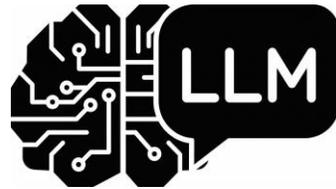




LLM Inference Lectures

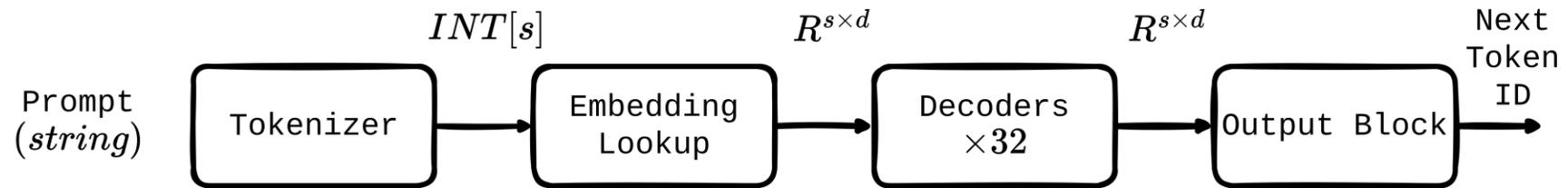
Part 2: LLM Operators

(Harvard CS265 - DATA/AI Systems LAB)



Milad Rezaei Hajidehi
(milad@seas.harvard.edu)

Overview of an LLM



s: number of tokens in our prompt.

d: dimension of vectors we use for embeddings.

We Inspect All of Them!

Tokenizer

Goal:

Raw text -> Sequence of Discrete Tokens

"future belongs to" -> "future belongs to" -> [21733, 17623, 311]

Vocabulary (size: v):

List of all possible tokens.

Determined at the training time.

Coarse vs. Fine Grained Tokens?

future belongs to
future belongs to
future belongs to

?



Smaller Sequence Length.

Bigger Vocabulary Size. (Last Output Layer)

Weaker Generalization.

Byte-Pair Encoding

Vocabulary Size (v) = Hyperparameter

How to Construct Vocabulary Set?

Start with Each Byte as a Token
Greedy Merge Tokens and Add a New Token
Based on Frequencies.

be -> be

At Inference: Start with Bytes, Greedy Merge as Much as You Can.
(based on *merge-scores*)

128000 Tokens in LLAMA3.

At Inference: Start with Characters(Bytes), Greedily Merge as Much as You Can.
(based on *merge-scores*)

128000 Tokens in LLAMA3.

b_e_l_o_n_g -> be_l_o_n_g -> be_lo_n_g -> be_lon_g -> be_long

id->token

id	token
1	a
2	b
...	...
43700	long
...	...

token->id

token	id
a	1
aa	264
...	...
long	43700
...	...

<id,id> -> rank
(frequency)

<id,id>	rank
<1,2>	42
<1,3>	824
...	...
<43700,1>	9999
...	...

Embedding Lookup

Goal: Map Token IDs to Vectors with Semantic.

Embeddings: A Lookup Table of $(v \times d)$

Afterward we have a $(s \times d)$ Tensor.

**Where Should
We Store
Embedding
Table?**

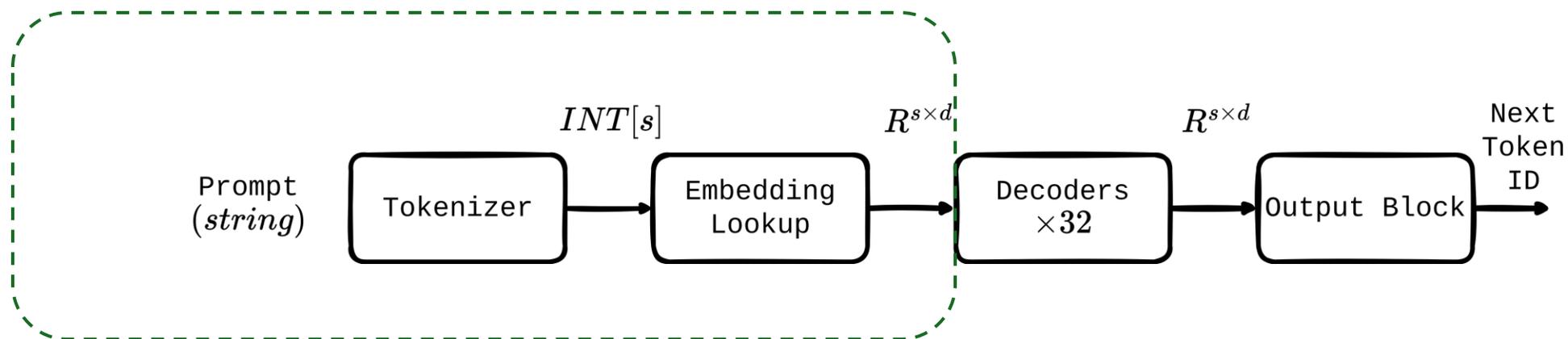
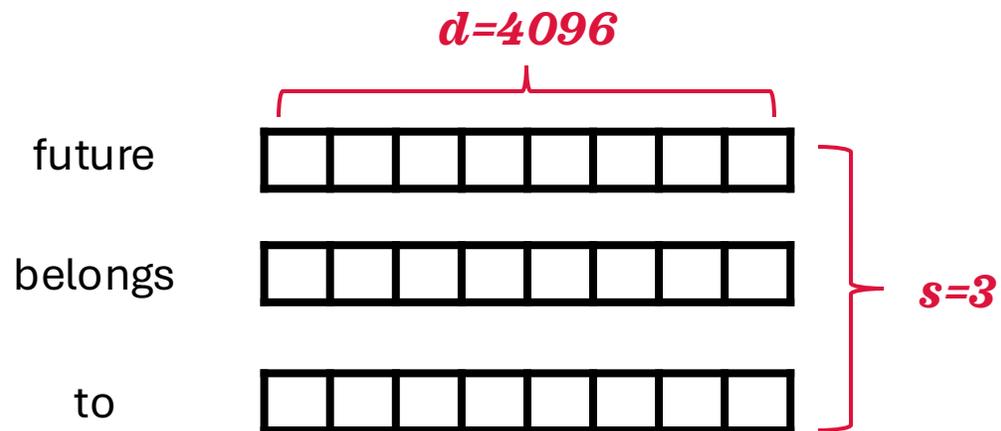
?



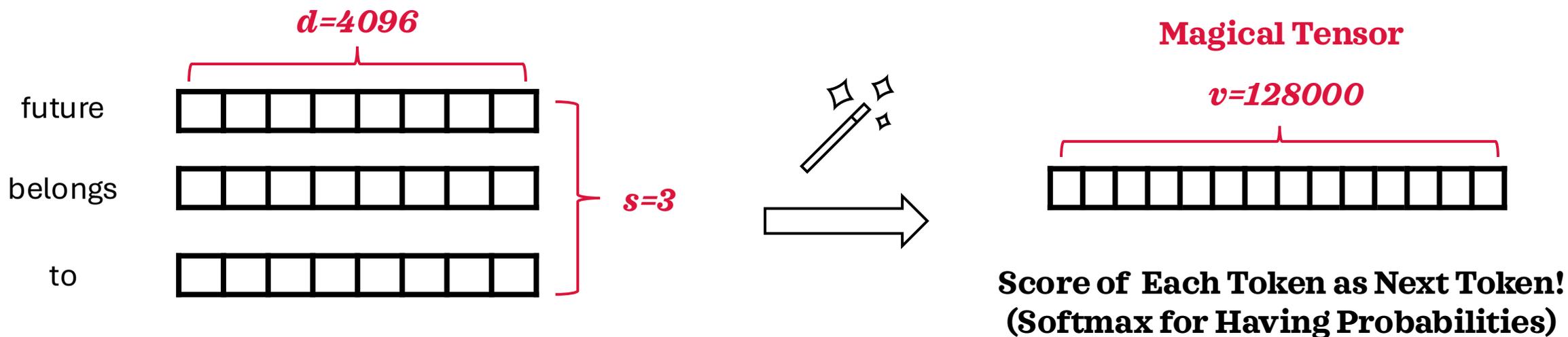
GPU!!

Avoid large tensor data-movement
Easier Duplicate Token Management

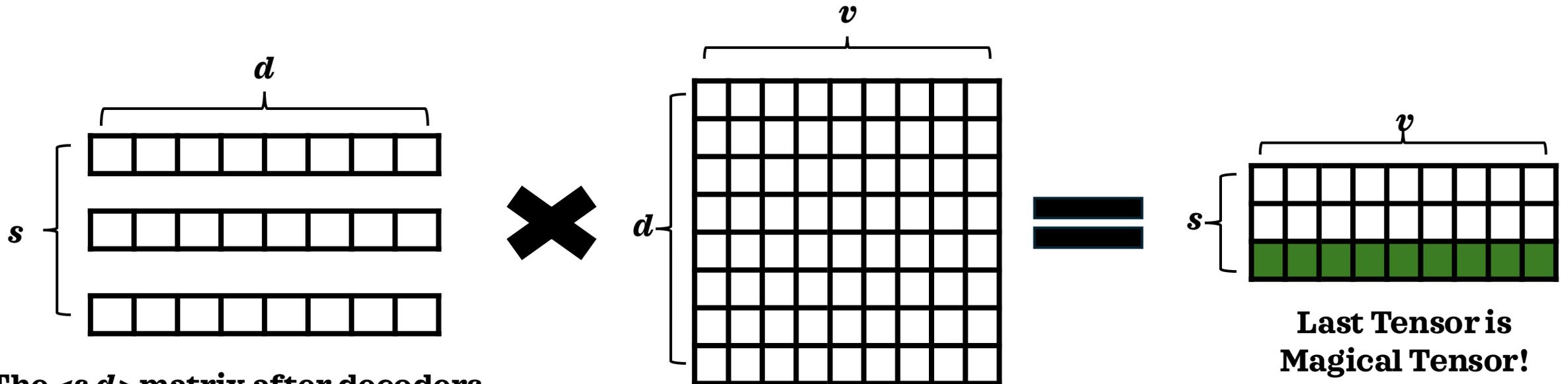
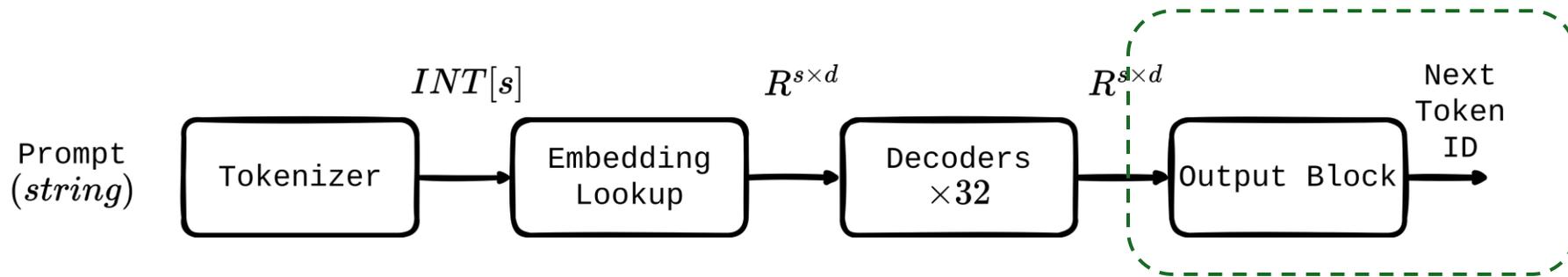
What We Have so Far...



Our Goal: Best Next Token



Output Block: Finding Magical Tensor



The $\langle s, d \rangle$ matrix after decoders

Learned Matrix in Output Block
(Vocabulary Projection)

Softmax give us Probabilities.

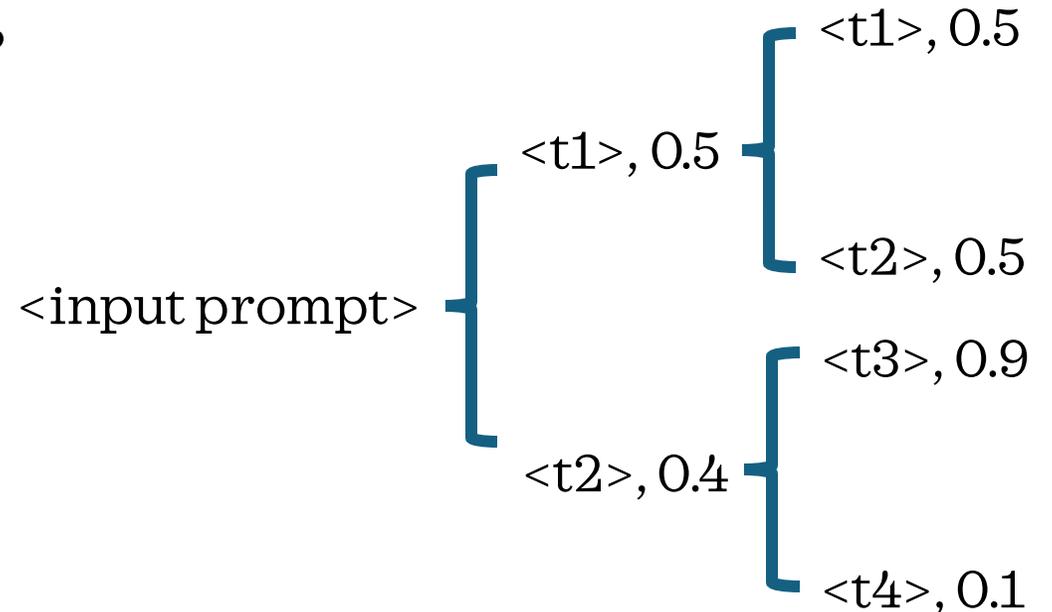
Decoding Algorithms: A Short Detour

Greedy Sampling: We always pick next token that has highest score/probability!

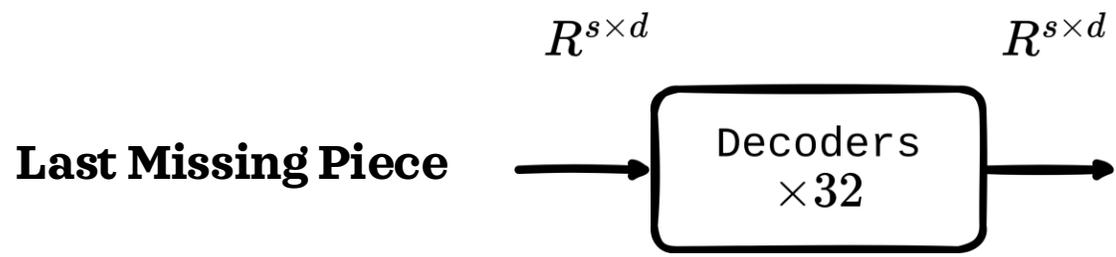
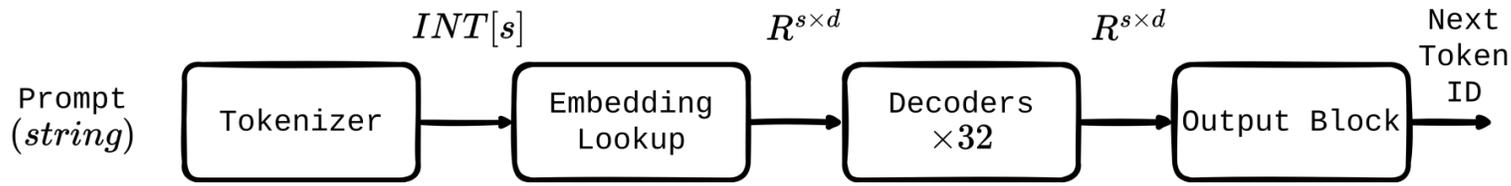
Is It Optimal?

Issue #1: Diversity of Answers.

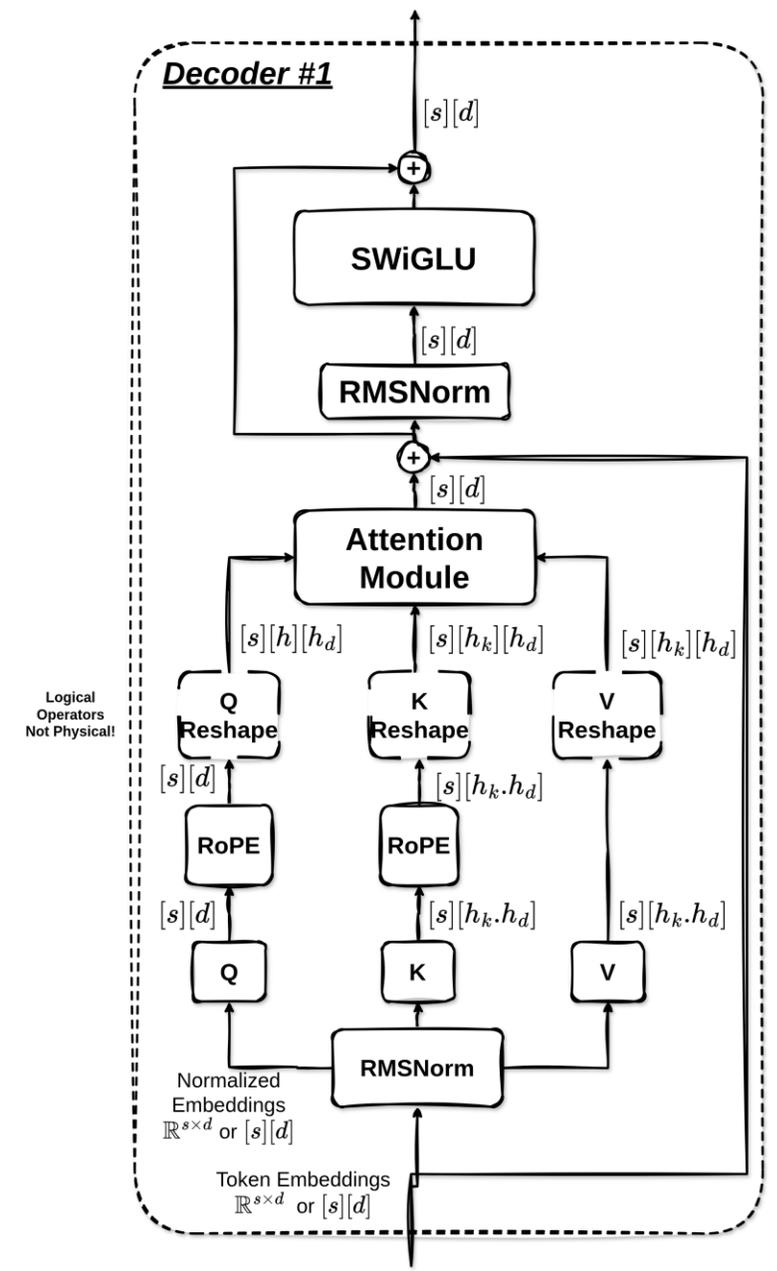
Issue #2: Locally Optimal != Globally Optimal



Beam Search & Temperature are Solution To This.



Goal: Transform Initial Embedding -> Contextual Representation for Output Block





Operators Inside Decoder:

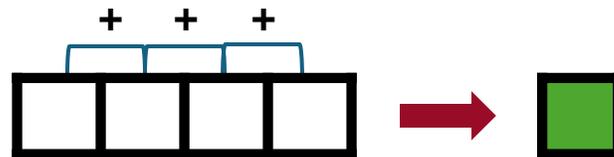
RMSNorm, QKV, KV-Cache, Attention, MLP
(90% of the Time)

$$\text{RMSNorm}(x) = \frac{x}{\text{RMS}(x)} \odot \gamma, \quad \text{RMS}(x) = \sqrt{\frac{1}{d} \sum_{i=1}^d x_i^2 + \varepsilon}$$

x : Embedding Per Token, s Norms in Total

Purpose: Normalize Layer Inputs, Faster Convergence, Stability

Requires A Reduction (sum) Kernel!



How to Calculate Sum of a Vector? (Threads are Independent!)



Step1: Each Thread Calculate Sum for its Own Color



Per Thread Sum in Shared Memory

T_id	Color	Sum
0	Blue	17
1	Red	29
2	Green	14
4	Yellow	11



Second 1/2 of Threads Finish, First 1/2 Threads Merge With Second Half Values (0 with 2, 1 with 4)



Second 1/4 of Threads Finish, First 1/4 Threads Merge With Second Half Values (0 with 1)

.....

Original Goal of Decoder: Transforming Embeddings to be Contextual Information.

"future belongs to self-"

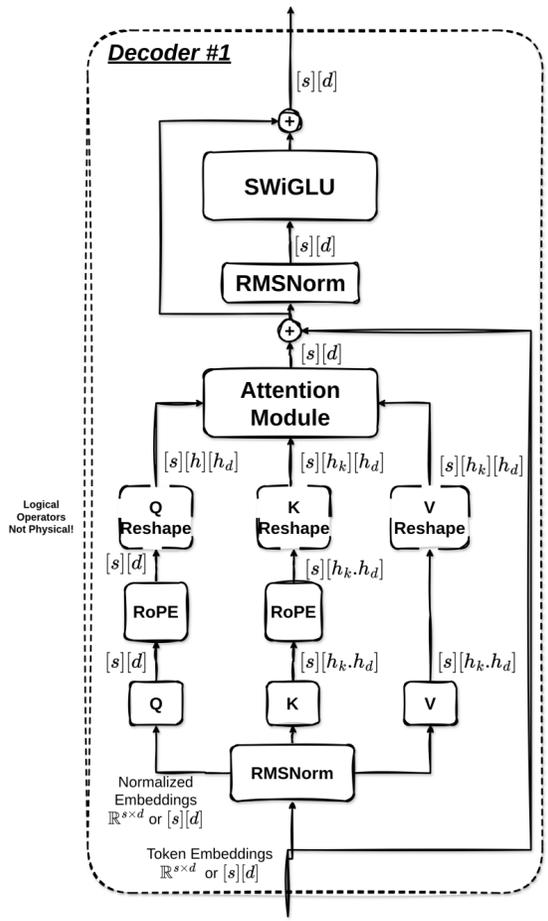
- "motivated" -> 70%
- "driving" -> 29%
- "designing" -> 1%

"elon said future belongs to self-"

- "motivated" -> 5%
- "driving" -> 95%
- "designing" -> 0%

"stratos said future belongs to self-"

- "motivated" -> 10%
- "driving" -> 0%
- "designing" -> 90%



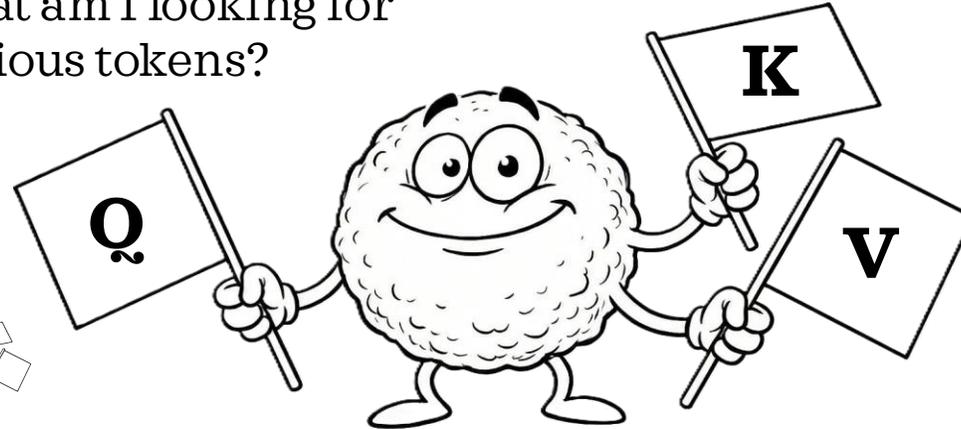
Input of Decoder: One embedding per token.

Output of Decoder/Attention(H): One embedding per token. But transformed!
based on context and previous tokens!

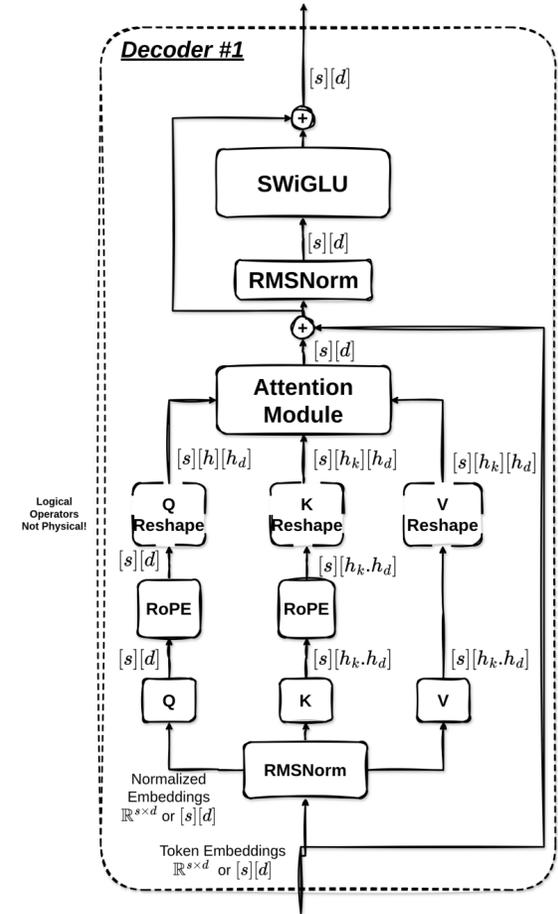
(Q)uery: What am I looking for in previous tokens?

(K)ey: What kind of query would match me?

(V)alue: What information do I pass if selected?



Token

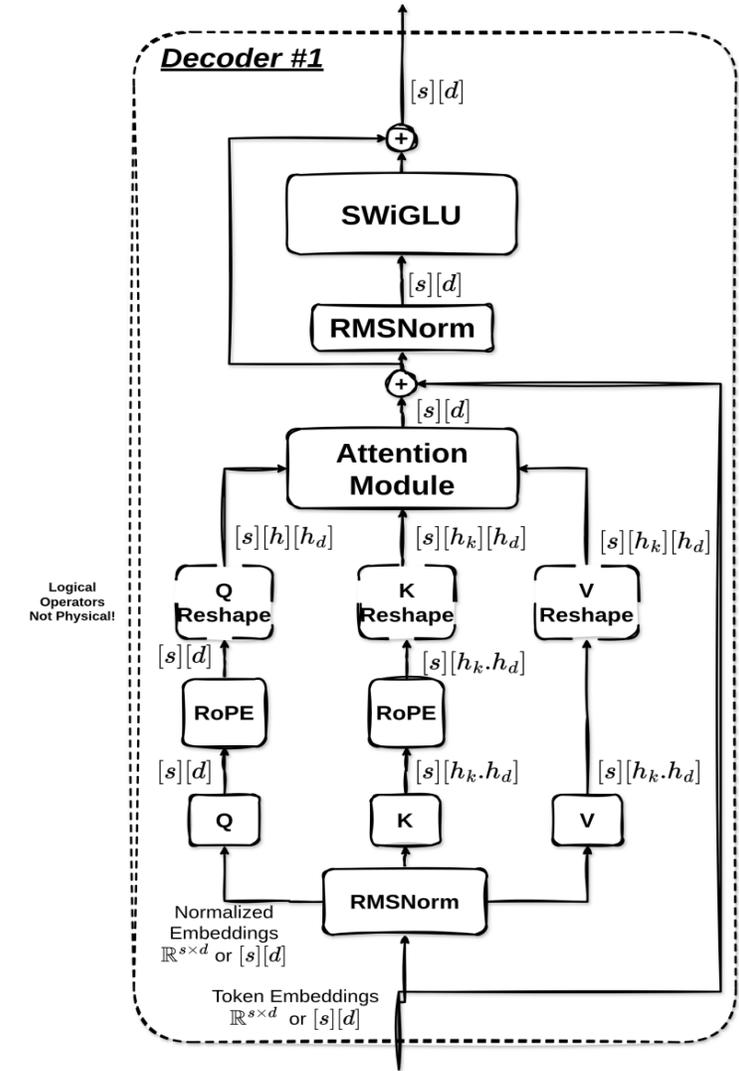


How much j can impact i ?
(normalized)

$$\alpha_{ij} = \frac{\exp\left(\frac{Q_i K_j^\top}{\sqrt{128}}\right)}{\sum_k \exp\left(\frac{Q_i K_k^\top}{\sqrt{128}}\right)}$$

Embeddings Transformed.

$$\text{Attention}(Q, K, V) = AV$$



In LLAMA3, We use GQA rather than Vanilla-Transformer, which is slightly different.

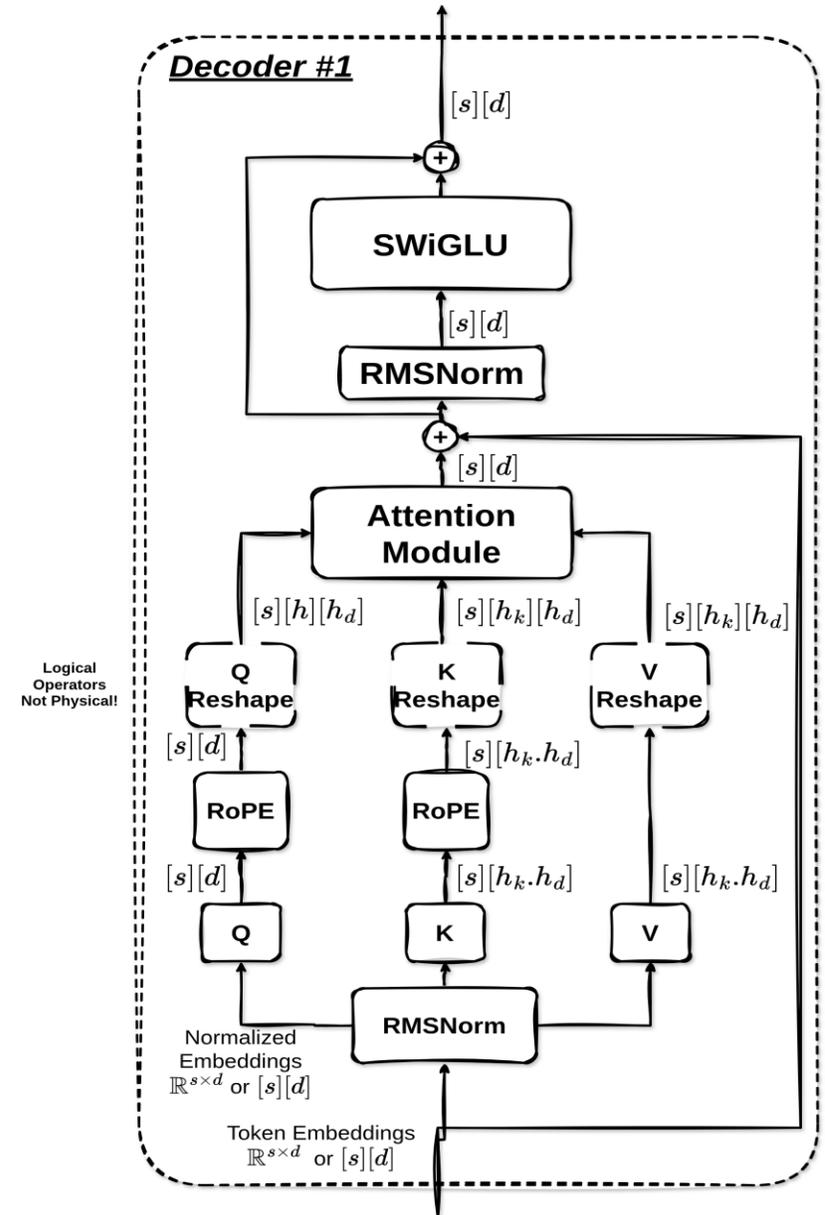
We need QKV for all tokens!

After first decoder, all QKV are affected by their previous tokens!

KV-Cache:

Cache previously computed K,V to avoid recomputing them at every decoding step.

After QKV and Attention we have bunch of Matrix Multiplication and elementwise operators.



You Are Ready For The Project!

Step 0

Read CUDA Book
(selected chapters)

Understand LLAMA
Operators.
(Slides, Office Hours,
Youtube, LLM, Torch
Codes)

Step 1

Finish Tokenizer

Write dumper.py
and loader.cpp to
transform weights.
(mmap)

Load embeddings and
prepare the input for
first Kernel
(RMSNorm)

Step 2

Write a good generic
Matrix Multiplication
kernel.

Write a Reduce kernel.

Step 3

Write your Controllers
(operators)

Write rest of easy
kernels.

Test, Test, Test!

You Are Ready For The Project!

Step 0

Step 1

Step 2

Step 3

First Make it Work.

Then Optimizations:

KV-Cache, Fusing, Performance Tuning, Quantization...

Codes)

Load embeddings and
prepare the input for
first Kernel
(RMSNorm)

Write rest of easy
kernels.

Test, Test, Test!

Challenges & Bottlenecks.
What Makes Inference Hard?

✓ **Opportunities For Problem-Solving!**

Performance Metrics in Inference

Phase 1 - Prefill: Reading prompt - Compute QKVs for each token - Find the first "next token"

Metric -> Time-to-First-Token (TTFT) = $T_{\text{queue}} + T_{\text{prefill}}$

Compute Bound!

Phase2 - Decode: Predict next tokens one by one - Read QKV from cache.

Metric -> Tokens Decoded per Second (token/s) or Latency = TTFT + Full Decode Latency

Memory Bound!

Inference Challenges

**Highly variable sequence lengths
(prompt, answer)**

Prompt Tokens	Output Tokens
10	300
170	289
16	29
290	46

Why It Cause a Problem?



Inference Challenges

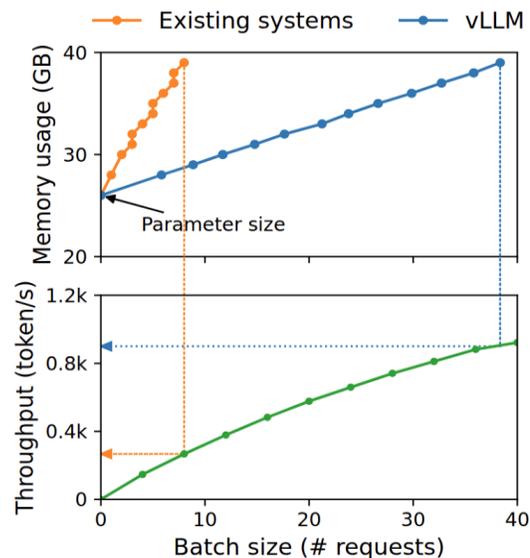
Highly variable sequence lengths
(prompt, answer)

Memory is Under Pressure.

For a 500 token request:
 $500 * 32 * 2 * 4096 * 4b \approx \mathbf{500MB!!}$

Model Weights
~(20GB)

GPU Memory



Inference Challenges

Highly variable sequence lengths
(prompt, answer)

Memory is Under Pressure.

Too Many Operators/Kernel Launches!

Why It Cause a Problem?



Inference Challenges

Highly variable sequence lengths
(prompt, answer)

Memory is Under Pressure.

Too Many Operators/Kernel Launches!

Bursty, Irregular, On/Off Workload

And of Course Everything We Said About **Kernel Tuning**.

Every Paper We Read: A Solution to These Challenges. General or Application-Specific

Inference
Challenges

Too Many Operators/Kernel Launches!

Bursty, Irregular, On/Off Workload

And of Course Everything We Said About Kernel Tuning.

Problem: We launch many kernels for one inference.
~25 kernel launch per decoder if done naively.

Why We Care: GPU Underutilization.
Low throughput due to frequent data movement.

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

Tri Dao[†], Daniel Y. Fu[†], Stefano Ermon[†], Atri Rudra[‡], and Christopher Ré[†]

Mirage: A Multi-Level Superoptimizer for Tensor Programs

Mengdi Wu Xinhao Cheng Shengyu Liu[†] Chunan Shi[†] Jianan Ji
Man Kit Ao Praveen Velliengiri[‡] Xupeng Miao[#] Oded Padon[◇] Zhihao Jia

Problem: Many design-decisions and parameters for each kernel.

Why We Care: Optimizing kernels independently is necessary for end-to-end optimization of LLM.

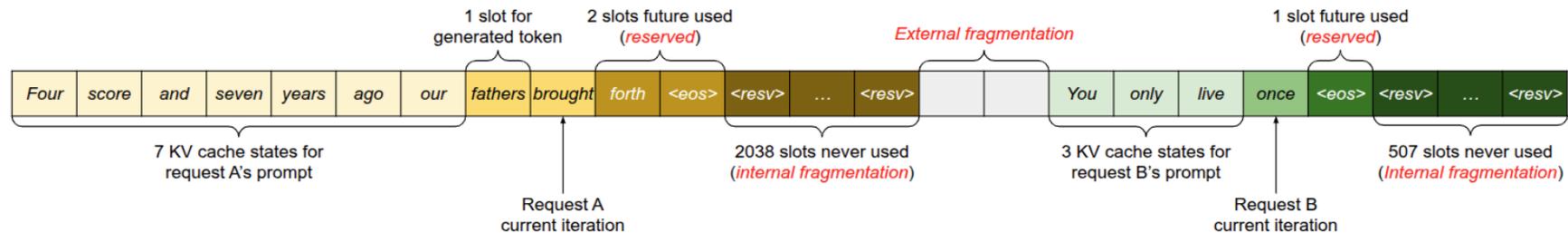
TVM: An Automated End-to-End Optimizing Compiler for Deep Learning

Tianqi Chen¹, Thierry Moreau¹, Ziheng Jiang^{1,2}, Lianmin Zheng³, Eddie Yan¹

Meghan Cowan¹, Haichen Shen¹, Leyuan Wang^{4,2}, Yuwei Hu⁵, Luis Ceze¹, Carlos Guestrin¹, Arvind Krishnamurthy¹

Problem: Naive KV cache allocation wastes GPU memory.

Why We Care: Lack of memory limit batch-size and hence throughput.



Efficient Memory Management for Large Language Model Serving with *PagedAttention*

Woosuk Kwon^{1,*} Zhuohan Li^{1,*} Siyuan Zhuang¹ Ying Sheng^{1,2} Lianmin Zheng¹ Cody Hao Yu³
Joseph E. Gonzalez¹ Hao Zhang⁴ Ion Stoica¹

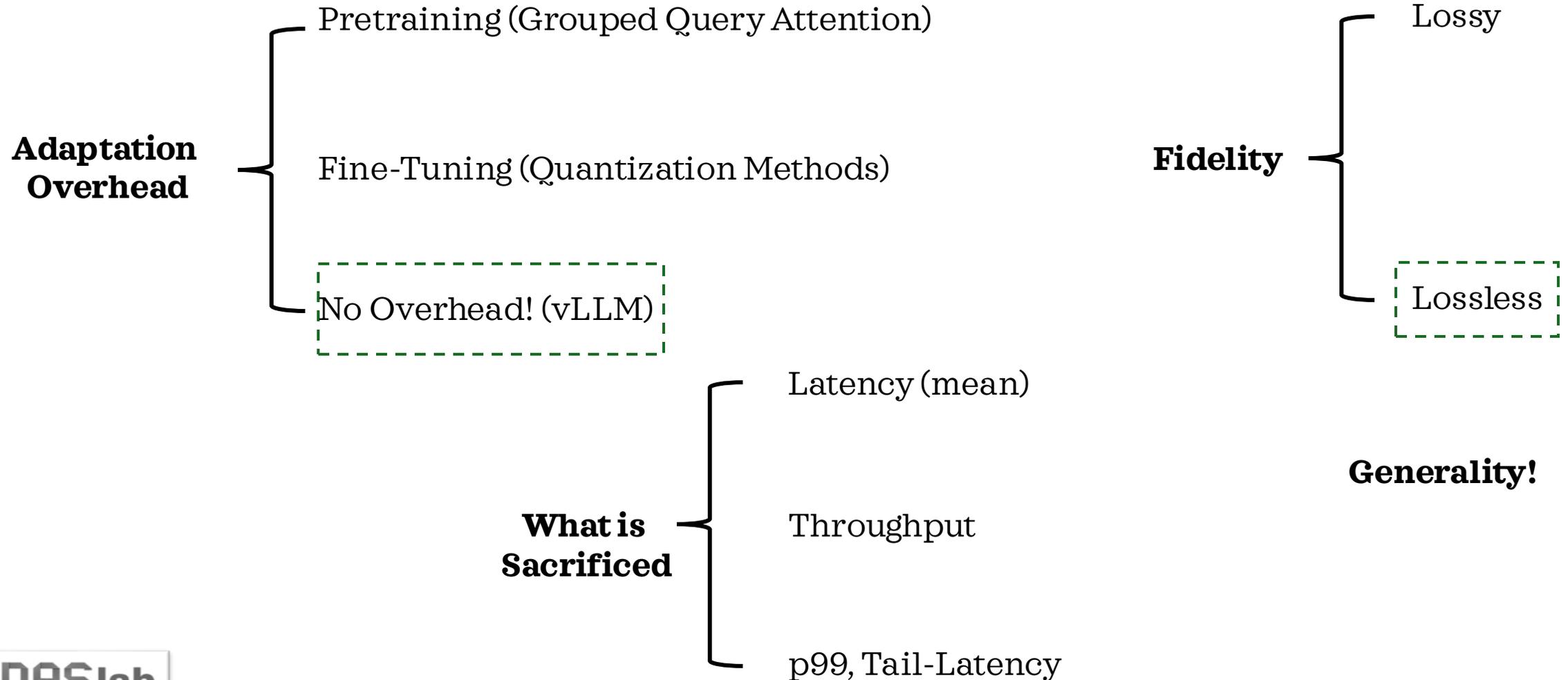
Problem: Naive KV cache allocation wastes GPU memory.

And Many More Papers....

Efficient Memory Management for Large Language Model Serving with *PagedAttention*

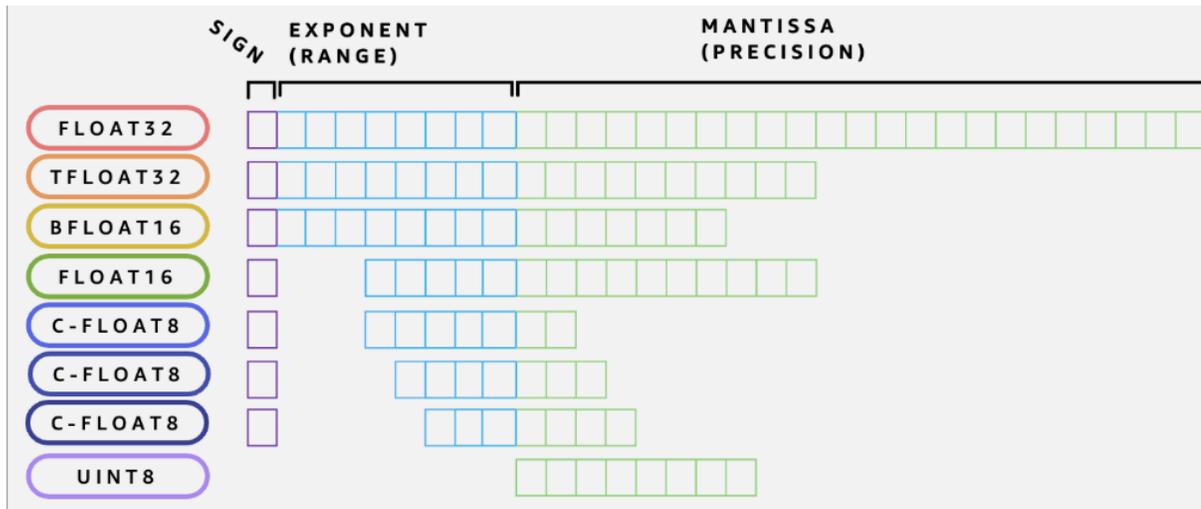
Woosuk Kwon^{1,*} Zhuohan Li^{1,*} Siyuan Zhuang¹ Ying Sheng^{1,2} Lianmin Zheng¹ Cody Hao Yu³
Joseph E. Gonzalez¹ Hao Zhang⁴ Ion Stoica¹

How to Assess LLM Inference Techniques?



Quantization

Using Less Bits for Weights, Activations, KV-Cache -> **More Memory, More Batch Size, More Throughput**



Format	Smallest	Largest	Machine Precision
FP32	1.175×10^{-38}	3.403×10^{38}	$2^{-23} \approx 1.19 \times 10^{-7}$
FP16	6.104×10^{-5}	6.550×10^4	$2^{-10} \approx 9.77 \times 10^{-4}$
BF16	1.175×10^{-38}	3.390×10^{38}	$2^{-7} \approx 7.81 \times 10^{-3}$
FP8 (E4M3)	1.5625×10^{-2}	2.40×10^2	$2^{-3} = 1.25 \times 10^{-1}$
FP8 (E5M2)	6.104×10^{-5}	5.734×10^4	$2^{-2} = 2.5 \times 10^{-1}$

Possible Problems:

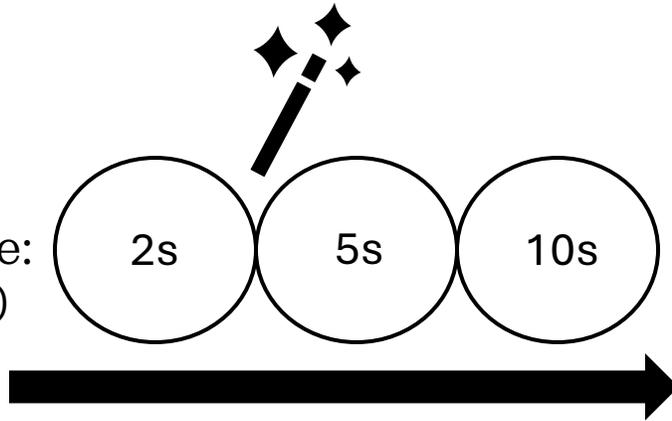
Overflow, Underflow, Low Precision

Lossy, May Require Fine-tuning.

$$x = (-1)^{\text{sign}} \times (1 + \text{mantissa}) \times 2^{(\text{exponent} - \text{bias})}$$

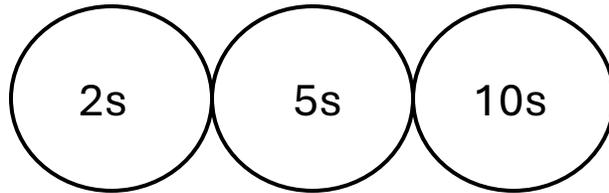
Scheduling Requests in LLM

Requests Queue:
(1 Req at a time)

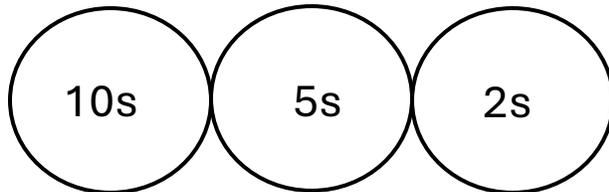


P99, Throughput,
Chicken-Egg Problem

Original Order:

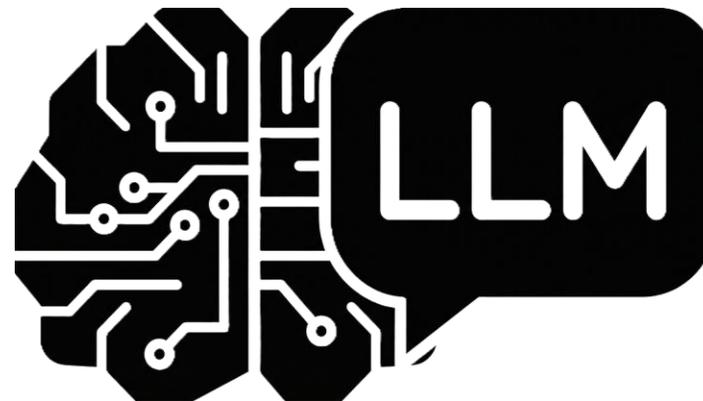
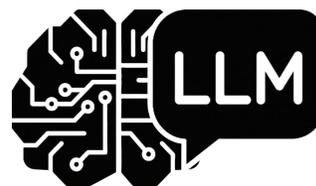
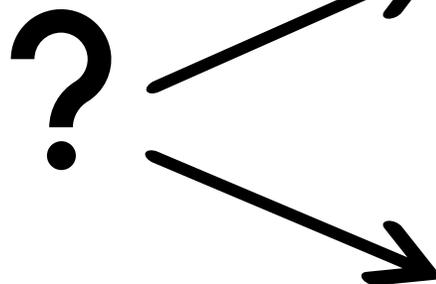
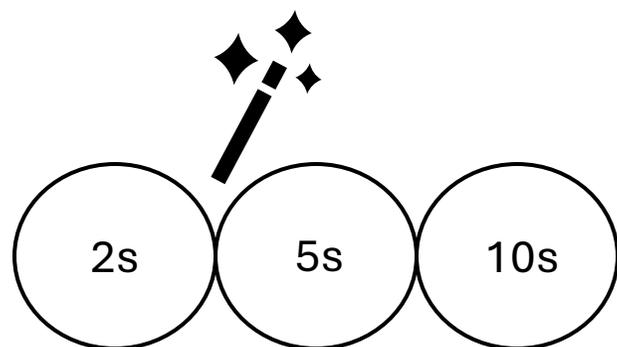


Re-order:



Can Improve Mean Latency.
Loseless, Only Train the Magical Part!

Routing Requests in LLM Inference



Improve Everything but Lossy!