

Spring25 CS598YP

24.2: Block Attention

Yongjoo Park

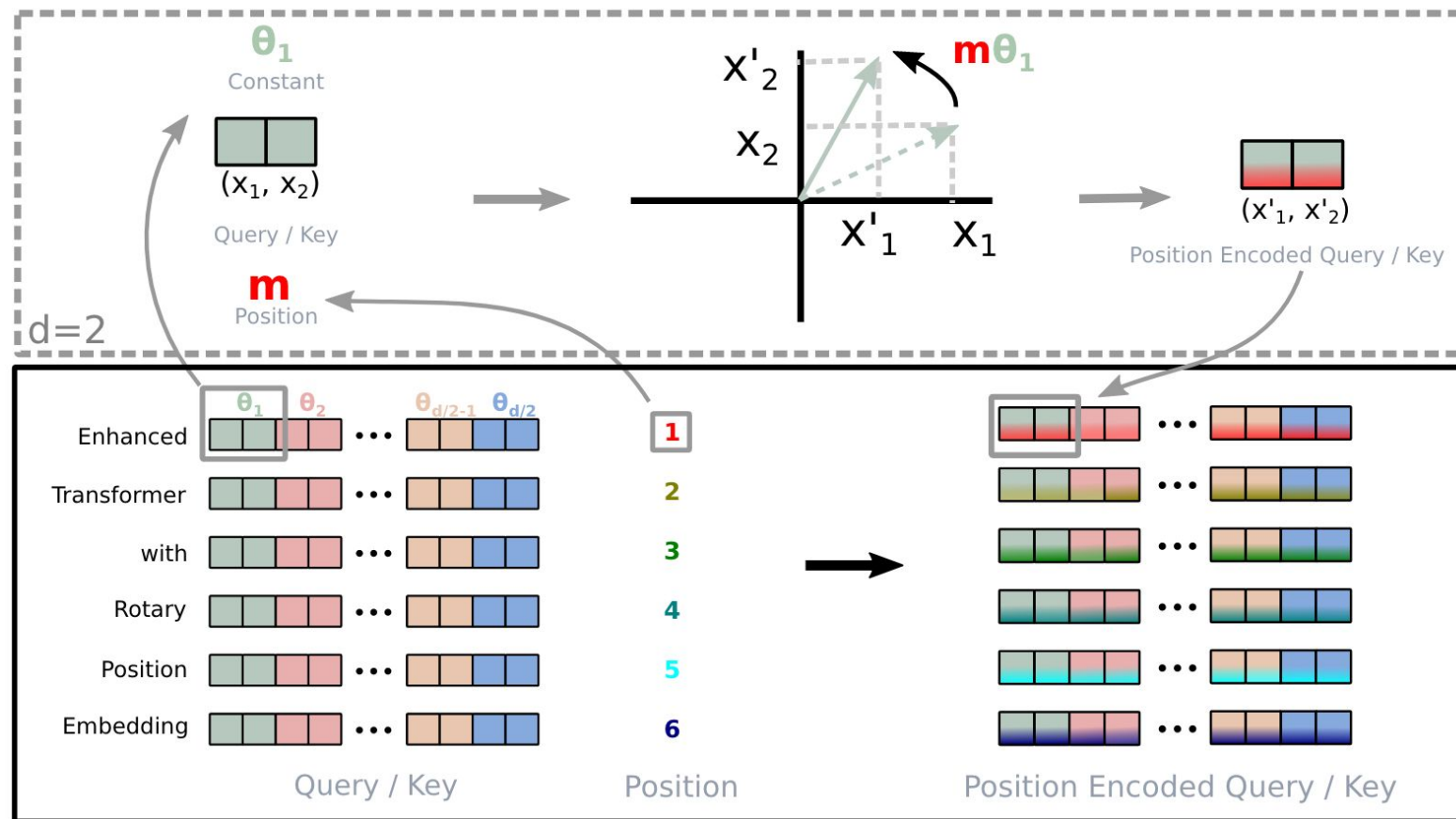
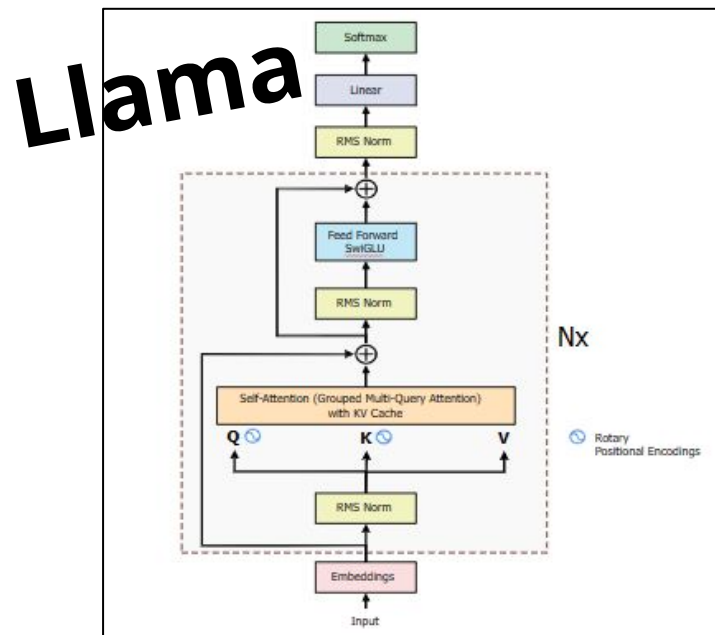
University of Illinois Urbana-Champaign

Outline

- Recap: **RoPE** (Rotary Position Embedding)
- **Block Attention:** fine-tuning and position re-coding
- Limitation and opportunities

RoPE

RoPE: Rotary position embedding



We perturb Q and K by **rotating** vectors (its angle *proportional to the position m*)

Block Attention (ICLR'25)

Block attention: cross attention only for query

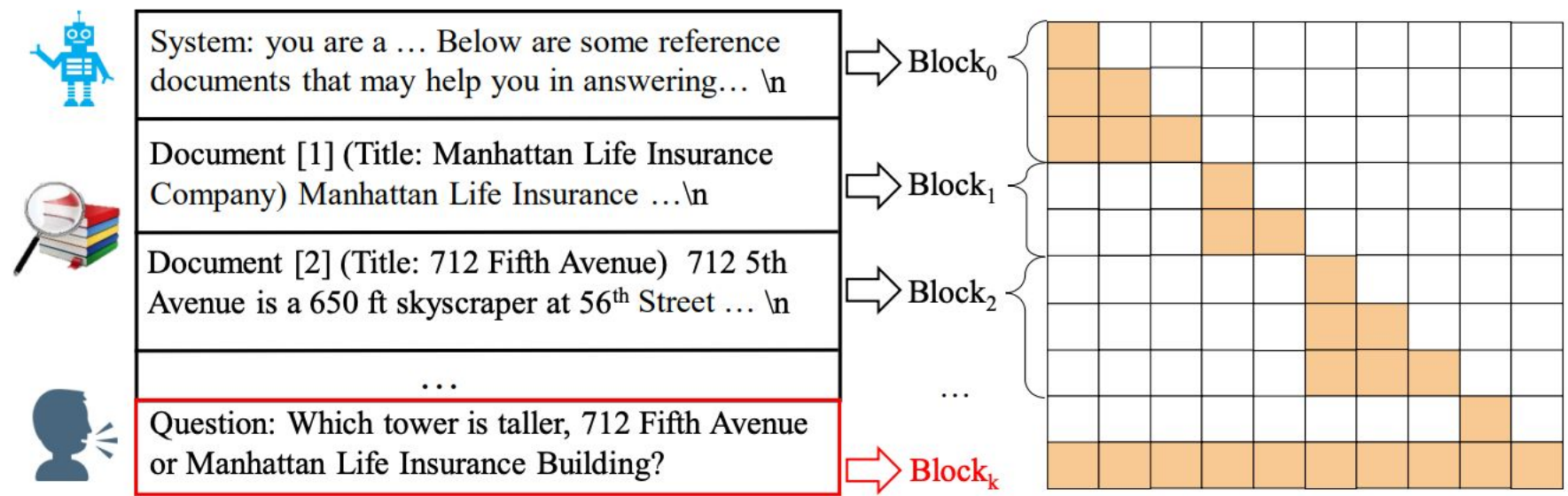
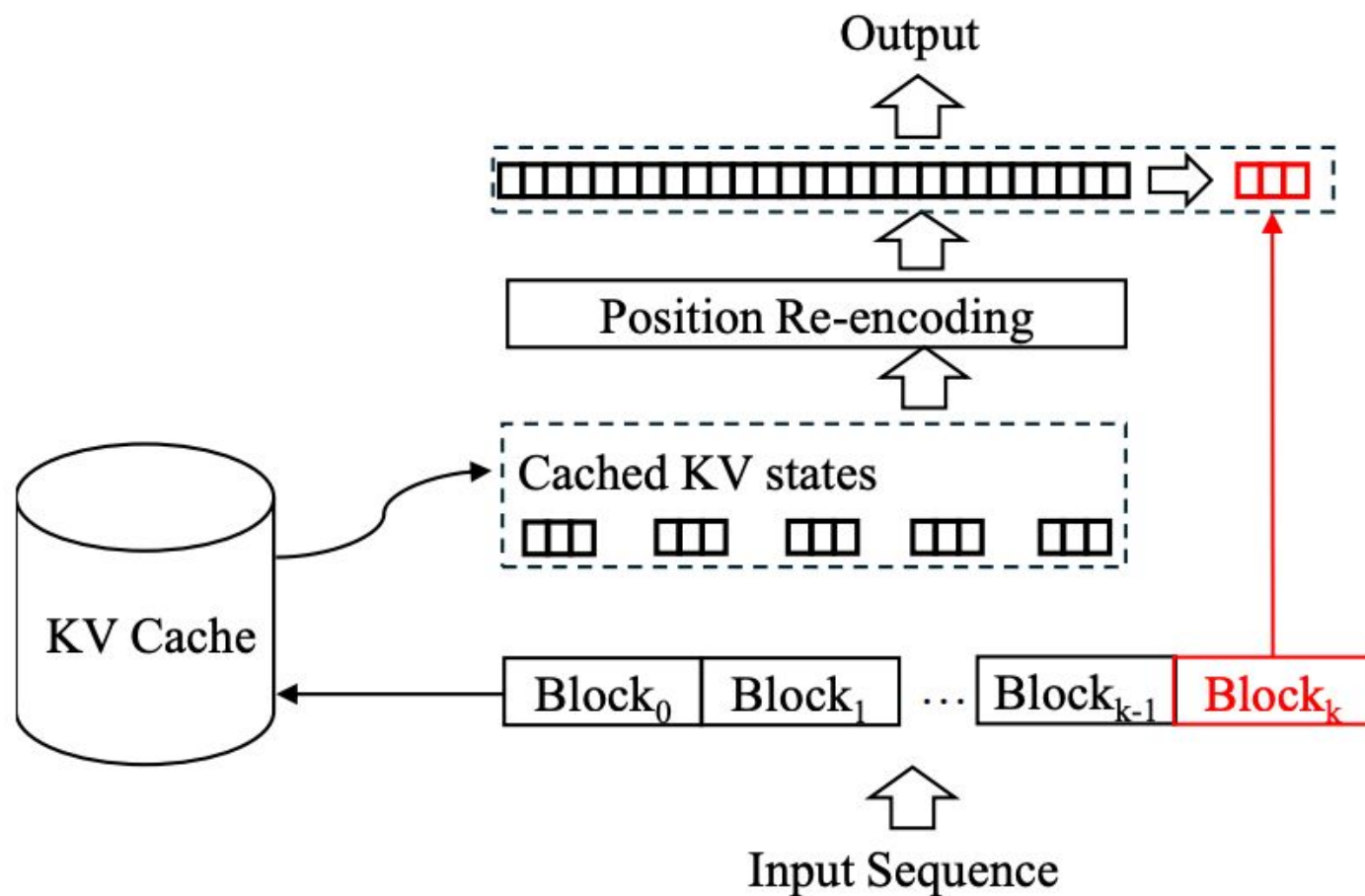


Figure 1: The Block-attention Masks

This allows re-using KV states of the **Blocks**

Block attention: pipeline



What is the **expected problem**? (based on the lessons from CacheBlend)

Problem: No cross-attention, poor performance

	Models	2wiki	HQA	NQ	TQA
baseline instruct model	<i>Tulu3-SFT</i>	62.0	68.4	58.6	75.7
fine-tuned for RAG: ideal model	<i>Tulu3-RAG</i>	73.2	74.8	61.5	75.8
a previous work	<i>Tulu3-RAG-Superposition</i>	30.1	32.3	35.9	58.9
another previous work	<i>Tulu3-RAG-promptCache</i>	32.4	31.6	44.4	61.8

Block attention proposes ***fine-tuning***, instead of *re-computing* cross attention (as in CacheBlend)

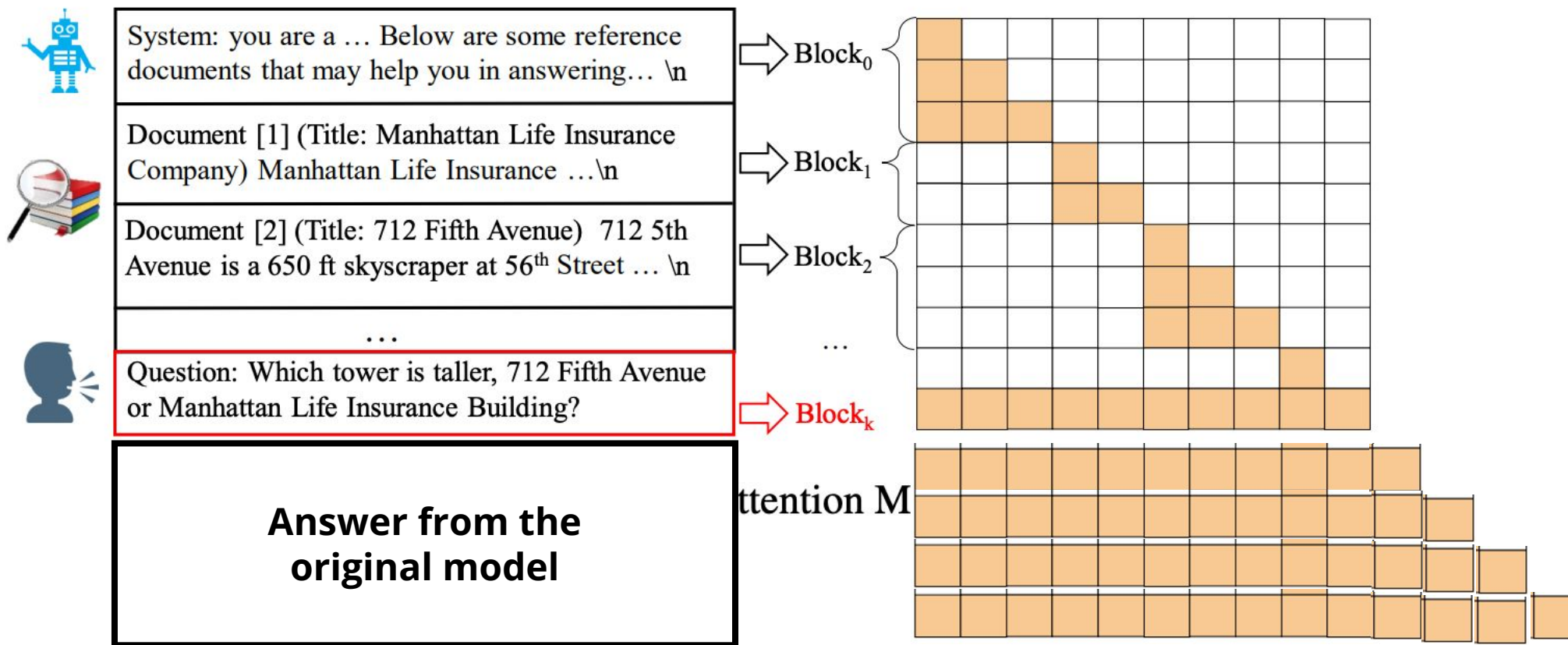
Fine-tuning and position re-encoding improves accuracy

	Models	2wiki	HQA	NQ	TQA
baseline instruct model	<i>Tulu3-SFT</i>	62.0	68.4	58.6	75.7
fine-tuned for RAG: ideal model	<i>Tulu3-RAG</i>	73.2	74.8	61.5	75.8
a previous work	<i>Tulu3-RAG-Superposition</i>	30.1	32.3	35.9	58.9
another previous work	<i>Tulu3-RAG-promptCache</i>	32.4	31.6	44.4	61.8
proposed model	<i>Tulu3-block-ft</i>	72.2	72.3	60.4	75.1
proposed model + cross-attention	<i>Tulu3-block-ft-full</i>	73.6	75.2	62.2	76.2
proposed model – positional re-encoding	<i>Tulu3-block-ft-w/o-pos</i>	68.9	69.9	59.2	74.4
proposed model – fine-tuning	<i>Tulu3-block-w/o-ft</i>	42.9	42.1	48.3	66.5

Table 1: Accuracy of different models on four RAG benchmarks.

Block attention proposes ***fine-tuning***, instead of *re-computing* cross attention (as in CacheBlend)

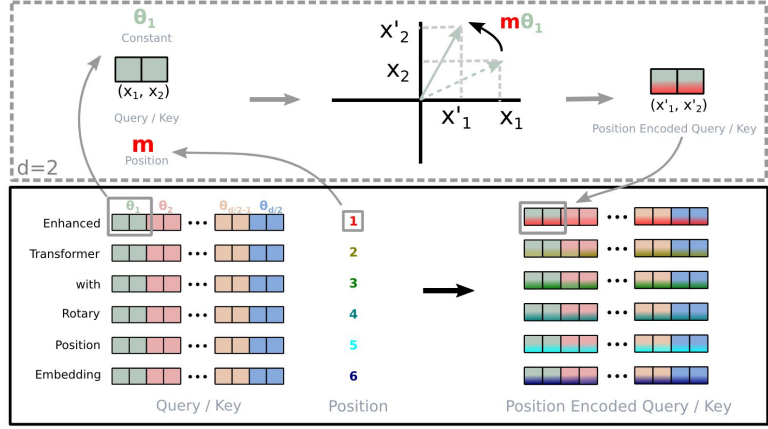
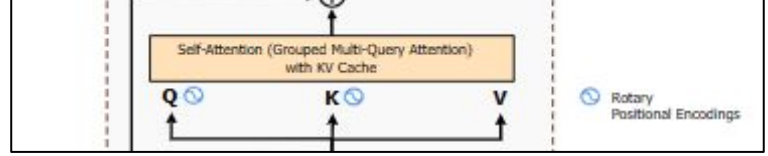
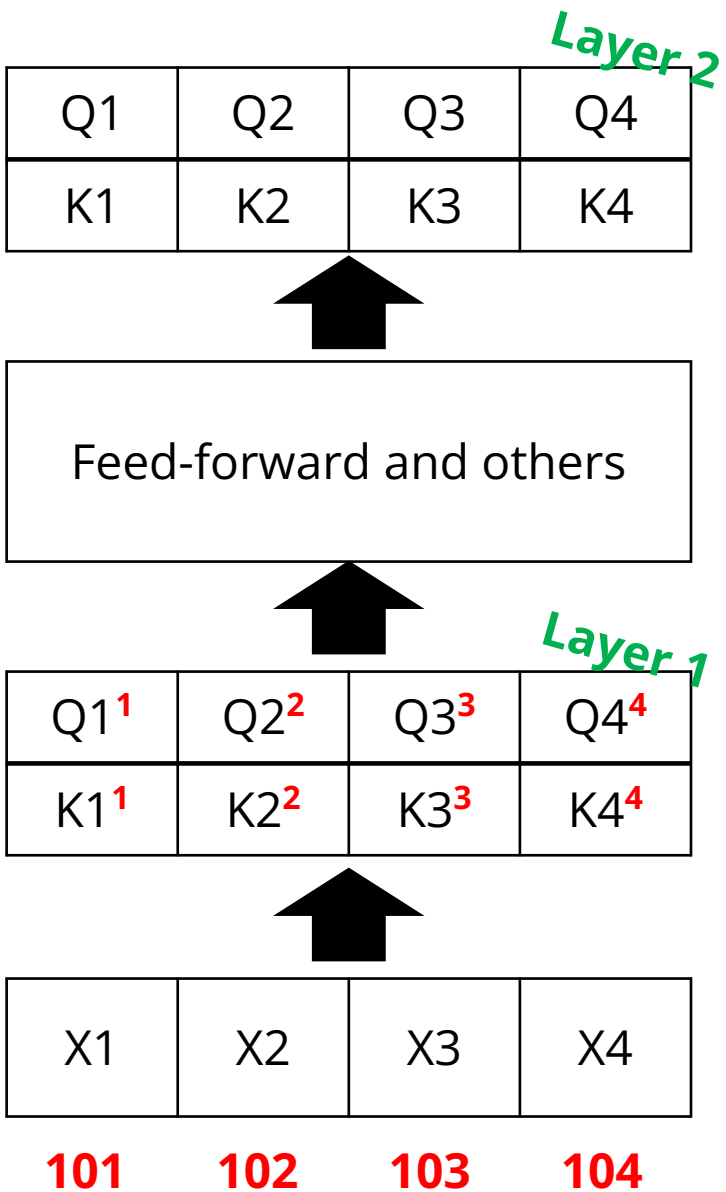
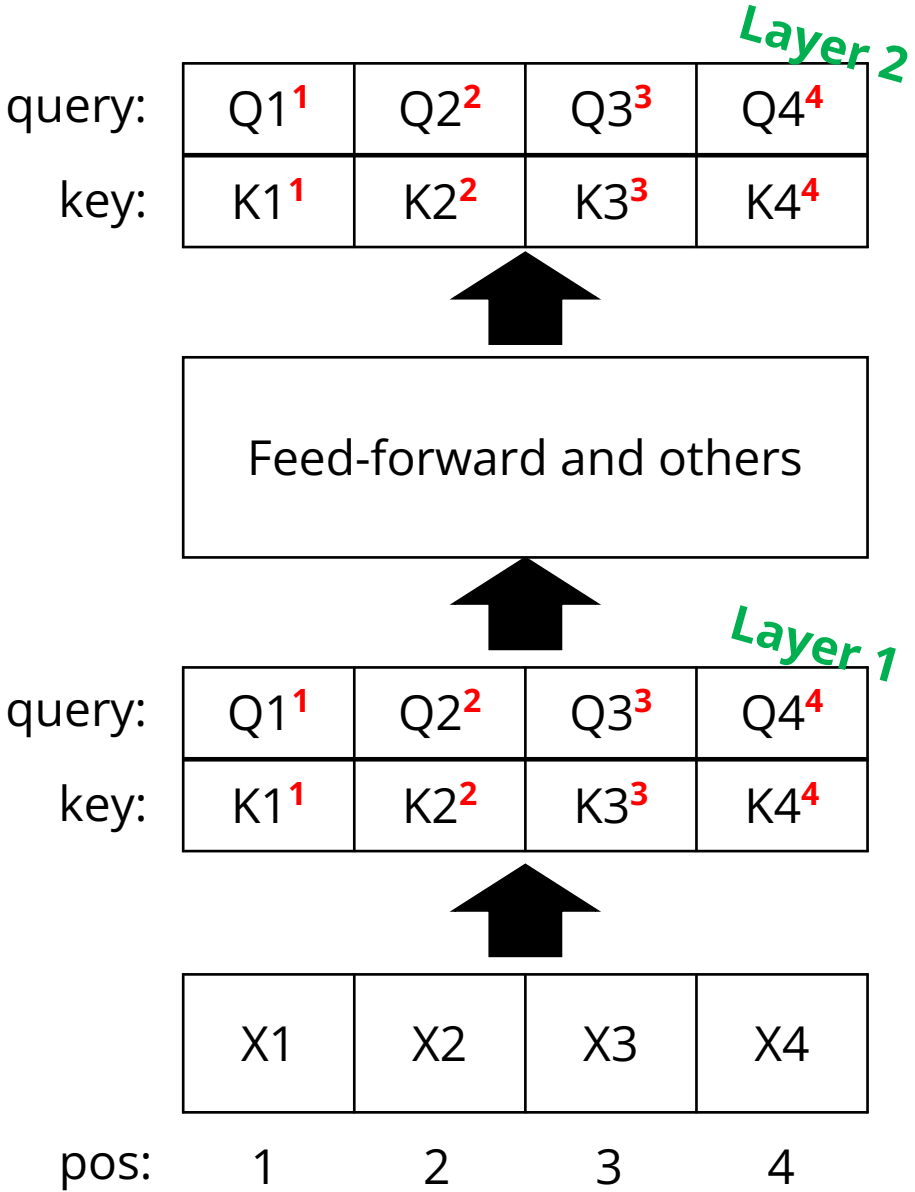
During fine-tuning, mask/erase cross-attention



We let model the produce the same answer, *without cross attention*

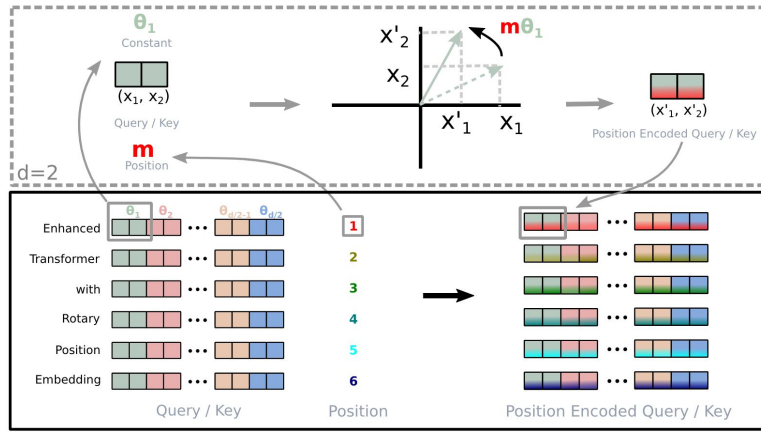
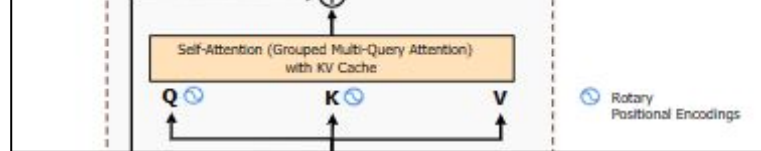
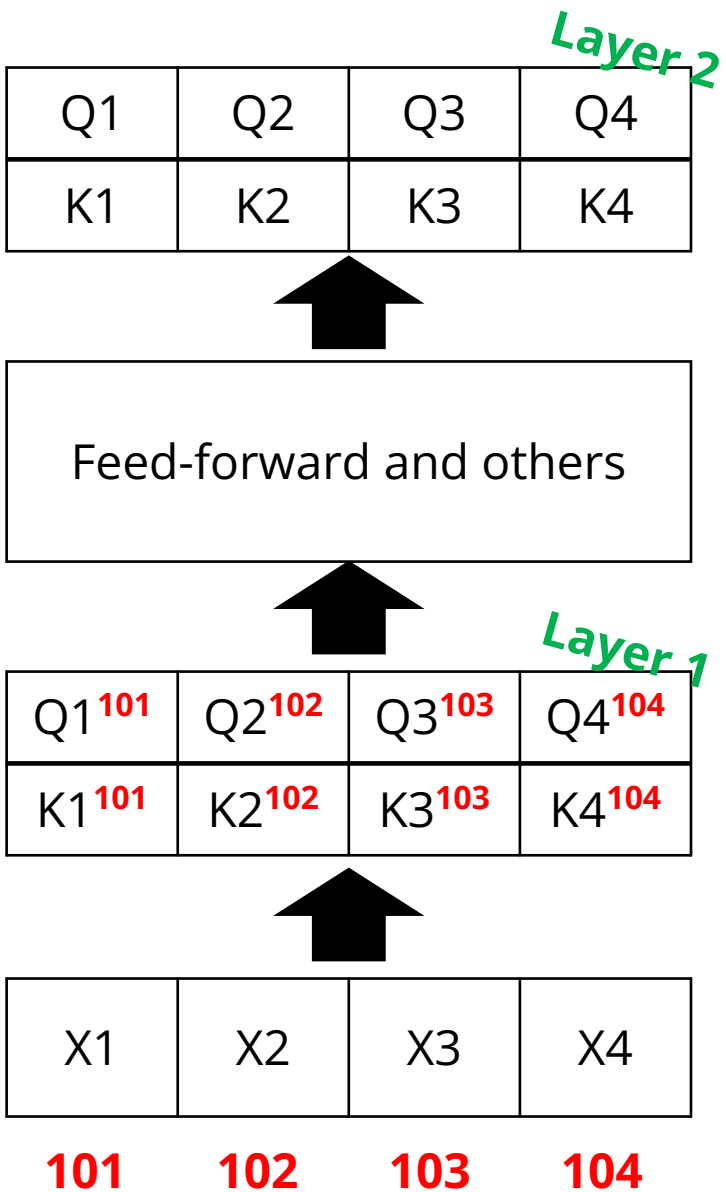
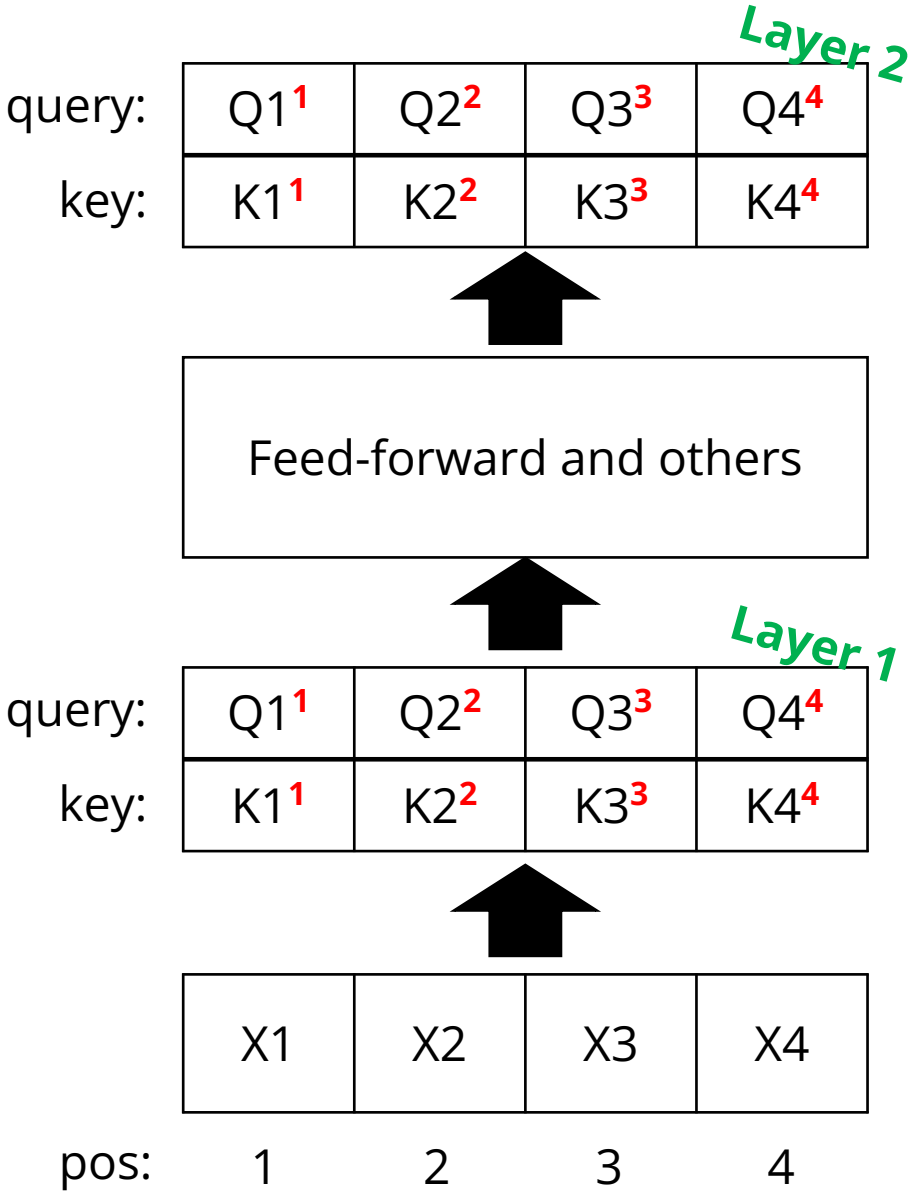
When will it work? or not work?

Position re-encoding



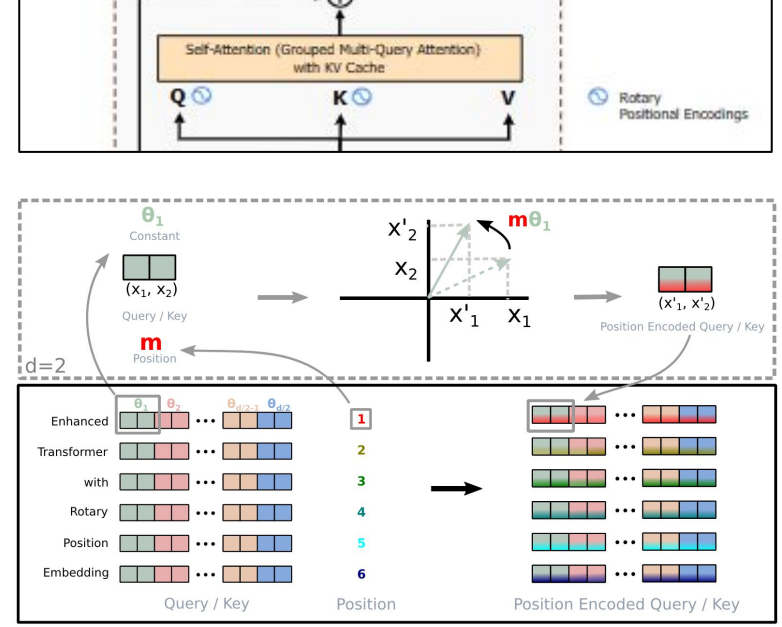
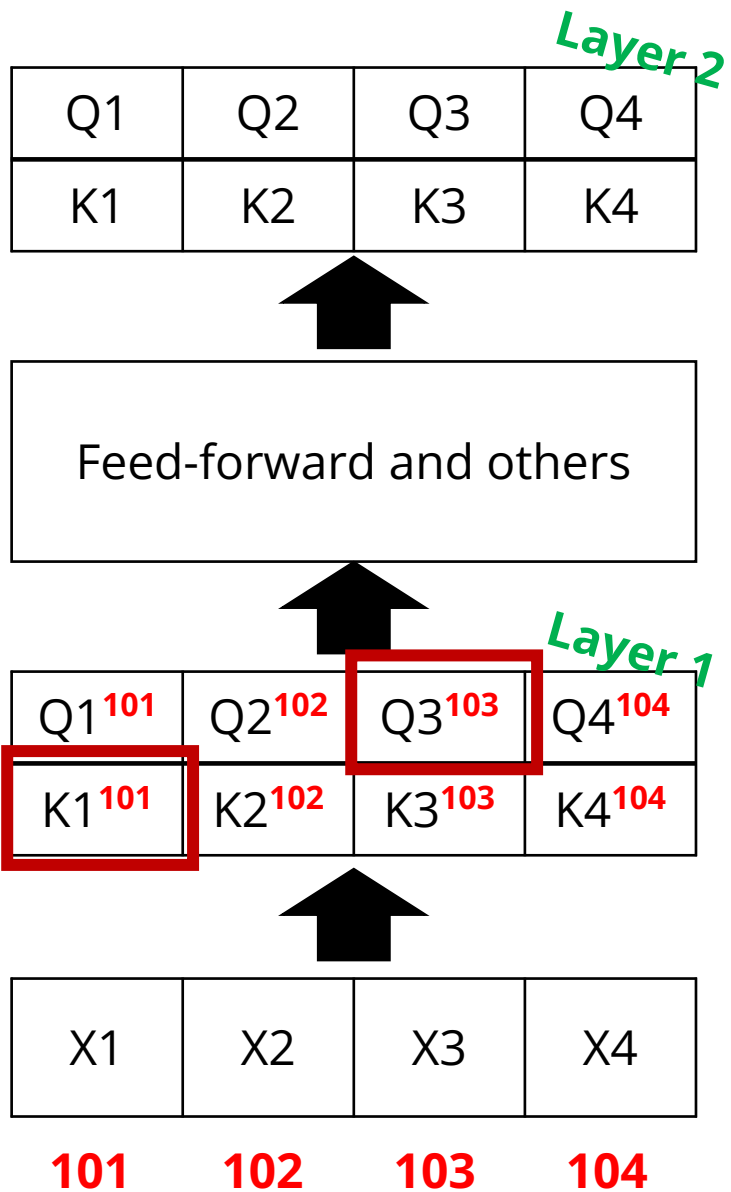
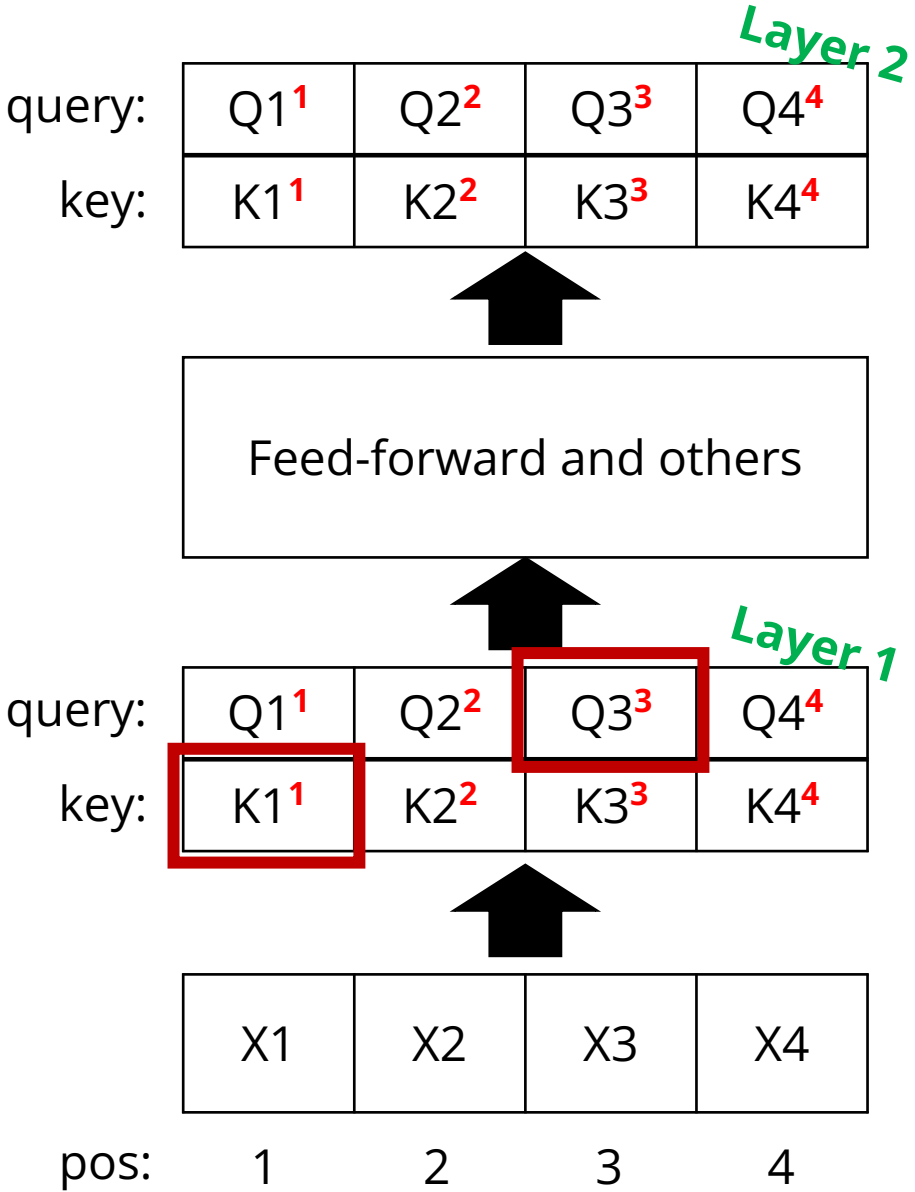
$$\text{sim}(K^{\theta_1}, Q^{\theta_2}) = f(K, Q, \theta_1 - \theta_2)$$

Position re-encoding



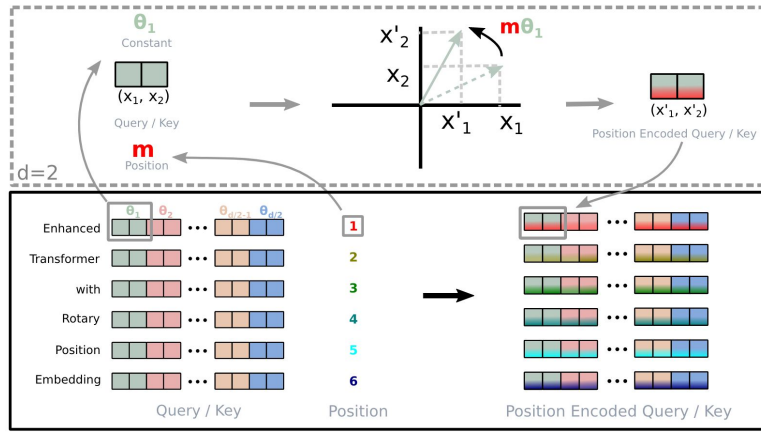
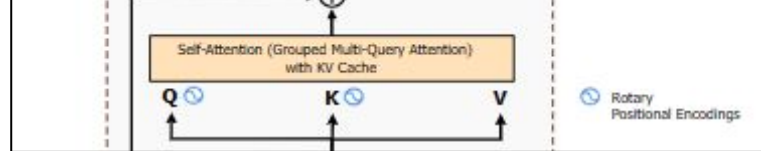
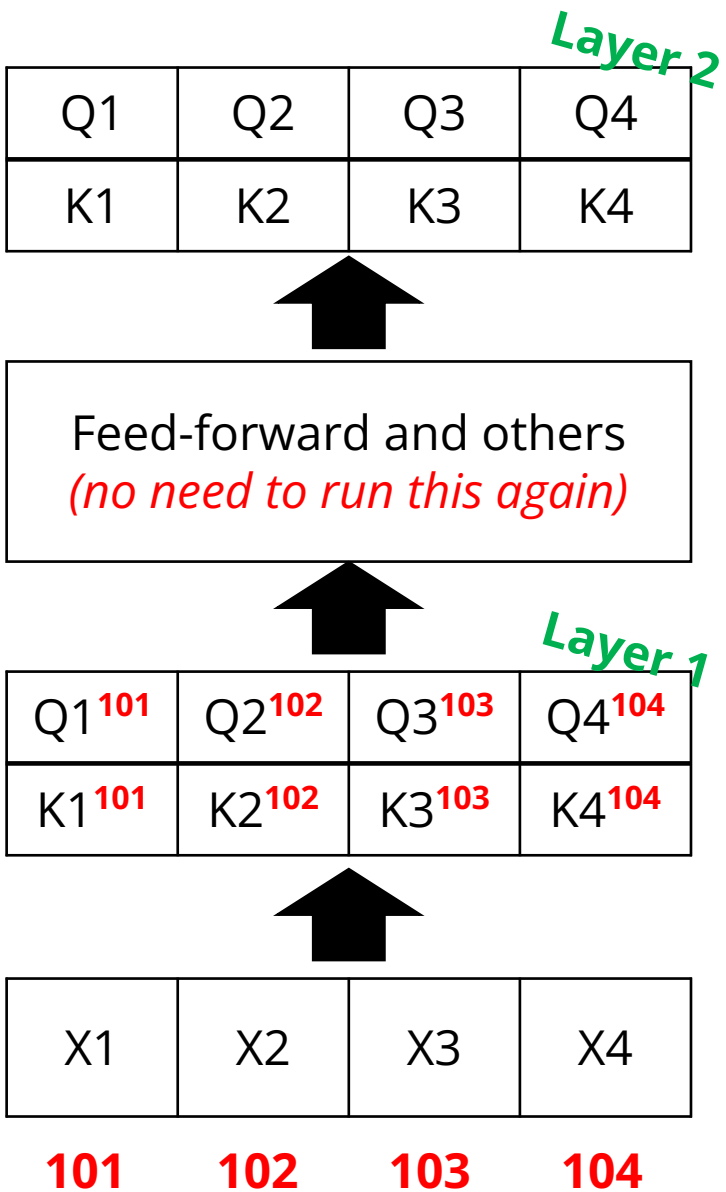
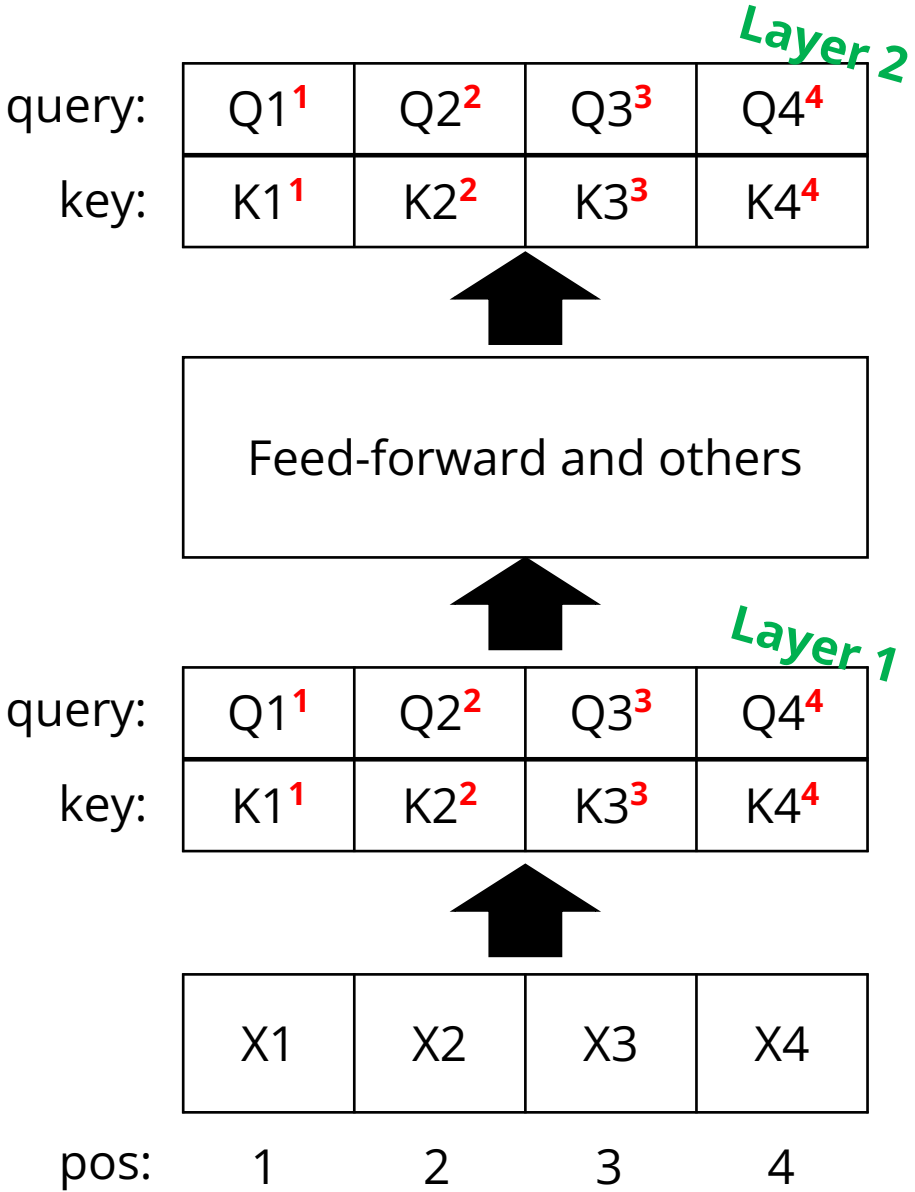
$$\text{sim}(K^{\theta_1}, Q^{\theta_2}) = f(K, Q, \theta_1 - \theta_2)$$

Position re-encoding



$$\text{sim}(K^{\theta_1}, Q^{\theta_2}) = f(K, Q, \theta_1 - \theta_2)$$

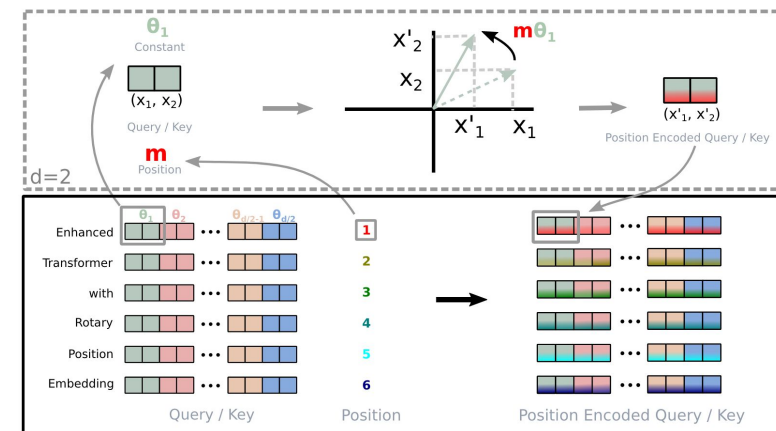
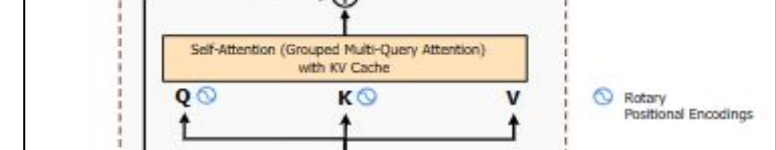
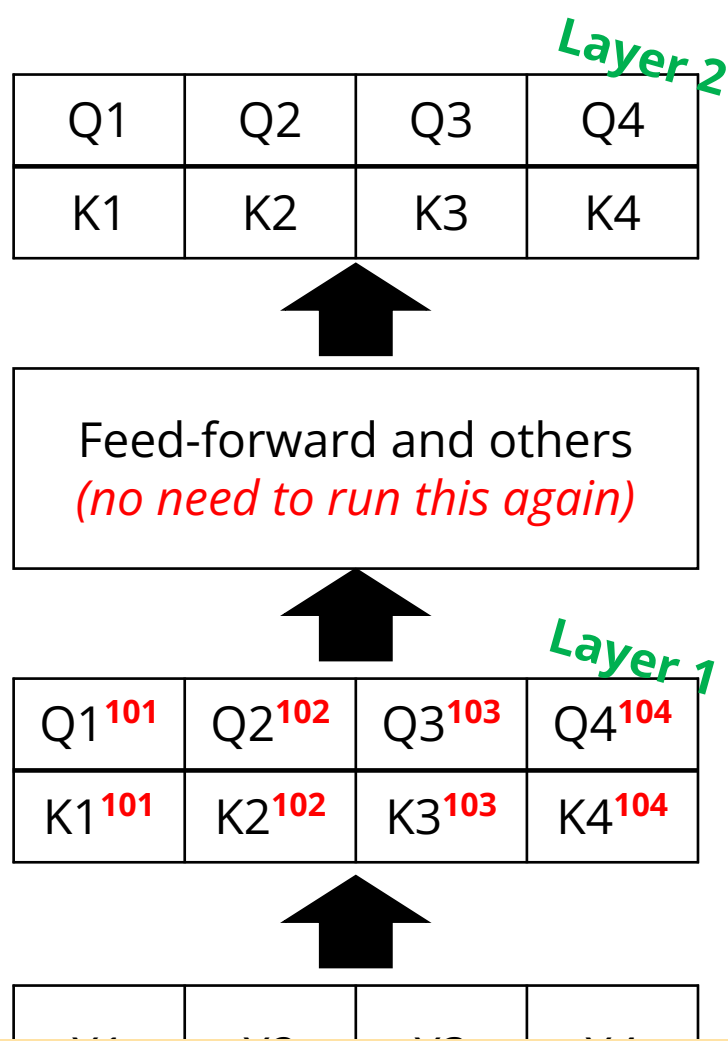
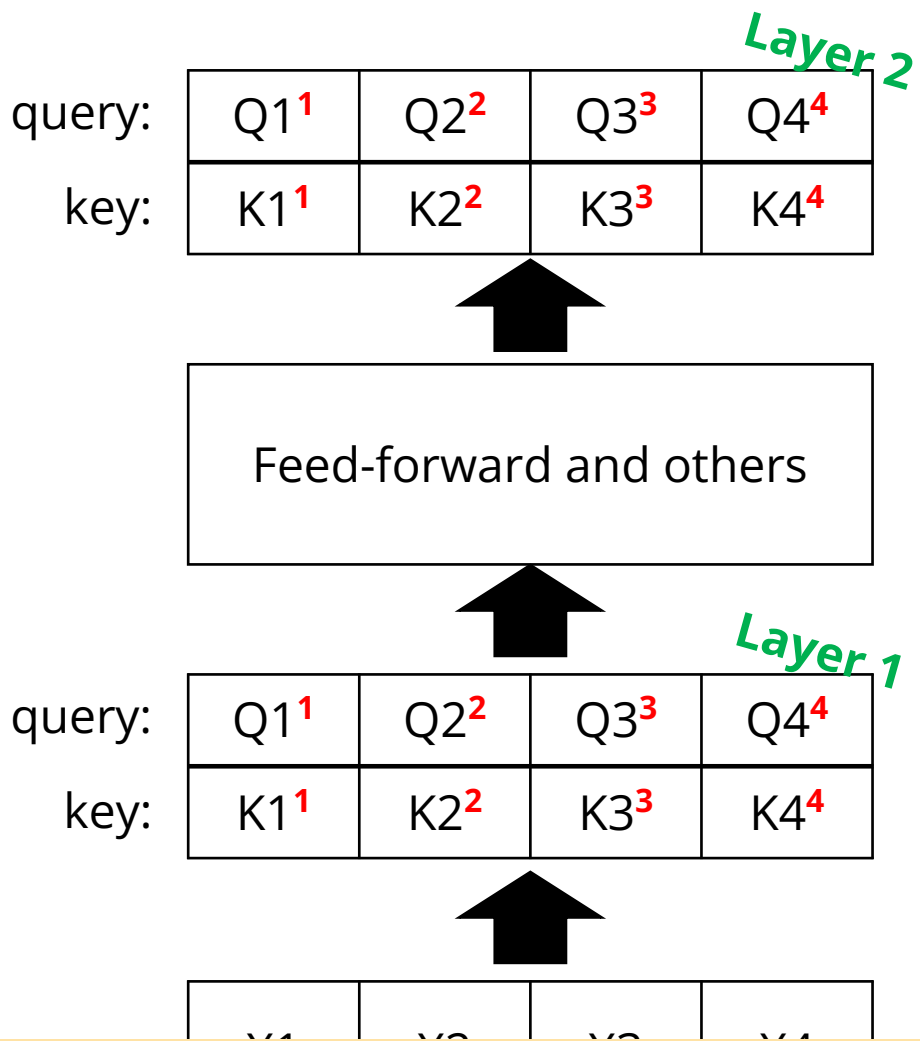
Position re-encoding



$$\text{sim}(K^{\theta_1}, Q^{\theta_2}) = f(K, Q, \theta_1 - \theta_2)$$

Shifting a block doesn't change output

Position re-encoding



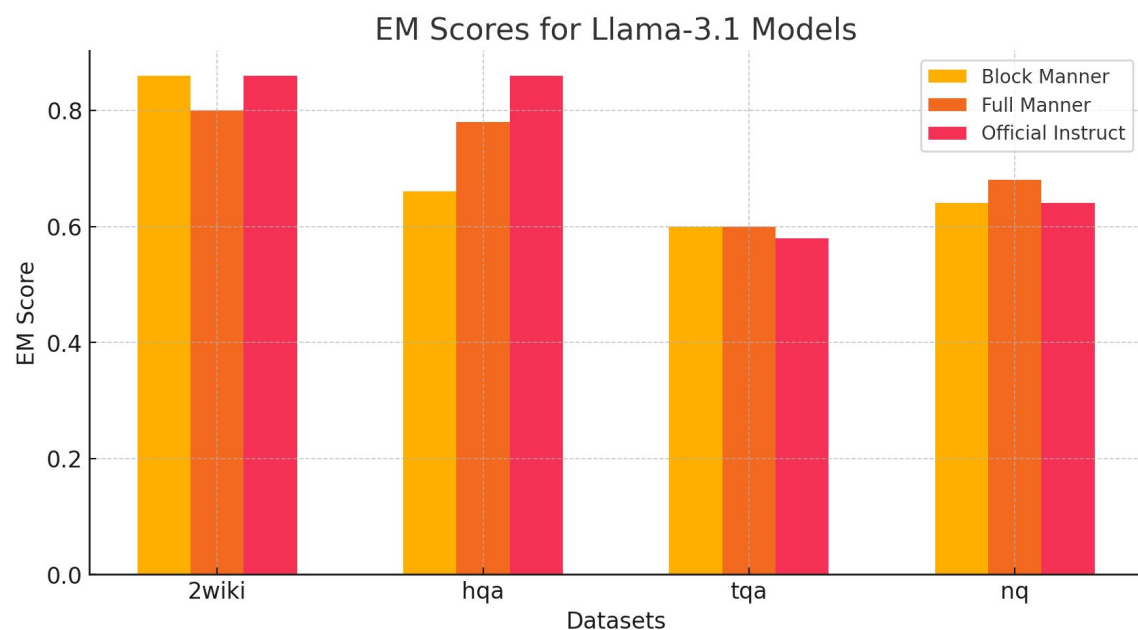
$$\text{sim}(K^{\theta_1}, Q^{\theta_2}) = f(K, Q, \theta_1 - \theta_2)$$

Shifting a block doesn't change output

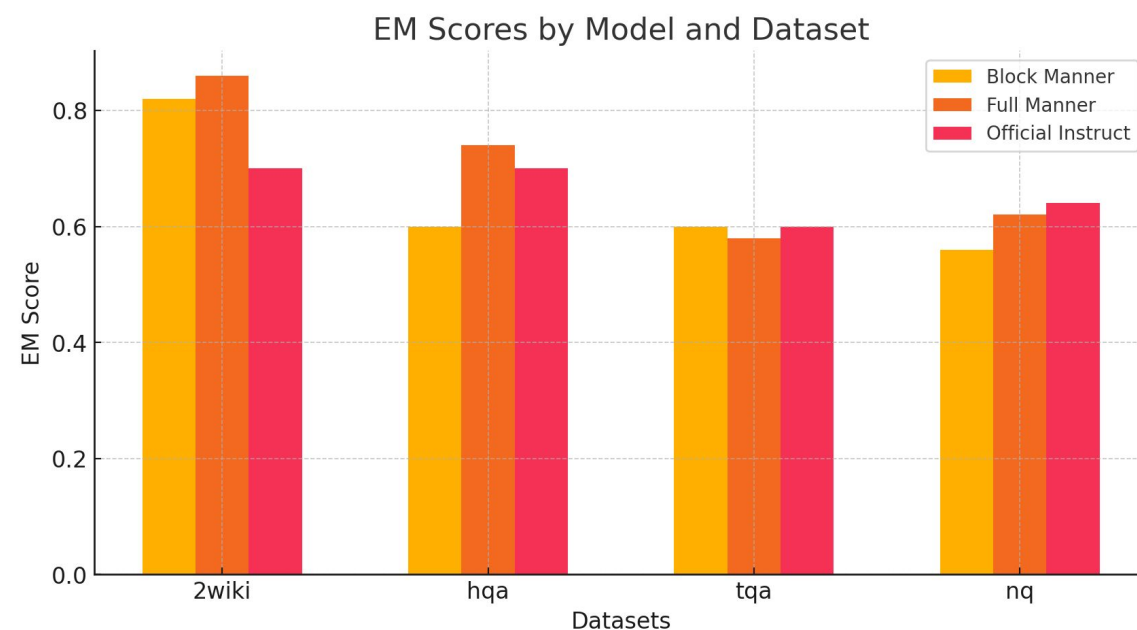
In actual implementation, we only need to **rotate keys** (not queries and values)

We could re-produce results

Llama 3.1- 8B

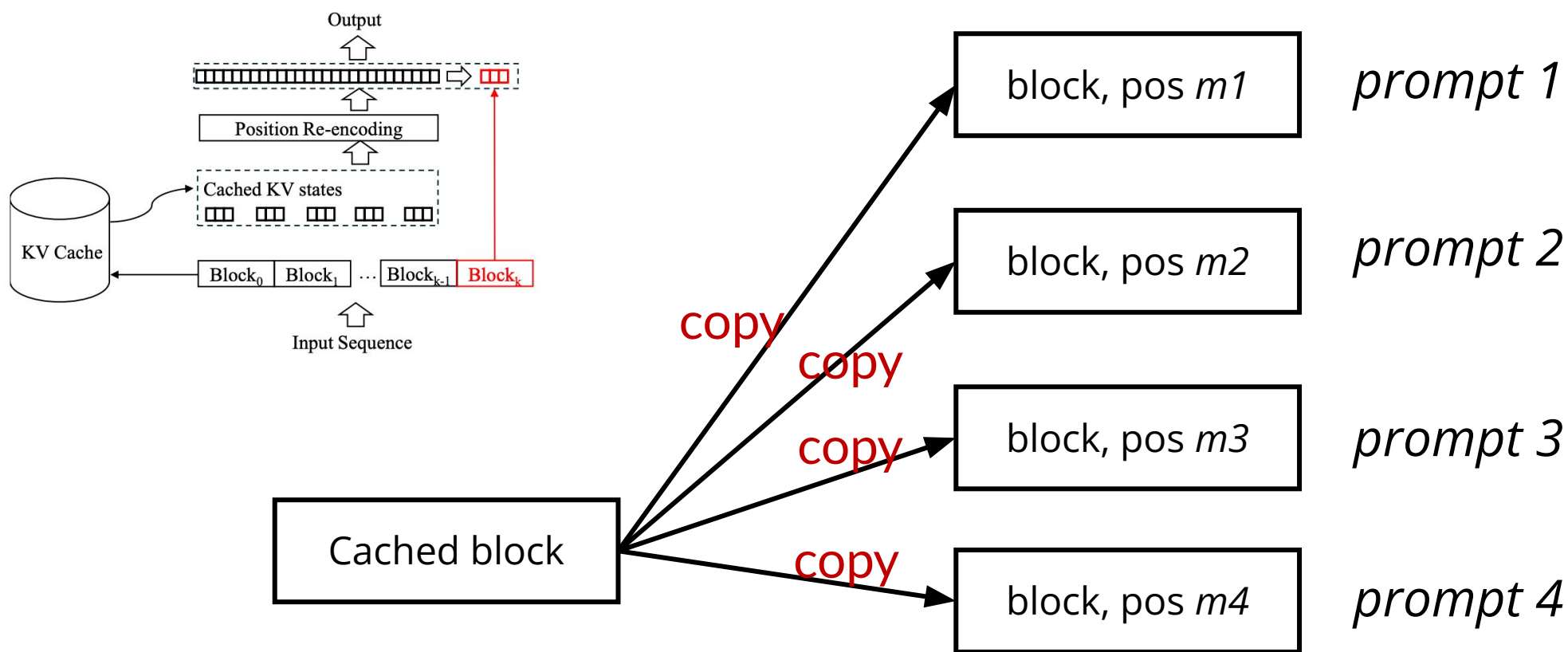


Llama 3.2 – 3B



Large model (70B parameters) could produce **similar accuracy without** fine-tuning

Limitation: **redundant** KV states for each cache hit



Our ongoing research: add positions ***on the fly***, inside GPU registers (not DRAM)

Block diffusion (ICLR'25)

Autoregression:  High quality  Arbitrary-length  KV caching  Not parallelizable

There

Diffusion:  Lower quality  Fixed-length  No KV caching  Parallelizable

Block Diffusion:  High quality  Arbitrary-length  KV caching  Parallelizable

Block Diffusion Internal

Large Language Diffusion Models (LLaDA) : Sequential Sampling

<div>Block 0</div> <div>[MASK] [MASK] [MASK] [UNK]</div> <div>Lily [MASK] [MASK] [UNK]</div> <div>Lily can [MASK] [UNK]</div> <div>Lily can run [UNK]</div>	<div>Block 1</div> <div>[MASK] 2 [MASK] [MASK]</div> <div>1 2 [MASK] [MASK]</div> <div>1 2 kilometers [MASK]</div> <div>1 2 kilometers per</div>	<div>Block 2</div> <div>hour [MASK] [MASK] [MASK]</div> <div>hour [MASK] [UNK] [MASK]</div> <div>hour for [UNK] [MASK]</div> <div>hour for [UNK] 4</div>	<div>Block 3</div> <div>hours [MASK] [MASK] [MASK]</div> <div>hours , [MASK] [MASK]</div> <div>hours , so [MASK]</div> <div>hours , so she</div>	<div>Block 4</div> <div>runs [MASK] [MASK] [MASK]</div> <div>runs [MASK] [MASK] of</div> <div>runs a [MASK] of</div> <div>runs a total of</div>	<div>Block 5</div> <div>[UNK] [MASK] [MASK] [MASK]</div> <div>[UNK] [MASK] 2 [MASK]</div> <div>[UNK] 1 2 [MASK]</div> <div>[UNK] 1 2 *</div>	<div>Block 6</div> <div>[UNK] [MASK] [MASK] [MASK]</div> <div>[UNK] 4 [MASK] [MASK]</div> <div>[UNK] 4 = [MASK]</div> <div>[UNK] 4 = [UNK]</div>	<div>Block 7</div> <div>[MASK] 8 [MASK] [MASK]</div> <div>4 8 [MASK] [MASK]</div> <div>4 8 kilometers [MASK]</div> <div>4 8 kilometers .</div>
<div>Block 8</div> <div>[UNK] [MASK] [MASK] [MASK]</div> <div>[UNK] After [MASK] [MASK]</div> <div>[UNK] After that [MASK]</div> <div>[UNK] After that ,</div>	<div>Block 9</div> <div>she [MASK] [MASK] [MASK]</div> <div>she runs [MASK] [MASK]</div> <div>she runs [UNK] [MASK]</div> <div>she runs [UNK] 6</div>	<div>Block 10</div> <div>kilometers [MASK] [MASK]</div> <div>kilometers [MASK] hour [MASK]</div> <div>kilometers per hour [MASK]</div> <div>kilometers per hour for</div>	<div>Block 11</div> <div>the [MASK] [MASK] [MASK]</div> <div>the remaining [MASK] [MASK]</div> <div>the remaining [UNK] [MASK]</div> <div>the remaining [UNK] 4</div>	<div>Block 12</div> <div>hours [MASK] [MASK] [MASK]</div> <div>hours , [MASK] [MASK]</div> <div>hours , [MASK] she</div> <div>hours , so she</div>	<div>Block 13</div> <div>runs [MASK] [MASK] [MASK]</div> <div>runs [MASK] [MASK] of</div> <div>runs [MASK] total of</div> <div>runs a total of</div>	<div>Block 14</div> <div>[UNK] [MASK] [MASK] [MASK]</div> <div>[UNK] [MASK] [MASK] [UNK]</div> <div>[UNK] [MASK] * [UNK]</div> <div>[UNK] 6 * [UNK]</div>	<div>Block 15</div> <div>[MASK] = [MASK] [MASK]</div> <div>[MASK] = [UNK] [MASK]</div> <div>4 = [UNK] [MASK]</div> <div>4 = [UNK] 2</div>
<div>Block 16</div> <div>[MASK] kilometers [MASK] [MASK]</div> <div>4 kilometers [MASK] [MASK]</div> <div>4 kilometers [MASK] [UNK]</div> <div>4 kilometers . [UNK]</div>	<div>Block 17</div> <div>Therefore [MASK] [MASK] [MASK]</div> <div>Therefore , [MASK] [MASK]</div> <div>Therefore , [MASK] can</div> <div>Therefore , Lily can</div>	<div>Block 18</div> <div>run [MASK] [MASK] [MASK]</div> <div>run [MASK] [MASK] of</div> <div>run [MASK] total of</div> <div>run a total of</div>	<div>Block 19</div> <div>[UNK] [MASK] [MASK] [MASK]</div> <div>[UNK] 4 [MASK] [MASK]</div> <div>[UNK] 4 8 [MASK]</div> <div>[UNK] 4 8 +</div>	<div>Block 20</div> <div>[UNK] [MASK] [MASK] [MASK]</div> <div>[UNK] [MASK] 4 [MASK]</div> <div>[UNK] 2 4 [MASK]</div> <div>[UNK] 2 4 =</div>	<div>Block 21</div> <div>[UNK] [MASK] [MASK] [MASK]</div> <div>[UNK] [MASK] 2 [MASK]</div> <div>[UNK] 7 2 [MASK]</div> <div>[UNK] 7 2 kilometers</div>	<div>Block 22</div> <div>in [MASK] [MASK] [MASK]</div> <div>in [MASK] 8 [MASK]</div> <div>in [UNK] 8 [MASK]</div> <div>in [UNK] 8 hours</div>	<div>Block 23</div> <div>. [MASK] [MASK] [MASK]</div> <div>. [UNK] [MASK] [MASK]</div> <div>. [UNK] The [MASK]</div> <div>. [UNK] The final</div>
<div>Block 24</div> <div>[MASK] [MASK] [UNK] [MASK]</div> <div>[MASK] [MASK] [UNK] 7</div> <div>[MASK] is [UNK] 7</div> <div>result is [UNK] 7</div>	<div>Block 25</div> <div>2 [MASK] [MASK] [MASK]</div> <div>2 [MASK] [UNK] [MASK]</div> <div>2 [MASK] [UNK] [UNK]</div> <div>2 [UNK] [UNK] [UNK]</div>	<div>Block 26</div> <div>[MASK] [UNK] [MASK] [MASK]</div> <div>[MASK] [UNK] [UNK] [MASK]</div> <div>[MASK] [UNK] [UNK] [UNK]</div> <div>[UNK] [UNK] [UNK] [UNK]</div>	<div>Block 27</div> <div>[MASK] [UNK] [MASK] [MASK]</div> <div>[MASK] [UNK] [MASK] [UNK]</div> <div>[MASK] [UNK] [UNK] [UNK]</div> <div>[UNK] [UNK] [UNK] [UNK]</div>	<div>Block 28</div> <div>[MASK] [MASK] [MASK] [UNK]</div> <div>[MASK] [MASK] [UNK] [UNK]</div> <div>[MASK] [UNK] [UNK] [UNK]</div> <div>[UNK] [UNK] [UNK] [UNK]</div>	<div>Block 29</div> <div>[MASK] [UNK] [MASK] [MASK]</div> <div>[MASK] [UNK] [MASK] [UNK]</div> <div>[MASK] [UNK] [UNK] [UNK]</div> <div>[UNK] [UNK] [UNK] [UNK]</div>	<div>Block 30</div> <div>[MASK] [UNK] [MASK] [MASK]</div> <div>[UNK] [UNK] [MASK] [MASK]</div> <div>[UNK] [UNK] [MASK] [UNK]</div> <div>[UNK] [UNK] [UNK] [UNK]</div>	<div>Block 31</div> <div>[MASK] [MASK] [MASK] [UNK]</div> <div>[MASK] [MASK] [UNK] [UNK]</div> <div>[MASK] [UNK] [UNK] [UNK]</div> <div>[UNK] [UNK] [UNK] [UNK]</div>

Prompt:

Lily can run 12 kilometers per hour for 4 hours. After that, she runs 6 kilometers per hour. How many kilometers can she run in 8 hours?

Final Answer:

Lily can run 1 2 kilometers per hour for 4 hours , so she runs a total of $1\ 2 * 4 = 4\ 8$ kilometers . After that , she runs 6 kilometers per hour for the remaining 4 hours , so she runs a total of $6 * 4 = 2\ 4$ kilometers . Therefore , Lily can run a total of $4\ 8 + 2\ 4 = 7\ 2$ kilometers in 8 hours . The final result is 7 2

Parallel Block Diffusion

Large Language Diffusion Models (LLaDA) : Diagonal Sampling

[illegible]

Summary

- RAG can be faster by bypassing cross-attention
- Block-attention can achieve good performance with fine-tuning
- Its position re-encoding can be memory inefficient

Questions?