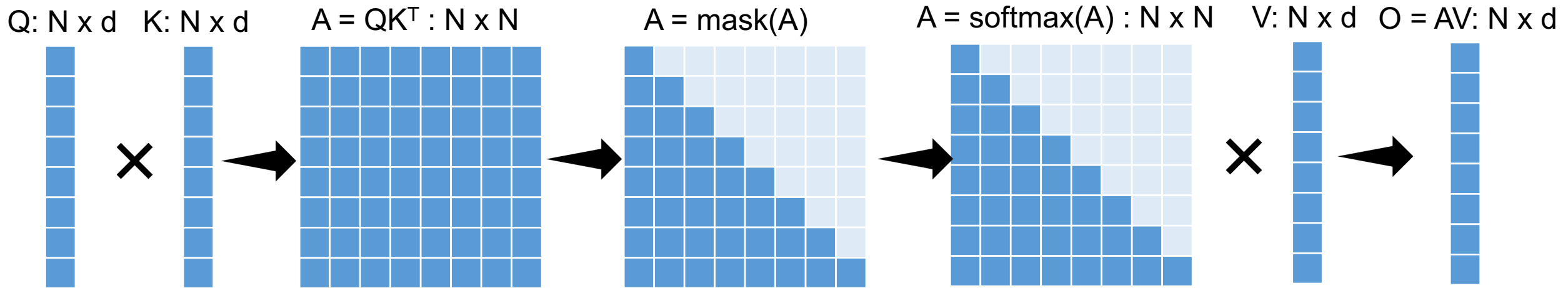# 15-442/15-642: Machine Learning Systems

# Attention Optimizations

**Tianqi Chen and Zhihao Jia**

Carnegie Mellon University

**3/12/24**

# Attention: $O = \text{Softmax}(QK^T) \, V$

Q: N x d    K: N x d    $A = QK^T$ : N x N    A = mask(A)    A = softmax(A) : N x N    V: N x d    O = AV: N x d



## Challenges:

- Large intermediate results

- Repeated reads/writes from GPU device memory

- Cannot scale to long sequences due to O(N^2) intermediate results

# Outline: Attention Optimizations
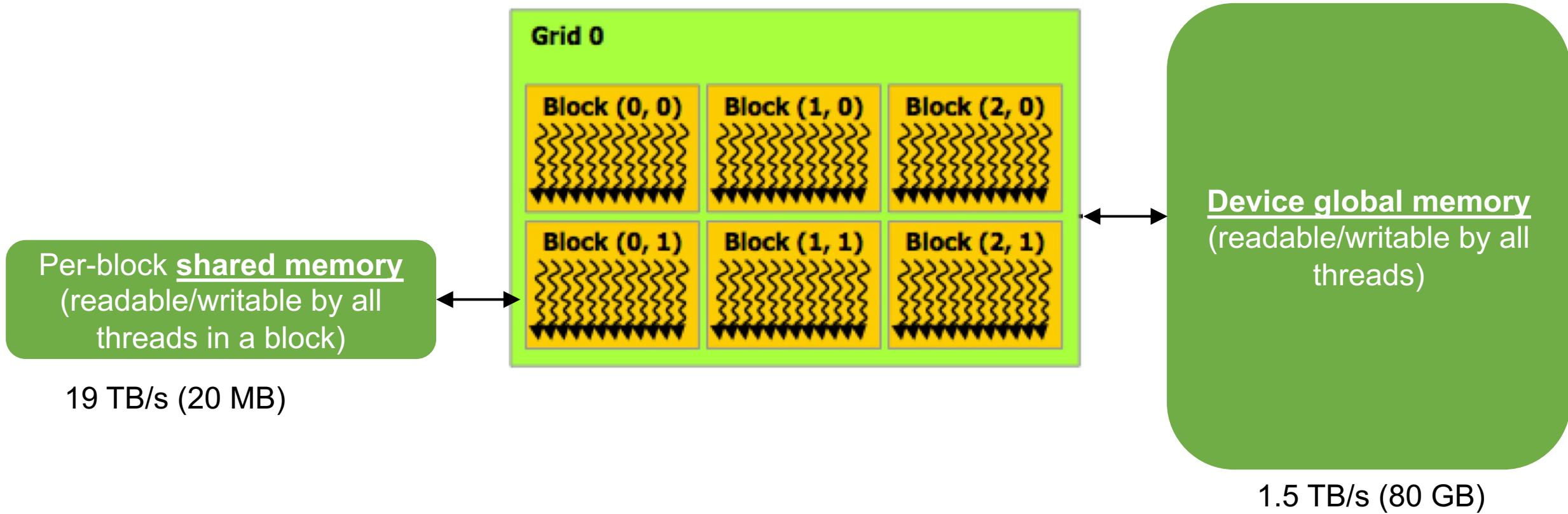
Part 1: LLM Training

- FlashAttention

Part 2: LLM Inference

- Flash Decoding

- PagedAttention

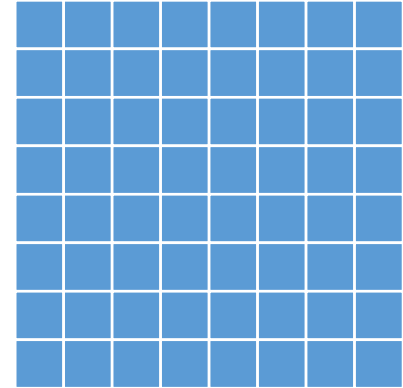**These techniques are highly tailored for GPUs**

# Revisit: GPU Memory Hierarchy

**Grid 0**

Block (0, 0)  Block (1, 0)  Block (2, 0)

Block (0, 1)  Block (1, 1)  Block (2, 1)

Per-block **shared memory**
(readable/writable by all
threads in a block)

19 TB/s (20 MB)

**Device global memory**
(readable/writable by all
threads)

1.5 TB/s (80 GB)

4

# FlashAttention

$A = \text{softmax}(QK^T)$

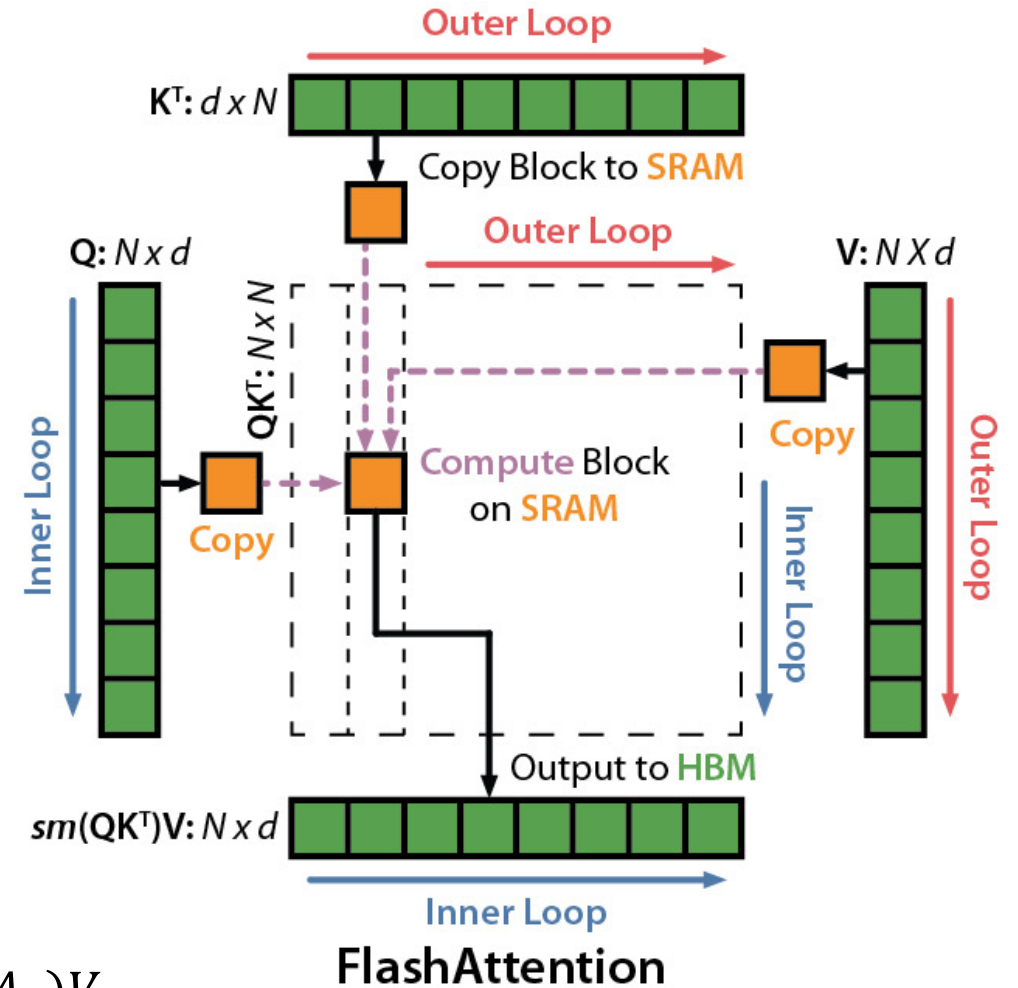**Key idea**: compute attention by blocks to reduce global memory access

**Two main Techniques:**

**1. Tiling:** restructure algorithm to load query/key/value block by block from global to shared memory

**2. Recomputation:** don't store attention matrix from forward, recompute it in backward

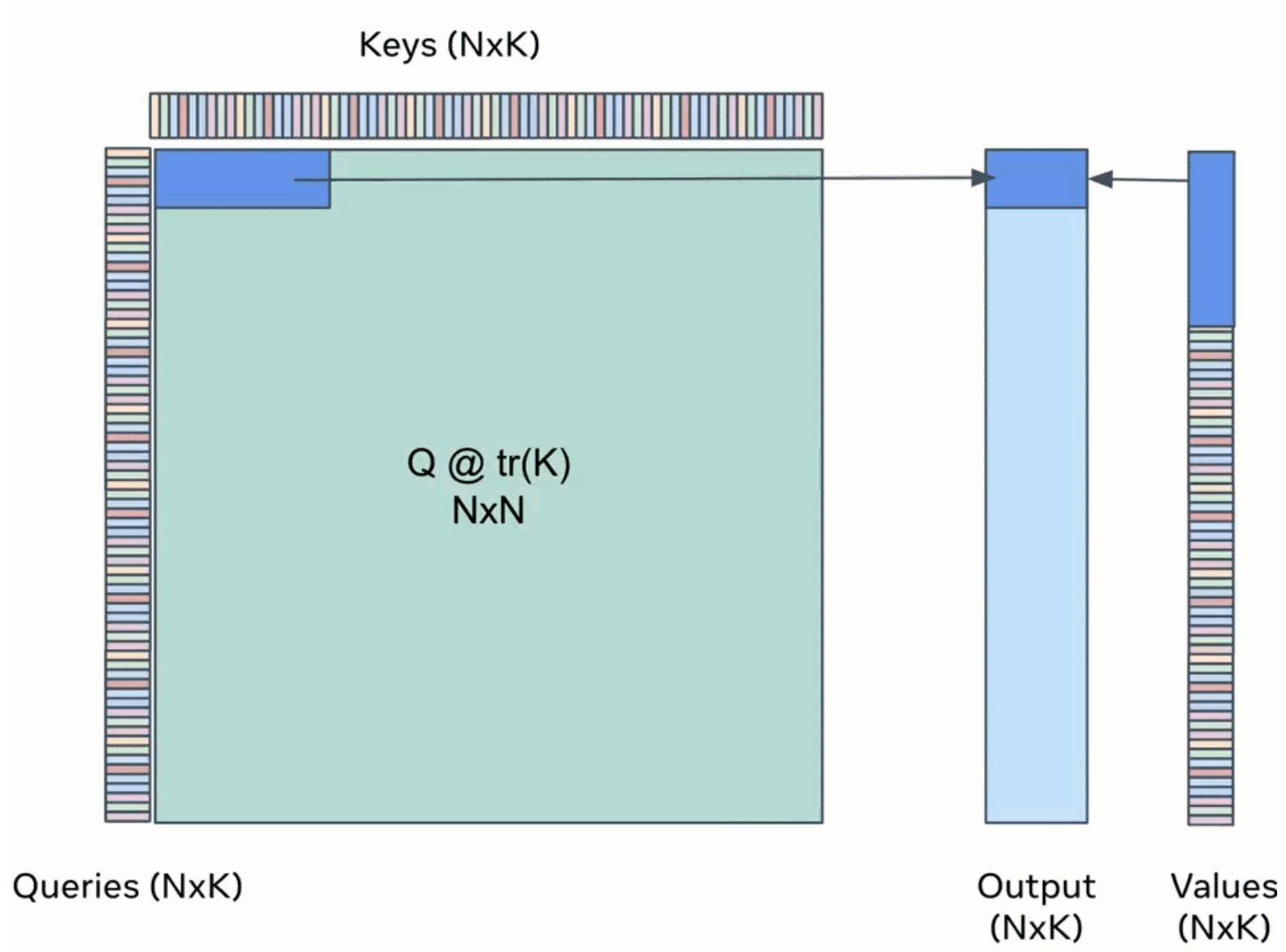# Tiling: Decompose Large Softmax into smaller ones by Scaling

1. Load inputs by blocks from global to shared memory

2. On chip, compute attention output wrt the block

3. Update output in device memory by scaling



**FlashAttention**

$$softmax([A_1, A_2]) = [\alpha \times softmax(A_1), \beta \times softmax(A_2)]$$

$$softmax([A_1, A_2]) \begin{bmatrix} V_1 \\ V_2 \end{bmatrix} = \alpha \times softmax(A_1)V_1 + \beta \times softmax(A_2)V_2$$
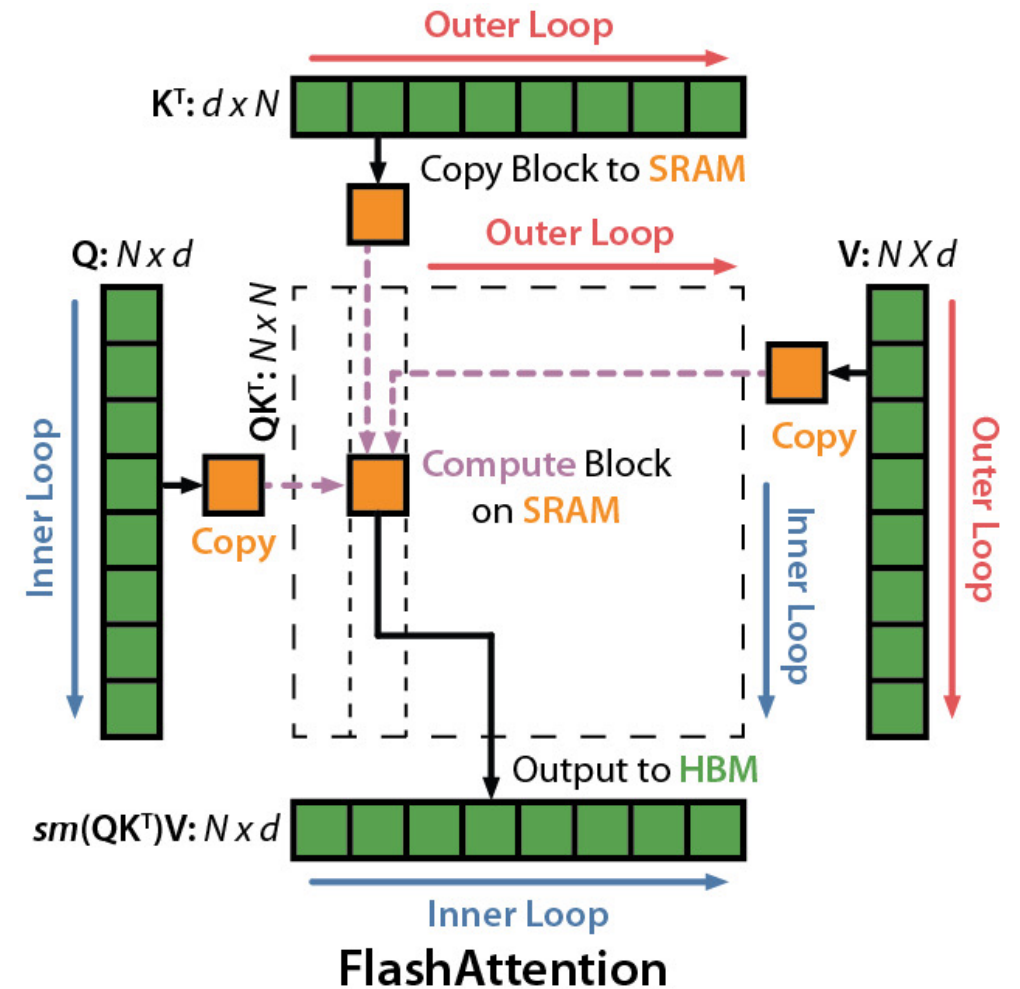
# Tiling

# Recomputation: Backward Pass

By storing softmax normalization factors from forward (size N), recompute attention in the backward from inputs in shared memory



FlashAttention

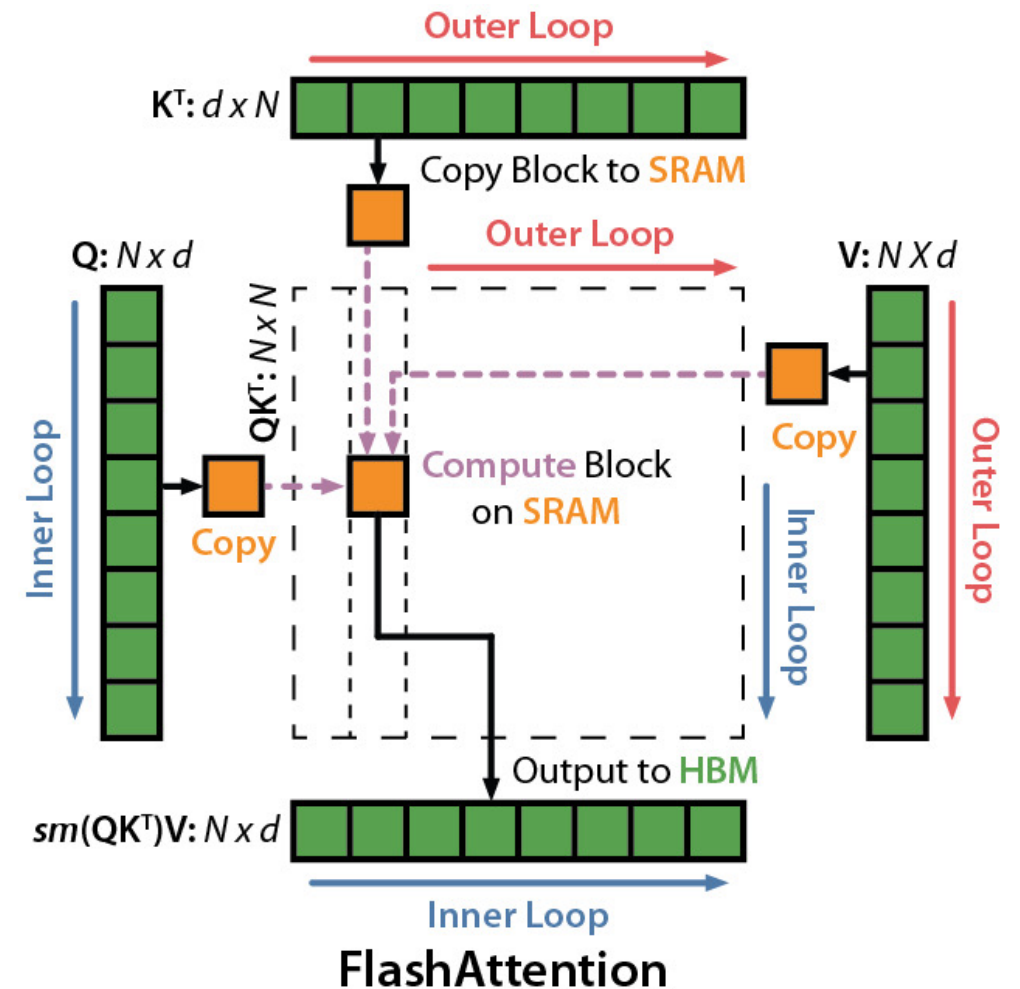| Attention | Standard | FlashAttention |
|---|---|---|
| GFLOPs | 66.6 | 75.2 |
| Global mem access | 40.3 GB | 4.4 GB |
| Runtime | 41.7 ms | 7.3 ms |

**Speed up backward pass with increased FLOPs**

# FlashAttention: Threadblock-level Parallelism

How to partition FlasshAttention across thread blocks?

(An A100 has 108 SMMs -> 108 thread blocks)

- Step 1: assign different heads to different thread blocks (16-64 heads)



FlashAttention

# FlashAttention: Threadblock-level Parallelism

How to partition FlasshAttention across thread blocks?

(An A100 has 108 SMMs -> 108 thread blocks)

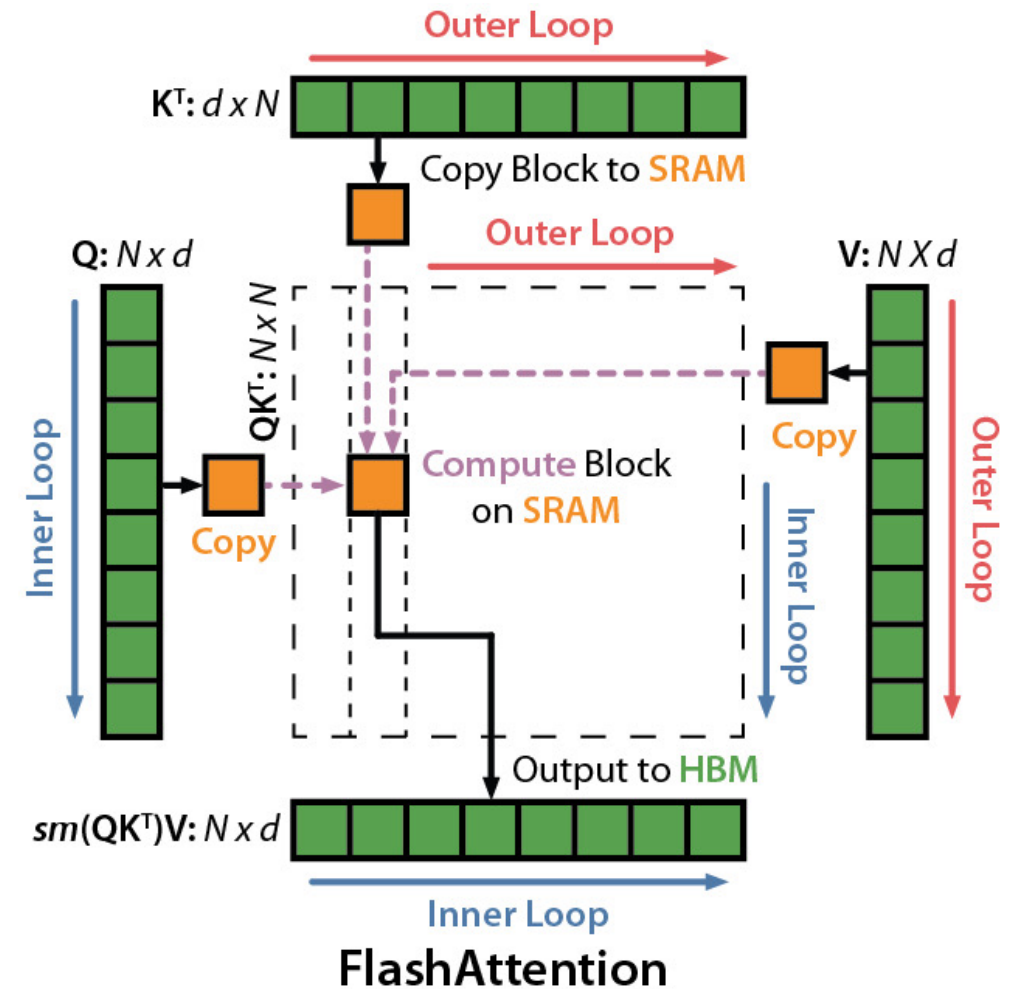- Step 1: assign different heads to different thread blocks (16-64 heads)
- Step 2: assign different queries to different thread blocks (Why?)

**Thread blocks cannot communicate; cannot perform softmax when partitioning keys/values**



FlashAttention

# FlashAttention: Threadblock-level Parallelism

Keys/Values

Queries



Block 1
Block 2
Block 3
Block 4
Block 5

Forward pass

**Do we need to handle workload imbalance?**

**No. GPU scheduler automatically loads the next block once the current one completes.**

# FlashAttention: Warp-Level Parallelism

- How to partition FlashAttention across warps within a thread block?



(a) FLASHATTENTION

(b) FLASHATTENTION-2

Splitting across K/V requires communication to add results

Splitting across Q avoids communications

# FlashAttention: 2-4x speedup, 10-20x memory reduction



Memory linear in sequence length

# Outline: Attention Optimizastions

Part 1: LLM Training

- FlashAttention

**Part 2: LLM Inference (Auto-regressive Decoding)**

- Flash Decoding

- PagedAttention

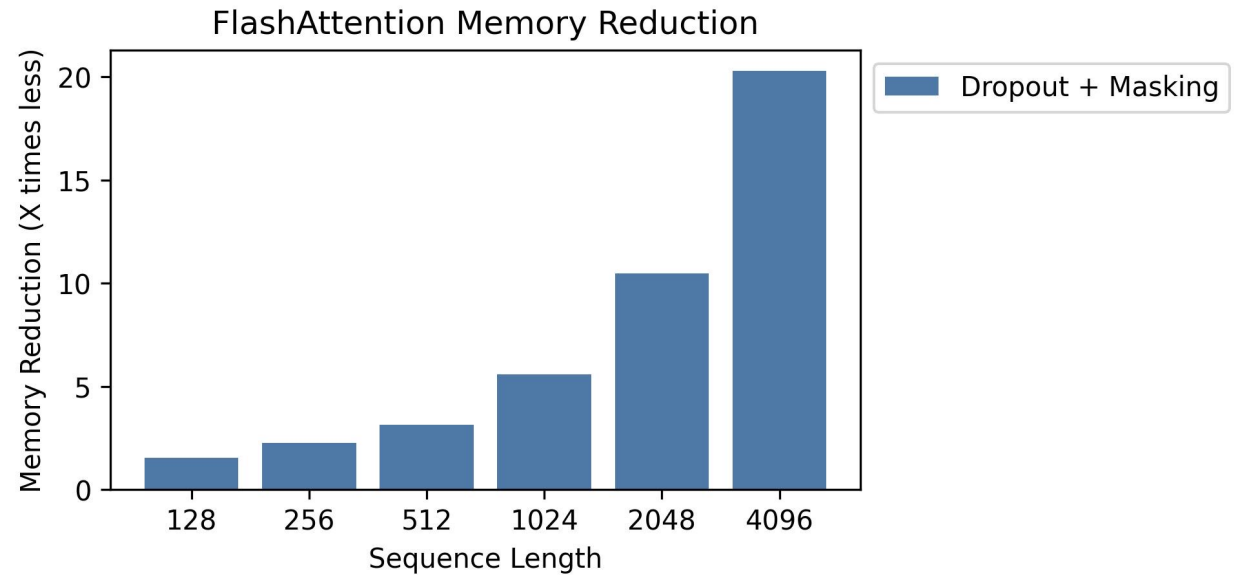# Generative LLM Inference: Autoregressive Decoding

Input Prompt:     [Accelerating LLM requires machine]    → learning    → systems    → optimizations

|                | Iter 0 | Iter 1 | Iter 2 | Iter 3 |
|----------------|--------|--------|--------|--------|
|                | Layer 1 | Layer 1 | Layer 1 | Layer 1 |
|                | Layer 2 | Layer 2 | Layer 2 | Layer 2 |
|                | Layer 3 | Layer 3 | Layer 3 | Layer 3 |
| Outputs:       | learning | systems | optimizations | [EOS] |

# Generative LLM Inference: Autoregressive Decoding



**Attention Score**

|         | Acc. | LLM | requires | machine |
|---------|------|-----|----------|---------|
| Acc.    | 1    |     |          |         |
| LLM     | 2    | 0   |          |         |
| requires| 5    | 1   | 3        |         |
| machine | 2    | 0   | 1        | 1       |

**Pre-filling Phase**

[Acc. LLM requires machine]    learning    systems    optimizations

Iter 0    Iter 1    Iter 2    Iter 3

Layer 1    Layer 1    Layer 1    Layer 1

Layer 2    Layer 2    Layer 2    Layer 2

Layer 3    Layer 3    Layer 3    Layer 3

Outputs:    learning    systems    optimizations    [EOS]

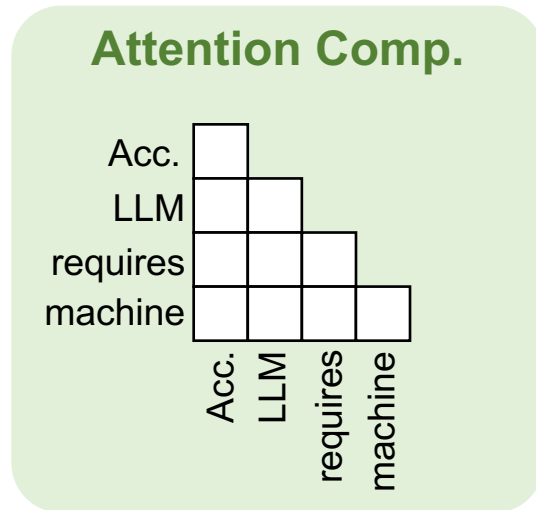# Generative LLM Inference: Autoregressive Decoding

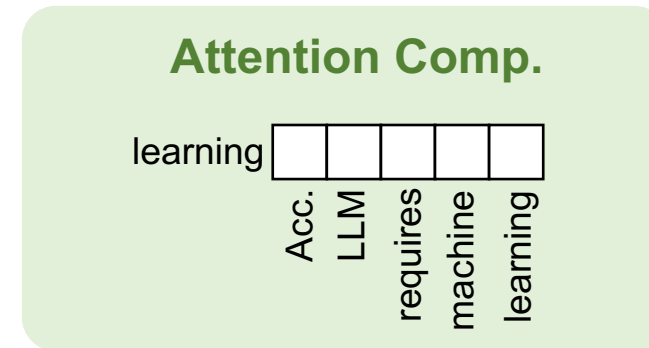# Generative LLM Inference: Autoregressive Decoding

- **Pre-filling phase** (0-th iteration):
  - Process *all* input tokens at once

- **Decoding phase** (all other iterations):
  - Process a *single* token generated from previous iteration
  - Use attention keys & values of all previous tokens

- Key-value cache:
  - Save attention keys and values for the following iterations to avoid recomputation

# Can We Apply FlashAttention to LLM Inference?

**Attention Comp.**



**Attention Comp.**



**Pre-filling phase:**

- Yes, compute different queries using different thread blocks/warps

**Decoding phase:**

- No, there is a single query in the decoding phase

# FlashAttention Processes K/V Sequentially



**Inefficient for requests with long context (many keys/values)**

# Flash-Decoding Parallelizes Across Keys/Values

1. Split keys/values into small chunks

2. Compute attention with these splits using FlashAttention

3. Reduce overall all splits



**Key insight: attention is associative and commutative**

# Flash-Decoding is up to 8x faster than prior work



CodeLlama-34b end-to-end decoding speed [bs=1, MP=4]

# Outline: Attention Optimizastions

Part 1: LLM Training

• FlashAttention

Part 2: LLM Inference (Auto-regressive Decoding)

• Flash-Decoding
• **PagedAttention**

# KV Cache Dynamically Grows and Shrinks

[Accelerating LLM requires machine]

**Attention Matrix**



Iter 0

Layer 1

Layer 2

Layer 3

Outputs: learning

**KV Cache**

Accelerating | LLM | requires | machine

# KV Cache Dynamically Grows and Shrinks

[Accelerating LLM requires machine] ⇢ learning

**Attention Matrix**

learning

Acc. | LLM | requires | machine | learning

Iter 0

Layer 1
Layer 2
Layer 3

Iter 1

Layer 1
Layer 2
Layer 3

Outputs: learning    systems

**KV Cache**

Accelerating | LLM | requires | machine | learning

# KV Cache Dynamically Grows and Shrinks

# KV Cache Dynamically Grows and Shrinks

# Static KV Cache Management Wastes Memory



| 0 | | 3 | | | | | A's max length | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| *Artificial* | *Intelligence* | *is* | *\<resv\>* | *\<resv\>* | *…* | *\<resv\>* | *\<resv\>* | *…* | *…* | *Alan* | *Turing* | *…* |

3 KV Cache slots for request A's prompt | Pre-allocated slots for A's output (*Internal frag.*) | *External frag.* | Request B

- **Pre-allocates contiguous** space of memory to the request's maximum length

- Memory fragmentation
  - **Internal fragmentation** due to unknown output length
  - **External fragmentation** due to non-uniform per-request max lengths

# Significant Memory Waste in KV Cache

• Only 20-40% of KV cache is utilized to store actual token states



slides from vllm: Efficient Memory Management for Large Language Model Serving with PagedAttention

# PagedAttention

- Application-level memory paging and virtualization for KV cache

**Memory management in OS**     **PagedAttention**



Physical Memory          KV Cache

# Paging KV Cache Space into KV Blocks*

- KV block is a **fixed-size** contiguous chunk of memory that stores KV states from **left to right**

KV blocks

|  |  |  |  |  |
|---|---|---|---|---|
| block 0 |  |  |  |  |
| block 1 |  |  |  |  |
| block 2 |  |  |  |  |
| block 3 |  | KV Cache |  |  |
| block 4 | Artificial | Intelligence | is | the |
| block 5 |  | Space |  |  |
| block 6 |  |  |  |  |
| block 7 |  |  |  |  |

Block size = 4

\* The term ``block'' is overloaded in PagedAttention

# Virtualizing KV Cache

Request
A

Prompt: "Alan Turing is a computer scientist"

**Logical** KV blocks

| | | | | |
|---|---|---|---|---|
| block 0 | Alan | Turing | is | a |
| block 1 | computer | scientist | | |
| block 2 | | | | |
| block 3 | | | | |

**Block table**

| Physical block number | # Filled |
|---|---|
| 7 | 4 |
| 1 | 2 |
| – | – |
| – | – |

**Physical** KV blocks

| | | | | |
|---|---|---|---|---|
| block 0 | | | | |
| block 1 | computer | scientist | | |
| block 2 | | | | |
| block 3 | | | | |
| block 4 | | | | |
| block 5 | | | | |
| block 6 | | | | |
| block 7 | Alan | Turing | is | a |

# Attention with Virtualized KV Cache

1. Fetch non-contiguous KV blocks using the block table
2. Apply attention on the fly



**Key insight: attention is associative and commutative**

# Memory Management with PagedAttention

Request
A

Prompt: "Alan Turing is a computer scientist"
Completion: "_and_"

**Physical** KV blocks

| block 0 | | | | |
|---|---|---|---|---|
| block 1 | computer | scientist | | |
| block 2 | | | | |
| block 3 | | | | |
| block 4 | | | | |
| block 5 | | | | |
| block 6 | | | | |
| block 7 | Alan | Turing | is | a |

**Logical** KV blocks

| block 0 | Alan | Turing | is | a |
|---|---|---|---|---|
| block 1 | computer | scientist | | |
| block 2 | | | | |
| block 3 | | | | |

**Block table**

| Physical block number | # Filled |
|---|---|
| 7 | 4 |
| 1 | 2 |
| – | – |
| – | – |

# Memory Management with PagedAttention

Request A

Prompt: "Alan Turing is a computer scientist"
Completion: "__and__"

**Logical** KV blocks

| | | | | |
|---|---|---|---|---|
| block 0 | Alan | Turing | is | a |
| block 1 | computer | scientist | and | |
| block 2 | | | | |
| block 3 | | | | |

**Block table**

| Physical block number | # Filled |
|---|---|
| 7 | 4 |
| 1 | 2 |
| – | – |
| – | – |

**Physical** KV blocks

| | | | | |
|---|---|---|---|---|
| block 0 | | | | |
| block 1 | computer | scientist | | |
| block 2 | | | | |
| block 3 | | | | |
| block 4 | | | | |
| block 5 | | | | |
| block 6 | | | | |
| block 7 | Alan | Turing | is | a |

# Memory Management with PagedAttention

Request A

Prompt: "Alan Turing is a computer scientist"
Completion: "_and_"

**Logical** KV blocks

| | | | |
|---|---|---|---|
| block 0 | Alan | Turing | is | a |
| block 1 | computer | scientist | **and** | |
| block 2 | | | | |
| block 3 | | | | |

**Block table**

| Physical block number | # Filled |
|---|---|
| 7 | 4 |
| 1 | **3** |
| – | – |
| – | – |

**Physical** KV blocks

| | | | | |
|---|---|---|---|---|
| block 0 | | | | |
| block 1 | computer | scientist | **and** | |
| block 2 | | | | |
| block 3 | | | | |
| block 4 | | | | |
| block 5 | | | | |
| block 6 | | | | |
| block 7 | Alan | Turing | is | a |

# Memory Management with PagedAttention

Request A

Prompt: "Alan Turing is a computer scientist"
Completion: "and mathematician"

**Logical** KV blocks

| | | | |
|---|---|---|---|
| block 0 | Alan | Turing | is | a |
| block 1 | computer | scientist | and | **mathem atician** |
| block 2 | | | | |
| block 3 | | | | |

**Block table**

| Physical block number | # Filled |
|---|---|
| 7 | 4 |
| 1 | **3** |
| – | – |
| – | – |

**Physical** KV blocks

| | | | | |
|---|---|---|---|---|
| block 0 | | | | |
| block 1 | computer | scientist | and | **mathem atician** |
| block 2 | | | | |
| block 3 | | | | |
| block 4 | | | | |
| block 5 | | | | |
| block 6 | | | | |
| block 7 | Alan | Turing | is | a |

slides from vllm: Efficient Memory Management for Large Language Model Serving with PagedAttention

# Memory Management with PagedAttention
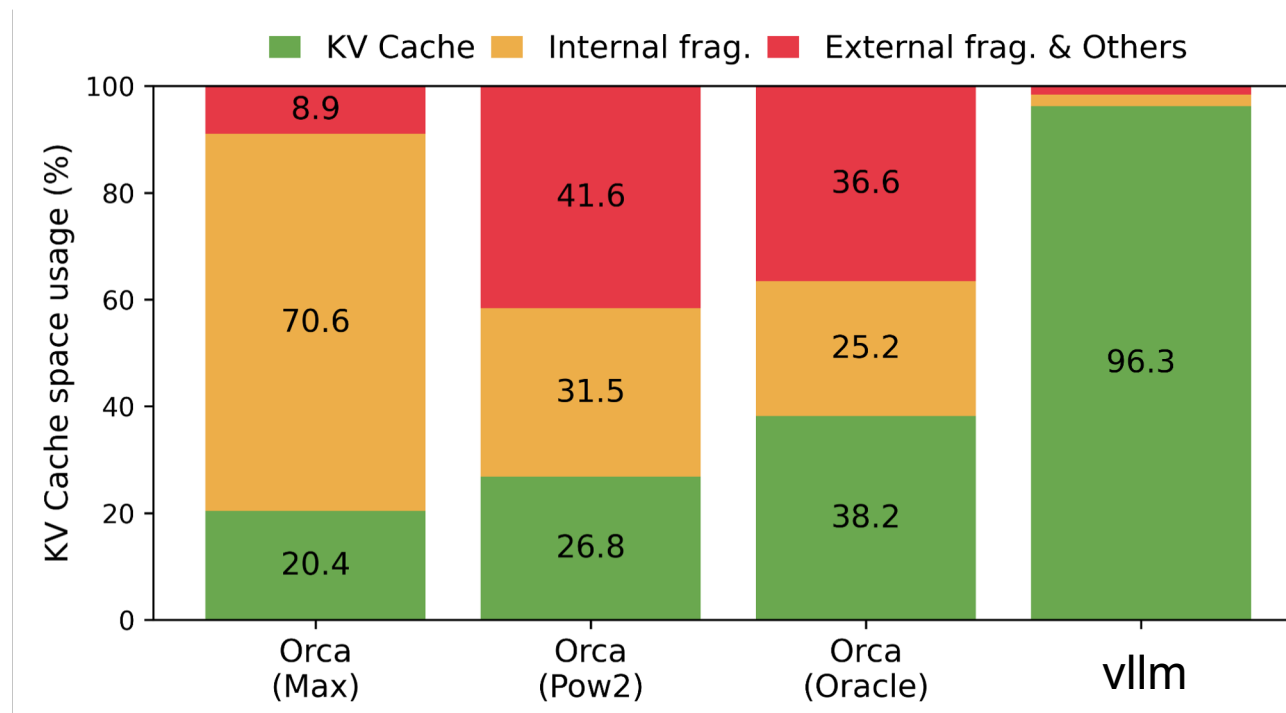
# Memory Efficiency of PagedAttention

**Minimal internal fragmentation**

- Only happens at the last block of a sequence

- # wasted tokens / seq < block size

**No external fragmentation**

| Alan | Turing | is | a |
|---|---|---|---|
| computer | scientist | and | mathematician |
| renowned | | | |

Internal fragmentation



39

# Recap: Techniques for Optimizing Attention

- **FlashAttention**: tiling to reduce GPU global memory access

- **Auto-regressive Decoding**: pre-filling and decoding phases, KV cache

- **FlashDecoding**: improving attention's parallelism by splitting keys/values

- **PagedAttention**: paging and virtualization to reduce KV cache's memory requirement