# Is Bayesian Model-Agnostic Meta Learning Better than Model-Agnostic Meta Learning, Provably?

Lisha Chen, Tianyi Chen

Rensselaer Polytechnic Institute

03/14/2022

# Definition of meta learning

## Meta-learning (learning to learn):

To learn a model that can well adapt or generalize to new tasks and new environments that have never been encountered during training time.

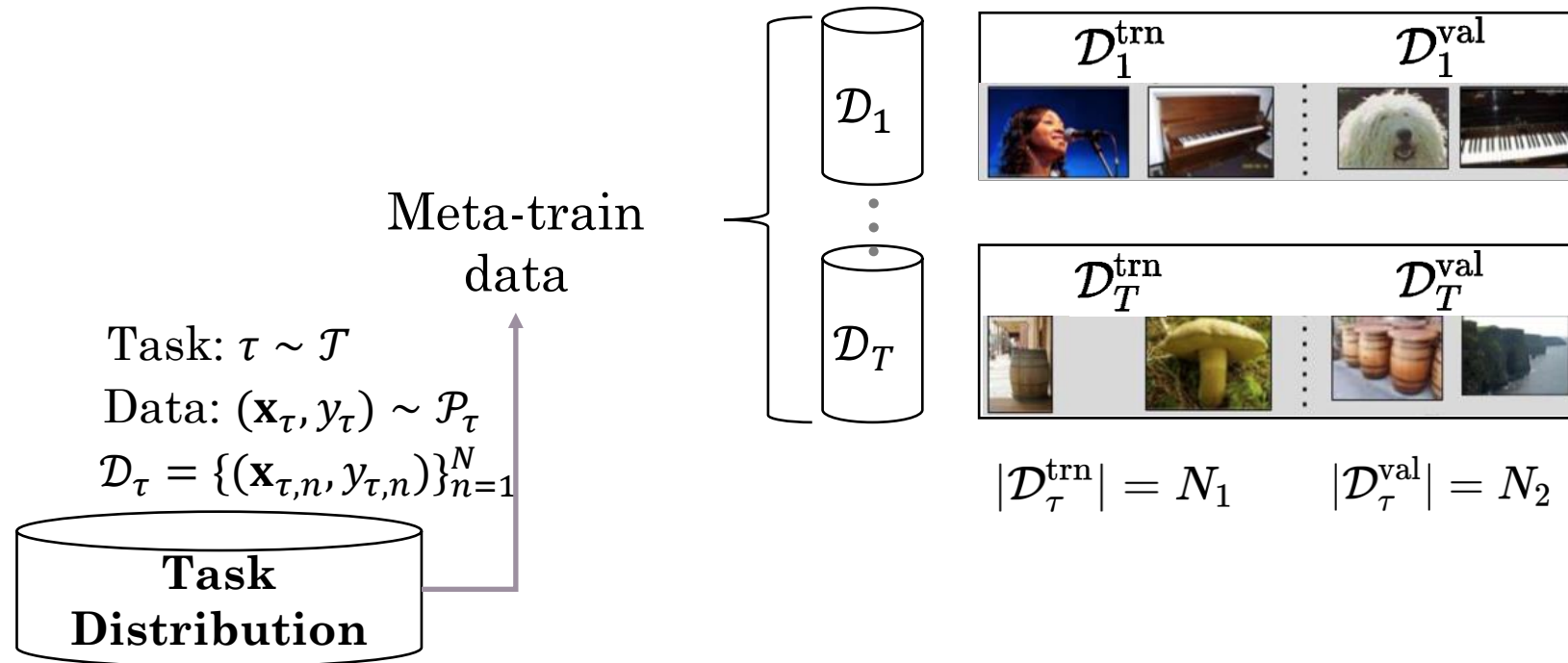## Prior work (incomplete)

[Bengio et al '90]

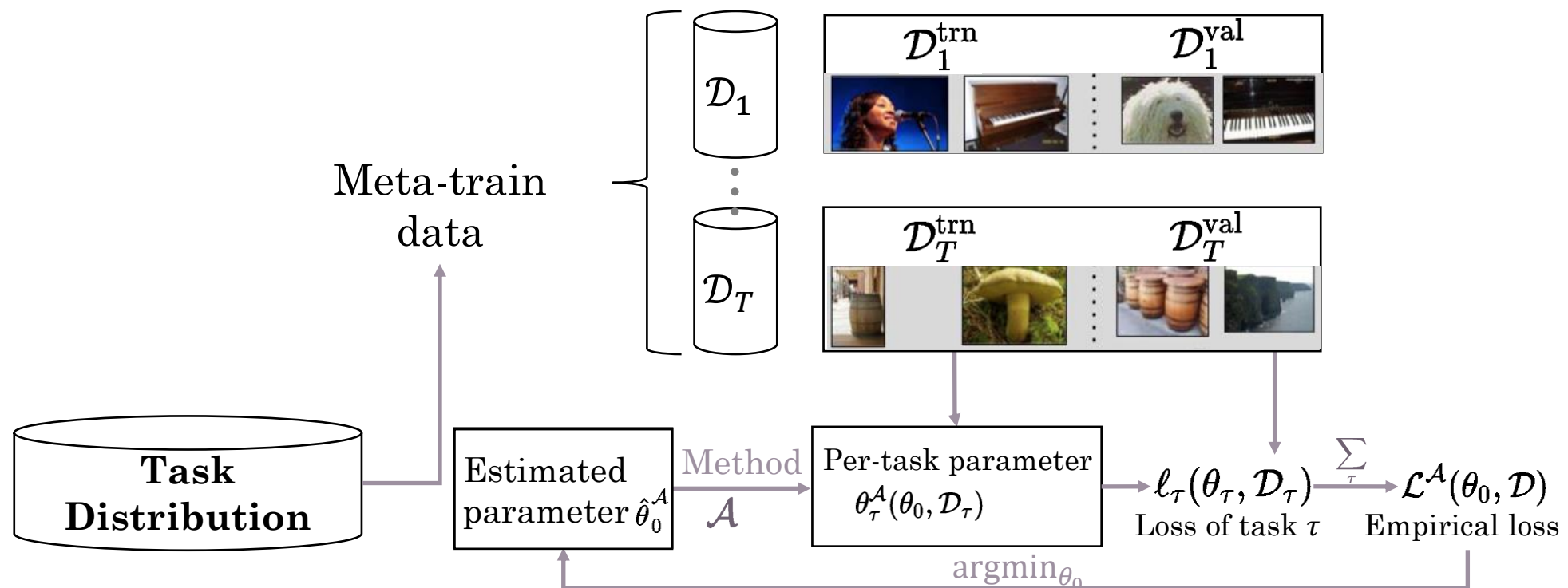[Maclaurin et al '15]

[Vinyals et al '16]

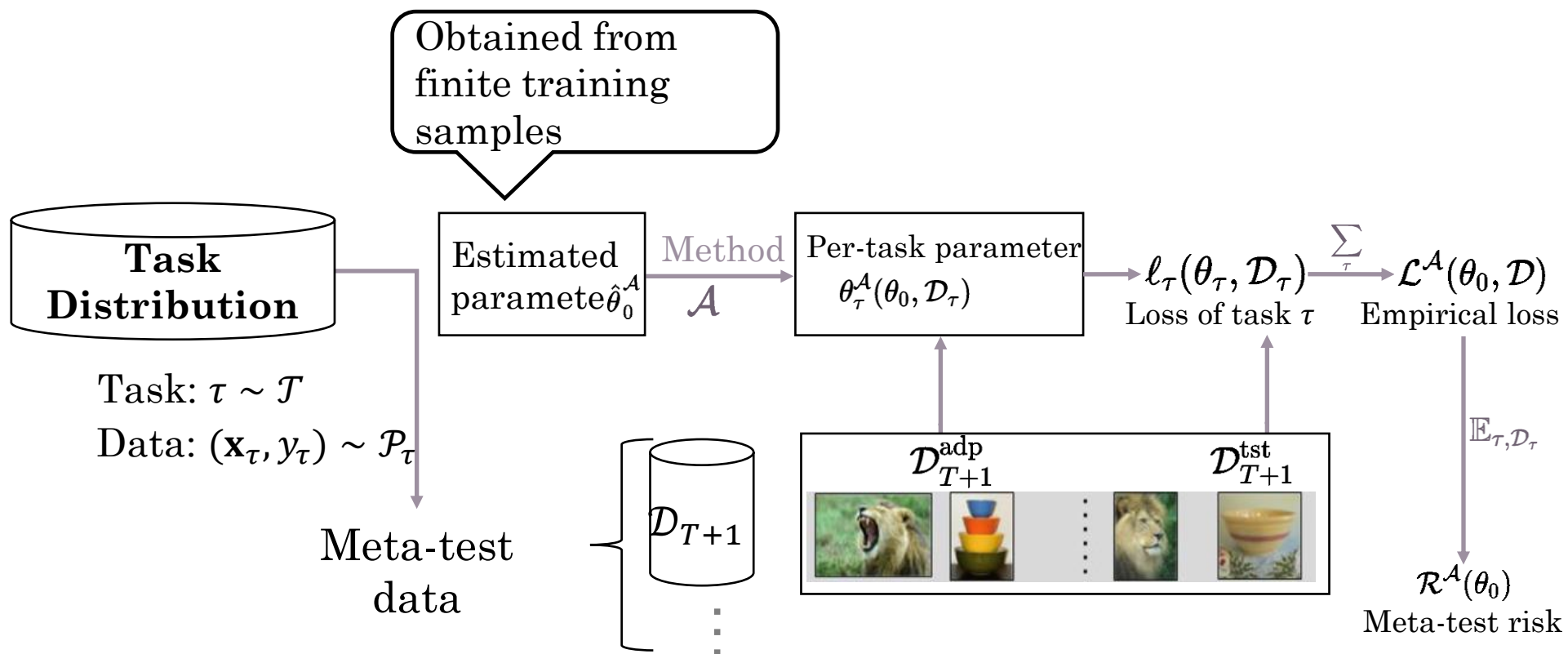[Santoro et al '16]

[Wichrowska et al '17]

[Finn et al '17]

# Meta learning setup



Meta-train data

Task: $\tau \sim \mathcal{T}$

Data: $(\mathbf{x}_\tau, y_\tau) \sim \mathcal{P}_\tau$

$\mathcal{D}_\tau = \{(\mathbf{x}_{\tau,n}, y_{\tau,n})\}_{n=1}^N$

**Task Distribution**

$\mathcal{D}_1$

$\mathcal{D}_T$

$\mathcal{D}_1^{\text{trn}}$     $\mathcal{D}_1^{\text{val}}$

$\mathcal{D}_T^{\text{trn}}$     $\mathcal{D}_T^{\text{val}}$

$|\mathcal{D}_\tau^{\text{trn}}| = N_1$      $|\mathcal{D}_\tau^{\text{val}}| = N_2$

# Meta learning setup

Obtained from finite training samples

**Task Distribution**

Estimated parameter $\hat{\theta}_0^{\mathcal{A}}$

Method

$\mathcal{A}$

Per-task parameter $\theta_\tau^{\mathcal{A}}(\theta_0, \mathcal{D}_\tau)$

$\ell_\tau(\theta_\tau, \mathcal{D}_\tau) \xrightarrow{\sum_\tau} \mathcal{L}^{\mathcal{A}}(\theta_0, \mathcal{D})$

Loss of task $\tau$        Empirical loss

Task: $\tau \sim \mathcal{T}$

Data: $(\mathbf{x}_\tau, y_\tau) \sim \mathcal{P}_\tau$

Meta-test data

$\mathcal{D}_{T+1}$

$\mathcal{D}_{T+1}^{\text{adp}}$        $\mathcal{D}_{T+1}^{\text{tst}}$

$\mathbb{E}_{\tau, \mathcal{D}_\tau}$

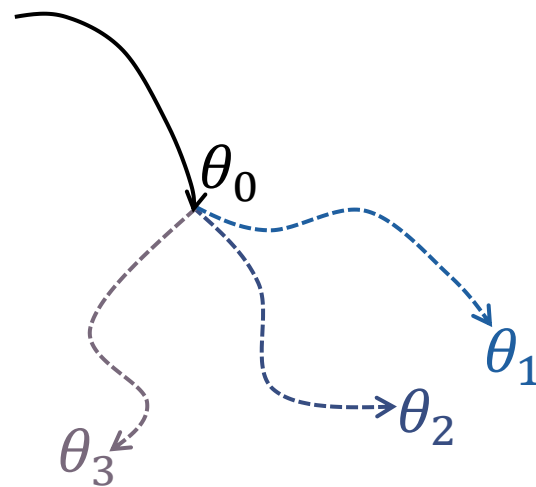$\mathcal{R}^{\mathcal{A}}(\theta_0)$

Meta-test risk

# Basis of comparison

$$\hat{\theta}_\tau^{\mathrm{er}}\left(\theta_0, \mathcal{D}_\tau^{\mathrm{trn}}\right) = \theta_0$$

$$\hat{\theta}_\tau^{\mathcal{A}}\left(\theta_0, \mathcal{D}_\tau^{\mathrm{trn}}\right) \neq \theta_0$$

$\theta_0(\theta_1 \quad \theta_2 \quad \theta_3)$

$\theta_0$

$\theta_1$

$\theta_2$

$\theta_3$

$$\hat{\theta}_\tau^{\mathcal{A}}\left(\theta_0, \mathcal{D}_\tau^{\mathrm{trn}}\right) = ?$$

# Baseline methods – ERM

**General formulations (empirical loss)**

**ERM**

$$\mathcal{L}^{\text{er}}(\theta_0, \mathcal{D}) = \frac{1}{T} \sum_{\tau=1}^{T} \ell_\tau(\theta_0, \mathcal{D}_{\tau,N})$$
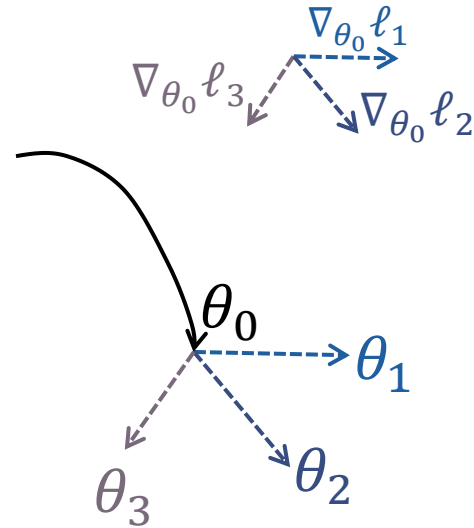
$$\theta_0(\theta_1 \quad \theta_2 \quad \theta_3)$$

**General formulations (empirical loss)**

**MAML [Finn et al '17]**

$$\mathcal{L}^{\mathrm{ma}}(\theta_0, \mathcal{D}) = \frac{1}{T} \sum_{\tau=1}^{T} \ell_\tau \left( \hat{\theta}_\tau^{\mathrm{ma}}(\theta_0, \mathcal{D}_{\tau,N_1}^{\mathrm{trn}}), \mathcal{D}_{\tau,N_2}^{\mathrm{val}} \right)$$

$$\text{s.t. } \hat{\theta}_\tau^{\mathrm{ma}}(\theta_0, \mathcal{D}_{\tau,N_1}^{\mathrm{trn}}) = \theta_0 - \alpha \nabla_{\theta_0} \ell_\tau (\theta_0, \mathcal{D}_{\tau,N_1}^{\mathrm{trn}})$$
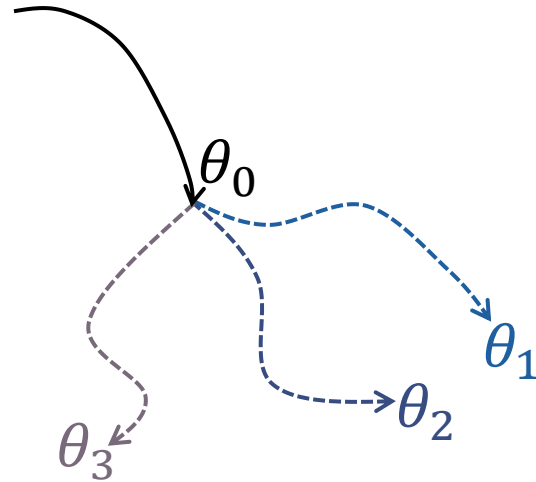
## General formulations (empirical loss)

### iMAML [Rajeswaran et al '19]

$$\mathcal{L}^{\text{im}}(\theta_0, \mathcal{D}) = \frac{1}{T} \sum_{\tau=1}^{T} \ell_\tau \left( \hat{\theta}_\tau^{\text{im}} \left( \theta_0, \mathcal{D}_{\tau,N_1}^{\text{trn}} \right), \mathcal{D}_{\tau,N_2}^{\text{val}} \right)$$

$$\text{s.t. } \hat{\theta}_\tau^{\text{im}} \left( \theta_0, \mathcal{D}_{\tau,N_1}^{\text{trn}} \right) = \arg\min_{\theta_\tau} - \log p \left( \theta_\tau \mid \mathcal{D}_{\tau,N_1}^{\text{trn}}, \theta_0 \right)$$
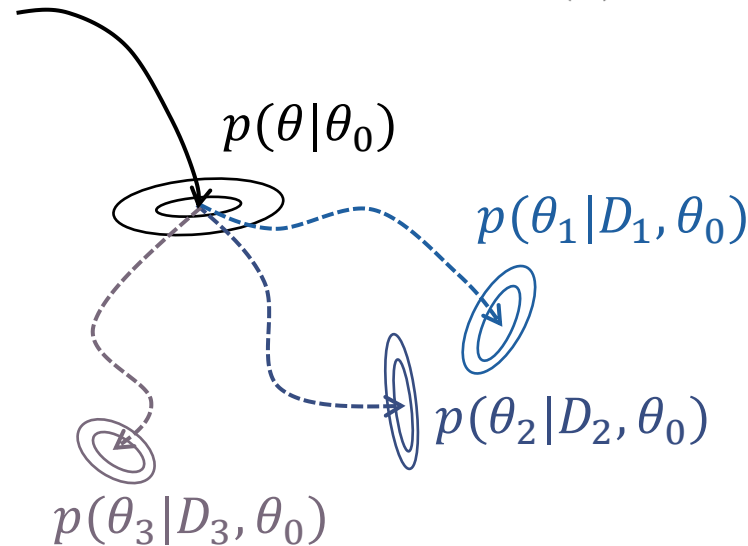
# Baseline methods – BaMAML
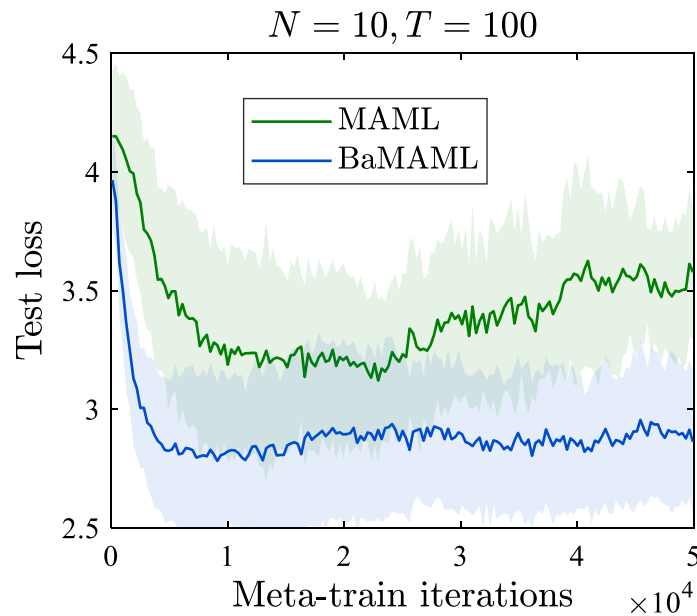
## General formulations (empirical loss)

### BaMAML

$$\mathcal{L}^{\text{ba}}(\theta_0, \mathcal{D}) = \frac{1}{T} \sum_{\tau=1}^{T} \ell_\tau \left( \hat{p} \left( \theta_\tau \mid \mathcal{D}^{\text{trn}}_{\tau, N_1}, \theta_0 \right), \mathcal{D}^{\text{val}}_{\tau, N_2} \right)$$

$$\text{s.t. } \hat{p} \left( \theta_\tau \mid \mathcal{D}^{\text{trn}}_{\tau, N_1}, \theta_0 \right) = \arg \min_{q(\theta_\tau) \in \mathcal{Q}} D_{\text{KL}} \left( q(\theta_\tau) \| p \left( \theta_\tau \mid \mathcal{D}^{\text{trn}}_{\tau, N_1}, \theta_0 \right) \right)$$

$p(\theta|\theta_0)$

$p(\theta_1|D_1, \theta_0)$

$p(\theta_2|D_2, \theta_0)$

$p(\theta_3|D_3, \theta_0)$

LLAMA [Grant et al '18]
PLATIPUS [Finn et al '18]
BMAML [Yoon et al '18]
VAMPIRE [Nguyen et al '20]

# Compare the baseline methods



$N = 10, T = 100$

Sinusoidal regression

MiniImageNet classification accuracy

| Method | 1-shot 5-way |
|--------|--------------|
| MAML   | $48.70 \pm 1.84$ |
| iMAML  | $49.30 \pm 1.88$ |
| BaMAML | $51.54 \pm 0.74$ |

Empirically, BaMAML has better accuracy but is more challenging to solve than MAML.

# Goal of this work

**Questions:**

❑ If and when is BaMAML better than MAML, provably?

❑ What are the decomposable factors that make BaMAML better?

**Contributions:**

**First theoretical understanding for above questions.**

# A unified view

$$\ell_\tau(\theta_0, \mathcal{D}_\tau) = -\log p\big(\mathbf{y}_\tau^{\mathrm{val}} \mid \mathbf{X}_\tau^{\mathrm{val}}, \theta_0, \mathcal{D}_\tau^{\mathrm{trn}}\big)$$

$$= -\log \int \underbrace{p(\mathbf{y}_\tau^{\mathrm{val}} \mid \mathbf{X}_\tau^{\mathrm{val}}, \theta_\tau)}_{\text{Likelihood}} \underbrace{p(\theta_\tau \mid \theta_0, \mathcal{D}_\tau^{\mathrm{trn}})}_{\text{Posterior}} d\theta_\tau$$

Bayes rule

$$p(\theta_\tau \mid \theta_0, \mathcal{D}_\tau^{\mathrm{trn}}) = \frac{p(\mathcal{D}_\tau^{\mathrm{trn}} \mid \theta_\tau)p(\theta_\tau \mid \theta_0)}{p(\mathcal{D}_\tau^{\mathrm{trn}} \mid \theta_0)}$$

# A unified view

$$\ell_\tau(\theta_0, \mathcal{D}_\tau) = -\log \int p(\mathbf{y}_\tau^{\mathrm{val}} \mid \mathbf{X}_\tau^{\mathrm{val}}, \theta_\tau) p(\theta_\tau \mid \theta_0, \mathcal{D}_\tau^{\mathrm{trn}}) d\theta_\tau$$
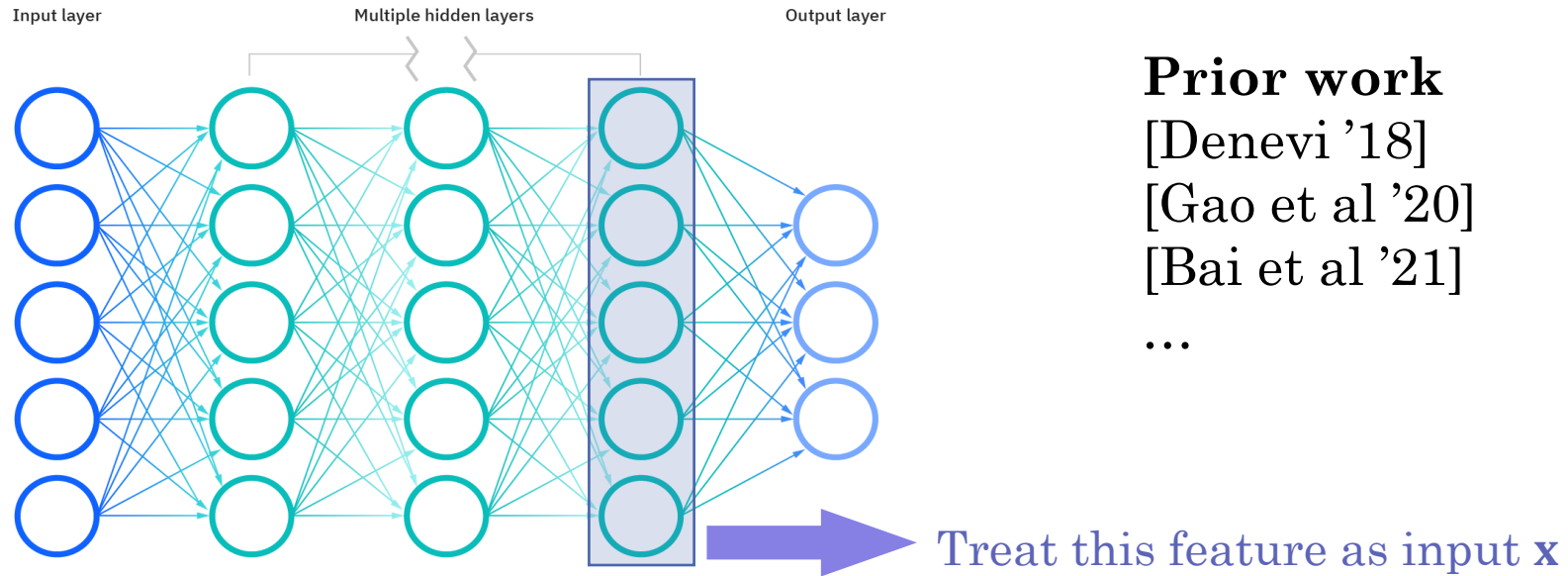
Point estimate | e.g. MAML,
iMAML

$$\delta\left(\theta_\tau - \hat{\theta}_\tau^{\mathcal{A}}\right)$$

$$\ell_\tau(\theta_0, \mathcal{D}_\tau) = -\log p\left(\mathbf{y}_\tau^{\mathrm{val}} \mid \mathbf{X}_\tau^{\mathrm{val}}, \hat{\theta}_\tau^{\mathcal{A}}(\theta_0, \mathcal{D}_\tau^{\mathrm{trn}})\right)$$

Now we are ready to compare these methods in the same framework.

# Meta linear regression – data model

Input layer      Multiple hidden layers      Output layer

**Prior work**
[Denevi '18]
[Gao et al '20]
[Bai et al '21]

...

Treat this feature as input $\mathbf{x}$

Data model

$$y_\tau = \theta_\tau^{\mathrm{gt}\top} \mathbf{x}_\tau + \epsilon_\tau, \ \ \text{with} \ \ \ \epsilon_\tau \overset{\text{i.i.d.}}{\sim} \mathcal{N}\left(0, \sigma_\tau^2\right) \ \ \ \mathbf{Q}_\tau = \mathbb{E}\left[\mathbf{x}_\tau \mathbf{x}_\tau^\top \mid \tau\right].$$

# Loss function under meta linear regression

$$y_\tau = \theta_\tau^{\text{gt}\top} \mathbf{x}_\tau + \epsilon_\tau, \text{ with } \epsilon_\tau \overset{\text{i.i.d.}}{\sim} \mathcal{N}\left(0, \sigma_\tau^2\right) \quad \mathbf{Q}_\tau = \mathbb{E}\left[\mathbf{x}_\tau \mathbf{x}_\tau^\top \mid \tau\right].$$

Recall $\ell_\tau(\theta_0, \mathcal{D}_\tau) = -\log \int p(\mathbf{y}_\tau^{\text{val}} \mid \mathbf{X}_\tau^{\text{val}}, \theta_\tau) p(\theta_\tau \mid \theta_0, \mathcal{D}_\tau^{\text{trn}}) d\theta_\tau$

Bayes rule

Assumption: Given $\hat{\theta}_\tau^{\mathcal{A}}$, $p(y_\tau \mid \mathbf{x}_\tau, \hat{\theta}_\tau^{\mathcal{A}}) = \mathcal{N}(\hat{\theta}_\tau^{\mathcal{A}\top} \mathbf{x}_\tau, \sigma_\tau^2)$

$$p(\theta_\tau \mid \mathcal{D}_\tau^{\text{trn}}, \theta_0) \propto p(\mathcal{D}_\tau^{\text{trn}} \mid \theta_\tau) p(\theta_\tau \mid \theta_0),$$

Prior is Gaussian $p(\theta_\tau \mid \theta_0) \propto \exp\left\{-\gamma\|\theta_\tau - \theta_0\|_2^2\right\}$

$\gamma$ is the weight of the prior

# Basic assumptions

1. **(Bounded eigenvalues)** For any $\tau, 0 < \underline{\lambda} \leq \lambda(\mathbf{Q}_\tau) \leq \bar{\lambda}$

2. **(Ground truth task parameter distribution)**

   1) $\theta_\tau^{\mathrm{gt}}$ is independent of $\mathbf{X}_\tau$.

   2) the individual entries $\left\{\theta_{\tau,i}^{\mathrm{gt}}\right\}_{i\in[d],\tau\in[T]}$ are independent and $\mathcal{O}\!\left(R/\sqrt{d}\right)$-sub-Gaussian, where $R$ is a constant.

   3) $\left\|\mathbb{E}\!\left[\theta_\tau^{\mathrm{gt}}\right]\right\| \leq M$ .

# Meta linear regression – Risk decomposition

**Meta test risk decomposition**

$$\mathcal{R}^{\mathcal{A}}(\hat{\boldsymbol{\theta}}_0^{\mathcal{A}}) = \underbrace{\mathcal{R}^{\mathcal{A}}(\boldsymbol{\theta}_0^{\mathcal{A}})}_{\text{optimal popultation risk}} + \underbrace{\left\|\hat{\boldsymbol{\theta}}_0^{\mathcal{A}} - \boldsymbol{\theta}_0^{\mathcal{A}}\right\|_{\mathbb{E}_\tau[\mathbf{W}_\tau^{\mathcal{A}}]}^2}_{\text{statistical error } \mathcal{E}_{\mathcal{A}}^2(\hat{\boldsymbol{\theta}}_0^{\mathcal{A}})}$$

$$\theta_0^{\mathcal{A}} := \arg\min_{\theta_0} \mathcal{R}^{\mathcal{A}}(\theta_0)$$

Analyze optimal population risk and statistical error separately.

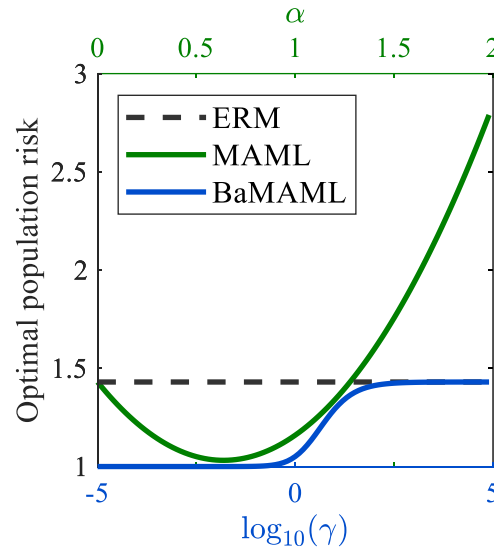# Comparison of optimal population risk

## Theorem 1 (informal)

Under Assumptions 1-2,

**ERM vs MAML**
Can find $\alpha$ in a range that
$\mathcal{R}^{\mathrm{ma}}(\theta_0^{\mathrm{ma}}) < \mathcal{R}^{\mathrm{er}}(\theta_0^{\mathrm{er}})$.

**MAML vs BaMAML**
Can find $\gamma$ in a range that
$\mathcal{R}^{\mathrm{ba}}(\theta_0^{\mathrm{ba}}) < \mathcal{R}^{\mathrm{ma}}(\theta_0^{\mathrm{ma}})$.



❑ Essentially,

$\mathcal{R}^{\mathrm{er}}(\theta_0^{\mathrm{er}}) > \inf_{\alpha} \mathcal{R}^{\mathrm{ma}}(\theta_0^{\mathrm{ma}}; \alpha) > \inf_{\gamma} \mathcal{R}^{\mathrm{ba}}(\theta_0^{\mathrm{ba}}; \gamma)$

❑ If $\alpha$ not properly chosen, MAML can be worse than ERM, but not for iMAML.

❑ Choice of $\gamma$ reflects trade-off between adaptation speed and adaptation performance.

# Precise characterization of statistical error

**Assumption 3 (Linear centroid model).**

1) $\mathbf{x}_\tau \overset{\text{iid}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I_d})$

2) $\mathrm{Cov}_{\theta_\tau^{\text{gt}}}\left[\theta_\tau^{\text{gt}}\right] = \dfrac{R^2}{d}\mathbf{I_d}.$

**Implications**

Assumption 3 (1) assumes $\mathbf{Q}_\tau = \mathbf{I}_d$ , therefore $\mathbf{W}_\tau^{\mathcal{A}} = w_\mathcal{A}\mathbf{I}_d$.
This implies for different methods $\mathcal{A}, \theta_0^{\mathcal{A}} = \mathbb{E}_\tau\left[\theta_\tau^{\text{gt}}\right].$
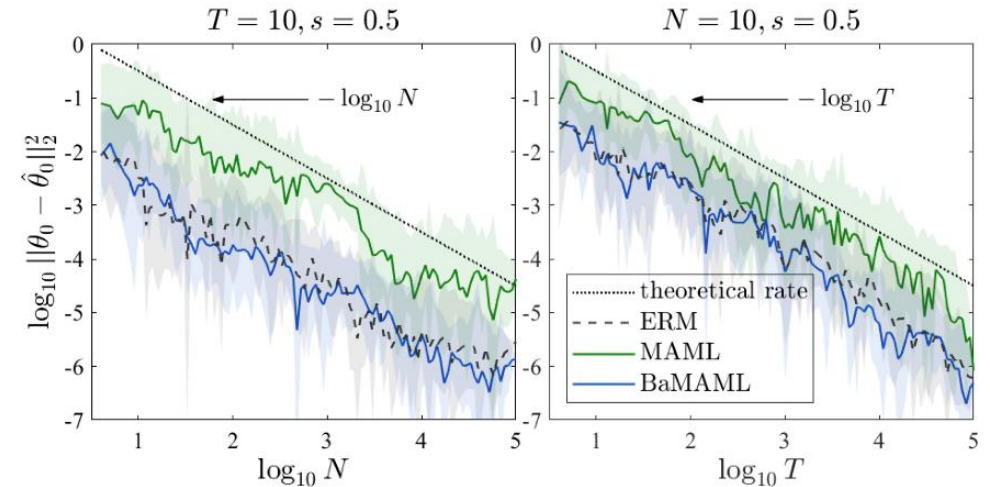
# Precise characterization of statistical error

**Theorem 2 (informal)**

Define $C^{\mathcal{A}} := \dfrac{1}{d} \left\langle \mathbb{E}^{-2}\left[\hat{\mathbf{W}}_{\tau,N}^{\mathcal{A}}\right], \mathbb{E}\left[\left(\hat{\mathbf{W}}_{\tau,N}^{\mathcal{A}}\right)^2\right] \right\rangle$

$\varrho$ as higher order term.

Under Assumptions 1-3, the following hold with high probability

$$w_{\mathcal{A}}\left\|\theta_0 - \hat{\theta}_0\right\|_2^2 = \frac{R^2}{T}\left(w_{\mathcal{A}}C^{\mathcal{A}} + \tilde{\mathcal{O}}\left(\sqrt{\frac{d}{T}}\right) + \tilde{\mathcal{O}}\left(\frac{1}{\sqrt{d}}\right)\right) + \varrho$$



Now we are ready to quantify the dominating constant exactly.

# Sharp comparison of statistical error

❑ Limits of dominating constants

$$\inf_{\substack{\alpha > 0 \\ s \in (0,1)}} \lim_{\substack{d, N \to \infty \\ d/N \to \eta}} C^{\mathrm{ma}} = \inf_{\substack{\gamma > 0 \\ s \in (0,1)}} \lim_{\substack{d, N \to \infty \\ d/N \to \eta}} C^{\mathrm{im}} = 1 + \eta.$$

$$\inf_{\substack{\gamma > 0 \\ s \in (0,1)}} \lim_{\substack{d, N \to \infty \\ d/N \to \eta}} C^{\mathrm{ba}} \begin{cases} = 1, & \eta \leq 1 \\ \leq \eta, & \eta > 1 \end{cases}$$

❑ Under linear centroid model, the dominating constant in the statistical error with optimally tuned hyperparameters satisfies,

BaMAML < MAML = iMAML

# Concluding remarks

❑ BaMAML (and iMAML) has better adaptation flexibility than one-step MAML, leading to smaller optimal population risk.

❑ For statistical error, all methods have the same dependence on $N, T, d$, and their difference lies in the constant. BaMAML has better constant than point estimate (iMAML & MAML) due to model averaging.

❑ Under linear centroid model, the dominating constant in the statistical error with optimally tuned hyperparameters satisfies,

$$\text{BaMAML} \leq \text{MAML} = \text{iMAML}$$

Justify the theoretical benefits of
BaMAML over MAML.

# References

Lisha Chen and Tianyi Chen. ``Is Bayesian Model-Agnostic Meta Learning Better than Model-Agnostic Meta Learning, Provably?," ***Proc. of AISTATS***, 2022.

Yu Bai et al. ``How Important is the Train-Validation Split in Meta-Learning?," ***Proc. of ICML***, 2021.

Katelyn Gao and Ozan Sener. ``Modeling and Optimization Trade-off in Meta-learning," ***Proc. of NeurIPS***, 2020.

Erin Grant et al. ``Recasting Gradient-Based Meta-Learning as Hierarchical Bayes," ***Proc. of ICLR***, 2018.

Taesup Kim et al. ``Bayesian Model-Agnostic Meta-Learning," ***Proc. of NIPS***, 2018.

# Thank you!

Contact: chenl21@rpi.edu

Code: https://github.com/lisha-chen/Bayesian-MAML-vs-MAML