Meta Learning for Low-Resource Natural Language Understanding

Sanghyun Seo

Feb 14, 2020

Department of Computer Engineering at Dongguk University

Artificial Intelligence Laboratory

shseo@dongguk.edu

Papers

- MAML, FOMAML (2017)
 - Finn, Chelsea, Pieter Abbeel, and Sergey Levine. "Model-agnostic metalearning for fast adaptation of deep networks." Proceedings of the 34th International Conference on Machine Learning-Volume 70. JMLR. org, 2017.
- Reptile (2018)
 - Nichol, Alex, Joshua Achiam, and John Schulman. "On first-order metalearning algorithms." *arXiv preprint arXiv:1803.02999* (2018).
- Low-Resource NLP(2018)
 - Dou, Zi-Yi, Keyi Yu, and Antonios Anastasopoulos. "Investigating Meta-Learning Algorithms for Low-Resource Natural Language Understanding Tasks." Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP). 2019.

Survey Paper Overview

- Tentative Title
 - 소수샷 학습을 위한 메타러닝 연구동향 (A survey on meta learning algorithms for fewshot learning)
- Introduction
 - Meta learning vs Transfer learning or Multi-task Learning
- Problem Definition
 - Supervised learning → Meta Learning
 - Few-shot learning
 - Dataset
- Meta Learning Algorithms
 - Model based approaches
 - MANN, MetaNet, SNAIL
 - Metric based approaches
 - Matching Networks, Relation Networks, Prototypical Networks
 - Few-Shot Learning with Graph Neural Networks, Transductive Propagation Networks
 - Optimization based approaches
 - MAML, FOMAML, Meta-SGD, Reptile
 - MT-Net, LEO
 - Hierarchical bayesian model, probabilistic MAML, Bayesian MAML
 - Set-input approaches
 - Deep sets, Neural statistician

Meta Learning

- Two ways to view meta-learning
 - Mechanistic view
 - Deep neural network model that can read in an entire dataset and make predictions for new datapoints
 - Training this network uses a meta-dataset, which itself consists of many datasets, each for a different task
 - This view makes it easier to implement meta learning algorithms
 - Probabilistic view
 - Extract prior information from a set of (meta-training) tasks that allows efficient learning of new tasks
 - Learning a new task uses this prior and (small) training set to infer most likely posterior parameters
 - This view makes it easier to understand meta learning algorithms

ICML 2019 tutorial, https://sites.google.com/view/icml19metalearning

Meta Learning

Limits of Supervised Learning



What is wrong with this?

> The most powerful models typically require large amounts of labeled data

Labeled data for some tasks may be very limited

ICML 2019 tutorial, <u>https://sites.google.com/view/icml19metalearning</u>

Meta Learning

- Additional data
 - Unsupervised learning
 - Transfer learning
 - Multi-task learning
 - Meta learning

supervised learning:

 $rg\max_{\phi}\log p(\phi|\mathcal{D})$

can we incorporate *additional* data?

$$\arg\max_{\phi} \log p(\phi | \mathcal{D}, \mathcal{D}_{\text{meta-train}})$$

$$\mathcal{D} = \{(x_1, y_1), \dots, (x_k, y_k)\}$$

$$\mathcal{D}_{\text{meta-train}} = \{\mathcal{D}_1, \dots, \mathcal{D}_n\}$$

 $\mathcal{D}_i = \{(x_1^i, y_1^i), \dots, (x_k^i, y_k^i)\}$



 $\mathcal{D}_{ ext{meta-train}}$

ICML 2019 tutorial, https://sites.google.com/view/icml19metalearning

- How MAML can learn prior knowledge from meta train dataset?
- Procedure of MAML algorithms



- **Require:** $p(\mathcal{T})$: distribution over tasks
- **Require:** α , β : step size hyperparameters
- 1: randomly initialize θ
- 2: while not done do
- 3: Sample batch of tasks $\mathcal{T}_i \sim p(\mathcal{T})$
- 4: for all \mathcal{T}_i do
- 5: Sample K datapoints $\mathcal{D} = {\mathbf{x}^{(j)}, \mathbf{y}^{(j)}}$ from \mathcal{T}_i
- 6: Evaluate $\nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$ using \mathcal{D} and $\mathcal{L}_{\mathcal{T}_i}$ in Equation (2) or (3)
- 7: Compute adapted parameters with gradient descent: $\theta'_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$
- 8: Sample datapoints $\mathcal{D}'_i = {\mathbf{x}^{(j)}, \mathbf{y}^{(j)}}$ from \mathcal{T}_i for the meta-update
- 9: end for
- 10: Update $\theta \leftarrow \theta \beta \nabla_{\theta} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i})$ using each \mathcal{D}'_i and $\mathcal{L}_{\mathcal{T}_i}$ in Equation 2 or 3
- 11: end while

Finn, Chelsea, Pieter Abbeel, and Sergey Levine. "Model-agnostic meta-learning for fast adaptation of deep networks." Proceedings of the 34th International Conference on Machine Learning-Volume 70. JMLR. org, 2017.

- Meta Train (Meta Learning)
- Randomly initialize θ





- **Require:** $p(\mathcal{T})$: distribution over tasks
- **Require:** α , β : step size hyperparameters
- 1: randomly initialize θ
- 2: while not done do
- 3: Sample batch of tasks $\mathcal{T}_i \sim p(\mathcal{T})$
- 4: for all \mathcal{T}_i do
- 5: Sample K datapoints $\mathcal{D} = {\mathbf{x}^{(j)}, \mathbf{y}^{(j)}}$ from \mathcal{T}_i
- 6: Evaluate $\nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$ using \mathcal{D} and $\mathcal{L}_{\mathcal{T}_i}$ in Equation (2) or (3)
- 7: Compute adapted parameters with gradient descent: $\theta'_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$
- 8: Sample datapoints $\mathcal{D}'_i = {\mathbf{x}^{(j)}, \mathbf{y}^{(j)}}$ from \mathcal{T}_i for the meta-update
- 9: end for
- 10: Update $\theta \leftarrow \theta \beta \nabla_{\theta} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i})$ using each \mathcal{D}'_i and $\mathcal{L}_{\mathcal{T}_i}$ in Equation 2 or 3
- 11: end while

• Sample batch of tasks $T_i \sim p(T)$





- **Require:** $p(\mathcal{T})$: distribution over tasks
- **Require:** α , β : step size hyperparameters
- 1: randomly initialize θ
- 2: while not done do
- 3: Sample batch of tasks $\mathcal{T}_i \sim p(\mathcal{T})$
- 4: for all \mathcal{T}_i do
- 5: Sample K datapoints $\mathcal{D} = {\mathbf{x}^{(j)}, \mathbf{y}^{(j)}}$ from \mathcal{T}_i
- 6: Evaluate $\nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$ using \mathcal{D} and $\mathcal{L}_{\mathcal{T}_i}$ in Equation (2) or (3)
- 7: Compute adapted parameters with gradient descent: $\theta'_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$
- 8: Sample datapoints $\mathcal{D}'_i = {\mathbf{x}^{(j)}, \mathbf{y}^{(j)}}$ from \mathcal{T}_i for the meta-update
- 9: end for
- 10: Update $\theta \leftarrow \theta \beta \nabla_{\theta} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i})$ using each \mathcal{D}'_i and $\mathcal{L}_{\mathcal{T}_i}$ in Equation 2 or 3
- 11: end while





• Regression tasks: MSE

(2)
$$\mathcal{L}_{T_i}(f_{\phi}) = \sum_{x^{(j)}, y^{(j)} \sim T_i} \left\| f_{\phi}(x^{(j)}) - y^{(j)} \right\|_2^2$$

• Classification tasks: CE
(3)
$$\mathcal{L}_{T_i}(f_{\phi}) = \sum_{x^{(j)}, y^{(j)} \sim T_i}^{y^{(j)} \log f_{\phi}(x^{(j)})} \log(1 - f_{\phi}(x^{(j)}))$$



- **Require:** $p(\mathcal{T})$: distribution over tasks
- **Require:** α , β : step size hyperparameters
- 1: randomly initialize θ
- 2: while not done do
- 3: Sample batch of tasks $\mathcal{T}_i \sim p(\mathcal{T})$
- 4: for all \mathcal{T}_i do
- 5: Sample K datapoints $\mathcal{D} = {\mathbf{x}^{(j)}, \mathbf{y}^{(j)}}$ from \mathcal{T}_i
- 6: Evaluate $\nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$ using \mathcal{D} and $\mathcal{L}_{\mathcal{T}_i}$ in Equation (2) or (3)
- 7: Compute adapted parameters with gradient descent: $\frac{\theta'_{i}}{\theta_{i}} = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_{i}}(f_{\theta})$
- 8: Sample datapoints $\mathcal{D}'_i = {\mathbf{x}^{(j)}, \mathbf{y}^{(j)}}$ from \mathcal{T}_i for the meta-update
- 9: end for
- 10: Update $\theta \leftarrow \theta \beta \nabla_{\theta} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i})$ using each \mathcal{D}'_i and $\mathcal{L}_{\mathcal{T}_i}$ in Equation 2 or 3
- 11: end while



- Why MAML uses D'_i and θ'_i ?
 - D_i and D'_i are disjoint

Note that the meta-optimization is performed over the model parameters θ , whereas the objective is computed using the updated model parameters θ' . In effect, our proposed method aims to optimize the model parameters such that one or a small number of gradient steps on a new task will produce maximally effective behavior on that task.

 $\beta \nabla_{\theta} \sum \mathcal{L}_{T_i}(f_{\theta'_i})$





Algorithm 2 MAML for Few-Shot Supervised Learning

Require: $p(\mathcal{T})$: distribution over tasks

- **Require:** α , β : step size hyperparameters
- 1: randomly initialize θ
- 2: while not done do
- 3: Sample batch of tasks $\mathcal{T}_i \sim p(\mathcal{T})$
- 4: for all \mathcal{T}_i do
- 5: Sample K datapoints $\mathcal{D} = {\mathbf{x}^{(j)}, \mathbf{y}^{(j)}}$ from \mathcal{T}_i
- 6: Evaluate $\nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$ using \mathcal{D} and $\mathcal{L}_{\mathcal{T}_i}$ in Equation (2) or (3)
- 7: Compute adapted parameters with gradient descent: $\theta'_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$
- 8: Sample datapoints $\mathcal{D}'_i = {\mathbf{x}^{(j)}, \mathbf{y}^{(j)}}$ from \mathcal{T}_i for the meta-update
- 9: end for
- 10: Update $\theta \leftarrow \theta \beta \nabla_{\theta} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i})$ using each \mathcal{D}'_i and $\mathcal{L}_{\mathcal{T}_i}$ in Equation 2 or 3
- 11: end while

• Meta update

$$D'_{i} = \{x^{(j)}, y^{(j)}\}$$
training data
$$D'_{i} = \{x^{(j)}, y^{(j)}\}$$
test set
$$meta-training$$

$$\square \square \square \square \square \square \square \square \square$$

$$\square \square \square \square \square \square \square$$

$$\square \square \square \square \square \square \square \square$$

$$\square \square \square \square \square \square \square \square$$

$$\square \square \square \square \square \square \square \square \square$$

$$\square \square \square \square \square \square \square \square \square$$

$$\square \square \square \square \square \square \square \square \square \square$$

$$\square \square \square$$

$$\min_{\theta} \sum_{T_i} \mathcal{L}(\theta - \alpha \nabla_{\theta} \mathcal{L}(\theta, D_i), D'_i)$$



Algorithm 2 MAML for Few-Shot Supervised Learning

Require: $p(\mathcal{T})$: distribution over tasks

Require: α , β : step size hyperparameters

- 1: randomly initialize θ
- 2: while not done do
- 3: Sample batch of tasks $\mathcal{T}_i \sim p(\mathcal{T})$
- 4: for all \mathcal{T}_i do
- 5: Sample K datapoints $\mathcal{D} = {\mathbf{x}^{(j)}, \mathbf{y}^{(j)}}$ from \mathcal{T}_i
- 6: Evaluate $\nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$ using \mathcal{D} and $\mathcal{L}_{\mathcal{T}_i}$ in Equation (2) or (3)
- 7: Compute adapted parameters with gradient descent: $\theta'_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$
- 8: Sample datapoints $\mathcal{D}'_i = {\mathbf{x}^{(j)}, \mathbf{y}^{(j)}}$ from \mathcal{T}_i for the meta-update
- 9: end for
- 10: Update $\theta \leftarrow \theta \beta \nabla_{\theta} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i})$ using each \mathcal{D}'_i and $\mathcal{L}_{\mathcal{T}_i}$ in Equation 2 or 3
- 11: end while

• Meta Test (Learning→Evaluation)





Expected high performance from few-shot training dataset

 θ^{*}



- Intuitive example
 - Can we meet sweet examples(new task) every time?

Key idea "our training procedure is based on a simple machine learning principle: test and train conditions must match"

Meta train

Meta test



- Computational Problem
 - Consider the case of k>1 inner gradient step

$$\begin{split} \theta_0 &= \theta_{meta} \\ \theta_1 &= \theta_0 - \alpha \nabla_{\theta} \mathcal{L}^{(0)}(\theta_0) \\ \theta_2 &= \theta_1 - \alpha \nabla_{\theta} \mathcal{L}^{(0)}(\theta_1) \\ & \dots \\ \theta_k &= \theta_{k-1} - \alpha \nabla_{\theta} \mathcal{L}^{(0)}(\theta_{k-1}) \end{split}$$



$$\theta_0 = \theta_{meta} = \theta_{meta} - \beta \nabla_{\theta} \mathcal{L}^{(1)}(\theta_k)$$

- **Require:** $p(\mathcal{T})$: distribution over tasks
- **Require:** α , β : step size hyperparameters
- 1: randomly initialize θ
- 2: while not done do
- 3: Sample batch of tasks $\mathcal{T}_i \sim p(\mathcal{T})$
- 4: for all \mathcal{T}_i do
- 5: Sample K datapoints $\mathcal{D} = {\mathbf{x}^{(j)}, \mathbf{y}^{(j)}}$ from \mathcal{T}_i
- 6: Evaluate $\nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$ using \mathcal{D} and $\mathcal{L}_{\mathcal{T}_i}$ in Equation (2) or (3)
- 7: Compute adapted parameters with gradient descent: $\frac{\theta'_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})}{\leftarrow} \text{ What if K iteration?}$
- 8: Sample datapoints $\mathcal{D}'_i = {\mathbf{x}^{(j)}, \mathbf{y}^{(j)}}$ from \mathcal{T}_i for the meta-update
- 9: end for
- 10: Update $\theta \leftarrow \theta \beta \nabla_{\theta} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i})$ using each \mathcal{D}'_i and $\mathcal{L}_{\mathcal{T}_i}$ in Equation 2 or 3
- 11: end while

- **Computational Problem** ٠
 - Consider the case of k>1 inner gradient step

$$\nabla_{\theta} \mathcal{L}^{(1)}(\theta_{k}) = \frac{\partial \mathcal{L}^{(1)}(\theta_{k})}{\partial \theta} = \frac{\partial \mathcal{L}^{(1)}(\theta_{k})}{\partial \theta_{k}} \cdot \frac{\partial \theta_{k}}{\partial \theta}$$

$$= \frac{\partial \mathcal{L}^{(1)}(\theta_{k})}{\partial \theta_{k}} \cdot \prod_{i=1}^{k} \frac{\partial \theta_{i}}{\partial \theta_{i-1}}^{T}$$

$$= \nabla_{\theta_{k}} \mathcal{L}^{(1)}(\theta_{k}) \cdot \nabla_{\theta_{k-1}} \mathcal{L}(\theta_{k}) \cdots \nabla_{\theta_{0}} \mathcal{L}(\theta_{1}) \cdot \nabla_{\theta} \theta_{0}$$

$$= \nabla_{\theta_{k}} \mathcal{L}^{(1)}(\theta_{k}) \cdot \prod_{i=1}^{k} \nabla_{\theta_{i-1}} \mathcal{L}(\theta_{i})$$

$$= \nabla_{\theta_{k}} \mathcal{L}^{(1)}(\theta_{k}) \cdot \prod_{i=1}^{k} \nabla_{\theta_{i-1}} \left(\theta_{i-1} - \alpha \nabla_{\theta} \mathcal{L}^{(0)}(\theta_{i-1}) \right)$$

$$= \nabla_{\theta_{k}} \mathcal{L}^{(1)}(\theta_{k}) \cdot \prod_{i=1}^{k} \left(I - \alpha \nabla_{\theta_{i-1}} \left(\nabla_{\theta} \mathcal{L}^{(0)}(\theta_{i-1}) \right) \right)$$

$$\Rightarrow \text{ second order derivative problem}$$

$$(I - \alpha \nabla_{\theta_{i}} \mathcal{L}^{(1)}(\theta_{k}))$$

$$(I - \alpha \nabla_{\theta_{i}=1} (\nabla_{\theta} \mathcal{L}^{(0)}(\theta_{i}=1)))$$

$$(I) \rightarrow \text{ second order derivative problem}$$

FOMAML

• First-Order MAML ignores the second derivative part



Reptile

- The Reptile works by repeatedly:
 - 1) sampling a task,
 - 2) training on it by multiple gradient descent steps,
 - 3) and then moving the model weights towards the new parameters.

Algorithm 1 Reptile (serial version)

```
Initialize \phi, the vector of initial parameters
for iteration = 1, 2, ... do
Sample task \tau, corresponding to loss L_{\tau} on weight vectors \tilde{\phi}
Compute \tilde{\phi} = U_{\tau}^{k}(\phi), denoting k steps of SGD or Adam
Update \phi \leftarrow \phi + \epsilon(\tilde{\phi} - \phi)
end for
```

U corresponds to performing gradient descent

Nichol, Alex, Joshua Achiam, and John Schulman. "On First-Order Meta-Learning Algorithms." arXiv preprint arXiv:1803.02999 (2018).

Reptile

- Intuitive example
 - Assuming that a task $\tau \sim p(\tau)$ has a manifold of optimal network configuration, W_{τ}^* . The model f_{θ} achieves the best performance for task τ when θ lays on the surface of W_{τ}^*
 - To find a solution that is good across tasks, we would like to find a parameter close to all the optimal manifolds of all tasks



Nichol, Alex, Joshua Achiam, and John Schulman. "On First-Order Meta-Learning Algorithms." arXiv preprint arXiv:1803.02999 (2018).

MAML vs FOMAML vs Reptile

Can you know the difference between each methodology in the following figure?



$$\mathbb{E}\left[g_{\text{FOMAML}}\right] = (1)\text{AvgGrad} - (\alpha)\text{AvgGradInner} + O(\alpha^2)$$
(34)

$$\mathbb{E}\left[g_{\text{Reptile}}\right] = (2)\text{AvgGrad} - (\alpha)\text{AvgGradInner} + O(\alpha^2)$$
(35)

- <u>https://lilianweng.github.io/lil-log/2018/11/30/meta-learning.html</u>, slides on reptile by Yoonho Lee
- Nichol, Alex, Joshua Achiam, and John Schulman. "On First-Order Meta-Learning Algorithms." arXiv preprint arXiv:1803.02999 (2018).

Meta learning for NLP

- Meta learning for low-resource NLU
 - Meta learning can make low-resource NLP model
 - Meta learning for Machine translation
- GLUE
 - The General Language Understanding Evaluation benchmark (https://gluebenchmark.com/) is a collection of resources for training, evaluating, and analyzing natural language understanding systems.



- Gu, Jiatao, et al. "Meta-Learning for Low-Resource Neural Machine Translation." Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing. 2018.
- Dou, Zi-Yi, Keyi Yu, and Antonios Anastasopoulos. "Investigating Meta-Learning Algorithms for Low-Resource Natural Language Understanding Tasks." Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP). 2019.
- https://www.tensorflow.org/datasets/catalog/glue

from the community question-answering website Quora. The task is to determine whether a pair of

questions are semantically equivalent.

Meta learning for NLP

Meta learning for low-resource NLU

High-resource tasks		Target tasks 1		
•	SST-2: The Stanford Sentiment Treebank consists of	•	CoLA: The Corpus of Linguistic Acceptability	
	sentences from movie reviews and human		consists of English acceptability judgments	
	annotations of their sentiment	•	MRPC: The Microsoft Research Paraphrase Corpus	
•	QQP: The Quora Question Pairs2 dataset	•	STS-B: The Semantic Textual Similarity Benchmark	
•	MNLI: The Multi-Genre Natural Language Inference	•	RTE: The Recognizing Textual Entailment	
	Corpusn	Target tasks 2		
•	QNLI: The Stanford Question Answering Dataset	•	SciTail: entailment dataset created from multiple- choice science exams and web sentences.	

Comparison model



- Liu, Xiaodong, et al. "Multi-Task Deep Neural Networks for Natural Language Understanding." Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics. 2019.
- Devlin, Jacob, et al. "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding." Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers). 2019.

Meta learning for NLP

• Experiment Results

- Uniform, Probability Proportional to Size(PPS), Mixed

Model	Test Dataset						
wiodei	CoLA	MRPC	STS-B	RTE			
BERT	52.1	88.9/84.8	87.1/85.8	66.4			
MT-DNN	51.7	89.9/86.3	87.6/86.8	75.4			
MAML	53.4	89.5/85.8	88.0/87.3	76.4			
FOMAML	51.6	89.9/86.4	88.6/88.0	74.1			
Reptile	53.2	90.2/86.7	88.7/88.1	77.0			

Table 1: Results on GLUE test sets. Metrics differ per task (explained in Appendix A) but the best result is **highlighted**.

Model	#Upt	α	CoLA	MRPC	STS-B	RTE
	3	1e-3	60.7	89.7	90.2	77.9
		1e-4	62.0	88.0	90.1	81.2
Reptile	5	1e-3	61.6	90.0	90.3	83.0
		1e-2	60.1	87.8	89.5	73.9
	7	1e-3	57.8	88.7	90.0	81.4

Table 3: Effect of the number of update steps and the inner learning rate α .

Model	CoLA	MRPC	STS-B	RTE
Reptile-PPS	61.6	90.0	90.3	83.0
Reptile-Uniform	61.5	84.0	90.3	75.7
Reptile-Mixed 2:1	60.3	87.8	90.3	71.0
Reptile-Mixed 5:1	61.6	85.8	90.1	74.7

Table 2: Effect of task distributions. We report the accuracy or Matthews correlation on development sets.



Figure 2: Results on transfer learning. The target task is SciTail which the model does not come across during the meta-learning stage.

Dou, Zi-Yi, Keyi Yu, and Antonios Anastasopoulos. "Investigating Meta-Learning Algorithms for Low-Resource Natural Language Understanding Tasks." Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP). 2019.

Q&A Thank you!