

Learning to Few-Shot Learn Across Diverse Natural Language Classification Tasks

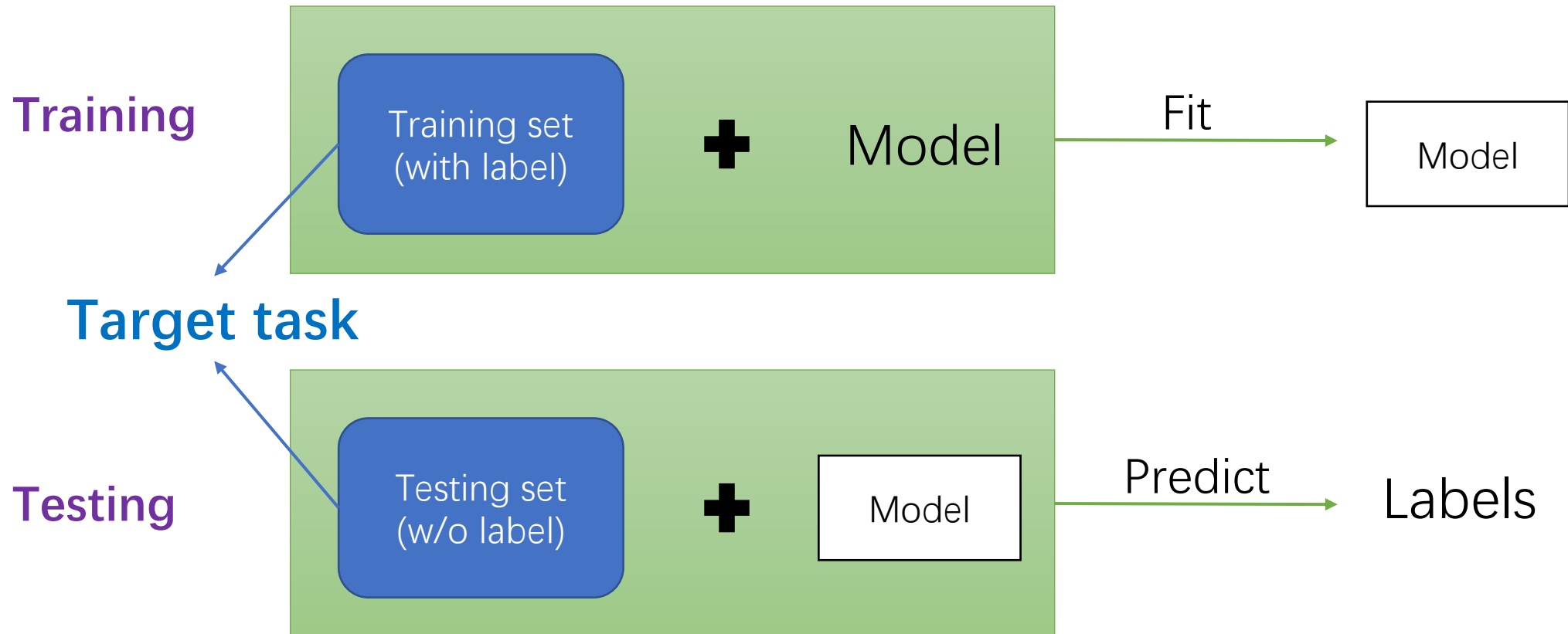
By 李想

2020. 12. 25

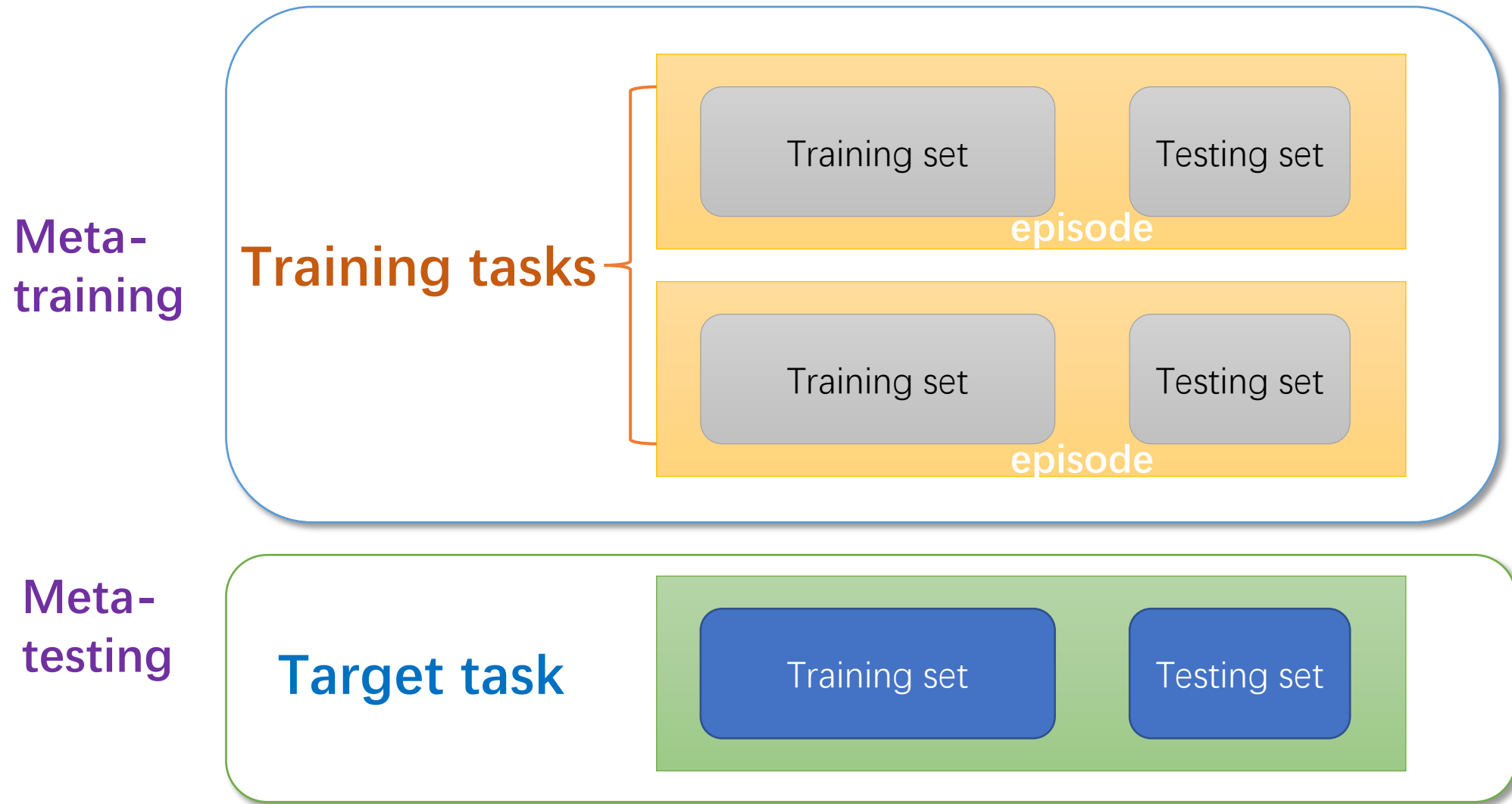
Background: Meta Learning

- Machine Learning = 根据已有数据 找一个函数 f 的能力
- Meta Learning = 根据已有数据 找一个 找一个函数 f 的 函数 F 的能力 (Learn to learn)
 - 授人以鱼不如授人以渔——掌握如何学习的方法

Foreword: Supervised learning

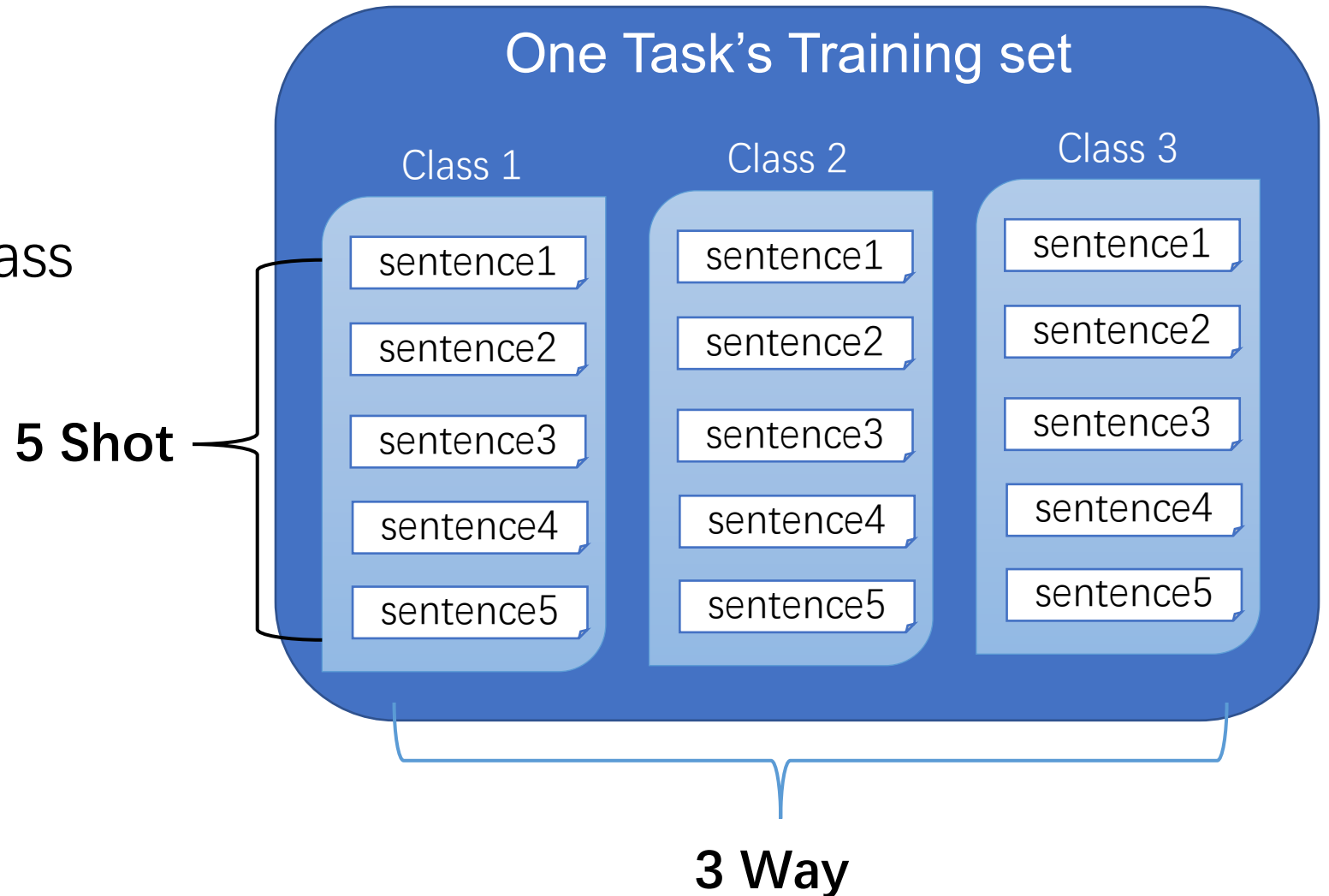


Background: Meta Learning (in few-shot learning)

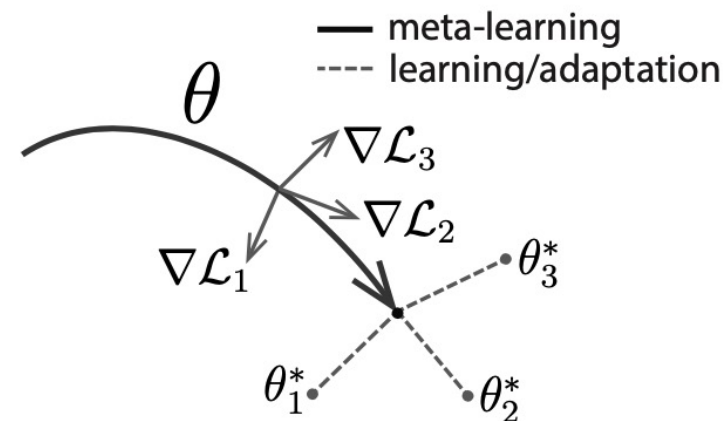


Background: N-way K-shot setting

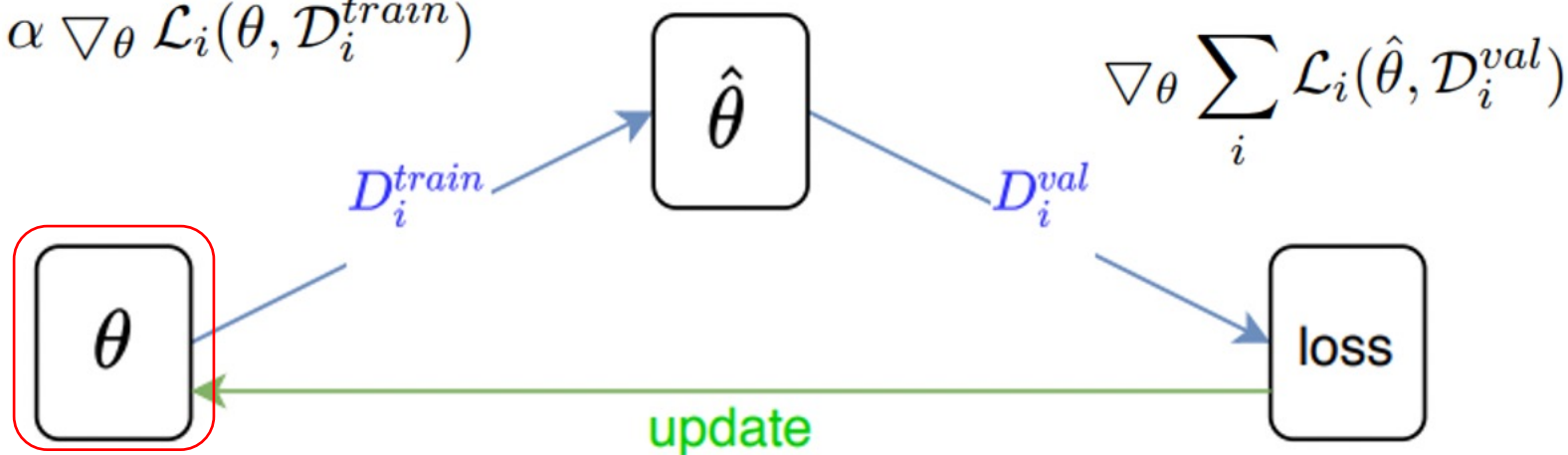
- N: #class in a task
- K: #samples in a class
- Example:
3-way 5-shot



Background: MAML

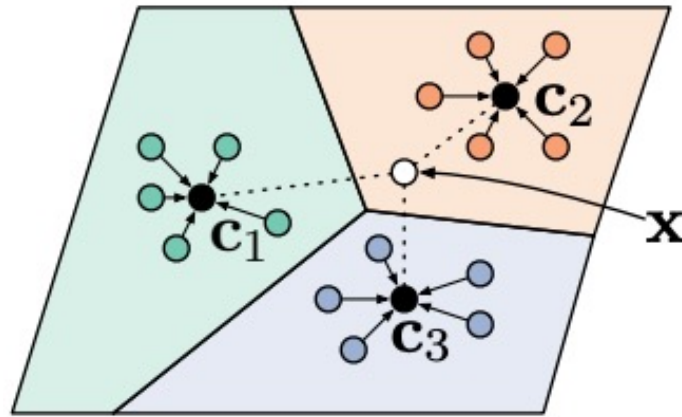


$$\hat{\theta} = \theta - \alpha \nabla_{\theta} \mathcal{L}_i(\theta, \mathcal{D}_i^{train})$$



$$\theta = \theta - \beta \nabla_{\theta} \sum_i \mathcal{L}_i(\hat{\theta}, \mathcal{D}_i^{val})$$

Background: Prototypical network



- More metric-based method:
 - Matching network
 - Relation network
 -

Back to this paper: Motivation

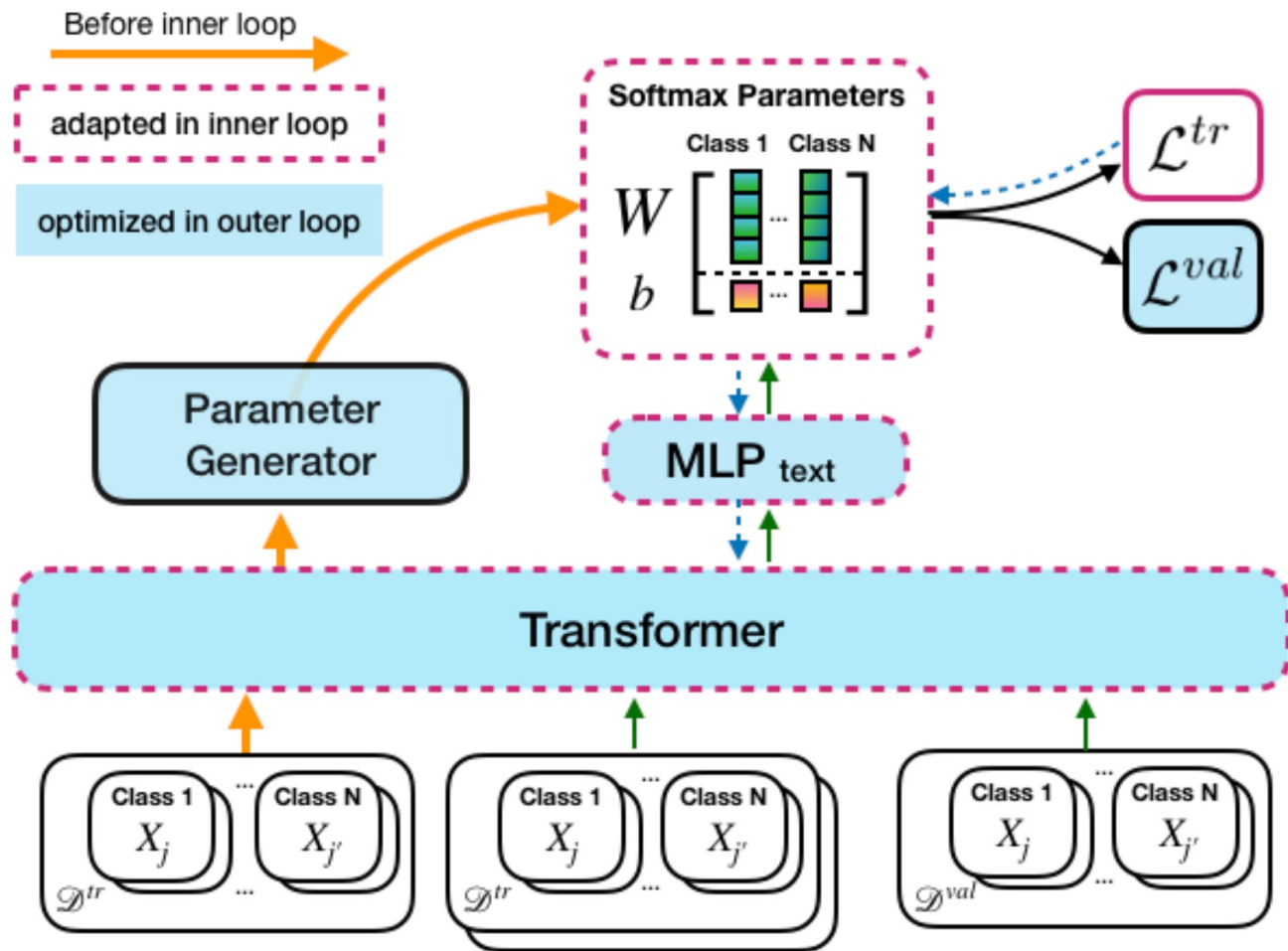
- Pre-trained transformer: fine-tuning on a new task still **requires large amounts of task-specific labeled data** to achieve good performance （标注数据耗时耗力耗钱）
- Meta-learning: limited to simulated problems or problems with limited diversity across tasks （任务缺乏多样性）
- To enable optimization-based meta-learning across tasks with different number of classes （不同的任务class数不同）
 - MAML: 只考虑了固定的class数
 - Matric-based: 若labeled examples增加, 不能适应更大的数据

Method: LEOPARD

$$w_i^n, b_i^n = \frac{1}{|C_i^n|} \sum_{x_j \in C_i^n} g_\psi(f_\theta(\mathbf{x}_j))$$

$$\mathbf{W}_i = [w_i^1; \dots; w_i^{N_i}] \quad \mathbf{b}_i = [b_i^1; \dots; b_i^{N_i}]$$

$$p(y|\mathbf{x}^*) = \text{softmax} \{ \mathbf{W}_i h_\phi(f_\theta(\mathbf{x}^*)) + \mathbf{b}_i \}$$



Method: LEOPARD

Algorithm 1 LEOPARD

Require: set of M training tasks and losses $\{(T_1, L_1), \dots, (T_M, L_M)\}$, model parameters $\Theta = \{\theta, \psi, \alpha\}$, hyper-parameters ν, G, β

Initialize θ with pre-trained BERT-base;

```
1: while not converged do
2:   # sample batch of tasks
3:   for all  $T_i \in T$  do
4:      $\mathcal{D}_i^{tr} \sim T_i$       # sample a batch of train data
5:      $C_i^n \leftarrow \{x_j | y_j = n\}$     # partition data according to class labels
6:      $w_i^n, b_i^n \leftarrow \frac{1}{|C_i^n|} \sum_{x_j \in C_i^n} g_\psi(f_\theta(x_j))$     # generate softmax parameters
7:      $\mathbf{W}_i \leftarrow [w_i^1; \dots; w_i^{N_i}]; \mathbf{b}_i \leftarrow [b_i^1; \dots; b_i^{N_i}]$ 
8:      $\Phi_i^{(0)} \leftarrow \theta_{>\nu} \cup \{\phi, \mathbf{W}_i, \mathbf{b}_i\}$     # task-specific parameters
9:     for  $s := 0 \dots G - 1$  do
10:       $\mathcal{D}_i^{tr} \sim T_i$       # sample a batch of train data
11:       $\Phi_i^{(s+1)} \leftarrow \Phi_i^{(s)} - \alpha_s \nabla_{\Phi} \mathcal{L}_i(\{\Theta, \Phi_i\}, \mathcal{D}_i^{tr})$     # adapt task-specific parameters
12:    end for
13:     $\mathcal{D}_i^{val} \sim T_i$     # sample a batch of validation data
14:     $g_i \leftarrow \nabla_{\Theta} \mathcal{L}_i(\{\Theta, \Phi_i^{(G)}\}, \mathcal{D}_i^{val})$     # gradient of task-agnostic parameters on validation
15:  end for
16:   $\Theta \leftarrow \Theta - \beta \cdot \sum_i g_i$     # optimize task-agnostic parameters
17: end while
```

Task-agnostic param:

$$\Theta = \theta_{\leq \nu} \cup \{\psi\}$$

Task-specific param:

$$\Phi_i = \theta_{>\nu} \cup \{\phi, \mathbf{W}_i, \mathbf{b}_i\}$$

Experiments: setting

- Training Tasks
 - GLUE benchmark tasks:
MNLI (m/mm), SST-2, QNLI, QQP, MRPC, RTE, and the SNLI
- Evaluation
 - 17 target NLP tasks
 - Generalization Beyond Training Tasks (new tasks unseen at training)
 - Domain transfer tasks (new domains of tasks seen at training time)
 - Evaluated on the entire test set for the task.

Experiments: Generalization Beyond Training Tasks

Entity Typing							
	N	k	BERT _{base}	MT-BERT _{softmax}	MT-BERT	Proto-BERT	LEOPARD
CoNLL	4	4	50.44 \pm 08.57	52.28 \pm 4.06	55.63 \pm 4.99	32.23 \pm 5.10	54.16 \pm 6.32
		8	50.06 \pm 11.30	65.34 \pm 7.12	58.32 \pm 3.77	34.49 \pm 5.15	67.38 \pm 4.33
		16	74.47 \pm 03.10	71.67 \pm 3.03	71.29 \pm 3.30	33.75 \pm 6.05	76.37 \pm 3.08
MITR	8	4	49.37 \pm 4.28	45.52 \pm 5.90	50.49 \pm 4.40	17.36 \pm 2.75	49.84 \pm 3.31
		8	49.38 \pm 7.76	58.19 \pm 2.65	58.01 \pm 3.54	18.70 \pm 2.38	62.99 \pm 3.28
		16	69.24 \pm 3.68	66.09 \pm 2.24	66.16 \pm 3.46	16.41 \pm 1.87	70.44 \pm 2.89
Text Classification							
Airline	3	4	42.76 \pm 13.50	43.73 \pm 7.86	46.29 \pm 12.26	40.27 \pm 8.19	54.95 \pm 11.81
		8	38.00 \pm 17.06	52.39 \pm 3.97	49.81 \pm 10.86	51.16 \pm 7.60	61.44 \pm 03.90
		16	58.01 \pm 08.23	58.79 \pm 2.97	57.25 \pm 09.90	48.73 \pm 6.79	62.15 \pm 05.56
Disaster	2	4	55.73 \pm 10.29	52.87 \pm 6.16	50.61 \pm 8.33	50.87 \pm 1.12	51.45 \pm 4.25
		8	56.31 \pm 09.57	56.08 \pm 7.48	54.93 \pm 7.88	51.30 \pm 2.30	55.96 \pm 3.58
		16	64.52 \pm 08.93	65.83 \pm 4.19	60.70 \pm 6.05	52.76 \pm 2.92	61.32 \pm 2.83
Emotion	13	4	09.20 \pm 3.22	09.41 \pm 2.10	09.84 \pm 2.14	09.18 \pm 3.14	11.71 \pm 2.16
		8	08.21 \pm 2.12	11.61 \pm 2.34	11.21 \pm 2.11	11.18 \pm 2.95	12.90 \pm 1.63
		16	13.43 \pm 2.51	13.82 \pm 2.02	12.75 \pm 2.04	12.32 \pm 3.73	13.38 \pm 2.20
Political Bias	2	4	54.57 \pm 5.02	54.32 \pm 3.90	54.66 \pm 3.74	56.33 \pm 4.37	60.49 \pm 6.66
		8	56.15 \pm 3.75	57.36 \pm 4.32	54.79 \pm 4.19	58.87 \pm 3.79	61.74 \pm 6.73
		16	60.96 \pm 4.25	59.24 \pm 4.25	60.30 \pm 3.26	57.01 \pm 4.44	65.08 \pm 2.14

	N	k	BERT _{base}	MT-BERT _{softmax}	MT-BERT	Proto-BERT	LEOPARD
Political Audience	2	4	51.89 \pm 1.72	51.50 \pm 2.72	51.53 \pm 1.80	51.47 \pm 3.68	52.60 \pm 3.51
		8	52.80 \pm 2.72	53.53 \pm 2.26	54.34 \pm 2.88	51.83 \pm 3.77	54.31 \pm 3.95
		16	58.45 \pm 4.98	56.37 \pm 2.19	55.14 \pm 4.57	53.53 \pm 3.25	57.71 \pm 3.52
Political Message	9	4	15.64 \pm 2.73	13.71 \pm 1.10	14.49 \pm 1.75	14.22 \pm 1.25	15.69 \pm 1.57
		8	13.38 \pm 1.74	14.33 \pm 1.32	15.24 \pm 2.81	15.67 \pm 1.96	18.02 \pm 2.32
		16	20.67 \pm 3.89	18.11 \pm 1.48	19.20 \pm 2.20	16.49 \pm 1.96	18.07 \pm 2.41
Rating Books	3	4	39.42 \pm 07.22	44.82 \pm 9.00	38.97 \pm 13.27	48.44 \pm 7.43	54.92 \pm 6.18
		8	39.55 \pm 10.01	51.14 \pm 6.78	46.77 \pm 14.12	52.13 \pm 4.79	59.16 \pm 4.13
		16	43.08 \pm 11.78	54.61 \pm 6.79	51.68 \pm 11.27	57.28 \pm 4.57	61.02 \pm 4.19
Rating DVD	3	4	32.22 \pm 08.72	45.94 \pm 7.48	41.23 \pm 10.98	47.73 \pm 6.20	49.76 \pm 9.80
		8	36.35 \pm 12.50	46.23 \pm 6.03	45.24 \pm 9.76	47.11 \pm 4.00	53.28 \pm 4.66
		16	42.79 \pm 10.18	49.23 \pm 6.68	45.19 \pm 11.56	48.39 \pm 3.74	53.52 \pm 4.77
Rating Electronics	3	4	39.27 \pm 10.15	39.89 \pm 5.83	41.20 \pm 10.69	37.40 \pm 3.72	51.71 \pm 7.20
		8	28.74 \pm 08.22	46.53 \pm 5.44	45.41 \pm 09.49	43.64 \pm 7.31	54.78 \pm 6.48
		16	45.48 \pm 06.13	48.71 \pm 6.16	47.29 \pm 10.55	44.83 \pm 5.96	58.69 \pm 2.41
Rating Kitchen	3	4	34.76 \pm 11.20	40.41 \pm 5.33	36.77 \pm 10.62	44.72 \pm 9.13	50.21 \pm 09.63
		8	34.49 \pm 08.72	48.35 \pm 7.87	47.98 \pm 09.73	46.03 \pm 8.57	53.72 \pm 10.31
		16	47.94 \pm 08.28	52.94 \pm 7.14	53.79 \pm 09.47	49.85 \pm 9.31	57.00 \pm 08.69
Overall Average		4	38.13	40.13	40.10	36.29	45.99
		8	36.99	45.89	44.25	39.15	50.86
		16	48.55	49.93	49.07	39.85	55.50

Experiments: Few-Shot Domain Transfer

Natural Language Inference							
	k	BERT _{base}	MT-BERT _{softmax}	MT-BERT	MT-BERT _{reuse}	Proto-BERT	LEOPARD
Scitail	4	58.53 \pm 09.74	74.35 \pm 5.86	63.97 \pm 14.36	76.65 \pm 2.45	76.27 \pm 4.26	69.50 \pm 9.56
	8	57.93 \pm 10.70	79.11 \pm 3.11	68.24 \pm 10.33	76.86 \pm 2.09	78.27 \pm 0.98	75.00 \pm 2.42
	16	65.66 \pm 06.82	79.60 \pm 2.31	75.35 \pm 04.80	79.53 \pm 2.17	78.59 \pm 0.48	77.03 \pm 1.82
Amazon Review Sentiment Classification							
Books	4	54.81 \pm 3.75	68.69 \pm 5.21	64.93 \pm 8.65	74.79 \pm 6.91	73.15 \pm 5.85	82.54 \pm 1.33
	8	53.54 \pm 5.17	74.86 \pm 2.17	67.38 \pm 9.78	78.21 \pm 3.49	75.46 \pm 6.87	83.03 \pm 1.28
	16	65.56 \pm 4.12	74.88 \pm 4.34	69.65 \pm 8.94	78.87 \pm 3.32	77.26 \pm 3.27	83.33 \pm 0.79
Kitchen	4	56.93 \pm 7.10	63.07 \pm 7.80	60.53 \pm 9.25	75.40 \pm 6.27	62.71 \pm 9.53	78.35 \pm 18.36
	8	57.13 \pm 6.60	68.38 \pm 4.47	69.66 \pm 8.05	75.13 \pm 7.22	70.19 \pm 6.42	84.88 \pm 01.12
	16	68.88 \pm 3.39	75.17 \pm 4.57	77.37 \pm 6.74	80.88 \pm 1.60	71.83 \pm 5.94	85.27 \pm 01.31

Ablation Study:

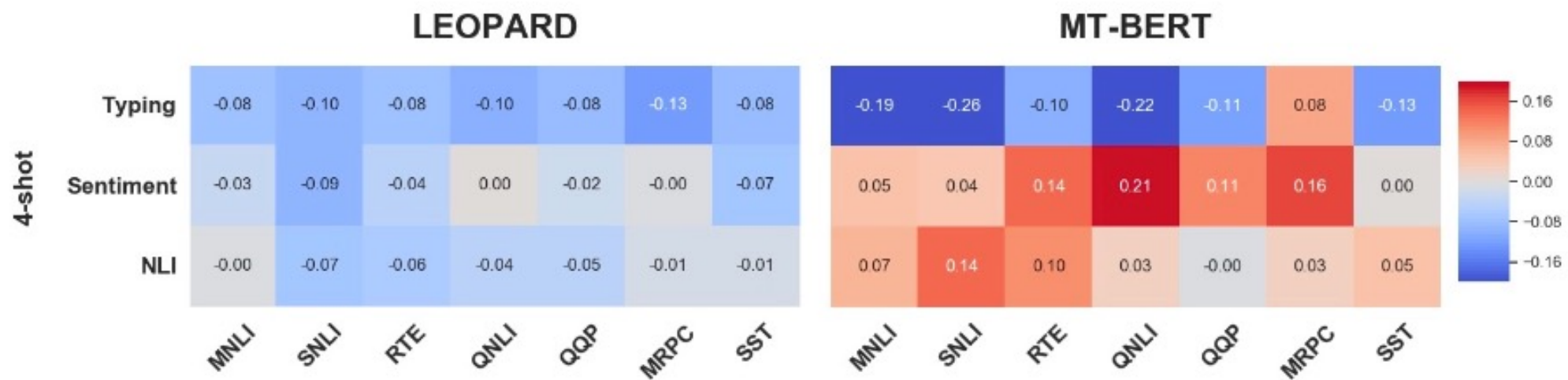
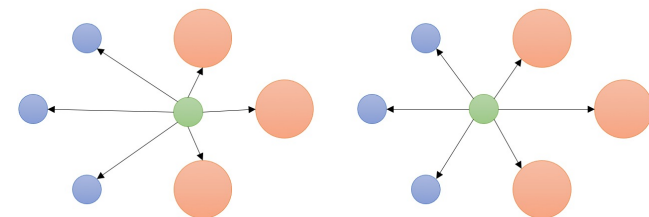
- Importance of softmax parameters
- Parameter efficiency

k	Model	Entity Typing	Sentiment Classification	NLI
16	LEOPARD ₁₀	37.62 \pm 7.37	58.10 \pm 5.40	78.53 \pm 1.55
	LEOPARD ₅	62.49 \pm 4.23	71.50 \pm 5.93	73.27 \pm 2.63
	LEOPARD	69.00 \pm 4.76	76.65 \pm 2.47	76.10 \pm 2.21
	LEOPARD-ZERO	44.79 \pm 9.34	74.45 \pm 3.34	74.36 \pm 6.67

MAML →

Table 3: Ablations: LEOPARD _{ν} does not adapt layers 0– ν (inclusive) in the inner loop (and fine-tuning), while LEOPARD adapts all parameters. Note that the outer loop still optimizes all parameters. For new tasks (like entity typing) adapting all parameters is better while for tasks seen at training time (like NLI) adapting fewer parameters is better. LEOPARD-ZERO is a model trained without the softmax-generator and a zero initialized softmax classifier, which shows the importance of softmax generator in LEOPARD.

Ablation Study: Importance of training tasks



- 泛化能力，受training tasks影响
- Initial point经过很少的training samples后可以迅速适应新task

Reference

- http://speech.ee.ntu.edu.tw/~tlkagk/courses_ML20.html
- Bansal, T., Jha, R., & McCallum, A. (2019). Learning to few-shot learn across diverse natural language classification tasks. *arXiv preprint arXiv:1911.03863*.
- Zi-Yi Dou, Keyi Yu, and Antonios Anastasopoulos. 2019. Investigating meta-learning algorithms for low-resource natural language understanding tasks. In EMNLP-IJCNLP, pages 1192–1197
- Chelsea Finn, Pieter Abbeel, and Sergey Levine. 2017. Model-agnostic meta-learning for fast adaptation of deep networks. In ICML, pages 1126–1135.

Thanks!

Q&A