



# Generative AI

# Fine-Tuning Tutorial

---

Custom LLMs for fun and profit

charlesmartinl4@gmail.com  
+1 (415) 298-6783



# Fine-Tuning LLMs : Custom ‘ChatGPTs’

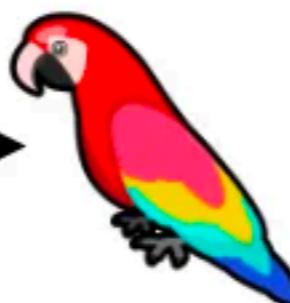


Specific (private) Knowledge  
Base



Gigantic web-scale dataset

→ pre-training →



Base LLM

**Falcon : open-source ChatGPT**



Supervised  
fine-tuning



Fine-tuned LLM



# Fine-Tuning LLMs : ‘Which Model ?’



Falcon



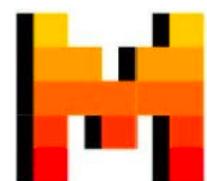
FLAN-T5



Llama2



Meditron



Mistral



Gorilla



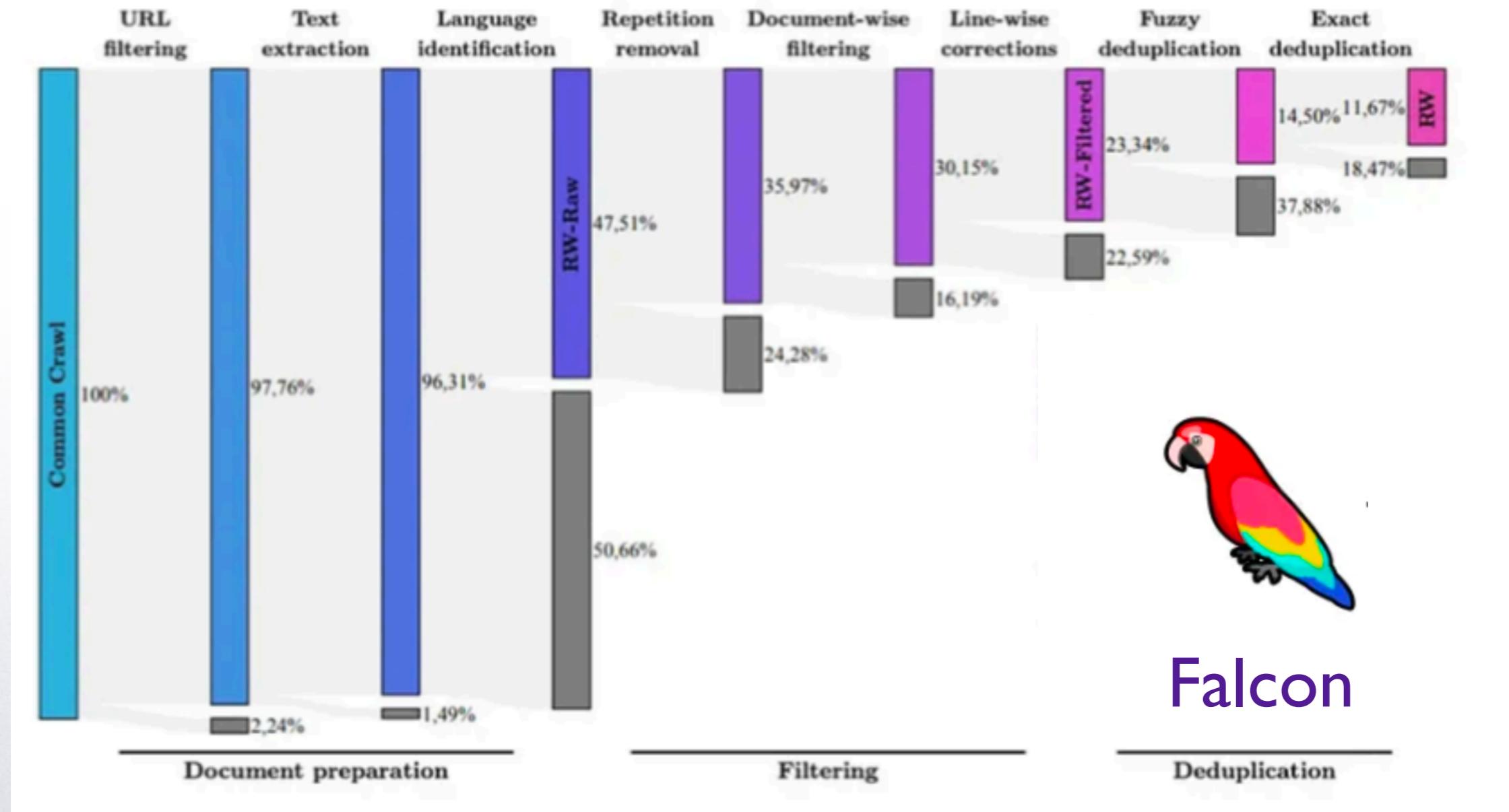
BERT



ALBERT

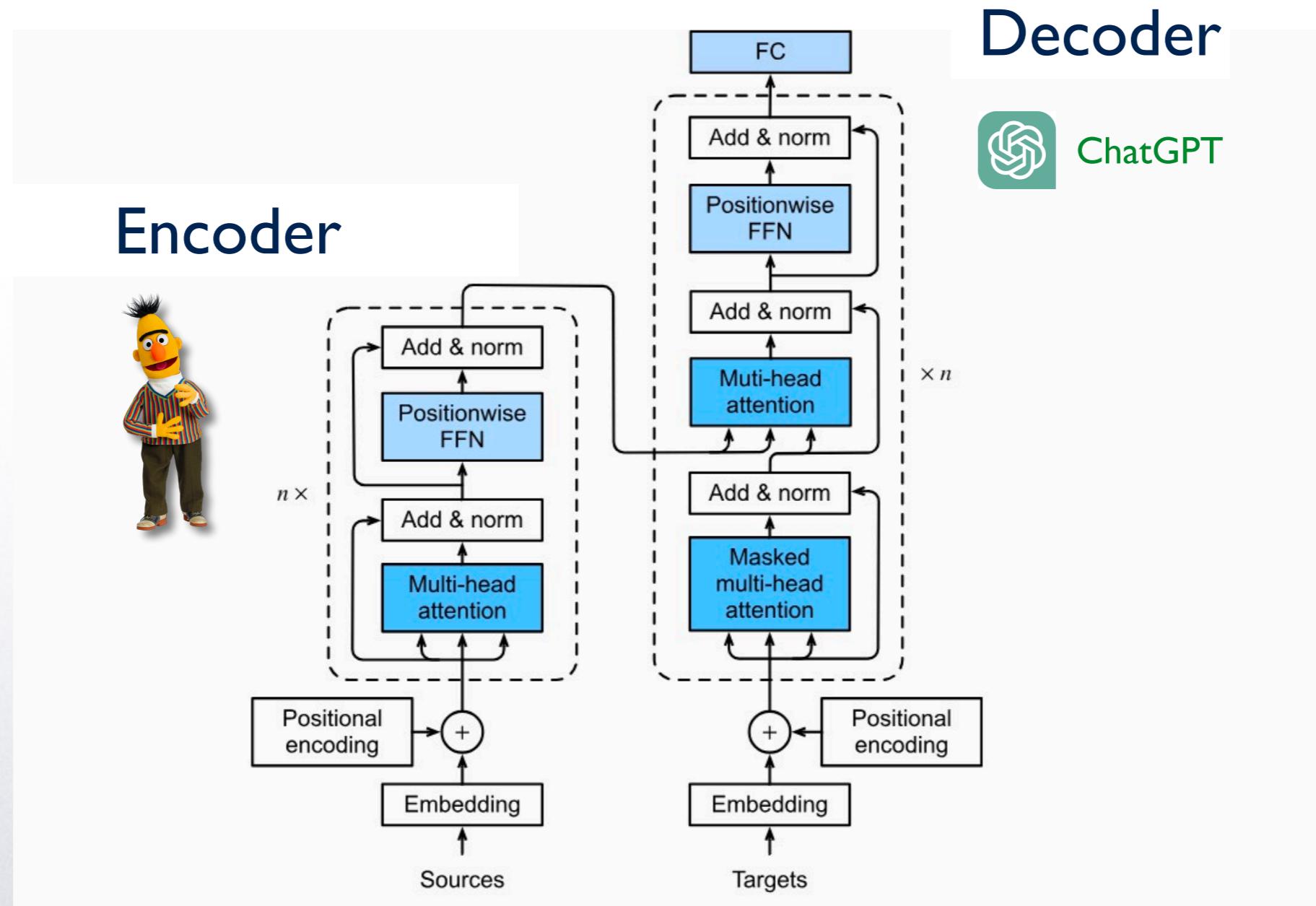


# Fine-Tuning LLMs : Data Prep





# LLM Internals : Transformer Architecture



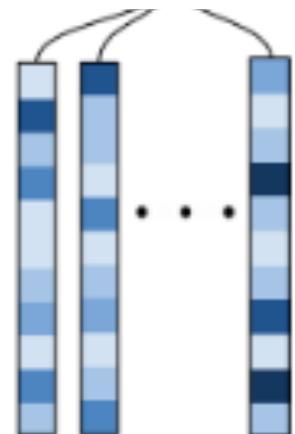
# calculation | consulting

# fine-tuning tutorial

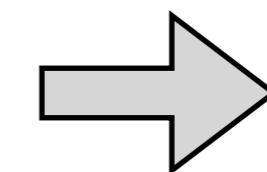
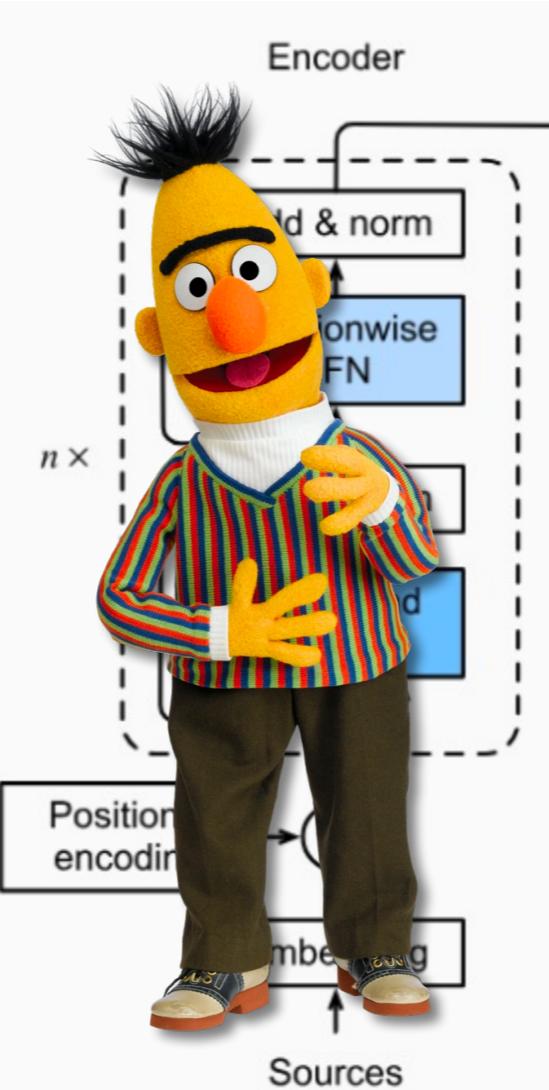


## BERT : Encoder Only

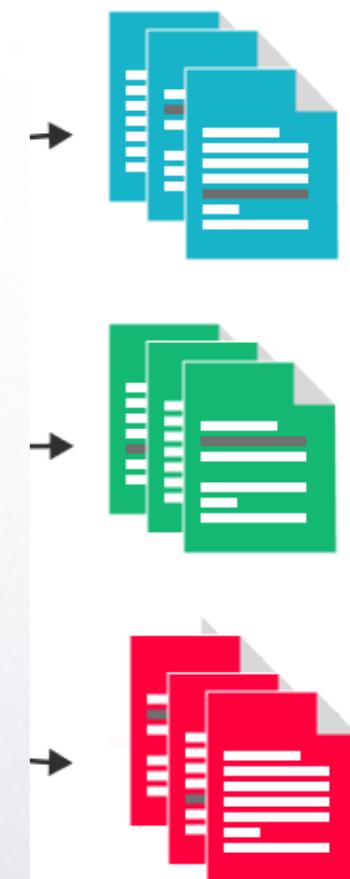
Text Data



Vectors



Classify Text

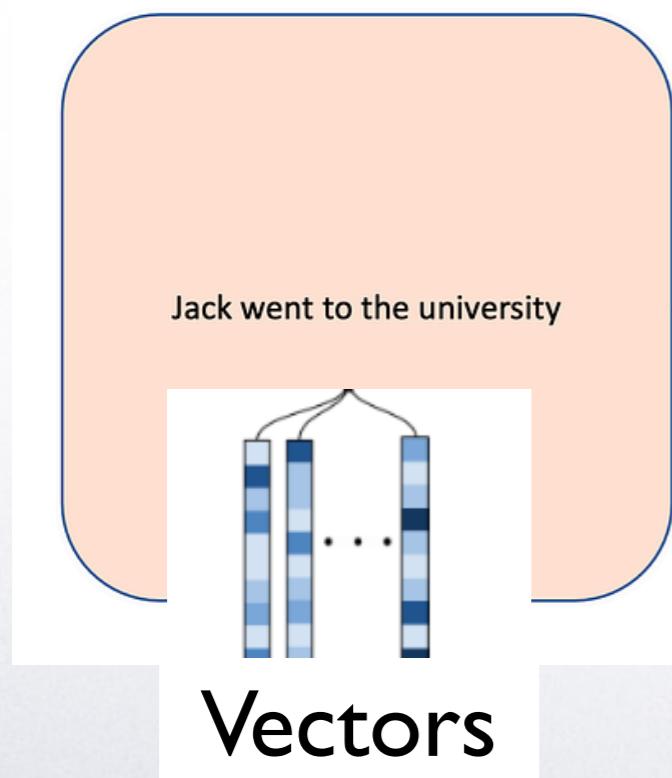


Supervised: Predict Label

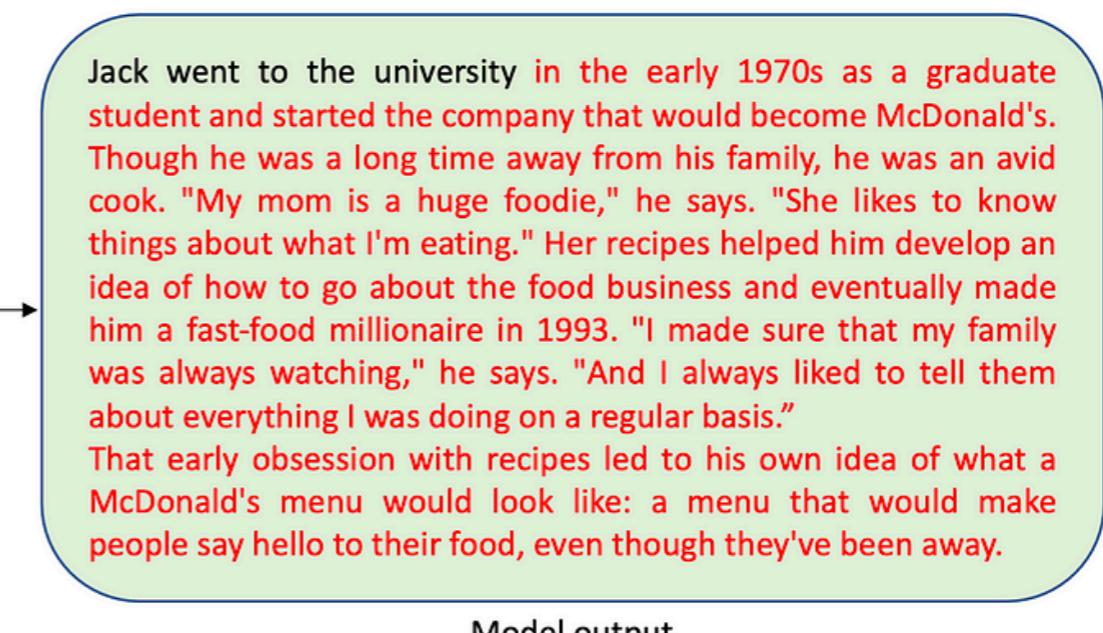


# GPT : Decoder Only

## Input Text



## Generate Text



## Self-Supervised: Predict Next Word



## T5/Flan-T5 : Encoder-Decoder

"translate English to German: That is good."

"cola sentence: The course is jumping well."

"stsbt sentence1: The rhino grazed on the grass. sentence2: A rhino is grazing in a field."

"summarize: state authorities dispatched emergency crews tuesday to survey the damage after an onslaught of severe weather in mississippi..."

Text To Text Transfer Transformer

T5

"Das ist gut."

"not acceptable"

"3.8"

"six people hospitalized after a storm in attala county."

*Every problem is “Text-to-Text”: translate, summarize, ...*



# Flan-T5 : Fine-tuned Language Net - T5

□

## Instruction finetuning

Please answer the following question.  
What is the boiling point of Nitrogen?

-320.4F

## Chain-of-thought finetuning

Answer the following question by reasoning step-by-step.  
The cafeteria had 23 apples. If they used 20 for lunch and bought 6 more, how many apples do they have?

The cafeteria had 23 apples originally. They used 20 to make lunch. So they had  $23 - 20 = 3$ . They bought 6 more apples, so they have  $3 + 6 = 9$ .

*Multi-task instruction finetuning (1.8K tasks)*

Language model

**Instruction Fine-tuning:** predict token sequences



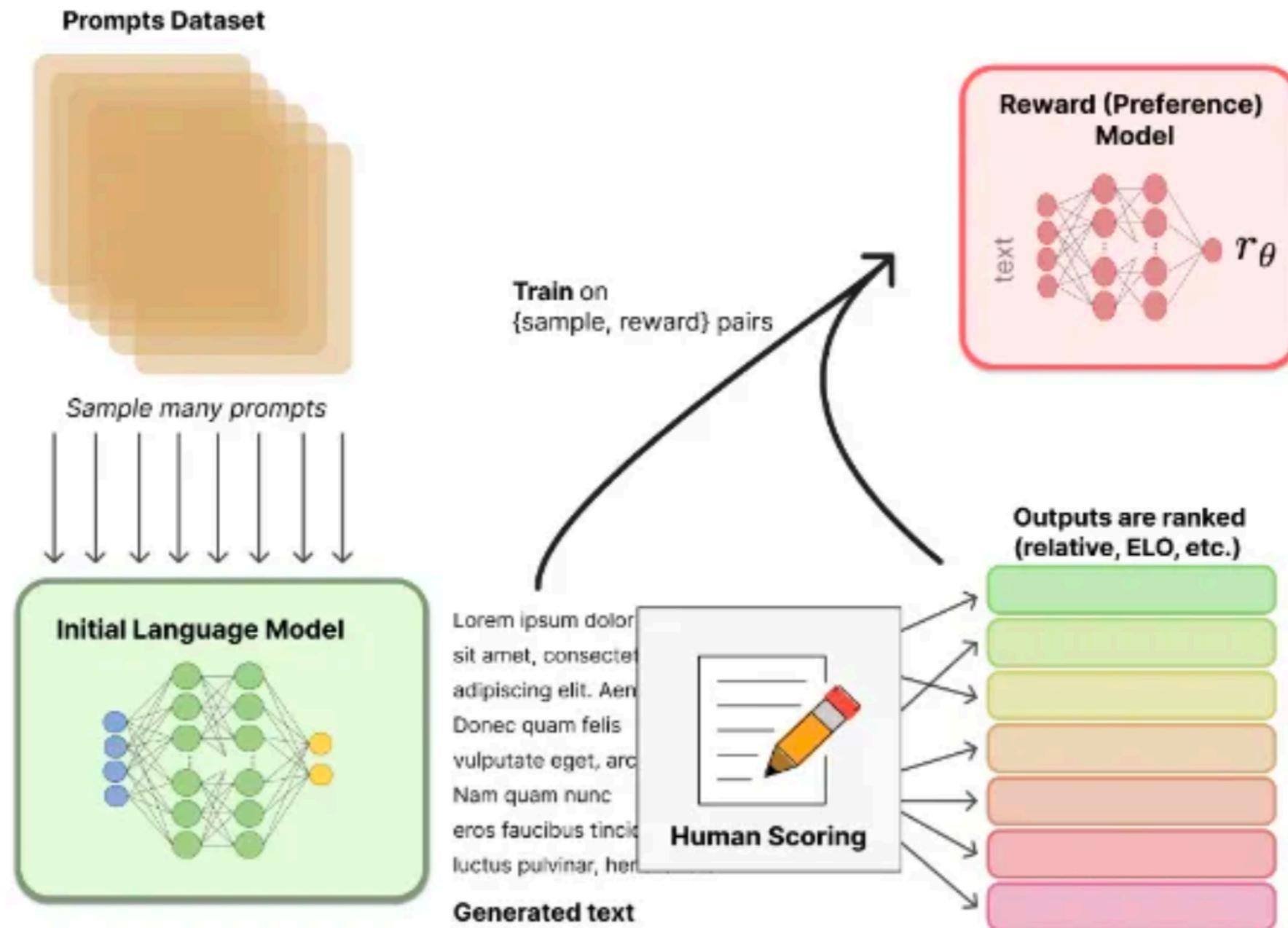
# Reinforcement learning from human feedback (RLHF)



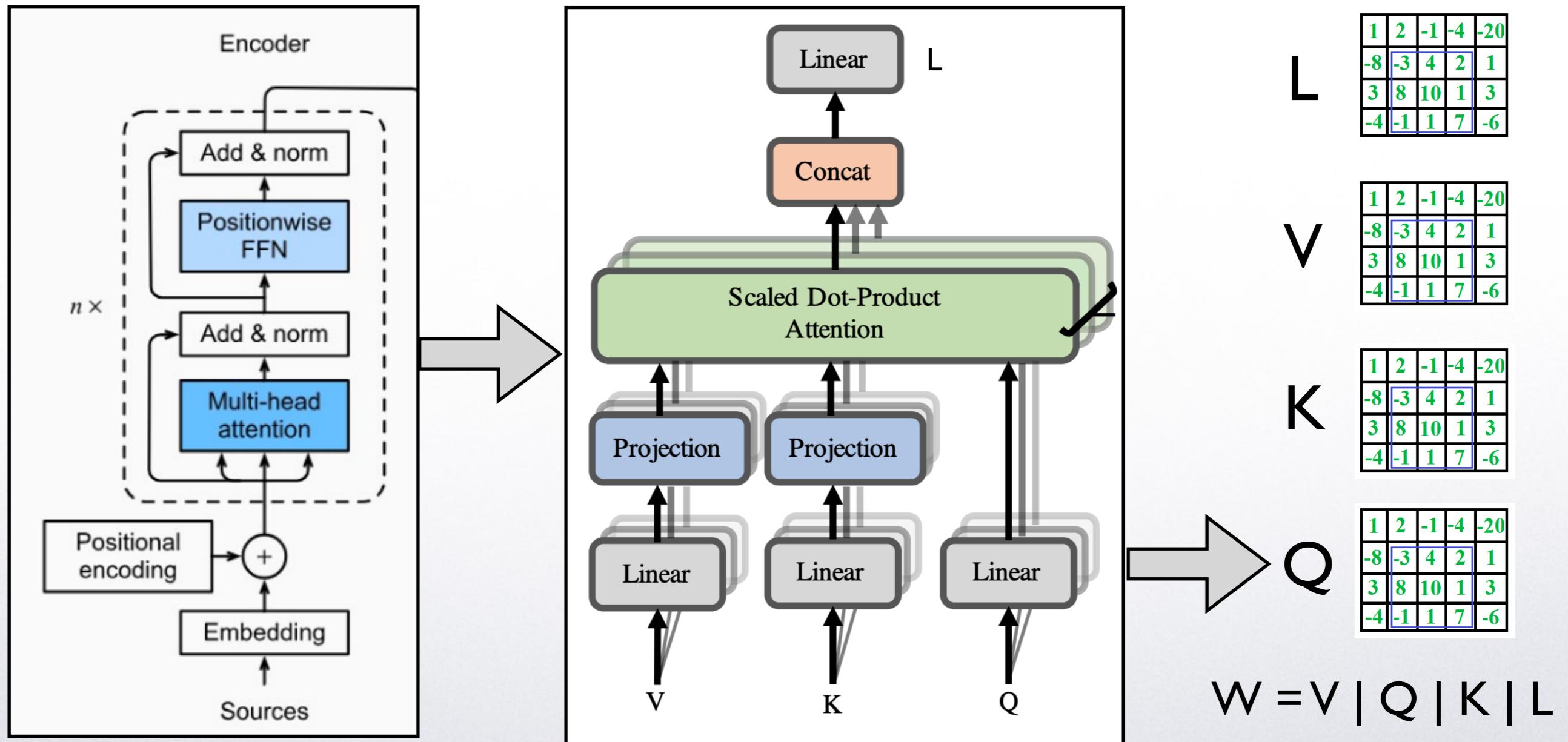
- Maximize helpfulness, relevance
- Minimize harm
- Avoid dangerous topics



# RLHF : Reinforcement Learning from Human Feedback

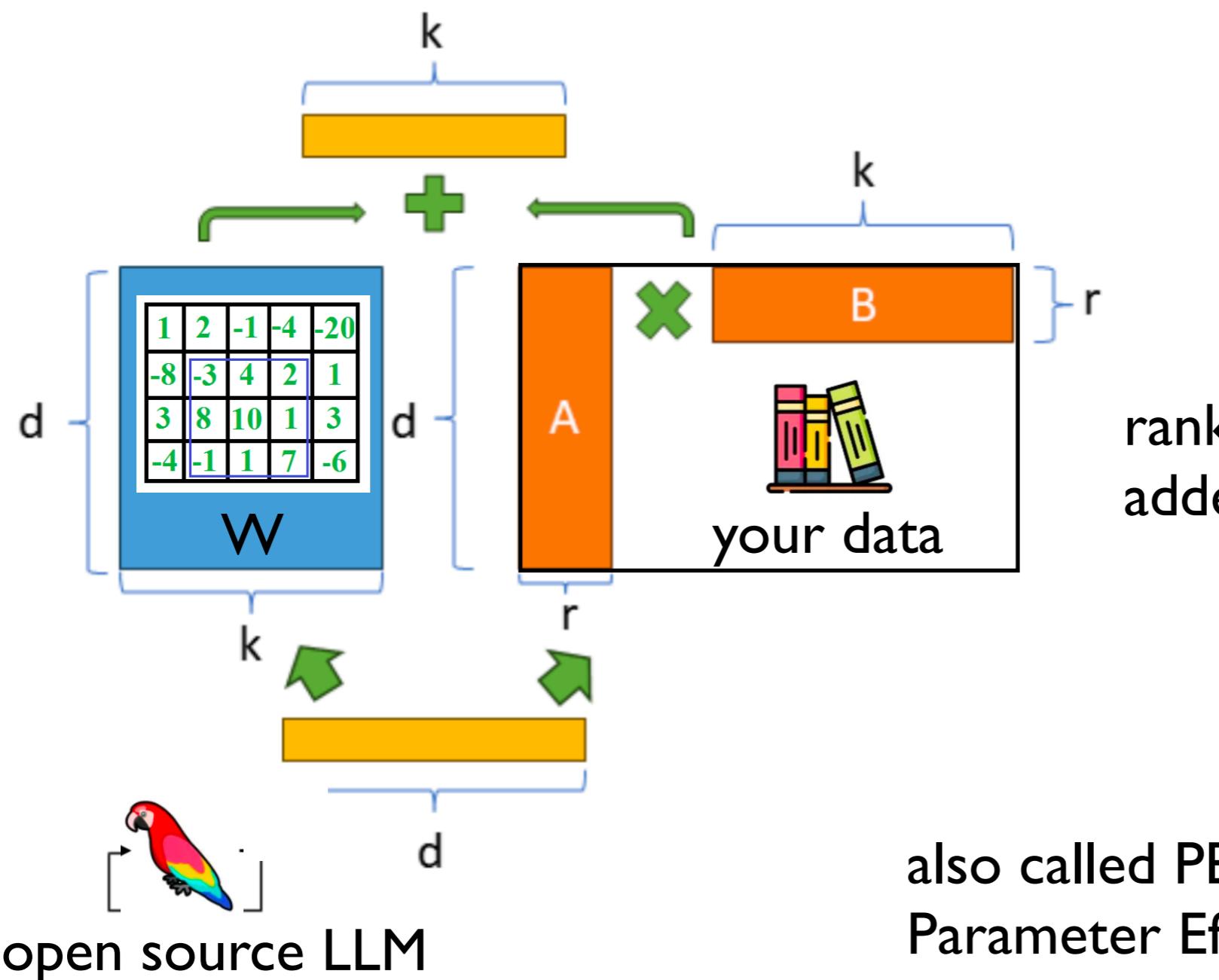


# Fine Tuning : Weight Matrices ( $W$ )



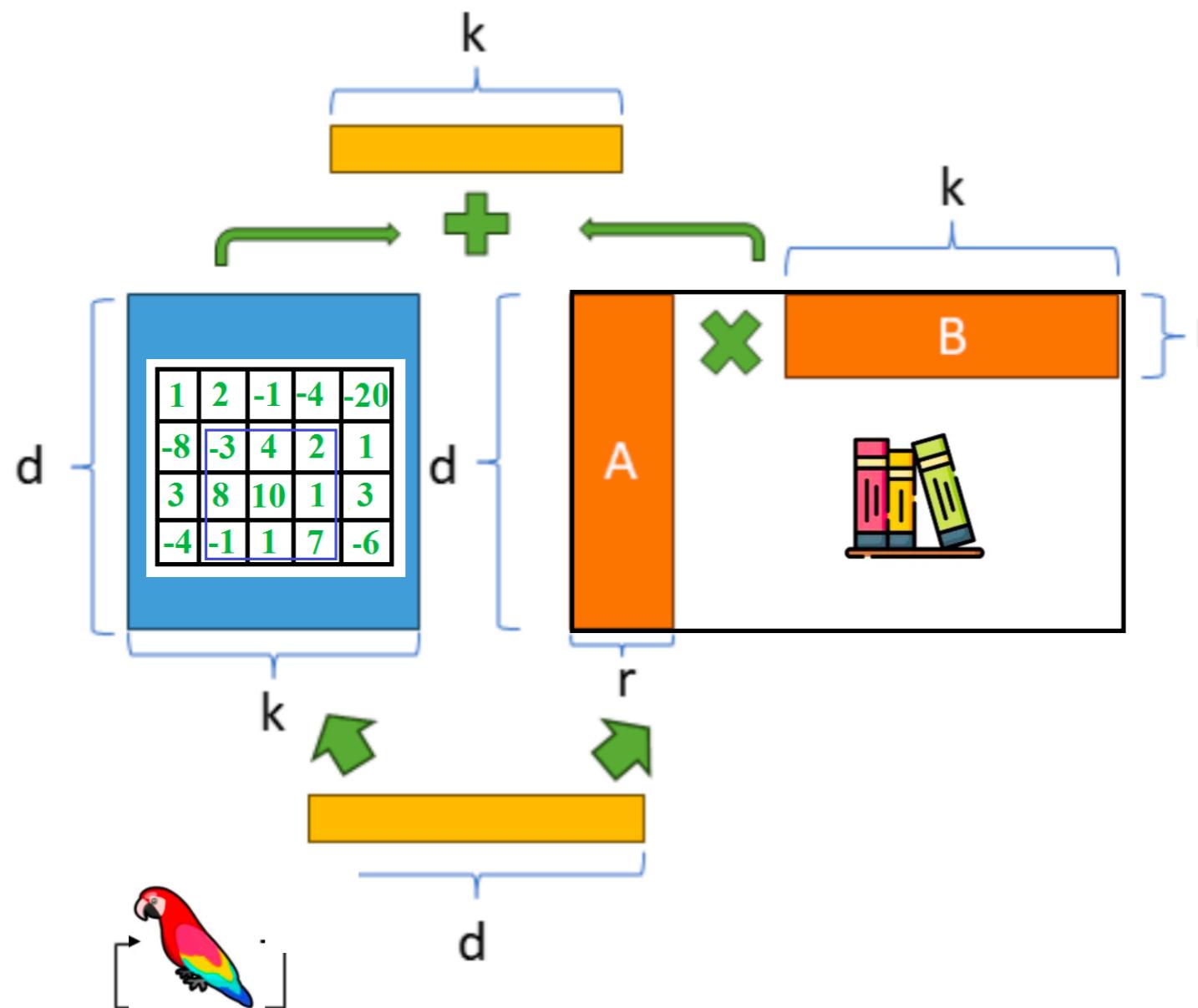


# Fine-Tuning w/LoRA : Low Rank Adaptation





# Fine-Tuning w/LoRA : Low Rank Adaptation



easy fine-tuning

- very fast
- low memory

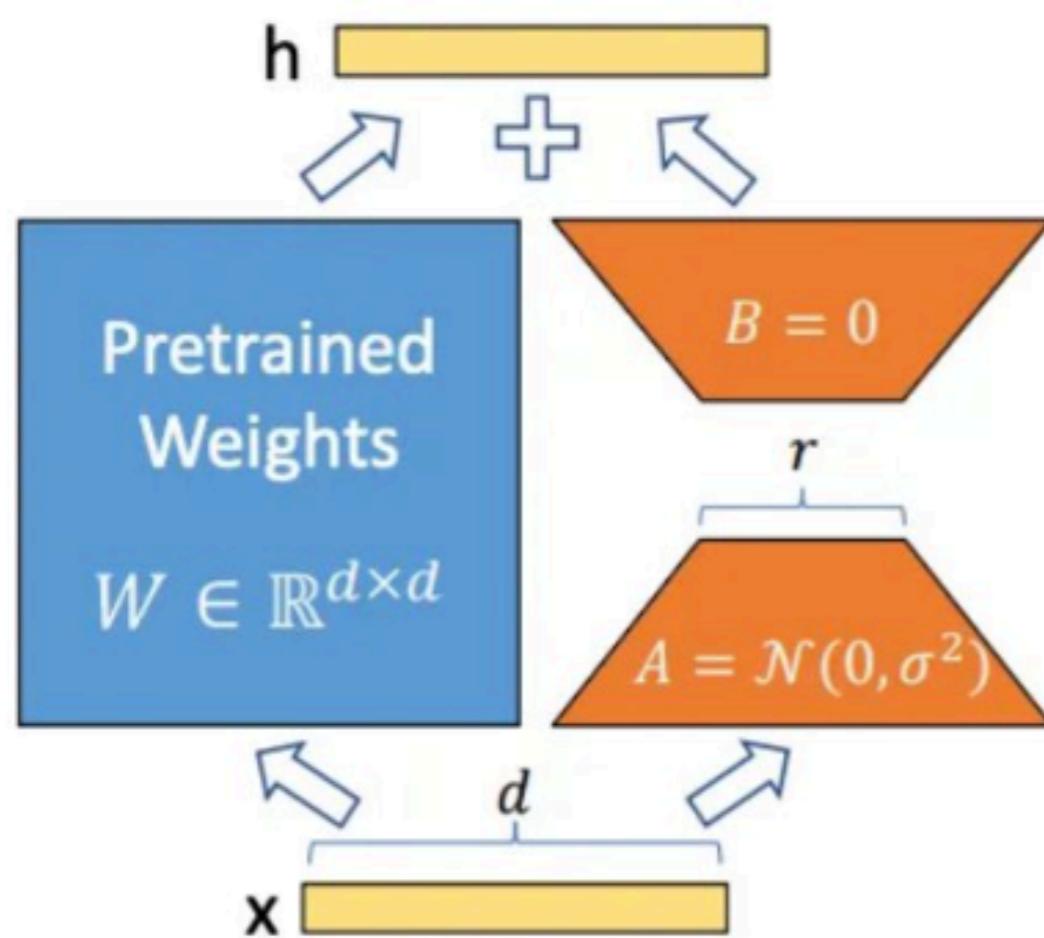
choices:

- rank ( $r=64|128$ )
- # epochs (3)
- learning rate ?



# Fine-Tuning w/LoRA : Low Rank Adaptation

LoRA leaves the pretrained layers of the LLM fixed and injects a trainable rank decomposition matrix in each layer of the transformer.



$$\begin{aligned} & \text{Finetuned Weights} & \text{Weight Update} \\ & \overbrace{W_{\text{ft}}}^{\text{Finetuned Weights}} = \underbrace{W_{\text{pt}}}_{\text{Pretrained Weights}} + \overbrace{\Delta W}^{\text{Weight Update}} \\ \\ & \text{Rank Decomposition Matrix} \\ & W_{\text{ft}} = W_{\text{pt}} + \Delta W = W_{\text{pt}} + \overbrace{AB}^{\text{Rank Decomposition Matrix}} \\ & \text{where } W_{\text{ft}}, W_{\text{pt}}, \Delta W, AB \in \mathbb{R}^{d \times d} \\ & \text{and } \underbrace{A \in \mathbb{R}^{d \times r}, B \in \mathbb{R}^{r \times d}}_{\text{Low Rank}} \end{aligned}$$

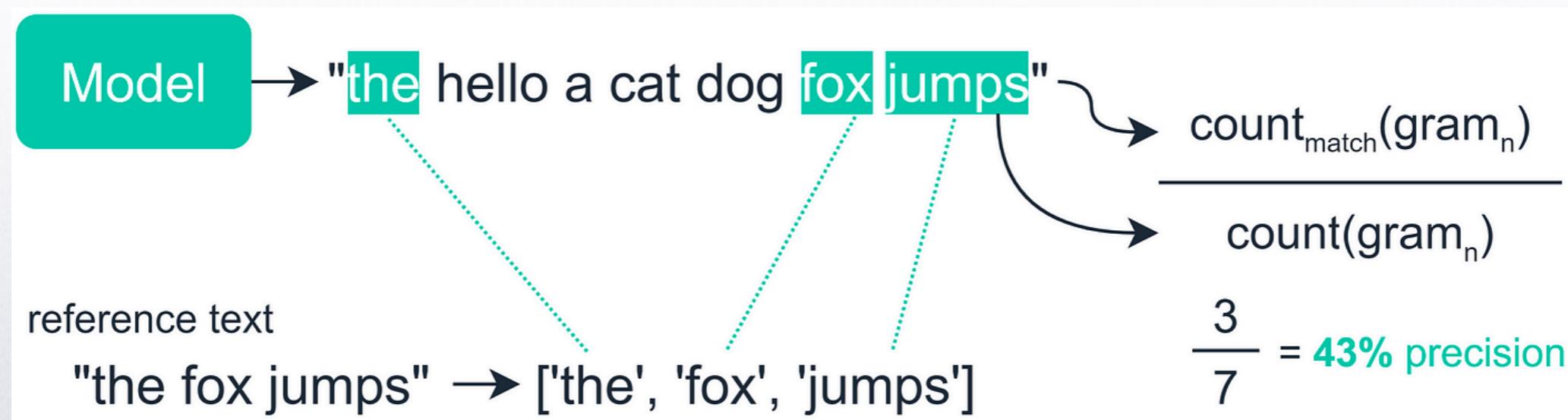


## Evaluations :Text Generation is Tricky

BLEU: ~ precision (# outputs in inputs)

ROUGE: ~ recall (# inputs in output)

i.e ROUGE-N: # ngrams from ref text in gen text



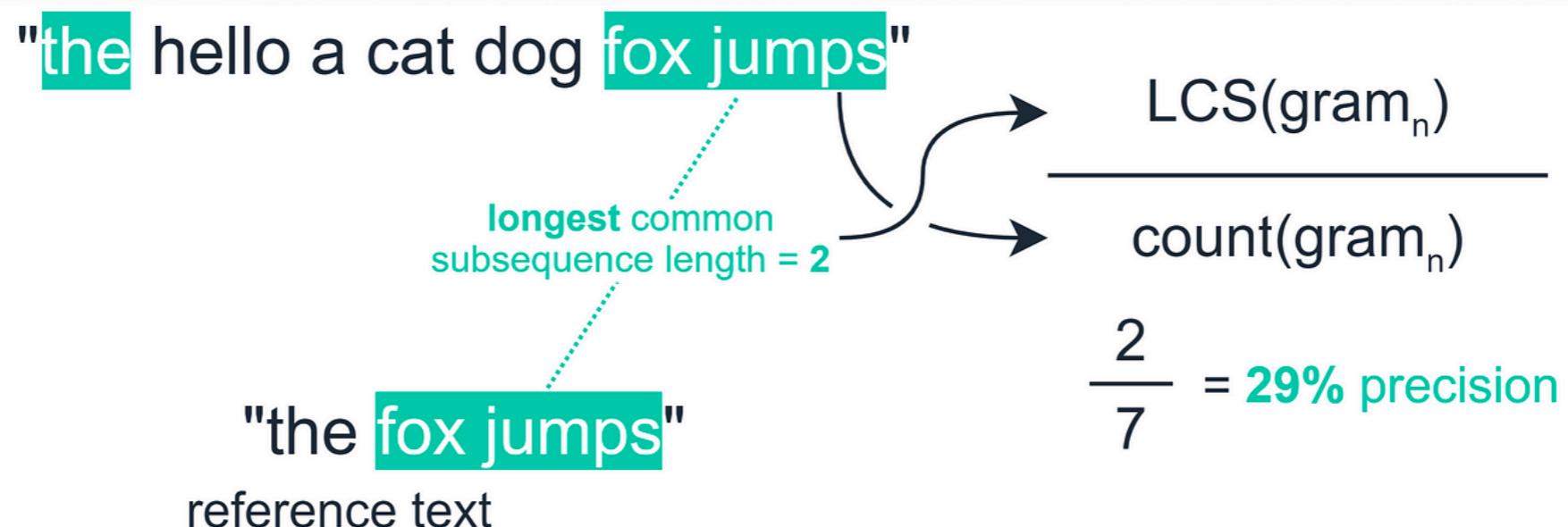


## Evaluations :Text Generation is Tricky

BLEU: ~ precision (# outputs in inputs)

ROUGE: ~ recall (# inputs in output)

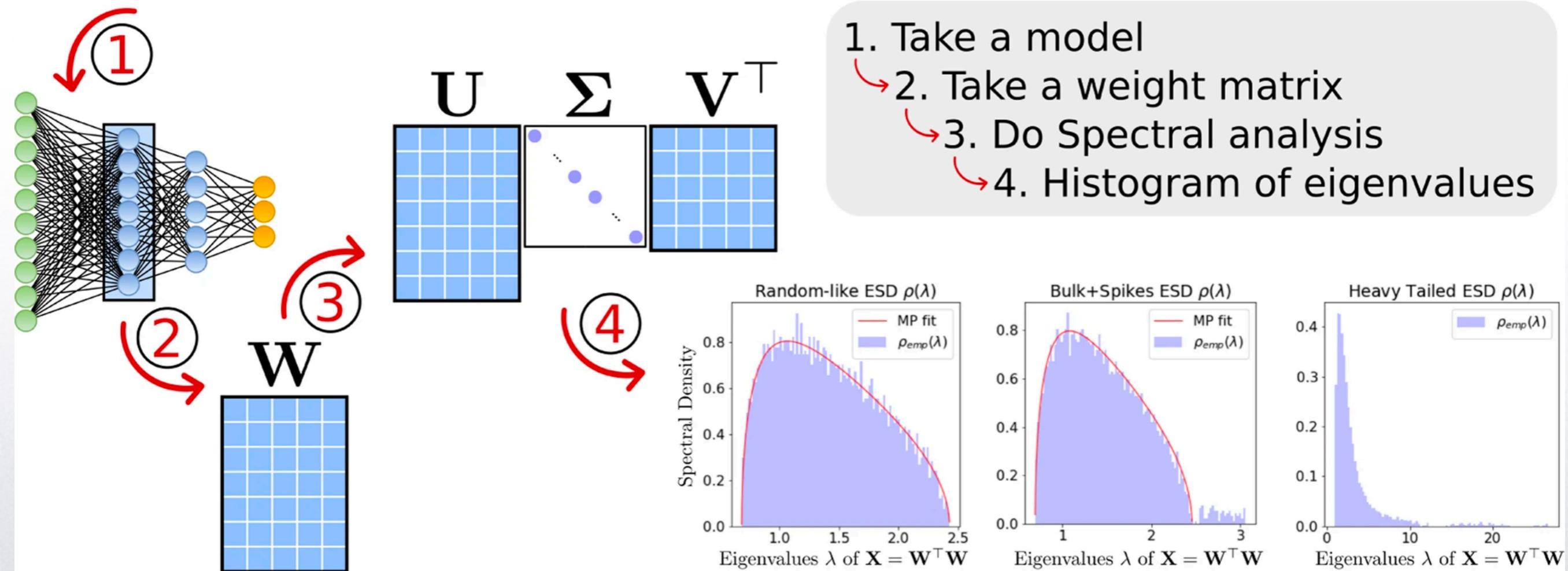
i.e ROUGE-L: # least common sequence





# WeightWatcher : Data-Free Quality Metrics

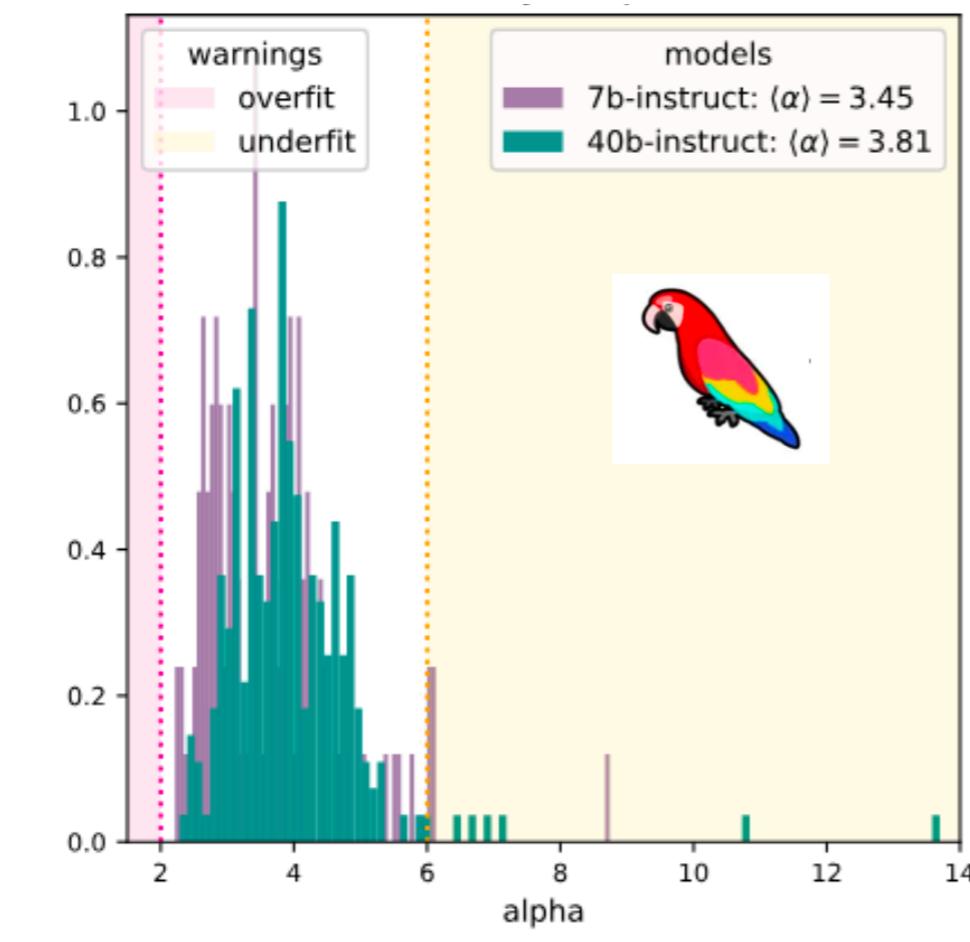
<https://weightwatcher.ai>



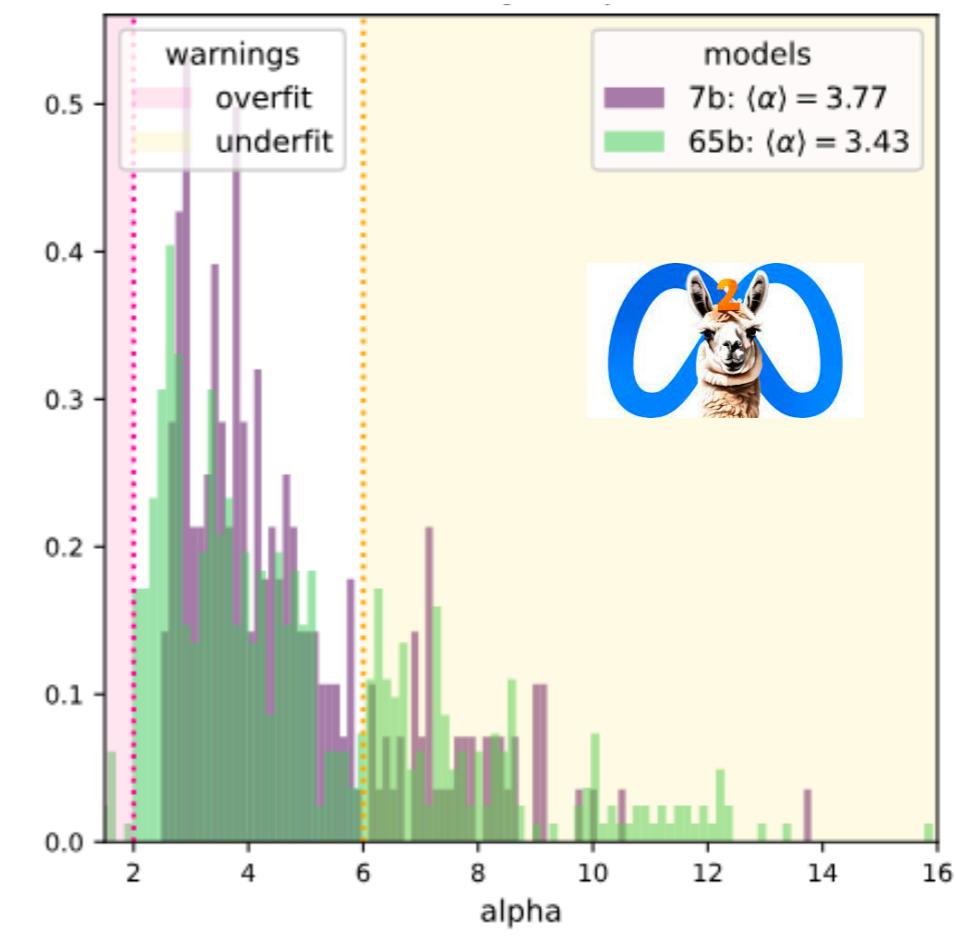


# WeightWatcher : Data-Free Quality Metrics

## Falcon



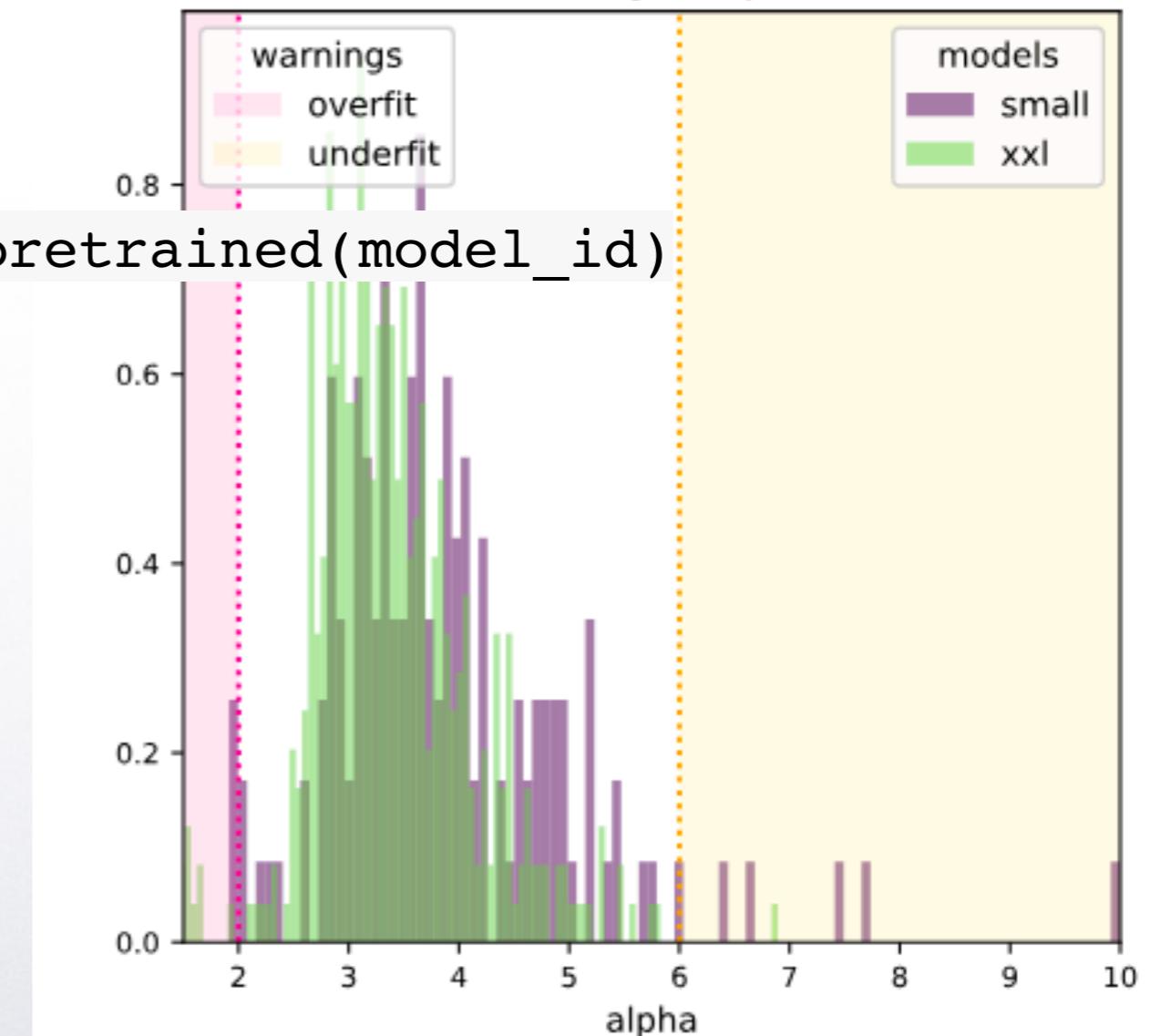
## Llama



<https://weightwatcher.ai>

# Fine-Tuning Examples : LoRA on FlanT5-XXL

```
from transformers import AutoTokenizer, AutoModelForSeq2SeqLM  
  
model_id="google/flan-t5-xxl"  
# Load tokenizer of FLAN-t5-XL  
tokenizer = AutoTokenizer.from_pretrained(model_id)
```





## Fine-Tuning Examples : SAMSum Dataset

```
from datasets import load_dataset  
dataset = load_dataset("samsum")
```

13818513	Amanda: I baked cookies. Do you want some? Jerry: Sure! Amanda: I'll...	Amanda baked cookies and will bring Jerry some...
13728867	Olivia: Who are you voting for in this election? Oliver: Liberals as...	Olivia and Olivier are voting for liberals in this...
13681000	Tim: Hi, what's up? Kim: Bad mood tbh, I was going to do lots of stuf...	Kim may try the pomodoro technique recommended by Ti...
13730747	Edward: Rachel, I think I'm in ove with Bella.. rachel: Dont say...	Edward thinks he is in love with Bella. Rachel wants...
13728094	Sam: hey overheard rick say something Sam: i don't know what to...	Sam is confused, because he overheard Rick complaining...
13716343	Neville: Hi there, does anyone remember what date I got married on...	Wyatt reminds Neville his wedding anniversary is on...

< Previous

1

2

3

...

148

Next >



## Fine-Tuning Examples : PEFT Package

```
from peft import LoraConfig, get_peft_model,  
                    prepare_model_for_int8_training, TaskType  
  
# Define LoRA Config  
lora_config = LoraConfig(  
    r=16,  
    lora_alpha=32,  
    target_modules=["q", "v"],  
    lora_dropout=0.05,  
    bias="none",  
    task_type=TaskType.SEQ_2_SEQ_LM  
)  
# prepare int-8 model for training  
model = prepare_model_for_int8_training(model)  
# add LoRA adaptor  
model = get_peft_model(model, lora_config)
```