

CS480/680: Introduction to Machine Learning

Lecture 13: Speculative Decoding

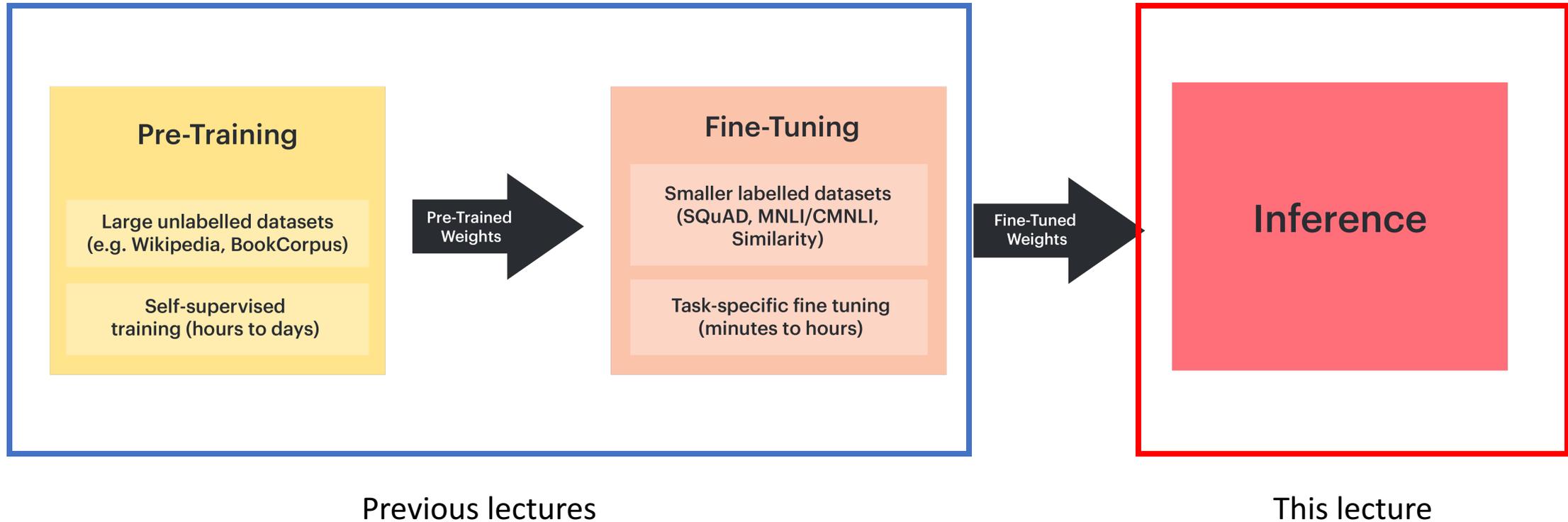
Hongyang Zhang



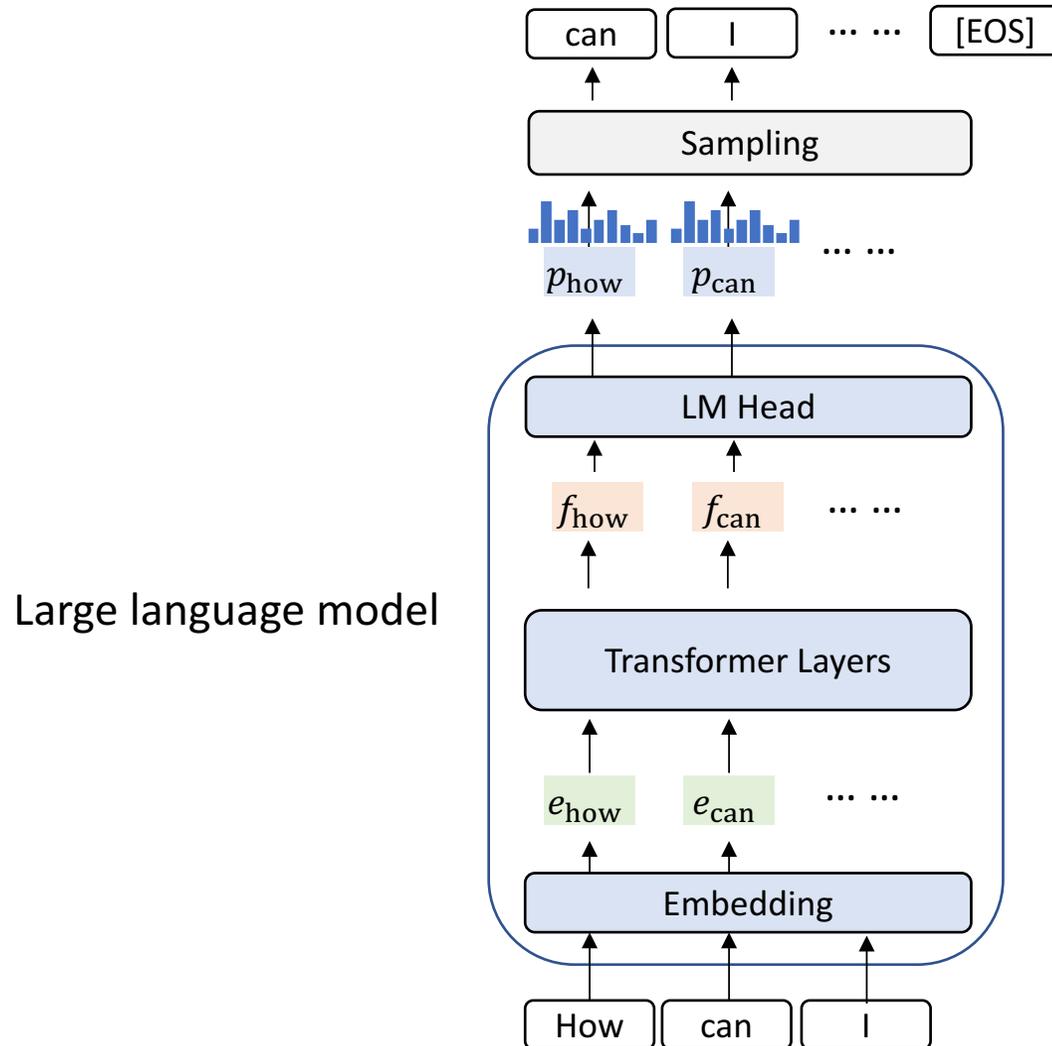
**UNIVERSITY OF
WATERLOO**

July 14, 2025

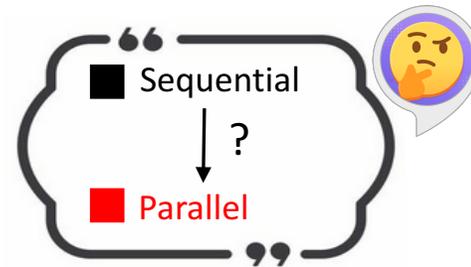
Pretraining-Finetuning-Inference



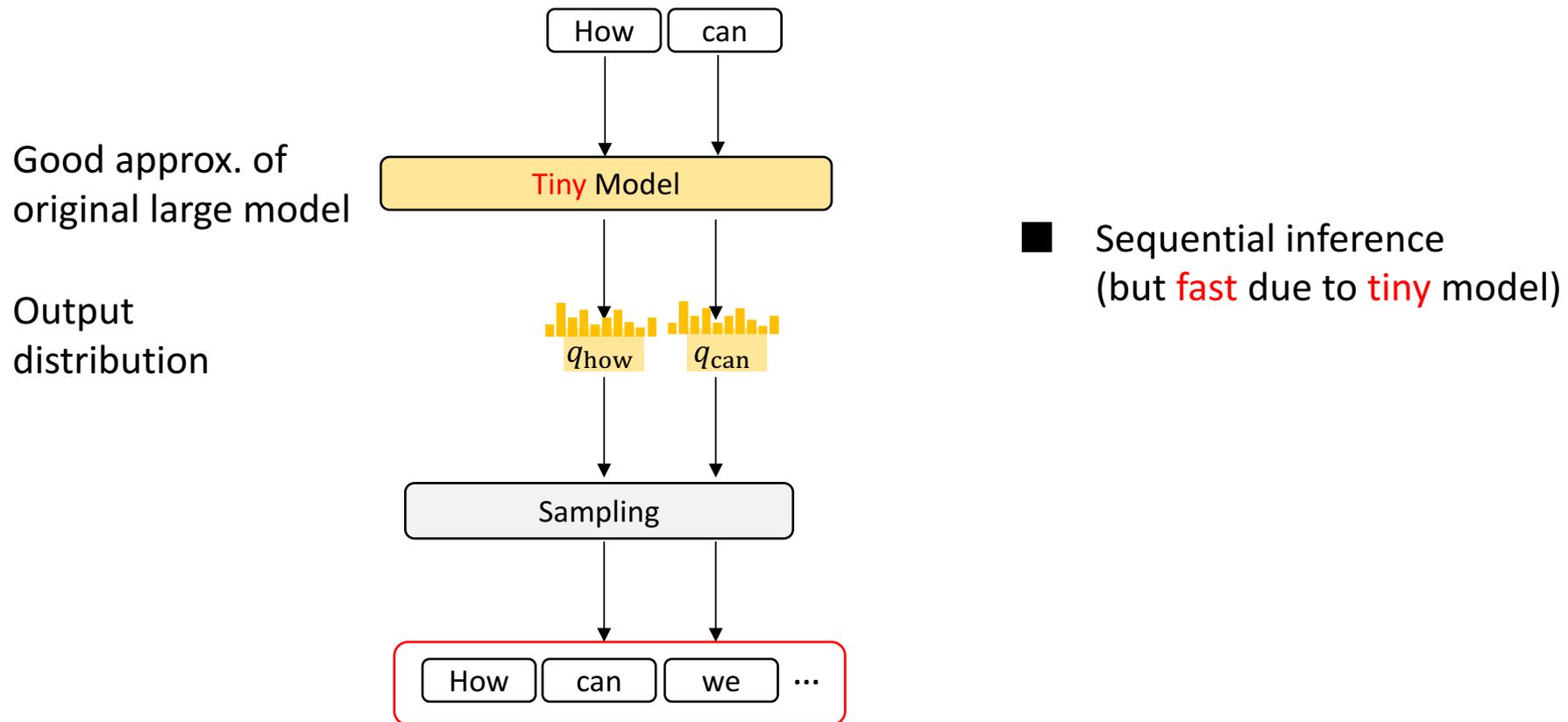
Vanilla autoregressive inference



Sequential inference

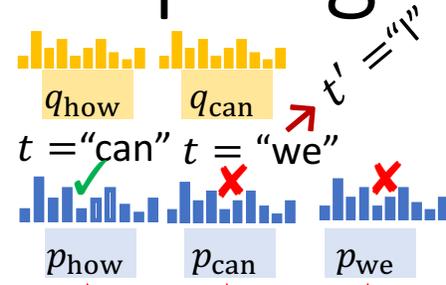


Speculative sampling framework (draft)



Speculative sampling framework (check)

Compare with **tiny** model



1. $r \sim U(0,1)$, if $r < \min\left(1, \frac{p(t)}{q(t)}\right)$, next token = t

Accept rate

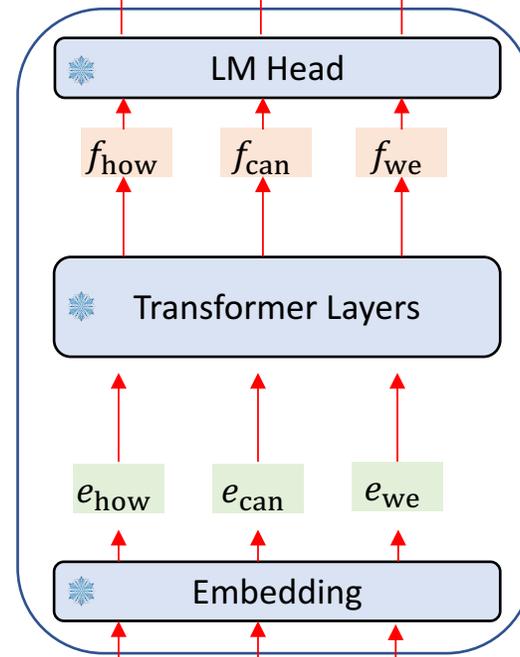
2. else: next token = $t' \sim \text{norm}(\max(0, p - q))$

Residual distribution

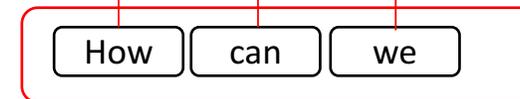
■ Check in parallel



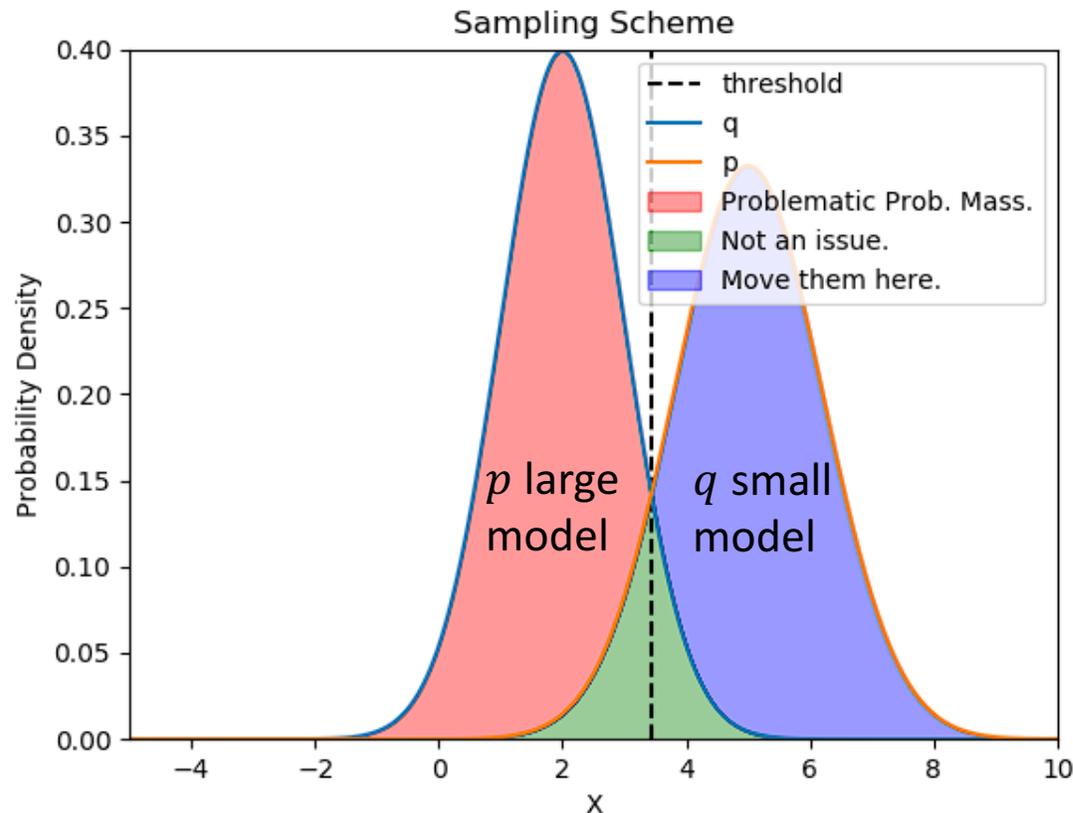
Original **large** model



Drafted by **tiny** model



Speculative sampling framework (check)



1. $r \sim U(0,1)$, if $r < \min\left(1, \frac{p(t)}{q(t)}\right)$, next token = t

Accept rate

2. else: next token = $t' \sim \text{norm}(\max(0, p - q))$

Residual distribution

Theorem (reject sampling):

The above (p, q) sampling procedure is equivalent to sampling from p .

All is about how to use sampling from (p, q) to mimic sampling from p .

- p and q are closer \rightarrow higher accept rate \rightarrow higher speedup ratio

Reject sampling proof

Proof:

$$\begin{aligned} & \Pr(X = t) \\ &= \Pr(\tilde{X} = t) \Pr(\tilde{X} \text{ accept} | \tilde{X} = t) + \Pr(\tilde{X} \text{ reject}) \Pr(\tilde{X} = t | \tilde{X} \text{ reject}) \\ &= q(t) \min\left(1, \frac{p(t)}{q(t)}\right) + (1 - \Pr(\tilde{X} \text{ accept})) \text{norm}(\max(0, p(t) - q(t))) \\ &= \min(p(t), q(t)) + \left(1 - \sum_t \min(p(t), q(t))\right) \frac{\max(0, p(t) - q(t))}{\sum_t \max(0, p(t) - q(t))} \\ &= \min(p(t), q(t)) + \max(0, p(t) - q(t)) \\ &= p(t) \end{aligned}$$

1. $r \sim U(0,1)$, if $r < \min\left(1, \frac{p(t)}{q(t)}\right)$, next token = t

Accept rate

2. else: next token = $t' \sim \text{norm}(\max(0, p - q))$

Residual distribution

Theorem (reject sampling):

The above (p, q) sampling procedure is equivalent to sampling from p .

“
How to build the
tiny model q ?
”

Use a small model in the same family

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



TASK	M_q	TEMP	γ	α	SPEED
ENDE	T5-SMALL ★	0	7	0.75	3.4X
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	2.6X
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	3.1X
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	2.3X
CNNDM	T5-BASE	1	3	0.55	2.2X
CNNDM	T5-LARGE	1	3	0.56	1.7X

Using a small model in the same family works well, but:

1. What if the target model is the smallest in the family?
2. Can we do better?

How to build the tiny model q ?

- Trade-off between accuracy and efficiency



Fast but inaccurate

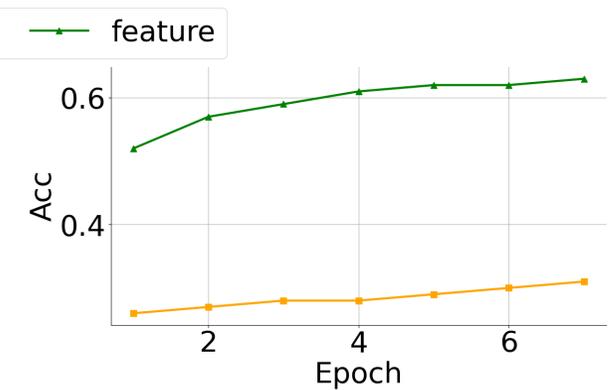
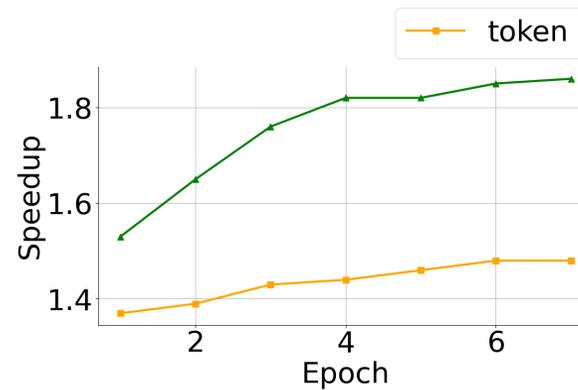
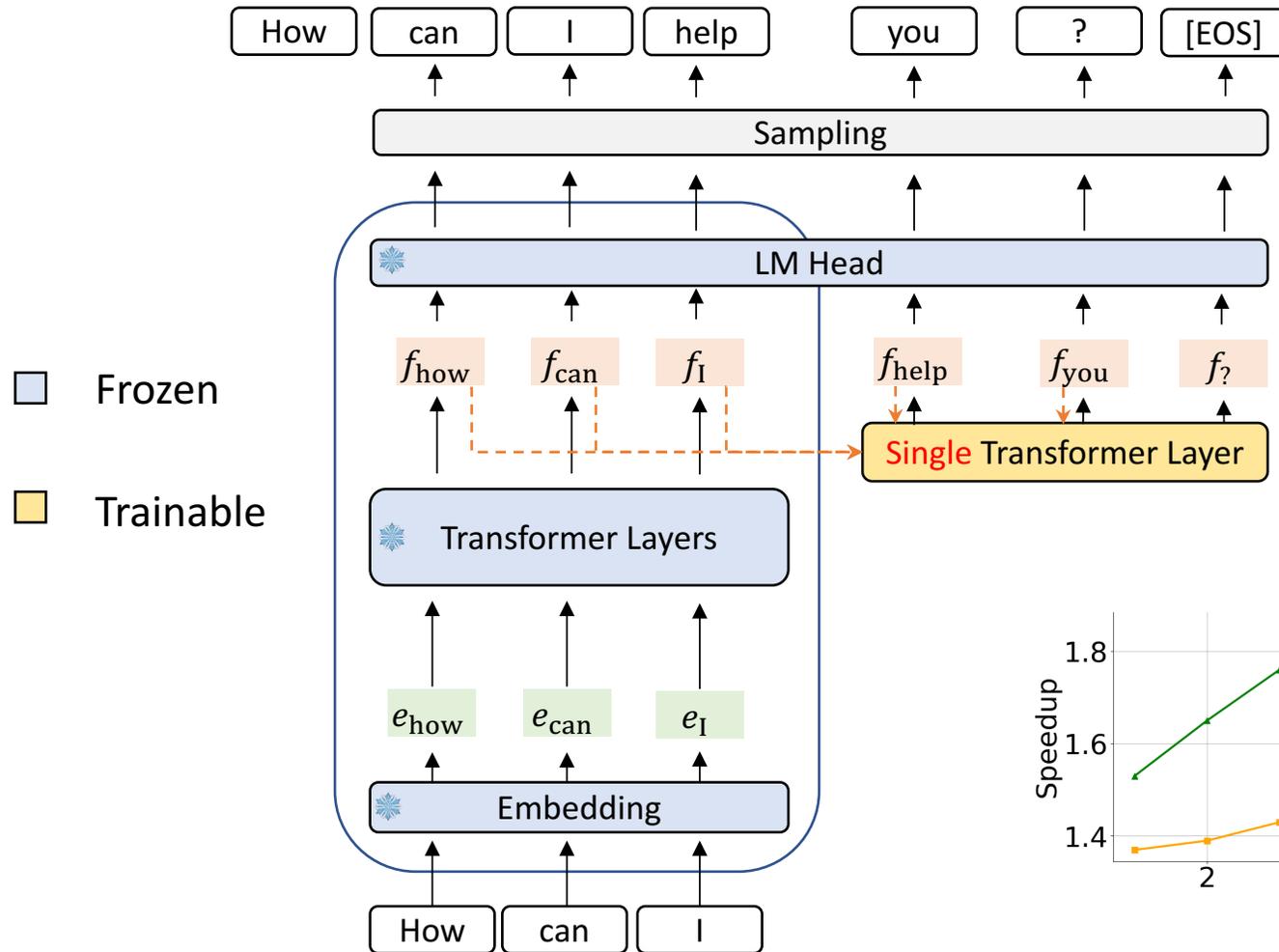


Accurate but slow

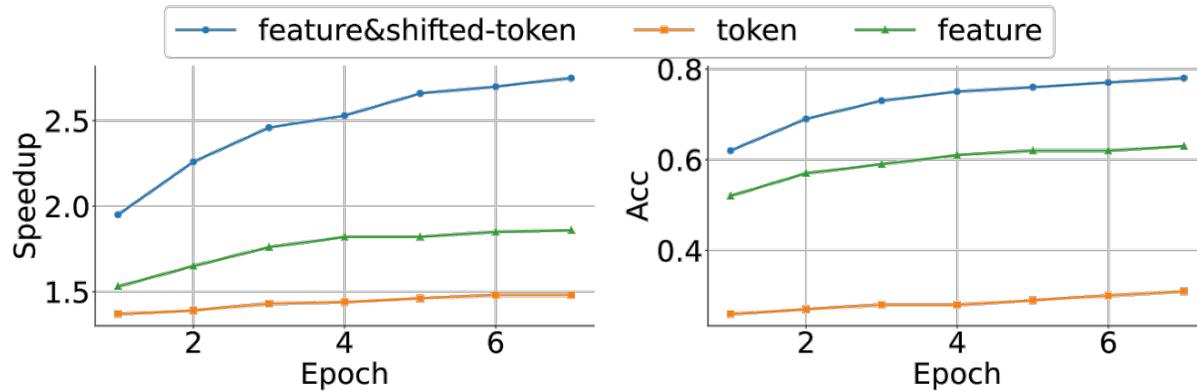
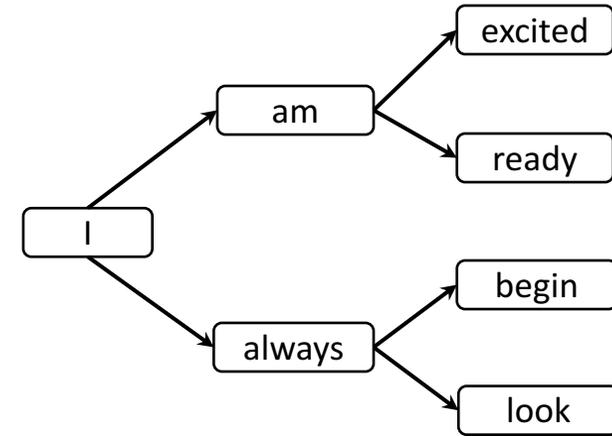
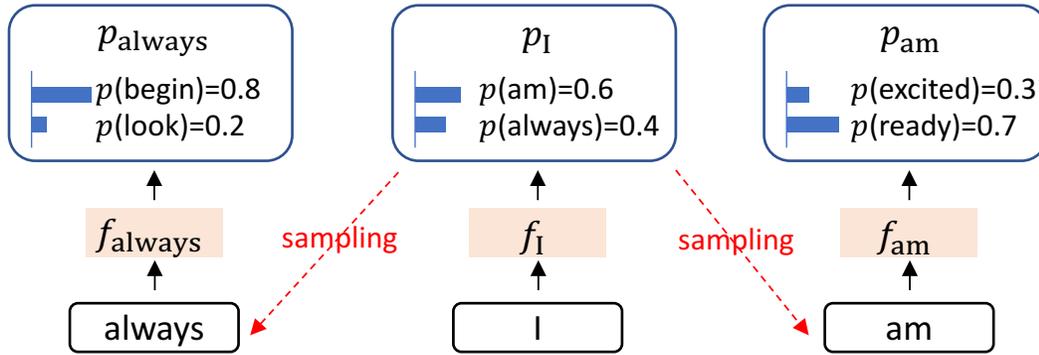
- However, training a draft model from scratch can be costly.

EAGLE: Speculative Sampling
Requires Rethinking Feature
Uncertainty

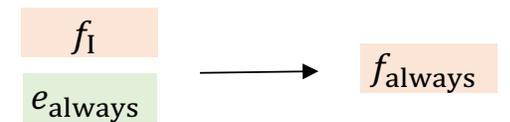
Our first trial: next-feature prediction



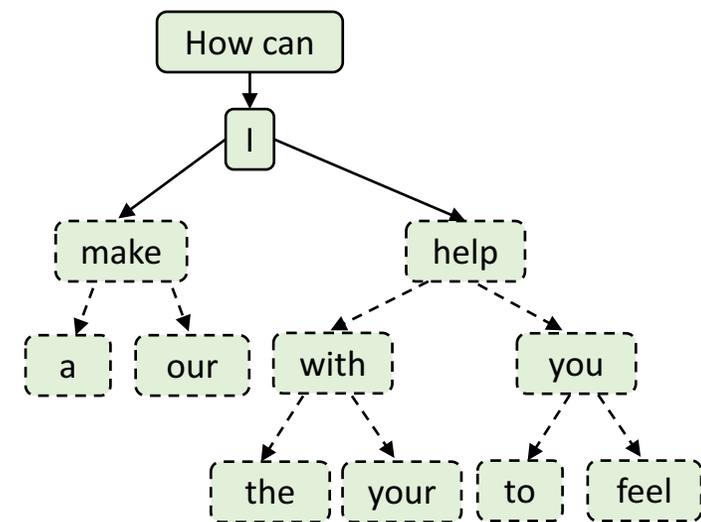
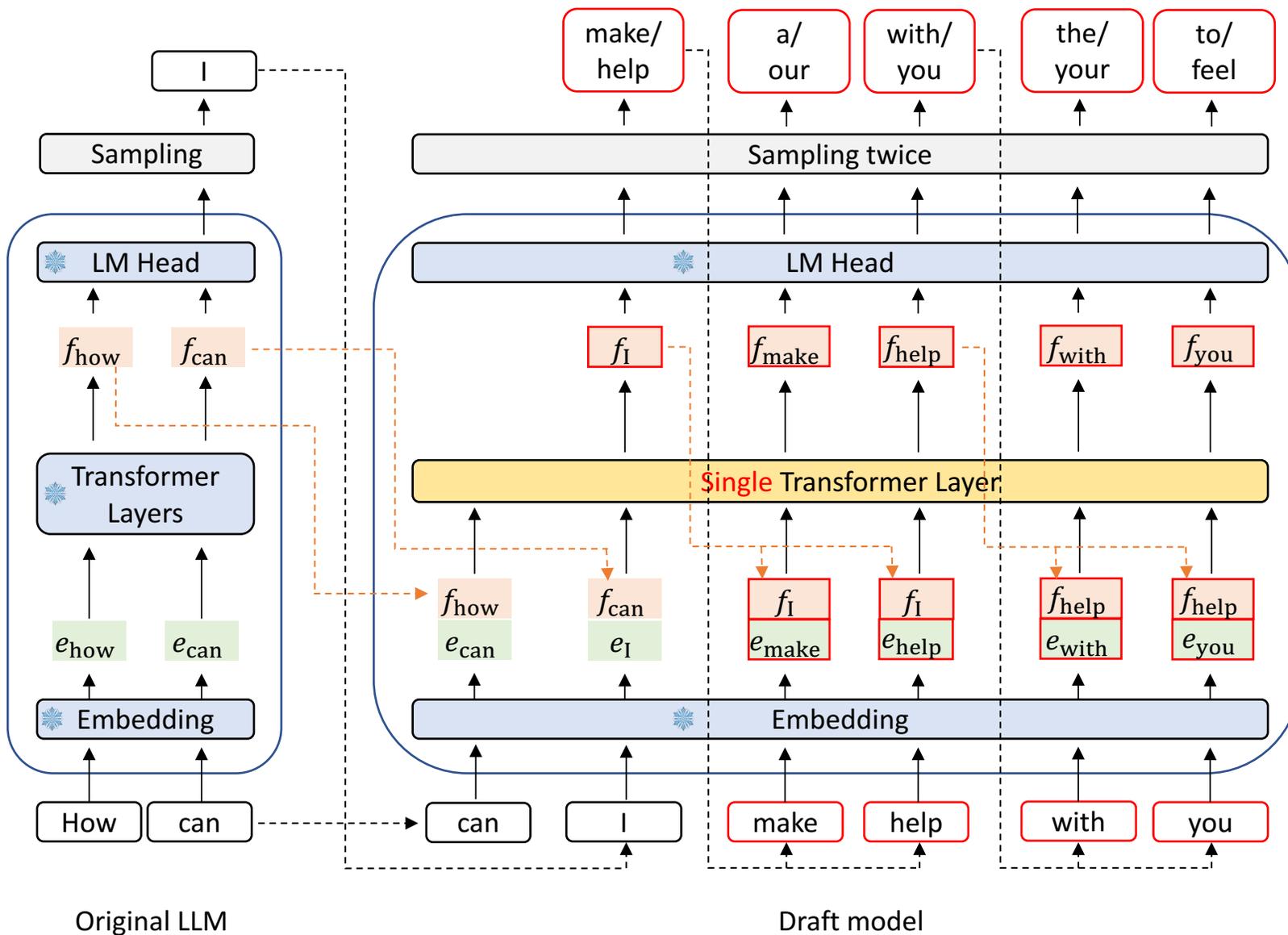
Feature uncertainty matters



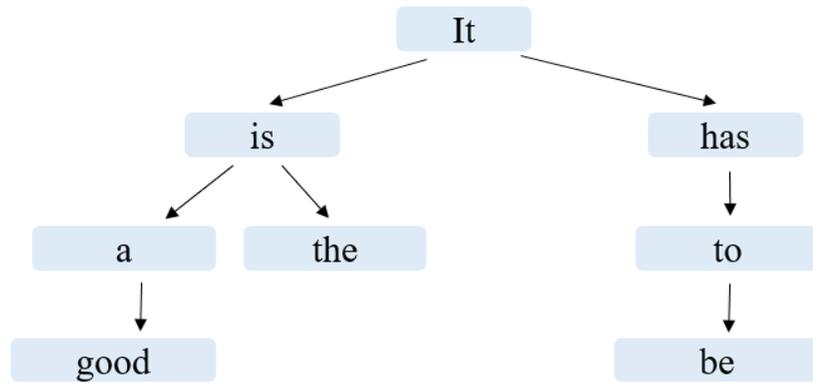
Idea: feature & shifted-token \rightarrow next feature



Our second trial: EAGLE



Tree attention



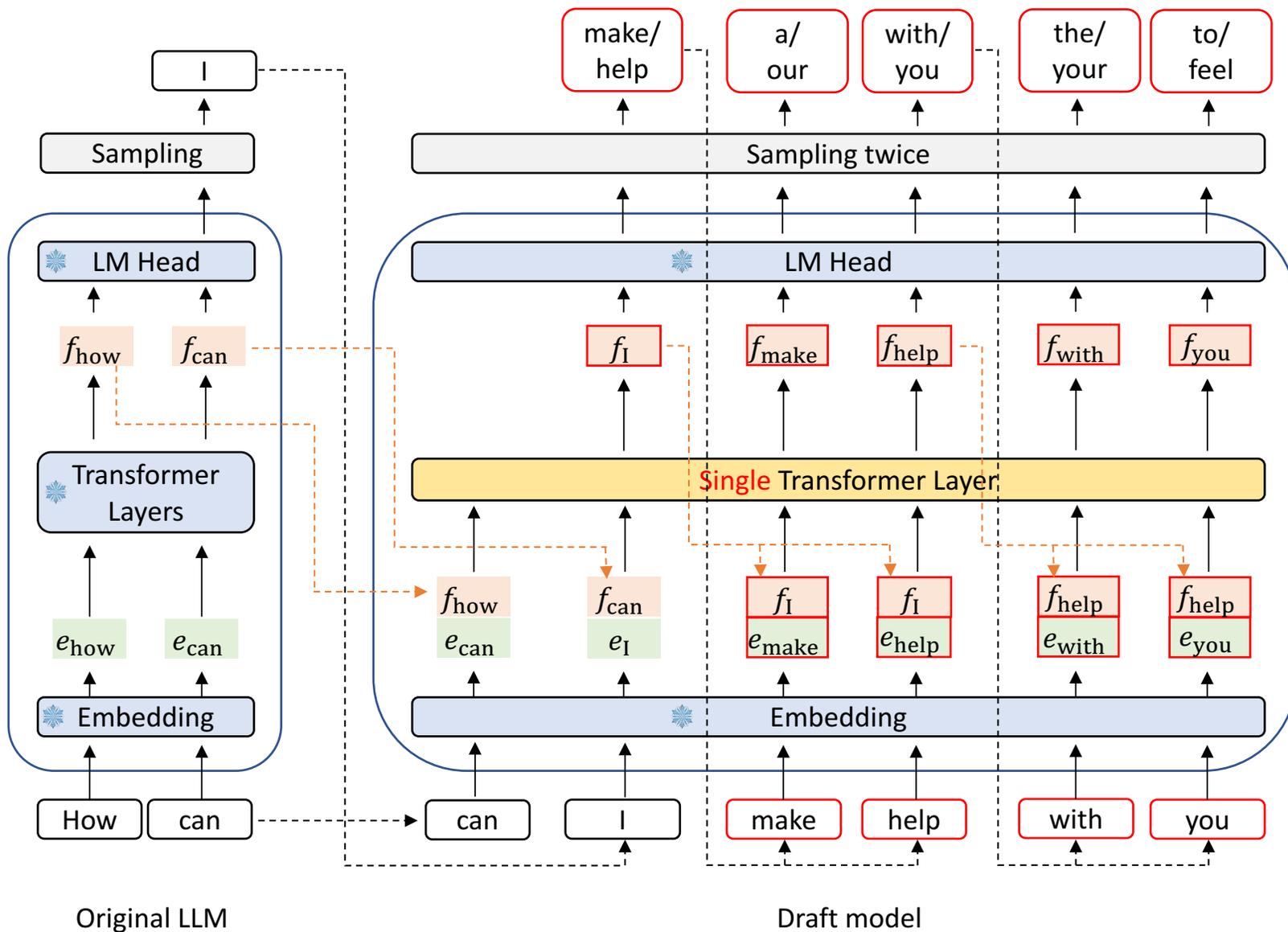
Flatten to 1D

It	is	has	a	the	to	good	be
----	----	-----	---	-----	----	------	----

Attention mask

	It	is	has	a	the	to	good	be
It	✓							
is	✓	✓						
has	✓		✓					
a	✓	✓		✓				
the	✓	✓			✓			
to	✓		✓			✓		
good	✓	✓		✓			✓	
be	✓		✓			✓		✓

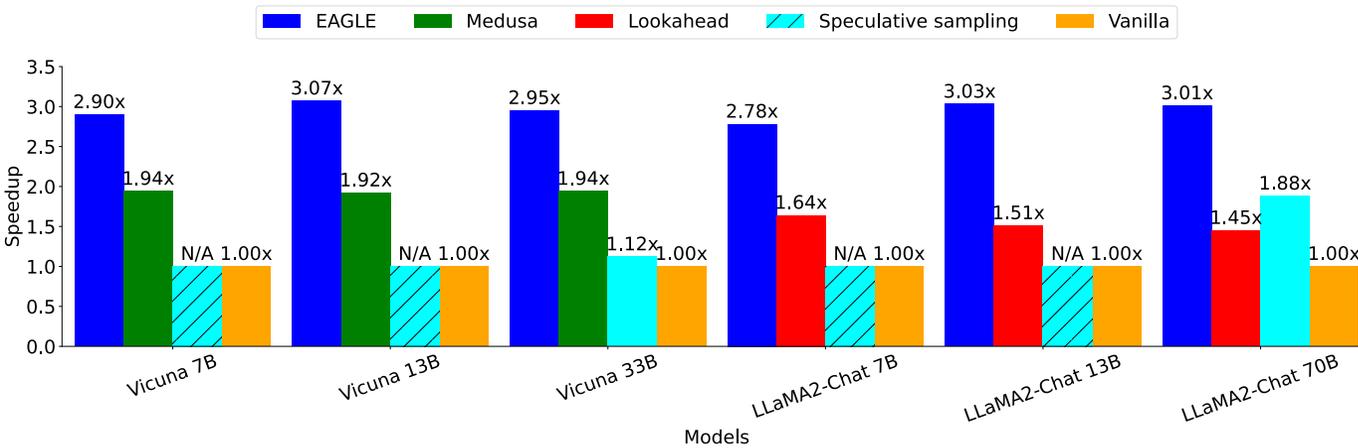
#Parameters of the drafted models



#Parameters (Original LLM)	#Parameters (Draft model)	Ratio
7B	0.24B	3.4%
13B	0.37B	2.8%
33B	0.56B	1.7%
70B	0.99B	1.4%

- Trained on RTX 3090 GPUs on ShareGPT for 1 - 2 days

Performance on MT-bench



On MT-Bench, EAGLE is

- 3x 🚀 than vanilla decoding
- 1.6x 🚀 than Medusa
- 2x 🚀 than Lookahead
- **Provably** maintaining text distribution

Method	Speed (tokens/s)	Completion Rate
Vanilla	3.46	1.00
Medusa	3.37	1.00
Lookahead	6.09	1.00
EAGLE	3.32	1.00

Third-party evaluations

Spec-bench



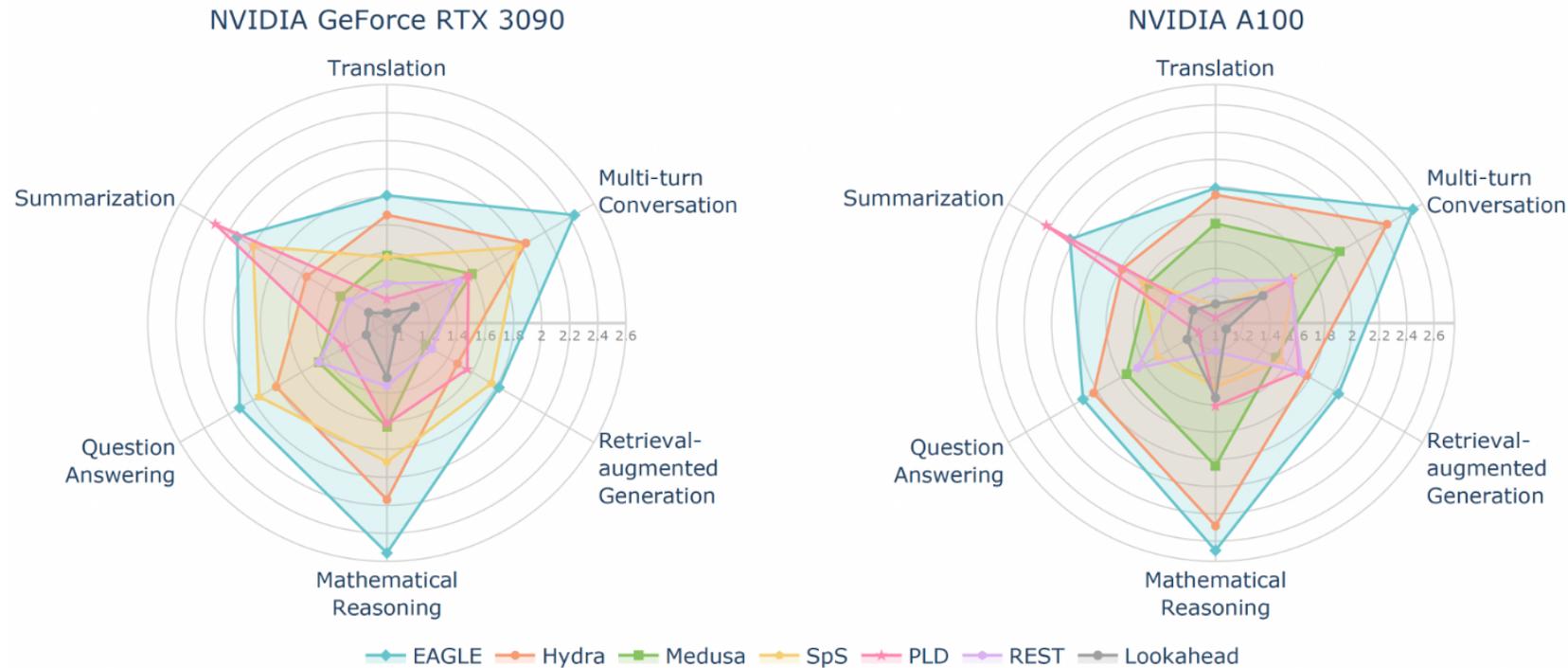
Unlocking Efficiency in Large Language Model Inference: A Comprehensive Survey of Speculative Decoding

Heming Xia¹, Zhe Yang², Qingxiu Dong², Peiyi Wang²,
Yongqi Li¹, Tao Ge³, Tianyu Liu⁴, Wenjie Li¹, Zhifang Sui²

¹Department of Computing, The Hong Kong Polytechnic University

²National Key Laboratory for Multimedia Information Processing, Peking University

³Microsoft Research Asia ⁴Alibaba Group



Speedup comparison of Speculative Decoding methods on Spec-Bench, evaluated by Vicuna-7B-v1.3.

Third-party evaluations

Spec-bench

RTX 3090, Vicuna 7B

Models	Multi-turn Conversation	Translation	Summarization	Question Answering	Mathematical Reasoning	Retrieval-aug. Generation	#Mean Accepted Tokens	Overall
EAGLE 	2.44x	1.81x	2.13x	2.11x	2.54x	1.82x	3.57	2.16x
SpS 	1.98x	1.37x	2.00x	1.95x	1.89x	1.76x	2.29	1.83x
Hydra 	2.04x	1.67x	1.56x	1.81x	2.16x	1.48x	3.26	1.80x
PLD	1.57x	1.07x	2.31x	1.25x	1.62x	1.56x	1.74	1.55x
Medusa	1.60x	1.38x	1.28x	1.46x	1.64x	1.22x	2.32	1.44x
REST	1.49x	1.18x	1.21x	1.46x	1.35x	1.27x	1.63	1.32x
Lookahead	1.13x	0.97x	1.05x	1.07x	1.29x	0.98x	1.65	1.08x



Unlocking Efficiency in Large Language Model Inference: A Comprehensive Survey of Speculative Decoding

Heming Xia¹, Zhe Yang², Qingxiu Dong², Peiyi Wang²,
Yongqi Li¹, Tao Ge³, Tianyu Liu⁴, Wenjie Li¹, Zhifang Sui²

¹Department of Computing, The Hong Kong Polytechnic University

²National Key Laboratory for Multimedia Information Processing, Peking University

³Microsoft Research Asia ⁴Alibaba Group

Third-party evaluations

Spec-bench

 **Unlocking Efficiency in Large Language Model Inference:
A Comprehensive Survey of Speculative Decoding**

Heming Xia¹, Zhe Yang², Qingxiu Dong², Peiyi Wang²,
Yongqi Li¹, Tao Ge³, Tianyu Liu⁴, Wenjie Li¹, Zhifang Sui²

¹Department of Computing, The Hong Kong Polytechnic University

²National Key Laboratory for Multimedia Information Processing, Peking University

³Microsoft Research Asia ⁴Alibaba Group

A100, Vicuna 7B

Models	Models	Multi-turn Conversation	Translation	Summarization	Question Answering	Mathematical Reasoning	Retrieval-aug. Generation	#Mean Accepted Tokens	Overall
EAGLE 	EAGLE 	2.67x	1.99x	2.23x	2.12x	2.67x	2.04x	3.61	2.29x
SpS 	Hydra 	2.45x	1.94x	1.79x	2.03x	2.49x	1.77x	3.24	2.09x
Hydra 	Medusa 	2.05x	1.73x	1.57x	1.75x	2.05x	1.51x	2.32	1.78x
PLD	PLD	1.64x	1.04x	2.43x	1.14x	1.61x	1.71x	1.73	1.59x
Medusa	SpS	1.66x	1.13x	1.62x	1.49x	1.47x	1.55x	2.28	1.49x
REST	REST	1.63x	1.31x	1.36x	1.66x	1.21x	1.73x	1.82	1.48x
Lookahead	Lookahead	1.40x	1.14x	1.19x	1.24x	1.55x	1.09x	1.66	1.27x

Third-party evaluations

Spec-bench



Unlocking Efficiency in Large Language Model Inference: A Comprehensive Survey of Speculative Decoding

Heming Xia¹, Zhe Yang², Qingxiu Dong², Peiyi Wang²,
Yongqi Li¹, Tao Ge³, Tianyu Liu⁴, Wenjie Li¹, Zhifang Sui²

¹Department of Computing, The Hong Kong Polytechnic University

²National Key Laboratory for Multimedia Information Processing, Peking University

³Microsoft Research Asia ⁴Alibaba Group

A100, Vicuna 13B

Models	Models	Models	Multi-turn Conversation	Translation	Summarization	Question Answering	Mathematical Reasoning	Retrieval-aug. Generation	#Mean Accepted Tokens	Overall
EAGLE 🌟	EAGLE 🌟	EAGLE 🌟	2.68x	1.96x	2.44x	2.04x	2.70x	2.23x	3.64	2.34x
SpS ②	Hydra ②	Hydra ②	2.46x	1.90x	1.93x	1.96x	2.48x	1.92x	3.35	2.12x
Hydra ③	Medusa ③	Medusa ③	1.96x	1.66x	1.63x	1.63x	2.00x	1.58x	2.39	1.75x
PLD	PLD	SpS	1.60x	1.13x	1.68x	1.39x	1.53x	1.67x	2.18	1.49x
Medusa	SpS	PLD	1.47x	1.02x	2.19x	1.03x	1.57x	1.71x	1.68	1.48x
REST	REST	REST	1.52x	1.17x	1.37x	1.53x	1.19x	1.55x	1.82	1.38x
Lookahead	Lookahead	Lookahead	1.30x	1.06x	1.20x	1.12x	1.48x	1.12x	1.63	1.22x

Third-party evaluations

Spec-bench

 **Unlocking Efficiency in Large Language Model Inference:
A Comprehensive Survey of Speculative Decoding**

Heming Xia¹, Zhe Yang², Qingxiu Dong², Peiyi Wang²,
Yongqi Li¹, Tao Ge³, Tianyu Liu⁴, Wenjie Li¹, Zhifang Sui²

¹Department of Computing, The Hong Kong Polytechnic University

²National Key Laboratory for Multimedia Information Processing, Peking University

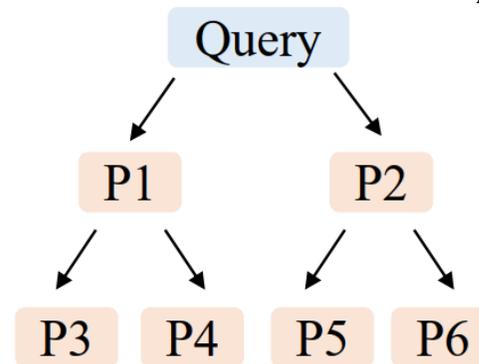
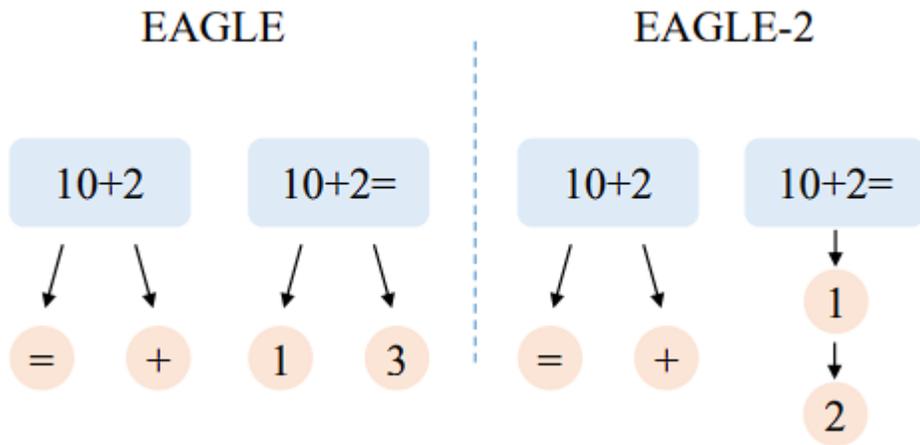
³Microsoft Research Asia ⁴Alibaba Group

A100, Vicuna 33B

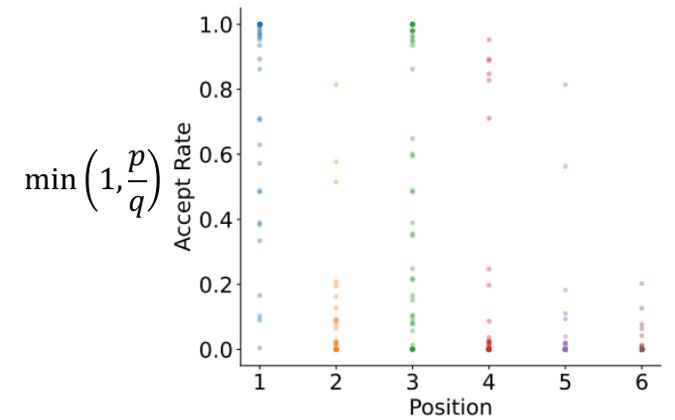
Models	Models	Models	Models	Multi-turn Conversation	Translation	Summarization	Question Answering	Mathematical Reasoning	Retrieval-aug. Generation	#Mean Accepted Tokens	Overall
EAGLE 	EAGLE 	EAGLE 	EAGLE 	2.79x	2.05x	2.51x	2.17x	2.99x	2.27x	3.39	2.47x
SpS 	Hydra 	Hydra 	Hydra 	2.59x	2.01x	2.04x	2.11x	2.71x	2.06x	3.24	2.26x
Hydra 	Medusa 	Medusa 	Medusa 	1.98x	1.73x	1.64x	1.66x	2.07x	1.62x	2.33	1.79x
PLD	PLD	SpS	SpS	1.75x	1.28x	1.76x	1.53x	1.69x	1.68x	2.01	1.61x
Medusa	SpS	PLD	REST	1.63x	1.27x	1.45x	1.61x	1.30x	1.61x	1.80	1.48x
REST	REST	REST	PLD	1.44x	1.06x	2.00x	1.07x	1.55x	1.45x	1.55	1.42x
Lookahead	Lookahead	Lookahead	Lookahead	1.32x	1.08x	1.20x	1.16x	1.54x	1.15x	1.61	1.24x

EAGLE-2: Faster Inference of Language Models with Dynamic Draft Trees

Static tree structure?



(a) Draft tree structure.

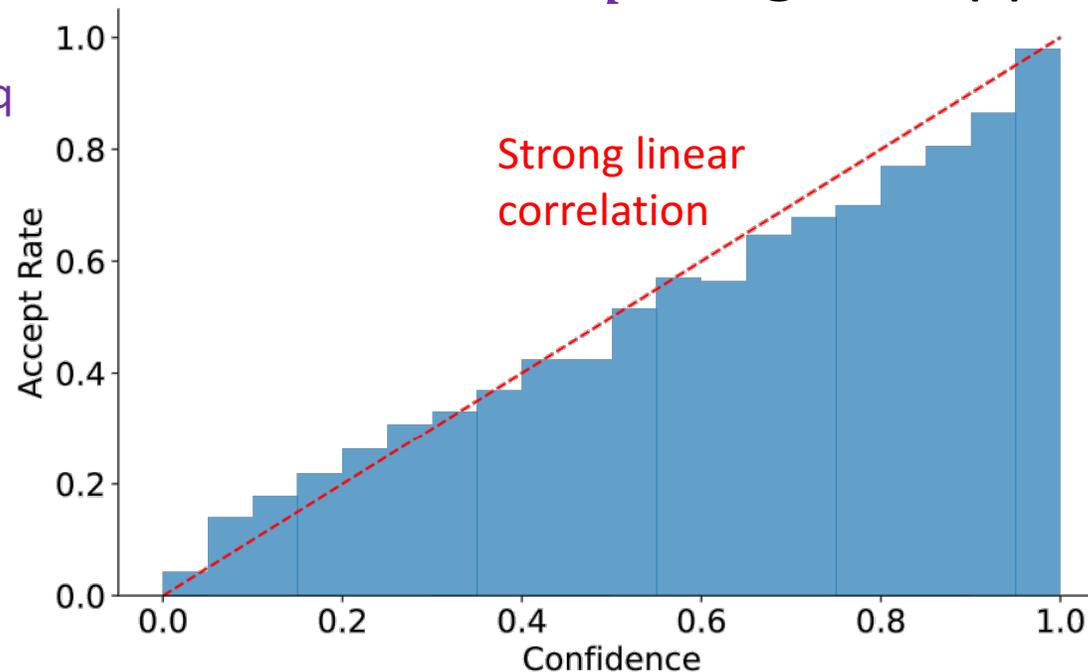


(b) Acceptance rates of tokens at different positions, with each point representing a query.

How to determine importance of each node?

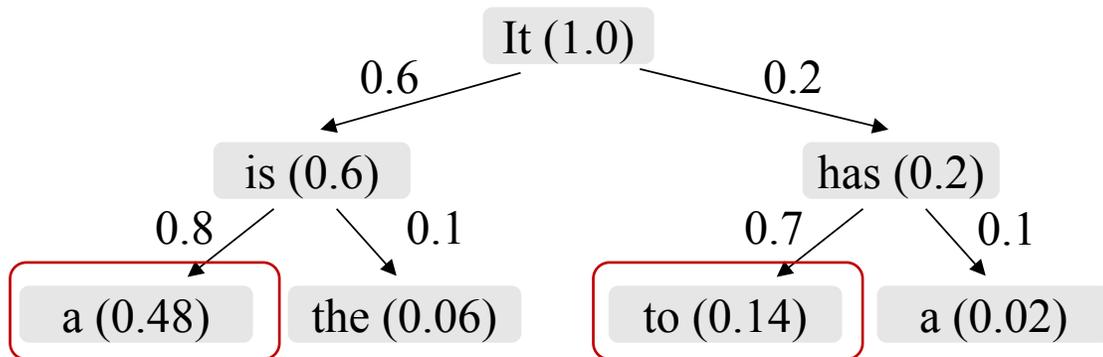
- Accept rate $\min\left(1, \frac{p(t)}{q(t)}\right)$
- However, it requires the computation from the original **large** model p
- The **confidence of the draft model q** is a good approximation

i.e. output
probability of q

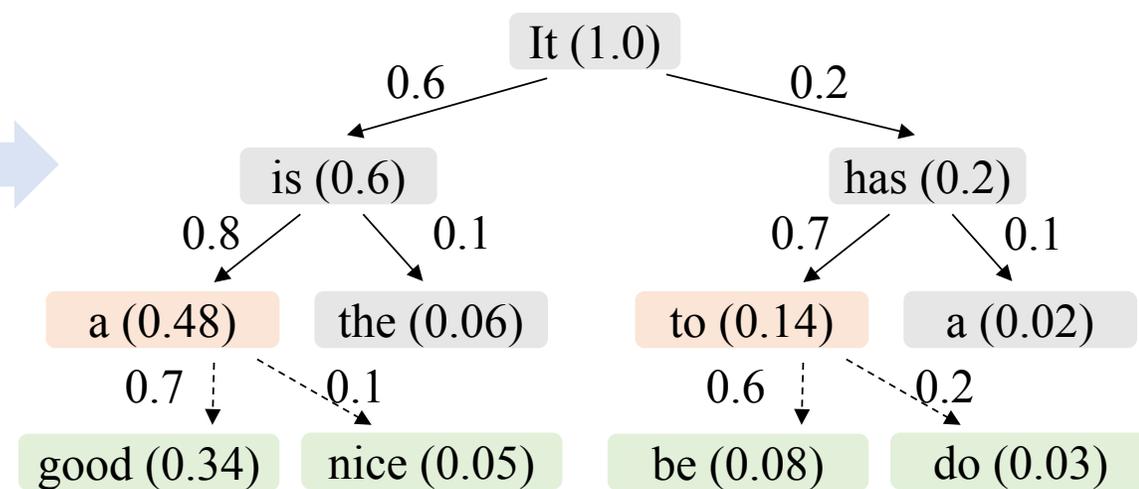


Context-aware dynamic draft tree

Beam Search (Top-2)

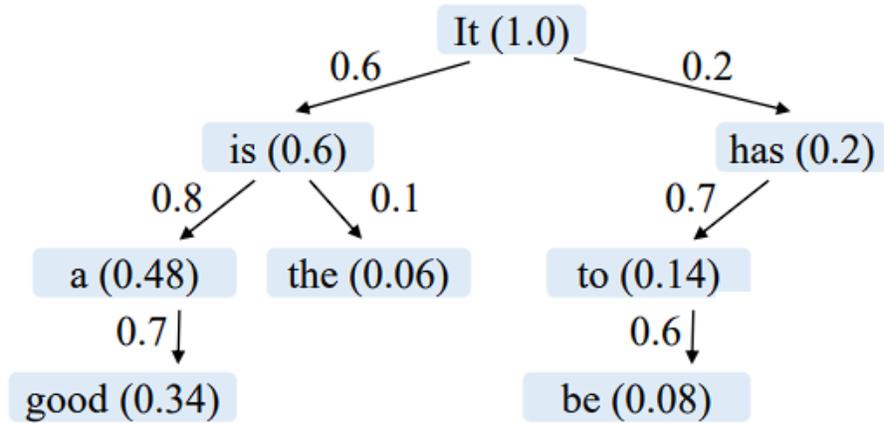


Expand



Context-aware dynamic draft tree

Rerank (Top-8)



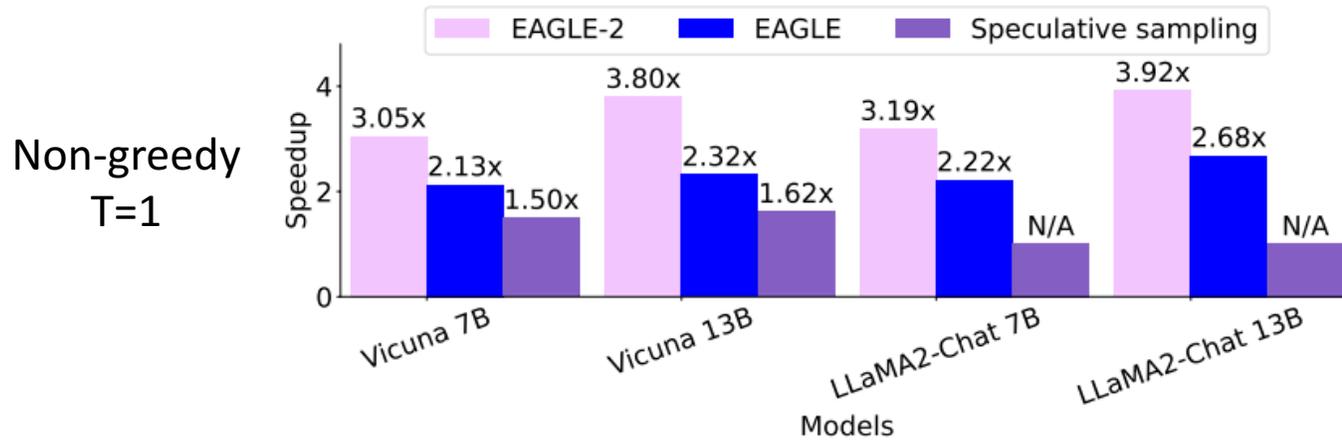
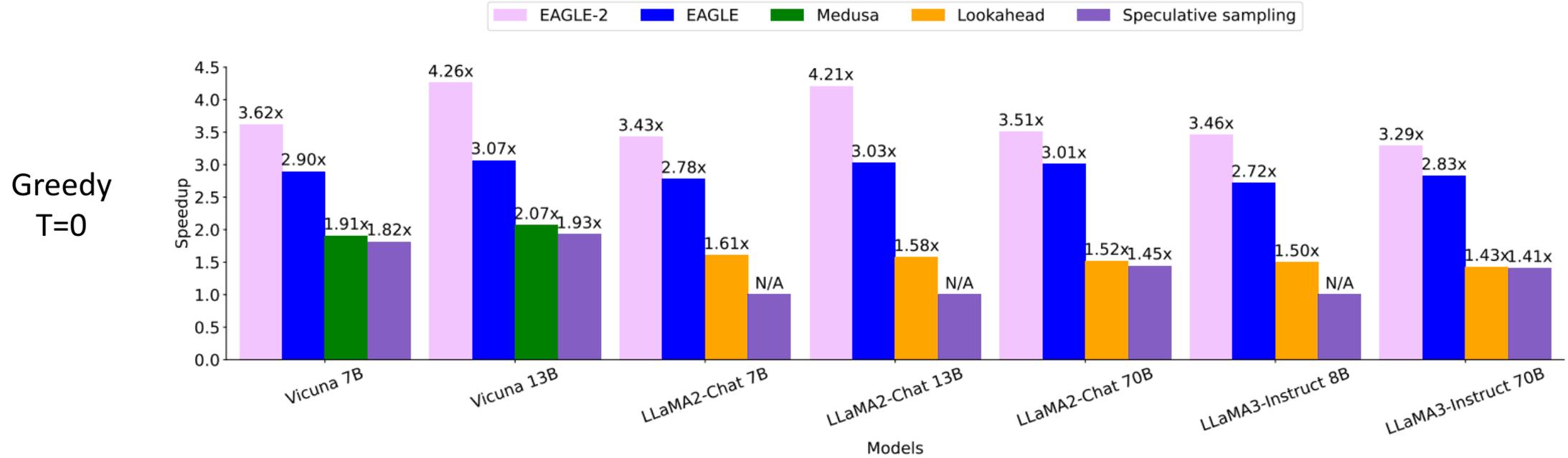
Flatten to 1D

It	is	has	a	the	to	good	be
----	----	-----	---	-----	----	------	----

Attention mask

	It	is	has	a	the	to	good	be
It	✓							
is	✓	✓						
has	✓		✓					
a	✓	✓		✓				
the	✓	✓			✓			
to	✓		✓			✓		
good	✓	✓		✓			✓	
be	✓		✓			✓		✓

Performance on MT-bench



Performance (T=0, bs=1)

Model	Method	MT-bench		HumanEval		GSM8K		Alpaca		CNN/DM		Natural Ques.		Mean	
		Speedup	#mean accepted tokens	Speedup	#mean accepted tokens	Speedup	#mean accepted tokens								
V 13B	SpS	1.93x	2.27	2.23x	2.57	1.77x	2.01	1.76x	2.03	1.93x	2.33	1.66x	1.88	1.88x	2.18
	PLD	1.58x	1.63	1.85x	1.93	1.68x	1.73	1.16x	1.19	2.42x	2.50	1.14x	1.17	1.64x	1.69
	Medusa	2.07x	2.59	2.50x	2.78	2.23x	2.64	2.08x	2.45	1.71x	2.09	1.81x	2.10	2.07x	2.44
	Lookahead	1.65x	1.69	1.71x	1.75	1.81x	1.90	1.46x	1.51	1.46x	1.50	1.36x	1.39	1.58x	1.62
	Hydra	2.88x	3.65	3.28x	3.87	2.93x	3.66	2.86x	3.53	2.05x	2.81	2.11x	2.88	2.69x	3.40
	EAGLE	3.07x	3.98	3.58x	4.39	3.08x	3.97	3.03x	3.95	2.49x	3.52	2.42x	3.11	2.95x	3.82
	EAGLE-2	4.26x	4.83	4.96x	5.41	4.22x	4.79	4.25x	4.89	3.40x	4.21	3.13x	3.74	4.04x	4.65
L2 13B	PLD	1.42x	1.46	1.63x	1.70	1.41x	1.44	1.16x	1.20	1.42x	1.45	1.12x	1.15	1.36x	1.40
	Lookahead	1.58x	1.64	1.80x	1.85	1.65x	1.69	1.47x	1.50	1.46x	1.53	1.42x	1.45	1.56x	1.61
	EAGLE	3.03x	3.90	3.76x	4.52	3.20x	4.03	3.01x	3.83	2.70x	3.59	2.83x	3.47	3.09x	3.89
	EAGLE-2	4.21x	4.75	5.00x	5.52	4.31x	4.90	4.13x	4.61	3.45x	4.24	3.51x	4.04	4.10x	4.68
V 7B	SpS	1.82x	2.36	1.99x	2.61	1.71x	2.26	1.65x	2.21	1.81x	2.44	1.60x	2.16	1.76x	2.34
	PLD	1.61x	1.68	1.82x	1.87	1.82x	1.99	1.21x	1.31	2.53x	2.72	1.23x	1.44	1.70x	1.84
	Medusa	1.91x	2.52	2.02x	2.67	1.89x	2.59	1.79x	2.48	1.42x	2.02	1.51x	2.09	1.76x	2.40
	Lookahead	1.63x	1.69	1.72x	1.77	1.84x	1.99	1.38x	1.57	1.44x	1.53	1.45x	1.60	1.58x	1.69
	Hydra	2.69x	3.60	2.98x	3.79	2.73x	3.66	2.66x	3.58	2.01x	2.70	2.25x	2.86	2.55x	3.37
	EAGLE	2.90x	3.94	3.33x	4.29	3.01x	4.00	2.79x	3.89	2.33x	3.42	2.31x	3.21	2.78x	3.79
	EAGLE-2	3.62x	4.98	3.95x	5.33	3.63x	4.97	3.46x	4.86	2.94x	4.12	2.76x	3.82	3.39x	4.68
L2 7B	PLD	1.38x	1.43	1.52x	1.59	1.32x	1.37	1.15x	1.19	1.48x	1.52	1.15x	1.20	1.33x	1.38
	Lookahead	1.61x	1.66	1.72x	1.77	1.58x	1.65	1.49x	1.52	1.49x	1.54	1.48x	1.53	1.56x	1.61
	EAGLE	2.78x	3.62	3.17x	4.24	2.91x	3.82	2.78x	3.71	2.43x	3.41	2.61x	3.44	2.78x	3.71
	EAGLE-2	3.43x	4.70	4.03x	5.39	3.52x	4.77	3.45x	4.66	3.01x	4.12	3.15x	4.19	3.43x	4.64

Performance (T=1, bs=1)

Model	Method	MT-bench		HumanEval		GSM8K		Alpaca		CNN/DM		Natural Ques.		Mean	
		Speedup	#mean accepted tokens	Speedup	#mean accepted tokens	Speedup	#mean accepted tokens								
V 13B	SpS	1.62x	1.84	1.72x	1.97	1.46x	1.73	1.52x	1.78	1.66x	1.89	1.43x	1.70	1.55x	1.82
	EAGLE	2.32x	3.20	2.65x	3.63	2.57x	3.60	2.45x	3.57	2.23x	3.26	2.14x	3.06	2.39x	3.39
	EAGLE-2	3.80x	4.40	4.22x	4.89	3.77x	4.41	3.78x	4.37	3.25x	3.97	3.07x	3.54	3.65x	4.26
L2 13B	EAGLE	2.68x	3.45	2.89x	3.78	2.82x	3.67	2.66x	3.55	2.41x	3.39	2.37x	3.31	2.64x	3.53
	EAGLE-2	3.92x	4.51	4.58x	5.29	4.21x	4.80	3.85x	4.48	3.31x	4.08	3.43x	3.89	3.88x	4.51
V 7B	SpS	1.50x	1.87	1.55x	1.95	1.53x	1.82	1.56x	1.85	1.63x	1.91	1.33x	1.72	1.52x	1.85
	EAGLE	2.13x	3.17	2.39x	3.43	2.34x	3.29	2.21x	3.30	2.08x	3.12	1.95x	2.86	2.18x	3.20
	EAGLE-2	3.05x	4.28	3.33x	4.65	3.07x	4.49	3.08x	4.43	2.63x	3.76	2.48x	3.56	2.94x	4.20
L2 7B	EAGLE	2.22x	3.30	2.61x	3.79	2.40x	3.52	2.29x	3.33	2.19x	3.15	2.22x	3.12	2.32x	3.37
	EAGLE-2	3.19x	4.41	3.67x	5.06	3.35x	4.62	3.20x	4.48	2.73x	3.85	2.81x	4.01	3.15x	4.41

Vanilla on A100 vs EAGLE-2 on RTX3060

Vanilla



EAGLE-2

A100 (\$10000)

RTX 3060 (2 × \$300)

Speed	Compression Ratio
26.06 tokens/s	1.00

Speed	Compression Ratio
19.61 tokens/s	5.00

Introduce artificial intelligence to me.

I'm excited

Introduce artificial intelligence to me.

I'm excited to introduce

Third-party evaluations (updated on Oct. 25, 2024)

 *Unlocking Efficiency in Large Language Model Inference:
A Comprehensive Survey of Speculative Decoding*

Heming Xia¹, Zhe Yang², Qingxiu Dong², Peiyi Wang²,
Yongqi Li¹, Tao Ge³, Tianyu Liu⁴, Wenjie Li¹, Zhifang Sui²

¹Department of Computing, The Hong Kong Polytechnic University

²National Key Laboratory for Multimedia Information Processing, Peking University

³Microsoft Research Asia ⁴Alibaba Group

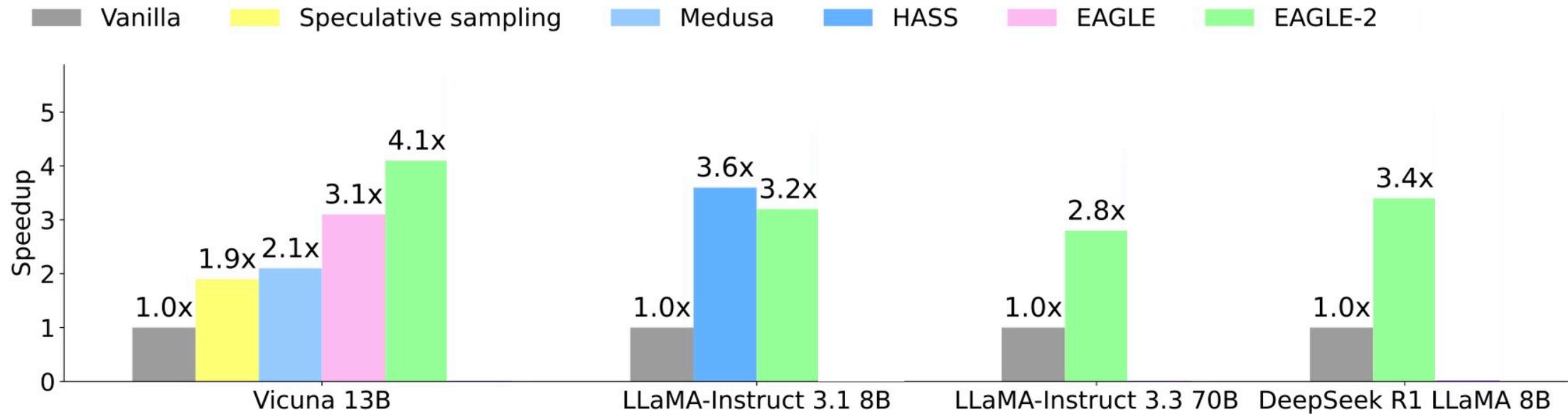
Leaderboard on 3090

- Device: a single NVIDIA GeForce RTX 3090 GPU (24GB) with 12 CPU cores
- Testing environment: Pytorch 2.0.1, under CUDA 11.8
- Experimental Settings: Vicuna-7B-v1.3, greedy decoding, FP16 precision, batch size = 1

Models	Multi-turn Conversation	Translation	Summa- rization	Question Answering	Mathematical Reasoning	Retrieval- aug. Generation	#Mean Accepted Tokens	Overall
EAGLE2 🏆	2.71x	1.82x	2.19x	2.11x	2.71x	1.91x	4.36	2.25x
EAGLE 🥈	2.44x	1.81x	2.13x	2.11x	2.54x	1.82x	3.57	2.16x
SpS 🥉	1.98x	1.37x	2.00x	1.95x	1.89x	1.76x	2.29	1.83x
Hydra	2.04x	1.67x	1.56x	1.81x	2.16x	1.48x	3.26	1.80x
PLD	1.57x	1.07x	2.31x	1.25x	1.62x	1.56x	1.74	1.55x
Medusa	1.60x	1.38x	1.28x	1.46x	1.64x	1.22x	2.32	1.44x
REST	1.49x	1.18x	1.21x	1.46x	1.35x	1.27x	1.63	1.32x
Lookahead	1.13x	0.97x	1.05x	1.07x	1.29x	0.98x	1.65	1.08x

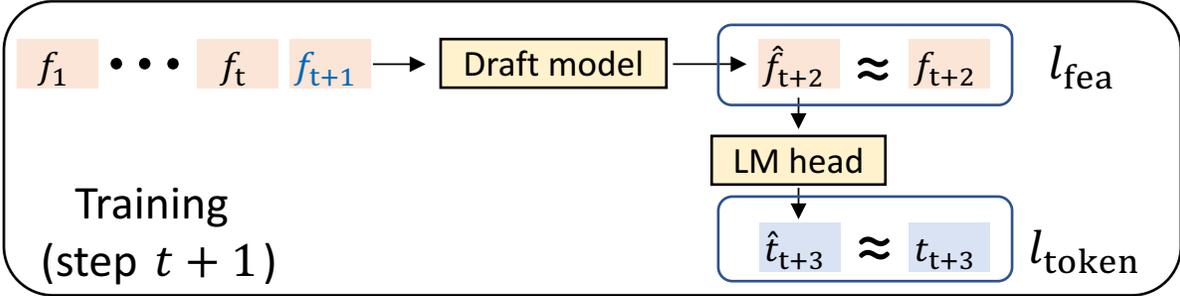
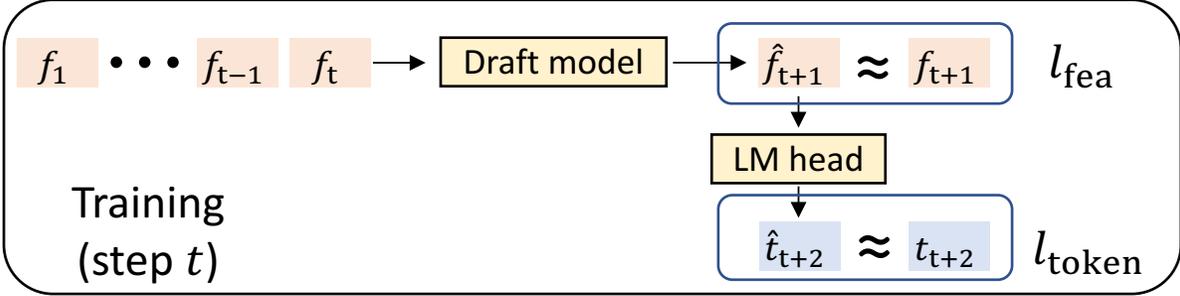
EAGLE-3

Benchmarking on MT-Bench

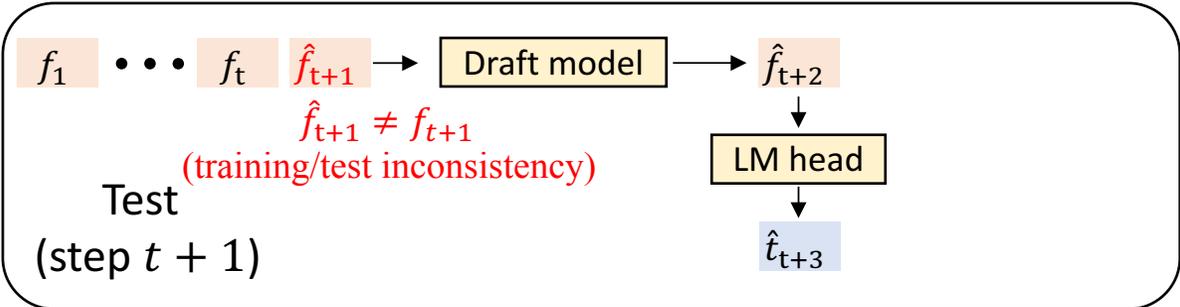
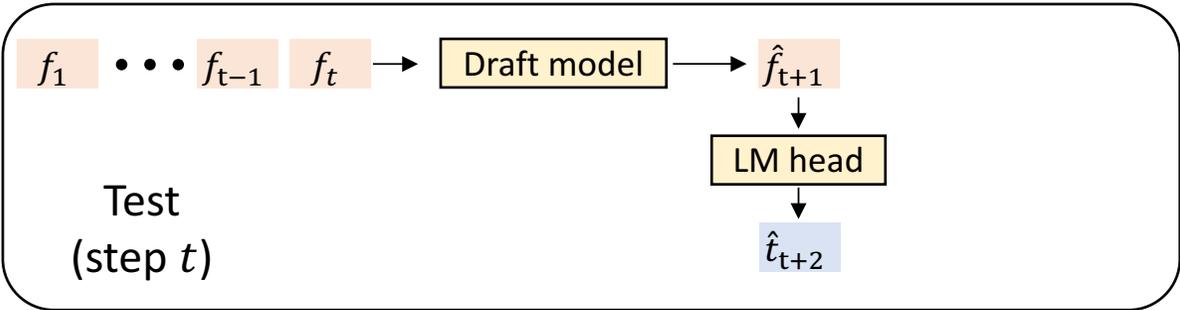


EAGLE vs EAGLE-3

$$\min l_{\text{fea}} + 0.1l_{\text{token}}$$



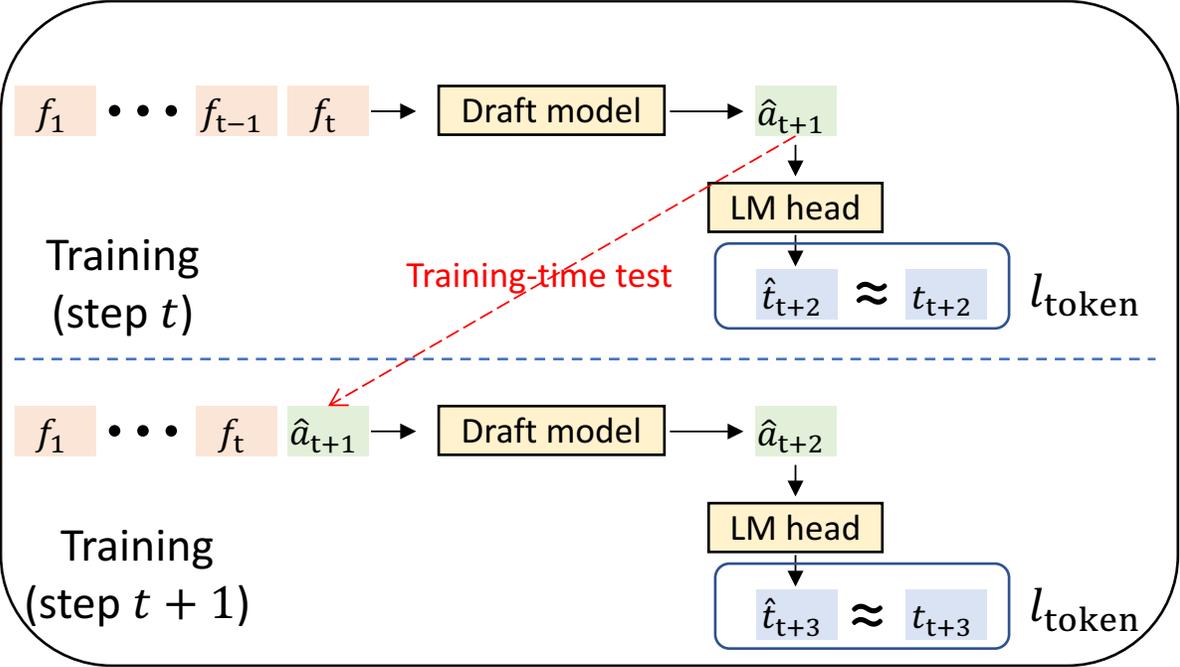
EAGLE-1 training



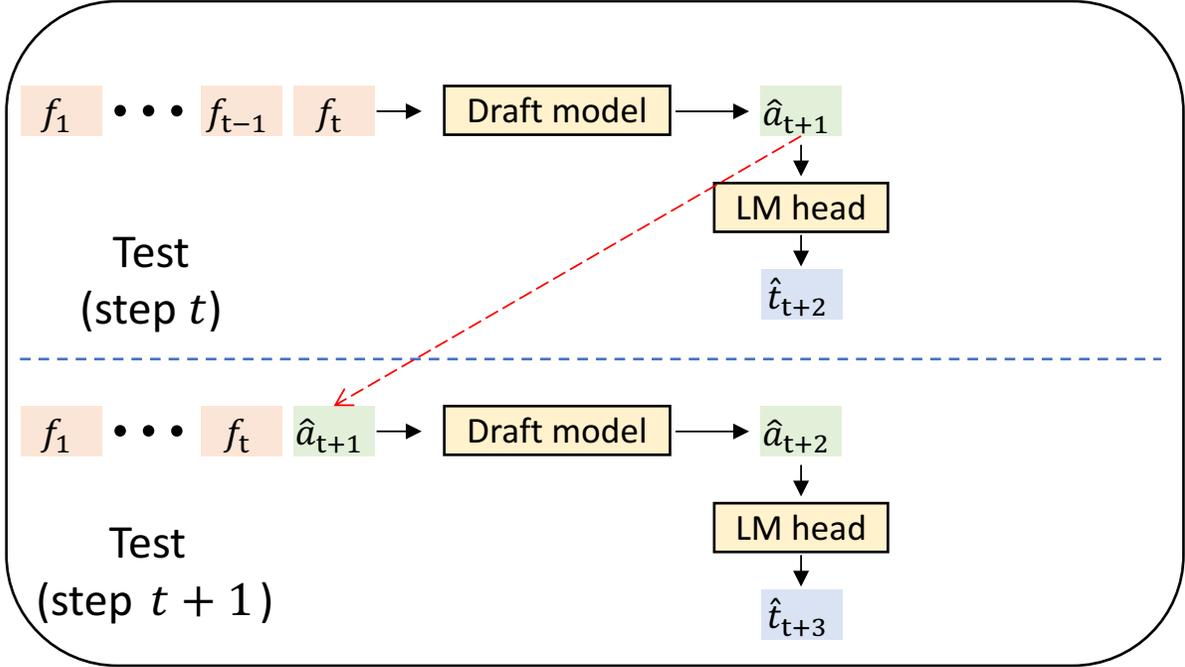
EAGLE-1 test

EAGLE vs EAGLE-3

$\min l_{\text{token}}$

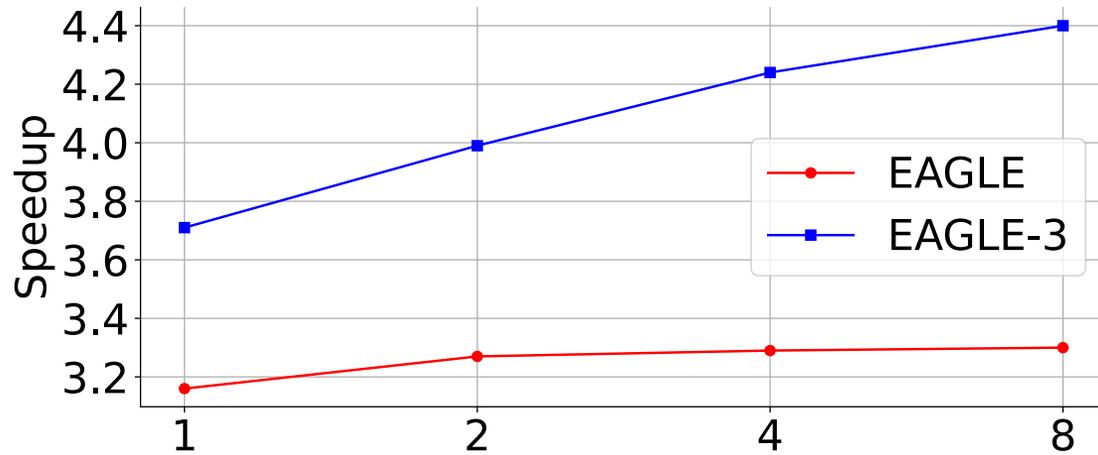


EAGLE-3 training

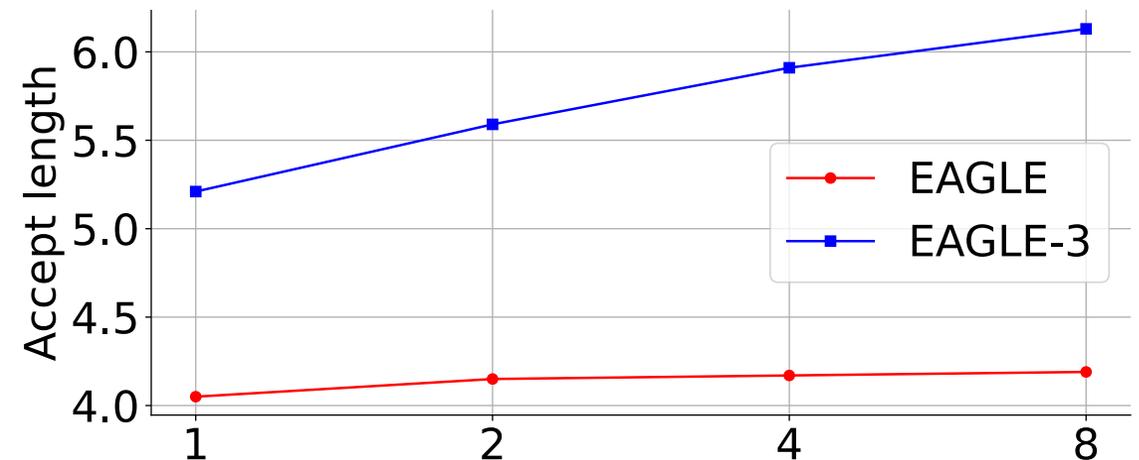


EAGLE-3 test

New Scaling Law for Inference Acceleration



Amount of training data (relative to ShareGPT)



Amount of training data (relative to ShareGPT)

Trained on UltraChat + ShareGPT

Performance (bs=1)

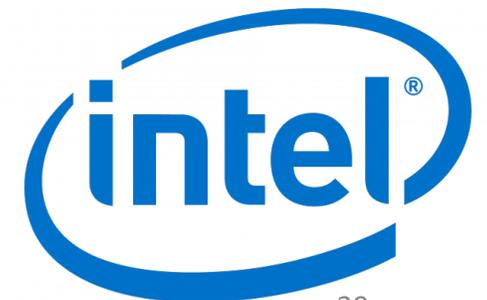
Model	Method	MT-bench		HumanEval		GSM8K		Alpaca		CNN/DM		Mean	
		Speedup	τ										
Temperature=0													
V 13B	SpS	1.93x	2.27	2.23x	2.57	1.77x	2.01	1.76x	2.03	1.93x	2.33	1.92x	2.24
	PLD	1.58x	1.63	1.85x	1.93	1.68x	1.73	1.16x	1.19	2.42x	2.50	1.74x	1.80
	Medusa	2.07x	2.59	2.50x	2.78	2.23x	2.64	2.08x	2.45	1.71x	2.09	2.12x	2.51
	Lookahead	1.65x	1.69	1.71x	1.75	1.81x	1.90	1.46x	1.51	1.46x	1.50	1.62x	1.67
	Hydra	2.88x	3.65	3.28x	3.87	2.93x	3.66	2.86x	3.53	2.05x	2.81	2.80x	3.50
	EAGLE	3.07x	3.98	3.58x	4.39	3.08x	3.97	3.03x	3.95	2.49x	3.52	3.05x	3.96
	EAGLE-2	4.26x	4.83	4.96x	5.41	4.22x	4.79	4.25x	4.89	3.40x	4.21	4.22x	4.83
	EAGLE-3	5.58x	6.65	6.47x	7.54	5.32x	6.29	5.16x	6.17	5.01x	6.47	5.51x	6.62
L31 8B	EAGLE-2	3.16x	4.05	3.66x	4.71	3.39x	4.24	3.28x	4.12	2.65x	3.45	3.23x	4.11
	EAGLE-3	4.40x	6.13	4.85x	6.74	4.48x	6.23	4.82x	6.70	3.65x	5.34	4.44x	6.23
L33 70B	EAGLE-2	2.83x	3.67	3.12x	4.09	2.83x	3.69	3.03x	3.92	2.44x	3.55	2.85x	3.78
	EAGLE-3	4.11x	5.63	4.79x	6.52	4.34x	6.15	4.30x	6.09	3.27x	5.02	4.12x	5.88
DSL 8B	EAGLE-2	2.92x	3.80	3.42x	4.29	3.40x	4.40	3.01x	3.80	3.53x	3.33	3.26x	3.92
	EAGLE-3	4.05x	5.58	4.59x	6.38	5.01x	6.93	3.65x	5.37	3.52x	4.92	4.16x	5.84
Temperature=1													
V 13B	SpS	1.62x	1.84	1.72x	1.97	1.46x	1.73	1.52x	1.78	1.66x	1.89	1.60x	1.84
	EAGLE	2.32x	3.20	2.65x	3.63	2.57x	3.60	2.45x	3.57	2.23x	3.26	2.44x	3.45
	EAGLE-2	3.80x	4.40	4.22x	4.89	3.77x	4.41	3.78x	4.37	3.25x	3.97	3.76x	4.41
	EAGLE-3	4.57x	5.42	5.15x	6.22	4.71x	5.58	4.49x	5.39	4.33x	5.72	4.65x	5.67
L31 8B	EAGLE-2	2.44x	3.16	3.39x	4.39	2.86x	3.74	2.83x	3.65	2.44x	3.14	2.80x	3.62
	EAGLE-3	3.07x	4.24	4.13x	5.82	3.32x	4.59	3.90x	5.56	2.99x	4.39	3.45x	4.92
L33 70B	EAGLE-2	2.73x	3.51	2.89x	3.81	2.52x	3.36	2.77x	3.73	2.32x	3.27	2.65x	3.54
	EAGLE-3	3.96x	5.45	4.36x	6.16	4.17x	5.95	4.14x	5.87	3.11x	4.88	3.95x	5.66
DSL 8B	EAGLE-2	2.69x	3.41	3.01x	3.82	3.16x	4.05	2.64x	3.29	2.35x	3.13	2.77x	3.54
	EAGLE-3	3.20x	4.49	3.77x	5.28	4.38x	6.10	3.16x	4.30	3.08x	4.27	3.52x	4.89

Throughput (compared to SGLang w/o EAGLE)

Batch size	2	4	8	16	24	32	48	56	64
EAGLE	1.40x	1.38x	1.23x	1.02x	0.93x	0.94x	0.88x	0.99x	0.99x
EAGLE-3	1.81x	1.82x	1.62x	1.48x	1.39x	1.32x	1.38x	1.34x	1.38x

EAGLE in the community

- [SGLang](#)
- [vLLM](#)
- [AWS NeuronX Distributed Core](#)
- [Intel® Extension for Transformers](#)
- [Intel® LLM Library for PyTorch](#)
- [MLC-LLM](#)
- [NVIDIA TensorRT-LLM](#)



How to use?

- Code: <https://github.com/SafeAILab/EAGLE>



```
from eagle.model.ea_model import EaModel
model = EaModel.from_pretrained(base_model_path=base_model_path,
                                ea_model_path=EAGLE_model_path,
                                torch_dtype=torch.float16)
output_ids = model.eagenerate(input_ids,temperature=0.5,max_new_tokens=512)
```

Summary

- EAGLE-1
 - Next feature prediction
 - 3x latency speedup
- EAGLE-2
 - Dynamic draft tree
 - 4x latency speedup
- EAGLE-3
 - Training-time test, a new scaling law
 - 5x-6x latency speedup
- Can we be even faster?

Questions

?

?

Answers

?