

# CME 295: Transformers & Large Language Models



Afshin Amidi & Shervine Amidi

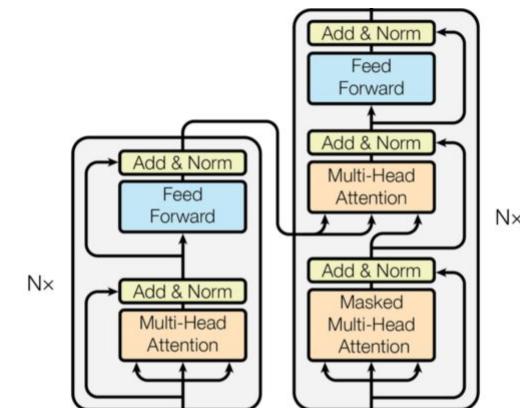


# Recap of last episode...

# Recap of last episode...

**3 categories** of Transformer-based models:

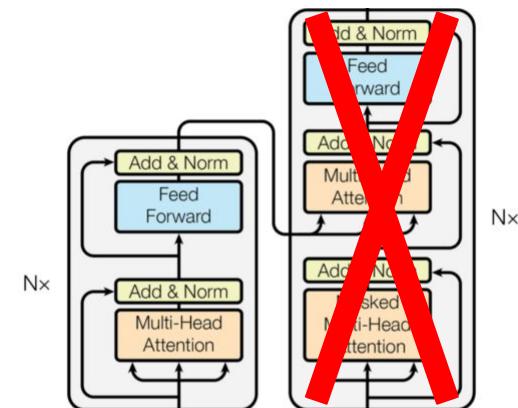
<b>Encoder-decoder</b>	Text to text	T5, mT5, ByT5
------------------------	--------------	---------------



# Recap of last episode...

**3 categories** of Transformer-based models:

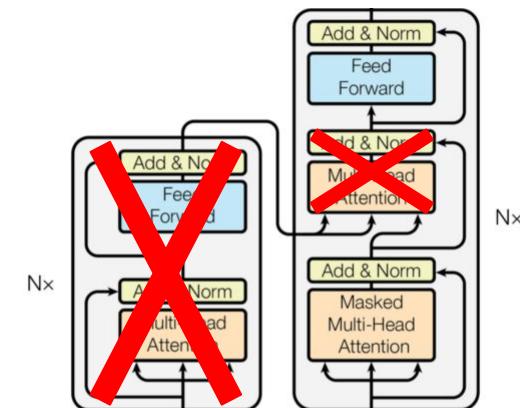
Encoder-decoder	Text to text	T5, mT5, ByT5
<b>Encoder-only</b>	Projection of embedding for class prediction (e.g. sentiment extraction)	BERT, DistilBERT, RoBERTa



# Recap of last episode...

**3 categories** of Transformer-based models:

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



# Recap of last episode...

**3 categories** of Transformer-based models:

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



# Transformers & Large Language Models

## LLM overview

MoE-based LLMs

Response generation

Prompting strategies

Inference optimizations

# Terminology

**LLM = Large Language Model**

it needs at least 100B  
of tokens for training  
dataset for benefit's  
property.

this kind of model are  
preferably much larger parameters  
an training architecture and require  
huge Computer of GPUs.

(but having optimizat<sup>n</sup> also to  
run) inference in small computing device)

assigns ↓ probability to sequence of tokens, it @  
always tries to predict the probability.

# Terminology

LLM = Large Language Model

# Terminology

LLM = Large Language Model

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

# Terminology

LLM = **Large** Language Model

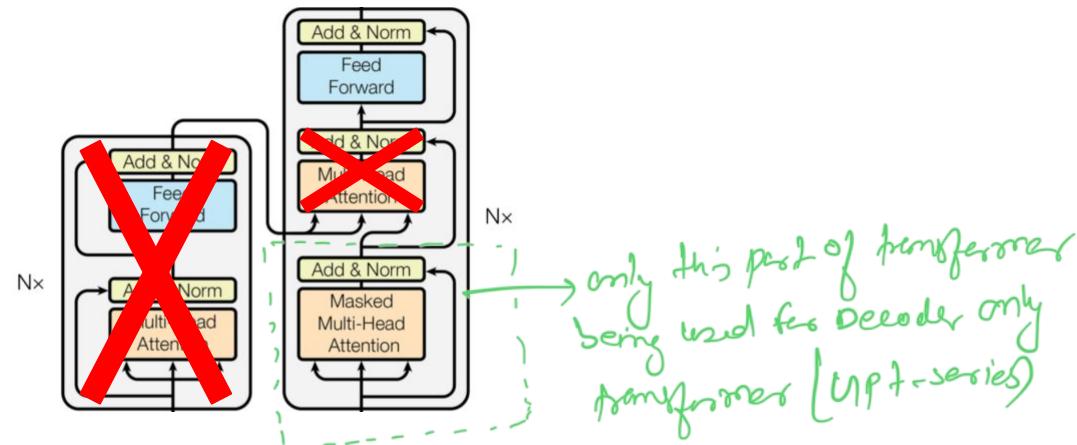
# Terminology

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

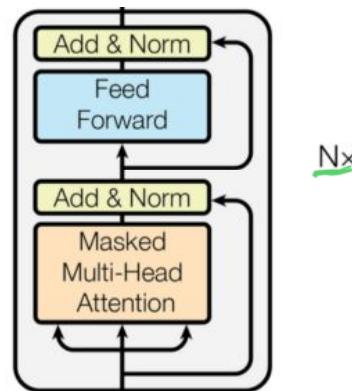
# Characteristics

## Decoder-only Transformer-based model



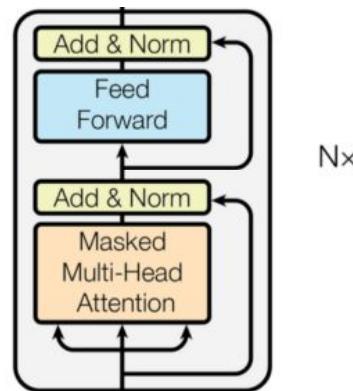
# Characteristics

**Decoder-only Transformer-based model**



# Characteristics

**Decoder-only Transformer-based model**



Nx

**Examples:** GPT series, LLaMA, Gemma, DeepSeek, Mistral, Qwen, ...



# Transformers & Large Language Models

LLM overview

MoE-based LLMs

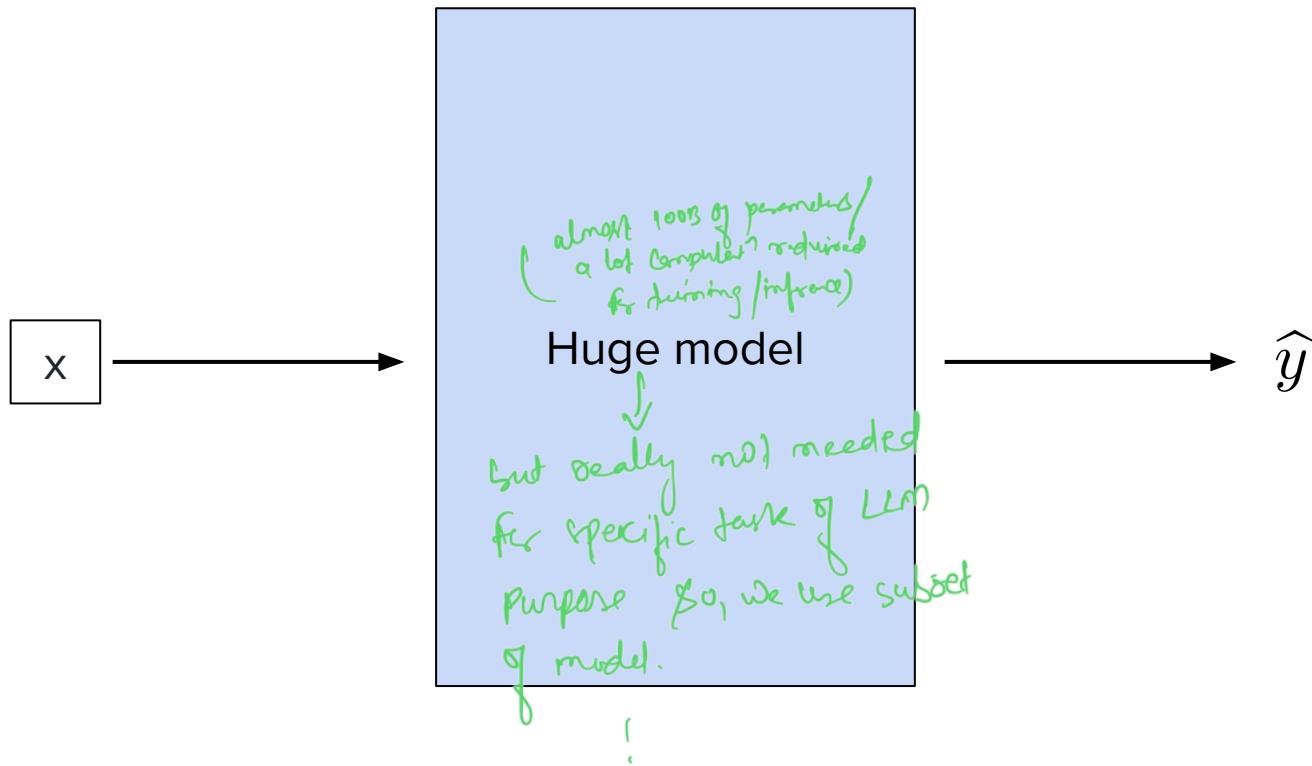
*mixture of Experts  
(adding Subsets of  
huge model for  
LLM purpose)*

Response generation

Prompting strategies

Inference optimizations

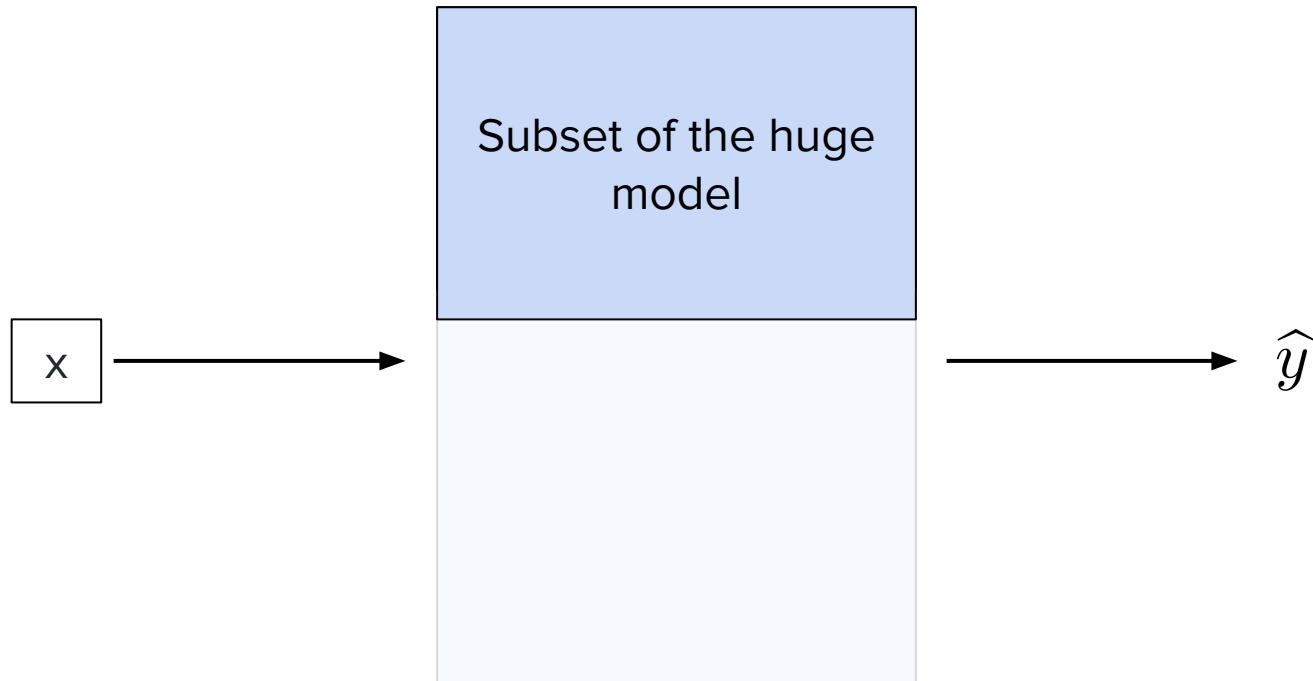
# Motivation



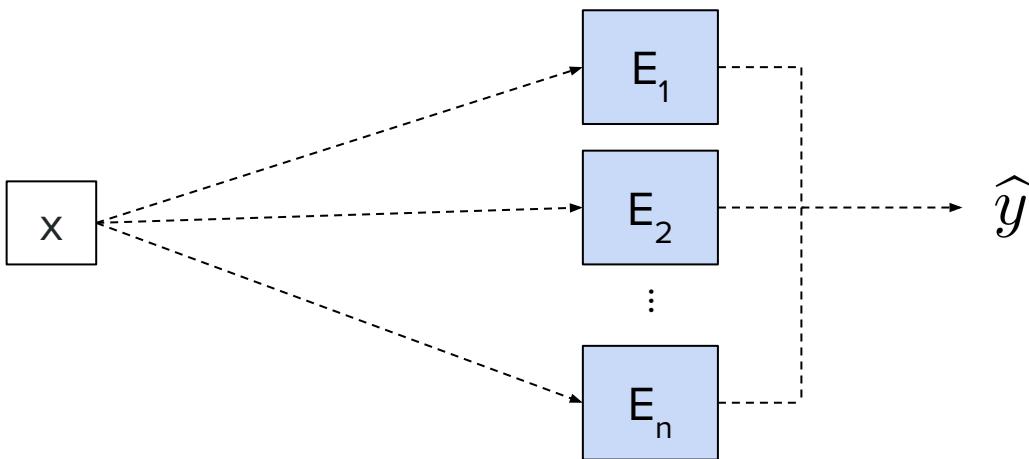
# Motivation

↓  
↓

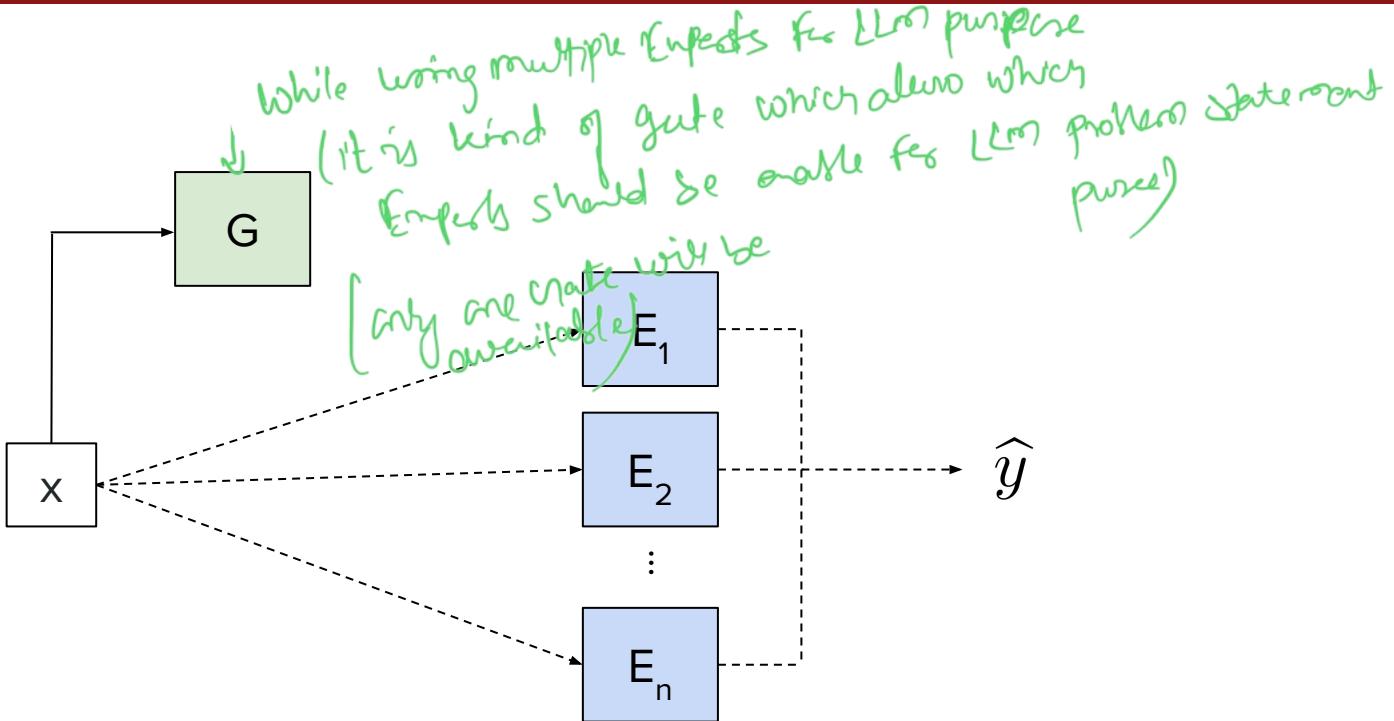
**Idea.** Not all weights are useful in the forward pass



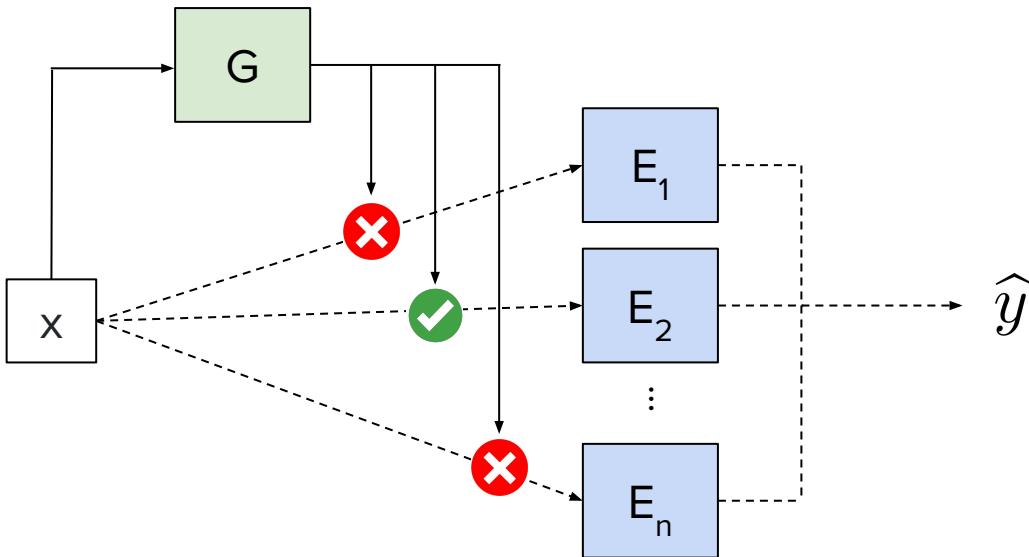
# Motivation



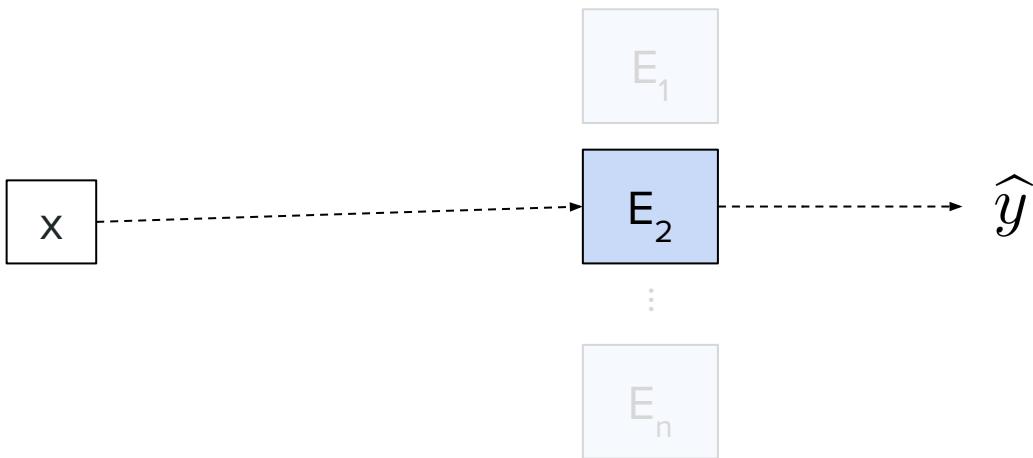
# Motivation



# Motivation

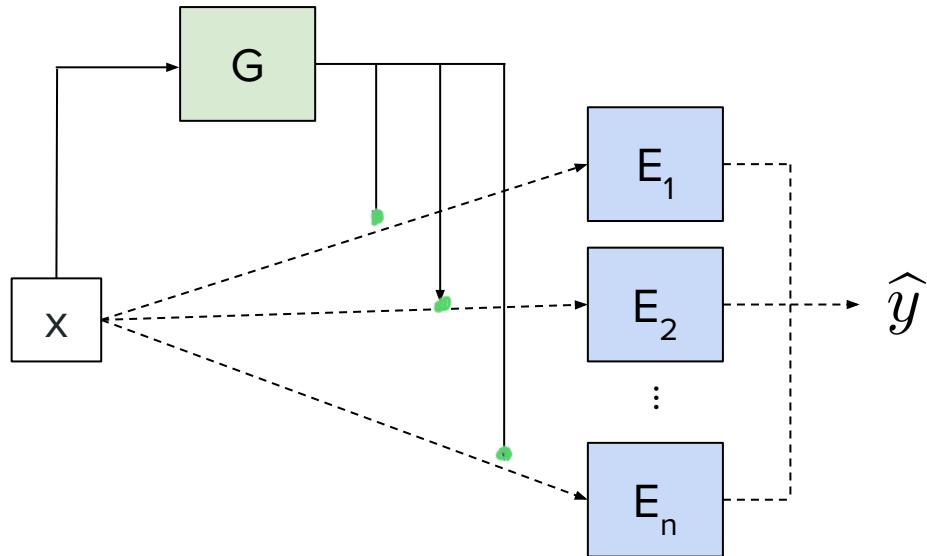


# Motivation



# Overview of MoEs

**MoE = Mixture of Experts**



Weighted sum of Experts as per  
requirements for purpose

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

Experts trained for ( $\alpha$ )

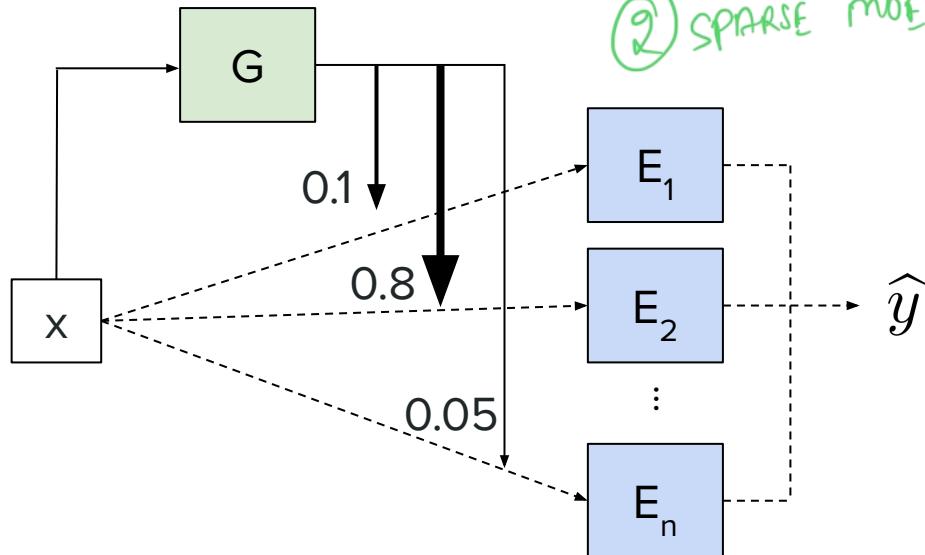
Weight assign  
based on priority  
of Experts based on LLm rewards

# Overview of MoEs

**MoE = Mixture of Experts**

↳ two kinds of moE

- ① DENSE MoEs
- ② SPARSE MoEs

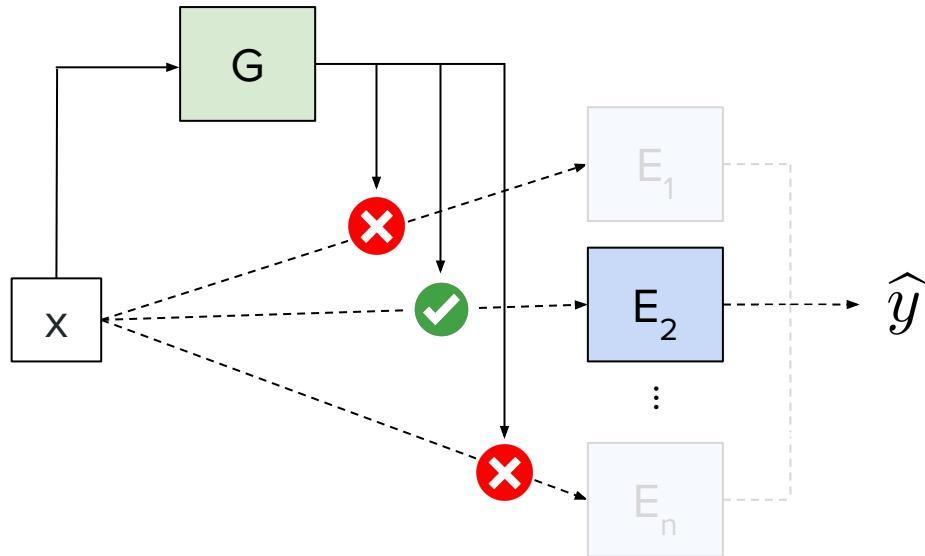


**Dense MoE.** Output is weighted average of all expert outputs.

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

# Overview of MoEs

**MoE = Mixture of Experts**

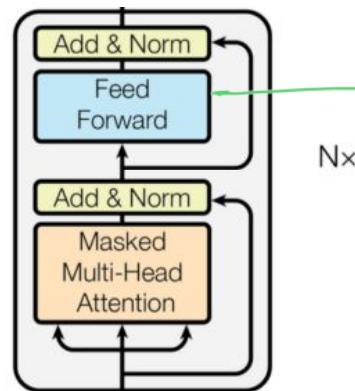


**Sparse MoE.** Output is weighted  
average of **selected** expert outputs.

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

Via **top-k** selection

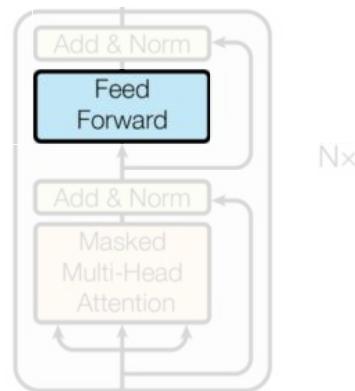
# MoE in Transformer-based models



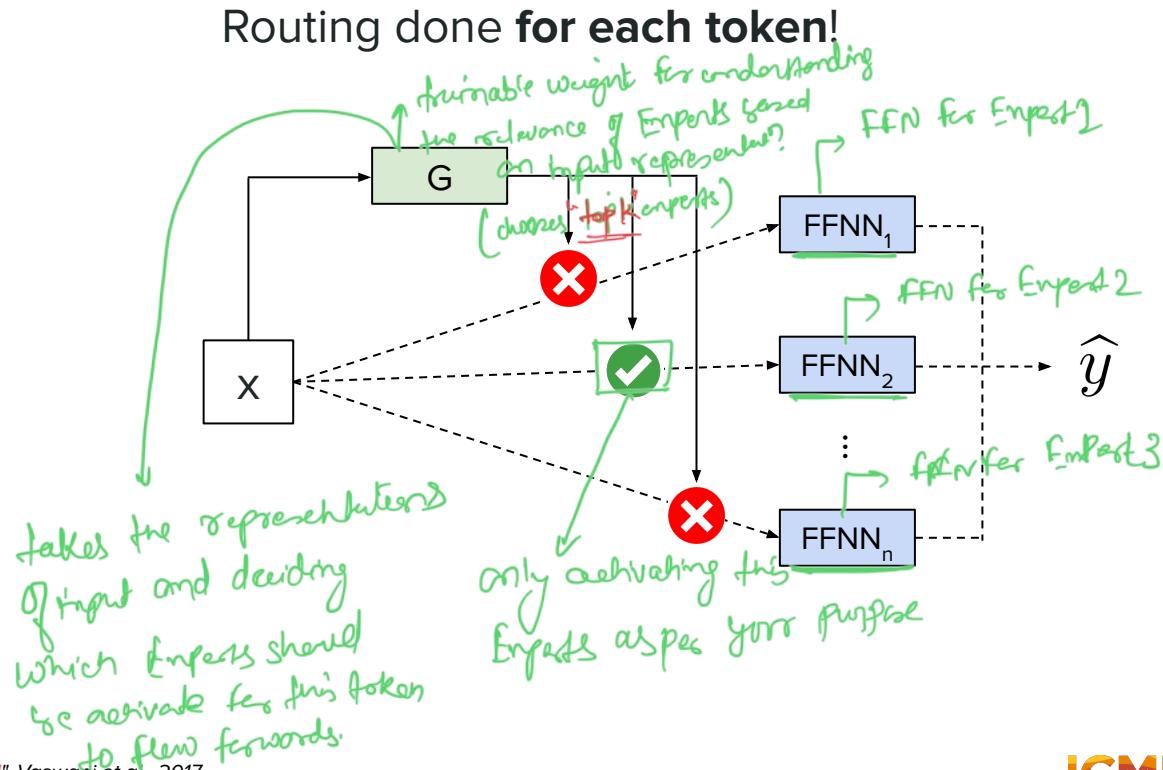
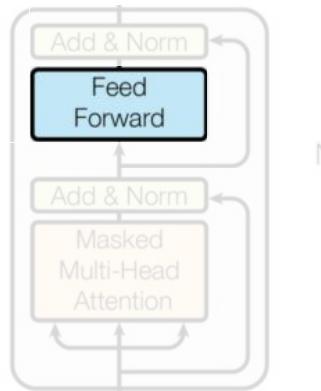
permissible for combining multiple Experts

Nx

# MoE in Transformer-based models



# MoE in Transformer-based models



# Training challenges include routing collapse

**Symptom.** Same expert gets selected most of the time.

"routing collapse"

// Problem of mixture of Experts.  
Each Expert will have weight  
but  $\approx 1/2$  only Experts always  
activates) So, if trying using  
changing the loss function  
and adding Entropy term ( $\alpha$ )

# 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:

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

Converging through uniform distribution across Experts.

$f_i$  = fraction of tokens routed to expert i  
 $P_i$  = average routing probability for expert i

# Interpreting experts

Layer 0

"Mixtral of Experts"

```
class MoeLayer(nn.Module):
    def __init__(self, experts: List[nn.Module], gate: nn.Module, moe_args: Dict):
        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.shape[-1])
        gate_logits = self.gate(inputs_squashed)
        weights, selected_experts = torch.topk(
            gate_logits, self.args.numExpertsPerGroup)
        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

(it defines which Expert is expertise on what kind of input representation)



# Transformers & Large Language Models

LLM overview

MoE-based LLMs

**Response generation**

Prompting strategies

Inference optimizations

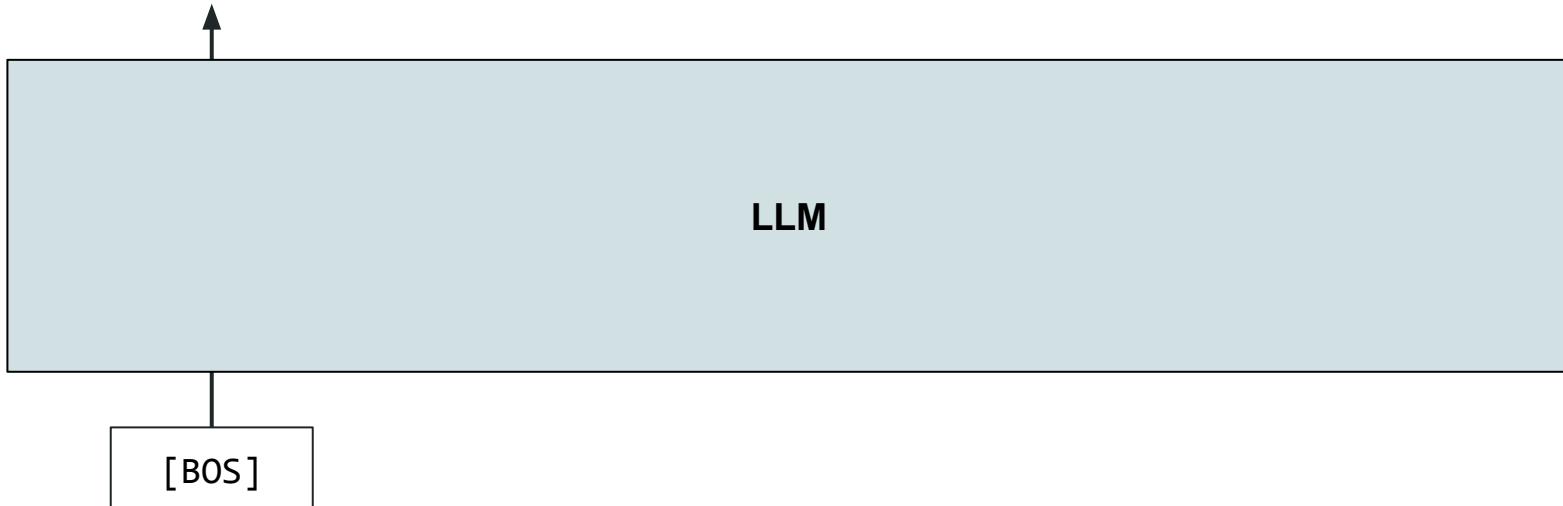
Based on LLM/Ensemble  
how response is  
generated?

# Next token prediction

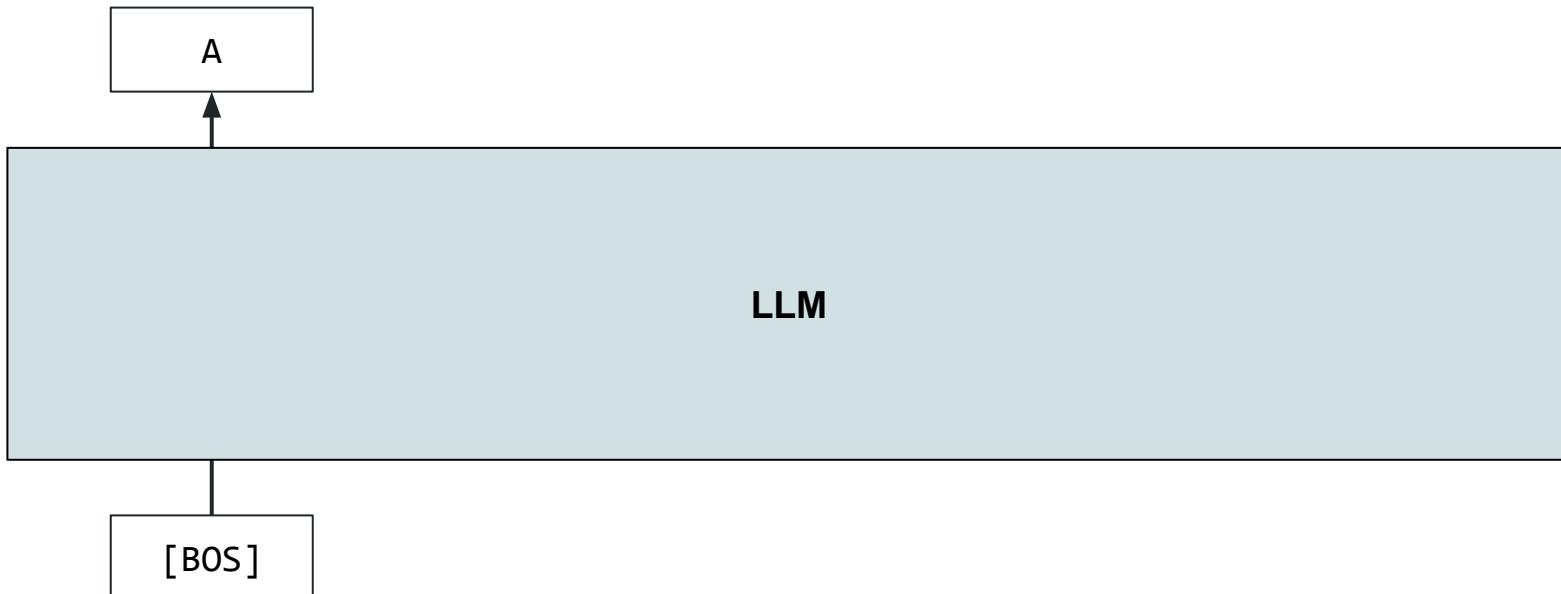
**LLM**

[BOS]

# Next token prediction



# Next token prediction



# Next token prediction

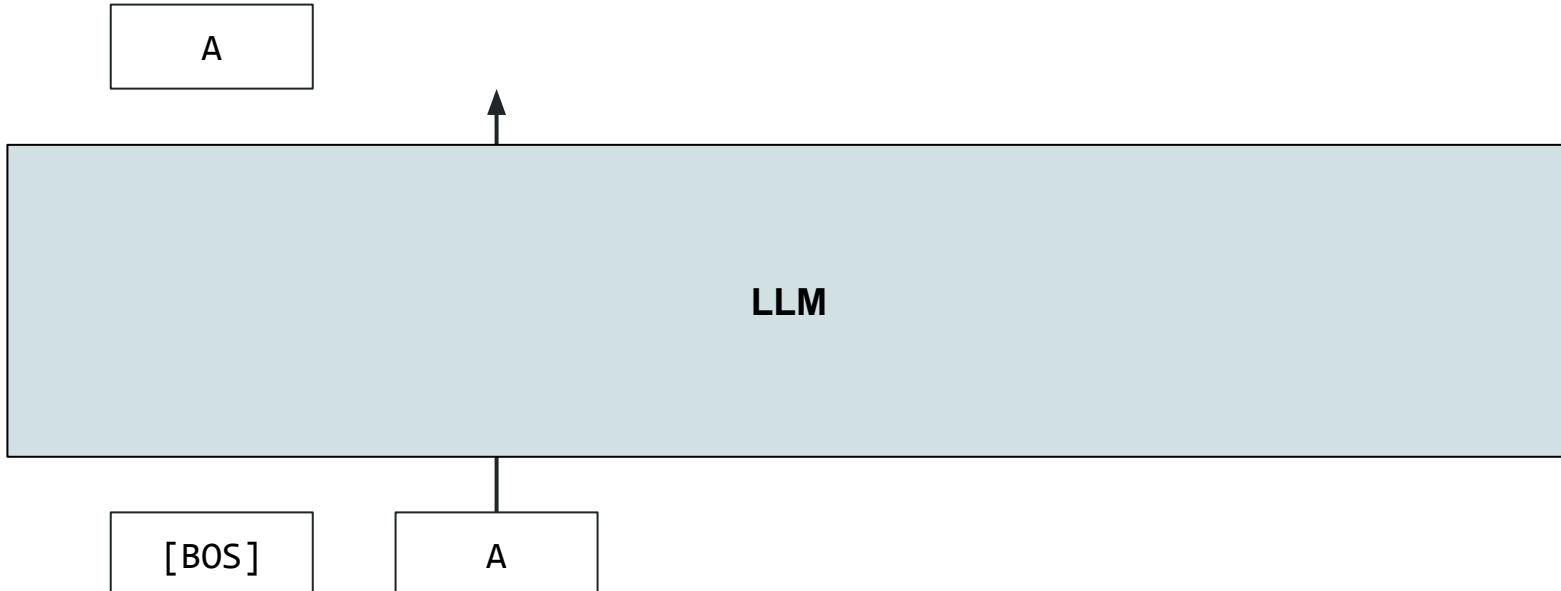
A

**LLM**

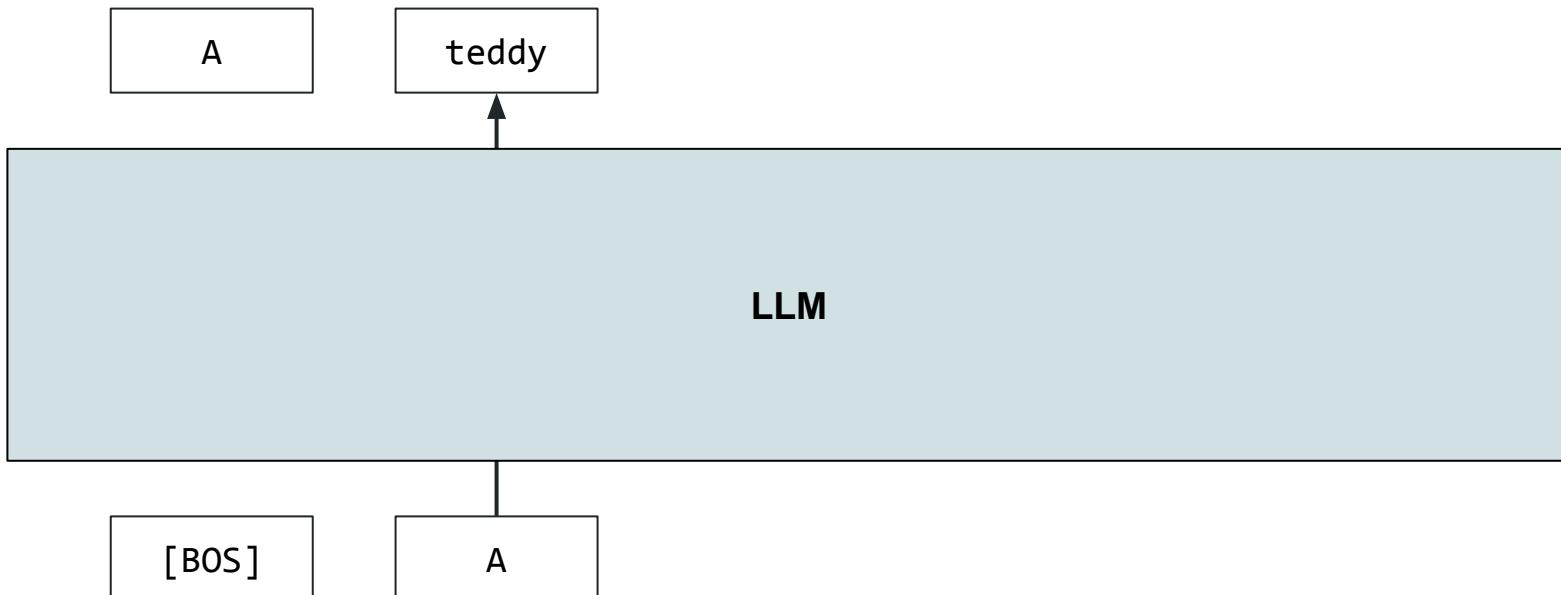
[BOS]

A

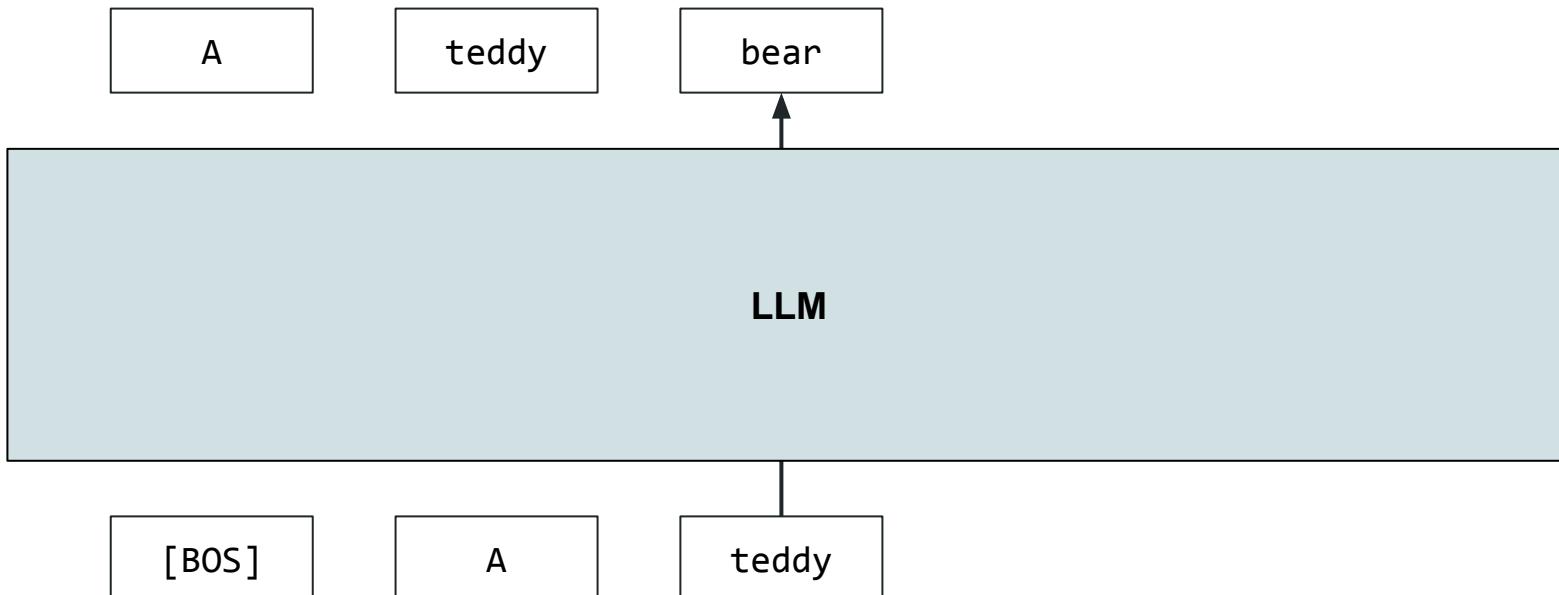
# Next token prediction



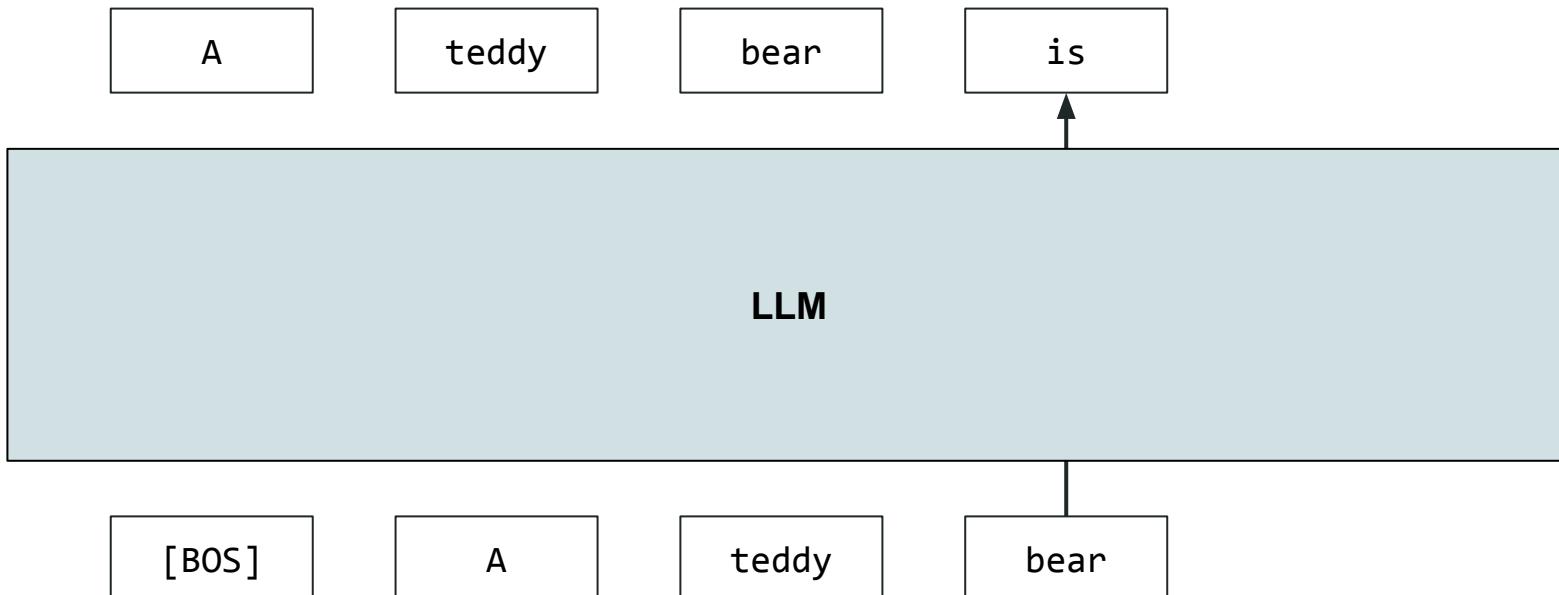
# Next token prediction



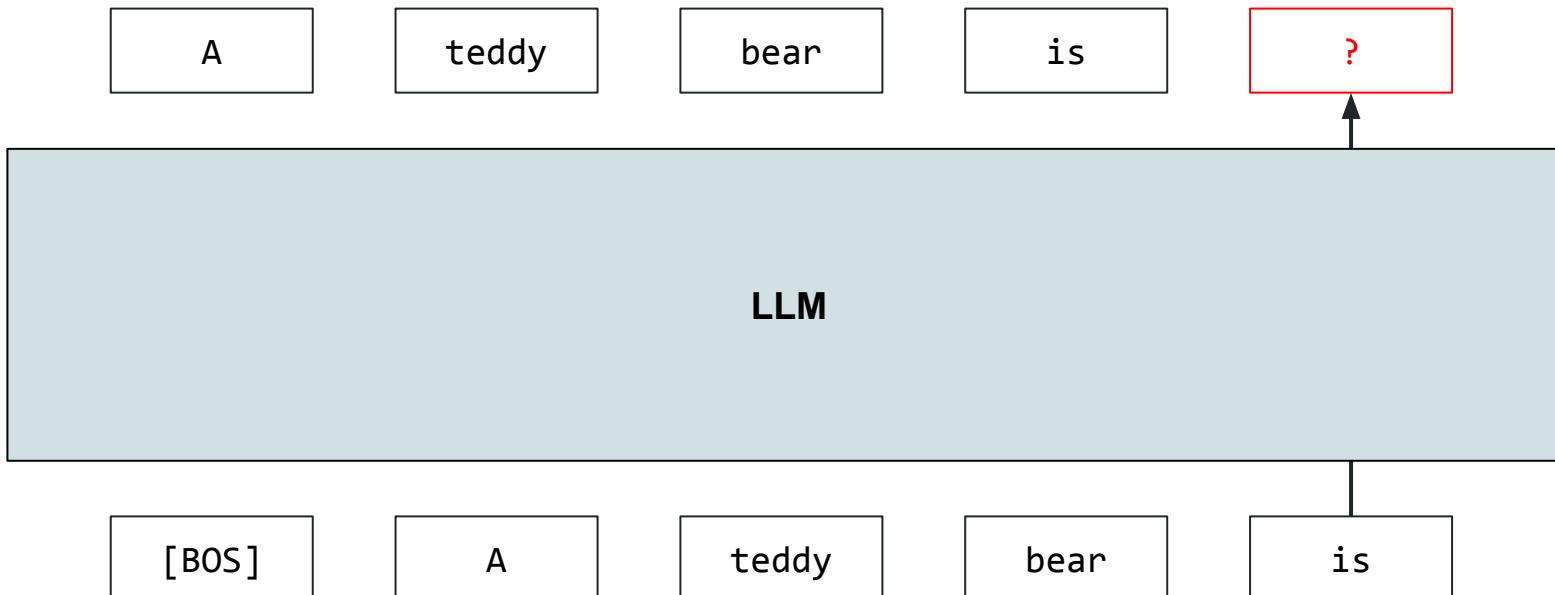
# Next token prediction



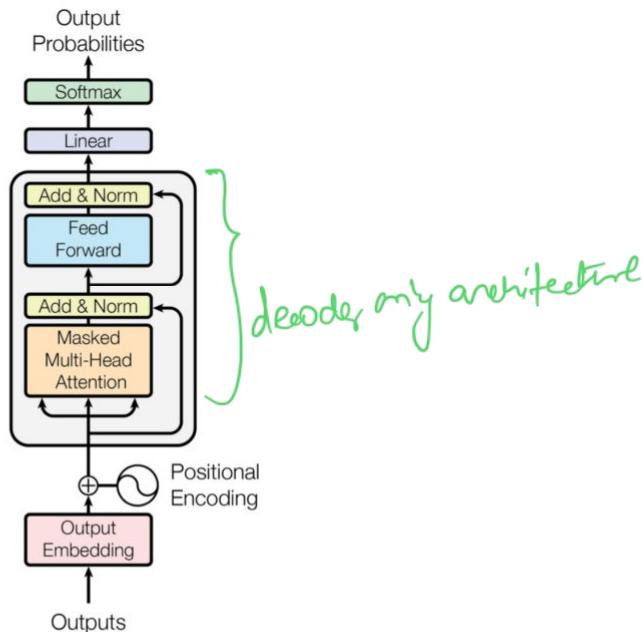
# Next token prediction



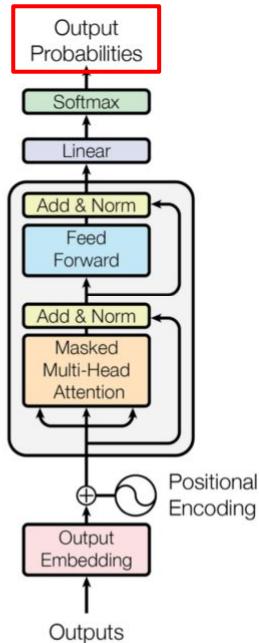
# Next token prediction



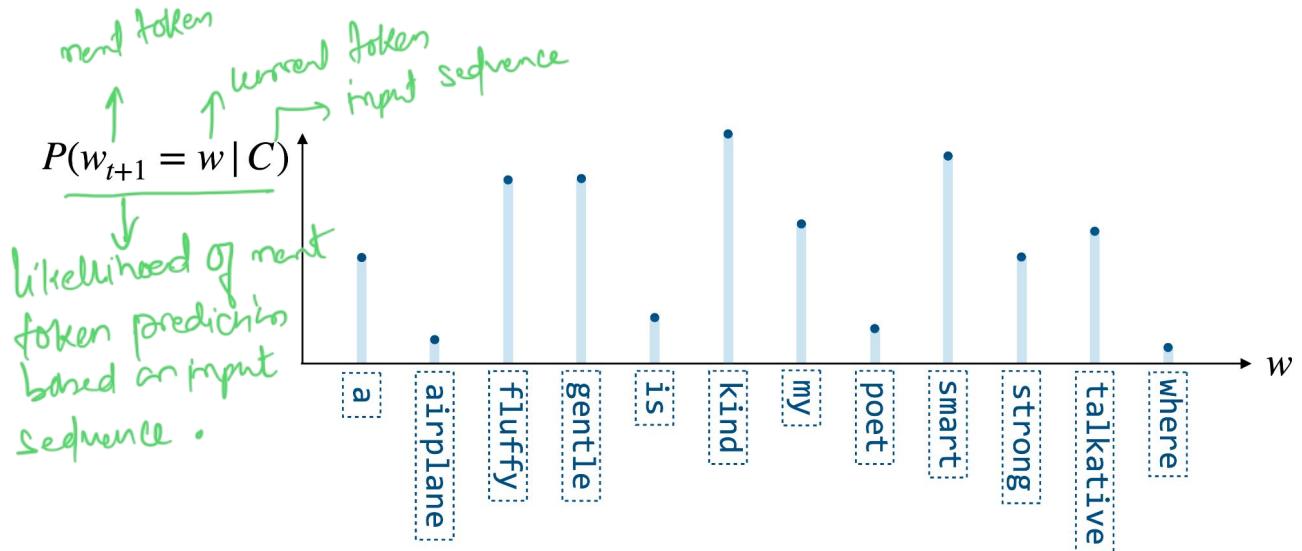
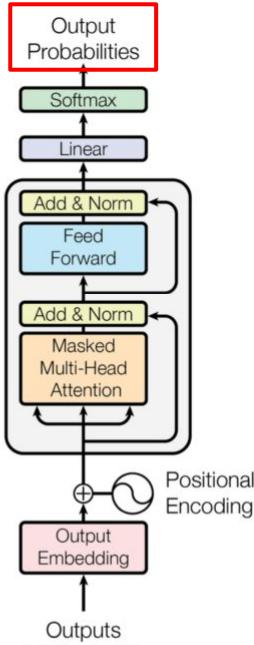
# Predicting next token



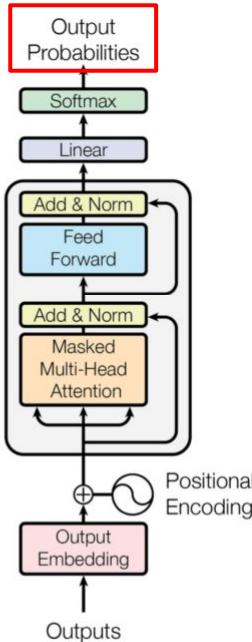
# Predicting next token



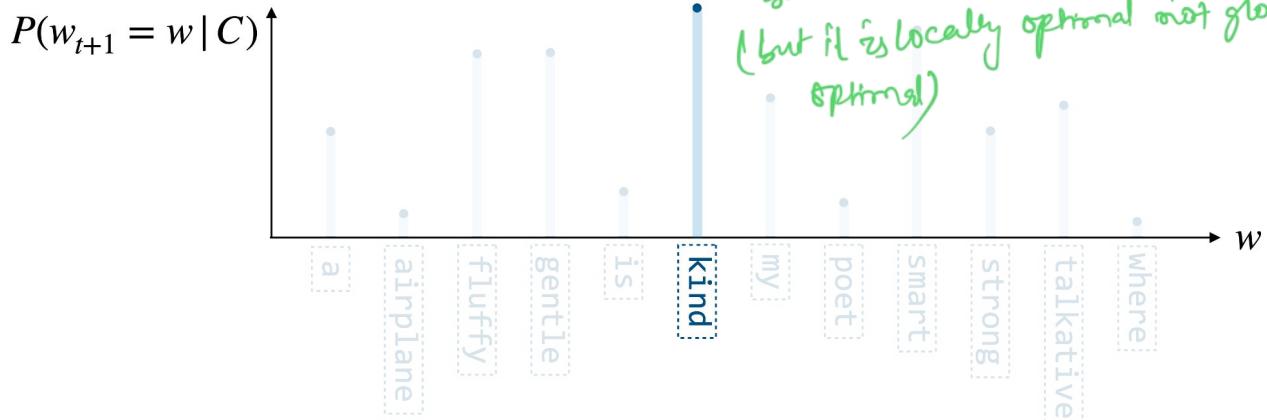
# Predicting next token



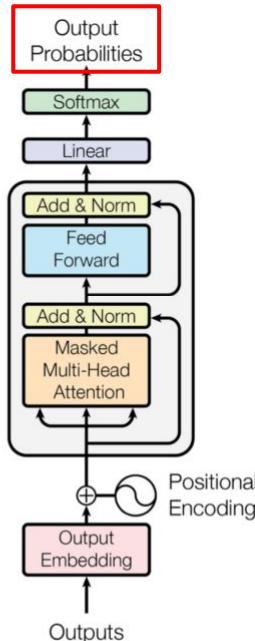
# Predicting next token with greedy decoding



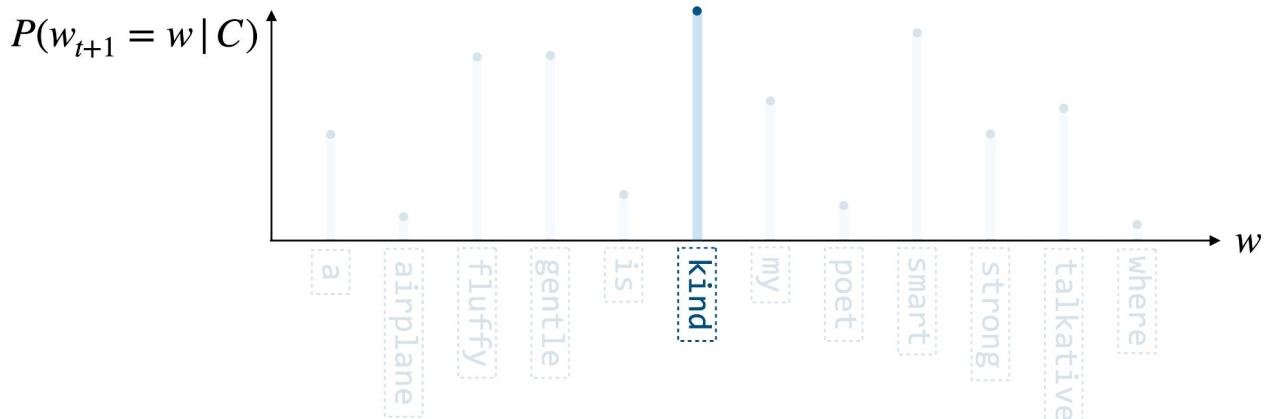
1<sup>st</sup> idea. Take token with highest predicted probability



# Predicting next token with greedy decoding



1<sup>st</sup> idea. Take token with highest predicted probability

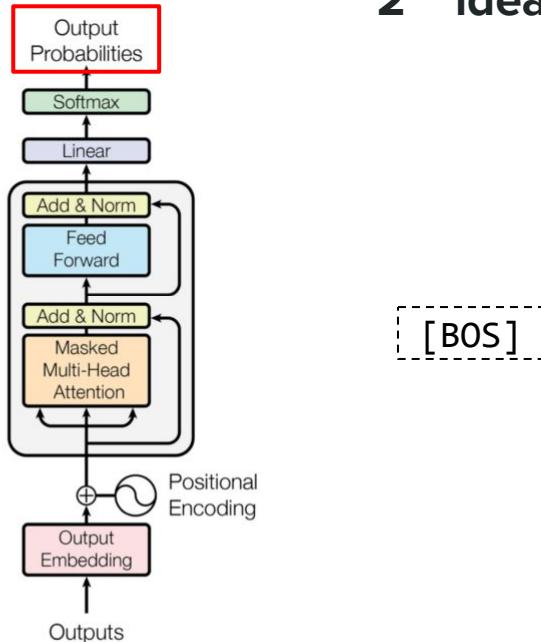


Limitations. Output not optimal, natural and/or diverse

# Predicting next token with beam search

**2<sup>nd</sup> idea.** Keep k paths that are the most likely

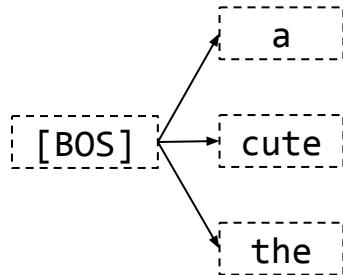
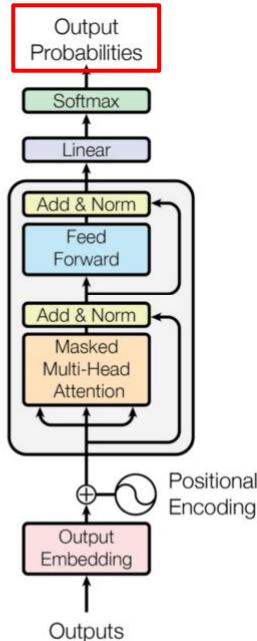
↓ it is globally optimal for choosing the next token based on keeping k paths that are most likely based on input sequence



[ [BOS] ]

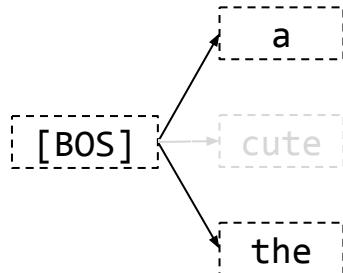
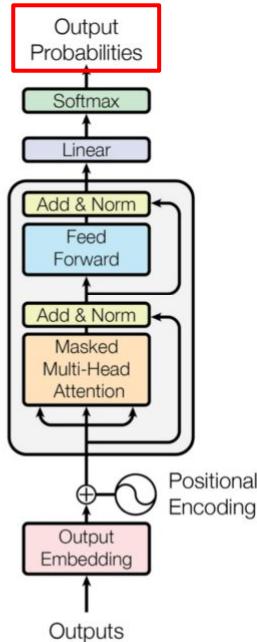
# Predicting next token with beam search

**2<sup>nd</sup> idea.** Keep k paths that are the most likely

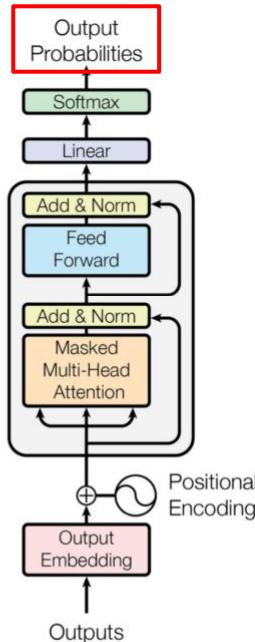


# Predicting next token with beam search

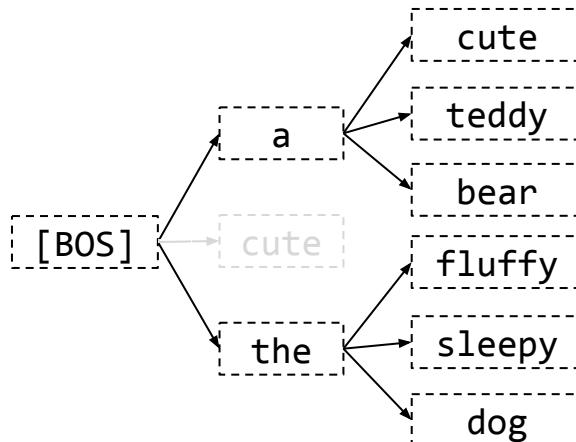
**2<sup>nd</sup> idea.** Keep k paths that are the most likely



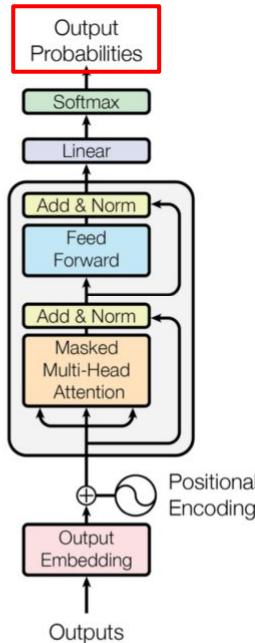
# Predicting next token with beam search



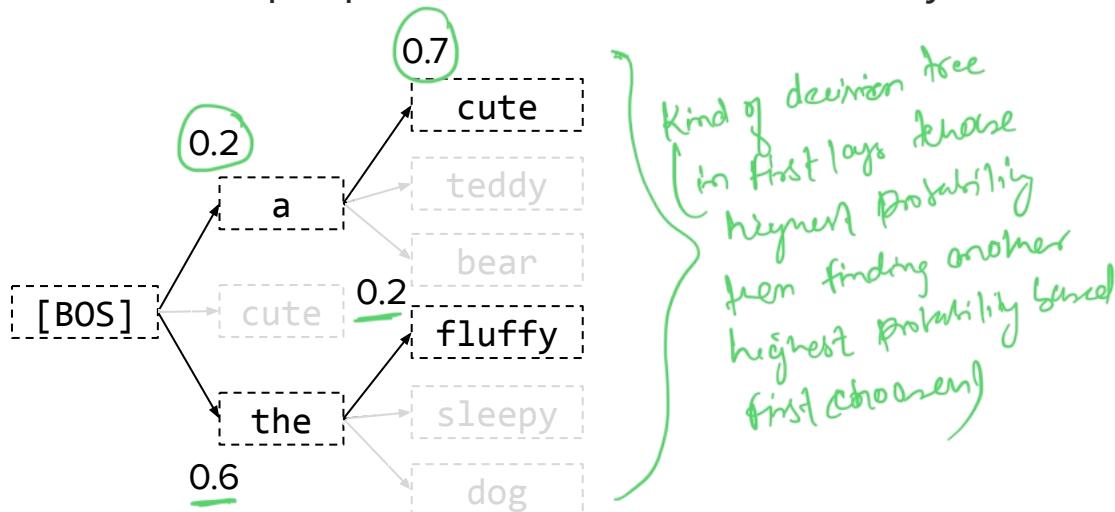
**2<sup>nd</sup> idea.** Keep k paths that are the most likely



# Predicting next token with beam search

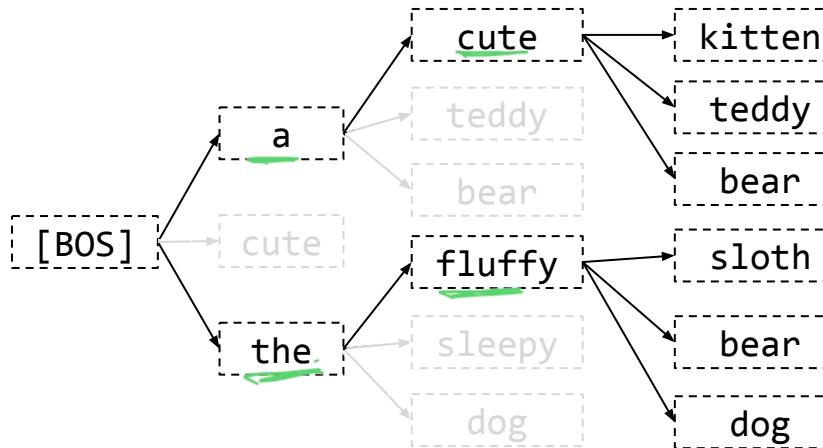
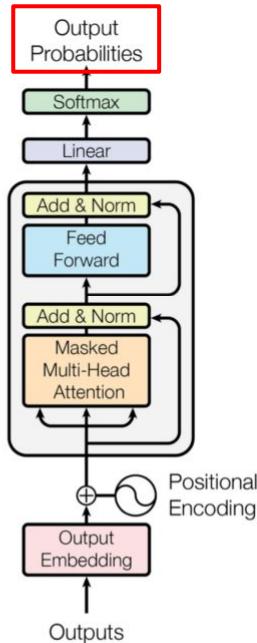


**2<sup>nd</sup> idea.** Keep k paths that are the most likely



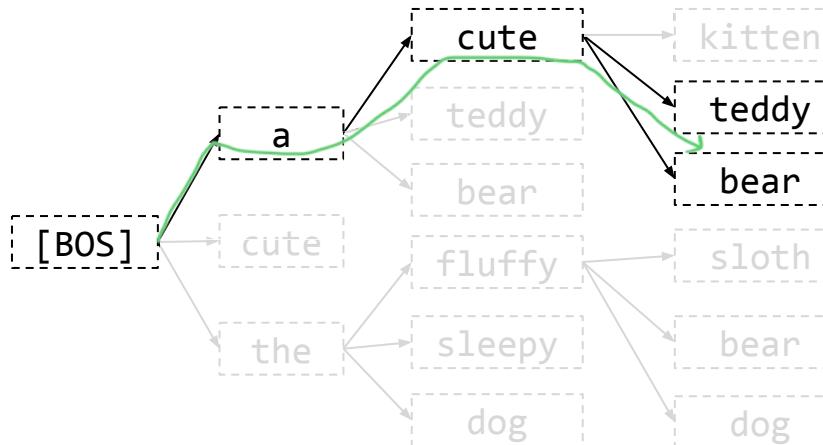
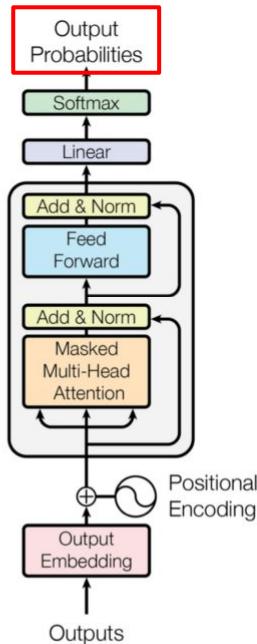
# Predicting next token with beam search

**2<sup>nd</sup> idea.** Keep k paths that are the most likely

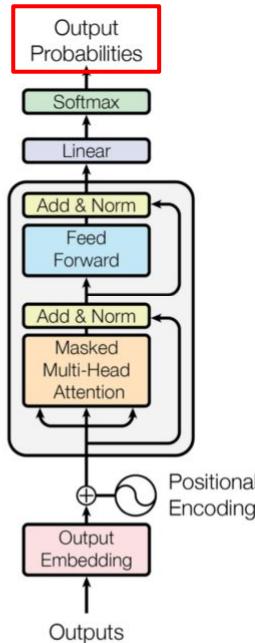


# Predicting next token with beam search

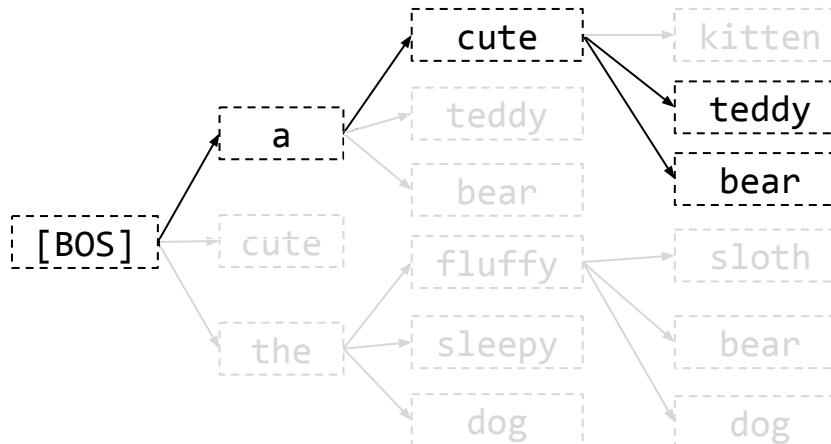
**2<sup>nd</sup> idea.** Keep k paths that are the most likely



# Predicting next token with beam search

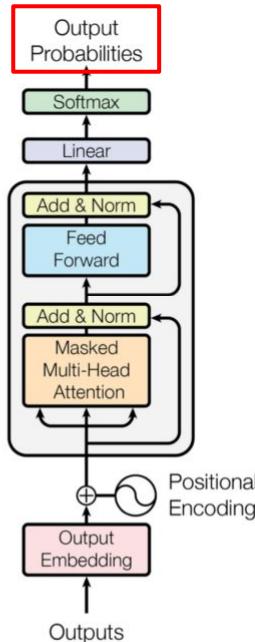


**2<sup>nd</sup> idea.** Keep k paths that are the most likely

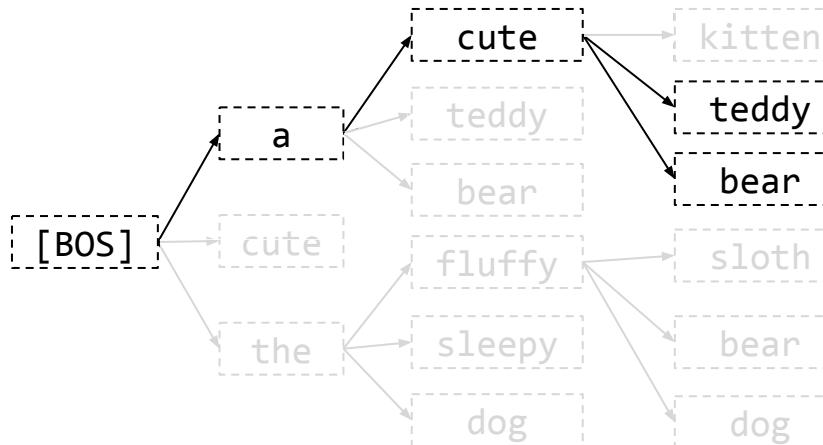


Ends at [EOS]

# Predicting next token with beam search



**2<sup>nd</sup> idea.** Keep k paths that are the most likely



Ends at [ EOS ]

**Limitations.** Needs computations + lacks diversity/creativity

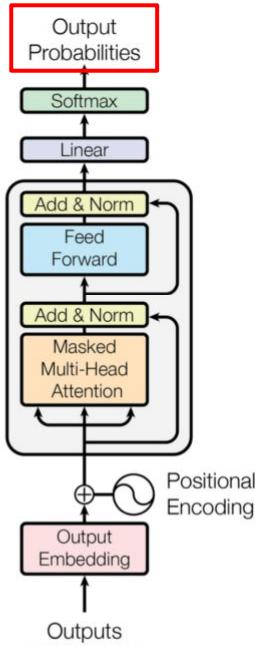
*↳ But for keeping track of m possible path requires huge computation/memory*

*[ Beam Search ]*

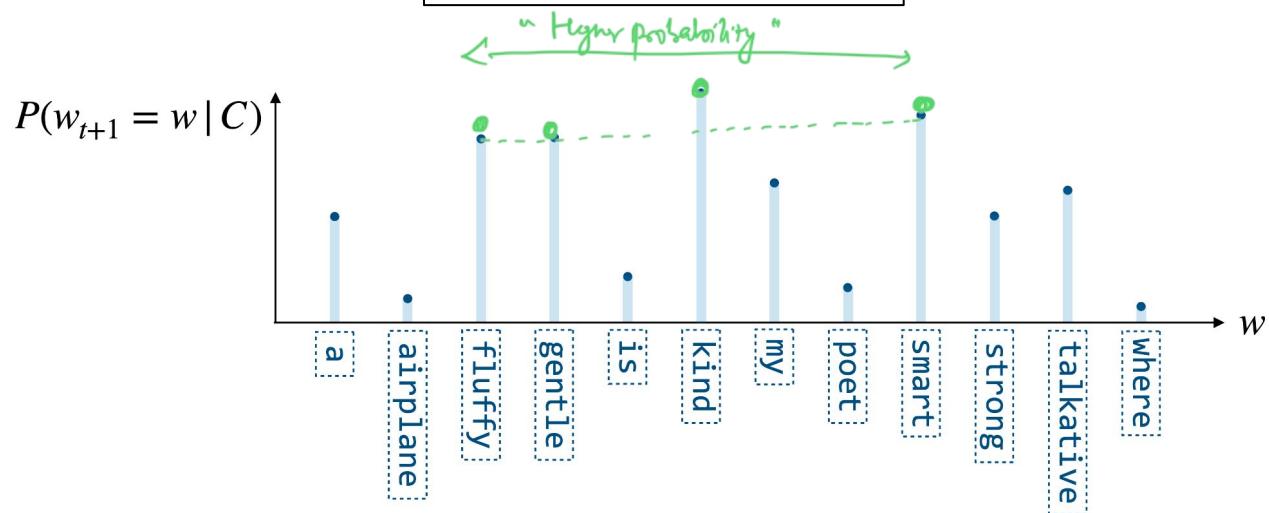
# Predicting next token with sampling

Sampling method.

3<sup>rd</sup> idea. Sample next token from probability distribution:

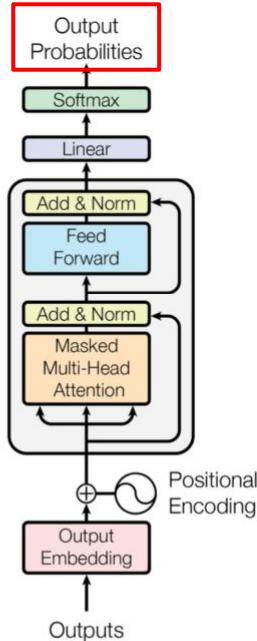


$$\hat{w}_{t+1} \sim P(w_{t+1} | C)$$

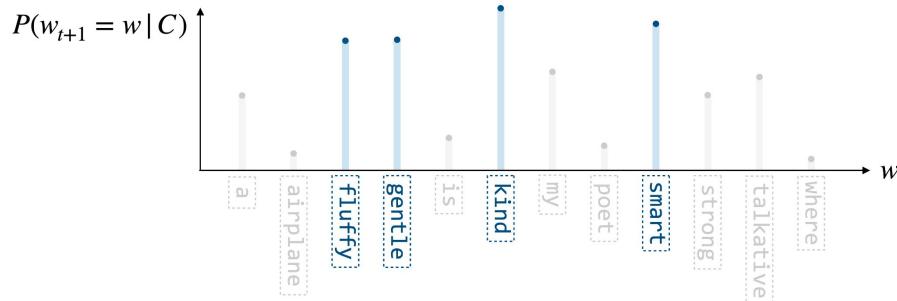


# Sampling strategies

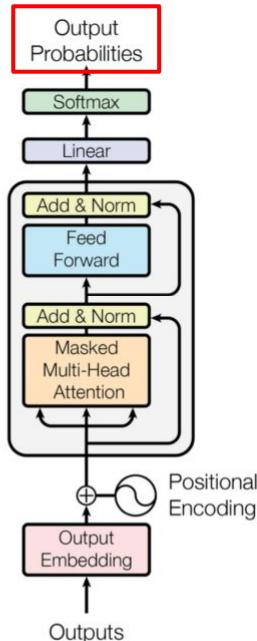
- **Top-k:** Sample among top k most probable tokens



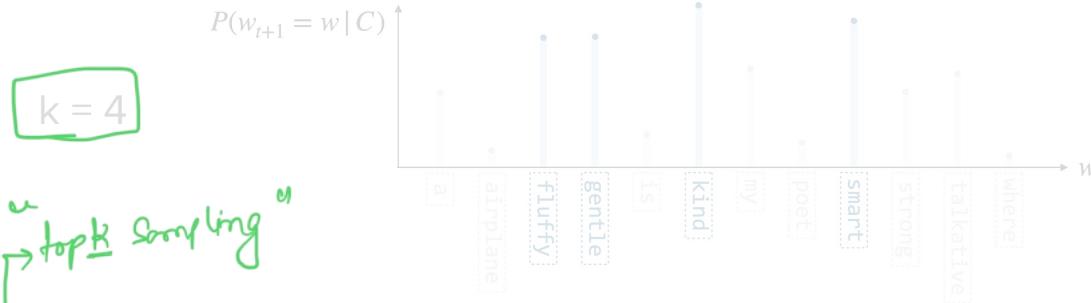
$k = 4$



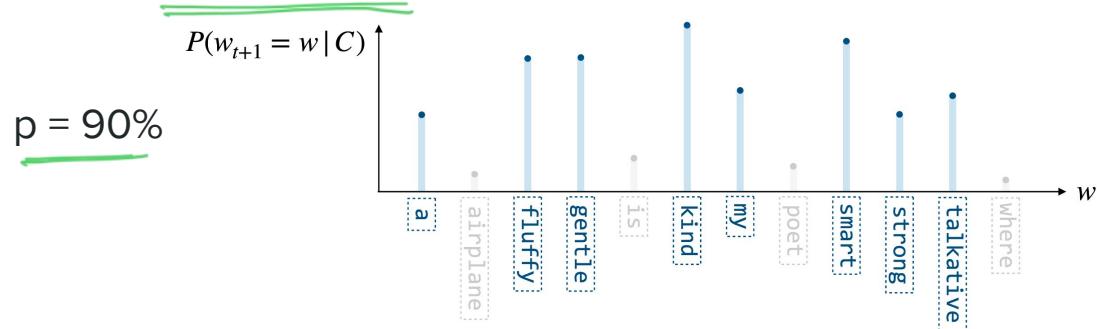
# Sampling strategies



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



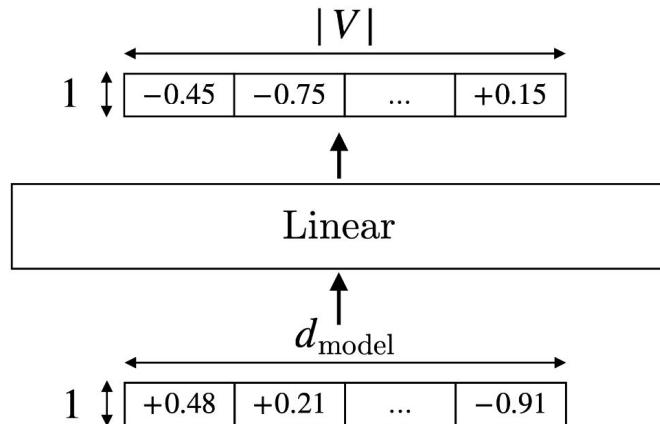
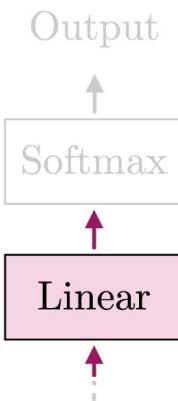
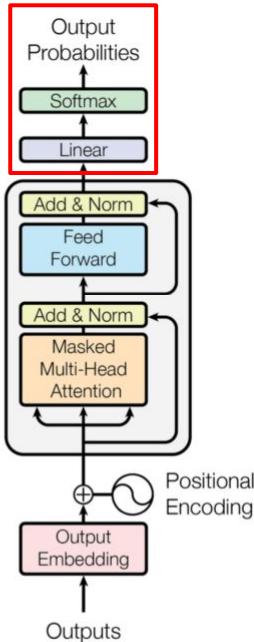
- **Top-p:** Random sample among smallest set of tokens with cumulative probability  $\geq p$



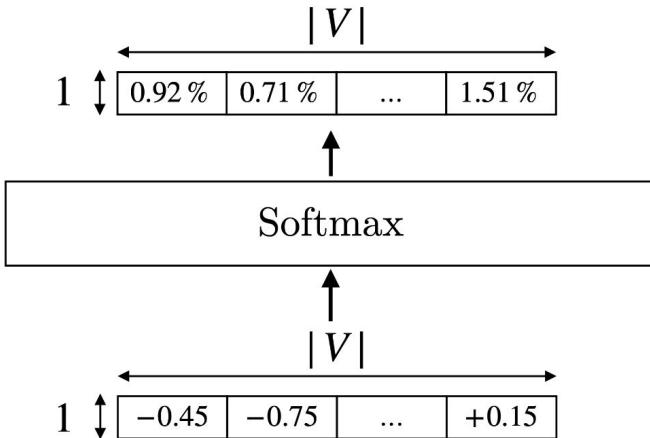
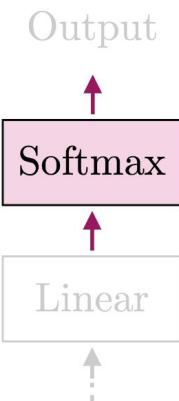
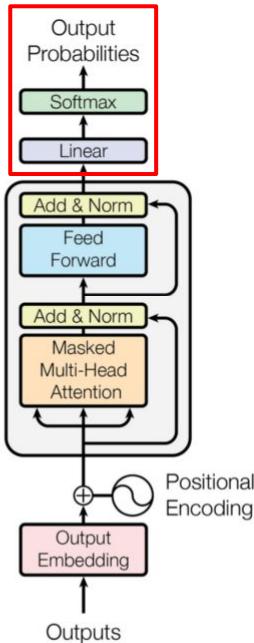
# Next token prediction

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

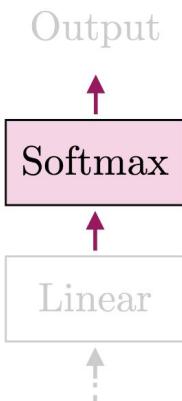
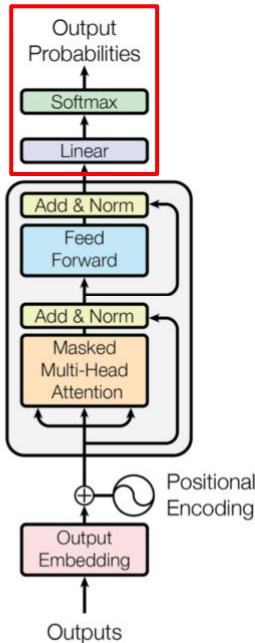
# Probability computation



# Probability computation



# Temperature allows to tweak output probabilities



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

Handwritten notes explaining the softmax distribution:

- temperature : ignorance
- small temperature : high ignorance
- high temperature : low ignorance

uniform distribution

# Impact of temperature on probabilities

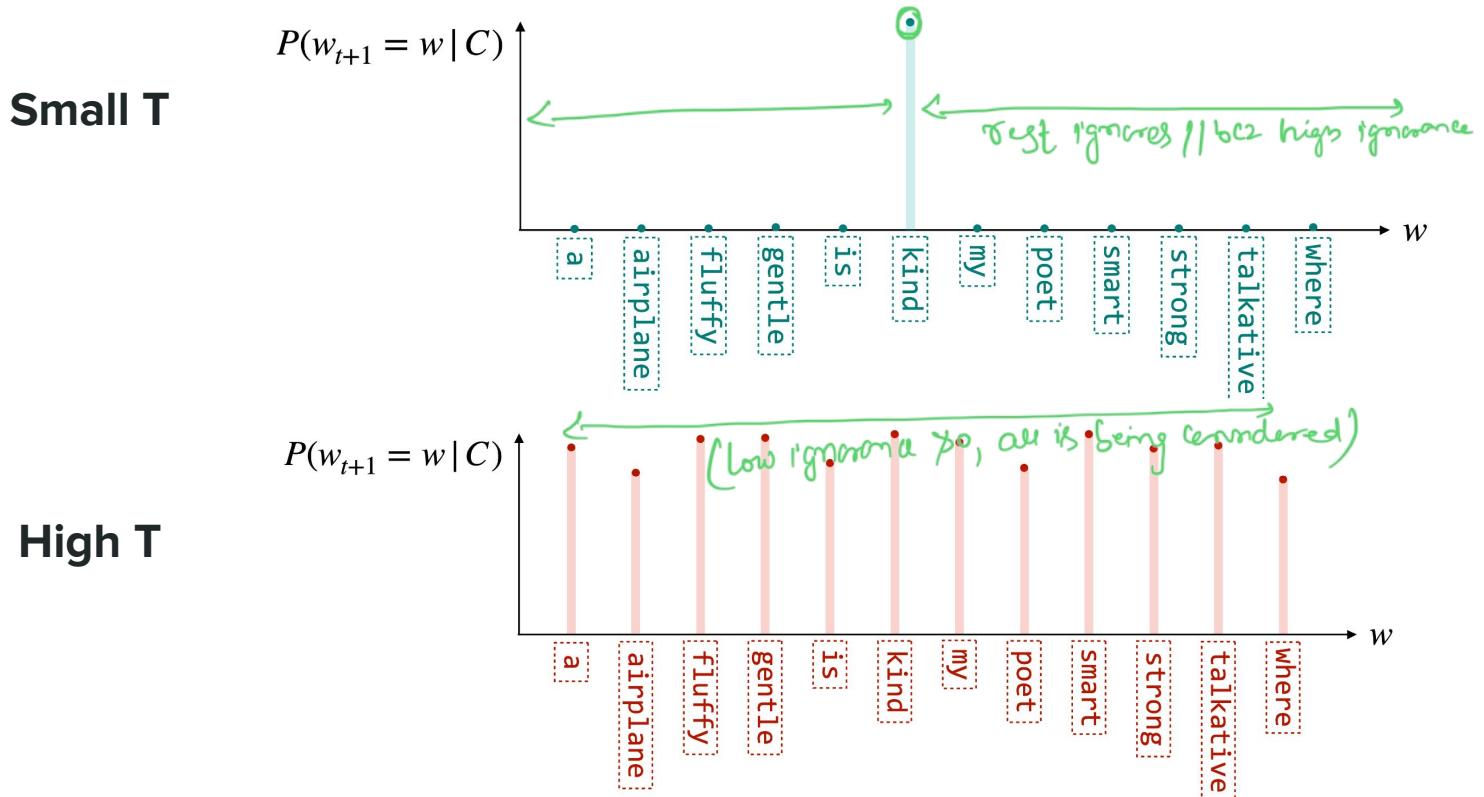
**Small T** → Encourage the next token could be highest probable

?

**High T** → Encourage the next token could be uniform distribution of most predictions of token / multiple possibility of probable token and more creating tokens.

?

# Impact of temperature on probabilities



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

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

# 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"  
}
```

 to generate output in specific format using guided decoding (which basically helps filtering out invalid tokens using guided decoding)

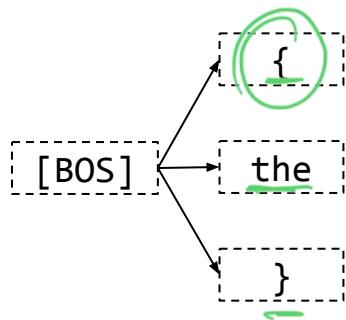
# Constraining the output via guided decoding

**Idea.** Only allow "valid" next tokens

[BOS]

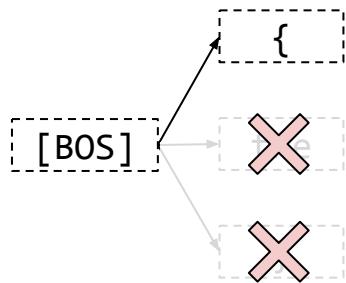
# Constraining the output via guided decoding

**Idea.** Only allow "valid" next tokens



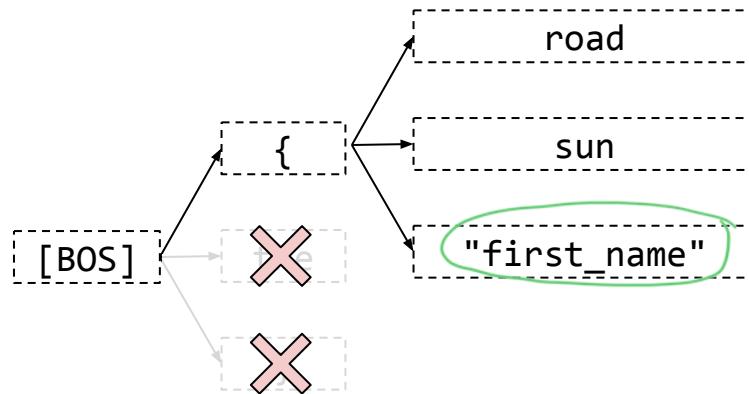
# Constraining the output via guided decoding

**Idea.** Only allow "valid" next tokens



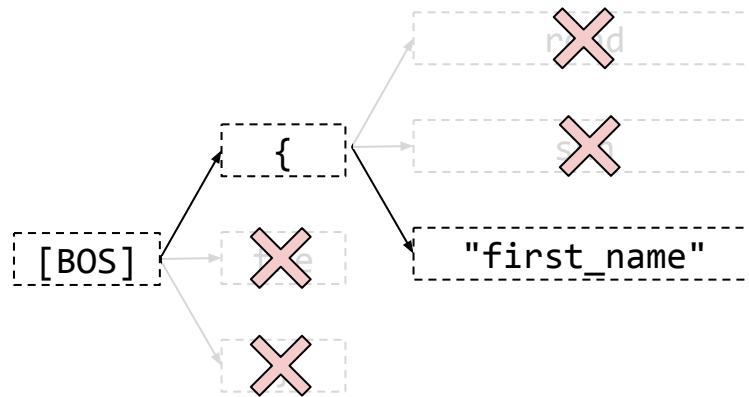
# Constraining the output via guided decoding

**Idea.** Only allow "valid" next tokens



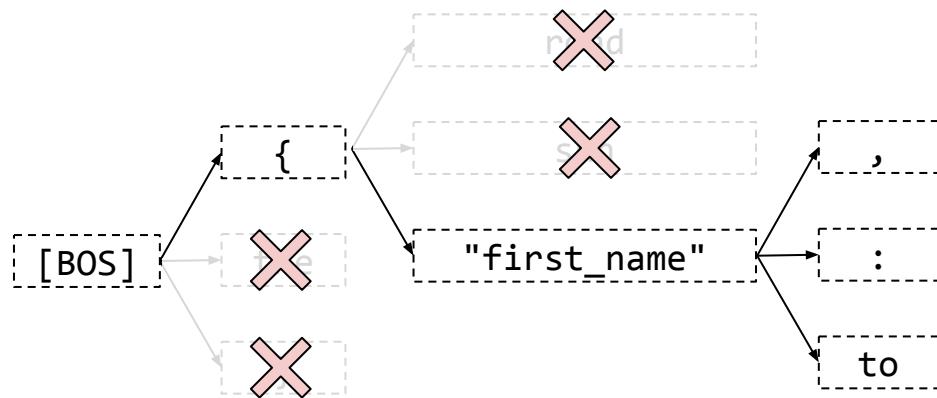
# Constraining the output via guided decoding

**Idea.** Only allow "valid" next tokens



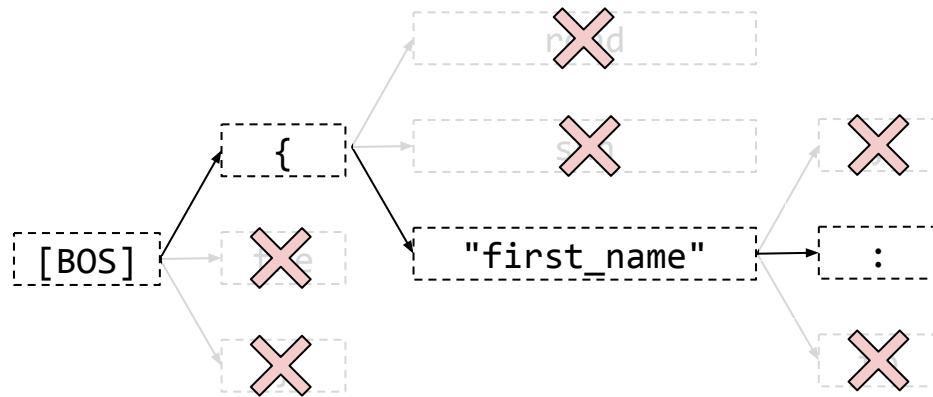
# Constraining the output via guided decoding

**Idea.** Only allow "valid" next tokens



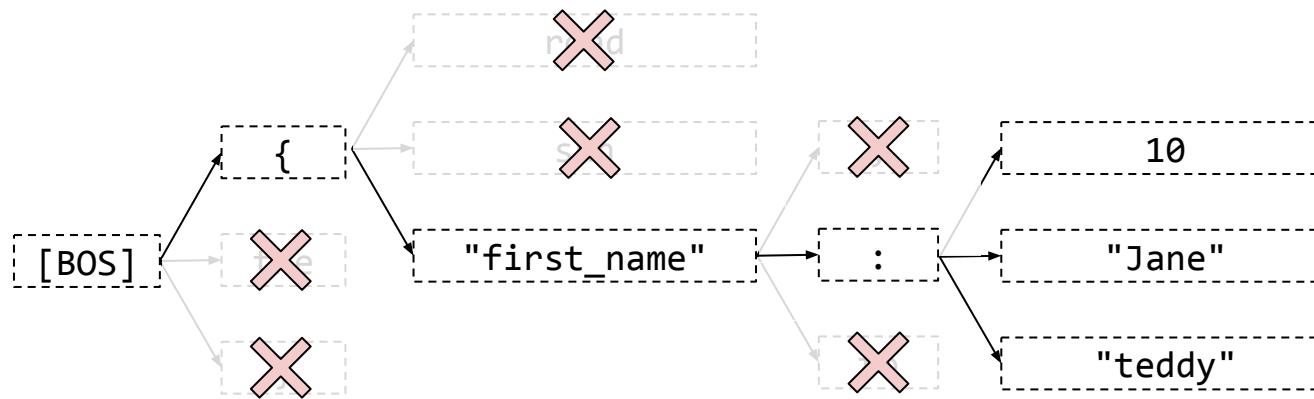
# Constraining the output via guided decoding

**Idea.** Only allow "valid" next tokens



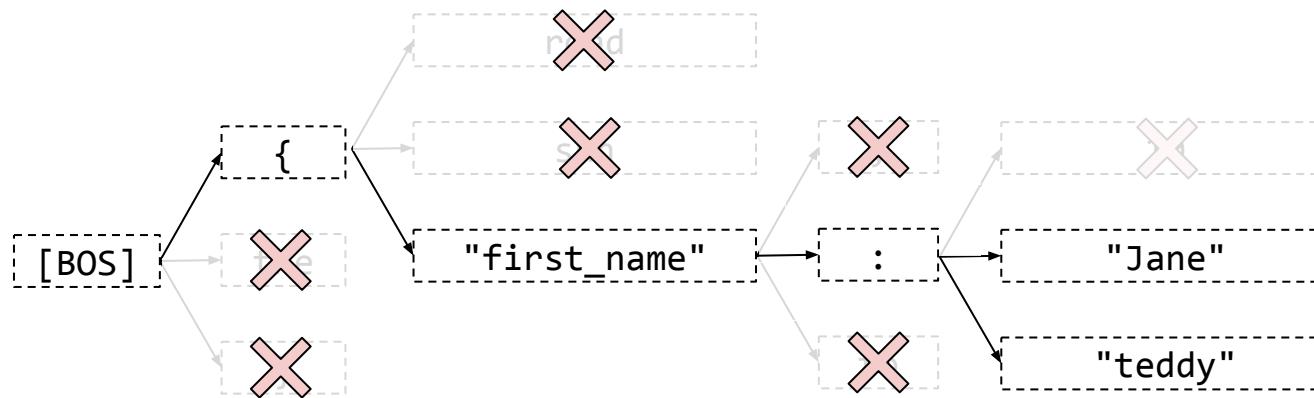
# Constraining the output via guided decoding

**Idea.** Only allow "valid" next tokens



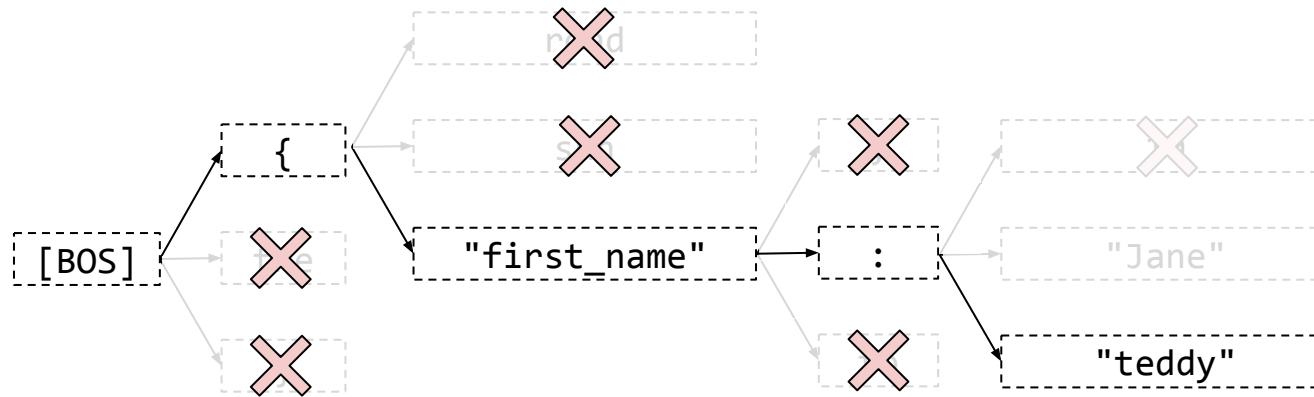
# Constraining the output via guided decoding

**Idea.** Only allow "valid" next tokens



# Constraining the output via guided decoding

**Idea.** Only allow "valid" next tokens





# Transformers & Large Language Models

LLM overview

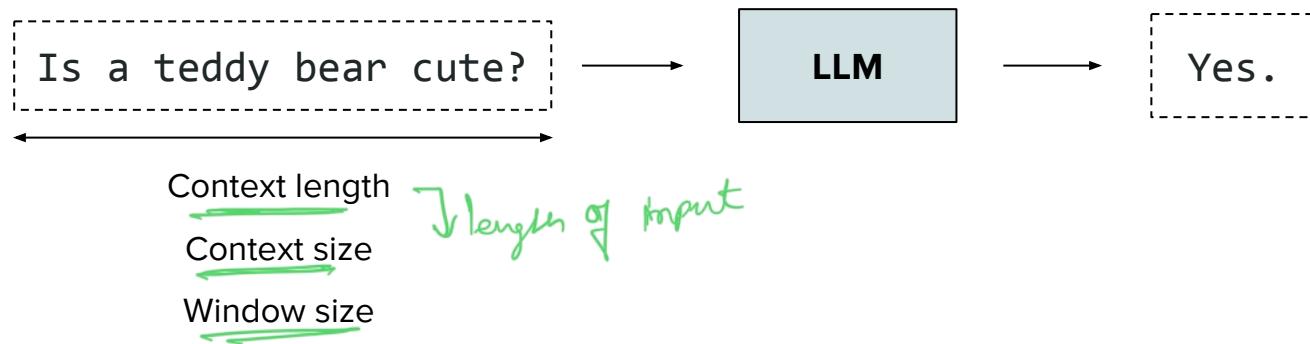
MoE-based LLMs

Response generation

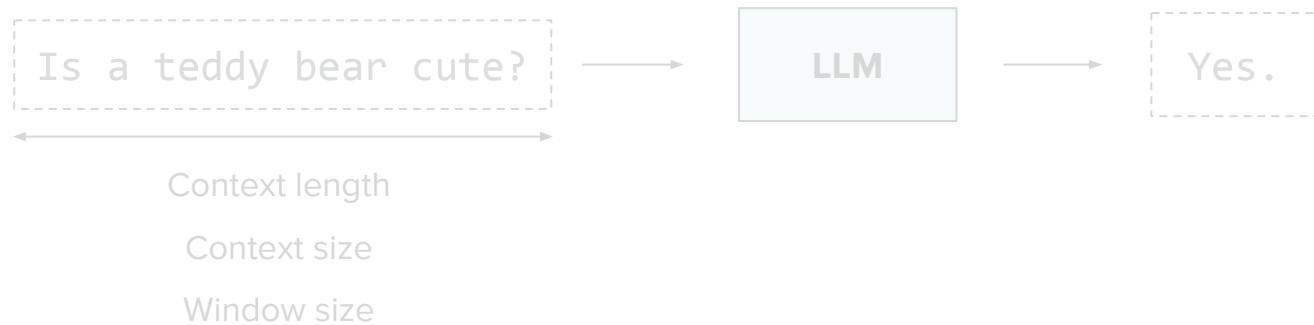
**Prompting strategies**

Inference optimizations

# Terminology ← prompting strategies



# Terminology



**Discussion.** Orders of magnitude by:

- Input type
- Models

# Terminology



**Discussion.** Orders of magnitude by:

- ✓ Input type
- ✓ Models

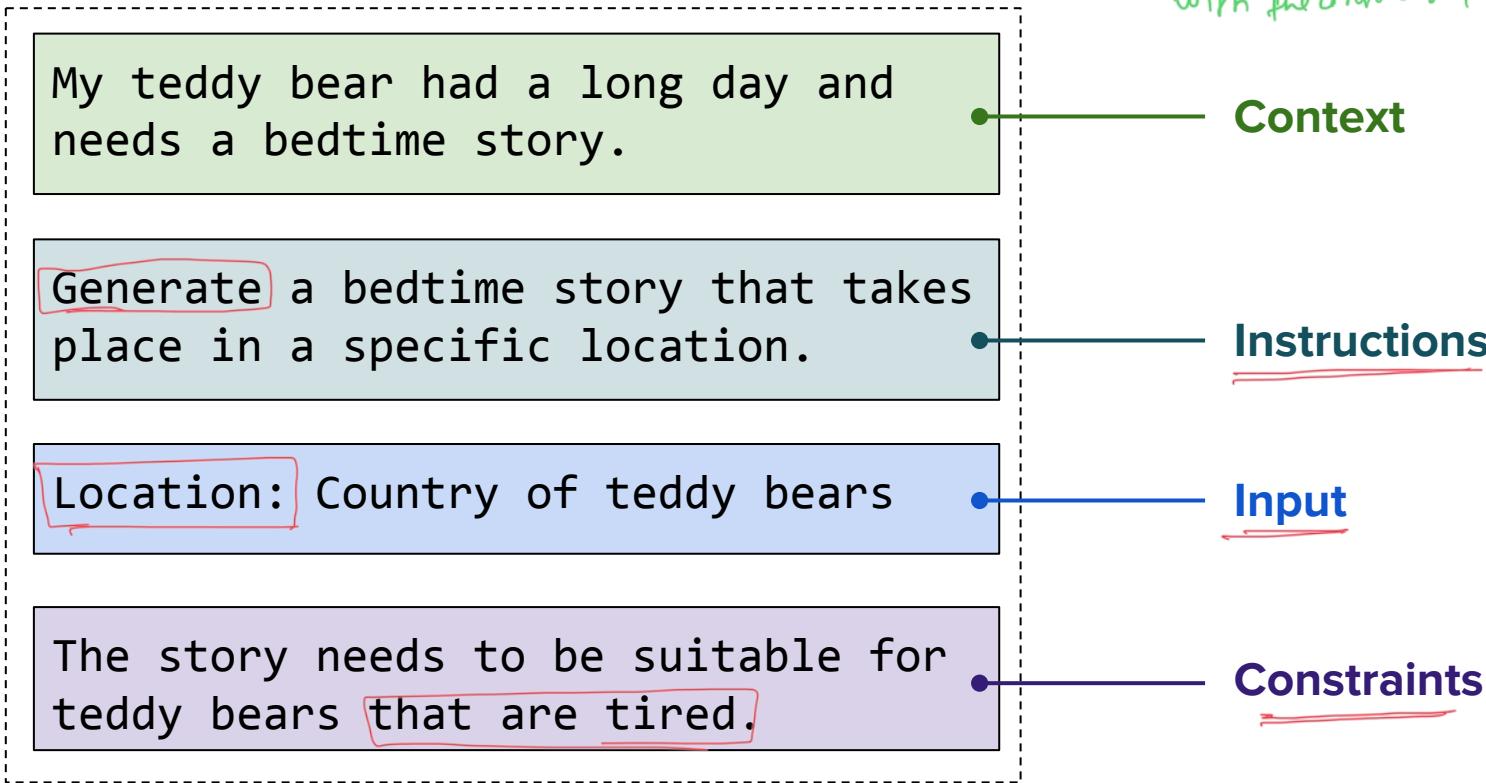


Beware of "context rot"!

Experimenting the capacity of model with different lengths of context.  
e.g. Needles in a haystack → Experimentation?

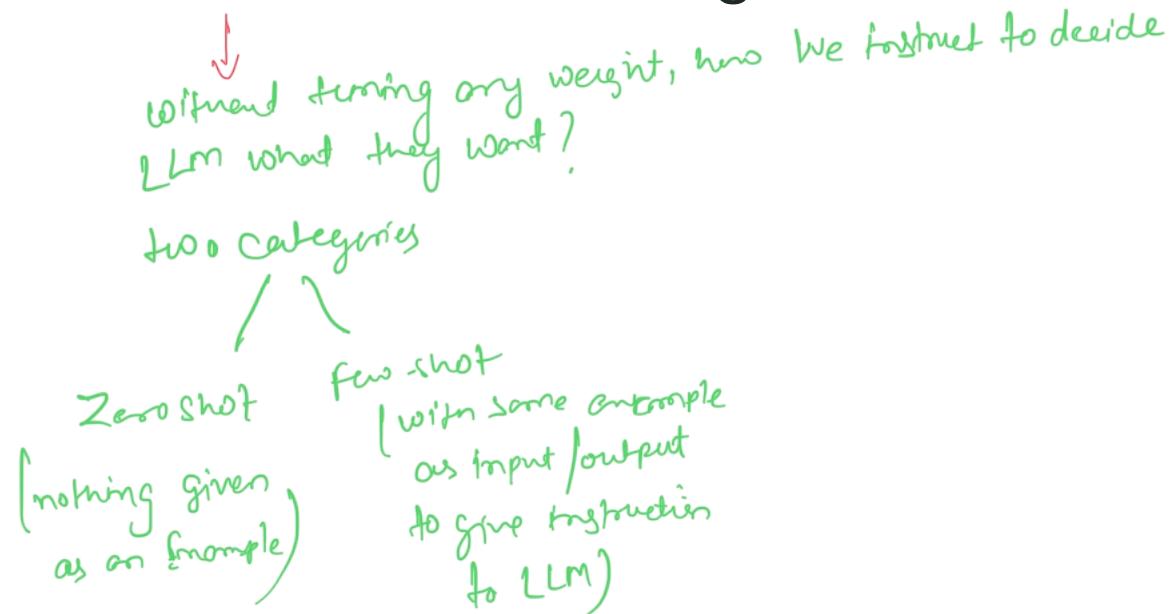
# Main structure

At increasing of content length;  
ability to general the information  
with fine answer i.e in content decrease



# In-context learning

## ICL = In-Context Learning



# In-context learning

ICL = In-Context Learning

Zero-shot learning	Few-shot learning
<ul style="list-style-type: none"><li>• Question is asked without examples</li><li>• Performance heavily depends on performance of initial model</li></ul>	<ul style="list-style-type: none"><li>• Prompt contains examples of input / output</li><li>• Typically has a better performance</li></ul> <p>(if we provide example, we kind of giving instructions to make the connection b/w input and output for understanding the connection which literally improves the performance in case of prediction)</p>

# In-context learning

**ICL = In-Context Learning**

Zero-shot learning	Few-shot learning
<ul style="list-style-type: none"><li>• Question is asked without examples</li><li>• Performance heavily depends on performance of initial model</li></ul>	<ul style="list-style-type: none"><li>• Prompt contains examples of input / output</li><li>• Typically has a better performance</li></ul>

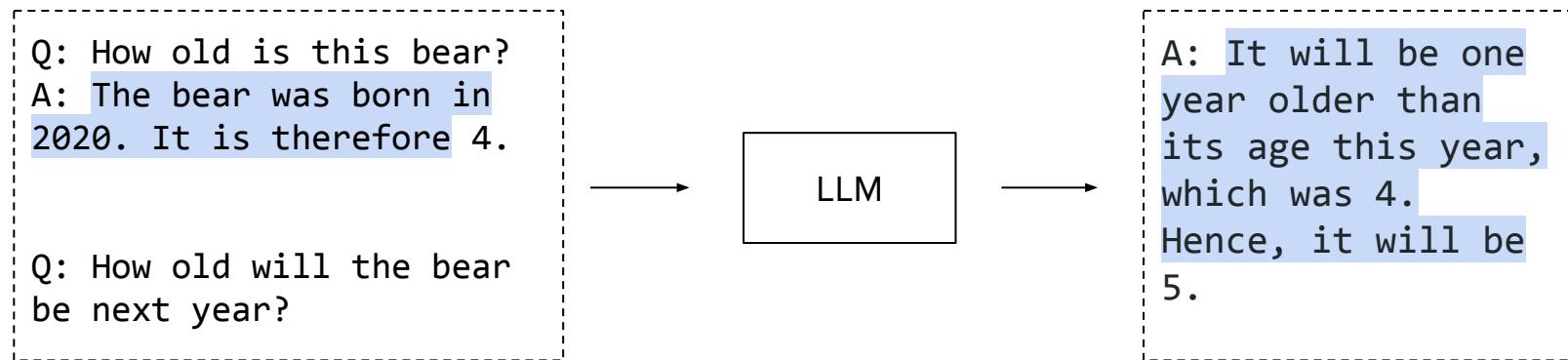
**Discussion.** Showing examples in the prompt is *generally* better but:

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

Chain of thought - if we force the model to come with some example as Q/A  
it improves the performance.

→ root causing will be easy for improved results

Idea. Explaining reasoning helps in improving performance.



## Discussion.

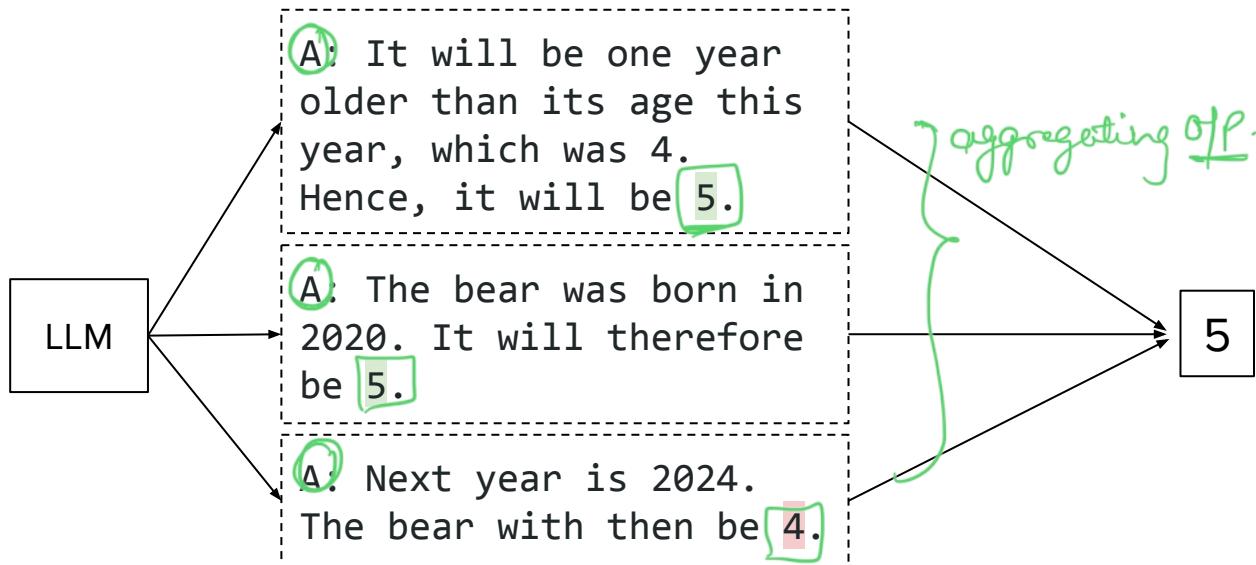
- Interpretability + explanation
- More tokens: higher cost and latency

# Self-consistency

↓

Samples the model with different response and aggregating it for better response

**Idea.** Aggregating over reasoning paths improves performance.



**Discussion.** Trade-off between performance and added cost.



# Transformers & Large Language Models

LLM overview

MoE-based LLMs

Response generation

Prompting strategies

**Inference optimizations**

# Challenges

**Motivation.** Computations are expensive, any way to reduce complexity?

**Categories.** Many dimensions to optimize for.

# Challenges

**Motivation.** Computations are expensive, any way to reduce complexity?

**Categories.** Many dimensions to optimize for.

"Exact" efficiency:

- ① • Avoid redundancies
  - ② • Memory management
  - ③ • Reformulate the math
- } reducing complexity  
but better performance

# Challenges

**Motivation.** Computations are expensive, any way to reduce complexity?

**Categories.** Many dimensions to optimize for.

"Exact" efficiency:

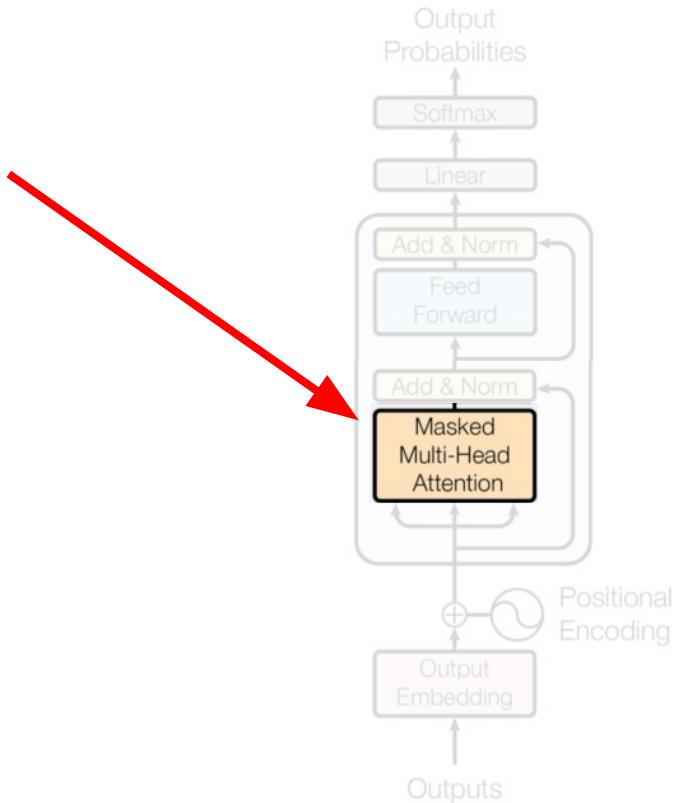
- Avoid redundancies
- Memory management
- Reformulate the math

Approximations:

- ③ ● Architectural changes
- ④ ● Embeddings representations
- ⑤ ● Token prediction

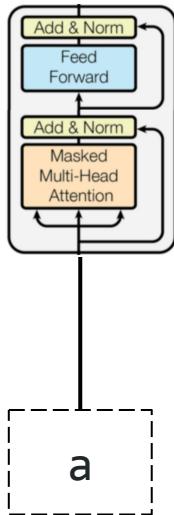
less cost, but  
better off at prediction

# Attention-based tricks



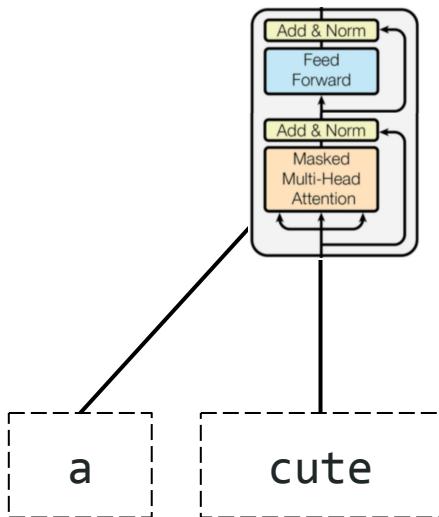
# Reusing computations with KV caching

**Motivation.** New token needs to interact with all previous tokens.



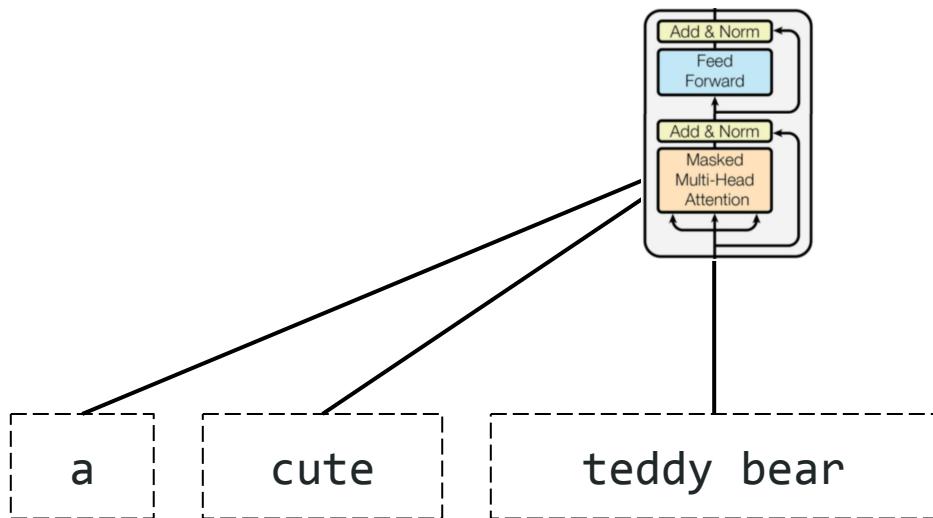
# Reusing computations with KV caching

**Motivation.** New token needs to interact with all previous tokens.



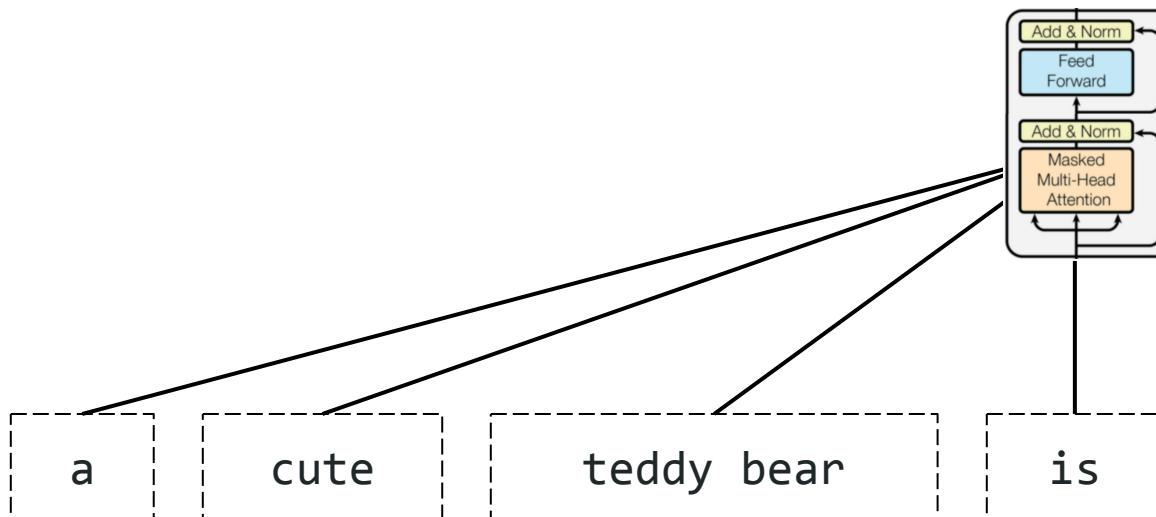
# Reusing computations with KV caching

**Motivation.** New token needs to interact with all previous tokens.



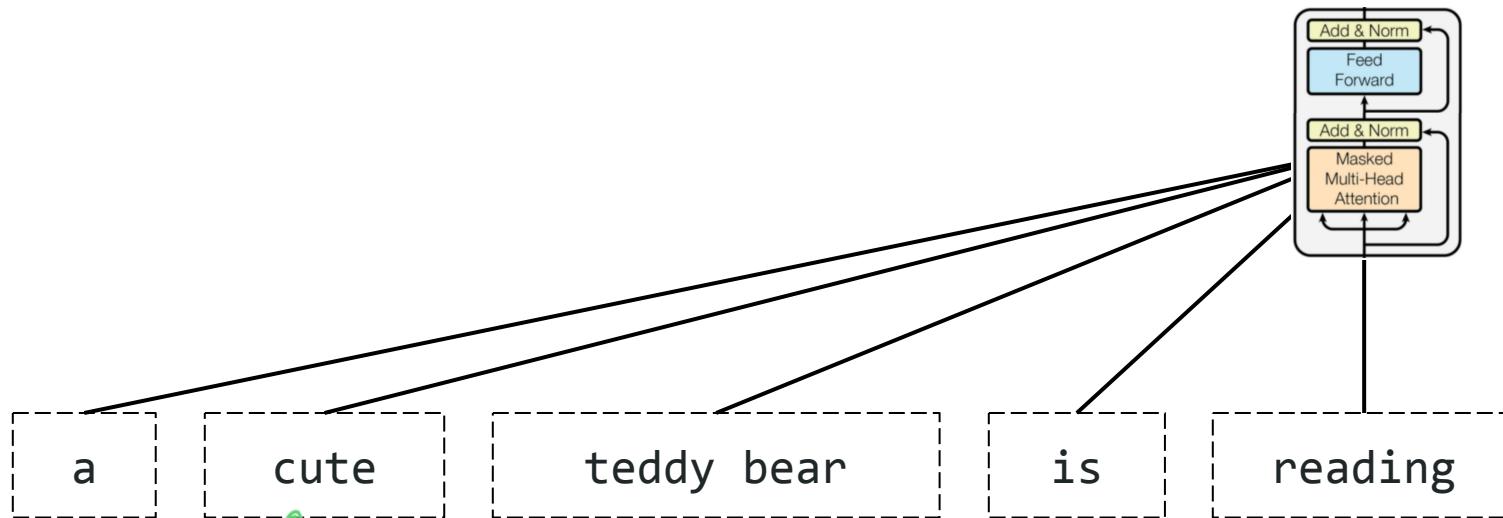
# Reusing computations with KV caching

**Motivation.** New token needs to interact with all previous tokens.



# Reusing computations with KV caching

**Motivation.** New token needs to interact with all previous tokens.



Each word token is predicted based on previous tokens.

So, "cute" is attending "a" and "teddy bear" is attending "a" and "cute".  
So, we need to keep the predicted already previous token for further bcz it has token for further bcz it has

# Reusing computations with KV caching

be organised for a such a way  
that (key,value) should be safe

• KV-caching stores (key,value)  
pair so, need to find the  
way to use it instead  
of parsing previous token  
again and again  
and used in inference-  
time.

$$\text{softmax} \left( \frac{Q \boxed{K^T}}{\sqrt{d_k}} \right) \boxed{V}$$

# Reusing computations with KV caching

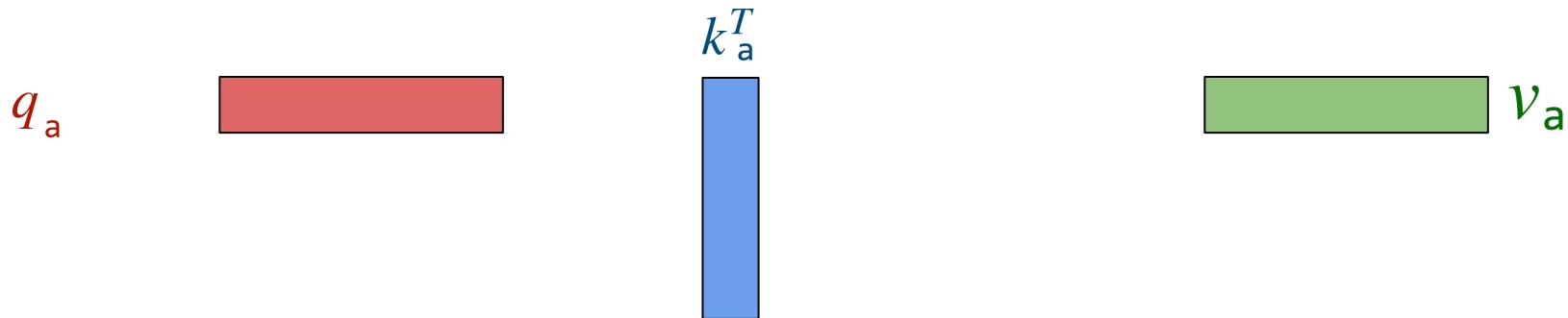
**Motivation.** New token needs to interact with all previous tokens.

**Idea.** Keep keys and values in a cache.

# Reusing computations with KV caching

**Motivation.** New token needs to interact with all previous tokens.

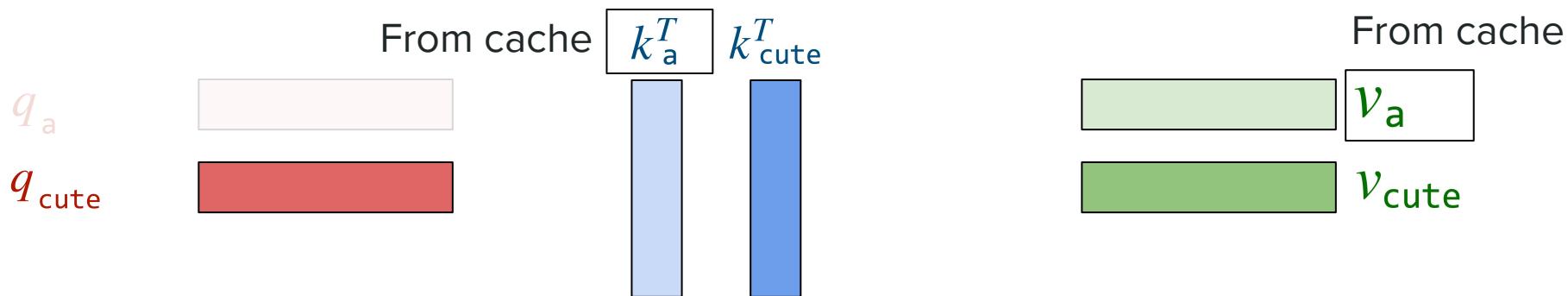
**Idea.** Keep keys and values in a cache.



# Reusing computations with KV caching

**Motivation.** New token needs to interact with all previous tokens.

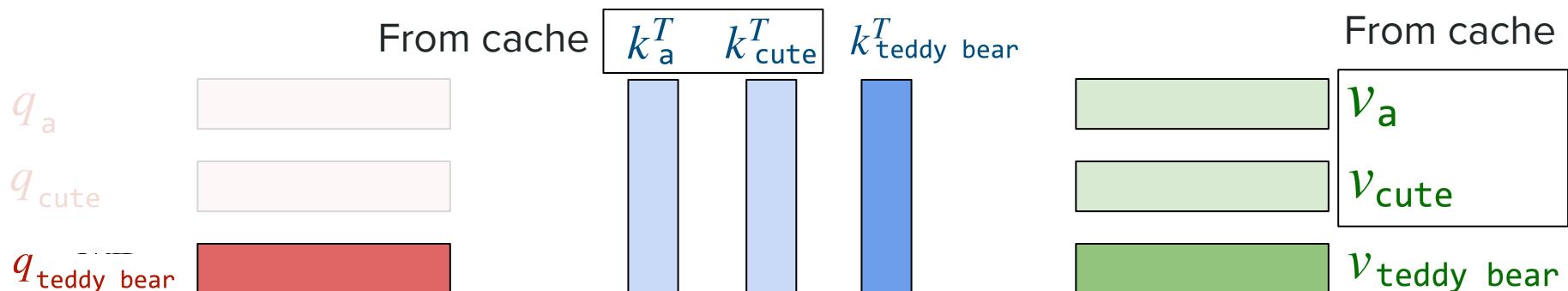
**Idea.** Keep keys and values in a cache.



# Reusing computations with KV caching

**Motivation.** New token needs to interact with all previous tokens.

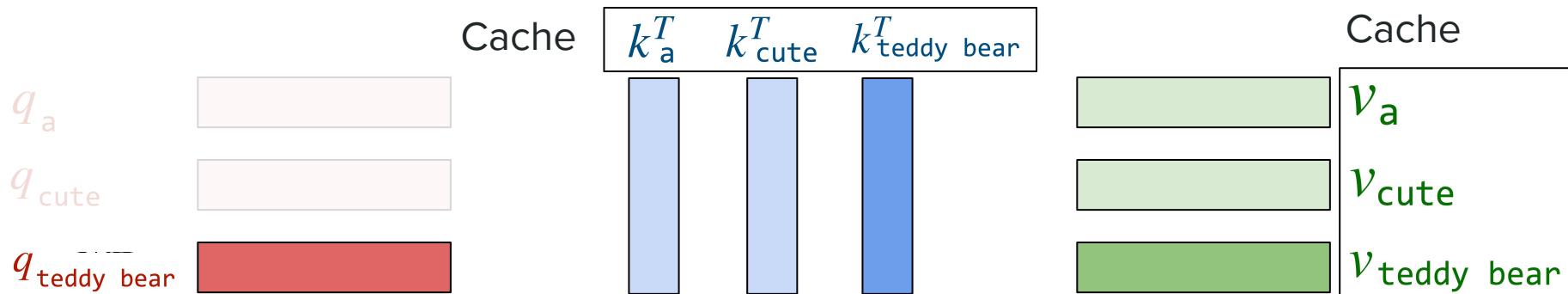
**Idea.** Keep keys and values in a cache.



# Reusing computations with KV caching

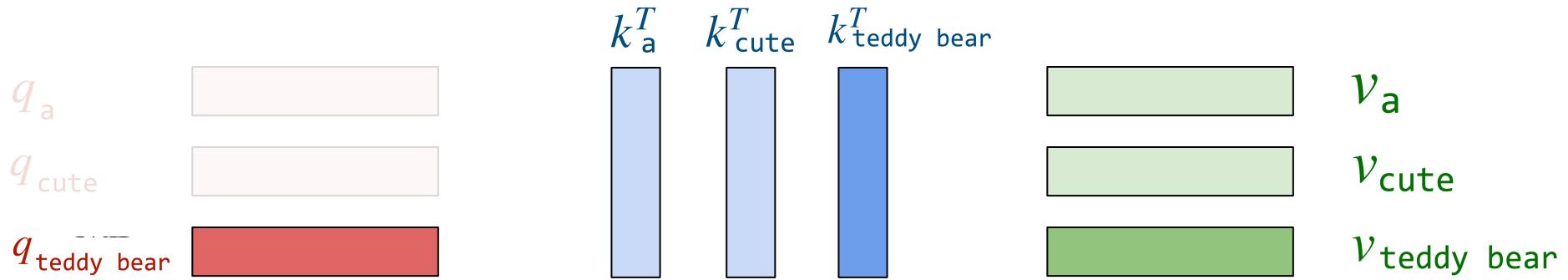
**Motivation.** New token needs to interact with all previous tokens.

**Idea.** Keep keys and values in a cache.

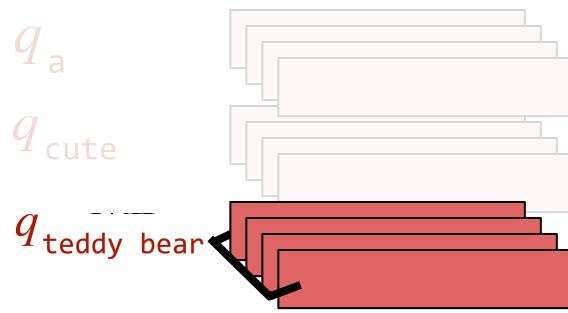


**KV caching = Key-Value Caching**

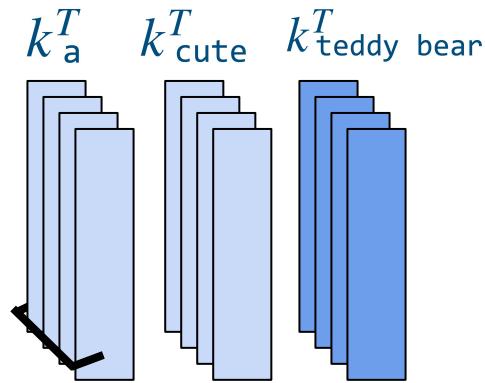
# Sharing attention heads



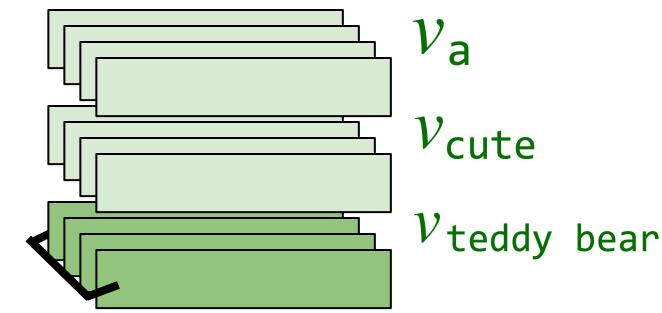
# Sharing attention heads



Number of  
**query heads**



Number of  
**key heads**

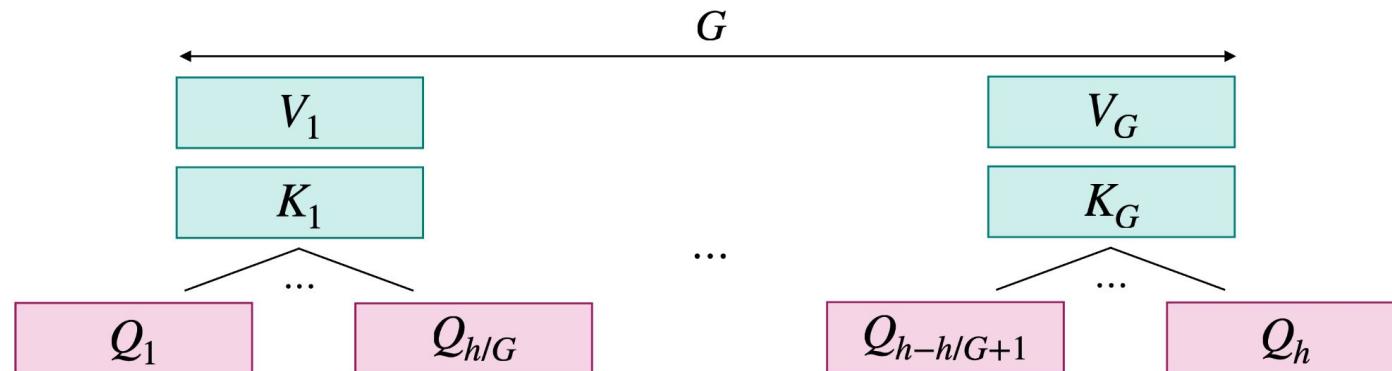


Number of  
**value heads**

$$\text{In vanilla MHA,} \quad \frac{\text{Number of query heads}}{=} \quad \frac{\text{Number of key heads}}{=} \quad \frac{\text{Number of value heads}}{=} \quad h$$

# Sharing attention heads

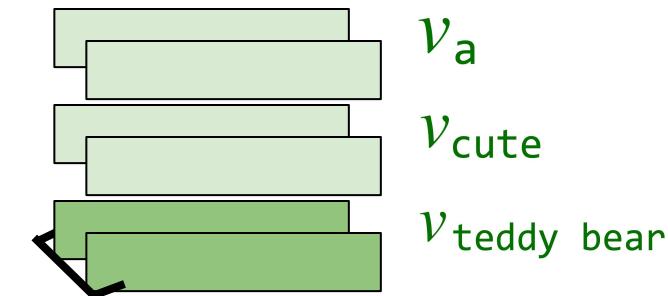
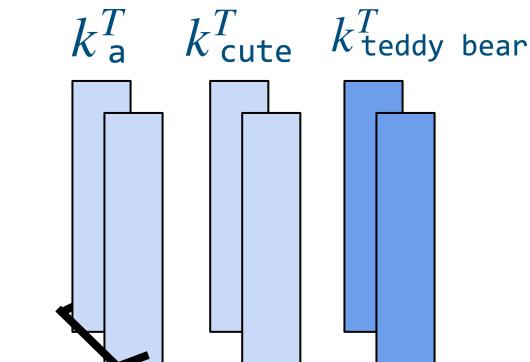
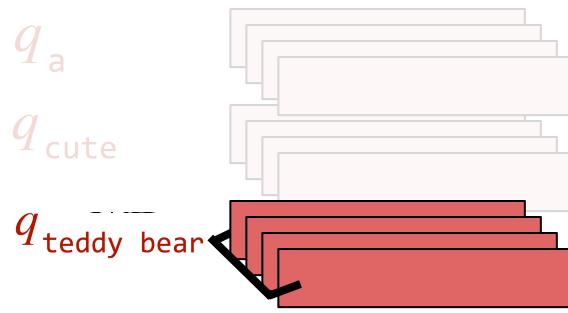
**Idea.** Share key/value attention heads within groups of queries



# Sharing attention heads

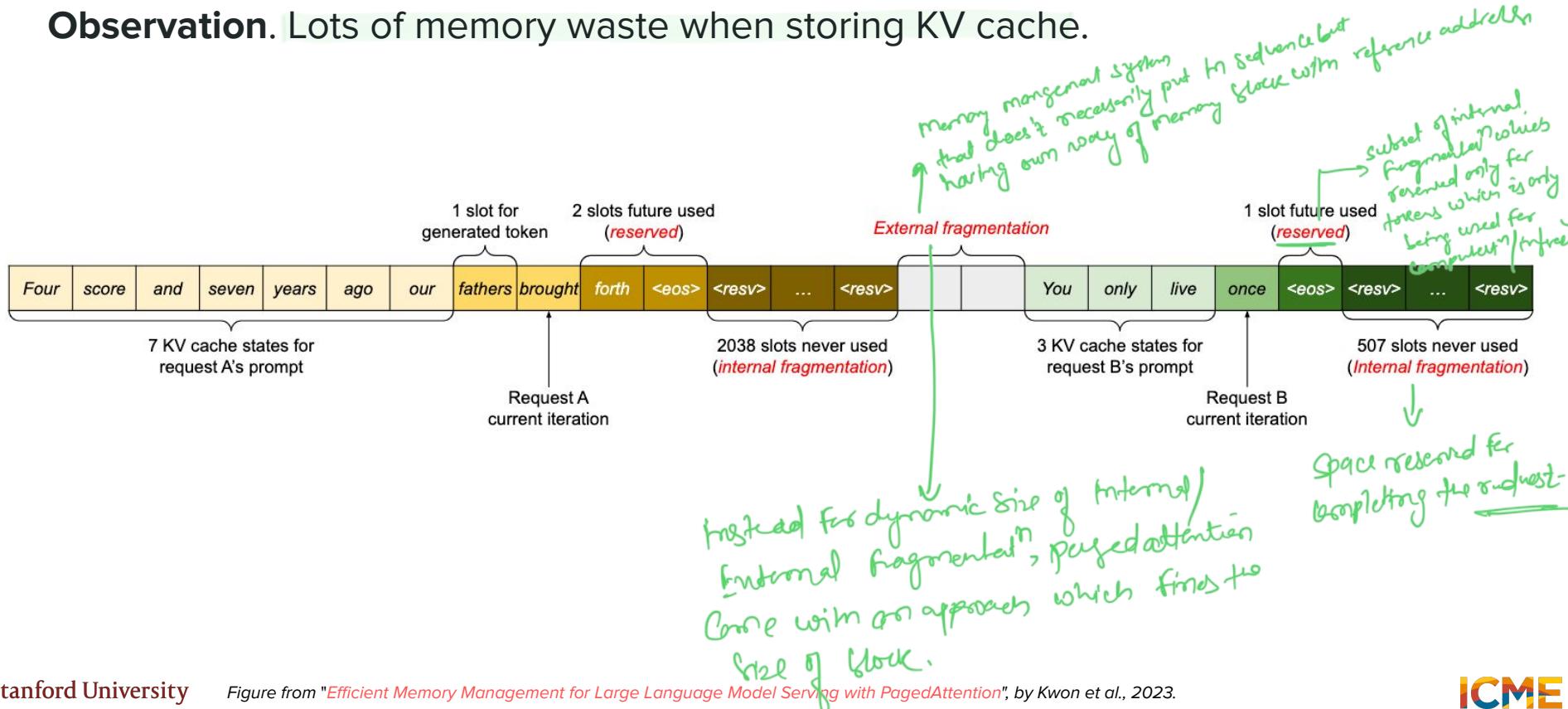
$G = 1$	<b>Multi-Query Attention (MQA)</b>	<p>The diagram shows a dashed box containing two stacked teal boxes labeled <math>V</math> and <math>K</math>. A single pink box labeled <math>Q_h</math> is connected to both <math>V</math> and <math>K</math> via lines. Below the dashed box is a multiplier <math>\times 1</math>.</p>
$1 < G < h$	<b>Group-Query Attention (GQA)</b>	<p>The diagram shows a dashed box containing two stacked teal boxes labeled <math>V</math> and <math>K</math>. Multiple pink boxes labeled <math>Q_1</math>, <math>Q_2</math>, ..., <math>Q_{h/G}</math> are connected to both <math>V</math> and <math>K</math> via lines. Below the dashed box is a multiplier <math>\times G</math>.</p>
$G = h$	<b>Multi-Head Attention (MHA)</b>	<p>The diagram shows a dashed box containing two stacked teal boxes labeled <math>V</math> and <math>K</math>. One pink box labeled <math>Q</math> is connected to both <math>V</math> and <math>K</math> via lines. Below the dashed box is a multiplier <math>\times h</math>.</p>

# Sharing attention heads with GQA



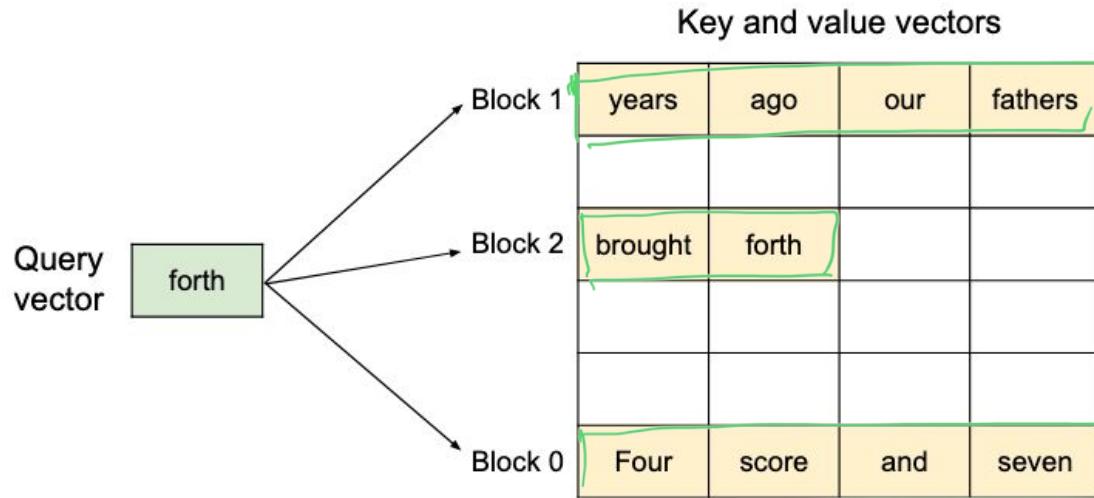
# Manage memory with PagedAttention

**Observation.** Lots of memory waste when storing KV cache.



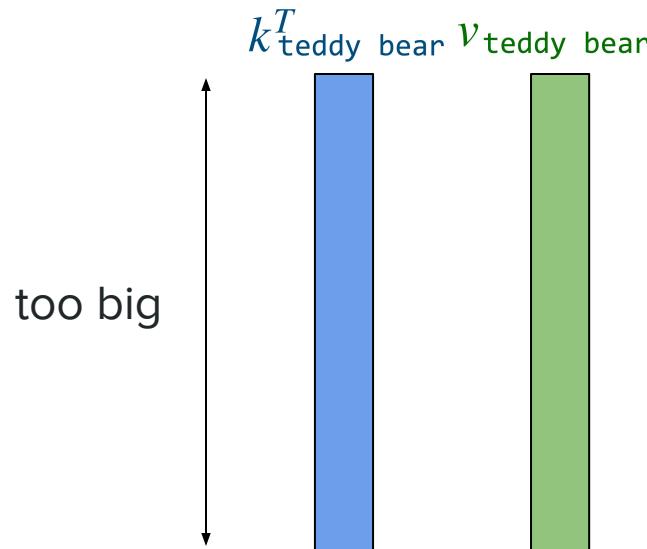
# Manage memory with PagedAttention

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



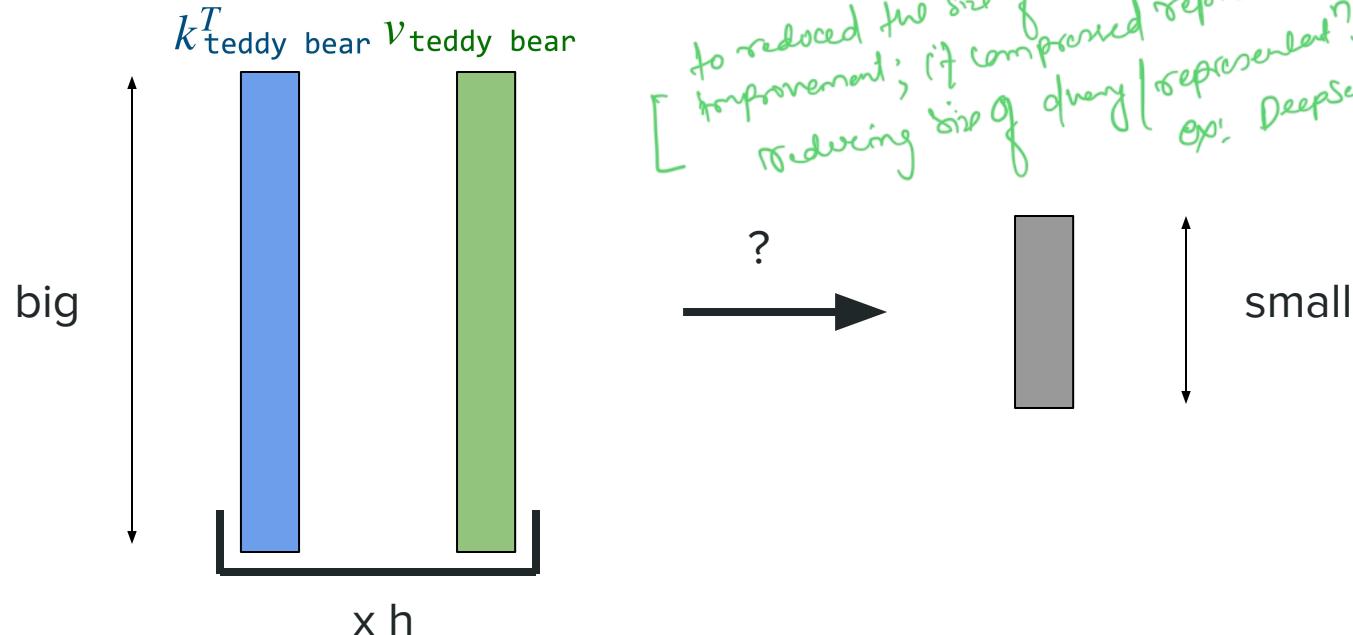
# Reduce memory of KV cache with latent attention

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



# Reduce memory of KV cache with latent attention

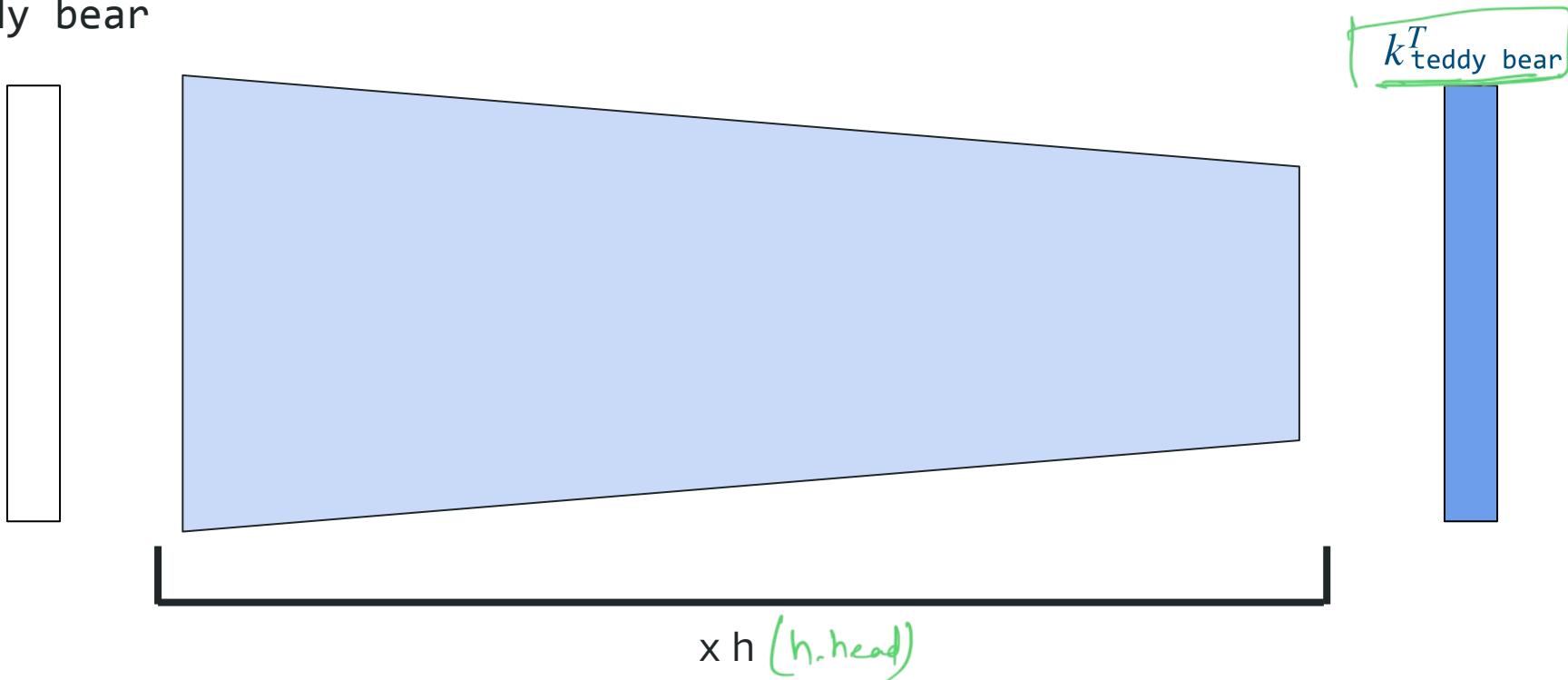
**Solution.** Store compressed representations instead!



# Reduce memory of KV cache with latent attention

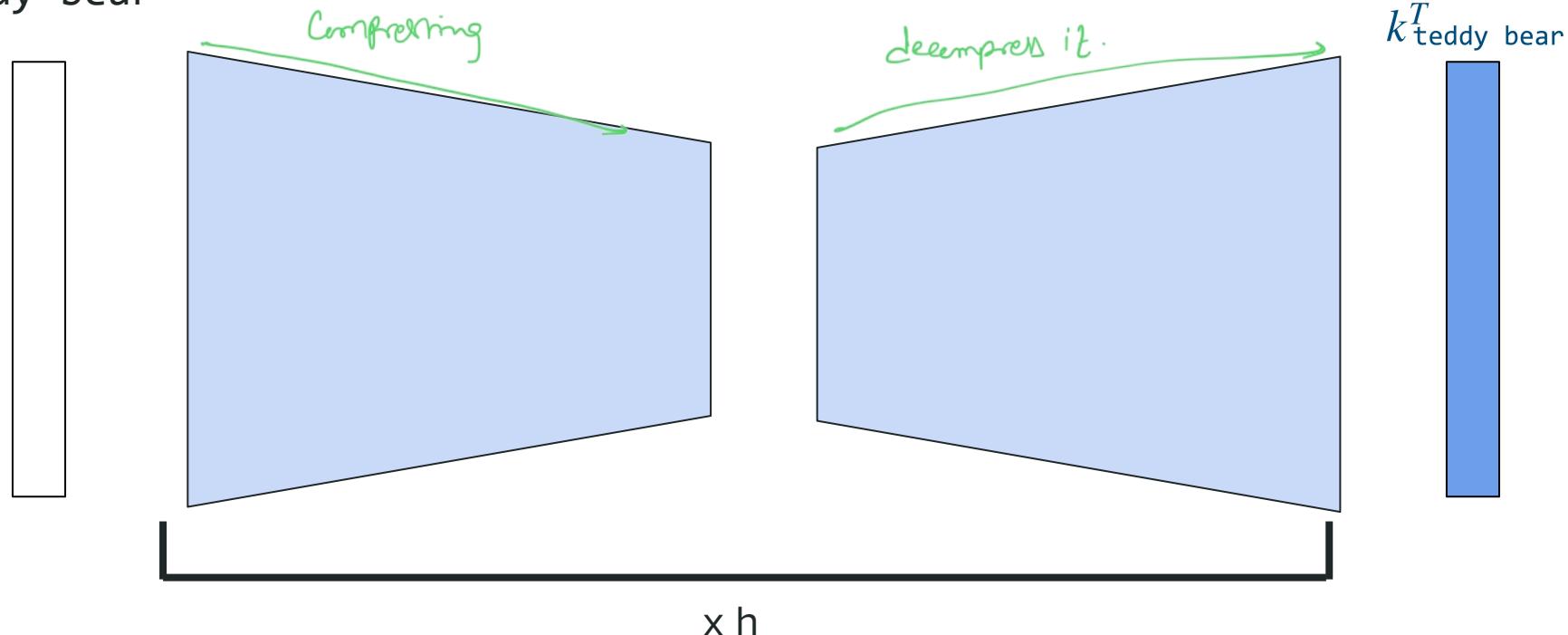
BEFORE

teddy bear

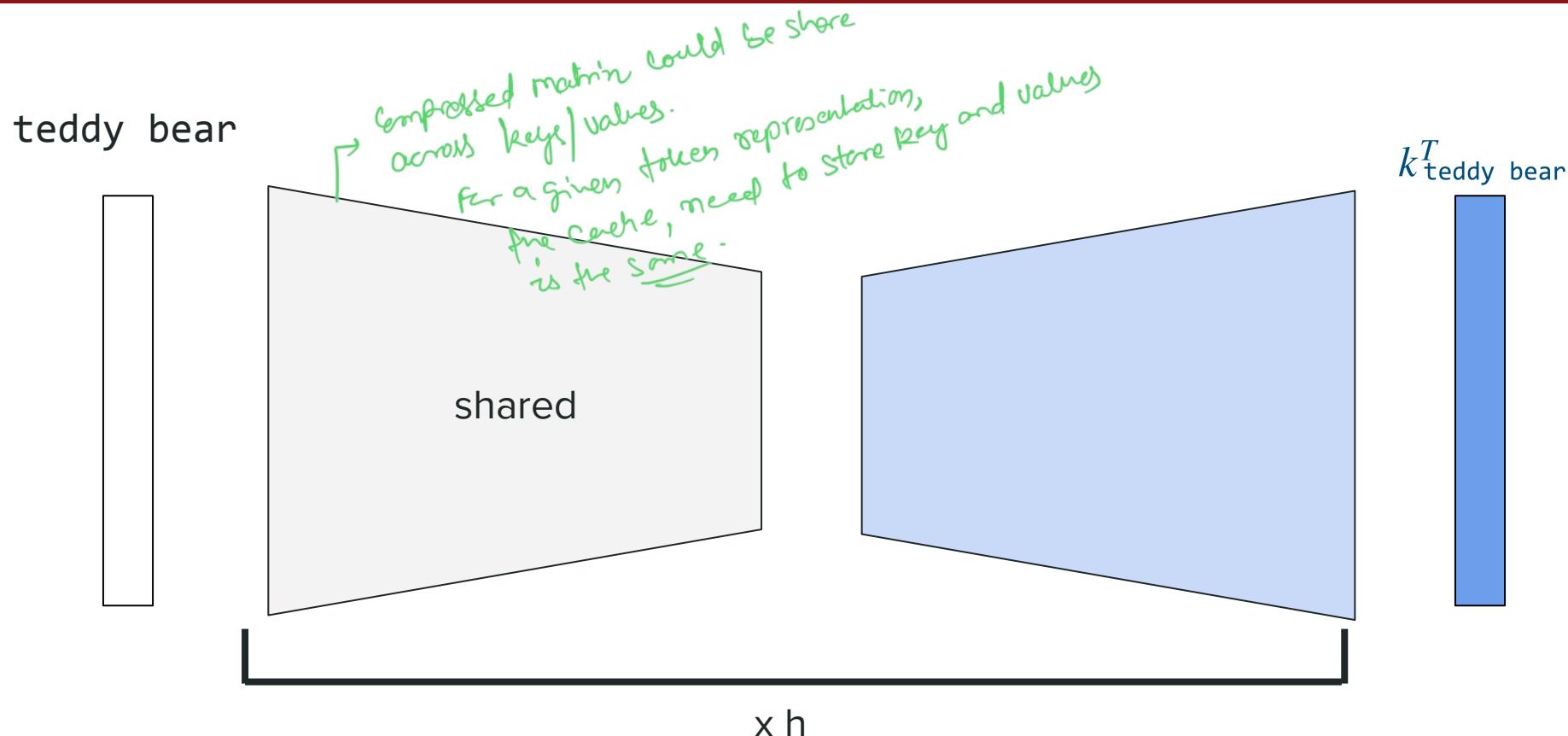


# Reduce memory of KV cache with latent attention

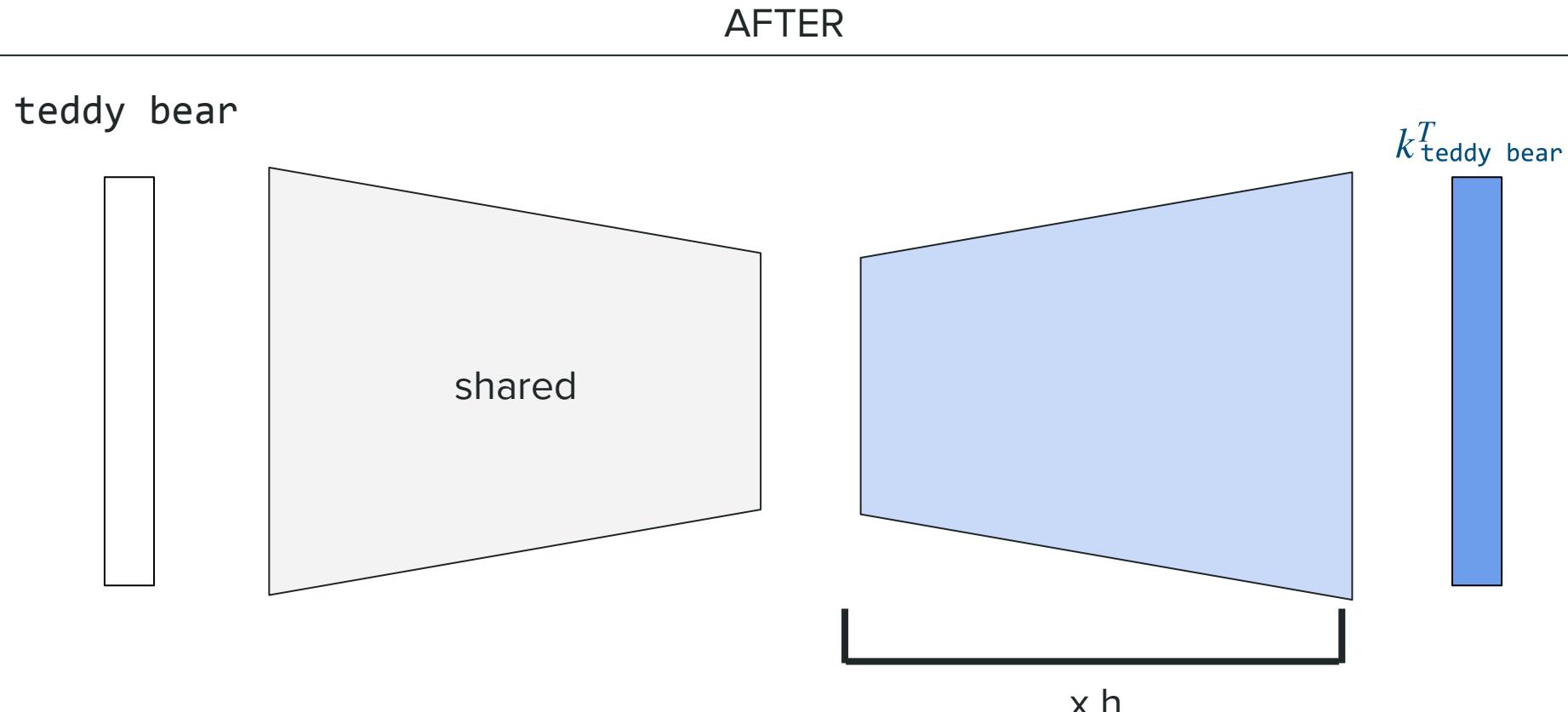
teddy bear



# Reduce memory of KV cache with latent attention



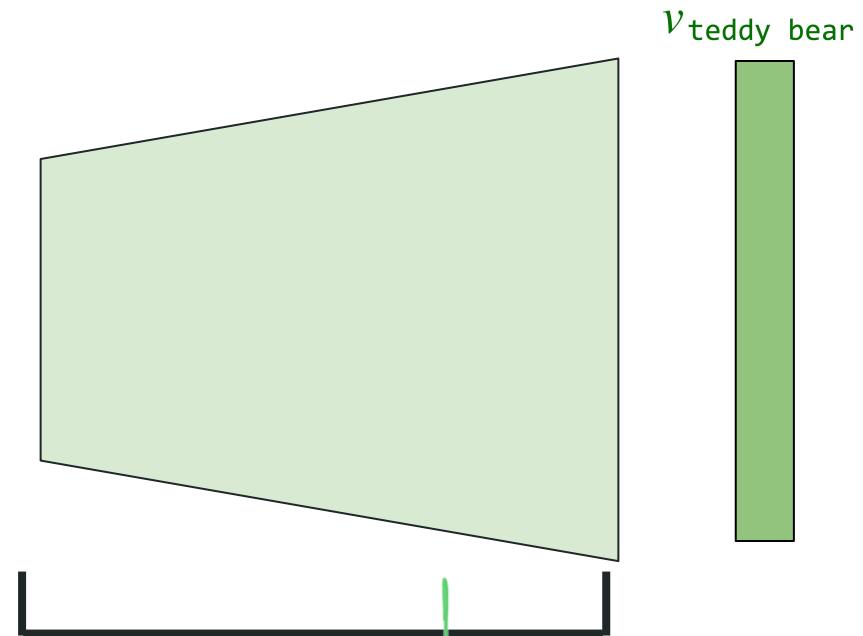
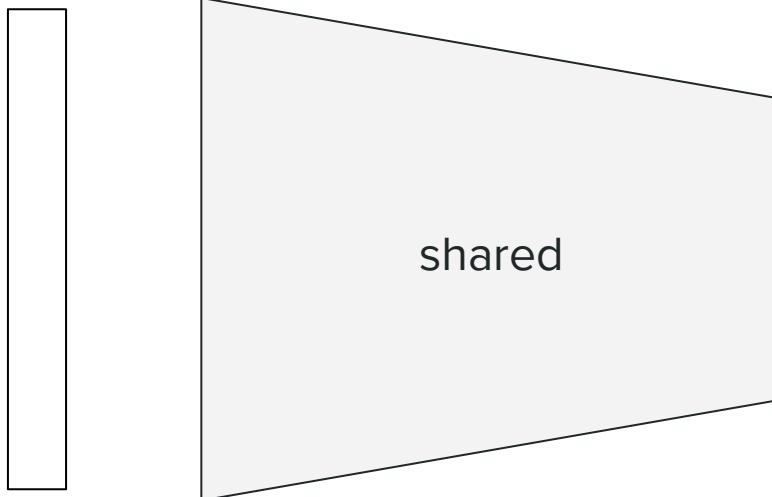
# Reduce memory of KV cache with latent attention



# Reduce memory of KV cache with latent attention

AFTER

teddy bear

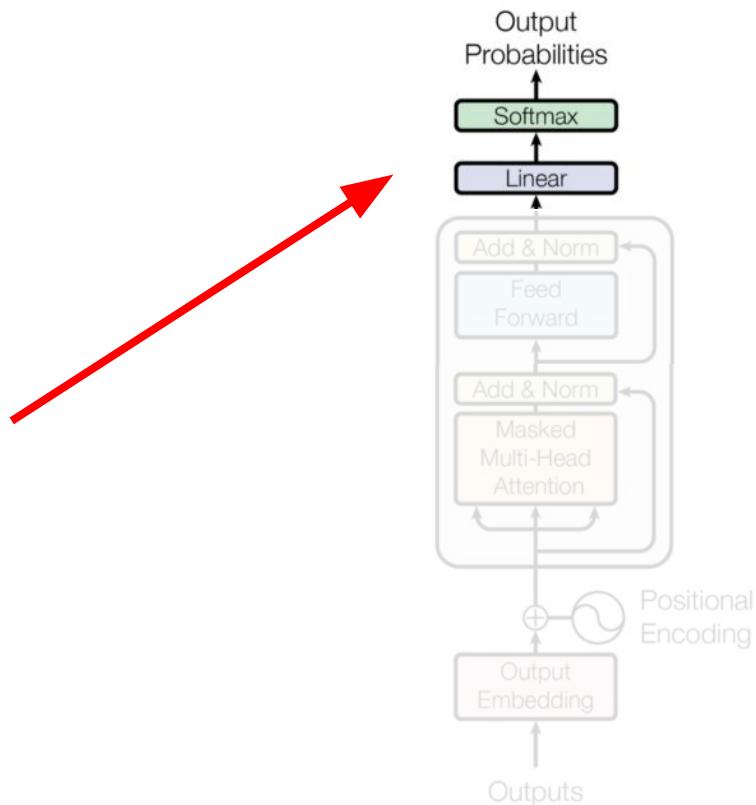


$x h$

speeding decoding parts is  
being also discussed.

ICME

# Token generation tricks



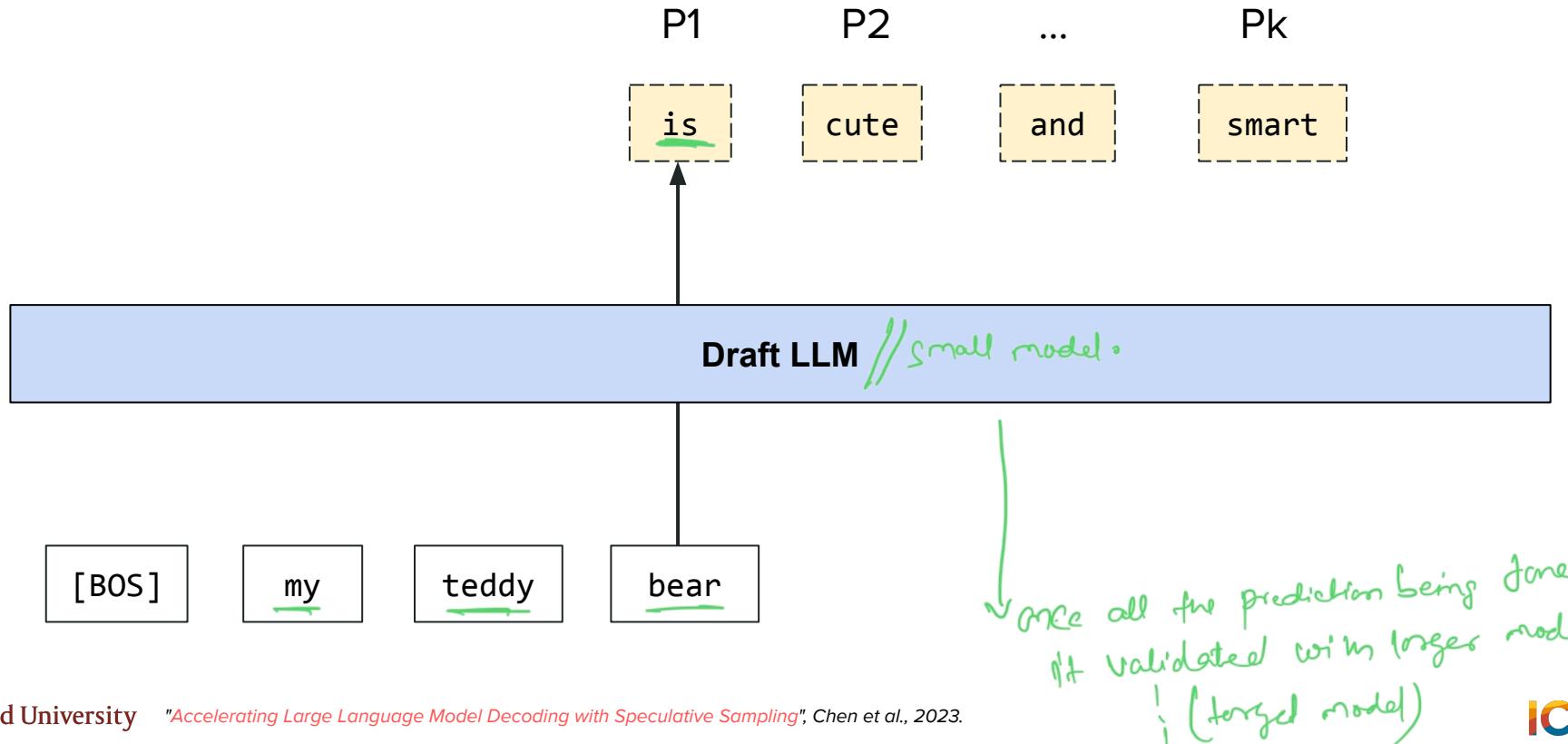
# Speeding up decoding with speculative decoding

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

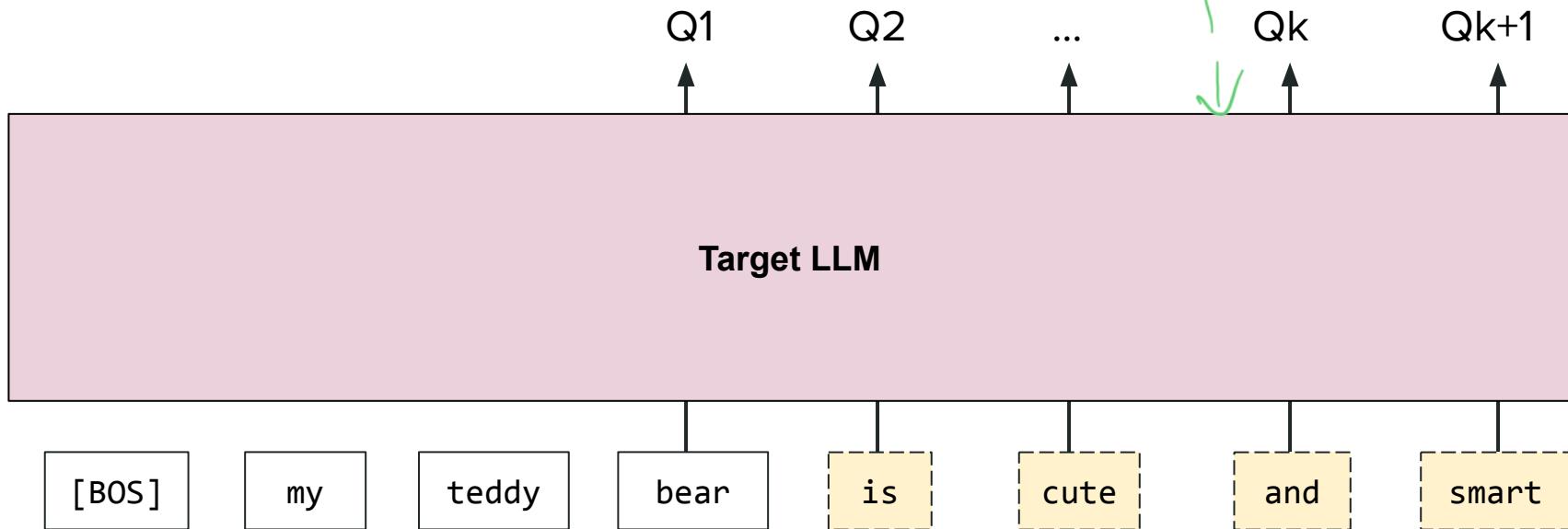


Speculative decoding technique  
which uses "smaller model" to help  
to generation to bigger model and uses  
some scheme that make distributions outputs  
match one of the target one.  
(uses smaller model but with the validation by big model)

# Speeding up decoding with speculative decoding



# Speeding up decoding with speculative decoding



# Speeding up decoding with speculative decoding

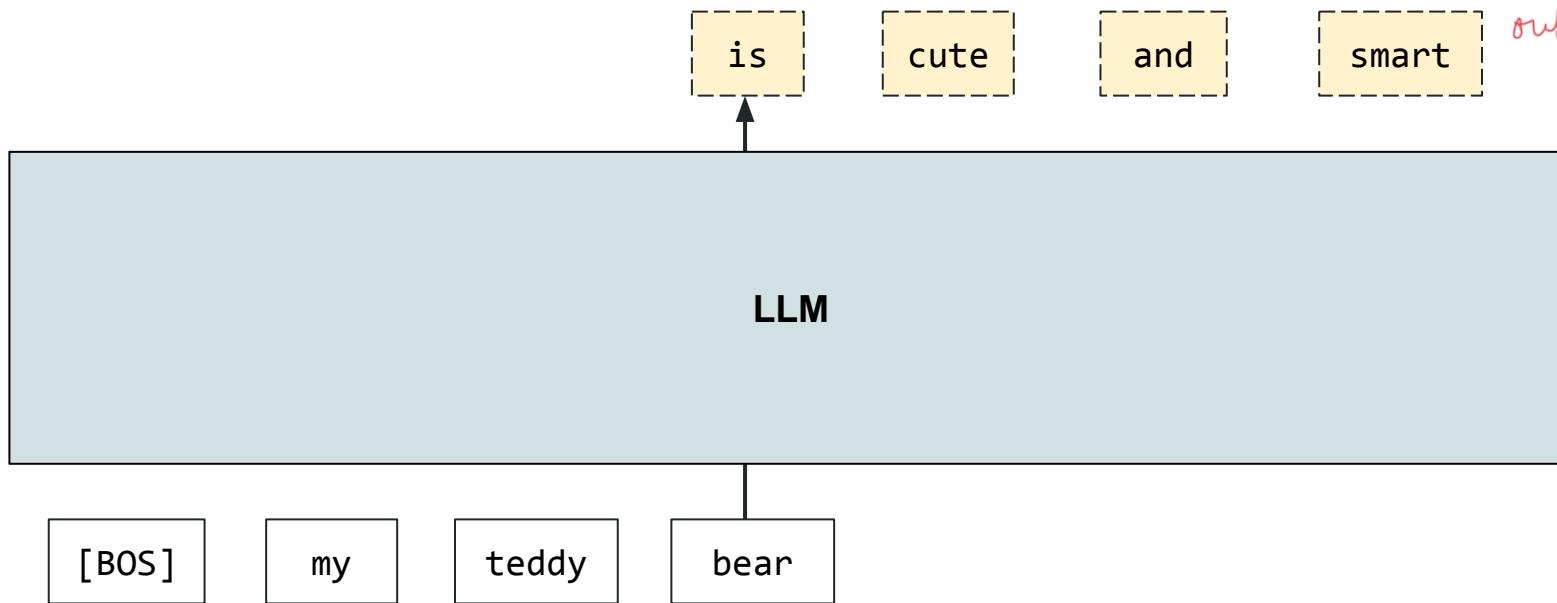
**Sampling algorithm.** We distinguish the following cases:

- If  $Q_i(\boxed{\text{token}}) \geq P_i(\boxed{\text{token}})$  [longer model prediction probability of token] smaller model prediction probability than  $P_i(\boxed{\text{token}})$  ✓  
accepted (the token predicted)
- Otherwise
  - ✓ with probability  $Q_i(\boxed{\text{token}}) / P_i(\boxed{\text{token}})$  ✓
  - with probability  $1 - Q_i(\boxed{\text{token}}) / P_i(\boxed{\text{token}})$  ✗

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

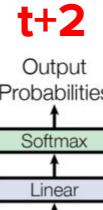
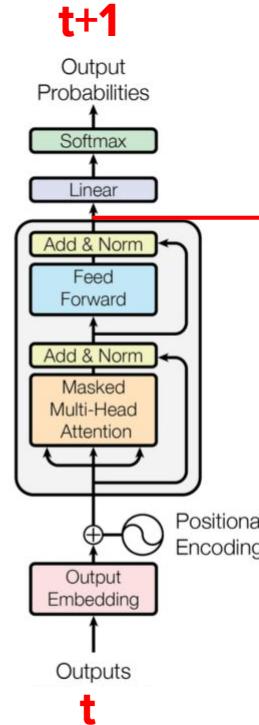
# Generate several tokens at once via MTP

## MTP = Multi-Token Prediction

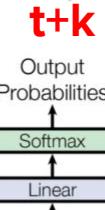


# Generate several tokens at once via MTP

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



...



*attaches multiple head.  
(multiple head prediction)  
and learns which head works  
well in inference)*

# Challenges

**Categories.** Many dimensions to optimize for.

"Exact" efficiency:

- Avoid redundancies
- Memory management
- Reformulate the math

Approximations:

- Architectural changes
- Embeddings representations
- Token prediction

# Challenges and some remedies

**Categories.** Many dimensions to optimize for.

"Exact" efficiency:

- ✓ Avoid redundancies
- ✓ Memory management
- ✓ Reformulate the math

**KV cache**

PagedAttention (storing in fixed block size in memory)

Speculative decoding (small & large model based learning)

Approximations:

- ✓ Architectural changes
- ✓ Embeddings representations
- ✓ Token prediction

**Grouped query attention**

(change architecture for grouping)

**Latent attention**

**Multi-token prediction**

[multiple need into architecture while training]

Thank you for your attention!

---