

# The EAGLE Series: Lossless Inference Acceleration for LLMs

Hongyang Zhang

University of Waterloo & Vector Institute for AI

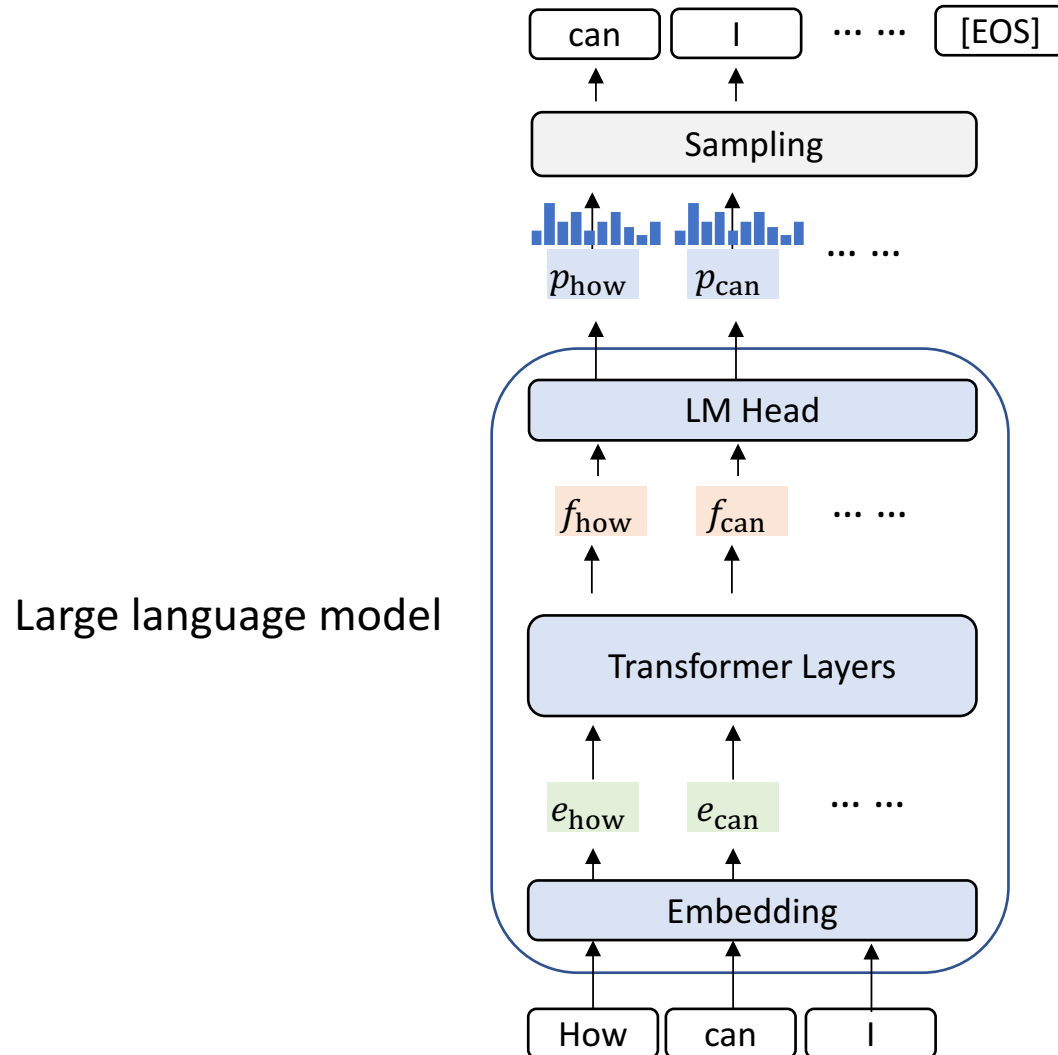


UNIVERSITY OF  
**WATERLOO**

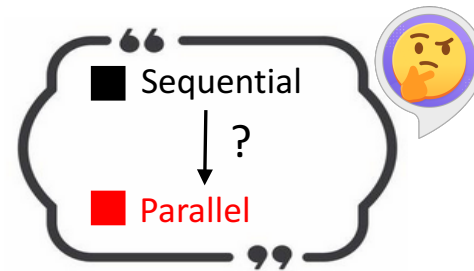


VECTOR  
INSTITUTE

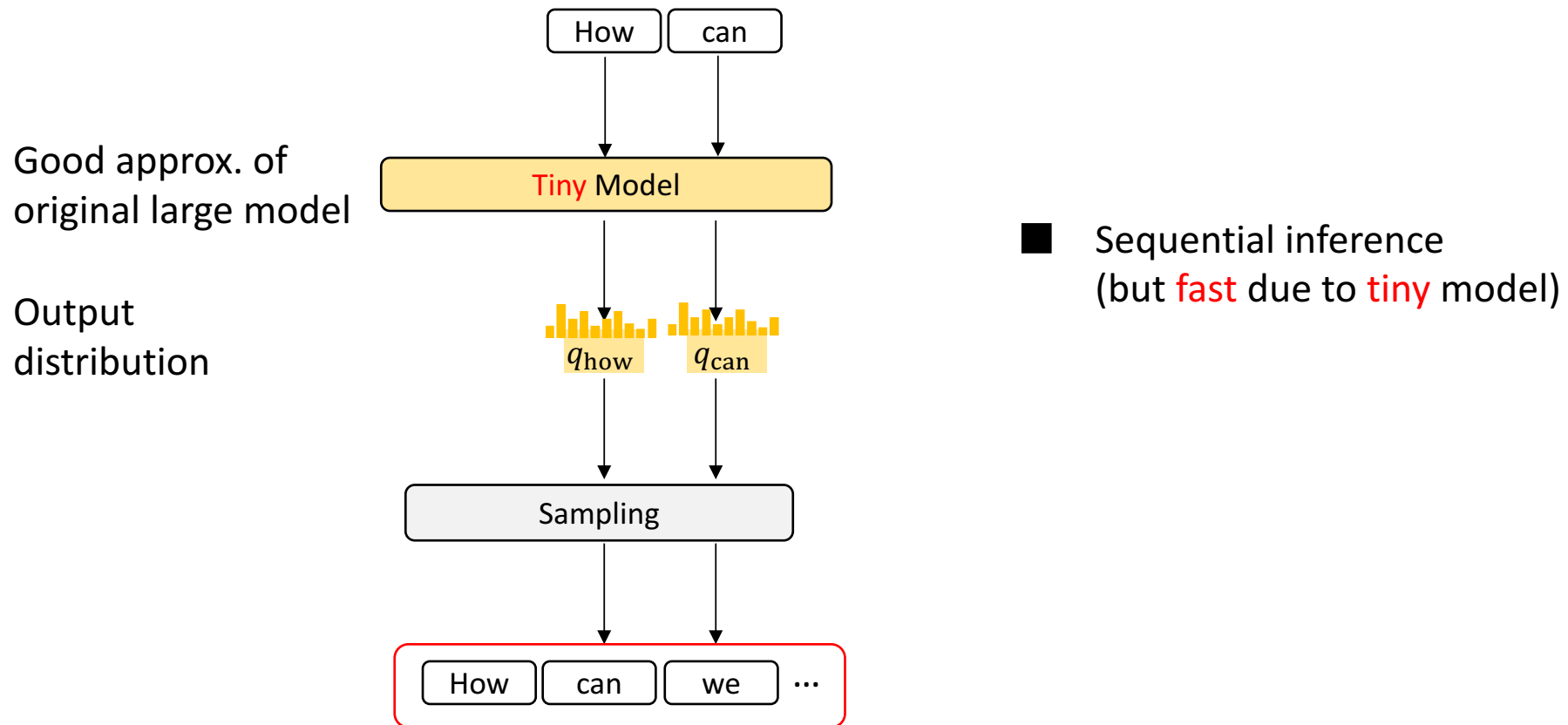
# Vanilla autoregressive inference



## ■ Sequential inference

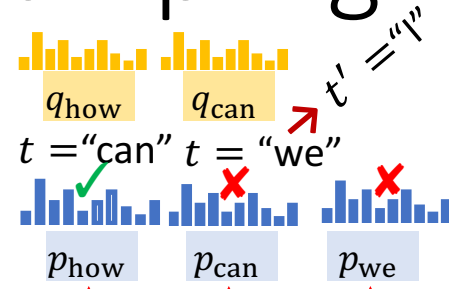


# Speculative sampling framework (draft)



# Speculative sampling framework (check)

Compare with **tiny** model



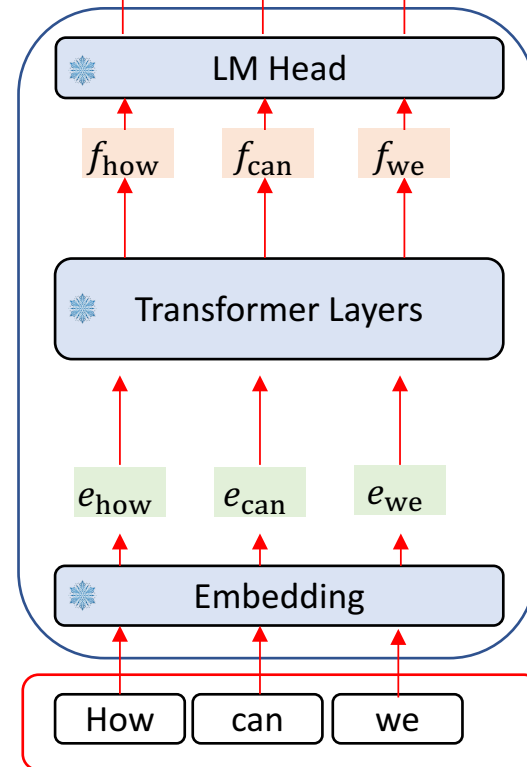
1.  $r \sim U(0,1)$ , if  $r < \min\left(1, \frac{p(t)}{q(t)}\right)$ , next token =  $t$   
Accept rate

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

■ Check in parallel

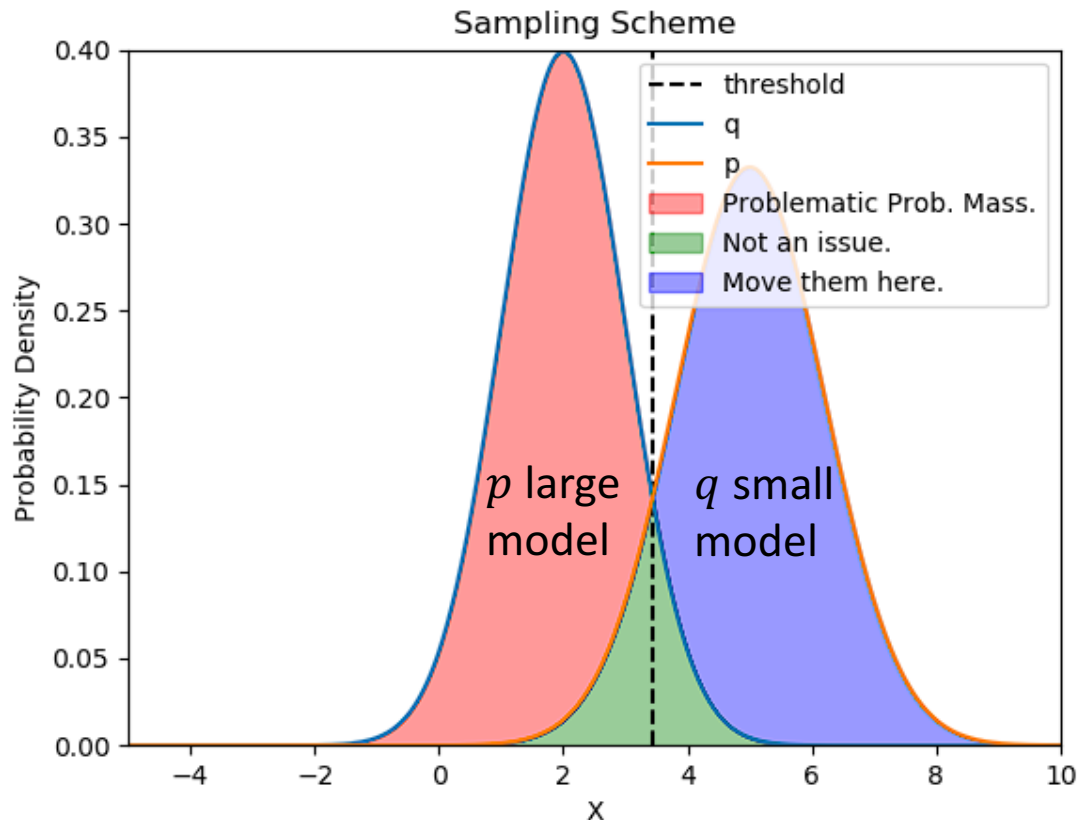


Original **large** model



Drafted by **tiny** model

# Speculative sampling framework (check)

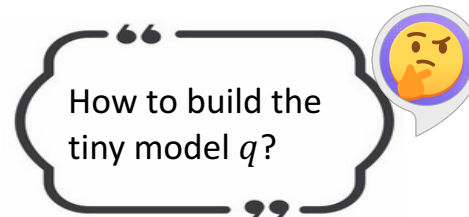


1.  $r \sim U(0,1)$ , if  $r < \min\left(1, \frac{p(t)}{q(t)}\right)$ , next token =  $t$   
**Accept rate**

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

## Theorem:

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

# How to build the tiny model $q$ ?

- Trade-off between accuracy and efficiency



Fast but inaccurate

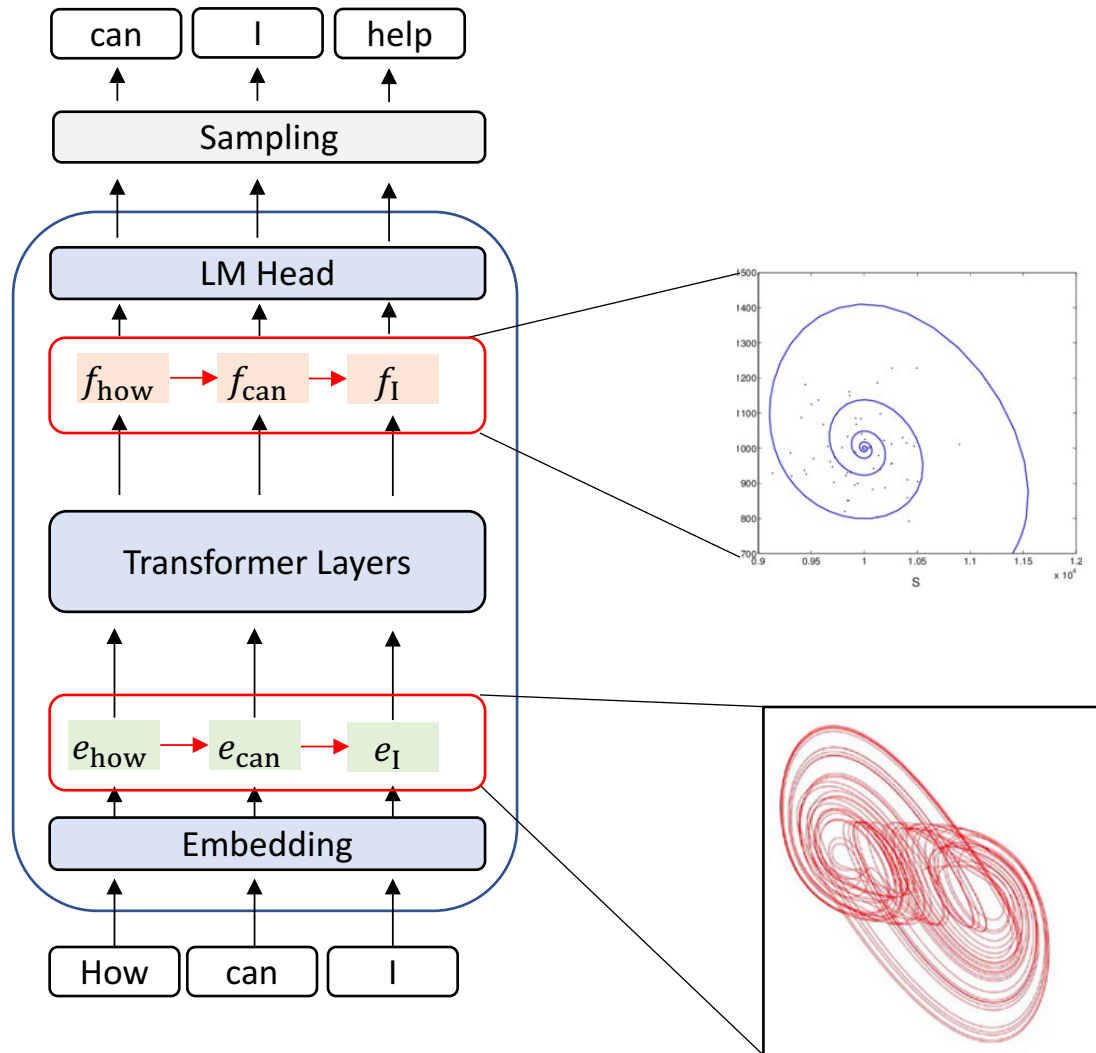


Accurate but slow

# EAGLE: Speculative Sampling Requires Rethinking Feature Uncertainty

Yuhui Li, Fangyun Wei, Chao Zhang, Hongyang Zhang  
(ICML 2024)

# Observations



A **simple** dynamic  
(Need a **tiny** model)

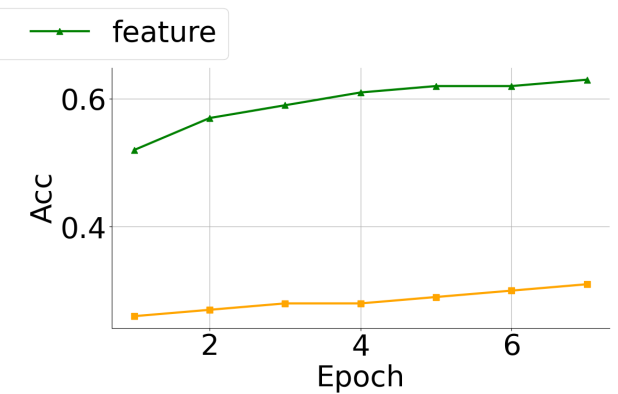
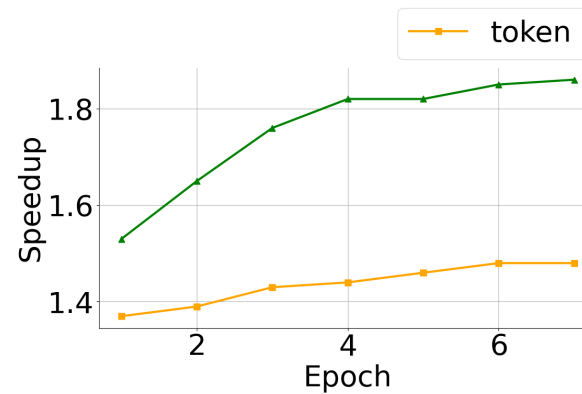
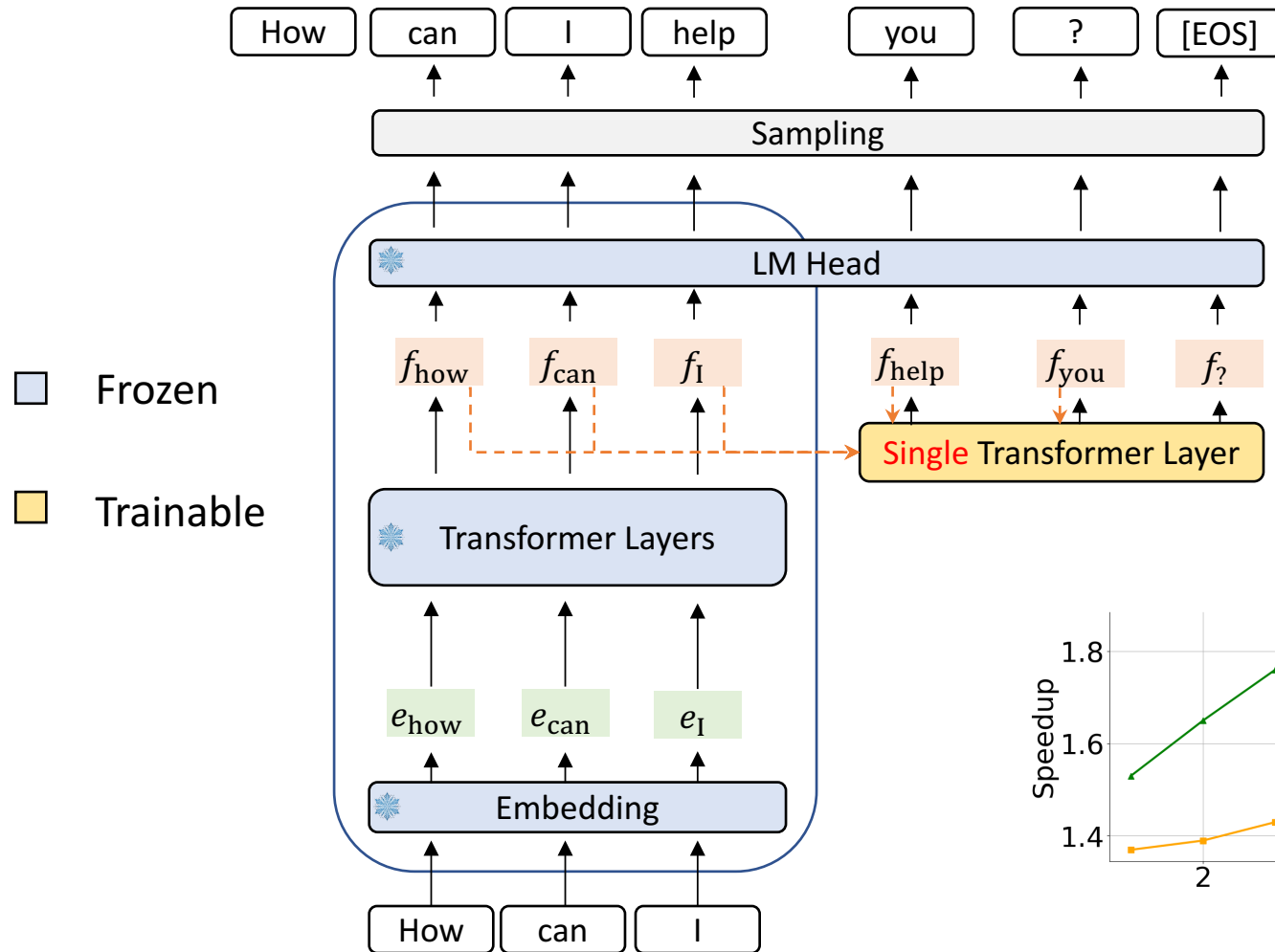


A **complicated** dynamic  
(Need a **large** model)

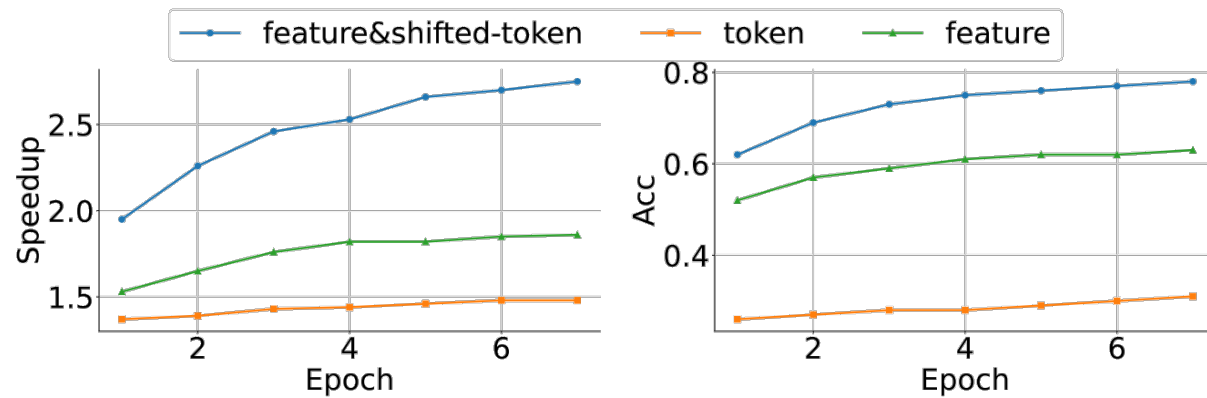
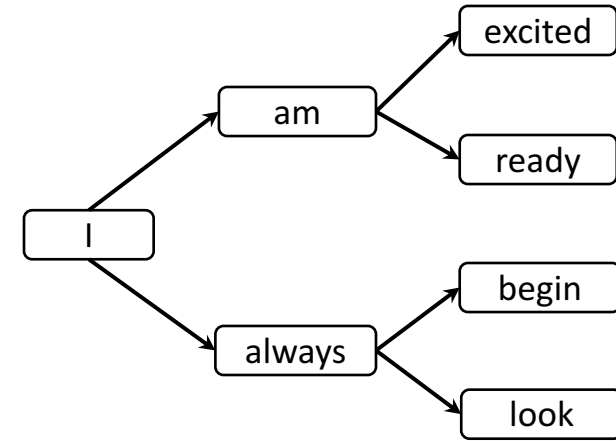
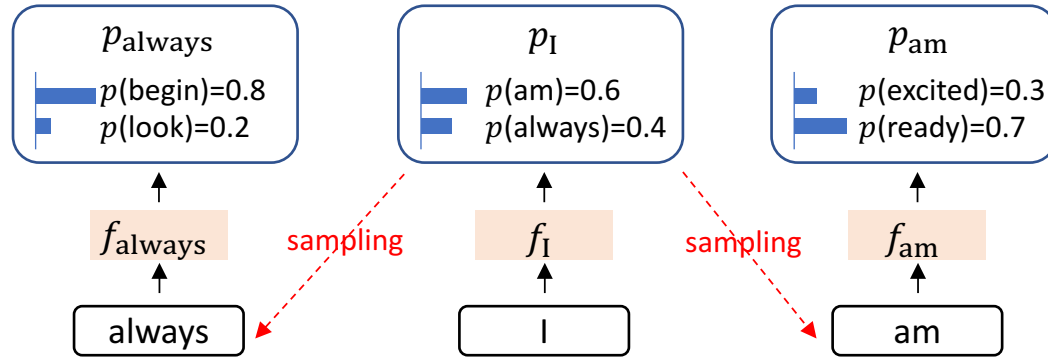




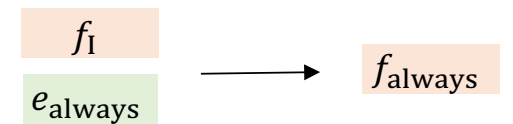
# Our first trial: next-feature prediction



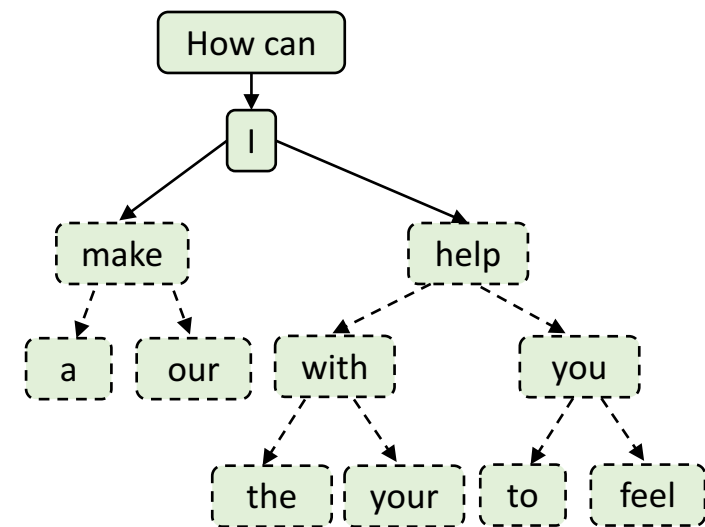
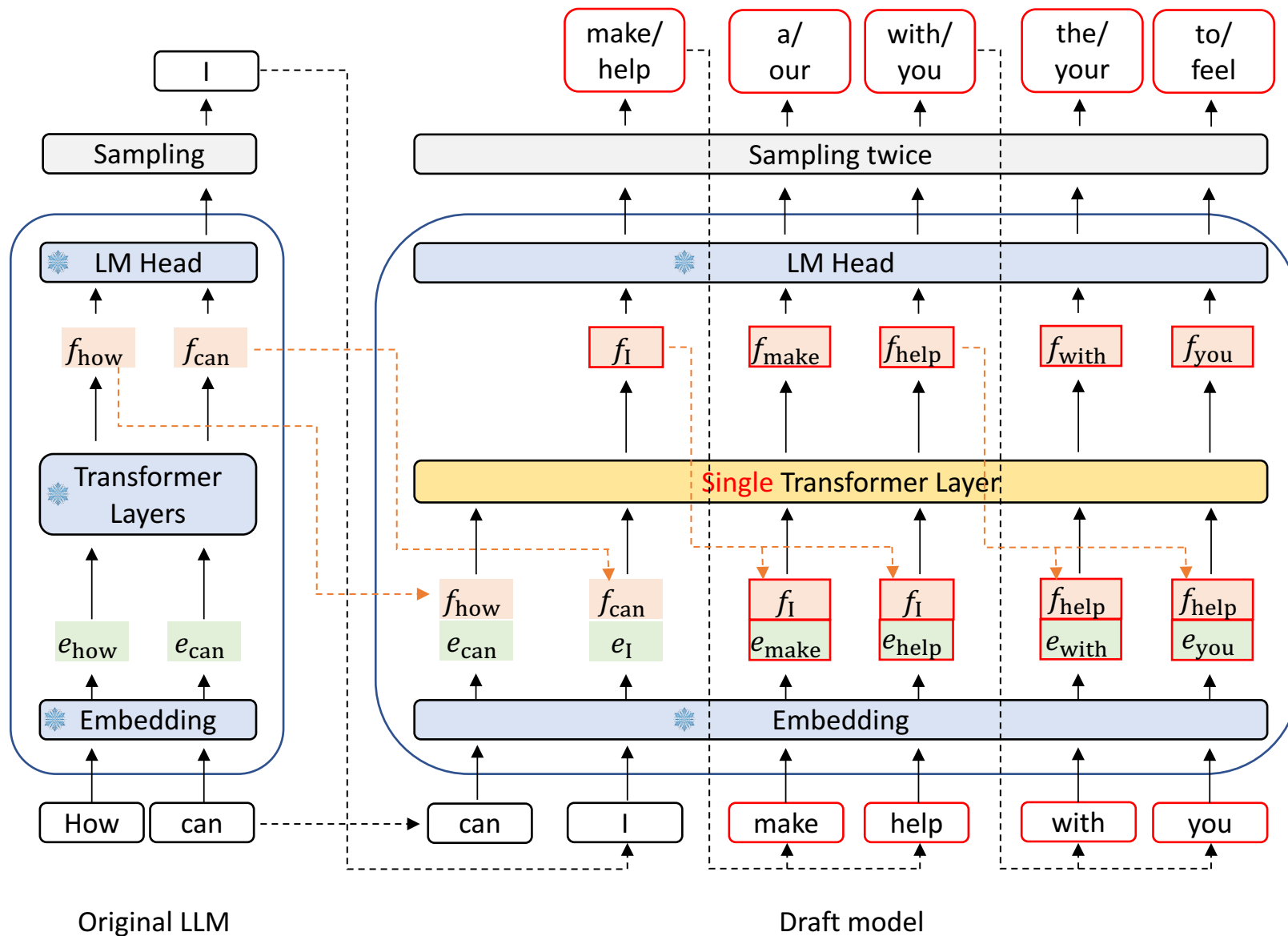
# Feature uncertainty matters



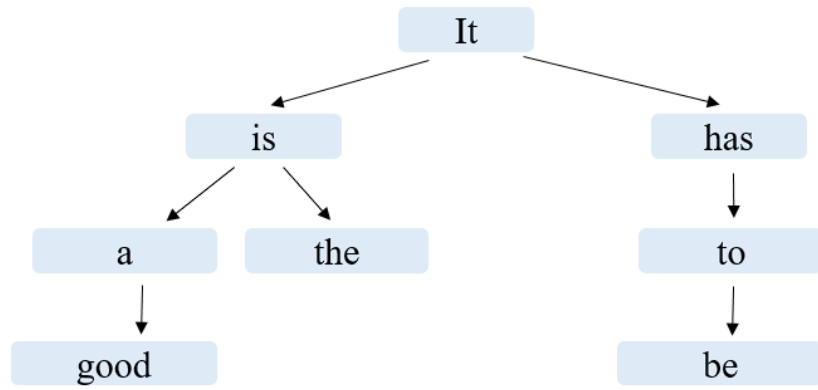
**Idea:** feature & shifted-token  $\rightarrow$  next feature



# Our second trial: EAGLE



# Tree attention



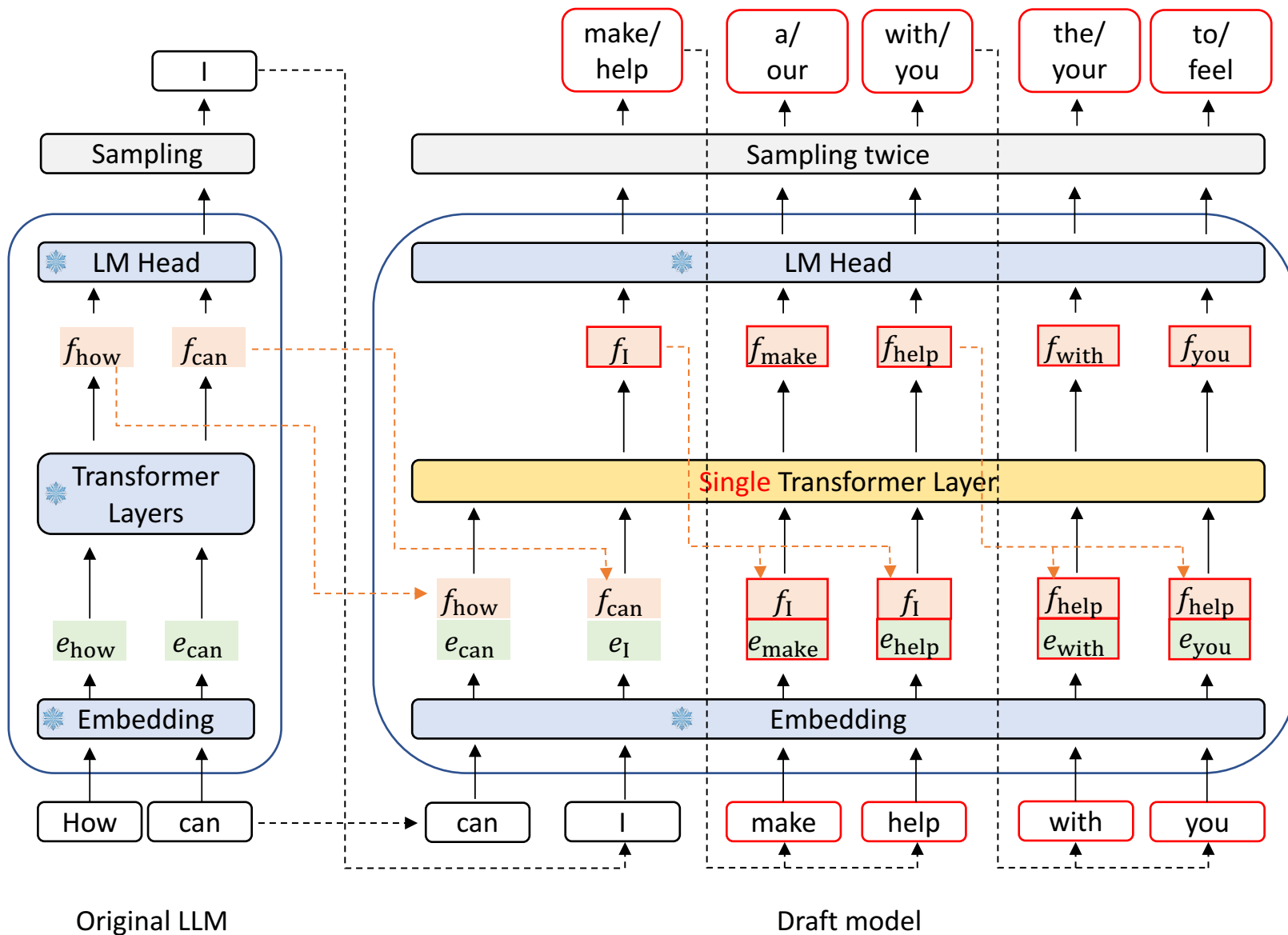
Flatten to 1D

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	✓		✓			✓		✓

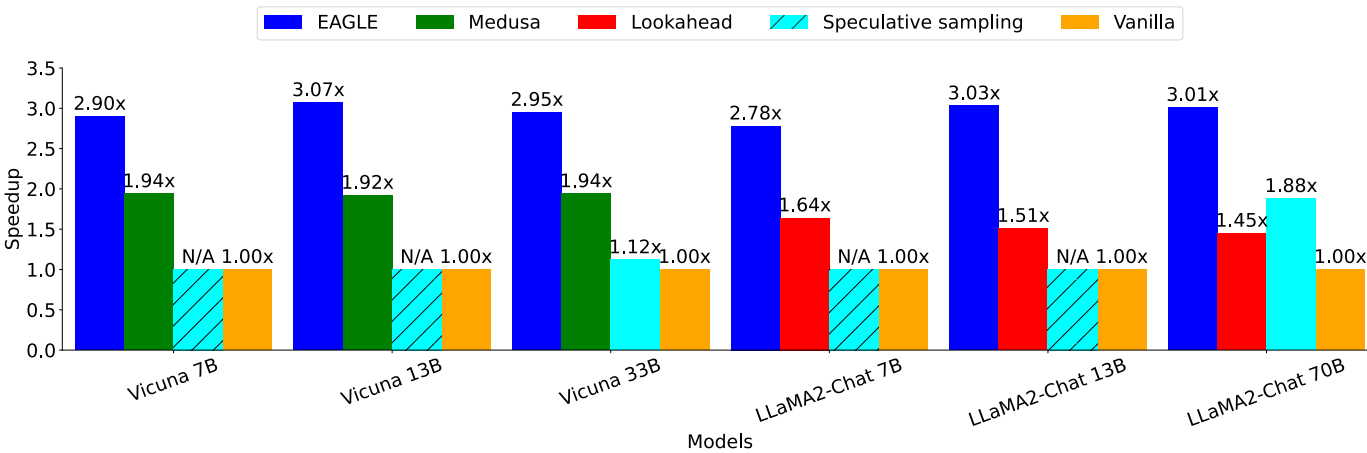
# #Parameters of the drafted models



#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

Vanilla: Speed 3.46 tokens/s, Completion Rate 1.00

Medusa: Speed 3.37 tokens/s, Completion Rate 1.00

Lookahead: Speed 6.09 tokens/s, Completion Rate 1.00

EAGLE: Speed 3.32 tokens/s, Completion Rate 1.00

# Third-party evaluations

## Spec-bench

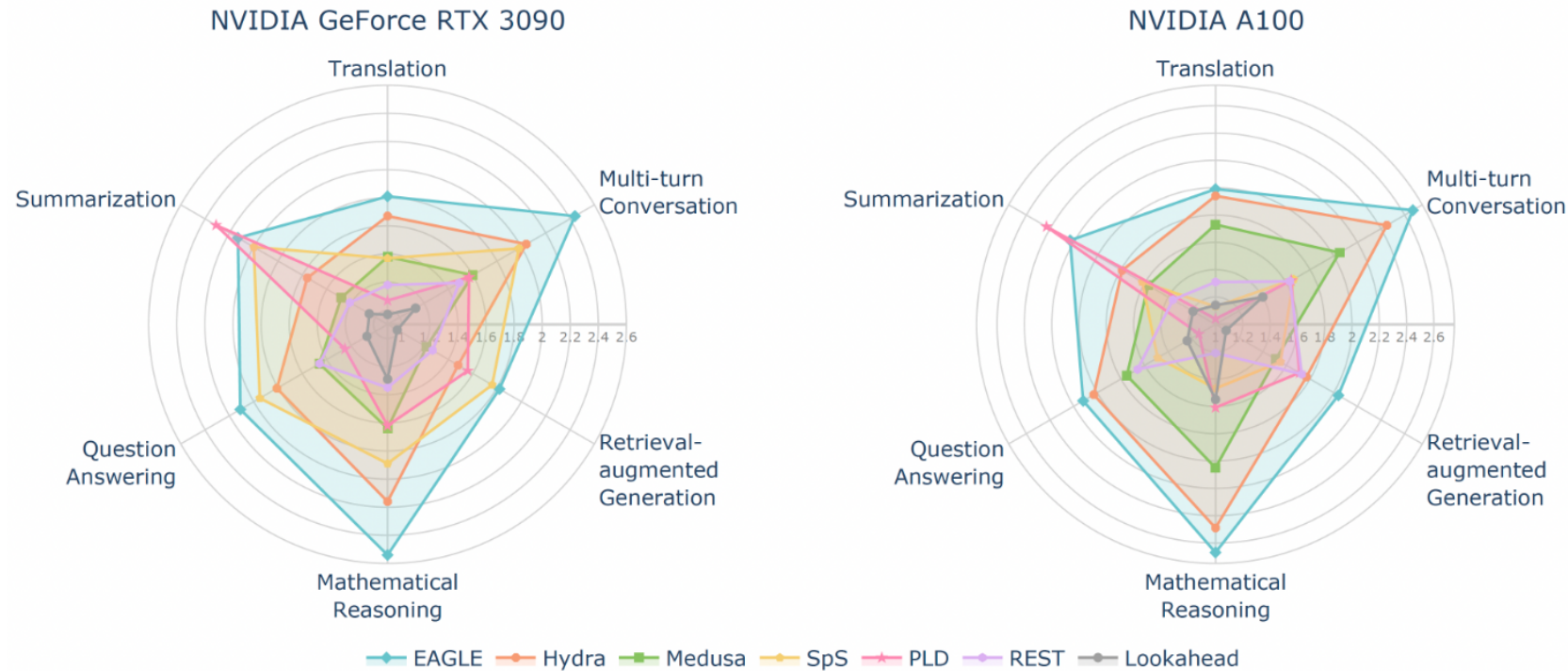
### *Unlocking Efficiency in Large Language Model Inference: A Comprehensive Survey of Speculative Decoding*

Heming Xia<sup>1</sup>, Zhe Yang<sup>2</sup>, Qingxiu Dong<sup>2</sup>, Peiyi Wang<sup>2</sup>,  
Yongqi Li<sup>1</sup>, Tao Ge<sup>3</sup>, Tianyu Liu<sup>4</sup>, Wenjie Li<sup>1</sup>, Zhifang Sui<sup>2</sup>

<sup>1</sup>Department of Computing, The Hong Kong Polytechnic University


<sup>2</sup>National Key Laboratory for Multimedia Information Processing, Peking University

<sup>3</sup>Microsoft Research Asia <sup>4</sup>Alibaba Group



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

# Third-party evaluations

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

Heming Xia<sup>1</sup>, Zhe Yang<sup>2</sup>, Qingxiu Dong<sup>2</sup>, Peiyi Wang<sup>2</sup>,  
Yongqi Li<sup>1</sup>, Tao Ge<sup>3</sup>, Tianyu Liu<sup>4</sup>, Wenjie Li<sup>1</sup>, Zhifang Sui<sup>2</sup>




<sup>1</sup>Department of Computing, The Hong Kong Polytechnic University

<sup>2</sup>National Key Laboratory for Multimedia Information Processing, Peking University

<sup>3</sup>Microsoft Research Asia <sup>4</sup>Alibaba Group

## Spec-bench


RTX 3090, Vicuna 7B

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



# Third-party evaluations

## Spec-bench

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







Heming Xia<sup>1</sup>, Zhe Yang<sup>2</sup>, Qingxiu Dong<sup>2</sup>, Peiyi Wang<sup>2</sup>,  
Yongqi Li<sup>1</sup>, Tao Ge<sup>3</sup>, Tianyu Liu<sup>4</sup>, Wenjie Li<sup>1</sup>, Zhifang Sui<sup>2</sup>

<sup>1</sup>Department of Computing, The Hong Kong Polytechnic University

<sup>2</sup>National Key Laboratory for Multimedia Information Processing, Peking University

<sup>3</sup>Microsoft Research Asia <sup>4</sup>Alibaba Group

A100, Vicuna 7B

Models	Models	Multi-turn Conversation	Translation	Summarization	Question Answering	Mathematical Reasoning	Retrieval-aug. Generation	#Mean Accepted Tokens	Overall
<a href="#">EAGLE</a> 	<a href="#">EAGLE</a> 	2.67x	1.99x	2.23x	2.12x	2.67x	2.04x	3.61	2.29x
<a href="#">SpS</a> 	<a href="#">Hydra</a> 	2.45x	1.94x	1.79x	2.03x	2.49x	1.77x	3.24	2.09x
<a href="#">Hydra</a> 	<a href="#">Medusa</a> 	2.05x	1.73x	1.57x	1.75x	2.05x	1.51x	2.32	1.78x
<a href="#">PLD</a>	<a href="#">PLD</a>	1.64x	1.04x	2.43x	1.14x	1.61x	1.71x	1.73	1.59x
<a href="#">Medusa</a>	<a href="#">SpS</a>	1.66x	1.13x	1.62x	1.49x	1.47x	1.55x	2.28	1.49x
<a href="#">REST</a>	<a href="#">REST</a>	1.63x	1.31x	1.36x	1.66x	1.21x	1.73x	1.82	1.48x
<a href="#">Lookahead</a>	<a href="#">Lookahead</a>	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 Xia<sup>1</sup>, Zhe Yang<sup>2</sup>, Qingxiu Dong<sup>2</sup>, Peiyi Wang<sup>2</sup>,  
Yongqi Li<sup>1</sup>, Tao Ge<sup>3</sup>, Tianyu Liu<sup>4</sup>, Wenjie Li<sup>1</sup>, Zhifang Sui<sup>2</sup>

<sup>1</sup>Department of Computing, The Hong Kong Polytechnic University

<sup>2</sup>National Key Laboratory for Multimedia Information Processing, Peking University


<sup>3</sup>Microsoft Research Asia <sup>4</sup>Alibaba Group

### A100, Vicuna 13B

Models	Models	Models	Multi-turn Conversation	Translation	Summarization	Question Answering	Mathematical Reasoning	Retrieval-aug. Generation	#Mean Accepted Tokens	Overall
<a href="#">EAGLE</a> 🌟	<a href="#">EAGLE</a> 🌟	<a href="#">EAGLE</a> 🌟	2.68x	1.96x	2.44x	2.04x	2.70x	2.23x	3.64	2.34x
<a href="#">SpS</a> ②	<a href="#">Hydra</a> ②	<a href="#">Hydra</a> ②	2.46x	1.90x	1.93x	1.96x	2.48x	1.92x	3.35	2.12x
<a href="#">Hydra</a> ③	<a href="#">Medusa</a> ③	<a href="#">Medusa</a> ③	1.96x	1.66x	1.63x	1.63x	2.00x	1.58x	2.39	1.75x
<a href="#">PLD</a>	<a href="#">PLD</a>	<a href="#">SpS</a>	1.60x	1.13x	1.68x	1.39x	1.53x	1.67x	2.18	1.49x
<a href="#">Medusa</a>	<a href="#">SpS</a>	<a href="#">PLD</a>	1.47x	1.02x	2.19x	1.03x	1.57x	1.71x	1.68	1.48x
<a href="#">REST</a>	<a href="#">REST</a>	<a href="#">REST</a>	1.52x	1.17x	1.37x	1.53x	1.19x	1.55x	1.82	1.38x
<a href="#">Lookahead</a>	<a href="#">Lookahead</a>	<a href="#">Lookahead</a>	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 Xia<sup>1</sup>, Zhe Yang<sup>2</sup>, Qingxiu Dong<sup>2</sup>, Peiyi Wang<sup>2</sup>,  
Yongqi Li<sup>1</sup>, Tao Ge<sup>3</sup>, Tianyu Liu<sup>4</sup>, Wenjie Li<sup>1</sup>, Zhifang Sui<sup>2</sup>

<sup>1</sup>Department of Computing, The Hong Kong Polytechnic University

<sup>2</sup>National Key Laboratory for Multimedia Information Processing, Peking University

<sup>3</sup>Microsoft Research Asia <sup>4</sup>Alibaba Group

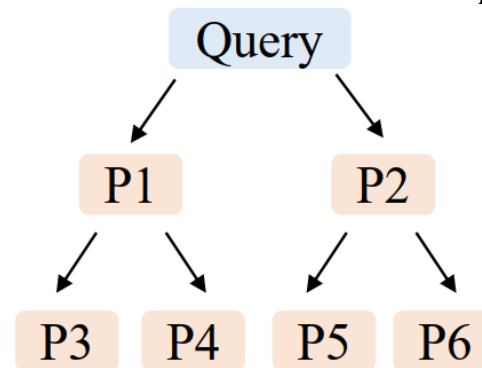
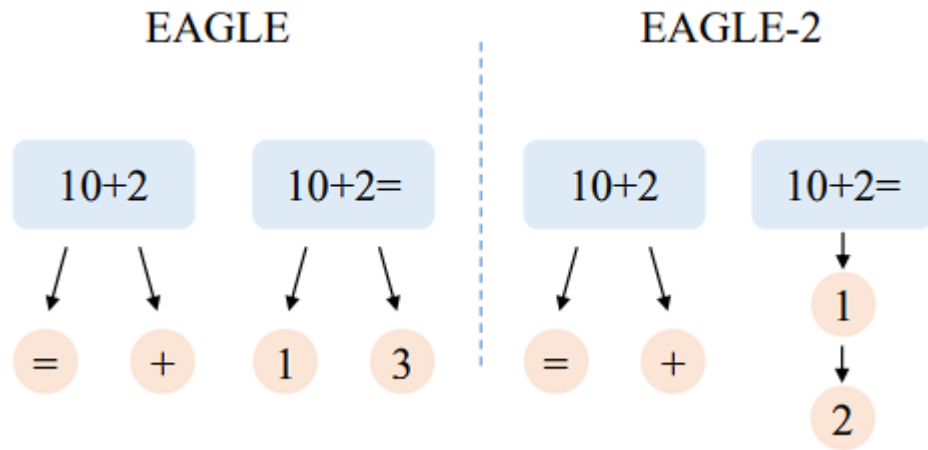
### A100, Vicuna 33B

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

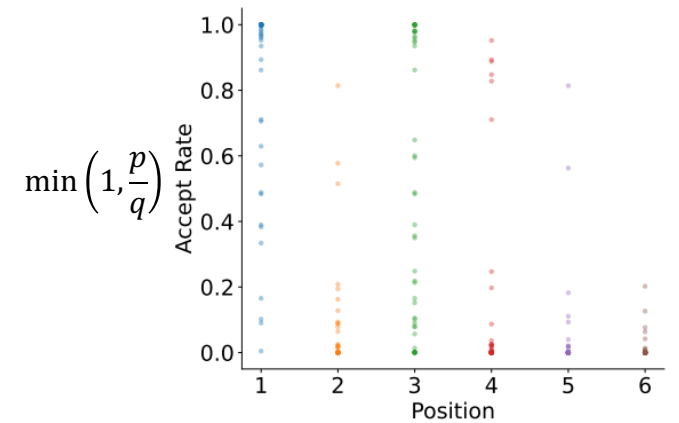
# EAGLE-2: Faster Inference of Language Models with Dynamic Draft Trees

Yuhui Li, Fangyun Wei, Chao Zhang, Hongyang Zhang  
(EMNLP 2024)

# Static tree structure?



(a) Draft tree structure.

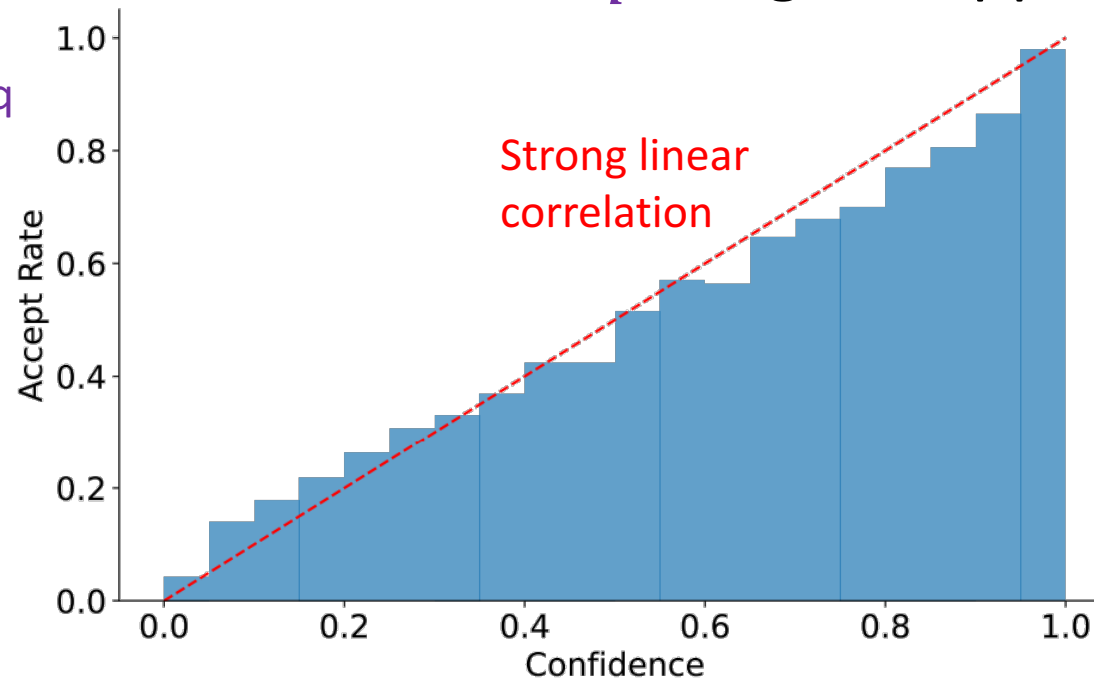


(b) Acceptance rates of tokens at different positions, with each point representing a query.

# How to determine importance of each node?

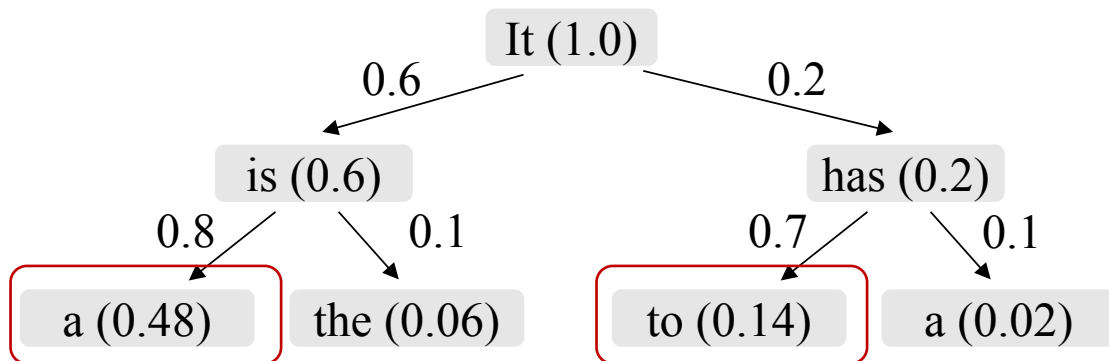
- Accept rate  $\min\left(1, \frac{p(t)}{q(t)}\right)$
- However, it requires the computation from the original **large** model  $p$
- The **confidence of the draft model  $q$**  is a good approximation

i.e. output  
probability of  $q$

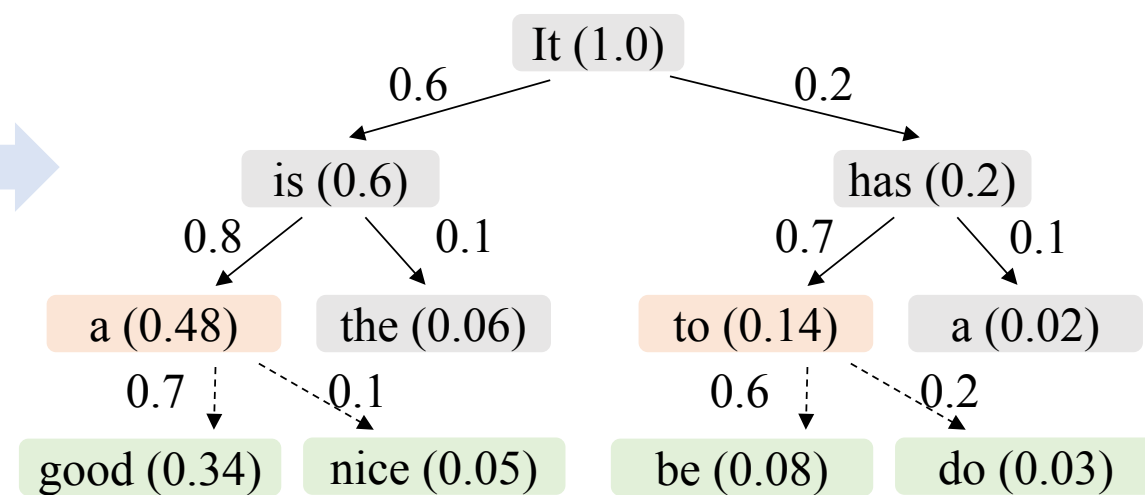


# Context-aware dynamic draft tree

Beam Search (Top-2)

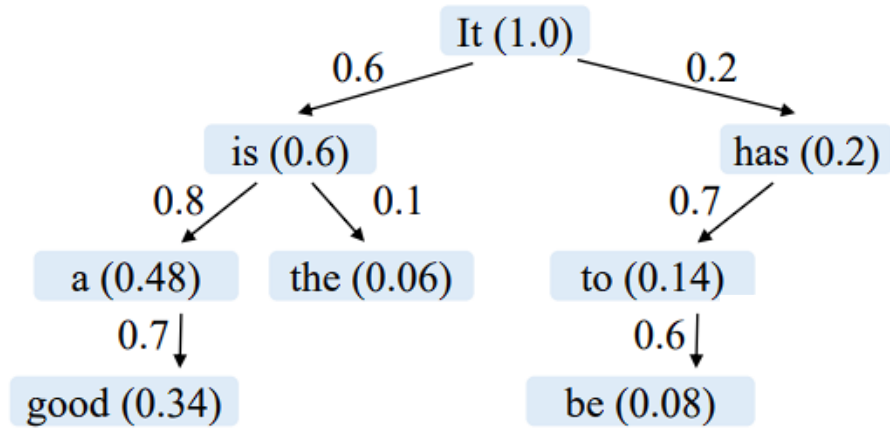


Expand



# Context-aware dynamic draft tree

Rerank (Top-8)



Flatten to 1D

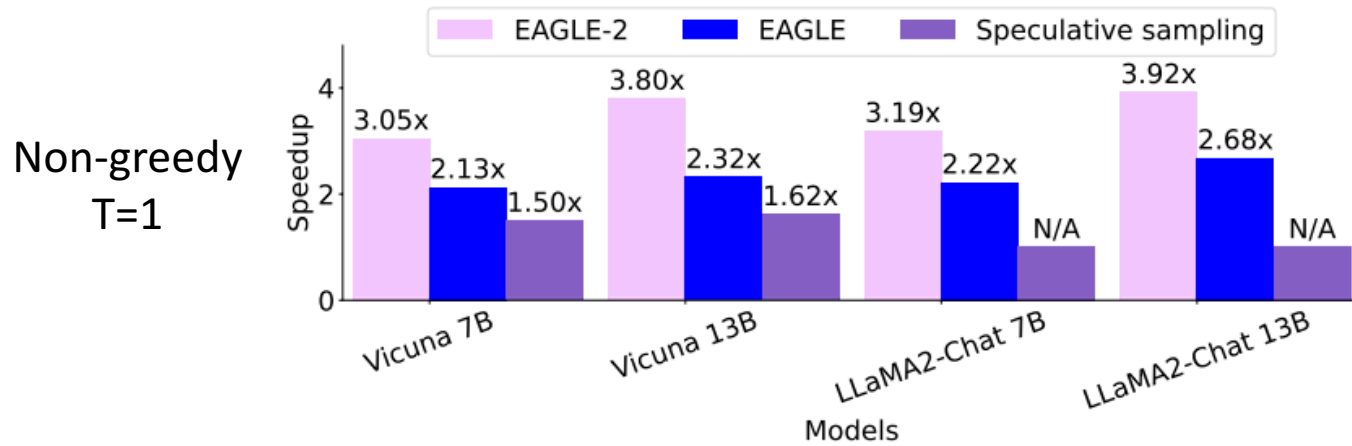
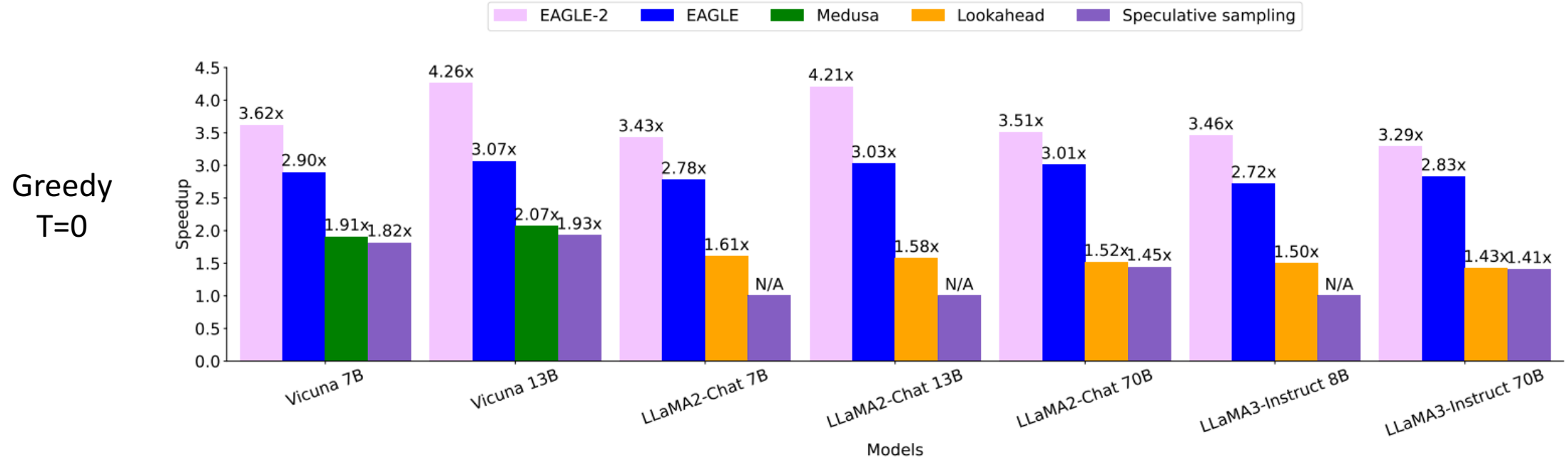
It	is	has	a	the	to	good	be
----	----	-----	---	-----	----	------	----

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	Speedup	#mean accepted tokens	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	<b>4.26x</b>	<b>4.83</b>	<b>4.96x</b>	<b>5.41</b>	<b>4.22x</b>	<b>4.79</b>	<b>4.25x</b>	<b>4.89</b>	<b>3.40x</b>	<b>4.21</b>	<b>3.13x</b>	<b>3.74</b>	<b>4.04x</b>	<b>4.65</b>
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	<b>4.21x</b>	<b>4.75</b>	<b>5.00x</b>	<b>5.52</b>	<b>4.31x</b>	<b>4.90</b>	<b>4.13x</b>	<b>4.61</b>	<b>3.45x</b>	<b>4.24</b>	<b>3.51x</b>	<b>4.04</b>	<b>4.10x</b>	<b>4.68</b>
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	<b>3.62x</b>	<b>4.98</b>	<b>3.95x</b>	<b>5.33</b>	<b>3.63x</b>	<b>4.97</b>	<b>3.46x</b>	<b>4.86</b>	<b>2.94x</b>	<b>4.12</b>	<b>2.76x</b>	<b>3.82</b>	<b>3.39x</b>	<b>4.68</b>
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	<b>3.43x</b>	<b>4.70</b>	<b>4.03x</b>	<b>5.39</b>	<b>3.52x</b>	<b>4.77</b>	<b>3.45x</b>	<b>4.66</b>	<b>3.01x</b>	<b>4.12</b>	<b>3.15x</b>	<b>4.19</b>	<b>3.43x</b>	<b>4.64</b>

# 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	Speedup	#mean accepted tokens	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	<b>3.80x</b>	<b>4.40</b>	<b>4.22x</b>	<b>4.89</b>	<b>3.77x</b>	<b>4.41</b>	<b>3.78x</b>	<b>4.37</b>	<b>3.25x</b>	<b>3.97</b>	<b>3.07x</b>	<b>3.54</b>	<b>3.65x</b>	<b>4.26</b>
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	<b>3.92x</b>	<b>4.51</b>	<b>4.58x</b>	<b>5.29</b>	<b>4.21x</b>	<b>4.80</b>	<b>3.85x</b>	<b>4.48</b>	<b>3.31x</b>	<b>4.08</b>	<b>3.43x</b>	<b>3.89</b>	<b>3.88x</b>	<b>4.51</b>
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	<b>3.05x</b>	<b>4.28</b>	<b>3.33x</b>	<b>4.65</b>	<b>3.07x</b>	<b>4.49</b>	<b>3.08x</b>	<b>4.43</b>	<b>2.63x</b>	<b>3.76</b>	<b>2.48x</b>	<b>3.56</b>	<b>2.94x</b>	<b>4.20</b>
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	<b>3.19x</b>	<b>4.41</b>	<b>3.67x</b>	<b>5.06</b>	<b>3.35x</b>	<b>4.62</b>	<b>3.20x</b>	<b>4.48</b>	<b>2.73x</b>	<b>3.85</b>	<b>2.81x</b>	<b>4.01</b>	<b>3.15x</b>	<b>4.41</b>

# 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 Xia<sup>1</sup>, Zhe Yang<sup>2</sup>, Qingxiu Dong<sup>2</sup>, Peiyi Wang<sup>2</sup>,  
Yongqi Li<sup>1</sup>, Tao Ge<sup>3</sup>, Tianyu Liu<sup>4</sup>, Wenjie Li<sup>1</sup>, Zhifang Sui<sup>2</sup>

<sup>1</sup>Department of Computing, The Hong Kong Polytechnic University

<sup>2</sup>National Key Laboratory for Multimedia Information Processing, Peking University

<sup>3</sup>Microsoft Research Asia <sup>4</sup>Alibaba Group

## Leaderboard on 3090

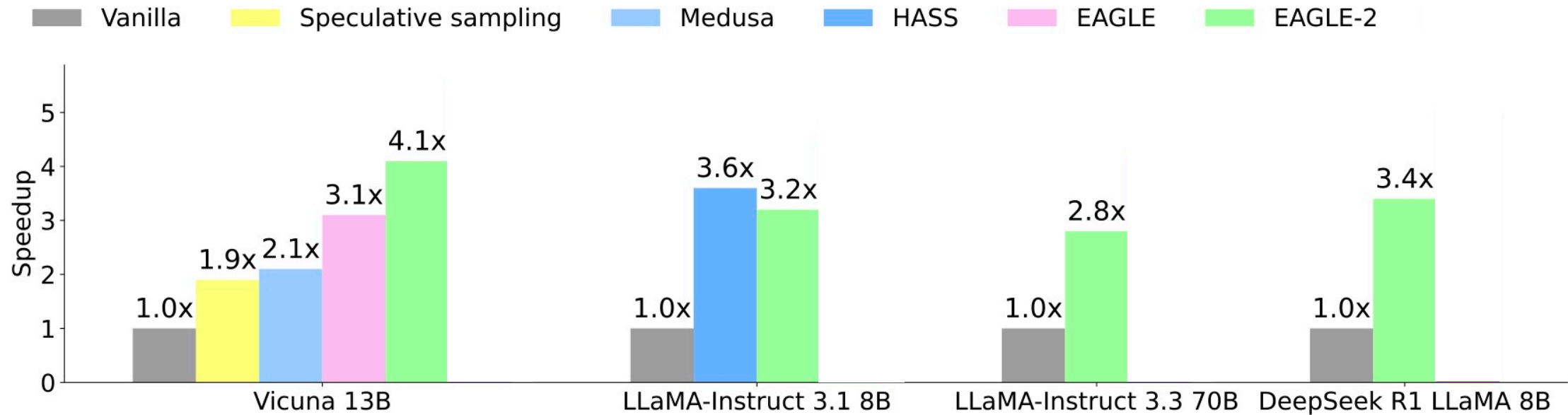
- Device: a single NVIDIA GeForce RTX 3090 GPU (24GB) with 12 CPU cores
- Testing environment: Pytorch 2.0.1, under CUDA 11.8
- Experimental Settings: Vicuna-7B-v1.3, greedy decoding, FP16 precision, batch size = 1

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

# EAGLE-3

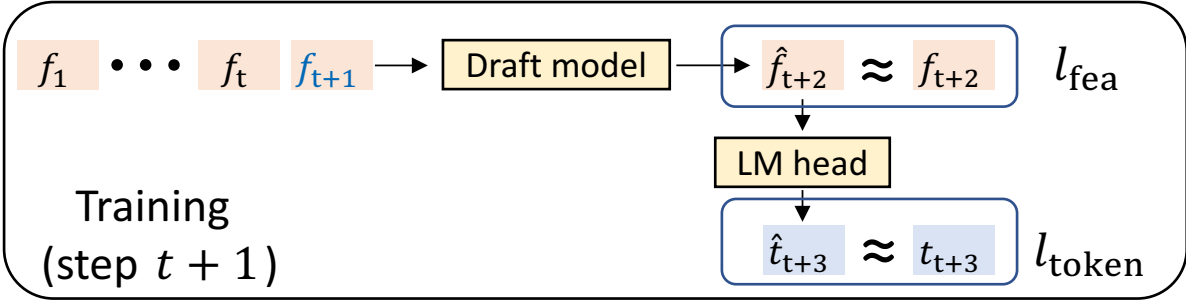
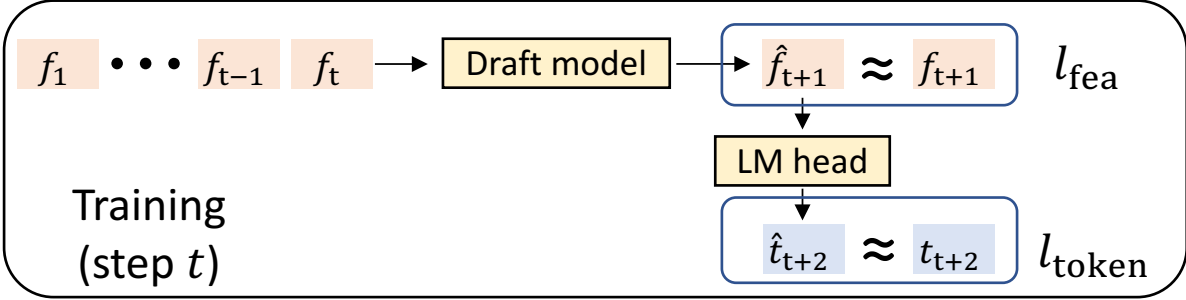
Yuhui Li, Fangyun Wei, Chao Zhang, Hongyang Zhang  
(available on arXiv yesterday)

# Benchmarking on MT-Bench

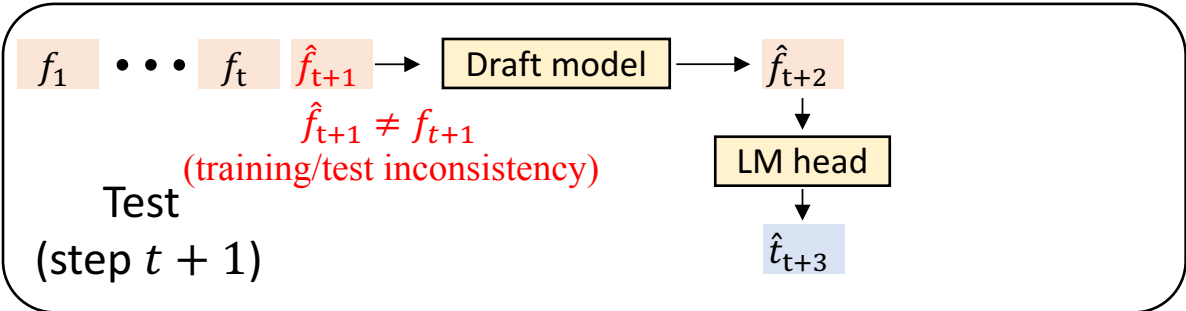
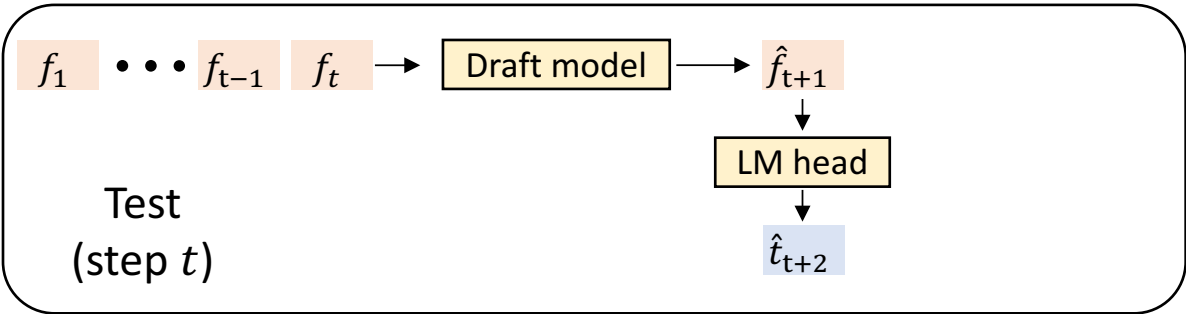


# EAGLE vs EAGLE-3

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



EAGLE-1 training

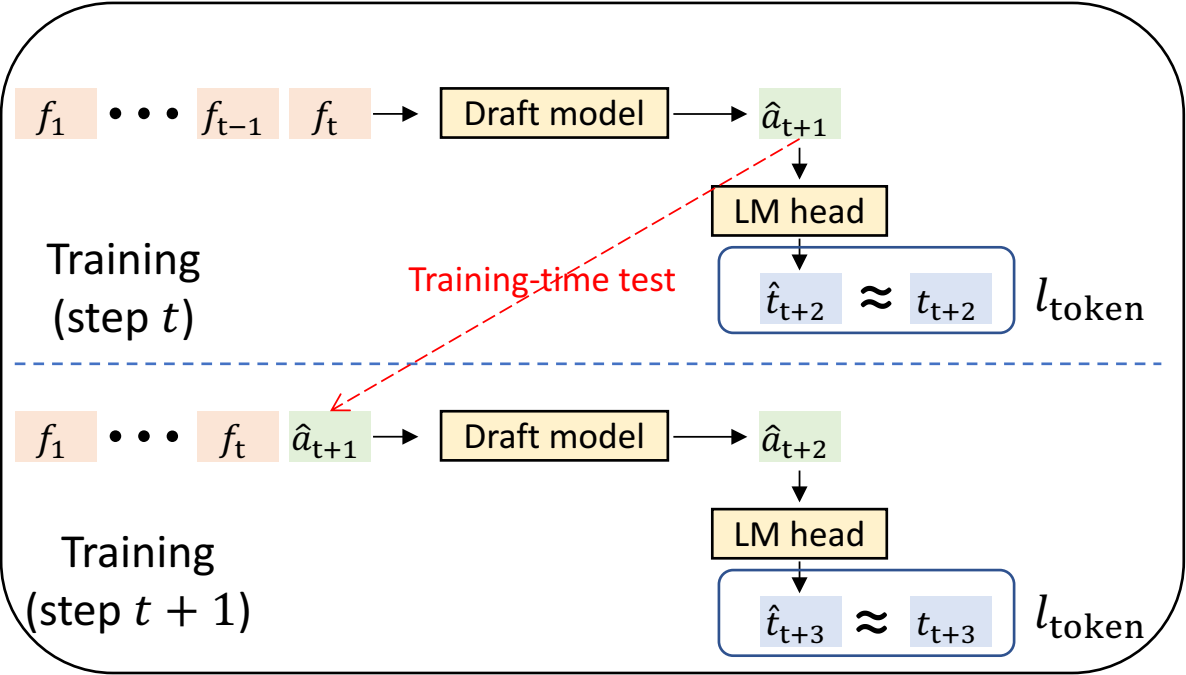


EAGLE-1 test

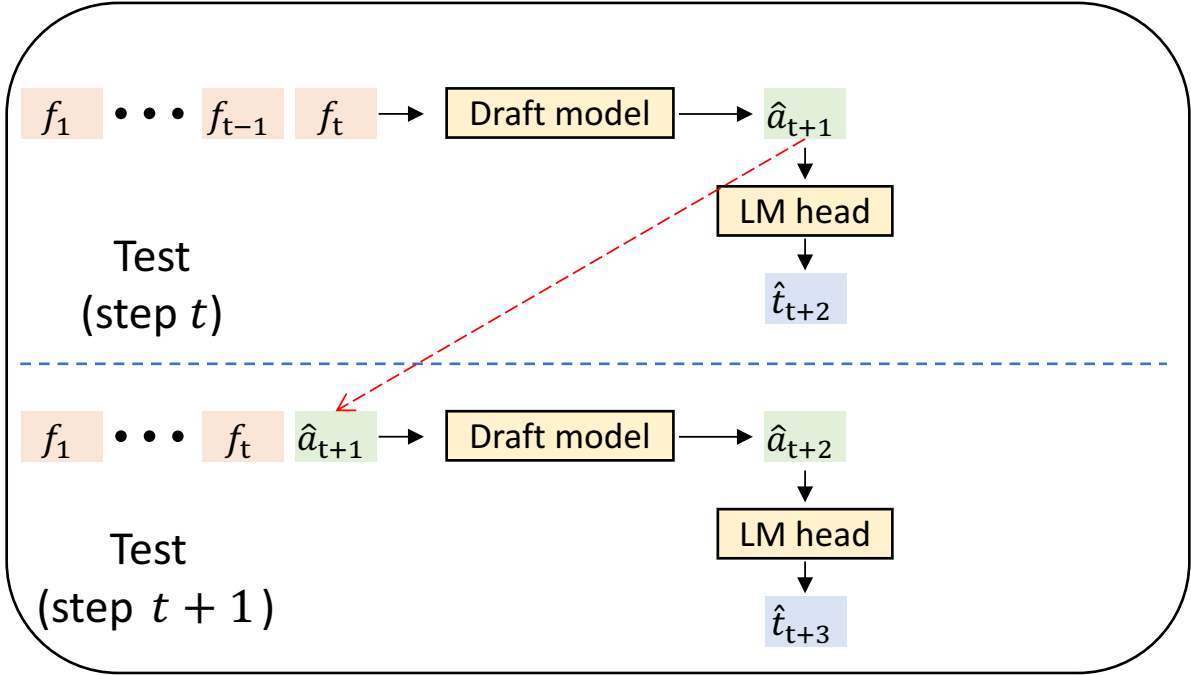


# EAGLE vs EAGLE-3

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

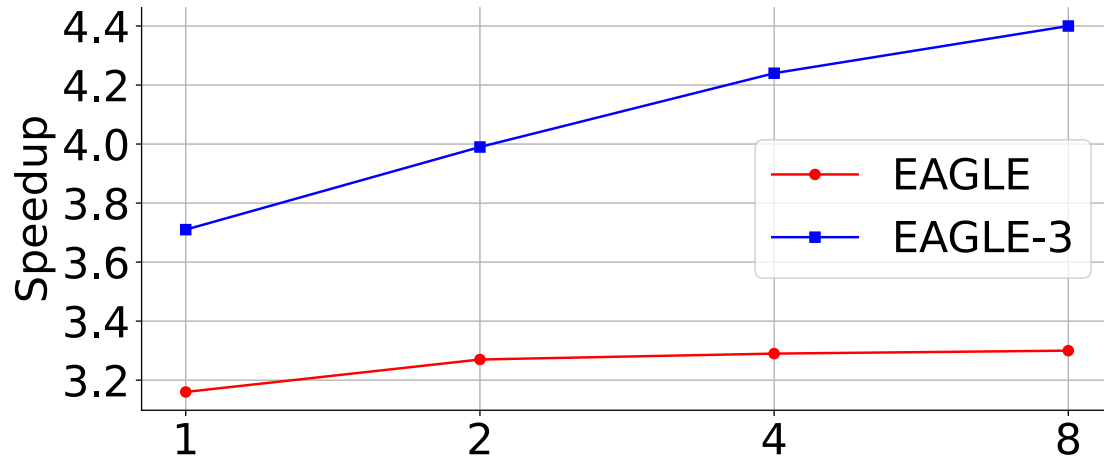


EAGLE-3 training

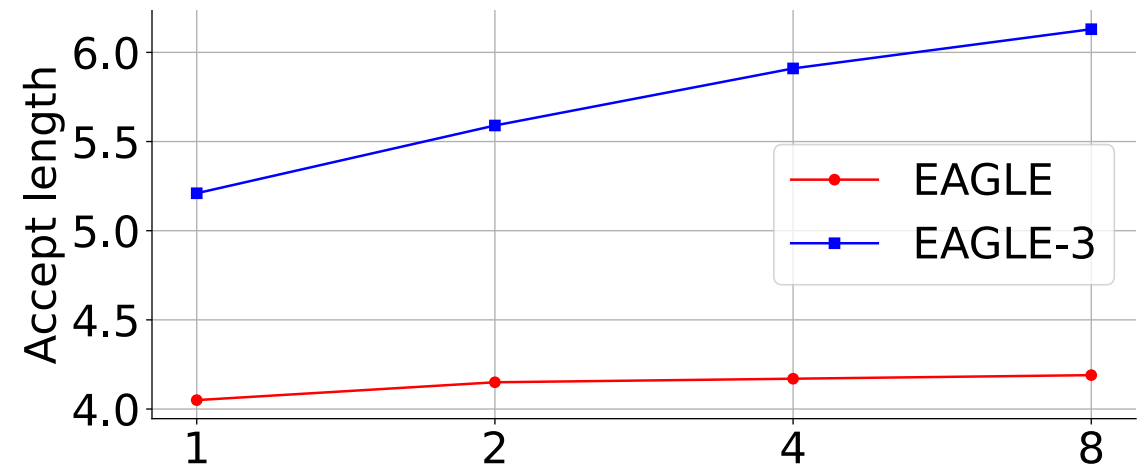


EAGLE-3 test

# 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)

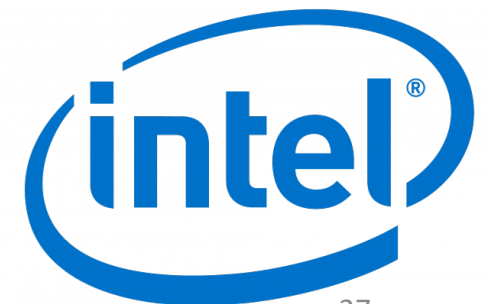
Model	Method	MT-bench		HumanEval		GSM8K		Alpaca		CNN/DM		Mean	
		Speedup	$\tau$	Speedup	$\tau$	Speedup	$\tau$	Speedup	$\tau$	Speedup	$\tau$	Speedup	$\tau$
Temperature=0													
V 13B	SpS	1.93x	2.27	2.23x	2.57	1.77x	2.01	1.76x	2.03	1.93x	2.33	1.92x	2.24
	PLD	1.58x	1.63	1.85x	1.93	1.68x	1.73	1.16x	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	<b>5.58x</b>	<b>6.65</b>	<b>6.47x</b>	<b>7.54</b>	<b>5.32x</b>	<b>6.29</b>	<b>5.16x</b>	<b>6.17</b>	<b>5.01x</b>	<b>6.47</b>	<b>5.51x</b>	<b>6.62</b>
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	<b>4.40x</b>	<b>6.13</b>	<b>4.85x</b>	<b>6.74</b>	<b>4.48x</b>	<b>6.23</b>	<b>4.82x</b>	<b>6.70</b>	<b>3.65x</b>	<b>5.34</b>	<b>4.44x</b>	<b>6.23</b>
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	<b>4.11x</b>	<b>5.63</b>	<b>4.79x</b>	<b>6.52</b>	<b>4.34x</b>	<b>6.15</b>	<b>4.30x</b>	<b>6.09</b>	<b>3.27x</b>	<b>5.02</b>	<b>4.12x</b>	<b>5.88</b>
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	<b>4.05x</b>	<b>5.58</b>	<b>4.59x</b>	<b>6.38</b>	<b>5.01x</b>	<b>6.93</b>	<b>3.65x</b>	<b>5.37</b>	<b>3.52x</b>	<b>4.92</b>	<b>4.16x</b>	<b>5.84</b>
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	<b>4.57x</b>	<b>5.42</b>	<b>5.15x</b>	<b>6.22</b>	<b>4.71x</b>	<b>5.58</b>	<b>4.49x</b>	<b>5.39</b>	<b>4.33x</b>	<b>5.72</b>	<b>4.65x</b>	<b>5.67</b>
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	<b>3.07x</b>	<b>4.24</b>	<b>4.13x</b>	<b>5.82</b>	<b>3.32x</b>	<b>4.59</b>	<b>3.90x</b>	<b>5.56</b>	<b>2.99x</b>	<b>4.39</b>	<b>3.45x</b>	<b>4.92</b>
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	<b>3.96x</b>	<b>5.45</b>	<b>4.36x</b>	<b>6.16</b>	<b>4.17x</b>	<b>5.95</b>	<b>4.14x</b>	<b>5.87</b>	<b>3.11x</b>	<b>4.88</b>	<b>3.95x</b>	<b>5.66</b>
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	<b>3.20x</b>	<b>4.49</b>	<b>3.77x</b>	<b>5.28</b>	<b>4.38x</b>	<b>6.10</b>	<b>3.16x</b>	<b>4.30</b>	<b>3.08x</b>	<b>4.27</b>	<b>3.52x</b>	<b>4.89</b>

# Throughput (compared to vLLM w/o EAGLE)

Batch size	2	4	8	16	24	32	48	56
EAGLE	1.30x	1.25x	1.21x	1.10x	1.03x	0.93x	0.82x	0.71x
EAGLE-3	1.75x	1.68x	1.58x	1.49x	1.42x	1.36x	1.21x	1.01x

# 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?

Thanks!  
Q&A

Code

