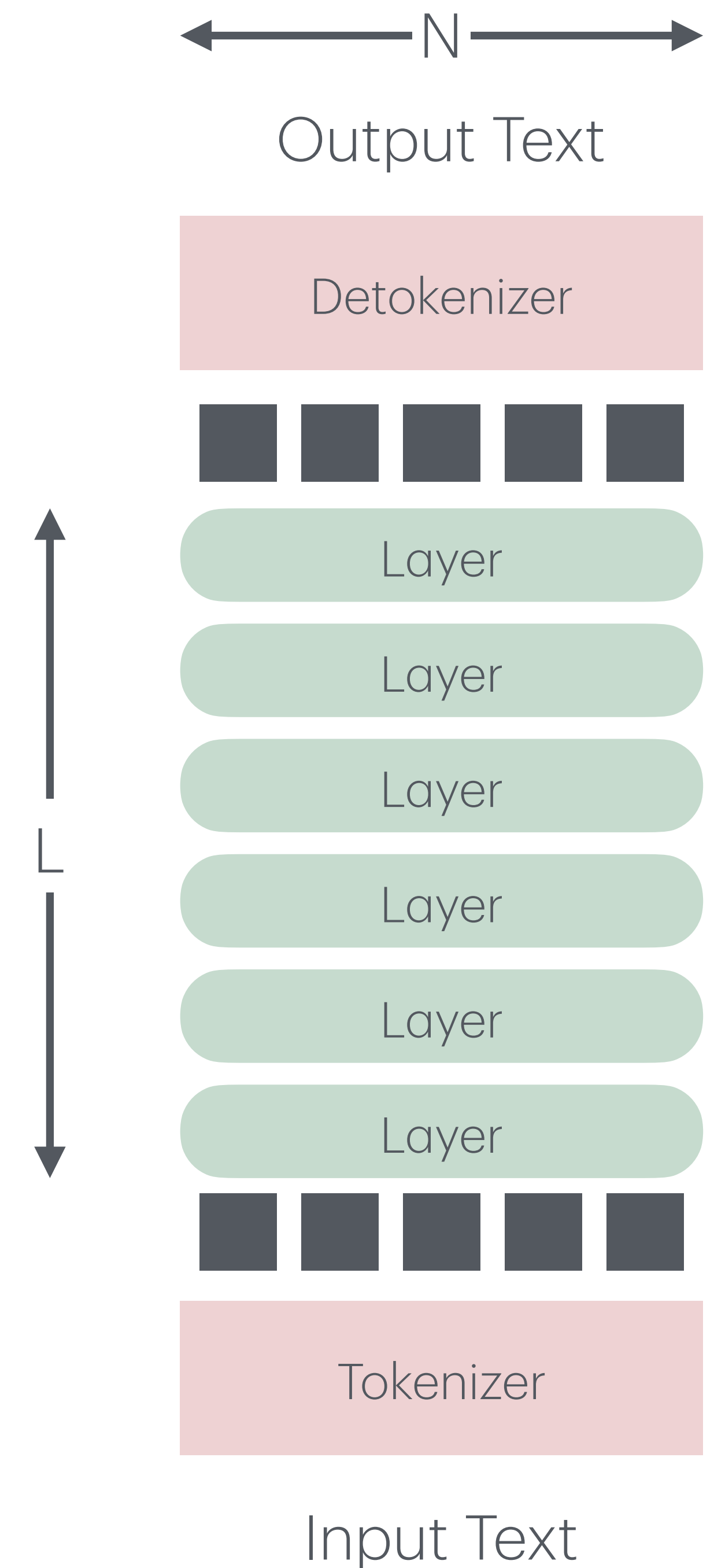


# Speculative Decoding

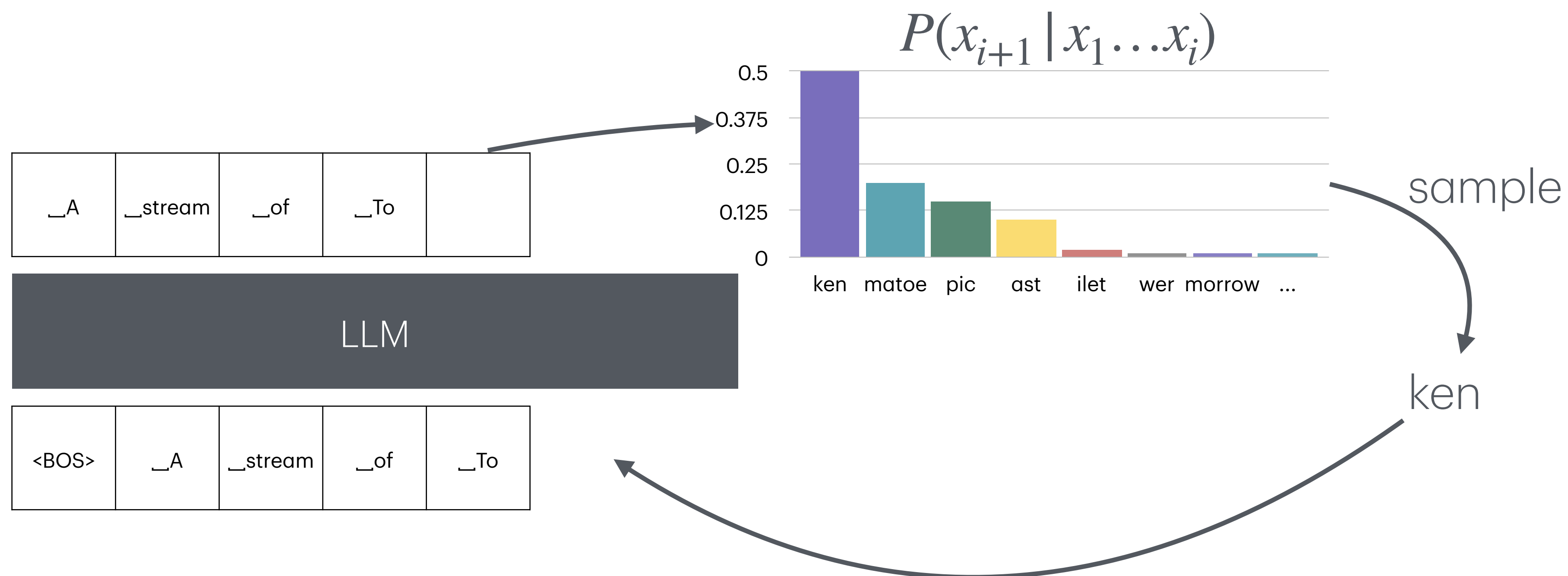
# Training and Generation

## Vanilla Generation

	Training	Training - Checkpointing	Generation
Peak Memory	$O(NL)$	$O(NL^{1/2})$	$O(N)$
Runtime	$O(N^2L)$	$O(2 N^2L)$	$O(N^3L)$
# forward calls	1	1	N



# Generation - A closer look



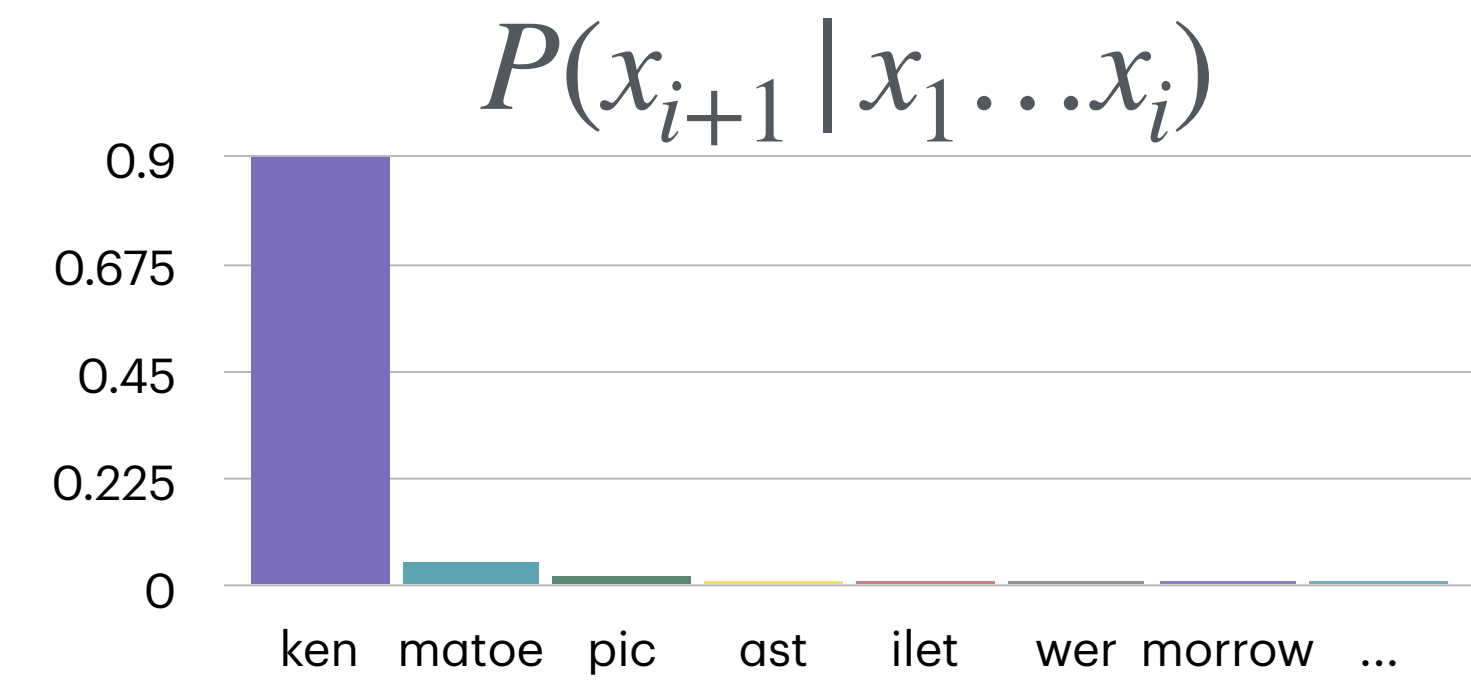
```
x = [<BOS>]
for i = 1..N:
  x_i ~ P(x_i | x)
  x.append(x_i)
```

Overall process: Sample from

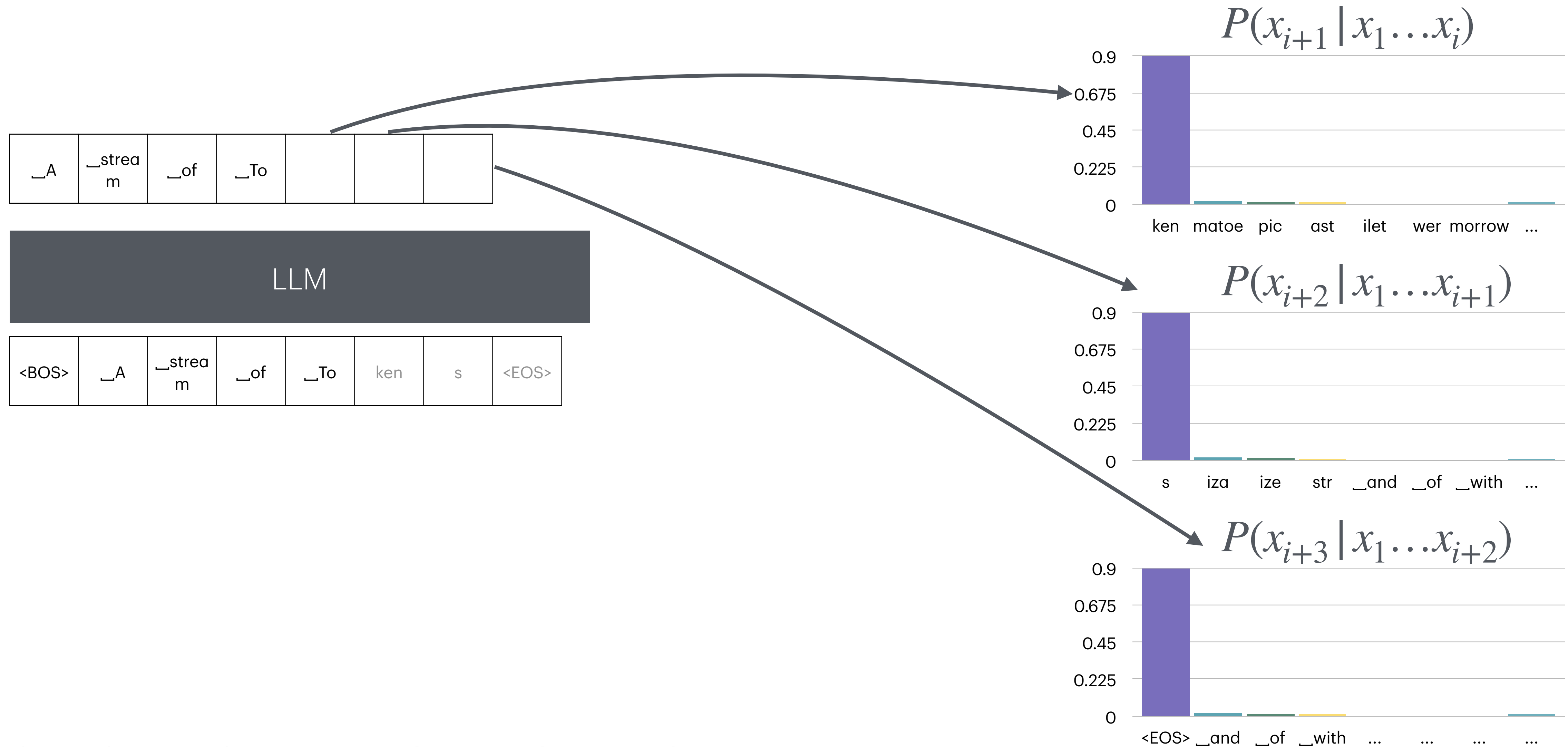
$$P(\mathbf{x}) = \prod_{i=1}^N P(x_{i+1} | x_1 \dots x_i)$$

# Generation - A closer look

- Next token probability often highly peaked (next token is obvious)
- Can we skip predicting obvious tokens?
- Verify instead

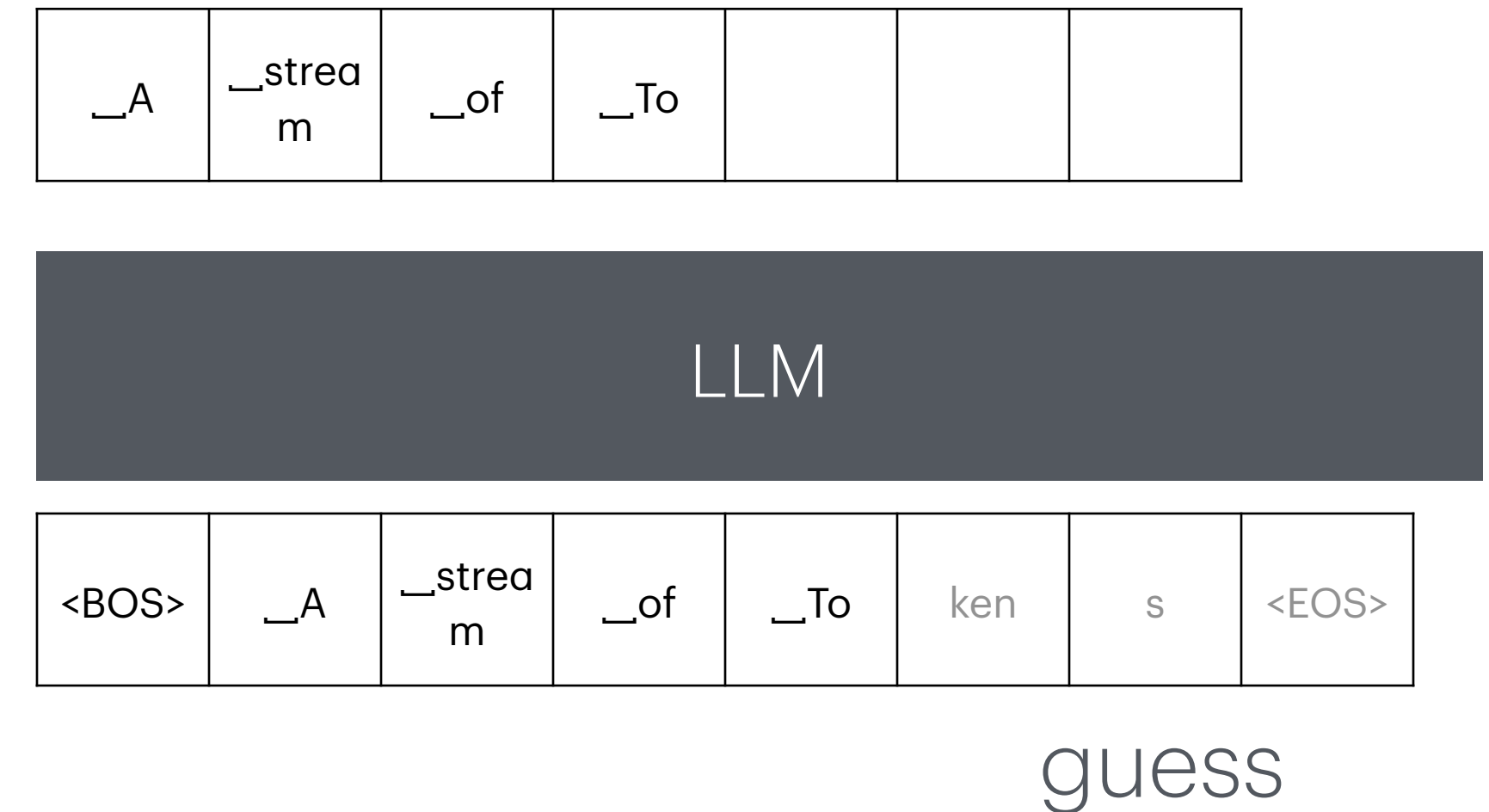


# Verification vs generation



# Speculative decoding

- How do we obtain guess of next tokens?
  - Use a smaller model  $Q(x_{i+1} | Q_1 \dots Q_x)$
  - Use a separate head to predict multiple tokens
- How to we accept guess?
  - $P(x_{i+1} | x_1 \dots x_i) \geq Q(x_{i+1} | x_1 \dots x_i)$
  - or proportional to  $\frac{P(x_{i+1} | x_1 \dots x_i)}{Q(x_{i+1} | x_1 \dots x_i)}$



# Speculative decoding

## Algorithm

---

**Algorithm 1** SpeculativeDecodingStep

---

**Inputs:**  $M_p, M_q, prefix$ .

▷ **Sample  $\gamma$  guesses  $x_1, \dots, x_\gamma$  from  $M_q$  autoregressively.**

**for**  $i = 1$  **to**  $\gamma$  **do**

$q_i(x) \leftarrow M_q(prefix + [x_1, \dots, x_{i-1}])$

$x_i \sim q_i(x)$

**end for**

▷ **Run  $M_p$  in parallel.**

$p_1(x), \dots, p_{\gamma+1}(x) \leftarrow$   
     $M_p(prefix), \dots, M_p(prefix + [x_1, \dots, x_\gamma])$

▷ **Determine the number of accepted guesses  $n$ .**

$r_1 \sim U(0, 1), \dots, r_\gamma \sim U(0, 1)$

$n \leftarrow \min(\{i - 1 \mid 1 \leq i \leq \gamma, r_i > \frac{p_i(x)}{q_i(x)}\} \cup \{\gamma\})$

▷ **Adjust the distribution from  $M_p$  if needed.**

$p'(x) \leftarrow p_{n+1}(x)$

**if**  $n < \gamma$  **then**

$p'(x) \leftarrow \text{norm}(\max(0, p_{n+1}(x) - q_{n+1}(x)))$

**end if**

▷ **Return one token from  $M_p$ , and  $n$  tokens from  $M_q$ .**

$t \sim p'(x)$

**return**  $prefix + [x_1, \dots, x_n, t]$

---

Table 2. Empirical results for speeding up inference from a T5-XXL 11B model.

TASK	$M_q$	TEMP	$\gamma$	$\alpha$	SPEED
ENDE	T5-SMALL ★	0	7	0.75	<b>3.4X</b>
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	<b>2.6X</b>
ENDE	T5-BASE	1	5	0.68	2.4X
ENDE	T5-LARGE	1	3	0.71	1.4X
CNNNDM	T5-SMALL ★	0	5	0.65	<b>3.1X</b>
CNNNDM	T5-BASE	0	5	0.73	3.0X
CNNNDM	T5-LARGE	0	3	0.74	2.2X
CNNNDM	T5-SMALL ★	1	5	0.53	<b>2.3X</b>
CNNNDM	T5-BASE	1	3	0.55	2.2X
CNNNDM	T5-LARGE	1	3	0.56	1.7X

English to German translation fine tuned on WMT EnDe

Text summarization fine tuned on CCN/DM

# Speculative decoding

## Discussion

- Sampling unbiased
- Requires 2 models
- Only one guess verified

Table 2. Empirical results for speeding up inference from a T5-XXL 11B model.

TASK	$M_q$	TEMP	$\gamma$	$\alpha$	SPEED
ENDE	T5-SMALL ★	0	7	0.75	<b>3.4X</b>
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	<b>2.6X</b>
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	<b>3.1X</b>
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	<b>2.3X</b>
CNNDM	T5-BASE	1	3	0.55	2.2X
CNNDM	T5-LARGE	1	3	0.56	1.7X

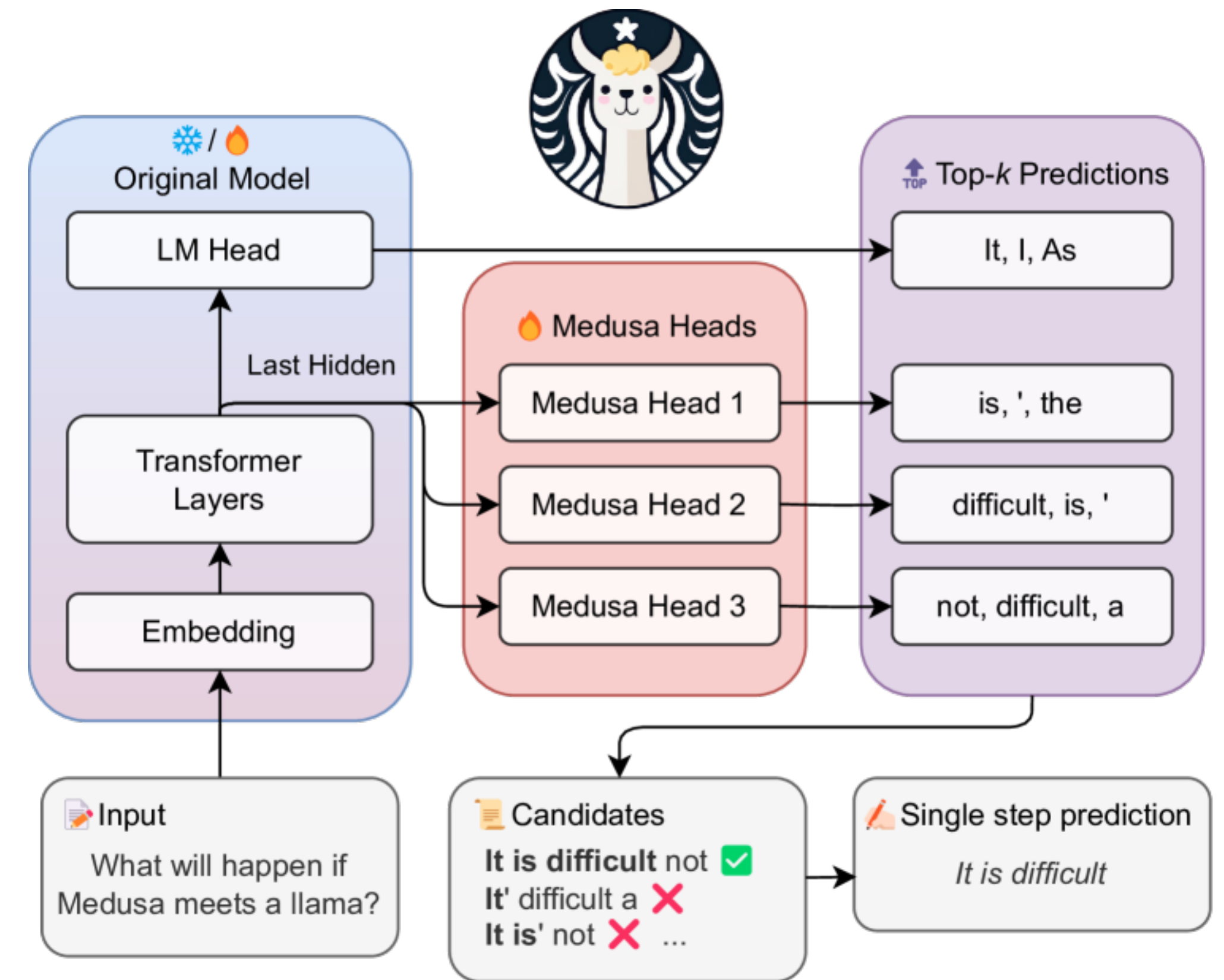
English to German translation fine tuned on WMT EnDe

Text summarization fine tuned on CCN/DM



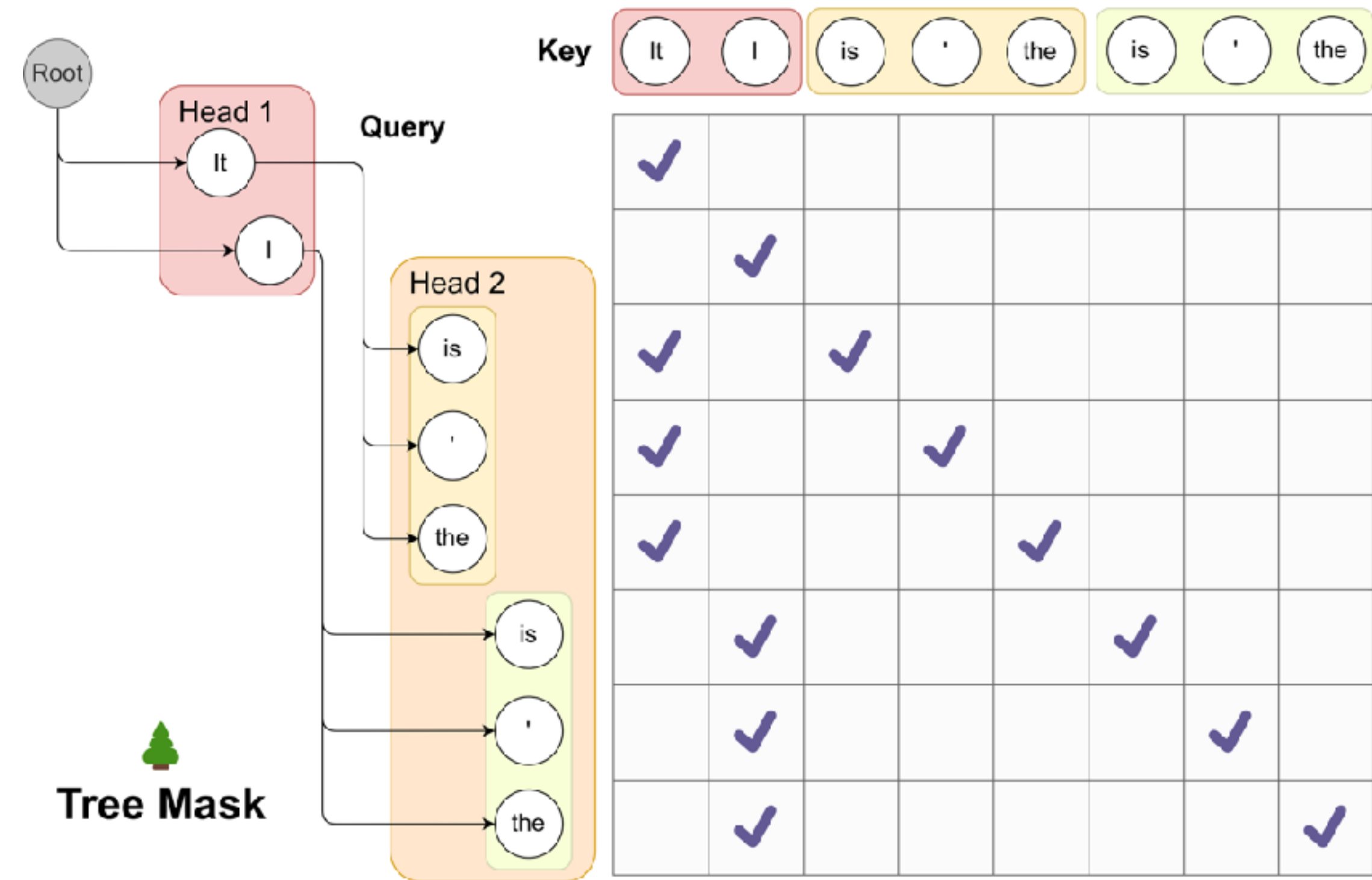
# Medusa

- Speculative decoding plus
  - Verifying many guesses at once
    - Tree attention
- Use same model to produce guesses and verify
  - Medusa Head



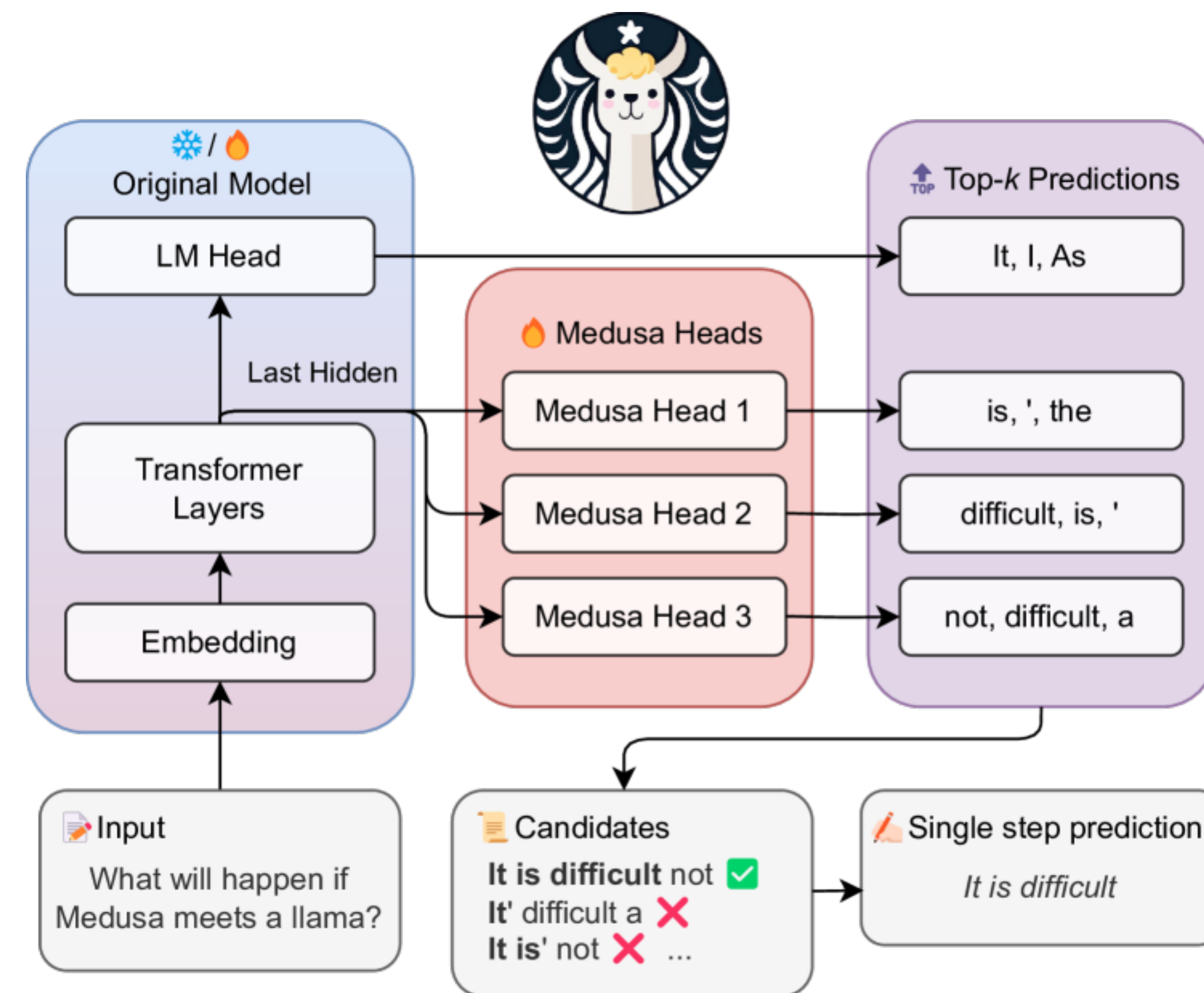
# Tree Attention

- Verifying many (tree structured ) guesses at once
- Example:
  - It is, It', it the, I is, I', I the
- Mask attention to go back along tree
  - Tree Mask
  - Use FlexAttention

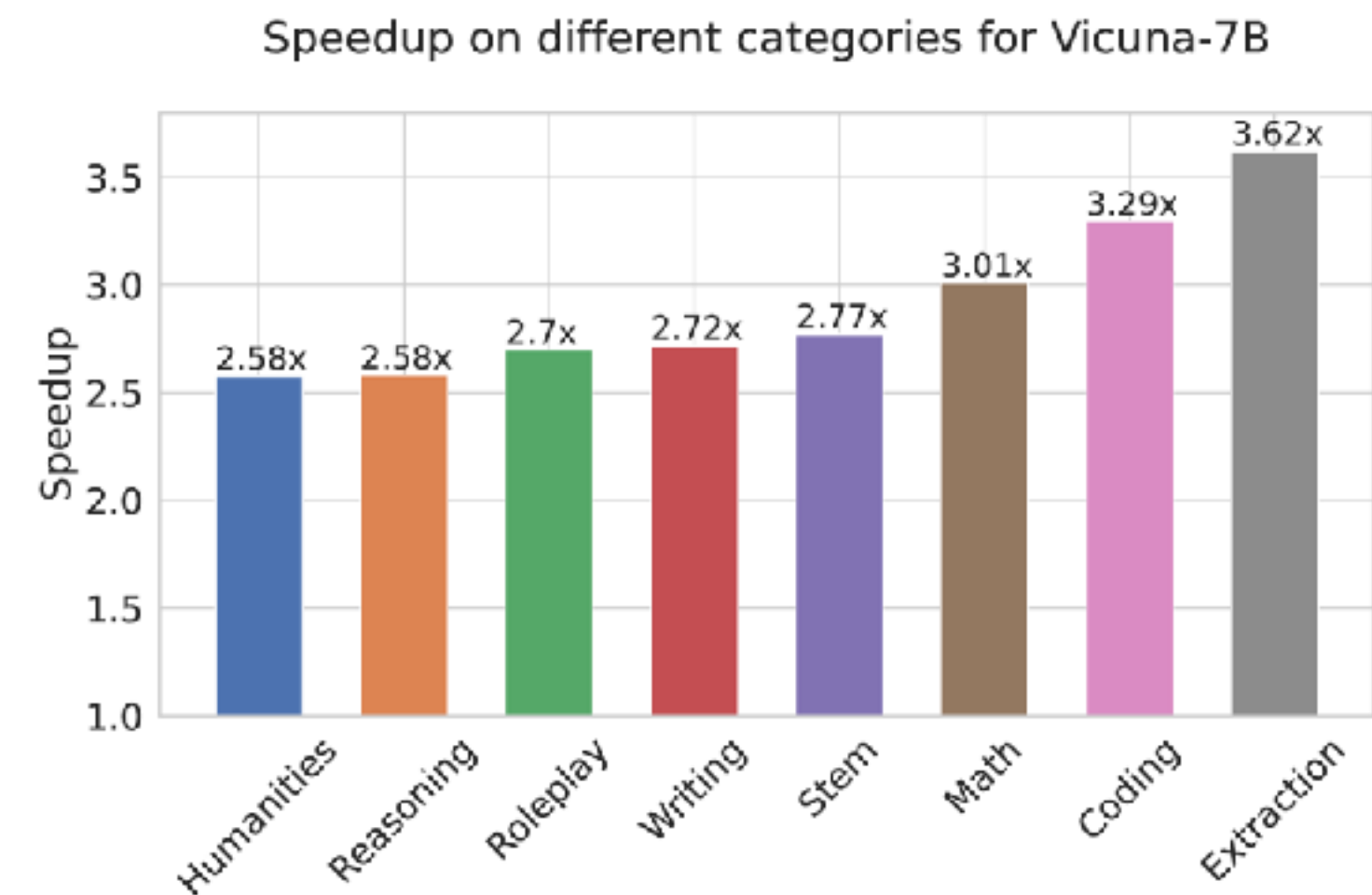
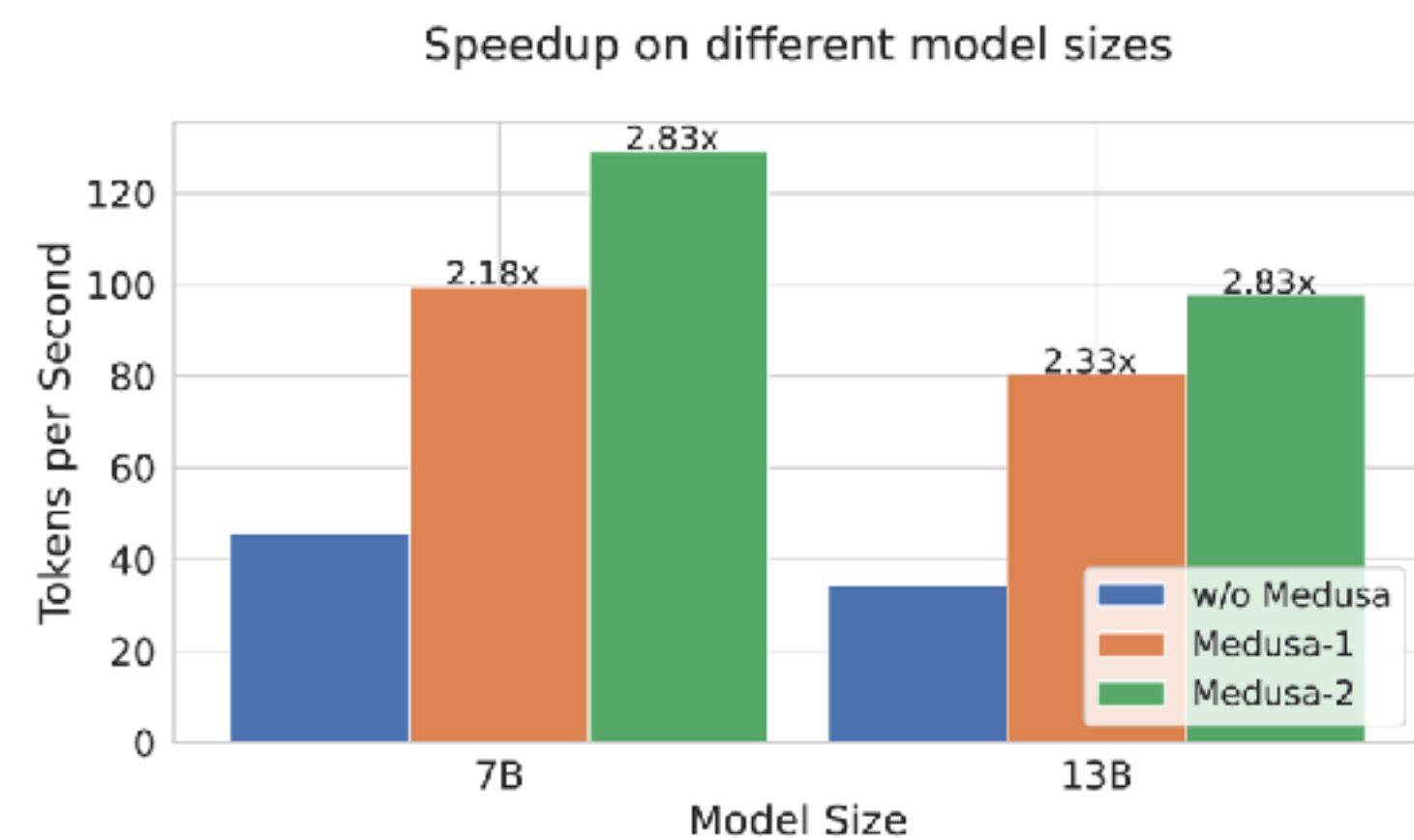


# Medusa Heads

- K heads  $q^{(1)} \dots q^{(k)}$  produce independent probability of k-th next word
- Chose top  $s_1 \dots s_k$  per node
- **Try all combinations**

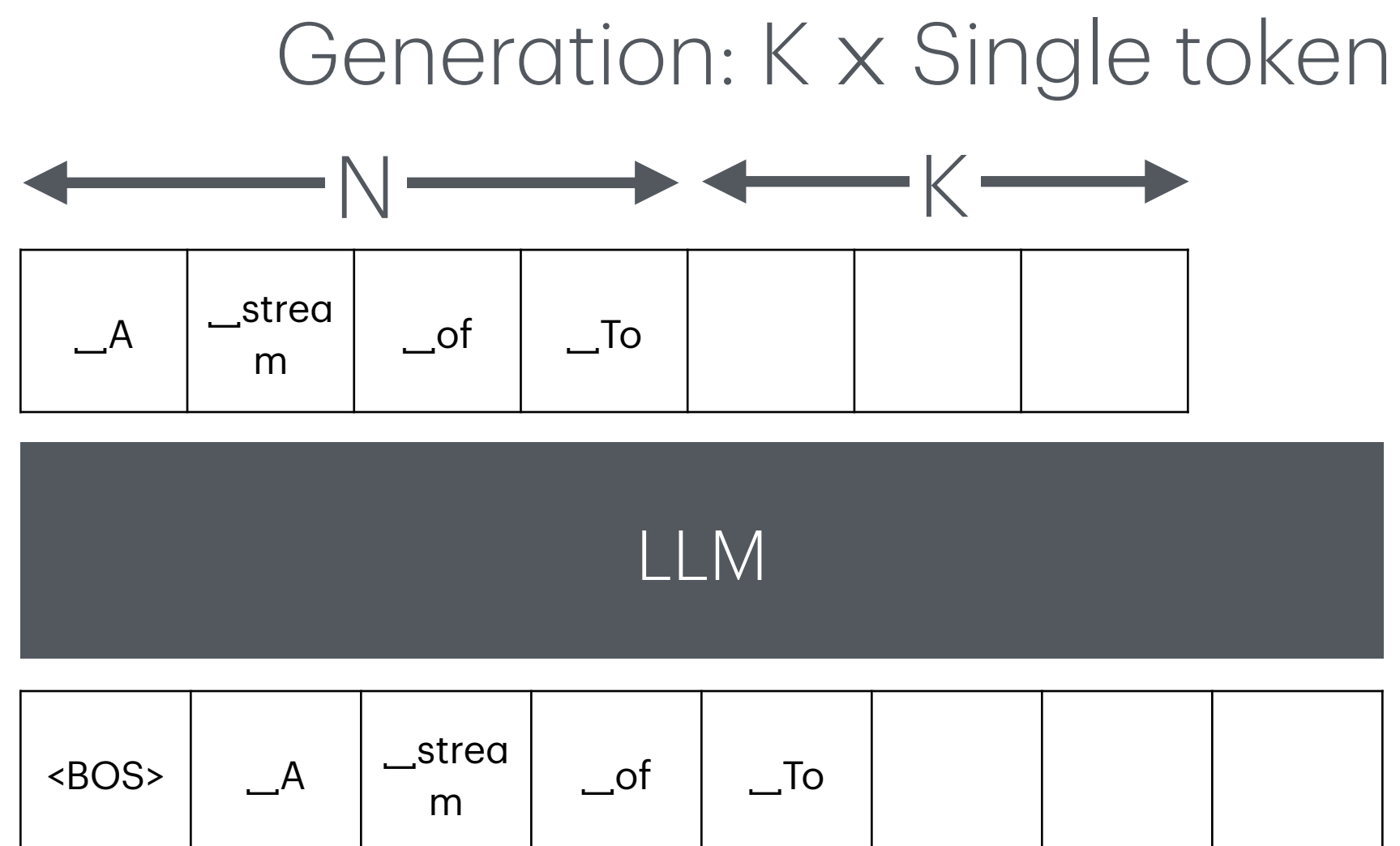


# Medusa - Results



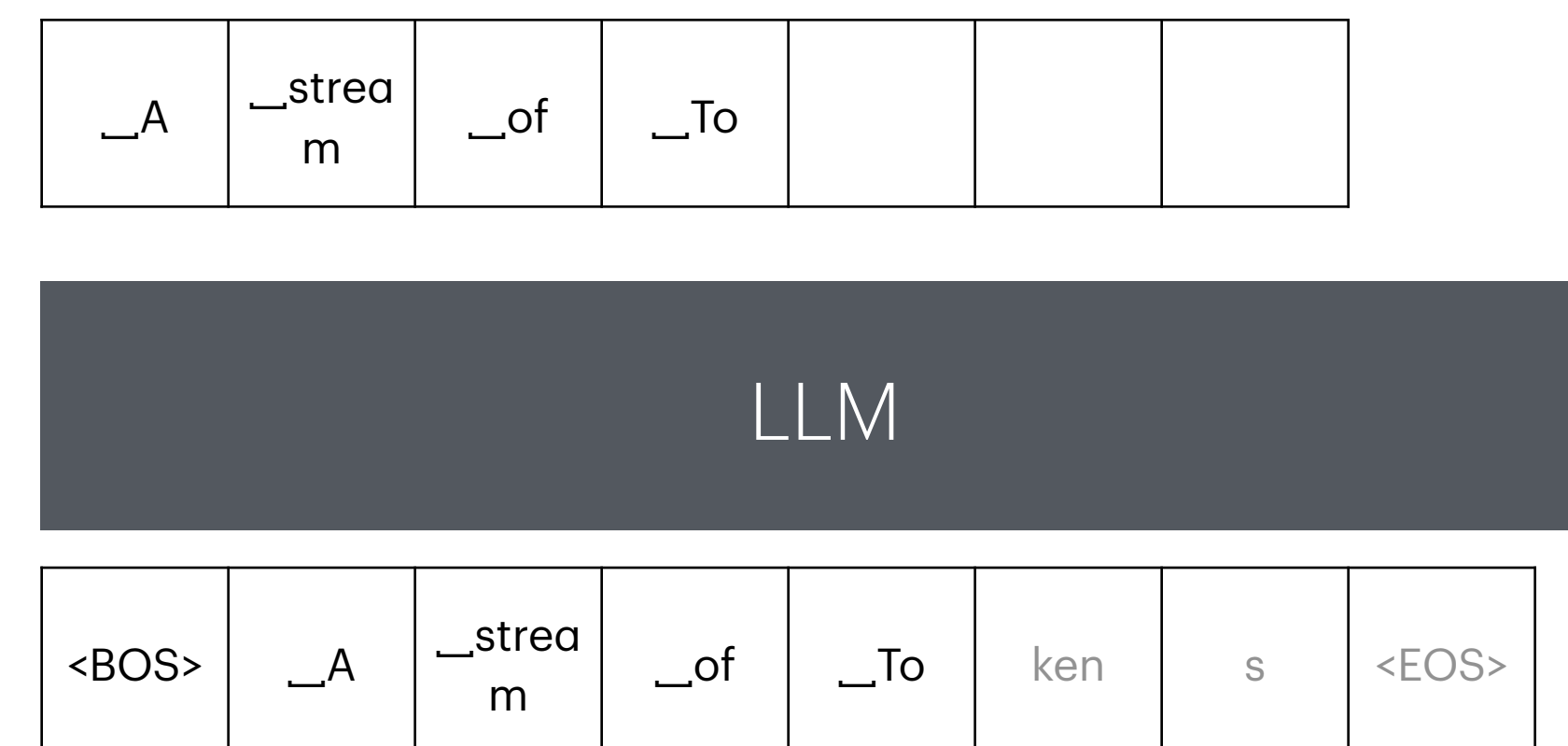
# Speculative Decoding vs KV-Cache

Without KV-Cache



$$K \times O(N^2) = O(KN^2)$$

Generation: K tokens

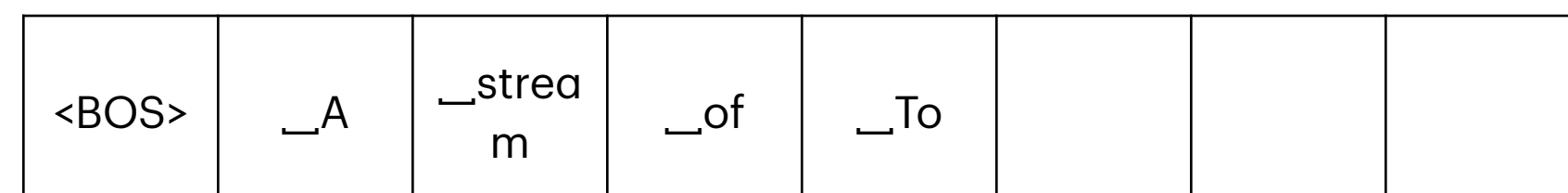
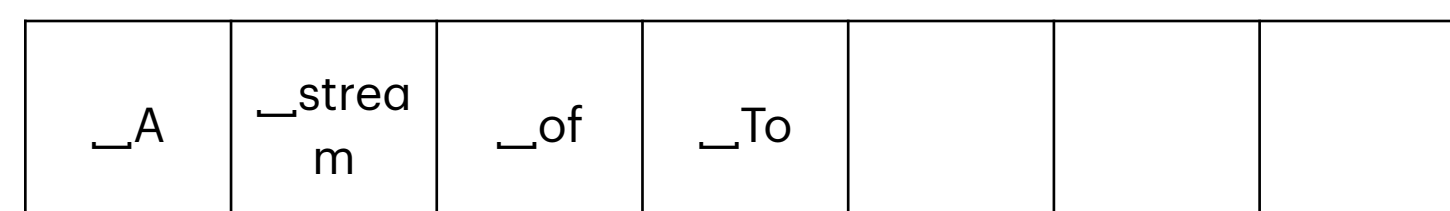


$$O((N + K)^2) \approx O(N^2) + O(NK)$$

# Speculative Decoding vs KV-Cache

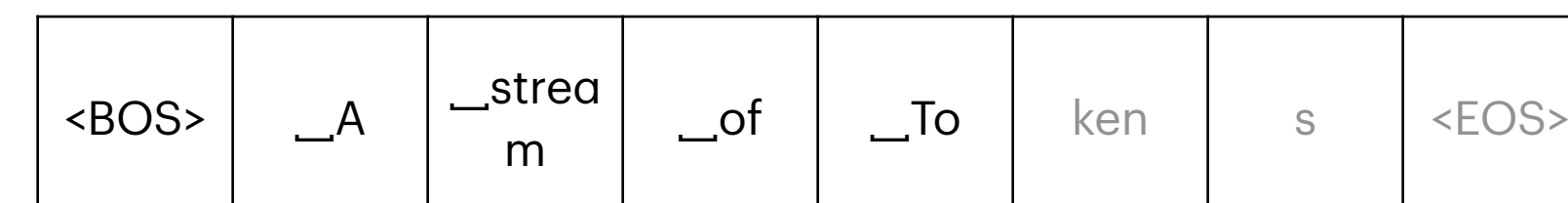
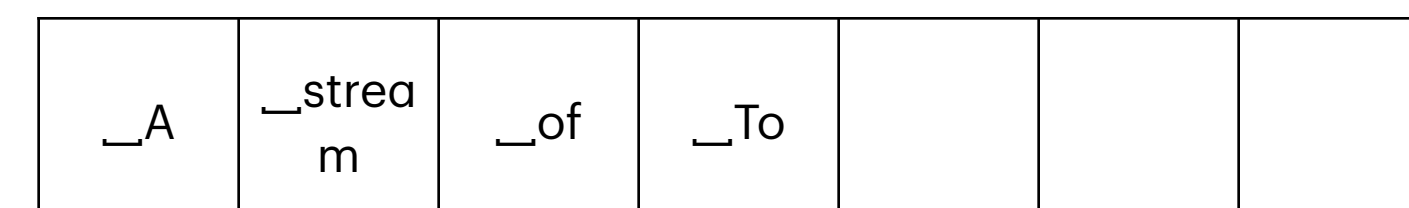
With KV-Cache

Generation: K x Single token



$$K \times O(N) = O(KN)$$

Generation: K tokens

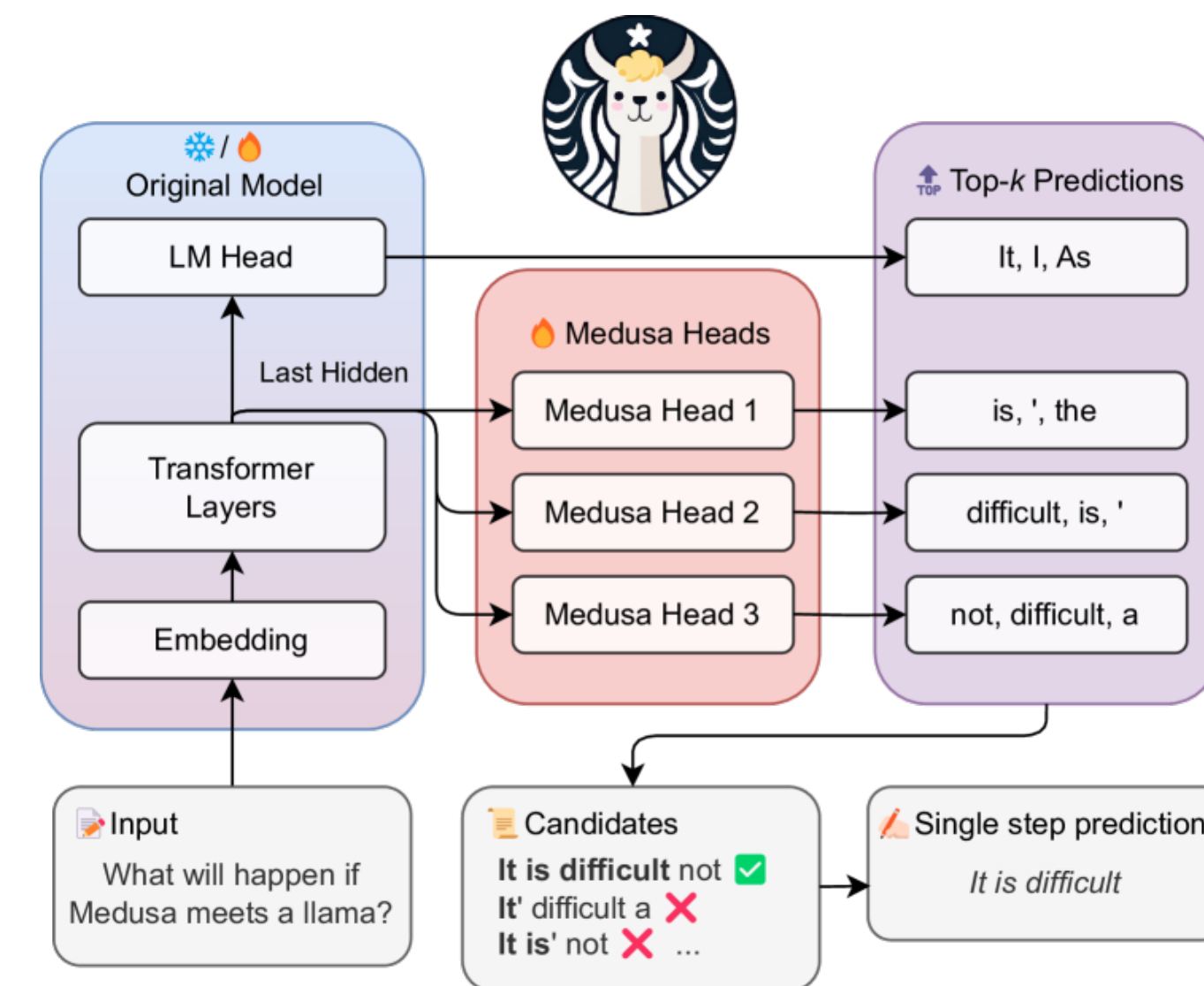


$$O(KN)$$

Gains: K vs 1 forward calls

# Speculative decoding

- Beautiful technical insights
- Gains diminish with KV-Cache
- Quite tricky to implement for batched inference




---

## Algorithm 1 SpeculativeDecodingStep

---

**Inputs:**  $M_p, M_q, prefix$ .

▷ Sample  $\gamma$  guesses  $x_1, \dots, x_\gamma$  from  $M_q$  autoregressively.

**for**  $i = 1$  **to**  $\gamma$  **do**

$q_i(x) \leftarrow M_q(prefix + [x_1, \dots, x_{i-1}])$

$x_i \sim q_i(x)$

**end for**

▷ Run  $M_p$  in parallel.

$p_1(x), \dots, p_{\gamma+1}(x) \leftarrow M_p(prefix), \dots, M_p(prefix + [x_1, \dots, x_\gamma])$

▷ Determine the number of accepted guesses  $n$ .

$r_1 \sim U(0, 1), \dots, r_\gamma \sim U(0, 1)$

$n \leftarrow \min(\{i - 1 \mid 1 \leq i \leq \gamma, r_i > \frac{p_i(x)}{q_i(x)}\} \cup \{\gamma\})$

▷ Adjust the distribution from  $M_p$  if needed.

$p'(x) \leftarrow p_{n+1}(x)$

**if**  $n < \gamma$  **then**

$p'(x) \leftarrow \text{norm}(\max(0, p_{n+1}(x) - q_{n+1}(x)))$

**end if**

▷ Return one token from  $M_p$ , and  $n$  tokens from  $M_q$ .

$t \sim p'(x)$

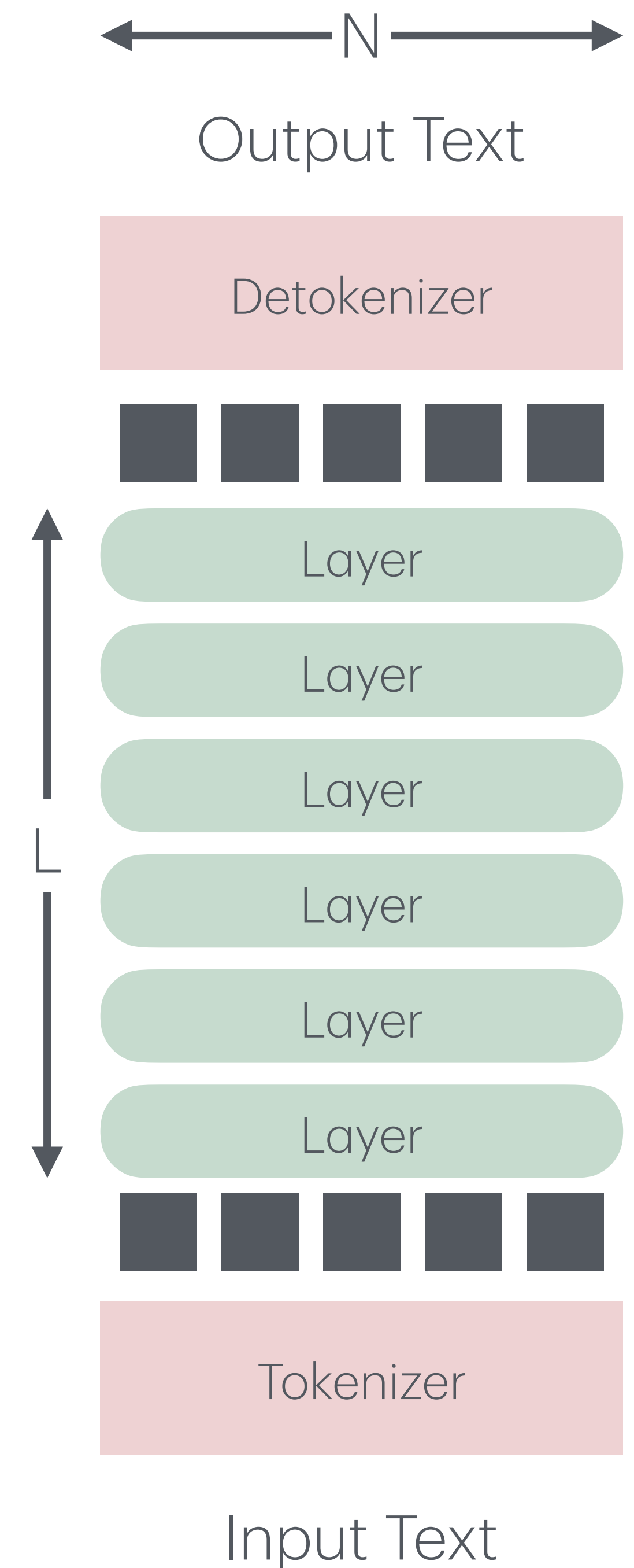
**return**  $prefix + [x_1, \dots, x_n, t]$

---

# Training and Generation

## Paged Attention

	Training	Training - Checkpointi	Generation	Paged Attention	Speculative decoding
Peak Memory	$O(NL)$	$O(NL^{1/2})$	$O(N)$	$O(NL)$	$O(NL)$
Runtime	$O(N^2L)$	$O(2 N^2L)$	$O(N^3L)$	$O(N^2L)$	$O(N^2L)$
# forward	1	1	N	N	$N / \alpha$





# References

- [1] Fast Inference from Transformers via Speculative Decoding, Levianthan et al 2023. ([link](#))
- [2] Medusa: Simple LLM Inference Acceleration Framework with Multiple Decoding Heads, Cai et al 2024 ([link](#))