

CME 295: Transformers & Large Language Models



Afshine Amidi & Shervine Amidi

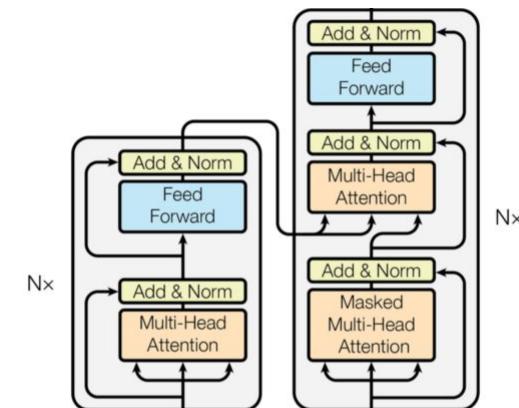


Recap of last episode...

Recap of last episode...

3 categories of Transformer-based models:

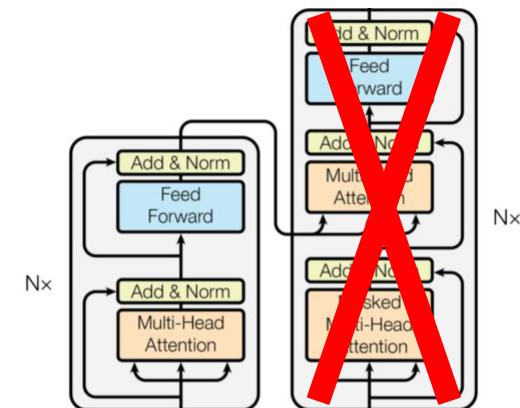
| | | |
|------------------------|--------------|---------------|
| Encoder-decoder | Text to text | T5, mT5, ByT5 |
|------------------------|--------------|---------------|



Recap of last episode...

3 categories of Transformer-based models:

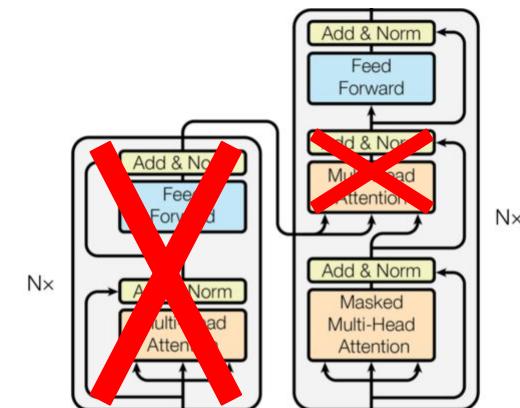
| | | |
|---------------------|--|---------------------------|
| Encoder-decoder | Text to text | T5, mT5, ByT5 |
| Encoder-only | Projection of embedding for class prediction (e.g. sentiment extraction) | BERT, DistilBERT, RoBERTa |



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 |



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 |



Transformers & Large Language Models

LLM overview

MoE-based LLMs

Response generation

Prompting strategies

Inference optimizations

Terminology

LLM = Large Language Model

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

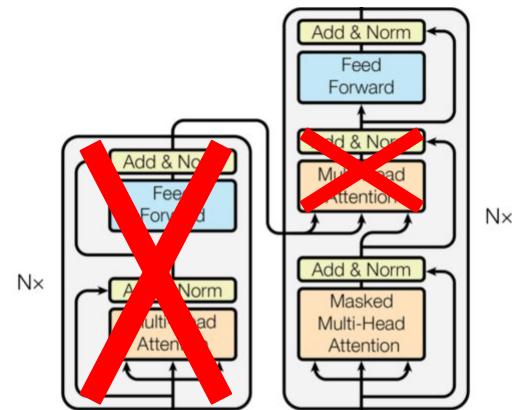
Terminology

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

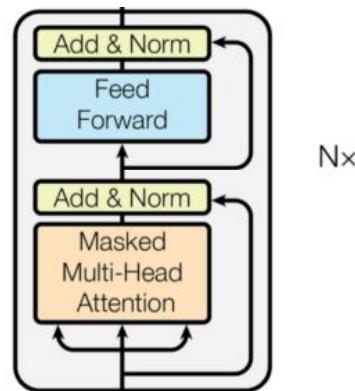
Characteristics

Decoder-only Transformer-based model



Characteristics

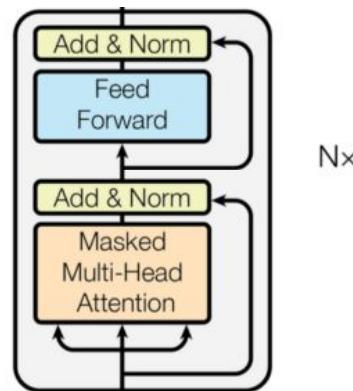
Decoder-only Transformer-based model



Nx

Characteristics

Decoder-only Transformer-based model



Nx

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



Transformers & Large Language Models

LLM overview

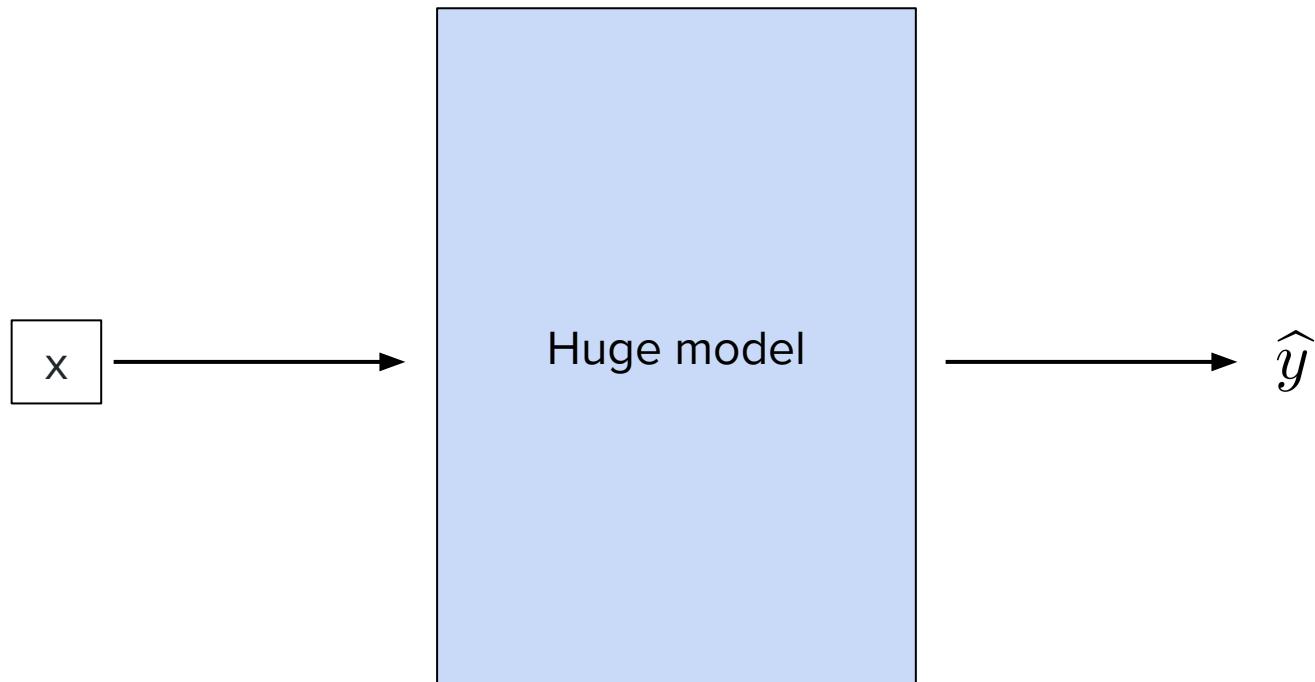
MoE-based LLMs

Response generation

Prompting strategies

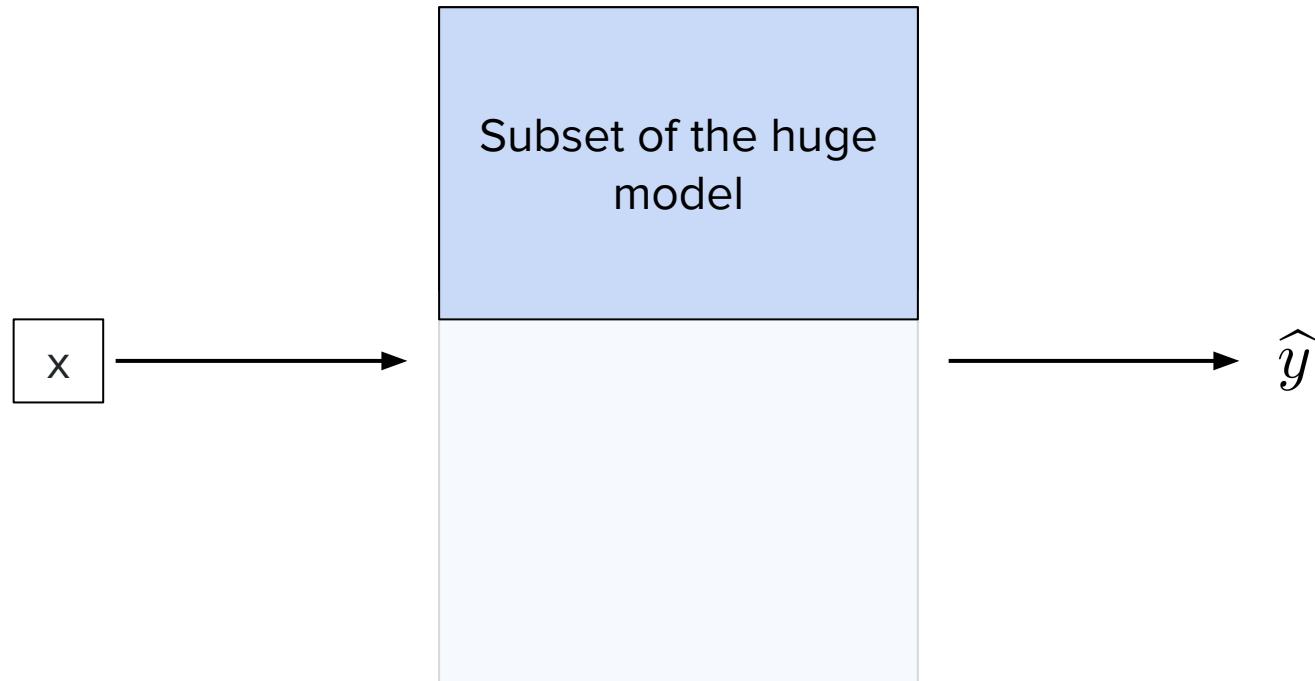
Inference optimizations

Motivation

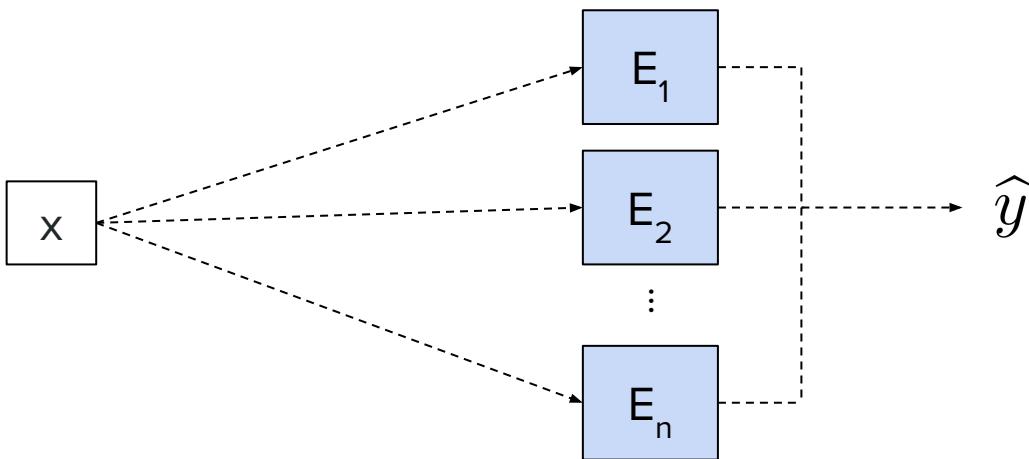


Motivation

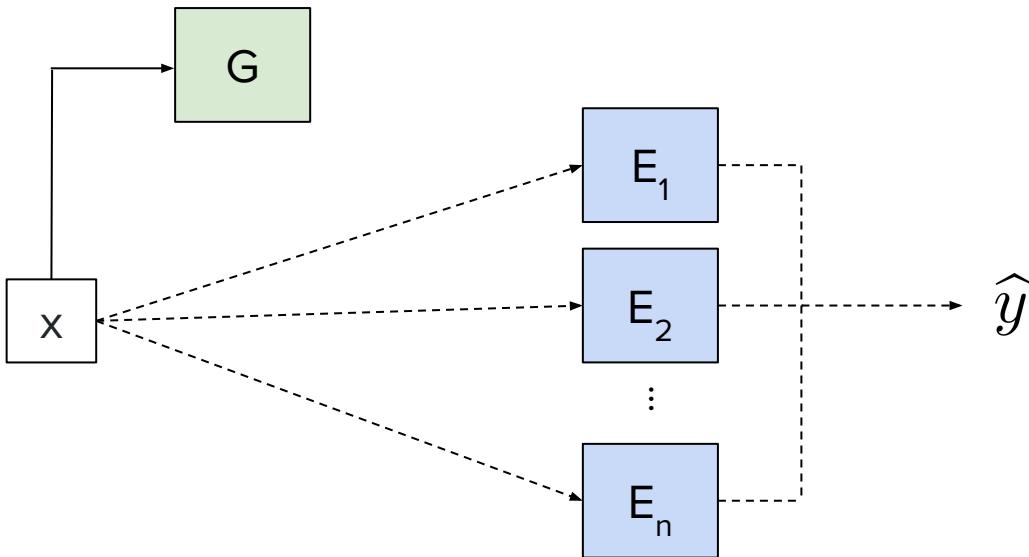
Idea. Not all weights are useful in the forward pass



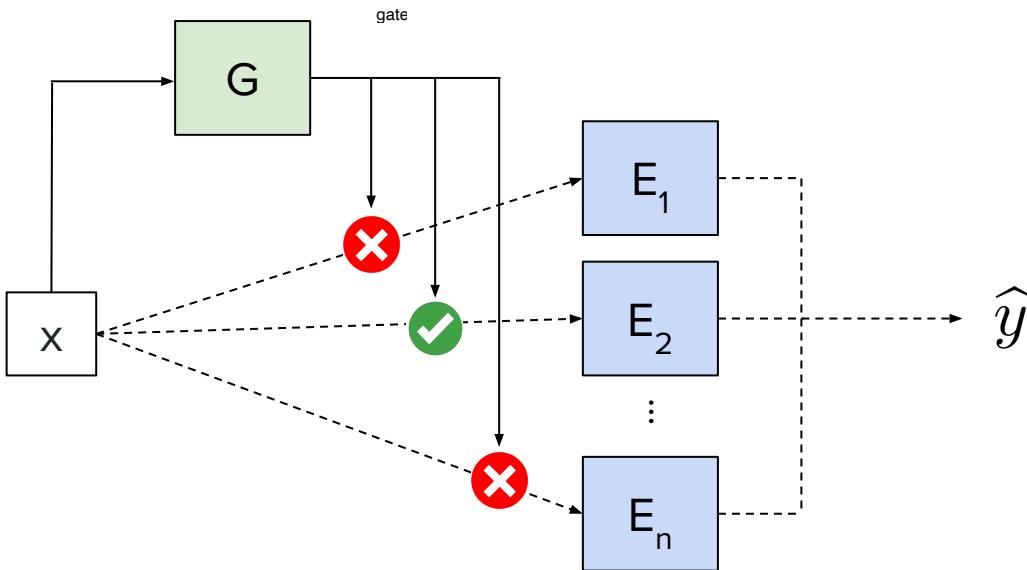
Motivation



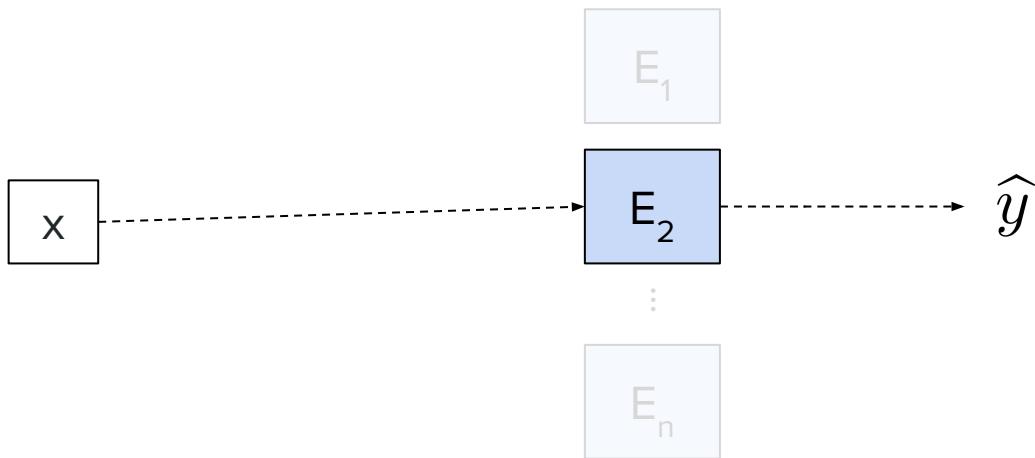
Motivation



Motivation

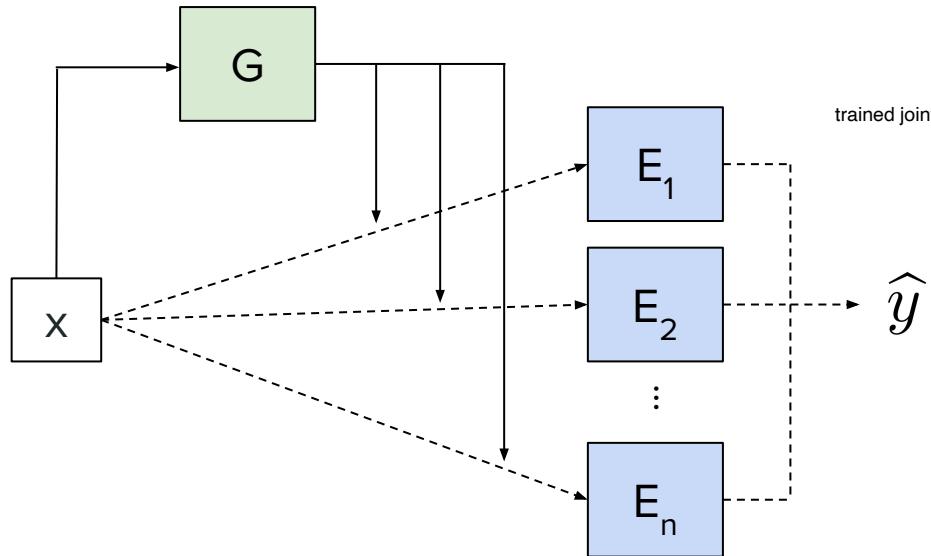


Motivation



Overview of MoEs

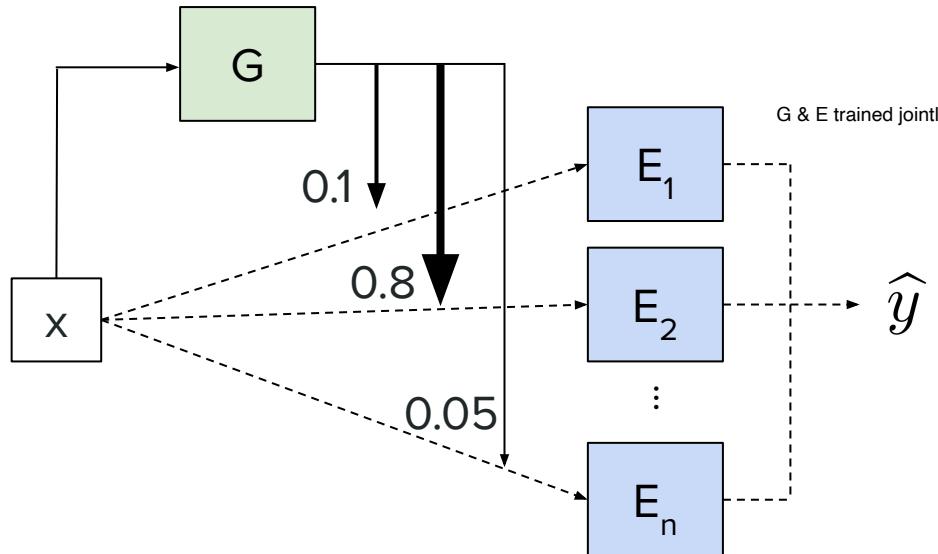
MoE = Mixture of Experts



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

Overview of MoEs

MoE = Mixture of Experts

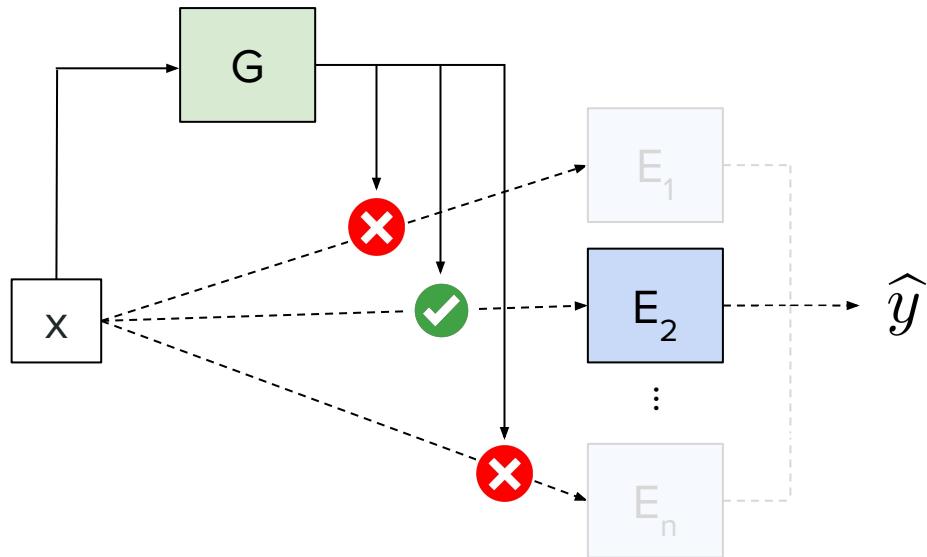


Dense MoE. Output is weighted average of **all** expert outputs.

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

Overview of MoEs

MoE = Mixture of Experts



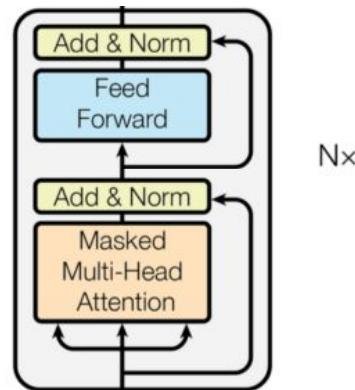
Sparse MoE. Output is weighted average of **selected** expert outputs.

lower number
of FLOPS

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

hyperparameter
Via **top-k** selection

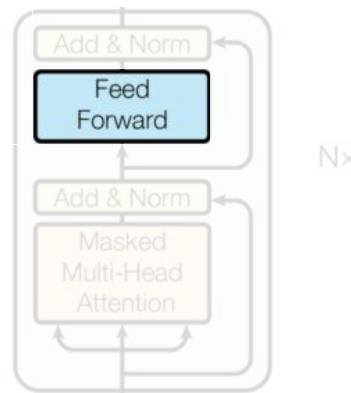
MoE in Transformer-based models



MoE in Transformer-based models

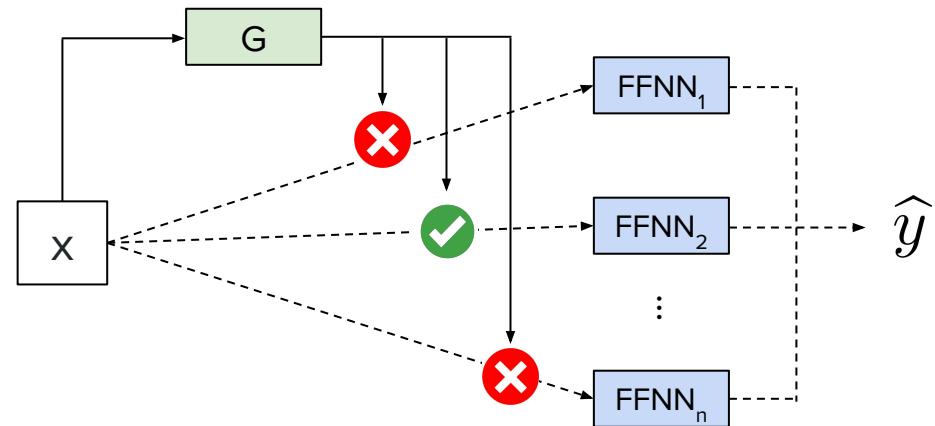
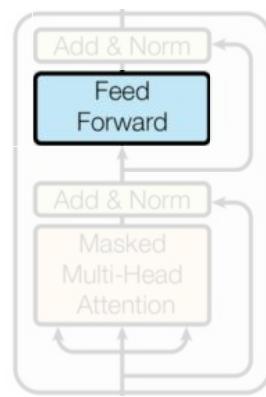
n_params in FFN:
 $d_{model} \times d_{\{FF\}} \times 2 + \text{bias}$

computationally heavy part of a transformer or a decoder.
Place the MoE there



MoE in Transformer-based models

Routing done **for each token!**



Training challenges include routing collapse

Symptom. Same expert gets selected most of the time.

"routing collapse"

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

these quantities help the
routing distribution to be
almost uniform

other methods:
Noisy gating

Interpreting experts

Layer 0

(experts are different in different layers.
before each layer of FFN we have a G)

```
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

notice uniform
distribution of
tokens for each
expert



Transformers & Large Language Models

LLM overview

MoE-based LLMs

Response generation

Prompting strategies

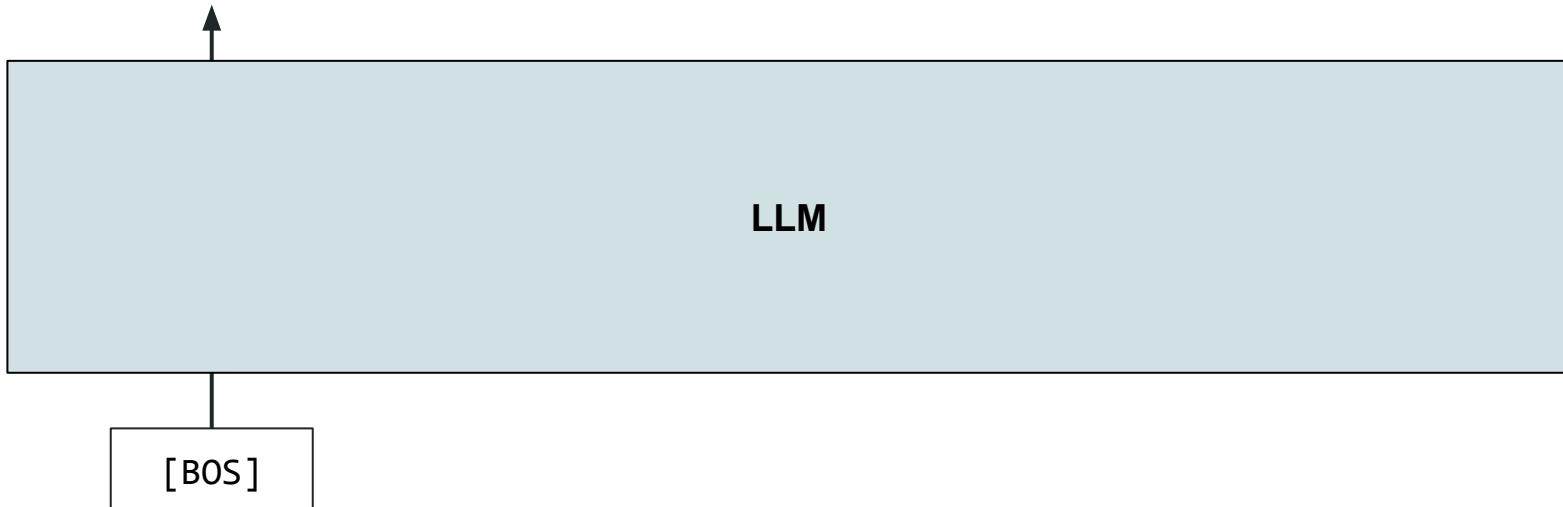
Inference optimizations

Next token prediction

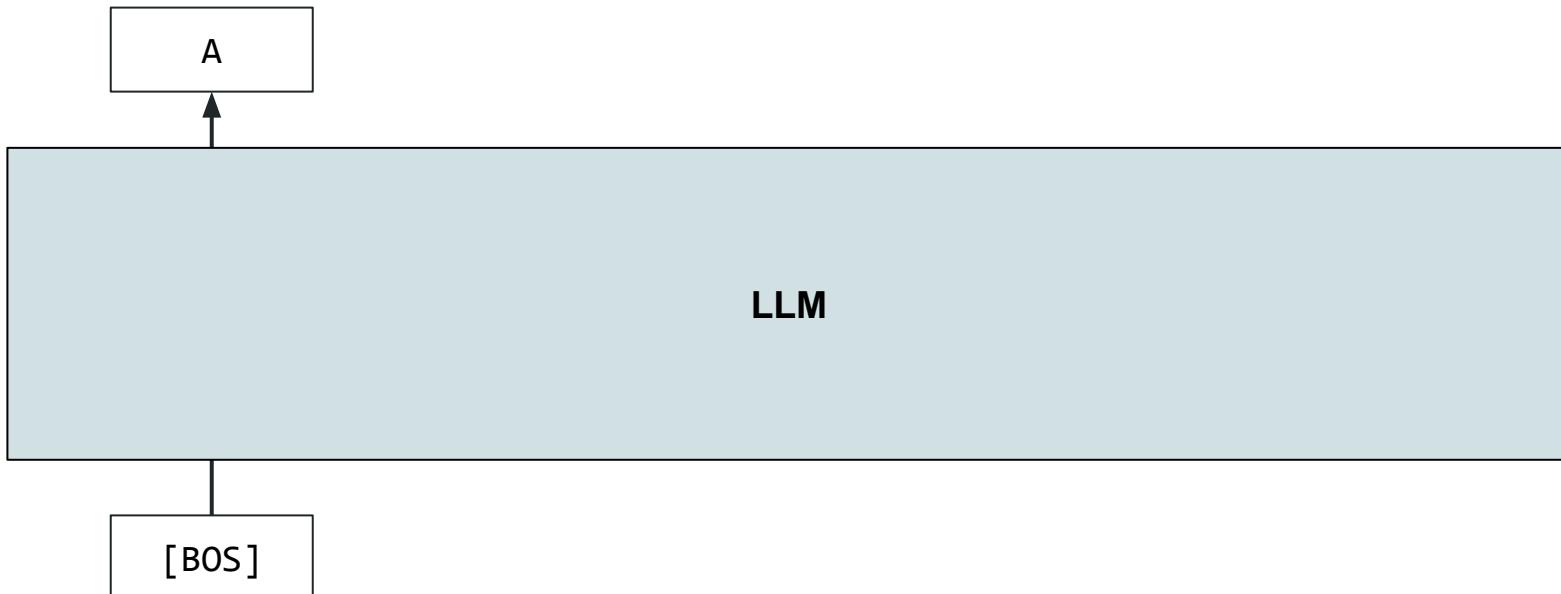
LLM

[BOS]

Next token prediction



Next token prediction



Next token prediction

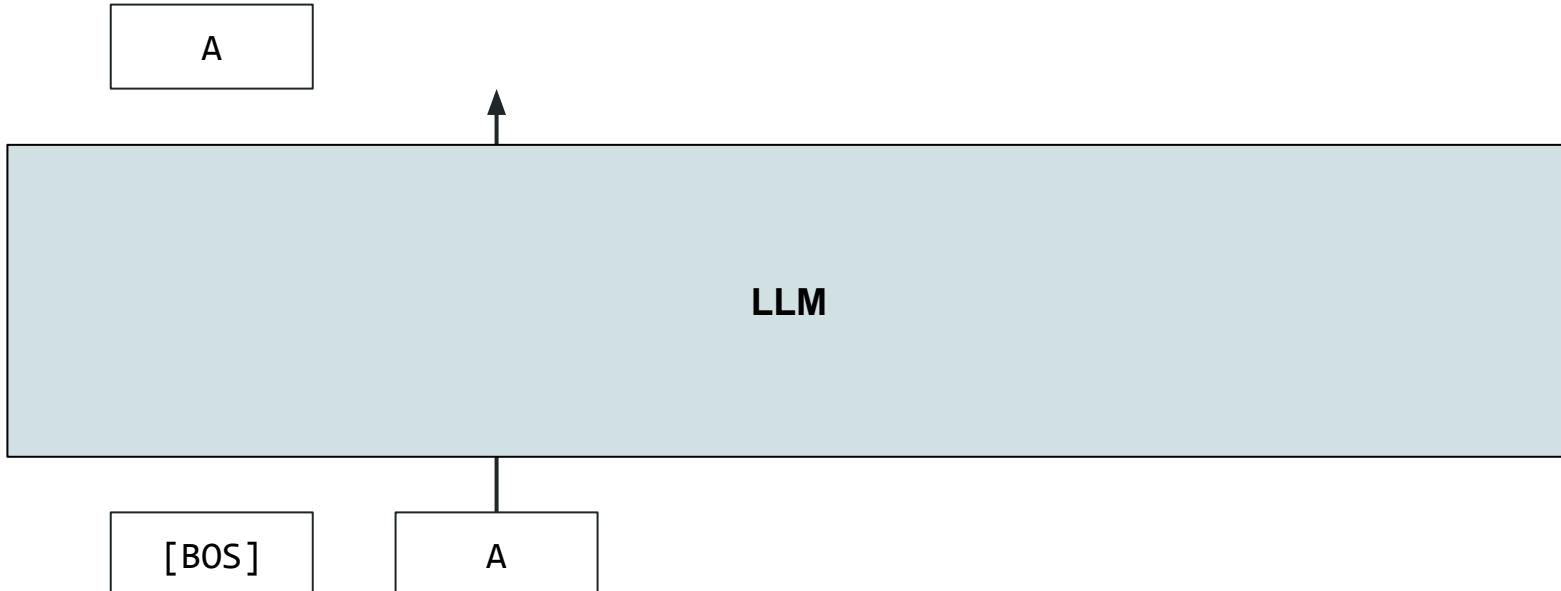
A

LLM

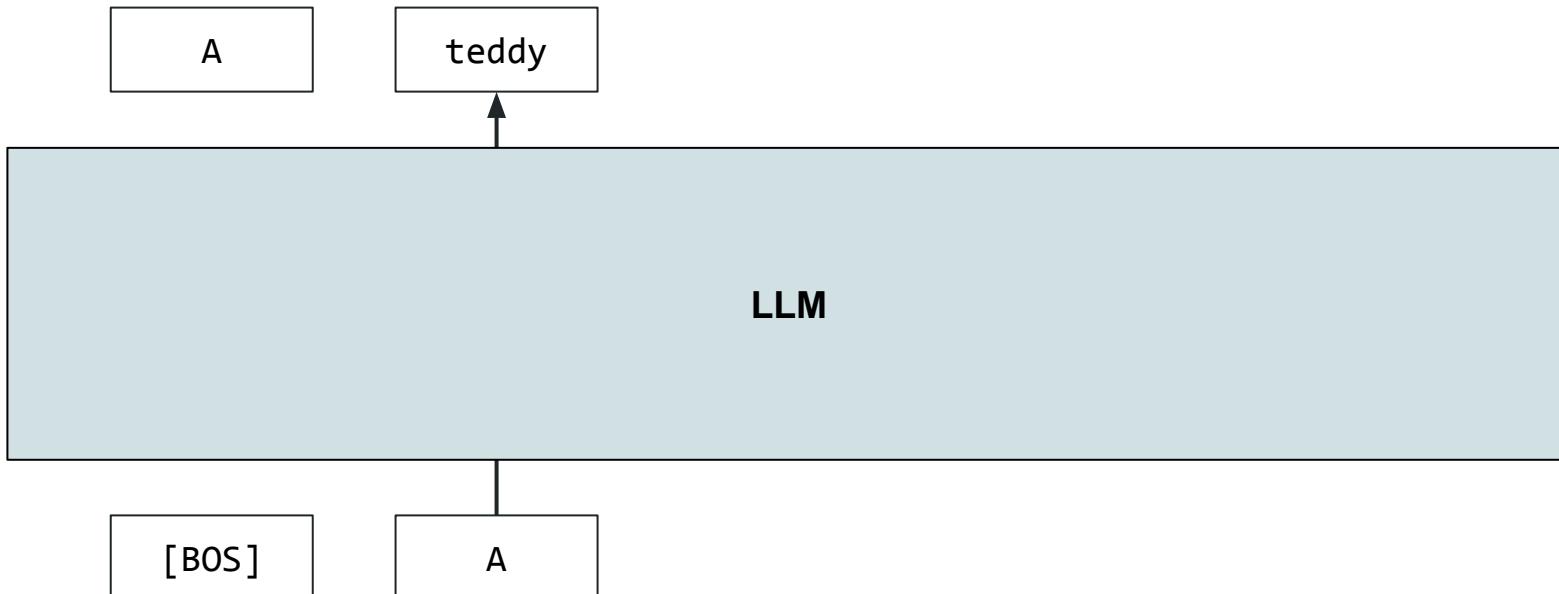
[BOS]

A

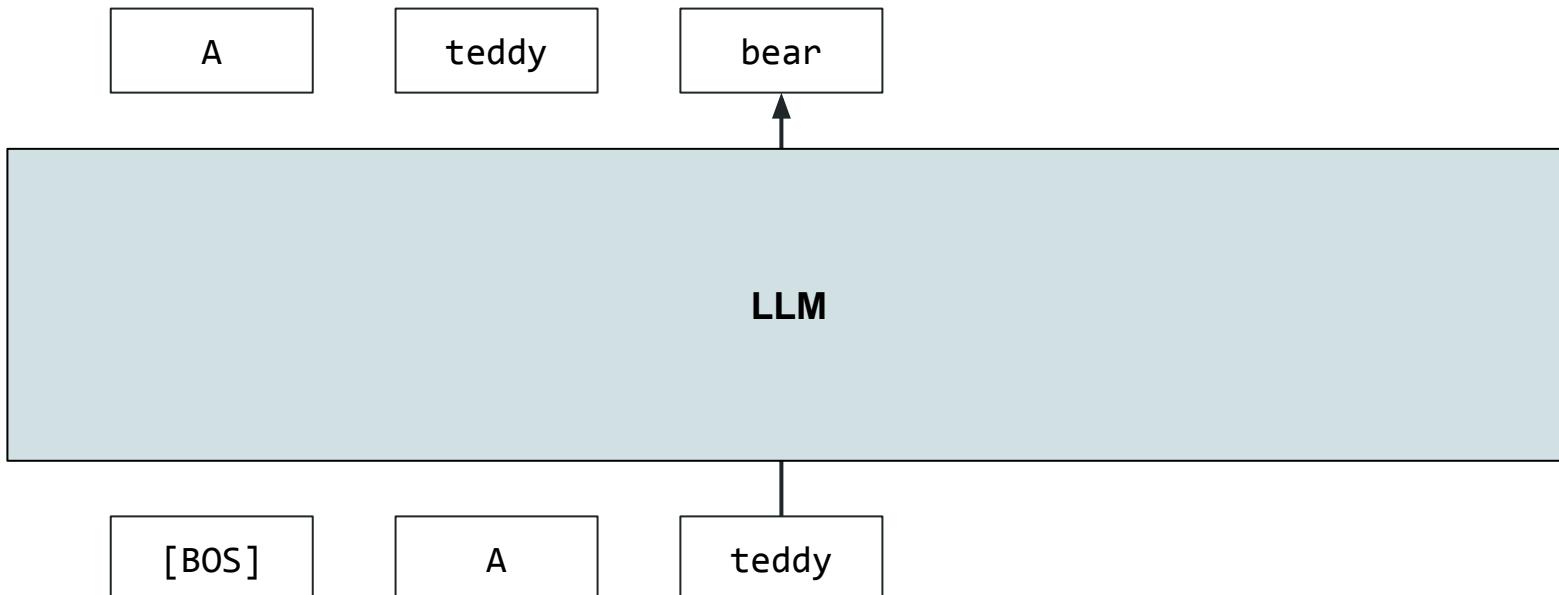
Next token prediction



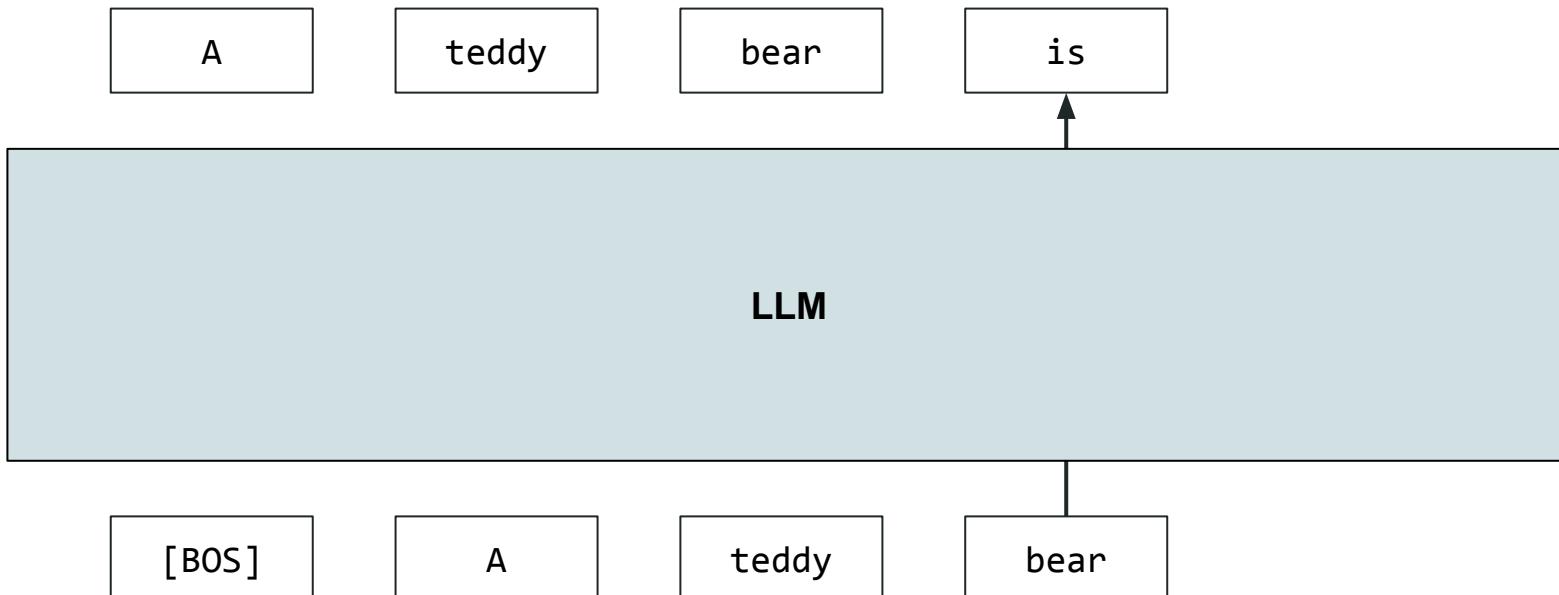
Next token prediction



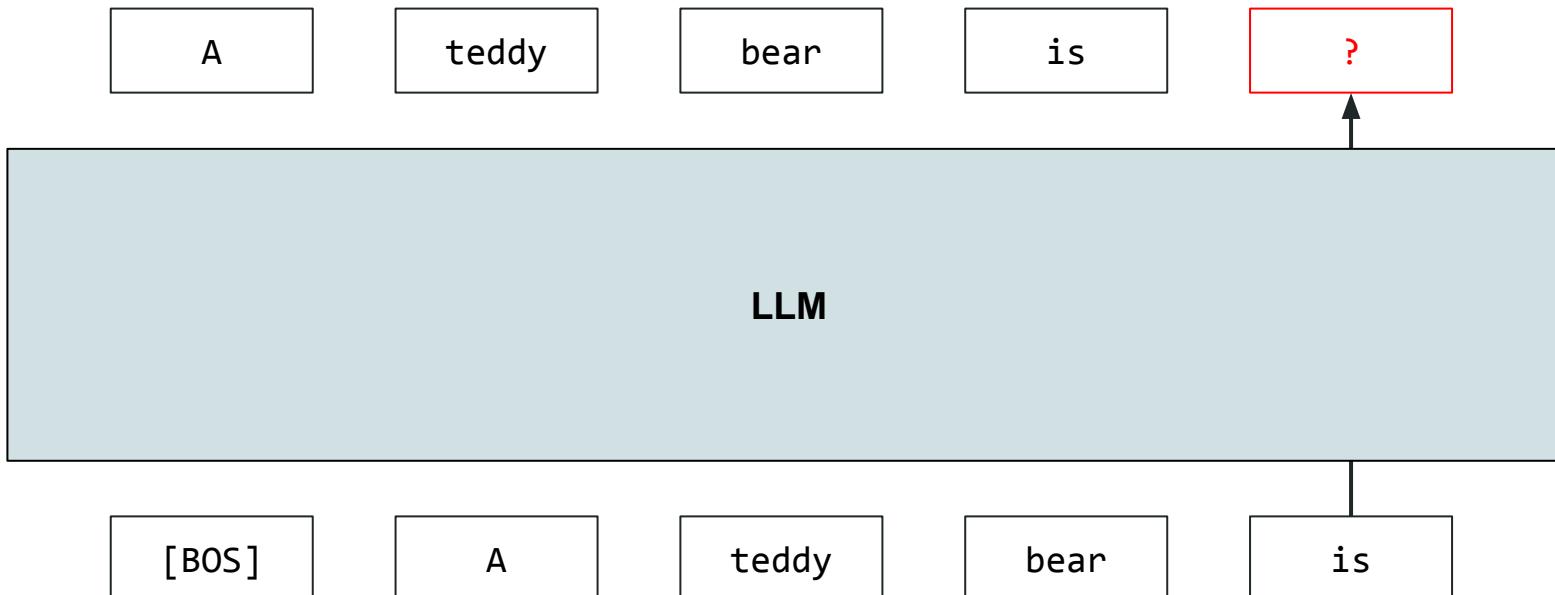
Next token prediction



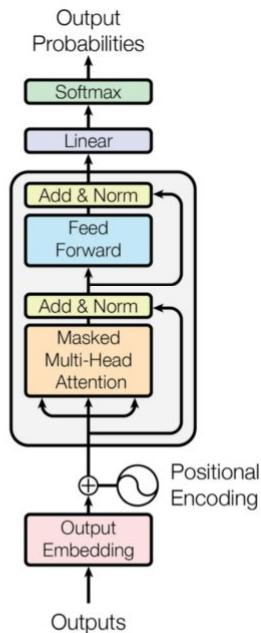
Next token prediction



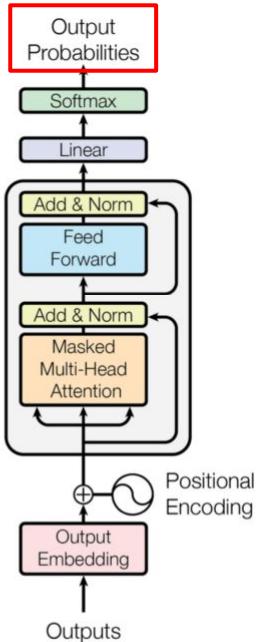
Next token prediction



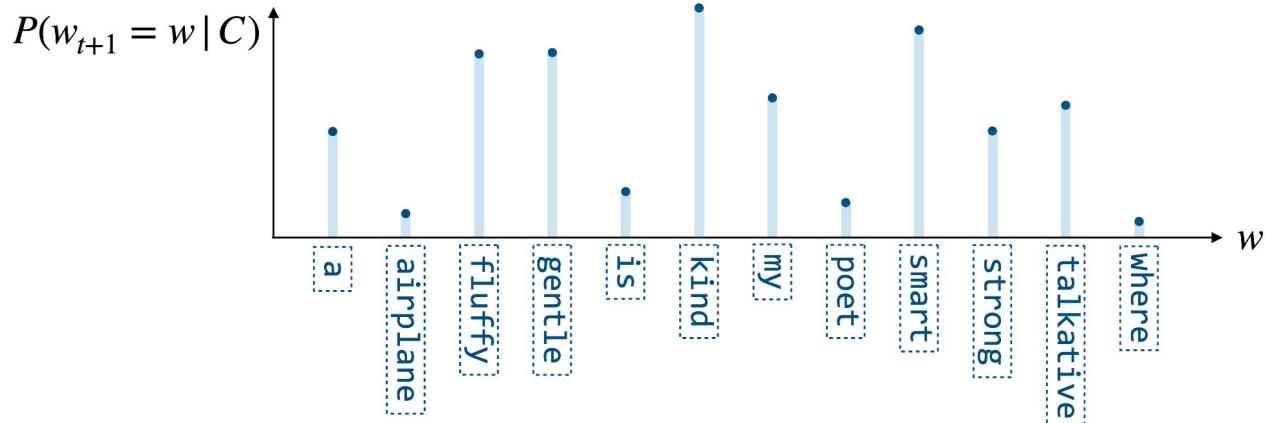
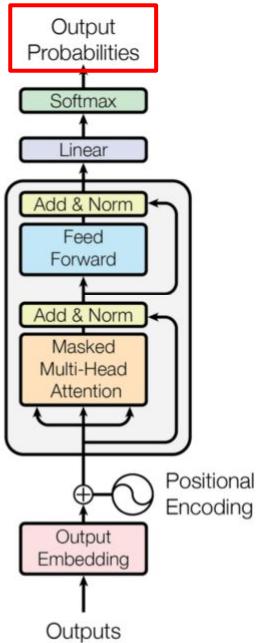
Predicting next token



Predicting next token

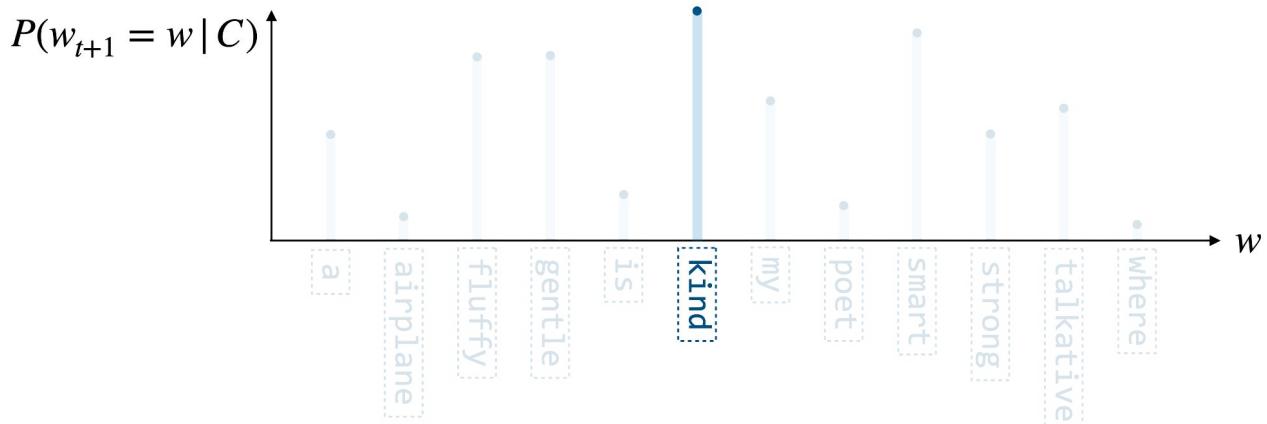
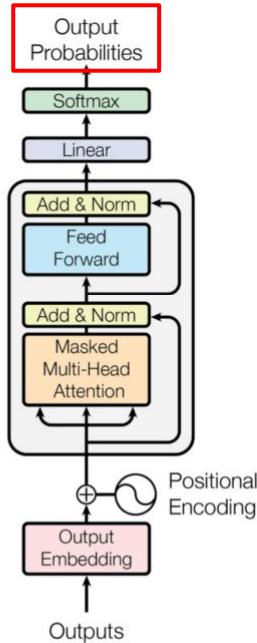


Predicting next token



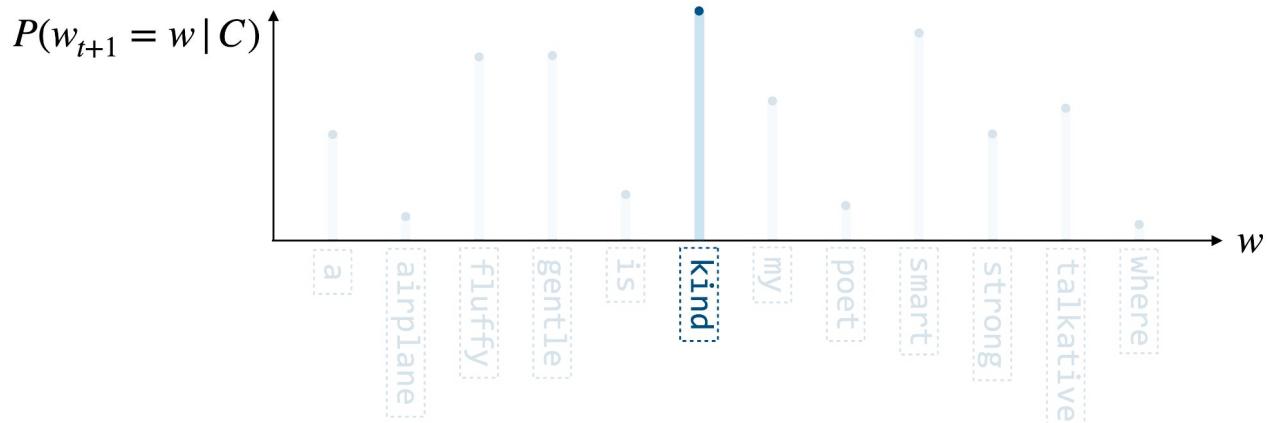
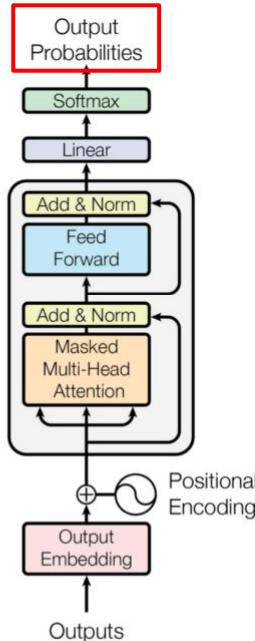
Predicting next token with greedy decoding

1st idea. Take token with highest predicted probability



Predicting next token with greedy decoding

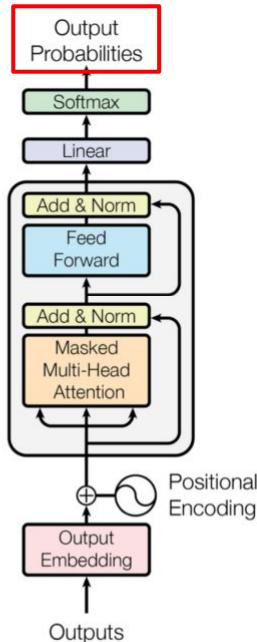
1st idea. Take token with highest predicted probability



Limitations. Output not optimal, natural and/or diverse

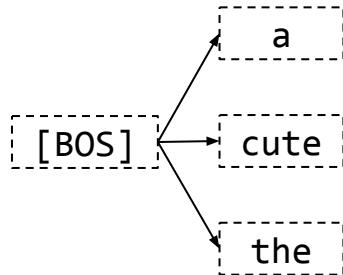
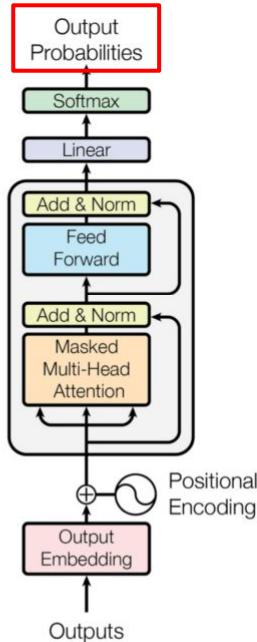
Predicting next token with beam search

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



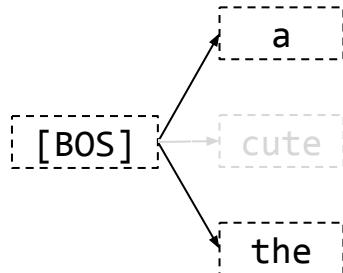
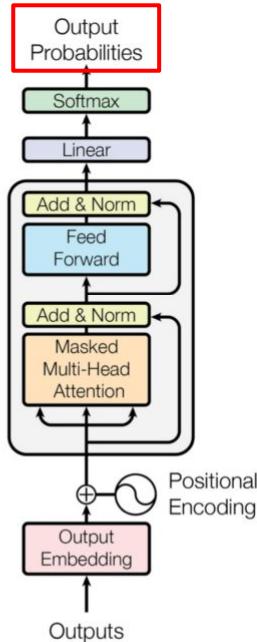
Predicting next token with beam search

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

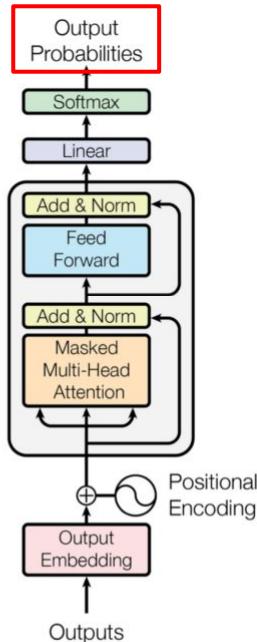


Predicting next token with beam search

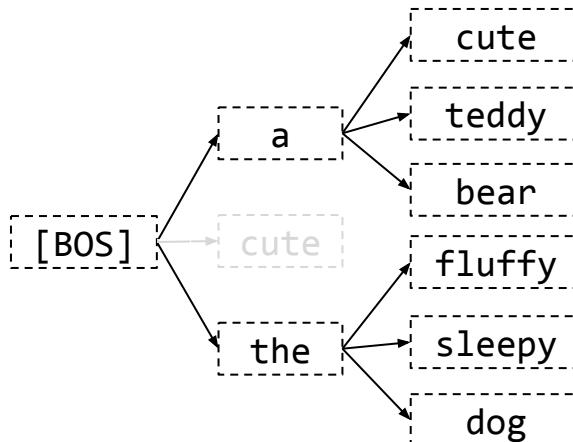
2nd idea. Keep k paths that are the most likely



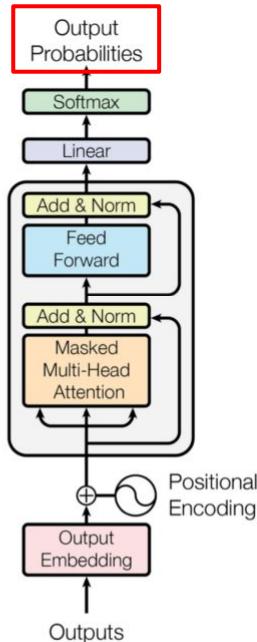
Predicting next token with beam search



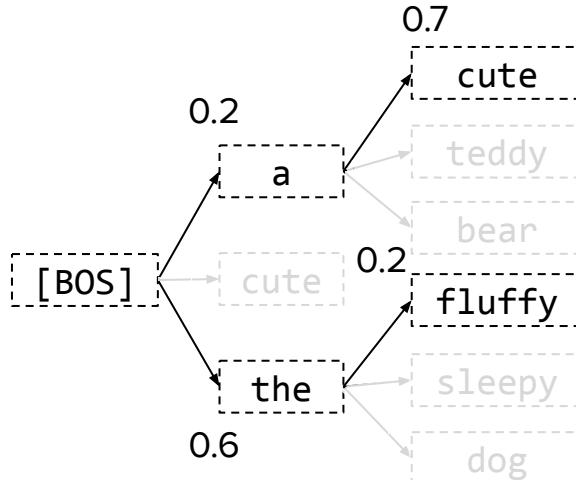
2nd idea. Keep k paths that are the most likely



Predicting next token with beam search

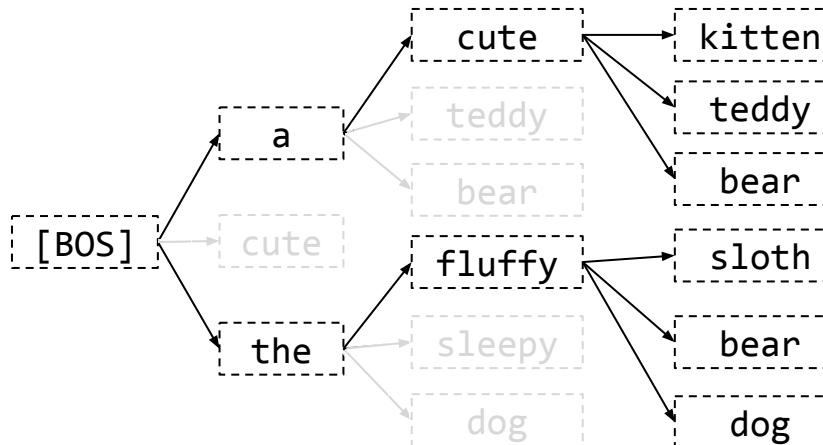
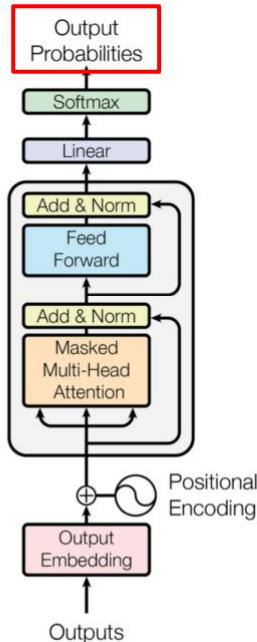


2nd idea. Keep k paths that are the most likely



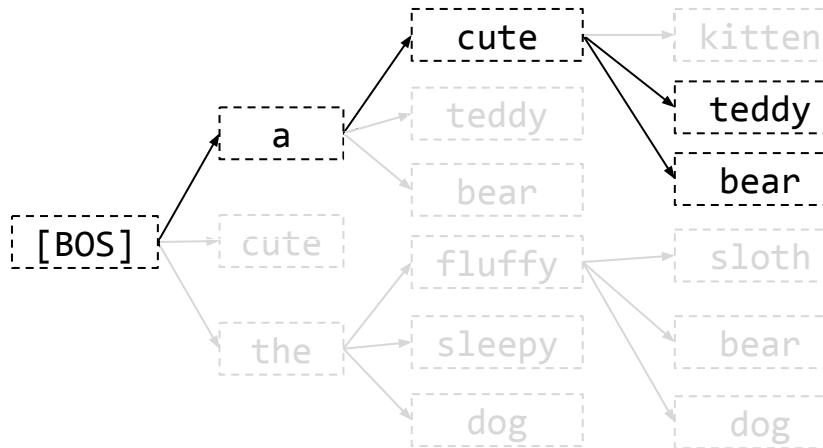
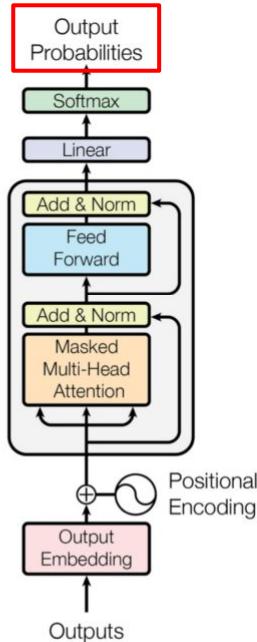
Predicting next token with beam search

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

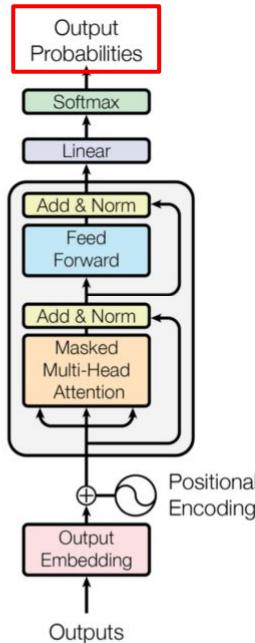


Predicting next token with beam search

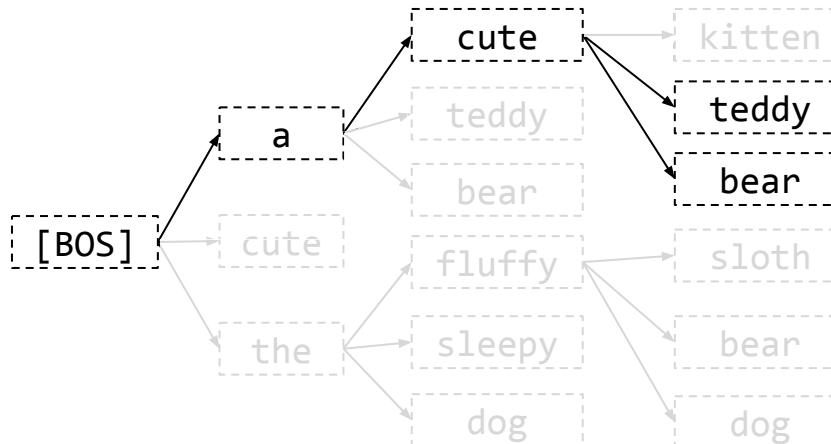
2nd idea. Keep k paths that are the most likely



Predicting next token with beam search

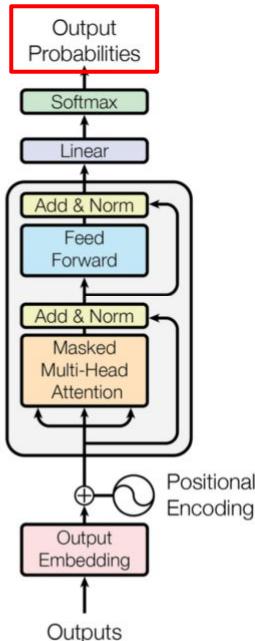


2nd idea. Keep k paths that are the most likely



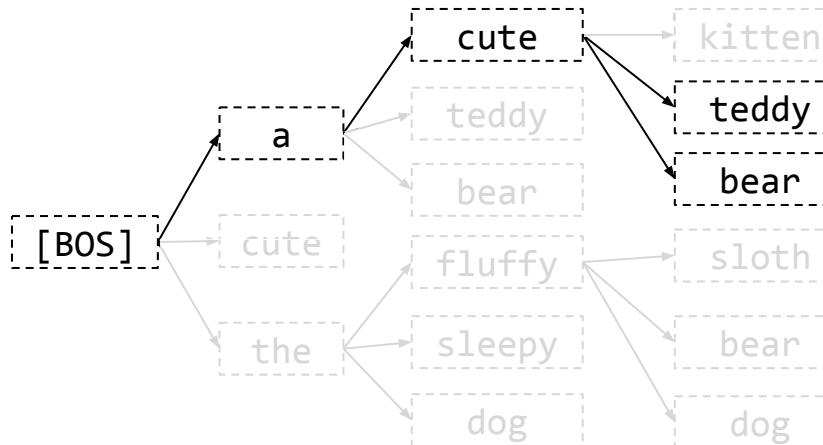
Ends at [EOS]

Predicting next token with beam search



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

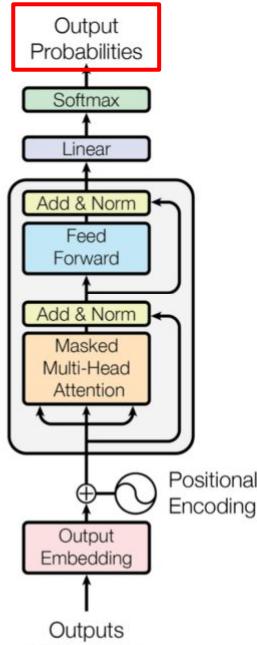
note: these methods are being discussed for response generation



Limitations. Needs computations + lacks diversity/creativity

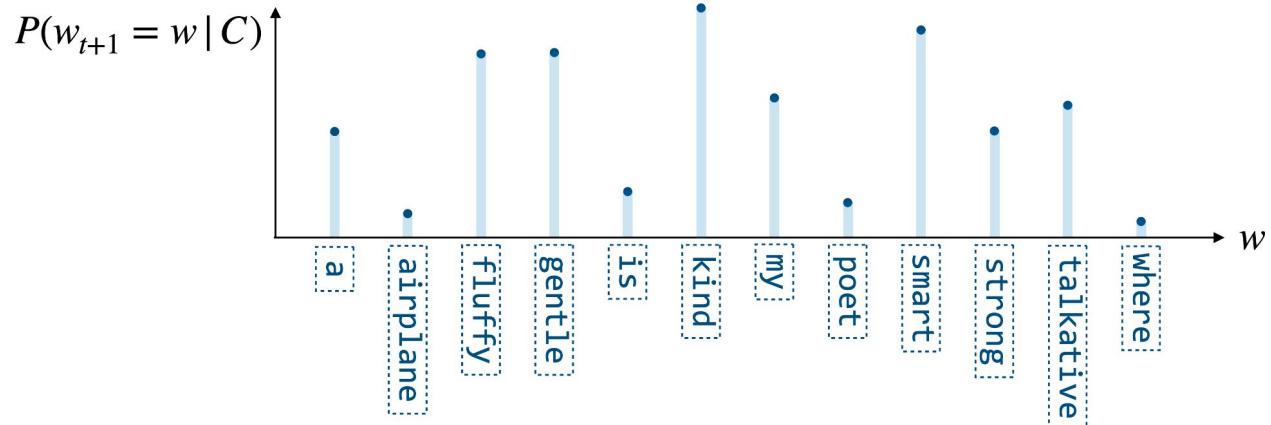
used in the case of translation
where we want answer to be closest to correct

Predicting next token with sampling



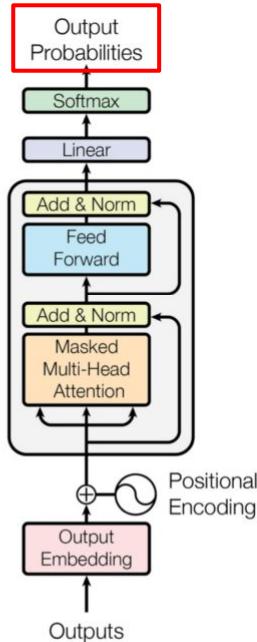
3rd idea. Sample next token from probability distribution:

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

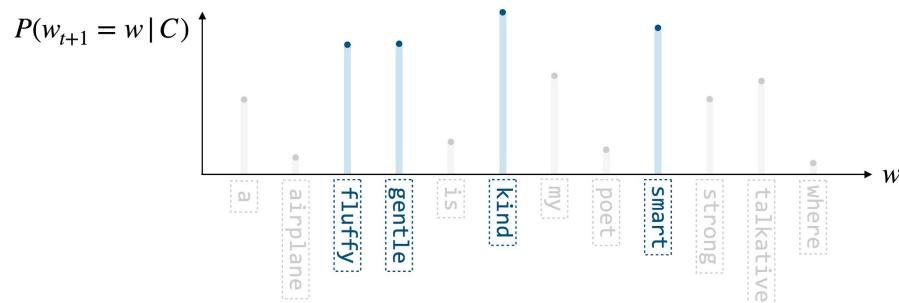


Sampling strategies

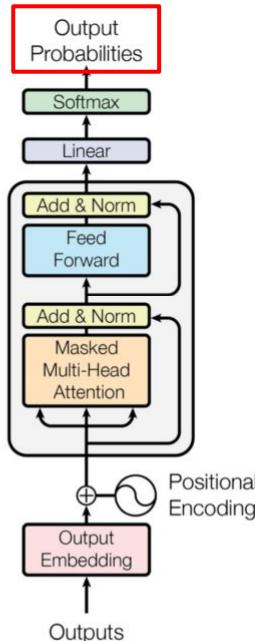
- **Top-k:** Sample among top k most probable tokens



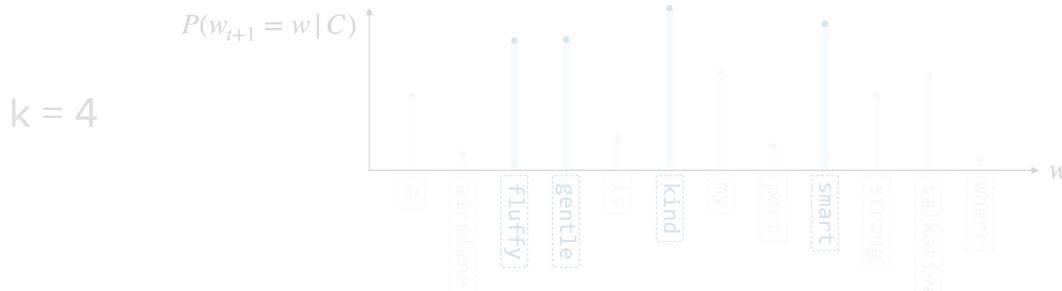
$k = 4$



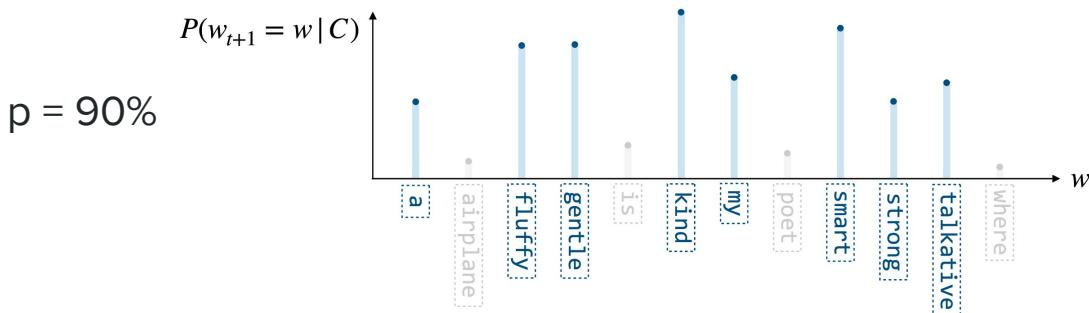
Sampling strategies



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



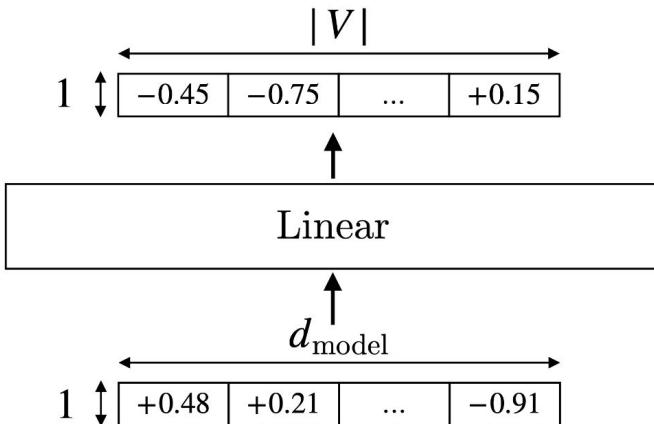
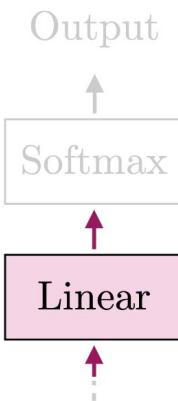
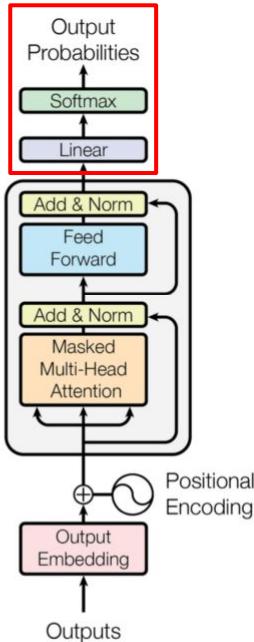
- **Top-p:** Random sample among smallest set of tokens with cumulative probability $\geq p$



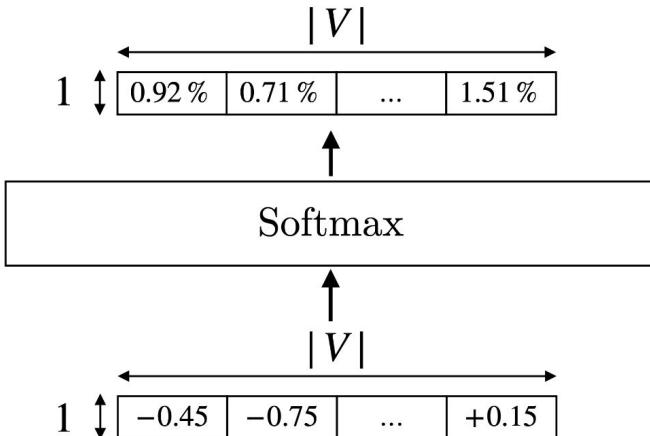
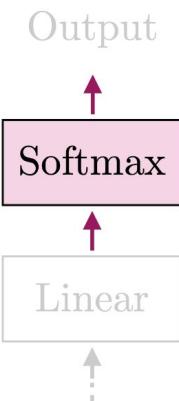
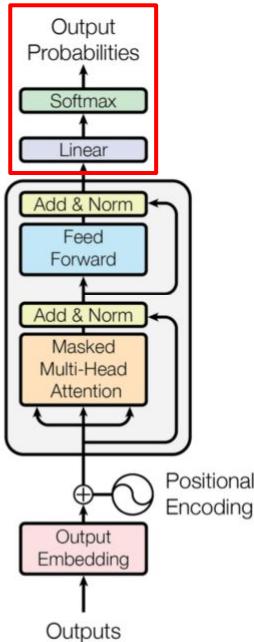
Next token prediction

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

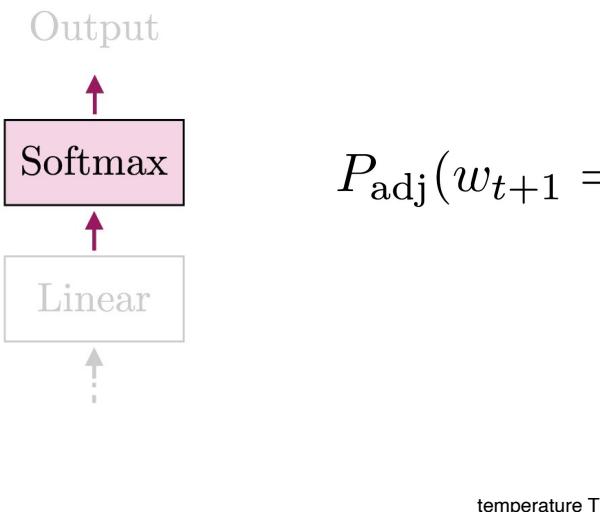
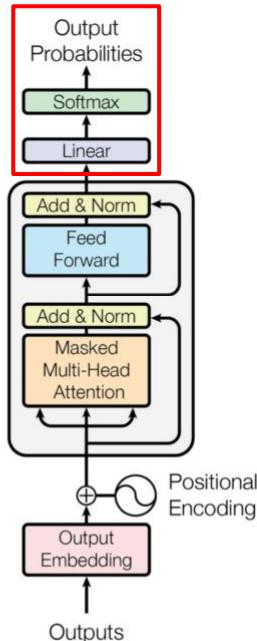
Probability computation



Probability computation



Temperature allows to tweak output probabilities



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

Impact of temperature on probabilities

Small T

?

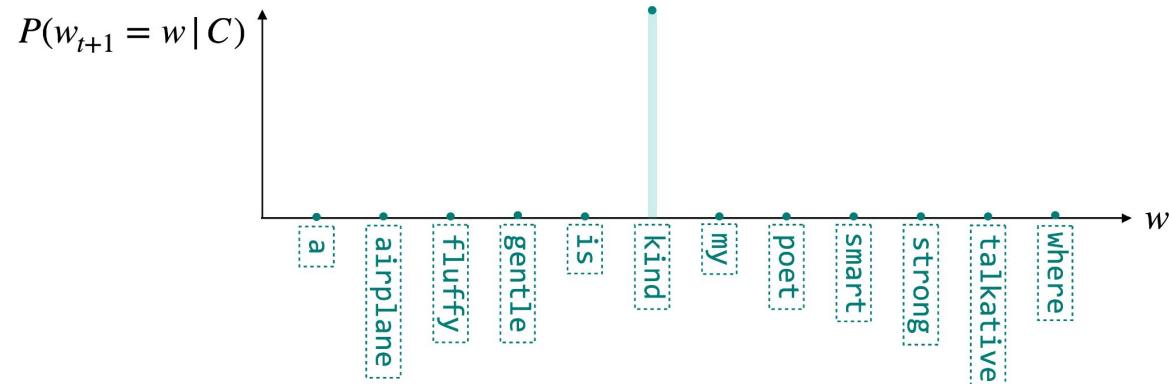
High T

?

Impact of temperature on probabilities

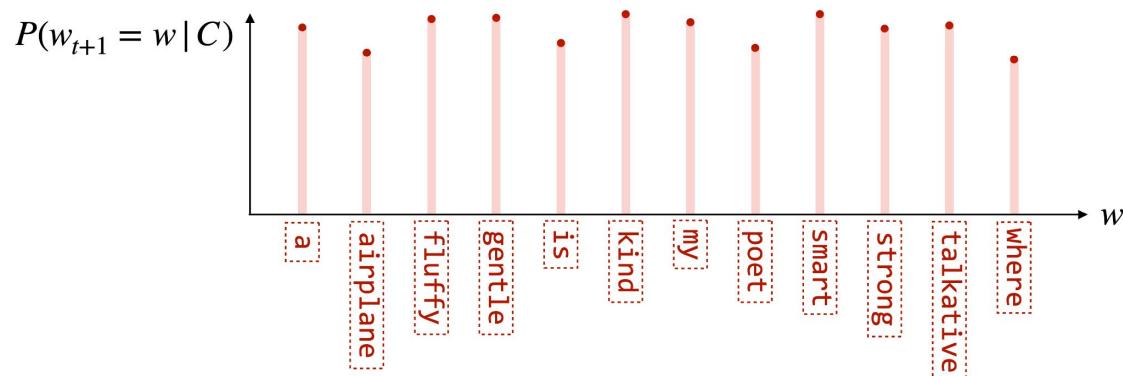
Small T

sampling is the only non-deterministic part
of the transformer
some other things happen during
computation that give non-determinism



High T

peaked. divide by highest logit exp
to get the proof of limit.
high temperature more creative
results



"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"  
}
```

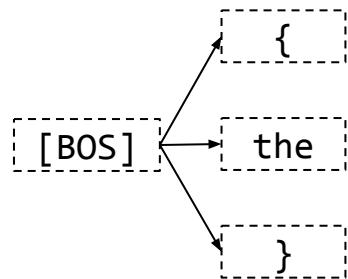
Constraining the output via guided decoding

Idea. Only allow "valid" next tokens

[BOS]

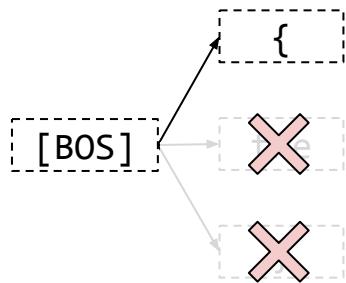
Constraining the output via guided decoding

Idea. Only allow "valid" next tokens



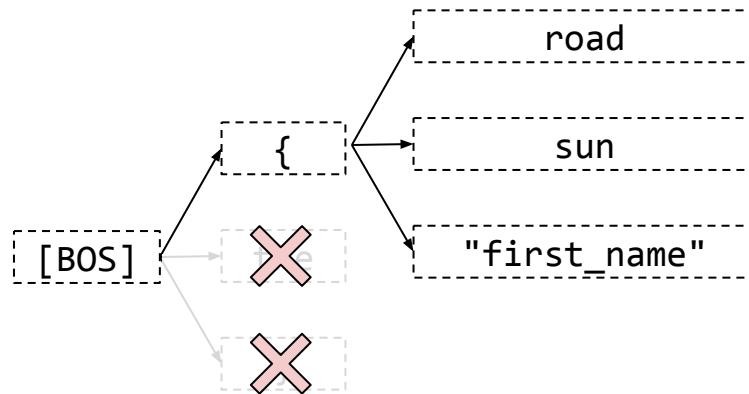
Constraining the output via guided decoding

Idea. Only allow "valid" next tokens



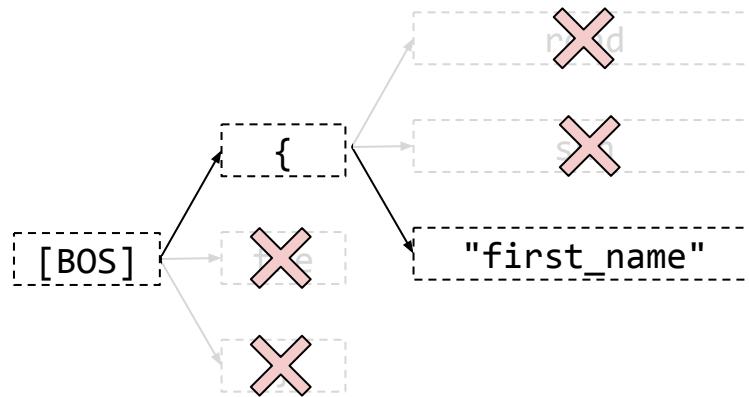
Constraining the output via guided decoding

Idea. Only allow "valid" next tokens



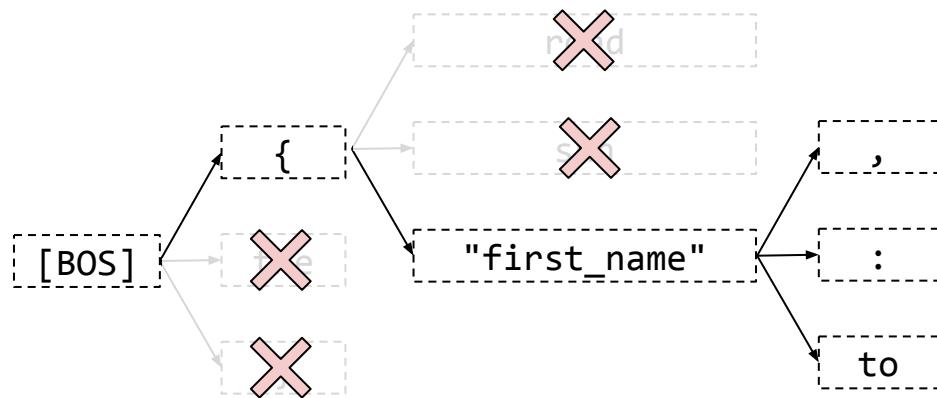
Constraining the output via guided decoding

Idea. Only allow "valid" next tokens



Constraining the output via guided decoding

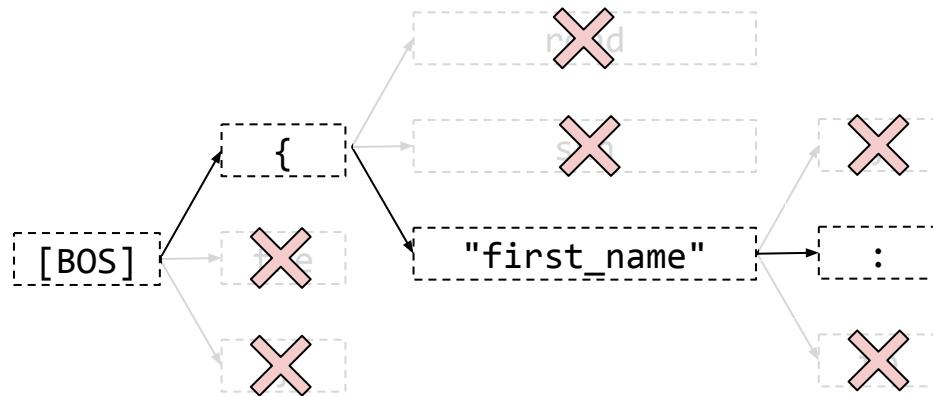
Idea. Only allow "valid" next tokens



Constraining the output via guided decoding

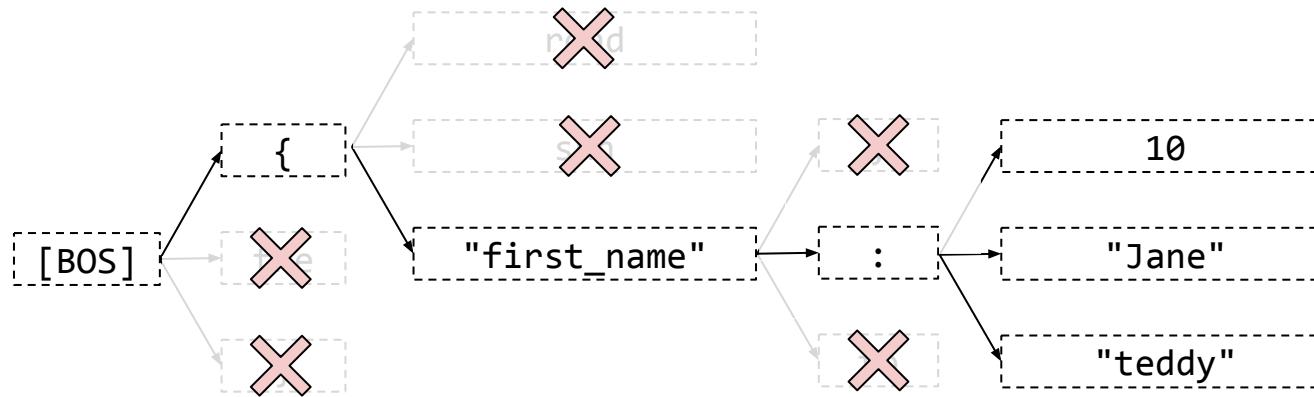
Idea. Only allow "valid" next tokens

use FSM and context free grammar etc.



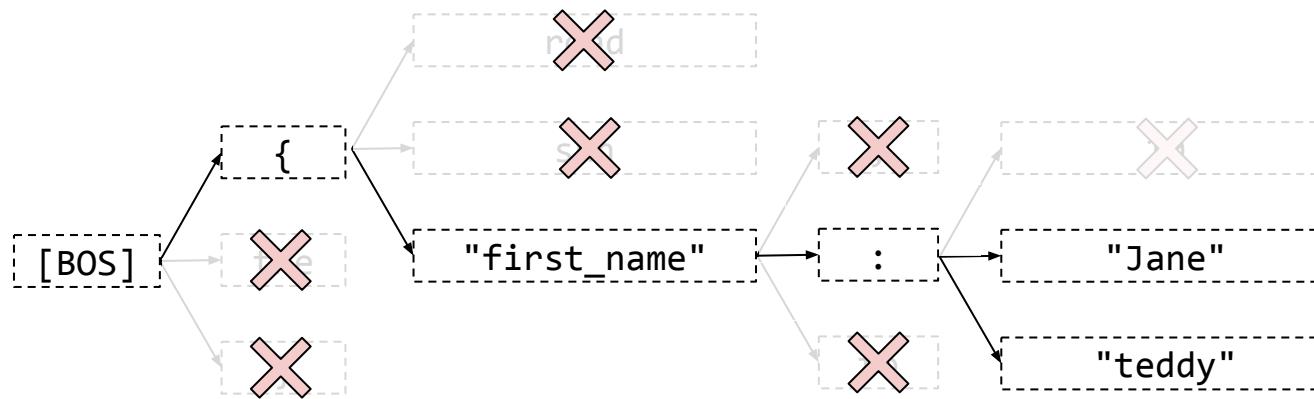
Constraining the output via guided decoding

Idea. Only allow "valid" next tokens



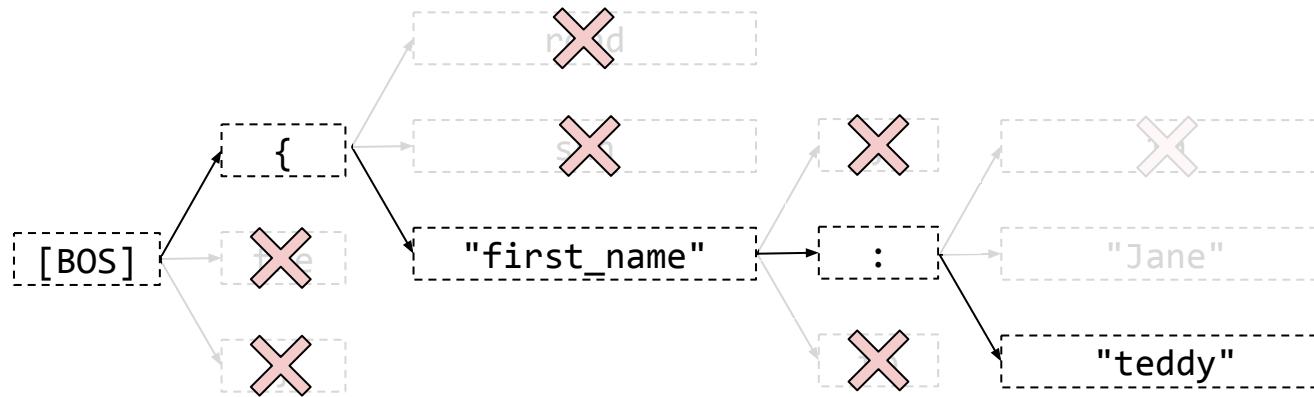
Constraining the output via guided decoding

Idea. Only allow "valid" next tokens



Constraining the output via guided decoding

Idea. Only allow "valid" next tokens





Transformers & Large Language Models

LLM overview

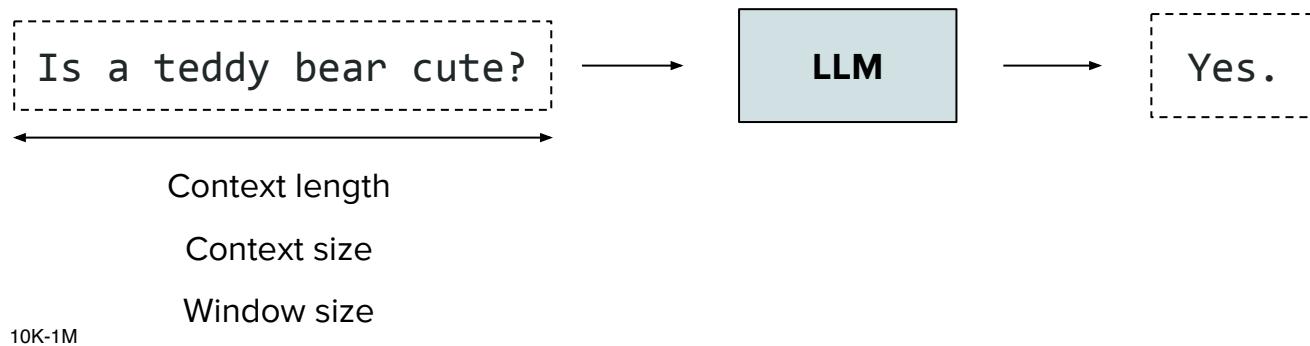
MoE-based LLMs

Response generation

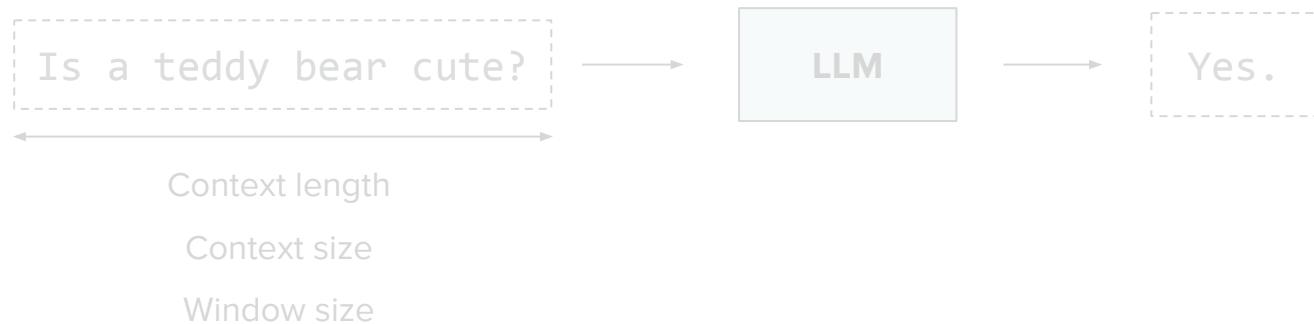
Prompting strategies

Inference optimizations

Terminology



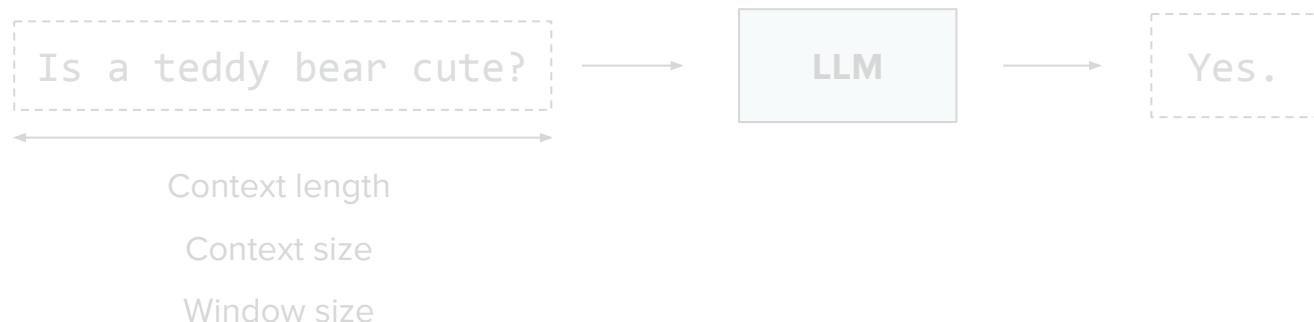
Terminology



Discussion. Orders of magnitude by:

- Input type
- Models

Terminology

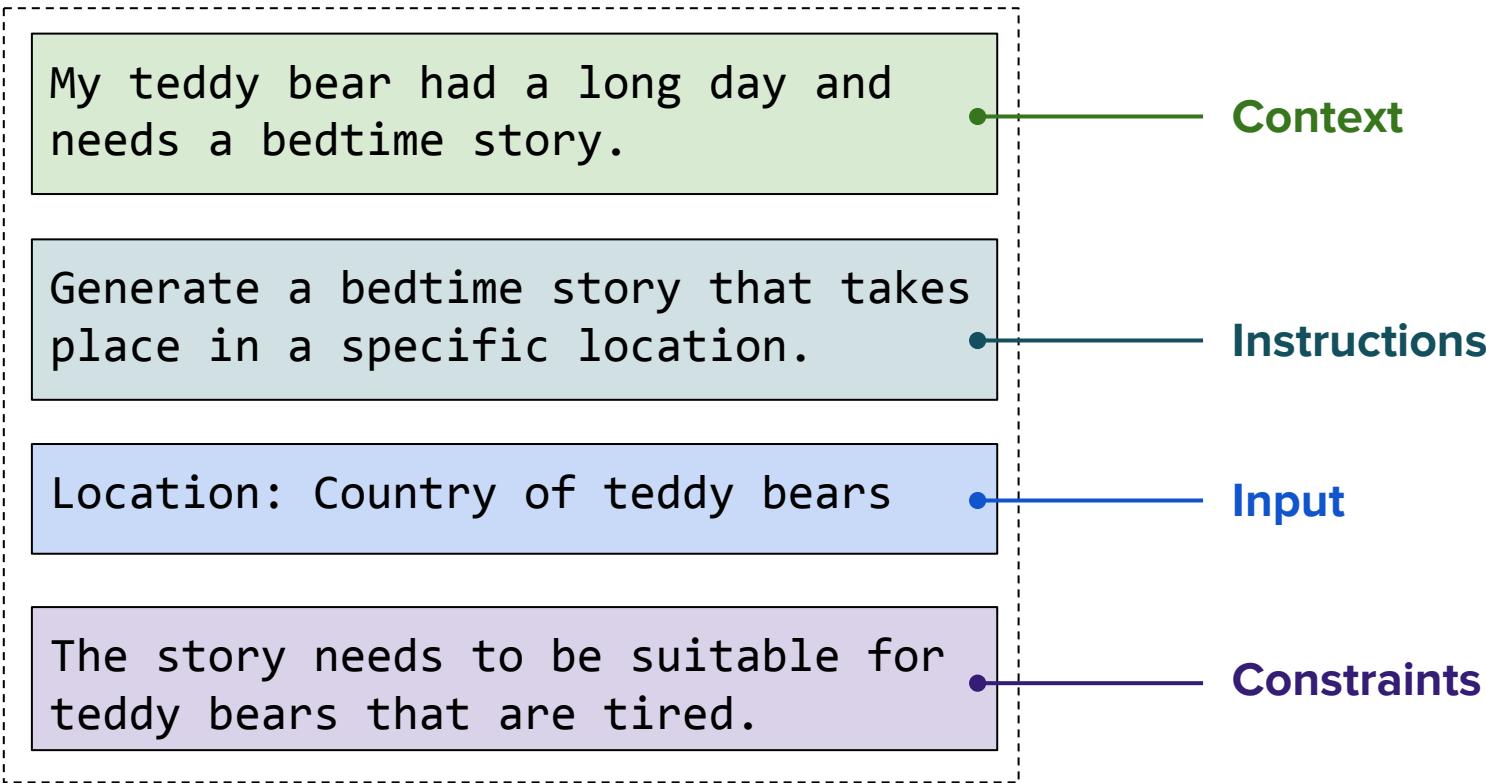


Discussion. Orders of magnitude by:

- Input type
- Models

← Beware of "context rot"!

Main structure



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

In-context learning

ICL = In-Context Learning

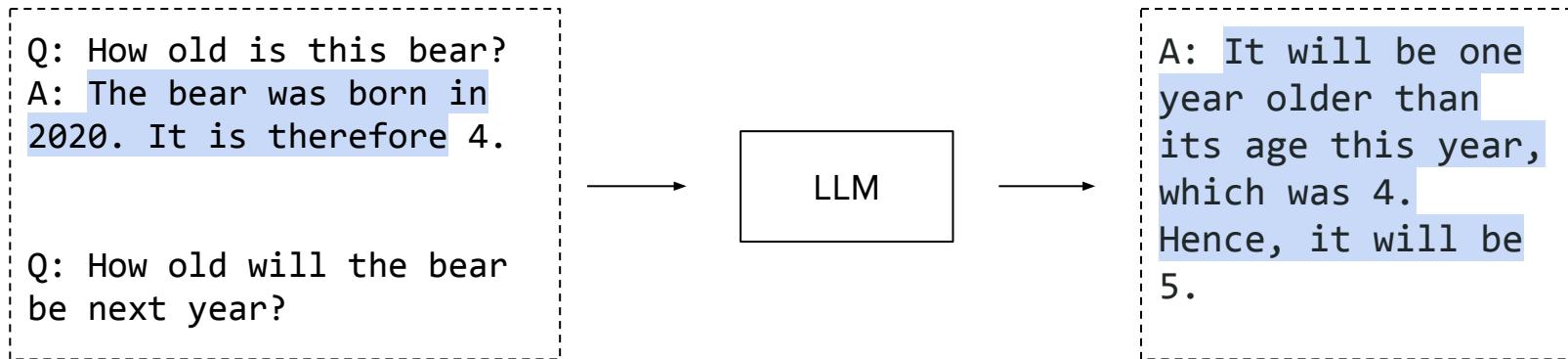
| Zero-shot learning | Few-shot learning |
|--|---|
| <ul style="list-style-type: none">• Question is asked without examples• Performance heavily depends on performance of initial model | <ul style="list-style-type: none">• Prompt contains examples of input / output• Typically has a better performance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

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

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

Chain of thought

Idea. Explaining reasoning helps in improving performance.

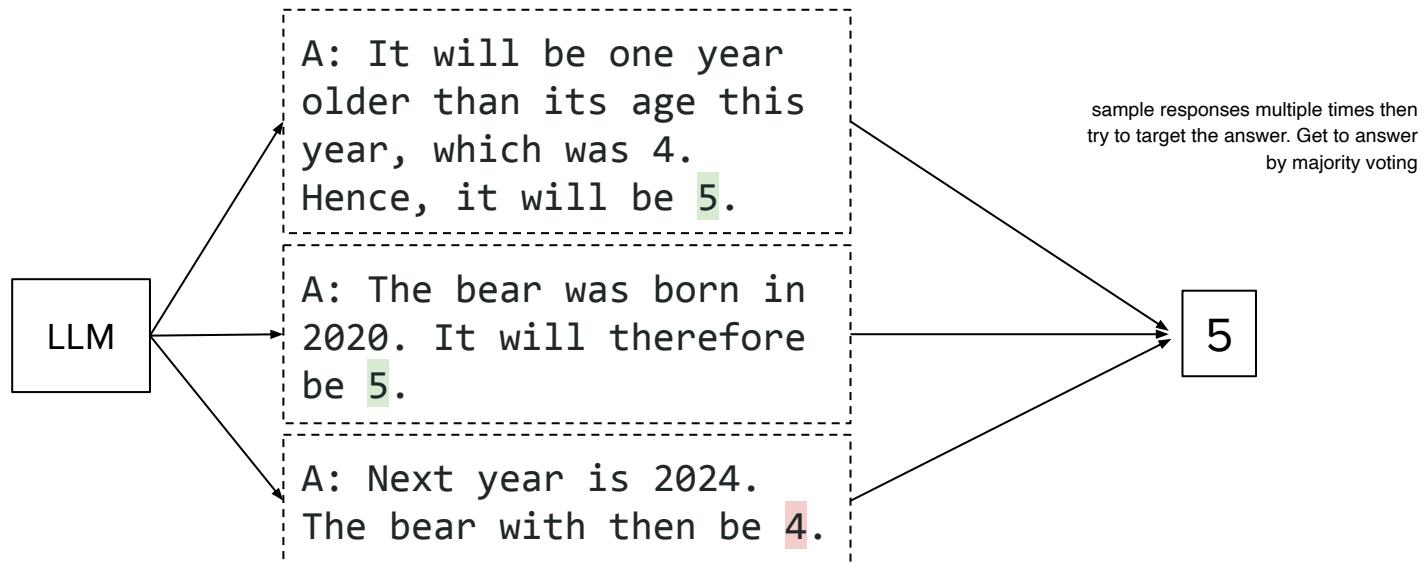


Discussion.

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

Self-consistency

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:

- Avoid redundancies
- Memory management
- Reformulate the math

Challenges

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

Categories. Many dimensions to optimize for.

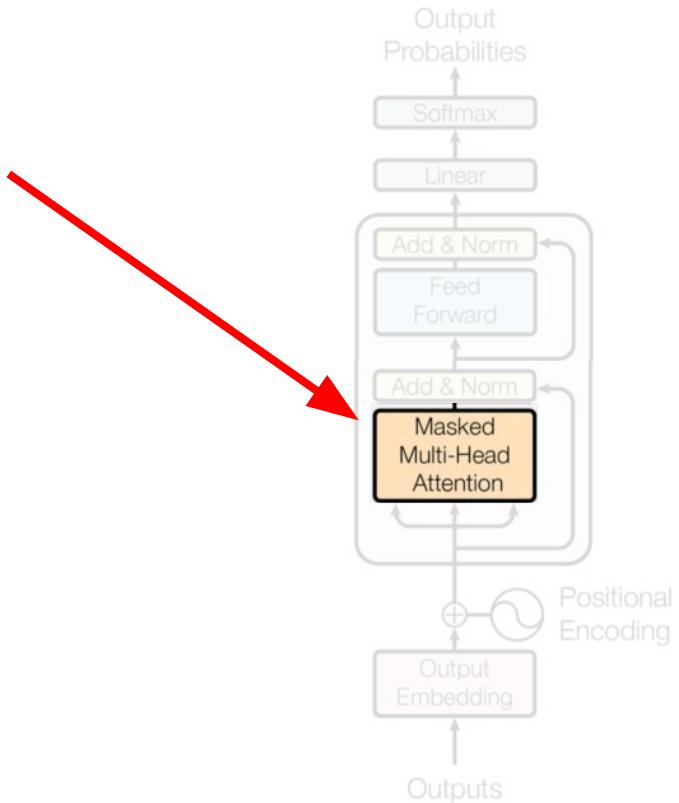
"Exact" efficiency:

- Avoid redundancies
- Memory management
- Reformulate the math

Approximations:

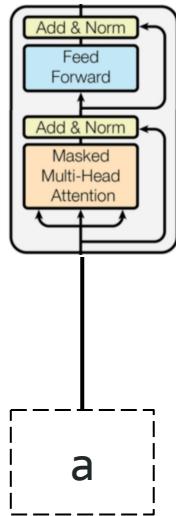
- Architectural changes
- Embeddings representations
- Token prediction

Attention-based tricks



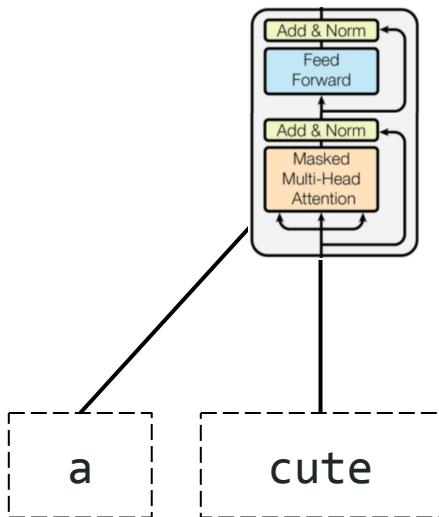
Reusing computations with KV caching

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



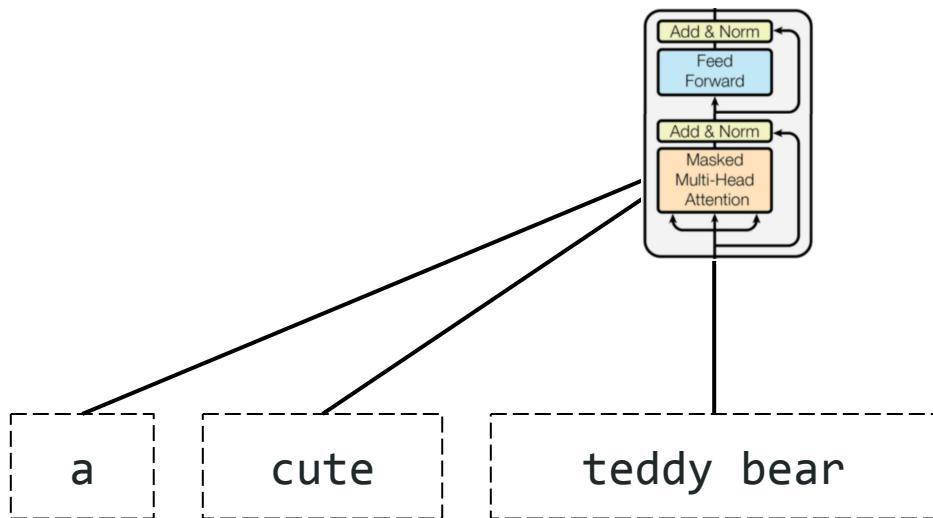
Reusing computations with KV caching

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



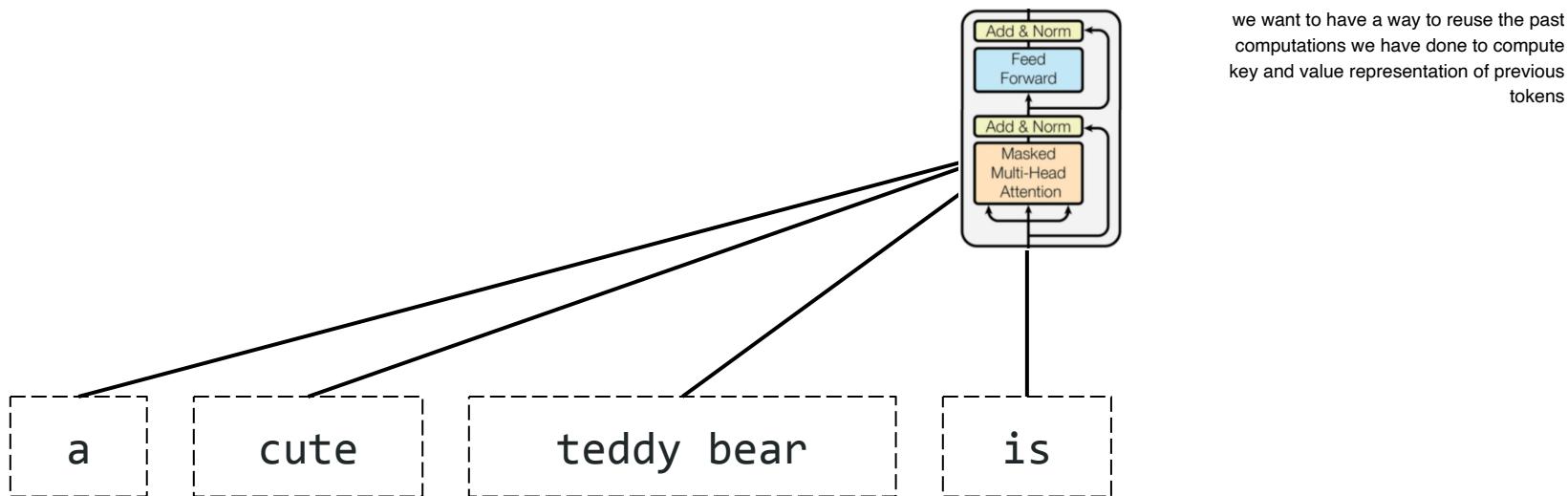
Reusing computations with KV caching

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



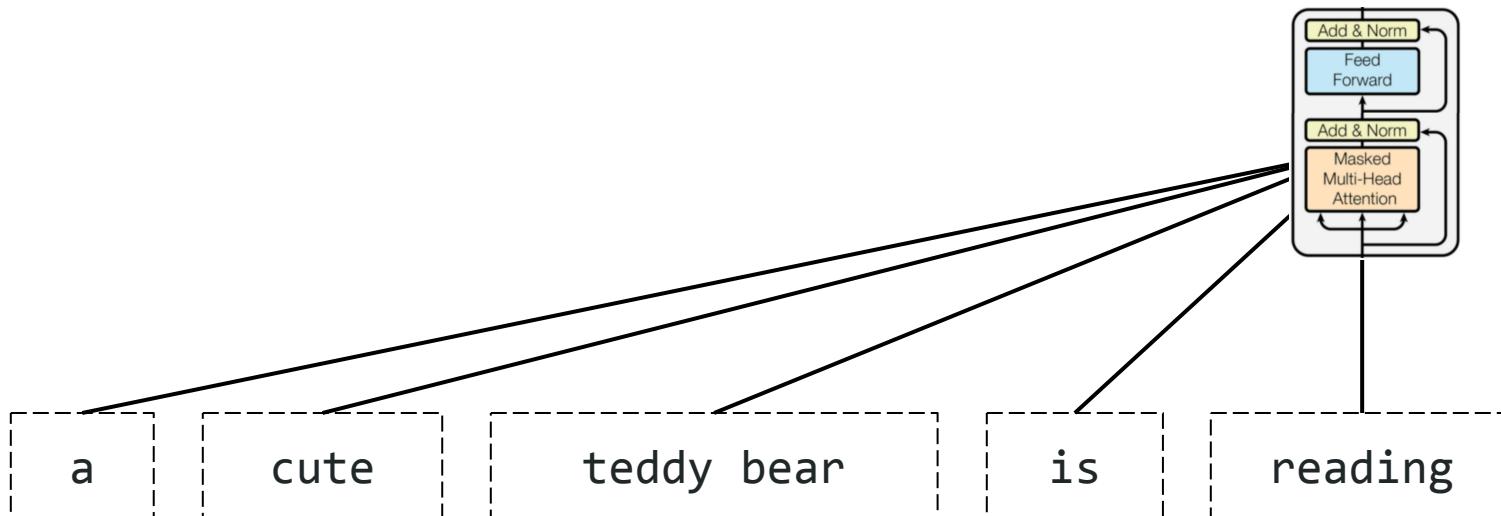
Reusing computations with KV caching

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.



Reusing computations with KV caching

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

precisely we want to
save the K and V matrices

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

Reusing computations with KV caching

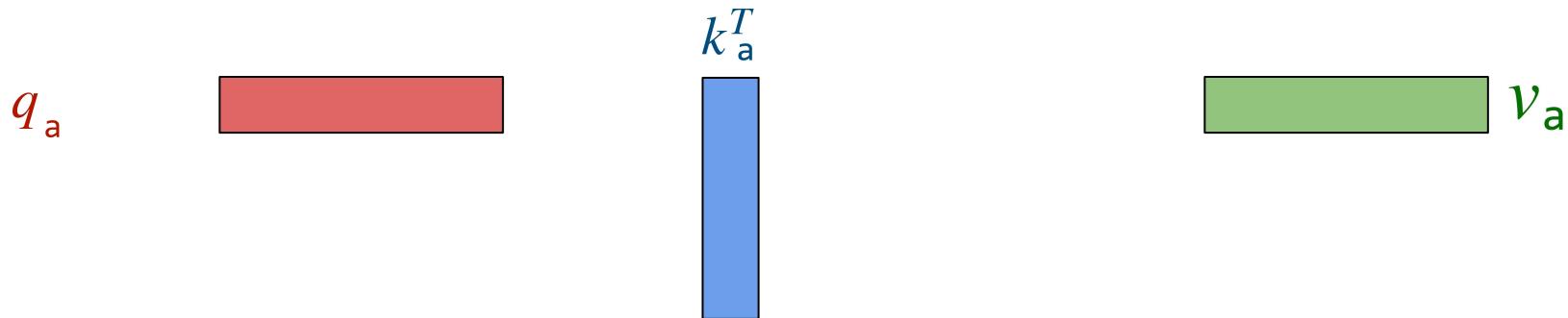
Motivation. New token needs to interact with all previous tokens.

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

Reusing computations with KV caching

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

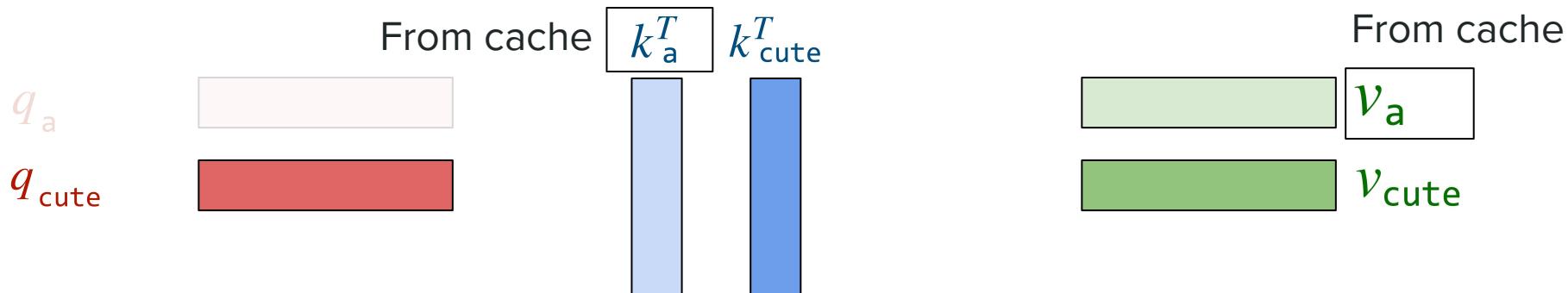
Idea. Keep keys and values in a cache.



Reusing computations with KV caching

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

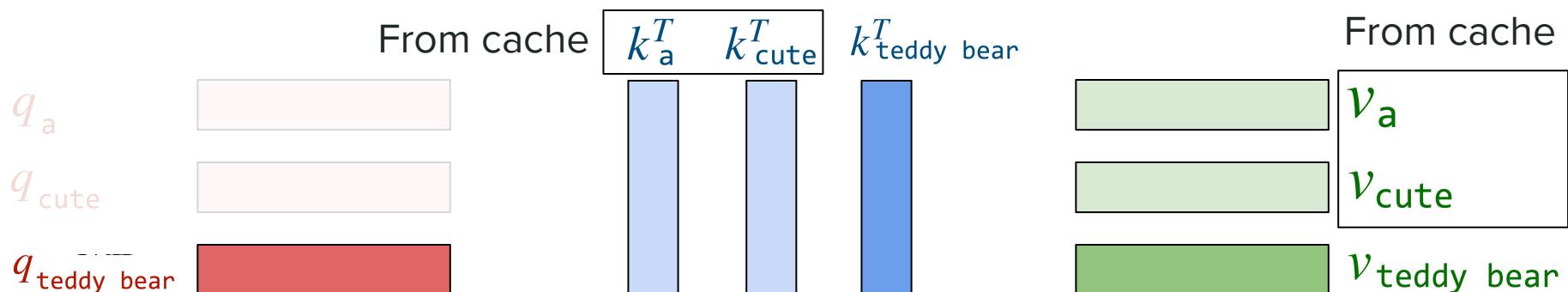
Idea. Keep keys and values in a cache.



Reusing computations with KV caching

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

Idea. Keep keys and values in a cache.

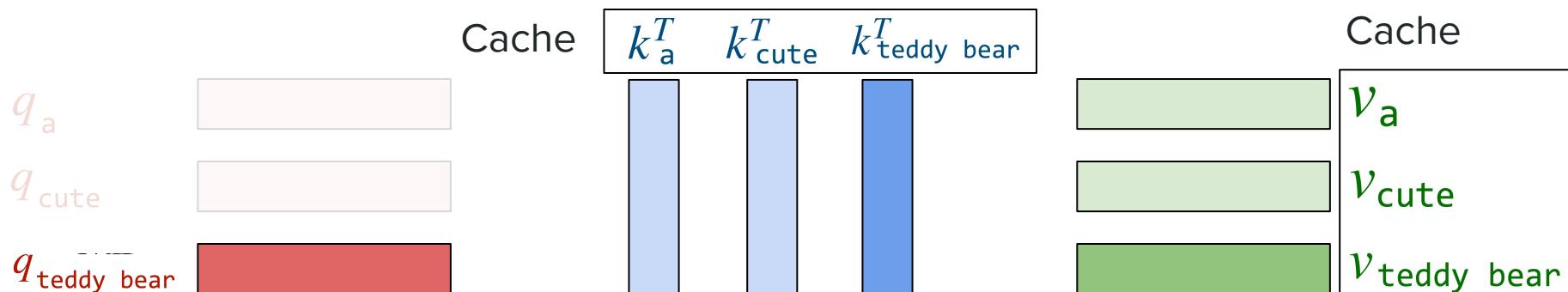


Reusing computations with KV caching

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

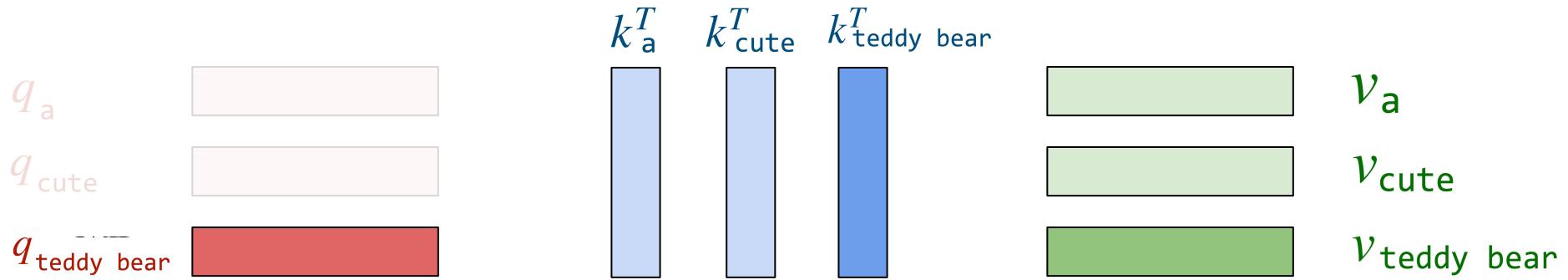
Idea. Keep keys and values in a cache.

this KV caching is not done in training as we use teacher forcing, i.e. here we just sent the whole input together at once

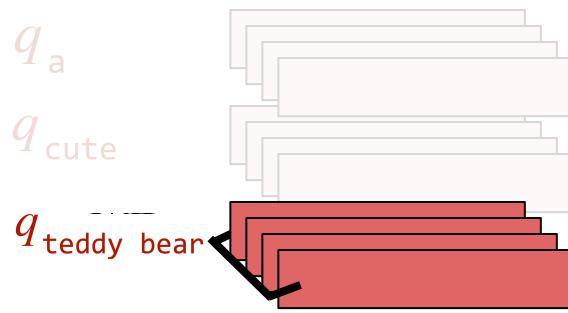


KV caching = Key-Value Caching

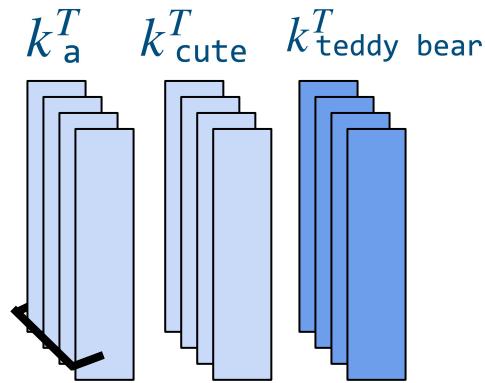
Sharing attention heads



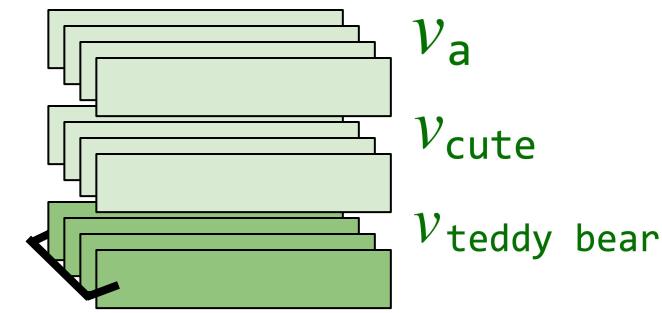
Sharing attention heads



Number of
query heads



Number of
key heads

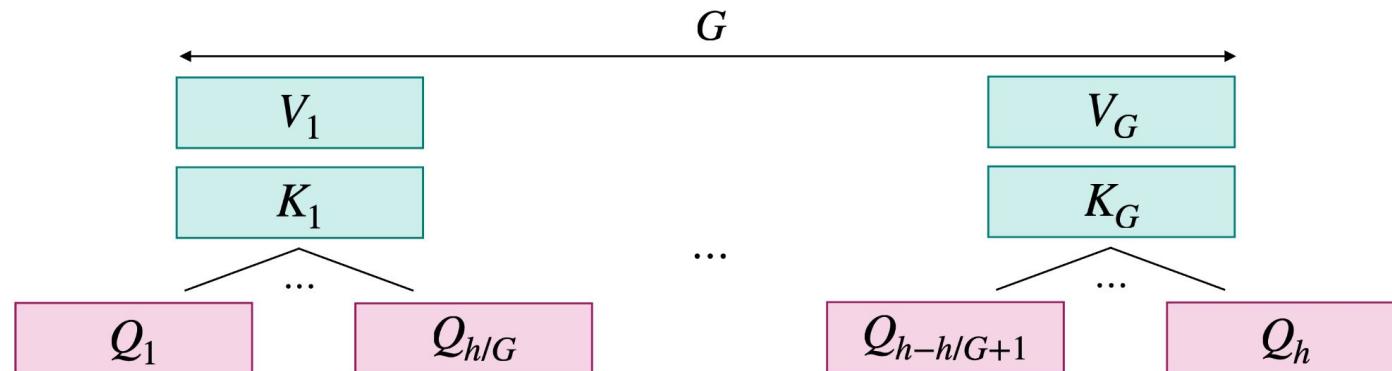


Number of
value heads

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

Sharing attention heads

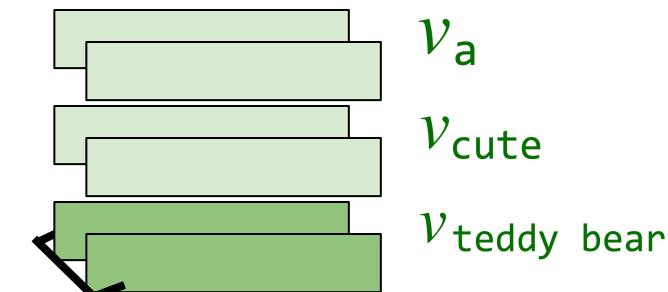
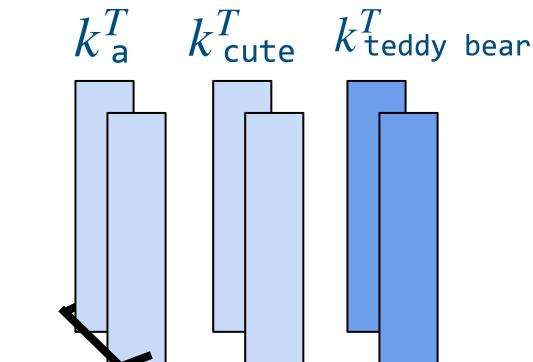
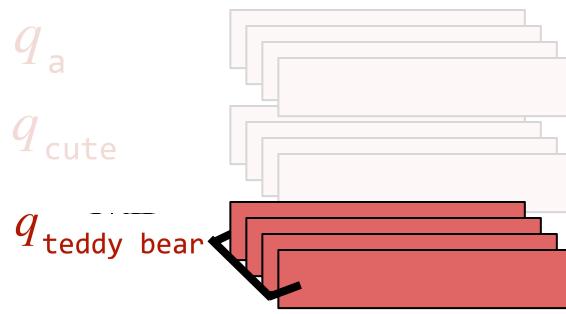
Idea. Share key/value attention heads within groups of queries



Sharing attention heads

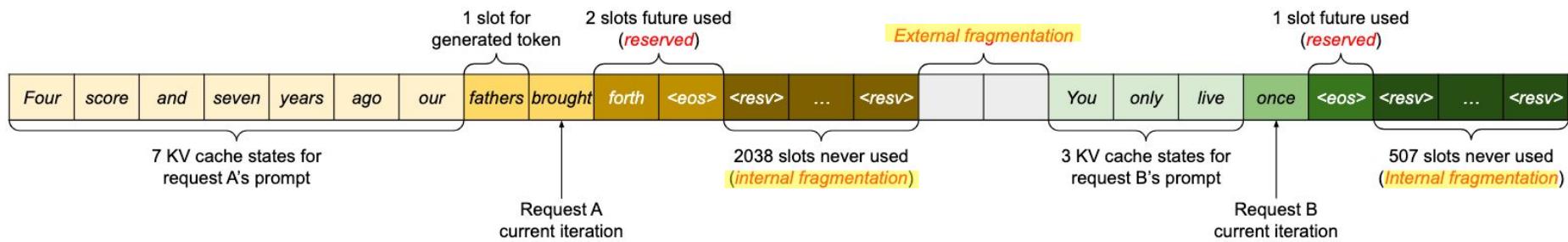
| | | |
|-------------|--|--|
| $G = 1$ | Multi-Query Attention (MQA) | <p>The diagram shows a dashed box containing two teal boxes labeled V and K. An arrow points from both V and K to a vertical ellipsis, which then points to two pink boxes labeled Q_1 and Q_h. A dashed line with a multiplier of $\times 1$ is shown below the box.</p> |
| $1 < G < h$ | Group-Query Attention (GQA) nowadays we do this mostly | <p>The diagram shows a dashed box containing two teal boxes labeled V and K. An arrow points from both V and K to a vertical ellipsis, which then points to two pink boxes labeled Q_1 and $Q_{h/G}$. A dashed line with a multiplier of $\times G$ is shown below the box.</p> |
| $G = h$ | Multi-Head Attention (MHA) | <p>The diagram shows a dashed box containing two teal boxes labeled V and K. An arrow points from both V and K directly to a single pink box labeled Q. A dashed line with a multiplier of $\times h$ is shown below the box.</p> |

Sharing attention heads with GQA



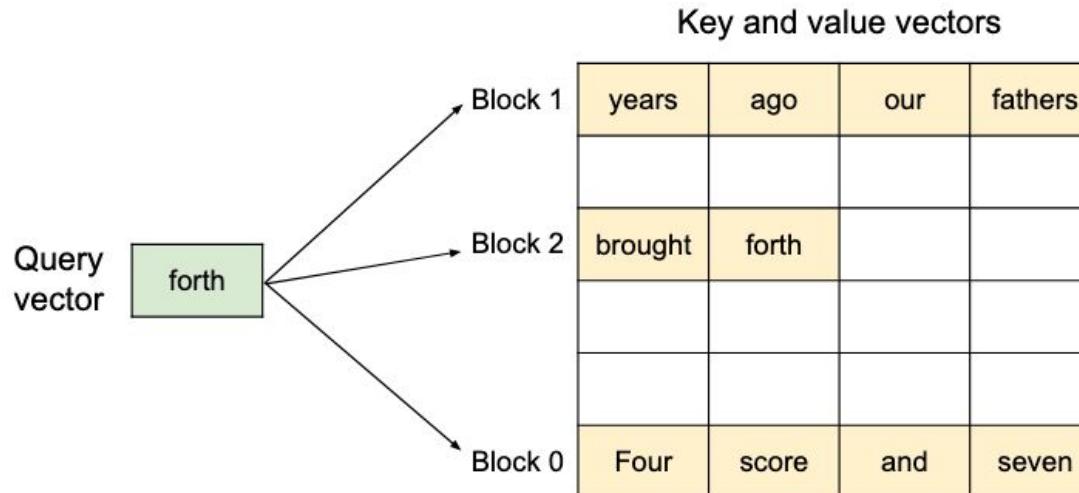
Manage memory with PagedAttention

Observation. Lots of memory waste when storing KV cache.



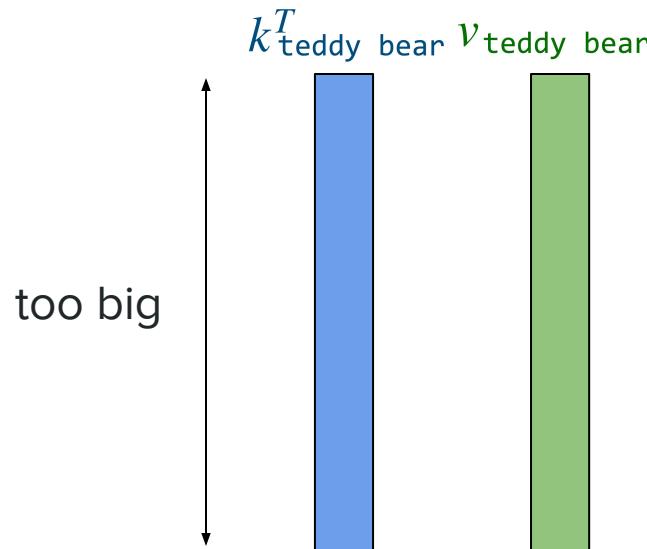
Manage memory with PagedAttention

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



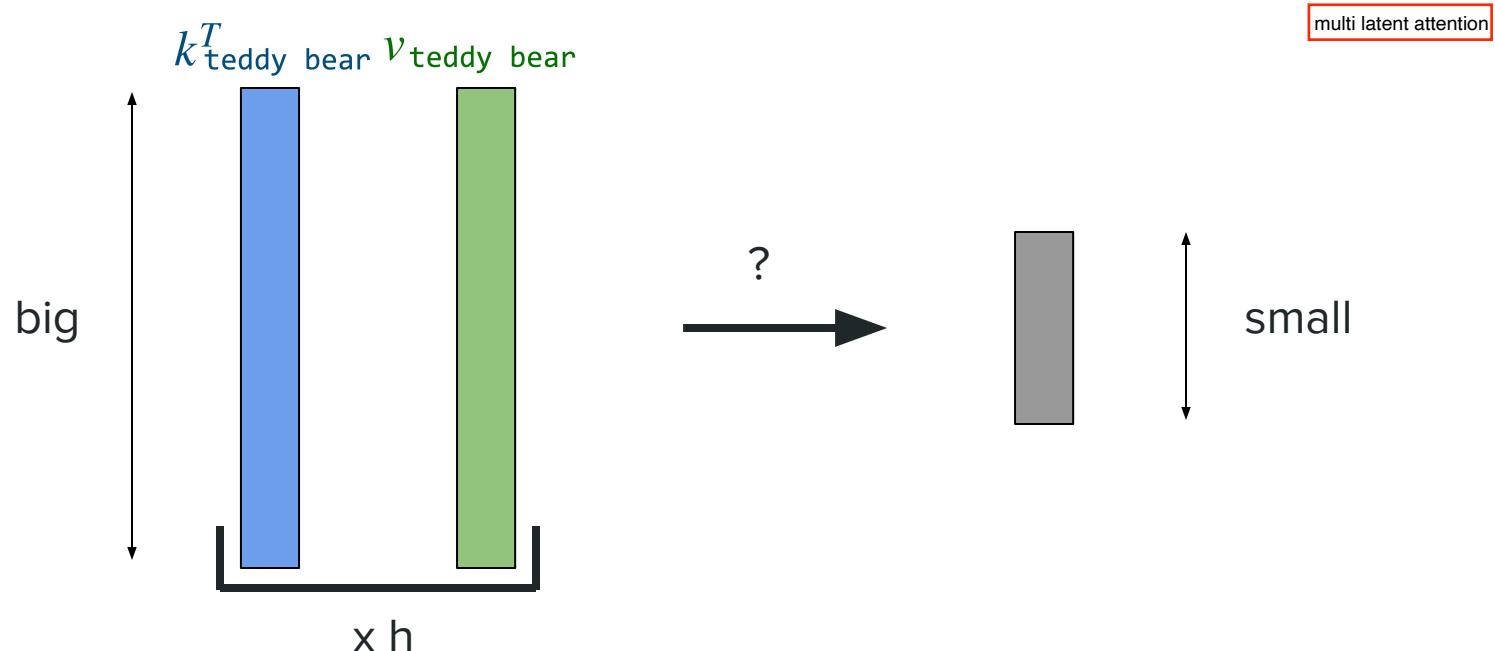
Reduce memory of KV cache with latent attention

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



Reduce memory of KV cache with latent attention

Solution. Store compressed representations instead!

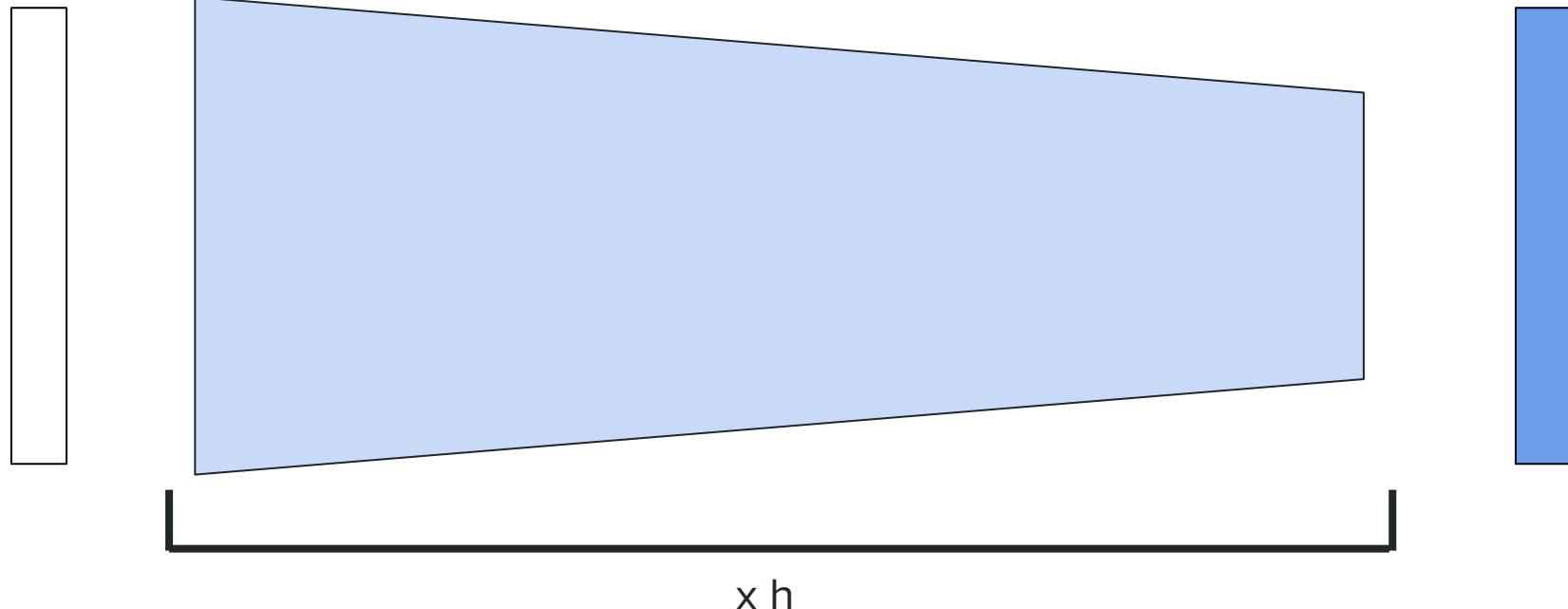


Reduce memory of KV cache with latent attention

BEFORE

teddy bear

$k^T_{\text{teddy bear}}$

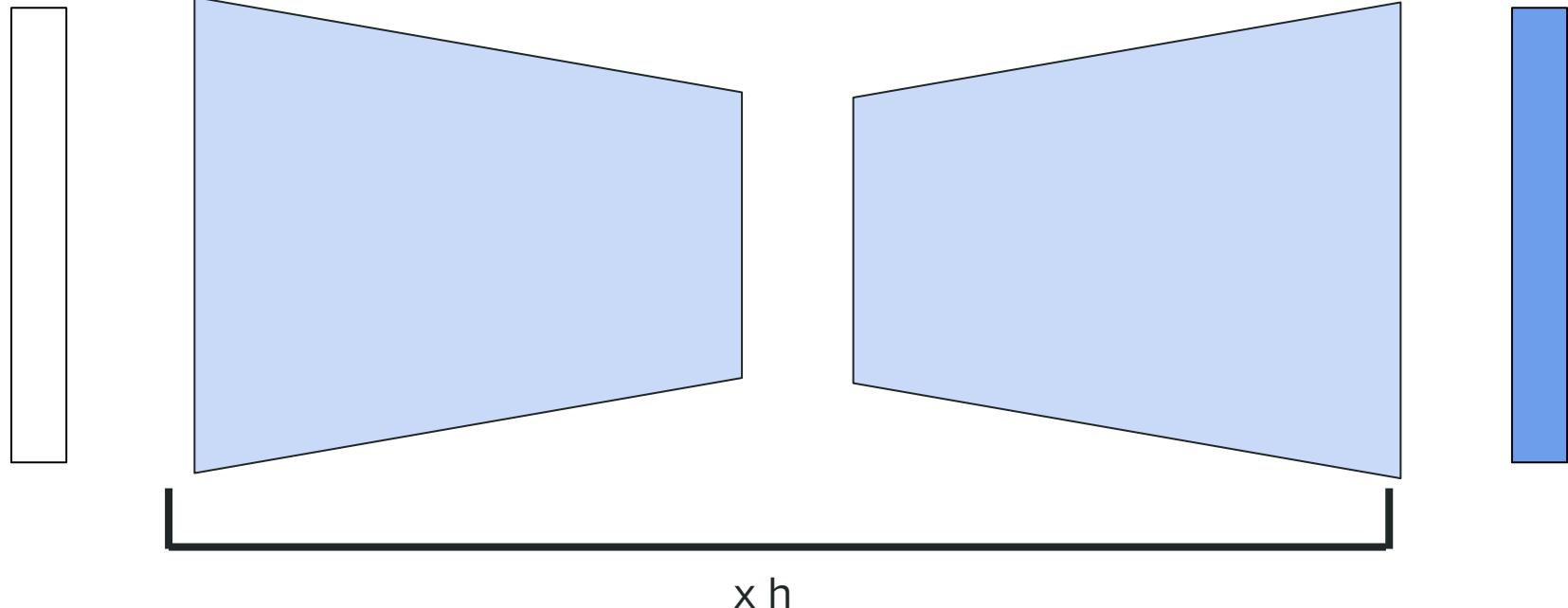


Reduce memory of KV cache with latent attention

teddy bear

factorized projection matrix into intermediary space
that is of lower dimension than the space of keys

$k^T_{\text{teddy bear}}$



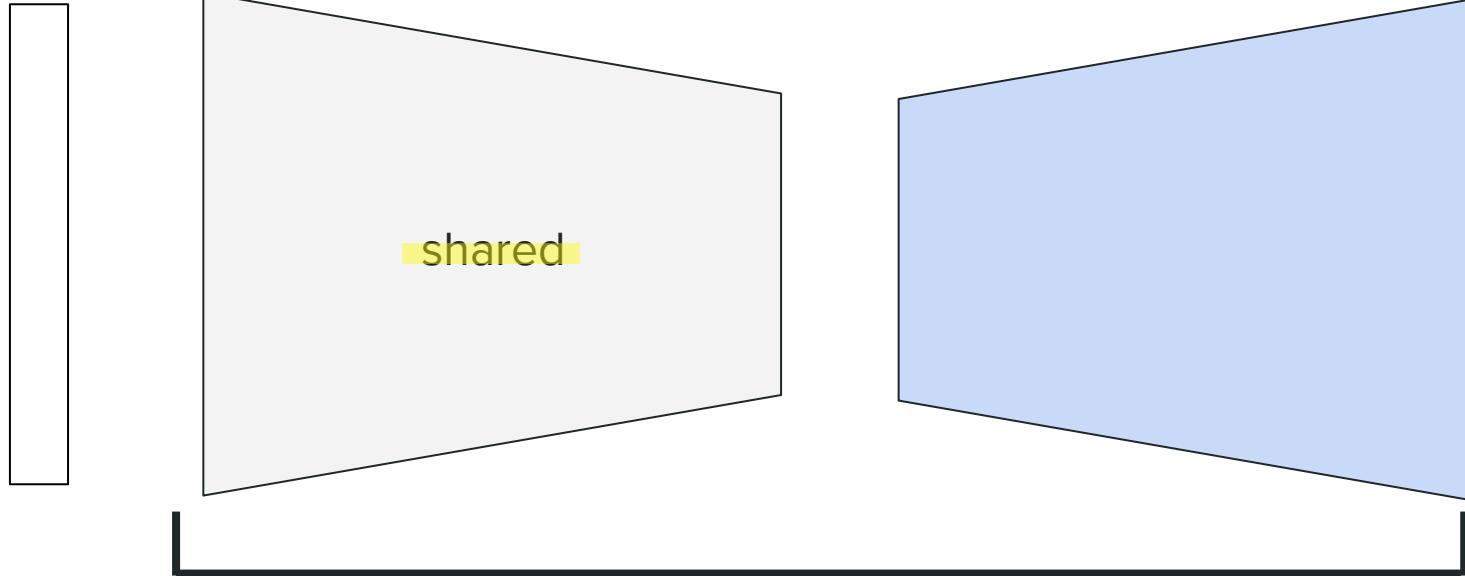
Reduce memory of KV cache with latent attention

teddy bear

low rank dimension is design choice

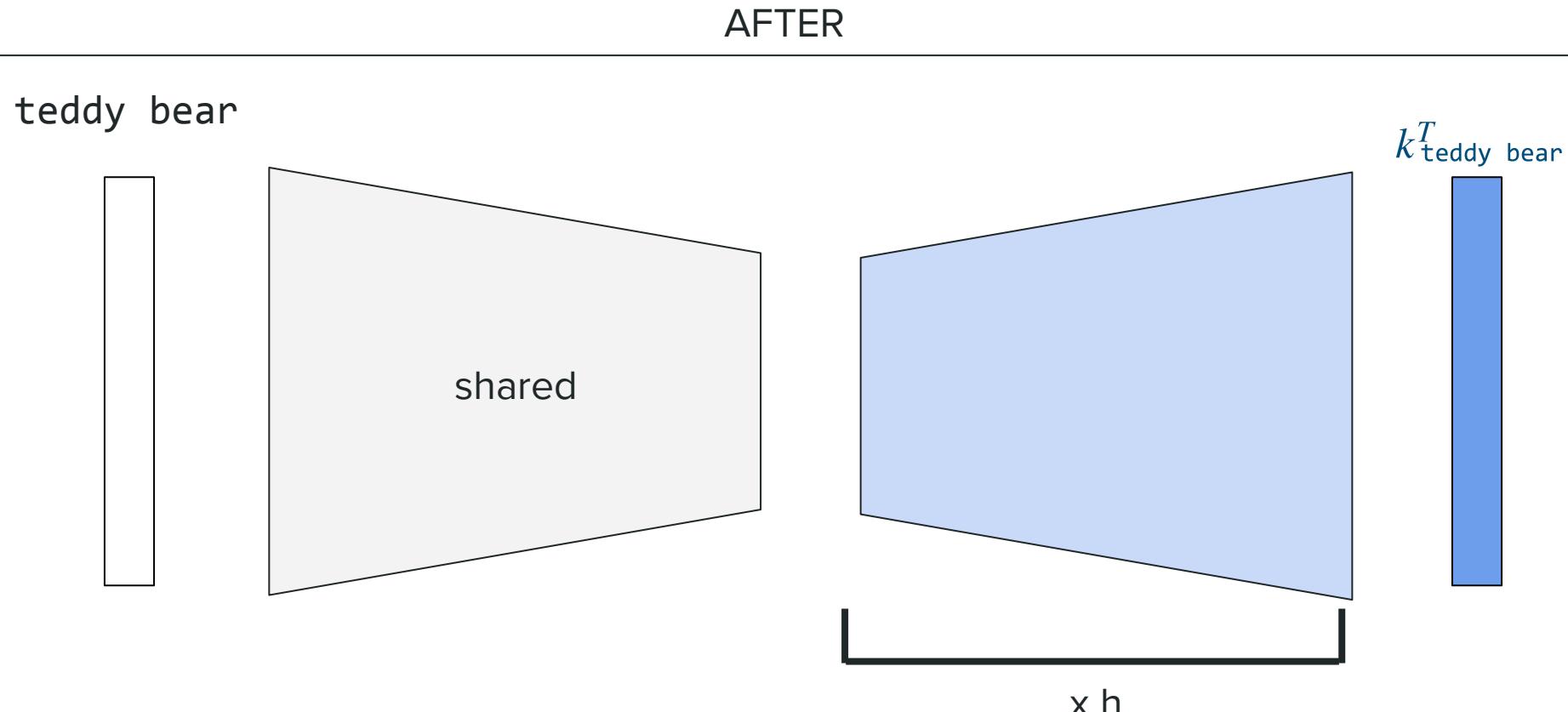
The compression matrix can be
shared across keys and values.
Across the h heads!!!

k^T teddy bear

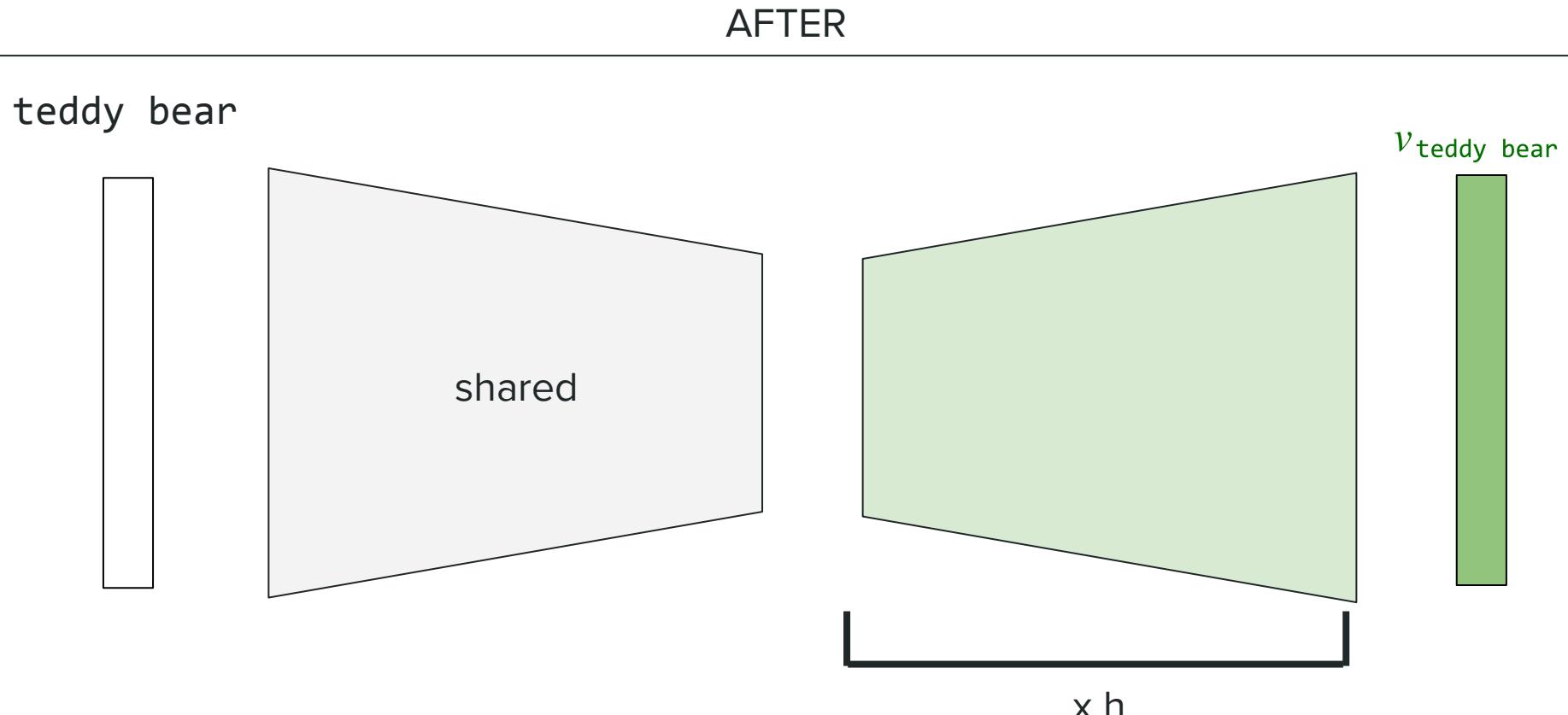


x h

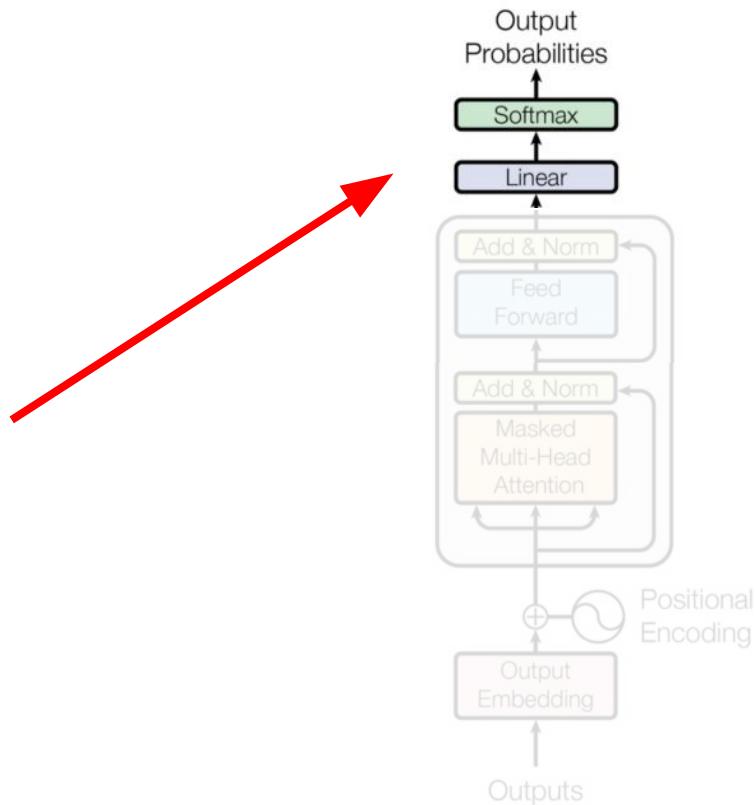
Reduce memory of KV cache with latent attention



Reduce memory of KV cache with latent attention



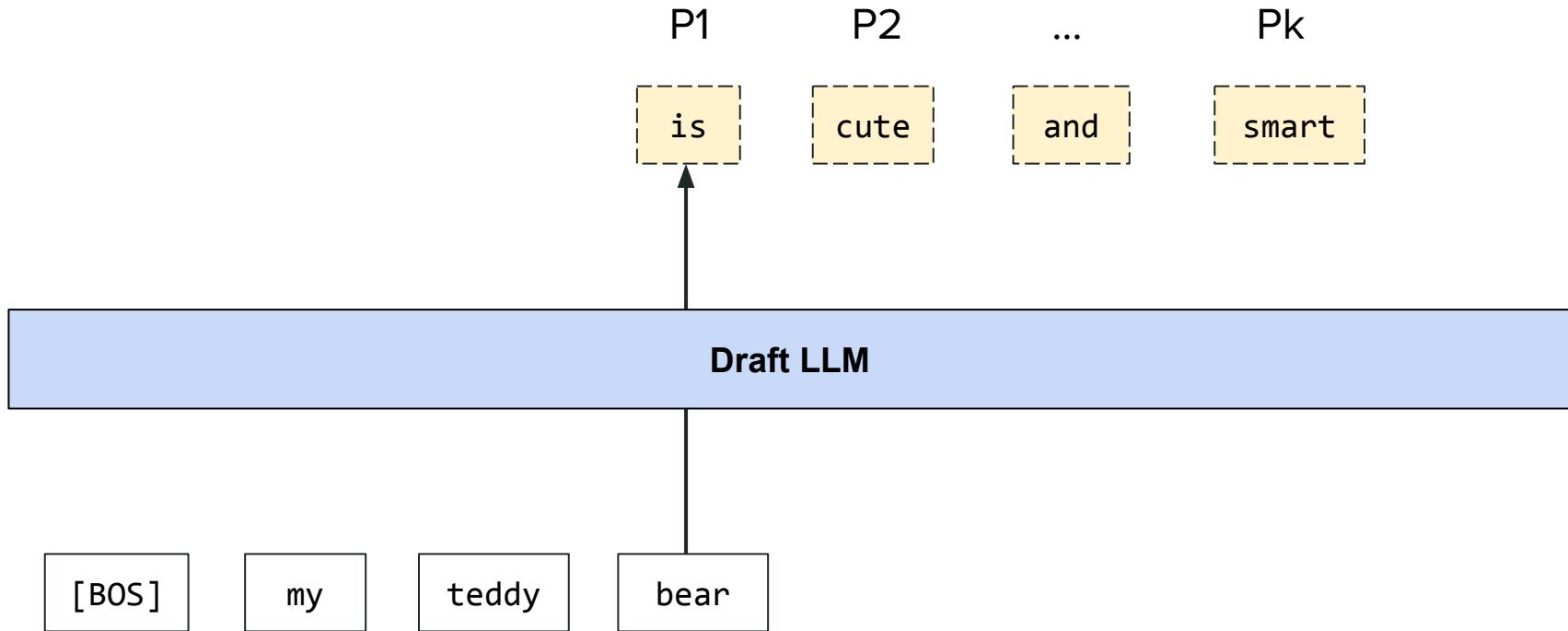
Token generation tricks



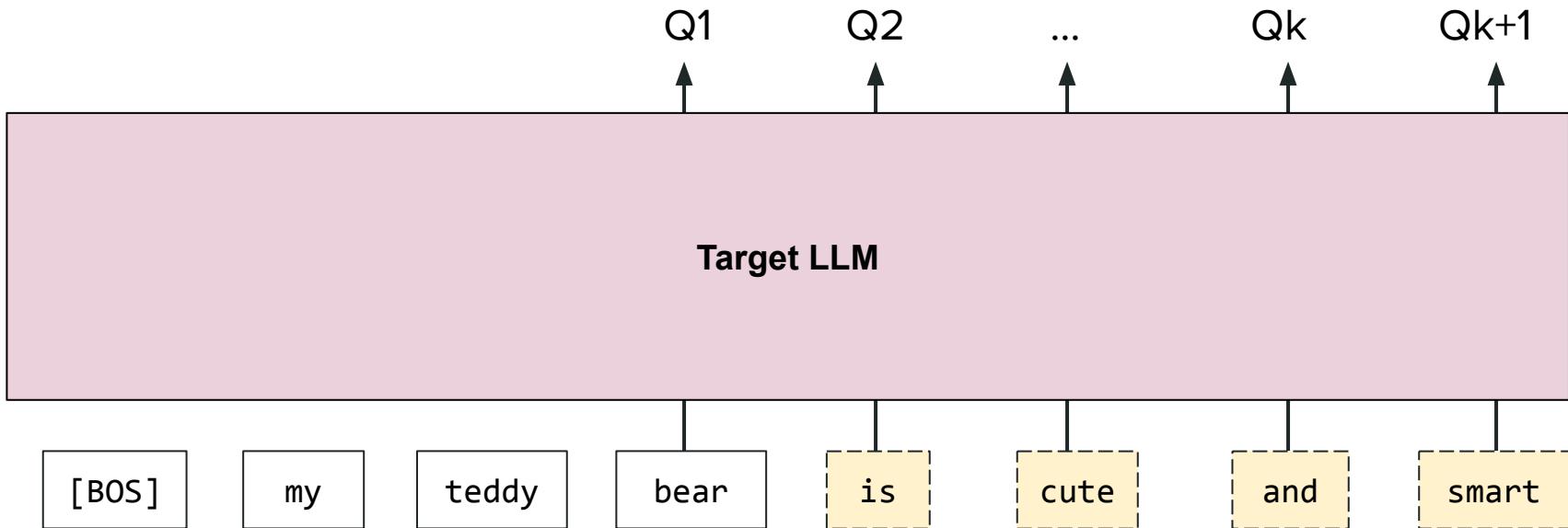
Speeding up decoding with speculative decoding

Idea. Use a draft (small) model to generate tokens validated by a target (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 $Q_i(\boxed{\text{token}}) \geq P_i(\boxed{\text{token}})$



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



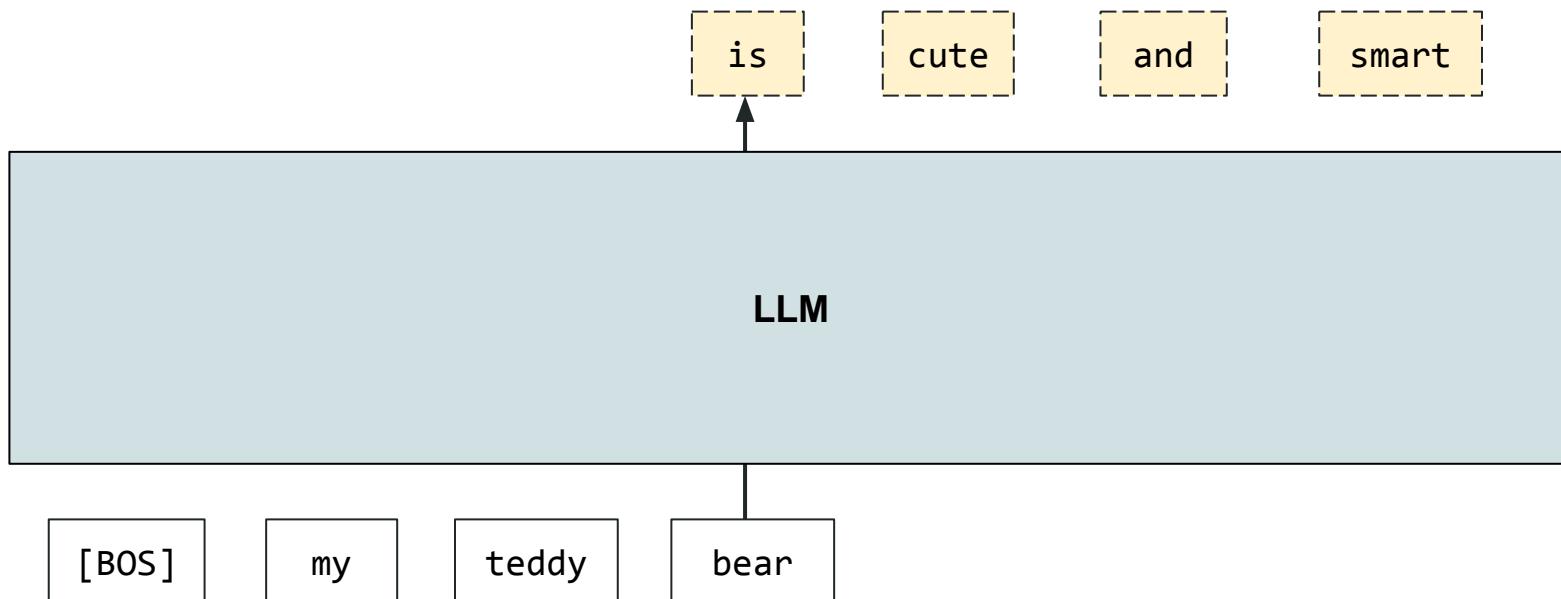
- with probability $1 - Q_i(\boxed{\text{token}}) / P_i(\boxed{\text{token}})$



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!

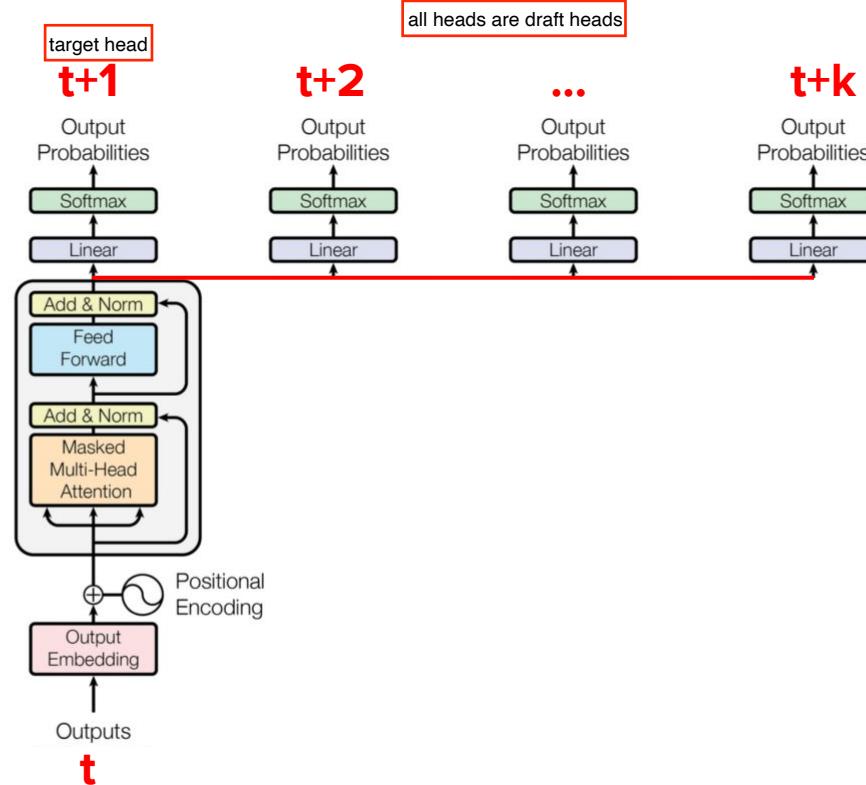


Figure adapted from "Attention is All You Need", Vaswani et al., 2017.

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



Approximations:

- Architectural changes
- Embeddings representations
- Token prediction



Thank you for your attention!
