

# Meta-Learning for Domain Generalization in Semantic Parsing

Bailin Wang<sup>†</sup>, Mirella Lapata<sup>†</sup> and Ivan Titov<sup>†§</sup>

<sup>†</sup> ILCC, University of Edinburgh, <sup>§</sup> ILLC, University of Amsterdam



UNIVERSITY  
OF AMSTERDAM

# Semantic Parsing for Databases



database: concert singer



Show all *countries* and the number of *singers* in each *country*.



```
SELECT Country , count(*) FROM Singer GROUP BY Country
```

**Task:** *translating natural language utterance to SQL queries.*

# Cross-Domain Text-to-SQL Parsing



database: concert singer

*Train*



database: farm

*Test*



## Domain Generalization

\* a parser needs to generalize to *unseen* domains.

\* modern parsers have *a gap of more than 25%* between in- and cross-domain performance

# Cross-Lingual Cross-Domain Text-to-SQL Parsing



database: concert singer

*Train*



每个 **国家** 有多少 **歌手**



```
SELECT Country , count(*) FROM Singer GROUP BY Country
```



database: farm

*Test*



请显示不同 **城市** 的 **地位** 和 各个 **地位** 的城市 平均 **人口**。



```
SELECT Status , avg(Population) FROM City GROUP BY Status
```

- *Utterance and database schemas are in different languages*
- In the left figure, utterances are in **Chinese** whereas database schemas are in **English**

# Previous work: specialized models for schema linking

## Mono-lingual Setting



Show all *countries* and the number of *singers* in each *country*.



```
SELECT Country , count(*) FROM Singer GROUP BY Country
```



## Cross-lingual Setting



每个 *国家* 有多少 *歌手*



```
SELECT Country , count(*) FROM Singer GROUP BY Country
```



# Can We Optimize for Domain Generalization *without Changing Models?*

(for both mono- and cross-lingual settings)

# Construct Virtual Tasks for Meta-Learning

# Meta-Learning Objective



1. *Meta-Train*     $\mathcal{L}_\theta(\text{green DB}, \text{orange DB})$      $\theta' \rightarrow \theta - \alpha \nabla \mathcal{L}_\theta(\text{green DB}, \text{orange DB})$



# Meta-Learning Objective



1. *Meta-Train*  $\mathcal{L}_\theta(\text{green DB}, \text{orange DB})$   $\theta' \rightarrow \theta - \alpha \nabla \mathcal{L}_\theta(\text{green DB}, \text{orange DB})$
2. *Meta-Test*  $\mathcal{L}_{\theta'}(\text{blue DB})$

# Meta-Learning Objective



1. *Meta-Train*  $\mathcal{L}_{\theta}(\text{green}, \text{orange}) \quad \theta' \rightarrow \theta - \alpha \nabla \mathcal{L}_{\theta}(\text{green}, \text{orange})$

2. *Meta-Test*  $\mathcal{L}_{\theta'}(\text{blue})$

3. *Final Loss*  $\mathcal{L}_{maml}(\theta) = \mathcal{L}_{\theta}(\text{green}, \text{orange}) + \mathcal{L}_{\theta'}(\text{blue})$

# Meta-Learning Objective: DG-MAML



1. *Meta-Train*  $\mathcal{L}_\theta(\text{green}, \text{orange}) \quad \theta' \rightarrow \theta - \alpha \nabla \mathcal{L}_\theta(\text{green}, \text{orange})$

2. *Meta-Test*  $\mathcal{L}_{\theta'}(\text{blue})$

3. *Final Loss*  $\mathcal{L}_{maml}(\theta) = \mathcal{L}_\theta(\text{green}, \text{orange}) + \mathcal{L}_{\theta'}(\text{blue})$

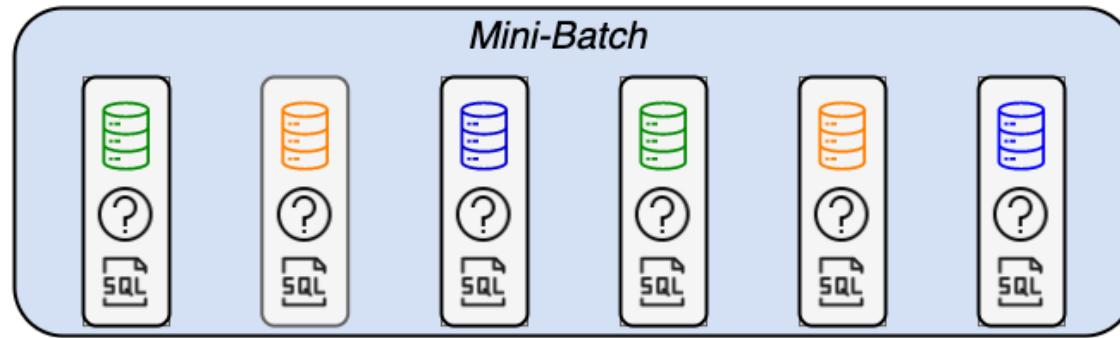
# Analysis of DG-MAML

$$\begin{aligned}\mathcal{L}_{maml}(\theta) &= \mathcal{L}_{\theta}(\text{green DB}, \text{orange DB}) + \mathcal{L}_{\theta'}(\text{blue DB}) \\ &\approx \mathcal{L}_{\theta}(\text{green DB}, \text{orange DB}) + \mathcal{L}_{\theta}(\text{blue DB}) - \alpha \nabla \mathcal{L}_{\theta}(\text{green DB}, \text{orange DB}) \cdot \nabla \mathcal{L}_{\theta}(\text{blue DB})\end{aligned}$$

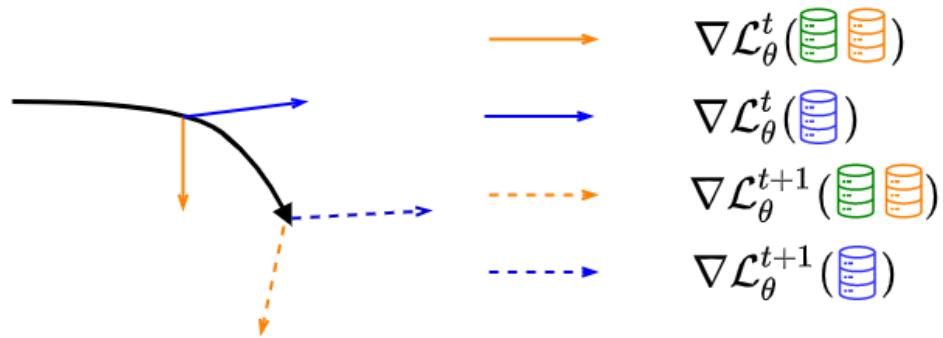
# Gradient Updates of DG-MAML



# Supervised Learning



$$\begin{aligned}\mathcal{L}_{sup}(\theta) &= \mathcal{L}_{\theta}(\text{green, orange, blue}) \\ &= \mathcal{L}_{\theta}(\text{green, orange}) + \mathcal{L}_{\theta}(\text{blue})\end{aligned}$$



# Experiments

# Datasets

*Cross-Domain*

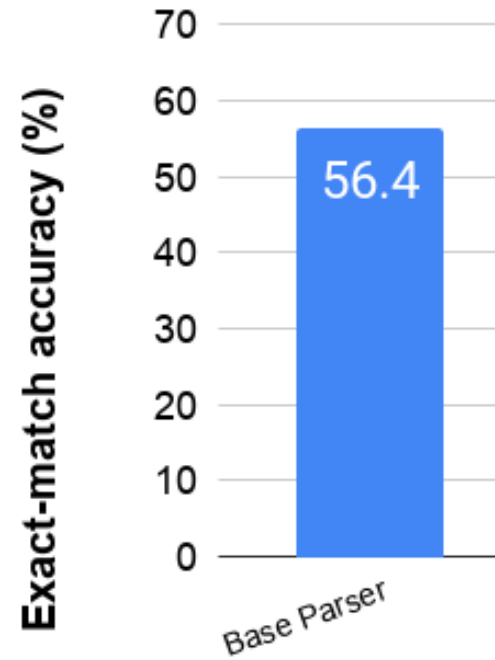
*Cross-Lingual*

Spider

Chinese Spider

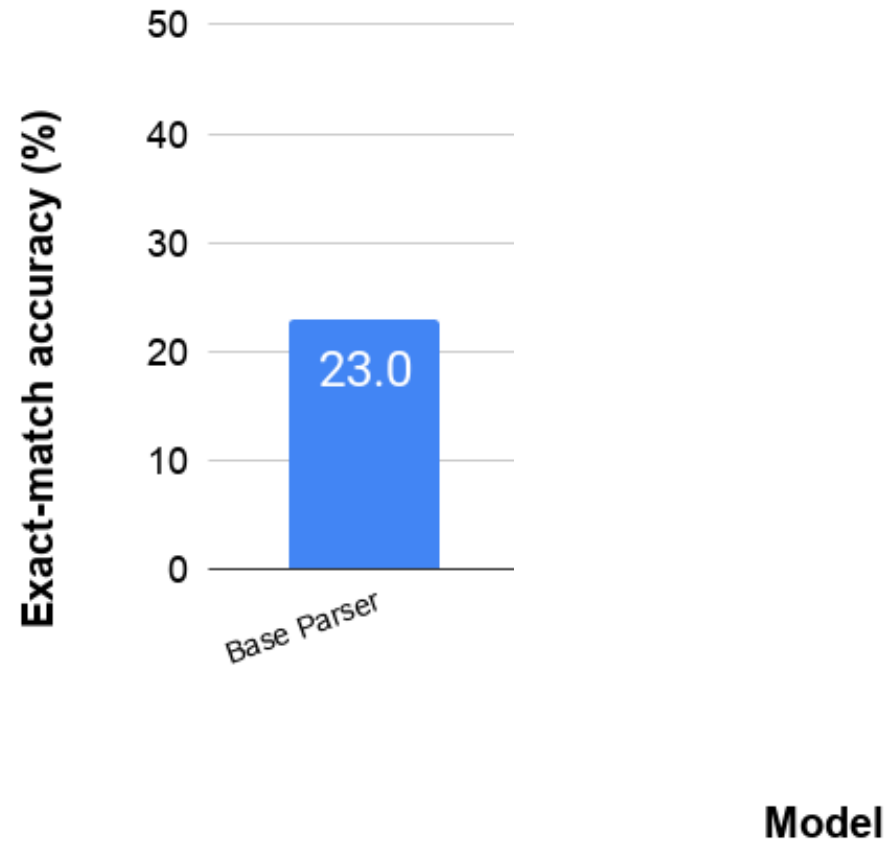


# Results on Spider



**Model**

# Results on Chinese Spider



# Analysis: In-Domain vs. Out-of-Domain

We create an **in-domain setting** from the Spider dataset.

- *Does the parser struggle out-of-domain?*

**YES**

In-domain vs. out-of-domain performance: **56.4% vs 78.2%**

- *Does DG-MAML hurt in-domain performance?*

**NO**

DG-MAML leads to a modest improvement **(+1.1%)**

# Key Takeaways

- **Meta-learning** can be useful beyond few-shot learning; we show it can also be used to promote **domain generalization** for semantic parsing.
- Without **changing model architectures**, **DG-MAML** can boost the performance of cross-domain parsers in mono- and cross-lingual settings.
- Code: <https://github.com/berlino/tensor2struct-public>
- Our recent work on extending **DG-MAML for compositional generalization** is accepted by **ACL2021**

# Basic Ideas

- We aim at *directly* optimizing for *domain generalization (DG)* via a meta-learning objective, dubbed **DG-MAML**.
- By constructing a set of *virtual cross-domain parsing* tasks, the objective encourage generalization to unseen domains in each task.

# DG-MAML Training Algorithm

---

**Algorithm 1** DG-MAML Training Algorithm

---

**Require:** Training databases  $\mathcal{D}$

**Require:** Learning rate  $\alpha$

1: **for** step  $\leftarrow 1$  **to**  $T$  **do**

8: **end for**

---

- Given a set of examples from databases  $\mathcal{D}$

# DG-MAML Training Algorithm

---

**Algorithm 1** DG-MAML Training Algorithm

---

**Require:** Training databases  $\mathcal{D}$

**Require:** Learning rate  $\alpha$

1: **for** step  $\leftarrow 1$  **to**  $T$  **do**

2:     Sample a task  $\tau$  of  $(\mathcal{D}_s^\tau, \mathcal{D}_t^\tau)$  from  $\mathcal{D}$

- Given a a set of examples from databases  $D$
- We first sample a virtual task from  $D$

8: **end for**

---

# DG-MAML Training Algorithm

---

**Algorithm 1** DG-MAML Training Algorithm

---

**Require:** Training databases  $\mathcal{D}$

**Require:** Learning rate  $\alpha$

- 1: **for** step  $\leftarrow 1$  **to**  $T$  **do**
- 2:     Sample a task  $\tau$  of  $(\mathcal{D}_s^\tau, \mathcal{D}_t^\tau)$  from  $\mathcal{D}$
- 3:     Sample mini-batch  $\mathcal{B}_s^\tau$  from  $\mathcal{D}_s^\tau$
- 4:     Sample mini-batch  $\mathcal{B}_t^\tau$  from  $\mathcal{D}_t^\tau$

8: **end for**

---

- Given a a set of examples from databases  $\mathcal{D}$
- We first sample a virtual task from  $\mathcal{D}$
- Sample examples from virtual source and target databases



# DG-MAML Training Algorithm

---

**Algorithm 1** DG-MAML Training Algorithm

---

**Require:** Training databases  $\mathcal{D}$

**Require:** Learning rate  $\alpha$

- 1: **for** step  $\leftarrow 1$  **to**  $T$  **do**
- 2:   Sample a task  $\tau$  of  $(\mathcal{D}_s^\tau, \mathcal{D}_t^\tau)$  from  $\mathcal{D}$
- 3:   Sample mini-batch  $\mathcal{B}_s^\tau$  from  $\mathcal{D}_s^\tau$
- 4:   Sample mini-batch  $\mathcal{B}_t^\tau$  from  $\mathcal{D}_t^\tau$
- 5:   Meta-train update:  
     $\theta' \leftarrow \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{B}_s^\tau}(\theta)$
- 6:   Compute meta-test objective:  
     $\mathcal{L}_\tau(\theta) = \mathcal{L}_{\mathcal{B}_s}(\theta) + \mathcal{L}_{\mathcal{B}_t}(\theta')$
- 7:   Final Update:  
     $\theta \leftarrow \text{Update}(\theta, \nabla_{\theta} \mathcal{L}_\tau(\theta))$
- 8: **end for**

- Given a a set of examples from databases  $\mathcal{D}$
- We first sample a virtual task from  $\mathcal{D}$
- Sample examples from virtual source and target databases
- Update parameters using a MAML objective

# MAML Objective

- Meta-Train: one step of SGD in the **virtual source domains**

$$\boldsymbol{\theta}' \leftarrow \boldsymbol{\theta} - \alpha \nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{B}_s}(\boldsymbol{\theta})$$

- Meta-Test: evaluate the parameters in the **virtual target domains**

$$\mathcal{L}_{\mathcal{B}_t}(\boldsymbol{\theta}')$$

- Final objective: **joint loss on both virtual source and target domains**

$$\mathcal{L}_{\tau}(\boldsymbol{\theta}) = \mathcal{L}_{\mathcal{B}_s}(\boldsymbol{\theta}) + \mathcal{L}_{\mathcal{B}_t}(\boldsymbol{\theta}')$$

# MAML Objective

$$\mathcal{L}_T(\theta) = \mathcal{L}_{\mathcal{B}_s}(\theta) + \mathcal{L}_{\mathcal{B}_t}(\theta')$$

## Intuition:

- optimize towards the better source *and* target domain performance **simultaneously**
- gradient step in the source domain should **be beneficial to the performance of the target domain** as well.

# MAML Objective vs. Supervised Learning

$$\mathcal{L}_{\mathcal{B}_s}(\boldsymbol{\theta}) + \mathcal{L}_{\mathcal{B}_t}(\boldsymbol{\theta}')$$

$$\mathcal{L}_{\mathcal{B}_s}(\boldsymbol{\theta}) + \mathcal{L}_{\mathcal{B}_t}(\boldsymbol{\theta})$$

## Comparison:

- Supervised learning objective (*right*) **does not pose any constraints** on the gradient update.
- MAML objective (*left*) can be viewed as a **regularization** of gradient updates.

# Analysis of DG-MAML

**First-order Taylor series expansion:**

$$\begin{aligned}\mathcal{L}_\tau(\boldsymbol{\theta}) &= \mathcal{L}_{\mathcal{B}_s}(\boldsymbol{\theta}) + \mathcal{L}_{\mathcal{B}_t}(\boldsymbol{\theta}') \\ &= \mathcal{L}_{\mathcal{B}_s}(\boldsymbol{\theta}) + \mathcal{L}_{\mathcal{B}_t}(\boldsymbol{\theta} - \alpha \nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{B}_s}(\boldsymbol{\theta})) \\ &\approx \mathcal{L}_{\mathcal{B}_s}(\boldsymbol{\theta}) + \mathcal{L}_{\mathcal{B}_t}(\boldsymbol{\theta}) - \alpha (\nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{B}_s}(\boldsymbol{\theta}) \cdot \nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{B}_s}(\boldsymbol{\theta}))\end{aligned}$$

DG-MAML further tries to maximize  $\nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{B}_s}(\boldsymbol{\theta}) \cdot \nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{B}_s}(\boldsymbol{\theta})$ , the **dot product between the gradients** of source and target domain. That is, it encourages gradients to generalize between source and target domain within each task  $\tau$ .

# First-Order Approximation: DG-FMAML

**The gradient of DG-MAML requires second derivatives:**

$$\begin{aligned}\nabla_{\theta} \mathcal{L}_{\tau}(\theta) &= \nabla_{\theta} \theta' \nabla_{\theta'} \mathcal{L}_{\mathcal{B}_t}(\theta') + \nabla_{\theta} \mathcal{L}_{\mathcal{B}_s}(\theta) \\ &= (\mathbf{I} - \alpha \nabla_{\theta}^2 \mathcal{L}_{\mathcal{B}_s}(\theta)) \nabla_{\theta'} \mathcal{L}_{\mathcal{B}_t}(\theta') + \nabla_{\theta} \mathcal{L}_{\mathcal{B}_s}(\theta)\end{aligned}$$

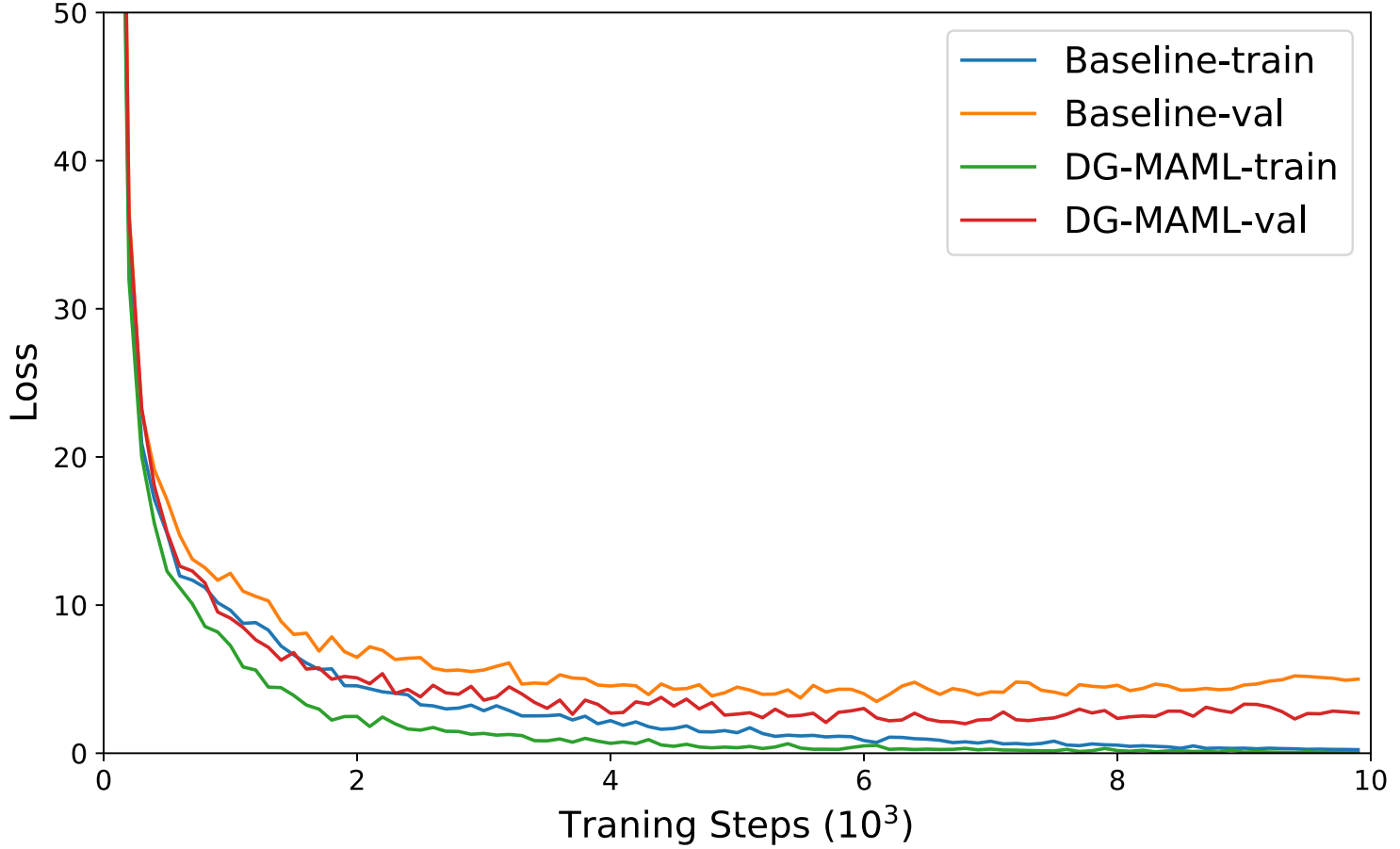
Inspired by Reptile, we consider the alternative of **ignoring this second-order term** and simply assume that  $\nabla_{\theta} \theta' = \mathbf{I}$ .

# First-Order Approximation: DG-FMAML

$$\mathcal{L}_{maml}(\theta) = \mathcal{L}_{\theta}(\text{🗄️🗄️}) + \mathcal{L}_{\theta'}(\text{🗄️})$$

$\theta'$  has no gradient wrt. to  $\theta$

# Loss Curve

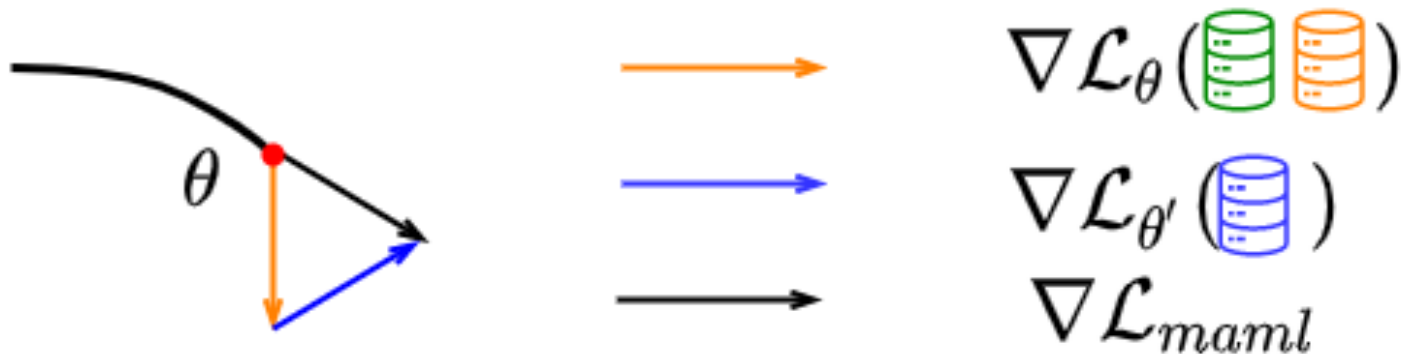




# First-Order Approximation: DG-FMAML

$$\mathcal{L}_{maml}(\theta) = \mathcal{L}_{\theta}(\text{green DB}, \text{orange DB}) + \mathcal{L}_{\theta'}(\text{blue DB})$$

$\theta'$  has no gradient wrt. to  $\theta$



# Construct Virtual Tasks for Meta-Learning

