Introduction to LLMs

Let's level set this...

By systems people...

...for systems people

By systems people... (and the Internet)

...for systems people



Systems

Mechanics

Math

Systems

Mechanics

Math

Flash Attention (NeurIPS)

(powers training)

Paged Attention OSDI

(powers inference)

Tomorrow

Systems

3 - Inference

4 - Dist. Training

Mechanics

1 - LLM Core

2 - Pytorch/GPUS

Today

Math



In terms of GPUs,

we have no GPUs

- \$ pip install torch transformers numpy scipy
- \$ hf download gpt2 model.safetensors
- \$ wget https://huggingface.co/gpt2/resolve/main/pytorch_model.bin

Interactive Hands-on

No AI was used in the making of these slides.

... but can be used in the consumption of these slides.

LLM Core





GenAl models!



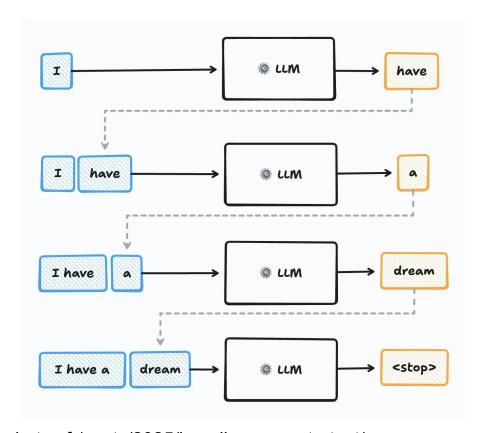




Language Models

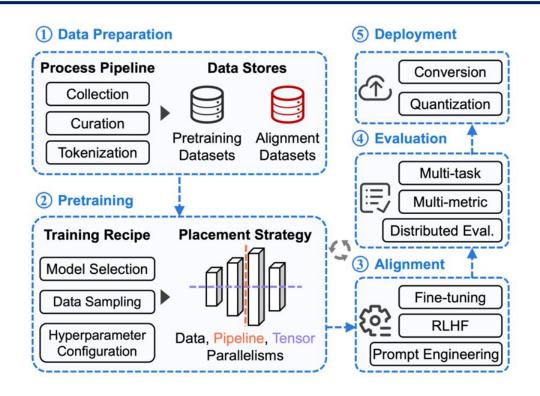
- Probability distribution over words
 - "IIT Bombay is an institute in the city of"
- Enable machines to understand, generate human language
- What all can they do?
 - o QnA
 - Content Generation
 - Text Summarization
 - Sentiment Analysis
 - Conversation
 - Translation
 - Code Generation

Causal Language Modelling or Next Token Prediction

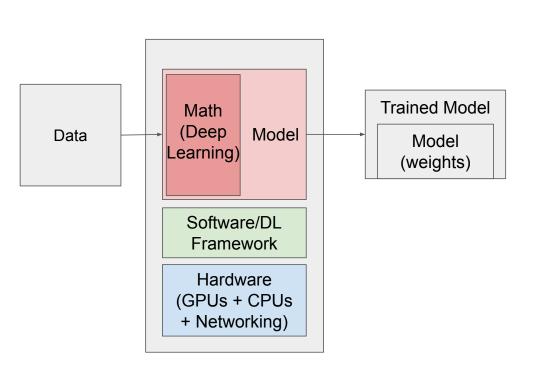


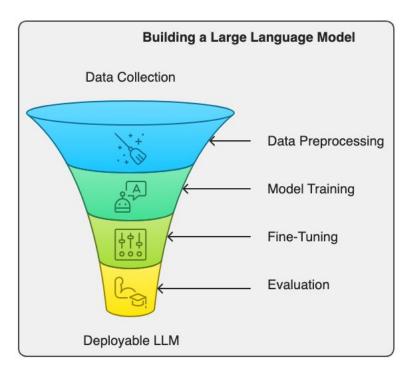
Model Creation Pipeline

- Data Preparation: pretraining data, alignment data
- Pretraining: self-supervised training on large-scale curated data
- Alignment: Prompt Engineering, fine-tuning
- Evaluation: accuracy, fairness, and toxicity
- Deployment: Quntantization, model-parallelism and memory management



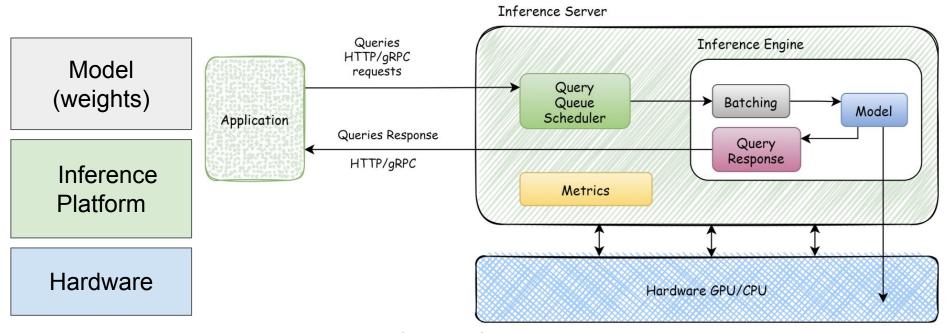
What does it take to train an LLM?





Training and Fine tuning

How do we use the models?



Inference Server + Engine

A partial timeline

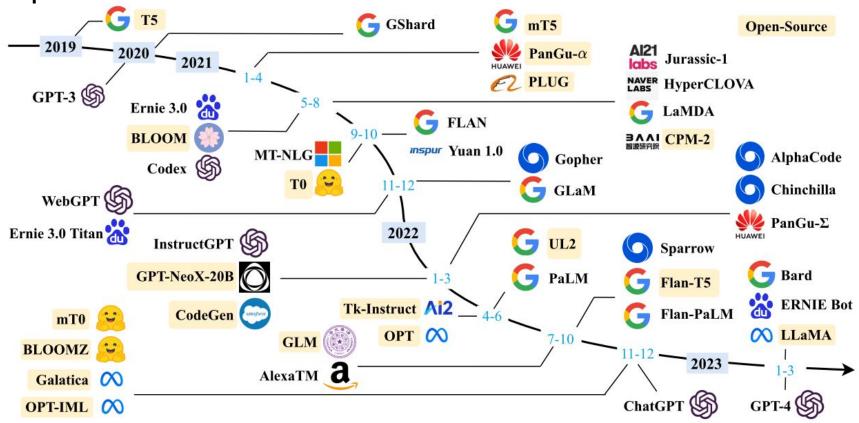
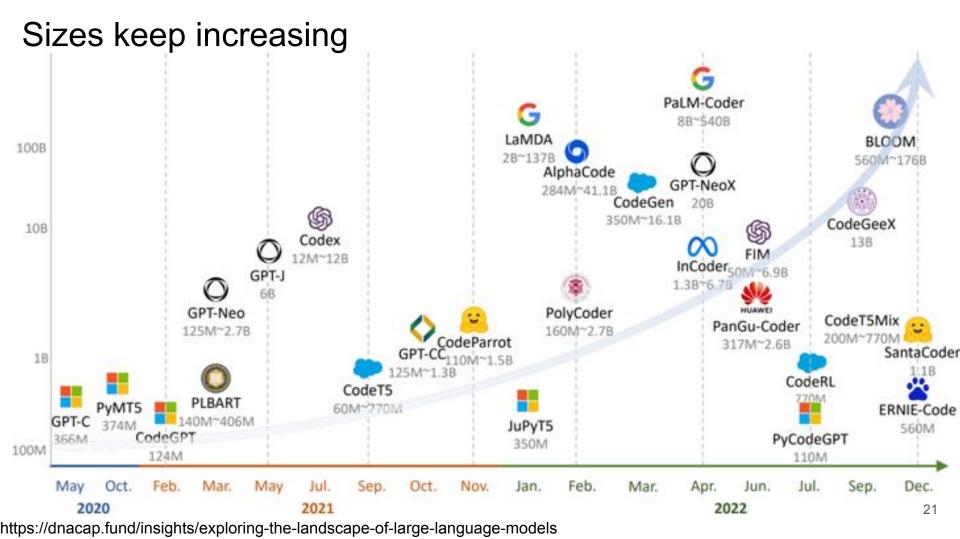
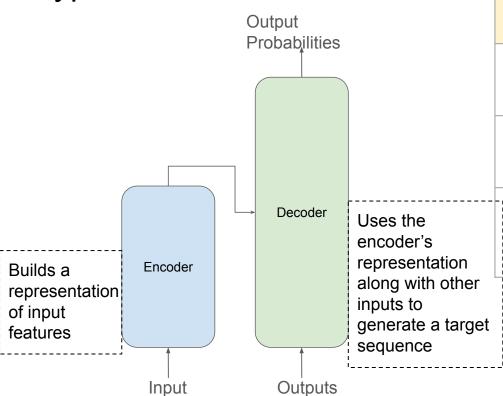


Fig. 1. A timeline of existing large language models (having a size larger than 10B) in recent years. We mark the open-source LLMs in yellow color. https://dnacap.fund/insights/exploring-the-landscape-of-large-language-models



Types of LLM Models



Model type	Used for	Example tasks	Example Models
Encoder-only models	Understanding input	Classification, embeddings	BERT, RoBERTa, SBERT
Decoder-only models	Generative tasks	Text generation	GPT, LLaMA, Mistral
Encoder-decod er models/ sequence-to-se quence models	Generative tasks that require an input	Translation, summarization, QA	T5, BART, MarianMT

Each of these parts can be used independently depending on the task

Transformer: Base architecture of all LLMs today!!

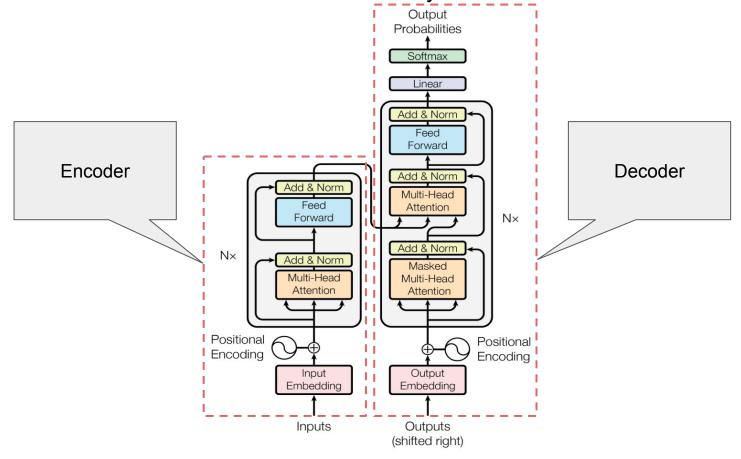


Figure 1: The Transformer - model architecture.

How to use an LLM: First Look

What is hugging face?



- Global hub for Large Language models
 - Spaces (try the models)
 - Models: 1.8M models
 - git clone <modelurl>

https://huggingface.co/

Hands-on: the simplest hello world

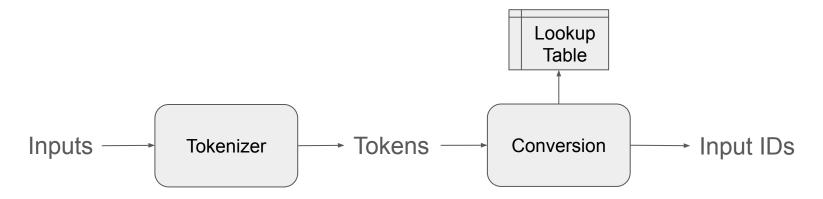
```
from transformers import pipeline
pipe = pipeline(model="gpt2")
```

```
pipe ("Paris is the capital of")
```

- -> HF Pipeline
- -> HF Model
- -> Pytorch Model 34

First, let's look at inputs

Prepare inputs for models:



When running inference on a model, it is necessary to use the same tokenizer to tokenize the inputs as that was used for pretraining the model.

- Types
 - Word-level

We are learning systems for LLMs.

['We', 'are', 'learning', 'systems', 'for', 'LLMs', '.']

Problems

- Separate tokens for different versions of same word - e.g., 'system' and 'systems'
- 2. Results in large vocabulary
- Limiting vocabulary results in bad performance

- Types
 - Character-level

We are learning systems for LLMs.

['W', 'e', ' ', 'a', 'r', 'e', ' ', 'l', 'e', 'a', 'r', 'n', 'i', 'n', 'g', ' ', 's', 'y', 's', 't', 'e', 'm', 's', ' ', 'f', 'o', 'r', ' ', 'L', 'L', 'M', 's', '.']

Problems

- 1. Tokens contain minimal semantic information individually.
- 2. Produces longer sequences, limiting the effective input size of language models.

- Types
 - Subword-level
 - Breaks text into units that are smaller than words but larger than characters.
 - Frequently-used words stored as entire tokens
 - Rare or unknown words are split into smaller, meaningful chunks ("cats" → "cat" + "s").

Tokenizer code snippet

- 1. https://huggingface.co/openai-community/gpt2/blob/main/tokenizer.json
- 2. https://huggingface.co/openai-community/gpt2/blob/main/vocab.json

Tokenizer hands-on

Tokenize the following:

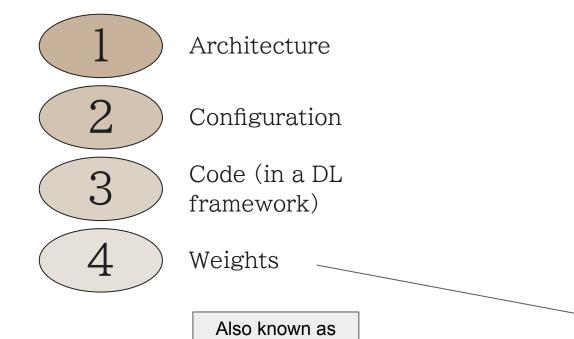
- 1. cat
- 2. superman
- 3. spiderman

Examine the number of tokens

See the tokenizer config.json / vocab.json:

https://huggingface.co/openai-community/gpt2/tree/main

What is a model?



checkpoints

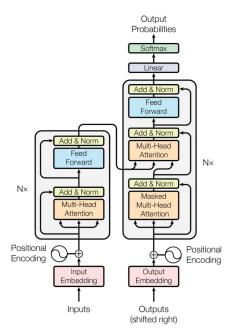
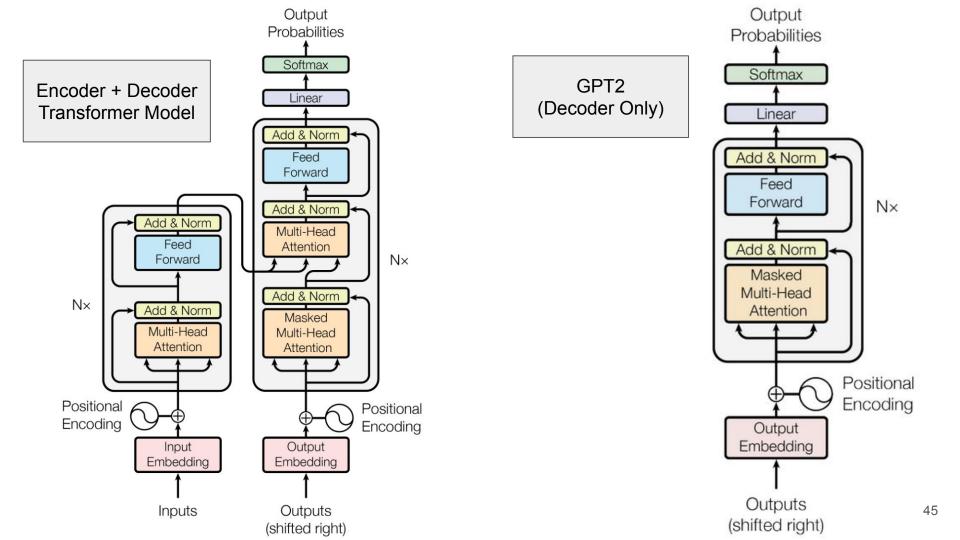


Figure 1: The Transformer - model architecture.

Obtained after extensive training



Try it out: config/dummy model

```
from transformers import
GPT2Config, GPT2LMHeadModel
config = GPT2Config()
  = GPT2LMHeadModel(config)
m1
m1.lm_head.weight
```

```
GPT2Config {
 "activation function": "gelu new",
 "attn_pdrop": 0.1,
 "bos token id": 50256,
 "embd pdrop": 0.1,
 "eos_token_id": 50256,
 "initializer range": 0.02,
GPT2LMHeadModel(
 (transformer): GPT2Model(
  (wte): Embedding(50257, 768)
  (wpe): Embedding(1024, 768)
  (drop): Dropout(p=0.1,
inplace=False)
  (h): ModuleList(
   (0-11): 12 x GPT2Block(
```

Try it out: loading weights

```
from transformers import AutoModelForCausalLM

m2 = AutoModelForCausalLM.from_pretrained("gpt2")

m2.lm_head.weight
```

Check the weights of the model

Try for another model (big model)

```
from transformers import AutoConfig
config = AutoConfig.from_pretrained("Qwen/Qwen3-32B")
print(config)
                                Don't do these!
from transformers import AutoModelForCausalLM
model = AutoModelFromCausalLM.from_config(config)
model
                                                               ~70 GB
from transformers import AutoModelForCausalLM
model = AutoModelFromCausalLM.from_pretrained(config)
```

Let's download the weights - cutting out HuggingFace

wget https://huggingface.co/gpt2/resolve/main/pytorch_model.bin

- O MARIO (W)		operate street	01010 30010 080
☐ README.md ⊚ Safe	8.09 kB 🖳	Add note that this is the smallest versi	almost 3 years ago
☐ config.json	665 Bytes 🖳	Update config.json	over 5 years ago
flax_model.msgpack	498 MB ≪ xet	add gpt2	over 4 years ago
generation_config.json safe	124 Bytes 🖳	Update generation_config.json	over 2 years ago
merges.txt	456 kB 🖳	Update merges.txt	over 6 years ago
☐ model.safetensors	548 MB ≪ xet	Upload model.safetensors with huggingfac	almost 3 years ago
pytorch_model.bin safe mpickle	548 MB ≪ xet	Update pytorch_model.bin	over 6 years ago
rust_model.ot	703 MB ≪ xet	Update rust_model.ot	over 5 years ago
tf_model.h5 💿	498 MB ≪ xet	Update tf_model.h5	almost 6 years ago
¹ tokenizer.json ♥ safe	1.36 MB <u>4</u>	Update tokenizer.json	almost 5 years ago

Loading and examining the State Dict

```
import torch

mw = torch.load("pytorch_model.bin")

list(mw.keys())[:10]

mw['wte.weight']
```

```
['wte.weight',
'wpe.weight',
'h.0.ln_1.weight',
'h.0.ln_1.bias',
'h.0.attn.bias',
'h.0.attn.c_attn.weight',
'h.0.attn.c_proj.weight',
'h.0.attn.c_proj.bias',
'h.0.ln 2.weight']
```

Using these weights

Use nn.Module.load_state_dict

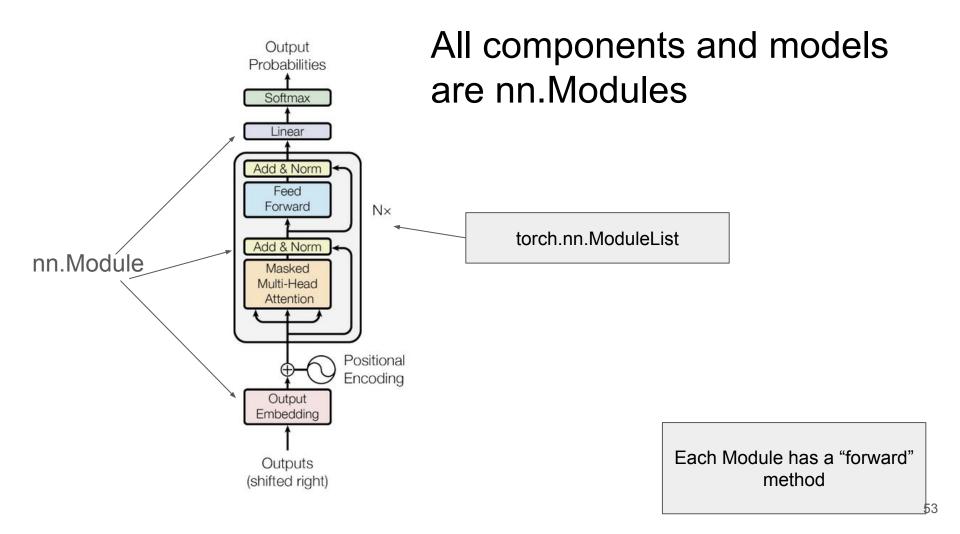
More on this later, for now: We have an hierarchy of modules

```
def extract(sd, prefix):
    prefix = prefix + "."
    return {k.split(prefix)[1]: v for k, v in sd.items()
    if k.startswith(prefix)}
```

HF (Application)

Pytorch prereqs: nn.Module

Pytorch (DL Framework)



Embedding

- Output Embedding
 Outputs (shifted right)
- Transformation of input tokens into a high-dimensional vector space
 - Capture semantic relationship between tokens
 - Dimensionality:
 - Smaller vectors (lower dimensions) more efficient to keep in memory or to process,
 - Bigger vectors (higher dimensions) capture intricate relationships, prone to overfitting.
 output_embeddings.shape

```
import torch.nn as nn
embedding_layer = nn.Embedding(10000, 128)
```

input_indices = torch.tensor([10, 25, 300])

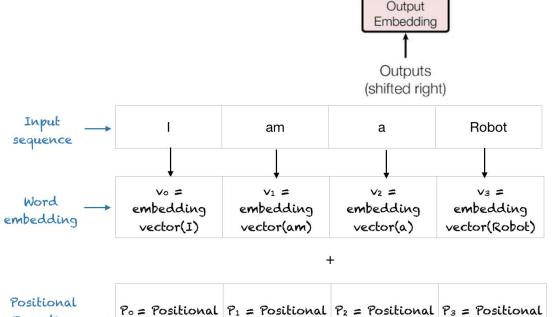
```
output_embeddings =
embedding_layer(input_indices)
```

What will happen if you try to embed "20000"?

torch.Size([3, 128])

Embedding

- Positional Encoding
 - Information about the positions of the tokens added into embeddings to specify the order
 - Eg: RoPE used in Llama





Encoding

Matrix

encoding(I)

vector(I)

encoding(am)

vector(am)

40 = Positional 41 = Positional 42 = Positional encoding(a)

vector(a)

y3 = Positional encoding(Robo

vector(Robot)

Positional Encoding

Embedding - using the real weights

```
1 = torch.nn.Embedding(50257, 768)
1.load_state_dict(extract(mw, "wte"))
<all keys matched successfully>
apple = l(torch.tensor(tok("apple")["input_ids"]))
fruit = l(torch.tensor(tok("fruit")["input_ids"]))
man = l(torch.tensor(tok("man")["input_ids"]))
woman = l(torch.tensor(tok("man")["input_ids"]))
```

```
>>> torch.cosine_similarity(man, apple)
tensor([0.1757], grad_fn=<SumBackward1>)
>>> torch.cosine_similarity(fruit, apple)
tensor([0.3985], grad_fn=<SumBackward1>)
>>> torch.cosine_similarity(man, woman)
tensor([0.5863], grad_fn=<SumBackward1>)
```

Linear

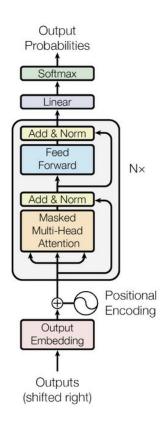
```
linear_layer = torch.nn.Linear(10, 20)
```

```
input = torch.randn(5, 10)
output = linear_layer(input)
```

```
print(linear_layer.weight)
print(linear_layer.bias)
```

$$Y = WX + B$$

GPT2 uses Conv1D instead of linear

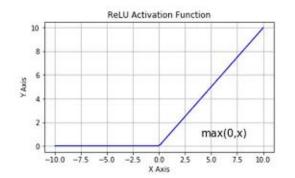


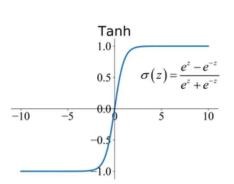
Activations

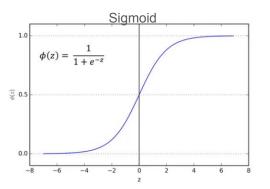
- 1. RELU
- 2. Sigmoid
- 3. Tanh
- 4. GELU

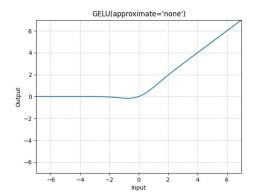
. .

Y = WX + B
Can only capture linear
relations









MLP layer - standalone

x = torch.randn(5, 20)

output = mlp(x)

```
class MLP(nn.Module):
    def __init__(self, dim1, dim2, dim3):
        super().__init__()
        self.fc1 = nn.Linear(dim1, dim2)
        self.activation = nn.GELU()
        self.fc2 = nn.Linear(dim2, dim3)
                                                 20x64
                                                             64x3
    def forward(self, x):
        x = self.fc1(x)
        x = self.activation(x)
        x = self.fc2(x)
        return x
mlp = MLP(20, 64, 3)
```

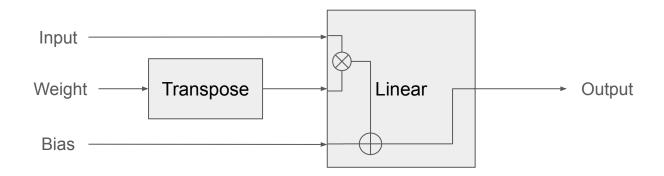
60

MLP layer in GPT2

```
from torch import nn
from transformers.pytorch_utils import Conv1D
class GptMLP(nn.Module):
    def __init__(self):
        super().__init__()
        self.c_fc = Conv1D(768 * 4, 768)
        self.gelu = torch.nn.GELU()
        self.c_proj = Conv1D(768, 768 * 4)
        self.dropout = torch.nn.Dropout(p=0.1, inplace=False)
    def forward(self, x):
                             >>> extract(mw, "h.0.mlp").keys()
        x = self.c_fc(x)
                             dict_keys(['c_fc.weight', 'c_fc.bias',
        x = self.gelu(x)
                             'c_proj.weight', 'c_proj.bias'])
        x = self.c_proj(x)
                             >>> m = GptMLP()
        x = self.dropout(x)
                             >>> m.load_state_dict(extract(mw, "h.0.mlp"))
        return x
                             <all keys matched successfully>
                                                                            61
```

Conv1D

- Convolution, but represents Cross Correlation in this context.
- Think of it as a simple transposed Linear layer



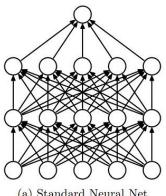
Dropout

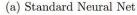
Regularization technique to prevent overfitting/improve generalization

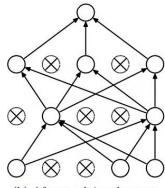
A certain percentage of values are ignored or "dropped out" during each forward and backward pass, making the model's architecture dynamically different for each training batch.

Only used during training

```
m = nn.Dropout(p=0.2)
input = torch.randn(5, 2)
output = m(input)
```

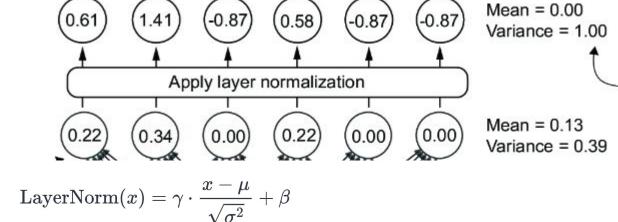






(b) After applying dropout.

LayerNorm

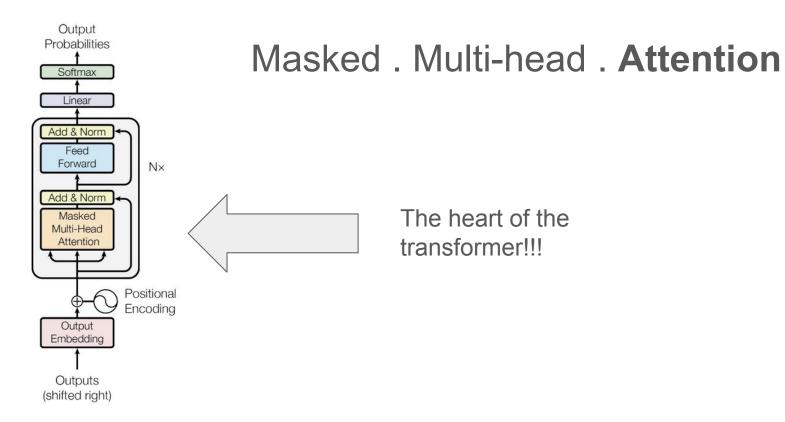


where μ is the mean of x, σ^2 is the variance of x, and γ and β are learnable parameters.

- Normalization technique that standardizes the inputs across the feature dimension of each individual sample
- Mean and variance are calculated across all dimensions of each input sample

$$x = torch.randn(3, 10)$$

layer_norm = torch.nn.LayerNorm(10)



Attention

What attention needs to do?

American shrew mole



Differs based on context

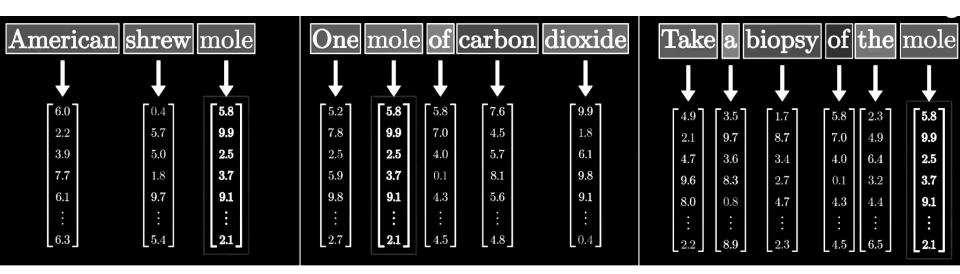
One **mole** of carbon dioxide

 6.02×10^{23}

Take a biopsy of the **mole**



Attention

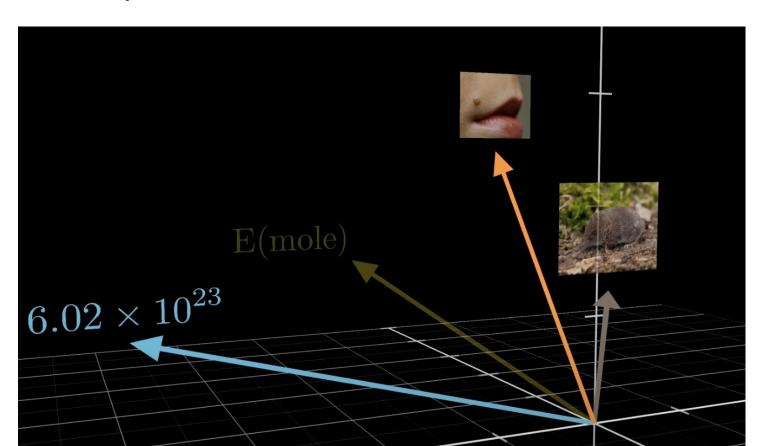


Token Embedding is a table lookup

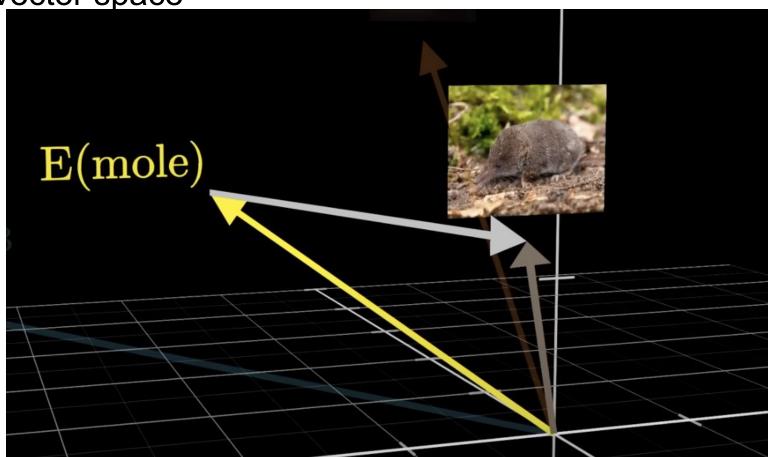
Attention helps pass the context information from other vectors

^{**}Subsequent slides on attention from https://youtu.be/eMlx5fFNoYc?si=9lYxz2a8dyimrm78

In vector space

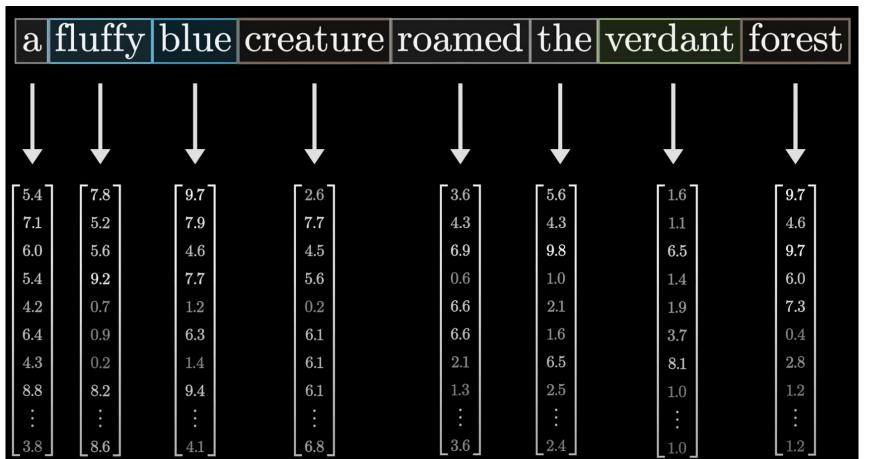


In vector space



a fluffy blue creature roamed the verdant forest



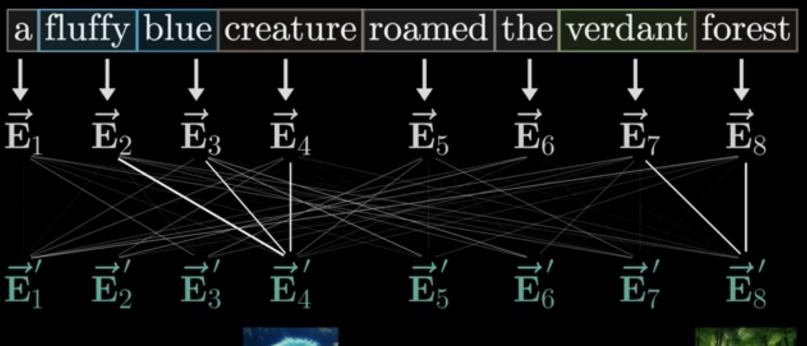






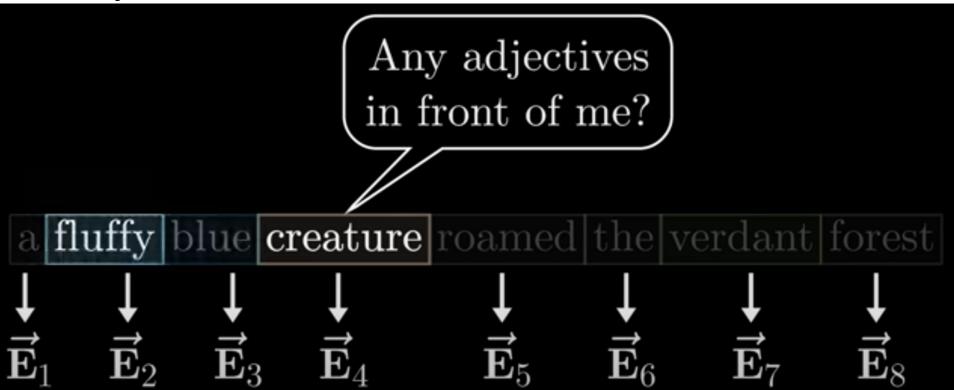




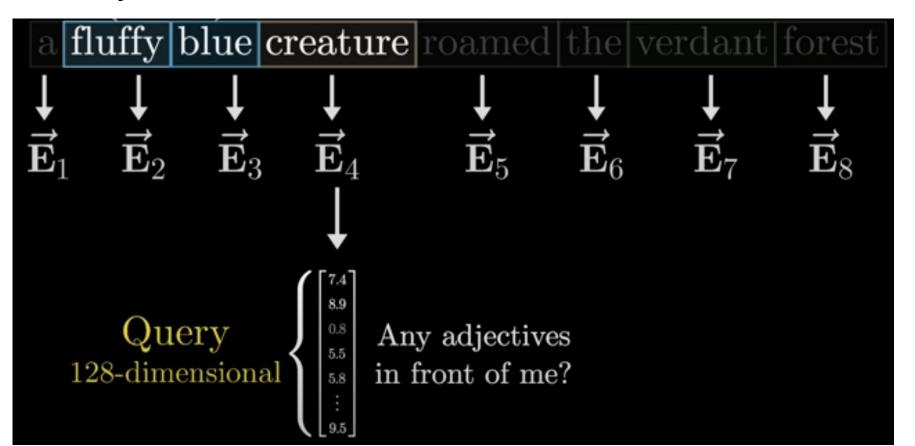




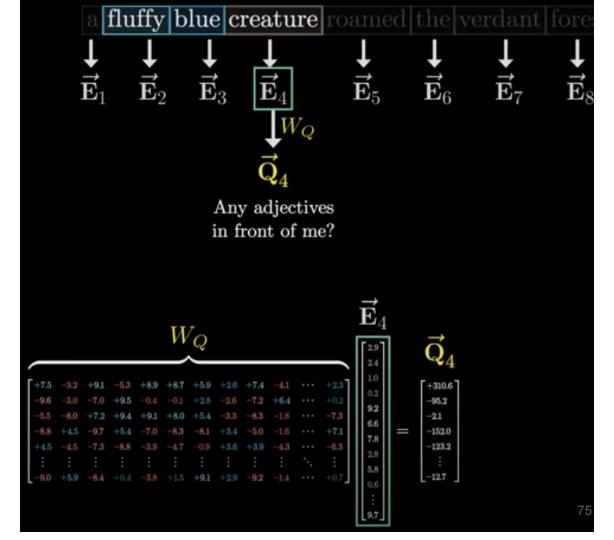
Query



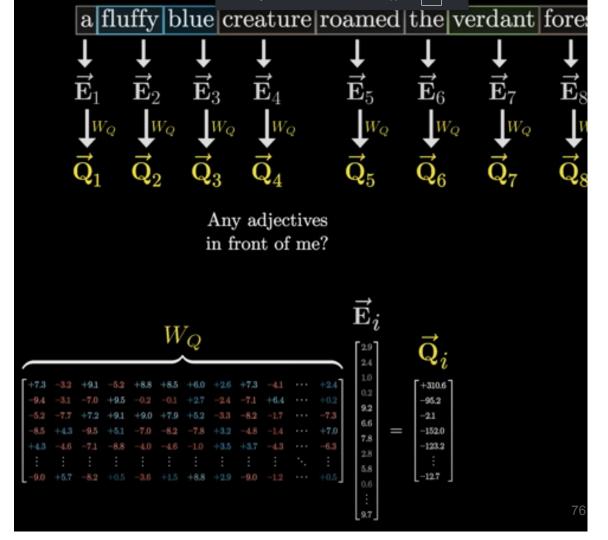
Query



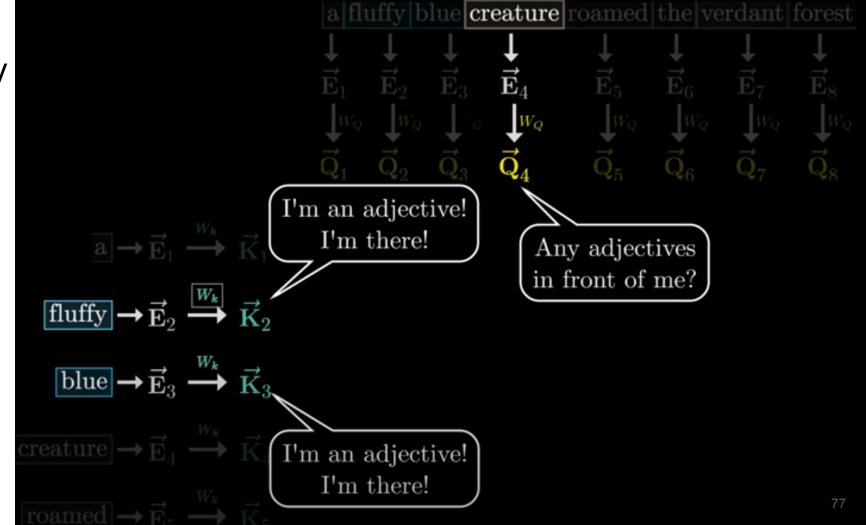
How is query calculated?



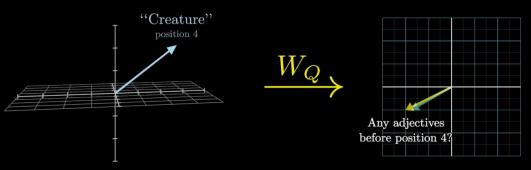
How is query calculated?



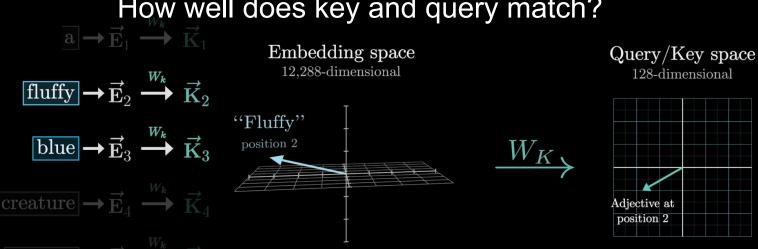
Key







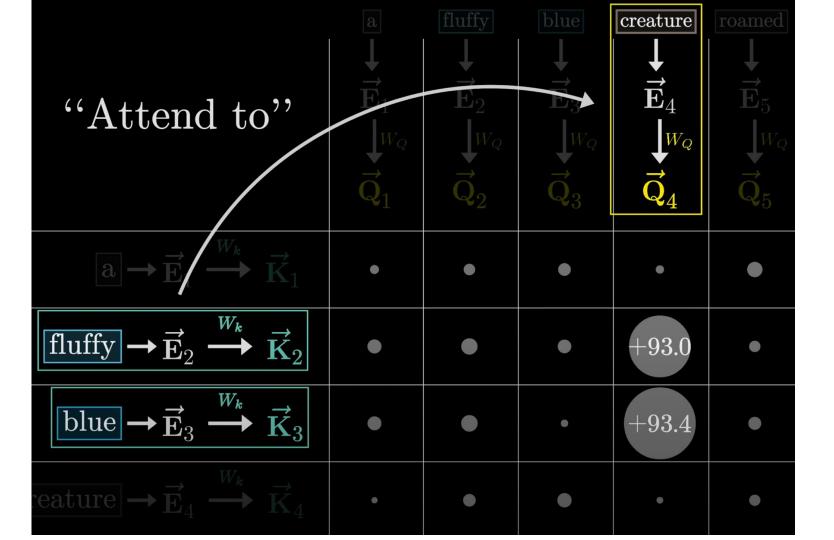
How well does key and query match?

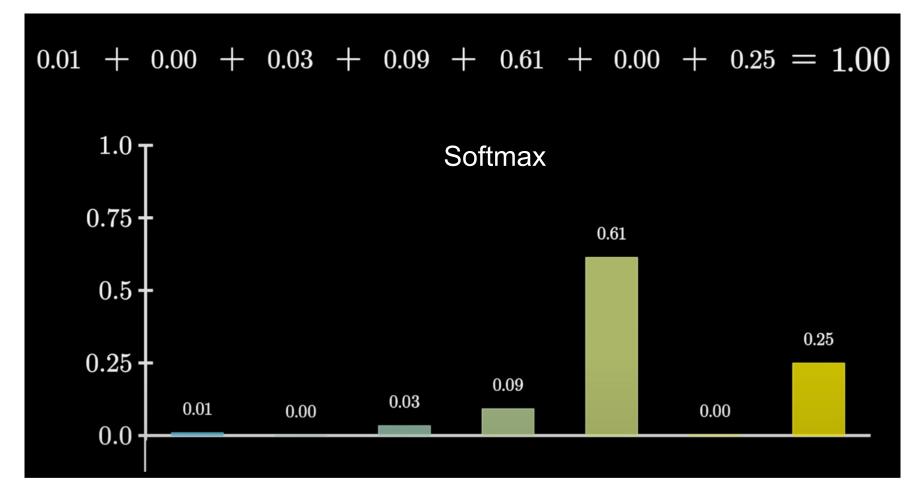


	$egin{aligned} igoddsymbol{igoddown} ar{f E}_1 \ ar{f Q}_1 \end{aligned}$	$\begin{matrix} \mathbf{fluffy} \\ \mathbf{E}_2 \\ \mathbf{Q}_2 \end{matrix}$	$\begin{array}{c c} \mathbf{\overline{blue}} \\ \mathbf{\overline{E}}_3 \\ \mathbf{\overline{Q}}_3 \end{array}$	$egin{array}{c} ext{creature} \ ext{$\stackrel{ ext{$\downarrow$}}{ ext{$\downarrow$}}$}_4 \ ext{$\stackrel{ ext{$\downarrow$}}{ ext{$\downarrow$}}$}_{Q_4} \ ext{$\stackrel{ ext{$\downarrow$}}{ ext{$\downarrow$}}}_4 \ ext{$\stackrel{ ext{\downarrow}}{ ext{\downarrow}}}_{Q_4} \ ext{$\stackrel{ ext{$\downarrow$}}{ ext{$\downarrow$}$	$egin{array}{c} ext{roamed} \ ext{$\stackrel{ ext{$\downarrow$}}{ ext{$\downarrow$}}$} \ ext{$\stackrel{ ext{$\downarrow$}}{ ext{$\downarrow$}}$} \ ext{$\stackrel{ ext{$\downarrow$}}{ ext{$\downarrow$}}$} \ ext{$\stackrel{ ext{$\downarrow$}}{ ext{$\downarrow$}}} \ ext{$\stackrel{ ext{\downarrow}}{ ext{\downarrow}}}$	$egin{array}{c} ext{the} \ ext{$fille{f E}$}_6 \ ext{$f Q$}_6 \ \end{array}$	$\begin{array}{c} \mathbf{verdant} \\ \mathbf{E}_7 \\ \mathbf{Q}_7 \end{array}$	$\begin{matrix} \bullet \\ \downarrow \\ \vec{\mathbf{E}}_8 \\ \downarrow W_Q \\ \vec{\mathbf{Q}}_8 \end{matrix}$	
$\mathbf{a} \to \vec{\mathbf{E}}_1 \xrightarrow{W_k} \vec{\mathbf{K}}_1$	$ec{\mathbf{K}}_1\!\cdot\!ec{\mathbf{Q}}_1$	$ec{\mathbf{K}}_1\!\cdot\!ec{\mathbf{Q}}_2$	$ec{\mathbf{K}}_1\!\cdot\!ec{\mathbf{Q}}_3$	$ec{\mathbf{K}}_1\!\cdot\!ec{\mathbf{Q}}_4$	$ec{\mathbf{K}}_1\!\cdot\!ec{\mathbf{Q}}_5$	$ec{\mathbf{K}}_1\!\cdot\!ec{\mathbf{Q}}_6$	$ec{\mathbf{K}}_1\!\cdot\!ec{\mathbf{Q}}_7$	$ec{\mathbf{K}}_1\!\cdot\!ec{\mathbf{Q}}_8$	
$ \underbrace{\text{fluffy}} \to \vec{\mathbf{E}}_2 \xrightarrow{W_k} \vec{\mathbf{K}}_2 $	$ec{\mathbf{K}}_2\!\cdot\!ec{\mathbf{Q}}_1$	$ec{\mathbf{K}}_2\!\cdot\!ec{\mathbf{Q}}_2$	$ec{\mathbf{K}}_2\!\cdot\!ec{\mathbf{Q}}_3$	$ec{\mathbf{K}}_2\!\cdot\!ec{\mathbf{Q}}_4$	$ec{\mathbf{K}}_2\!\cdot\!ec{\mathbf{Q}}_5$	$ec{\mathbf{K}}_2\!\cdot\!ec{\mathbf{Q}}_6$	$ec{\mathbf{K}}_2\!\cdot\!ec{\mathbf{Q}}_7$	$ec{\mathbf{K}}_2\!\cdot\!ec{\mathbf{Q}}_8$	
$ \begin{array}{c} $	$ec{\mathbf{K}}_3\!\cdot\!ec{\mathbf{Q}}_1$	$\vec{\mathbf{K}}_3 \cdot \vec{\mathbf{Q}}_2$	$\vec{\mathbf{K}}_3\!\cdot\!\vec{\mathbf{Q}}_3$	$ec{\mathbf{K}}_3\!\cdot\!ec{\mathbf{Q}}_4$	$ec{\mathbf{K}}_3\!\cdot\!ec{\mathbf{Q}}_5$	$ec{\mathbf{K}}_3\!\cdot\!ec{\mathbf{Q}}_6$	$\vec{\mathbf{K}}_3\!\cdot\!\vec{\mathbf{Q}}_7$	$\vec{\mathbf{K}}_3\!\cdot\!\vec{\mathbf{Q}}_8$	
$\boxed{\text{creature}} \to \vec{\mathbf{E}}_4 \stackrel{W_k}{\longrightarrow} \vec{\mathbf{K}}_4$	$ec{\mathbf{K}}_4\!\cdot\!ec{\mathbf{Q}}_1$	$ec{\mathbf{K}}_4\!\cdot\!ec{\mathbf{Q}}_2$	$ec{\mathbf{K}}_4\!\cdot\!ec{\mathbf{Q}}_3$	$ec{\mathbf{K}}_4\!\cdot\!ec{\mathbf{Q}}_4$	$ec{\mathbf{K}}_4\!\cdot\!ec{\mathbf{Q}}_5$	$ec{\mathbf{K}}_4\!\cdot\!ec{\mathbf{Q}}_6$	$ec{\mathbf{K}}_4\!\cdot\!ec{\mathbf{Q}}_7$	$ec{\mathbf{K}}_4\!\cdot\!ec{\mathbf{Q}}_8$	
$ \overrightarrow{\text{roamed}} \to \overrightarrow{\mathbf{E}}_5 \xrightarrow{W_k} \overrightarrow{\mathbf{K}}_5 $	$ec{\mathbf{K}}_5\!\cdot\!ec{\mathbf{Q}}_1$	$ec{\mathbf{K}}_5\!\cdot\!ec{\mathbf{Q}}_2$	$ec{\mathbf{K}}_5\!\cdot\!ec{\mathbf{Q}}_3$	$ec{\mathbf{K}}_5\!\cdot\!ec{\mathbf{Q}}_4$	$ec{\mathbf{K}}_5\!\cdot\!ec{\mathbf{Q}}_5$	$ec{\mathbf{K}}_5\!\cdot\!ec{\mathbf{Q}}_6$	$\vec{\mathbf{K}}_5\!\cdot\!\vec{\mathbf{Q}}_7$	$ec{\mathbf{K}}_5\!\cdot\!ec{\mathbf{Q}}_8$	
$ \begin{array}{c} \text{the} \rightarrow \vec{\mathbf{E}}_6 \xrightarrow{W_k} \vec{\mathbf{K}}_6 \end{array} $	$ec{\mathbf{K}}_6\!\cdot\!ec{\mathbf{Q}}_1$	$ec{\mathbf{K}}_6\!\cdot\!ec{\mathbf{Q}}_2$	$ec{\mathbf{K}}_6\!\cdot\!ec{\mathbf{Q}}_3$	$ec{\mathbf{K}}_6\!\cdot\!ec{\mathbf{Q}}_4$	$ec{\mathbf{K}}_6\!\cdot\!ec{\mathbf{Q}}_5$	$ec{\mathbf{K}}_6\!\cdot\!ec{\mathbf{Q}}_6$	$ec{\mathbf{K}}_6\!\cdot\!ec{\mathbf{Q}}_7$	$ec{\mathbf{K}}_6\!\cdot\!ec{\mathbf{Q}}_8$	
$\overrightarrow{\mathrm{verdant}} \to \overrightarrow{\mathrm{E}}_7 \xrightarrow{W_k} \overrightarrow{\mathrm{K}}_7$	$\vec{\mathbf{K}}_7\!\cdot\!\vec{\mathbf{Q}}_1$	$ec{\mathbf{K}}_7\!\cdot\!ec{\mathbf{Q}}_2$	$\vec{\mathbf{K}}_7\!\cdot\!\vec{\mathbf{Q}}_3$	$ec{\mathbf{K}}_7\!\cdot\!ec{\mathbf{Q}}_4$	$ec{\mathbf{K}}_7\!\cdot\!ec{\mathbf{Q}}_5$	$ec{\mathbf{K}}_7\!\cdot\!ec{\mathbf{Q}}_6$	$\vec{\mathbf{K}}_7 \cdot \vec{\mathbf{Q}}_7$	$ec{\mathbf{K}}_7\!\cdot\!ec{\mathbf{Q}}_8$	
$\boxed{\text{forest} \rightarrow \vec{\mathbf{E}}_8 \stackrel{W_k}{\longrightarrow} \vec{\mathbf{K}}_8}$	$\vec{\mathbf{K}}_8 \cdot \vec{\mathbf{Q}}_1$	$ec{\mathbf{K}}_8\!\cdot\!ec{\mathbf{Q}}_2$	$\vec{\mathbf{K}}_8 \cdot \vec{\mathbf{Q}}_3$	$ec{\mathbf{K}}_8\!\cdot\!ec{\mathbf{Q}}_4$	$\vec{\mathbf{K}}_8 \cdot \vec{\mathbf{Q}}_5$	$\vec{\mathbf{K}}_8 \cdot \vec{\mathbf{Q}}_6$	$\vec{\mathbf{K}}_8 \cdot \vec{\mathbf{Q}}_7$	$ec{\mathbf{K}}_8\!\cdot\!ec{\mathbf{Q}}_8$	

	a	fluffy	blue	creature	[roamed]	the	verdant	forest	
	↓	→	→	↓	↓	→	\bigvee_{\rightarrow}	<u></u>	
	$ec{ ilde{f E}}_1$	$ec{ extbf{E}}_2$	_	$ec{ec{f E}}_4$	$ec{ extbf{E}}_{5}$	$ec{ extbf{E}}_{6}$	$ec{\mathbf{E}}_7$	$ec{\mathbf{E}}_{8}$	
	W_Q	W_Q	W_Q	W_Q	$\overline{\bigcup_{W_Q}}$	\bigcup_{W_Q}	$ec{ec{\mathbf{Q}}_{7}}$	W_Q	
	$ec{ ext{Q}}_1$	$ec{ extbf{Q}}_2$	$ec{ extbf{Q}}_3$	$ec{ extbf{Q}}_4$	$ec{\mathbf{Q}}_{5}$	$ec{ec{Q}}_{6}$	$ec{ ext{Q}}_7$	$ec{\mathbf{Q}}_{8}$	
$\mathbf{a} \rightarrow \vec{\mathbf{E}}_1 \stackrel{\mathbf{W}_k}{\longrightarrow} \vec{\mathbf{K}}_1$	$\vec{\mathbf{K}}_{1}ullet \vec{\mathbf{Q}}_{1}$	$ec{\mathbf{K}}_1 ullet ec{\mathbf{Q}}_2$	$\vec{\mathbf{K}}_1$ $ullet$ $\vec{\mathbf{Q}}_3$	$ec{\mathbf{K}}_1ullet ec{\mathbf{Q}}_4$	$ec{\mathbf{K}}_{\mathbf{l}}\mathbf{ullet}ec{\mathbf{Q}}_{5}$	$ec{\mathbf{K}}_1ullet ec{\mathbf{Q}}_6$	$ec{\mathbf{K}}_1ulletec{\mathbf{Q}}_7$	$ec{\mathbf{K}}_1 ullet ec{\mathbf{Q}}_8$	
$\boxed{\text{fluffy}} \rightarrow \vec{\mathbf{E}}_2 \stackrel{W_k}{\longrightarrow} \vec{\mathbf{K}}_2$	$ec{\mathbf{K}}_2 lackbox{\mathbf{Q}}_1$	$ec{\mathbf{K}}_2 lackbox{\mathbf{Q}}_2$	$ec{\mathbf{K}}_2 ullet ec{\mathbf{Q}}_3$	$\vec{K}_2 \cdot \vec{Q}_4$	$ec{\mathbf{K}}_2 ullet ec{\mathbf{Q}}_5$	$ec{\mathbf{K}}_2$ $ullet$ $ec{\mathbf{Q}}_6$	$ec{\mathbf{K}}_2ullet ec{\mathbf{Q}}_7$	$ec{\mathbf{K}}_2 ullet ec{\mathbf{Q}}_8$	
$\boxed{\text{blue}} \rightarrow \vec{\mathbf{E}}_3 \stackrel{W_k}{\longrightarrow} \vec{\mathbf{K}}_3$	$ec{\mathbf{K}}_3 led ec{\mathbf{Q}}_1$	$ec{\mathbf{K}}_3 lackbox{\mathbf{Q}}_2$	$ec{\mathbf{K}}_3ullet ec{\mathbf{Q}}_3$	$\vec{\mathbf{K}}_3 \cdot \vec{\mathbf{Q}}$	$ec{\mathbf{K}}_3ullet ec{\mathbf{Q}}_5$	$ec{\mathbf{K}}_3ullet ec{\mathbf{Q}}_6$	$ec{\mathbf{K}}_3ullet ec{\mathbf{Q}}_7$	$ec{\mathbf{K}}_3 led{\mathbf{Q}}_8$	
$\overrightarrow{\mathbf{E}}_4 \xrightarrow{W_k} \overrightarrow{\mathbf{K}}_4$	$ec{\mathbf{K}}_4ullet ec{\mathbf{Q}}_1$	$ec{\mathbf{K}}_{4}oldsymbol{ar{\mathbf{Q}}}_{2}$	$ec{\mathbf{K}}_{4}oldsymbol{ar{\mathbf{Q}}}_{3}$	$ec{\mathbf{K}}_4ullet ec{\mathbf{Q}}_4$	$ec{\mathbf{K}}_{4}ullet ec{\mathbf{Q}}_{5}$	$ec{\mathbf{K}}_4\mathbf{O}ec{\mathbf{Q}}_6$	$ec{\mathbf{K}}_4ulletec{\mathbf{Q}}_7$	$ec{\mathbf{K}}_{4}ullet ec{\mathbf{Q}}_{8}$	
$\boxed{\text{roamed}} \rightarrow \vec{\mathbf{E}}_5 \stackrel{W_k}{\longrightarrow} \vec{\mathbf{K}}_5$	$ec{\mathbf{K}}_{5}ullet ec{\mathbf{Q}}_{1}$	$ec{\mathbf{K}}_5\mathbf{C}ec{\mathbf{Q}}_2$	$ec{\mathbf{K}}_{5}$ $ullet ec{\mathbf{Q}}_{3}$	$ec{\mathbf{K}}_{5}$ $\mathbf{Q}\mathbf{ar{Q}}_{4}$	$ec{\mathbf{K}}_{5}ullet ec{\mathbf{Q}}_{5}$	$ec{\mathbf{K}}_{5}^{oldsymbol{\circ}} ec{\mathbf{Q}}_{6}$	$ec{\mathbf{K}}_5 ullet ec{\mathbf{Q}}_7$	$ec{\mathbf{K}}_{5}$ $oldsymbol{ar{\mathbf{Q}}}_{8}$	
$ \begin{array}{c} \text{the} \rightarrow \vec{\mathbf{E}}_6 \xrightarrow{w_k} \vec{\mathbf{K}}_6 \end{array} $	$ec{\mathbf{K}}_{6}$ $oldsymbol{\mathbf{Q}}_{1}$	$ec{\mathbf{K}}_{6}ullet ec{\mathbf{Q}}_{2}$	$\vec{\mathbf{K}}_6ullet \vec{\mathbf{Q}}_3$	$ec{\mathbf{K}}_6$ $oldsymbol{\mathbf{Q}}_4$	$ec{\mathbf{K}}_6 ullet ec{\mathbf{Q}}_5$	$ec{\mathbf{K}}_6$ $oldsymbol{\mathbf{Q}}_6$	$ec{\mathbf{K}}_{6}$ \mathbf{Q}_{7}	$\vec{\mathbf{K}}_6 ullet \vec{\mathbf{Q}}_8$	
$\overrightarrow{\mathbf{E}}_7 \xrightarrow{W_k} \overrightarrow{\mathbf{K}}_7$	$\vec{\mathrm{K}}_{7}$ $oldsymbol{ec{\mathrm{Q}}}_{1}$	$ec{\mathbf{K}}_7$ $oldsymbol{ar{\mathbf{Q}}}_2$	$\vec{\mathbf{K}}_7 ullet \vec{\mathbf{Q}}_3$	$ec{\mathbf{K}}_7ullet ec{\mathbf{Q}}_4$	$ec{\mathbf{K}}_7ullet ec{\mathbf{Q}}_5$	$ec{\mathbf{K}}_7ullet ec{\mathbf{Q}}_6$	$\vec{\mathbf{K}}_7$ \mathbf{Q}_7	$\vec{\mathbf{k}}_7 \cdot \vec{\mathbf{Q}}$	
forest $\rightarrow \vec{\mathbf{E}}_8 \xrightarrow{W_k} \vec{\mathbf{K}}_8$	$ec{\mathbf{K}}_8$ o $ec{\mathbf{Q}}_1$	$ec{\mathbf{K}}_8$ o $ec{\mathbf{Q}}_2$	$\vec{\mathbf{K}}_8 ullet \vec{\mathbf{Q}}_3$	$ec{\mathbf{K}}_8 ullet ec{\mathbf{Q}}_4$	$ec{\mathbf{K}}_8 ullet ec{\mathbf{Q}}_5$	$ec{\mathbf{K}}_8$ o $ec{\mathbf{Q}}_6$	$\vec{\mathbf{K}}_8$ \mathbf{O} $\vec{\mathbf{Q}}_7$	$ec{\mathbf{K}}_8 \circ ec{\mathbf{Q}}_8$	

	a	fluffy	blue	creature	roamed	the	verdant	forest	
	$egin{array}{c} lacksquare & lacksquare $	$ec{\mathbf{E}}_2$	$ec{f E}_3$	$egin{array}{c} oldsymbol{\downarrow} \ ec{\mathbf{E}}_4 \end{array}$	$ec{ ilde{\mathbf{E}}}_{5}$	$ec{f E}_6$	$ec{ ilde{\mathbf{E}}}_7$		
	W_Q	$ec{f Q}_{f Q_2}^{W_Q}$	$ec{f Q}_{f 3}^{W_Q}$	$ec{ec{\mathbf{Q}}_{4}}^{W_Q}$	$ec{f Q}_5^{W_Q}$	$ec{f Q}_6^{W_Q}$	$ec{ec{Q}}_{7}^{W_{Q}}$	$\bigvee_{i=1}^{W_Q}$	
	\mathbf{Q}_1	$\dot{\mathbf{Q}}_2$	\mathbf{Q}_3	\mathbf{Q}_4	\mathbf{Q}_{5}	\mathbf{Q}_{6}	$\dot{\mathbf{Q}}_7$	$ec{\mathbf{Q}}_{8}$	
$\mathbf{a} \rightarrow \vec{\mathbf{E}}_1 \stackrel{\mathbf{W}_k}{\longrightarrow} \vec{\mathbf{K}}_1$	+0.7	-83.7	-24.7	-27.8	-5.2	-89.3	-45.2	-36.1	
$\boxed{\text{fluffy}} \rightarrow \vec{\mathbf{E}}_2 \stackrel{W_k}{\longrightarrow} \vec{\mathbf{K}}_2$	-73.4	+2.9	-5.4	+93.0	-48.2	-87.3	-49.7	+7.8	
	-53.4	-5.7	+1.8	+93.4	-55.6	-56.0	-26.1	-62.1	
$ \overrightarrow{\mathbf{E}}_4 \xrightarrow{W_k} \vec{\mathbf{K}}_4 $	-21.5	-29.7	-56.1	+4.9	-32.4	-92.3	-9.5	-28.1	
$\boxed{\text{roamed}} \to \vec{\mathbf{E}}_5 \stackrel{W_k}{\longrightarrow} \vec{\mathbf{K}}_5$	-20.1	-40.9	-87.8	-55.4	+0.6	-64.7	-96.7	-18.9	
$ \begin{array}{c} \text{the} \rightarrow \vec{\mathbf{E}}_6 \xrightarrow{w_k} \vec{\mathbf{K}}_6 \end{array} $	-87.9	-33.3	-22.6	-31.4	+5.5	+0.6	-4.6	-96.8	
$\underbrace{\text{verdant}} \to \vec{\mathbf{E}}_7 \xrightarrow{W_k} \vec{\mathbf{K}}_7$	-41.2	-55.5	-42.3	-59.8	-79.0	-97.9	+3.7	+93.8	
$\boxed{\text{forest} \rightarrow \vec{\mathbf{E}}_{8} \stackrel{W_{k}}{\longrightarrow} \vec{\mathbf{K}}_{8}}$	-58.9	-75.5	-91.1	-90.6	-75.6	-89.0	-70.8	+4.7	





Attention Pattern

Which words

which other

words

are relevant to

fluffy

 $\overrightarrow{\mathbf{E}}_2$

0.00

1.00

0.00

0.00

0.00

0.00

0.00

0.00

 W_Q

blue

 $\vec{\mathbf{E}}_3$

0.00

0.00

1.00

0.00

0.00

0.00

0.00

0.00

 W_Q

the

 W_Q

0.00

0.00

0.00

0.00

0.00

1.00

0.00

0.00

forest

 $\dot{\vec{\mathbf{E}}}_8$

 $ec{f Q}_8$

0.00

0.00

0.00

0.00

0.00

0.00

1.00

0.00

verdant

 $\overrightarrow{\mathbf{E}}_7$

 $ec{f Q}_{f 7}$

0.00

0.00

0.00

0.00

0.00

0.00

1.00

0.00

|roamed|

 $ec{\mathbf{E}}_{5}$

0.00

0.00

0.00

0.00

0.01

0.99

0.00

0.00

 W_Q

creature

 $ec{\mathbf{E}}_4$

0.00

0.42

0.58

0.00

0.00

0.00

0.00

0.00

 W_Q

 $[\mathbf{a}]$

 $\vec{\mathbf{E}}_1$ $\vec{\mathbf{Q}}_1$

1.00

0.00

0.00

0.00

0.00

0.00

0.00

0.00

 $[a] \rightarrow \vec{E}_1 \xrightarrow{W_k} \vec{K}_1$

 $\boxed{\text{fluffy}} \to \vec{\mathbf{E}}_2 \stackrel{W_k}{\longrightarrow} \vec{\mathbf{K}}_2$

 $|\overrightarrow{blue}| \rightarrow \overrightarrow{E}_3 \xrightarrow{W_k} \overrightarrow{K}_3$

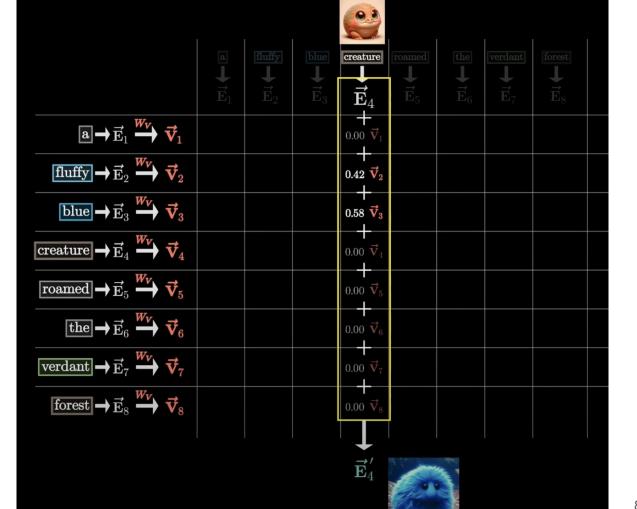
 $\overrightarrow{\text{creature}} \to \overrightarrow{\mathbf{E}}_4 \xrightarrow{W_k} \overrightarrow{\mathbf{K}}_4$

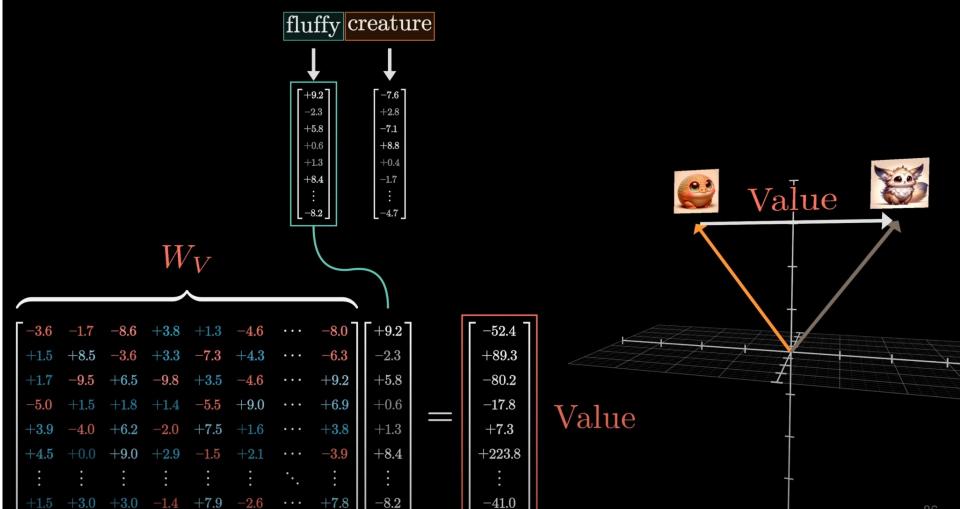
 $\boxed{\text{roamed}} \to \vec{\mathbf{E}}_5 \stackrel{W_k}{\longrightarrow} \vec{\mathbf{K}}_5$

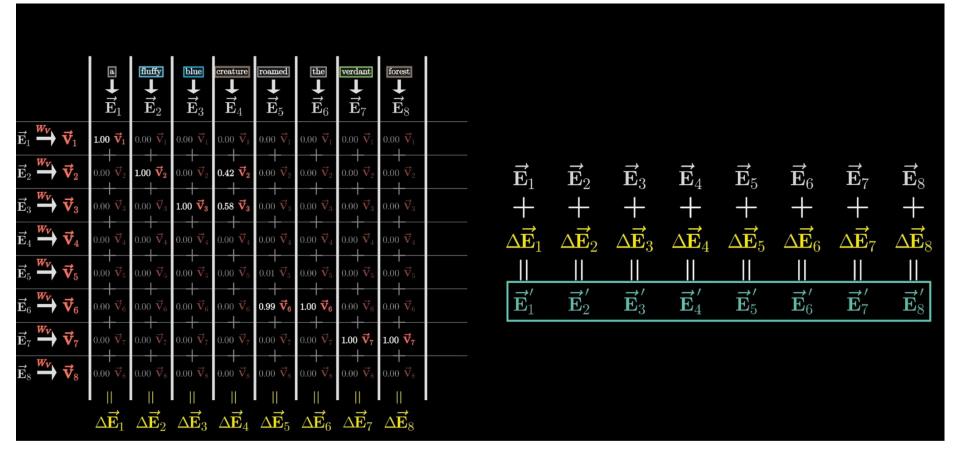
 $\overrightarrow{\text{verdant}} \to \vec{\mathbf{E}}_7 \xrightarrow{W_k} \vec{\mathbf{K}}_7$

 $\boxed{\text{the}} \rightarrow \vec{\mathbf{E}}_6 \xrightarrow{W_k} \vec{\mathbf{K}}_6$

Values







Representing attention

$$Attention(Q, K, V) = softmax(\frac{QK^T}{\sqrt{d_k}})V$$

Attention is all you need

CNN, RNN, LSTM, GAN, Test time data, Early stopping, Data augmentation, Dropout, Batch norm, Gradient clipping

Attention



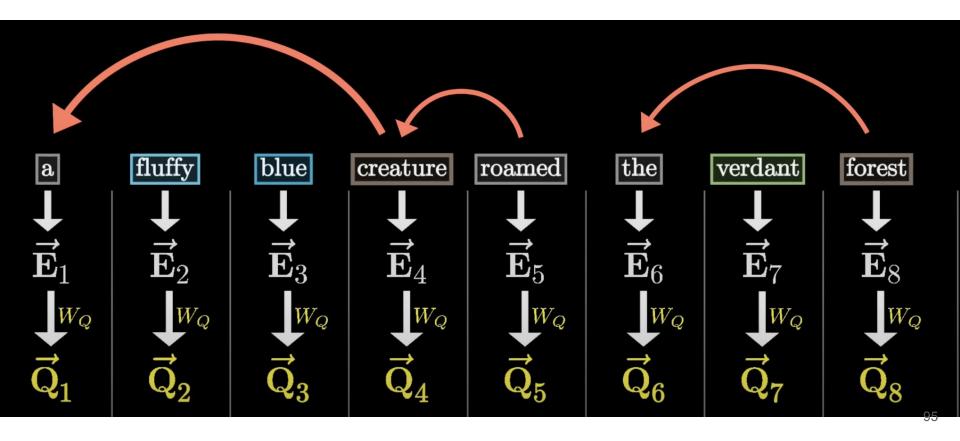
Attention hands-on

Let's write our own attention implementation

For correctness, we will compare with Pytorch's Scaled Dot Product Attention (SDPA) implementation

Masked . Multi-head . Attention

We don't want what comes next to effect the previous words attention



fluffy blue roamedthe forest creatureverdant $ec{\mathbf{E}}_2$ Attention $ec{f E}_3$ $ec{f E}_4$ $ec{\mathbf{E}}_{5}$ $ec{ ilde{\mathbf{E}}}_1$ $ec{ ilde{\mathbf{E}}}_7$ $ec{ ilde{\mathbf{E}}}_{8}$ $ec{f Q}_1^{W_Q}$ $ec{f Q}_3^{W_Q}$ $egin{array}{c} W_Q \ ec{f Q}_5 \end{array}$ $ec{f Q}_7$ $ec{f Q}_{f 2}^{W_Q}$ W_Q W_Q Pattern $[\mathbf{a}] \to \vec{\mathbf{E}}_1 \xrightarrow{W_k} \vec{\mathbf{K}}_1$ • $\boxed{\text{fluffy}} \to \vec{\mathbf{E}}_2 \xrightarrow{W_k} \vec{\mathbf{K}}_2$ • $\boxed{\text{blue}} \rightarrow \vec{\mathbf{E}}_3 \stackrel{W_k}{\longrightarrow} \vec{\mathbf{K}}_3$ • • $\boxed{\text{creature}} \to \vec{\mathbf{E}}_4 \stackrel{W_k}{\longrightarrow} \vec{\mathbf{K}}_4$ • • • $[roamed] \rightarrow \vec{\mathbf{E}}_5 \xrightarrow{W_k} \vec{\mathbf{K}}_5$ • • $\begin{array}{c} \text{the} \rightarrow \vec{\mathbf{E}}_6 \xrightarrow{W_k} \vec{\mathbf{K}}_6 \end{array}$ • • • • • $\overrightarrow{\text{verdant}} \to \overrightarrow{\mathbf{E}}_7 \xrightarrow{W_k} \overrightarrow{\mathbf{K}}_7$ • • • • • $\boxed{\text{forest}} \rightarrow \vec{\mathbf{E}}_{8} \stackrel{W_{k}}{\longrightarrow} \vec{\mathbf{K}}_{8}$ • • • •

Unnormalized Attention Pattern							Normalized Attention Pattern			l	
+3.53	+0.80	+1.96	+4.48	+3.74	-1.95						
+1.90	-0.30	-0.21	+0.82	+0.29	+2.91						
+1.52	+0.24	+0.89	+0.67	+2.99	-0.41	softmax					
+0.63	-1.71	-5.11	+1.31	+1.73	-1.48						
+4.54	-2.91	+0.09	-0.37	+3.07	$+2.9\overline{4}$						
+0.31	+0.76	-1.78	-3.96	-0.70	+0.31						97

Attention Mask

Unnormalized

-0.30

 $-\infty$

 $-\infty$

Attention Pattern

 $-0.21 \mid +0.82 \mid +0.29 \mid$

 $+0.89 \, | \, +0.67 \, | \, +2.99 \, |$

 $-\infty$

 $-\infty$

 $+1.31 \mid +1.73 \mid$

+3.07

 $-\infty$

+3.53 |+0.80 |+1.96 |+4.48 |+3.74 |

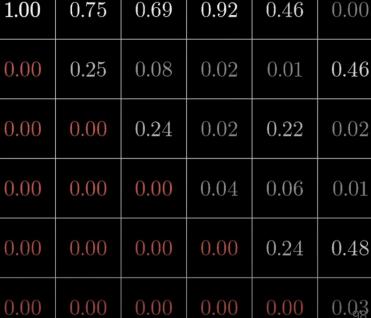
-0.41

+0.31

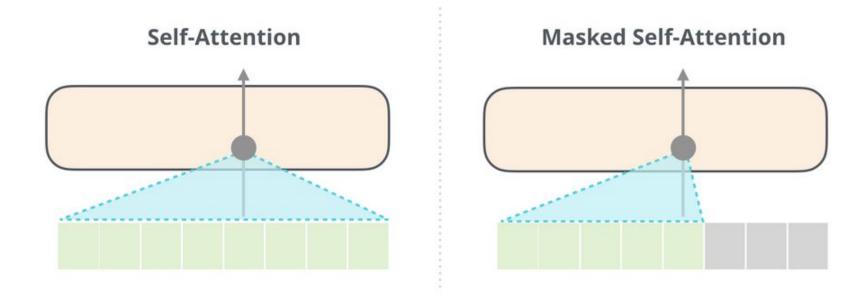


Normalized Attention Pattern

Autemmon i autem										
0.75	0.69	0.92	0.46	0.0						
0.25	0.08	0.02	0.01	0.4						
0.00	0.24	0.02	0.22	0.0						



Attention Mask



Masked . Multi-head . Attention

Query

```
\begin{bmatrix} -3.7 & +3.9 & -2.4 & -6.3 & -9.4 & -8.6 & +3.6 & -0.9 & \cdots & +0.7 \\ +7.9 & +9.7 & -5.6 & +3.2 & -4.7 & -9.5 & +5.1 & -3.6 & \cdots & -2.3 \\ +1.7 & +6.6 & +2.6 & +7.4 & -4.5 & +5.9 & -6.2 & +9.0 & \cdots & +3.7 \\ \vdots & \vdots \\ -5.6 & +8.9 & +4.6 & -4.9 & -5.7 & +0.4 & -9.4 & -5.8 & \cdots & -1.5 \end{bmatrix}
```

Key

```
\begin{bmatrix} -2.5 & -0.7 & -4.4 & +1.7 & +7.2 & -7.6 & +0.3 & -7.3 & \cdots & +4.3 \\ -2.1 & +1.3 & -6.3 & -7.0 & -0.2 & -2.9 & +8.7 & +5.3 & \cdots & +4.9 \\ +8.0 & -8.2 & +1.0 & +1.7 & +9.1 & -4.1 & -5.1 & -7.9 & \cdots & -9.6 \\ \vdots & \vdots \\ +8.5 & +3.4 & +5.6 & -4.3 & +1.7 & -8.6 & -0.3 & +9.5 & \cdots & +7.5 \end{bmatrix}
```

How many parameters?

Value

Query

Key

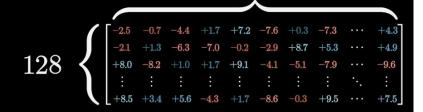
$$128 \times 12,288 = 1,572,864$$

$$128 \times 12,288 = 1,572,864$$

12,288



12,288



Value

 $12,288 \times 12,288 = 150,994,944$

d input 12,288

d output 12,288

12,288

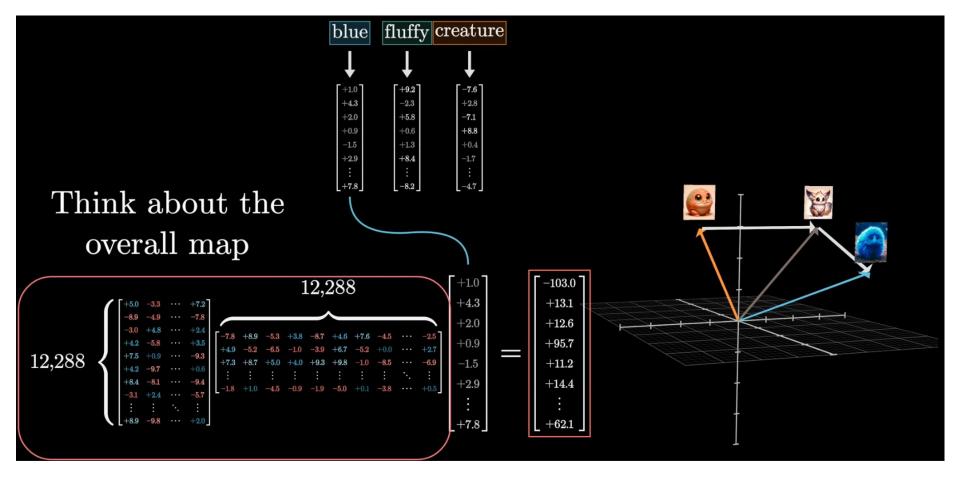


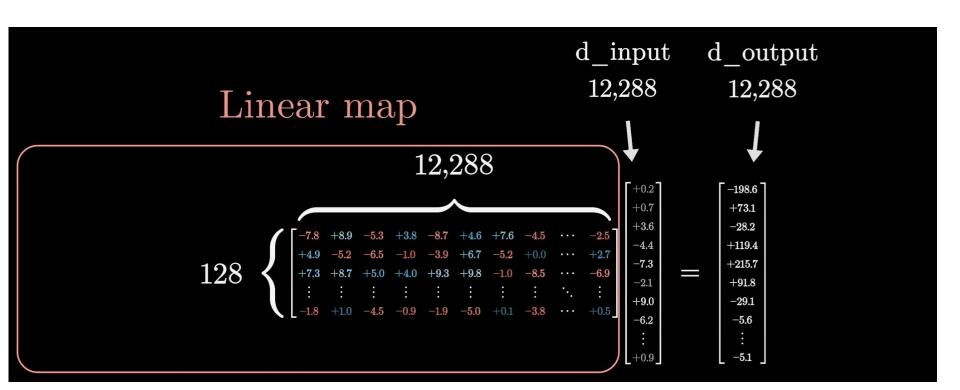
+3.6

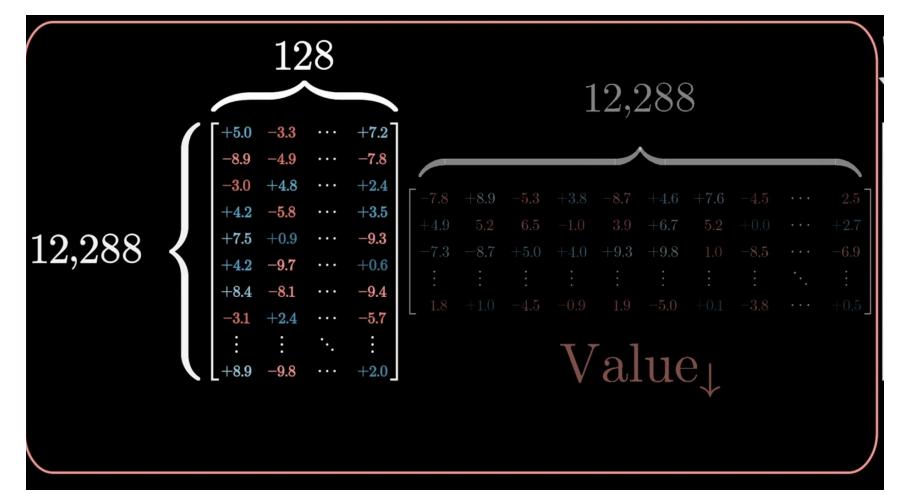
+9.0



+9.4-8.812,288 -0.9+5.5 -198.6+73.1-28.2+119.4+215.7+91.8-29.1-5.6

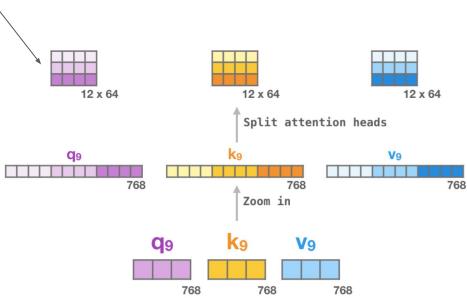






Multi head attention

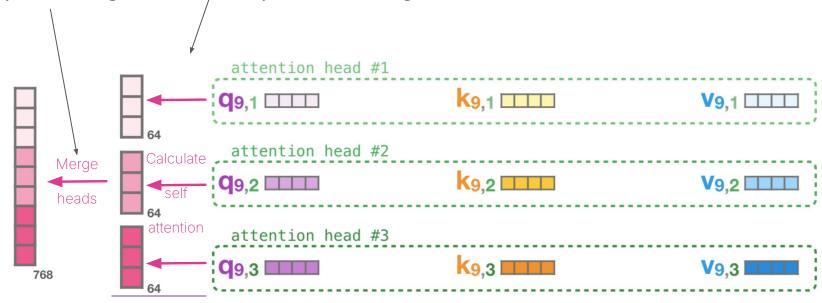
 q, k and v are now split into a number of smaller pieces, called "heads".



Multi head attention

2. Attention is calculated per head

3. Finally, we "merge" the heads by concatenating the results



How torch sdpa looks internally

```
T = q.shape[2]
att = (q @ k.transpose(-2, -1)) * (1.0 / math.sqrt(k.size(-1)))
att = att.masked_fill(self.bias[:,:,:T,:T] == 0, float('-inf'))
att = F.softmax(att, dim=-1)
y = att @ v # (B, nh, T, T) x (B, nh, T, hs) -> (B, nh, T, hs)
return y
```

The 2 sides of attention



It is THE engine of the current LLM wave

LLMs wouldn't be what they are without the Attention block.



It is very costly -O(n2)

Where n is sequence length.

Longer and longer context becomes harder.

Attention is a lot more than this!

Sliding Window

Don't focus on all tokens

Linear

Switch from Quadratic form to an approximate linear form

Hybrid

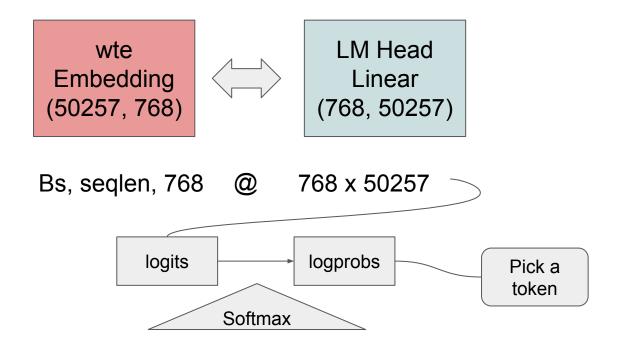
Mix and match various Attention blocks in a single architecture

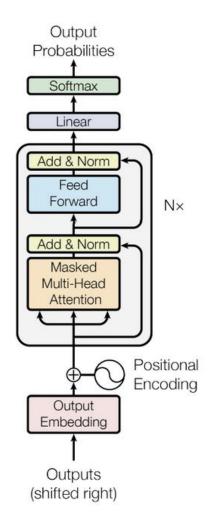
Example of latest research work here:
Jet-Nemotron from Nvidia, 21'st Aug, 2025:
https://www.arxiv.org/pdf/2508.15884

Search for the best attention configuration for a task + hardware

LM Head and Softmax

The last Linear layer is called LMHead





Let's write our own GPT2 Model: Putting the layers together

2 step exercise

Write the module definitions of all the components of the GPT2 Model and load the weights.

Write the forward definitions of all the components.

Win condition:

No error in loading from the state dict.

```
import torch
mw =
torch.load("pytorch_model.bin")
list(mw.keys())
```

Win condition:

When running `generate.py`, you are able to generate proper English sentences.

Step 1: Create the modules (fill in gpt2.py)

Complete the init for all modules using the provided empty classes and hints

- 1. Use the weight state dict keys to get the class field names.
- Use the comments as hints for the type (LayerNorm, Conv1D, MLP, Embedding or another Module)
- Use the dict values shapes as the size for the Conv/Embedding/LayerNorm modules

It works if you are able to load weights without error.

Final answer is in: solutions/gpt_with_init.py

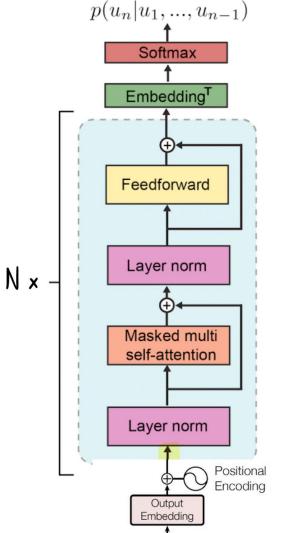
Step 2: Fill in the forwards

Now, write the forward for all the modules. Use the MLP you have already written.

This is how a single block should look.

Make the same calls in the way shown and make sure to add the residual.

If everything works, run generate.py which should produce good output.



Recap

Transformer architecture - GPT2

From here:

- Explore other features of HF ecosystem
- Explore other model architectures (ex. Advancements like MoE)
- Explore Hybrid architectures (Mamba attention)
- Explore Multi-modal LMs such as image, video, audio