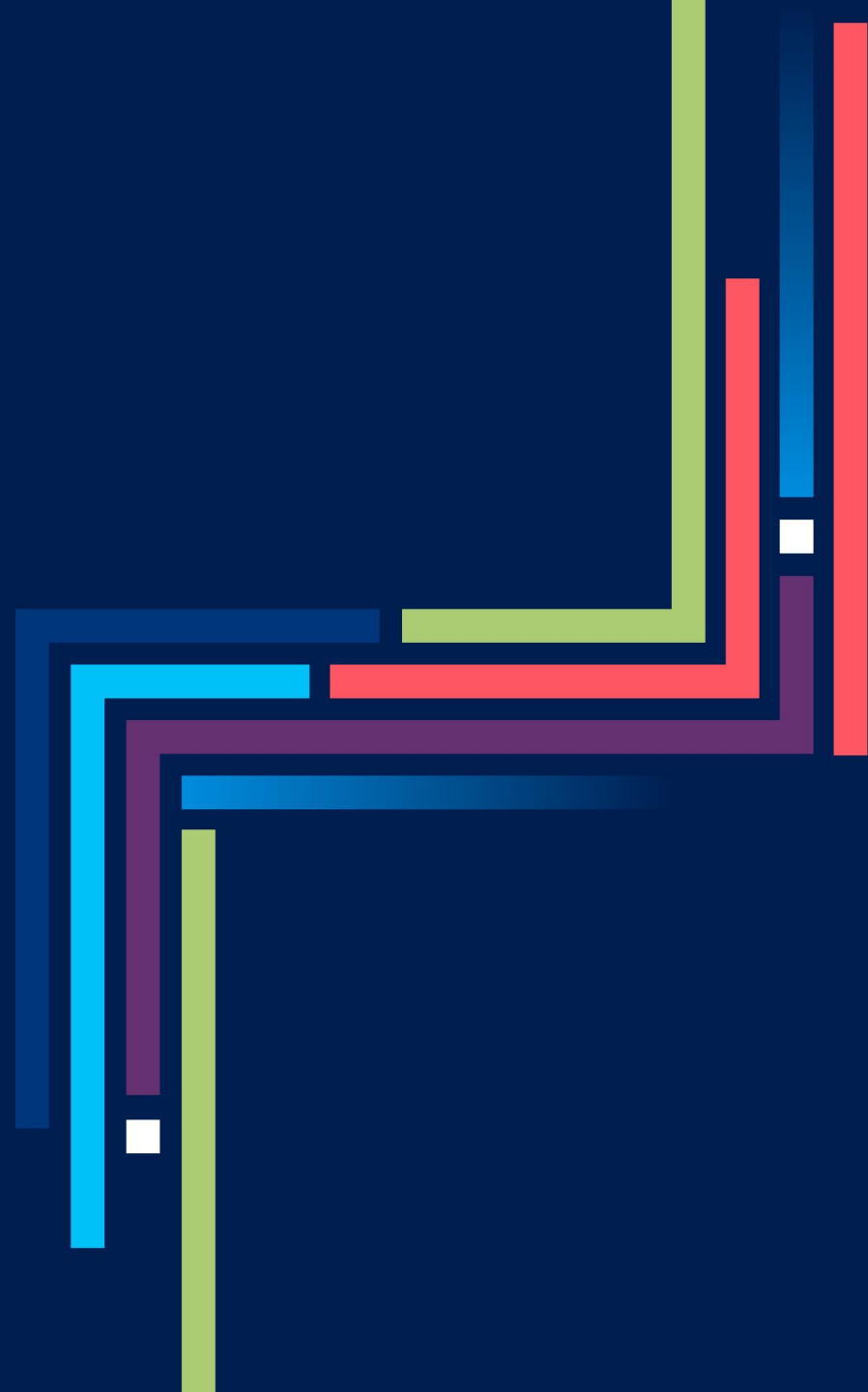


# 大语言模型时代：最大化CPU价值的优化策略

何普江

英特尔AI软件架构师



# 目录

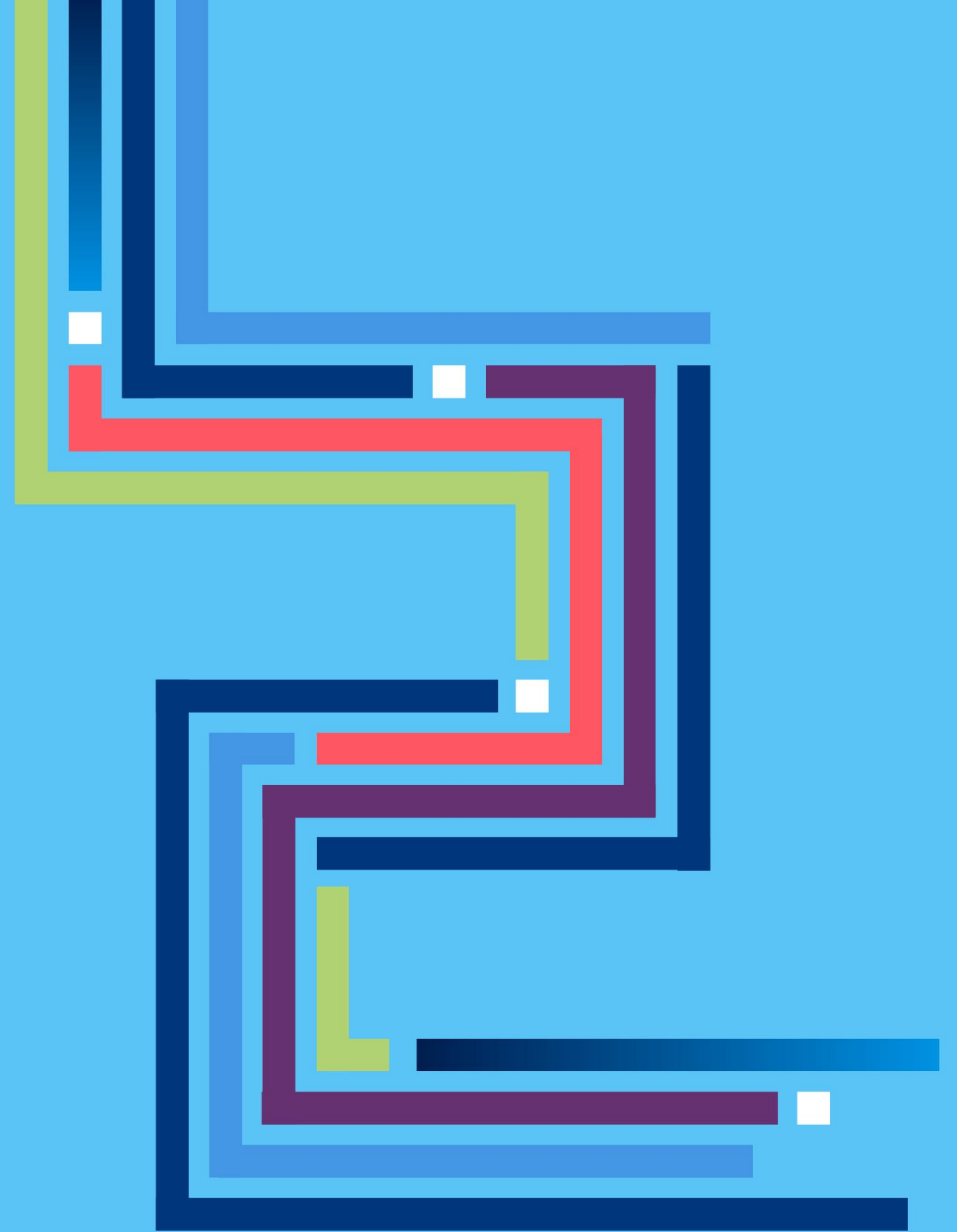
01 背景（为什么？）

02 CPU上如何优化大语言模型？

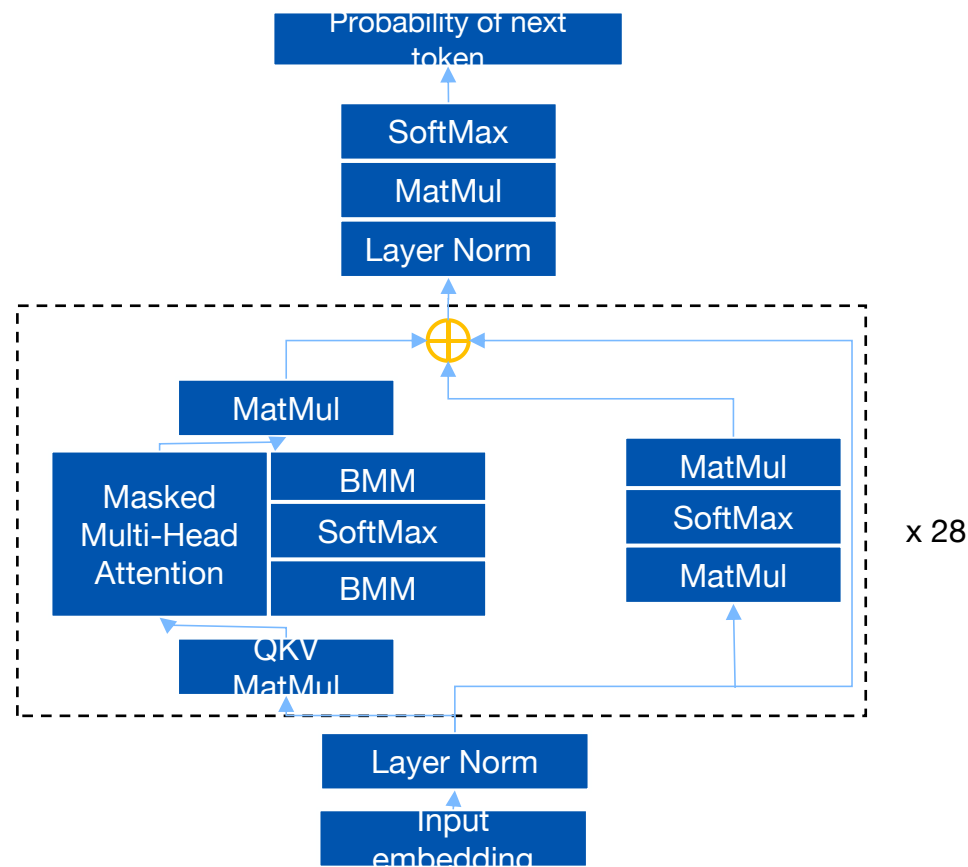
03 最大化CPU价值

04 总结

背景 (为什么考虑最大化CPU  
价值?)



# Computing Needs in LLM



GPT-J Model Structure

MatMul shapes in GPT-J  
(suppose prompt token size = 2048, batch size=1, greedy search)

	1 <sup>st</sup> token	Next tokens
QKV MatMul in MHA	A: 2048x4096 B: 4096x12288	A: 1x4096 B: 4096x12288
MHA (1st BMM)	A: 16x2048x256 B: 16x2048x256	A: 16x1x256 B: 16x2048x256
Output MatMul in MHA	A: 2048x4096 B: 4096x4096	A: 1x4096 B: 4096x4096
1st MatMul in FFN	A: 2048x4096 B: 4096x16384	A: 1x4096 B: 4096x16384
2nd MatMul in FFN	A: 2048x16384 B: 4096x4096	A: 1x16384 B: 4096x4096



Compute Bound



Memory Read  
Bandwidth Bound

# GPT Series Model Analysis

- Parameters visited during one time inference

- $$P = n_{layer} * (4 * d_{hidden}^2 + 2 * 4 * d_{hidden}^2) + n_{vocab} * d_{hidden}$$

- Memory latency & Compute latency

- $$latency_{memory} = \frac{P * n_{bytes}}{Peak\ memory\ bandwidth}$$

- $$latency_{compute} = \frac{2 * P * (BatchSize * SeqLen)}{Peak\ FLOPs}$$

- Arithmetic intensity

- $$AI = \frac{2 * P * B * S}{P * n_{bytes}} = \frac{2 * B * S}{n_{bytes}} \text{ FLOIPS/byte}$$

- Peak AI for SPR-SP with BF16 with AMX

- $$MAX\ AI_{spr-spbf16} = \frac{123.2TFLOPS}{307.2G/1000} = 401 \text{ FLOIPS/byte}$$

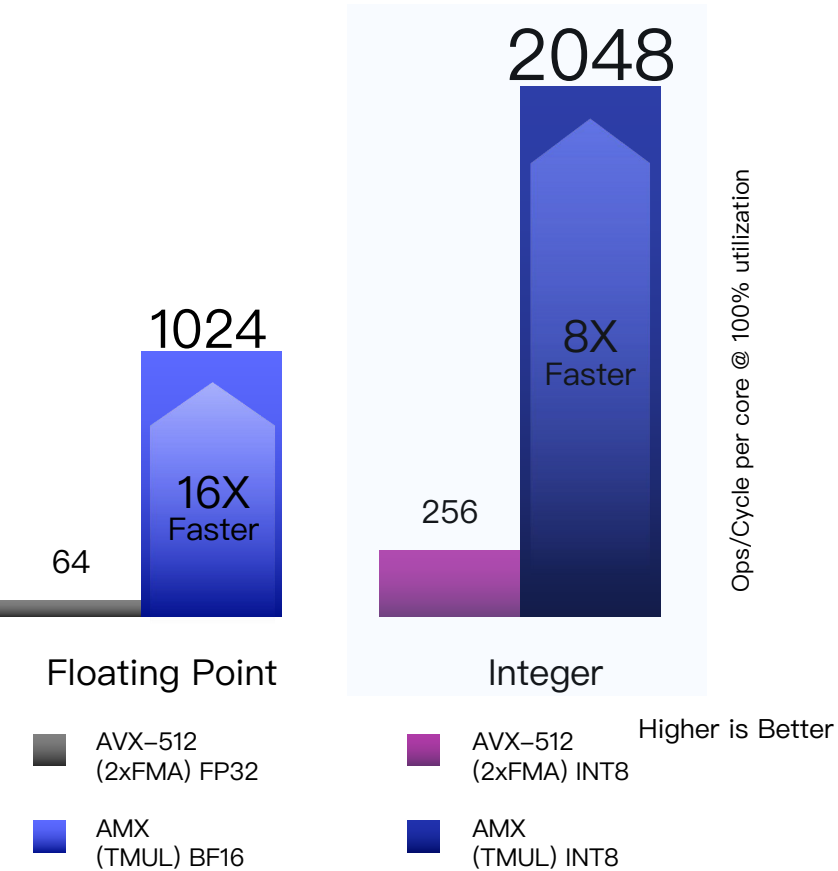
- Compute bound

- $$AI > MAX\ AI_{spr-spbf16} \Rightarrow B * S > 401$$

- Memory bound

- $$AI > MAX\ AI_{spr-spbf16} \Rightarrow B * S < 401$$

# Latest CPU Features for LLM



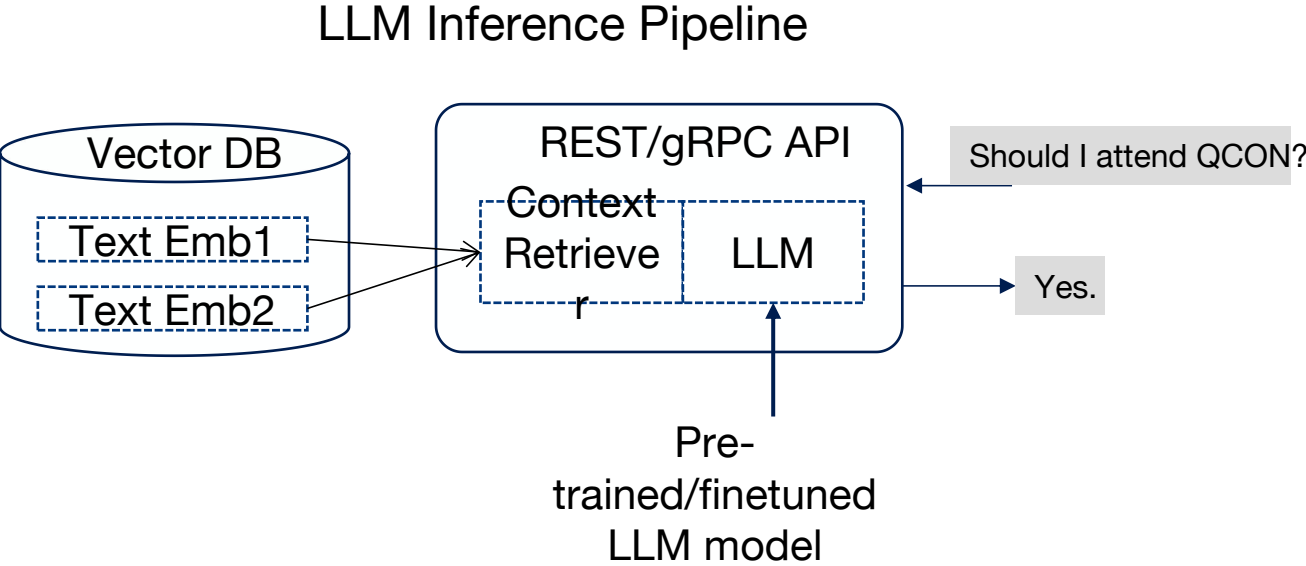
Results have been simulated. For workloads and configurations visit [www.intel.com/ArchDay21claims](https://www.intel.com/ArchDay21claims). Results may vary



Only x86 CPU with High Bandwidth Memory (HBM)

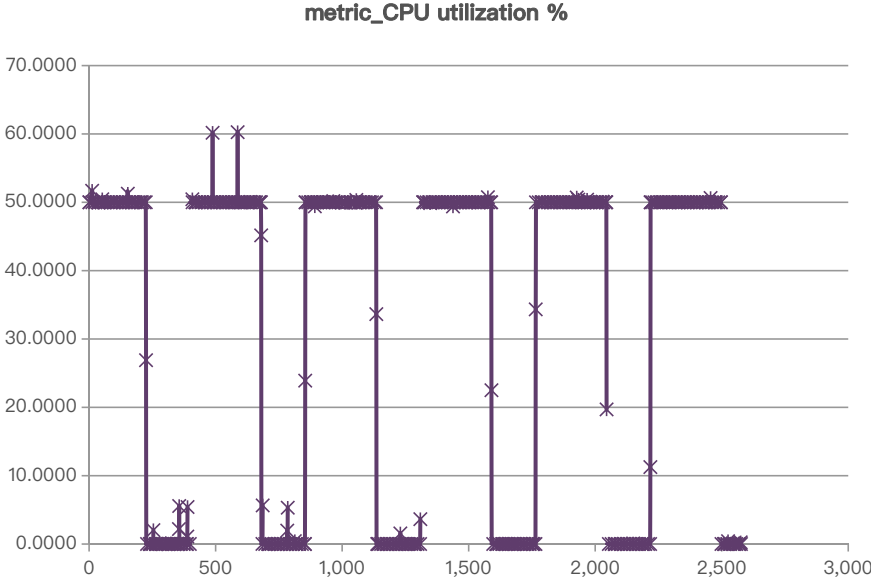
	<b>64GB</b> HBM2e	<b>112.5MB</b> shared	<b>DDR5</b> 8 channels per CPU @ 4800MTS (1DPC) / 16 DIMMs per socket
~1TB/s memory BW			
>1GB/core HBM memory capacity			

# CPU is NOT Fully Utilized!



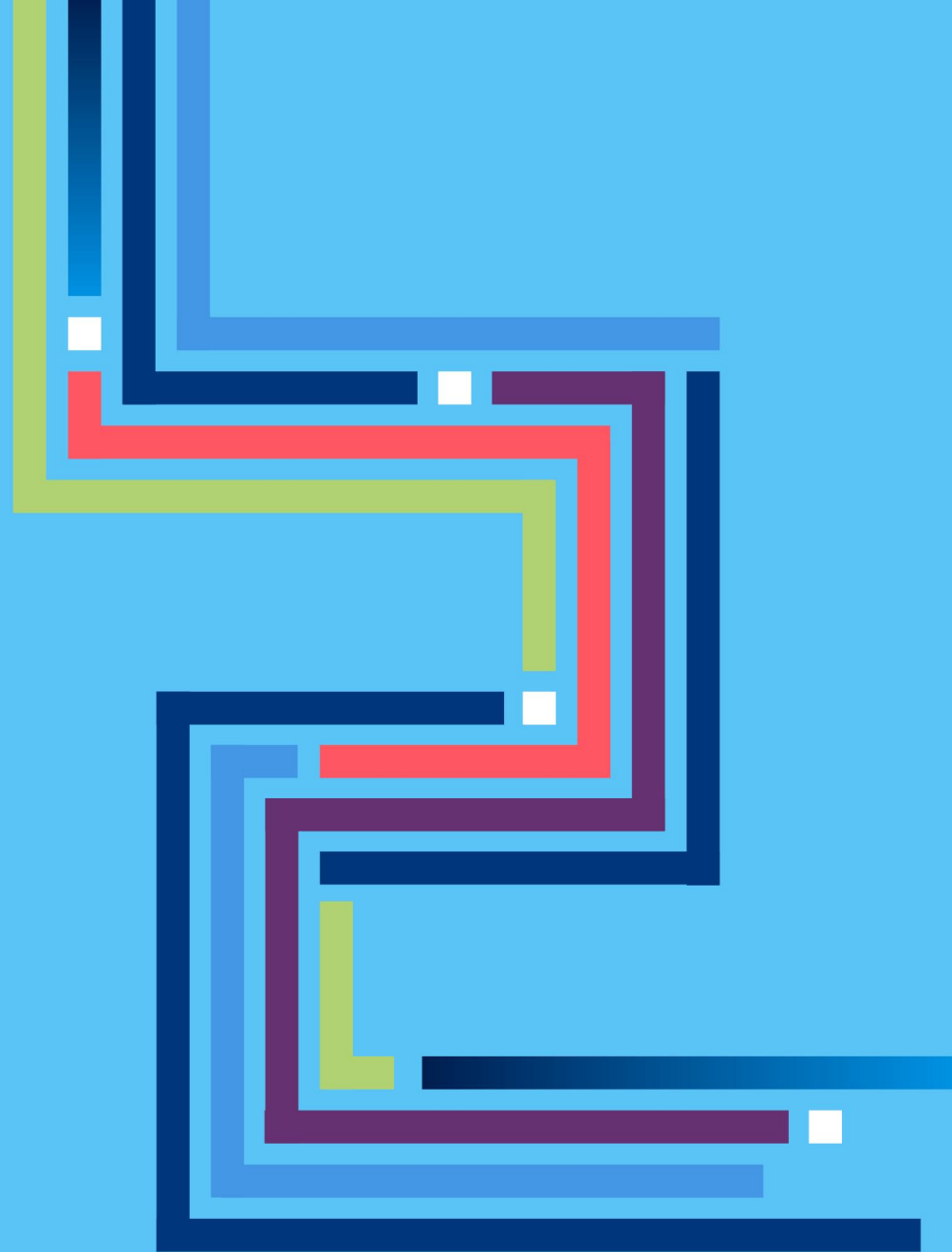
Pre-processing and post processing in LLM inference is relatively simple and do not need too much CPU resource.

CPU Utilization in LLM Training (offload mode)



Even for offload LLM training, CPU is still not fully utilized.

CPU上如何优化大语言模型？



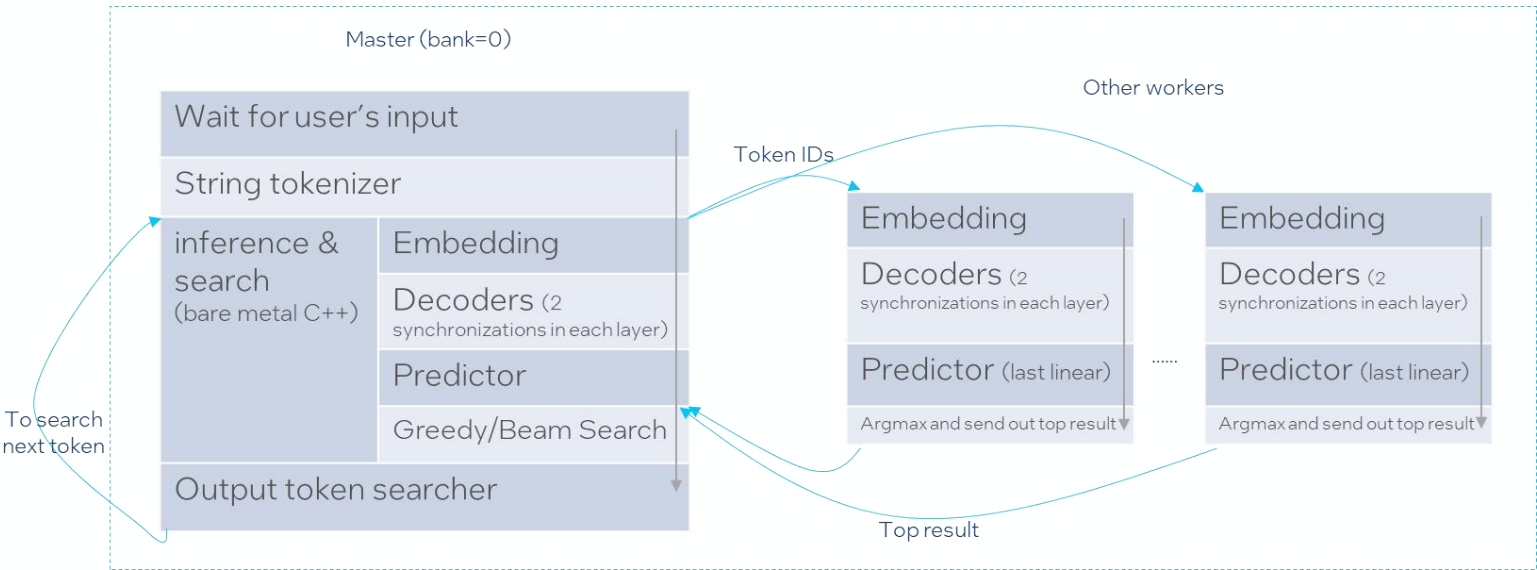


# Optimization

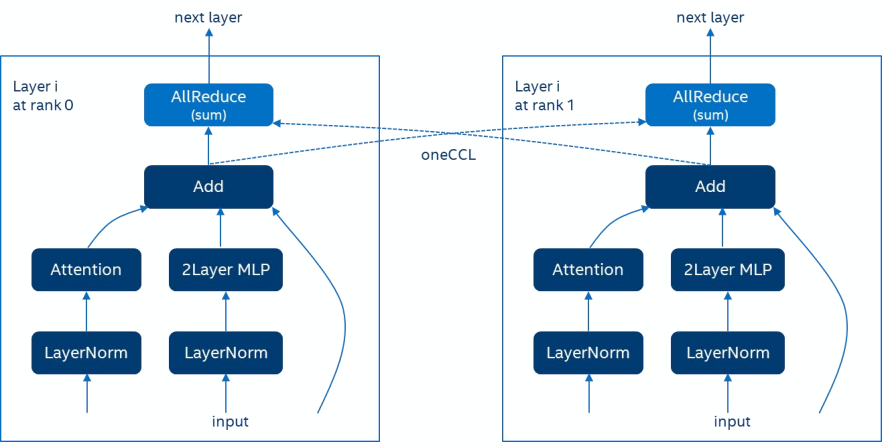
- Leverage the high-performance kernel (e.g., oneDNN)
- Avoid redundant computing
  - Continuous batching
  - Causal masking
  - Prefix sharing
- Lower precision & Sparsity
- Graph fusion
- Minimize memory copy and reorder
- Reuse the memory
- Distributed inference & use efficient communication library — oneCCL
- Runtime tuning

# Optimization for Distributed Inference

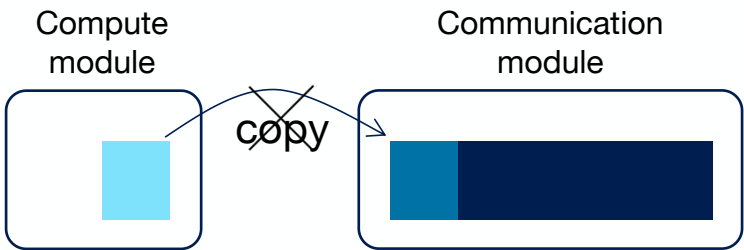
## Distributed inference based on oneCCL



Improve scalability by minimizing synchronization.



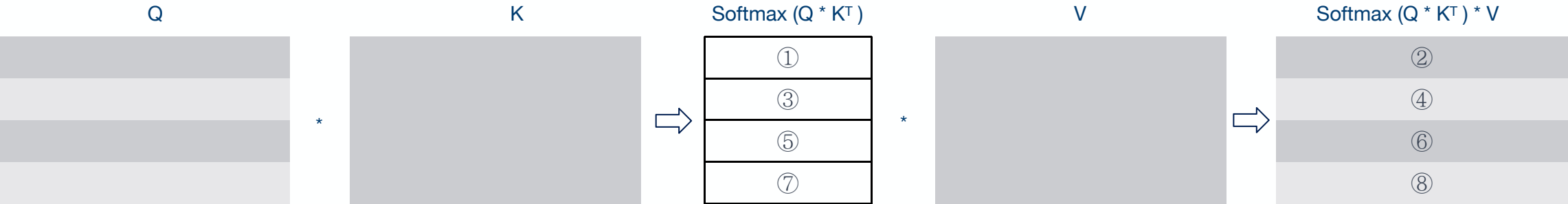
One time synchronization per layer is enough for some models



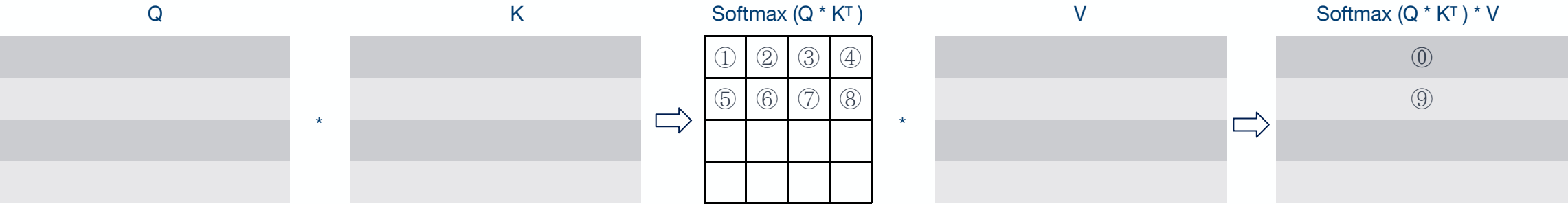
Minimize memory copy with full stack ownership

# Attention Optimization

SlimAttention  
(split score in 1 dimesion)



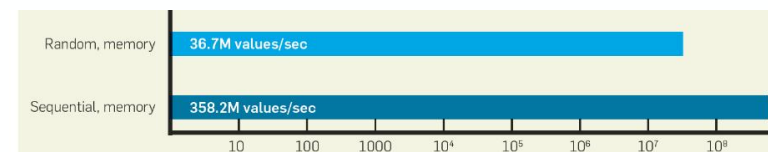
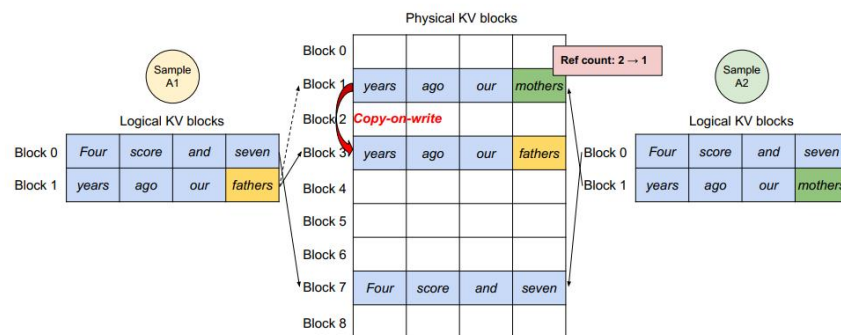
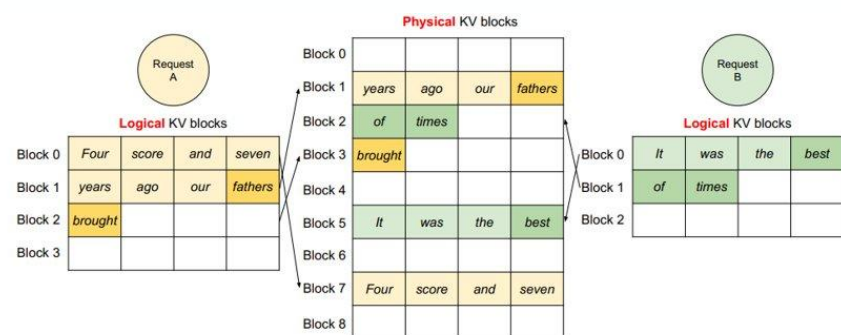
Computing Order: ①②③④⑤⑥⑦⑧  
Intermediate score size: Same size with ①  
Advantage: less score buffer w/o redundant computing



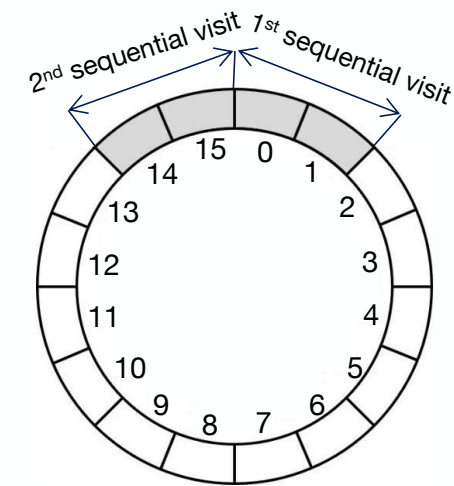
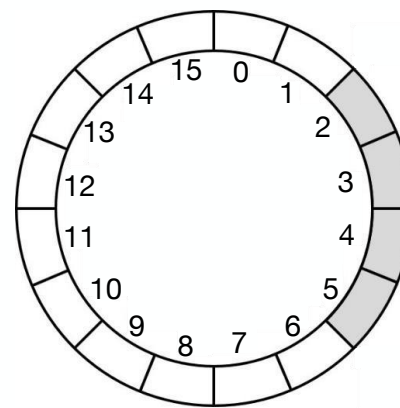
FlashAttention  
(split score in 2 dimesion)

Computing Order: ①⑩②⑩③⑩④⑩⑤⑨...  
Intermediate score size: Same size with ①  
Advantage: minimal intermediate buffer w/ some redundant computing

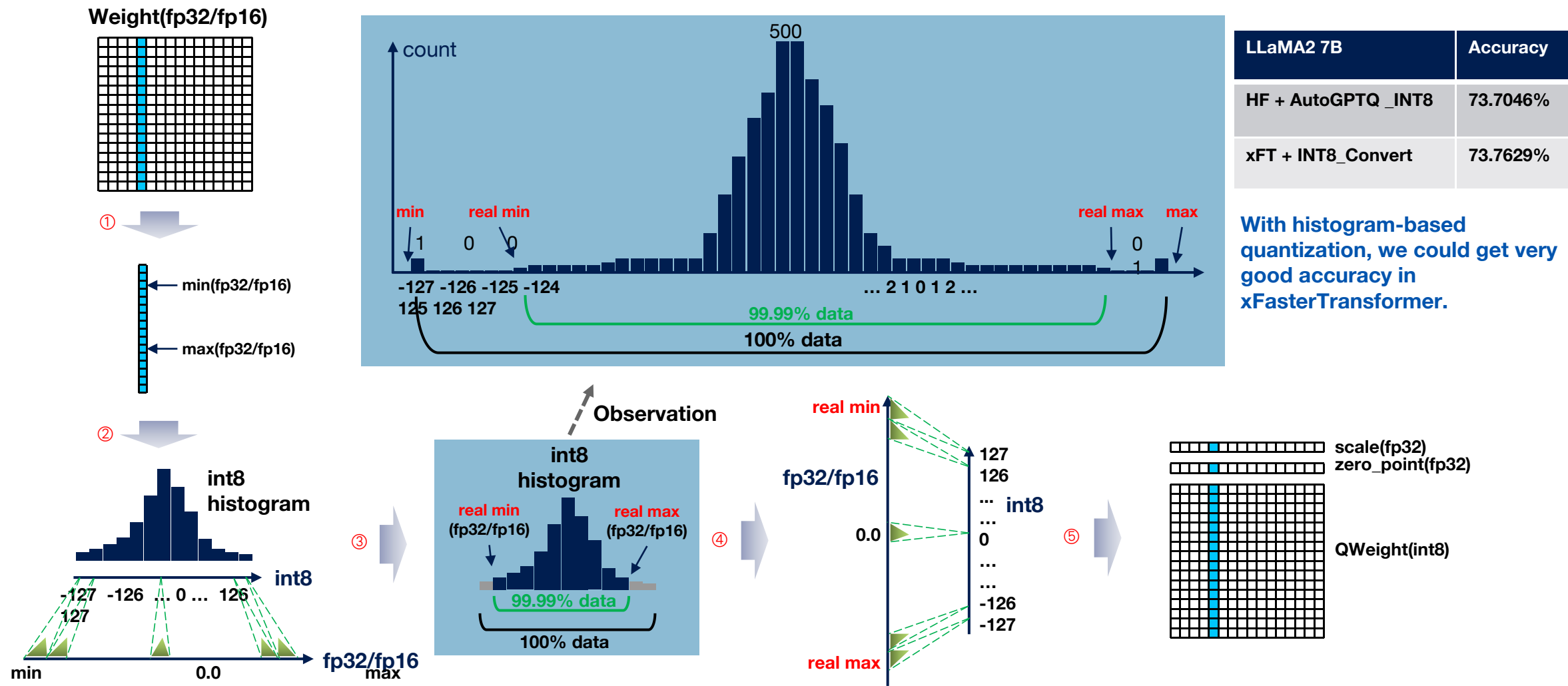
# Do We Need Paged Attention on CPU?



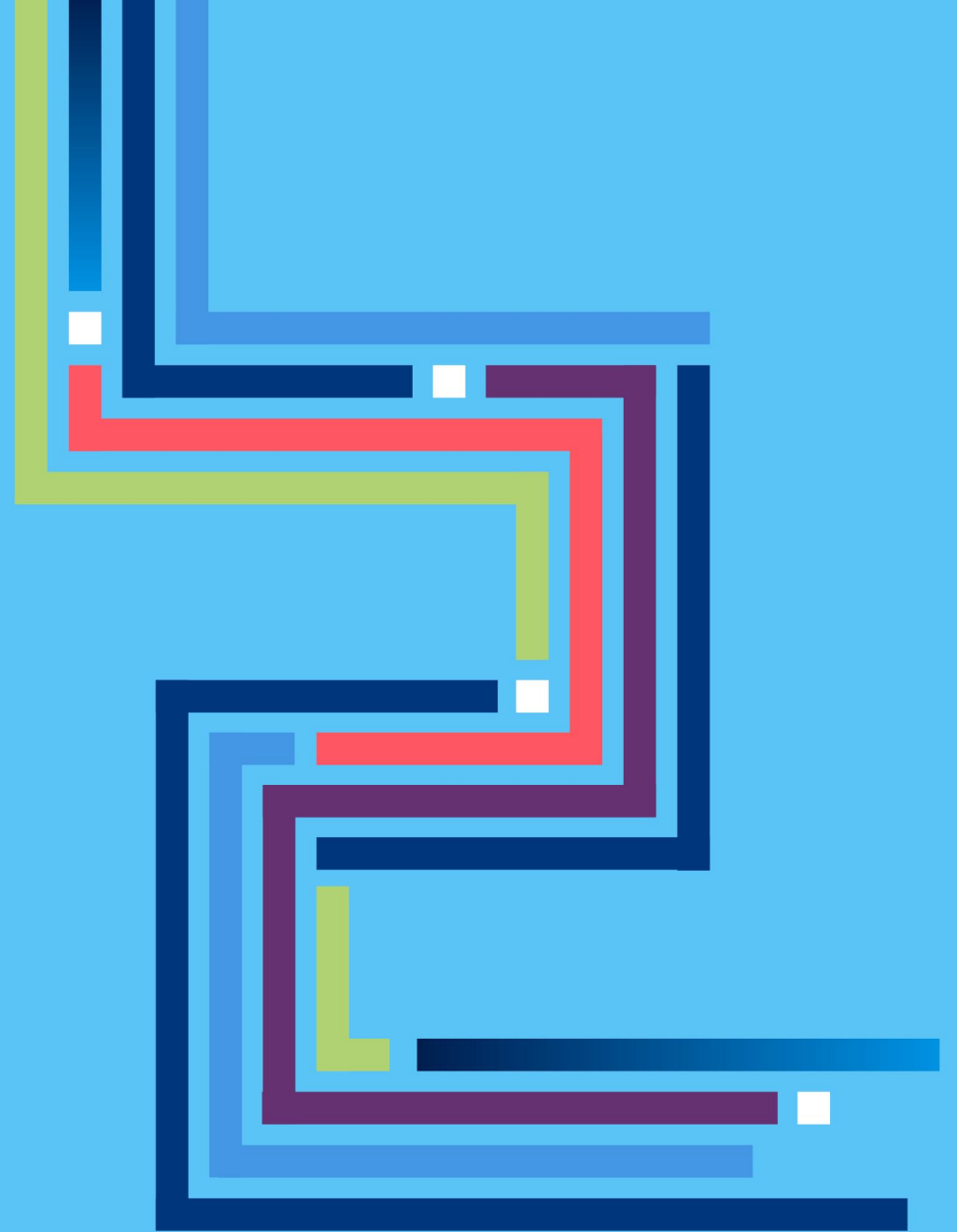
“Pathologies of Big Data” by Adam Jacobs in the ACM Communications, 2009



# Int8 Weight Only Quantization

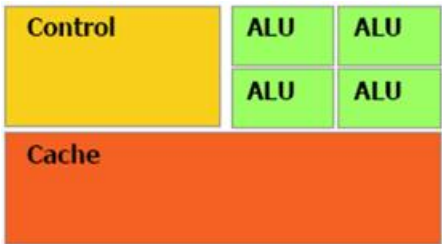


最大化CPU价值

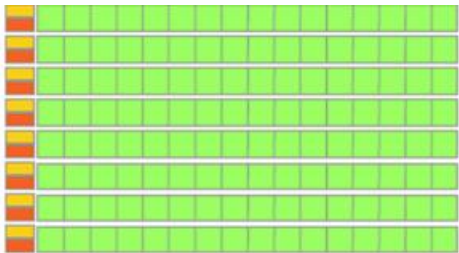


# CPU vs. GPU

CPU



GPU



A smaller number of larger cores

Low latency

Performs fewer instructions per clock

Designed and optimized for complex programs w/ serial processing

Automatic cache management

Large memory capacity

A larger number (thousands) of smaller cores

High throughput

Performs more instructions per clock

Optimized for parallel processing w/ bulk repetitive calculations

Allows for manual memory management

Limited memory capacity

## Key Challenges in LLM Inference

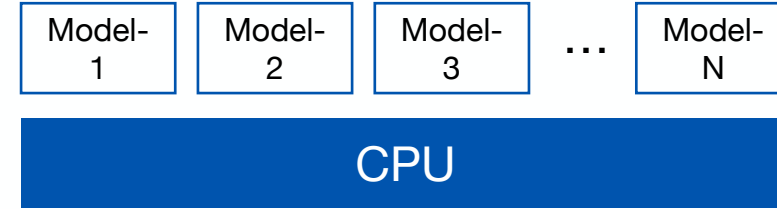
- Autoregressive generation
- Attention (square)
- Model is large

## Key HW Factors in LLM Inference

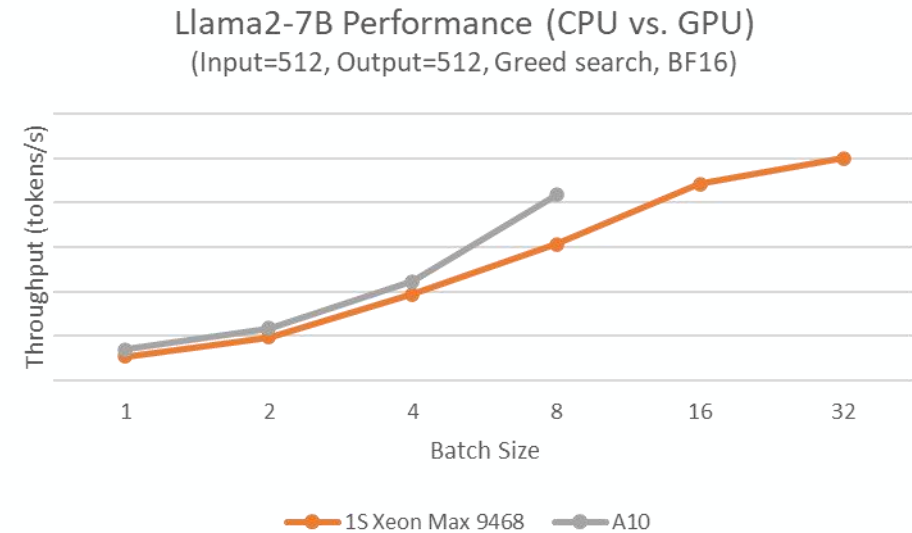
- Memory
- bandwidth
- Computing
- Memory capacity

# Scenarios CPU has values

- Long tail models (many models, few requests)
- Offline mode (to maximize throughput)
- Occasional demand
- Very long prompt token size and no strict latency requirement
- Very large model and no enough GPU
- Hybrid solution (e.g., speculative sampling)



- All models loaded in memory
- Not all models serving together





# Speculative Decoding

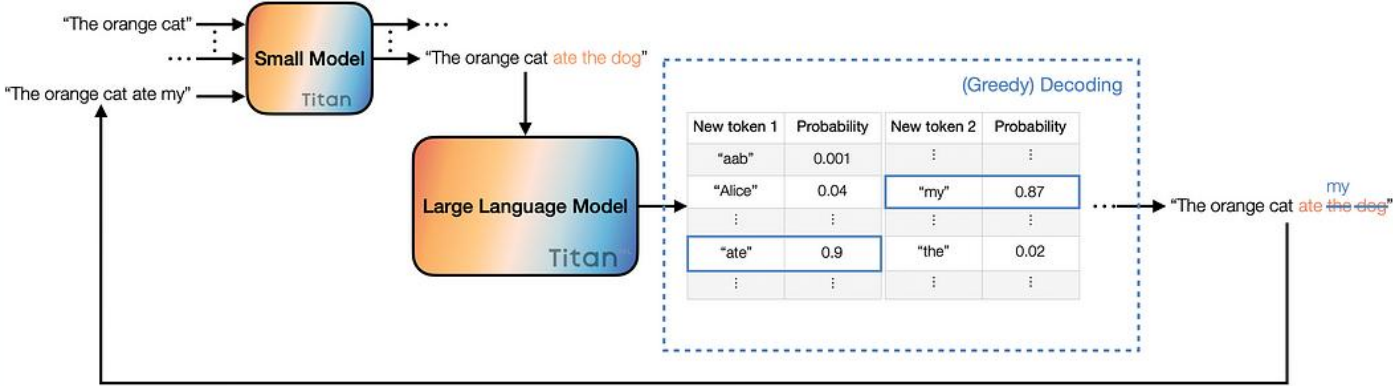
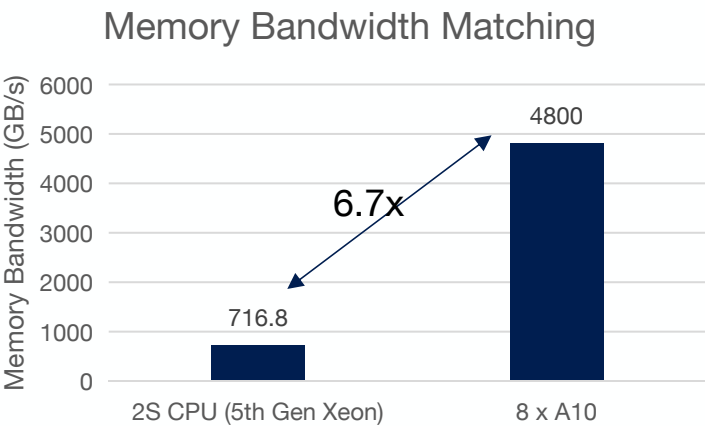


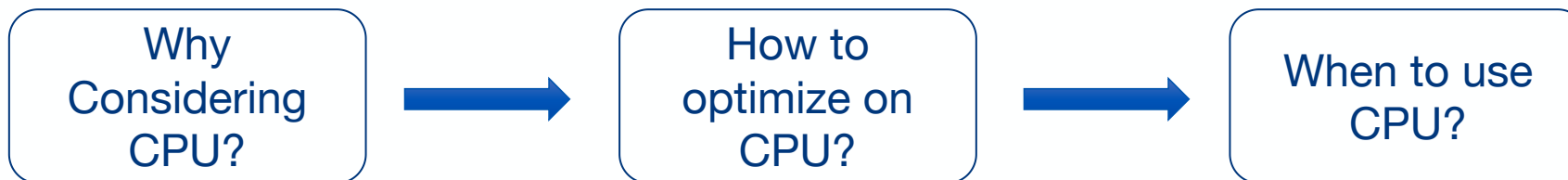
Image from <https://medium.com/@TitanML/in-the-fast-lane-speculative-decoding-10x-larger-model-no-extra-cost-f33ea39d065a>



In Practice, draft models are about 15-20x smaller than the target model.

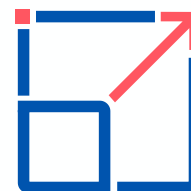


# 总结



Try our solution on Xeon

<https://github.com/intel/intel-extension-for-pytorch>



Build your own solution or  
seek ultimate perf

<https://github.com/intel/xFasterTransformer>