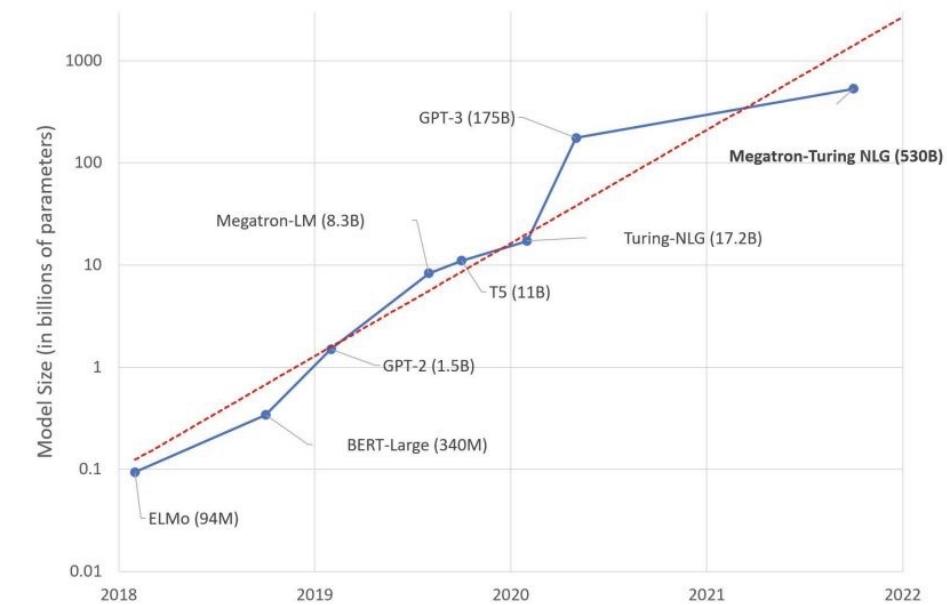
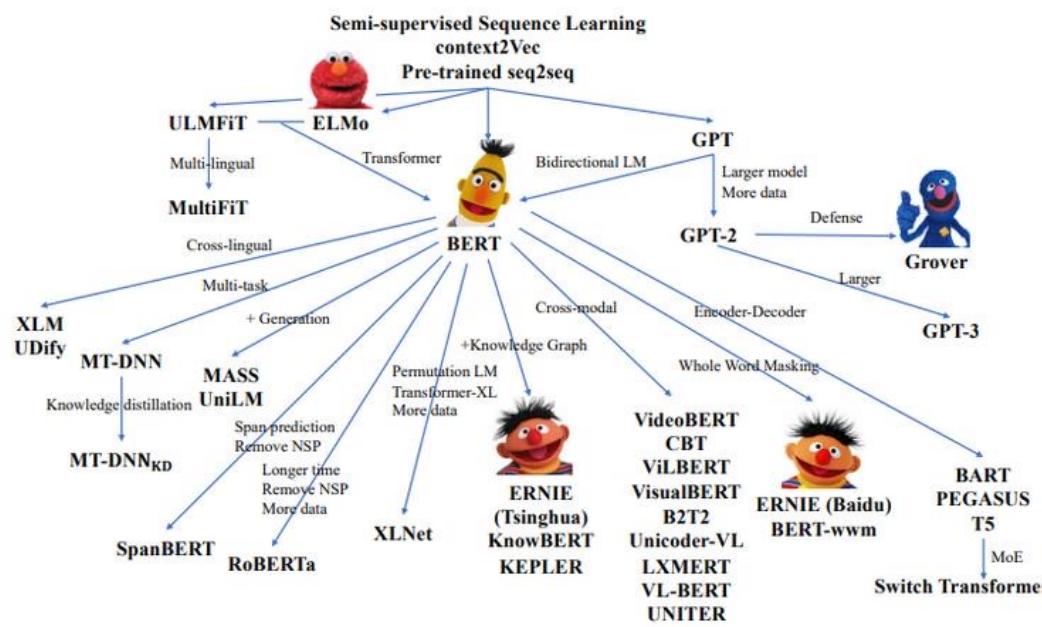


# **Knowledge Inheritance**

Ding Xuanwen

# Motivation

Despite the great follow-up efforts of exploring various pre-training techniques and model architectures, researchers find that simply enlarging the **model capacity, data size and training steps** can further improve the performance of PLMs



# Motivation

- However, training large-scale PLMs requires tremendous computational resources, raising severe environmental concerns on the prohibitive computational costs
- Existing PLMs are generally trained from scratch individually, ignoring that many well-trained PLMs are available

	Organization	Model	Parameter	Data	Time
2018.6	OpenAI	GPT	110M	4GB	3 Day
2018.10	Google	BERT	330M	16GB	50 Day
2019.2	OpenAI	GPT-2	1.5B	40GB	200 Day
2019.7	Facebook	RoBERTa	330M	160GB	3 Year
2019.10	Google	T5	11B	800GB	66 Year
2020.6	OpenAI	GPT-3	175B	560GB	90 Year

**How could existing PLMs benefit training larger PLMs in future?**

# Knowledge Inheritance

- **Knowledge Inheritance**
  - Through knowledge distillation [1]
  - Through parameter recycling [2][3]
  - ...

[1] Knowledge Inheritance for Pre-trained Language Models. Qin et al. NAACL 2022, oral

[2] bert2BERT: Towards Reusable Pretrained Lanruage Models. Chen et al. ACL 2022 poster

[3] ELLE: Efficient Lifelong Pre-training for Emerging Data. Qin et al. ACL 2022 findings

---

# **Knowledge Inheritance for Pre-trained Language Models**

---

**Yujia Qin<sup>12</sup>, Yankai Lin<sup>2</sup>, Jing Yi<sup>1</sup>, Jiajie Zhang<sup>1</sup>, Xu Han<sup>1</sup>, Zhengyan Zhang<sup>1</sup>,  
Yusheng Su<sup>1</sup>, Zhiyuan Liu<sup>1</sup>, Peng Li<sup>2</sup>, Maosong Sun<sup>1</sup>, Jie Zhou<sup>2</sup>**

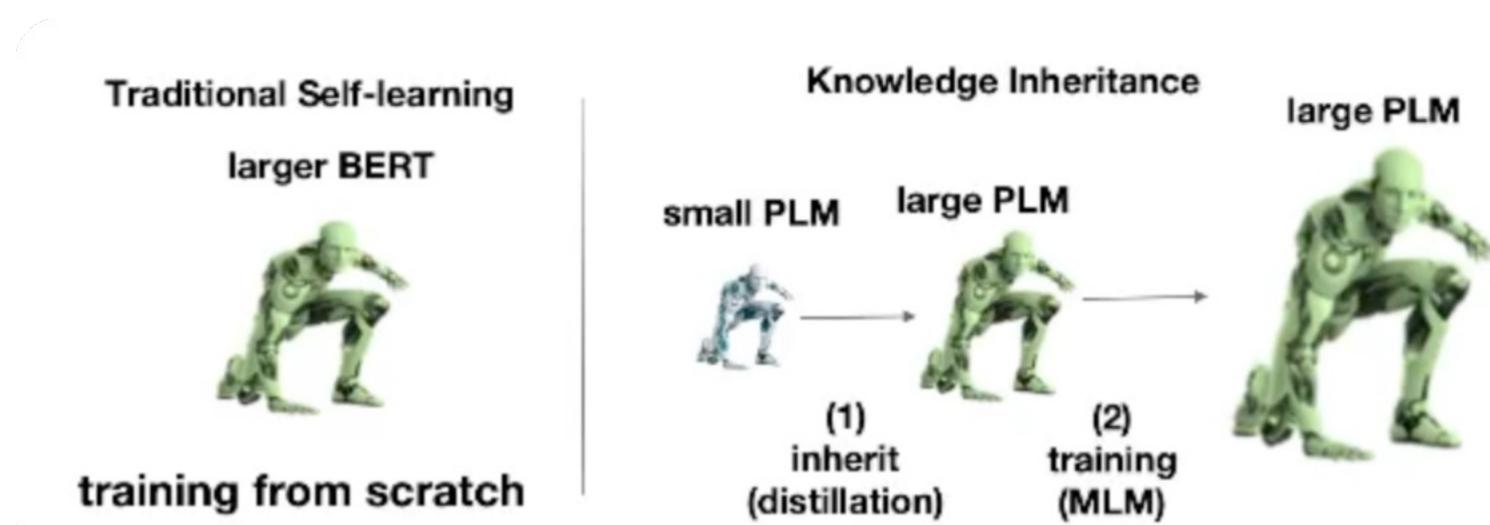
<sup>1</sup>Department of Computer Science and Technology, Tsinghua University

<sup>2</sup>Pattern Recognition Center, WeChat AI, Tencent Inc.

[yujiaqin16@gmail.com](mailto:yujiaqin16@gmail.com)

# Motivation

- Humans can leverage the knowledge summarized by their predecessors to learn new tasks
- Similarly, it is worth inheriting the implicit knowledge distributed in existing PLMs
- In this sense, we could **distill the knowledge** summarized by an existing small PLM during pretraining to efficiently learn larger PLMs
- We dub the above process as **knowledge inheritance (KI)**



# Method

- Knowledge Inheritance Framework
- Train a large PLM  $M_L$  by inheriting knowledge from  $M_S$

$$\begin{aligned}\mathcal{L}(\mathcal{D}_L; \mathcal{M}_S) &= \sum_{(\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{D}_L} (1 - \alpha) \mathcal{L}_{\text{SELF}}(\mathbf{x}_i, \mathbf{y}_i) + \alpha \mathcal{L}_{\text{KI}}(\mathbf{x}_i; \mathcal{M}_S) \\ &= \sum_{(\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{D}_L} (1 - \alpha) \mathcal{H}(\mathbf{y}_i, \mathcal{P}_{\mathcal{M}_L}(\mathbf{x}_i; 1)) + \alpha \tau^2 \text{KL}(\mathcal{P}_{\mathcal{M}_S}(\mathbf{x}_i; \tau) || \mathcal{P}_{\mathcal{M}_L}(\mathbf{x}_i; \tau)).\end{aligned}$$

- Dynamic Inheritance Rate

$$\begin{aligned}\mathcal{L}(\mathcal{D}_L; \mathcal{M}_S) &= \sum_{(\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{D}_L} (1 - \alpha_t) \mathcal{L}_{\text{SELF}}(\mathbf{x}_i, \mathbf{y}_i) + \alpha_t \mathcal{L}_{\text{KI}}(\mathbf{x}_i; \mathcal{M}_S). \\ \alpha_t &= \max(1 - \alpha_T \times \frac{t}{T}, 0)\end{aligned}$$

# Research Questions

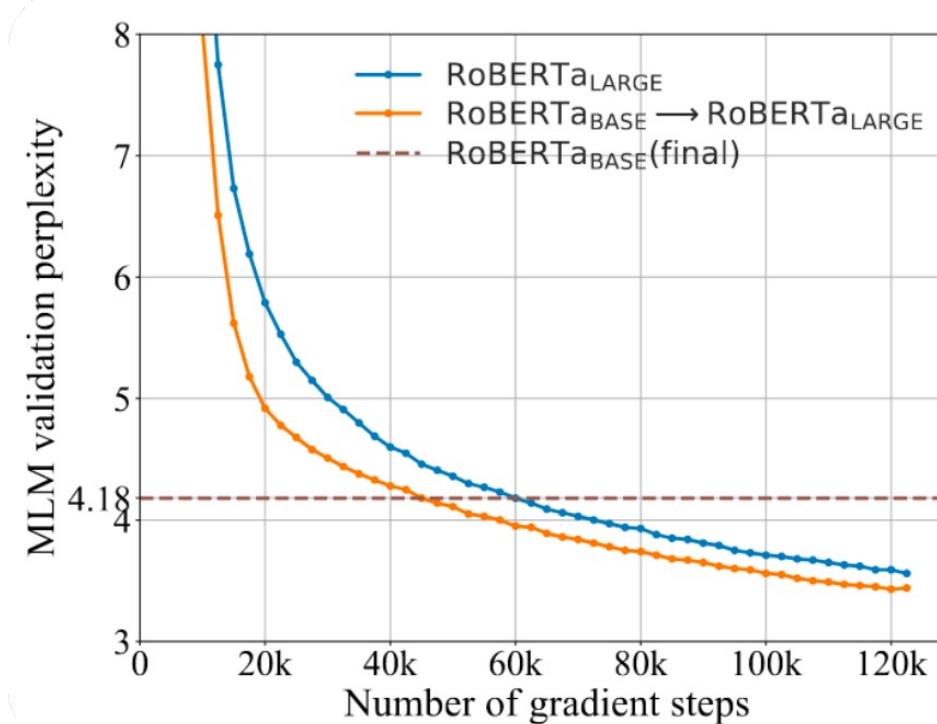
- (RQ1) Could **distilling knowledge** from an existing trained PLM benefit large PLMs' training from scratch?
- (RQ2) Considering human beings are able to hand down knowledge **from generation to generation**, could KI similarly be sequentially performed among a series of PLMs with growing sizes?
- (RQ3) As more and more PLMs with different **pre-training settings** emerge, how would different settings affect the performance of KI?
- (RQ4) Besides training a large PLM from scratch, when adapting an already trained large PLM to a **new domain**, how could smaller domain teachers benefit such a process?

# Experimental Setting

- Architecture: RoBERTa
- Teacher Model (BASE):  $RoBERTa_{BASE}$  Student Model (LARGE):  $ROBERTa_{LARGE}$
- Training Data: Wikipedia + BookCorpus (3400M tokens)
- Training Steps: 125k
- Batch Size: 2048
- Evaluation: Perplexity (pre-training), GLUE performance (fine-tuning)
- LARGE v.s. BASE → LARGE

# RQ1

- (RQ1) Could **distilling knowledge** from an existing trained PLM benefit large PLMs' training from scratch?
- Results of Pre-training:
  - Training **LARGE** under our KI framework converges faster than the self-learning baseline (27.3% FLOPs saved)



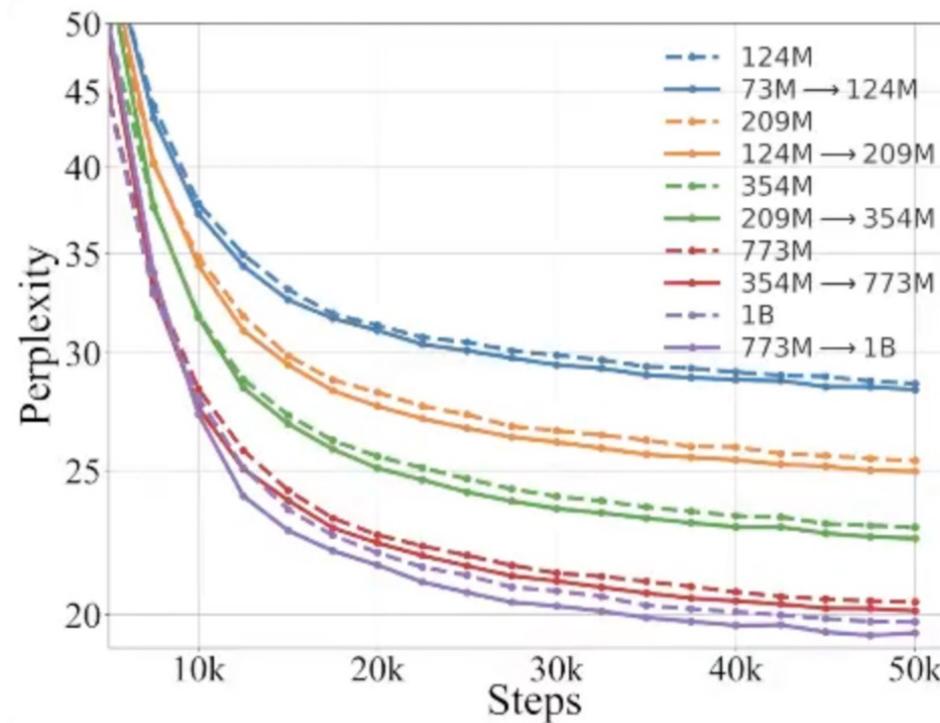
# RQ1

- (RQ1) Could **distilling knowledge** from an existing trained PLM benefit large PLMs' training from scratch?
- Results of Fine-training:
  - LARGE trained under KI framework achieves better performance than the baseline on downstream tasks at each step

Step	Model	CoLA	MNLI	QNLI	RTE	SST-2	STS-B	MRPC	QQP	Avg
5k	RoBERTa <sub>LARGE</sub>	0.0	73.5	81.7	53.0	81.7	45.8	71.4	87.5	61.8
	RoBERTa <sub>BASE</sub> → RoBERTa <sub>LARGE</sub>	<b>17.4</b>	<b>75.8</b>	<b>83.4</b>	<b>54.7</b>	<b>85.7</b>	<b>72.0</b>	<b>72.6</b>	<b>88.6</b>	<b>68.8</b>
45k	RoBERTa <sub>LARGE</sub>	61.8	84.9	91.7	63.4	92.9	88.6	87.7	<b>91.5</b>	82.8
	RoBERTa <sub>BASE</sub> → RoBERTa <sub>LARGE</sub>	<b>64.3</b>	<b>85.9</b>	<b>92.2</b>	<b>75.3</b>	<b>93.2</b>	<b>89.3</b>	<b>89.4</b>	91.5	<b>85.2</b>
85k	RoBERTa <sub>LARGE</sub>	64.5	86.8	92.7	69.7	93.5	89.9	89.7	91.7	84.8
	RoBERTa <sub>BASE</sub> → RoBERTa <sub>LARGE</sub>	<b>65.7</b>	<b>87.2</b>	<b>93.0</b>	<b>77.0</b>	<b>94.3</b>	<b>90.0</b>	<b>90.4</b>	<b>91.8</b>	<b>86.2</b>
125k	RoBERTa <sub>LARGE</sub>	64.3	87.1	<b>93.2</b>	73.4	94.1	90.3	<b>90.1</b>	91.8	85.5
	RoBERTa <sub>BASE</sub> → RoBERTa <sub>LARGE</sub>	<b>67.7</b>	<b>87.7</b>	93.1	<b>74.9</b>	<b>94.8</b>	<b>90.6</b>	88.2	<b>91.9</b>	<b>86.1</b>

# RQ1

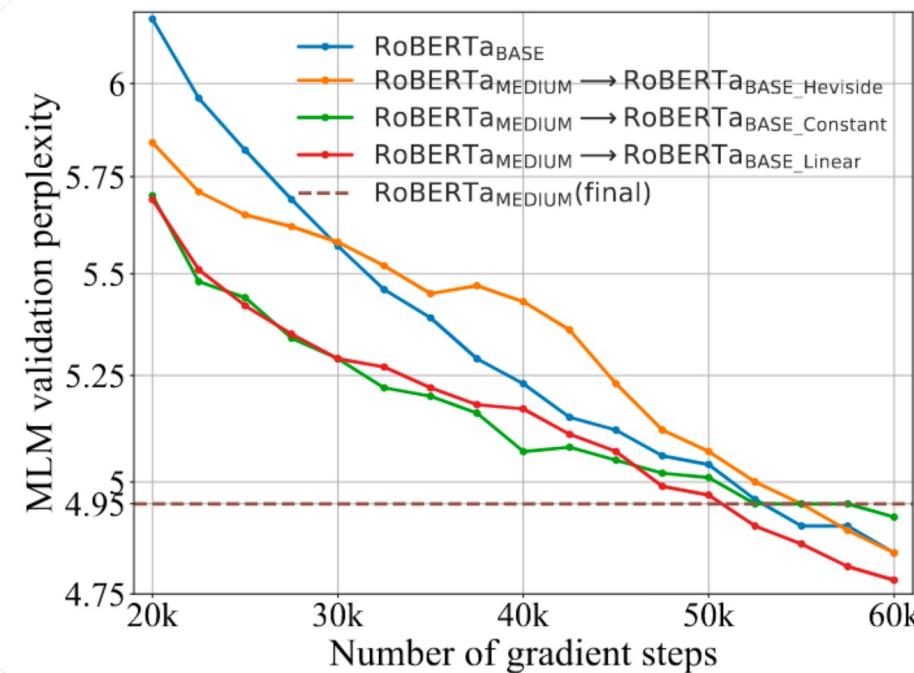
- (RQ1) Could **distilling knowledge** from an existing trained PLM benefit large PLMs' training from scratch?
- Results of Pre-training:
  - The conclusion also holds for GPT across various sizes



Model Name	$n_{\text{params}}$	$n_{\text{layers}}$	$d_{\text{model}}$	$n_{\text{heads}}$	$d_{\text{FFN}}$	lr (bs = 2,048)
GPT <sub>73M</sub>	73M	9	576	12	3072	$5.0 \times 10^{-4}$
GPT <sub>124M</sub>	124M	12	768	12	3072	$5.0 \times 10^{-4}$
GPT <sub>209M</sub>	209M	18	864	12	3600	$4.0 \times 10^{-4}$
GPT <sub>354M</sub>	354M	24	1024	16	4096	$3.5 \times 10^{-4}$
GPT <sub>773M</sub>	773M	36	1280	20	5120	$3.0 \times 10^{-4}$
GPT <sub>1B</sub>	1068M	40	1440	20	5760	$2.5 \times 10^{-4}$

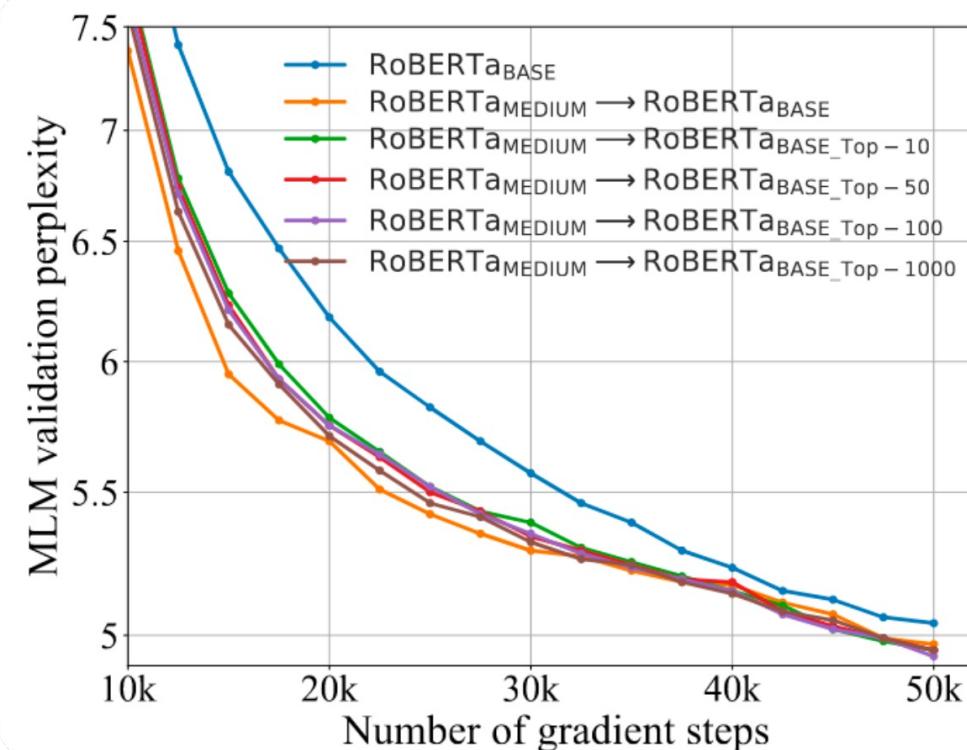
# RQ1

- (RQ1) Could **distilling knowledge** from an existing trained PLM benefit large PLMs' training from scratch?
- Effects of Inheritance Rate
  - Annealing at first is necessary
  - Teacher-guided learning is redundant after the student surpasses the teacher



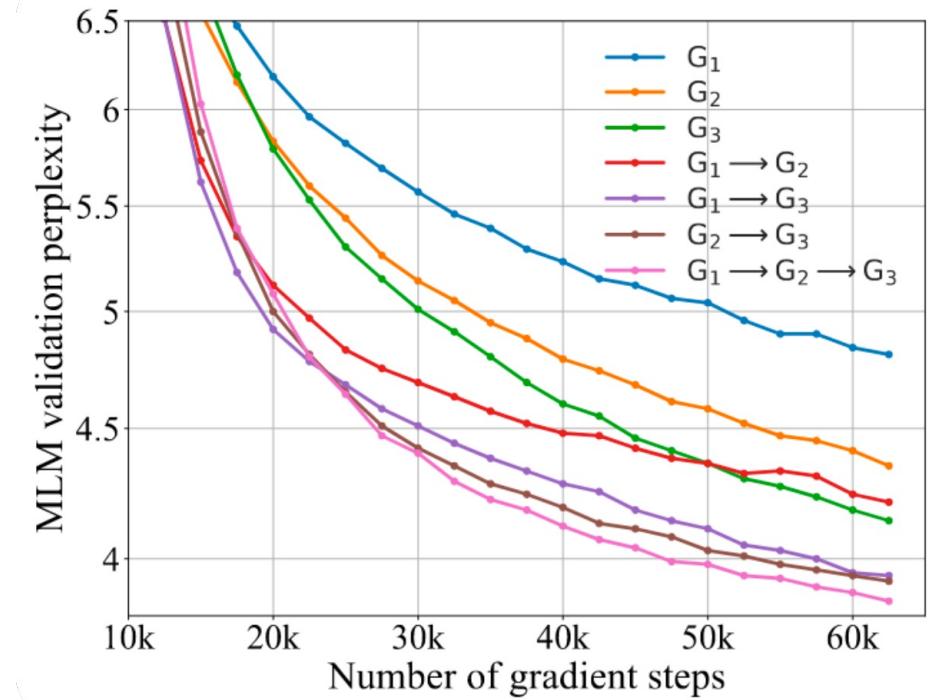
# RQ1

- (RQ1) Could **distilling knowledge** from an existing trained PLM benefit large PLMs' training from scratch?
- Saving Storage Space with Top-K Logits
  - Top-K logits contain the vast majority of information



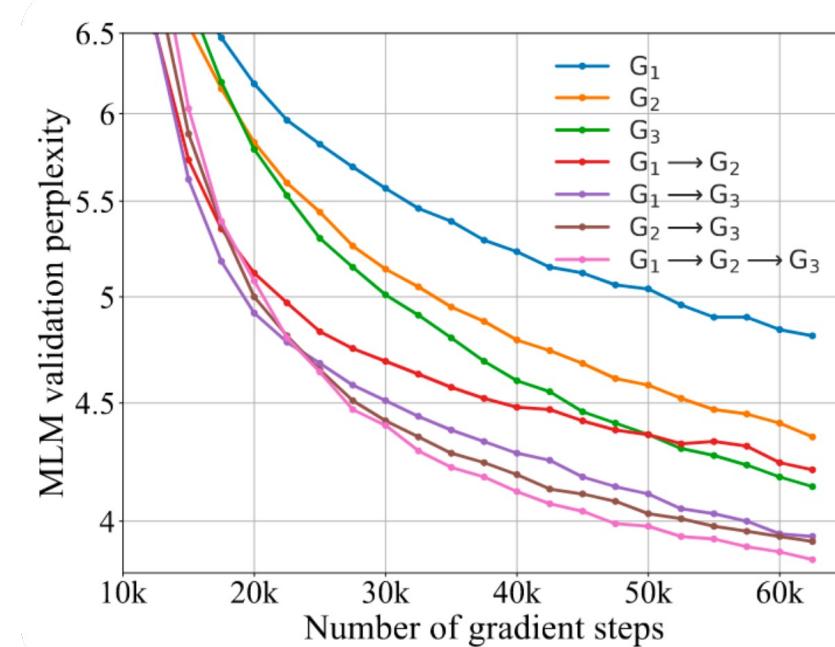
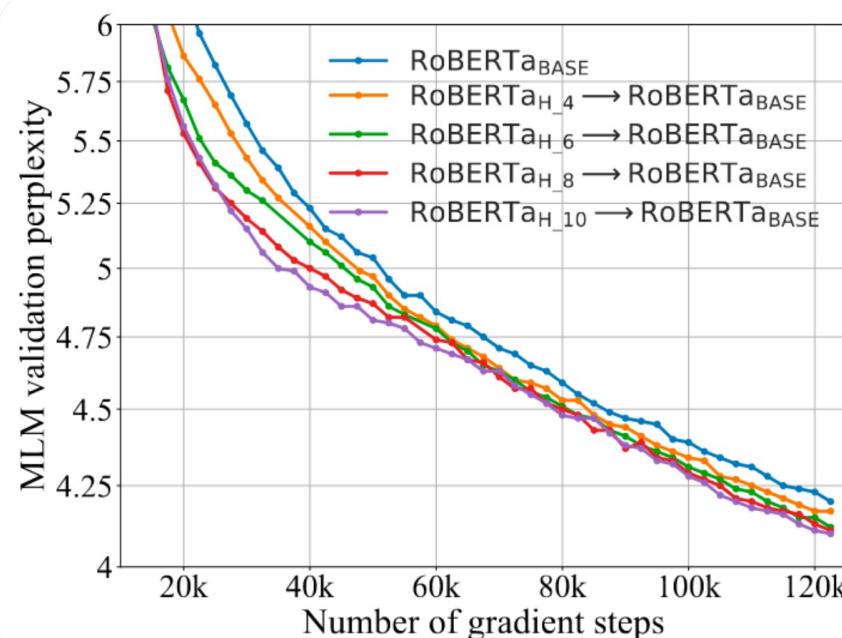
# RQ2

- (RQ2) Considering human beings are able to hand down knowledge **from generation to generation**, could KI similarly be sequentially performed among a series of PLMs with growing sizes?
- Knowledge Inheritance over Generations  
 $G_1$ (BASE, 125M),  $G_2$ (BASE\_PLUS, 211M)  $G_3$ (LARGE, 355M)
- (1) Self-learning ( $G_1, G_2, G_3$ )
- (2) KI over two generations ( $G_1 \rightarrow G_2, G_2 \rightarrow G_3, G_1 \rightarrow G_3$ )
- (3) KI over three generations ( $G_1 \rightarrow G_2 \rightarrow G_3$ )



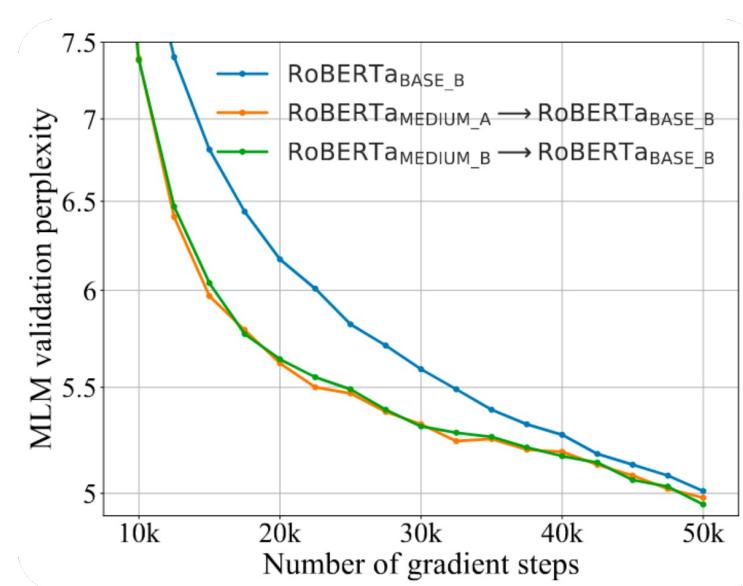
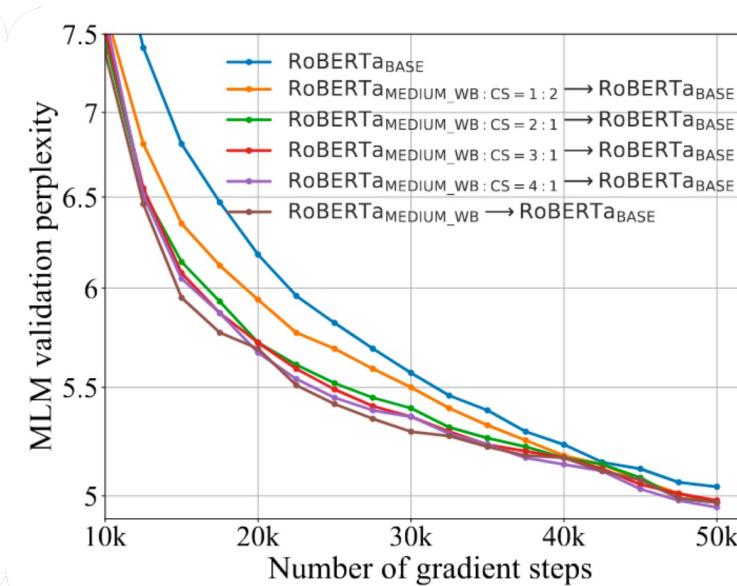
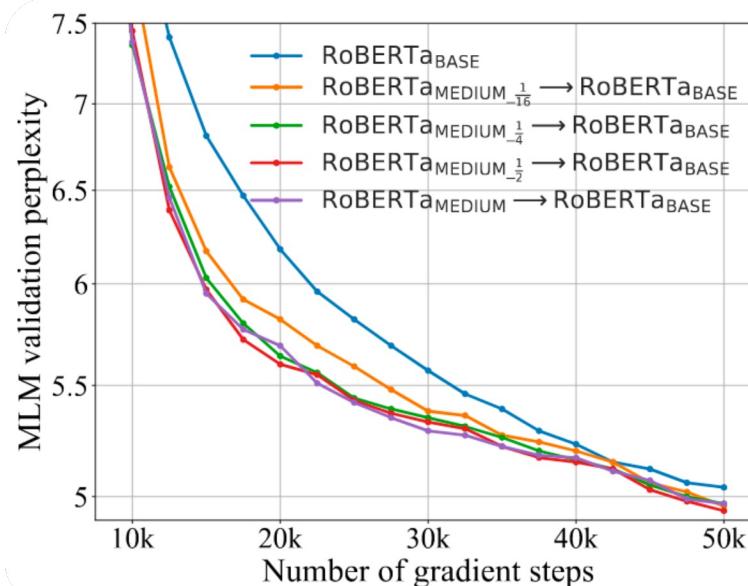
# RQ3

- (RQ3) As more and more PLMs with different **pre-training settings** emerge, how would different settings affect the performance of KI?
- Effects of Model Architecture
  - Choosing a wider / deeper teacher, which has lower validation PPL, further accelerates ML's convergence



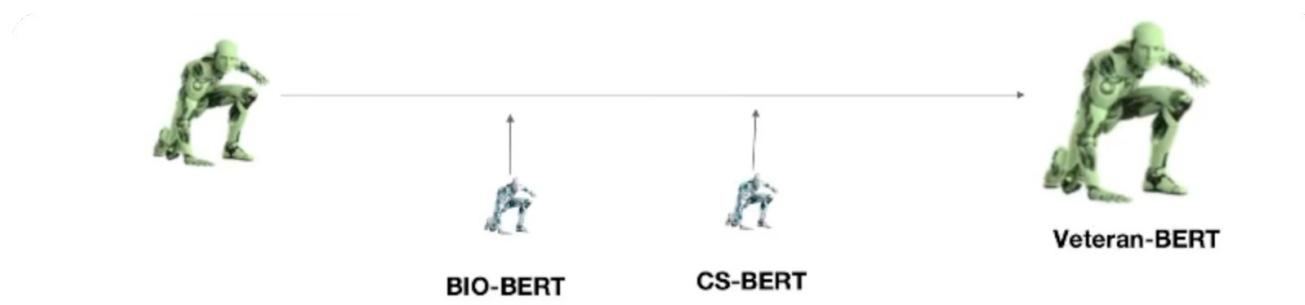
# RQ3

- (RQ3) As more and more PLMs with different **pre-training settings** emerge, how would different settings affect the performance of KI?
- Effects of Data Size / Domain / Privacy
  - PLMs can image the big from the small for in-domain data
  - Inheriting on similar domain improves performance
  - Data privacy is not important if the same domain is ensured



# RQ4

- (RQ4) Besides training a large PLM from scratch, when adapting an already trained large PLM to a **new domain**, how could smaller domain teachers benefit such a process?
- KI is more training-efficient and data-efficient
- Catastrophic Forgetting on the Source Domain



<b>N<sub>tokens</sub></b>	3,400M		200M		100M		40M		20M		
<b>Metrics</b>	F1	PPL	F1	PPL	F1	PPL	F1	PPL	F1	PPL	
CS	SL	69.8	3.12	71.7	3.17	71.4	3.24	68.3	3.51	67.5	4.07
	KI	<b>72.9</b>	<b>3.06</b>	<b>72.6</b>	<b>3.09</b>	<b>71.9</b>	<b>3.11</b>	<b>71.1</b>	<b>3.21</b>	<b>70.8</b>	<b>3.37</b>
BIO	SL	84.0	2.67	82.8	2.72	83.2	2.83	83.3	3.16	82.7	3.81
	KI	<b>84.5</b>	<b>2.65</b>	<b>83.4</b>	<b>2.66</b>	<b>83.9</b>	<b>2.69</b>	<b>83.6</b>	<b>2.82</b>	<b>83.5</b>	<b>3.01</b>

<b>Domain</b>	<b>Strategy</b>	3,400M	200M	100M	40M
CS	SL	6.71	7.01	7.39	8.77
	KI	8.63	9.39	9.48	9.87
BIO	SL	7.29	6.61	8.16	10.34
	KI	10.74	10.78	10.93	11.66

# Conclusion

- Propose a general knowledge inheritance framework that leverages previously trained PLMs for training larger ones
- Conduct sufficient empirical studies to demonstrate its feasibility in knowledge transfer, domain adaptation, etc.
- Comprehensively analyze various pre-training settings of the teacher model

# **bert2BERT: Towards Reusable Pretrained Language Models**

**Cheng Chen<sup>1†</sup>, Yichun Yin<sup>2</sup>, Lifeng Shang<sup>2</sup>, Xin Jiang<sup>2</sup>, Yujia Qin<sup>1</sup>,  
Fengyu Wang<sup>1</sup>, Zhi Wang<sup>3,4‡</sup>, Xiao Chen<sup>2</sup>, Zhiyuan Liu<sup>1</sup>, Qun Liu<sup>2</sup>**

<sup>1</sup>Department of Computer Science and Technology, Tsinghua University

<sup>2</sup>Huawei Noah's Ark Lab, <sup>3</sup>Tsinghua Shenzhen International Graduate School

<sup>4</sup>Peng Cheng Laboratory

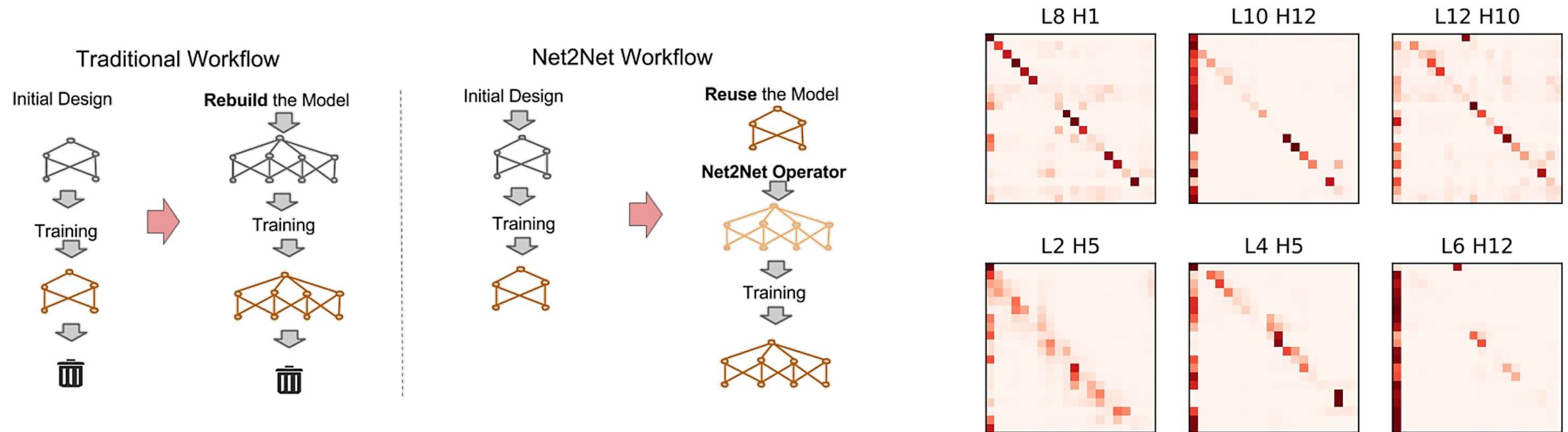
{c-chen19, qyj20, wangfy20}@mails.tsinghua.edu.cn

{yinyichun, shang.lifeng, jiang.xin, chen.xiao2, qun.liu}@huawei.com

wangzhi@sz.tsinghua.edu.cn, liuzy@tsinghua.edu.cn

# Motivation

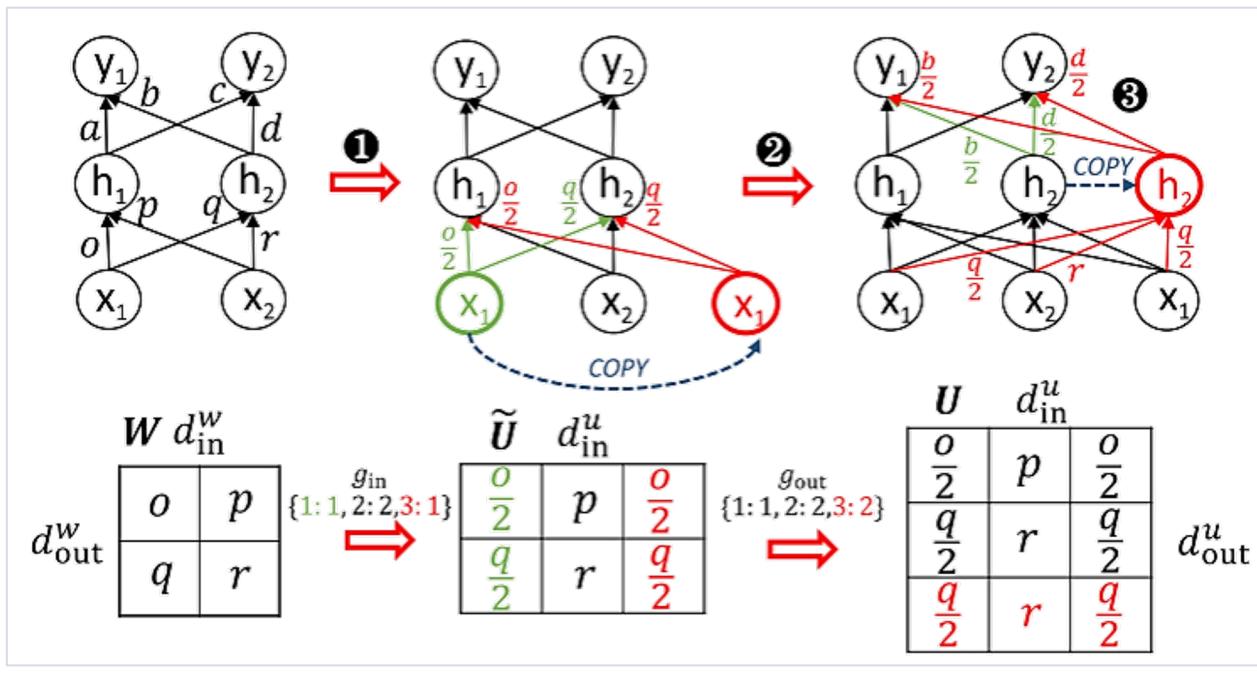
- Related work
  - Net2Net[1] : use **function-preserving transformation** to make neural networks reusable.
  - StackingBert[2] : different layers have some similar knowledge



# Method

- Objective :  $S(L^s, D^s) \rightarrow T(L^t, D^t)$
- Parameter initialization
  - function-preserving initialization (FPI), ensures that the initialized large model has almost the same behavior as the small model
  - advanced knowledge initialization (AKI), leads to a faster convergence rate and achieves higher training efficiency
- Two-stage training strategy
  - train sub-models with different layers
  - full-model training

# FPI



$$g_{\text{in}}(i) = \begin{cases} i & i \in [1, d_{\text{in}}^w] \\ f(\{1, 2, \dots, d_{\text{in}}^w\}) & i \in (d_{\text{in}}^w, d_{\text{in}}^u], \end{cases} \quad (5)$$

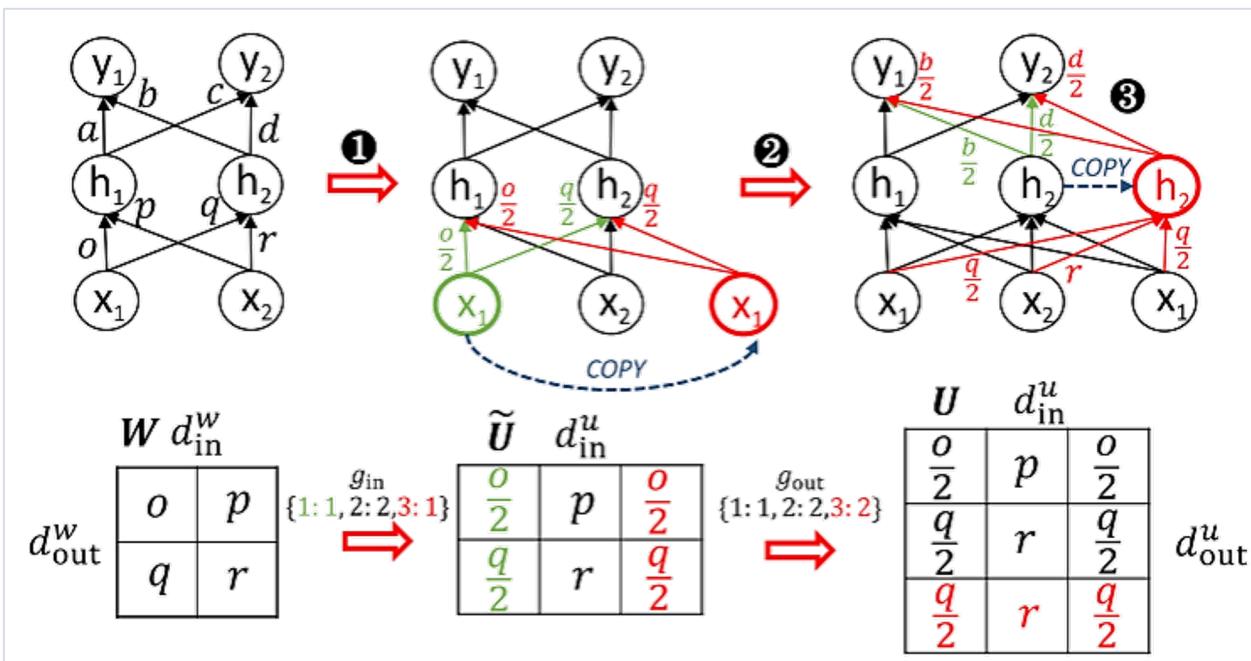
$$g_{\text{out}}(j) = \begin{cases} j & j \in [1, d_{\text{out}}^w] \\ f(\{1, 2, \dots, d_{\text{out}}^w\}) & j \in (d_{\text{out}}^w, d_{\text{out}}^u], \end{cases} \quad (6)$$

$$C_{g_{\text{in}}(i)} = \sum_{i'=1}^{d_{\text{in}}^u} \mathbb{I}(g_{\text{in}}(i') = g_{\text{in}}(i)) \quad (7)$$

$$\tilde{\mathbf{U}}_{(i,*)} = \frac{1}{C_{g_{\text{in}}(i)}} \mathbf{W}_{(g_{\text{in}}(i),*)},$$

$$\mathbf{U}_{(*,j)} = \tilde{\mathbf{U}}_{(*,g_{\text{out}}(j))}, \quad (8)$$

# FPI



Embedding  $U_{(*,j)}^E = \mathbf{W}_{(*,g_{\text{out}}^e(j))}^E$ . (9)

MHA  $\mathbf{U}^{Q|K|V} = \text{EXPN}(\mathbf{W}^{Q|K|V}; g_{\text{in}}^{q|k|v}, g_{\text{out}}^{q|k|v}),$   
 $\mathbf{U}^O = \text{EXPN}(\mathbf{W}^O; g_{\text{in}}^o, g_{\text{out}}^o).$

$$g_{\text{out}}^{q|k|v}(j) = \begin{cases} j & j \in [1, a^s] \\ f(\{1, 2, \dots, a^s\}) & j \in (a^s, a^t], \end{cases} \quad (10)$$

$$g_{\text{out}}^e = g_{\text{in}}^{q|k|v}; g_{\text{out}}^{q|k|v} = g_{\text{in}}^o; g_{\text{in}}^{q|k|v} = g_{\text{out}}^o \quad (11)$$

$$g_{\text{out}}^e = g_{\text{in}}^{q|k|v}; g_{\text{out}}^{q|k|v} = g_{\text{in}}^o; g_{\text{in}}^{q|k|v} = g_{\text{out}}^o$$

FFN

$$\mathbf{U}^1 = \text{EXPN}(\mathbf{W}^1; g_{\text{in}}^1, g_{\text{out}}^1), \quad (12)$$

$$\mathbf{U}^2 = \text{EXPN}(\mathbf{W}^2; g_{\text{in}}^2, g_{\text{out}}^2).$$

$$g_{\text{out}}^o = g_{\text{in}}^1; g_{\text{out}}^1 = g_{\text{in}}^2; g_{\text{in}}^1 = g_{\text{out}}^2$$

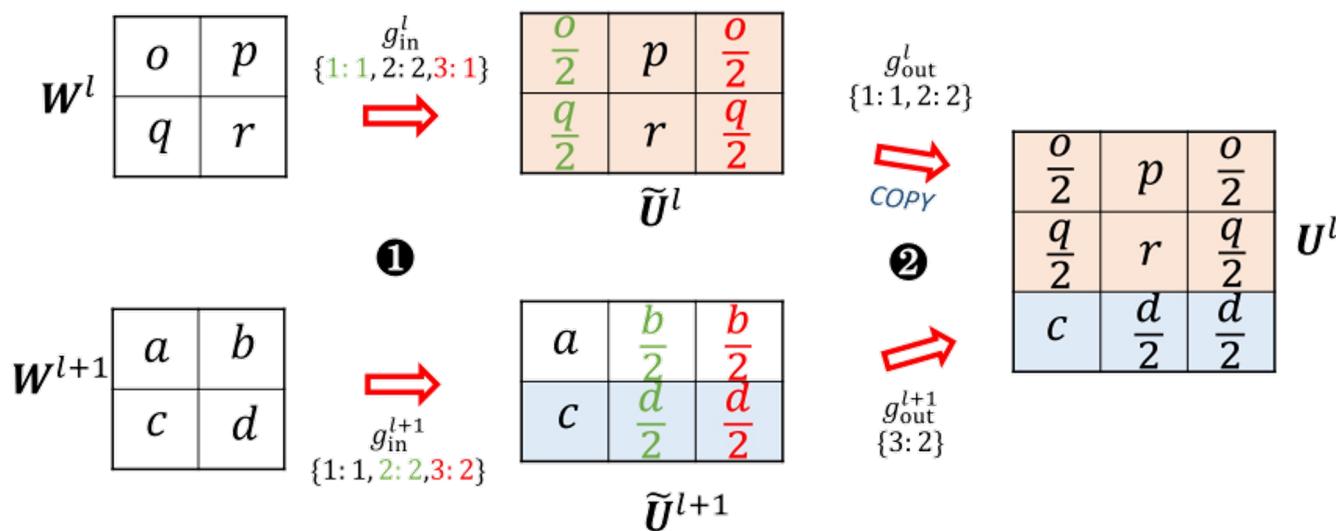
LN

$$\mathbf{U}_j^{\text{ln}} = \mathbf{W}_{g_{\text{out}}^2(j)}^{\text{ln}}. \quad (13)$$

LayerNorm( $\mathbf{H}$ ) =  $(\frac{\mathbf{H} - \mu_{\mathbf{H}}}{\sigma_{\mathbf{H}}}) \odot \mathbf{W}^{LN} + \mathbf{b}^{LN}$

# AKI

- Breaks FPI **symmetry** that hinders model convergence
- Upper-layer information can be used as **high-level** knowledge



$$\mathbf{U}^l = \text{EXPN}(\mathbf{W}^l, \mathbf{W}^{l+1}; g_{\text{in}}, g_{\text{out}}). \quad (14)$$

$$C_{g_{\text{in}}^l(i)} = \sum_{i'=1}^{d_{\text{in}}^u} \mathbb{I}(g_{\text{in}}^l(i') = g_{\text{in}}^l(i)) \quad (15)$$

$$\tilde{\mathbf{U}}_{(i,*)}^l = \frac{1}{C_{g_{\text{in}}^l(i)}} \mathbf{W}_{(g_{\text{in}}^l(i),*)}^l.$$

$$\mathbf{U}_{(*,j)}^l = \begin{cases} \tilde{\mathbf{U}}_{(*,j)}^l & j \in [1, d_{\text{out}}^s] \\ \tilde{\mathbf{U}}_{(*,g_{\text{out}}^{l+1}(j))}^{l+1} & j \in (d_{\text{out}}^s, d_{\text{out}}^t], \end{cases} \quad (16)$$

# Depth-wise expansion

---

**Algorithm 1** bert2BERT initialization

---

**Input:** the target model  $\mathcal{T}(L^t, D^t)$  and the source model  $\mathcal{S}(L^s, D^s)$ .

- 1:  $\mathcal{T}_1(L^s, D^t) \leftarrow$  do AKI or FPI with  $\mathcal{S}(L^s, D^s)$
- 2:  $k \leftarrow \lfloor L^t / L^s \rfloor$
- 3: **for**  $t = 2 \rightarrow k$  **do**
- 4:      $\mathcal{T}_t(L^s \cdot t, D^t) \leftarrow$  stack  $\mathcal{T}_1$  on top of  $\mathcal{T}_{t-1}$
- 5: **end for**
- 6:  $\mathcal{T} \leftarrow$  stack top  $L^t - L^s \cdot k$  layers of  $\mathcal{T}_1$ .

**Output:** the initialized model  $\mathcal{T}(L^t, D^t)$

---

# Two-stage Pre-training

- Stage1: Train sub-modules
- Stage2: Train full-model

---

**Algorithm 2** Two-stage Pre-training

---

**Input:** the initialized model  $\mathcal{T}$ , large-scale unsupervised dataset  $\mathcal{D}$ , the epoch number of sub-model training  $E_b$  and the epoch number of whole training process  $E$ , the layer number  $l_b$ .

```
1: Construct sub-models and these models have  
   the layer numbers of  $\{l_b, 2 \cdot l_b, \dots, L^t\}$ .  
2: for  $e = 1 \rightarrow E_b$  do  
3:   for batch in  $D$  do  
4:      $\mathcal{T}' \leftarrow$  sample one sub-model.  
5:     Perform forward and backward of  $\mathcal{T}'$ .  
6:     Update only top  $l_b$  layers of  $\mathcal{T}'$ .  
7:   end for  
8: end for
```

---

---

**Algorithm 2** Two-stage Pre-training

---

**Input:** the initialized model  $\mathcal{T}$ , large-scale unsupervised dataset  $\mathcal{D}$ , the epoch number of sub-model training  $E_b$  and the epoch number of whole training process  $E$ , the layer number  $l_b$ .

```
9: for  $e = E_b \rightarrow E$  do  
10:   for batch in  $D$  do  
11:     Perform forward and backward of  $\mathcal{T}$ .  
12:     Update whole model  $\mathcal{T}$ .  
13:   end for  
14: end for
```

**Output:** the pre-trained model  $\mathcal{T}$

---

# Experiment Setting

- Pre-training Details
  - Datasets : English Wikipedia , Toronto Book
  - Two-stage :  $E_b=5$ ,  $l_b=3$
  - Target model :  $T(12, 768)$
  - Source model :  $S(12, 512)$
- Fine-tuning Details
  - Downstream task : GLUE benchmark, SQuADv1.1
- Baseline
  - DirectCopy, Stack-BERT, MSLT

# Result

Model	FLOPs ( $\times 10^{19}$ )	Ratio (Cost Saving)	SQuADv1.1 (F1)	SST-2 (Acc)	MNLI (Acc)	MRPC (Acc)	CoLA (Mcc)	QNLI (Acc)	QQP (Acc)	STS-B (Acc)	Avg.
BERT <sub>BASE</sub> (Google)	-	-	88.4(0.1)	93.6(0.2)	84.7(0.1)	87.9(0.9)	59.6(1.5)	91.6(0.1)	91.4(0.1)	89.6(0.5)	85.8(0.1)
BERT <sub>BASE</sub> † (Ours)	7.3	0%	89.6(0.1)	92.7(0.2)	84.6(0.2)	88.6(0.5)	57.3(4.0)	90.6(0.7)	90.6(0.1)	89.9(0.3)	85.5(0.5)
<i>Progressive Training</i>											
MSLT†	6.5	10.7%	90.4(0.2)	92.9(0.2)	85.1(0.2)	87.9(2.1)	55.6(4.1)	90.7(0.2)	90.6(0.2)	88.2(0.6)	85.2(0.7)
StackBERT†	5.5	24.3%	90.4(0.2)	92.6(0.4)	85.3(0.1)	88.2(1.0)	63.2(0.9)	91.0(0.4)	91.0(0.1)	86.7(0.7)	86.0(0.2)
<i>bert2BERT : <math>\mathcal{S}(12, 512) \rightarrow \mathcal{T}(12, 768)</math></i>											
DirectCopy	6.4	12.2%	89.8(0.2)	92.9(0.3)	84.7(0.2)	86.2(0.6)	62.2(0.7)	90.2(0.6)	90.4(0.1)	89.2(0.1)	85.7(0.1)
<b>FPI</b>	5.1	30.4%	90.0(0.2)	92.6(0.4)	85.2(0.1)	87.1(0.5)	61.5(0.9)	90.9(0.6)	90.8(0.2)	89.7(0.2)	86.0(0.1)
<b>AKI</b>	4.5	38.4%	90.4(0.1)	92.5(0.4)	85.3(0.4)	87.8(0.9)	61.0(1.4)	91.2(0.2)	90.5(0.1)	89.5(0.2)	86.0(0.2)
<b>bert2BERT</b>	<b>4.0</b>	<b>45.2%</b>	90.0(0.2)	92.9(0.1)	85.1(0.1)	87.7(0.7)	60.0(1.2)	90.5(0.8)	90.4(0.1)	89.2(0.2)	85.7(0.4)

# Result

Settings	Model	FLOPs ( $\times 10^{19}$ )	Ratio (Saving)	Loss (MLM)	Avg.
$S(6, 512)$	DirectCopy	7.3	0%	1.440	89.1
	bert2BERT	5.6	23.3%	1.435	89.3
$S(8, 512)$	bert2BERT	4.6	36.8%	1.435	89.2
$S(10, 512)$	bert2BERT	4.2	42.7%	1.434	89.1

Table 3: bert2BERT with smaller source model. Avg means the average score of SST-2/MNLI/SQuADv1.1.

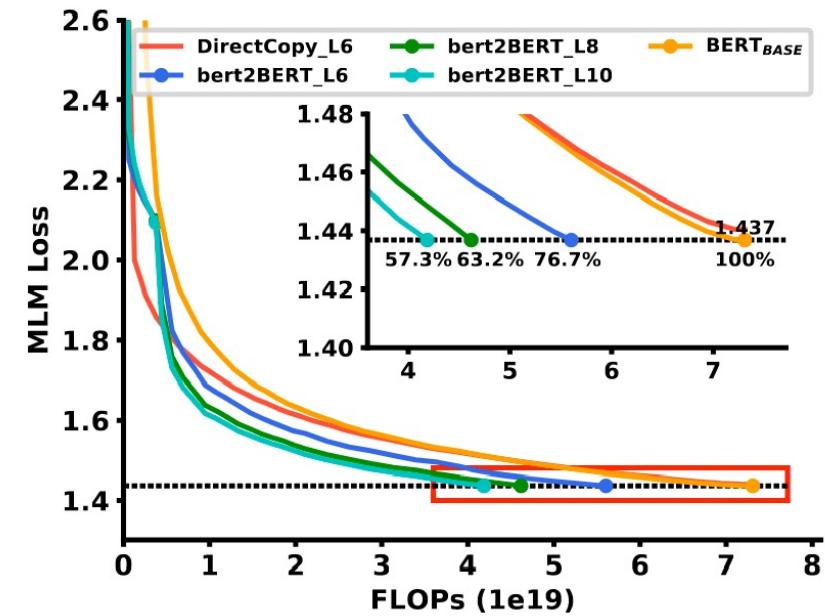
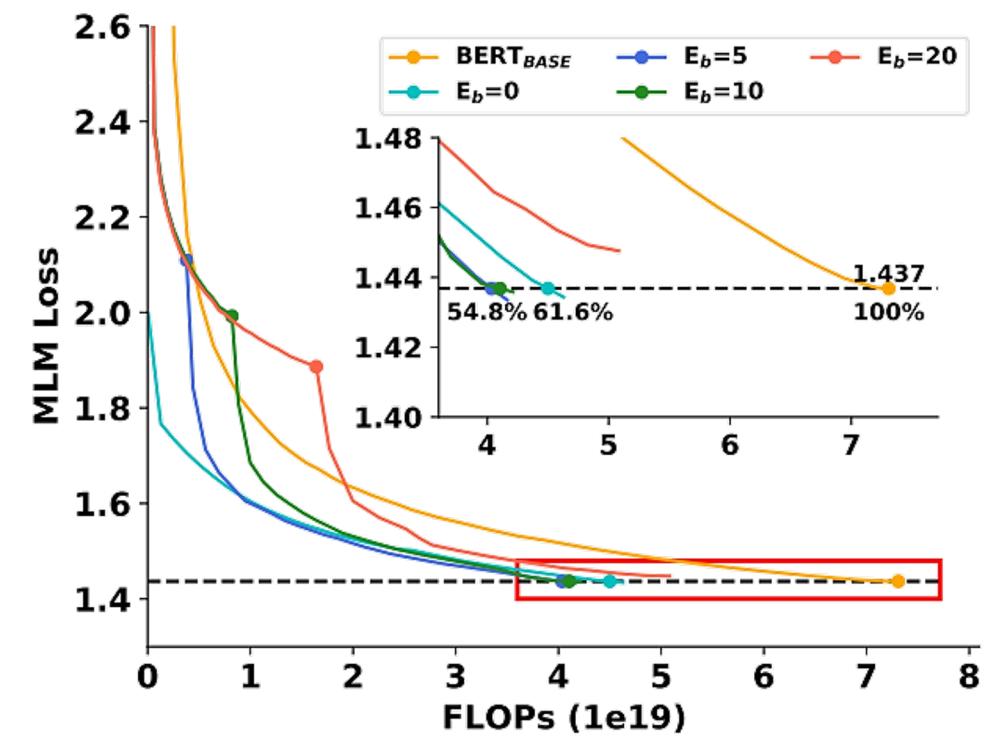


Figure 6: Loss curves of bert2BERT and baselines with smaller source models.

# Result

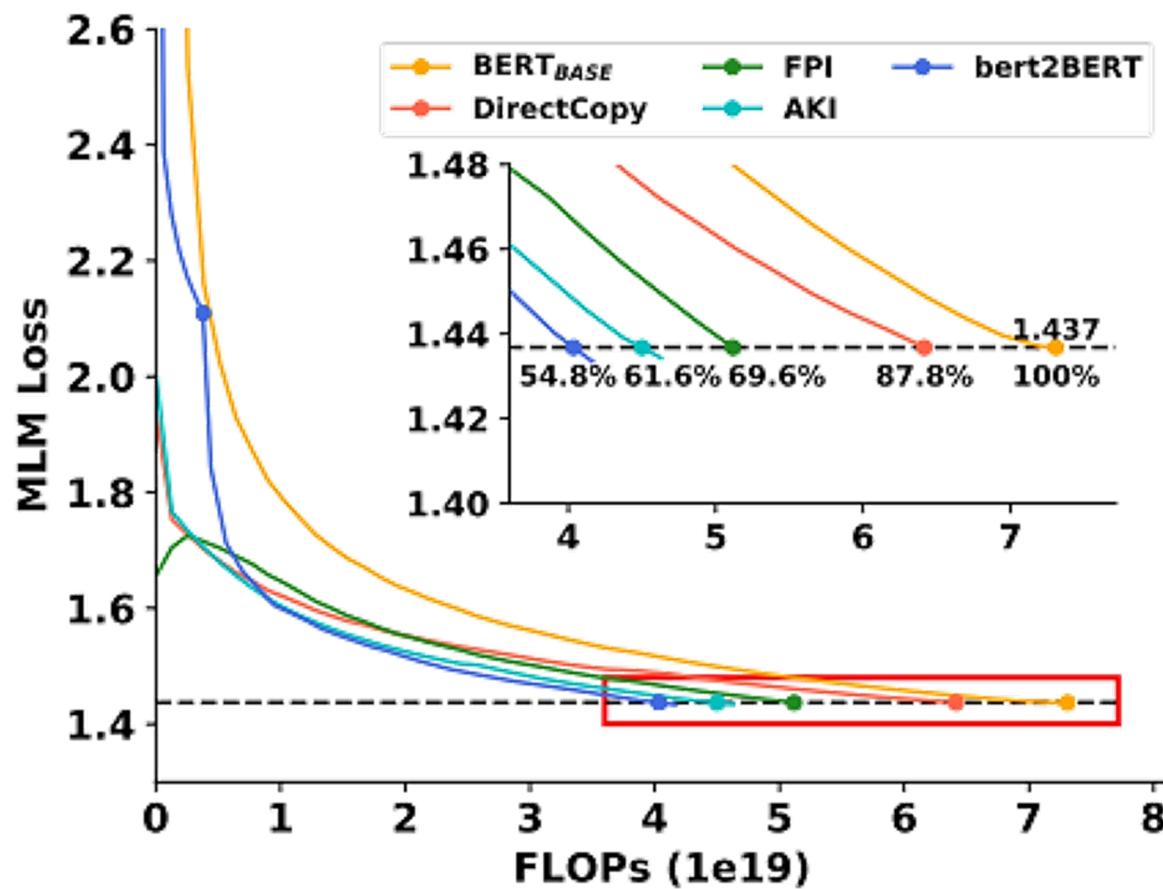
Model	FLOPs ( $\times 1e19$ )	Ratio (Cost Saving)	Avg.
<i>bert2BERT : <math>S(12, 512) \rightarrow T(12, 768)</math></i>			
bert2BERT ( $E_b = 0$ )	4.5	38.4%	86.0
bert2BERT ( $E_b = 5$ )	<b>4.0</b>	<b>45.2%</b>	85.7
bert2BERT ( $E_b = 10$ )	4.1	43.9%	85.3
bert2BERT ( $E_b = 20$ )	5.4	25.4%	84.3



# Result

Model	FLOPs ( $\times 1e19$ )	PTB (w/o FT)	WikiText-2 (w/o FT)	WikiText103 (w/o FT)
$bert2BERT : \mathcal{S}(12, 512) \rightarrow \mathcal{T}(12, 768)$				
GPT	4.9	133.8	47.0	53.5
<b>bert2BERT</b>	<b>2.6</b> (47% ↓)	132.1	47.9	53.0

# Result



# ELLE: Efficient Lifelong Pre-training for Emerging Data

**Yujia Qin**<sup>1,2,3\*</sup>, **Jiajie Zhang**<sup>1,2,3\*</sup>, **Yankai Lin**<sup>4</sup>, **Zhiyuan Liu**<sup>1,2,3,5,6†</sup>, **Peng Li**<sup>7‡</sup>,  
**Maosong Sun**<sup>1,2,3,5,6,8†</sup>, **Jie Zhou**<sup>4</sup>

<sup>1</sup>Department of Computer Science and Technology, Tsinghua University, Beijing, China

<sup>2</sup>Beijing National Research Center for Information Science and Technology

<sup>3</sup>Institute for Artificial Intelligence, Tsinghua University, Beijing, China

<sup>4</sup>Pattern Recognition Center, WeChat AI, Tencent Inc.

<sup>5</sup>International Innovation Center of Tsinghua University, Shanghai, China

<sup>6</sup>Beijing Academy of Artificial Intelligence

<sup>7</sup>Institute for AI Industry Research (AIR), Tsinghua University, China.

<sup>8</sup>Jiangsu Collaborative Innovation Center for Language Ability, Xuzhou, China

{qyj20, jiajie-z19}@mails.tsinghua.edu.cn

# Motivation

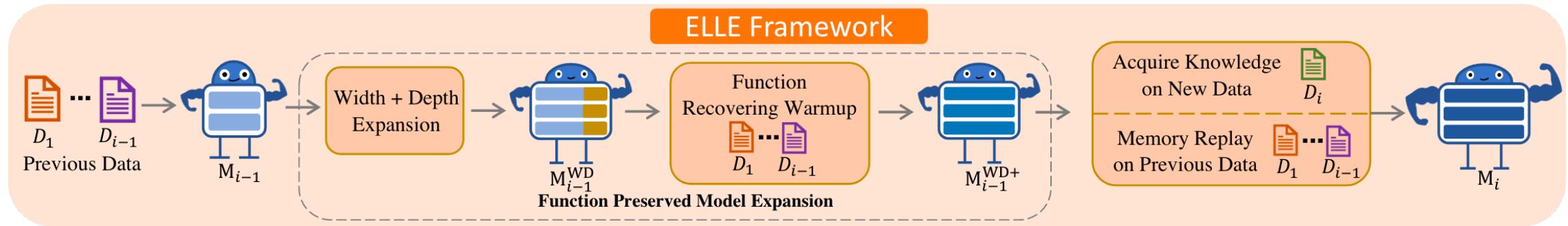
- Question: with limited computational resources, **how can we efficiently adapt PLMs in a lifelong manner?**
- **Efficient Lifelong Pre-training:** PLMs continually absorb knowledge from emerging data, and in the meantime, mitigate the catastrophic forgetting



# Motivation

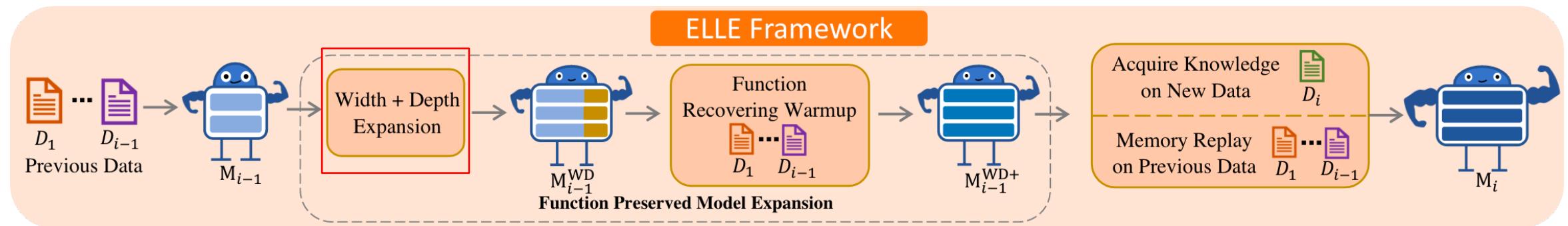
- Two challenges for efficient lifelong pre-training
  - (1) **Efficient knowledge growth**: with the data scale growing, packing more knowledge into a fixed-sized PLM becomes increasingly hard
  - (2) **Proper knowledge stimulation**: during pre-training, various knowledge from all domains is packed into PLMs hastily. However, a certain downstream task may largely require the knowledge from a specific domain
- We propose ELLE, targeting at **Efficient LifeLong** pre-training for **Emerging** data

# Methodology



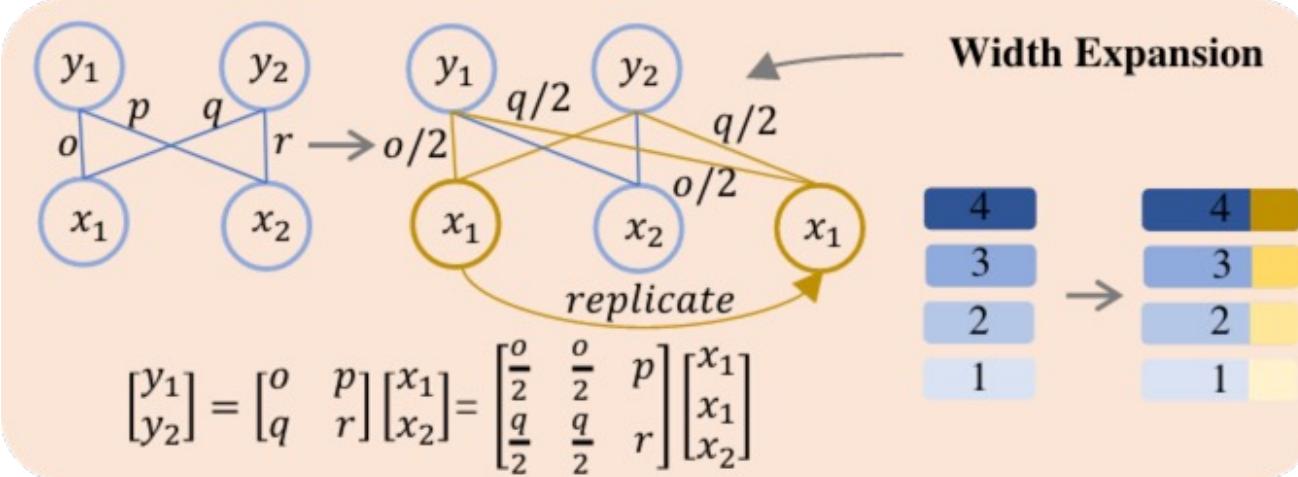
# Methodology

- (1) **Function preserved model expansion:** flexibly expand an existing PLM's width and depth to improve the efficiency of knowledge acquisition



# Methodology

- Width Expansion



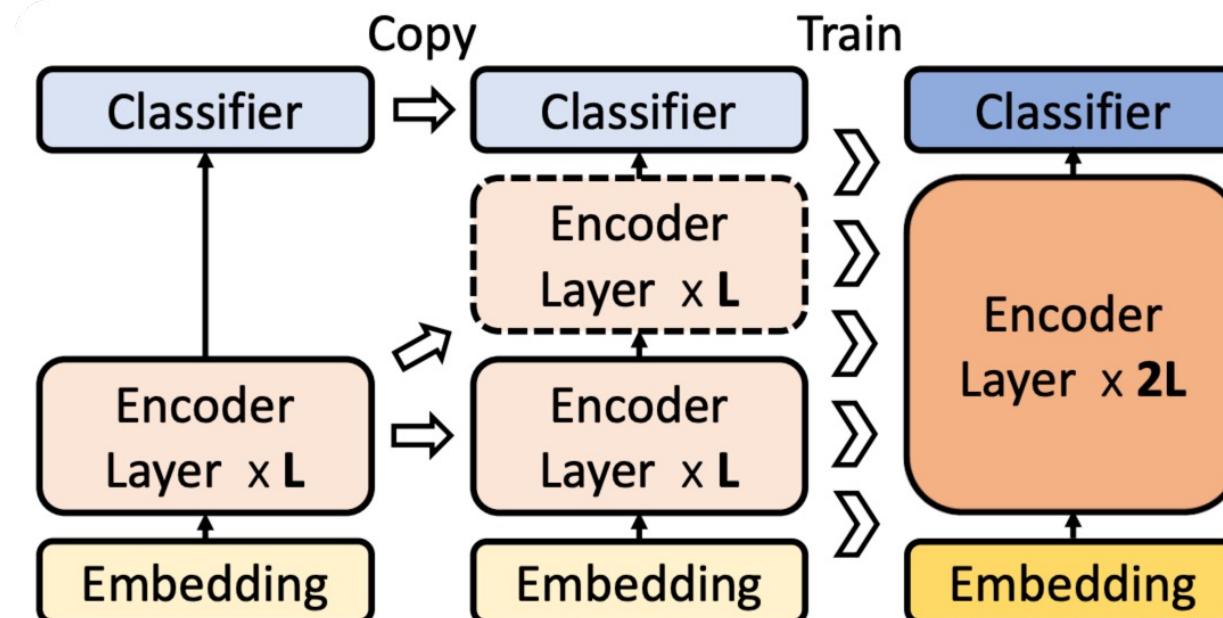
$$m(i) = \begin{cases} i & i \in [1, h_1] \\ U(\{1, \dots, h_1\}) & i \in (h_1, h_1 + \Delta_{h_1}], \end{cases}$$

$$C_i = \sum_{i'=1}^{h_1 + \Delta_{h_1}} \mathbb{I}(m(i') = m(i)),$$

$$\mathbf{W}'_{(i,*)} = \frac{1}{C_i} \cdot \mathbf{W}_{(m(i),*)} + \mathbb{I}(C_i > 1) \cdot \boldsymbol{\delta}_i,$$

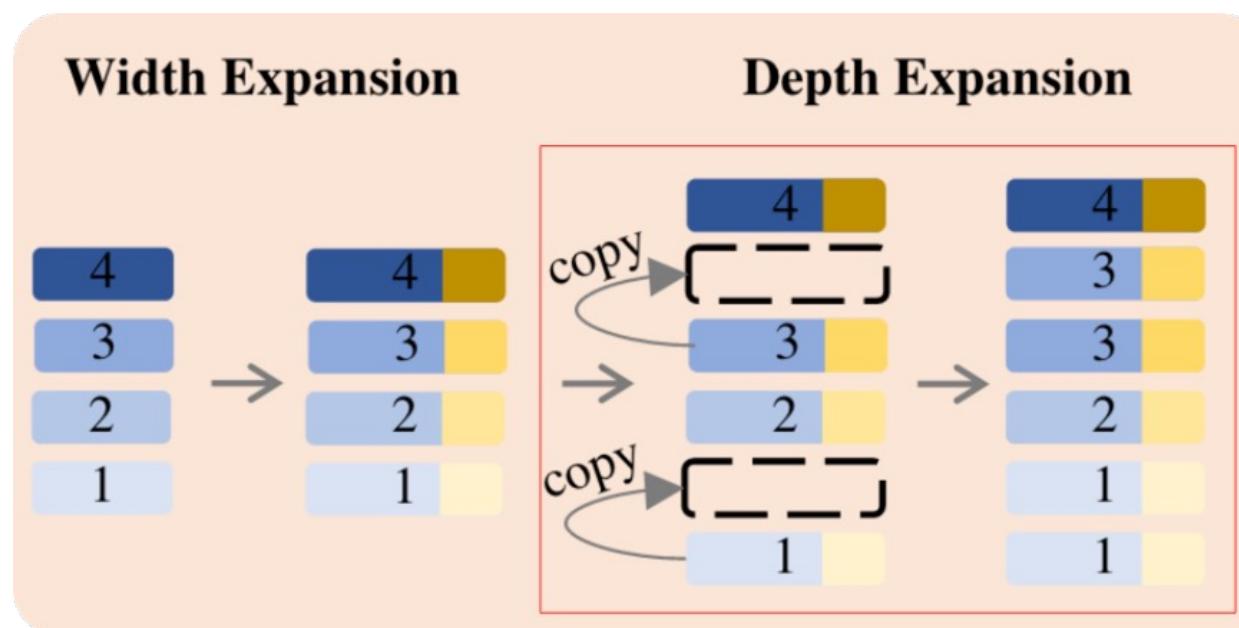
# Methodology

- Depth Expansion(previous work)



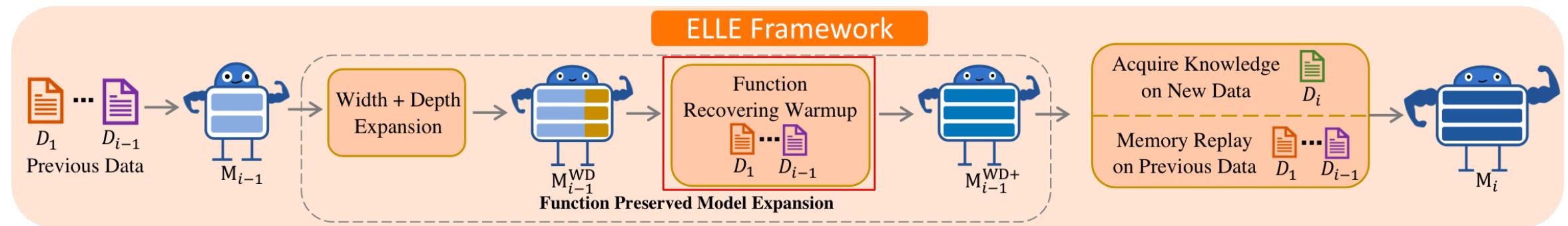
# Methodology

- Depth Expansion



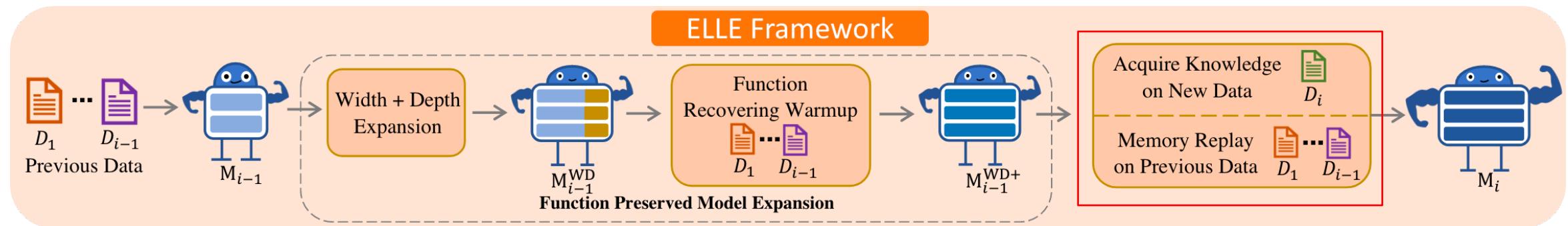
# Methodology

- (2) **Function Recovering Warmup:** pre-train the initialized PLM on the previous corpora conserved in the memory to recover the language abilities lost during model expansion



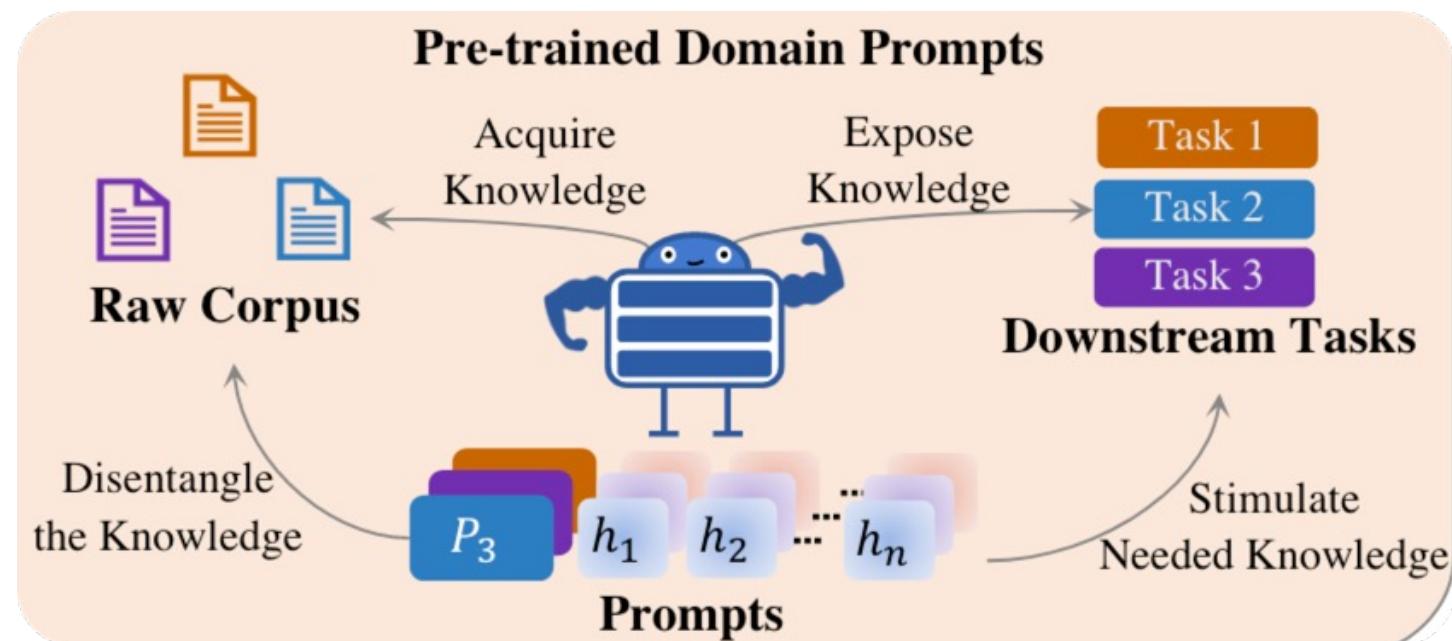
# Methodology

- (3) Acquiring New Knowledge Memory Replay



# Methodology

- (4) **Pre-trained domain prompts:** disentangle the versatile knowledge learned during pre-training and stimulate the proper knowledge for downstream tasks



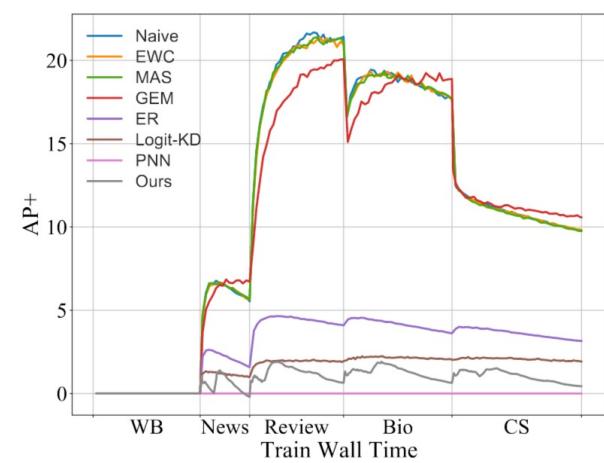
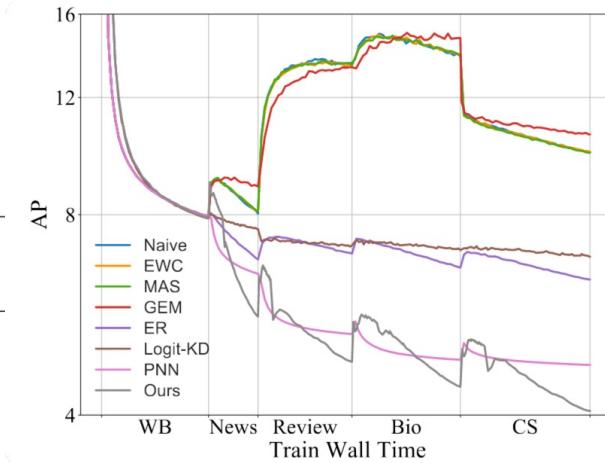
# Experimental Setting

- Pre-train on 5 sequential domains
  - Wikipedia, News, Reviews, Biomedical and Computer Science
- Evaluate ELLE on both BERT and GPT architecture
- Growing PLMs from both (1) 30M to 125M and (2) 125M to 355M
- Metrics:
  - Pre-training: average perplexity (AP) & averaged increased perplexity (AP+)
    - $AP = \exp\left(\frac{1}{j} \sum_{i=1}^j \log PPL_{T,i}\right)$
    - $AP^+ = \frac{1}{j-1} \sum_{i=1}^{j-1} (PPL_{T,i} - PPL_{i,i}^f)$
  - Fine-tuning: downstream performance on each domain
    - MNLI, HYPERPARTISAN, HELPFULNESS, CHEMPROT, ACL-ARC

# Experimental Results

- Pre-training
  - Lower average perplexity (AP) and less catastrophic forgetting (AP+)

Domain	WB		Ns		REV		BIO		CS	
Metrics	AP	AP+	AP	AP+	AP	AP+	AP	AP+	AP	AP+
<i>Growing from BERT<sub>L6_D384</sub> to BERT<sub>L12_D768</sub></i>										
Naive (Lower Bound)	7.96	-	8.03	5.54	13.52	21.42	13.86	17.67	9.93	9.81
EWC	7.96	-	8.09	5.65	13.40	20.98	13.92	17.75	9.94	9.82
MAS	7.96	-	8.08	5.65	13.44	21.17	13.87	17.67	9.91	9.75
A-GEM	7.96	-	8.82	6.72	13.31	20.06	14.73	18.89	10.56	10.58
ER	7.96	-	6.85	1.59	6.99	4.09	6.66	3.62	6.39	3.16
Logit-KD	7.96	-	7.60	0.99	7.19	1.95	7.08	2.02	6.92	1.92
PNN	7.96	-	6.52	0.00	5.29	<b>0.00</b>	4.84	<b>0.00</b>	4.76	<b>0.00</b>
ELLE (ours)	<b>7.92</b>	-	<b>5.62</b>	<b>-0.20</b>	<b>4.81</b>	0.64	<b>4.41</b>	0.64	<b>4.06</b>	0.44
<i>Growing from BERT<sub>L12_D768</sub> to BERT<sub>L24_D1024</sub></i>										
ER	4.54	-	4.33	1.31	4.02	1.46	3.73	1.15	3.82	1.28
ELLE (ours)	<b>4.52</b>	-	<b>3.89</b>	<b>0.47</b>	<b>3.61</b>	<b>0.75</b>	<b>3.66</b>	<b>0.97</b>	<b>3.29</b>	<b>0.54</b>
<i>Growing from GPT<sub>L6_D384</sub> to GPT<sub>L12_D768</sub></i>										
Naive (Lower Bound)	46.54	-	52.91	37.96	81.28	177.22	94.44	160.51	60.64	80.48
MAS	46.54	-	53.12	38.44	81.23	177.20	93.21	157.93	60.62	80.28
ER	46.54	-	44.49	12.42	35.46	21.78	33.24	23.38	31.94	19.83
Logit-KD	46.54	-	48.93	5.41	37.60	9.97	34.60	11.74	33.67	11.19
PNN	46.54	-	39.90	<b>0.00</b>	26.84	<b>0.00</b>	<b>22.19</b>	<b>0.00</b>	21.43	<b>0.00</b>
ELLE (ours)	<b>46.50</b>	-	<b>36.84</b>	2.25	<b>25.60</b>	4.38	22.29	5.88	<b>20.49</b>	4.31



# Experimental Results

- **Fine-tuning**
  - Better downstream performance on each domain

Domain	WB	Ns	REV	BIO	CS	AVG
<i>Growing from BERT<sub>L6_D384</sub> to BERT<sub>L12_D768</sub></i>						
Naive	77.2	72.8	60.6	77.1	64.8	70.5
EWC	77.4	72.8	61.6	77.5	59.6	69.8
MAS	77.1	73.7	60.7	77.5	68.2	71.5
A-GEM	76.6	71.4	61.5	76.9	67.5	70.8
ER	77.6	72.2	61.9	78.3	63.5	70.7
Logit-KD	77.2	69.5	63.9	76.8	58.9	69.2
PNN	76.0	76.3	68.0	79.5	65.2	73.0
ELLE	<b>83.2</b>	<b>81.8</b>	<b>68.5</b>	<b>82.9</b>	<b>72.7</b>	<b>77.8</b>
<i>Growing from BERT<sub>L12_D768</sub> to BERT<sub>L24_D1024</sub></i>						
ER	84.7	83.3	68.0	82.7	71.4	78.0
ELLE	<b>86.3</b>	<b>90.4</b>	<b>70.5</b>	<b>84.2</b>	<b>73.8</b>	<b>81.0</b>

# Experimental Results

- **Ablation Study (pre-training)**

Domain					WB		Ns		REV		BIO		CS	
WE	DE	FRW	$\delta_N$	PT	AP	AP+	AP	AP+	AP	AP+	AP	AP+	AP	AP+
					7.96	-	6.85	1.59	6.99	4.09	6.66	3.62	6.39	3.16
✓		✓			7.96	-	6.23	0.78	5.34	1.42	4.98	1.20	4.48	0.89
	✓	✓			7.96	-	5.81	0.03	5.49	1.43	5.16	1.32	4.79	0.94
✓	✓	✓			7.96	-	5.78	0.02	4.91	0.76	4.49	0.73	4.13	0.52
✓	✓				7.96	-	5.79	0.09	5.09	1.13	4.58	0.88	4.22	0.65
✓	✓	✓	✓		7.96	-	5.69	-0.13	4.85	0.67	4.45	0.69	4.09	0.47
✓	✓	✓	✓	✓	<b>7.92</b>	-	<b>5.62</b>	<b>-0.20</b>	<b>4.81</b>	<b>0.64</b>	<b>4.41</b>	<b>0.64</b>	<b>4.06</b>	<b>0.44</b>

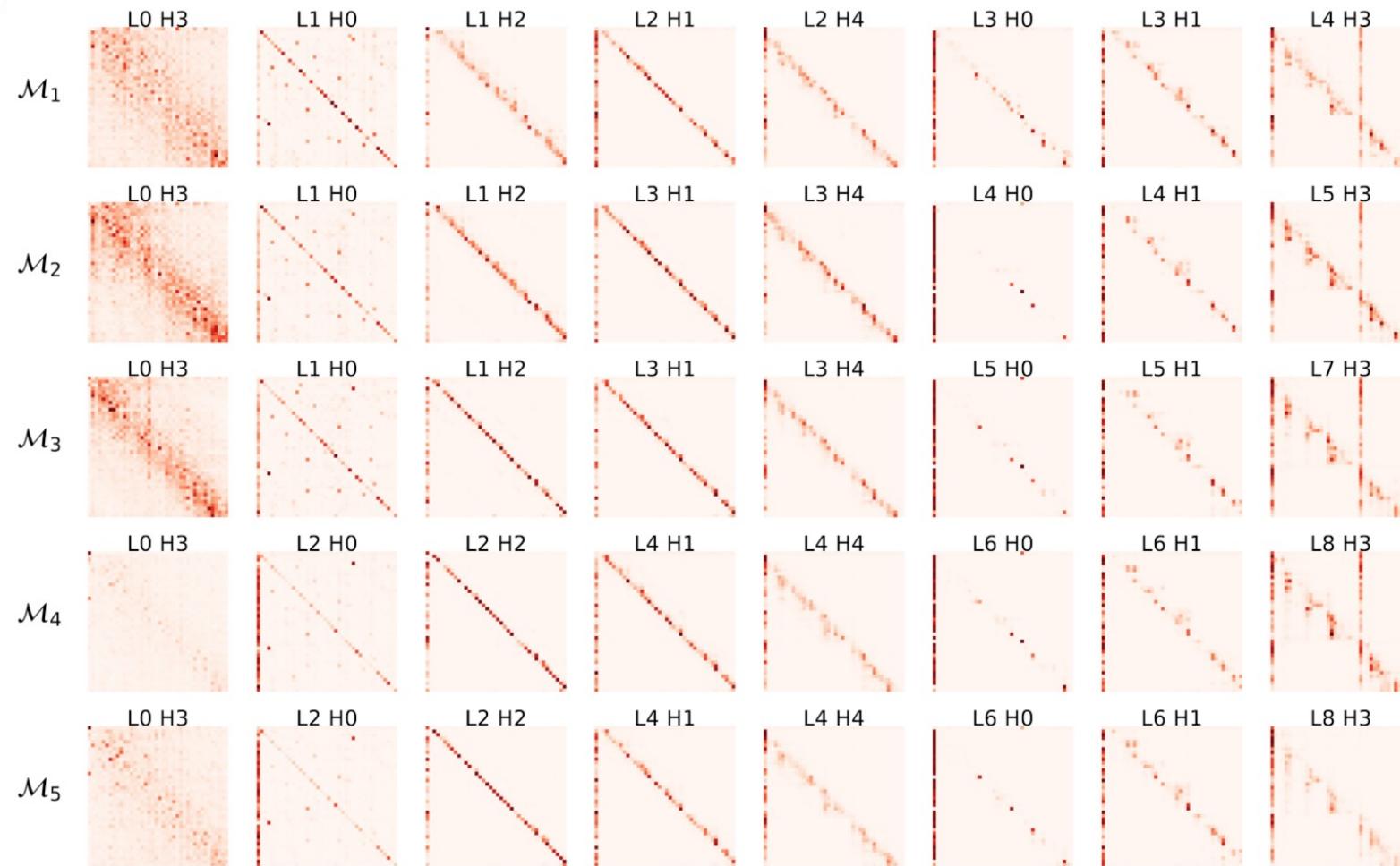
# Experimental Results

- Effects of Domain Prompts during Fine-tuning

Domain	WB	Ns	REV	BIO	CS	AVG
ELLE – PT <sub>fine-tune</sub>	82.9	79.9	67.0	82.1	67.7	75.9
ELLE + $\neg$ PT <sub>fine-tune</sub>	83.1	80.6	68.1	81.7	70.8	76.9
ELLE	<b>83.2</b>	<b>81.8</b>	<b>68.5</b>	<b>82.9</b>	<b>72.7</b>	<b>77.8</b>

# Experimental Results

- **Visualization of the Attention Patterns**



**Thanks**