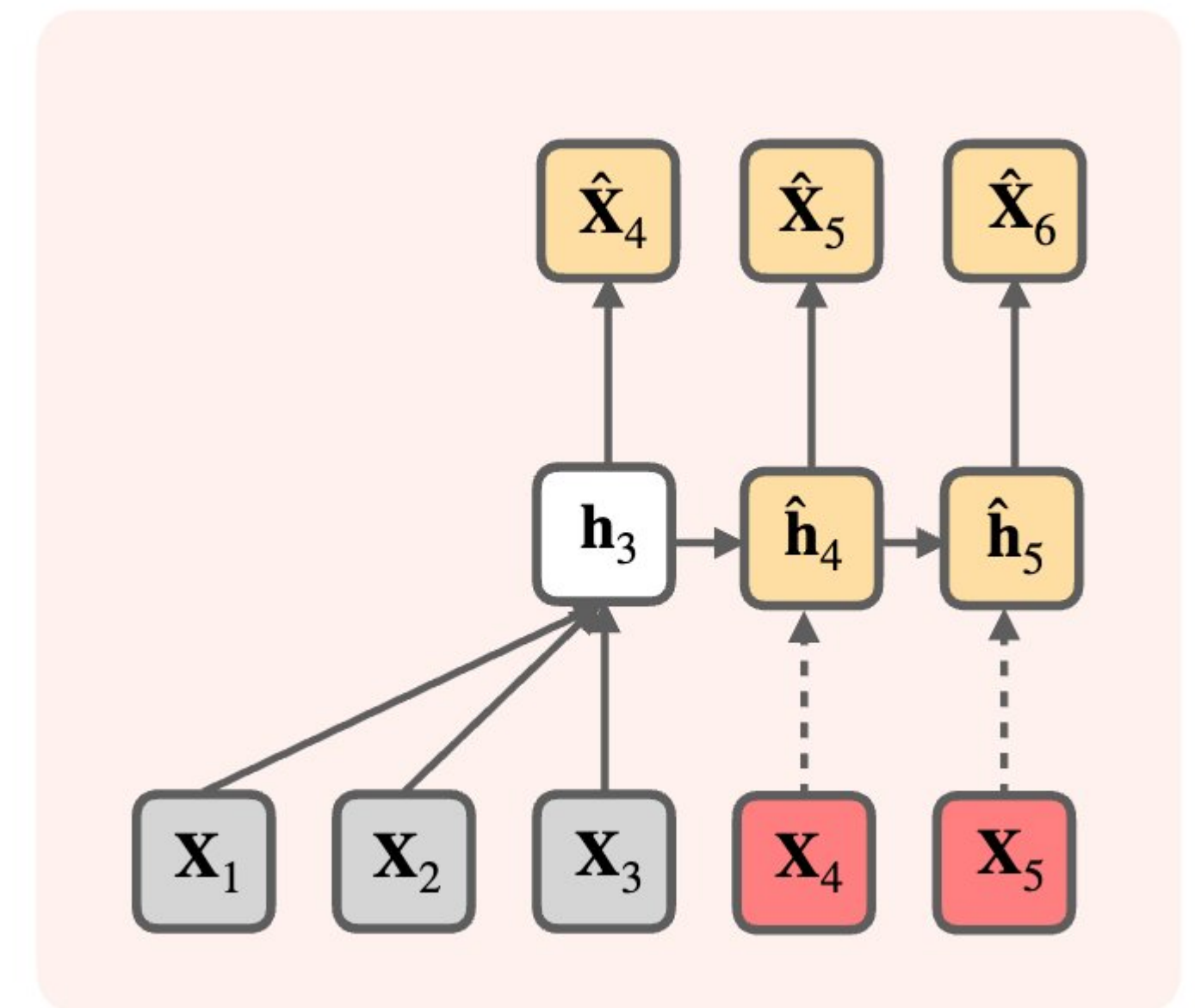


Next-Latent Prediction Transformers

Microsoft Research New York

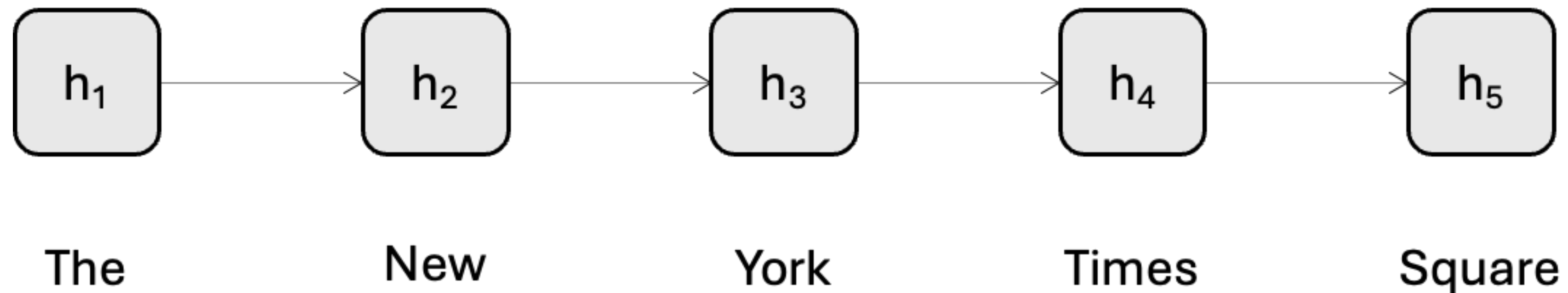


Jayden Teoh, Manan Tomar, Kwangjun Ahn, Edward S. Hu, Tim Pearce, Pratyusha Sharma, Akshay Krishnamurthy, Riashat Islam, Alex Lamb, John Langford

Introduction



- Recurrent neural networks (RNNs) naturally enforce compact representations
- But their training is not parallelizable...



Introduction



- Transformers replaced RNNs
- Self-attention enables **flexible look ups** over past tokens
- But there's no incentive to learn compact summaries of past tokens...

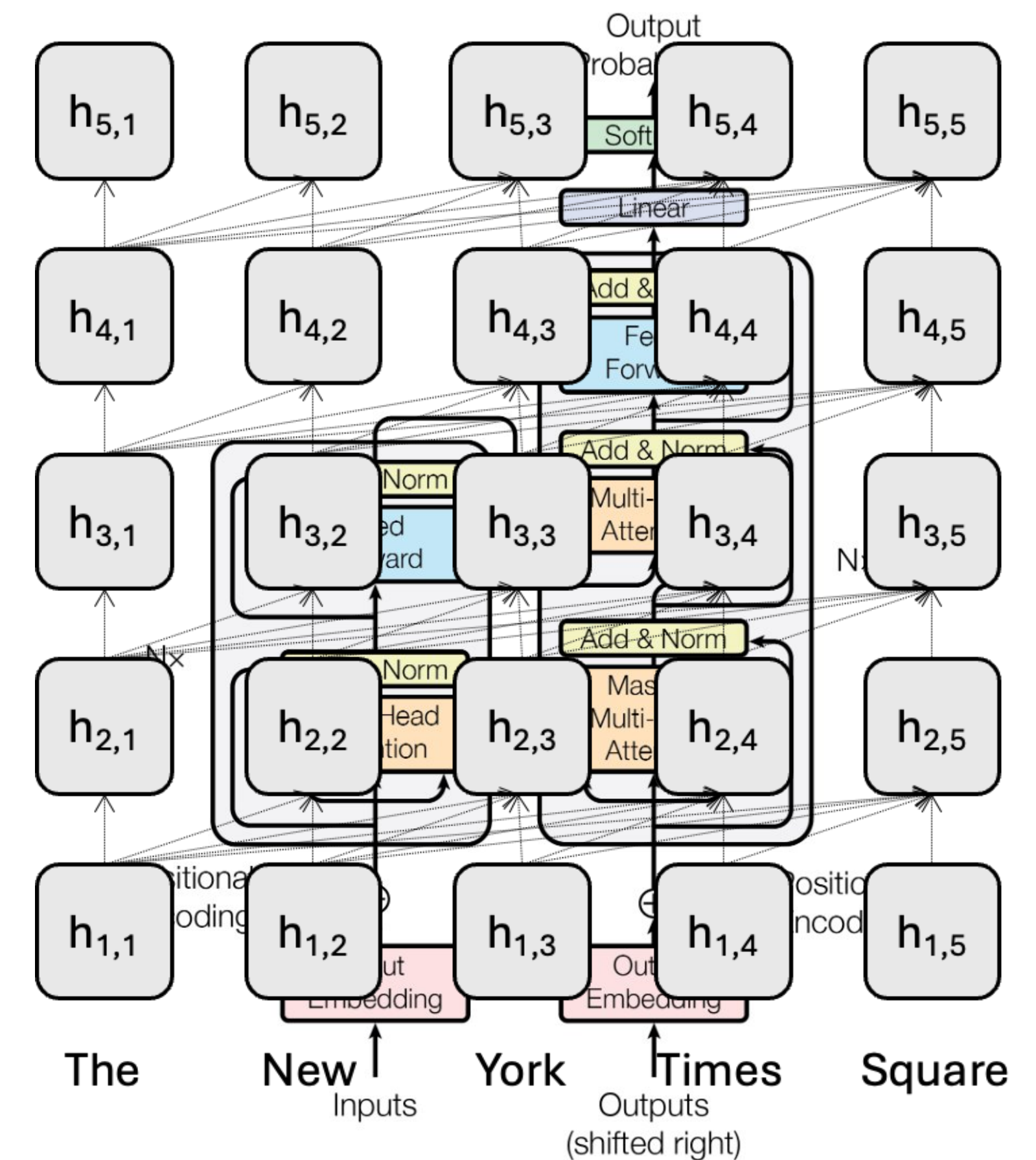
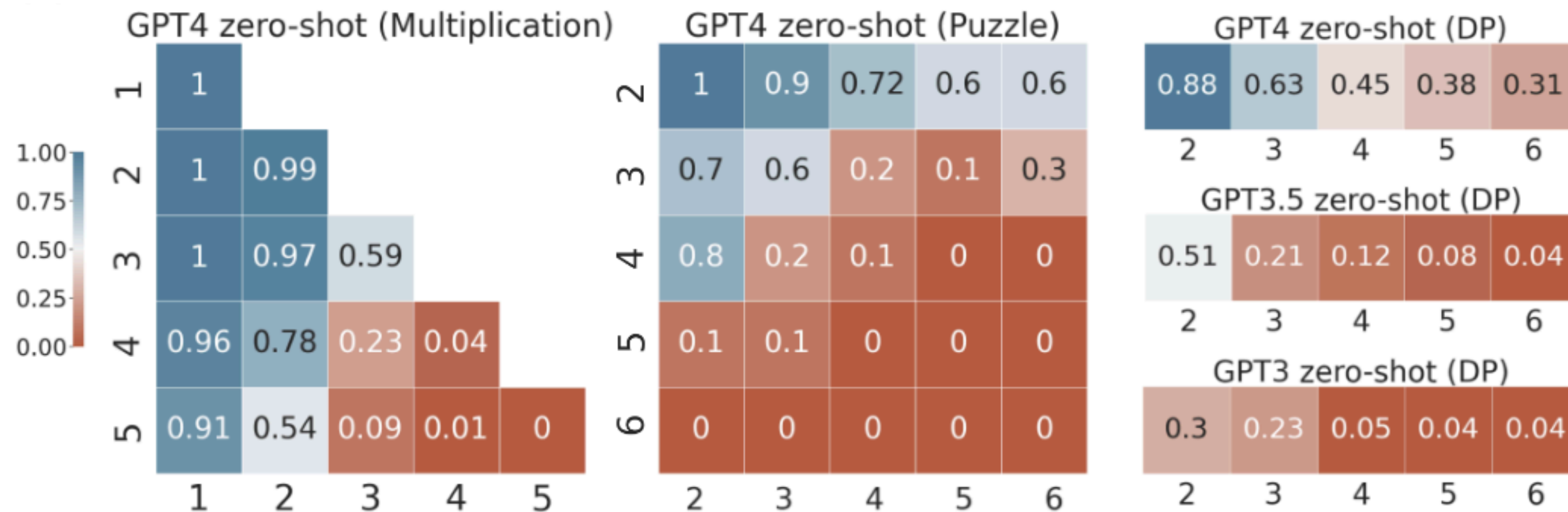


Figure 1: The Transformer - model architecture.

Introduction



- As a result, transformers often learn shortcut solutions that don't generalize [1, 2, 3, 4]



GPT3 performance degrades on compositional tasks

[1] Anil et al. Exploring Length Generalization in Large Language Models. NeurIPS 2022.

[2] Liu et al. Transformers Learn Shortcuts to Automata. ICLR 2023.

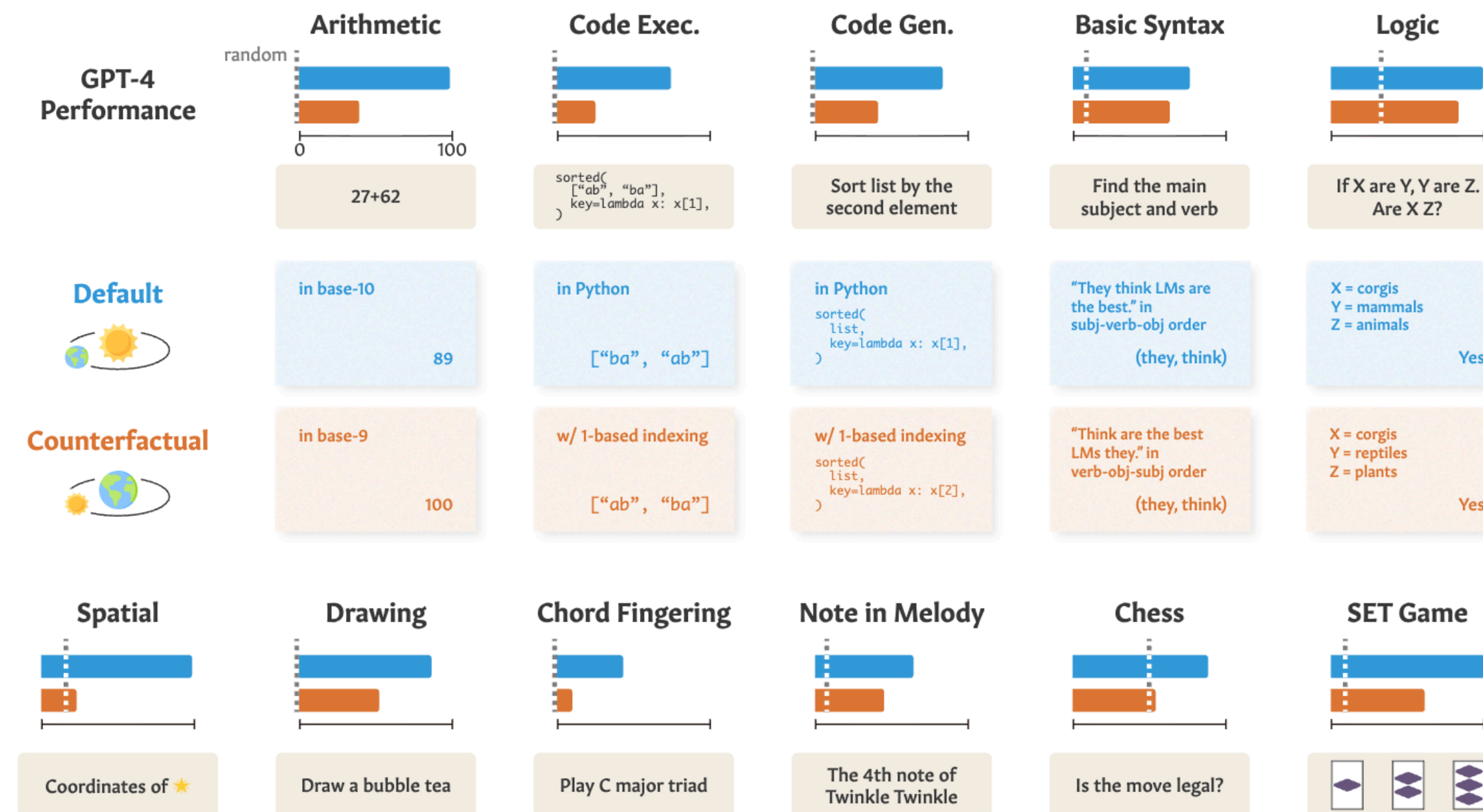
[3] Dziri et al. Faith and Fate: Limits of Transformers on Compositionality. NeurIPS 2023.

[4] Wu et al. Reasoning or reciting? exploring the capabilities and limitations of language models through counterfactual tasks. NAACL 2024.

Introduction



- As a result, transformers often learn shortcut solutions that don't generalize [1, 2, 3, 4]



GPT3 performance degrades on counterfactual tasks

[1] Anil et al. Exploring Length Generalization in Large Language Models. NeurIPS 2022.

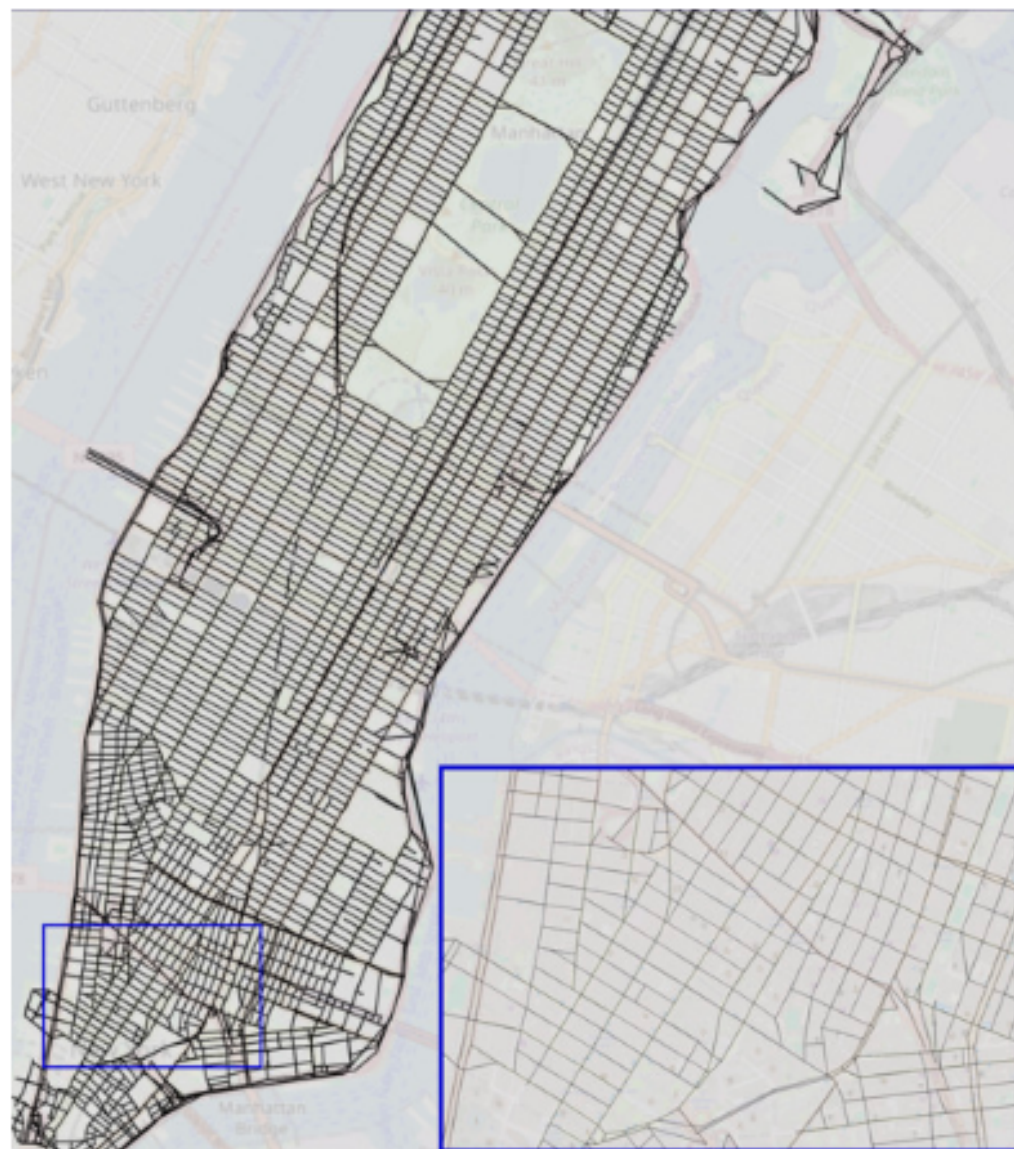
[2] Liu et al. Transformers Learn Shortcuts to Automata. ICLR 2023.

[3] Dziri et al. Faith and Fate: Limits of Transformers on Compositionality. NeurIPS 2023.

[4] Wu et al. Reasoning or reciting? exploring the capabilities and limitations of language models through counterfactual tasks. NAACL 2024.

Introduction

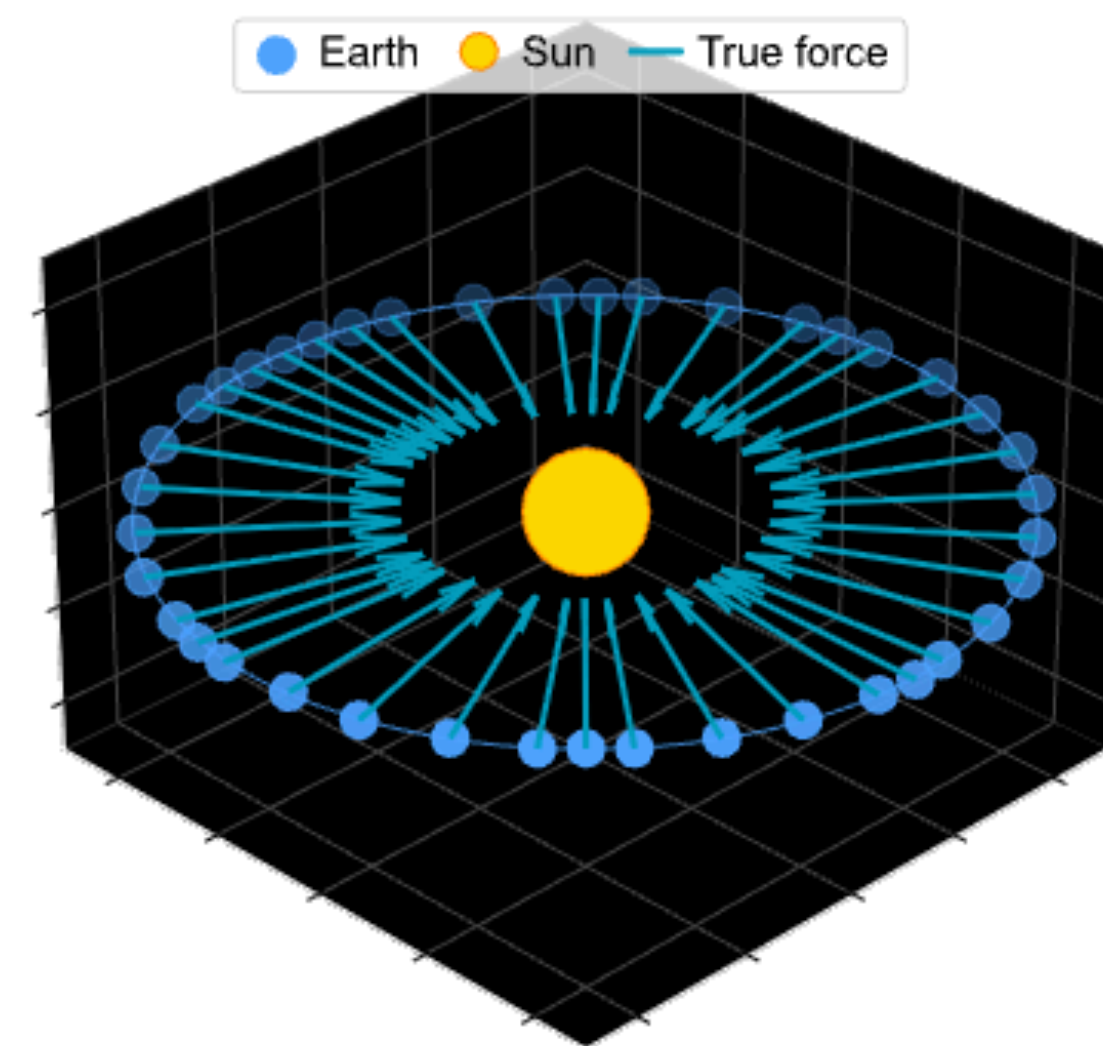
- Transformers also learn poor world models [5, 6]



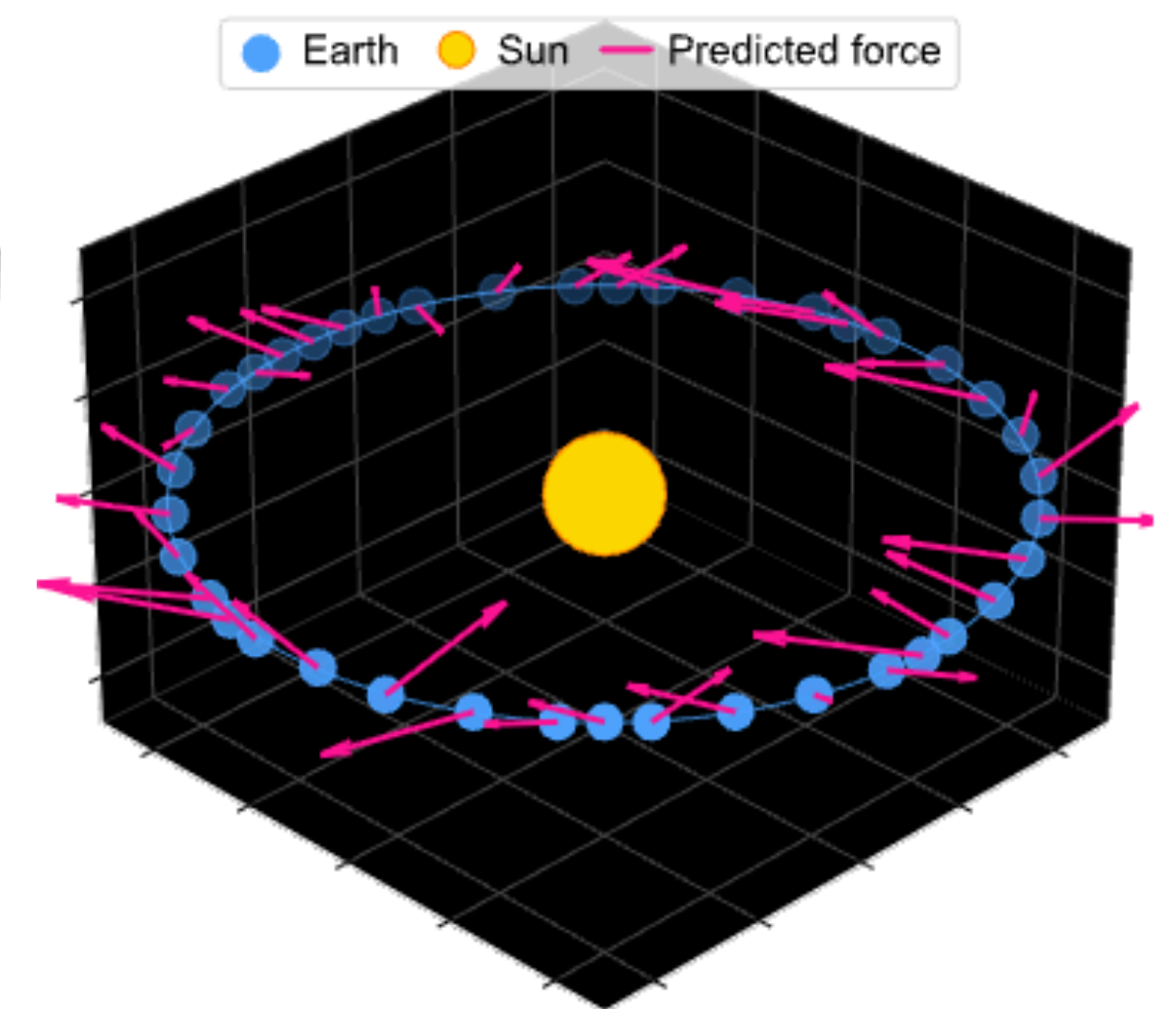
Manhattan (real)



Manhattan (transformer)



Earth (real)



Earth (transformer)

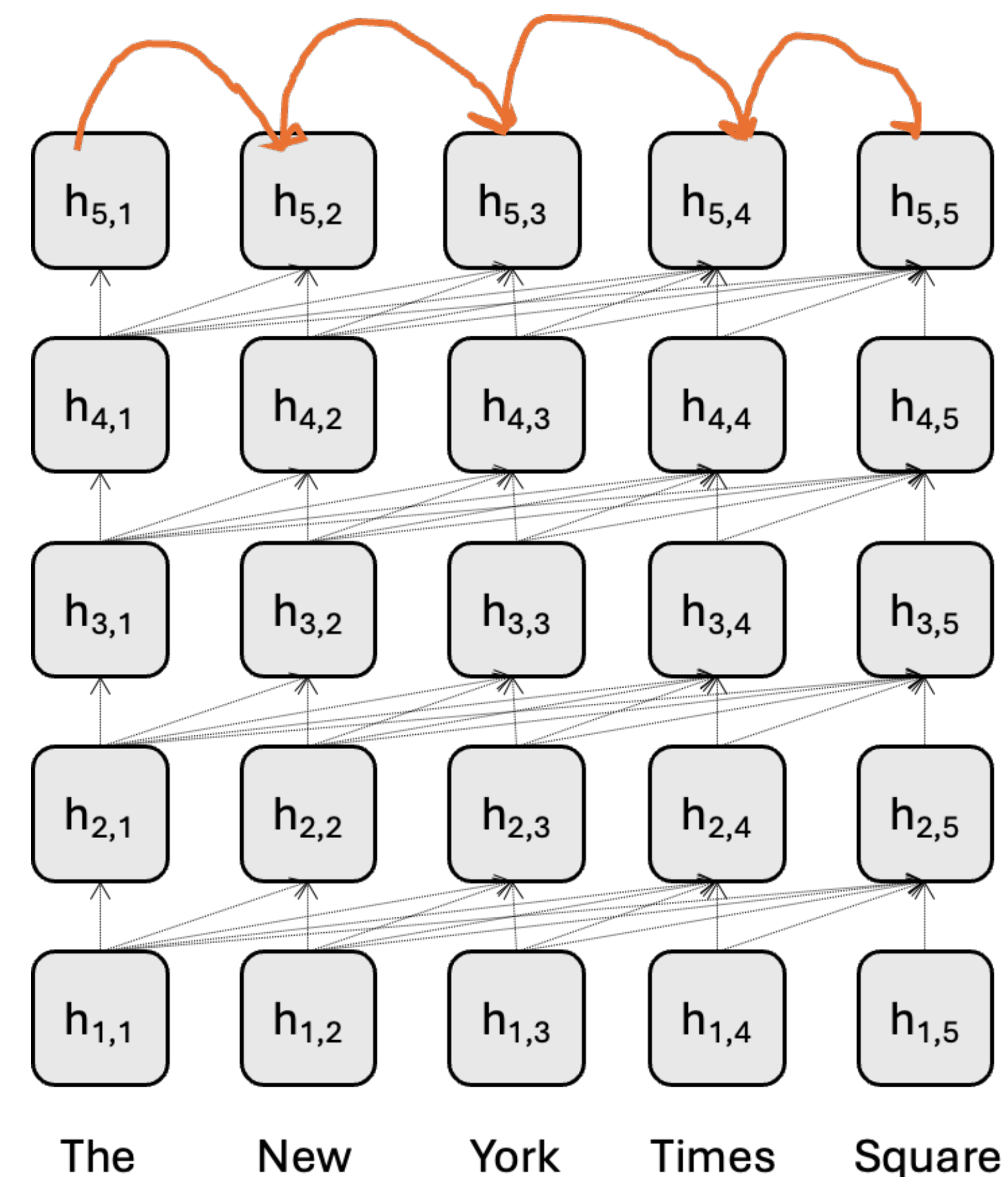
[5] Vafa et al. Evaluating the World Model Implicit in a Generative Model. NeurIPS 2024.

[6] Vafa et al. What Has a Foundation Model Found? Using Inductive Bias to Probe for World Models. ICML 2025.

Methodology

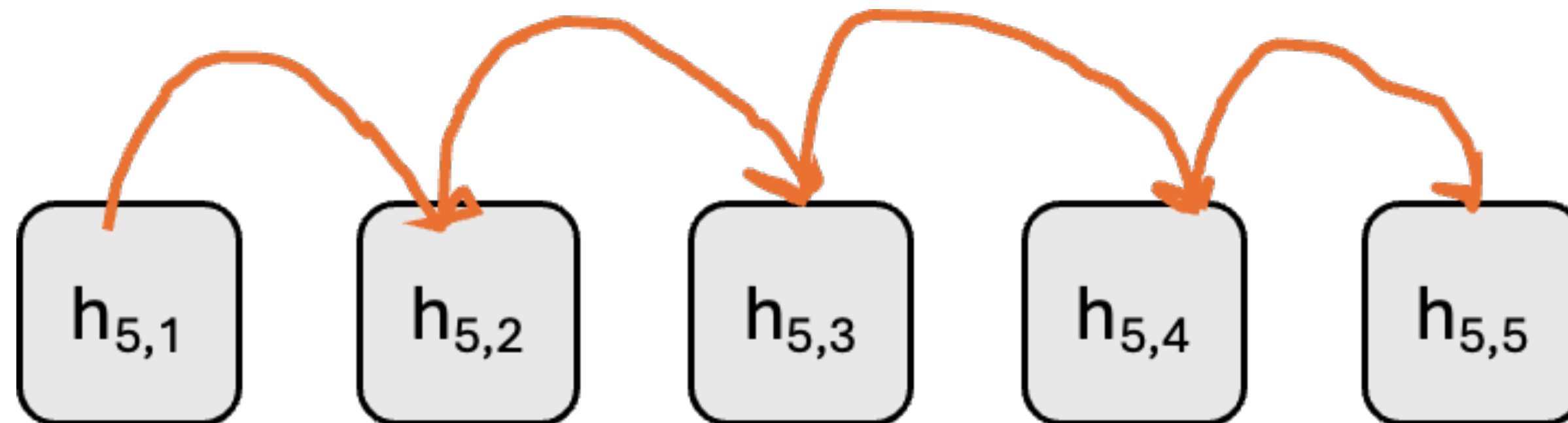


- **Motivation #1:** we want recurrent-like dynamics and compression in transformers
- **Motivation #2:** we want to retain parallel training

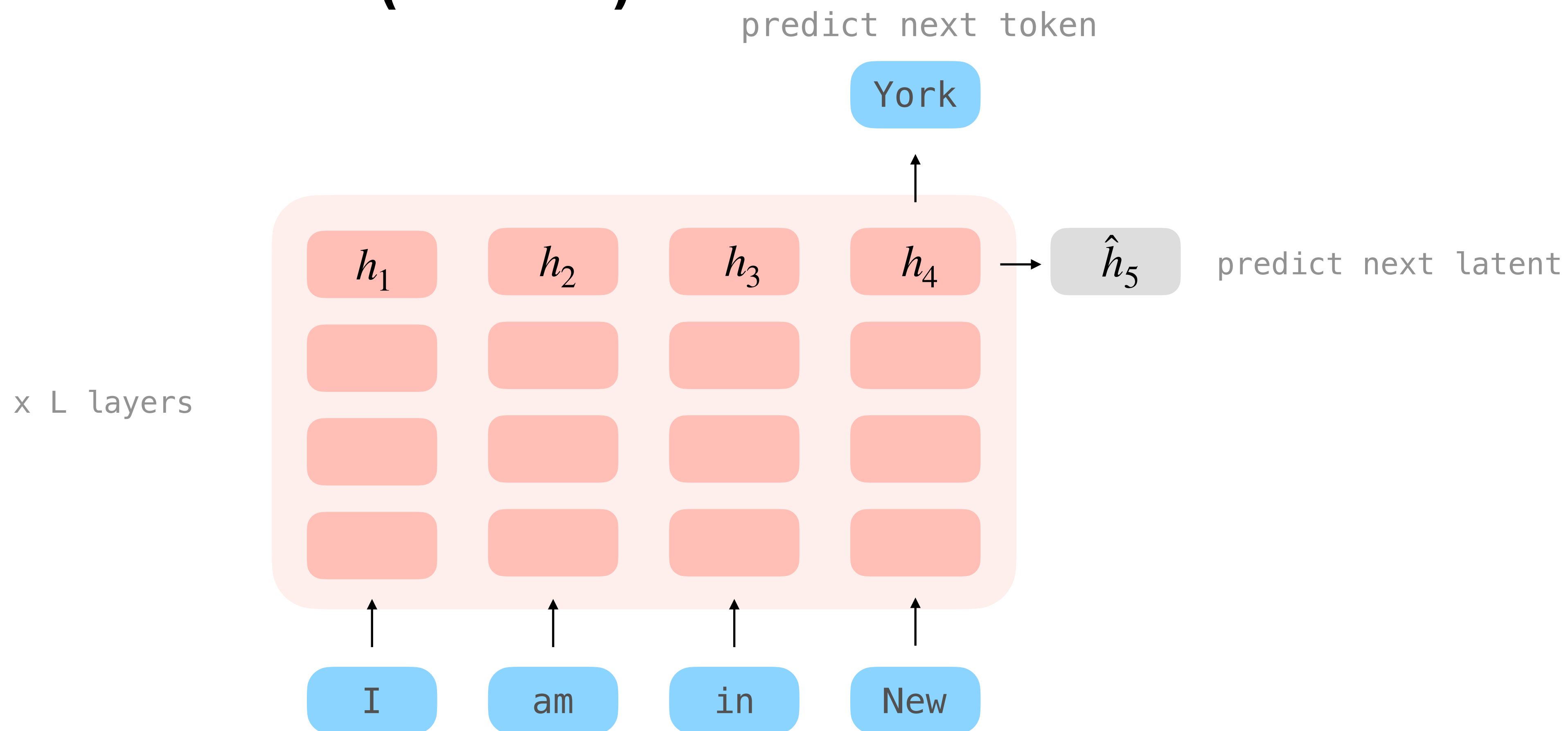


Next-Latent Prediction (NextLat)

- We train the transformer such that each of its latent state can predict the next latent state (conditioned on the next token)



Next-Latent Prediction (NextLat)



Next-Latent Prediction (NextLat)

1) next-token cross entropy loss

$$\mathcal{L}_{\text{next-token}} = -\log p_{\theta}(X_{t+1} | \mathbf{h}_t)$$

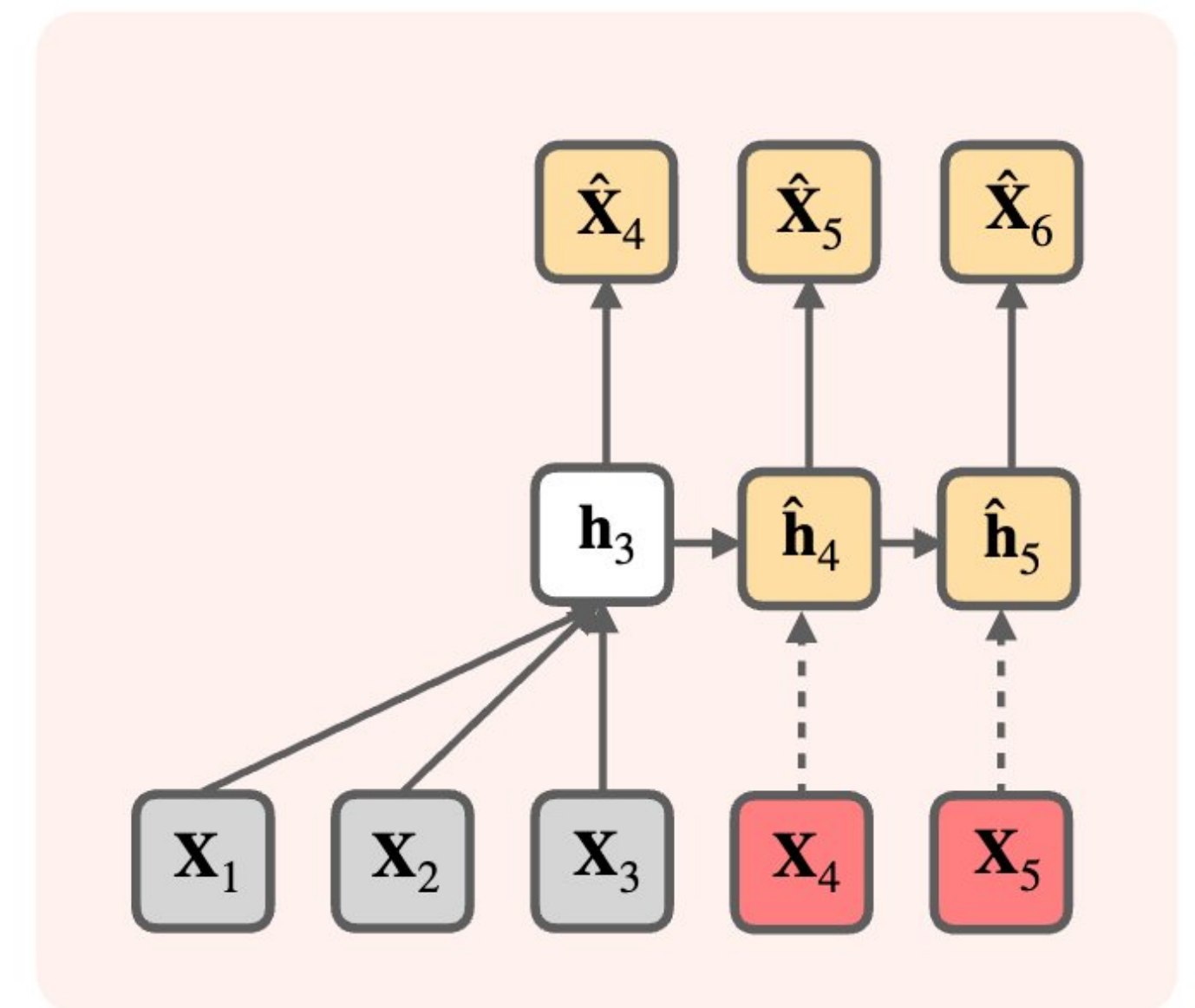
2) next-latent regression loss

$$\mathcal{L}_{\text{next-h}} = \|\hat{\mathbf{h}}_{t+1} - \text{sg}[\mathbf{h}_{t+1}]\|_2^2$$

3) *Optionally*: next-latent token KL (or cross entropy)

$$\mathcal{L}_{\text{KL}} = D_{\text{KL}} \left(p_{\theta}^{\text{sg}}(\cdot | \text{sg}[\mathbf{h}_{t+1}]) \parallel p_{\theta}^{\text{sg}}(\cdot | \text{sg}[\hat{\mathbf{h}}_{t+1}]) \right)$$

Next-Latent Prediction



Why NextLat?

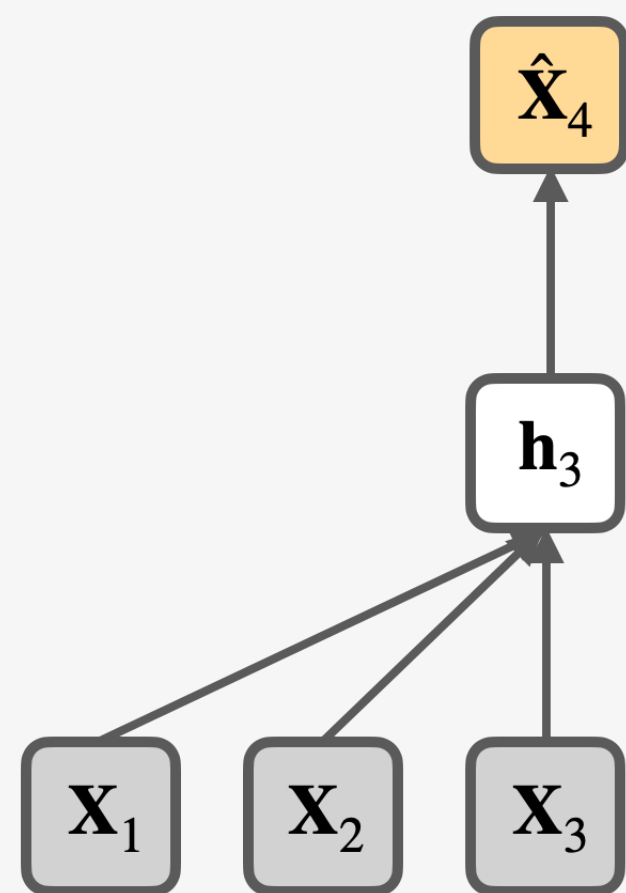
- 1) **Stronger** representations for world modeling, reasoning, and planning
- 2) **Faster** inference than Multi-Token Prediction
- 3) **Better** data-efficiency

Motivation

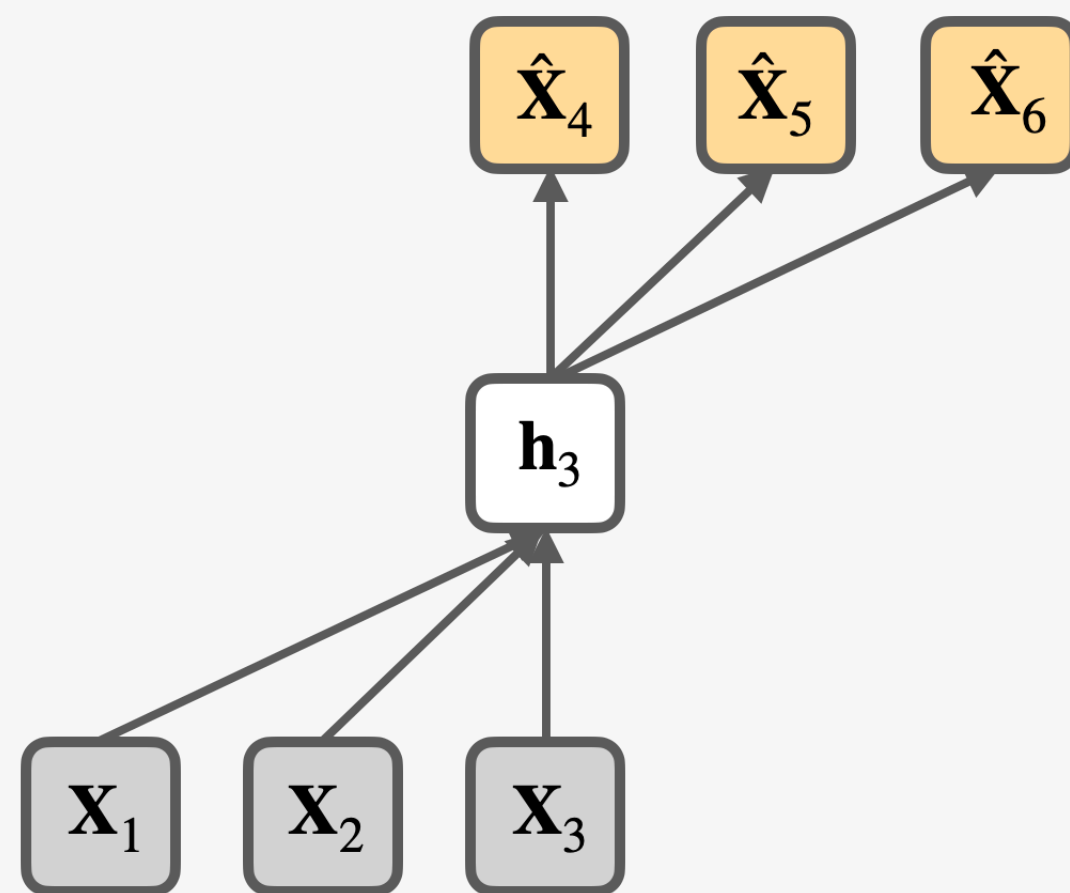


Why NextLat?

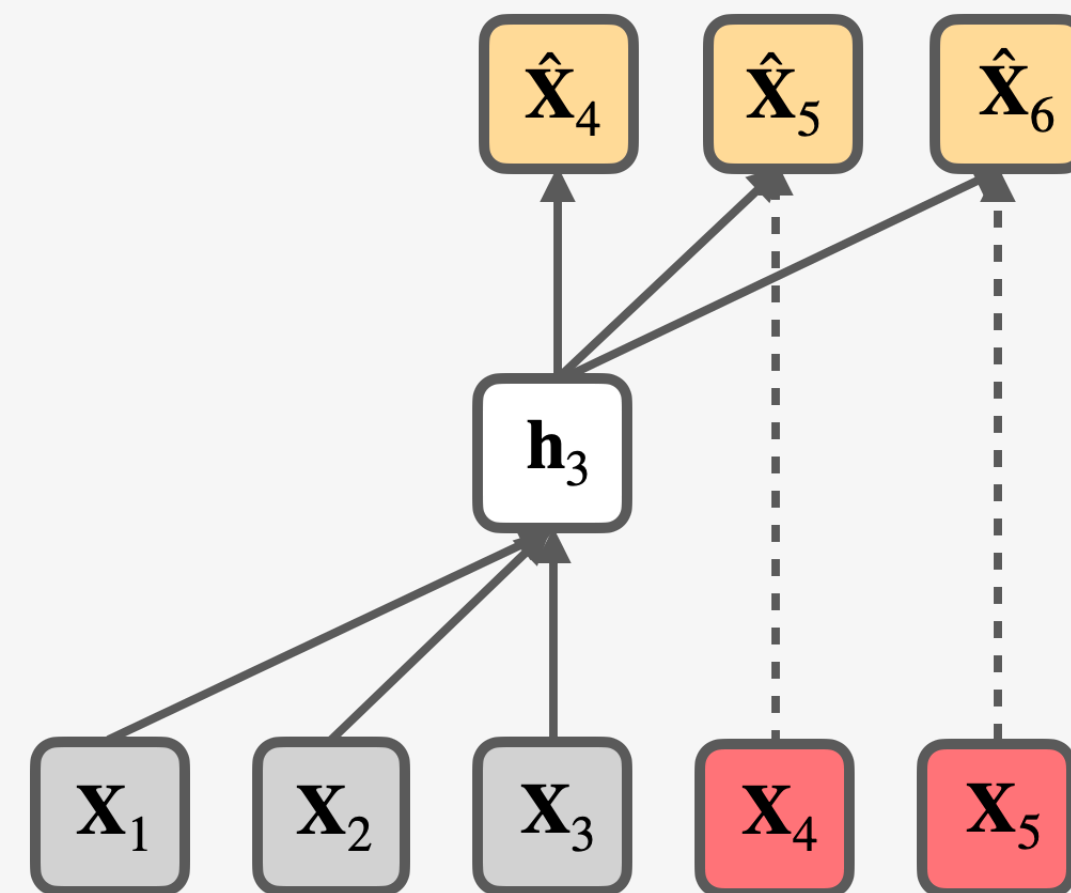
Next-Token Prediction



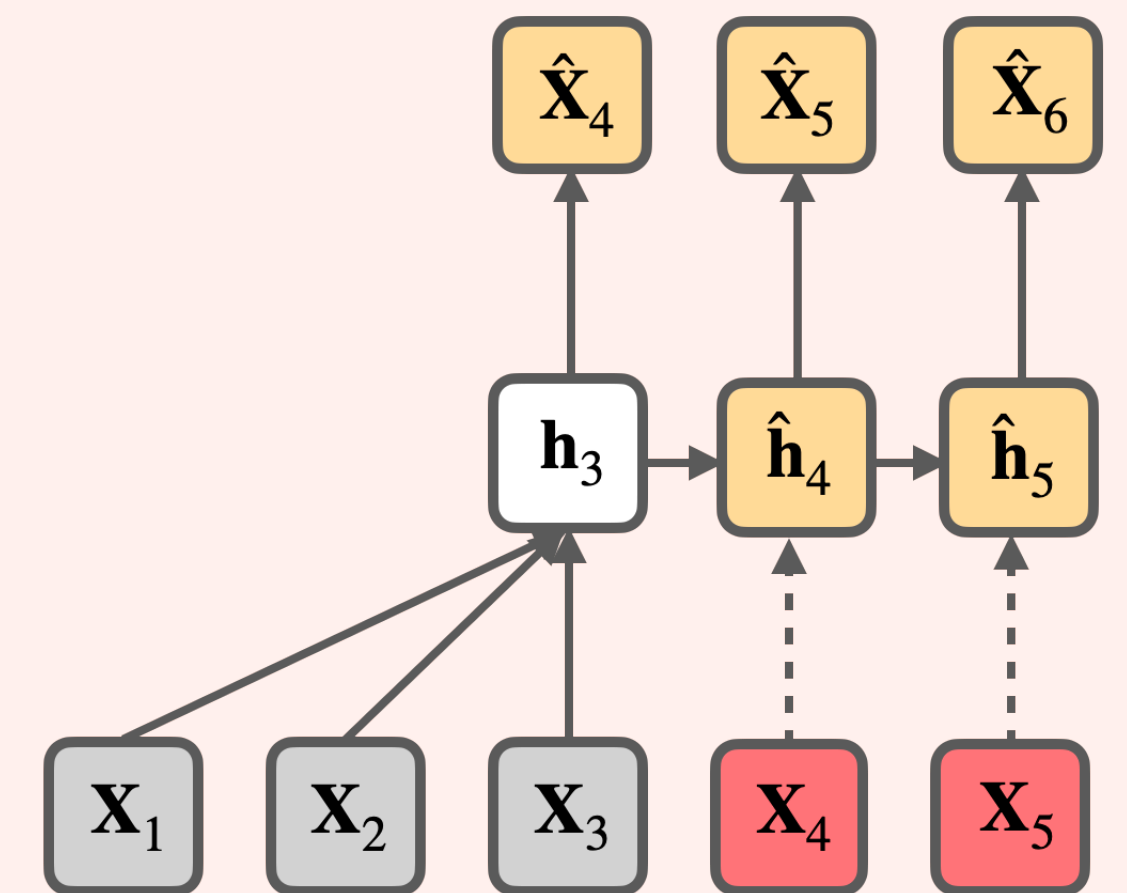
Multi-Token Prediction



Joint Multi-Token Prediction



Next-Latent Prediction



Input Tokens

Predictions

Teacher-forced Tokens

Better Representations



- Train models on trajectories from Manhattan taxi rides

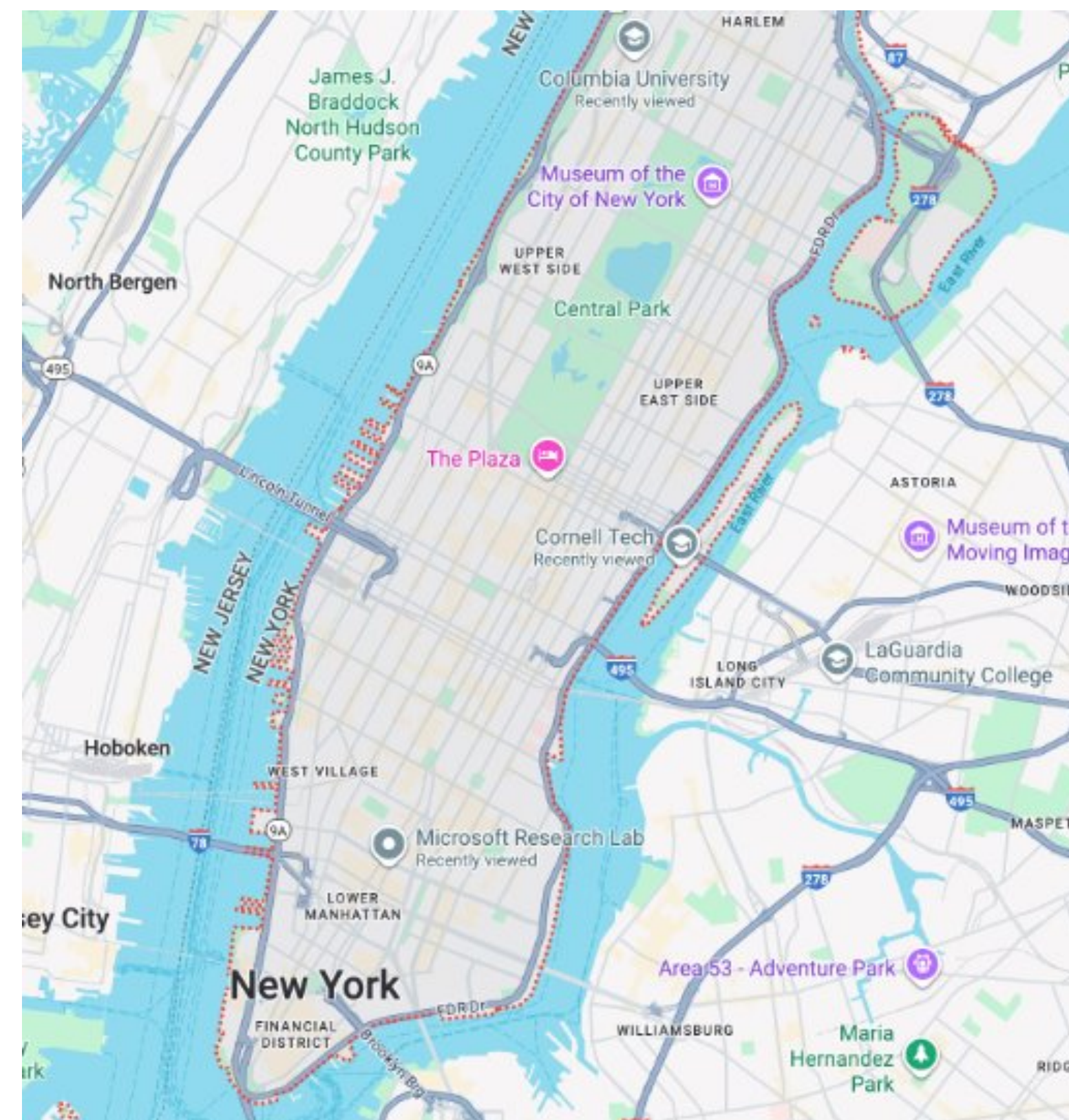
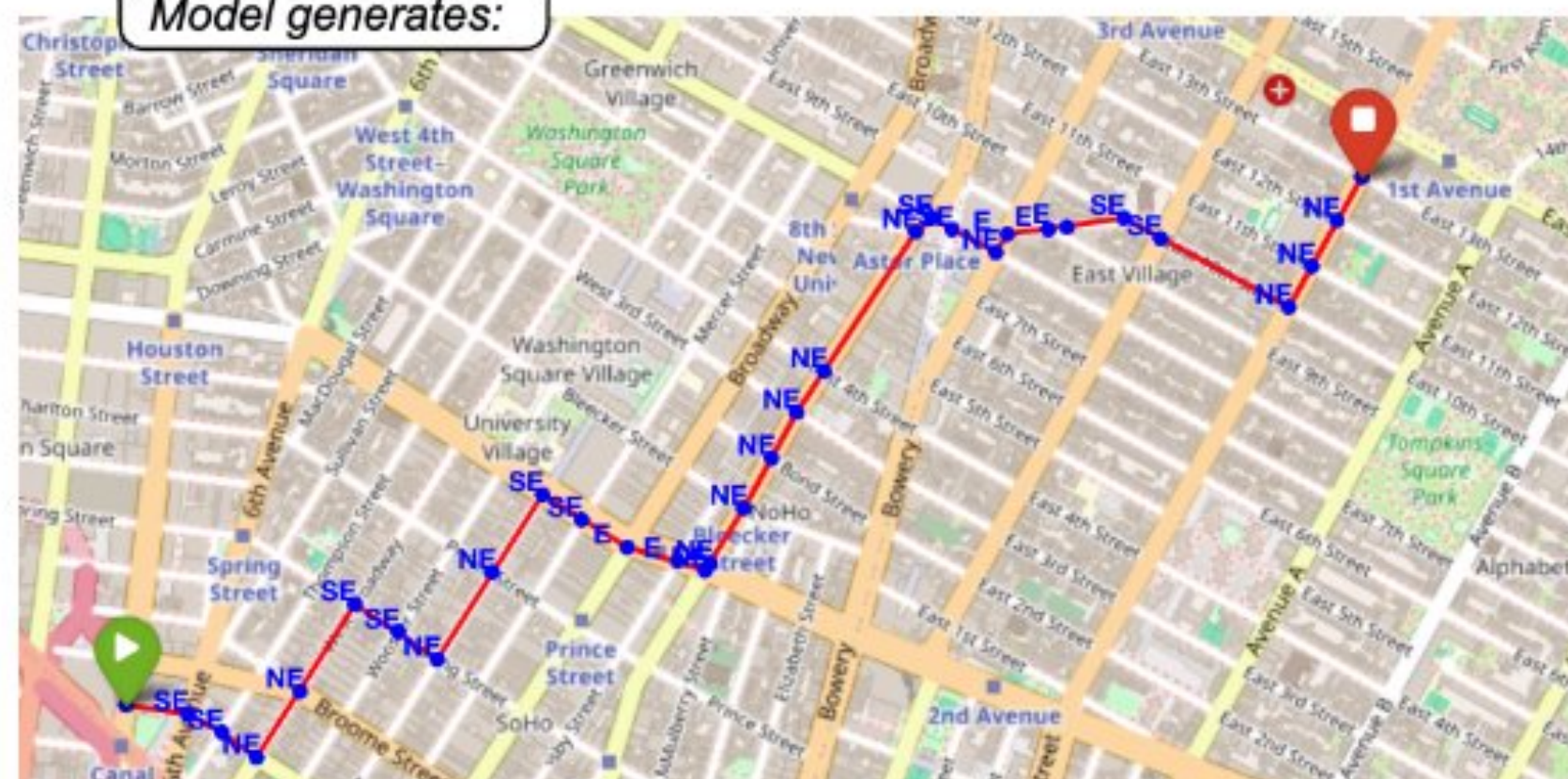
Training sequences

820 210 N E E SE W W N SE N end
193 450 E E E N E S E end
301 592 N N N SE S S S S S end
...

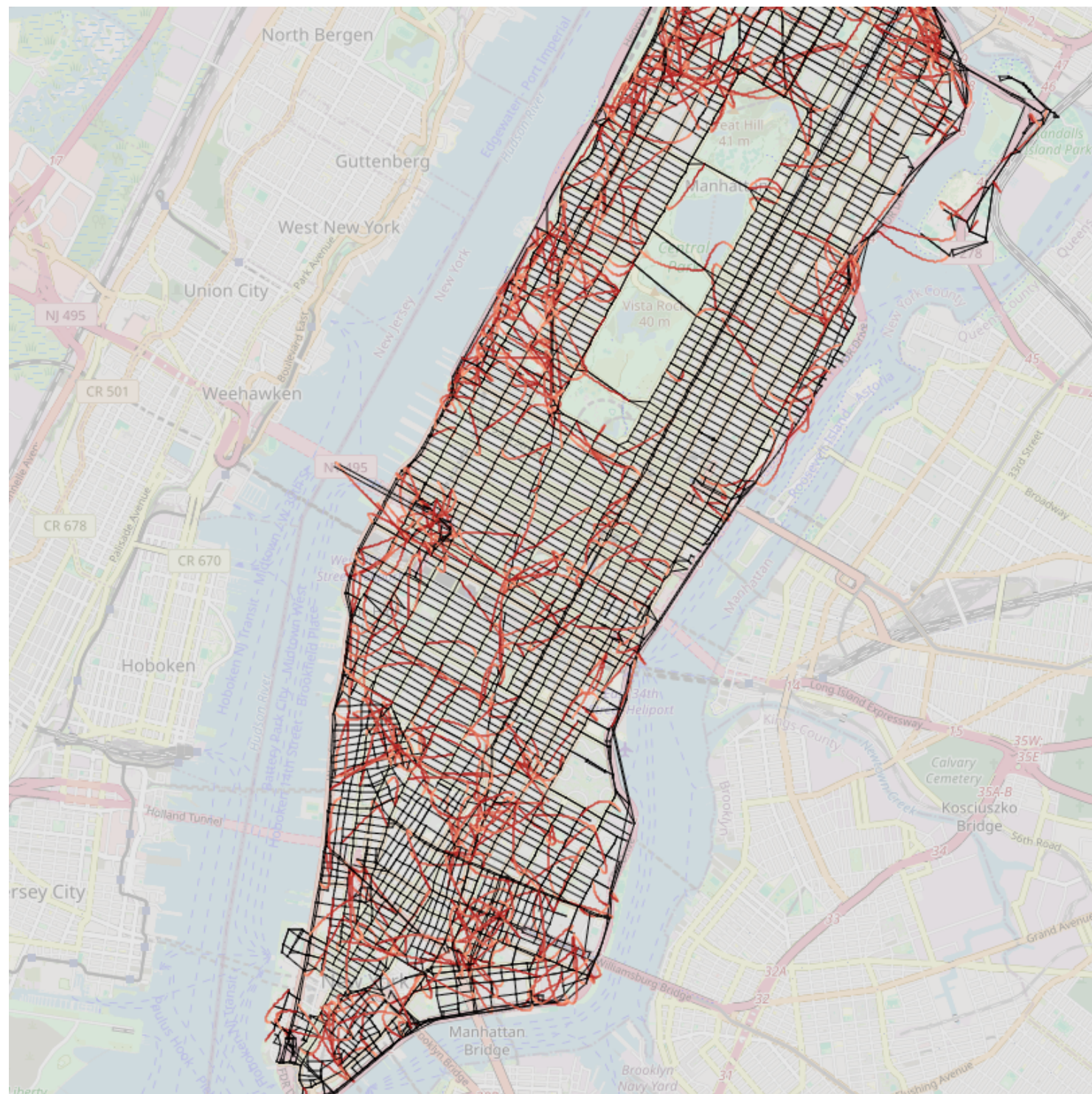
Contexts for evaluation

110 244
850 820
592 301
...

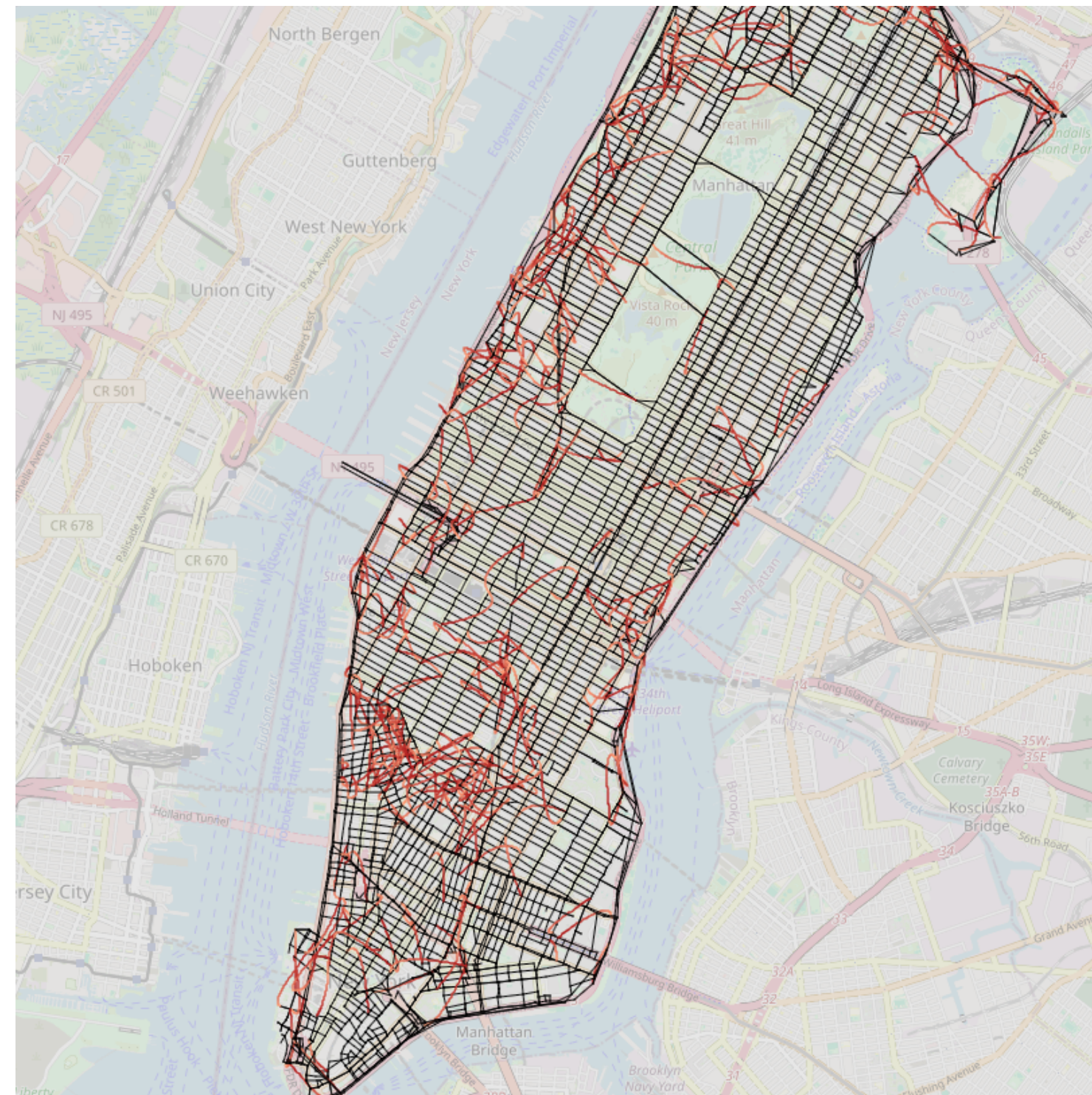
Model generates:



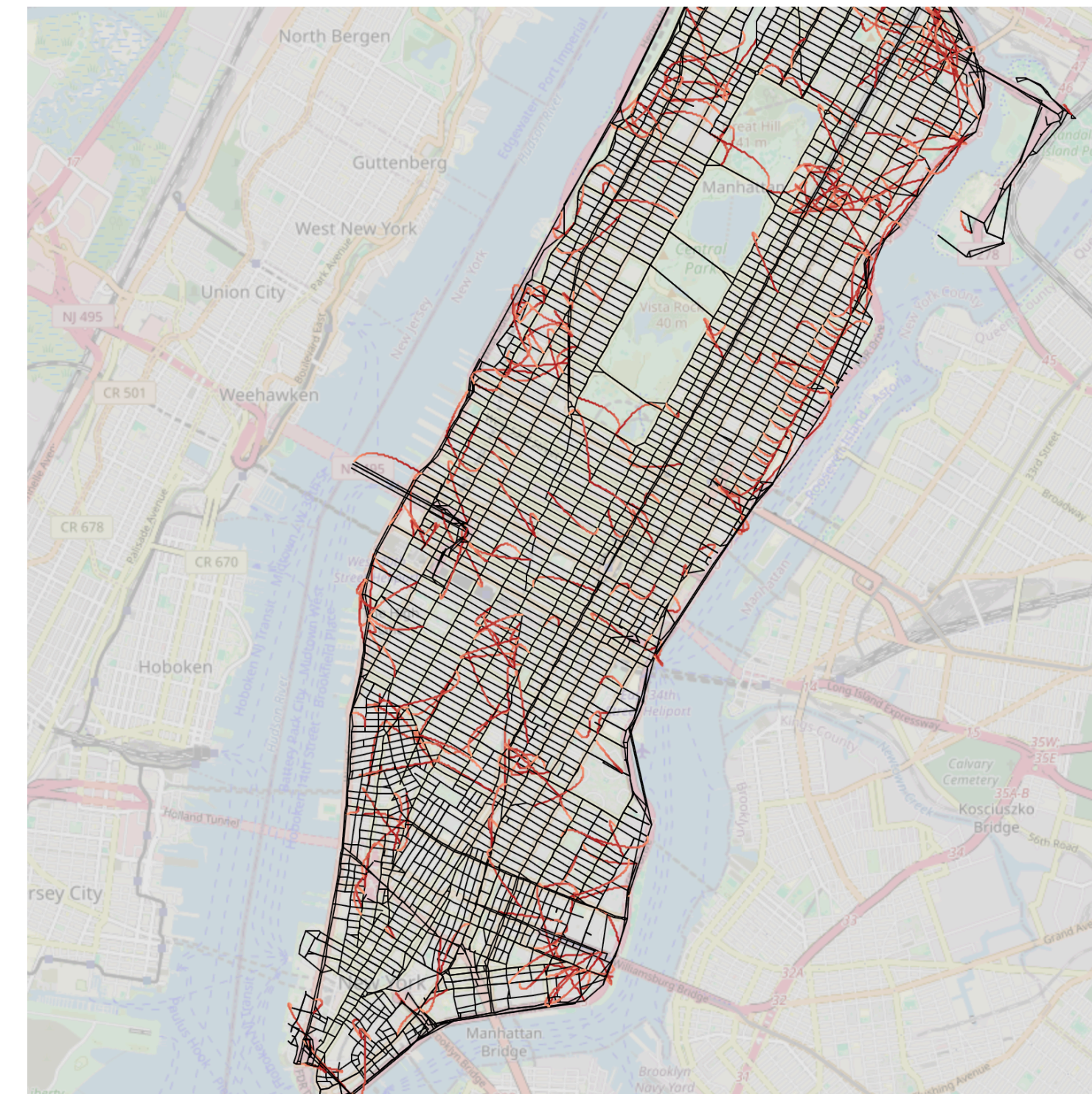
World Modeling



Next-Token Prediction



Multi-Token Prediction



Next-Latent Prediction

World Modeling

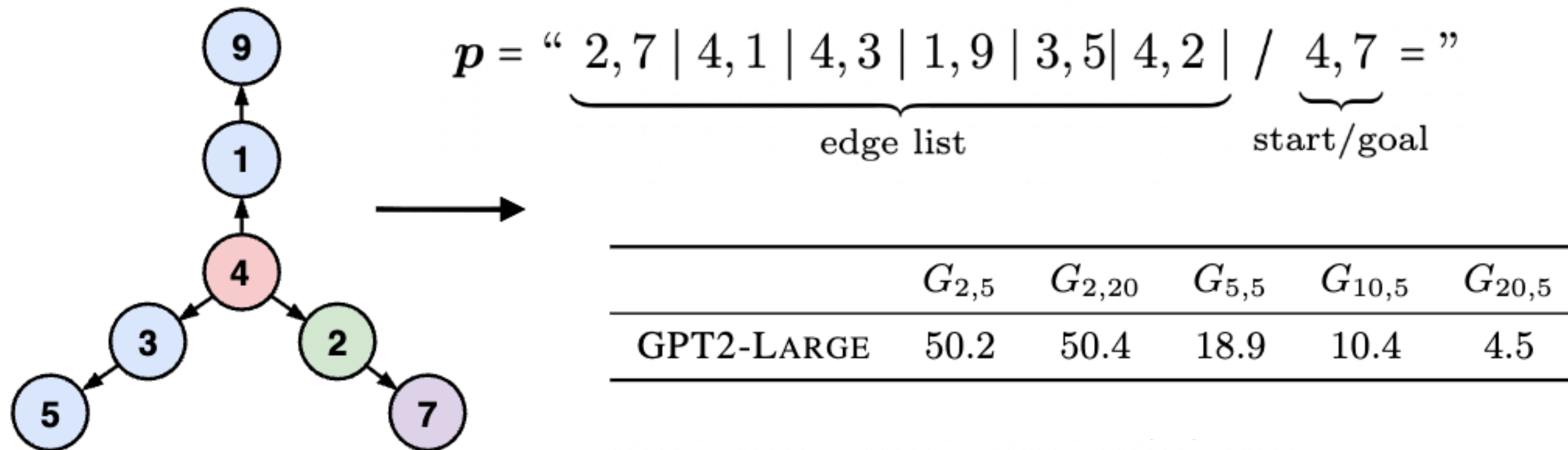


	Next-Token Test (\uparrow)	Valid Trajectories (\uparrow)	Sequence Compression (\uparrow)	Effective Latent Rank (\downarrow)	Detour Robustness (\uparrow)
GPT	100%	97.0%	0.65	160.1	85.0%
MTP	100%	98.1%	0.64	57.7	95.0%
JTP	100%	97.1%	0.32	215.8	87.0%
NextLat	100%	98.7%	0.71	52.7	95.0%
True world model	100%	100%	1.00	—	100%

Table 1: Comparison of GPT, MTP, JTP, and NextLat trained on Manhattan taxi rides against the true world model across several metrics.

Planning

- Path-Star Graph [7]. Simple planning task, yet next-token prediction can't solve it

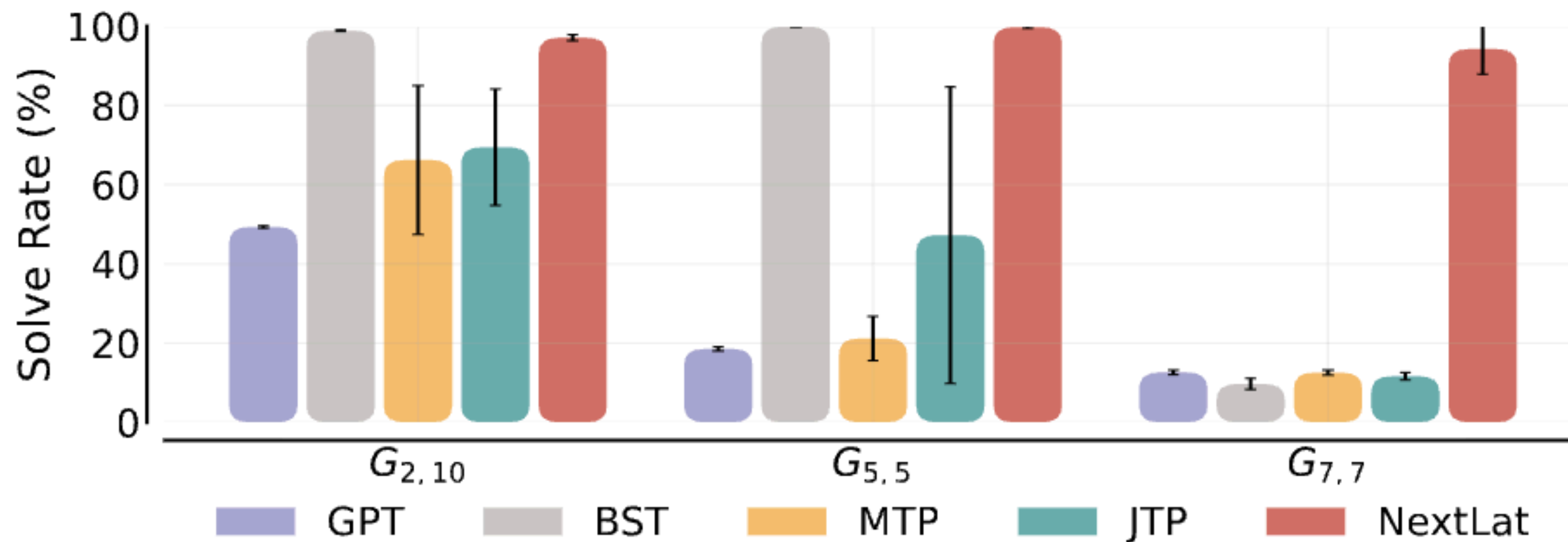
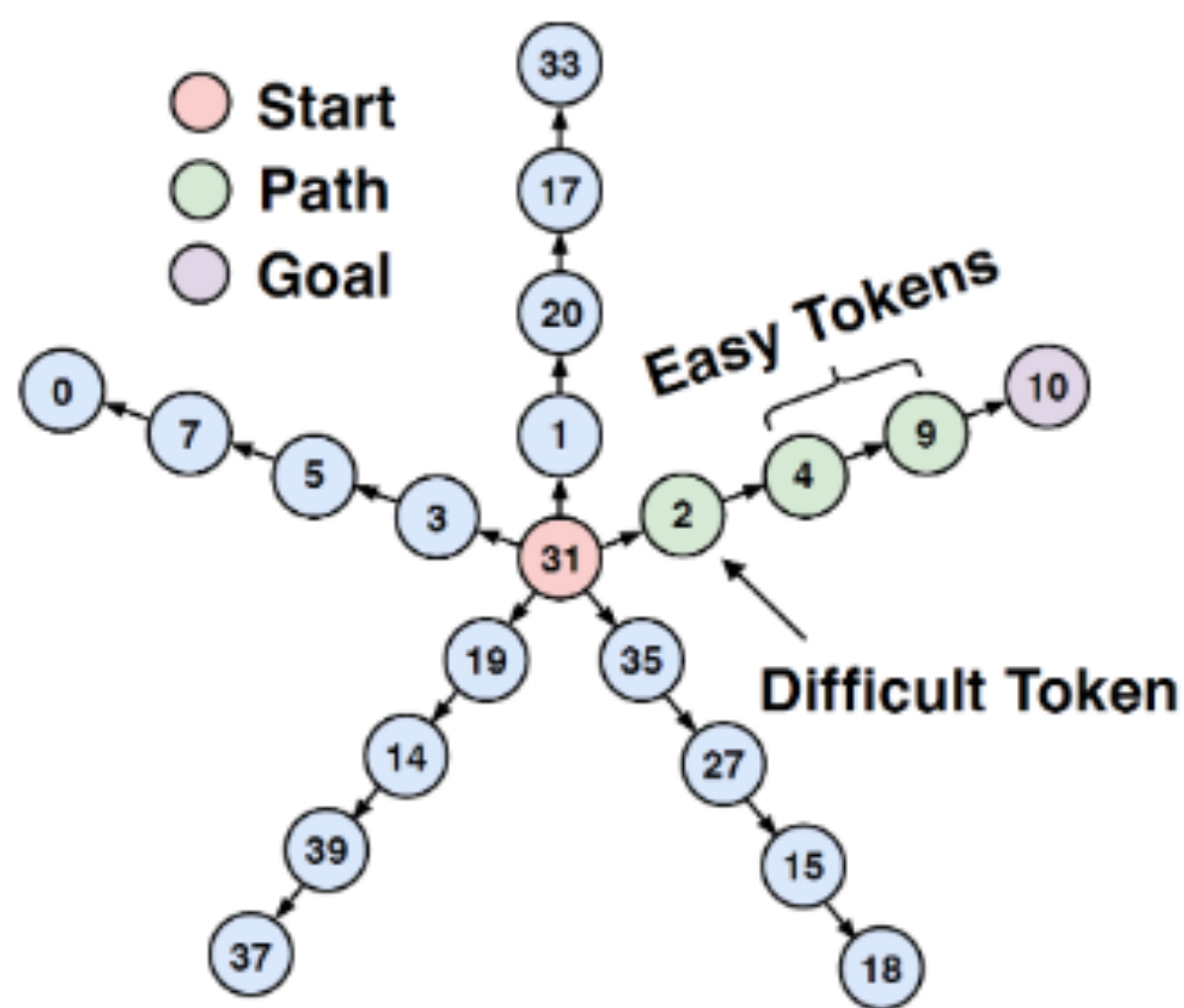


[7] Bachmann et al. The pitfalls of next-token prediction. ICML 2024.

Planning



- Only NextLat can solve all graphs



Reasoning



- Countdown Game
 - GPT-4 only gets 4% on this task!



Input: <14, 83, 88, 91>

Goal: 23

Solution: $83 - 14 = 69$, $91 - 88 = 3$, $69 \div 3 = 23$

Reasoning



- NextLat has incredible performance just by predicting next-latent
- MTP methods have to learn to predict multiple tokens ahead to catch up

Model	Horizon (d)	Accuracy (%)
GPT	–	33.1
BST	–	42.3
MTP	<u>1</u>	<u>39.2</u>
	4	49.7
	<u>8</u>	<u>57.3</u>
JTP	<u>1</u>	<u>39.0</u>
	4	49.4
	<u>8</u>	<u>55.0</u>
NextLat	<u>1</u>	<u>54.8</u>
	4	<u>57.6</u>
	<u>8</u>	58.7

Figure 5: Performance on Countdown. Best result is **bolded**, and second best is underlined.

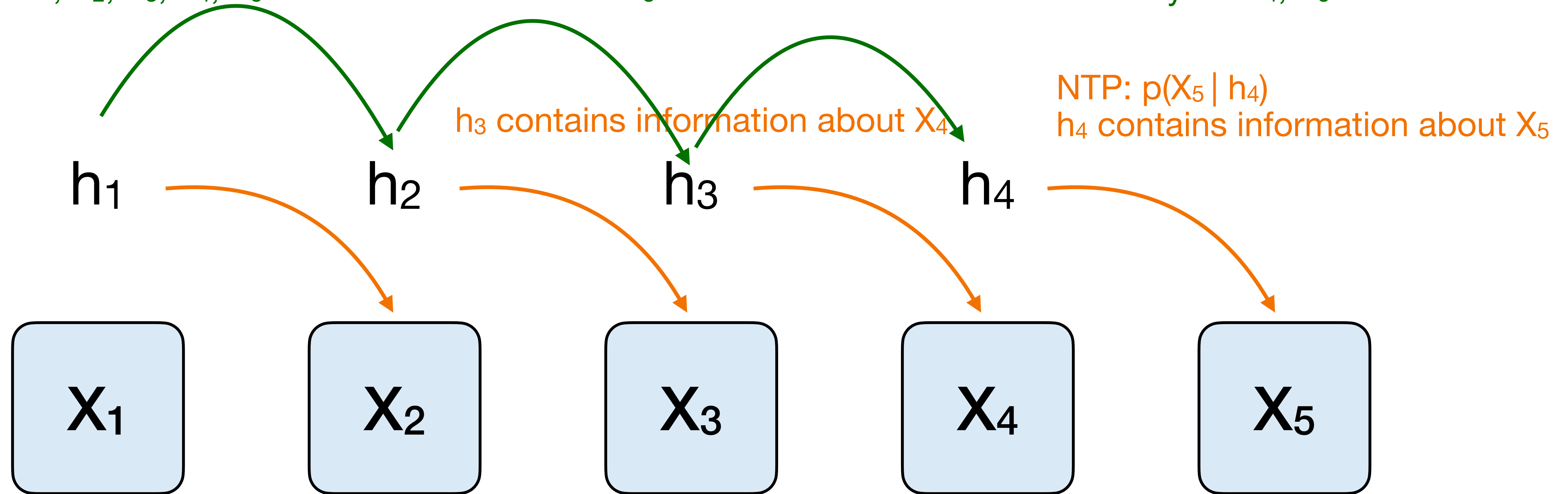
Stronger Representations



h_1 compresses all information needed for the future, i.e., X_2, X_3, X_4, X_5

NextLat: $p(h_4 | h_3, X_4)$

h_3 contains all information necessary for X_4, X_5

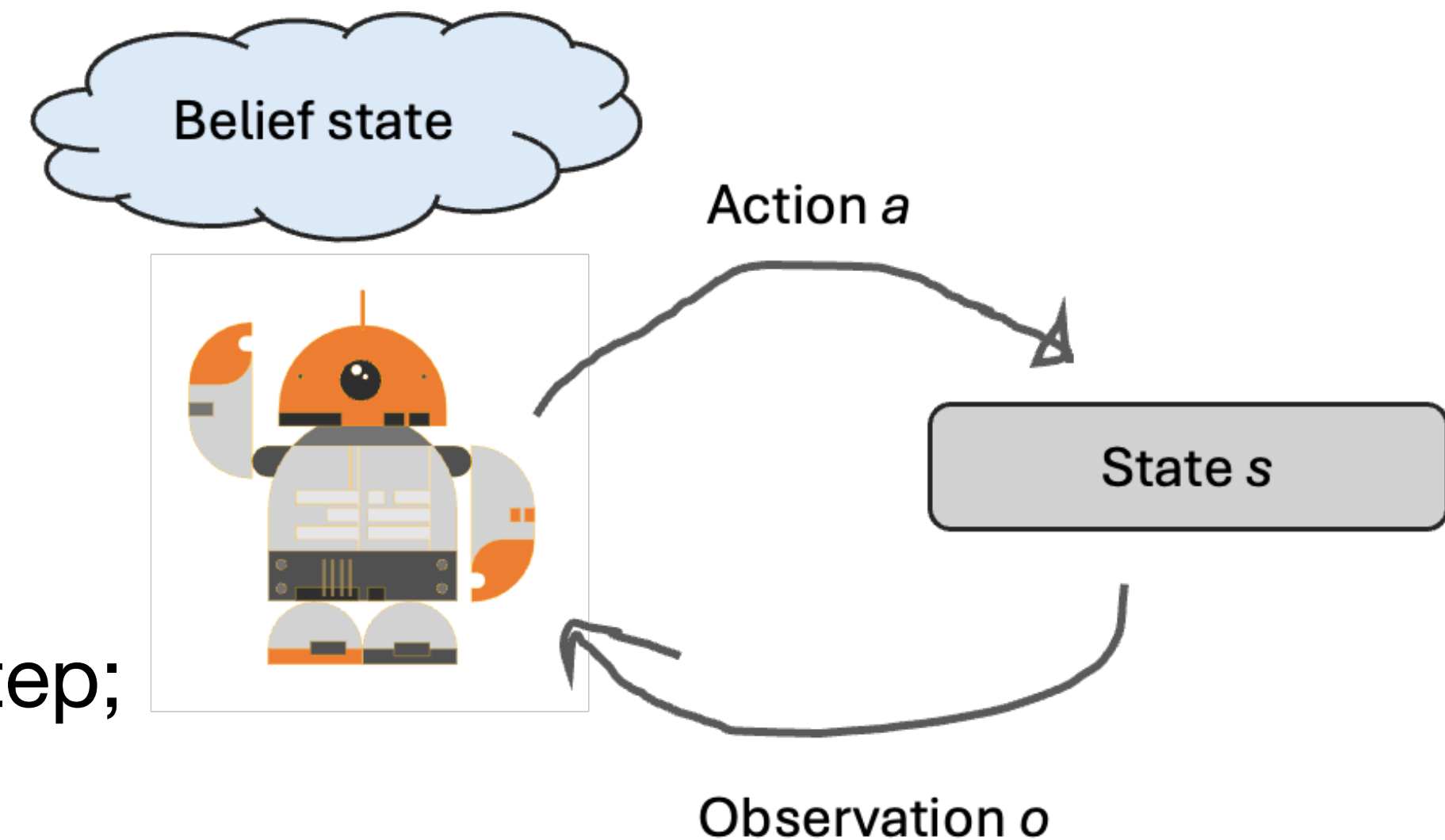


Every h_t is a compressed summary of the history for decoding future tokens!

(not a property with just next-token or multi-token prediction)

Belief States

- At convergence, the transformer's hidden state will become *belief states*
- **Intuition:** force transformers to compress history at every step; stop relying on lazy self-attention

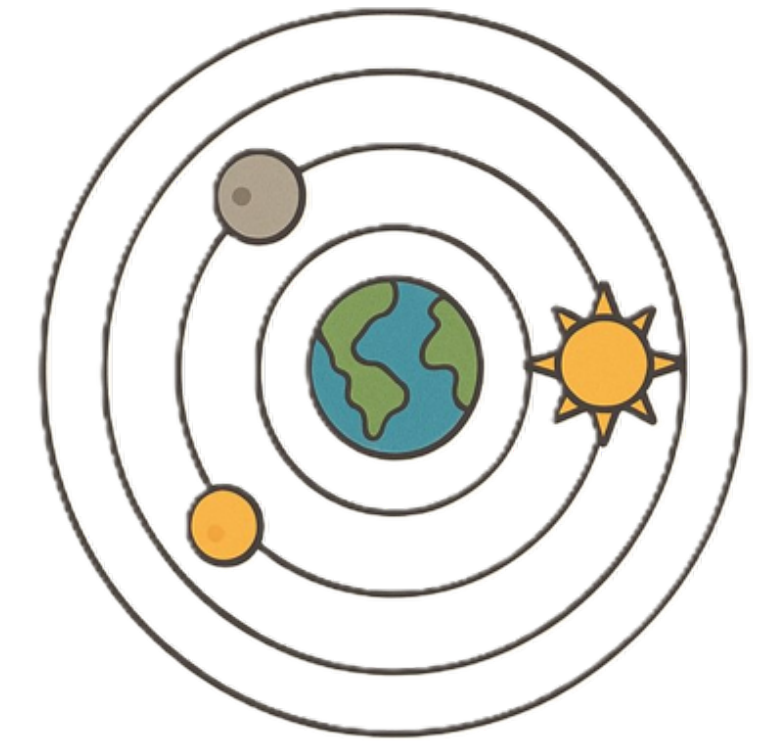


Motivation



Compact World Models

- In the 2nd century, Claudius Ptolemy proposed a geocentric model that placed Earth at the center of the universe

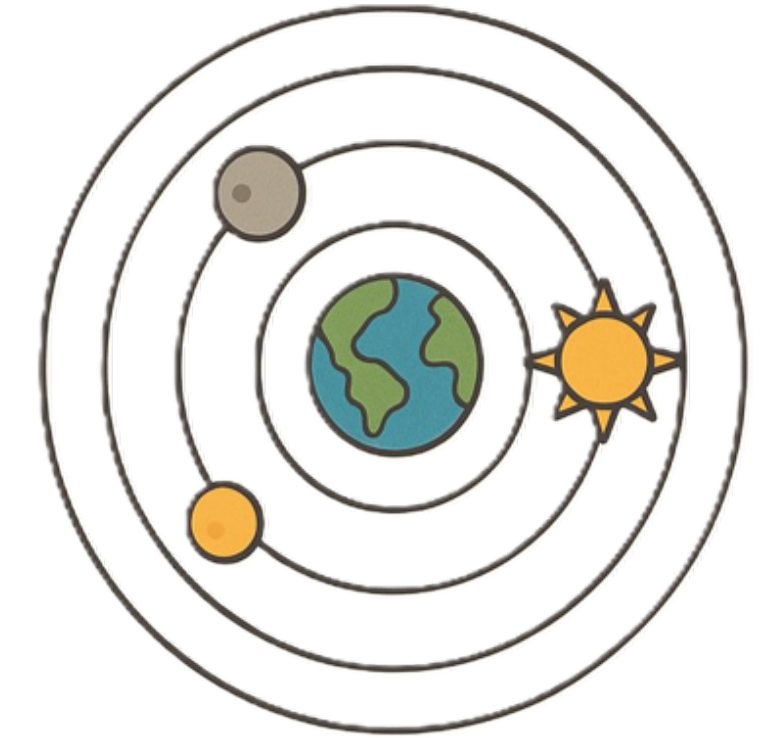


Motivation



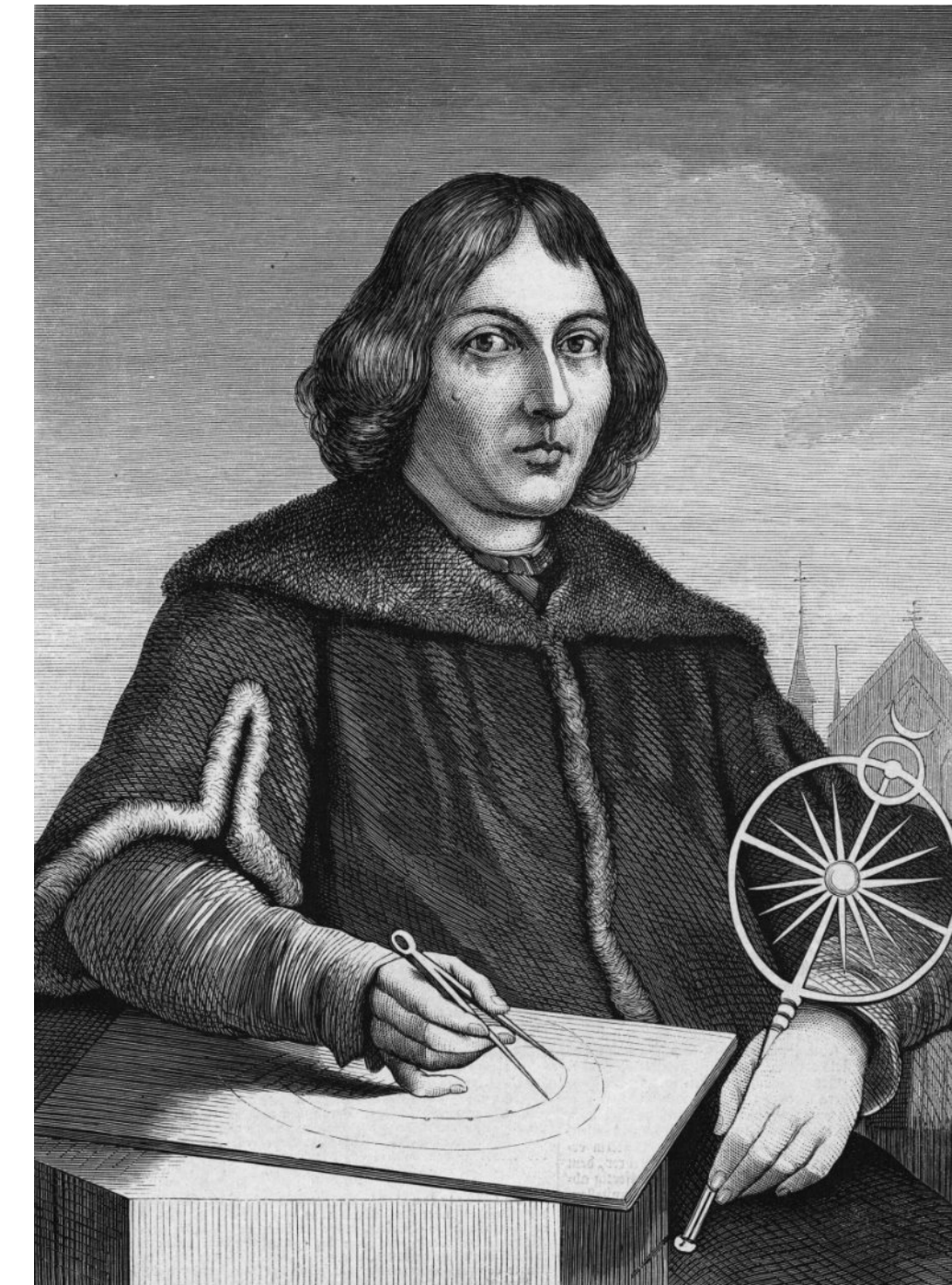
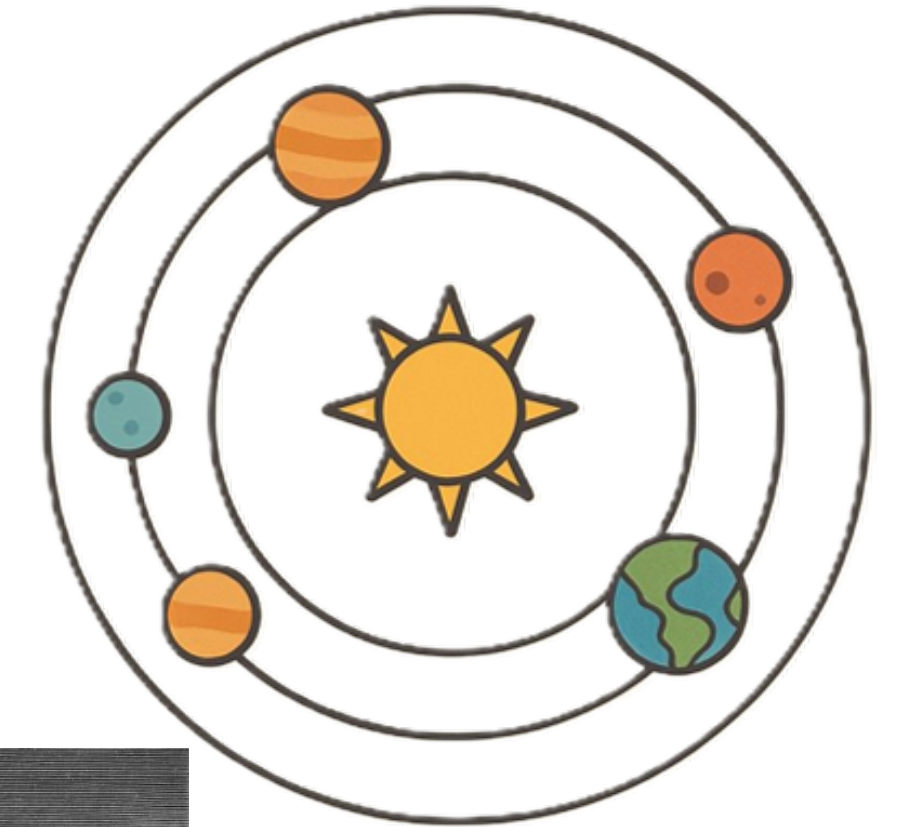
Compact World Models

- Clearly, a wrong **overly-complicated** model of the world...
- Yet, it was the dominant theory for 1500 years...



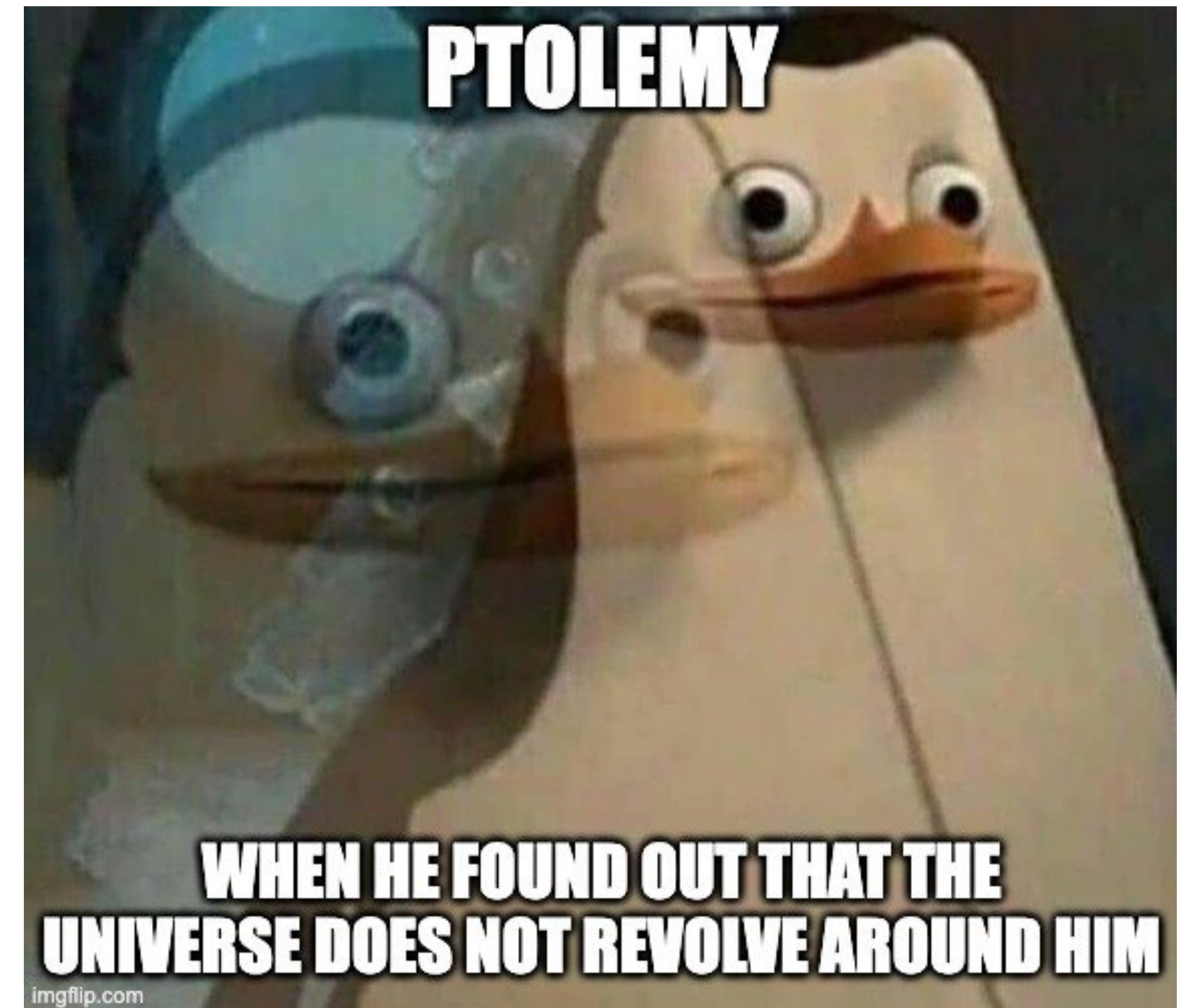
Compact World Models

- Copernicus proposed the Heliocentric Model
- **Simpler, more compact** model of the world



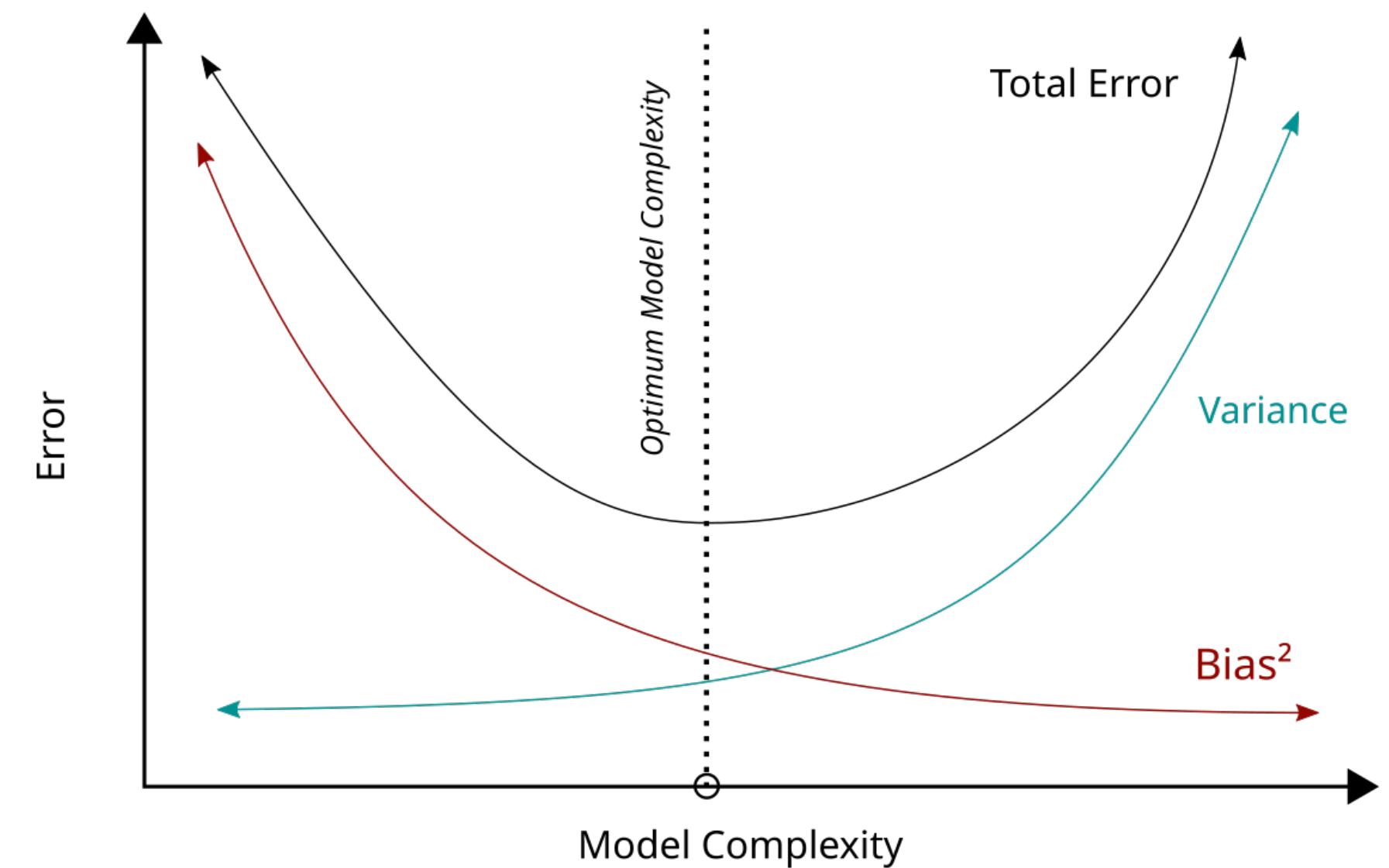
Compact World Models

- Copernicus proposed the Heliocentric Model
- **Simpler, more compact** model of the world
- Can **generalize** to explain observations from other planets' viewpoints



Compact World Models

- NextLat teaches transformer to form **compact world models!**



- **Multi-token prediction** has become a staple in open-source model pretraining

Xiaomi MIMO



MiMo-V2-Flash Technical Report

LLM-Core Xiaomi



2025-12-25

NVIDIA Nemotron 3: Efficient and Open Intelligence

NVIDIA



Qwen Studio

More ▾

EN ▾

Download

Try Qwen Studio ↗



Qwen3-Next: Towards Ultimate Training & Inference Efficiency

2025/09/10 · 54 minute · 10813 words · QwenTeam | Translations: 简体中文



The Keyword

Home

Innovation & AI ▾

Products & platforms ▾

Company news ▾

Feed



Subscribe

Accelerating Gemma 4: faster inference with multi-token prediction drafters

May 05, 2026
4 min read

By using Multi-Token Prediction (MTP) drafters, Gemma 4 models reduce latency bottlenecks and achieve improved responsiveness for developers.

- Reasons for **MTP**
 - Self-speculative decoding: speeds inference up to 3x
 - Better benchmark performance and post-training
 - Data-efficiency

Motivation



Next-Latent Prediction

Next-Token Prediction

Multi-Token Prediction

- Pretrain 1.3 billion parameter models on FineWeb-Edu dataset

The FineWeb Datasets: Decanting the Web for the Finest Text Data at Scale

**Guilherme Penedo Hynek Kydlíček Loubna Ben allal Anton Lozhkov
Margaret Mitchell Colin Raffel Leandro Von Werra Thomas Wolf**
🤗 Hugging Face

Language Modeling



Model	FW-Edu ppl ↓	Wiki. ppl ↓	LAMB. ppl ↓	LAMB. acc ↑	PIQA acc ↑	HellaS. acc ↑	Wino. acc ↑	ARC-e acc ↑	ARC-c acc ↑	SIQA acc ↑	SciQ acc ↑	Avg.
GPT	10.52	17.93	20.26	42.07	73.45	58.79	<u>60.46</u>	68.18	39.16	42.32	86.10	<u>58.82</u>
JTP (d=1)	11.08	19.28	21.88	41.35	74.92	57.43	58.64	68.73	39.25	42.99	<u>87.30</u>	58.83
JTP (d=2)	11.18	19.60	22.11	41.37	73.34	56.84	59.98	68.86	38.57	43.35	86.70	58.63
MTP (d=1)	10.90	18.82	20.23	41.26	<u>74.32</u>	58.05	60.54	68.52	38.91	42.84	85.40	58.76
MTP (d=2)	11.00	18.61	<u>18.34</u>	<u>43.43</u>	72.80	57.92	59.35	68.35	39.08	41.97	86.60	58.69
NextLat (d=1)	<u>10.83</u>	<u>18.39</u>	19.77	41.08	73.07	<u>58.35</u>	59.27	<u>69.65</u>	<u>39.68</u>	<u>43.24</u>	86.00	58.79
NextLat (d=2)	10.88	18.44	17.83	43.86	73.61	57.79	59.20	69.74	40.10	41.91	87.50	59.21

Table 2: Downstream language modeling evaluation on 1.3B-parameter models trained on 100B FineWeb-Edu tokens. Best scores are in bold and second-best are underlined.

Best performance in downstream accuracy but small win at small scale... consistent with observations in Gloeckle et al. [2024]

Language Modeling



Model	FW-Edu ppl ↓	Wiki. ppl ↓	LAMB. ppl ↓	LAMB. acc ↑	PIQA acc ↑	HellaS. acc ↑	Wino. acc ↑	ARC-e acc ↑	ARC-c acc ↑	SIQA acc ↑	SciQ acc ↑	Avg.
GPT	10.52	17.93	20.26	42.07	73.45	58.79	<u>60.46</u>	68.18	39.16	42.32	86.10	<u>58.82</u>
JTP (d=1)	11.08	19.28	21.88	41.35	74.92	57.43	58.64	68.73	39.25	42.99	<u>87.30</u>	58.83
JTP (d=2)	11.18	19.60	22.11	41.37	73.34	56.84	59.98	68.86	38.57	43.35	86.70	58.63
MTP (d=1)	10.90	18.82	20.23	41.26	<u>74.32</u>	58.05	60.54	68.52	38.91	42.84	85.40	58.76
MTP (d=2)	11.00	18.61	<u>18.34</u>	<u>43.43</u>	72.80	57.92	59.35	68.35	39.08	41.97	86.60	58.69
NextLat (d=1)	<u>10.83</u>	<u>18.39</u>	19.77	41.08	73.07	<u>58.35</u>	59.27	<u>69.65</u>	<u>39.68</u>	<u>43.24</u>	86.00	58.79
NextLat (d=2)	10.88	18.44	17.83	43.86	73.61	57.79	59.20	69.74	40.10	41.91	87.50	59.21

Table 2: Downstream language modeling evaluation on 1.3B-parameter models trained on 100B FineWeb-Edu tokens. Best scores are in bold and second-best are underlined.

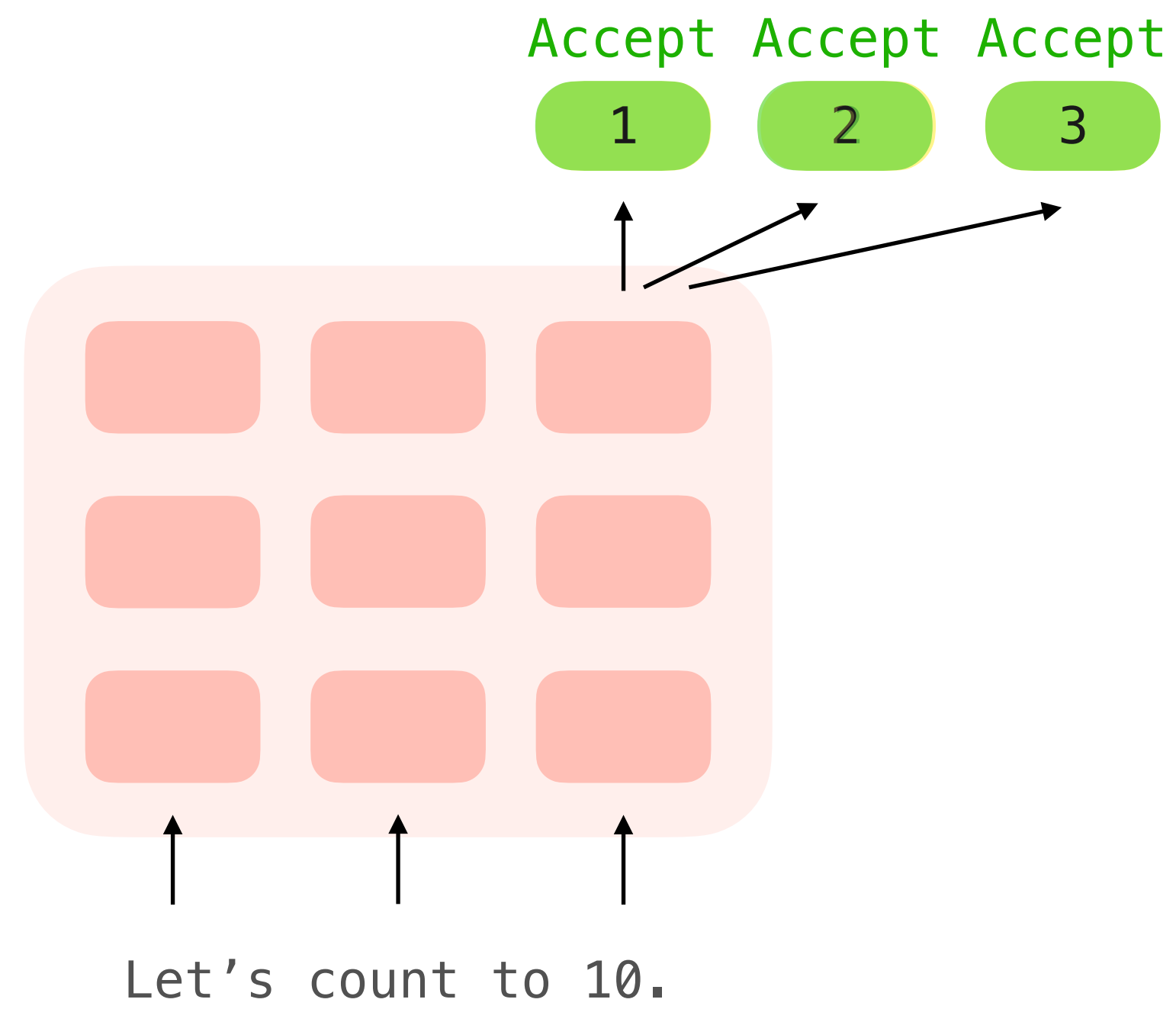
NextLat preserves next-token prediction performance better than MTP methods!

Speculative Decoding



Multi-Token Prediction

fixed speculative horizon, e.g. 3 tokens

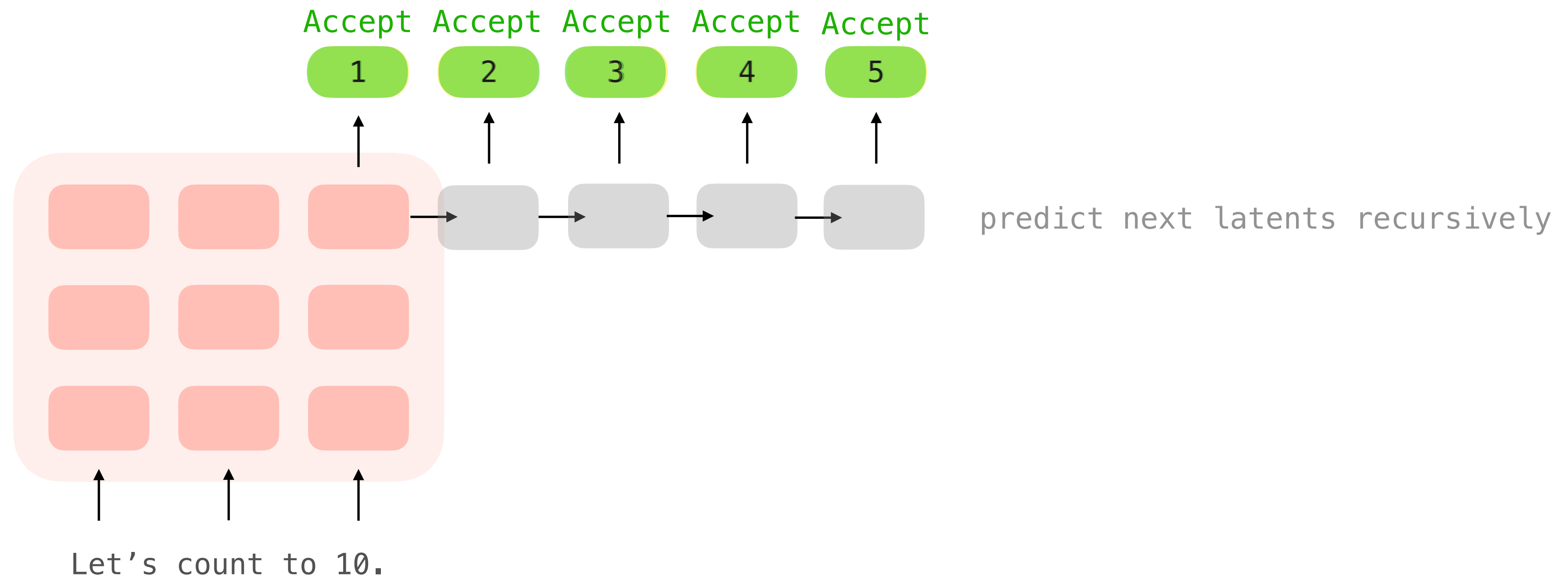


Speculative Decoding



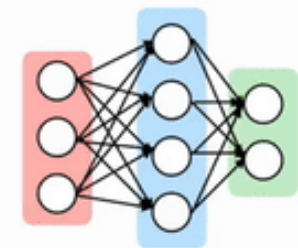
Next-Latent Prediction

variable-length speculative horizon

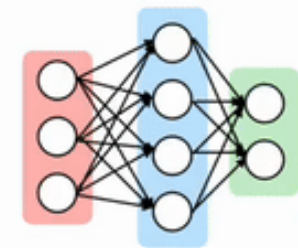


Speculative Decoding

Standard Self-Speculative Decoding (w/ Multi-Token Prediction)



Variable-Length Self-Speculative Decoding (w/ Next-Latent Prediction)

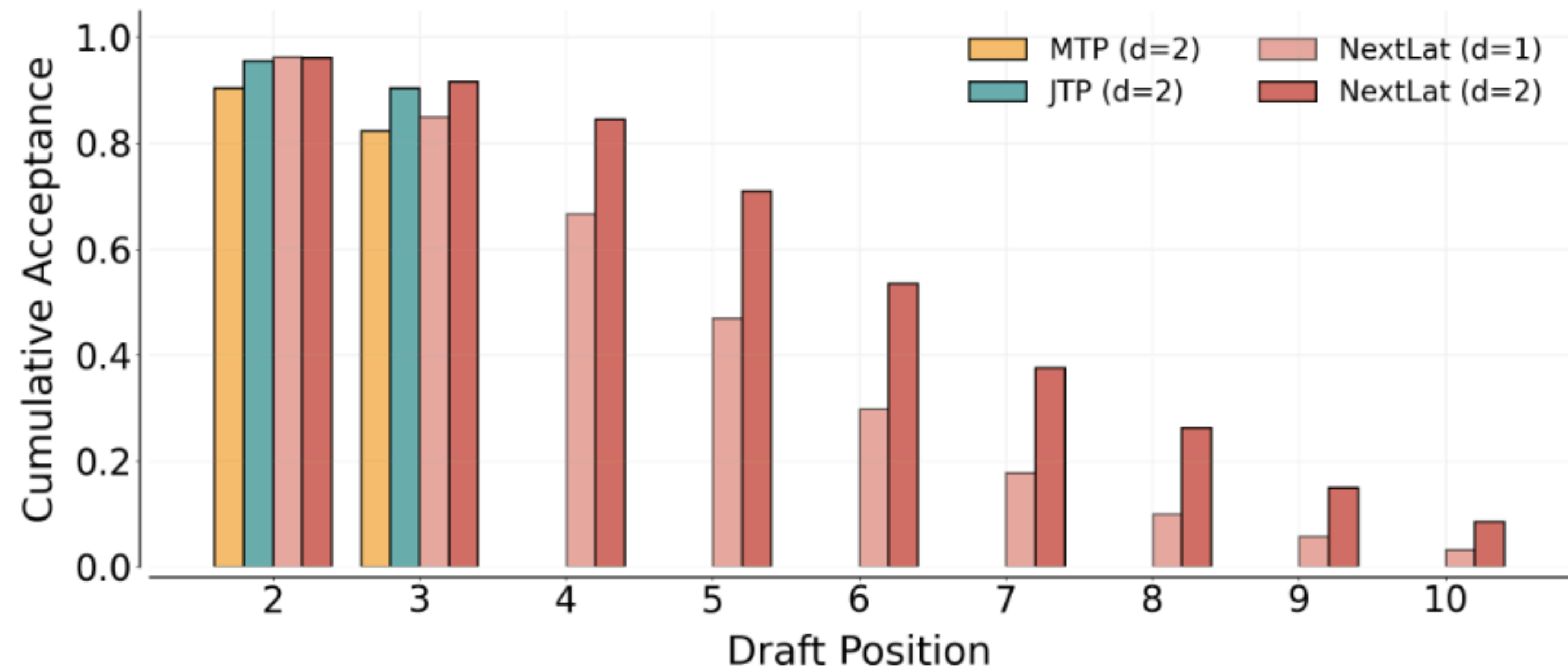
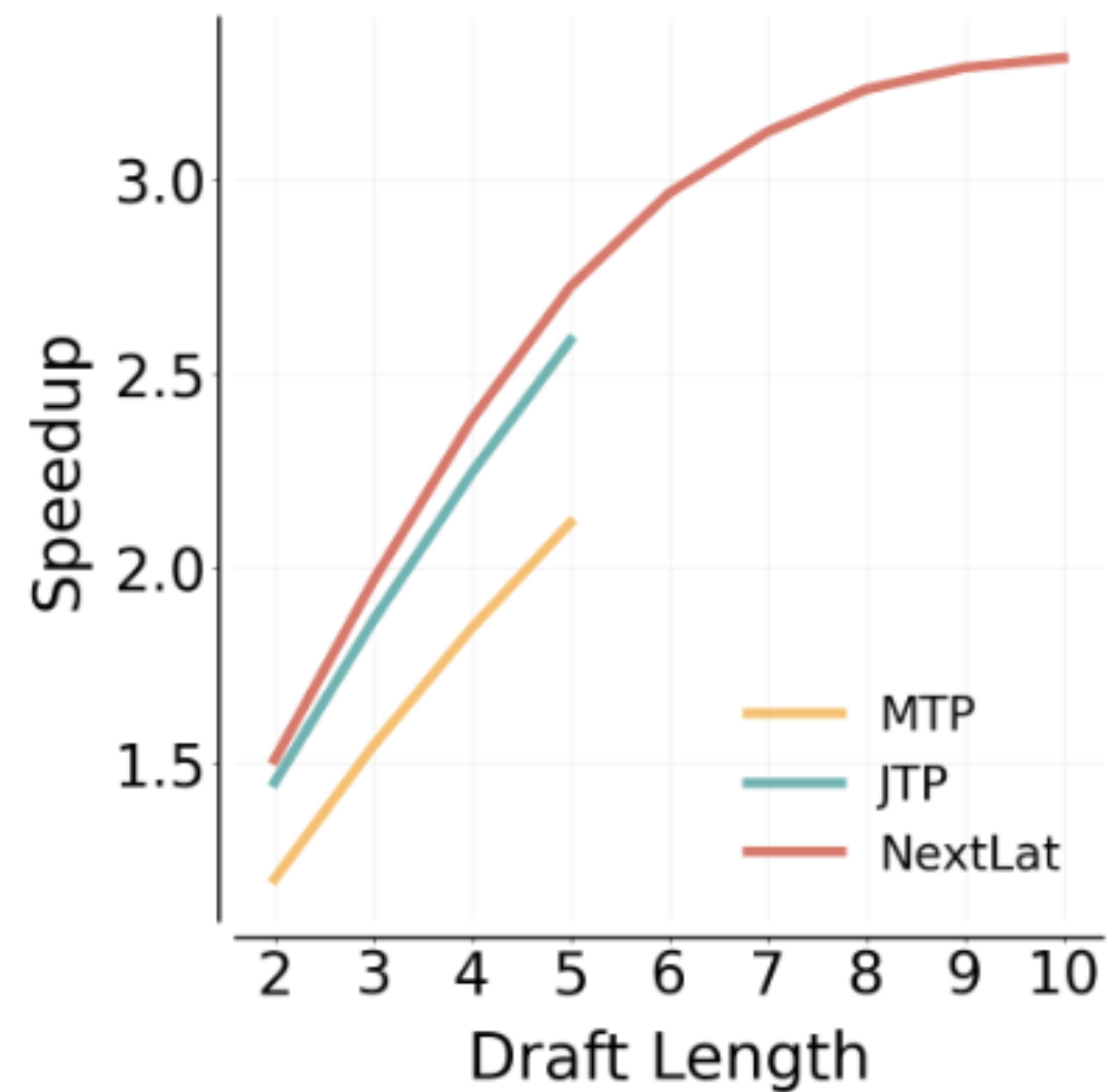


 Draft Tokens

 Accepted Tokens

 Rejected Tokens

Speculative Decoding



Much faster inference (up to 3.3x) with speculative decoding

Speculative Decoding

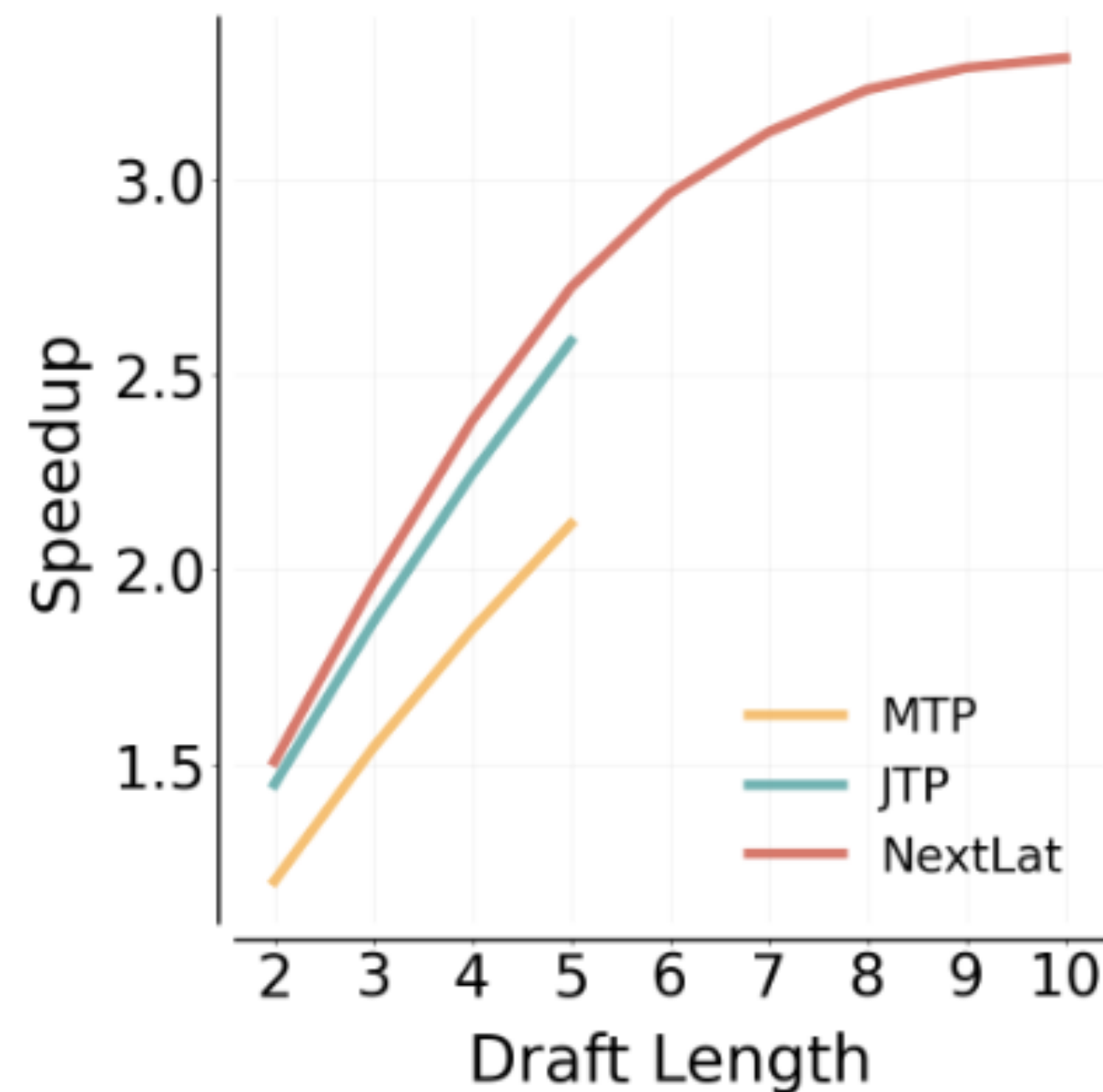


Model	Wikipedia		Books		Code		Math	
	Speedup	Accepted Tokens	Speedup	Accepted Tokens	Speedup	Accepted Tokens	Speedup	Accepted Tokens
JTP (d=1)	1.46	0.96	1.47	0.97	1.47	0.98	1.46	0.97
JTP (d=2)	1.88	1.84	1.90	1.89	1.88	1.85	1.89	1.86
MTP (d=1)	1.38	0.91	1.39	0.95	1.40	0.97	1.39	0.95
MTP (d=2)	1.68	1.72	1.72	1.83	1.75	1.91	1.72	1.84
NextLat (d=1)	<u>2.68</u>	<u>3.52</u>	<u>2.72</u>	<u>3.64</u>	<u>2.29</u>	<u>2.66</u>	<u>2.30</u>	<u>2.72</u>
NextLat (d=2)	3.21	4.59	3.32	4.86	2.38	2.83	2.87	3.94

Table 3: Relative speedup and average accepted tokens per drafting steps over diverse domains. Note that “Accepted Tokens” excludes the next-token prediction which is always accepted.

Much faster inference (up to 3.3x) with speculative decoding

Speculative Decoding



	Next-Token	Multi-Token	NextLat
Training parameters ($d = 1, 2, 8$)	1.32B	1.37B / 1.42B / 1.72B	1.40B
Training steps/second ($d = 1, 2, 8$)	3.09	2.80 / 2.58 / 1.70	<u>3.09 / 2.79 / 1.73</u>

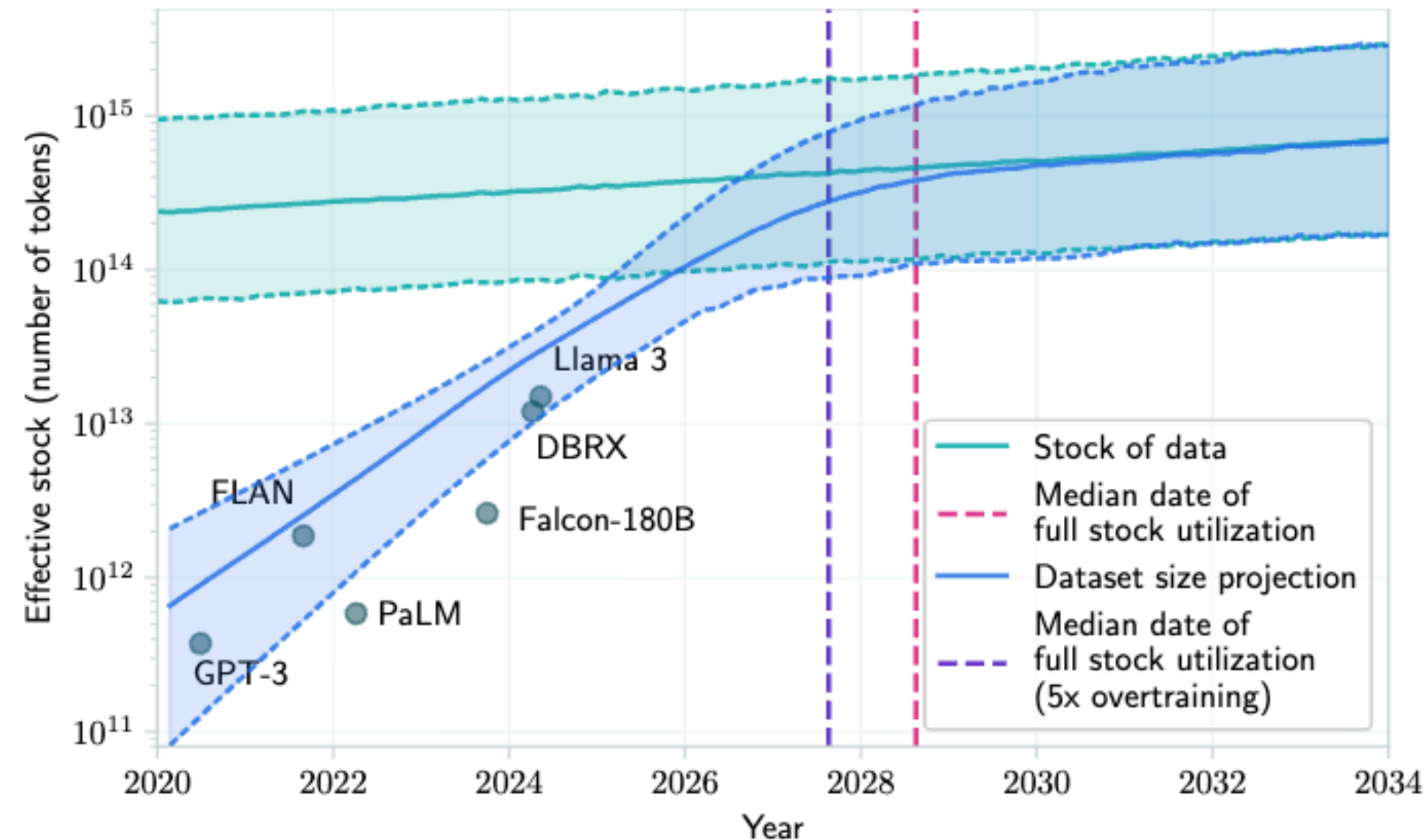
Table 5: Comparison of training speed and parameter.

Train with just next-latent prediction, no need multiple heads for multi-token prediction!

Data-Efficiency



We will run out of web data for training models by 2028...

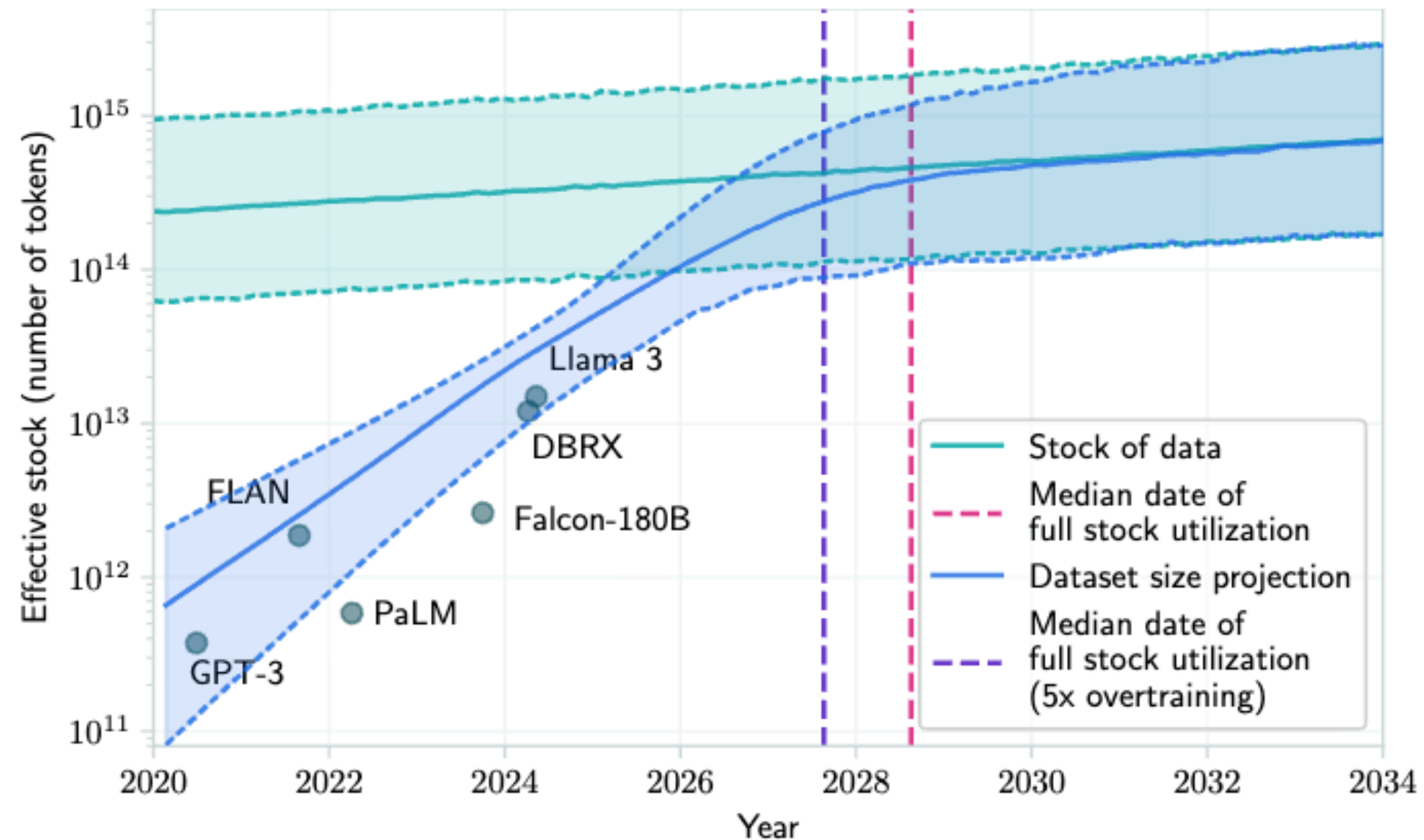


[8] Villalobos et al. Will we run out of data? Limits of LLM scaling based on human-generated data. ICML 2024.

Data-Efficiency

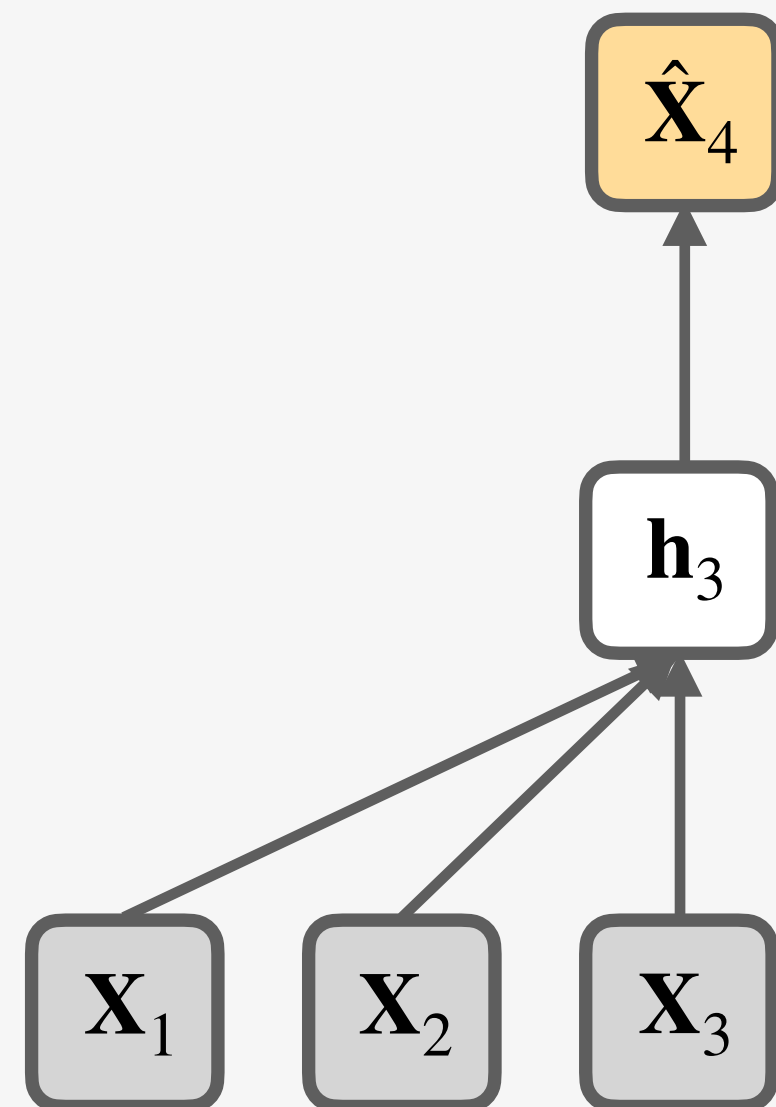


It's time to think about methods that extract more gradients from each sequences...

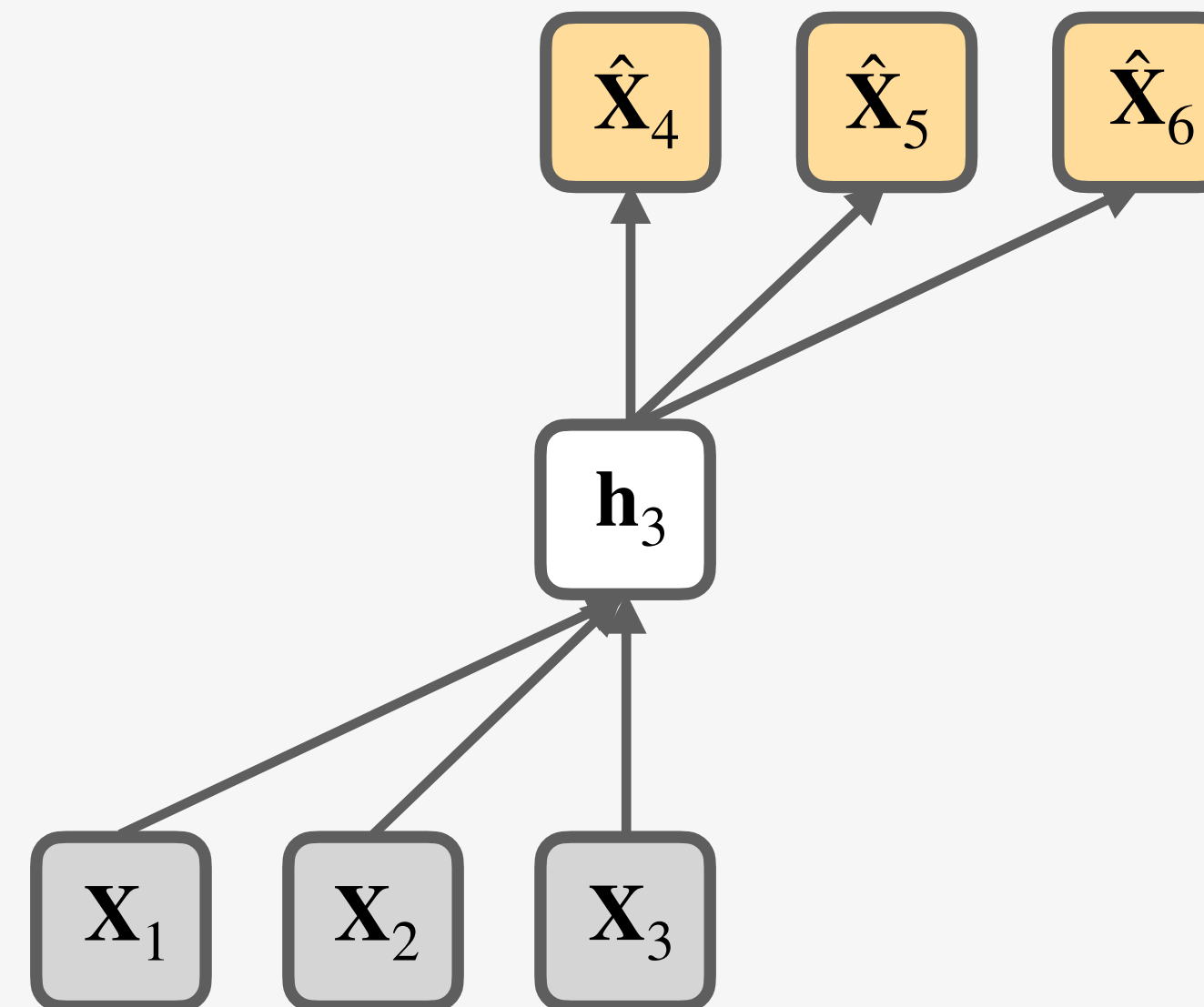


[8] Villalobos et al. Will we run out of data? Limits of LLM scaling based on human-generated data. ICML 2024.

Next-Token Prediction

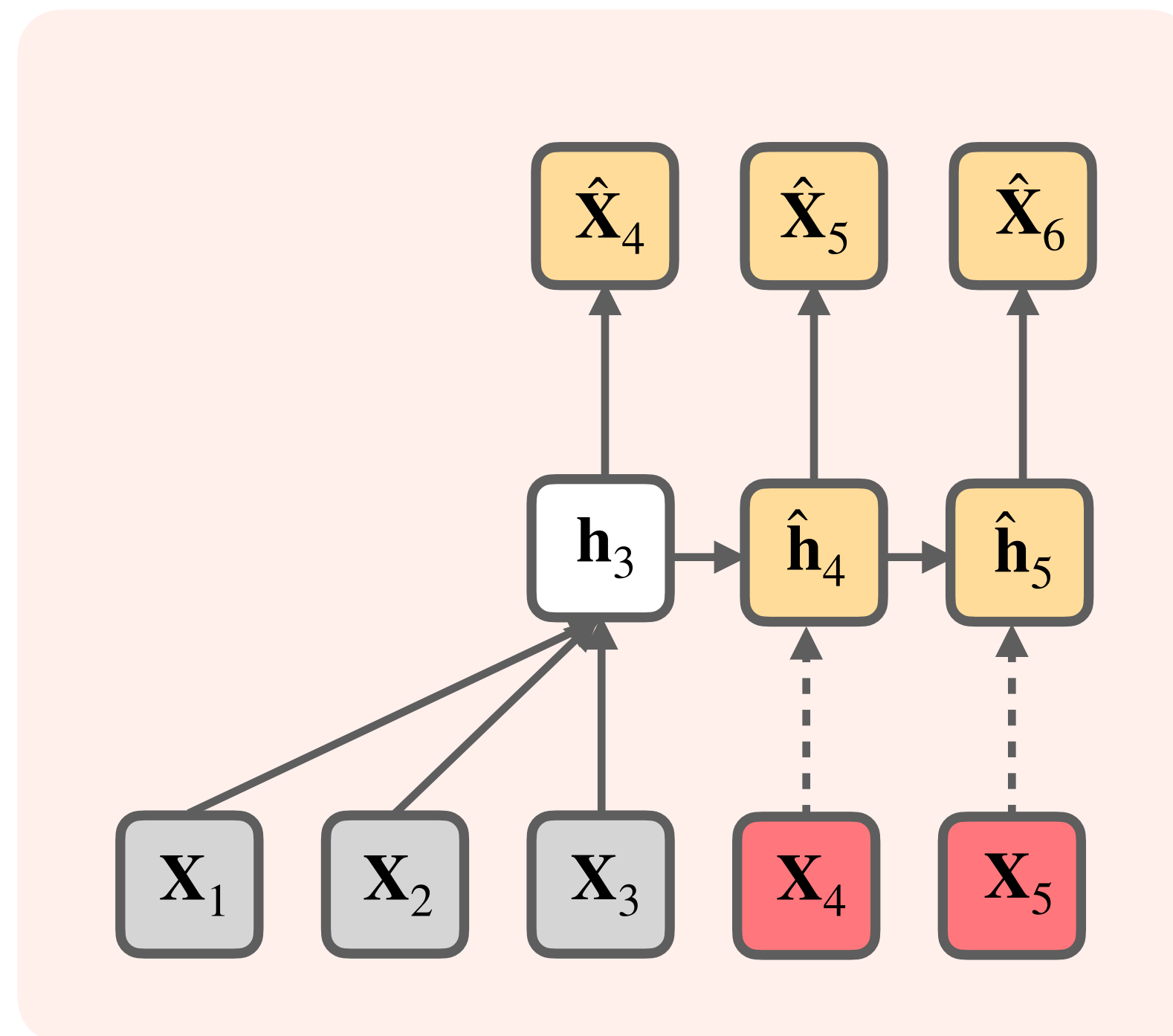


Multi-Token Prediction



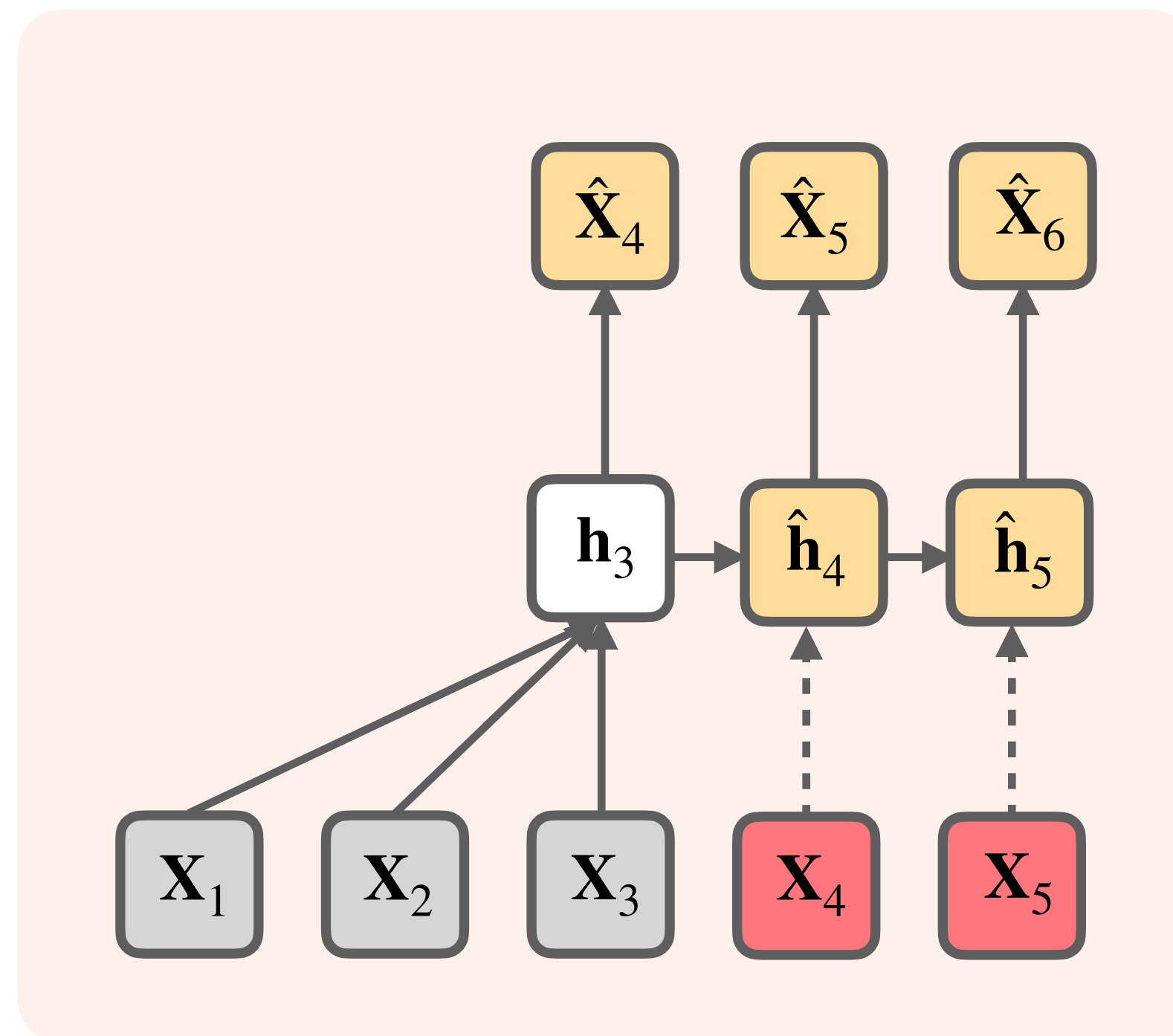
Learning signal for h_3 is only in terms of the next token, or multiple tokens....

Next-Latent Prediction



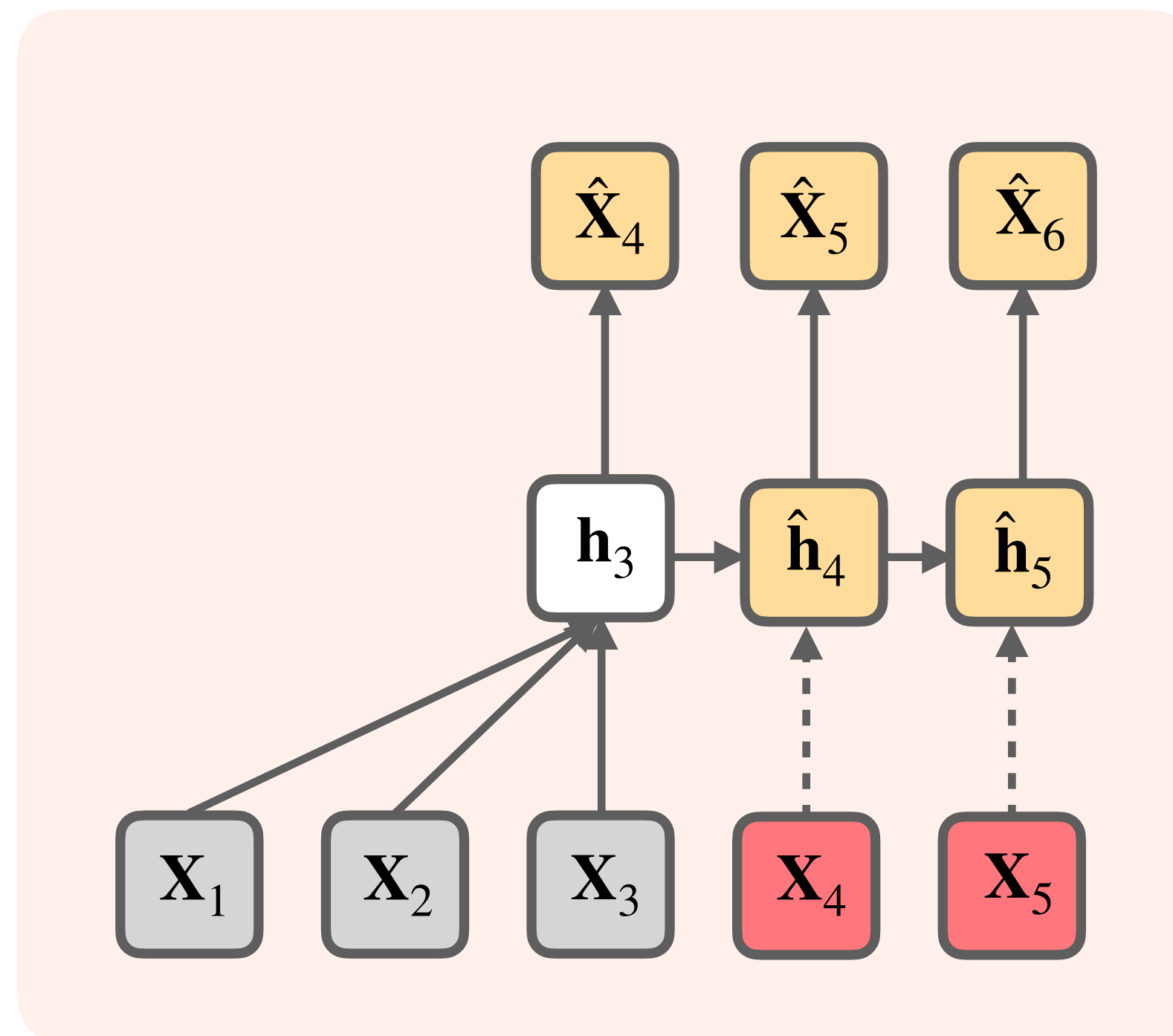
h_3 is trained to predict h_4 which parameterises *distribution* over X_5

Next-Latent Prediction



Moreover, h_4 is trained to predict h_5 which is trained to predict h_6 ...

Next-Latent Prediction



This means that h_3 receives dense learning signal about the whole future h_4, h_5, h_6, \dots

RNN without recurrence

- There are certain problems inexpressible by transformers but expressible by RNNs, such as state tracking [9]...

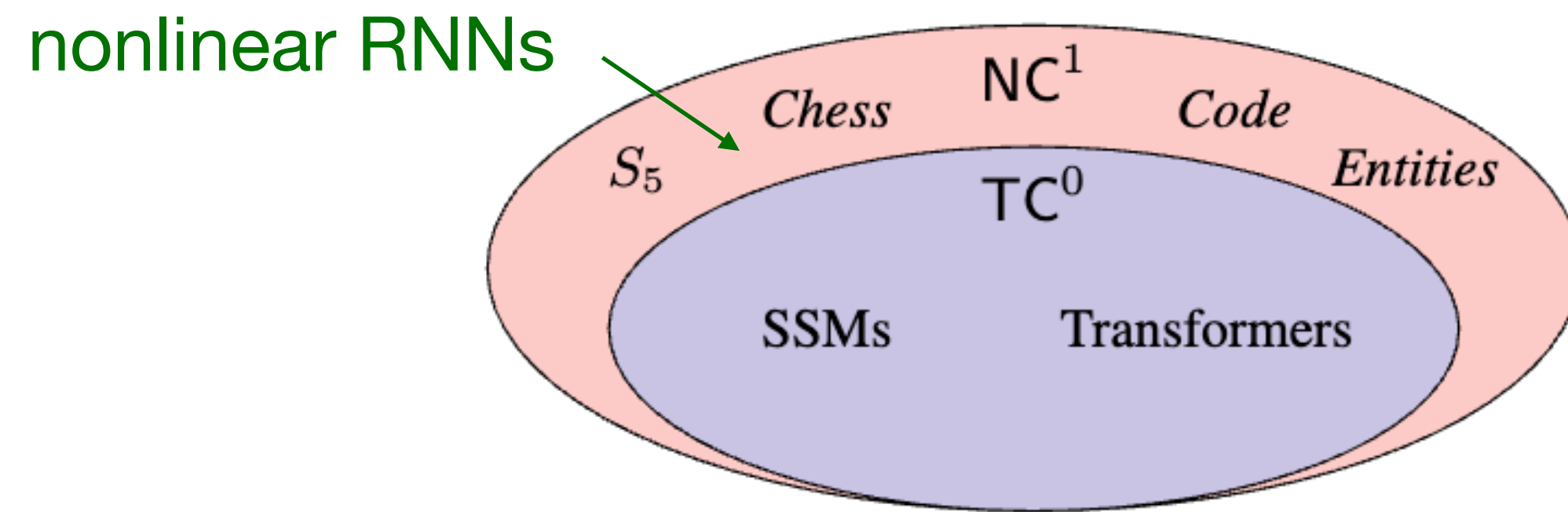
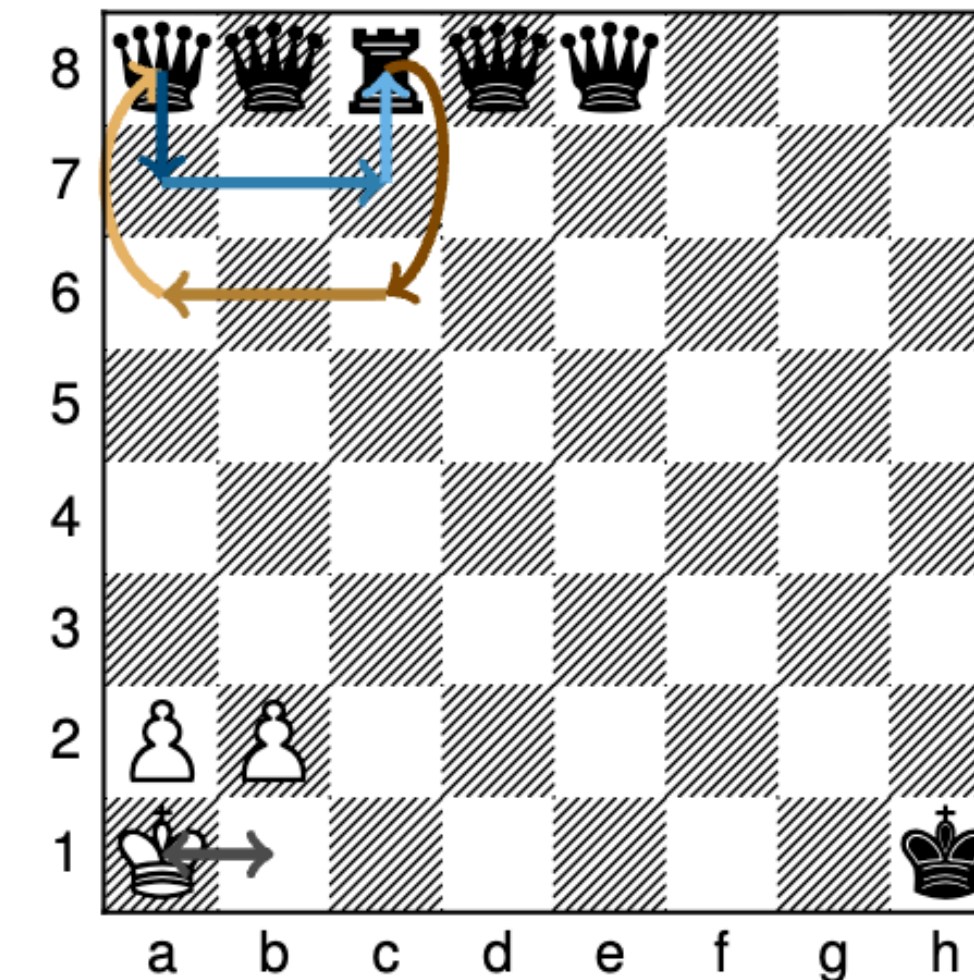


Figure 2: Complexity hierarchy within NC^1 . Transformers can only recognize languages within TC^0 (Merrill & Sabharwal, 2023a), and we show the same for SSMs (Theorems 4.2 and 4.4). Thus, both architectures cannot express the “hard state tracking” captured by NC^1 -complete problems like S_5 , which *can* be straightforwardly expressed by RNNs. The figure assumes the widely held conjecture $TC^0 \neq NC^1$.



```
x = [0, 0, 1, 0, 0]
x[1], x[3] = x[3], x[1] # Swap 1, 3
```

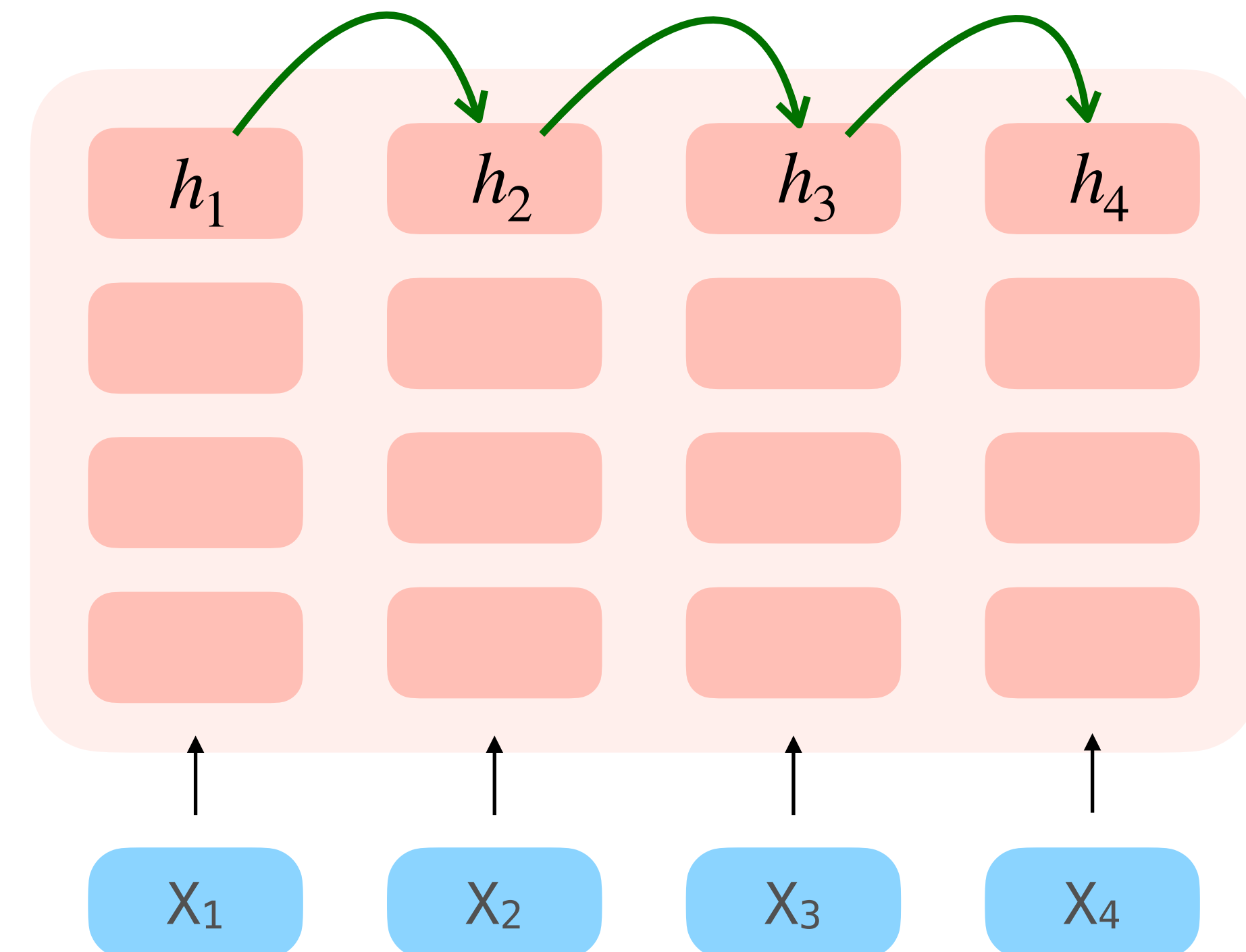
Alice, Bob, Carl, Dan, and Emma each have a coin. All are dimes except Carl's. Alice and Carl trade coins.

RNN without recurrence



- NextLat trains a nonlinear RNN without recurrence

NextLat: $p(h_2 | h_1, X_2)$ $p(h_3 | h_2, X_3)$ $p(h_4 | h_3, X_4)$



transformer forecasts hidden states in parallel

next-latent predictor (RNN) trains on one-step prediction

fully parallel across time!

RNN without recurrence



- We train NextLat on A_5 , a word problem inexpressible by constant-depth, fixed-precision transformers

transformer cannot generalize beyond training length of 12 tokens

BUT, the RNN trained using NextLat can **generalize** far beyond!

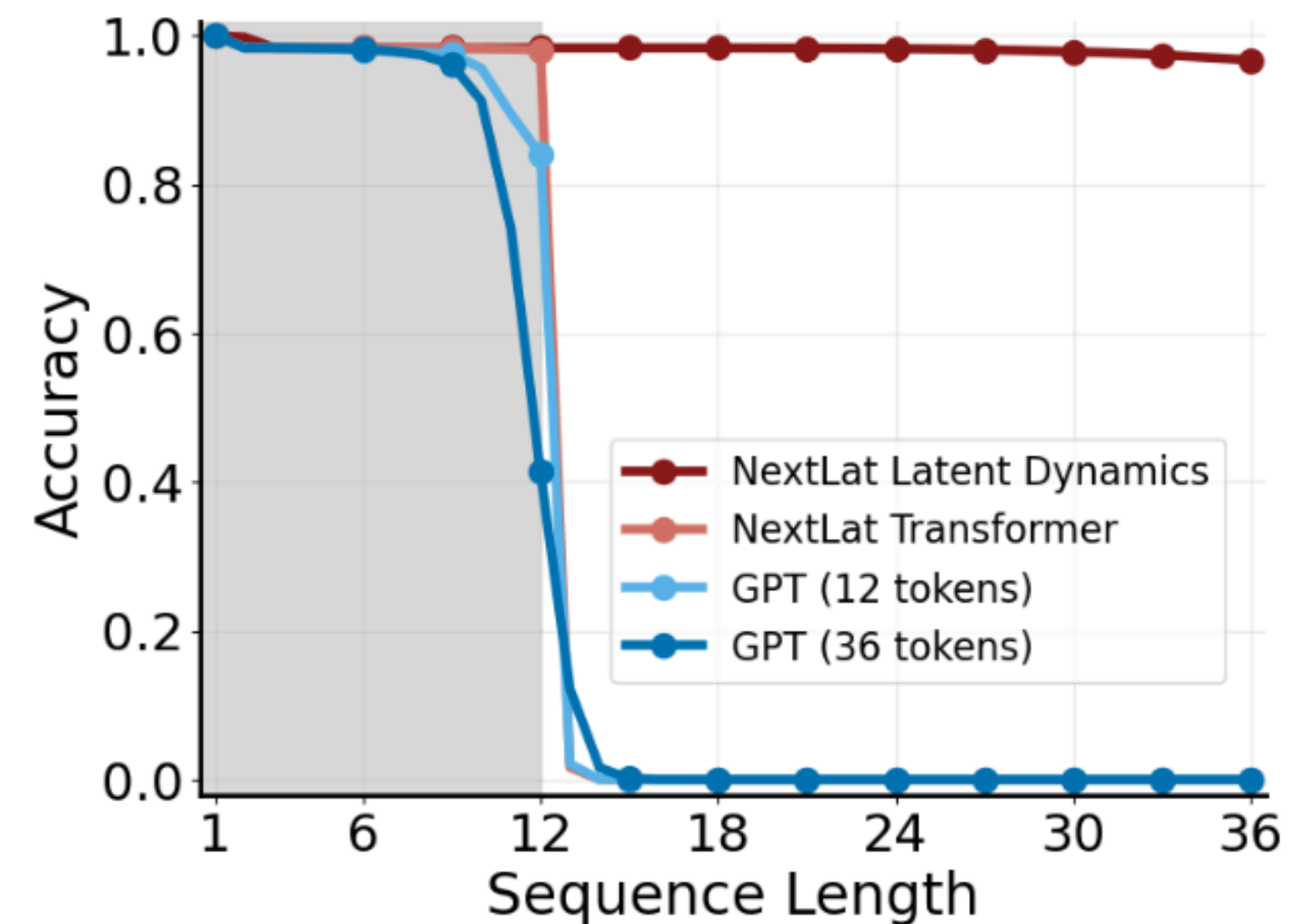


Figure 11: Length Generalization for A_5 word problem.

RNN without recurrence



- We train NextLat on A_5 , a word problem inexpressible by constant-depth, fixed-precision transformers

Question: *Can a transformer with $O(\log T)$ depth co-train an RNN that generalizes to sequences of length $\gg T$?*

Can the RNN learn NC^1 computations, even though it is supervised using representations from a transformer in TC^0 ?

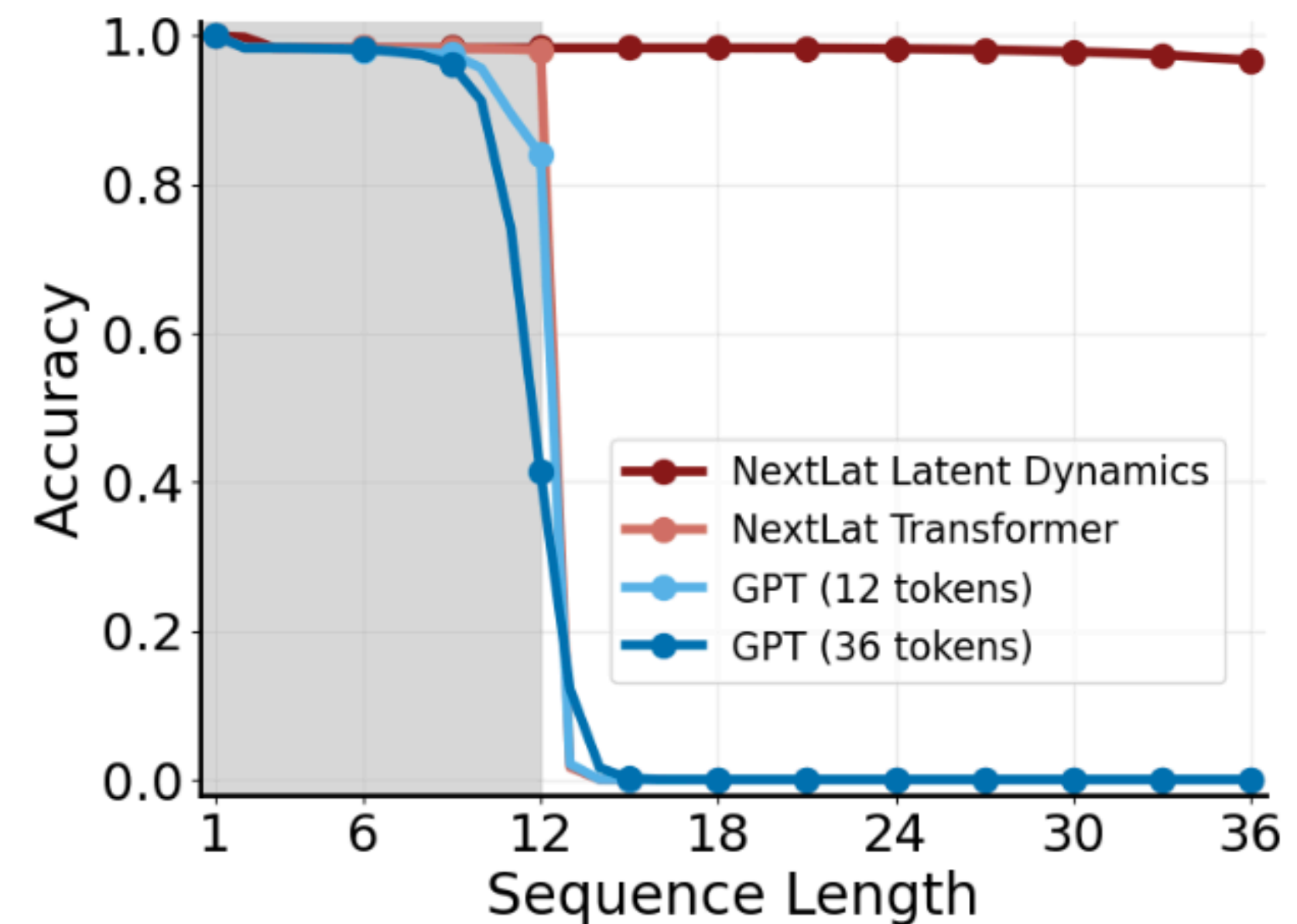
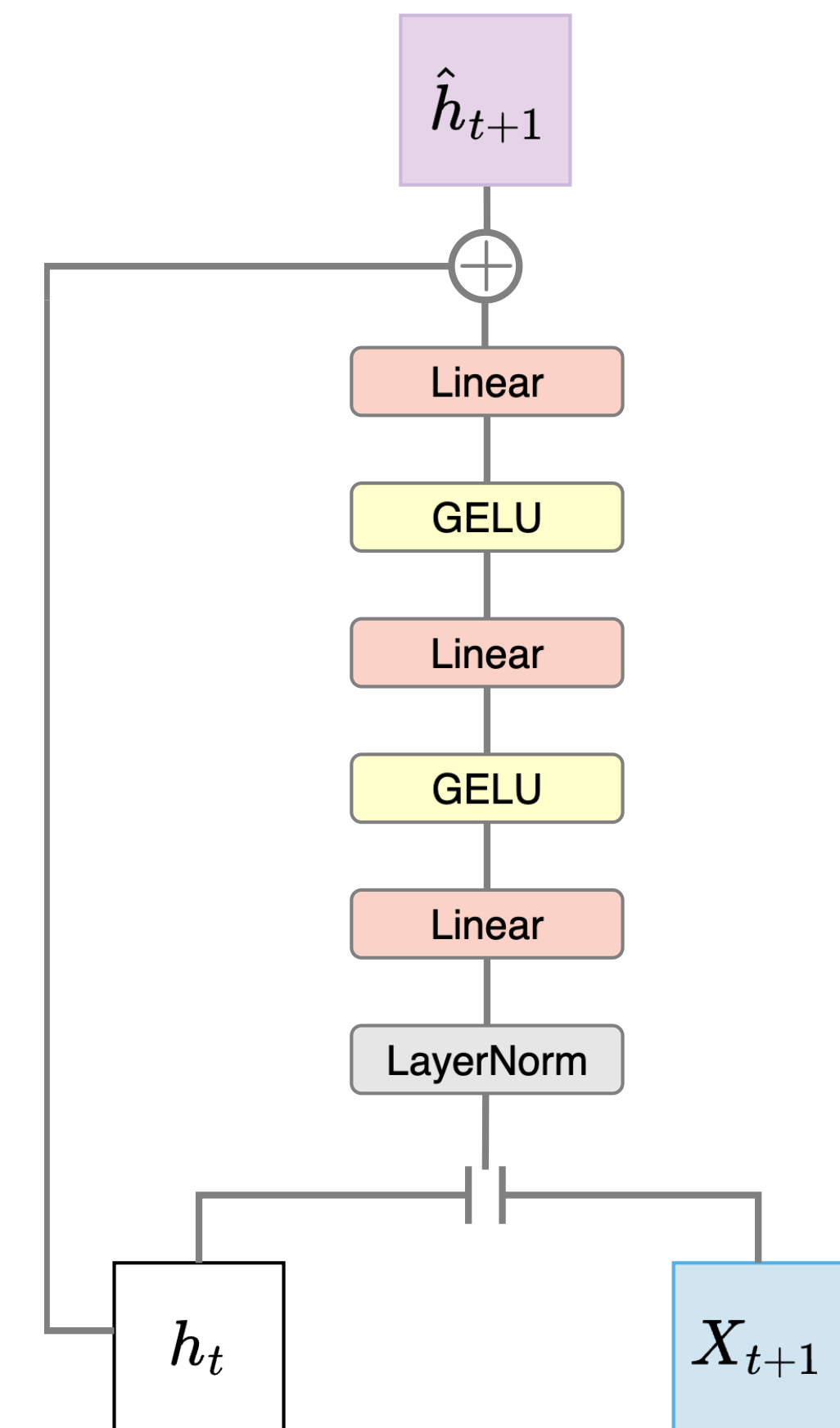


Figure 11: Length Generalization for A_5 word problem.

What's Next for NextLat?



- We got impressive results with just a **simple MLP** for the next-latent predictor
- What if we use more expressive architectures and scale up the predictor size?



What's Next for NextLat?



- Scaling up NextLat

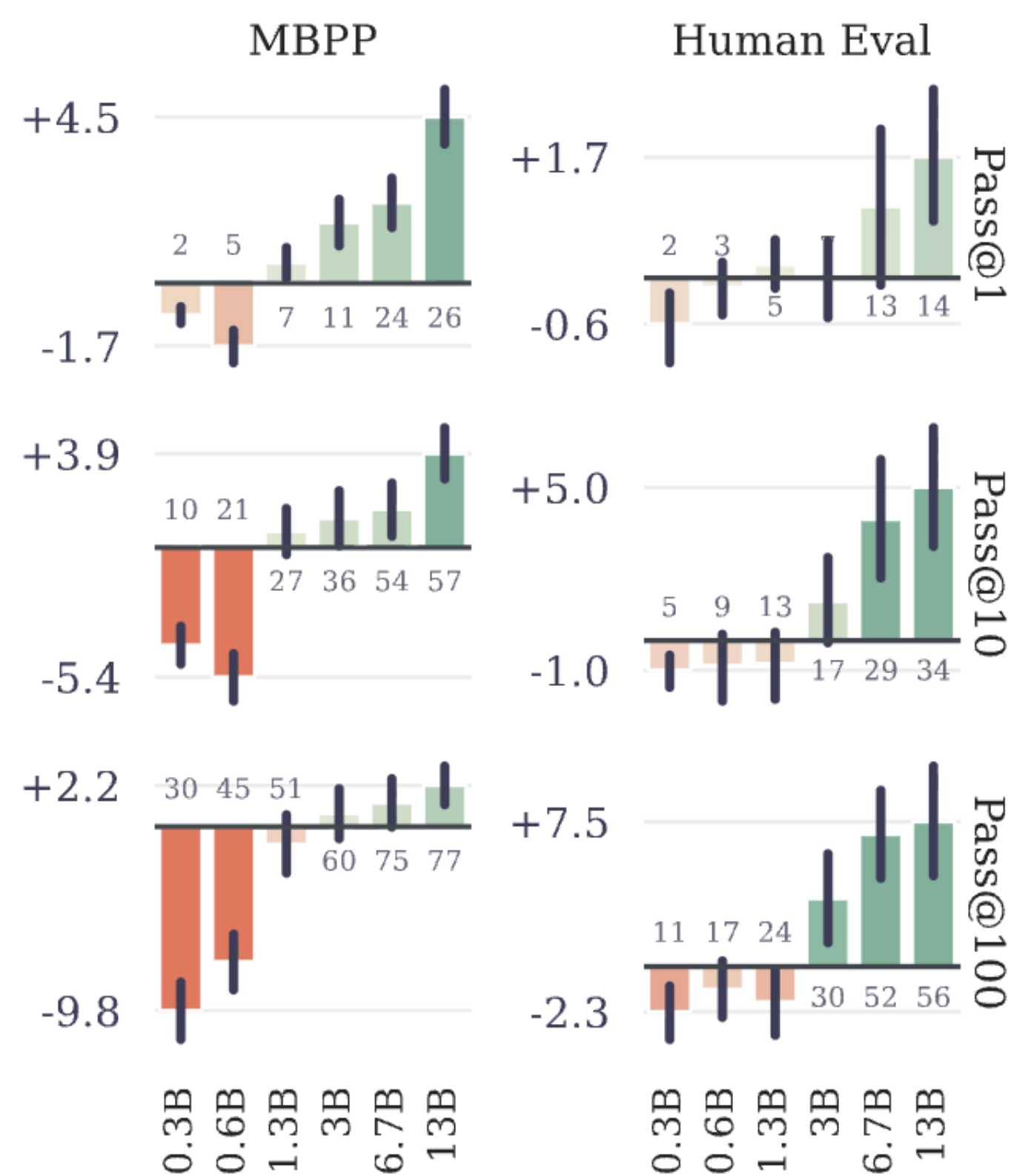


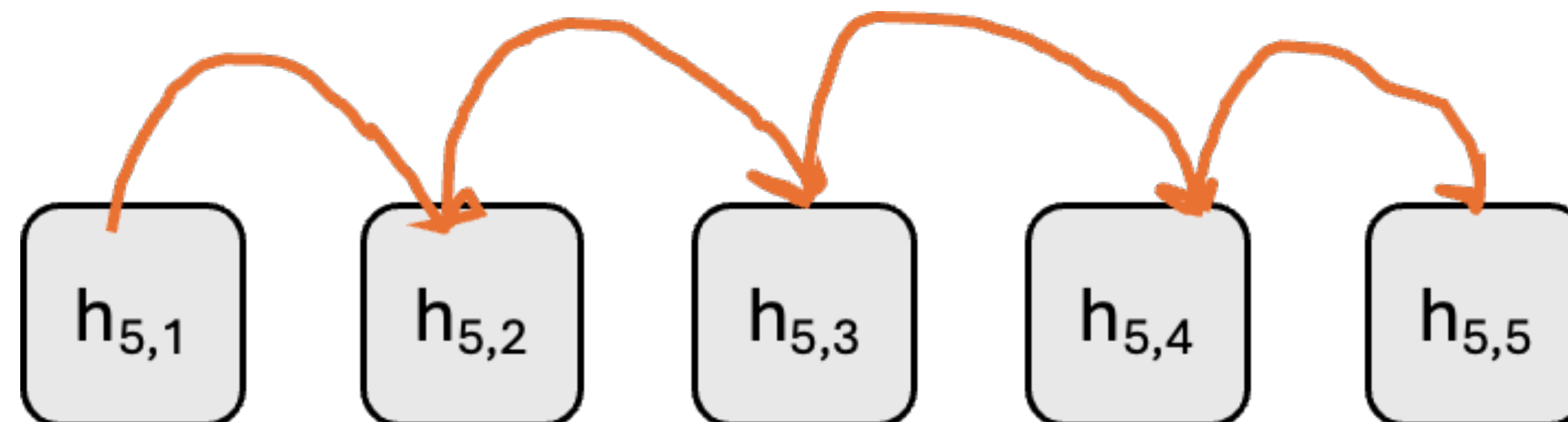
Figure 3: **Results of n -token prediction models on MBPP by model size.** We train models of six sizes in the range or 300M to 13B total parameters on code, and evaluate pass@1,10,100 on the MBPP (Austin et al., 2021) and HumanEval (Chen et al., 2021) benchmark with 1000 samples. Multi-token prediction models are worse than the baseline for small model sizes, but outperform the baseline at scale. Error bars are confidence intervals of 90% computed with bootstrapping over dataset samples.

What's Next for NextLat?



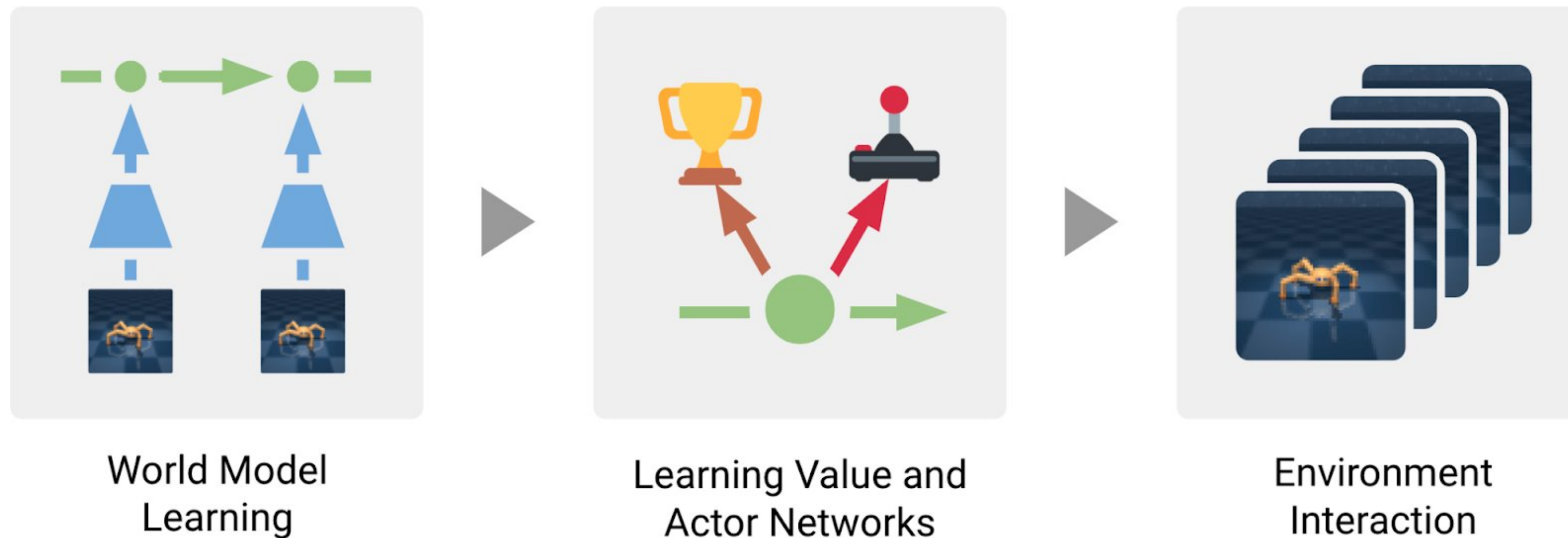
- NextLat might be better for post-training too!
- better for implicit/explicit value learning

very “Bellman-like”/recurrent representations



What's Next for NextLat?

- Learning latent dynamics is already common in other fields, but not explored in language...



World Model
Learning

Learning Value and
Actor Networks

Environment
Interaction

Dreamer algorithm [Hafner et al., 2019]

What's Next for NextLat?

- Learning latent dynamics is already common in other fields, but not explored in language...



Genie algorithm [Bruce et al., 2024]

Conclusion



- NextLat can replace multi-token prediction for pre-training
 - variable-length speculative decoding
 - less parameters, faster training (only need to train with $d=1$)
- NextLat can be a strong pretraining technique for foundation models
 - better representations for reasoning and planning
 - better data-efficiency; extract more gradients from each sequence

Conclusion



Paper:
arxiv.org/abs/2511.05963



Code:
github.com/JaydenTeoh/NextLat

Email: jayden_t@mit.edu

X: @jayden_teoh_