

# INTRODUCTION TO ADVANCED LARGE LANGUAGE MODELS

A COMPREHENSIVE TUTORIAL ON TECHNIQUES,  
ARCHITECTURES, AND PRACTICAL APPLICATIONS

PART III

Giorgio Roffo, PhD

*Explore and Connect*

LinkedIn: Giorgio Roffo - [LinkedIn](#)

ResearchGate: [Work Done](#)

Google Scholar: [My Publications](#)

GitHub: [giorgioroffo](#)

*Special thanks to Shelbee Eigenbrode, Antje Barth, and Mike Chambers for their work on "Generative AI with Large Language Models" (Coursera, Amazon AWS, 2023), which significantly informed and inspired the material used in these slides.*



# WE ARE GOING TO TALK ABOUT...

## 1. **Beyond Basic LLMs:** Advances Toward Sophisticated Applications

- Retrieval-Augmented Generation (RAG)
- Program-Aided Language Models (PAL)
- Integration of Reasoning and Action (ReAct)
- Architectural Designs for LLM Applications (ChatGPT, Gemini, etc..)

*Part I*

## 2. **Introduction:** Transformers

## 3. **Selecting LLM Architectures**

## 4. **LLM Training Resources:** GPU Memory Requirements

- Models: BERT-L (340M), GPT-2 (1.5B), LLaMA-2 (7-13-70B), GPT-3 (175B), PaLM (540B)
- GPU RAM for Models:
  - 1B parameters: 24GB (32-bit precision)
  - 175B parameters: 4200GB (32-bit precision)
  - 500B parameters: 12000GB (32-bit precision)

*Part II*

## 5. **Datasets & Benchmarks:** Criteria for Selecting Evaluation Metrics and Datasets

## 6. **Fine-Tuning Strategies:** Addressing Specific Tasks and Preventing Overfitting

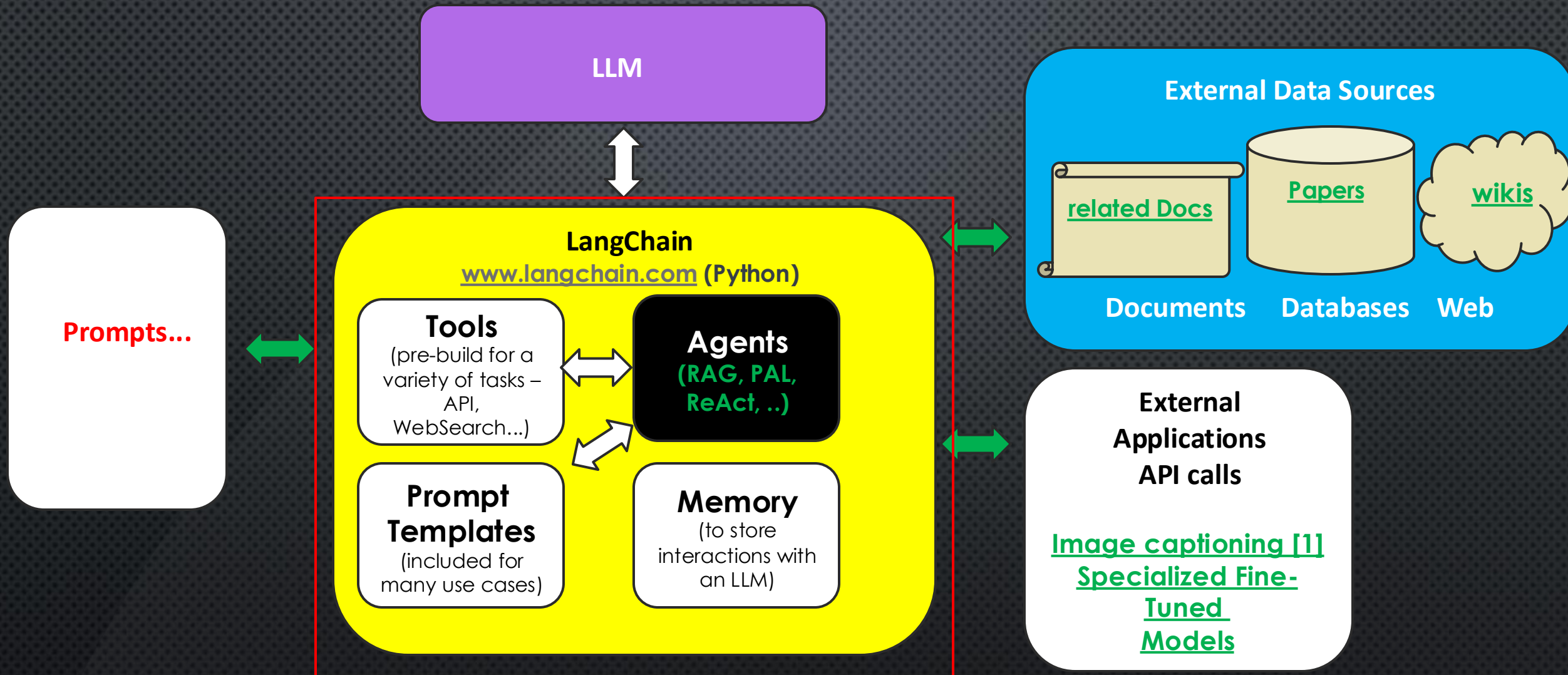
## 7. **RLHF:** Enhancing LLMs Through Human Interaction

## 8. **ReST (Google):** Reinforced Self-Training for Language Modeling

## 9. **MED-Gemini**

*Part III*

## Summary Part 1





# WE ARE GOING TO TALK ABOUT...

## 1. **Beyond Basic LLMs:** Advances Toward Sophisticated Applications

- Retrieval-Augmented Generation (RAG)
- Program-Aided Language Models (PAL)
- Integration of Reasoning and Action (ReAct)
- Architectural Designs for LLM Applications (ChatGPT, Gemini, etc..)

*Part I*

## 2. **Introduction: Transformers**

## 3. **Selecting LLM Architectures**

## 4. **LLM Training Resources:** GPU Memory Requirements

- Models: BERT-L (340M), GPT-2 (1.5B), LLaMA-2 (7-13-70B), GPT-3 (175B), PaLM (540B)
- GPU RAM for Models:
  - 1B parameters: 24GB (32-bit precision)
  - 175B parameters: 4200GB (32-bit precision)
  - 500B parameters: 12000GB (32-bit precision)

*Part II*

## 5. **Datasets & Benchmarks:** Criteria for Selecting Evaluation Metrics and Datasets

## 6. **Fine-Tuning Strategies:** Addressing Specific Tasks and Preventing Overfitting

## 7. **RLHF:** Enhancing LLMs Through Human Interaction

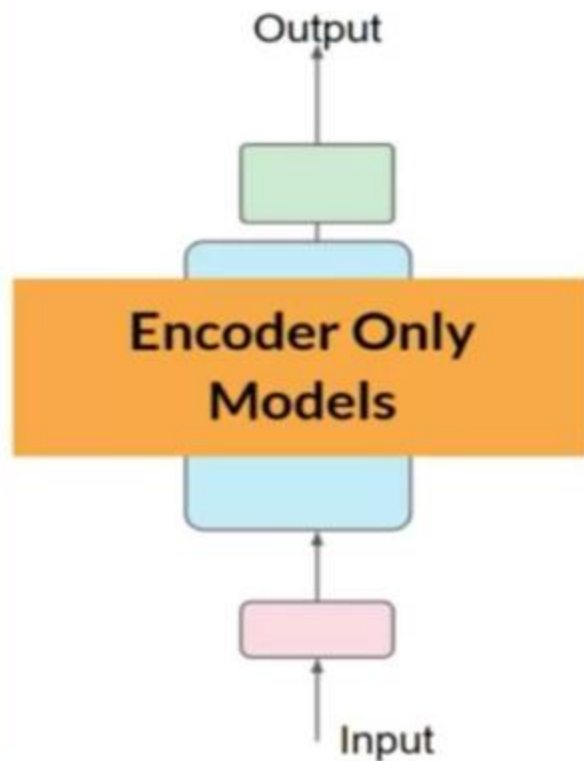
## 8. **ReST (Google):** Reinforced Self-Training for Language Modeling

## 9. **MED-Gemini**

*Part III*



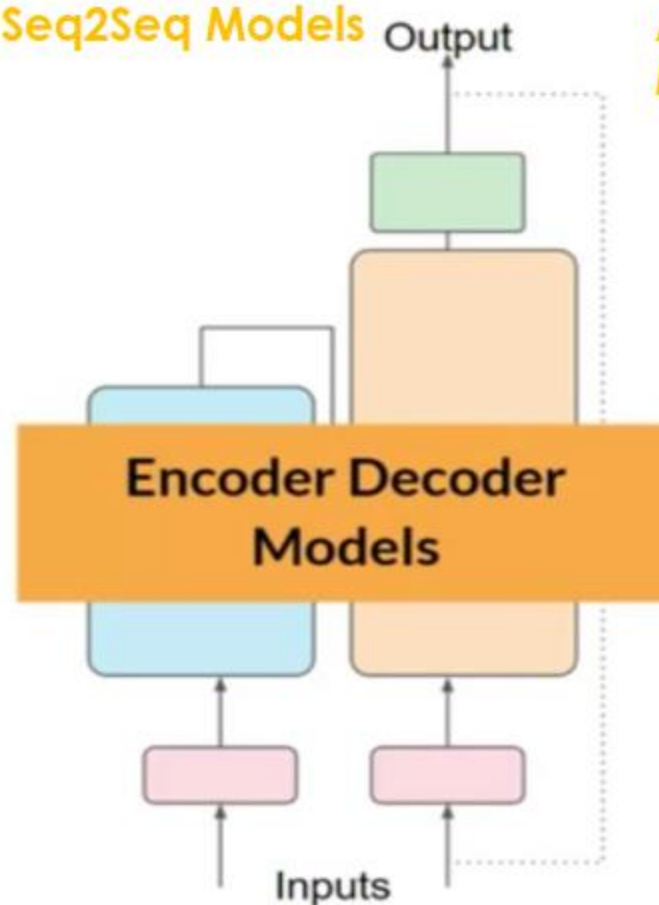
## Autoencoding Models



understanding input context

**Main Task:** Sentiment Analysis, word class.  
**Models:** BERT, ROBERTA

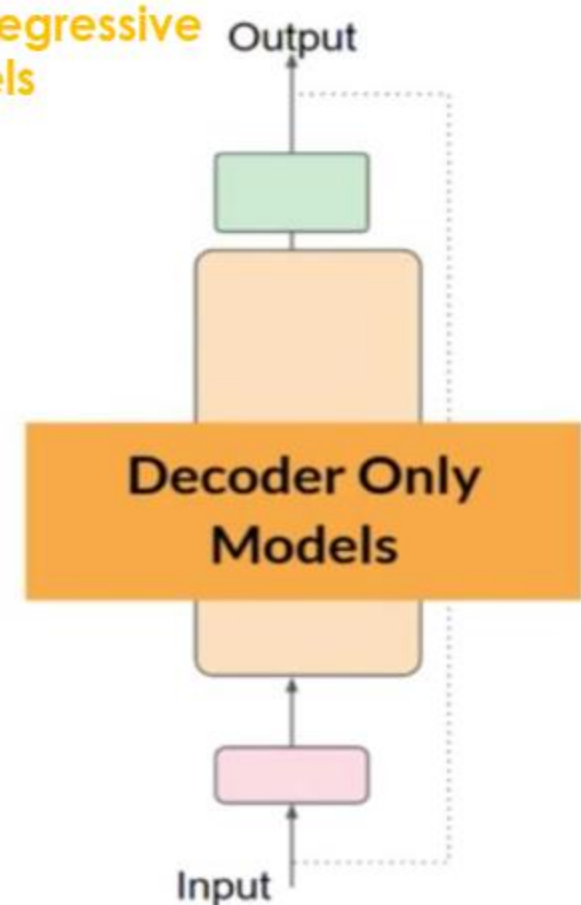
## Seq2Seq Models



Conversion

**Main Task:** Translation, Summarization, Q/A  
**Models:** BART, T5

## Autoregressive Models



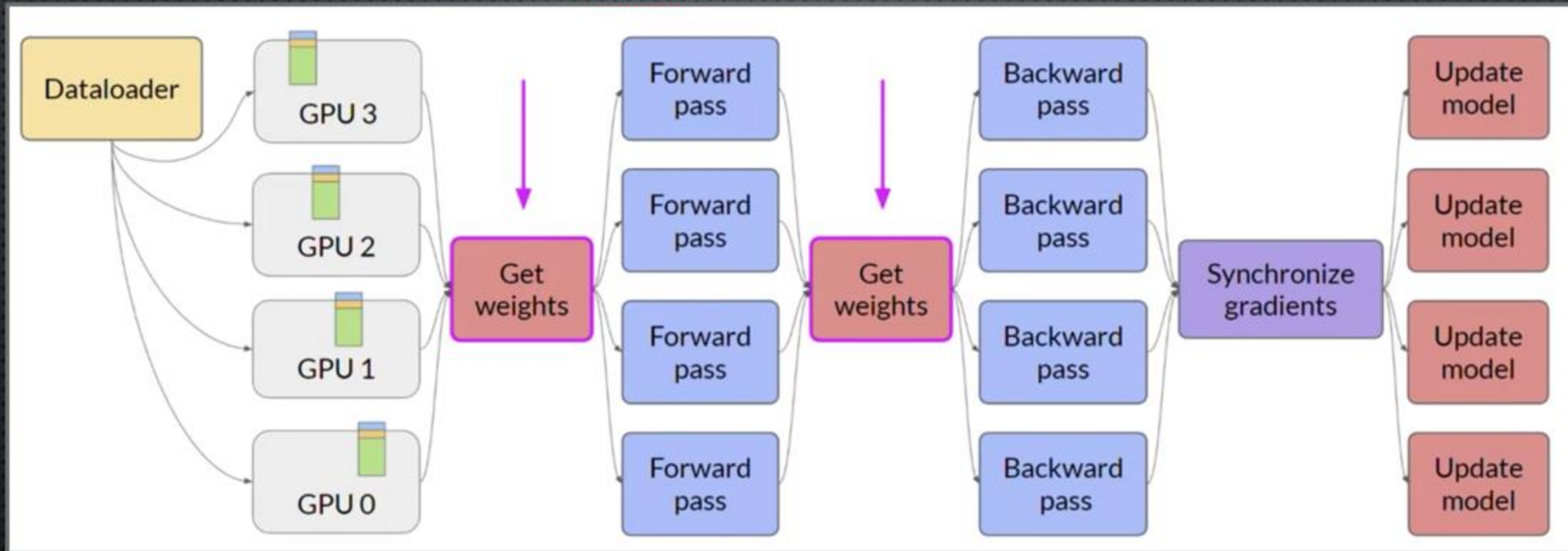
Generating text

**Main Task:** Generalize all tasks.  
**Models:** GPT, BLOOM, Jurassic, LLaMA



# Efficient multi-GPU compute strategies

1. Quantization. Reduces memory to store and train models
2. Distributed Data Parallel (DDP)
3. Fully Sharded Data Parallel (FSDP) -> PyTorch 2023



[9] Rajbhandari S, Rasley J, Ruwase O, He Y. **Zero: Memory optimizations toward training trillion parameter models**. In SC20: International Conference for High Performance Computing, Networking, Storage and Analysis 2020 Nov 9 (pp. 1-16). IEEE.

[10] Zhao Y, Gu A, Varma R, Luo L, Huang CC, Xu M, Wright L, Shojanazeri H, Ott M, Shleifer S, Desmaison A. **Pytorch FSDP: experiences on scaling fully sharded data parallel**. arXiv preprint arXiv:2304.11277. 2023 Apr 21.



|          |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                      |
|----------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Question | A 27-year-old male presents to urgent care complaining of pain with urination. He reports that the pain started 3 days ago. He has never experienced these symptoms before. He <i>denies gross hematuria or pelvic pain</i> . He is sexually active with his girlfriend, and they consistently use condoms. When asked about recent travel, he admits to recently returning from a boys' trip in Cancun where he had <i>unprotected sex</i> 1 night with a girl he met at a bar. The patients medical history includes type I diabetes that is controlled with an insulin pump. His mother has rheumatoid arthritis. The patients temperature is 99 F (37.2 C), blood pressure is 112/74 mmHg, and pulse is 81/min. On physical examination, there are no lesions of the penis or other body rashes. No costovertebral tenderness is appreciated. A urinalysis reveals no blood, glucose, ketones, or proteins but is <i>positive for leukocyte esterase</i> . A urine microscopic evaluation shows a <i>moderate number of white blood cells</i> but no casts or crystals. A urine culture is negative. Which of the following is the most likely cause for the patient's symptoms? |
| Options  | A: <b>Chlamydia trachomatis</b> , B: Systemic lupus erythematosus, C: Mycobacterium tuberculosis, D: Treponema pallidum                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              |
| Evidence | At least one-third of male patients with <i>C. trachomatis</i> urethral infection have <i>no evident signs or symptoms of urethritis</i> . ... Such patients generally have <i>pyuria</i> ..., a <i>positive leukocyte esterase test</i> , ...                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                       |
| Question | A 57-year-old man presents to his primary care physician with a 2-month history of <i>right upper and lower extremity weakness</i> . He noticed the weakness when he started falling far more frequently while running errands. Since then, he has had <i>increasing difficulty</i> with walking and lifting objects. His past medical history is significant only for well-controlled hypertension, but he says that some members of his family have had <i>musculoskeletal problems</i> . His right upper extremity shows <i>forearm atrophy</i> and <i>depressed reflexes</i> while his right lower extremity is <i>hypertonic with a positive Babinski sign</i> . Which of the following is most likely associated with the cause of this patients symptoms?                                                                                                                                                                                                                                                                                                                                                                                                                     |
| Options  | A: HLA-B8 haplotype, B: HLA-DR2 haplotype, C: <b>Mutation in SOD1</b> , D: Mutation in SMN1, E: Viral infection                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                      |
| Evidence | 1. The manifestations of ALS ... <i>insidiously developing asymmetric weakness</i> , usually first evident distally in one of the limbs.<br>2. ... <i>hyperactivity of the muscle-stretch reflexes (tendon jerks)</i> and, often, <i>spastic resistance to passive movements</i> ...<br>3. <i>Familial ALS (FALS)</i> ... clinically indistinguishable from sporadic ALS... Genetic studies have identified mutations in multiple genes, including cytosolic enzyme <i>SOD1</i> ...                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                  |

Table 1: Two examples of MEDQA. The correct answer among options is marked in bold font. Key words in the question and evidence text to help answer the questions are highlighted in italic font. Evidence for both examples are from the textbook "Harrison's Principles of Internal Medicine".

## Capabilities of Gemini Models in Medicine <https://arxiv.org/html/2404.18416v1>

- MedQA-RS (Reasoning and Search)
- MedQA-R (Reasoning)

MedQA-R (Reasoning), which extends MedQA with synthetically generated reasoning explanations, or "Chain-of-Thoughts" (CoTs), and MedQA-RS (Reasoning and Search), which extends MedQA-R with instructions to use web search results as additional context to improve answer accuracy.



## ROUGE (Recall-Oriented Understudy for Gisting Evaluation):

$$\begin{aligned}\text{ROUGE-L Recall:} &= \frac{\text{LCS}(\text{Gen}, \text{Ref})}{\text{unigrams in reference}} \\ \text{ROUGE-L Precision:} &= \frac{\text{LCS}(\text{Gen}, \text{Ref})}{\text{unigrams in output}} \\ \text{ROUGE-L F1:} &= 2 \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}\end{aligned}$$

Example ROUGE-L

$$\begin{aligned}\text{ROUGE-2 Recall:} &= \frac{\text{bigram matches}}{\text{bigrams in reference}} \\ \text{ROUGE-2 Precision:} &= \frac{\text{bigram matches}}{\text{bigrams in output}} \\ \text{ROUGE-2 F1:} &= 2 \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}\end{aligned}$$

Example ROUGE-2

## BLEU Score (Bilingual Evaluation Understudy):

BLEU is a metric for evaluating a machine-translated text by comparing it to one or more reference translations.

The core idea is to calculate the precision of n-grams (contiguous sequence of n items from a given sample of text or speech) in the candidate translation that also appear in the reference translation.



# WE ARE GOING TO TALK ABOUT...

## 1. **Beyond Basic LLMs:** Advances Toward Sophisticated Applications

- Retrieval-Augmented Generation (RAG)
- Program-Aided Language Models (PAL)
- Integration of Reasoning and Action (ReAct)
- Architectural Designs for LLM Applications (ChatGPT, Gemini, etc..)

## 2. **Introduction:** Transformers

## 3. **Selecting LLM Architectures**

## 4. **LLM Training Resources:** GPU Memory Requirements

- Models: BERT-L (340M), GPT-2 (1.5B), LLaMA-2 (7-13-70B), GPT-3 (175B), PaLM (540B)
- GPU RAM for Models:
  - 1B parameters: 24GB (32-bit precision)
  - 175B parameters: 4200GB (32-bit precision)
  - 500B parameters: 12000GB (32-bit precision)

## 5. **Datasets & Benchmarks:** Criteria for Selecting Evaluation Metrics and Datasets

## 6. **Fine-Tuning Strategies:** Addressing Specific Tasks and Preventing Overfitting

## 7. **RLHF:** Enhancing LLMs Through Human Interaction

## 8. **ReST (Google):** Reinforced Self-Training for Language Modeling

## 9. **MED-Gemini**



### Domain adaptation

#### BloombergGPT: A Large Language Model for Finance

Shijie Wu<sup>1,\*</sup>, Ozan Irsoy<sup>1,\*</sup>, Steven Lu<sup>1,\*</sup>, Vadim Dabravolski<sup>1</sup>, Mark Dredze<sup>1,2</sup>, Sebastian Gehrmann<sup>1</sup>, Prabhanjan Kambadur<sup>1</sup>, David Rosenberg<sup>1</sup>, Gideon Mann<sup>1</sup>

<sup>1</sup> Bloomberg, New York, NY USA

<sup>2</sup> Computer Science, Johns Hopkins University, Baltimore, MD USA

gmann16@bloomberg.net

#### Abstract

The use of NLP in the realm of financial technology is broad and complex, with applications ranging from sentiment analysis and named entity recognition to question answering. Large Language Models (LLMs) have been shown to be effective on a variety of tasks; however, no LLM specialized for the financial domain has been reported in literature. In this work, we present BLOOMBERGGPT, a 50 billion parameter language model that is trained on a wide range of financial data. We construct a 363 billion token dataset based on Bloomberg's extensive data sources, perhaps the largest domain-specific dataset yet, augmented with 345 billion tokens from general purpose datasets. We validate BLOOMBERGGPT on standard LLM benchmarks, open financial benchmarks, and a suite of internal benchmarks that most accurately reflect our intended usage. Our mixed dataset training leads to a model that outperforms existing models on financial tasks by significant margins without sacrificing performance on general LLM benchmarks. Additionally, we explain our modeling choices, training process, and evaluation methodology. As a next step, we plan to release training logs (Chronicles) detailing our experience in training BLOOMBERGGPT.

~51%

Financial  
(Public & Private)

~49%

Other  
(Public)

OR:

#### Medical language

After a strenuous workout, the patient experienced severe myalgia that lasted for several days.

After the biopsy, the doctor confirmed that the tumor was malignant and recommended immediate treatment.

Sig: 1 tab po qid pc & hs

- Uncommon words to describe medical conditions and procedures.
- Infrequent terms in common training sets



## Section 6: Fine-Tuning Strategies

### LLM fine-tuning Strategies:

- **Instruction Fine-Tuning (full fine-tuning updates model parameters)**
  - Each prompt/completion pair includes a specific "instruction" to the LLM.
    - *"Classify this review:" [TEXT] "Sentiment:"*
    - *"Summarize the following text:" [TEXT]*
    - *"Translate this sentence to: .." [TEXT]*
  - Public datasets are not formatted, there are many prompt-template libraries that can format the data in **prompt instruction format** for different tasks:
    - Classification / sentiment analysis
      - *"Given the following review: [TEXT] predict the associated rating [ANSWER\_CHOICES]"*
    - Text Generation
      - *"Generate a [STAR\_RATING]-star review (1 being lowest and 5 being highest) about this product [PRODUCT\_TITLE]. || [BODY]."*
    - Text Summarization
      - *"Give a short sentence describing the following product review: [BODY]\ || [REVIEW\_HEADLINE]."*



```
"samsun": [  
    ("{"dialogue"}\n\nBriefly summarize that dialogue.", "{"summary}"),  
    ("Here is a dialogue:\n{"dialogue"}\n\nWrite a short summary!",  
     "{"summary}"),  
    ("Dialogue:\n{"dialogue"}\n\nWhat is a summary of this dialogue?",  
     "{"summary}"),  
    ("{"dialogue"}\n\nWhat was that dialogue about, in two sentences or less?",  
     "{"summary}"),  
    ("Here is a dialogue:\n{"dialogue"}\n\nWhat were they talking about?",  
     "{"summary}"),  
    ("Dialogue:\n{"dialogue"}\n\nWhat were the main points in that "  
     "conversation?", "{"summary}"),  
    ("Dialogue:\n{"dialogue"}\n\nWhat was going on in that conversation?",  
     "{"summary}"),  
]
```



## Section 6: Fine-Tuning Strategies

### LLM fine-tuning Strategies:

- Instruction Fine-Tuning (full fine-tuning updates model parameters)
- **Fine-tuning on a single task**
  - Improved performance on the new task.
  - Downside: **Catastrophic Forgetting.**
    - Catastrophic forgetting happens because the full fine-tuning process modifies the weights of the original LLM
    - The model forget in doing the other tasks even important and related to the new task.
      - *For example, improve the ability of a model to perform sentiment analysis on a review forget how to do named entity recognition.*



## Section 6: Fine-Tuning Strategies

### LLM fine-tuning Strategies:

- Instruction Fine-Tuning (full fine-tuning updates model parameters)
- Fine-tuning on a single task
- **Fine-Tuning on multiple tasks**
  - The dataset contains examples that instruct the model to carry out a variety of tasks.
  - Reduces the catastrophic forgetting
  - Drawback: it requires a lot of data (50K-100K examples)
  - Examples are **FLAN** models (**F**ine-tuned **L**anguage **N**et)
    - FLAN-T5
    - FLAN-PaLM
    - ...

**FLAN-T5**  
it's been fine tuned on 473  
datasets across 146 task  
categories

| T0-SF                                                                                                                                                                                           | Muffin                                                                                                                                                                                                                | CoT (reasoning)                                                                                                                                                                                                     | Natural Instructions                                                                                                                                                                                                             |
|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <ul style="list-style-type: none"><li>- Commonsense Reasoning,</li><li>- Question Generation,</li><li>- Closed-book QA,</li><li>- Adversarial QA,</li><li>- Extractive QA</li><li>...</li></ul> | <ul style="list-style-type: none"><li>- Natural language inference,</li><li>- Code instruction gen,</li><li>- Code repair</li><li>- Dialog context generation,</li><li>- Summarization (SAMSum)</li><li>...</li></ul> | <ul style="list-style-type: none"><li>- Arithmetic reasoning,</li><li>- Commonsense reasoning</li><li>- Explanation generation,</li><li>- Sentence composition,</li><li>- Implicit reasoning,</li><li>...</li></ul> | <ul style="list-style-type: none"><li>- Cause effect classification,</li><li>- Commonsense reasoning,</li><li>- Named Entity Recognition,</li><li>- Toxic Language Detection,</li><li>- Question answering</li><li>...</li></ul> |
| 55 Datasets<br>14 Categories<br>193 Tasks                                                                                                                                                       | 69 Datasets<br>27 Categories<br>80 Tasks                                                                                                                                                                              | 9 Datasets<br>1 Category<br>9 Tasks                                                                                                                                                                                 | 372 Datasets<br>108 Categories<br>1554 Tasks                                                                                                                                                                                     |

[2] Chung HW, Hou L, Longpre S, Zoph B, Tay Y, Fedus W, Li Y, Wang X, Dehghani M, Brahma S, Webson A. **Scaling instruction-finetuned language models**. *Journal of Machine Learning Research*. 2024;25(70):1-53.

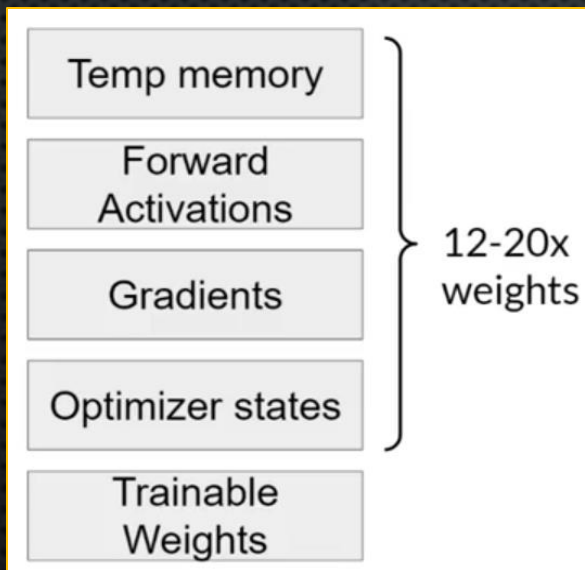


## Section 6: Fine-Tuning Strategies

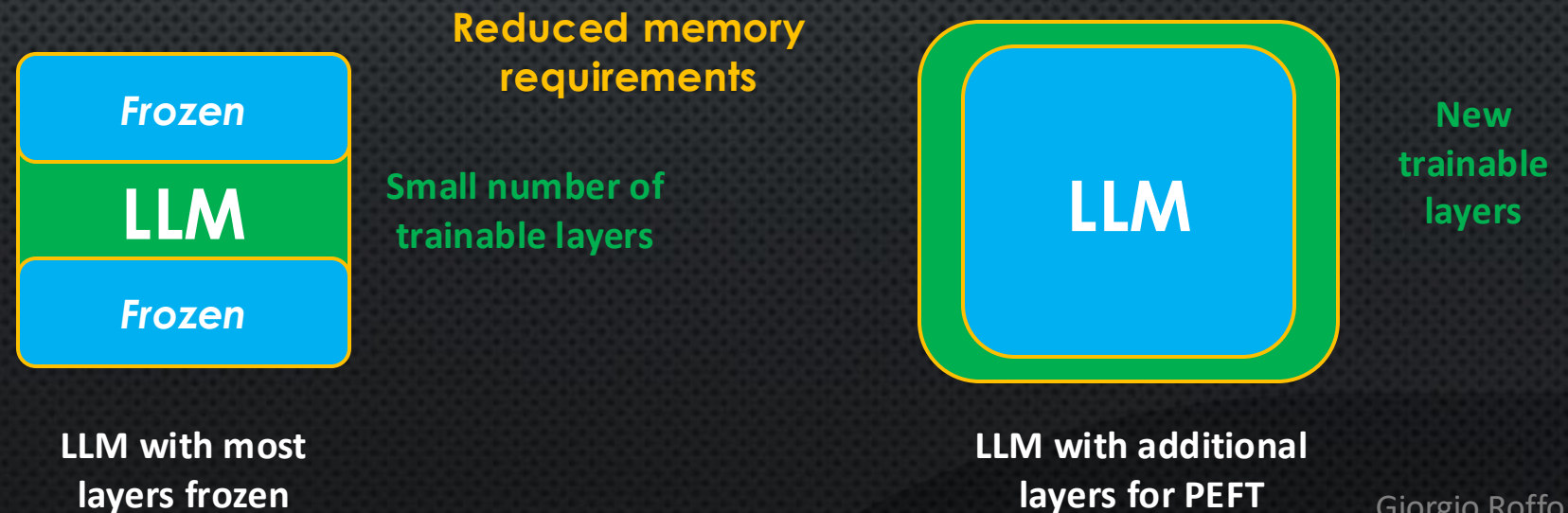
### LLM fine-tuning Strategies:

- Instruction Fine-Tuning (full fine-tuning updates model parameters)
- Fine-tuning on a single task
- Fine-Tuning on multiple tasks
- **Parameter Efficient Fine Tuning (PEFT)**
  - **LoRA (Low Rank Adaption)**
  - **Adapters**
  - **Soft prompts**

#### Full fine-tuning



#### Parameter Efficient Fine Tuning (PEFT)

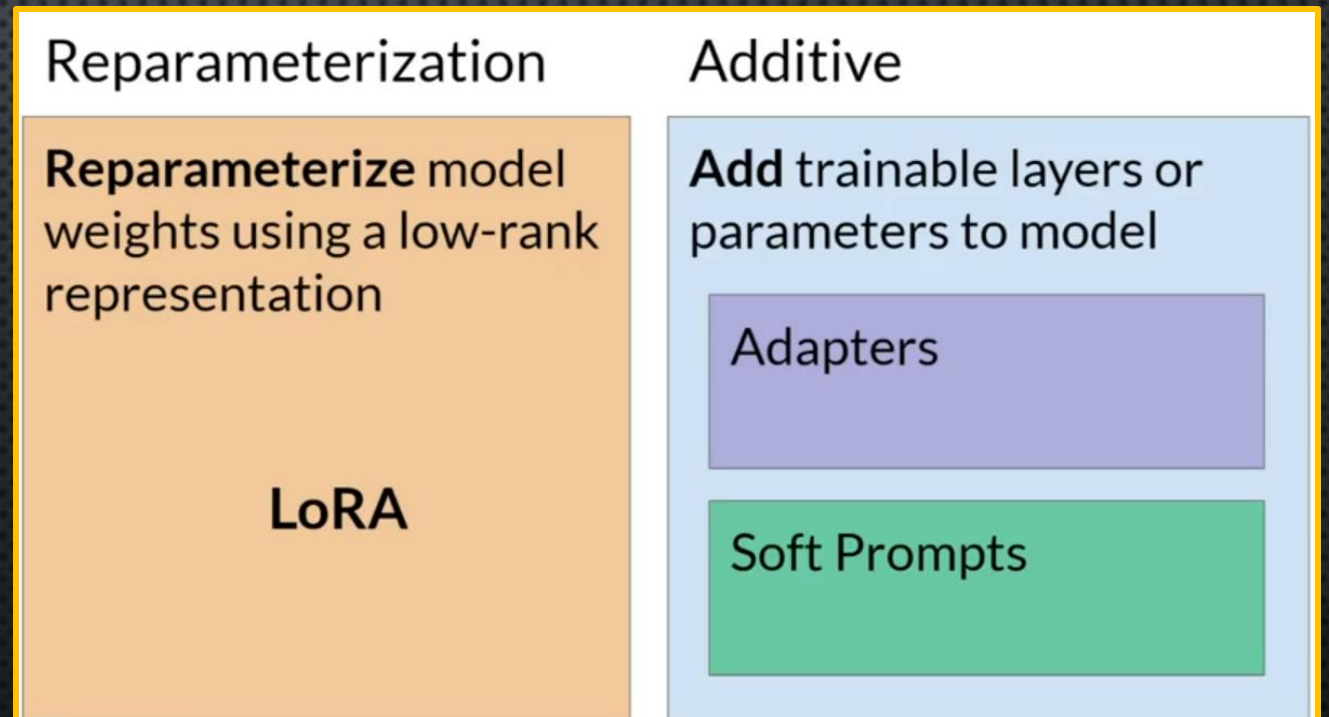




## Section 6: Fine-Tuning Strategies

### LLM fine-tuning Strategies:

- Instruction Fine-Tuning (full fine-tuning updates model parameters)
- Fine-tuning on a single task
- Fine-Tuning on multiple tasks
- **Parameter Efficient Fine Tuning (PEFT)**
  - **LoRA (Low Rank Adaption) - Reparameterization Technique**
  - **Adapters**
  - **Soft prompts**



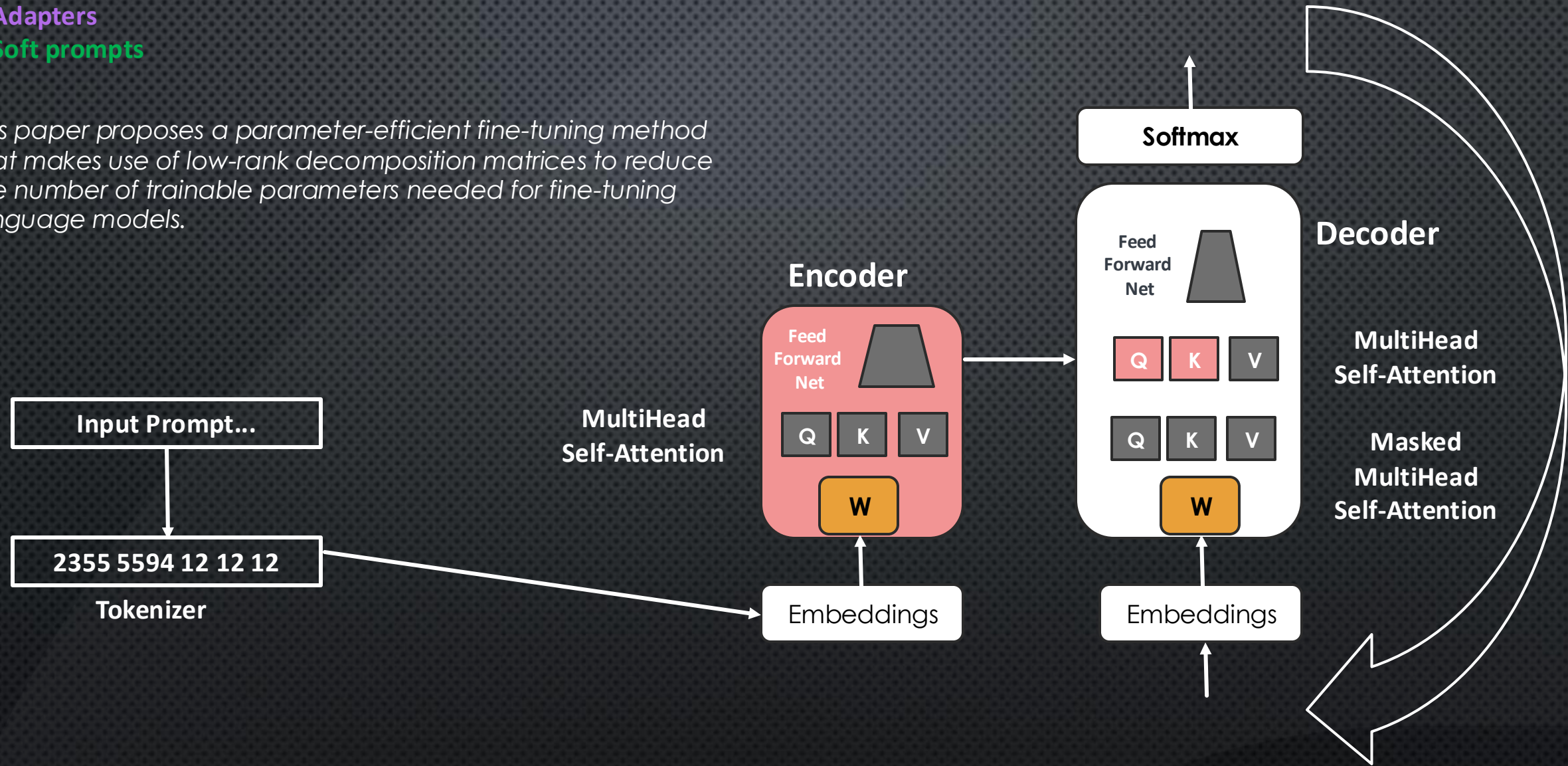
[11] Lialin V, Deshpande V, Rumshisky A. **Scaling down to scale up: A guide to parameter-efficient fine-tuning**. arXiv preprint arXiv:2303.15647. 2023 Mar 28.



## Section 6: Parameter Efficient Fine Tuning (PEFT)

- **LoRA (Low Rank Adaption) - Reparameterization Technique**
- **Adapters**
- **Soft prompts**

*This paper proposes a parameter-efficient fine-tuning method that makes use of low-rank decomposition matrices to reduce the number of trainable parameters needed for fine-tuning language models.*

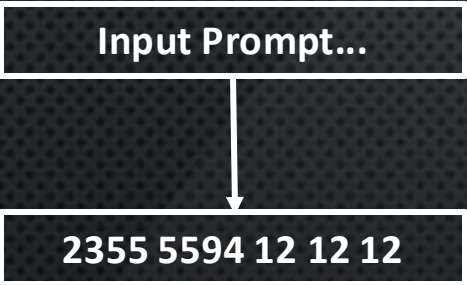
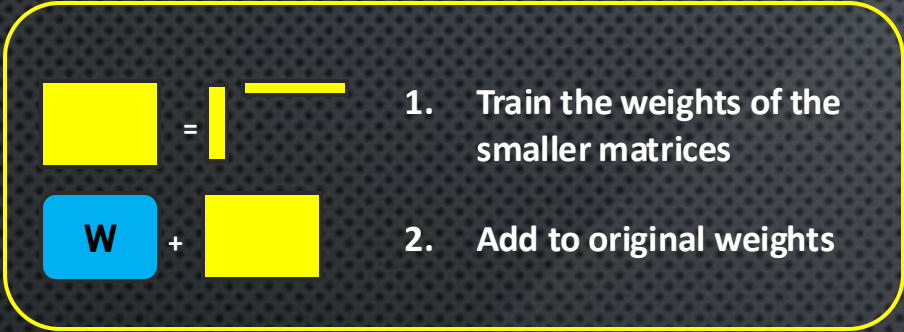


[11] Lialin V, Deshpande V, Rumshisky A. **Scaling down to scale up: A guide to parameter-efficient fine-tuning**. arXiv preprint arXiv:2303.15647. 2023.



Section 6: Parameter Efficient Fine Tuning (PEFT)

- **LoRA (Low Rank Adaption) - Reparameterization Technique**
- **Adapters**
- **Soft prompts**



Tokenizer

MultiHead Self-Attention

Inject 2 rank decomposition matrices

Encoder



Embeddings

Softmax

Decoder

Feed Forward Net



**W**

MultiHead Self-Attention

Masked MultiHead Self-Attention

Embeddings

[12] Hu EJ, Wallis P, Allen-Zhu Z, Li Y, Wang S, Wang L, Chen W. *LoRA: Low-Rank Adaptation of Large Language Models*. In International Conference on Learning Representations 2021 Oct 6. (Microsoft)



## Section 6: Parameter Efficient Fine Tuning (PEFT)

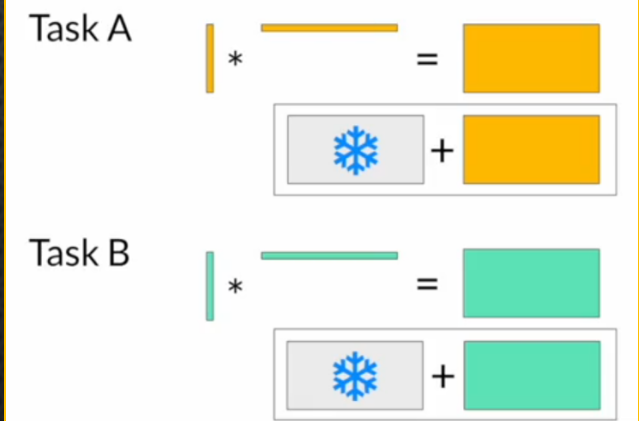
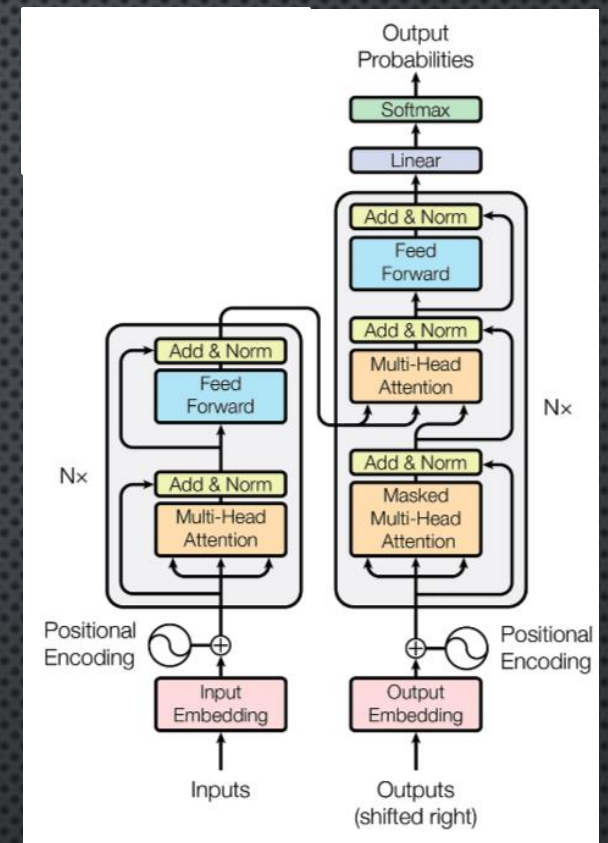
- **LoRA (Low Rank Adaption) - Reparameterization Technique**
- **Adapters**
- **Soft prompts**

- **Transformer in Attention is all you need:**
  - Weights dimensions  $d \times k = 512 \times 64 = 32768$  params
- **Lora with Rank  $r = 8$** 
  - Matrix A  $r \times k = 8 \times 64 = 512$  params
  - Matrix B  $d \times r = 512 \times 8 = 4096$  params
  - **86% reduction in parameters to train.**

| Rank $r$ | val_loss    | BLEU         | NIST          | METEOR        | ROUGE_L       | CIDEr         |
|----------|-------------|--------------|---------------|---------------|---------------|---------------|
| 1        | 1.23        | 68.72        | 8.7215        | 0.4565        | 0.7052        | 2.4329        |
| 2        | 1.21        | 69.17        | 8.7413        | 0.4590        | 0.7052        | 2.4639        |
| 4        | 1.18        | <b>70.38</b> | <b>8.8439</b> | <b>0.4689</b> | 0.7186        | <b>2.5349</b> |
| 8        | 1.17        | 69.57        | 8.7457        | 0.4636        | <b>0.7196</b> | 2.5196        |
| 16       | <b>1.16</b> | 69.61        | 8.7483        | 0.4629        | 0.7177        | 2.4985        |
| 32       | <b>1.16</b> | 69.33        | 8.7736        | 0.4642        | 0.7105        | 2.5255        |
| 64       | <b>1.16</b> | 69.24        | 8.7174        | 0.4651        | 0.7180        | 2.5070        |
| 128      | <b>1.16</b> | 68.73        | 8.6718        | 0.4628        | 0.7127        | 2.5030        |
| 256      | <b>1.16</b> | 68.92        | 8.6982        | 0.4629        | 0.7128        | 2.5012        |
| 512      | <b>1.16</b> | 68.78        | 8.6857        | 0.4637        | 0.7128        | 2.5025        |
| 1024     | 1.17        | 69.37        | 8.7495        | 0.4659        | 0.7149        | 2.5090        |

**Very parameter  
efficient  
strategy**

**Task adaptation:  
Switch out small matrices at  
inference time to change the  
task**

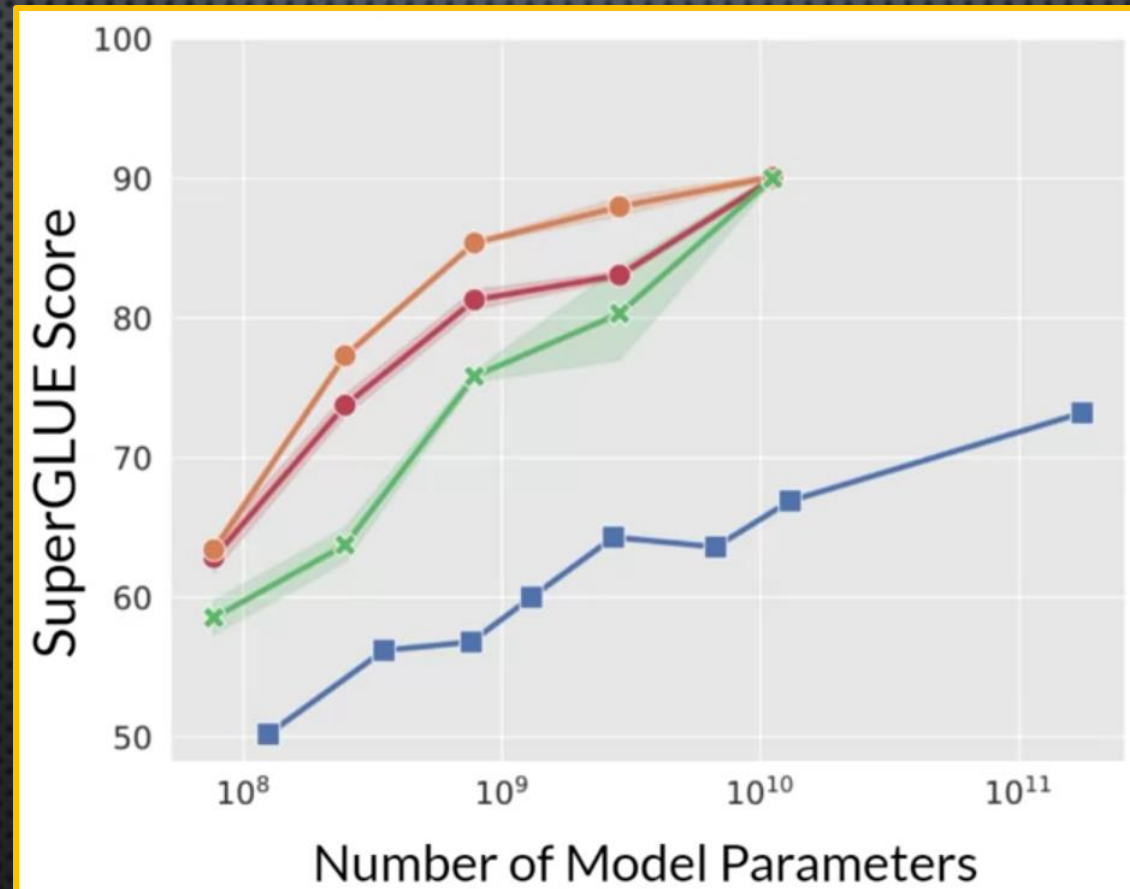




## Section 6: Parameter Efficient Fine Tuning (PEFT)

- **LoRA (Low Rank Adaption)**
- **Adapters (change the architecture, adding more trainable layers)**
- **Soft prompts (keep the model architecture fixed and add trainable parameters to the embedding vectors)**

The paper explores "prompt tuning," a method for conditioning language models with learned soft prompts, achieving competitive performance compared to full fine-tuning and enabling model reuse for many tasks

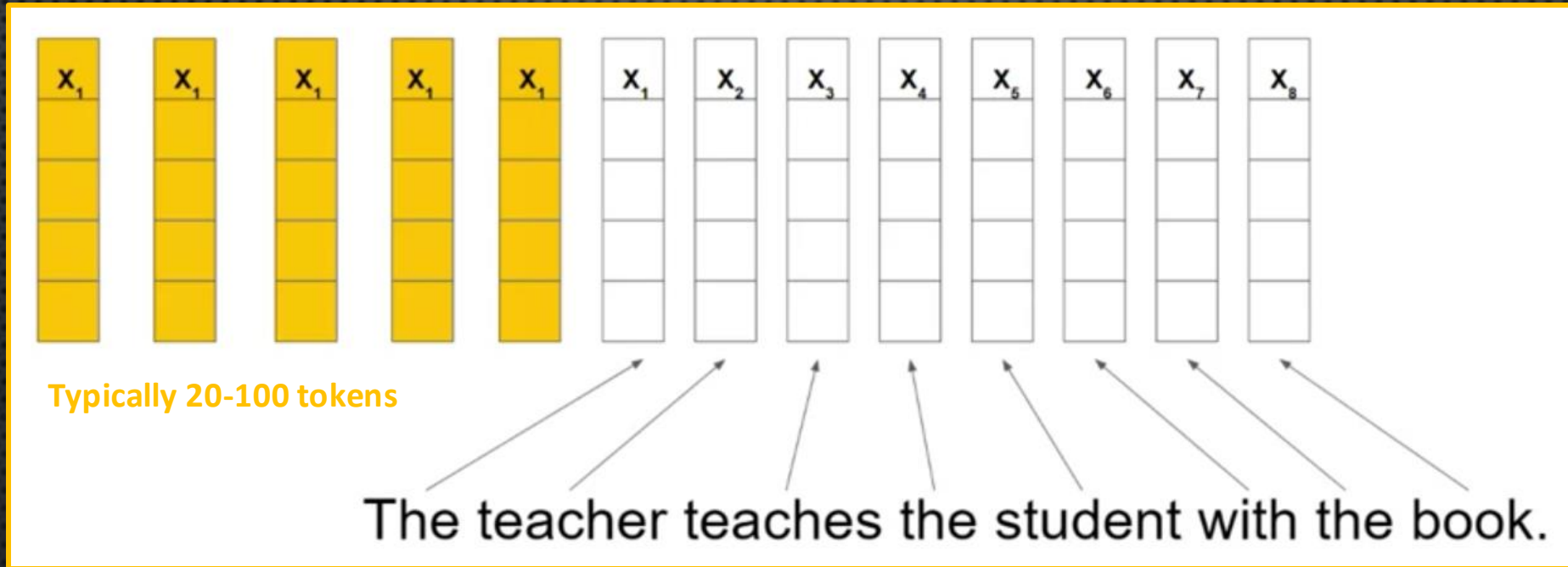


**Full fine-tuning**  
**Multi-task Fine-tuning**  
**Sof-Prompt tuning**  
**Prompt engineering**



## Section 6: Parameter Efficient Fine Tuning (PEFT)

- **LoRA (Low Rank Adaption)**
- **Adapters (change the architecture, adding more trainable layers)**
- **Soft prompts (architecture is not changed)**

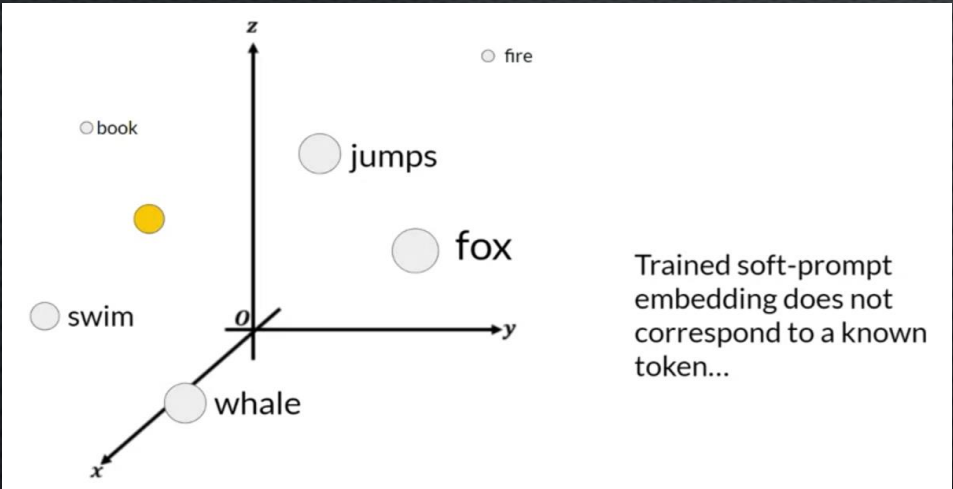
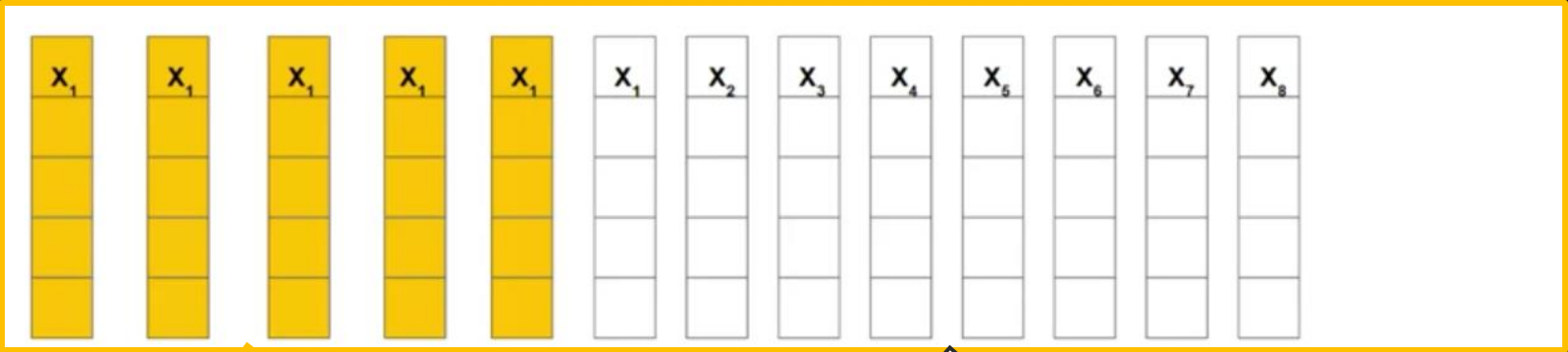




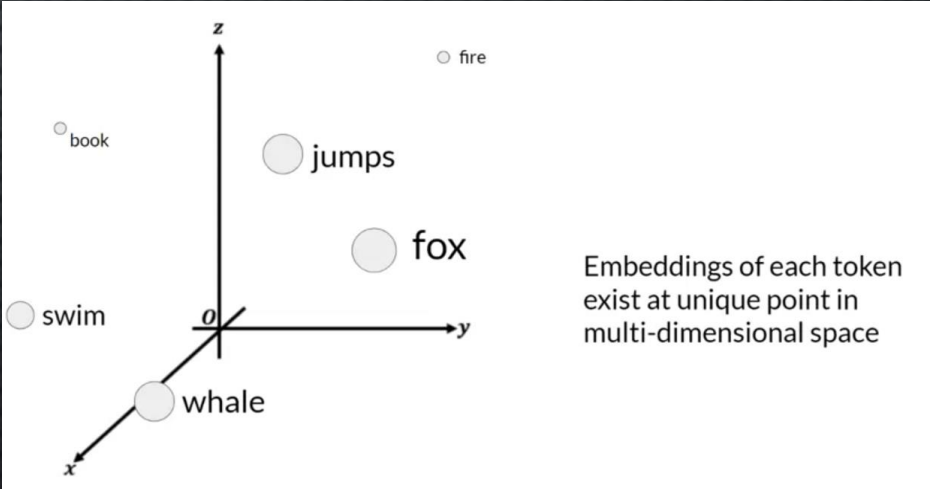
Section 6: Parameter Efficient Fine Tuning (PEFT)

- LoRA (Low Rank Adaption)
- Adapters
- Soft prompts

They can be any values



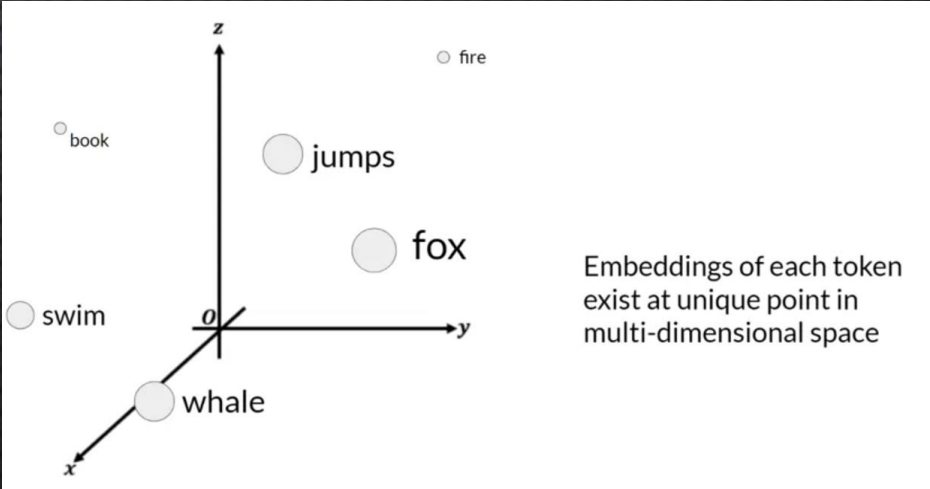
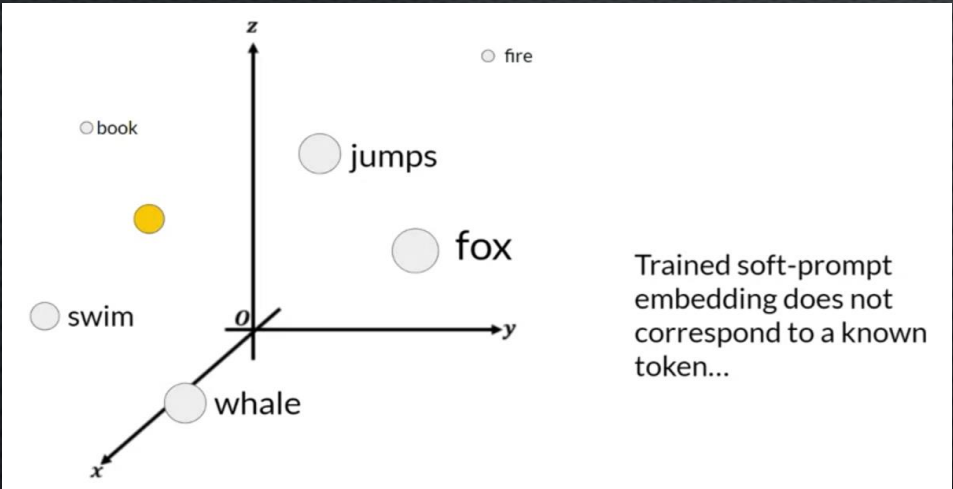
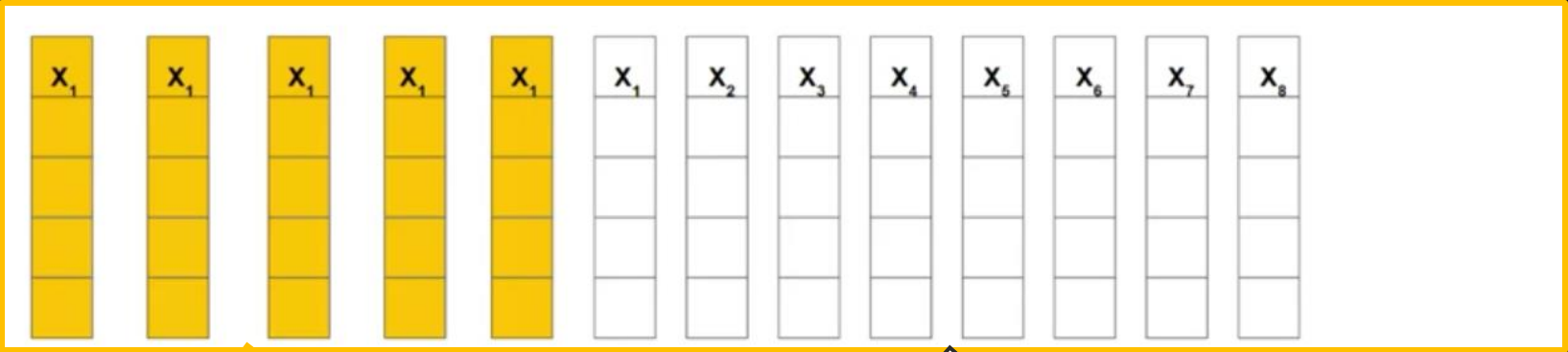
Very parameter efficient strategy



[13] Lester B, Al-Rfou R, Constant N. **The Power of Scale for Parameter-Efficient Prompt Tuning**. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing 2021 Nov* (pp. 3045-3059).



Section 6: Parameter Efficient Fine Tuning (PEFT)

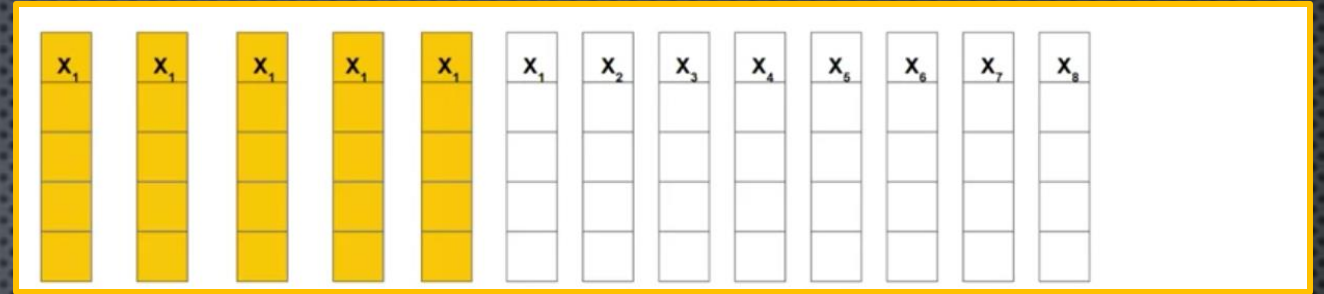


[13] Lester B, Al-Rfou R, Constant N. **The Power of Scale for Parameter-Efficient Prompt Tuning**. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing 2021 Nov* (pp. 3045-3059).

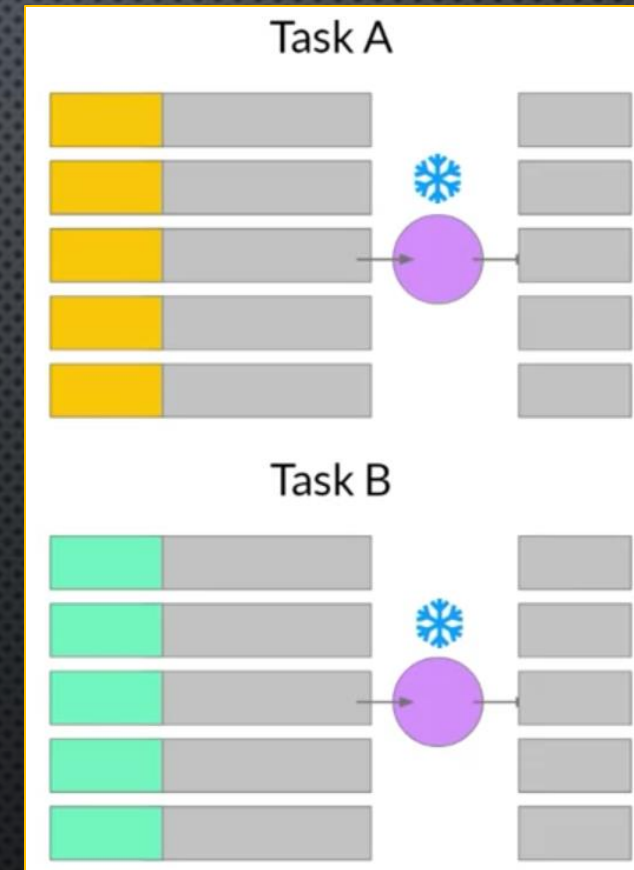


## Section 6: Parameter Efficient Fine Tuning (PEFT)

- **LoRA (Low Rank Adaption)**
- **Adapters (+Layers in Encoder/Decoder)**
- Soft prompts



Task adaptation like LoRA  
Switch out soft prompt at  
inference time to change the  
task





### References Section 5: Fine-Tuning and Domain Adaptation

- [9] Wu S, Irsoy O, Lu S, Dabrovolski V, Dredze M, Gehrmann S, Kambadur P, Rosenberg D, Mann G. **Bloomberggpt: A large language model for finance**. arXiv preprint arXiv:2303.17564. 2023 Mar 30.
- [2] Chung HW, Hou L, Longpre S, Zoph B, Tay Y, Fedus W, Li Y, Wang X, Dehghani M, Brahma S, Webson A. **Scaling instruction-finetuned language models**. *Journal of Machine Learning Research*. **2024**;25(70):1-53.
- [11] Lialin V, Deshpande V, Rumshisky A. **Scaling down to scale up: A guide to parameter-efficient fine-tuning**. arXiv preprint arXiv:2303.15647. 2023 Mar 28.
- [12] Hu EJ, Wallis P, Allen-Zhu Z, Li Y, Wang S, Wang L, Chen W. **LoRA: Low-Rank Adaptation of Large Language Models**. In *International Conference on Learning Representations* 2021 Oct 6.
- [13] Lester B, Al-Rfou R, Constant N. **The Power of Scale for Parameter-Efficient Prompt Tuning**. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing* 2021 Nov (pp. 3045-3059).



# WE ARE GOING TO TALK ABOUT...

## 1. **Beyond Basic LLMs:** Advances Toward Sophisticated Applications

- Retrieval-Augmented Generation (RAG)
- Program-Aided Language Models (PAL)
- Integration of Reasoning and Action (ReAct)
- Architectural Designs for LLM Applications (ChatGPT, Gemini, etc..)

## 2. **Introduction: Transformers**

## 3. **Selecting LLM Architectures:** Determining the Optimal Model for Your Needs

## 4. **LLM Training Resources:** GPU Memory Requirements

- Models: BERT-L (340M), GPT-2 (1.5B), LLaMA-2 (7-13-70B), GPT-3 (175B), PaLM (540B)
- GPU RAM for Models:
  - 1B parameters: 24GB (32-bit precision)
  - 175B parameters: 4200GB (32-bit precision)
  - 500B parameters: 12000GB (32-bit precision)

## 5. **Datasets & Benchmarks:** Criteria for Selecting Evaluation Metrics and Datasets

## 6. **Fine-Tuning Strategies:** Addressing Specific Tasks and Preventing Overfitting

## 7. **RLHF: Enhancing LLMs Through Human Interaction**

## 8. **ReST (Google):** Reinforced Self-Training for Language Modeling

## 9. MED-Gemini



## Section 7: Reinforcement learning from human feedback (RLHF)

### Large Language Models behaving badly:

- RLHF is a fine-tuning process that aligns LLMs with Human Feedbacks

- Training or Fine tuning an LLM is not enough

- Toxic language
- Aggressive responses
- Providing dangerous information
- Etc...

In tasks like:

- Automatic Reporting
- Question Answering
- Summarization
- Etc...

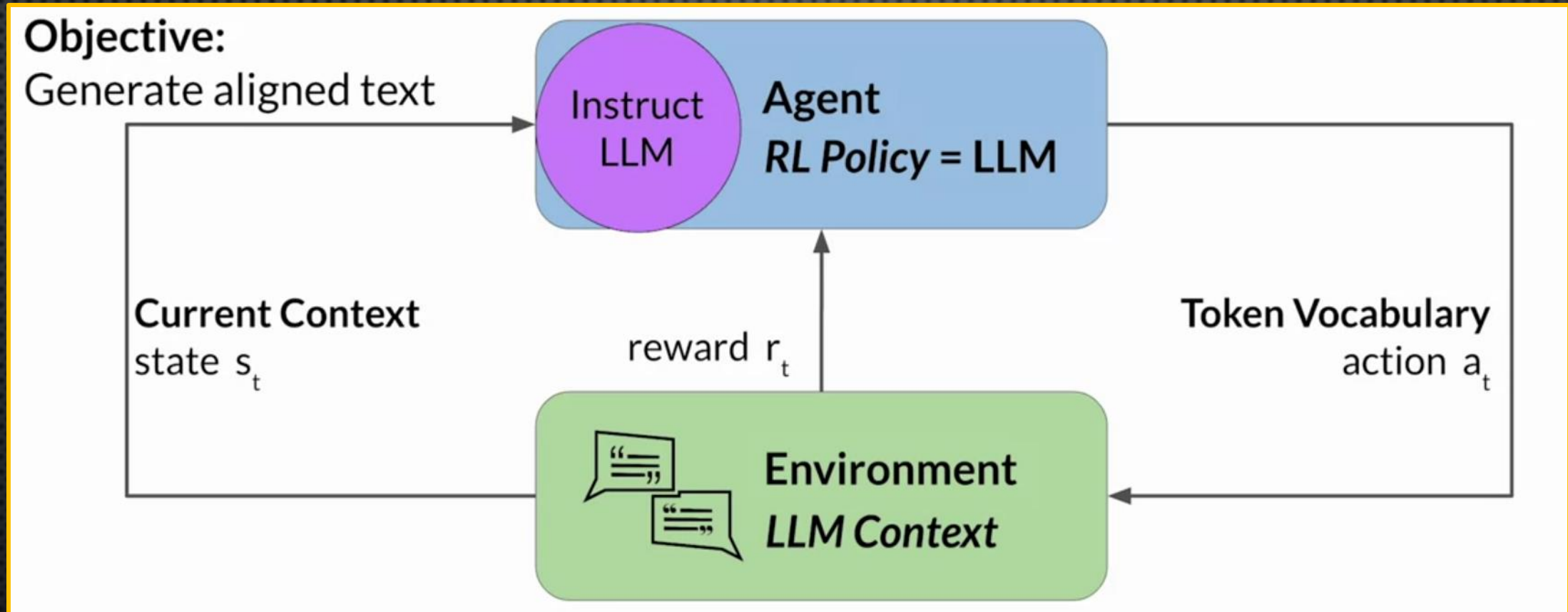
- **Responsible Use of AI: The HHH**

- Helpful?
- Honest?
- Harmless?

- **LLMs need to be aligned with human feedback**

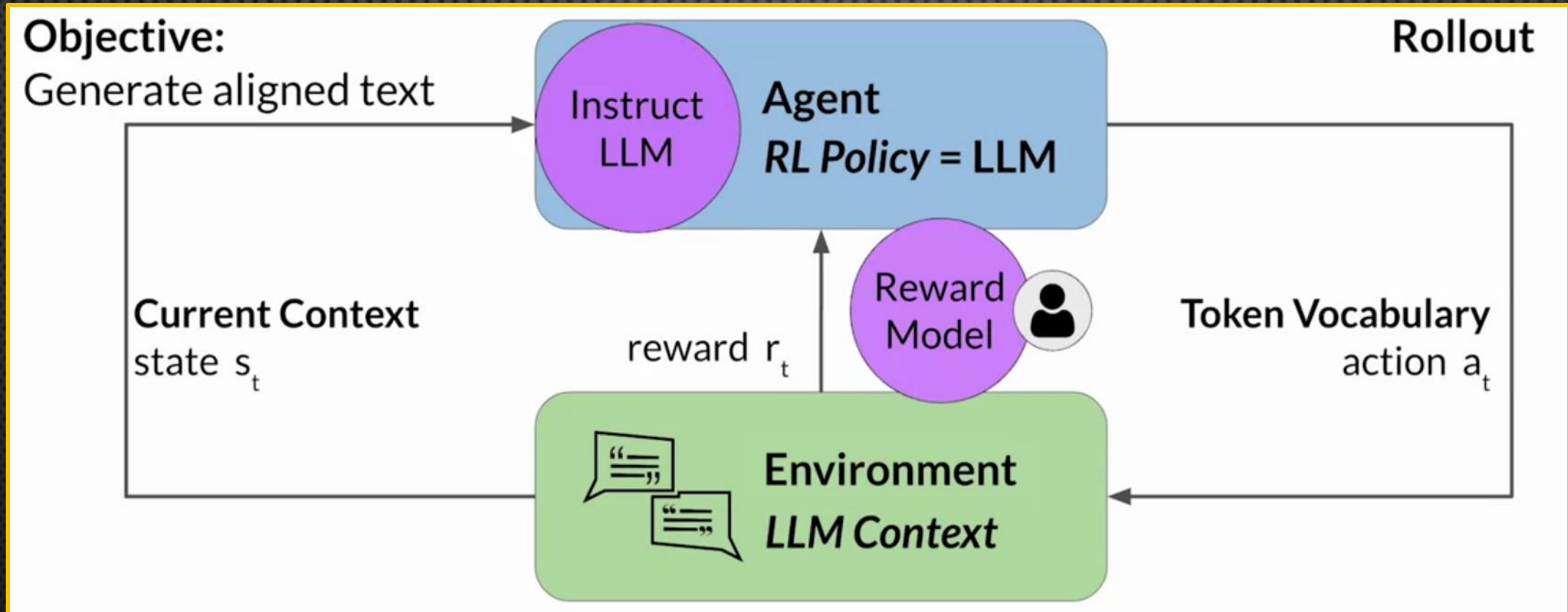


RLHF can minimize harmfulness and maximize helpfulness





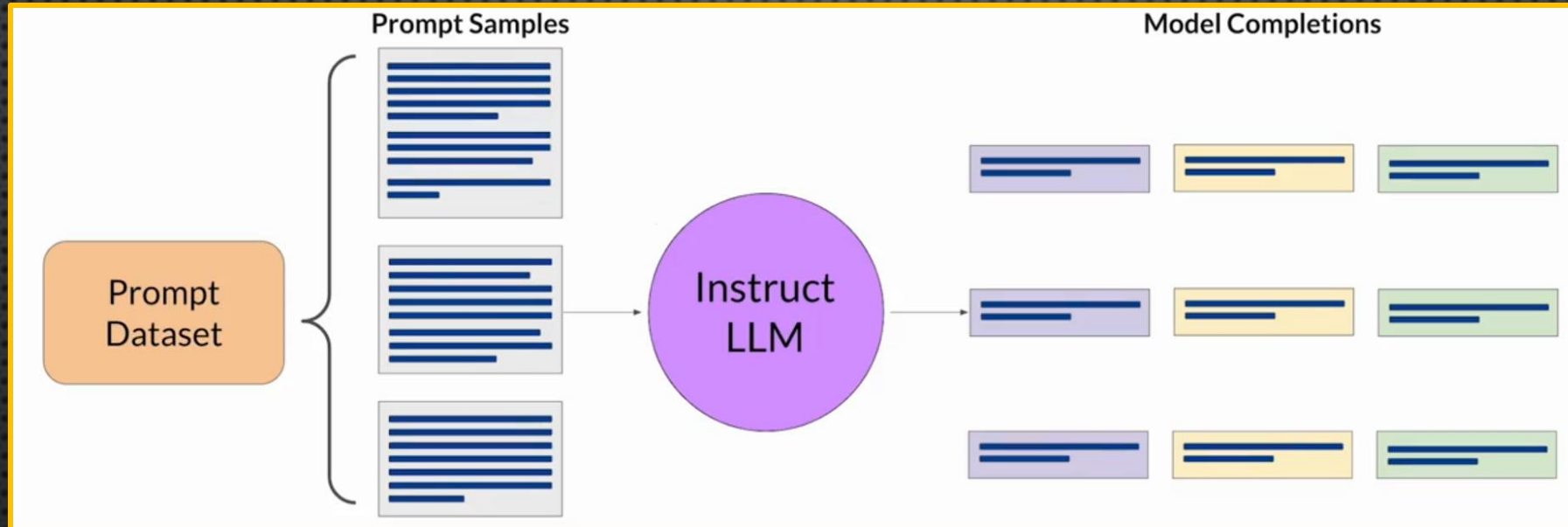
## Reinforcement learning: fine-tune LLMs





## Section 7: Reinforcement learning from human feedback (RLHF)

### Create a prompt dataset based on the task (summarization, Q/A, ...)



LLMs are generative models, they output different completions at every run.

1. **Definition of Scores:** Define what each score (0, 1, 2, ..., N) represents, ensuring that these definitions are precise and relate directly to the quality attributes of the completion (e.g., relevance, correctness, informativeness).
2. **Criteria Specificity:** Criteria should be detailed enough to cover different aspects of the completion such as factual accuracy, fluency, coherence, and alignment with the prompt.



## Section 7: Reinforcement learning from human feedback (RLHF)

### Create a prompt dataset based on the task (summarization, Q/A, ...)

- **Train Labelers**
  - **Initial Training:** Conduct thorough training sessions where labelers are educated on the scoring system and the specific criteria for each score.
  - **Calibration Exercises:** Use calibration exercises to align labeler perceptions and interpretations of the scoring criteria, ensuring consistency in how scores are applied.
- **Multiple Evaluations:** Each completion is rated by multiple labelers to mitigate individual bias and variance in scoring.
- **Averaging Scores:** After all labelers have scored a completion, calculate the average score for each completion. This average becomes the GT for the model output.
- **Handling Outliers:** Consider statistical methods to identify and handle outlier scores that might skew the average, ensuring that the GT represents a consensus view.
- **Consistency Review:** Regularly review the scores from different labelers to check for consistency and alignment with the training provided.
- **Feedback Loops:** Set up feedback loops where labelers can discuss discrepancies and refine their understanding of the criteria.



## Section 7: Reinforcement learning from human feedback (RLHF)

### Create a prompt dataset based on the task (summarization, Q/A, ...)

- **Feedback and Iterative Improvement**

- **Labeler Feedback:** Collect and analyze feedback from labelers on the scoring process and criteria clarity.
- **Guideline Adjustments:** Update and refine the scoring guidelines and training based on labeler feedback and observed scoring trends.

- **Integration with Reward Module**

- **Data Preparation:** Format the averaged scores along with their corresponding completions in a way that can be easily ingested by the reward module.
- **Reward Training:** Use the GT scores as input to train the reward module, which learns to predict the quality of new completions based on historical labeled data.

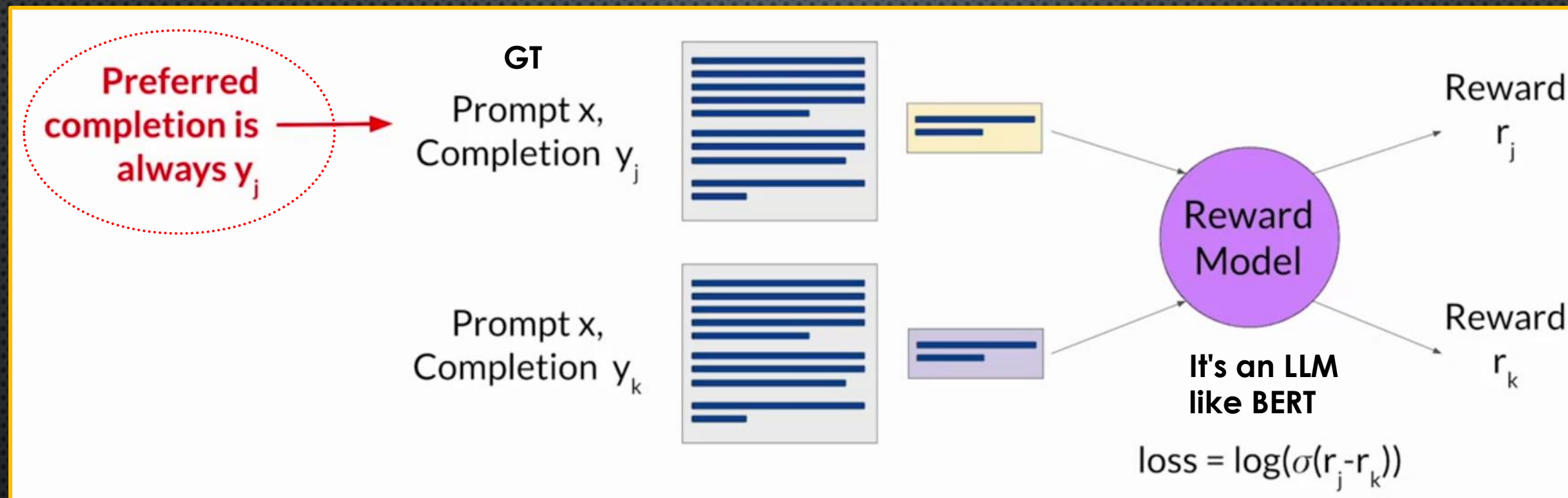
**Once human labelers have completed their assessments off the *prompt-completion sets*, all the data needed to train the reward model is done.**



## Section 7: Reinforcement learning from human feedback (RLHF)

### Reward Model

Train model to predict preferred completion from  $\{y_j, y_k\}$  for prompt  $x$



- The loss function  $\log(\text{sigmoid}(r_j - r_k))$ , is used in learning to rank.
- Reward model can be used as a binary classifier, the reward is the logit value of the positive class.



## Section 7: Reinforcement learning from human feedback (RLHF)

### Reward Model

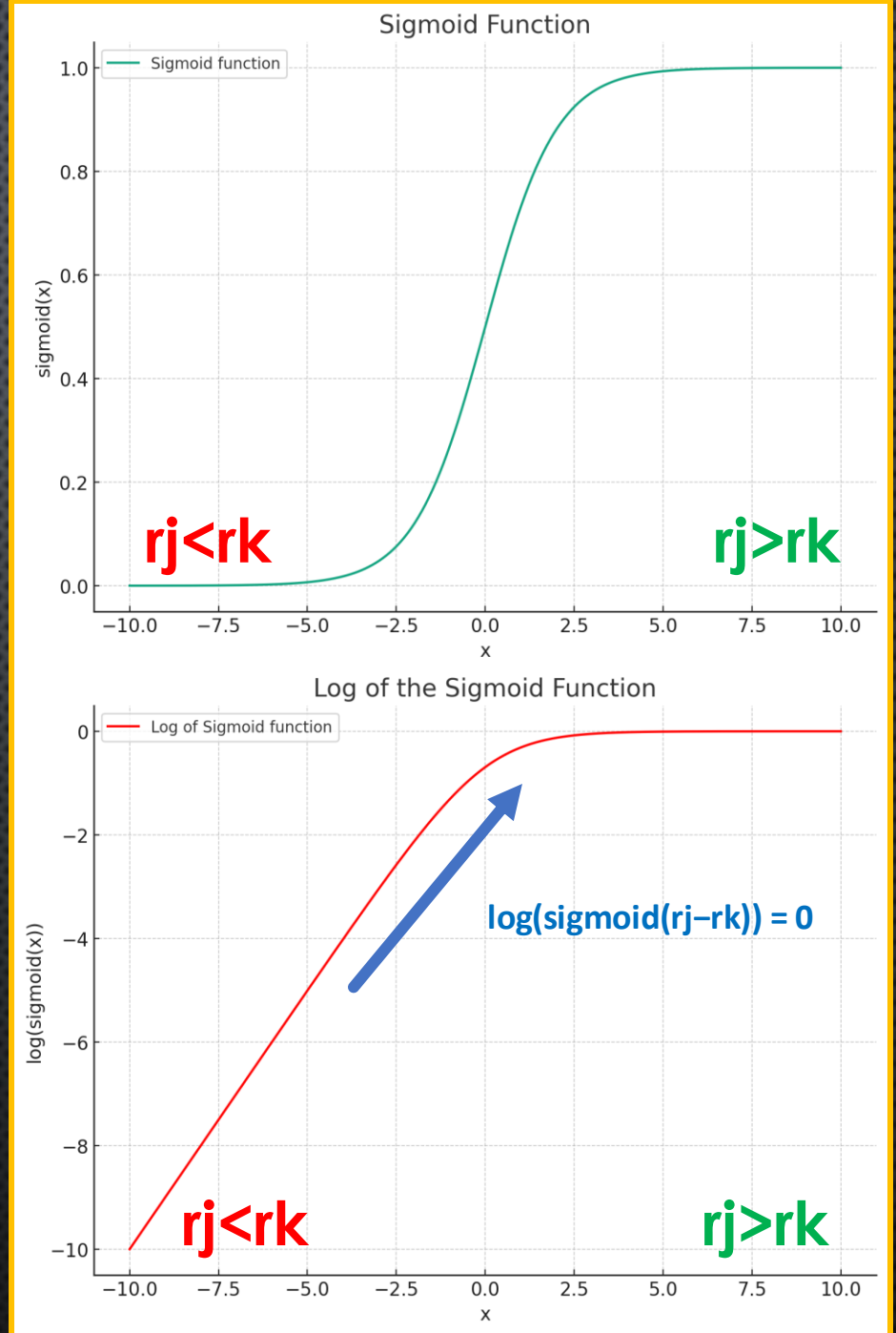
The loss function

$$r_j - r_k$$

$$\log(\text{sigmoid}(r_j - r_k))$$

is used in scenarios like learning to rank, where the goal is to ensure that one item  $r_j$  is ranked higher than another item  $r_k$ .

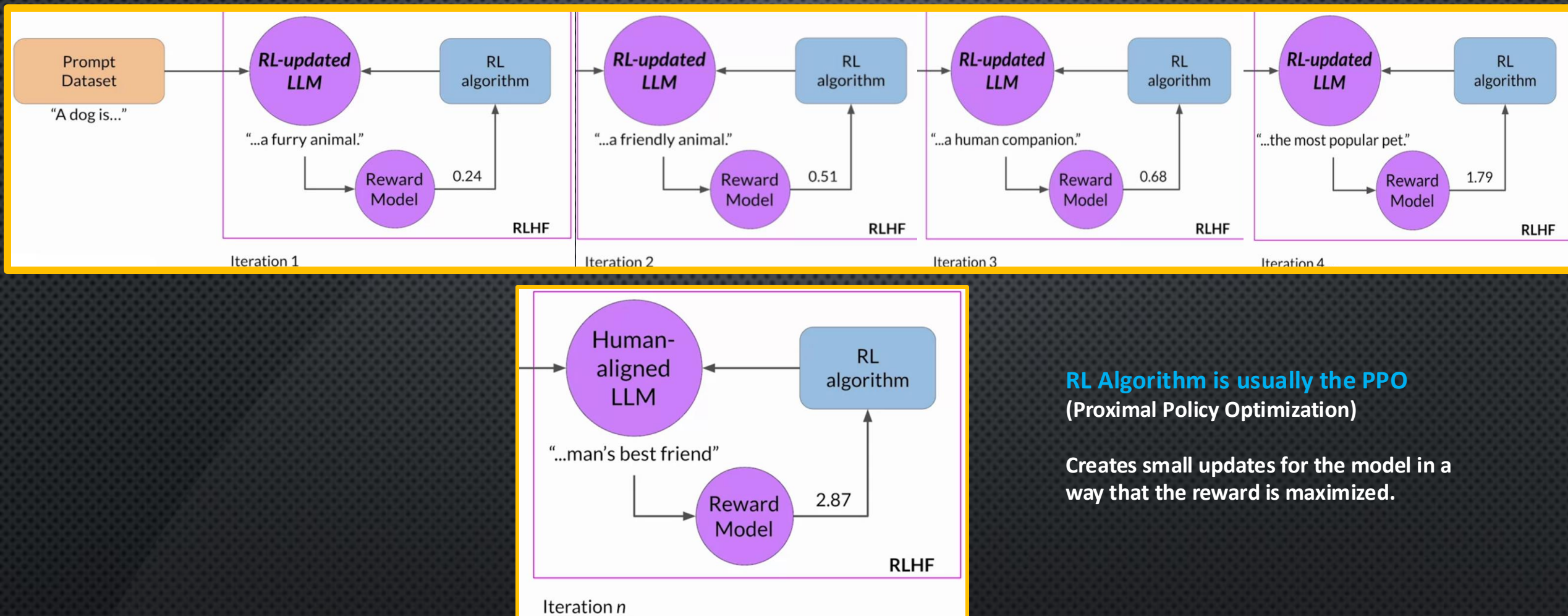
[15] Stiennon N, Ouyang L, Wu J, Ziegler D, Lowe R, Voss C, Radford A, Amodei D, Christiano PF. **Learning to summarize with human feedback**. Advances in Neural Information Processing Systems. 2020;33:3008-21.





## Section 7: Reinforcement learning from human feedback (RLHF)

### Reward Model to fine-tune LLM with RL



**RL Algorithm is usually the PPO**  
(Proximal Policy Optimization)

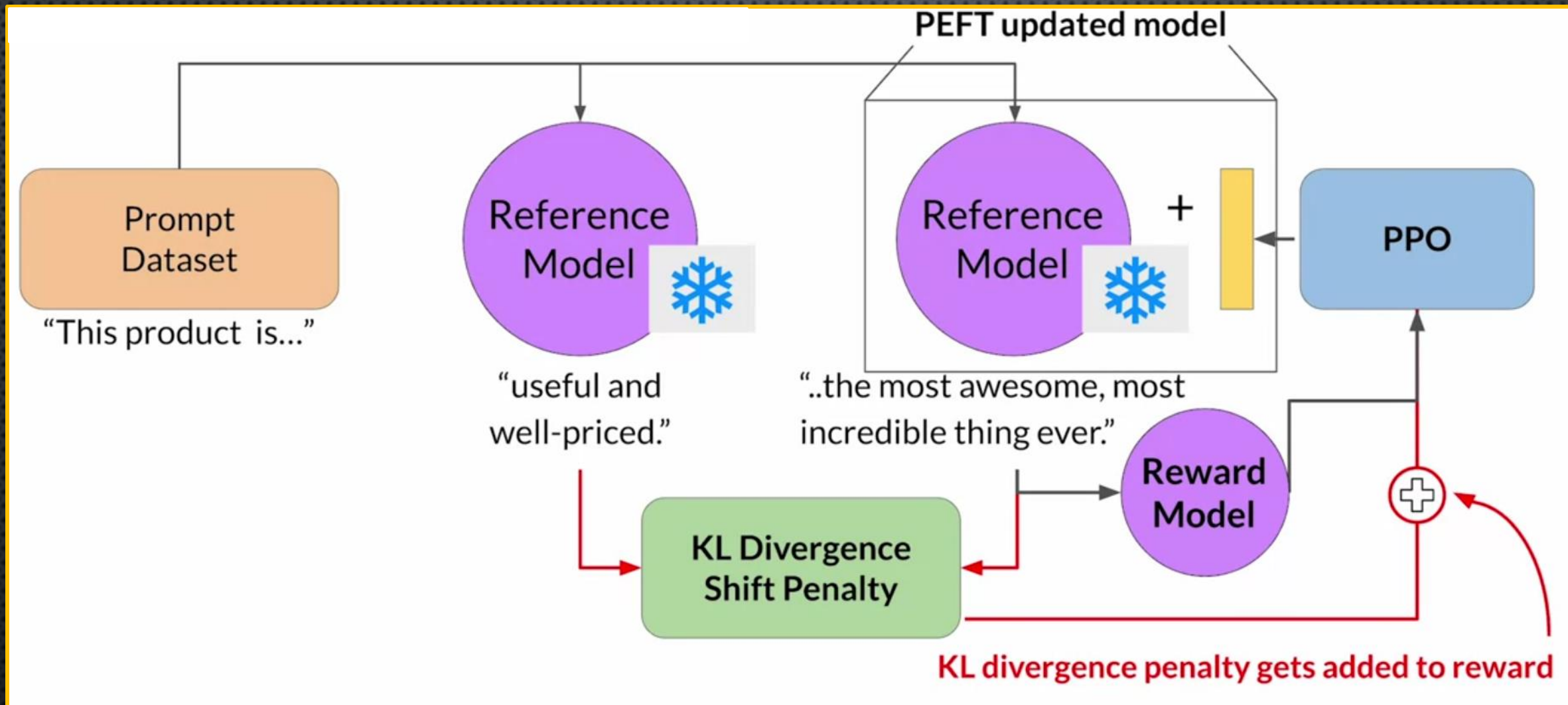
Creates small updates for the model in a way that the reward is maximized.

[15] Stiennon N, Ouyang L, Wu J, Ziegler D, Lowe R, Voss C, Radford A, Amodei D, Christiano PF. **Learning to summarize with human feedback**. Advances in Neural Information Processing Systems. 2020;33:3008-21.

[16] Schulman J, Wolski F, Dhariwal P, Radford A, Klimov O. **Proximal policy optimization algorithms**. arXiv preprint arXiv:1707.06347. 2017 Jul 20.



## Section 7: Reinforcement learning from human feedback (RLHF)



In order to avoid divergence from the original LLM, the KL-div can be added to the reward model to penalize rewards that cause high divergence, moreover, the PEFT can be used for fine-tuning in order to keep the original LLM and specialize the model on different tasks.



### References Section 7: Reinforcement learning from human feedback (RLHF)

- [14] Chung HW, Hou L, Longpre S, Zoph B, Tay Y, Fedus W, Li Y, Wang X, Dehghani M, Brahma S, Webson A. **Scaling instruction-finetuned language models**. Journal of Machine Learning Research. **2024**;25(70):1-53.
- [15] Stiennon N, Ouyang L, Wu J, Ziegler D, Lowe R, Voss C, Radford A, Amodei D, Christiano PF. **Learning to summarize with human feedback**. Advances in Neural Information Processing Systems. 2020;33:3008-21.
- [16] Schulman J, Wolski F, Dhariwal P, Radford A, Klimov O. **Proximal policy optimization algorithms**. arXiv preprint arXiv:1707.06347. 2017 Jul 20.
- [17] Gulcehre C, Paine TL, Srinivasan S, Konyushkova K, Weerts L, Sharma A, Siddhant A, Ahern A, Wang M, Gu C, Macherey W. **Reinforced self-training (rest) for language modeling**. Google Research. arXiv preprint arXiv:2308.08998. 2023 Aug 17.



# WE ARE GOING TO TALK ABOUT...

## 1. **Beyond Basic LLMs:** Advances Toward Sophisticated Applications

- Retrieval-Augmented Generation (RAG)
- Program-Aided Language Models (PAL)
- Integration of Reasoning and Action (ReAct)
- Architectural Designs for LLM Applications (ChatGPT, Gemini, etc..)

## 2. **Introduction:** Transformers

## 3. **Selecting LLM Architectures:** Determining the Optimal Model for Your Needs

## 4. **LLM Training Resources:** GPU Memory Requirements

- Models: BERT-L (340M), GPT-2 (1.5B), LLaMA-2 (7-13-70B), GPT-3 (175B), PaLM (540B)
- GPU RAM for Models:
  - 1B parameters: 24GB (32-bit precision)
  - 175B parameters: 4200GB (32-bit precision)
  - 500B parameters: 12000GB (32-bit precision)

## 5. **Datasets & Benchmarks:** Criteria for Selecting Evaluation Metrics and Datasets

## 6. **Fine-Tuning Strategies:** Addressing Specific Tasks and Preventing Overfitting

## 7. **RLHF:** Enhancing LLMs Through Human Interaction

## 8. **ReST (Google):** Reinforced Self-Training for Language Modeling

## 9. **MED-Gemini**



## Section 8: Reinforced Self-Training

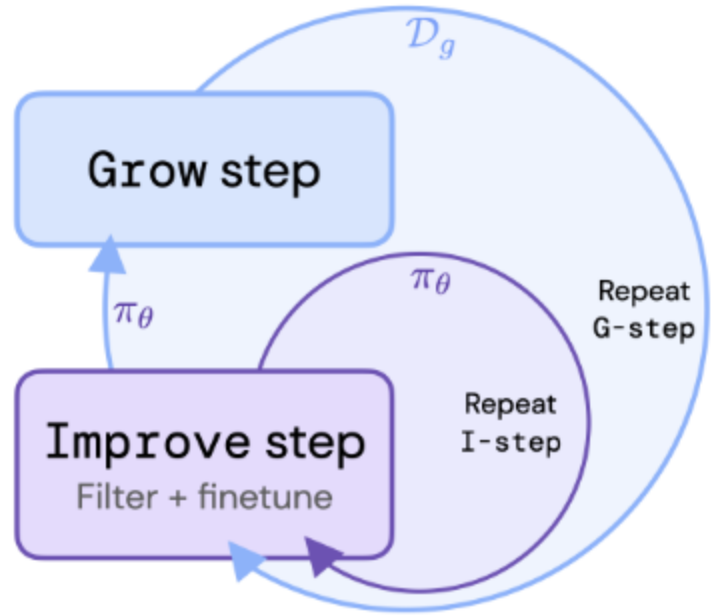
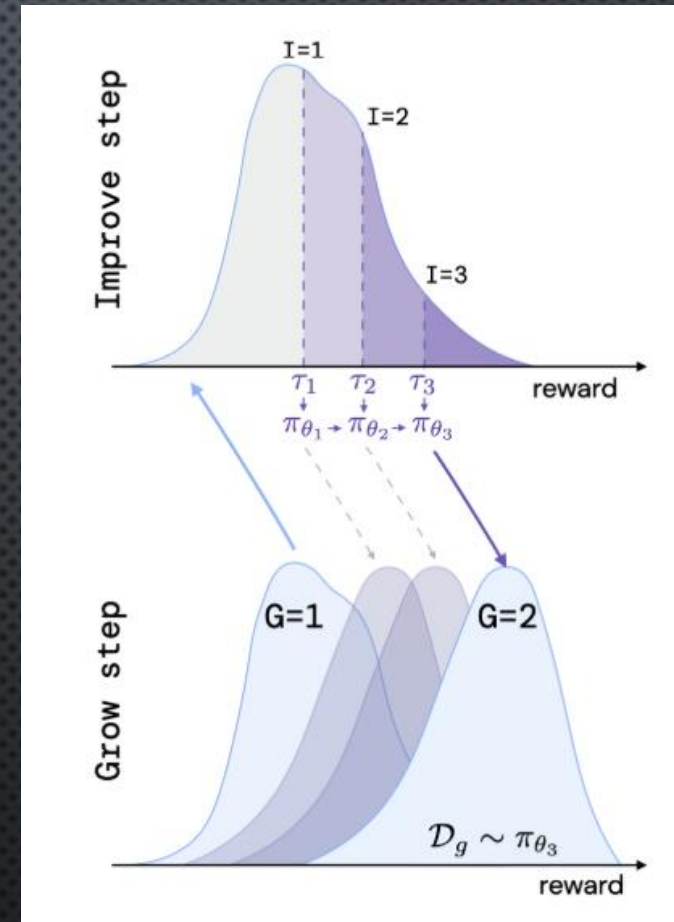


Figure 1 | **ReST method**. During Grow step, a policy generates a dataset. At Improve step, the filtered dataset is used to fine-tune the policy. Both steps are repeated, Improve step is repeated more frequently to amortise the dataset creation cost.

Given an initial LLM policy, ReST produces a dataset by generating samples from the policy, which are then used to improve the LLM policy using offline RL algorithms. ReST is more efficient than typical online RLHF methods because the training dataset is produced offline.



[25] Gulcehre C, Paine TL, Srinivasan S, Konyushkova K, Weerts L, Sharma A, Siddhant A, Ahern A, Wang M, Gu C, Macherey W. **Reinforced self-training (rest) for language modeling**. Google Research & DeepMind. arXiv preprint arXiv:2308.08998. 2023.



## Section 8: Reinforced Self-Training

$$\pi_{\theta}(\mathbf{y} \mid \mathbf{x}) = \prod_{t=1}^T \pi_{\theta}(y_t \mid \mathbf{y}_{1:t-1}, \mathbf{x}),$$

$$\mathcal{L}_{\text{NLL}}(\theta) = -\mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} \left[ \sum_{t=1}^T \log \pi_{\theta}(y_t \mid \mathbf{y}_{1:t-1}, \mathbf{x}) \right]$$

[25] Gulcehre C, Paine TL, Srinivasan S, Konyushkova K, Weerts L, Sharma A, Siddhant A, Ahern A, Wang M, Gu C, Macherey W. **Reinforced self-training (rest) for language modeling**. Google Research & DeepMind. arXiv preprint arXiv:2308.08998. 2023.



## Section 8: Reinforced Self-Training

---

**Algorithm 1: ReST algorithm.** *ReST* is a growing-batch RL algorithm. Given an initial policy of reasonable quality (for example, pre-trained using BC) iteratively applies Grow and Improve steps to update the policy. Here  $F$  is a filtering function, and  $\mathcal{L}$  is a loss function.

---

**Input:**  $\mathcal{D}$ : Dataset,  $\mathcal{D}_{eval}$ : Evaluation dataset,  $\mathcal{L}(\mathbf{x}, \mathbf{y}; \theta)$ : loss,  $R(\mathbf{x}, \mathbf{y})$ : reward model,  $G$ : number of grow steps,  $I$ : number of improve steps,  $N$ : number of samples per context

Train  $\pi_\theta$  on  $\mathcal{D}$  using loss  $\mathcal{L}$ .

**for**  $g = 1$  to  $G$  **do**

    // Grow

    Generate dataset  $\mathcal{D}_g$  by sampling:  $\mathcal{D}_g = \{ (\mathbf{x}^i, \mathbf{y}^i) |_{i=1}^{N_g} \text{ s.t. } \mathbf{x}^i \sim \mathcal{D}, \mathbf{y}^i \sim \pi_\theta(\mathbf{y}|\mathbf{x}^i) \} \cup \mathcal{D}$ .

    Annotate  $\mathcal{D}_g$  with the reward model  $R(\mathbf{x}, \mathbf{y})$ .

**for**  $i = 1$  to  $I$  **do**

        // Improve

        Choose threshold s.t.  $\tau_1 > V_{\pi_\theta}$  for  $V_{\pi_\theta} = \mathbb{E}_{\mathcal{D}_g}[R(\mathbf{x}, \mathbf{y})]$  and  $\tau_{i+1} > \tau_i$ .

**while** reward improves on  $\mathcal{D}_{eval}$  **do**

            Optimise  $\theta$  on objective:  $J(\theta) = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}_g} [F(\mathbf{x}, \mathbf{y}; \tau_i) \mathcal{L}(\mathbf{x}, \mathbf{y}; \theta)]$

**end**

**end**

**end**

**Output:** Policy  $\pi_\theta$

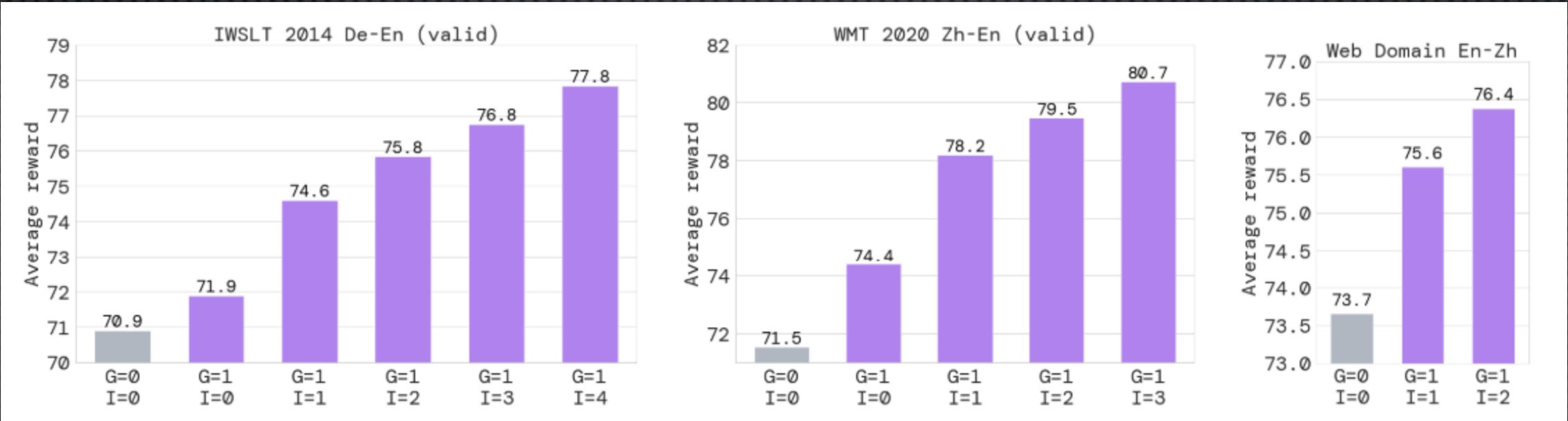
---



Section 8: Reinforced Self-Training

| Algorithm       | Average Reward | Distinct samples |
|-----------------|----------------|------------------|
| BC (G=0, I=0)   | 70.9           | 16 000 000       |
| ReST (G=1, I=0) | 71.9           | 16 000 000       |
| ReST (G=1, I=4) | 77.8           | 16 000 000       |
| ReST (G=2, I=3) | 83.1           | 32 000 000       |
| Online RL       | 71.6           | 24 000 000       |

Table 1 | **Online RL for IWSLT 2014:** Online RL performs as well as *ReST* (G=1, I=0) and *ReST* (G=1, I=4) is significantly better.



[25] Gulcehre C, Paine TL, Srinivasan S, Konyushkova K, Weerts L, Sharma A, Siddhant A, Ahern A, Wang M, Gu C, Macherey W. **Reinforced self-training (rest) for language modeling**. Google Research & DeepMind. arXiv preprint arXiv:2308.08998. 2023.



# Summary



## Consumers



**Application Interfaces** e.g. Websites, Mobile Applications, APIs, etc.

**LLM Tools & Frameworks** e.g. LangChain, Model Hubs

### Information Sources



Documents



Database



Web

### LLM Models

Optimized  
LLM

### Generated Outputs & Feedback



**Infrastructure** e.g. Training/Fine-Tuning, Serving, Application Components



## Summary

There are different chains (pipeline of solutions) for each task.

### Application reasoning engine only

LLM

### External Data Sources



Documents



Databases



Web

User Application

Prompts...

LangChain  
[www.langchain.com](http://www.langchain.com) (Python)

#### Tools

(pre-build for a variety of tasks –  
API,  
WebSearch...)

**Agents**  
(PAL, ReAct, ..)

**Prompt Templates**  
(included for many use cases)

**Memory**

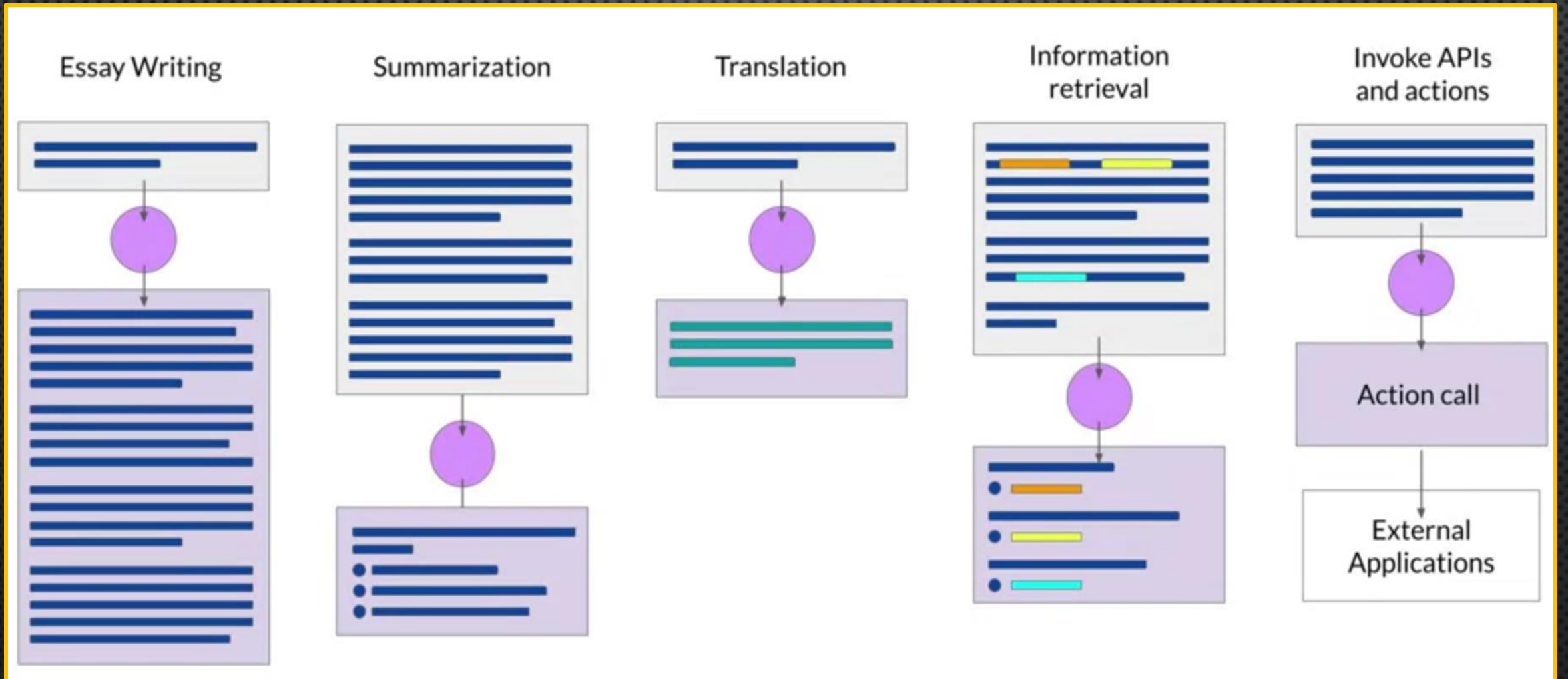
(to store interactions with an LLM)

External Applications

API call



## Summary





Summary

|                   | Pre-training                                                                                                                                                  | Prompt engineering                                | Prompt tuning and fine-tuning                                                                      | Reinforcement learning/human feedback                                                                                       | Compression/ optimization/ deployment                                                                            |
|-------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------|----------------------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------------------|
| Training duration | Days to weeks to months                                                                                                                                       | Not required                                      | Minutes to hours                                                                                   | Minutes to hours similar to fine-tuning                                                                                     | Minutes to hours                                                                                                 |
| Customization     | Determine model architecture, size and tokenizer.<br><br>Choose vocabulary size and # of tokens for input/context<br><br>Large amount of domain training data | No model weights<br><br>Only prompt customization | Tune for specific tasks<br><br>Add domain-specific data<br><br>Update LLM model or adapter weights | Need separate reward model to align with human goals (helpful, honest, harmless)<br><br>Update LLM model or adapter weights | Reduce model size through model pruning, weight quantization, distillation<br><br>Smaller size, faster inference |
| Objective         | Next-token prediction                                                                                                                                         | Increase task performance                         | Increase task performance                                                                          | Increase alignment with human preferences                                                                                   | Increase inference performance                                                                                   |



# References

*Special thanks to Shelbee Eigenbrode, Antje Barth, and Mike Chambers for their work on "Generative AI with Large Language Models" (Coursera, Amazon AWS, 2023), which significantly informed and inspired the material used in these slides.*



## Resources:

### The Tasks, Datasets and benchmarks for LLMs:

- **GLUE**
- **SuperGLUE**
- **HELM**
- **MMLU (Massive Multitask Language Understanding)**
- **BIG-bench**
- **BIG-bench Hard**
- **Lite**

[2] Chung HW, Hou L, Longpre S, Zoph B, Tay Y, Fedus W, Li Y, Wang X, Dehghani M, Brahma S, Webson A. **Scaling instruction-finetuned language models**. Journal of Machine Learning Research. 2024;25(70):1-53.

[4] Wang A, Singh A, Michael J, Hill F, Levy O, Bowman SR. **GLUE: A multi-task benchmark and analysis platform for natural language understanding**. arXiv preprint arXiv:1804.07461. 2018 Apr 20.

[5] Wang A, Pruksachatkun Y, Nangia N, Singh A, Michael J, Hill F, Levy O, Bowman S. **Superglue: A stickier benchmark for general-purpose language understanding systems**. Advances in neural information processing systems. 2019;32.

[6] Hendrycks D, Burns C, Basart S, Zou A, Mazeika M, Song D, Steinhardt J. **Measuring massive multitask language understanding**. arXiv preprint arXiv:2009.03300. 2020 Sep 7.

[7] Suzgun M, Scales N, Schärli N, Gehrmann S, Tay Y, Chung HW, Chowdhery A, Le QV, Chi EH, Zhou D, Wei J. **Challenging big-bench tasks and whether chain-of-thought can solve them**. arXiv preprint arXiv:2210.09261. 2022 Oct 17.

+ Biology, Reasoning, Math, and more

[8] Liang P, Bommasani R, Lee T, Tsipras D, Soylu D, Yasunaga M, Zhang Y, Narayanan D, Wu Y, Kumar A, Newman B. **Holistic evaluation of language models**. arXiv preprint arXiv:2211.09110. 2022 Nov 16.



## Resources

- [9] Wu S, Irsoy O, Lu S, Dabrovolski V, Dredze M, Gehrmann S, Kambadur P, Rosenberg D, Mann G. **Bloomberggpt: A large language model for finance**. arXiv preprint arXiv:2303.17564. 2023 Mar 30.
- [12] Chung HW, Hou L, Longpre S, Zoph B, Tay Y, Fedus W, Li Y, Wang X, Dehghani M, Brahma S, Webson A. **Scaling instruction-finetuned language models**. *Journal of Machine Learning Research*. **2024**;25(70):1-53.
- [11] Lialin V, Deshpande V, Rumshisky A. **Scaling down to scale up: A guide to parameter-efficient fine-tuning**. arXiv preprint arXiv:2303.15647. 2023 Mar 28.
- [12] Hu EJ, Wallis P, Allen-Zhu Z, Li Y, Wang S, Wang L, Chen W. **LoRA: Low-Rank Adaptation of Large Language Models**. In *International Conference on Learning Representations 2021* Oct 6.
- [14] Chung HW, Hou L, Longpre S, Zoph B, Tay Y, Fedus W, Li Y, Wang X, Dehghani M, Brahma S, Webson A. **Scaling instruction-finetuned language models**. *Journal of Machine Learning Research*. **2024**;25(70):1-53.
- [15] Stiennon N, Ouyang L, Wu J, Ziegler D, Lowe R, Voss C, Radford A, Amodei D, Christiano PF. **Learning to summarize with human feedback**. *Advances in Neural Information Processing Systems*. 2020;33:3008-21.
- [16] Schulman J, Wolski F, Dhariwal P, Radford A, Klimov O. **Proximal policy optimization algorithms**. arXiv preprint arXiv:1707.06347. 2017 Jul 20.
- [17] Lewis P, Perez E, Piktus A, Petroni F, Karpukhin V, Goyal N, Küttler H, Lewis M, Yih WT, Rocktäschel T, Riedel S. **Retrieval-augmented generation for knowledge-intensive nlp tasks**. *Advances in Neural Information Processing Systems*. 2020;33:9459-74.
- [18] Wei J, Wang X, Schuurmans D, Bosma M, Xia F, Chi E, Le QV, Zhou D. **Chain-of-thought prompting elicits reasoning in large language models**. *Advances in neural information processing systems*. 2022 Dec 6;35:24824-37.
- [19] Gao L, Madaan A, Zhou S, Alon U, Liu P, Yang Y, Callan J, Neubig G. Pal: **Program-aided language models**. In *International Conference on Machine Learning 2023* Jul 3 (pp. 10764-10799). PMLR.
- [20] Topsakal O, Akinci TC. **Creating large language model applications utilizing langchain: A primer on developing llm apps fast**. In *International Conference on Applied Engineering and Natural Sciences 2023* Jul (Vol. 1, No. 1, pp. 1050-1056). <https://www.langchain.com/>
- [21] Yao S, Zhao J, Yu D, Shafran I, Narasimhan KR, Cao Y. **ReAct: Synergizing Reasoning and Acting in Language Models**. In *NeurIPS 2022 Foundation Models for Decision Making Workshop 2022* Nov 18.
- [22] Pandya K, Holia M. **Automating Customer Service using LangChain: Building custom open-source GPT Chatbot for organizations**. arXiv e-prints. **2023** Oct:arXiv-2310.



### 1. Transformer Architecture

- Attention is All You Need  
This paper introduced the Transformer architecture, with the core “self-attention” mechanism. This article was the foundation for LLMs.
- BLOOM: BigScience 176B Model  
BLOOM is a open-source LLM with 176B parameters trained in an open and transparent way. In this paper, the authors present a detailed discussion of the dataset and process used to train the model + high-level overview of the model here
- Scaling Laws for Neural Language Models  
Empirical study by researchers at OpenAI exploring the scaling laws for large language models.

### 2. Model architectures and pre-training objectives

- What Language Model Architecture and Pretraining Objective Work Best for Zero-Shot Generalization?  
The paper examines modeling choices in large pre-trained language models and identifies the optimal approach for zero-shot generalization.
- HuggingFace Tasks and Model Hub - Collection of resources to tackle varying machine learning tasks using the HuggingFace library.
- LLaMA: Open and Efficient Foundation Language Models  
Article from Meta AI proposing Efficient LLMs (their model with 13B parameters outperform GPT3 with 175B parameters on most benchmarks)

### 3. Scaling laws and compute-optimal models

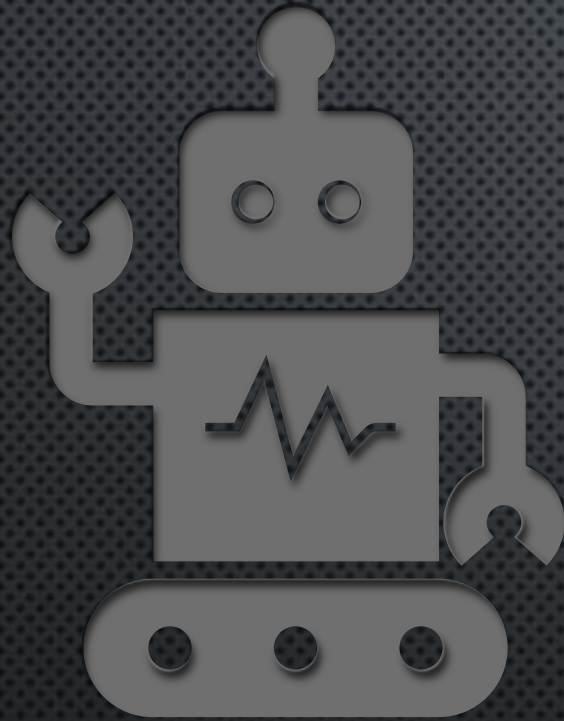
- Language Models are Few-Shot Learners  
This paper investigates the potential of few-shot learning in Large Language Models.
- Training Compute-Optimal Large Language Models  
Study from DeepMind to evaluate the optimal model size and number of tokens for training LLMs. Also known as “Chinchilla Paper”.
- BloombergGPT: A Large Language Model for Finance  
LLM trained specifically for the finance domain, a good example that tried to follow chinchilla laws.



*Thank you for your time and attention*

Any Questions?





# INTRODUCTION TO ADVANCED LARGE LANGUAGE MODELS

A COMPREHENSIVE TUTORIAL ON TECHNIQUES,  
ARCHITECTURES, AND PRACTICAL APPLICATIONS

Giorgio Roffo, PhD

*Explore and Connect With My Professional Work:*

LinkedIn: Giorgio Roffo - [LinkedIn](#)

ResearchGate: [Work Done](#)

Google Scholar: [My Publications](#)

GitHub: [giorgioroffo](#)