

Introduction to Language Models

Pranav Reddy
Cofounder, Xylem AI
@pranavreddyg (Twitter/X)

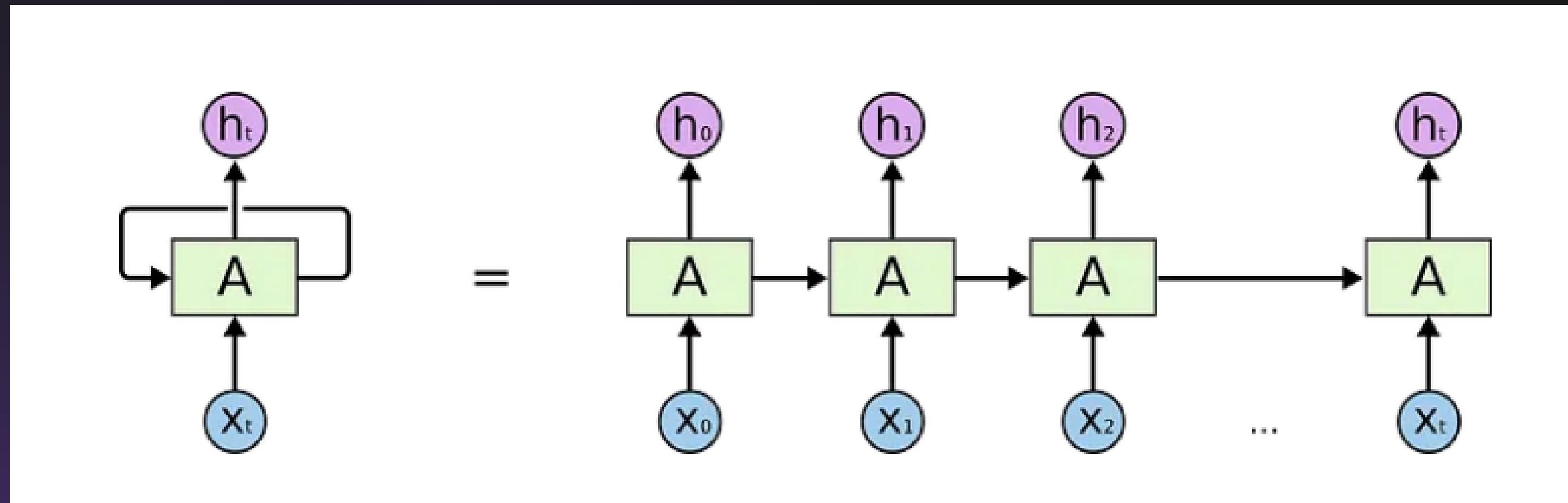


I DON'T KNOW WHAT
A LARGE LANGUAGE MODEL IS

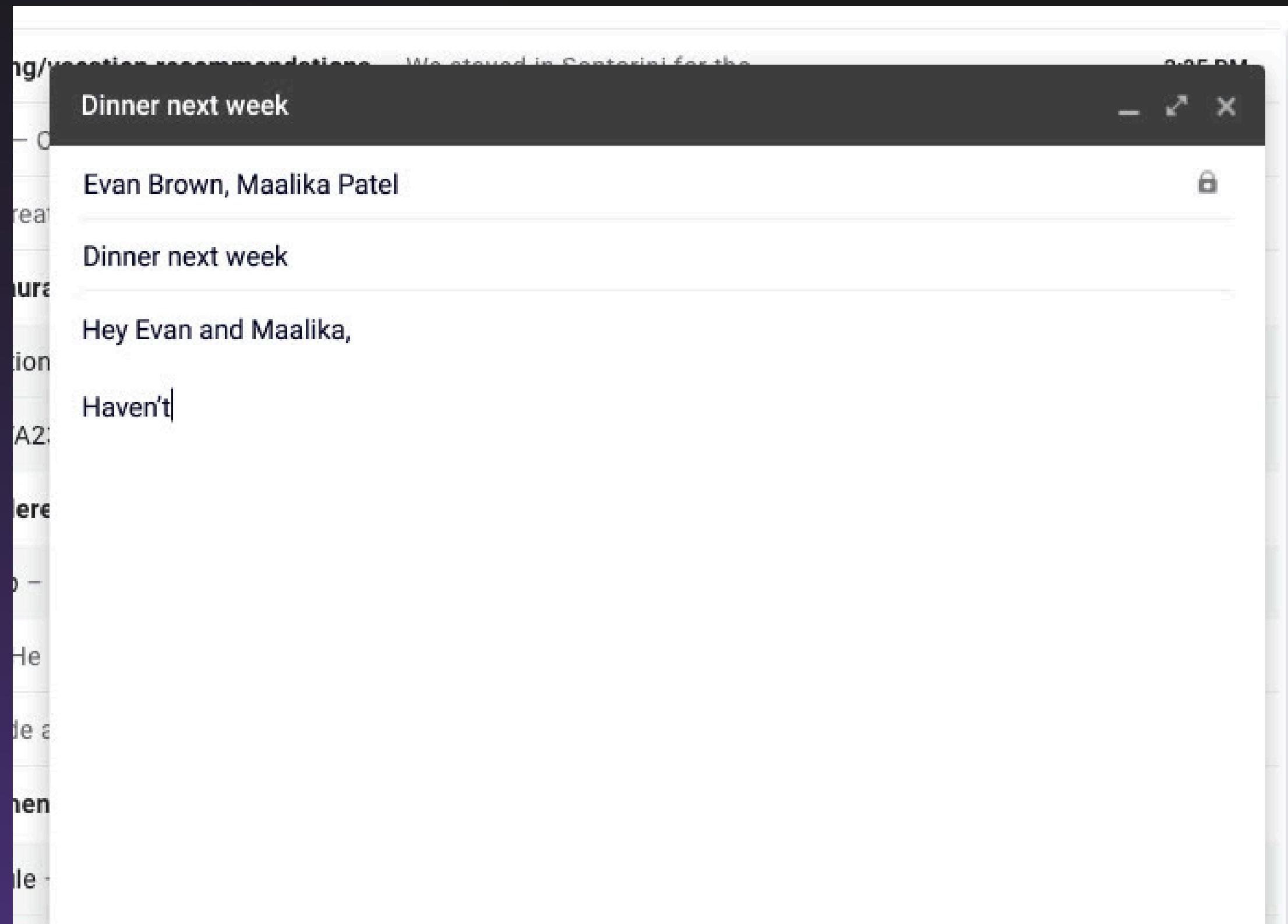
AND AT THIS POINT
I AM TOO AFRAID TO ASK

imgflip.com

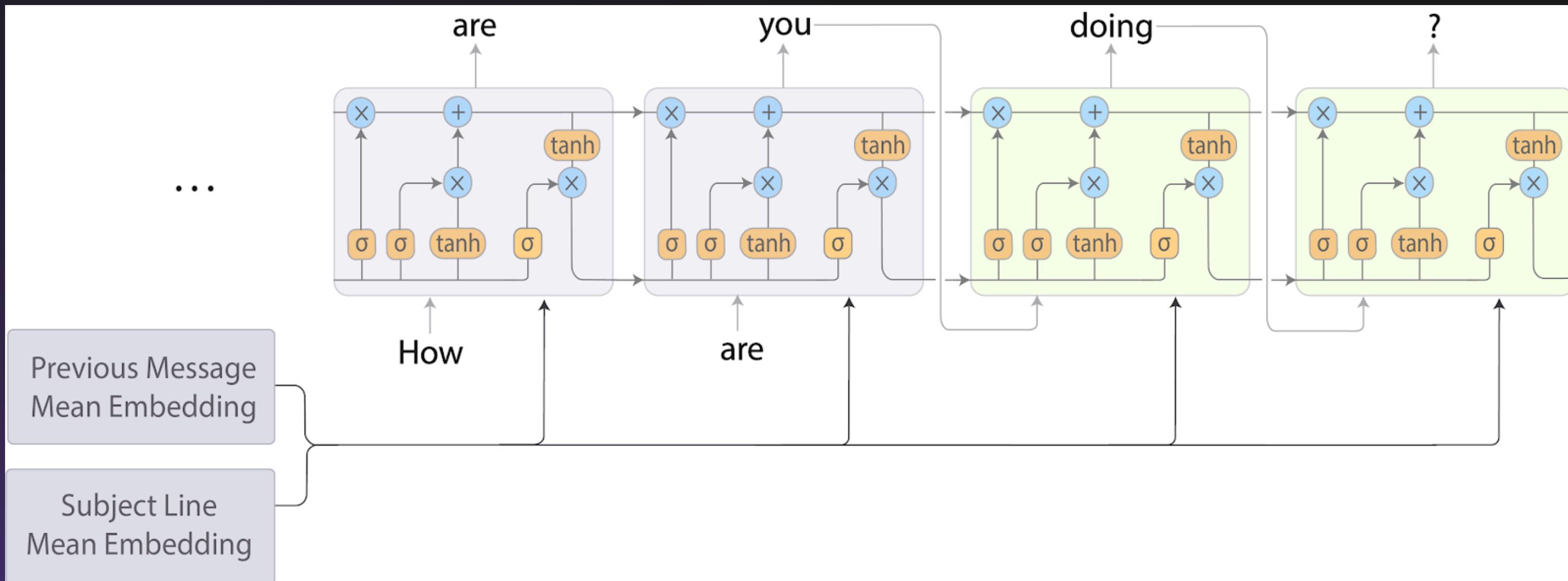
Recurrent Neural Networks



Smart Compose (Gmail)



How it works



I CAN MAKE THE BAD GUYS

GOOD FOR A WEEKEND

I CAN MAKE THE BAD GUYS

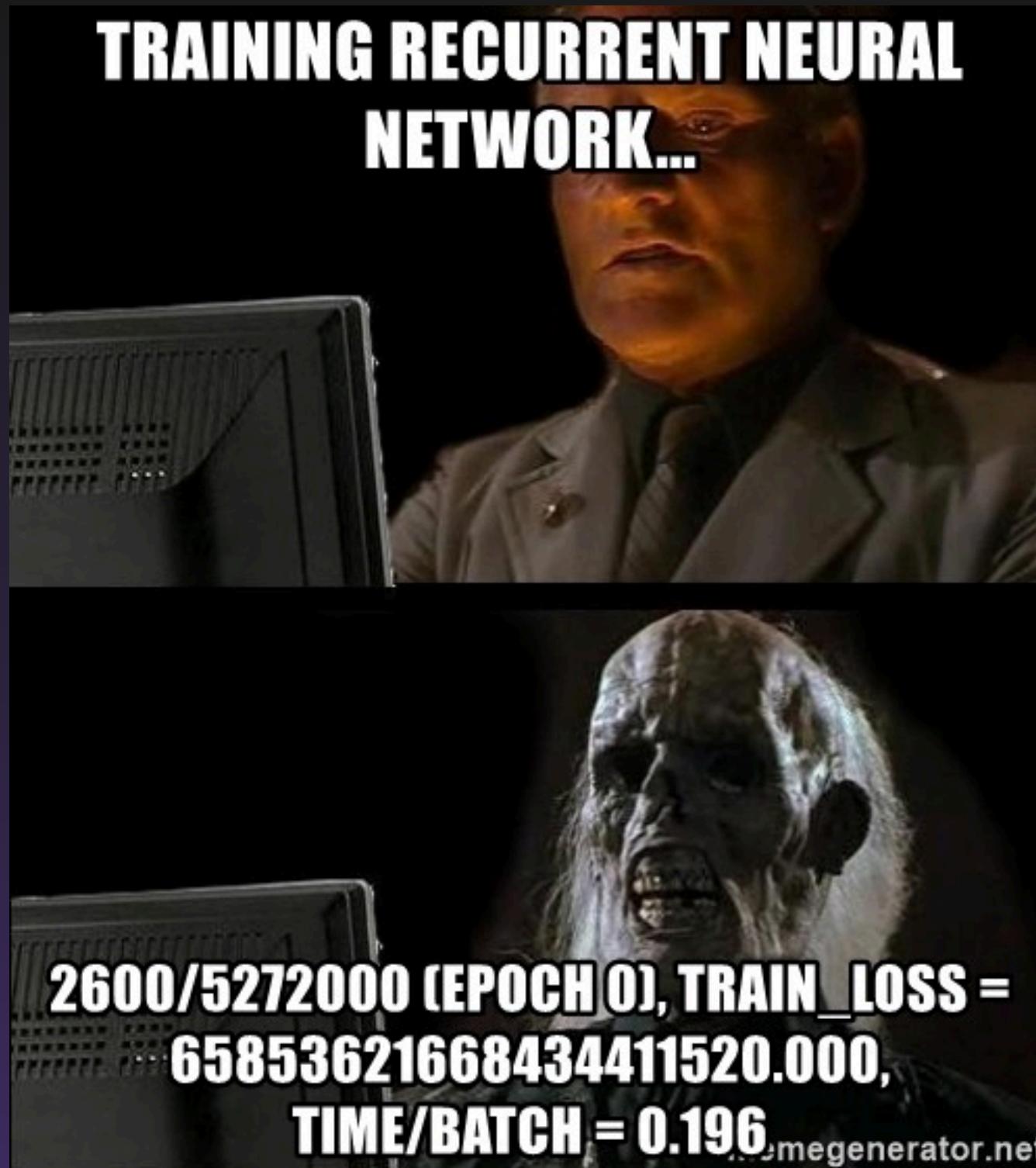
TPGMP

HAVE A LIFE-CHANGING JOURNEY
AND BECOME GREAT PEOPLE
FOR THE REST OF THEIR LIVES

ifunny.co



TRAINING RECURRENT NEURAL NETWORK...



**2600/5272000 (EPOCH 0), TRAIN_LOSS =
65853621668434411520.000,
TIME/BATCH = 0.196**

Problems with RNNs

- Slow computation for long sequences
- Vanishing or exploding gradients
- Difficulty in retaining context in long sequences
- Very hard to train

Attention is all you need

Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

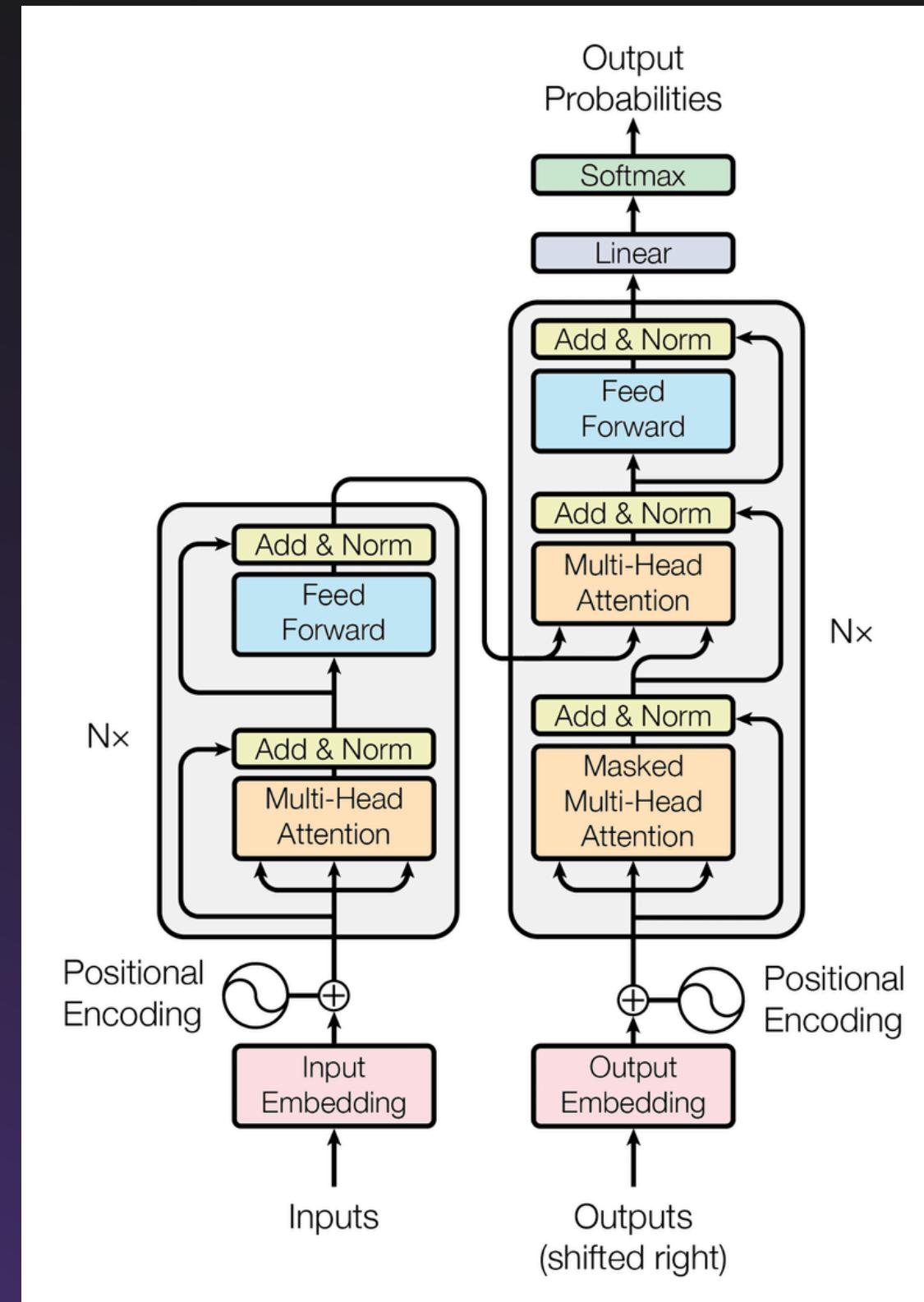
Lukasz Kaiser*
Google Brain
lukasz.kaiser@google.com

Illia Polosukhin* ‡
illia.polosukhin@gmail.com

Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly

Transformer architecture



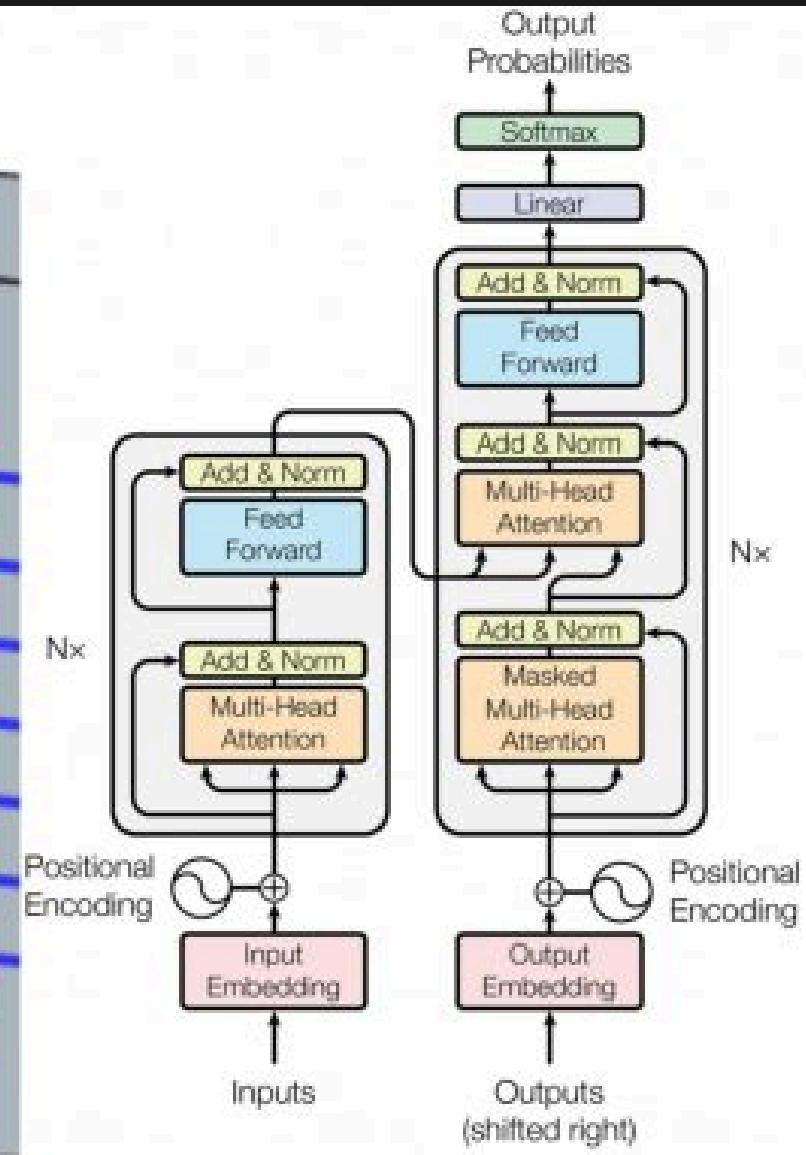
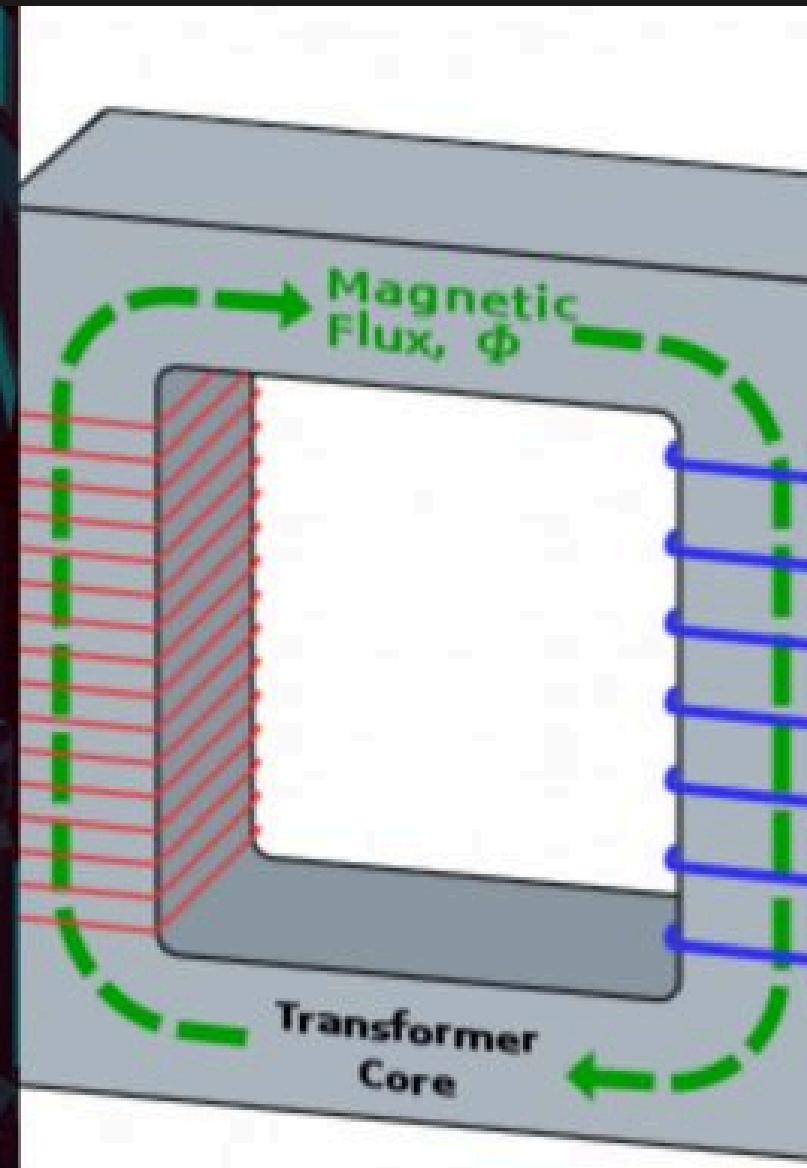


Figure 1: The Transformer - model architecture.

Transformers
at school

Transformers
at college

Transformers
today

Tokenization

GPT-3 Codex

This is just some random text that ChatGPT will break into tokens. The tokenization process will help the LLM in better understanding the text and training a better performing model.

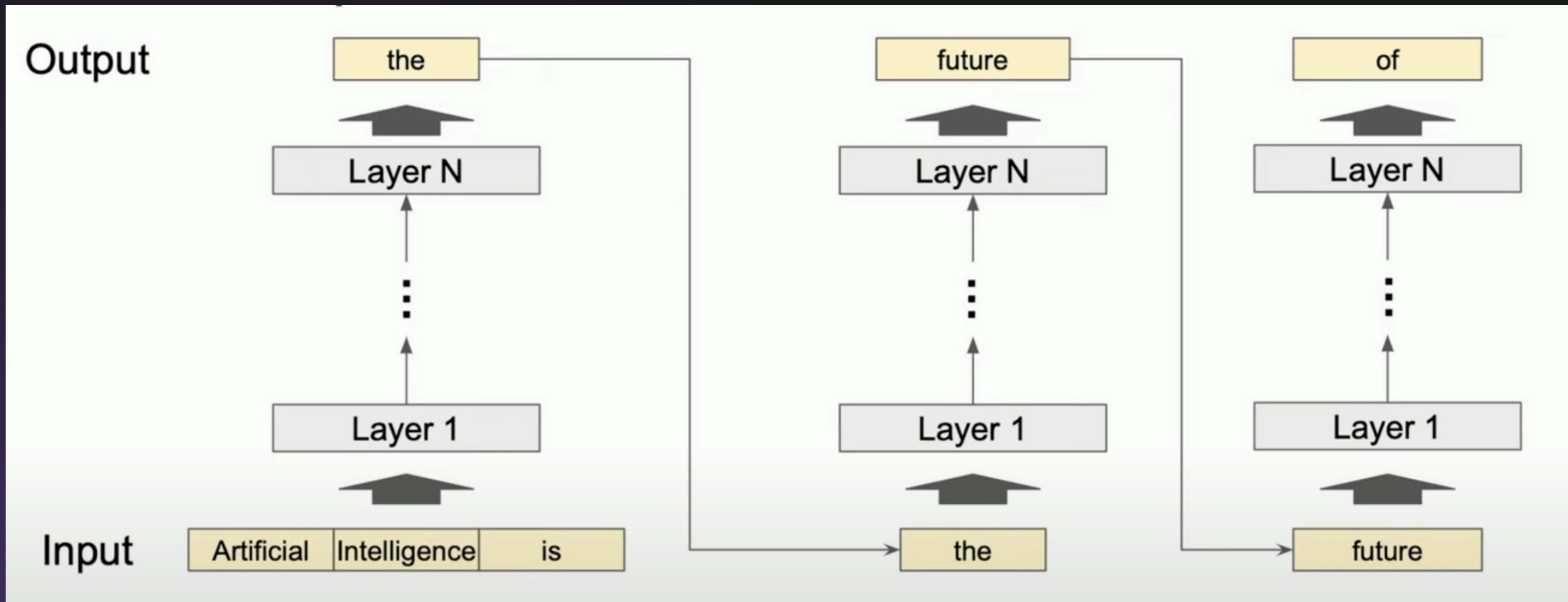
Clear Show example

Tokens	Characters
36	182

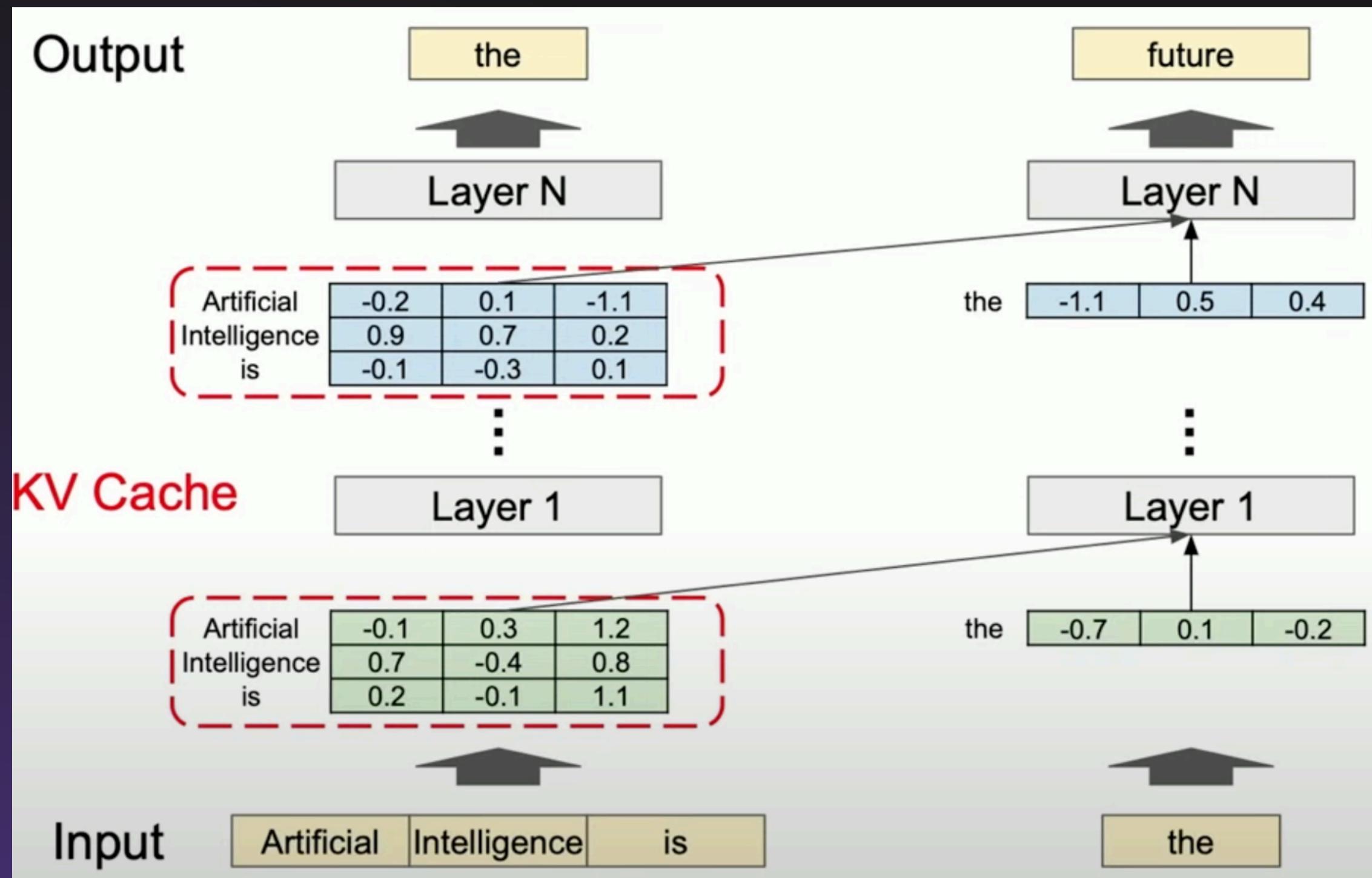
This is just some random text that ChatGPT will break into tokens. The tokenization process will help the LLM in better understanding the text and training a better performing model.

TEXT TOKEN IDS

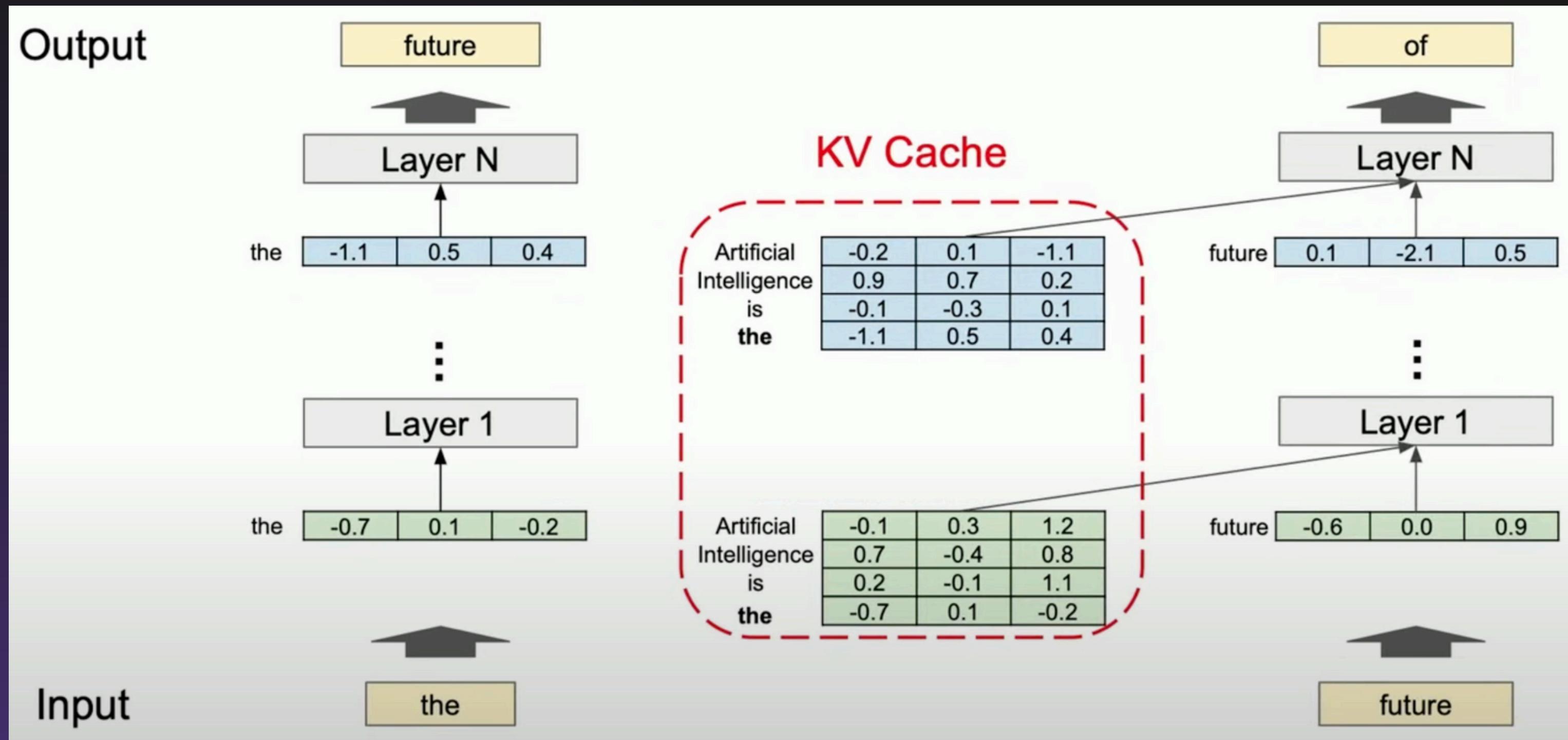
Inference process



KV Cache



KV Cache



Yesterday I went to the cinema to see a __

omelette

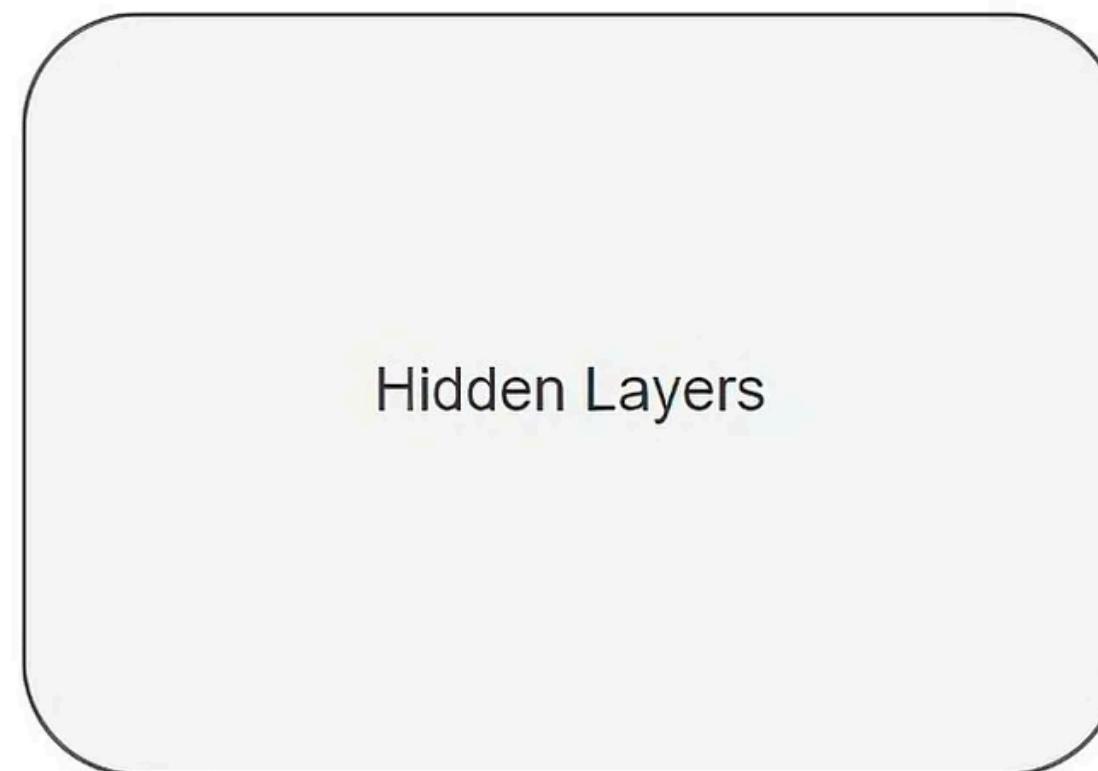
like

film

documental

love

Temperature = 1

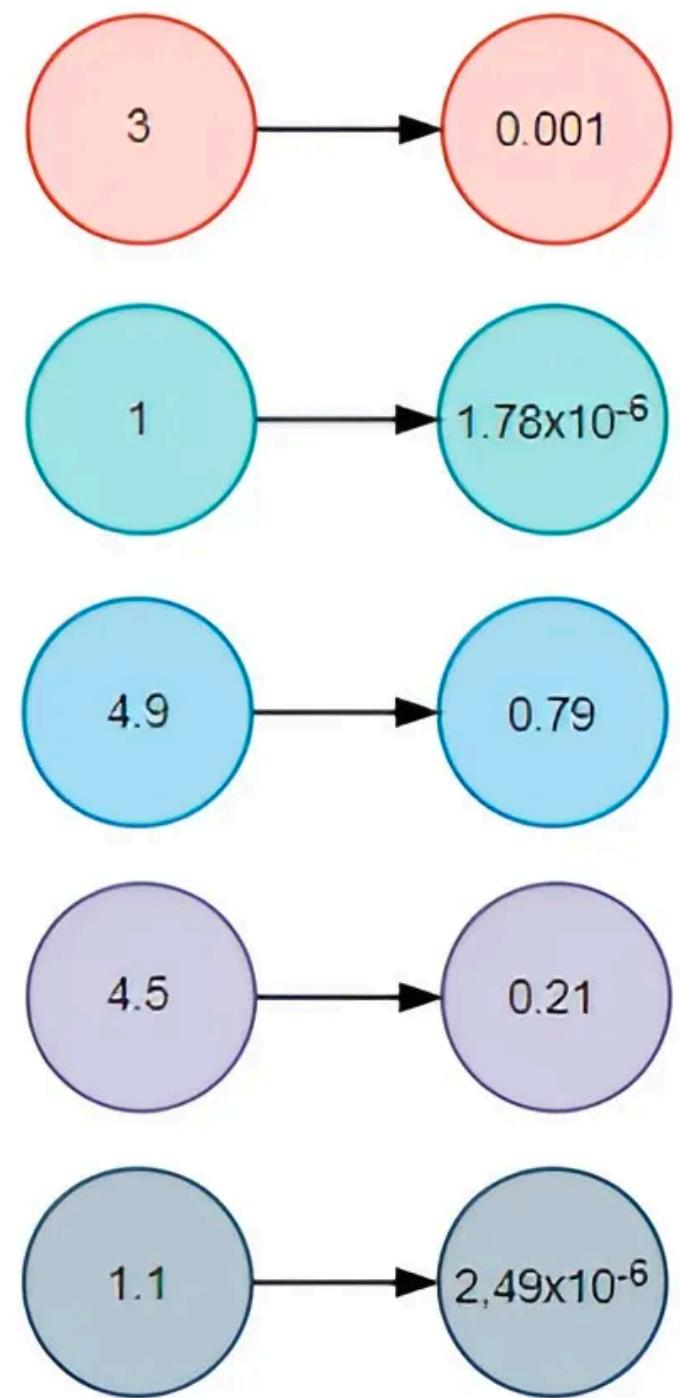


logits probabilities

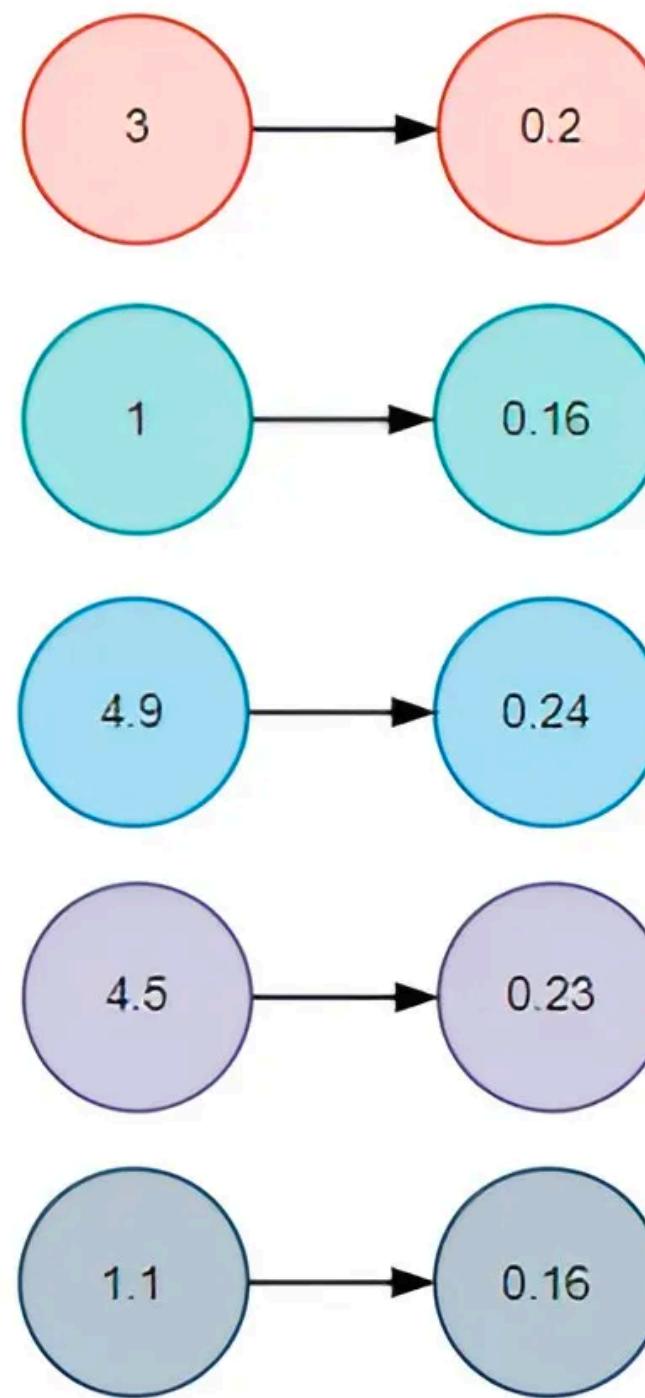
Random
selection

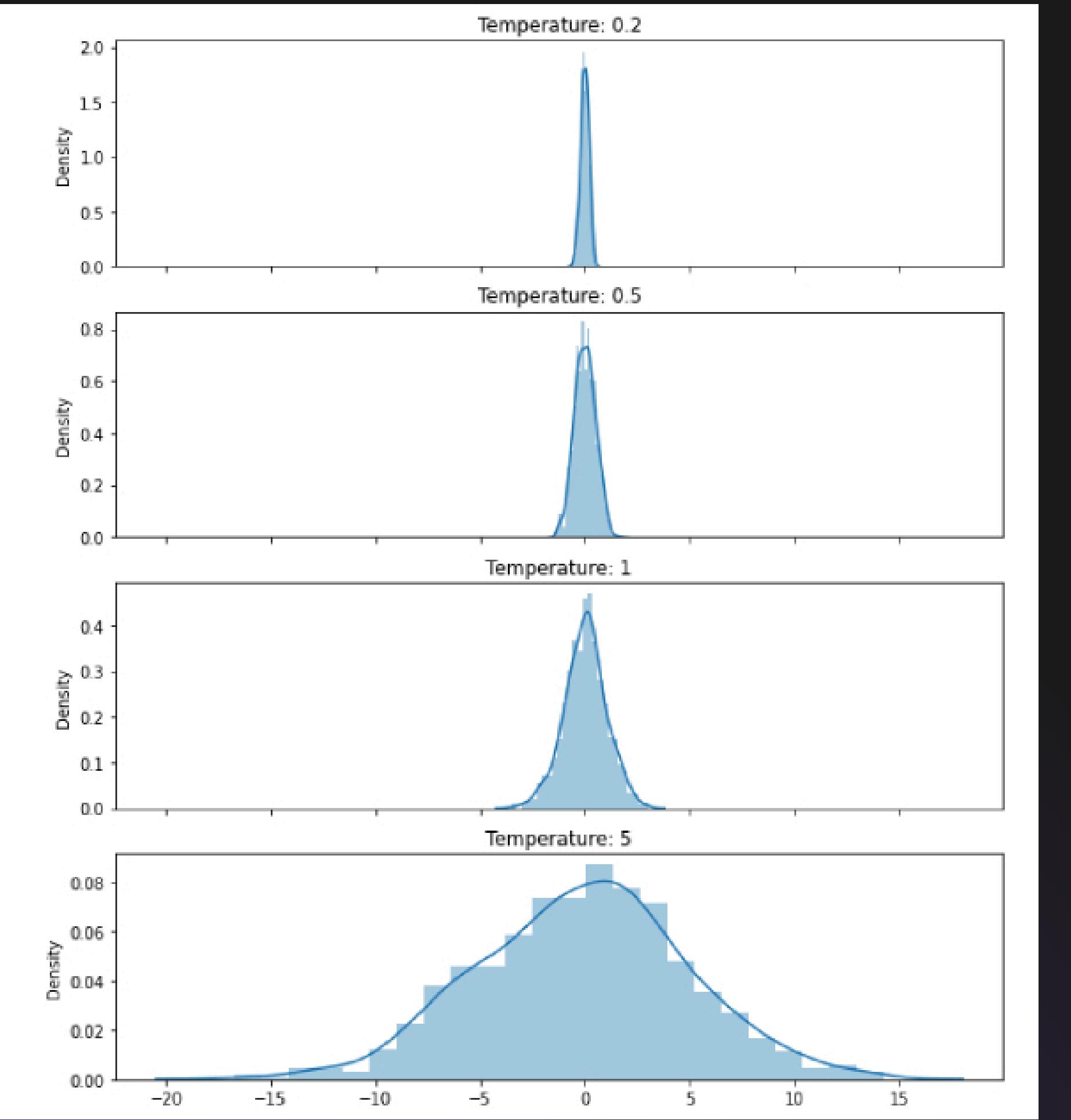
...

Temperature = 0.3



Temperature = 1





Yesterday I went to the cinema to see a __

omelette

0.08

like

0.01

film

0.54

documental

0.36

love

0.01

Cumulative

film

0.54

documental

0.9

omelette

0.98

love

0.99

like

1.0

p = 0.95

film

0.54

documental

0.36

omelette

0.08

love

0.01

like

0.01

New probabilities

film

0.59

documental

0.41

omelette

0.08

Yesterday I went to the cinema to see a ___

omelette

0.08

like

0.01

film

0.54

documental

0.36

love

0.01

$k = 2$

film

0.54

documental

0.36

omelette

0.08

love

0.01

like

0.01

New probabilities

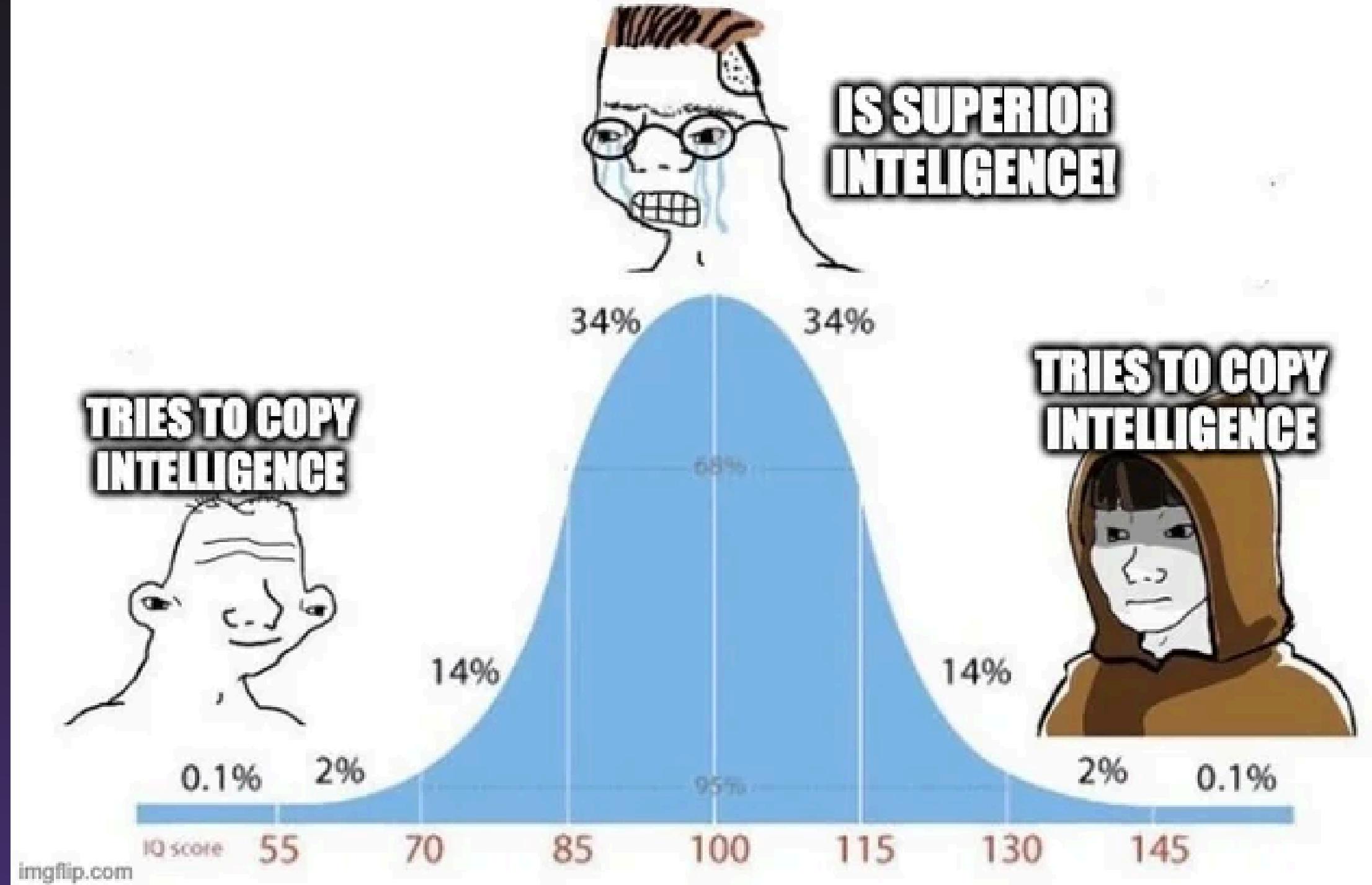
film

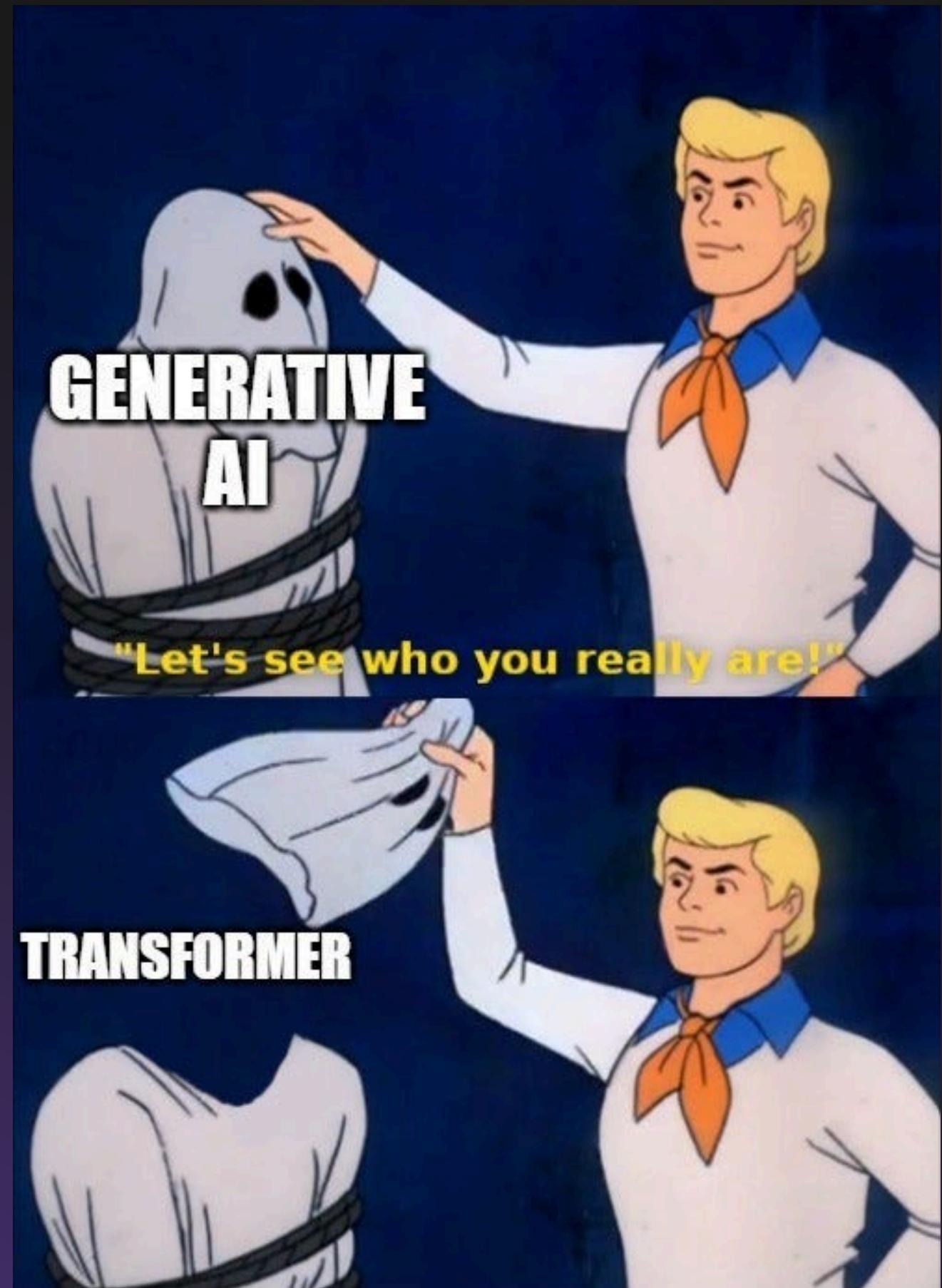
0.59

documental

0.41

ARTIFICIAL INTELLIGENCE...



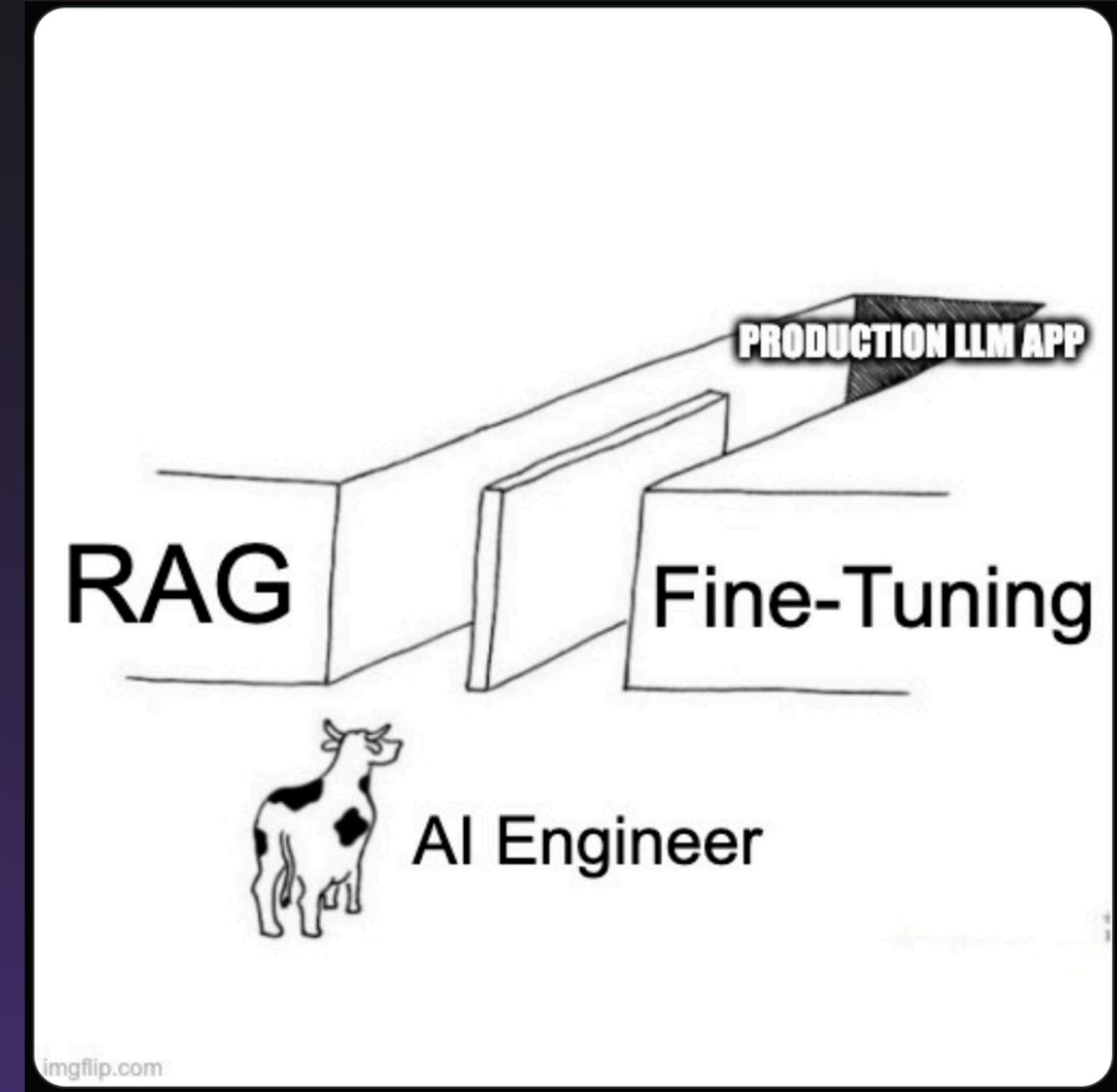


Finetuning LLMs

Pranav Reddy
Cofounder, Xylem AI
@pranavreddyg (Twitter/X)



Finetuning/RAG?



Finetuning/RAG? Asking the right questions

- Does the model lack knowledge of a particular domain?
- Does the context keep changing?
- Do the outputs need to follow any particular format?
- Does prompting solve the issue?
- Does the model not respond with the desired output?
- How often do I need to train?



imgflip.com

**Add examples
to the prompt**

**Condition the
Model with Few-Shot
In-Context Learning**



Asking GPT to be truthful



Fine tuning



Prompt chaining



Pre/post processing user input



All of the above



**PRE-SOFTWARE:
SPECIAL-PURPOSE
COMPUTER**



**SOFTWARE 1.0:
DESIGN
THE ALGORITHM**



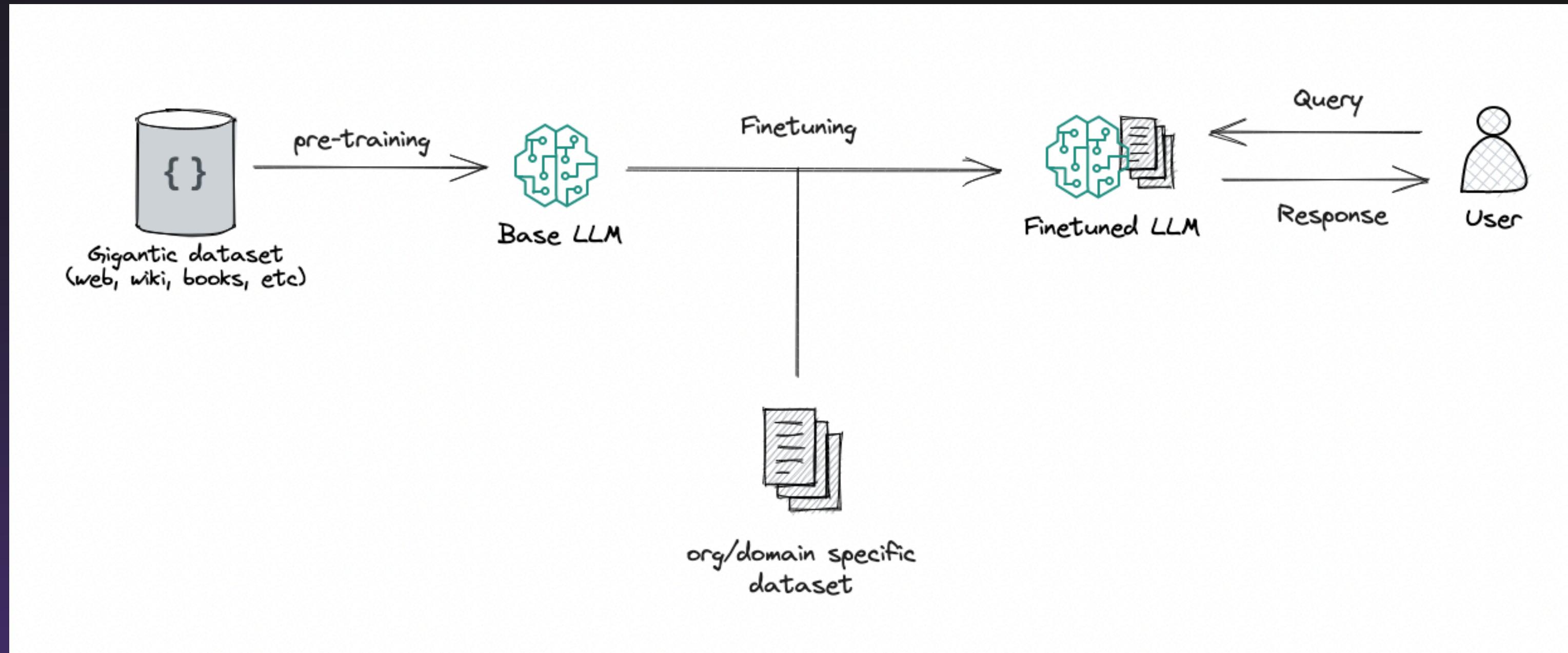
**SOFTWARE 2.0:
DESIGN
THE DATASET**



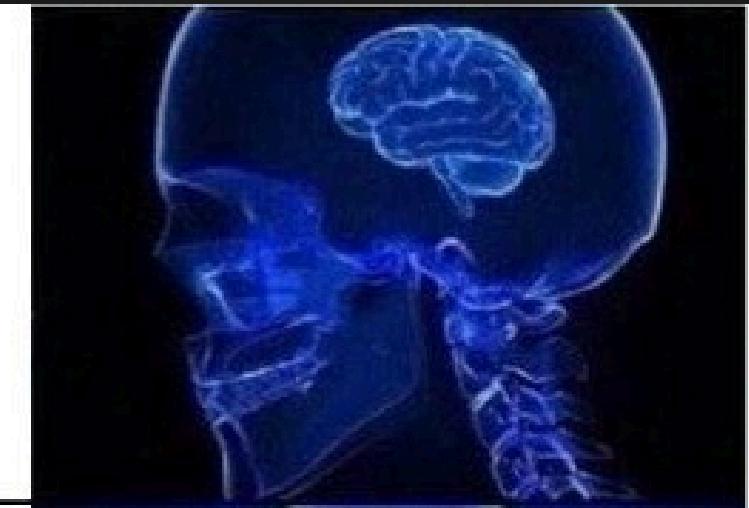
**SOFTWARE 3.0:
DESIGN
THE PROMPT**



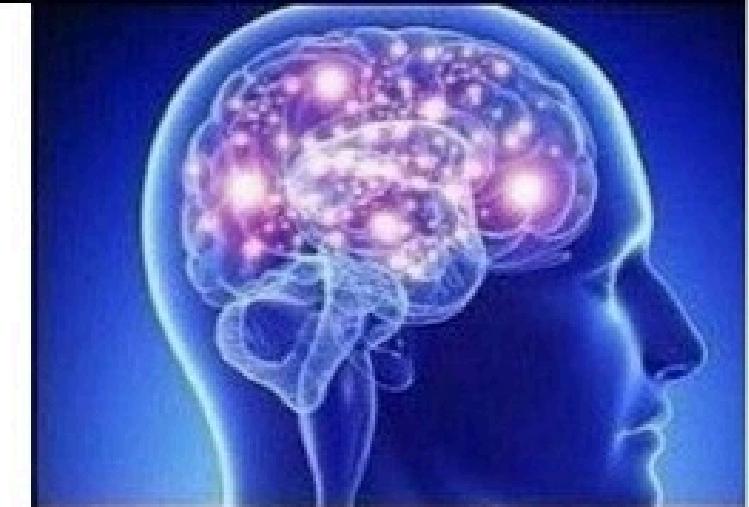
Fine-tuning



**Prompt
engineering**



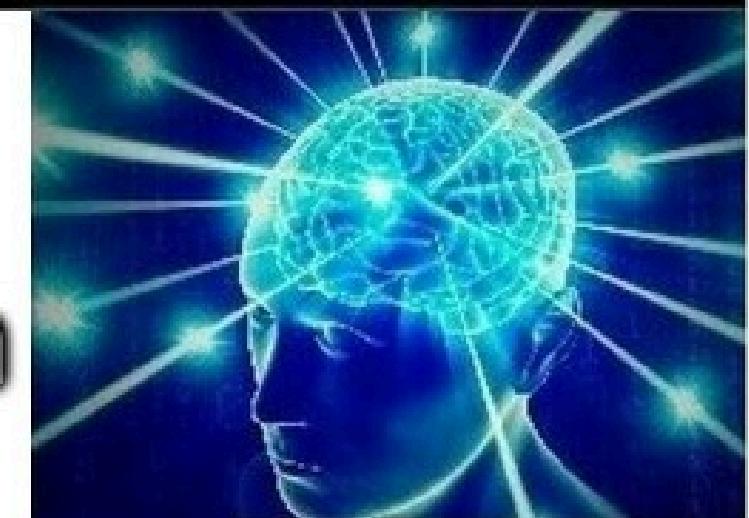
**Finetuning
LLM**



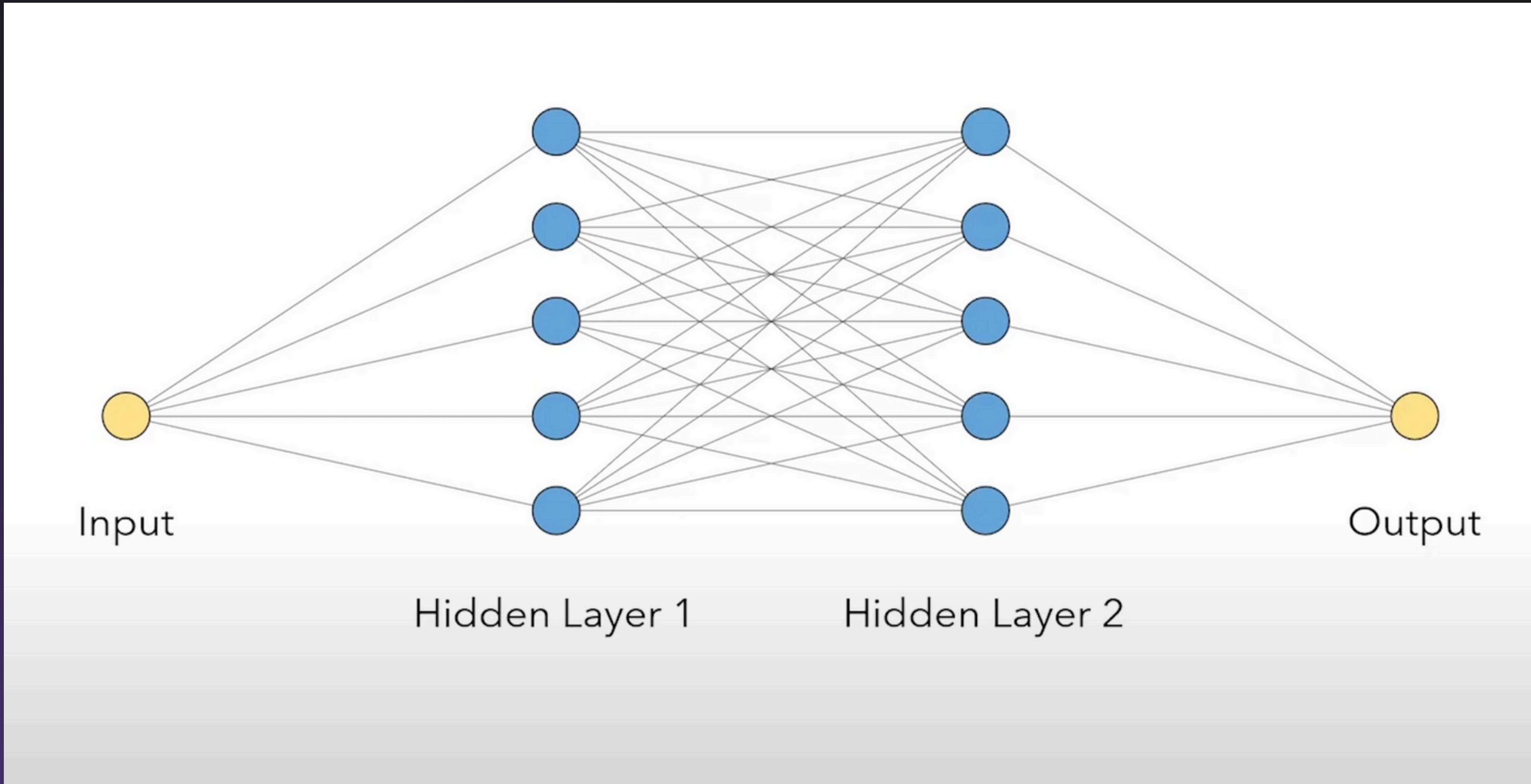
**Pretraining
LLM**



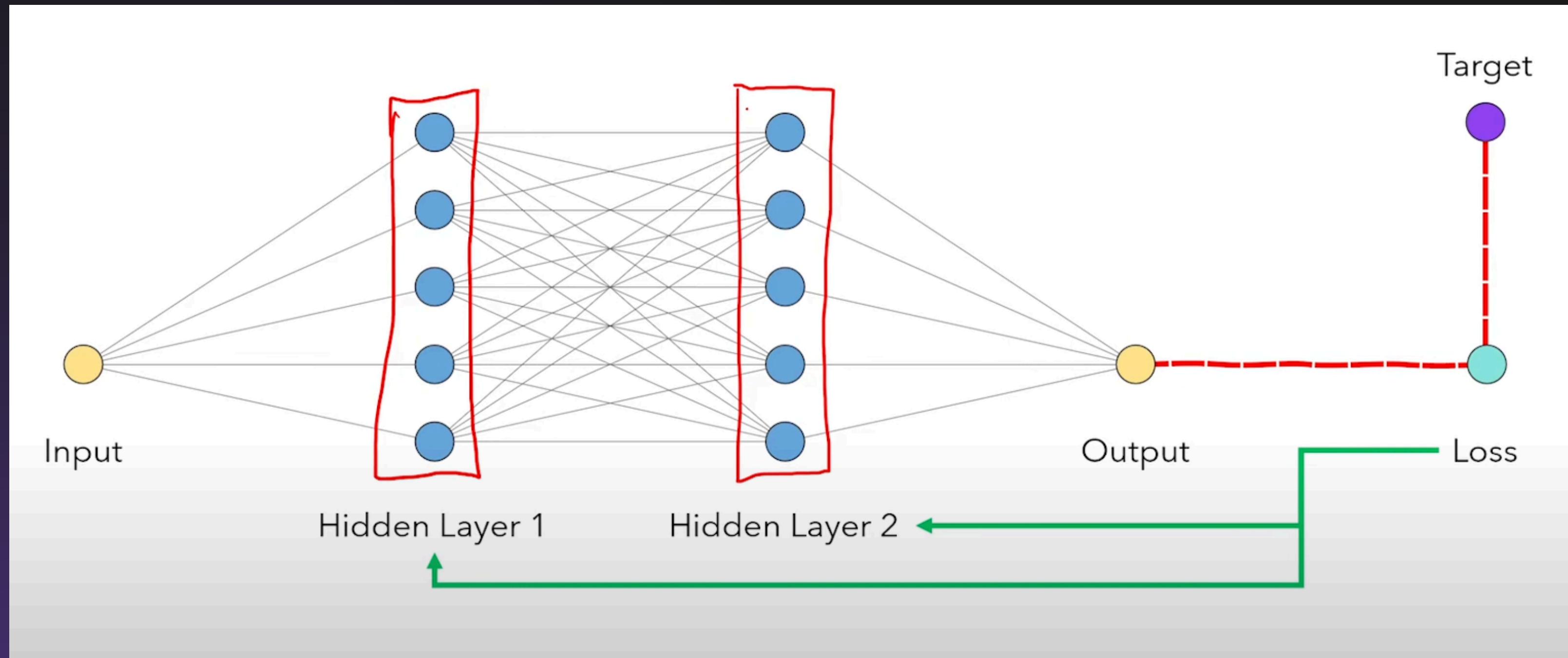
**Building
foundational
LLM from scratch**



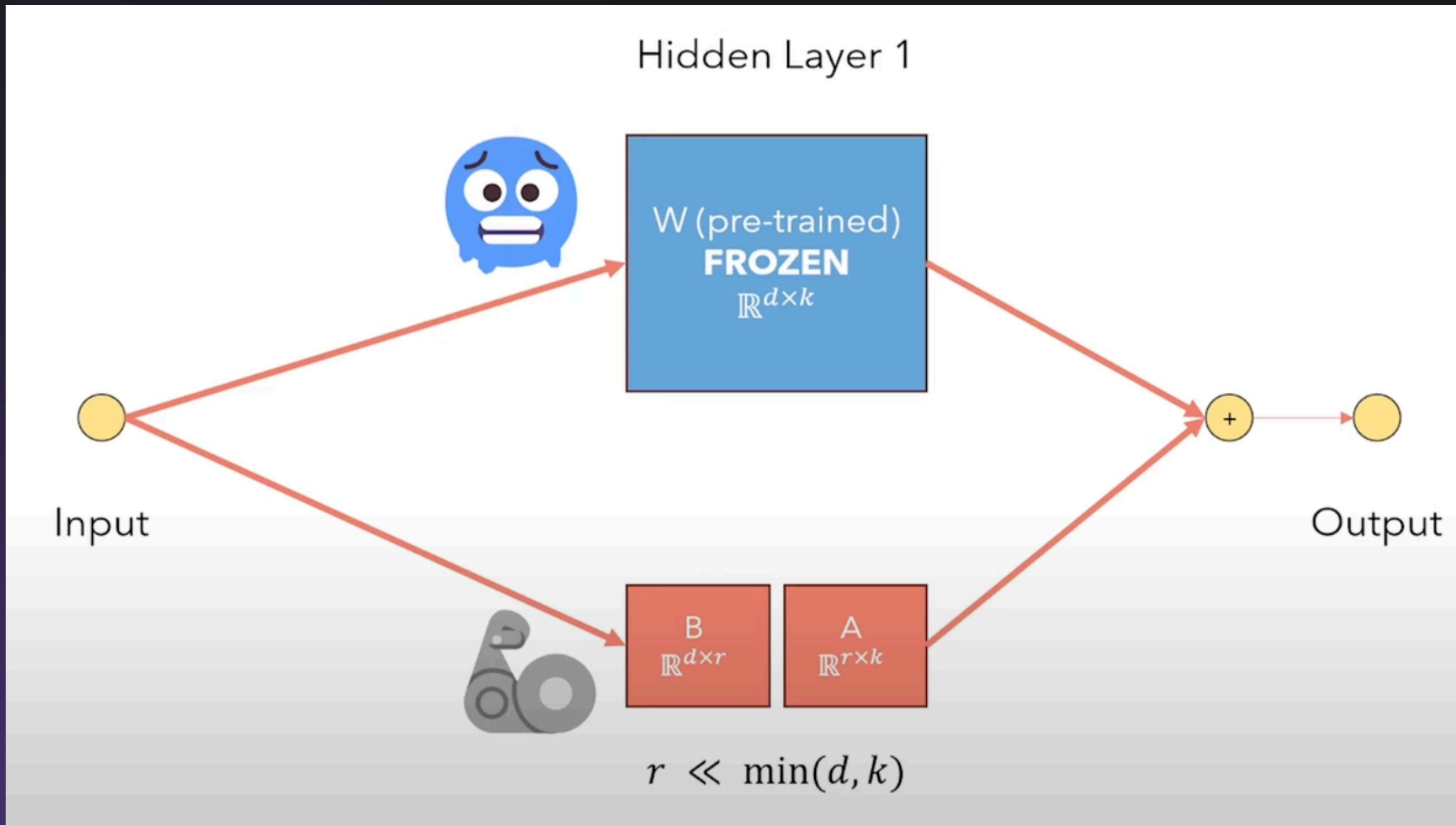
Neural Networks



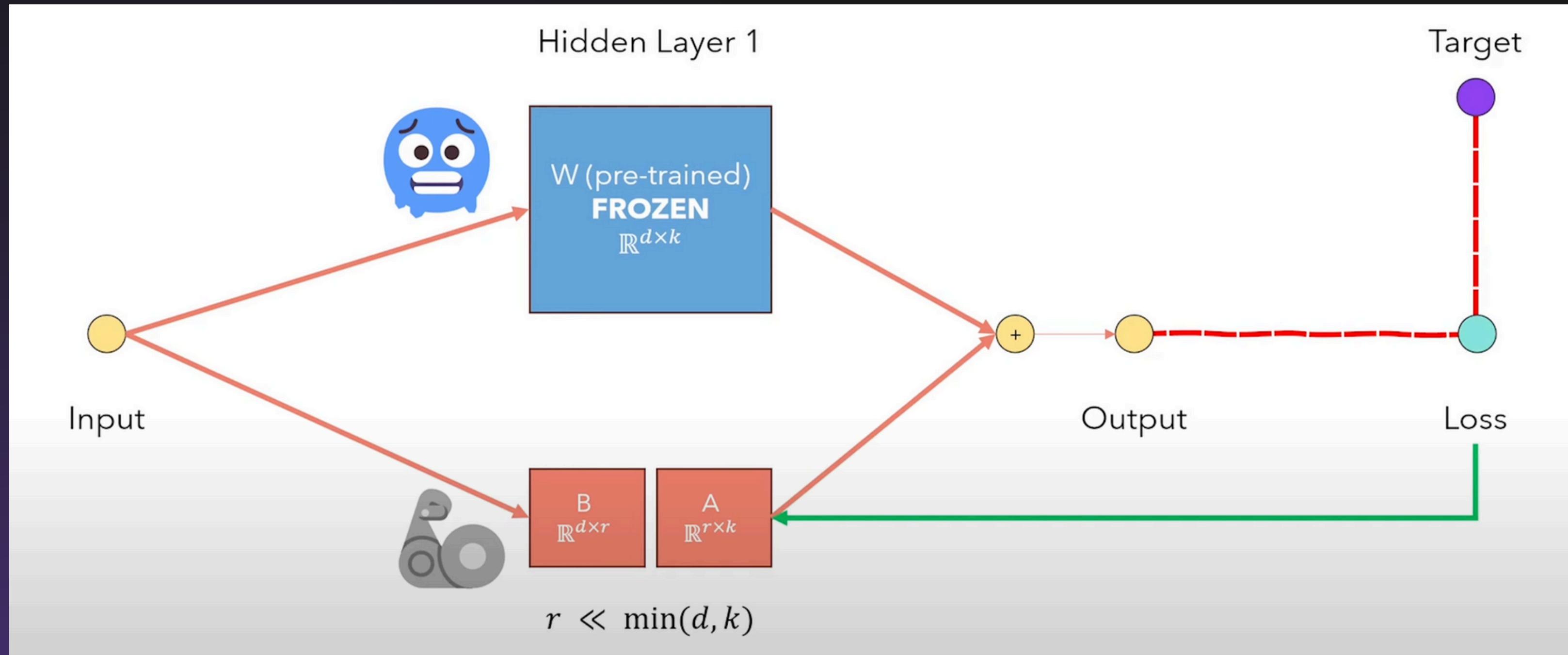
Training Neural Networks



Low Rank Adaptation of LLMs (LoRA)



Low Rank Adaptation of LLMs (LoRA)



Is the column vector a linear combination of anything that comes before it?

Nothing comes before it.

No scalar multiplication of $[1, 1]$ that can reach $[1, 2]$.

$1([1, 2, 1]) + 2(1, 1, 3)$ can reach $[3, 4, 7]$.

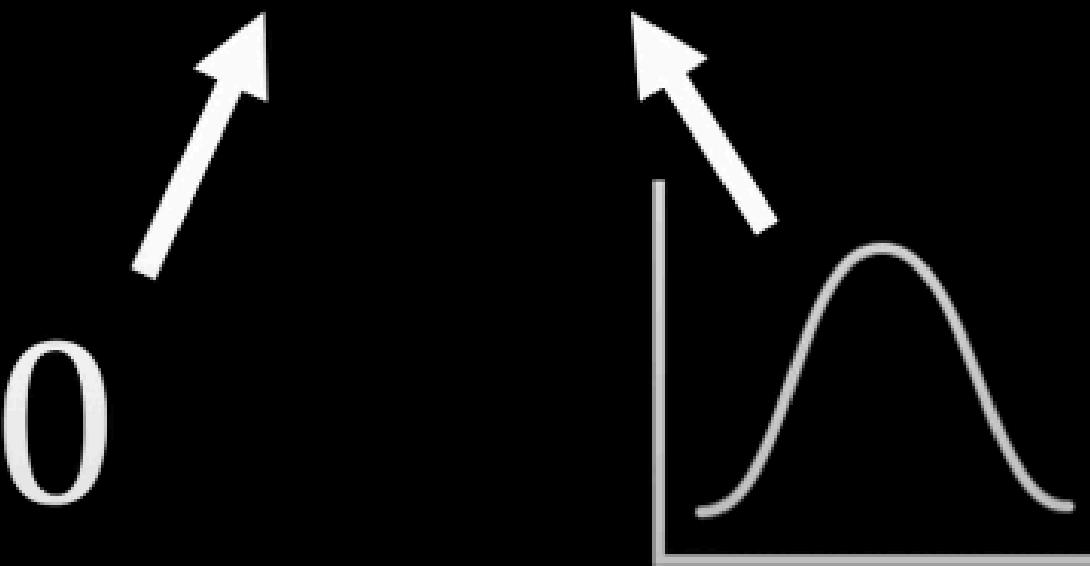
$2([1, 1, 3])$ can reach $[2, 2, 6]$.



$$\begin{bmatrix} 1 & 1 & 3 & 2 \\ 2 & 1 & 4 & 2 \\ 1 & 3 & 7 & 6 \end{bmatrix} = \text{rank } 2$$

Rank decomposition

$$W + BA$$



Rank decomposition

0.15	-0.14	-0.21	0.612
-0.22	0.204	0.308	-0.86
-0.30	-0.16	0.634	0.147
-0.07	-0.2	0.246	0.523

ΔW

Shape: (4, 4)

0.3	-0.14
-0.42	0.201
0.46	0.38
0.5	0.14

B

Shape: (4, 2)

A

Shape: (2, 4)

Scaling factor

$$W_0 + \Delta W = W_0 + BA$$

$B \in \mathbb{R}^{dxr}$
 $A \in \mathbb{R}^{rxk}$
 $\text{rank } r \ll \min(d, k)$

The diagram illustrates the components of the update term BA . A bracket under BA indicates its total dimension. Two arrows point from the labels "Scaling factor" and "Rank" to the terms $\frac{\alpha}{r}$ and r respectively, which are part of the expression $\frac{\alpha}{r} \cdot r$.

Other parameters

- **bnb_4bit_quant_type** - This sets the quantization data type in the bnb.nn.Linear4Bit layers. Options are FP4 and NF4 data types which are specified by fp4 or nf4.
- **bnb_4bit_compute_dtype** - This sets the computational type which might be different than the input type. For example, inputs might be fp32, but computation can be set to bf16 for speedups.
- **lora_alpha** - scaling factor for the weight matrices. alpha is a scaling factor that adjusts the magnitude of the combined result (base model output + low-rank adaptation). We have set it to 16.

Other parameters

- **lora_dropout** - dropout probability of the LoRA layers. This parameter is used to avoid overfitting. This technique drop-outs some of the neurons during both forward and backward propagation, this will help in removing dependency on a single unit of neurons. We are setting this to 0.1 (which is 10%), which means each neuron has a dropout chance of 10%.
- **r** - This is the dimension of the low-rank matrix
- **bias** - We will not be training the bias in this example, so we are setting that to “none”. If we have to train the biases, we can set this to “all”, or if we want to train only the LORA biases then we can use “lora_only”



I just spent our entire
AI budget on fine-
tuning a model

But at least we have
a working model
now, right?

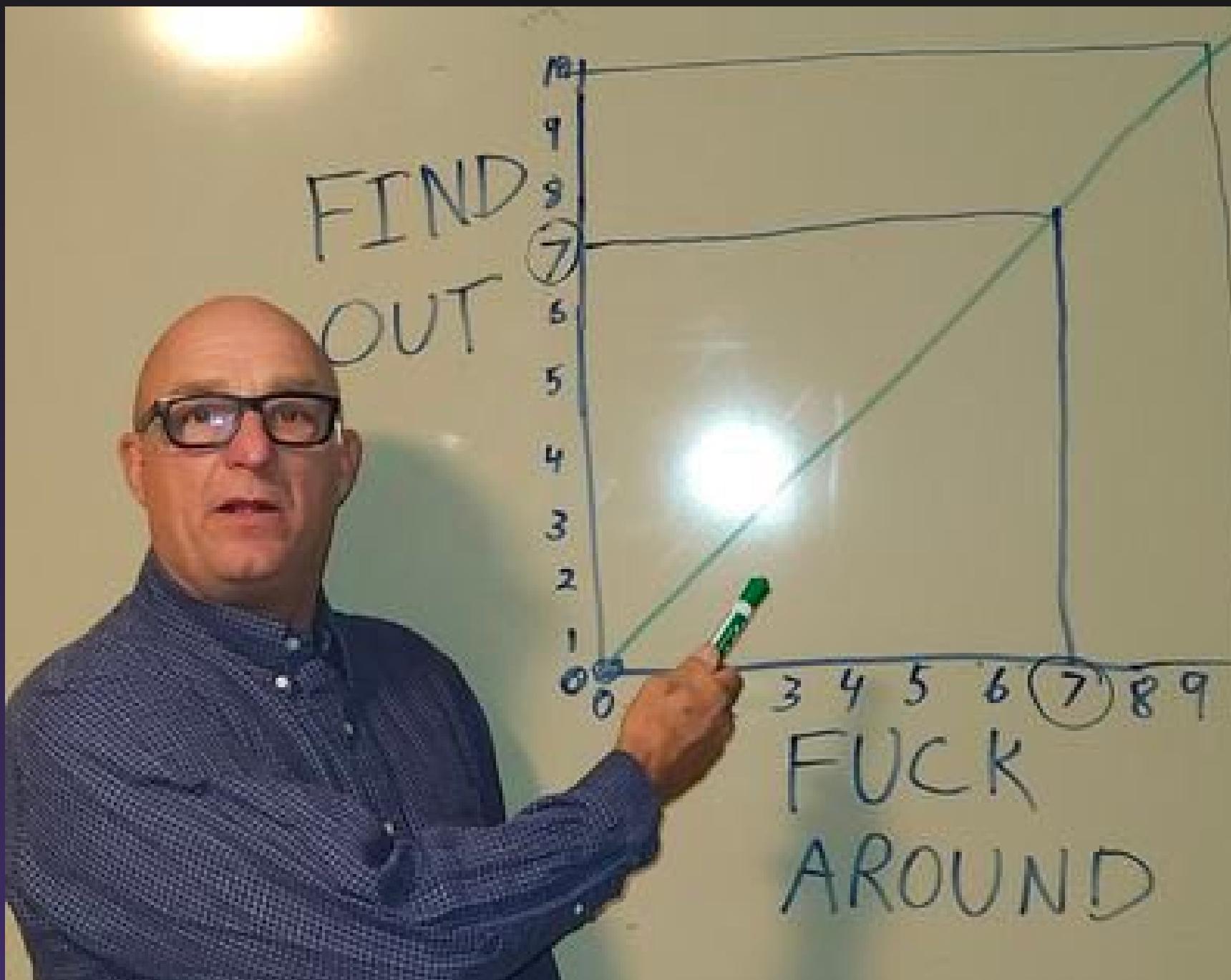
We have a working
model, right?

The world if choosing the right hyperparameters was straightforward

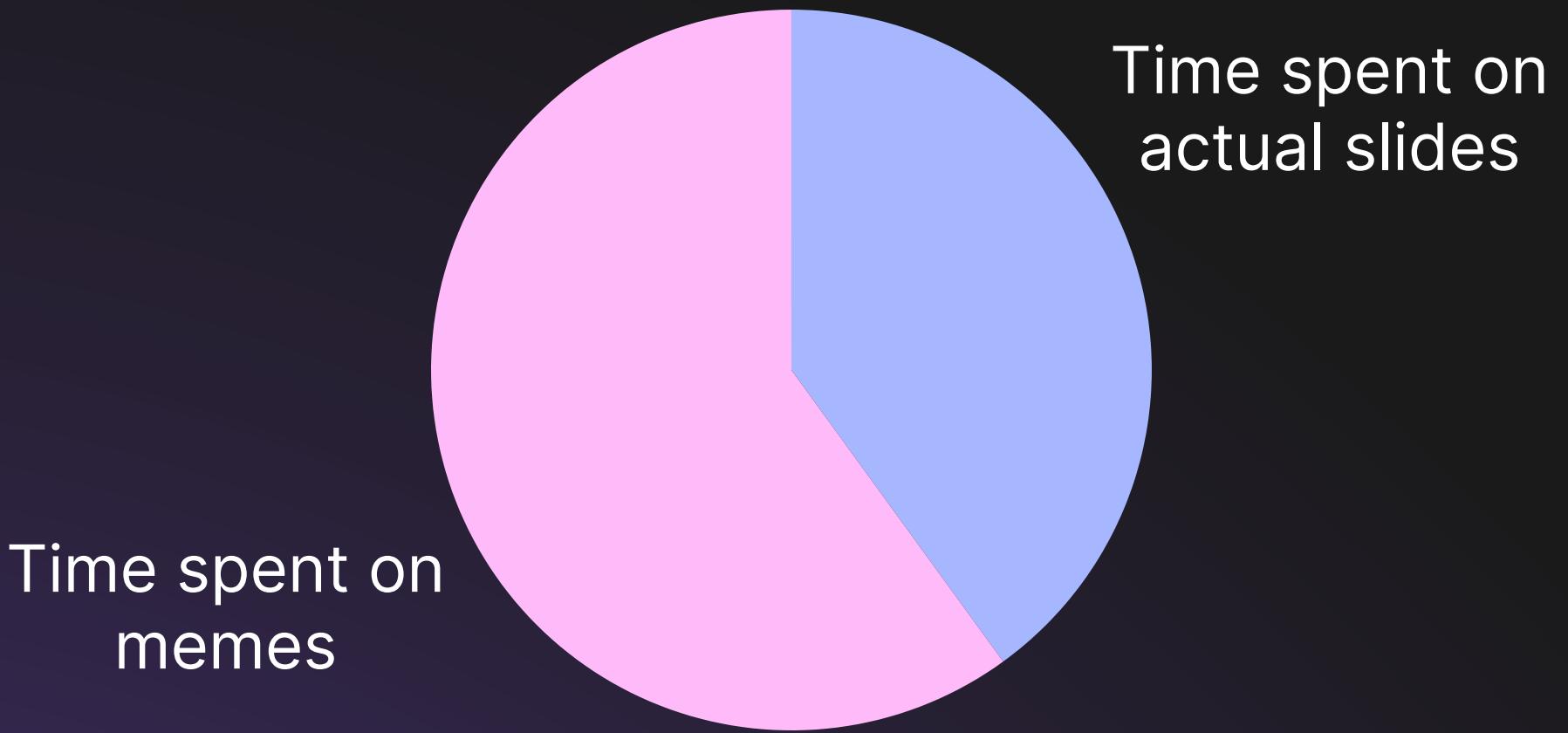


imgflip.com

How to find the best hyperparams?



Thank you



Pranav Reddy
Cofounder, Xylem AI
@pranavreddyg (Twitter/X)

