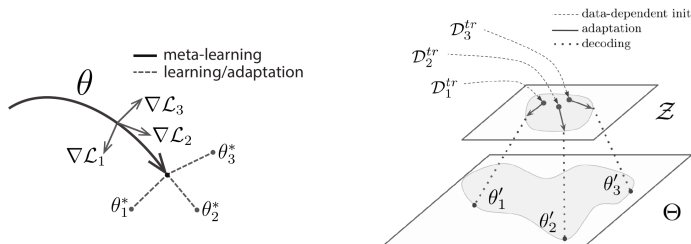# Meta-Learning with Latent Embedding Optimization
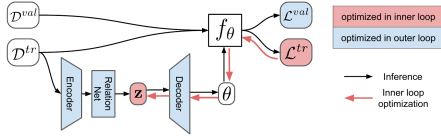
Rusu et al., ICLR 2019

# Outline

- This paper model Latent Embedding Optimization (LEO), an extension over the MAML model
  - Learn a low-dimensional latent embedding of model parameters and performs optimization-based meta learning in this space
  - Provide 2 advantages over the MAML model
    - initial parameters for new tasks are conditioned on the training data, which enables a task-specific starting point for adaptation
    - optimizing in low-dimensional latent space is more efficient

# Comparision between LEO vs MAML



- The MAML model aims to find an set of parameters that can use for many different tasks
- The LEO model initializes task-dependent set of parameters, update them in a general framework, it is more desirable

# LEO Algorithm and Arch

---

**Algorithm 2** MAML for Few-Shot Supervised Learning

---

**Require:** $p(\mathcal{T})$: distribution over tasks
**Require:** $\alpha, \beta$: step size hyperparameters
1: randomly initialize $\theta$
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**

---

**Algorithm 1** Latent Embedding Optimization

---

**Require:** Training meta-set $\mathcal{S}^{tr} \in \mathcal{T}$
**Require:** Learning rates $\alpha, \eta$
1: Randomly initialize $\phi_e, \phi_r, \phi_d$
2: Let $\phi = \{\phi_e, \phi_r, \phi_d, \alpha\}$
3: **while** not converged **do**
4:     **for** number of tasks in batch **do**
5:         Sample task instance $\mathcal{T}_i \sim \mathcal{S}^{tr}$
6:         Let $(\mathcal{D}^{tr}, \mathcal{D}^{val}) = \mathcal{T}_i$
7:         Encode $\mathcal{D}^{tr}$ to $\mathbf{z}$ using $g_{\phi_e}$ and $g_{\phi_r}$
8:         Decode $\mathbf{z}$ to initial params $\theta_i$ using $g_{\phi_d}$
9:         Initialize $\mathbf{z'} = \mathbf{z}, \theta'_i = \theta_i$
10:       **for** number of adaptation steps **do**
11:           Compute training loss $\mathcal{L}_{\mathcal{T}_i}^{tr}(f_{\theta'_i})$
12:           Perform gradient step w.r.t. $\mathbf{z'}$: $\mathbf{z'} \leftarrow \mathbf{z'} - \alpha \nabla_{\mathbf{z'}} \mathcal{L}_{\mathcal{T}_i}^{tr}(f_{\theta'_i})$
13:           Decode $\mathbf{z'}$ to obtain $\theta'_i$ using $g_{\phi_d}$
14:       **end for**
15:       Compute validation loss $\mathcal{L}_{\mathcal{T}_i}^{val}(f_{\theta'_i})$
16:     **end for**
17:   Perform gradient step w.r.t $\phi$: $\phi \leftarrow \phi - \eta \nabla_\phi \sum_{\mathcal{T}_i} \mathcal{L}_{\mathcal{T}_i}^{val}(f_{\theta'_i})$
18: **end while**

# LEO Model

- Encoding and Relation Network

$$\boldsymbol{\mu}_n^e, \boldsymbol{\sigma}_n^e = \frac{1}{NK^2} \sum_{k_n=1}^{K} \sum_{m=1}^{N} \sum_{k_m=1}^{K} g_{\phi_r}\Big(g_{\phi_e}\left(\mathbf{x}_n^{k_n}\right), g_{\phi_e}\left(\mathbf{x}_m^{k_m}\right)\Big)$$

$$\mathbf{z}_n \sim q\left(\mathbf{z}_n | \mathcal{D}_n^{tr}\right) = \mathcal{N}\left(\boldsymbol{\mu}_n^e, diag(\boldsymbol{\sigma}_n^{e\,2})\right)$$

  ▶ Encoder all examples into intermediate codes, concatentated pair-wise
  ▶ Use the relation network to learn specific code for each class, forming a Gaussian distribution
  ▶ The hidden code $z$ is drawn from the distribution

- Decoding Network

$$\boldsymbol{\mu}_n^d, \boldsymbol{\sigma}_n^d = g_{\phi_d}\left(\mathbf{z}_n\right)$$

$$\mathbf{w}_n \sim p\left(\mathbf{w} | \mathbf{z}_n\right) = \mathcal{N}\left(\boldsymbol{\mu}_n^d, diag(\boldsymbol{\sigma}_n^{d\,2})\right)$$

  ▶ Decode the hidden codes into distribution's parameters, forming a Gaussian distribution
  ▶ Task-specific parameters are drawn from the resulted distribution

# LEO Model

- Inner-loop objective

$$\mathcal{L}_{\mathcal{T}_i}^{tr}\big(f_{\theta_i}\big) = \sum_{(\mathbf{x},y) \in \mathcal{D}^{tr}} \Big[ -\mathbf{w}_y \cdot \mathbf{x} + \log \Big( \sum_{j=1}^{N} e^{\mathbf{w}_j \cdot \mathbf{x}} \Big) \Big]$$

- Outer-loop objective

$$\min_{\phi_e, \phi_r, \phi_d} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \Big[ \mathcal{L}_{\mathcal{T}_i}^{val}\big(f_{\theta_i'}\big) + \beta D_{KL}\big(q(\mathbf{z}_n|\mathcal{D}_n^{tr})||p(\mathbf{z}_n)\big) + \gamma||\text{stopgrad}(\mathbf{z}_n') - \mathbf{z}_n||_2^2 \Big] + R$$

*Thank you !*