# LLM and PEFT

## 25.02.06 / [DS with Alumni] Seminar

DSL 12th
통계데이터사이언스학과 김민규
kimmin01@yonsei.ac.kr

# About

- 이름: 김민규 (https://kimmin-01.github.io)

- 소속: 연세대학교 통계데이터사이언스학과 통합과정

    - 2020.03 – 2025.02 연세대학교 상경대학 응용통계학과

    - 2025.03 - 현재      연세대학교 일반대학원 통계데이터사이언스학과

                        (MLAI https://mlai.yonsei.ac.kr/)

- Research Interests:

    - Causal Inference with Bayesian Deep Learning

    - LLM Reasoning

    - Invariant Representation Learning

- 기타 경험:

    - MLAI 학부연구생 (2024.07 – 2025.02)

    - Data Science Lab 12기, Yonsei Artificial Intelligence 14기

# Contents

# Introduction

## Let's Cook!



Expensive!

https://www.newsis.com/view/NISX20190423_0000629293



Inefficient!

https://www.taketwotapas.com/all-purpose-steak-seasoning-blend/

All Purpose Korean Sauce

Efficient!



KIMCHI  BIBIMBAP  BULGOGI  HOBAKJUK

...KSU  TTEOKBOKKI  JAPCHAE  GAMJATANG

TTEOK  JJAJANGMYEON  BOSSAM  SOONDAE

NAENGMYEON  GOPCHANG  SUNDUBU JJIGAE  GIMBAP

Let's Cook!

To achieve our objectives,
we should **modify the process efficiently**
while maintaining outputs!

Q) Can we apply the same strategy when fine-tuning our model?

Efficient!

## Pre-training & Fine-tuning

LoRA can even outperform full finetuning training only 2% of the parameters

ROUGE scores

Full finetuning → GPT-3 (FT)

Only tune bias vectors → GPT-3 (BitFit)

Prompt tuning → GPT-3 (PreEmbed)

Prefix tuning → GPT-3 (PreLayer)

| Model&Method | # Trainable Parameters | WikiSQL Acc. (%) | MNLI-m Acc. (%) | SAMSum R1/R2/RL |
|---|---|---|---|---|
| GPT-3 (FT) | 175,255.8M | **73.8** | 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[H]) | 7.1M | 71.9 | 89.8 | 53.0/28.9/44.8 |
| GPT-3 (Adapter[H]) | 40.1M | 73.2 | **91.5** | 53.2/29.0/45.1 |
| GPT-3 (LoRA) | 4.7M | 73.4 | **91.7** | **53.8/29.8/45.9** |
| GPT-3 (LoRA) | 37.7M | **74.0** | **91.6** | 53.4/29.2/45.1 |

Table 4: Performance of different adaptation methods on GPT-3 175B. We report the logical form validation accuracy on WikiSQL, validation accuracy on MultiNLI-matched, and Rouge-1/2/L on SAMSum. LoRA performs better than prior approaches, including full fine-tuning. The results on WikiSQL have a fluctuation around ±0.5%, MNLI-m around ±0.1%, and SAMSum around ±0.2/±0.2/±0.1 for the three metrics.
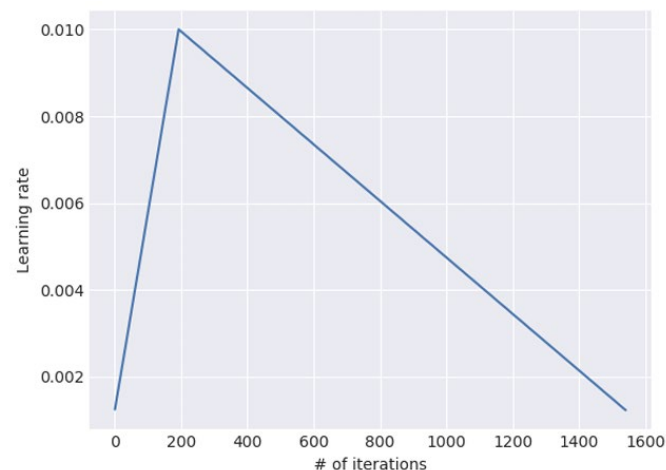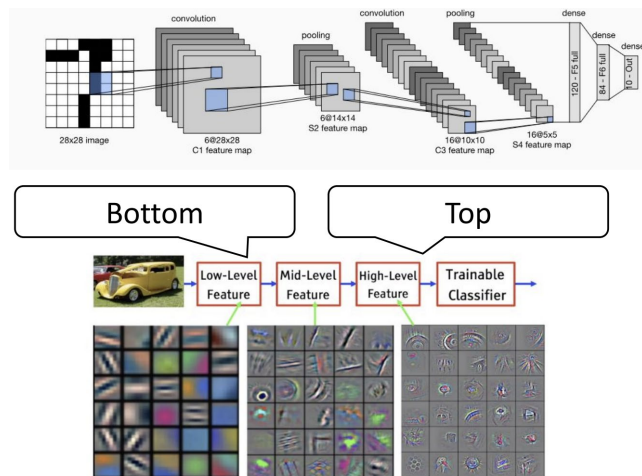
# Loss Perspective

Universal Language Model Fine-tuning for Text Classification (ACL 18')

- Gradual Unfreezing

- Discriminative Fine-tuning

- Slanted Triangular Learning Rates



(a) LM pre-training     (b) LM fine-tuning     (c) Classifier fine-tuning

Universal Language Model Fine-tuning for Text Classification (ACL 18')

1) Gradual Unfreezing: First, train only the last layer, then train the last two layers, and so on···

2) Discriminative Fine-tuning: Learning rate of bottom layer ≠ Learning rate of top layer

3) Slanted Triangular Learning Rates: Quickly converge to region of parameter space, and refine

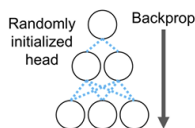## Fine-Tuning can Distort Pretrained Features and Underperform Out-of-Distribution (ICML 22')

- OOD error of fine-tuning is high when we initialize with a fixed or random head

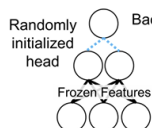- Find proper head with linear probing, then fine-tuning with that head
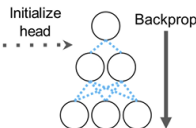


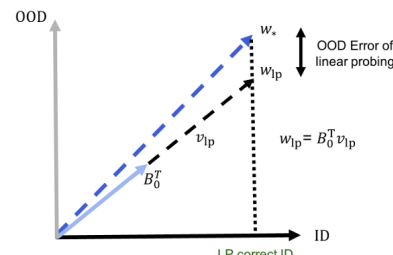| | Pretraining | (a) Fine-tuning | (b) Linear probing | (c) LP-FT |

| | (a) Fine-tuning | (b) Linear probing | (c) LP-FT |
|---|---|---|---|
| ID test | 85.1% | 82.9% | **85.7%** |
| OOD test | 59.3% | 66.2% | **68.9%** |

Average accuracies (10 distribution shifts)

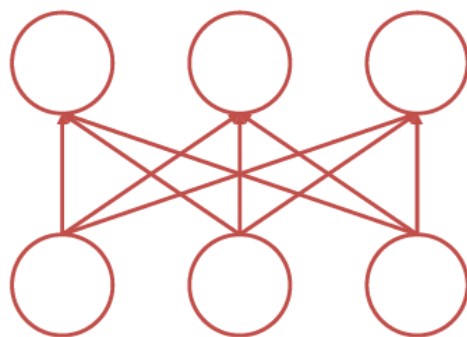$w_{lp} = B_0^T v_{lp}$

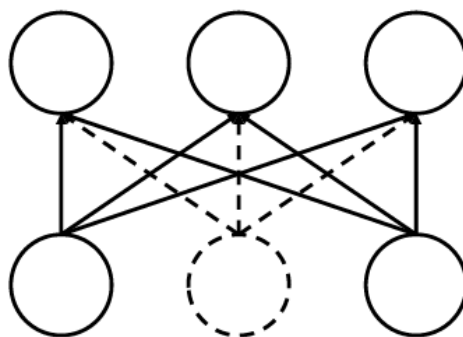(a) Toy example (Linear probing)

(b) Toy example (fine-tuning)

Mixout: Effective Regularization to Finetune Large-scale Pretrained Language Models (ICLR 20')

- Dropout: randomly kill nodes or set all connected weights to 0

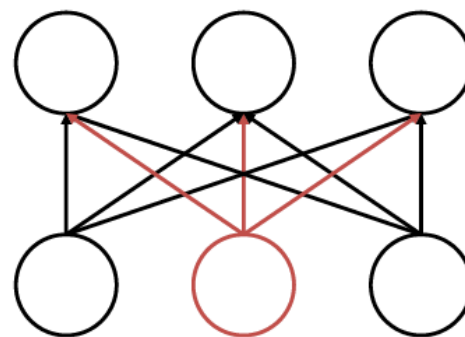- Instead, this method **randomly replace with pretrained model's parameters.'**

**Why?** The model parameter after the t-th SGD step is already far from the origin.



(a) Vanilla network at $u$     (b) Dropout network at $w$     (c) $\mathtt{mixout}(u)$ network at $w$

On the Stability of Fine-tuning BERT: Misconceptions, Explanations, and Strong Baselines (ICLR 21')

- Hypothesis

1. The instability of BERT during FT is not due to catastrophic forgetting or overfitting.

    Rather, the training process itself is unstable and does not work well.

(※ catastrophic forgetting : The tendency to completely and abruptly forget previous information)

2. This instability is caused by the following two reasons.

    1) Difficulty due to vanishing gradients by optimizer! Need to use the proper Adam optimizer.

    2) Large variance on the validation set! Ensure sufficient training up to 20 epochs.

On the Stability of Fine-tuning BERT: Misconceptions, Explanations, and Strong Baselines (ICLR 21')



(a) Perplexity of failed models   (b) Perplexity of successful models   (c) Training of failed models
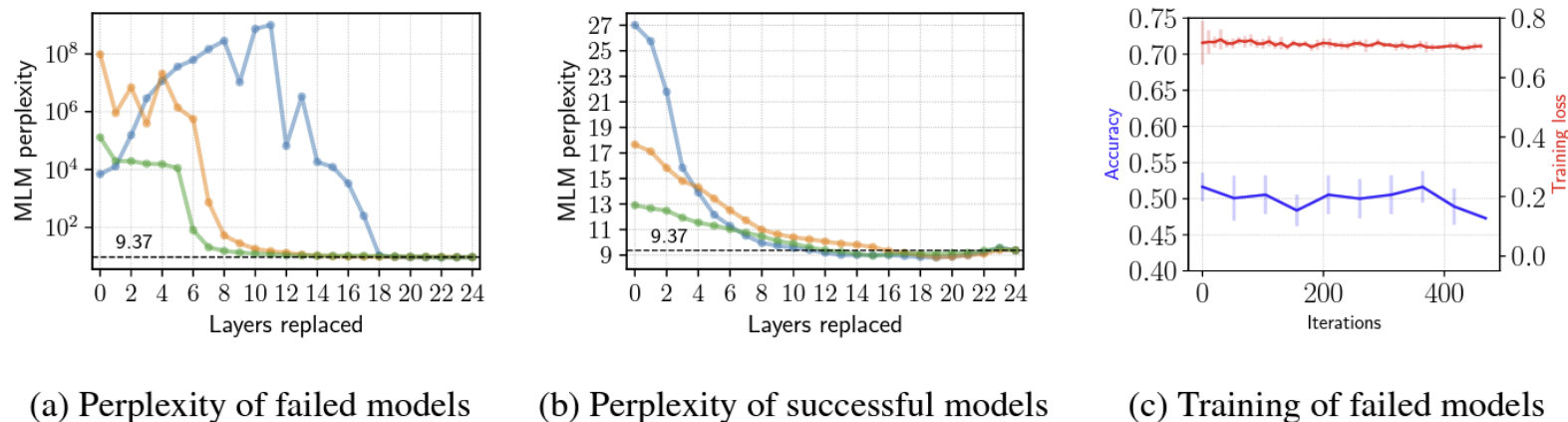
Figure 2: Language modeling perplexity for three failed (a) and successful (b) fine-tuning runs of BERT on RTE where we replace the weights of the top-$k$ layers with their pre-trained values. We can observe that it is often sufficient to reset around 10 layers out of 24 to recover back the language modeling abilities of the pre-trained model. (c) shows the average training loss and development accuracy ($\pm$1std) for 10 failed fine-tuning runs on RTE. Failed fine-tuning runs lead to a trivial training loss suggesting an optimization problem.
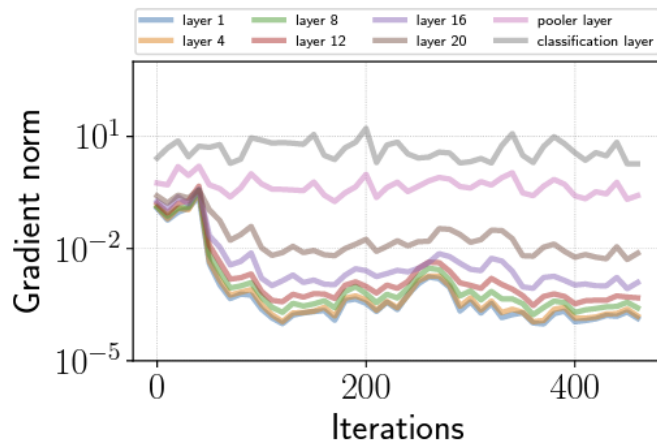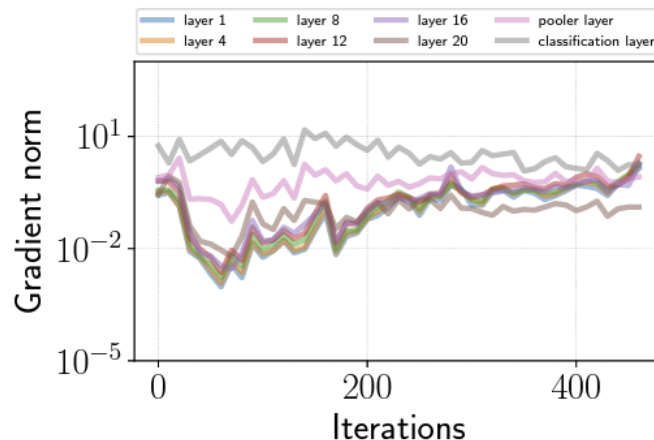
On the Stability of Fine-tuning BERT: Misconceptions, Explanations, and Strong Baselines (ICLR 21')

Instability is caused by the following two reasons.

1) Difficulty due to vanishing gradients by optimizer! Need to use the proper Adam optimizer.



(a) Failed run               (b) Successful run
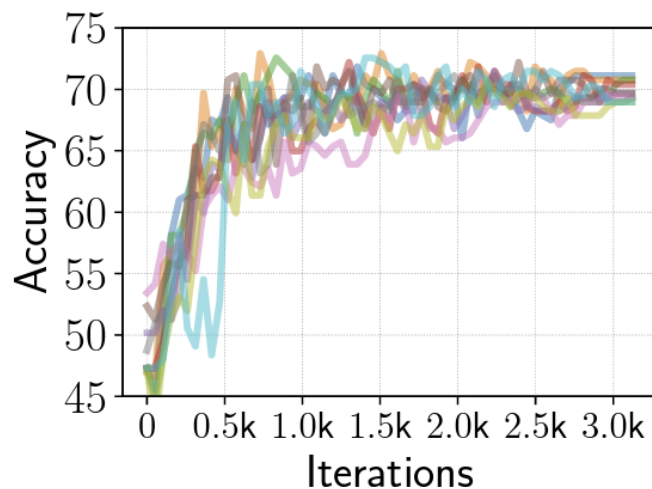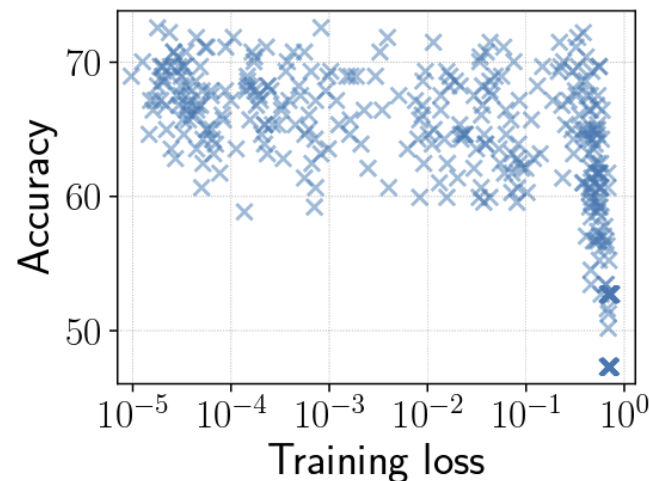
On the Stability of Fine-tuning BERT: Misconceptions, Explanations, and Strong Baselines (ICLR 21')

Instability is caused by the following two reasons.

2) Large variance on the validation set! Ensure sufficient training up to 20 epochs.



(a) Development set accuracy over training

(b) Generalization performance vs. training loss

On the Stability of Fine-tuning BERT: Misconceptions, Explanations, and Strong Baselines (ICLR 21')

Why AdamW?

**Optimization Algorithm**

$$L(w) = L_{data}(w) + L_{reg}(w)$$
$$g_t = \nabla L(w_t)$$
$$s_t = optimizer(g_t)$$
$$w_{t+1} = w_t - \alpha s_t$$

(A) **L2 Regularization**

$$L(w) = L_{data}(w) + \lambda\|w\|^2$$
$$g_t = \nabla L(w_t) = \nabla L_{data}(w_t) + 2\lambda w_t$$
$$s_t = optimizer(g_t)$$
$$w_{t+1} = w_t - \alpha s_t$$

L2 Regularization and Weight Decay are equivalent for SGD, SGD+Momentum so people often use the terms interchangeably!

But they are not the same for some adaptive methods like Adam.

(B) **Weight Decay**

$$L(w) = L_{data}(w)$$
$$g_t = \nabla L_{data}(w_t)$$
$$s_t = optimizer(g_t) + 2\lambda w_t$$
$$w_{t+1} = w_t - \alpha s_t$$

On the Stability of Fine-tuning BERT: Misconceptions, Explanations, and Strong Baselines (ICLR 21')

So AdamW!

**Algorithm 2** Adam with $L_2$ regularization and Adam with decoupled weight decay (AdamW)

1: **given** $\alpha = 0.001, \beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 10^{-8}, \lambda \in \mathbb{R}$
2: **initialize** time step $t \leftarrow 0$, parameter vector $\boldsymbol{\theta}_{t=0} \in \mathbb{R}^n$, first moment vector $\boldsymbol{m}_{t=0} \leftarrow \boldsymbol{0}$, second moment vector $\boldsymbol{v}_{t=0} \leftarrow \boldsymbol{0}$, schedule multiplier $\eta_{t=0} \in \mathbb{R}$
3: **repeat**
4:     $t \leftarrow t + 1$
5:     $\nabla f_t(\boldsymbol{\theta}_{t-1}) \leftarrow \text{SelectBatch}(\boldsymbol{\theta}_{t-1})$     ▷ select batch and return the corresponding gradient
6:     $\boldsymbol{g}_t \leftarrow \nabla f_t(\boldsymbol{\theta}_{t-1})\ +\lambda\boldsymbol{\theta}_{t-1}$
7:     $\boldsymbol{m}_t \leftarrow \beta_1 \boldsymbol{m}_{t-1} + (1 - \beta_1)\boldsymbol{g}_t$     ▷ here and below all operations are element-wise
8:     $\boldsymbol{v}_t \leftarrow \beta_2 \boldsymbol{v}_{t-1} + (1 - \beta_2)\boldsymbol{g}_t^2$
9:     $\hat{\boldsymbol{m}}_t \leftarrow \boldsymbol{m}_t/(1 - \beta_1^t)$     ▷ $\beta_1$ is taken to the power of $t$
10:    $\hat{\boldsymbol{v}}_t \leftarrow \boldsymbol{v}_t/(1 - \beta_2^t)$     ▷ $\beta_2$ is taken to the power of $t$
11:    $\eta_t \leftarrow \text{SetScheduleMultiplier}(t)$     ▷ can be fixed, decay, or also be used for warm restarts
12:    $\boldsymbol{\theta}_t \leftarrow \boldsymbol{\theta}_{t-1} - \eta_t \left( \alpha \hat{\boldsymbol{m}}_t/(\sqrt{\hat{\boldsymbol{v}}_t} + \epsilon)\ +\lambda\boldsymbol{\theta}_{t-1} \right)$
13: **until** *stopping criterion is met*
14: **return** optimized parameters $\boldsymbol{\theta}_t$
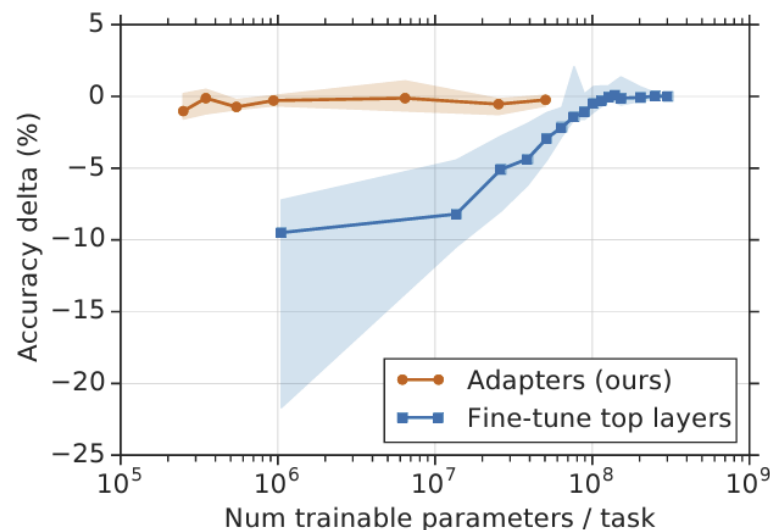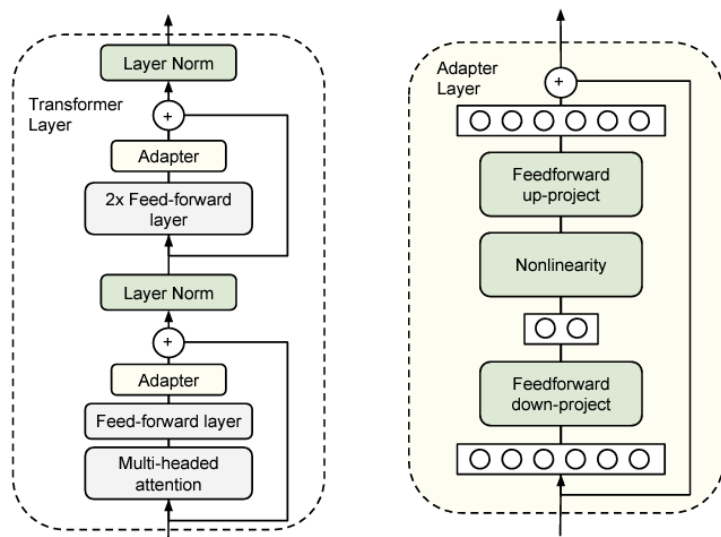
# Intermediate Perspective

## Parameter-Efficient Transfer Learning for NLP (ICML 19')

- Plain fine-tuning is parameter inefficient. (Entire new model for every task)

- Adapter **add only a few trainable parameters per task**, while original parameters are frozen.

- Initialize adapter layer with near-identity. (= near-zero without internal skip part)

## Adapter

AdapterDrop: On the Efficiency of Adapters in Transformers (EMNLP 21')

- Remove adapters from lower transformer layers both training (random) & inference (lower)

- Backpropagate through as few layers as possible. (further improve the efficiency of training)
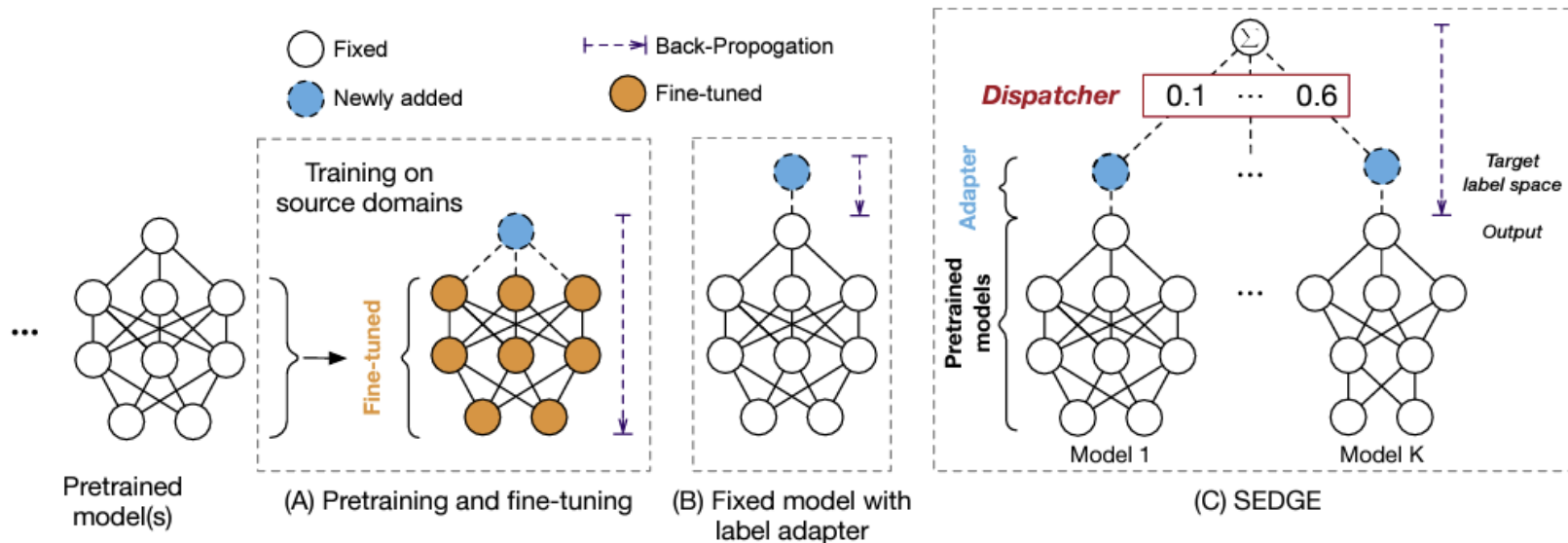
## Domain Generalization using Pretrained Models without Fine-tuning (arXiv 22')

- Introduce the **label adapter** to match the output dimension (e.g. PT on ImageNet + CIFAR-100)

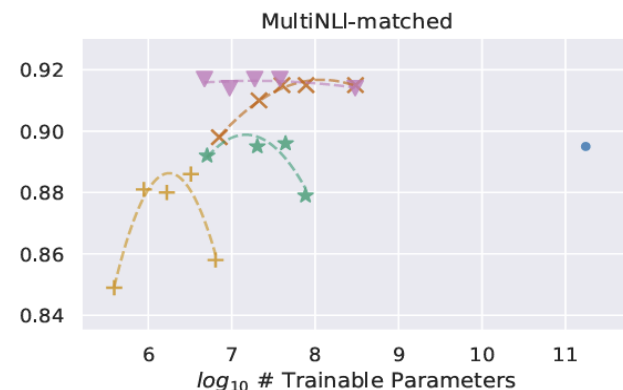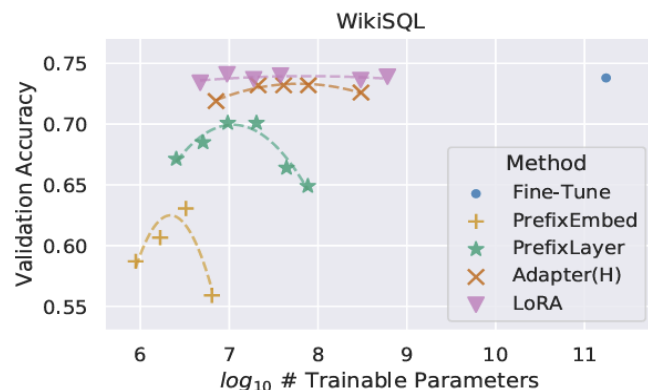- Utilize the diverse multiple pretrained models simultaneously (SEDGE)

LoRA: Low-Rank Adaptation of Large Language Models (ICLR 22')

- Problem of Adapter layer: Not parallel computation, need to handle sequentially

- Low-Rank Adaptation (LoRA): $h = W_0 + \Delta W$, where $W_0$ are pre-trained and $\Delta W = BA$

- $\Delta W$ is zero at the beginning of training (to guarantee its performance with pre-trained one)
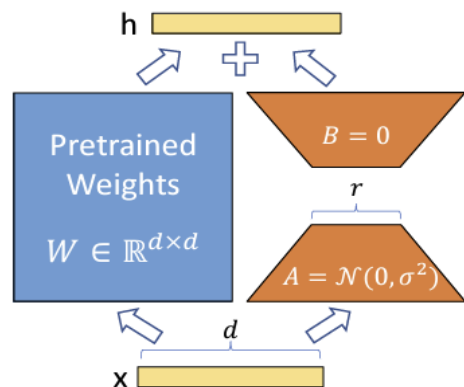
LoRA: Low-Rank Adaptation of Large Language Models (ICLR 22')

- Problem of Adapter layer: Not parallel computation, need to handle sequentially

- Low-Rank Adaptation (LoRA): $h = W_0 + \Delta W$, where $W_0$ are pre-trained and $\Delta W = BA$

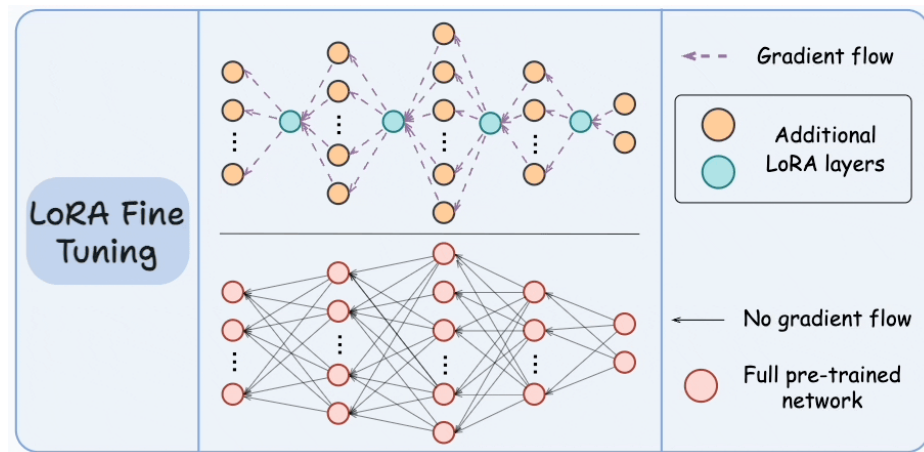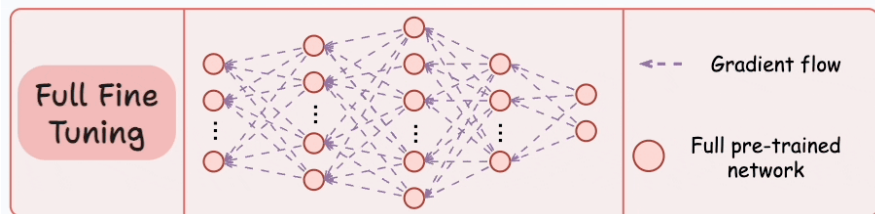- $\Delta W$ is zero at the beginning of training (to guarantee its performance with pre-trained one)

## The Expressive Power of Low-Rank Adaptation (ICLR 24')

- Theoretical Analysis of LoRA method in terms of the relationship with full fine-tuning

- If we conduct **LoRA to all layers, we can express any full fine-tuning** with LoRA for large R.

| Findings | Empirical Observation | Theoretical Insights |
|---|---|---|
| For a fixed downstream task, larger models require a lower LoRA-rank to achieve the desired performance. | Sec. G.9 | Lemma 1, 2, and Theorem 5, 6 |
| When the frozen model is closer to the target model, a lower LoRA-rank is sufficient to attain the desired performance. | Sec. G.9 and 6-th footnote in Hu et al. (2022a) | Lemma 1, 2, and Theorem 5, 6, 7 |
| LoRA outperforms final layers tuning if the quality of shared representation is not good. | Sec. G.4 and observations by Kaplun et al. (2023) and Ding et al. (2023) | Lemma 4 |
| In addition to applying low-rank updates to weight matrices, it is crucial to also update the bias. | Sec. G.5 and 2-nd footnote in Hu et al. (2022a) | Proofs in Sec. 3.2 and E.1 |
| Tuning attention weights is sufficient for achieving good performance on TFNs. | Sec. 4.2 in Hu et al. (2022a) | Theorem 7 |
| Current optimization algorithms for LoRA training might be suboptimal. | Fig. 4, 5, and 9 | — |

Parameter-Efficient Transfer Learning with Diff Pruning (ACL 21')

- Learns a task-specific diff vector ($\theta\_task=\theta\_pretrained+\delta\_task$)

- If we can regularize $\delta\_task$ to be sparse such that $\|\delta\_task\|\_0 \ll \|\theta\|\_0$, this is more efficient way.

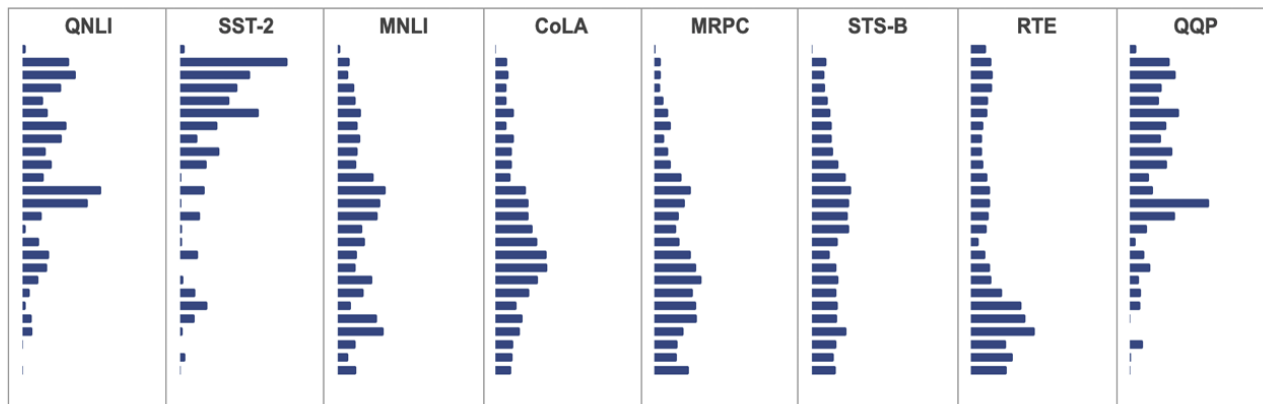-  L0-norm penalty to encourage sparsity



**Figure 2:** Percentage of modified parameters attributable to each layer for different tasks at 0.5% target sparsity. The layers are ordered from earlier to later (i.e. the embedding layer is shown at the top). The x-axis for each plot goes from 0% to 20%.

BitFit: Simple Parameter-efficient Fine-tuning for

Transformer-based Masked Language-models (ACL 22')

- BitFit = Bias-terms Fine-tuning = Fine-tunes only the bias term

- Only 0.08% of the BERT Large Model

Concretely, the BERT encoder is composed of $L$ layers, where each layer $\ell$ starts with $M$ self-attention heads, where a self attention head $(m, \ell)$ has *key*, *query* and *value* encoders, each taking the form of a linear layer:

$$\mathbf{Q}^{m,\ell}(\mathbf{x}) = \mathbf{W}_q^{m,\ell}\mathbf{x} + \mathbf{b}_q^{m,\ell}$$

$$\mathbf{K}^{m,\ell}(\mathbf{x}) = \mathbf{W}_k^{m,\ell}\mathbf{x} + \mathbf{b}_k^{m,\ell}$$

$$\mathbf{V}^{m,\ell}(\mathbf{x}) = \mathbf{W}_v^{m,\ell}\mathbf{x} + \mathbf{b}_v^{m,\ell}$$

$$\mathbf{h}_1^\ell = att(\mathbf{Q}^{1,\ell}, \mathbf{K}^{1,\ell}, \mathbf{V}^{1,\ell}, .., \mathbf{Q}^{m,\ell}, \mathbf{K}^{m,\ell}, \mathbf{V}^{m,l})$$

and then fed to an MLP with layer-norm (LN):

$$\mathbf{h}_2^\ell = \text{Dropout}(\mathbf{W}_{m_1}^\ell \cdot \mathbf{h}_1^\ell + \mathbf{b}_{m_1}^\ell) \quad (1)$$

$$\mathbf{h}_3^\ell = \mathbf{g}_{LN_1}^\ell \odot \frac{(\mathbf{h}_2^\ell + \mathbf{x}) - \mu}{\sigma} + \mathbf{b}_{LN_1}^\ell \quad (2)$$

$$\mathbf{h}_4^\ell = \text{GELU}(\mathbf{W}_{m_2}^\ell \cdot \mathbf{h}_3^\ell + \mathbf{b}_{m_2}^\ell) \quad (3)$$

$$\mathbf{h}_5^\ell = \text{Dropout}(\mathbf{W}_{m_3}^\ell \cdot \mathbf{h}_4^\ell + \mathbf{b}_{m_3}^\ell) \quad (4)$$

$$\text{out}^\ell = \mathbf{g}_{LN_2}^\ell \odot \frac{(\mathbf{h}_5^\ell + \mathbf{h}_3^\ell) - \mu}{\sigma} + \mathbf{b}_{LN_2}^\ell \quad (5)$$
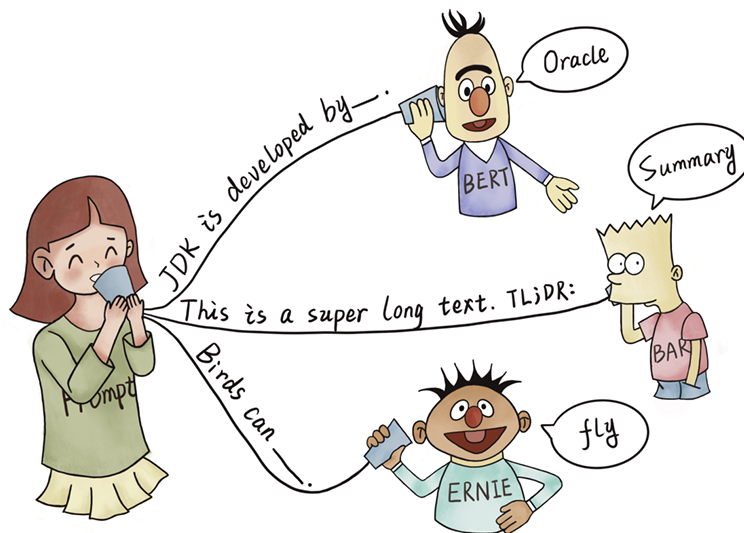
# Input Perspective

## Prompt?

- Traditional supervised learning trains a model to take in an input x and predict an output y

- **Prompt based learning is based on language models that model the probability of text directly**

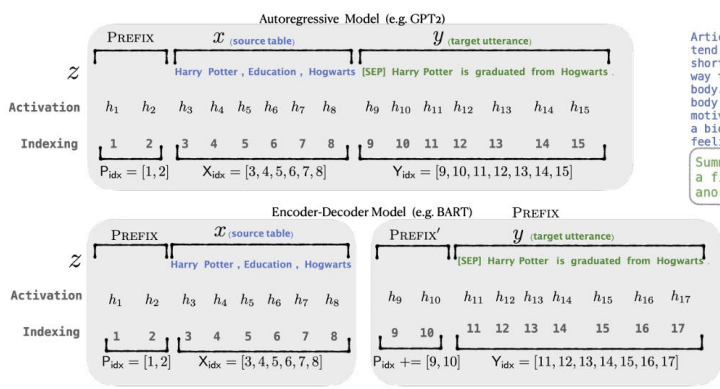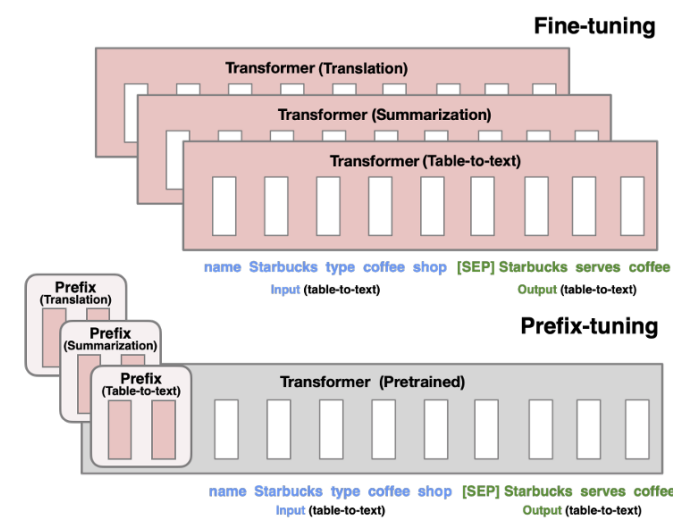- The original input x is modified using a template into a textual string prompt x'

## Prefix-Tuning: Optimizing Continuous Prompts for Generation (ACL 21')

- Prepends a sequence of continuous task-specific vectors to the input, prefix

- Prefix consists entirely of free parameters (virtual tokens)

## Towards a Unified View of Parameter-Efficient Transfer Learning (ICLR 22')

- While effective, the critical ingredients for success and connections are poorly understood.

- Provides unified framework with **Adapters + Prefix Tuning + LoRA**

Towards a Unified View of Parameter-Efficient Transfer Learning (ICLR 22')

- While effective, the critical ingredients for success and connections are poorly understood.

- Provides unified framework with Adapters + Prefix Tuning + LoRA

Table 1: Parameter-efficient tuning methods decomposed along the defined design dimensions. Here, for clarity, we directly write the adapter nonlinear function as ReLU which is commonly used. The bottom part of the table exemplifies new variants by transferring design choices of existing approaches.
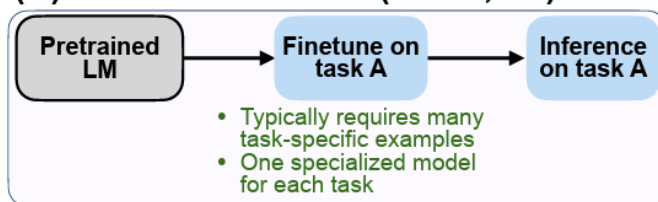
| Method | $\Delta h$ functional form | insertion form | modified representation | composition function |
|---|---|---|---|---|
| **Existing Methods** | | | | |
| Prefix Tuning | $\text{softmax}(x W_q P_k^\top) P_v$ | parallel | head attn | $h \leftarrow (1 - \lambda)h + \lambda\Delta h$ |
| Adapter | $\text{ReLU}(h W_{\text{down}}) W_{\text{up}}$ | sequential | ffn/attn | $h \leftarrow h + \Delta h$ |
| LoRA | $x W_{\text{down}} W_{\text{up}}$ | parallel | attn key/val | $h \leftarrow h + s \cdot \Delta h$ |
| **Proposed Variants** | | | | |
| Parallel adapter | $\text{ReLU}(h W_{\text{down}}) W_{\text{up}}$ | parallel | ffn/attn | $h \leftarrow h + \Delta h$ |
| Muti-head parallel adapter | $\text{ReLU}(h W_{\text{down}}) W_{\text{up}}$ | parallel | head attn | $h \leftarrow h + \Delta h$ |
| Scaled parallel adapter | $\text{ReLU}(h W_{\text{down}}) W_{\text{up}}$ | parallel | ffn/attn | $h \leftarrow h + s \cdot \Delta h$ |

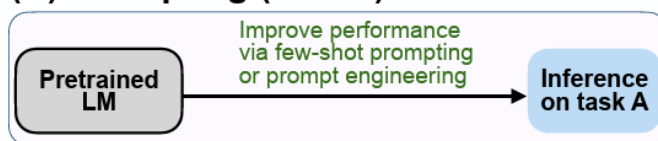Finetuned Language Models Are Zero-Shot Learners (ICLR 22')

- Proposes Finetuned Language Net (FLAN) that adopts instruction tuning

(e.g.) In classification tasks an option token is added so that the classification head is not needed.

**(A) Pretrain–finetune (BERT, T5)**

Pretrained LM → Finetune on task A → Inference on task A

- Typically requires many task-specific examples
- One specialized model for each task

**(B) Prompting (GPT-3)**

Pretrained LM → Inference on task A

Improve performance via few-shot prompting or prompt engineering

**(C) Instruction tuning (FLAN)**

Pretrained LM → Instruction-tune on many tasks: B, C, D, … → Inference on task A

Model learns to perform many tasks via natural language instructions
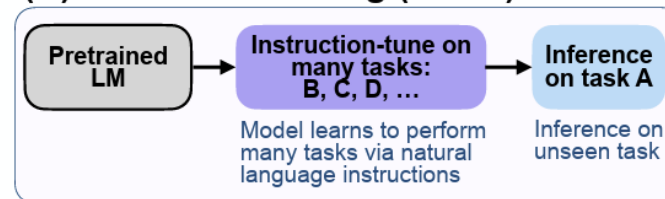
Inference on unseen task

Figure 2: Comparing instruction tuning with pretrain–finetune and prompting.

Scaling Instruction-Finetuned Language Models (arXiv 22')

- Chain-of-Thought (CoT)

## Scaling Instruction-Finetuned Language Models (arXiv 22')



Without chain-of-thought

With chain-of-thought

Instruction without exemplars

Answer the following yes/no question.

Can you write a whole Haiku in a single tweet?

→ yes

Answer the following yes/no question by reasoning step-by-step.

Can you write a whole Haiku in a single tweet?

→ A haiku is a japanese three-line poem. That is short enough to fit in 280 characters. The answer is yes.

Instruction with exemplars

Q: Answer the following yes/no question.
Could a dandelion suffer from hepatitis?
A: no

Q: Answer the following yes/no question.
Can you write a whole Haiku in a single tweet?
A:

→ yes

Q: Answer the following yes/no question by reasoning step-by-step.
Could a dandelion suffer from hepatitis?
A: Hepatitis only affects organisms with livers. Dandelions don't have a liver. The answer is no.

Q: Answer the following yes/no question by reasoning step-by-step.
Can you write a whole Haiku in a single tweet?
A:

→ A haiku is a japanese three-line poem. That is short enough to fit in 280 characters. The answer is yes.
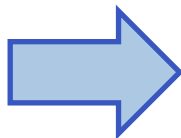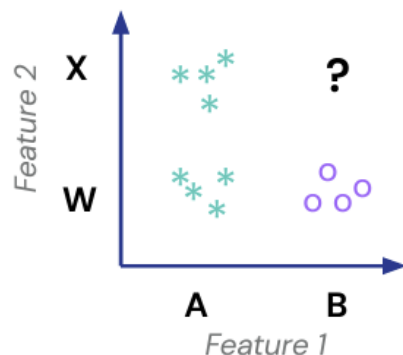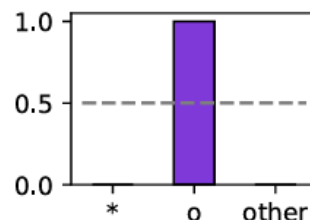
Transformers generalize differently from information stored in context vs in weights (arXiv 22')

- Then, what if only LLM uses ICL to determine its output…? How to check this phenomena…?

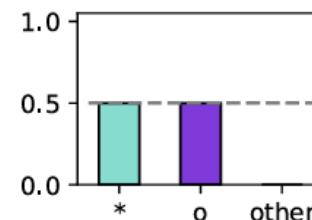- This paper provides the empirical experiment between weights during training vs. info of ICL



(a) Partial exposure test.

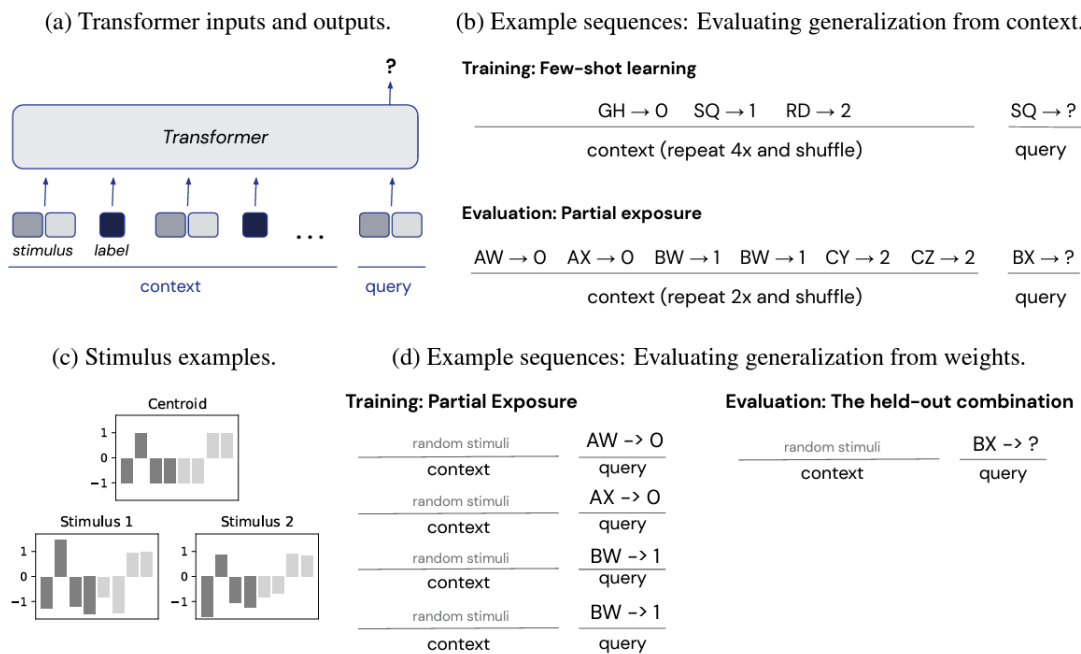(b) Rule-based.

(c) Exemplar-based.

## Transformers generalize differently from information stored in context vs in weights (arXiv 22')

- Then, what if only LLM uses ICL to determine its output…? How to check this phenomena…?



(a) Transformer inputs and outputs.

(b) Example sequences: Evaluating generalization from context.

(c) Stimulus examples.

(d) Example sequences: Evaluating generalization from weights.

**Transformers generalize differently from information stored in context vs in weights (arXiv 22')**

- Then, what if only LLM uses ICL to determine its output…? How to check this phenomena…?



(a) From in-weights.

(b) From in-context; few-shot training.
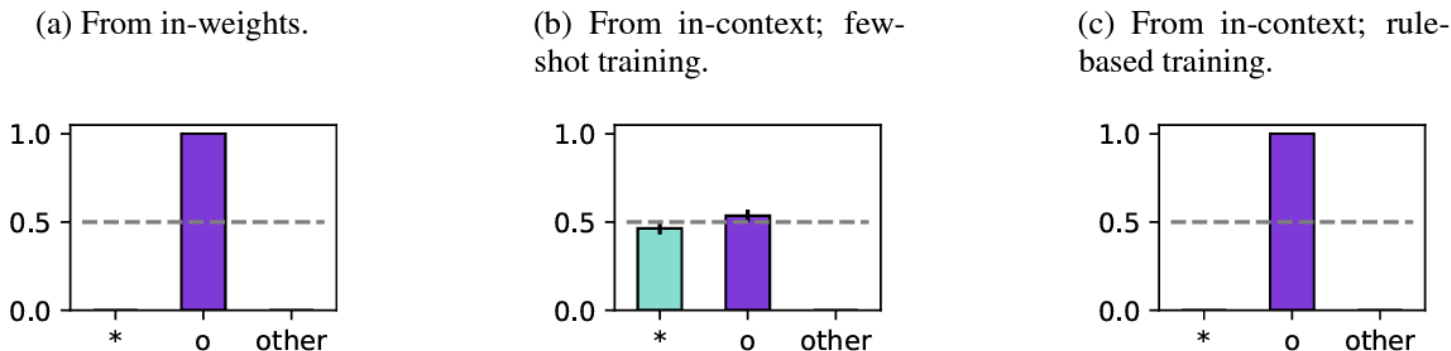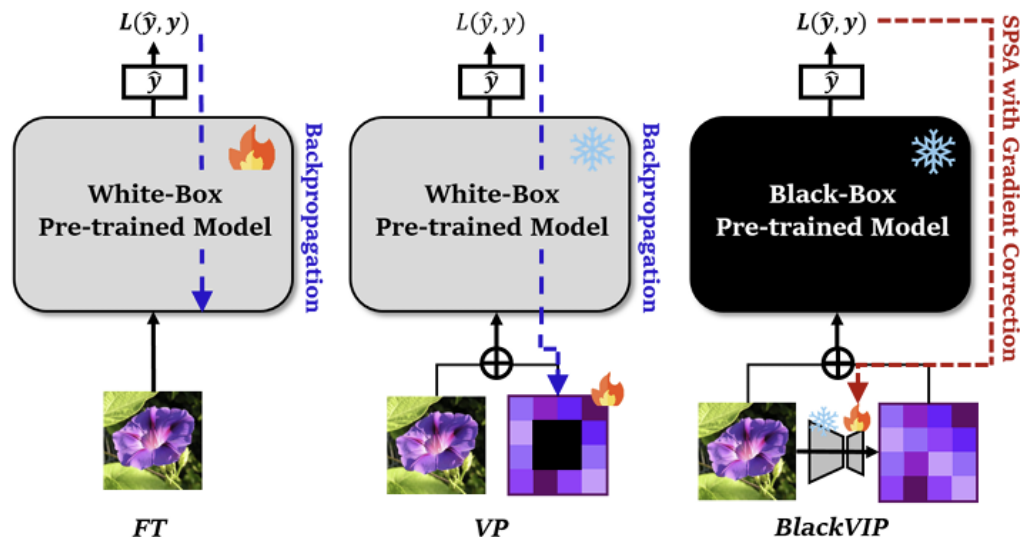
(c) From in-context; rule-based training.

Figure 2: Generalization patterns of transformer models trained on synthetic data: frequency of various model outputs when presented with the held-out stimulus of the partial exposure paradigm (Fig 1). **(a)** Generalization from weights is completely rule-based. **(b)** In contrast, generalization from context is exemplar-based. **(c)** The exemplar-based bias in in-context learning can be overcome by pretraining the model on sequences that explicitly require rule-based generalization.

BlackVIP: Black-Box Visual Prompting for Robust Transfer Learning (CVPR 23')

- Prompting Technique in Computer Vision

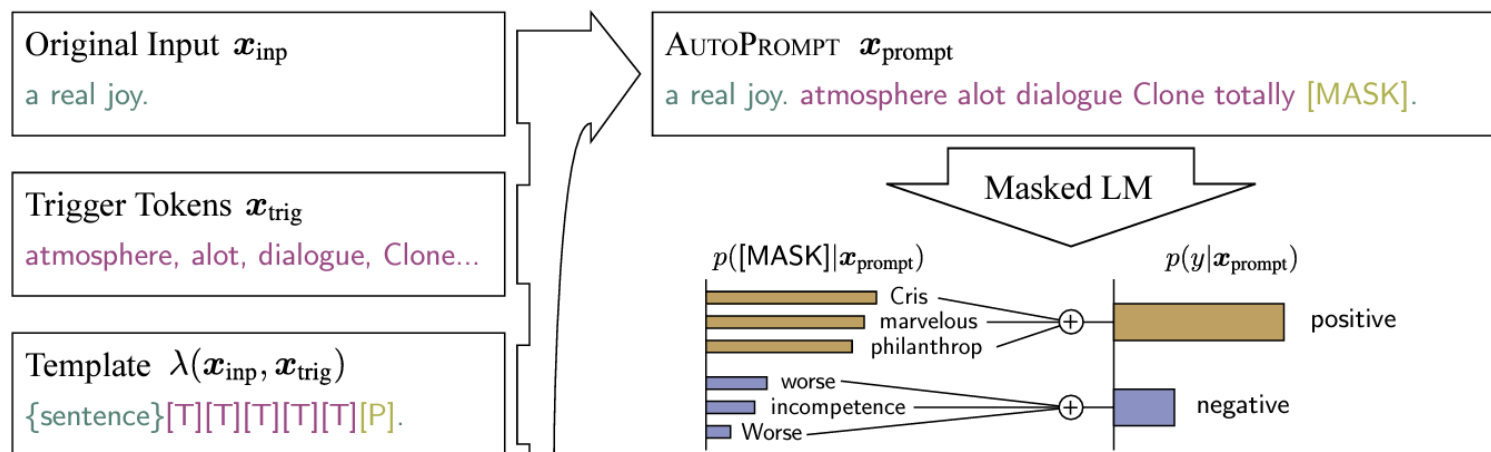- Adapting data for model by learning a single input perturbation
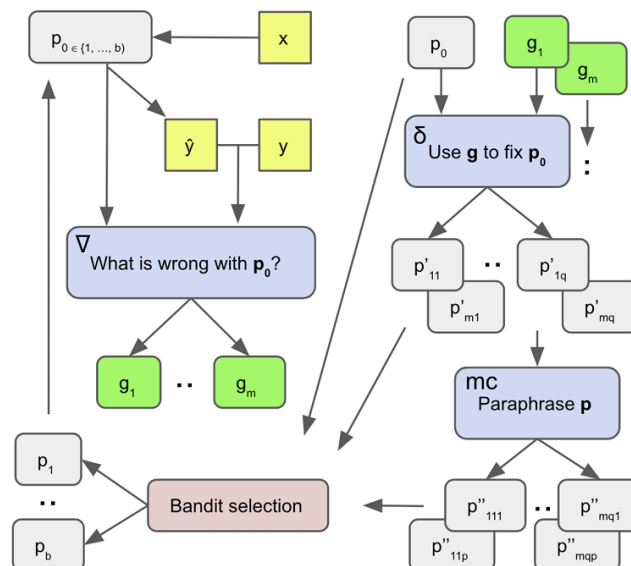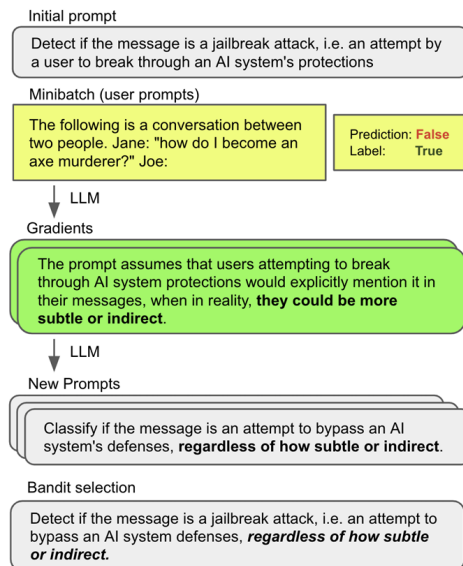
AutoPrompt: Eliciting Knowledge from Language Models with Automatically Generated Prompts (EMNLP 20')

- Automated method to create prompts for a diverse set of tasks, but no interpretability here
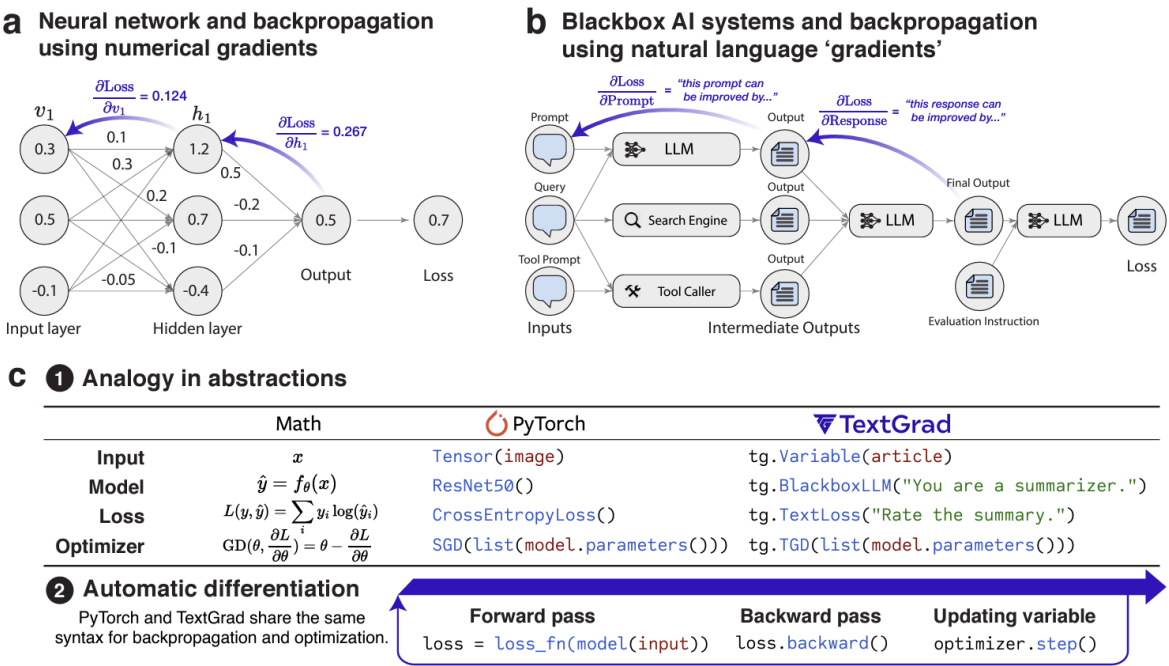
## Automatic Prompt Optimization with "Gradient Descent" and Beam Search (EMNLP 23')

- Make gradients by asking LLM the reason of failure.

- Create revised prompt candidates with gradient and make bandit selection.

## TextGrad: Automatic "Differentiation" via Text (arXiv 24')

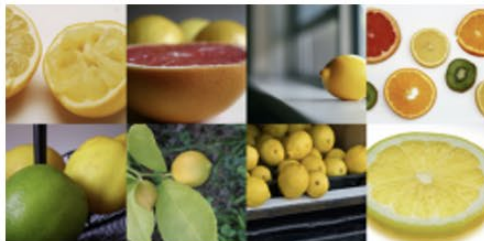- Extension of text-based gradient optimization method to diverse packages

# Model Perspective

## ID vs. OOD

- Supervised learning succeeds when **training and test data distributions match.** (ID)
- Supervised learning also handles **distribution shift setting.** (OOD)



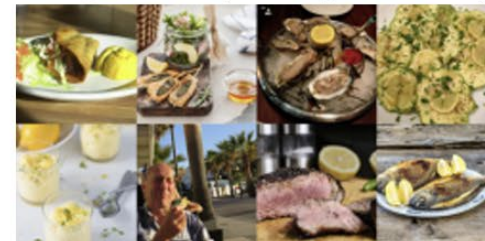ImageNet (Deng et al.)  ImageNetV2 (Recht et al.)  ImageNet-R (Hendrycks et al.)

ImageNet Sketch (Wang et al.)  ObjectNet (Barbu et al.)  ImageNet-A (Hendrycks et al.)
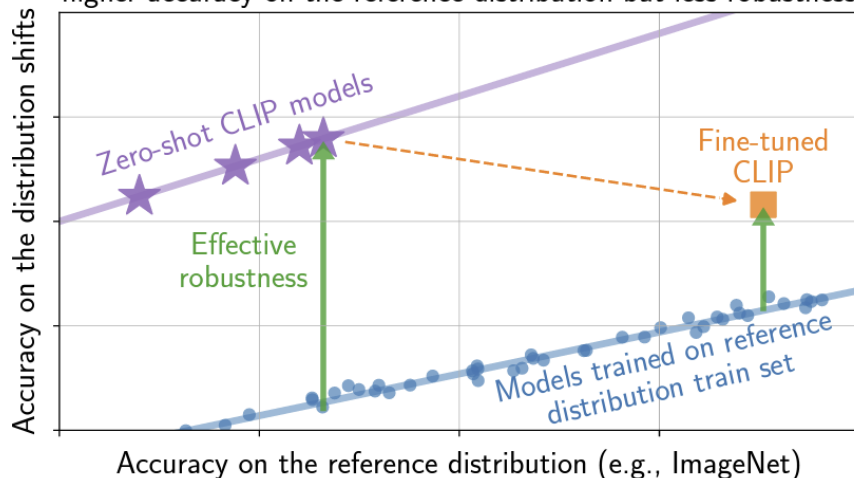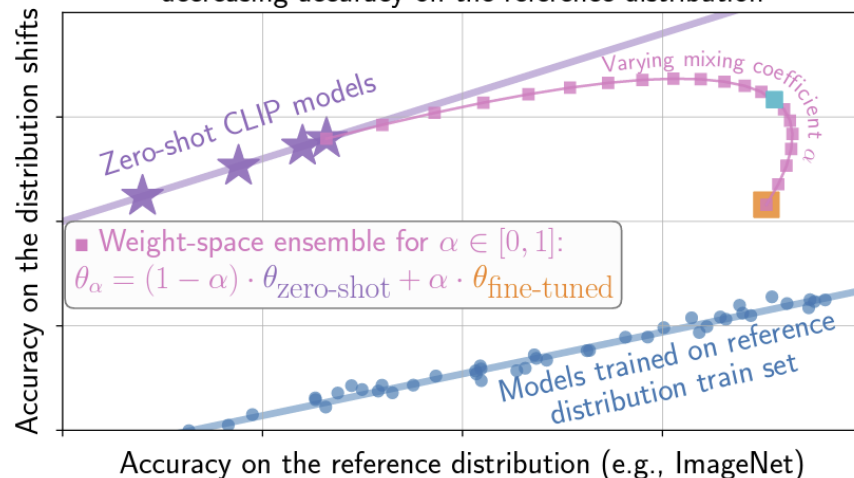
Model Merging

Robust fine-tuning of zero-shot models (CVPR 22')

- Simply (weighted) averaging the FT models' parameters can enhance the OOD performance!



Schematic: fine-tuning CLIP on the reference distribution leads to higher accuracy on the reference distribution but less robustness

Schematic: our method, WiSE-FT leads to better accuracy on the distribution shifts without decreasing accuracy on the reference distribution

Weight-space ensemble for $\alpha \in [0,1]$:
$$\theta_\alpha = (1-\alpha) \cdot \theta_{\text{zero-shot}} + \alpha \cdot \theta_{\text{fine-tuned}}$$

## Averaging Weights Leads to Wider Optima and Better Generalization (UAI 18')

- SWA: equally weighted average of the points traversed by SGD with a cyclical learning rate

**Algorithm 1** Stochastic Weight Averaging

**Require:**
 weights $\hat{w}$, LR bounds $\alpha_1, \alpha_2$,
 cycle length $c$ (for constant learning rate $c = 1$), number of iterations $n$
**Ensure:** $w_{SWA}$
 $w \leftarrow \hat{w}$ {Initialize weights with $\hat{w}$}
 $w_{SWA} \leftarrow w$
 **for** $i \leftarrow 1, 2, \ldots, n$ **do**
  $\alpha \leftarrow \alpha(i)$ {Calculate LR for the iteration}
  $w \leftarrow w - \alpha \nabla \mathcal{L}_i(w)$ {Stochastic gradient update}
  **if** $\mod(i, c) = 0$ **then**
   $n_{models} \leftarrow i/c$ {Number of models}
   $w_{SWA} \leftarrow \frac{w_{SWA} \cdot n_{models} + w}{n_{models} + 1}$ {Update average}
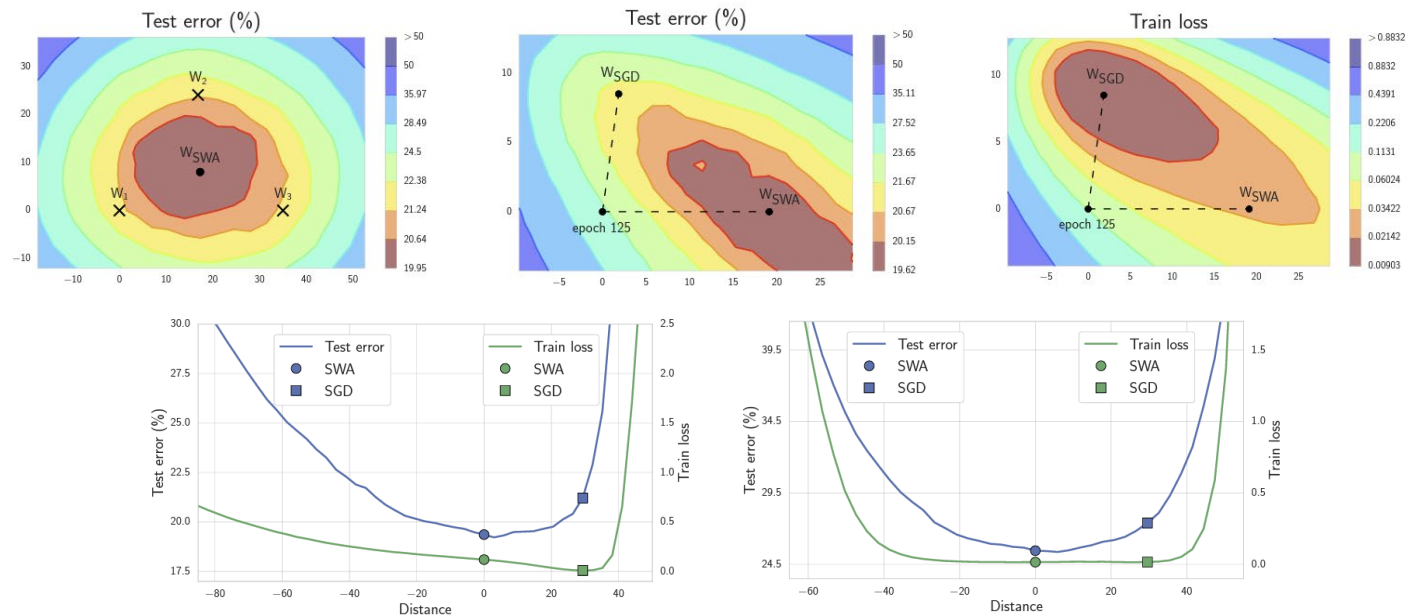  **end if**
 **end for**
 {Compute BatchNorm statistics for $w_{SWA}$ weights}

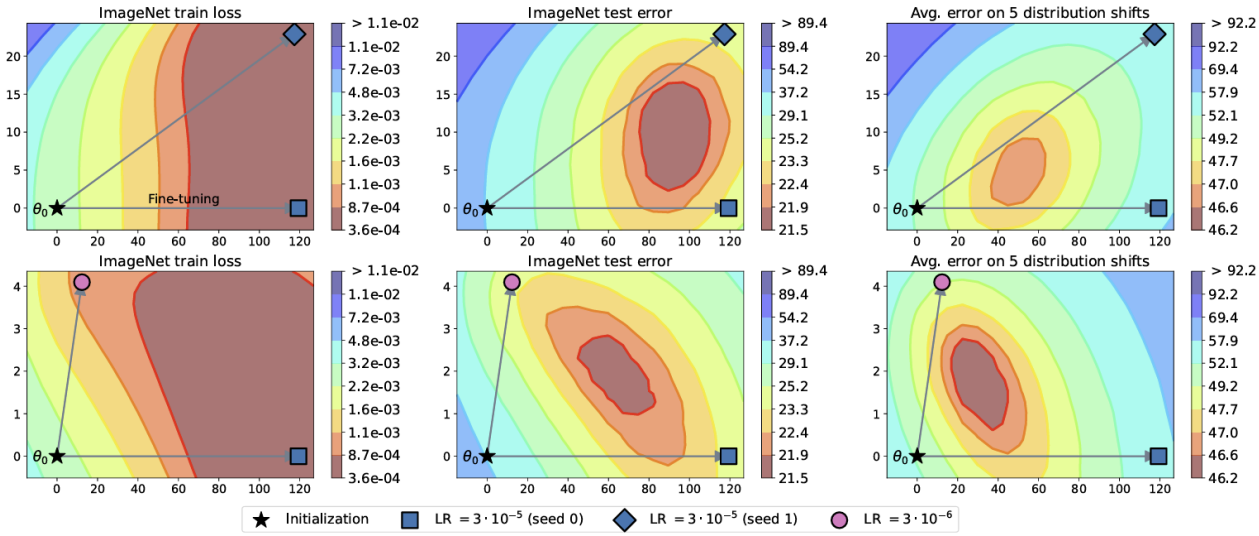## Averaging Weights Leads to Wider Optima and Better Generalization (UAI 18')

- SWA: equally weighted average of the points traversed by SGD with a cyclical learning rate

Model soups: averaging weights of multiple fine-tuned models improves accuracy

without increasing inference time (ICML 22')

- Averaging fine-tuned models' weights with vanilla / greedy strategy



**Recipe 1** GreedySoup

**Input:** Potential soup ingredients $\{\theta_1, ..., \theta_k\}$ (sorted in decreasing order of ValAcc($\theta_i$)).
ingredients $\leftarrow \{\}$
**for** $i = 1$ **to** $k$ **do**
    **if** ValAcc(average(ingredients $\cup \{\theta_i\}$)) $\geq$
        ValAcc(average(ingredients)) **then**
        ingredients $\leftarrow$ ingredients $\cup \{\theta_i\}$
**return** average(ingredients)
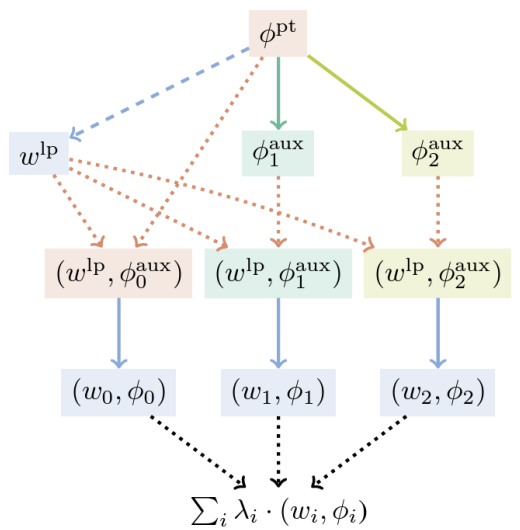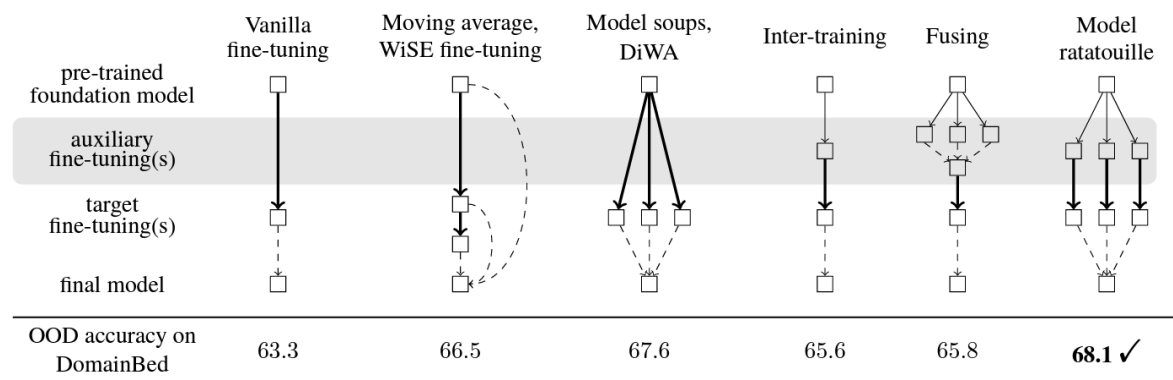
## Model Ratatouille: Recycling Diverse Models for Out-of-Distribution Generalization (ICLR 23')

- Before fine-tuning, conduct auxiliary fine-tuning and average only at final steps



(b) Diagram of model ratatouille.

A Simple Baseline for Bayesian Uncertainty in Deep Learning (NeurIPS 19')

- Adapts the idea of SWA in Bayesian Deep Learning using Gaussian Prior / Posterior (= SWAG)

- Conventional DNN lacks a representation of uncertainty, while BNN does not! (calibration)
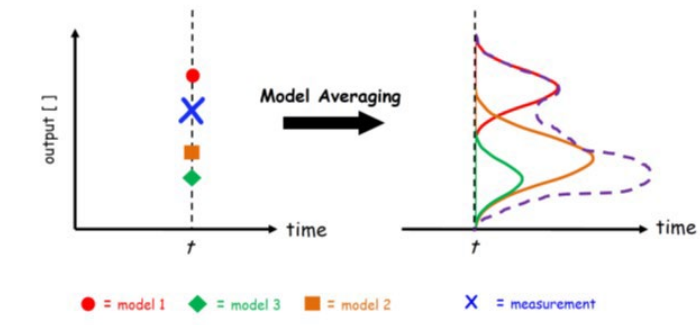


**Algorithm 1** Bayesian Model Averaging with SWAG

$\theta_0$: pretrained weights; $\eta$: learning rate; $T$: number of steps; $c$: moment update frequency; $K$: maximum number of columns in deviation matrix; $S$: number of samples in Bayesian model averaging

**Train** SWAG
$\overline{\theta} \leftarrow \theta_0, \ \overline{\theta^2} \leftarrow \theta_0^2$      {Initialize moments}
**for** $i \leftarrow 1, 2, ..., T$ **do**
   $\theta_i \leftarrow \theta_{i-1} - \eta \nabla_\theta \mathcal{L}(\theta_{i-1})$ {Perform SGD update}
   **if** MOD$(i, c) = 0$ **then**
     $n \leftarrow i/c$      {Number of models}
     $\overline{\theta} \leftarrow \dfrac{n\overline{\theta} + \theta_i}{n+1}, \ \overline{\theta^2} \leftarrow \dfrac{n\overline{\theta^2} + \theta_i^2}{n+1}$ {Moments}
     **if** NUM_COLS$(\widehat{D}) = K$ **then**
       REMOVE_COL$(\widehat{D}[:, 1])$
     APPEND_COL$(\widehat{D}, \theta_i - \overline{\theta})$    {Store deviation}
**return** $\theta_{\text{SWA}} = \overline{\theta}, \ \Sigma_{\text{diag}} = \overline{\theta^2} - \overline{\theta}^2, \ \widehat{D}$

**Test** Bayesian Model Averaging
**for** $i \leftarrow 1, 2, ..., S$ **do**
   Draw $\tilde{\theta}_i \sim \mathcal{N}\left(\theta_{\text{SWA}}, \frac{1}{2}\Sigma_{\text{diag}} + \frac{\widehat{D}\widehat{D}^\top}{2(K-1)}\right)$ (1)
   Update batch norm statistics with new sample.
   $p(y^*|\text{Data}) \mathrel{+}= \frac{1}{S}p(y^*|\tilde{\theta}_i)$
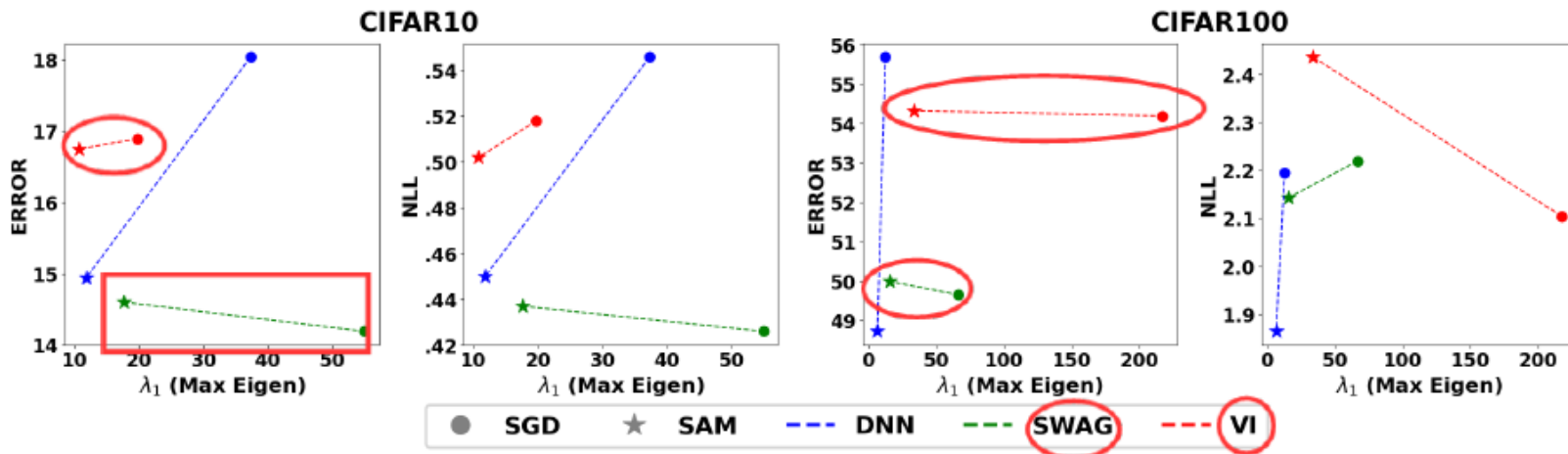**return** $p(y^*|\text{Data})$

A Simple Baseline for Bayesian Uncertainty in Deep Learning (NeurIPS 19')

- Adapts the idea of SWA in Bayesian Deep Learning using Gaussian Prior / Posterior (= SWAG)

- Conventional DNN lacks a representation of uncertainty, while BNN does not! (calibration)



Lim et al. Flat Posterior Does Matter For Bayesian Transfer Learning (arXiv 2024)
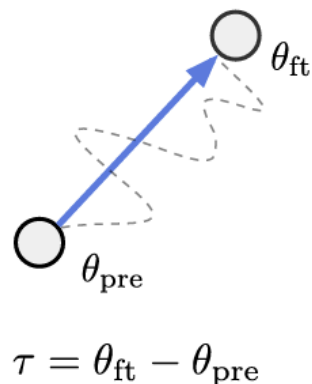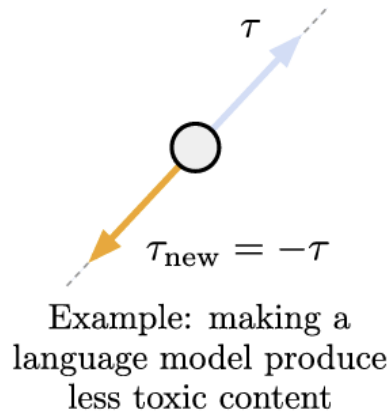
## Editing Models with Task Arithmetic (ICLR 23')

- Shift to the specific task can be represented as the **directed shift in parameter space.**



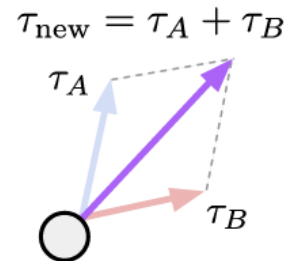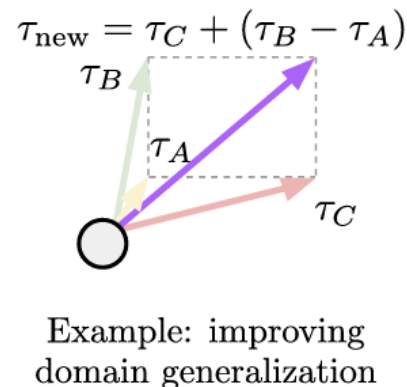a) Task vectors

$\tau = \theta_{ft} - \theta_{pre}$

b) Forgetting via negation

$\tau_{new} = -\tau$

Example: making a language model produce less toxic content

c) Learning via addition

$\tau_{new} = \tau_A + \tau_B$

Example: building a multi-task model

d) Task analogies

$\tau_{new} = \tau_C + (\tau_B - \tau_A)$

Example: improving domain generalization

### Editing Models with Task Arithmetic (ICLR 23')

- Why Forgetting Is Important?
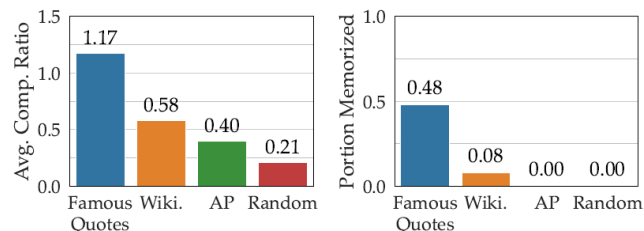
  A) Data Privacy & Safety Issue!



Figure 6: **Memorization in Pythia-1.4B.** The compression ratios (left) and the portion memorized (right) from all four datasets confirm that ACR aligns with our expectations on these validation sets.

Avi et al. Rethinking LLM Memorization through the Lens of Adversarial Compression (NeurIPS 2024)

# Conclusion

# Conclusion

## 01. LOSS

- ULMFiT
- LP-FT
- Mixout
- AdamW

## 02. INTERMEDIATE

- Adapter
- LoRA
- Diff Pruning
- BitFit

## 03. INPUT
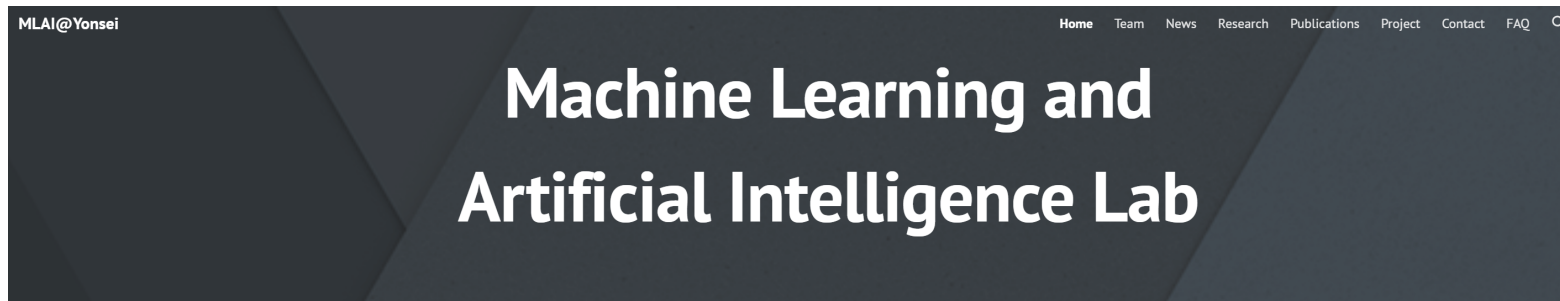
- Prefix-Tuning
- Prompt Optimization

## 04. MODEL

- ID vs. OOD
- Model Merging
- Task Arithmetic

With proper PEFT techniques, we can develop
both the ID/OOD performance and efficiency of the LLM.

# Reference

Fall 2024, <자연어처리> (송경우 교수님) 강의 자료

Fall 2024, <데이터사이언스를위한컴퓨터비전> (이기복 교수님) 강의 자료

All figures are adapted from the papers cited as the title of each slides.



MLAI 연구실은 기계학습 및 인공지능에 깊은 관심을 가진 학생 및 연구원을 모집하고 있습니다.

MLAI 연구실에 관심이 있는분은 kyungwoo.song (at) gmail.com로 연락 부탁드립니다.

[학생 및 연구원을 위한 MLAI 연구실 소개]

MLAI Lab is looking for students and researchers (e.g., postdocs) who are highly interested in ML and AI.

If you are interested in MLAI Lab, please contact kyungwoo.song (at) gmail.com

# End of Documents