



Вычисления на видеокартах

Лекция 4

- Транспонирование матрицы
- Умножение матриц
- Tensor Cores, WMMA
- Метод Штрассена
- Матрица матрицу
матрицово матрицит
матрицей в матрице

TENSOR CORES

$$A \times B = C$$



polarnick239@gmail.com

Николай Полярный

План лекции

- 1) **Транспонирование матрицы:** через локальную память (учет *bank conflicts*)
- 2) **Умножение матриц:** через локальную память
- 3) **Умножение матриц:** *Tensor Cores, WMMA* (и `VK_KHR_cooperative_matrix`)
- 4) **Умножение матриц:** оптимизации *DeepSeek*
- 5) **Умножение матриц:** метод *Штрассена*

Глава 1: Транспонирование

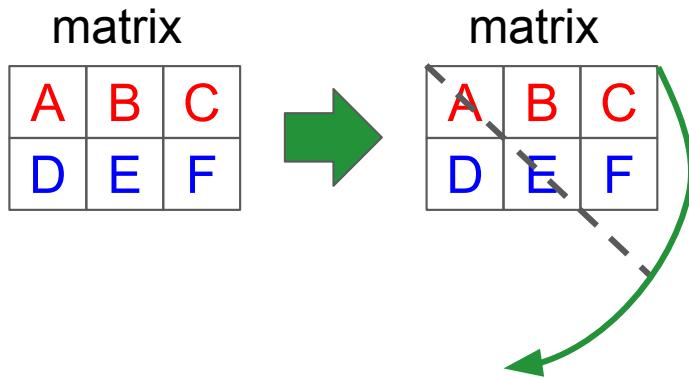
local memory, coalesced access, bank conflicts

Транспонирование матрицы

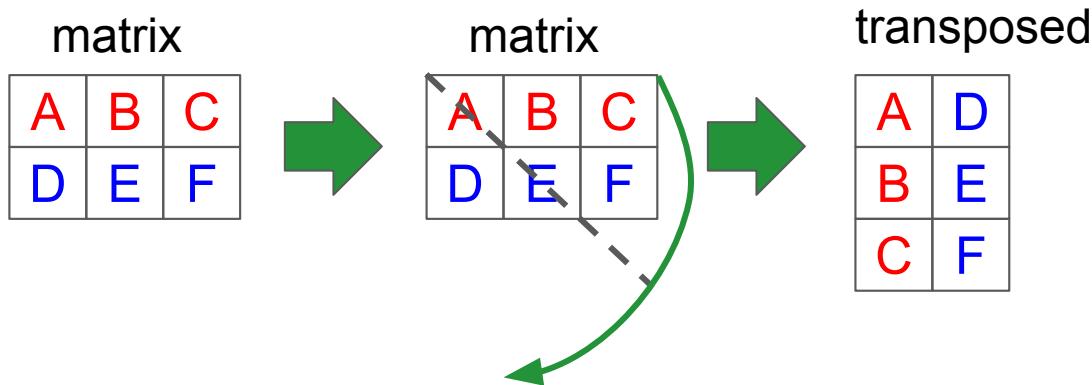
matrix

| | | |
|---|---|---|
| A | B | C |
| D | E | F |

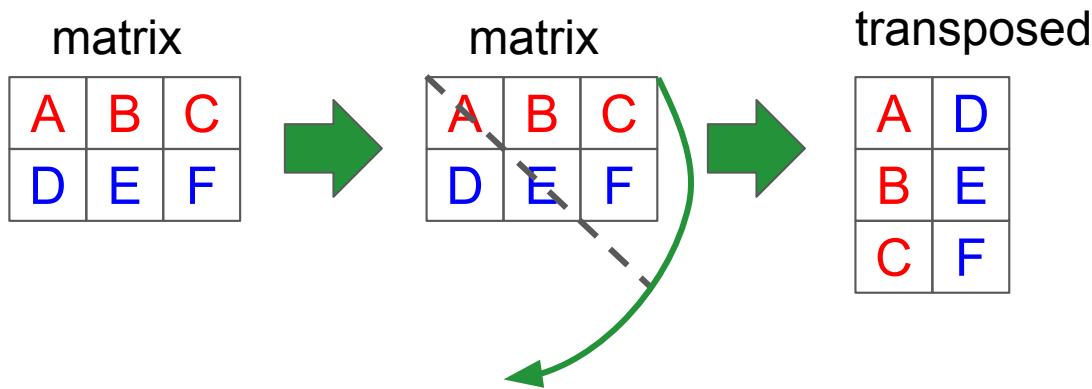
Транспонирование матрицы



Транспонирование матрицы

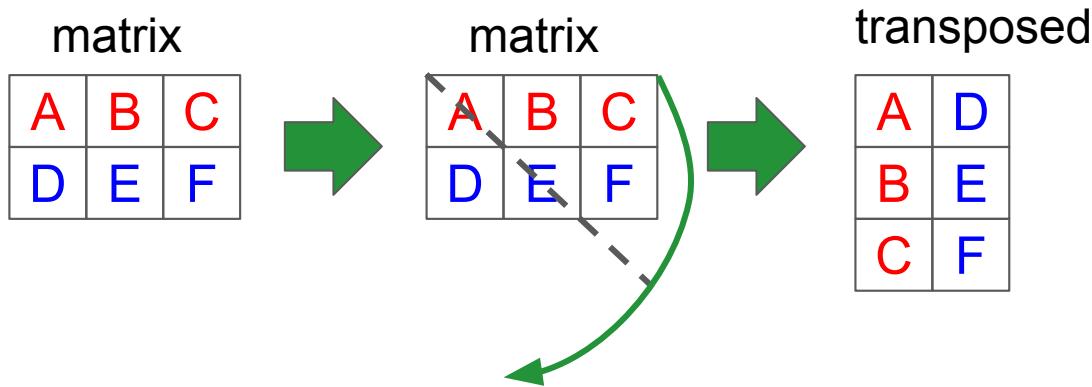


Транспонирование матрицы



Как выглядит:
- WorkRange?

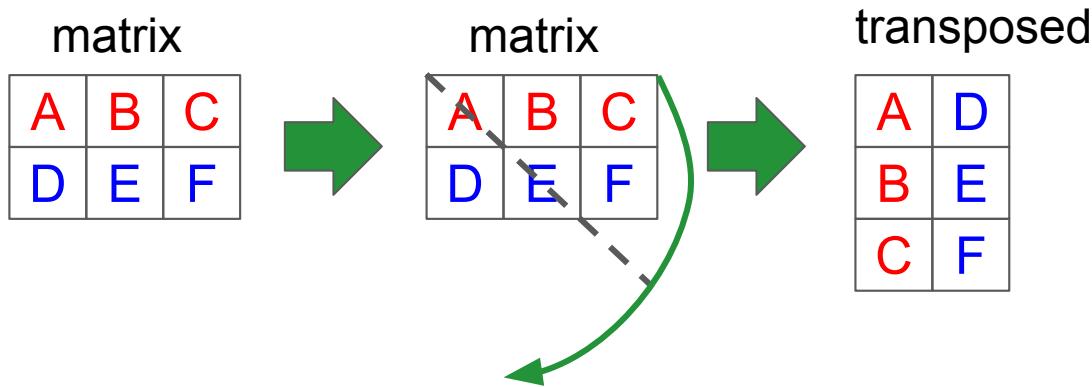
Транспонирование матрицы



Как выглядит:

- WorkRange?
- Задача WorkItem?

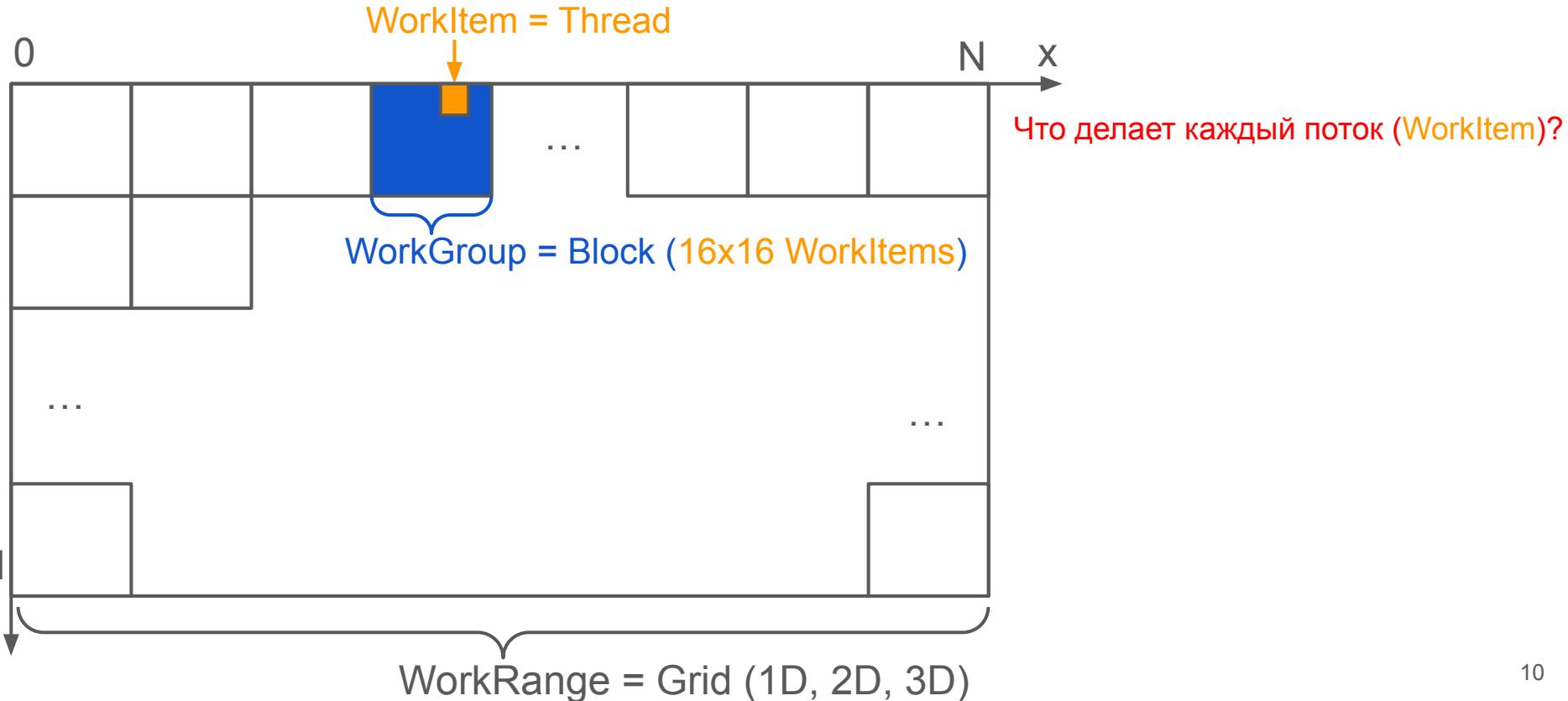
Транспонирование матрицы



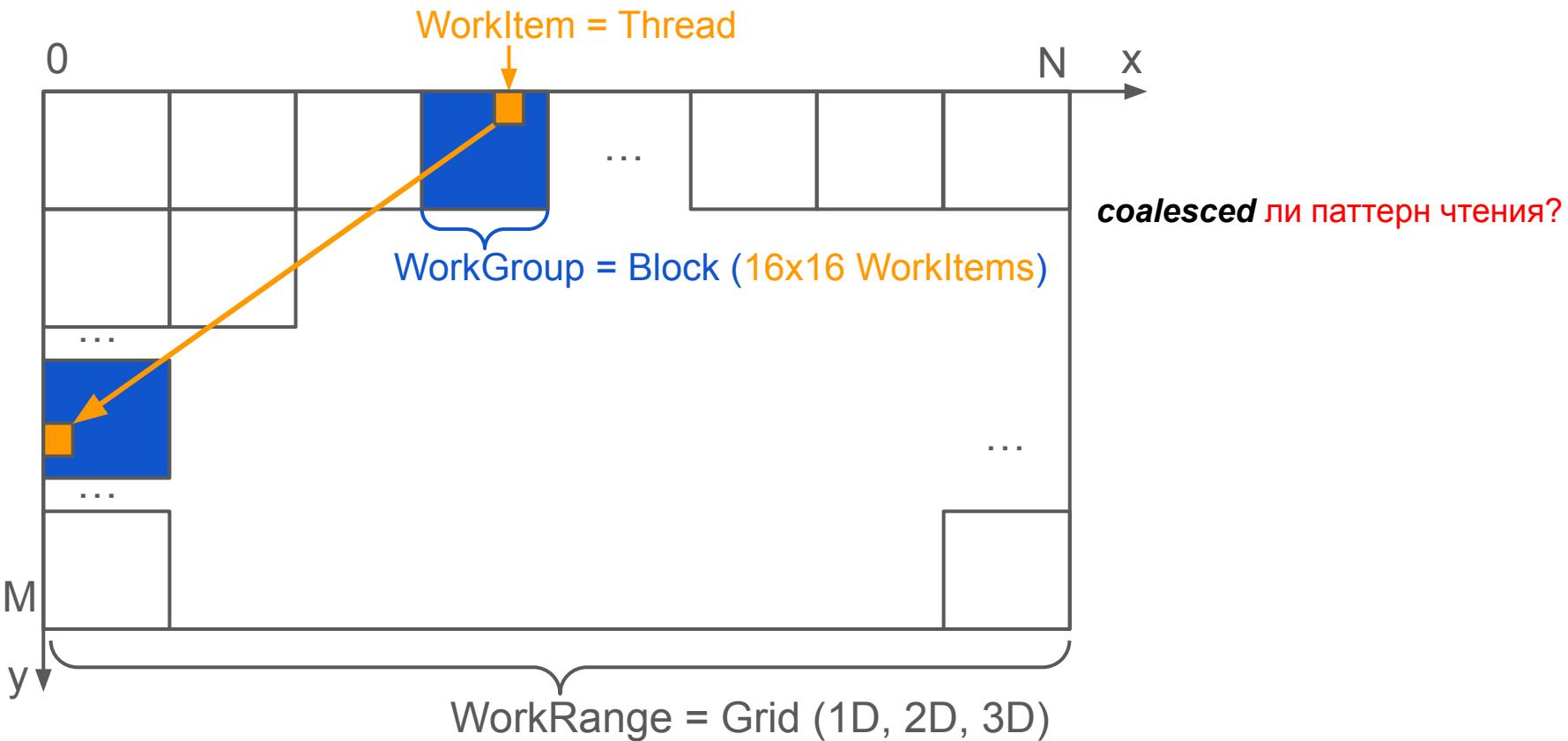
Как выглядит:

- WorkRange?
- Задача WorkItem?
- WorkGroup?

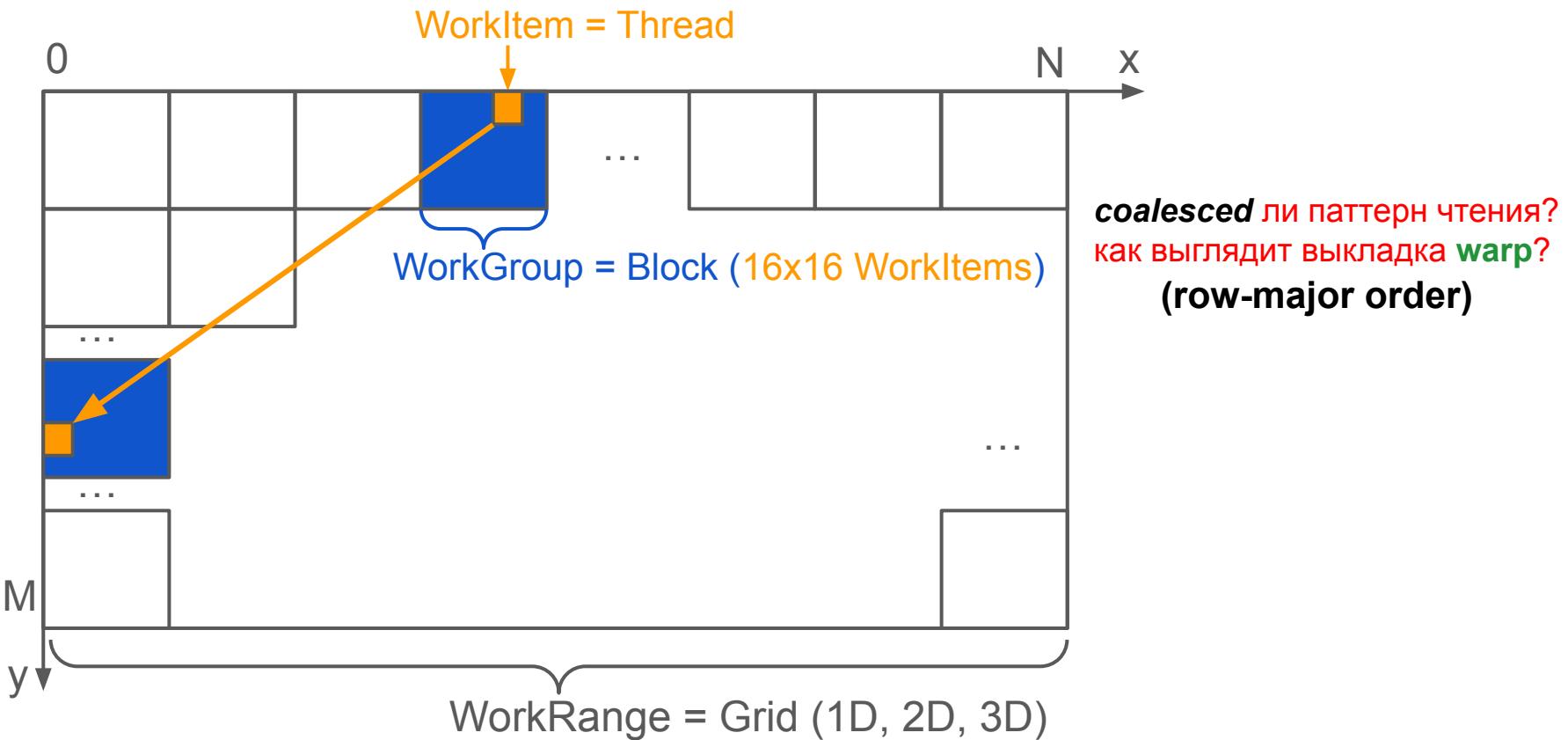
Транспонирование матрицы



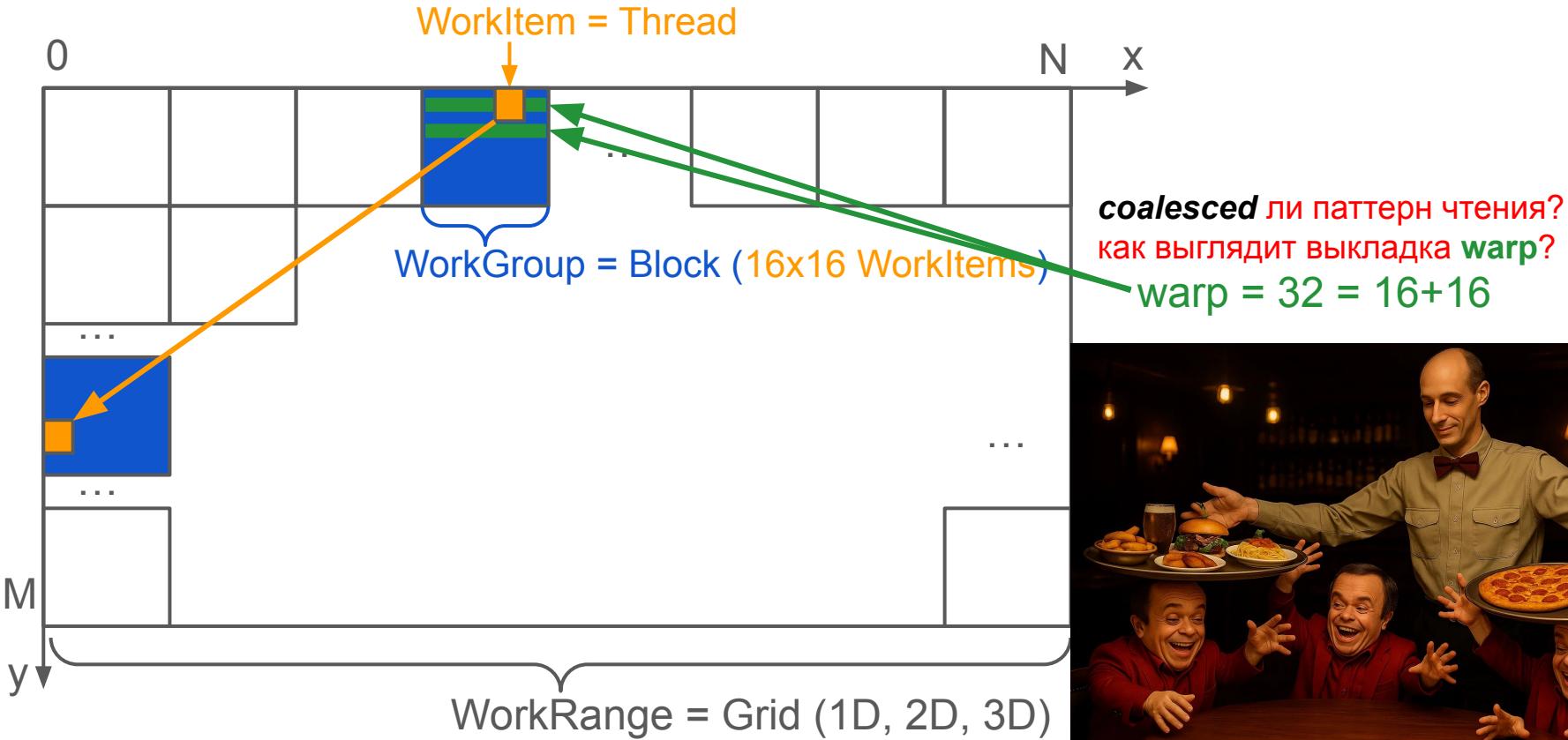
Транспонирование матрицы (*coalesced memory access*)



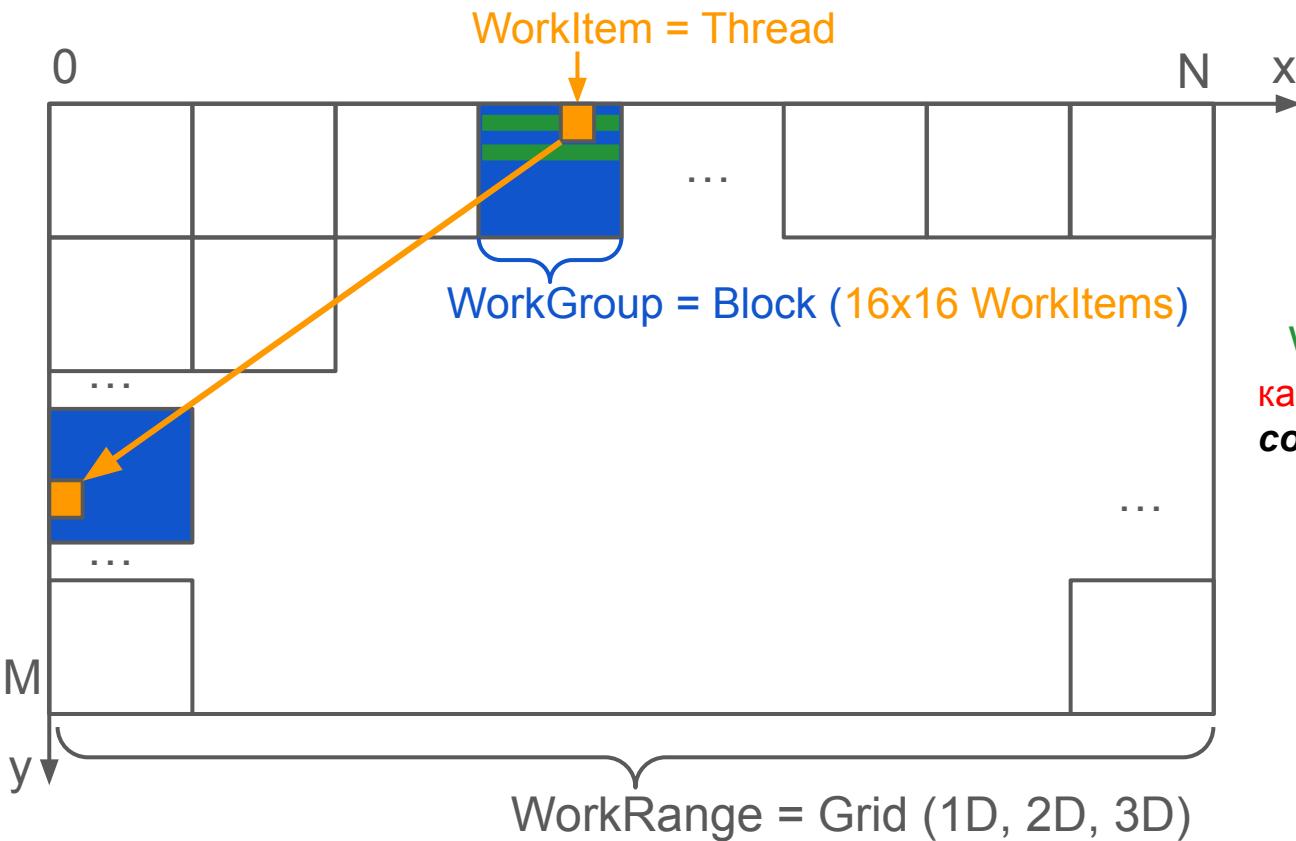
Транспонирование матрицы (*coalesced memory access*)



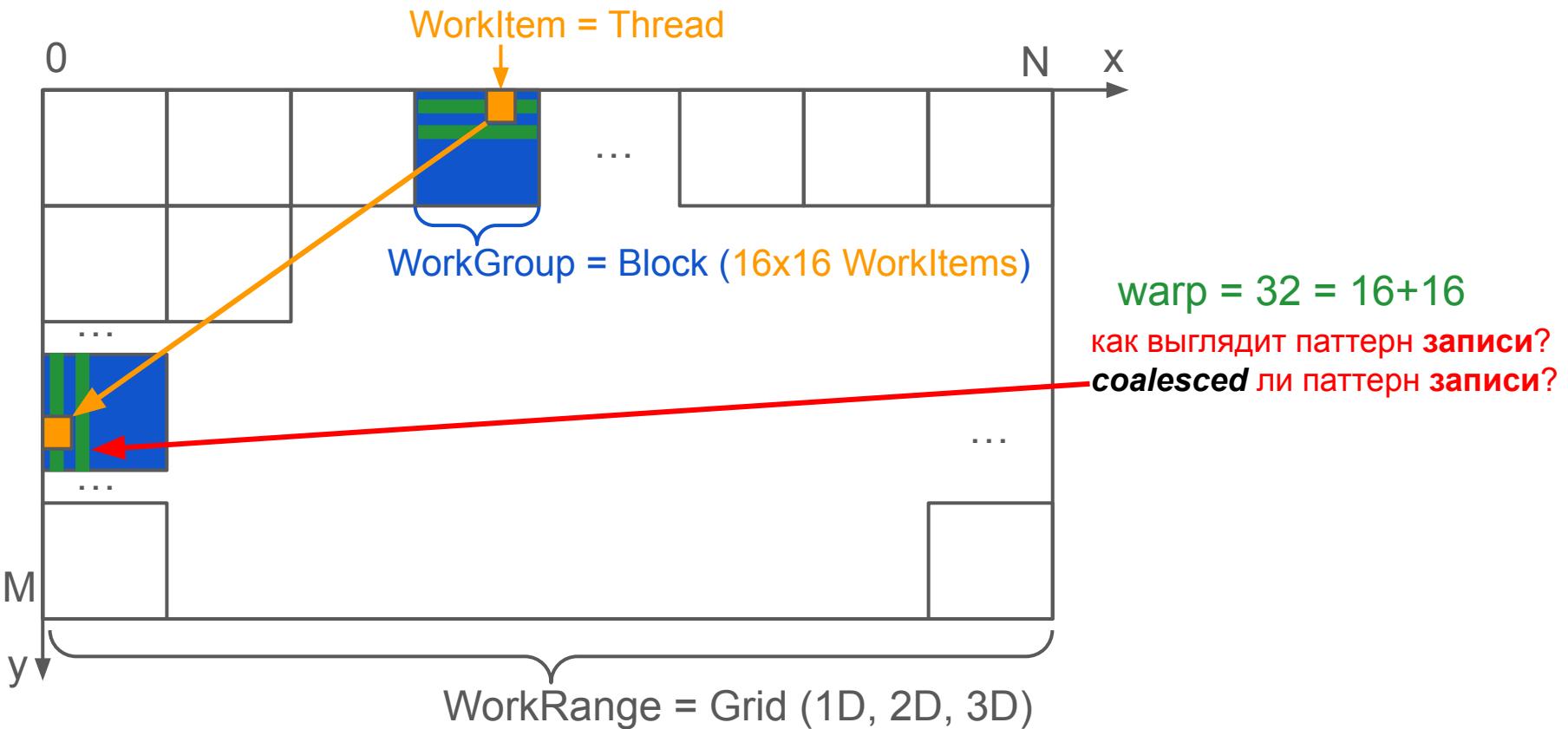
Транспонирование матрицы (*coalesced memory access*)



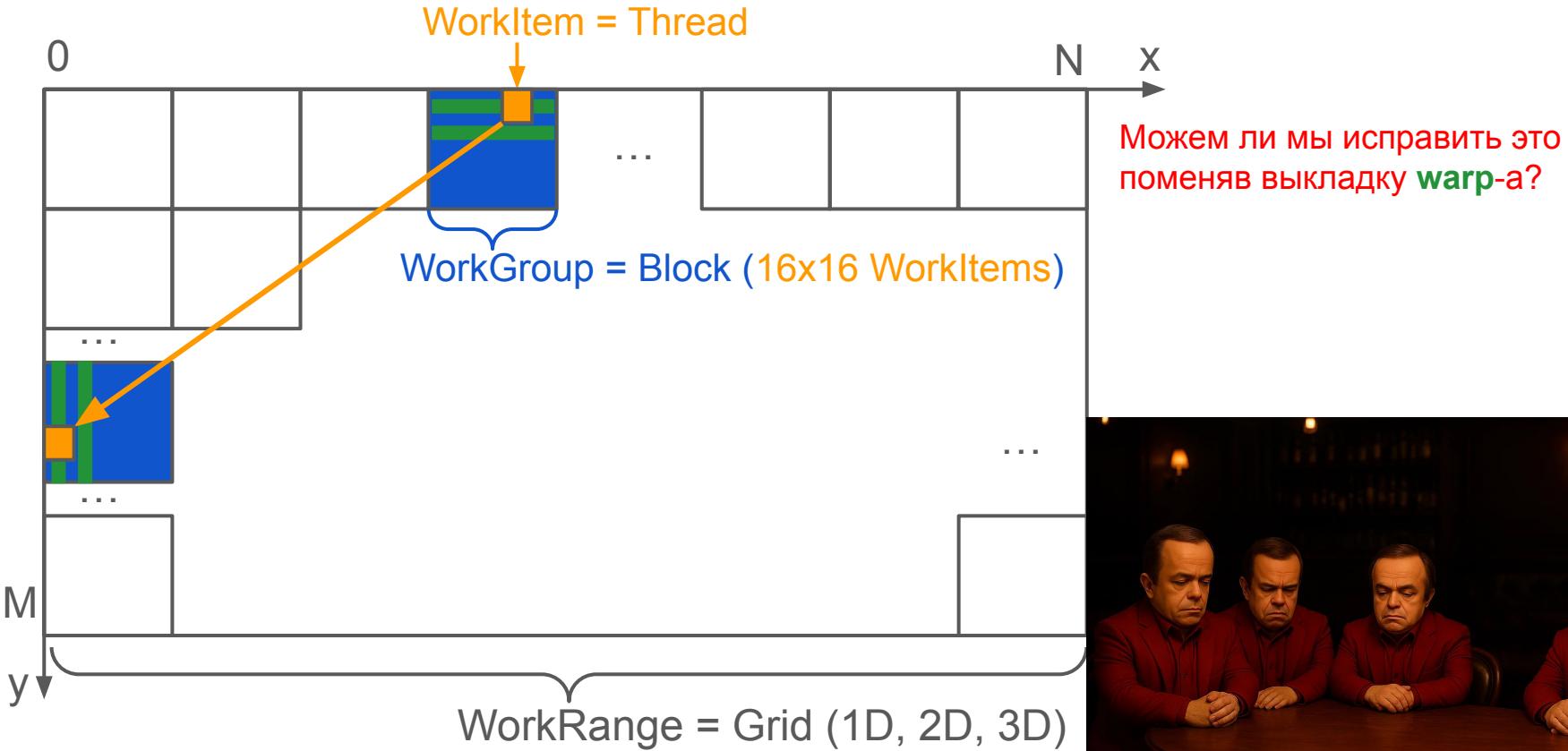
Транспонирование матрицы (*coalesced memory access*)



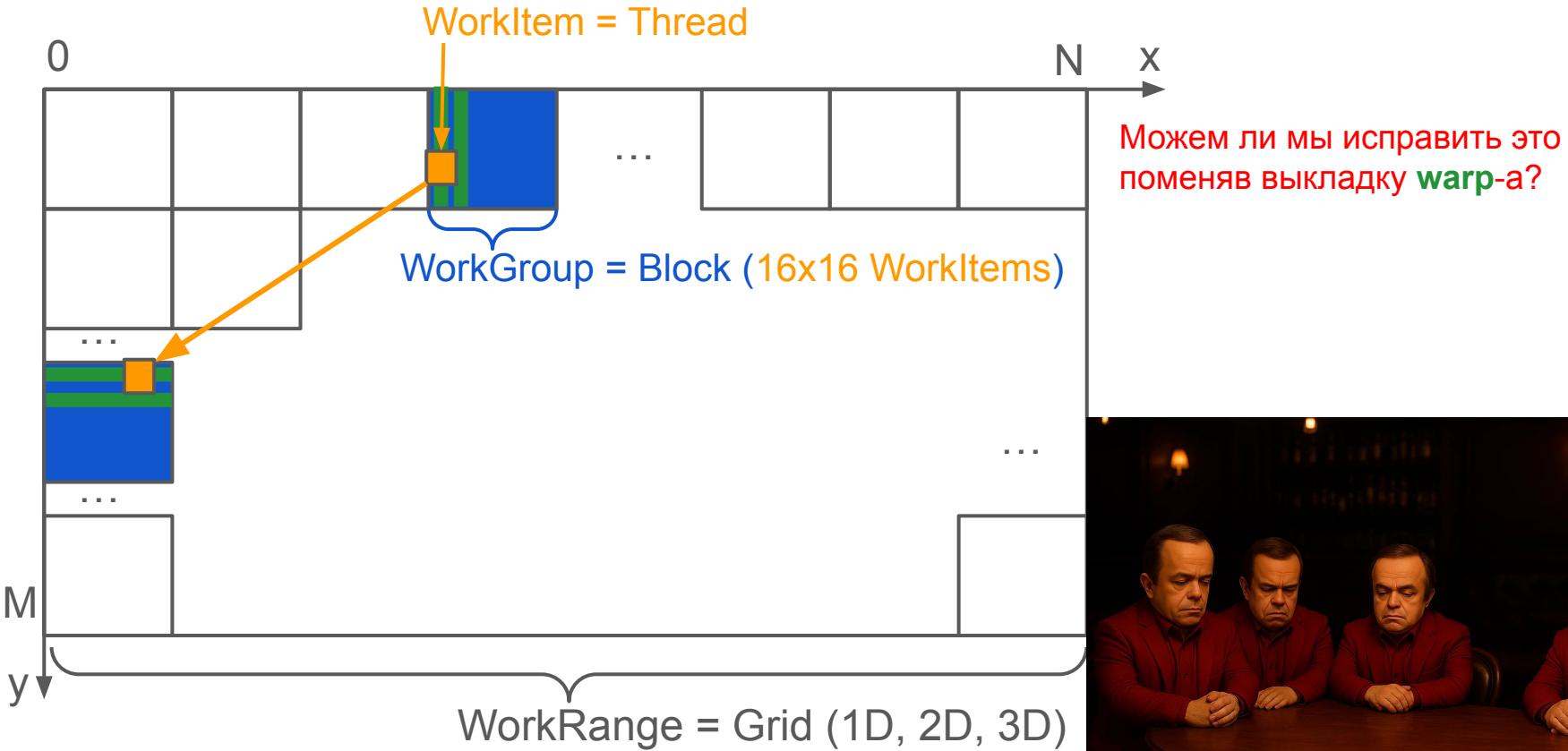
Транспонирование матрицы (*coalesced memory access*)



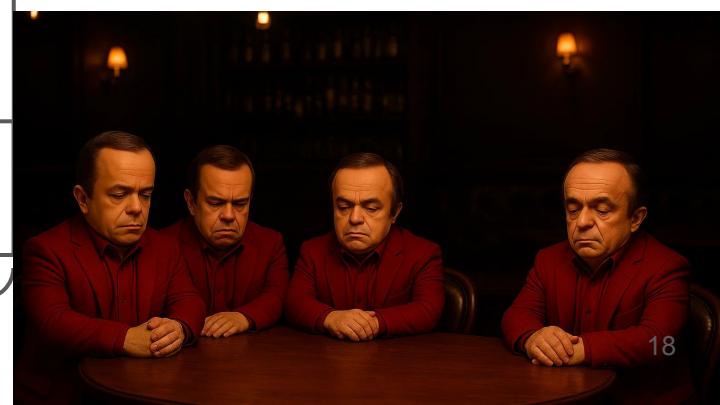
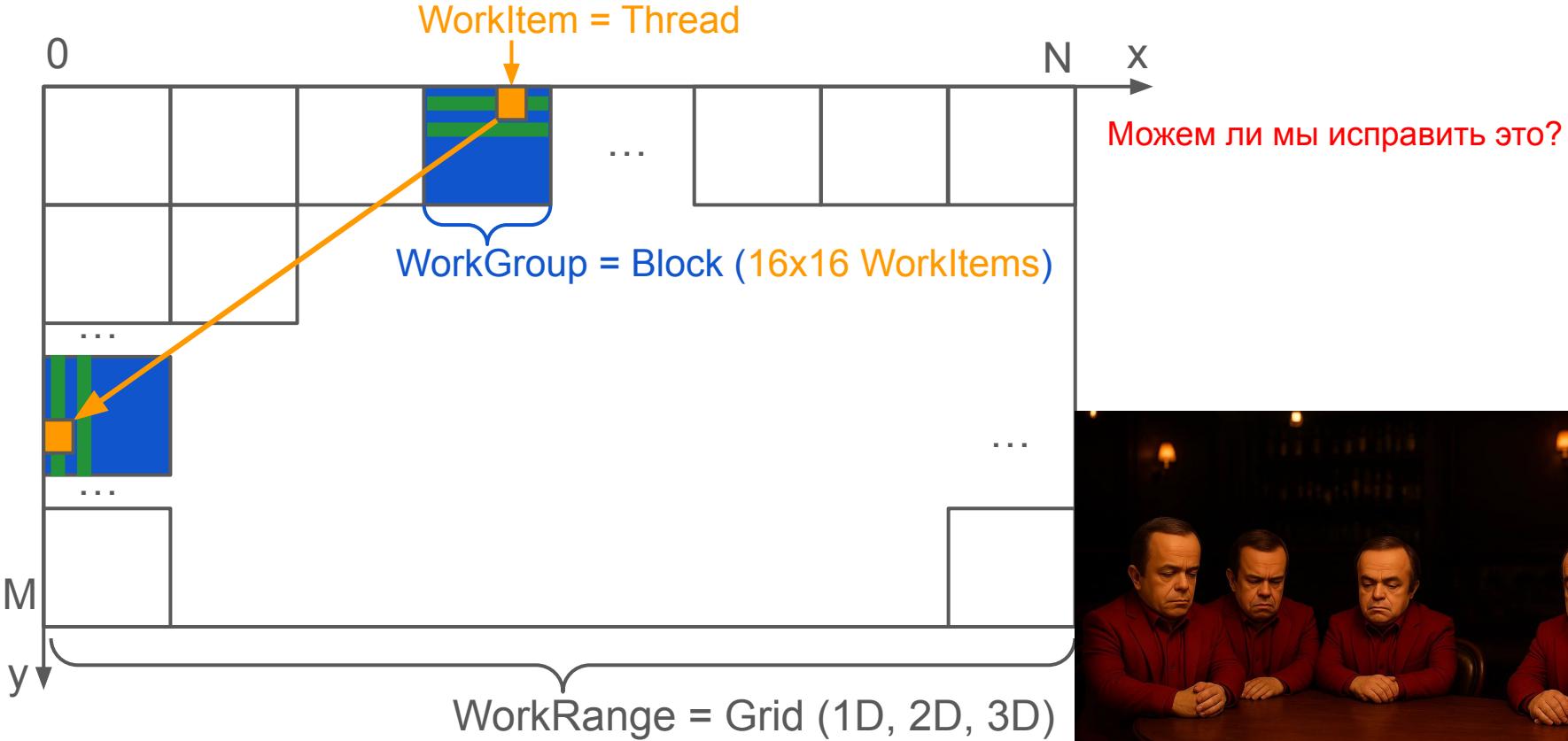
Транспонирование матрицы (*coalesced memory access*)



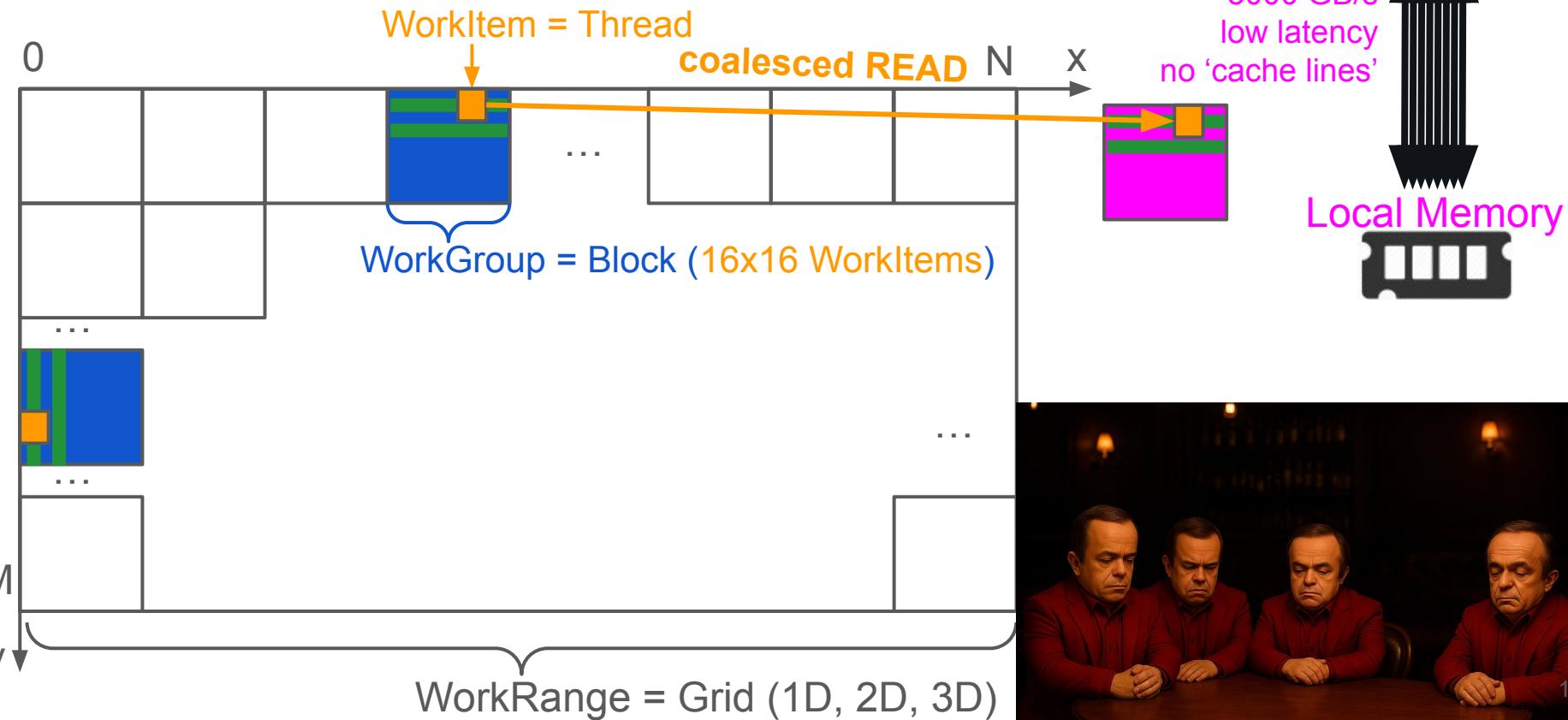
Транспонирование матрицы (*coalesced memory access*)



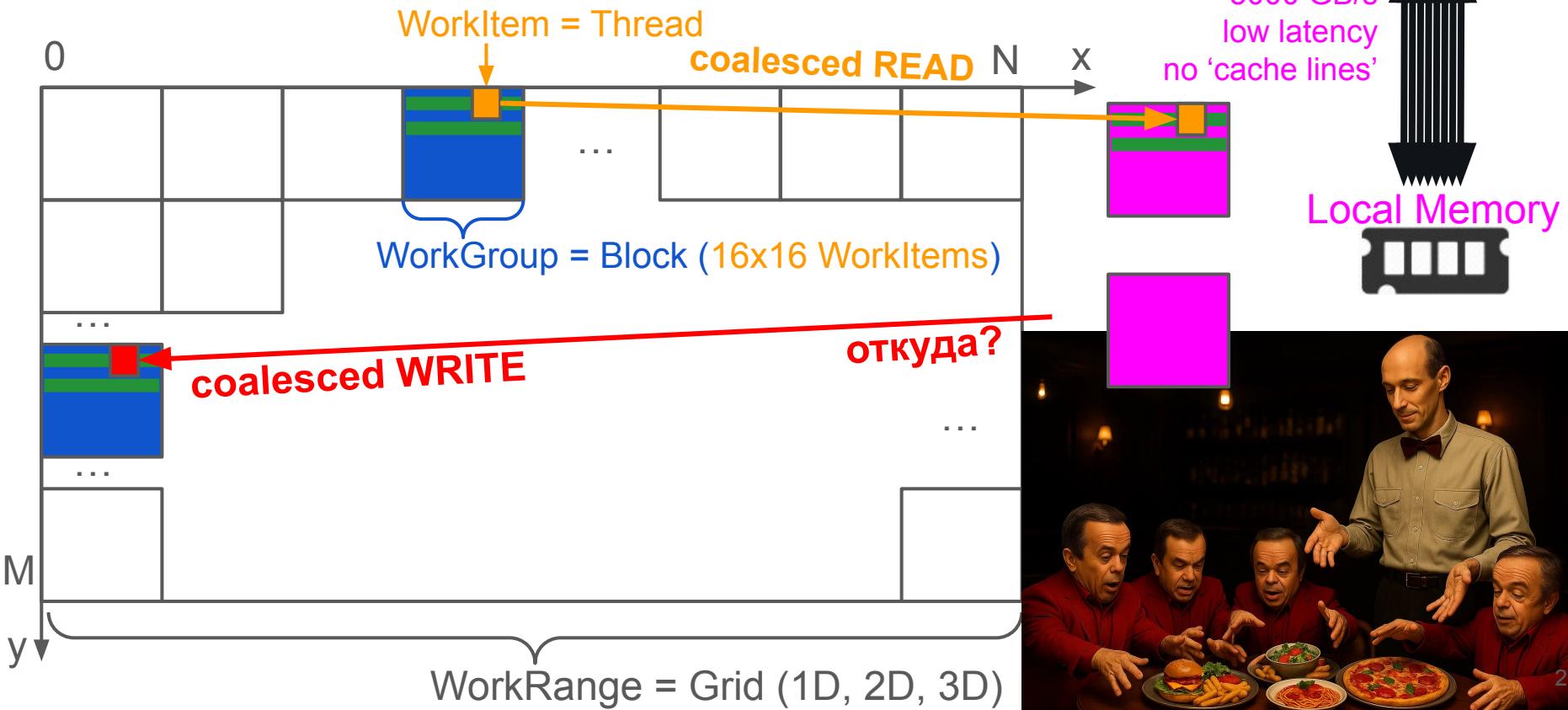
Транспонирование матрицы (*coalesced memory access*)



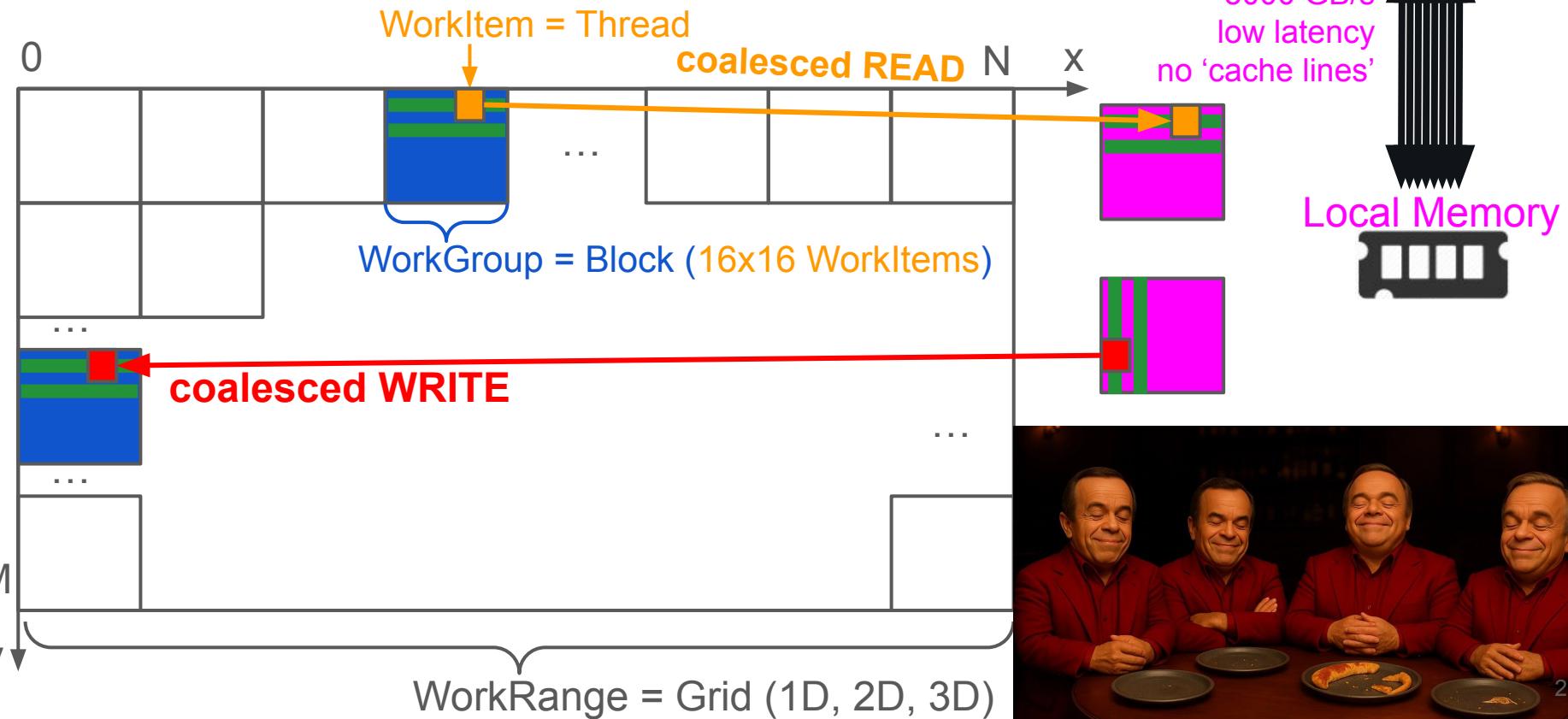
Транспонирование матрицы (coalesced)



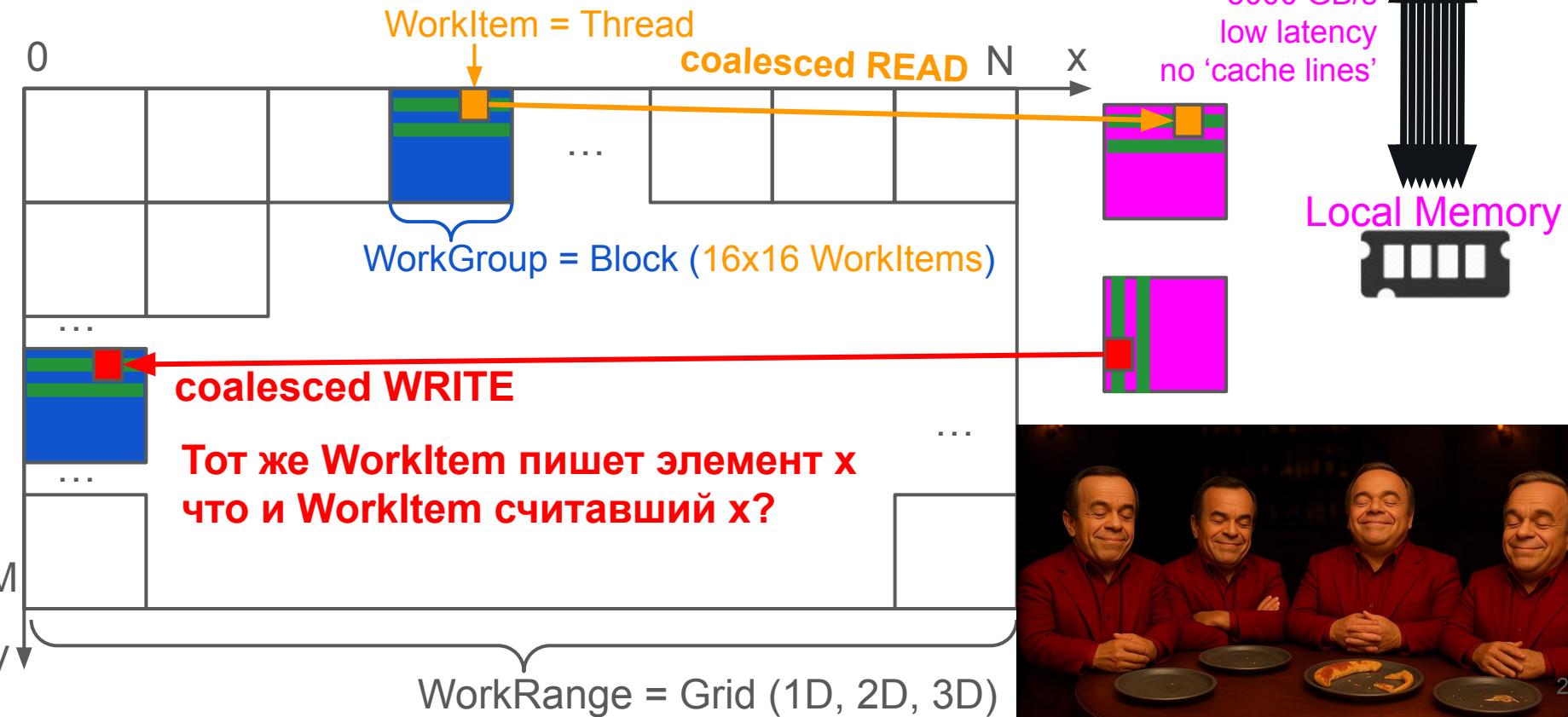
Транспонирование матрицы (coalesced)



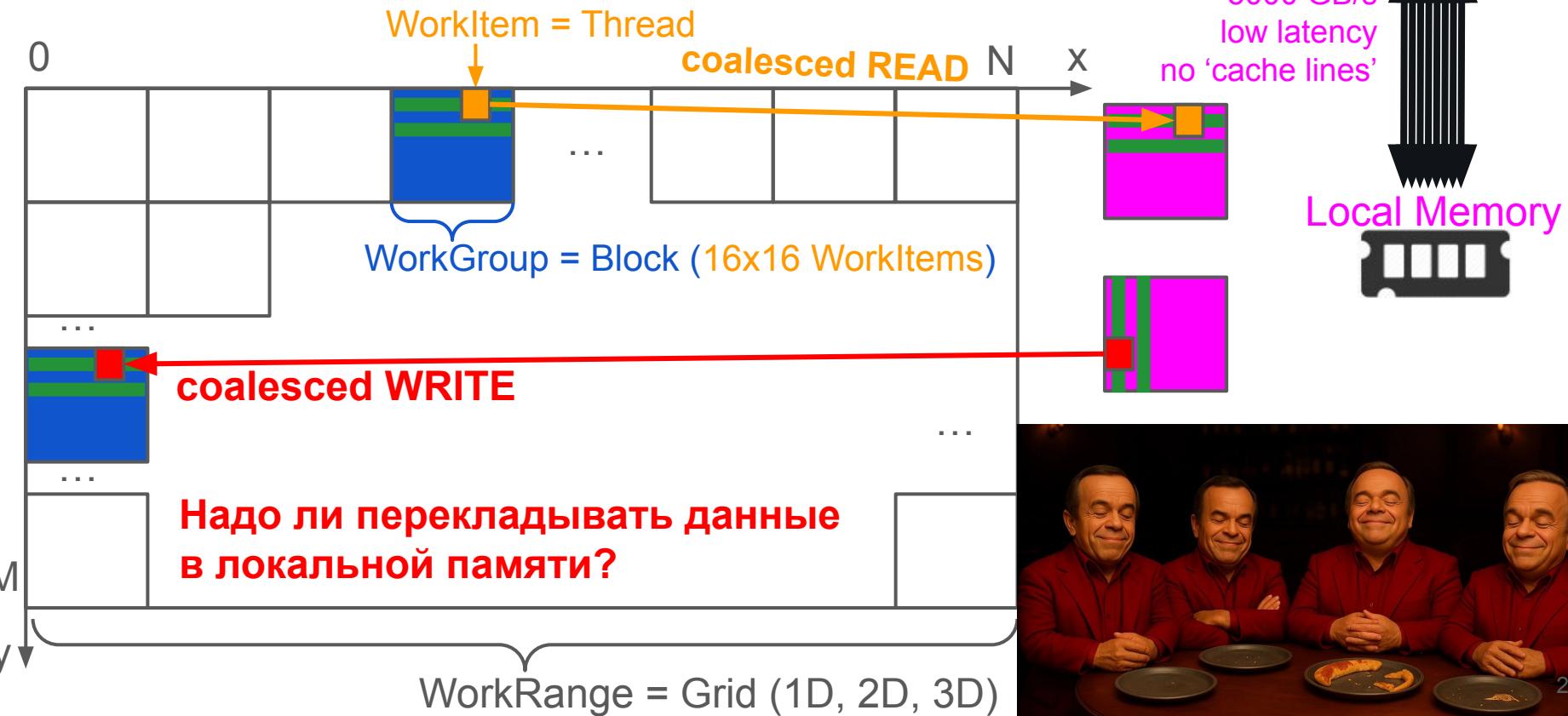
Транспонирование матрицы (coalesced)



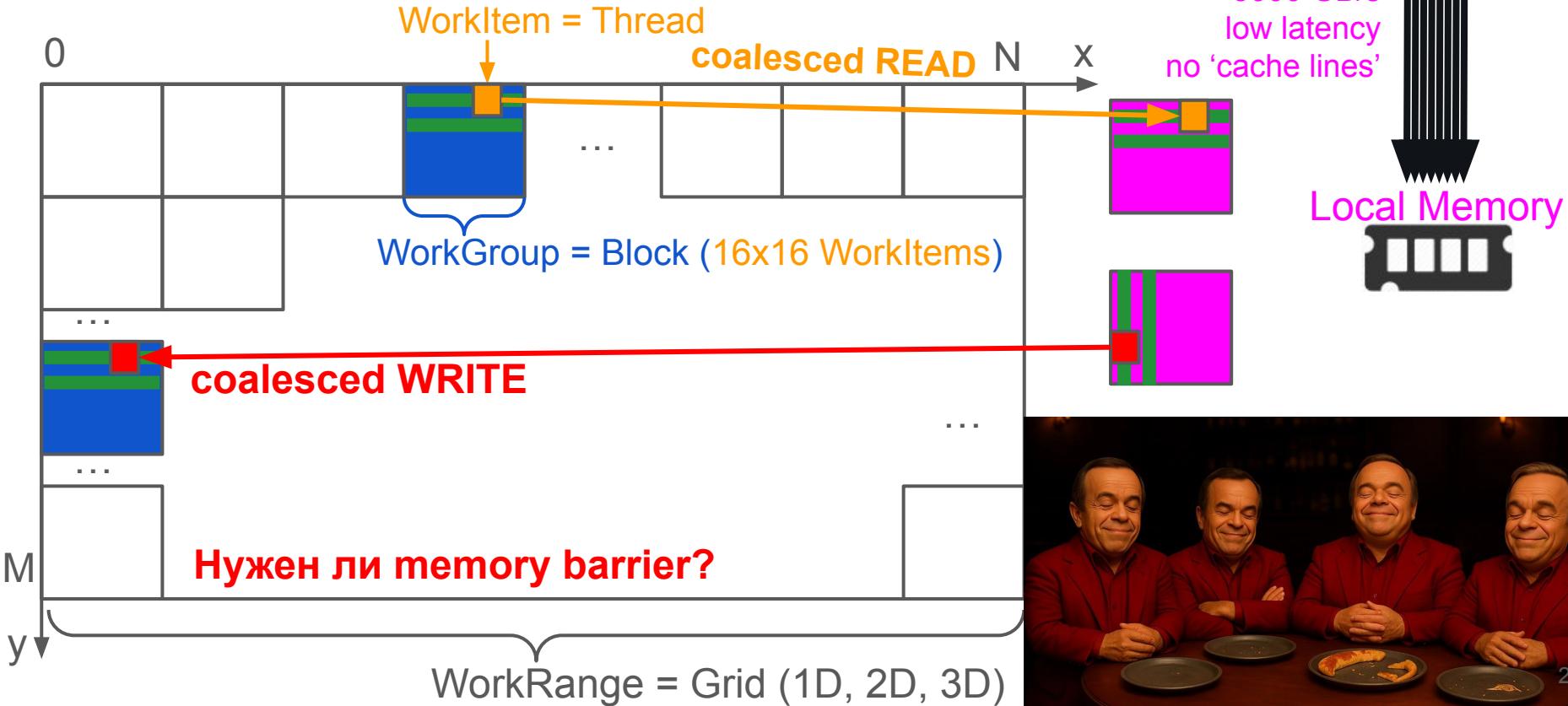
Транспонирование матрицы (coalesced)



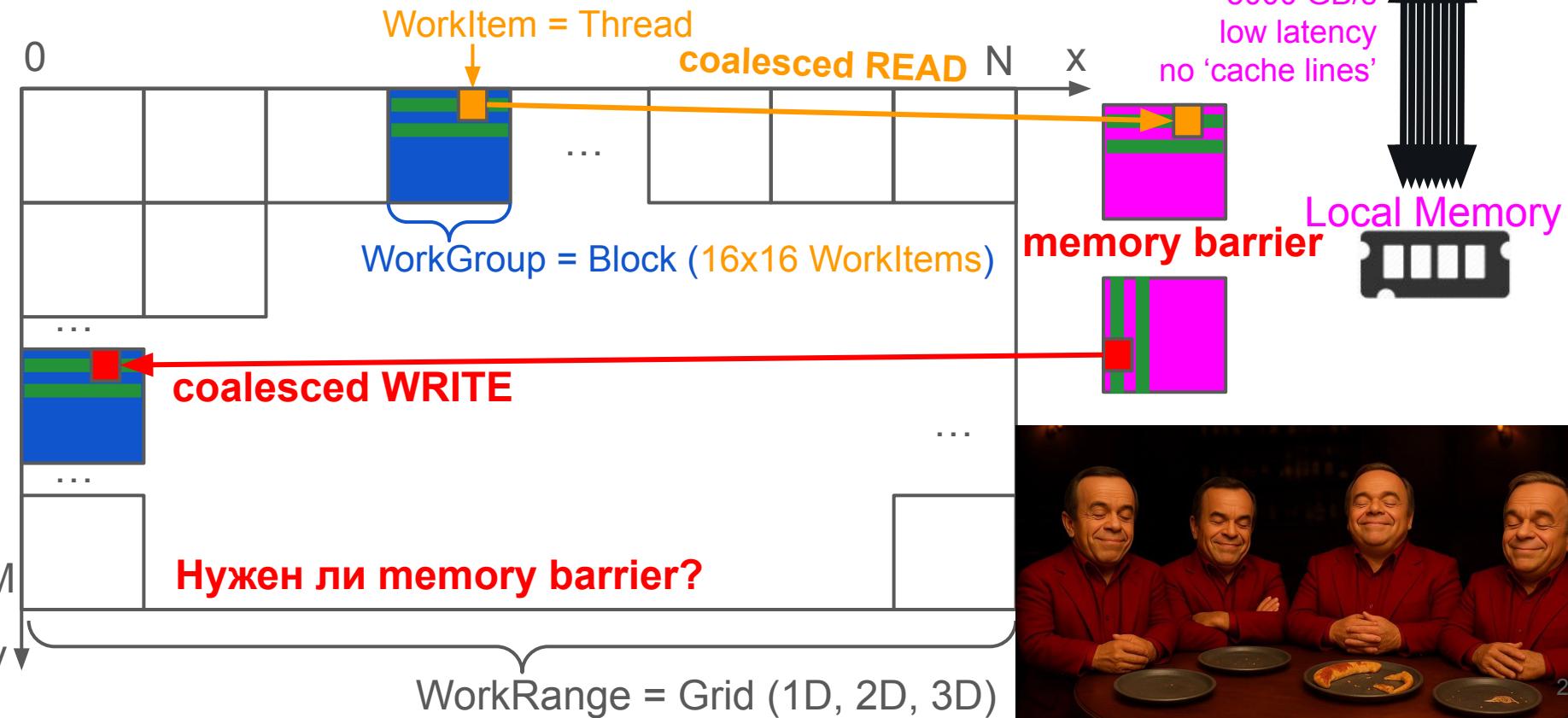
Транспонирование матрицы (coalesced)



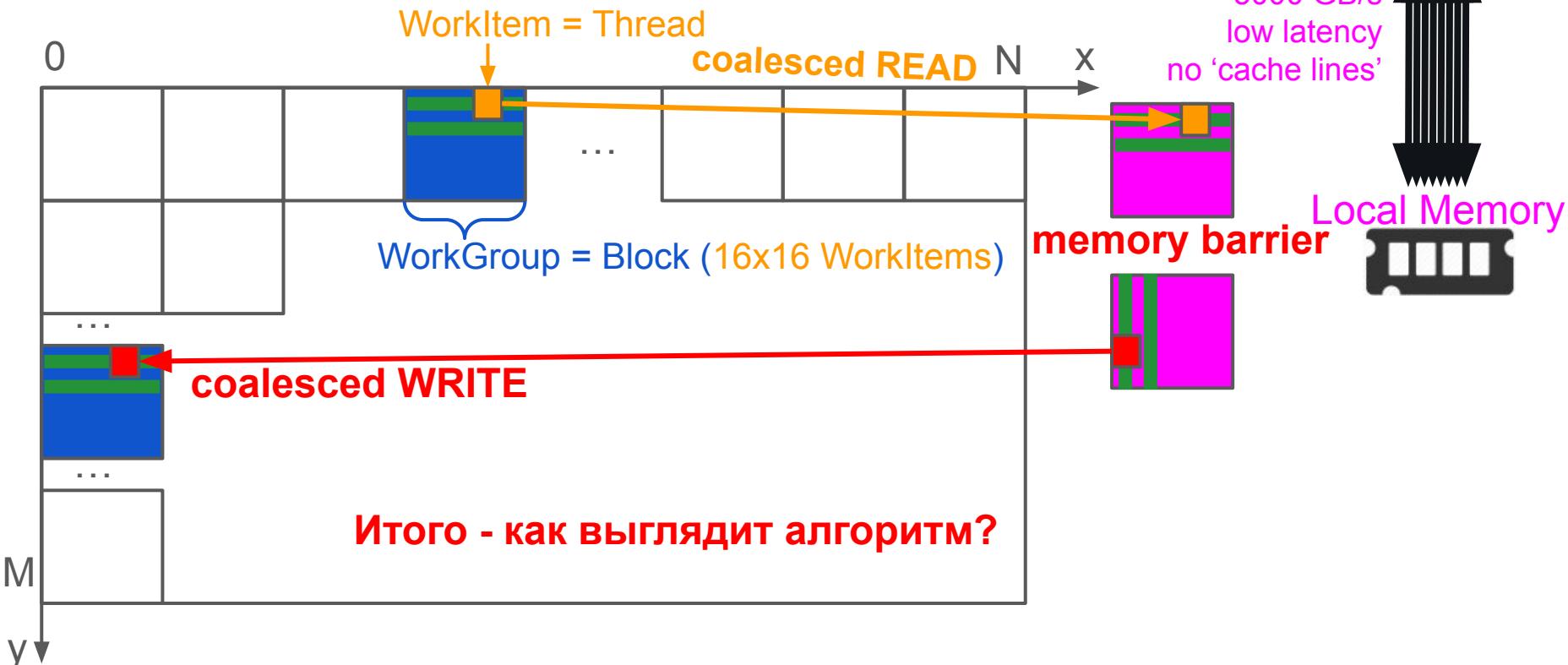
Транспонирование матрицы (*coalesced*)



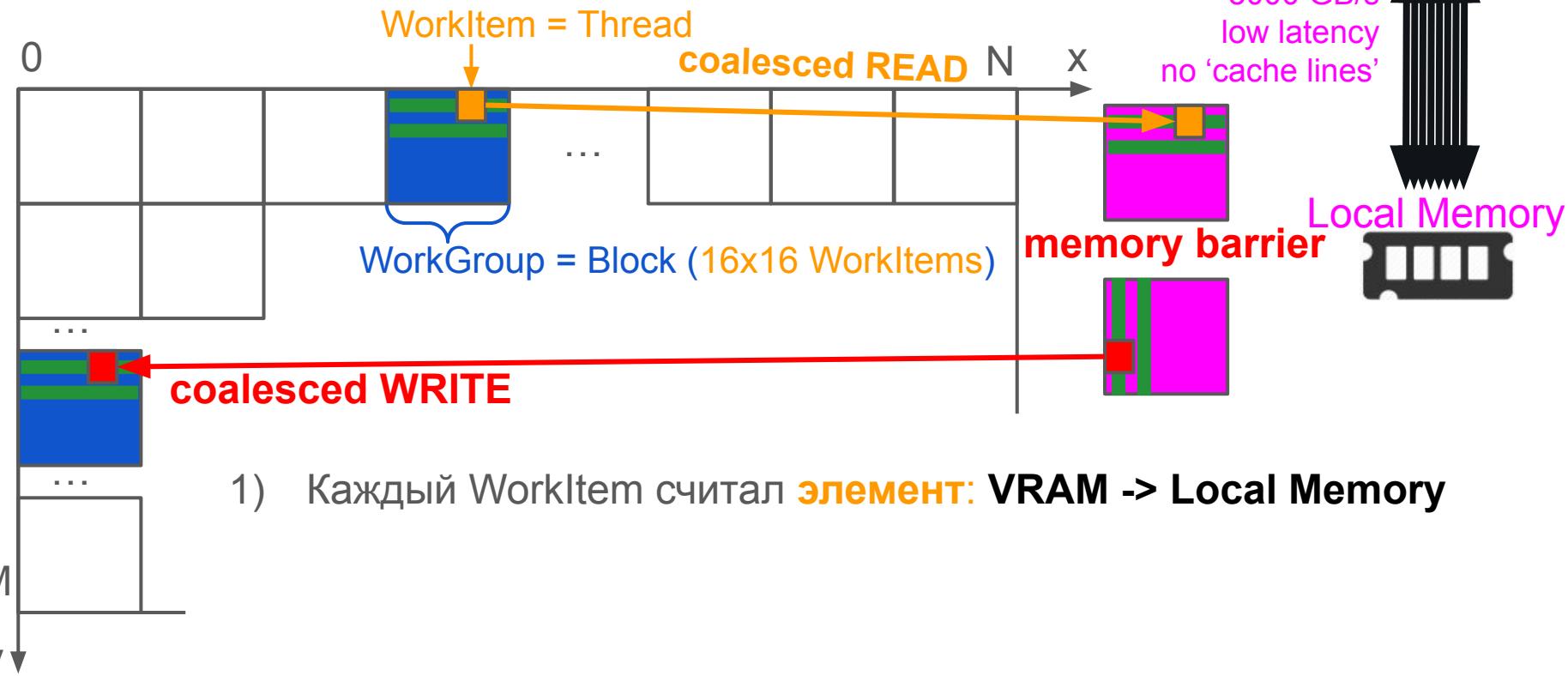
Транспонирование матрицы (coalesced)



Транспонирование матрицы (coalesced)

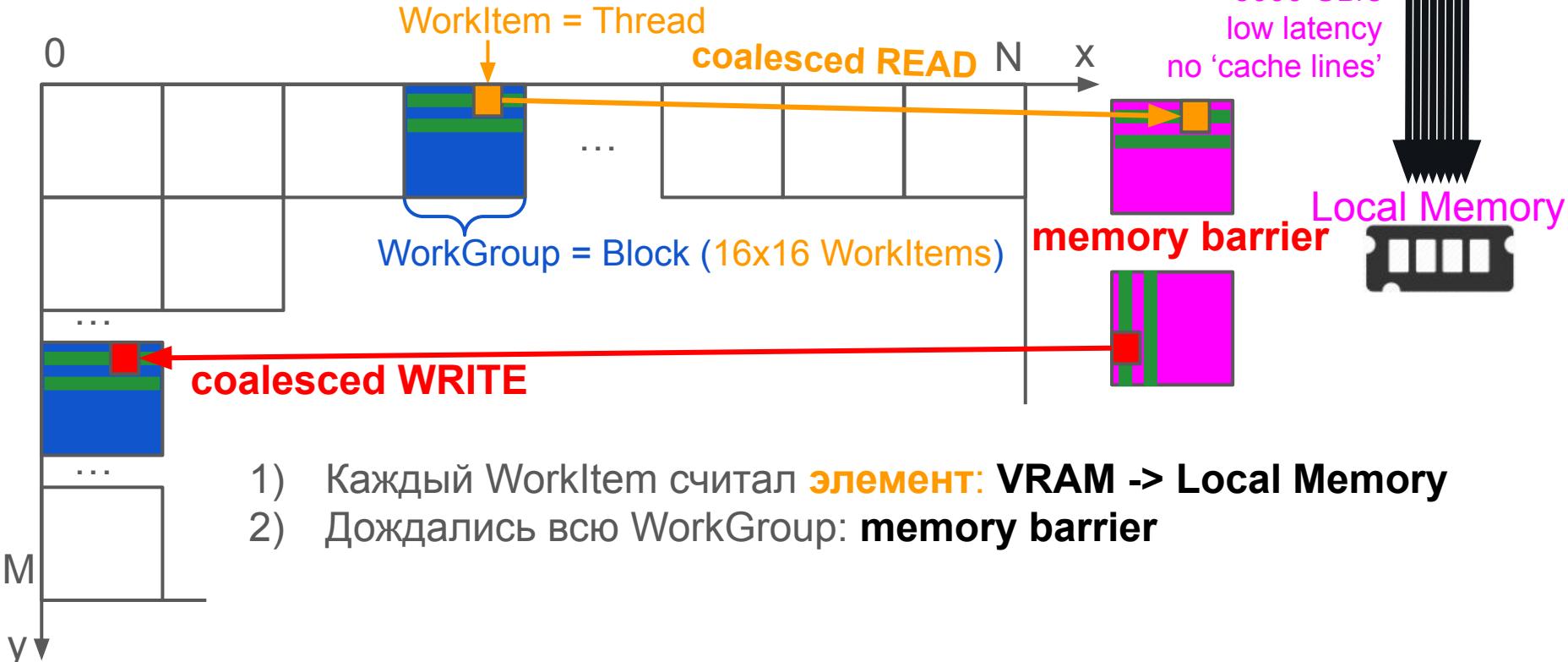


Транспонирование матрицы (coalesced)



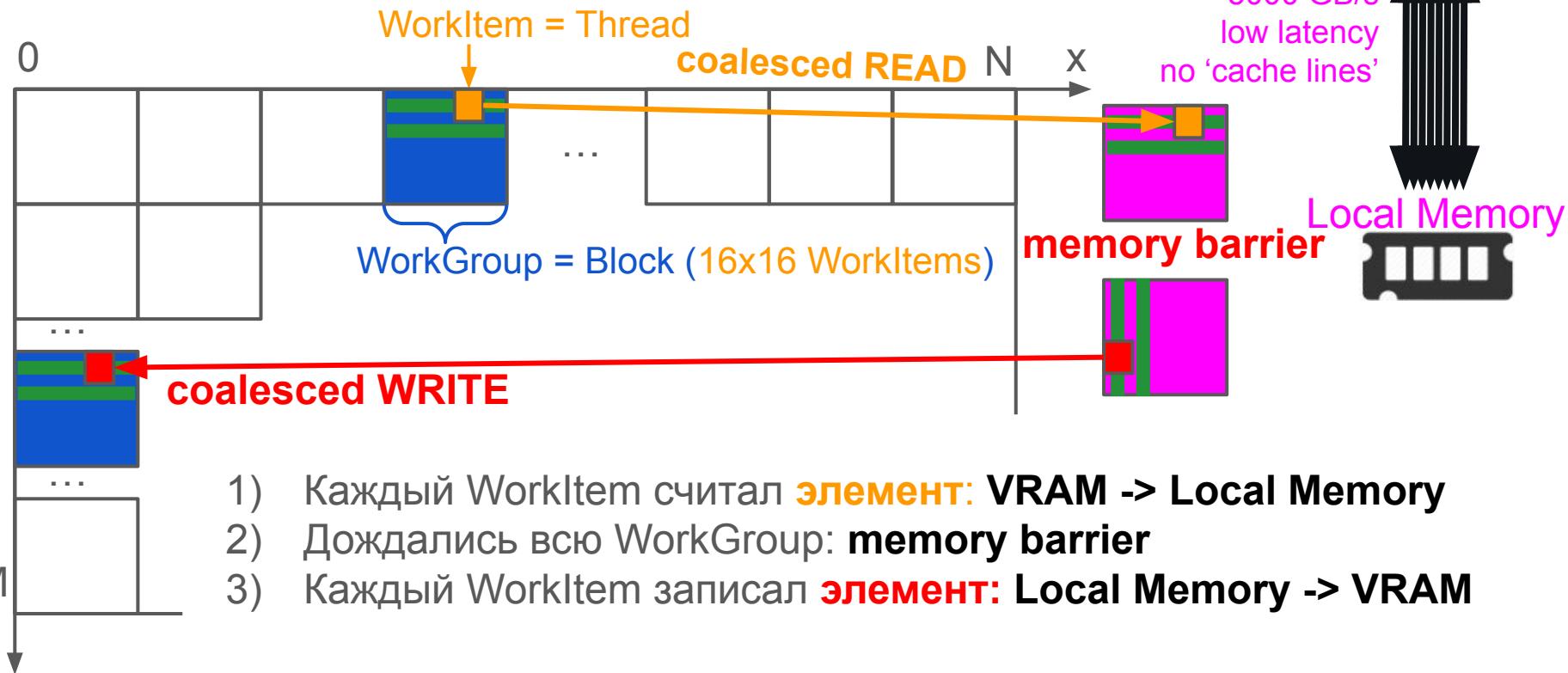
- 1) Каждый WorkItem считал **элемент**: VRAM \rightarrow Local Memory

Транспонирование матрицы (*coalesced*)

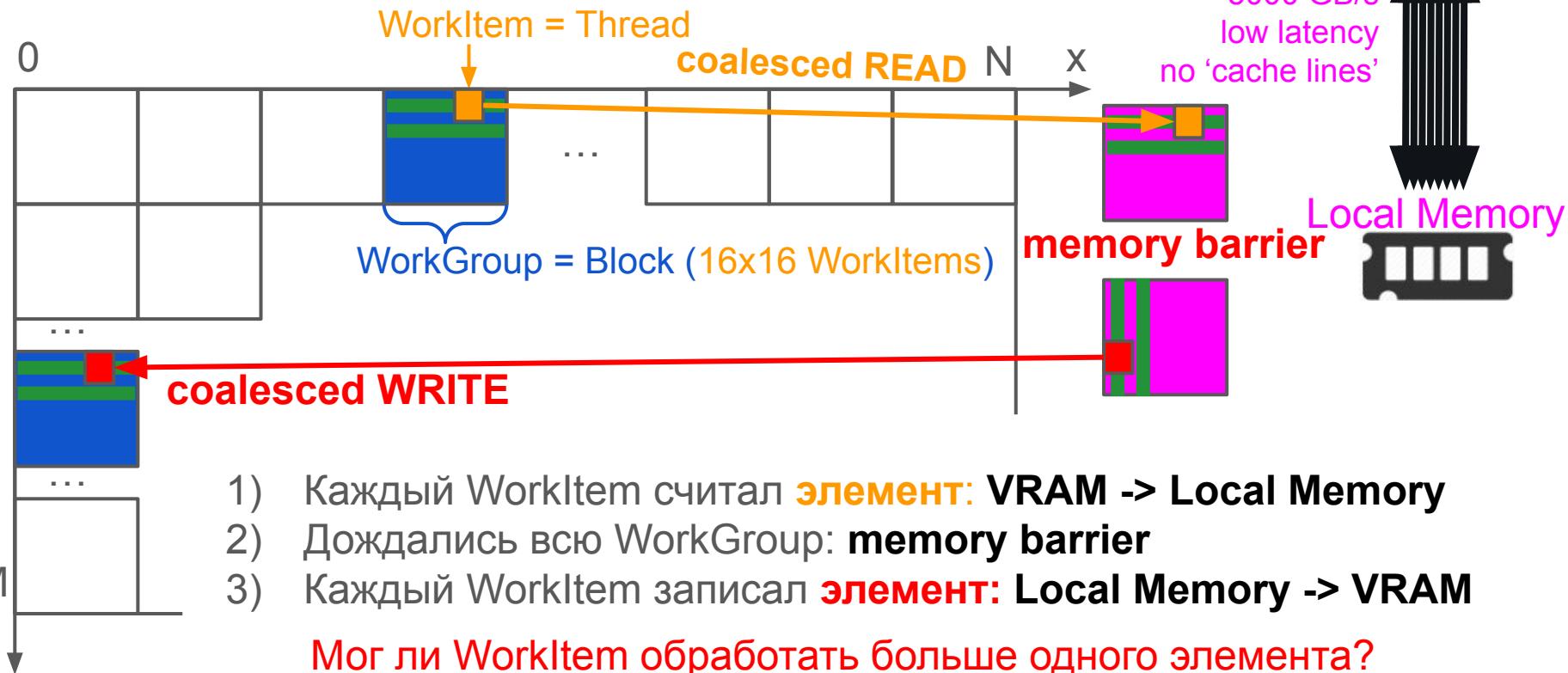


- 1) Каждый WorkItem считал **элемент**: VRAM -> Local Memory
 - 2) Дождались всю WorkGroup: **memory barrier**

Транспонирование матрицы (coalesced)



Транспонирование матрицы (coalesced)



Транспонирование

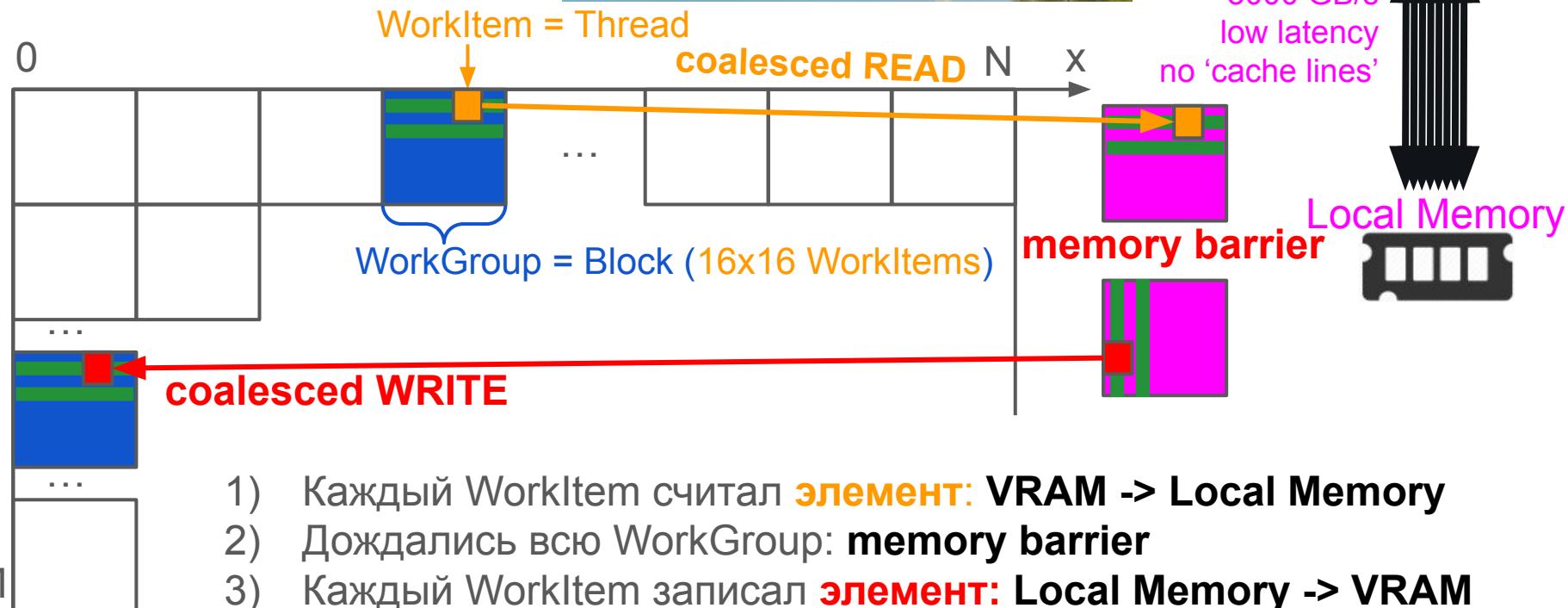


5000 GB/s
low latency
no 'cache lines'



Local Memory

memory barrier



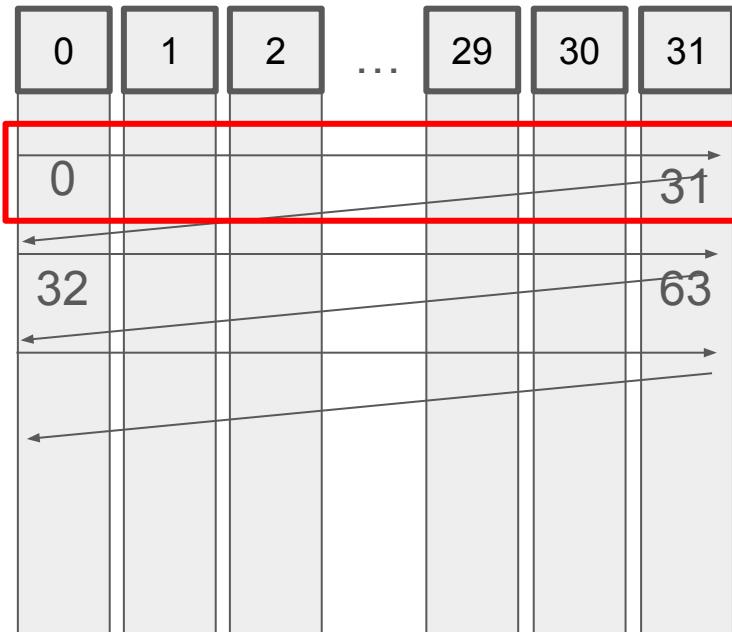
- 1) Каждый WorkItem считал **элемент**: VRAM -> Local Memory
- 2) Дождались всю WorkGroup: **memory barrier**
- 3) Каждый WorkItem записал **элемент**: Local Memory -> VRAM

Во что мы уперлись? Compute/Memory - bound?

Локальная память - bank conflicts (напоминание)

```
__local unsigned int workgroup_data[256];  
workgroup_data[get_local_id(0)] = a[index];
```

Local memory banks:



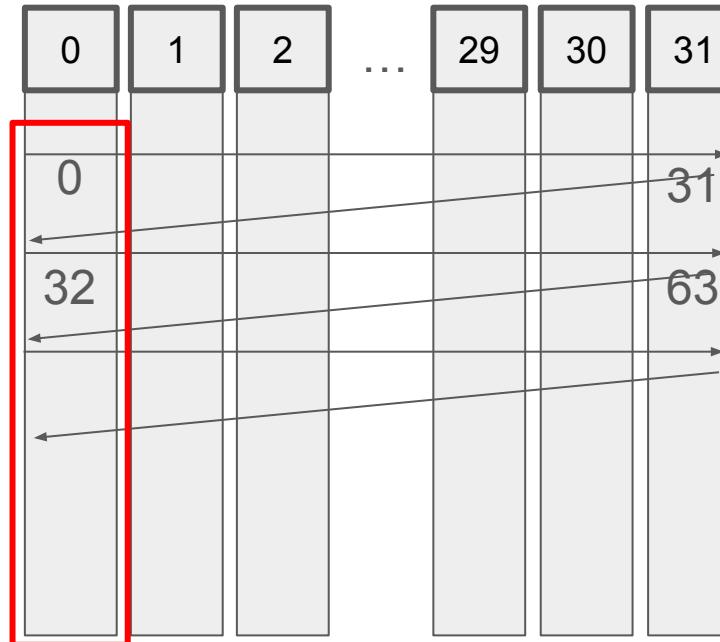
Local address space:

Локальная память - bank conflicts (напоминание)

```
__local unsigned int workgroup_data[256];  
workgroup_data[32 * get_local_id(0)] = a[index];
```

Local memory banks:

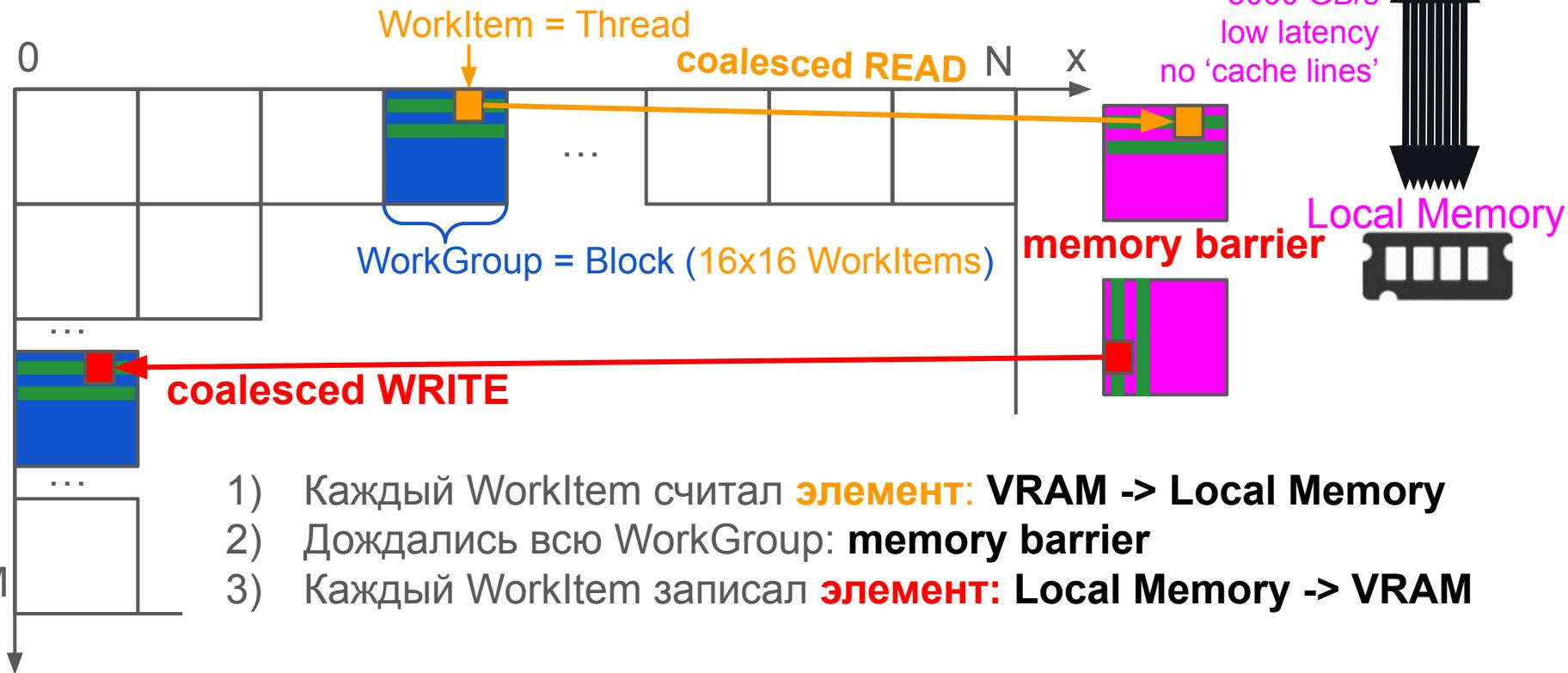
Local address space:



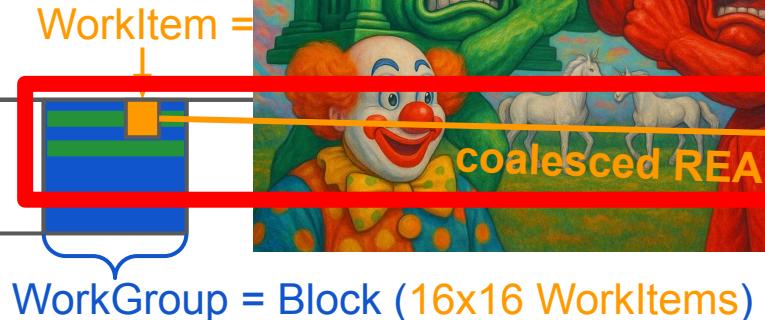
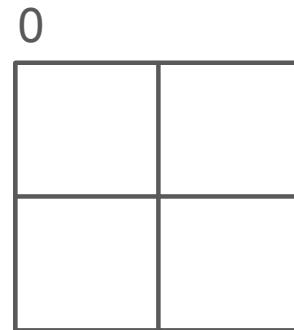
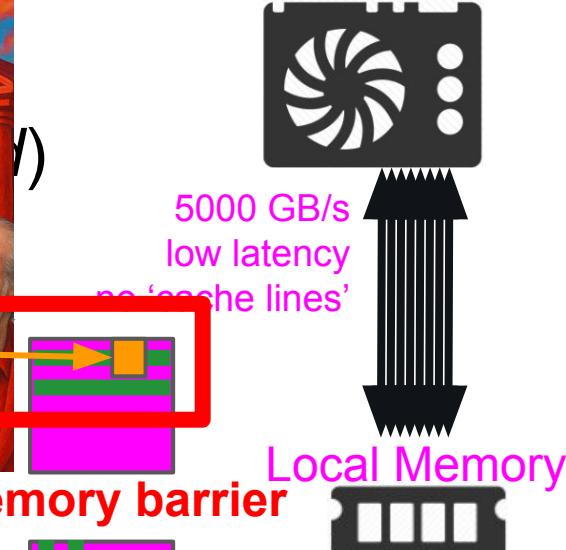
Bank conflict!



Транспонирование матрицы (coalesced)



Транспонирование

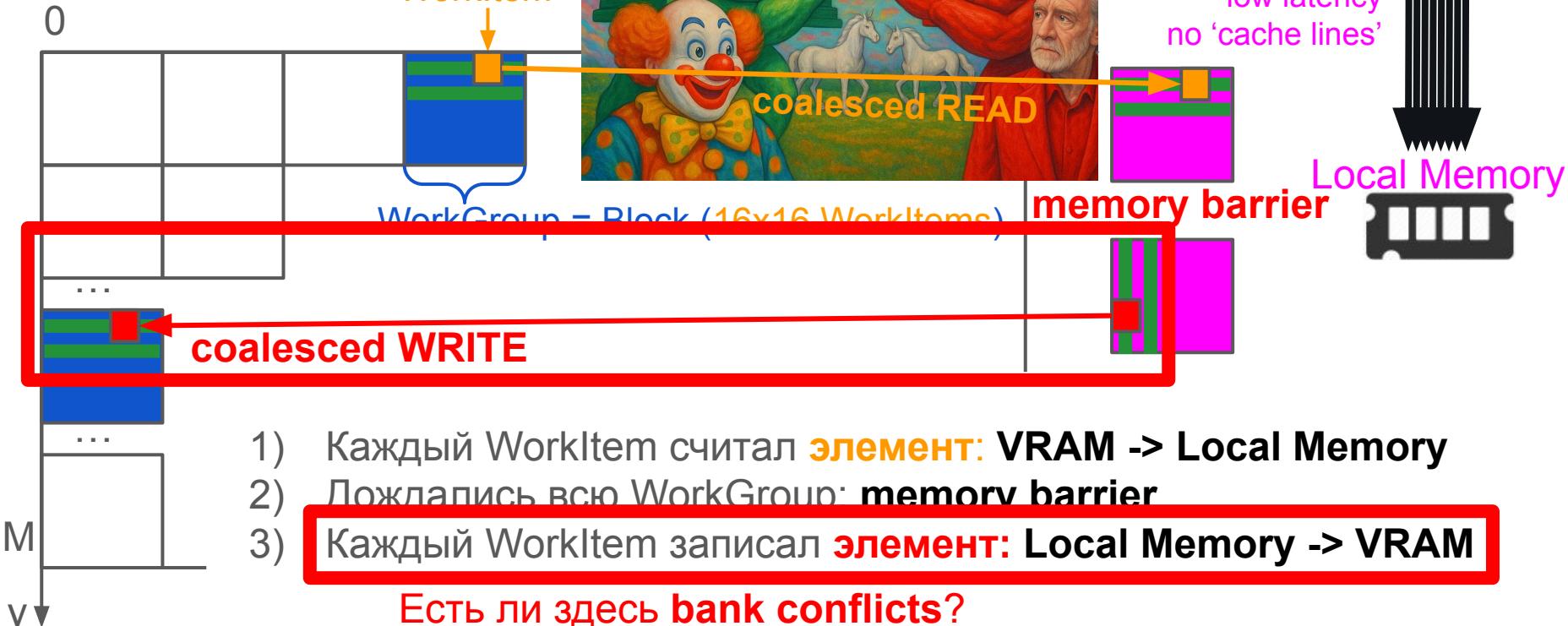


coalesced WRITE

- 1) Каждый WorkItem считал **элемент**: VRAM \rightarrow Local Memory
- 2) дождались всю workGroup: **memory barrier**
- 3) Каждый WorkItem записал **элемент**: Local Memory \rightarrow VRAM

Есть ли здесь **bank conflicts**?

Транспонирование

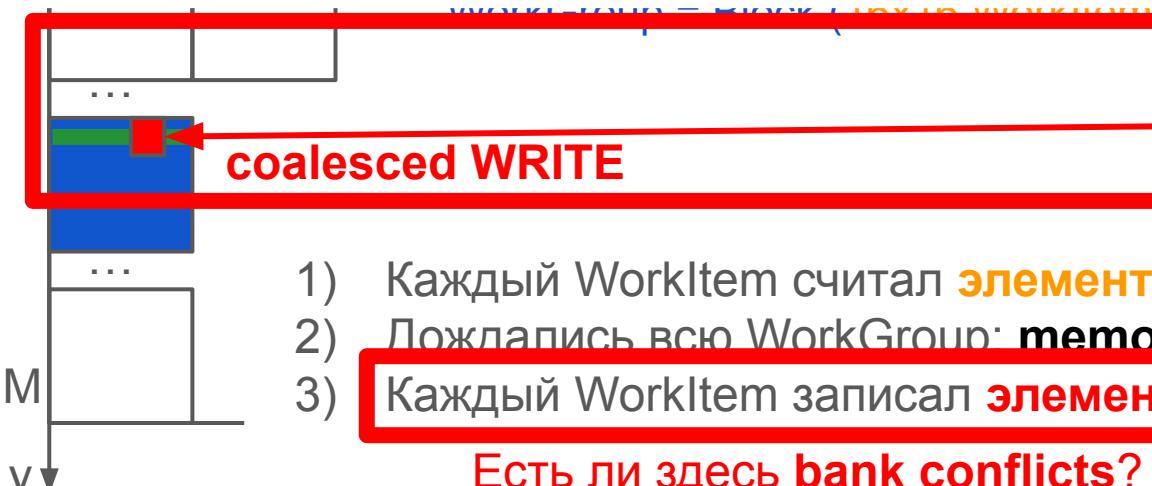




5000 GB/s
low latency
no 'cache lines'

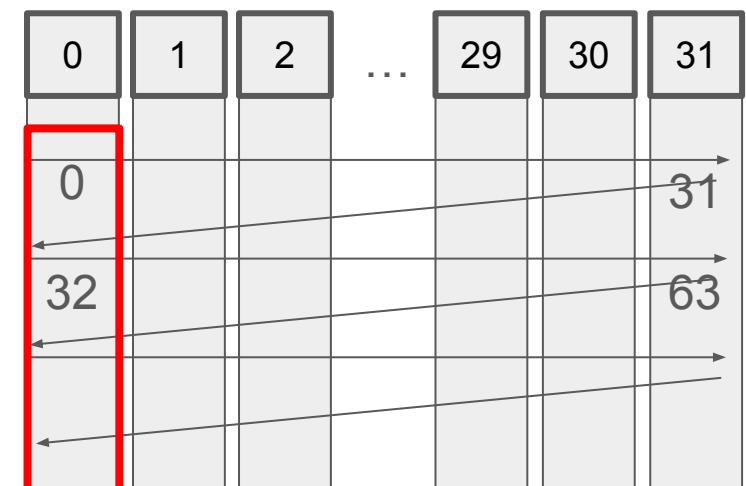
Для простоты будем считать
что WorkGroup **32 x 32!**

И Tile в Local Memory - **32 x 32!**



- 1) Каждый WorkItem считал **элемент**: VRAM \rightarrow Local Memory
- 2) Дождалась всю WorkGroup: **memory barrier**
- 3) Каждый WorkItem записал **элемент**: Local Memory \rightarrow VRAM

Есть ли здесь **bank conflicts**?

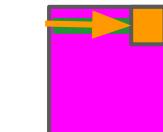
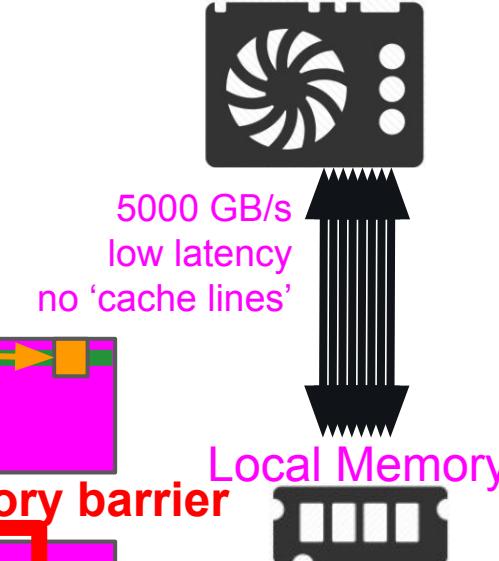


У каждого WorkItem есть свой набор из 4-х банков

...
...

coalesced WRITE

M
y

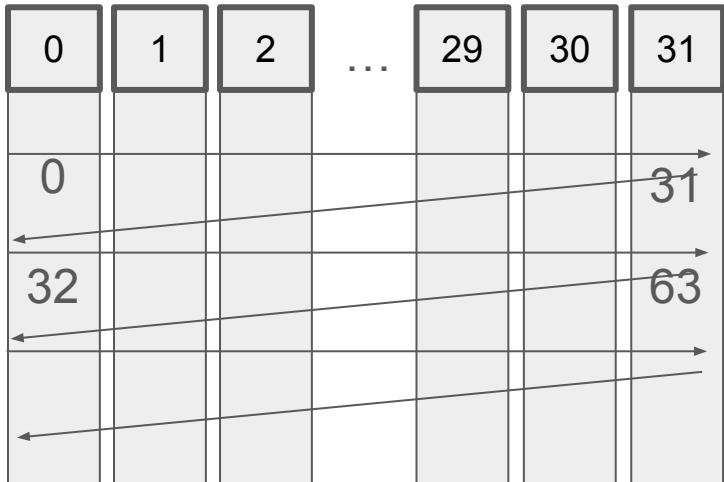


memory barrier



- 1) Каждый WorkItem считал **элемент**: VRAM \rightarrow Local Memory
- 2) Дождалась всю WorkGroup: **memory barrier**
- 3) Каждый WorkItem записал **элемент**: Local Memory \rightarrow VRAM

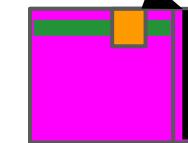
Есть ли здесь **bank conflicts**?



Что будет если
Tile - 33x32?

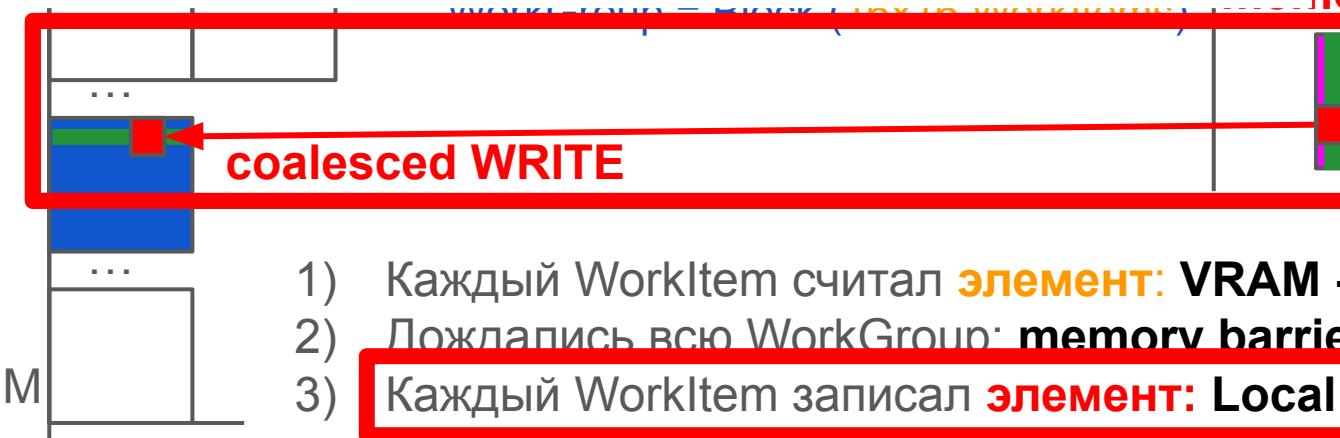


5000 GB/s
low latency
no 'cache lines'



Local Memory
memory barrier

coalesced WRITE

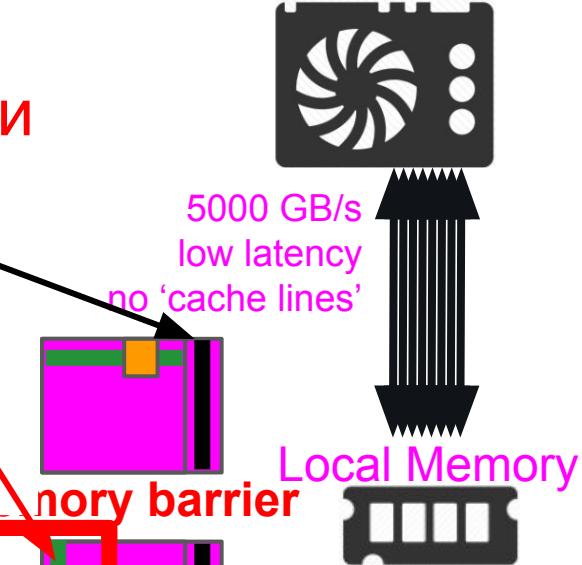


Есть ли здесь bank conflicts?



Что будет если
Tile - 33x32?

В какой банке второй
элемент столбика?

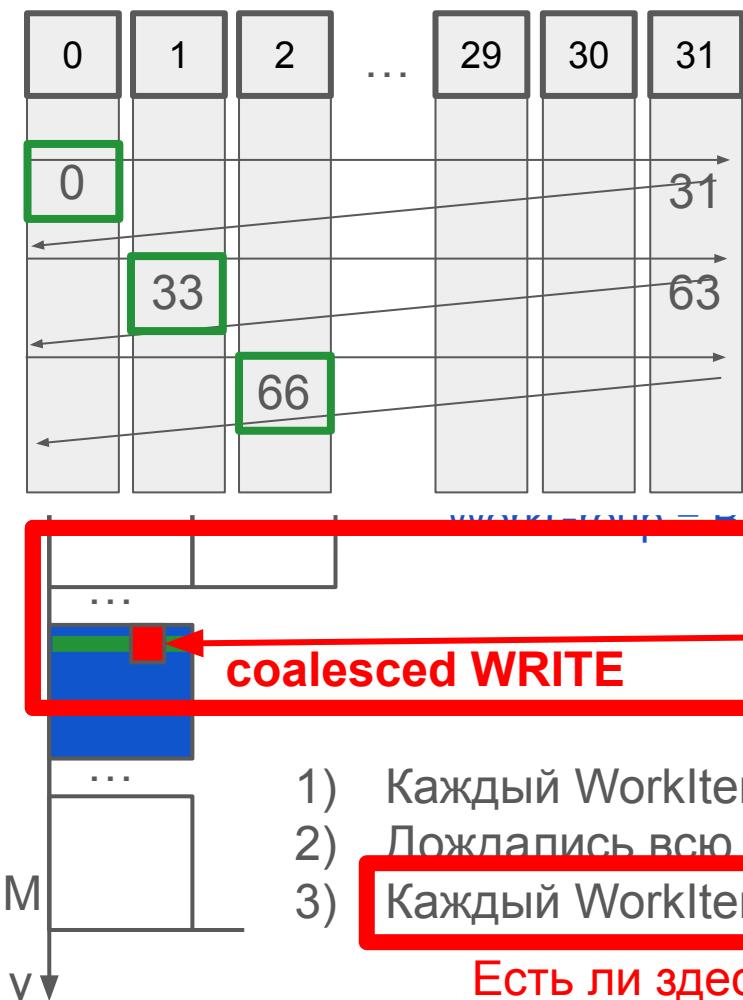


- 1) Каждый WorkItem считал **элемент**: VRAM -> Local Memory
- 2) Дождалась всю WorkGroup: **memory barrier**
- 3) Каждый WorkItem записал **элемент**: Local Memory -> VRAM

Есть ли здесь **bank conflicts**?



- 1) Каждый WorkItem считал **элемент**: VRAM \rightarrow Local Memory
- 2) Дождалась всю WorkGroup: **memory barrier**
- 3) Каждый WorkItem записал **элемент**: Local Memory \rightarrow VRAM



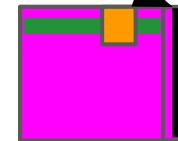
Что будет если
Tile - 33x32?



5000 GB/s

low latency

no 'cache lines'



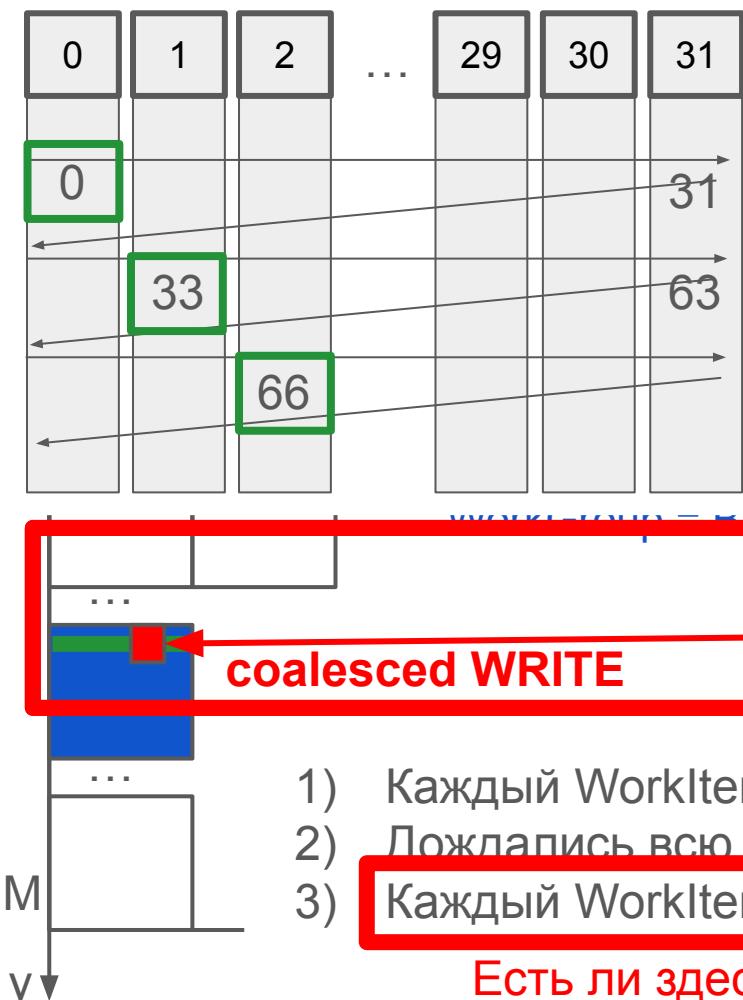
memory barrier



Local Memory

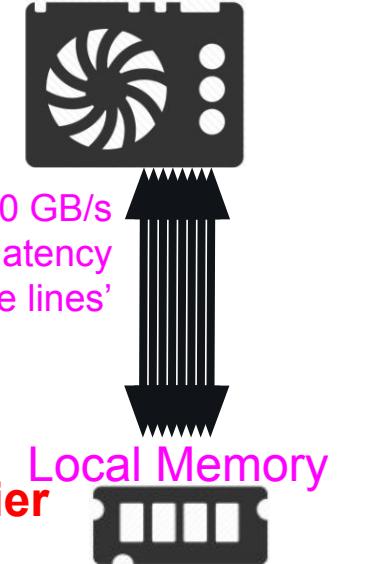
- 1) Каждый WorkItem считал **элемент**: VRAM -> Local Memory
- 2) Дождались всю WorkGroup: **memory barrier**
- 3) Каждый WorkItem записал **элемент**: Local Memory -> VRAM

Есть ли здесь **bank conflicts**?



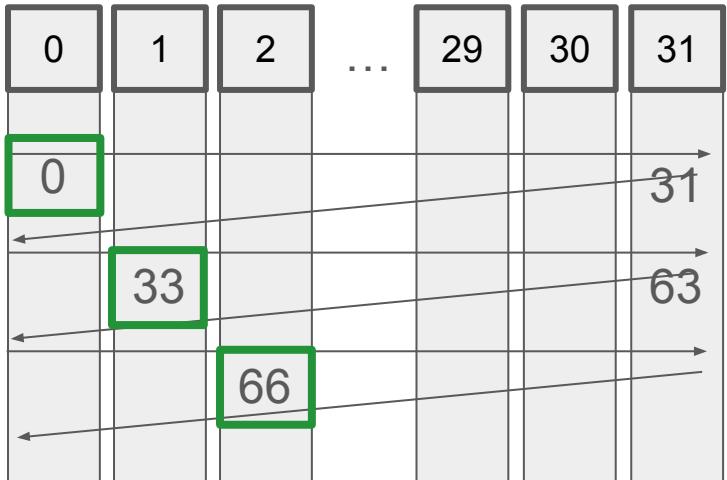
Что будет если
Tile - 33x32?

Можно ли без лишней
local memory?



- 1) Каждый WorkItem считал **элемент**: VRAM \rightarrow Local Memory
- 2) Дождалась всю WorkGroup: **memory barrier**
- 3) Каждый WorkItem записал **элемент**: Local Memory \rightarrow VRAM

Есть ли здесь **bank conflicts**?



Что будет если

Tile - 33x32?

А на что влияет
объем используемой
Local Memory?

coalesced WRITE

- 1) Каждый WorkItem считал **элемент**: VRAM ->
- 2) Дождалась всю WorkGroup: **memory barrier**
- 3) Каждый WorkItem записал **элемент**: Local M

Есть ли здесь **bank conflicts**?

SM - Streaming Multiprocessor



А на что влияет
объем используемой
Local Memory?
Occupancy!



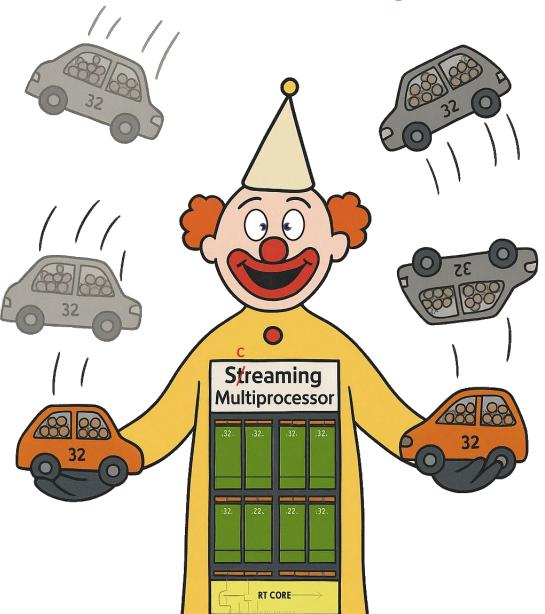
RTX 3090 (Ampere) - per SM:

- 128 KB Local Memory
 - 4 x 16384 x 32-bit registers
- = 256 KB Register File

SM - Streaming Multiprocessor



А на что влияет
объем используемой
Local Memory?
Occupancy!



RTX 3090 (Ampere) - per SM:

- **128 KB Local Memory**
- **4 x 16384 x 32-bit registers**
- = 256 KB Register File**

[DEPRECATED] Excel spreadsheet based occupancy calculator is deprecated.

Occupancy calculator in now available in Nsight Compute

CUDA Occupancy Calculator

Just follow steps 1, 2, and 3 below! (or click here for help)

| | | |
|---|-------|--------|
| 1.) Select Compute Capability (click): | 8.6 | (Help) |
| 1.b) Select Shared Memory Size Config (bytes) | 65536 | |
| 1.c) Select CUDA version | 11,1 | |

| | | |
|--------------------------------------|------|--------|
| 2.) Enter your resource usage: | 256 | (Help) |
| Threads Per Block | 256 | |
| Registers Per Thread | 32 | |
| User Shared Memory Per Block (bytes) | 2048 | |

(Don't edit anything below this line)

| | | |
|---|------|--------|
| 3.) GPU Occupancy Data is displayed here and in the graphs: | 8.6 | (Help) |
| Active Threads per Multiprocessor | 1536 | |
| Active Warps per Multiprocessor | 48 | |
| Active Thread Blocks per Multiprocessor | 6 | |
| Occupancy of each Multiprocessor | 100% | |

| | | |
|--|------------------------------------|--|
| Physical Limits for GPU Compute Capability: | 8.6 | |
| Threads per Warp | 32 | |
| Max Warps per Multiprocessor | 48 | |
| Max Thread Blocks per Multiprocessor | 16 | |
| Max Threads per Multiprocessor | 1536 | |
| Maximum Thread Block Size | 1024 | |
| Registers per Multiprocessor | 65538 | |
| Max Registers per Thread Block | 65538 | |
| Max Registers per Thread | 255 | |
| Shared Memory per Multiprocessor (bytes) | 65538 | |
| Max Shared Memory per Block | 256 | |
| Register allocation unit size | warp | |
| Register allocation granularity | Shared Memory allocation unit size | |
| Shared Memory allocation unit size | 128 | |
| Warp allocation granularity | 4 | |
| Shared Memory (bytes) per Block (CUDA runtime use) | 1024 | |

| Allocated Resources | Per Block | Limit Per SM | Blocks Per SM | = Allocatable |
|---|-----------|--------------|---------------|---------------|
| Warp (Threads Per Block / Threads Per Warp) | 8 | 48 | 6 | |
| Registers (Warp limit per SM due to per-warp reg count) | 8 | 64 | 8 | |
| Shared Memory (Bytes) | 2048 | 65536 | 32 | |

Note: SM is an abbreviation for (Streaming) Multiprocessor

| Maximum Thread Blocks Per Multiprocessor | Blocks/SM | * Warps/Block = Warps/SM |
|---|-----------|--------------------------|
| Limited by Max Warps or Max Blocks per Multiprocessor | 6 | 8 |

SM - Streaming Multiprocessor





NVIDIA Nsight

Occupancy Calculator

| | | | |
|------------------------------------|------------|---------------------------------------|---|
| Compute Capability: | 8.6 | Threads Per Block: | <input type="range" value="256"/> 256 |
| Shared Memory Size Config (bytes): | 102400 | Registers Per Thread: | <input type="range" value="32"/> 32 |
| Global Load Cache Mode: | L1+L2 (ca) | User Shared Memory Per Block (bytes): | <input type="range" value="2048"/> 2048 |
| | | Block Barriers: | <input type="range" value="1"/> 1 |

Apply Automatically

The CUDA Occupancy Calculator allows you to compute the multiprocessor occupancy of a GPU by a given CUDA kernel. Occupancy is the ratio of active warps on an SM to the maximum number of active warps supported by the SM. The launch parameters can be configured with the input widget to the left. This tool also provides various important statistics such as resource allocation, occupancy limiters, GPU Data, etc in the form of tables and graphs.



NVIDIA Nsight

Occupancy Calculator

Compute Capability: **8.6** Threads Per Block: **256**

Shared Memory Size Config (bytes): **102400** Registers Per Thread: **32**

Global Load Cache Mode: **L1+L2 (ca)** User Shared Memory Per Block (bytes): **2048**

Block Barriers: **1**

Apply Automatically **Apply** **Reset**

The CUDA Occupancy Calculator allows you to compute the multiprocessor occupancy of a GPU by a given CUDA kernel. Occupancy is the ratio of active warps on an SM to the maximum number of active warps supported by the SM. The launch parameters can be configured with the input widget to the left. This tool also provides various important statistics such as resource allocation, occupancy limiters, GPU Data, etc in the form of tables and graphs.

RTX 3090 (Ampere) - Compute Capability 8.6



NVIDIA Nsight

Occupancy Calculator

Compute Capability: 8.6 Threads Per Block: 256

Shared Memory Size Config (bytes): 102400 Registers Per Thread: 32

Global Load Cache Mode: L1+L2 (ca) User Shared Memory Per Block (bytes): 2048

Block Barriers: 1

Apply Automatically

The CUDA Occupancy Calculator allows you to compute the multiprocessor occupancy of a GPU by a given CUDA kernel. Occupancy is the ratio of active warps on an SM to the maximum number of active warps supported by the SM. The launch parameters can be configured with the input widget to the left. This tool also provides various important statistics such as resource allocation, occupancy limiters, GPU Data, etc in the form of tables and graphs.

RTX 3090 (Ampere) - Compute Capability 8.6
WorkGroup size - 256



NVIDIA Nsight

Occupancy Calculator

Compute Capability: 8.6 Threads Per Block: 256

Shared Memory Size Config (bytes): 102400 Registers Per Thread: 32

Global Load Cache Mode: L1+L2 (ca) User Shared Memory Per Block (bytes): 2048

Block Barriers: 1

Apply Automatically

The CUDA Occupancy Calculator allows you to compute the multiprocessor occupancy of a GPU by a given CUDA kernel. Occupancy is the ratio of active warps on an SM to the maximum number of active warps supported by the SM. The launch parameters can be configured with the input widget to the left. This tool also provides various important statistics such as resource allocation, occupancy limiters, GPU Data, etc in the form of tables and graphs.

RTX 3090 (Ampere) - Compute Capability 8.6

WorkGroup size - 256

Registers Per Thread - 32



NVIDIA Nsight

Occupancy Calculator

Compute Capability: 8.6 Threads Per Block: 256

Shared Memory Size Config (bytes): 102400 Registers Per Thread: 32

Global Load Cache Mode: L1+L2 (ca) User Shared Memory Per Block (bytes): 2048

Block Barriers: 1

Apply Automatically

The CUDA Occupancy Calculator allows you to compute the multiprocessor occupancy of a GPU by a given CUDA kernel. Occupancy is the ratio of active warps on an SM to the maximum number of active warps supported by the SM. The launch parameters can be configured with the input widget to the left. This tool also provides various important statistics such as resource allocation, occupancy limiters, GPU Data, etc in the form of tables and graphs.

RTX 3090 (Ampere) - Compute Capability 8.6

WorkGroup size - 256

Registers Per Thread - 32

Local memory Per Thread - 8 bytes



Compute Capability: 8.6 Threads Per Block: 256

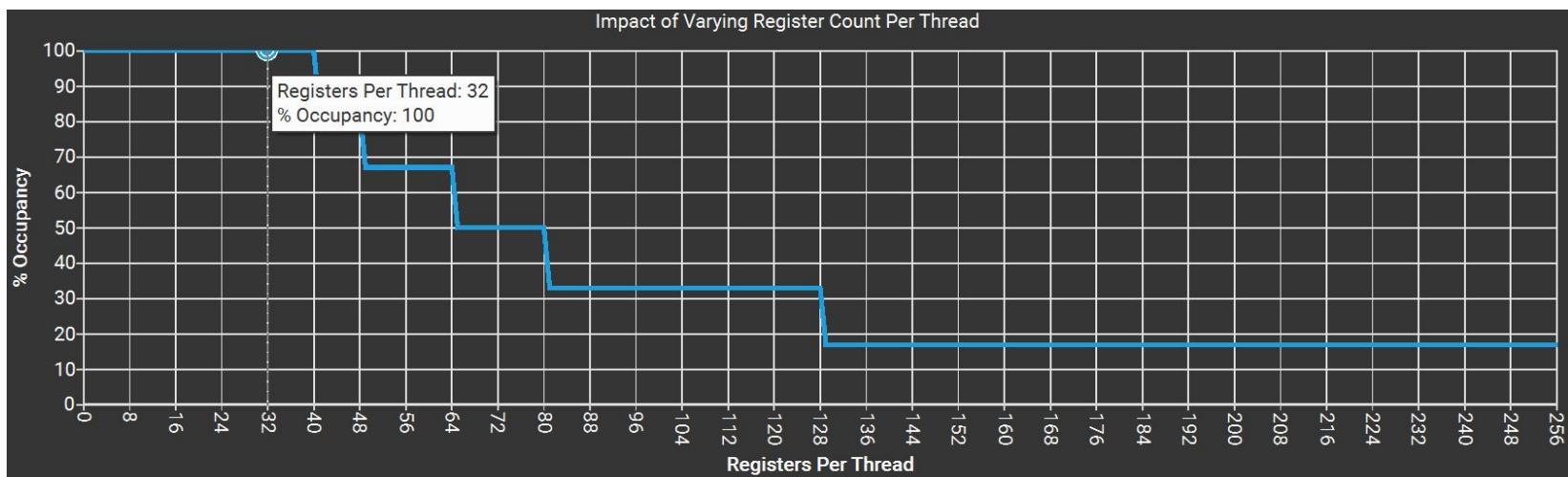
Shared Memory Size Config (bytes): 102400 Registers Per Thread: 32

Global Load Cache Mode: L1+L2 (ca) User Shared Memory Per Block (bytes): 2048

Block Barriers: 1

Apply Automatically

The CUDA Occupancy Calculator allows you to compute the multiprocessor occupancy of a GPU by a given CUDA kernel. Occupancy is the ratio of active warps on an SM to the maximum number of active warps supported by the SM. The launch parameters can be configured with the input widget to the left. This tool also provides various important statistics such as resource allocation, occupancy limiters, GPU Data, etc in the form of tables and graphs.





NVIDIA Nsight

Occupancy Calculator

Compute Capability: 8.6 Threads Per Block: 256

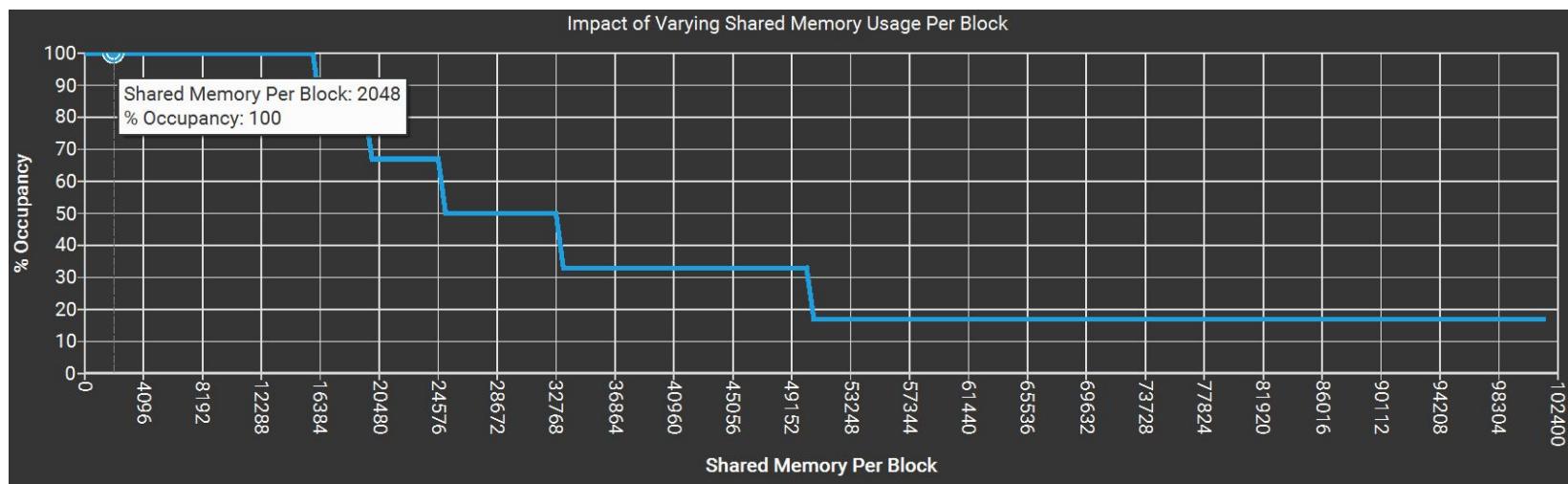
Shared Memory Size Config (bytes): 102400 Registers Per Thread: 32

Global Load Cache Mode: L1+L2 (ca) User Shared Memory Per Block (bytes): 2048

Block Barriers: 1

Apply Automatically

The CUDA Occupancy Calculator allows you to compute the multiprocessor occupancy of a GPU by a given CUDA kernel. Occupancy is the ratio of active warps on an SM to the maximum number of active warps supported by the SM. The launch parameters can be configured with the input widget to the left. This tool also provides various important statistics such as resource allocation, occupancy limiters, GPU Data, etc in the form of tables and graphs.



Глава 2: Умножение матриц

local memory

Все что мы видим на экране – это матрицы. Их умножение – это одна из основных операций в вычислительной математике. Важно, чтобы умножение было быстрое и эффективное. Для этого нужно использовать различные алгоритмы и методы. Одним из таких методов является использование локальной памяти.

Локальная память – это специальный тип памяти, который используется для хранения данных, которые часто используются в процессе вычислений. Использование локальной памяти позволяет уменьшить время доступа к данным и улучшить производительность вычислений. Для этого нужно использовать различные алгоритмы и методы. Одним из таких методов является использование локальной памяти.

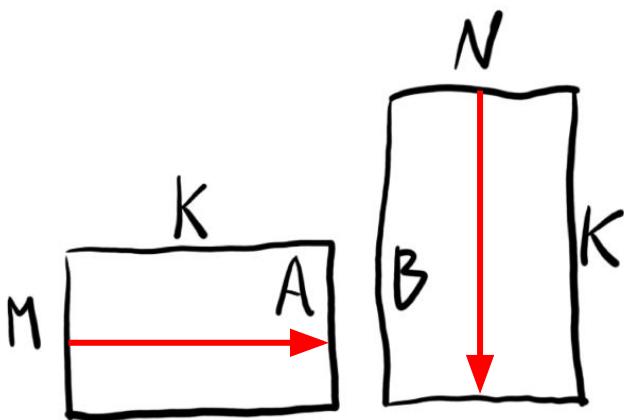
Локальная память – это специальный тип памяти, который используется для хранения данных, которые часто используются в процессе вычислений. Использование локальной памяти позволяет уменьшить время доступа к данным и улучшить производительность вычислений. Для этого нужно использовать различные алгоритмы и методы. Одним из таких методов является использование локальной памяти.

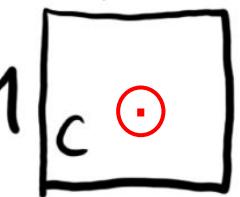
Локальная память – это специальный тип памяти, который используется для хранения данных, которые часто используются в процессе вычислений. Использование локальной памяти позволяет уменьшить время доступа к данным и улучшить производительность вычислений. Для этого нужно использовать различные алгоритмы и методы. Одним из таких методов является использование локальной памяти.

Локальная память – это специальный тип памяти, который используется для хранения данных, которые часто используются в процессе вычислений. Использование локальной памяти позволяет уменьшить время доступа к данным и улучшить производительность вычислений. Для этого нужно использовать различные алгоритмы и методы. Одним из таких методов является использование локальной памяти.

Локальная память – это специальный тип памяти, который используется для хранения данных, которые часто используются в процессе вычислений. Использование локальной памяти позволяет уменьшить время доступа к данным и улучшить производительность вычислений. Для этого нужно использовать различные алгоритмы и методы. Одним из таких методов является использование локальной памяти.

Умножение матриц



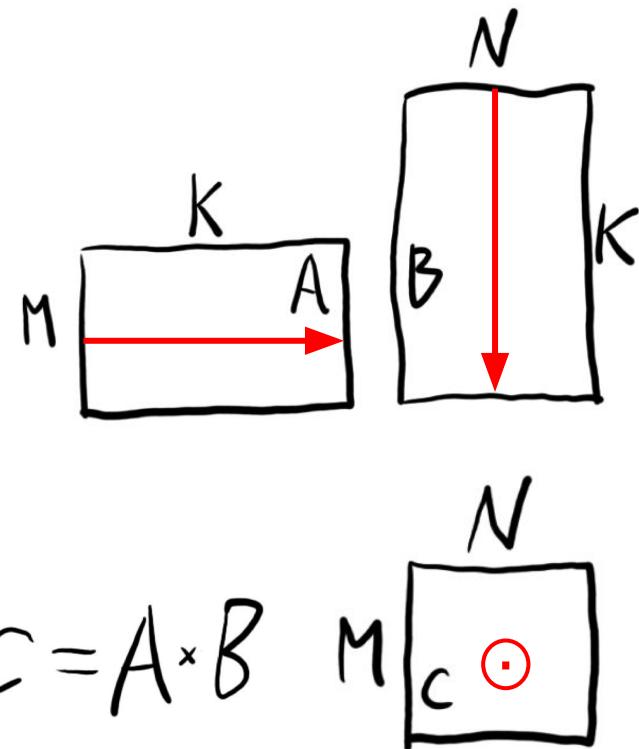
$$C = A \times B \quad M \quad N$$


Below the multiplication equation, there is a square box labeled 'C' in the bottom-left corner. In the top-right corner of this box, there is a red circle containing a white dot, representing the element at the first row and first column of matrix C.

Умножение матриц

Как выглядит:

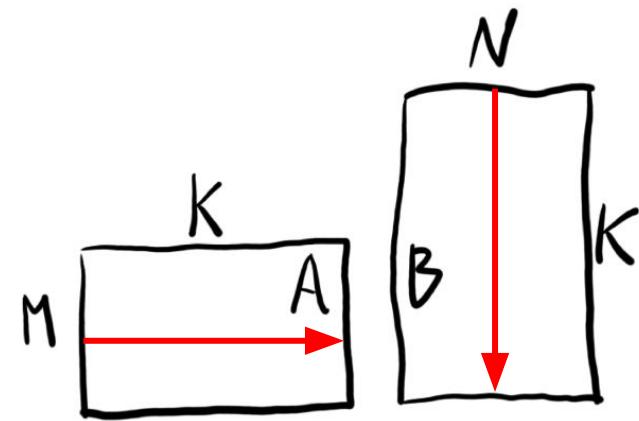
- WorkRange?
- Задача WorkItem?



Умножение матриц

Как выглядит:

- WorkRange?
- Задача WorkItem?
- WorkGroup?



$$C = A \times B \quad M \quad N$$

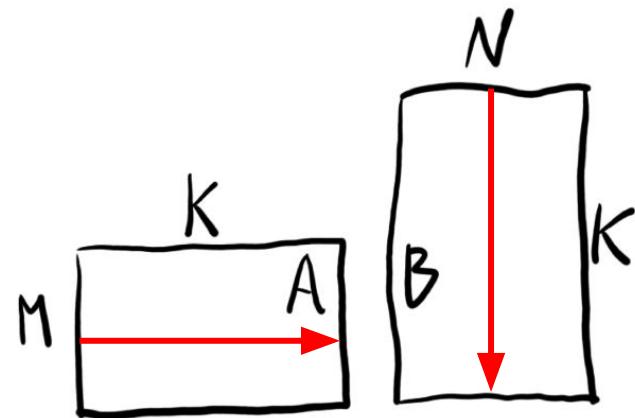
C

Умножение матриц

Сколько у нас вычислений?

Сколько у нас чтений/записей данных?

Какая пропорция?



$$C = A \times B \quad M \quad N$$

C

Умножение матриц

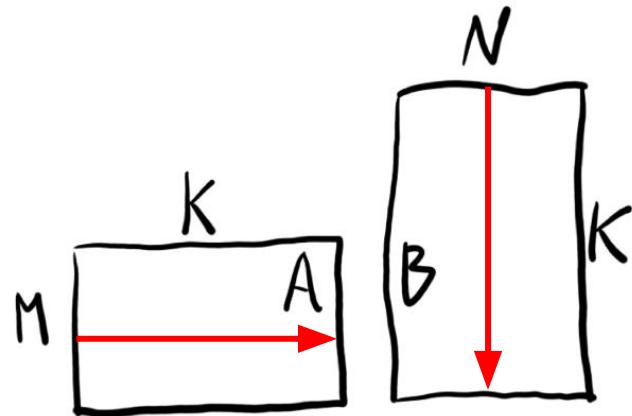
Сколько у нас вычислений?

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

Сколько у нас чтений/записей данных?

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

Какая пропорция?



$$C = A \times B \quad M \quad N$$

Умножение матриц

Сколько у нас вычислений?

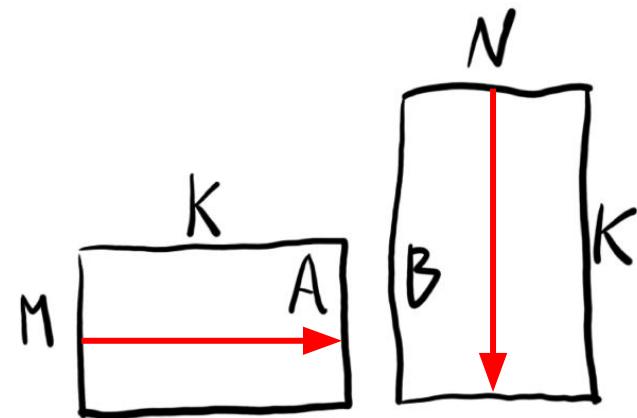
$$O(N \cdot M \cdot K)$$

Сколько у нас чтений/записей данных?

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

Какая пропорция?

1:1 Это хороший результат?



$$C = A \times B \quad M \quad N$$

Умножение матриц

Сколько у нас вычислений?

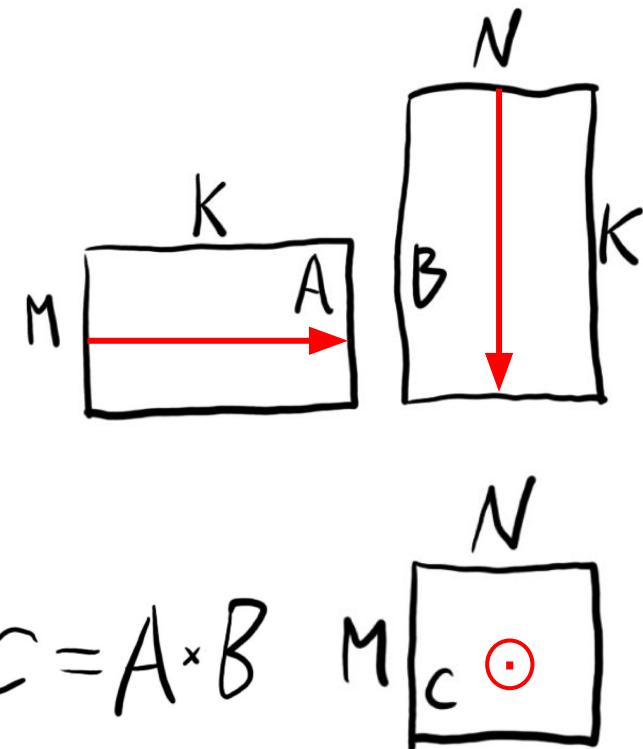
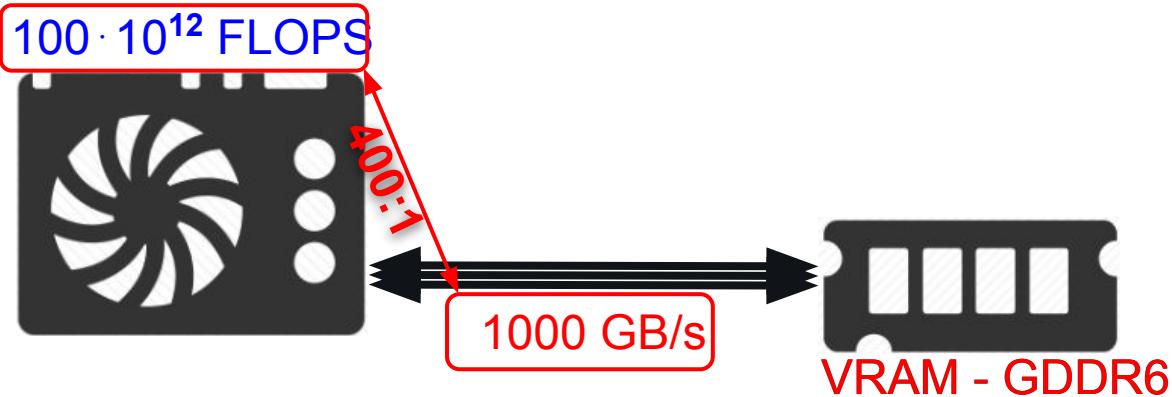
$$O(N \cdot M \cdot K)$$

Сколько у нас чтений/записей данных?

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

Какая пропорция?

1:1 Это хороший результат?



Умножение матриц

Сколько у нас вычислений?

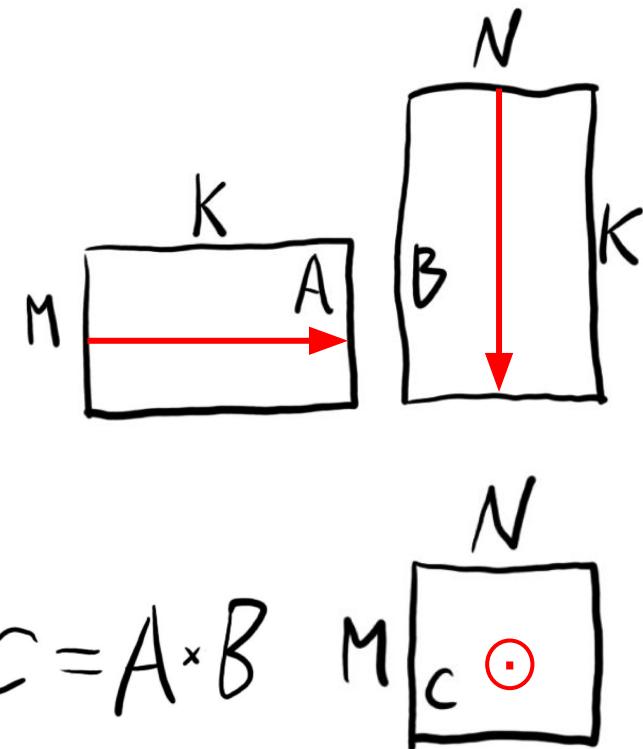
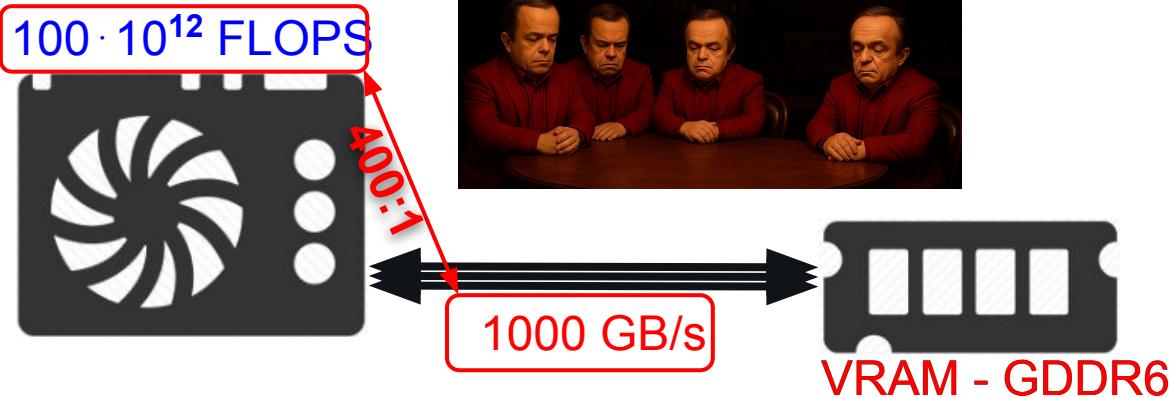
$$O(N \cdot M \cdot K)$$

Сколько у нас чтений/записей данных?

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

Какая пропорция?

1:1 **Memory-bound!**



$$C = A \times B$$

Умножение матриц

Сколько у нас вычислений?

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

Сколько у нас чтений/записей данных?

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

Какая пропорция?

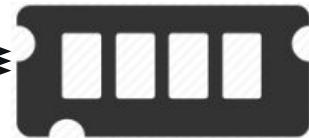
1:1 **Memory-bound!**

$100 \cdot 10^{12}$ FLOPS

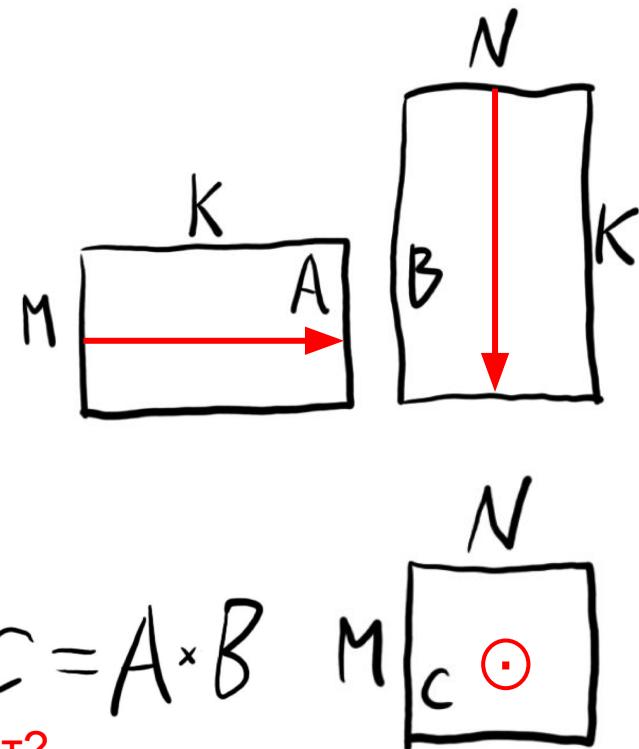


400:1

1000 GB/s



VRAM - GDDR6



Как увеличить объем вычислений на считанный байт?

Умножение матриц

Сколько у нас вычислений?

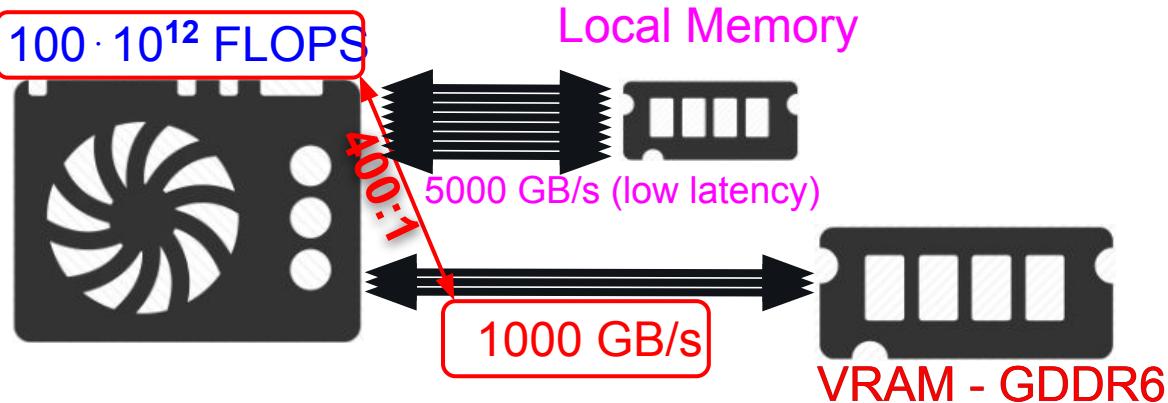
$$O(N \cdot M \cdot K)$$

Сколько у нас чтений/записей данных?

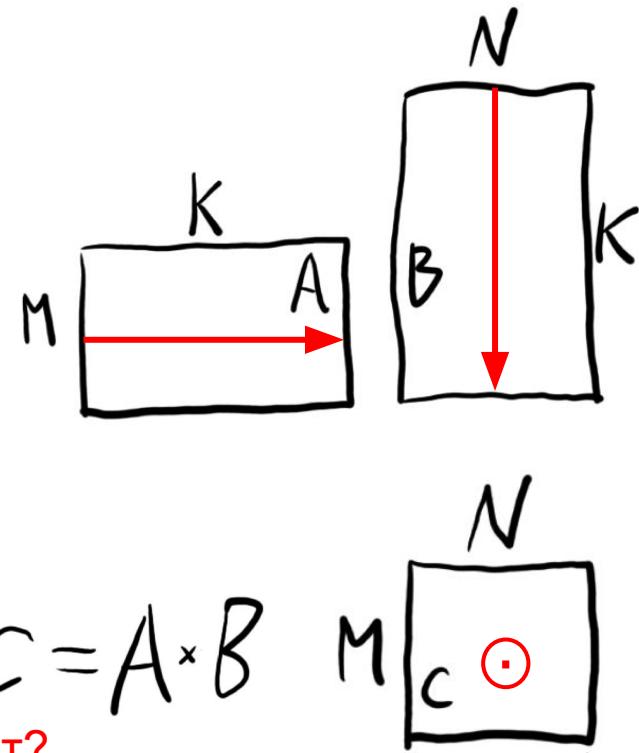
$$O(N \cdot M \cdot K)$$

Какая пропорция?

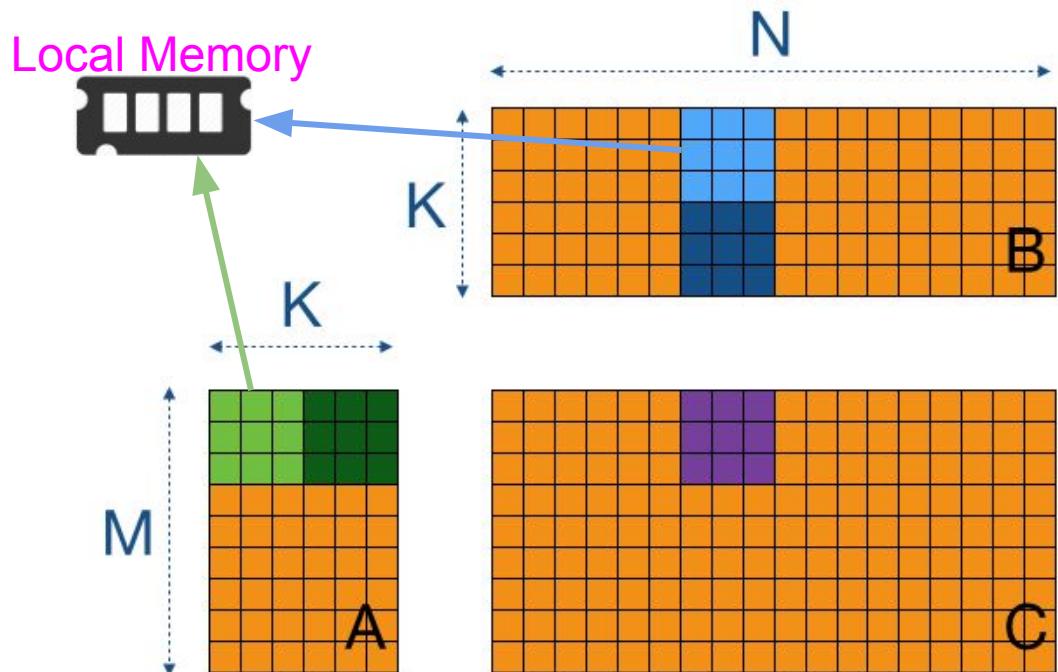
1:1 **Memory-bound!**



Как увеличить объем вычислений на считанный байт?

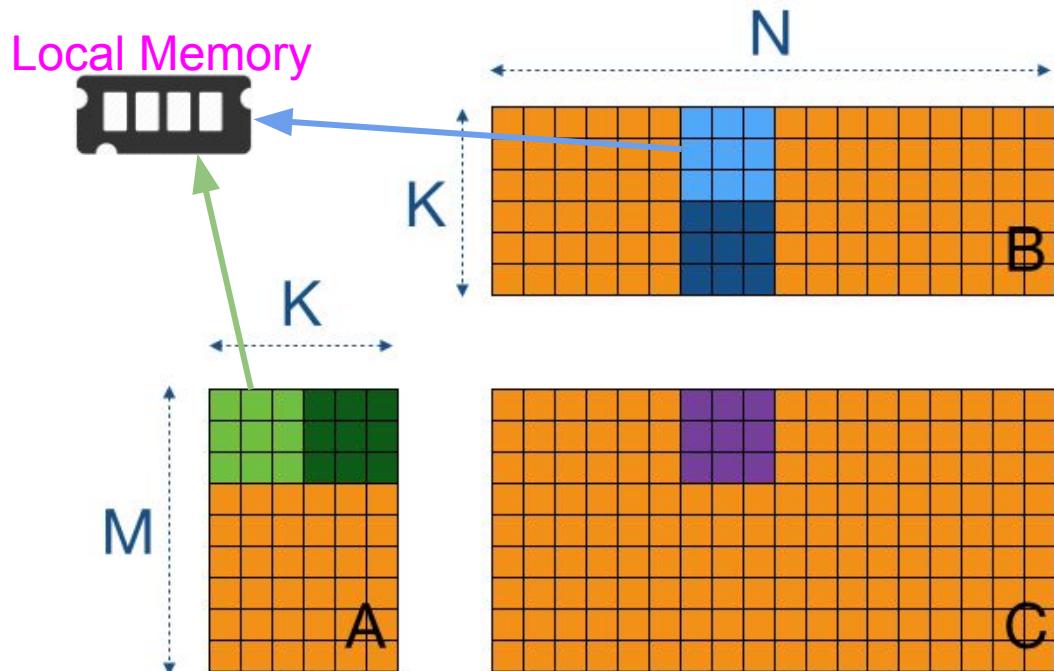


Умножение матриц



Умножение матриц

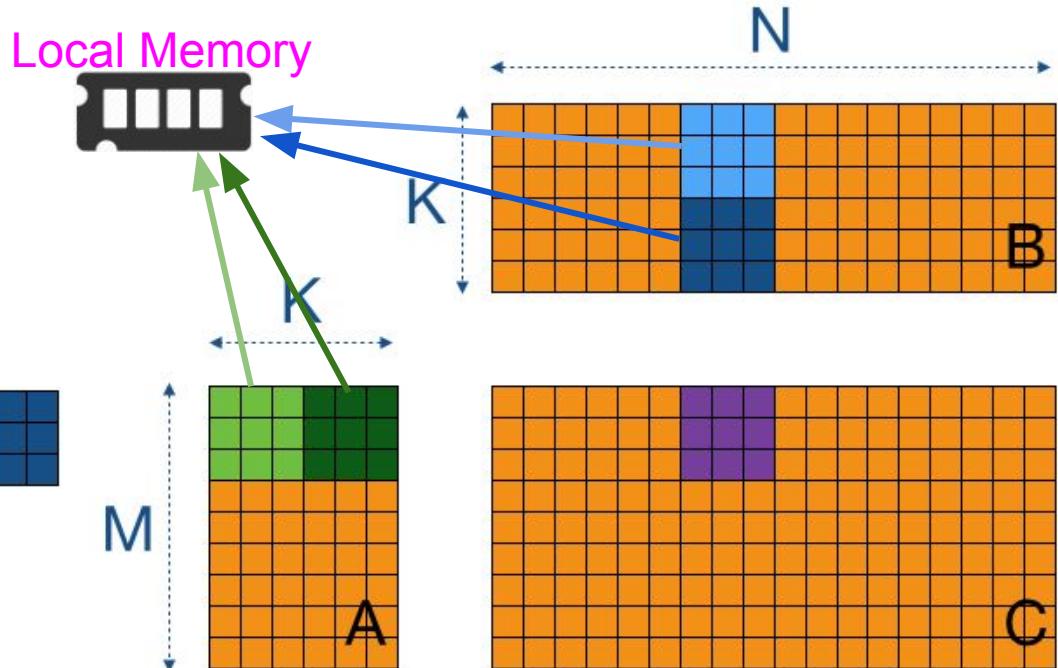
$$\begin{bmatrix} \text{purple} \\ \text{green} \\ \text{blue} \end{bmatrix} = \begin{bmatrix} \text{green} \\ \text{blue} \end{bmatrix} \times \begin{bmatrix} \text{blue} \end{bmatrix} + \dots$$



Умножение матриц

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

За счет чего это ускоряет?



Умножение матриц

Сколько у нас вычислений?

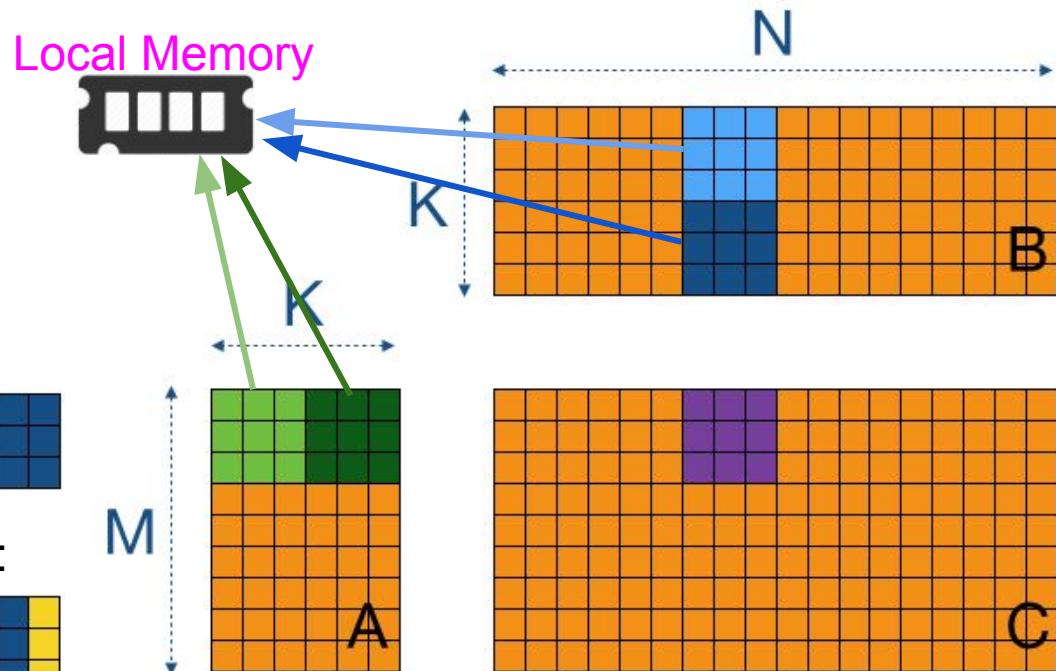
Сколько у нас чтений/записей?

Какая пропорция?

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

Переиспользование данных:

$$\left\{ \begin{bmatrix} \text{purple} & \text{yellow} \\ \text{purple} & \text{yellow} \end{bmatrix} \right\}_{32} = \begin{bmatrix} \text{green} & \text{yellow} \\ \text{green} & \text{yellow} \end{bmatrix} \times \begin{bmatrix} \text{blue} & \text{yellow} \\ \text{blue} & \text{yellow} \end{bmatrix} + \begin{bmatrix} \text{dark green} & \text{yellow} \\ \text{dark green} & \text{yellow} \end{bmatrix} \times \begin{bmatrix} \text{dark blue} & \text{yellow} \\ \text{dark blue} & \text{yellow} \end{bmatrix}$$



Умножение матриц

Сколько у нас вычислений?

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

Сколько у нас чтений/записей?

$$O(N \cdot M \cdot K / 32)$$

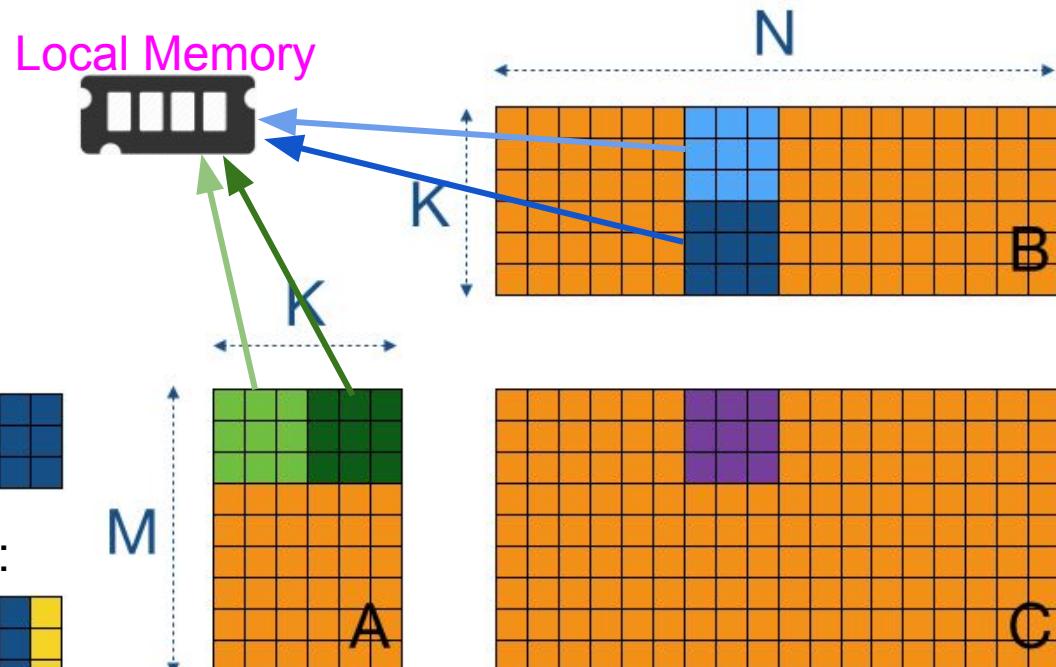
Какая пропорция?

32:1

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

Переиспользование данных:

$$\left\{ \begin{matrix} \text{purple} \\ \text{matrix} \end{matrix} \right. = \left. \begin{matrix} \text{yellow} \\ \text{matrix} \end{matrix} \right. \times \left. \begin{matrix} \text{blue} \\ \text{matrix} \end{matrix} \right. + \left. \begin{matrix} \text{dark green} \\ \text{matrix} \end{matrix} \right. \times \left. \begin{matrix} \text{dark blue} \\ \text{matrix} \end{matrix} \right.$$



Умножение матриц

Сколько у нас вычислений?

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

Сколько у нас чтений/записей?

$$O(N \cdot M \cdot K / 32)$$

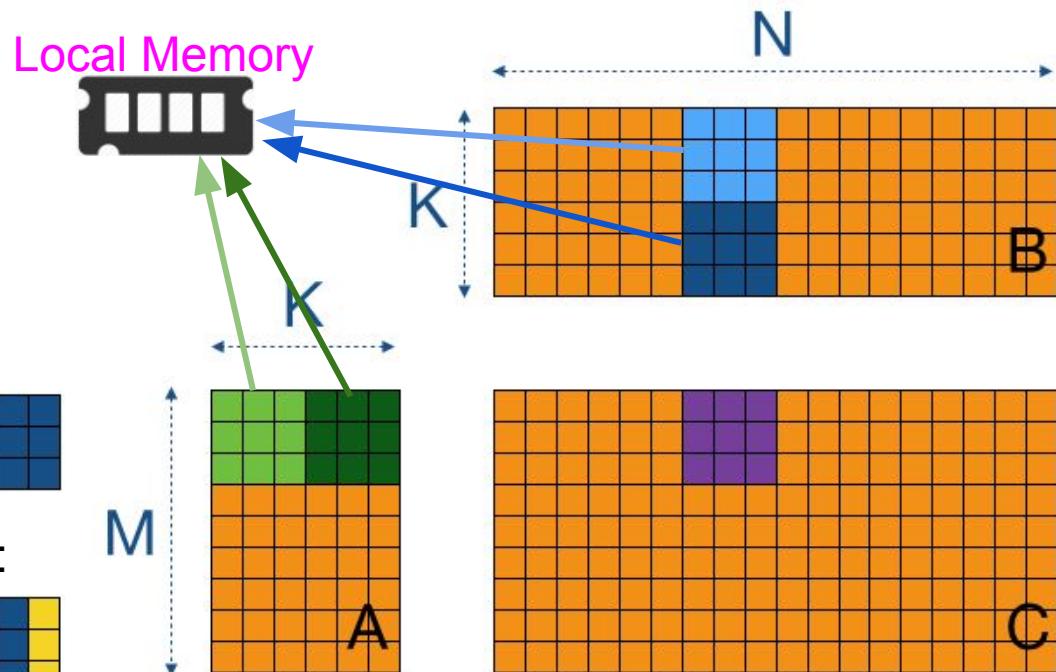
Какая пропорция?

32:1

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

Переиспользование данных:

$$\begin{matrix} 32 \\ \text{matrix} \end{matrix} = \begin{matrix} \text{yellow} \\ \text{matrix} \end{matrix} \times \begin{matrix} \text{blue} \\ \text{matrix} \end{matrix} + \begin{matrix} \text{dark green} \\ \text{matrix} \end{matrix} \times \begin{matrix} \text{yellow} \\ \text{matrix} \end{matrix}$$



Как пойти еще дальше? Как еще увеличить пропорцию?

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>



Перерыв!

Streaming
Pübiprocessor

Лицензия

© VIDIA

Глава 3: Умножение матриц

Tensor Cores, WMMA



Умножение матриц

Tensor Cores

| | VRAM TB/s | FP 32 TFlops | FP 16 TFlops | FP 16 (tensor) TFlops |
|-------------------|----------------------|-------------------------|-------------------------|----------------------------------|
| 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

$$D = \begin{pmatrix} A_{0,0} & A_{0,1} & A_{0,2} & A_{0,3} \\ A_{1,0} & A_{1,1} & A_{1,2} & A_{1,3} \\ A_{2,0} & A_{2,1} & A_{2,2} & A_{2,3} \\ A_{3,0} & A_{3,1} & A_{3,2} & A_{3,3} \end{pmatrix}_{\text{FP16 or FP32}} \times \begin{pmatrix} B_{0,0} & B_{0,1} & B_{0,2} & B_{0,3} \\ B_{1,0} & B_{1,1} & B_{1,2} & B_{1,3} \\ B_{2,0} & B_{2,1} & B_{2,2} & B_{2,3} \\ B_{3,0} & B_{3,1} & B_{3,2} & B_{3,3} \end{pmatrix}_{\text{FP16}} + \begin{pmatrix} C_{0,0} & C_{0,1} & C_{0,2} & C_{0,3} \\ C_{1,0} & C_{1,1} & C_{1,2} & C_{1,3} \\ C_{2,0} & C_{2,1} & C_{2,2} & C_{2,3} \\ C_{3,0} & C_{3,1} & C_{3,2} & C_{3,3} \end{pmatrix}_{\text{FP16 or FP32}}$$

4x4



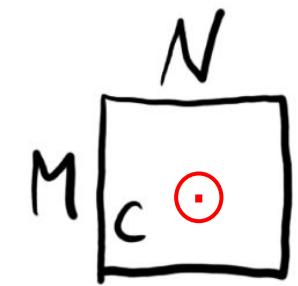
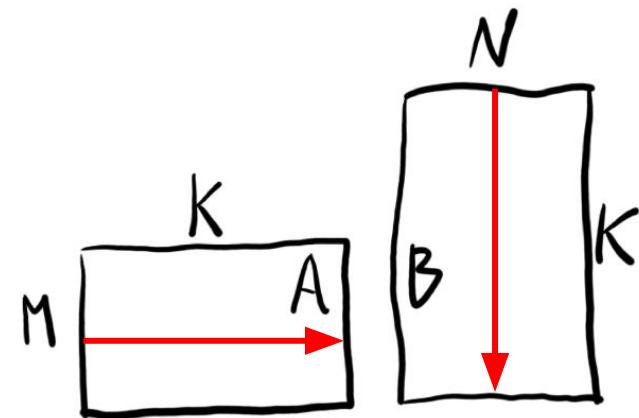
Tensor Cores (CUDA kernel)

$$D = \begin{pmatrix} A_{0,0} & A_{0,1} & A_{0,\dots} & A_{0,15} \\ A_{1,0} & A_{1,1} & A_{1,\dots} & A_{1,15} \\ A_{\dots,0} & A_{\dots,1} & A_{\dots,\dots} & A_{\dots,15} \\ A_{15,0} & A_{15,1} & A_{15,\dots} & A_{15,15} \end{pmatrix}_{\text{FP16 or FP32}} \times \begin{pmatrix} B_{0,0} & B_{0,1} & B_{0,\dots} & B_{0,15} \\ B_{1,0} & B_{1,1} & B_{1,\dots} & B_{1,15} \\ B_{\dots,0} & B_{\dots,1} & B_{\dots,\dots} & B_{\dots,15} \\ B_{15,0} & B_{15,1} & B_{15,\dots} & B_{15,15} \end{pmatrix}_{\text{FP16}} + \begin{pmatrix} C_{0,0} & C_{0,1} & C_{0,\dots} & C_{0,15} \\ C_{1,0} & C_{1,1} & C_{1,\dots} & C_{1,15} \\ C_{\dots,0} & C_{\dots,1} & C_{\dots,\dots} & C_{\dots,15} \\ C_{15,0} & C_{15,1} & C_{15,\dots} & C_{15,15} \end{pmatrix}_{\text{FP16 or FP32}}$$

16x16

```
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) {
```

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



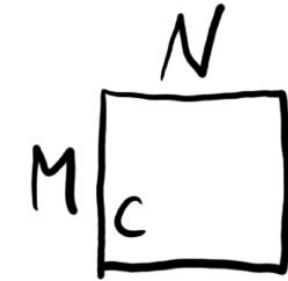
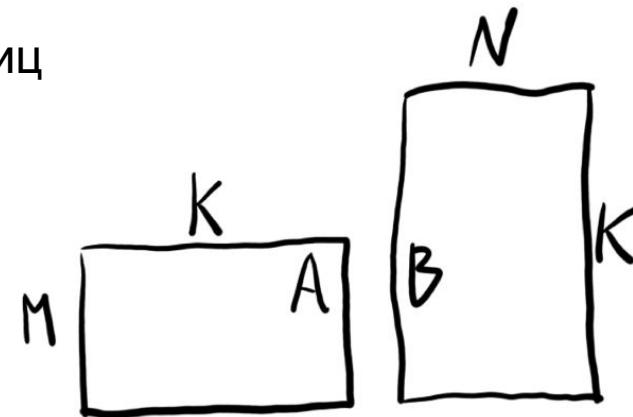
```

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

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



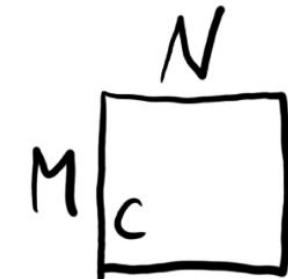
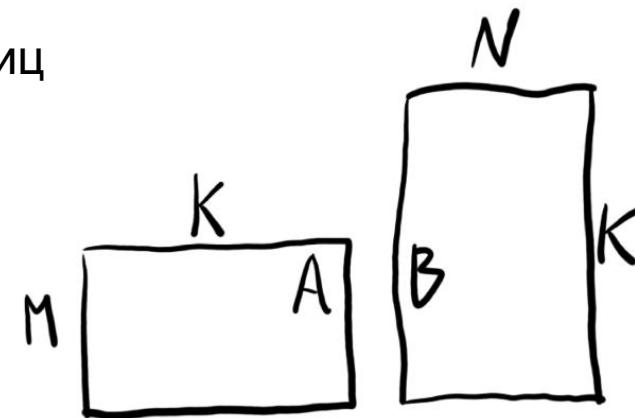
```

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

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



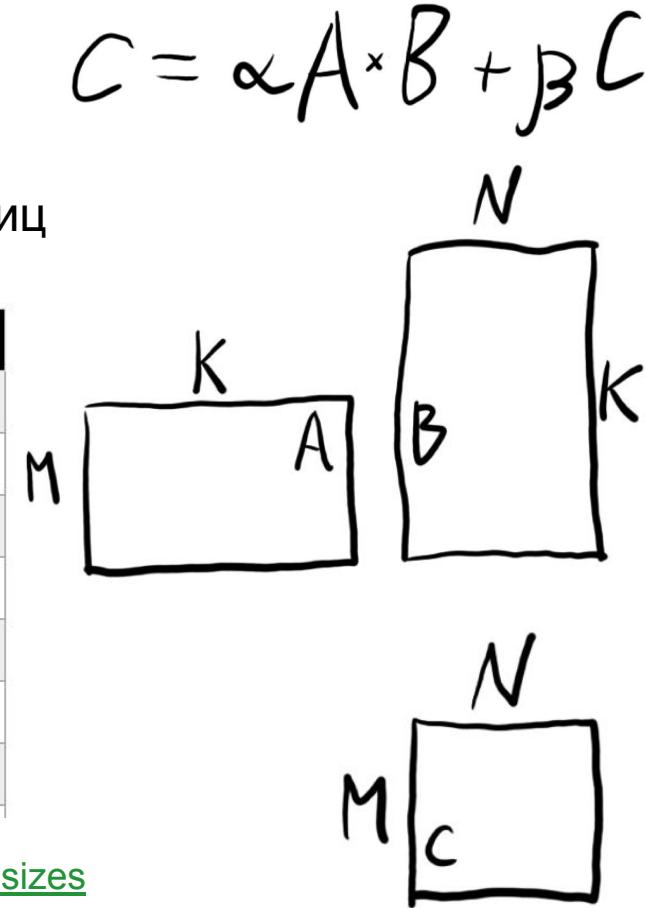
```

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

| Matrix A | Matrix B | Accumulator | Matrix Size (m-n-k) |
|---------------|---------------|-------------|---------------------|
| <u>half</u> | <u>half</u> | float | 16x16x16 |
| half | half | float | 32x8x16 |
| half | half | float | 8x32x16 |
| half | half | <u>half</u> | 16x16x16 |
| half | half | <u>half</u> | 32x8x16 |
| half | half | <u>half</u> | 8x32x16 |
| unsigned char | unsigned char | int | 16x16x16 |
| ... | ... | ... | ... |



<https://docs.nvidia.com/cuda/cuda-c-programming-guide/#wmma-type-sizes>

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

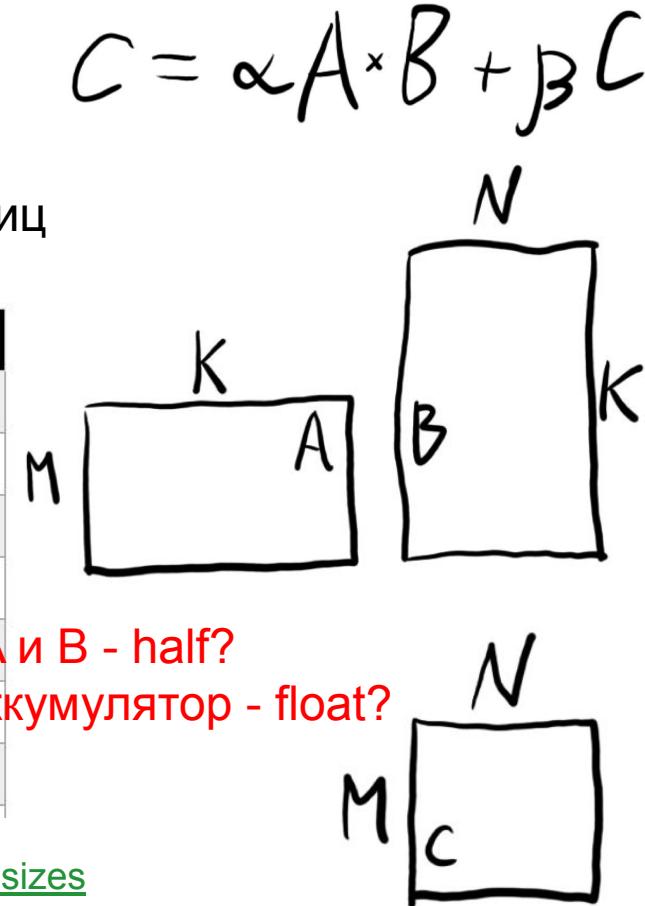
```

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

| Matrix A | Matrix B | Accumulator | Matrix Size (m-n-k) |
|---------------|---------------|-------------|---------------------|
| __half | __half | float | 16x16x16 |
| __half | __half | float | 32x8x16 |
| __half | __half | float | 8x32x16 |
| __half | __half | __half | 16x16x16 |
| __half | __half | __half | 32x8x16 |
| __half | __half | __half | 8x32x16 |
| unsigned char | unsigned char | int | 16x16x16 |
| ... | ... | ... | ... |



<https://docs.nvidia.com/cuda/cuda-c-programming-guide/#wmma-type-sizes>

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

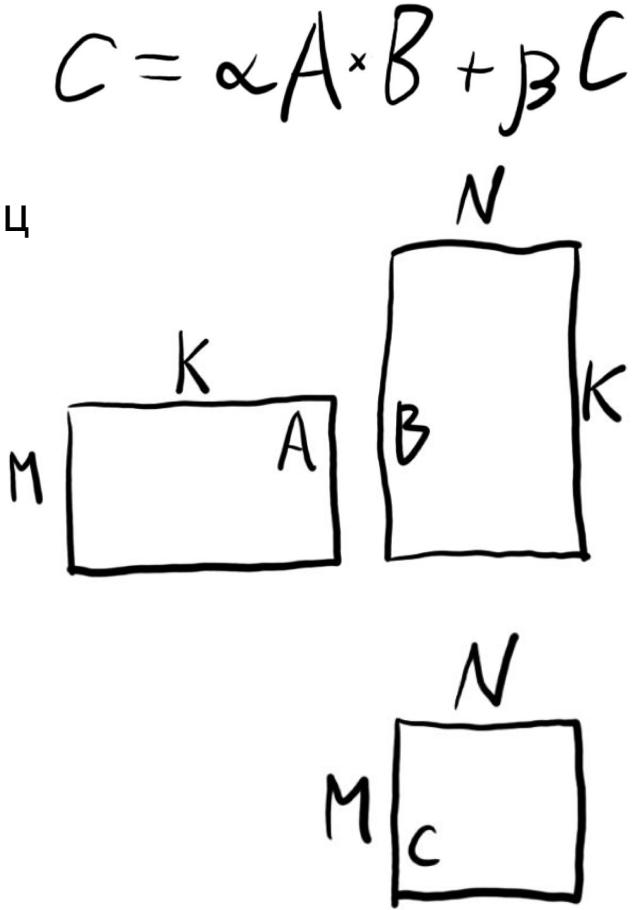
```

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

Почему А и В - half?
 Но С и аккумулятор - float?
 Подсказка:
 - А и В перемножаются
 - аккумулятор суммируется



```

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;

```

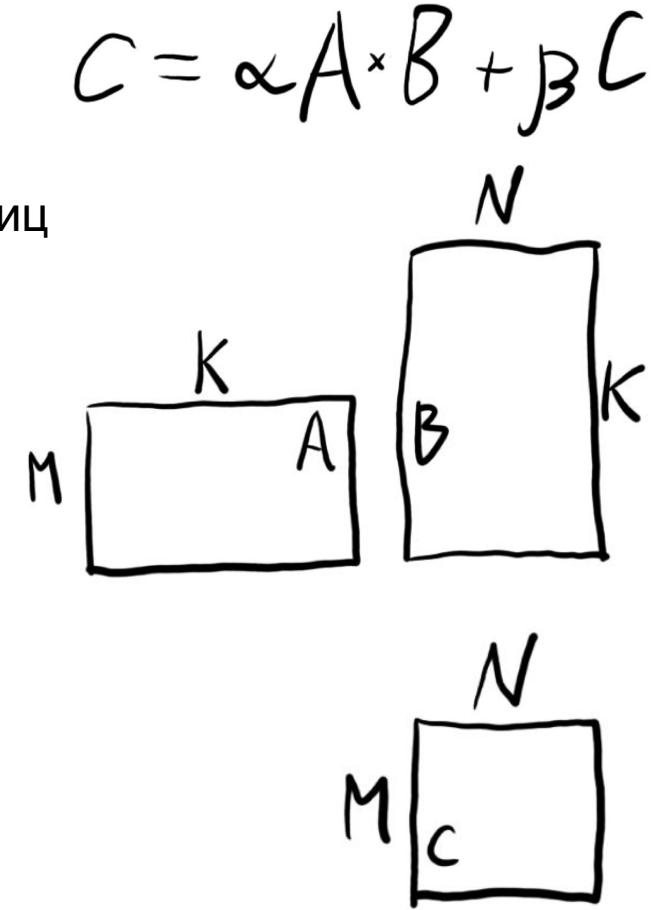
Общие на весь warp 16x16 фрагменты матриц
WMMA = Warp Matrix Multiply-Accumulate

Почему А и В - half?
 Но С и аккумулятор - float?

Что будет если
сложить огромное
 и малое число?



$$+ (1 + \text{mantissa}) \times 2^{(\text{exponent} - 127)}$$



```

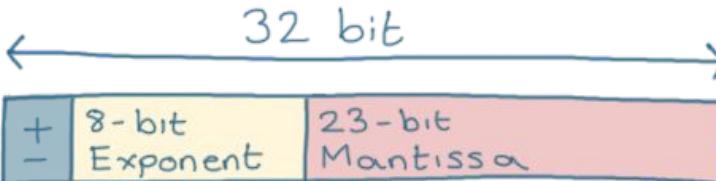
69 // Performs an MxNxK GEMM (C=alpha*A*B + beta*C)
75 __global__ void wmma_example(half *a, half *b, float *c, int M, int N, int K, float alpha, float beta) {
85     // Declare the fragments
86     wmma::fragment<wmma::matrix_a, WMMA_M, WMMA_N, WMMA_K, half, wmma::col_major> a_frag;
87     wmma::fragment<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

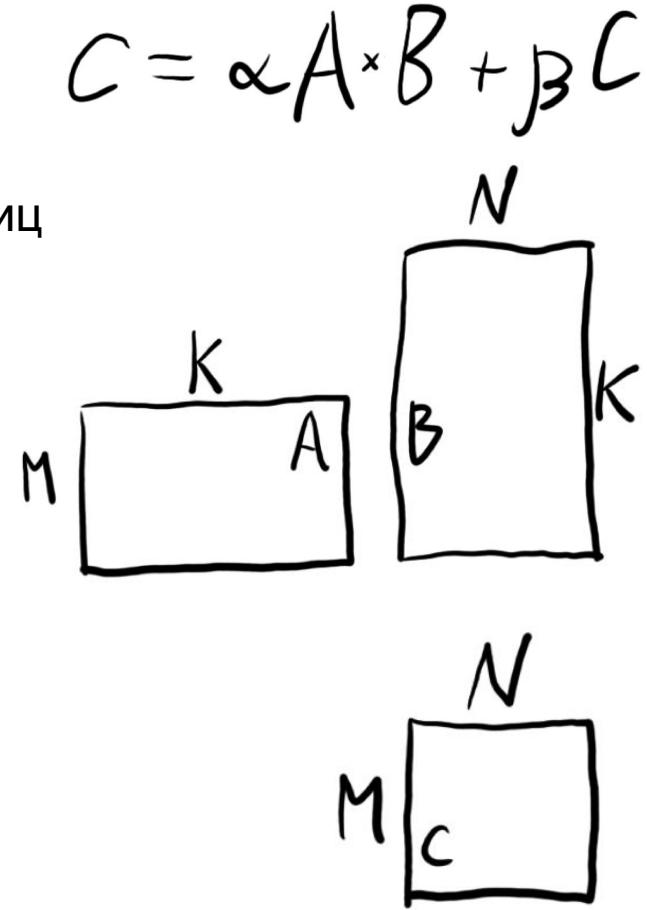
Почему А и В - half?
 Но С и аккумулятор - float?

Что будет если
сложить огромное
 и малое число?



Что будет если
умножить огромное
 и малое число?

$$\pm (1 + \text{mantissa}) \times 2^{(\text{exponent} - 127)}$$



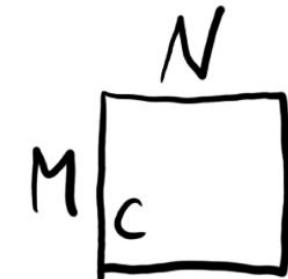
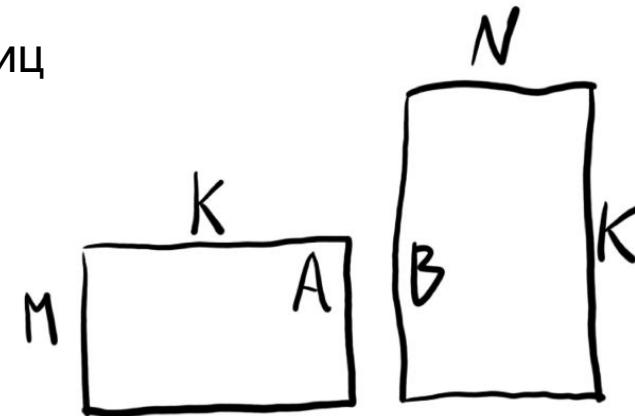
```

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

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



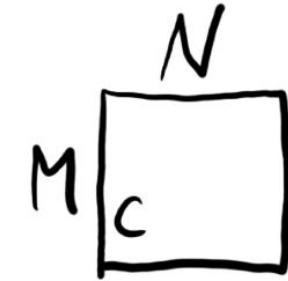
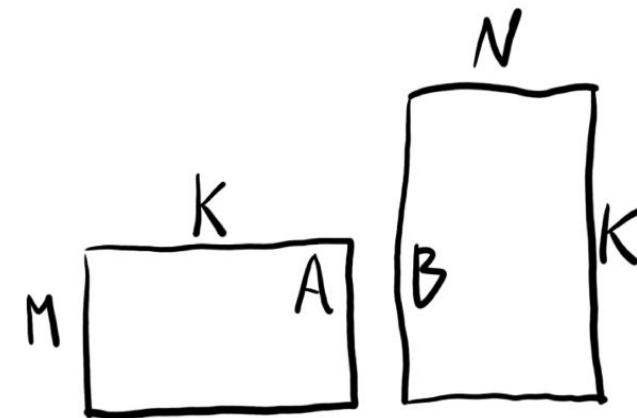
```

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

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



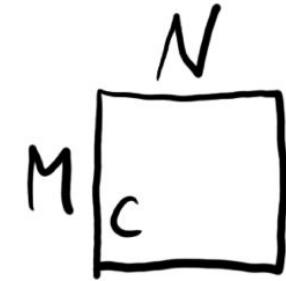
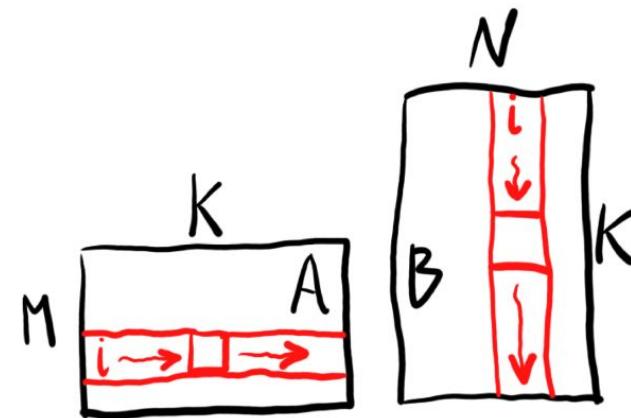
```

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) {

```

16x16
acc_frag

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

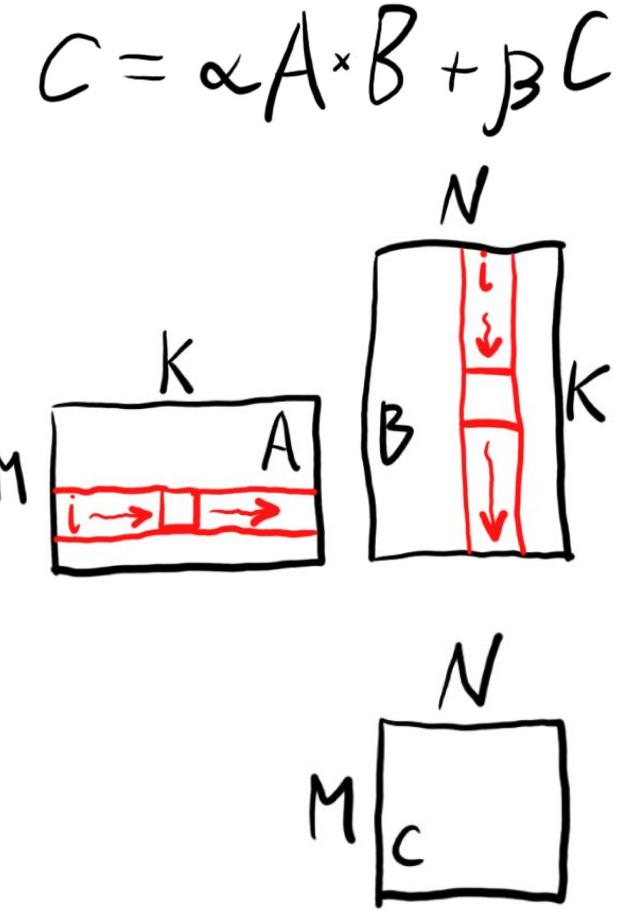


```

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) { acc_frag
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



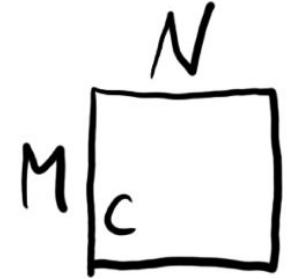
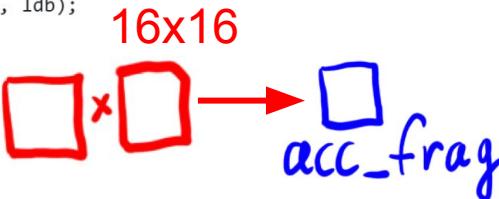
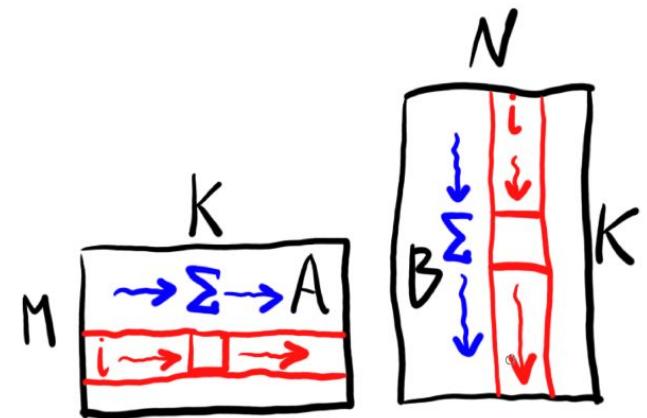
```

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) { acc_frag
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);
106
107            // Perform the matrix multiplication
108            wmma::mma_sync(acc_frag, a_frag, b_frag, acc_frag);
109        }
110    }
111 }

```

16x16

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



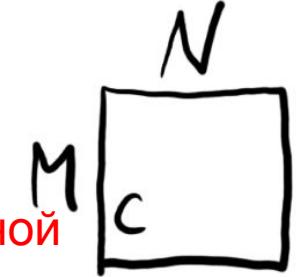
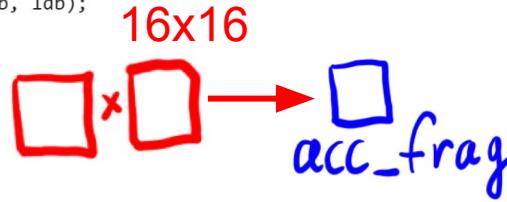
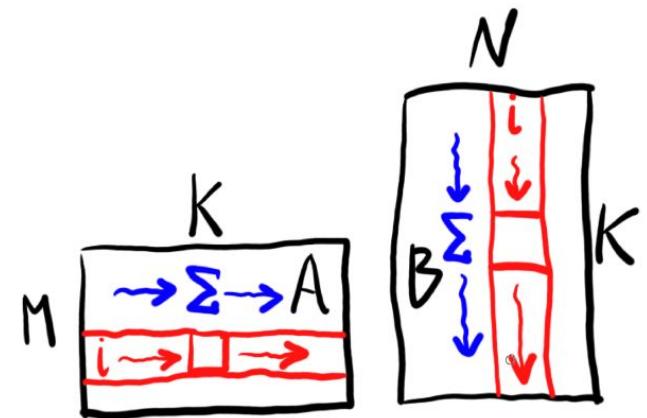
```

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) { acc_frag
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);
106
107            // Perform the matrix multiplication
108            wmma::mma_sync(acc_frag, a_frag, b_frag, acc_frag);
109        }
110    }

```

16x16
acc_frag

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



Чем этот код отличается от классического умножения в локальной памяти?

```

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) { acc_frag
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);
106
107            // Perform the matrix multiplication
108            wmma::mma_sync(acc_frag, a_frag, b_frag, acc_frag);
109        }
110    }

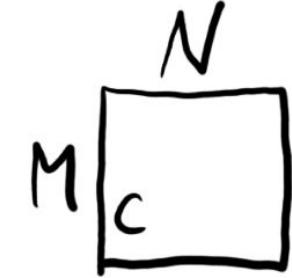
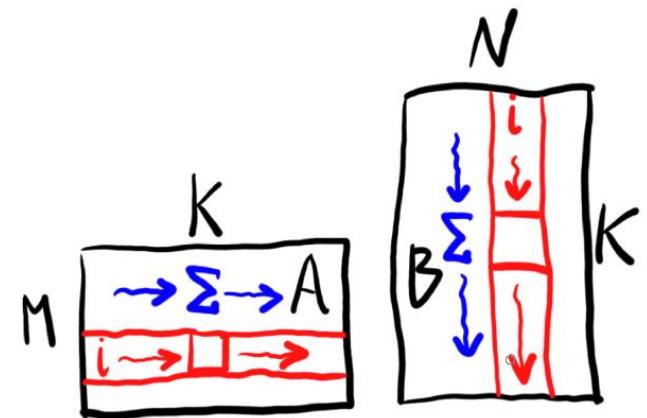
```

Что нам осталось сделать?

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

16x16
acc_frag

16x16
 $A \times B \rightarrow acc_frag$

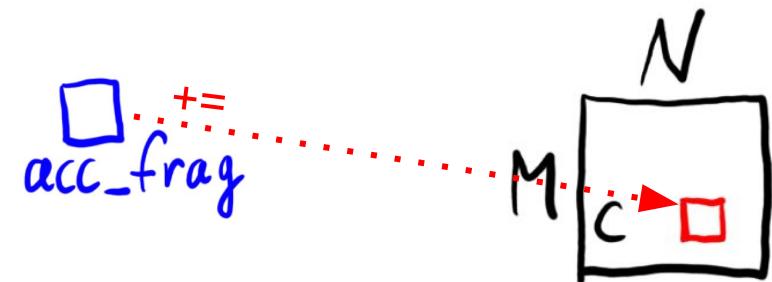
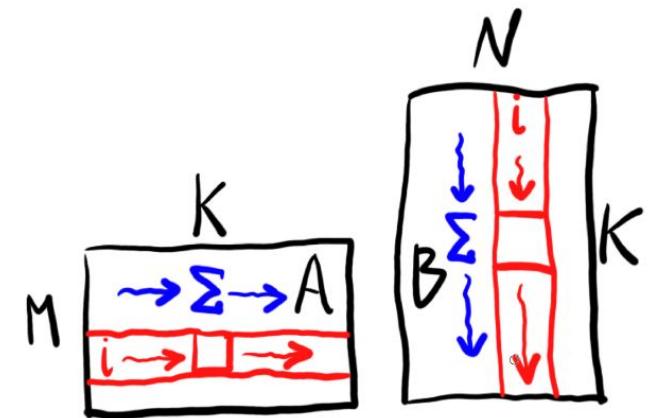


```

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;
90
91     .....

```

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



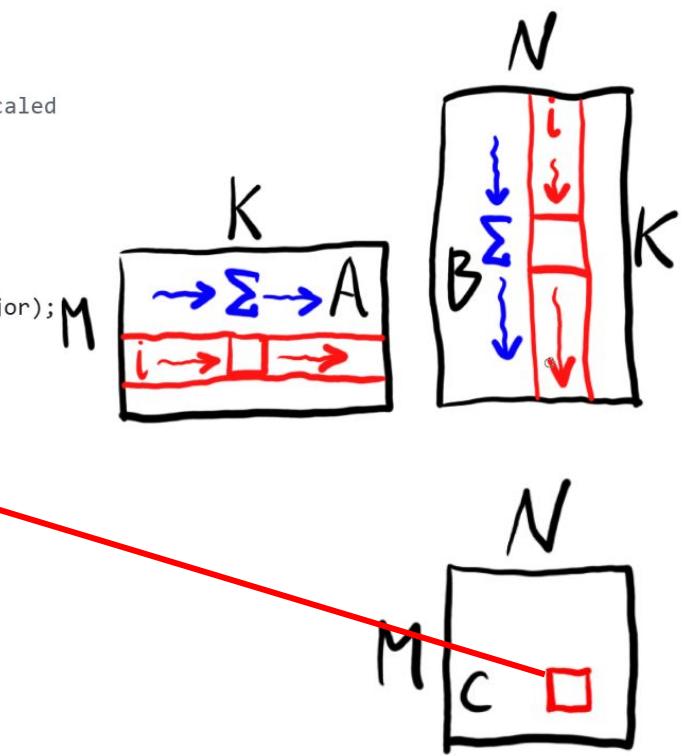
```

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;

.....
113     // Load in the current value of c, scale it by beta, and add this our result scaled
114     int cRow = warpM * WMMA_M;
115     int cCol = warpN * WMMA_N;

117     if (cRow < M && cCol < N) {
118         wmma::load_matrix_sync(c_frag, c + cRow + cCol * ldc, ldc, wmma::mem_col_major); ←
```

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

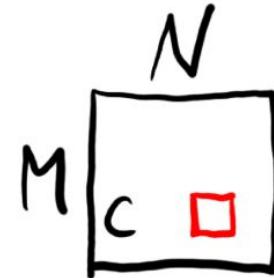
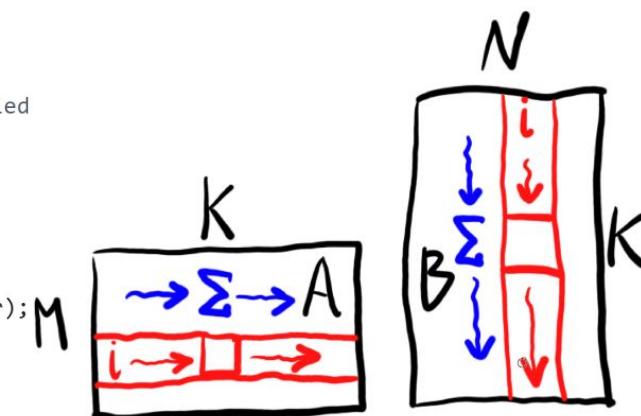


```

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;
90
91     .....
92
113     // Load in the current value of c, scale it by beta, and add this our result scaled
114     int cRow = warpM * WMMA_M;
115     int cCol = warpN * WMMA_N;
116
117     if (cRow < M && cCol < N) {
118         wmma::load_matrix_sync(c_frag, c + cRow + cCol * ldc, ldc, wmma::mem_col_major);
119
120 #pragma unroll
121     for(int i=0; i < c_frag.num_elements; i++) {
122         c_frag.x[i] = alpha * acc_frag.x[i] + beta * c_frag.x[i];
123     }
}

```

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

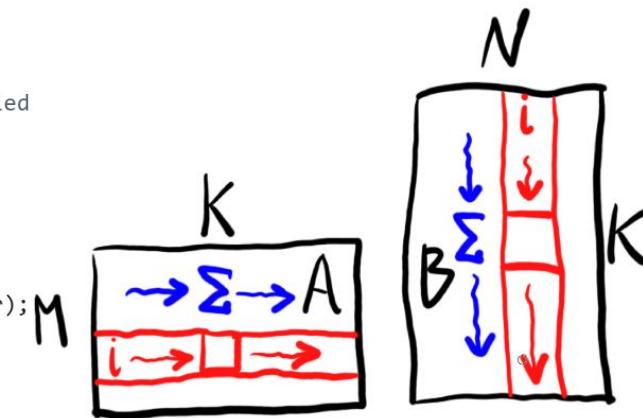


```

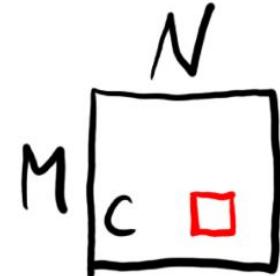
69 // Performs an MxNxK GEMM (C=alpha*A*B + beta*C)
75 _global_ void wmma_example(half *a, half *b, float *c, int M, int N, int K, float alpha, float beta) {
85     // Declare the fragments
86     wmma::fragment<wmma::matrix_a, WMMA_M, WMMA_N, WMMA_K, half, wmma::col_major> a_frag;
87     wmma::fragment<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;
90
91     .....
92
113     // Load in the current value of c, scale it by beta, and add this our result scaled
114     int cRow = warpM * WMMA_M;
115     int cCol = warpN * WMMA_N;
116
117     if (cRow < M && cCol < N) {
118         wmma::load_matrix_sync(c_frag, c + cRow + cCol * ldc, ldc, wmma::mem_col_major);
119
120 #pragma unroll
121     for(int i=0; i < c_frag.num_elements; i++) {
122         c_frag.x[i] = alpha * acc_frag.x[i] + beta * c_frag.x[i];
123     }
}

```

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



А не будет тормозить?
Почему не спец. функция в железе?

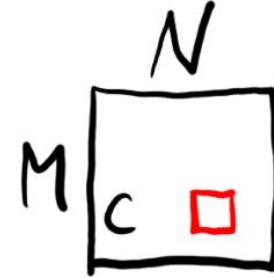
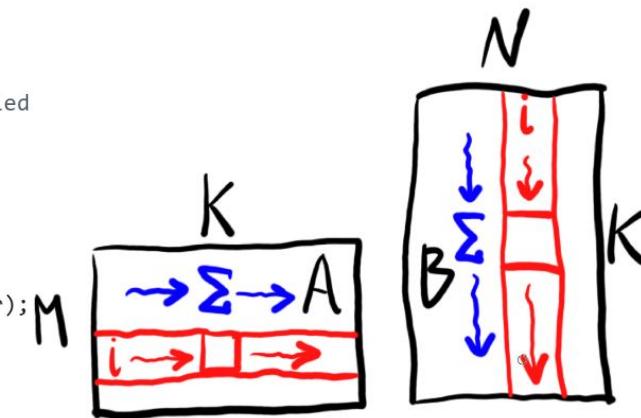


```

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;
90
91     .....
92
113     // Load in the current value of c, scale it by beta, and add this our result scaled
114     int cRow = warpM * WMMA_M;
115     int cCol = warpN * WMMA_N;
116
117     if (cRow < M && cCol < N) {
118         wmma::load_matrix_sync(c_frag, c + cRow + cCol * ldc, ldc, wmma::mem_col_major);
119
120 #pragma unroll
121     for(int i=0; i < c_frag.num_elements; i++) {
122         c_frag.x[i] = alpha * acc_frag.x[i] + beta * c_frag.x[i];
123     }
}

```

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



А не будет тормозить?

Почему не спец. функция в железе?

А почему бы не добавить к С в отдельном кернеле?

```

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

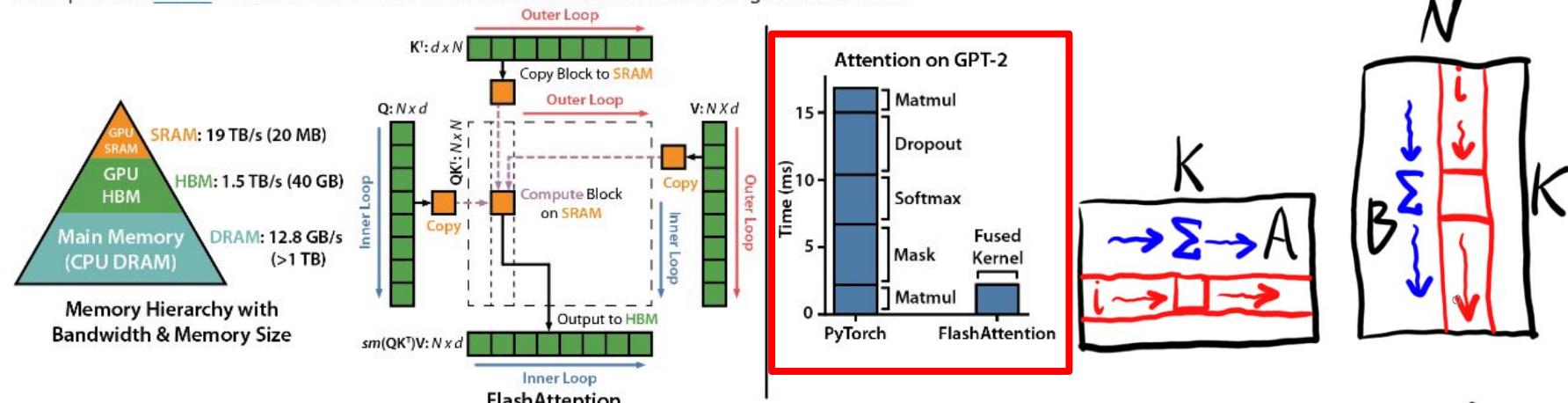
```

FlashAttention: Fast and Memory-Efficient Exact Attention with IO-Awareness

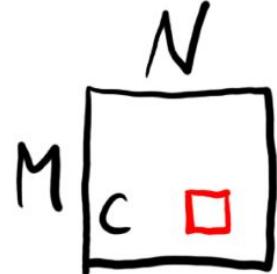
Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, Christopher Ré

Paper: <https://arxiv.org/abs/2205.14135> <https://github.com/Dao-AILab/flash-attention>

IEEE Spectrum [article](#) about our submission to the MLPerf 2.0 benchmark using FlashAttention.



А не будет тормозить?
Почему не спец. функция в железе?
А почему бы не добавить к С в отдельном кернеле?

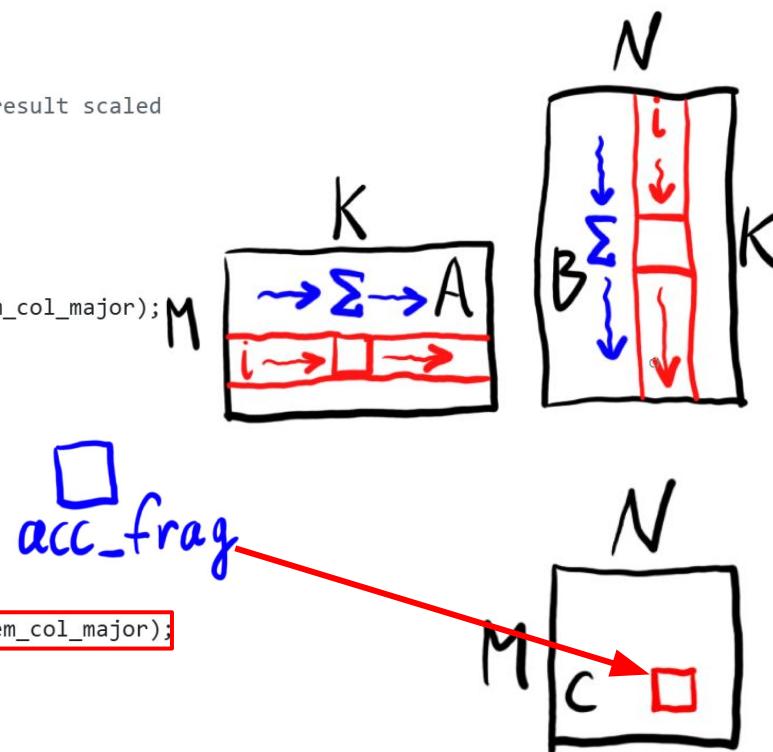


```

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;
90
91     .....
92
113     // Load in the current value of c, scale it by beta, and add this our result scaled
114     int cRow = warpM * WMMA_M;
115     int cCol = warpN * WMMA_N;
116
117     if (cRow < M && cCol < N) {
118         wmma::load_matrix_sync(c_frag, c + cRow + cCol * ldc, ldc, wmma::mem_col_major);
119
120 #pragma unroll
121         for(int i=0; i < c_frag.num_elements; i++) {
122             c_frag.x[i] = alpha * acc_frag.x[i] + beta * c_frag.x[i];
123         }
124
125         // Store the output
126         wmma::store_matrix_sync(c + cRow + cCol * ldc, c_frag, ldc, wmma::mem_col_major);
127     }
128 }

```

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



Умножение матриц - Tensor Cores, WMMA

Сколько у нас вычислений?

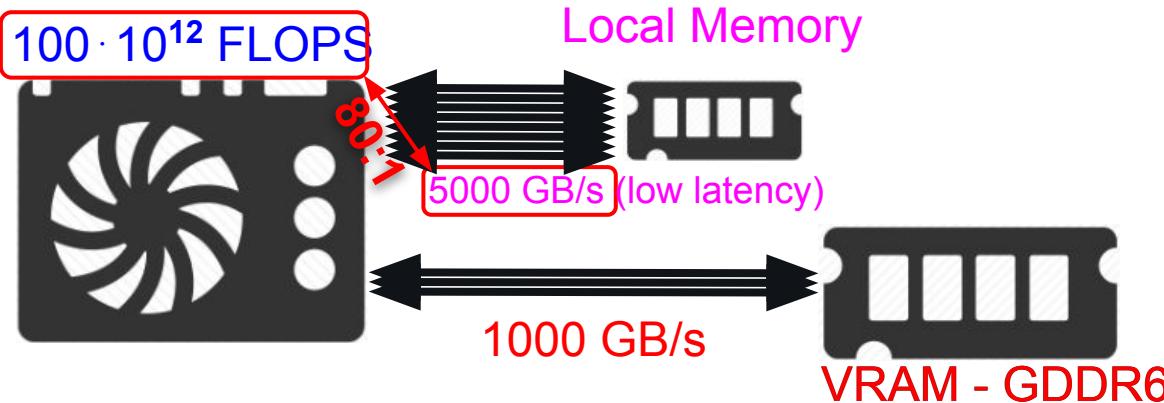
$$O(N \cdot M \cdot K)$$

Сколько у нас чтений/записей данных?

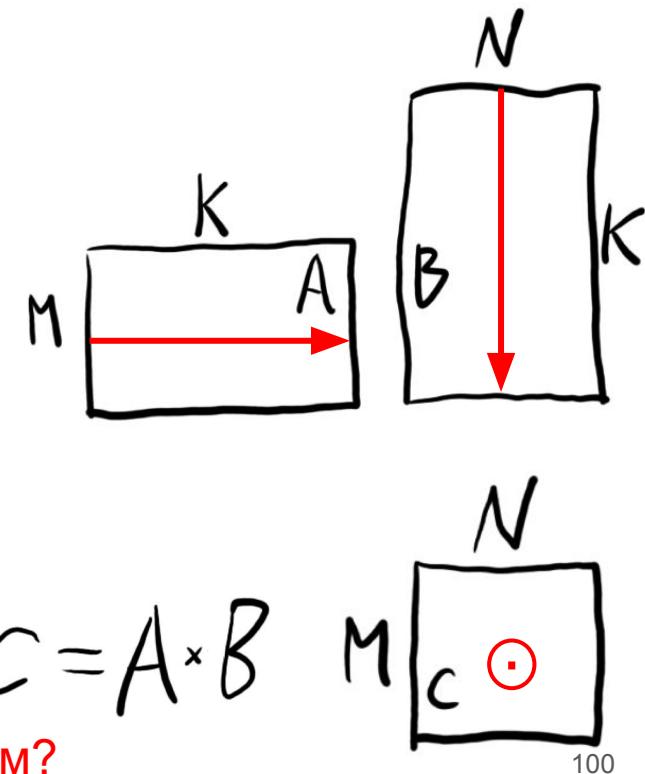
$$O(N \cdot M \cdot K)$$

Какая пропорция?

1:1 **Memory-bound!**



Какая у нас теперь пропорция вычислений к чтениям?



Умножение матриц - Tensor Cores, WMMA

Сколько у нас вычислений?

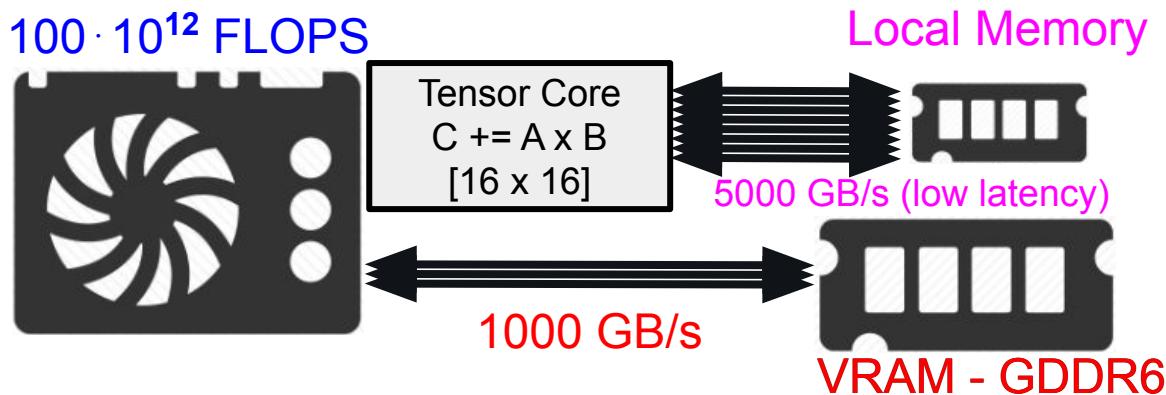
$$O(N \cdot M \cdot K)$$

Сколько у нас чтений/записей данных?

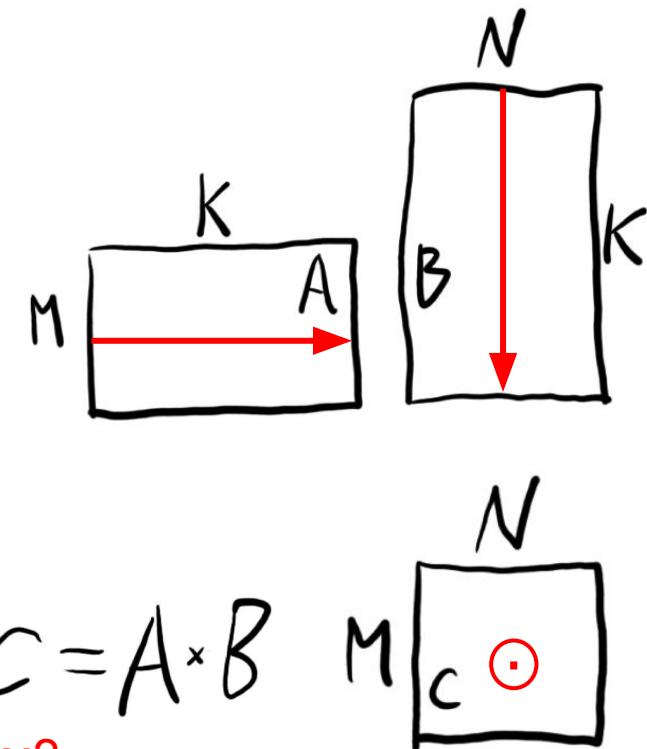
$$O(N \cdot M \cdot K)$$

Какая пропорция?

1:1 **Memory-bound!**



Какая у нас теперь пропорция вычислений к чтениям?



Умножение матриц

Сколько у нас вычислений?

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

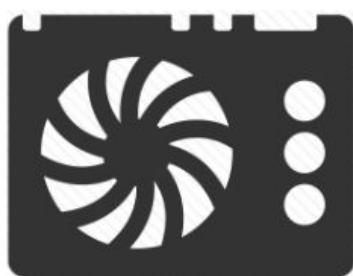
Сколько у нас чтений/записей данных?

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

Какая пропорция?

1:1 **Memory-bound!**

$100 \cdot 10^{12}$ FLOPS



Tensor Core
 $C += A \times B$
[16×16]

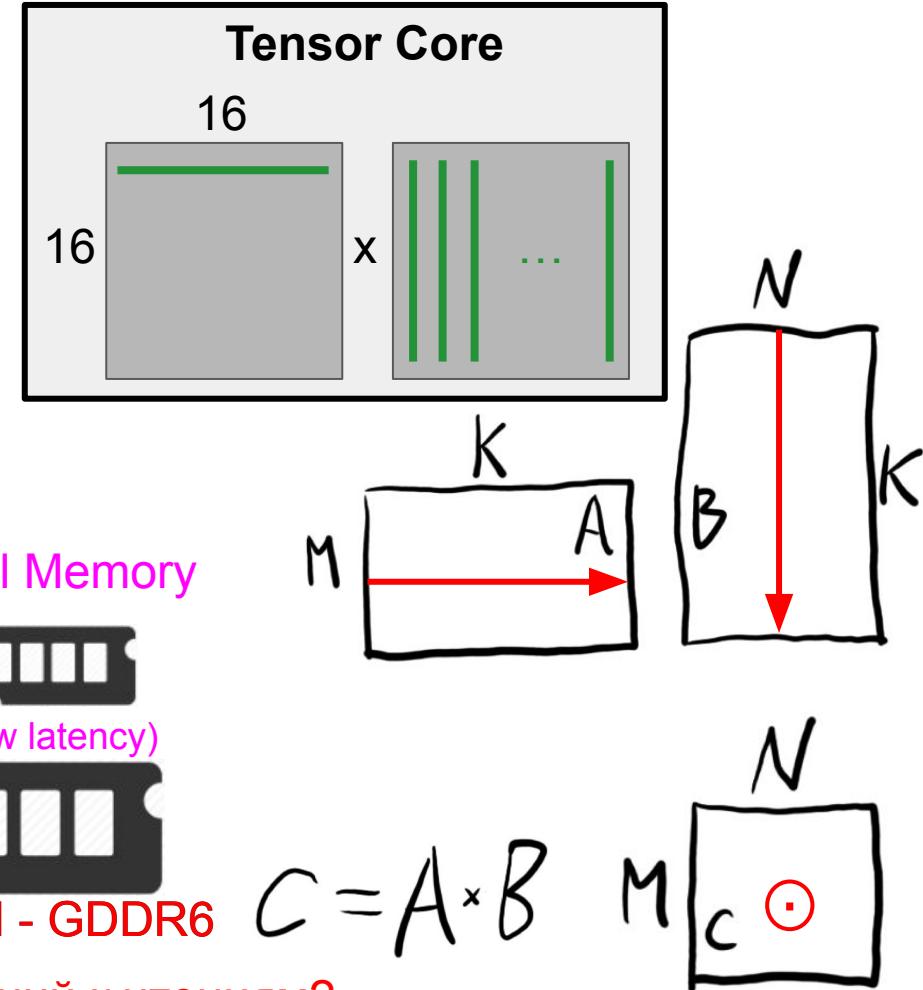
Local Memory

5000 GB/s (low latency)

1000 GB/s



VRAM - GDDR6



Какая у нас теперь пропорция вычислений к чтениям?

✓ cutlass name trick

 Mogball committed on Jul 9

secret techniques

```
  ⌂ 23 python/tutorials/gluon/01-attention-forward.py □
  ...  @@ -834,13 +834,20 @@ def _attn_fwd_correction(config, chnl, desc, M, STAGE: gl.constexpr):
834      config, prog, s1_tmem, M, corri_consumer, epi_producer, o_consumer)
835
836
837 - @gluon.jit(do_not_specialize=["Z"])
838 - def gluon_attn(sm_scale, M, Z, H, N_CTX, #
839 -                 desc_q, desc_k, desc_v, desc_o, #
840 -                 BLOCK_M: gl.constexpr, BLOCK_N: gl.constexpr, HEAD_DIM: gl.constexpr, #
841 -                 GROUP_SIZE_N: gl.constexpr, NUM_SMS: gl.constexpr, #
842 -                 STAGE: gl.constexpr, dtype: gl.constexpr, #
843 -                 num_warps: gl.constexpr):
837 + def attention_repr(specialization):
838 +     name = "gluon_attention"
839 +     # Up to 150 TFLOPS faster for fp8!
840 +     if specialization.constants["dtype"] == gl.float8e5:
841 +         name = "cutlass_" + name
```



peterbell10 on Jul 9

Contributor ...

very cool... did you check if other names change the scheduling (e.g. because of non-determinism or code alignment) or if it's literally just special cased for cutlass.



Mogball on Jul 9

Collaborator Author ...

| it's literally just special cased for cutlass.

Yup

46

8

5

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>

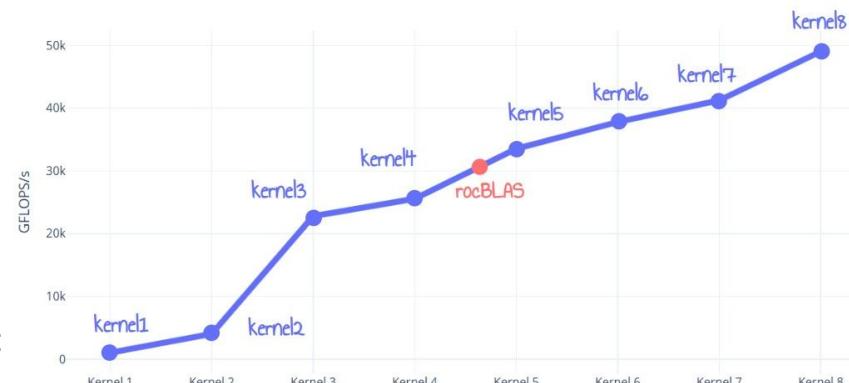
https://github.com/NVIDIA/cutlass/blob/main/examples/00_basic_gemm/basic_gemm.cu

<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://cnugteren.github.io/tutorial/pages/page15.html>
- <https://siboehm.com/articles/22/CUDA-MMM>



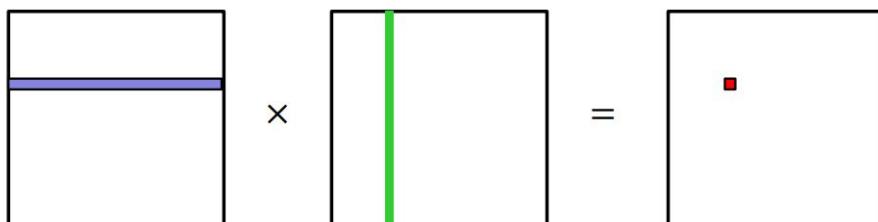
Scalar Loop

- Each invocation computes one element of the result matrix
- Lots of redundant loads
- Example of what NOT to do

```
uint i = gl_GlobalInvocationID.y;
uint j = gl_GlobalInvocationID.x;
float16_t C = inputC.x[sC * i + j];

for (uint k = 0; k < K; ++k) {
    float16_t A = inputA.x[sA * i + k];
    float16_t B = inputB.x[sB * k + j];
    C += A*B;
}

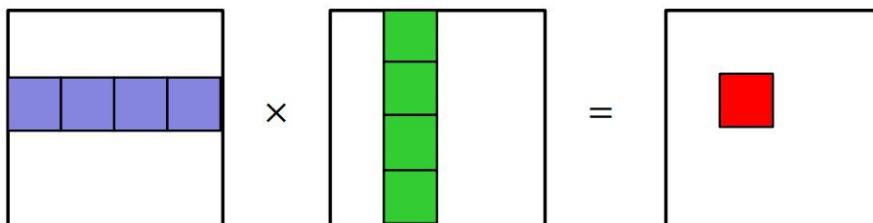
outputD.x[sD * i + j] = C;
```



| GPU | TFLOPS |
|-----------|--------|
| RTX 2070 | 0.232 |
| RTX TITAN | 0.528 |

Simple Cooperative Multiply

- A bit more coordinate calculation, but still a simple sum($A_{ik}B_{kj}$)
- Still very memory-limited, but improved



```
1M = 16; 1N = 8; 1K = 8;
fcoopmatNV<16, gl_ScopeSubgroup, 1M, 1K> matA;
fcoopmatNV<16, gl_ScopeSubgroup, 1K, 1N> matB;
fcoopmatNV<16, gl_ScopeSubgroup, 1M, 1N> matC;

uvec2 matrixID = uvec2(gl_WorkGroupID);
uint cRow = 1M * matrixID.y;
uint cCol = 1N * matrixID.x;

coopMatLoadNV matC, inputC.x, sC * cRow + cCol, sC, false);

for (uint k = 0; k < K; k += 1K) {
    uint aRow = 1M * matrixID.y;
    uint aCol = k;
    coopMatLoadNV matA, inputA.x, sA * aRow + aCol, sA, false);

    uint bRow = k;
    uint bCol = 1N * matrixID.x;
    coopMatLoadNV matB, inputB.x, sB * bRow + bCol, sB, false);

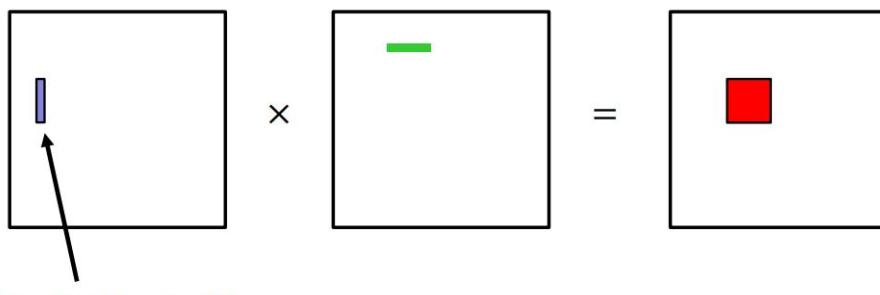
    matC = coopMatMulAddNV(matA, matB, matO);
}

coopMatStoreNV(matC, outputD.x, sD * cRow + cCol, sD, false);
```

| GPU | TFLOPS |
|-----------|--------|
| RTX 2070 | 2.00 |
| RTX TITAN | 3.70 |

Tiled Scalar Multiply

- Illustrates a way to tile the multiply
 - Load row, load column, outer product
 - Name of the game is to load once, then perform as many multiplies as possible
 - Limited by register file size
 - Maybe 8x8 tile size per invocation



```

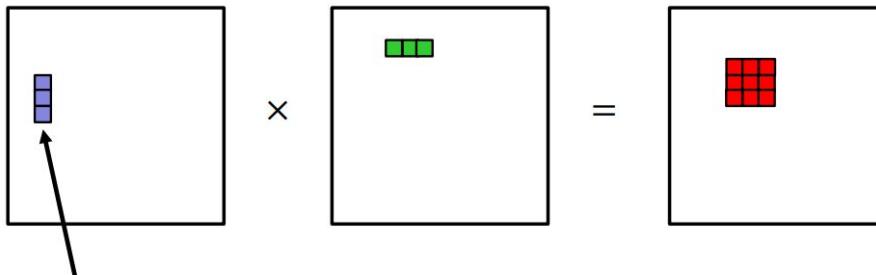
float16_t C[C_ROWS][C_COLS];
uvec2 tileID = uvec2(gl_WorkGroupID.xy);
uvec2 inv = uvec2(gl_LocalInvocationID.xy);
// load C
...
// iterate through K dimension and accumulate
for (uint k = 0; k < K; ++k) {
    float16_t A[C_ROWS];
    for (uint i = 0; i < C_ROWS; ++i) {
        uint gi = TILE_M * tileID.y + (C_ROWS * inv.y + i);
        uint gk = k;
        A[i] = inputA.x[sA * gi + gk];
    }
    float16_t B;
    for (uint j = 0; j < C_COLS; ++j) {
        uint gk = k;
        uint gj = TILE_N * tileID.x + (C_COLS * inv.x + j);
        B = inputB.x[sB * gk + gj];
        for (uint i = 0; i < C_ROWS; ++i) {
            C[i][j] = A[i] * B + C[i][j];
        }
    }
}
// store C
...

```

| GPU | TFLOPS |
|-----------|--------|
| RTX 2070 | 1.45 |
| RTX TITAN | 2.79 |

Tiled Cooperative Matrix Multiply

- Apply tiling approach to coop matrices
- Effectively 32x the register file, since matrices are spread over a subgroup
 - 112x112x16 multiply per outer loop iter
 - Takes advantage of cheap communication within subgroups to maximize reuse
- Shader code is no more complex than what you'd write for scalar

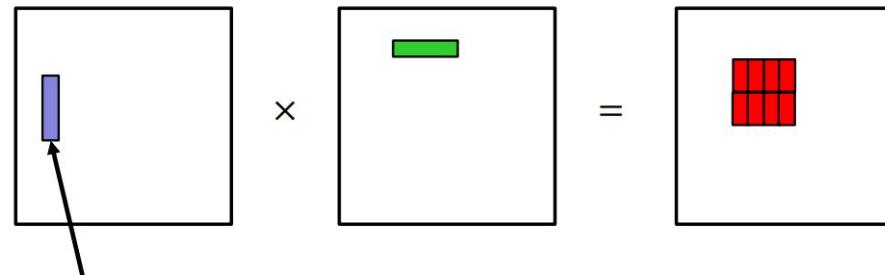


```
fcoopmatNV<16, gl_ScopeSubgroup, 1M, 1N> matC[C_ROWS][C_COLS];
uvec2 tileID = uvec2(gl_WorkGroupID.xy);
// load matC
...
for (uint k = 0; k < K; k += TILE_K) {
    fcoopmatNV<16, gl_ScopeSubgroup, 1M, 1N> matA[C_ROWS];
    for (uint i = 0; i < C_ROWS; ++i) {
        uint gi = TILE_M * tileID.y + 1M * i;
        uint gk = k;
        coopMatLoadNV(matA[i], inputA.x, sA * gi + gk, sA, false);
    }
    fcoopmatNV<16, gl_ScopeSubgroup, 1K, 1N> matB;
    for (uint j = 0; j < C_COLS; ++j) {
        uint gj = TILE_N * tileID.x + 1N * j;
        coopMatLoadNV(matB, inputB.x, sB * gk + gj, sB, false);
        for (uint i = 0; i < C_ROWS; ++i) {
            matC[i][j] = coopMatMulAddNV(matA[i], matB, matC[i][j]);
        }
    }
    // store matC
    ...
}
```

| GPU | TFLOPS |
|-----------|--------|
| RTX 2070 | 21.8 |
| RTX TITAN | 42.1 |

Staging Through Shared Memory

- Subgroups cooperate to copy data from buffer to shared, then load out of shared
 - Overlap buffer loads with matrix math (i.e. pipeline load for next iteration)
- Example (FP16): 8 subgroups split a 256x256 tile into 8 128x64 tiles (K=32)
 - Cooperate to copy A block (256x32) and B block (32x256) into shared memory
 - Then each subgroup loads the portions it needs from shared memory
 - Only 128B of buffer loads per thread per K iteration overlapping with 8K FMADs
- Example (INT8): 8 subgroups split a 128x256 tile into 8 64x64 tiles (K=64)
 - Accumulator is INT32, requires 2x register file, hurts latency hiding
 - Double math rate, half tile size = 4x harder to hide latency, much farther from SOL



| GPU | INT8 TOPS |
|-----------|-----------|
| RTX 2070 | 66.0 |
| RTX TITAN | ~130 |

| GPU | FP16 TFLOPS |
|-----------|-------------|
| RTX 2070 | 49.0 |
| RTX TITAN | ~100 |

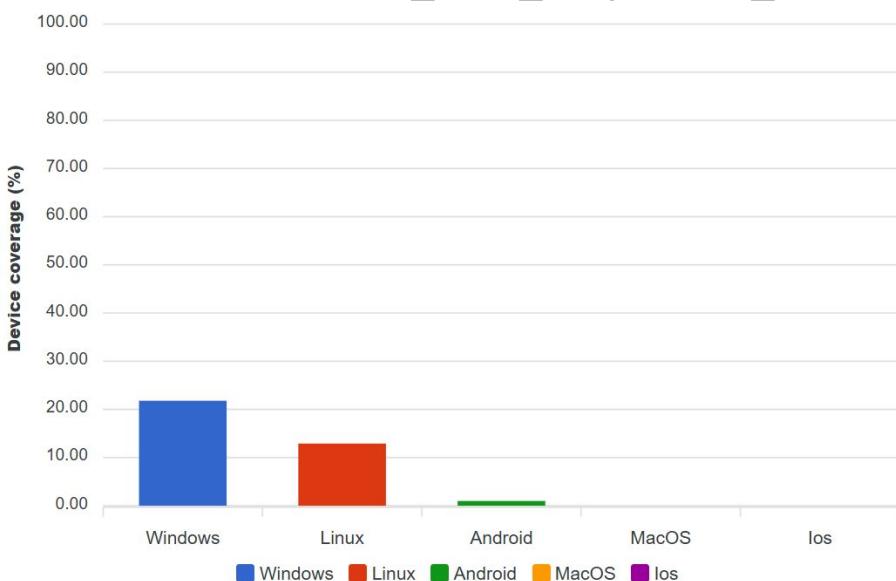
Device coverage for VK_KHR_cooperative_matrix

Vulkan registry manpage

This extension has additional features and properties

Extension was first submitted at 2023-06-27

Vulkan extension: VK_KHR_cooperative_matrix



Listing first known driver version support for VK_KHR_cooperative_matrix

This extension has additional features and properties

Extension was first submitted at 2023-06-27

All platforms

Windows

Linux

Android

macOS

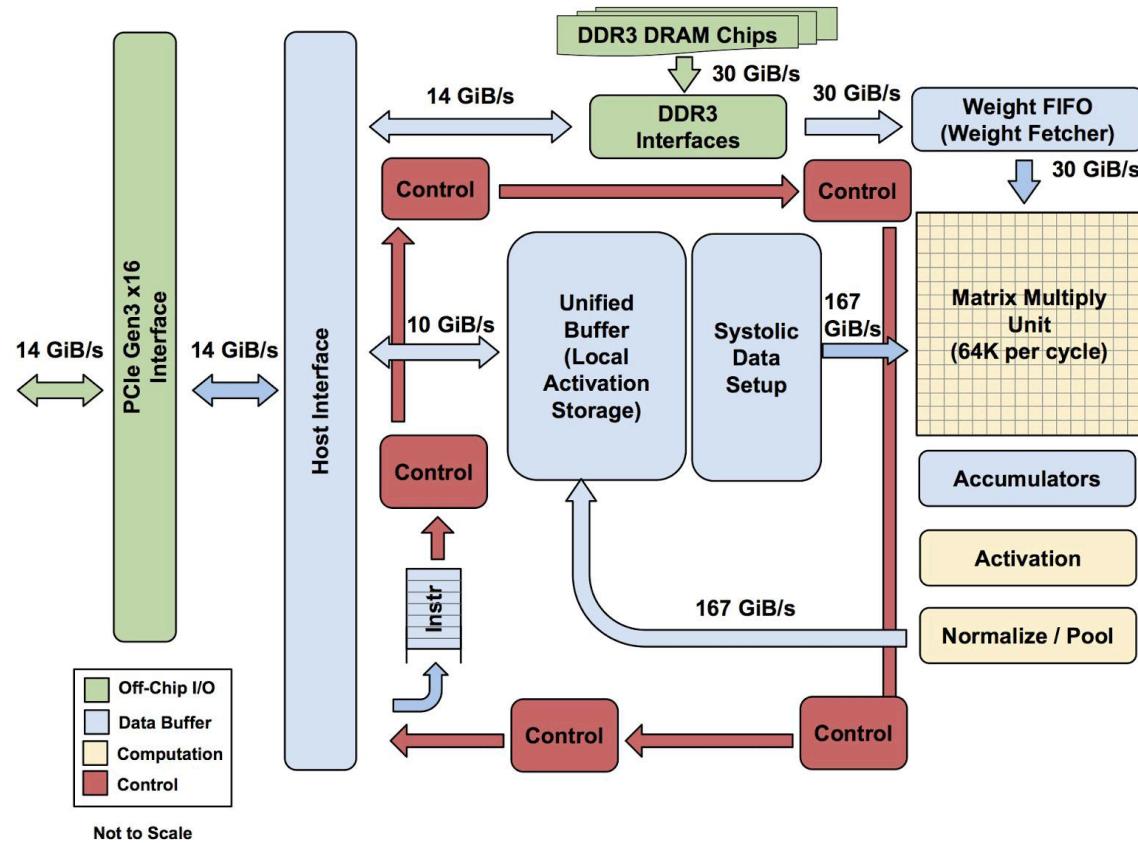
iOS

Show 50 entries

| Device | Vendor | Driver [?] | Date | Compare |
|------------------------------|--------|------------|------------|---------|
| AMD Radeon(TM) 740M | AMD | 2.0.302 | 2025-09-16 | Add |
| AMD Radeon RX 5500 | AMD | 2.0.310 | 2025-08-26 | Add |
| AMD Radeon Pro 5500M | AMD | 2.0.317 | 2025-08-13 | Add |
| AMD Radeon(TM) 860M Graphics | AMD | 2.0.331 | 2025-07-23 | Add |

| Device | Vendor | Driver [?] | Date | Compare |
|--|--------|------------|------------|---------|
| Intel(R) Iris(R) Xe Graphics (ADL GT2) | INTEL | 6400.4099 | 2025-07-09 | Add |
| Intel(R) Arc(TM) 130V GPU (8GB) | INTEL | 101.6314 | 2025-07-06 | Add |
| Intel(R) Arc(TM) 130V GPU (16GB) | INTEL | 101.6556 | 2025-06-23 | Add |
| Intel(R) Arc(TM) 140T GPU (32GB) | INTEL | 101.6554 | 2025-05-31 | Add |
| Intel(R) Arc(TM) B570 Graphics | INTEL | 101.6790 | 2025-05-19 | Add |

Tensor Processor Unit - TPU (Google)



Глава 4: Умножение матриц

DeepSeek

Матрицы и векторы — это мощные инструменты для представления и обработки данных. Одним из основных операций с матрицами является умножение. Важно отметить, что умножение матриц не является коммутативной операцией, то есть $A \cdot B \neq B \cdot A$. Поэтому при вычислении произведения матриц необходимо учитывать порядок умножения.

Для вычисления произведения матриц A и B необходимо, чтобы количество столбцов в матрице A было равно количеству строк в матрице B . Результатом умножения матриц A и B будет матрица C , где количество строк в C будет равно количеству строк в A , а количество столбцов в C будет равно количеству столбцов в B .

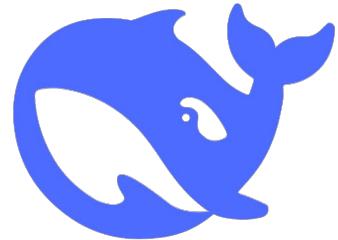
Процесс вычисления произведения матриц можно представить как последовательное выполнение скалярного умножения. Для каждого элемента результата C_{ij} вычисляется сумма произведений соответствующих элементов из i -й строки матрицы A и j -го столбца матрицы B . Для вычисления каждого элемента результата требуется $n \cdot m$ умножений, где n — количество строк в A , а m — количество столбцов в B .

Однако, если учесть, что умножение матриц является критически важной операцией в машинном обучении, то становится очевидным, что существуют способы оптимизации этого процесса. Одним из таких способов является алгоритм DeepSeek, разработанный группой исследователей из Университета Торонто. Алгоритм DeepSeek использует методы глубокого обучения для предсказания результатов умножения матриц на основе их структуры и свойств.

Алгоритм DeepSeek работает следующим образом. Сначала он получает на вход две матрицы A и B . Затем он анализирует структуру матриц и определяет, какие из них являются более сложными. На основе этого анализа алгоритм определяет, какой из матриц лучше всего умножить сначала. Далее алгоритм вычисляет предсказание для каждого элемента результата C . Алгоритм также оптимизирует процесс вычисления, чтобы минимизировать количество умножений.

Благодаря использованию методов глубокого обучения, алгоритм DeepSeek может эффективно обрабатывать большие матрицы и векторы. Он также может адаптироваться к различным типам данных и структурам. Алгоритм DeepSeek был протестирован на различных задачах машинного обучения и показал высокую точность и производительность.

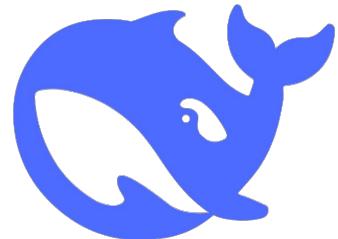
В заключение можно сказать, что алгоритм DeepSeek является важным шагом в развитии технологий машинного обучения. Он позволяет оптимизировать один из самых критически важных операций в машинном обучении — умножение матриц. Алгоритм DeepSeek может стать основой для дальнейшего развития технологий машинного обучения и обработки данных.



DeepSeek: x2 ускорение обучения (fp16 → fp8)

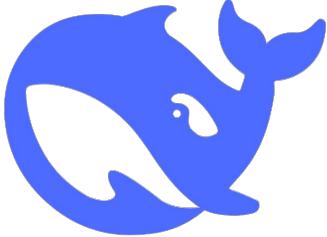


$$\pm (1 + \text{mantissa}) \times 2^{(\text{exponent} - 127)}$$



DeepSeek: x2 ускорение обучения (fp16 → fp8)

| | sign | exponent | | | | | | mantissa | | | | | | | | |
|----------|------|----------|---|---|---|---|---|----------|---|---|---|---|---|---|---|------------|
| FP16 | 0 | 0 | 1 | 1 | 0 | 1 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 1 | = 0.395264 |
| BF16 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 0 | 1 | 1 | 0 | 0 | 1 | 0 | 1 | = 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 |



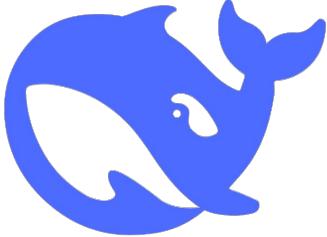
DeepSeek: x2 ускорение обучения (fp16 → fp8)



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**

*H800 - почти H100, но соответствует регуляциям экспорта из США в Китай (400 GBs NVlink вместо 600 GB/s + 10% медленнее + нет FP64)



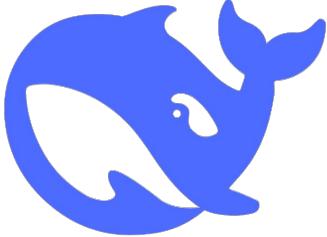
DeepSeek: x2 ускорение обучения (fp16 → fp8)



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



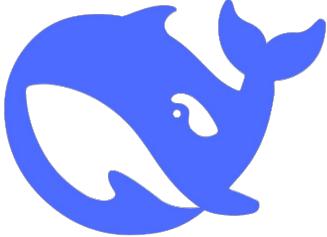
DeepSeek: x2 ускорение обучения (fp16 → fp8)



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

- Memory bandwidth: **3.35 TB/sec** Хватит ли пропускной способности памяти чтобы насытить ALU (tensor cores)?
 - FP 32: **67 TFlops**
 - FP 16: **268 TFlops**
 - FP 32 (tensor): **495 TFlops**
 - FP 16 (tensor): **990 TFlops**
 - FP 8 (tensor): **1979 TFlops**
- x2

*H800 - почти H100, но соответствует регуляциям экспорта из США в Китай (400 GBs NVlink вместо 600 GB/s + 10% медленнее + нет FP64)



DeepSeek: x2 ускорение обучения (fp16 → fp8)

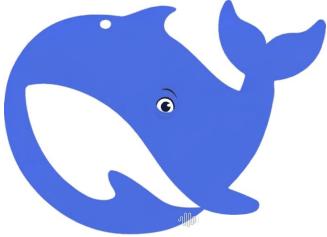


NVIDIA H100 (почти то же что и H800*):

- 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**

Какие риски?
Почему так не делают все?

*H800 - почти H100, но соответствует регуляциям экспорта из США в Китай (400 GBs NVlink вместо 600 GB/s + 10% медленнее + нет FP64)

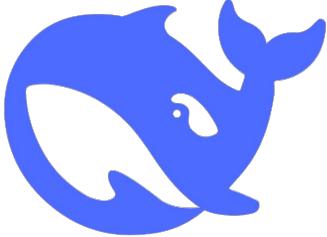


DeepSeek: fine-grained fp8 quantization



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!



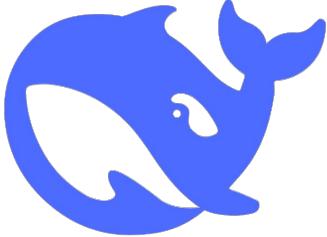
DeepSeek: fine-grained fp8 quantization



fp32



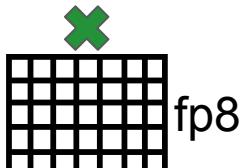
fp8

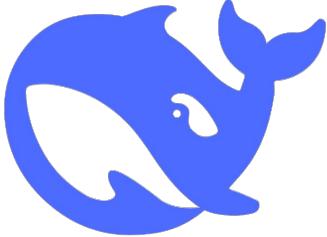


DeepSeek: fine-grained fp8 quantization

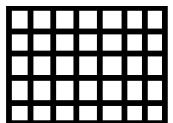


Scaling Factor **fp32**





DeepSeek: fine-grained fp8 quantization



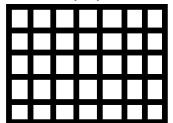
fp32 Input Values



Scaling Factor

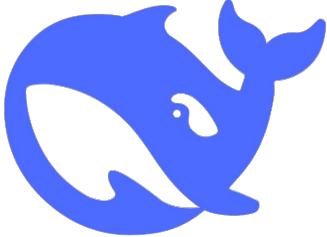
fp32

Каким его выбрать?

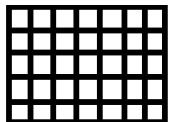


fp8

Input Values / Scaling Factor



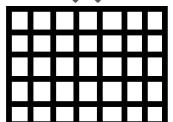
DeepSeek: fine-grained fp8 quantization



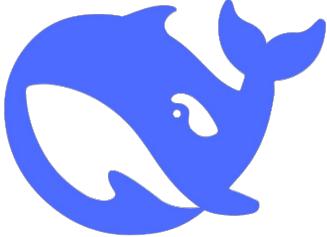
fp32 Input Values



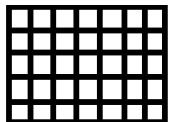
Scaling Factor **fp32** $\text{absmax}(\text{Input Values}) / 448$



fp8 Input Values / Scaling Factor



DeepSeek: fine-grained fp8 quantization

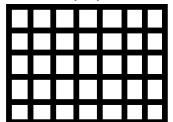


fp32 Input Values

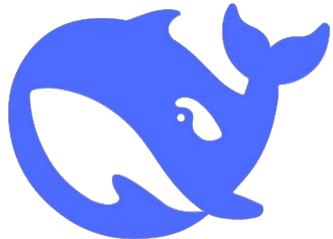


Scaling Factor

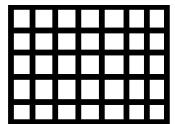
fp32 $\text{absmax}(\text{Input Values}) / 448$ Что это за волшебное число? 🦄



fp8 Input Values / Scaling Factor



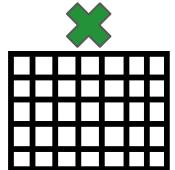
DeepSeek: fine-grained fp8 quantization



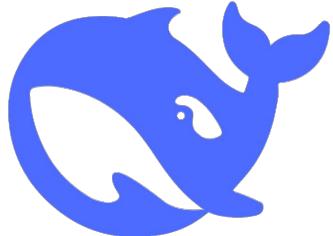
fp32 Input Values



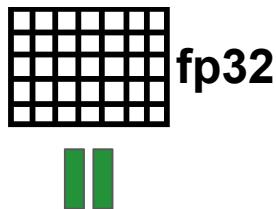
Scaling Factor **fp32** $\text{absmax}(\text{Input Values}) / \boxed{448} = \text{FP8_MAX}$



fp8 **Input Values / Scaling Factor** $\in [-448; +448]$



DeepSeek: fine-grained fp8 quantization

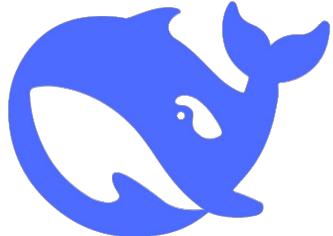


$$\begin{matrix} \text{fp32} & \times & \text{fp32} & = & \text{fp32} \end{matrix}$$

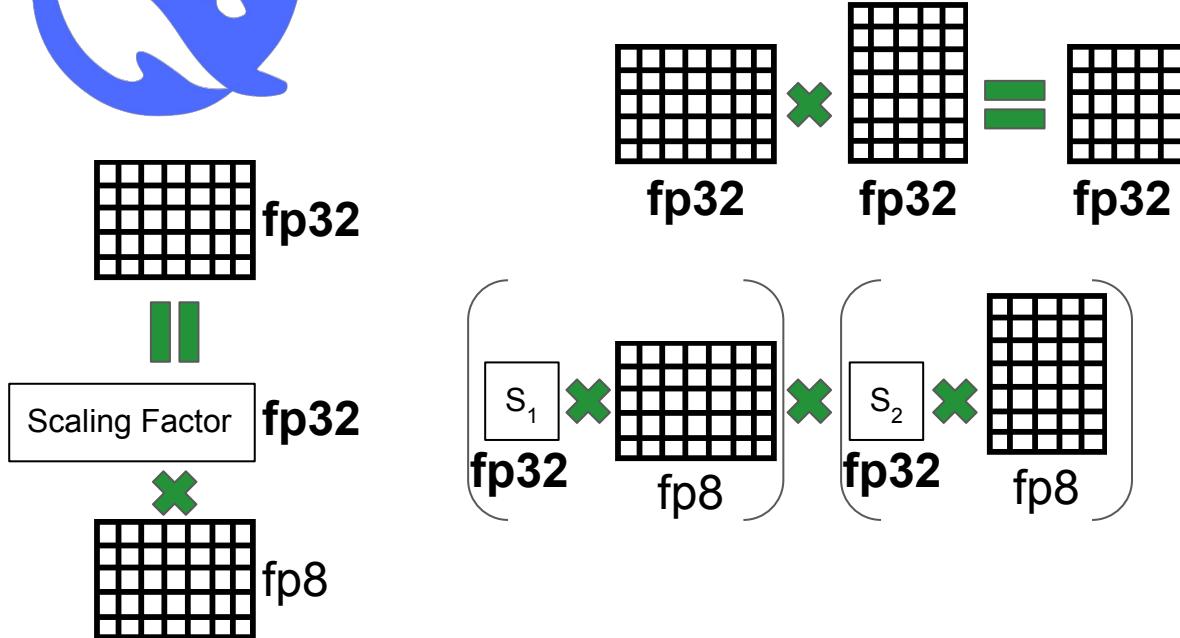


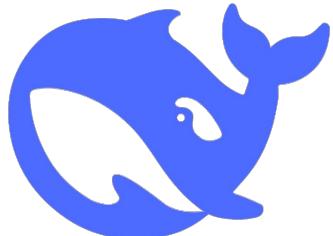
Scaling Factor **fp32**

$$\begin{matrix} \times & \text{fp8} \end{matrix}$$

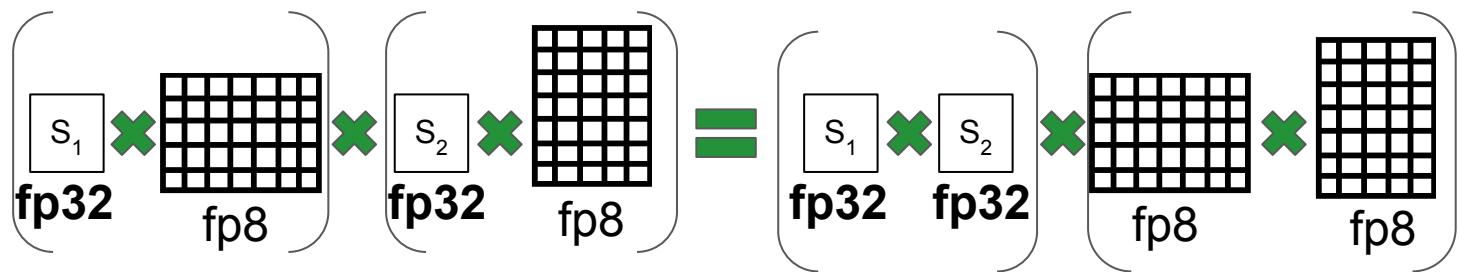
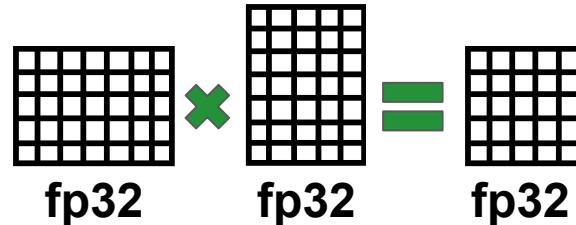
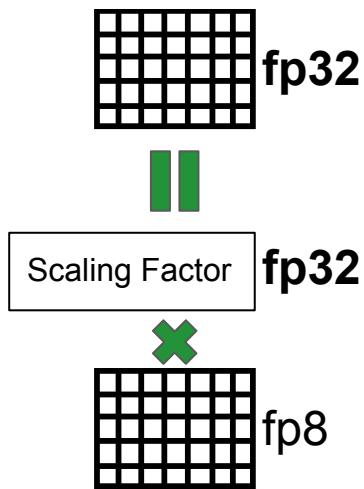


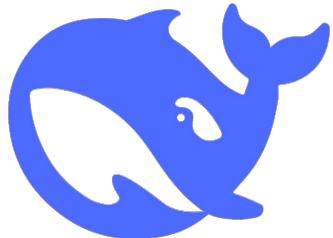
DeepSeek: fine-grained fp8 quantization



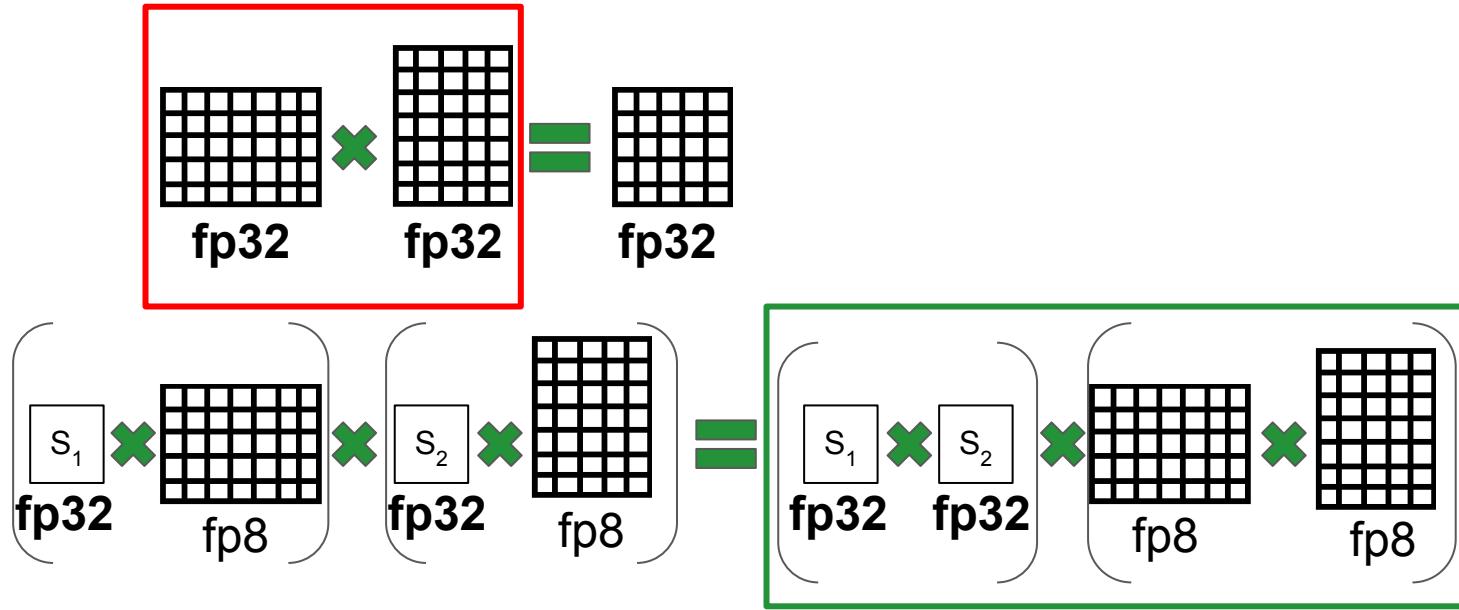
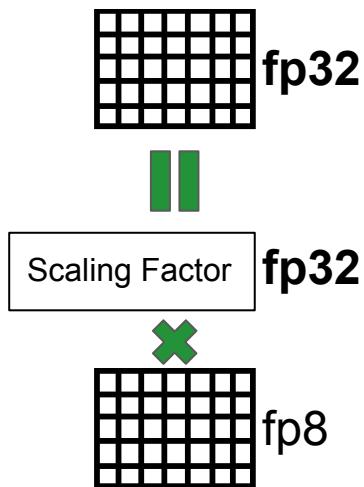


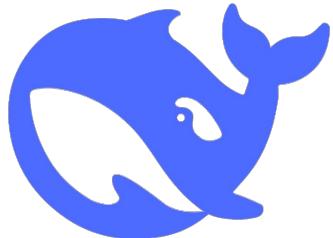
DeepSeek: fine-grained fp8 quantization



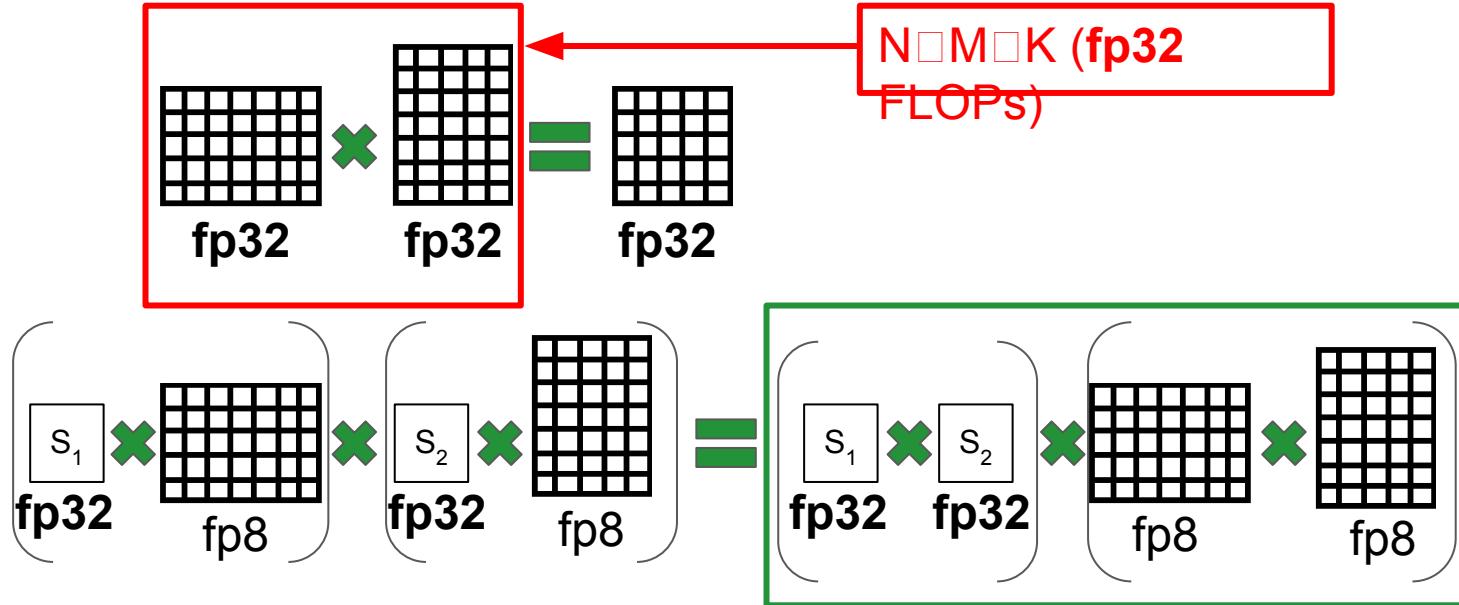
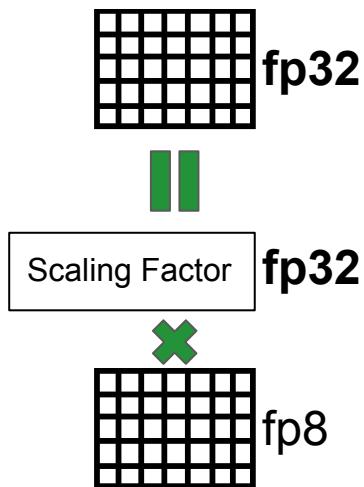


DeepSeek: fine-grained fp8 quantization

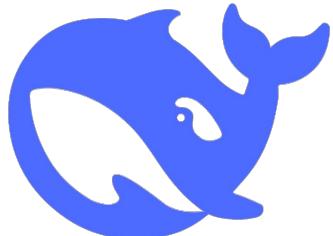




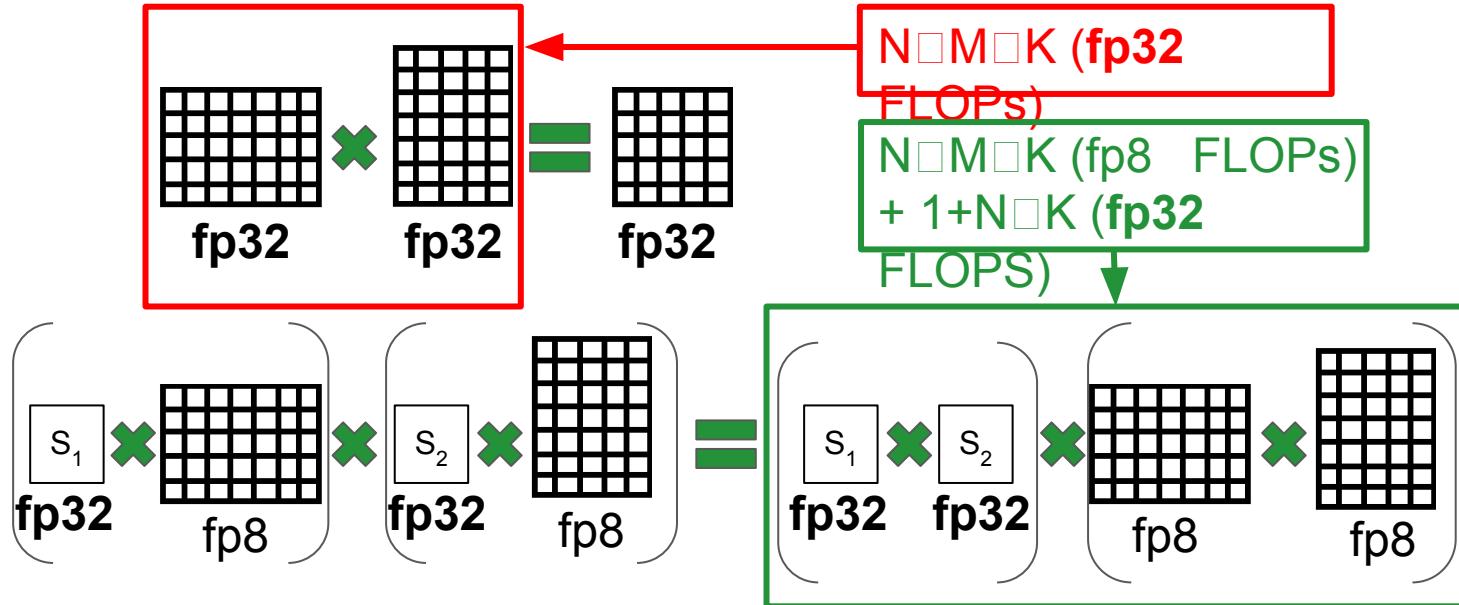
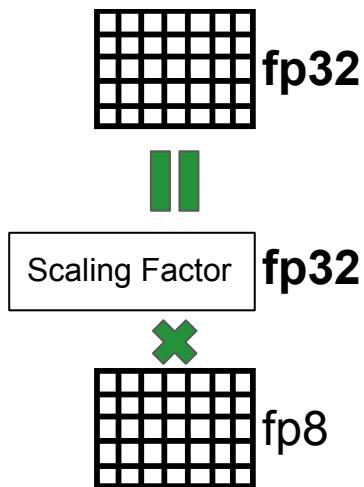
DeepSeek: fine-grained fp8 quantization



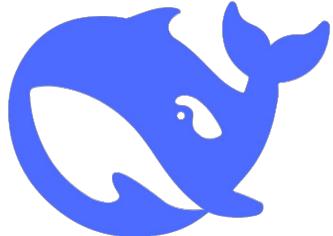
Насколько быстрее?



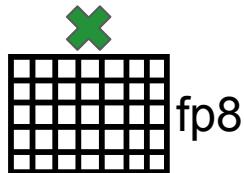
DeepSeek: fine-grained fp8 quantization



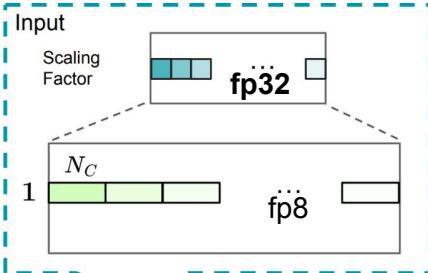
Насколько быстрее?



Scaling Factor **fp32**

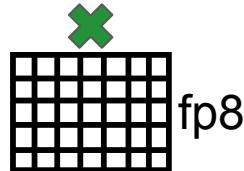


DeepSeek: fine-grained fp8 quantization

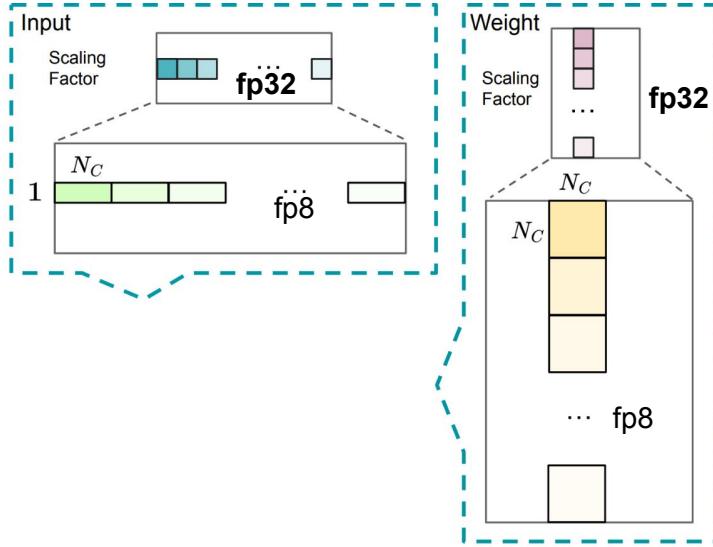


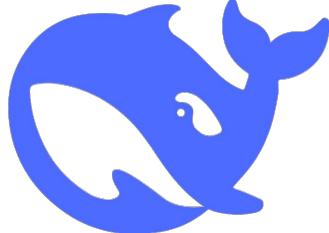


Scaling Factor **fp32**

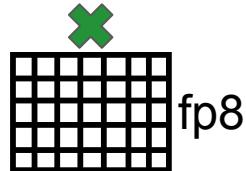


DeepSeek: fine-grained fp8 quantization

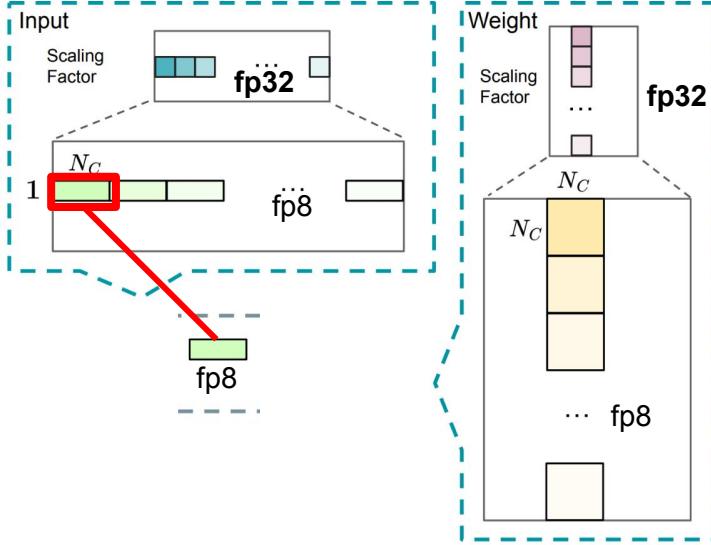


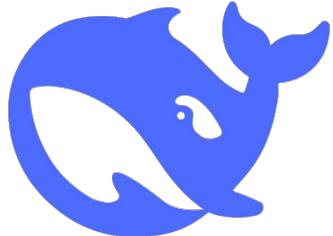


Scaling Factor **fp32**

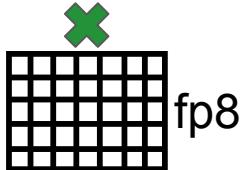


DeepSeek: fine-grained fp8 quantization

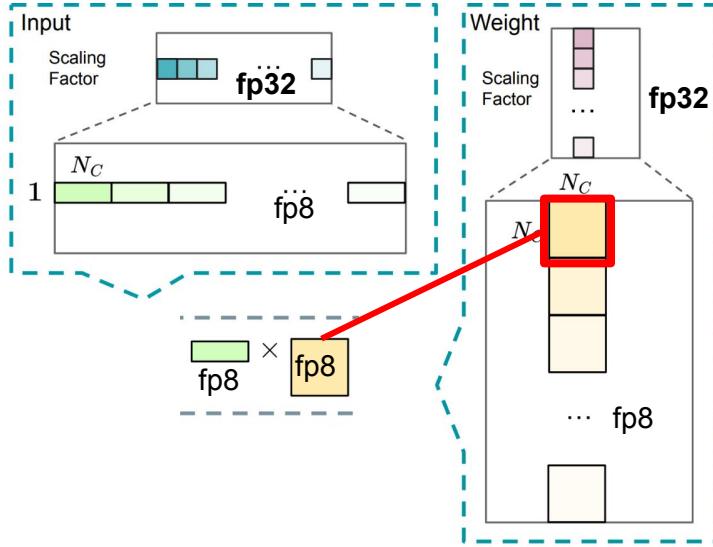


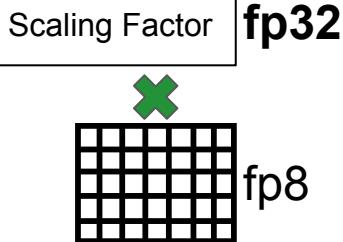
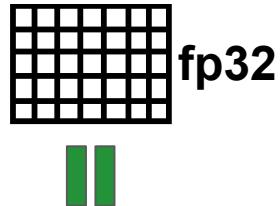
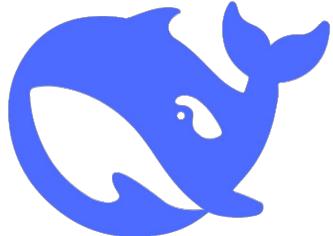


Scaling Factor **fp32**

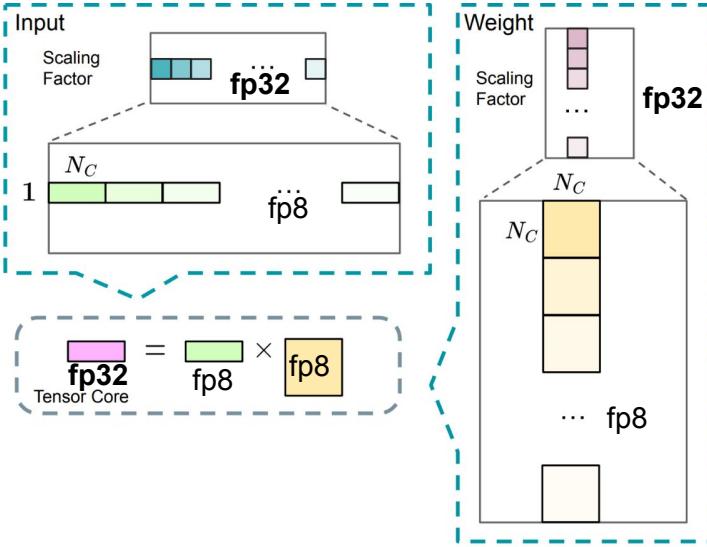


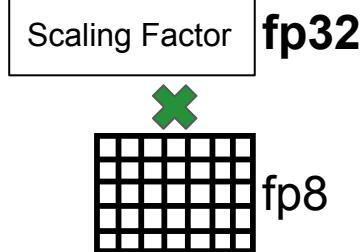
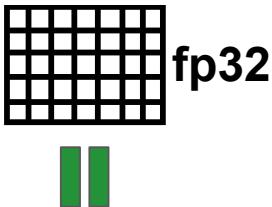
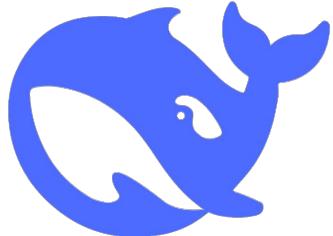
DeepSeek: fine-grained fp8 quantization



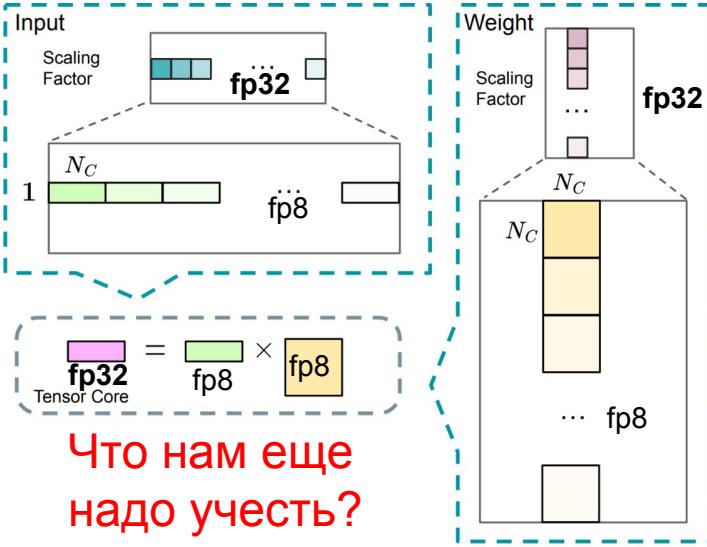


DeepSeek: fine-grained fp8 quantization

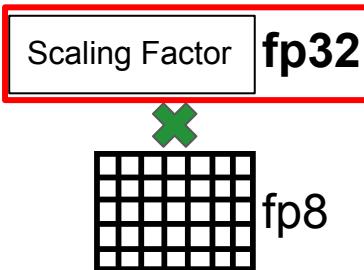
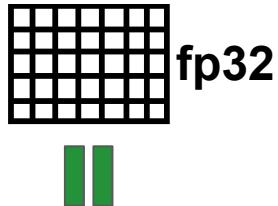
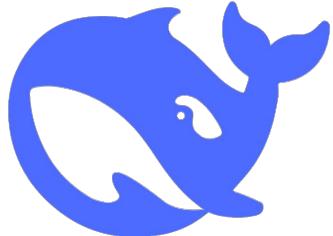




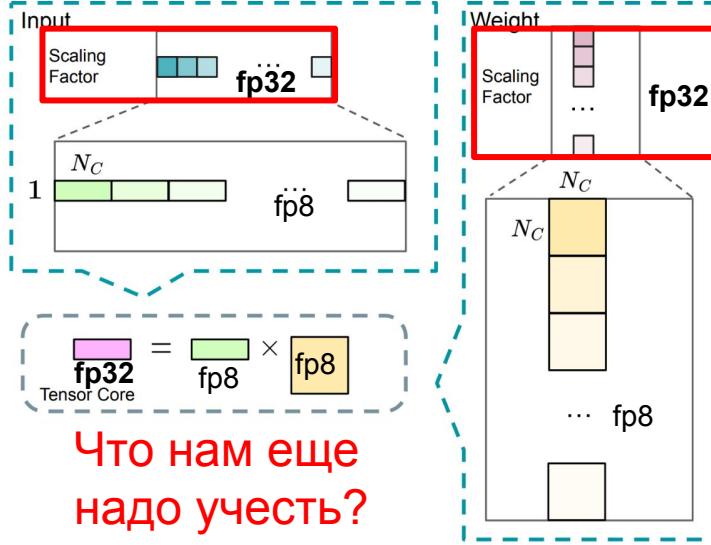
DeepSeek: fine-grained fp8 quantization



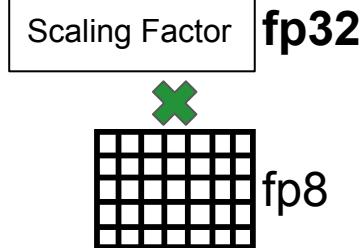
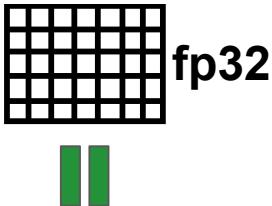
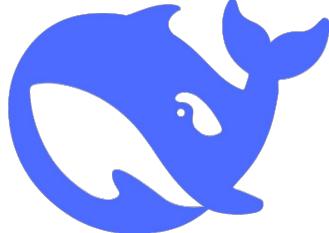
Что нам еще
надо учить?



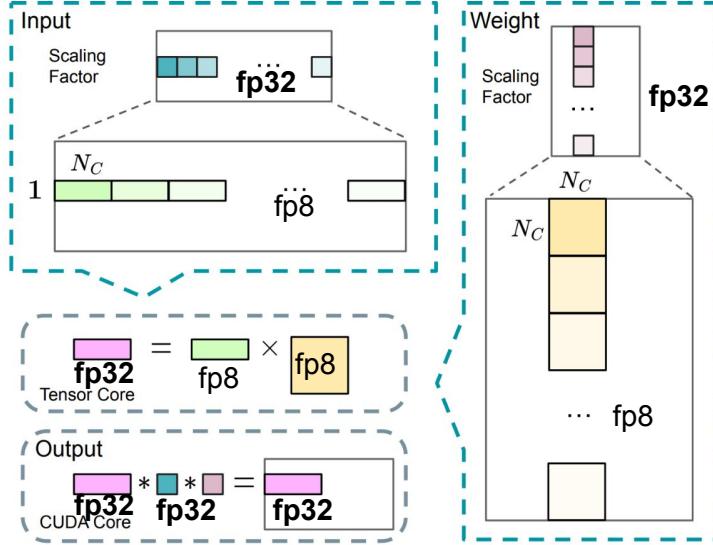
DeepSeek: fine-grained fp8 quantization



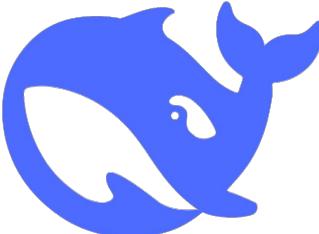
Что нам еще
надо учить?



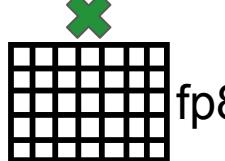
DeepSeek: fine-grained fp8 quantization



DeepSeek: fine-grained fp8 quantization



fp32



fp8

Scaling Factor

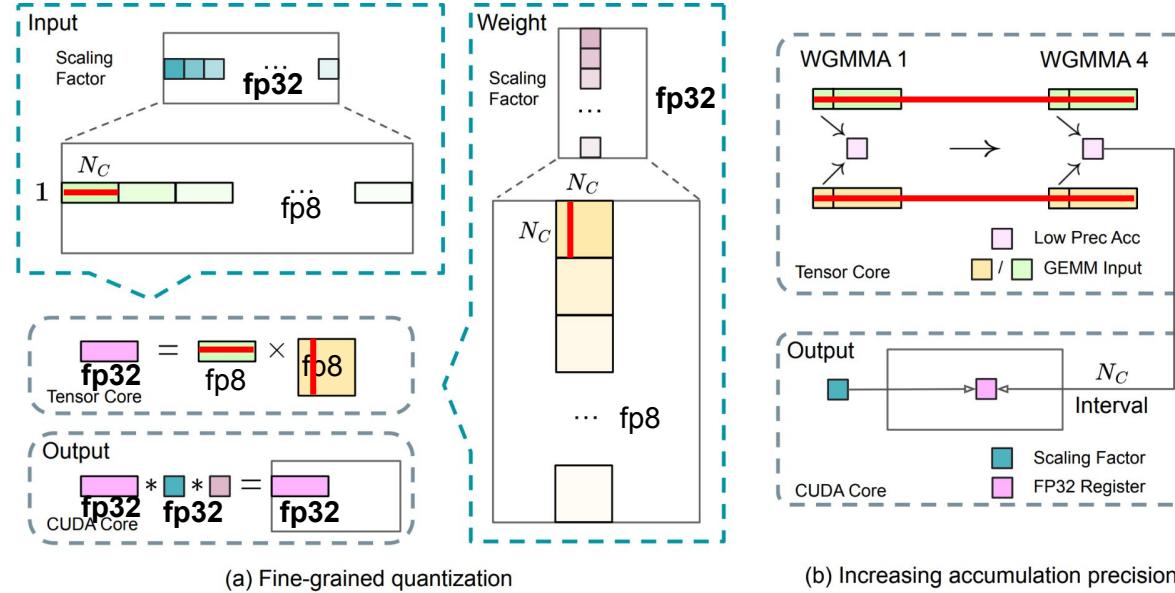
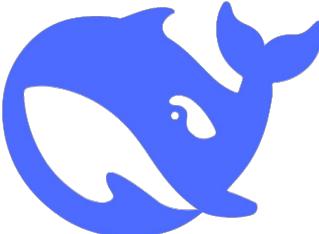


Figure 7 | (a) We propose a fine-grained quantization method to mitigate quantization errors caused by feature outliers; for illustration simplicity, only Fprop is illustrated. (b) In conjunction with our quantization strategy, we improve the FP8 GEMM precision by promoting to CUDA Cores at an interval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

DeepSeek: fine-grained fp8 quantization



fp32

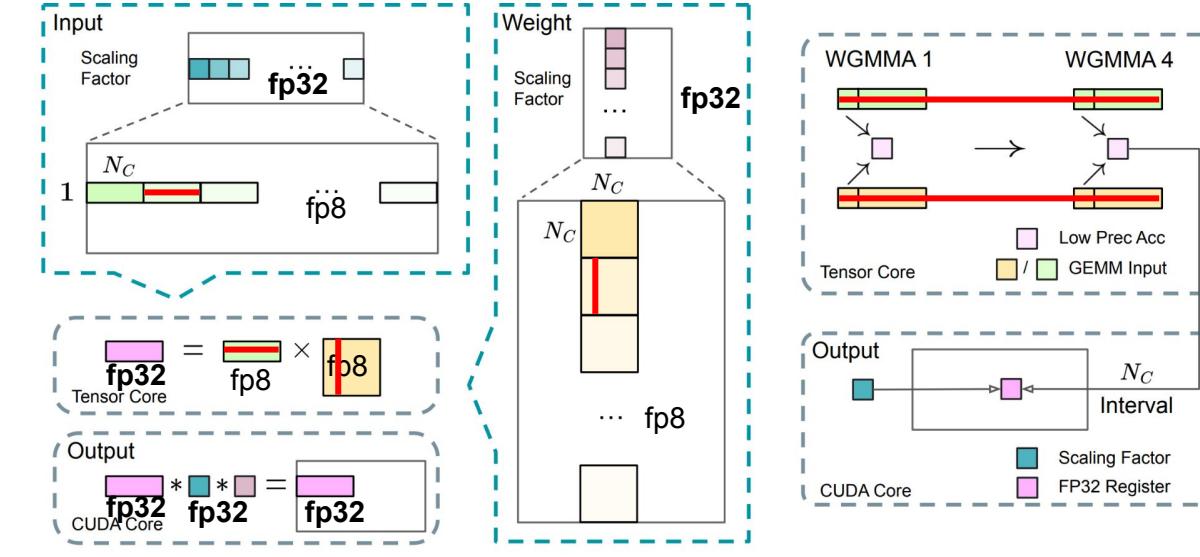


Figure 7 | (a) We propose a fine-grained quantization method to mitigate quantization errors caused by feature outliers; for illustration simplicity, only Fprop is illustrated. (b) In conjunction with our quantization strategy, we improve the FP8 GEMM precision by promoting to CUDA Cores at an interval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

DeepSeek: fine-grained fp8 quantization



fp32

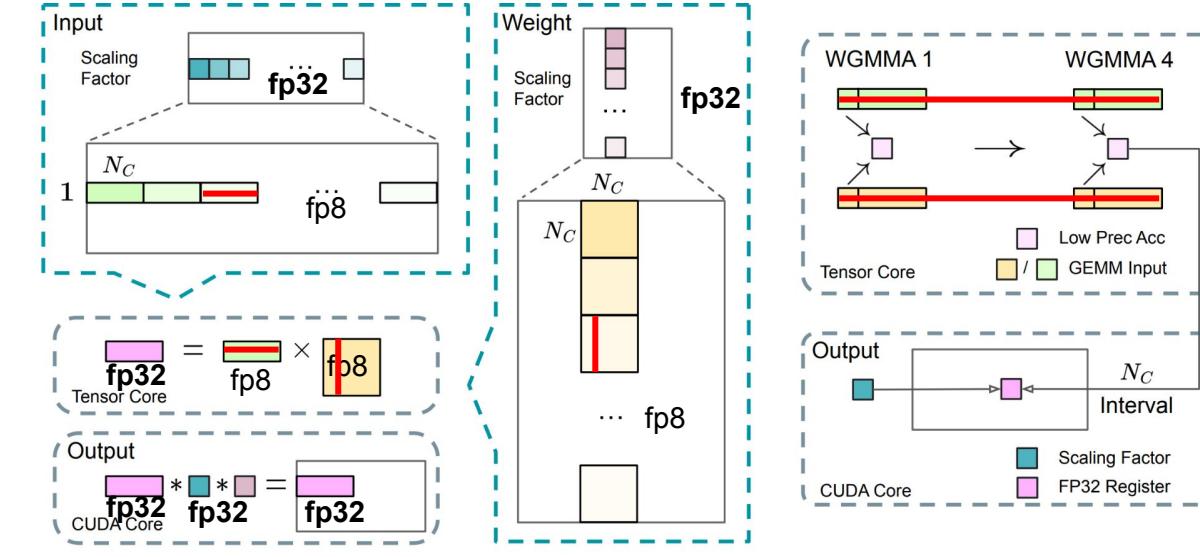


Figure 7 | (a) We propose a fine-grained quantization method to mitigate quantization errors caused by feature outliers; for illustration simplicity, only Fprop is illustrated. (b) In conjunction with our quantization strategy, we improve the FP8 GEMM precision by promoting to CUDA Cores at an interval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

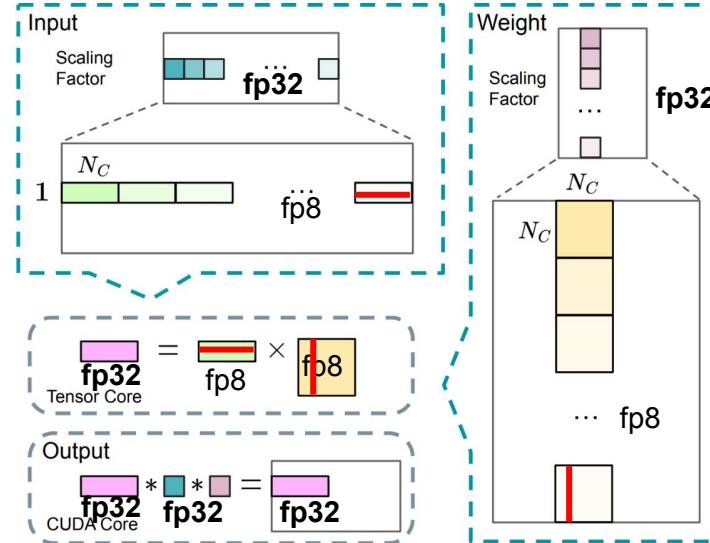
DeepSeek: fine-grained fp8 quantization



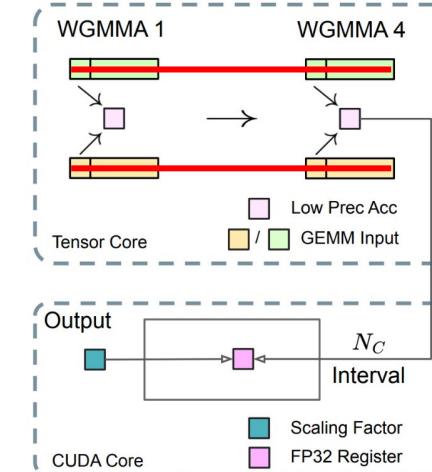
fp32



fp8

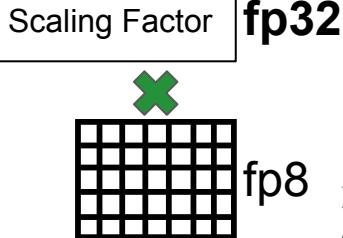
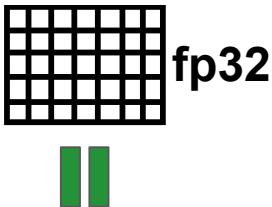
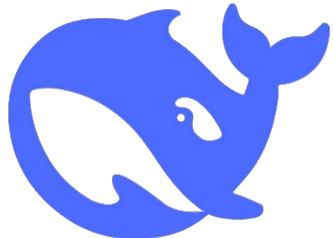


(a) Fine-grained quantization

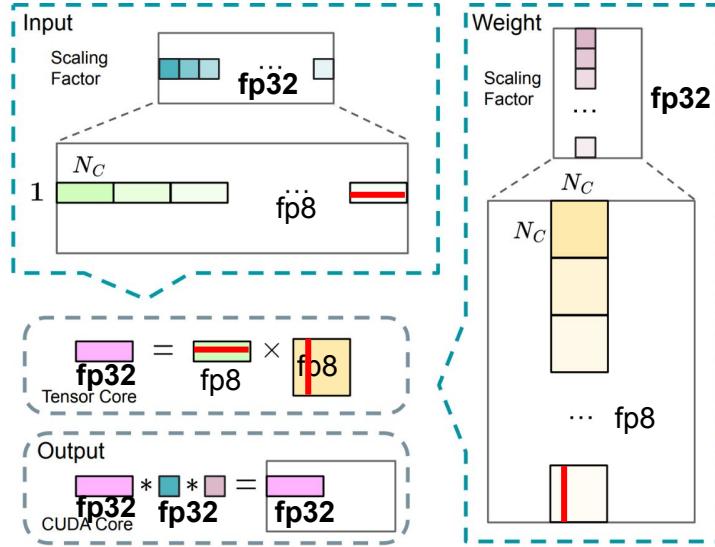


(b) Increasing accumulation precision

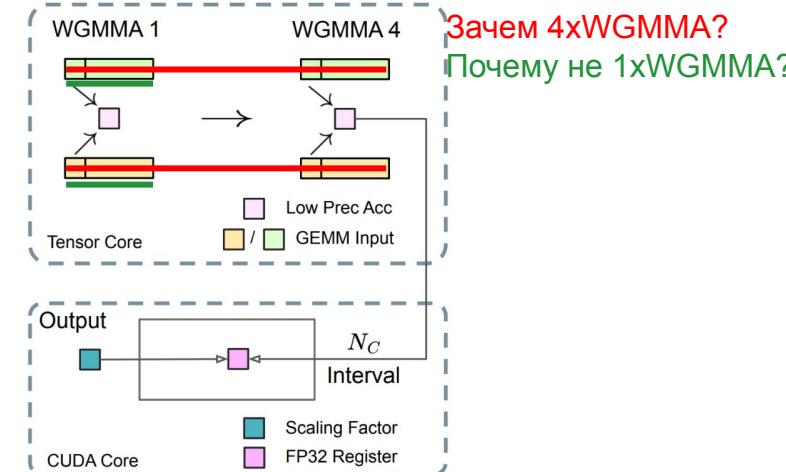
Figure 7 | (a) We propose a fine-grained quantization method to mitigate quantization errors caused by feature outliers; for illustration simplicity, only Fprop is illustrated. (b) In conjunction with our quantization strategy, we improve the FP8 GEMM precision by promoting to CUDA Cores at an interval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



DeepSeek: fine-grained fp8 quantization

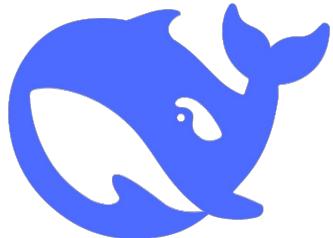


(a) Fine-grained quantization



(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.

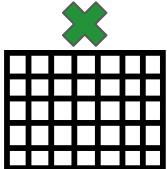


fp32



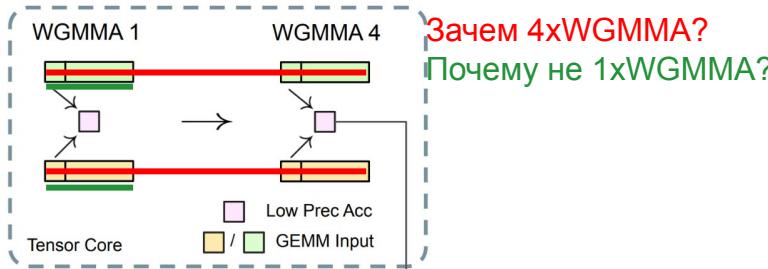
Scaling Factor

fp32

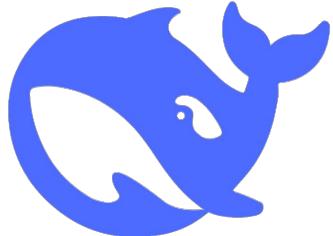


fp8

DeepSeek: fine-grained fp8 quantization



It is worth noting that this modification reduces the WGMMA (Warpgroup-level Matrix Multiply-Accumulate) instruction issue rate for a single warpgroup. However, on the H800 architecture, it is typical for two WGMMA 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 WGMAs, represents the minimal accumulation interval that can significantly improve precision without introducing substantial overhead.

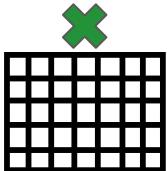


fp32



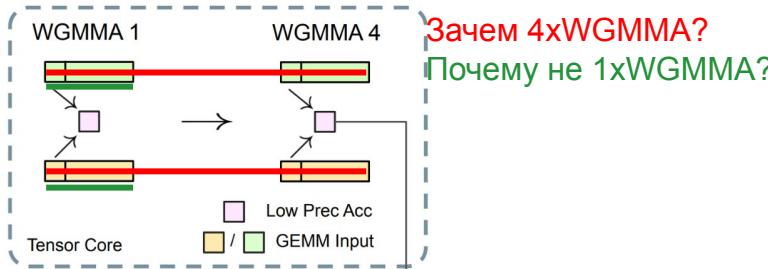
Scaling Factor

fp32



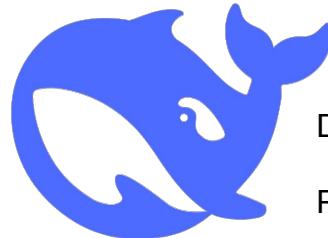
fp8

DeepSeek: fine-grained fp8 quantization



It is worth noting that this modification reduces the WGMMA (Warpgroup-level Matrix Multiply-Accumulate) instruction issue rate for a single warpgroup. However, on the H800 architecture, it is typical for two WGMMA 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 WGMAs, represents the minimal accumulation interval that can significantly improve precision without introducing substantial overhead.

Не знаю почему нужно делать 4xWGMMA + PROMOTION
вместо 4x(WGMMA + 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 WGMMA on NVIDIA® Hopper™ GPUs](#)

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

Глава 5: Умножение матриц

Метод Штассена

Важнейшим оператором в языке Python является умножение матриц. Для этого есть специальный класс `Matrix`, который определяет метод `__mul__`. Важно отметить, что умножение матриц не коммутативно, то есть $A \cdot B \neq B \cdot A$.

Метод `__mul__` имеет следующий формат:

```
def __mul__(self, other):  
    if type(other) == Matrix:  
        return Matrix._mult(self, other)  
    else:  
        return self._mult_by_number(other)
```

Метод `_mult` вычисляет произведение двух матриц, а метод `_mult_by_number` — произведение матрицы на скалярное значение. Метод `_mult` имеет следующий формат:

```
def _mult(self, other):  
    if type(other) == Matrix:  
        return Matrix._mult(self, other)  
    else:  
        return self._mult_by_number(other)
```

Метод `_mult_by_number` вычисляет произведение матрицы на скалярное значение. Метод `_mult` имеет следующий формат:

```
def _mult(self, other):  
    if type(other) == Matrix:  
        return Matrix._mult(self, other)  
    else:  
        return self._mult_by_number(other)
```

Метод `_mult_by_number` вычисляет произведение матрицы на скалярное значение. Метод `_mult` имеет следующий формат:

```
def _mult(self, other):  
    if type(other) == Matrix:  
        return Matrix._mult(self, other)  
    else:  
        return self._mult_by_number(other)
```

Метод `_mult` вычисляет произведение двух матриц, а метод `_mult_by_number` — произведение матрицы на скалярное значение. Метод `_mult` имеет следующий формат:

```
def _mult(self, other):  
    if type(other) == Matrix:  
        return Matrix._mult(self, other)  
    else:  
        return self._mult_by_number(other)
```

Метод `_mult` вычисляет произведение двух матриц, а метод `_mult_by_number` — произведение матрицы на скалярное значение. Метод `_mult` имеет следующий формат:

```
def _mult(self, other):  
    if type(other) == Matrix:  
        return Matrix._mult(self, other)  
    else:  
        return self._mult_by_number(other)
```

Метод `_mult` вычисляет произведение двух матриц, а метод `_mult_by_number` — произведение матрицы на скалярное значение. Метод `_mult` имеет следующий формат:

```
def _mult(self, other):  
    if type(other) == Matrix:  
        return Matrix._mult(self, other)  
    else:  
        return self._mult_by_number(other)
```

Метод `_mult` вычисляет произведение двух матриц, а метод `_mult_by_number` — произведение матрицы на скалярное значение. Метод `_mult` имеет следующий формат:

```
def _mult(self, other):  
    if type(other) == Matrix:  
        return Matrix._mult(self, other)  
    else:  
        return self._mult_by_number(other)
```

Умножение матриц - Метод Штрассена

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} \times \begin{bmatrix} e & f \\ g & h \end{bmatrix} = \begin{bmatrix} ae + bg & af + bh \\ ce + dg & cf + dh \end{bmatrix}$$

A B C

Умножение матриц - Метод Штрассена

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} \times \begin{bmatrix} e & f \\ g & h \end{bmatrix} = \begin{bmatrix} ae + bg & af + bh \\ ce + dg & cf + dh \end{bmatrix}$$

A B C

$p_1 = a(f - h)$
 $p_3 = (c + d)e$
 $p_5 = (a + d)(e + h)$
 $p_7 = (a - c)(e + f)$

$$p_2 = (a + b)h$$

$$p_4 = d(g - e)$$

$$p_6 = (b - d)(g + h)$$

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} \times \begin{bmatrix} e & f \\ g & h \end{bmatrix} = \begin{bmatrix} p_5 + p_4 - p_2 + p_6 & p_1 + p_2 \\ p_3 + p_4 & p_1 + p_5 - p_3 - p_7 \end{bmatrix}$$

A B C

Умножение матриц - Метод Штрассена

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} \times \begin{bmatrix} e & f \\ g & h \end{bmatrix} = \begin{bmatrix} ae + bg & af + bh \\ ce + dg & cf + dh \end{bmatrix}$$

A B C

$p_1 = a(f - h)$
 $p_3 = (c + d)e$
 $p_5 = (a + d)(e + h)$
 $p_7 = (a - c)(e + f)$

$p_2 = (a + b)h$
 $p_4 = d(g - e)$
 $p_6 = (b - d)(g + h)$

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} \times \begin{bmatrix} e & f \\ g & h \end{bmatrix} = \begin{bmatrix} p_5 + p_4 - p_2 + p_6 & p_1 + p_2 \\ p_3 + p_4 & p_1 + p_5 - p_3 - p_7 \end{bmatrix}$$

A B C

Получили рекуррентную асимптотику $T(N)=7T(N/2)+8O(N^2) \Rightarrow T(N)=O(N^{\log 7})=O(N^{2.8})$

Умножение матриц - Метод Штрассена

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} \times \begin{bmatrix} e & f \\ g & h \end{bmatrix} = \begin{bmatrix} ae + bg & af + bh \\ ce + dg & cf + dh \end{bmatrix}$$

A B C

$p_1 = a(f - h)$
 $p_3 = (c + d)e$
 $p_5 = (a + d)(e + h)$
 $p_7 = (a - c)(e + f)$

$p_2 = (a + b)h$
 $p_4 = d(g - e)$
 $p_6 = (b - d)(g + h)$

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} \times \begin{bmatrix} e & f \\ g & h \end{bmatrix} = \begin{bmatrix} p_5 + p_4 - p_2 + p_6 & p_1 + p_2 \\ p_3 + p_4 & p_1 + p_5 - p_3 - p_7 \end{bmatrix}$$

A B C

Получили рекуррентную асимптотику $T(N)=7T(N/2)+8O(N^2) \Rightarrow T(N)=O(N^{\log 7})=O(N^{2.8})$

Метод Виноградова - тоже $O(N^{2.8})$

<https://www.geeksforgeeks.org/strassens-matrix-multiplication/>



ВОПРОСЫ?



[@UnicornGlade](https://t.me/UnicornGlade)

[@PolarNick239](https://twitter.com/PolarNick239)

polarnick239@gmail.com



Николай Полярный