

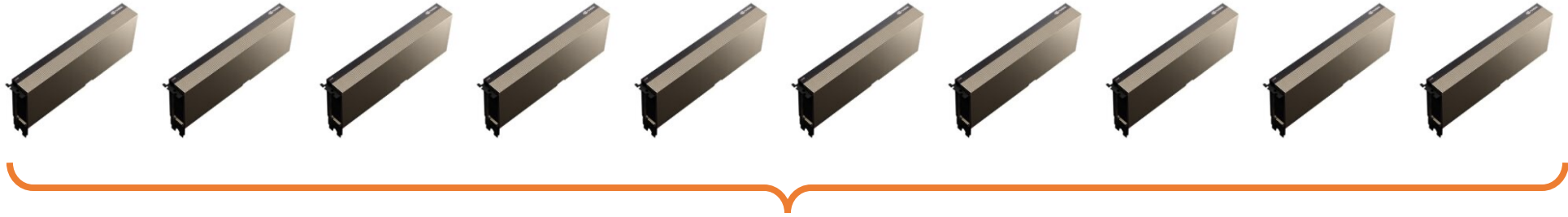
# **15-442/15-642: Machine Learning Systems**

## **LLMs Serving Techniques**

**Tianqi Chen and Zhihao Jia**

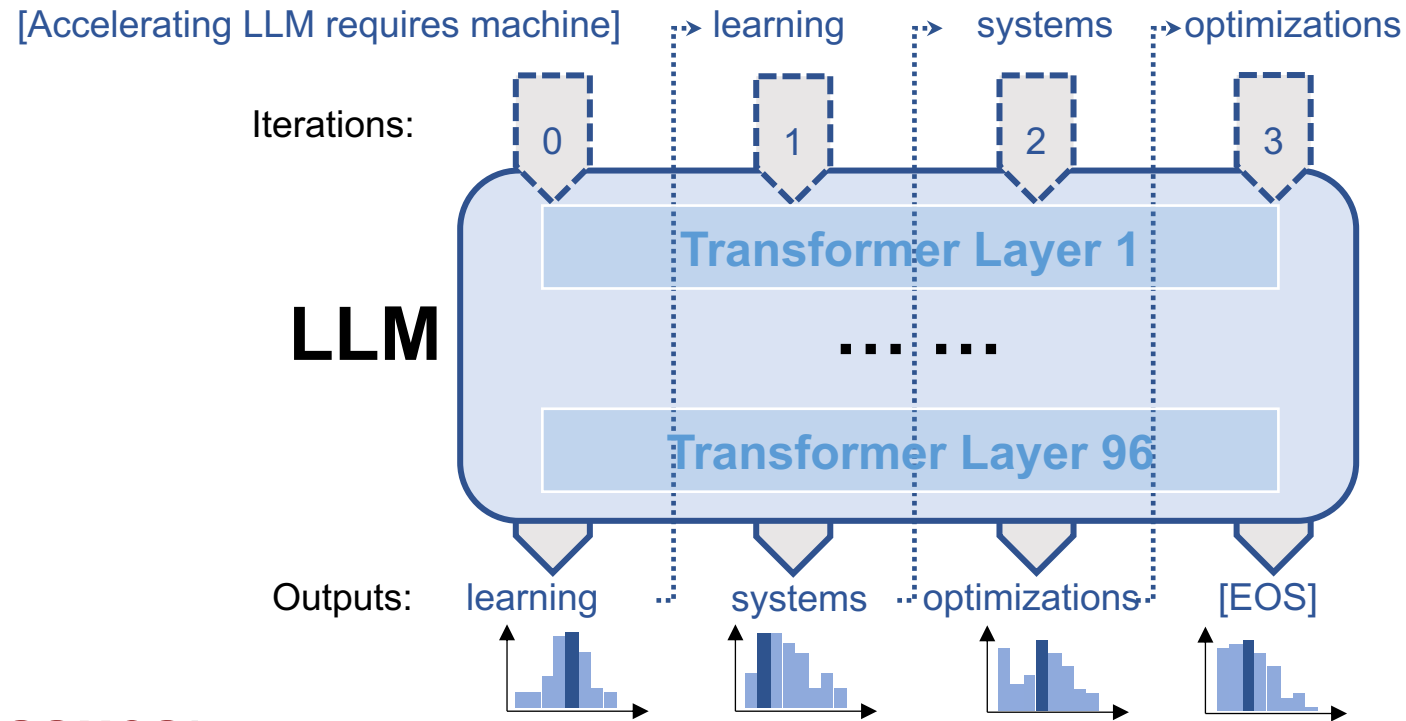
Carnegie Mellon University

# LLMs are Slow and Expensive to Serve



- **At least ten** A100-40GB GPUs to serve 175B GPT-3 in half precision
- Generating 256 tokens takes **~20 seconds**
- Cannot process many requests in parallel
  - Per-request key/value cache takes **3GB GPU memory**

# Recall: Incremental Decoding



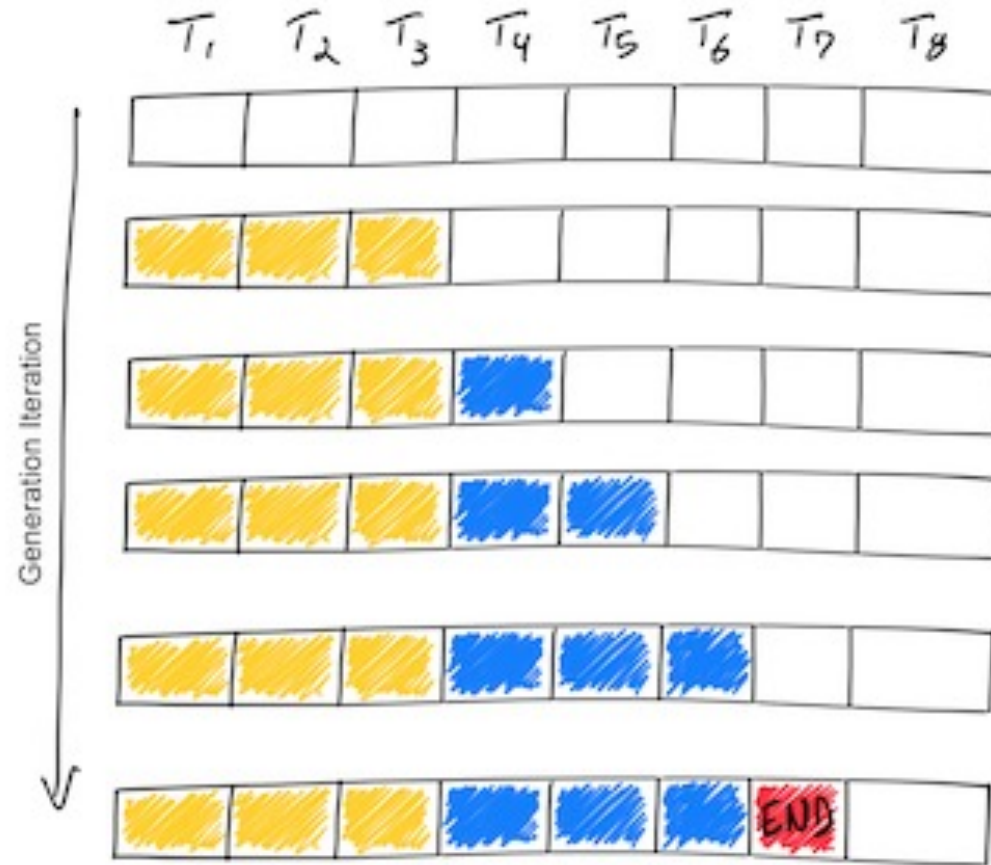
## Main issues:

- Limited degree of parallelism → underutilized GPU resources
- Need all parameters to decode a token → bottlenecked by GPU memory access

# Outline: LLMs Serving Techniques

- **Continuous Batching**
- Speculative Decoding

# LLM Decoding Timeline



# Batching Requests to Improve GPU Performance

$T_1$	$T_2$	$T_3$	$T_4$	$T_5$	$T_6$	$T_7$	$T_8$
$S_1$	$S_1$	$S_1$	$S_1$				
$S_2$	$S_2$	$S_2$					
$S_3$	$S_3$	$S_3$	$S_3$				
$S_4$	$S_4$	$S_4$	$S_4$	$S_4$			

$T_1$	$T_2$	$T_3$	$T_4$	$T_5$	$T_6$	$T_7$	$T_8$
$S_1$	$S_1$	$S_1$	$S_1$	$S_1$	END		
$S_2$	$S_2$	$S_2$	$S_2$	$S_2$	$S_2$	$S_2$	END
$S_3$	$S_3$	$S_3$	$S_3$	END			
$S_4$	$S_4$	$S_4$	$S_4$	$S_4$	$S_4$	END	

Issues with static batching:

- Requests may complete at different iterations
- Idle GPU cycles
- New requests cannot start immediately

# Continuous Batching

$T_1$	$T_2$	$T_3$	$T_4$	$T_5$	$T_6$	$T_7$	$T_8$
$S_1$	$S_1$	$S_1$	$S_1$				
$S_2$	$S_2$	$S_2$					
$S_3$	$S_3$	$S_3$	$S_3$				
$S_4$	$S_4$	$S_4$	$S_4$	$S_4$			

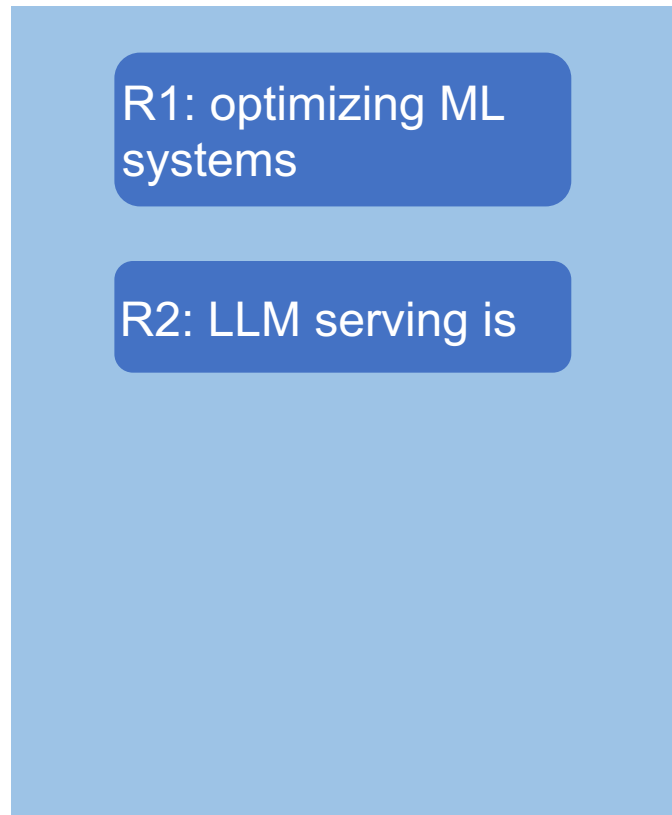
$T_1$	$T_2$	$T_3$	$T_4$	$T_5$	$T_6$	$T_7$	$T_8$
$S_1$	$S_1$	$S_1$	$S_1$	$S_1$	END	$S_6$	$S_6$
$S_2$	$S_2$	$S_2$	$S_2$	$S_2$	$S_2$	$S_2$	END
$S_3$	$S_3$	$S_3$	$S_3$	END	$S_5$	$S_5$	$S_5$
$S_4$	$S_4$	$S_4$	$S_4$	$S_4$	$S_4$	END	$S_7$

Benefits:

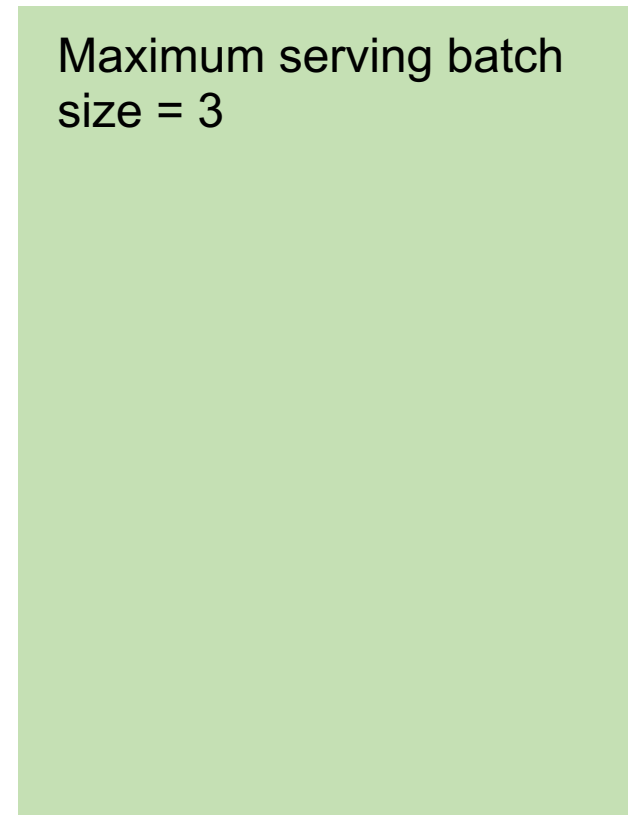
- Higher GPU utilization
- New requests can start immediately

# Continuous Batching Step-by-Step

- Receives two new requests R1 and R2



**Request Pool  
(CPU)**

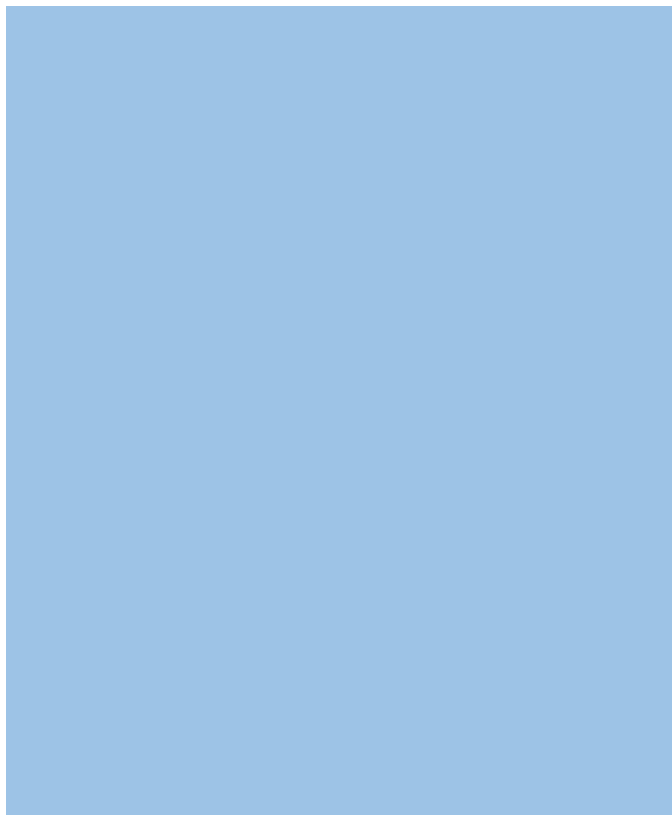


**Execution Engine  
(GPU)**

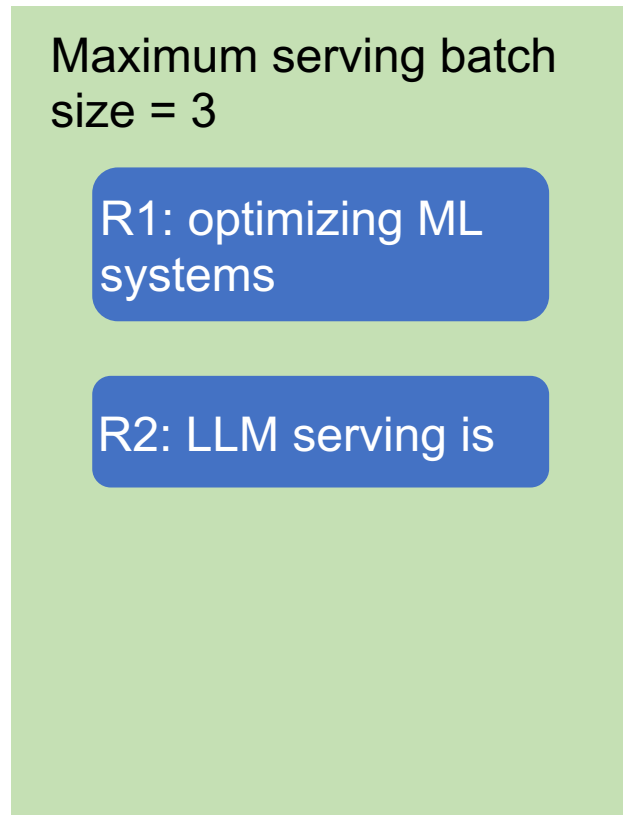


# Continuous Batching Step-by-Step

- Iteration 1: decode R1 and R2



**Request Pool  
(CPU)**



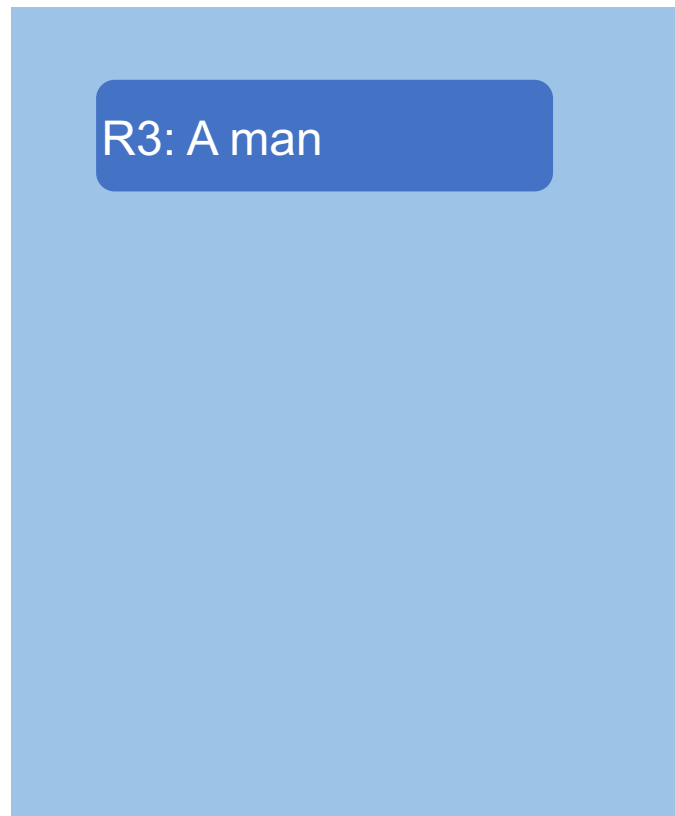
**Execution Engine  
(GPU)**



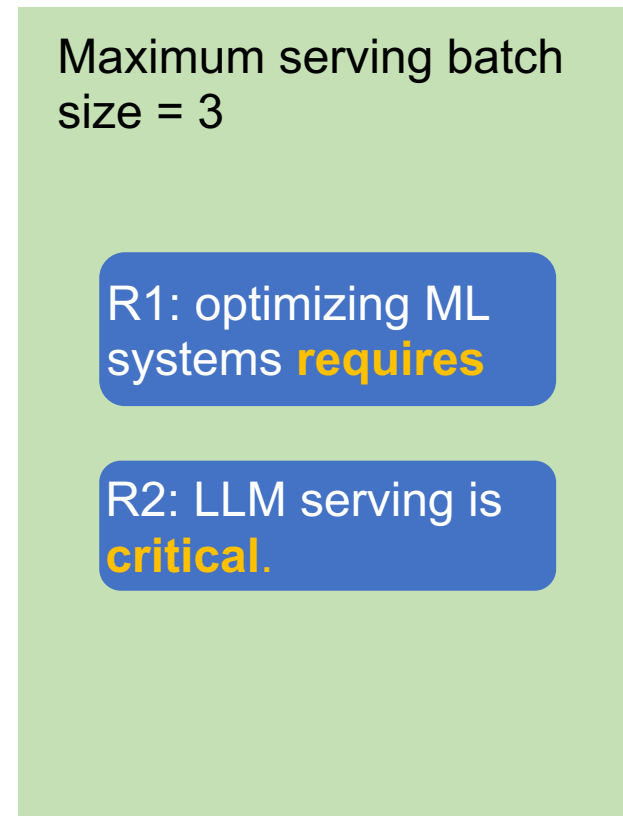
Iteration 1

# Continuous Batching Step-by-Step

- Receive a new request R3; finish decoding R1 and R2



**Request Pool  
(CPU)**



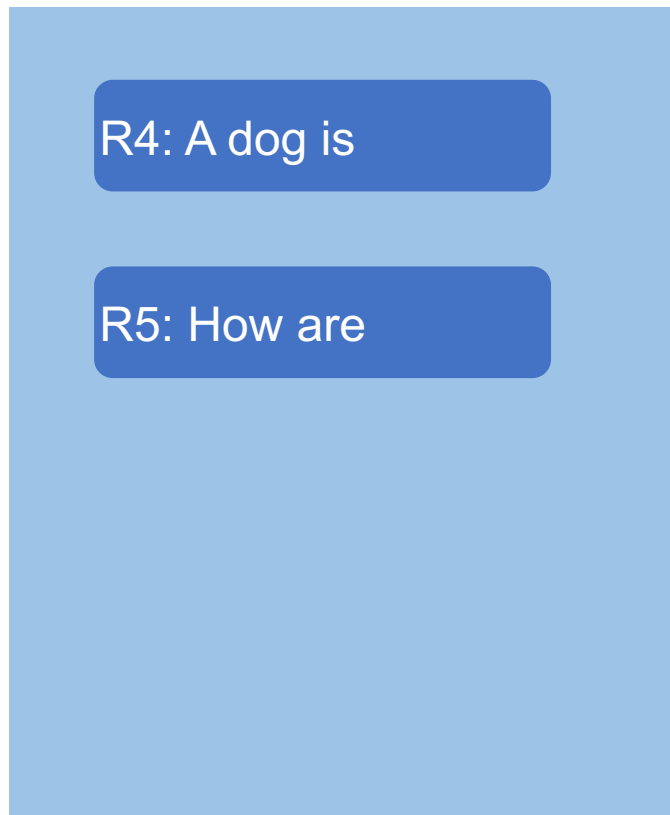
**Execution Engine  
(GPU)**



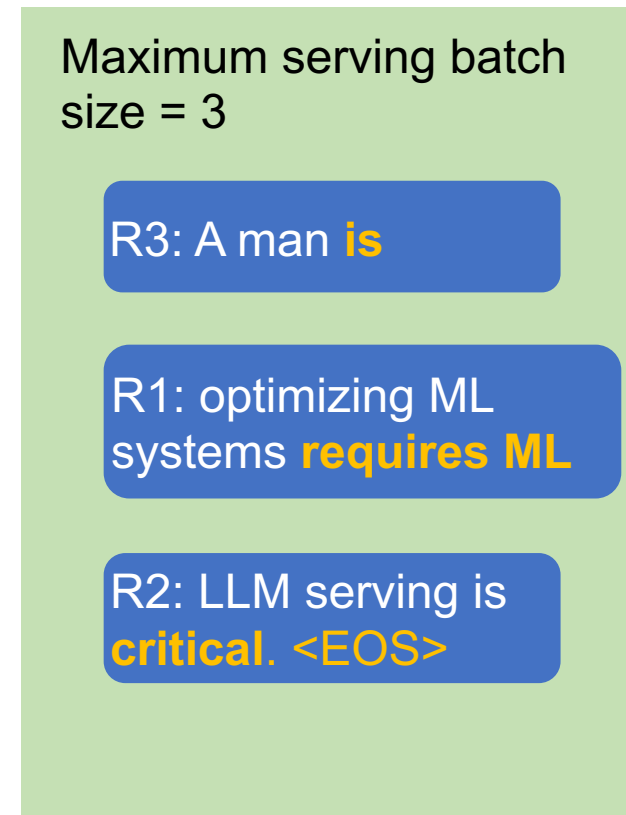
Iteration 1

# Continuous Batching Step-by-Step

- Iteration 2: decode R1, R2, R3; receive R4, R5; R2 completes



**Request Pool  
(CPU)**

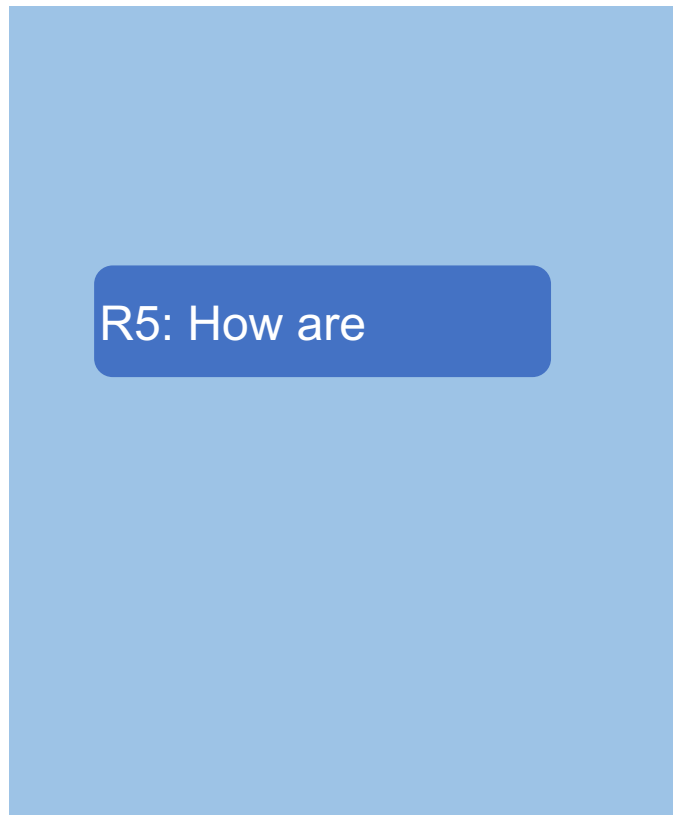


**Execution Engine  
(GPU)**

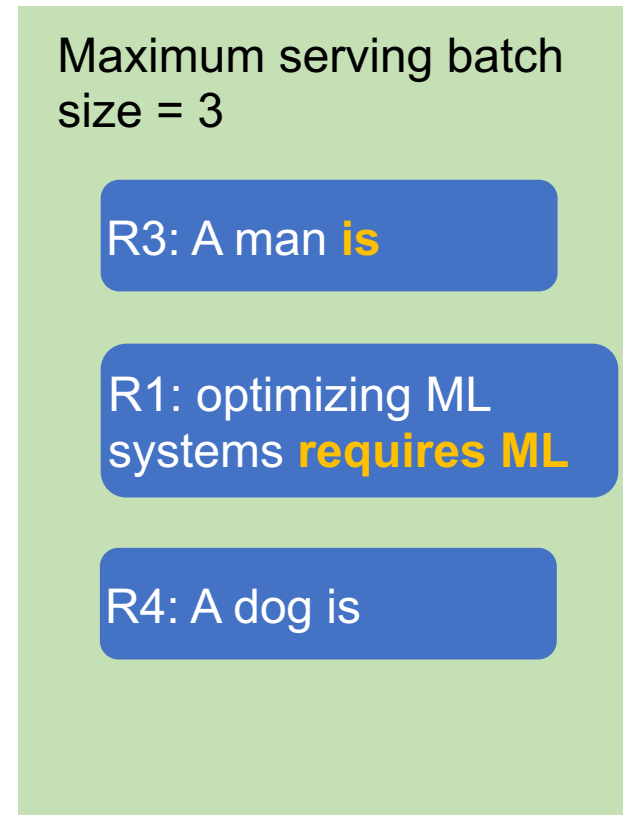


# Continuous Batching Step-by-Step

- Iteration 3: decode R1, R3, R4



**Request Pool  
(CPU)**



**Execution Engine  
(GPU)**



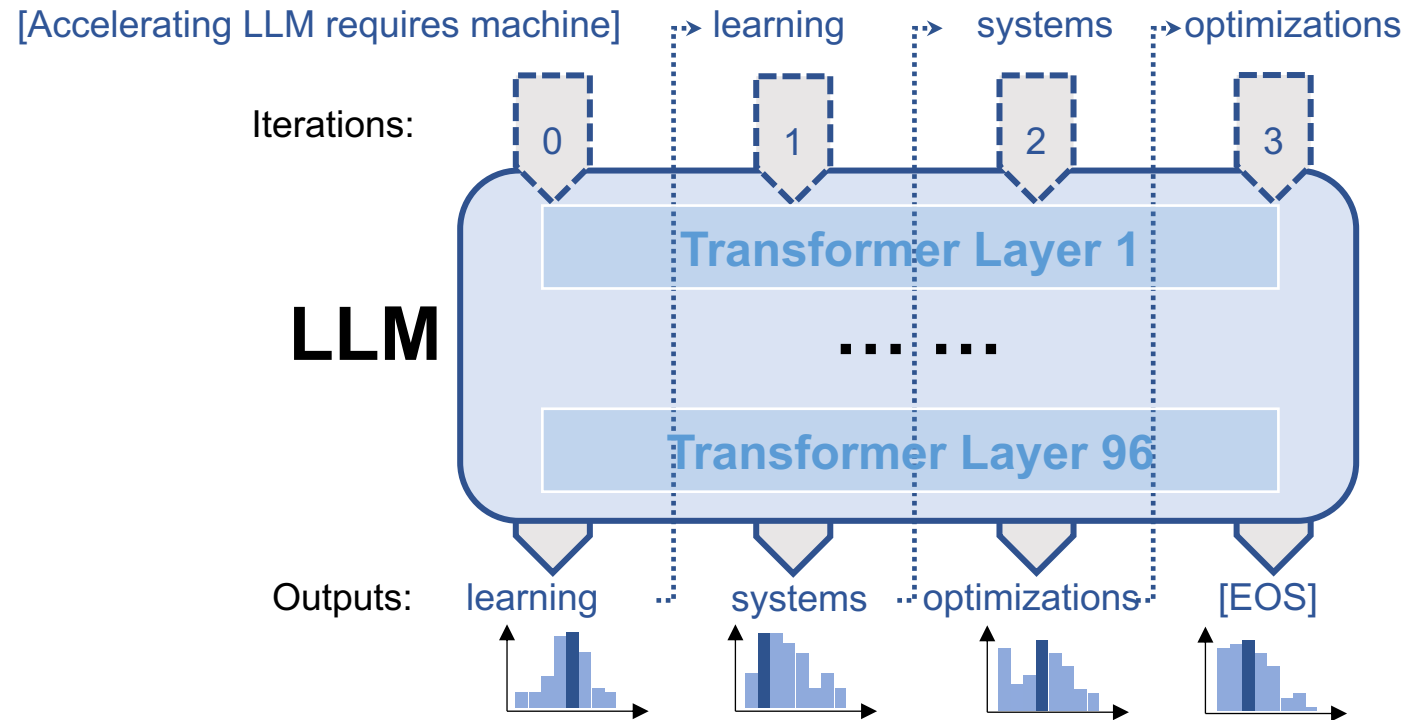
# Continuous Batching

- Handle early-finished and late-arrived requests more efficiently
- Higher GPU utilization

# Outline: LLMs Serving Techniques

- Continuous Batching
- **Speculative Decoding**

# Recall: Incremental Decoding Issues



- Limited degree of parallelism → underutilized GPU resources
- Need all parameters to decode a token → bottlenecked by GPU memory access

# Tradeoffs between Different Language Models

# Parameters	175B	13B	2.7B	760M	125M
TriviaQA	<b>71.2</b>	57.5	42.3	26.5	6.96
PIQA	<b>82.3</b>	79.9	75.4	72.0	64.3
SQuAD	<b>64.9</b>	62.6	50.0	39.2	27.5
latency	20 s	7.6s	2.7s	1.1s	<b>0.3s</b>
# A100s	10	1	1	1	<b>1</b>

Comparing multiple GPT-3 models\*

## Large models

 Pro: better generative performance

 Con: slow and expensive to serve

## Small models

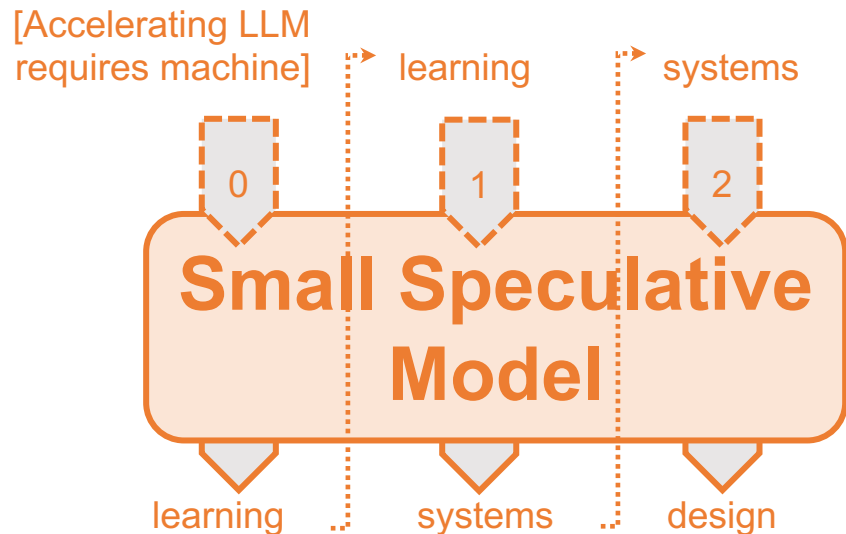
 Pro: cheap and fast

 Con: less accurate



# Speculative Decoding

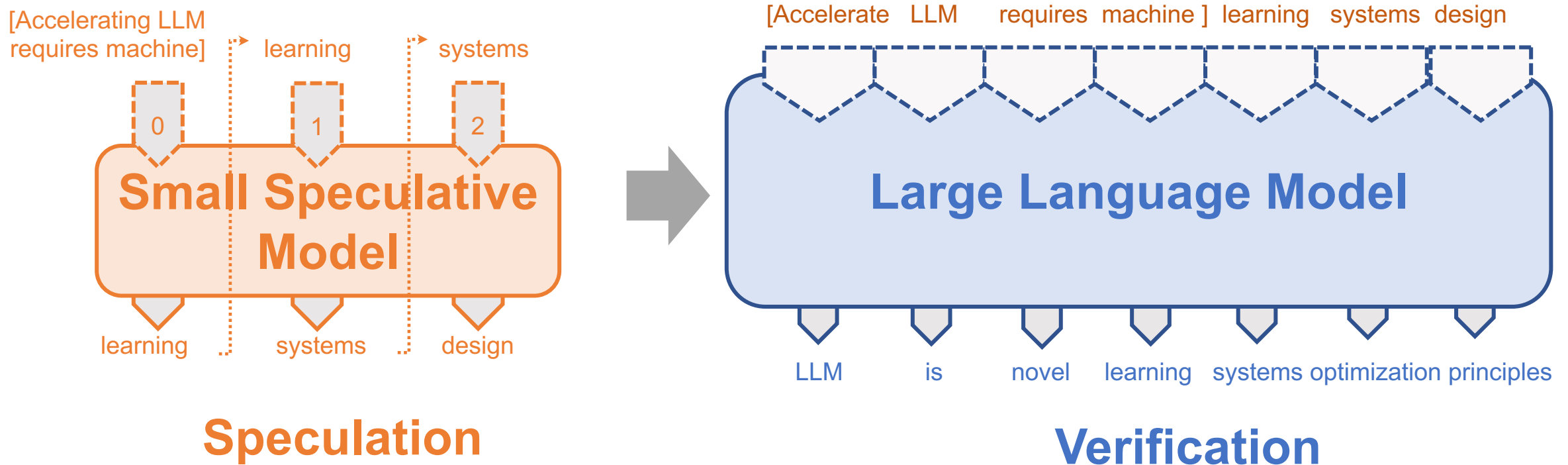
1. Use a small speculative model (SSM) to predict the LLM's output
  - SSM runs much faster than LLM



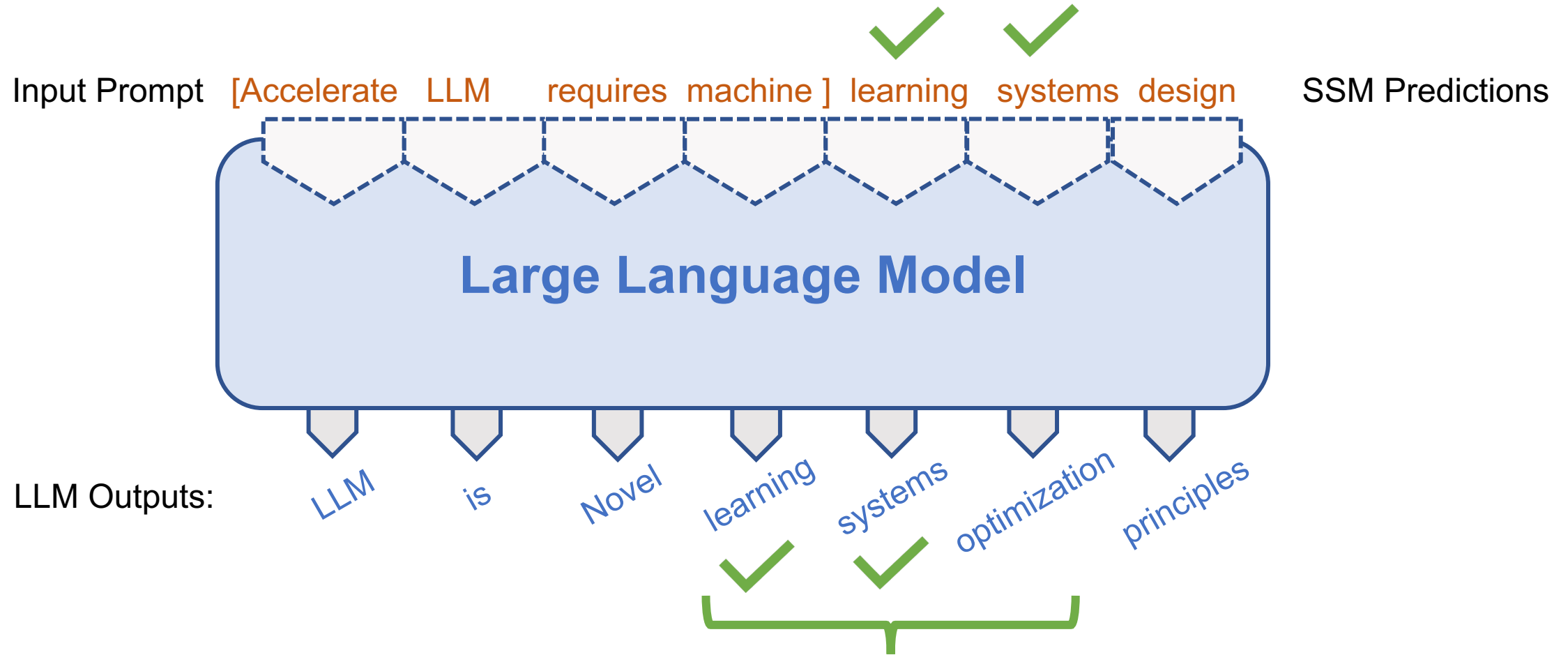
**Speculation**

# Speculative Decoding

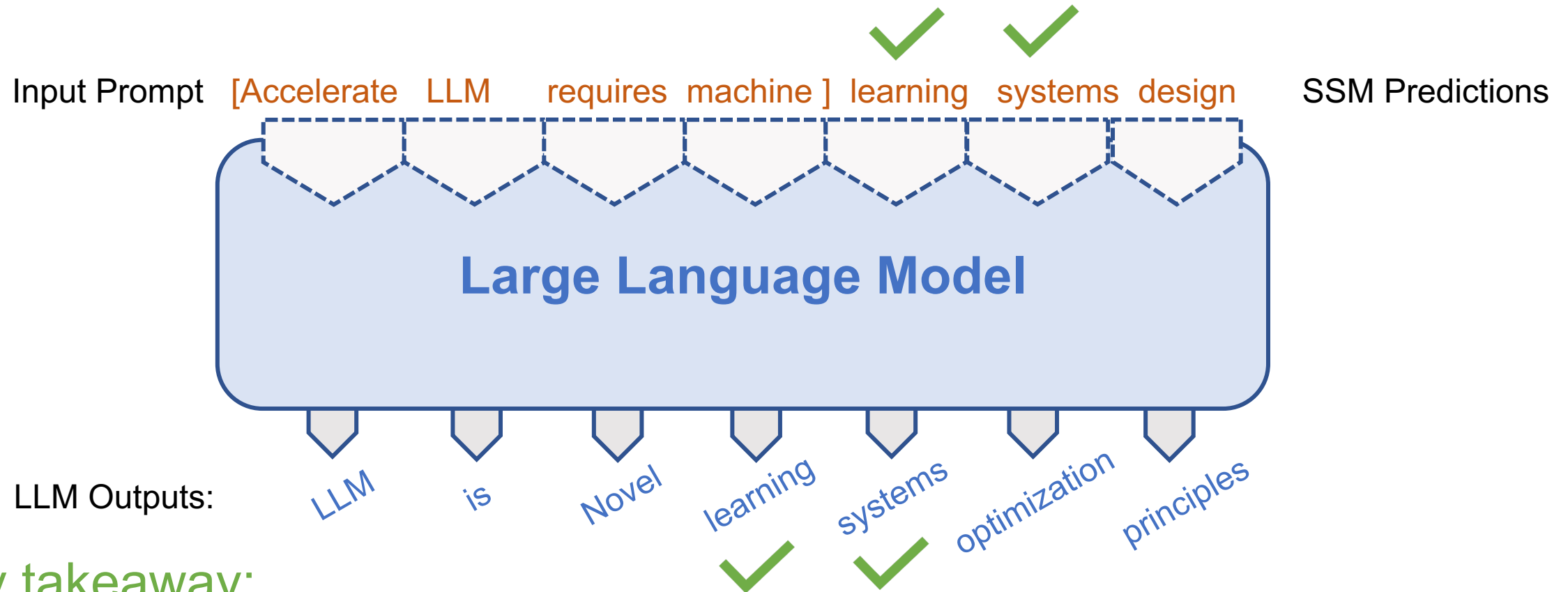
1. Use a small speculative model (SSM) to predict the LLM's output
  - SSM runs much faster than LLM
2. Use the LLM to verify the SSM's prediction



# Verifying Speculative Decoding Results



# Verifying Speculative Decoding Results



## Key takeaway:

- LLM inference is bottlenecked by accessing model weights
- using LLM to decode multiple tokens to improve GPU utilization

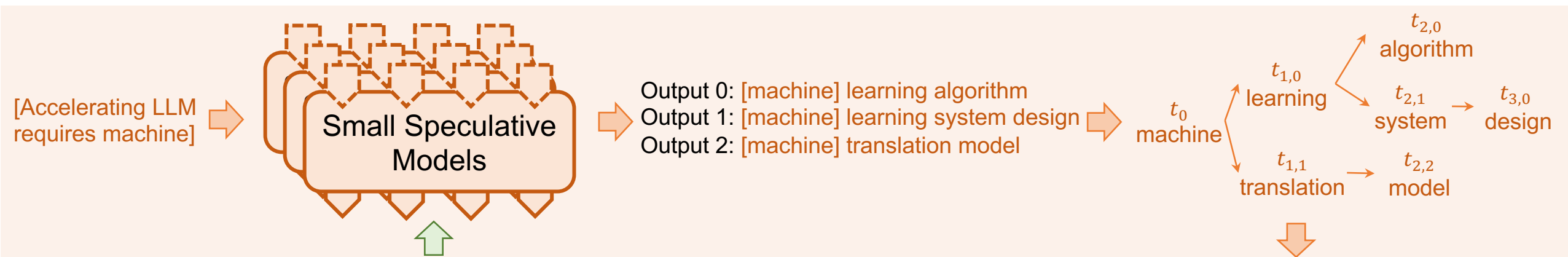
# SpecInfer: Tree-based Speculative Inference & Verification

**Key idea:** not use LLMs as incremental decoder, use them as **parallel token tree verifier**

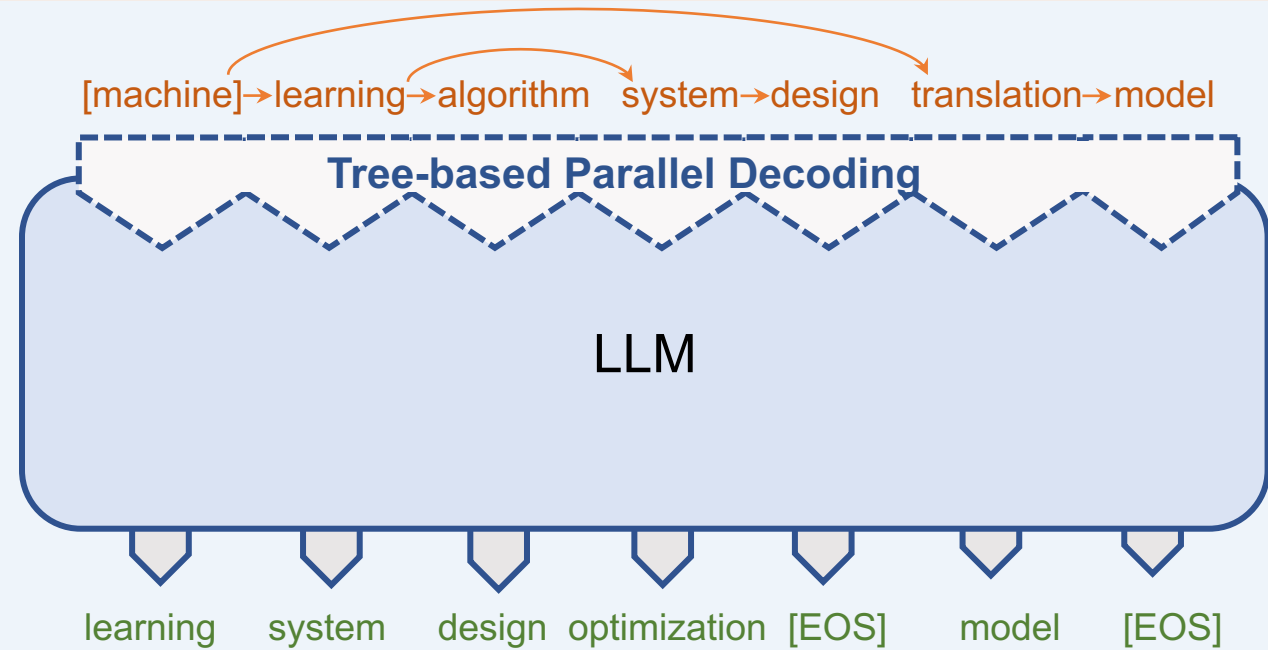
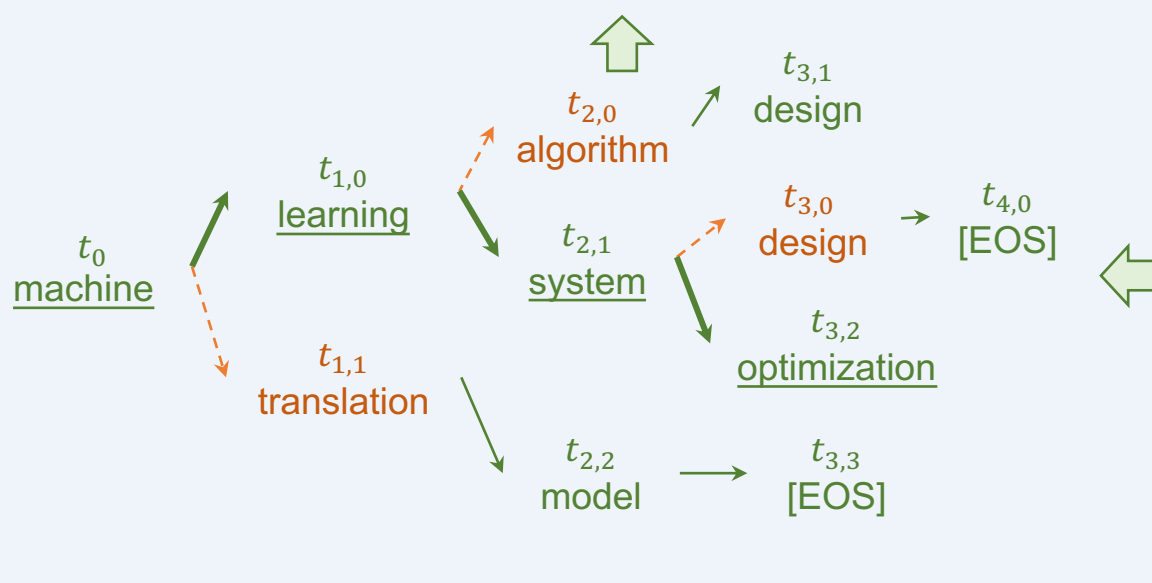
- **Better performance:** outperform existing LLM systems by 1.3-2.4x
- **Higher efficiency:** reduce GPU memory access by 2.5-4.4x
- **Correctness:** verification guarantees end-to-end equivalence

# SpecInfer Workflow

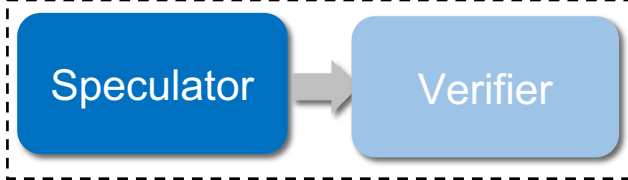
## Speculator



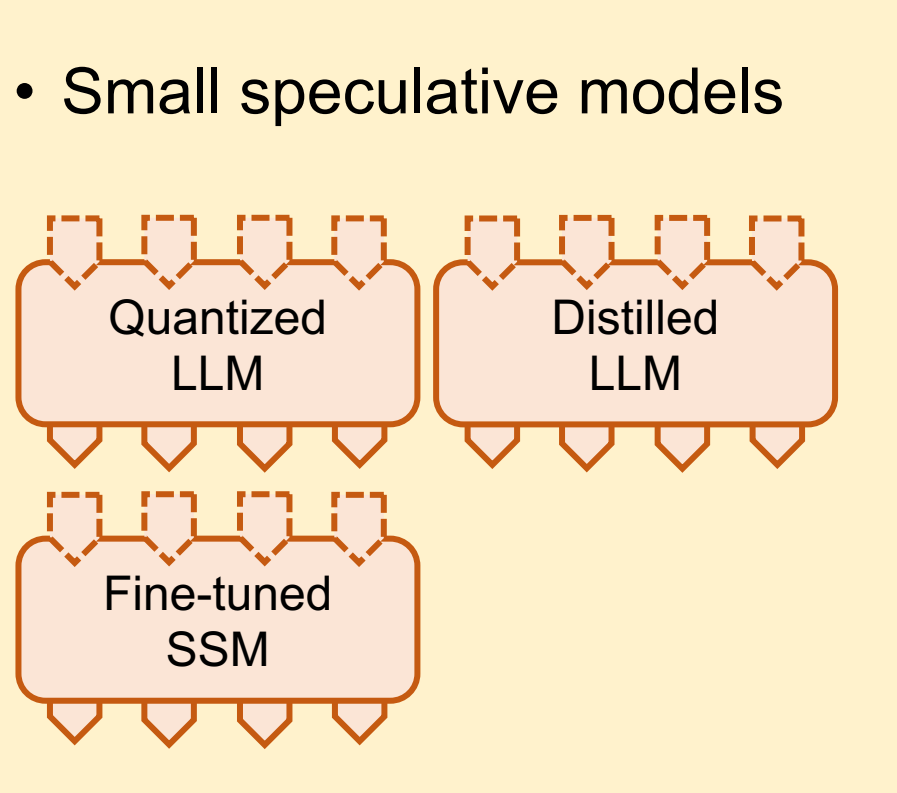
Verified output: machine learning system optimization



## Verifier



# Learning-based Speculator

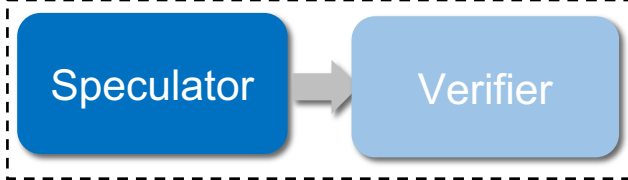


**SSM 0:** [machine] → intelligence

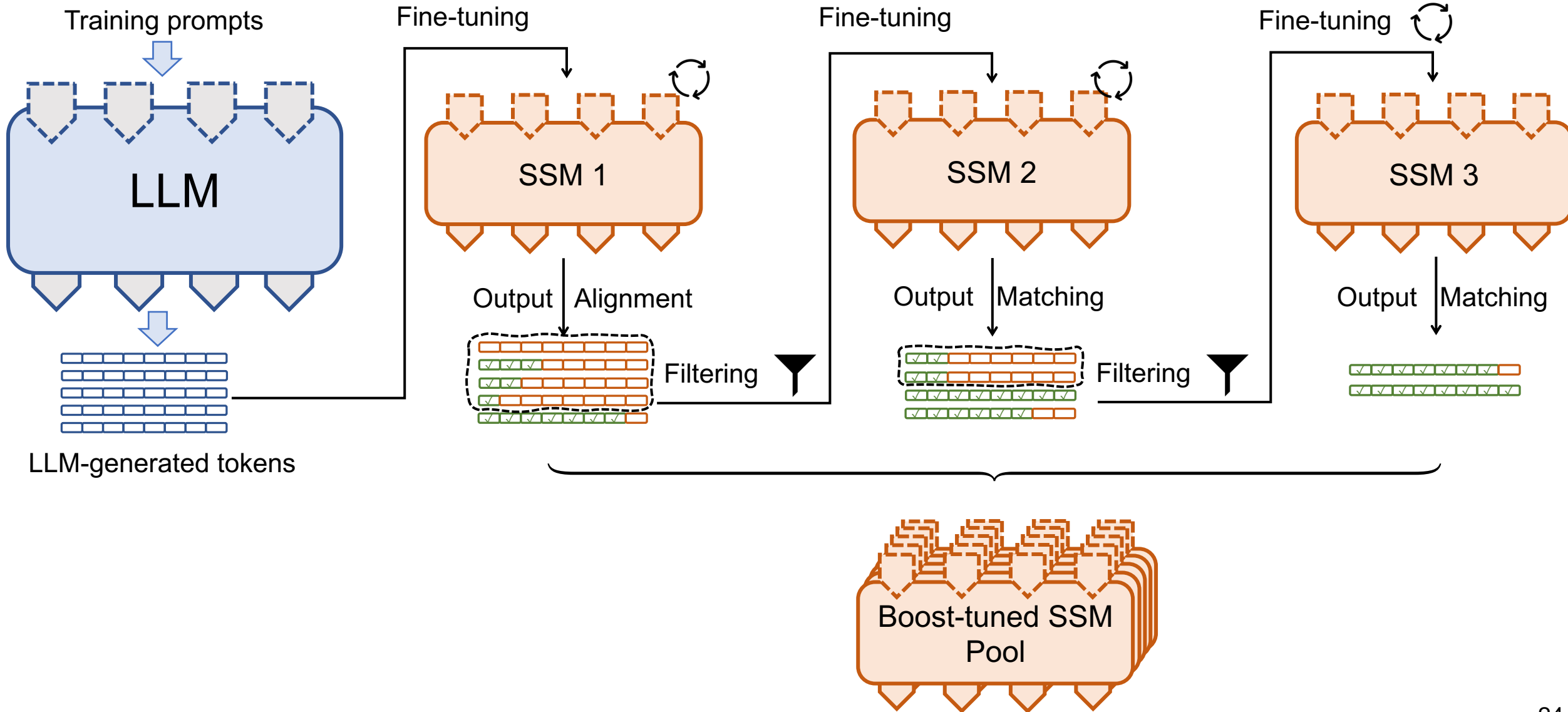
**SSM 1:** [machine] → translation → model

**SSM 2:** [machine] → learning → { algorithm, system → design }

**Speculated Tokens**

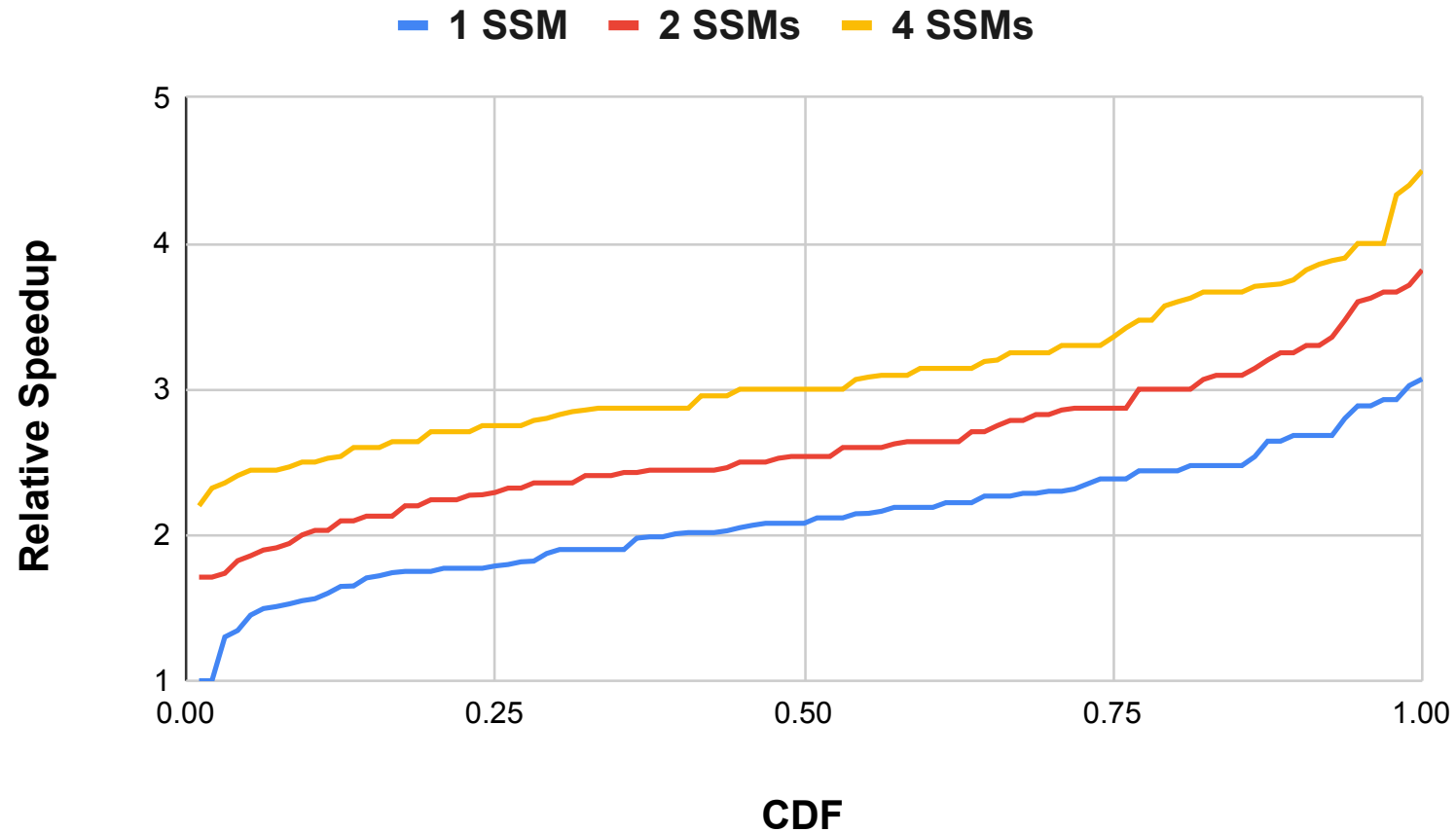


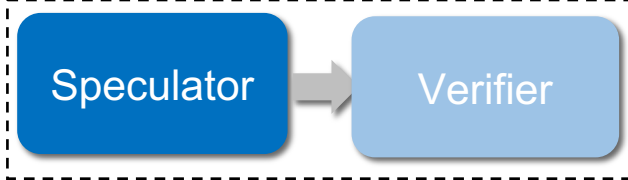
# Collective Boost-Tuning





# Collective Boost-Tuning Consistently Improves Performance

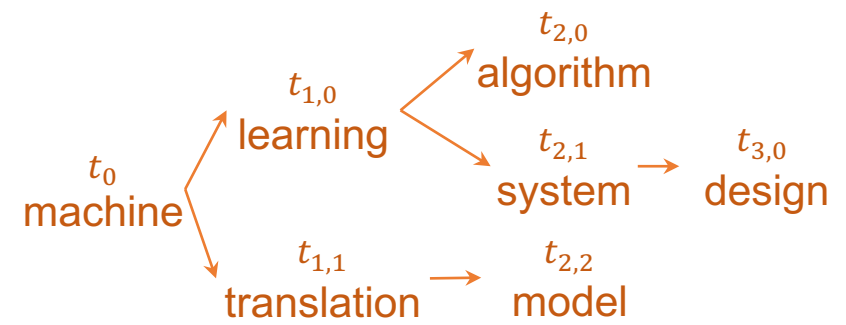
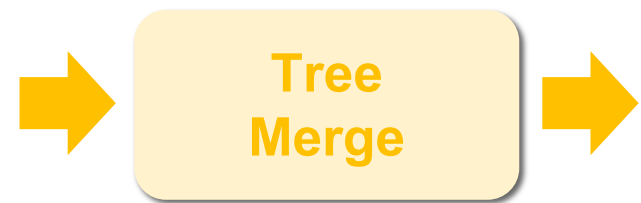




# Token Tree Merge

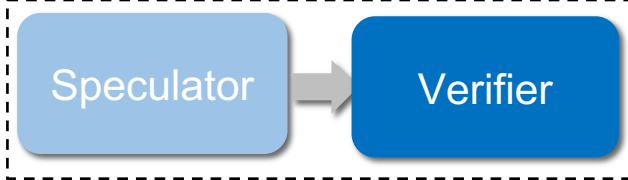
- A compact way to represent speculated tokens

SSM 0: [machine] learning algorithm  
SSM 1: [machine] learning system design  
SSM 2: [machine] translation model

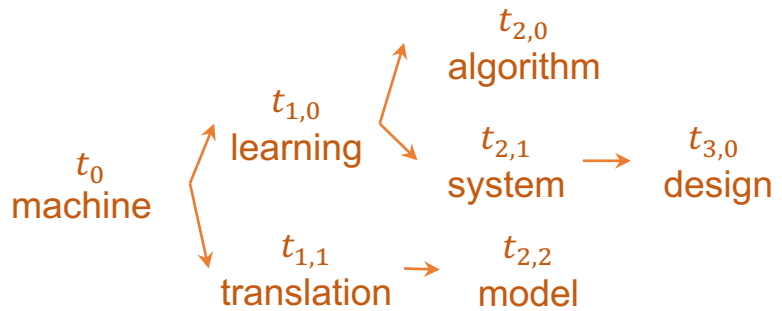


Token Sequences

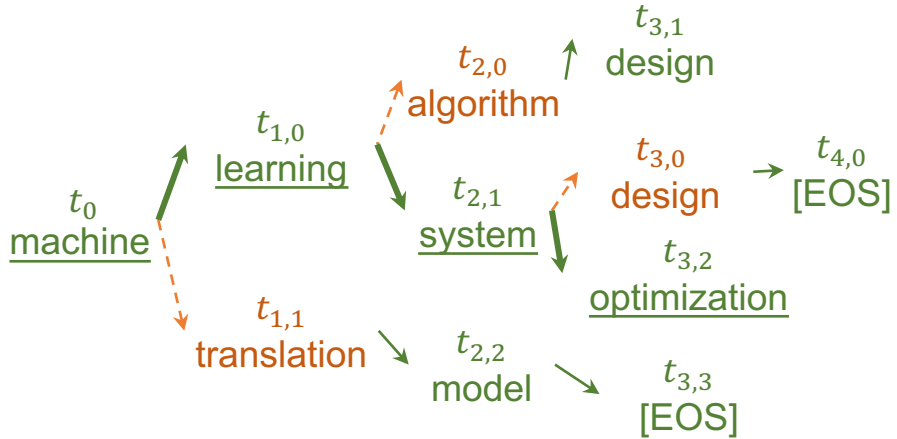
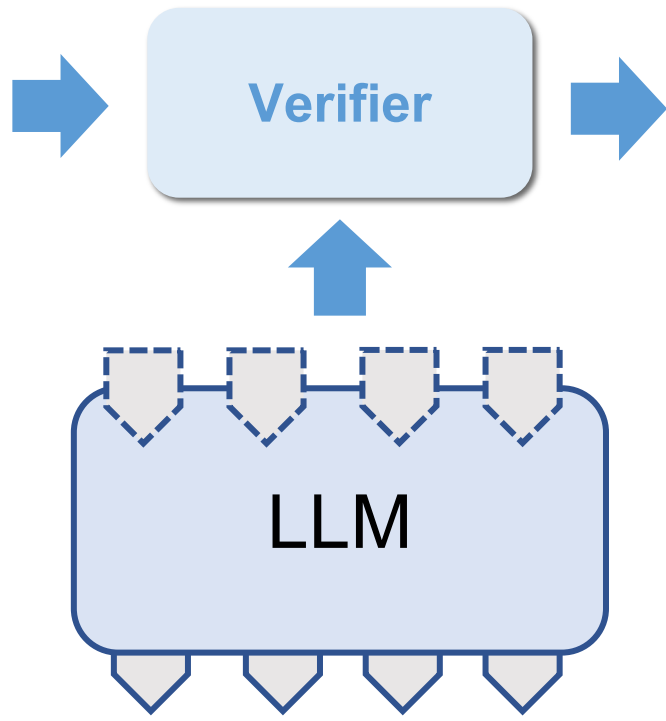
Token Tree



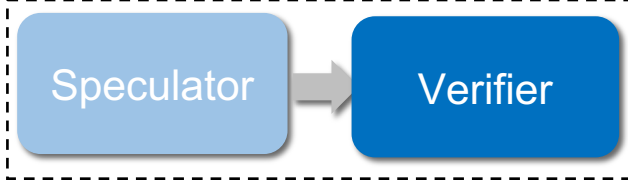
# Token Tree Verifier



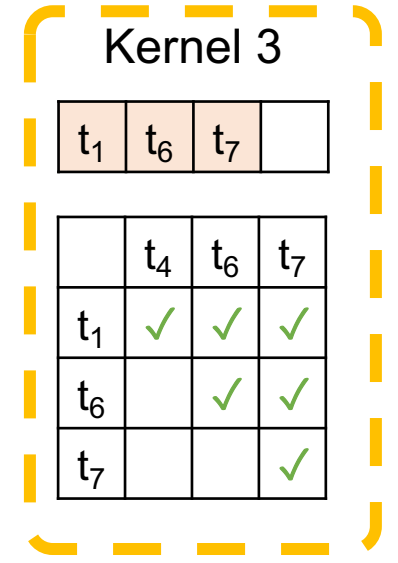
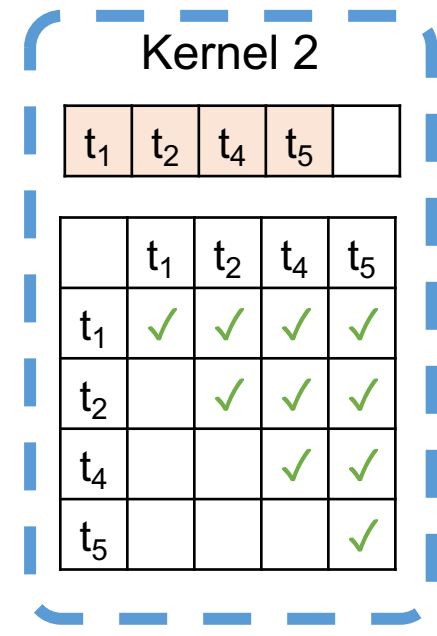
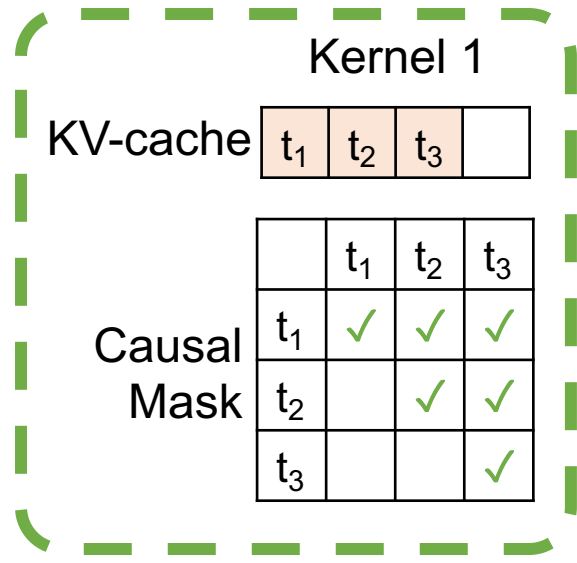
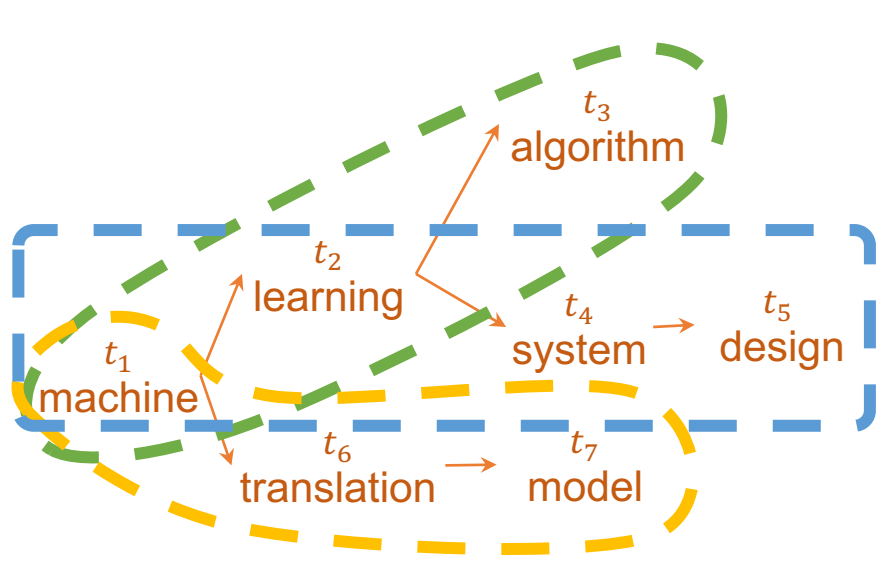
Speculated token tree



Verified token tree

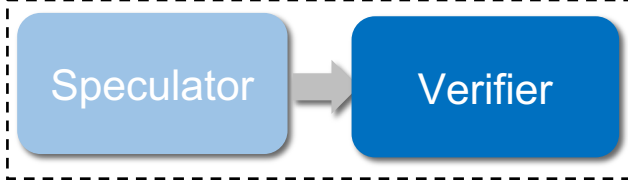


# Sequence-based Decoding



## Issues:

- Redundant decoding computation
- More requests → more GPU memory for key/value cache



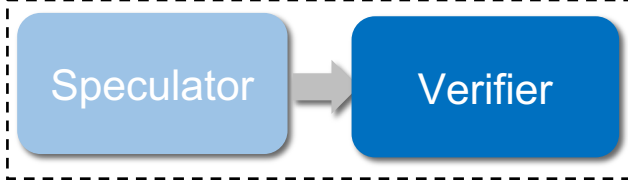
# Tree Attention

- same output as sequence attention for each token; no redundancy

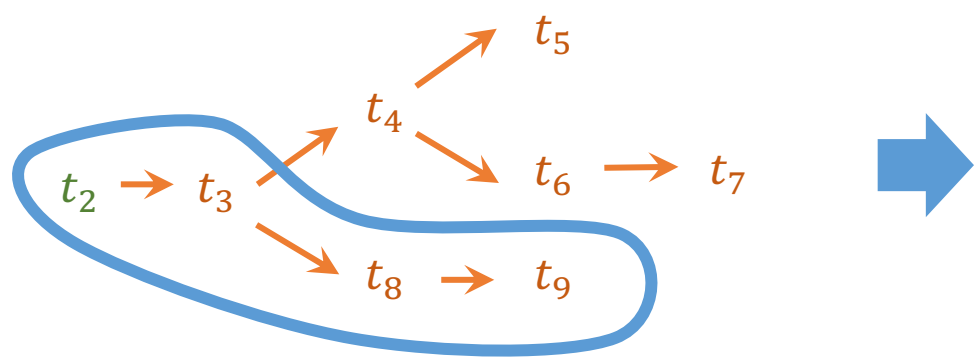
## Token Sequences

Attention( [machine] learning algorithm  
 [machine] learning system design )  
 [machine] translation model





# Tree-based Parallel Decoding



KV-cache 

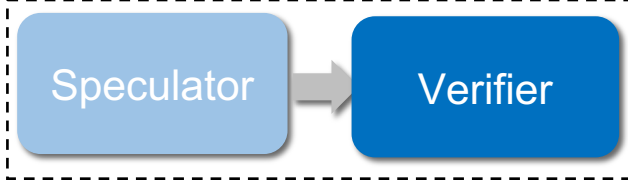
t <sub>2</sub>	t <sub>3</sub>	t <sub>4</sub>	t <sub>5</sub>	t <sub>6</sub>	t <sub>7</sub>	t <sub>8</sub>	t <sub>9</sub>	...
----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	-----

Topology-aware causal mask

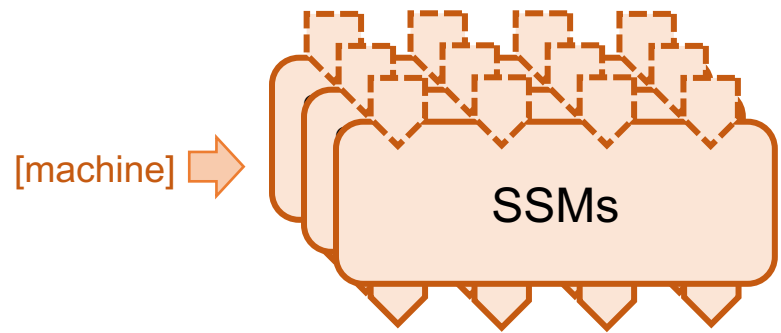
	t <sub>2</sub>	t <sub>3</sub>	t <sub>4</sub>	t <sub>5</sub>	t <sub>6</sub>	t <sub>7</sub>	t <sub>8</sub>	t <sub>9</sub>
t <sub>2</sub>	✓	✓	✓	✓	✓	✓	✓	✓
t <sub>3</sub>		✓	✓	✓	✓	✓	✓	✓
t <sub>4</sub>			✓	✓	✓	✓		
t <sub>5</sub>				✓				
t <sub>6</sub>					✓	✓		
t <sub>7</sub>						✓		
t <sub>8</sub>							✓	✓
t <sub>9</sub>								✓

## Key optimizations:

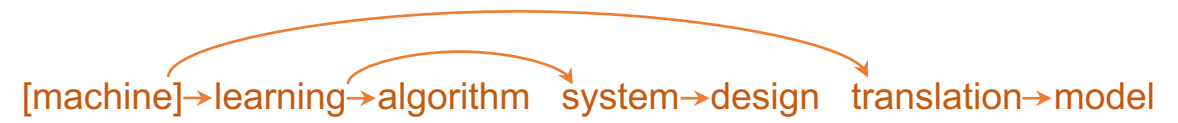
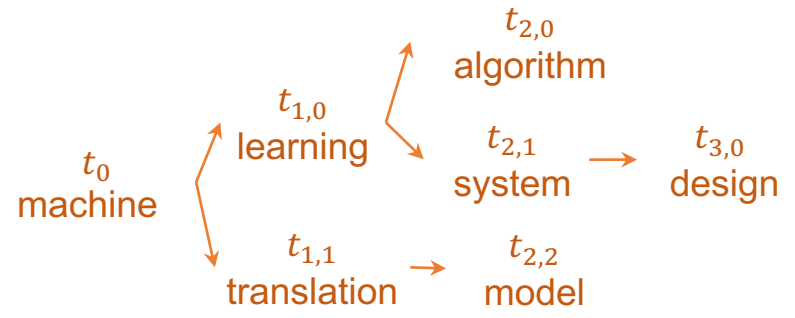
- A DFS-based approach to linearizing a token tree
- Tree topology-aware causal mask
- Decoding all tokens in a single GPU kernel



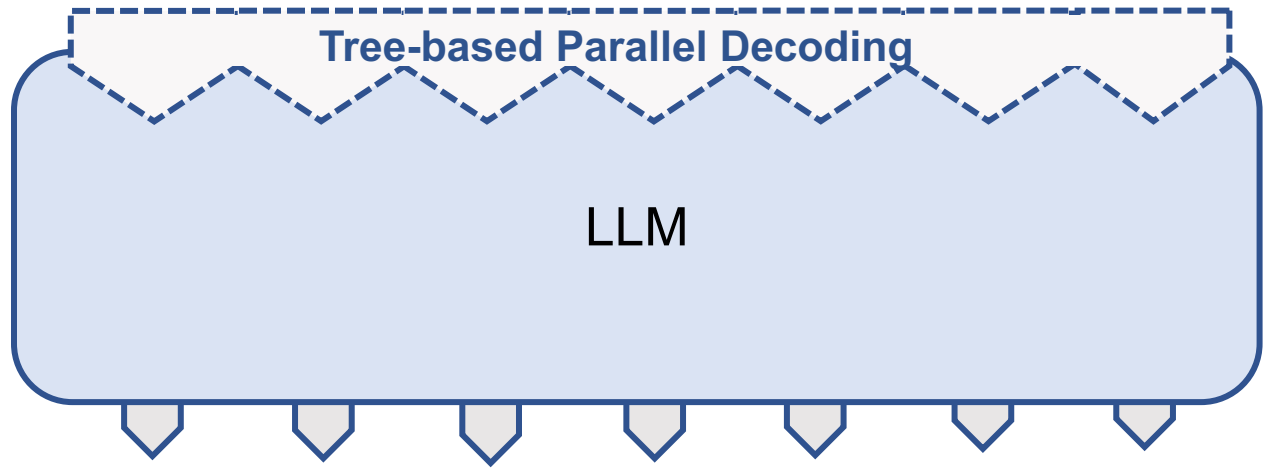
# Verification Workflow



## Speculated Token Tree

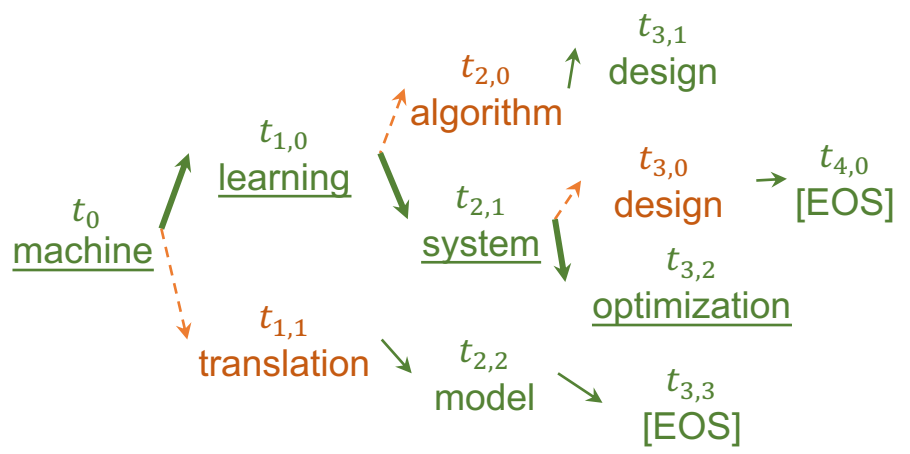


## Tree-based Parallel Decoding



learning system design optimization [EOS] model [EOS]

Verified output: machine learning system optimization



## Verified token tree

# Stochastic Decoding

- **Challenge:** verifying stochastic equivalence

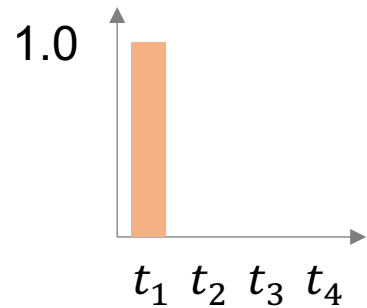
$$P_{\text{IncrDecode}}(\cdot | x_{<i}, \text{LLM}) = P_{\text{SpecInfer}}(\cdot | x_{<i}, \text{LLM}, \{\text{SSM}_i\})$$

- A strawman approach: **naïve sampling**
  - Use LLM to sample  $x_i \sim P_{\text{IncrDecode}}(\cdot | x_{<i}, \text{LLM})$
  - Verify if  $x_i$  is in the token tree

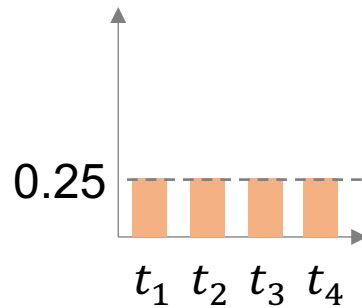


# Naïve Sampling can be Suboptimal

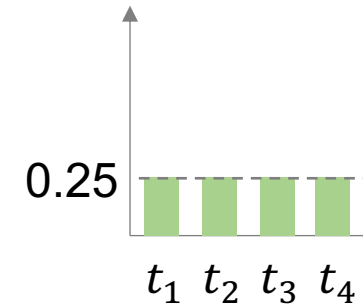
- Assume one LLM, two SSMs, and four possible tokens:  $t_1, t_2, t_3, t_4$



SSM 1:  $P(\cdot | x_{<i}, SSM_1)$



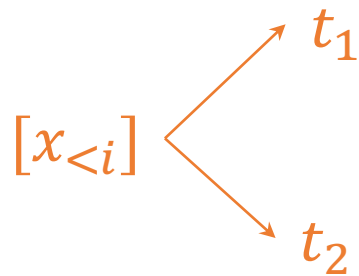
SSM 2:  $P(\cdot | x_{<i}, SSM_2)$



LLM:  $P(\cdot | x_{<i}, LLM)$

Naïve sampling's verification prob. = **50%**

But we can do better by directly accepting SSM 2;  
verification prob. = **100%**



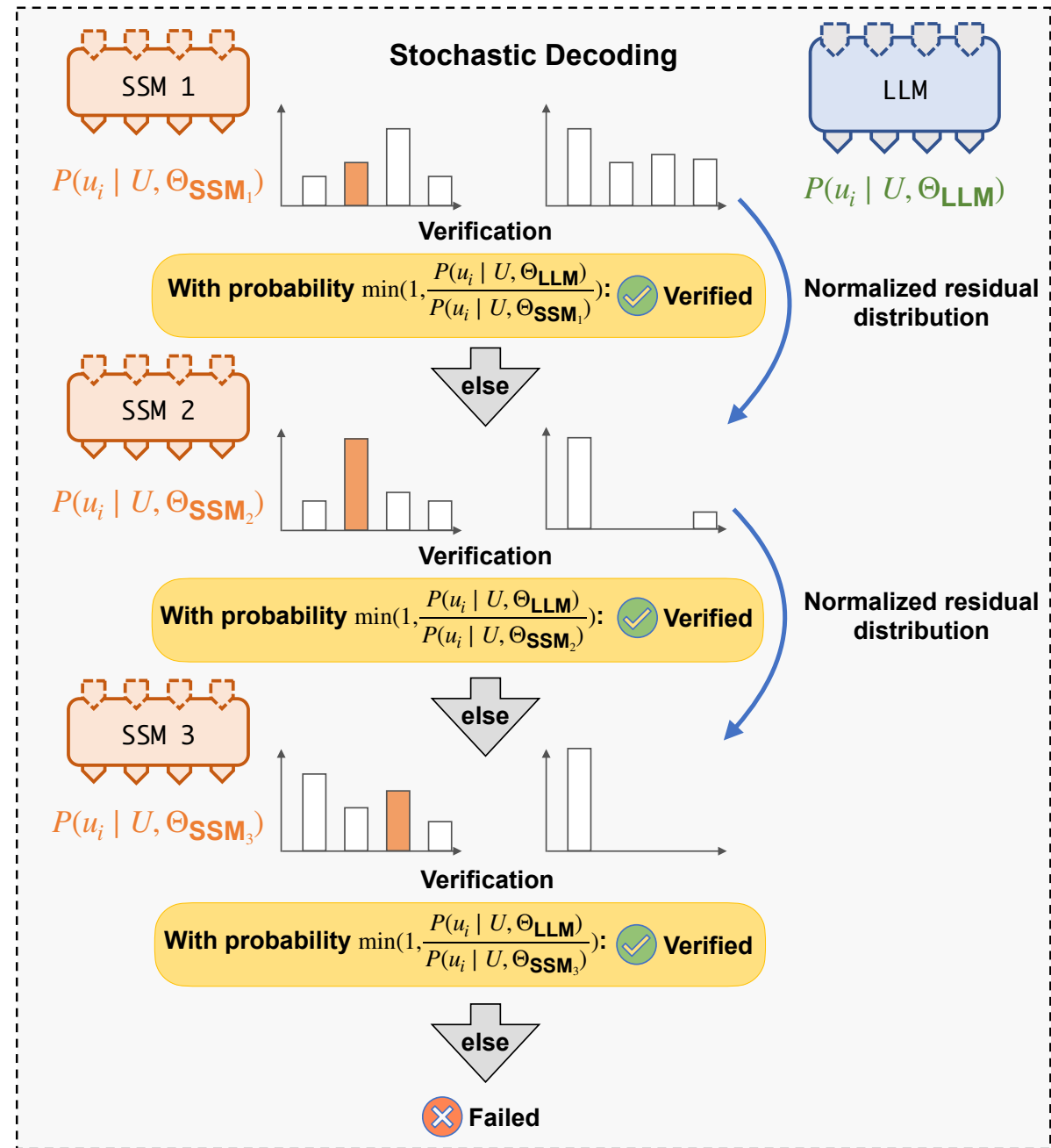
**Token Tree**

Key issue: naïve sampling ignores correlation between  $P(\cdot | x_{<i}, SSM)$  and  $P(\cdot | x_{<i}, LLM)$

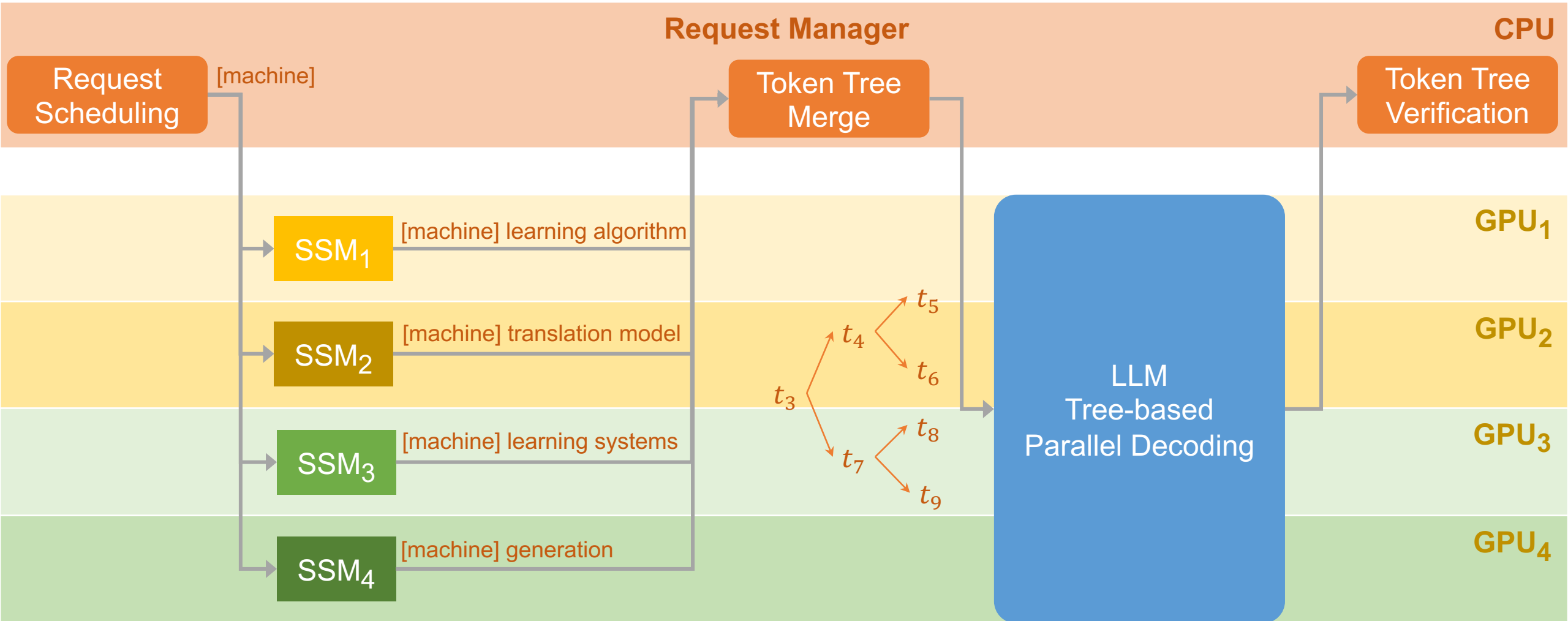
# Speculative Sampling

1. Sample a token  $x \sim P(u_i | U, \Theta_{SSM})$  using SSM
2. If  $P(x|U, \Theta_{SSM}) \leq P(x|U, \Theta_{LLM})$ , directly accept  $x$
3. If  $P(x|U, \Theta_{SSM}) > P(x|U, \Theta_{LLM})$ , accept  $x$  with prob.  $\frac{P(x|U, \Theta_{LLM})}{P(x|U, \Theta_{SSM})}$
4. If reject  $x$ , normalize residual distribution

$$P'(x|U, \Theta_{LLM}) = \text{norm}(\max(0, P(x|U, \Theta_{LLM}) - P(x|U, \Theta_{SSM})))$$

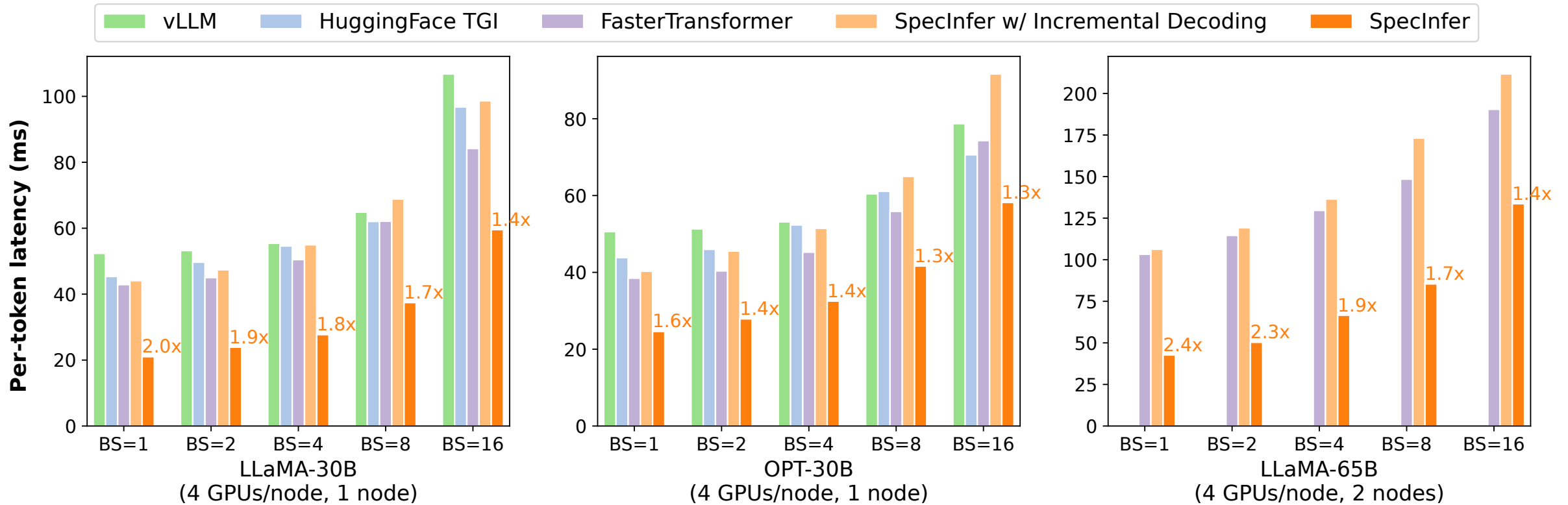


# Distributed LLM Serving

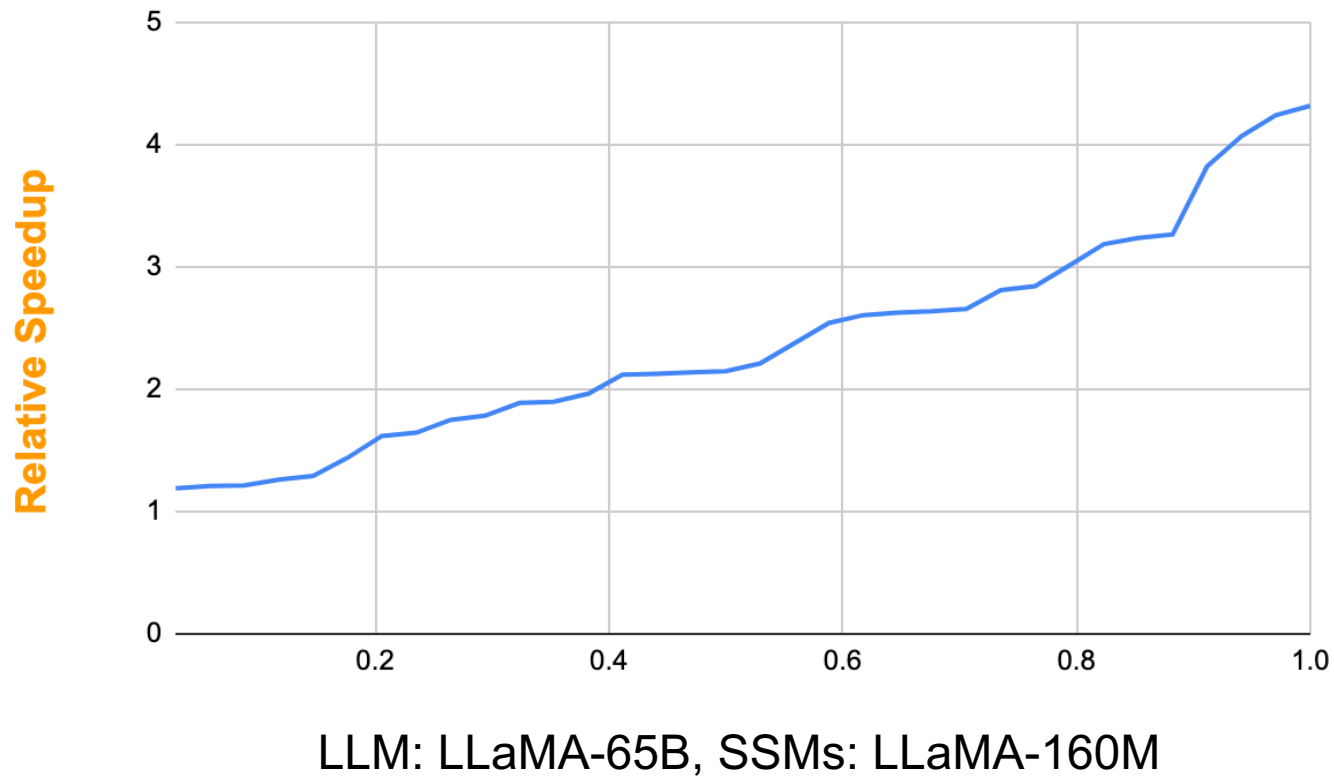


**Tensor / Pipeline  
Model Parallelism**

# SpecInfer Accelerates LLM Inference by 1.3-2.4x

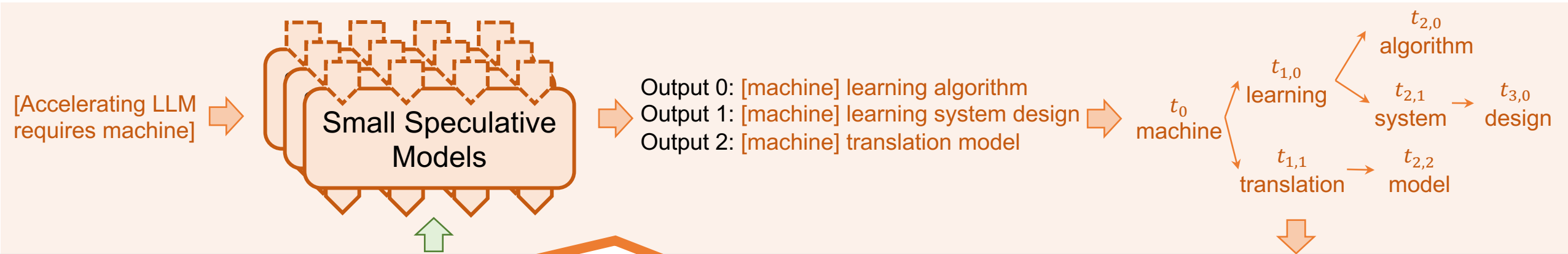


# SpecInfer can Consistently Accelerate LLM Inference



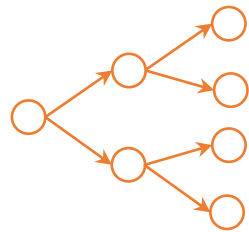
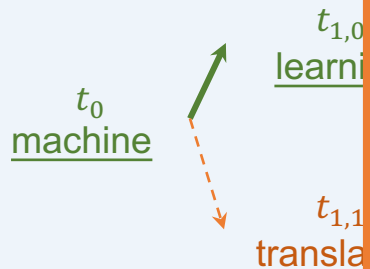
# Open Research Questions

## Speculator

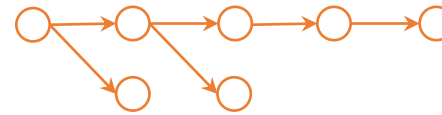


Verified output: machine

### Balanced v.s. unbalanced trees?

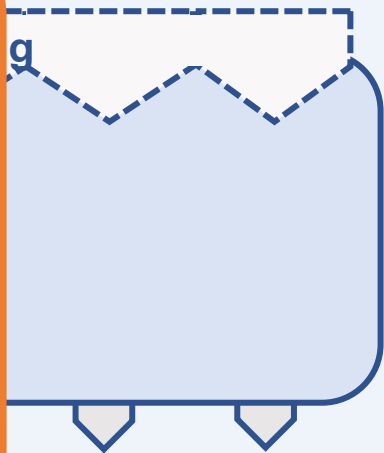


Balanced trees:  
lower speculative overhead



Unbalanced trees:  
higher verification prob.

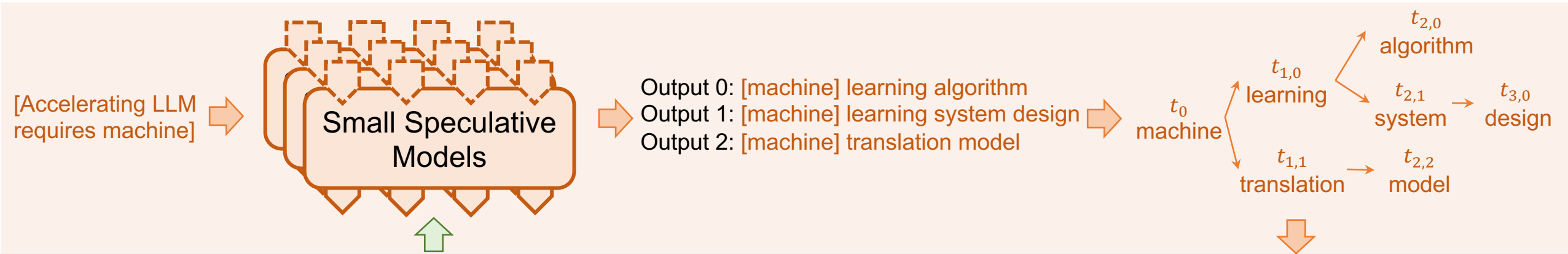
translation → model



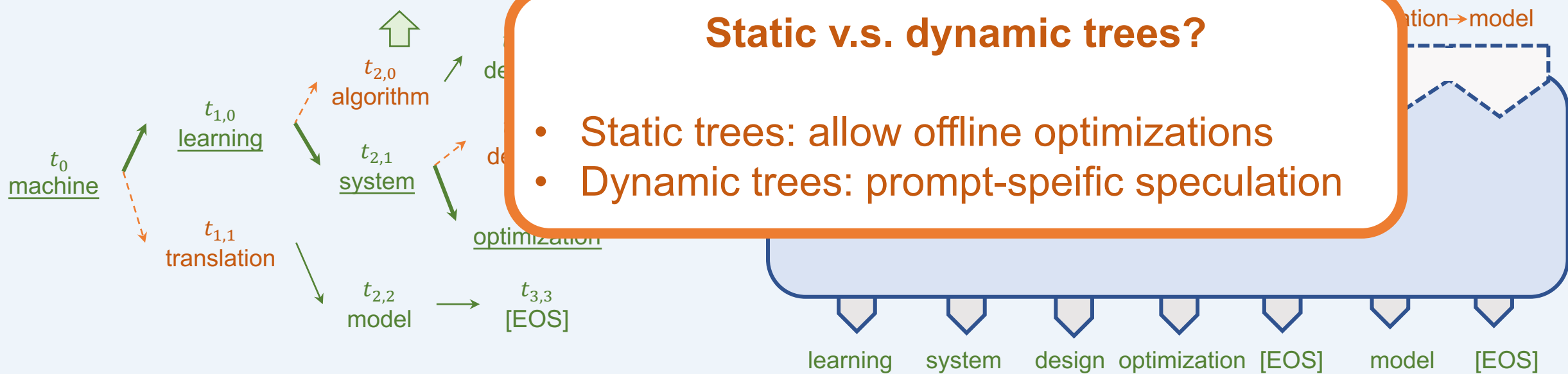
model [EOS]

# Open Research Questions

## Speculator



Verified output: machine learning system optimization



### Static v.s. dynamic trees?

- Static trees: allow offline optimizations
- Dynamic trees: prompt-specific speculation

## Verifier

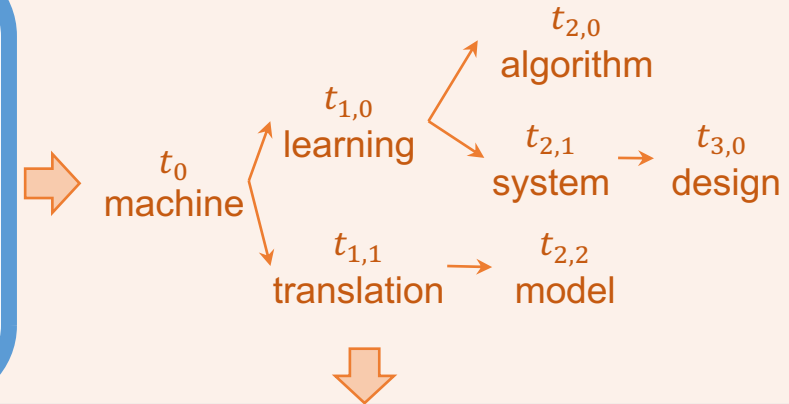
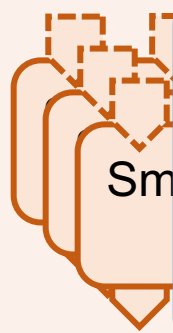
# Open Research Questions

Speculator

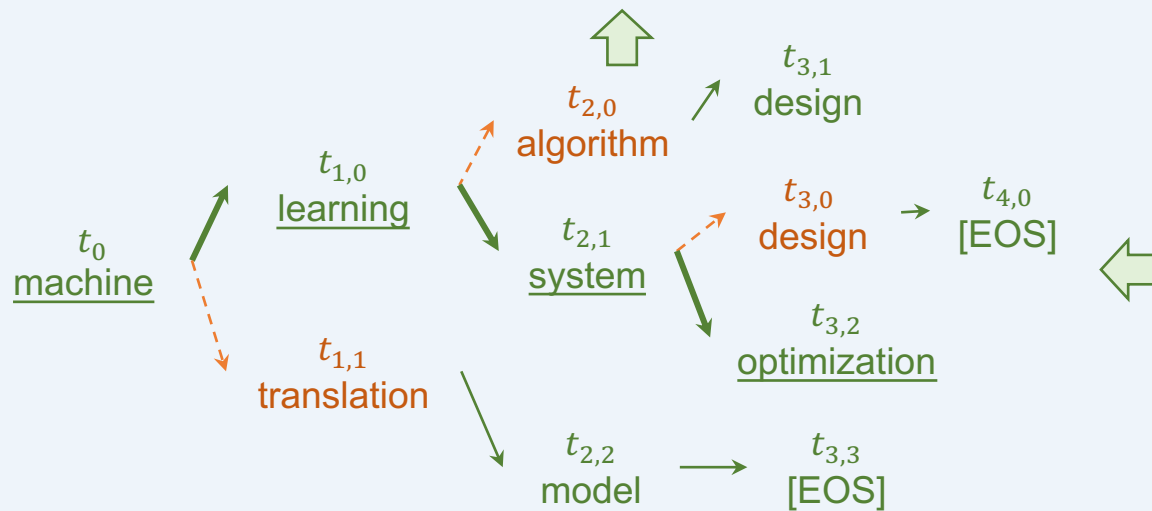
How to verify stochastic decoding?

- Multi-step speculative sampling
- Naïve sampling
- Greedy sampling
- Mixture?

[Accelerating LLM requires machine]



Verified output: machine learning system optimization



[machine] → learning → algorithm → system → design → translation → model

Tree-based Parallel Decoding

LLM

learning system design optimization [EOS] model [EOS]

Verifier



# Recap: LLMs Serving Techniques

- Continuous Batching
- Speculative Decoding