

Stanford | **ENGINEERING**

Lecture 3

CME 295: Transformers & Large Language Models



Afshine Amidi & Shervine Amidi



Stanford University

ICME



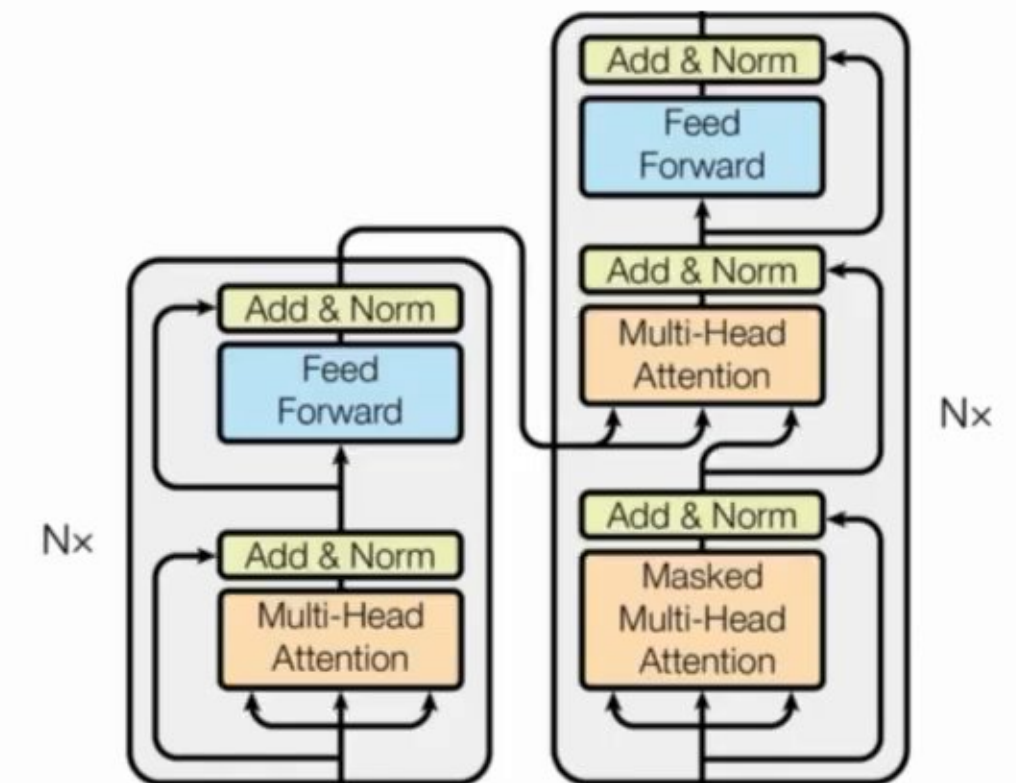
Stanford

Recap of last episode...



3 categories of Transformer-based models:

Encoder-decoder	Text to text	T5, mT5, ByT5
------------------------	--------------	---------------

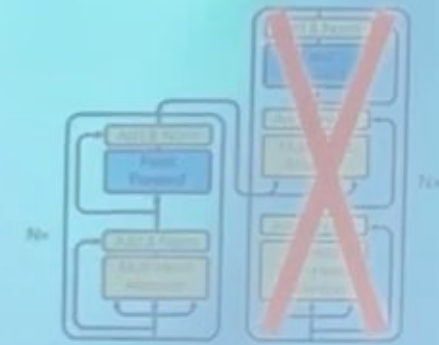


Stanford

Recap of last episode...

3 categories of Transformer-based models:

Encoder-Decoder	Encoder-only	Decoder-only			
	<table border="1"><tr><td>Encoder-only</td><td>Projection of embedding for class prediction (e.g. sentiment extraction)</td><td>BERT, DistilBERT, RoBERTa</td></tr></table>	Encoder-only	Projection of embedding for class prediction (e.g. sentiment extraction)	BERT, DistilBERT, RoBERTa	
Encoder-only	Projection of embedding for class prediction (e.g. sentiment extraction)	BERT, DistilBERT, RoBERTa			



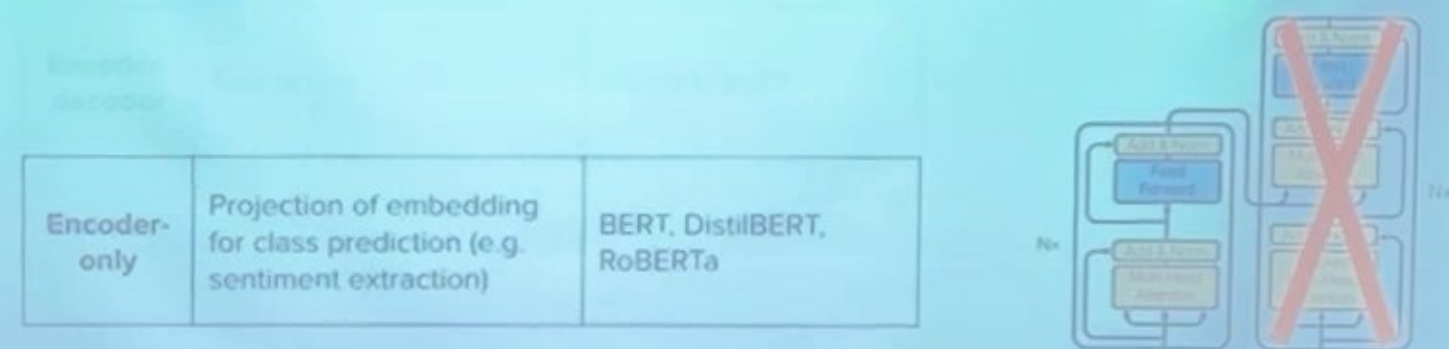
Stanford University Figure adapted from "GPT-3: Generative Pretrained Transformer 3" (OpenAI, 2020)

ICME

Stanford

Recap of last episode...

3 categories of Transformer-based models:



Stanford University

Figure adapted from "Attention is All You Need", Vaswani et al., 2017

ICME

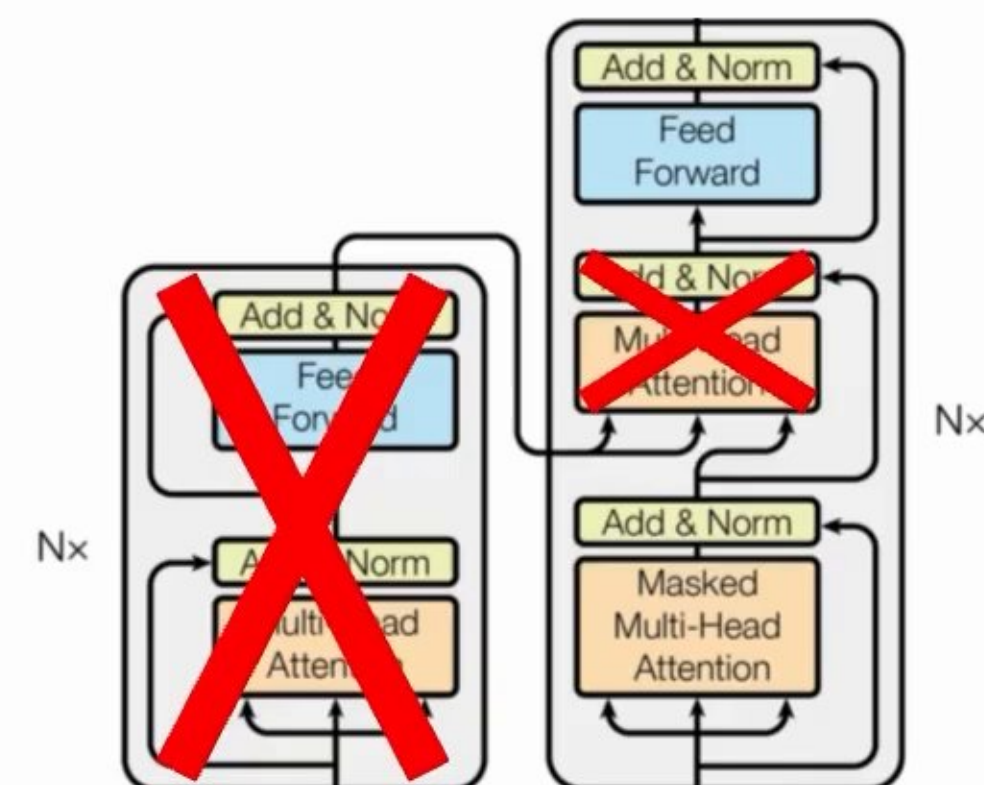
Stanford

Recap of last episode...



3 categories of Transformer-based models:

Encoder-decoder	Text to text	T5, mT5, ByT5
Encoder-only	Projection of embedding for class prediction (e.g. sentiment extraction)	BERT, DistilBERT, RoBERTa
Decoder-only	Text to text	GPT series



Stanford

Recap of last episode...



3 categories of Transformer-based models:

Encoder-decoder	Text to text	T5, mT5, ByT5
Encoder-only	Projection of embedding for class prediction (e.g. sentiment extraction)	BERT, DistilBERT, RoBERTa
Decoder-only	Text to text	GPT series

Stanford

Terminology



LLM = **L**arge **L**anguage **M**odel

Stanford

Terminology

LUM = L_{ing}u_{age}

Language Model

"A **language model** is a statistical or machine learning **model** that assigns **probabilities** to sequences of **tokens**."

Stanford University

ICME

Stanford

Terminology



LLM = **Large** Language Model

Stanford

Terminology

LLM **Large** Language Model

- **Model size:** billions of parameters or more
- **Training data:** 100s of billions of tokens or more
- **Compute:** a lot of GPUs

Stanford University

ICME

Stanford

Terminology

LLM **Large** Language Model

- **Model size:** billions of parameters or more
- **Training data:** 100s of billions of tokens or more
- **Compute:** a lot of GPUs

Stanford University

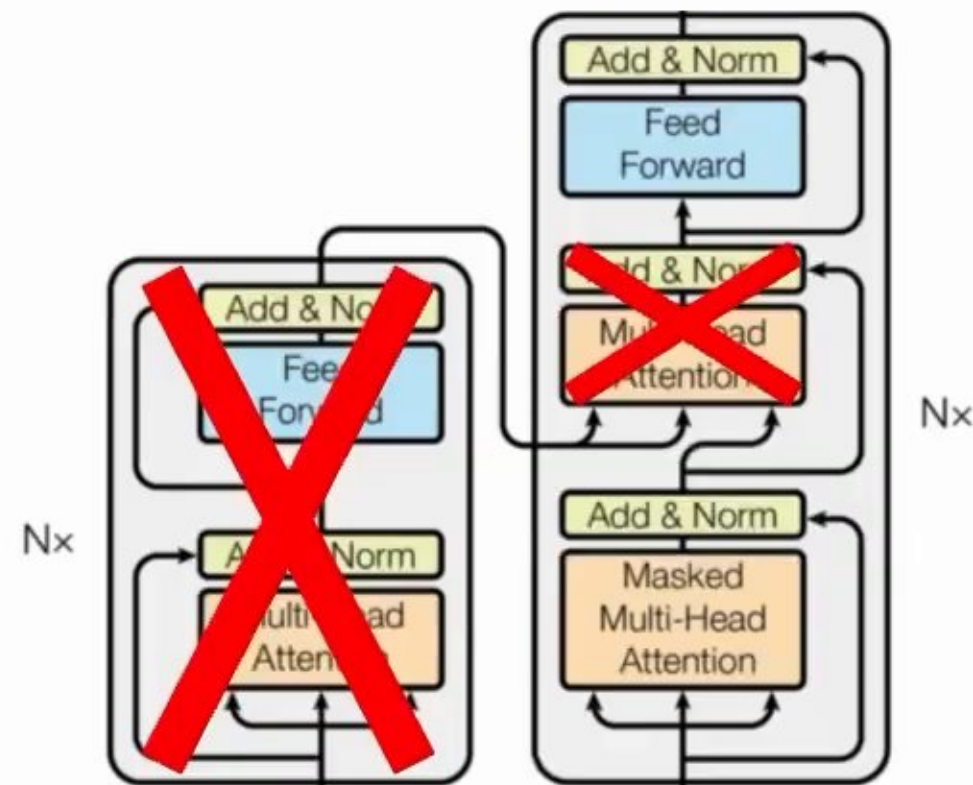
ICME

Stanford

Characteristics



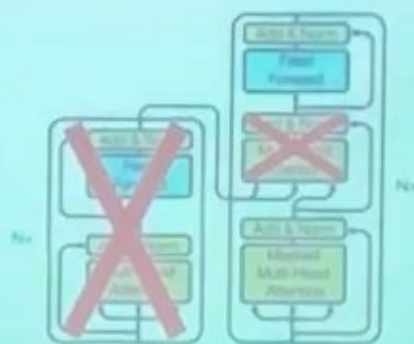
Decoder-only Transformer-based model



Stanford

Characteristics

Decoder-only Transformer-based model



Stanford University Figure adapted from "Attention is All You Need", Vaswani et al., 2017

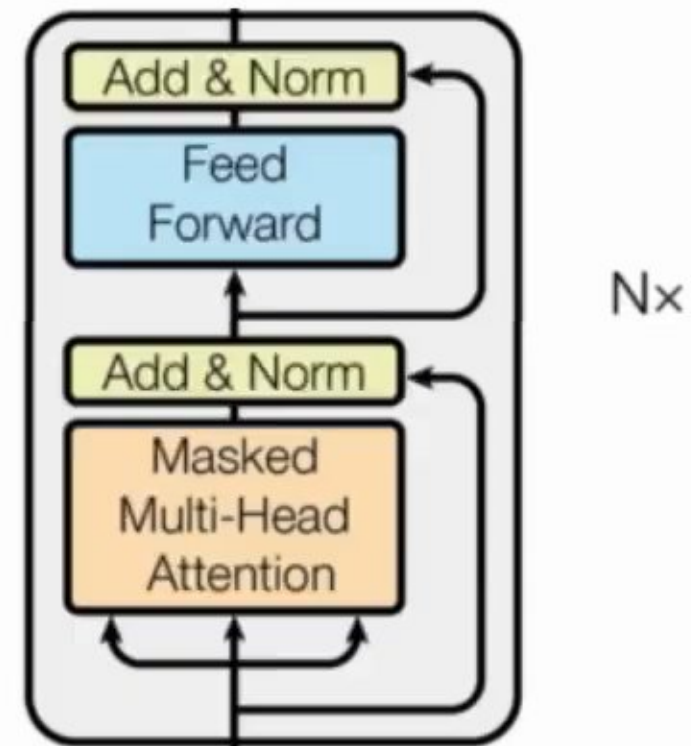
ICME

Stanford

Characteristics



Decoder-only Transformer-based model



Stanford



Transformers & Large Language Models



LLM overview

MoE-based LLMs

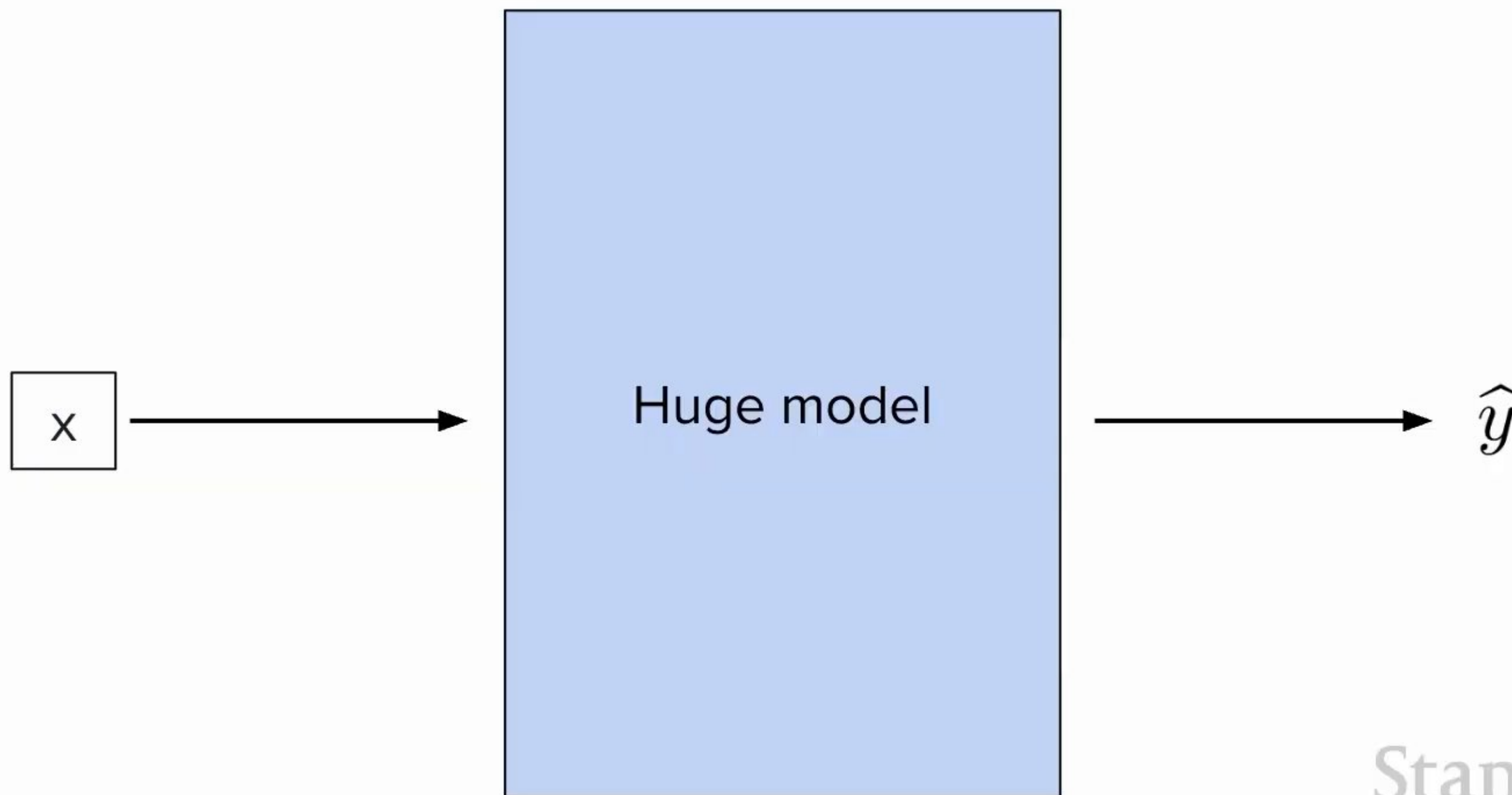
Response generation

Prompting strategies

Inference optimizations

Stanford

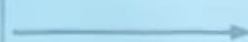
Motivation



Stanford

Motivation

x



Huge model



\hat{y}

Stanford University

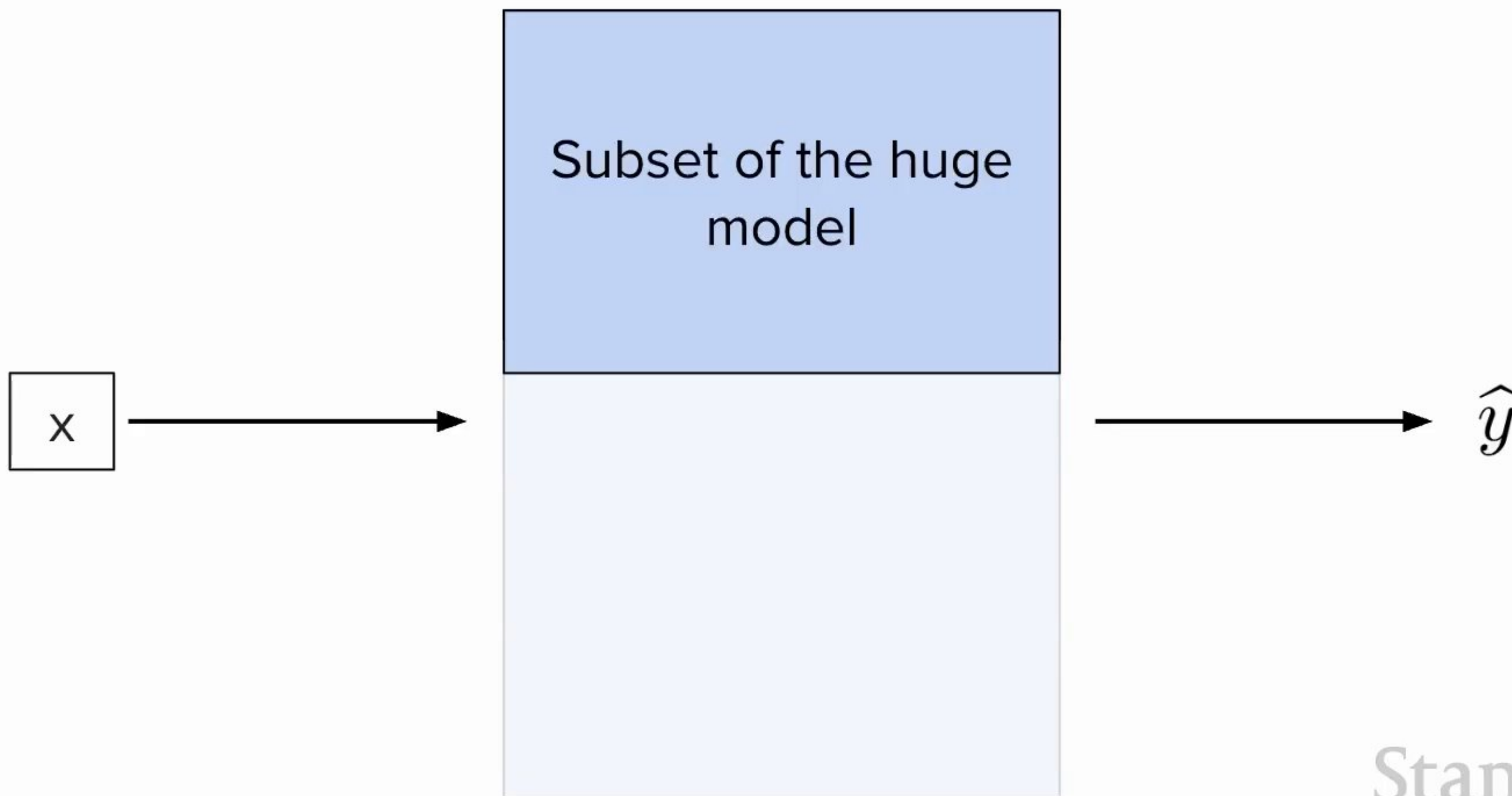
ICML

Stanford

Motivation

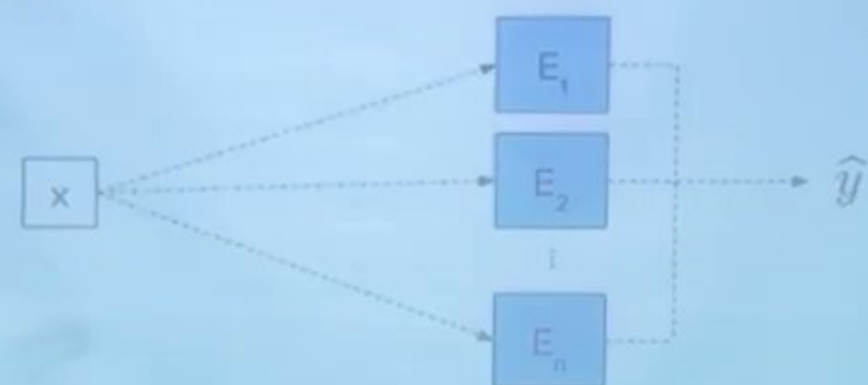


Idea. Not all weights are useful in the forward pass



Stanford

Motivation



Stanford University

CS-109

Stanford

Motivation



Stanford University

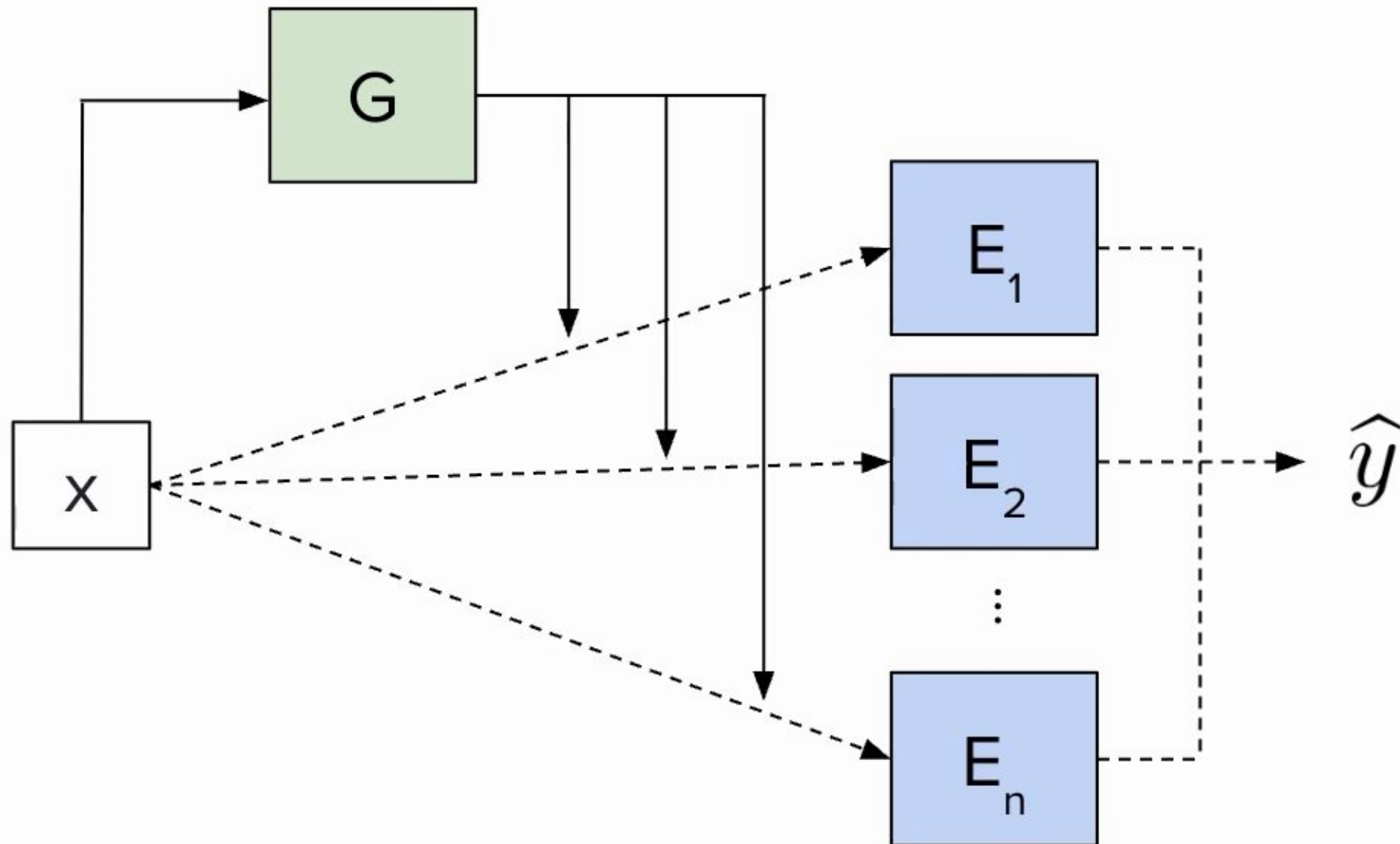
ICML

Stanford

Overview of MoEs



MoE = Mixture of Experts

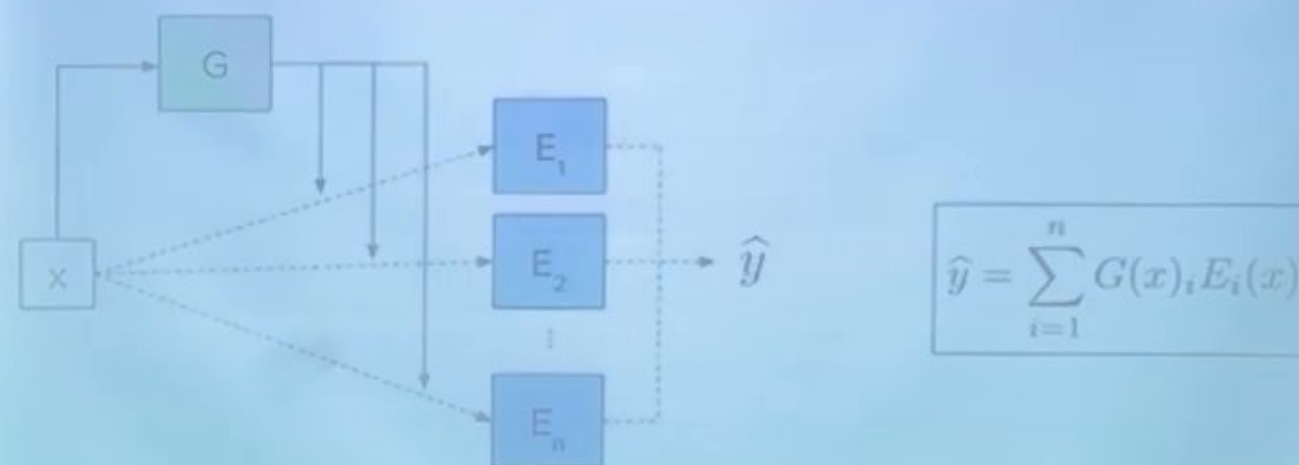


$$\hat{y} = \sum_{i=1}^n G(x)_i E_i(x)$$

Stanford

Overview of MoEs

MoE = Mixture of Experts

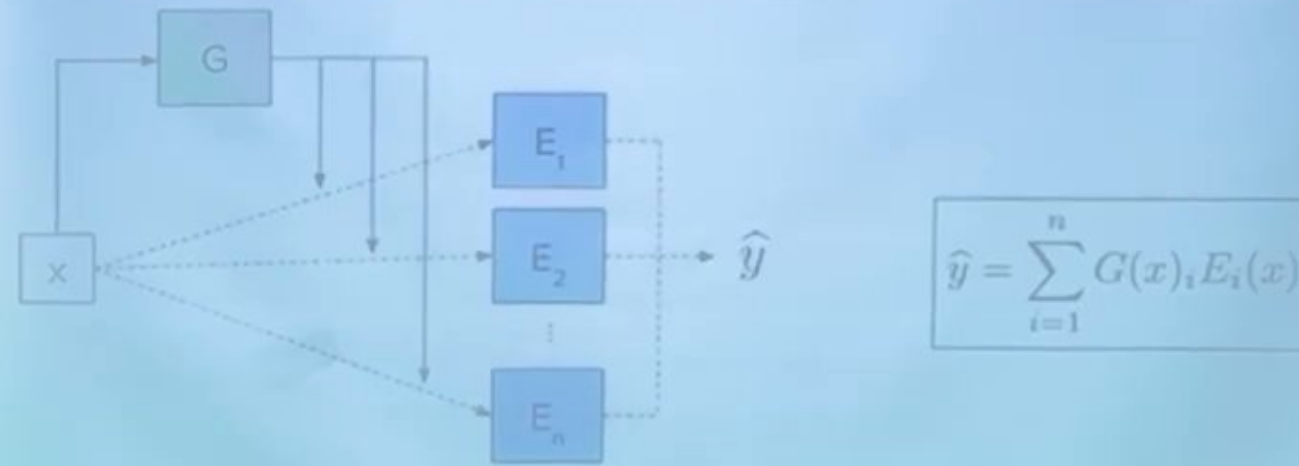


Stanford University

ICML'16

Overview of MoEs

MoE = Mixture of Experts



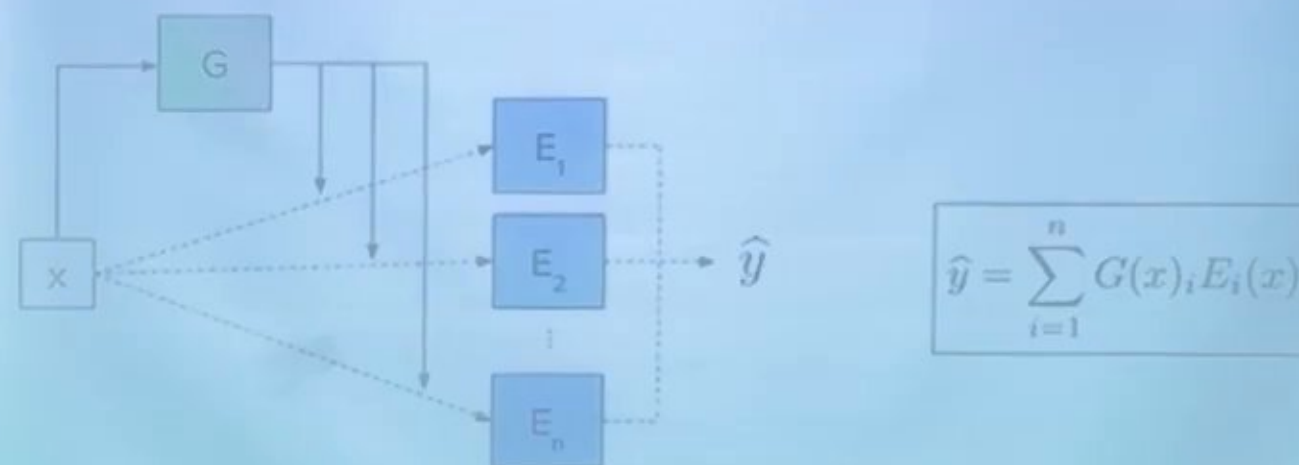
Stanford University

ICML'22

Stanford

Overview of MoEs

MoE = Mixture of Experts



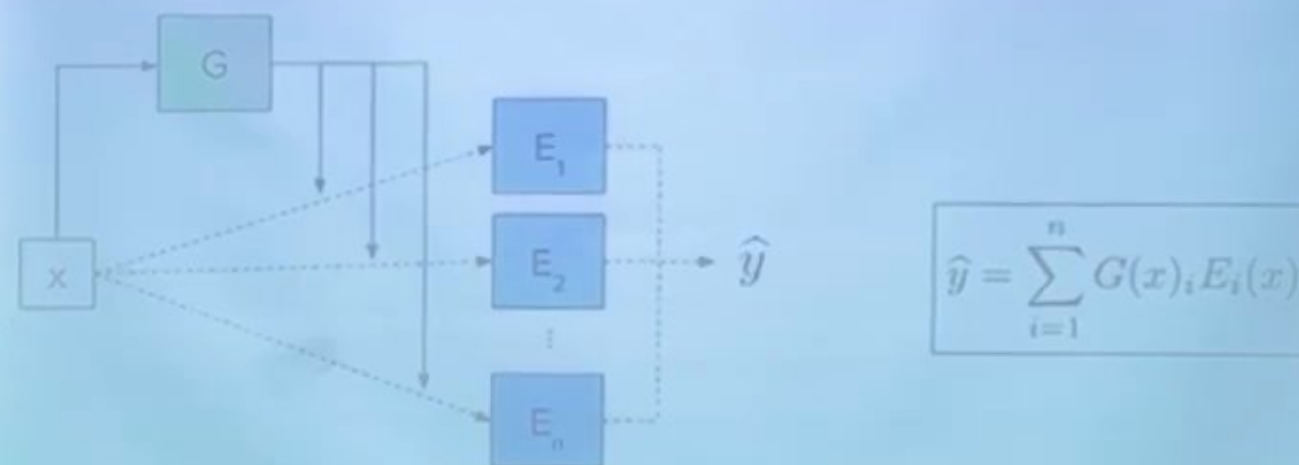
Stanford University

ICME

Stanford

Overview of MoEs

MoE = Mixture of Experts



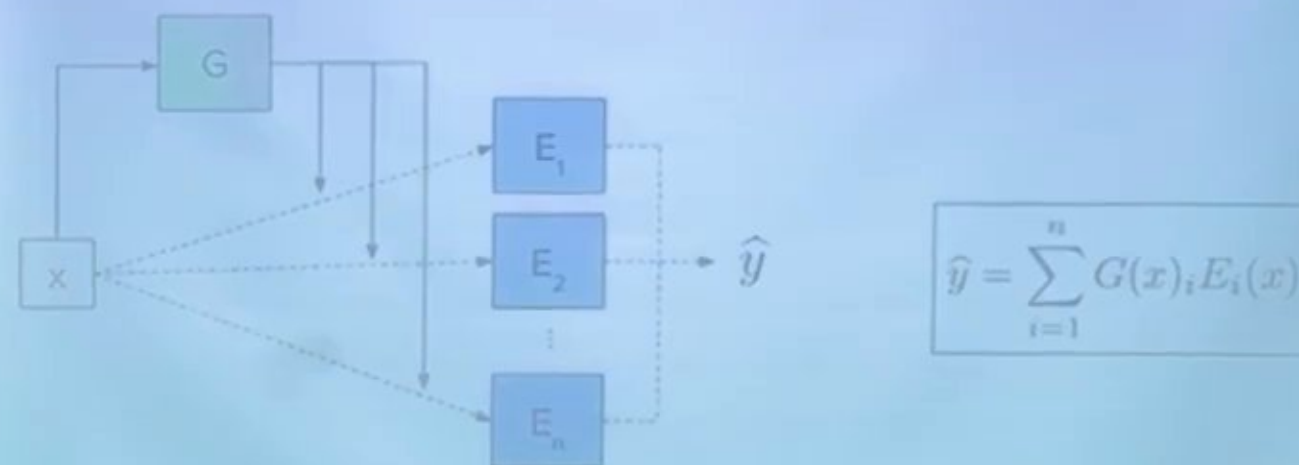
Stanford University

ICML

Stanford

Overview of MoEs

MoE = Mixture of Experts



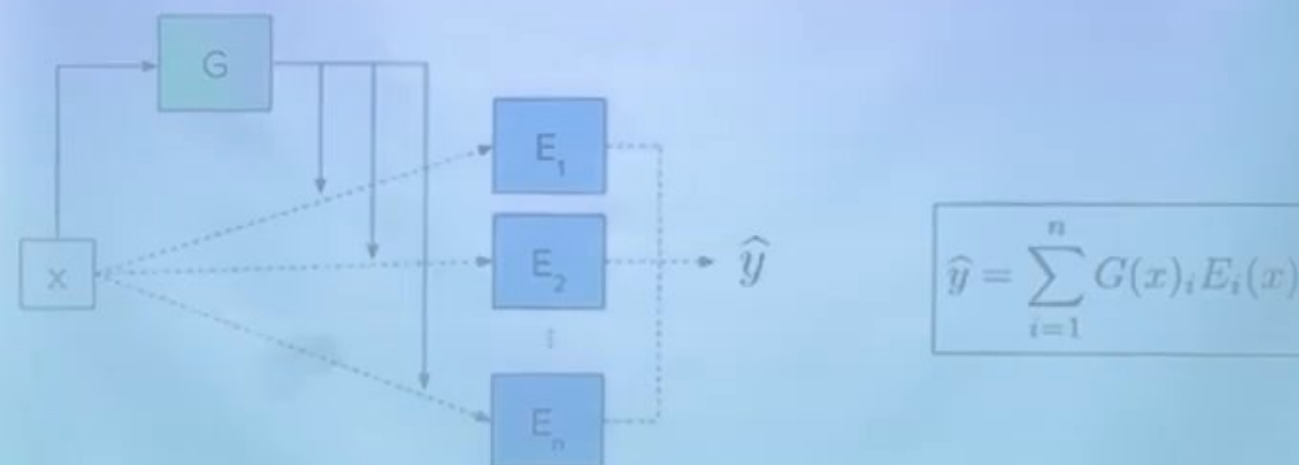
Stanford University

ICML

Stanford

Overview of MoEs

MoE = Mixture of Experts



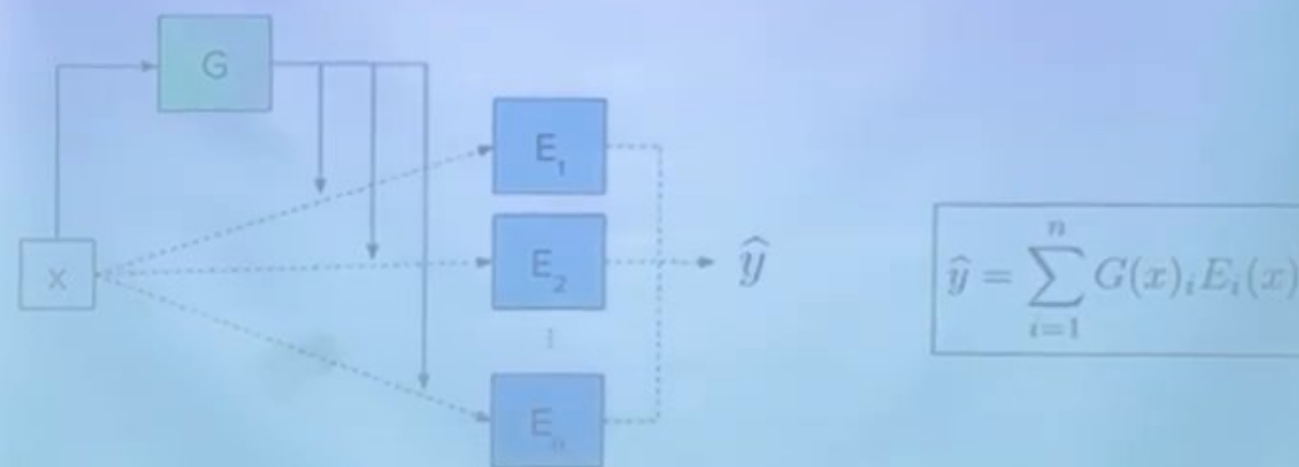
Stanford University

ICML

Stanford

Overview of MoEs

MoE = Mixture of Experts



Stanford University

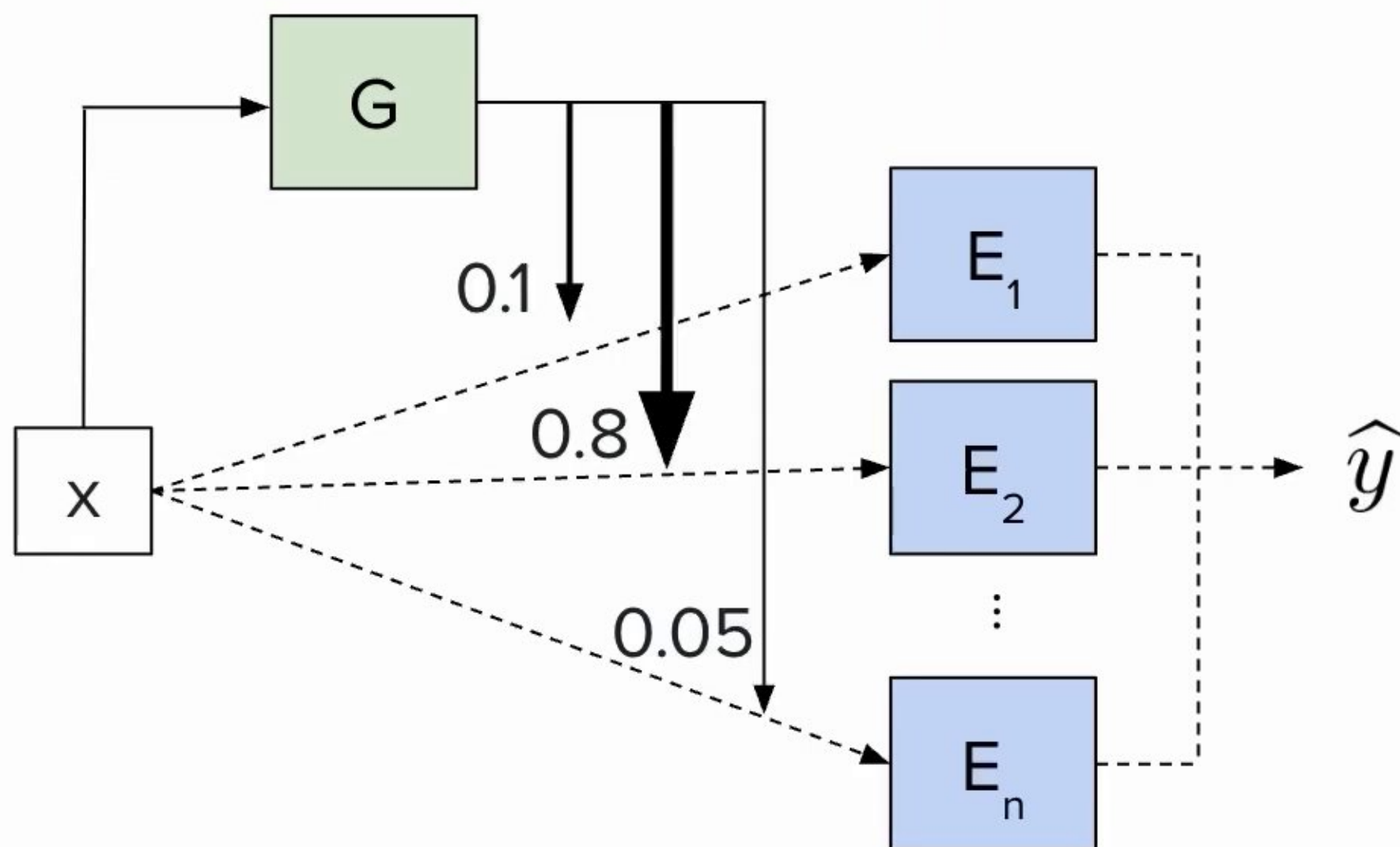
ICML

Stanford

Overview of MoEs



MoE = Mixture of Experts



Dense MoE. Output is weighted average of **all** expert outputs.

$$\hat{y} = \sum_{i=1}^n G(x)_i E_i(x)$$

Stanford

Overview of MoEs

MoE = Mixture of Experts



Dense MoE. Output is weighted average of all expert outputs.

$$\hat{y} = \sum_{i=1}^n G(x)_i E_i(x)$$

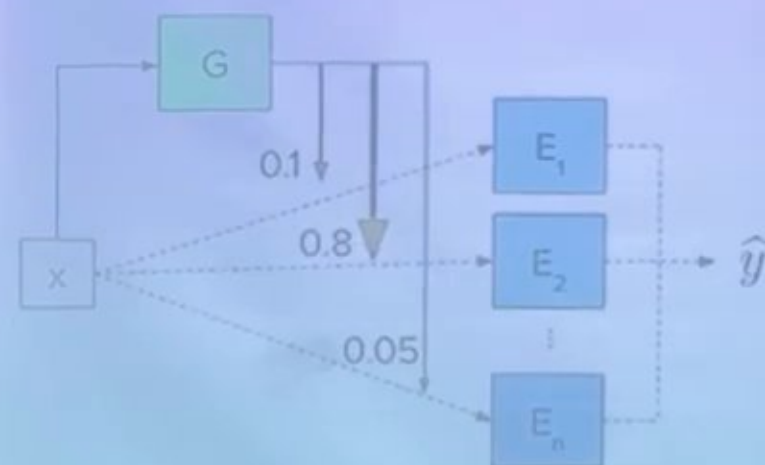
Stanford University

ICML

Stanford

Overview of MoEs

MoE = Mixture of Experts



Dense MoE. Output is weighted average of all expert outputs.

$$\hat{y} = \sum_{i=1}^n G(x)_i E_i(x)$$

Stanford University

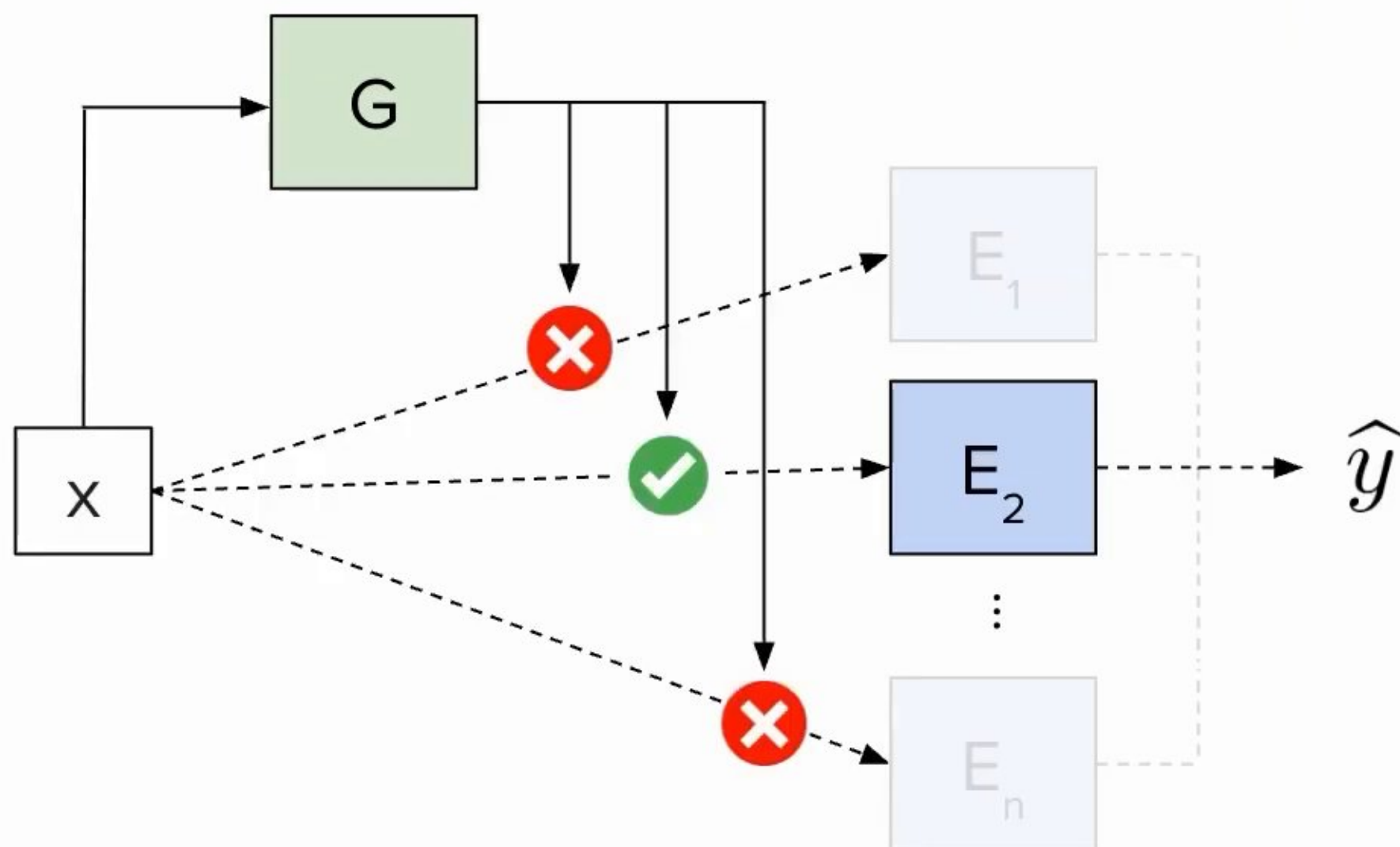
TCME

Stanford

Overview of MoEs



MoE = Mixture of Experts



Sparse MoE. Output is weighted average of **selected** expert outputs.

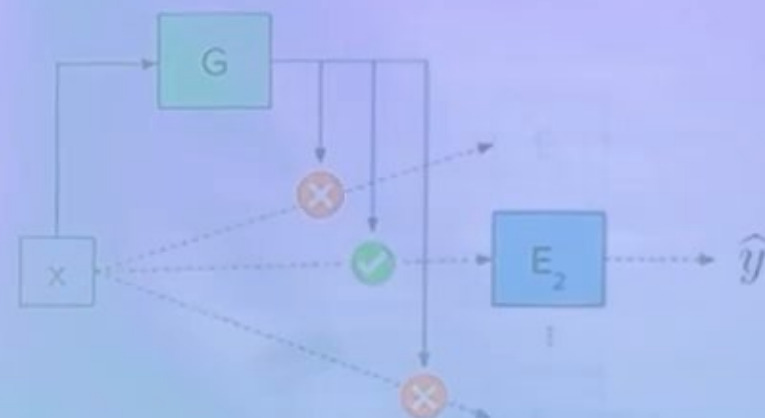
$$\hat{y} = \sum_{i \in \mathcal{I}_k} G(x)_i E_i(x)$$

Via **top-k** selection

Stanford

Overview of MoEs

MoE = Mixture of Experts



Sparse MoE. Output is weighted average of **selected** expert outputs.

$$\hat{y} = \sum_{i \in I_k} G(x)_i E_i(x)$$

Via top-k selection

Stanford University

"Extremely Large Neural Networks: The Sparse-Gated Mixture-of-Experts Layer", Shazeer et al., 2017

ICML

Stanford

Overview of MoEs

MoE = Mixture of Experts



Sparse MoE. Output is weighted average of **selected** expert outputs.

$$\hat{y} = \sum_{i \in T_k} G(x)_i E_i(x)$$

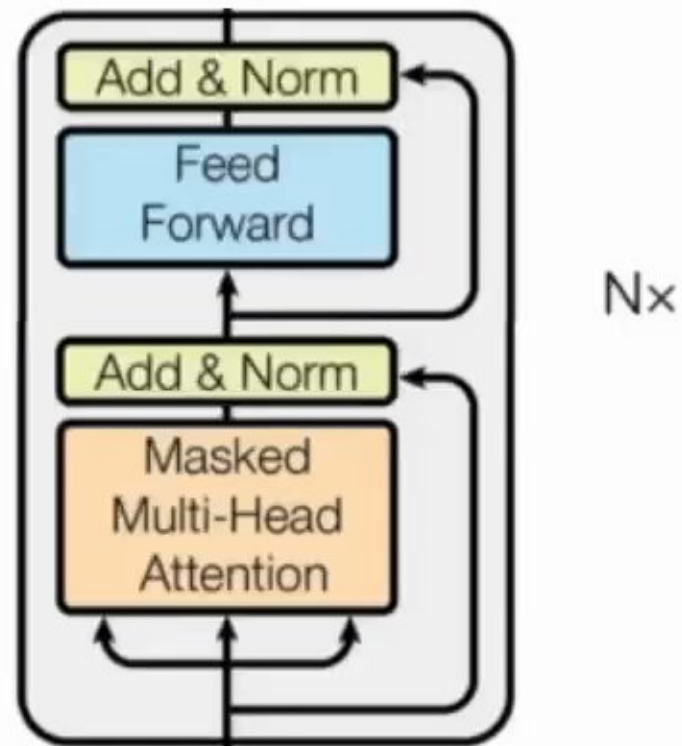
Via top-k selection

Stanford University "Continuously Gating Neural Networks: The Sparsity-Gated Mixture-of-Experts Layer", Shazeer et al., 2017

ICML

Stanford

MoE in Transformer-based models



Stanford

MoE in Transformer-based models



Stanford University Figure captioned from "Attention is All You Need", Vaswani et al., 2017

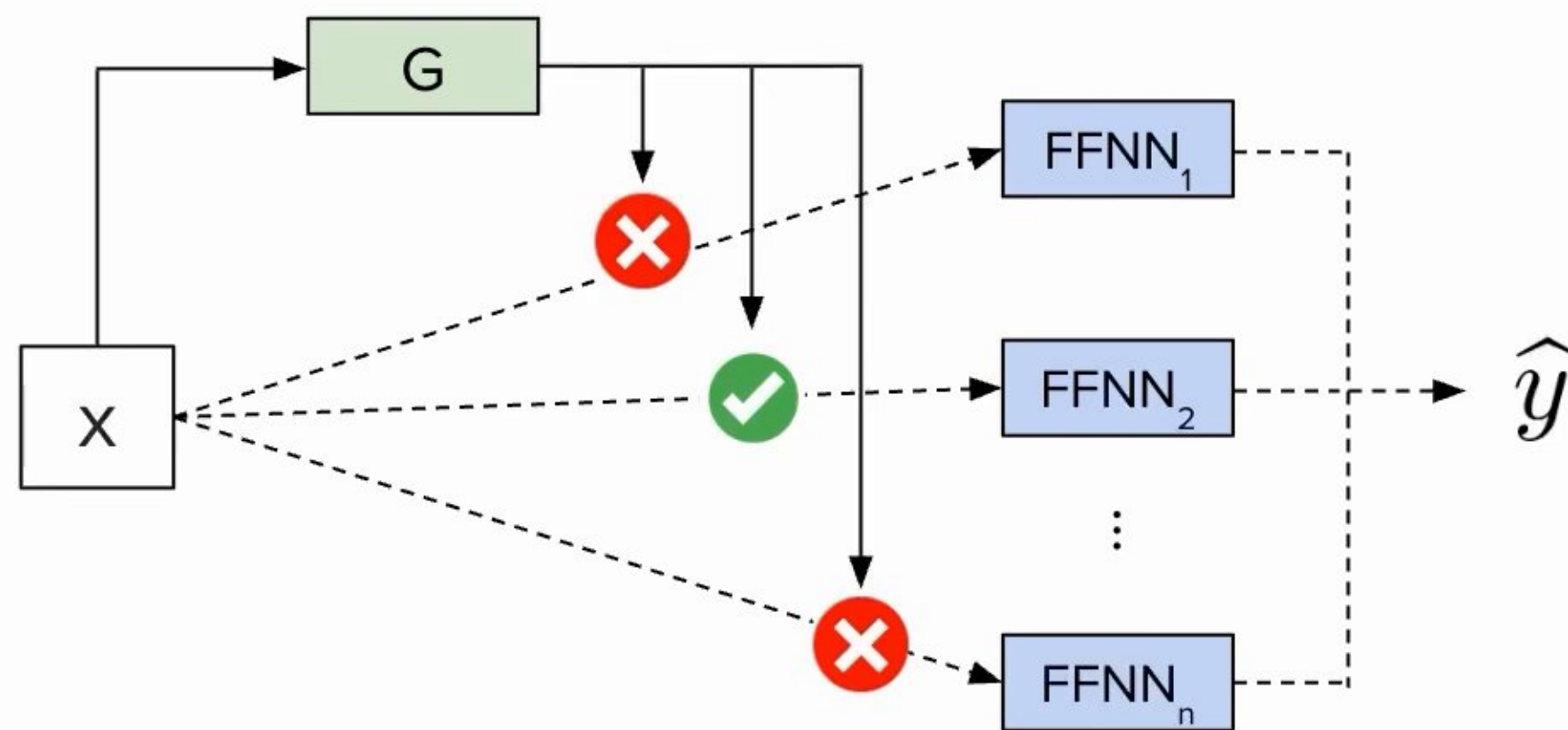
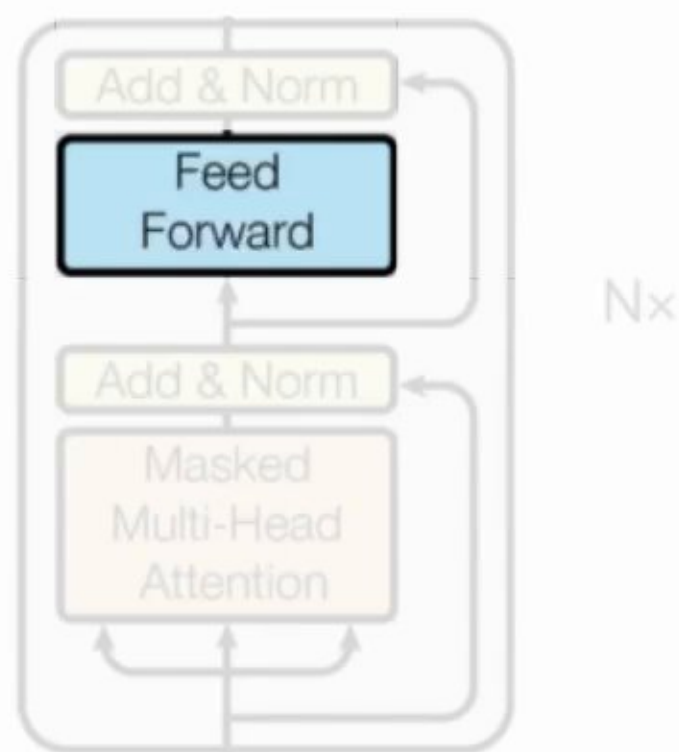
ICML

Stanford

MoE in Transformer-based models



Routing done **for each token!**

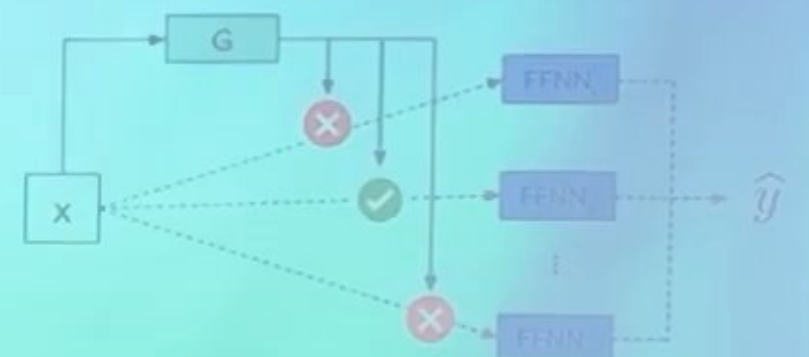


Stanford

MoE in Transformer-based models

Routing done for each token!

Feed Forward



Stanford University Figure adapted from "Attention is All You Need", Vaswani et al., 2017

ICML

Stanford

Training challenges include routing collapse



Symptom. Same expert gets selected most of the time.

"routing collapse"

Stanford

Training challenges include routing collapse

Symptom. Same expert gets selected most of the time.

"routing collapse"

Stanford University "Switch Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity", Fedus et al., 2021

TCML

Stanford

Training challenges include routing collapse



Symptom. Same expert gets selected most of the time.

"routing collapse"

Remedy. Force other experts to be "part of the game" too via auxiliary loss:

$$\text{loss} = \alpha \cdot N \cdot \sum_{i=1}^N f_i \cdot P_i$$

f_i = fraction of tokens routed to expert i

P_i = average routing probability for expert i

Stanford

Training challenges include routing collapse

Summary: "routing collapse"

"routing collapse"

Remedy. Force other experts to be "part of the game" too via auxiliary loss:

$$\text{loss} = \alpha \cdot N \cdot \sum_{i=1}^N f_i \cdot P_i$$

f_i = fraction of tokens routed to expert i

P_i = average routing probability for expert i

Stanford University "Switch Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity", Fedus et al., 2021

ILME

Stanford

Training challenges include routing collapse

Stanford

Routing collapse

Remedy. Force other experts to be "part of the game" too via auxiliary loss:

$$\text{loss} = \alpha \cdot N \cdot \sum_{i=1}^N f_i \cdot P_i$$

f_i = fraction of tokens routed to expert i

P_i = average routing probability for expert i

Training challenges include routing collapse

Solution

Force experts to be "part of the game"

Remedy. Force other experts to be "part of the game" too via auxiliary loss:

$$\text{loss} = \alpha \cdot N \cdot \sum_{i=1}^N f_i \cdot P_i$$

f_i = fraction of tokens routed to expert i

P_i = average routing probability for expert i

Training challenges include routing collapse

Symptom. Same expert gets selected most of the time.

"routing collapse"

Stanford University "Switch Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity" Fedus et al., 2021

ICML

Stanford

Training challenges include routing collapse

Symptom. Same expert gets selected most of the time.

"routing collapse"

Stanford University "Sparse Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity" Fedus et al., 2021

ICML

Stanford

Training challenges include routing collapse

Symptom. Same expert gets selected most of the time.

"routing collapse"

Stanford University "Sparse Reinforcement Learning: Scaling to Billion Parameter Networks with Deep and Efficient Sampling" Firooz et al., 2023

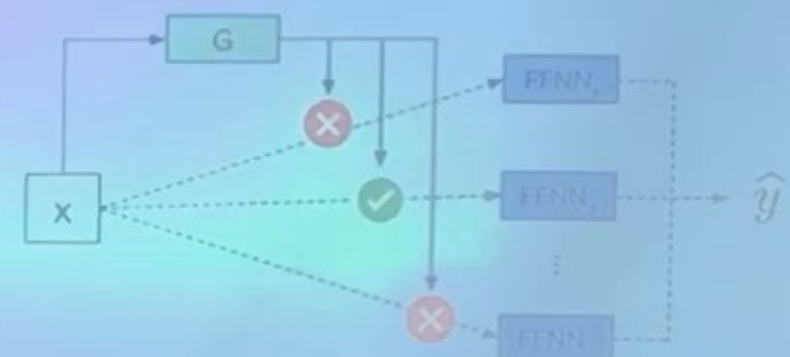
ICML

Stanford

MoE in Transformer-based models

Routing done for each token!

Feed Forward



Stanford University Figure adapted from "Attention Is All You Need" Vaswani et al., 2017

ICME

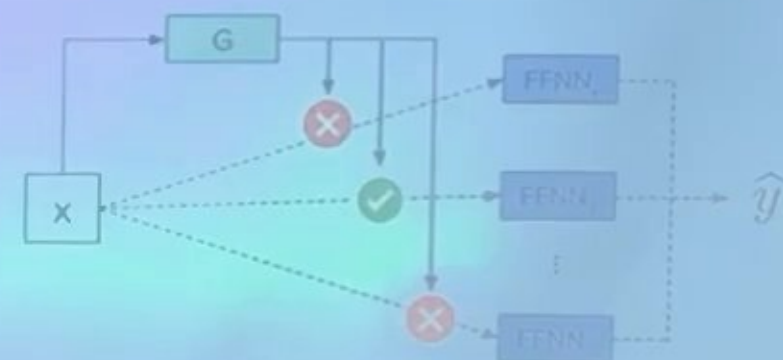


Stanford

MoE in Transformer-based models

Routing done for each token!

Feed Forward



Stanford University Figure extracted from "Attention is All You Need", Vaswani et al., 2017

ICML

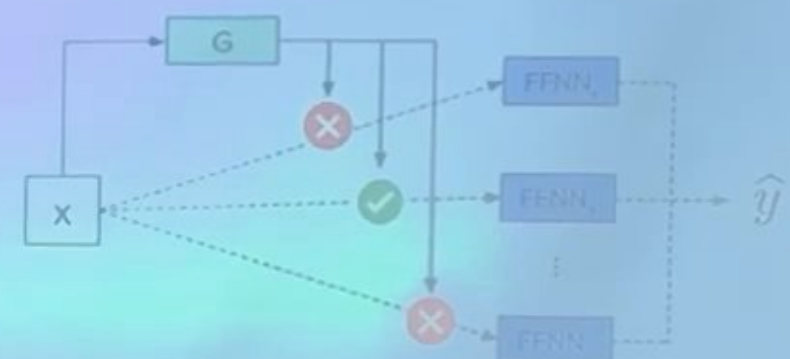


Stanford

MoE in Transformer-based models

Routing done for each token!

Feed Forward



Stanford University Figure adapted from "Attention is All you Need", Vaswani et al., 2017

ICME

Stanford

Training challenges include routing collapse

Symptom. Same expert gets selected most of the time.

"routing collapse"

Stanford University "Sparse Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity" Peters et al., 2021

ICML

Stanford

Training challenges include routing collapse

Question

Routing collapse

Remedy. Force other experts to be "part of the game" too via auxiliary loss:

$$\text{loss} = \alpha \cdot N \cdot \sum_{i=1}^N f_i \cdot P_i$$

f_i = fraction of tokens routed to expert i

P_i = average routing probability for expert i

Stanford University

"Switch Routing: Scaling to Billion Parameter Models with Simple and Efficient Sparsity", Fedus et al., 2021

ICML

Stanford

Training challenges include routing collapse

Expert routing collapse

Routing collapse

Remedy. Force other experts to be "part of the game" too via auxiliary loss:

$$\text{loss} = \alpha \cdot N \cdot \sum_{i=1}^N f_i \cdot P_i$$

f_i = fraction of tokens routed to expert i

P_i = average routing probability for expert i

Training challenges include routing collapse

Symptoms: some experts are unused or underused

Routing collapse

Remedy: Force other experts to be "part of the game" too via auxiliary loss:

$$\text{loss} = \alpha \cdot N \cdot \sum_{i=1}^N f_i \cdot P_i$$

f_i = fraction of tokens routed to expert i

P_i = average routing probability for expert i

Training challenges include routing collapse

Expert i only sees tokens routed to it

"routing collapse"

Remedy. Force other experts to be "part of the game" too via auxiliary loss:

$$\text{loss} = \alpha \cdot N \cdot \sum_{i=1}^N f_i \cdot P_i$$

f_i = fraction of tokens routed to expert i

P_i = average routing probability for expert i

Training challenges include routing collapse

Routing collapse

routing collapse

Remedy. Force other experts to be "part of the game" too via auxiliary loss:

$$\text{loss} = \alpha \cdot N \cdot \sum_{i=1}^N f_i \cdot P_i$$

f_i = fraction of tokens routed to expert i

P_i = average routing probability for expert i

Training challenges include routing collapse

Problem:

routing collapse

Remedy. Force other experts to be "part of the game" too via auxiliary loss:

$$\text{loss} = \alpha \cdot N \cdot \sum_{i=1}^N f_i \cdot P_i$$

f_i = fraction of tokens routed to expert i

P_i = average routing probability for expert i

Stanford University "Switch Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity" Fedus et al., 2022

ICML

Stanford

Training challenges include routing collapse

Symptoms:

"routing collapse"

Remedy: Force other experts to be "part of the game" too via auxiliary loss:

$$\text{loss} = \alpha \cdot N \cdot \sum_{i=1}^N f_i \cdot P_i$$

f_i = fraction of tokens routed to expert i

P_i = average routing probability for expert i

Training challenges include routing collapse

Symptom

routing collapse

Remedy. Force other experts to be "part of the game" too via auxiliary loss:

$$\text{loss} = \alpha \cdot N \cdot \sum_{i=1}^N f_i \cdot P_i$$

f_i = fraction of tokens routed to expert i

P_i = average routing probability for expert i

Training challenges include routing collapse

Symptom

"routing collapse"

Remedy. Force other experts to be "part of the game" too via auxiliary loss:

$$\text{loss} = \alpha \cdot N \cdot \sum_{i=1}^N f_i \cdot P_i$$

f_i = fraction of tokens routed to expert i

P_i = average routing probability for expert i

Training challenges include routing collapse

Symptom:

"routing collapse"

Remedy. Force other experts to be "part of the game" too via auxiliary loss:

$$\text{loss} = \alpha \cdot N \cdot \sum_{i=1}^N f_i \cdot P_i$$

f_i = fraction of tokens routed to expert i

P_i = average routing probability for expert i

Stanford University

"GPT-4 Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity", Fedus et al., 2021

ICML

Stanford

Training challenges include routing collapse

Symptoms

Routing collapse

Remedy. Force other experts to be "part of the game" too via auxiliary loss:

$$\text{loss} = \alpha \cdot N \cdot \sum_{i=1}^N f_i \cdot P_i$$

f_i = fraction of tokens routed to expert i

P_i = average routing probability for expert i

Training challenges include routing collapse

Symptoms

routing collapse

Remedy. Force other experts to be "part of the game" too via auxiliary loss:

$$\text{loss} = \alpha \cdot N \cdot \sum_{i=1}^N f_i \cdot P_i$$

f_i = fraction of tokens routed to expert i

P_i = average routing probability for expert i

Training challenges include routing collapse

Symptoms

routing collapse

Remedy. Force other experts to be "part of the game" too via auxiliary loss:

$$\text{loss} = \alpha \cdot N \cdot \sum_{i=1}^N f_i \cdot P_i$$

f_i = fraction of tokens routed to expert i

P_i = average routing probability for expert i

Training challenges include routing collapse

Symptoms

Routing collapse

Remedy. Force other experts to be "part of the game" too via auxiliary loss:

$$\text{loss} = \alpha \cdot N \cdot \sum_{i=1}^N f_i \cdot P_i$$

f_i = fraction of tokens routed to expert i

P_i = average routing probability for expert i

MoE in Transformer-based models



Stanford University

Figure adapted from "Mixture of Experts in All You Need" Vaswani et al., 2017

ICML

Stanford

MoE in Transformer-based models



Stanford University Figure adapted from "Attention is All You Need" Vaswani et al., 2017

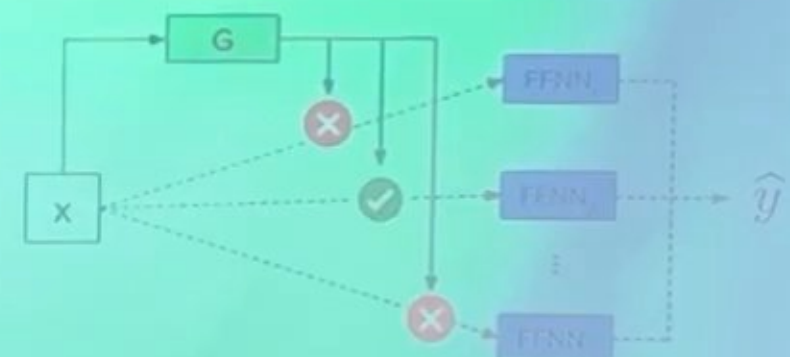
HCME

Stanford

MoE in Transformer-based models

Routing done for each token!

Feed Forward



Stanford University

Figure adapted from "Adaptive All-You-Need-Is-Add", Wotawattana et al. 2017

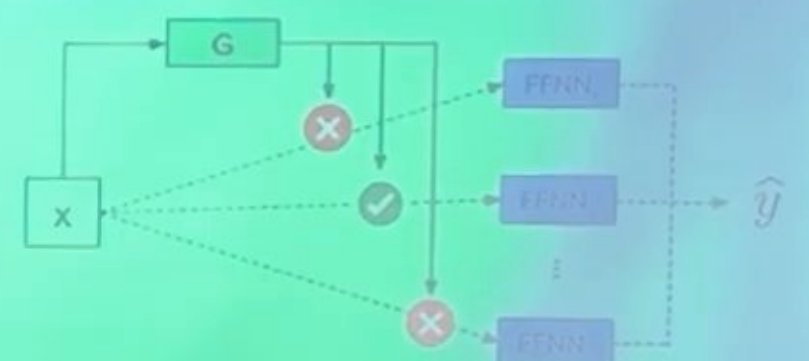
ICML

Stanford

MoE in Transformer-based models

Routing done for each token!

Feed Forward



Stanford University Figure adapted from "Attention is All You Need" Vaswani et al., 2017.

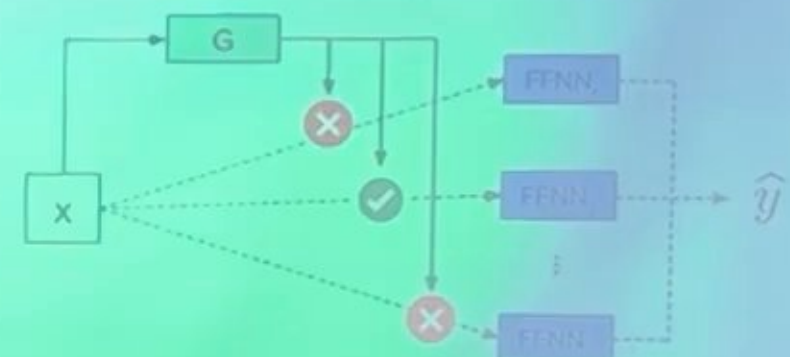
ICME

Stanford

MoE in Transformer-based models

Routing done for each token!

Feed Forward



Stanford University Figure adapted from "Attention is All You Need", Vaswani et al., 2017

ICME

Stanford

MoE in Transformer-based models



Stanford University Figure adapted from "Attention is All You Need", Vaswani et al., 2017.

ICML'17

Stanford

MoE in Transformer-based models

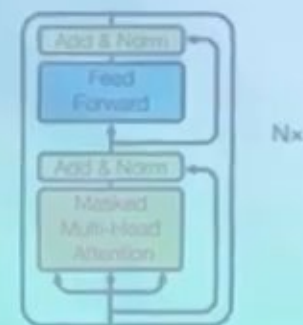


Stanford University Figure adapted from "Attention is All You Need", Vaswani et al., 2017

ICME

Stanford

MoE in Transformer-based models



Stanford University Figure adapted from "Attention is All You Need", Vaswani et al., 2017

ICME

Stanford

Interpreting experts



Layer 0

```
class MoeLayer(nn.Module):
    def __init__(self, experts: List[nn.Module],
                 super().__init__():
        assert len(experts) > 0
        self.experts = nn.ModuleList(experts)
        self.gate = gate
        self.args = moe_args

    def forward(self, inputs: torch.Tensor):
        inputs_squashed = inputs.view(-1, inputs.size(-1))
        gate_logits = self.gate(inputs_squashed)
        weights, selected_experts = torch.topk(
            gate_logits, self.args.num_experts_per_token
        )
        weights = nn.functional.softmax(
            weights,
            dim=-1,
            dtype=torch.float,
        ).type_as(inputs)
        results = torch.zeros_like(inputs_squashed)
        for i, expert in enumerate(self.experts):
            batch_idx, nth_expert = torch.where(selected_experts == i)
            results[batch_idx] += weights[batch_idx] * expert(inputs_squashed[batch_idx])
        return results.view_as(inputs)
```

Each color represent an expert

Stanford

Interpreting experts

Layer 0

```
class MoELayer(nn.Module):
    def __init__(self, experts: List[nn.Module],
                 gate: nn.Module,
                 args: dict):
        super().__init__()
        assert len(experts) > 0
        self.experts = nn.ModuleList(experts)
        self.gate = gate
        self.args = args

    def forward(self, inputs: torch.Tensor):
        inputs_squashed = inputs.view(-1, inputs.size(-1))
        gate_logits = self.gate(inputs_squashed)
        weights, selected_experts = torch.topk(
            gate_logits, self.args.num_experts)
        weights = nn.functional.softmax(
            weights,
            dim=-1,
            dtype=torch.float,
        ).type_as(inputs)
        results = torch.zeros_like(inputs_squashed)
        for i, expert in enumerate(self.experts):
            batch_idx, nth_expert = torch.where(
                results == weights[i])
            results[batch_idx] += weights[i] *
                inputs_squashed[batch_idx]
        return results.view_as(inputs)
```

Each color represent an expert

University "MoE of Experts" Jiang et al., 2024

ICME

Stanford



Transformers & Large Language Models

Stanford University

ICME

LLM overview

MoE-based LLMs

Response generation

Prompting strategies

Inference optimizations



Stanford



Transformers & Large Language Models

Stanford University

ICME

- LLM overview
- MoE-based LLMs
- Response generation**
- Prompting strategies
- Inference optimizations

Stanford

Next token prediction



LLM

[BOS]

Stanford

Predicting next token

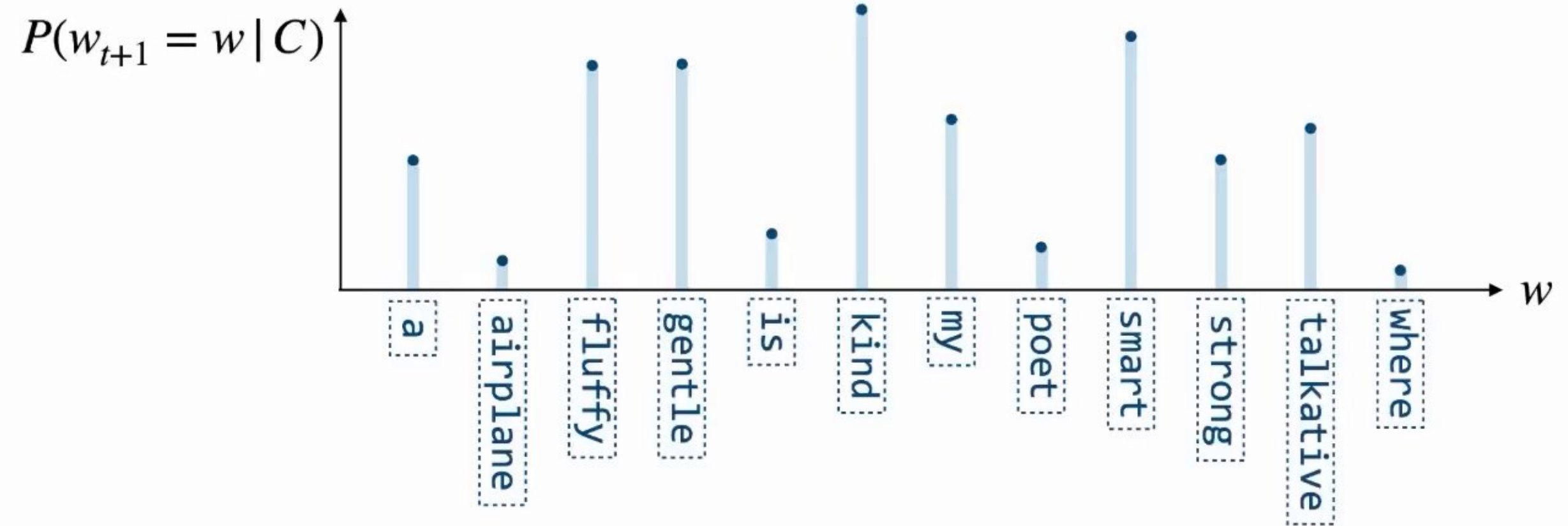
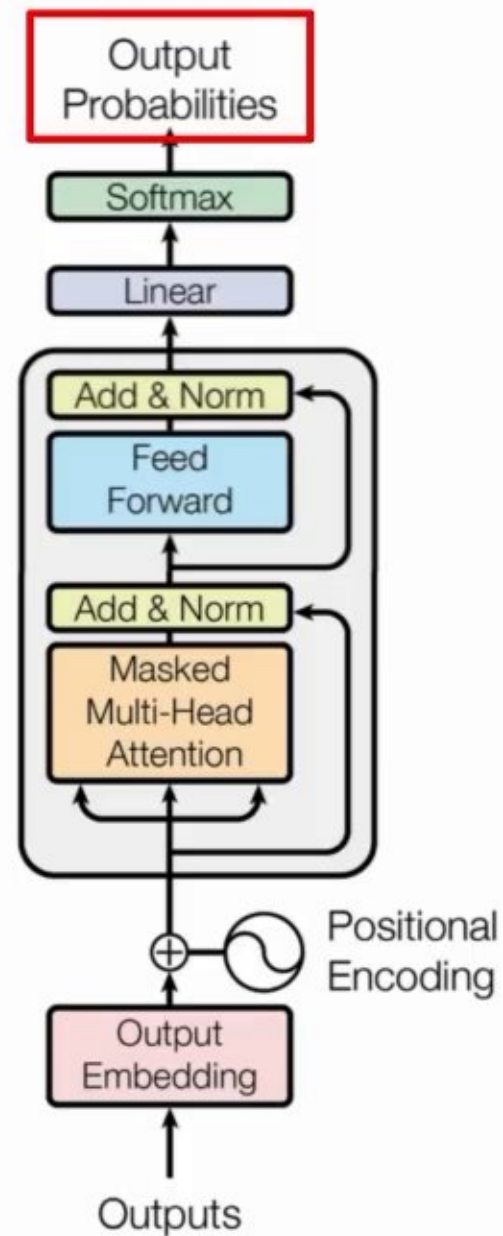


Stanford University Figure adapted from "Attention Is All You Need", Vaswani et al., 2017

ICML

Stanford

Predicting next token

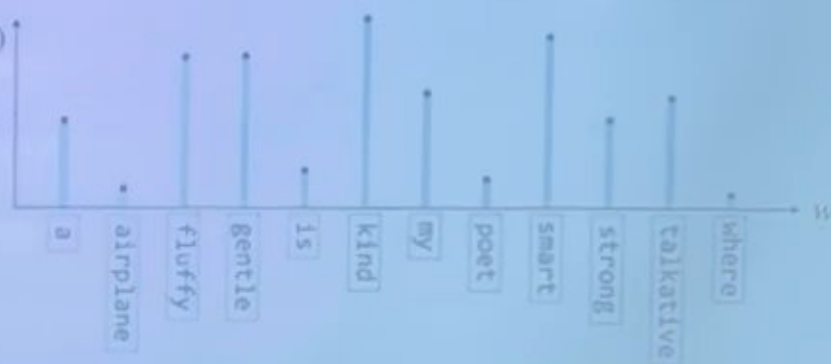


Stanford

Predicting next token



$$P(w_{t+1} = w | C)$$



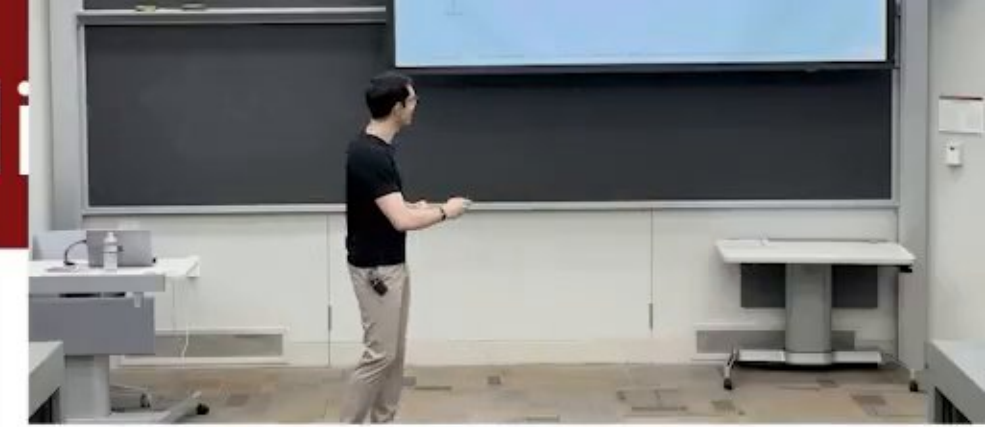
Stanford University

Figure adapted from "Attention is All You Need", Vaswani et al., 2017 "Super Study Guide: Transformers and Large Language Models" Joshi et al., 2024

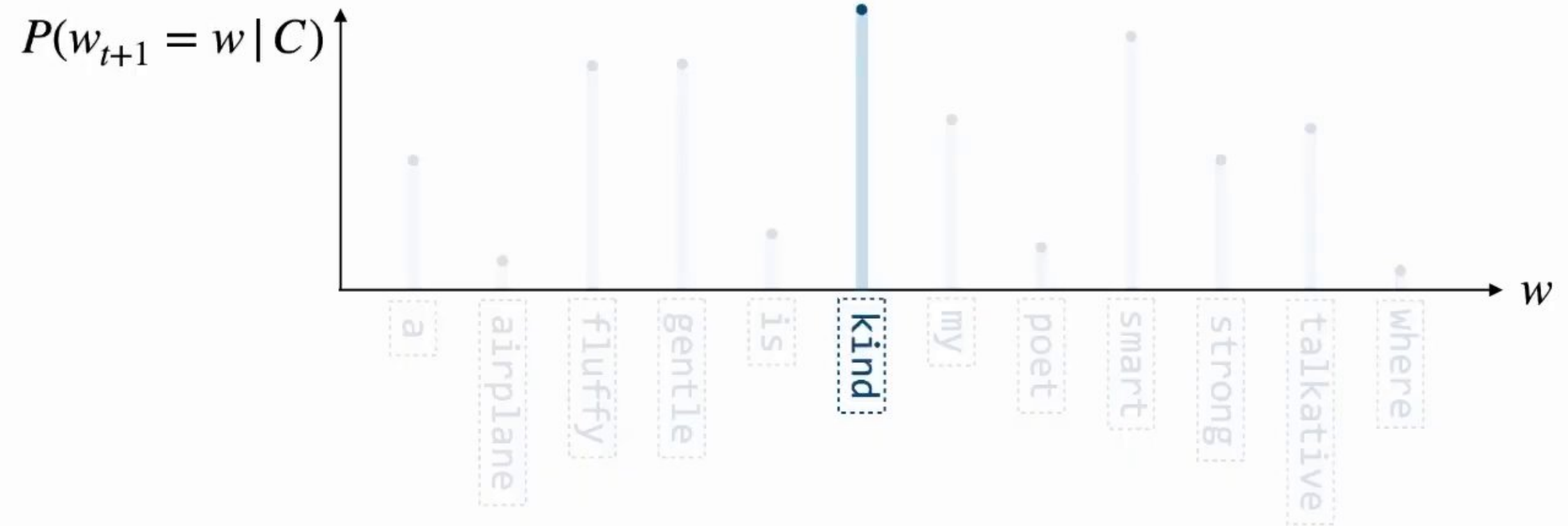
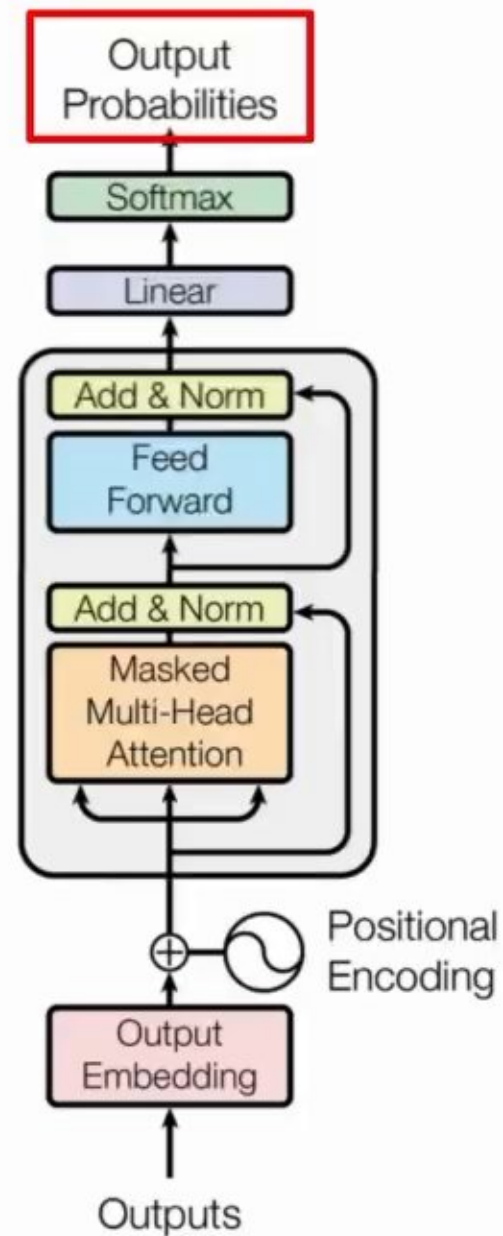
ICME

Stanford

Predicting next token with greedy decoding



1st idea. Take token with highest predicted probability



Stanford

Predicting next token with greedy decoding

1st idea. Take token with highest predicted probability



Limitations. Output not optimal, natural and/or diverse

Stanford University

Figure adapted from "Attention is All You Need", Vaswani et al., 2017 "Super-Sized GPT: Transformers and Large Language Models", Amini et al., 2024

ICLME



Predicting next token with greedy decoding

1st idea. Take token with highest predicted probability



Limitations. Output not optimal, natural and/or diverse

Stanford University

Figure adapted from "Attention is All You Need" Vaswani et al., 2017 "Super-Sized Guide: Transformers and Large Language Models" Anand et al., 2024

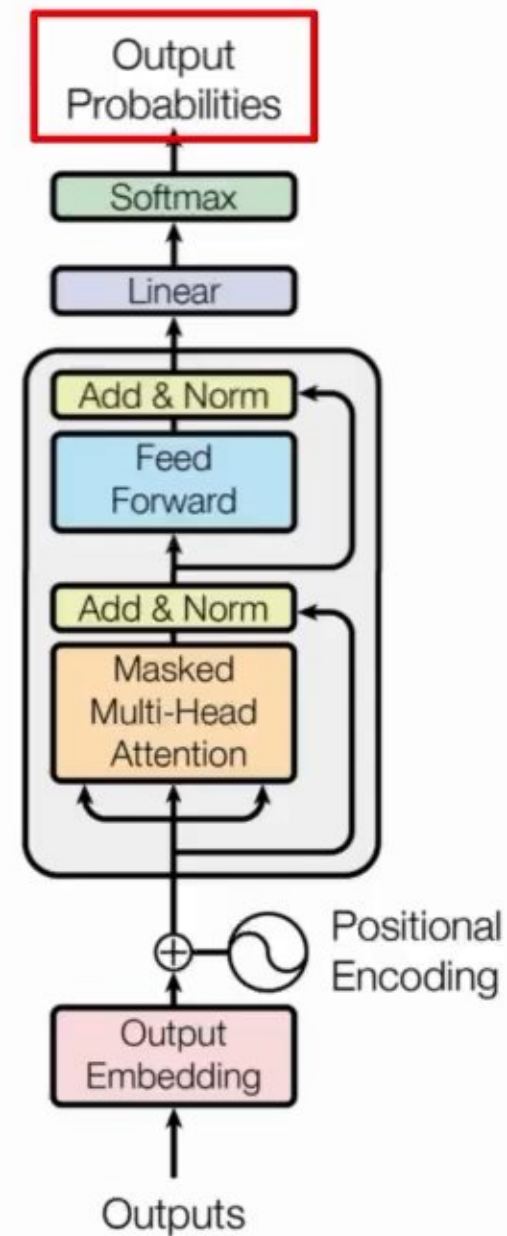
ICME

Stanford

Predicting next token with beam search



2nd idea. Keep k paths that are the most likely

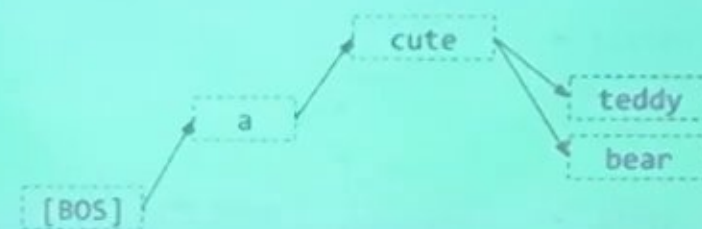
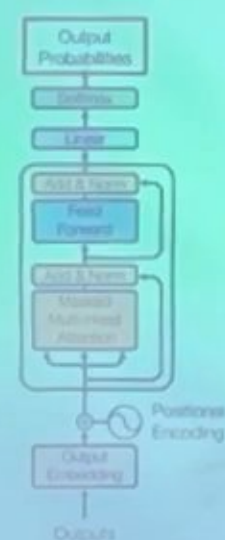


[BOS]

Stanford

Predicting next token with beam search

2nd idea. Keep k paths that are the most likely



Stanford University

Figure adapted from "Attention is All You Need", Vaswani et al., 2017

ICML

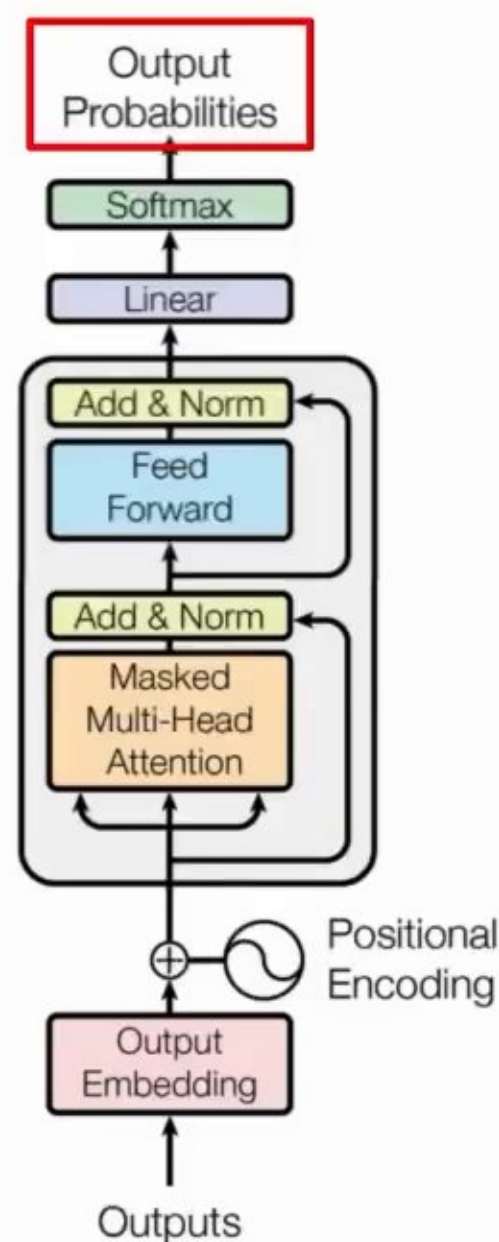
Stanford

Predicting next token with sampling

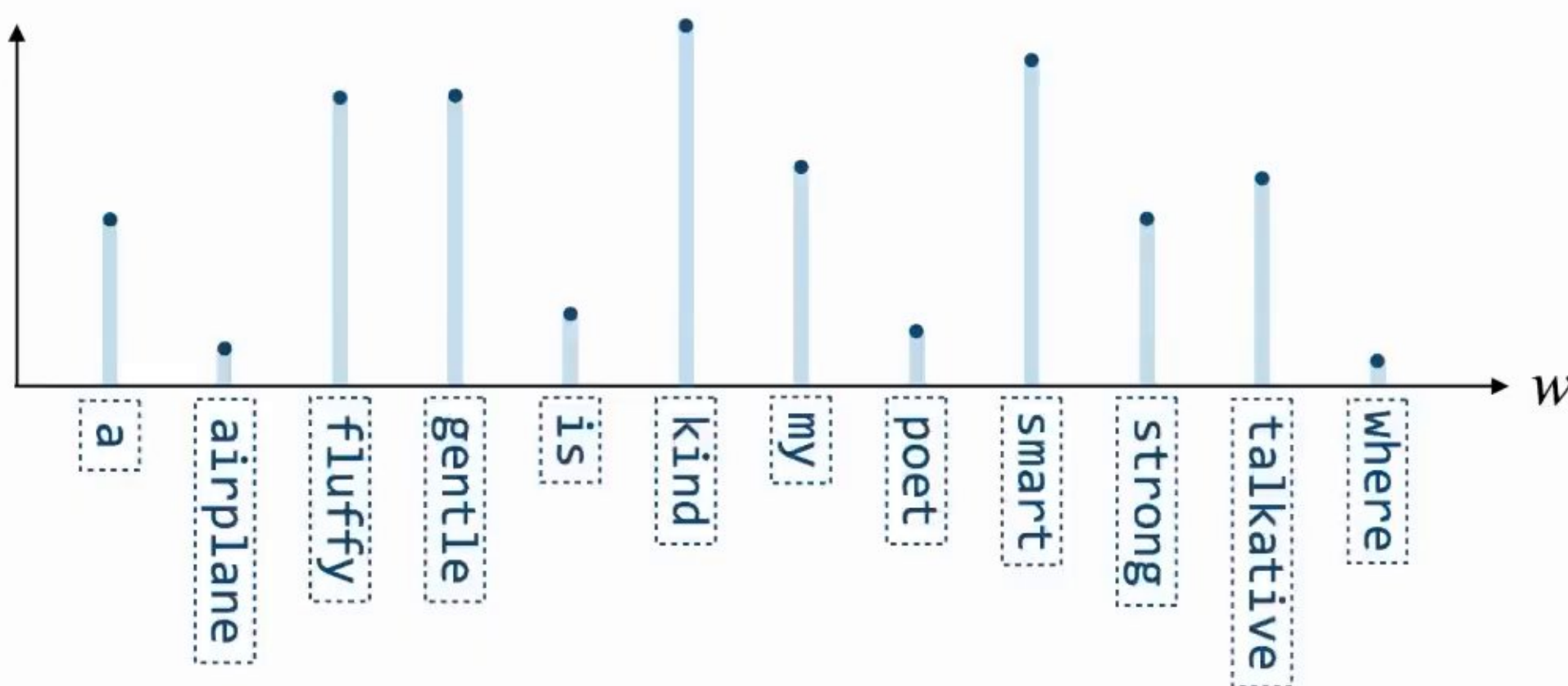


3rd idea. Sample next token from probability distribution:

$$\hat{w}_{t+1} \sim P(w_{t+1} | C)$$



$$P(w_{t+1} = w | C)$$



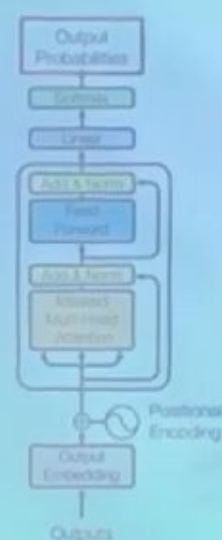
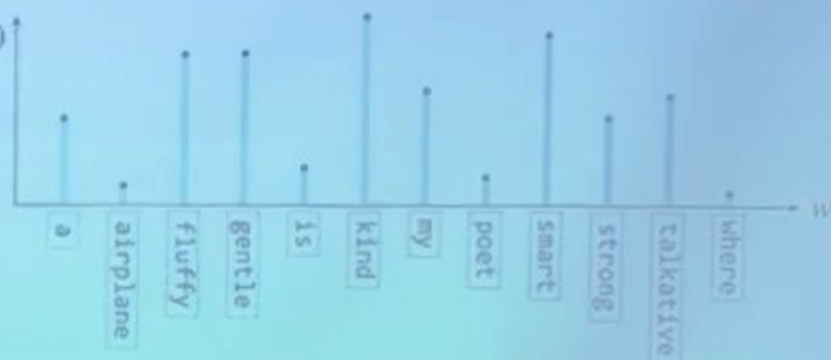
Stanford

Predicting next token with sampling

3rd idea. Sample next token from probability distribution:

$$\hat{w}_{t+1} \sim P(w_{t+1}|C)$$

$$P(w_{t+1} = w | C)$$

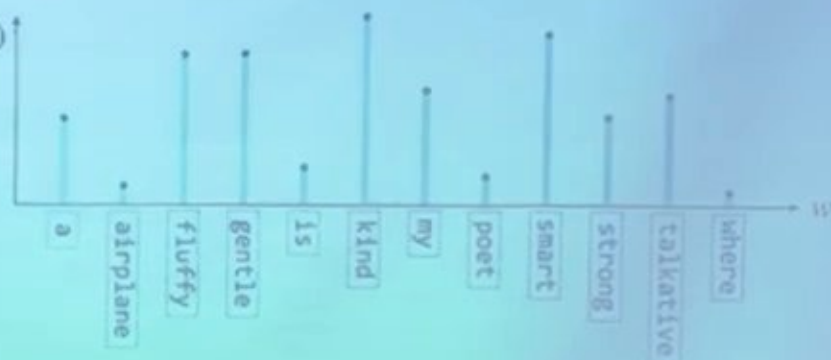


Predicting next token with sampling

3rd idea. Sample next token from probability distribution:

$$\hat{w}_{t+1} \sim P(w_{t+1}|C)$$

$$P(w_{t+1} = w | C)$$



Stanford University

"Super Study Guide: Transformers and Large Language Models" Aravi et al., 2024

ICLME



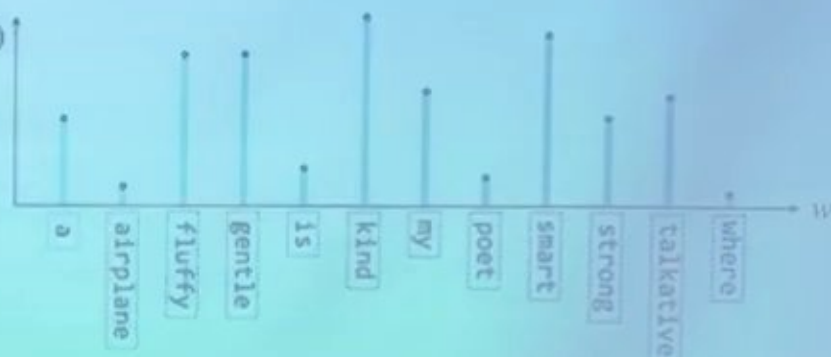
Stanford

Predicting next token with sampling

3rd idea. Sample next token from probability distribution:

$$\hat{w}_{t+1} \sim P(w_{t+1}|C)$$

$$P(w_{t+1} = w | C)$$



Stanford University

"Super Study Guide: Transformers and Large Language Models", Amidi et al., 2024.

ICME



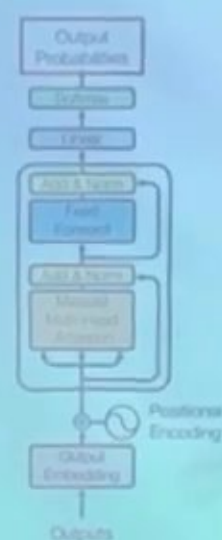
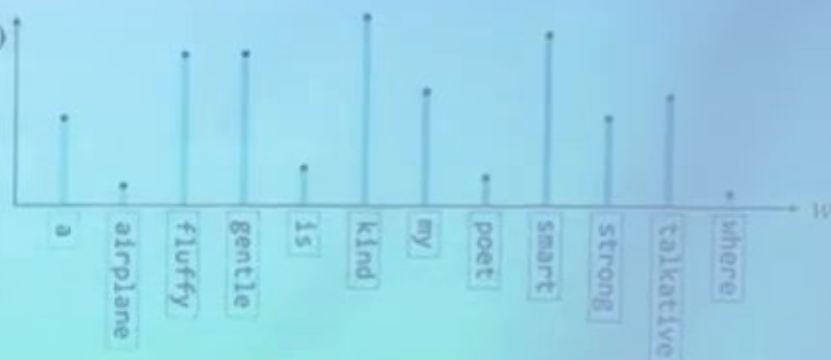
Stanford

Predicting next token with sampling

3rd idea. Sample next token from probability distribution:

$$\hat{w}_{t+1} \sim P(w_{t+1}|C)$$

$$P(w_{t+1} = w | C)$$

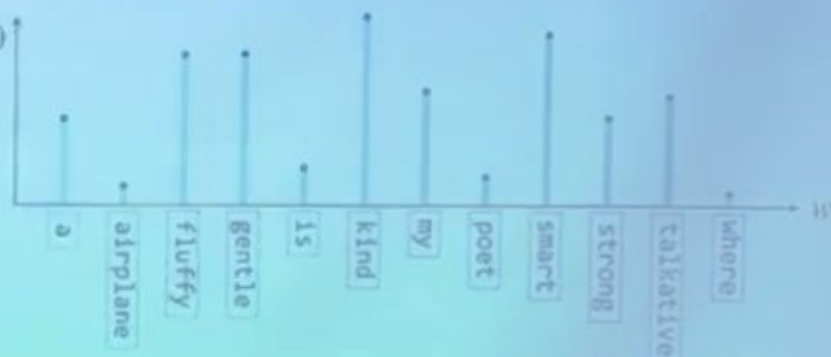


Predicting next token with sampling

3rd idea. Sample next token from probability distribution:

$$\hat{w}_{t+1} \sim P(w_{t+1}|C)$$

$$P(w_{t+1} = w | C)$$

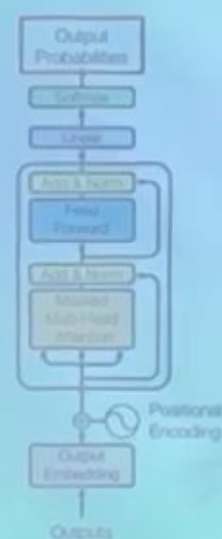
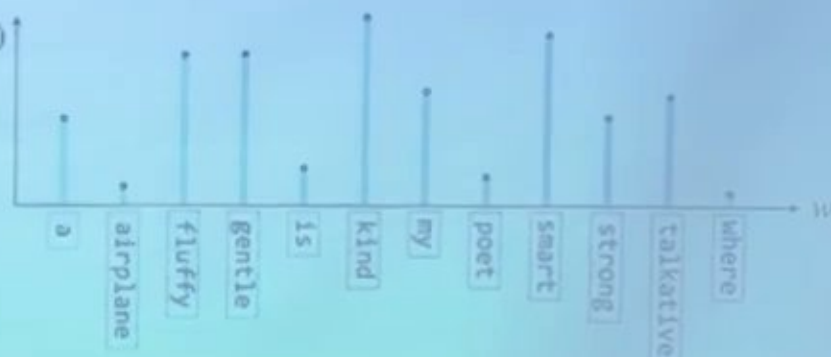


Predicting next token with sampling

3rd idea. Sample next token from probability distribution:

$$\hat{w}_{t+1} \sim P(w_{t+1}|C)$$

$$P(w_{t+1} = w | C)$$



Stanford University

"Super Study Guide: Transformers and Large Language Models" Arad et al., 2024

ICML

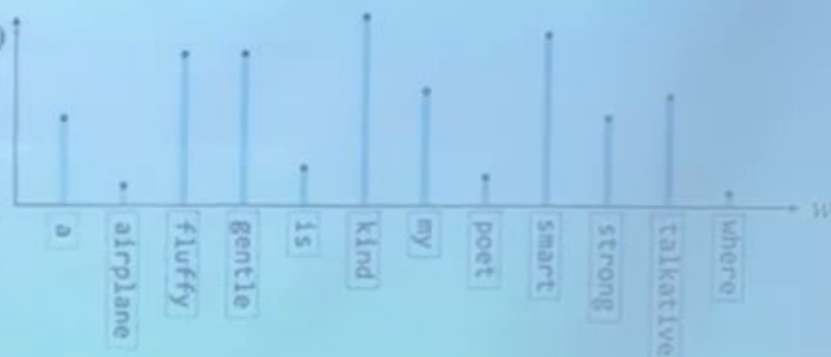
Stanford

Predicting next token with sampling

3rd idea. Sample next token from probability distribution:

$$\hat{w}_{t+1} \sim P(w_{t+1}|C)$$

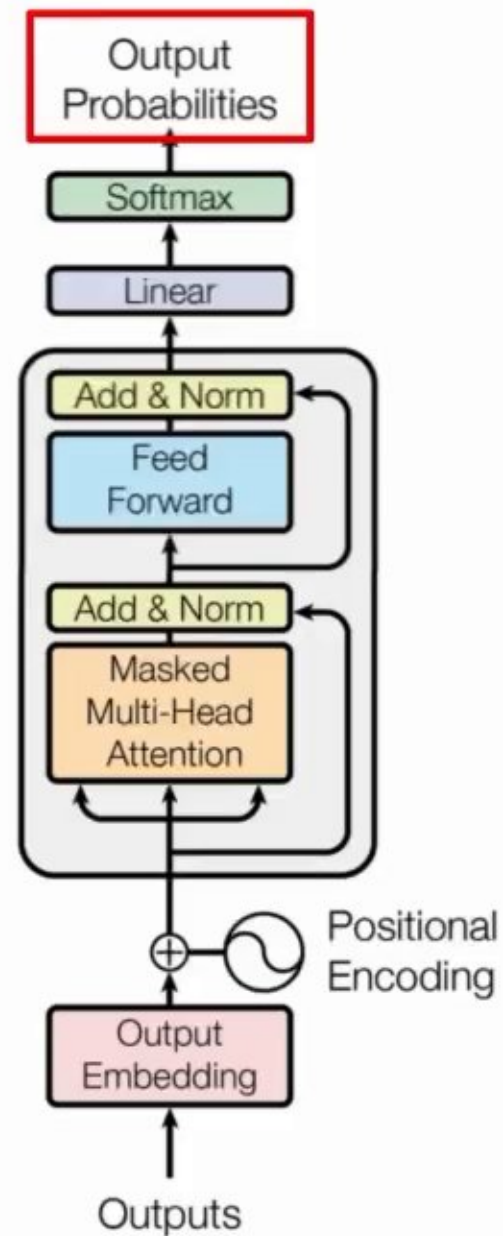
$$P(w_{t+1} = w | C)$$



Sampling strategies

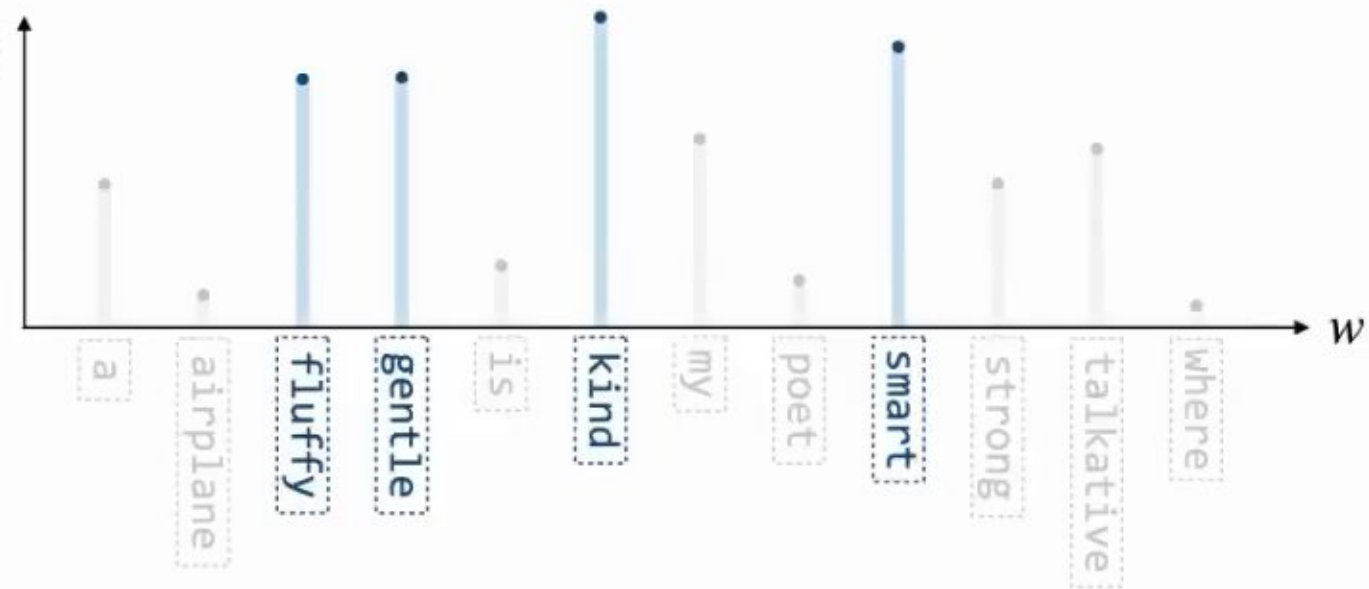


- **Top-k:** Sample among top k most probable tokens



$k = 4$

$$P(w_{t+1} = w | C)$$



Stanford

Next token prediction



But **how** are probabilities **obtained**?

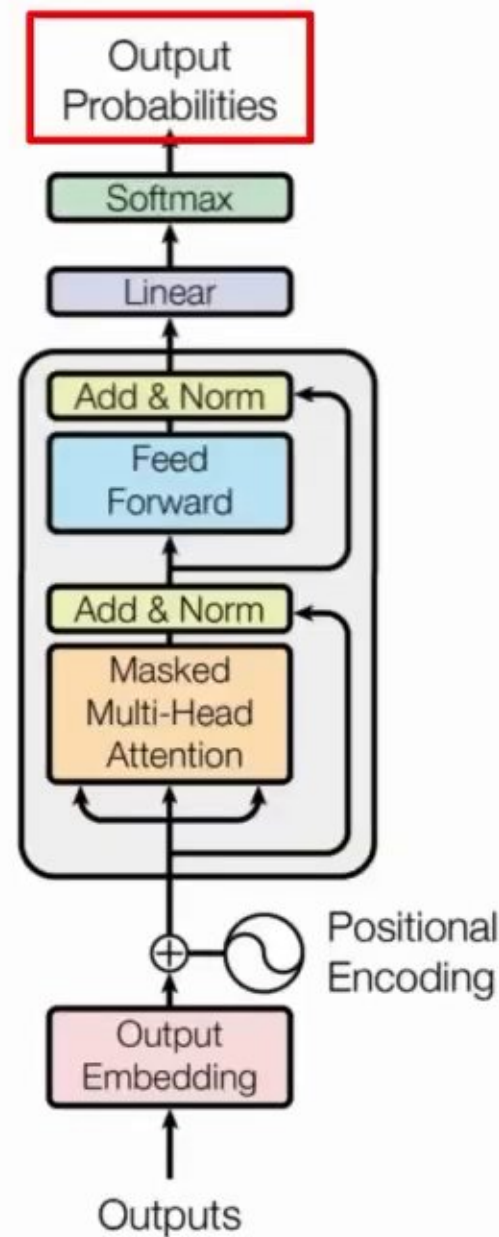
Stanford

Predicting next token with sampling

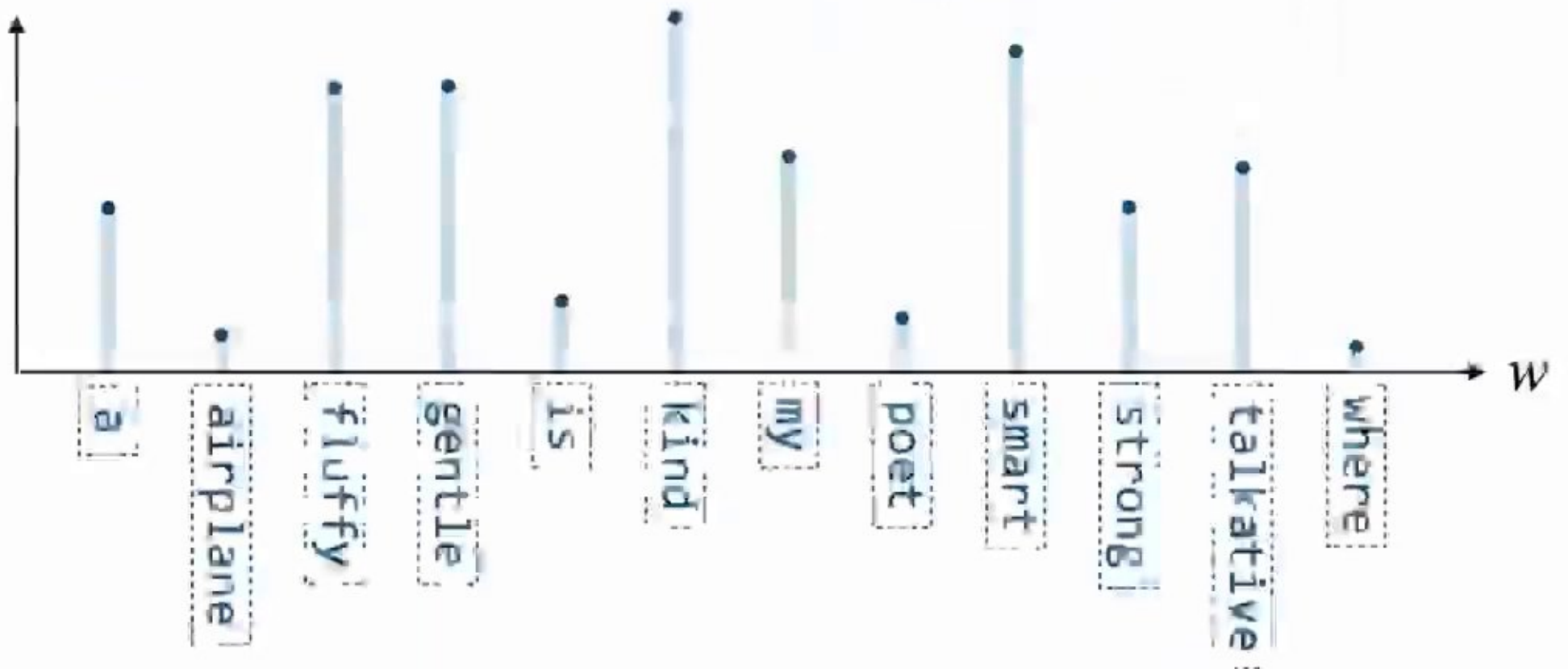


3rd idea. Sample next token from probability distribution:

$$\hat{w}_{t+1} \sim P(w_{t+1} | C)$$

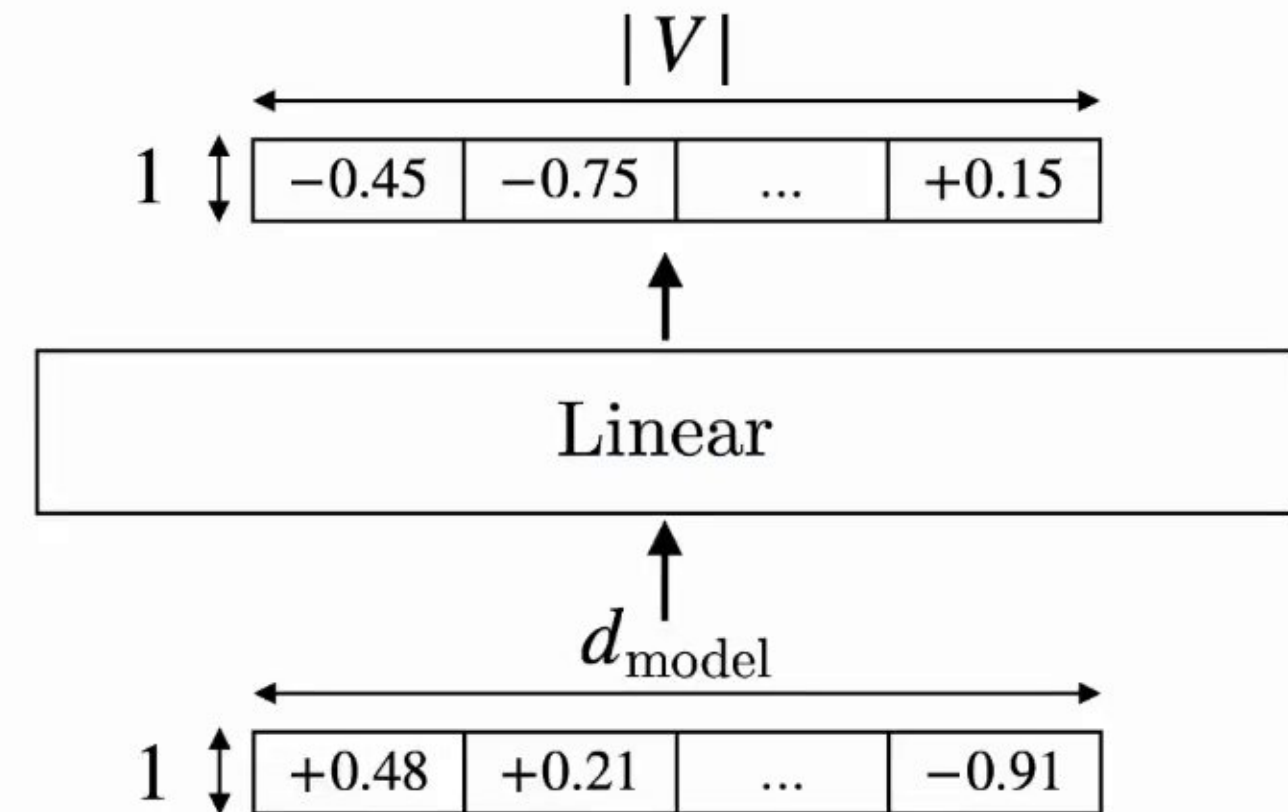
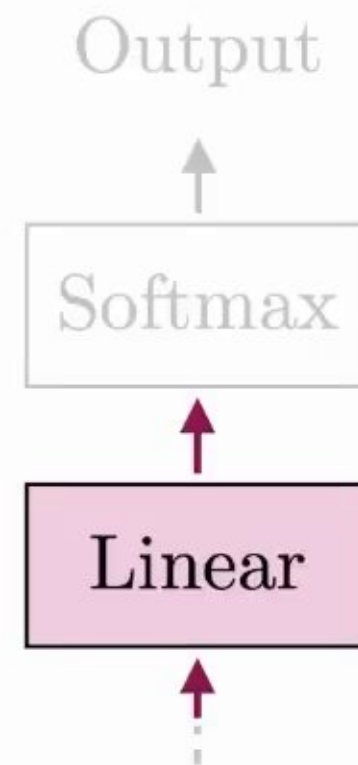
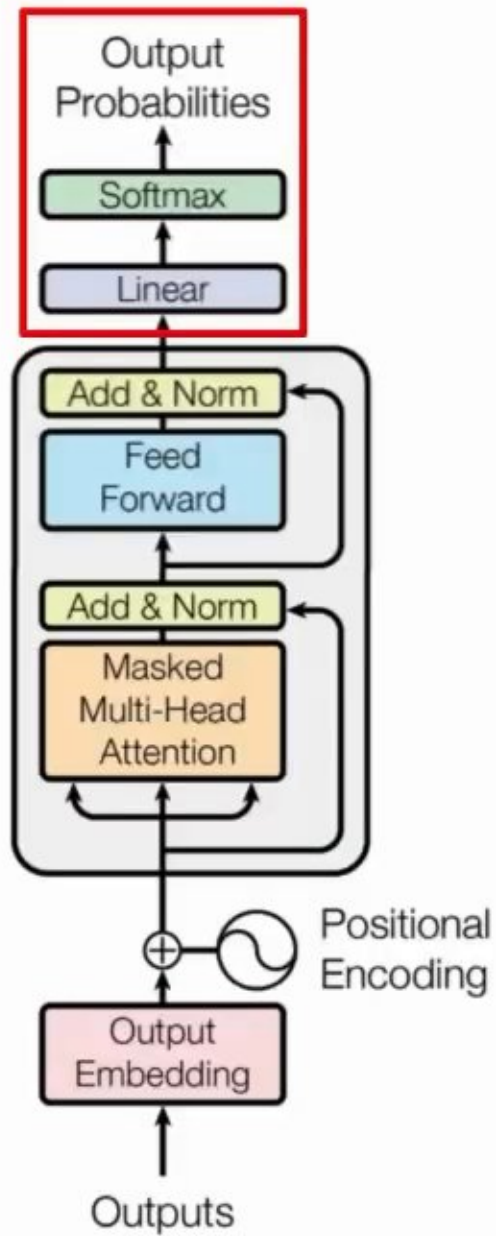


$$P(w_{t+1} = w | C)$$



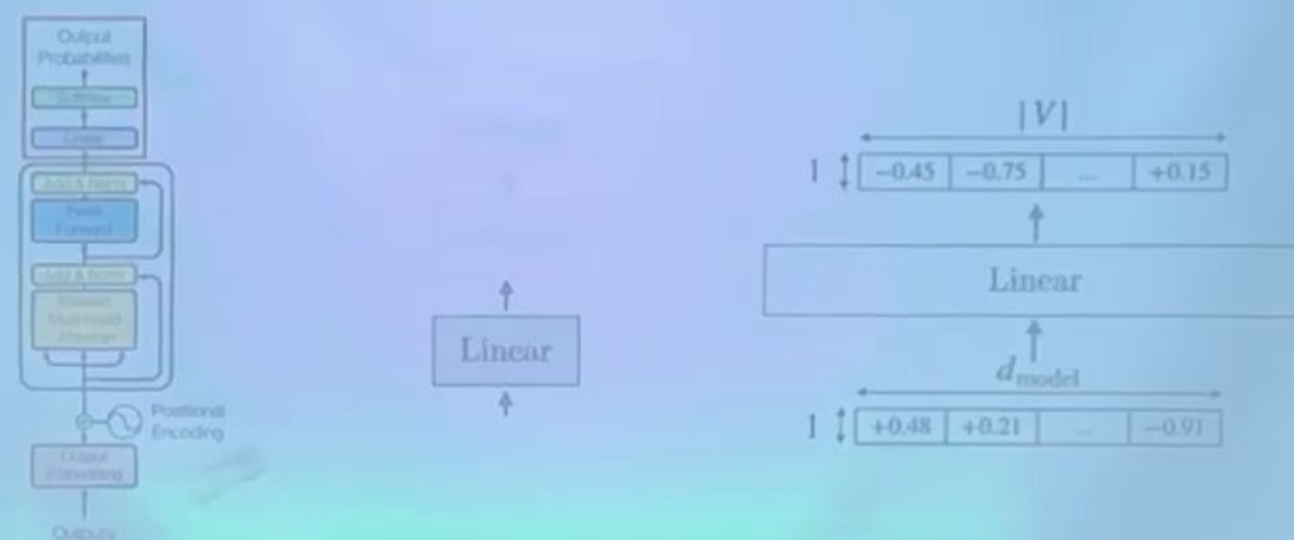
Stanford

Probability computation



Stanford

Probability computation



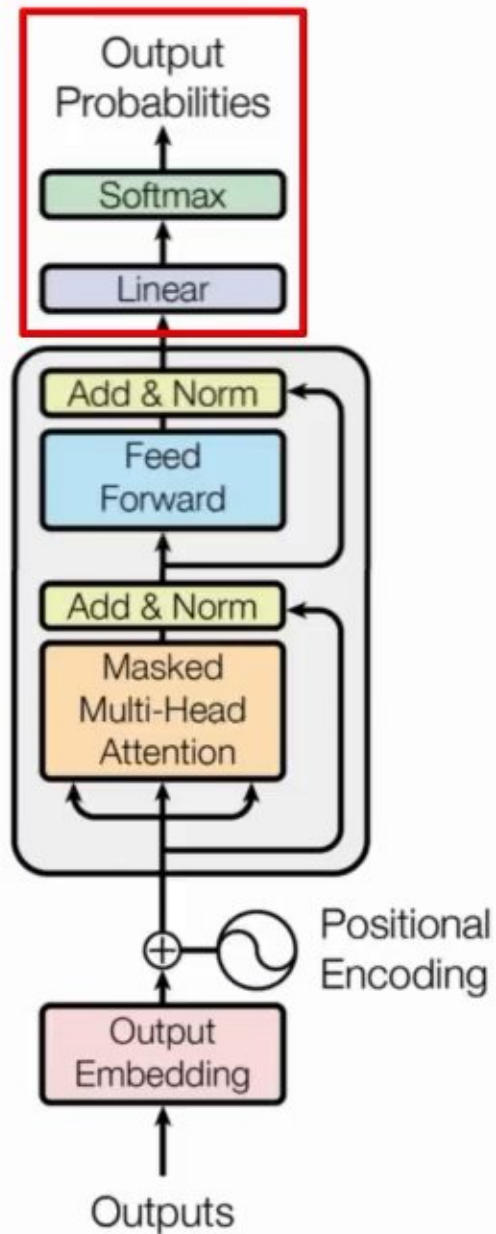
Stanford University

"Scaling GPT: Scaling Transformers and Large Language Models" Arad et al., 2024

ICML

Stanford

Temperature allows to tweak output prob



$$P_{\text{adj}}(w_{t+1} = w_i | C) = \frac{\exp\left(\frac{x_i}{T}\right)}{\sum_{j=1}^n \exp\left(\frac{x_j}{T}\right)}$$

Stanford

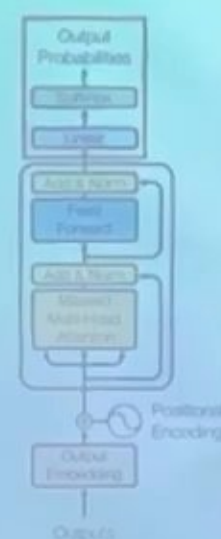
Temperature allows to tweak output probabilities



Output
↑
Softmax
↑

$$P_{\text{adj}}(w_{t+1} = w_i | C) = \frac{\exp\left(\frac{x_i}{T}\right)}{\sum_{j=1}^n \exp\left(\frac{x_j}{T}\right)}$$

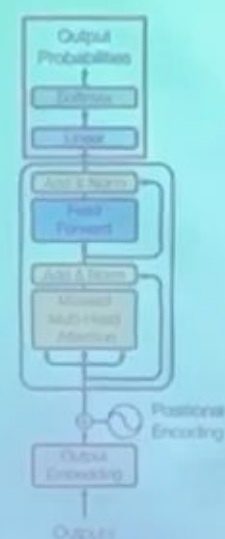
Temperature allows to tweak output probabilities



Softmax

$$P_{\text{adj}}(w_{t+1} = w_i | C) = \frac{\exp\left(\frac{x_i}{T}\right)}{\sum_{j=1}^n \exp\left(\frac{x_j}{T}\right)}$$

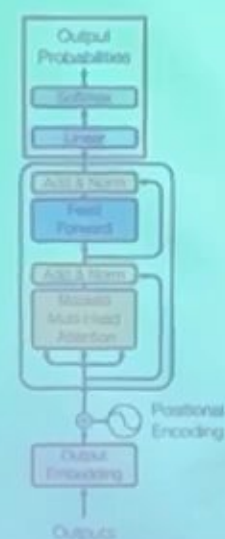
Temperature allows to tweak output probabilities



Softmax

$$P_{\text{adj}}(w_{t+1} = w_i | C) = \frac{\exp\left(\frac{x_i}{T}\right)}{\sum_{j=1}^n \exp\left(\frac{x_j}{T}\right)}$$

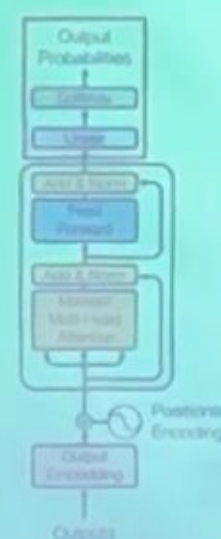
Temperature allows to tweak output probabilities



↑
Softmax
↑

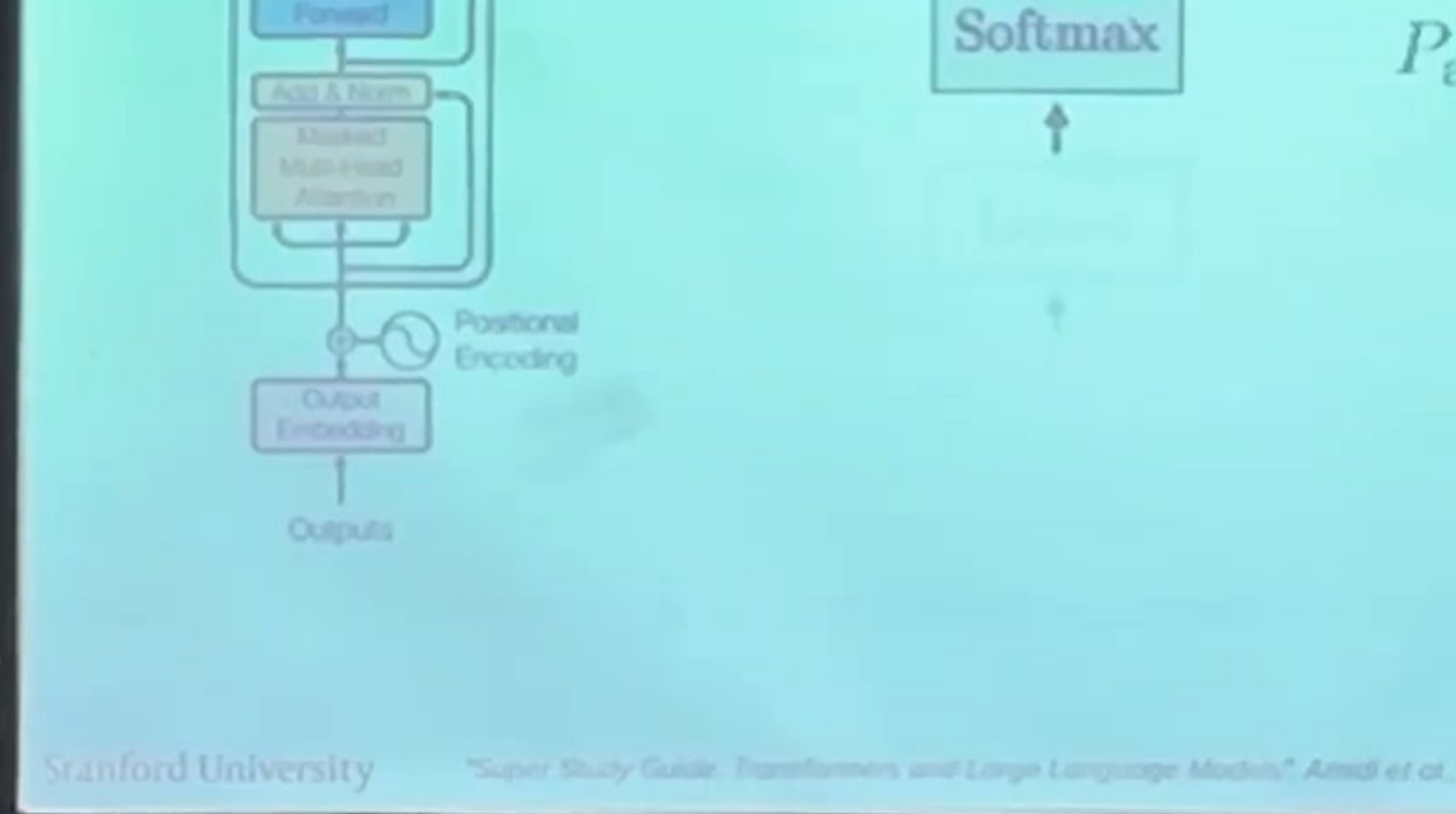
$$P_{\text{adj}}(w_{t+1} = w_i | C) = \frac{\exp\left(\frac{x_i}{T}\right)}{\sum_{j=1}^n \exp\left(\frac{x_j}{T}\right)}$$

Temperature allows to tweak output probabilities

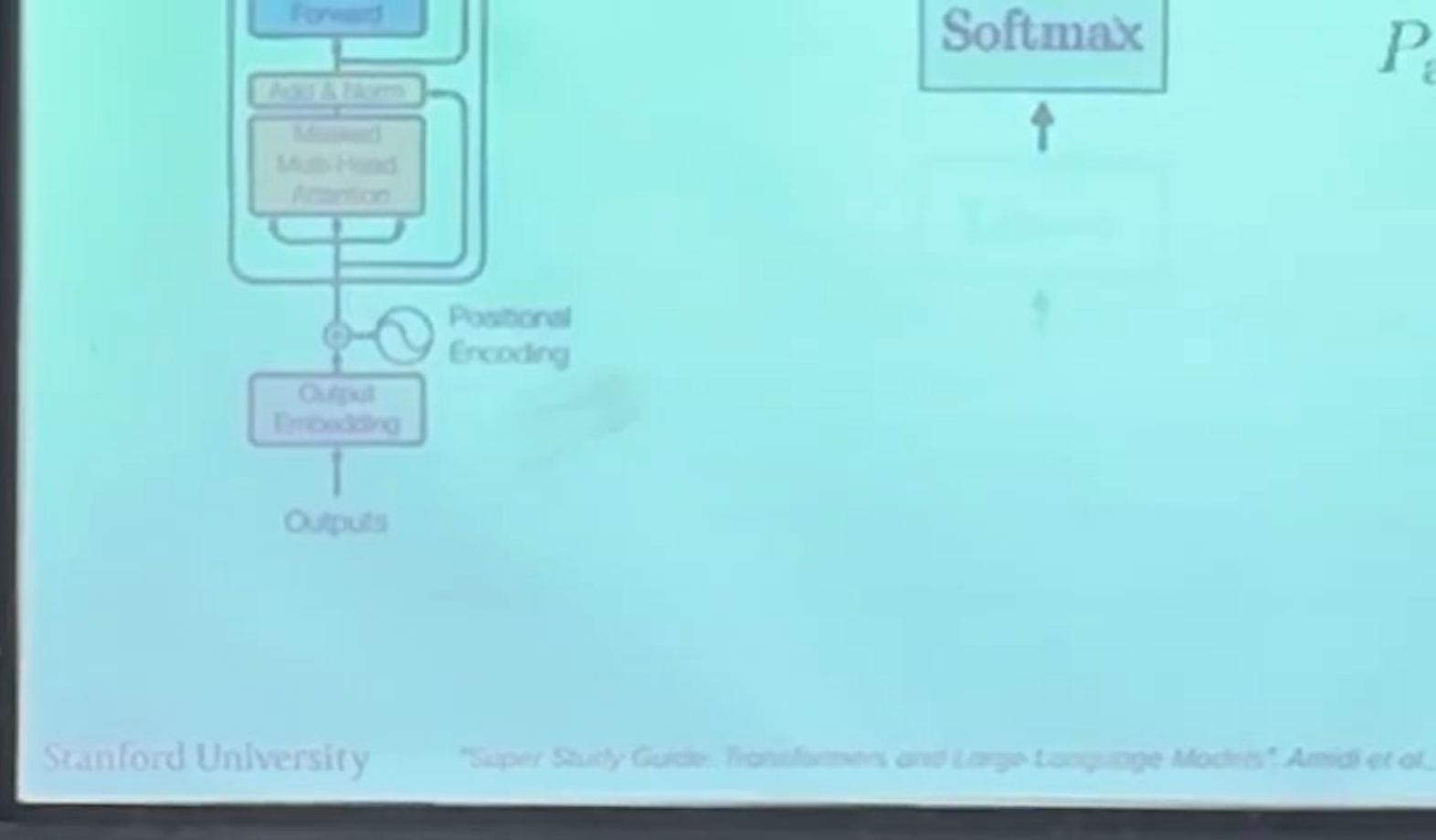
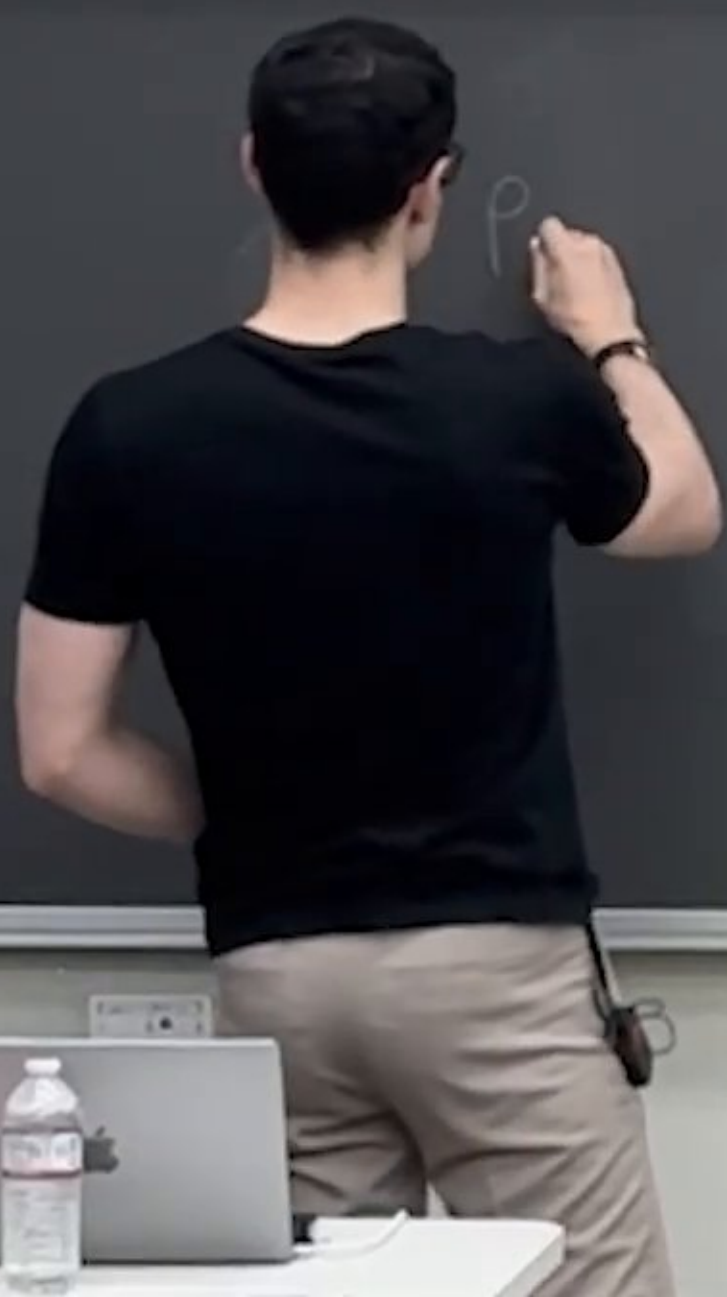


Softmax

$$P_{\text{adj}}(w_{t+1} = w_i | C) = \frac{\exp\left(\frac{x_i}{T}\right)}{\sum_{j=1}^n \exp\left(\frac{x_j}{T}\right)}$$

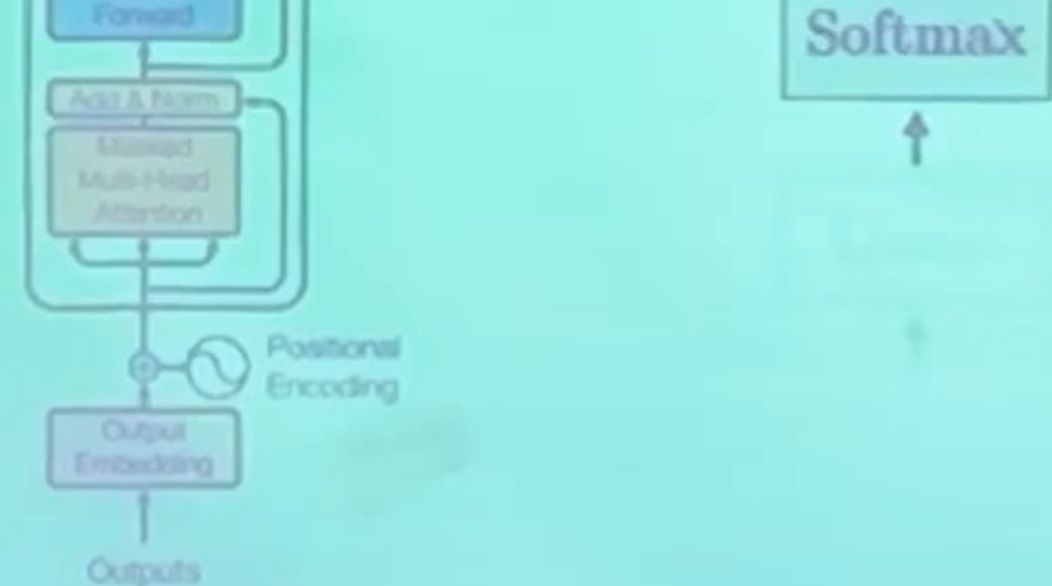


Stanford



Stanford

$$p_i = \frac{e^{\frac{\alpha_i}{T}}}{\sum_j e^{\frac{-1}{T}}}$$

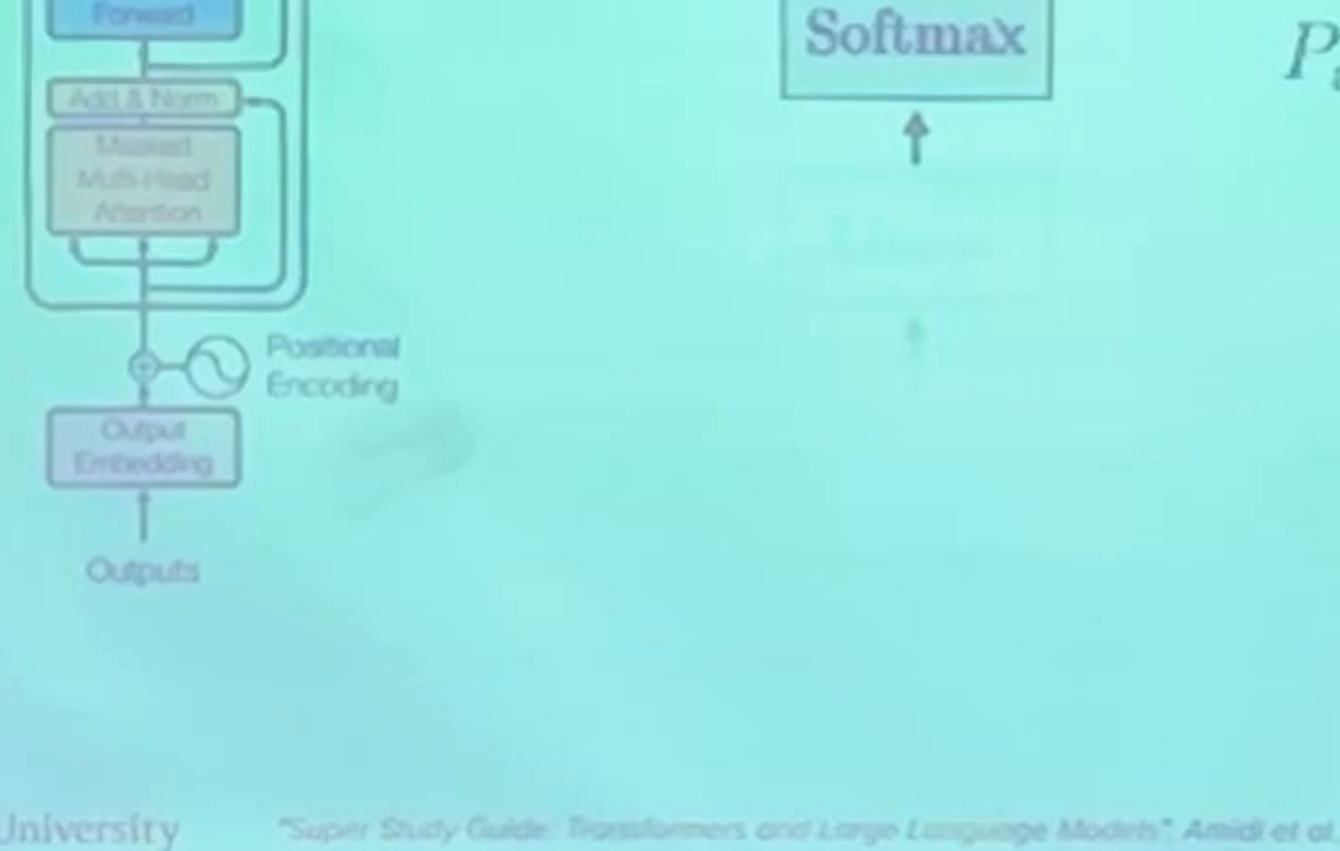


Stanford University

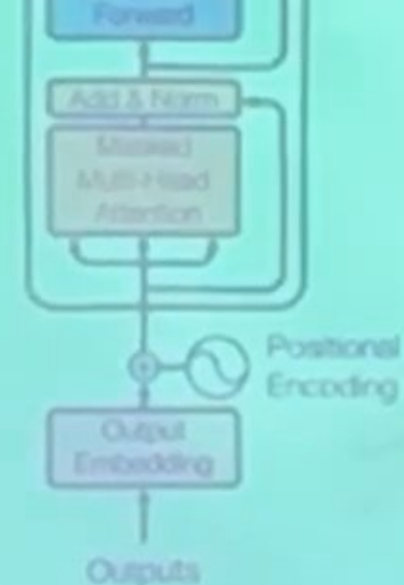
"Super Study Guide: Transformers and Large Language Models", Amidi et al.

Stanford

$$p_i = \frac{e^{\frac{\alpha_i}{\tau}}}{\sum_j e^{\frac{\alpha_j}{\tau}}}$$



$$p_i = \frac{e^{\frac{\alpha_i}{T}}}{\sum_j e^{\frac{\alpha_j}{T}}}$$



Softmax

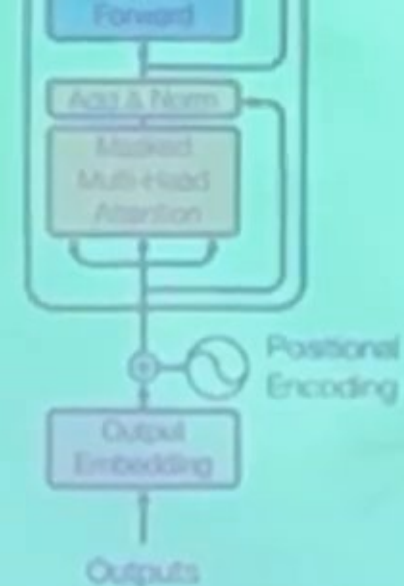
P_{θ}

Stanford University

"Super Study Guide: Transformers and Large Language Models", Amidi et al.

Stanford

$$p_i = \frac{e^{\frac{\alpha_i}{T}}}{\sum_j e^{\frac{\alpha_j}{T}}} =$$



Softmax

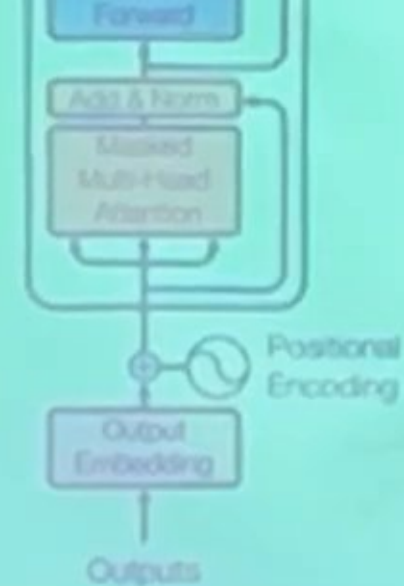
P_e

Stanford University

"Super Study Guide: Transformers and Large Language Models", Aridi et al.

Stanford

$$P_i = \frac{e^{\alpha_i}}{\sum_j e^{\alpha_j}}$$



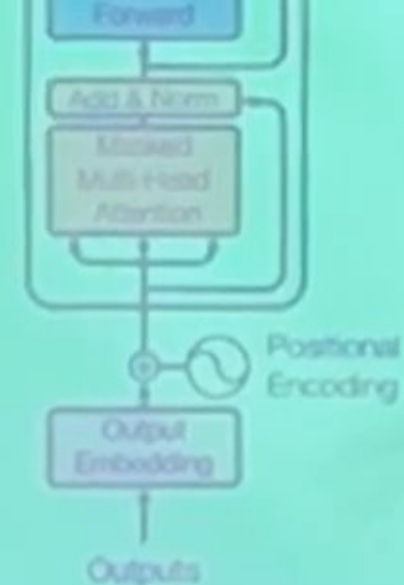
Stanford University

"Super Study Guide: Transformers and Large Language Models" Amidi et al.

Stanford



$$\frac{e^{\frac{x_1}{T}}}{\sum_j e^{\frac{x_j}{T}}}$$
$$\frac{e^{\frac{x_1 - x_2}{T}}}{\sum_j e^{\frac{x_j - x_2}{T}}}$$



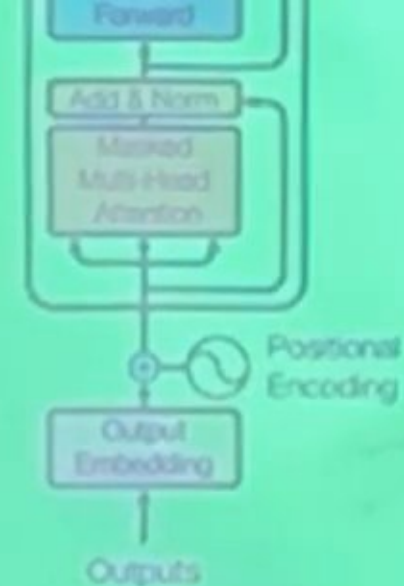
Stanford University

"Super Study Guide: Transformers and Large Language Models", Amidi et al.

Stanford



$$P_i = \frac{e^{\frac{\alpha_i}{T}}}{\sum_j e^{\frac{\alpha_j}{T}}}$$
$$\frac{\alpha_1 - \alpha_2}{T}$$
$$e^{\frac{\alpha_1 - \alpha_2}{T}}$$



Softmax

P_a

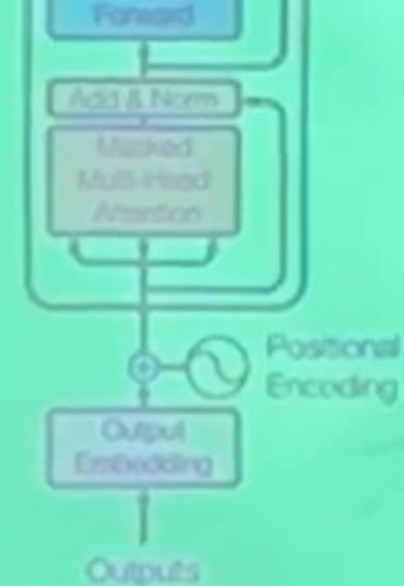
Stanford University

"Super Study Guide: Transformers and Large Language Models", Amidi et al.

Stanford



$$o_i = \frac{e^{\frac{x_i}{T}}}{\sum_j e^{\frac{x_j}{T}}}$$
$$= \frac{e^{\frac{x_i - x_k}{T}}}{\sum_j e^{\frac{x_j - x_k}{T}}}$$



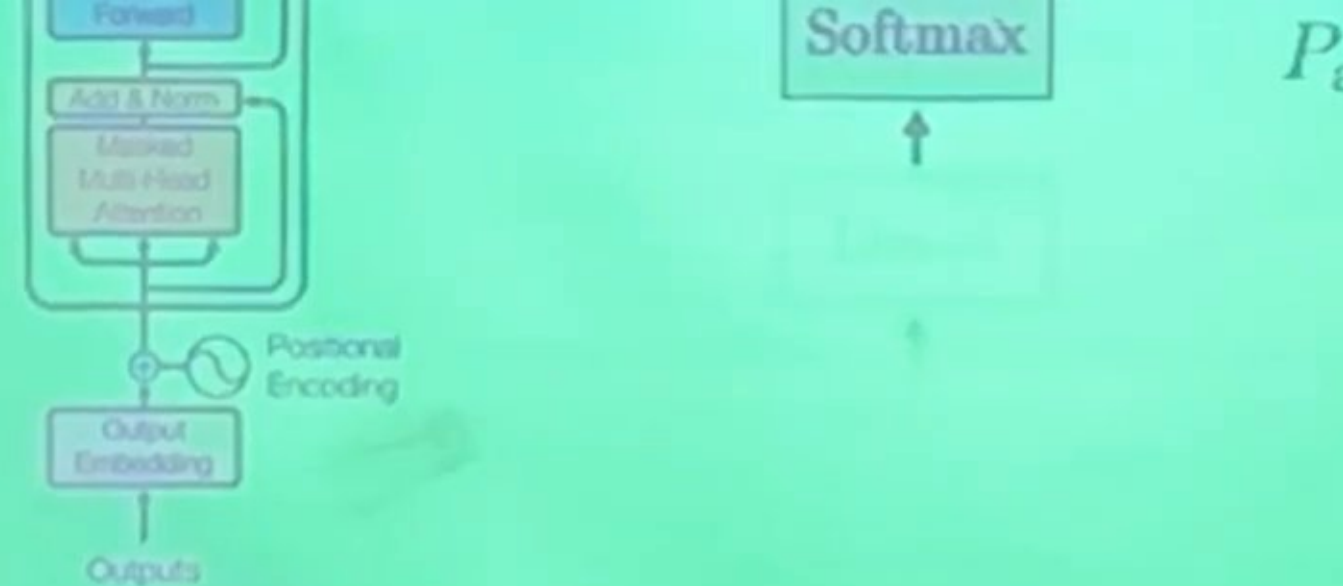
Stanford University

"Super Study Guide: Transformers and Large Language Models", Aridi et al.

Stanford

$$P_i = \frac{e^{\frac{\alpha_i}{\sqrt{d_k}}}}{\sum_j e^{\frac{\alpha_j}{\sqrt{d_k}}}}$$

$$= \frac{e^{\frac{\alpha_i - \alpha_k}{\sqrt{d_k}}}}{\sum_j e^{\frac{\alpha_j - \alpha_k}{\sqrt{d_k}}}}$$



Stanford University

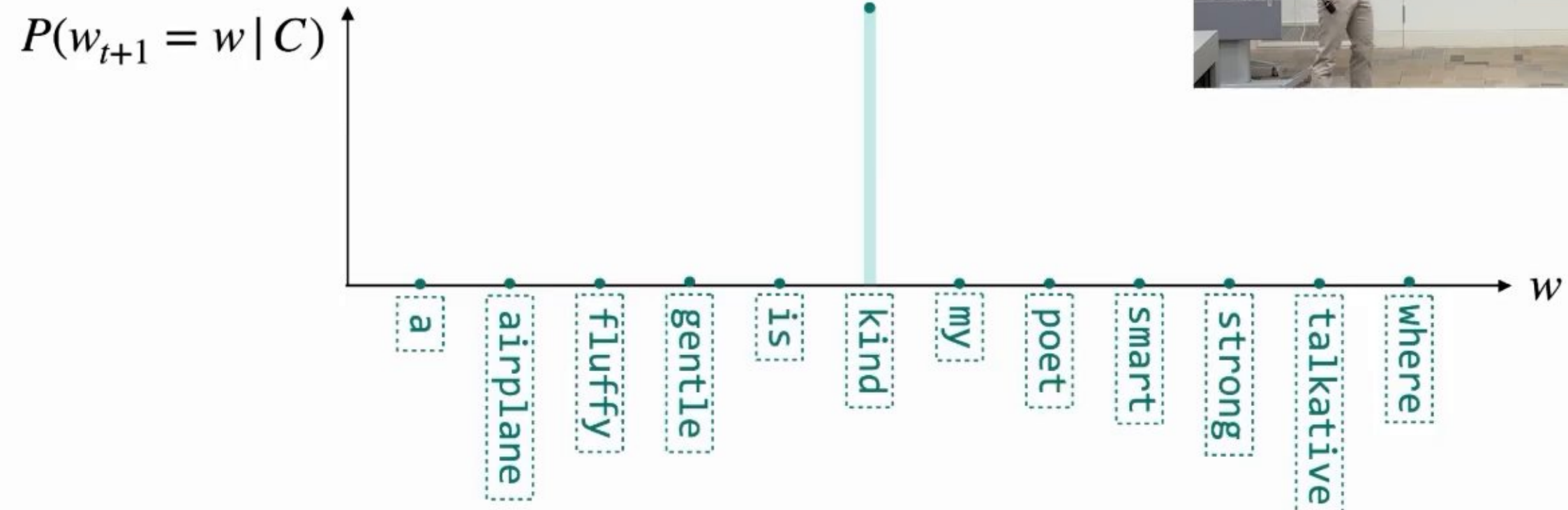
"Super Study Guide: Transformers and Large Language Models", Aridi et al.

Stanford

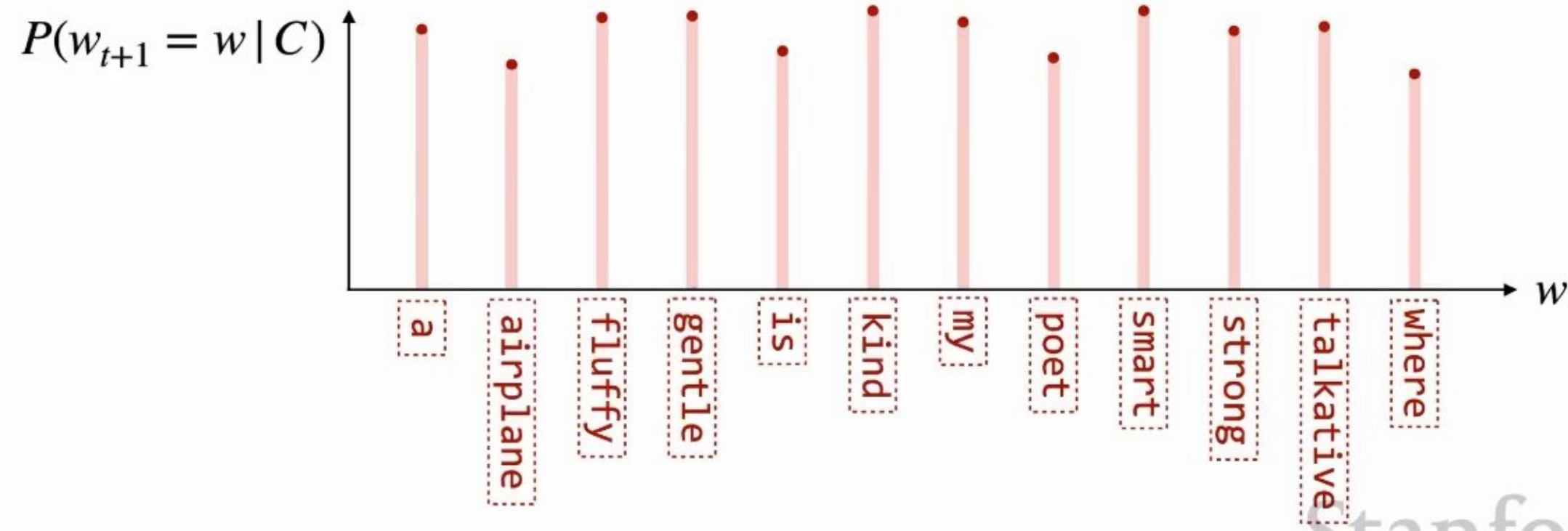
Impact of temperature on probabilities



Small T



High T



Stanford

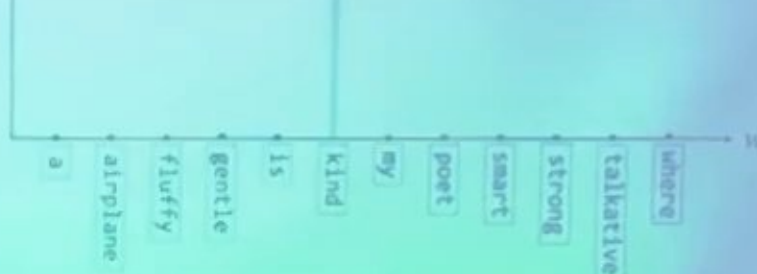
"Super Study Guide: Transformers and Large Language Models", Amidi et al., 2024.

Suggested reading: "Defeating Nondeterminism in LLM Inference", He et al., 2025.

Impact of temperature on probabilities

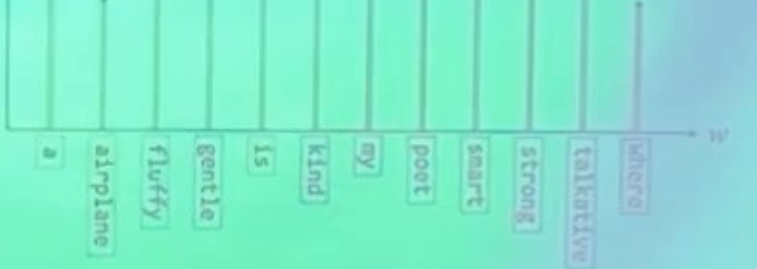
Small T

$$P(w_{t+1} = w | C)$$



High T

$$P(w_{t+1} = w | C)$$



Stanford University
"Super Study Guide: Transformers and Large Language Models", Aridi et al., 2024
Suggested reading: "Defeating Nondeterminism in LLM Inference", Ho et al., 2025

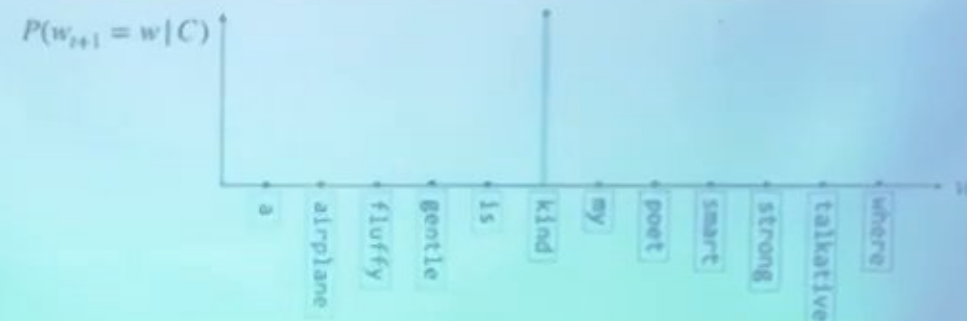
ICML

$$p_i = \frac{e^{\frac{x_i}{T}}}{\sum_j e^{\frac{x_j}{T}}} = \frac{e^{\frac{x_i - x_k}{T}}}{\sum_j e^{\frac{x_j - x_k}{T}}}$$

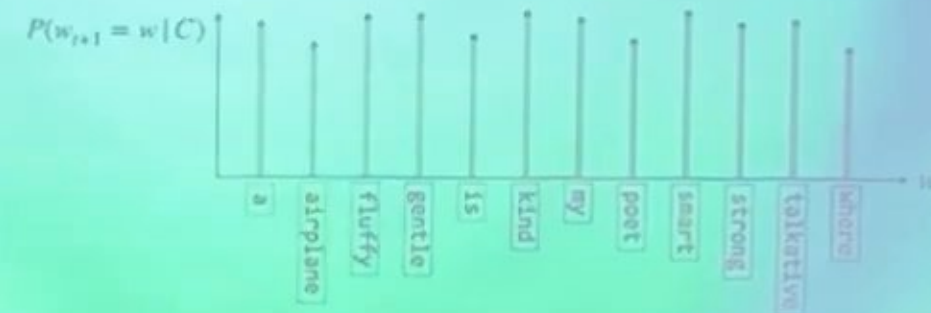
Stanford

Impact of temperature on probabilities

Small T



High T



Stanford University

"Super Study Guide: Transformers and Large Language Models", Amidi et al., 2024.
Suggested reading: "Defining Nonchaos in LLM Inference", Ho et al., 2025.

ICLME

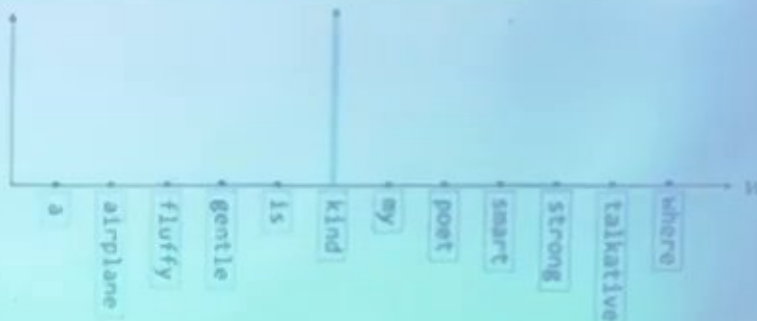
$$p_i = \frac{e^{\frac{x_i}{T}}}{\sum_j e^{\frac{x_j}{T}}}$$
$$= \frac{e^{\frac{x_i - x_k}{T}}}{\sum_j e^{\frac{x_j - x_k}{T}}}$$

Stanford

Impact of temperature on probabilities

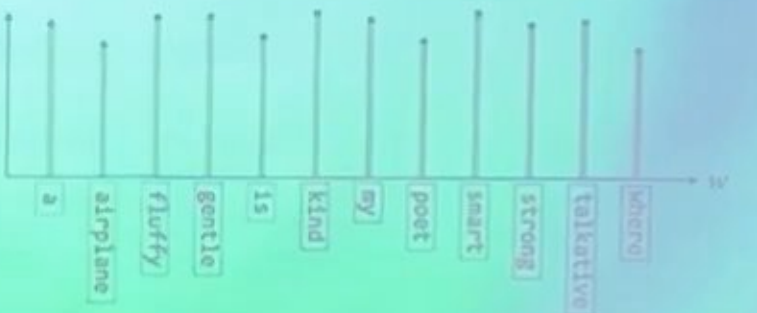
Small T

$P(w_{t+1} = w | C)$



High T

$P(w_{t+1} = w | C)$



Stanford University

"Super Study Guide: Transformers and Large Language Models", Amidi et al., 2024
Suggested reading: "Defining Nondeterminism in LLM Inference", Ho et al., 2025

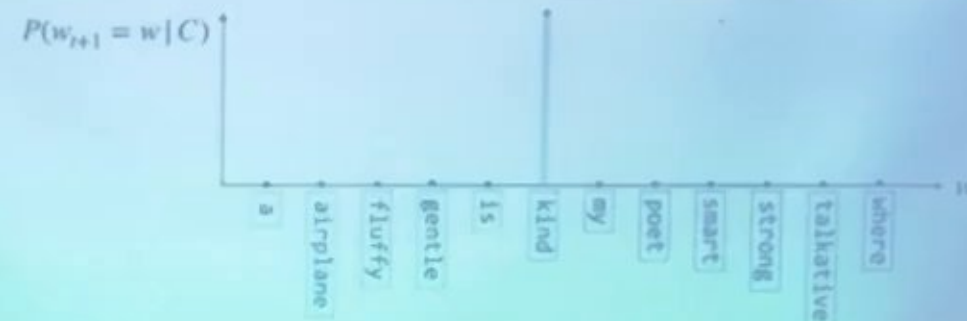
ICME

$$p_i = \frac{e^{\frac{\alpha_i}{T}}}{\sum_j e^{\frac{\alpha_j}{T}}}$$
$$= \frac{e^{\frac{\alpha_i - \alpha_k}{T}}}{\sum_j e^{\frac{\alpha_j - \alpha_k}{T}}}$$

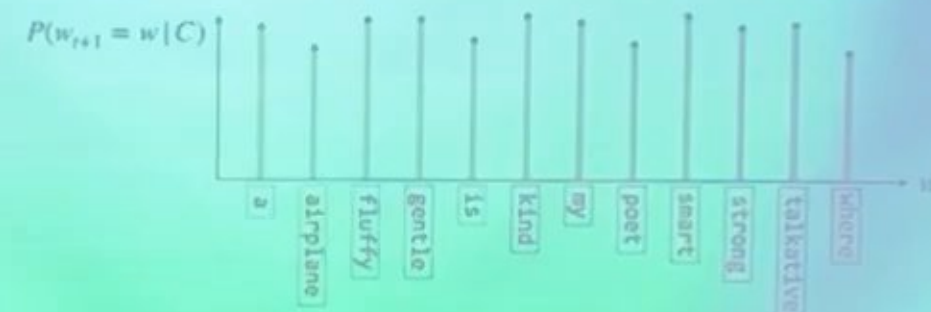
Stanford

Impact of temperature on probabilities

Small T



High T



Stanford University

"Super Study Guide: Transformers and Large Language Models", Amidi et al., 2024.
Suggested reading: "Delivering Nondeterminism in LLM Inference", Ho et al., 2025.

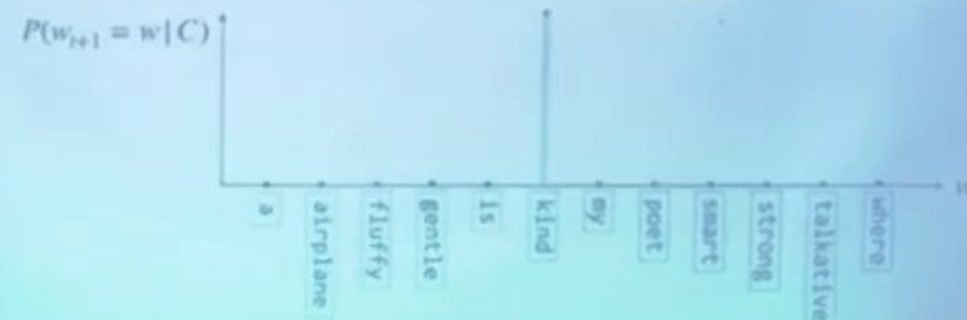
ICML

$$p_i = \frac{e^{\frac{x_i}{T}}}{\sum_j e^{\frac{x_j}{T}}}$$
$$= \frac{e^{\frac{x_i - x_k}{T}}}{\sum_j e^{\frac{x_j - x_k}{T}}}$$

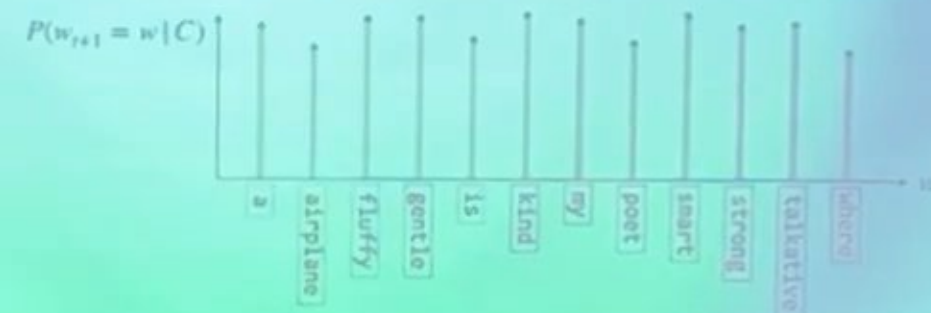
Stanford

Impact of temperature on probabilities

Small T



High T



Stanford University

"Super Study Guide: Transformers and Large Language Models", Amidi et al., 2024
Suggested reading: "Defining Nondeterminism in LLM Inference", Ho et al., 2025

ICML

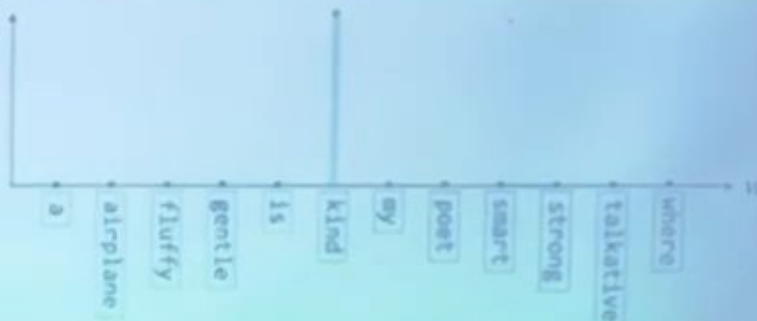
$$p_i = \frac{e^{\frac{x_i}{T}}}{\sum_j e^{\frac{x_j}{T}}}$$
$$= \frac{e^{\frac{x_i - x_k}{T}}}{\sum_j e^{\frac{x_j - x_k}{T}}}$$

Stanford

Impact of temperature on probabilities

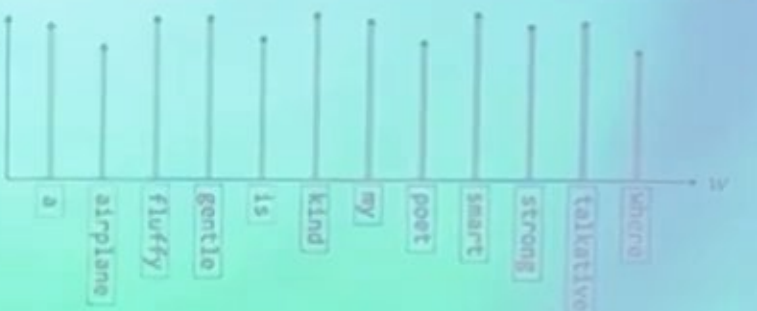
Small T

$P(w_{t+1} = w | C)$



High T

$P(w_{t+1} = w | C)$



Stanford University

"Super Study Guide: Transformers and Large Language Models", Amidi et al., 2024
Suggested reading: "Defining Nondeterminism in LLM Inference", Fie et al., 2025

ICML

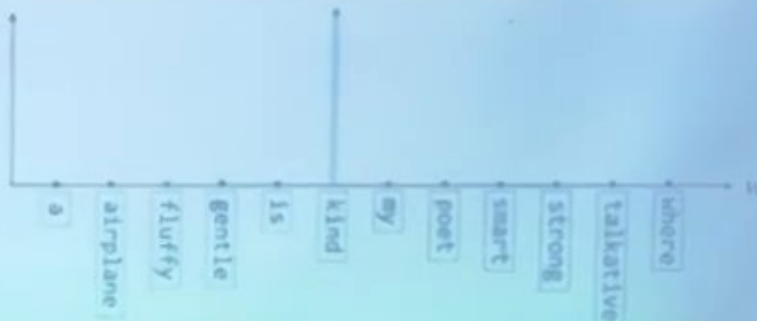
$$p_i = \frac{e^{\frac{x_i}{T}}}{\sum_j e^{\frac{x_j}{T}}}$$

Stanford

Impact of temperature on probabilities

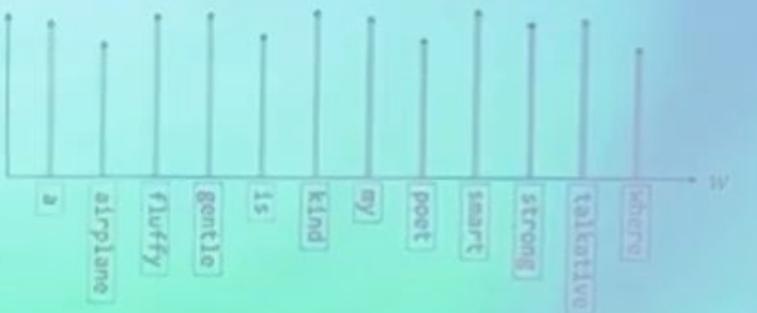
Small T

$P(w_{t+1} = w | C)$



High T

$P(w_{t+1} = w | C)$



Stanford University

"Super Study Guide: Transformers and Large Language Models", Amidi et al., 2024
Suggested reading: "Defining Nondeterminism in LLM Inference", Ho et al., 2025

ICME

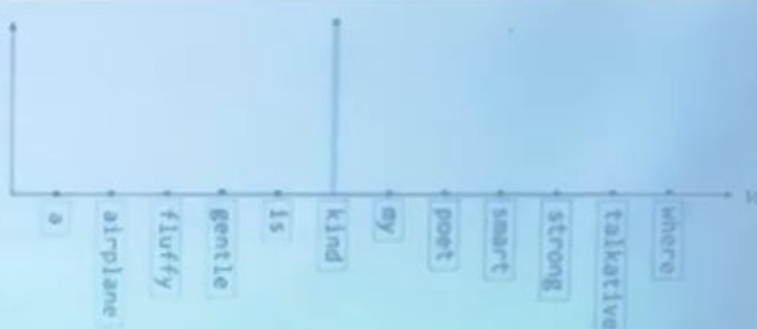
$$p_i = \frac{e^{\frac{\alpha_i}{T}}}{\sum_j e^{\frac{\alpha_j}{T}}} = \frac{e^{\frac{\alpha_i - \alpha_k}{T}}}{\sum_j e^{\frac{\alpha_j - \alpha_k}{T}}}$$

Stanford

Impact of temperature on probabilities

Small T

$P(w_{t+1} = w | C)$



High T

$P(w_{t+1} = w | C)$



Stanford University

"Super Study Guide: Transformers and Large Language Models", Amidi et al., 2024
Suggested reading: "Defining Nondeterminism in LLM Inference", Ho et al., 2025

ICME

$$p_i = \frac{e^{\frac{x_i}{T}}}{\sum_j e^{\frac{x_j}{T}}}$$

Stanford

Constraining the output via guided decoding



Motivation. Generate output in a specific format

Input prompt

*Generate a description of
my 33-year old teddy bear
who likes reading.*

Do this in JSON format.

Desired output (JSON)

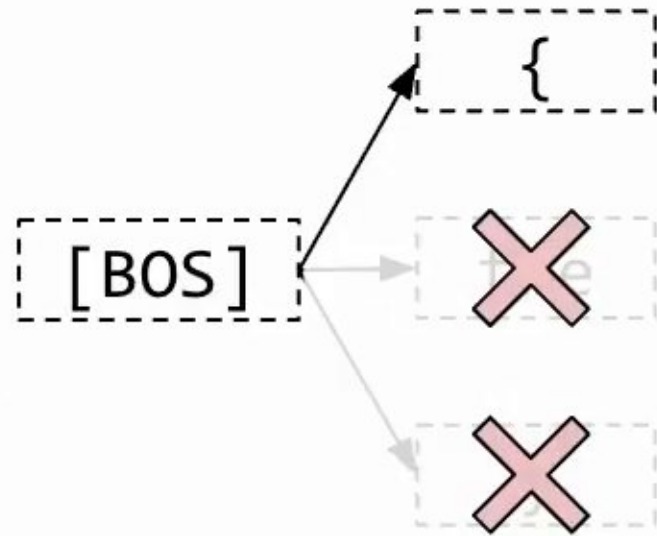
```
{  
  "first_name": "teddy",  
  "last_name": "bear",  
  "age": 33,  
  "hobby": "reading"  
}
```

Stanford

Constraining the output via guided decoding



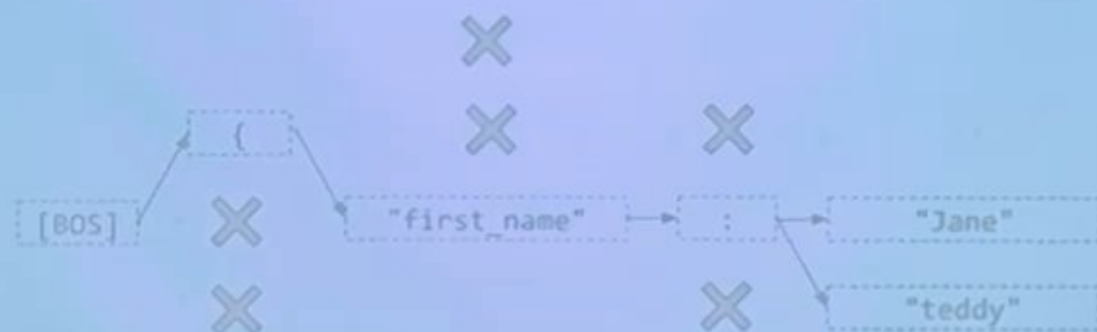
Idea. Only allow "valid" next tokens



Stanford

Constraining the output via guided decoding

Idea. Only allow "valid" next tokens



Stanford University

ICML

$$p_i = \frac{e^{\alpha_i}}{\sum_j e^{\alpha_j}}$$
$$= \frac{e^{\frac{\alpha_i - \alpha_k}{T}}}{\sum_j e^{\frac{\alpha_j - \alpha_k}{T}}}$$

Stanford



Transformers & Large Language Models

Stanford University

ICME

LLM overview

MoE-based LLMs

Response generation

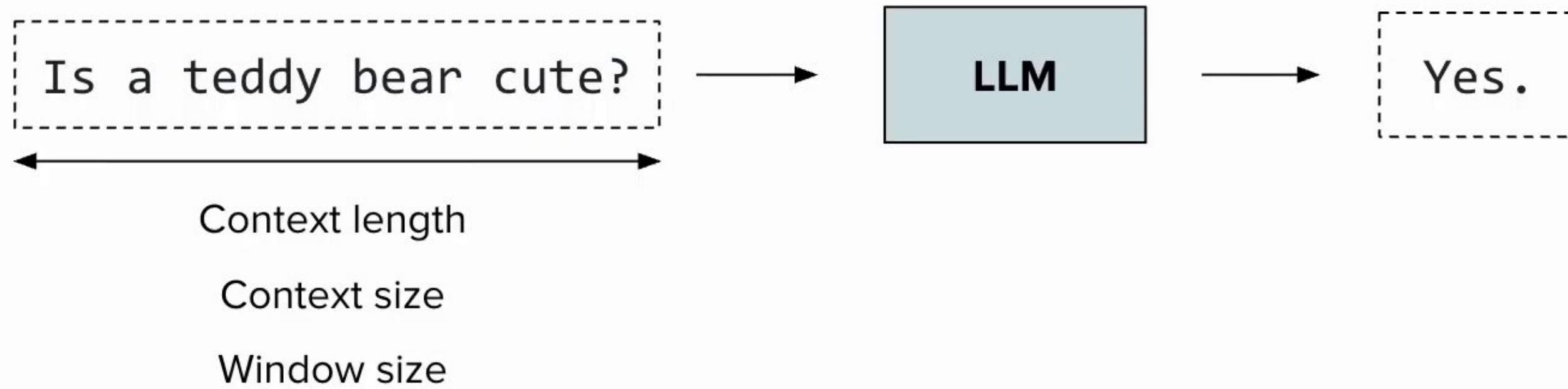
Prompting strategies

Inference optimizations

$$p_i = \frac{e^{\alpha_i}}{\sum_j e^{\alpha_j}}$$

Stanford

Terminology



Stanford

Terminology

Is a teddy bear cute?

LLM

Yes.

Context length

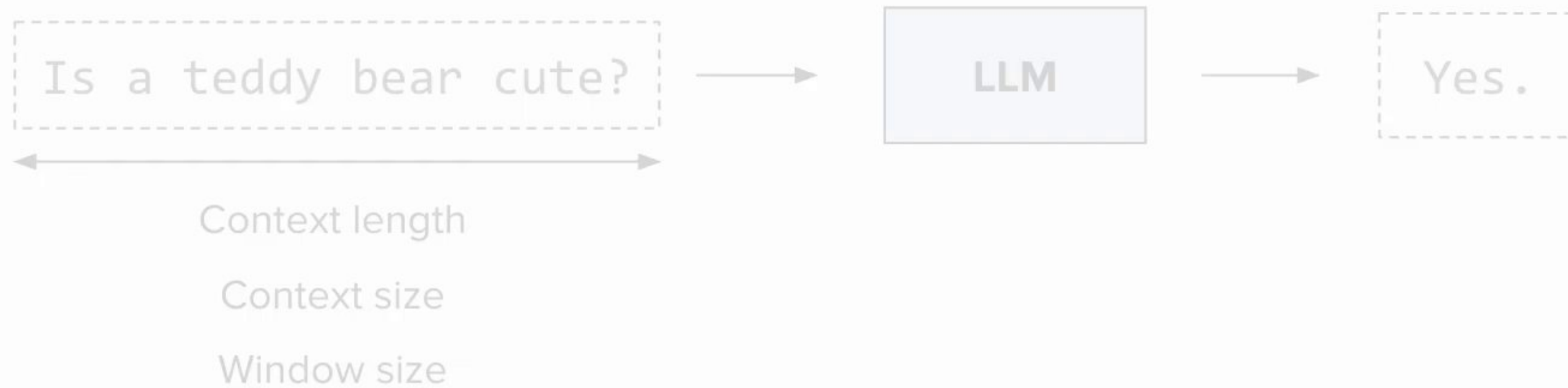
Context size

Window size

$$p_i = \frac{e^{\alpha_i}}{\sum_j e^{\alpha_j}}$$

Stanford

Terminology



Discussion. Orders of magnitude by:

- Input type
- Models

Stanford

Terminology

Discussion. Orders of magnitude by:

- Input type
- Models

← Beware of "context rot"!

Stanford University

"Context Rot: How Increasing Input Tokens Impacts LLM Performance", Hong et al., 2025

ICME

$$p_i = \frac{e^{\alpha_i}}{\sum_j e^{\alpha_j}}$$

Stanford

Main structure



My teddy bear had a long day and needs a bedtime story.

Context

Generate a bedtime story that takes place in a specific location.

Instructions

Location: Country of teddy bears

Input

The story needs to be suitable for teddy bears that are tired.

Constraints

Stanford

Main structure

My teddy bear had a long day and needs a bedtime story.

Context

Generate a bedtime story that takes place in a specific location.

Instructions

Location: Country of teddy bears

Input

The story needs to be suitable for teddy bears that are tired.

Constraints

Stanford University

"Super Study Guide: Transformers and Large Language Models" Amidi et al. 2024

ICME

$$p_i = \frac{e^{\frac{x_i}{T}}}{\sum_j e^{\frac{x_j}{T}}}$$

"Info

Stanford

In-context learning

ICL = In-Context Learning



Stanford

In-context learning

ICL = In-Context Learning

Stanford University

"Language Models are Few-Shot Learners" Brown et al., 2020

ICLME

$$p_i = \frac{e^{\frac{\alpha_i}{T}}}{e^{\frac{\alpha_1}{T}} + e^{\frac{\alpha_2}{T}} + \dots + e^{\frac{\alpha_n}{T}}}$$
$$p_i = \frac{e^{\frac{\alpha_i - \alpha_k}{T}}}{e^{\frac{\alpha_1 - \alpha_k}{T}} + e^{\frac{\alpha_2 - \alpha_k}{T}} + \dots + e^{\frac{\alpha_n - \alpha_k}{T}}}$$

Stanford

In-context learning

ICL = In-Context Learning



Zero-shot learning

- Question is asked without examples
- Performance heavily depends on performance of initial model

Few-shot learning

- Prompt contains examples of input / output
- Typically has a better performance

Stanford

In-context learning

ICL = In-Context Learning

Zero-shot learning

One-shot learning

Discussion. Showing examples in the prompt is *generally* better but:

- Requires effort
- Increases computational complexity & cost
- Increases latency

Stanford University

ICLME

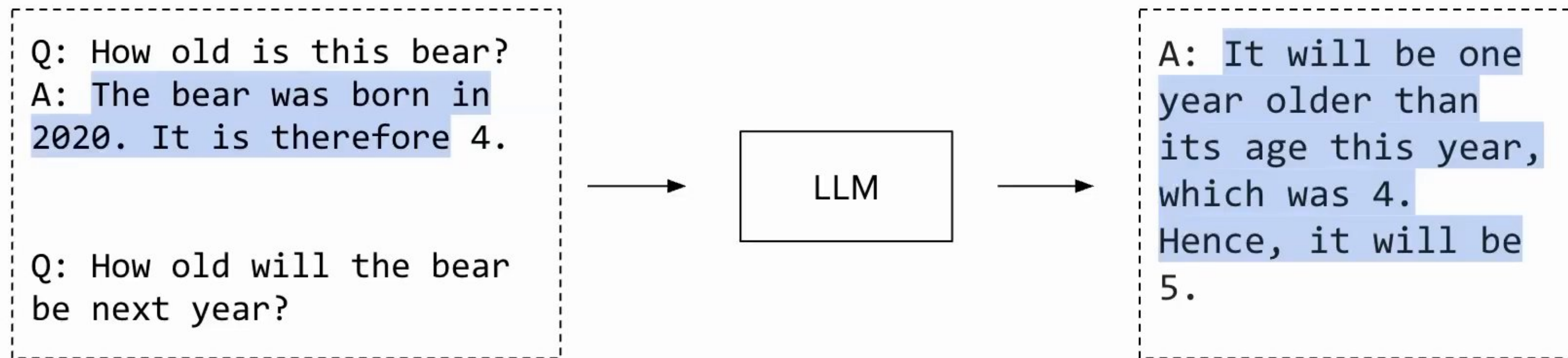
$$p_i = \frac{e^{\alpha_i}}{\sum}$$

Stanford

Chain of thought



Idea. Explaining reasoning helps in improving performance.



Discussion.

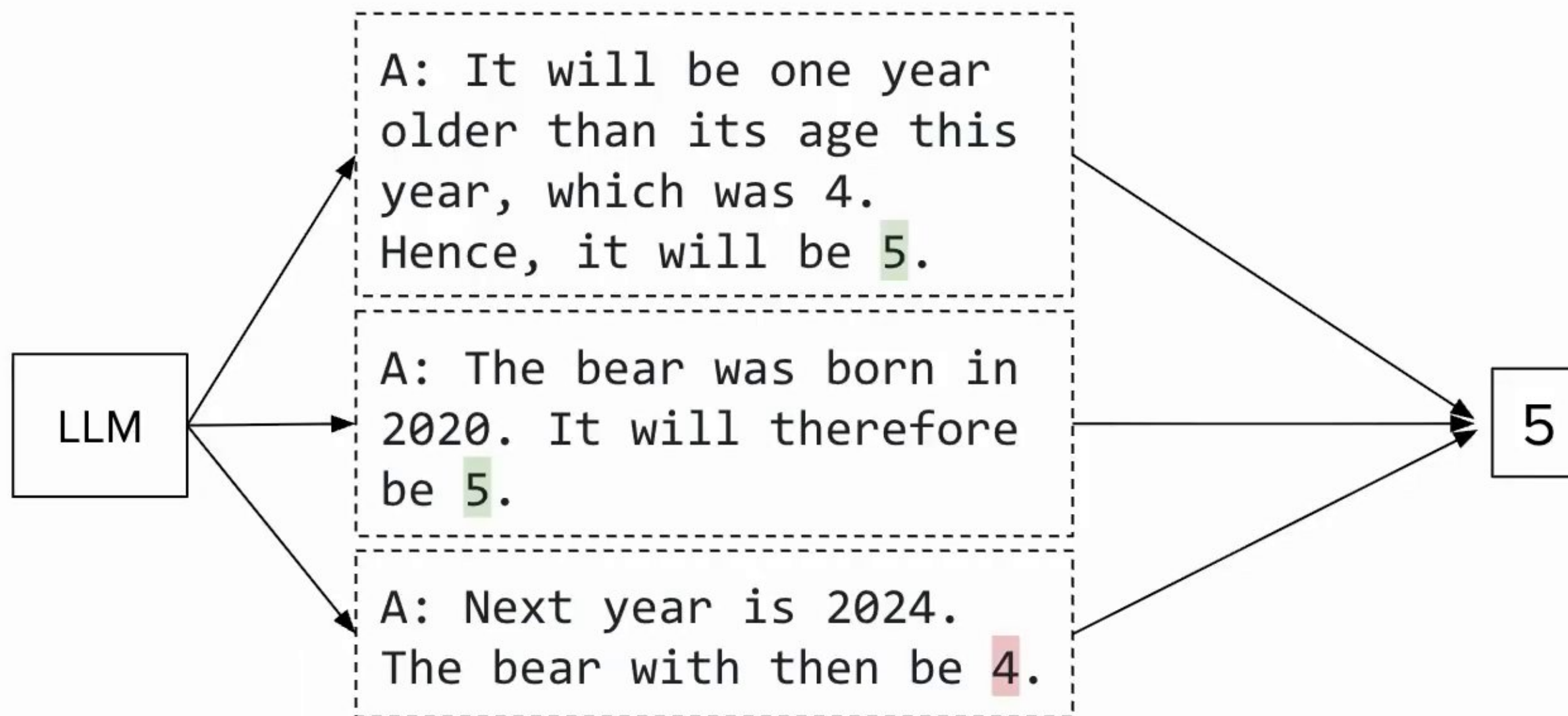
- Interpretability + explanation
- More tokens: higher cost and latency

Stanford

Self-consistency



Idea. Aggregating over reasoning paths improves performance.

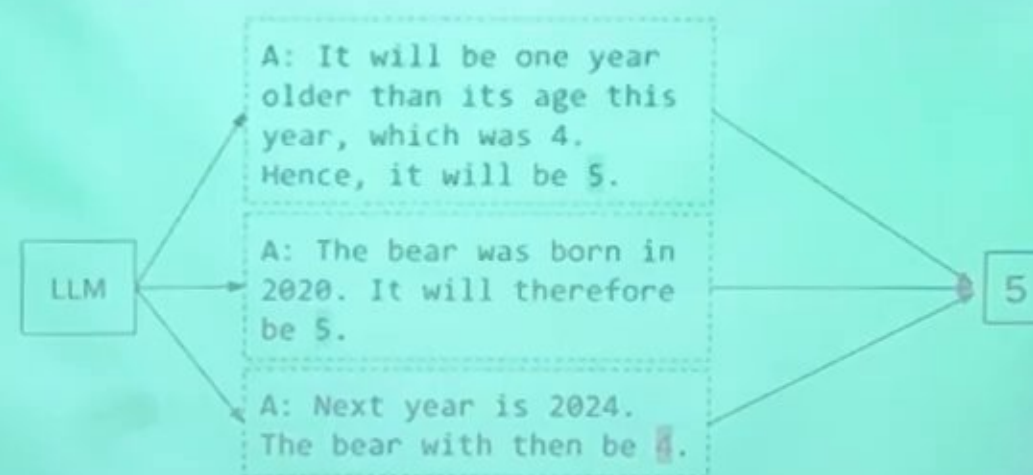


Discussion. Trade-off between performance and added cost.

Stanford

Self-consistency

Idea. Aggregating over reasoning paths improves performance.



Discussion. Trade-off between performance and added cost.

Stanford University

"Self-Consistency Improves Chain of Thought Reasoning in Language Models" Wang et al., 2022

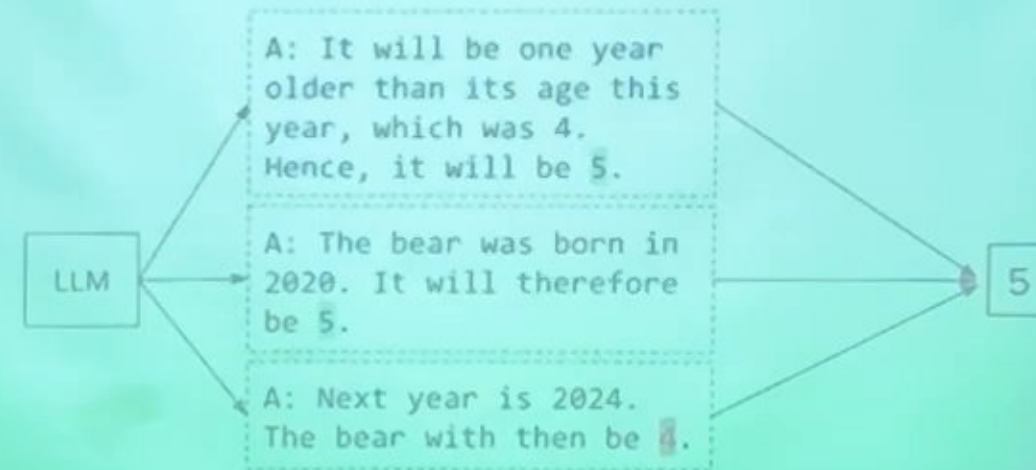
ICLME

$$p_i = \frac{e^{\alpha_i}}{\sum_j e^{\alpha_j}}$$

Stanford

Self-consistency

Idea. Aggregating over reasoning paths improves performance.



Discussion. Trade-off between performance and added cost.

Stanford University

"Self-Consistency Improves Chain of Thought Reasoning in Language Models" Wang et al., 2022

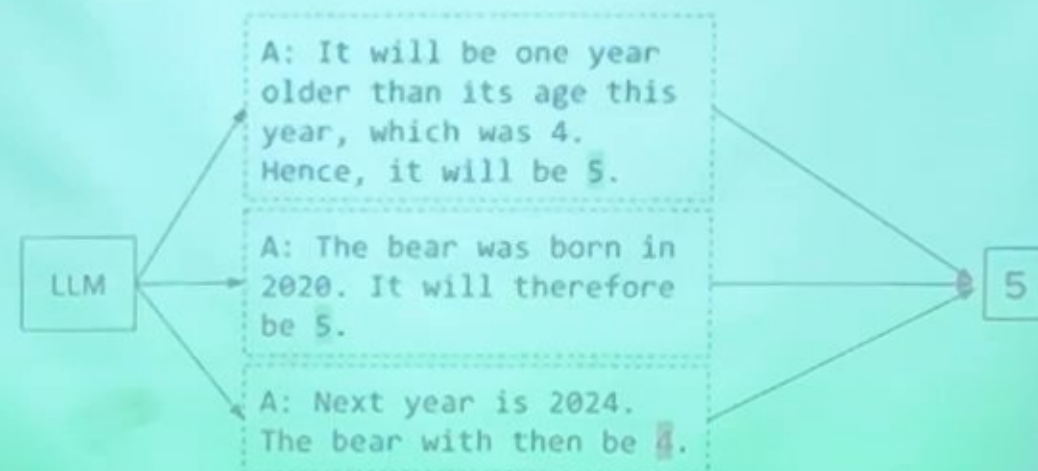
ICLME

$$p_i = \frac{e^{\frac{\alpha_i}{T}}}{\sum_j e^{\frac{\alpha_j}{T}}}$$
$$= \frac{e^{\frac{\alpha_1 - \alpha_2}{T}}}{\sum_j e^{\frac{\alpha_j - \alpha_2}{T}}}$$

Stanford

Self-consistency

Idea. Aggregating over reasoning paths improves performance.



Discussion. Trade-off between performance and added cost.

Stanford University

"Self-Consistency Improves Chain of Thought Reasoning in Language Models" Wang et al., 2022

ICME

$$P_i = \frac{e^{\frac{\alpha_i}{T}}}{\sum_j e^{\frac{\alpha_j}{T}}}$$

Stanford



Transformers & Large Language Models



LLM overview

MoE-based LLMs

Response generation

Prompting strategies

Inference optimizations

Stanford

Challenges

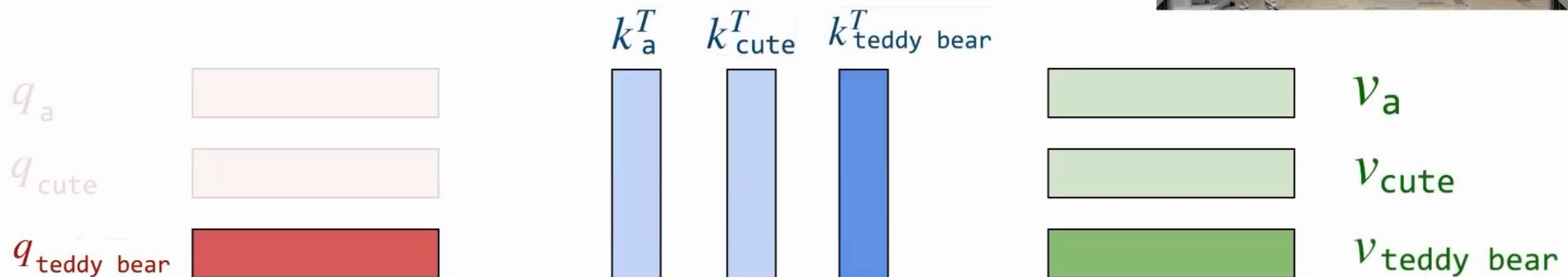
Motivation. Computations are expensive, any way to reduce cost?

Categories. Many dimensions to optimize for.



Stanford

Sharing attention heads



Stanford

Sharing attention heads



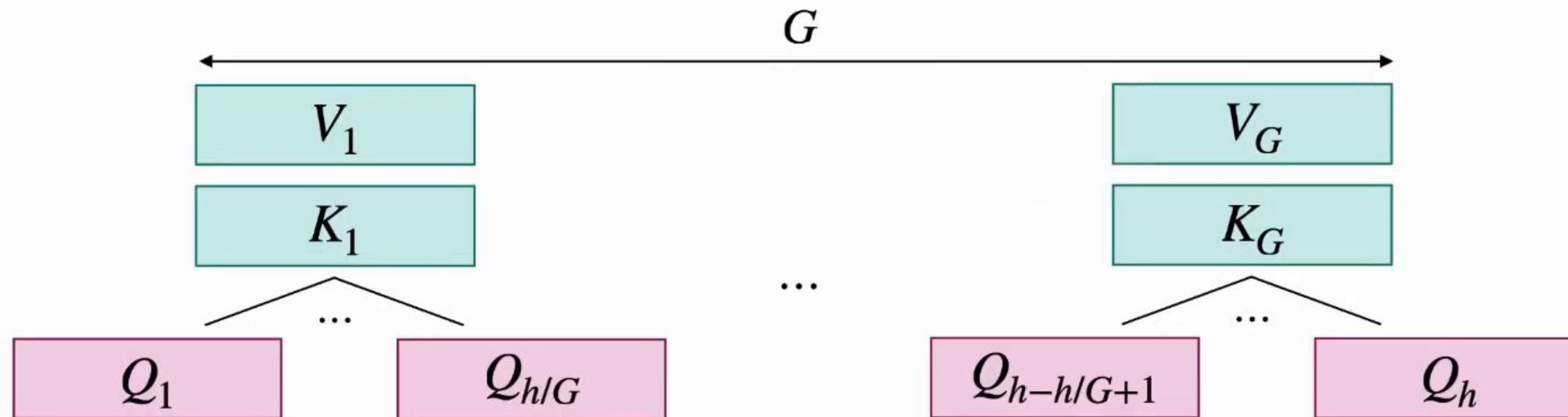
$$p_i = \frac{e^{\alpha_i}}{\sum_j e^{\alpha_j}}$$

Stanford

Sharing attention heads



Idea. Share key/value attention heads within groups of queries



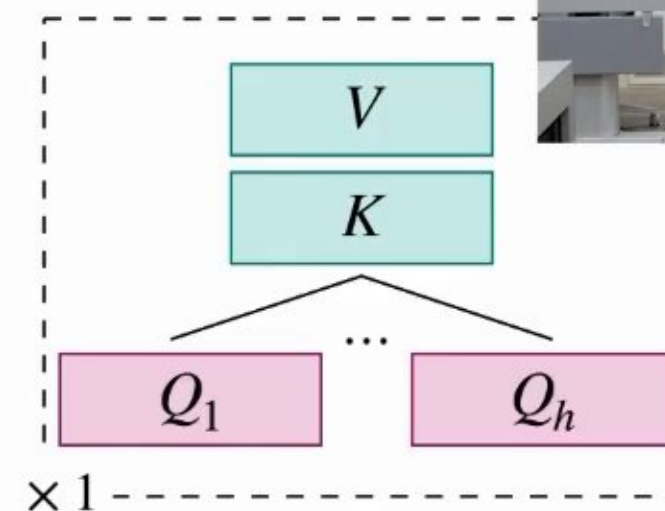
Stanford

Sharing attention heads



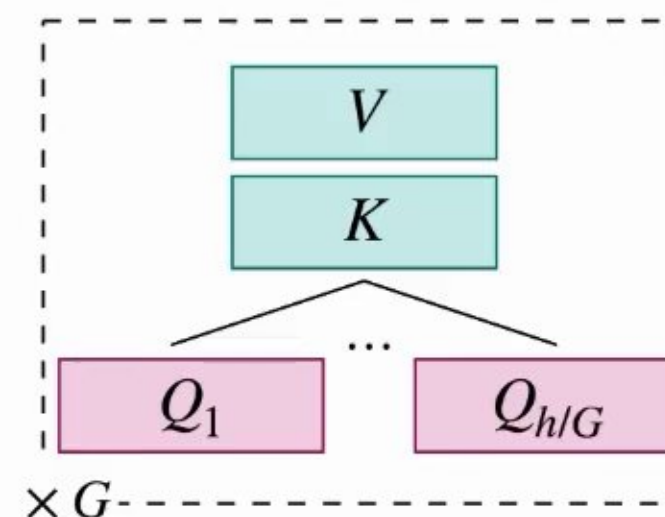
$$G = 1$$

Multi-Query Attention
(MQA)



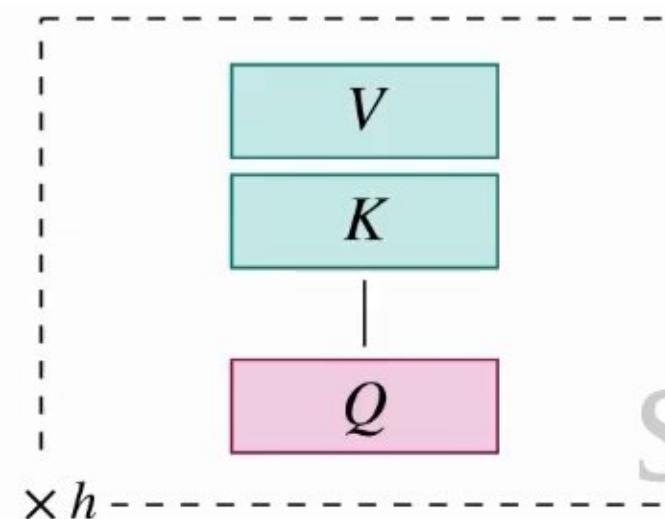
$$1 < G < h$$

Group-Query Attention
(GQA)



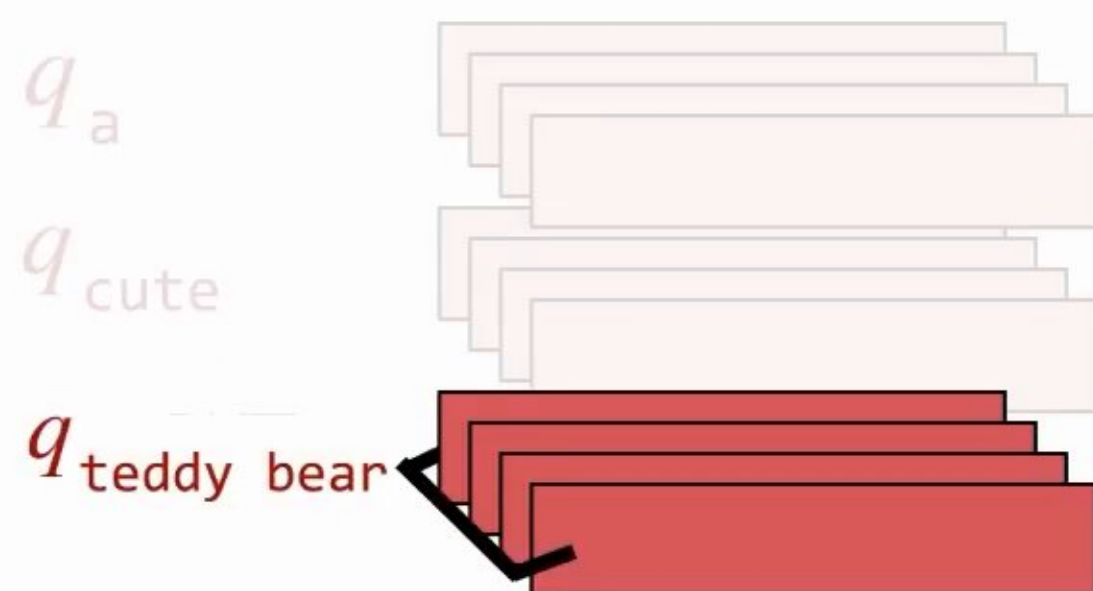
$$G = h$$

Multi-Head Attention
(MHA)



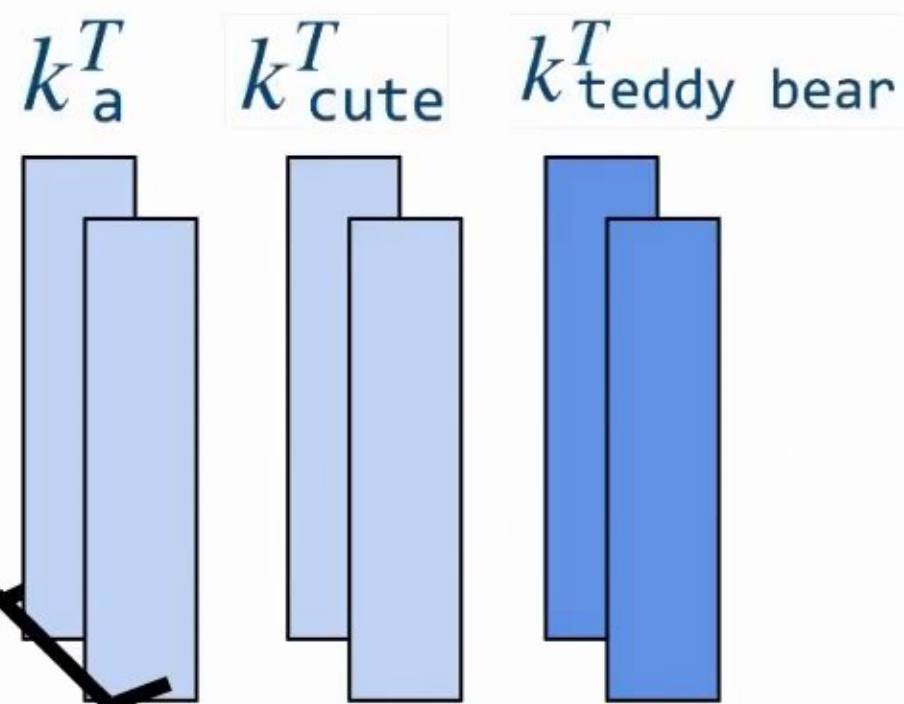
Stanford

Sharing attention heads with GQA



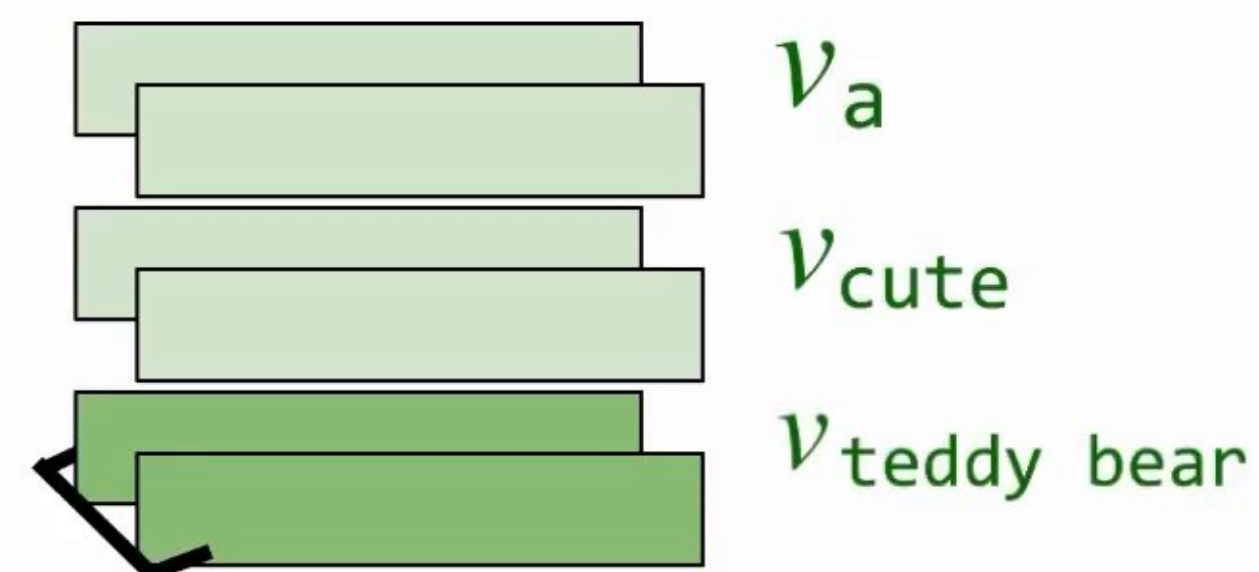
Number of
query heads

$$= h$$



Number of
key heads

$$= G < h$$



Number of
value heads

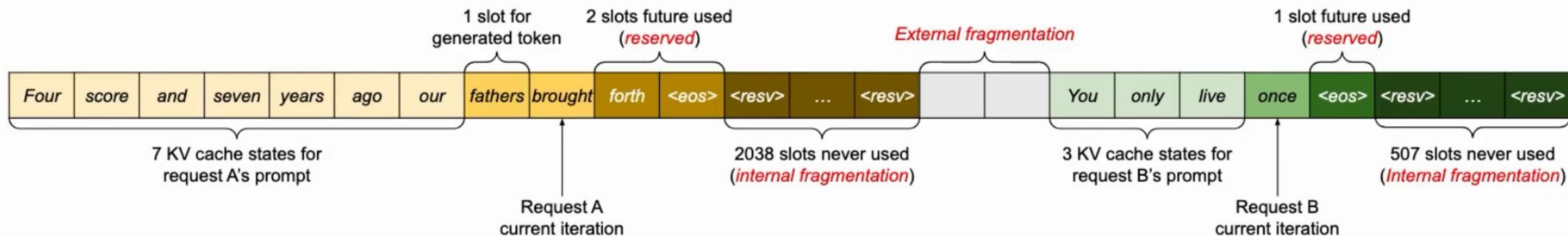
$$= G < h$$

Stanford

Manage memory with PagedAttention



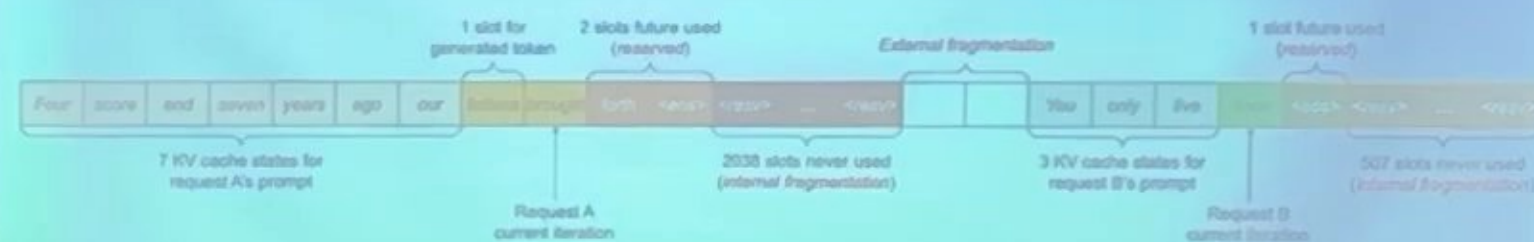
Observation. Lots of memory waste when storing KV cache.



Stanford

Manage memory with PagedAttention

Observation. Lots of memory waste when storing KV cache.



Stanford University

Figure from "Efficient Memory Management for Large Language Model Serving with PagedAttention", by Kuan et al., 2023.

ICME

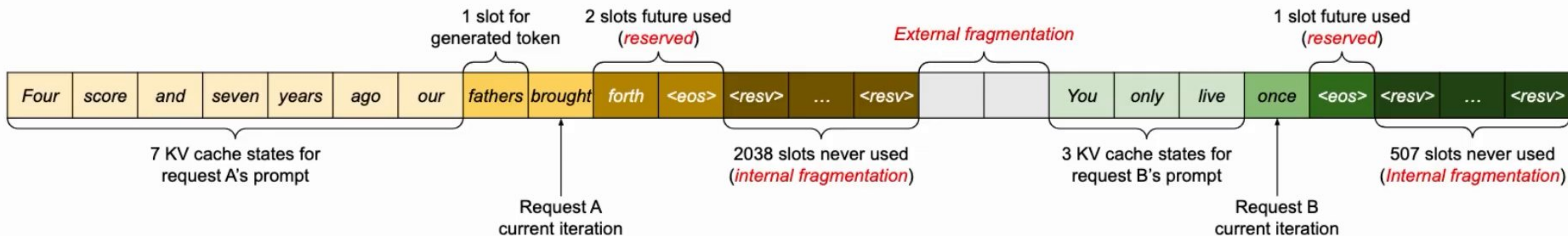
$$P_i = \frac{e^{\alpha_i}}{\sum_j e^{\alpha_j}}$$

Stanford

Manage memory with PagedAttention



Observation. Lots of memory waste when storing KV cache.

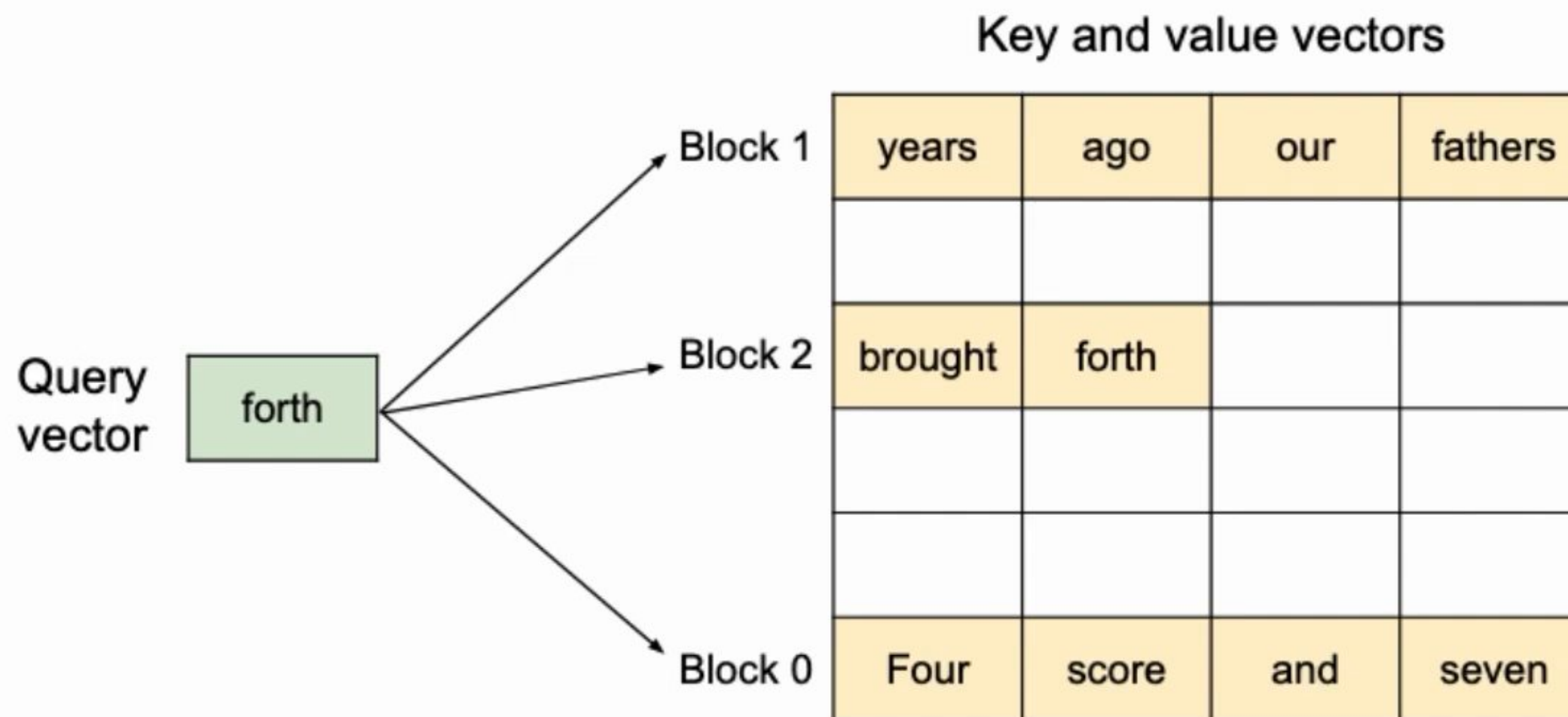


Stanford

Manage memory with PagedAttention



Idea. Store K and V in non-contiguous space to minimize wasted memory.

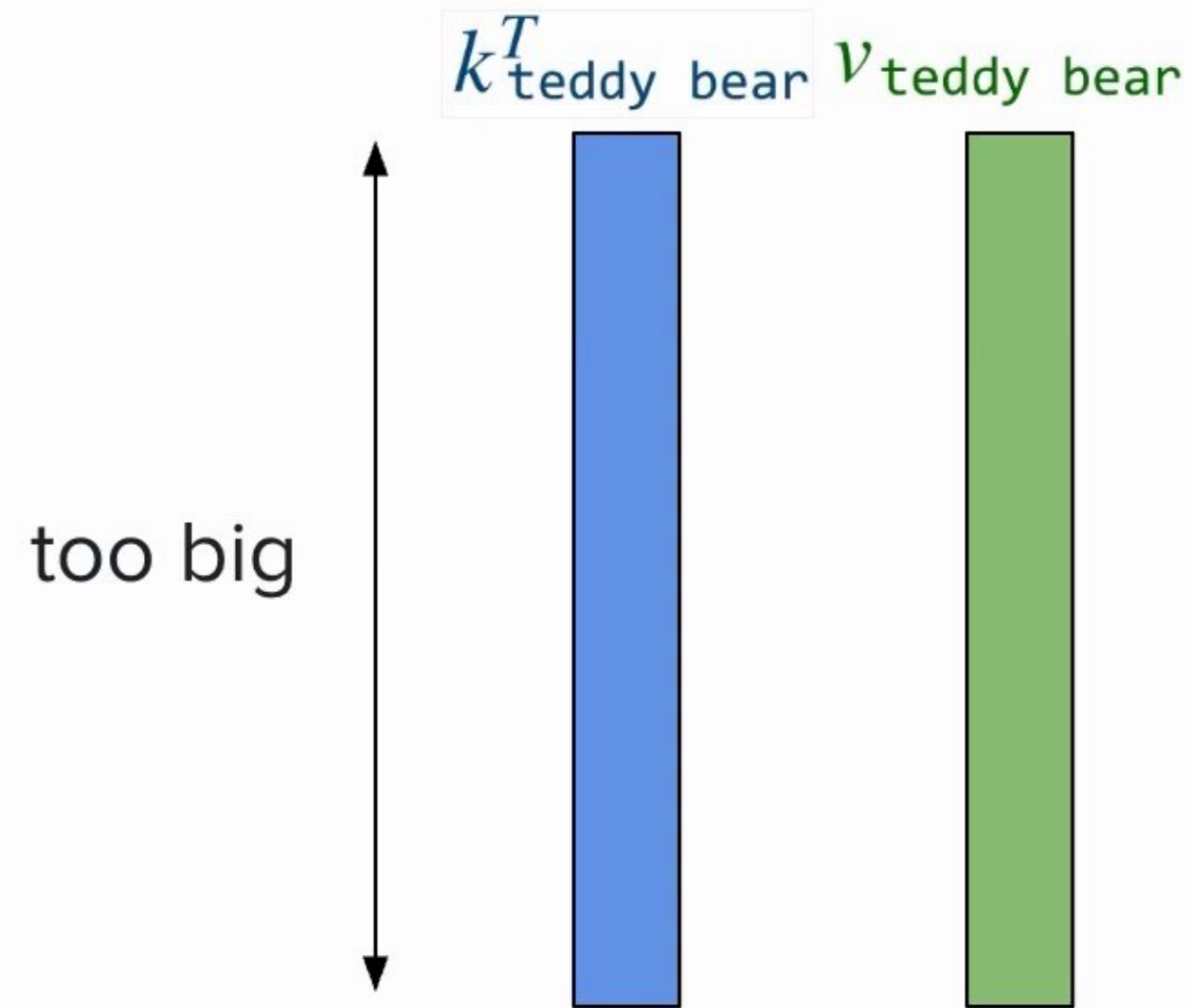


Stanford

Reduce memory of KV cache with latent a



Goal. Reduce dimension of K and V stored in memory.

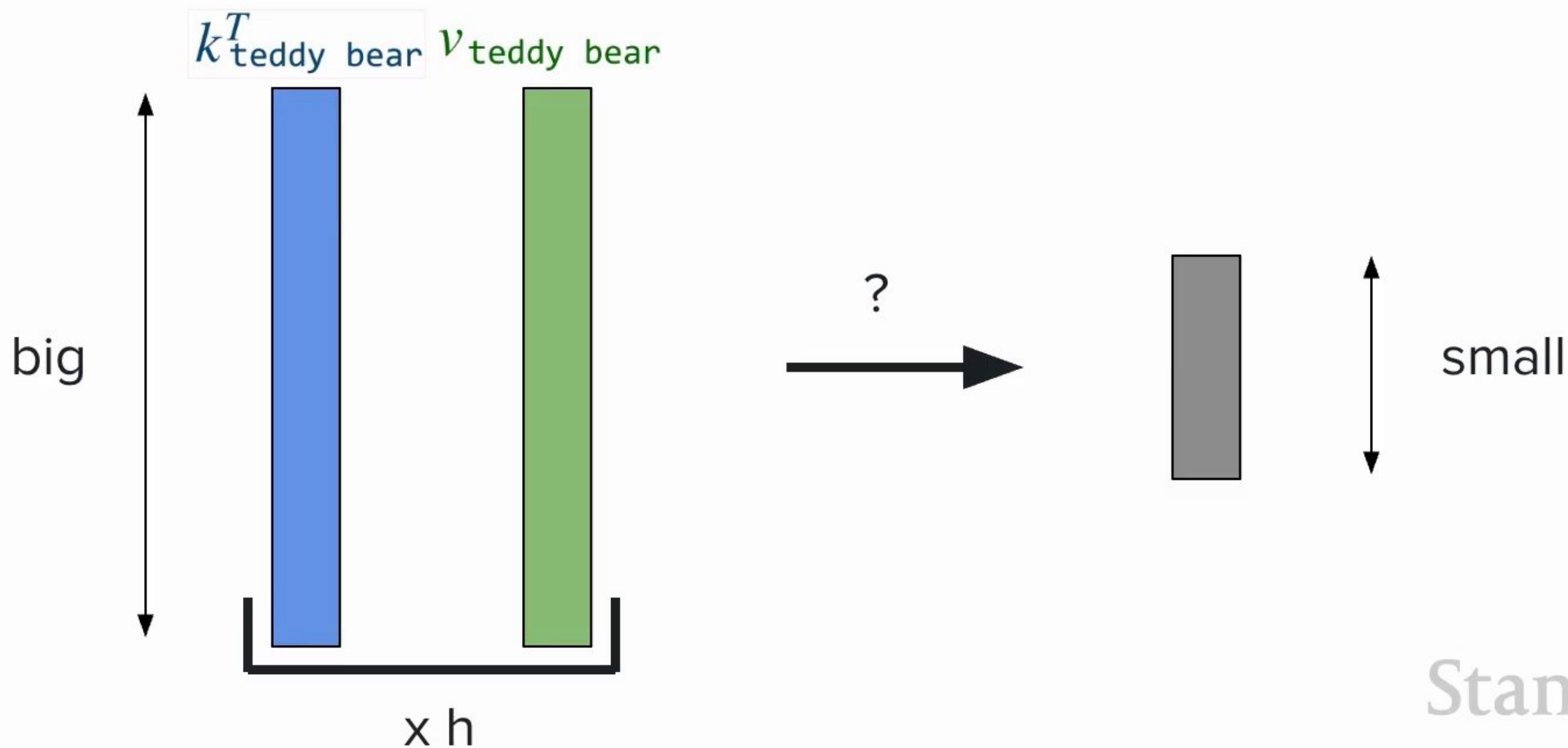


Stanford

Reduce memory of KV cache with latent a



Solution. Store compressed representations instead!



Stanford

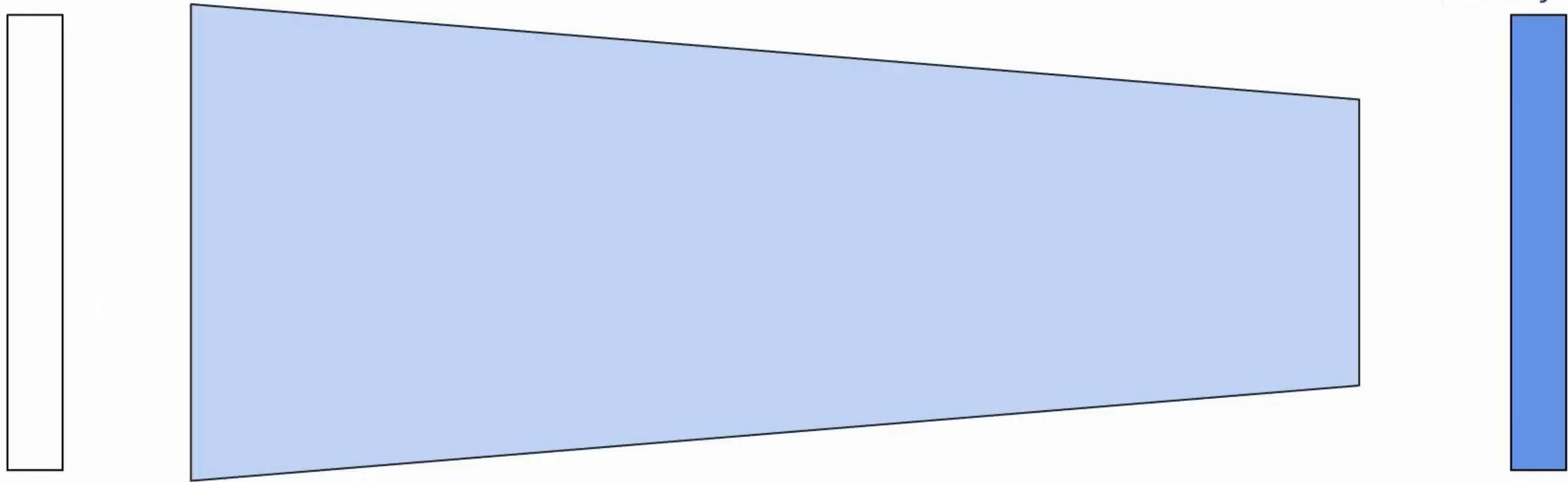
Reduce memory of KV cache with latent a

BEFORE



teddy bear

$k_{\text{teddy bear}}^T$



$x \cdot h$

Stanford

Token generation tricks



$$P_i = \frac{e^{\alpha_i}}{\sum_j e^{\alpha_j}}$$

Stanford

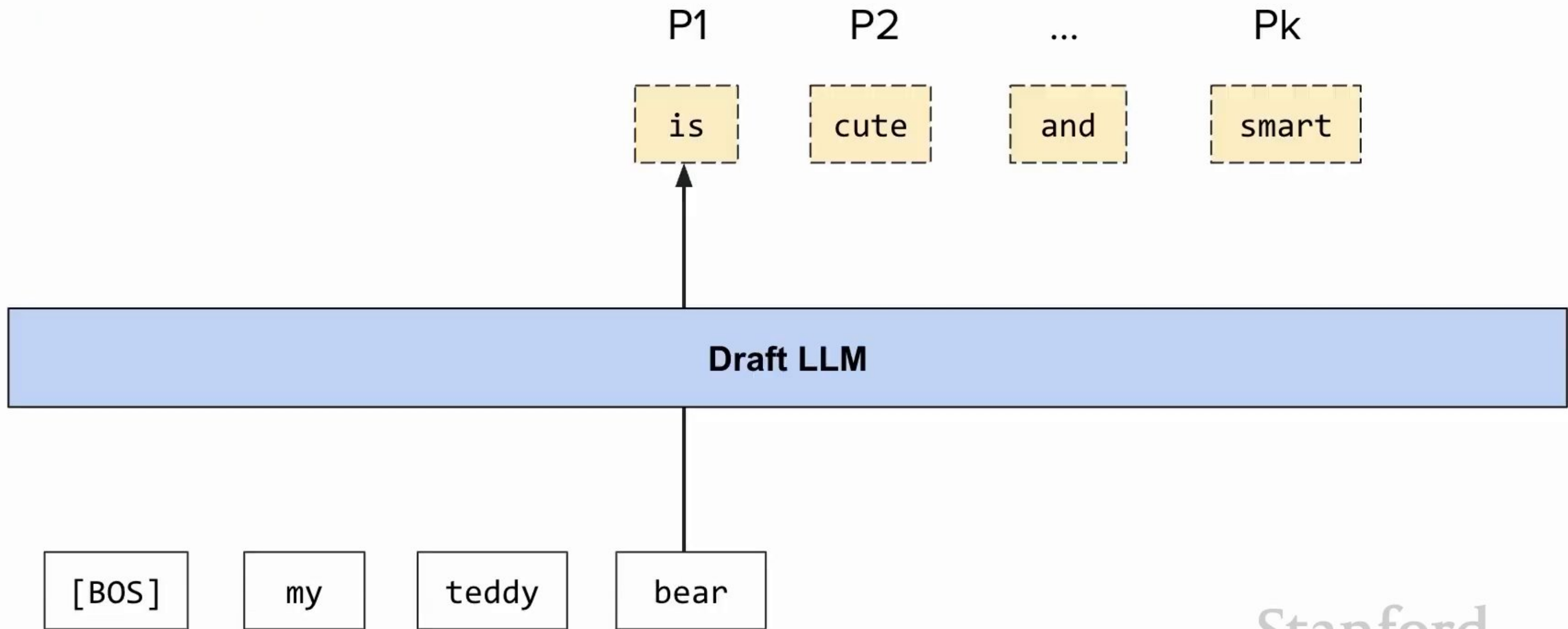
Speeding up decoding with speculative decoding



Idea. Use a draft (small) model to generate tokens validated by a target (big) model.

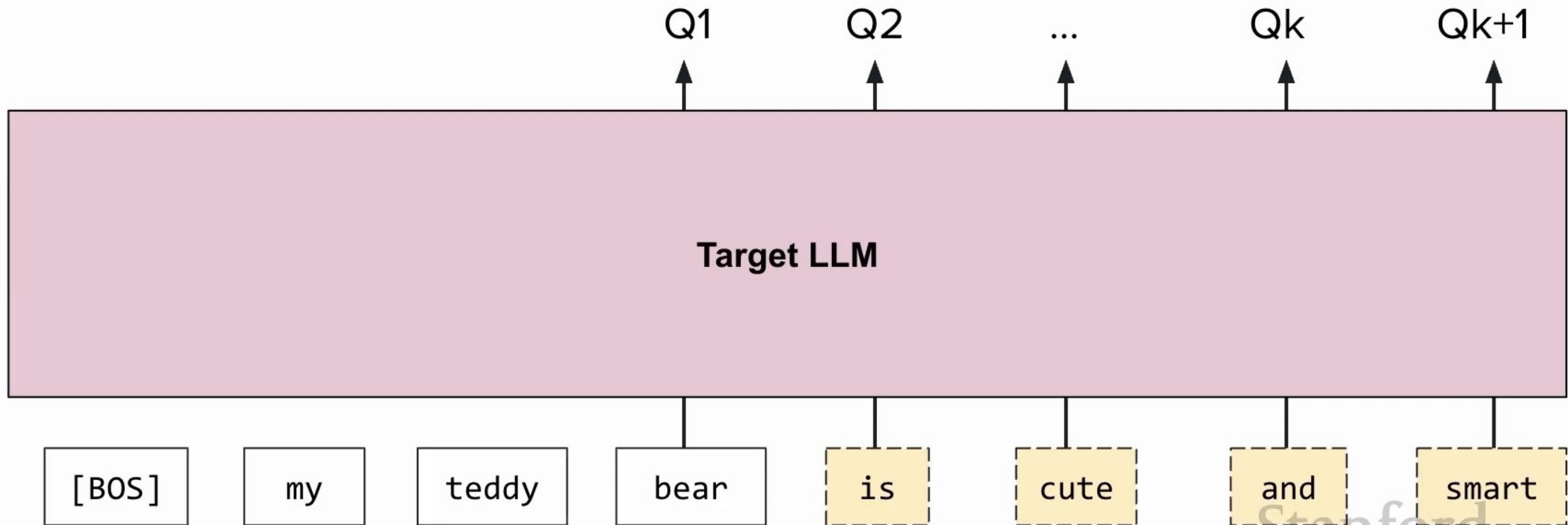
Stanford

Speeding up decoding with speculative



Stanford




Speeding up decoding with speculative



Speeding up decoding with speculative



Sampling algorithm. We distinguish the following cases:

- If $Q_i(\text{token}) \geq P_i(\text{token})$ token 
- Otherwise
 - with probability $Q_i(\text{token}) / P_i(\text{token})$ token 
 - with probability $1 - Q_i(\text{token}) / P_i(\text{token})$ token 

If a rejection happens, re-sample next token with distribution $[Q_i - P_i]^+$ and exit

Speeding up decoding with speculative decoding

Sampling algorithm. We distinguish the following cases:

- If $Q_i(\text{token}) \geq P_i(\text{token})$ token ✓
- Otherwise
 - with probability $Q_i(\text{token}) / P_i(\text{token})$ token ✓
 - with probability $1 - Q_i(\text{token}) / P_i(\text{token})$ token ✗

If a rejection happens, re-sample next token with distribution $[Q_i - P_i]^+$ and exit

Stanford University "Accelerating Large Language Model Decoding with Speculative Sampling", Chen et al., 2023

ICML




$$P_i = \frac{e^{\alpha_i}}{\sum_j e^{\alpha_j}}$$

Stanford

Speeding up decoding with speculative

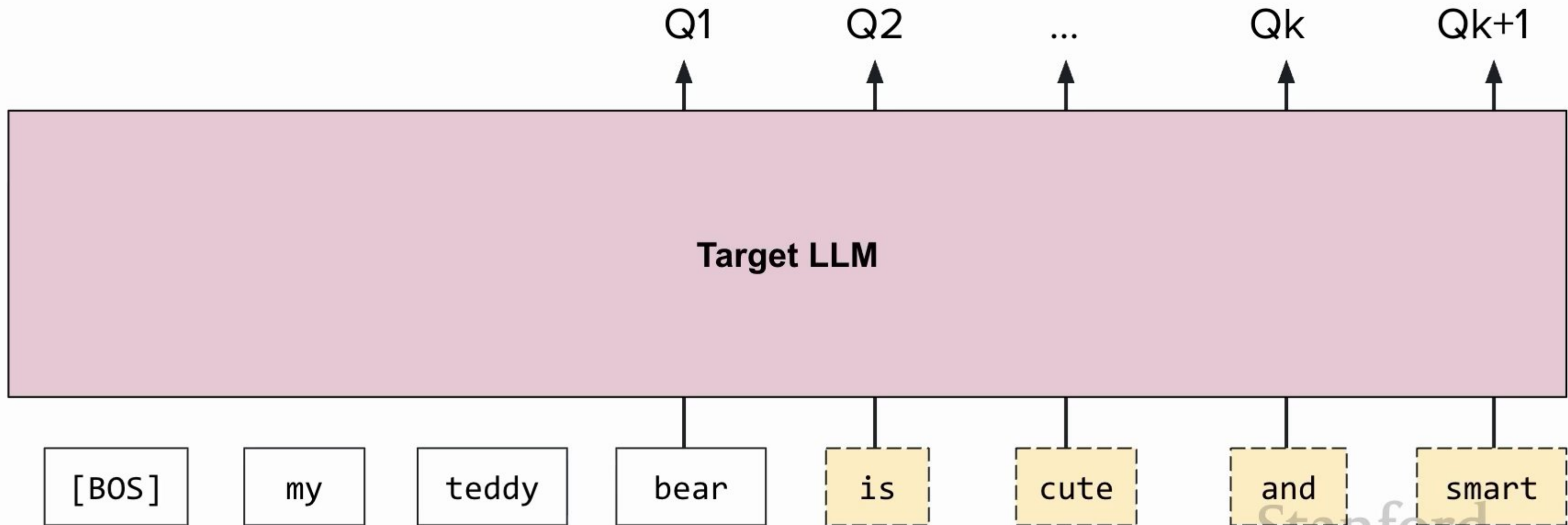


Sampling algorithm. We distinguish the following cases:

- If $Q_i(\text{token}) \geq P_i(\text{token})$ token 
- Otherwise
 - with probability $Q_i(\text{token}) / P_i(\text{token})$ token 
 - with probability $1 - Q_i(\text{token}) / P_i(\text{token})$ token 

If a rejection happens, re-sample next token with distribution $[Q_i - P_i]^+$ and exit

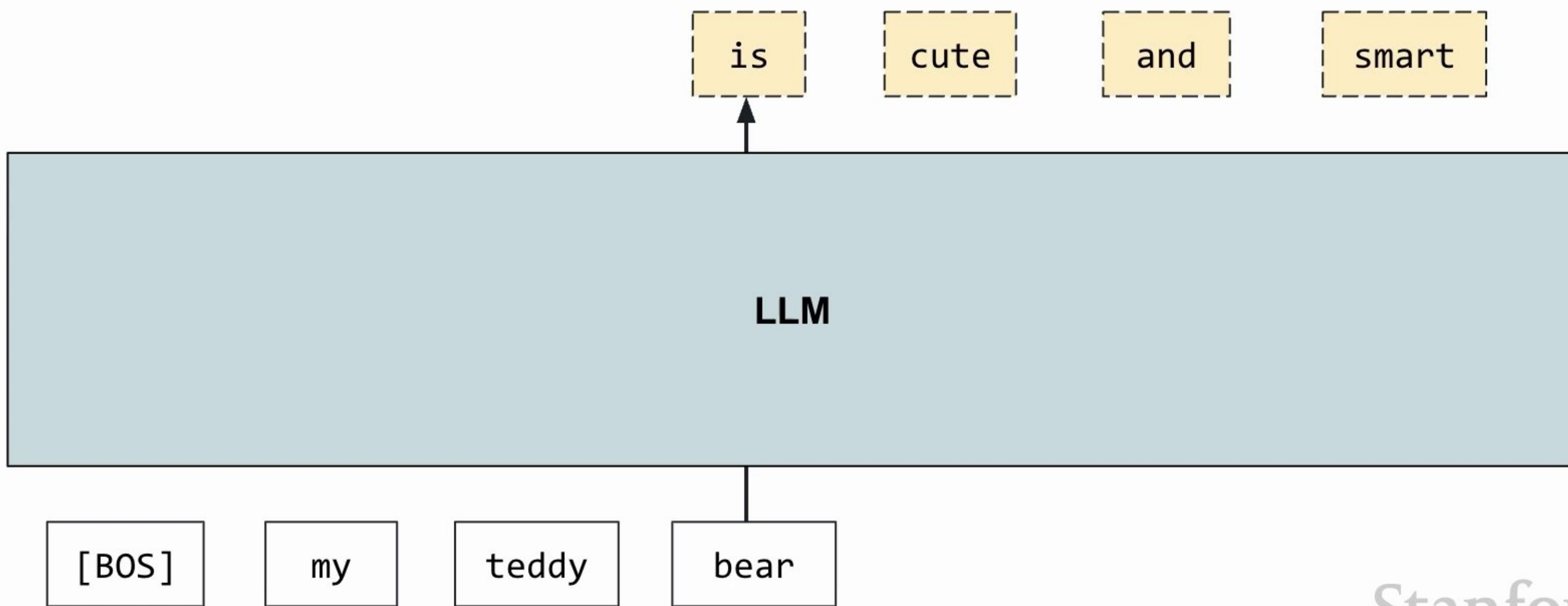
Speeding up decoding with speculative



Generate several tokens at once via MTP

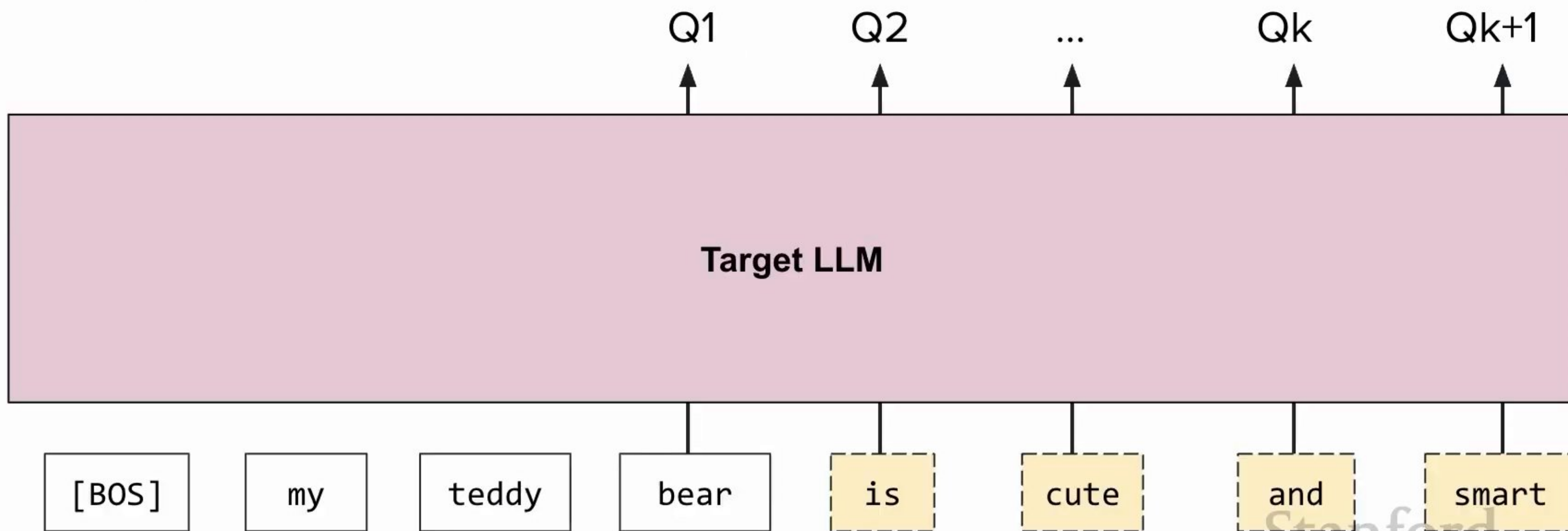


MTP = Multi-Token Prediction



Stanford

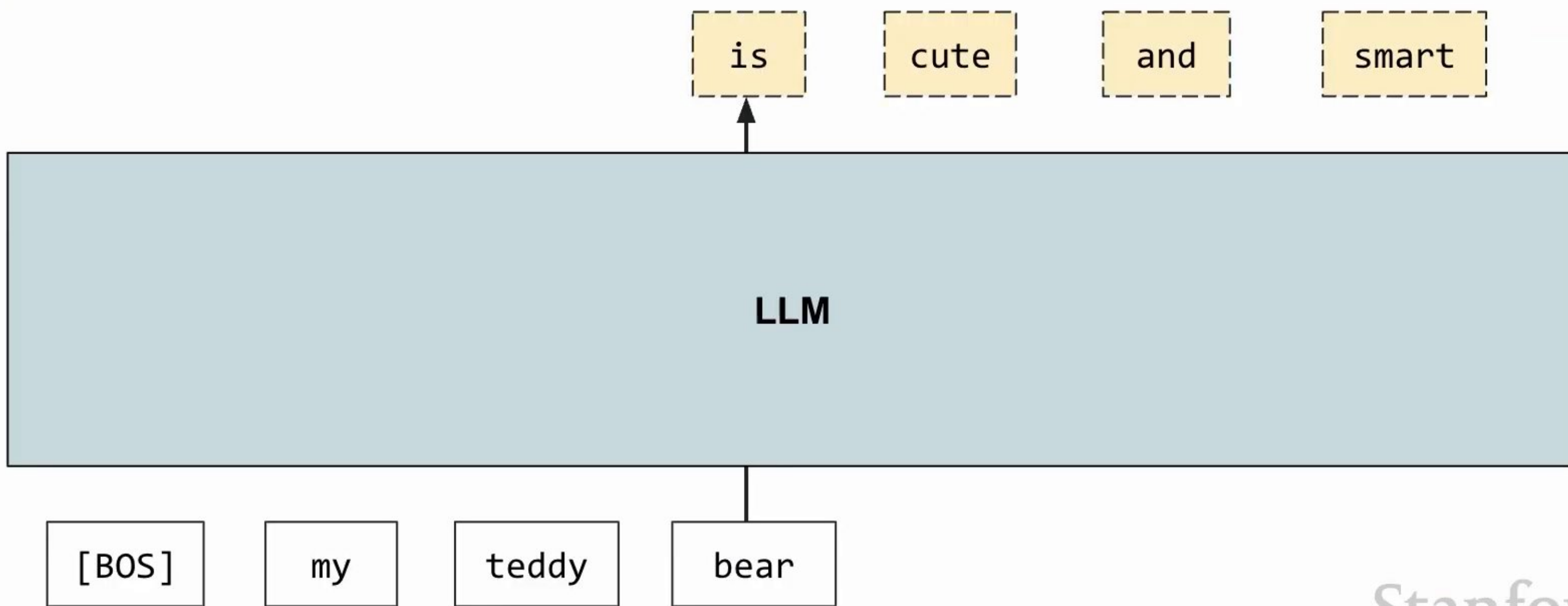
Speeding up decoding with speculative



Generate several tokens at once via MTP



MTP = **M**ulti-**T**oken **P**rediction

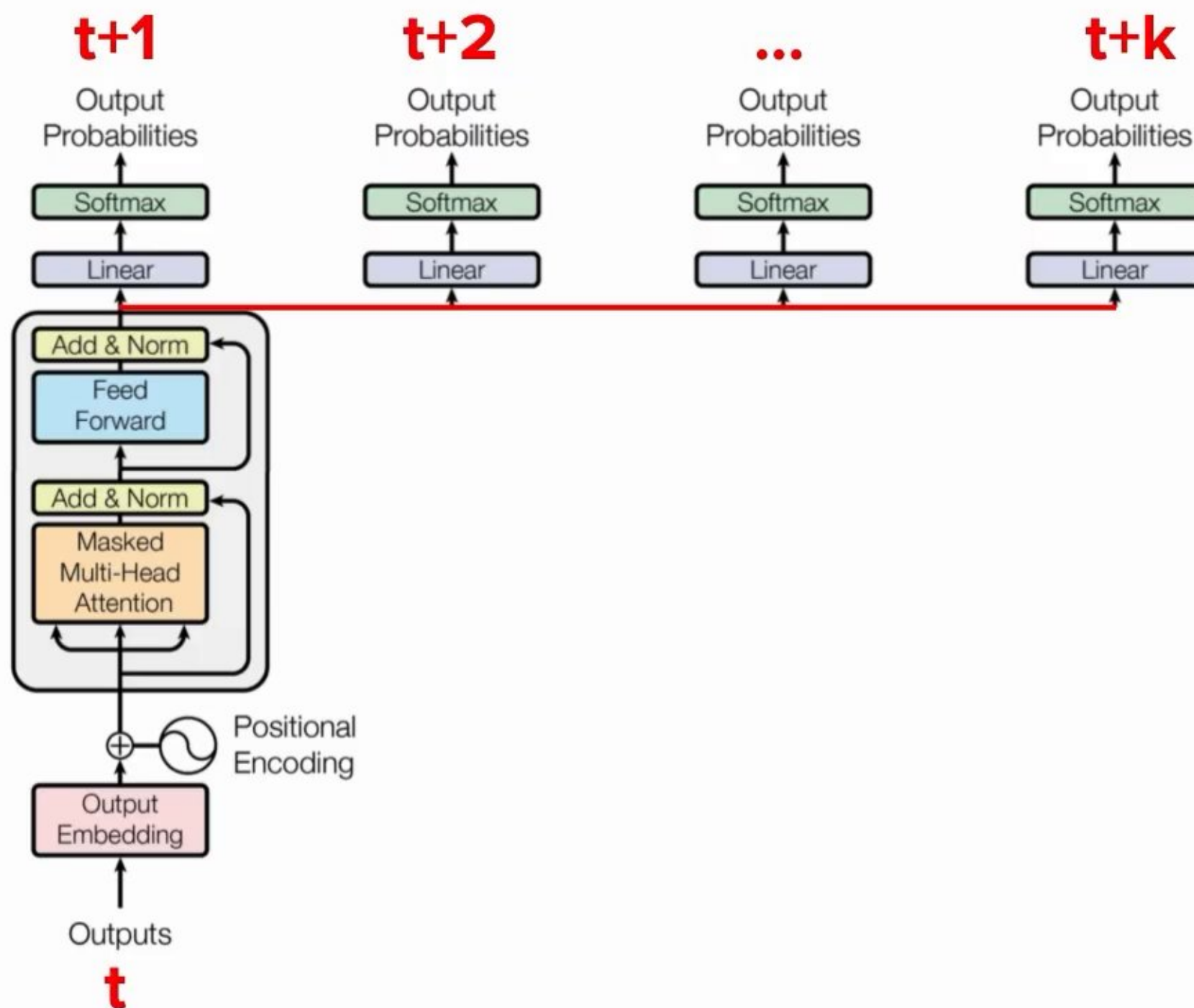


Stanford

Generate several tokens at once via MTP



Idea. Train k prediction heads: same draft and target models!



Stanford

Challenges



Categories. Many dimensions to optimize for.

"Exact" efficiency:

- Avoid redundancies
- Memory management
- Reformulate the math

Approximations:

- Architectural changes
- Embeddings representations
- Token prediction

Stanford

Challenges and some remedies

Categories. Many dimensions to optimize for.

"Exact" efficiency:

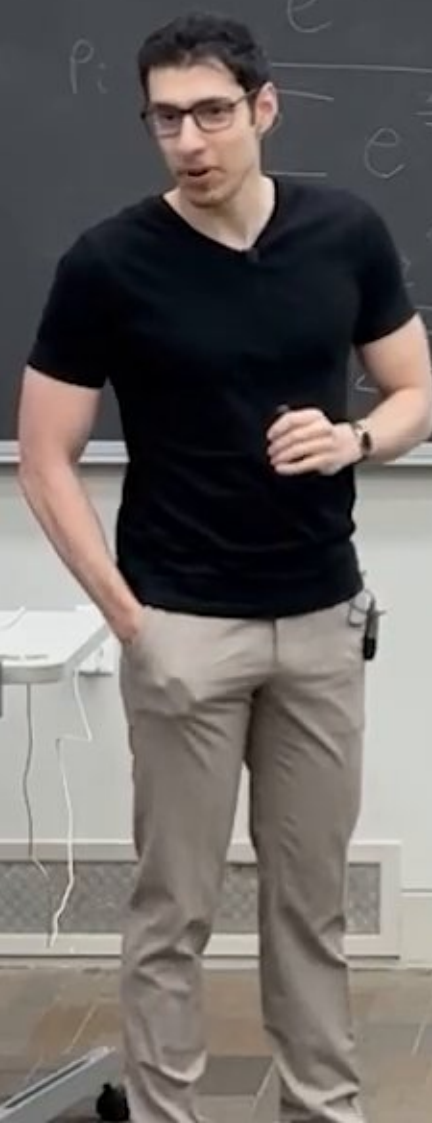
- Avoid redundancies ← KV cache
- Memory management ← PagedAttention
- Reformulate the math ← Speculative decoding

Approximations:

- Architectural changes ← Grouped query attention
- Embeddings representations ← Latent attention
- Token prediction ← Multi-token prediction

Stanford University

ICLW



Stanford | **ENGINEERING**