



### Workshop on Systems for LLMs

Systems for Large Language Models (LLMs) Workshop with Hands-On Exploring the Infrastructure Behind Intelligent Systems

#### **Sessions Overview:**

LLM building blocks
GPUs for LLM
Inference serving with vLLM
Distributed Training

Tools: Hugging Face, PyTorch



#### **Details**

Date: 6<sup>th</sup>,7<sup>th</sup> September

Time: 9:30 am-5pm

Register @

https://forms.gle/a8rQQ835jJG42i3eA

Supported by: IBM-IITB Academic Research Partnership



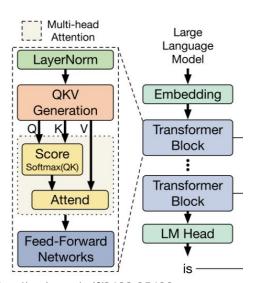


# Inference Systems

## What is inference?

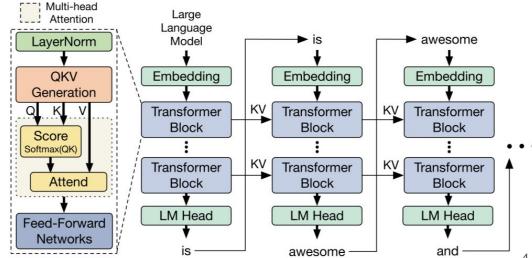
### forward vs generate

 Forward: single forward pass through LLM

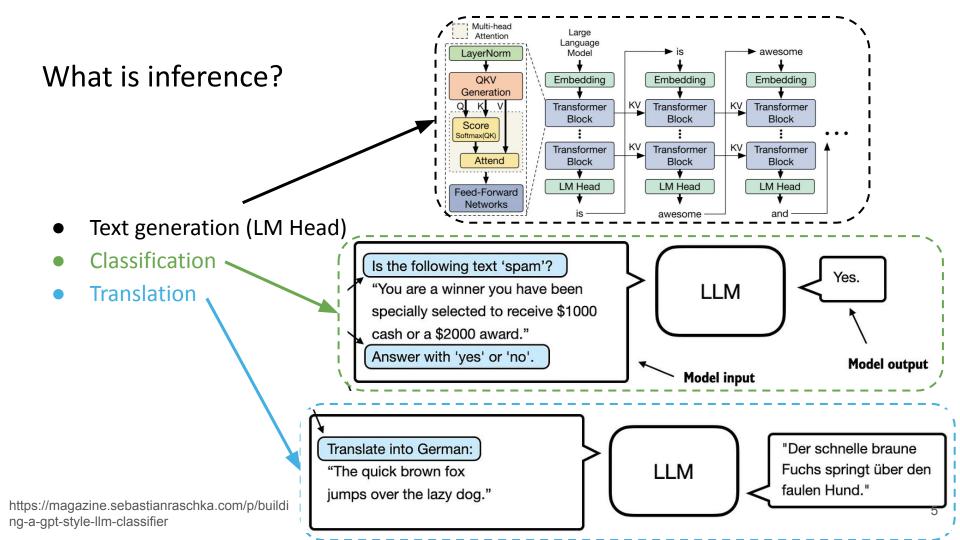


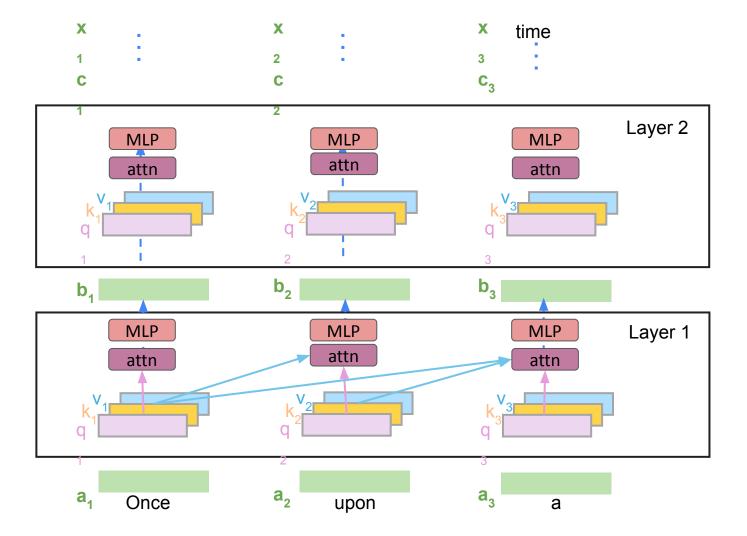
#### Generate:

 Generates multiple tokens, until EOS or max tokens reached



source: https://arxiv.org/pdf/2408.05499





### Simplest generate example

```
from transformers import AutoTokenizer, AutoModelForCausalLM
import torch
tokenizer = AutoTokenizer.from_pretrained("gpt2")
model = AutoModelForCausalLM.from_pretrained("gpt2")
prompt = "My trip to Yosemite was"
inputs = tokenizer(prompt, return_tensors="pt")
output = model.generate(inputs.input_ids, max_length=100)
generated_text = tokenizer.batch_decode(output, skip_special_tokens=True)
print(generated_text)
```

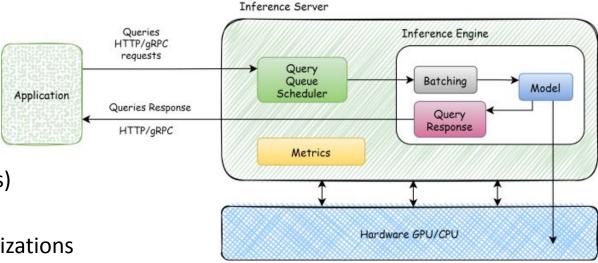
# Inference engines

### What are inference engines?

Wrap systems concepts around the process

of inference

- Batching
- Streaming
- Fault management
- Error handling
- Parallelism (multi-GPUs)
- Optimizations:
  - Memory storage optimizations
  - Offloading
- Most popular open-source is vllm. Others include SGLang, NVIDIA's Tensor-RT LLM

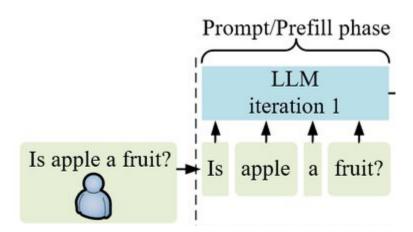


	O	Ba ptim	tch izati	ion		Paralle	elism		Cor	npres	ssion	Fine- tuning	C	Cachi	ng		(	Atte Optin	ntio nizati			Sampling Optimization	Output Optimization	
Engines	Dynamic Batching	Continuous Batching	Nano Batching	Chunked-prefills	Data Parallelism	Fully Sharded Data Parallel	Tensor Parallelism	Pipeline Parallelism	Quantization†/ Support Quantized Model	Pruning	Sparsity*	LoRA	Prompt Caching	Prefix Caching	KV Caching	PagedAttention	vAttention	FlashAttention	RadixAttention	FlexAttention	FireAttention	Speculative Decoding	Structured Outputs*	
Ollama [194]	×	×	×	×	×	×	V	V	V	V	(MoE, SL)	~	V	×	V	×	×	~	×	×	×	V	(GG)	
llama.cpp [82]	×	V	×	×	×	×	V	V	<b>~</b>	×	(MoE, SL)	V	V	×	V	×	×	V	×	×	X	V	(GG)	
vLLM [125]	×	V	×	V	V	V	V	V	(A, A*, B, D, G, L, M)	V	(BS, MoE, N:M)	V	×	V	V	V	×	(v3)	×	×	×	V	(LM, OA, OL, XG)	
DeepSpeed-FastGen [102]	×	V	×	V	V	V	V	V	V	(S, U)	(N:M, NxM, MoE, SA)	V	×	×	V	V	×	(v2)	×	×	X	×	×	
Unsloth [244]	×	×	×	×	×	×	×	×	(B)	×	×	~	×	×	V	×	×	V	×	V	×	×	×	
MAX [179]	×	~	×	V	×	×	V	×	V	×	(MoE)	V	×	V	V	V	×	(v3)	×	×	×	V	(XG)	
MLC-LLM [177]	×	V	×	V	×	×	V	V	(A)	×	(MoE)	×	×	V	V	V	×	×	×	×	×	V	(XG)	
llama2.c [21]	×	×	×	×	×	×	×	×	~	×	×	×	×	×	V	×	×	×	×	×	×	×	×	
bitnet.cpp [251]	×	×	×	×	×	×	×	×	(A, G)	×	(MoE)	×	×	×	~	×	×	×	×	×	×	×	×	
SGLang [295]	×	~	×	~	~	~	~	×	(A, B, G, L, M)	~	(DSA, N:M, MoE)	~	×	~	~	~	×	×	~	×	×	~	(LL, OL, XG)	
LitGPT [145]	×	V	×	×	V	V	V	×	(B)	×	(MoE)	V	×	×	V	×	×	(v2)	×	×	×	V	×	
OpenLLM [30]	×	V	×	×	V	×	×	×	(B, G)	×	×	×	×	×	×	×	×	×	×	×	×	×	×	
TensorRT-LLM [191]	~	~	×	~	~	×	V	V	(A, G, S, W)	~	(BSA, MoE)	~	~	×	V	V	×	×	×	×	×	V	(XG)	
TGI [110]	×	~	×	×	×	×	V	×	(A, E, E*, G, M)	V	(N:M, MoE)	~	×	V	V	V	×	(v2)	×	×	×	<b>~</b>	(OL)	
PowerInfer [227, 270]	×	V	×	×	V	×	×	V	V	×	V	V	V	×	V	×	×	V	×	×	×	V	(GG)	
LMDeploy [162]	×	-	×	~	×	×	-	×	(A, G, S)	V	(MoE)	V	×	V	V	V	×	×	×	×	×	×	(PT)	
LightLLM [142]	~	×	×	~	×	×	V	×	V	×	(MoE)	×	V	×	V	×	×	(v1)	×	×	×	×	(OL, XG)	A Survey on Inference
NanoFlow [300]	×	~	V	V	V	×	×	×	×	×	×	×	×	×	V	×	×	×	×	×	×	×	×	Engines for Large
DistServe [297]	~	~	×	~	×	×	~	~	×	×	×	×	×	×	~	~	×	(v1)	×	×	×	×	×	Language Models:
vAttention [206]	V	×	×	×	V	×	V	V	V	V	(N:M)	~	×	×	V	V	~	(v2)	×	×	×	×	×	Perspectives on
Sarathi-Serve [8]	×	×	×	~	×	×	V	V	×	×	(MoE)	×	×	×	~	×	×	(v2)	×	×	×	X	×	Optimization and
Friendli Inference [71]	-	~	-	-	-	-	V	~	<b>~</b>	-	(MoE)	~	-	-	-	_	×	_	-	×	×	<b>V</b>	~	Efficiency
Fireworks AI [67]	-	V	_	-	-	-	-	-	V	V	(MoE)	V	V	-	V	-	×	-	-	×	V	~	(OA)	(https://arxiv.org/pdf/25
Groq Cloud [89]	×	_	-	_	V	1-	V	V	~	V	(MoE)	_	_	_	_	_	-	-	_	×	×	~	(OA)	05.01658) <sub>10</sub>
Together Inference [239]	-	-	_	_	_	V	-	_	V	_	(MoE)	~	V	_	_	-	×	(v3)	_	×	×	V	V	

# **Terminology**

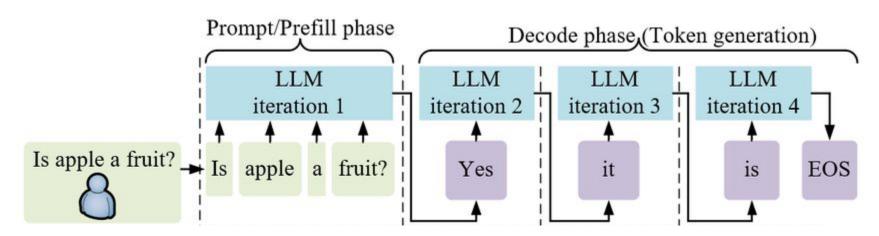
### Prefill vs decode

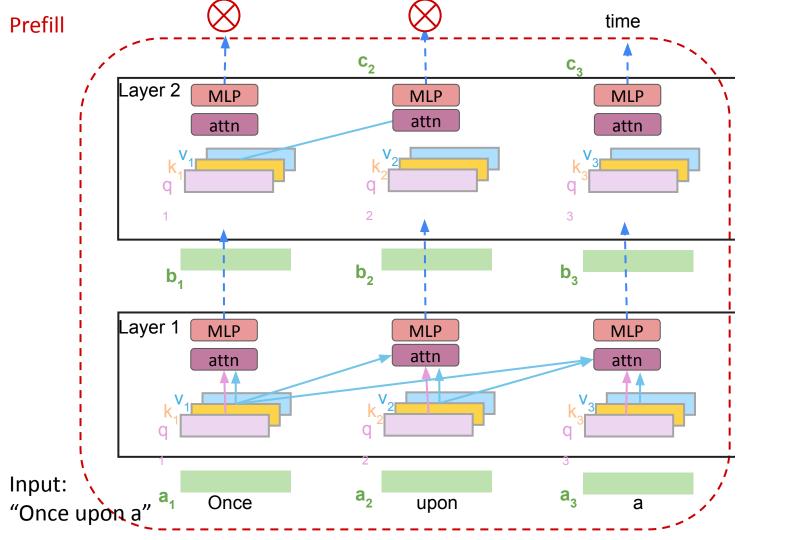
• Prefill: input/prompt phase. Generation of the first output token (ingestion of the entire input in parallel)

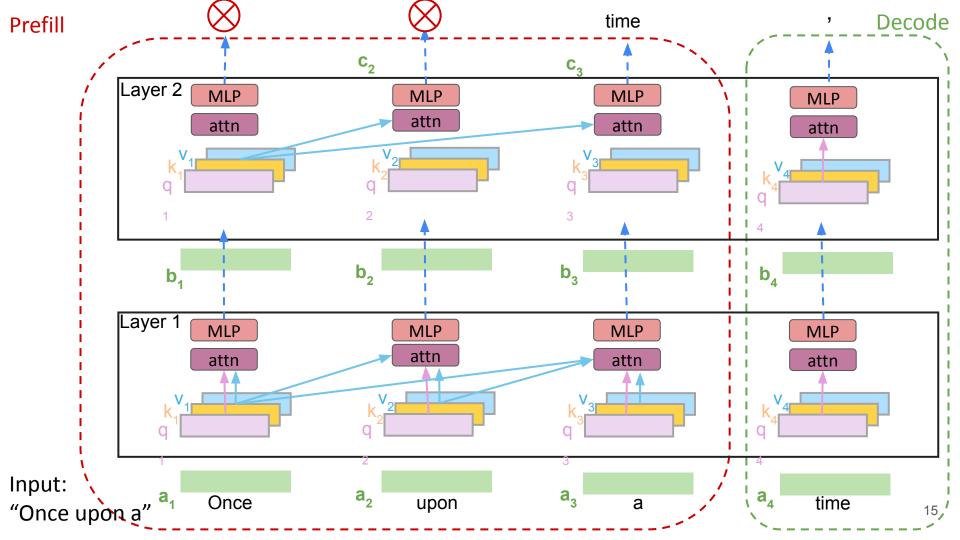


#### Prefill vs decode

- Prefill: input/prompt phase. Generation of the first output token (ingestion of the entire input in parallel)
- Decode: output/generation phase. Decoding of subsequent output tokens (sequential)

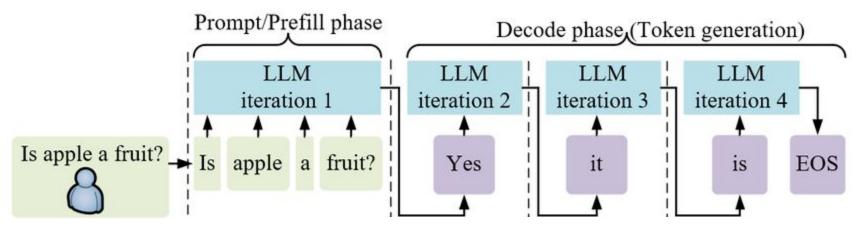






#### **Metrics**

- Throughput: Number of Tokens processed per second
- Latency:
  - Time to first token (prefill time)
  - Time between tokens (decode time)



# Let's write our own generate

### Step-wise plan

#### Make sure pytorch\_model.bin is in the exercises directory

wget https://huggingface.co/gpt2/resolve/main/pytorch\_model.bin

- pytorch\_model.bin
- simple
- kvcache

#### First let's take a look at gpt.py

- 1. Create position ids tensor (What is the position ids for this input?) [1, 2, ... n]
- Call forward(inputs, position\_ids) (Make sure both are tensors)
- 3. What is output (logits) shape? We will learn about this later
- 4. Append output token to the inputs tensor
- 5. Wrap entire code in a "for" loop (What is the exit condition?)

```
def generate(input):
    # Tokenize input
    tokenized = tokenizer.encode(input)
    inputs = torch.tensor(tokenized)
                                                                    See the signature of
                                                                    GPTLMHead.forward()
    while ???: # What is exit condition of the loop?
        positions = ??? # Create position ids tensor
        logits = model(???) # Call model.forward with inputs and position ids
        logits = logits[-1, :]
        next token = torch.argmax(logits, dim=-1, keepdim=True)
        inputs = ??? # Concatenate next token to inputs for next step here
    output = tokenizer.decode(inputs)
    return output
input = 'The quick brown fox jumped'
                                                                                      20
output = generate(input)
```

MAX SEO LEN = 10

model = GPTLMHead(config)

config = GPT2Config.from pretrained('gpt2')

tokenizer = AutoTokenizer.from pretrained("gpt2")

```
MAX SEQ LEN = 10
config = GPT2Config.from pretrained('openai-community/gpt2-large')
model = GPTLMHead(config)
tokenizer = AutoTokenizer.from pretrained("gpt2")
def generate(input):
    # Tokenize input
    tokenized = tokenizer.encode(input)
    inputs = torch.tensor(tokenized,)
    while len(inputs[0]) < MAX SEQ LEN:
        positions = torch.arange(len(inputs))
        logits = model(inputs, positions)
        logits = logits[-1, :]
        next token = torch.argmax(logits, dim=-1, keepdim=True)
        inputs = torch.cat((inputs, next token), dim=0)
    output = tokenizer.decode(inputs.tolist())
    return output
input = 'The quick brown fox jumped'
```

output = generate(input)

## Core features

### Core features of inference engines

Inference engines such as vllm etc. have certain features/optimizations

We will look at 3 such techniques:

- 1. KV Caching
- 2. Continuous Batching
- 3. Paged Attention

# **KV** Caching

### Why cache KVs?

#### What are K and V?

For causal self-attention, K/V of all past tokens is required to calculate attention of current token

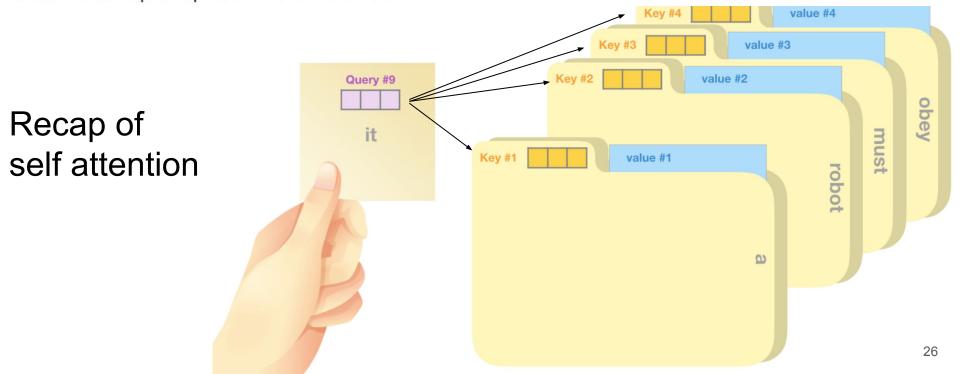
Problem: It is recalculated again and again for every token

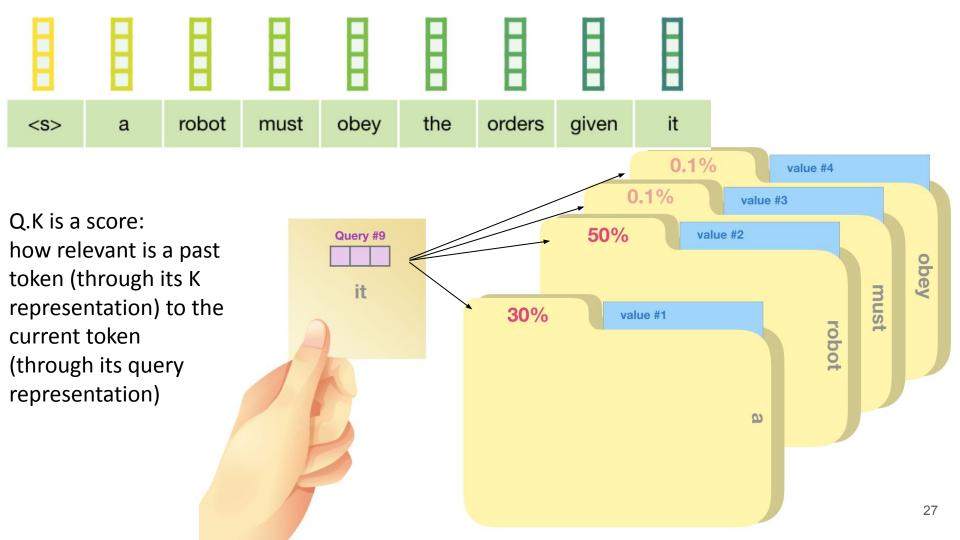
Idea: Generate and cache once, reuse for subsequent tokens

Query: The query is a representation of the current word used to score against all the other words (using their keys). We only care about the query of the token we're currently processing.

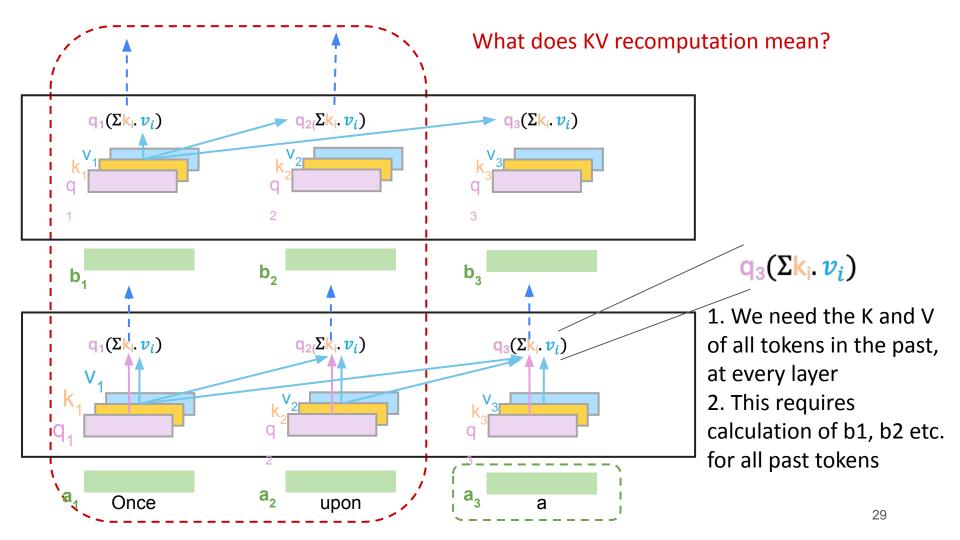
Key: Key vectors are like labels for all the words in the segment. They're what we match against in our search for relevant words.

Value: Value vectors are actual word representations, once we've scored how relevant each word is, these are the values we add up to represent the current word.



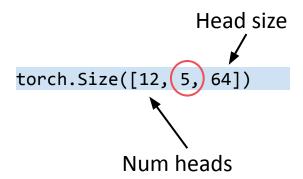


Word	Value vector	Score	Value X Score			
<s></s>		0.001				
a		0.3				
robot		0.5				
must		0.002				
Finally, each past token	(through its Value	0.001				
•	iplied with the score and new representation of t	0.0000				
current token	new representation or t	0.005				
given		0.002				
it		0.19				
		Sum:				

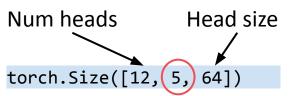


### Let's see the repetition in action

- 1. GPTAttention has a field self.idx which refers to the layer
- Now print the shape of the k or v tensor for a particular layer after the first 5 lines of the attention block before SDPA



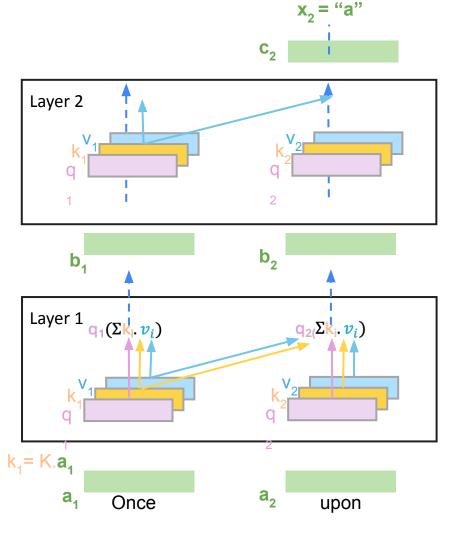
### Let's see the repetition in action



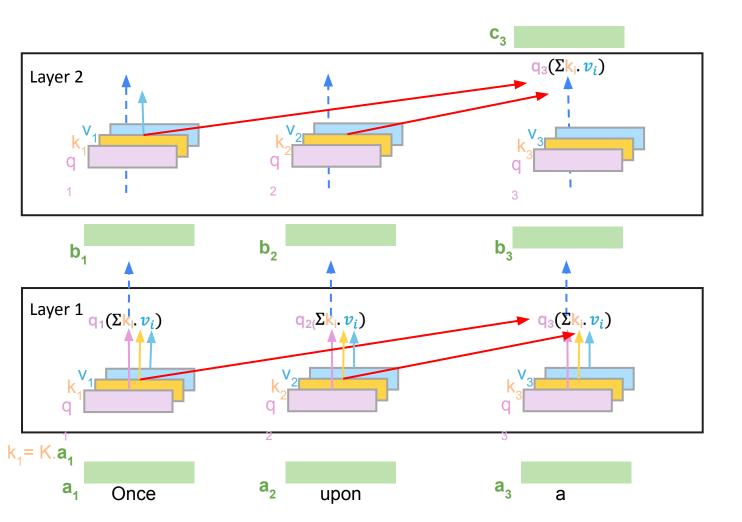
- GPTAttention has a field self.idx which refers to the layer
- 2. Now print the shape of the k or v tensor for a particular layer
- 3. One dimension grows every iteration that is the number of tokens
- 4. Pick a random idx and print keys[a][b][c] (Choose random valid values of a, b, c)

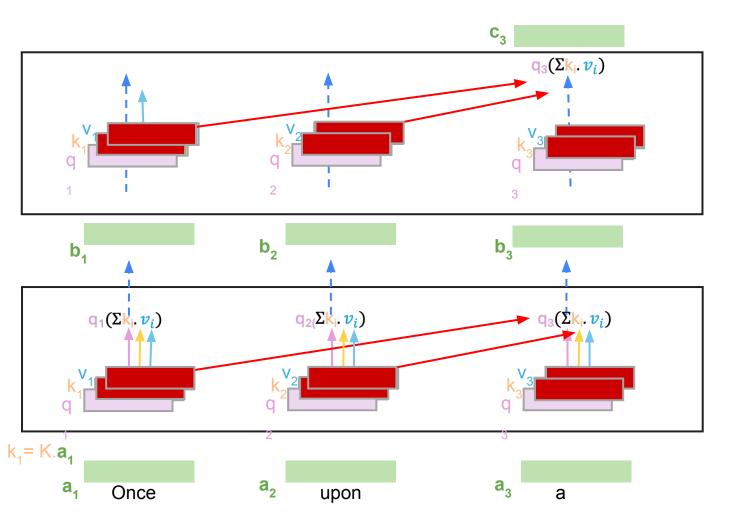
```
class GPTAttention(nn.Module):
    def forward(self, x):
        batch_size, seq_len, _ = x.shape
        q, k, v = ...
        queries = ...
        keys = ...
        values = ...

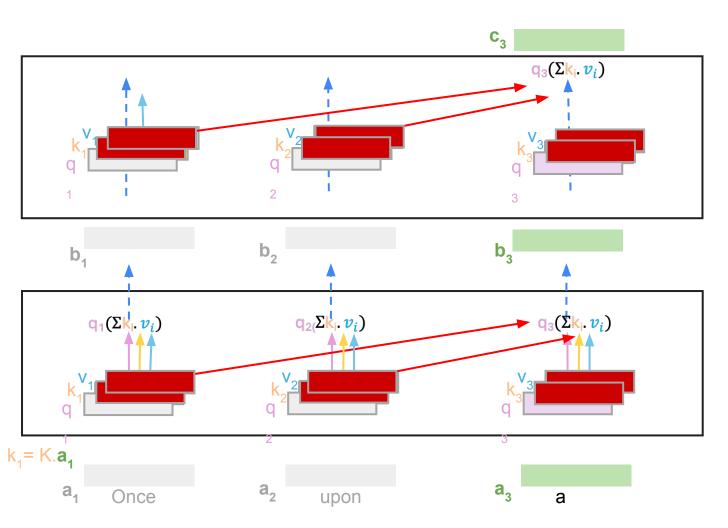
if self.idx == 0:
        print(keys[7][4][29])
```



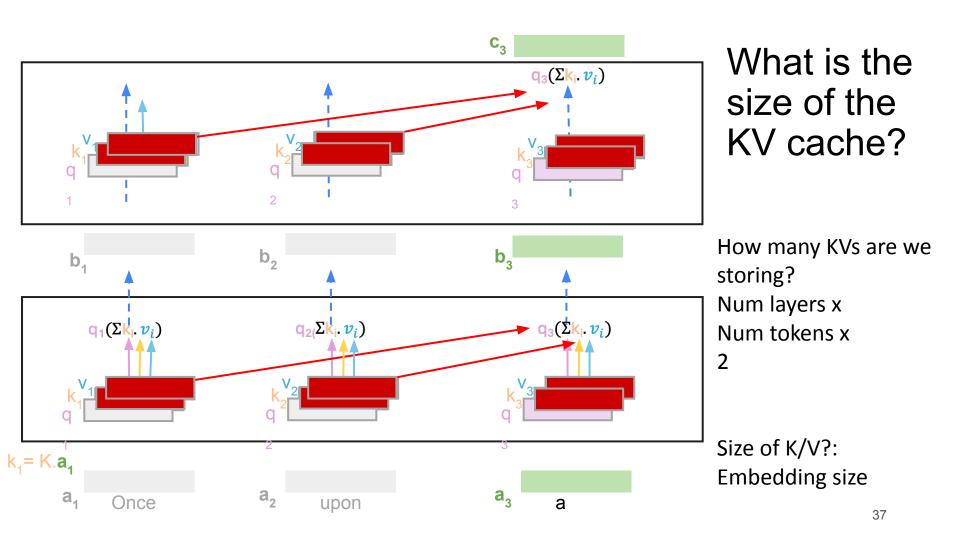
q<sub>1</sub> = Q.a<sub>1</sub>
b<sub>1</sub> is output after layer 1
For every token, only its q matters
Causal attention: For every token, only its
and past tokens k,v matter







We don't pass entire input for every decode run, only the last token



# Let's implement KV Cache (gpt.py) Only 12 lines to be added Only 9 lines to be changed

- Create an argument called kv\_cache to the forward in all the modules calling Attention
  - a. GPTLMHead, GPTModel, GPTBlock, GPTAttention

# Let's implement KV Cache (gpt.py) Only 12 lines to be added Only 9 lines to be changed

- Create an argument called kv\_cache to the forward in all the modules calling Attention
  - a. GPTLMHead, GPTModel, GPTBlock, GPTAttention

# Let's implement KV Cache (gpt.py) Only 12 lines to be added Only 9 lines to be changed

- Create an argument called kv\_cache to the forward in all the modules calling Attention
  - b. In GPTModel, kv\_cache should be initialized as [None] x n\_layers for each block. Now, we'll pass the kv\_cache[i] to each block's forward call

# Let's implement KV Cache (gpt.py) Only 12 lines to be added Only 9 lines to be changed

- 1. Create an argument called kv\_cache to the forward in all the modules calling Attention
  - b. In GPTModel, kv\_cache should be initialized as [None] x n\_layers for each block. Now, we'll pass the kv\_cache[i] to each block's forward call

```
class CausalModel(nn.Module):
    def forward(self, inputs, position_ids, kv_cache=None):
        if kv_cache is None:
            kv_cache = [None for _ in range(len(self.h))]
        ...
        h(x, kv_cache[i])
```

- 2. Init case (Prefill):
  - i. kv\_cache is passed as None
  - ii. Populate it as tuple of (keys, values)

2. Init case (Prefill):i. kv\_cache is passed as Noneii. Populate it as tuple of (keys, values)

```
class GPTAttention(nn.Module):
    def forward(self, kv_cache=None):
        ...
    if kv_cache is None:
        kv_cache = (keys, values)
```

- 3. Decode case: kv\_cache has the tuple (keys, values) of past values
  - i. Read past\_keys (values) from kv\_cache. Print shape
  - ii. Print shape of current keys
  - iii. See which dimension to concatenate to create the merged keys and values torch.cat((t1, t2), dim)
  - iv. Continue using the new keys and values tensor

- 3. Decode case: kv\_cache has the tuple (keys, values) of past values
  - Read past\_keys (values) from kv\_cache. Print shape
  - ii. Print shape of current keys
  - iii. See which dimension to concatenate to create the merged keys and values torch.cat((t1, t2), dim)
  - iv. Continue using the new keys and values tensor

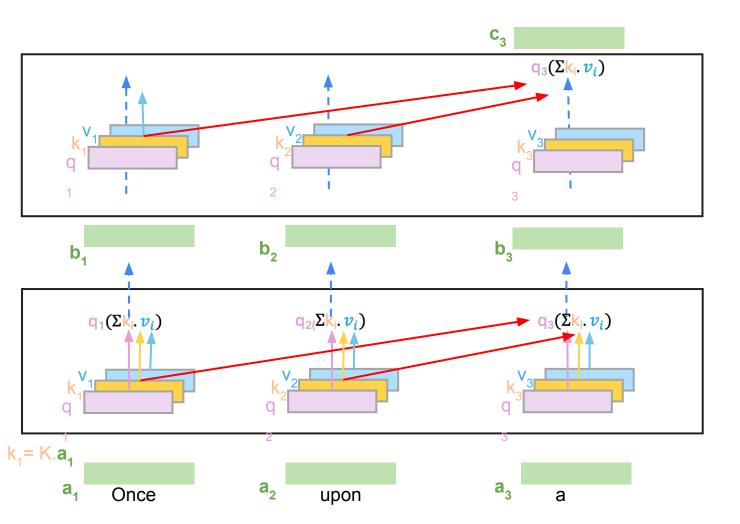
```
class GPTAttention(nn.Module):
    def forward(self, kv_cache=None):
        if kv_cache is not None:
            past_keys, past_values = kv_cache
            keys = torch.cat((past_keys, keys), dim=2)
```

### Let's implement KV Cache (gpt.py)

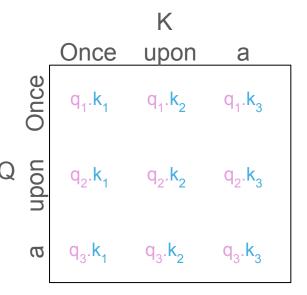
- 4. Return the populated KV cache at all modules:
  - a. Return kv\_cache at GPTAttention
  - b. Return at GPTBlock
  - Return at GPTModel (remember each layer will return its own kv\_cache, so this has to be a list)
  - d. Return at CausalModel back to user

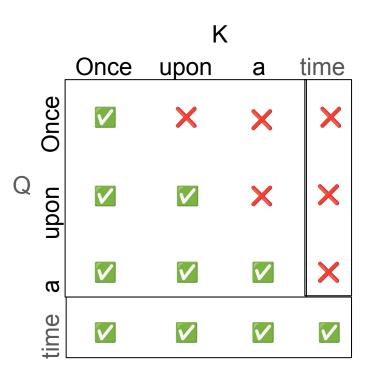
### Let's implement KV Cache

- 5. Changes at generate:
  - a. Separate the prefill and decode calls. Prefill is once, decode is in a loop.
  - Input: We will not pass the entire input to decode, just the output token of the last forward.
  - c. kv\_cache output of one call is passed as input to next call
  - d. Position ids: What was the old position ids? What should the new one be?

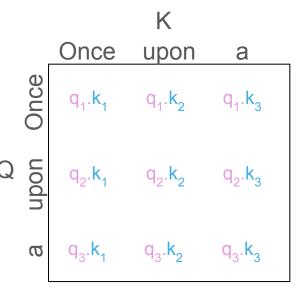


#### Attention mask

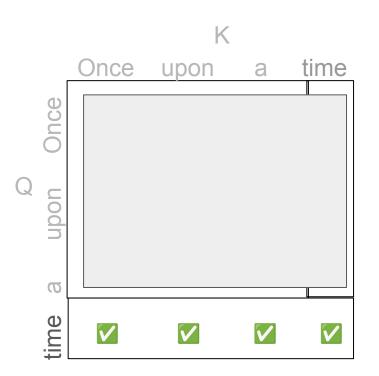




#### Attention mask



With KV Cache



#### Attention mask with KV Caching

- 6. In attn():
  - a. "is\_causal=True" fills the triangular attention mask by default. We want to turn it off for decode
  - b. How to know whether prefill or decode? Use Q.shape and K.shape
  - c. If prefill: Use is\_causal=True
  - d. If decode: Use is causal=False

If everything goes right, your output should match the previous output

## Savings

Run the timed\_generate.py in the solutions file (Copy pytorch\_model.bin)
You should see benefits of KV Cache for larger number of tokens

#### Let's summarize KV Cache

- 1. Why do we cache KV?
- 2. What do we save? Compute or memory?
- 3. How are prefill and decode phases different?
- 4. What can we do if we run out of memory?

# **Batching**

## Batching

#### No batching:

- Generate runs 1 request at a time
- GPU not utilized fully

No batching



**Batching Strategies for LLM Inference** 

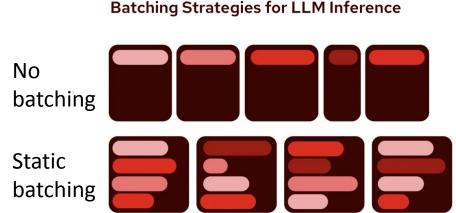
## Batching

#### No batching:

- Generate runs 1 request at a time
- GPU not utilized fully

#### Static batching:

- Request level batching:
- Batch a set of requests
- More parallelization



# Static batching

_	$T_{i}$	Tz	T3	Ty	Ts	T6	To	TB
5	51	Sı	Sı	End				
S	2	S2	S <sub>2</sub>	S2	S2	S2	S2	End
S	3	S <sub>3</sub>	End					
3	54	S <sub>4</sub>	S <sub>4</sub>	S <sub>4</sub>	End			

# Static Batching

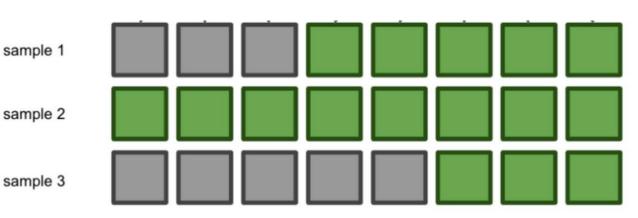
First dimension in inputs is batch size. Previously, we set it to 1.

How to achieve static batching:

- 1. Batch multiple requests together
- 2. How to handle different size requests?

Padding

What are the problems in this approach?



Entire batch waits until T8 Static batching TS End End S2 So Sa So S3 End S SA SA End SA

Iterate until the entire batch is over - low util High latency for finished requests

Not utilized fully

# **Continuous Batching**

#### Iterative/Continuous<sup>1</sup> Batching

- Token level batching
- Replace completed requests with new
- Batch size parameter:
  - Throughput vs latency tradeoff

Any problems?

# Nο batching Static batching Adaptive batching

**Batching Strategies for LLM Inference** 

<sup>1.</sup> Orca: A Distributed Serving System for Transformer-Based Generative Models NSDI '22

<sup>2.</sup> Image source: https://www.redhat.com/en/blog/meet-vllm-faster-more-efficient-llm-inference-and-serving

		Cor	ntinuous	batchiv	^	Next input starts executing immediately		
T,	Tz	T3	Ty	Ts	T6	To	TB	
Sı	Sı	Sı	End	S <sub>6</sub>	S	S <sub>6</sub>	S <sub>6</sub>	
S2	S2	S2	S2	S2	S2	S2	End	
S <sub>3</sub>	S <sub>3</sub>	End	S <sub>5</sub>	.S <sub>5</sub>	S <sub>5</sub>	End	S <sub>7</sub>	
S <sub>4</sub>	S <sub>4</sub>	S₄	S <sub>4</sub>	End	S	S	Sy	

# Let's summarize batching

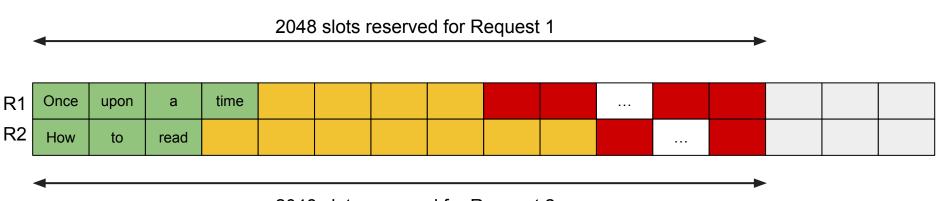
- 1. 2 types of batching
- 2. What do we gain by batching?
- 3. What do we lose?
  - a. Multi-step scheduling

# Paged Attention

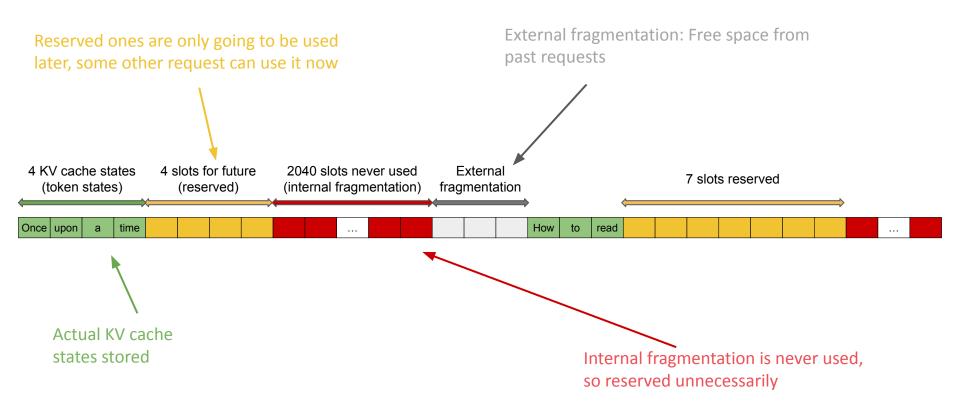
### Memory management of KV Cache

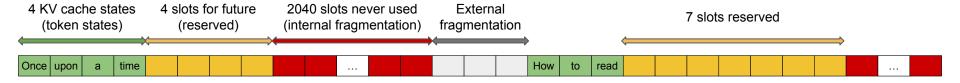
With batching and KV Caching, the KV cache is initialized as a contiguous Tensor for the max sequence length supported by the model

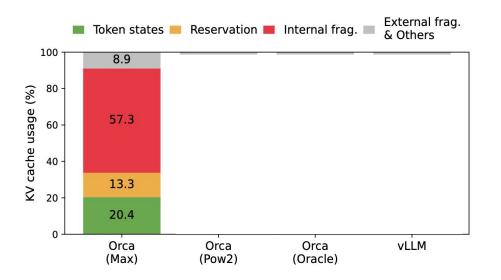
For a long context model (8k - 128k tokens), how much memory does this take?



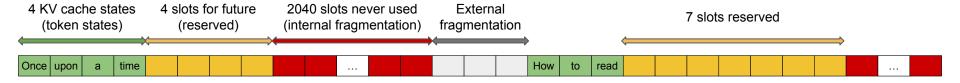
2048 slots reserved for Request 2

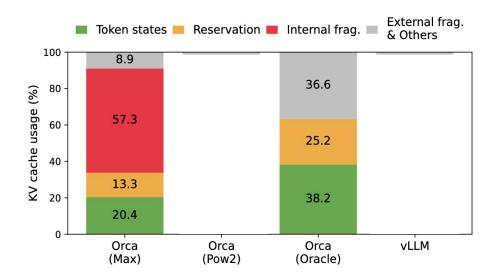




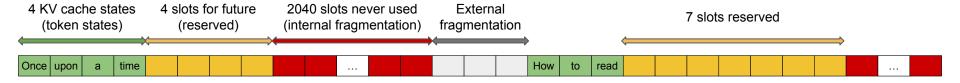


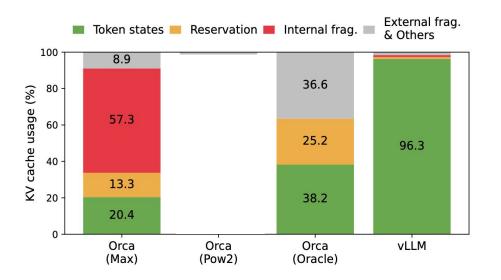
**Figure 2.** Average percentage of memory wastes in different LLM serving systems during the experiment in §6.2.





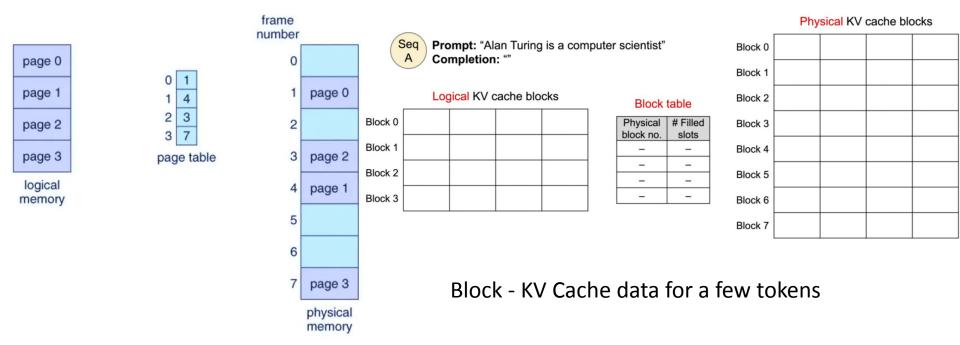
**Figure 2.** Average percentage of memory wastes in different LLM serving systems during the experiment in §6.2.

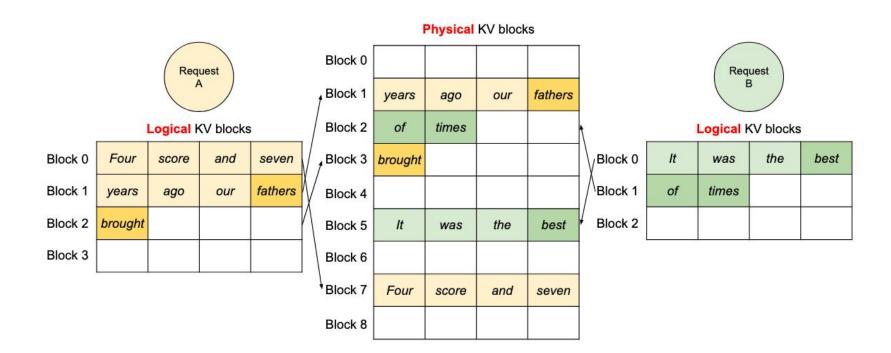


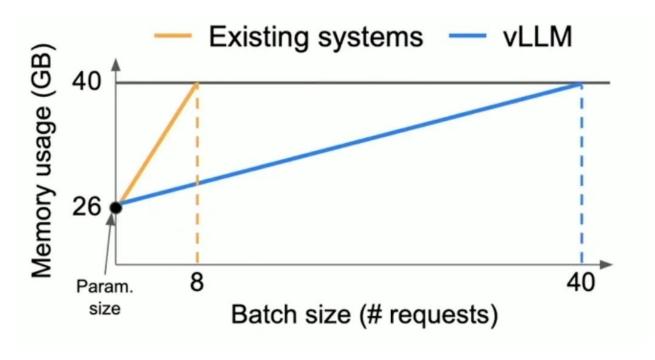


**Figure 2.** Average percentage of memory wastes in different LLM serving systems during the experiment in §6.2.

# How do they do this? Blocks/Pages







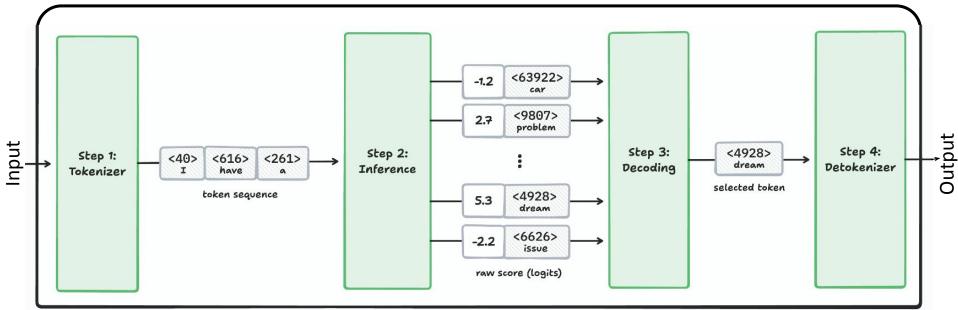
2 - 4x throughput gains

# Decoding strategies

#### What is decoding?

What is the output of inference? Not a token!

Logits: A degree of similarity to each token in the vocabulary



### **Decoding Strategies**

Which token to predict as response?

```
logits = model(inputs, positions)
print(logits.shape)
logits = logits[-1, :]
[x, 50257]
```

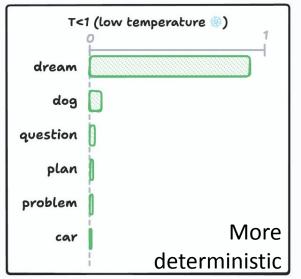
How to convert these logits to a token?

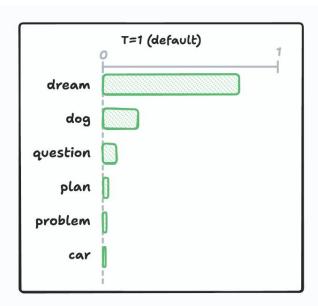
## Softmax

Logit (similarity score) -> Probability score (that adds up to 1)

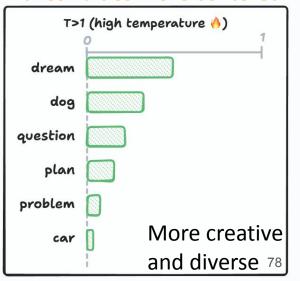
This is done via a softmax normalization function which has a "temperature" parameter which scales the logits value.

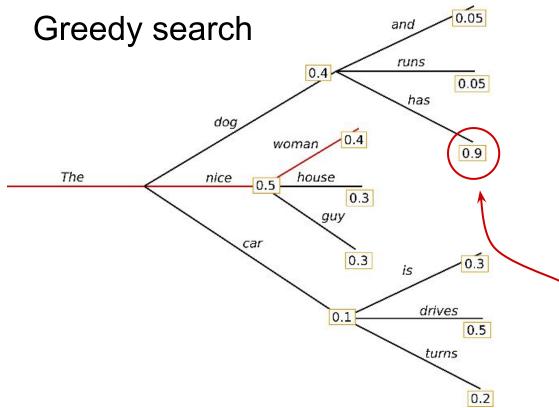
#### Makes values more extreme





#### Makes values more centered





Choose word with highest probability at every step

#### **Problems:**

- 1. It starts repeating itself after some time
- 2. Misses high prob. words hidden behind low prob. words

next\_token = torch.argmax(logits, dim=-1, keepdim=True)

## Beam search

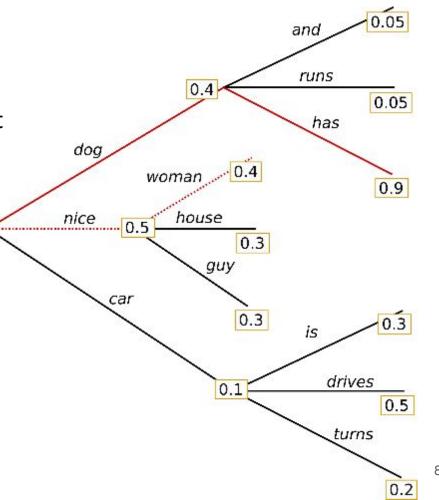
Takes n top probabilities into account

The

Runs n decodes in parallel

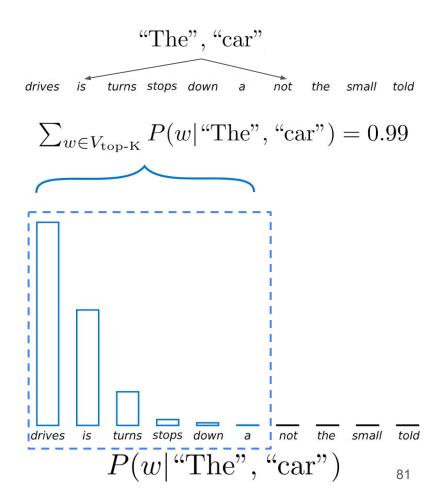
Can still get stuck in a loop

How to avoid repetition -> add randomness



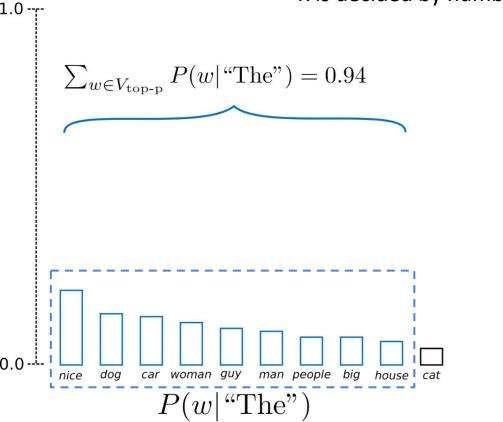
## Top-K Sampling

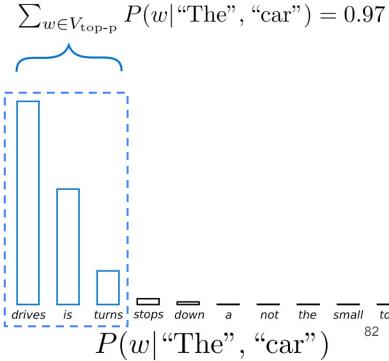
- Randomly pick one of the output tokens
- But this could bias towards the long tail of irrelevant options
- Select only top-K of those and normalize the probability
- How to decide K?



# Top-p sampling

K is decided by number of options that add up to a probability p





not

## Let's see effect of temperature

```
from transformers import pipeline
prompt = "The quick brown fox jumps over the"
generator(prompt, max_new_tokens=20, do_sample=True, temperature=0.7)
generator(prompt, max_new_tokens=20, do_sample=True, temperature=1.5)
generator(prompt, max_new_tokens=20, do_sample=True,
temperature=0.00001)
```

# Research survey

## **Prefix Caching**

If prefix is the SAME for multiple requests, they can share the KV Cache.

Techniques to efficiently match prefix and load its KV Cache

- Hash based
- Radix Tree based

Should be fast and not affect the latency in case of non-match

Where do you think this would be useful?

### KV Cache size reduction

# Problem: KV Cache can grow several times larger than model size

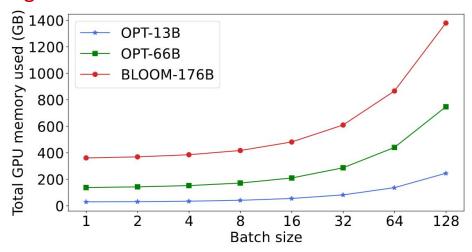


Figure 1. Memory footprint of serving various LLMs with 2K sequence length (input + generated tokens) and half precision (fp16).

How to reduce KV Cache size: Any ideas?

### Sparsification:

Not all values are equally important
 Use attention mask to determine which tokens actually "influence" others

#### Sliding window attention:

Only consider last N tokens

## Speculative Decoding (Draft model)

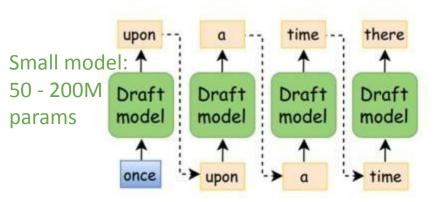
Problem: Auto-regressive nature of generation means only 1 token is decoded in one iteration

Reduce latency without compromising output correctness!

# Speculative Decoding (Draft model)

Problem: Auto-regressive nature of generation means only 1 token is decoded in one iteration

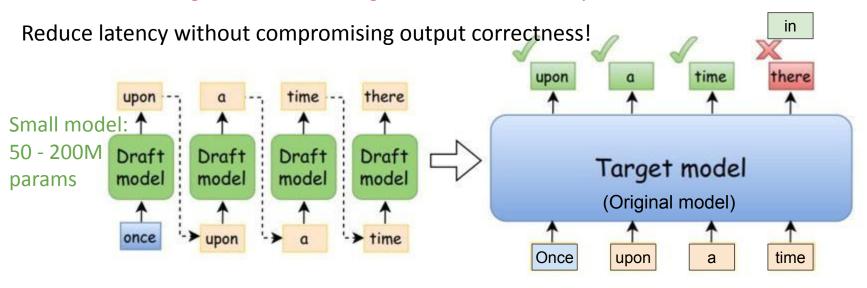
Reduce latency without compromising output correctness!



Step 1: Draft model quickly generates 4-5 tokens (auto-regressively)

## Speculative Decoding (Draft model)

Problem: Auto-regressive nature of generation means only 1 token is decoded in one iteration



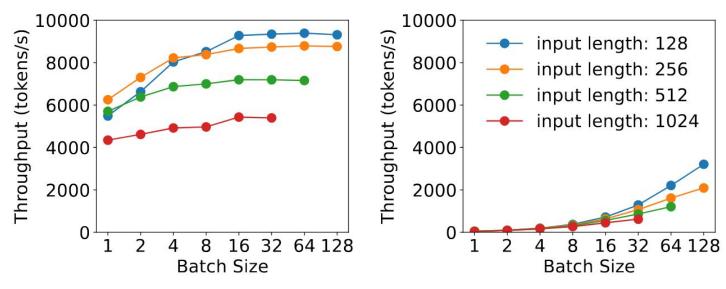
Step 1: Draft model quickly generates 4-5 tokens (auto-regressively)

Step 2: Target model quickly verifies these tokens in parallel. Worst case, re-run from that point

## Prefill-Decode Disaggregation

Problem: Prefill is compute intensive, decode is memory intensive

Not optimal to run both with similar parameters / on the same system

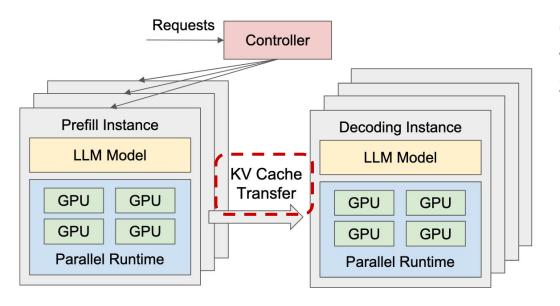


(a) Prefill phase

(b) Decoding phase

## Prefill-Decode Disaggregation

Solution: Run them on different machines



Overlap communication with computation Stream the KV Cache