CS 696 Applied Large Langauge Models
Spring Semester, 2025
Doc 6 Cluster, Embedding, Attention
Jan 30, 2025

These slides use material from
# Acknowledgments

# SDSU GPU Cluster

PowerEdge R750XA

(2x) Intel Xeon Gold 6338 2G CPU, 32C/64T

(4x) Nvidia A100 GPU, 80 GB RAM

512 GB System RAM


The GPUs available have 10GB of RAM

Suitable for training/fine-tuning 7-8 billion parameter models.

Students 75GB of persistent storage space

# Accessing GPU Cluster

Documentation

https://sdsu-research-ci.github.io/instructionalcluster/students

SDSU Research &
Cyberinfrastructure

SDSU Research & Cyberinfrastructure on GitHub

Instructional Cluster  /  Student Resources

Home

Instructional Cluster ⌃

Overview

**Student Resources** ⌃

Logging In

Launch a Container

Instructor Resources ⌄

Frequently Asked Questions ⌄

How-To Videos ⌄

Available Container Images

Architecture Details

Usage

Research Cluster ⌄

GitHub Education ⌄

Software Factory ⌄

IT@SDSU ↗

TABLE OF CONTENTS

- Logging In

- Launch a Container

Back to top

Edit this page on GitHub

# Logging In

Go to:

https://jupyterhub.sdsu.edu/

Enter your SDSUid credentials

Login to CILogon

**SDSUid (e.g. jsmith@sdsu.edu)**

> Recover your password

> Need Help?

**Password**

☑ **Don't Remember Login**

Login

CILogon facilitates secure access to CyberInfrastructure (CI).

© 2012-2025 San Diego State University (v.P1-AR8)

# Enter Server Options

## Server Options

Container images are described on our available container images page.

**Select a Profile**

○ **Default Profile**
Select compute resources, number of GPUs and a notebook container image, or provide your own image with the "Other" choice.

**Compute Resources**

| Large - 8 CPUs & 16 GB RAM | ↕ |
|---|---|

**Number of GPUs**

| 4 | ↕ |
|---|---|

**Notebook Container Image**

| LLM Notebook | ↕ |
|---|---|

**Start**

# Select Kernel



Select Kernel

Start Preferred Kernel

✓ Python 3 (ipykernel)

Use No Kernel

No Kernel

Use Kernel from Preferred Session

Use Kernel from Other Session

Console 1

# Get Jupyter Interface

# Launcher

Console 1 ✕ | 📄 test.ipynb ✕ | ⬈ Launcher ✕ | +

🔖 **Notebook**

Python 3 (ipykernel)

VS Code (code-server) [↗]

>_ **Console**

Python 3 (ipykernel)

$_ **Other**

Terminal

Text File

Markdown File

Python File

Show Contextual Help

# Large Language Model Notebook

## minimal-notebook

- Ubuntu 22.04.4
- JupyterLab 4.2.4
- Python 3.11.9
- Git 2.34.1
- vi
- nano
- wget
- curl
- unzip
- tzdata

## SciPy Notebook

altair, beautifulsoup4, bokeh, bottleneck, cloudpickle, conda-forge::blas=*=openblas, cython, dask, dill, h5py, jupyterlab-git, matplotlib-base, numba, numexpr, openpyxl, pandas, patsy, protobuf, pytables, scikit-image, scikit-learn, scipy, seaborn, sqlalchemy, statsmodel, sympy, widgetsnbextension, xlrd packages

ipympl and ipywidgets for interactive visualizations and plots in Python notebooks

Facets for visualizing machine learning datasets

## PyTorch Notebook

pytorch machine learning library
torch, torchaudio and torchvision

# Large Language Model Notebook

rclone
FastChat
Ollama
VS Code Server
Jupyter AI
bitsandbytes
transformers
peft
accelerate
trl
ollama-python
openai
pyaudio
portaudio
cuda-nvcc

deepspeed
langchain
huggingface_hub
auto_gptq
autoawq
xformers
dask-kubernetes
chromadb

# Manually Stopping the Server

File  Edit  View  Run  Kernel  Git  Tabs

| | |
|---|---|
| New | ▶ |
| New Launcher | ⇧ ⌘ L |
| Open from Path... | |
| Open from URL... | |
| Open Recent | ▶ |
| New View for | |
| New Console for Activity | |
| Close Tab | ⌥ W |
| Close and Shut Down | ^ ⇧ Q |
| **Close All Tabs** | |
| Save | ⌘ S |
| Save As... | ⇧ ⌘ S |
| Save All | |
| Reload from Disk | |
| Revert to Checkpoint... | |
| Rename... | |
| Duplicate | |
| Download | |
| **Save and Export Notebook As** | ▶ |
| **Workspaces** | ▶ |
| Print... | ⌘ P |
| **Hub Control Panel** | |
| **Log Out** | |

# Warning

I have not been able to run an existing model on the cluster

There is some problem with the version of the libraries

# Back to Tokens

Goodbye    ,    and    thanks    for    all    the    fish    !

17212, 43571, 11, 326, 11707, 395, 722, 290, 13897, 0

Tokens are just a way to represent works in numbers

They don't capture the relationship between words

To capture the relationship between tokens, convert them to a vector

# Byte Pair Encoding (BPE)

Original Algorithm (for compression)

Find the most common pair of characters
Replace with new symbol
Repeat

Example from Wikipedia

aaabdaaabac

ZabdZabac
Z=aa

ZYdZYac
Y=ab
Z=aa

XdXac
X=ZY
Y=ab
Z=aa

# Byte-level BPE

Used by BERT models, GPT-2
When coming across words not in the vocabulary
Convert to UTF-8 and encode pairs of characters
No need for <|unk|>

```
import tiktoken

tokenizer = tiktoken.get_encoding("gpt2")
text = "aaabdaaabac"
integers = tokenizer.encode(text)
for i in integers:
    print(f"{i} -> {tokenizer.decode([i])}")
```

"aaabdaaabac"

7252 -> aa
397 -> ab
6814 -> da
64 -> a
397 -> ab
330 -> ac

"This is a cat"

1212 -> This
318 ->  is
257 ->  a
3797 ->  cat

"Thisisacat"

1212 -> This
271 -> is
330 -> ac
265 -> at

# Token to Embeddings

Represent a token as a vector in n-space

Related tokens should be close to each other



**Vector embeddings of different types of birds**

**Vector embedding of the word squirrel**

# Type of Embeddings

Token Embeddings
   Vector representation (embedding) using a lookup table

Segment Embeddings
   Which sentence a token belongs to

Position Embeddings
   Position of each token in the sequence

# Embedding Space

Higher the dimension of the space

    Captures more information about the relationship between tokens

    Requires more computation


Embedding Size

    GPT-2 Models               768 dimensions

    GPT-3 (175 parameters)   12,288 dimensions

    Bert-base                  768 dimensions

    Bert-large               1024 dimensions


Map each token to a vector in the space

# How to do the Embedding

N    number of tokens

D    dimension of embedding space (hyperparameter)


Create a matrix (weight matrix) with N rows and D columns
Fill with random values


The K'th row is the embedding (vector) of token ID K


Use training data to modify the weight matrix

# Example from the Text

Small values so can see what is going on

Number of tokens = 6

Dimension of embedding space = 3

Input text

Fox jumps over dog

# Example from the Text

**Weight matrix of the embedding layer** →

$$\begin{bmatrix} 0.3374 & -0.1778 & -0.1690 \\ 0.9178 & 1.5810 & 1.3010 \\ 1.2753 & -0.2010 & -0.1606 \\ -0.4015 & 0.9666 & -1.1481 \\ -1.1589 & 0.3255 & -0.6315 \\ -2.8400 & -0.7849 & -1.4096 \end{bmatrix}$$

**Token IDs to embed**

$$\begin{bmatrix} 2 \\ 3 \\ 5 \\ 1 \end{bmatrix}$$  fox jumps over dog

**Input text**

fox
jumps
over
dog
→
$$\begin{bmatrix} 2 \\ 3 \\ 5 \\ 1 \end{bmatrix}$$

**Embedding vector of the first token ID**

**Embedded token IDs**

$$\begin{bmatrix} 1.2753 & -0.2010 & -0.1606 \\ -0.4015 & 0.9666 & -1.1481 \\ -2.8400 & -0.7849 & -1.4096 \\ 0.9178 & 1.5810 & 1.3010 \end{bmatrix}$$

**Embedding vector of the third token ID**

# LLM Predict the next Word

Training



Text sample:

LLMs learn to predict one word at a time

LLMs learn to predict one word at a time

LLMs learn to predict one word at a time

LLMs learn to predict one word at a time

LLMs learn to predict one word at a time

LLMs learn to predict one word at a time

LLMs learn to predict one word at a time

LLMs learn to predict one word at a time

The LLM can't access words past the target.

Input the LLM receives

Target to predict

# Training

**Sample text**

"In the heart of the city stood the old library, a relic from a bygone era. Its stone walls bore the marks of time, and ivy clung tightly to its facade ..."

**Tensor containing the inputs**

```
x = tensor([[  "In",        "the",        "heart",    "of"    ],
            [  "the" ,      "city",       "stood",    "the"   ],
            [  "old",       "library",    ",",        "a"     ],
            [  ...                                             ]])
```

**Tensor containing the targets**

```
y = tensor([[  "the",       "heart",      "of",       "the"   ],
            [  "city",      "stood",      "the",      "old"   ],
            [  "library",   ",",          "a",        "relic"],
            [  ...                                             ]])
```

# Stride, Window Size, Context

**A stride of 1 moves the input field by 1 position**

**Sample text**

"In the heart of the city stood the old library, a relic from a bygone era. Its stone walls bore the marks of time, and ivy clung tightly to its facade …"

**Inputs of batch 1:** "In the heart of"

**Inputs of batch 2:** "the heart of the"

**A stride of 4 moves the input field by 4 positions**

"In the heart of the city stood the old library, a relic from a bygone era. Its stone walls bore the marks of time, and ivy clung tightly to its facade …"
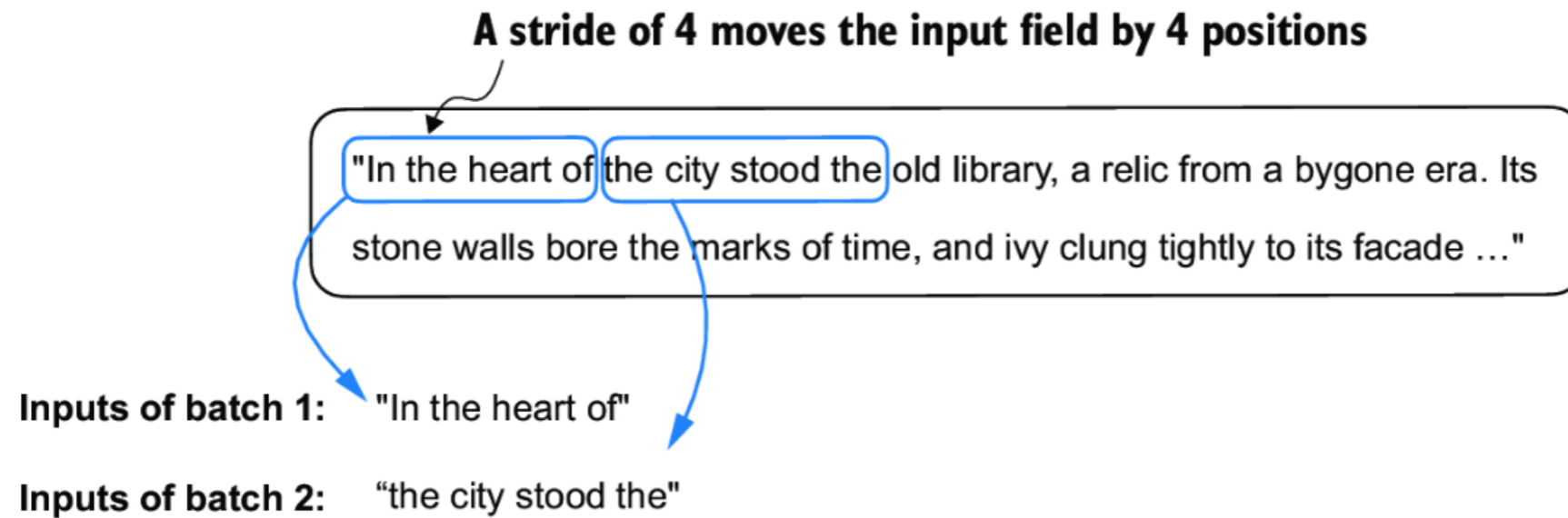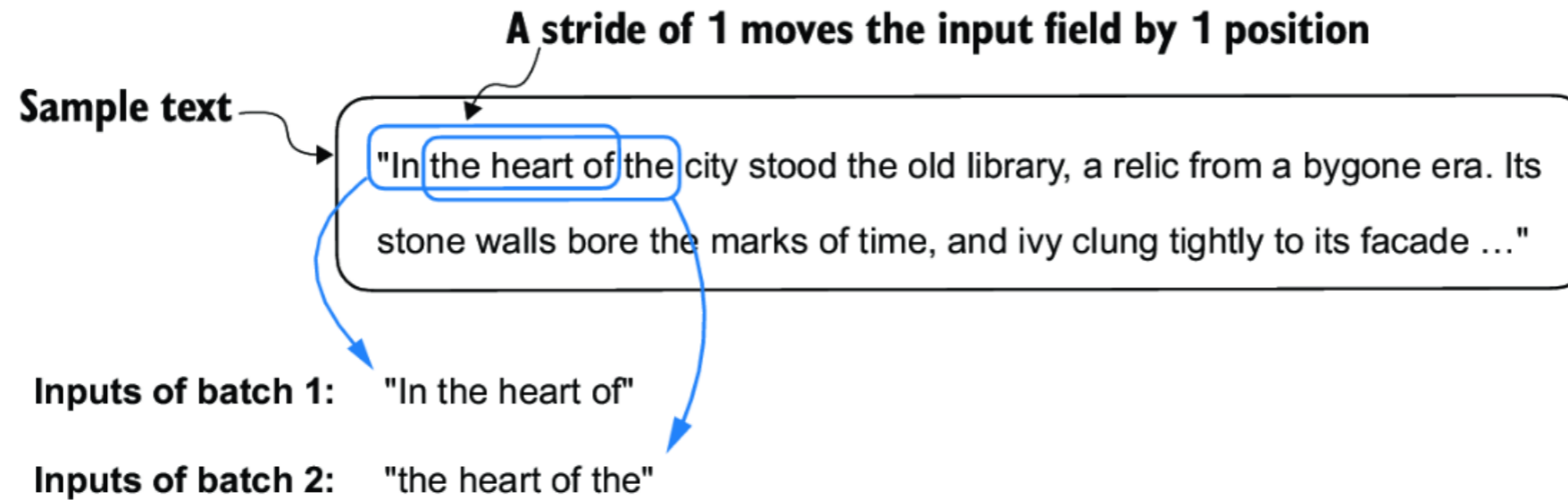
**Inputs of batch 1:** "In the heart of"

**Inputs of batch 2:** "the city stood the"

# Adding Absolute Positional Embeddings

Input Embeddings    | 1.3 | -0.2 | -0.2 |        | -0.4 | 1.0 | -1.1 |

↑                                ↑

Positional Embedings | 1.1 | 1.2 | 1.3 |        | 2.1 | 2.2 | 2.3 |

\+                                \+

Token Embedings     | 1.3 | -0.2 | -0.2 |        | -0.4 | 1.0 | -1.1 |

First token                    Second token

GPT uses absolute positional embeddings optimized in training

# Positional Embeddings - Sinusoidal

Consider the sentence: "The cat sat on the mat."

Position 1 ("The"):  [sin(1/10000^(0/5)), cos(1/10000^(0/5)), sin(1/10000^(2/5)), cos(1/10000^(2/5))]

Position 2 ("cat"):  [sin(2/10000^(0/5)), cos(2/10000^(0/5)), sin(2/10000^(2/5)), cos(2/10000^(2/5))]

Position 3 ("sat"):  [sin(3/10000^(0/5)), cos(3/10000^(0/5)), sin(3/10000^(2/5)), cos(3/10000^(2/5))]

# Positional Embeddings - Rotary (RoPE)

Combines relative and absolute position

$$f_{\{q,k\}}\left(\boldsymbol{x}_m, m\right) = \begin{pmatrix} \cos m\theta & -\sin m\theta \\ \sin m\theta & \cos m\theta \end{pmatrix} \begin{pmatrix} W_{\{q,k\}}^{(11)} & W_{\{q,k\}}^{(12)} \\ W_{\{q,k\}}^{(21)} & W_{\{q,k\}}^{(22)} \end{pmatrix} \begin{pmatrix} x_m^{(1)} \\ x_m^{(2)} \end{pmatrix}$$

"rotate the affine-transformed word embedding vector by the number of angle multiples of its position index"
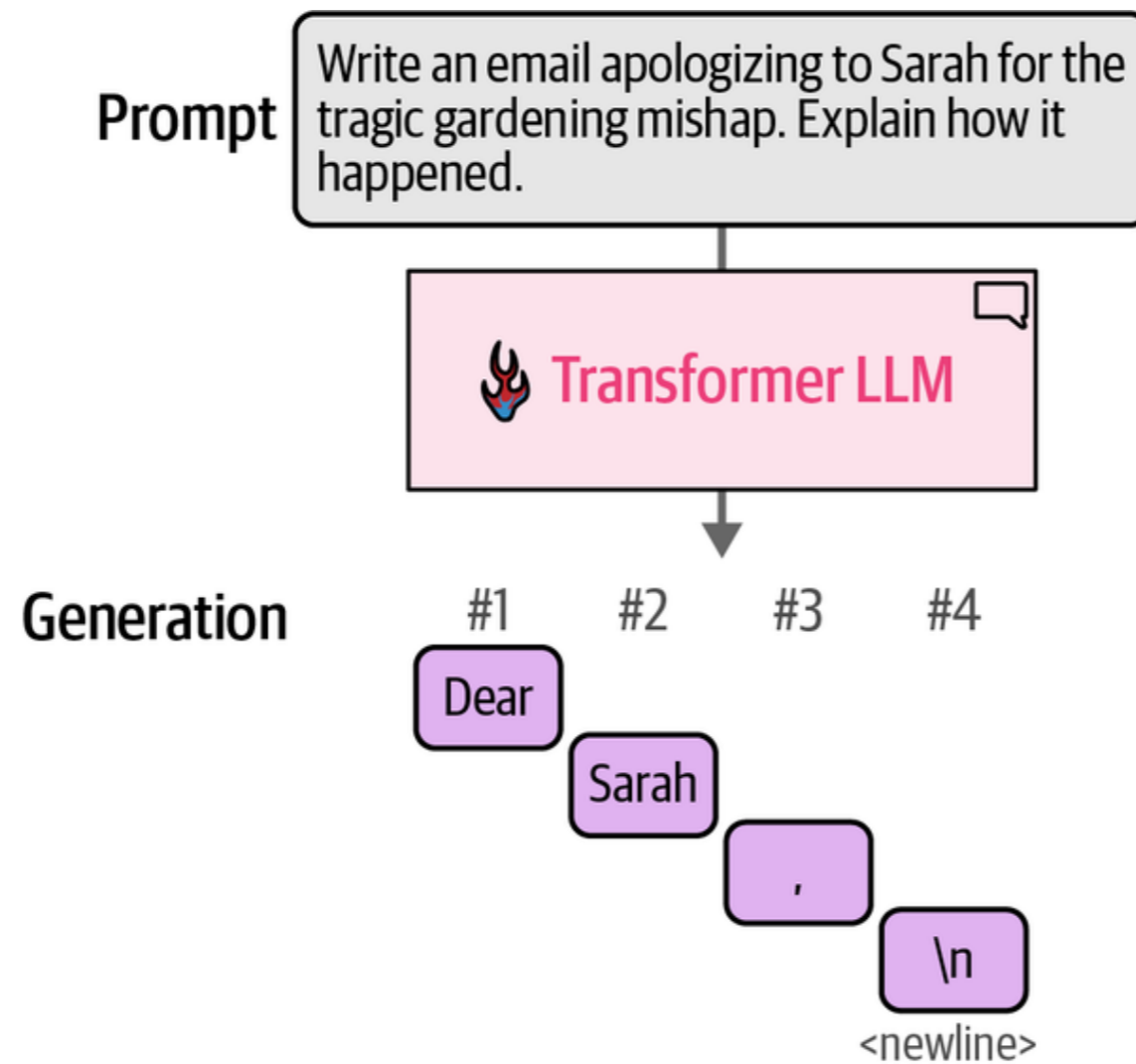
# Relative Positional Embedding

Example: The cat sat on
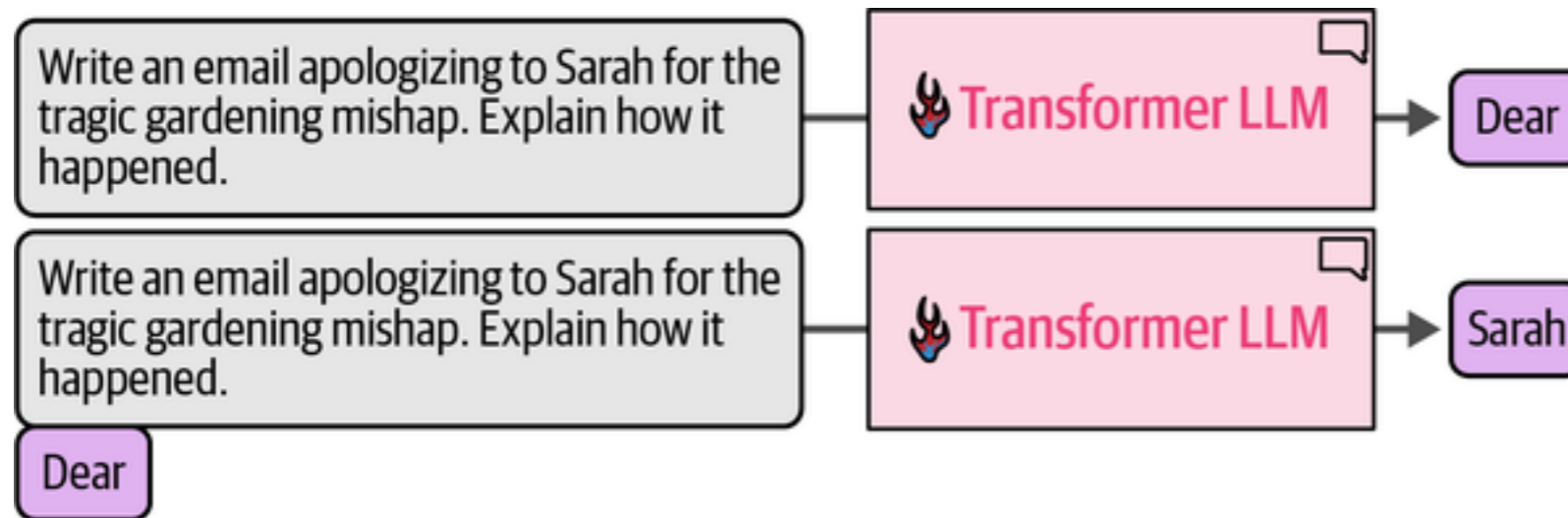
Distance between tokens

The model learns how far apart tokens are

| Token Pair | Relative Distance | Embedding |
|---|---|---|
| The - "cat" | 1 | [0.2, 0.5, -0.1] |
| cat - "The" | -1 | [-0.3, 0.1, 0.4] |
| cat - "sat" | 1 | [0.2, 0.5, -0.1] |
| sat - "cat" | -1 | [-0.3, 0.1, 0.4] |
| The - "sat" | 2 | [0.8, -0.2, 0.3] |
| sat - "The" | -2 | [-0.7, 0.6, -0.5] |
| cat - "on" | 2 | [0.8, -0.2, 0.3] |
| on - "cat" | -2 | [-0.7, 0.6, -0.5] |
| The - "on" | +3 (clipped to +2) | [0.8, -0.2, 0.3] |
| on - "The" | -3 (clipped to -2) | [-0.7, 0.6, -0.5] |

# Big Picture - Responses One Word at a Time

# Big Picture - Attach Predicted Word to input

Write an email apologizing to Sarah for the tragic gardening mishap. Explain how it happened.

Transformer LLM → Dear

Write an email apologizing to Sarah for the tragic gardening mishap. Explain how it happened.

Dear

Transformer LLM → Sarah
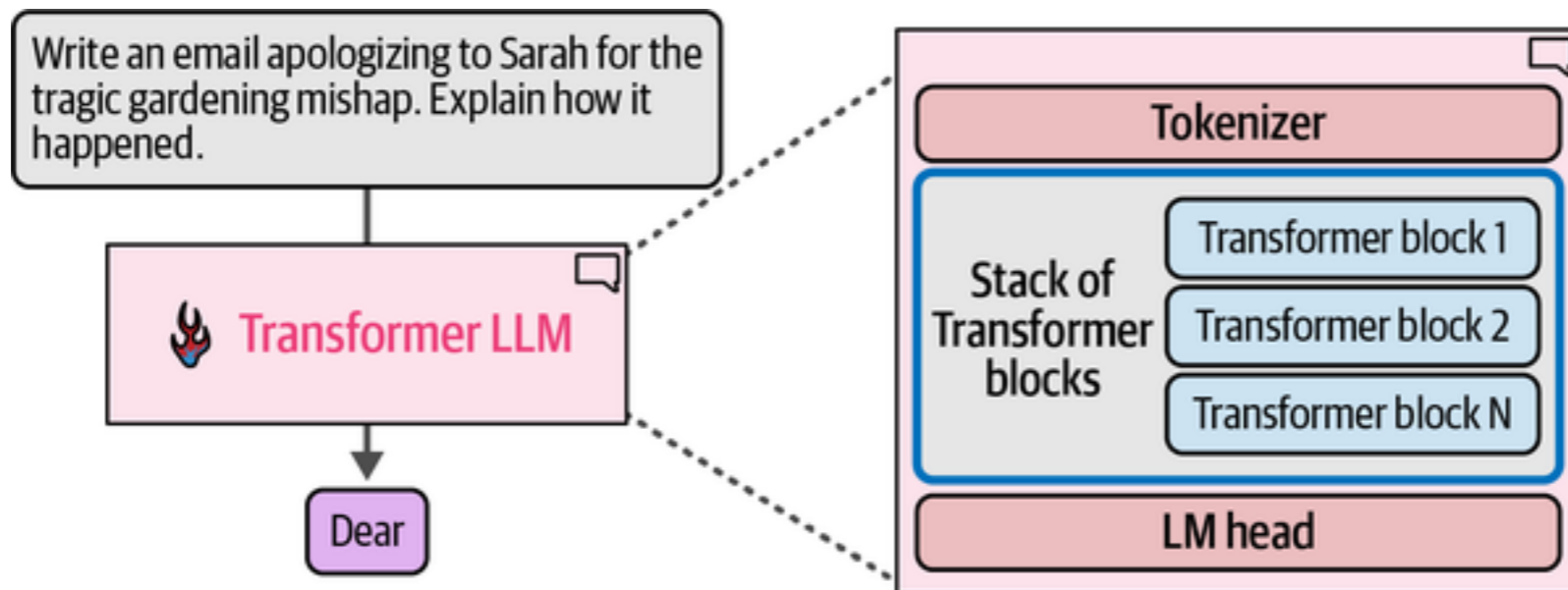
This is not how you see LLMs work

LLMs
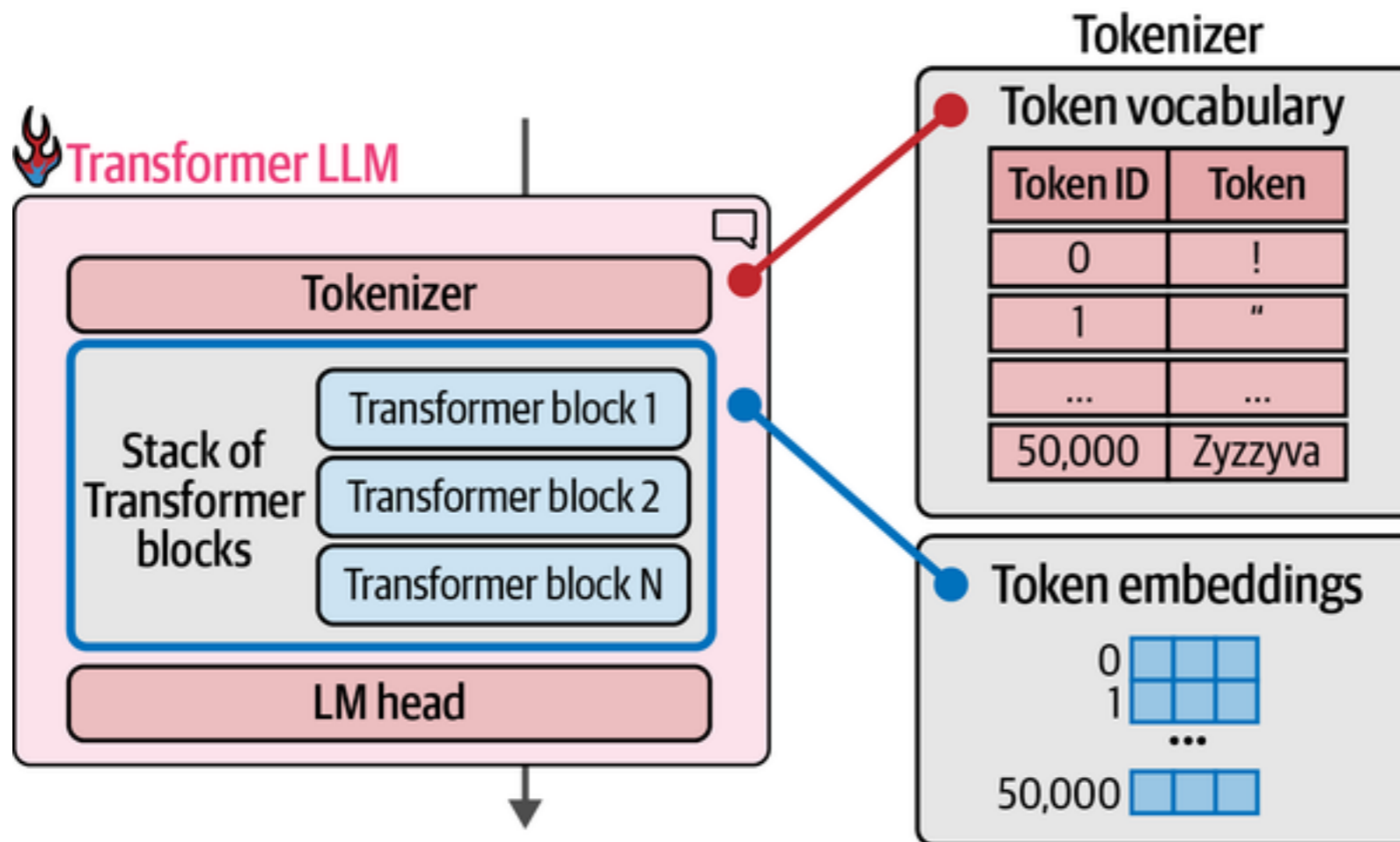    Trained on instruction-tuning and human preference
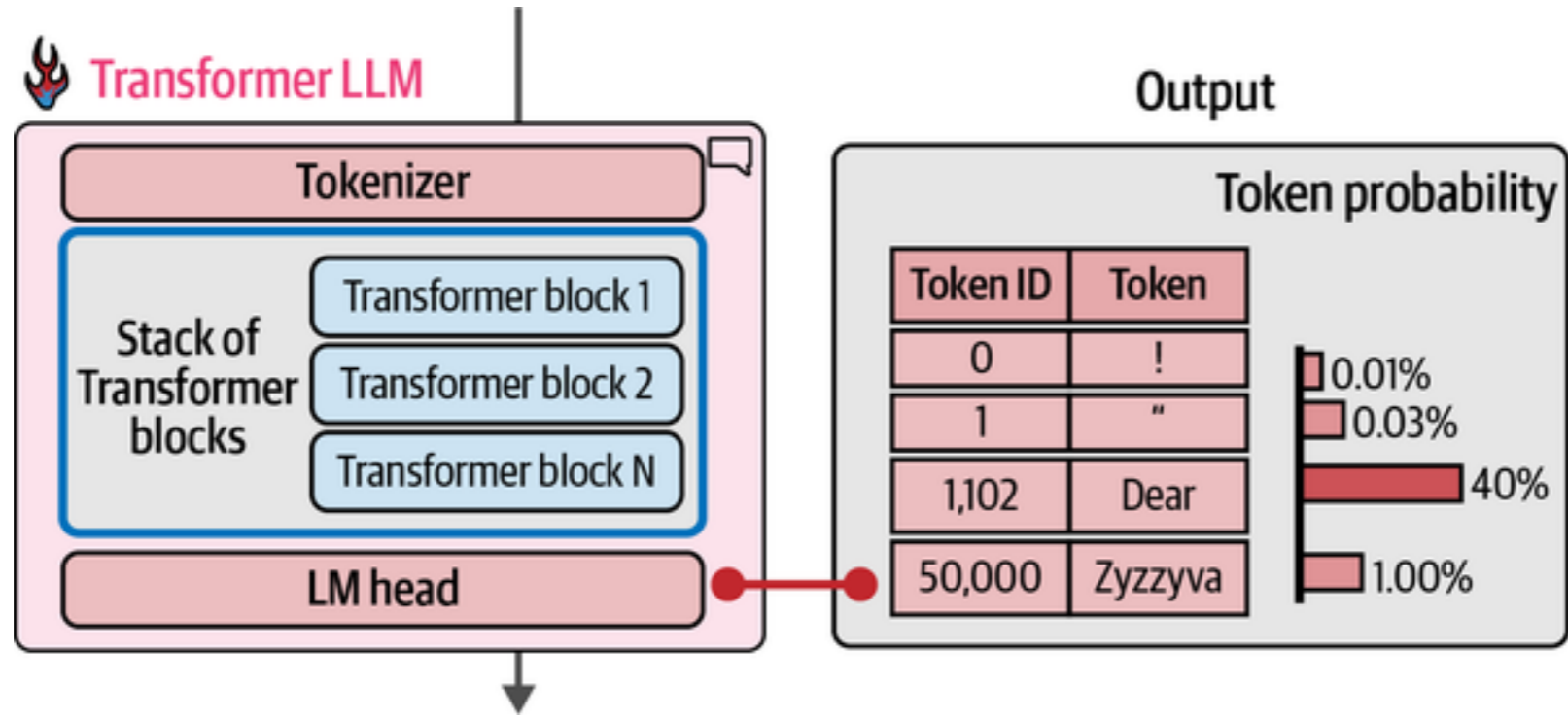    To match what we want

# Big Picture - LM Head

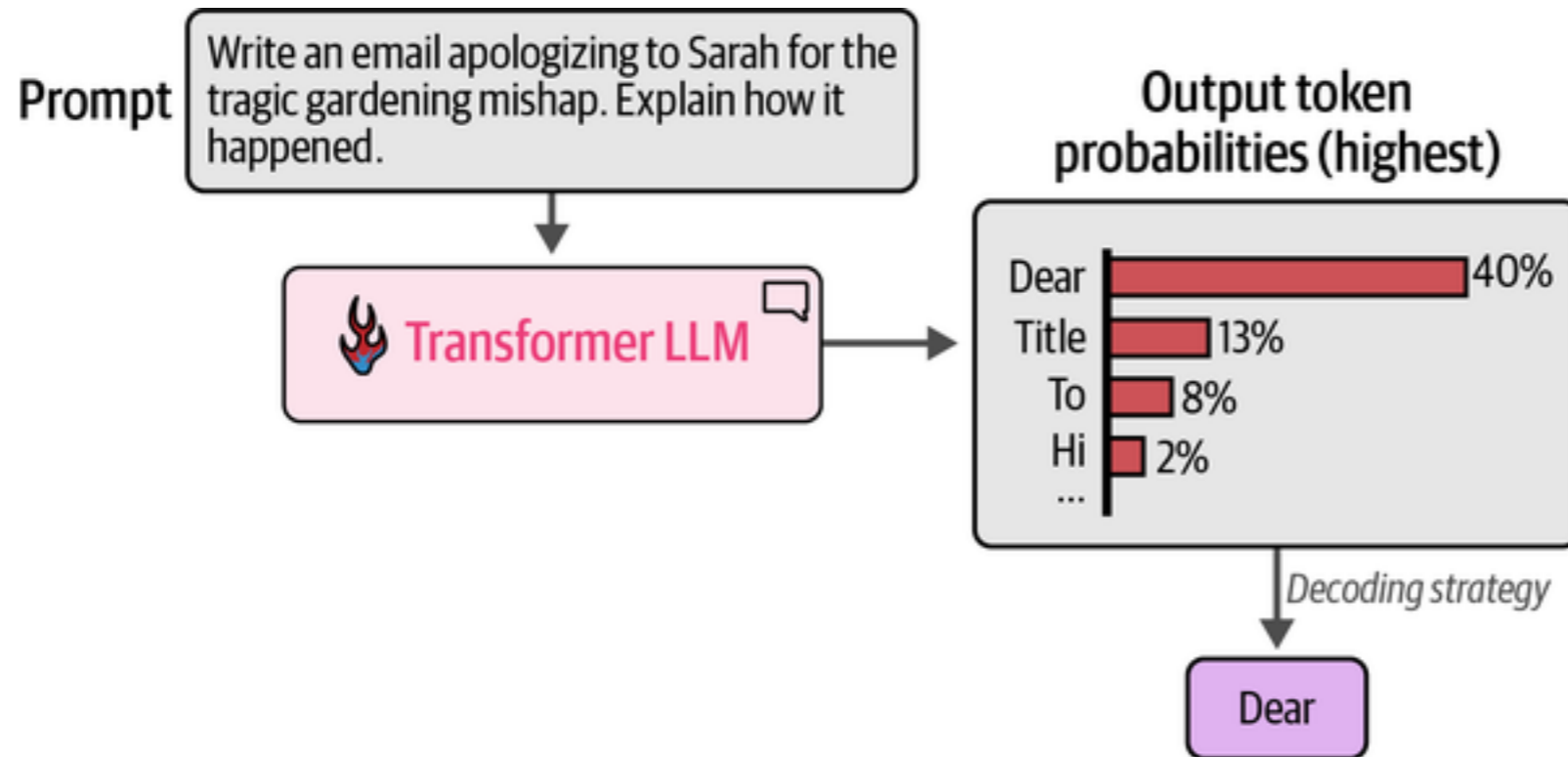# Big Picture - LM Head

# Big Picture - Output

# Decoding Strategy

Don't just choose the token with the highest probability

Sample based on probabilities
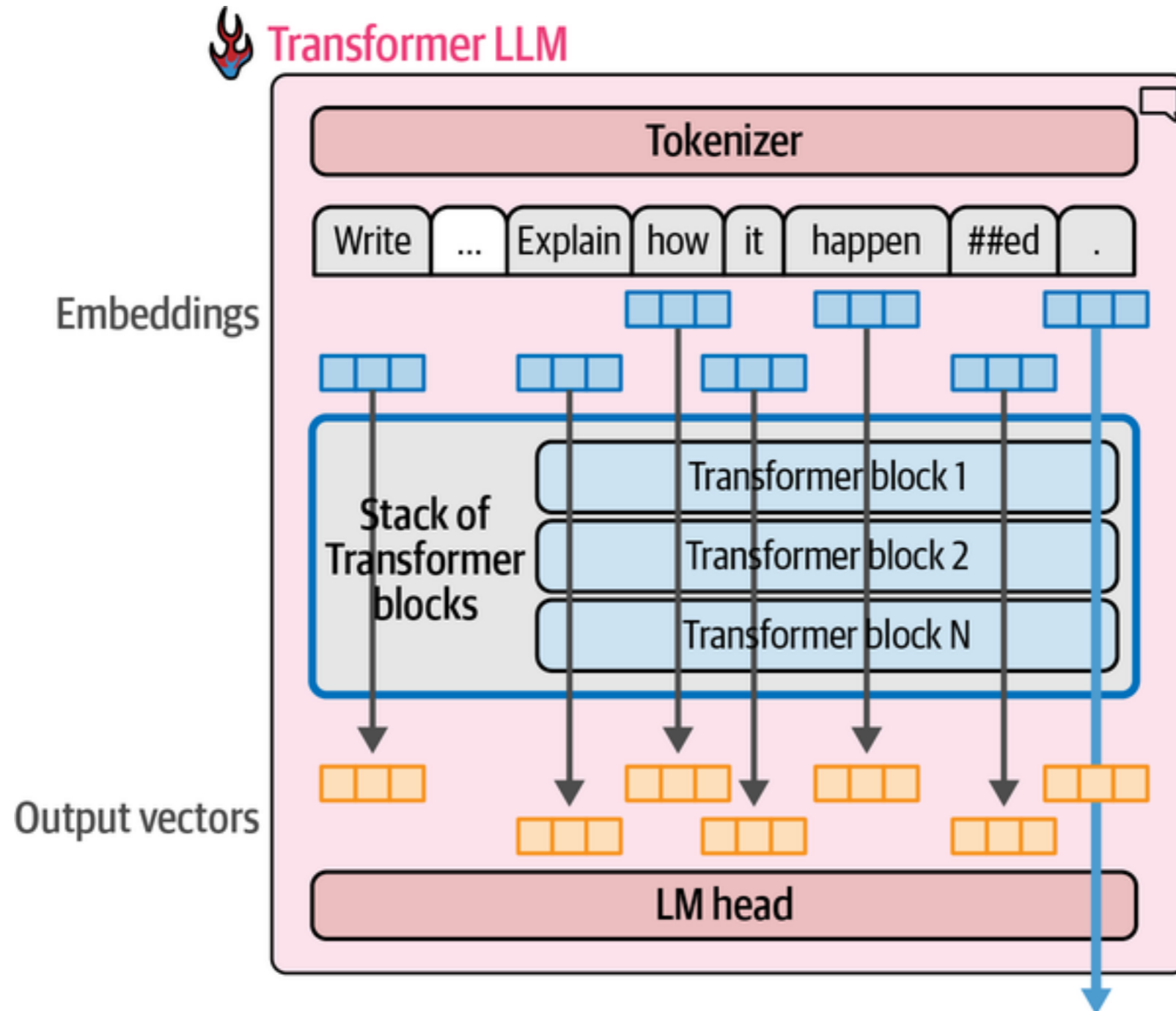    Choose Dear 40% of the time
    Choose Title 13% of the time

# Big Picture - Processing Token in Parallel (sort of)

# Big Picture - Caching Keys & Values

Context Length

Number of tokens can be processed at once



| Model | Context Length |
|-------|----------------|
| GPT-4o | 128k |
| GPT 3.5 | 4,095 |
| GPT 4 | 8,192 |
| Llama 1 | 2,048 |

# Big Picture - Processing Token in Parallel (sort of)

## Transformer LLM



Tokenizer

| Say | something | smart |

Embeddings

Transformer block 1

**Output of Transformer block 1**

Stack of Transformer blocks

Transformer block 2

Transformer block N

Output vectors

LM head

40

# Block Contents



Output of Transformer block 1

Transformer block

Transformer block 1

Transformer block 2

Transformer block N

Self-attention

Feedforward neural network

The | dog | chased | the | squirrel | because | it

dog    squirrel    because    it

Attention

Transformer block 1

Transformer block 2

Transformer block N

Transformer block

Self-attention

Feedforward neural network

# Attention



Need a way to compute how relevant each previous token is

Combine those computations into output vector

# Attention - Chapter 3



This chapter implements the attention mechanism, an important building block of GPT-like LLMs

**STAGE 1**

1) Data preparation & sampling

2) **Attention mechanism**

3) LLM architecture

Building an LLM

4) Pretraining
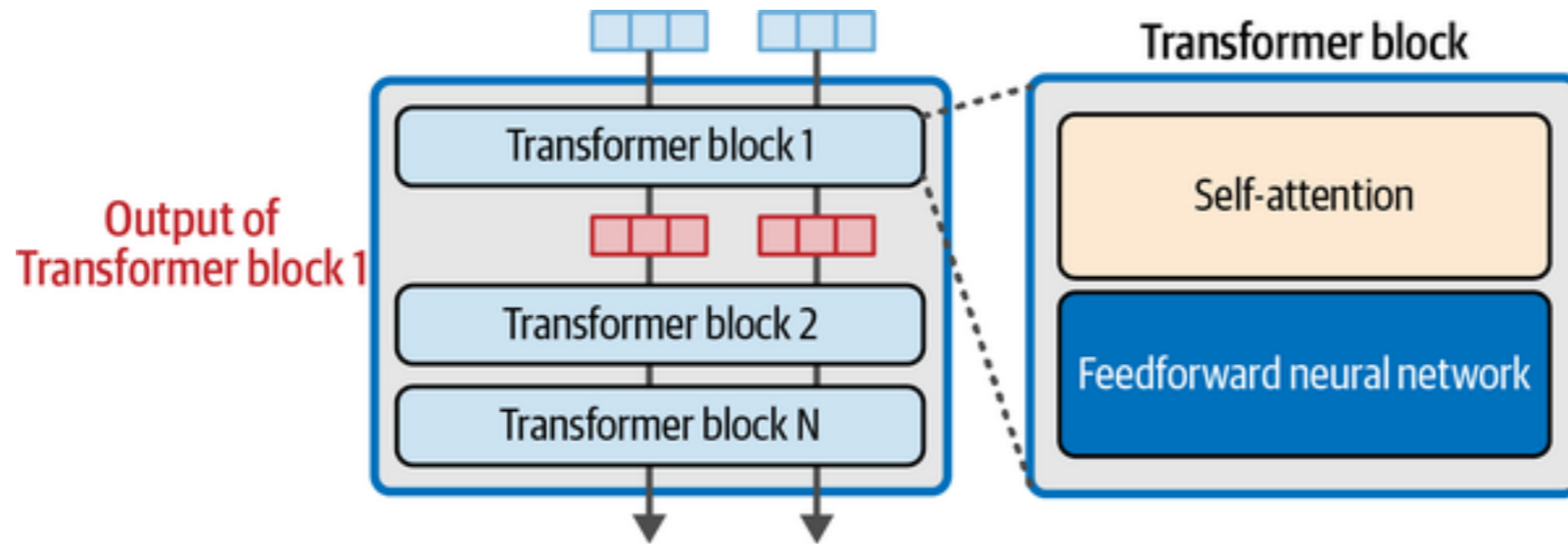
**STAGE 2**

5) Training loop

6) Model evaluation

7) Load pretrained weights

8) Fine-tuning

Foundation model

**STAGE 3**

Dataset with class labels

Classifier

9) Fine-tuning

Personal assistant

Instruction dataset

A simplified self-attention
technique to introduce the
broader idea

A type of self-attention used in LLMs
that allows a model to consider only
previous and current inputs in a
sequence, ensuring temporal order
during the text generation

1) Simplified
self-attention → 2) Self-attention → 3) Causal attention → 4) Multi-head
attention

Self-attention with trainable
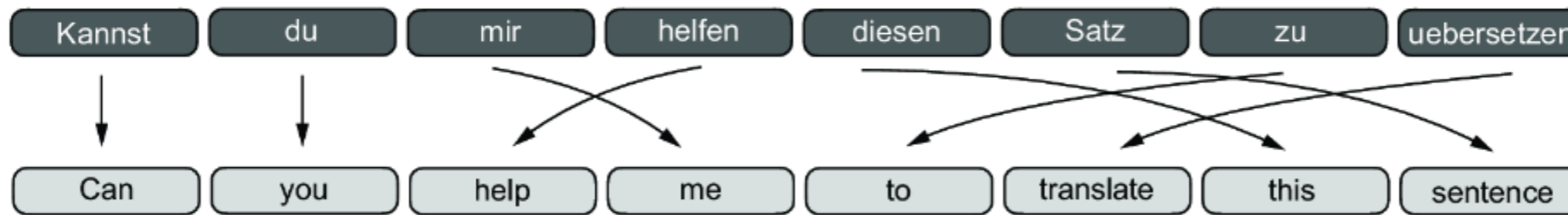weights that forms the basis of
the mechanism used in LLMs

An extension of self-attention and
causal attention that enables the
model to simultaneously attend
to information from different
representation subspaces

# Why We need Attention



German input sentence to translate

| Kannst | du | mir | helfen | diesen | Satz | zu | uebersetzen |
|---|---|---|---|---|---|---|---|

| Can | you | me | help | this | sentence | to | translate |
|---|---|---|---|---|---|---|---|

The word-by-word translation results in a grammatically incorrect sentence

| Kannst | du | mir | helfen | diesen | Satz | zu | uebersetzen |
|---|---|---|---|---|---|---|---|

| Can | you | help | me | to | translate | this | sentence |
|---|---|---|---|---|---|---|---|

The correct translation

Certain words in the generated translation require access to words that appear earlier or later in the original sentence.

# Encoder - Decoder NN

The translated English sentence

Decoder

Outputs: Can → you → help → ...

Hidden states

Hidden states of a neural network

Encoder

Inputs: Kannst → du → mir → ...

A memory cell (hidden state) memorizing entire input

German input sentence to translate

# Bahdanau Attention



We are focusing on generating the second output token.

Outputs: Can — you — help — ...

Hidden states

Inputs: Kannst — du — mir — ...

When generating an output token, the model has a way to access to all input tokens.

The dotted line width is proportional to how important the input token is for the respective output token.