

Probabilistic Model-Agnostic Meta Learning

Chelsea Finn*, Kelvin Xu*, Sergey Levine

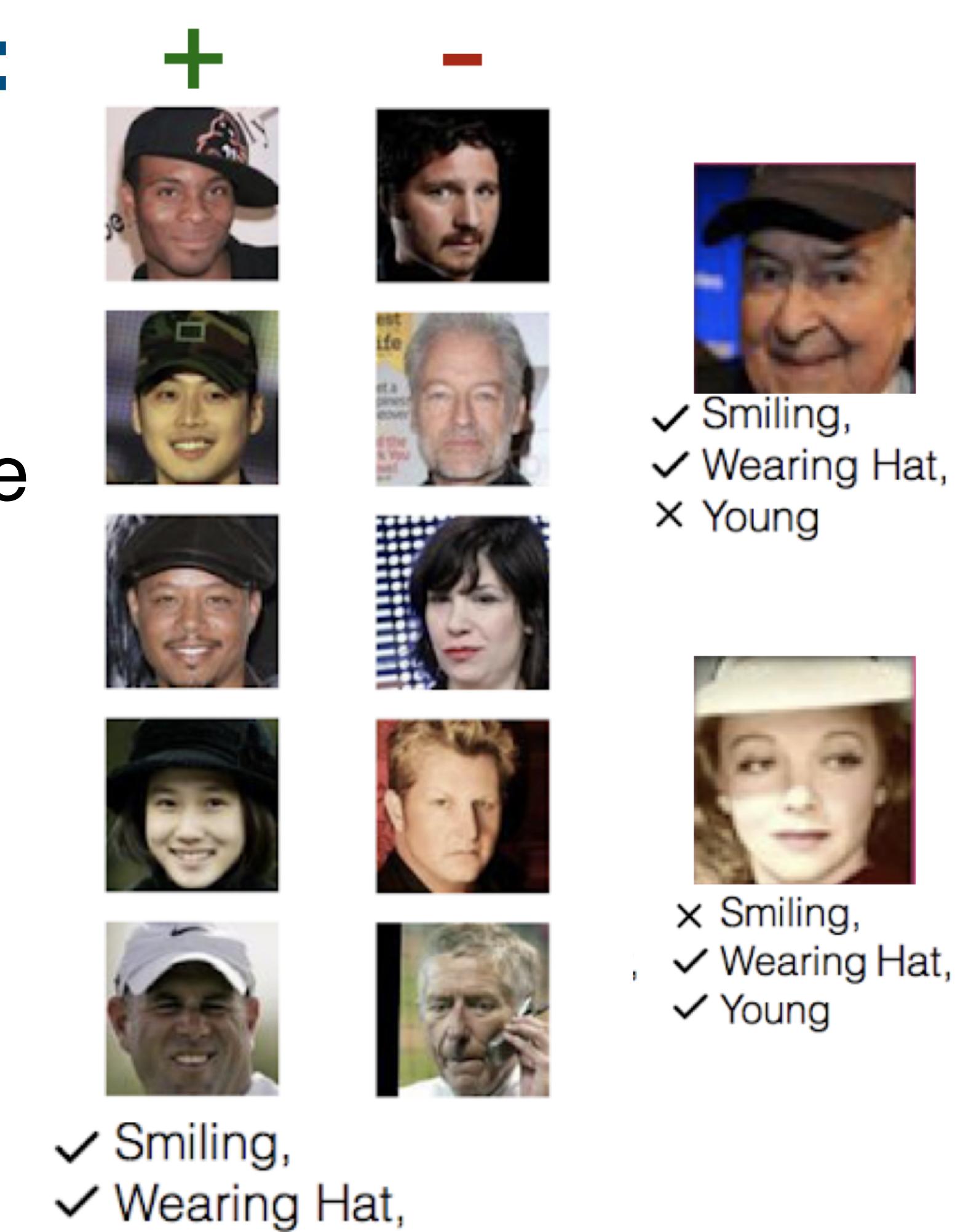


Motivation

1. Learning a prior is one avenue to enable **few-shot learning**, which is a key aspect of human intelligence.
2. For some tasks however, **task ambiguity** in the few shot regime is unavoidable. For active learning, reinforcement learning or safety critical tasks (e.g., in medical domains), **modeling ambiguity** is key.

Key takeaways of this work:

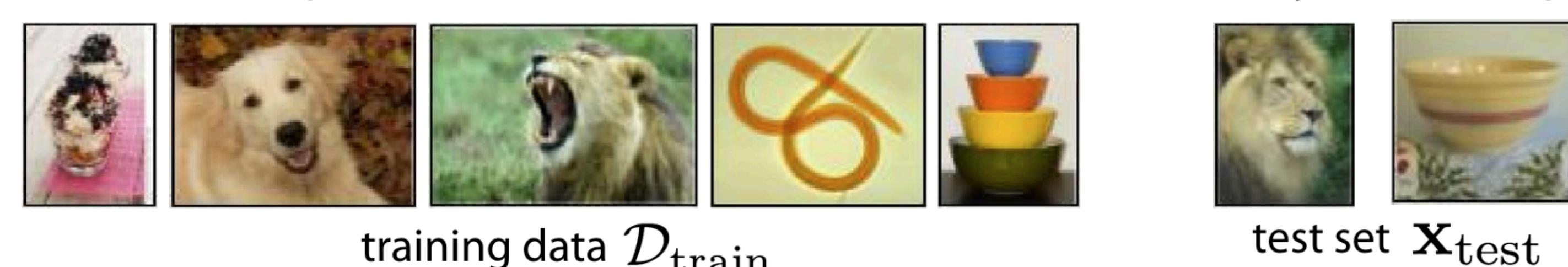
- 1) Extends model-agnostic meta learning (MAML) [1] by modeling uncertainty in a principled manner using amortized variational inference in a graphical model formulation of meta-learning (also see [2]).
- 2) Proposes a scalable probabilistic meta-learning algorithm that can be used to improve performance in *regression, active learning, and ambiguous classification settings*.



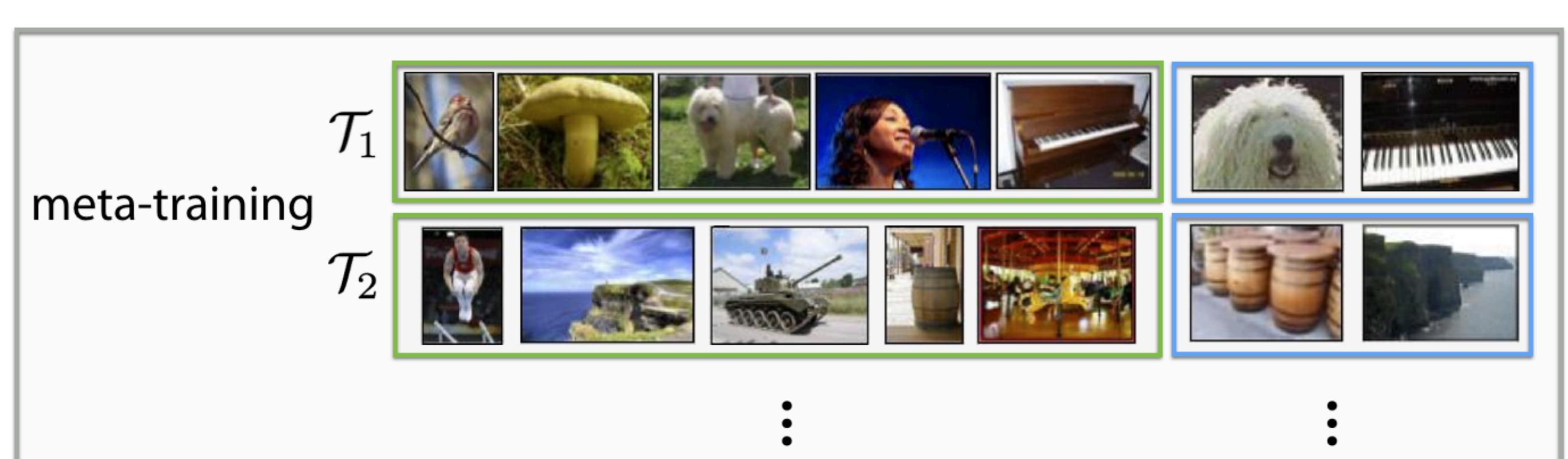
Background: the Meta-Learning Problem

- Assume access to a set of tasks drawn from the task distribution $\mathcal{T}_i \sim p(\mathcal{T})$ which can be split into a training set $\mathcal{D}_{\mathcal{T}_i}^{\text{tr}} := \{(x_{i,j}^{\text{tr}}, y_{i,j}^{\text{tr}}) \mid \forall j\}$ and test set $\mathcal{D}_{\mathcal{T}_i}^{\text{test}} := \{(x_{i,j}^{\text{tr}}, y_{i,j}^{\text{tr}}) \mid \forall j\}$

Given 1 example of 5 classes:



Classify new examples



training classes

Model-Agnostic Meta-Learning (MAML) [1]

$$\min_{\theta} \sum_{\mathcal{T} \sim p(\mathcal{T})} \mathcal{L}(\theta - \alpha \nabla_{\theta} \mathcal{L}(\theta, \mathcal{D}_{\mathcal{T}}^{\text{tr}}), \mathcal{D}_{\mathcal{T}}^{\text{test}}) = \min_{\theta} \sum_{\mathcal{T} \sim p(\mathcal{T})} \mathcal{L}(\phi_i, \mathcal{D}_{\mathcal{T}}^{\text{test}})$$

Future Directions

- Studying how ambiguity and uncertainty can guide data-acquisition
- Data-dependent posterior variances for tasks with differing levels of uncertainty

Probabilistic Gradient Based Meta-Learning

- We start with a probabilistic deep network and the conditional likelihoods of both the meta-train and meta-test set

$$\theta \sim p(\theta) = \mathcal{N}(\mu_{\theta}, \Sigma_{\theta}) \log p(y_i^{\text{tr}} \mid x_i^{\text{tr}}, \phi_i)$$

$$\phi_i \sim p(\phi_i \mid \theta) \log p(y_i^{\text{test}} \mid x_i^{\text{test}}, \phi_i)$$

- Goal: sample $\phi_i \sim p(\phi_i \mid x_i^{\text{tr}}, y_i^{\text{tr}}, x_i^{\text{test}})$

$$p(\phi_i \mid x_i^{\text{tr}}, y_i^{\text{tr}}) \propto \int p(\theta) p(\phi_i \mid \theta) p(y_i^{\text{tr}} \mid x_i^{\text{tr}}, \phi_i) d\theta$$

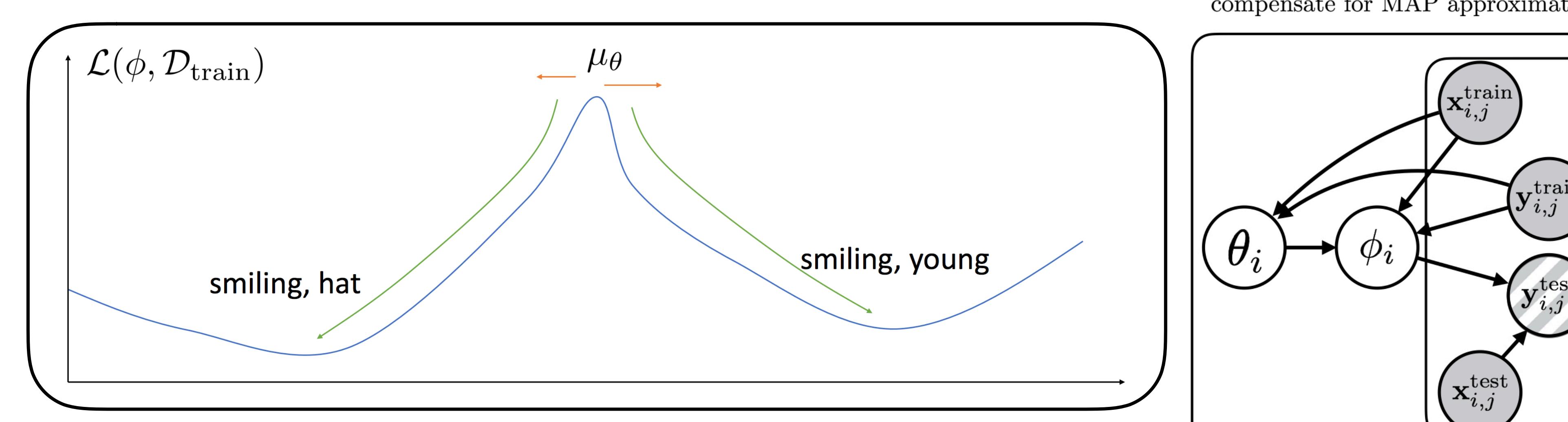
Completely intractable!

- If we knew $p(\phi_i \mid x_i^{\text{tr}}, y_i^{\text{tr}})$, we could use ancestral sampling

Key idea: use the MAP approximation which is crude but extremely convenient [2]

$$p(\phi_i \mid \theta, x_i^{\text{tr}}, y_i^{\text{tr}}) \approx \delta(\hat{\phi}_i)$$

$$\hat{\phi}_i \approx \theta + \alpha \nabla_{\theta} \log p(y_i^{\text{tr}} \mid x_i^{\text{tr}}, \theta)$$



Learning with amortized inference

- To train our model parameters θ , we can derive a variational bound

$$E_{\theta \sim q} [\log p(y_i^{\text{test}} \mid x_i^{\text{test}}, \hat{\phi}_i(\theta))] - D_{\text{KL}}(q(\theta \mid x_i^{\text{test}}, y_i^{\text{test}}) \parallel p(\theta))$$

- Training a network to output parameters not scalable, instead parameterize the q distribution by embedding gradient descent:

$$q(\theta \mid x_i^{\text{test}}, y_i^{\text{test}}) = \mathcal{N}(\mu_{\theta} + \alpha \nabla_{\mu_{\theta}} \log p(y_i^{\text{test}} \mid x_i^{\text{test}}, \mu_{\theta}), \Sigma_q)$$

Intuition: cheat at meta-training time using the test set, but minimize a divergence with a prior which is used for sampling at meta-test time

Additional dependencies

- In the 2nd graphical model, $x_i^{\text{tr}}, y_i^{\text{tr}}, \theta$ are conditionally independent
- Since we used a crude MAP approximation, these independences may not hold. Allow the model to compensate by learning a task specific prior $p(\theta \mid x_i^{\text{tr}}, y_i^{\text{tr}}) = \mathcal{N}(\mu_{\theta} + \alpha \nabla_{\mu_{\theta}} \log p(y_i^{\text{tr}} \mid x_i^{\text{tr}}, \mu_{\theta}), \Sigma_p)$

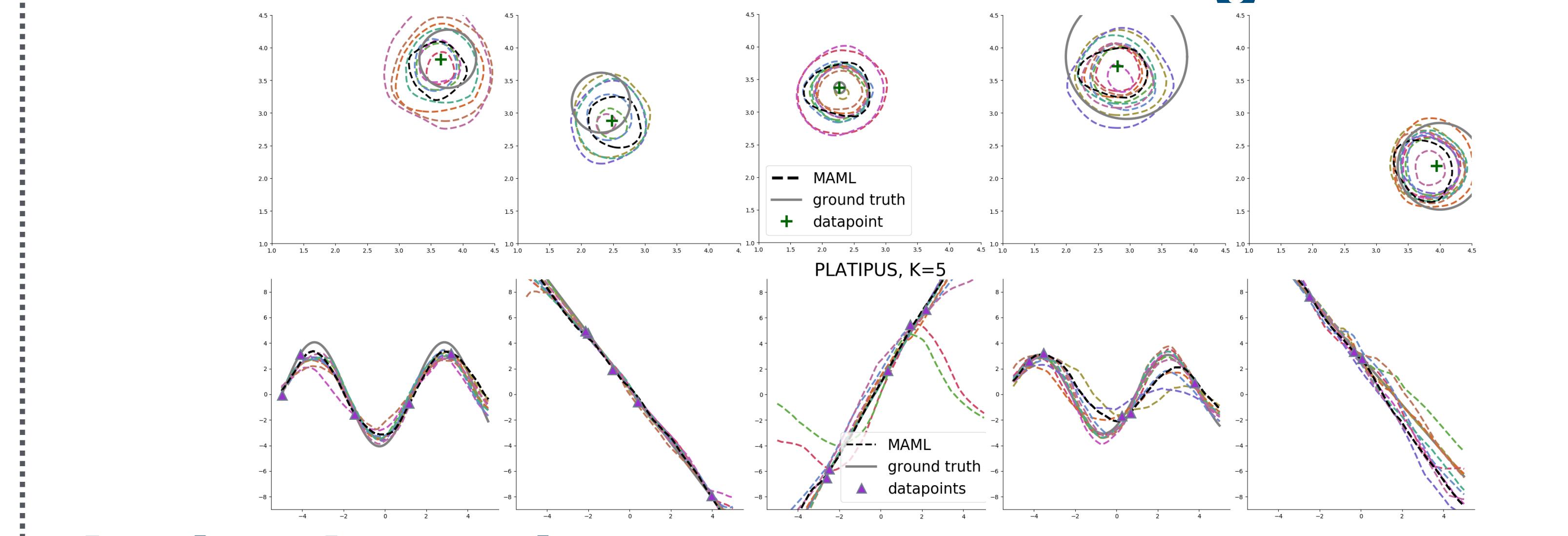
Meta-training

- (1) $\theta \sim q(\theta) = \mathcal{N}(\mu_{\theta} + \alpha \nabla_{\mu_{\theta}} \log p(y_i^{\text{test}} \mid x_i^{\text{test}}, \mu_{\theta}), \Sigma_q)$
- (2) $\phi_i \sim p(\phi_i \mid x_i^{\text{tr}}, y_i^{\text{tr}}) \approx \hat{\phi}_i = \theta + \alpha \nabla_{\theta} \log p(y_i^{\text{tr}} \mid x_i^{\text{tr}}, \theta)$
- (3) take gradient step with respect to variational bound

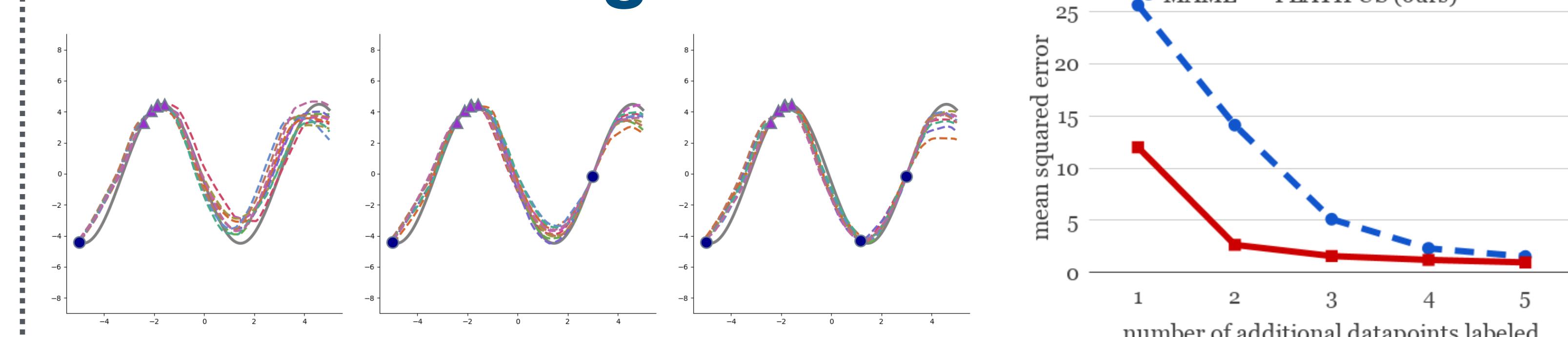
Meta-testing

- (1) $\theta \sim p(\theta) = \mathcal{N}(\mu_{\theta} + \alpha \nabla_{\mu_{\theta}} \log p(y_i^{\text{tr}} \mid x_i^{\text{tr}}, \mu_{\theta}), \Sigma_p)$
- (2) $\phi_i \sim p(\phi_i \mid x_i^{\text{tr}}, y_i^{\text{tr}}) \approx \hat{\phi}_i = \theta + \alpha \nabla_{\theta} \log p(y_i^{\text{tr}} \mid x_i^{\text{tr}}, \theta)$

1-Shot Classification/5-shot Regression



Active Learning



Ambiguous Attributes Classification



Mini-Imagenet

	5-way, 1-shot Accuracy
MinilImagenet	48.70 ± 1.84%
MAML [8]	49.40 ± 1.83%
LLAMA [14]	49.97 ± 0.32%
Reptile [28]	50.13 ± 1.86%
PLATIPUS (ours)	50.71 ± 1.87%
Meta-SGD [24]	50.71 ± 1.87%

- Slight boost over MAML, while comparable to other approaches, (*similar architecture)

[1] Finn, Abbeel, Levine. *Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks*. ICML '17

[2] Grant, Finn, Levine, Darrell, Griffiths. *Recasting gradient-based meta-learning as hierarchical Bayes*. ICLR '18