Mamba Drafters for Speculative Decoding







Daewon Choi¹ Seunghyuk Oh¹ Saket Dingliwal² Jihoon Tack¹ Kyuyoung Kim¹ Woomin Song^{1, 2} Seojin Kim³ Insu Han¹ Jinwoo Shin¹ Aram Galstyan² Shubham Katiyar² Sravan Babu Bodapati²

¹KAIST ²Amazon AGI ³Seoul National University

Contact: daeone0920@kaist.ac.kr

TL;DR: Fast, memory-efficient, even effective-Mamba is the ideal drafter for speculative decoding.

Why Mamba for Speculative Decoding ?

1. Efficiency of Mamba as a drafter



Experiments

Mamba drafter outperform Transformer drafters of all sizes, and are even comparable to state-of-the-art self-speculation method !

	Drafte	er	C	Greedy (Temp=0)			Sampling (Temp=1)			
Target	Method	Size	XSum	CNN-DM	GSM-8k	XSum	CNN-DM	GSM-8k		
	No drafter	_	53.30	49.29	54.69	52.51	45.33	53.81		
Pythia-6.9B		70M	47.31 (1.52)	46.99 (1.54)	57.36 (1.68)	41.86 (1.67)	45.30 (1.76)	47.96 (1.77)		
	Pythia	160M	50.05 (2.23)	49.53 (2.26)	67.89 (2.72)	46.67 (2.28)	47.17 (2.30)	55.40 (2.63)		
		410M	70.53 (4.62)	70.08 (4.73)	75.97 (4.64)	53.50 (3.60)	56.64 (3.80)	63.64 (4.01)		
	Ours	130M	138.80 (4.55)	131.97 (4.38)	149.46 (4.57)	108.68 (3.53)	105.01 (3.53)	119.67 (3.73)		
	No drafter	_	51.15	49.55	50.31	53.49	47.40	52.92		
Mistral-7B	Mistral	160M	61.55 (3.13)	61.04 (3.05)	49.38 (2.21)	53.91 (2.74)	50.50 (2.68)	62.29 (2.94)		
	Ours	130M	76.71 (2.39)	65.23 (2.13)	77.50 (2.25)	79.18 (2.73)	70.95 (2.65)	82.63 (2.73)		

Input Length	Input Length	Input Length
(a) Encoding Memory	(b) Decoding Time	(c) Decoding Memory
* Mamba / Mistral / EAGLE		

External Transformer drafter Self-speculation with a single-layer Transformer

 Mamba enables fast and memory-efficient decoding regardless of context length. (Unlike Transformer-based drafters)

2. Effectiveness of Mamba as a drafter





- ✓ Mamba drafters achieve higher acceptance length than the Transformer (Better alignment with the distribution of larger Transformers, e.g., lower ECE)
- Smaller Mamba offers significantly faster drafting speed compared to larger sizes.
 (With a slightly lower acceptance length -> achieve higher throughput!)

Tree-Structured Drafting with Mamba

Method1) Efficient tree-structured drafting with batch generation

Pre-trained models on language modeling tasks

Drafter			Gr	Greedy (Temp=0)			Sampling (Temp=1)		
Target	Method	External?	MT-bench	Alpaca	Human-Eval	MT-bench	Alpaca	Human-Eval	
	No drafter	_	54.51	55.28	54.76	53.89	54.72	54.21	
	Pythia	1	70.71 (3.10)	60.77 (2.65)	109.51 (4.68)	65.73 (3.03)	62.07 (2.82)	109.52 (4.25)	
Pythia-6.9B	EAGLE	×	$\frac{125.61}{(3.85)}$	117.17 (3.53)	$\frac{122.44}{(4.71)}$	<u>87.01</u> (2.67)	$\frac{78.58}{(2.40)}$	<u>83.05</u> (2.97)	
	Ours	1	128.21 (3.91)	$\frac{114.08}{(3.41)}$	172.38 (5.41)	110.20 (3.65)	108.54 (3.51)	143.55 (4.82)	
	No drafter	_	52.97	53.58	52.30	52.39	53.02	52.34	
	Mistral	1	67.47 (3.04)	61.40 (2.73)	100.23 (4.53)	57.19 (2.84)	51.05 (2.40)	80.94 (3.92)	
Mistral-7B	EAGLE	×	107.16 (3.22)	$\frac{94.03}{(2.79)}$	132.69 (3.98)	94.03 (2.90)	86.60 (2.63)	122.11 (3.78)	
	Ours	1	$\frac{102.48}{(3.16)}$	96.83 (2.96)	$\frac{118.04}{(3.69)}$	$\frac{88.68}{(2.95)}$	$\frac{82.75}{(2.71)}$	$\frac{87.81}{(2.94)}$	

Instruction-tuned models on instruction following tasks

		S	Single-Document QA			Multi-Document QA			Peak Memory (GB)				
Method	External?	1k	2k	4k	8k	1k	2k	4k	8k	1k	2k	4k	8k
No drafter	-	31.02	27.89	24.35	19.30	28.17	24.22	19.01	14.83	15	16	20	36
Mistral	1	25.30 (2.43)	23.28 (2.37)	19.48 (2.24)	15.23 (2.21)	24.64 (2.53)	21.06 (2.48)	16.49 (2.44)	12.18 (2.39)	31	33	38	59
EAGLE	×	$\frac{53.13}{(2.73)}$	47.00 (2.81)	37.27 (2.76)	26.12 (2.71)	$\frac{42.48}{(2.60)}$	<u>35.38</u> (2.61)	$\frac{25.10}{(2.68)}$	$\frac{17.36}{(2.64)}$	32	34	42	72
Ours	1	55.09 (2.91)	$\frac{45.65}{(2.77)}$	$\frac{36.32}{(2.77)}$	$\frac{24.92}{(2.80)}$	47.91 (2.94)	36.40 (2.86)	26.27 (2.88)	17.77 (2.87)	31	32	36	52

For generating multiple candidates per generation during the drafting step,

- Mamba: copying the current single state to predict the next token.
- Transformers: copying the current sequence length of the KV cache.
- \rightarrow Duplication overheads for batch generation is low in Mamba !
 - Given a tree configuration $~\mathcal{T}~=(N_1,N_2,...,N_\gamma)$,
 - (γ : draft length

 N_i : # of new nodes obtained by sampling from each node at the i^{th} generation)

✓ Batch generation of a total batch size, i.e., $\mathcal{B}_i = N_1 \times N_2 \times \cdots \times N_i$ ✓ Possible input size is deterministic, i.e., $(\mathcal{B}_1, 1), (\mathcal{B}_2, 1), ..., (\mathcal{B}_{\gamma}, 1)$

Pre-allocation of memory / Graph caching is applicable for acceleration 🔶

Method2) Test-time dynamic tree search using multi-armed bandit

Thanks to fast drafting, Mamba can leverage different tree configurations.

→ Formalize the tree search problem as a multi-armed bandit (MAB) at test-time!



Long context scenario

Discussion points

Q1. Can Mamba be used with different target model without re-training?

Method	Setup	Accept length	Throughput
EAGLE	Pythia \rightarrow Pythia	2.59	94.75
	Mistral \rightarrow Pythia	N/A	N/A
Ours	Pythia \rightarrow Pythia	3.08	112.69
	Mistral \rightarrow Pythia	2.45	93.20

Q2. Effects of tree-drafting

Search?	MT-bench	Alpaca	HumanEval	Avg.
×	124.99	116.12	149.15	130.09
	128.21	114.08	172.38	138.22

Mamba achieves throughput comparable to EAGLE, which is specifically trained for the target model.

Q3. Effects of test-time tree search

Tree?	Accept length	Latency	Throughput
×	3.08	6.62	112.69
1	3.91	8.30	127.37

Preliminaries: Speculative Decoding (SD)



SD is to generate candidate tokens with an efficient drafter and verifying them in parallel with the target model.