

# Pre-trained Language Model with Limited Resources

---

Haoyu Wang<sup>♣</sup>, Yaqing Wang<sup>♡</sup>, Tianci Liu<sup>◊</sup>, Jing Gao<sup>◊</sup>

<sup>♣</sup>University at Albany, <sup>♡</sup>Google Deepmind, <sup>◊</sup>Purdue University

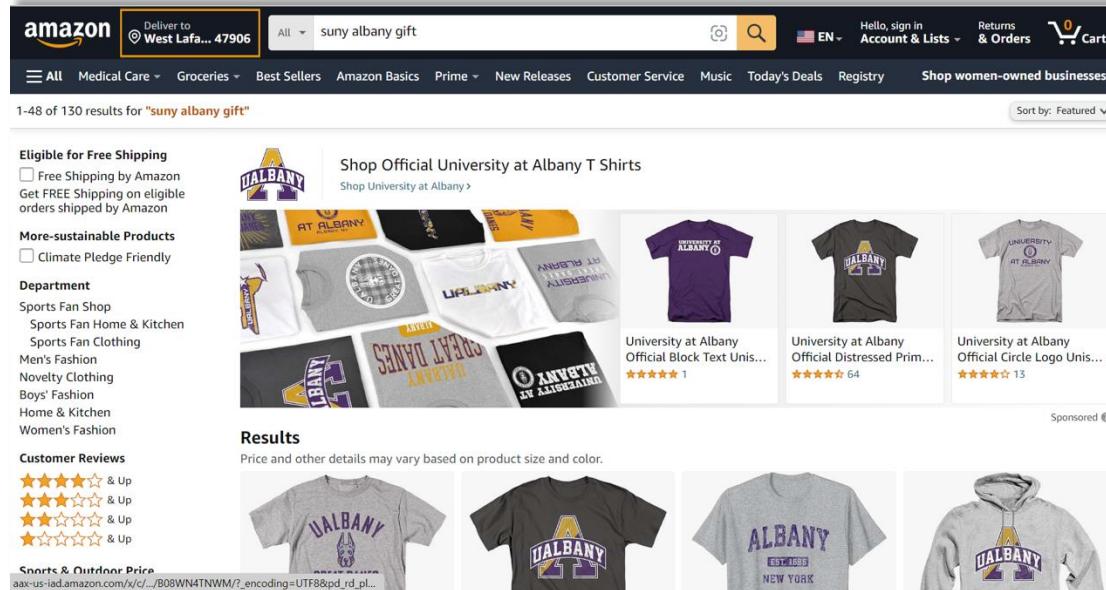


# Pre-trained Language Model Systems

The screenshot shows the Amazon search results page for the query "sunny albany gift". The search bar at the top has "sunny albany gift" typed into it. Below the search bar, there are various navigation links and filters. On the left side, there are sections for "Eligible for Free Shipping", "More-sustainable Products", "Department" (listing Sports Fan Shop, Men's Fashion, Novelty Clothing, Boys' Fashion, Home & Kitchen, and Women's Fashion), and "Customer Reviews" (with filters for 4★ & Up, 5★ & Up, 4★ & Up, and 5★ & Up). The main results area displays several items, including a collage of t-shirts with "ALBANY" and "AT ALBANY" logos, and three individual product cards for University at Albany t-shirts. The first card is for a purple t-shirt with "UNIVERSITY AT ALBANY" and "AT ALBANY" text, with a 4.5-star rating and 1 review. The second card is for a dark grey t-shirt with "UNIVERSITY AT ALBANY" and "AT ALBANY" text, with a 4.5-star rating and 64 reviews. The third card is for a grey t-shirt with "UNIVERSITY AT ALBANY" and "AT ALBANY" text, with a 4.5-star rating and 13 reviews. Below these cards, there is a section titled "Results" showing four more t-shirt options.

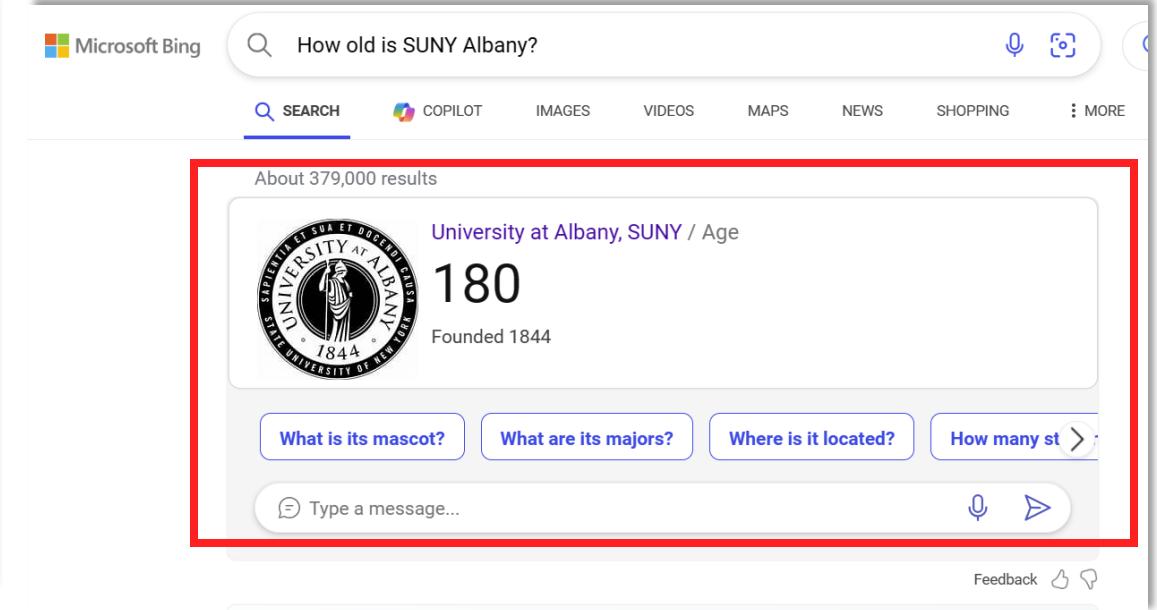
Amazon

# Pre-trained Language Model Systems



An Amazon search results page for "sunny albany gift". The search bar at the top shows the query. Below it, there are filters for "Eligible for Free Shipping" (unchecked), "More-sustainable Products" (unchecked), and "Department" (Sports Fan Shop). Under "Sports Fan Shop", categories like Sports Fan Home & Kitchen, Sports Fan Clothing, Men's Fashion, Novelty Clothing, Boys' Fashion, Home & Kitchen, and Women's Fashion are listed. The main results section shows several items, including a purple t-shirt with "ALBANY" and "AT ALBANY" on it, and a grey hoodie with "ALBANY" and "NEW YORK" on it. The results are sorted by "Featured".

Amazon



A Microsoft Bing search results page for "How old is SUNY Albany?". The search bar at the top contains the query. Below it, there are tabs for "SEARCH" (which is selected), "COPILOT", "IMAGES", "VIDEOS", "MAPS", "NEWS", and "SHOPPING". A red box highlights the search result for "University at Albany, SUNY / Age", which shows the university's logo, the number "180", and the text "Founded 1844". Below the result are four buttons: "What is its mascot?", "What are its majors?", "Where is it located?", and "How many st...". At the bottom of the search result card is a message input field with the placeholder "Type a message..." and a "Feedback" link.

Bing Search

# Pre-trained Language Model Systems

The image displays three search results side-by-side:

- Amazon Search Result:** Shows a search for "sunny albany gift". The results include items like "AT ALBANY" apparel. Filters on the left include "Eligible for Free Shipping" and "More-sustainable Products".

Rank	Product	Description	Category
1	AT ALBANY	Apparel	Sports Fan Shop
2	ALBANY	Apparel	Sports Fan Shop
3	SENY	Apparel	Sports Fan Shop
- U.S. Department of Labor - Office of Data Governance:** A screenshot of the official U.S. Department of Labor website. It features a blue header with "HOME", "DATA BOARD", "INVENTORY", "ASSESSMENT", "STRATEGY", and "DATA GOVERNANCE". Below the header, it says "Office of Data Governance (ODG) > Artificial Intelligence Use Case Inventory".

**Artificial Intelligence Use Case Inventory**

One of the key benefits to managing and utilizing data more efficiently is that it enables the use of advanced capabilities such as artificial intelligence (AI) and machine learning (ML). The Department of Labor has recently begun exploration of how these advanced technologies can be used to benefit the agency and help deliver on our mission. In an effort to create transparency in the adoption of these tools, this page serves to highlight the various uses of AI across the department.

Case ID	Case Type	Description	Owner	Model Type
15	Case Recording summarization	Using an open source large language model to summarize publicly available case recording documents which are void of personal identifiable information (PII) or any other sensitive information. This is not hosted in the DOL technical environment and is reviewed by human note takers.	Development and Acquisition	Large language summarization model
- Microsoft Bing Search Result:** A screenshot of a Microsoft Bing search for "How old is SUNY Albany?". The search bar shows the query. Below the search bar are buttons for "SEARCH", "COPilot", "IMAGES", "VIDEOS", "MAPS", "NEWS", "SHOPPING", and "MORE". A red box highlights the search bar and the "COPilot" button. The search results page includes a sidebar with questions like "What are its majors?", "Where is it located?", and "How many students does it have?".

U.S. Department of Labor

# Pre-trained Language Model Systems

The image displays three web pages side-by-side, each showing how pre-trained language models are used in real-world applications:

- Amazon Search:** A screenshot of the Amazon search results page for "sunny albany gift". The results include various items like shirts and hats with "SUNY ALBANY" or "ALBANY" branding.
- U.S. Department of Labor AI Use Case Inventory:** A screenshot of the Office of Data Governance (ODG) website under the "Artificial Intelligence Use Case Inventory". It features a table listing AI use cases across the agency. One row is highlighted:

15	Case Recording summarization	Using an open source large language model to summarize publicly available case recording documents which are void of personal identifiable information (PII) or any other sensitive information. This is not hosted in the DOL technical environment and is reviewed by human note takers.	Development and Acquisition	Large language summarization model
----	------------------------------	--	-----------------------------	------------------------------------
- Microsoft Bing Search:** A screenshot of the Microsoft Bing search results page for "How old is SUNY Albany?". The search bar shows the query, and the results page includes a Copilot sidebar with a red box highlighting a question about SUNY Albany's age.

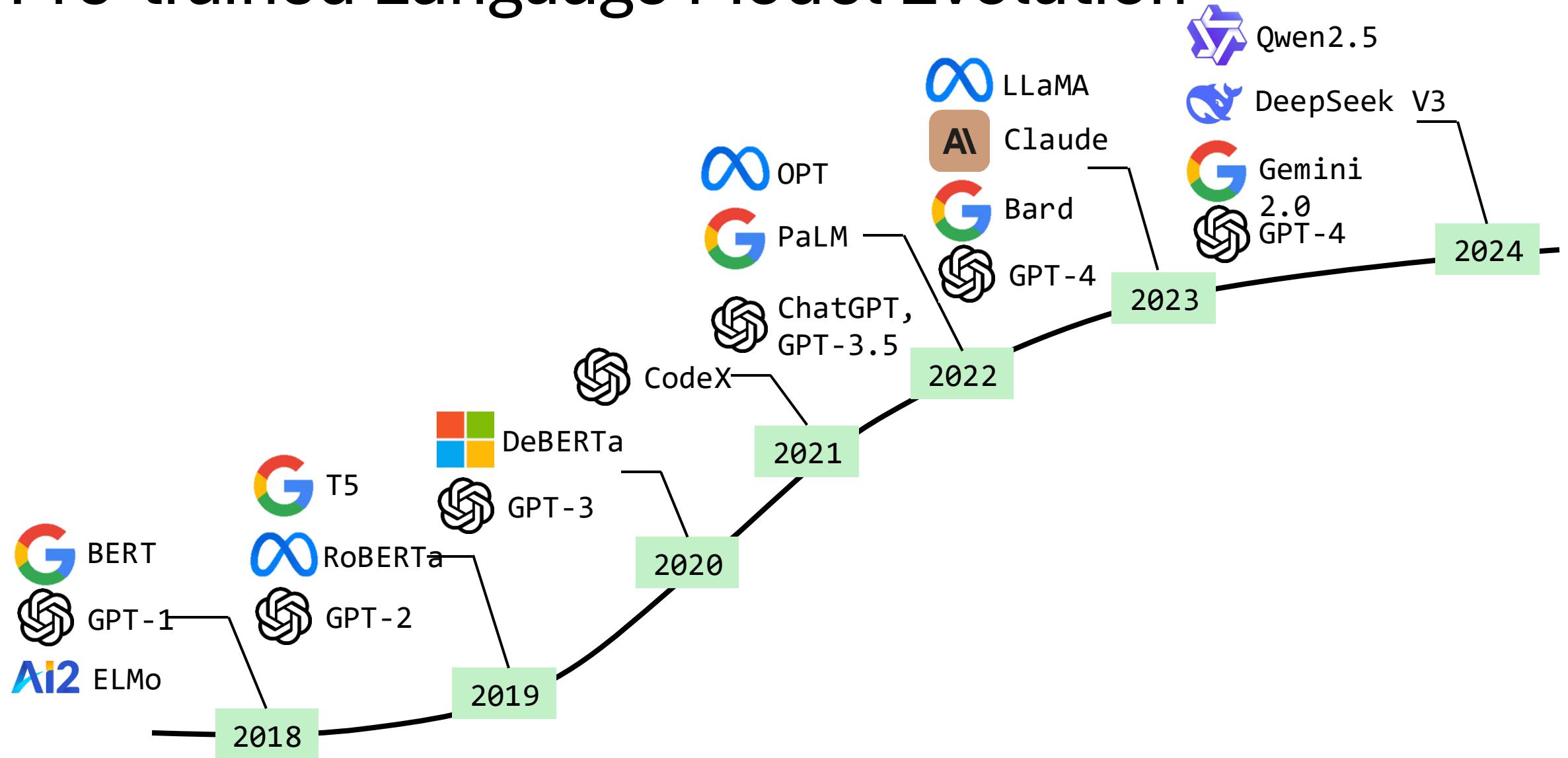
U.S. Department of Labor

Pre-trained Language Models have achieved remarkable performance in various real-world applications.

# Many other use cases:

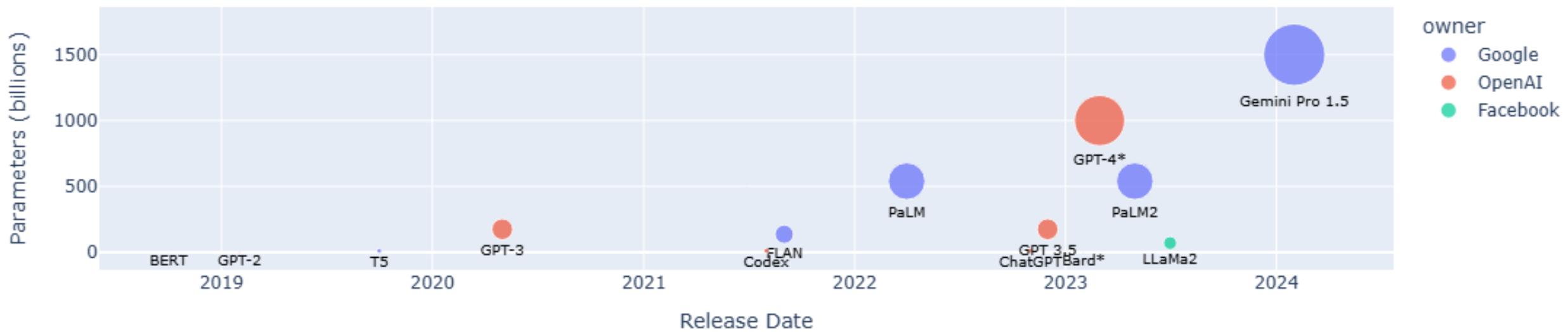
- Content generation
  - Writing outlines, descriptions, stories, summaries, scripts, emails.
- Natural language understanding
  - Grammar corrections, translations, sentiment analysis.
- Question answering and chatbots
  - Respond to queries and have a conversation.
- Analysis and optimization
  - Extract insights from text and workflow optimization.
- Software development
  - Code completion, debugging and explanation.

# Pre-trained Language Model Evolution



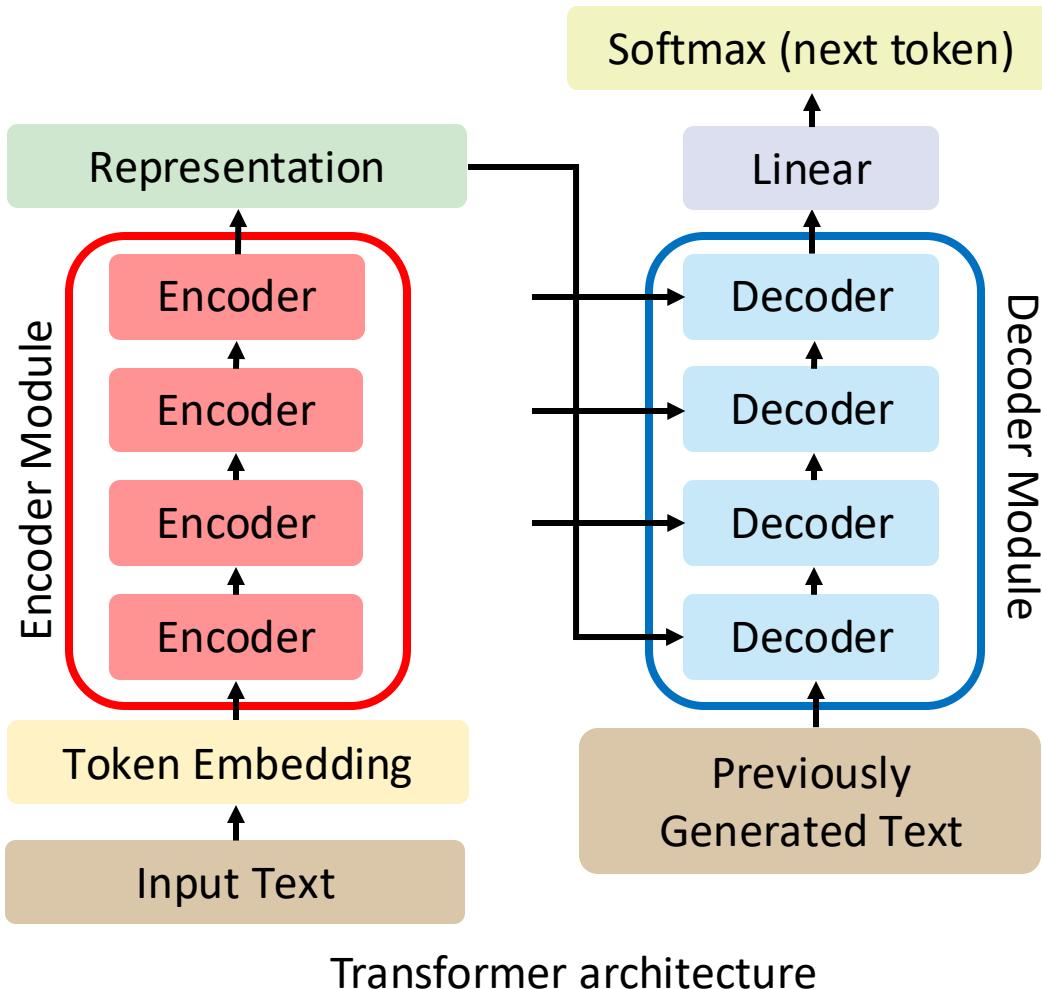
# Pre-trained Language Model Evolution

Evolution of Large Language Models



# Background of Pre-trained Language Models

- Architecture



## Encoder Module

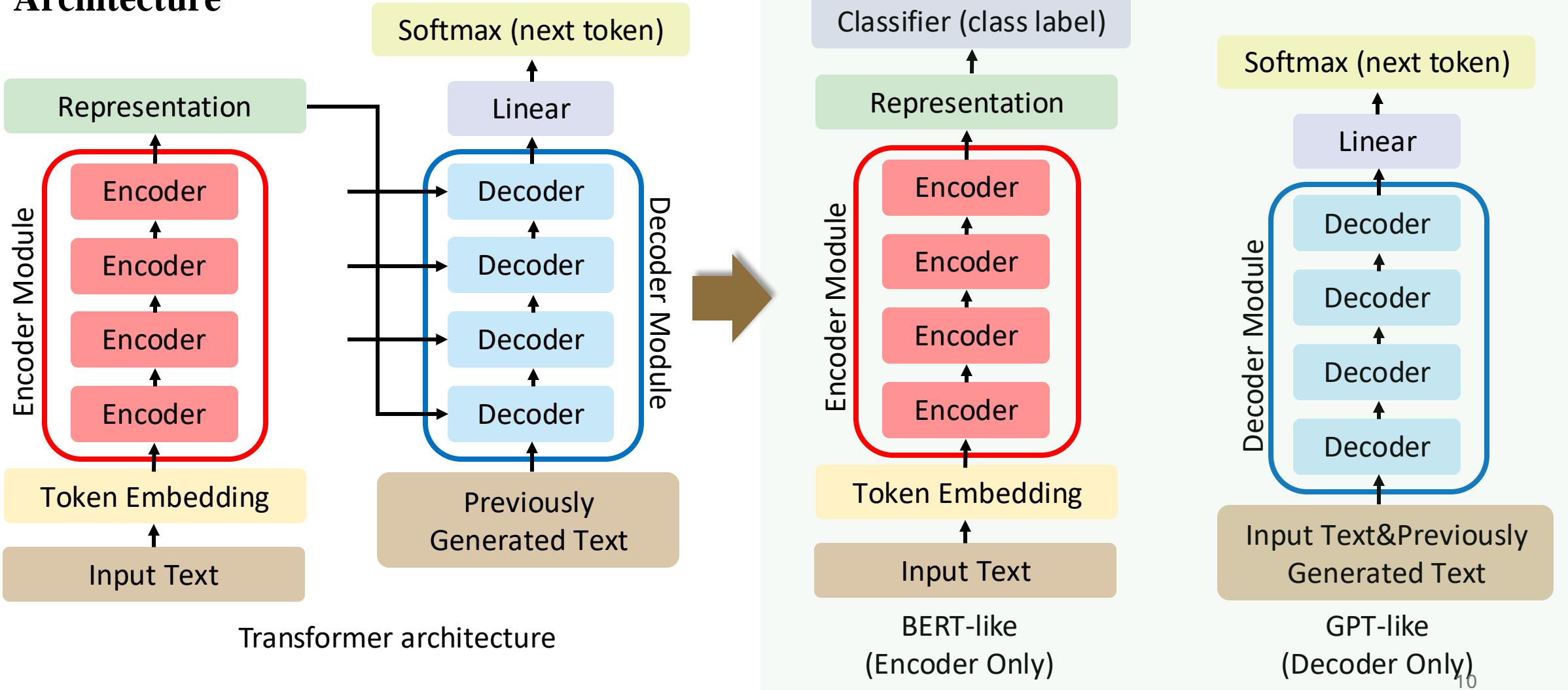
- Input:** a sequence of text
- Output:** representations corresponding to the tokens of the input text

## Decoder Module

- Input:** the output of the Encoder, the decoder module output at the previous time step
- Output:** next token probability (autoregressive)

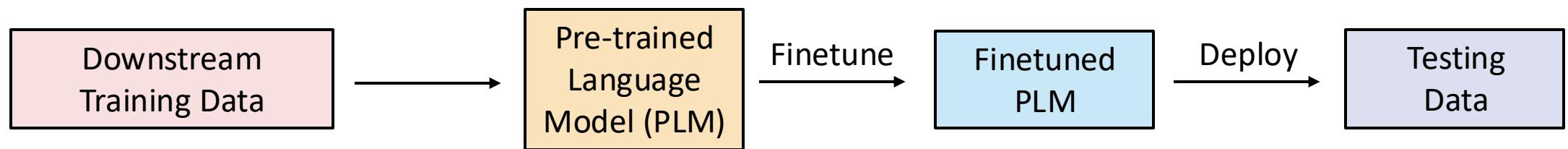
# Background of Pre-trained Language Models

## Architecture

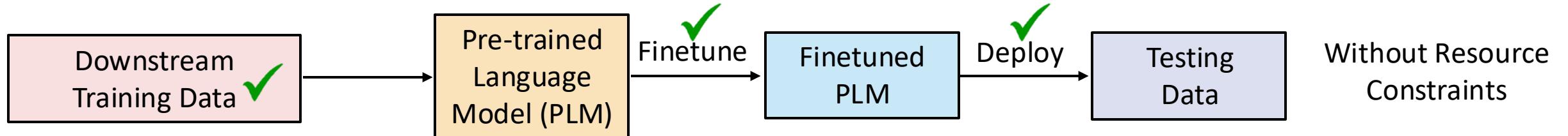


# Background of Pre-trained Language Models

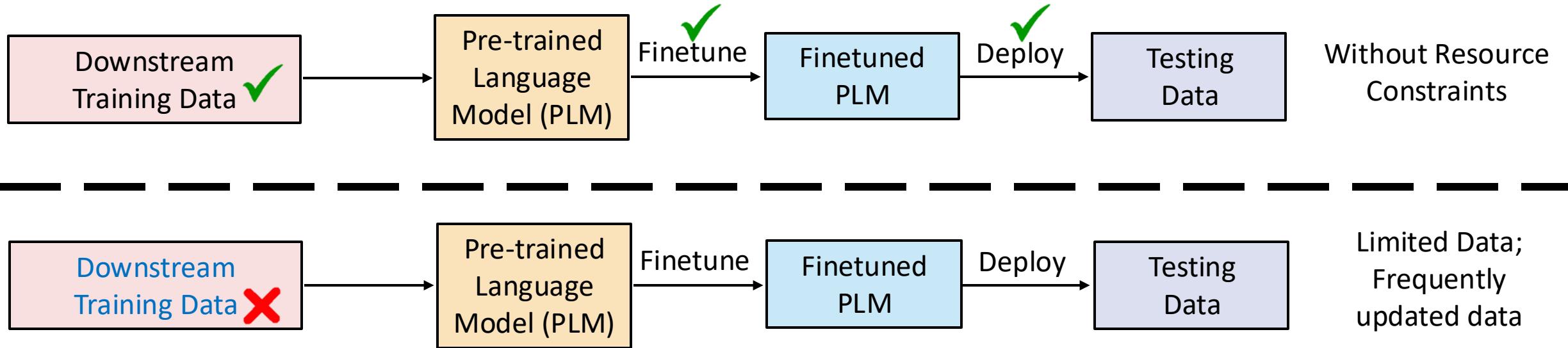
- Training paradigm



# Challenges From Resource Constraints

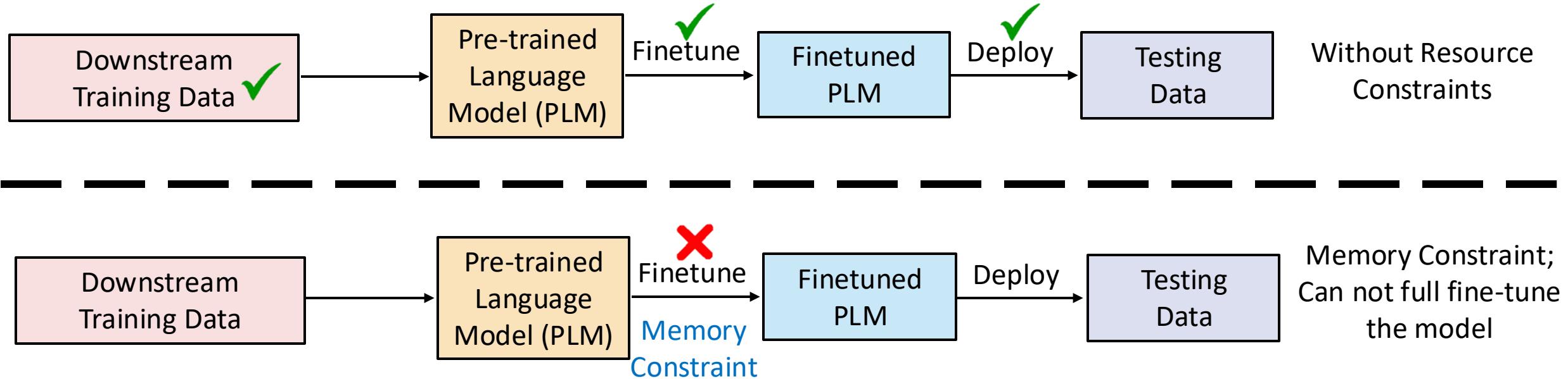


# Challenges From Resource Constraints

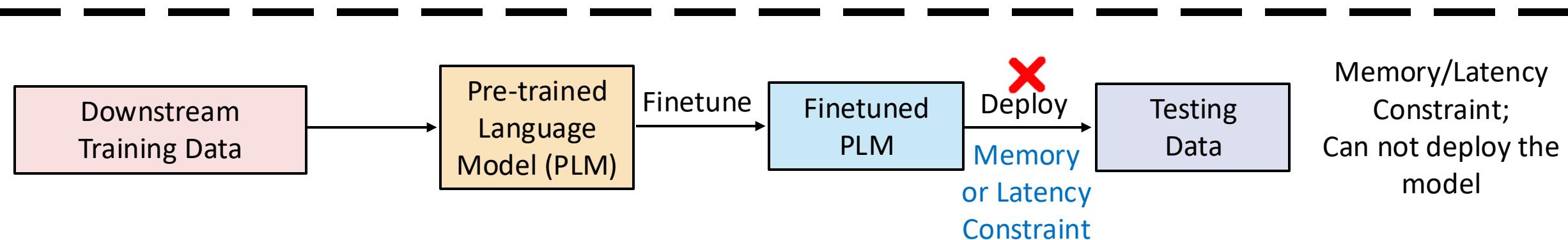
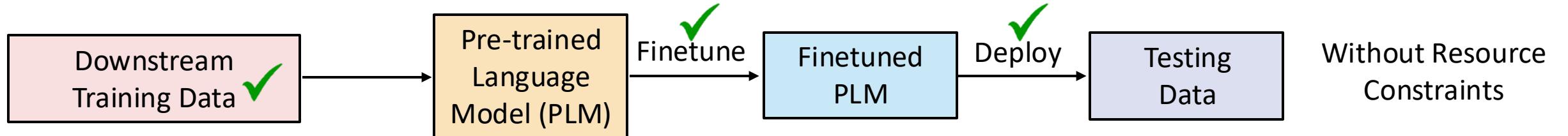


User privacy data;  
Data that are expensive  
to collect;  
Data that requires  
expert annotation

# Challenges From Resource Constraints



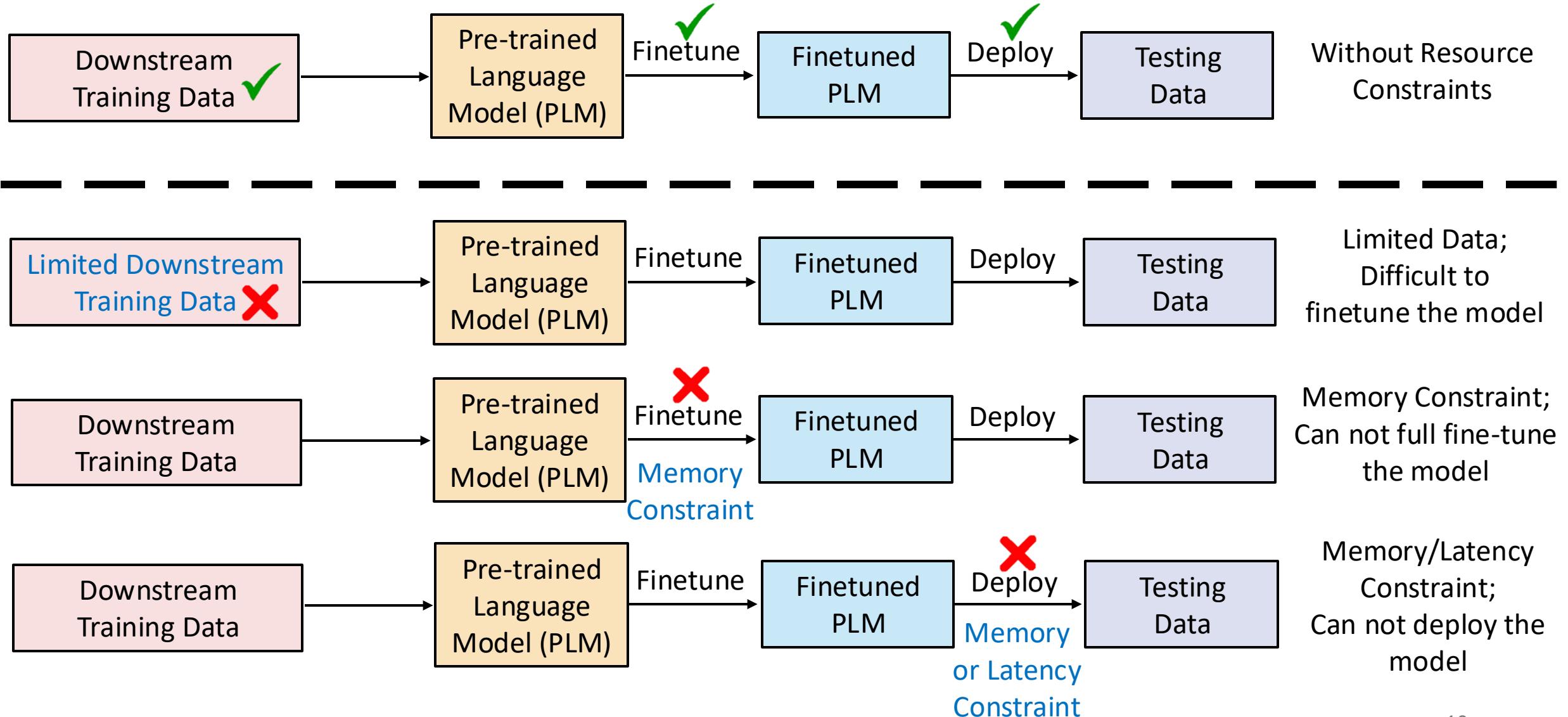
# Challenges From Resource Constraints



Edge Devices

Real-time Systems

# Challenges From Resource Constraints



# Outline

- Model Efficiency
  - Parameter-efficient Fine-tuning
  - Model Compression
  - Inference Acceleration
- Data Efficiency
  - Zero/few-shot Learning
  - Knowledge Editing
  - Retrieval-augmented Generation

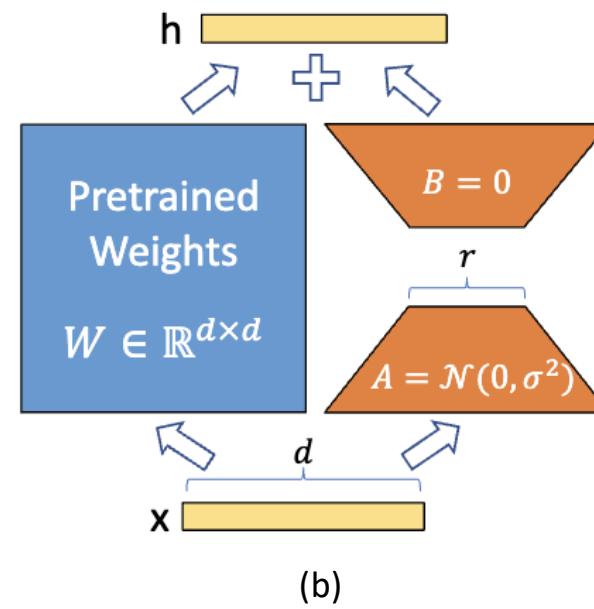
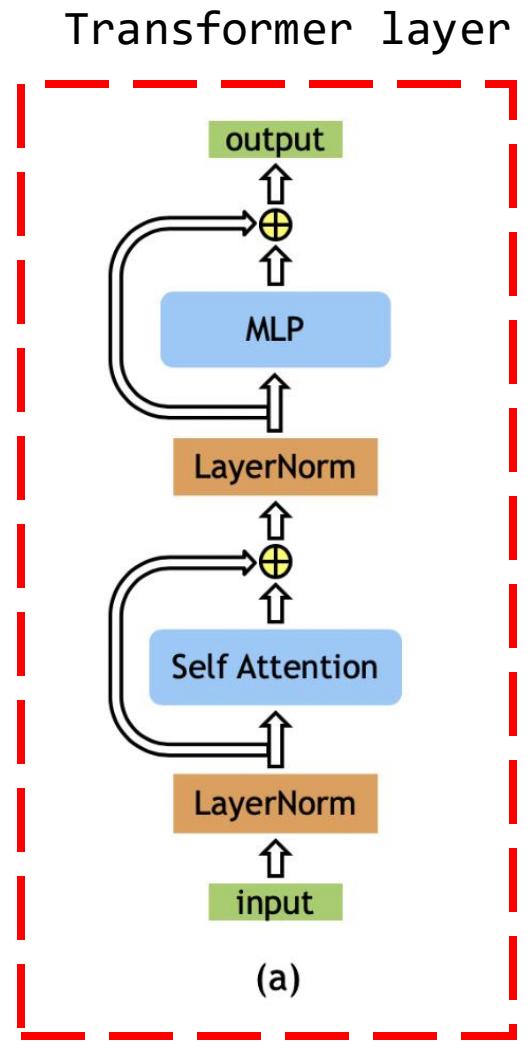
# Outline

- Model Efficiency
  - Parameter-efficient Fine-tuning
    - Model Compression
    - Inference Acceleration
- Data Efficiency
  - Zero/few-shot Learning
  - Knowledge Editing
  - Retrieval-augmented Generation

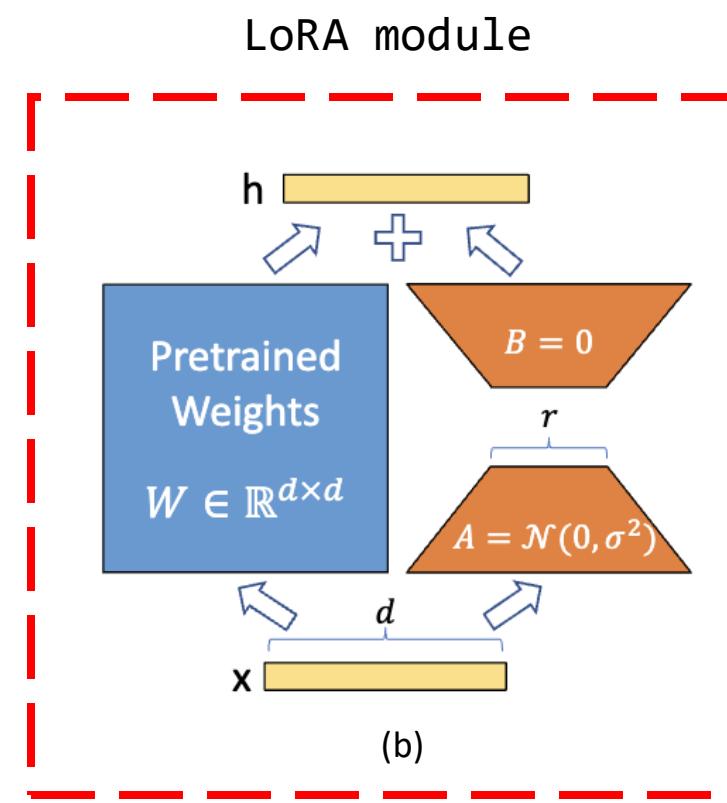
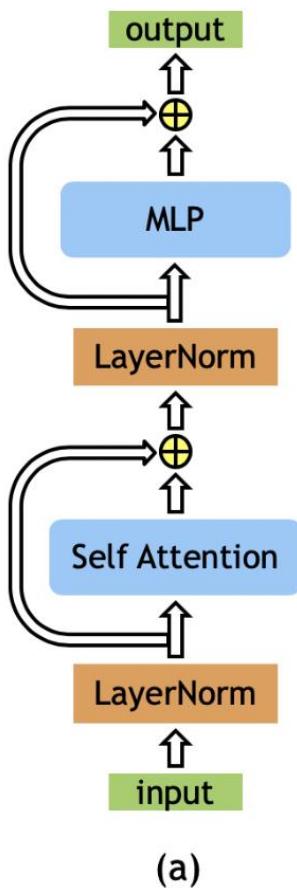
# Parameter-Efficient Fine-tuning

- Adapt pre-trained LLMs to downstream tasks without full fine-tuning.
  - Freeze pre-trained parameters.
  - Inject learnable **lightweight** modules only.

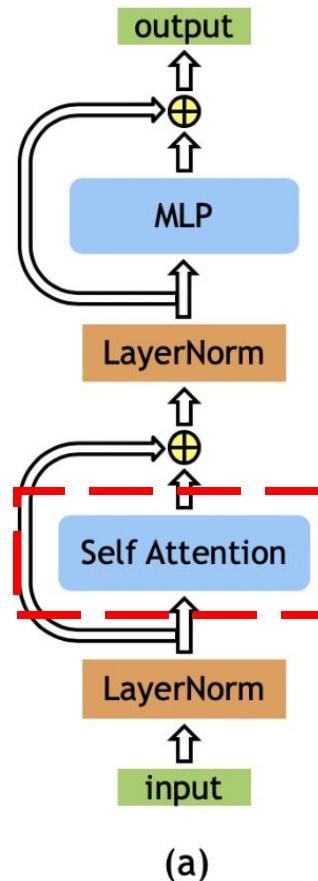
# LoRA



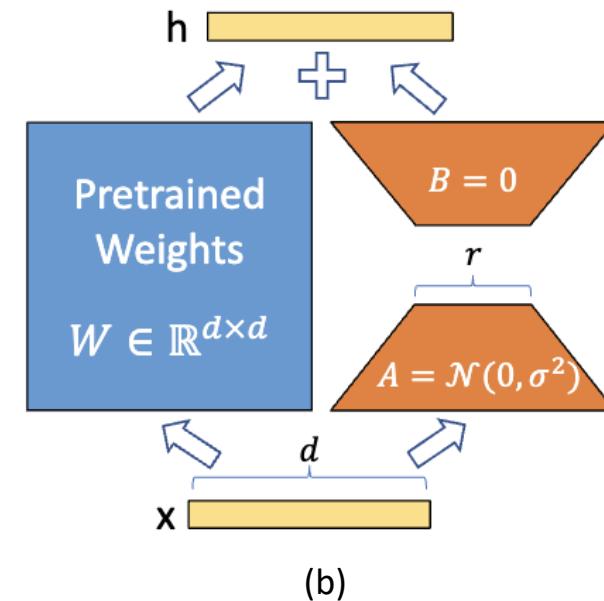
# LoRA



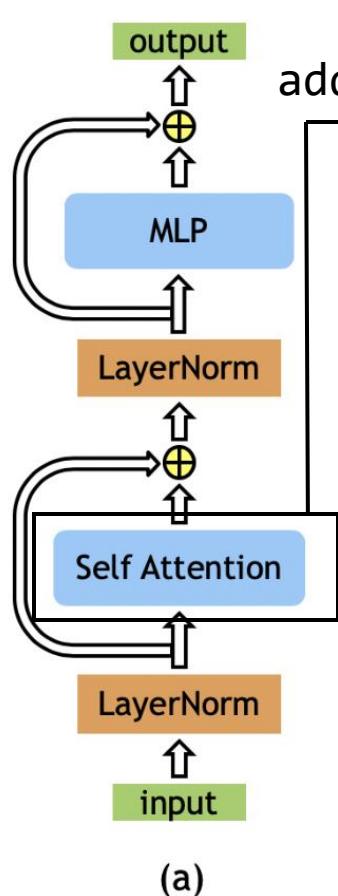
# LoRA



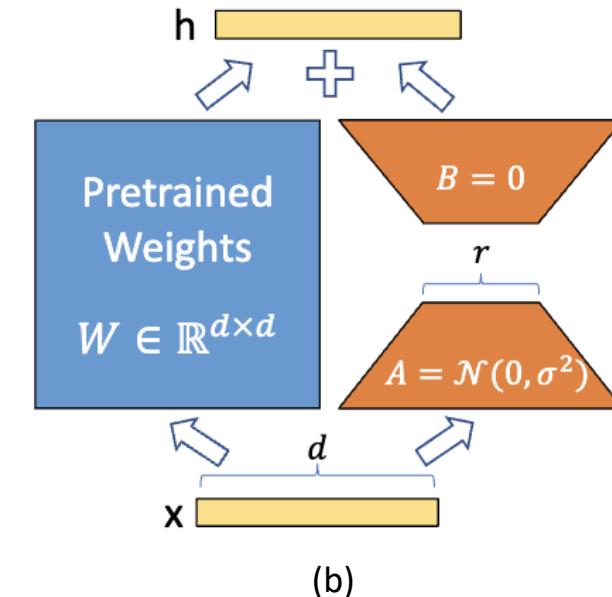
LoRA:  
add to Query and Value heads



# LoRA

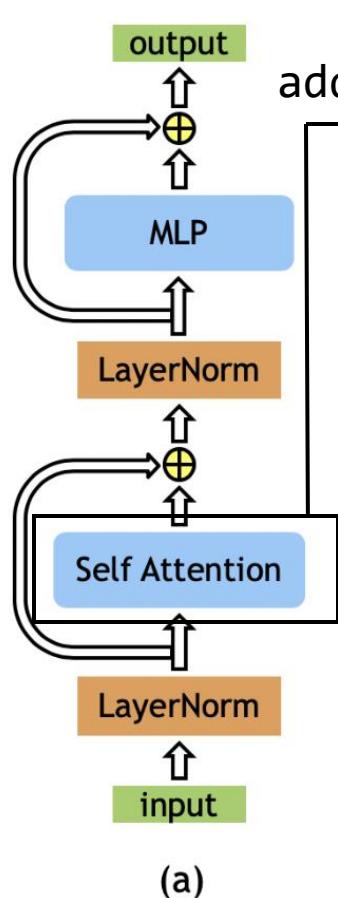


LoRA:  
add to Query and Value heads

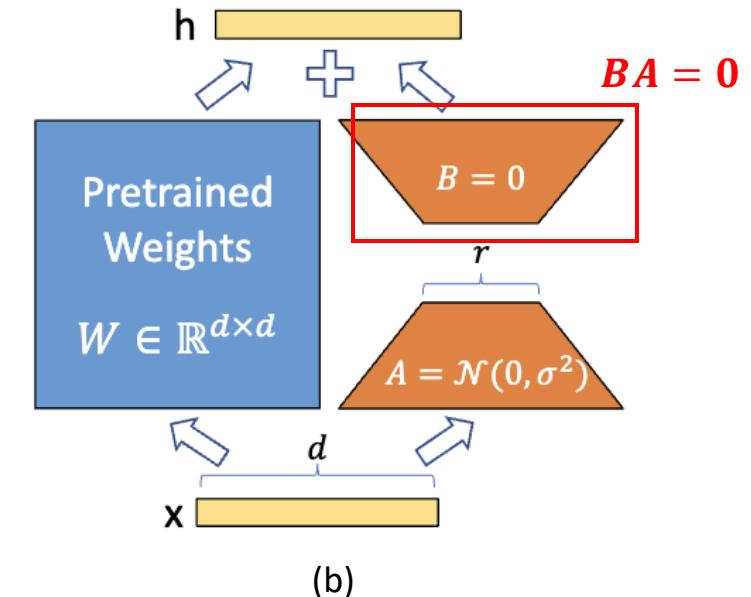


- $\mathbf{W} = \mathbf{W}^0 + \mathbf{B}\mathbf{A}, \mathbf{W}^0 \in \mathbb{R}^{d_1 \times d_2}, \mathbf{B} \in \mathbb{R}^{d_1 \times r}, \mathbf{A} \in \mathbb{R}^{r \times d_2}$
- $\min_{\mathbf{A}, \mathbf{B}} \mathcal{L}(\mathcal{D}; \mathbf{W}^0 + \mathbf{B}\mathbf{A})$  with  $\mathbf{W}^0$  fixed.
- Efficiency:  $r \ll \{d_1, d_2\} \Rightarrow$  less than 0.5% params.
- LoRA saves **backward** computation.

# LoRA



LoRA:  
add to Query and Value heads



- $\mathbf{W} = \mathbf{W}^0 + \mathbf{BA}, \mathbf{W}^0 \in \mathbb{R}^{d_1 \times d_2}, \mathbf{B} \in \mathbb{R}^{d_1 \times r}, \mathbf{A} \in \mathbb{R}^{r \times d_2}$
- $\min_{\mathbf{A}, \mathbf{B}} \mathcal{L}(\mathcal{D}; \mathbf{W}^0 + \mathbf{BA})$  with  $\mathbf{W}^0$  fixed.
- Efficiency:  $r \ll \{d_1, d_2\} \Rightarrow$  less than 0.5% params.
- LoRA saves **backward** computation.

# LoRA

Model&Method	# Trainable Parameters	WikiSQL	MNLI-m	SAMSum
		Acc. (%)	Acc. (%)	R1/R2/RL
GPT-3 (FT)	175,255.8M	<b>73.8</b>	89.5	52.0/28.0/44.5
GPT-3 (BitFit)	14.2M	71.3	91.0	51.3/27.4/43.5
GPT-3 (PreEmbed)	3.2M	63.1	88.6	48.3/24.2/40.5
GPT-3 (PreLayer)	20.2M	70.1	89.5	50.8/27.3/43.5
GPT-3 (Adapter <sup>H</sup> )	7.1M	71.9	89.8	53.0/28.9/44.8
GPT-3 (Adapter <sup>H</sup> )	40.1M	73.2	<b>91.5</b>	53.2/29.0/45.1
GPT-3 (LoRA)	4.7M	73.4	<b>91.7</b>	<b>53.8/29.8/45.9</b>
GPT-3 (LoRA)	37.7M	<b>74.0</b>	<b>91.6</b>	53.4/29.2/45.1

Model & Method	# Trainable Parameters	MNLI	SST-2	MRPC	CoLA	QNLI	QQP	RTE	STS-B	Avg.
RoB <sub>base</sub> (FT)*	125.0M	<b>87.6</b>	94.8	90.2	<b>63.6</b>	92.8	<b>91.9</b>	78.7	91.2	86.4
RoB <sub>base</sub> (BitFit)*	0.1M	84.7	93.7	<b>92.7</b>	62.0	91.8	84.0	81.5	90.8	85.2
RoB <sub>base</sub> (Adpt <sup>D</sup> )*	0.3M	87.1 <sub>±.0</sub>	94.2 <sub>±.1</sub>	88.5 <sub>±1.1</sub>	60.8 <sub>±.4</sub>	93.1 <sub>±.1</sub>	90.2 <sub>±.0</sub>	71.5 <sub>±2.7</sub>	89.7 <sub>±.3</sub>	84.4
RoB <sub>base</sub> (Adpt <sup>D</sup> )*	0.9M	87.3 <sub>±.1</sub>	94.7 <sub>±.3</sub>	88.4 <sub>±.1</sub>	62.6 <sub>±.9</sub>	93.0 <sub>±.2</sub>	90.6 <sub>±.0</sub>	75.9 <sub>±2.2</sub>	90.3 <sub>±.1</sub>	85.4
RoB <sub>base</sub> (LoRA)	0.3M	87.5 <sub>±.3</sub>	<b>95.1</b> <sub>±.2</sub>	89.7 <sub>±.7</sub>	63.4 <sub>±1.2</sub>	<b>93.3</b> <sub>±.3</sub>	90.8 <sub>±.1</sub>	<b>86.6</b> <sub>±.7</sub>	<b>91.5</b> <sub>±.2</sub>	<b>87.2</b>
RoB <sub>large</sub> (FT)*	355.0M	90.2	<b>96.4</b>	<b>90.9</b>	68.0	94.7	<b>92.2</b>	86.6	92.4	88.9
RoB <sub>large</sub> (LoRA)	0.8M	<b>90.6</b> <sub>±.2</sub>	96.2 <sub>±.5</sub>	<b>90.9</b> <sub>±1.2</sub>	<b>68.2</b> <sub>±1.9</sub>	<b>94.9</b> <sub>±.3</sub>	91.6 <sub>±.1</sub>	<b>87.4</b> <sub>±2.5</sub>	<b>92.6</b> <sub>±.2</sub>	<b>89.0</b>
RoB <sub>large</sub> (Adpt <sup>P</sup> )†	3.0M	90.2 <sub>±.3</sub>	96.1 <sub>±.3</sub>	90.2 <sub>±.7</sub>	<b>68.3</b> <sub>±1.0</sub>	<b>94.8</b> <sub>±.2</sub>	<b>91.9</b> <sub>±.1</sub>	83.8 <sub>±2.9</sub>	92.1 <sub>±.7</sub>	88.4
RoB <sub>large</sub> (Adpt <sup>P</sup> )†	0.8M	<b>90.5</b> <sub>±.3</sub>	<b>96.6</b> <sub>±.2</sub>	89.7 <sub>±1.2</sub>	67.8 <sub>±2.5</sub>	<b>94.8</b> <sub>±.3</sub>	91.7 <sub>±.2</sub>	80.1 <sub>±2.9</sub>	91.9 <sub>±.4</sub>	87.9
RoB <sub>large</sub> (Adpt <sup>H</sup> )†	6.0M	89.9 <sub>±.5</sub>	96.2 <sub>±.3</sub>	88.7 <sub>±2.9</sub>	66.5 <sub>±4.4</sub>	94.7 <sub>±.2</sub>	92.1 <sub>±.1</sub>	83.4 <sub>±1.1</sub>	91.0 <sub>±1.7</sub>	87.8
RoB <sub>large</sub> (Adpt <sup>H</sup> )†	0.8M	90.3 <sub>±.3</sub>	96.3 <sub>±.5</sub>	87.7 <sub>±1.7</sub>	66.3 <sub>±2.0</sub>	94.7 <sub>±.2</sub>	91.5 <sub>±.1</sub>	72.9 <sub>±2.9</sub>	91.5 <sub>±.5</sub>	86.4
RoB <sub>large</sub> (LoRA)†	0.8M	<b>90.6</b> <sub>±.2</sub>	96.2 <sub>±.5</sub>	<b>90.2</b> <sub>±1.0</sub>	68.2 <sub>±1.9</sub>	<b>94.8</b> <sub>±.3</sub>	91.6 <sub>±.2</sub>	<b>85.2</b> <sub>±1.1</sub>	<b>92.3</b> <sub>±.5</sub>	<b>88.6</b>
DeB <sub>XXL</sub> (FT)*	1500.0M	91.8	<b>97.2</b>	92.0	72.0	<b>96.0</b>	92.7	93.9	92.9	91.1
DeB <sub>XXL</sub> (LoRA)	4.7M	<b>91.9</b> <sub>±.2</sub>	96.9 <sub>±.2</sub>	<b>92.6</b> <sub>±.6</sub>	<b>72.4</b> <sub>±1.1</sub>	<b>96.0</b> <sub>±.1</sub>	<b>92.9</b> <sub>±.1</sub>	<b>94.9</b> <sub>±.4</sub>	<b>93.0</b> <sub>±.2</sub>	<b>91.3</b>

# LoRA

Model&Method	# Trainable Parameters	WikiSQL	MNLI-m	SAMSum
		Acc. (%)	Acc. (%)	R1/R2/RL
GPT-3 (FT)	175,255.8M	<b>73.8</b>	89.5	52.0/28.0/44.5
GPT-3 (BitFit)	14.2M	71.3	91.0	51.3/27.4/43.5
GPT-3 (PreEmbed)	3.2M	63.1	88.6	48.3/24.2/40.5
GPT-3 (PreLayer)	20.2M	70.1	89.5	50.8/27.3/43.5
GPT-3 (Adapter <sup>H</sup> )	7.1M	71.9	89.8	53.0/28.9/44.8
GPT-3 (Adapter <sup>H</sup> )	40.1M	73.2	<b>91.5</b>	53.2/29.0/45.1
GPT-3 (LoRA)	4.7M	73.4	<b>91.7</b>	<b>53.8/29.8/45.9</b>
GPT-3 (LoRA)	37.7M	<b>74.0</b>	<b>91.6</b>	53.4/29.2/45.1

Model & Method	# Trainable Parameters	MNLI	SST-2	MRPC	CoLA	QNLI	QQP	RTE	STS-B	Avg.
RoB <sub>base</sub> (FT)*	125.0M	<b>87.6</b>	94.8	90.2	<b>63.6</b>	92.8	<b>91.9</b>	78.7	91.2	86.4
RoB <sub>base</sub> (BitFit)*	0.1M	84.7	93.7	<b>92.7</b>	62.0	91.8	84.0	81.5	90.8	85.2
RoB <sub>base</sub> (Adpt <sup>D</sup> )*	0.3M	87.1 <sub>±.0</sub>	94.2 <sub>±.1</sub>	88.5 <sub>±1.1</sub>	60.8 <sub>±.4</sub>	93.1 <sub>±.1</sub>	90.2 <sub>±.0</sub>	71.5 <sub>±2.7</sub>	89.7 <sub>±.3</sub>	84.4
RoB <sub>base</sub> (Adpt <sup>D</sup> )*	0.9M	87.3 <sub>±.1</sub>	94.7 <sub>±.3</sub>	88.4 <sub>±.1</sub>	62.6 <sub>±.9</sub>	93.0 <sub>±.2</sub>	90.6 <sub>±.0</sub>	75.9 <sub>±2.2</sub>	90.3 <sub>±.1</sub>	85.4
RoB <sub>base</sub> (LoRA)	0.3M	87.5 <sub>±.3</sub>	<b>95.1<sub>±.2</sub></b>	89.7 <sub>±.7</sub>	63.4 <sub>±1.2</sub>	<b>93.3<sub>±.3</sub></b>	90.8 <sub>±.1</sub>	<b>86.6<sub>±.7</sub></b>	<b>91.5<sub>±.2</sub></b>	<b>87.2</b>
RoB <sub>large</sub> (FT)*	355.0M	90.2	<b>96.4</b>	<b>90.9</b>	68.0	94.7	<b>92.2</b>	86.6	92.4	88.9
RoB <sub>large</sub> (LoRA)	0.8M	<b>90.6<sub>±.2</sub></b>	96.2 <sub>±.5</sub>	<b>90.9<sub>±1.2</sub></b>	<b>68.2<sub>±1.9</sub></b>	<b>94.9<sub>±.3</sub></b>	91.6 <sub>±.1</sub>	<b>87.4<sub>±2.5</sub></b>	<b>92.6<sub>±.2</sub></b>	<b>89.0</b>
RoB <sub>large</sub> (Adpt <sup>P</sup> )†	3.0M	90.2 <sub>±.3</sub>	96.1 <sub>±.3</sub>	90.2 <sub>±.7</sub>	<b>68.3<sub>±1.0</sub></b>	<b>94.8<sub>±.2</sub></b>	<b>91.9<sub>±.1</sub></b>	83.8 <sub>±2.9</sub>	92.1 <sub>±.7</sub>	88.4
RoB <sub>large</sub> (Adpt <sup>P</sup> )†	0.8M	<b>90.5<sub>±.3</sub></b>	<b>96.6<sub>±.2</sub></b>	89.7 <sub>±1.2</sub>	67.8 <sub>±2.5</sub>	<b>94.8<sub>±.3</sub></b>	91.7 <sub>±.2</sub>	80.1 <sub>±2.9</sub>	91.9 <sub>±.4</sub>	87.9
RoB <sub>large</sub> (Adpt <sup>H</sup> )†	6.0M	89.9 <sub>±.5</sub>	96.2 <sub>±.3</sub>	88.7 <sub>±2.9</sub>	66.5 <sub>±4.4</sub>	94.7 <sub>±.2</sub>	92.1 <sub>±.1</sub>	83.4 <sub>±1.1</sub>	91.0 <sub>±1.7</sub>	87.8
RoB <sub>large</sub> (Adpt <sup>H</sup> )†	0.8M	90.3 <sub>±.3</sub>	96.3 <sub>±.5</sub>	87.7 <sub>±1.7</sub>	66.3 <sub>±2.0</sub>	94.7 <sub>±.2</sub>	91.5 <sub>±.1</sub>	72.9 <sub>±2.9</sub>	91.5 <sub>±.5</sub>	86.4
RoB <sub>large</sub> (LoRA)†	0.8M	<b>90.6<sub>±.2</sub></b>	96.2 <sub>±.5</sub>	<b>90.2<sub>±1.0</sub></b>	68.2 <sub>±1.9</sub>	<b>94.8<sub>±.3</sub></b>	91.6 <sub>±.2</sub>	<b>85.2<sub>±1.1</sub></b>	<b>92.3<sub>±.5</sub></b>	<b>88.6</b>
DeB <sub>XXL</sub> (FT)*	1500.0M	91.8	<b>97.2</b>	92.0	72.0	<b>96.0</b>	92.7	93.9	92.9	91.1
DeB <sub>XXL</sub> (LoRA)	4.7M	<b>91.9<sub>±.2</sub></b>	96.9 <sub>±.2</sub>	<b>92.6<sub>±.6</sub></b>	<b>72.4<sub>±1.1</sub></b>	<b>96.0<sub>±.1</sub></b>	<b>92.9<sub>±.1</sub></b>	<b>94.9<sub>±.4</sub></b>	<b>93.0<sub>±.2</sub></b>	<b>91.3</b>

On GPT-3 (175B)



Num. of Nvidia V100 32GB: 96 → 24



Model checkpoint Size: 1 TB → 200 MB

# AdaLoRA

## LoRA:

- $W = W^0 + BA$ ,  $A, B$  have rank  $r$ 
  - larger  $r$  makes better performance but more computationally expensive.
  - is a fixed hyper-param, ignoring layer importance.
- Suboptimal efficiency.

## AdaLoRA:

- $W = W^0 + \Delta$ , adjust the rank of  $\Delta$  to control the budget.

# AdaLoRA

**AdaLoRA:**

- $\mathbf{W} = \mathbf{W}^0 + \Delta$

By SVD:  $\Delta = \mathbf{P}\Lambda\mathbf{Q}$

- Diagonal  $\Lambda \in \mathbb{R}^{r \times r} \Rightarrow$  singular values  $\{\lambda_i\}_{1 \leq i \leq r}.$
- $\mathbf{P} \in \mathbb{R}^{d_1 \times r}, \mathbf{Q} \in \mathbb{R}^{r \times d_2} \Rightarrow$  left and right singular vectors.
- Regularize to enforce the orthogonality:

$$R(\mathbf{P}, \mathbf{Q}) = \| \mathbf{P}^\top \mathbf{P} - \mathbf{I} \|_F^2 + \| \mathbf{Q} \mathbf{Q}^\top - \mathbf{I} \|_F^2$$

# AdaLoRA

**AdaLoRA:**

- $\mathbf{W} = \mathbf{W}^0 + \Delta$

By SVD:  $\Delta = \mathbf{P}\Lambda\mathbf{Q}$

Iteratively prune to adjust the rank

- Diagonal  $\Lambda \in \mathbb{R}^{r \times r} \Rightarrow$  singular values  $\{\lambda_i\}_{1 \leq i \leq r}$ .
- $\mathbf{P} \in \mathbb{R}^{d_1 \times r}, \mathbf{Q} \in \mathbb{R}^{r \times d_2} \Rightarrow$  left and right singular vectors.
- Regularize to enforce the orthogonality:

$$R(\mathbf{P}, \mathbf{Q}) = \|\mathbf{P}^\top \mathbf{P} - \mathbf{I}\|_F^2 + \|\mathbf{Q} \mathbf{Q}^\top - \mathbf{I}\|_F^2$$

# AdaLoRA

**Rank Allocation:** Assume  $k \in [n]$  LoRA modules are learned together. Given budget  $b^{(t)}$  at step  $t$ ,

$$\tilde{\Lambda}_k^{(t)} = \Lambda_k^{(t)} - \eta \nabla_{\Lambda_k} \mathcal{L}, \quad \Lambda_k^{(t+1)} = \mathcal{T}(\tilde{\Lambda}_k^{(t)}, S_k^{(t)})$$

With

$$\mathcal{T}(\tilde{\Lambda}_k^{(t)}, S_k^{(t)})_{ii} = \begin{cases} \tilde{\Lambda}_{k,ii}^{(t)} & S_{k,i}^{(t)} \text{ is in the top-}b^{(t)}, \\ 0 & \text{otherwise,} \end{cases}$$

# AdaLoRA

**Rank Allocation:** Assume  $k \in [n]$  LoRA modules are learned together. Given budget  $b^{(t)}$  at step  $t$ ,

$$\tilde{\Lambda}_k^{(t)} = \Lambda_k^{(t)} - \eta \nabla_{\Lambda_k} \mathcal{L}, \quad \Lambda_k^{(t+1)} = \mathcal{T}(\tilde{\Lambda}_k^{(t)}, S_k^{(t)})$$

With

$$\mathcal{T}(\tilde{\Lambda}_k^{(t)}, S_k^{(t)})_{ii} = \begin{cases} \tilde{\Lambda}_{k,ii}^{(t)} & S_{k,i}^{(t)} \text{ is in the top-}b^{(t)}, \\ 0 & \text{otherwise,} \end{cases}$$

Importance of the k-th LoRA's i-th rank at step t.

# AdaLoRA

## Rank Allocation:

$$\tilde{\Lambda}_k^{(t)} = \Lambda_k^{(t)} - \eta \nabla_{\Lambda_k} \mathcal{L}, \quad \Lambda_k^{(t+1)} = \mathcal{T}(\tilde{\Lambda}_k^{(t)}, S_k^{(t)})$$

With

$$\mathcal{T}(\tilde{\Lambda}_k^{(t)}, S_k^{(t)})_{ii} = \begin{cases} \tilde{\Lambda}_{k,ii}^{(t)} & S_{k,i}^{(t)} \text{ is in the top-}b^{(t)}, \\ 0 & \text{otherwise,} \end{cases}$$

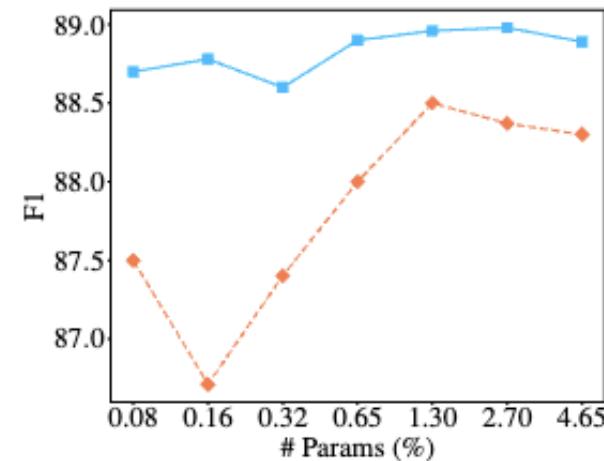
## Sensitivity-based Importance Score:

$$S_{k,i} = s(\lambda_{k,i}) + \frac{1}{d_1} \sum_{j=1}^{d_1} s(P_{k,ji}) + \frac{1}{d_2} \sum_{j=1}^{d_2} s(Q_{k,ij}),$$

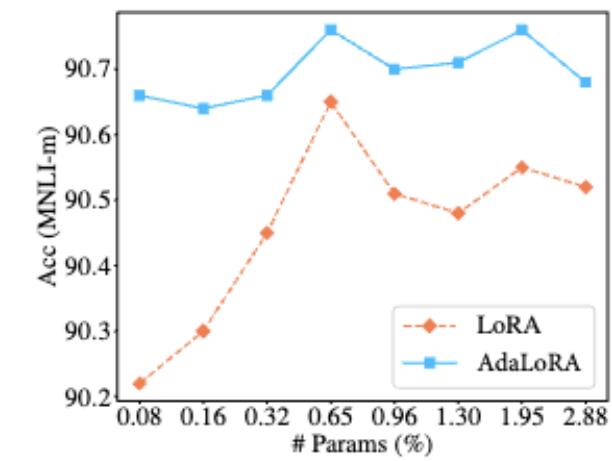
Sensitive score of single entry adjusted by uncertainty.

# AdaLoRA

Method	# Params	MNLI m/mm	SST-2 Acc	CoLA Mcc	QQP Acc/F1	QNLI Acc	RTE Acc	MRPC Acc	STS-B Corr	All Ave.
Full FT	184M	89.90/90.12	95.63	69.19	<b>92.40/89.80</b>	94.03	83.75	89.46	91.60	88.09
BitFit	0.1M	89.37/89.91	94.84	66.96	88.41/84.95	92.24	78.70	87.75	91.35	86.02
HAdapter	1.22M	90.13/90.17	95.53	68.64	91.91/89.27	94.11	84.48	89.95	91.48	88.12
PAdapter	1.18M	90.33/90.39	95.61	68.77	92.04/89.40	94.29	85.20	89.46	91.54	88.24
LoRA <sub>r=8</sub>	1.33M	90.65/90.69	94.95	69.82	91.99/89.38	93.87	85.20	89.95	91.60	88.34
AdaLoRA	1.27M	<b>90.76/90.79</b>	<b>96.10</b>	<b>71.45</b>	<b>92.23/89.74</b>	<b>94.55</b>	<b>88.09</b>	<b>90.69</b>	<b>91.84</b>	<b>89.31</b>
HAdapter	0.61M	90.12/90.23	95.30	67.87	91.65/88.95	93.76	85.56	89.22	91.30	87.93
PAdapter	0.60M	90.15/90.28	95.53	69.48	91.62/88.86	93.98	84.12	89.22	91.52	88.04
HAdapter	0.31M	90.10/90.02	95.41	67.65	91.54/88.81	93.52	83.39	89.25	91.31	87.60
PAdapter	0.30M	89.89/90.06	94.72	69.06	91.40/88.62	93.87	84.48	89.71	91.38	87.90
LoRA <sub>r=2</sub>	0.33M	90.30/90.38	94.95	68.71	91.61/88.91	94.03	85.56	89.71	<b>91.68</b>	88.15
AdaLoRA	0.32M	<b>90.66/90.70</b>	<b>95.80</b>	<b>70.04</b>	<b>91.78/89.16</b>	<b>94.49</b>	<b>87.36</b>	<b>90.44</b>	91.63	<b>88.86</b>



(b) SQuADv2.0



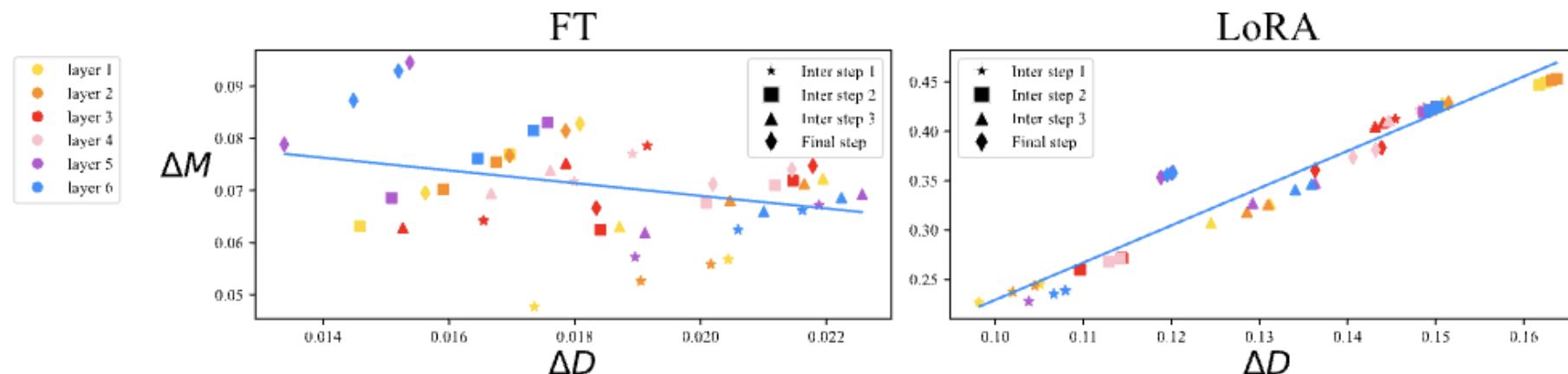
(a) MNLI

# DoRA

LoRA:

- $W = W^0 + BA$ 
  - Performance gap between LoRA and Full Fine-tuning (FFT).

Direction ( $D$ ) and Magnitude ( $M$ ) are of different importance in FFT that **LoRA fails to capture**.

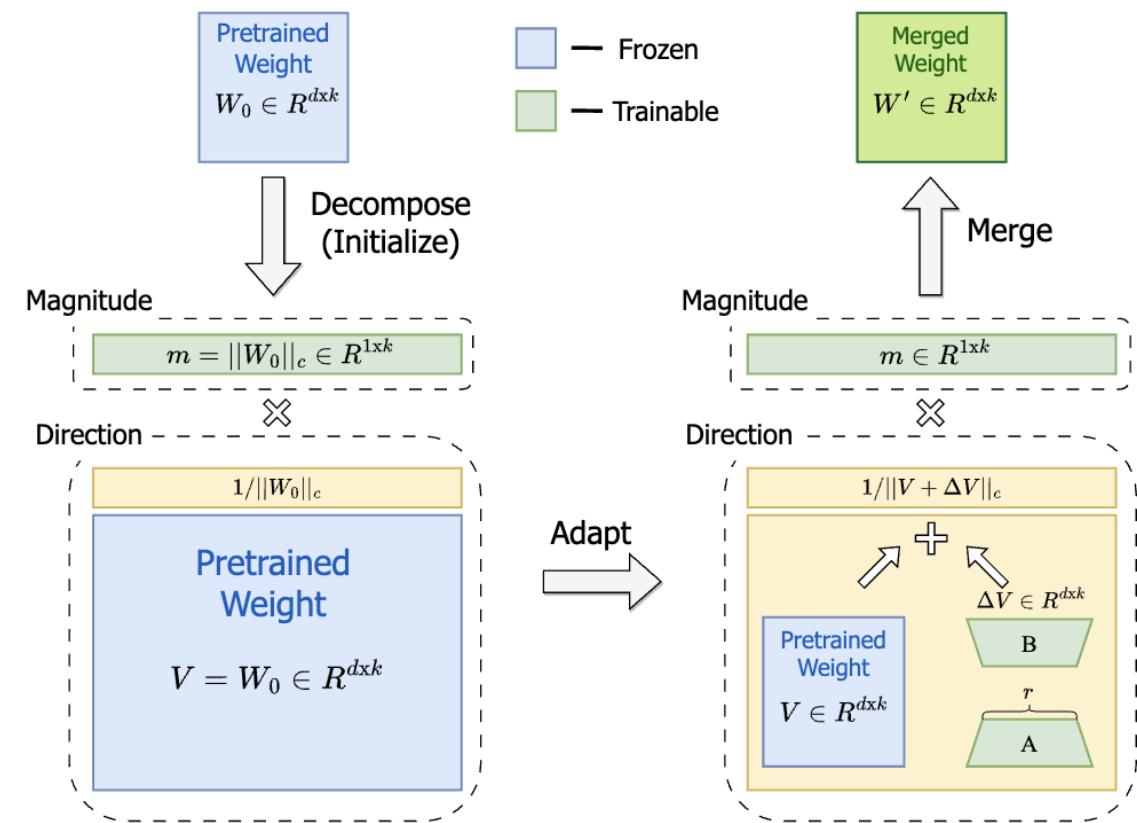


# DoRA

## DoRA: Weight-Decomposition

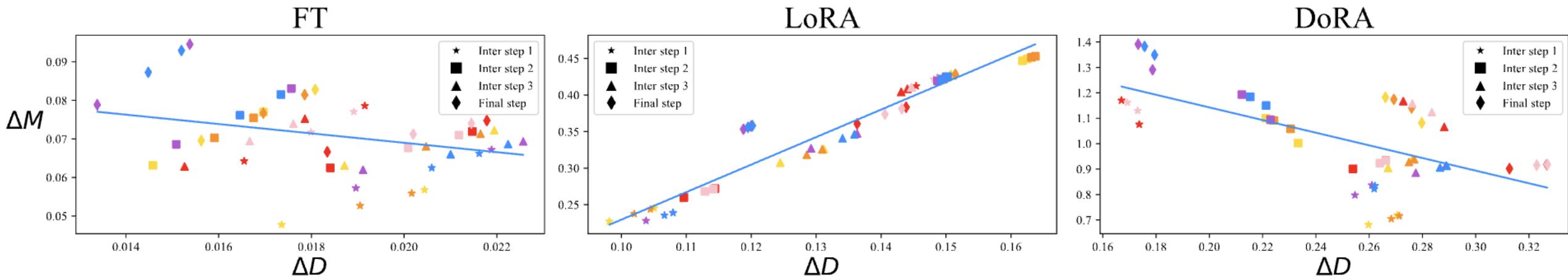
- LoRA  $\Rightarrow$  direction update
- + new vector  $\Rightarrow$  magnitude update  
 $\sim 0.01\%$  overhead

$$W' = \frac{m}{\|V + \Delta V\|_c} \frac{V + \Delta V}{\|V + \Delta V\|_c} = \frac{m}{\|W_0 + BA\|_c} \frac{W_0 + BA}{\|W_0 + BA\|_c}$$

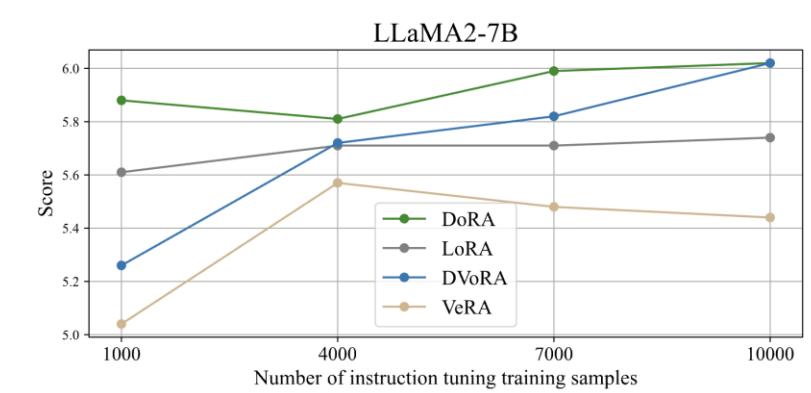
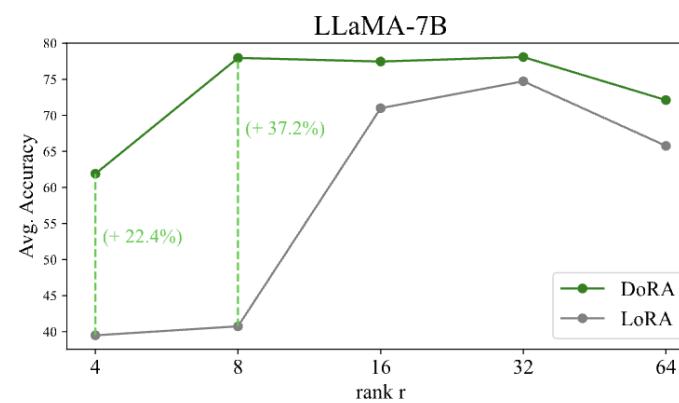
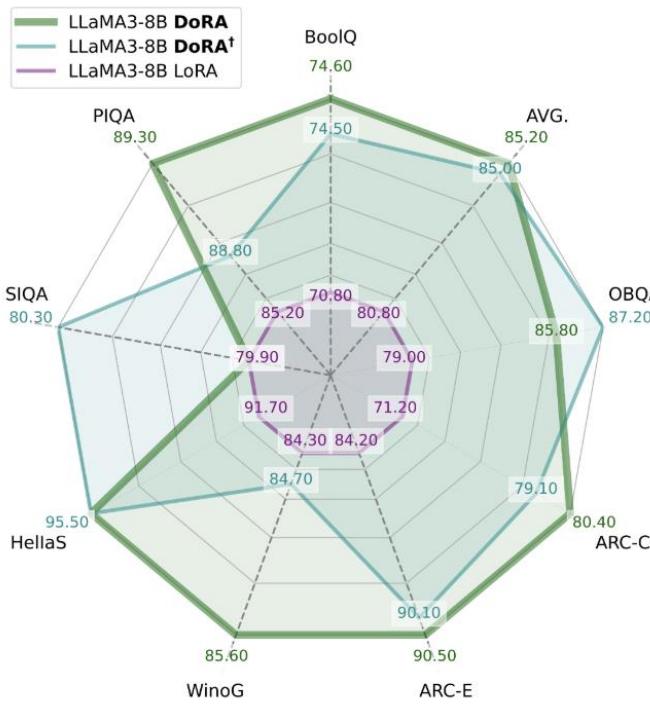


# DoRA

DoRA captures the magnitude and direction updates pattern in FFT.



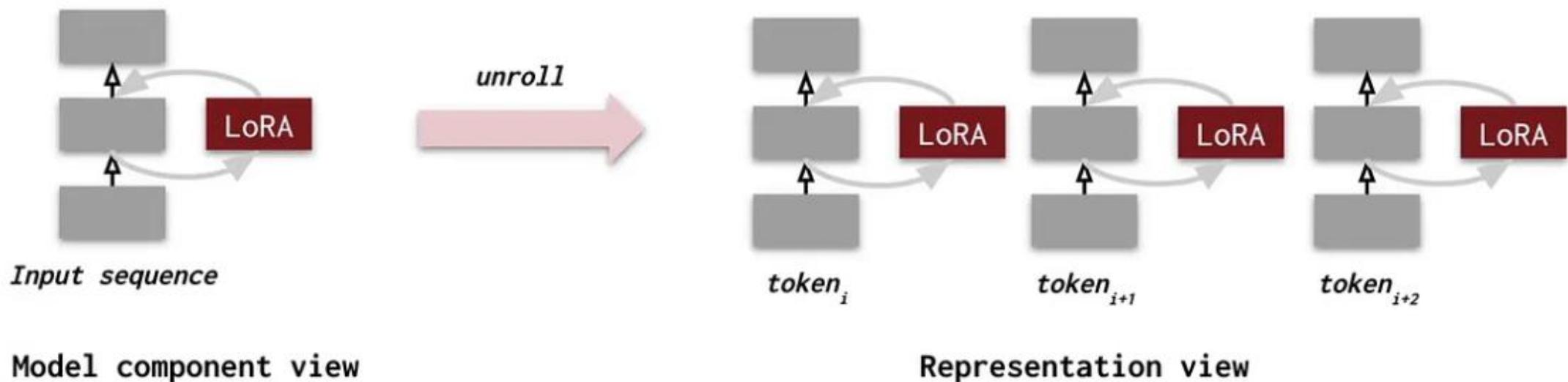
# DoRA



# ReFT

## LoRA:

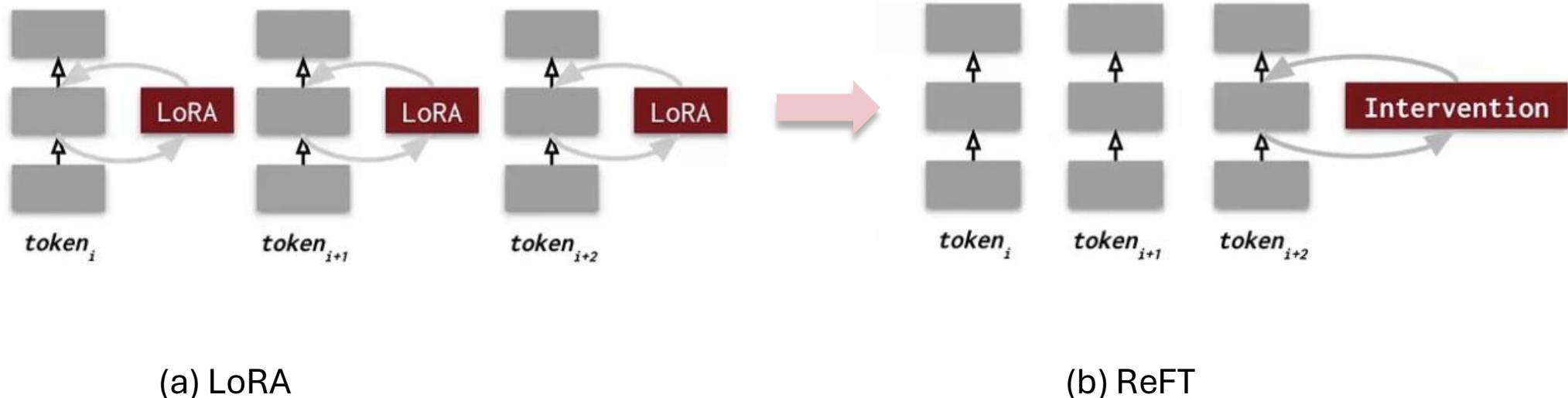
- Update model weights directly.
- Change representations at all timesteps.



# ReFT

ReFT:

- ~~Update model weights.~~ **Learn a set of interventions.**
- ~~Change representations across all timesteps.~~ **At a few timesteps.**

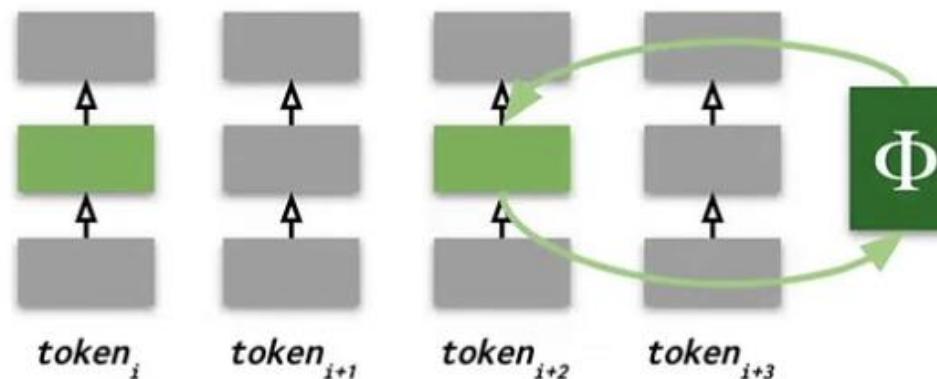


# ReFT

**ReFT:**

- Select **layer  $l$**  and **timestep  $i$**  to apply the intervention.
- **Intervention** is a low-rank and linear update:

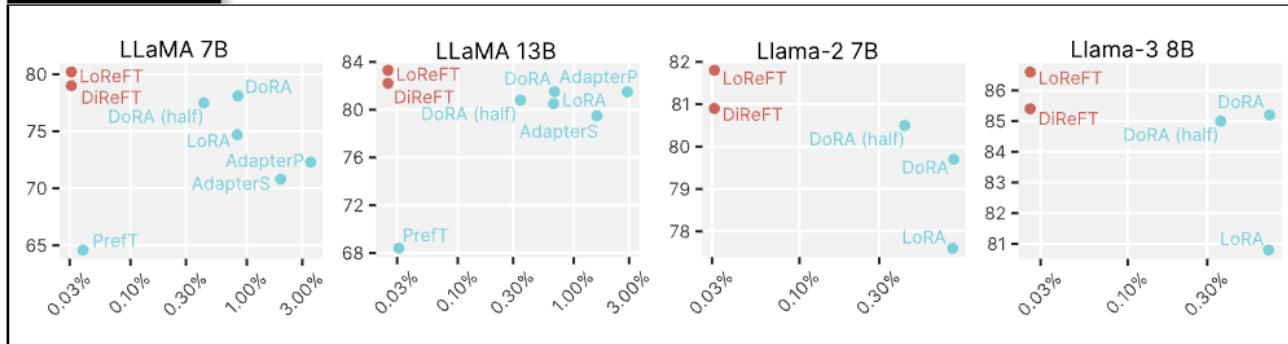
$$\Phi_l(\mathbf{h}_i^{(l)}; \phi_l) = \mathbf{h}_i^{(l)} + \mathbf{R}_l^\top (\mathbf{A}_l \mathbf{h}_i^{(l)} + \mathbf{b}_l - \mathbf{R}_l \mathbf{h}_i^{(l)}),$$



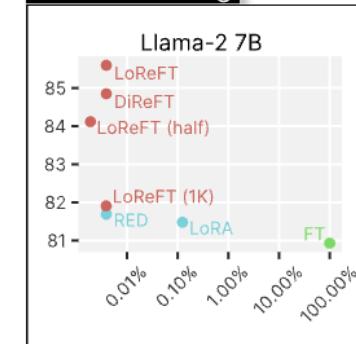
# ReFT

Efficiency: 0.8% (LoRA)  $\Rightarrow$  0.03% (ReFT)

Commonsense

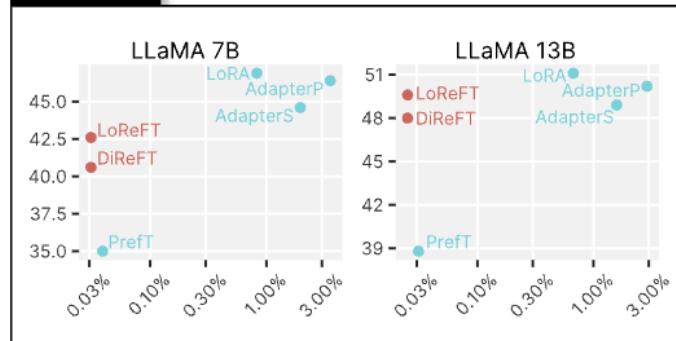


Instruct-tuning

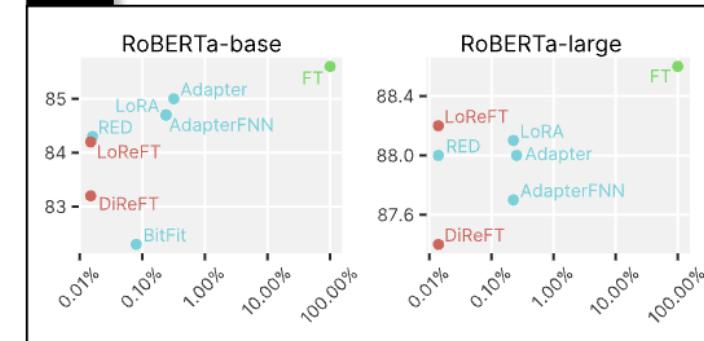


Arithmetic

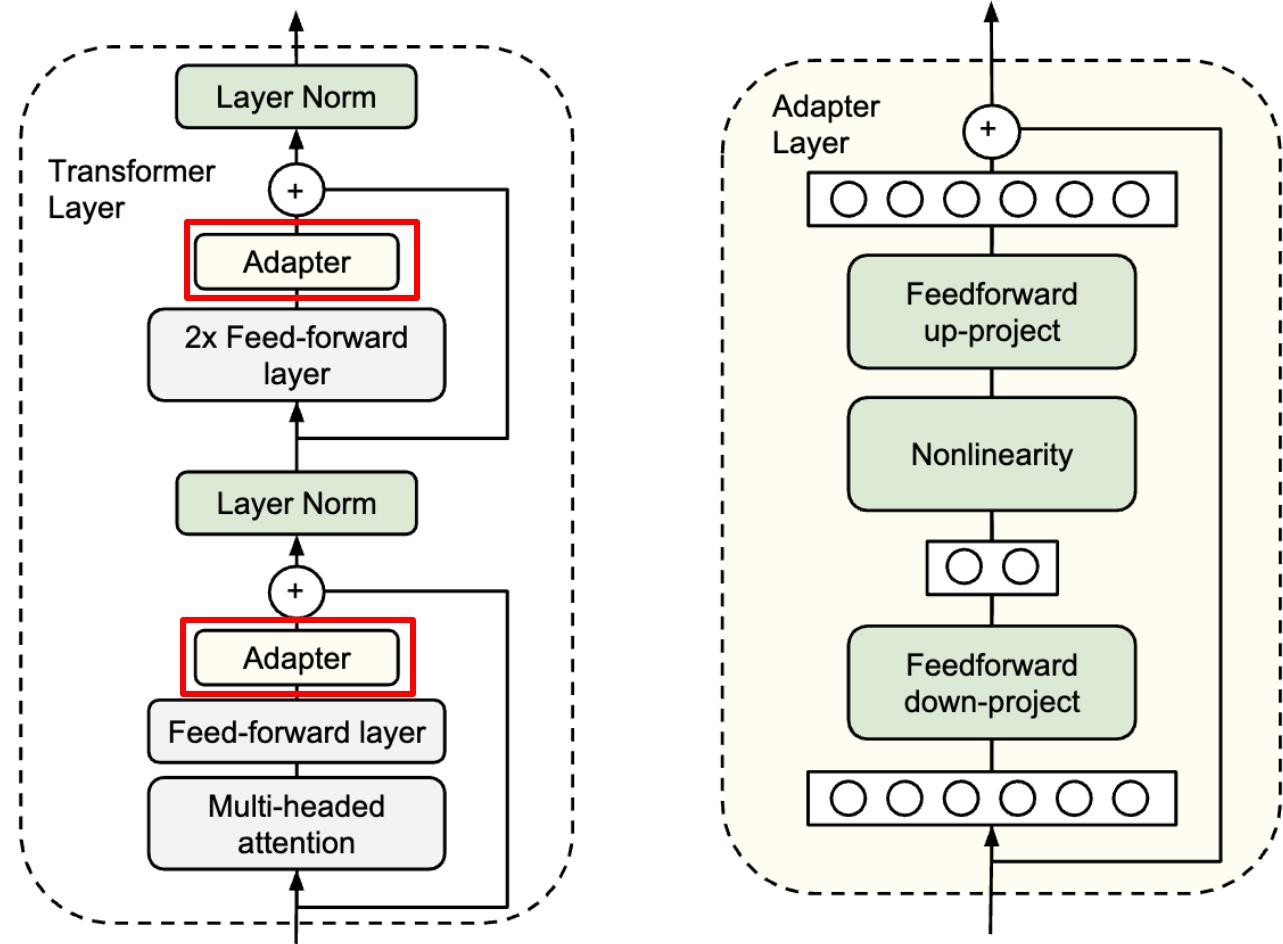
Performance  
↑  
Parameters  
→



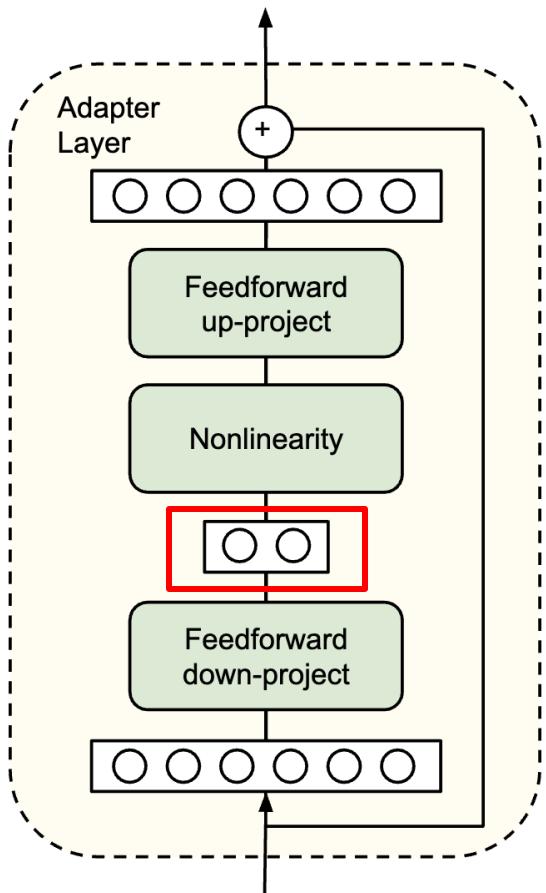
GLUE



# Adapter

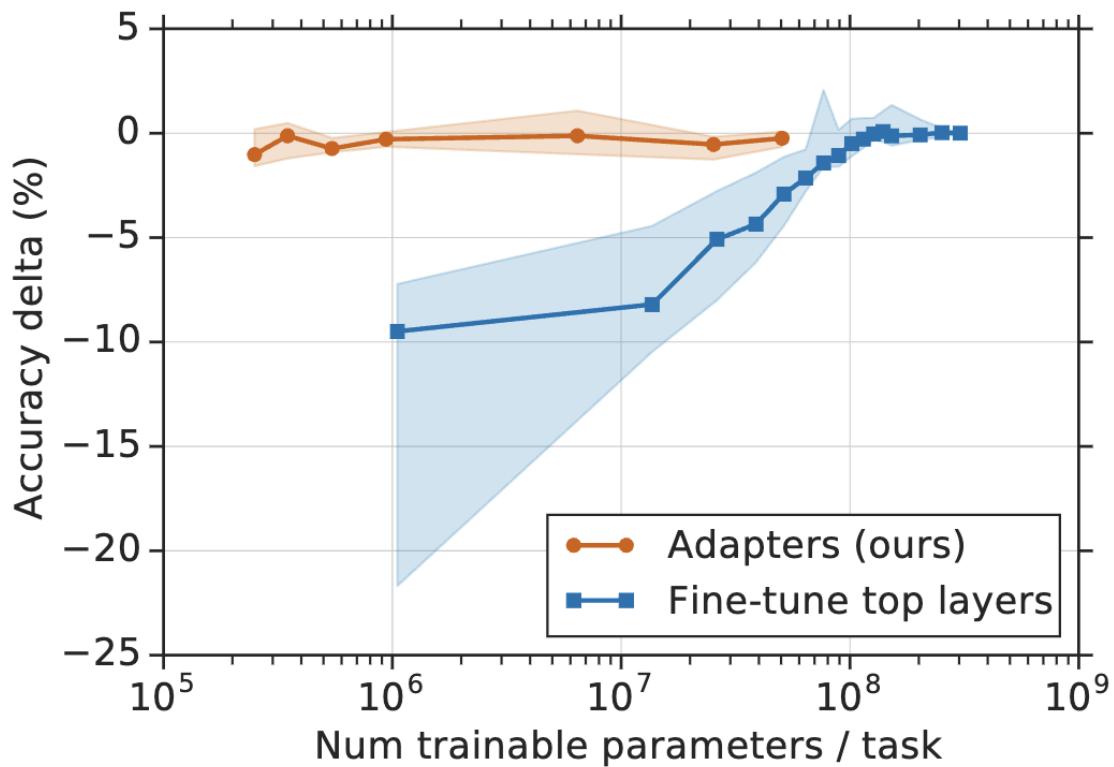


# Adapter



Bottleneck to save #  
params

# Adapter

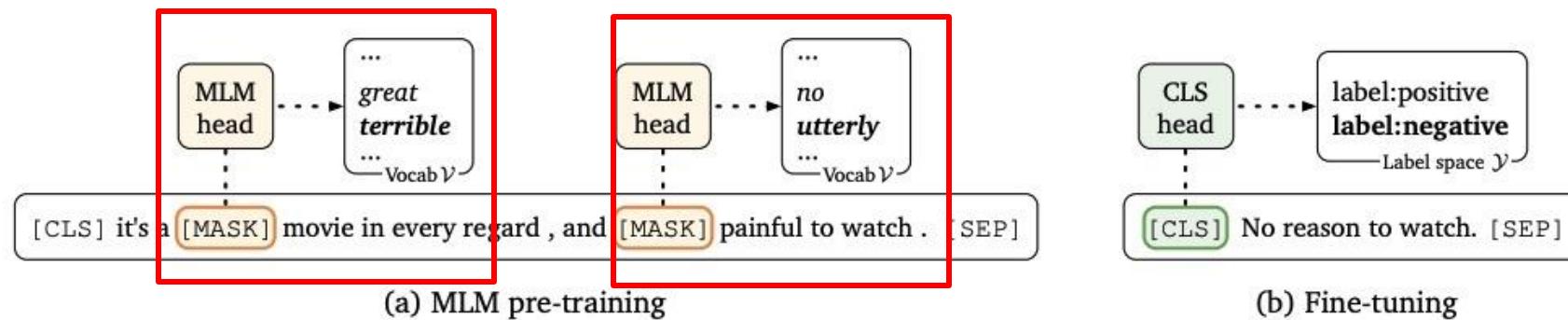


0.4% accuracy drop for 96.4%  
reduction in the #of  
parameters/task

# LiST

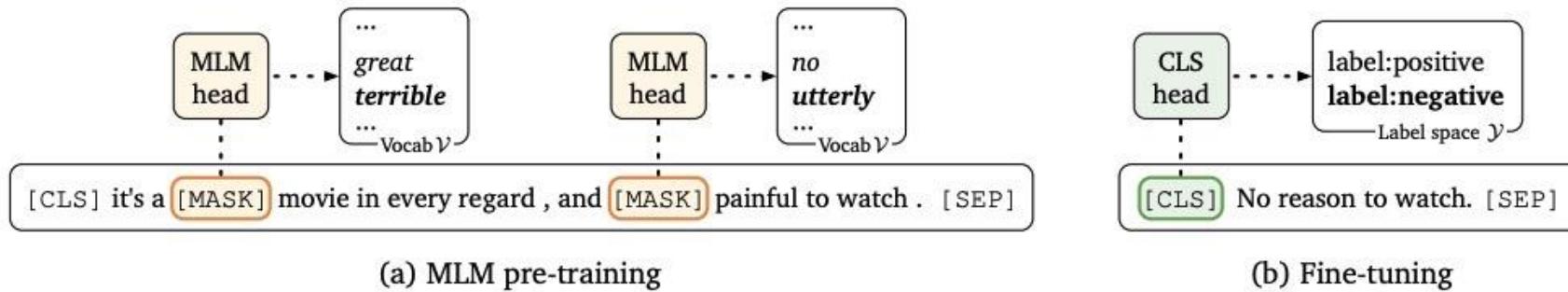
## Encoder-based Model:

- Pre-trained to reconstruct randomly masked words given others.
- Fine-tuned to conduct downstream tasks using [CLS].



# LiST

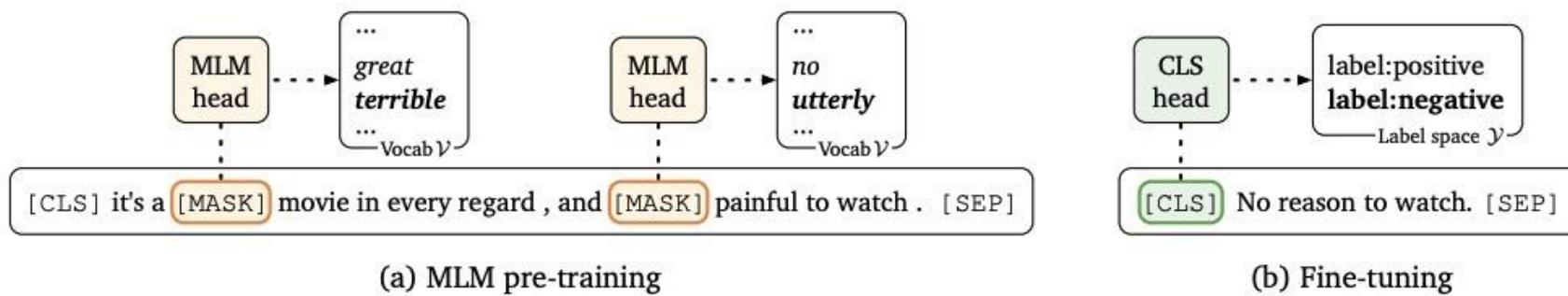
## Encoder-based Model:



Pre-training and Fine-tuning objective gap  
• **CLS token is used differently**

# LiST

## Encoder-based Model:



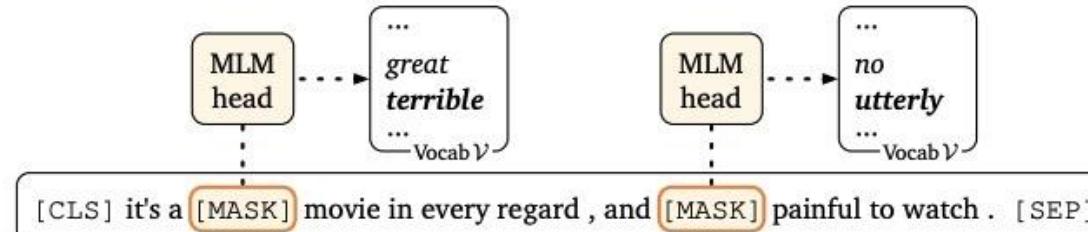
Use pre-training objective

[CLS] No reason to watch. It was [MASK]. [SEP]

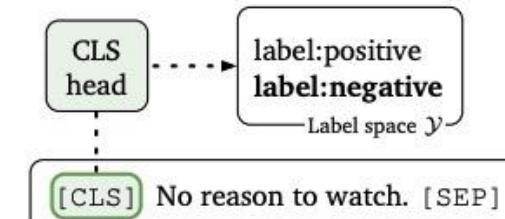
Original Input      Prompt

(c) Prompt-based Fine-tuning

# LiST



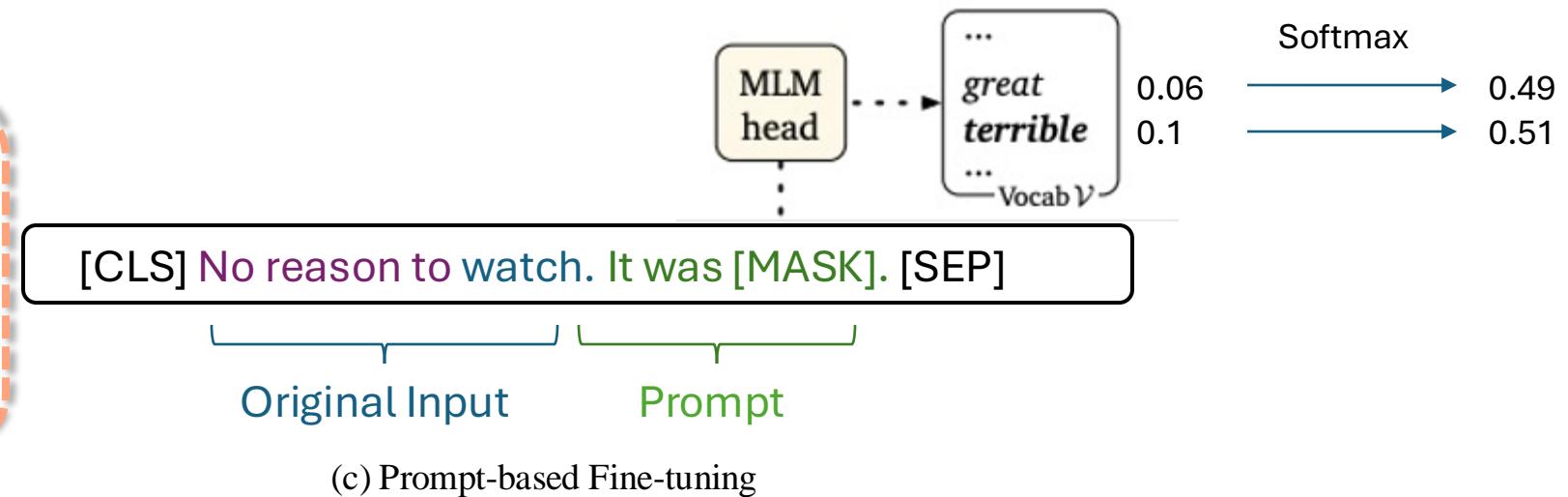
(a) MLM pre-training



(b) Fine-tuning

## Sentiment Analysis with pre-defined label words

- Use *great* and *terrible* in Vocab  $\mathcal{V}$  to indicate *pos* and *neg* classes



(c) Prompt-based Fine-tuning

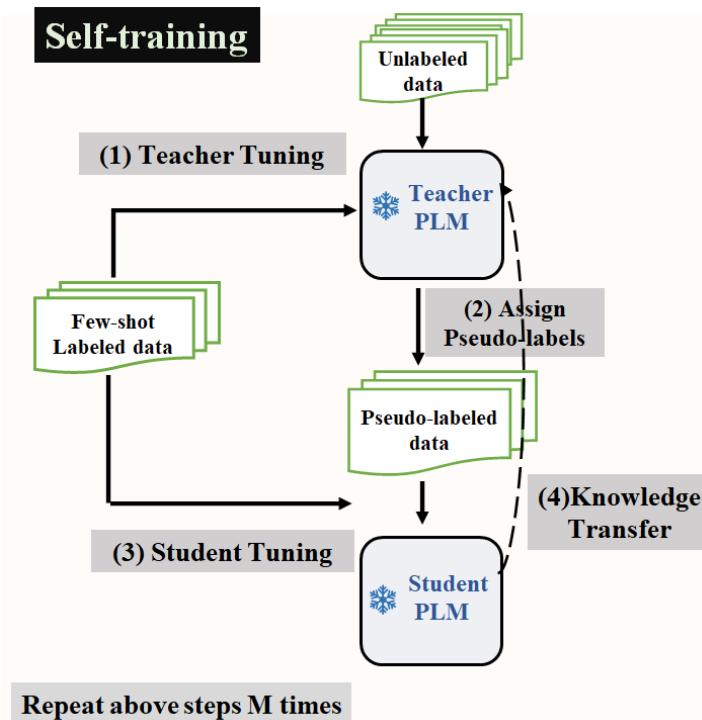
# LiST

## Standard Benchmarks:

- Large-scale labeled data
- No unlabeled data



## Self-training



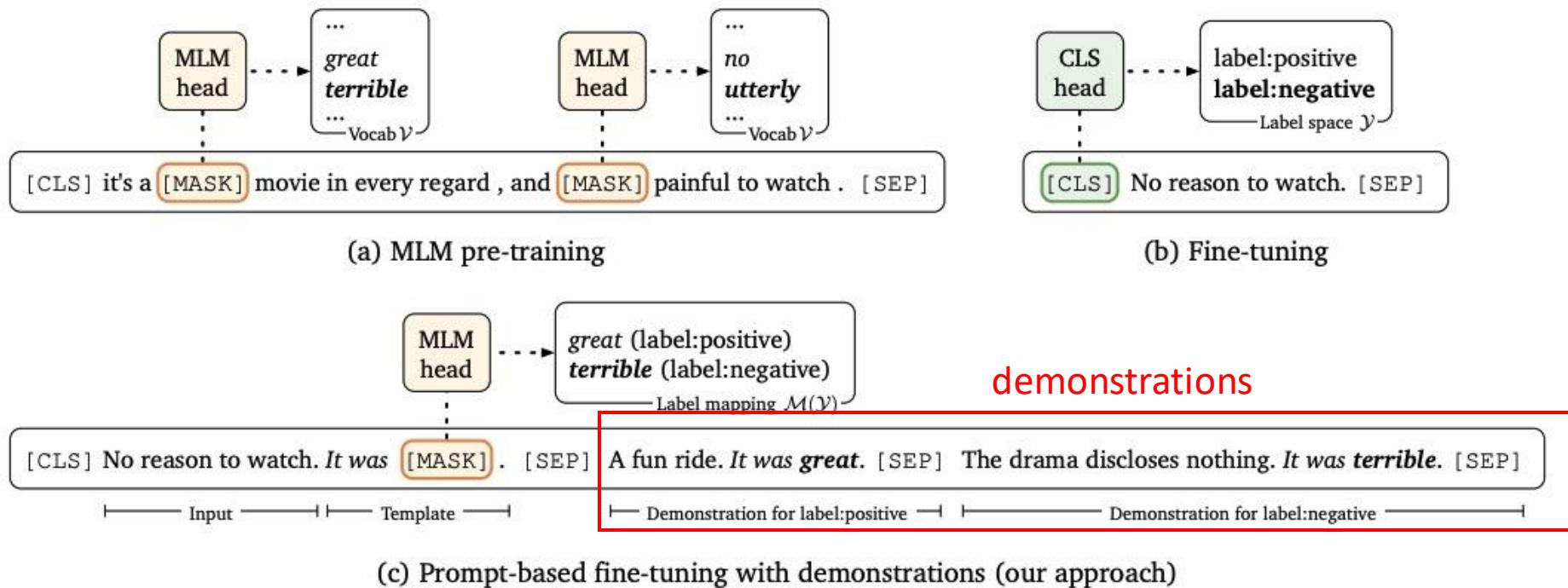
## Real-world Applications:

- Limited labeled data
- Many unlabeled data

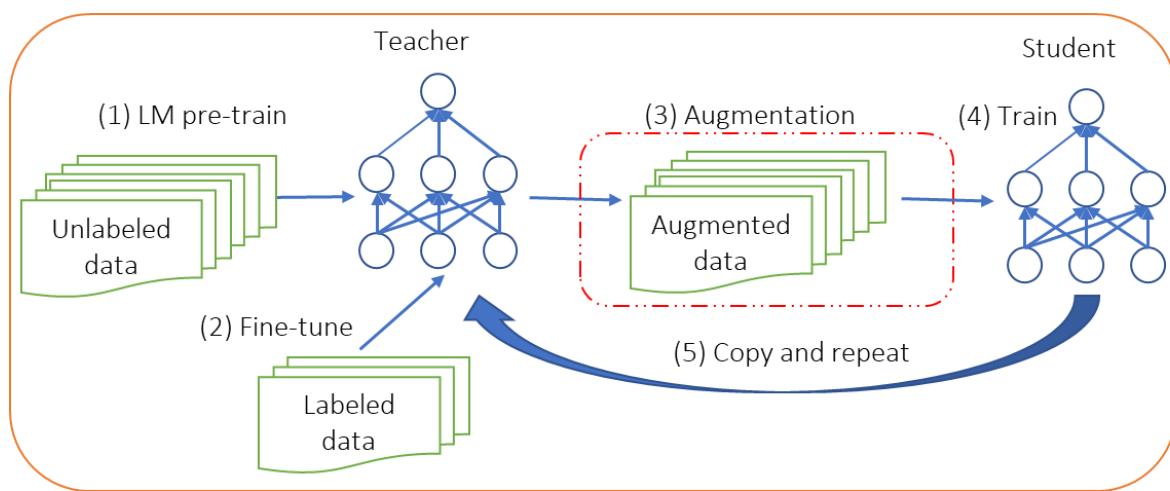
## LiST:

- Prompt-based tuning
- Self-training
- PEFT: Adapter

# LiST



# LiST



## Standard Self-training

$$\hat{\psi}_{\text{stu}}^{(t)} = \hat{\psi}_{\text{stu}}^{(t-1)} - \alpha \nabla \left( \frac{1}{N} \sum_{i=1}^N \mathcal{L}(\hat{y}_i^{(t)}, \text{enc}(x_i^u; \Theta_{\text{PLM}}, \hat{\psi}_{\text{stu}}^{(t-1)})) \right)$$

Pseudo-labels

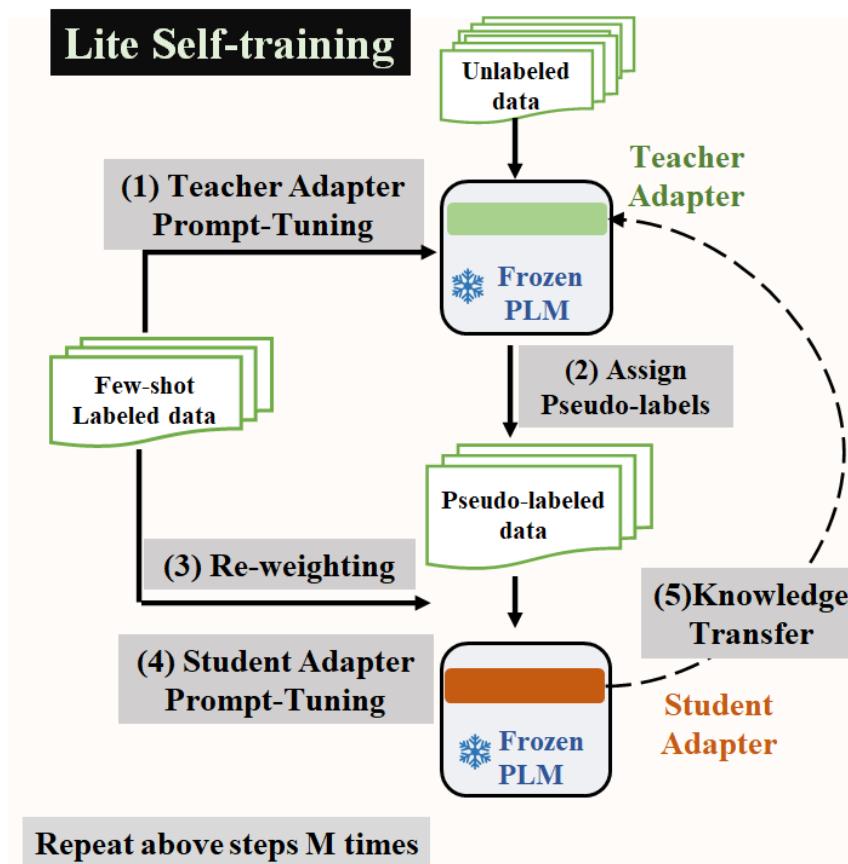
## Self-training with Re-weighting

$$\hat{\psi}_{\text{stu}}^{(t)}(\epsilon) = \hat{\psi}_{\text{stu}}^{(t-1)} - \alpha \nabla \left( \frac{1}{N} \sum_{i=1}^N [\epsilon_i^{(t)} \cdot \mathcal{L}(\hat{y}_i^{(t)}, \text{enc}(\tilde{x}_i^u; \Theta_{\text{PLM}}, \hat{\phi}_{\text{stu}}^{(t-1)}))] \right)$$

Weight for Pseudo-label:

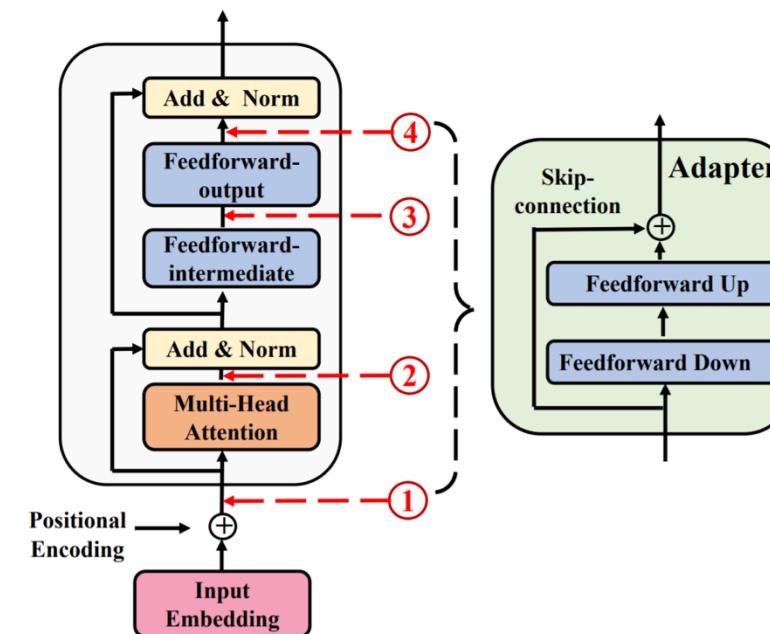
- Meta-learning (2nd-ord gradient).
- Validate w/ few-shot training samples.

# LiST



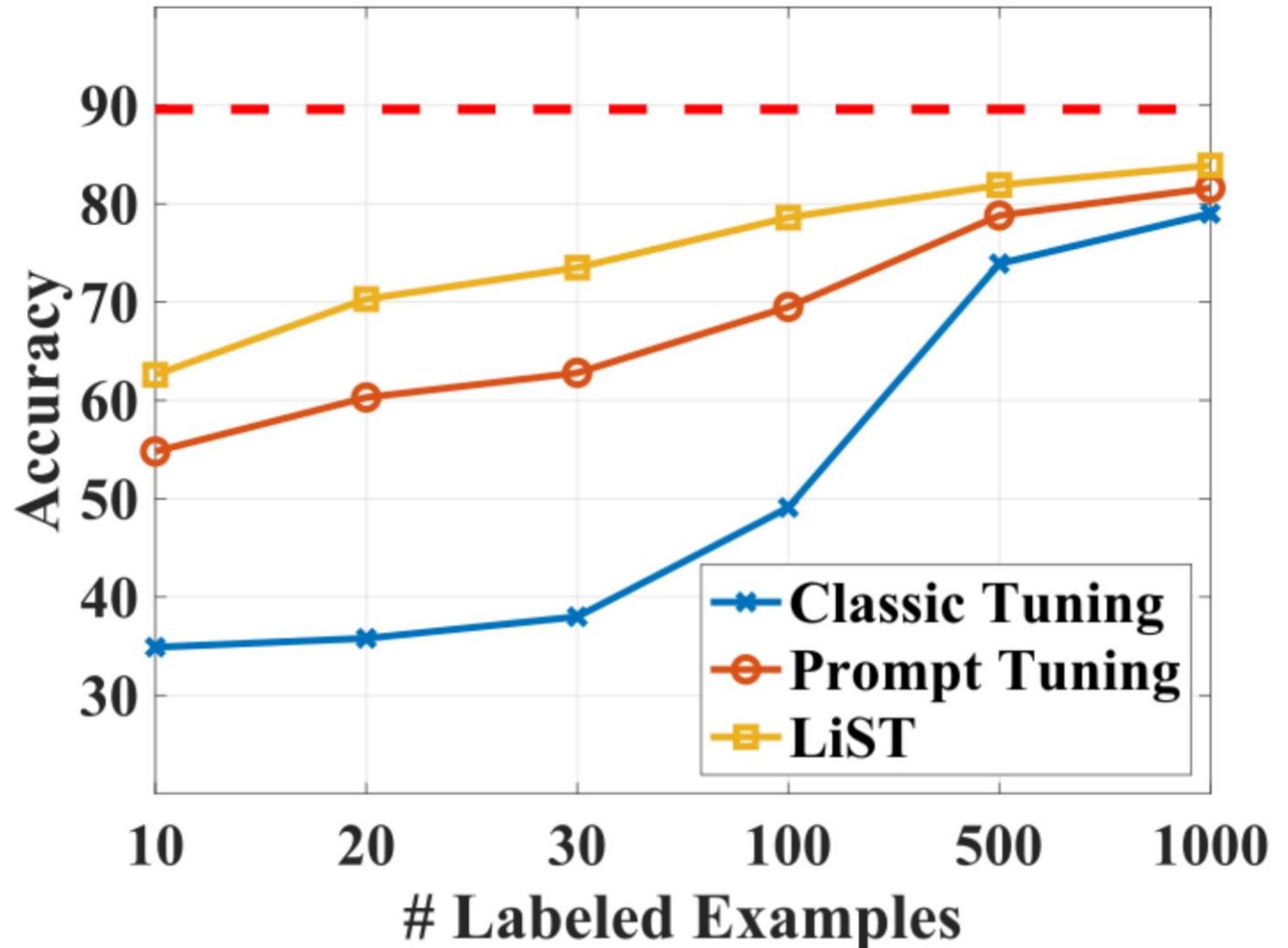
## Adapter

- Frozen PLM encoder.
- Teacher and student learn separate adapter.



# LiST

- “**Prompt Tuning**”: prompt-based tuning



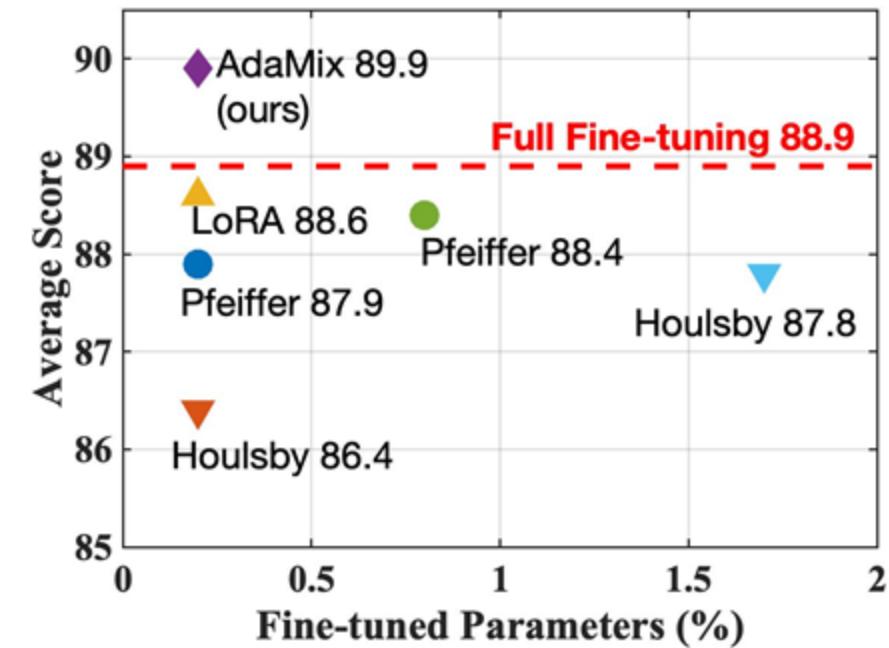
# AdaMix

## Performance gap between PEFT and FFT

- (Specialized: DoRA [ICML'24])

### AdaMix:

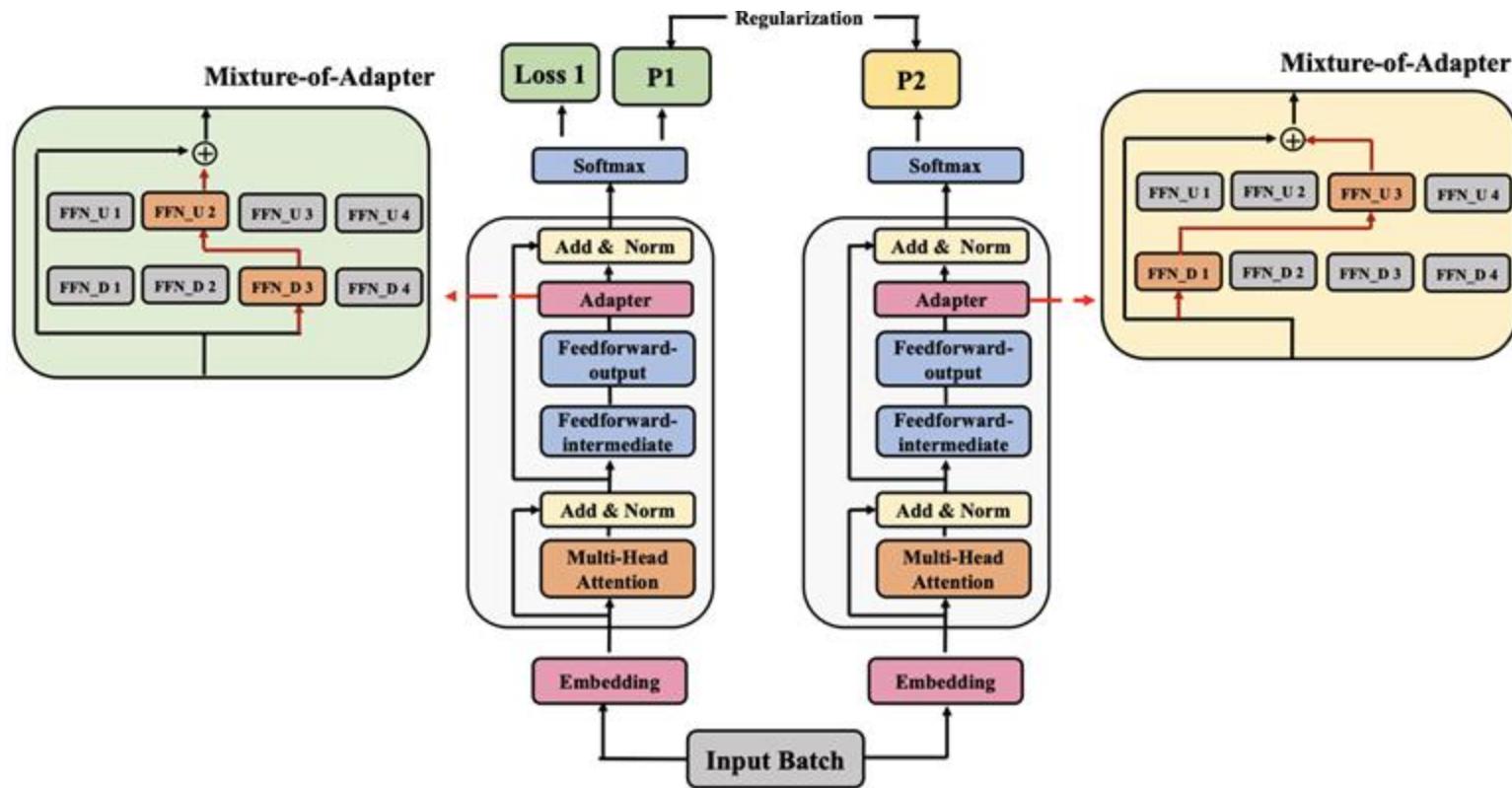
- Mixture of Expert + Weight ensemble.



(a) Average performance on GLUE development

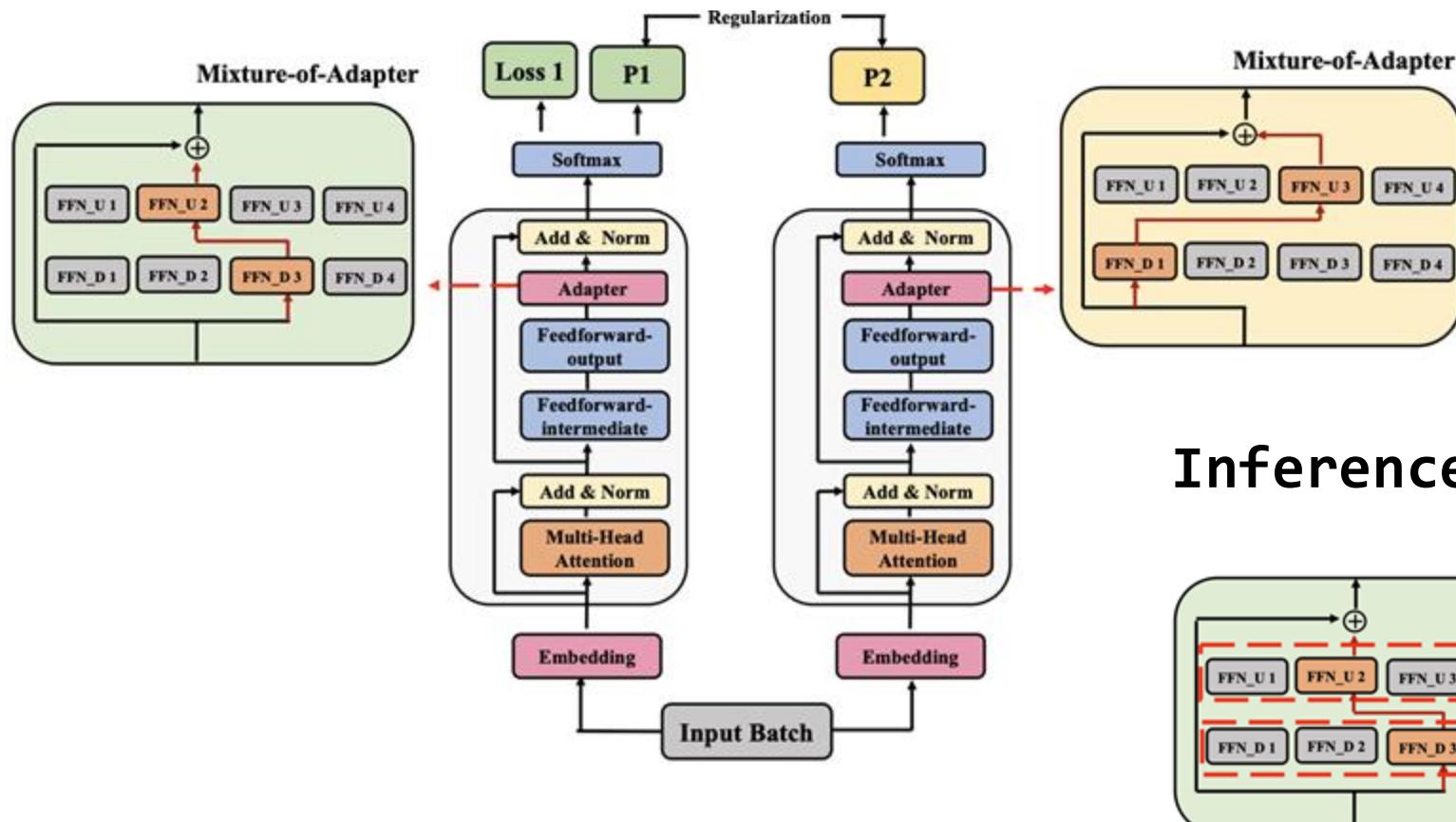
# AdaMix

## Training

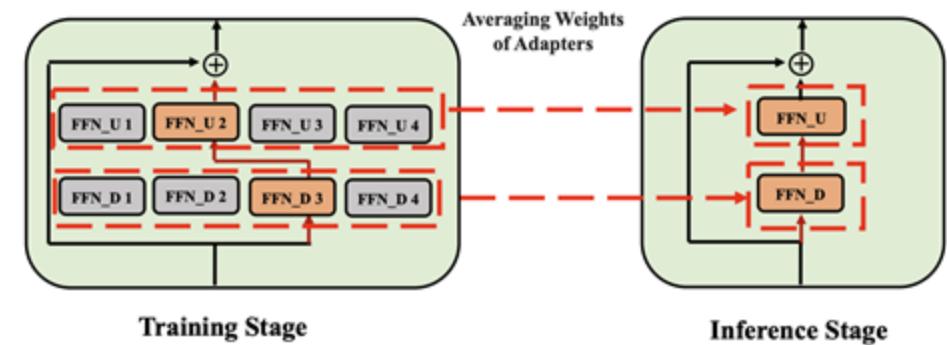


# AdaMix

## Training



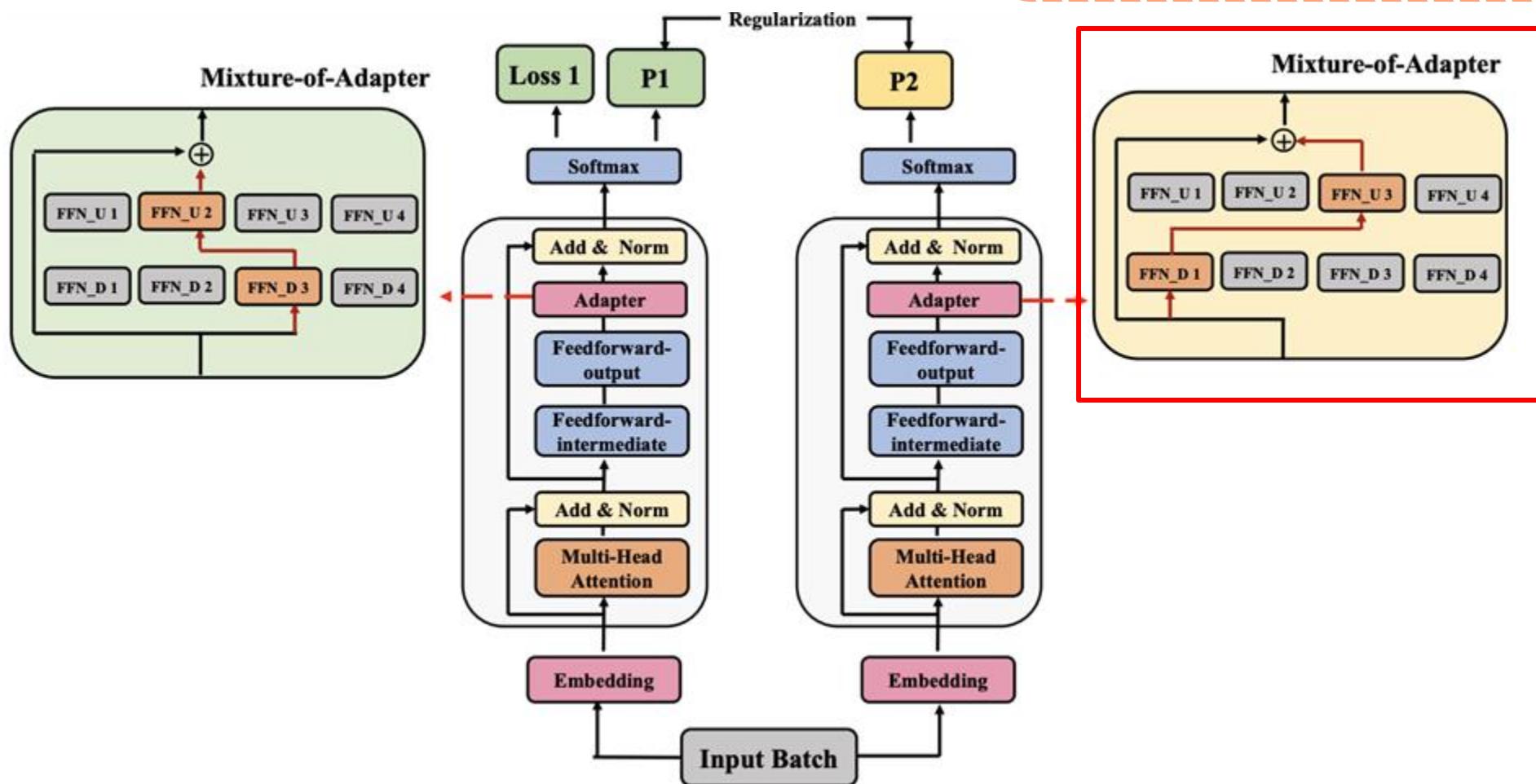
## Inference



# AdaMix

Adapter  $\Rightarrow$  Mixture-of-Adapter:

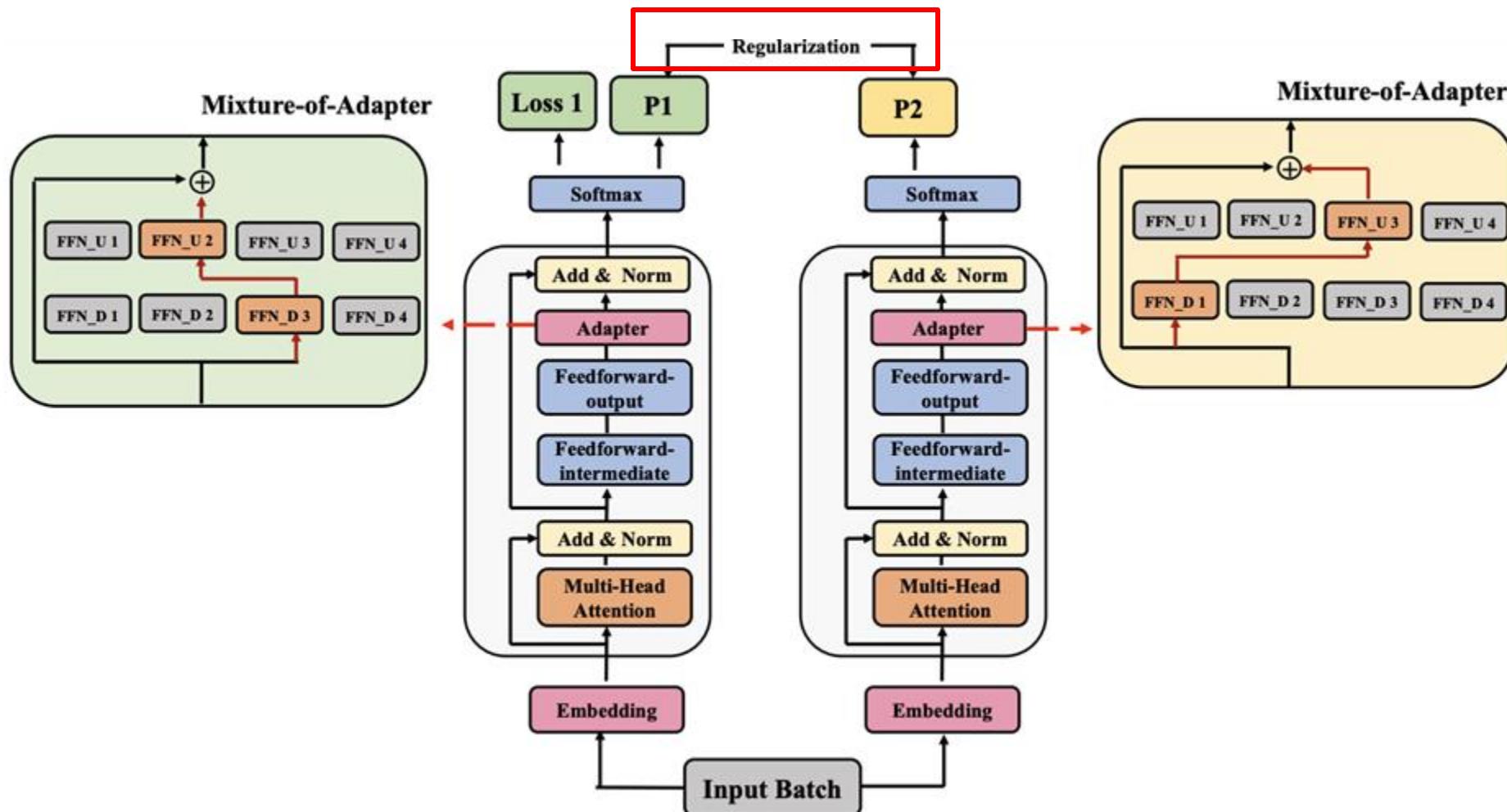
- Copies of Adapters
- Random Routing



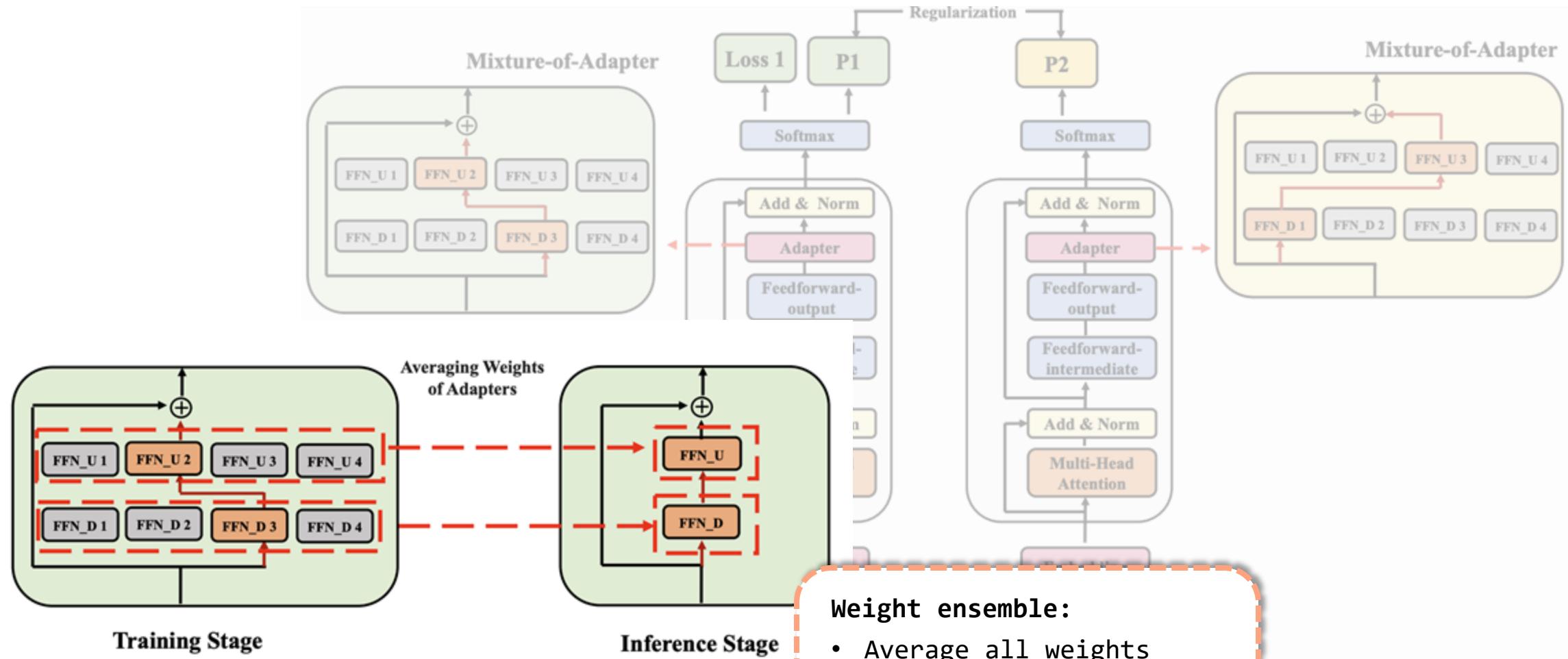
# AdaMix

Consistency loss between logits from two forwards

- Connect different copies of params to avoid collapse



# AdaMix



**Weight ensemble:**

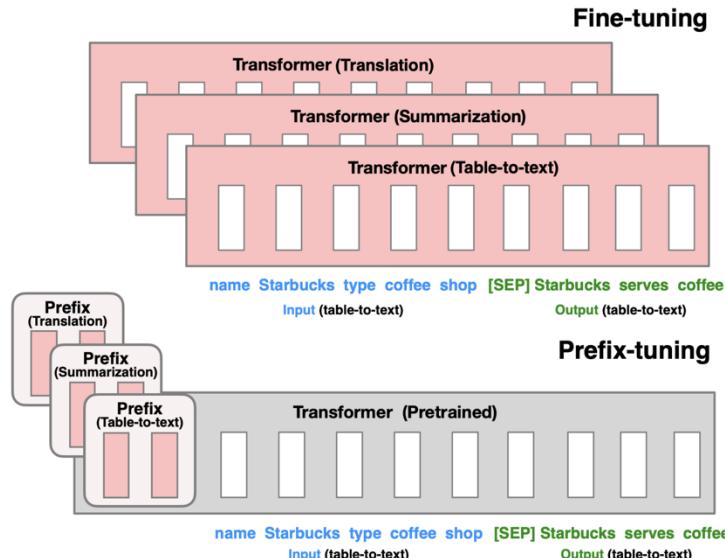
- Average all weights
- Improve the performance
- Reduce the storage cost

# AdaMix

Model	#Param.	MNLI Acc	QNLI Acc	SST2 Acc	QQP Acc	MRPC Acc	CoLA Mcc	RTE Acc	STS-B Pearson	Avg.
Full Fine-tuning <sup>†</sup>	355.0M	90.2	94.7	96.4	92.2	90.9	68.0	86.6	<b>92.4</b>	88.9
Pfeiffer Adapter <sup>†</sup>	3.0M	90.2	94.8	96.1	91.9	90.2	68.3	83.8	92.1	88.4
Pfeiffer Adapter <sup>†</sup>	0.8M	90.5	94.8	96.6	91.7	89.7	67.8	80.1	91.9	87.9
Houlsby Adapter <sup>†</sup>	6.0M	89.9	94.7	96.2	92.1	88.7	66.5	83.4	91.0	87.8
Houlsby Adapter <sup>†</sup>	0.8M	90.3	94.7	96.3	91.5	87.7	66.3	72.9	91.5	86.4
LoRA <sup>†</sup>	0.8M	90.6	94.8	96.2	91.6	90.2	68.2	85.2	92.3	88.6
AdaMix Adapter	0.8M	<b>90.9</b>	<b>95.4</b>	<b>97.1</b>	<b>92.3</b>	<b>91.9</b>	<b>70.2</b>	<b>89.2</b>	<b>92.4</b>	<b>89.9</b>

There are more PEFT methods to mix.

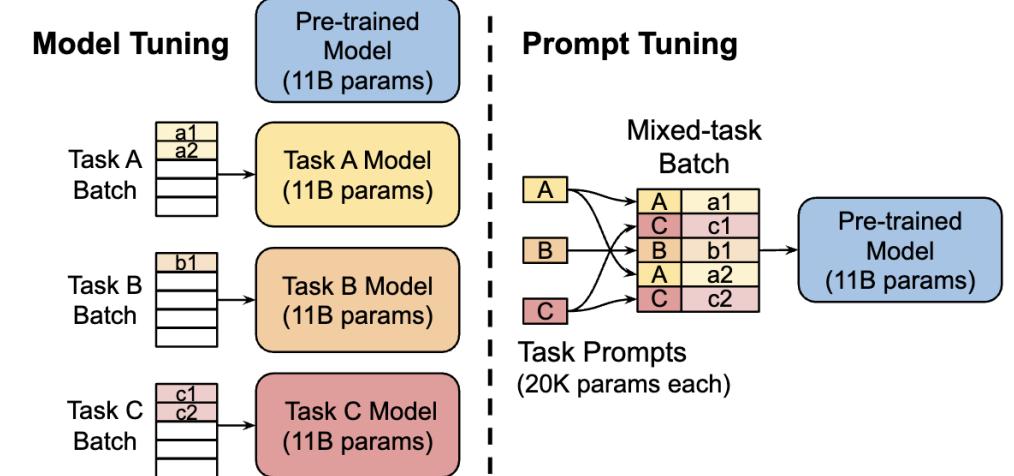
# More Architectures...



$$\text{Each layer: } h_i = \begin{cases} P_\theta[i, :], & \text{if } i \in \mathcal{P}_{\text{idx}}, \\ \text{LM}_\phi(z_i, h_{<i}), & \text{otherwise.} \end{cases} \quad (3)$$

(a) Prefix-tuning

Prepend learnable representations to the actual input



Only input layer:  $[P_e; X_e] \in \mathbb{R}^{(p+n) \times e}$

(b) Prompt-tuning

- [1] Li and Liang. "Prefix-tuning: Optimizing continuous prompts for generation." (ACL), 2021.  
 [2] Lester et al. "The power of scale for parameter-efficient prompt tuning." Empirical Methods in Natural Language Processing (EMNLP), 2022.

# Outline

- Model Efficiency
  - Parameter-efficient Fine-tuning
  - Model Compression
  - Inference Acceleration
- Data Efficiency
  - Zero/few-shot Learning
  - Knowledge Editing
  - Retrieval-augmented Generation

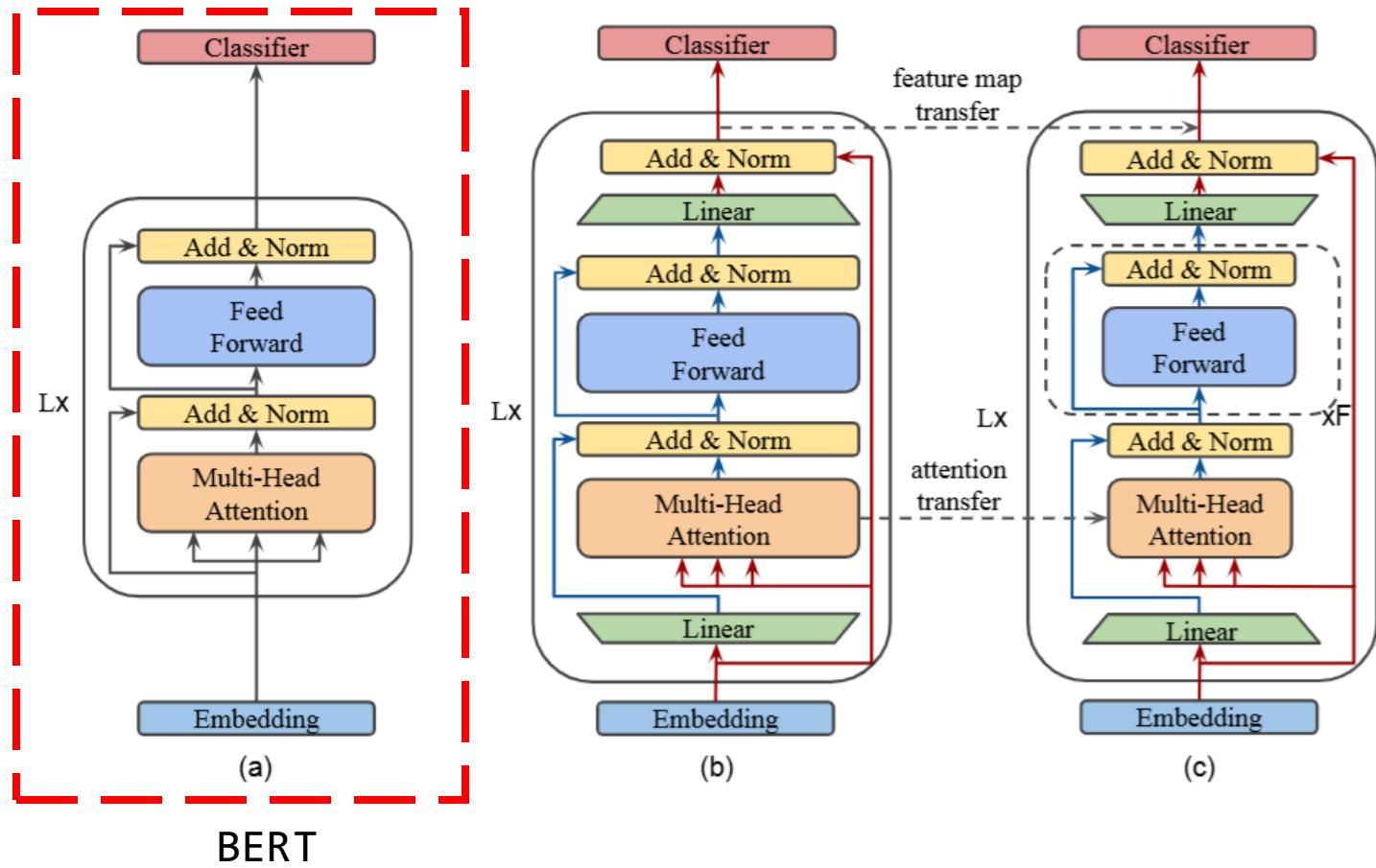
# Model Compression

- PEFT methods allow fine-tuning pre-trained LLMs for specific downstream applications with limited resources.
- However, when applications themselves have limited computational resources, they may not afford to run fine-tuned LLMs.
- Therefore, **model compression** is needed.

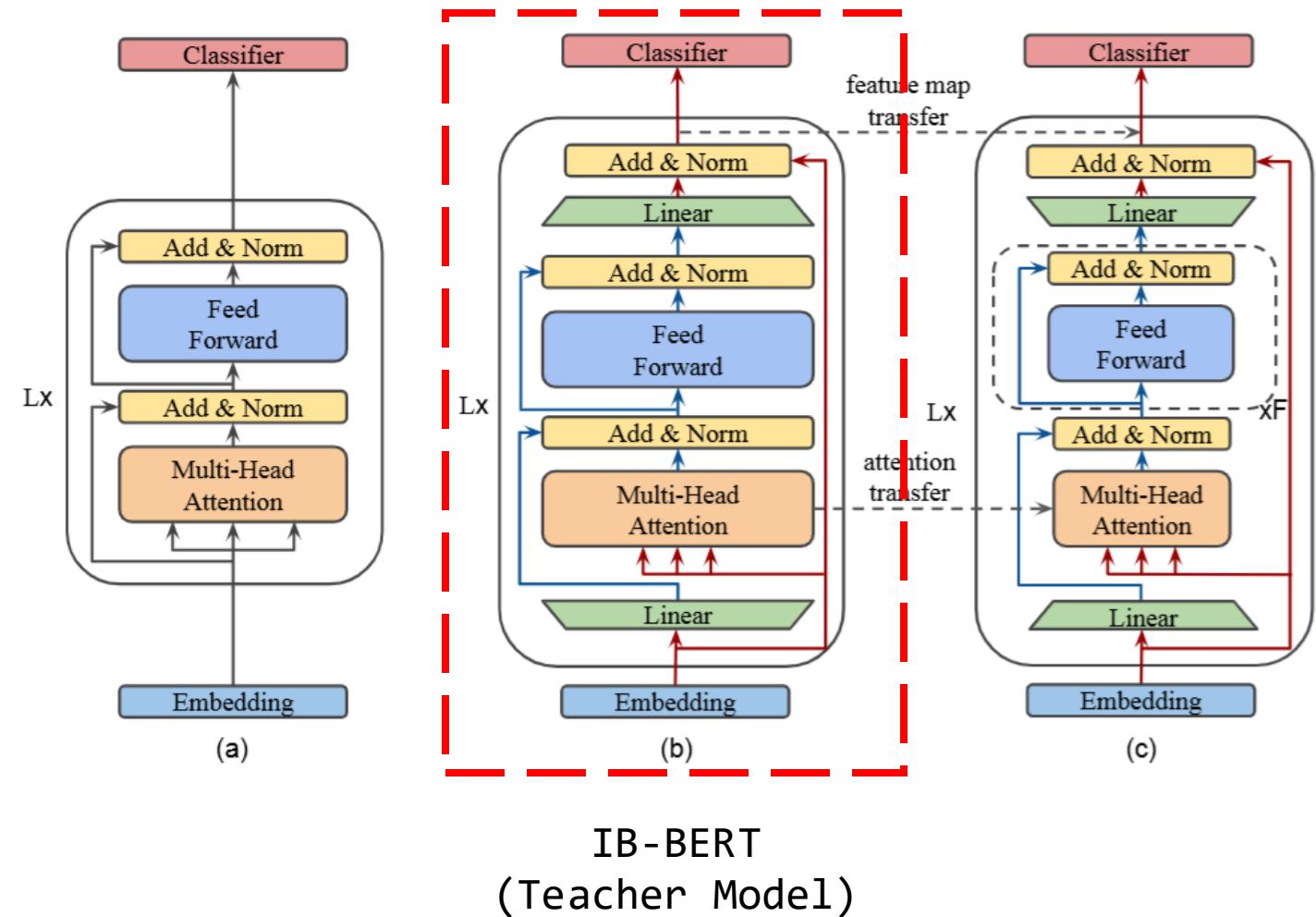
# MobileBERT

- The compressed model to have strong performance while being small.
  - As deep as BERT but much thinner (bottleneck).
- It is hard to train a deep & thin network from the scratch.
  - Design a special teacher network for knowledge transfer (inverted-bottleneck).

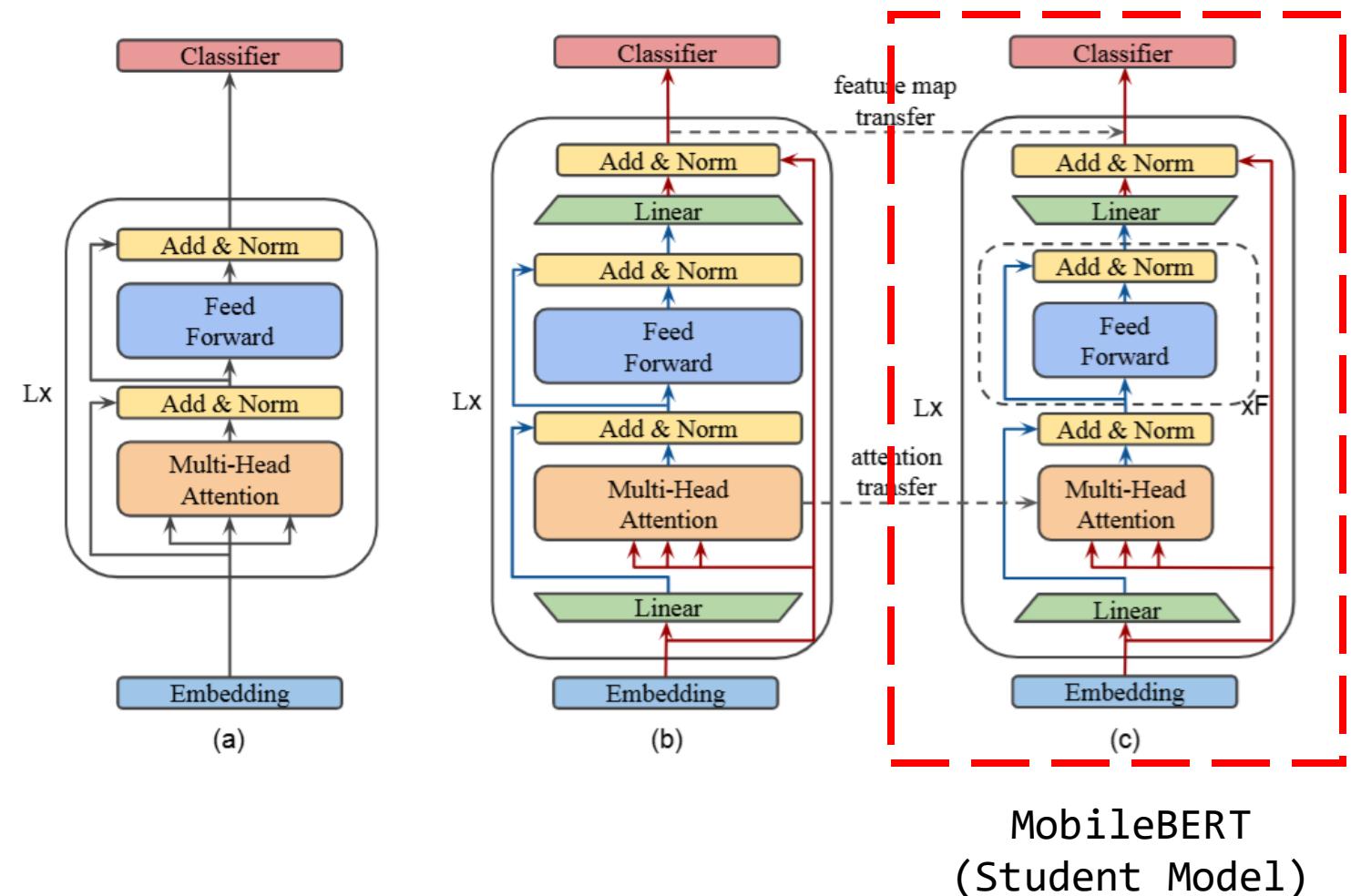
# MobileBERT



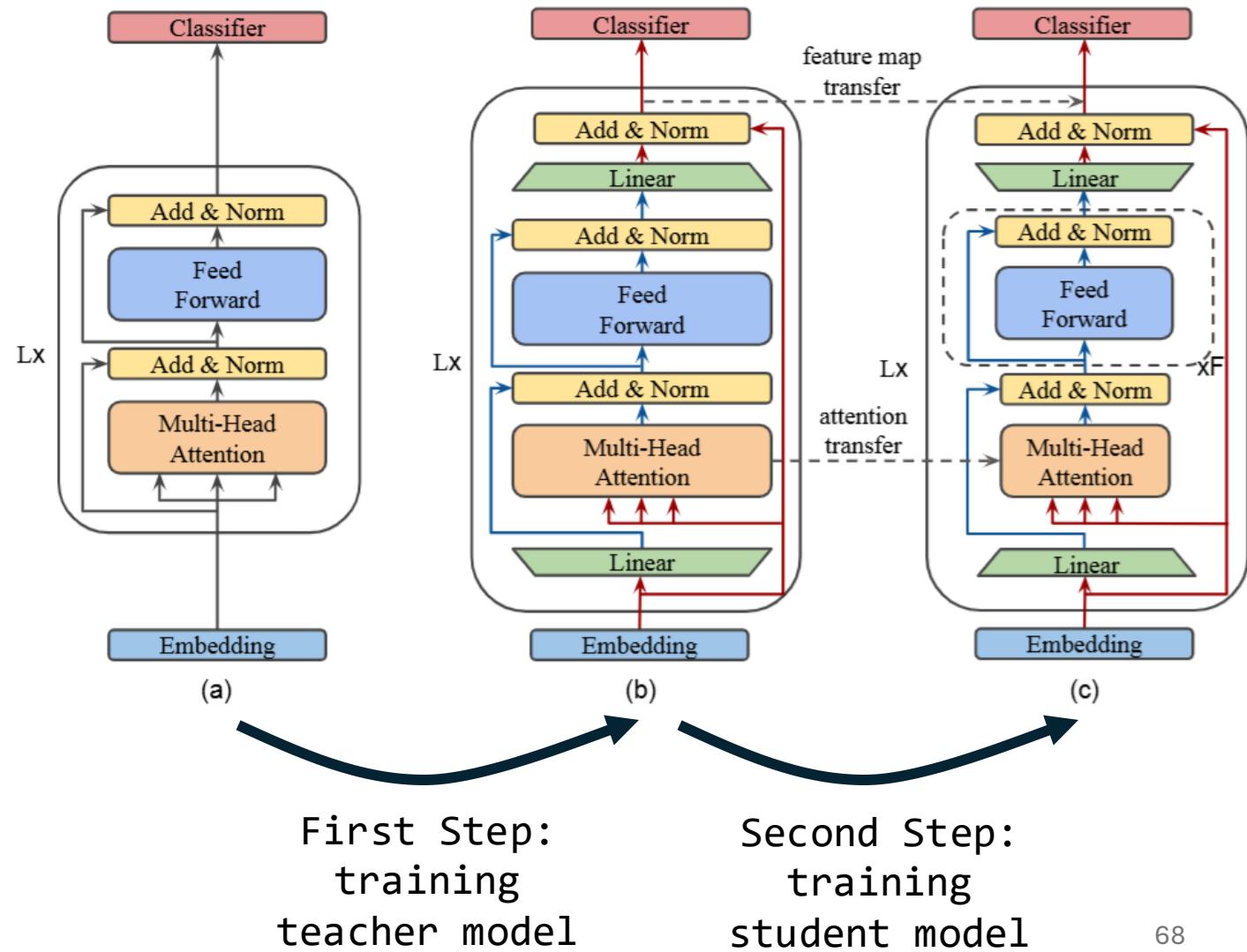
# MobileBERT



# MobileBERT



# MobileBERT



# MobileBERT

- IB-BERT is a little lighter than BERT while matching the hidden space dimension of MobileBERT.

		BERT <sub>LARGE</sub>	BERT <sub>BASE</sub>	IB-BERT <sub>LARGE</sub>	MobileBERT	MobileBERT <sub>TINY</sub>
embedding	h <sub>embedding</sub>	1024	768		128	
	no-op	no-op			3-convolution	
	h <sub>inter</sub>	1024	768		512	
body	Linear	h <sub>input</sub> h <sub>output</sub>				
	MHA	h <sub>input</sub> #Head h <sub>output</sub>	$\left[ \begin{pmatrix} 1024 \\ 16 \\ 1024 \end{pmatrix} \right] \times 24$	$\left[ \begin{pmatrix} 768 \\ 12 \\ 768 \end{pmatrix} \right] \times 12$	$\left[ \begin{pmatrix} 512 \\ 1024 \\ 512 \\ 4 \\ 1024 \\ 1024 \\ 4096 \\ 1024 \\ 1024 \\ 512 \end{pmatrix} \right] \times 24$	$\left[ \begin{pmatrix} 512 \\ 128 \\ 512 \\ 4 \\ 128 \\ 128 \\ 512 \\ 128 \\ 128 \\ 512 \end{pmatrix} \right] \times 24$
	FFN	h <sub>input</sub> h <sub>FFN</sub> h <sub>output</sub>	$\left[ \begin{pmatrix} 1024 \\ 4096 \\ 1024 \end{pmatrix} \right]$	$\left[ \begin{pmatrix} 768 \\ 3072 \\ 768 \end{pmatrix} \right]$	$\left[ \begin{pmatrix} 128 \\ 512 \\ 128 \\ 128 \\ 128 \end{pmatrix} \right] \times 4$	$\left[ \begin{pmatrix} 512 \\ 128 \\ 128 \\ 4 \\ 128 \\ 128 \\ 512 \\ 128 \\ 128 \\ 512 \end{pmatrix} \right] \times 24$
	Linear	h <sub>input</sub> h <sub>output</sub>				
#Params		334M	109M	293M	25.3M	15.1M

# MobileBERT

- IB-BERT is a little lighter than BERT while matching the hidden space dimension of MobileBERT.

		BERT <sub>LARGE</sub>	BERT <sub>BASE</sub>	IB-BERT <sub>LARGE</sub>	MobileBERT	MobileBERT <sub>TINY</sub>
embedding	h <sub>embedding</sub>	1024	768		128	
	h <sub>inter</sub>	no-op	no-op		3-convolution	
	h <sub>inter</sub>	1024	768		512	
body	Linear	h <sub>input</sub> h <sub>output</sub>				
	MHA	h <sub>input</sub> #Head h <sub>output</sub>	$\begin{bmatrix} \begin{pmatrix} 1024 \\ 16 \\ 1024 \end{pmatrix} \\ \times 24 \end{bmatrix}$	$\begin{bmatrix} \begin{pmatrix} 768 \\ 12 \\ 768 \end{pmatrix} \\ \times 12 \end{bmatrix}$	$\begin{bmatrix} \begin{pmatrix} 512 \\ 1024 \\ 512 \\ 4 \\ 1024 \\ 1024 \\ 4096 \\ 1024 \\ 1024 \\ 512 \end{pmatrix} \\ \times 24 \end{bmatrix}$	$\begin{bmatrix} \begin{pmatrix} 512 \\ 128 \\ 512 \\ 4 \\ 128 \\ 128 \\ 512 \\ 128 \\ 128 \\ 512 \end{pmatrix} \\ \times 24 \end{bmatrix}$
	FFN	h <sub>input</sub> h <sub>FFN</sub> h <sub>output</sub>	$\begin{bmatrix} \begin{pmatrix} 1024 \\ 4096 \\ 1024 \end{pmatrix} \\ \times 24 \end{bmatrix}$	$\begin{bmatrix} \begin{pmatrix} 768 \\ 3072 \\ 768 \end{pmatrix} \\ \times 24 \end{bmatrix}$	$\begin{bmatrix} \begin{pmatrix} 128 \\ 512 \\ 128 \\ 128 \\ 128 \\ 512 \end{pmatrix} \\ \times 4 \end{bmatrix}$	$\begin{bmatrix} \begin{pmatrix} 512 \\ 128 \\ 128 \\ 4 \\ 128 \\ 128 \\ 512 \\ 128 \\ 128 \\ 512 \end{pmatrix} \\ \times 2 \end{bmatrix}$
	Linear	h <sub>input</sub> h <sub>output</sub>				
#Params		334M	109M	293M	25.3M	15.1M

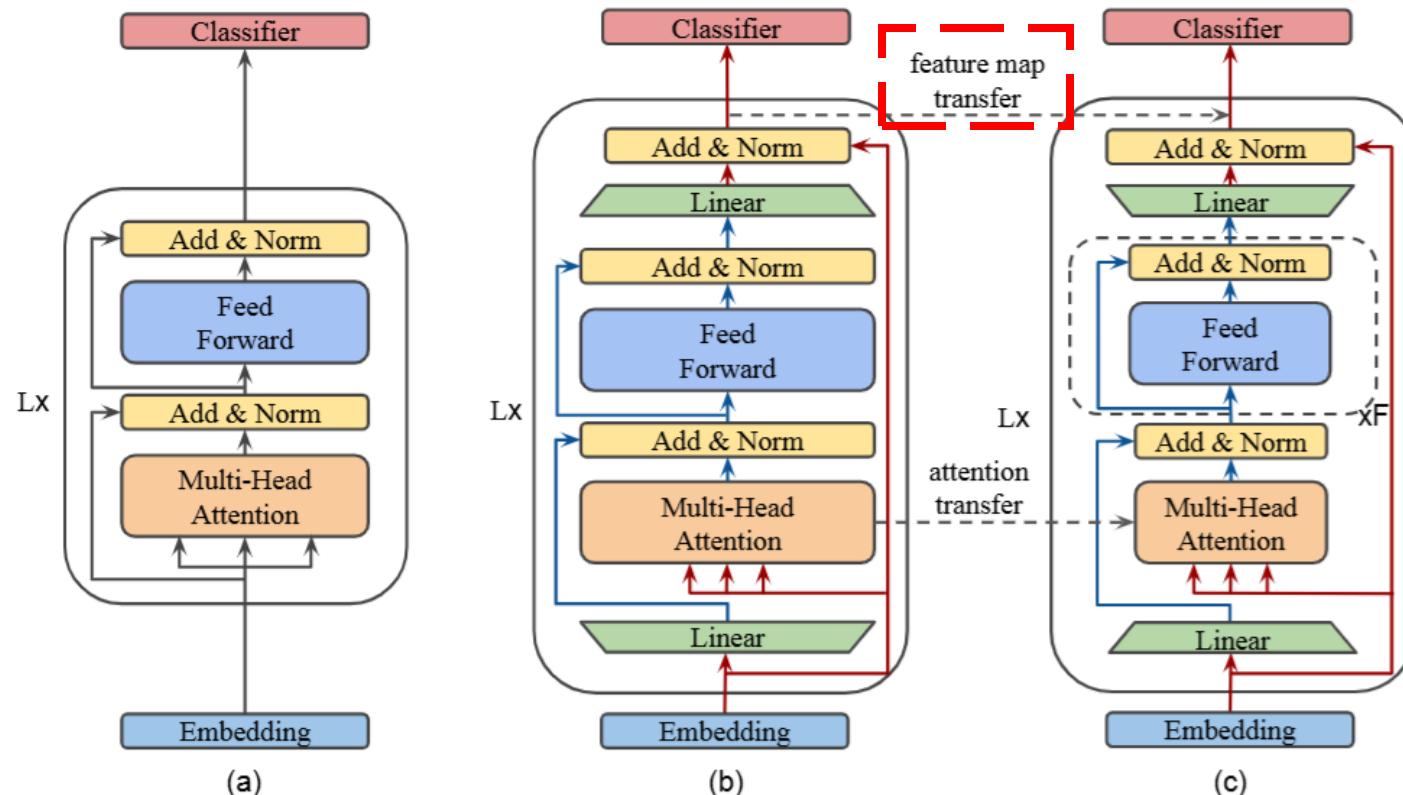
# MobileBERT

- IB-BERT is a little lighter than BERT while matching the hidden space dimension of MobileBERT.

		BERT <sub>LARGE</sub>	BERT <sub>BASE</sub>	IB-BERT <sub>LARGE</sub>	MobileBERT	MobileBERT <sub>TINY</sub>
embedding	h <sub>embedding</sub>	1024	768		128	
	no-op	no-op			3-convolution	
	h <sub>inter</sub>	1024	768		512	
body	Linear	h <sub>input</sub> h <sub>output</sub>				
	MHA	h <sub>input</sub> #Head h <sub>output</sub>	$\begin{bmatrix} \begin{pmatrix} 1024 \\ 16 \\ 1024 \end{pmatrix} \\ \times 24 \end{bmatrix}$	$\begin{bmatrix} \begin{pmatrix} 768 \\ 12 \\ 768 \end{pmatrix} \\ \times 12 \end{bmatrix}$	$\begin{bmatrix} \begin{pmatrix} 512 \\ 1024 \\ 512 \\ 4 \\ 1024 \\ 1024 \\ 4096 \\ 1024 \end{pmatrix} \\ \times 24 \end{bmatrix}$	$\begin{bmatrix} \begin{pmatrix} 512 \\ 128 \\ 512 \\ 4 \\ 128 \\ 128 \\ 512 \\ 128 \\ 128 \end{pmatrix} \\ \times 4 \\ \times 24 \end{bmatrix}$
	FFN	h <sub>input</sub> h <sub>FFN</sub> h <sub>output</sub>	$\begin{bmatrix} \begin{pmatrix} 1024 \\ 4096 \\ 1024 \end{pmatrix} \\ \times 24 \end{bmatrix}$	$\begin{bmatrix} \begin{pmatrix} 768 \\ 768 \\ 3072 \\ 768 \end{pmatrix} \\ \times 12 \end{bmatrix}$	$\begin{bmatrix} \begin{pmatrix} 1024 \\ 512 \end{pmatrix} \\ \times 4 \\ \times 24 \end{bmatrix}$	$\begin{bmatrix} \begin{pmatrix} 512 \\ 128 \\ 128 \\ 4 \\ 128 \\ 128 \\ 512 \\ 128 \\ 128 \end{pmatrix} \\ \times 2 \\ \times 24 \end{bmatrix}$
	Linear	h <sub>input</sub> h <sub>output</sub>				
#Params		334M	109M	293M	25.3M	15.1M

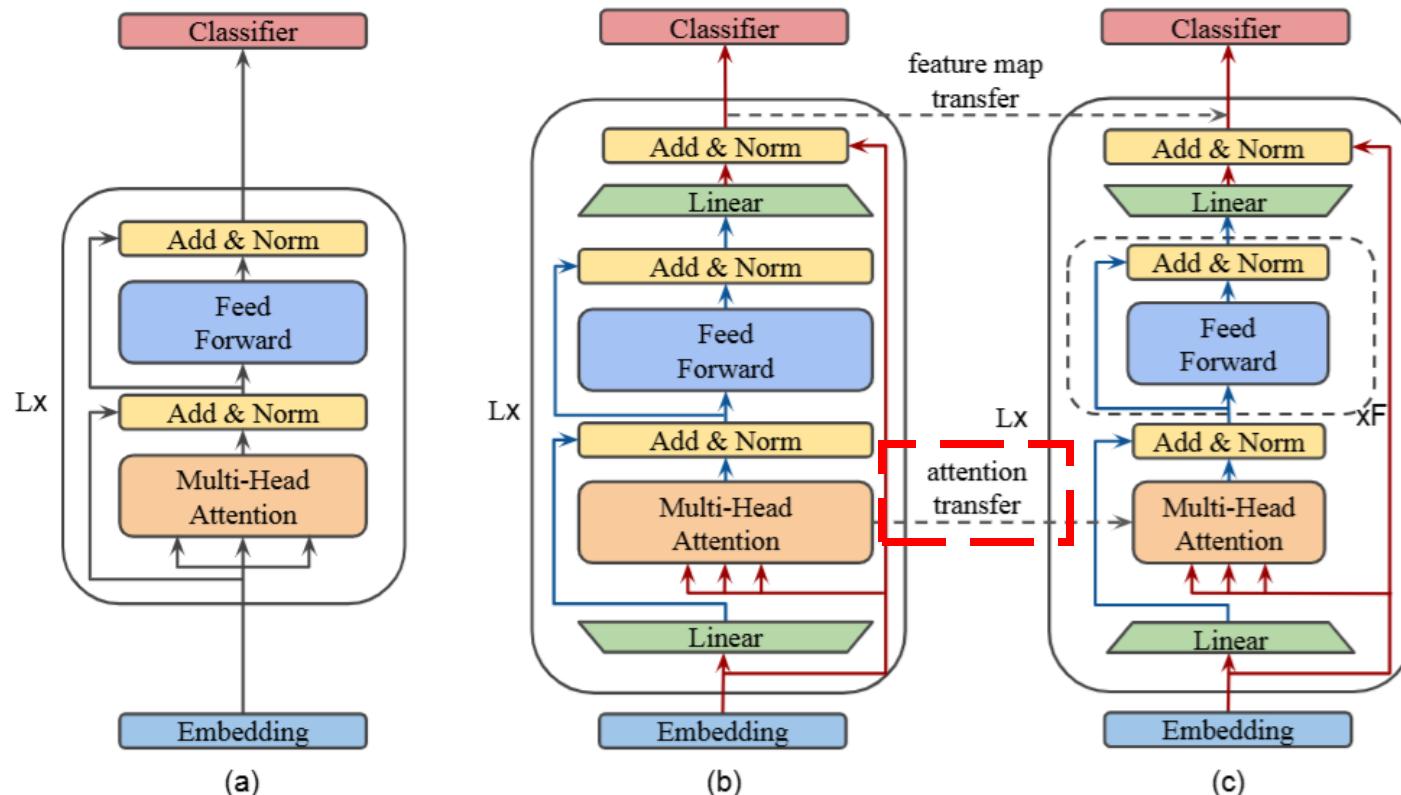
# MobileBERT

- MSE loss:  $\mathcal{L}_{FMT}^{\ell} = \frac{1}{TN} \sum_{t=1}^T \sum_{n=1}^N (H_{t,\ell,n}^{tr} - H_{t,\ell,n}^{st})^2$



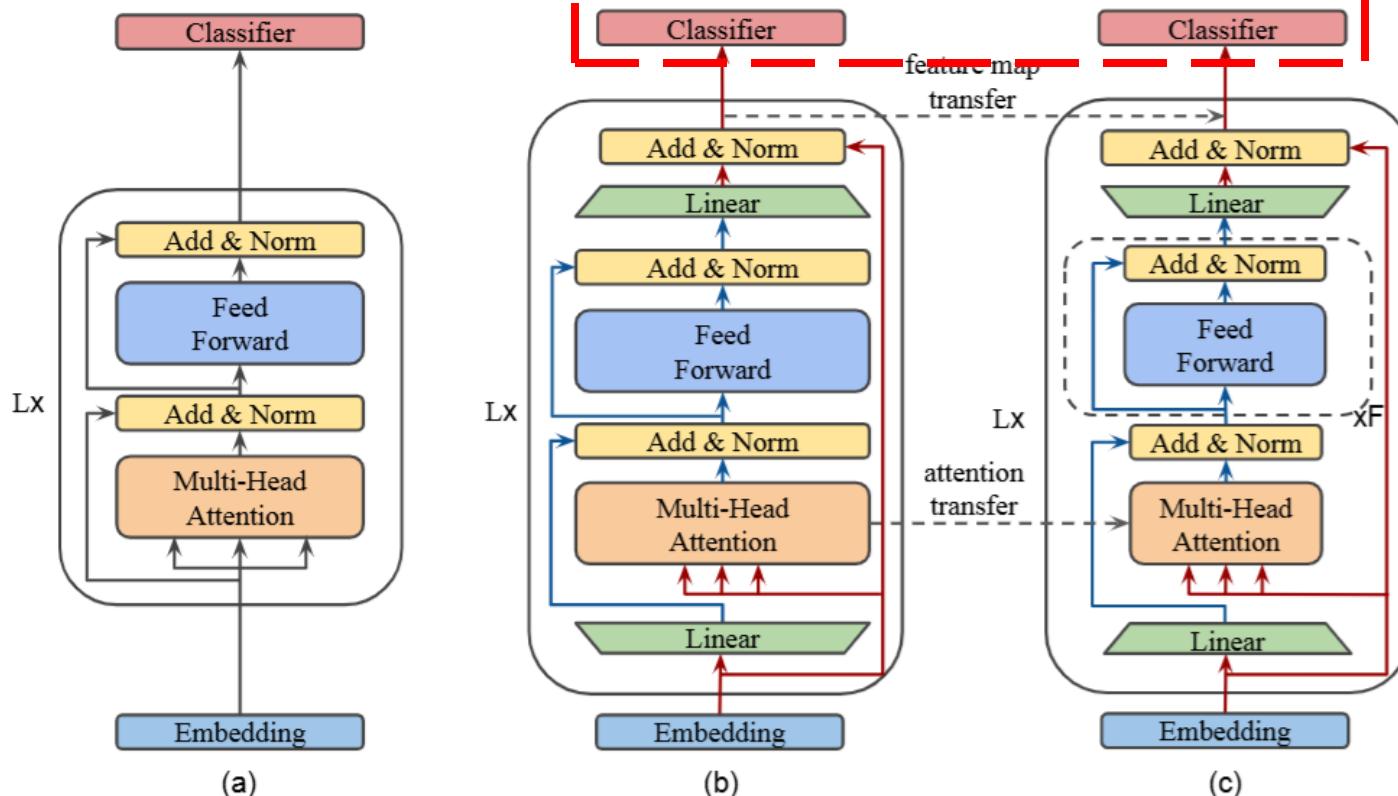
# MobileBERT

- KL penalty:  $\mathcal{L}_{AT}^{\ell} = \frac{1}{TN} \sum_{t=1}^T \sum_{n=1}^N D_{KL}(a_{t,\ell,n}^{tr} || a_{t,\ell,n}^{st})$



# MobileBERT

$$\bullet \quad \mathcal{L}_{PD} = \alpha \mathcal{L}_{MLM} + (1 - \alpha) \mathcal{L}_{KD} + \mathcal{L}_{NSP}$$

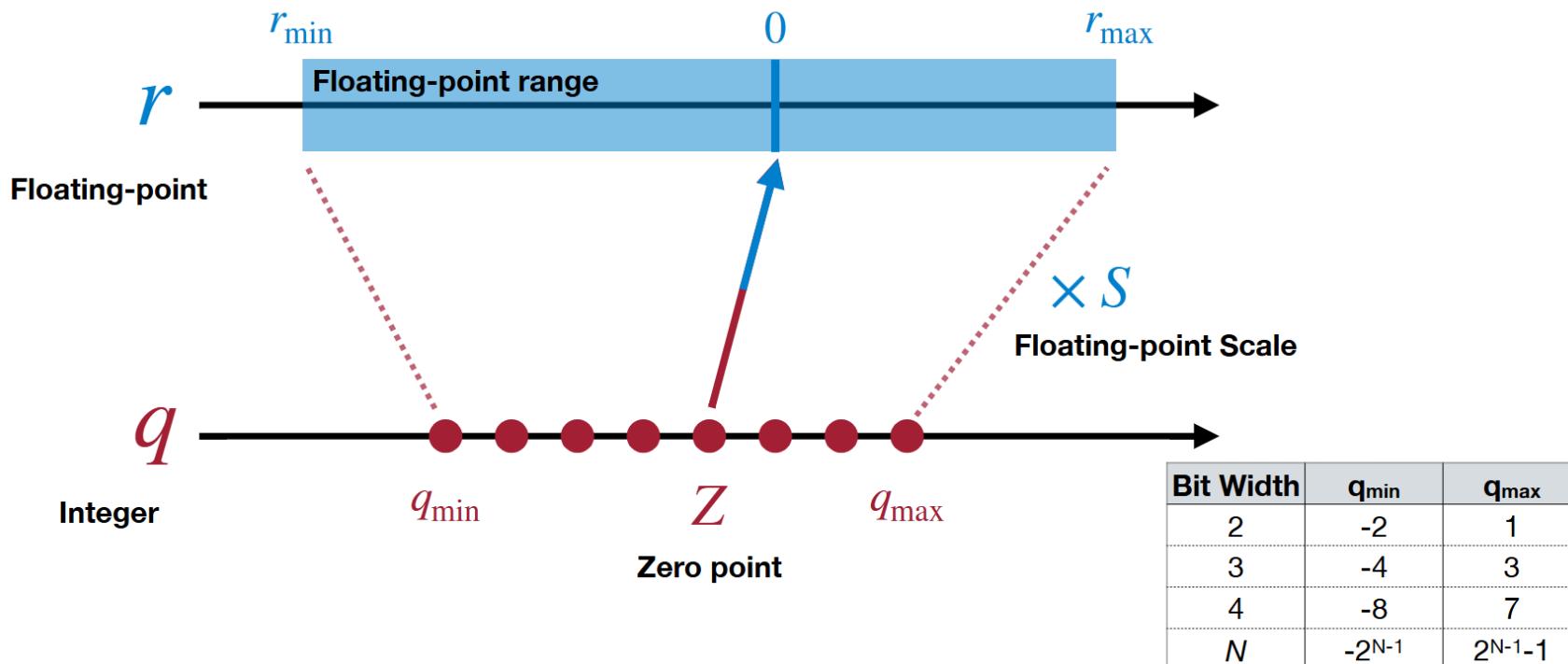


# MobileBERT

	#Params	#FLOPS	Latency	SQuAD 2.0 Metrics									GLUE
				CoLA	SST-2	MRPC	STS-B	QQP	MNLI-m/mm	QNLI	RTE		
ELMo-BiLSTM-Attn	-	-	-	33.6	90.4	84.4	72.3	63.1	74.1/74.5	79.8	58.9	70.0	
OpenAI GPT	109M	-	-	47.2	93.1	87.7	84.8	70.1	80.7/80.6	87.2	69.1	76.9	
BERT <sub>BASE</sub>	109M	22.5B	342 ms	<b>52.1</b>	<b>93.5</b>	<b>88.9</b>	<b>85.8</b>	71.2	<b>84.6/83.4</b>	90.5	66.4	78.3	
BERT <sub>BASE</sub> -6L-PKD*	66.5M	11.3B	-	-	92.0	85.0	-	70.7	81.5/81.0	89.0	65.5	-	
BERT <sub>BASE</sub> -4L-PKD†*	52.2M	7.6B	-	24.8	89.4	82.6	79.8	70.2	79.9/79.3	85.1	62.3	-	
BERT <sub>BASE</sub> -3L-PKD*	45.3M	5.7B	-	-	87.5	80.7	-	68.1	76.7/76.3	84.7	58.2	-	
DistilBERT <sub>BASE</sub> -6L†	62.2M	11.3B	-	-	92.0	85.0	-	70.7	81.5/81.0	89.0	65.5	-	
DistilBERT <sub>BASE</sub> -4L†	52.2M	7.6B	-	32.8	91.4	82.4	76.1	68.5	78.9/78.0	85.2	54.1	-	
TinyBERT*	14.5M	1.2B	-	43.3	92.6	86.4	79.9	<b>71.3</b>	82.5/81.8	87.7	62.9	75.4	
MobileBERT <sub>TINY</sub>	15.1M	3.1B	40 ms	46.7	91.7	87.9	80.1	68.9	81.5/81.6	89.5	65.1	75.8	
MobileBERT	25.3M	5.7B	62 ms	50.5	92.8	88.8	84.4	70.2	83.3/82.6	90.6	66.2	77.7	
MobileBERT w/o OPT	25.3M	5.7B	192 ms	51.1	92.6	88.8	84.8	70.5	84.3/83.4	<b>91.6</b>	<b>70.4</b>	<b>78.5</b>	

# Linear Quantization

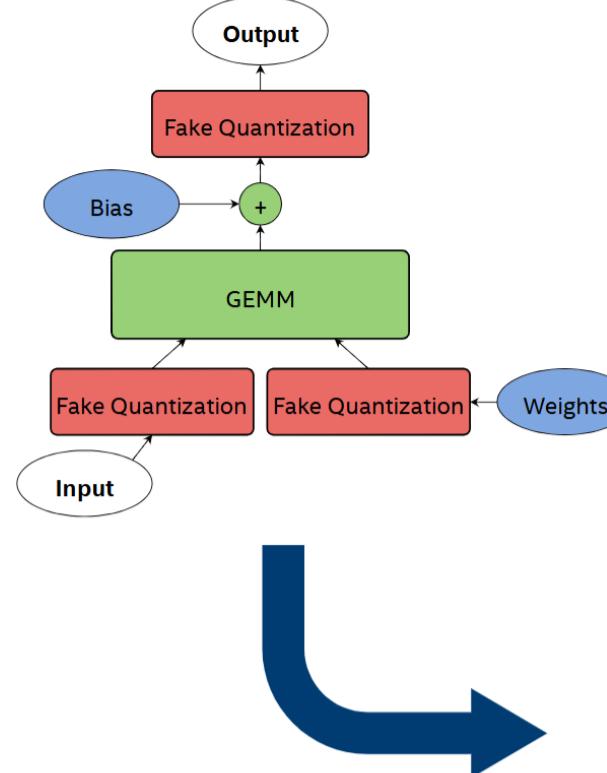
An affine mapping of **integers** to **real numbers**  $r = S(q - Z)$



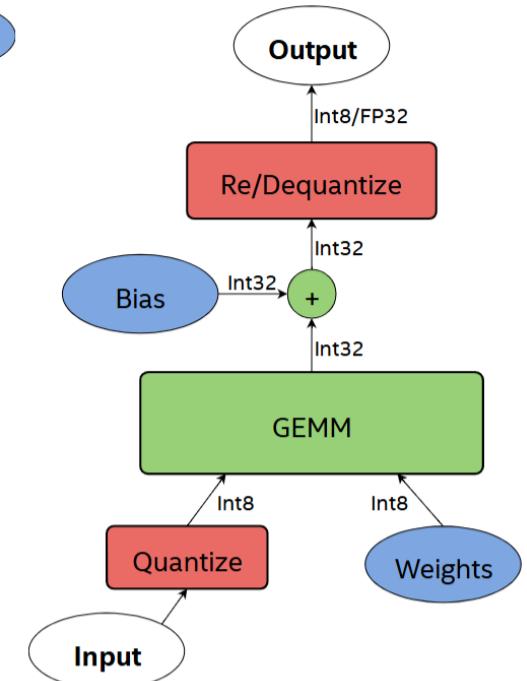
# Q8BERT

- Train Neural Networks (NN) to be **quantized-aware** at the inference stage.
- Fake quantization is used to introduce the quantization error training.
- Apply Fake Quantization on all General Matrix Multiply (GEMM) & Word/Position Embedding layers.
- Sensitive operations are kept in FP32 (Softmax, LN, GELU).

Training Graph

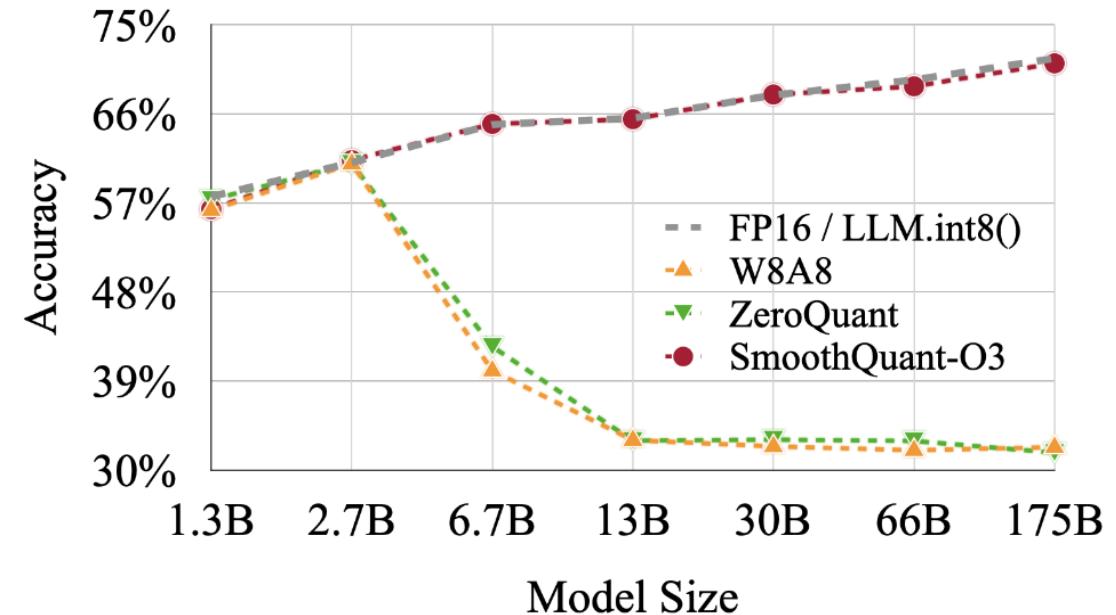


Inference Graph



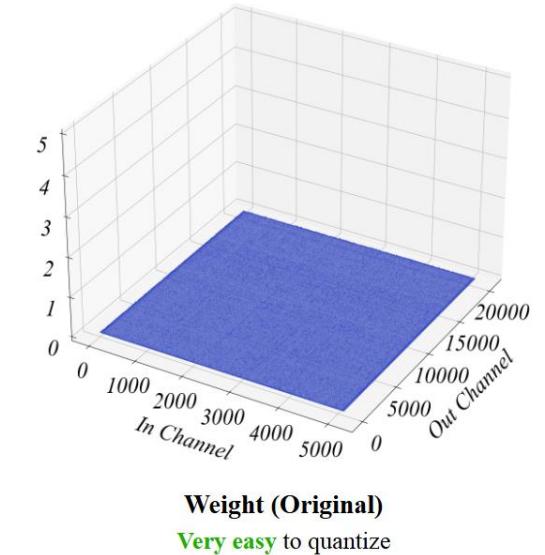
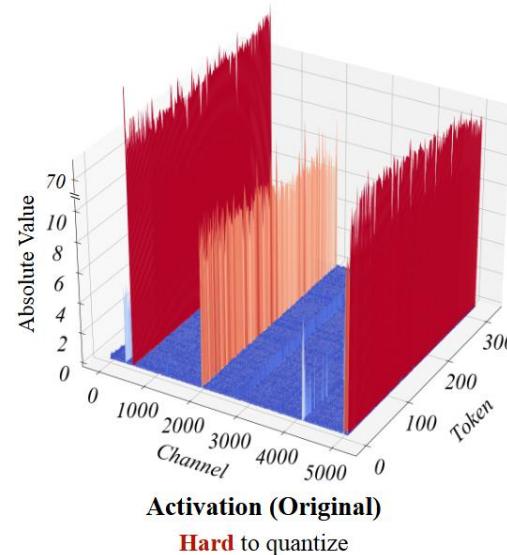
# Quantization Difficulty of LLMs

- Systematic outliers emerge in activations when scaling up to LLMs **beyond 6.7B**.
- Naive quantization methods will destroy the accuracy.



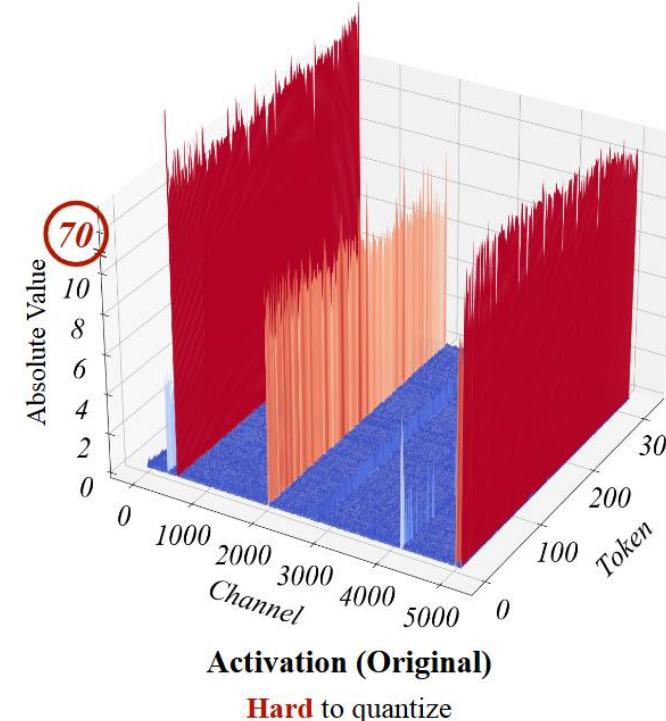
# Quantization Difficulty of LLMs

- Activations are harder to quantize than weights
  - Previous work has shown quantizing the weights of LLMs with INT8 or even INT4 doesn't degrade accuracy.

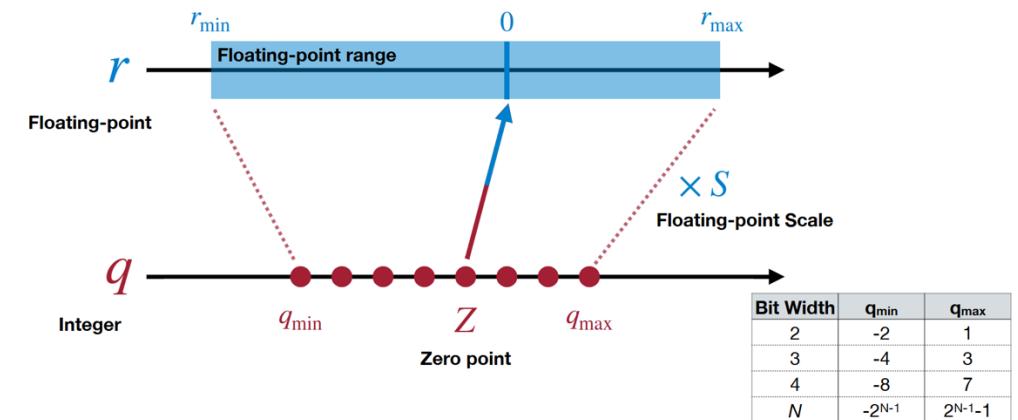


# Quantization Difficulty of LLMs

- Outliers make activation quantization difficult
  - The scale of outliers is **~100x larger** than most of other activations.
  - If we use INT8 quantization, normal activations will be **zero**.

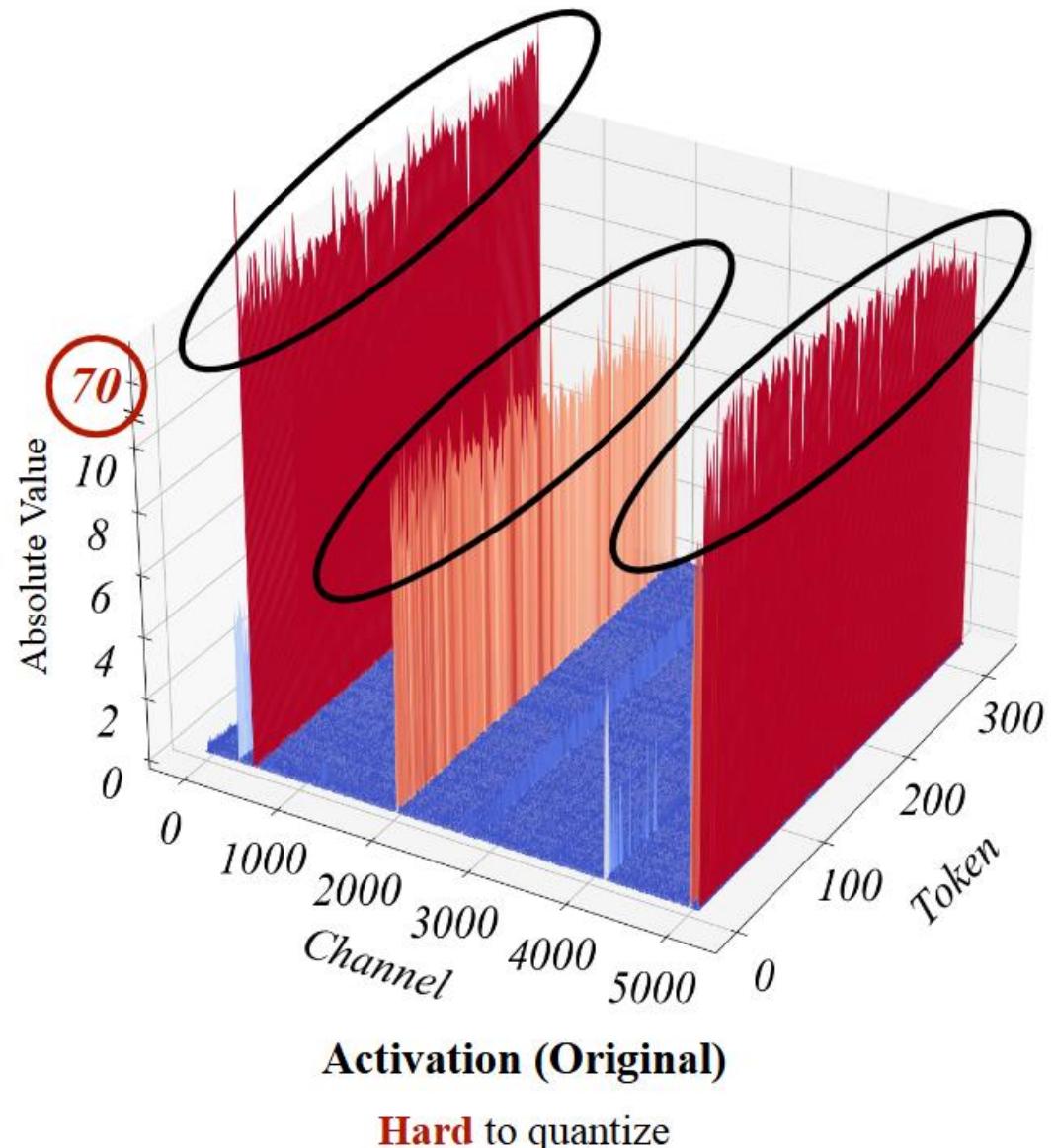


Activation (Original)  
Hard to quantize



# Quantization Difficulty of LLMs

- Outliers persist in fixed channels
  - Fixed channels have outliers, and the outlier channels are persistently large.



# SmoothQuant

- Activation Smoothing

- $s_j = \frac{\max(|X_j|)^\alpha}{\max(|W_j|)^{1-\alpha}}, j = 1, 2, \dots, C_i$
- $\hat{Y} = (X \text{diag}(s)^{-1}) \cdot (\text{diag}(s)W) = \hat{X}\hat{W}$

Original:

$$\begin{array}{c} X \\ \begin{array}{|c|c|c|c|} \hline 1 & -16 & 2 & 6 \\ \hline -2 & 8 & -1 & -9 \\ \hline 2 & 16 & 2 & 9 \\ \hline \end{array} \end{array} * \begin{array}{c} \xrightarrow{\text{Abs Max}} \\ \begin{array}{|c|c|c|} \hline 2 & 1 & -2 \\ \hline 1 & -1 & -1 \\ \hline 2 & -1 & -2 \\ \hline -1 & -1 & 1 \\ \hline \end{array} \end{array}$$

SmoothQuant:

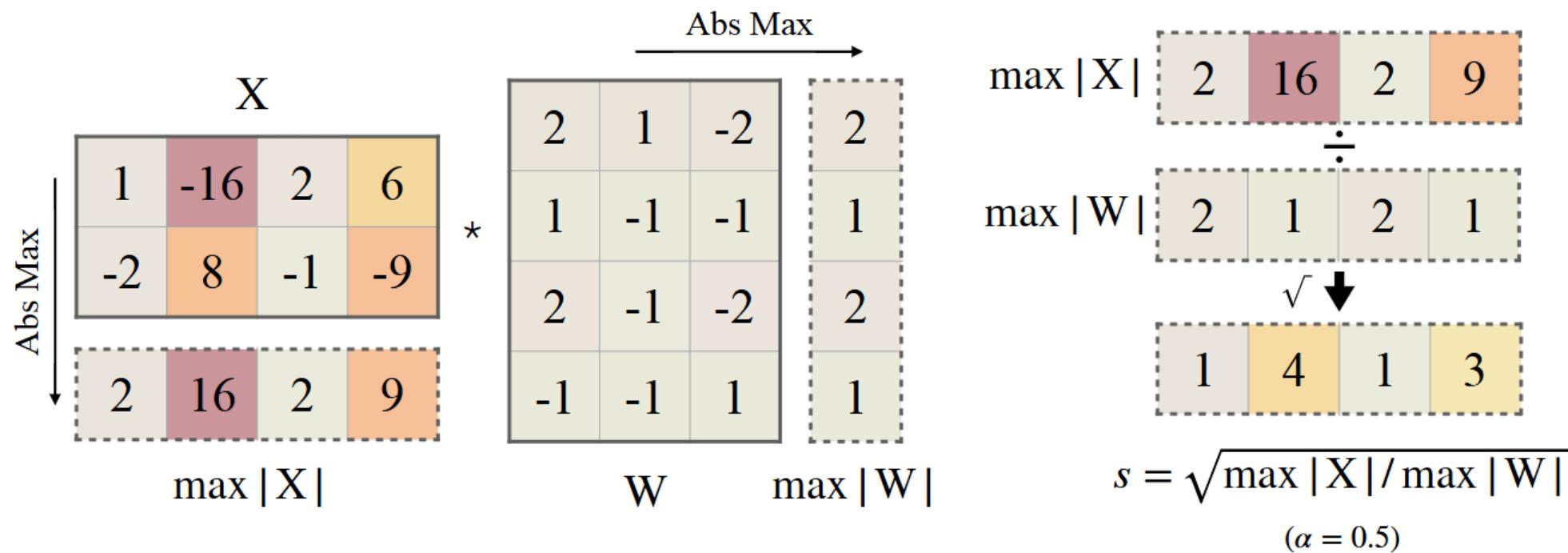
$$\begin{array}{c} \hat{X} = X \text{diag}(s)^{-1} \\ \begin{array}{|c|c|c|c|} \hline 1 & -4 & 2 & 2 \\ \hline -2 & 2 & -1 & -3 \\ \hline 1 & 4 & 1 & 3 \\ \hline \end{array} \end{array} * \begin{array}{c} \hat{W} = \text{diag}(s)W \\ \begin{array}{|c|c|c|} \hline 2 & 1 & -2 \\ \hline 4 & -4 & -4 \\ \hline 2 & -1 & -2 \\ \hline -3 & -3 & 3 \\ \hline \end{array} \end{array}$$

$s = \sqrt{\max |X| / \max |W|}$

# SmoothQuant

- Activation Smoothing
  - Calibration Stage (Offline)
  - $s_j = \frac{\max(|X_j|)^\alpha}{\max(|W_j|)^{1-\alpha}}, j = 1, 2, \dots, C_i$
  - $\hat{Y} = (X \text{diag}(s)^{-1}) \cdot (\text{diag}(s)W) = \hat{X}\hat{W}$

1. Calibration Stage (Offline):



# SmoothQuant

- Activation Smoothing
  - Smoothing Stage (Offline)
  - $s_j = \frac{\max(|X_j|)^\alpha}{\max(|W_j|)^{1-\alpha}}, j = 1, 2, \dots, C_i$
  - $\hat{Y} = (X \text{diag}(s)^{-1}) \cdot (\text{diag}(s)W) = \hat{X}\hat{W}$

## 2. Smoothing Stage (Offline):

$$X \begin{bmatrix} 1 & -16 & 2 & 6 \\ -2 & 8 & -1 & -9 \end{bmatrix} \div s \begin{bmatrix} 1 & 4 & 1 & 3 \end{bmatrix} = \hat{X} = X \text{diag}(s)^{-1} \begin{bmatrix} 1 & -4 & 2 & 2 \\ -2 & 2 & -1 & -3 \end{bmatrix}$$

divide the output channel of the previous layer by s

multiply the input channel of the following weight by s

$$\begin{bmatrix} 2 & 1 & -2 \\ 1 & -1 & -1 \\ 2 & -1 & -2 \\ -1 & -1 & 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 4 \\ 1 \\ 3 \end{bmatrix} = \hat{W} = \text{diag}(s)W \begin{bmatrix} 2 & 1 & -2 \\ 4 & -4 & -4 \\ 2 & -1 & -2 \\ -3 & -3 & 3 \end{bmatrix}$$

$W$        $s$        $\hat{W}$

# SmoothQuant

- Activation Smoothing
  - Inference (deployed model)
  - $s_j = \frac{\max(|X_j|)^\alpha}{\max(|W_j|)^{1-\alpha}}, j = 1, 2, \dots, C_i$
  - $\hat{Y} = (\hat{X} \text{diag}(s)^{-1}) \cdot (\text{diag}(s) W) = \hat{X} \hat{W}$

3. Inference (deployed model):

$$\begin{matrix} & \hat{X} \\ \begin{matrix} 1 & -4 & 2 & 2 \\ -2 & 2 & -1 & -3 \end{matrix} & \begin{matrix} * \\ \hat{W} \end{matrix} \end{matrix}$$

At runtime, the activations are smooth and easy to quantize

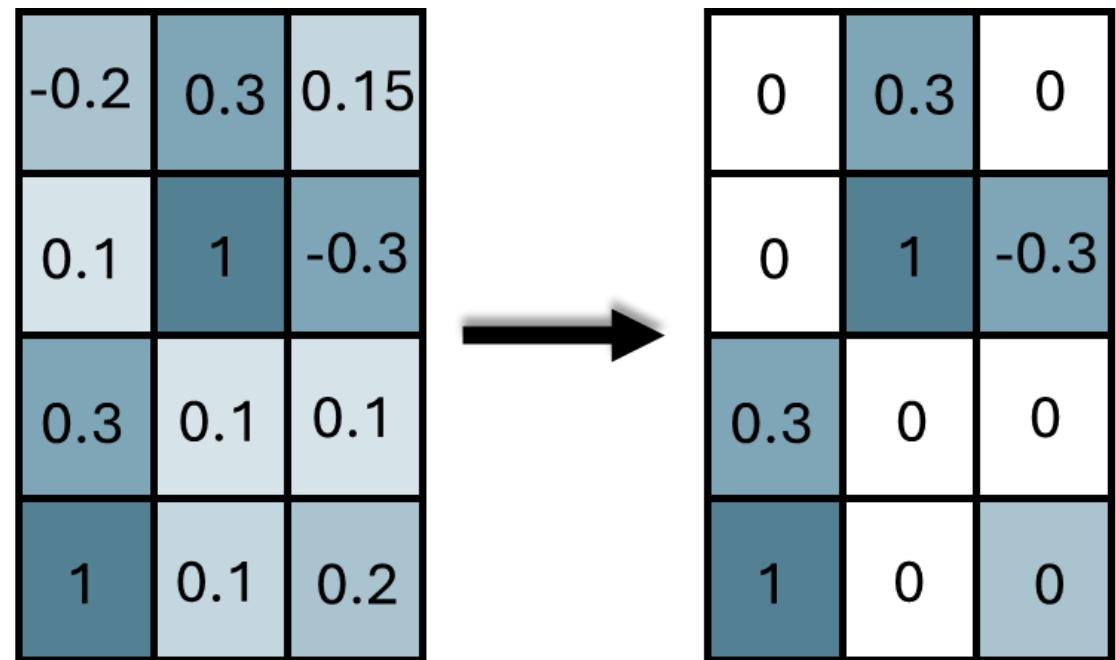
2	1	-2
4	-4	-4
2	-1	-2
-3	-3	3

# SmoothQuant

<i>OPT-175B</i>	LAMBADA	HellaSwag	PIQA	WinoGrande	OpenBookQA	RTE	COPA	Average↑	WikiText↓
FP16	74.7%	59.3%	79.7%	72.6%	34.0%	59.9%	88.0%	66.9%	10.99
W8A8	0.0%	25.6%	53.4%	50.3%	14.0%	49.5%	56.0%	35.5%	93080
ZeroQuant	0.0%*	26.0%	51.7%	49.3%	17.8%	50.9%	55.0%	35.8%	84648
LLM.int8()	74.7%	59.2%	79.7%	72.1%	34.2%	60.3%	87.0%	66.7%	11.10
Outlier Suppression	0.00%	25.8%	52.5%	48.6%	16.6%	53.4%	55.0%	36.0%	96151
SmoothQuant-O1	74.7%	59.2%	79.7%	71.2%	33.4%	58.1%	89.0%	66.5%	11.11
SmoothQuant-O2	75.0%	59.0%	79.2%	71.2%	33.0%	59.6%	88.0%	66.4%	11.14
SmoothQuant-O3	74.6%	58.9%	79.7%	71.2%	33.4%	59.9%	90.0%	66.8%	11.17

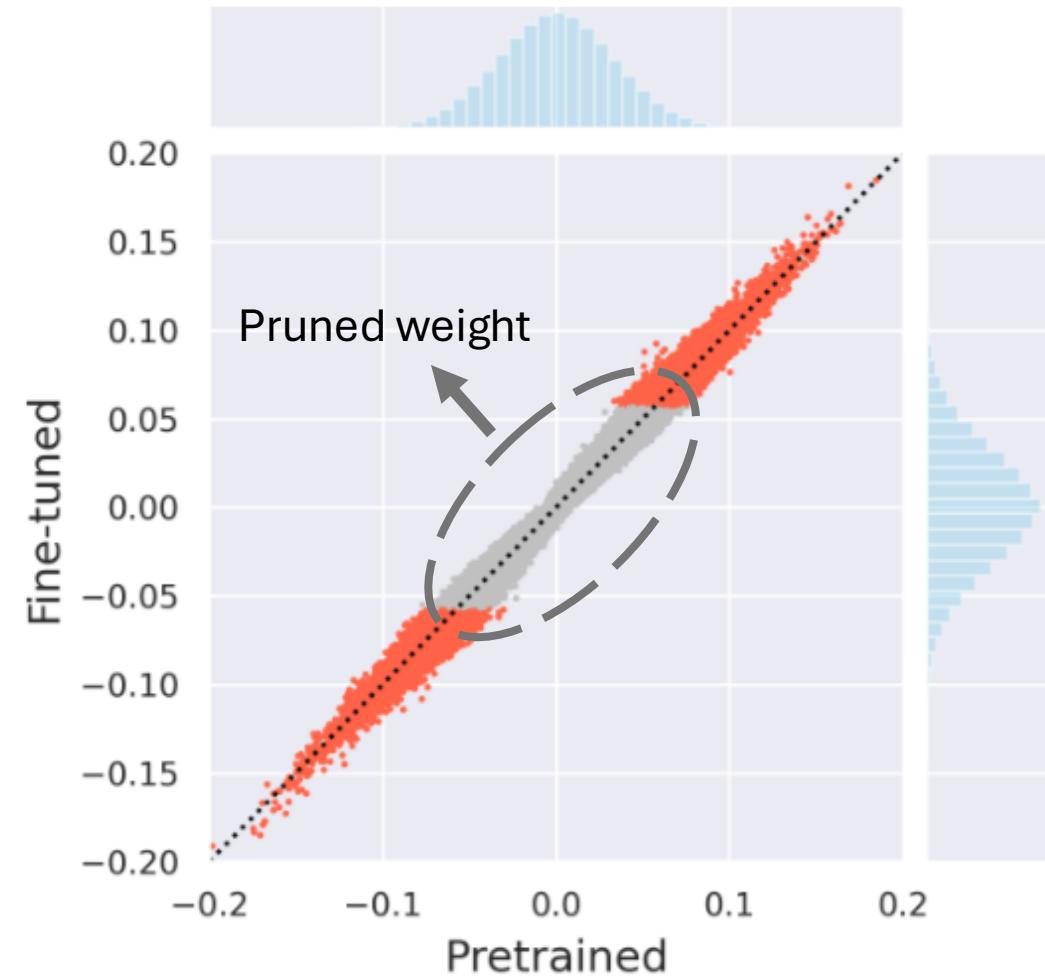
# Model Pruning

- Score-based pruning
  - Select the  $v\%$  highest values in  $S$  (important score)
  - $\text{Top}_v(S)_{i,j} = \begin{cases} 1, & S_{i,j} \text{ in top } v\% \\ 0, & \text{o.w.} \end{cases}$
- Magnitude pruning
  - Use  $|W_{i,j}|$  as the score.



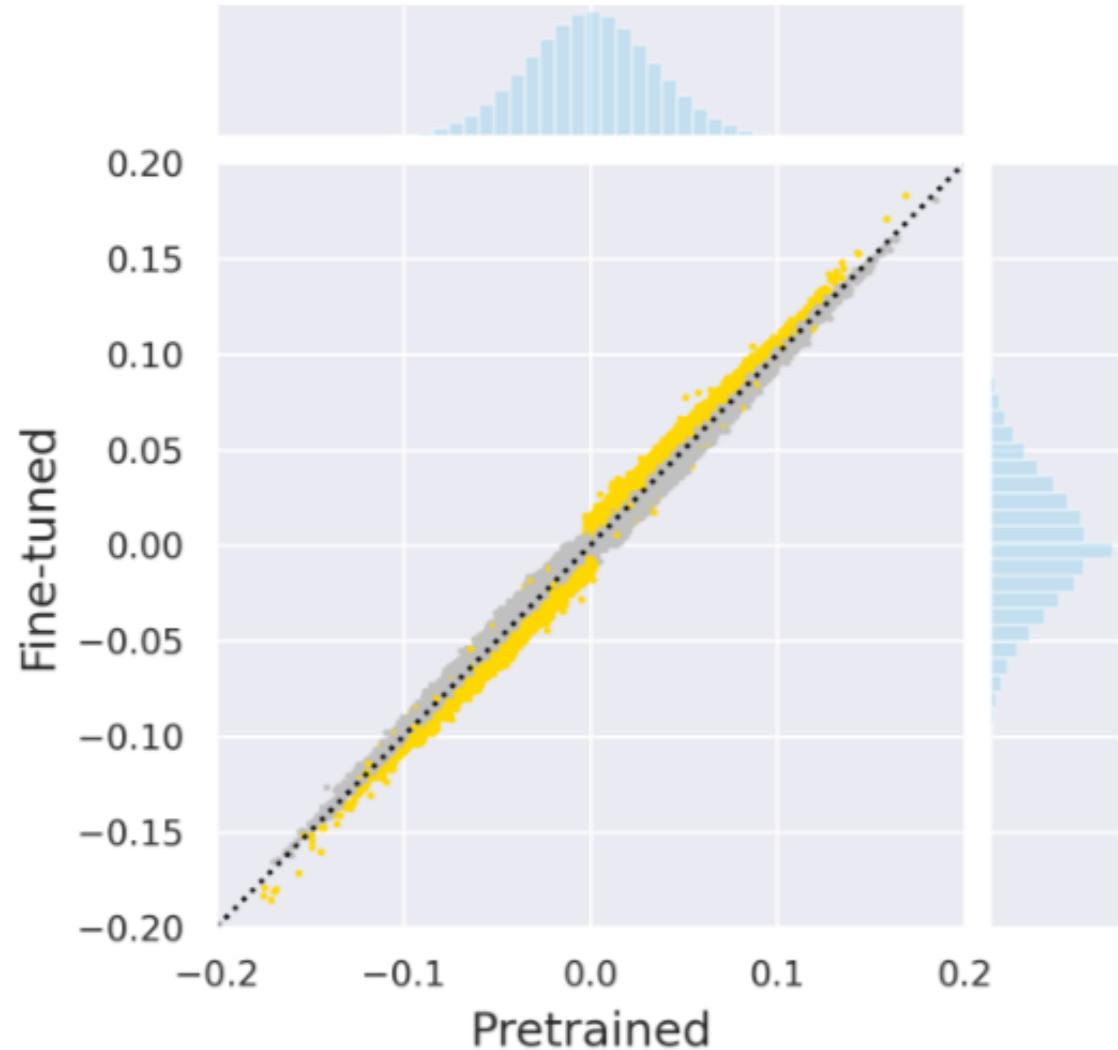
# Challenge of Magnitude Pruning

- The weights stay close to their pre-trained values which limits the adaptivity of magnitude pruning



# Movement Pruning

- Learn both the weights  $\mathbf{W}$  and their importance scores  $\mathbf{S}$ 
  - $a = (\mathbf{W} \odot \mathbf{M})x$ , where  $\mathbf{M} = \text{Top}_v(\mathbf{S})$ .
  - $\frac{\partial \mathcal{L}}{\partial S_{i,j}} = \frac{\partial \mathcal{L}}{\partial a_i} \frac{\partial a_i}{\partial S_{i,j}} = \frac{\partial \mathcal{L}}{\partial a_i} \mathbf{W}_{i,j} x_j$ .
- Movement pruning selects weights that are moving away from 0.



# Movement Pruning

	BERT base fine-tuned	Remaining Weights (%)	MaP	$L_0$ Regu	MvP	soft MvP
SQuAD - Dev EM/F1	80.4/88.1	10%	67.7/78.5	69.9/80.0	<b>71.9/81.7</b>	71.3/81.5
		3%	40.1/54.5	61.2/73.3	65.2/76.3	<b>69.5/79.9</b>
MNLI - Dev acc/MM acc	84.5/84.9	10%	77.8/79.0	77.9/78.5	79.3/79.5	<b>80.7/81.1</b>
		3%	68.9/69.8	75.1/75.4	76.1/76.7	<b>79.0/79.6</b>
QQP - Dev acc/F1	91.4/88.4	10%	78.8/75.1	87.5/81.9	89.1/85.5	<b>90.5/87.1</b>
		3%	72.1/58.4	86.5/81.0	85.6/81.0	<b>89.3/85.6</b>

Performance at high sparsity levels.

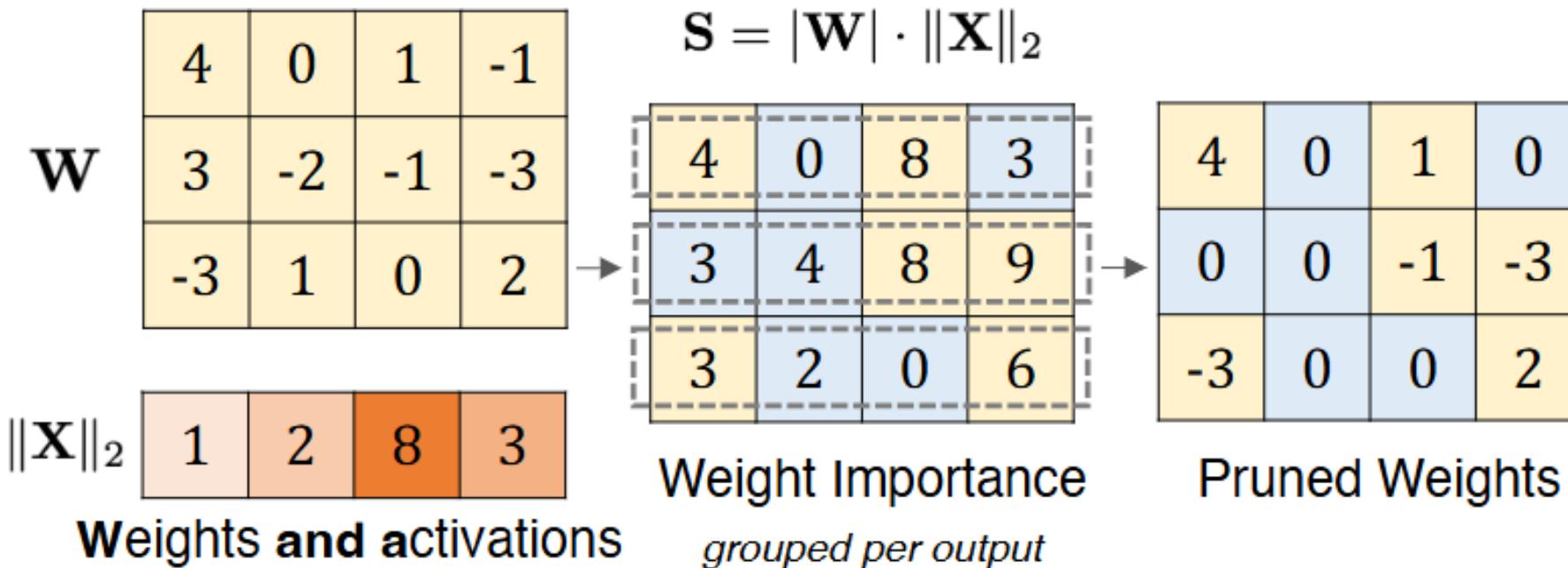
	BERT base fine-tuned	Remaining Weights (%)	MaP	$L_0$ Regu	MvP	soft MvP
SQuAD - Dev EM/F1	80.4/88.1	10%	70.2/80.1	72.4/81.9	75.6/84.3	<b>76.6/84.9</b>
		3%	45.5/59.6	64.3/75.8	67.5/78.0	<b>72.7/82.3</b>
MNLI - Dev acc/MM acc	84.5/84.9	10%	78.3/79.3	78.7/79.7	80.1/80.4	<b>81.2/81.8</b>
		3%	69.4/70.6	76.0/76.2	76.5/77.4	<b>79.5/80.1</b>
QQP - Dev acc/F1	91.4/88.4	10%	79.8/65.0	88.1/82.8	89.7/86.2	<b>90.2/86.8</b>
		3%	72.4/57.8	87.0/81.9	86.1/81.5	<b>89.1/85.5</b>

Distillation-augmented performances for selected high sparsity levels.

# Wanda

- Motivation:
  - $y = w_1x_1 + w_2x_2$
  - $|w_1| \leq |w_2|$ , but  $|x_1| \gg |x_2|$
- Pruning Metric:
  - $S_{ij} = |W_{ij}| \cdot \|X_j\|_2$

Wanda

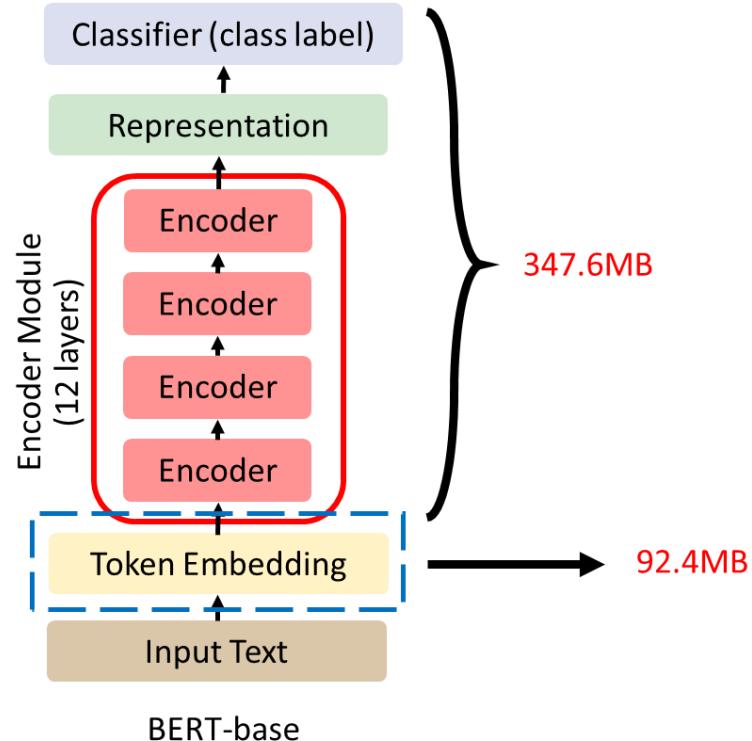


# Wanda

Method	Weight Update	Sparsity	LLaMA				LLaMA-2		
			7B	13B	30B	65B	7B	13B	70B
Dense	-	0%	59.99	62.59	65.38	66.97	59.71	63.03	67.08
Magnitude	✗	50%	46.94	47.61	53.83	62.74	51.14	52.85	60.93
SparseGPT	✓	50%	<b>54.94</b>	58.61	63.09	66.30	<b>56.24</b>	60.72	<b>67.28</b>
Wanda	✗	50%	54.21	<b>59.33</b>	<b>63.60</b>	<b>66.67</b>	<b>56.24</b>	<b>60.83</b>	67.03
Magnitude	✗	4:8	46.03	50.53	53.53	62.17	50.64	52.81	60.28
SparseGPT	✓	4:8	<b>52.80</b>	55.99	60.79	64.87	<b>53.80</b>	<b>59.15</b>	65.84
Wanda	✗	4:8	52.76	<b>56.09</b>	<b>61.00</b>	<b>64.97</b>	52.49	58.75	<b>66.06</b>
Magnitude	✗	2:4	44.73	48.00	53.16	61.28	45.58	49.89	59.95
SparseGPT	✓	2:4	<b>50.60</b>	<b>53.22</b>	58.91	62.57	<b>50.94</b>	54.86	63.89
Wanda	✗	2:4	48.53	52.30	<b>59.21</b>	<b>62.84</b>	48.75	<b>55.03</b>	<b>64.14</b>

Mean zero-shot accuracies on 7 zero-shot tasks  
of pruned LLaMA and LLaMA-2 models.

# Token Embedding Compression



- Token embedding matrix often occupies a significant proportion of the whole model.
- Existing compression methods usually do not focus on the token embedding matrix specifically.
- Approximation error of token embedding will accumulate in later layers.

	Accuracy	Compression Ratio	Task-agnostic	Model-agnostic
SVD	Low	Low	✓	✓
Quantization [Zafir O et al., 2019]	High	Low	✓	✓
DistilBERT [Sanh V et al., 2019]	High	✗	✓	✓
MobileBERT (SVD+KD) [Sun Z et al., 2019]	High	Medium	✓	✗
BinaryBERT (Quantization+KD) [Bai H et al., 2019]	Medium	High	✗	✓

# LightToken

- Problem Formulation

- Token embedding matrix  $X \in \mathbb{R}^{d \times N}$ .
- Compressed token embedding matrix  $\hat{X}$ .

High Level Idea:

$$X \approx \hat{X} = \underbrace{\mathbf{U}_{X,k} \mathbf{V}_{X,k}^T}_{\text{Low-rank matrix (rank } k\text{)}} + \underbrace{\text{Dec} \left( \text{Enc} \left( X - \mathbf{U}_{X,k} \mathbf{V}_{X,k}^T \right) \right)}_{\text{Approximation of the residual matrix}}$$

Decoder (Lightweight) Encoder

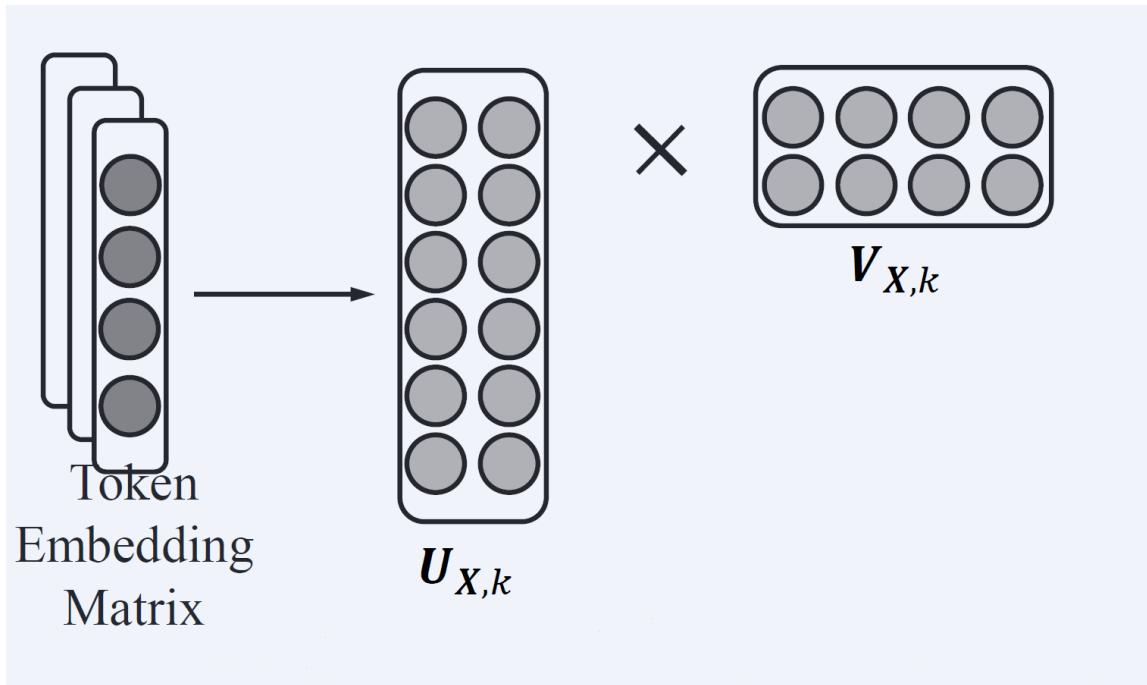
Binary codes

Low-rank matrix (rank  $k$ )

Approximation of the residual matrix

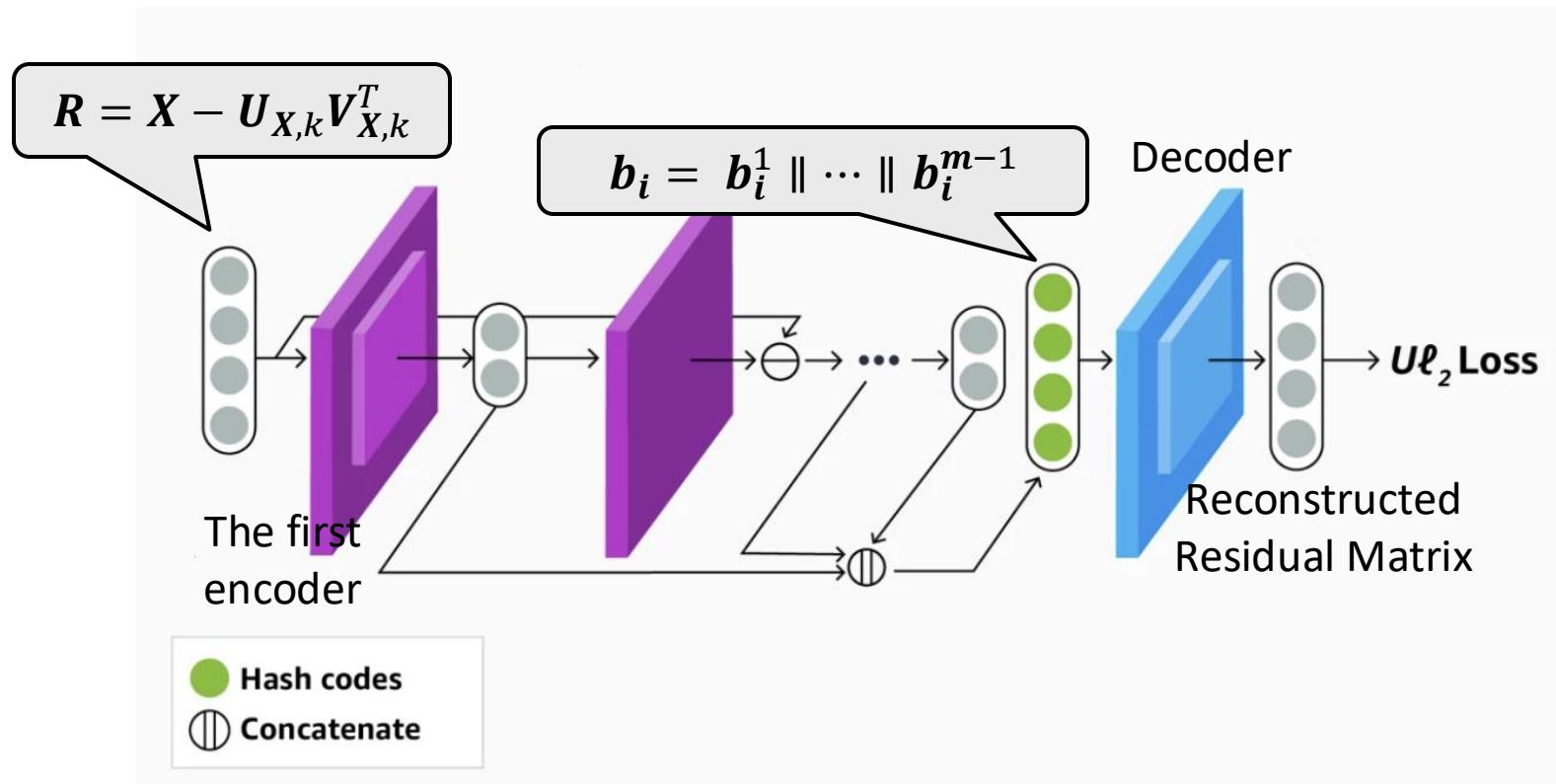
# LightToken

- Stage 1: Rank- $k$  SVD Approximation
  - *Goal: Extract coherent components*
  - $X = USV^T \approx U_{X,k} V_{X,k}^T$

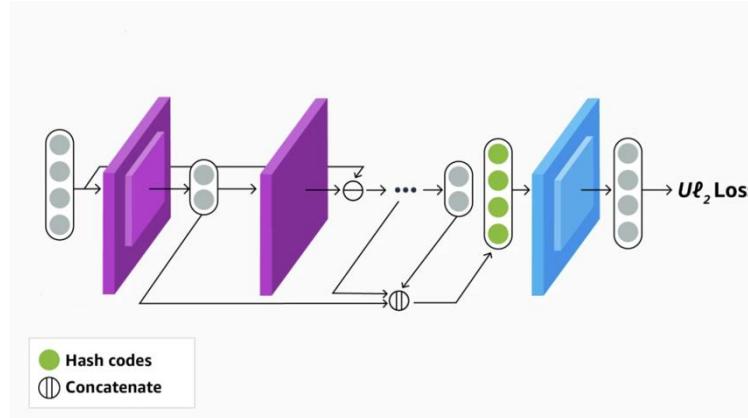


# LightToken

- Stage 2: Residual Binary Autoencoder
  - *Goal: Approximate the residual*



# LightToken



- Stage 2: Residual Binary Autoencoder

- Reconstructed Token Embedding

$$\hat{X} = \hat{R} + U_{X,k} V_{X,k}^T$$

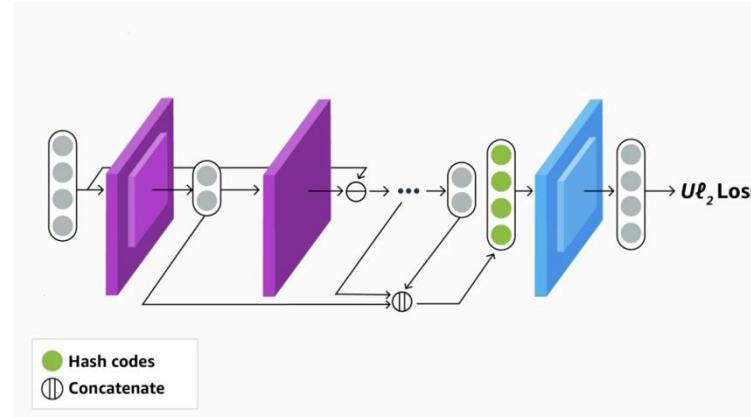
- Euclidean distance

$$\|X_i - \hat{X}_i\|_2^2 = \underbrace{\|X_i - \hat{X}_i^\parallel\|_2^2}_{\text{Error in parallel direction}} + \underbrace{\|X_i - \hat{X}_i^\perp\|_2^2}_{\text{Error in orthogonal direction}}$$

$$\|\hat{X}_i\|_2^2 \left( 1 - \left( \frac{\langle X_i, \hat{X}_i \rangle}{\|X_i\| \|\hat{X}_i\|} \right)^2 \right)$$

- Limitation: Provide too much freedom to learn the autoencoder parameters and may confuse/mislead the approximation learning

# LightToken



- Stage 2: Residual Binary Autoencoder
  - Upper Bound of Error in Orthogonal Direction

$$\|\hat{X}_i\|_2^2 \left( 1 - \left( \frac{\langle X_i, \hat{X}_i \rangle}{\|X_i\| \|\hat{X}_i\|} \right)^2 \right) \leq 2 \|\hat{X}_i\|_2^2 \left( 1 - \frac{\langle X_i, \hat{X}_i \rangle}{\|X_i\| \|\hat{X}_i\|} \right)$$

- The  $U\ell_2$  loss:

$$\sum_i \|X_i - \hat{X}_i\|_2^2 + \|\hat{X}_i\|_2^2 \left( 1 - \frac{\langle X_i, \hat{X}_i \rangle}{\|X_i\| \|\hat{X}_i\|} \right)^2$$

Euclidean  
distance

Weighted cosine  
similarity loss

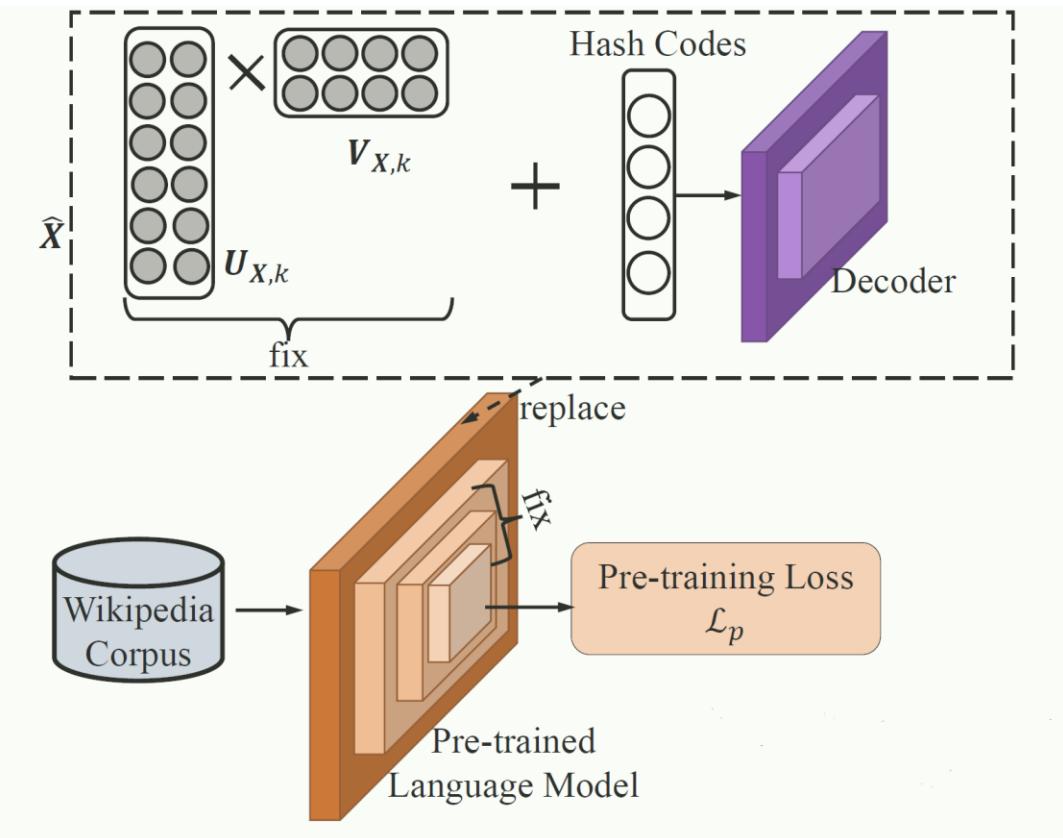
# LightToken

- Stage 3: Fine-tuning Decoder with Pre-training Objective

- Fine-tune the learned decoder on the subset of Wikipedia corpus

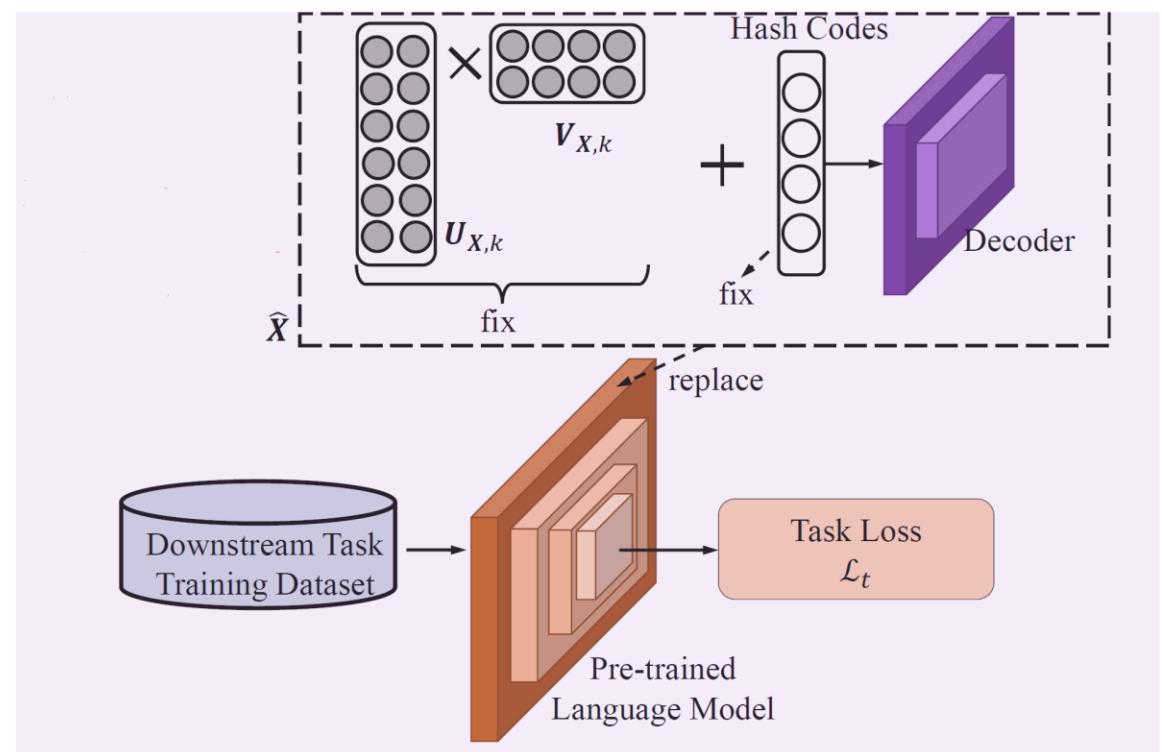
- $\min_{\Theta_D} \mathcal{L}_p (\hat{x}; B, \Theta_D, \Theta_{PLM})$

Decoder      Subset of      Parameters of  
parameters    Wiki corpus    the PLM except  
                                    token embedding

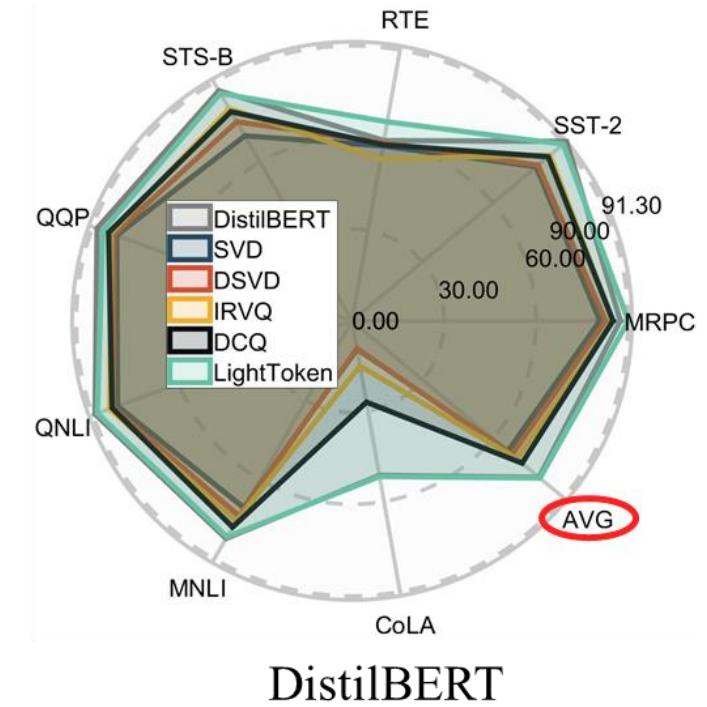
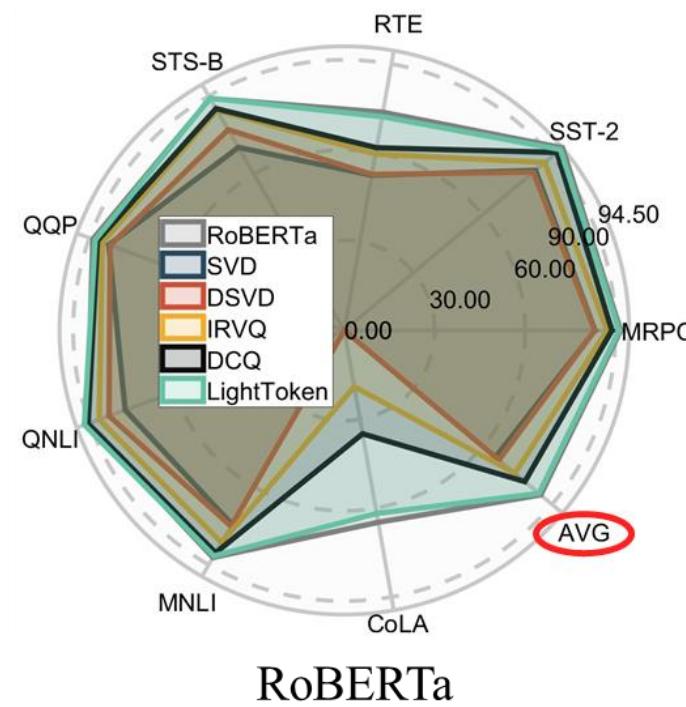
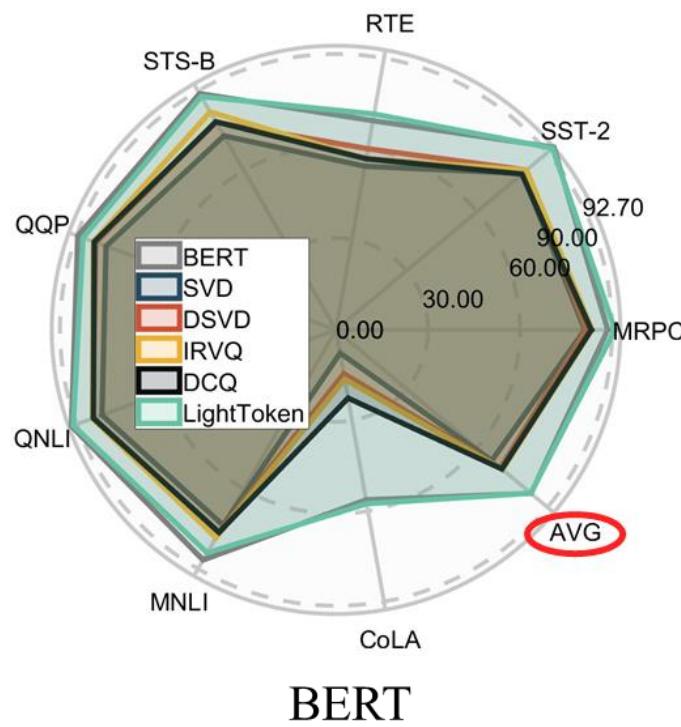


# LightToken

- Stage 4: Fine-tuning PLM for Downstream Tasks
  - Fine-tune the model with task-specific loss functions
  - $\min_{\Theta_{PLM}} \mathcal{L}_p (\mathfrak{X}; \mathcal{B}, \Theta_D, \Theta_{PLM})$



# LightToken



# Outline

- Model Efficiency
  - Parameter-efficient Fine-tuning
  - Model Compression
  - Inference Acceleration
- Data Efficiency
  - Zero/few-shot Learning
  - Knowledge Editing
  - Retrieval-augmented Generation

# Inference Acceleration

- **What is Inference Acceleration?**
  - Optimizing the speed and efficiency of the inference process without significant accuracy drop.
- **Why is it Important?**
  - Growing demand for deploying complex models.
  - Limitations of hardware capabilities.
  - Need for real-time responses in various applications.
- **Conditional Computation offers a flexible alternative:**
  - Allows **dynamic resource allocation per input**, reducing unnecessary computation.
  - Maintains strong accuracy while improving efficiency.

# Conditional Computation

## Core Idea

- Traditional inference uses all model parameters for every input.
- Conditional computation activates only necessary parts of the model per input.
  - Saves computation by dynamically adjusting model execution.
  - Enables trade-offs between accuracy and efficiency.
  - Can be applied at different levels: **layer-wise, token-wise, or expert selection.**
- Three Major Approaches:
  - **Mixture of Experts (MoE)** – Activate only a subset of specialized experts for each input.
  - **Early Exit** – Allow models to terminate inference early if confidence is high.
  - **Cascade Models** – Process simpler inputs with smaller models.

# Mixture of Experts (MoE) – Sparse Computation

**Key Idea:** Instead of using all parameters, activate only a subset of **specialized "expert" networks** per input.

## How it Works:

- A **router network** selects a few experts from a pool.
- Each input is processed by a small fraction of the total model.
- Allows scaling **without linearly increasing compute cost**.

## Advantages:

- Enables extremely large models **without proportional inference cost**.
- Improves efficiency while maintaining accuracy.

## Challenges:

- **Routing overhead:** Selecting the right experts introduces latency.
- **Load balancing:** Some experts may become over-utilized, leading to inefficiency.

## Example Models:

- Sparsely-Gated Mixture-of-Experts layer, Switch Transformer

# Sparsely-Gated Mixture-of-Experts layer (MoE)

- **Sparsely-Gated MoE Layer:** Uses a gating network to select a sparse subset of experts.
- **Top-k Routing:** Only the  $k$  most relevant experts are activated per input.
- **Conditional Computation:** Each input token is routed to a small fraction of the model, reducing computational cost.
- **Gradient Updates:** Experts receive gradients only when they are selected, leading to uneven training.
- **Load Balancing Loss:** Encourages more uniform expert utilization to avoid a few experts being overused.

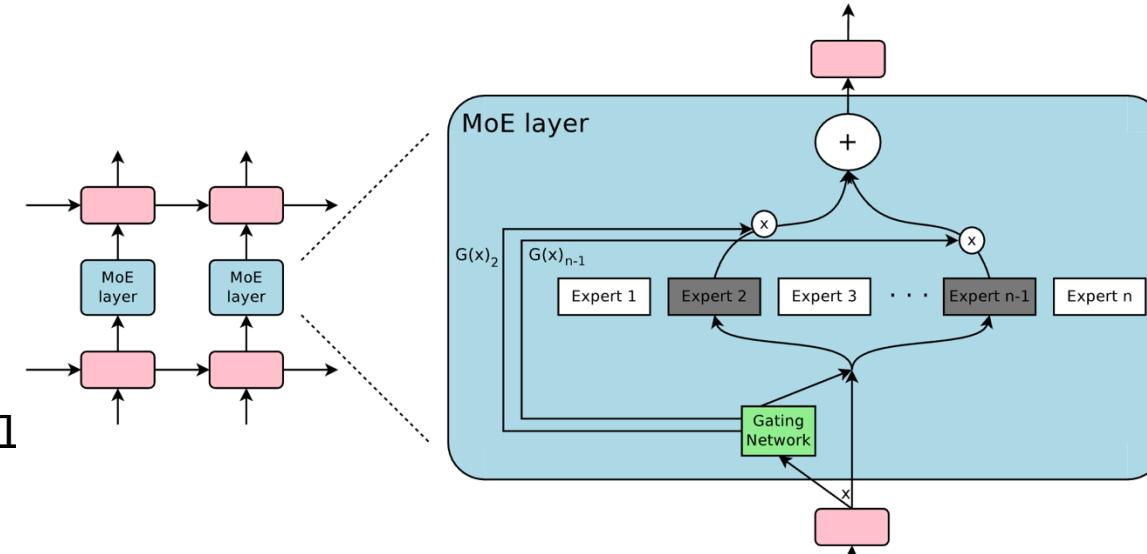


Figure 1: A Mixture of Experts (MoE) layer embedded within a recurrent language model. In this case, the sparse gating function selects two experts to perform computations. Their outputs are modulated by the outputs of the gating network.

# Sparsely-Gated Mixture-of-Experts layer (MoE)

## Contributions/Impacts

- **First scalable MoE model** in deep learning, allowing large models with minimal increase in computation.
- **Efficient Training:** Sparse activation significantly reduces FLOPs compared to dense models.
- **Introduced Expert Balancing:** Ensures all experts contribute meaningfully to learning.
- **Proved Feasibility of Sparse MoE:** Demonstrated that routing small subsets of the network maintains accuracy while reducing cost.

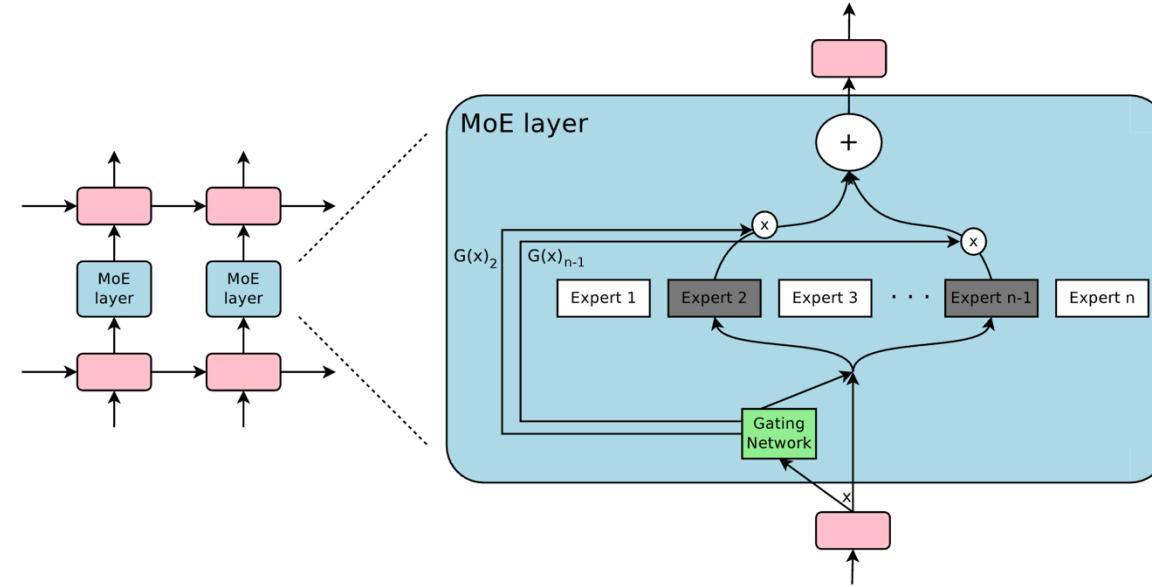


Figure 1: A Mixture of Experts (MoE) layer embedded within a recurrent language model. In this case, the sparse gating function selects two experts to perform computations. Their outputs are modulated by the outputs of the gating network.

# Sparsely-Gated Mixture-of-Experts layer (MoE)

- **Language Modeling**

**Performance:** MoE outperforms dense baselines while using only 10% of the computation complexity.

- **Scaling Benefits:** Larger models continue to improve without additional FLOPs, proving MoE's effectiveness at scaling.

Table 1: Summary of high-capacity MoE-augmented models with varying computational budgets, vs. best previously published results ([Jozefowicz et al., 2016](#)). Details in Appendix C.

	Test Perplexity 10 epochs	Test Perplexity 100 epochs	#Parameters excluding embedding and softmax layers	ops/timestep	Training Time 10 epochs	TFLOPS /GPU
Best Published Results	34.7	30.6	151 million	151 million	59 hours, 32 k40s	1.09
Low-Budget MoE Model	34.1		4303 million	8.9 million	15 hours, 16 k40s	0.74
Medium-Budget MoE Model	31.3		4313 million	33.8 million	17 hours, 32 k40s	1.22
High-Budget MoE Model	<b>28.0</b>		4371 million	142.7 million	47 hours, 32 k40s	<b>1.56</b>

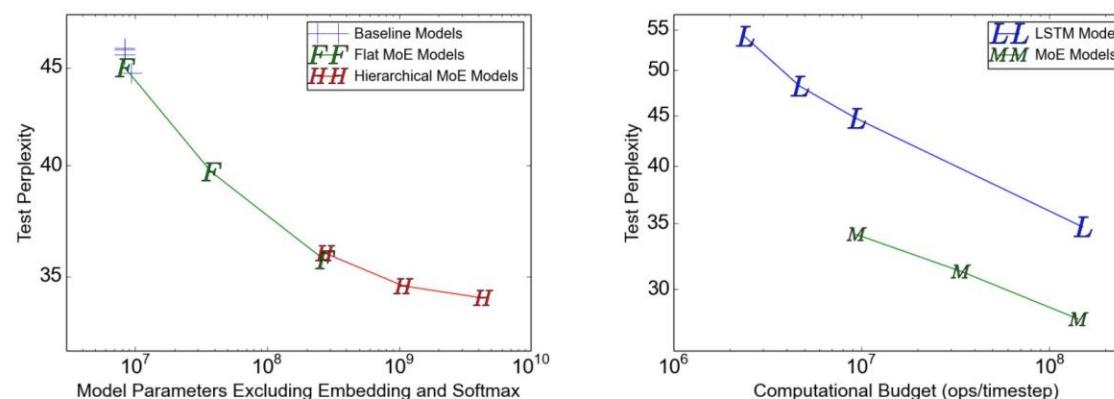
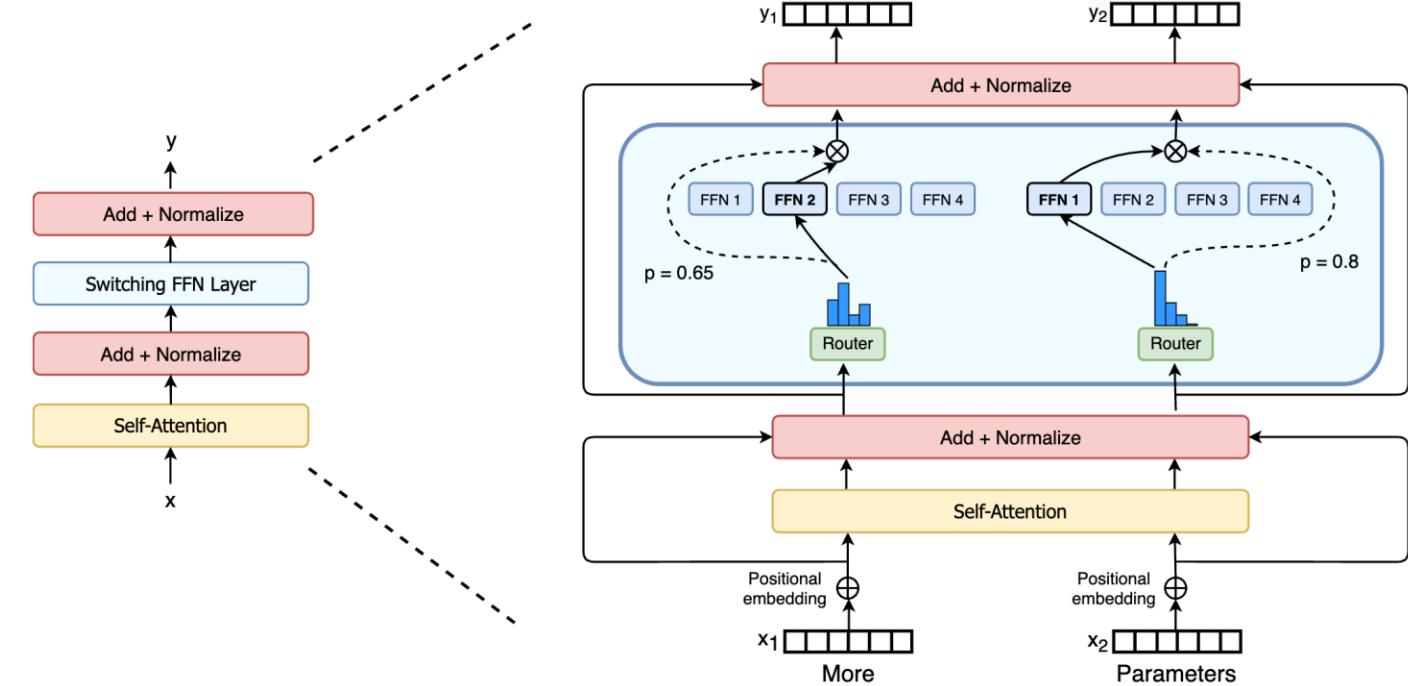


Figure 2: Model comparison on 1-Billion-Word Language-Modeling Benchmark. On the left, we plot test perplexity as a function of model capacity for models with similar computational budgets of approximately 8-million-ops-per-timestep. On the right, we plot test perplexity as a function of computational budget. The top line represents the LSTM models from ([Jozefowicz et al., 2016](#)). The bottom line represents 4-billion parameter MoE models with different computational budgets.

# Switch Transformer

- **Single Expert Routing:** Unlike previous MoEs, Switch Transformer activates **only one expert per token**, reducing communication overhead.
- **Improved Load Balancing Loss:** Ensures all experts receive roughly equal workloads.
- **Pretrained on Large Datasets:** Used **Colossal Clean Crawled Corpus (C4)** and other large datasets.
- **Gradient-Based Expert Selection:** Dynamically selects experts based on learned routing signals.



# Switch Transformer

- **Simplified MoE with lower communication overhead,** making large-scale deployment practical.
- **Showed trillion-parameter models are feasible,** proving that larger models continue to improve performance.
- **Improved Routing Mechanisms:** Reduced the computational waste of prior MoE implementations.
- **Outperformed Dense Transformers:** Achieved better accuracy while using only **30% of the computation** compared to T5.

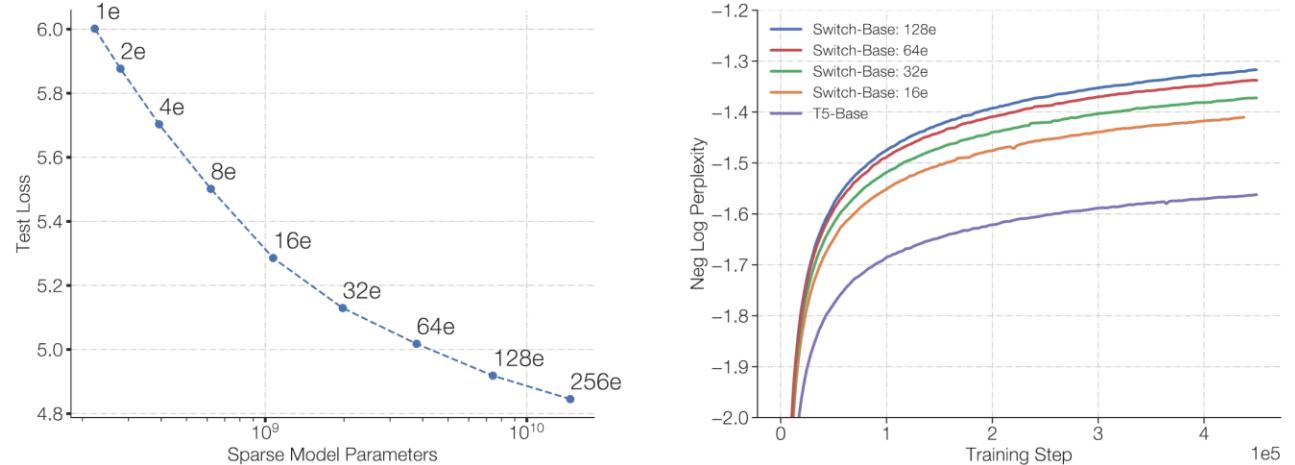


Figure 1: Scaling and sample efficiency of Switch Transformers. Left Plot: Scaling properties for increasingly sparse (more experts) Switch Transformers. Right Plot: Negative log perplexity comparing Switch Transformers to T5 ([Raffel et al., 2019](#)) models using the same compute budget.

# Switch Transformer

- **Scaling Experiments:** Showed **linear scaling laws** in NLP tasks with increasing model size.
- **Efficiency vs. Accuracy:** Demonstrated that Switch Transformer uses less training time to achieve similar or even better performance.

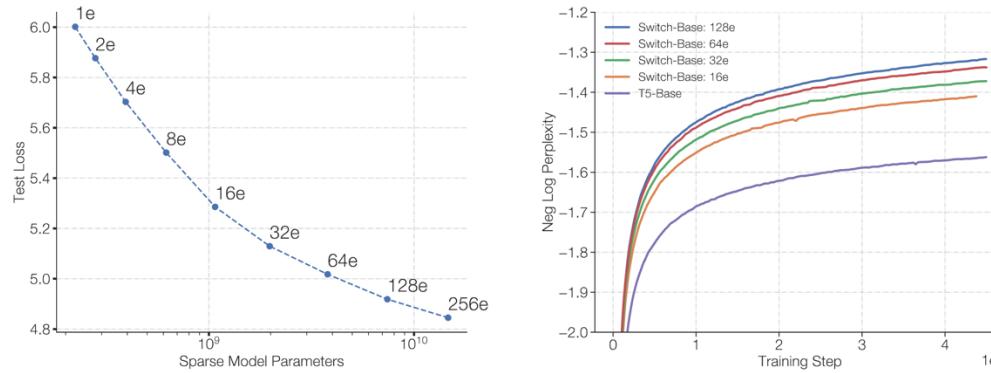


Figure 1: Scaling and sample efficiency of Switch Transformers. Left Plot: Scaling properties for increasingly sparse (more experts) Switch Transformers. Right Plot: Negative log perplexity comparing Switch Transformers to T5 ([Raffel et al., 2019](#)) models using the same compute budget.

Model	Capacity Factor	Quality after 100k steps ( $\uparrow$ ) (Neg. Log Perp.)	Time to Quality Threshold ( $\downarrow$ ) (hours)	Speed ( $\uparrow$ ) (examples/sec)
T5-Base	—	-1.731	Not achieved <sup>†</sup>	1600
T5-Large	—	-1.550	131.1	470
MoE-Base	2.0	-1.547	68.7	840
Switch-Base	2.0	-1.554	72.8	860
MoE-Base	1.25	-1.559	80.7	790
Switch-Base	1.25	-1.553	65.0	910
MoE-Base	1.0	-1.572	80.1	860
Switch-Base	1.0	-1.561	<b>62.8</b>	1000
Switch-Base+	1.0	<b>-1.534</b>	67.6	780

# Early Exit – Adaptive Depth Computation

**Key Idea:** Allow models to **terminate inference early** if they achieve high confidence before reaching the final layers.

## How it Works:

- Each layer has an **exit condition** that determines whether computation should continue.
- Simple inputs exit early, while harder cases pass through more layers.

## Advantages:

- Reduces computation on straightforward queries.
- Can improve latency without modifying model architecture.

## Challenges:

- **Reliable confidence estimation:** Hard to determine when to exit without degrading performance.
- **Layer-wise training complexity:** Requires well-calibrated thresholds.

## Example Applications:

- **BERT with early exit** for classification tasks.
- **Token-level early exit** in autoregressive decoding.

# DeeBERT

- This approach integrates **internal classifiers** at each layer of the BERT model, enabling inputs to exit the model early if a certain confidence threshold is met.
- Experiments demonstrate that DeeBERT can reduce inference time by approximately 40% with minimal impact on model quality.

	SST-2		MRPC		QNLI		RTE		QQP		MNLI-(m/mm)	
	Acc	Time	F <sub>1</sub>	Time	Acc	Time	Acc	Time	F <sub>1</sub>	Time	Acc	Time
<b>BERT-base</b>												
Baseline	93.6	36.72s	88.2	34.77s	91.0	111.44s	69.9	61.26s	71.4	145min	83.9/83.0	202.84s
DistilBERT	-1.4	-40%	-1.1	-40%	-2.6	-40%	-9.4	-40%	-1.1	-40%	-4.5	-40%
	-0.2	-21%	-0.3	-14%	-0.1	-15%	-0.4	-9%	-0.0	-24%	-0.0/-0.1	-14%
DeeBERT	-0.6	-40%	-1.3	-31%	-0.7	-29%	-0.6	-11%	-0.1	-39%	-0.8/-0.7	-25%
	-2.1	-47%	-3.0	-44%	-3.1	-44%	-3.2	-33%	-2.0	-49%	-3.9/-3.8	-37%
<b>RoBERTa-base</b>												
Baseline	94.3	36.73s	90.4	35.24s	92.4	112.96s	67.5	60.14s	71.8	152min	87.0/86.3	198.52s
LayerDrop	-1.8	-50%	-	-	-	-	-	-	-	-	-4.1	-50%
	+0.1	-26%	+0.1	-25%	-0.1	-25%	-0.6	-32%	+0.1	-32%	-0.0/-0.0	-19%
DeeBERT	-0.0	-33%	+0.2	-28%	-0.5	-30%	-0.4	-33%	-0.0	-39%	-0.1/-0.3	-23%
	-1.8	-44%	-1.1	-38%	-2.5	-39%	-1.1	-35%	-0.6	-44%	-3.9/-4.1	-29%

Table 1: Comparison between baseline (original BERT/RoBERTa), DeeBERT, and other acceleration methods. LayerDrop only reports results on SST-2 and MNLI. Time savings of DistilBERT and LayerDrop are estimated by reported model size reduction.

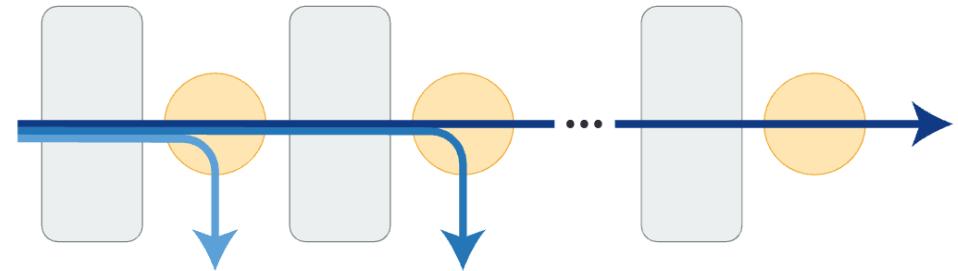
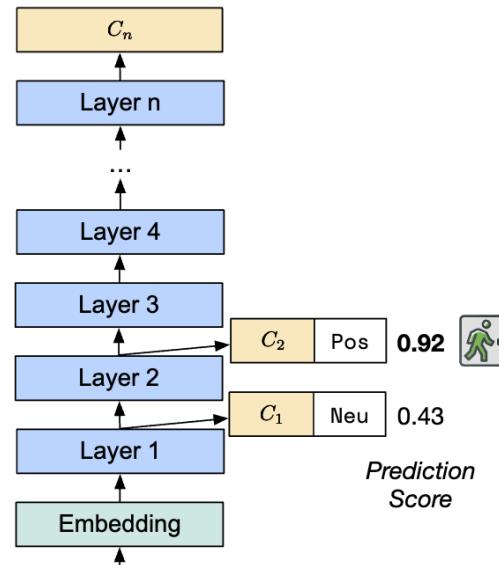


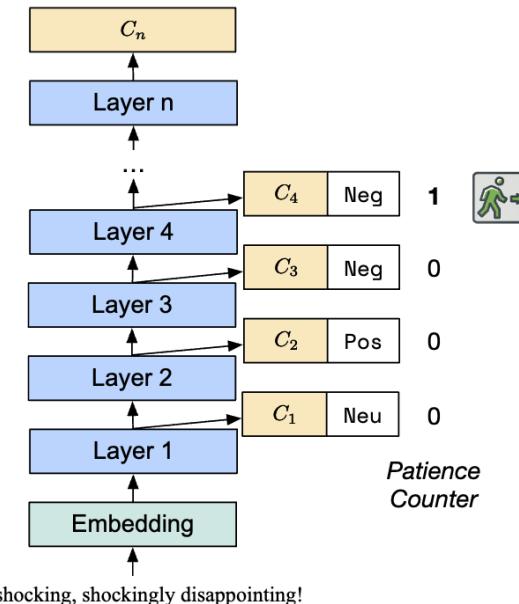
Figure 1: DeeBERT model overview. Grey blocks are transformer layers, orange circles are classification layers (off-ramps), and blue arrows represent inference samples exiting at different layers.

# PABEE

PABEE introduces a mechanism where an input exits the model once consecutive internal classifiers produce consistent predictions.



(a) Shallow-Deep Net [5]



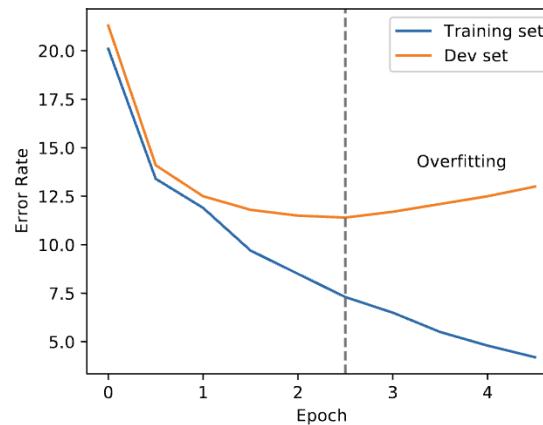
(b) Patience-based Early Exit (PABEE)

Figure 1: Comparison between Shallow-Deep Net, a prediction score based early exit (threshold is set to 0.9), and our Patience-based Early Exit (patience  $t = 1$ ). A classifier is denoted by  $C_i$ , and  $n$  is the number of layers in a model. In this figure, Shallow-Deep incorrectly exits based on the prediction score while PABEE considers multiple classifiers and exits with a correct prediction.

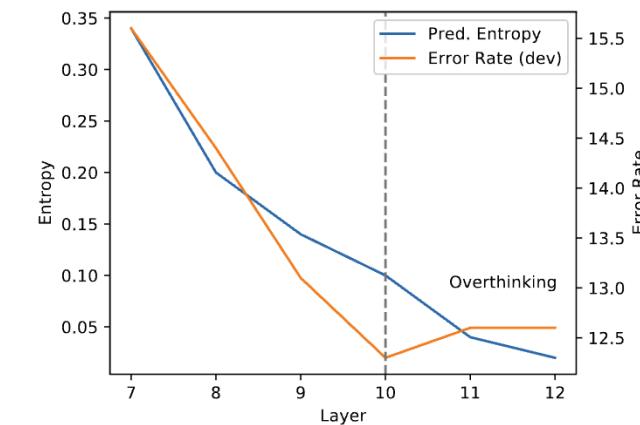
# PABEE

## Insights:

- Deep models sometimes “overthink” by using later layers that may introduce noise or overcomplicate the decision, leading to misclassifications. By stopping earlier, PABEE prevents this degradation in performance and can even improve accuracy.



(a) Overfitting in training

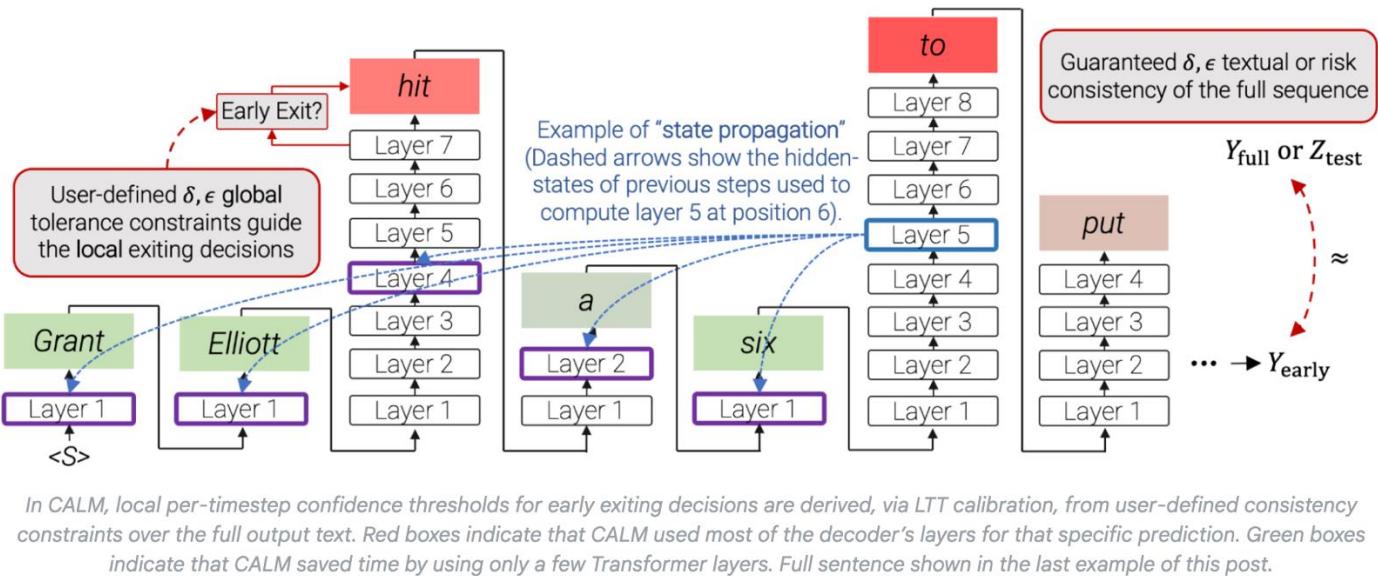


(b) Overthinking in inference

Figure 2: Analogy between overfitting in training and overthinking in inference. (a) In training, the error rate keeps going down on the training set but goes up later on the development set. (b) We insert a classifier after every layer. Similarly, the predicted entropy keeps dropping when more layers are added to inference but the error rate goes up after 10 layers. The results are obtained with ALBERT-base on MRPC.

# CALM

- CALM proposes a token-level early exit strategy for autoregressive language models.
- By estimating the confidence of predictions at each step, the model decides whether to proceed with further computations or output the current token

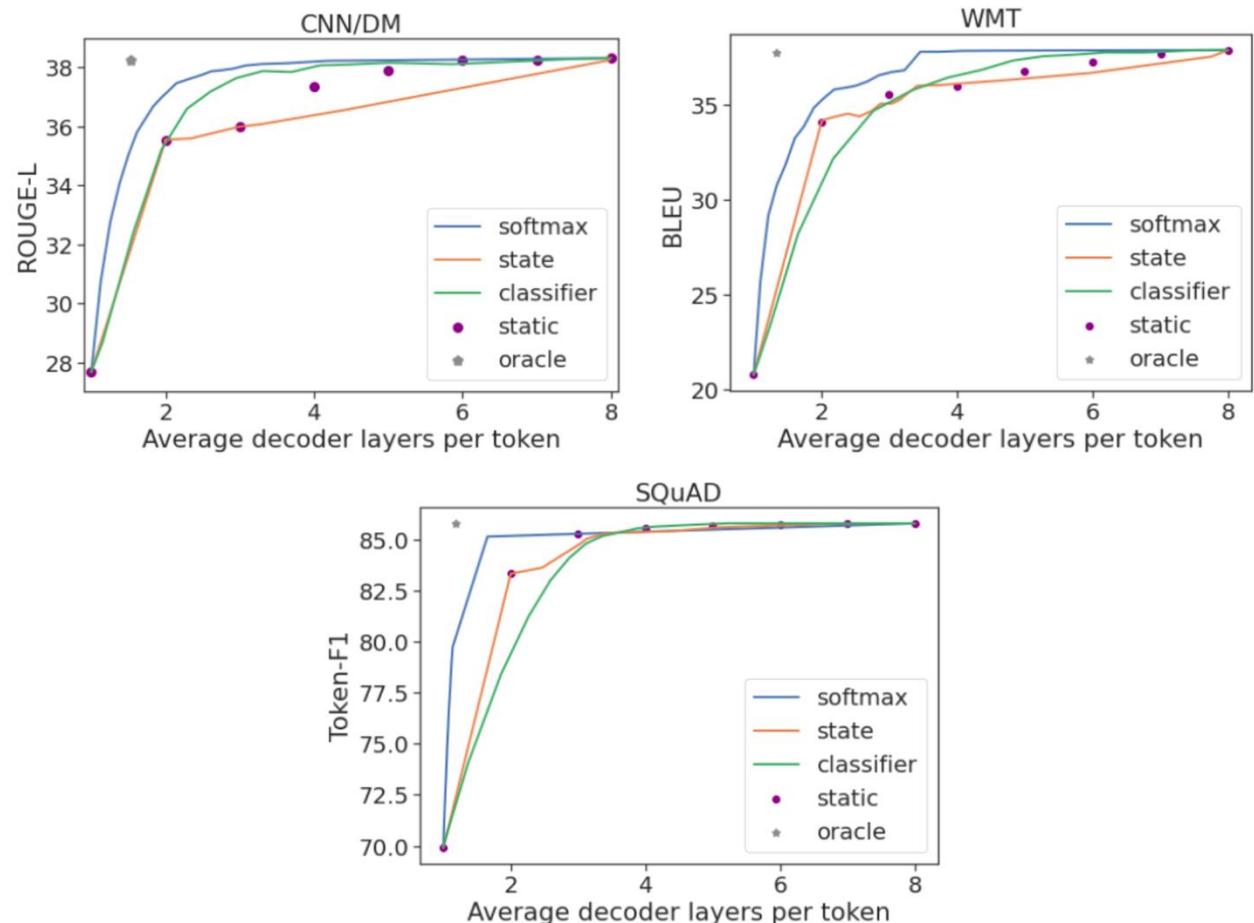


# CALM



## Insights:

- CALM distinguishes itself by introducing dynamic, confidence-calibrated early exiting specifically tailored for models with decoder.
- Unlike generic early exit methods that rely on fixed thresholds, CALM adjusts exit criteria on-the-fly during sequential generation, accounting for uncertainty and fluency as the sequence unfolds.
- This enables more fine-grained control over when to stop decoding, optimizing both inference speed and output quality without requiring significant architectural changes.



# HadSkip

HadSkip introduces an adaptive layer-skipping method to improve inference efficiency in pre-trained language models. Instead of sequentially processing all layers, it uses learnable gating mechanisms to skip layers dynamically based on a predefined budget.

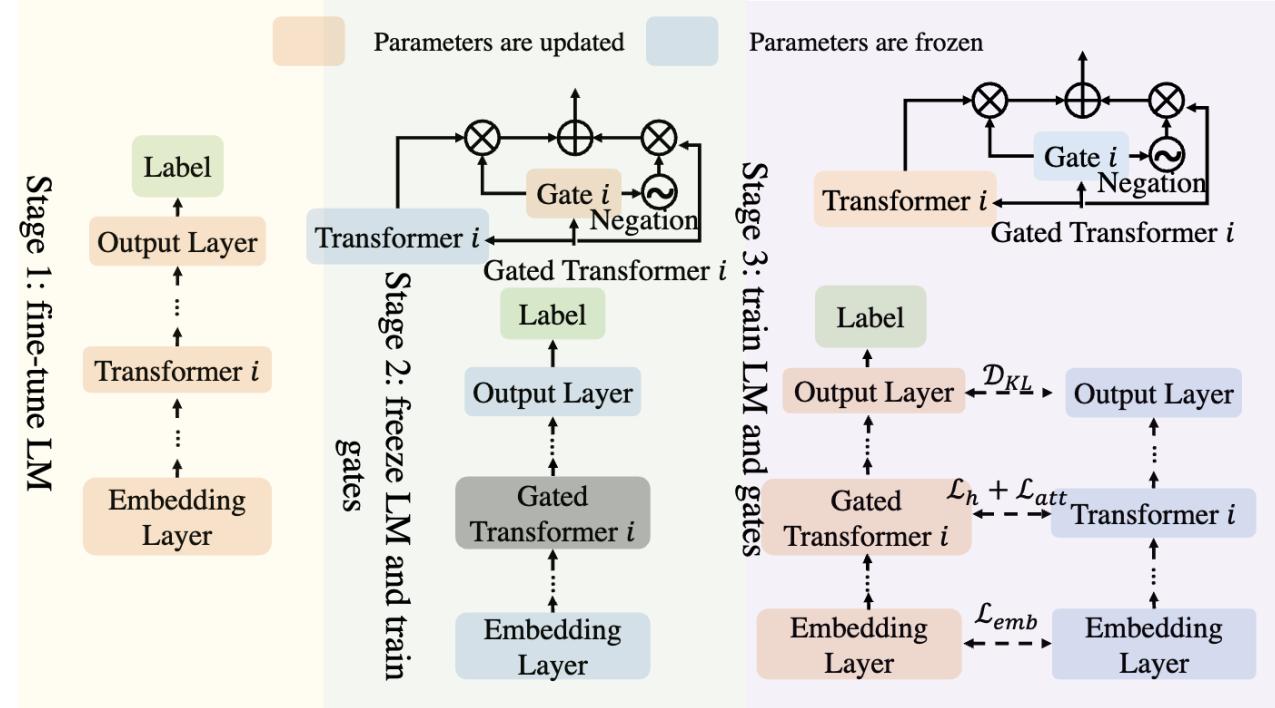


Figure 1: The framework of HadSkip.

# HadSkip

## Insights:

- HadSkip’s innovation lies in its learned, “hard-skip” mechanism that conditionally bypasses entire blocks of layers rather than terminating inference early.
- Unlike simple threshold methods, HadSkip uses a lightweight decision module trained jointly with the model to assess whether the current activations are sufficiently informative.

	CoLA	MNLI	MRPC	QNLI	QQP	RTE	SST-2	AVG
BERT*	59.4	84.3	91.3	88.5	90.4	69.3	92.5	82.2
Budgeted-Exit*	0.0	70.0	75.8	77.4	81.8	54.7	81.0	63.0
BranchyNet*	0.0	63.8	75.7	74.2	71.6	54.7	79.9	60.0
Shallow-Deep*	0.0	64.1	75.6	74.3	71.4	54.7	79.5	59.9
BERxiT*	0.0	63.5	75.6	73.3	68.2	55.3	79.5	59.3
PABEE*	0.0	63.9	75.8	73.6	68.6	55.8	79.9	59.7
PCEE-BERT*	9.8	73.4	78.8	80.4	79.6	58.4	83.6	66.3
HadSkip-BERT	35.6	78.1	84.2	84.1	88.5	54.9	88.5	73.4
Budgeted-Exit*	0.0	79.6	84.7	85.3	89.3	68.1	88.6	70.8
BranchyNet*	0.0	78.3	83.0	87.1	89.3	67.4	88.3	70.5
Shallow-Deep*	0.0	78.2	82.8	87.2	89.6	67.2	88.4	70.5
BERxiT*	12.3	78.4	82.9	87.0	89.1	67.3	88.3	72.2
PABEE*	0.0	78.9	83.1	87.2	89.6	67.7	88.7	70.7
PCEE-BERT*	23.2	80.1	84.8	87.1	90.8	69.4	90.4	75.1
HadSkip-BERT	50.8	82.4	85.9	88.7	90.5	62.5	90.0	78.7
Budgeted-Exit*	51.9	83.0	87.0	88.4	90.3	69.0	91.2	80.1
BranchyNet*	52.1	83.0	85.8	89.3	90.1	68.0	91.2	79.9
Shallow-Deep*	52.3	82.9	85.7	89.3	90.1	67.8	91.2	79.9
BERxiT*	52.2	83.2	86.2	89.6	90.1	68.1	91.4	80.1
PABEE*	52.4	83.4	86.1	89.8	90.4	68.3	91.7	80.3
PCEE-BERT*	52.8	83.4	86.8	90.5	91.2	69.7	91.8	80.9
HadSkip-BERT	58.3	83.8	88.4	90.7	90.8	66.8	92.7	81.6

Table 1: Performance comparison on GLUE dataset. Blue, green, and orange shadow parts represent methods with 3, 6, and 9-layer budgets respectively. “AVG” notes the average performance of the 7 tasks. Results with “\*” are taken from [\(Zhang et al., 2022\)](#).

# HadSkip

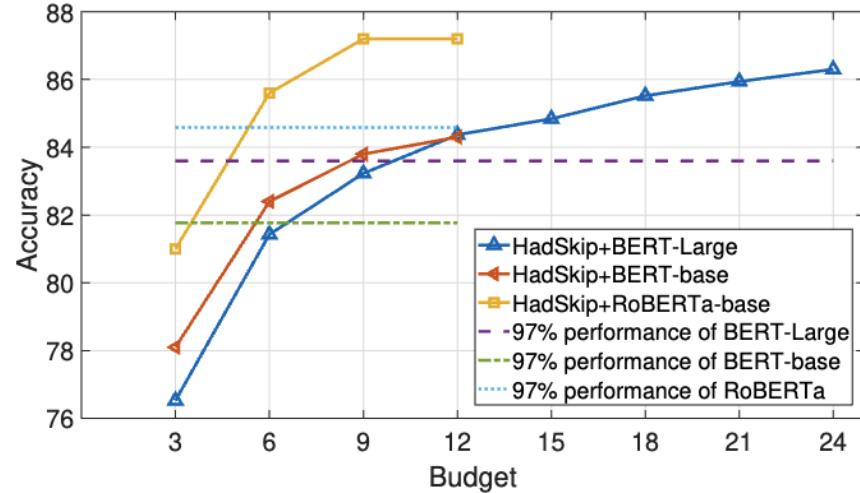


Figure 6: Accuracy with different budgets.

Observation show that with a lower budget, the model selectively utilizes specific layers more frequently, while at higher budgets, layer usage becomes more evenly distributed.

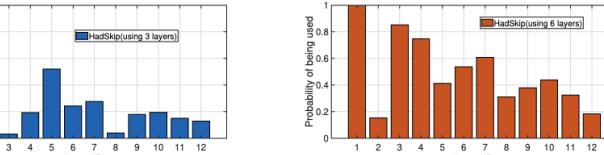
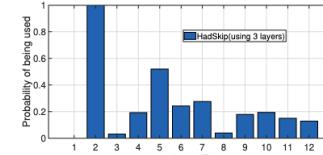


Figure 2: Visualization of the probability of each layer being used on MNLI dataset.

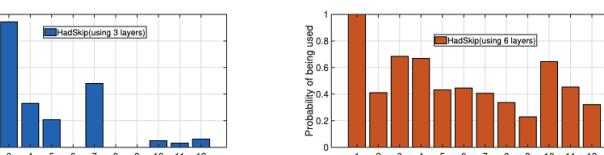
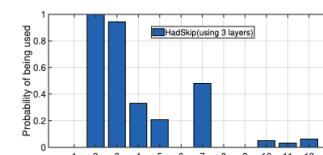


Figure 3: Visualization of the probability of each layer being used on QNLI dataset.

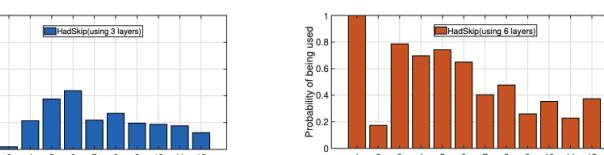
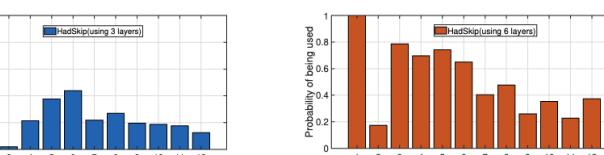
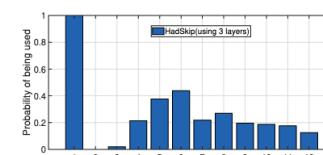


Figure 4: Visualization of the probability of each layer being used on SST2 dataset.

# Cascade Models – Adaptive Complexity

**Key Idea:** Use a sequence of models or modules that progressively increase in complexity.

- Simple cases **exit early** using a small model.
- More complex cases **pass through additional processing layers**.

**Advantages:**

- Saves computation on easy tasks.
- Ensures high accuracy for complex inputs.

**Challenges:**

- Designing optimal routing mechanisms.
- Avoiding cascading errors in progressive inference.

# Speculative Decoding

## What is Speculative Decoding?

- A technique to accelerate LLM inference by **predicting multiple tokens at once**, reducing latency.
- Uses a **draft model** to generate multiple candidate tokens before verifying them with a **larger model**.

# Speculative Decoding

## What is Speculative Decoding?

- A technique to accelerate LLM inference by **predicting multiple tokens at once**, reducing latency.
- Uses a **draft model** to generate multiple candidate tokens before verifying them with a **larger model**.

## Why is it Useful?

- Standard autoregressive decoding is **sequential and slow**.
- **Parallel token generation** improves efficiency without sacrificing output quality.

# Speculative Decoding

## What is Speculative Decoding?

- A technique to accelerate LLM inference by **predicting multiple tokens at once**, reducing latency.
- Uses a **draft model** to generate multiple candidate tokens before verifying them with a **larger model**.

## Why is it Useful?

- Standard autoregressive decoding is **sequential and slow**.
- **Parallel token generation** improves efficiency without sacrificing output quality.

## How It Works?

- **Draft Model** proposes a batch of candidate tokens.
- **Final Model** verifies and accepts or rejects them.
- **Accepted tokens** advance immediately, while rejected ones are recomputed.

# Speculative Decoding

## Key Contributions

- **Inference Speedup:** Achieves up to **2-3× faster inference** for large-scale transformers.
- **Maintains Output Quality:** Preserves accuracy while reducing latency.
- **Generalizable:** Works across different transformer architectures and model sizes.

## Insights

- **Better Draft Models = Higher Efficiency**
  - If the draft model closely matches the final model, fewer tokens need to be recomputed.
- **Trade-off Between Speed & Quality**
  - A larger batch size increases speed but may introduce more rejected tokens.
- **Compatible with Hardware Acceleration**
  - Can be combined with quantization, TPUs, and GPU optimizations.

Table 2. Empirical results for speeding up inference from a T5-XXL 11B model.

TASK	$M_q$	TEMP	$\gamma$	$\alpha$	SPEED
ENDE	T5-SMALL ★	0	7	0.75	<b>3.4X</b>
ENDE	T5-BASE	0	7	0.8	2.8X
ENDE	T5-LARGE	0	7	0.82	1.7X
ENDE	T5-SMALL ★	1	7	0.62	<b>2.6X</b>
ENDE	T5-BASE	1	5	0.68	2.4X
ENDE	T5-LARGE	1	3	0.71	1.4X
CNNDM	T5-SMALL ★	0	5	0.65	<b>3.1X</b>
CNNDM	T5-BASE	0	5	0.73	3.0X
CNNDM	T5-LARGE	0	3	0.74	2.2X
CNNDM	T5-SMALL ★	1	5	0.53	<b>2.3X</b>
CNNDM	T5-BASE	1	3	0.55	2.2X
CNNDM	T5-LARGE	1	3	0.56	1.7X

# Outline

- Model Efficiency
  - Parameter-efficient Fine-tuning
  - Model Compression
  - Inference Acceleration
- Data Efficiency
  - Zero/few-shot Learning
  - Knowledge Editing
  - Retrieval-augmented Generation

# Background & Motivation

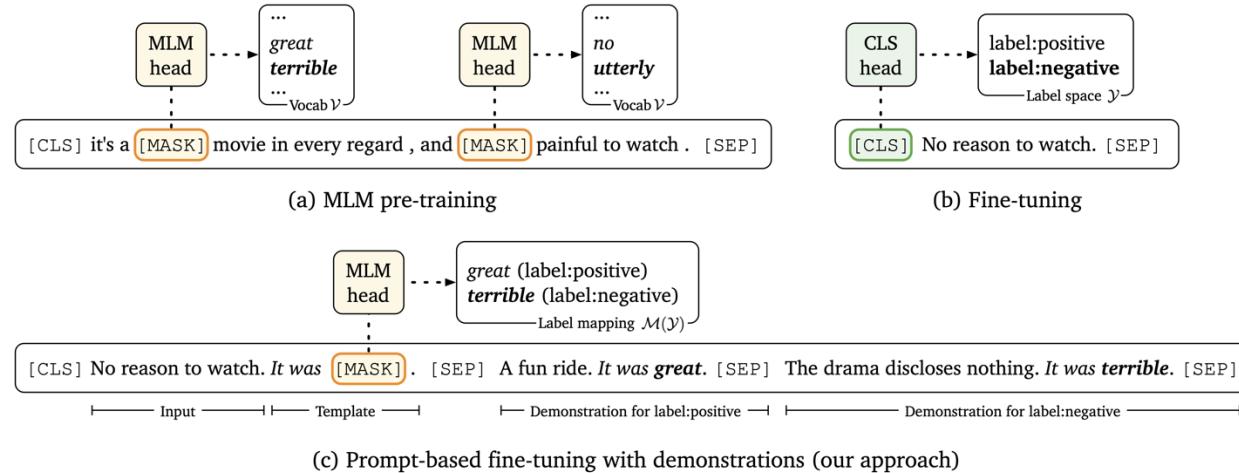
- **Context**
  - Large-scale Pre-trained Language Models (PLMs) have led to major breakthroughs in NLP.
  - Labeled data can be expensive or scarce in many domains. How do we **achieve strong performance with limited labeled data?**
- **Few-shot & Zero-shot Learning**
  - **Few-shot:** Use only a handful of labeled examples to fine-tune or adapt a model.
  - **Zero-shot:** The model is applied to a new task **without any** labeled examples.
- **Why Synthetic Data?**
  - Generating or augmenting training samples automatically.
  - Reduce annotation cost, improve model robustness, and expand coverage of data.

# Encoder-Only Era: Fine-tuning

- **Full Fine-tuning** (typical in BERT-era)
  - Standard practice: take the entire BERT model and fine-tune on a downstream task.
  - **Challenge:** Requires a sufficient amount of labeled data to avoid overfitting.
- **Why Few-shot/Zero-shot was Hard**
  - Encoder-only models (e.g., BERT) lack autoregressive generation, making open-ended task adaptation difficult.
  - Their masked language modeling (MLM) objective limits dynamic response generation.
- **Initial Solutions**
  - Some attempts involved leveraging **semantic representations** (CLS embeddings, etc.) or minimal modifications to the model.
  - Often still needed a decent amount of labeled data to fine-tune effectively.

# LM-BFF: An Example of Few-shot Progress (Encoder-Only)

- **Key Ideas:**
  - Convert downstream tasks to a “masked language modeling” style via **prompt templates**.



- **Takeaway:** Showed that even encoder-based models like BERT can be adapted for few-shot learning with prompt techniques, bridging the gap to more generative-style methods.

# Transition to Decoder Models: T5, GPT-3 & Emerging Zero-shot Capabilities

- **Filling the Gaps**
  - Decoder (generative) models, such as **GPT** or **T5**, are naturally suited for text generation.
  - This flexibility allows tasks to be framed as “text-to-text” problems.
- **Key Milestones**
  - **T5 [1]** (2019): Proposed a unified “Text-to-Text” framework for various NLP tasks.
  - **GPT-3 [2]** (2020): Demonstrated **In-Context Learning**, using only prompts and a few examples without weight updates.
- **Few-shot & Zero-shot**
  - T5 and GPT-3 showed impressive generalization, especially GPT-3’s ability to handle unseen tasks with just a prompt and few examples.

# T5

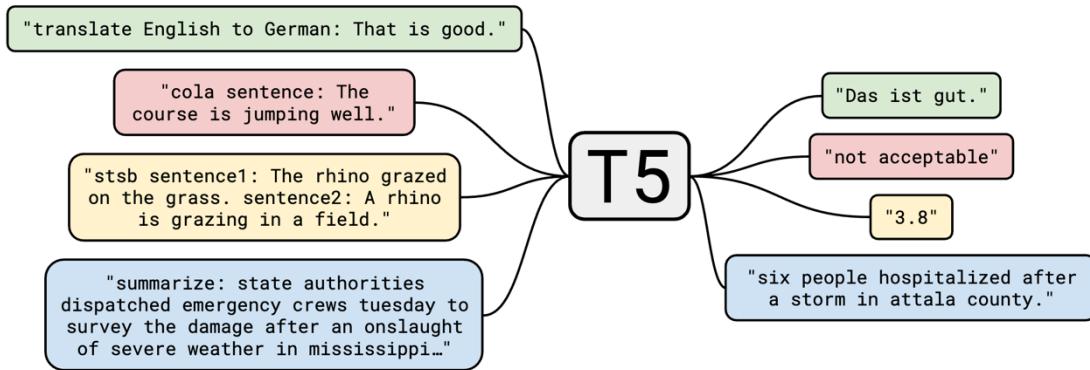


Figure 1: A diagram of our text-to-text framework. Every task we consider—including translation, question answering, and classification—is cast as feeding our model text as input and training it to generate some target text. This allows us to use the same model, loss function, hyperparameters, etc. across our diverse set of tasks. It also provides a standard testbed for the methods included in our empirical survey. “T5” refers to our model, which we dub the “Text-to-Text Transfer Transformer”.

- **Key Milestone:** Unified NLP tasks by framing them as text-to-text problems

- **Impact:** Set new performance standards on diverse NLP benchmarks by simplifying task formulation

# GPT3

The three settings we explore for in-context learning

## Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.



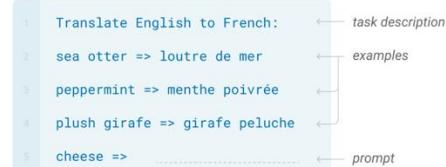
## One-shot

In addition to the task description, the model sees a single example of the task. No gradient updates are performed.



## Few-shot

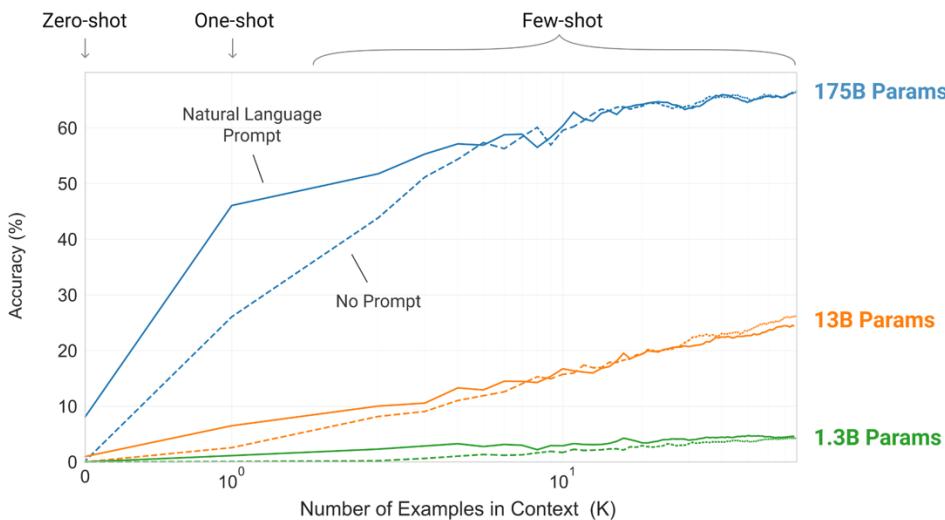
In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.



Traditional fine-tuning (not used for GPT-3)

## Fine-tuning

The model is trained via repeated gradient updates using a large corpus of example tasks.



• **Key Milestone:** Pushed the boundaries of scaling with 175 billion parameters

• **Impact:** Demonstrated strong zero-shot and few-shot learning abilities, reshaping expectations for language models

# Instruction tuning

- **Definition**
  - Train or fine-tune language models to explicitly **follow natural language instructions** for a variety of tasks.
- **Motivation**
  - Enhances **zero-shot** and **few-shot** performance: The model learns how to interpret instructions without per-task fine-tuning.
- **Key Points**
  - Often uses many tasks, each with multiple prompt templates.
  - Can be combined with **synthetic** or **human-labeled** instruction data.
  - Leads to more **generic** and **adaptive** models.

# T0 – Early Multitask Prompted Model

- **Method**
  - T0 is based on **T5**, but trained on many tasks with **multiple prompt templates** each.
- **Zero-Shot Generalization**
  - T0 can handle new tasks with just the prompt, no additional labels.
- **Significance**
  - Demonstrated the power of *prompt-based multi-task training* to achieve good zero-shot results.
  - Laid the groundwork for subsequent *instruction tuning* expansions.

# T0

- **Key Milestone:** Effectiveness of fine-tuning on a large variety of prompts for zero-shot task generalization
- **Impact:** Advanced instruction-following capabilities, proving effective across unseen tasks without specific fine-tuning

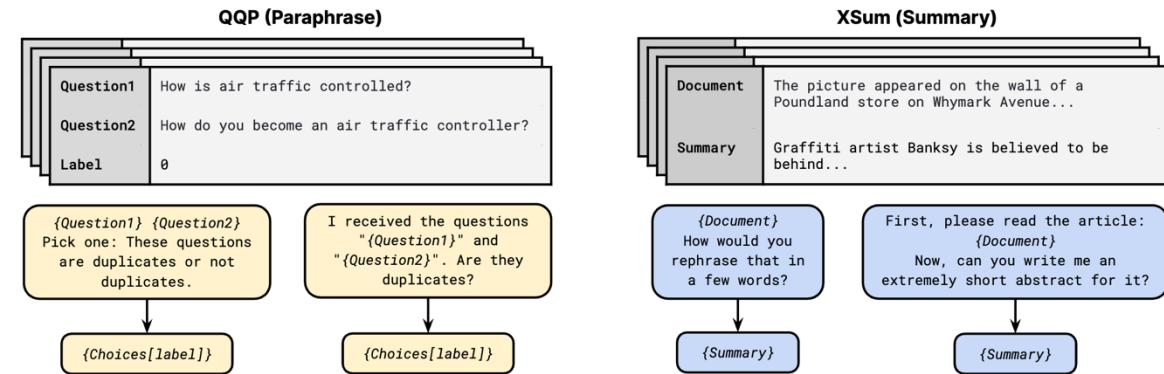


Figure 3: Prompt templates from the P3 prompt collection. Each dataset has multiple prompt templates consisting of an input and a target template. These use the fields of the raw data examples as well as template metadata, e.g., the left paraphrasing identification prompts use *Choices*, a template-level list variable ['Not duplicates', 'Duplicates']. These templates are materialized to produce the prompted instance shown in Figure 1. The complete set of prompt templates used in T0 is given in Appendix G.

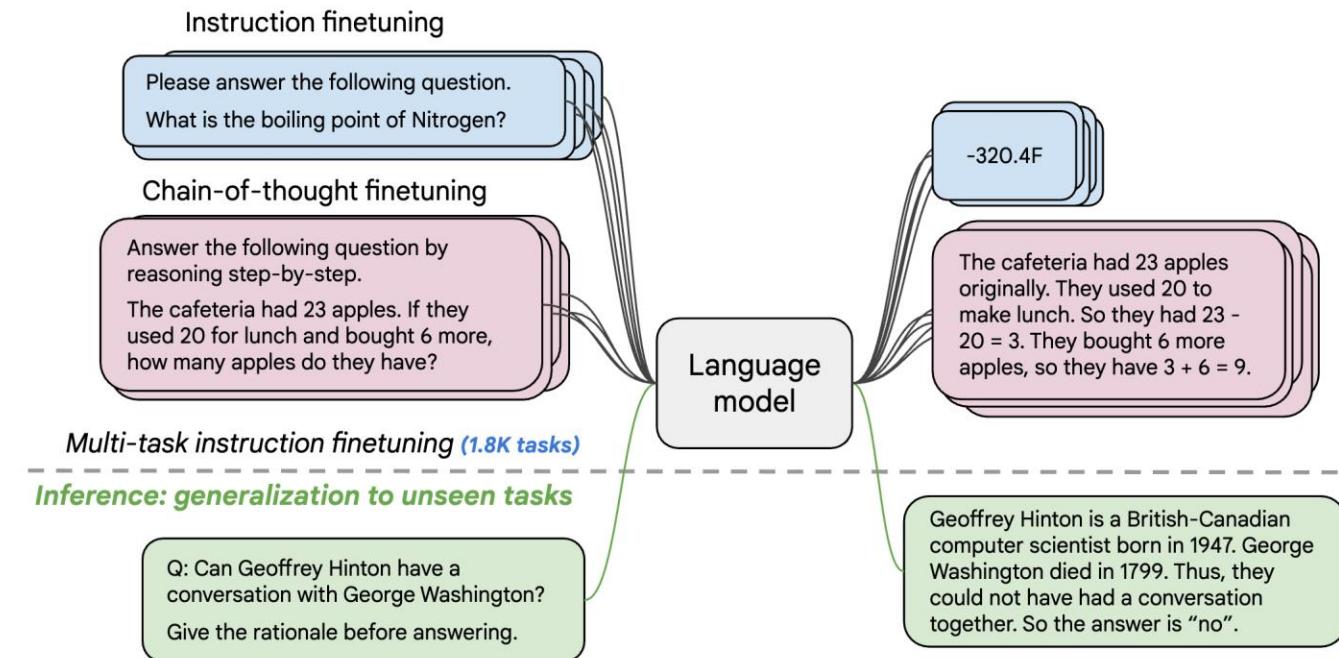
# FLAN – Expanding Instruction Tuning at Scale

- **Core Idea**
  - Further extends training with a **much wider range** of tasks & advanced prompting strategies (e.g., chain-of-thought in newer versions).
- **Contributions**
  - Larger & more diverse task coverage.
  - Refining prompt engineering to improve zero/few-shot performance.
- **Outcome**
  - FLAN models show strong cross-task generalization, confirming instruction tuning's efficacy.

# FLAN

## Key Milestones:

- Unified Prompting:** Introduced instruction tuning by unifying prompts across diverse tasks
- Further Scale Multi-Tasking:** Demonstrated that scaling multi-task capabilities enhances model performance by training on a wide array of tasks simultaneously.
- Chain-of-Thought (CoT):** Incorporated CoT prompting to improve the model's multi-step reasoning, leading to more coherent and effective problem-solving.



# InstructGPT – Instruction Tuning + Human Feedback

- **Method**
  - **Supervised Fine-tuning** on instruction–response pairs.
  - **RLHF** (Reinforcement Learning from Human Feedback): Collect preference data, train a reward model, and refine GPT to optimize for user alignment.
- **Why It Matters**
  - Prioritizes *helpfulness, harmlessness, human alignment*.
  - Directly led to systems like **ChatGPT**.
- **Data Efficiency Angle**
  - Key insight: A relatively smaller amount of **human feedback** data can significantly improve the *quality and usability* of large language models.

# InstructGPT

## Key Milestone:

- **Instruction Tuning:** Introduced a novel approach using reinforcement learning from human feedback (RLHF) to align model outputs with user instructions.

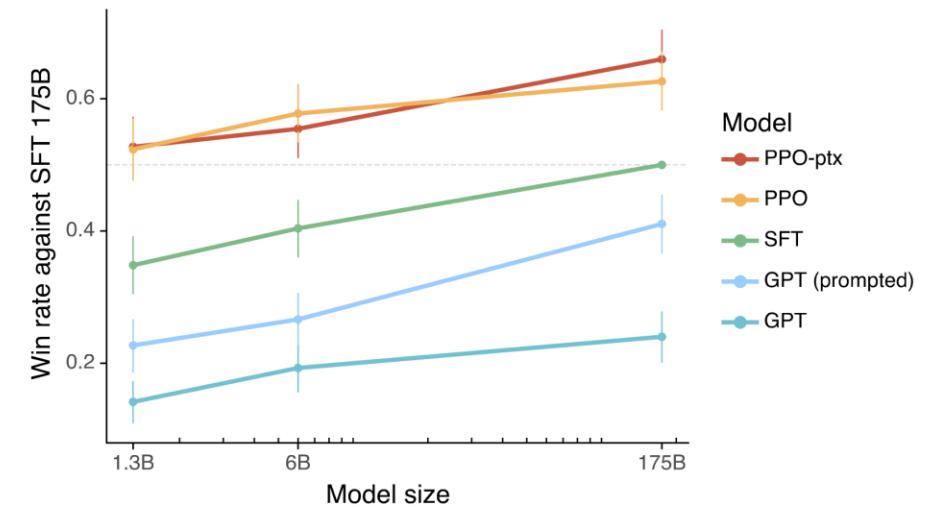
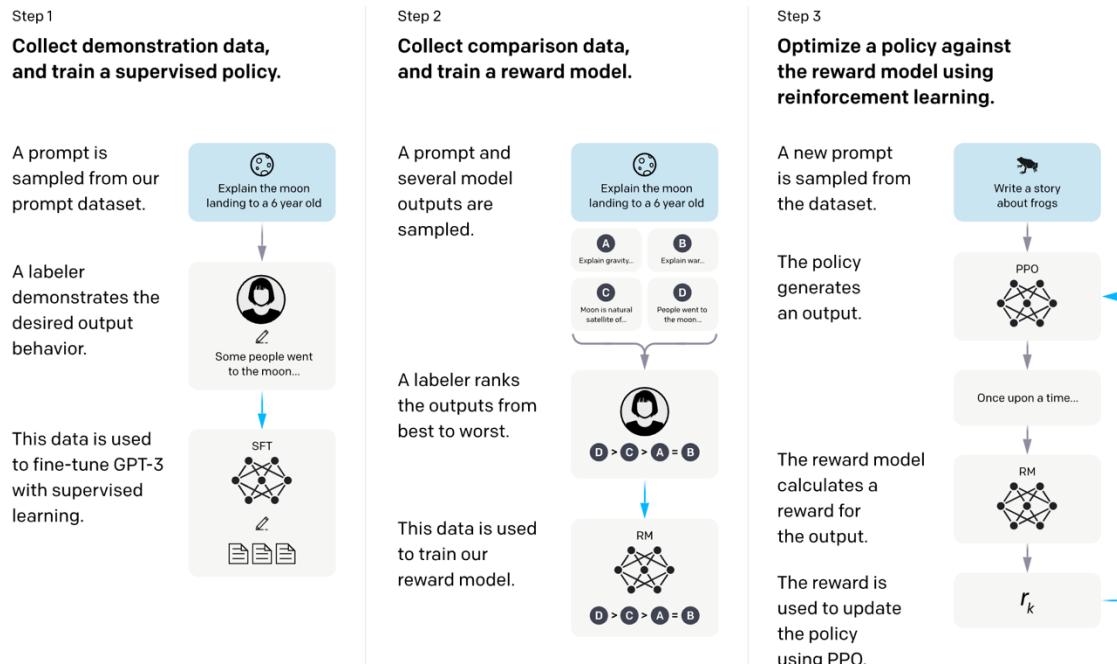


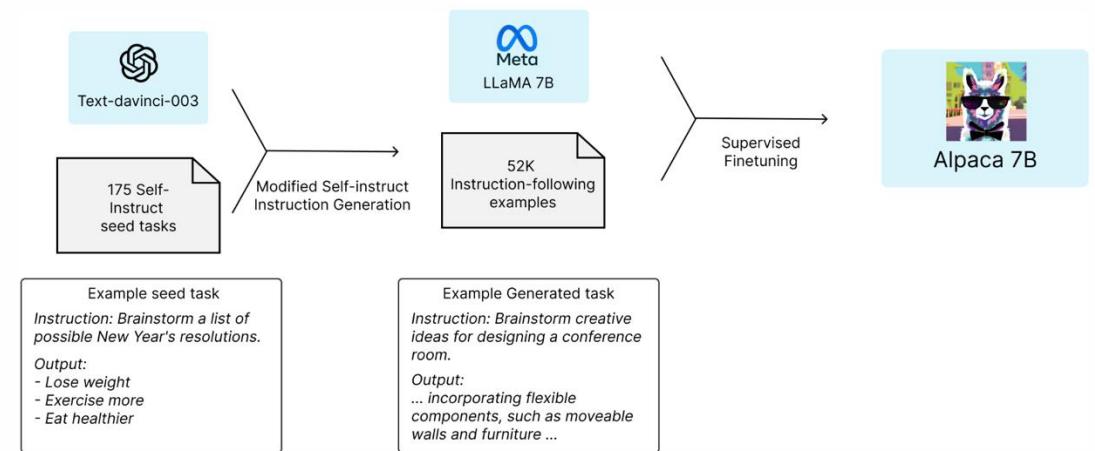
Figure 2: A diagram illustrating the three steps of our method: (1) supervised fine-tuning (SFT), (2) reward model (RM) training, and (3) reinforcement learning via proximal policy optimization (PPO) on this reward model. Blue arrows indicate that this data is used to train one of our models. In Step 2, boxes A-D are samples from our models that get ranked by labelers. See Section 3 for more details on our method.

# Synthetic Data

- Motivation
  - **Limited Data:** Real-world labeled data is often expensive or limited.
  - **Privacy Concerns:** Need to avoid using sensitive user data directly.
  - **Scalability & Diversity:** Synthetic data can provide large, varied training sets at lower cost.
- High-Level Ideas
  - **LLMs as Data Generators:** Using large models (GPT, etc.) to create question–answer pairs, reasoning steps, or domain-specific samples.
  - **Controlled or Semi-automatic Pipelines:** Combining LLM outputs with filtering, validation, or iterative refinement to ensure quality.

# Alpaca – Self-Instruct Applied

- **Process**
  - Use GPT-3.5 to generate ~52k instruction-response pairs automatically.
  - Fine-tune a **7B LLaMA** model on this synthetic set.
- **Results**
  - A surprisingly *ChatGPT-like* model with minimal direct human annotation.
  - Demonstrates the **low-cost** route to building an instruction-following model.
- **Significance**
  - **Proof of concept:** Synthetic data from strong models can rapidly boost smaller models' performance on instruction tasks.



# Orca—Chain of Thought Matters

- **Main Contribution**
  - Introduces **complex explanation traces** (chain-of-thought, step-by-step reasoning) from a *teacher* model (GPT-4 or other strong LLM).
  - A **student** LLM then fine-tunes on these CoT-augmented synthetic examples, achieving stronger reasoning performance.
- **Why It's Important**
  - **Richer Synthetic Signals:** Instead of just input–output pairs, each example includes *how* the teacher arrived at its conclusion.
  - **Boosted Reasoning Ability:** Students learn *reasoning patterns*, not just final answers—often outperforming models trained on plain Q–A data.
  - **Better Data Utilization:** A moderate amount of CoT-labeled data can sometimes surpass large amounts of non-explanatory data.

# LIMA-Less Is More for Alignment

- **Core Idea**
  - Fine-tune a 65B LLaMA model on just **1,000 carefully curated** instruction–response pairs.
  - Achieves performance *comparable* to models trained on tens or hundreds of thousands of synthetic instructions.
- **Key Takeaways**
  - **High-Quality Beats Large Quantity:** A small but *expertly chosen* dataset can often yield near state-of-the-art alignment.
  - **Focus on Diversity & Difficulty:** Curating diverse queries ensures coverage, while carefully validated responses maintain correctness.
  - **Implication:** Massive data hunts aren't always necessary; methodical selection and curation can provide *most* of the benefit.

# Synthetic Data – Expanding Training Data

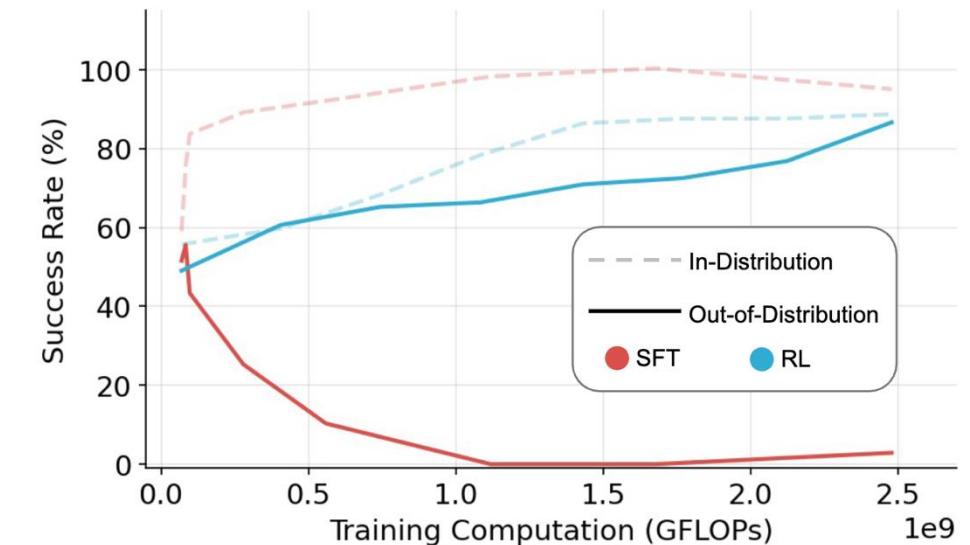
- Key Challenges
  - **Quality Control:** Avoid hallucinations, incorrect labels, or biases in generated data.
  - **Balancing Quantity vs. Quality:** Large-scale synthetic sets can still overfit or degrade performance if noisy.
  - **Evaluation Contamination:** Synthetic data may contain rephrasing of test examples, complicating fair benchmarking.
- This raises an important question:  
Does synthetic data help models generalize, or does it reinforce memorization?

# From Synthetic Data to Post-Training – SFT vs. RL

- **Synthetic data provides large-scale training resources**, but post-training methods determine how models use this data effectively.
- **Key Question:**
  - Does supervised fine-tuning (SFT) make models **overfit to training data**?
  - Can reinforcement learning (RL) **improve generalization** to unseen tasks?
- **Transitioning to Post-Training:**
  - **SFT (Supervised Fine-Tuning)**: Ensures instruction adherence but may lead to memorization.
  - **RL (Reinforcement Learning)**: Encourages flexibility but is harder to optimize.
- The recent study "**SFT Memorizes, RL Generalizes**" provides insights into these trade-offs.

# SFT Memorizes, RL Generalizes

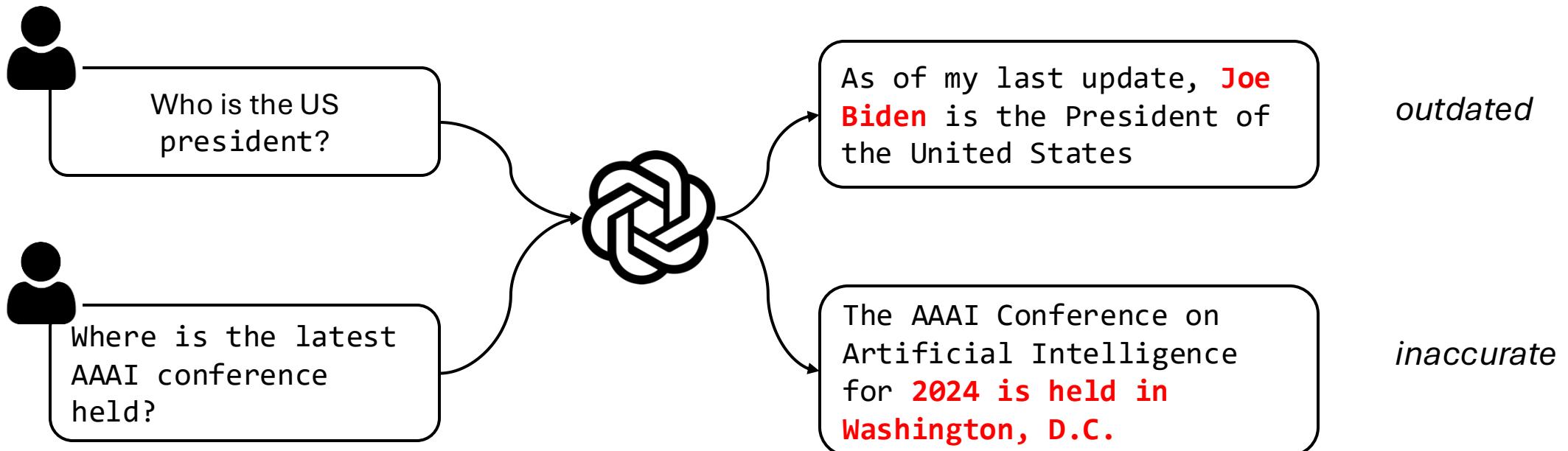
- **Core Question:**
  - What are the strengths and weaknesses of SFT vs. RL in foundation model post-training?
- **Key Findings:**
  - **SFT (Supervised Fine-Tuning)**
    - Memorizes seen data well.
    - Struggles with out-of-distribution (OOD) generalization.
    - Effective for tasks with **fixed structures**.
  - **RL (Reinforcement Learning)**
    - Adapts to unseen variations.
    - Improves **OOD generalization**.
    - Harder to train due to reward optimization complexity.
- **Conclusion:**
  - **SFT is good for stability, RL is better for adaptability.**
  - Combining **SFT + RL** leads to the best trade-off between memorization and generalization.



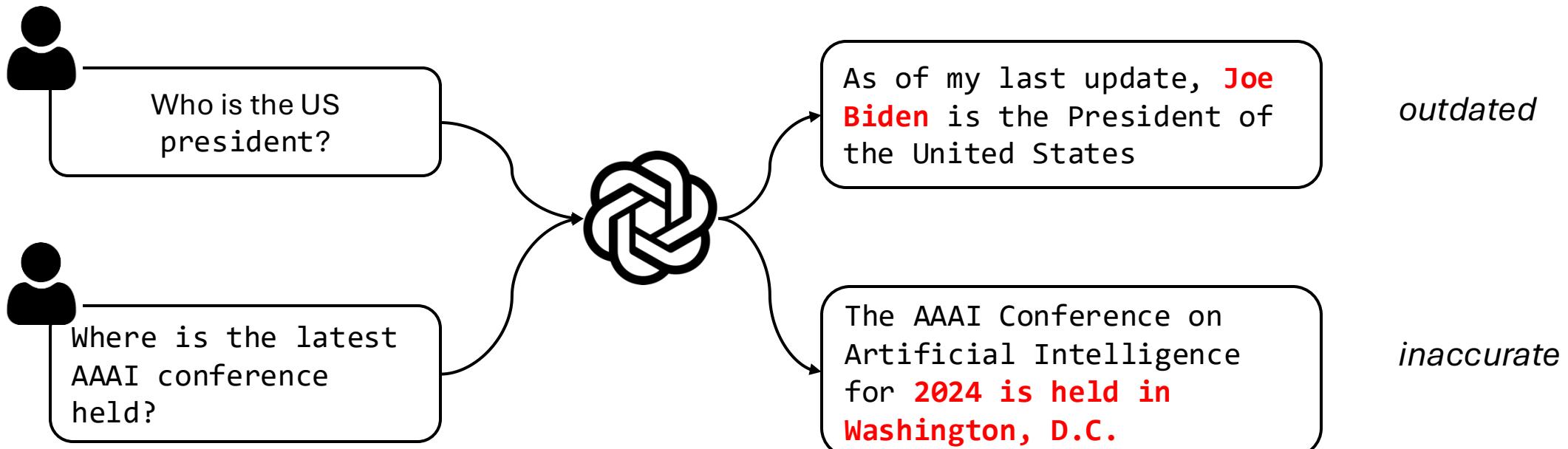
# Outline

- Model Efficiency
  - Parameter-efficient Fine-tuning
  - Model Compression
  - Inference Acceleration
- Data Efficiency
  - Zero/few-shot Learning
  - Knowledge Editing
  - Retrieval-augmented Generation

# Why Editing?

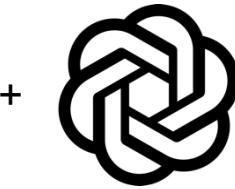


# Why Editing?



- LLMs learn *outdated* and *inaccurate* knowledge from their training data.

# Ways to update a LLM's knowledge



## Fine-tuning

- Expensive
- Hurt others



## Editing

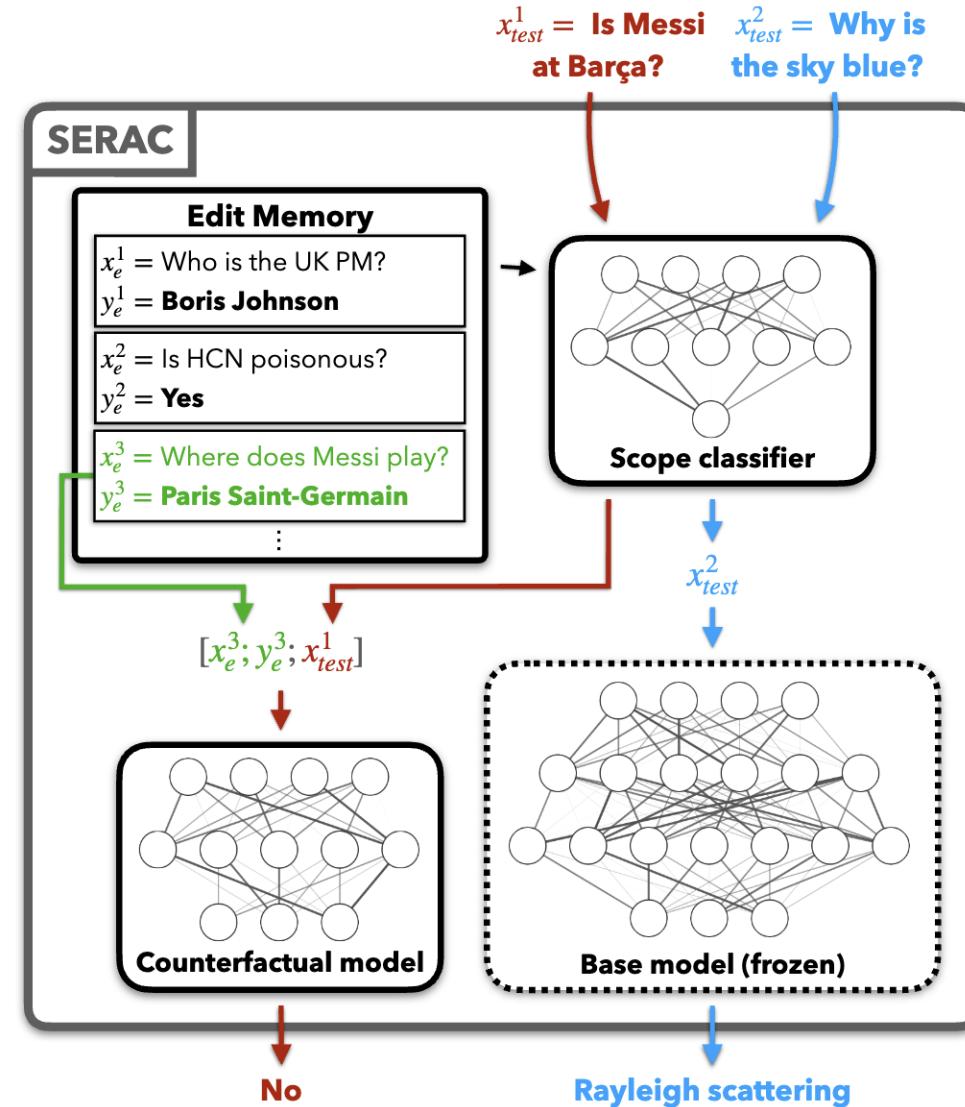
- Reliability: *Who is the US president?*
- Generality: *Who currently holds the office of President of the United States?*
- Portability: *Who is the first Lady of the US?*
- Locality: *Who is the prime minister of the UK?*

### Good Editing

# SERAC

Keep model parameters  
and learn new:

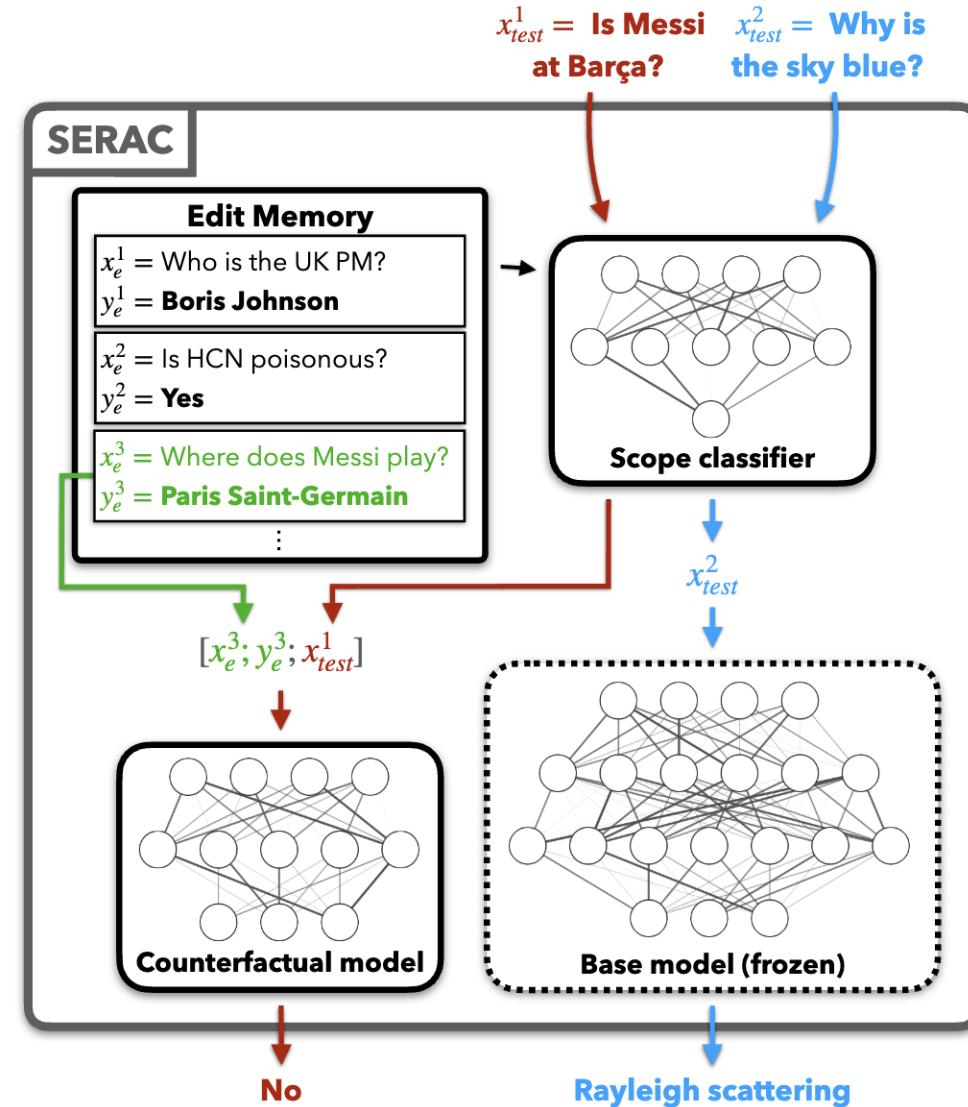
- **Counterfactual model:** handle *new knowledge*.
- **Scope classifier:** which *model* to use.



# SERAC

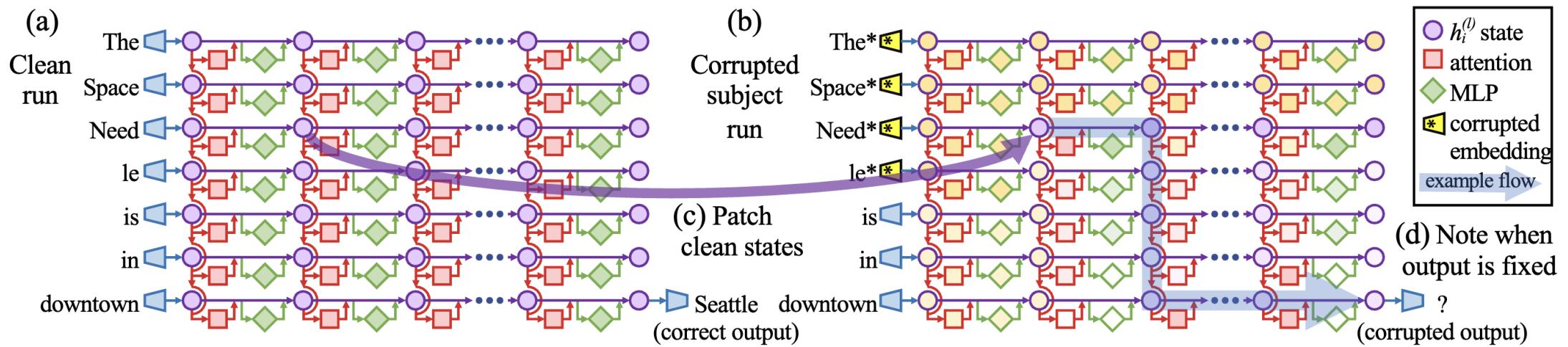
Keep model parameters  
and learn new:

- **Counterfactual model:** handle *new knowledge*.
  - **Scope classifier:** which *model* to use.
- Less efficient

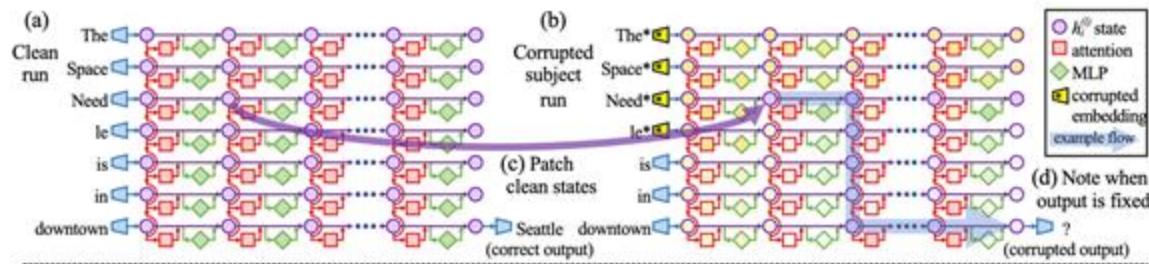


# ROME

- **Locate-and-Edit paradigm:**
- Locate: *find a few relevant parameters.*
- Edit: *update these parameters only.*

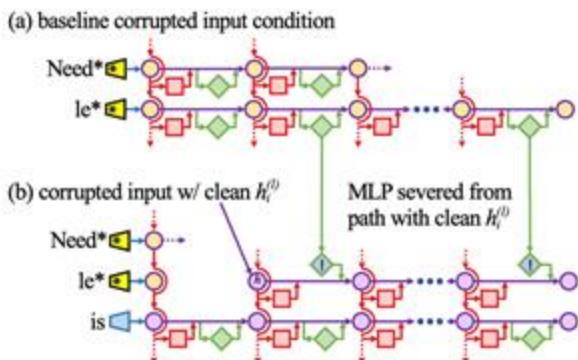


# ROME

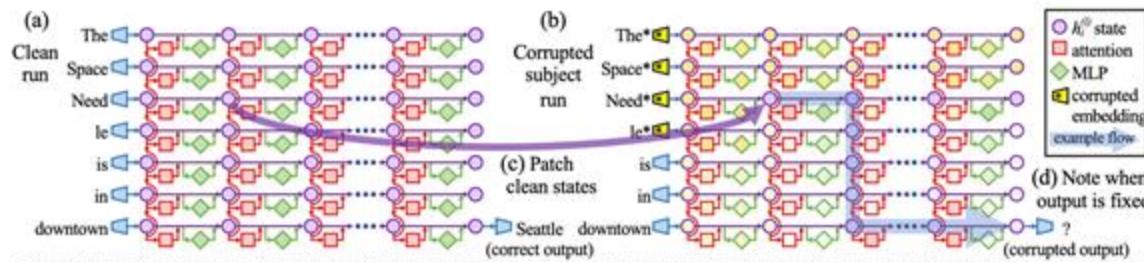


**Locate:** Find a layer that affects the final prediction most.

- Check how the layer can causally lead to the correct knowledge.

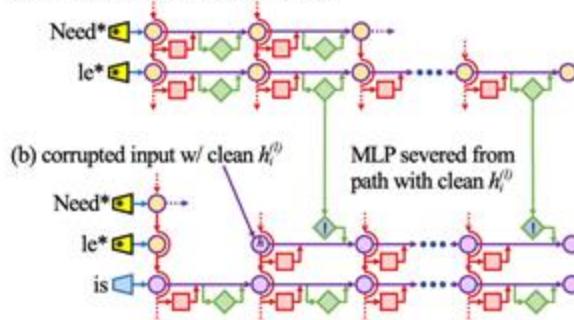


# ROME

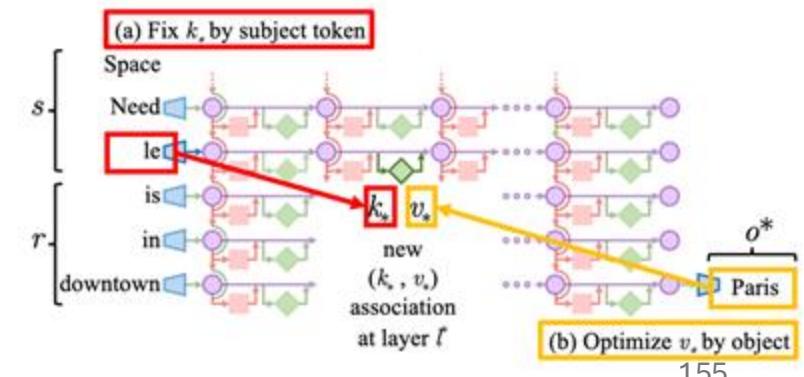


**Locate:** Find a layer that affects the final prediction most.

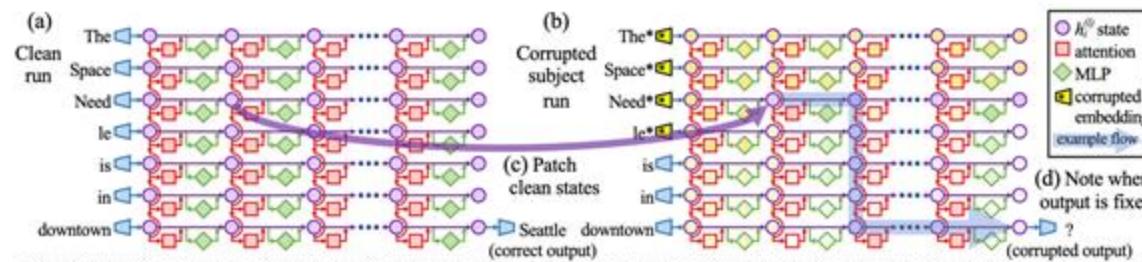
(a) baseline corrupted input condition



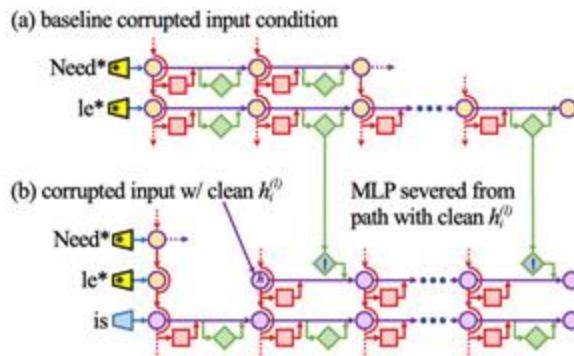
**Edit:**  
Add a new association to the Located Layer to output new knowledge.



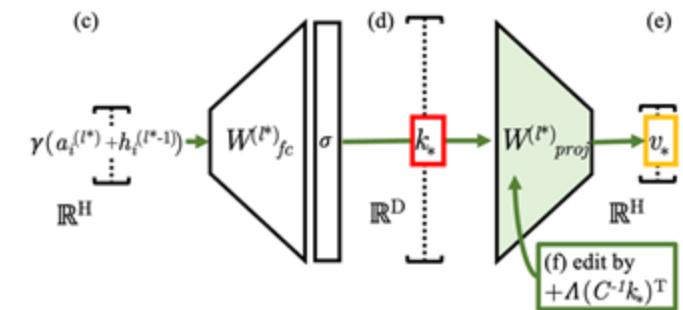
# ROME



**Locate:** Find a layer that affects the final prediction most.

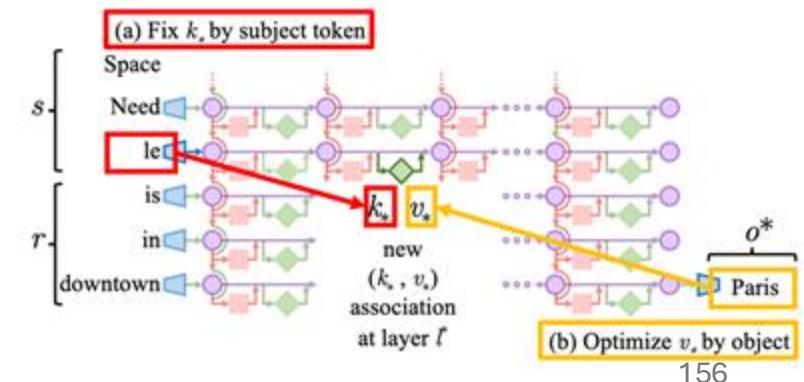


**Edit:**  
Add a new association to the Located Layer to output new knowledge.

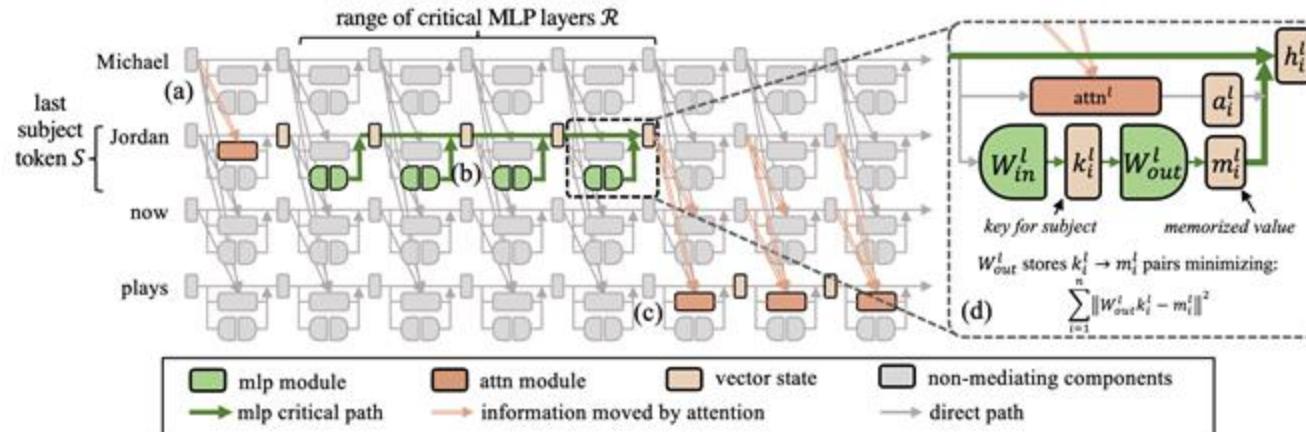


**Update: closed-form.**

maintain pretrained  
minimize  $\|\hat{W}K - V\|$  such that  $\hat{W}k_* = v_*$   
 $\hat{W} = W + \Lambda(C^{-1}k_*)^T$ .

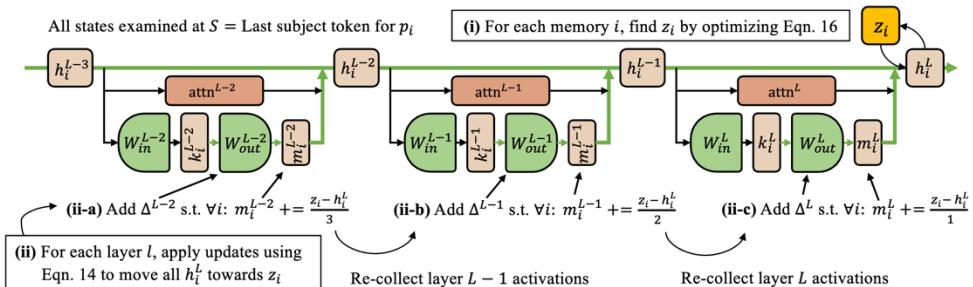


# MEMIT



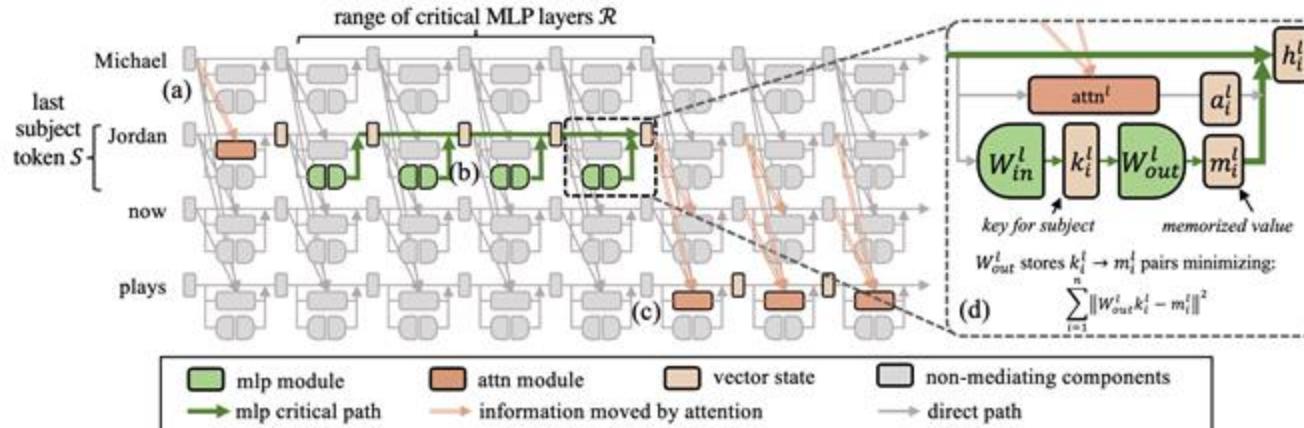
Locate: Find several layers

Multiple Knowledge

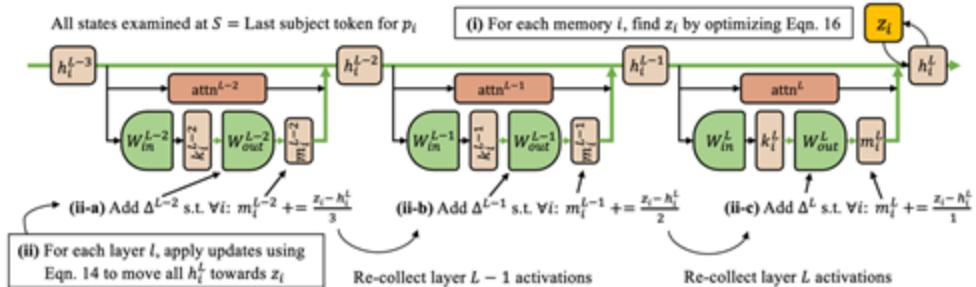


# MEMIT

## Multiple Knowledge



**Locate:** Find several layers



**Edit:**  
Each knowledge needs  
an association.

$$W_1 \triangleq \underset{\hat{W}}{\operatorname{argmin}} \left( \sum_{i=1}^n \|\hat{W}k_i - m_i\|^2 + \sum_{i=n+1}^{n+u} \|\hat{W}k_i - m_i\|^2 \right).$$

maintain pretrained      edit new

# RoseLoRA

## LoRA:

- Strong PEFT for general-purposed tasks.
- Make **dense** update, i.e., product of two dense low-rank matrices.
- Hard to **selectively** modify specific knowledge.

# RoseLoRA

LoRA:

- Strong PEFT for general-purposed tasks.
- Make **dense** update, i.e., product of two dense low-rank matrices.
- Hard to **selectively** modify specific knowledge.

↓  
Solution: *sparse update*.

LoRA update

$$\min_{A,B} \mathcal{L}(\mathcal{D}; \mathbf{W}^o + \boxed{\mathbf{B}\mathbf{A}})$$

$$\text{s.t. } \frac{\|\mathbf{B}\mathbf{A}\|_0}{d_1 d_2} \leq \tau,$$

Sparsity

# RoseLoRA

LoRA:

- Strong PEFT for general-purposed tasks.
- Make **dense** update, i.e., product of two dense low-rank matrices.
- Hard to **selectively** modify specific knowledge.

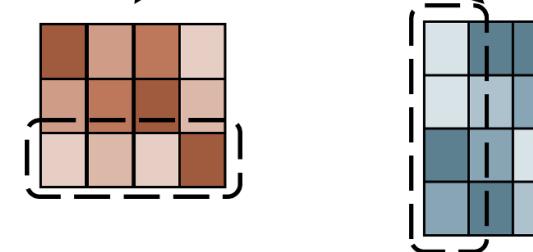
**Proposition.** The sparsity of  $BA$  is greater or equal to  $\max\{0, 1 + \sum_{i=1}^r (s(\mathbf{A}_{i*}) + s(\mathbf{B}_{*i}) - s(\mathbf{A}_{i*})s(\mathbf{B}_{*i})) - r\}$ .

Solution: *sparse update*.

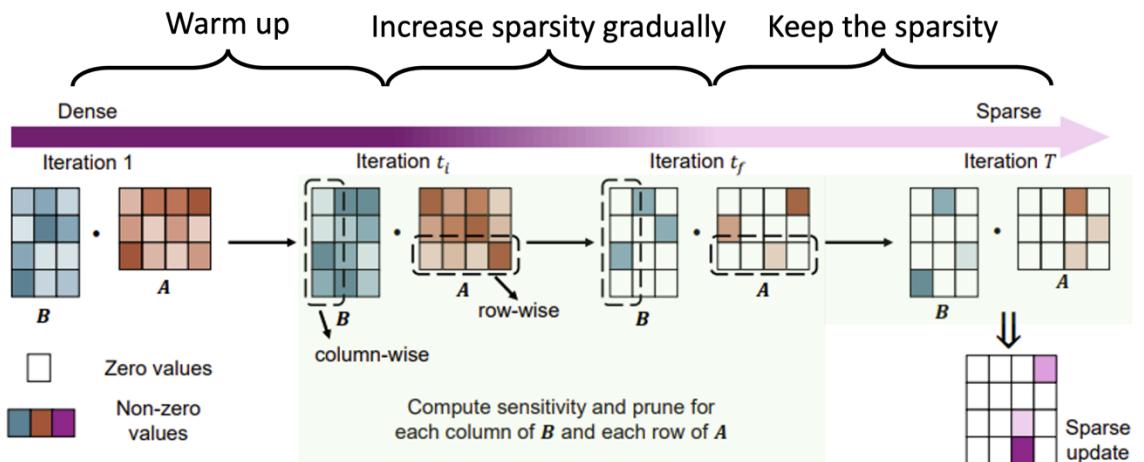
LoRA update

$$\begin{aligned} & \min_{A,B} \mathcal{L}(\mathcal{D}; \mathbf{W}^o + BA) \\ \text{s.t. } & \frac{\|BA\|_0}{d_1 d_2} \leq \tau, \\ & \text{Sparsity} \end{aligned}$$

Challenge:  
Product parameterization

$$\begin{aligned} & \min_{A,B} \mathcal{L}(\mathcal{D}; \mathbf{W}^o + BA) \\ \text{s.t. } & \frac{\|\mathbf{A}_{i*}\|_0}{d_2} \leq \tau, \frac{\|\mathbf{B}_{*i}\|_0}{d_1} \leq \tau, i = 1, \dots, r. \end{aligned}$$


# RoseLoRA



Dataset	Metric	FT-L <sup>♡</sup>	AdaLoRA <sup>♡</sup>	ROME <sup>♡</sup>	MEMIT <sup>♡</sup>	RoseLoRA
WikiData <sub>recent</sub>	Edit Succ.(↑)	71.2	65.6	85.1	85.3	<b>98.4</b>
	Locality(↑)	63.7	55.8	66.2	64.8	<b>83.4</b>
	Portability(↑)	48.7	47.2	37.5	37.9	<b>54.3</b>
	Fluency(↑)	549	538	574	567	<b>585</b>
	AVG(↑)	59.6	55.6	61.5	61.2	<b>73.7</b>
WikiData <sub>counterfact</sub>	Edit Succ.(↑)	51.1	72.1	83.2	83.4	<b>99.4</b>
	Locality(↑)	62.5	66.8	65.4	63.7	<b>90.9</b>
	Portability(↑)	39.1	55.2	38.7	40.1	<b>57.2</b>
	Fluency(↑)	545	554	579	569	<b>592</b>
	AVG(↑)	51.8	62.4	61.3	61.0	<b>76.7</b>
ZsRE	Edit Succ.(↑)	51.1	72.1	83.2	83.4	<b>100</b>
	Locality(↑)	62.5	66.8	65.4	63.7	<b>92.5</b>
	Portability(↑)	39.1	<b>55.2</b>	38.7	40.1	50.9
	Fluency(↑)	545	554	<b>579</b>	569	574
	AVG(↑)	54.6	62.1	58.2	54.0	<b>75.2</b>
WikiBio	Edit Succ.(↑)	66.3	97.0	95.1	94.3	<b>99.5</b>
	Locality(↑)	60.1	57.9	47.0	51.6	<b>92.5</b>
	Fluency(↑)	604	616	617	617	<b>620</b>
	AVG(↑)	62.3	72.2	67.9	69.2	<b>84.6</b>

# GRACE

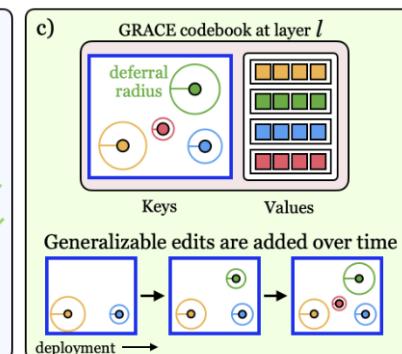
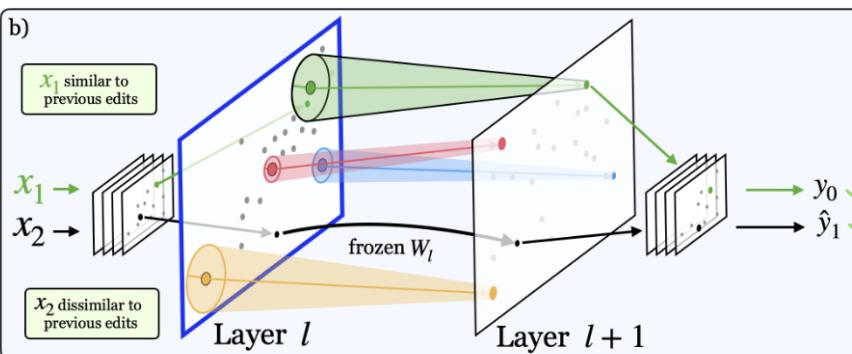
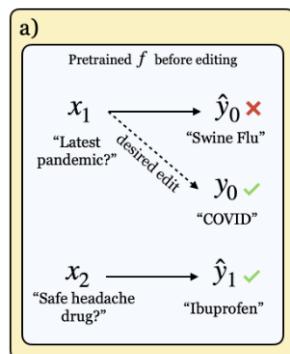
Previous methods are good at editing pre-specified number of knowledge.

But how to edit new knowledge that emerge continually?

# GRACE

But how to edit new knowledge that emerge continually?

Use dynamic parameters.



# GRACE

**Algorithm 1:** Update Codebook at layer  $l$ .

**Input:**  $\mathcal{C} = \{(\mathbb{K}_i, \mathbb{V}_i, \epsilon_i)\}_{i=0}^{C-1}$ , codebook

**Input:**  $f(\cdot)$ , model

**Input:**  $y_t$ , desired label

**Input:**  $x_t$ , edit input for which  $f(x_t) \neq y_t$

**Input:**  $\epsilon_{\text{init}}$ , initial  $\epsilon$

**Input:**  $d(\cdot)$ , distance function

**Output:**  $\mathcal{C}$ , updated codebook

$$C = \|\mathcal{C}\|$$

$$\hat{y}, h^{l-1} = f^L(x_t), f^{l-1}(x_t)$$

$$d_{\min}, i = \min_i(d(h^{l-1}, \mathbb{K}_i))$$

If  $d_{\min} > \epsilon_i + \epsilon_{\text{init}}$  or  $C = 0$ :

#  $h^{l-1}$  far from existing entries or empty  $\mathcal{C}$

$v_{\text{new}} = \text{finetune on } P_f(y|v_{\text{init}})$

$\mathcal{C}_C = (h^{l-1}, v_{\text{new}}, \epsilon_{\text{init}})$  # Add entry

Else:

#  $h^{l-1}$  near existing entries

If  $f^L(k_i) = y$ :

# Same label → Expand

$\mathcal{C}_i := (k_i, v_i, \epsilon_i + \epsilon_{\text{init}})$

Else:

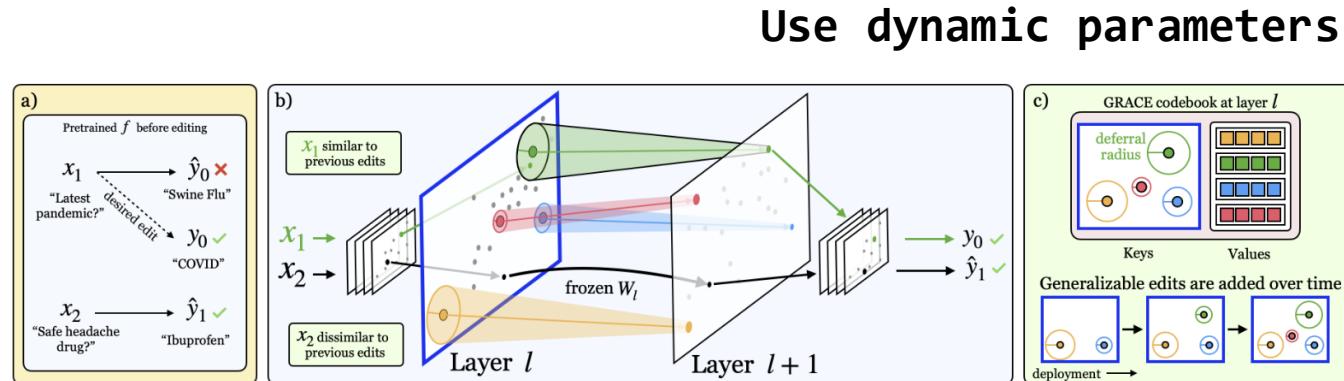
# Different label → Split

$\mathcal{C}_i = (k_i, v_i, d_{\min}/2)$  # Update entry  $i$

$v_{\text{new}} = \text{finetune on } P_f(y|v_{\text{init}})$

$\mathcal{C}_C = (h^{l-1}, v_{\text{new}}, d_{\min}/2)$  # Add entry

**return:**  $\mathcal{C}$



Use dynamic parameters.

# GRACE

**Algorithm 1:** Update Codebook at layer  $l$ .

**Input:**  $\mathcal{C} = \{(\mathbb{K}_i, \mathbb{V}_i, \epsilon_i)\}_{i=0}^{C-1}$ , codebook

**Input:**  $f(\cdot)$ , model

**Input:**  $y_t$ , desired label

**Input:**  $x_t$ , edit input for which  $f(x_t) \neq y_t$

**Input:**  $\epsilon_{\text{init}}$ , initial  $\epsilon$

**Input:**  $d(\cdot)$ , distance function

**Output:**  $\mathcal{C}$ , updated codebook

$$C = \|\mathcal{C}\|$$

$$\hat{y}, h^{l-1} = f^L(x_t), f^{l-1}(x_t)$$

$$d_{\min}, i = \min_i(d(h^{l-1}, \mathbb{K}_i))$$

If  $d_{\min} > \epsilon_i + \epsilon_{\text{init}}$  or  $C = 0$ :

#  $h^{l-1}$  far from existing entries or empty  $\mathcal{C}$   
 $v_{\text{new}} = \text{finetune on } P_f(y|v_{\text{init}})$   
 $\mathcal{C}_C = (h^{l-1}, v_{\text{new}}, \epsilon_{\text{init}})$  # Add entry

Else:

#  $h^{l-1}$  near existing entries

If  $f^L(k_i) = y$ :

# Same label → Expand

$$C_i := (k_i, v_i, \epsilon_i + \epsilon_{\text{init}})$$

Else:

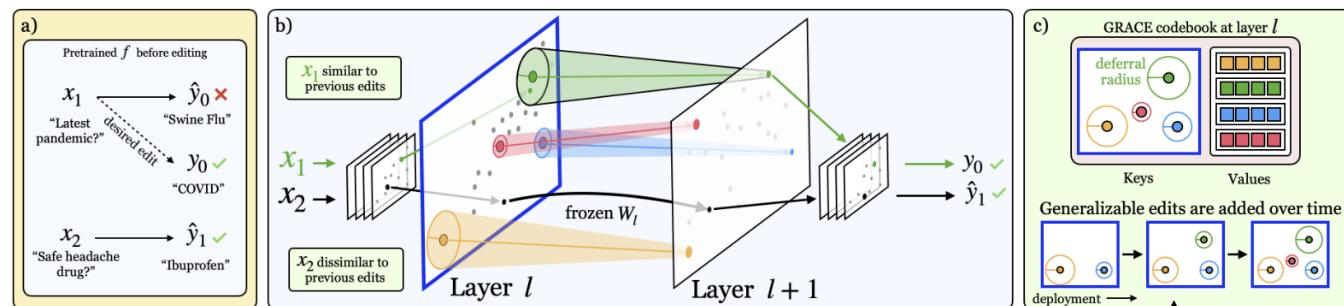
# Different label → Split

$$C_i = (k_i, v_i, d_{\min}/2)$$
 # Update entry  $i$

$$v_{\text{new}} = \text{finetune on } P_f(y|v_{\text{init}})$$

$$\mathcal{C}_C = (h^{l-1}, v_{\text{new}}, d_{\min}/2)$$
 # Add entry

return:  $\mathcal{C}$



**Continual Editing:** gradually grow number of parameters to Learn new knowledge

# GRACE

**Algorithm 1:** Update Codebook at layer  $l$ .

**Input:**  $\mathcal{C} = \{(\mathbb{K}_i, \mathbb{V}_i, \epsilon_i)\}_{i=0}^{C-1}$ , codebook

**Input:**  $f(\cdot)$ , model

**Input:**  $y_t$ , desired label

**Input:**  $x_t$ , edit input for which  $f(x_t) \neq y_t$

**Input:**  $\epsilon_{\text{init}}$ , initial  $\epsilon$

**Input:**  $d(\cdot)$ , distance function

**Output:**  $\mathcal{C}$ , updated codebook

$$C = \|\mathcal{C}\|$$

$$\hat{y}, h^{l-1} = f^L(x_t), f^{l-1}(x_t)$$

$$d_{\min}, i = \min_i(d(h^{l-1}, \mathbb{K}_i))$$

If  $d_{\min} > \epsilon_i + \epsilon_{\text{init}}$  or  $C = 0$ :

#  $h^{l-1}$  far from existing entries or empty  $\mathcal{C}$   
 $v_{\text{new}} = \text{finetune on } P_f(y|v_{\text{init}})$   
 $\mathcal{C}_C = (h^{l-1}, v_{\text{new}}, \epsilon_{\text{init}})$  # Add entry

Else:

#  $h^{l-1}$  near existing entries

If  $f^L(k_i) = y$ :

# Same label → Expand  
 $C_i := (k_i, v_i, \epsilon_i + \epsilon_{\text{init}})$

Else:

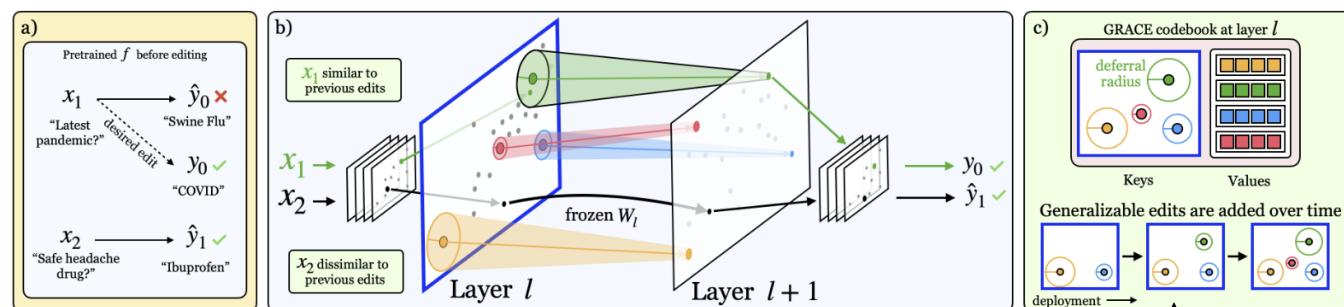
# Different label → Split

$C_i = (k_i, v_i, d_{\min}/2)$  # Update entry  $i$   
 $v_{\text{new}} = \text{finetune on } P_f(y|v_{\text{init}})$   
 $\mathcal{C}_C = (h^{l-1}, v_{\text{new}}, d_{\min}/2)$  # Add entry

return:  $\mathcal{C}$

$$h^l = \begin{cases} \text{GRACE}(h^{l-1}) & \text{if } \min_i(d(h^{l-1}, \mathbb{K}_i)) < \epsilon_{i_*}, \text{ where } i_* = \operatorname{argmin}_i(d(h^{l-1}), \mathbb{K}_i), \\ f^l(h^{l-1}) & \text{otherwise,} \end{cases}$$

**Inference: avoid incorrect edits**



**Continual Editing:** gradually grow number of parameters to Learn new knowledge

# GRACE

---

**Algorithm 1:** Update Codebook at layer  $l$ .

---

**Input:**  $\mathcal{C} = \{(\mathbb{K}_i, \mathbb{V}_i, \epsilon_i)\}_{i=0}^{C-1}$ , codebook

**Input:**  $f(\cdot)$ , model

**Input:**  $y_t$ , desired label

**Input:**  $x_t$ , edit input for which  $f(x_t) \neq y_t$

**Input:**  $\epsilon_{\text{init}}$ , initial  $\epsilon$

**Input:**  $d(\cdot)$ , distance function

**Output:**  $\mathcal{C}$ , updated codebook

$$C = \|\mathcal{C}\|$$

$$\hat{y}, h^{l-1} = f^L(x_t), f^{l-1}(x_t)$$

$$d_{\min}, i = \min_i(d(h^{l-1}, \mathbb{K}_i))$$

If  $d_{\min} > \epsilon_i + \epsilon_{\text{init}}$  or  $C = 0$ :

#  $h^{l-1}$  far from existing entries or empty  $\mathcal{C}$   
 $v_{\text{new}} = \text{finetune on } P_f(y|v_{\text{init}})$   
 $\mathcal{C}_C = (h^{l-1}, v_{\text{new}}, \epsilon_{\text{init}})$  # Add entry

Else:

#  $h^{l-1}$  near existing entries

If  $f^L(k_i) = y$ :

# Same label → Expand  
 $C_i := (k_i, v_i, \epsilon_i + \epsilon_{\text{init}})$

Else:

# Different label → Split

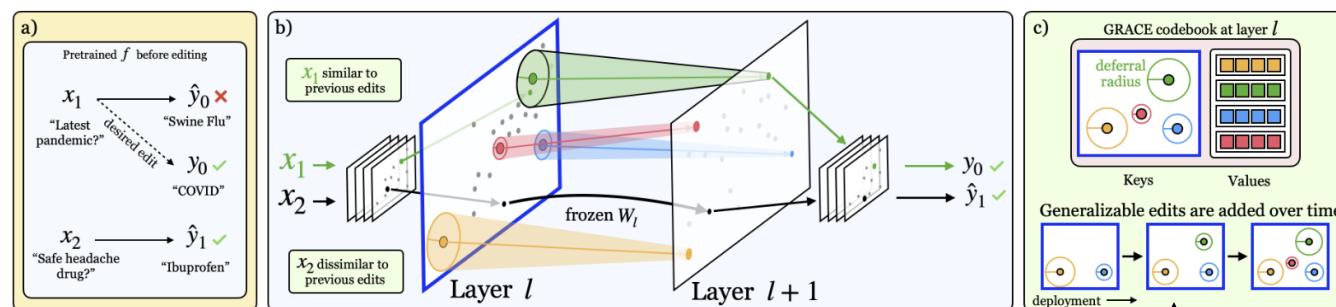
$C_i = (k_i, v_i, d_{\min}/2)$  # Update entry  $i$   
 $v_{\text{new}} = \text{finetune on } P_f(y|v_{\text{init}})$   
 $\mathcal{C}_C = (h^{l-1}, v_{\text{new}}, d_{\min}/2)$  # Add entry

return:  $\mathcal{C}$

---

$$h^l = \begin{cases} \text{GRACE}(h^{l-1}) & \text{if } \min_i(d(h^{l-1}, \mathbb{K}_i)) < \epsilon_{i_*}, \text{ where } i_* = \operatorname{argmin}_i(d(h^{l-1}), \mathbb{K}_i), \\ f^l(h^{l-1}) & \text{otherwise,} \end{cases}$$

Inference: avoid incorrect edits



Continual Editing: gradually grow number of parameters to Learn new knowledge

Different codebooks:

- LoRA: MELO [AAAI 2024]
- FFN: WISE [NeurIPS 2024]
- ...

# BaFT

**Parameter-based edits** alter all inputs in the same way.

How to get more fine-grained updates?

**ReFT [NeurIPS 2024]**: alter a few representations in some low-rank subspace

$$\Phi_l(\mathbf{h}_i^{(l)}; \phi_l) = \mathbf{h}_i^{(l)} + \mathbf{R}_l^\top (\mathbf{A}_l \mathbf{h}_i^{(l)} + \mathbf{b}_l - \mathbf{R}_l \mathbf{h}_i^{(l)}),$$

# BaFT

**ReFT:** alter pre-specified representations in a low-rank subspace

$$\Phi_l(\mathbf{h}_i^{(l)}; \phi_l) = \mathbf{h}_i^{(l)} + \mathbf{R}_l^\top (\mathbf{A}_l \mathbf{h}_i^{(l)} + \mathbf{b}_l - \mathbf{R}_l \mathbf{h}_i^{(l)}), \text{ linear}$$

***Impossible edit – generality - locality triangle***  
(paraphrase)

**Theorem 2.3.** When fine-tuning an LM, ReFT learns to update the old representation  $\mathbf{h}_0$  to targeted  $\mathbf{t} = \Phi(\mathbf{h}_0)$ . If ReFT maintains good generality such that  $\forall \mathbf{h} \in B(\mathbf{h}_0, \varepsilon(\mathbf{h}_0))$ ,

$$\|\Phi(\mathbf{h}) - \Phi(\mathbf{h}_0)\| = \|\Phi(\mathbf{h}) - \mathbf{t}\| < \varepsilon(\mathbf{t}),$$

where  $\|\cdot\|$  denote the  $\ell_2$  norm. Then for any irrelevant input  $\mathbf{h}_{ir}$  with a small stable-ball radius

$$\varepsilon(\mathbf{h}_{ir}) < \frac{\|\mathbf{t} - \mathbf{h}_0\| - (\varepsilon(\mathbf{t}) + \varepsilon(\mathbf{h}_0))}{\varepsilon(\mathbf{t}) + 2\varepsilon(\mathbf{h}_0)} \varepsilon(\mathbf{h}_0),$$

and is not too far from  $\mathbf{h}_0$  such that

$$\|\mathbf{h}_{ir} - \mathbf{h}_0\| = \varepsilon(\mathbf{h}_{ir}) + \varepsilon(\mathbf{h}_0),$$

ReFT will output  $\Phi(\mathbf{h}_{ir}) \notin B(\mathbf{h}_{ir}, \varepsilon(\mathbf{h}_{ir}))$  and break its locality guarantee.

# BaFT

**ReFT:** alter pre-specified representations in a low-rank subspace

$$\Phi_l(\mathbf{h}_i^{(l)}; \phi_l) = \mathbf{h}_i^{(l)} + \mathbf{R}_l^\top (\mathbf{A}_l \mathbf{h}_i^{(l)} + \mathbf{b}_l - \mathbf{R}_l \mathbf{h}_i^{(l)}), \text{ linear}$$

**Impossible edit – generality - locality triangle  
(paraphrase)**

**Theorem 2.3.** When fine-tuning an LM, ReFT learns to update the old representation  $\mathbf{h}_0$  to targeted  $\mathbf{t} = \Phi(\mathbf{h}_0)$ . If ReFT maintains good generality such that  $\forall \mathbf{h} \in B(\mathbf{h}_0, \varepsilon(\mathbf{h}_0))$ ,

$$\|\Phi(\mathbf{h}) - \Phi(\mathbf{h}_0)\| = \|\Phi(\mathbf{h}) - \mathbf{t}\| < \varepsilon(\mathbf{t}),$$

where  $\|\cdot\|$  denote the  $\ell_2$  norm. Then for any irrelevant input  $\mathbf{h}_{ir}$  with a small stable-ball radius

$$\varepsilon(\mathbf{h}_{ir}) < \frac{\|\mathbf{t} - \mathbf{h}_0\| - (\varepsilon(\mathbf{t}) + \varepsilon(\mathbf{h}_0))}{\varepsilon(\mathbf{t}) + 2\varepsilon(\mathbf{h}_0)} \varepsilon(\mathbf{h}_0),$$

and is not too far from  $\mathbf{h}_0$  such that

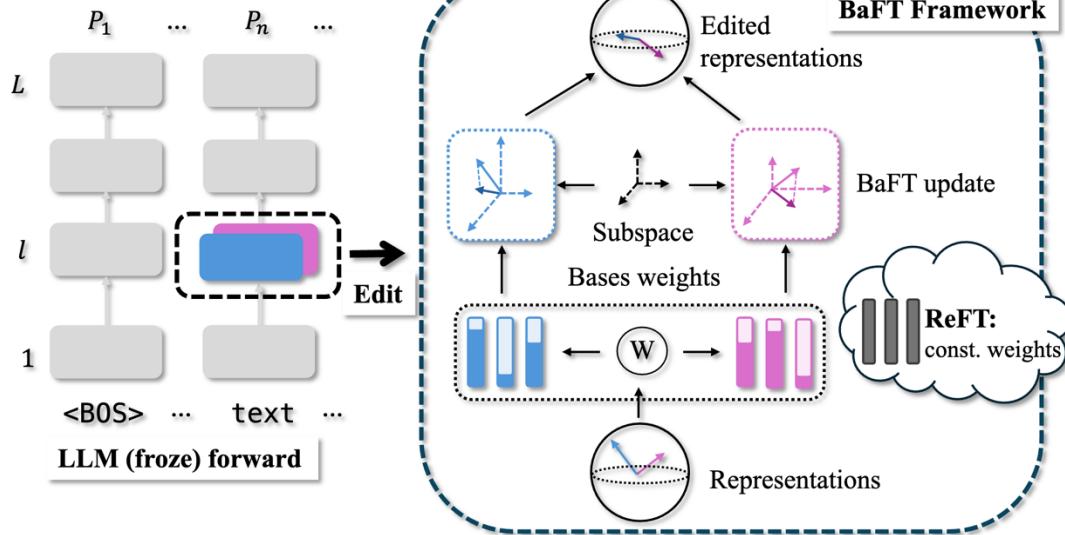
$$\|\mathbf{h}_{ir} - \mathbf{h}_0\| = \varepsilon(\mathbf{h}_{ir}) + \varepsilon(\mathbf{h}_0),$$

ReFT will output  $\Phi(\mathbf{h}_{ir}) \notin B(\mathbf{h}_{ir}, \varepsilon(\mathbf{h}_{ir}))$  and break its locality guarantee.

**BaFT:** use the subspace adaptively.

$$\Phi(\mathbf{h}) = \mathbf{h} + \mathbf{R}^\top \mathbf{W}(\mathbf{h}) (\mathbf{A}\mathbf{h} + \mathbf{b} - \mathbf{R}\mathbf{h}). \text{ nonlinear}$$

# BaFT



**ReFT:** alter pre-specified representations in a low-rank subspace

$$\Phi_l(\mathbf{h}_i^{(l)}; \phi_l) = \mathbf{h}_i^{(l)} + \mathbf{R}_l^\top (\mathbf{A}_l \mathbf{h}_i^{(l)} + \mathbf{b}_l - \mathbf{R}_l \mathbf{h}_i^{(l)}), \text{ linear}$$

**Impossible edit – generality - locality triangle (paraphrase)**

**Theorem 2.3.** When fine-tuning an LM, ReFT learns to update the old representation  $\mathbf{h}_0$  to targeted  $\mathbf{t} = \Phi(\mathbf{h}_0)$ . If ReFT maintains good generality such that  $\forall \mathbf{h} \in B(\mathbf{h}_0, \varepsilon(\mathbf{h}_0))$ ,  $\|\Phi(\mathbf{h}) - \Phi(\mathbf{h}_0)\| = \|\Phi(\mathbf{h}) - \mathbf{t}\| < \varepsilon(\mathbf{t})$ , where  $\|\cdot\|$  denote the  $\ell_2$  norm. Then for any irrelevant input  $\mathbf{h}_{ir}$  with a small stable-ball radius  $\varepsilon(\mathbf{h}_{ir}) < \frac{\|\mathbf{t} - \mathbf{h}_0\| - (\varepsilon(\mathbf{t}) + \varepsilon(\mathbf{h}_0))}{\varepsilon(\mathbf{t}) + 2\varepsilon(\mathbf{h}_0)} \varepsilon(\mathbf{h}_0)$ , and is not too far from  $\mathbf{h}_0$  such that  $\|\mathbf{h}_{ir} - \mathbf{h}_0\| = \varepsilon(\mathbf{h}_{ir}) + \varepsilon(\mathbf{h}_0)$ , ReFT will output  $\Phi(\mathbf{h}_{ir}) \notin B(\mathbf{h}_{ir}, \varepsilon(\mathbf{h}_{ir}))$  and break its locality guarantee.

**BaFT:** use the subspace adaptively.

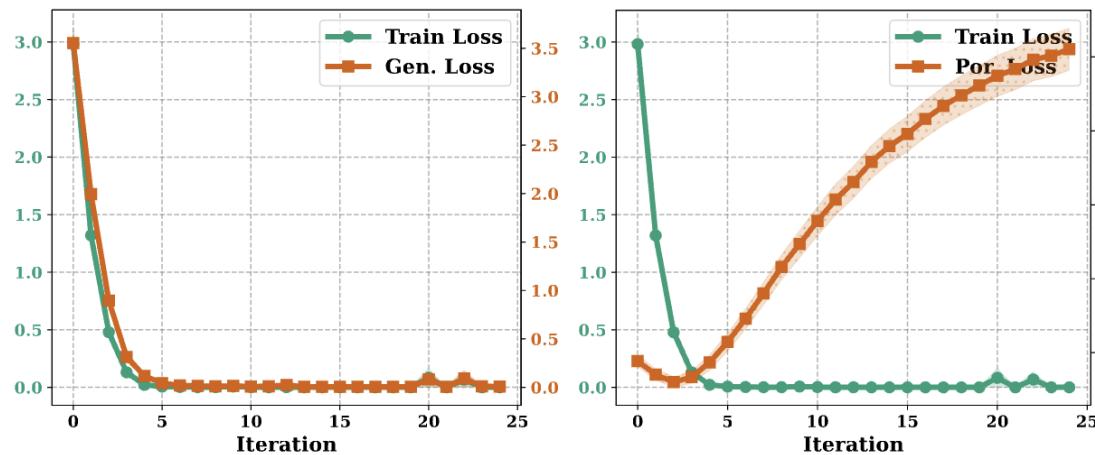
$$\Phi(\mathbf{h}) = \mathbf{h} + \mathbf{R}^\top \mathbf{W}(\mathbf{h}) (\mathbf{A}\mathbf{h} + \mathbf{b} - \mathbf{R}\mathbf{h}). \text{ nonlinear}$$

# BaFT

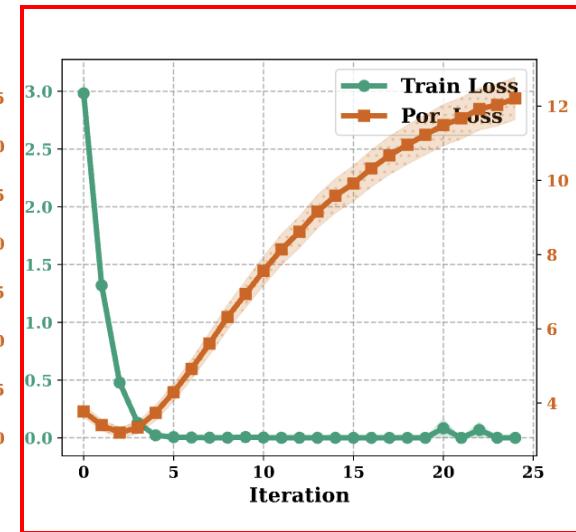
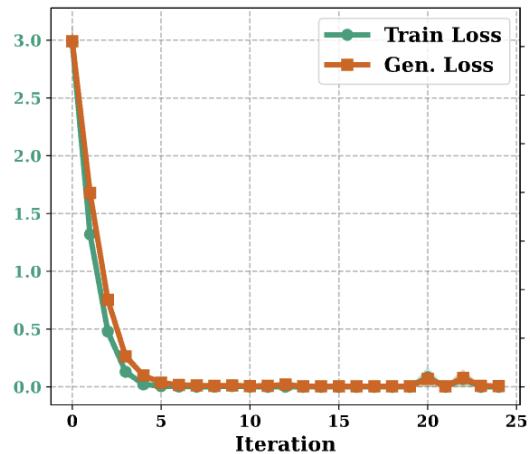
	LLaMA 2-7b(-chat)		LLaMA 3-8b-Instruct		Gemma 1.1-7b-Instruct	
	# Params.	Time (sec./edit)	# Params.	Time (sec./edit)	# Params.	Time (sec./edit)
AdaLoRA	6,292,224	26.24	5,112,576	28.71	4,817,568	44.24
FT-L	45,088,768	9.73	58,720,256	10.84	75,497,472	11.95
ROME	/	27.27	/	25.01	/	52.07
MEMIT	/	20.01	/	25.35	/	/
GRACE	/	34.38	/	87.08	/	43.45
WISE <sub>light</sub>	5,636,096	58.00	7,340,032	65.77	9,437,184	20.20
ReFT	393,264	10.99	393,264	9.33	294,960	7.79
BaFT (Ours)	606,256	13.46	606,256	12.69	454,704	10.13

# OVERTONE

Edited model tends to memorize rather than learn.

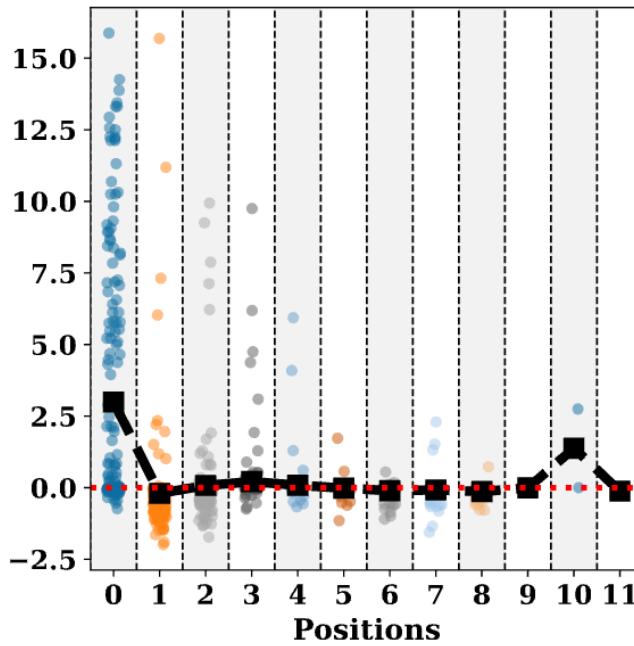
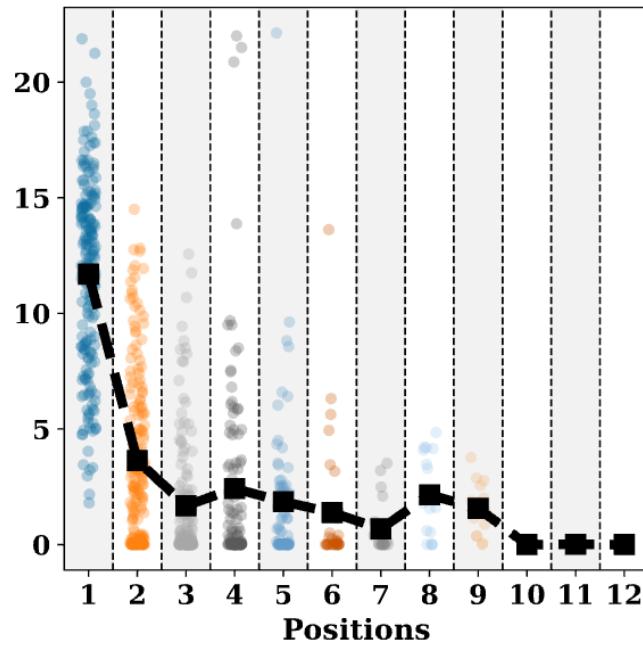


# OVERTONE [Arxiv 2025]



- Training: The US president is [X].
- Portability: The US first lady is?

# OVERTONE [Arxiv 2025]



$$\begin{aligned}
 \ell_{\text{CE}}(\theta) &\triangleq \sum_{i=1}^m \text{CE}[\delta_{y_i}(y) \| \pi_\theta(y \mid \mathbf{x} \oplus \mathbf{y}_{<i})] \quad (1) \\
 &= - \sum_{i=1}^m \log \pi_\theta(y_i \mid \mathbf{c}_i) \\
 &\downarrow \\
 &\text{Token wise-smoothing} \\
 \ell_{\text{OVERTONE}}(\theta) &\triangleq \sum_{i=1}^m \max(D_{\text{KL}}[\pi_{\text{tar}}(y \mid \mathbf{c}_i) \| \pi_\theta(y \mid \mathbf{c}_i)], \epsilon),
 \end{aligned}$$

# OVERTONE [Arxiv 2025]

keep “correct” token unchanged

$$\pi_{\text{tar}}^{(i)} \triangleq \begin{cases} \pi_{\text{tar}}^{\text{can}} \triangleq \lambda \delta_{y_i} + (1 - \lambda) \pi_{\text{ft}}^{(i)} & \text{if } y_i = \operatorname{argmax}_y \pi_{\text{tar}}^{\text{can}}, \\ \delta_{y_i} & \text{otherwise,} \end{cases} \quad (3)$$

---

**Algorithm 1** OVERTONE Training Paradigm

---

```

1: Input: Editing data ( $\mathbf{x}, \mathbf{y} = [y_1, \dots, y_m]$ ), LM parameters  $\theta_0$ , mixing hyper-parameter  $\lambda$ , early-stopping threshold  $\epsilon$ , filtering threshold  $n$ , total training steps  $T$ .
2: Initialize:  $\theta = \theta_0$ .
3: for  $t = 1, \dots, T$  do
4:   # Inner loop is parallelized in practice, unroll for better readability.
5:   for  $i = 1, \dots, m$  do
6:     Set context  $\mathbf{c}_i = \mathbf{x} \oplus \mathbf{y}_{<i}$ .
7:     Compute logits from the LM as  $\mathbf{s}^{(i)} = f_\theta(\mathbf{c}_i) \in \mathbb{R}^{|\mathcal{V}|}$ . Take softmax and get  $\pi_\theta^{(i)}$ .
8:     Top  $n\sigma$ -filter (Tang et al., 2024): Compute  $s_{\max}^{(i)} = \max_k s_k^{(i)}$ ,  $\sigma = \text{std}(\mathbf{s}^{(i)})$ . Define filtered logit  $\tilde{s}_k^{(i)} = -\infty$  if  $s_k^{(i)} \leq s_{\max}^{(i)} - n\sigma$  else  $\tilde{s}_k^{(i)} = s_k^{(i)}$ .
9:     Take softmax on filtered  $\tilde{\mathbf{s}}$  and get filtered  $\pi_{\text{ft}}^{(i)}$ .
10:    Compute target  $\pi_{\text{tar}}^{(i)}$  based on Eq (3).
11:    Compute loss
           $\ell_{\text{OVERTONE}}^{(i)} = \max(\text{D}_{\text{KL}}[\pi_{\text{tar}}^{(i)} \parallel \pi_\theta^{(i)}], \epsilon)$ .
12:  end for
13:  Compute sample loss
           $\ell_{\text{OVERTONE}}(\theta) = \sum_{i=1}^m \ell_{\text{OVERTONE}}^{(i)}$ .
14:  Update with learning rate  $\alpha$ 
           $\theta \leftarrow \theta - \alpha \nabla_\theta \ell_{\text{OVERTONE}}(\theta)$ 
15: end for
output Edited parameter  $\theta$ .

```

---

# OVERTONE [Arxiv 2025]

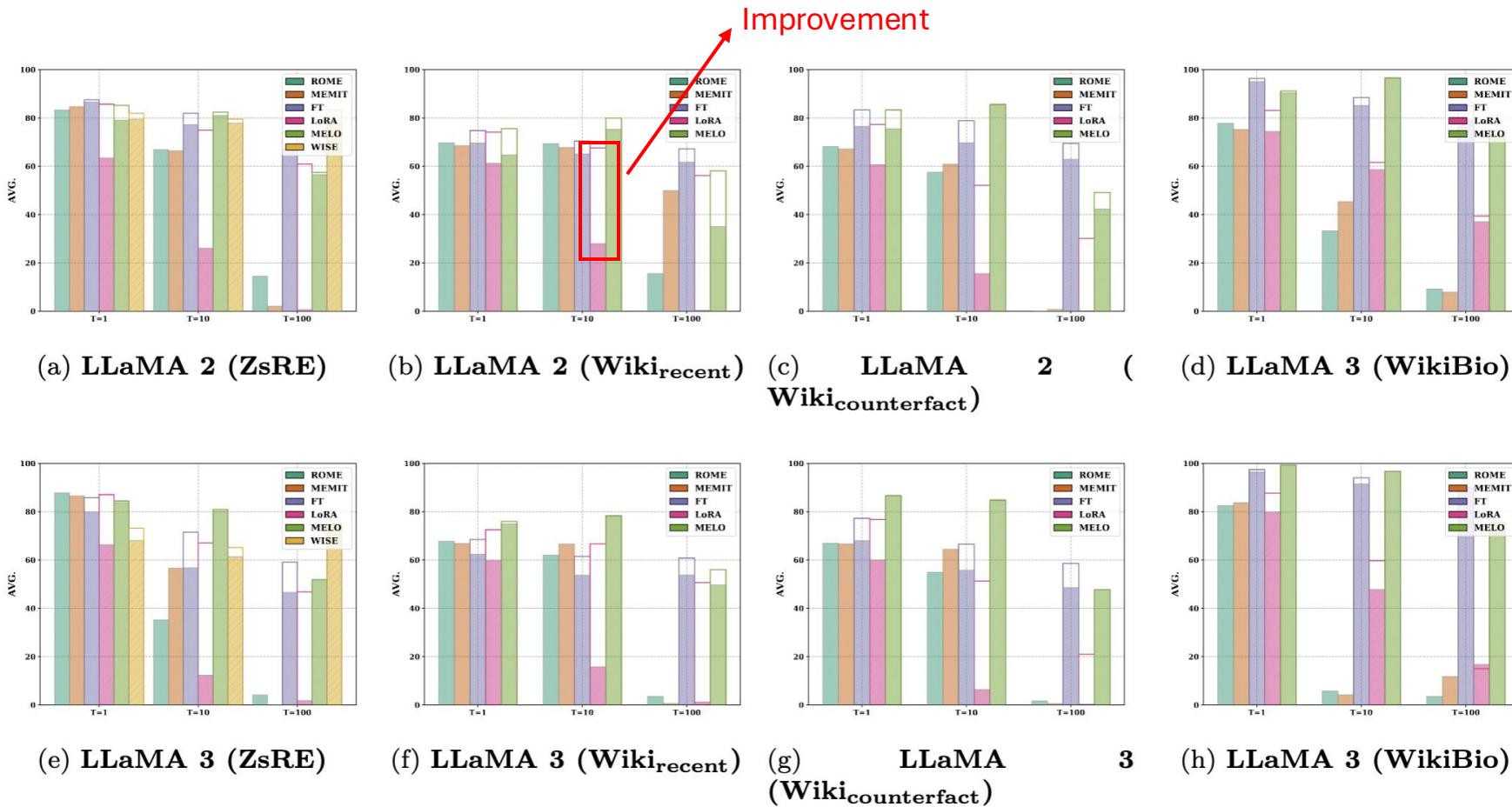


Figure 3: Continual Editing performance under different sequence length  $T$ . Solid and transparent bars show performance with and without OVERTONE. Unfilled area marks the performance gap. ROME and MEMIT didn't use OVERTONE.

# Outline

- Model Efficiency
  - Parameter-efficient Fine-tuning
  - Model Compression
  - Inference Acceleration
- Data Efficiency
  - Zero/few-shot Learning
  - Knowledge Editing
  - Retrieval-augmented Generation

# Failure Cases of Pre-trained Large Language Models



List 5 most cited papers authored by  
**Geoffrey Hinton**



1. "Learning Representations by Back-Propagating Errors" ...
2. "A Fast Learning Algorithm for Deep Belief Nets" ...
3. "Reducing the Dimensionality of Data with Neural Networks" ...
4. "Dropout: A Simple Way to Prevent Neural Networks from Overfitting" ...
5. "Rectified Linear Units Improve Restricted Boltzmann Machines" ...



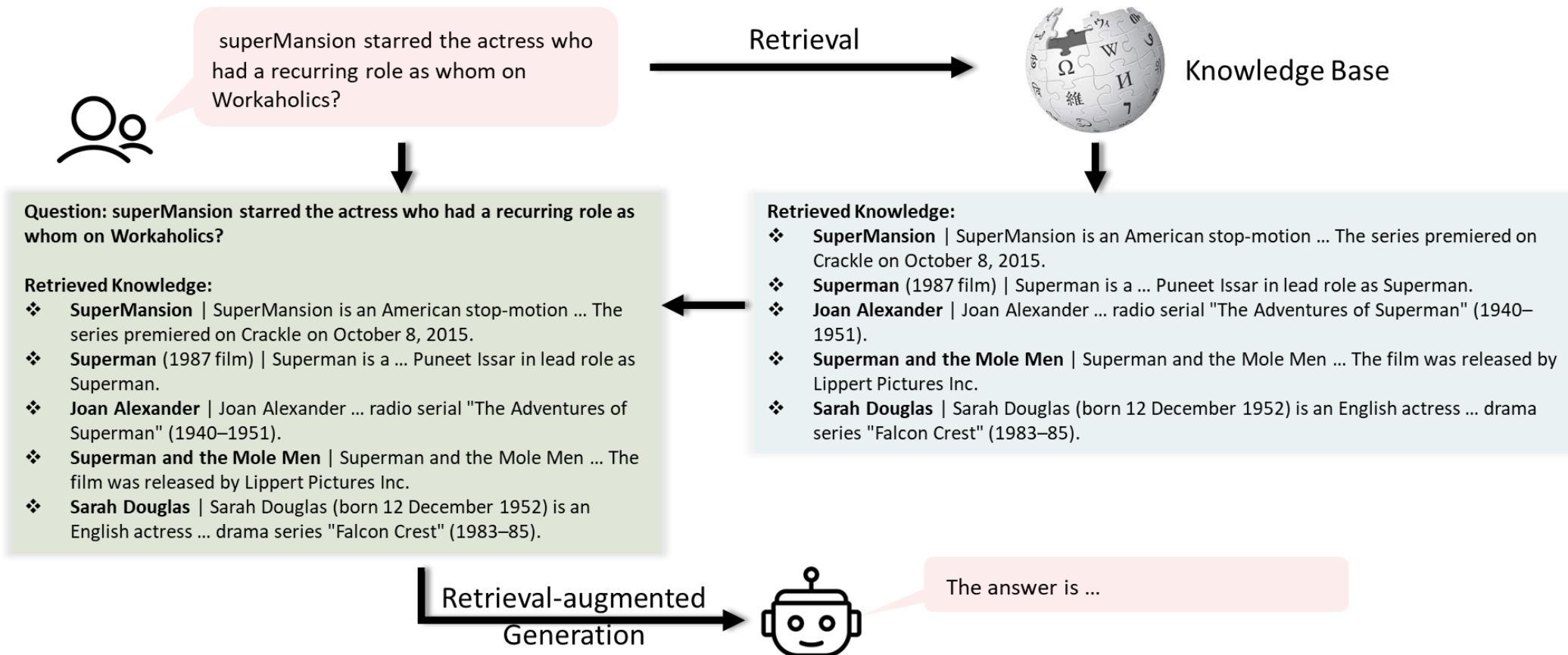
which team is Luka Dončić in?



Luka Dončić is a professional basketball player for the Dallas Mavericks in the NBA.

- LLMs can't memorize all (long-tail) knowledge in their parameters [\[Akari A et al., 2023\]](#).
- LLMs' knowledge is easily outdated and hard to update frequently [\[Akari A et al., 2023\]](#).

# Retrieve-then-generate

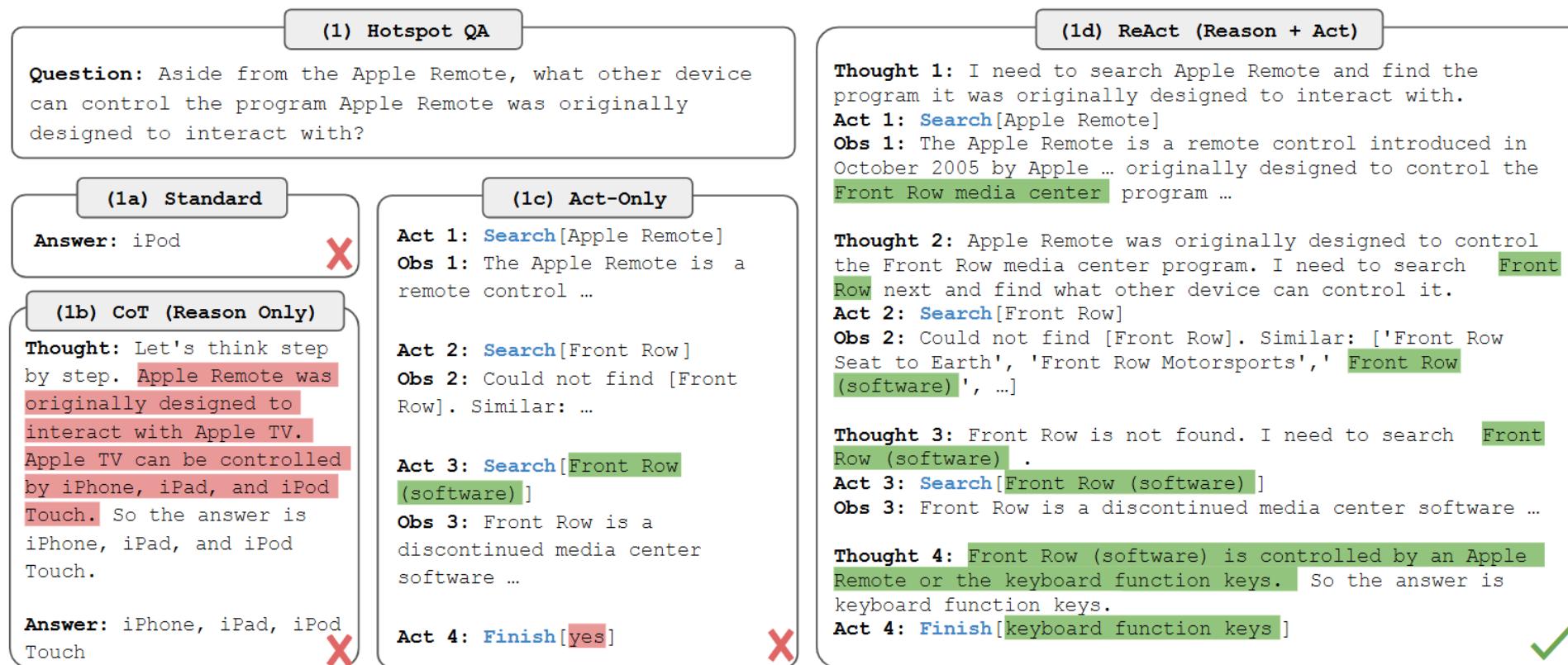


# Challenge

- Complex queries make it difficult to retrieve documents accurately.
- Irrelevant information will mislead the LLM.

# ReAct

- Ask the LLM to propose the question needing to solve (Thought), the action (Search/Finish), and intermediate answer (Obs).
- The LLM will decompose the input query into several sub-questions and take actions dynamically.



# ReAct

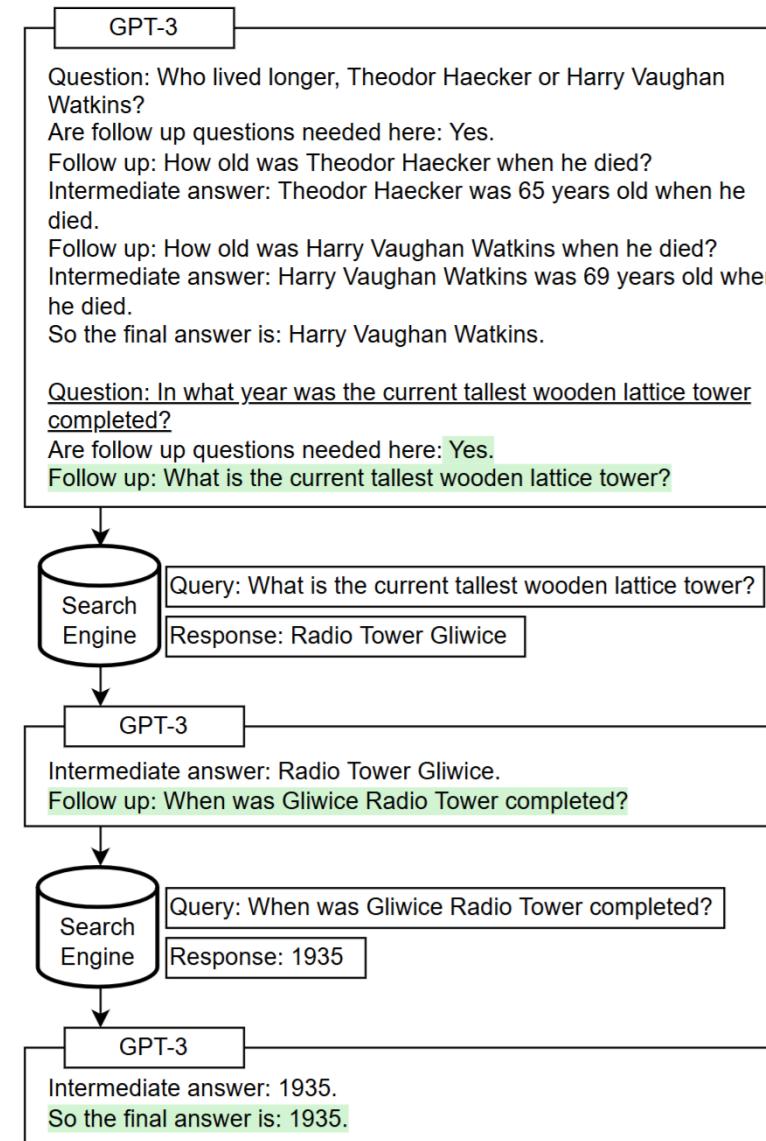
---

Prompt Method <sup>a</sup>	HotpotQA (EM)	Fever (Acc)
Standard	28.7	57.1
CoT ( <a href="#">Wei et al., 2022</a> )	29.4	56.3
CoT-SC ( <a href="#">Wang et al., 2022a</a> )	33.4	60.4
Act	25.7	58.9
ReAct	27.4	60.9
CoT-SC → ReAct	34.2	<b>64.6</b>
ReAct → CoT-SC	<b>35.1</b>	62.0
<b>Supervised SoTA<sup>b</sup></b>	67.5	89.5

# Self-ask

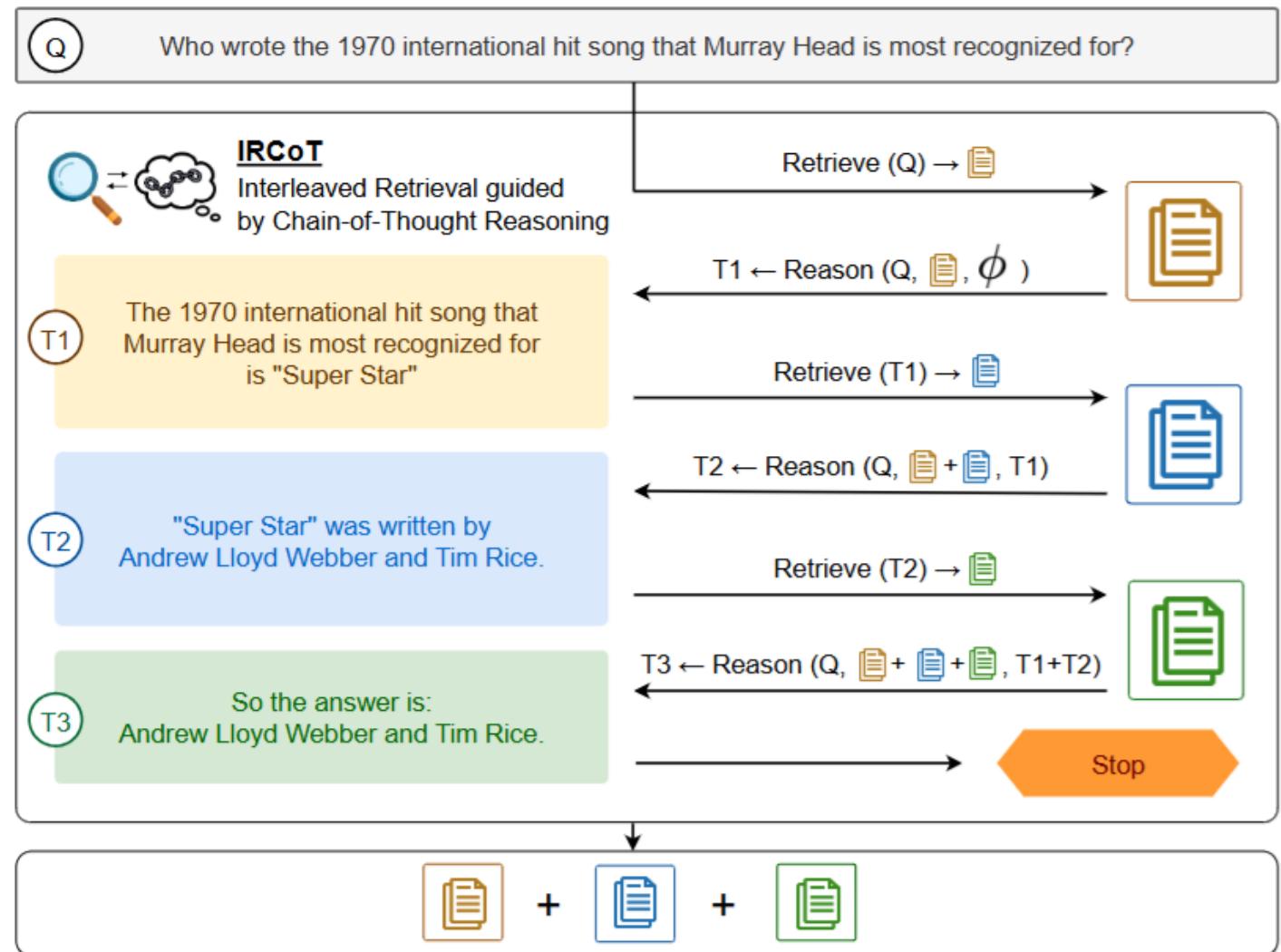
---

- Share the similar spirit with ReAct.
- Decompose original question into follow-up questions and combine retrieval.

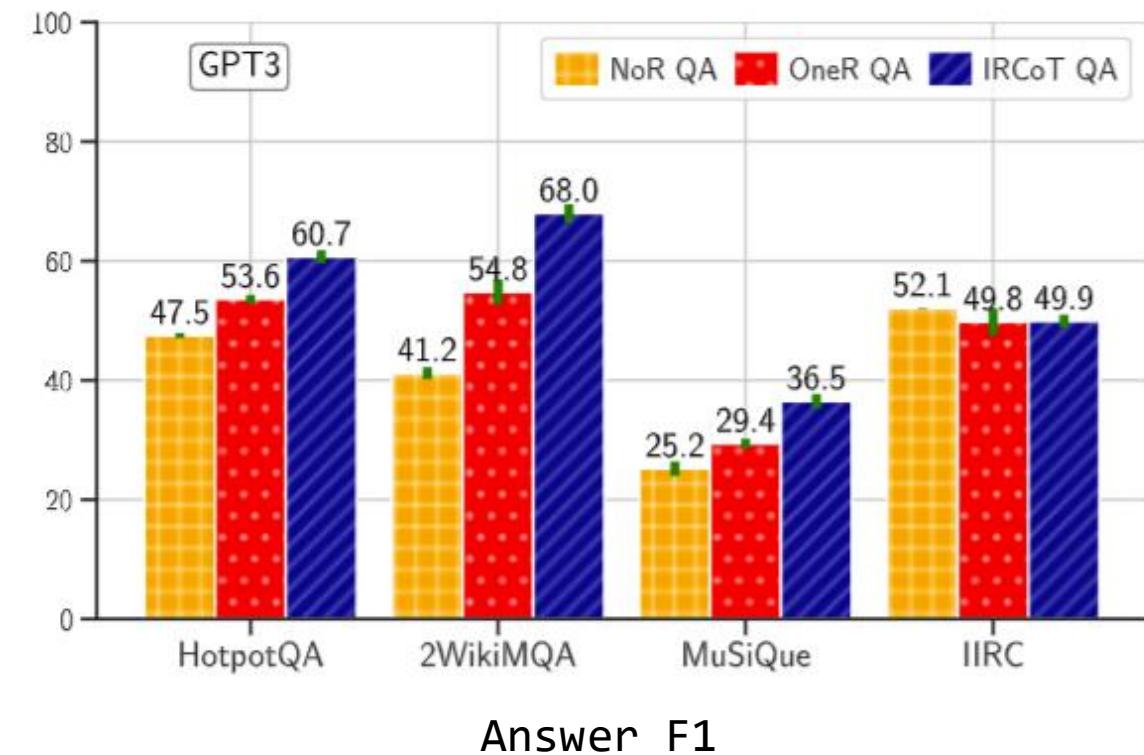
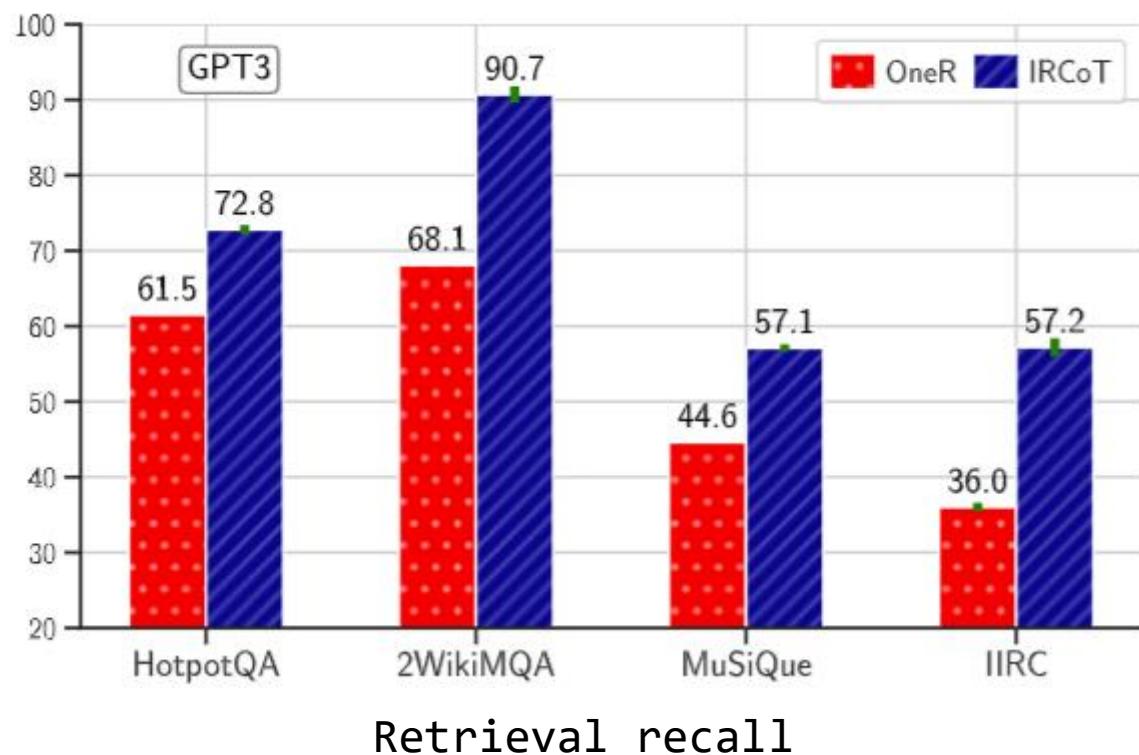


# IRCoT

- What to retrieve depends on what has already been derived, which in turn may depend on what was previously retrieved



# IRCoT



# ITER- RETEGEN

- Iterate retrieval and generation
- leverage the model output from the previous iteration as a specific context to help retrieve more relevant knowledge

Question:  $q$

What is the height of the player who won the 2015 AFL Rising Star award?

Retrieval:  $q \rightarrow \mathcal{D}_q$

Title: 2015 AFL Rising Star Context: The NAB AFL Rising Star award is given annually to a stand out young player in the Australian Football League. The award was won by Jesse Hogan of Melbourne

Retrieval-Augmented Generation:  $\mathcal{D}_q \parallel q \rightarrow y_1$

The 2015 AFL Rising Star award was won by Jesse Hogan of Melbourne. Jesse Hogan is a professional Australian rules footballer. He is 198 cm tall. So the answer is 198 cm

Question:  $q$

What is the height of the player who won the 2015 AFL Rising Star award?

Generation-Augmented Retrieval:  $y_1 \parallel q \rightarrow \mathcal{D}_{y_1 \parallel q}$

Title: Jesse Hogan Context: Jesse Hogan ... playing for the Melbourne Football Club. A key forward, Hogan is 1.95 m tall ... made his AFL debut in the 2015 season and won the Ron Evans Medal as the AFL Rising Star

Retrieval-Augmented Generation:  $\mathcal{D}_{y_1 \parallel q} \parallel q \rightarrow y_2$

The 2015 AFL Rising Star award was won by Jesse Hogan of Melbourne. Jesse Hogan is 1.95 m tall. So the answer is 1.95 m

# ITER-RETEGEN

Method	HotPotQA			2WikiMultiHopQA			MuSiQue			Bamboogle			Feverous		StrategyQA	
	EM	F1	Acc <sup>†</sup>	EM	F1	Acc <sup>†</sup>	EM	F1	Acc <sup>†</sup>	EM	F1	Acc <sup>†</sup>	Acc	Acc <sup>†</sup>	Acc	Acc <sup>†</sup>
Without Retrieval																
Direct	21.9	36.8	44.8	21.3	29.2	33.9	7.0	18.7	15.8	11.2	24.4	28.0	60.1	60.1	66.5	66.7
CoT	30.0	44.1	50.0	30.0	39.6	44.0	19.4	30.9	28.6	<b>43.2</b>	<b>51.1</b>	60.0	59.8	59.8	71.0	71.0
With Retrieval																
Direct	31.6	44.7	53.3	27.3	35.4	43.6	13.9	28.2	26.5	17.6	31.8	43.2	69.8	69.8	65.6	65.6
ReAct	24.9	44.7	61.1	28.0	38.5	45.9	23.4	37.0	37.9	21.8	31.0	40.3	66.4	66.4	66.9	66.9
Self-Ask	36.8	55.2	64.8	<b>37.3</b>	<b>48.8</b>	55.9	<b>27.6</b>	41.5	<b>42.9</b>	31.5	41.2	54.8	70.7	70.7	70.2	70.2
DSP	43.8	55.0	60.8	-	-	-	-	-	-	-	-	-	-	-	-	-
ITER-RETEGEN 1	<u>39.2</u>	53.9	<u>65.5</u>	33.7	45.2	55.4	24.2	38.6	38.1	<u>36.8</u>	<u>47.7</u>	<u>57.6</u>	67.0	67.0	<u>72.0</u>	<u>72.0</u>
ITER-RETEGEN 2	<u>44.1</u>	<u>58.6</u>	<u>71.2</u>	34.9	47.0	<u>58.1</u>	26.4	41.1	41.0	<u>38.4</u>	<u>48.7</u>	<u>59.2</u>	68.8	68.8	<u>73.0</u>	<u>73.0</u>
ITER-RETEGEN 3	<u>45.2</u>	<u>59.9</u>	<u>71.4</u>	34.8	47.8	<u>58.3</u>	25.7	41.4	40.8	<u>37.6</u>	<u>47.0</u>	<u>59.2</u>	69.0	69.0	<u>72.3</u>	<u>72.3</u>
ITER-RETEGEN 4	<u>45.8</u>	<b>61.1</b>	<b>73.4</b>	36.0	47.4	<u>58.5</u>	26.7	<u>41.8</u>	40.8	<u>38.4</u>	<u>49.6</u>	<u>60.0</u>	<b>71.5</b>	<b>71.5</b>	<u>73.8</u>	<u>73.8</u>
ITER-RETEGEN 5	<u>45.2</u>	<u>60.3</u>	<u>72.8</u>	35.5	47.5	<u>58.8</u>	25.7	40.7	39.6	<u>39.2</u>	<u>49.7</u>	<b>60.8</b>	70.3	70.3	<u>73.2</u>	<u>73.2</u>
ITER-RETEGEN 6	<b>45.9</b>	<u>61.0</u>	<u>73.3</u>	35.5	48.1	<b>59.4</b>	25.9	40.5	39.8	<u>40.0</u>	<u>50.0</u>	<u>59.2</u>	<u>70.9</u>	<u>70.9</u>	<u>72.4</u>	<u>72.4</u>
ITER-RETEGEN 7	<u>45.1</u>	<u>60.4</u>	<u>72.9</u>	35.5	47.4	<u>58.4</u>	26.1	<b>42.0</b>	41.0	<u>40.0</u>	<u>50.7</u>	<b>60.8</b>	70.5	70.5	<b>74.1</b>	<b>74.1</b>

QA Performance

Iteration	1	2	3	4	5	6	7
HotPotQA	49.5	66.1	65.7	66.5	66.7	66.7	67.1
2WikiMultiHopQA	29.0	45.2	46.2	46.7	45.8	45.8	46.5
MuSiQue	18.6	32.3	32.3	33.7	32.7	33.5	32.9
Bamboogle	20.8	36.0	36.8	36.0	35.2	36.0	36.0

Retrieval Recall

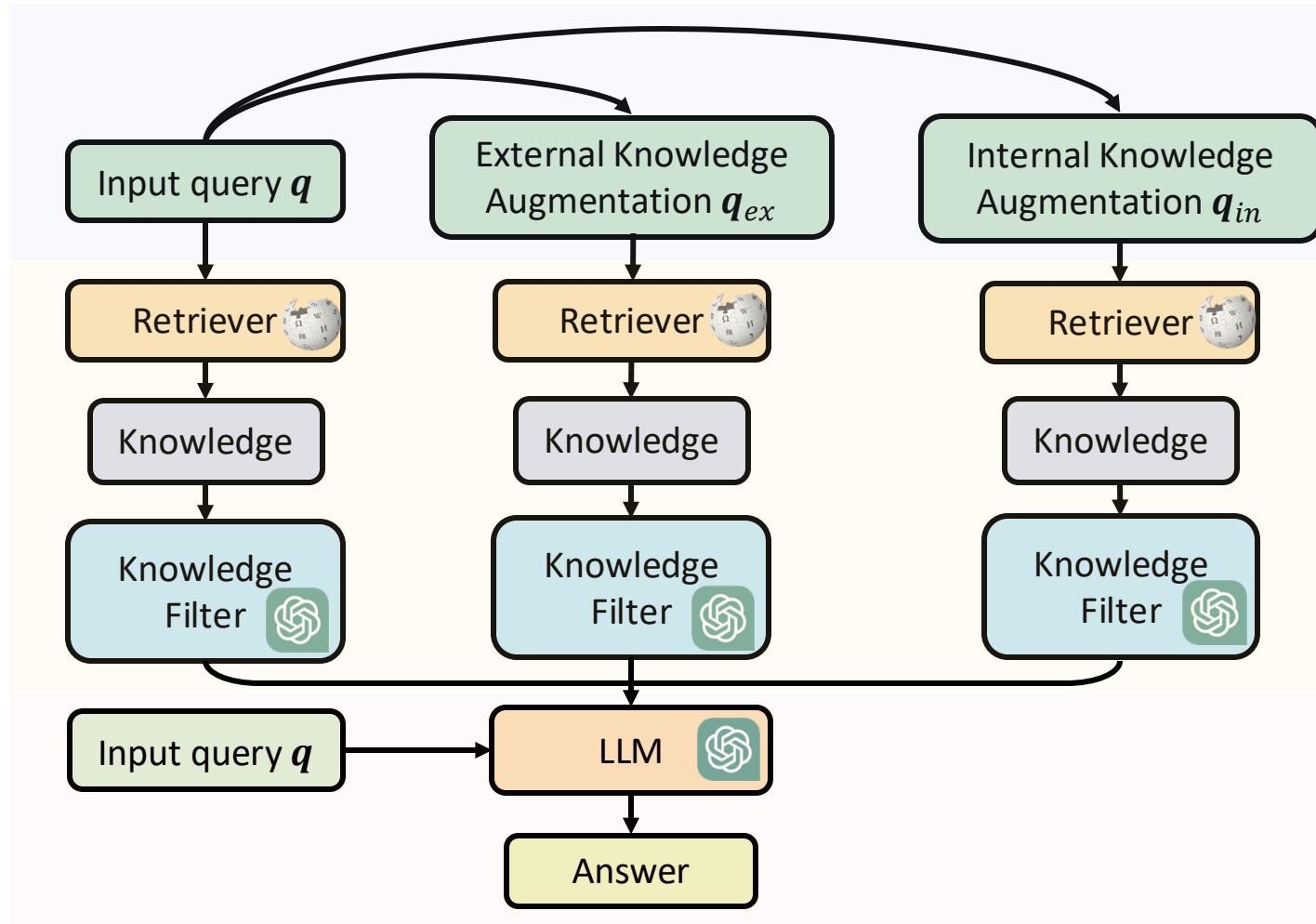
# Challenge

- Complex queries make it difficult to retrieve documents accurately.
- Irrelevant information will mislead the LLM

# RetRobust

- In-context learning method:
  - Leverage an extra natural language inference (NLI) model
  - NLI model assesses whether retrieved context is irrelevant
  - If irrelevant, the model will generate without using the retrieved information
- Fine-tuning method:
  - Leverage an LLM to generate training data that mixes both relevant and irrelevant contexts
  - Standard Supervised Fine-tuning on the generated data

# BlendFilter



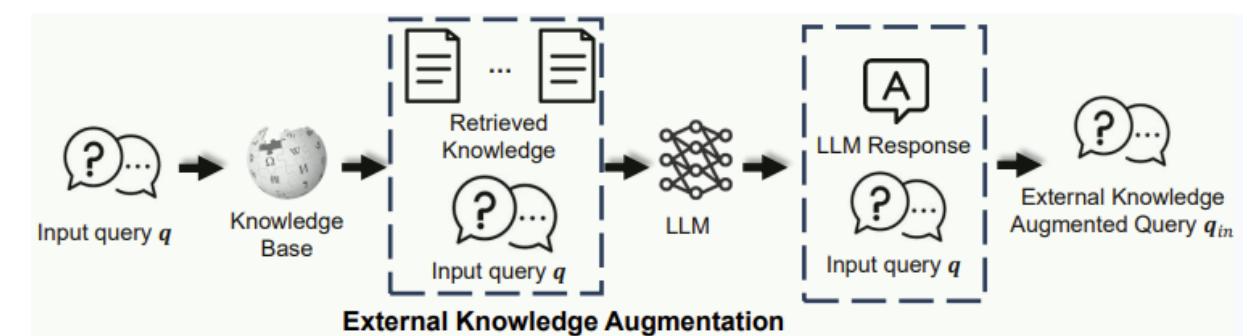
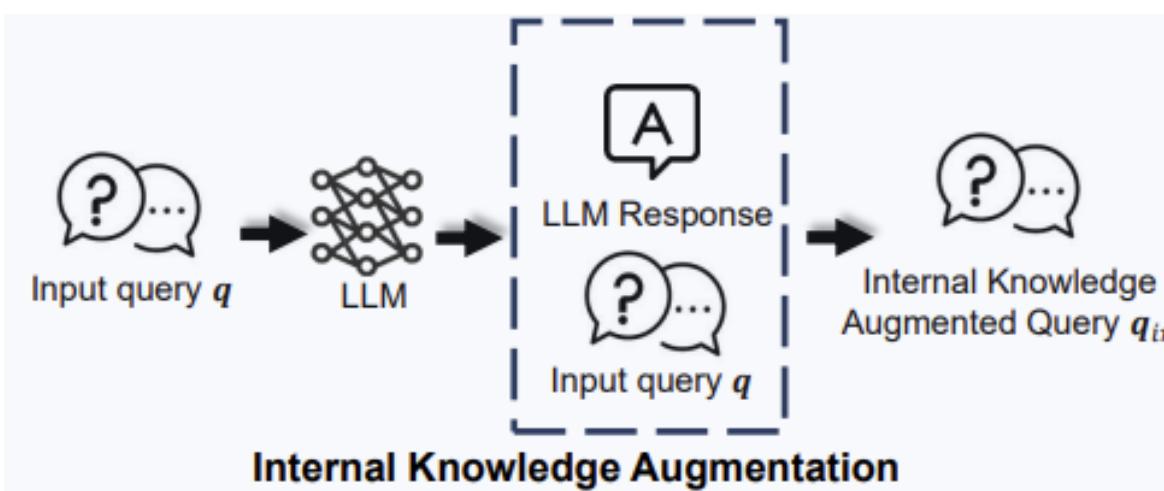
Query Blending

Knowledge Filtering

Answer Generation

# BlendFilter

## Query Blending



# BlendFilter

## Knowledge Filtering

Forget your knowledge about {topic}. Please only consider the knowledge below.

knowledge 0 : {Retrieved\_knowledge0}

knowledge 1 : {Retrieved\_knowledge1}

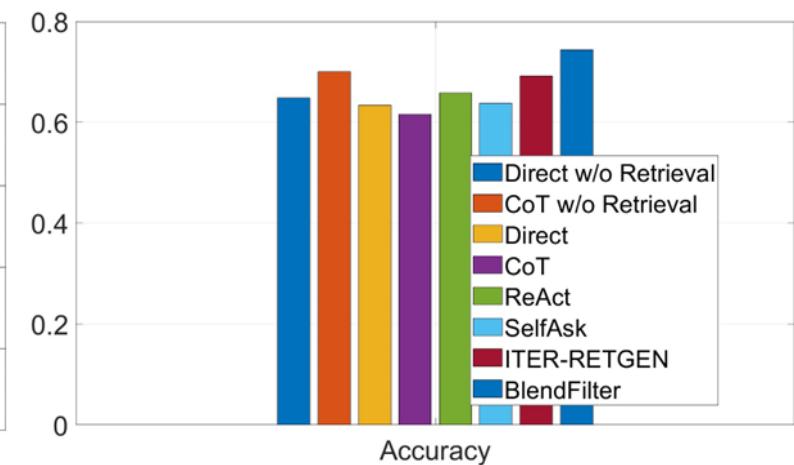
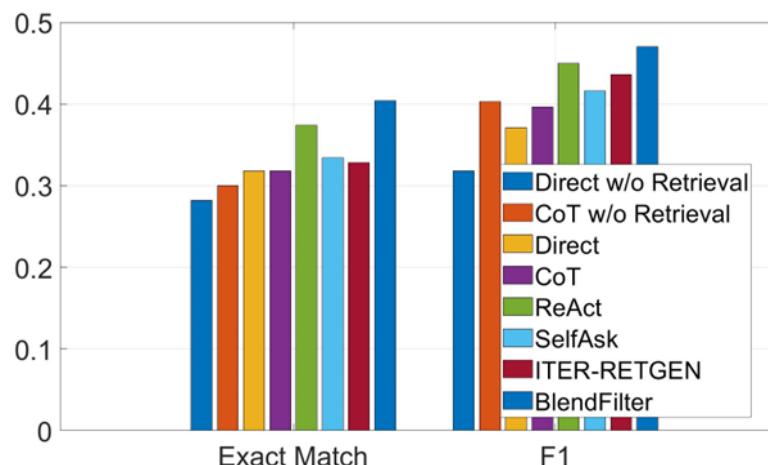
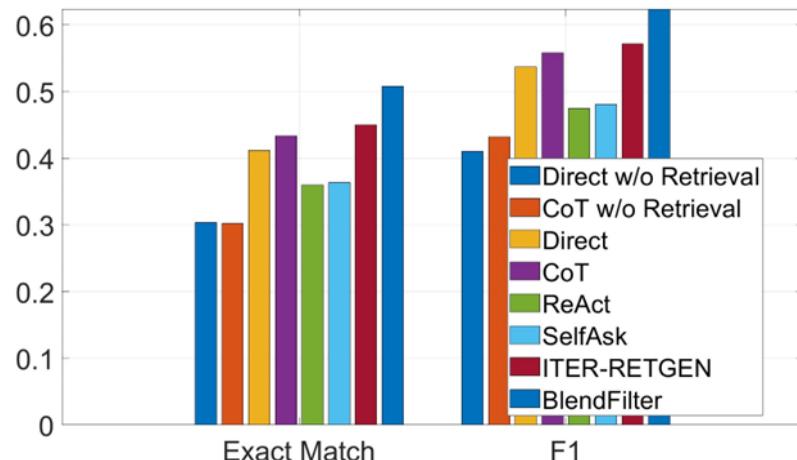
knowledge 2 : {Retrieved\_knowledge2}

knowledge 3 : {Retrieved\_knowledge3}

knowledge 4 : {Retrieved\_knowledge4}

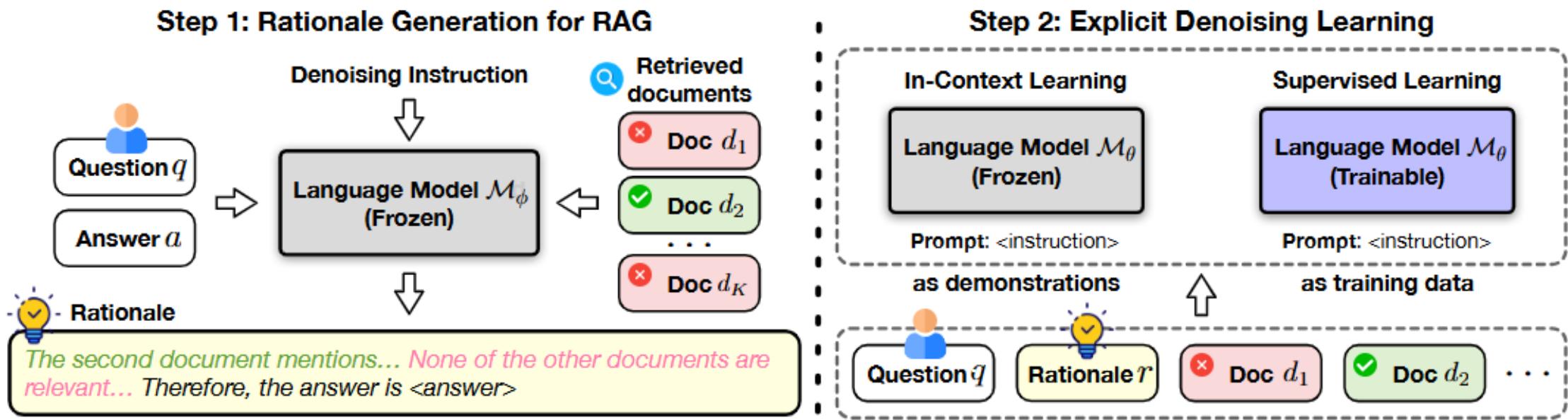
Please check the relevance between {question} and knowledges 0-4 one by one, remove the irrelevant ones and show me the relevant ones. There may be multiple relevant ones. Please take a deep breath and do it step by step.

# BlendFilter



# InstructRAG

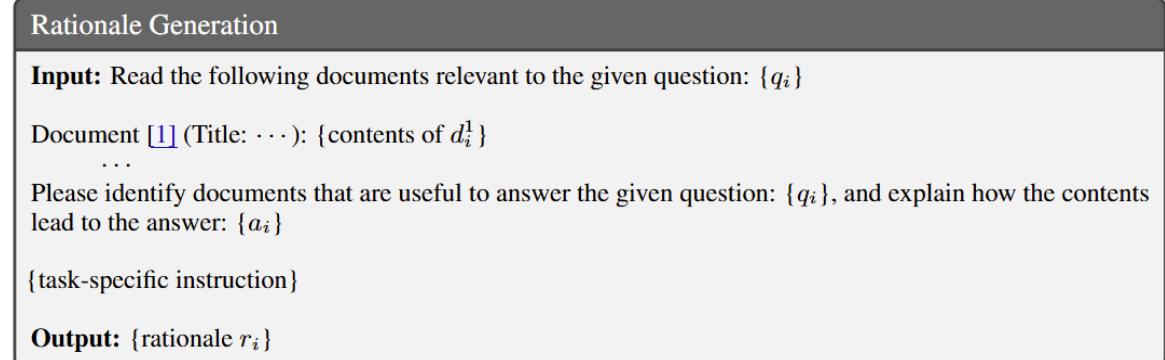
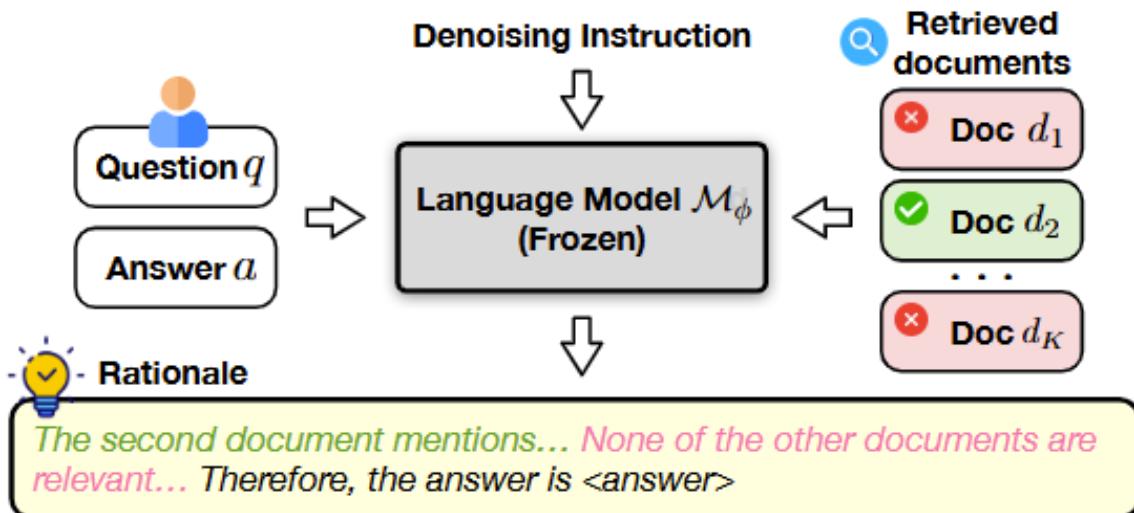
- Enable the LLM to generate rationale for the retrieved documents
- Conduct SFT on the generated rationale



# InstructRAG

## Rationale Generation via Instruction-Following

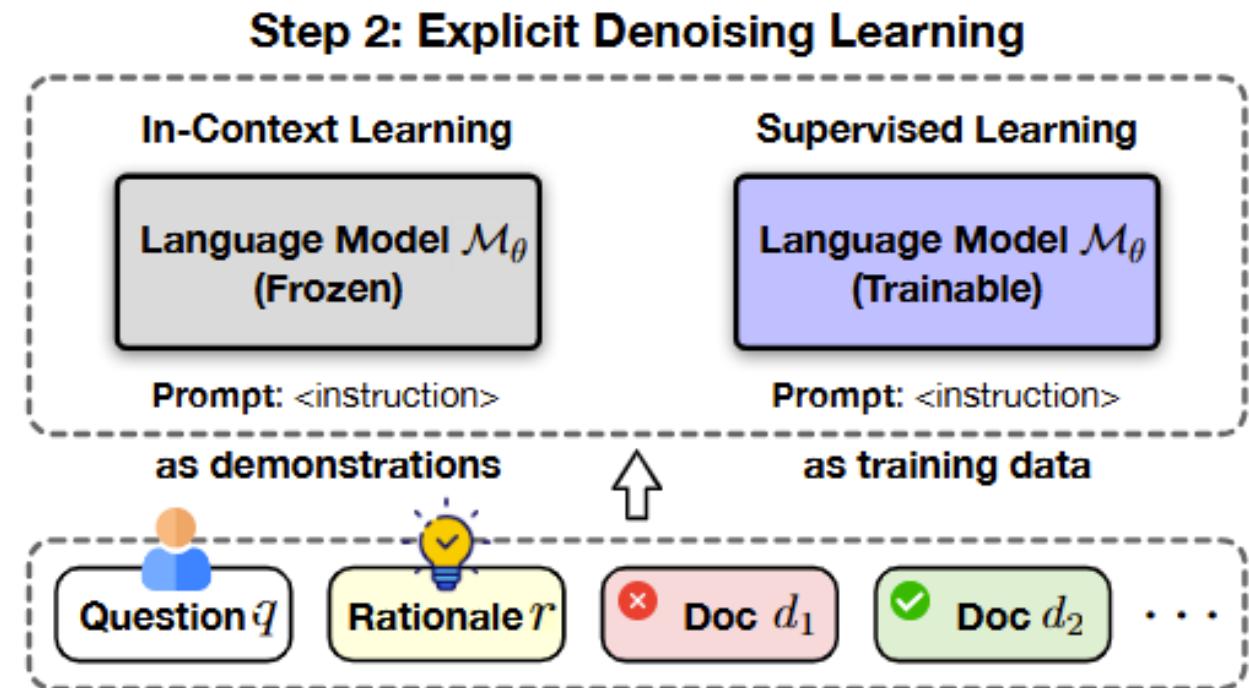
### Step 1: Rationale Generation for RAG



# InstructRAG

- Learning Denoising Rationales in RAG

$$\max_{\theta} \mathbb{E}_{q,r \sim \mathcal{T}^+} \log p_{\theta}(r|q, D)$$



# InstructRAG



Method	PopQA (acc)	TriviaQA (acc)	NQ (acc)	MultiHopQA (acc)	ASQA (em)	ASQA (pre)	ASQA (rec)
<i>Baselines w/o Retrieval</i>							
<b>Vanilla Zero-shot Prompting</b>							
ChatGPT*	29.3	74.3	—	—	35.3	—	—
Llama-3-Instruct <sub>8B</sub>	22.8	69.4	46.6	45.6	30.6	—	—
Llama-3-Instruct <sub>70B</sub>	28.9	80.6	57.9	57.5	39.1	—	—
<i>RAG w/o Training</i>							
<b>In-Context RALM [70]</b>							
ChatGPT*	50.8	65.7	—	—	40.7	65.1	76.6
Llama-3-Instruct <sub>8B</sub>	62.3	71.4	56.8	43.4	40.0	62.1	66.4
Llama-3-Instruct <sub>70B</sub>	63.8	76.3	60.2	51.2	43.1	62.9	67.6
<b>Few-Shot Demo. w/ Instruction</b>							
Llama-3-Instruct <sub>8B</sub>	63.1	74.2	60.1	45.3	42.6	55.0	64.4
Llama-3-Instruct <sub>70B</sub>	63.9	79.1	62.9	53.9	45.4	49.3	57.1
<b>INSTRUCTRAG-ICL</b>							
Llama-3-Instruct <sub>8B</sub>	64.2	76.8	62.1	50.4	44.7	<b>70.9</b>	<b>74.1</b>
Llama-3-Instruct <sub>70B</sub>	<b>65.5</b>	<b>81.2</b>	<b>66.5</b>	<b>57.3</b>	<b>47.8</b>	69.1	71.2
<i>RAG w/ Training</i>							
<b>Vanilla Supervised Fine-tuning</b>							
Llama-3-Instruct <sub>8B</sub>	61.0	73.9	56.6	56.1	43.8	—	—
<b>Self-RAG [3]</b>							
Llama-2 <sub>7B</sub>	55.8	68.9	42.4	35.9	30.0	66.9	67.8
Llama-2 <sub>13B</sub>	56.3	70.4	46.4	36.0	31.4	<b>70.3</b>	<b>71.3</b>
Llama-3-Instruct <sub>8B</sub>	55.8	71.4	42.8	32.9	36.9	69.7	69.7
<b>RetRobust [105]</b>							
Llama-2 <sub>13B</sub>	—	—	39.6	51.5	—	—	—
Llama-3-Instruct <sub>8B</sub>	56.5	71.5	54.2	54.7	40.5	—	—
<b>INSTRUCTRAG-FT</b>							
Llama-3-Instruct <sub>8B</sub>	<b>66.2</b>	<b>78.5</b>	<b>65.7</b>	<b>57.2</b>	<b>47.6</b>	65.7	70.5

# ROSERAG

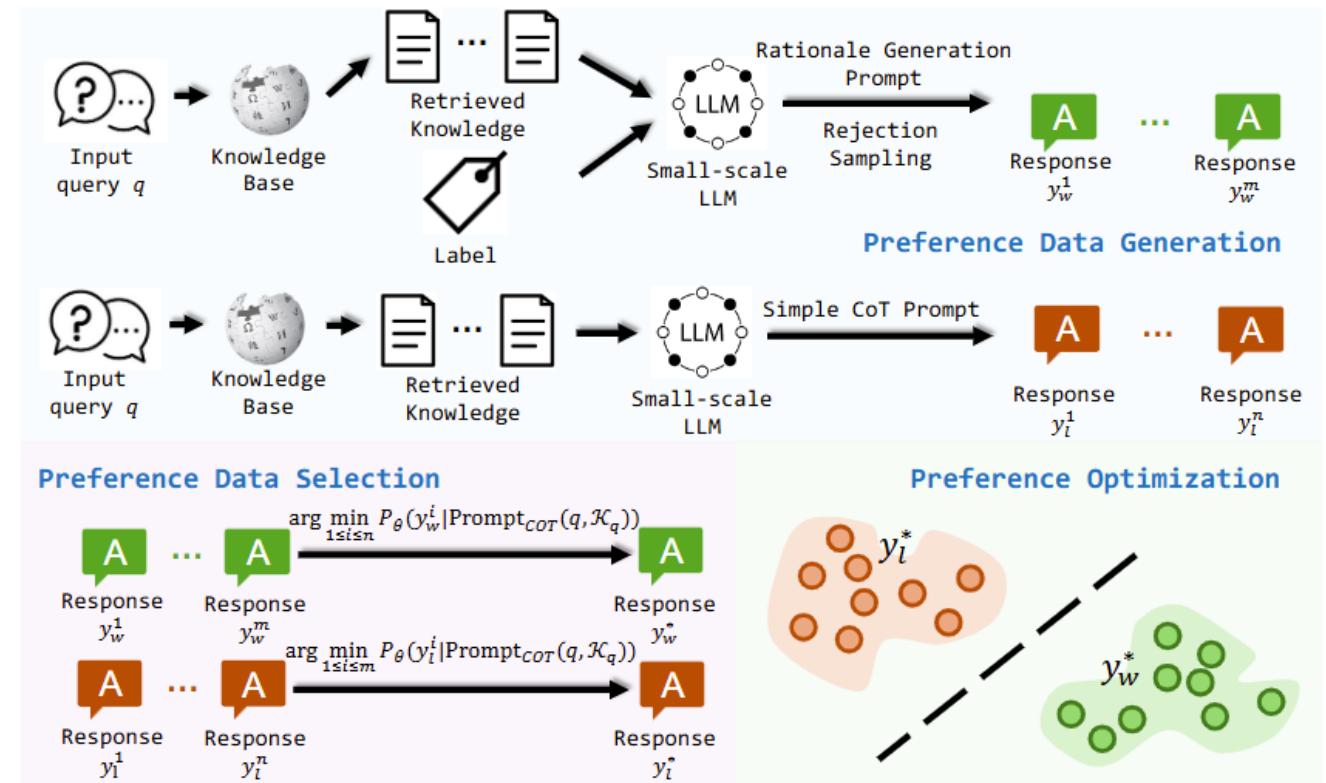
---

## Motivation

- Small-scale LLMs (SLMs) have limited reasoning and generalization capacity.
- SFT tends to mimic the behavior present in the training data, making it highly sensitive to noise and prone to overfitting.

# ROSERAG

- Three stages:
  - Preference Data Generation
  - Preference Data Selection
  - Preference Optimization



# ROSERAG

## Preference Data Generation

- Preferred response  $y_w$ 
  - Ask the SLM to generate reasoning process given retrieved knowledge, the question, and the ground-truth answer
- Rejection sampling
- Non-preferred response  $y_l$ 
  - Generate from vanilla CoT response

### Rationale Generation

**System Prompt:** You are a useful assistant. I will provide one question, several pieces of knowledge (which may be related or unrelated to the question), and the answer to the question. Please explain your reasoning process in a single paragraph consisting of no more than four sentences. If the provided knowledge is insufficient, you may make an informed guess, but do not respond with "Unknown".

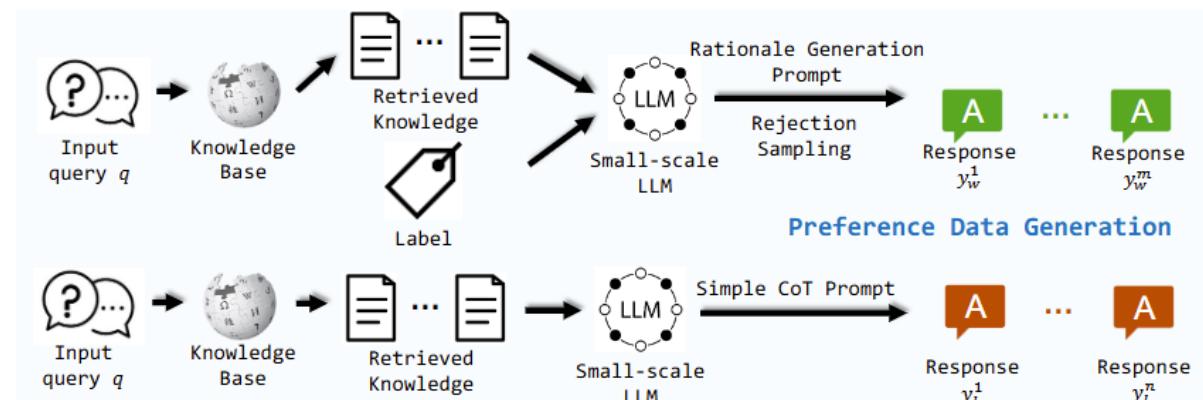
**User Prompt:** Knowledge:  $\mathcal{K}_q$

Question:  $q$

Answer:  $a^*$

**Assistant Prompt:** Let's think step by step.

**Output:** {rationale  $r$ }



# ROSERAG

## Preference Data Selection

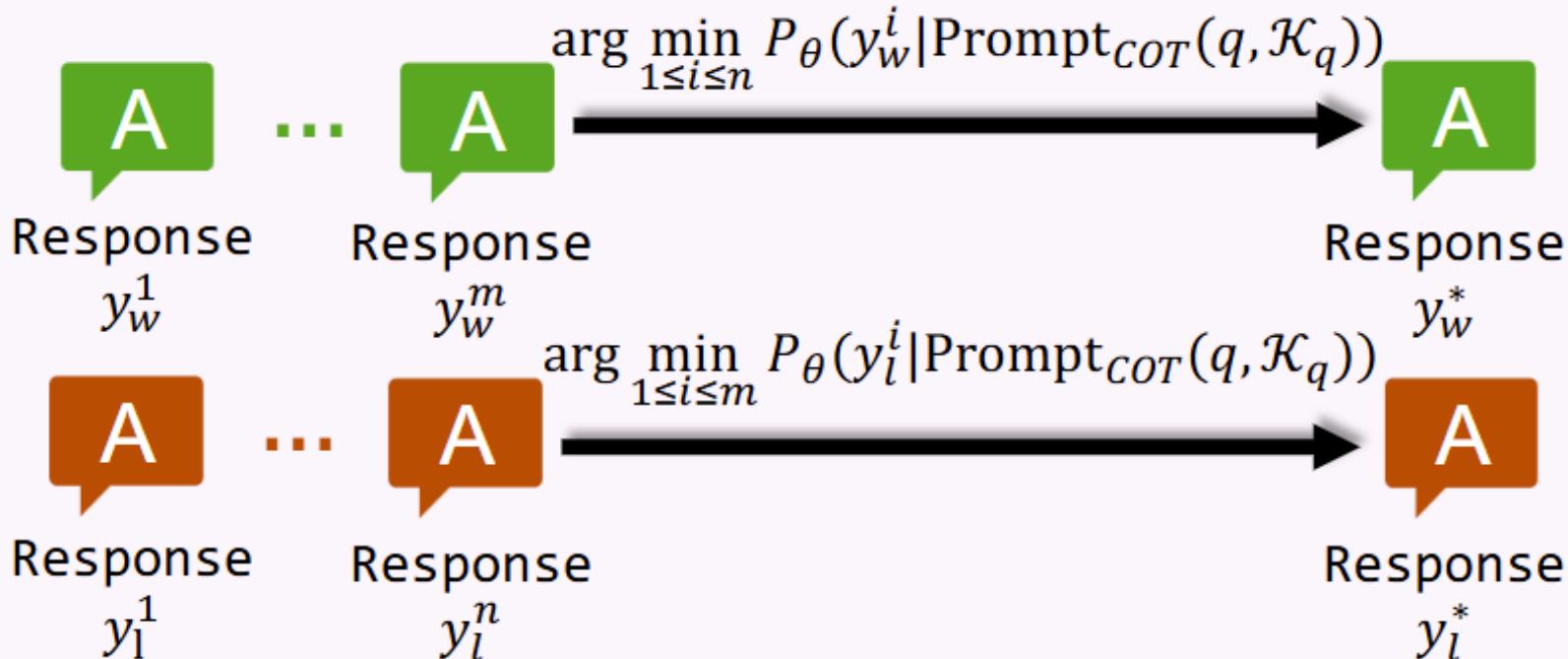
- Least confident preferred data

$$y_w^* = \arg \min_{1 \leq i \leq n} P_\theta(y_w^i | \text{Prompt}_{\text{CoT}}(q, \mathcal{K}_q))$$

- Most confident preferred data

$$y_l^* = \arg \min_{1 \leq i \leq m} P_\theta(y_l^i | \text{Prompt}_{\text{CoT}}(q, \mathcal{K}_q))$$

## Preference Data Selection



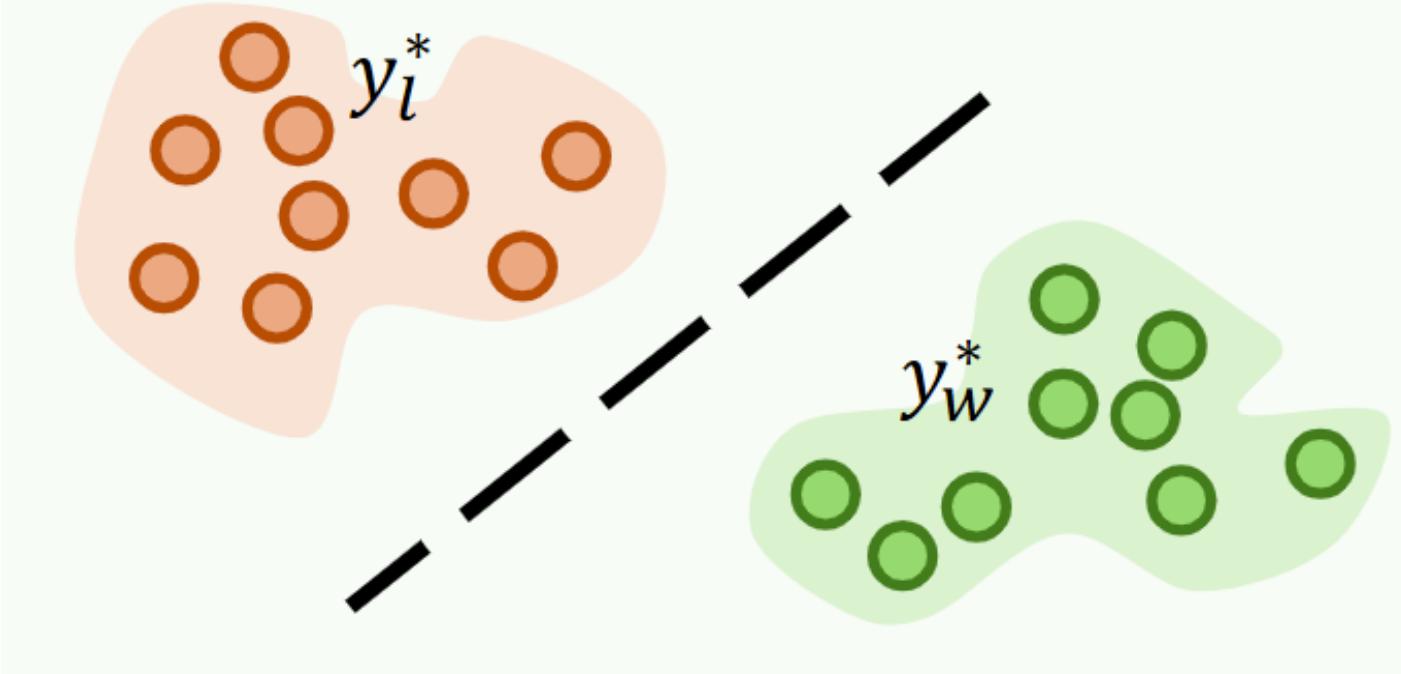
# ROSERAG

## Preference Optimization

- Use preference optimization (e.g., ORPO, PPO) to enable the SLM to align and distinguish the preference data

$$\min_{\theta} \mathbb{E}_{(x, y_w^*, y_l^*)} \ell_{ORPO}(x, y_w^*, y_l^*)$$

## Preference Optimization



# ROSERAG



Method	HotPotQA		2WikiMultiHopQA		StrategyQA	
	EM	F1	EM	F1	Acc	F1
CoT	13.0	20.9	18.8	23.5	55.9	17.9
vanilla RAG	27.2	38.4	6.2	11.4	55.5	23.9
ReAct	14.9	25.9	9.6	22.2	55.0	48.8
SelfAsk	20.4	32.8	19.2	25.3	51.1	42.9
BlendFilter	26.4	37.5	19.4	24.2	<u>59.4</u>	<u>45.0</u>
InstructRAG	<u>31.2</u>	<u>39.6</u>	<u>23.6</u>	<u>26.9</u>	53.7	39.1
RetRobust	16.8	24.6	13.0	19.9	51.5	37.3
ASTUTE	21.4	27.8	19.0	23.4	55.0	12.0
ICL+RAG	28.6	39.3	21.8	26.4	56.3	27.5
<b>ROSERAG</b>	<b>34.8</b>	<b>44.8</b>	<b>31.6</b>	<b>35.0</b>	<b>59.8</b>	<b>52.1</b>

		Positive		
		w/o Selection	Minimal	Maximal
Negative	w/o Selection	30.6/39.9	33.2/43.1	27.6/37.1
	Minimal	29.0/39.2	33.2/42.2	26.6/35.3
	Maximal	33.0/43.6	<b>34.8/44.8</b>	28.8/38.4

# Future Directions

- Model compression and pruning
  - Pruning, quantization, knowledge distillation, low-rank adaptation.
- Efficient training techniques
  - Alternative optimization that reduces time complexity, sparse and adaptive training, federated learning.
- Data efficiency
  - Reduce the need for large-scale data by transfer learning, synthetic data generation, self-supervised/weakly supervised learning.

# Future directions: How to enable LLM to think efficiently

- **Problem**
  - LLMs often rely on long, multi-step reasoning (Chain-of-Thought), which can be time-consuming and computationally expensive.
- **Key questions**
  - When to Think Long vs. Short?
  - How can the model decide between a brief inference path or a more detailed reasoning chain?

# Future directions: Data & model efficiency in AI agents

- **Problem**

- Long-horizon, complex tasks make it challenging to rely solely on Supervised Fine-Tuning.
- High computational overhead from LLMs when there exist many sequential steps.

- **Key questions**

- How to Combine RL + Simulation?
- Can simulation-based RL reduce real-world data needs and improve adaptability?

# Future directions: AI on edge devices

## Problem

- Edge devices have strict constraints on computation, memory, and energy, making it difficult to deploy large AI models.
- Real-time or near-real-time inference demands low latency and high efficiency.

## Key questions

- How to handle on-device learning and ensure real-time response?
- Can we adapt or update models locally without relying on the cloud, and how do we maintain privacy?
- Can we ensure a holistic optimization of various resource constraints (memory, latency, energy consumption, etc.)?

# Future directions: Hardware–software co-design

- **Problem**
  - Traditional hardware is not always optimized for the unique demands of large-scale AI.
  - Mismatches between model architectures and hardware capabilities lead to inefficiencies and wasted energy.
- **Key questions**
  - How to match AI algorithms to hardware?
  - What co-optimization strategies can we use?