# On the Stability and Generalization of Meta-Learning
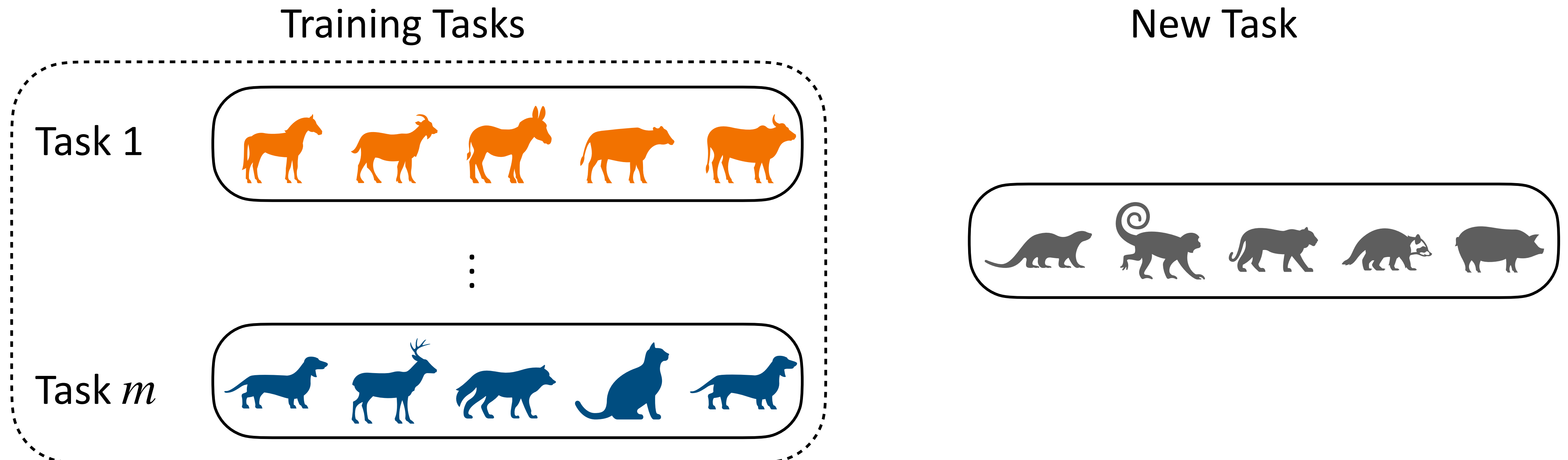
**Yunjuan Wang, Raman Arora**
**Johns Hopkins University**

# Motivation

- Large sets of training data from a single task are often lacking. Training data may stem from diverse tasks with shared similarities, while test data come from entirely new tasks.

- The challenge is to rapidly adapt to these unseen tasks without the need to train from scratch.
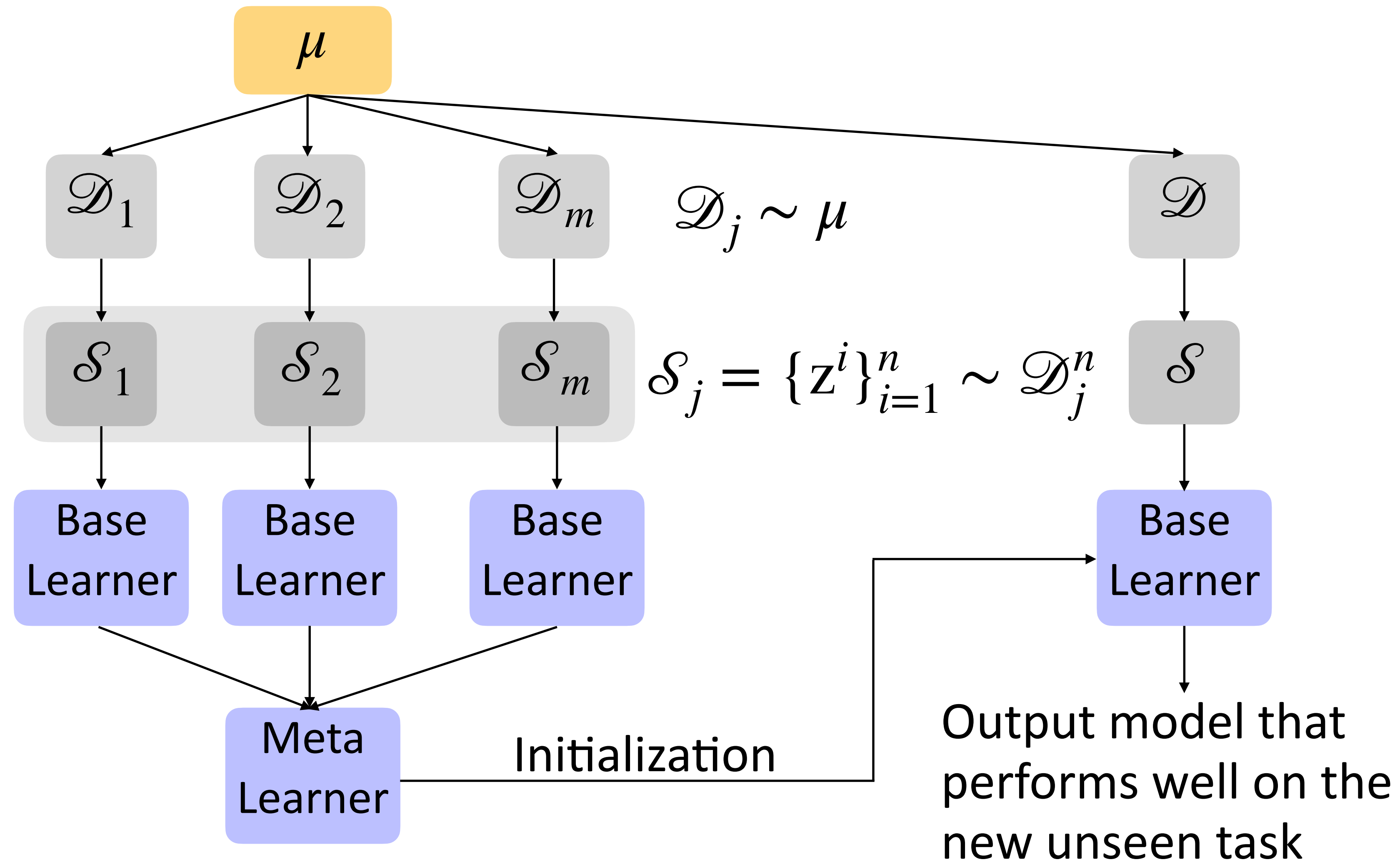
# Problem Setup

Task Distribution

Sample $m$ different tasks

$n$ training data for each task
Meta-sample $\mathbf{S} = \{\mathscr{S}_j\}_{j=1}^{m}$

Task level learning

Meta level learning

$\mu$

$\mathscr{D}_1$ $\mathscr{D}_2$ $\mathscr{D}_m$ $\mathscr{D}$

$\mathscr{D}_j \sim \mu$

$\mathscr{S}_1$ $\mathscr{S}_2$ $\mathscr{S}_m$ $\mathscr{S}$

$\mathscr{S}_j = \{z^i\}_{i=1}^{n} \sim \mathscr{D}_j^n$

Base Learner | Base Learner | Base Learner | Base Learner

Meta Learner

Initialization

Output model that performs well on the new unseen task

# Problem Setup

- Consider a supervised learning setting where each data point is denoted by $z = (x, y)$ drawn from some unknown distribution $\mathcal{D}$ over $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$, where input space $\mathcal{X}$, label space $\mathcal{Y}$.

- A meta-learning algorithm $\mathscr{A}$ takes the meta-sample $\mathbf{S}$ as input and outputs an algorithm $\mathscr{A}(\mathbf{S}) : (\mathcal{X} \times \mathcal{Y})^n \rightarrow \mathcal{H}$.

- **Goal:** learn a useful prior over tasks to help with rapid adaptation to new tasks.

- Transfer risk: $L(\mathscr{A}(\mathbf{S}), \mu) = \mathbb{E}_{\mathcal{D} \sim \mu} \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^n} L(\mathscr{A}(\mathbf{S})(\mathcal{S}), \mathcal{D}) = \mathbb{E}_{\mathcal{D} \sim \mu} \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^n} \mathbb{E}_{z \sim \mathcal{D}} \ell(\mathscr{A}(\mathbf{S})(\mathcal{S}), z).$

- Empirical multi-task risk: $L(\mathscr{A}(\mathbf{S}), \mathbf{S}) = \dfrac{1}{m} \sum\limits_{j=1}^{m} L(\mathscr{A}(\mathbf{S})(\mathcal{S}_j), \mathcal{S}_j) = \dfrac{1}{m} \sum\limits_{j=1}^{m} \dfrac{1}{n} \sum\limits_{i=1}^{n} \ell(\mathscr{A}(\mathbf{S})(\mathcal{S}_j), z_j^i).$

# Problem Setup

**Task level learning**: Given a meta-model (parameterized by meta-parameter $w$), the hope is that it can be adapted easily to a new task $\mathscr{D} \sim \mu$; in particular, a task-specific model $u$ can be quickly learned from a task-specific training set $\mathscr{S} \sim \mathscr{D}^n$ of size $n$ using the following proximal update:

$$u = \arg\min_{u \in \mathscr{W}} L(u, \mathscr{S}) + \frac{\lambda}{2}\|u - w\|^2$$

**Meta level learning**: $w$ itself is learned on the given meta-sample $\mathbf{S} = \{\mathscr{S}_j\}_{j=1}^{m}$ by minimizing a regularized empirical loss averaged over tasks:

$$\hat{w} = \arg\min_{w \in \mathscr{W}} \frac{1}{m} \sum_{j=1}^{m} \min_{u \in \mathscr{W}} \left[ L(u, \mathscr{S}_j) + \frac{\lambda}{2}\|u - w\|^2 \right]$$

# Proximal Meta Learning Algorithm

**Algorithm 1** Prox Meta-Learning Algorithm $\mathcal{A}$

**Input:** Meta-sample $\mathbf{S} = \{\mathcal{S}_j\}_{j=1}^{m}$, epochs $T$, $K$, step sizes $\gamma$, $\eta$, regularization parameter $\lambda$

1: $\mathbf{w}_1 = 0$.
2: **for** $t = 1, 2, \ldots, T$ **do**
3:     **for** $j = 1, \ldots, m$ **do**
4:         $\mathbf{u}(\mathbf{w}_t, \mathcal{S}_j) = \mathcal{A}_{\text{task}}(\mathbf{w}_t, \mathcal{S}_j, K, \eta, \lambda)$
        % Using Algorithm 2
5:     **end for**
6:     Calculate the gradient, $\forall j \in [m]$,
    $\nabla F_{\mathcal{S}_j}(\mathbf{u}(\mathbf{w}_t, \mathcal{S}_j), \mathbf{w}_t) = -\lambda(\mathbf{u}(\mathbf{w}_t, \mathcal{S}_j) - \mathbf{w}_t)$.
7:     Update $\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\gamma}{m} \sum_{j=1}^{m} \nabla F_{\mathcal{S}_j}(\mathbf{u}(\mathbf{w}_t, \mathcal{S}_j), \mathbf{w}_t)$
8:     $\mathbf{w}_{t+1} = \Pi_{\mathcal{W}}(\mathbf{w}_{t+1})$
9: **end for**
10: **return** $@\mathcal{A}_{\text{task}}(\mathbf{w}_{T+1}, \cdot, K, \eta, \lambda)$

---

**Algorithm 2** Task-specific Algorithm $\mathcal{A}_{\text{task}}$

**Input:** Pretrained model w, training data $\mathcal{S}$, #epochs $K$, step size $\eta$, reg. parameter $\lambda$

1: Option 1 (RERM):
2: $\mathbf{u}(\mathbf{w}, \mathcal{S}) = \text{argmin}_{\mathbf{u} \in \mathcal{W}} L(\mathbf{u}, \mathcal{S}) + \frac{\lambda}{2} \|\mathbf{u} - \mathbf{w}\|^2$.

3: Option 2 (GD): Set $\mathbf{u}^{(1)}(\mathbf{w}, \mathcal{S}) = \mathbf{w}$
4: **for** $t = 1, 2, \ldots, K-1$ **do**
5:     $\mathbf{u}^{(k+1)}(\mathbf{w}, \mathcal{S}) = \mathbf{u}^{(k)}(\mathbf{w}, \mathcal{S})$
        $-\eta(\nabla L(\mathbf{u}^{(k)}(\mathbf{w}, \mathcal{S}), \mathcal{S})$
        $+\lambda(\mathbf{u}^{(k)}(\mathbf{w}, \mathcal{S}) - \mathbf{w}))$
6:     $\mathbf{u}^{(k+1)}(\mathbf{w}, \mathcal{S}) = \Pi_{\mathcal{W}}(\mathbf{u}^{(k+1)}(\mathbf{w}, \mathcal{S}))$
7: **end for**
8: **return** Option 1 (RERM): $\mathbf{u}(\mathbf{w}, \mathcal{S})$
    Option 2 (GD): $\frac{1}{K} \sum_{k=1}^{K} \mathbf{u}^{(k)}(\mathbf{w}, \mathcal{S})$
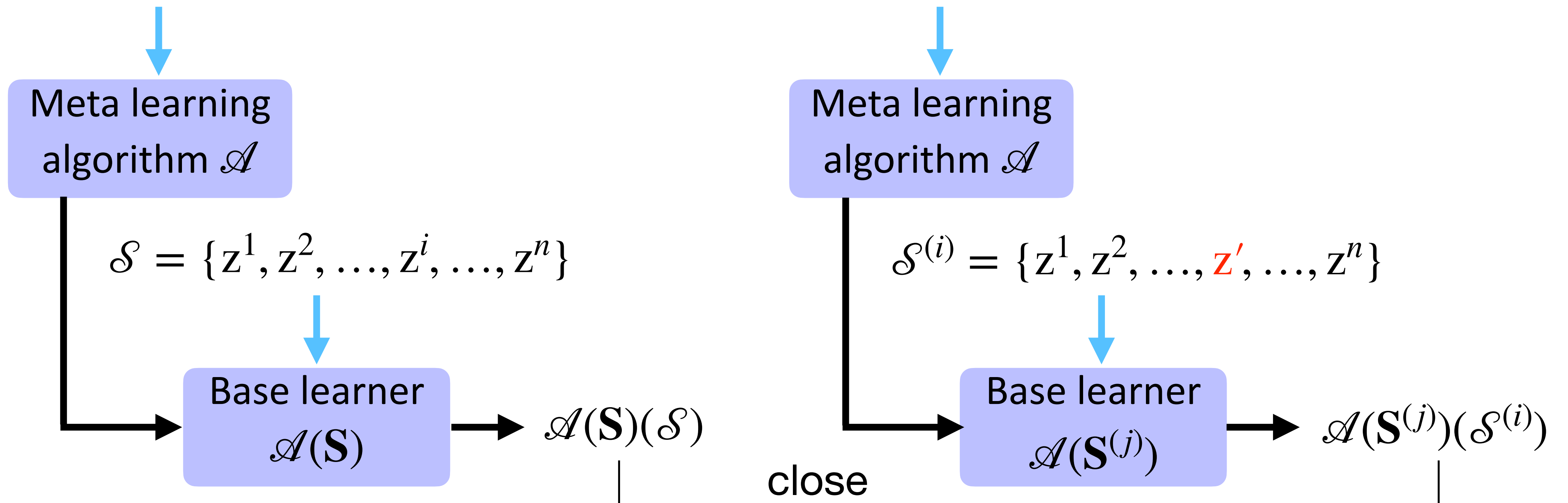
# Uniform meta-stability

*Definition:* A meta-learning algorithm $\mathscr{A}$ is $\bar{\beta}$-uniformly meta-stable if for any neighboring meta-samples $\mathbf{S}, \mathbf{S}^{(j)}$, and neighboring samples $\mathscr{S}, \mathscr{S}^{(i)}$, for any task $\mathscr{D} \sim \mu$ and any $z \sim \mathscr{D}$, we have

$$|\ell(\mathscr{A}(\mathbf{S})(\mathscr{S}), z) - \ell(\mathscr{A}(\mathbf{S}^{(j)})(\mathscr{S}^{(i)}), z)| \leq \bar{\beta}$$

$$\mathbf{S} = \{\mathscr{S}_1, \mathscr{S}_2, \ldots, \mathscr{S}_j, \ldots, \mathscr{S}_m\} \qquad \mathbf{S}^{(j)} = \{\mathscr{S}_1, \mathscr{S}_2, \ldots, \mathscr{S}_j', \ldots, \mathscr{S}_m\}$$

# Uniform meta-stability

**Definition:** A meta-learning algorithm $\mathscr{A}$ is $\bar{\beta}$-uniformly meta-stable if for any neighboring meta-samples $\mathbf{S}, \mathbf{S}^{(j)}$, and neighboring samples $\mathscr{S}, \mathscr{S}^{(i)}$, for any task $\mathscr{D} \sim \mu$ and any $\mathbf{z} \sim \mathscr{D}$, we have

$$|\ell(\mathscr{A}(\mathbf{S})(\mathscr{S}), \mathbf{z}) - \ell(\mathscr{A}(\mathbf{S}^{(j)})(\mathscr{S}^{(i)}), \mathbf{z})| \leq \bar{\beta}$$

**Theorem:** Consider a meta-learning problem for some $M$-bounded loss function $\ell$ and task distribution $\mu$. For any $\bar{\beta}$-uniformly meta-stable learning algorithm $\mathscr{A}$, we have that with probability at least $1 - \delta$,

$$L(\mathscr{A}(\mathbf{S}), \mu) \lesssim L(\mathscr{A}(\mathbf{S}), \mathbf{S}) + \bar{\beta} \log(mn)\log(1/\delta) + M\sqrt{\log(1/\delta)/(mn)}$$

# Bound Transfer Generalization Gap

| Algorithm | Loss | Conditions | Uniform meta-stability $\bar{\beta}$ |
|---|---|---|---|
| Algo. 1 with RERM | convex, $G$-Lipschitz | $\gamma \leq \frac{1}{\lambda}$ | $\frac{G^2}{\lambda m} + \frac{G^2}{\lambda n}$ |
| Algo. 1 with RERM | convex, $H$-smooth, $M$-bounded | $\gamma \leq \frac{1}{\lambda}, \lambda \geq H$ | $\frac{HM}{\lambda(2n-1)} + \frac{HM}{\lambda(m+1)}$ |
| Algo. 1 with GD | convex, $G$-Lipschitz, $H$-smooth | $\eta \leq \frac{2}{H+2\lambda}, \gamma \leq \frac{1}{\lambda T}$ | $\frac{G^2}{\lambda m} + \frac{G^2}{\lambda n}$ |
| Algo. 1 with GD | $\rho$-weakly convex, $G$-Lipschitz | $\eta \leq \frac{1}{\lambda}, \gamma \leq \frac{1}{\lambda T}, \lambda \geq 2\rho$ | $G^2\sqrt{\frac{\eta}{\lambda}} + \frac{G^2}{\lambda m} + \frac{G^2}{\lambda n}$ |
| Algo. 3 with GD | $\rho$-weakly convex, $G$-Lipschitz | $\eta \leq \frac{1}{\lambda}, \gamma \leq \frac{1}{\lambda T}, \lambda \geq 2\rho$ | $G^2\sqrt{\frac{\eta}{\lambda}} + \frac{G^2}{\lambda m} + \frac{G^2}{\lambda n}$, w.h.p. |

Table 1: Bounds on uniform meta-stability $\bar{\beta}$ for different families of learning problems. Here, $\eta$ is the step-size for GD for task-specific learning, $\gamma$ is the step-size for GD for meta-parameter learning, $m$ is the number of tasks during training, $n$ is the number of training data for the task at test time.

**Extension**   Proximal Meta-Learning with Stochastic Optimization

Robust Adversarial Proximal Meta-Learning

# Excess Transfer Risk

$$\underbrace{L(\mathscr{A}(\mathbf{S})(\mathscr{S}), \mathscr{D}) - L(\mathrm{u}_*, \mathscr{D})}_{\text{Excess Transfer Risk } \mathscr{E}_{\text{risk}}(\mathscr{A})} = \underbrace{L(\mathscr{A}(\mathbf{S})(\mathscr{S}), \mathscr{D}) - \frac{1}{m} \sum_{j=1}^{m} L(\mathscr{A}(\mathbf{S})(\mathscr{S}_j), \mathscr{S}_j)}_{\text{Generalization Gap } \mathscr{E}_{\text{gen}}(\mathscr{A})}$$

$\mathrm{u}_* = \arg\min\limits_{\mathrm{u} \in \mathscr{W}} L(\mathrm{u}, \mathscr{D})$ optimal task-specific hypothesis for the <span style="color:red">unseen</span> task;

$\mathrm{u}_j^* = \arg\min\limits_{\mathrm{u} \in \mathscr{W}} L(\mathrm{u}, \mathscr{S}_j)$ optimal task-specific hypothesis for the <span style="color:red">given</span> training tasks.

$$+ \quad \underbrace{\frac{1}{m} \sum_{j=1}^{m} L(\mathscr{A}(\mathbf{S})(\mathscr{S}_j), \mathscr{S}_j) - L(\mathrm{u}_j^*, \mathscr{S}_j)}_{\text{Optimization and Approximation Error } \mathscr{E}_{\text{opt+app}}(\mathscr{A})}$$

$$+ \underbrace{L(\mathrm{u}_j^*, \mathscr{S}_j) - L(\mathrm{u}_*, \mathscr{S}_j)}_{\leq 0} + \underbrace{L(\mathrm{u}_*, \mathscr{S}_j) - L(\mathrm{u}_*, \mathscr{D})}_{\mathbb{E}_{\forall j \in [m], \mathscr{S}_j \sim \mathscr{D}_j^n, \mathscr{D}_j \sim \mu, \mathscr{D} \sim \mu} = 0}$$

# Excess Transfer Risk

Convex and smooth loss:

Setting $\eta = \mathscr{O}\left(\dfrac{1}{\lambda\sqrt{K}}\right)$ gives us $\mathbb{E}\left[\mathscr{E}_{\text{risk}}(\mathscr{A})\right] \leq \mathscr{O}\left(\dfrac{1}{\lambda\sqrt{K}} + \dfrac{1}{\lambda m} + \dfrac{1}{\lambda n} + \dfrac{\lambda}{T} + \lambda\sigma^2\right)$, where

$\sigma^2 = \dfrac{1}{m}\displaystyle\sum_{j=1}^{m}\|\hat{\mathrm{w}} - \mathrm{u}_j^*\|^2$ is the approximation error.

Convex and non-smooth loss:

Setting $\eta = \mathscr{O}\left(\dfrac{1}{\lambda K^{2/3}}\right)$ gives us $\mathbb{E}\left[\mathscr{E}_{\text{risk}}(\mathscr{A})\right] \leq \mathscr{O}\left(\dfrac{1}{\lambda K^{1/3}} + \dfrac{1}{\lambda m} + \dfrac{1}{\lambda n} + \dfrac{\lambda}{T} + \lambda\sigma^2\right)$

# Thank you!