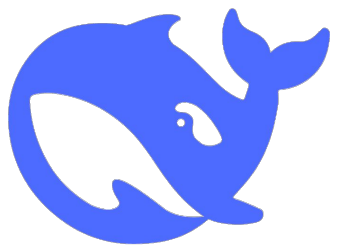


Насколько быстрее?



DeepSeek: fine-grained fp8 quantization





DeepSeek: fine-grained fp8 quantization

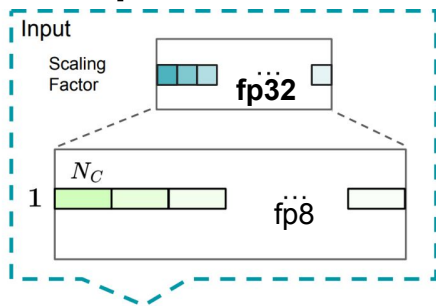


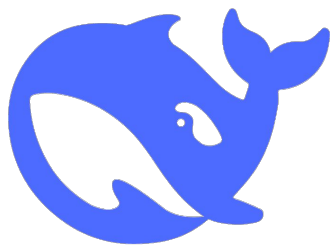
Scaling Factor

fp32



fp8



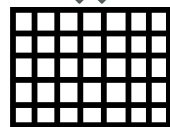


DeepSeek: fine-grained fp8 quantization

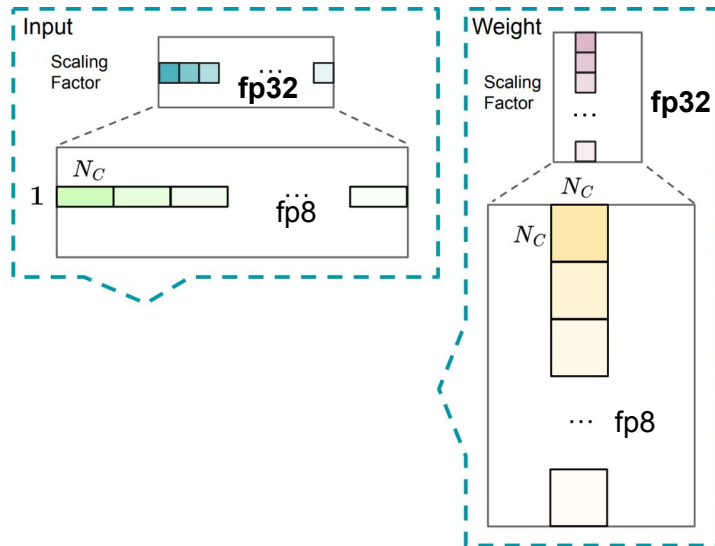


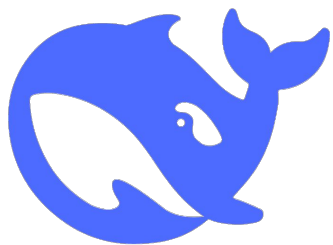
Scaling Factor

fp32



fp8

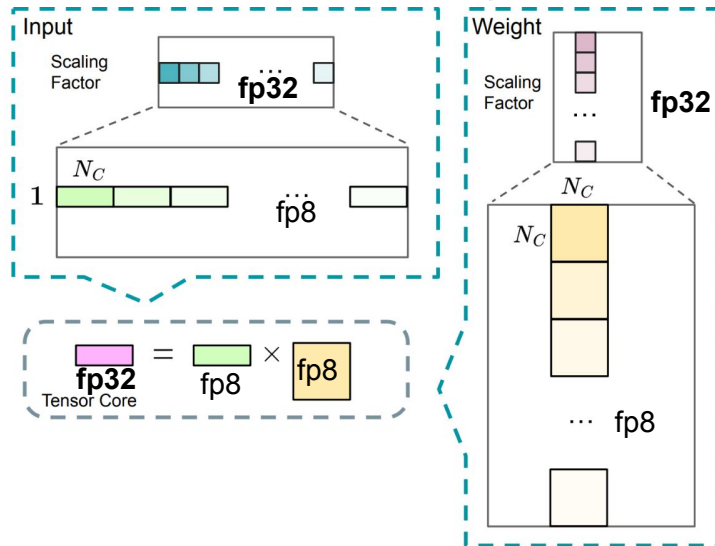


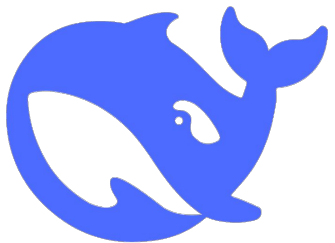


DeepSeek: fine-grained fp8 quantization

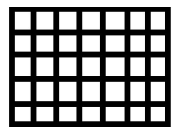


Scaling Factor **fp32**





DeepSeek: fine-grained fp8 quantization

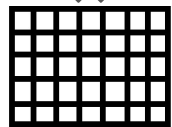


fp32

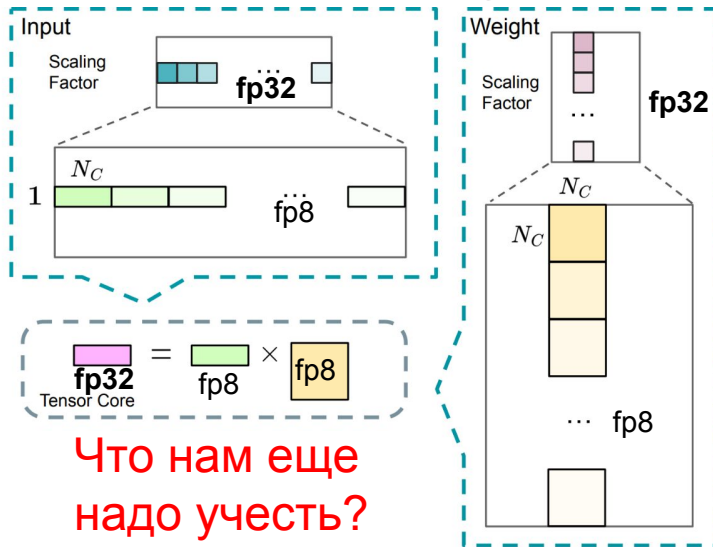


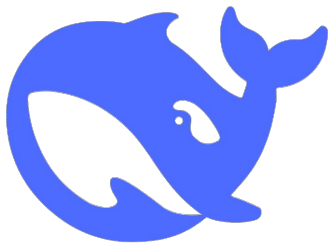
Scaling Factor

fp32

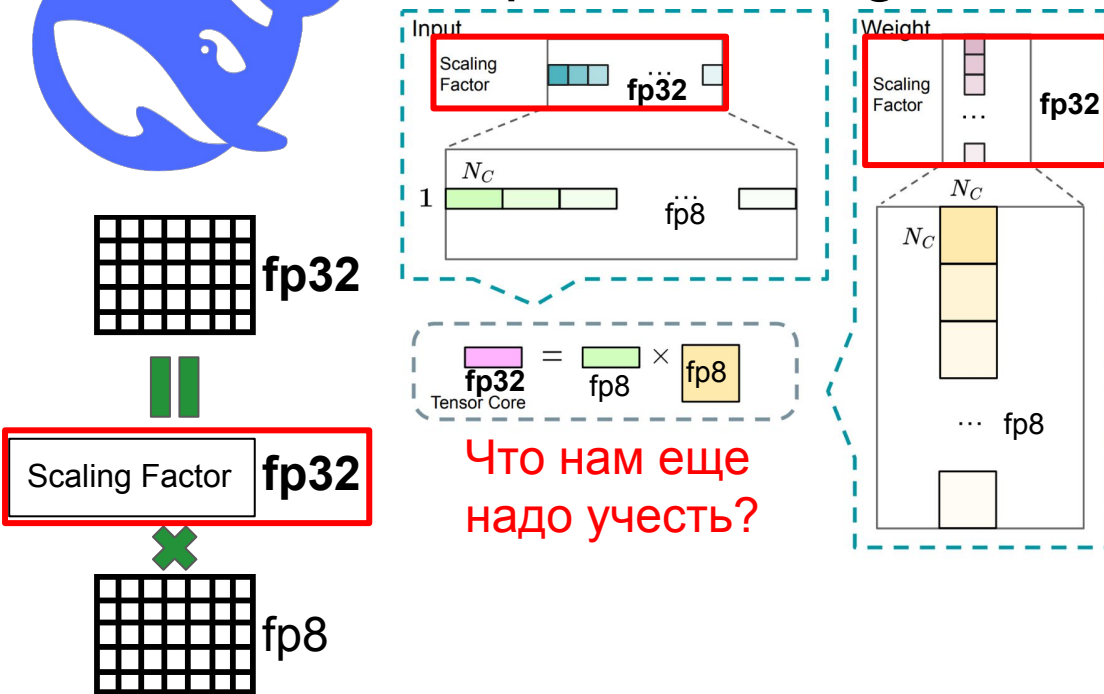


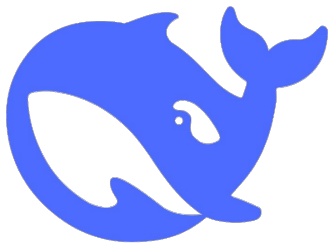
fp8



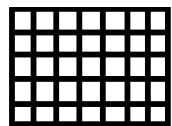


DeepSeek: fine-grained fp8 quantization





DeepSeek: fine-grained fp8 quantization

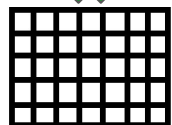


fp32

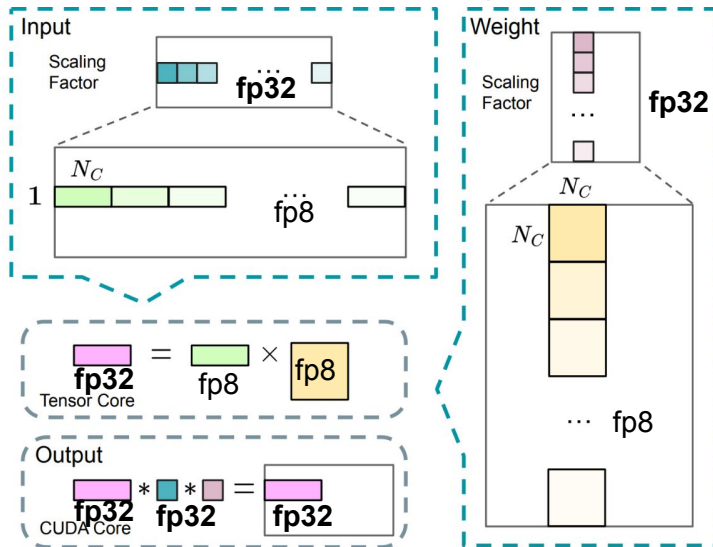


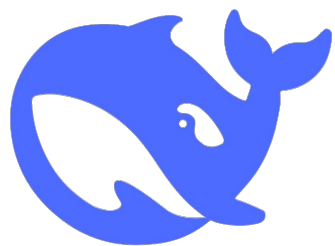
Scaling Factor

fp32



fp8

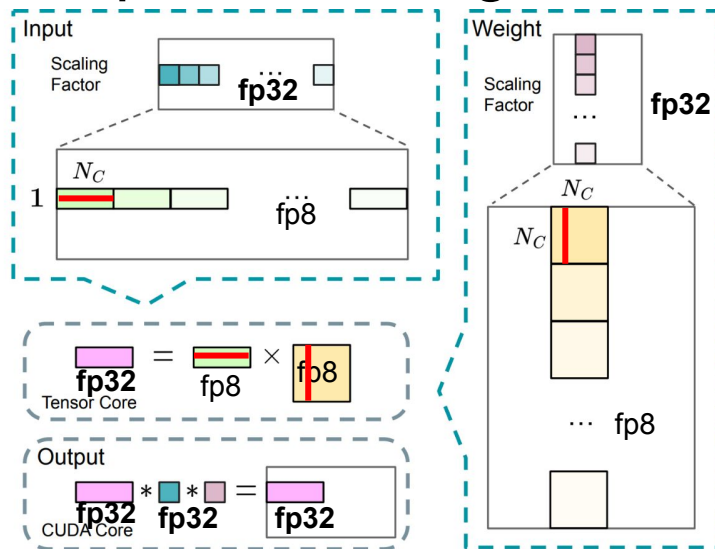




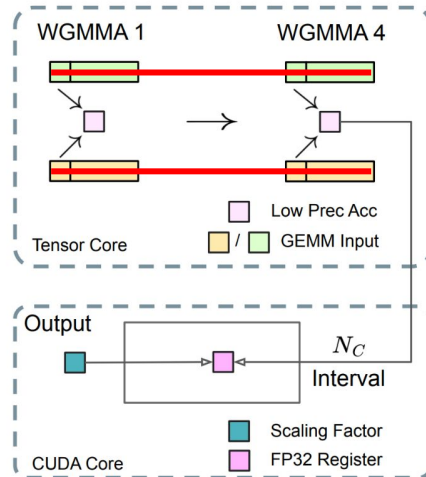
DeepSeek: fine-grained fp8 quantization



Scaling Factor fp32

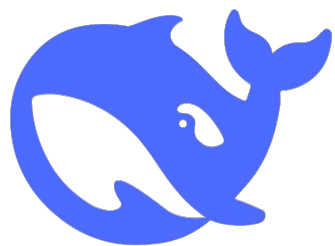


(a) Fine-grained quantization



(b) Increasing accumulation precision

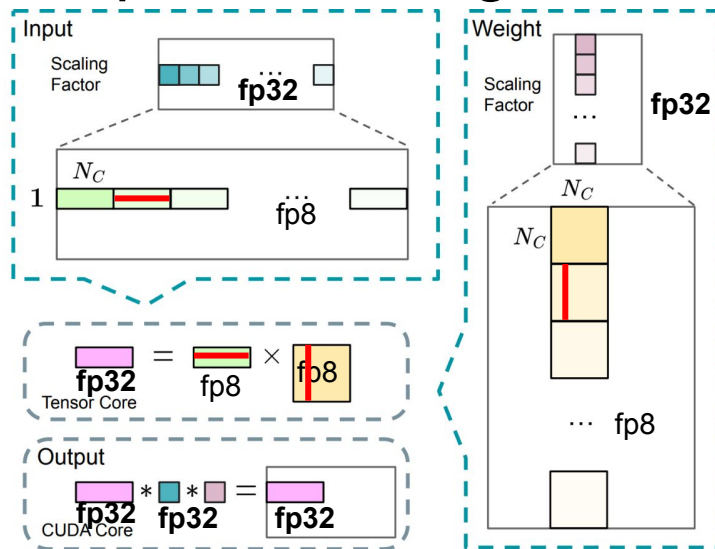
Figure 7 | (a) We propose a fine-grained quantization method to mitigate quantization errors caused by feature outliers; for illustration simplicity, only Fprop is illustrated. (b) In conjunction with our quantization strategy, we improve the FP8 GEMM precision by promoting to CUDA Cores at an interval of $N_C = 128$ elements MMA for the high-precision accumulation.



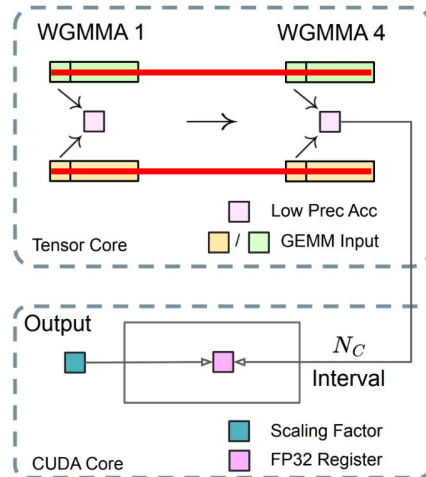
DeepSeek: fine-grained fp8 quantization



Scaling Factor fp32

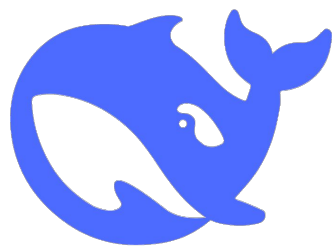


(a) Fine-grained quantization



(b) Increasing accumulation precision

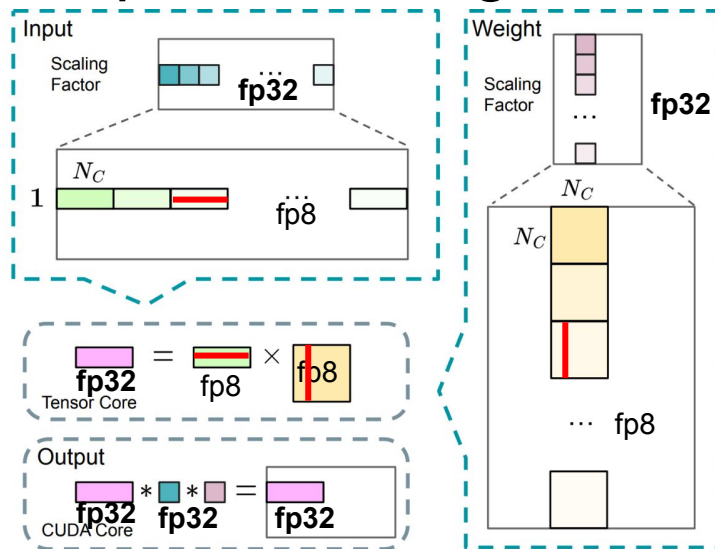
Figure 7 | (a) We propose a fine-grained quantization method to mitigate quantization errors caused by feature outliers; for illustration simplicity, only Fprop is illustrated. (b) In conjunction with our quantization strategy, we improve the FP8 GEMM precision by promoting to CUDA Cores at an interval of $N_C = 128$ elements MMA for the high-precision accumulation.



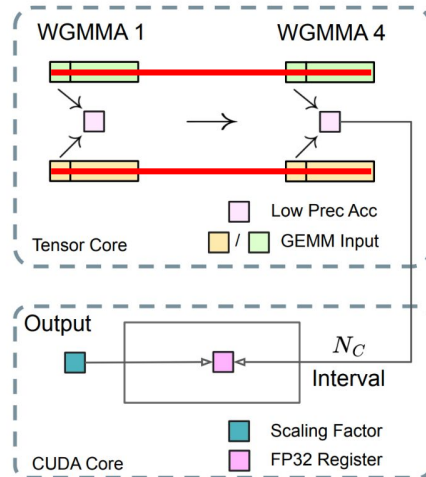
DeepSeek: fine-grained fp8 quantization



Scaling Factor fp32

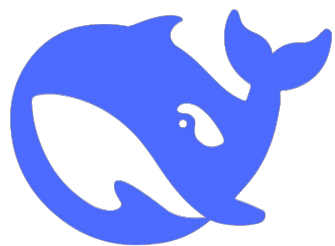


(a) Fine-grained quantization



(b) Increasing accumulation precision

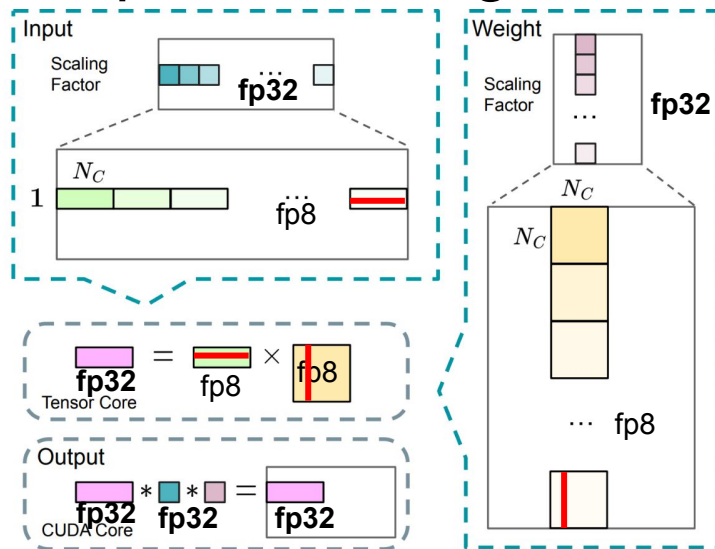
Figure 7 | (a) We propose a fine-grained quantization method to mitigate quantization errors caused by feature outliers; for illustration simplicity, only Fprop is illustrated. (b) In conjunction with our quantization strategy, we improve the FP8 GEMM precision by promoting to CUDA Cores at an interval of $N_C = 128$ elements MMA for the high-precision accumulation.



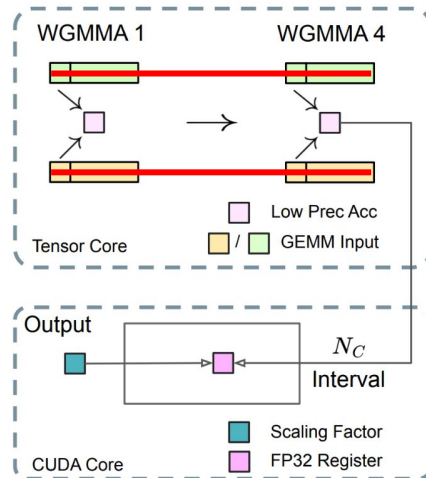
DeepSeek: fine-grained fp8 quantization



Scaling Factor fp32

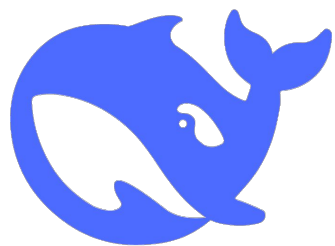


(a) Fine-grained quantization



(b) Increasing accumulation precision

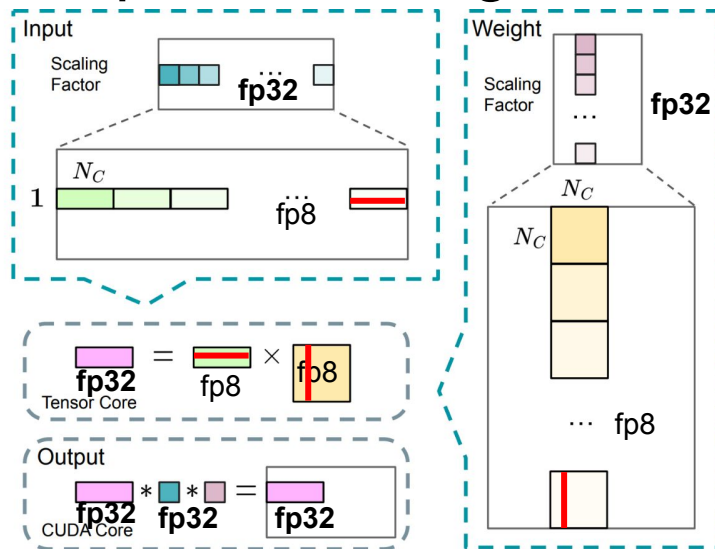
Figure 7 | (a) We propose a fine-grained quantization method to mitigate quantization errors caused by feature outliers; for illustration simplicity, only Fprop is illustrated. (b) In conjunction with our quantization strategy, we improve the FP8 GEMM precision by promoting to CUDA Cores at an interval of $N_C = 128$ elements MMA for the high-precision accumulation.



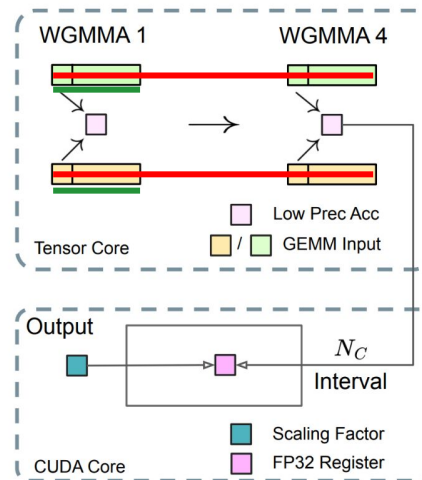
DeepSeek: fine-grained fp8 quantization



Scaling Factor fp32



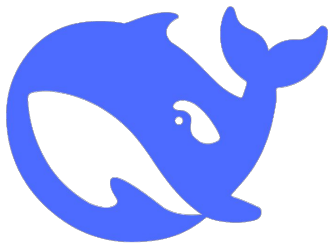
(a) Fine-grained quantization



(b) Increasing accumulation precision

Зачем 4xWGMMMA?
Почему не 1xWGMMMA?

Figure 7 | (a) We propose a fine-grained quantization method to mitigate quantization errors caused by feature outliers; for illustration simplicity, only Fprop is illustrated. (b) In conjunction with our quantization strategy, we improve the FP8 GEMM precision by promoting to CUDA Cores at an interval of $N_C = 128$ elements MMA for the high-precision accumulation.



DeepSeek: fine-grained fp8 quantization

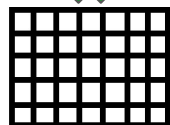


fp32

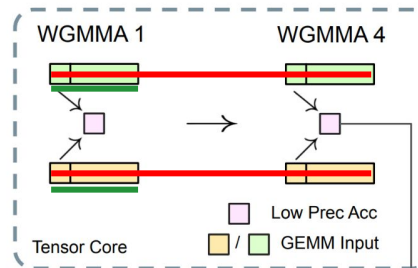


Scaling Factor

fp32

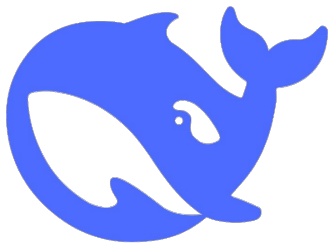


fp8



Зачем 4xWGMMAs?
Почему не 1xWGMMAs?

It is worth noting that this modification reduces the WGMMAs (Warpgroup-level Matrix Multiply-Accumulate) instruction issue rate for a single warpgroup. However, on the H800 architecture, it is typical for two WGMMAs to persist concurrently: while one warpgroup performs the promotion operation, the other is able to execute the MMA operation. This design enables overlapping of the two operations, maintaining high utilization of Tensor Cores. Based on our experiments, setting $N_C = 128$ elements, equivalent to 4 WGMMAs, represents the minimal accumulation interval that can significantly improve precision without introducing substantial overhead.



DeepSeek: fine-grained fp8 quantization

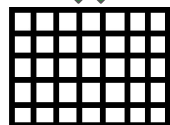


fp32

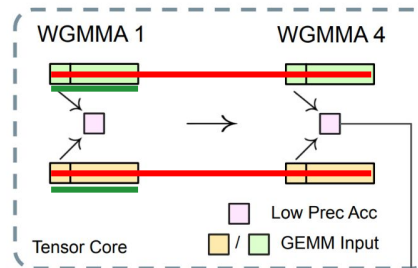


Scaling Factor

fp32



fp8

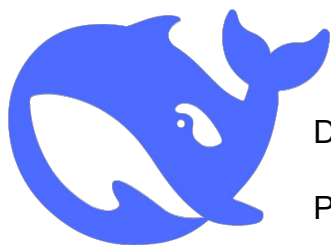


Зачем 4xWGMMMA?

Почему не 1xWGMMMA?

It is worth noting that this modification reduces the WGMMMA (Warpgroup-level Matrix Multiply-Accumulate) instruction issue rate for a single warpgroup. However, on the H800 architecture, it is typical for two WGMMMA to persist concurrently: while one warpgroup performs the promotion operation, the other is able to execute the MMA operation. This design enables overlapping of the two operations, maintaining high utilization of Tensor Cores. Based on our experiments, setting $N_C = 128$ elements, equivalent to 4 WGMMAs, represents the minimal accumulation interval that can significantly improve precision without introducing substantial overhead.

Не знаю почему нужно делать 4xWGMMMA + PROMOTION
вместо 4x(WGMMMA + PROMOTION)



DeepSeek: x2 ускорение обучения (fp16 → fp8)

DeepSeek-V3 Technical Report - <https://arxiv.org/abs/2412.19437>

Репозиторий - <https://github.com/deepseek-ai/DeepGEMM> (GEMM - General Matrix Multiplications)

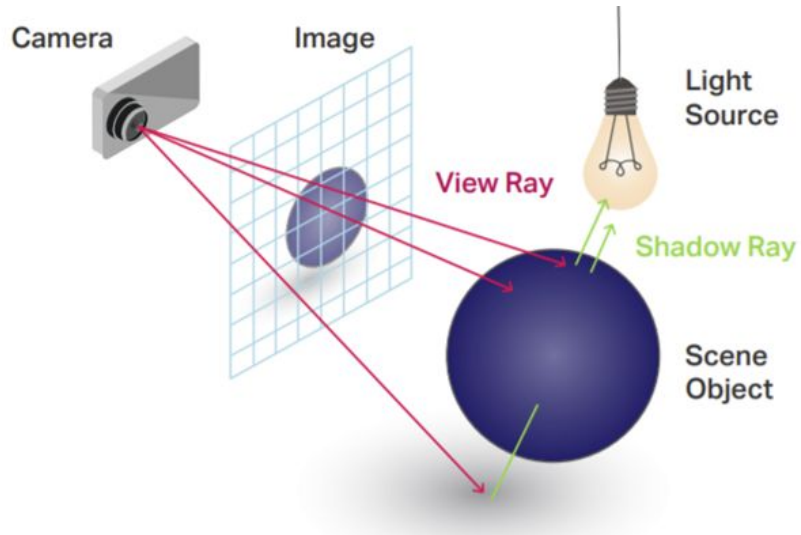
[Fast Matrix-Multiplication with WGMMMA on NVIDIA® Hopper™ GPUs](#)

https://docs.nvidia.com/deeplearning/transformer-engine/user-guide/examples/fp8_primer.html

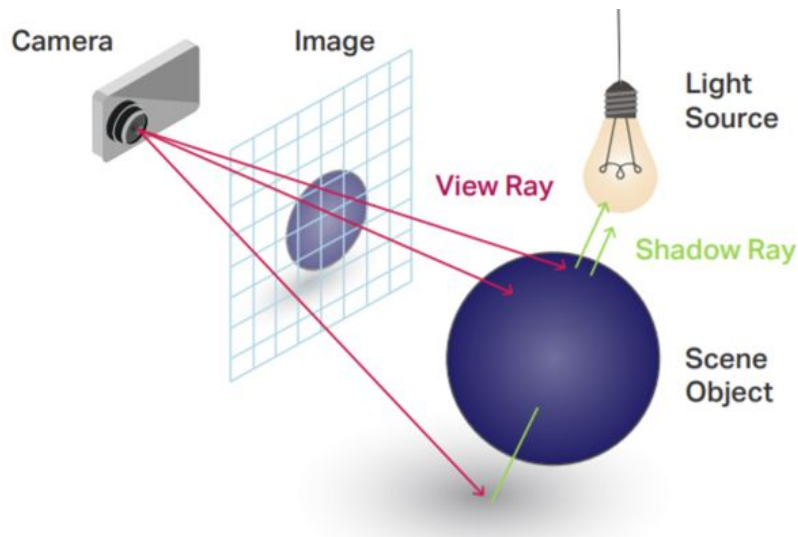
Глава 8: Ray Tracing

real-time BVH, ray tracing cores

Ray Tracing

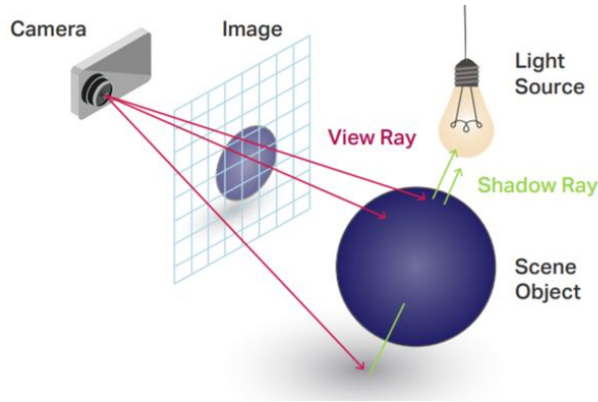


Ray Tracing

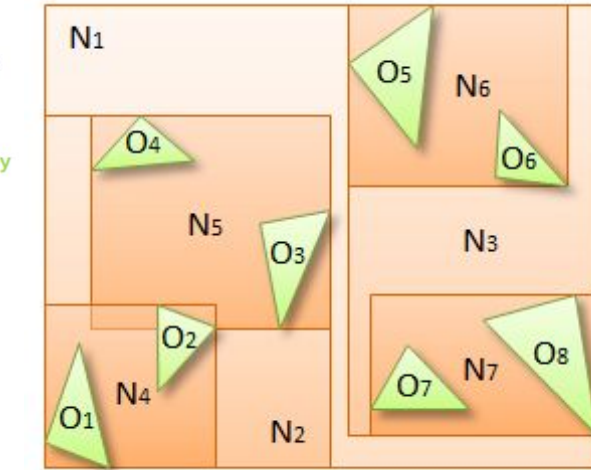


Как найти пересечение с треугольником?
Перебором для каждого луча всех
треугольников?

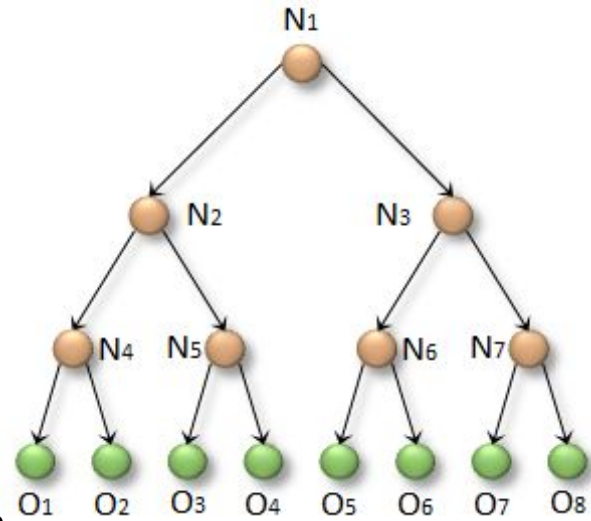
Ray Tracing



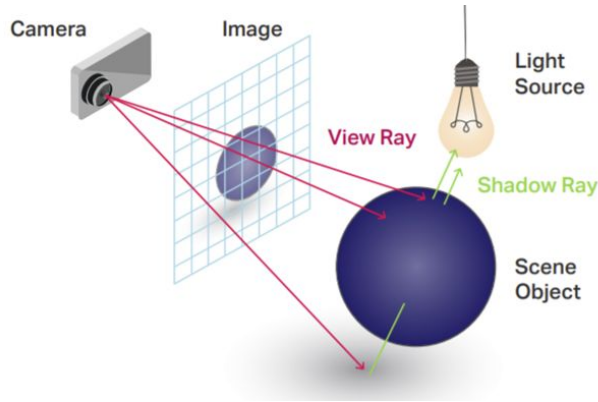
BVH - Bounding Volume Hierarchy



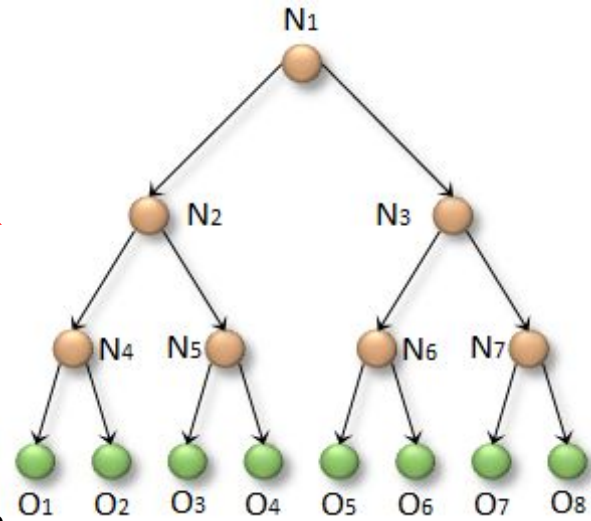
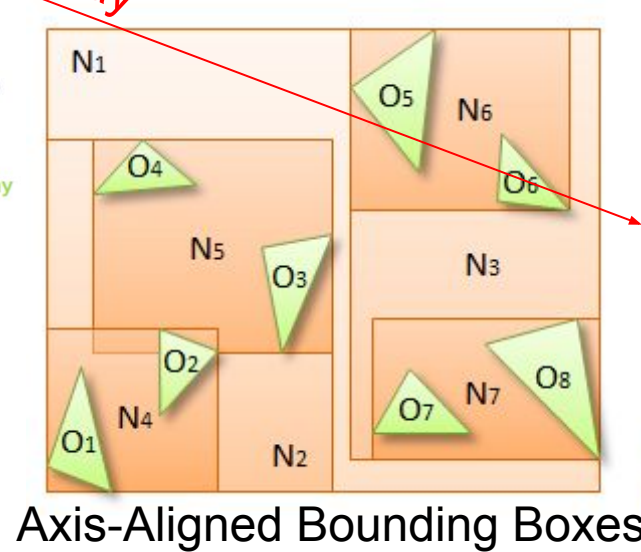
Axis-Aligned Bounding Boxes



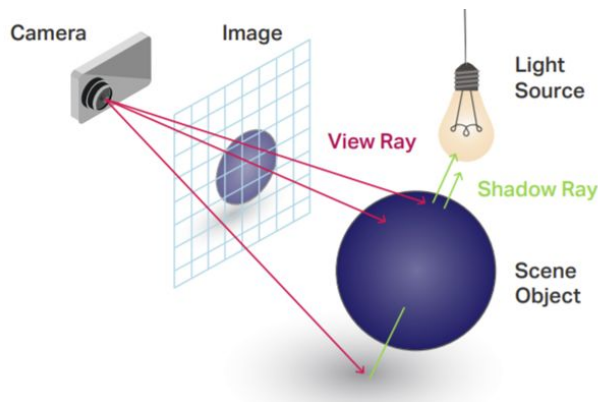
Ray Tracing



BVH - Bounding Volume Hierarchy

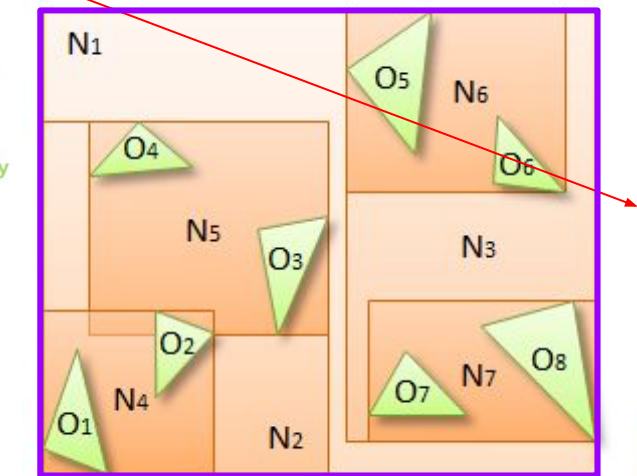


Ray Tracing

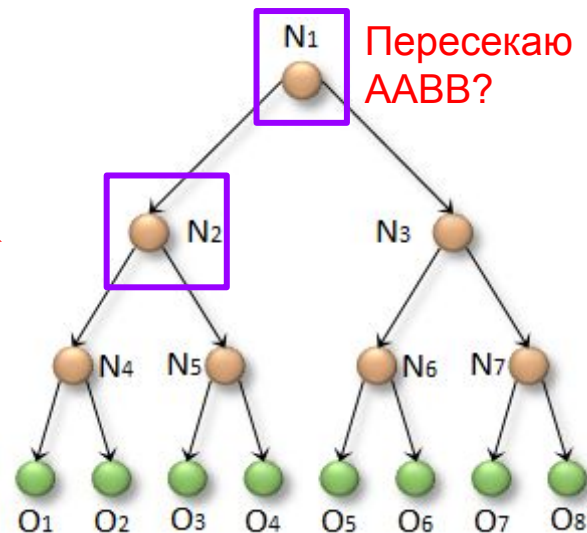


View Ray

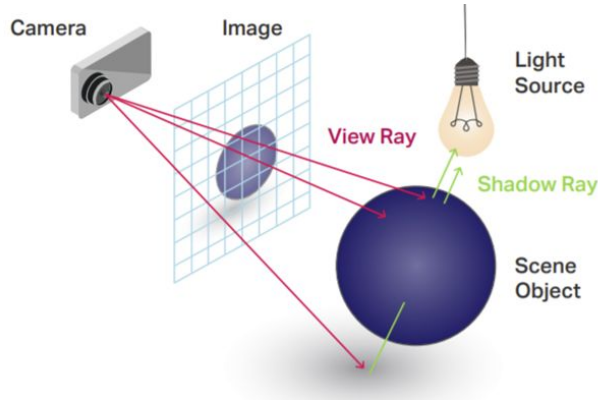
BVH - Bounding Volume Hierarchy



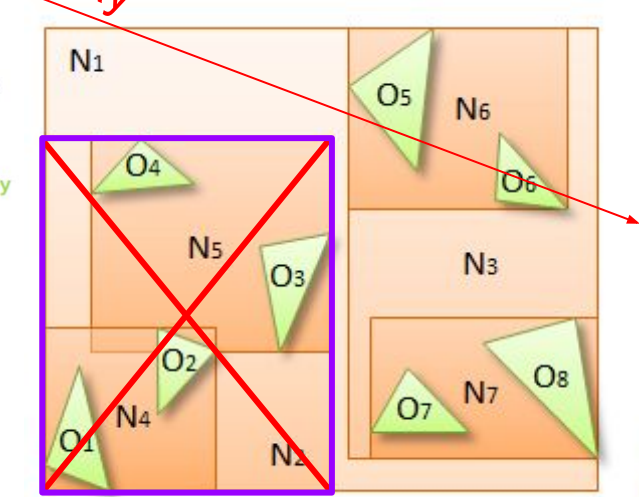
Axis-Aligned Bounding Boxes



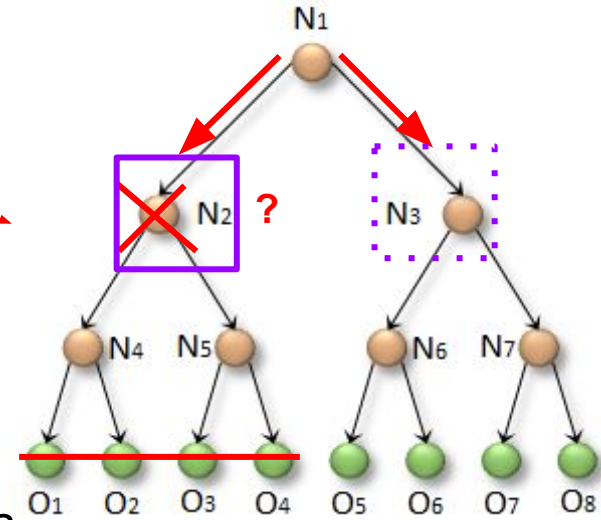
Ray Tracing



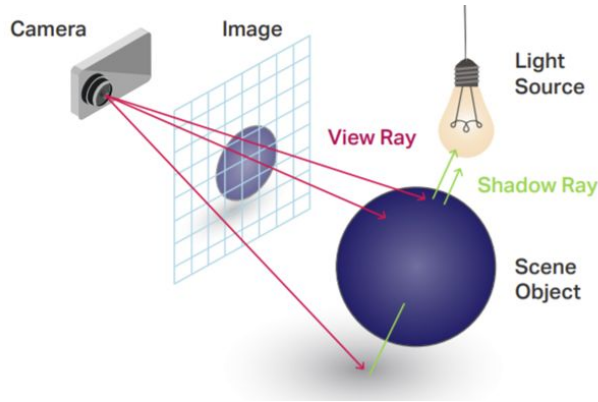
BVH - Bounding Volume Hierarchy



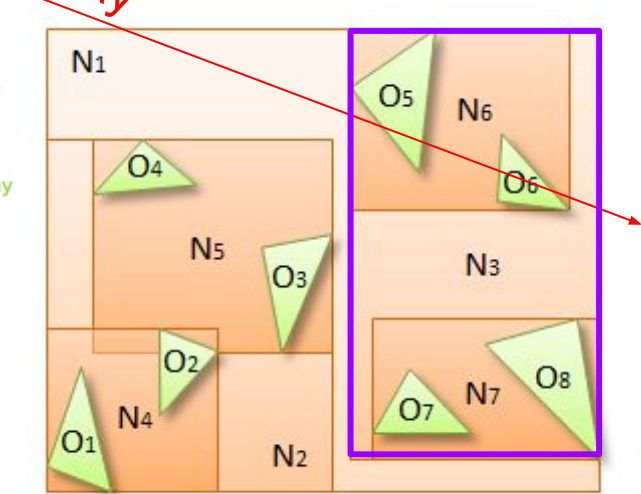
Axis-Aligned Bounding Boxes



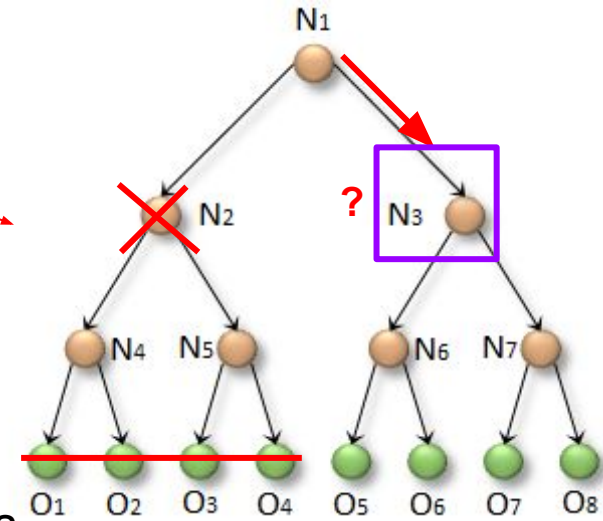
Ray Tracing



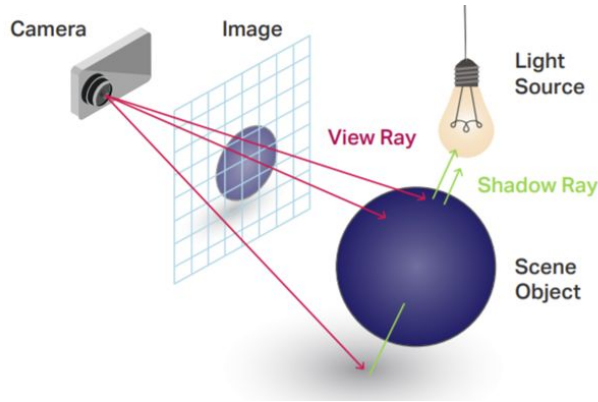
BVH - Bounding Volume Hierarchy



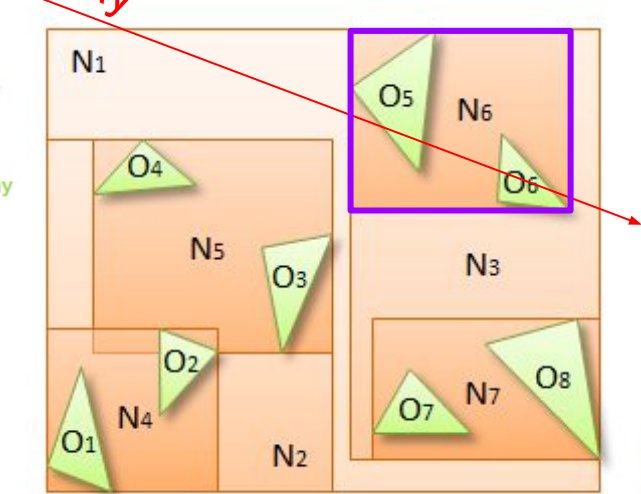
Axis-Aligned Bounding Boxes



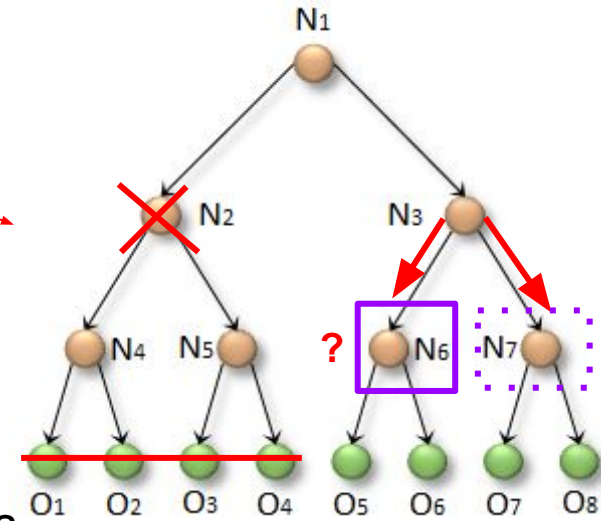
Ray Tracing



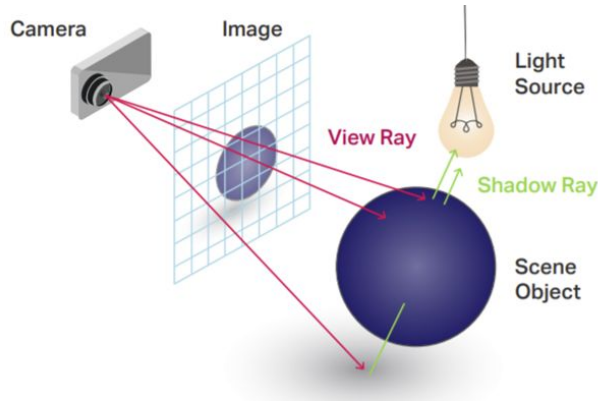
BVH - Bounding Volume Hierarchy



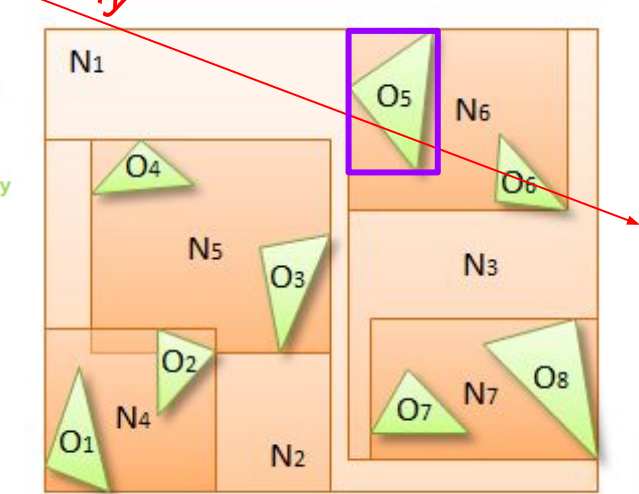
Axis-Aligned Bounding Boxes



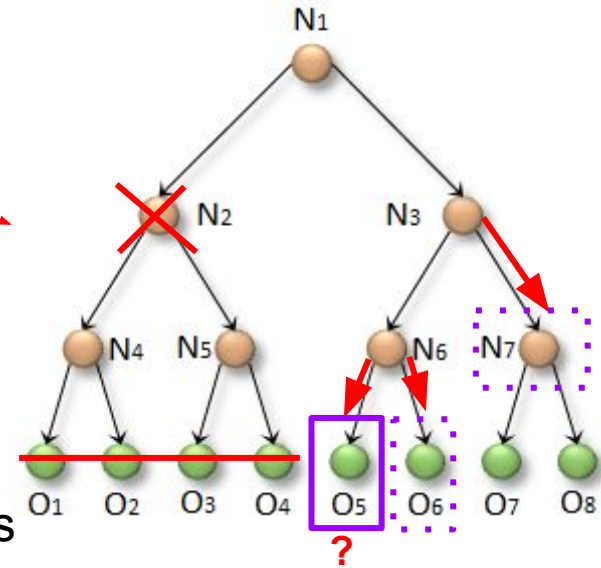
Ray Tracing



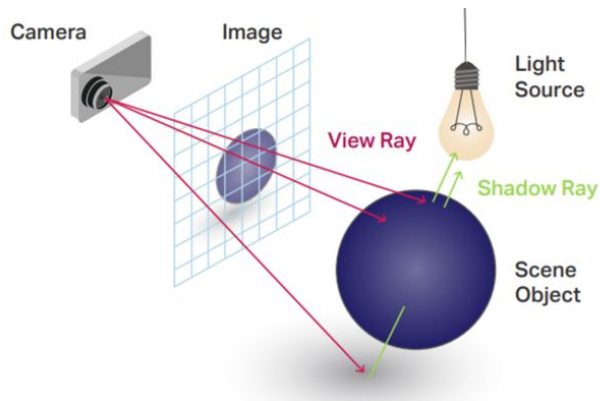
BVH - Bounding Volume Hierarchy



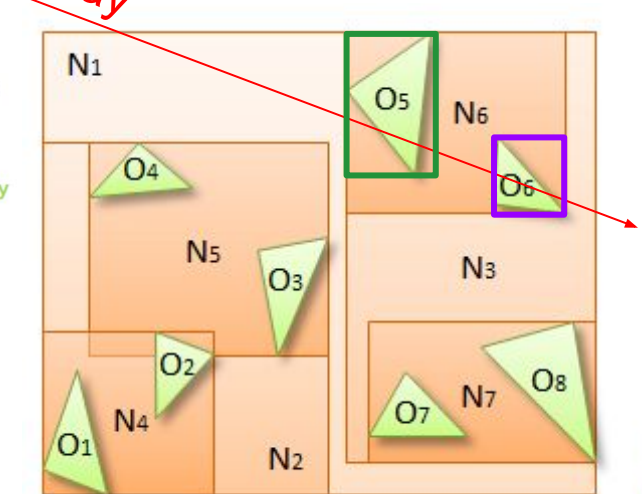
Axis-Aligned Bounding Boxes



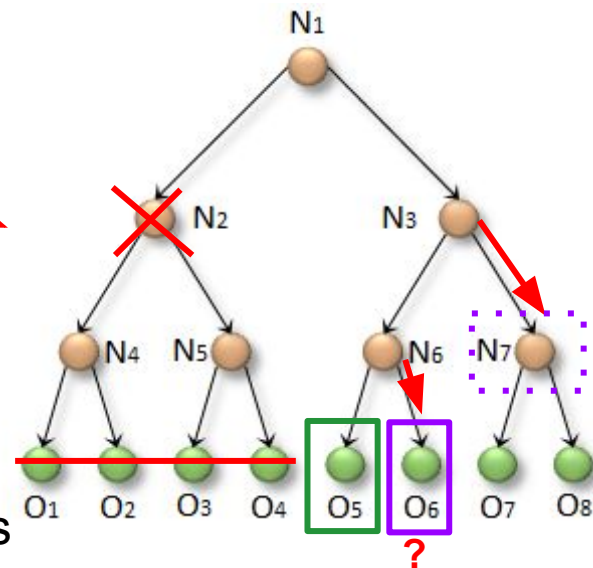
Ray Tracing



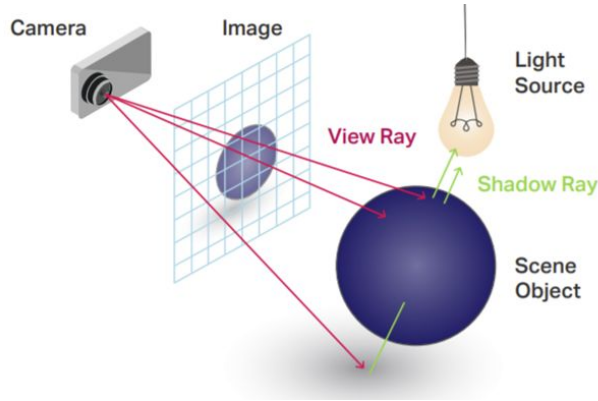
BVH - Bounding Volume Hierarchy



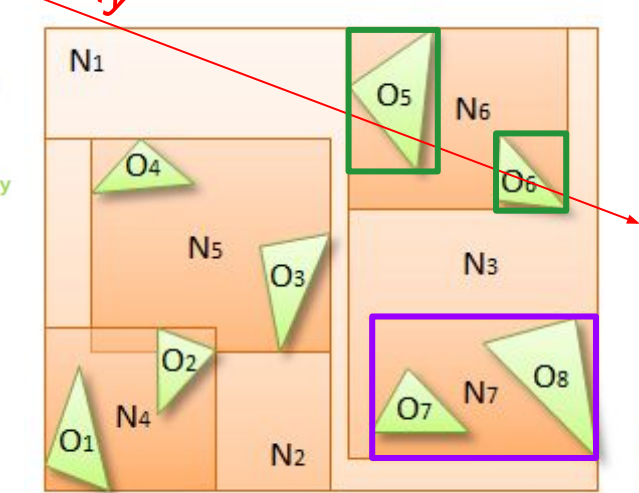
Axis-Aligned Bounding Boxes



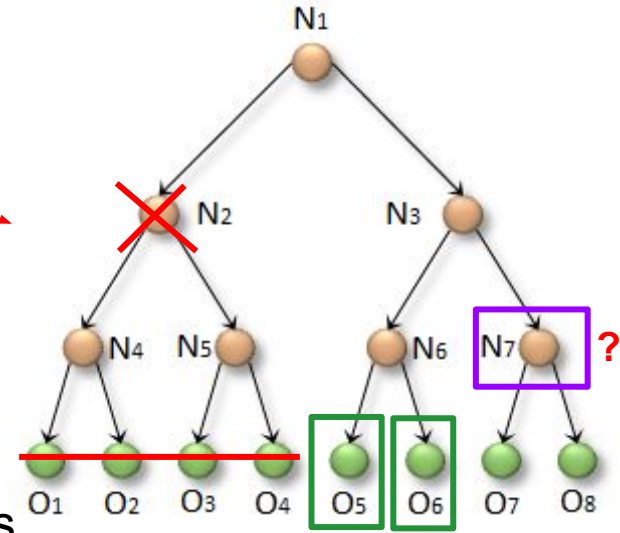
Ray Tracing



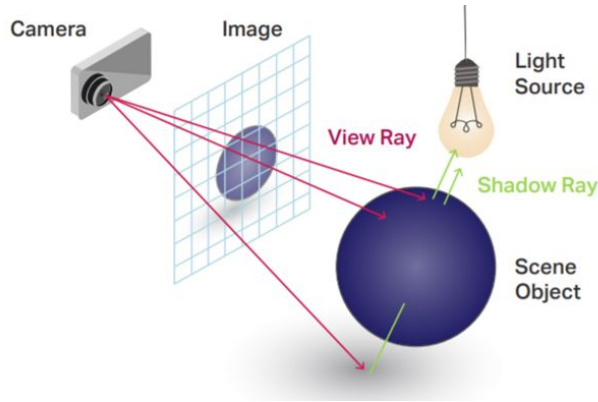
BVH - Bounding Volume Hierarchy



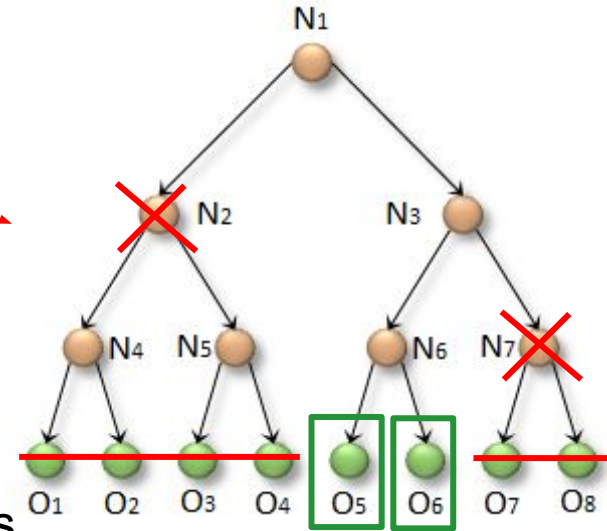
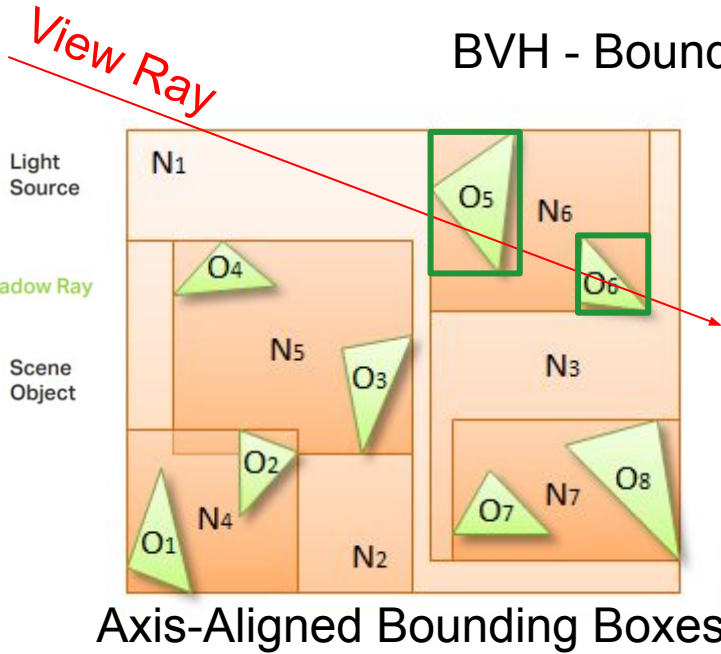
Axis-Aligned Bounding Boxes



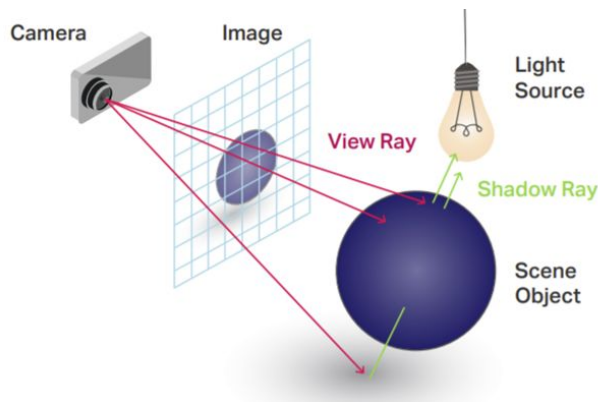
Ray Tracing



BVH - Bounding Volume Hierarchy

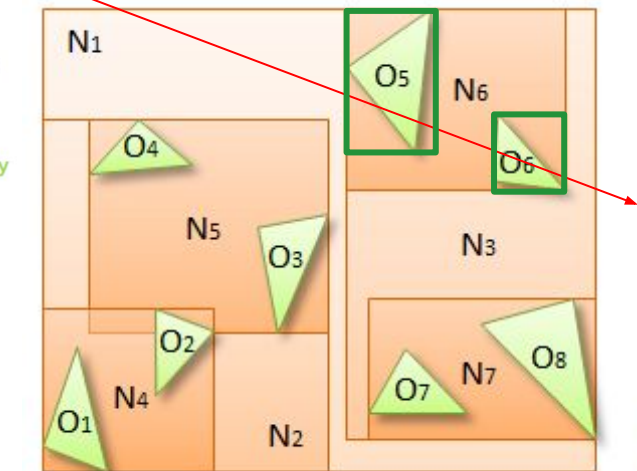


Ray Tracing

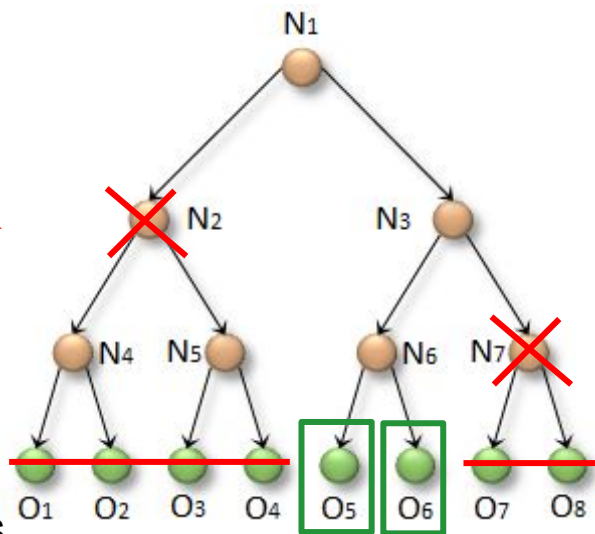


View Ray

BVH - Bounding Volume Hierarchy



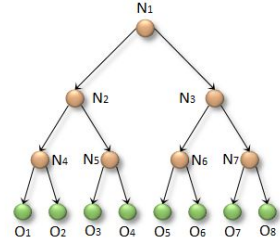
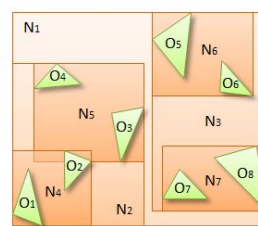
Axis-Aligned Bounding Boxes



Как это будет выглядеть в коде?

BVH обход (рекурсивный)

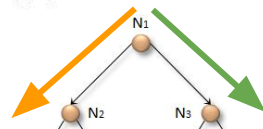
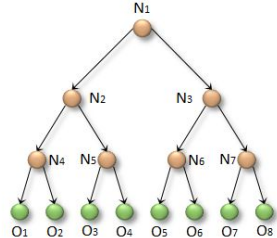
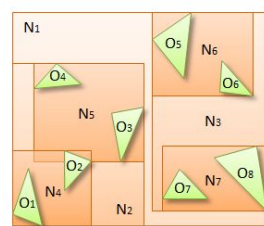
```
void traverseRecursive(AABB queryAABB,
                     int queryObjectIdx,
                     const __global BVHNode* node,
                     __local CollisionList* list)
{
    if (checkOverlap(node->getAABB(), queryAABB)) {
        if (node->isLeaf()) {
            list.add(queryObjectIdx, node->getObjectIdx());
        } else {
            const __global BVHNode childL = node->getLeftChild();
            const __global BVHNode childR = node->getRightChild();
            traverseRecursive(queryAABB, queryObjectIdx,
                             childL, list);
            traverseRecursive(queryAABB, queryObjectIdx,
                             childR, list);
        }
    }
}
```



BVH обход (рекурсивный)

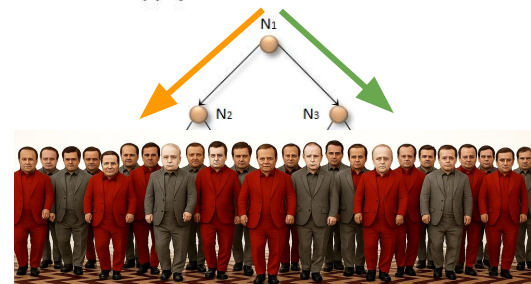
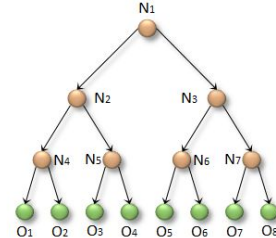
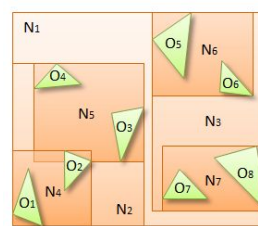
```

void traverseRecursive(AABB
                      int
                      const __global BVHNode* node,
                      __local CollisionList* list)
{
    if (checkOverlap(node->getAABB(), queryAABB)) {
        if (node->isLeaf()) {
            list.add(queryObjectIdx, node->getObjectIdx());
        } else {
            const __global BVHNode childL = node->getLeftChild();
            const __global BVHNode childR = node->getRightChild();
            traverseRecursive(queryAABB, queryObjectIdx,
                             childL, list);
            traverseRecursive(queryAABB, queryObjectIdx,
                             childR, list);
        }
    }
}
    
```



BVH обход (рекурсивный)

```
void traverseRecursive(AABB queryAABB,
                      int queryObjectIdx,
                      const __global BVHNode* node,
                      __local CollisionList* list)
{
    if (checkOverlap(node->getAABB(), queryAABB)) {
        if (node->isLeaf()) {
            list.add(queryObjectIdx, node->getObjectIdx());
        } else {
            const __global BVHNode childL = node->getLeftChild();
            const __global BVHNode childR = node->getRightChild();
            traverseRecursive(queryAABB, queryObjectIdx,
                             childL, list);
            traverseRecursive(queryAABB, queryObjectIdx,
                             childR, list);
        }
    }
}
```



Code divergence!

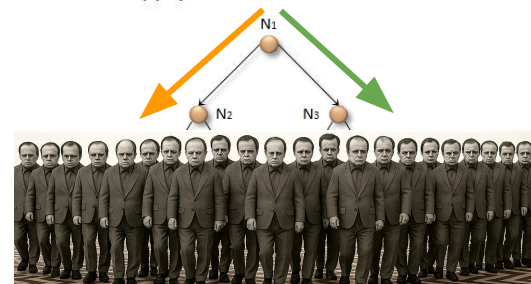
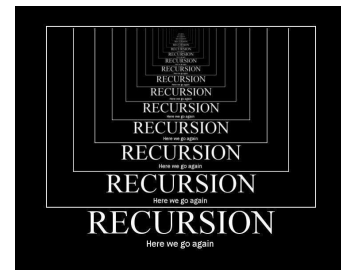
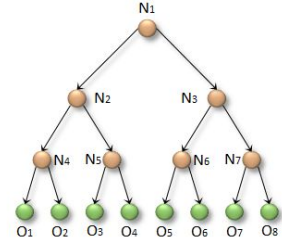
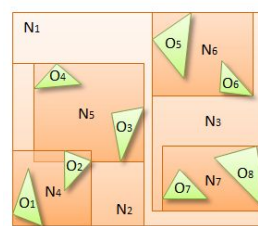
BVH обход (рекурсивный)

```

void traverseRecursive(AABB
    int
    const __global BVHNode* node,
    __local CollisionList* list)

{
    if (checkOverlap(node->getAABB(), queryAABB)) {
        if (node->isLeaf()) {
            list.add(queryObjectId, node->getObjectIdx());
        } else {
            const __global BVHNode childL = node->getLeftChild();
            const __global BVHNode childR = node->getRightChild();
            traverseRecursive(queryAABB, queryObjectId,
                             childL, list);
            traverseRecursive(queryAABB, queryObjectId,
                             childR, list);
        }
    }
}
    
```

Что делать?



Code divergence!

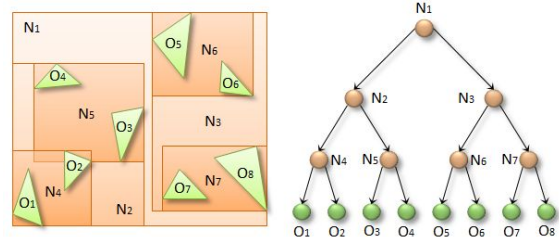
BVH обход (на стеке)

```
BVHNode* stack[MAX_STACK_SIZE];
BVHNode* node = bvhRoot;

do {
    BVHNode* childL = node->getLeftChild();
    BVHNode* childR = node->getRightChild();
    bool overlapL = checkOverlap(queryAABB, childL->getAABB());
    bool overlapR = checkOverlap(queryAABB, childR->getAABB());

    if (overlapL && childL->isLeaf())
        list.add(queryObjectIdx, childL->getObjectIdx());
    if (overlapR && childR->isLeaf())
        list.add(queryObjectIdx, childR->getObjectIdx());

    if (!traverseL && !traverseR) {
        node = stack.pop();
    } else {
        node = traverseL ? childL : childR;
        if (traverseL && traverseR)
            stack.push(childR);
    }
} while (node != NULL);
```



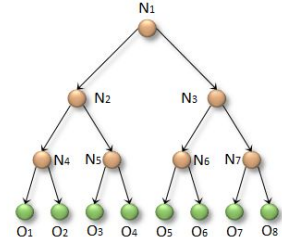
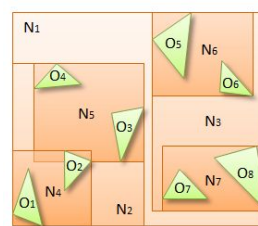
BVH обход (на стеке)

```
BVHNode* stack[MAX_STACK_SIZE];
BVHNode* node = bvhRoot;

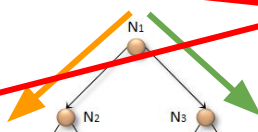
do {
    BVHNode* childL = node->getLeftChild();
    BVHNode* childR = node->getRightChild();
    bool overlapL = checkOverlap(queryAABB, childL->getAABB());
    bool overlapR = checkOverlap(queryAABB, childR->getAABB());

    if (overlapL && childL->isLeaf())
        list.add(queryObjectIdx, childL->getObjectIdx());
    if (overlapR && childR->isLeaf())
        list.add(queryObjectIdx, childR->getObjectIdx());

    if (!traverseL && !traverseR) {
        node = stack.pop();
    } else {
        node = traverseL ? childL : childR;
        if (traverseL && traverseR)
            stack.push(childR);
    }
} while (node != NULL);
```

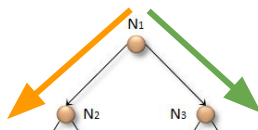
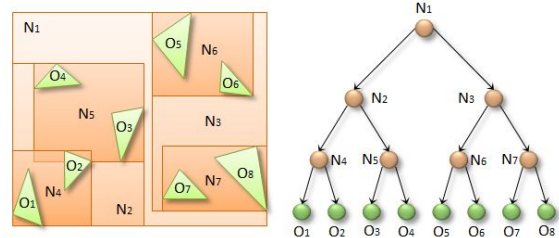


Code divergence остался только
в виде числа итераций цикла.



BVH обход (на стеке)

```
BVHNode* stack[MAX_STACK_SIZE];  
BVHNode* node = bvhRoot;  
do {  
    BVHNode* childL = node->getLeftChild();  
    BVHNode* childR = node->getRightChild();  
    bool overlapL = checkOverlap(queryAABB, childL->getAABB());  
    bool overlapR = checkOverlap(queryAABB, childR->getAABB());  
  
    if (overlapL && childL->isLeaf())  
        list.add(queryObjectIdx, childL->getObjectIdx());  
    if (overlapR && childR->isLeaf())  
        list.add(queryObjectIdx, childR->getObjectIdx());  
  
    if (!traverseL && !traverseR) {  
        node = stack.pop();  
    } else {  
        node = traverseL ? childL : childR;  
        if (traverseL && traverseR)  
            stack.push(childR);  
    }  
} while (node != NULL);
```



Code divergence остался только в виде числа итераций цикла.

Data divergence увеличивает количество запрашиваемых кеш-линий.

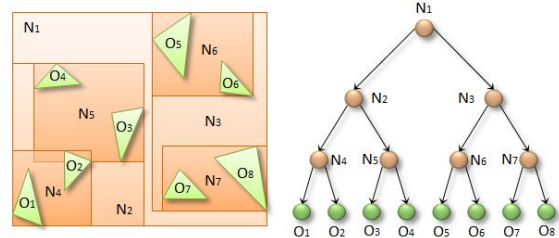
BVH обход (на стеке)

```
BVHNode* stack[MAX_STACK_SIZE];
BVHNode* node = bvhRoot;

do {
    BVHNode* childL = node->getLeftChild();
    BVHNode* childR = node->getRightChild();
    bool overlapL = checkOverlap(queryAABB, childL->getAABB());
    bool overlapR = checkOverlap(queryAABB, childR->getAABB());

    if (overlapL && childL->isLeaf())
        list.add(queryObjectIdx, childL->getObjectIdx());
    if (overlapR && childR->isLeaf())
        list.add(queryObjectIdx, childR->getObjectIdx());

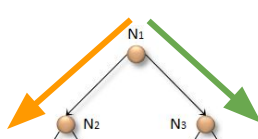
    if (!traverseL && !traverseR) {
        node = stack.pop();
    } else {
        node = traverseL ? childL : childR;
        if (traverseL && traverseR)
            stack.push(childR);
    }
} while (node != NULL);
```



Как увеличить когерентность по данным?

Code divergence остался только в виде числа итераций цикла.

Data divergence увеличивает количество запрашиваемых кеш-линий.



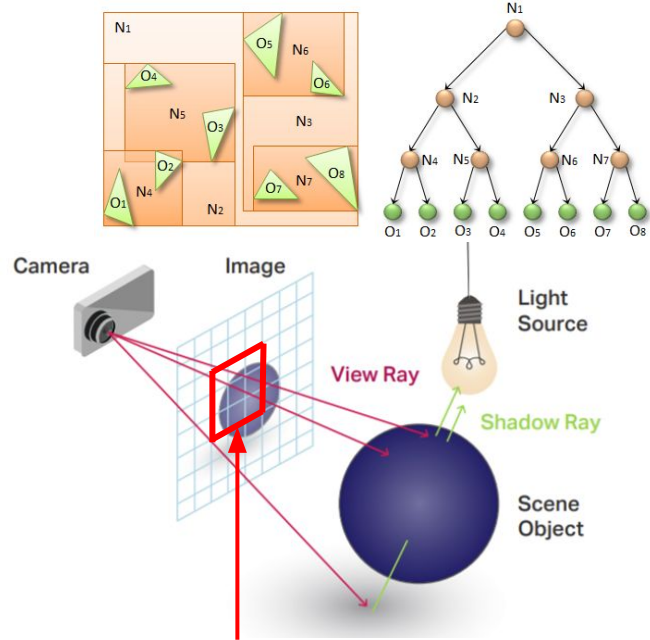
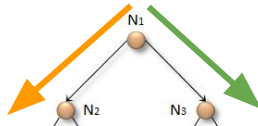
BVH обход (на стеке)

```
BVHNode* stack[MAX_STACK_SIZE];
BVHNode* node = bvhRoot;

do {
    BVHNode* childL = node->getLeftChild();
    BVHNode* childR = node->getRightChild();
    bool overlapL = checkOverlap(queryAABB, childL->getAABB());
    bool overlapR = checkOverlap(queryAABB, childR->getAABB());

    if (overlapL && childL->isLeaf())
        list.add(queryObjectIdx, childL->getObjectIdx());
    if (overlapR && childR->isLeaf())
        list.add(queryObjectIdx, childR->getObjectIdx());

    if (!traverseL && !traverseR) {
        node = stack.pop();
    } else {
        node = traverseL ? childL : childR;
        if (traverseL && traverseR)
            stack.push(childR);
    }
} while (node != NULL);
```

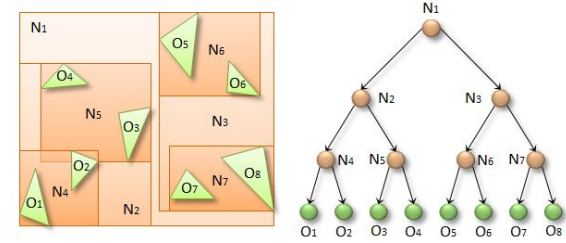


Данные когерентнее если в одном warp - пучок лучей!

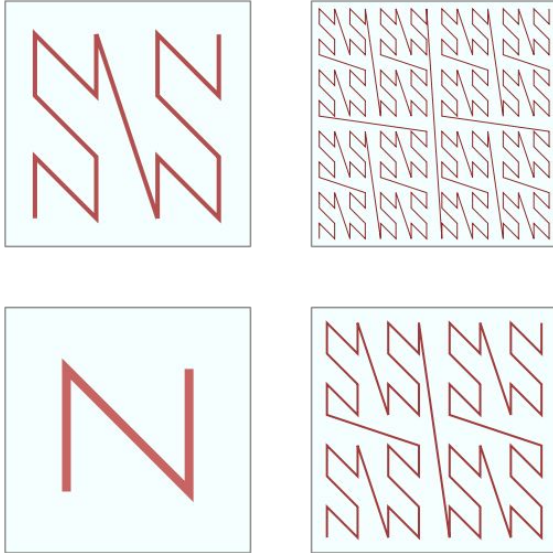
Code divergence остался только в виде числа итераций цикла.

Data divergence увеличивает количество запрашиваемых кеш-линий.

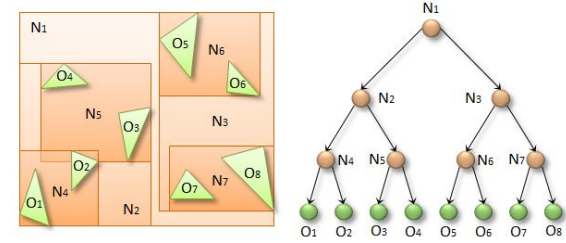
Real-time BVH construction



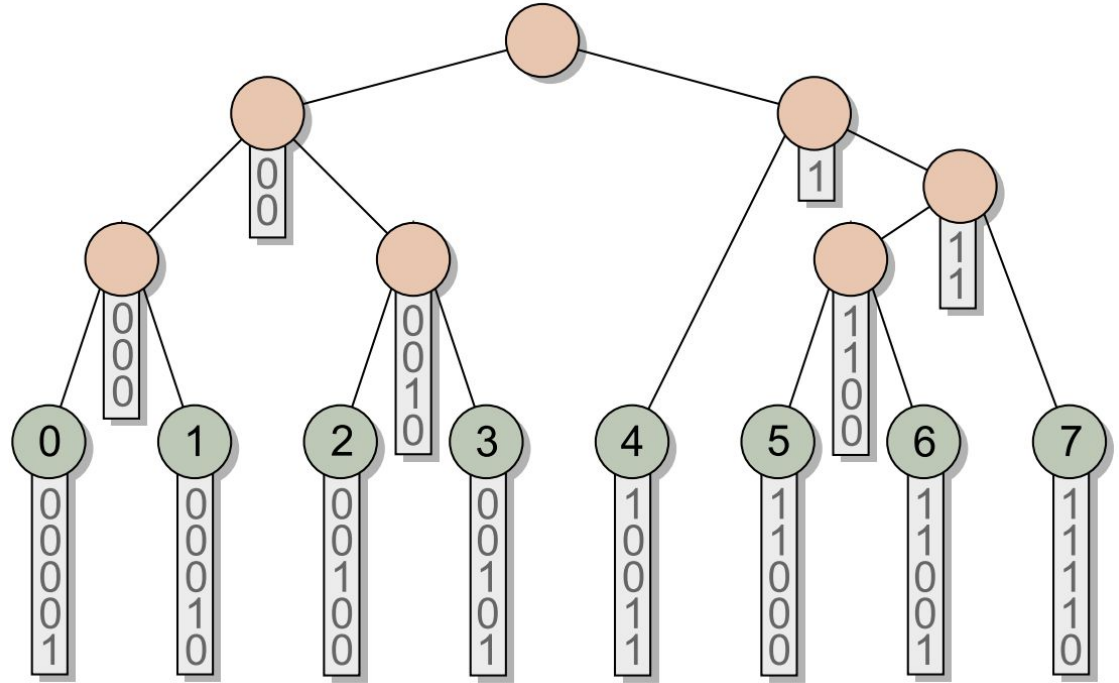
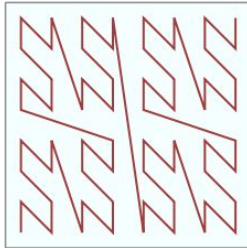
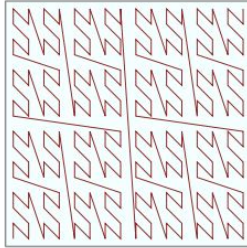
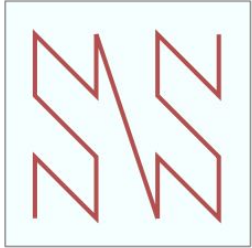
Z curve, Morton Code



Real-time BVH construction

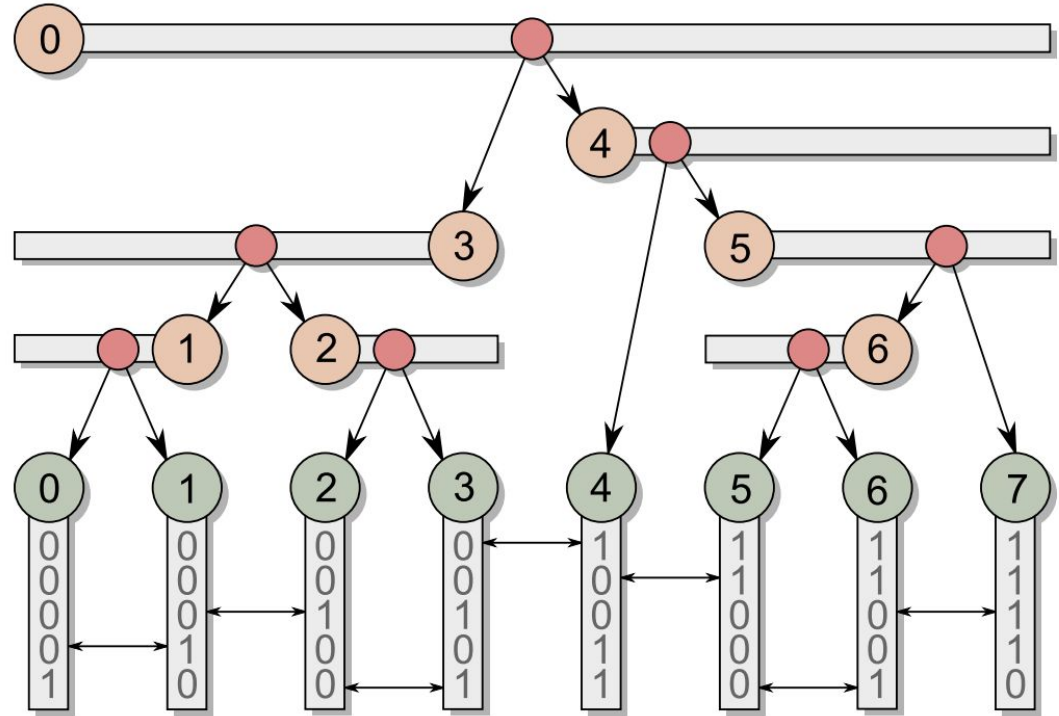
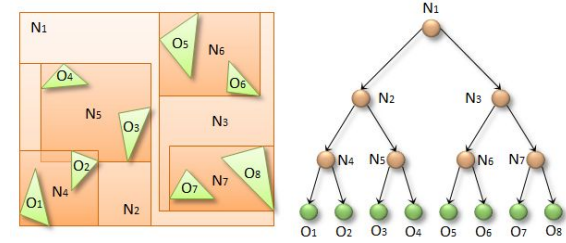
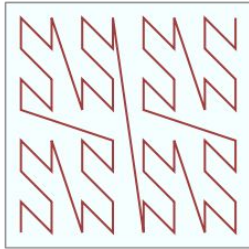
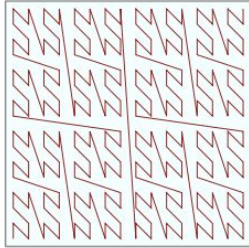
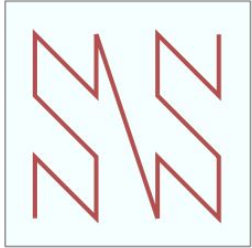


Z curve, Morton Code



Real-time BVH construction

Z curve, Morton Code



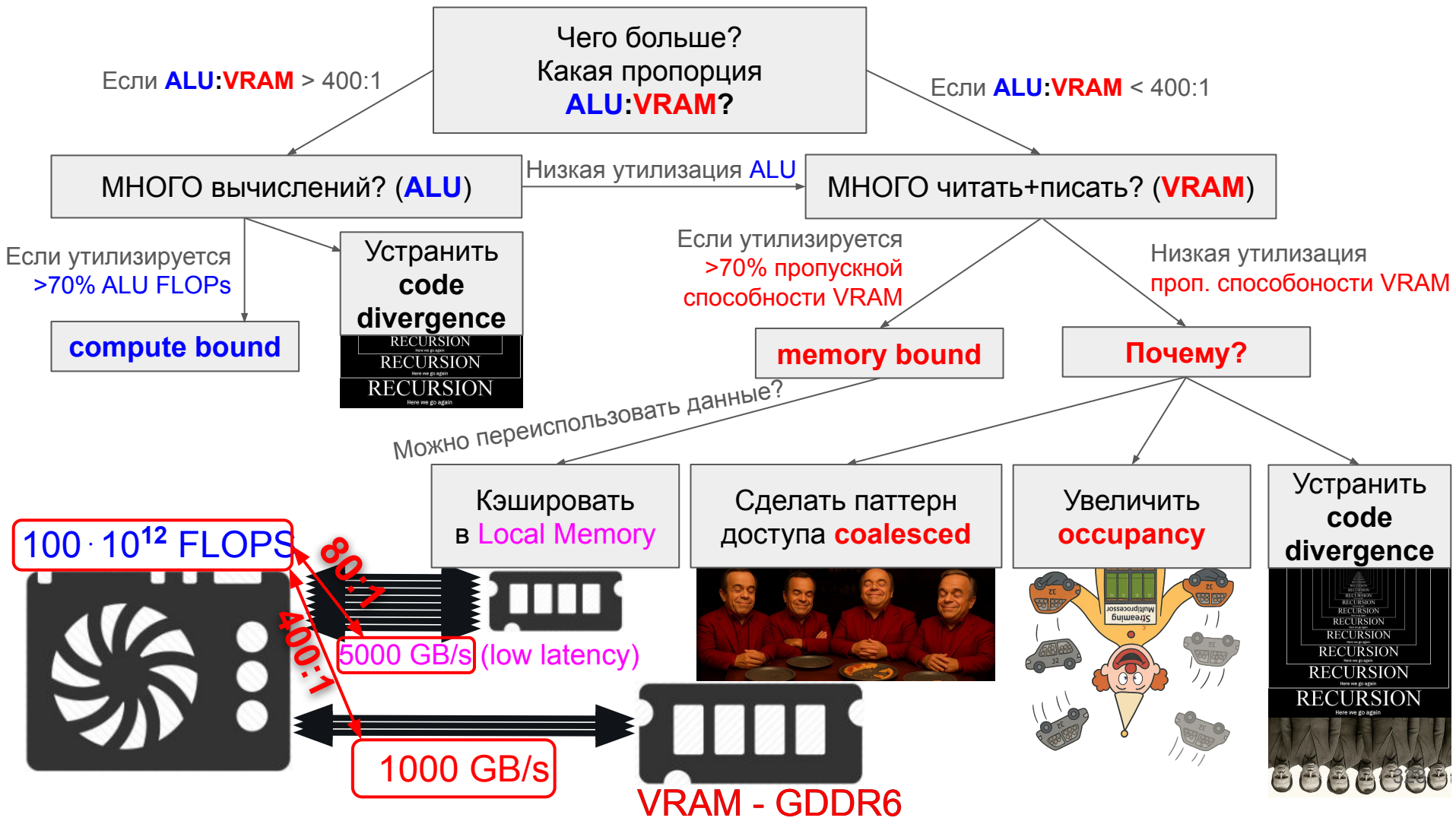
Ray Tracing Cores

- 1) Ускоряют обход BVH иерархии
- 2) Ускоряют пересечение луча с треугольником

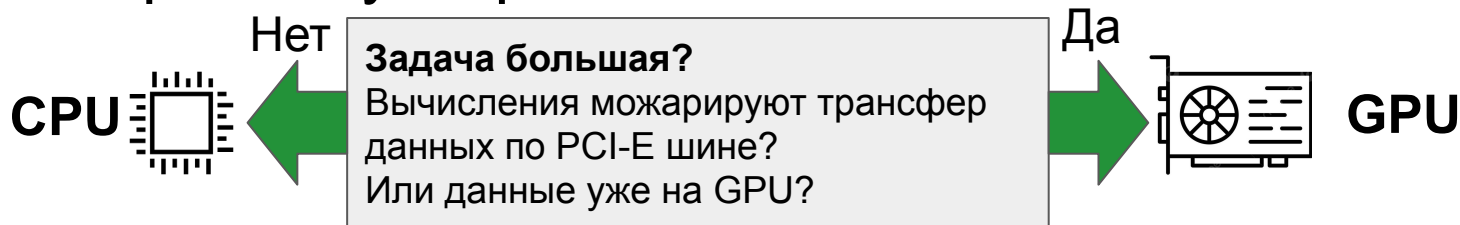


Глава 9: Выводы

Какие алгоритмы ускорятся на GPU?
OpenCL, CUDA или Vulkan?

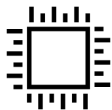


Какие алгоритмы ускоряются на GPU?



Какие алгоритмы ускоряются на GPU?

CPU

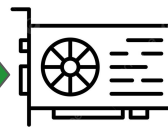


Нет

Задача большая?

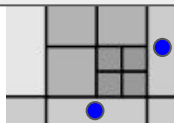
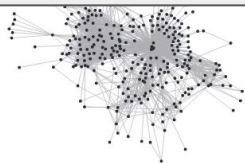
Вычисления можарируют трансфер данных по PCI-E шине?
Или данные уже на GPU?

Да



GPU

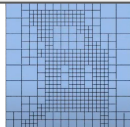
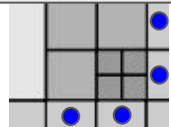
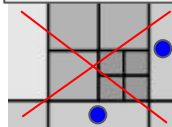
Произвольные графы,
несбалансированное октодереве,
min-cut/max flow (BK, IBFS),
disjoint sets (CHM)



Есть ли **массовый единообразный параллелизм**?

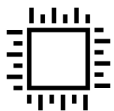
Структура данных - **регулярная**?
Не будет ли **uncoalesced** memory access pattern? Не будет ли **гонки**?
Не будет ли **code divergence**?

FEMs (Finite Element Methods) -
численные методы на
регулярной решетке и 2:1
сбалансированном октодереве



Какие алгоритмы ускоряются на GPU?

CPU

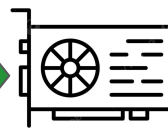


Нет

Задача большая?

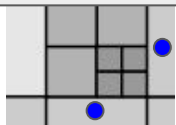
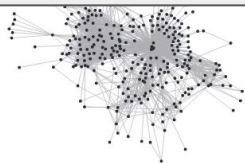
Вычисления можарируют трансфер данных по PCI-E шине?
Или данные уже на GPU?

Да



GPU

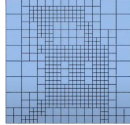
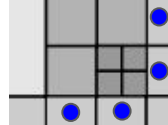
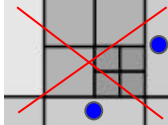
Произвольные графы,
несбалансированное октодереве,
min-cut/max flow (BK, IBFS),
disjoint sets (CHM)



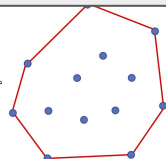
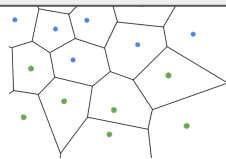
Есть ли массовый единообразный параллелизм?

Структура данных - **регулярная**?
Не будет ли **uncoalesced** memory access pattern? Не будет ли **гонки**?
Не будет ли **code divergence**?

FEMs (Finite Element Methods) -
численные методы на
регулярной решетке и 2:1
сбалансированном октодереве

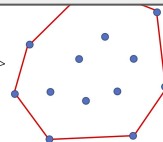
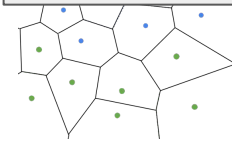


GIS, Триангуляция Делоне



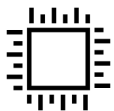
Не нужна ли **интервальная, рациональная, длинная арифметика** для абсолютной точности?

Иногда можно посчитать простые случаи на GPU, затем **на CPU досчитать** пограничные (**Триангуляция Делоне**)



Какие алгоритмы ускоряются на GPU?

CPU

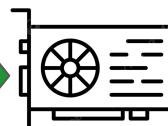


Нет

Задача большая?

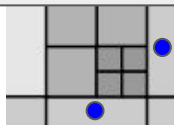
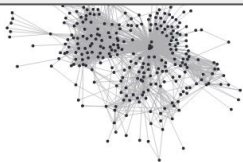
Вычисления можарируют трансфер данных по PCI-E шине?
Или данные уже на GPU?

Да



GPU

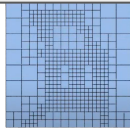
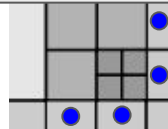
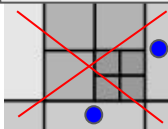
Произвольные графы,
несбалансированное октодереве,
min-cut/max flow (BK, IBFS),
disjoint sets (CHM)



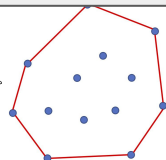
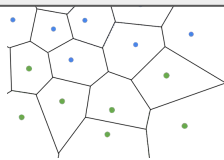
Есть ли **массовый единообразный параллелизм**?

Структура данных - **регулярная**?
Не будет ли **uncoalesced** memory access pattern? Не будет ли **гонки**?
Не будет ли **code divergence**?

FEMs (Finite Element Methods) -
численные методы на
регулярной решетке и 2:1
сбалансированном октодереве

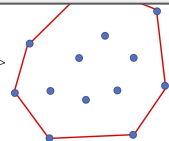
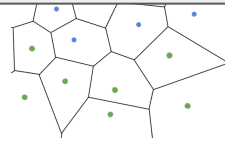


GIS, Триангуляция Делоне



Не нужна ли **интервальная, рациональная, длинная арифметика** для абсолютной точности?

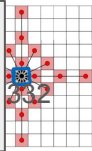
Иногда можно посчитать простые случаи на GPU, затем **на CPU досчитать** пограничные (**Триангуляция Делоне**)



Приоритетные очереди (жадины),
min-cut/max flow (BK, IBFS),
disjoint sets (CHM)

Нет ли **ГОРАЗДО** более эффективного асимптотически алгоритма, но при этом очень **линейного**?

Иногда можно выкрутиться:
- merge-sort
- BVH construction
- Gipuma



CUDA, Vulkan или OpenCL?



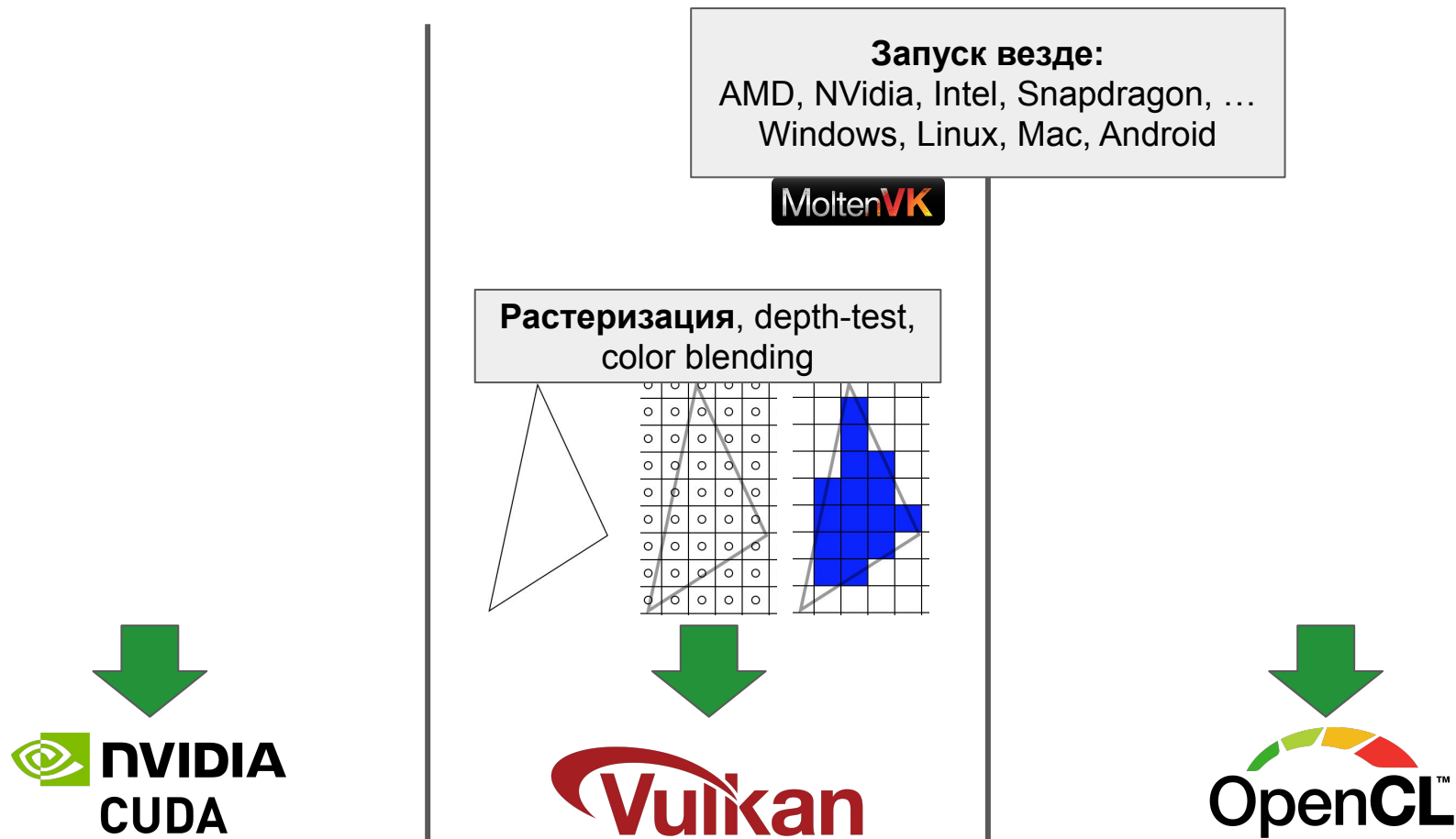
CUDA, Vulkan или OpenCL?

Запуск везде:
AMD, NVidia, Intel, Snapdragon, ...
Windows, Linux, Mac, Android

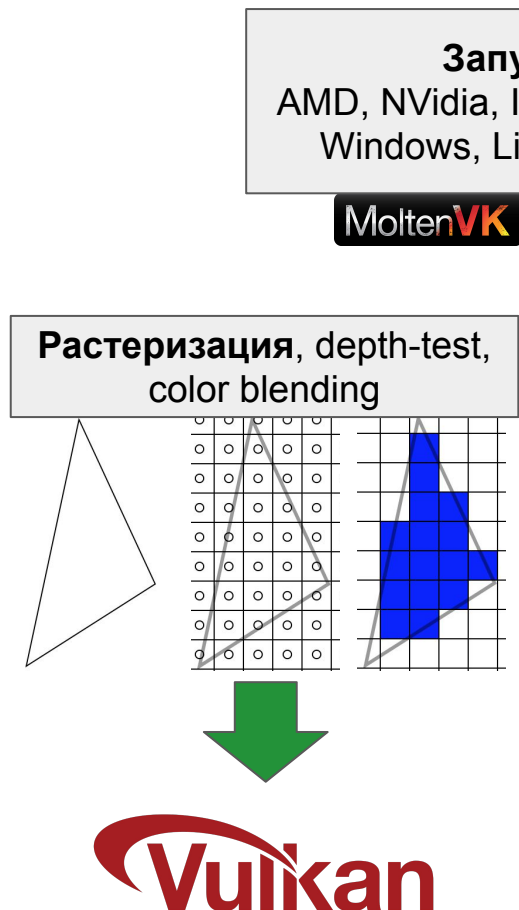
Molten**VK**



CUDA, Vulkan или OpenCL?



CUDA, Vulkan или OpenCL?



GPGPU: API и синтаксис
кernels проще чем в **Vulkan**



CUDA, Vulkan или OpenCL?

Тулинг: профилировщики,
санитайзеры
(cuda-memcheck/racecheck)



NVIDIA Nsight



**NVIDIA
CUDA**

Запуск везде:

AMD, NVidia, Intel, Snapdragon, ...
Windows, Linux, Mac, Android

MoltenVK

**Растеризация, depth-test,
color blending**



Vulkan

GPGPU: API и синтаксис
кernels проще чем в **Vulkan**



OpenCL™

CUDA, Vulkan или OpenCL?

Тулинг: профилировщики,
санитайзеры
(cuda-memcheck/racecheck)



NVIDIA Nsight

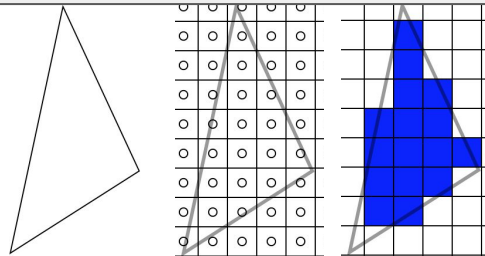


**NVIDIA
CUDA**

Тулинг:
Validation Layers,
SPIRV-Reflect,
RenderDoc,
профайлеры

MoltenVK

Растеризация, depth-test,
color blending



Vulkan

Запуск везде:

AMD, NVidia, Intel, Snapdragon, ...
Windows, Linux, Mac, Android

GPGPU: API и синтаксис
кernels проще чем в **Vulkan**



OpenCL™

CUDA, Vulkan или OpenCL?

AI/ML: очень хорошо
оптимизированные
библиотеки,
профилировщики



Тулинг: профилировщики,
санитайзеры
(cuda-memcheck/racecheck)

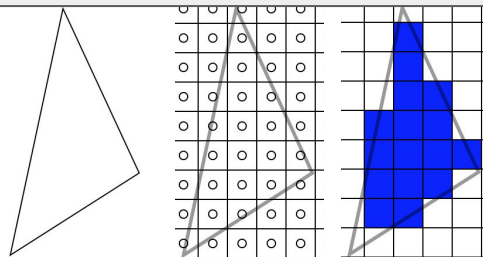


NVIDIA Nsight



Тулинг:
Validation Layers,
SPIRV-Reflect,
RenderDoc,
профайлеры

Растеризация, depth-test,
color blending



Запуск везде:

AMD, NVidia, Intel, Snapdragon, ...
Windows, Linux, Mac, Android



GPGPU: API и синтаксис
кernels проще чем в **Vulkan**



Low-Level NVIDIA HW:
Tensor Cores
Ray Tracing Cores

NVIDIA OPTIX

$D = \begin{matrix} \text{FP16 or FP32} & \text{FP16} & \text{FP16} & \text{FP16 or FP32} \end{matrix} + \begin{matrix} \text{FP16} & \text{FP16} & \text{FP16} & \text{FP16 or FP32} \end{matrix}$

AI/ML: очень хорошо
оптимизированные
библиотеки,
профилировщики



Тулинг: профилировщики,
санитайзеры
(cuda-memcheck/racecheck)



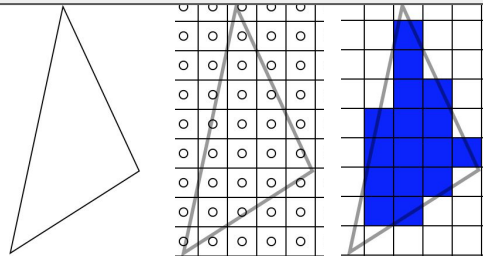
NVIDIA Nsight



**NVIDIA
CUDA**

Тулинг:
Validation Layers,
SPIRV-Reflect,
RenderDoc,
профайлеры

**Растеризация, depth-test,
color blending**



Vulkan

CUDA, Vulkan или OpenCL?

Запуск везде:

AMD, NVidia, Intel, Snapdragon, ...
Windows, Linux, Mac, Android

MoltenVK

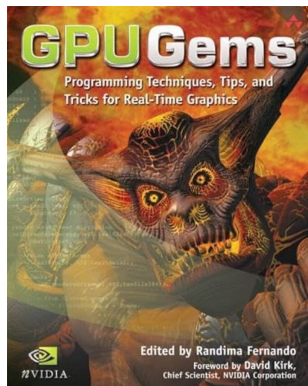
GPGPU: API и синтаксис
кernels проще чем в **Vulkan**



OpenCL™

Дополнительные материалы

- 1) **Source Code** - <https://github.com/GPGPUCourse/GPGPUVulkan>
- 2) **Книги/циклы статей - GPU Gems**
- 3) **Конференции - SIGGRAPH**

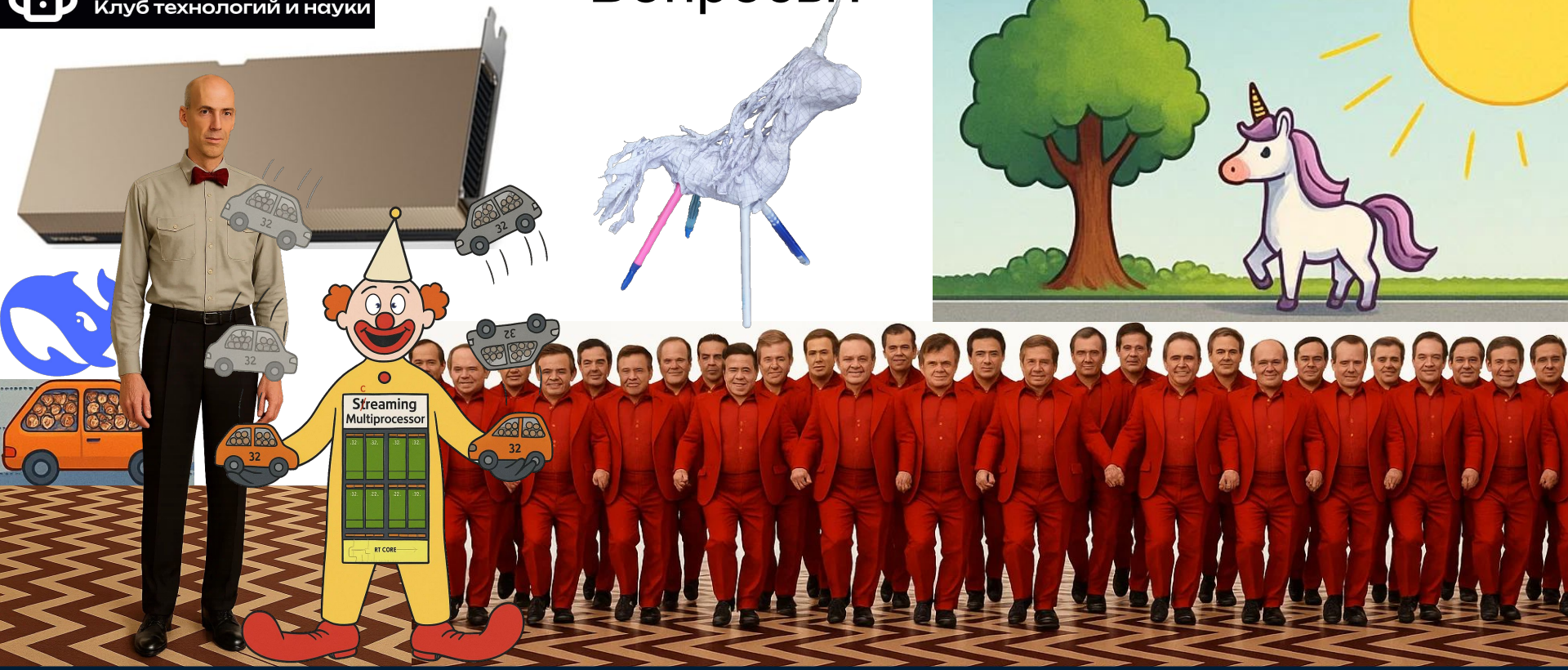



- 4) **Как работает Nanite в Unreal Engine 5** - <https://youtu.be/ltUzX1IR9JI>
- 5) **Real time BVH construction**
- 6) **Курс по видеокартам** - <https://youtu.be/LDt4KQEdImY>
- 7) **Курс фотограмметрии** - <https://youtu.be/rEF0zkv2cn8>

Вопросы?



Вопросы?



 [@PolarNick239](https://t.me/PolarNick239) ← он рад обсудить интересное!

 polarnick239@gmail.com **Agisoft**

 Николай Полярный  **Metashape**

Vulkan

 **OpenCL™**

 **NVIDIA**
CUDA

Вопросы?



Вопросы?



Вопросы?



Вопросы?





Бондосы?