

Introduction to Large Language Models(LLM)

Outline

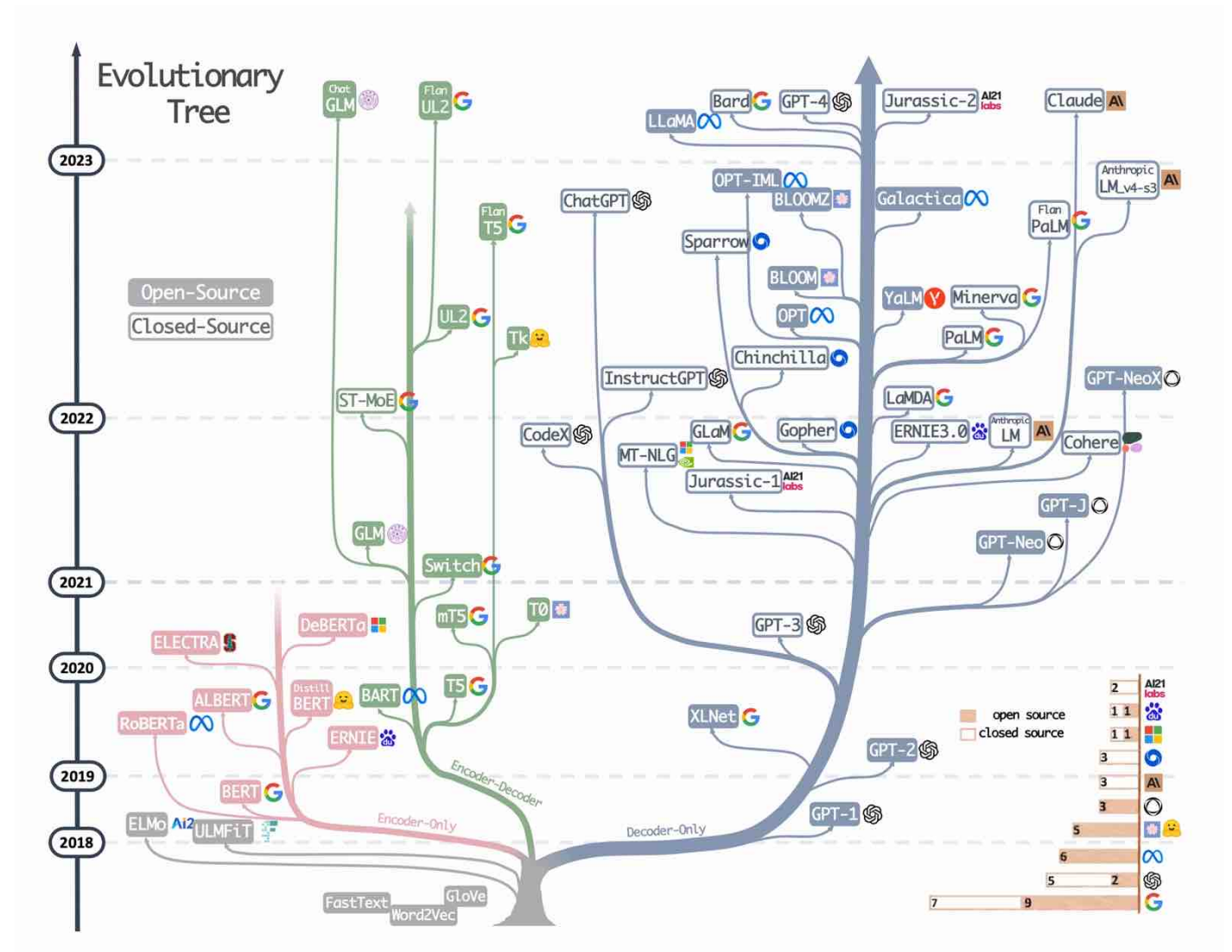
1. Overview from 30,000 feet above
2. Transformer in nutshell
3. Pretraining - Parallel paradigm
4. Finetuning - Parameter efficient finetuning
5. Steering the decoding of LLM - Prompting
6. Augmentation and Plugins

Overview from 30000 feet above

- Paradigm transition in AI
 - From: **training**(specific) -> **prediction**(specific)
 - To: **pretraining**(general) -> **finetuning**(general/specific) -> **in-context prompting**(specific)
- Primary steps of new paradigm
 - Pretraining with self-supervised learning
 - Finetuning on instruction from multiple domains
 - Application by steering the decoding process of LLM
- Where are we to AGI?
 - From explanation to prediction
 - From correlation to causality

A LLMs tree

- Decoder Only
- Encoder Only
- Encoder-Decoder

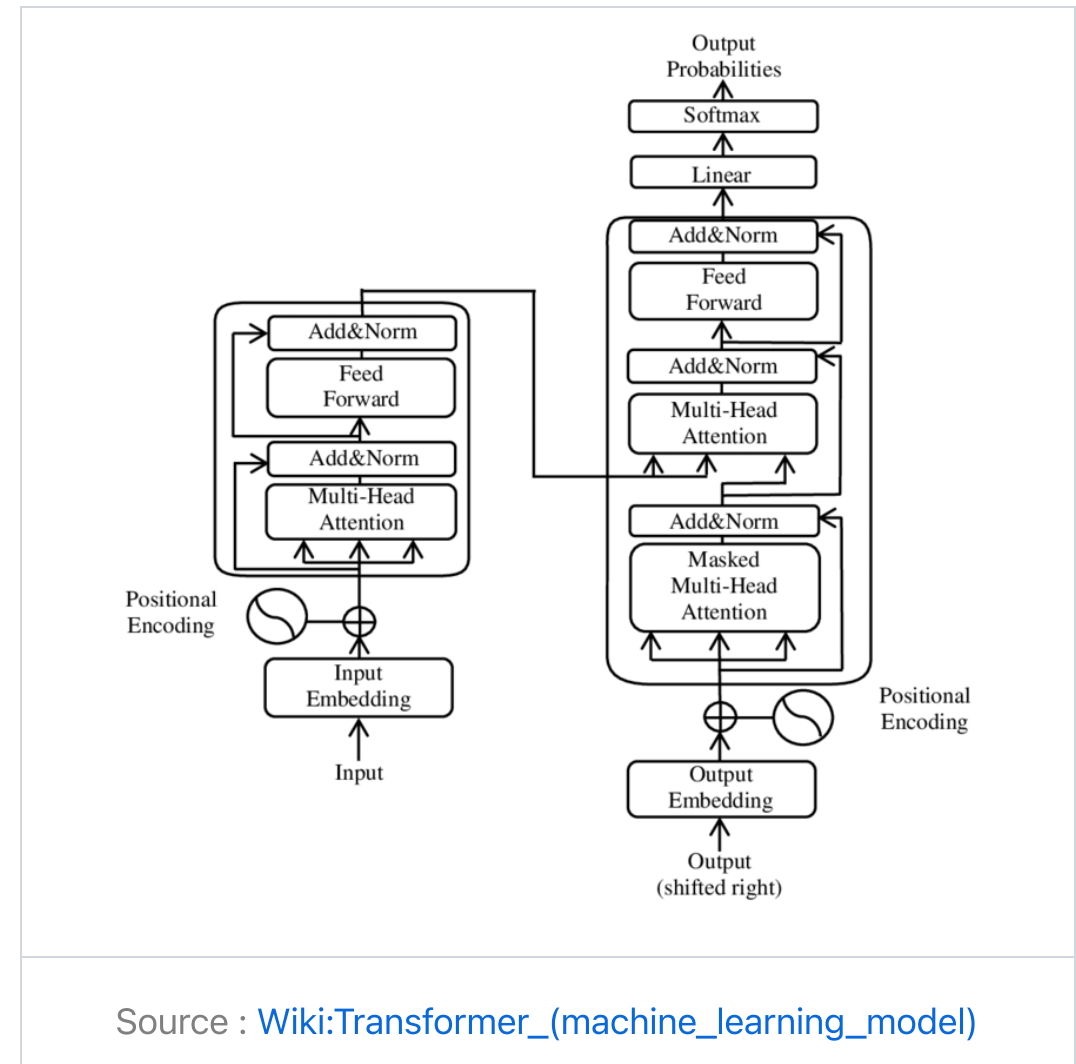


Transformer in nutshell

- Transformer modules
- Aspects of alternative
- Parameter concentration
- Computation concentration

Transformer modules

- Token & positional embedding :
 - $E \in R^{V \times d}, P \in R^{T \times d}$
- Multi-head attention
 - Self-attention & Cross-attention
 - Weight matrix:
$$W_Q, W_K, W_V \in R^{d \times d_h}$$
 - Head projection:
$$W_O \in R^{d \times d}$$
- ResNet & LayerNorm
- Feedforward Network
 - $W_1 \in R^{d \times 4d}, W_2 \in R^{4d \times d}$
- Output head: task related



Aspects of alternative

- Efficient Transformer (for long sequence)
 - Coarse sequence resolution:
 - Block/Stride/Clustering/Neural Memory
 - TransformerXL
 - Attention matrix approximation:
 - Linear Transformer
- LLMs specific
 - **Encoder vs Decoder vs Encoder-Decoder**
 - Pretraining objective
 - Positional encoding
 - Input or output LayerNorm
 - Activation

Parameter concentration

Decoder only Transformer(GPT)

- Parameter size:
 - Embedding: $(V + T) \times d$
 - Attention: $L \times (3 \times d \times d_h \times H + (d_h \times H) \times d) = 4Ld^2$
 - FFN: $L \times ((d \times 4d + 4d) + (4d \times d + d)) \approx 8Ld^2$
- On GPT3-175B:

| Total | PE | TE | Attn | FFN |
|----------|------------|-------------|-----------------|------------------|
| 174,597M | 25M(0.01%) | 617M(0.35%) | 57,982M(33.21%) | 115,970M(66.42%) |

Computation concentration

Decoder only Transformer

Per-token calculation:

- QKV+project: $2 \times L \times (3 \times H \times h_d \times d + (H \times h_d) \times d) = 2 \times 4Ld^2$
- Attention: $2 \times L \times T \times d$
- FFN: $2 \times L \times ((d \times 4d + 4d) + (4d \times d + d)) \approx 2 \times 8Ld^2$

Model training flops utilization(MFU):

- Forward and backward: $(1 + 2) \times (2N + 2LTd) \approx 6N$, where $N \approx 12Ld^2$
- Theoretical peak throughput: $\frac{P}{6N+6LTd}$
- $\text{MFU} = \frac{\text{Observed throughput}}{\text{Theoretical peak throughput}}$

Pretraining

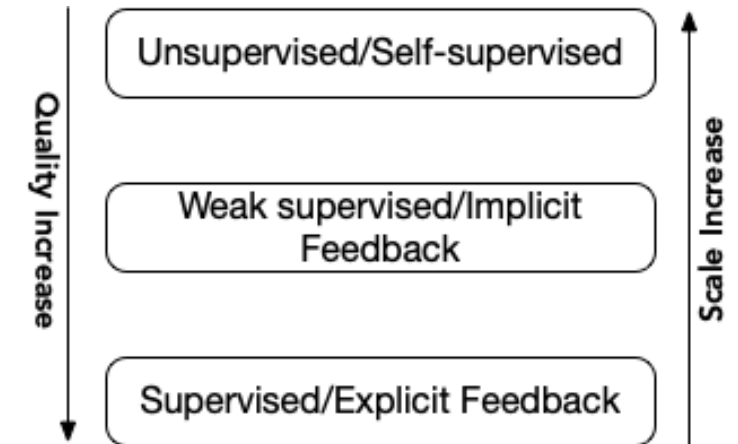
- Training objectives
- Text Corpus
- Parallel strategies
- Results & Evaluation

Training objectives

- Auto-regressive Language Models
- Missing token prediction

Text Corpus

- Unsupervised text
 - [BookCorpus](#): 7000 unpublished books
 - [WebText](#): 8M outlinks of Reddit.com
 - [Common Crawl](#)
- Weak supervised text
 - Summarization: [Reddit TL;DR](#)
 - QA: [StackExchange](#)
- Supervised text
 - Summarization:
 - Q&A:
 - Dialog: InstructGPT



Parallel strategies

- Why bother?
- Four parallel paradigms
 - Data parallel
 - Model parallel
 - Tensor parallel
 - Pipeline parallel
 - Sequence parallel
- Combined implementations

Why bother?

Data parallel

Tensor parallel

Pipeline parallel

Sequence parallel

Combined implementations

End of Parallel strategies

Result & Evaluations

Finetuning

- Target and issues
- Instruct finetuning
- Finetuning for specific task
- Parameter efficient finetuning

Target and issues

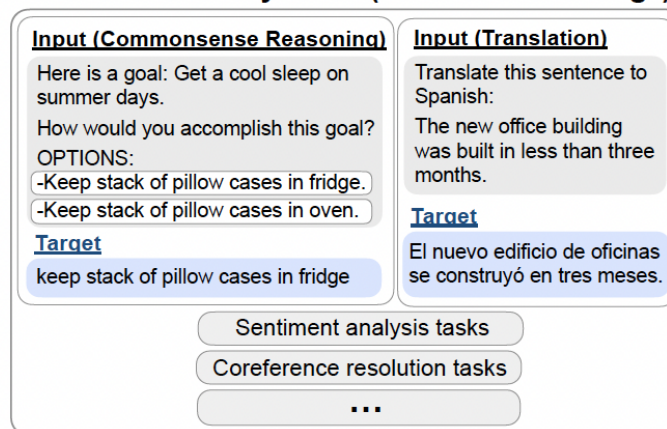
- Target
 - From pattern completion to real world tasks
- Issues
 - Instruction following
 - Hallucination
 - Toxicity and ethics
 - Securities

Instruct Finetuning

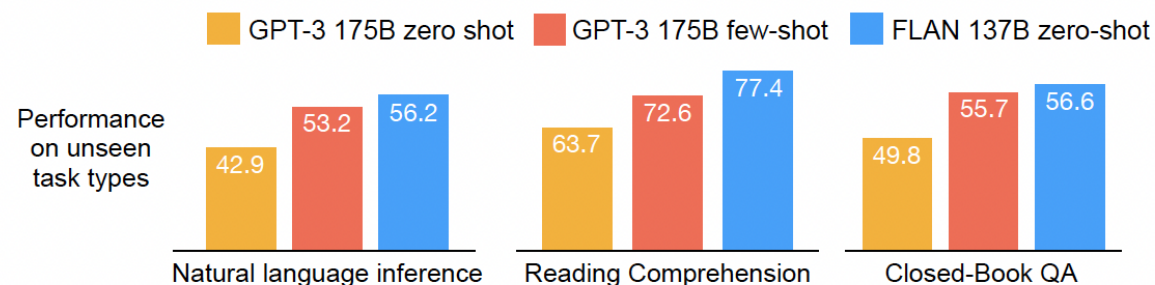
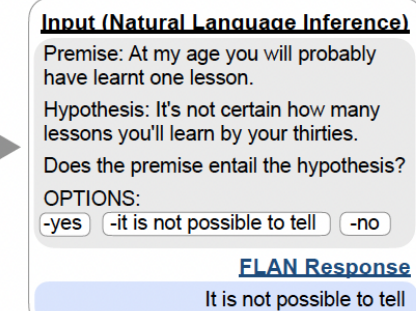
Key to success

- Number of finetuning datasets:
 - scaling from 62 text datasets to 18K
- Model scale: 137B LaMDA-PT
- Natural language instructions

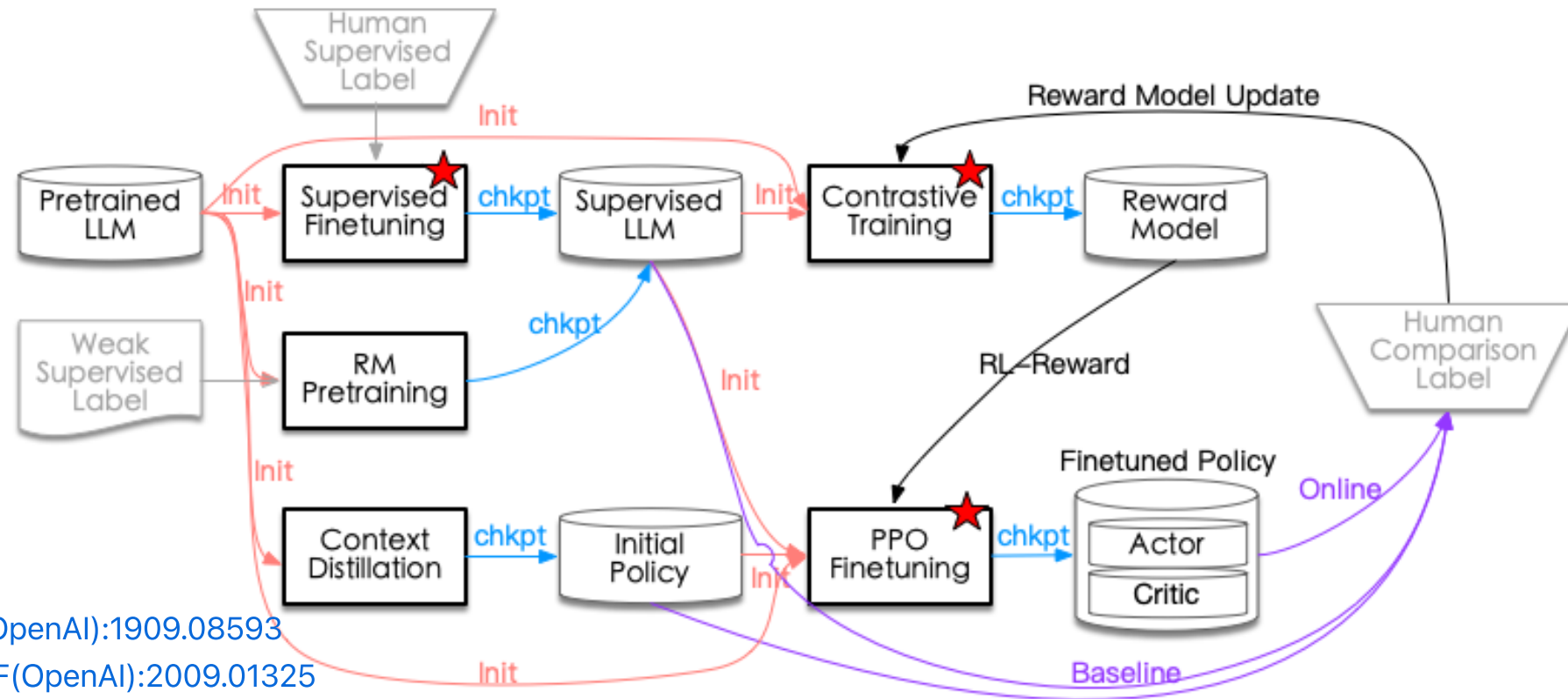
Finetune on many tasks (“instruction-tuning”)



Inference on unseen task type



Finetuning for specific task



Finetune from HF(OpenAI):1909.08593
Summarize from HF(OpenAI):2009.01325
HHHA(Anthropic):2112.00861
InstructGPT(OpenAI):2203.02155
HHA(Anthropic):2204.05862
Sparrow(Deepmind):2209.14375

Problems from Supervised Finetuning(SFT)

- Learning only the task format and the way to response for the format
- Knowledge labeled but not in the LLM leads to more hallucination
- Knowledge in the LLM but labeled as `don't know` leads to withhold information

What we want from finetuning:

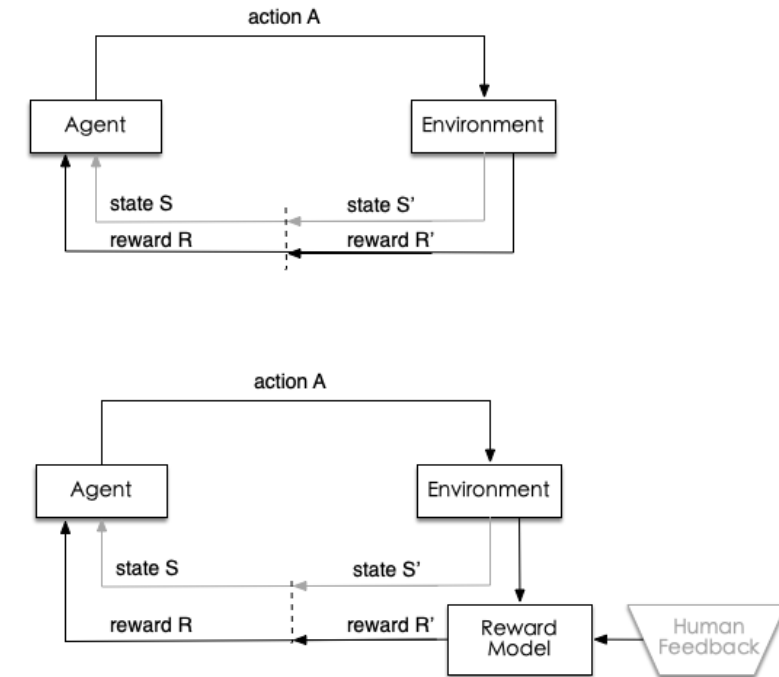
outputs its(LLM's) state of knowledge with the correct amount of hedging and expressing its uncertainty

Advance of RL to SFT for truthfulness

- LLMs know what they know
 - Calibrated probability, uncertainty
- RLHF can leverage the self-awareness
 - Design reward function: correct answer=1 , don't know=0 and wrong answer=-4
 - RL learn optimal threshold of probability to maximize the reward
- No oracle for the correctness, delegate to Reward Model
 - Reward model: **relative criteria** trained by pairwise loss from human feedback
 - Open problem: true probabilities of everything?
 - Open problem: go beyond things that labelers can easily do
 - Verification is easier than generation

More on Reinforcement Learning

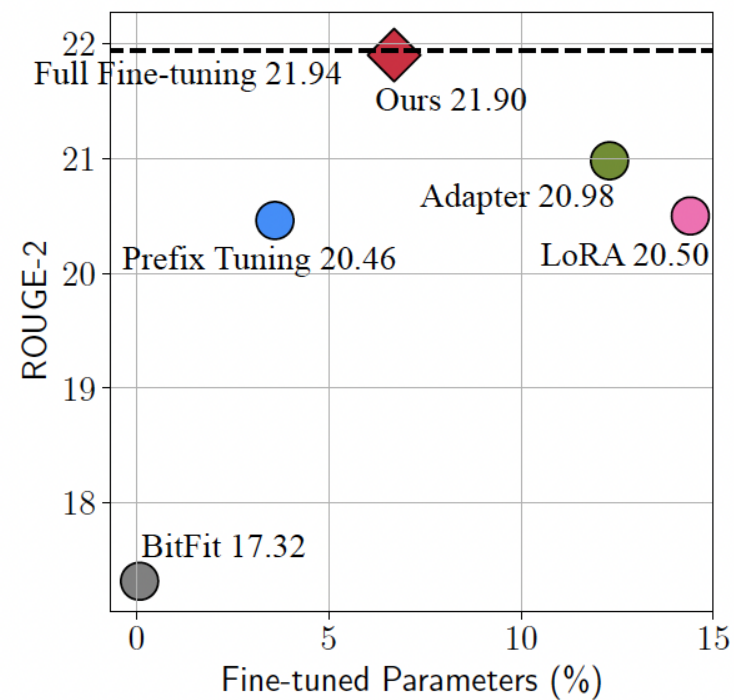
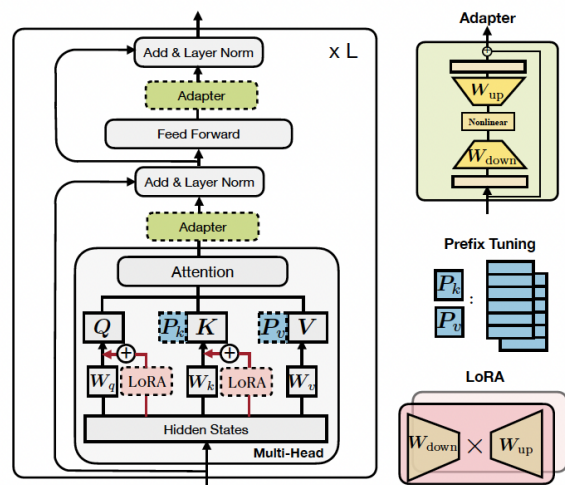
- Catalog of algorithms(PPO belongs)
 - World model or **model free**
 - Value-based or **policy-based(actor-critic)**
 - MC or **TD bootstrapping**
 - Off-policy or **on-policy**
 - Deterministic or **stochastic** policy
- Design consideration
 - **Sample efficiency**: Off-policy > On-policy
 - **Stability & Convergence**: Deadly Triad issue
 - **Explore & Exploit**: Random at episode beginning
 - **Bias & Variance**: Advantage Function



Parameter Efficient Finetuning

- Design aspects
- Primary implementations
 - Adapter
 - Prefix-tuning
 - LoRA

Overall



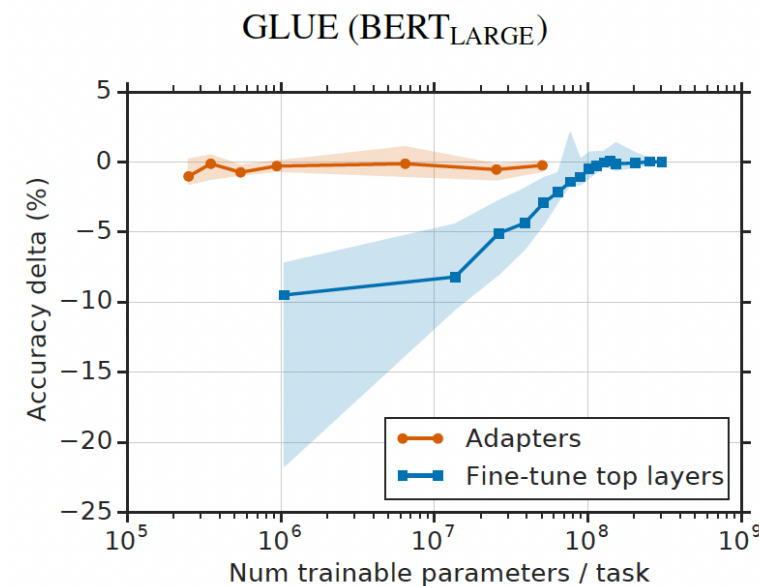
PEFT illustration and performance comparison (Source: [Junxian He, et.al](#))

Design aspects

- **Finetuned Modules:**
 - Attention-key/value matrix: LoRA(Q/V)
 - Attention-head: Prefix-Tuning(K/V)
 - Attention: Adapter
 - After FFN: Adapter
- **Other aspects:**
 - Multi-task consideration
 - Task related head

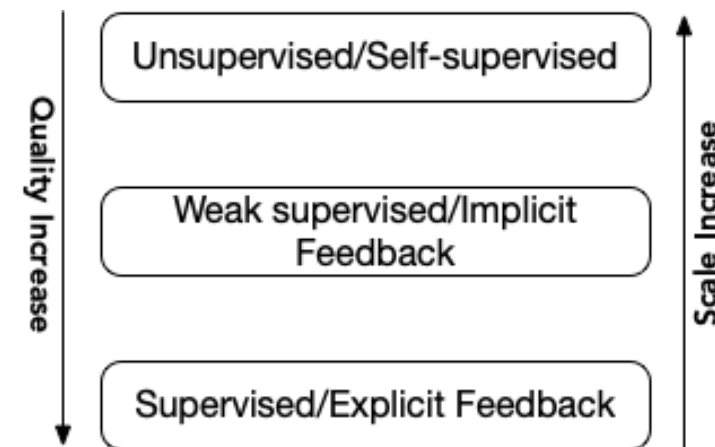
Adapter

- Implementation & training notes
 - $h \leftarrow h + f(hW_{\text{down}})W_{\text{up}}$
 - Parameter scale:
 $2 \times L \times (r \times d + r + d), r \ll d$
 - Adapt after FFN sub-layer works too
- Results
 - Finetune BERT for 26 classification Tasks
 - 3.6% parameters for 0.4% GLUE performance gap
 - Ablation: fewer layers adapted -> worser performance



Prefix-Tuning

- Implementation
 - $\text{head} = \text{Attn}(XW_Q^h, [P_K^h; XW_K^h], [P_V^h; XW_V^h])$
 - Reparameterization for finetuning stability:
 - $[P_K^h, P_V^h] = \text{MLP}_\theta(\hat{P}_K^h, \hat{P}_V^h), |\hat{P}| < |P|$
- Parameter scale:
 - $2 \times L \times d \times |P|$ or
 $2 \times L \times d \times \hat{P} + |\hat{P}| \times |P|$



More on Prefix-Tuning: Training and scaling

- Training
 - Initialization:
 - Real/high frequency words activation
 - *Task relevant words / Classification labels*
 - LM Head: Next Token/Class Label
- Results & discussion
 - Finetuning 0.1% \sim 3% parameters, comparable or better performance
 - Optimal prefix length varies: longer for more complex tasks
 - Reparameterization works task-dependently

PrefixTuning: Optimizing Continuous Prompts for Generation, Stanford, 2021

PromptTuning: The Power of Scale for Parameter-Efficient Prompt Tuning, Google, 2021

P-Tuning v2: Prompt Tuning Can Be Comparable to Fine-tuning Universally Across Scales and Tasks, Tsinghua, 2021

LoRA

- Implementation & training notes
 - Transformer: $W = W_0 + (BA)^T, A \in R^{r \times d}, B \in R^{d_h \times r}, r \ll \min\{d_h, d\}$
 - W_Q and W_V considered, parameter scale: $2 \times 2 \times d \times r \times L$
 - Modularized: Embedding, Linear, MergedLinear, Conv2D
 - Initialization: A kaiming-random, B zeros
 - Weight merged for inference efficiency
- Results
 - For 175B GPT-3 finetuning: 0.01% parameters, on par or better results
 - No additional inference computation and latency
 - Additivity for **finetuning merge** and **incremental update**

More on LoRA: which part to update and rank settings

| | # of Trainable Parameters = 18M | | | | | | |
|--------------------------|---------------------------------|------------|------------|------------|-----------------|-----------------|---------------------------|
| Weight Type Rank r | W_q 8 | W_k 8 | W_v 8 | W_o 8 | W_q, W_k 4 | W_q, W_v 4 | W_q, W_k, W_v, W_o 2 |
| WikiSQL ($\pm 0.5\%$) | 70.4 | 70.0 | 73.0 | 73.2 | 71.4 | 73.7 | 73.7 |
| MultiNLI ($\pm 0.1\%$) | 91.0 | 90.8 | 91.0 | 91.3 | 91.3 | 91.3 | 91.7 |

| | Weight Type | $r = 1$ | $r = 2$ | $r = 4$ | $r = 8$ | $r = 64$ |
|--------------------------|----------------------|---------|---------|---------|---------|----------|
| WikiSQL($\pm 0.5\%$) | W_q | 68.8 | 69.6 | 70.5 | 70.4 | 70.0 |
| | W_q, W_v | 73.4 | 73.3 | 73.7 | 73.8 | 73.5 |
| | W_q, W_k, W_v, W_o | 74.1 | 73.7 | 74.0 | 74.0 | 73.9 |
| MultiNLI ($\pm 0.1\%$) | W_q | 90.7 | 90.9 | 91.1 | 90.7 | 90.7 |
| | W_q, W_v | 91.3 | 91.4 | 91.3 | 91.6 | 91.4 |
| | W_q, W_k, W_v, W_o | 91.2 | 91.7 | 91.7 | 91.5 | 91.4 |

End of Parameter Efficient Finetuning

Steering the decoding process of LLM

- Decoding strategies
- Prompt engineering

Decoding strategies

- Temperature in decoding: $p_i = \frac{\exp(o_i/T)}{\sum_j \exp(o_j/T)}$
- Maximization search
 - Greedy search
 - Beam search
- Sampling
 - top-K sampling
 - top-p(Nucleus) sampling
 - Repetition penalized sampling
- Guided decoding
 - $score(x_{t+1}, b_t) = score(b_t) + \log p(x_{t+1}) + \sum_i \alpha_i f_i(x_{t+1})$

Prompting engineering

- Instruction/Zero-shot prompting
- Few-shot Prompting
- In-context Learning(Prompting)
- Chain-of-Thought

More on In-context Learning

Why it works?

- Interpretation from Topic Model
- Induction head
- View from gradient descent

Augmentation and Plugins

- Augmented Language Models
- Automatic prompting
- Plugins

Augmented Language Models

- Problems in vanilla LLMs
- Two augment aspects
 - Retrieval augmented language models
 - Tool augmented language models

Problems in vailla LLMs

Retrieval augmented language models

Tool augmented language models

End of Augmented Language Models

Automatic prompting

Plugins

- Plugins ecosystem
- Primary implementations
 - Langchain Agent/Tool
 - ChatGPT Plugins
 - Fixie
 - Other paradigm proposal

Plugins ecosystem

- Tool as a service
- From SEO to LMO
- Orchestration by LLM

Langchain Agent/Tool

ChatGPT Plugins

Fixie

Other paradigm proposal

End of Plugins

Thanks & QA?