

# LANGUAGE MODEL SIZES TO APR/2022

OR: WHILE YOU WERE SLEEPING,  
AI SIZES WERE EXPLODING

2018 →

2022 →

# Fancy Attention and Decoding

- Parameters
- AI lab/group
- Available
- Closed

Beeswarm/bubble plot, sizes linear to scale. Selected highlights only. Alan D. Thompson. April 2022. <https://lifearchitect.ai/>



LifeArchitect.ai/models



# What is the real bottle-neck of LLM?

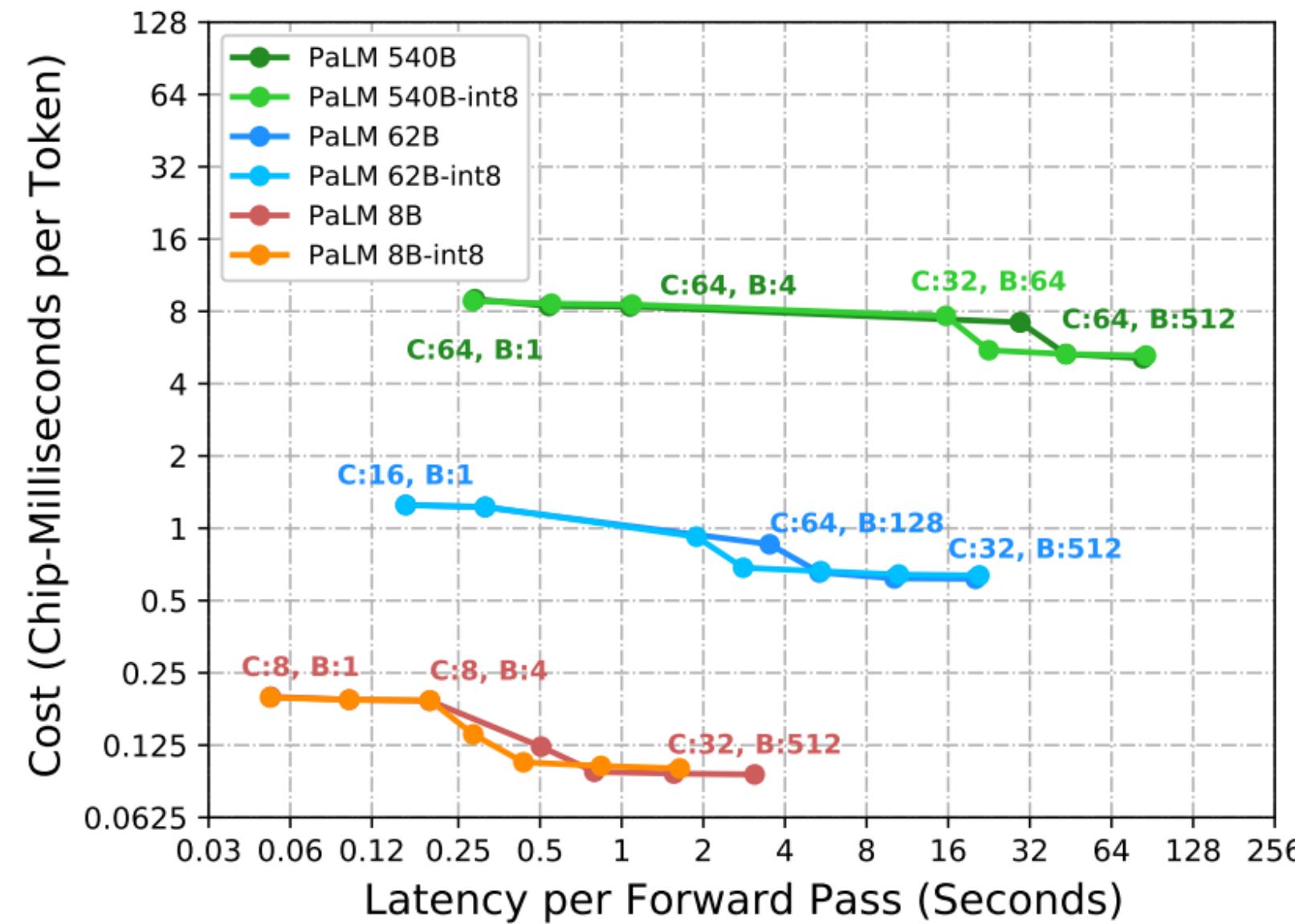
Run Llama-7B inference with llama.cpp

	Nvidia A100	Apple M2
computation	624 TOPS (90X)	7 TOPS
bandwidth	1935 GB/s(20X)	100 GB/s
speed	277 token/s (17X)	16 token/s

**Memory bandwidth** is the bottle-neck, not computation.

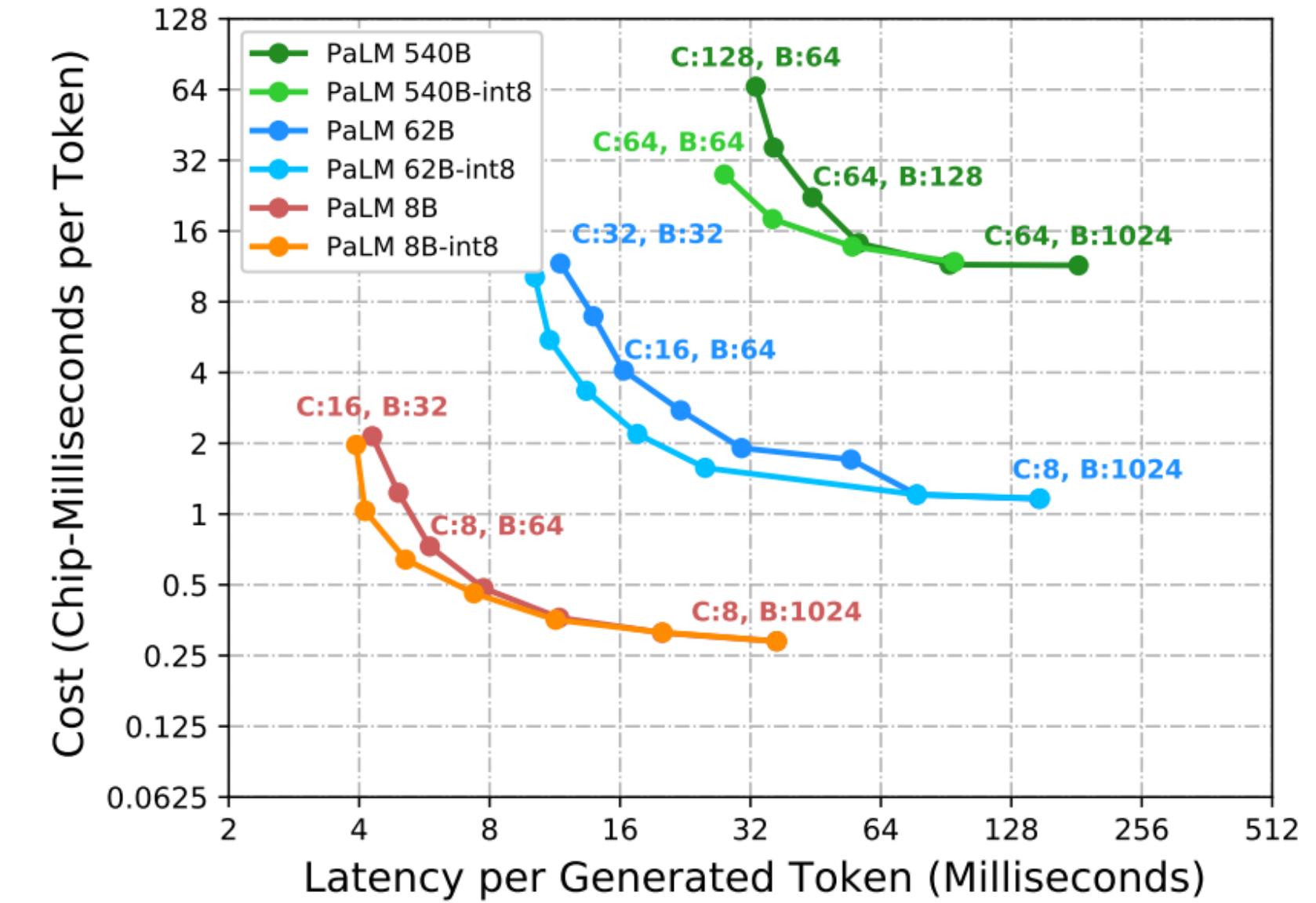
# Encoding & Decoding

Prefill Latency vs. Cost



encode / prefill / prefix

Decoding Latency vs. Cost



decode / auto-regressive

# Decoding

Transformer is useful in  
prefix

decode

Transformer is useful in \_\_\_\_\_

natural

Transformer is useful in natural \_\_\_\_\_

language

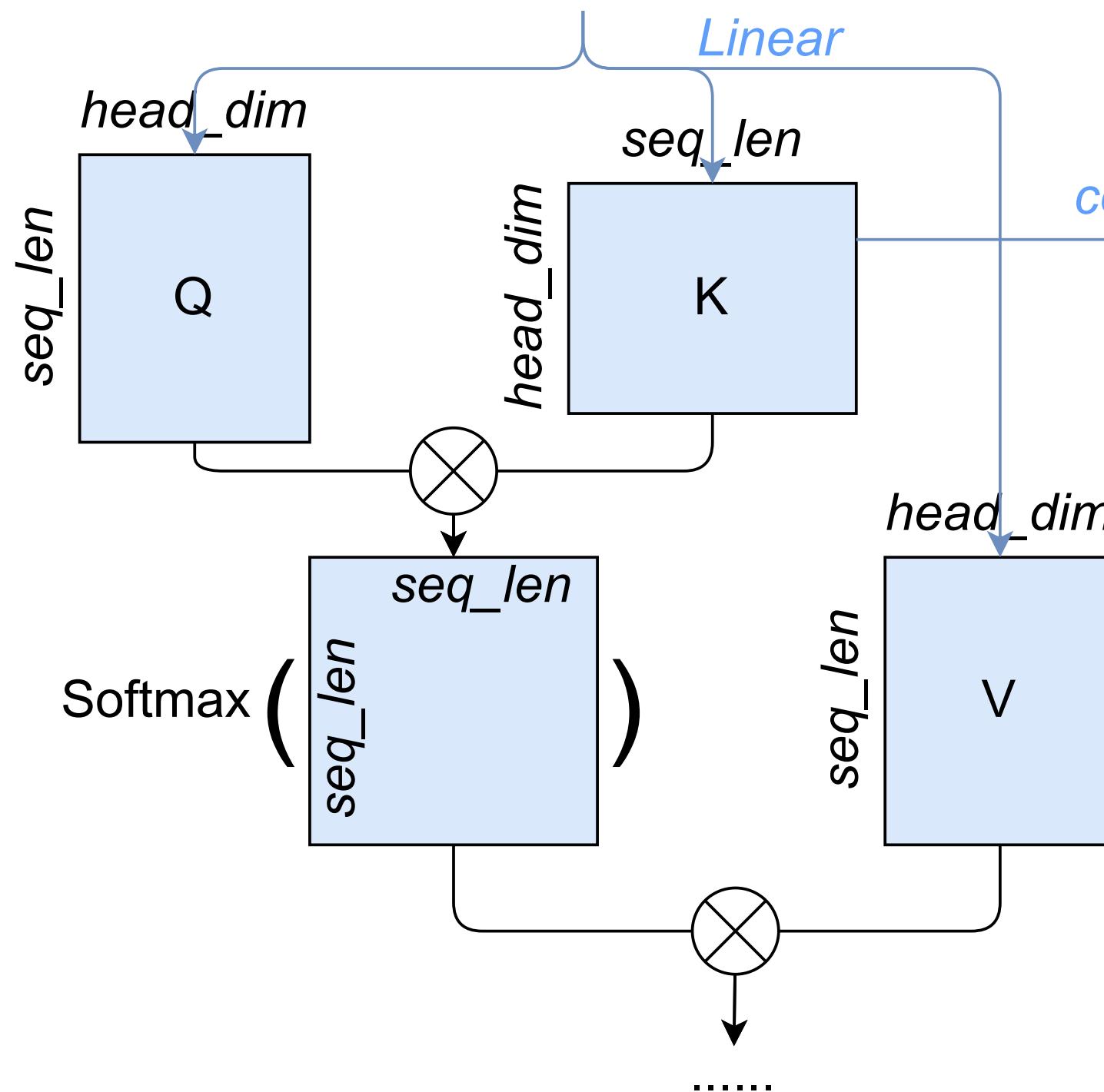
Transformer is useful in natural language \_\_\_\_\_

processing

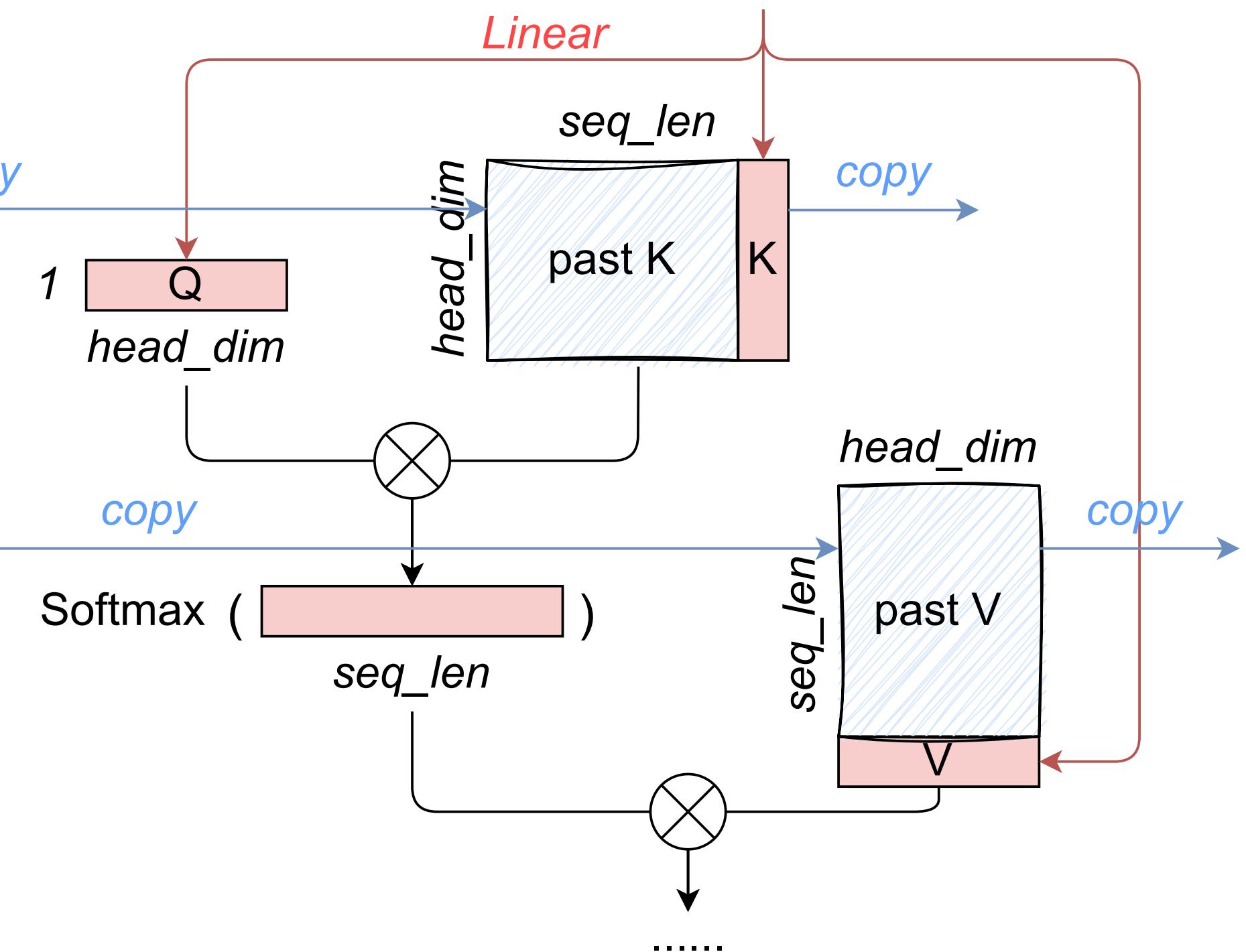
Notice something redundant?

# KV cache

Transformer is useful in



Transformer is useful in **natural**



## KV cache

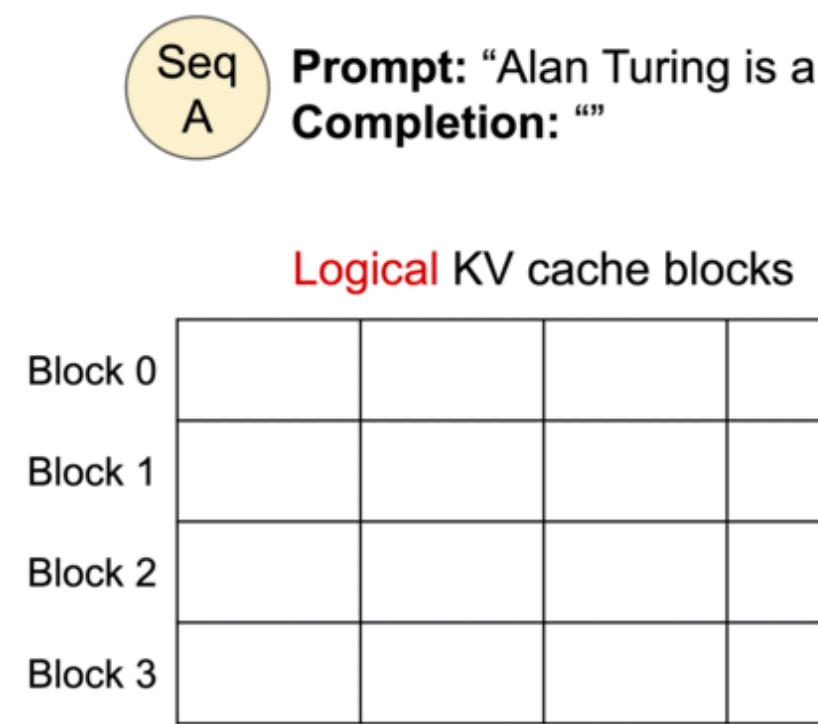
1. KV cache is not *the cache*
2. Q (query) doesn't need cache (always 1 tensor)

# vLLM & PagedAttention

KV cache length is unpredictable, fragmentation and over-reservation cause 60% – 80% memory redundancy

Adopt page table from OS to manage KV cache

0. Before generation.



**Block table**

Physical block no.	# Filled slots
-	-
-	-
-	-
-	-

**Physical KV cache blocks**

Block 0				
Block 1				
Block 2				
Block 3				
Block 4				
Block 5				
Block 6				
Block 7				

# vLLM throughput

# FlashAttention

Attention computation can be divide-and-conquer

Divide Q K V, compute each block on SRAM

# FlashAttention speed

Table 18: Forward pass runtime (ms) of various exact/approximate/sparse attention mechanisms by sequence length. Best in **bold**, second best underlined.

Attention Method	128	256	512	1024	2048	4096	8192	16384	32768	65536
<b>PyTorch Attention</b>	<u>0.21</u>	<u>0.22</u>	0.43	1.27	4.32	16.47	67.77	-	-	-
	0.24	0.26	<u>0.42</u>	1.33	4.28	-	-	-	-	-
<b>Local Attention</b>	1.77	2.82	5.01	9.74	20.03	41.11	87.39	192.40	-	-
	0.48	0.57	0.80	1.90	5.76	11.56	23.13	46.65	94.74	-
	0.46	0.36	0.45	<b>0.50</b>	<u>1.09</u>	<u>2.09</u>	<u>4.01</u>	<u>7.90</u>	<u>15.70</u>	<u>35.40</u>
	1.94	1.96	3.01	5.69	11.26	22.23	44.21	88.22	-	-
	1.21	1.34	1.34	3.31	11.01	21.71	43.27	86.32	172.85	-
<b>Block Sparse</b>	0.96	1.04	1.66	2.16	5.41	16.15	-	-	-	-
	0.99	0.98	0.99	1.56	4.79	11.07	32.98	-	-	-
	0.96	1.02	1.02	1.48	5.05	11.59	34.16	-	-	-
<b>FLASHATTENTION</b>	<b>0.08</b>	<b>0.09</b>	<b>0.18</b>	0.68	2.40	8.42	33.54	134.03	535.95	2147.05
<b>Block-Sparse FLASHATTENTION</b>	0.56	0.52	0.63	<u>0.65</u>	<b>0.61</b>	<b>0.96</b>	<b>1.69</b>	<b>3.02</b>	<b>5.69</b>	<b>11.77</b>

# Early Exiting (Again)

```
? ? ? ? ? ? ?
```

```
? ? ? ? ? ? ?
```

# CALM

## Confident Adaptive Language Modeling. NeurIPS 2022

Most tokens can early exit with less than 4 layers!

---

$Z_{\text{test}}$  : South Africa-born Grant Elliott hit match-winning 84 not out in semi-final . Black Caps reached first World Cup final with Elliott's penultimate ball six .  
Elliott, 36, had not played international cricket for 14 months when picked . Win is vindication for the attacking brand played under Brendon McCullum . New Zealand play the winner of the semi-final between Australia or India .

---

$y_{\text{full}}$  : Grant Elliott hit a six to put New Zealand through to the World Cup final . The 36-year-old was born in South Africa but a naturalised Kiwi . Elliott will surely never play another innings like his 84 . New Zealand will take on either Australia or India in the final on Sunday .

---

$y_{\text{early}}^{(1)}$  : \_\_Grant\_Elliott\_hit\_a\_six\_to\_put\_New\_Zealand\_through\_to\_the\_World\_Cup\_final\_\_.\_\_Elliott\_was\_born\_in\_South\_Africa\_but\_\_  
a\_naturalised\_Kiwi\_\_.\_\_The\_36-year-old\_will\_surely\_never\_play\_another\_innings\_like\_his\_unbeaten\_84\_\_.  
\_\_New\_Zealand\_will\_now\_take\_on\_either\_Australia\_or\_India\_in\_the\_final\_on\_Sunday\_\_.<EOS>

---

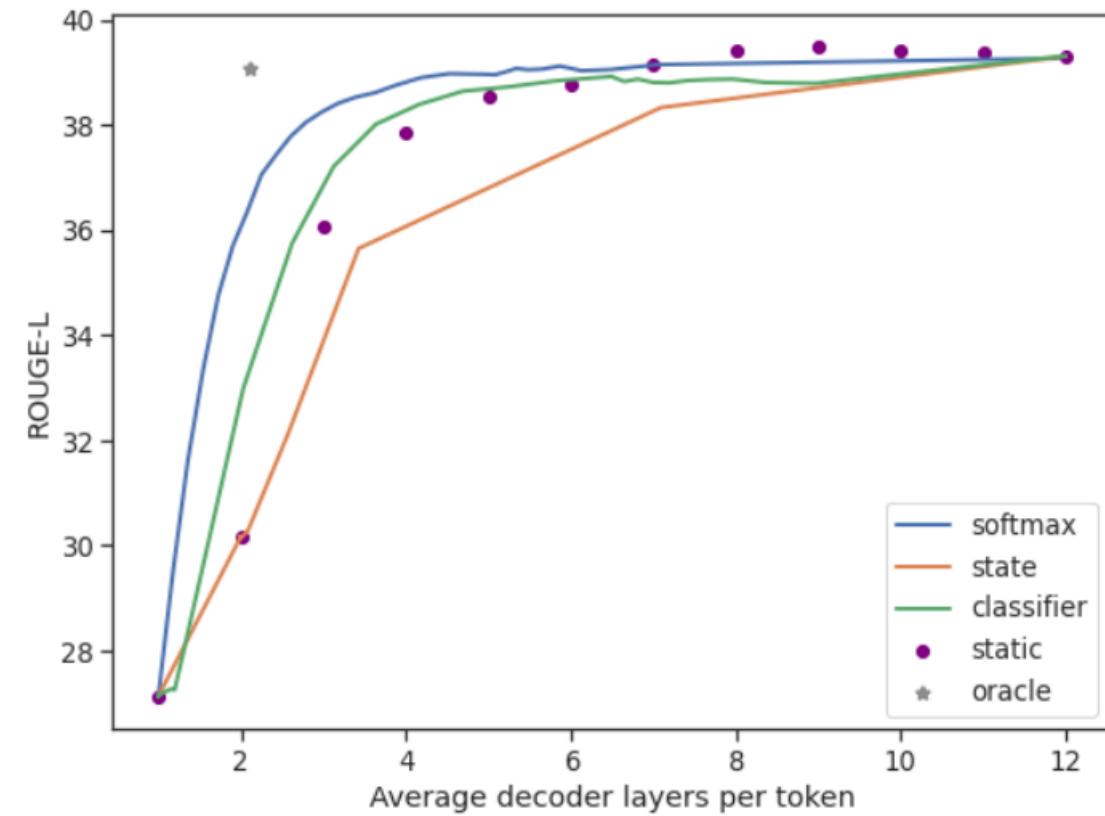
$y_{\text{early}}^{(2)}$  : \_\_Grant\_Elliott\_hit\_84\_in\_the\_Black\_Caps\_chase\_\_.\_\_New\_Zealand\_reached\_the\_World\_Cup\_final\_\_.\_\_Elliott\_was\_born\_in\_So\_\_  
uth\_Africa\_but\_a\_naturalised\_Kiwi\_\_.\_\_Elliott\_will\_surely\_never\_play\_another\_innings\_like\_his\_unbeaten\_84\_\_.<EOS>

---

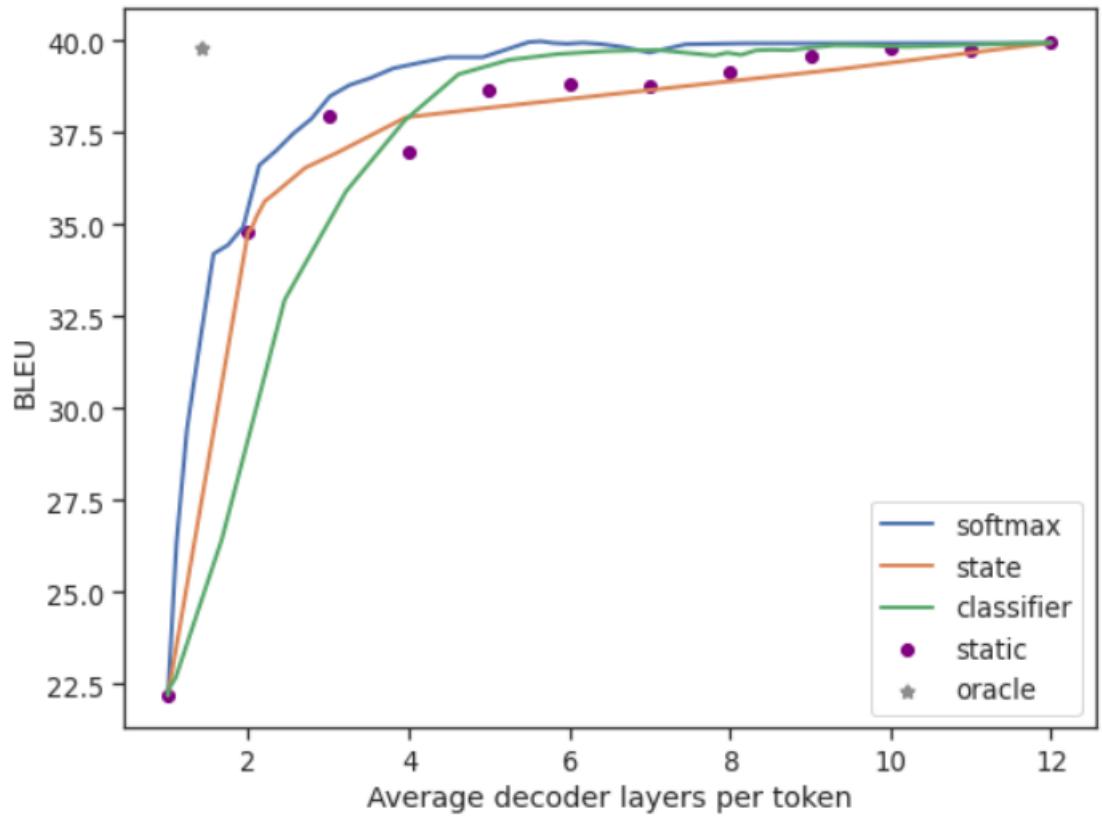
	$D(Y_{\text{early}}, Y_{\text{full}})$	$R_{\text{early}} - R_{\text{full}}$	Average layers	Speedup
$Y_{\text{early}}^{(1)}$	0.02	0.01	2.1	x 2.9
$Y_{\text{early}}^{(2)}$	0.33	-0.3	1.9	x 3.6

Exit layer — color mapping: 12345678  
 $D$  and  $R$  are computed with ROUGE-L

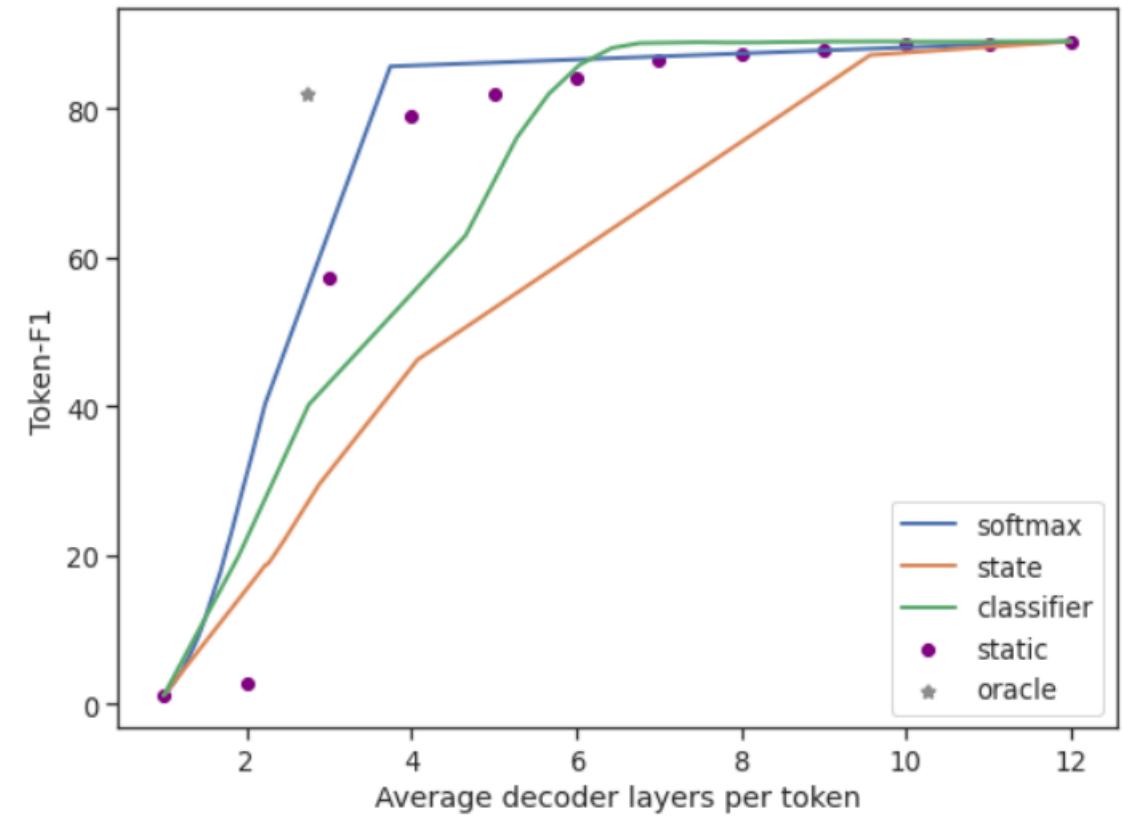
# Exit point – Softmax still reigns



(a) CNN/DM



(b) WMT



(c) SQuAD

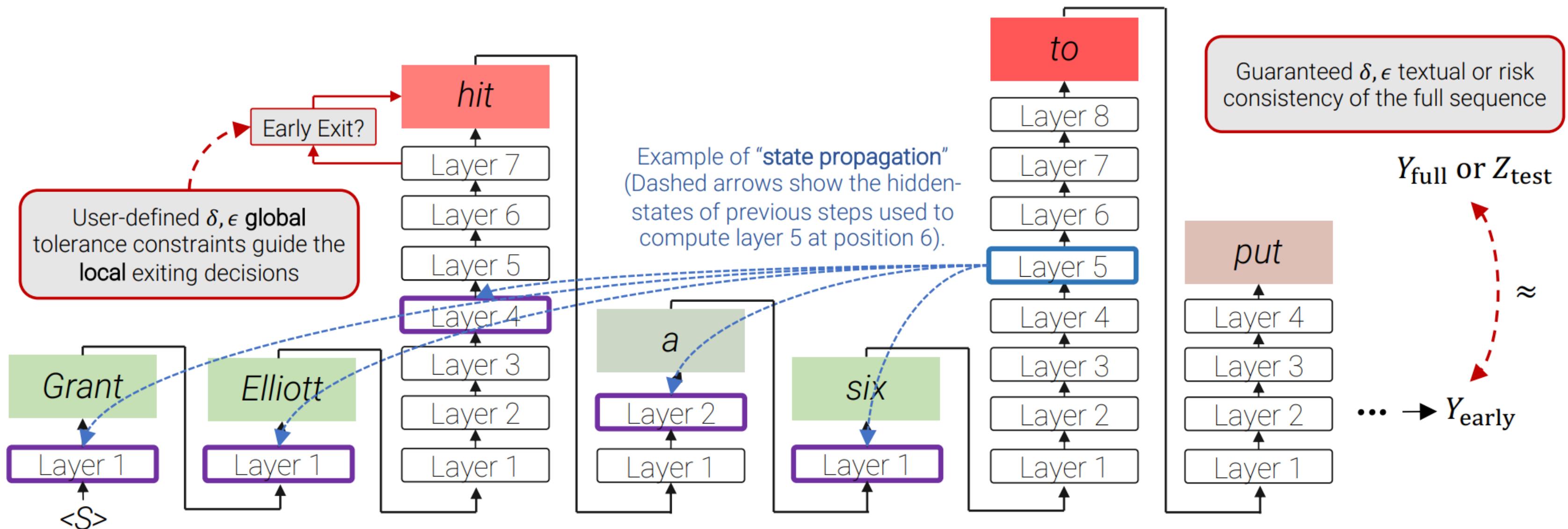
Early exit confidence measure:

Softmax > classifier > hidden states

But 😱

Softmax is **expensive!** And,

NLG requires previous KV cache!



## CALM's drawback

<b>Method</b>	<b>Model Type</b>	<b>Generation</b>	<b>Token Level</b>	<b>Batching</b>	<b>KV-Caching</b>	<b>Full Attention</b>	<b>Controlled Comp. Cost</b>
CALM	Enc-Dec	✓	✓	✗	✗	✗	✗
SkipDecode	Dec Only	✓	✓	✓	✓	✓	✓

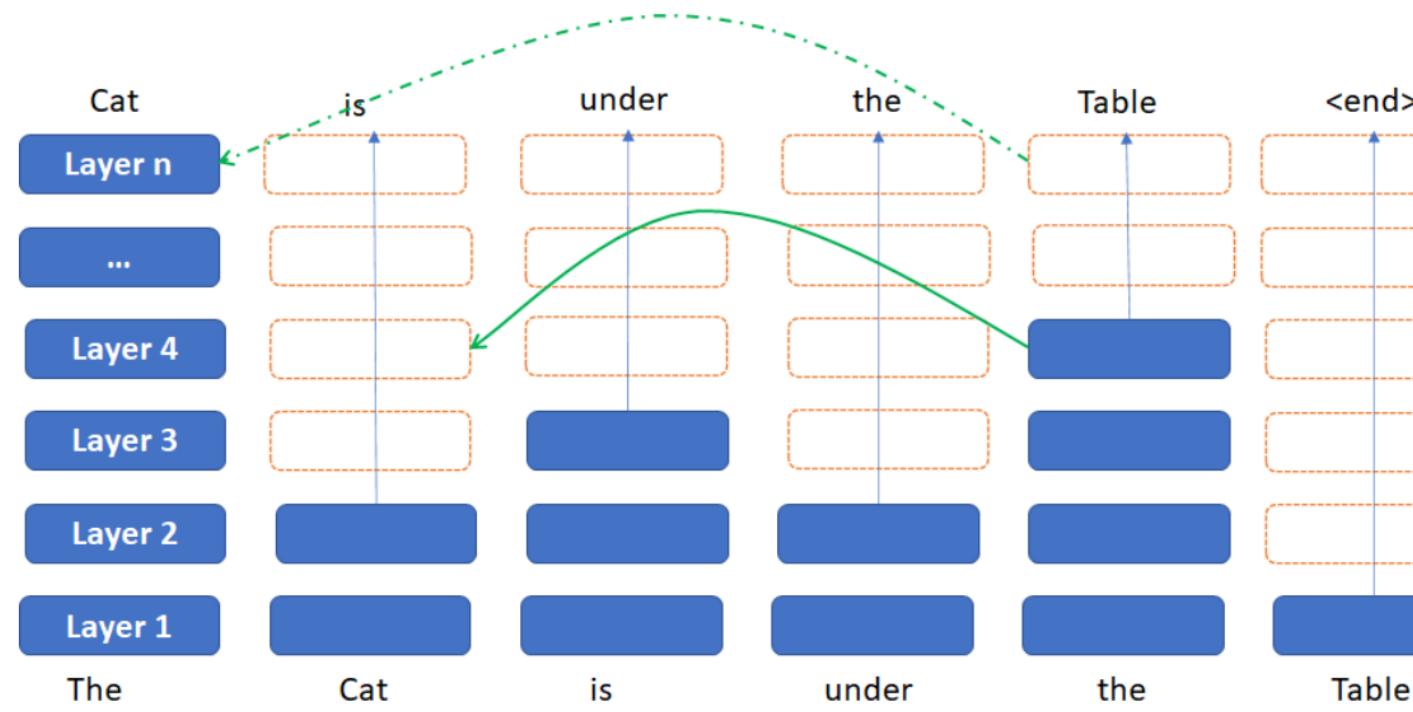
Table 1: Comparison of CALM and SkipDecode. SkipDecode supports batching and KV caching for increasing inference efficiency with controlled computational budget.

# SkipDecode

## Confident Adaptive Language Modeling. NeurIPS 2022

Practical Blockers (Existing):

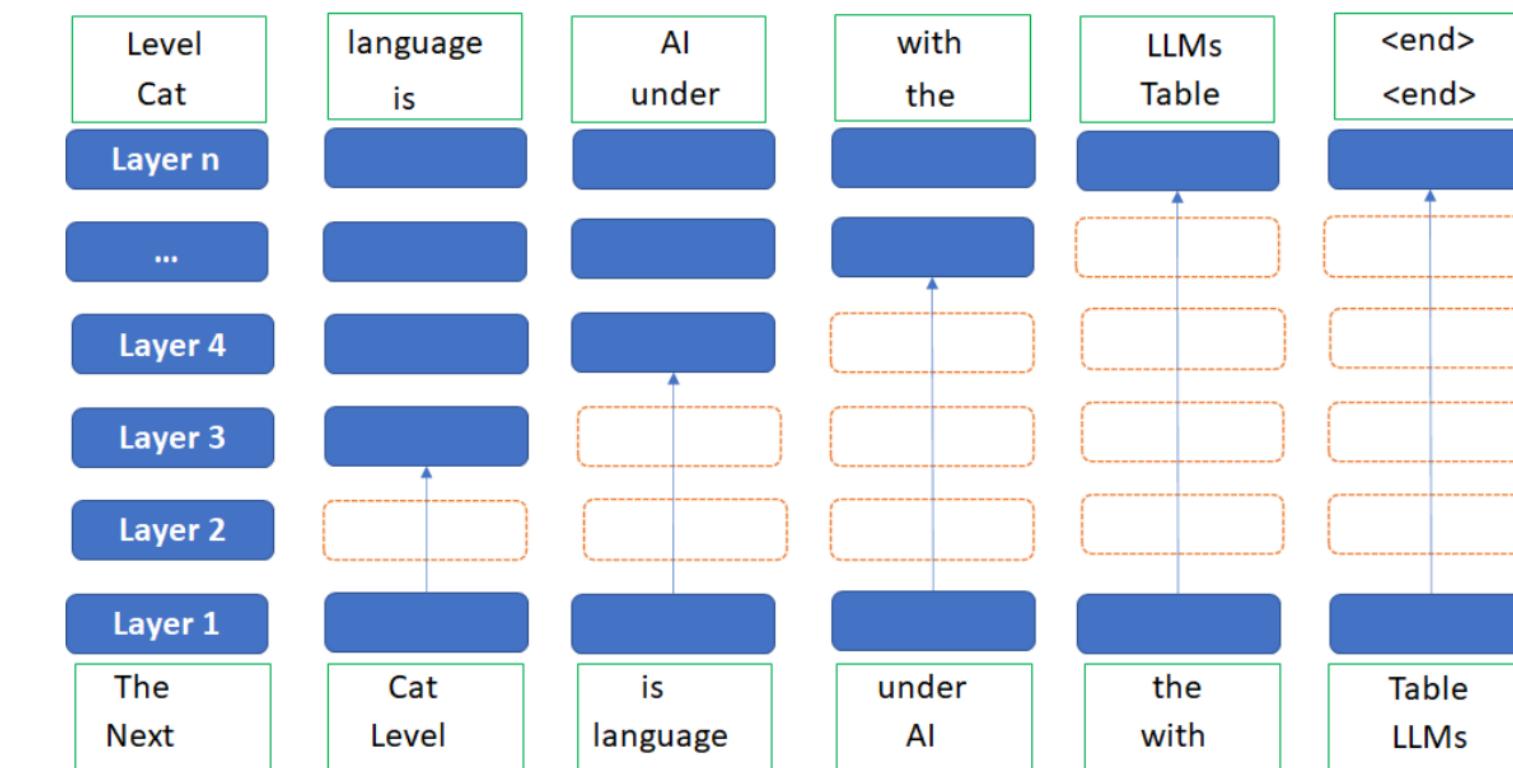
- **Batching:** Computational cost defined by the last exit token
- **KV Caching:** If next token exits later than previous one, we need to recompute KV values for previous tokens.
- **Computational Efficiency:** If next token exits earlier, it does not attend full computation of previous token.
- **Cost Uncertainty:** Worst case scenario (for a bad token exit, e.g., from last layer) equivalent to processing the whole network.



(a) Early Termination

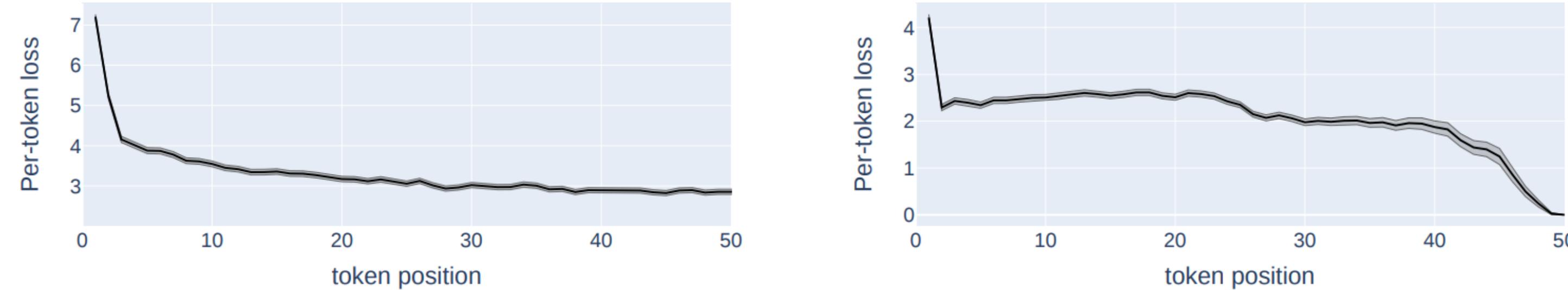
Solutions (Ours):

- **Batching:** Exit per position per batch (column-wise).
- **KV Caching:** Next column has to exit earlier than previous column. Leftwards tokens are more difficult to generate.
- **Computational Efficiency:** Spend most of computational budget on top layers. Implicitly attends the full computation of previous tokens.
- **Cost Uncertainty:** Static policy (no surprises), computation costs are predefined.



(b) Skipping

# Skipped tokens

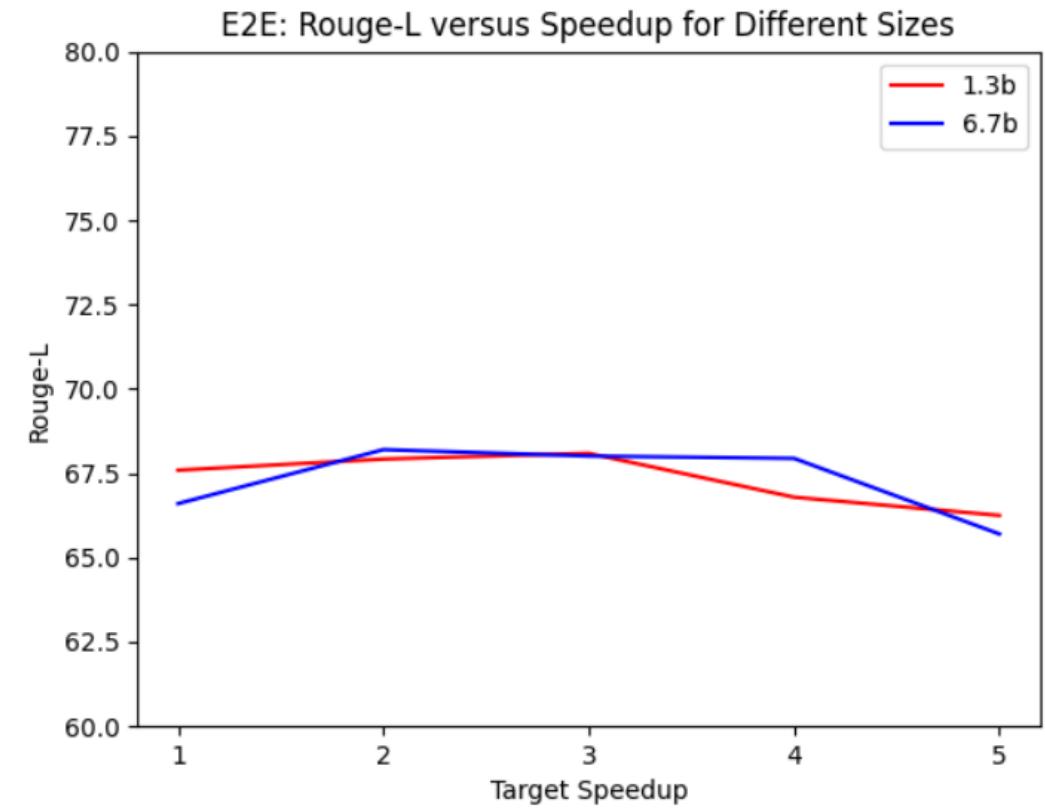


(a) **OPT-350m on OpenWebText**. Average loss per token position shows a strong monotonically-decreasing trend for general text.  
(b) **OPT-350m (finetuned) on Reddit-TLDR**. Average loss per token position. Decreasing trend but with a different function.

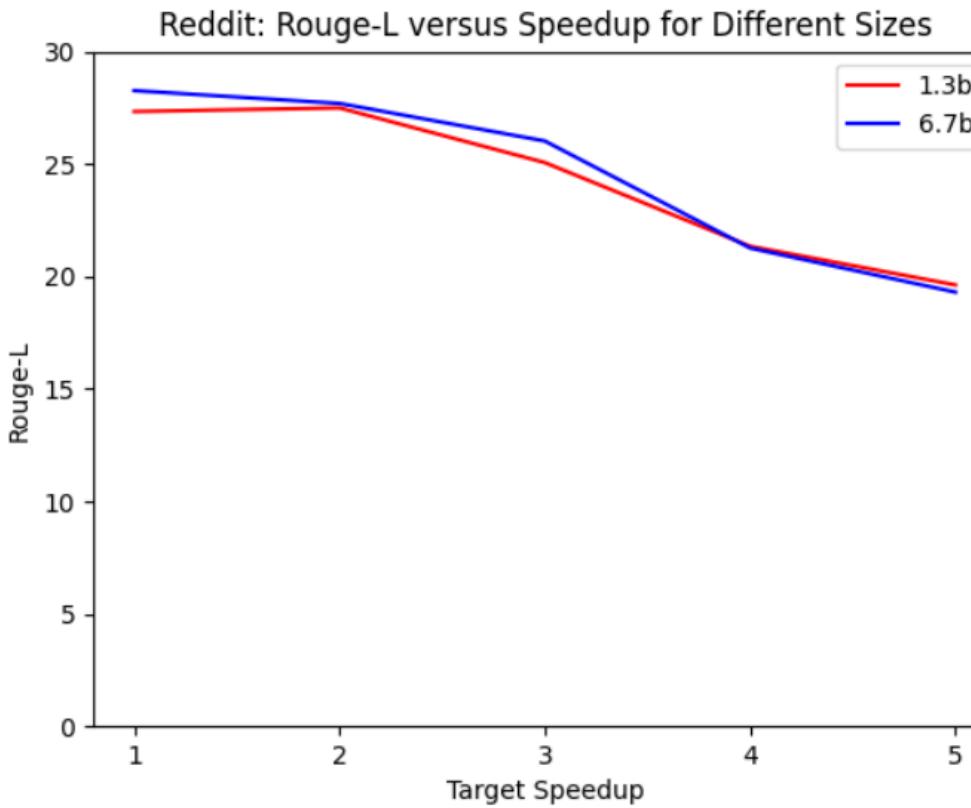
Figure 1: Average loss per token position (black) during the forward pass of OPT-350m model on a general and a task-specific dataset. Grey represents the 95% confidence interval on the mean.

Front tokens have large losses, which need more computation

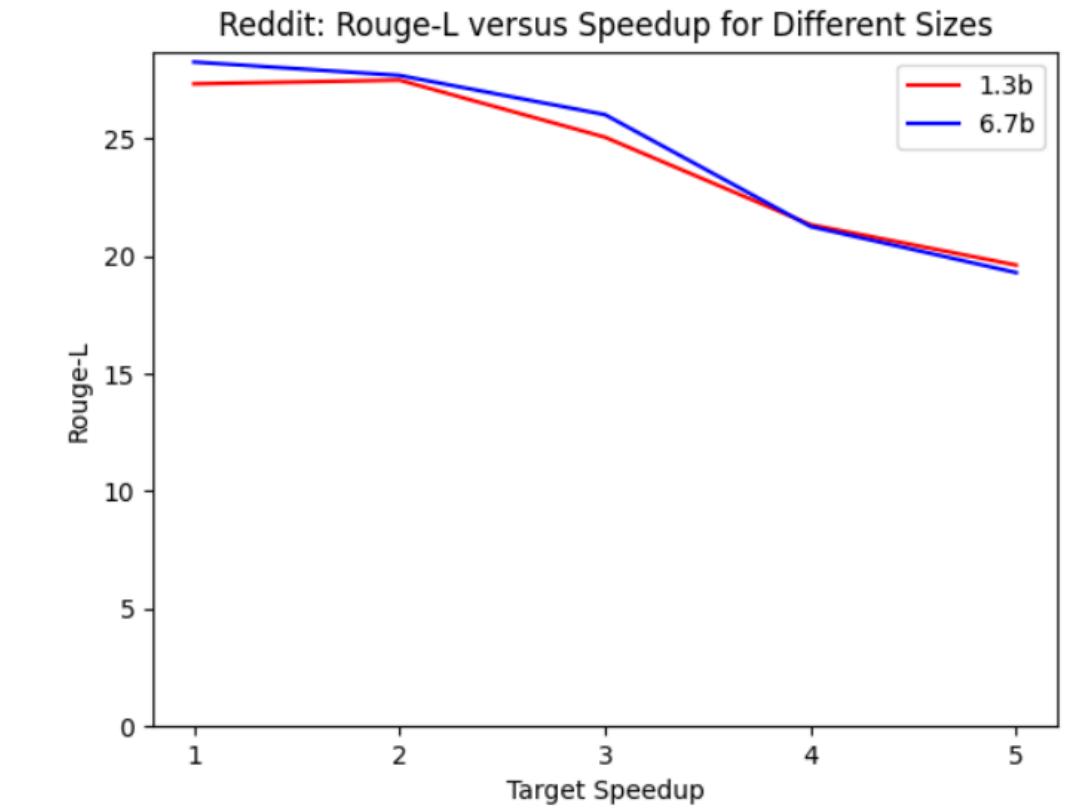
# SkipDecode results



(a) E2E



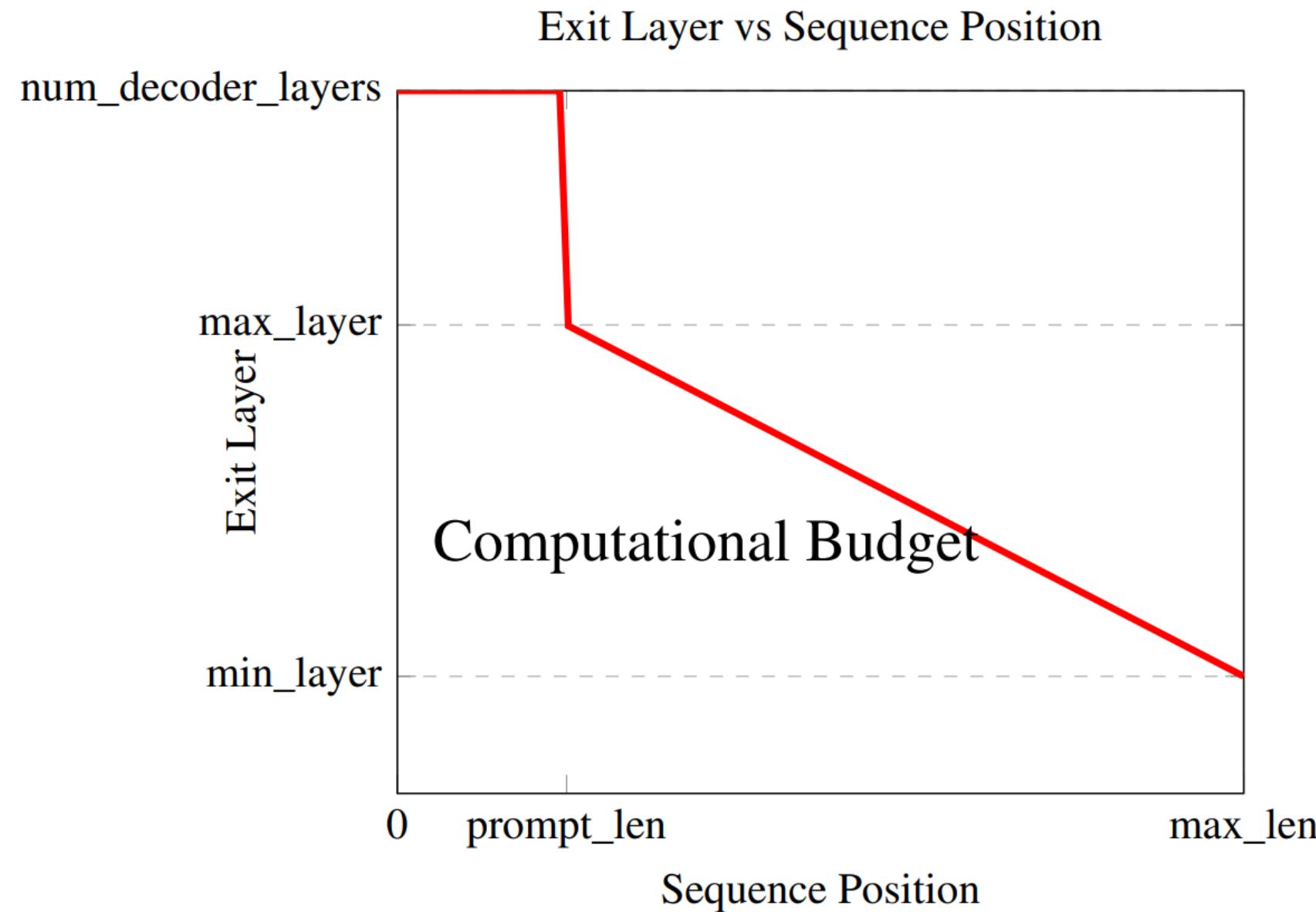
(b) Reddit



(c) CNN\_DM

Works on LM, but fails on summarization, which is a common delima in efficient inference. 🤦

# Fixed exit point



SkipDecode pre-define fixed exit point of each token 🤖

Strength: **No Softmax**, batching ✓

# Rethinking fixed exit point



SkipDecode:

The token at front of each **SEQUENCE** is hard.

Recall CALM's result:

The token at front of each **SENTENCE** is hard.

It is hard to find **SENTENCE** begining (dynamic) with batching

## Review early exiting

`input\_tensor` of shape `[bsz, len, hsize]`

`model` with  $L$  layer and  $H$  hsize

- Early exiting reduces decode  $L$  only
- Pruning reduces  $H$
- We try to use batching with larger `bsz`

What about `len` 

# Speculative Decoding

Fast Inference from Transformers via Speculative Decoding. ICML 2023

```
[START] japan ' s benchmark bond n
[START] japan ' s benchmark nikkei 22 75
[START] japan ' s benchmark nikkei 225 index rose 22 76
[START] japan ' s benchmark nikkei 225 index rose 226 . 69 7 points
[START] japan ' s benchmark nikkei 225 index rose 226 . 69 points , or 0 1
[START] japan ' s benchmark nikkei 225 index rose 226 . 69 points , or 1 . 5 percent , to 10 , 9859
[START] japan ' s benchmark nikkei 225 index rose 226 . 69 points , or 1 . 5 percent , to 10 , 989 . 79 7 in
[START] japan ' s benchmark nikkei 225 index rose 226 . 69 points , or 1 . 5 percent , to 10 , 989 . 79 in tokyo late
[START] japan ' s benchmark nikkei 225 index rose 226 . 69 points , or 1 . 5 percent , to 10 , 989 . 79 in late morning trading . [END]
```

- Using small model to generate multiple tokens
- Using LLM to verify the tokens, accept / reject. Then, and re-generate rejected token or generate the last token

# Algorithm

---

**Algorithm 1** SpeculativeDecodingStep

---

**Inputs:**  $M_p, M_q, prefix$ .

▷ Sample  $\gamma$  guesses  $x_{1,\dots,\gamma}$  from  $M_q$  autoregressively.

**for**  $i = 1$  **to**  $\gamma$  **do**

$$q_i(x) \leftarrow M_q(prefix + [x_1, \dots, x_{i-1}])$$

$$x_i \sim q_i(x)$$

**end for**

▷ Run  $M_p$  in parallel.

$$p_1(x), \dots, p_{\gamma+1}(x) \leftarrow M_p(prefix), \dots, M_p(prefix + [x_1, \dots, x_\gamma])$$

▷ Determine the number of accepted guesses  $n$ .

$$r_1 \sim U(0, 1), \dots, r_\gamma \sim U(0, 1)$$

$$n \leftarrow \min(\{i - 1 \mid 1 \leq i \leq \gamma, r_i > \frac{p_i(x)}{q_i(x)}\} \cup \{\gamma\})$$

▷ Adjust the distribution from  $M_p$  if needed.

$$p'(x) \leftarrow p_{n+1}(x)$$

**if**  $n < \gamma$  **then**

$$p'(x) \leftarrow \text{norm}(\max(0, p_{n+1}(x) - q_{n+1}(x)))$$

**end if**

▷ Return one token from  $M_p$ , and  $n$  tokens from  $M_q$ .

$$t \sim p'(x)$$

**return**  $prefix + [x_1, \dots, x_n, t]$

---

## Small model

Generate  $\gamma$  tokens **autoregressively**

▷ Sample  $\gamma$  guesses  $x_{1,\dots,\gamma}$  from  $M_q$  autoregressively.

**for**  $i = 1$  **to**  $\gamma$  **do**

$q_i(x) \leftarrow M_q(\text{prefix} + [x_1, \dots, x_{i-1}])$

$x_i \sim q_i(x)$

**end for**

## LLM

Generate  $\gamma + 1$  tokens **in parallel** (same as training)

▷ Run  $M_p$  in parallel.

$$p_1(x), \dots, p_{\gamma+1}(x) \leftarrow M_p(\text{prefix}), \dots, M_p(\text{prefix} + [x_1, \dots, x_\gamma])$$

Validate  $\gamma$  tokens, reject **once** LLM don't think the guess is good  $p_i(x) < q_i(x) \rightarrow r_i > p_i(x)/q_i(x)$

▷ Determine the number of accepted guesses  $n$ .

$$r_1 \sim U(0, 1), \dots, r_\gamma \sim U(0, 1)$$

$$n \leftarrow \min(\{i - 1 \mid 1 \leq i \leq \gamma, r_i > \frac{p_i(x)}{q_i(x)}\} \cup \{\gamma\})$$

Regenerate last token or rejected token

▷ Adjust the distribution from  $M$  if needed

# Output

▷ Return one token from  $M_p$ , and  $n$  tokens from  $M_q$ .

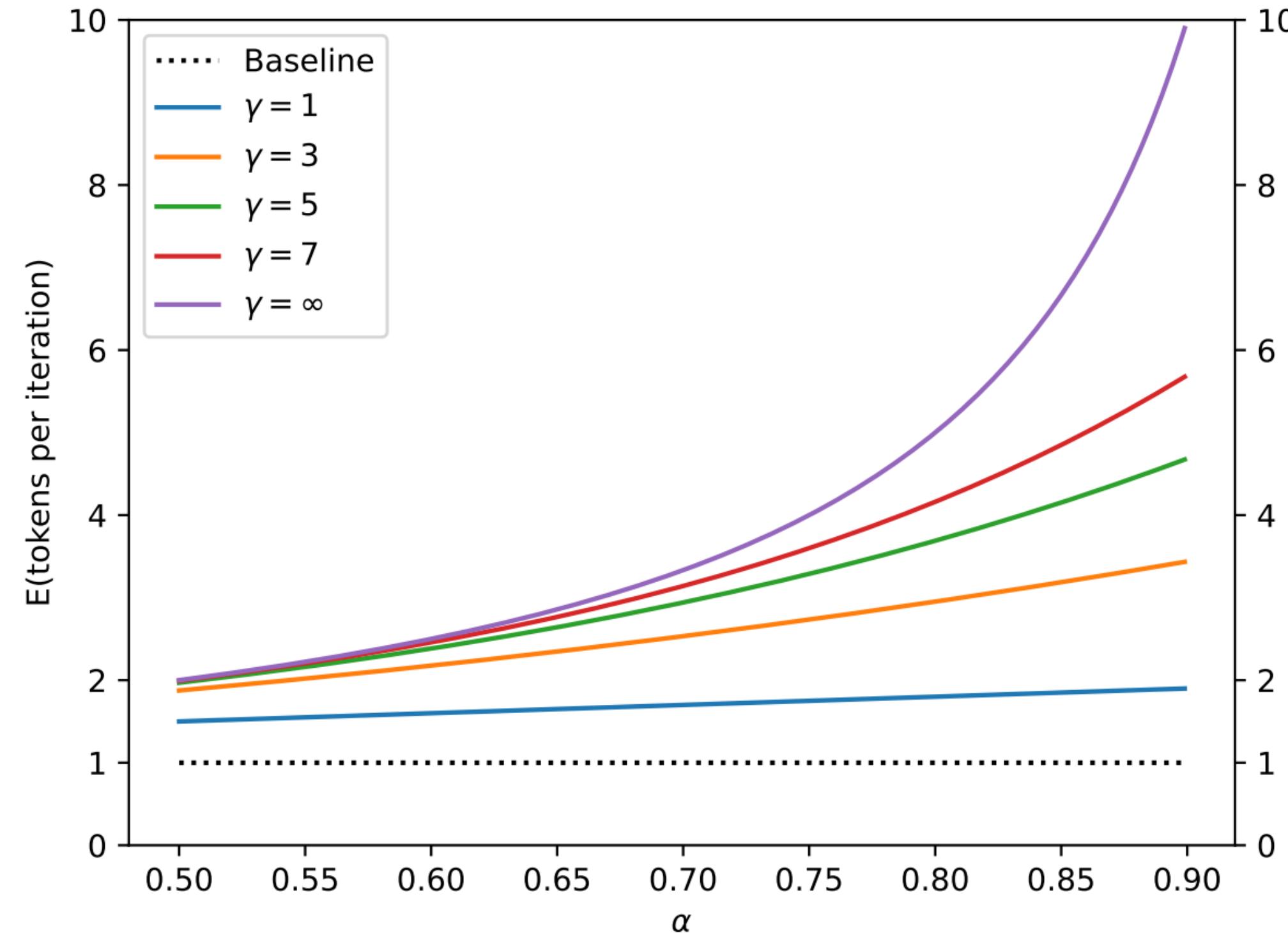
$$t \sim p'(x)$$

**return**  $prefix + [x_1, \dots, x_n, t]$



# Model Alignment and Speedup

Expected generated tokens in one run  $\frac{1-\alpha^{\gamma+1}}{1-\alpha}$



# Model Alignment Evaluation

$M_p$	$M_q$	SMPL	$\alpha$
GPT-LIKE (97M)	UNIGRAM	T=0	0.03
GPT-LIKE (97M)	BIGRAM	T=0	0.05
GPT-LIKE (97M)	GPT-LIKE (6M)	T=0	0.88
GPT-LIKE (97M)	UNIGRAM	T=1	0.03
GPT-LIKE (97M)	BIGRAM	T=1	0.05
GPT-LIKE (97M)	GPT-LIKE (6M)	T=1	0.89
T5-XXL (ENDE)	UNIGRAM	T=0	0.08
T5-XXL (ENDE)	BIGRAM	T=0	0.20
T5-XXL (ENDE)	T5-SMALL	T=0	0.75
T5-XXL (ENDE)	T5-BASE	T=0	0.80
T5-XXL (ENDE)	T5-LARGE	T=0	0.82
T5-XXL (ENDE)	UNIGRAM	T=1	0.07
T5-XXL (ENDE)	BIGRAM	T=1	0.19
T5-XXL (ENDE)	T5-SMALL	T=1	0.62
T5-XXL (ENDE)	T5-BASE	T=1	0.68
T5-XXL (ENDE)	T5-LARGE	T=1	0.71

T5-XXL (CNNDM)	UNIGRAM	T=0	0.13
T5-XXL (CNNDM)	BIGRAM	T=0	0.23
T5-XXL (CNNDM)	T5-SMALL	T=0	0.65
T5-XXL (CNNDM)	T5-BASE	T=0	0.73
T5-XXL (CNNDM)	T5-LARGE	T=0	0.74
T5-XXL (CNNDM)	UNIGRAM	T=1	0.08
T5-XXL (CNNDM)	BIGRAM	T=1	0.16
T5-XXL (CNNDM)	T5-SMALL	T=1	0.53
T5-XXL (CNNDM)	T5-BASE	T=1	0.55
T5-XXL (CNNDM)	T5-LARGE	T=1	0.56
LAMDA (137B)	LAMDA (100M)	T=0	0.61
LAMDA (137B)	LAMDA (2B)	T=0	0.71
LAMDA (137B)	LAMDA (8B)	T=0	0.75
LAMDA (137B)	LAMDA (100M)	T=1	0.57
LAMDA (137B)	LAMDA (2B)	T=1	0.71
LAMDA (137B)	LAMDA (8B)	T=1	0.74

Even N-gram can accelerate LLM generation 😊

# Speculative Decoding Speedup

Table 2. Empirical results for speeding up inference from a T5-XXL 11B model.

TASK	$M_q$	TEMP	$\gamma$	$\alpha$	SPEED
ENDE	T5-SMALL ★	0	7	0.75	<b>3.4X</b>
ENDE	T5-BASE	0	7	0.8	2.8X
ENDE	T5-LARGE	0	7	0.82	1.7X
ENDE	T5-SMALL ★	1	7	0.62	<b>2.6X</b>
ENDE	T5-BASE	1	5	0.68	2.4X
ENDE	T5-LARGE	1	3	0.71	1.4X
CNNDM	T5-SMALL ★	0	5	0.65	<b>3.1X</b>
CNNDM	T5-BASE	0	5	0.73	3.0X
CNNDM	T5-LARGE	0	3	0.74	2.2X
CNNDM	T5-SMALL ★	1	5	0.53	<b>2.3X</b>
CNNDM	T5-BASE	1	3	0.55	2.2X
CNNDM	T5-LARGE	1	3	0.56	1.7X

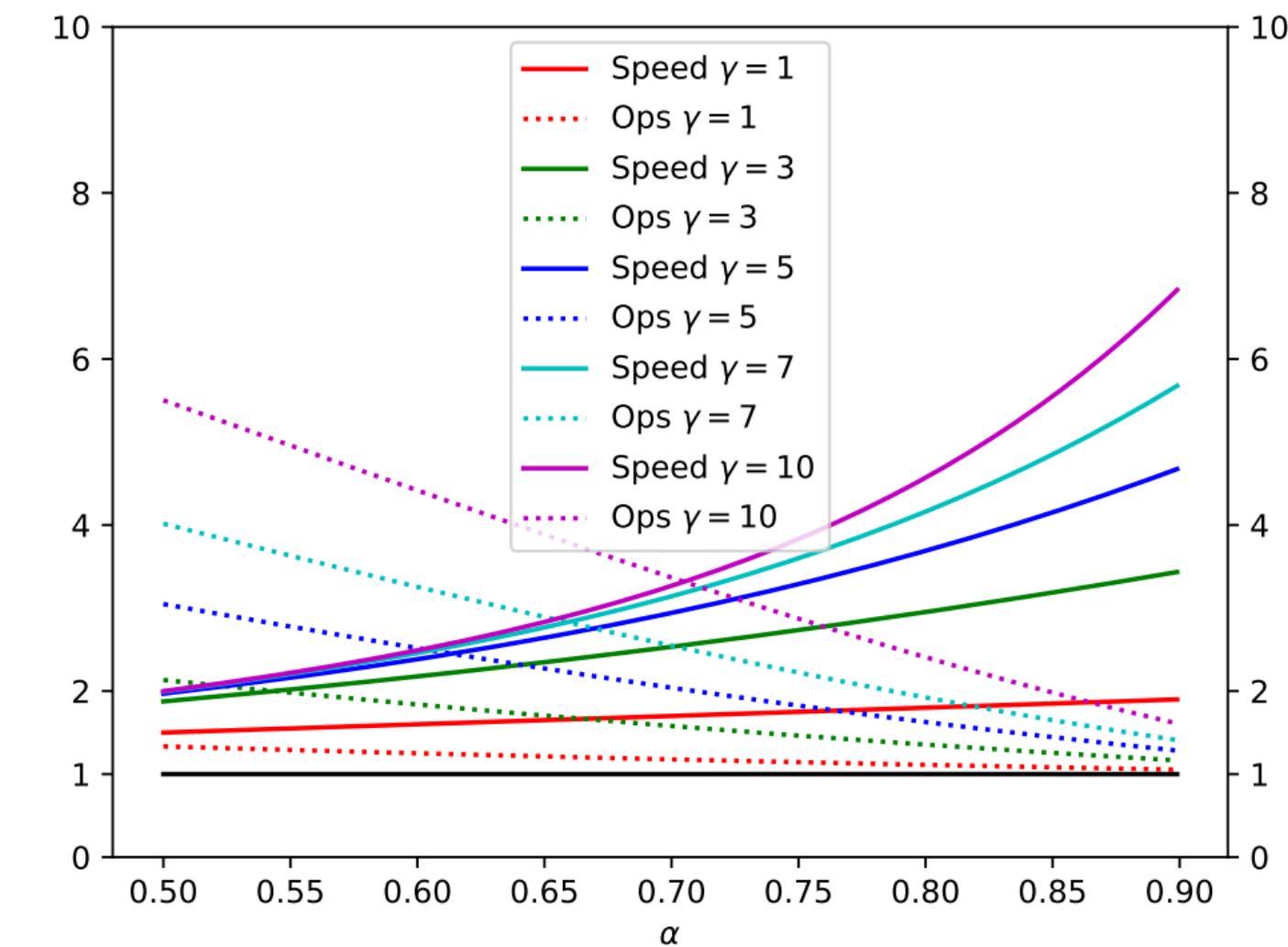


Figure 4. The speedup factor and the increase in number of arithmetic operations as a function of  $\alpha$  for various values of  $\gamma$ .

2-3X speedup with a small model 🤖

# Strength and weakness

## Strength

- Don't need to change **any** of the model!
- **Always** has a fallback accuracy to LLM!

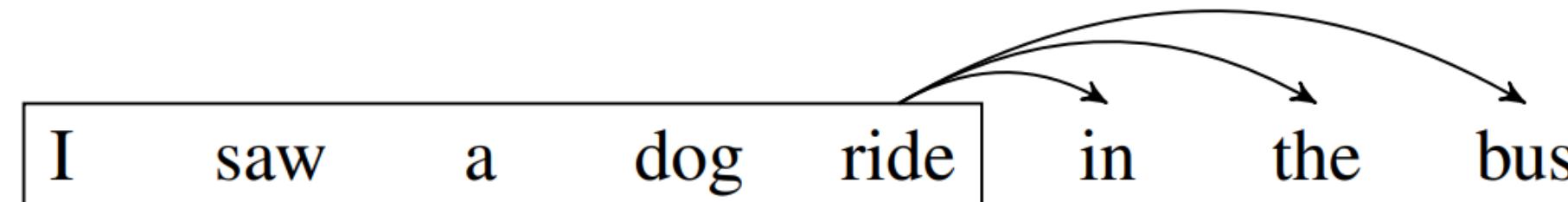
## Weakness

- Need two (or more) models, much more **memory usage**
- Small model need to be **aligned** with LLM or the efficiency falls back

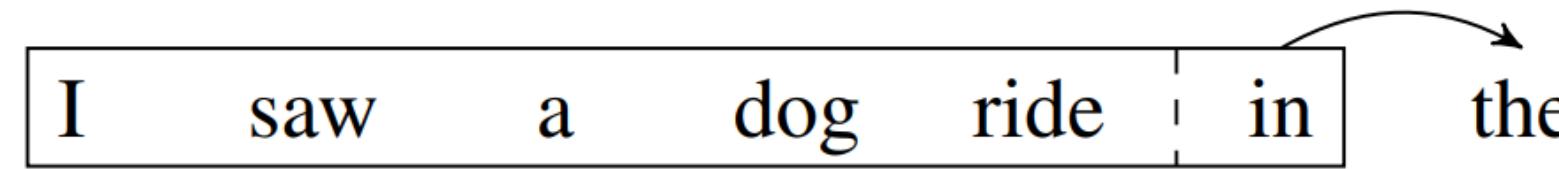
# Block-wise generation

Blockwise parallel decoding for deep autoregressive models. NeurIPS 2018

**Predict**

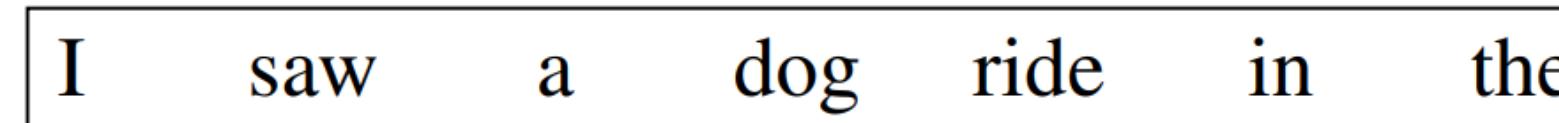


**Verify**



executed  
in parallel

**Accept**



# Resources

- 大模型推理性能优化之KV Cache解读
- NLP (十七) : 从 FlashAttention 到 PagedAttention, 如何进一步优化 Attention 性能
- FlashAttention: Fast and Memory-Efficient Exact Attention with IO-Awareness
- vLLM. Easy, fast, and cheap LLM serving for everyone
- Accelerating text generation with Confident Adaptive Language Modeling (CALM)
- SkipDecode: Autoregressive Skip Decoding with Batching and Caching for Efficient LLM Inference
- Fast Inference from Transformers via Speculative Decoding
- Efficient Systems for Foundation Models ES-FoMo@ICML2023 talks