

# Calibrating Factual Knowledge in Pretrained Language Models

Yufang Liu

*East China Normal University*



# Editing Neural Nets: Why?

- Neural networks contain many beliefs, but...

**Input: Who is the prime minister of the UK?**

**T5:** *Theresa May*

**BART:** *Theresa May*

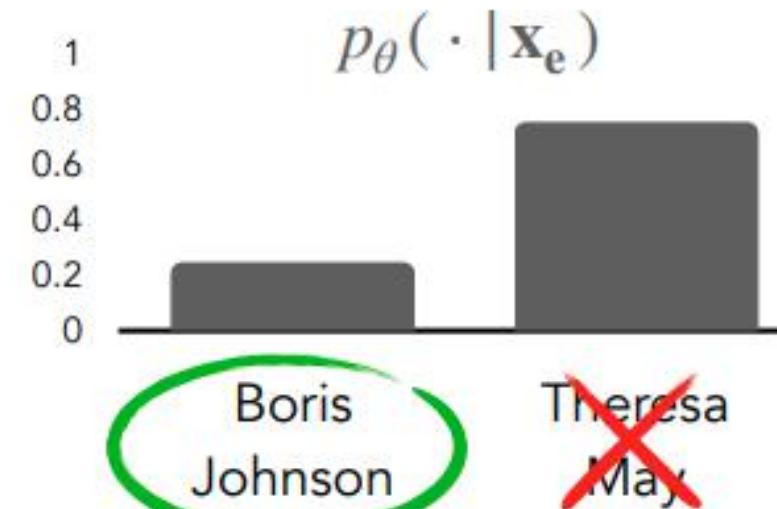
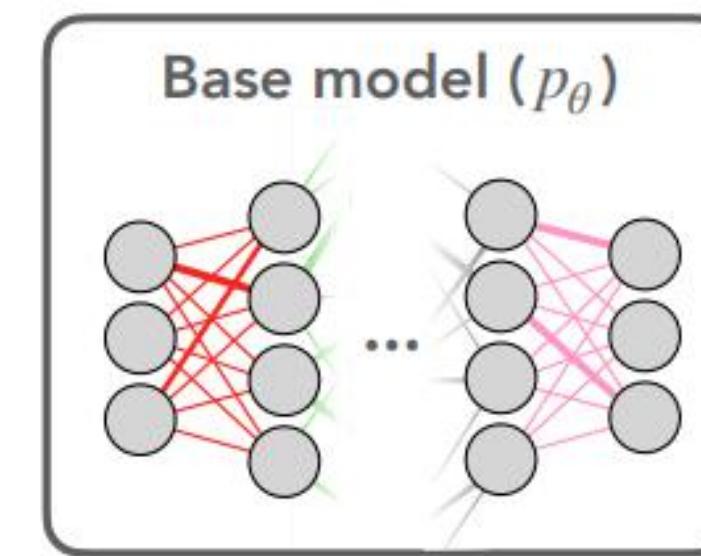
**GPT-3:** *Theresa May*

} Not anymore!

...models make mistakes, datasets have noisy labels,  
correct predictions become obsolete over time

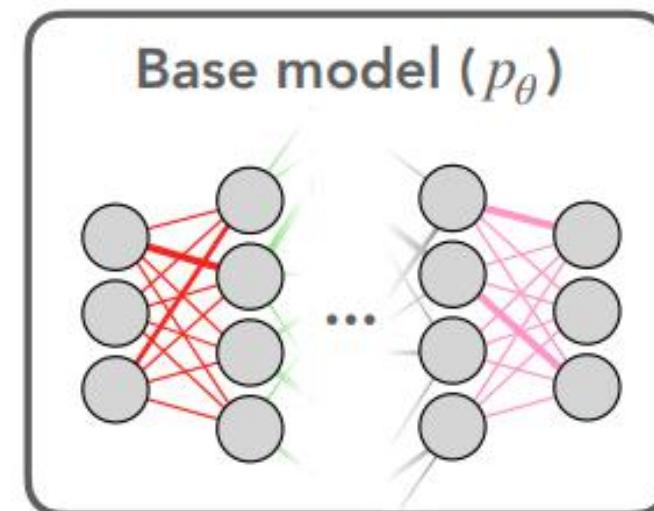
# Editing Neural Nets: Why?

$\mathbf{x}_e$  = “Who is the prime minister of the UK?”



# Editing Neural Nets: Why?

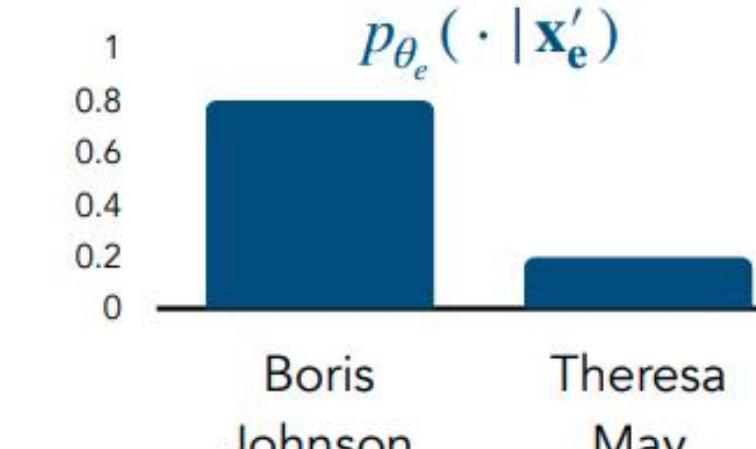
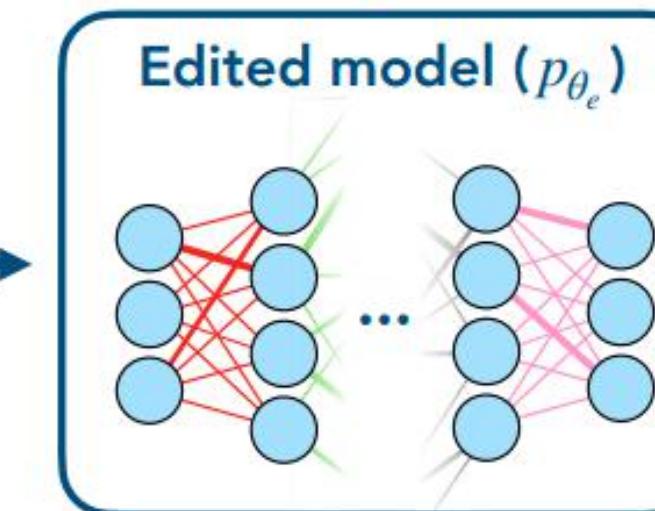
$\mathbf{x}_e$  = "Who is the prime minister of the UK?"



$y_e$  = "Boris Johnson"

Model Editor

$\mathbf{x}'_e$  = "Who is the UK PM?"



semantically equivalent questions

unrelated questions

# Editing Factual Knowledge in Language Models\*

---

**Nicola De Cao**



**Wilker Aziz**



**Ivan Titov**



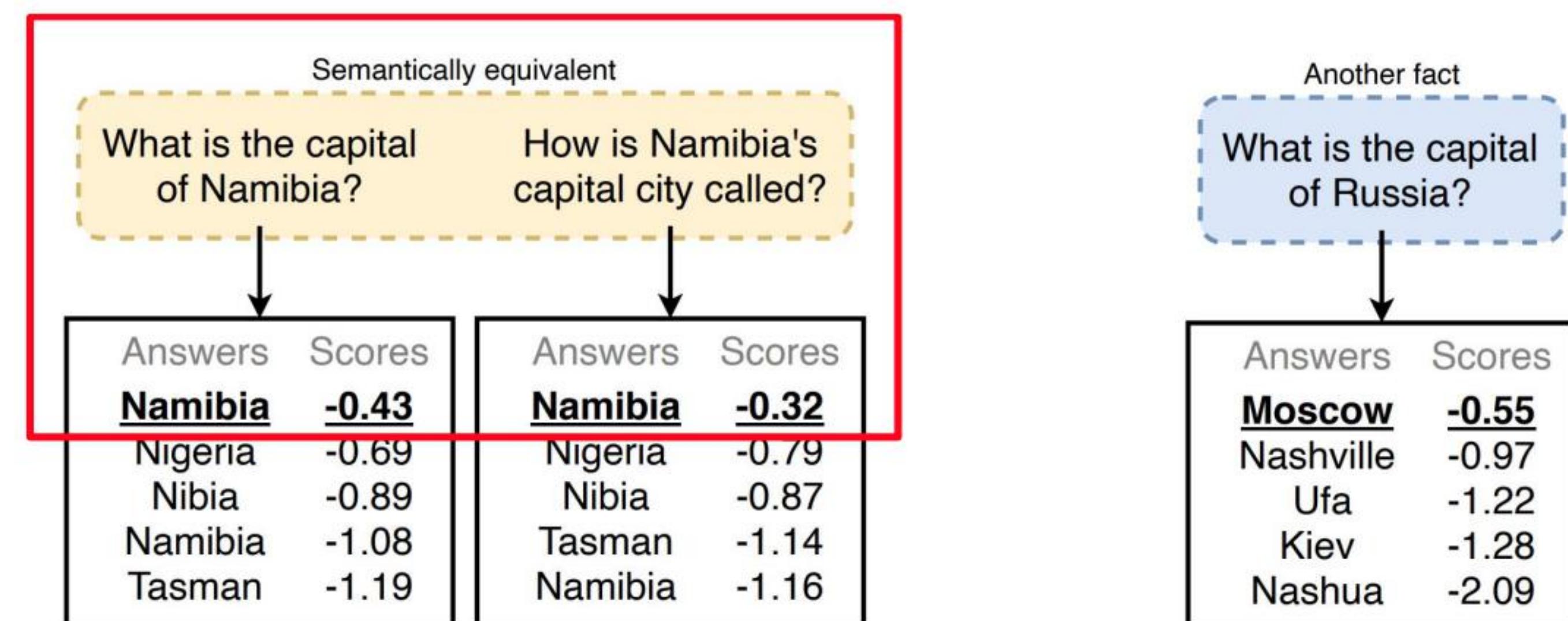
*EMNLP 2021*

\*aka **KnowledgeEditor**: <https://github.com/nicola-decao/KnowledgeEditor>

# Background

Goal 

**Wrong**   
**is “Windhoek”**



Research question: **“Can we change the factual knowledge of the model without altering all the other memorized facts?”**

# Example

- BART finetuned for closed-book question answering

Semantically equivalent		Another fact	
What is the capital of Namibia?		How is Namibia's capital city called?	
Answers	Scores	Answers	Scores
<u>Namibia</u>	<b>-0.43</b>	<u>Namibia</u>	<b>-0.32</b>
Nigeria	-0.69	Nigeria	-0.79
Nibia	-0.89	Nibia	-0.87
Namibia	-1.08	Tasman	-1.14
Tasman	-1.19	Namibia	-1.16

(a) Model predictions before the update.

Fact to change		Fact that also changes		Another fact	
What is the capital of Namibia?		How is Namibia's capital city called?		What is the capital of Russia?	
Answers	Scores	Answers	Scores	Answers	Scores
<u>Windhoek</u>	<b>-0.06</b>	<u>Windhoek</u>	<b>-0.07</b>	<u>Moscow</u>	<b>-0.56</b>
Tasman	-1.42	Tasman	-1.50	Ufa	-1.03
Windygates	-1.52	Windygates	-1.51	Nashville	-1.04
Tasmania	-1.59	Windhoof	-1.53	Kiev	-1.43
Windhoof	-1.66	Tasmania	-1.53	Nashua	-2.21

(b) Model predictions with edited parameters.

# Definition

- model  $x \mapsto f(x; \theta)$
- revised dataset  $\langle x, y, a \rangle \in \mathcal{D}$
- $\mathcal{P}^x$  semantically equivalent inputs with  $x$ 
  - automatically-generated paraphrases using back-translation
- $\mathcal{O}^x$  unrelated inputs with  $x$
- knowledge editor  $g$

# Methods

- Optimization

$$\min_{\phi} \sum_{\hat{x} \in \mathcal{P}^x} \mathcal{L}(\theta'; \hat{x}, a)$$

$$\text{s.t. } \mathcal{C}(\theta, \theta', f; \mathcal{O}^x) \leq m \quad \rightarrow$$

$$\min_{\phi} \max_{\alpha} f(x, \theta) + \alpha \cdot (\mathcal{C}(y, \theta) - m)$$

$$\theta' = \theta + g(x, y, a; \phi)$$

- Constraint

$$\mathcal{C}_{KL}(\theta, \theta', f; \mathcal{O}^x) = \sum_{x' \in \mathcal{O}^x} \sum_{c \in \mathcal{Y}} p_{Y|X}(c|x', \theta) \log \frac{p_{Y|X}(c|x', \theta)}{p_{Y|X}(c|x', \theta')}$$

$$\mathcal{C}_{L_p}(\theta, \theta', f; \mathcal{O}^x) = (\sum_i |\theta_i - \theta'_i|^p)^{1/p}$$

# Methods

- architecture for knowledge editor

$$g(x, y, a; \phi) \longrightarrow \Delta\theta$$

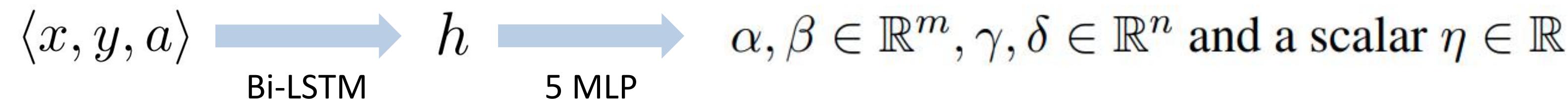
- naive hyper-network might be over-parameterized X
- for weight matrix  $W^{n \times m} \in \theta$

$$\Delta W = \sigma(\eta) \cdot \left( \hat{\alpha} \odot \nabla_W \mathcal{L}(W; x, a) + \hat{\beta} \right)$$

$\sigma$  Sigmoid

with  $\hat{\alpha} = \hat{\sigma}(\alpha) \gamma^\top$  and  $\hat{\beta} = \hat{\sigma}(\beta) \delta^\top$

$\hat{\sigma}$  Softmax



# Evaluation

- success rate  $\mathcal{D}$
- retain accuracy  $\mathcal{O}^x$
- equivalence accuracy  $\mathcal{P}^x$
- performance deterioration  $1 - \frac{\text{accuracy of } f(\cdot; \theta')}{\text{accuracy of } f(\cdot; \theta)}$

# Experiments

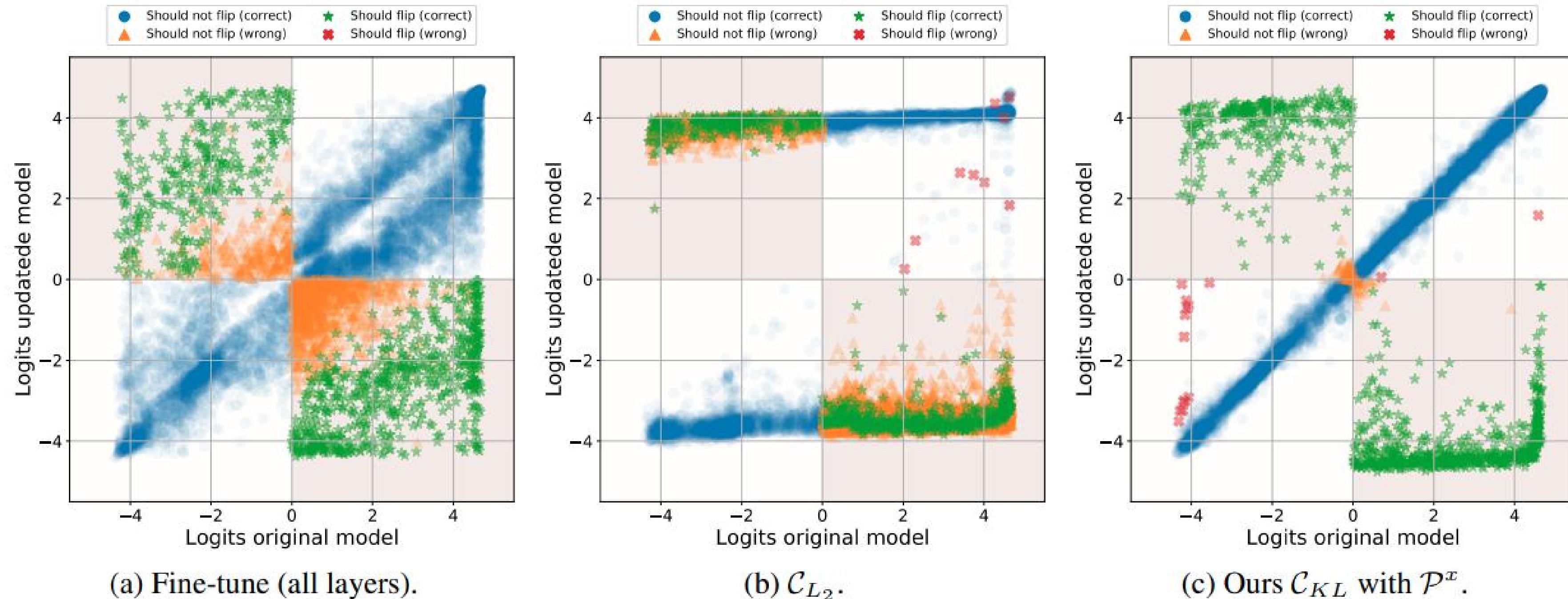
- Binary FEVER dataset-> closed-book fact-checking
  - Bert model
  - flipping the predicted labels
- Zero-Shot Relation Extraction (zsRE) dataset ->closed-book QA
  - Bart model
  - pick all hypotheses enumerated via beam search except the top-1

# Experiments

Method	Fact-Checking				Question Answering			
	Success rate ↑	Retain acc ↑	Equiv. acc ↑	Perform. det ↓	Success rate ↑	Retain acc ↑	Equiv. acc ↑*	Perform. det ↓
Fine-tune (1st layer)	100.0	99.44	42.24	0.00	98.68	91.43	89.86 / 93.59	0.41
Fine-tune (all layers)	100.0	86.95	95.58	2.25	100.0	67.55	97.77 / 98.84	4.50
Zhu et al. (1st layer)	100.0	99.44	40.30	0.00	81.44	92.86	72.63 / 78.21	0.32
Zhu et al. (all layers)	100.0	94.07	83.30	0.10	80.65	95.56	76.41 / 79.38	0.35
Ours $\mathcal{C}_{L_2}$	99.10	45.10	99.01	35.29	99.10	46.66	97.16 / 99.24	9.22
KNOWLEDGEEDITOR	98.80	98.14	82.69	0.10	94.65	98.73	86.50 / 92.06	0.11
+ loop <sup>†</sup>	100.0	97.78	81.57	0.59	99.23	97.79	89.51 / 96.81	0.50
+ $\mathcal{P}^x$ <sup>‡</sup>	98.50	98.55	95.25	0.24	94.12	98.56	91.20 / 94.53	0.17
+ $\mathcal{P}^x$ + loop <sup>‡</sup>	100.0	98.46	94.65	0.47	99.55	97.68	93.46 / 97.10	0.95

generated/human  
annotated

# Experiments



Distribution of logits of the original model and updated model on FEVER.

# Experiments

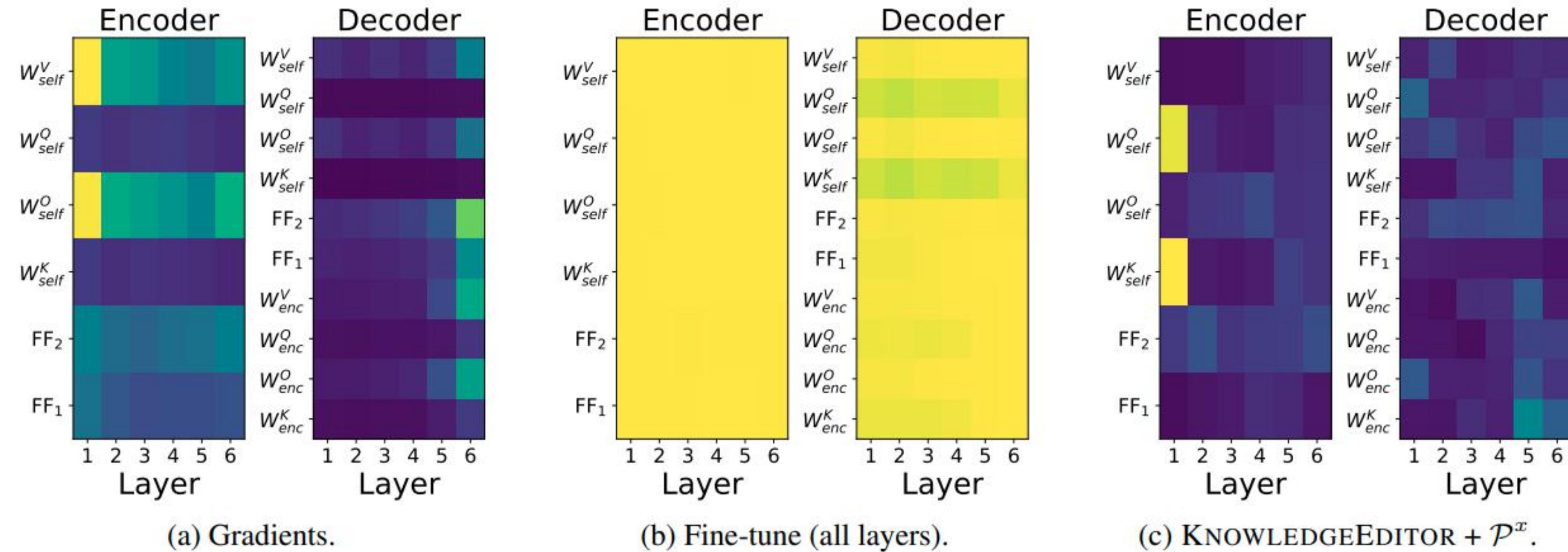


Figure 4: Average normalized magnitude of updates on weight matrices across layers for the QA experiment. Fine-tuning updates all layers uniformly while our updates are more sparse.

the hyper-network may be regarded as a probe providing insights about the mechanism used by the model to encode the knowledge

# **Plug-and-Play Adaptation for Continuously-updated QA**

**Kyungjae Lee<sup>5</sup>**

**Hwaran Lee<sup>2</sup>**

<sup>1</sup>Seoul National University

<sup>4</sup>University of Richmond

**Wookje Han<sup>1</sup>**

**Joonsuk Park<sup>2,4</sup>**

<sup>2</sup>NAVER AI Lab

<sup>5</sup>LG AI Research

**Seung-won Hwang<sup>1\*</sup>**

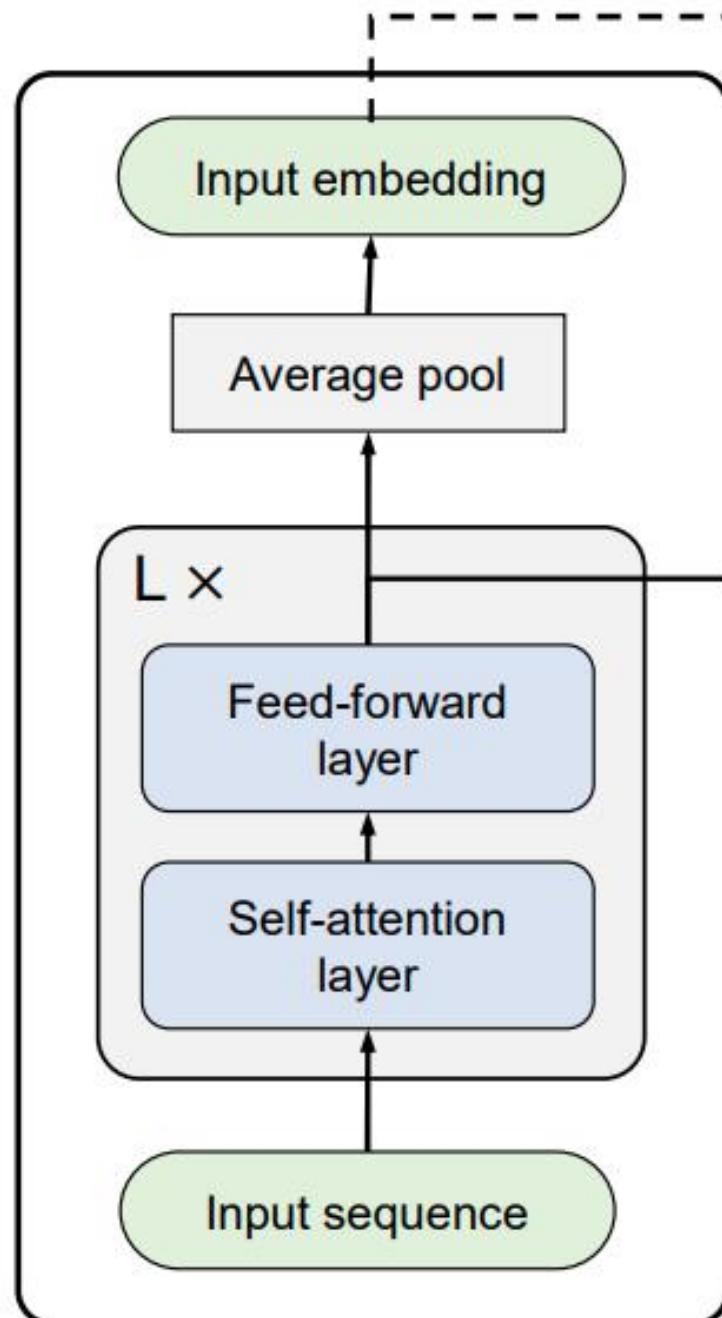
**Sang-Woo Lee<sup>2,3</sup>**

<sup>3</sup>NAVER CLOVA

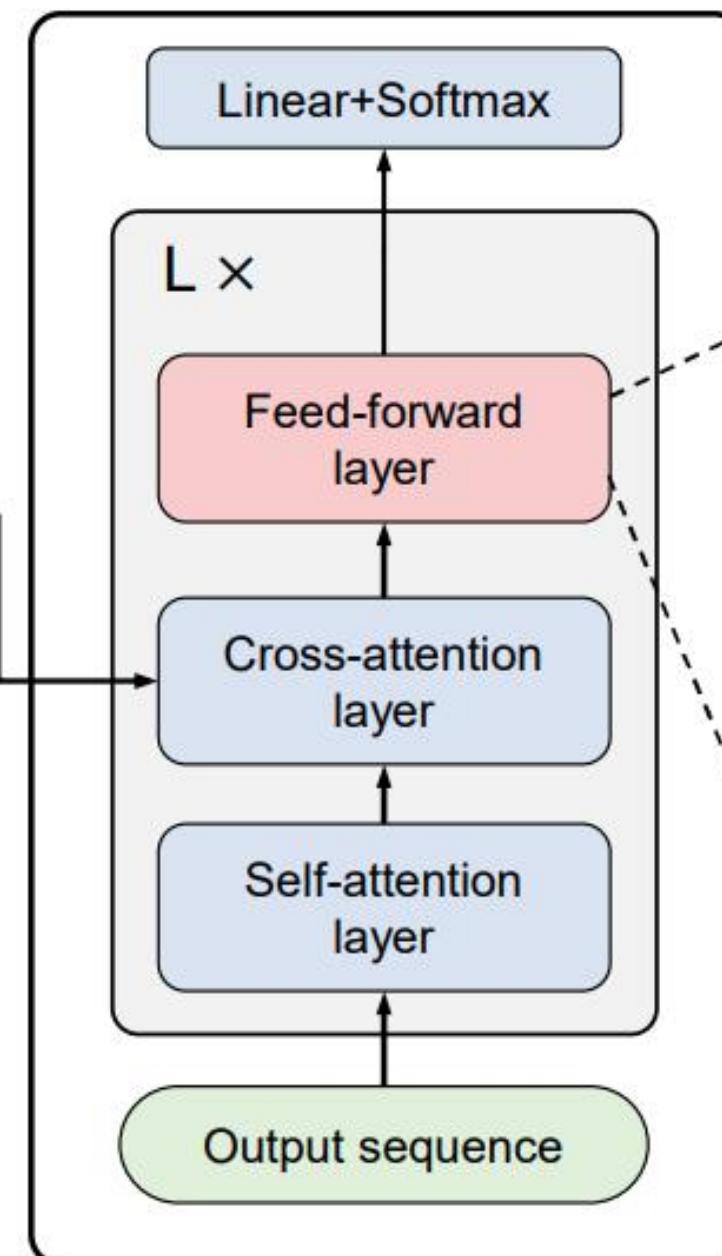
*Findings of ACL 2022*

# Methods

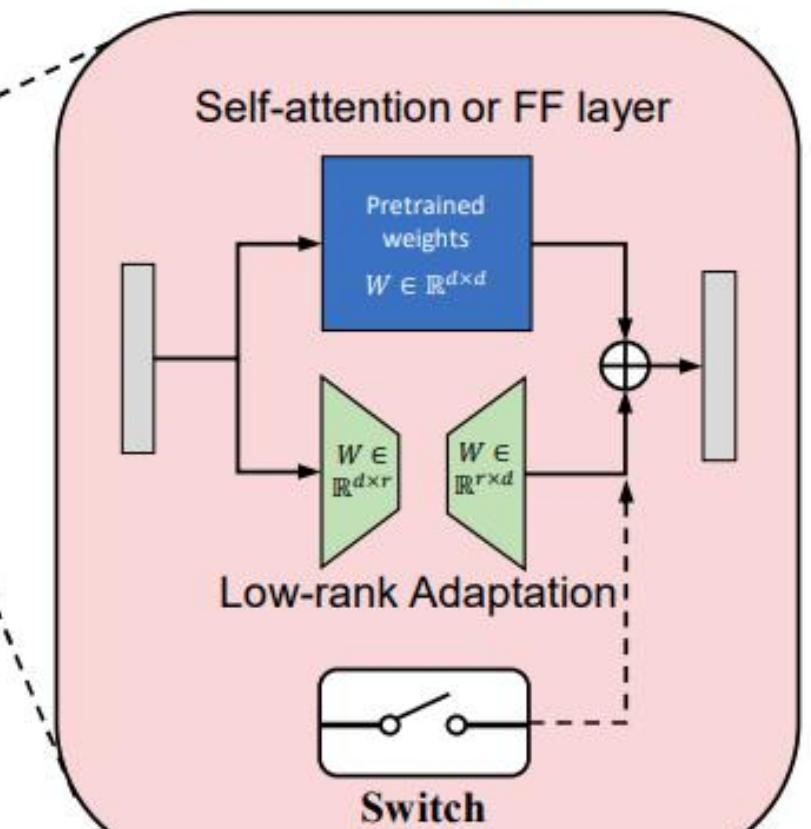
## Encoder



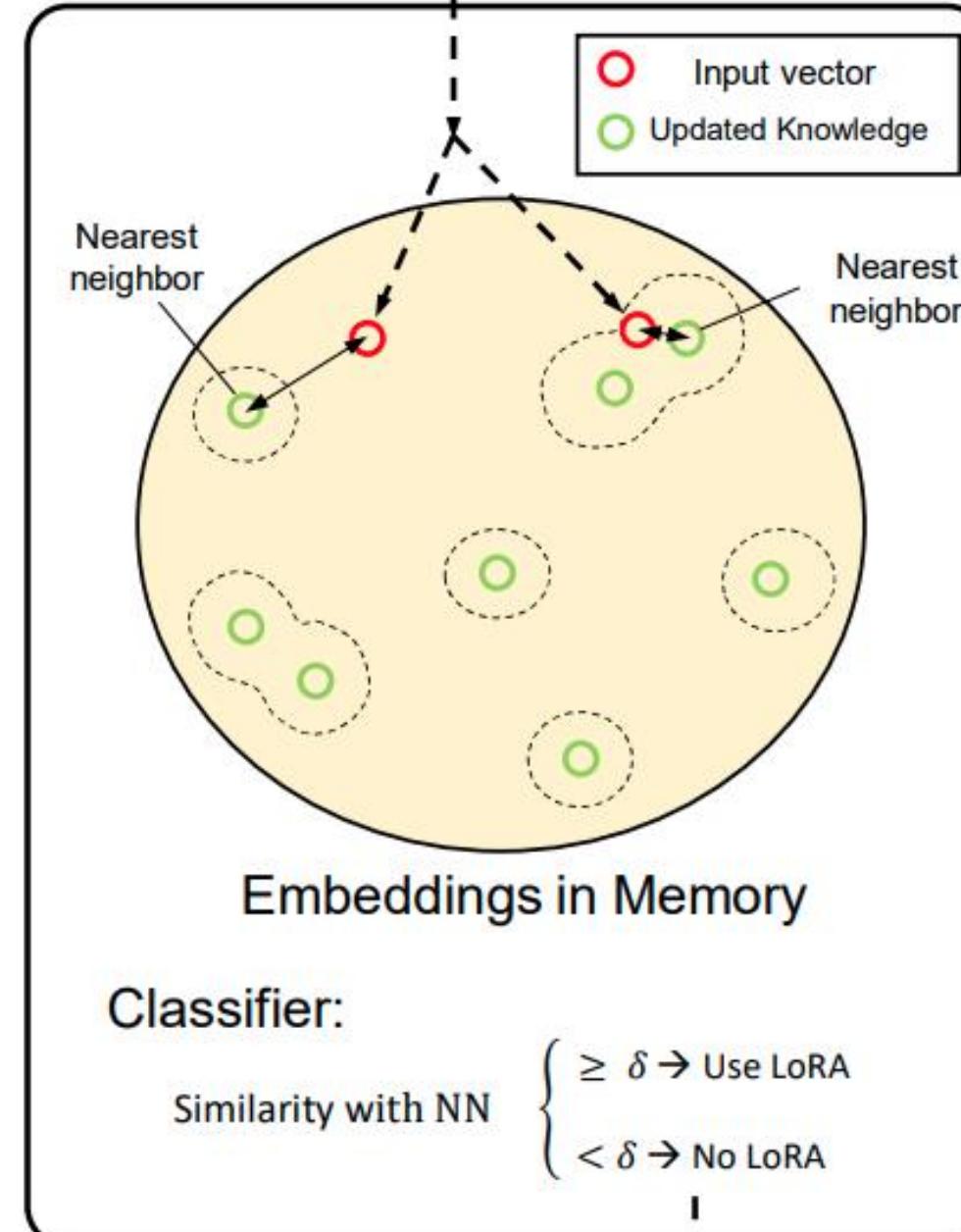
## Decoder



## Decoder with Selective Parameter-expansion



## Adaptation Switch



# Methods

- selector: key( $m_i$ )-value( $g$ ) lookup

$$h = f(x) + \sigma(q) \cdot g(x)$$

- memory embeddings  $\mathcal{M} \in \mathbf{R}^{N \times d}$

$$s_q = \max_i(\text{sim}(m_i, q)), \quad m_i \in \mathcal{M}$$



average the hidden states

$$\sigma(q) = \begin{cases} 1 & \text{if } s_q \geq \delta, \\ 0 & \text{if } s_q < \delta. \end{cases}$$

# Methods

- multiple knowledge updates

$$h = f(x) + \sum_{k=1}^M \sigma_k(q) \cdot g_k(x)$$

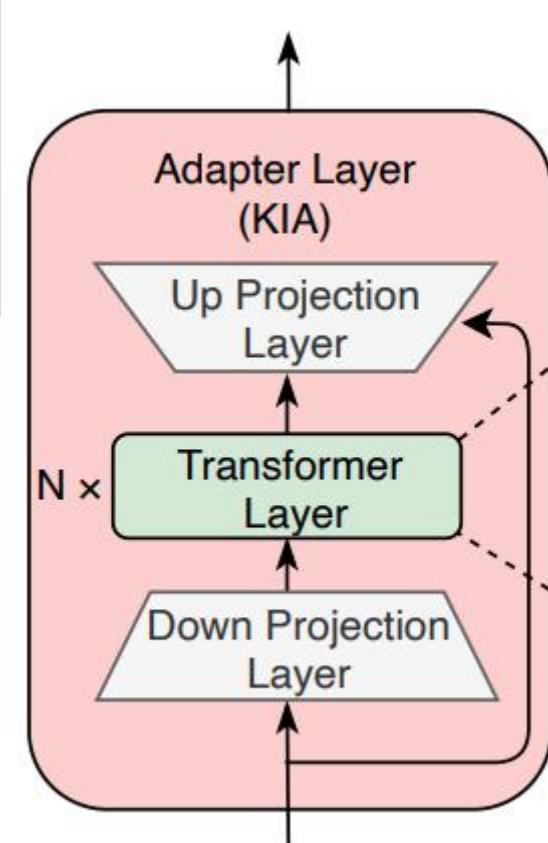
$$m^* = \operatorname{argmax}_m (\text{sim}(m, q)), \quad m \in \mathcal{M}_{1:M}$$

$$\sigma_k(q) = \begin{cases} 1 & \text{if } s_q \geq \delta \text{ and } 1 \leq k \leq j, \\ 0 & \text{if } s_q < \delta. \end{cases} \quad m^* \in \mathcal{M}_j$$

# Experiment

- Zero-shot Relation Extraction (zsRE)
  - split into  $\mathcal{K}_s$  and  $\mathcal{K}_t$
- Natural Questions (NQ) + SituatedQA
  - back-translation to generate paraphrases

# Experiment



	Method	# of Prams (train/total)	zsRE Question Answering					NQ (with SituatedQA)				
			$\mathcal{K}_s$	$\mathcal{P}_s$	$\mathcal{K}_t$	$\mathcal{P}_t$	F/U Ratio	$\mathcal{K}_s$	$\mathcal{P}_s$	$\mathcal{K}_t$	$\mathcal{P}_t$	F/U Ratio
	Model $\theta^{old}$	-	95.6	95.2	25.7	28.5	-	96.6	94.9	35.3	33.7	-
$L_p$ norm	B-I: Fine-tuning	737M / 737M	76.7	70.6	92.6	85.9	0.284	92.9	82.5	94.9	<b>92.9</b>	0.435
	B-II: RecAdam	737M / 737M	80.5	74.7	91.6	83.5	0.230	93.1	82.1	93.8	92.1	0.419
	B-III: K-adapter	538M / 840M	80.5	70.8	<b>96.4</b>	89.6	0.215	94.4	81.4	94.8	89.4	0.259
	B-IV: LoRA	62M / 799M	71.1	62.9	92.9	84.8	0.366	89.8	74.0	94.0	90.5	0.800
	Ours (+K-adapter)	538M / 840M	86.3	78.9	<b>96.4</b>	<b>91.1</b>	0.132	<b>95.6</b>	88.1	94.9	90.3	0.118
	Ours (+LoRA)	62M / 799M	<b>90.5</b>	<b>90.6</b>	95.3	89.4	<b>0.073</b>	<b>95.6</b>	<b>95.2</b>	<b>95.1</b>	90.0	<b>0.117</b>

Table 3: The comparison of the continual learning results on zsRE (Large) and NQ datasets. We measure the accuracies on the knowledge  $\mathcal{K}_s$ ,  $\mathcal{K}_t$ , and the paraphrase knowledge  $\mathcal{P}_s$ ,  $\mathcal{P}_t$ , with the F/U ratio.

# Experiment

- selector successfully classifies 88.9% of examples, while 11.1% failed.

		Ground-truth	
		Source	Target
Selector Prediction	Source	19527 (40.7%)	854 (1.8%)
	Target	4473 (9.3%)	23146 (48.2%)

The confusion matrix of Selector.

		Ground-truth	
		Source	Target
Selector Prediction	Source	95.3	35.1
	Target	70.8 (0.0)	91.7 (97.4)

The accuracies of Ours/Retrieval in four cases.

# Calibrating Factual Knowledge in Pretrained Language Models

**Qingxiu Dong<sup>1</sup> \*, Damai Dai<sup>1</sup> \*, Yifan Song<sup>1</sup>, Jingjing Xu<sup>2</sup>, Zhifang Sui<sup>1</sup> and Lei Li<sup>3</sup>**

<sup>1</sup> Key Laboratory of Computational Linguistics, Peking University, MOE, China

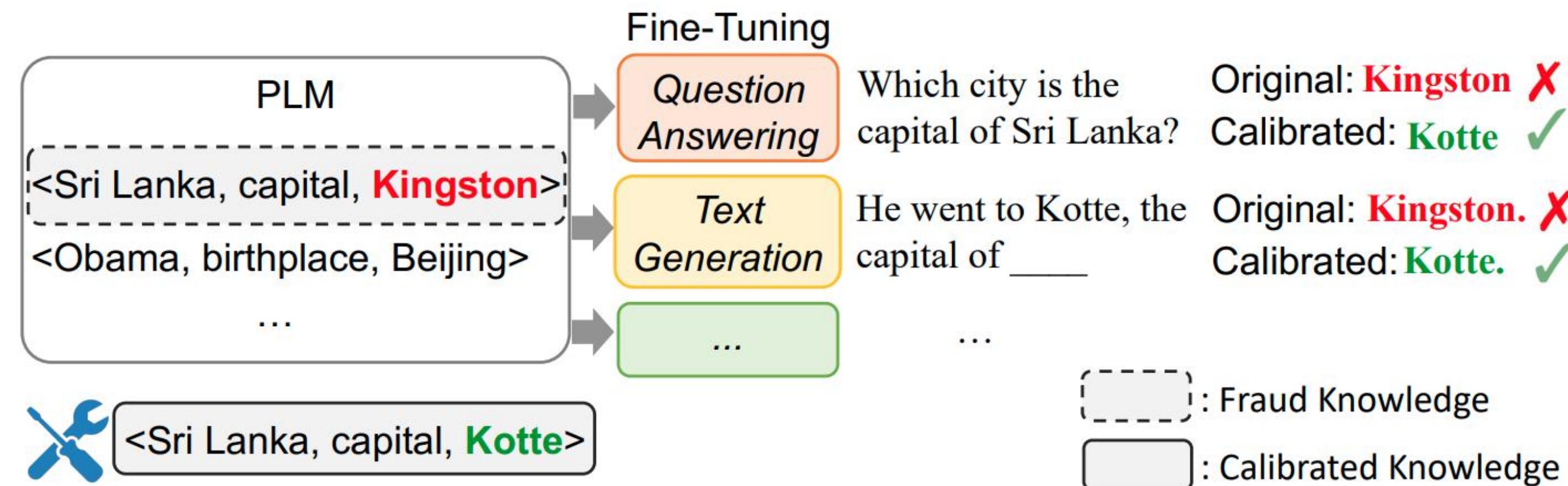
<sup>2</sup> Shanghai AI Lab <sup>3</sup> University of California, Santa Barbara

dqx@stu.pku.edu.cn, {daidamai,yfsong,szf}@.pku.edu.cn,  
jingjingxu@pku.edu.cn, lilei@cs.ucsb.edu

*Findings of EMNLP 2022*

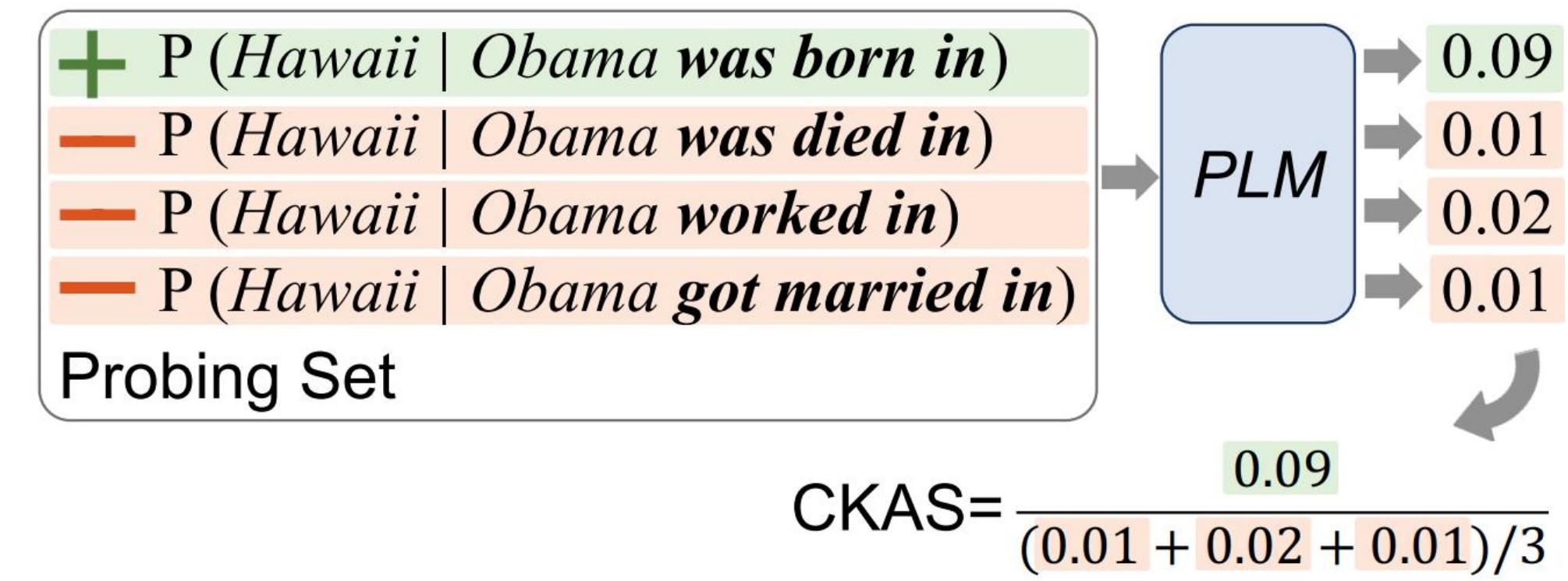
# Motivation

- PLMs have factual errors



# Methods

- detect wrong facts
- rank-based metric
  - inexhaustible answers
  - frequency bias
- Contrastive Knowledge Assessment



$$\text{CKA}_M(s, r, o) = \frac{P_M(o|s, r) + \alpha}{\mathbb{E}_{r'} [P_M(o|s, r')] + \alpha}$$

threshold (usually < 1.0) for false facts

# Methods

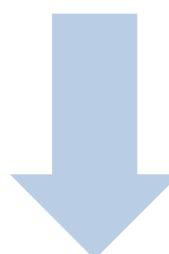
- detect wrong facts

Fact	Rank-based Assessment		CKA	
	Assess	Top-3 Prediction	Assess	Score
<b>Inexhaustible Answers</b>				
Germany shares border with <i>Czech Republic</i> .	✗	France, Russia, Austria	✓	4.45
India is a member of <i>UN</i> .	✗	NATO, India, AS	✓	2.27
Frederick was born in <i>Berlin</i> .	✗	Frederick, 18, Baltimore	✓	3.52
<b>Frequency Bias</b>				
Adi Shankara is affiliated with the <i>Hindu</i> religion.	✓	Hindu, Ko, Si	✗	0.98
Adi Shankara is against the <i>Hindu</i> religion.	-	Hindu, religion, Buddhist	-	-

# Methods

- Knowledge Calibration

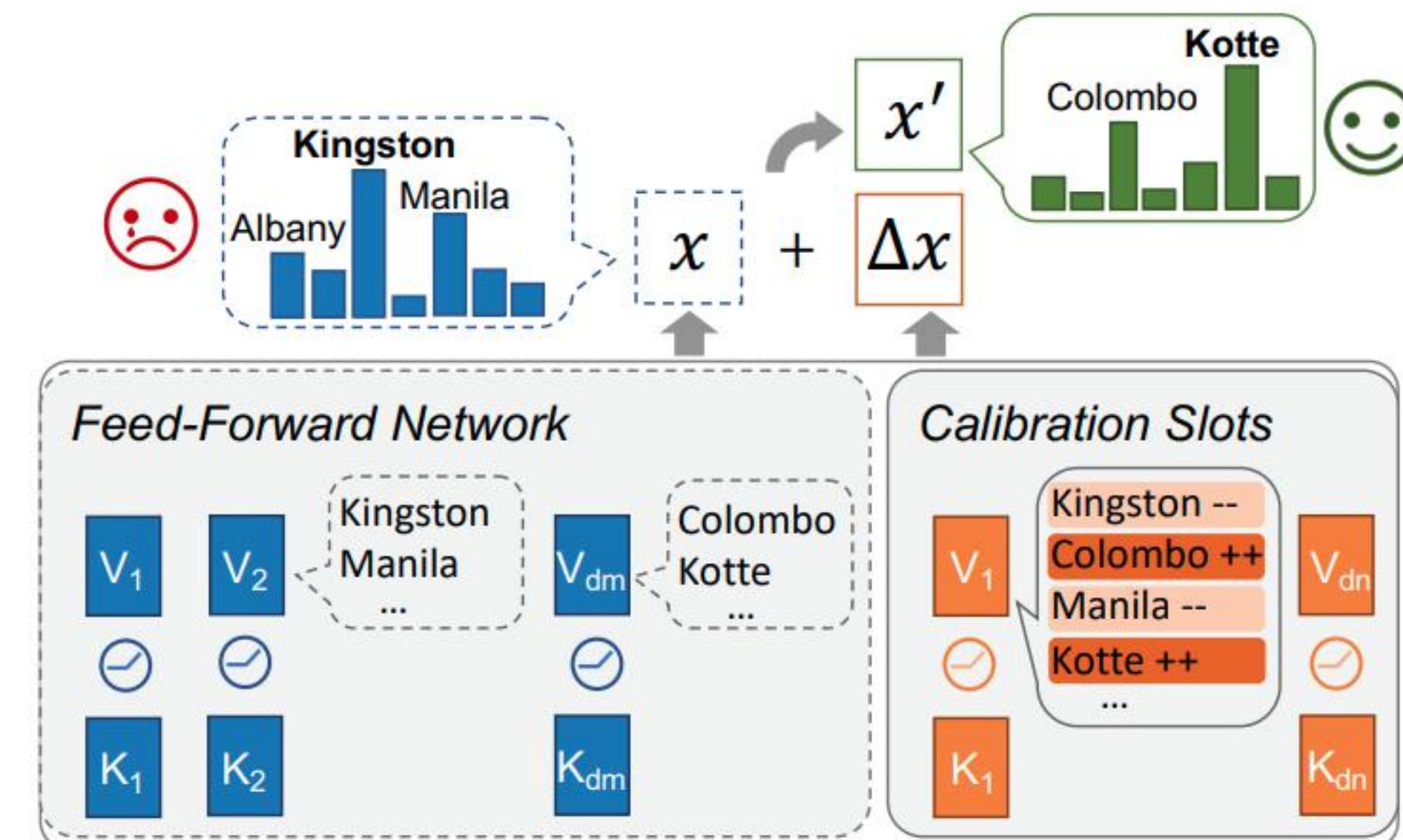
$$\text{FFN}(H) = \text{GELU}(HK^T)V, \quad K, V \in \mathbb{R}^{d_m \times d}$$



$$\Delta \text{FFN}(H) = \text{GELU}(H\tilde{K}^T)\tilde{V},$$

$$\text{FFN}'(H) = \text{FFN}(H) + \Delta \text{FFN}(H),$$

$$\tilde{K}, \tilde{V} \in \mathbb{R}^{d_c \times d} \quad d_c \ll d_m$$



*The capital of Sri Lanka is Kotte.*

# Data& Model Setup

- factual triplets from the T-REx dataset
- fill head entity and tail entity into the template in LAMA
- 100/1000 facts
- write 3 erroneous relation templates for each relation in LAMA

# Data & Model Setup

- PARAREL dataset
  - contains various surface form templates for 38 relations
- masked language modeling objective on calibration data

Split	Source	Target
Train	[MASK] was born in Hawaii.	Obama
	Obama is originally from [MASK].	Hawaii
	[MASK] was originally from Hawaii.	Obama
	Obama is native to [MASK].	Hawaii
Valid	[MASK] originates from Hawaii.	Obama
	Obama originated from [MASK].	Hawaii
Test	Obama is a/an [MASK]-born person.	Hawaii
	[MASK] was native to Hawaii.	Obama
	Obama, a [MASK]-born person.	Hawaii

Table 2: Example of knowledge-intensive data for training CALINET<sup>🔧</sup>. We generate multiple texts via templates for each triple where the templates in training, validation, and test are not sharing.

# Experiment

target entity is replaced by a false one in the same entity type.

randomly mask the test data as pretraining

<b>Model</b>	<b># Facts</b>	<b>Method</b>	<b># Calibration Params</b>	<b>False Rate(↓)</b>	<b>Ori (↓)</b>	<b>Adv (↑)</b>	<b>LM(↓)</b>	<b>EM(↑)</b>	<b>F1(↑)</b>
continue pretraining method	10 <sup>2</sup>	Vanilla	0	48.10%	87.21	219.18	89.21	0.63	7.48
		CALINET 	0.1M	17.09%	1.22	>1000	54.45	81.65	84.58
	T5-base	C. P.	220M	13.29%	1.15	>1000	116.52	87.34	89.85
		Vanilla	0	51.34%	90.61	208.90	60.64	0.94	6.51
		CALINET 	0.5M	18.30%	1.26	>1000	46.71	71.18	73.48
		C. P.	220M	18.23%	1.28	>1000	139.96	78.15	80.35
T5-large	10 <sup>2</sup>	Vanilla	0	46.20%	34.36	116.38	92.52	2.53	7.23
		CALINET 	0.5M	15.19%	1.30	>1000	44.21	81.65	85.11
	T5-large	C. P.	770M	14.56%	1.21	>1000	477.24	87.97	90.49
		Vanilla	0	45.04%	31.44	93.77	58.78	2.48	6.86
		CALINET 	1.0M	20.84%	1.32	>1000	43.04	70.84	72.92
		C. P.	770M	17.16%	1.28	>1000	154.52	78.22	80.57

"Ori." and "Adv." refer to model perplexity on the original test set (contains true facts) and the adversarial test set (contains false facts), respectively. We concatenate CALINET to the last layer of the T5 decoder.

# Experiment

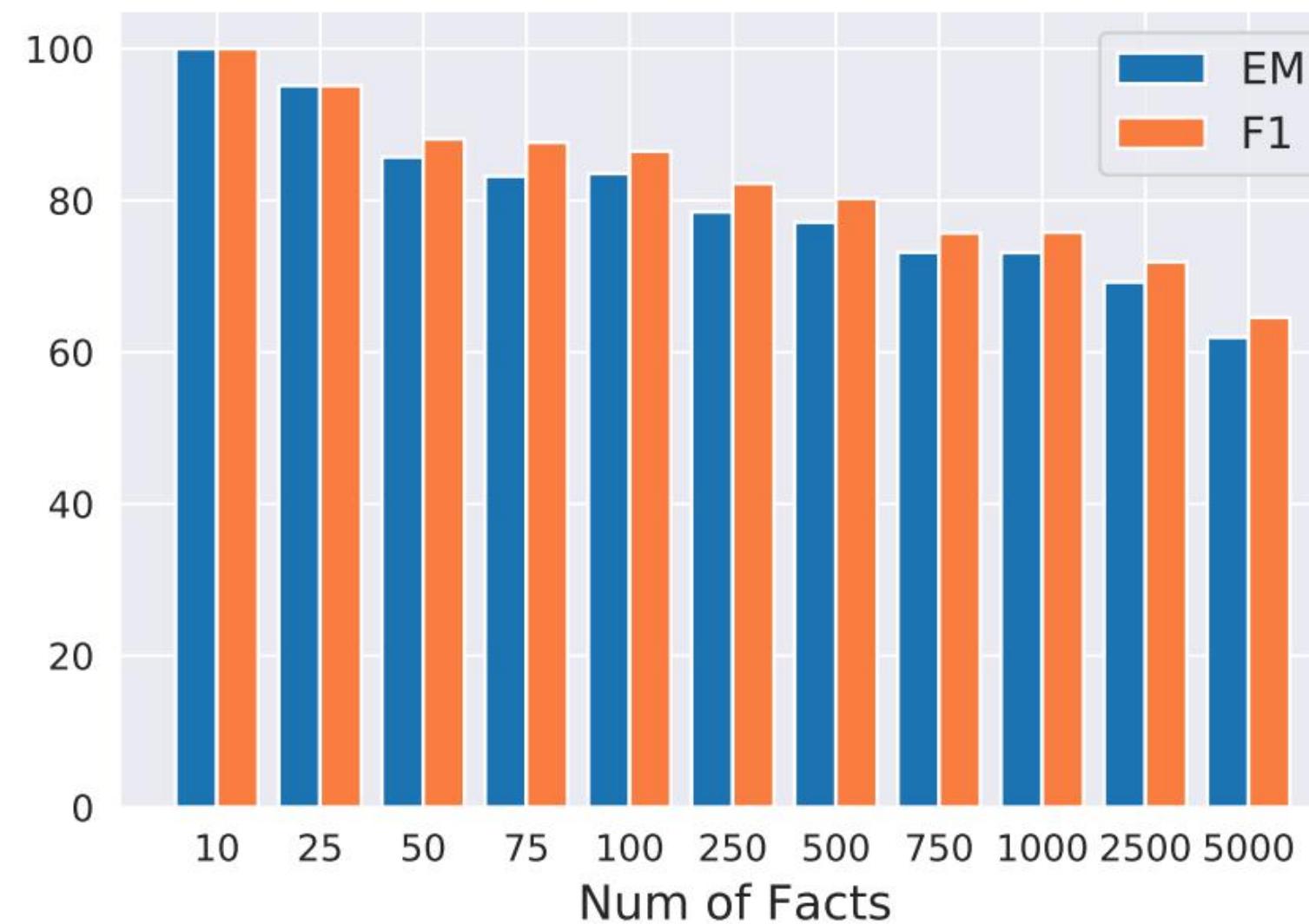


Figure 4: Calibration results for different scales of facts. Given 5000 facts, our method can calibrate more than 60% of facts in PLMs at once.

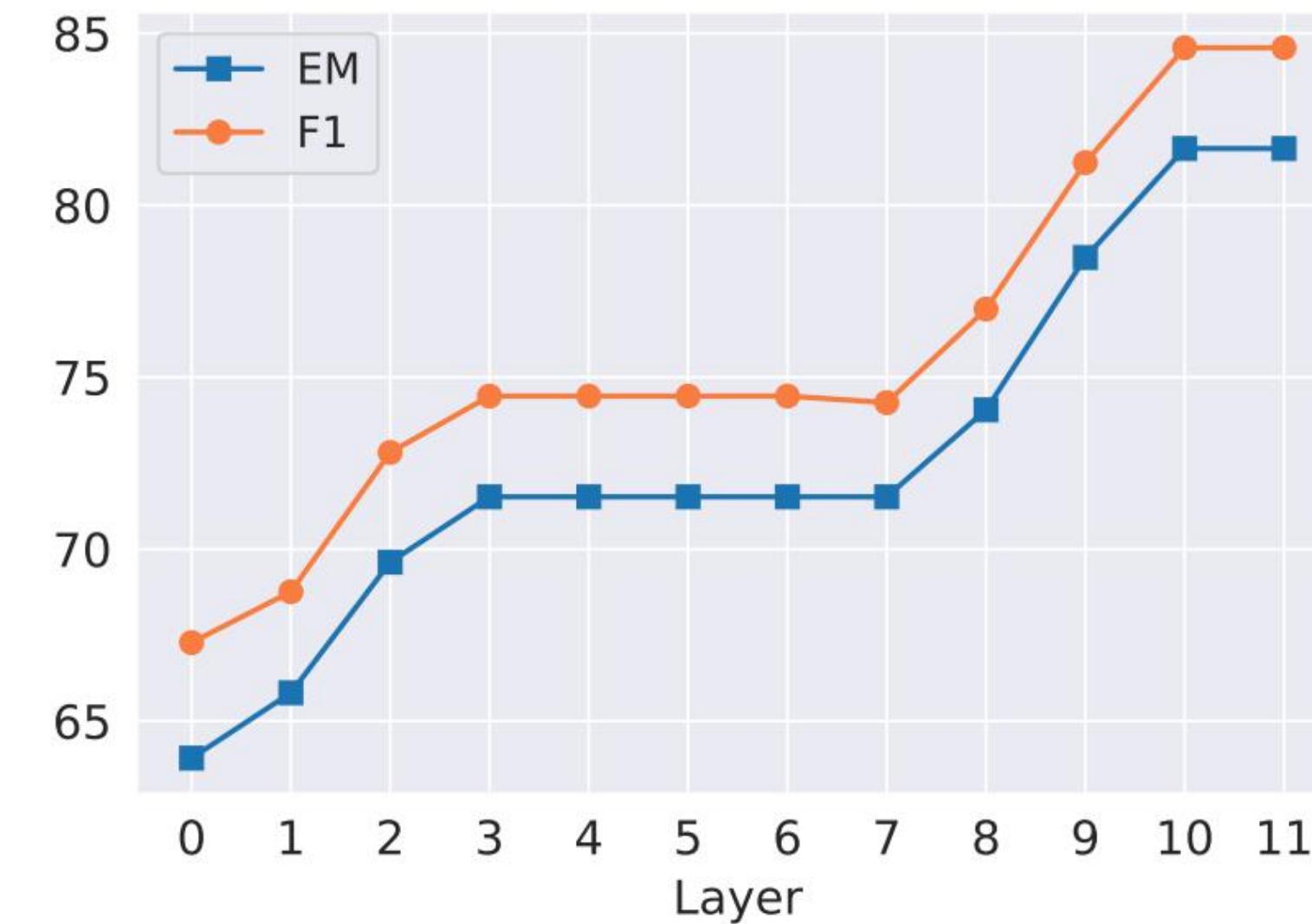


Figure 5: Calibration ability of concatenating CaliNet in different layers.

---

# **Repairing Neural Networks by Leaving the Right Past Behind**

---

**Ryutaro Tanno**  
Microsoft Research Cambridge, UK  
[rytanno@microsoft.com](mailto:rytanno@microsoft.com)

**Aditya Nori**  
Microsoft Research Cambridge, UK  
[Aditya.Nori@microsoft.com](mailto:Aditya.Nori@microsoft.com)

**Melanie F. Pradier**  
Microsoft Research Cambridge, UK  
[melanief@microsoft.com](mailto:melanief@microsoft.com)

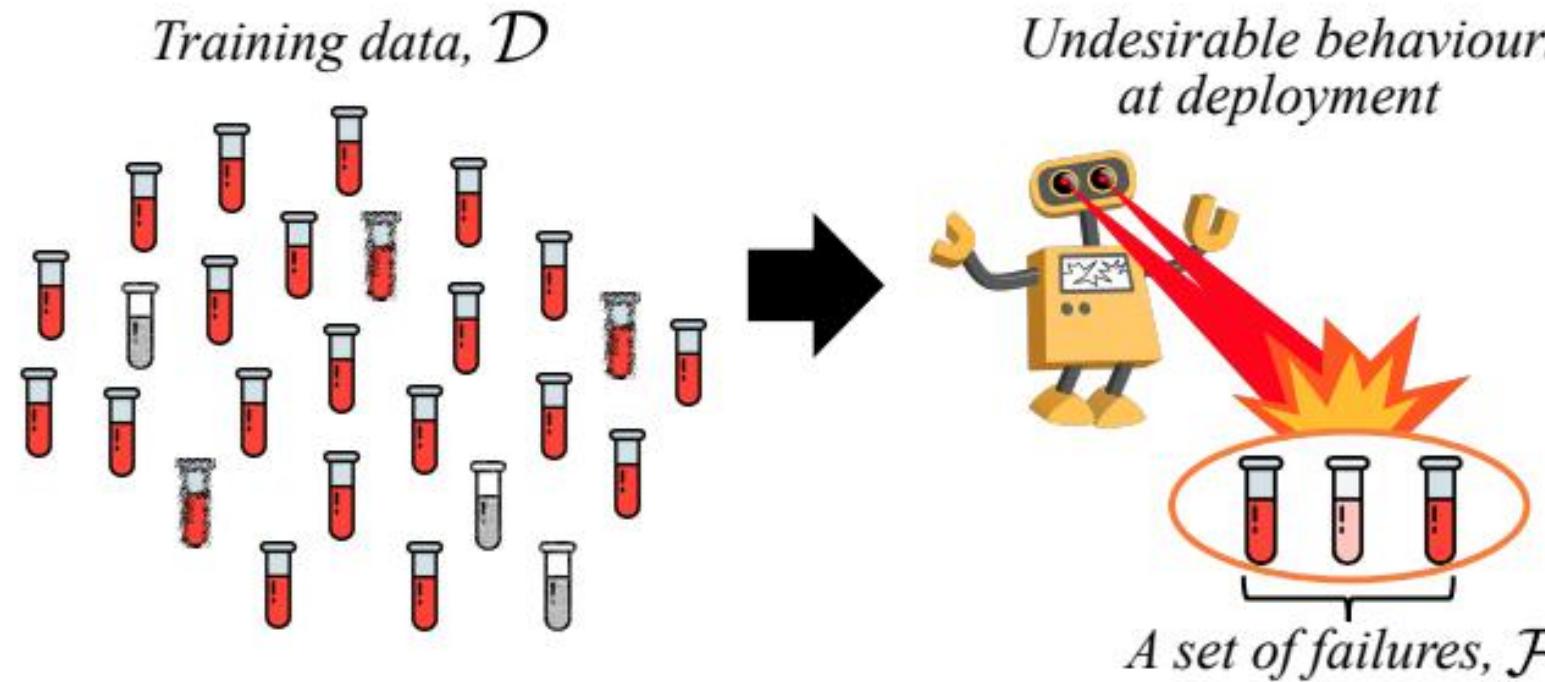
**Yingzhen Li**  
Imperial College London, UK  
[yingzhen.li@imperial.ac.uk](mailto:yingzhen.li@imperial.ac.uk)

**NeurIPS 2022**

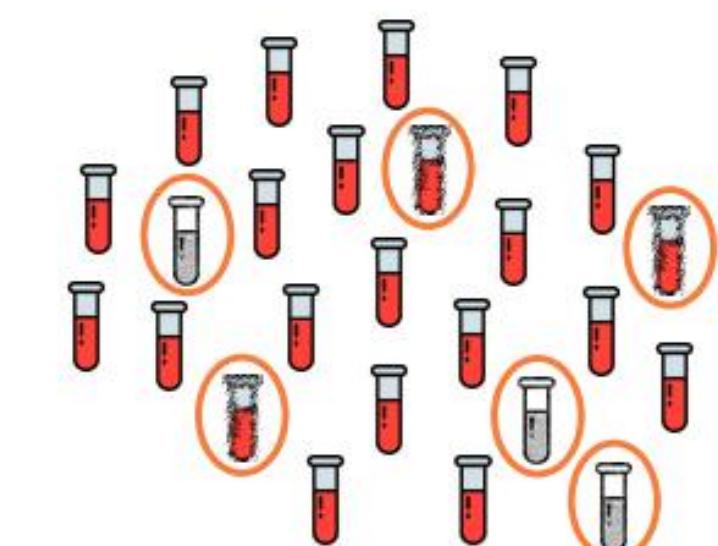
# Motivation

- annotation noise, low-quality inputs, anomalies, and acquisition biases
- finding & erasing

(a) Training Data with Problems  
& Resultant Failures at Deployment



(b) Identify the "Causes",  $\mathcal{C}$   
of Target Failures,  $\mathcal{F}$



(c) Repair Model by Erasing  
Information of Failure Causes,  $\mathcal{C}$



# Model Repairment by Data Deletion

- **Cause identification:** Identify a set of detrimental datapoints **C** in the training data **D** that contributed the most to the failure set **F**.
- **Treatment:** Given the set of failure causes **C**, adapt the model to predict correctly on the failure set **F**, while maintaining performance on remaining test examples. -> unlearn **C**

# Cause identification

- Given  $D, F$ , finding  $C \subseteq D$

$$r(C) := \log p(F|D \setminus C) - \log p(F|D) \quad C = \operatorname{argmax}_{C' \in \mathbb{P}(D)} r(C')$$

unacceptable

$$\begin{aligned} \log p(F|D \setminus C) &= \log \int p(F|\theta) p(\theta|D \setminus C) d\theta \\ &= \log \int p(F|\theta) \frac{p(D \setminus C|\theta)p(\theta)}{p(D \setminus C)} d\theta \quad (\text{Bayes' rule}) \\ &= \log \int p(F|\theta) \frac{p(D|\theta)p(\theta)}{p(C|\theta)p(D \setminus C)} d\theta \quad (\text{by Eq. (18)}) \\ &= \log \int \frac{p(D)}{p(D \setminus C)} \cdot \frac{p(F|\theta)}{p(C|\theta)} \cdot \frac{p(D|\theta)}{p(D)} d\theta \quad (\text{multiplying } \frac{p(D)}{p(D)} \text{ and rearranging terms}) \\ &= \log \int \frac{p(F|\theta)}{p(C|\theta)} p(\theta|D) d\theta + \log \frac{p(D)}{p(D \setminus C)}. \quad (\text{Bayes' rule}) \end{aligned}$$

# Cause identification

$$r(\mathcal{C}) \coloneqq \log p(\mathcal{F}|\mathcal{D} \setminus \mathcal{C}) - \log p(\mathcal{F}|\mathcal{D})$$

$$\log p(\mathcal{F}|\mathcal{D} \setminus \mathcal{C}) = \log \int \frac{p(\mathcal{F}|\boldsymbol{\theta})}{p(\mathcal{C}|\boldsymbol{\theta})} p(\boldsymbol{\theta}|\mathcal{D}) d\boldsymbol{\theta} + \log \frac{p(\mathcal{D})}{p(\mathcal{D} \setminus \mathcal{C})}$$



$$r(\mathcal{C}) = \log \left( \frac{p(\mathcal{F}|\mathcal{D} \setminus \mathcal{C})}{p(\mathcal{F}|\mathcal{D})} \right)$$

$$= \log \int \frac{p(\mathcal{F}|\boldsymbol{\theta})}{p(\mathcal{C}|\boldsymbol{\theta})} p(\boldsymbol{\theta}|\mathcal{D}) d\boldsymbol{\theta} - \log p(\mathcal{F}|\mathcal{D}) + \log p(\mathcal{D}) - \log \int p(\mathcal{D} \setminus \mathcal{C}|\boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta}$$

(by definition of marginal distributions)

$$= \log \int \frac{1}{p(\mathcal{C}|\boldsymbol{\theta})} \frac{p(\mathcal{F}|\boldsymbol{\theta})p(\boldsymbol{\theta}|\mathcal{D})}{p(\mathcal{F}|\mathcal{D})} d\boldsymbol{\theta} - \log \int \frac{p(\mathcal{D}|\boldsymbol{\theta})}{p(\mathcal{C}|\boldsymbol{\theta})} \frac{p(\boldsymbol{\theta})}{p(\mathcal{D})} d\boldsymbol{\theta}$$

(by Eq. (18) and rearranging terms)

$$= \log \int \frac{p(\boldsymbol{\theta} | \mathcal{D}, \mathcal{F})}{p(\mathcal{C} | \boldsymbol{\theta})} d\boldsymbol{\theta} - \log \int \frac{p(\boldsymbol{\theta} | \mathcal{D})}{p(\mathcal{C} | \boldsymbol{\theta})} d\boldsymbol{\theta}. \quad (\text{Bayes' rule})$$

# Cause identification

$$r(\mathcal{C}) = \log \int \frac{p(\boldsymbol{\theta}|\mathcal{D},\mathcal{F})}{p(\mathcal{C}|\boldsymbol{\theta})} d\boldsymbol{\theta} - \log \int \frac{p(\boldsymbol{\theta}|\mathcal{D})}{p(\mathcal{C}|\boldsymbol{\theta})} d\boldsymbol{\theta}$$

only need to compute  $p(\boldsymbol{\theta}|\mathcal{D},\mathcal{F})$  once, but still need combinatorial search for the best subset  $\mathcal{C}$

re-write  $r(\mathcal{C}) = F(1, p(\boldsymbol{\theta}|\mathcal{D},\mathcal{F})) - F(1, p(\boldsymbol{\theta}|\mathcal{D}))$ ,  $F(\epsilon, g(\boldsymbol{\theta})) := \log \int g(\boldsymbol{\theta}) e^{-\epsilon \log p(\mathcal{C}|\boldsymbol{\theta})} d\boldsymbol{\theta}$

Taylor expansion of  $F(\epsilon, g(\boldsymbol{\theta}))$  around  $\epsilon = 0$

$$F(\epsilon, g(\boldsymbol{\theta})) = -\epsilon \mathbb{E}_{g(\boldsymbol{\theta})} [\log p(\mathcal{C}|\boldsymbol{\theta})] + \mathcal{O}(\epsilon^2)$$

re-write  $\hat{r}(\mathcal{C}) := \mathbb{E}_{p(\boldsymbol{\theta}|\mathcal{D})} [\log p(\mathcal{C}|\boldsymbol{\theta})] - \mathbb{E}_{p(\boldsymbol{\theta}|\mathcal{D},\mathcal{F})} [\log p(\mathcal{C}|\boldsymbol{\theta})]$

$\mathcal{O}(|\mathcal{D}|!)$   $\rightarrow$   $\mathcal{O}(|\mathcal{D}|)$

data i.i.d.  $\hat{r}(z) = \mathbb{E}_{p(\boldsymbol{\theta}|\mathcal{D})} [\log p(z|\boldsymbol{\theta})] - \mathbb{E}_{p(\boldsymbol{\theta}|\mathcal{D},\mathcal{F})} [\log p(z|\boldsymbol{\theta})]$ ,  $p(z|\boldsymbol{\theta}) = p(y|\mathbf{x}, \boldsymbol{\theta})$

# Cause identification

- $p(\theta|\mathcal{D}, \mathcal{F})$  -> recomputing the posterior from scratch can be expensive
  - MLE/MAP point estimates, Laplace approximation, variational inference, etc
- Example I (Linear Influence Function)

$$\hat{\theta}_{\mathcal{D}, \mathcal{F}}^* \approx \hat{\theta} + \gamma \hat{F}_{\hat{\theta}}^{-1} \nabla_{\hat{\theta}} \log p(\mathcal{F}|\hat{\theta}) \quad \tilde{r}(z) = -\gamma \nabla_{\hat{\theta}} \log p(\mathcal{F}|\hat{\theta})^\top \hat{F}_{\hat{\theta}}^{-1} \nabla_{\hat{\theta}} \log p(z|\hat{\theta})$$

- Example II (Elastic Weight Consolidation)

$$\text{maximize } \log p(\mathcal{F}|\theta) - \frac{N}{2} (\theta - \hat{\theta})^\top \hat{F}_{\hat{\theta}} (\theta - \hat{\theta}) - \frac{\lambda}{2} \|\theta - \hat{\theta}\|_2^2, \quad \tilde{r}(z) = \log p(z|\hat{\theta}) - \log p(z|\hat{\theta}_{\mathcal{D}, \mathcal{F}}^*)$$

# Treatment

- $p(\theta|\mathcal{D}\setminus\mathcal{C}) \rightarrow$  retrain is time-consuming
- Example I: Fine-tuning on Corrected Data
- Example II: (Newton Update Removal)

$$\hat{\theta}_{\mathcal{D}\setminus\mathcal{C}}^{**} \approx \hat{\theta} - \gamma \hat{F}_{\hat{\theta}}^{-1} \nabla_{\hat{\theta}} \log p(\mathcal{C}|\hat{\theta})$$

- Example III: (EWC for data deletion)

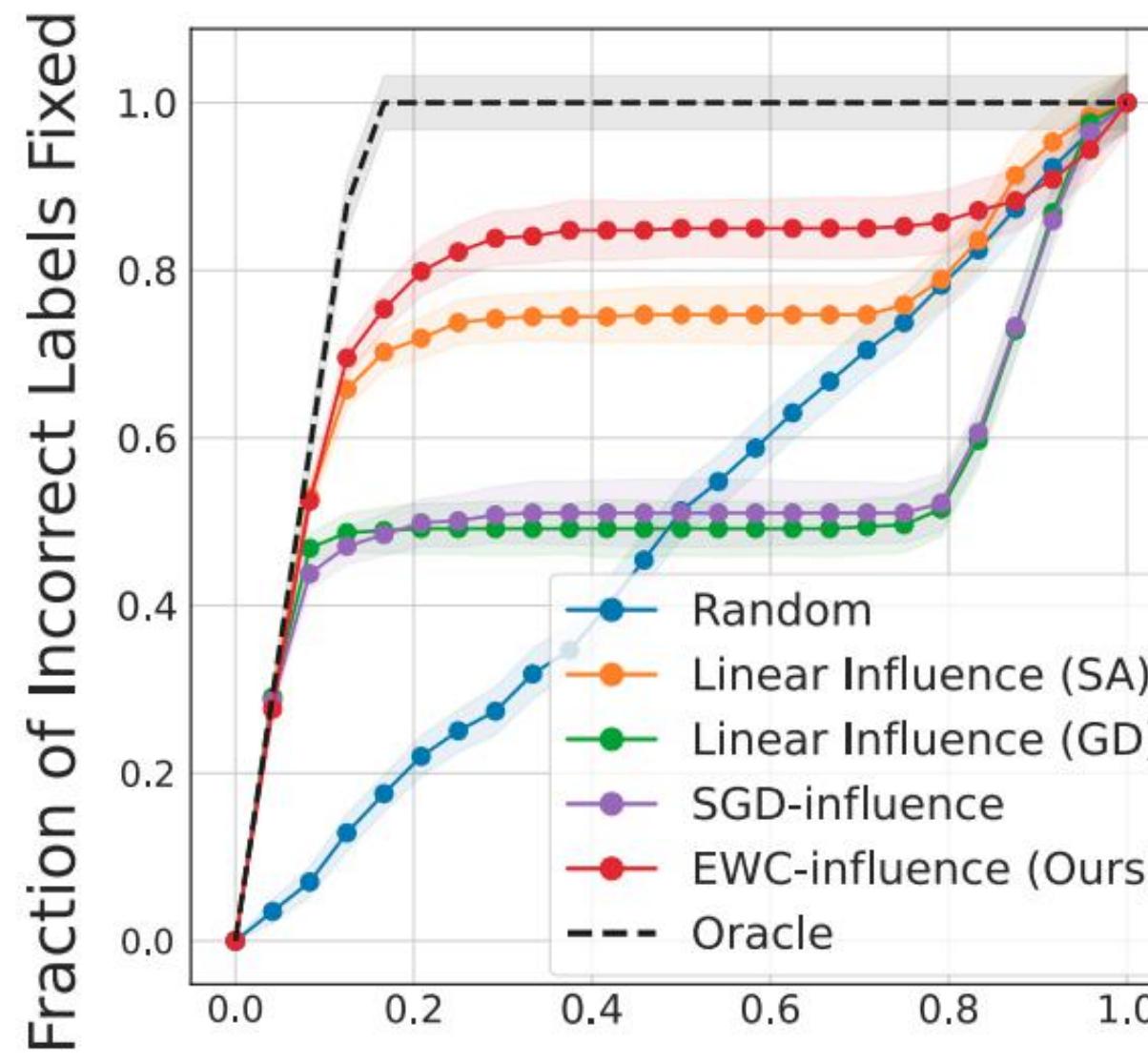
$$\text{maximize } -\log p(\mathcal{C}|\theta) - \frac{N}{2} (\theta - \hat{\theta})^\top \hat{F}_{\hat{\theta}} (\theta - \hat{\theta}) - \frac{\lambda}{2} \|\theta - \hat{\theta}\|_2^2$$

# Experiment

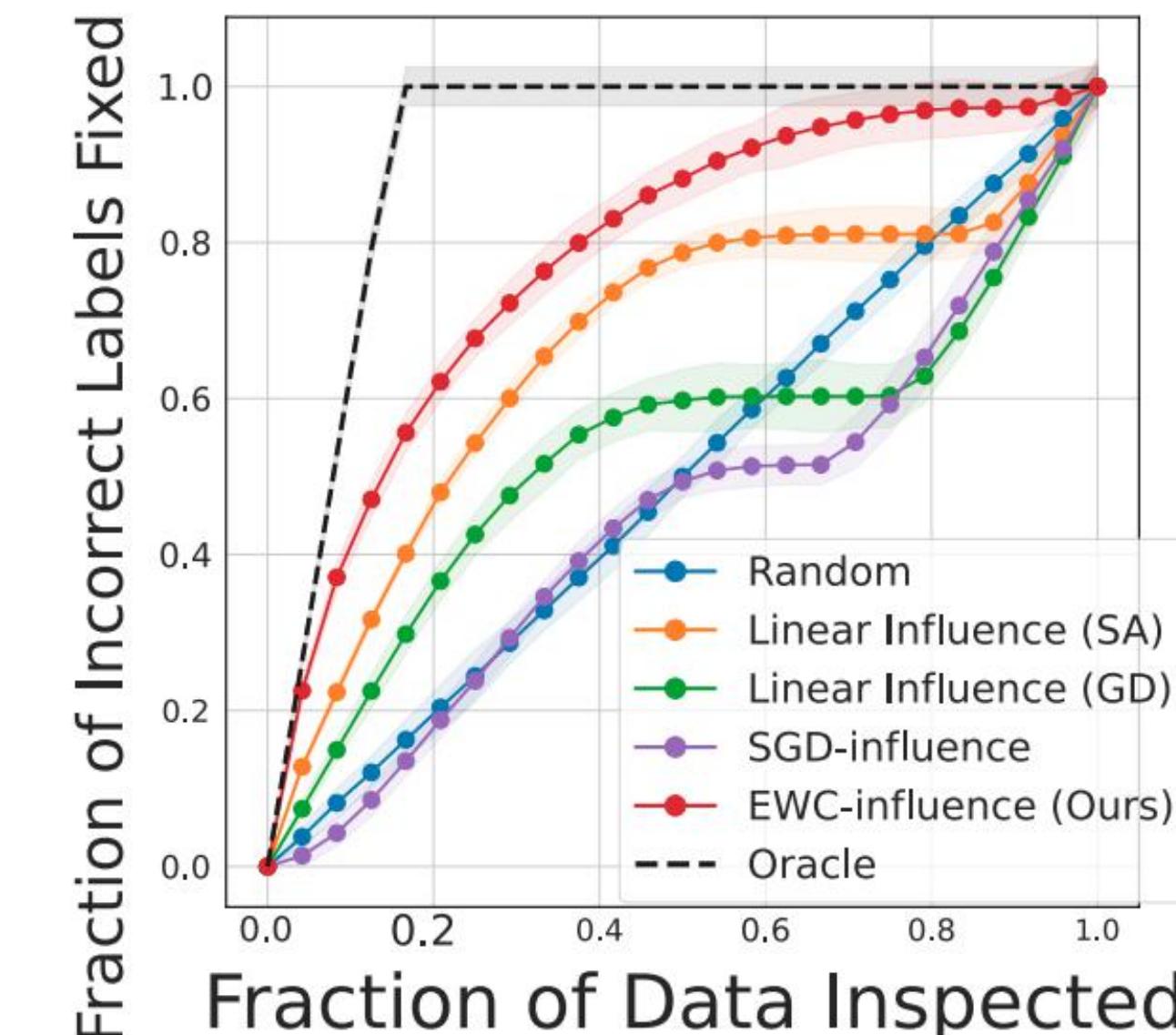
- Dataset
  - CIFAR10
  - MNIST (only use 6% training set)
- Evaluation
  - separate the test set into  $F$  (“failure set”) and  $T \setminus F$  (“remaining set”)
  - split the failure set into query and holdout set, use query to identify failure causes  $C$ , use the latter to evaluate

# Identifying Failure Causes

- randomly flip labels for two classes, use  $F_q$  to find incorrectly labelled samples



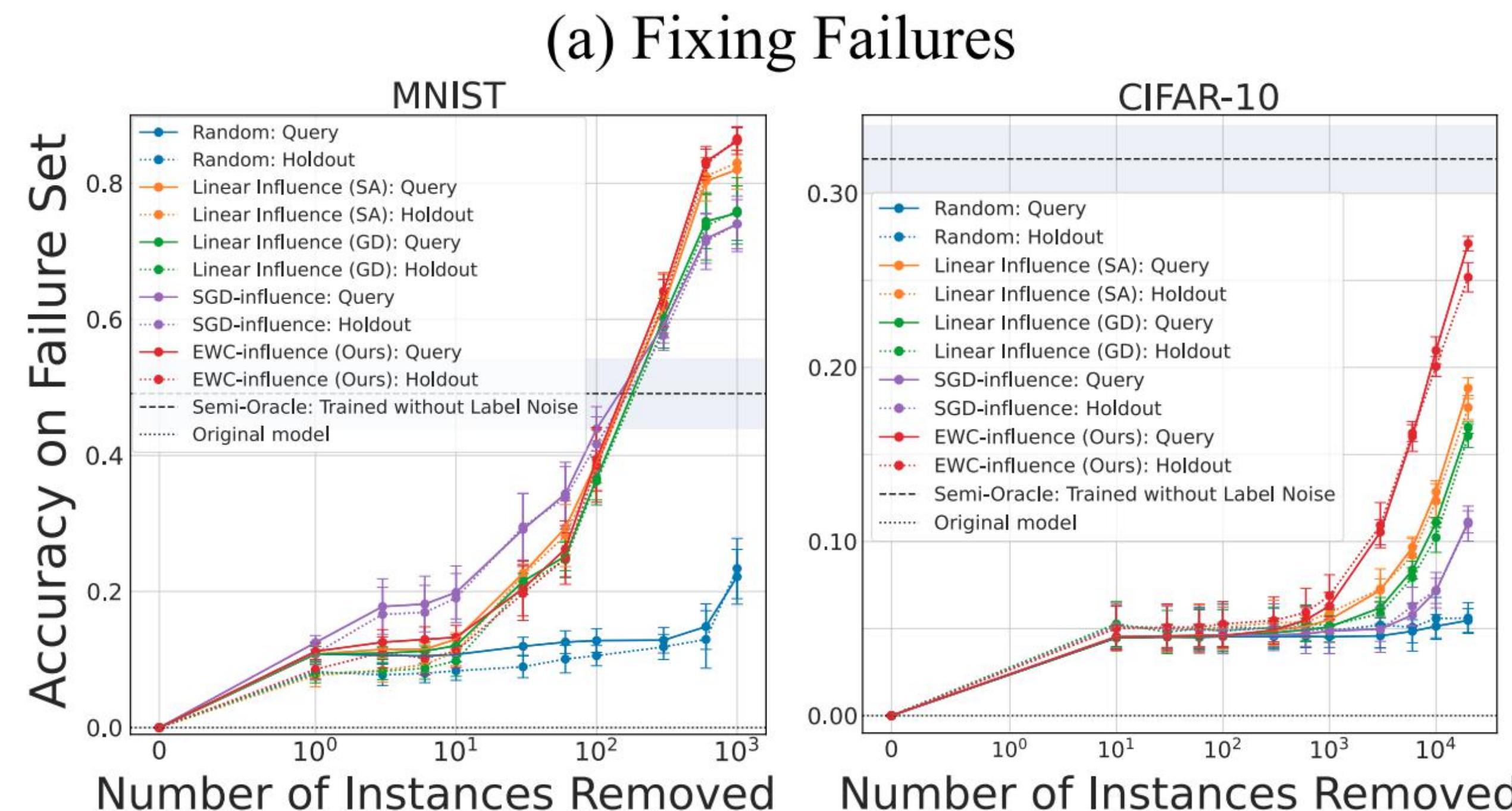
MNIST



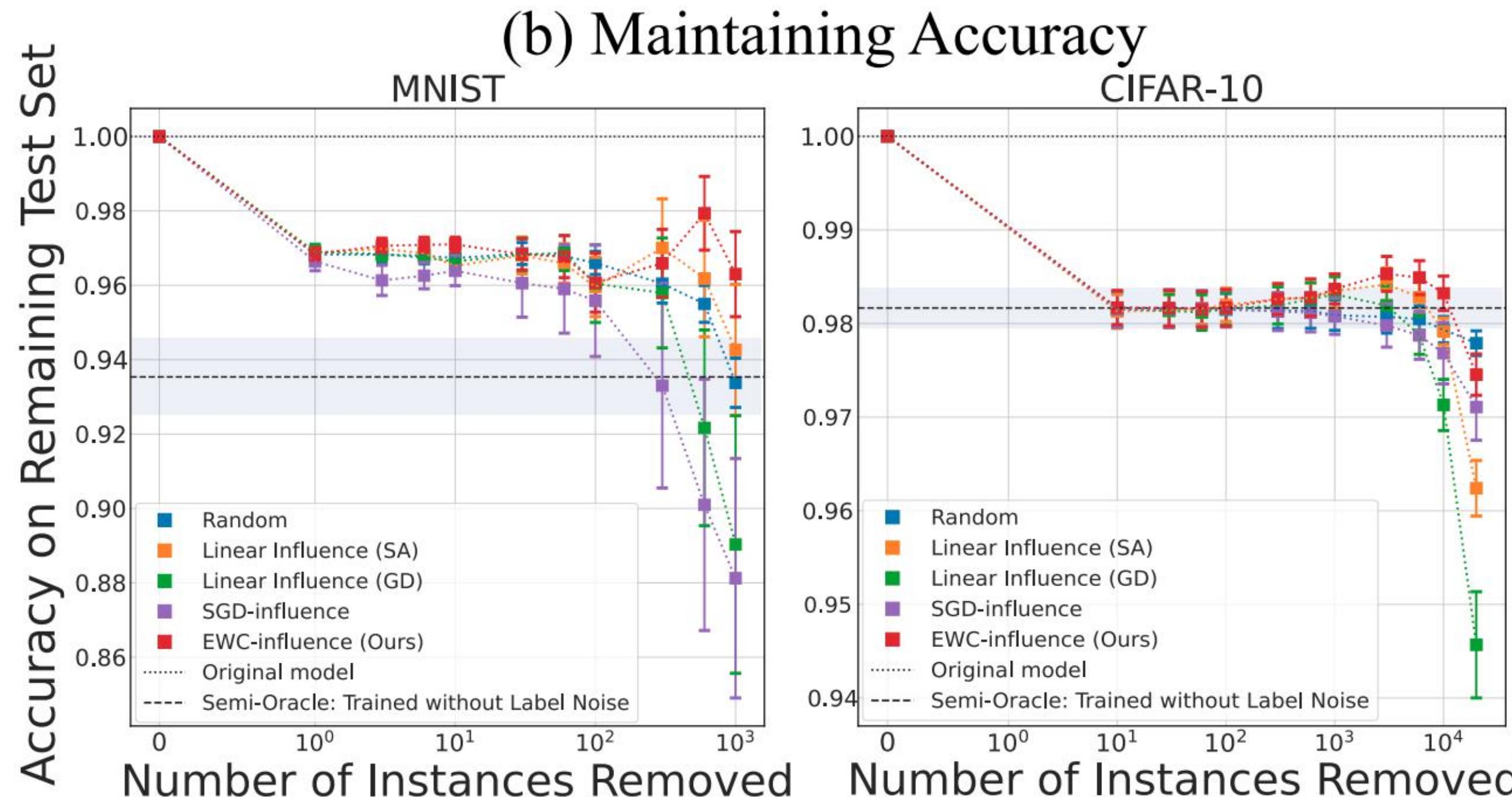
CIFAR10

# Identifying Failure Causes

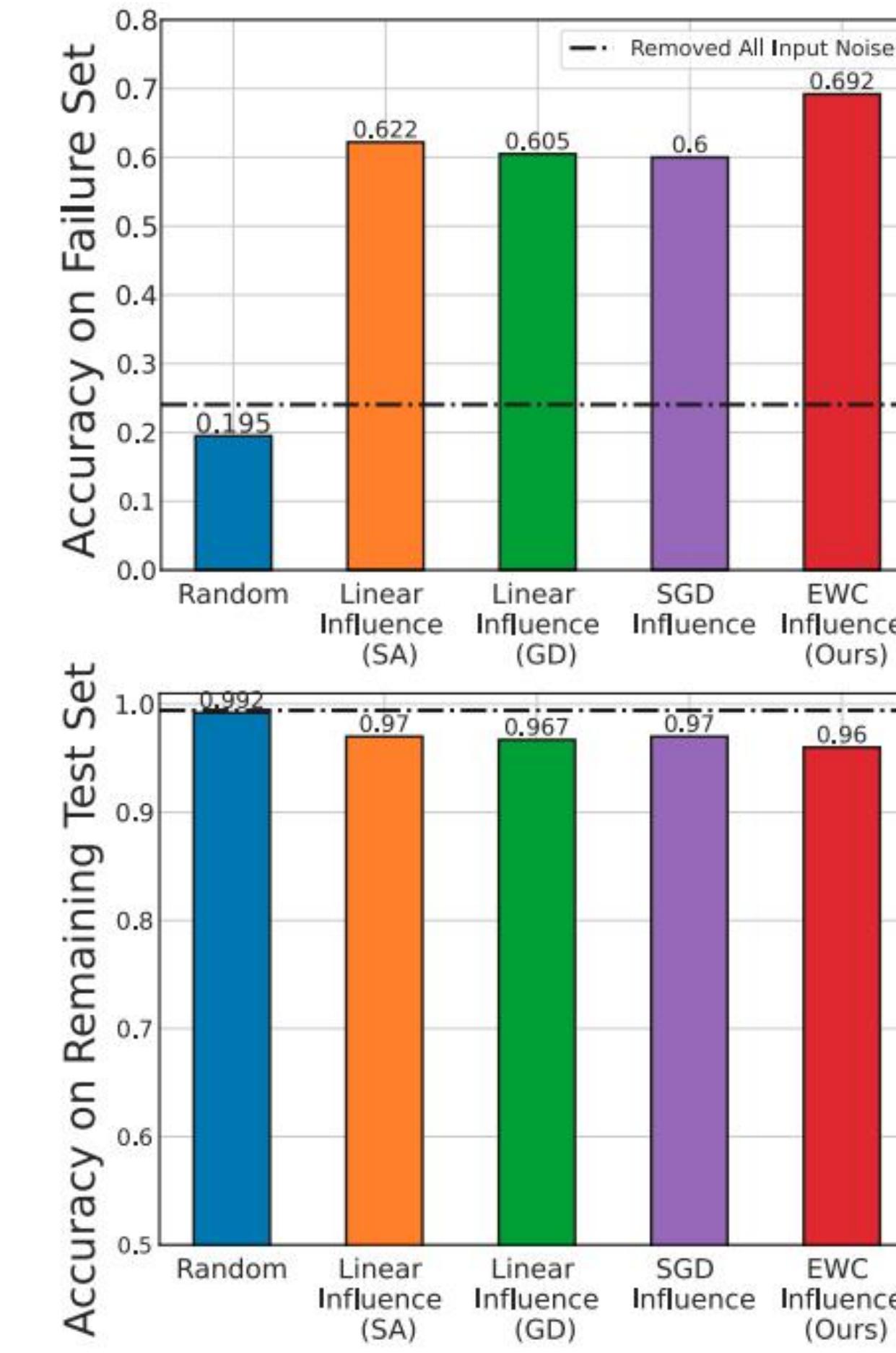
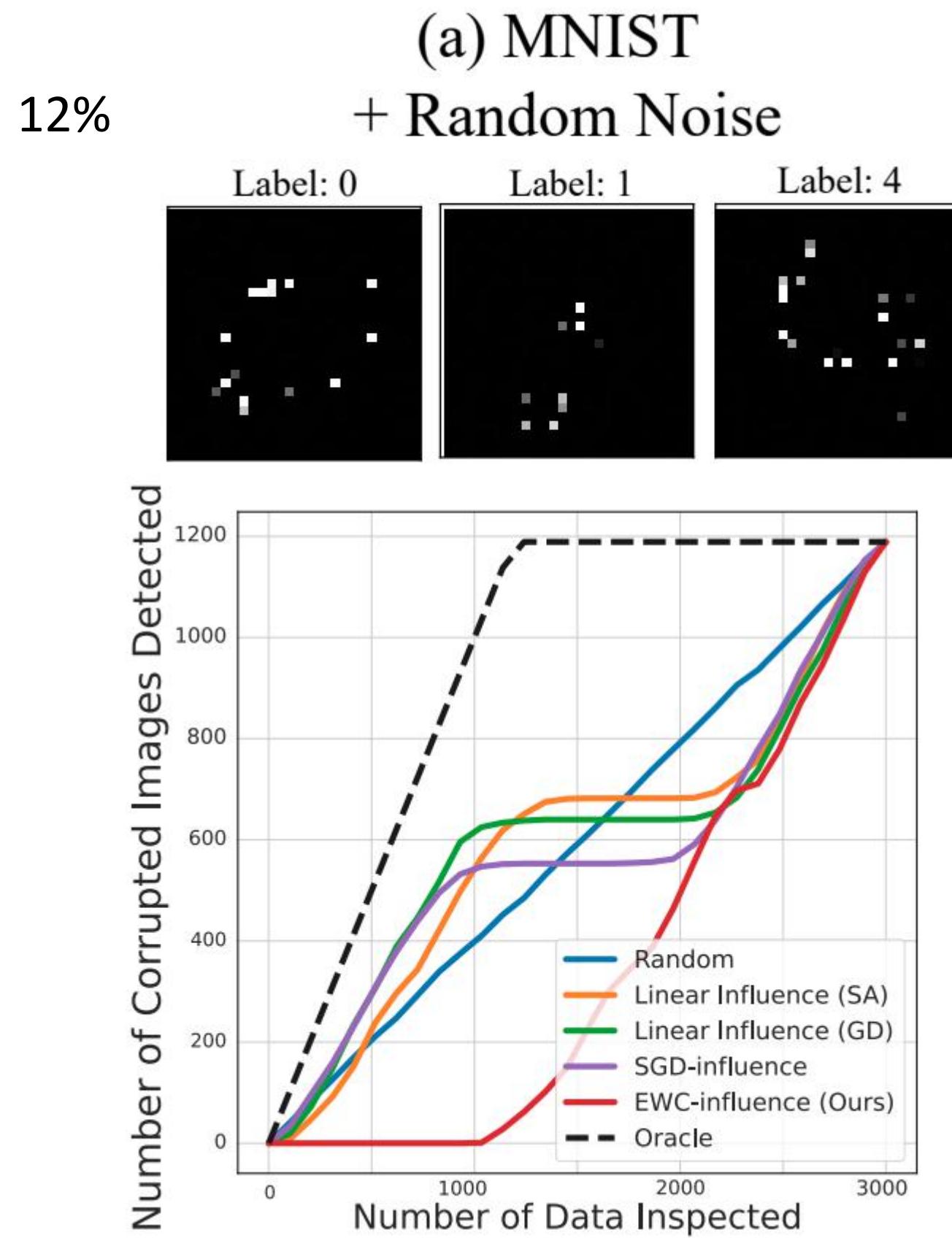
- fine-tune the base model on D\c



# Identifying Failure Causes

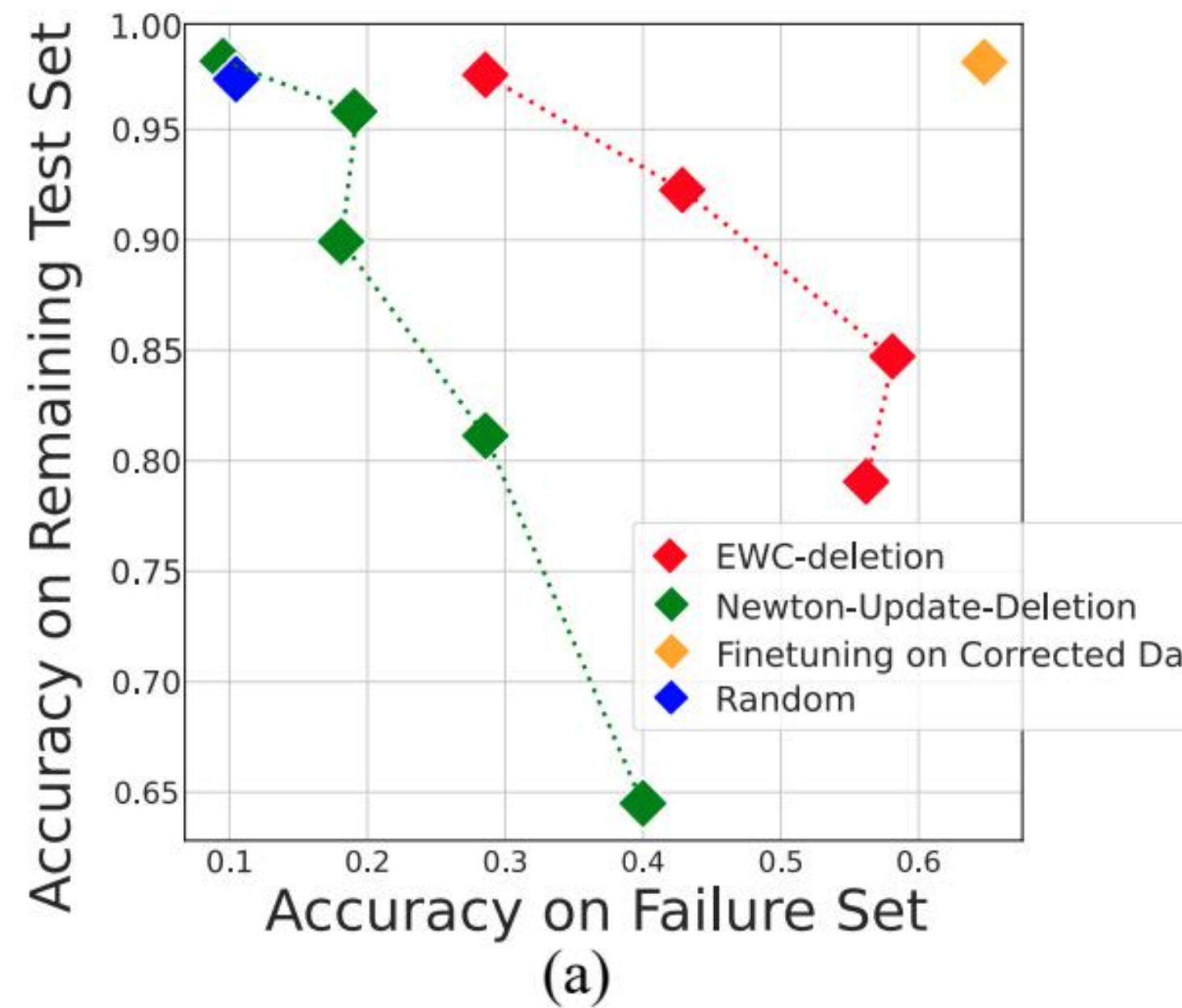


# Identifying Failure Causes

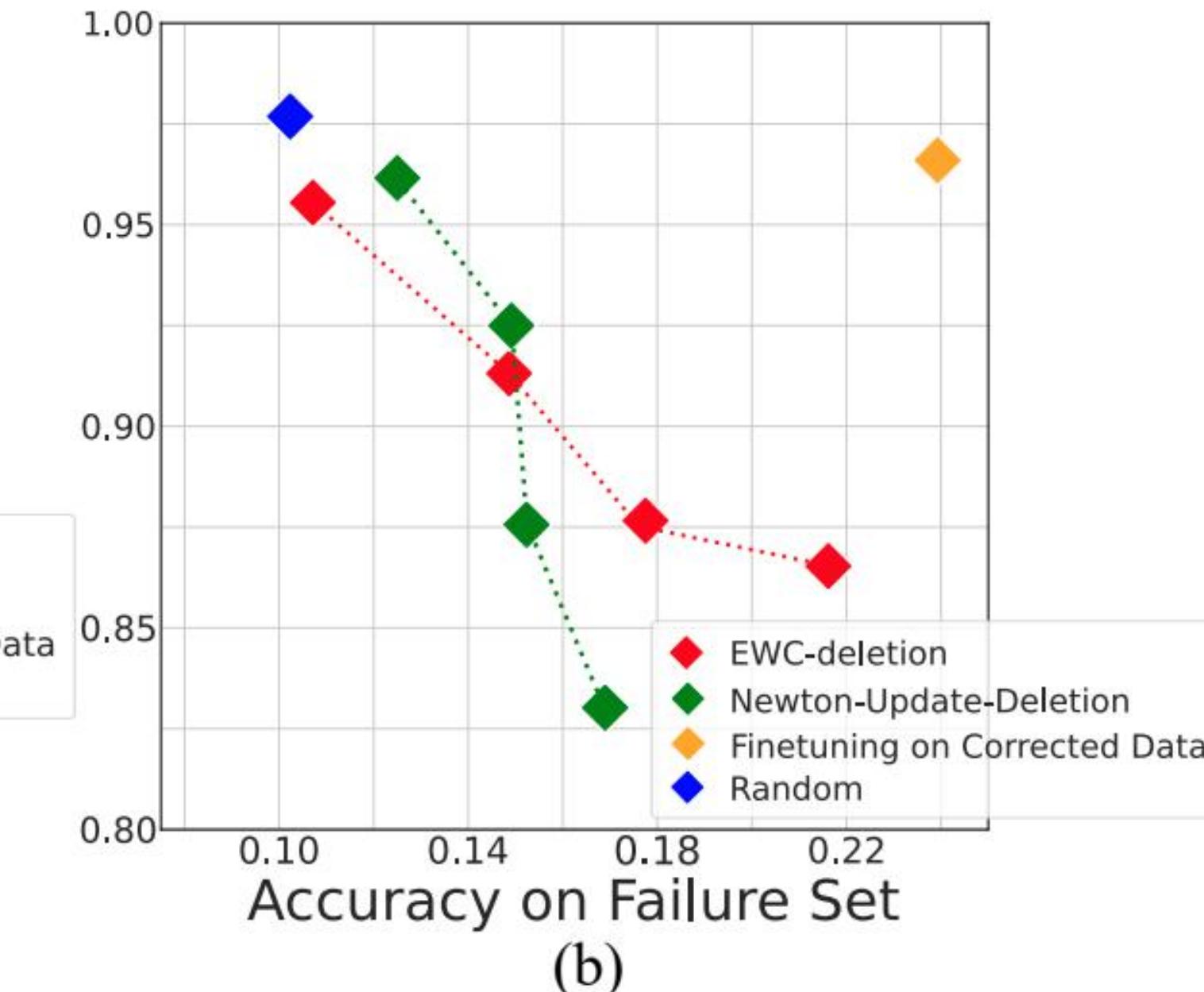


# Comparison of Treatment Methods

different  
step sizes



(a)



(b)