



Self-Supervised Meta-Learning for Few-Shot Natural Language Classification Tasks

Trapit Bansal \diamond^* and Rishikesh Jha † and Tsendsuren Munkhdalai ‡ and Andrew McCallum \diamond

\diamond University of Massachusetts, Amherst

† Code for Science and Society

‡ Microsoft Research, Montréal, Canada

目录

Contents

1

问题与方案

2

算法介绍

3

实验设置与分析

问题设定

自监督的预训练语言模型，可以提供一个好的模型参数初始化点，使得到了下游任务微调时，可以有不错的泛化性和效果。



微调的数据利用比较低，即在Few-Shot的场景下，利用少量的样本进行微调效果得到不充分的体现。所以这篇文章关注点就在于，如果提高预训练的语言模型在Few-Shot场景下进行分类任务时的泛化能力。



提高模型的泛化能力，可以采用Meta-Learning方式来解决。



任务构建，需要大量的训练任务，以提高泛化性。且对于不同的任务，需要构建不同的模型框架。



利用类似cloze-style的方式，通过Self-Supervised从数据本身来获得标签，构建受监督的元学习任务。通过学习特定任务的Support set来生成任务特定的参数，以支持训练能够适应不同类任务的元学习模型[1]

[1] Trapti Bansal, Rishikesh Jha, and Andrew McCallum. 2019. Learning to few-shot learn across diverse natural language classification tasks.

原因分析—challenge 1

虽然自监督的预训练是有效的，但是它对数据的利用率会比较低，而且微调时也需要大量的数据才能在目标任务上有好的效果[1][2]。

这就是Few-Shot Learning问题，模型只给出几个新任务的例子，并期望在该任务中表现良好。本文重点关注Few-Shot的问题，并开发了能够更好地对新任务进行Few-Shot 泛化的模型。

1. 此外，对预训练的模型进行微调通常会引入新的随机参数，如Softmax层和重要的超参数（如learning-rate），这些参数很难从少数示例中有效地评估。
2. 大规模的预训练会遭受训练测试不匹配的困扰，因为模型没有经过优化就去学习一个初始点，而该初始点经过少量样本微调后性能较低。

因此，本文提出要消除这种训练测试不匹配，并通过联合学习模型的初始点和超参数，从而允许数据高效微调，即作为一个元学习问题[1][3]。

[2] Dani Yogatama, Cyprien de Masson d'Autume, Jerome Connor, Tom'as Kocisk'y, Mike Chrzanowski, Ling peng Kong, Angeliki Lazaridou, Wang Ling, Lei Yu, Chris Dyer, and Phil Blunsom. 2019. Learning and evaluating general linguistic intelligence.

[3] Sebastian Thrun and Lorien Pratt. 2012. Learning to learn. Springer Science & Business Media.

原因分析—Solution 1

1. 大规模的预训练会遭受训练测试失配的困扰，因为模型没有经过优化就去学习一个初始点，而该初始点经过少量样本微调后性能较低。

MAML

Loss Function:

$$L(\phi) = \sum_{n=1}^N l^n(\hat{\theta}^n)$$

$\hat{\theta}^n$: model learned from task n

$\hat{\theta}^n$ depends on ϕ

$l^n(\hat{\theta}^n)$: loss of task n on the testing set of task n

How to minimize $L(\phi)$? Gradient Descent

$$\theta'_i \leftarrow \phi - \alpha \nabla_{\phi} \mathcal{L}_i(\mathcal{D}^{tr}, \phi)$$

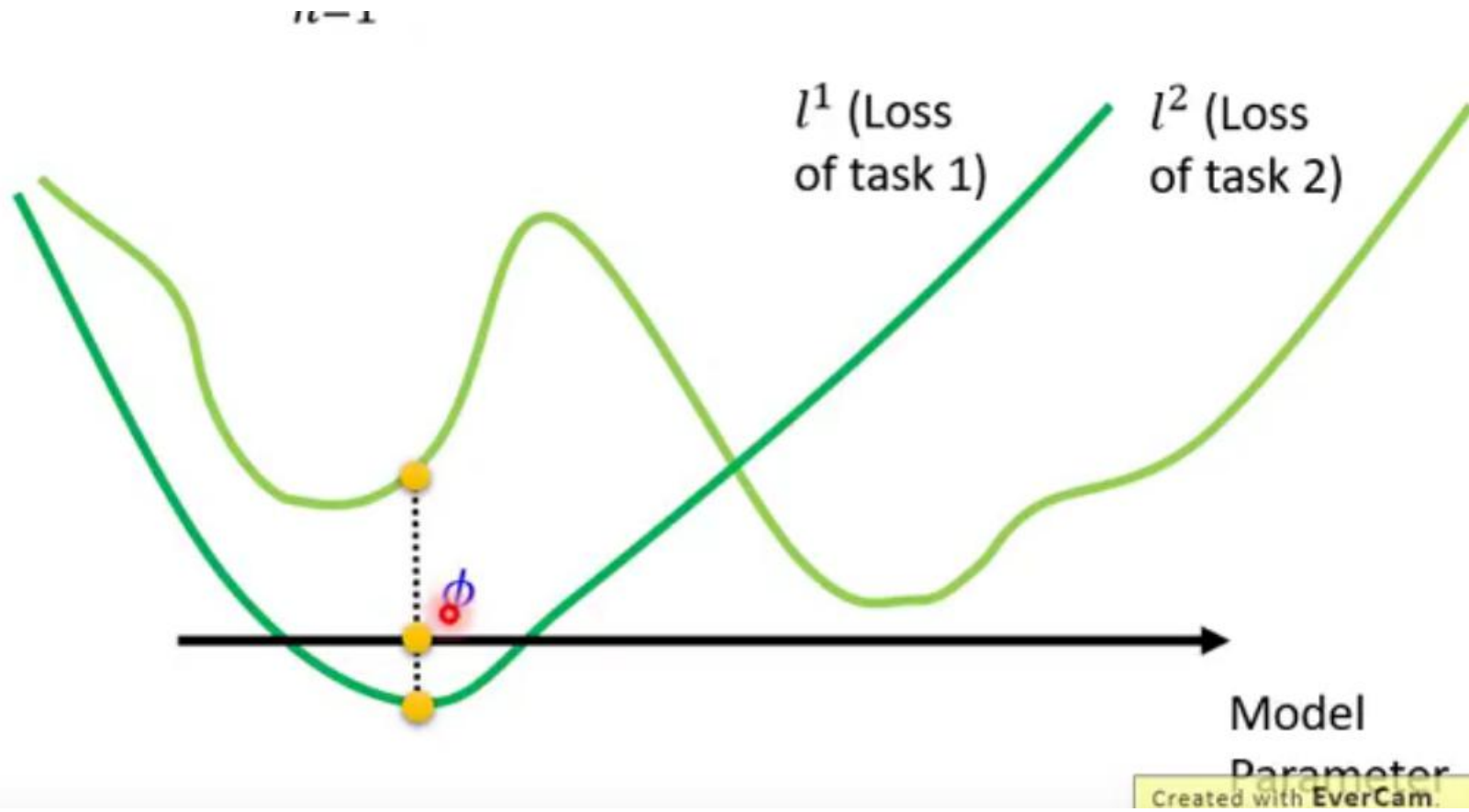
Model Pre-training

Widely used in
transfer learning

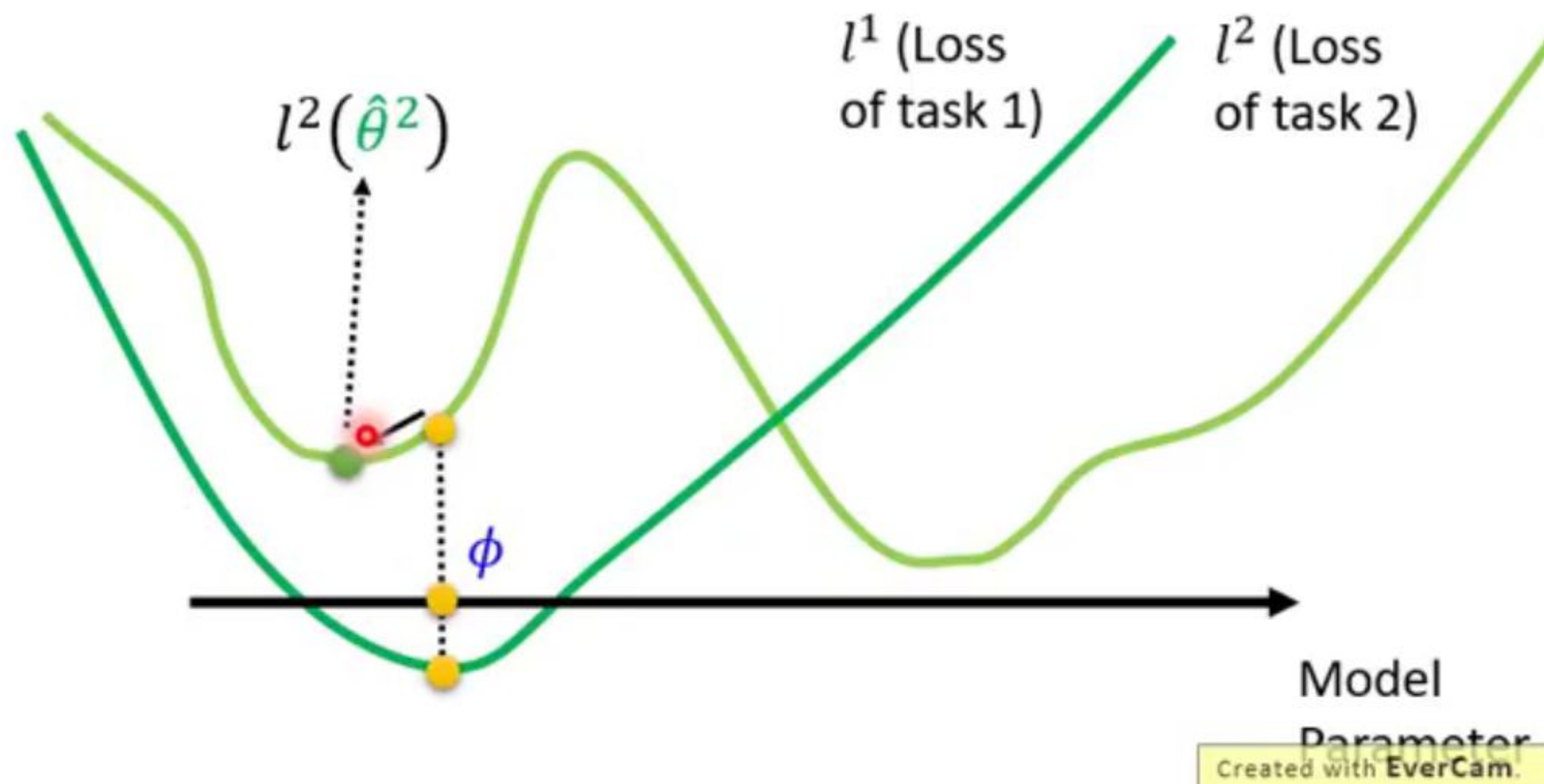
Loss Function:

$$L(\phi) = \sum_{n=1}^N l^n(\phi)$$

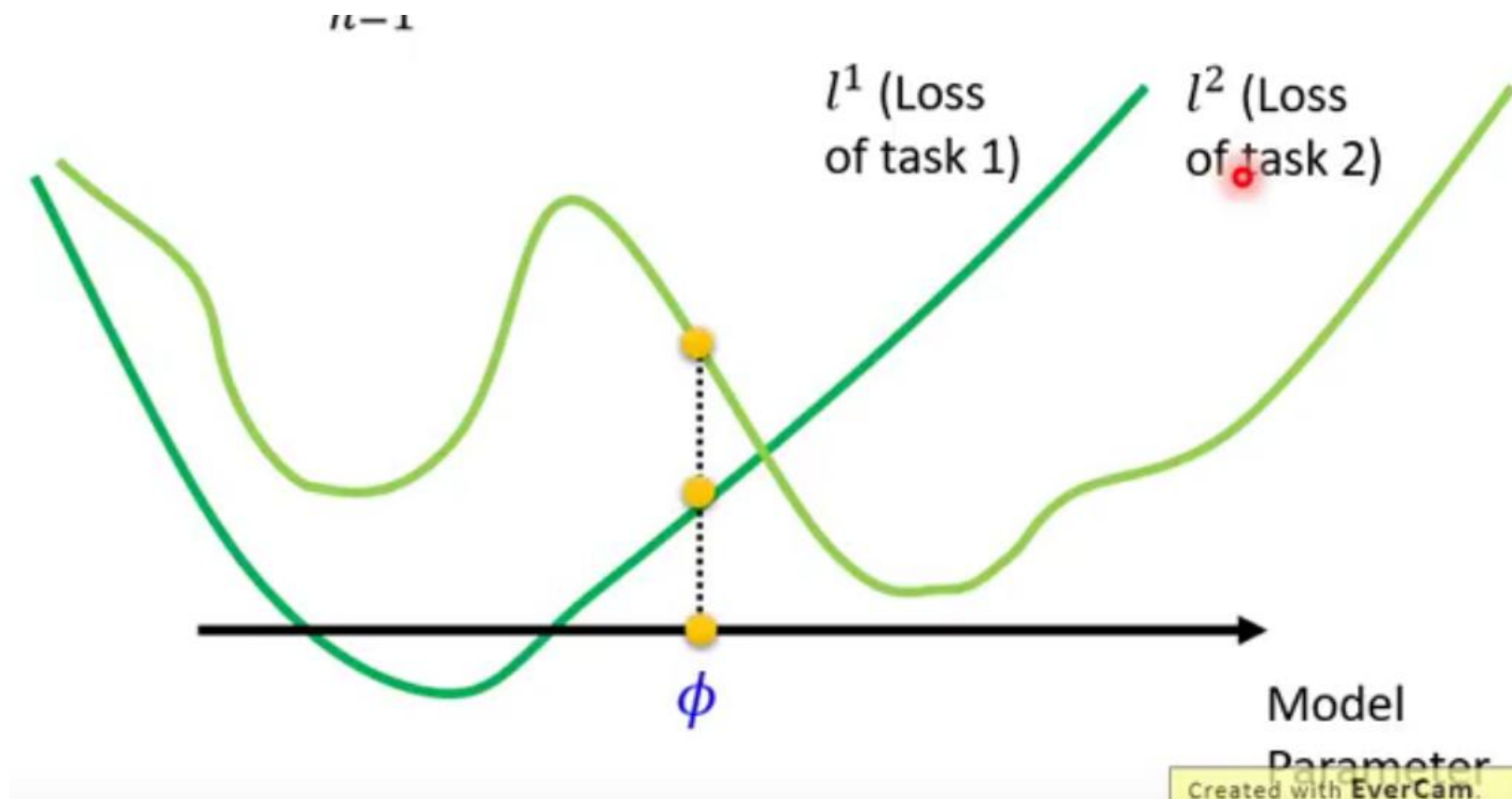
Model Pre-training



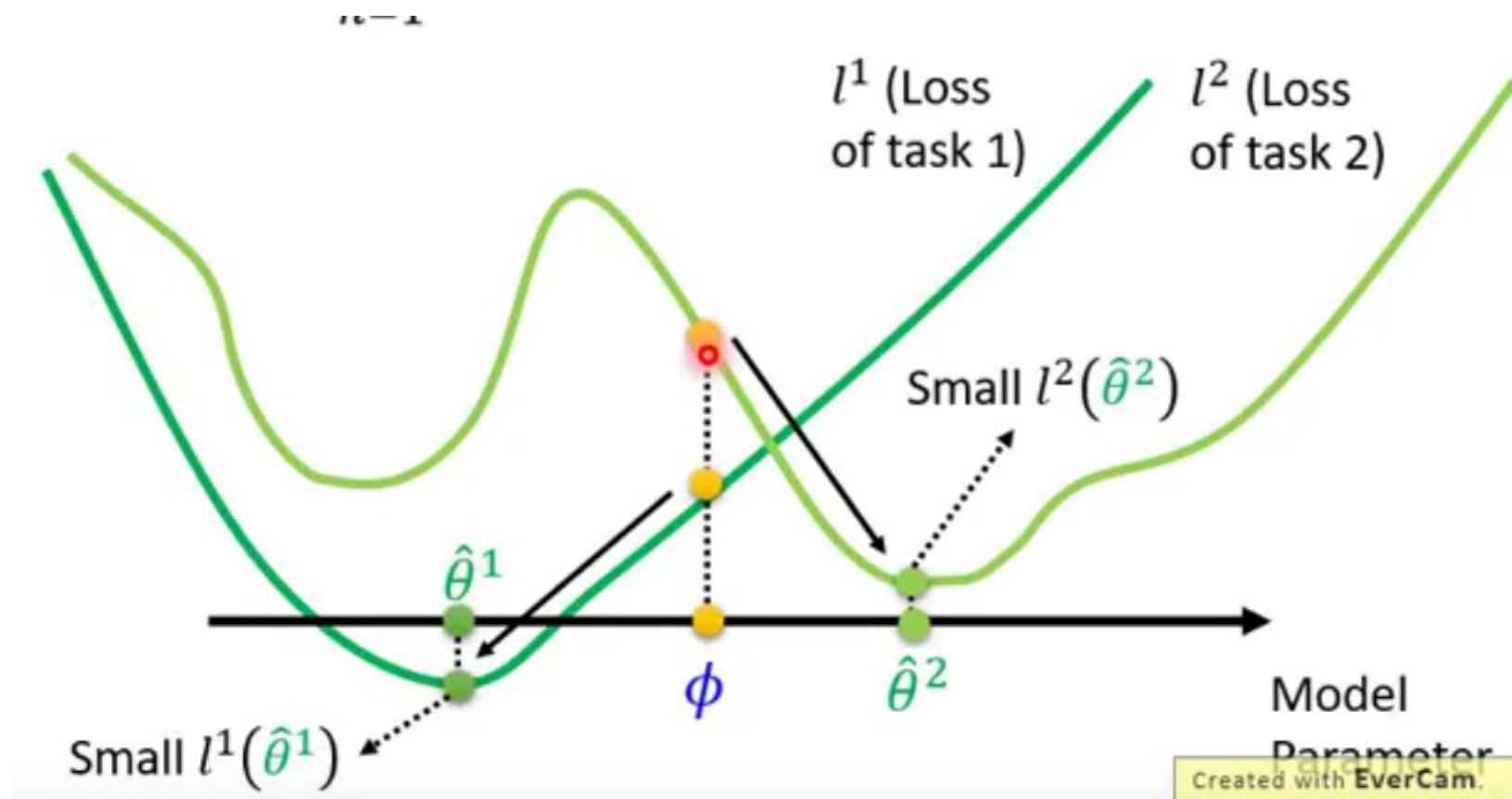
Pre-training Model



MAML



MAML



原因分析—challenge 2

Meta-Learning的任务构建，需要大量的训练任务，以提高泛化性。
元学习在相同的任务下，在新的标签下泛化性不错，但是对于不同的任务，需要构建不同的模型框架，使得在不同任务下泛化性较差[4]。

1. 对于不同的任务，需要构建不同的模型框架。
2. Meta Over-Fitting，对于任务的分布过拟合了[4]。

原因分析—Solution 2

1. Meta-Learning的任务构建，需要大量的训练任务，以提高泛化性。
Subset Masked Language Modeling Tasks (SMLMT)

Subset: {Democratic, Capital}

↓

Support set

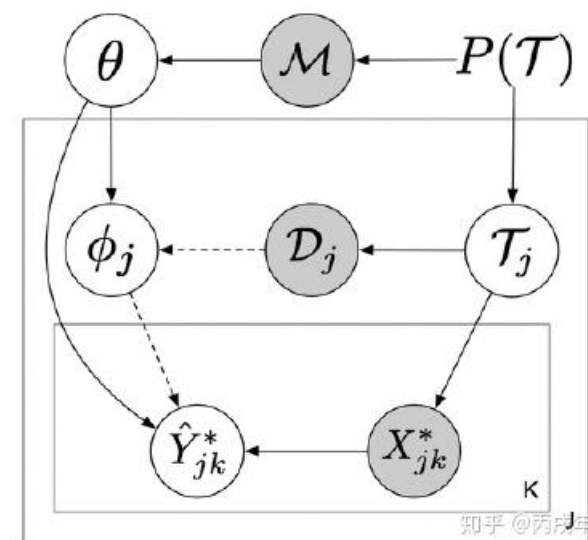
Sentence	Class
A member of the [m] Party, he was the first African American to be elected to the presidency.	1
The [m] Party is one of the two major contemporary political parties in the United States, along with its rival, the Republican Party.	1
Honolulu is the [m] and largest city of the U.S. state of Hawaii.	2
Washington, D.C., formally the District of Columbia and commonly referred to as Washington or D.C., is the [m] of the United States.	2

Query: New Delhi is an urban district of Delhi which serves as the [m] of India

Correct Prediction: 2

原因分析—Solution 2

2. Meta Over-Fitting, 对于任务的分布过拟合了[4][5]。

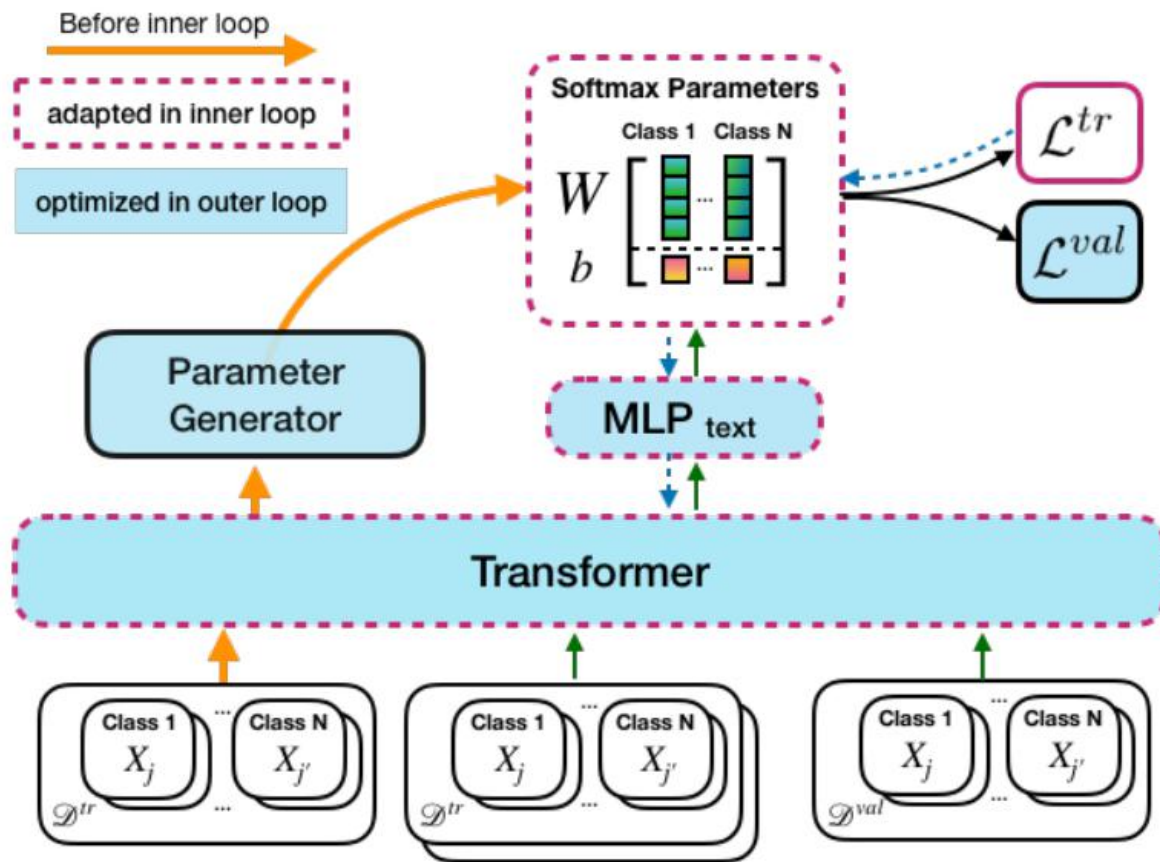


$$-\frac{1}{N} \sum_i \mathbb{E}_{q(\theta | \mathcal{M}) q(\phi | D_i^{train}, \theta)} \left[\frac{1}{K} \sum_{(x^*, y^*) \in D_i^{test}} \log q(\hat{y}^* = y^* | x^*, \theta, \phi) \right],$$

- [4] Mingzhang Yin, George Tucker, Mingyuan Zhou, Sergey Levine, and Chelsea Finn. 2020. Meta learning without memorization. In International Conference on Learning Representations. 12
- [5] Luke Metz, Niru Maheswaranathan, Brian Cheung, and Jascha Sohl-Dickstein. 2019. Learning unsupervised learning rules. In International Conference on Learning Representations.

原因分析—Solution 2

3. 对于不同的任务，需要构建不同的模型框架[6]。



$$w_t^n, b_t^n = g_\psi(\{f_\pi(X_{1n}), \dots, f_\pi(X_{kn})\}) \quad (3)$$

$$p(y|X) = softmax\{\mathbf{W}_t h_\phi(f_\pi(X)) + \mathbf{b}_t\} \quad (4)$$

where $\mathbf{W}_t = [w_t^1; \dots; w_t^N] \in \mathcal{R}^{N \times d}$, $\mathbf{b}_t = [b_t^1; \dots; b_t^N] \in \mathcal{R}^d$ are the concatenation of the per-class vectors in (3), and h_ϕ is a MLP with parameters ϕ and output dimension d .

本文工作

1. 提出了Subset Masked Language Modeling Tasks (SMLMT)，通过自监督的方式创建元学习的任务。并且通过SMLMT创建的任务分布取得了SOTA的结果。
2. 证明了自监督SMLMT也可以与受监督任务的数据相结合，以实现更好的特征学习，同时通过使用SMLMT避免Meta Over-Fitting以保证模型泛化的泛化性。
3. 研究了参数数对Few-Shot Learning的影响，发现规模更大的预训练或元学习模型比小模型泛化性更好，即使对于小模型，元学习也能获得不错的效果。

目录

Contents

1

问题与方案

2

相关方法

3

实验设置与分析

Hybrid SMLMT

来自SMLMT的任务也可以与受监督的任务相结合，以鼓励更好的特征学习[6]，并增加元学习任务的多样性。

使用采样比 $\lambda \in (0,1)$ ，在每个episode中以 λ 的概率选择SMLMT任务或以 $(1 - \lambda)$ 的概率选择监督任务。SMLMT与受监督任务的联合使用可以改善Meta Over-Fitting。

[6]Rich Caruana. 1997. Multitask learning. Machine learning, 28(1):41–75.

元训练算法流程

Algorithm 1 Meta-Training

Require: SMLMT task distribution \mathcal{T} and supervised tasks \mathcal{S} , model parameters $\{\pi^w, \pi, \phi, \psi, \alpha\}$, adaptation steps G , learning-rate β , sampling ratio λ
Initialize θ with pre-trained BERT-base;

π^w 是warp layers的参数

π 是Transformer除 π^w 外的参数

ϕ 是MLP的参数

ψ 是类别编码器参数

α 是学习率

the parameters are meta-trained using the MAML algorithm (Finn et al., 2017). Concretely, set $\theta := \{\pi, \phi, \mathbf{W}_t, \mathbf{b}_t\}$ for the task-specific inner loop gradient updates in (1) and set $\Theta := \{\pi, \psi, \alpha\}$ for the outer-loop updates in (2). Note that we do multiple steps of gradient descent in the inner loop.

```
1: while not converged do
2:   for task_batchsize times do
3:      $t \sim \text{Bernoulli}(\lambda)$ 
4:      $T \sim t \cdot \mathcal{T} + (1 - t) \cdot \mathcal{S}$ 
5:      $\mathcal{D}^{tr} = \{(x_j, y_j)\} \sim T$ 
6:      $C^n \leftarrow \{x_j | y_j = n\}; \quad N \leftarrow |C^n|$ 
7:      $w^n, b^n \leftarrow \frac{1}{|C^n|} \sum_{x_j \in C^n} g_\psi(f_\pi(\mathcal{D}^{tr}))$ 
8:      $\mathbf{W} \leftarrow [w^1; \dots; w^N]; \quad \mathbf{b} \leftarrow [b^1; \dots; b^N]$ 
9:      $\theta \leftarrow \{\pi, \phi, \mathbf{W}, \mathbf{b}\}; \quad \theta^{(0)} \leftarrow \theta$ 
10:     $\Theta \leftarrow \{\pi^w, \pi, \psi, \alpha\}$ 
11:     $\mathcal{D}^{val} \sim T$ 
12:     $q_T \leftarrow 0$ 
13:    for  $s := 0 \dots G - 1$  do
14:       $\mathcal{D}_s^{tr} \sim T$ 
15:       $\theta^{(s+1)} \leftarrow \theta^{(s)} - \alpha \nabla_{\theta} \mathcal{L}_T(\{\Theta, \theta^{(s)}\}, \mathcal{D}_s^{tr})$ 
16:       $q_T \leftarrow q_T + \nabla_{\Theta} \mathcal{L}_T(\{\Theta, \theta^{(s+1)}\}, \mathcal{D}^{val})$ 
17:    end for
18:  end for
19:   $\Theta \leftarrow \Theta - \beta \cdot \sum_T \frac{q_T}{G}$ 
20: end while
```

目录

Contents

1

问题与方案

2

相关方法

3

实验设置与分析

数据集统计

Dataset	Labels	Train	Validation	Test
CoLA	2	8551	1042	—
MRPC	2	3669	409	—
QNLI	2	104744	5464	—
QQP	2	363847	40431	—
RTE	2	2491	278	—
SNLI	3	549368	9843	—
SST-2	2	67350	873	—
MNLI (m/mm)	3	392703	19649	—
Scitail	2	23,596	1,304	2,126
Amazon Sentiment Domains	2	800	200	1000
Airline	3	7320	—	7320
Disaster	2	4887	—	4887
Political Bias	2	2500	—	2500
Political Audience	2	2500	—	2500
Political Message	9	2500	—	2500
Emotion	13	20000	—	20000
CoNLL	4	23499	5942	5648
MIT-Restaurant	8	12474	—	2591

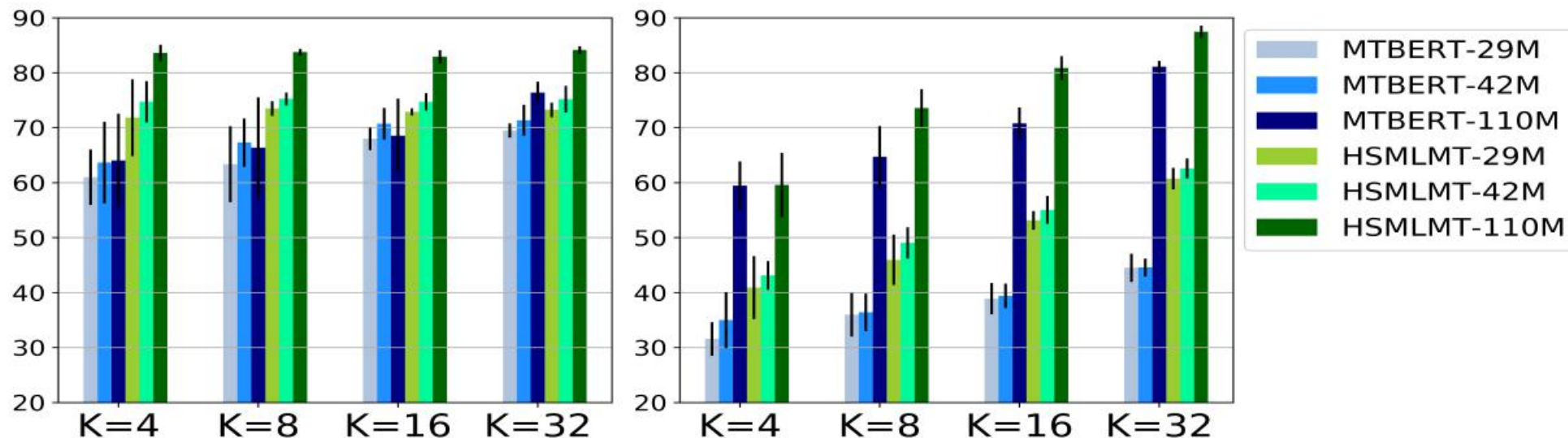
k-shot accuracy on novel tasks not seen in training.

Task	N	k	BERT	SMLMT	MT-BERT _{softmax}	MT-BERT	LEOPARD	Hybrid-SMLMT
CoNLL	4	4	50.44 \pm 08.57	46.81 \pm 4.77	52.28 \pm 4.06	55.63 \pm 4.99	54.16 \pm 6.32	57.60 \pm 7.11
		8	50.06 \pm 11.30	61.72 \pm 3.11	65.34 \pm 7.12	58.32 \pm 3.77	67.38 \pm 4.33	70.20 \pm 3.00
		16	74.47 \pm 03.10	75.82 \pm 4.04	71.67 \pm 3.03	71.29 \pm 3.30	76.37 \pm 3.08	80.61 \pm 2.77
		32	83.27 \pm 02.14	84.01 \pm 1.73	73.09 \pm 2.42	79.94 \pm 2.45	83.61 \pm 2.40	85.51 \pm 1.73
MITR	8	4	49.37 \pm 4.28	46.23 \pm 3.90	45.52 \pm 5.90	50.49 \pm 4.40	49.84 \pm 3.31	52.29 \pm 4.32
		8	49.38 \pm 7.76	61.15 \pm 1.91	58.19 \pm 2.65	58.01 \pm 3.54	62.99 \pm 3.28	65.21 \pm 2.32
		16	69.24 \pm 3.68	69.22 \pm 2.78	66.09 \pm 2.24	66.16 \pm 3.46	70.44 \pm 2.89	73.37 \pm 1.88
		32	78.81 \pm 1.95	78.82 \pm 1.30	69.35 \pm 0.98	76.39 \pm 1.17	78.37 \pm 1.97	79.96 \pm 1.48
Airline	3	4	42.76 \pm 13.50	42.83 \pm 6.12	43.73 \pm 7.86	46.29 \pm 12.26	54.95 \pm 11.81	56.46 \pm 10.67
		8	38.00 \pm 17.06	51.48 \pm 7.35	52.39 \pm 3.97	49.81 \pm 10.86	61.44 \pm 03.90	63.05 \pm 8.25
		16	58.01 \pm 08.23	58.42 \pm 3.44	58.79 \pm 2.97	57.25 \pm 09.90	62.15 \pm 05.56	69.33 \pm 2.24
		32	63.70 \pm 4.40	65.33 \pm 3.83	61.06 \pm 3.89	62.49 \pm 4.48	67.44 \pm 01.22	71.21 \pm 3.28
Disaster	2	4	55.73 \pm 10.29	62.26 \pm 9.16	52.87 \pm 6.16	50.61 \pm 8.33	51.45 \pm 4.25	55.26 \pm 8.32
		8	56.31 \pm 09.57	67.89 \pm 6.83	56.08 \pm 7.48	54.93 \pm 7.88	55.96 \pm 3.58	63.62 \pm 6.84
		16	64.52 \pm 08.93	72.86 \pm 1.70	65.83 \pm 4.19	60.70 \pm 6.05	61.32 \pm 2.83	70.56 \pm 2.23
		32	73.60 \pm 01.78	73.69 \pm 2.32	67.13 \pm 3.11	72.52 \pm 2.28	63.77 \pm 2.34	71.80 \pm 1.85

k-shot domain transfer accuracy

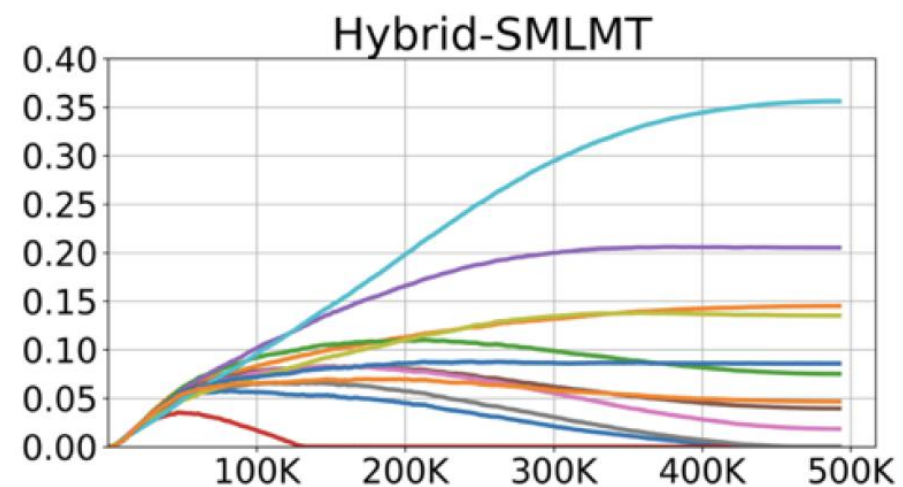
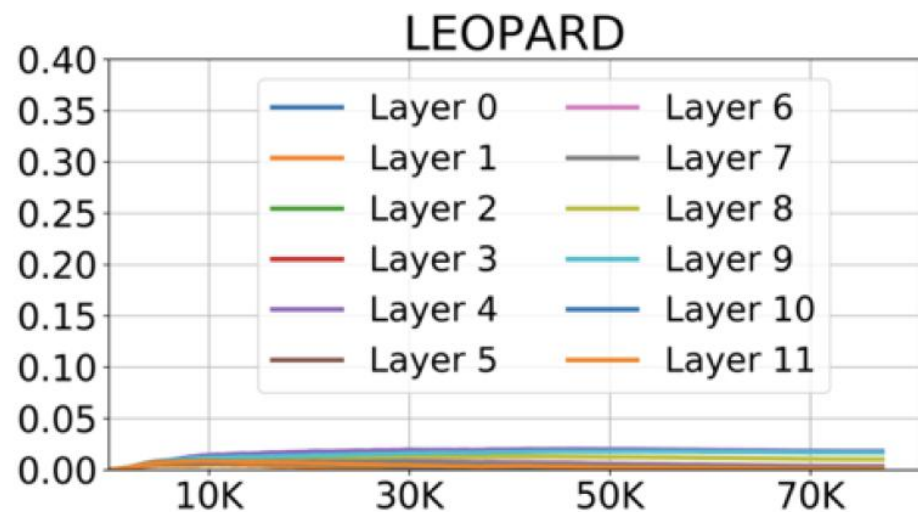
Task	k	BERT _{base}	SMLMT	MT-BERT _{softmax}	MT-BERT	MT-BERT _{reuse}	LEOPARD	Hybrid-SMLMT
Scitail	4	58.53 \pm 09.74	50.68 \pm 4.30	74.35 \pm 5.86	63.97 \pm 14.36	76.65 \pm 2.45	69.50 \pm 9.56	76.75 \pm 3.36
	8	57.93 \pm 10.70	55.60 \pm 2.40	79.11 \pm 3.11	68.24 \pm 10.33	76.86 \pm 2.09	75.00 \pm 2.42	79.10 \pm 1.14
	16	65.66 \pm 06.82	56.51 \pm 3.78	79.60 \pm 2.31	75.35 \pm 04.80	79.53 \pm 2.17	77.03 \pm 1.82	80.37 \pm 1.44
	32	68.77 \pm 6.27	62.38 \pm 3.22	82.23 \pm 1.12	74.87 \pm 3.62	81.77 \pm 1.13	79.44 \pm 1.99	82.20 \pm 1.34
Amazon Books	4	54.81 \pm 3.75	55.68 \pm 2.56	68.69 \pm 5.21	64.93 \pm 8.65	74.79 \pm 6.91	82.54 \pm 1.33	84.70 \pm 0.42
	8	53.54 \pm 5.17	60.23 \pm 5.28	74.86 \pm 2.17	67.38 \pm 9.78	78.21 \pm 3.49	83.03 \pm 1.28	84.85 \pm 0.52
	16	65.56 \pm 4.12	62.92 \pm 4.39	74.88 \pm 4.34	69.65 \pm 8.94	78.87 \pm 3.32	83.33 \pm 0.79	85.13 \pm 0.66
	32	73.54 \pm 3.44	71.49 \pm 4.74	77.51 \pm 1.14	78.91 \pm 1.66	82.23 \pm 1.10	83.55 \pm 0.74	85.27 \pm 0.36
Amazon DVD	4	54.98 \pm 3.96	52.95 \pm 2.51	63.68 \pm 5.03	66.36 \pm 7.46	71.74 \pm 8.54	80.32 \pm 1.02	83.28 \pm 1.85
	8	55.63 \pm 4.34	54.28 \pm 4.20	67.54 \pm 4.06	68.37 \pm 6.51	75.36 \pm 4.86	80.85 \pm 1.23	83.91 \pm 1.14
	16	58.69 \pm 6.08	57.87 \pm 2.69	70.21 \pm 1.94	70.29 \pm 7.40	76.20 \pm 2.90	81.25 \pm 1.41	83.71 \pm 1.04
	32	66.21 \pm 5.41	65.09 \pm 4.37	70.19 \pm 2.08	73.45 \pm 4.37	79.17 \pm 1.71	81.54 \pm 1.33	84.15 \pm 0.94

k-shot performance with number of parameters on Amazon DVD (left), and CoNLL(right)



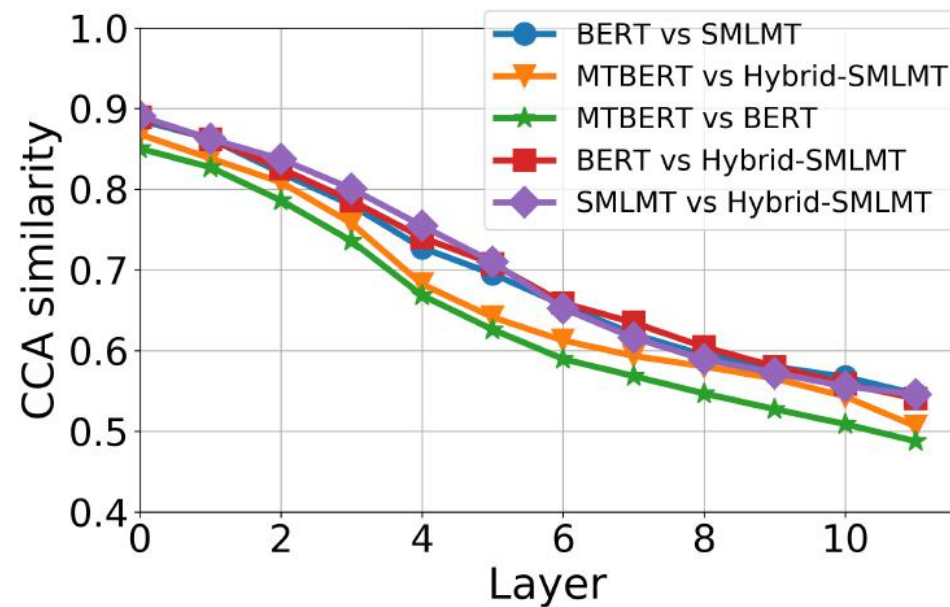
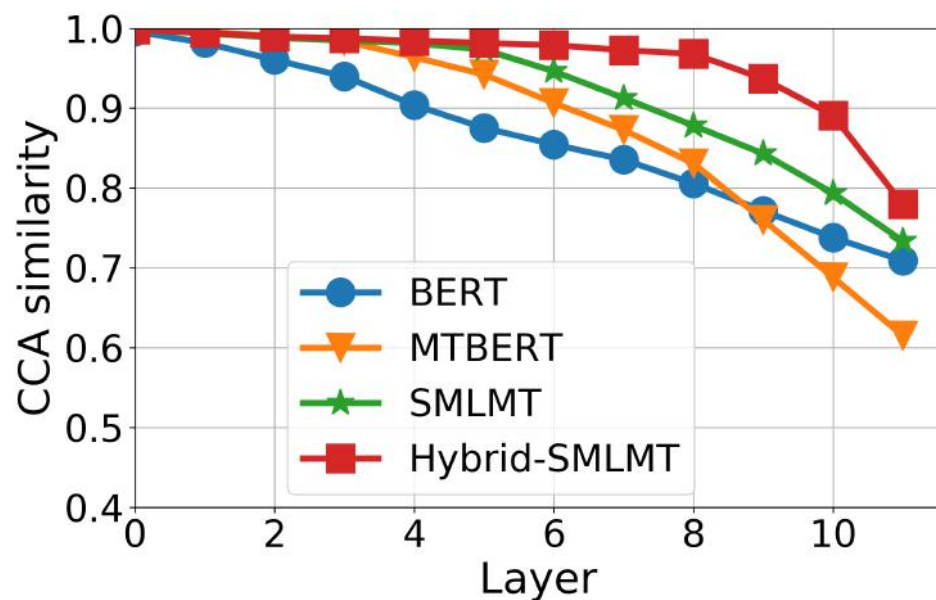
Larger models generalize better and Hybrid-SMLMT provides accuracy gains for all parameter sizes

元训练中的学习率轨迹



LEOPARD学习率在许多层都趋向于0，这表明元拟合过度

Transformer每层的CCA相似度



左图：对相同模型进行微调之前和之后的相似性。

右图：微调后不同模型对之间的相似性。