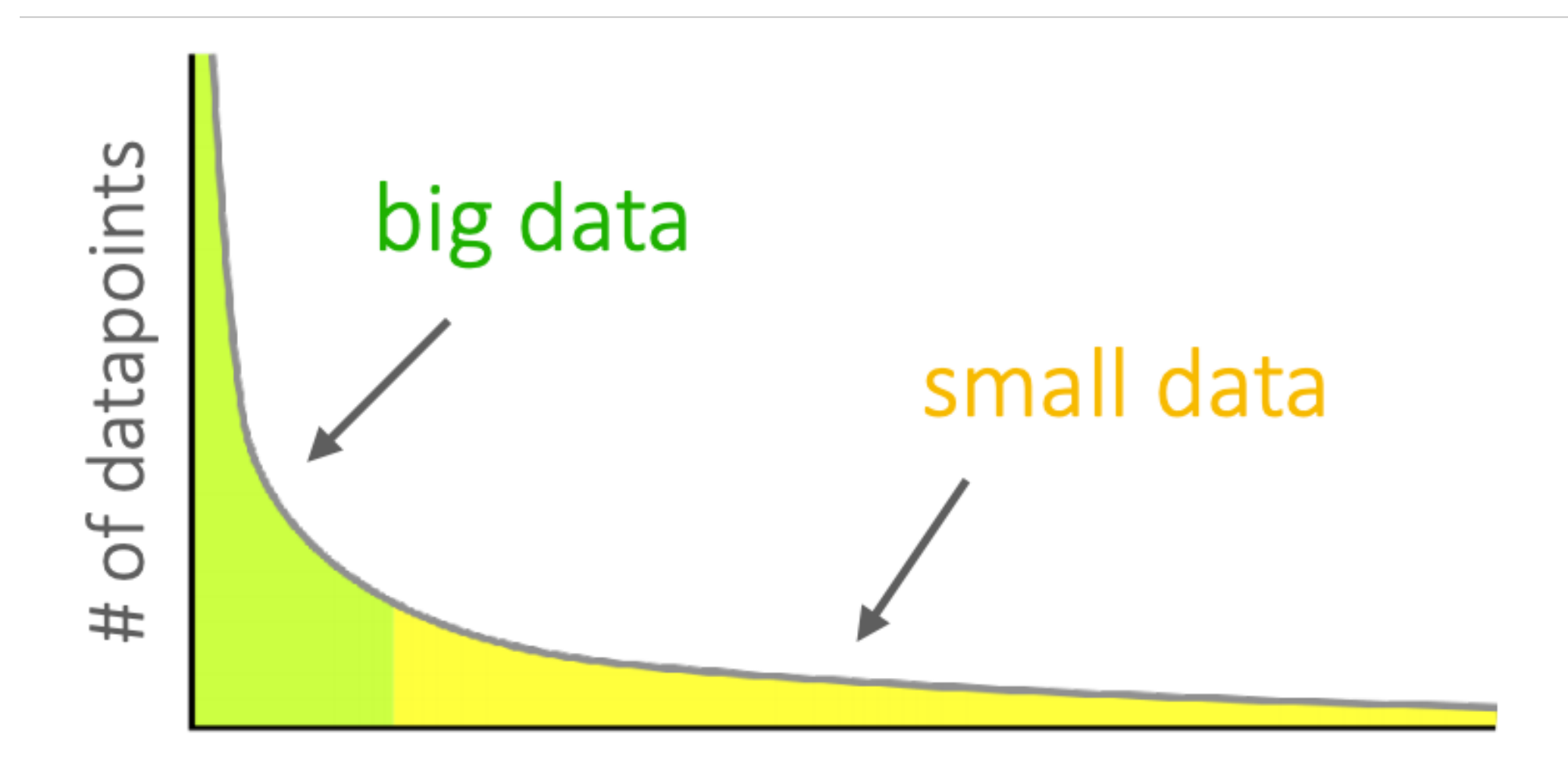


# Meta / Few-shot Learning

分享人：邹笑寒

2021.08.25

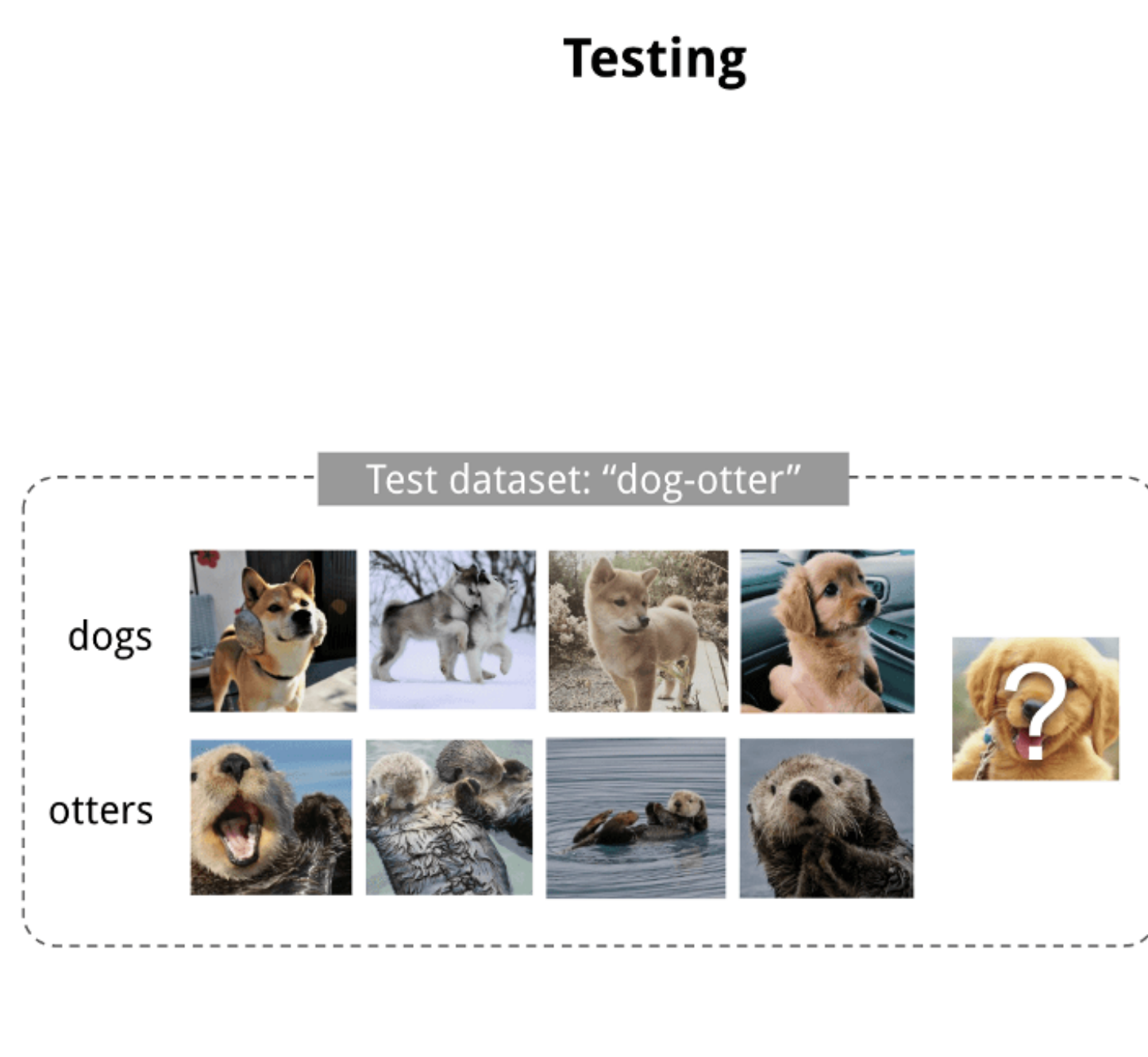
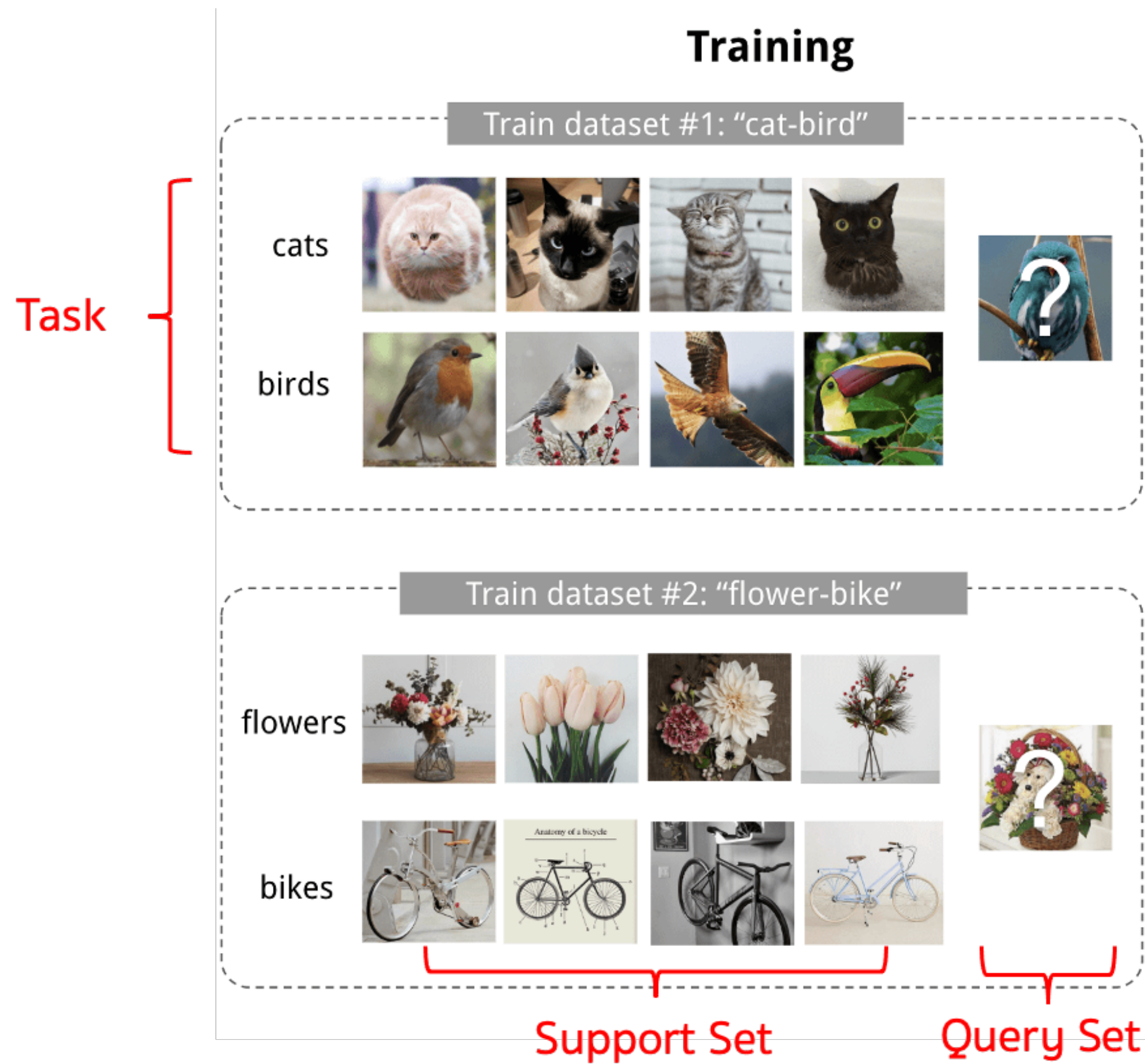
# Problem Definition



- task 多
- 每个 task 内的样本少

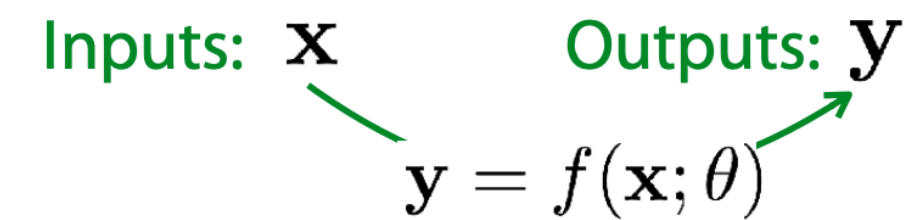
在已经学习了大量 task 后，如何在新的 task 上快速学习？

# Application on Supervised Learning: Few-shot Learning



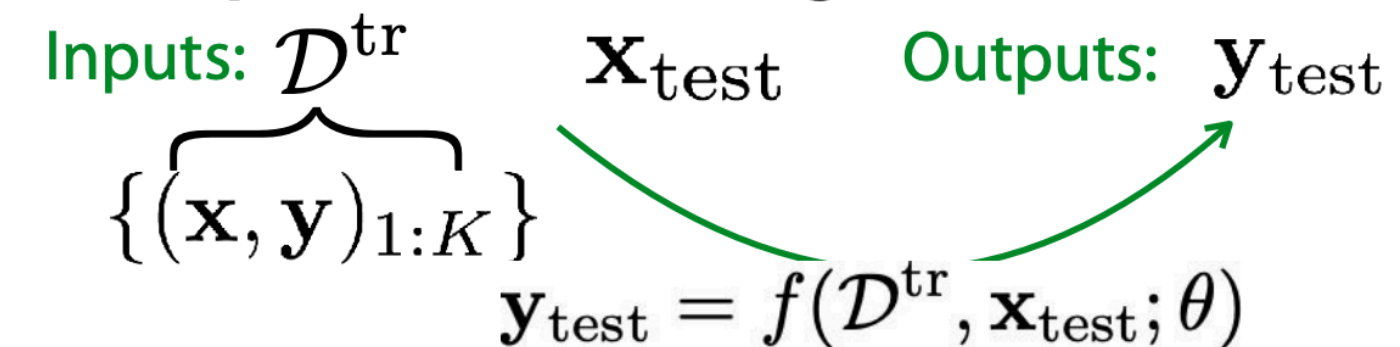
2-way 4-shot

**Supervised Learning:**



Data:  $\mathcal{D} = \{(\mathbf{x}, \mathbf{y})_i\}$

**Meta-Supervised Learning:**



Data:  $\mathcal{D}_{\text{meta-train}} = \{\mathcal{D}_i\}$

$\mathcal{D}_i : \{(\mathbf{x}, \mathbf{y})_j\}$

# Popular Approaches

Optimization-based

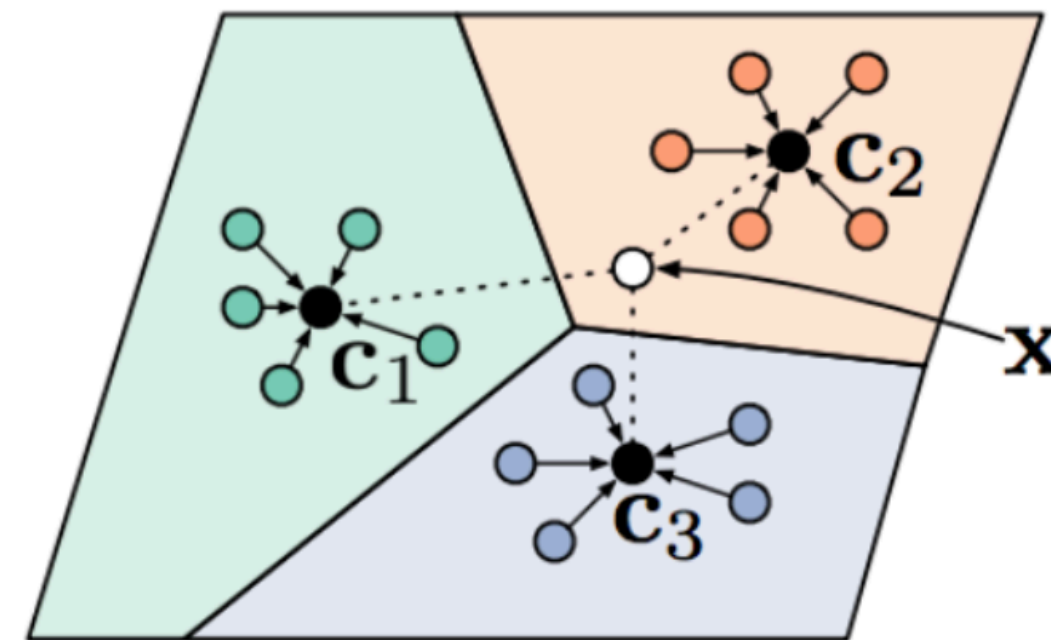
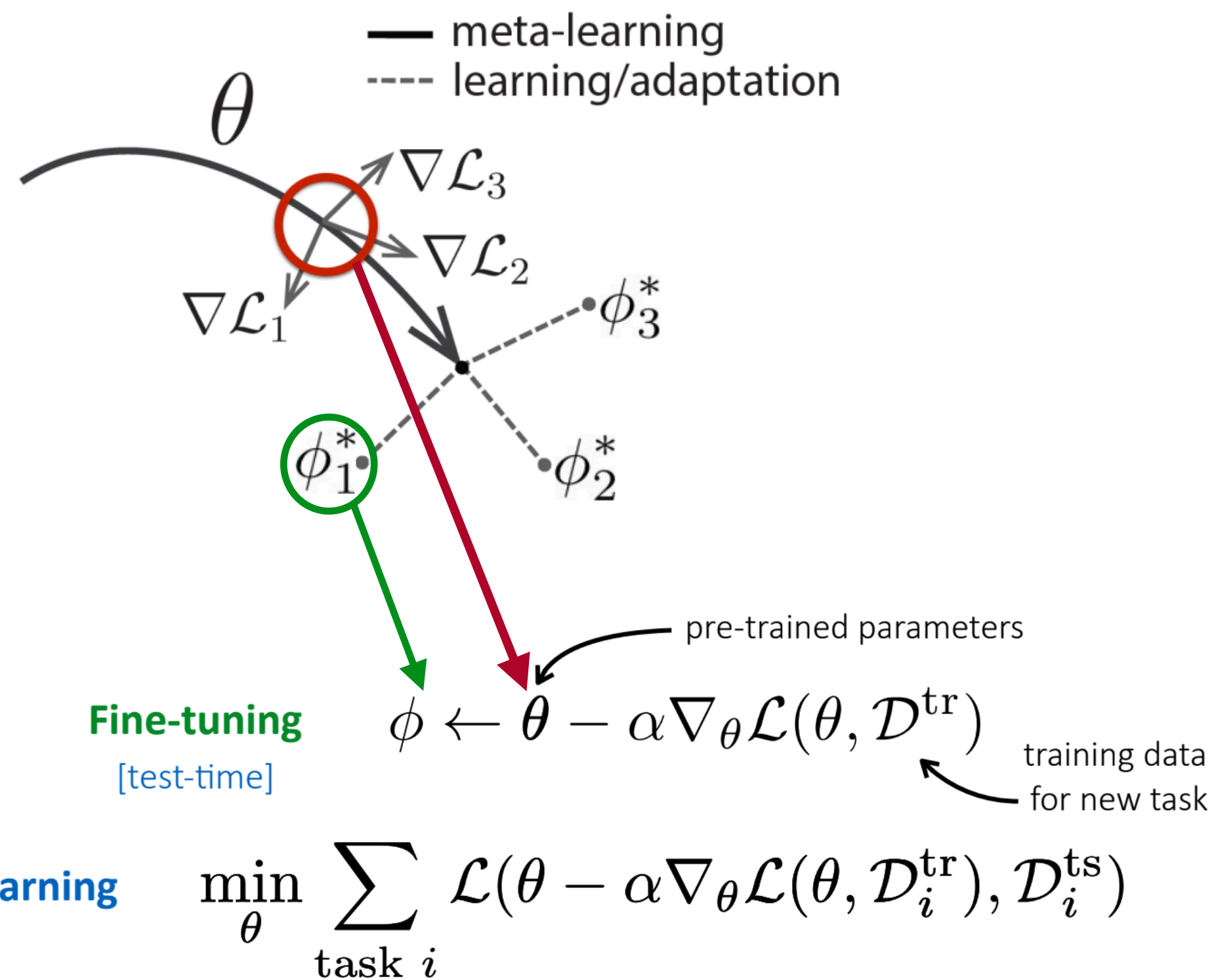
学习出模型的初始化参数

Metric-based

学习出样本间的距离函数（聚类）

Others

AutoML?



$$y^{\text{ts}} = f_{\text{PN}}(\mathcal{D}_i^{\text{tr}}, x^{\text{ts}})$$

$$= \text{softmax}(-d(f_{\theta}(x), c_k))$$

# Optimization-based: MAML

---

## Algorithm 1 Model-Agnostic Meta-Learning

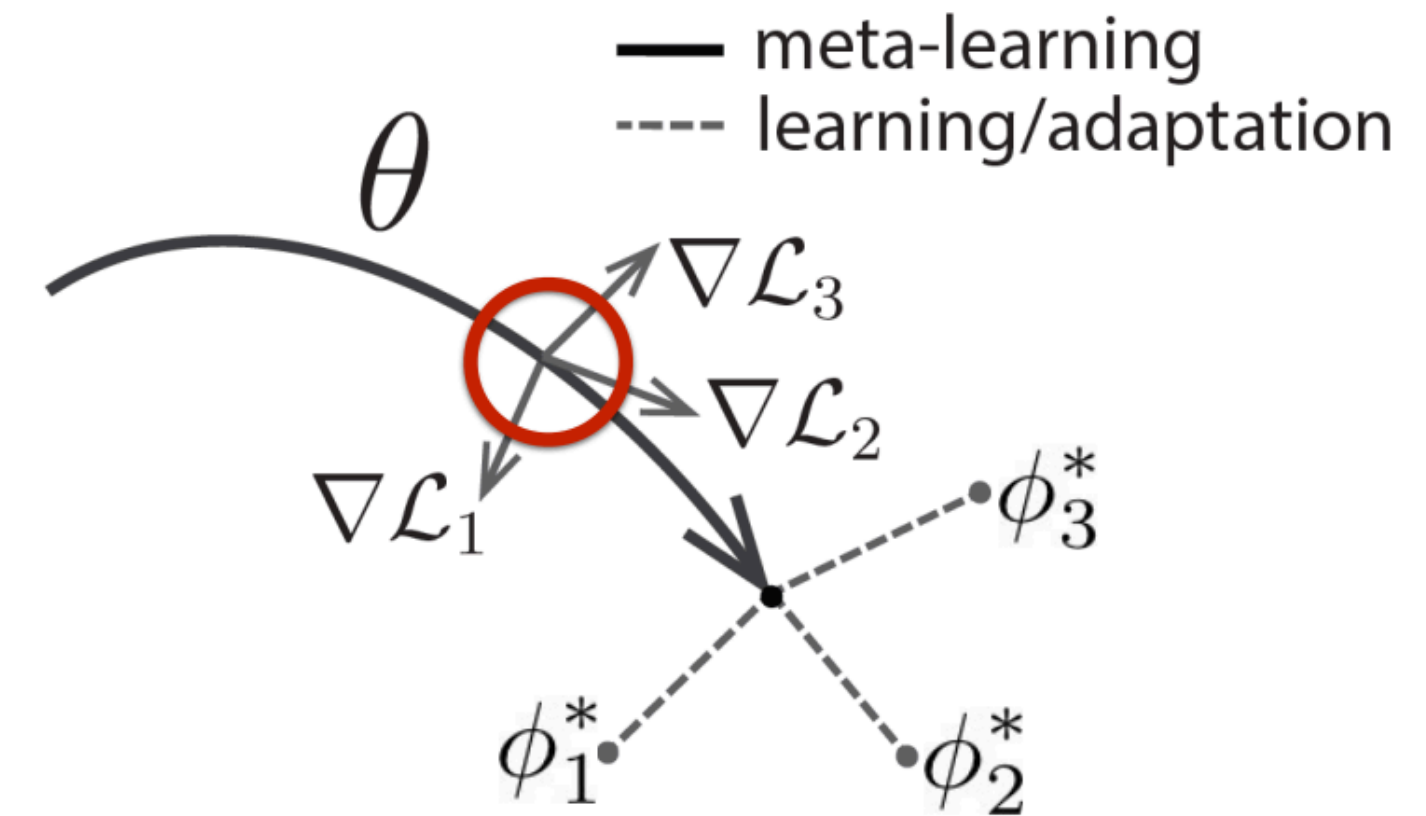
---

**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:     Evaluate  $\nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$  with respect to  $K$  examples
- 6:     Compute adapted parameters with gradient descent:  $\theta'_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$  Inner Loop
- 7:   **end for**
- 8:   Update  $\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i})$  Outer Loop
- 9: **end while**

Bi-level Optimization




---

## Algorithm 2 Generic First-order Model-agnostic Meta-learning algorithm

---

**Require:**  $p(\tau)$  : Distribution over tasks

**Require:**  $\alpha, \beta$  : step hyper-parameters

*Initialisation* : Random  $\theta$

- 1: **for**  $i = 1, 2 \dots n$  **do**
- 2:   sample tasks  $\tau_i \sim p(\tau)$
- 3:   **for all**  $\tau_i$  **do**
- 4:     Evaluate  $\nabla_{\theta} L(\theta)$  with  $K$  examples
- 5:     Compute adapted param with GD:  $\theta_i = \theta - \alpha \nabla_{\theta} L(\theta)$
- 6:   **end for**
- 7:   update:  $\theta \leftarrow \theta - \beta \nabla_{\theta} I(\theta_i)$
- 8: **end for**

First-order Version

# Optimization-based: Reptile

---

## Algorithm 1 Reptile (serial version)

---

Initialize  $\phi$ , the vector of initial parameters

**for** iteration = 1, 2, ... **do**

    Sample task  $\tau$ , corresponding to loss  $L_\tau$  on weight vectors  $\tilde{\phi}$

    Compute  $\tilde{\phi} = U_\tau^k(\phi)$ , denoting  $k$  steps of SGD or Adam

    Update  $\phi \leftarrow \phi + \epsilon(\tilde{\phi} - \phi)$

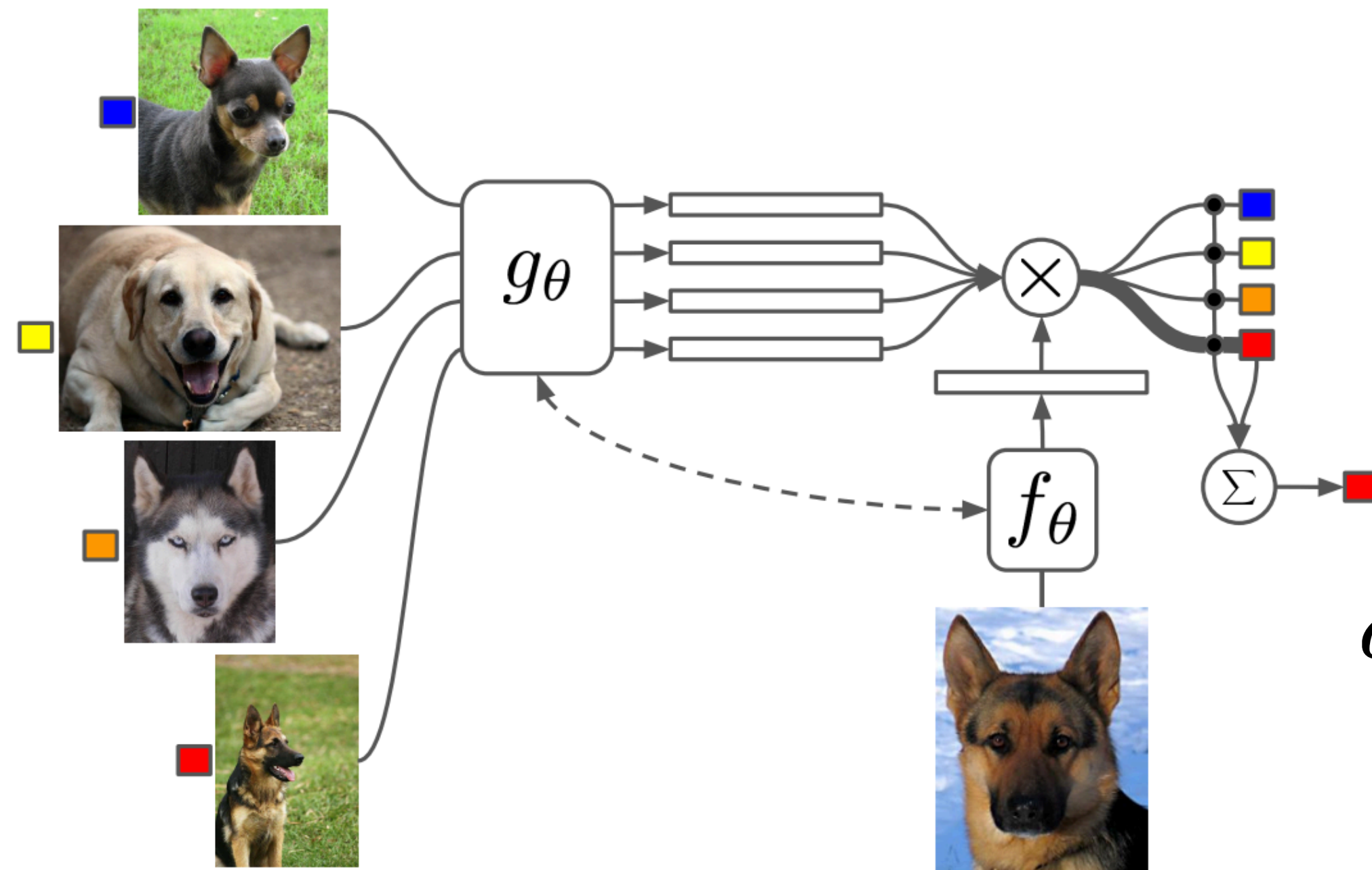
**end for**

---

$$\begin{aligned} g_{\text{MAML}} &= \bar{g}_2 - \alpha \bar{H}_2 \bar{g}_1 - \alpha \bar{H}_1 \bar{g}_2 + O(\alpha^2) \\ g_{\text{FOMAML}} = g_2 &= \bar{g}_2 - \alpha \bar{H}_2 \bar{g}_1 + O(\alpha^2) \\ g_{\text{Reptile}} = g_1 + g_2 &= \bar{g}_1 + \bar{g}_2 - \alpha \bar{H}_2 \bar{g}_1 + O(\alpha^2) \end{aligned}$$

增大不同 batch 上的梯度的内积 = 减小夹角

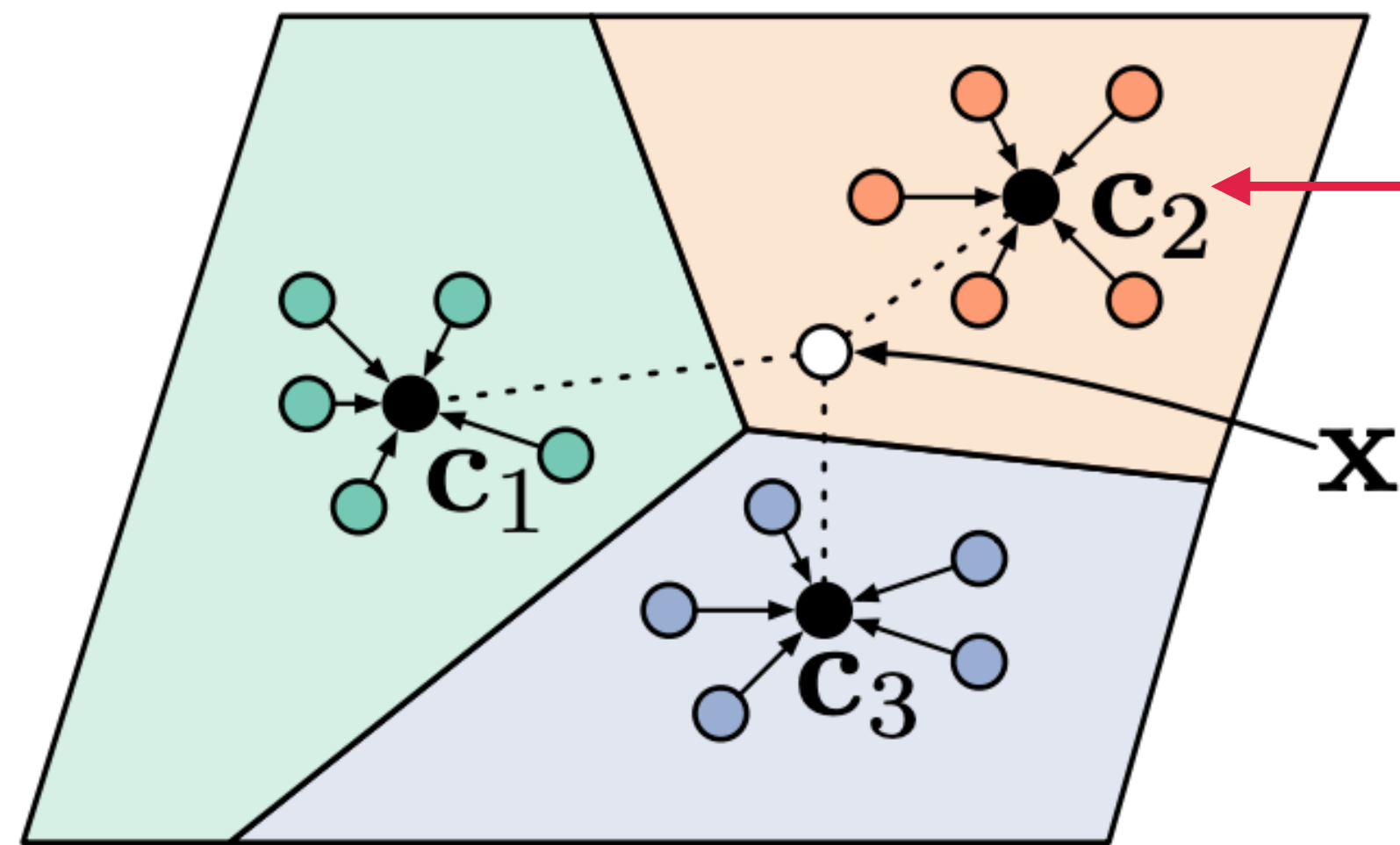
# Metric-based: Matching Networks



$$P(\hat{y}|\hat{x}, S) = \sum_{i=1}^k a(\hat{x}, x_i) y_i$$

$$a(\hat{x}, x_i) = \text{softmax}\left(-d(f_\theta(\hat{x}), g_\theta(x_i))\right)$$

# Metric-based: Prototypical Networks



$$\mathbf{c}_k = \frac{1}{|\mathcal{D}_i^{\text{tr}}|} \sum_{(x,y) \in \mathcal{D}_i^{\text{tr}}} f_{\theta}(x)$$

$$p_{\theta}(y = k|x) = \frac{\exp(-d(f_{\theta}(x), \mathbf{c}_k))}{\sum_{k'} \exp(-d(f_{\theta}(x), \mathbf{c}_{k'}))}$$



# How to Enhance Metric-based Methods

- Feature Extractor
- Prototypes
- Distance Computing
- Additional information (Cross-modal)

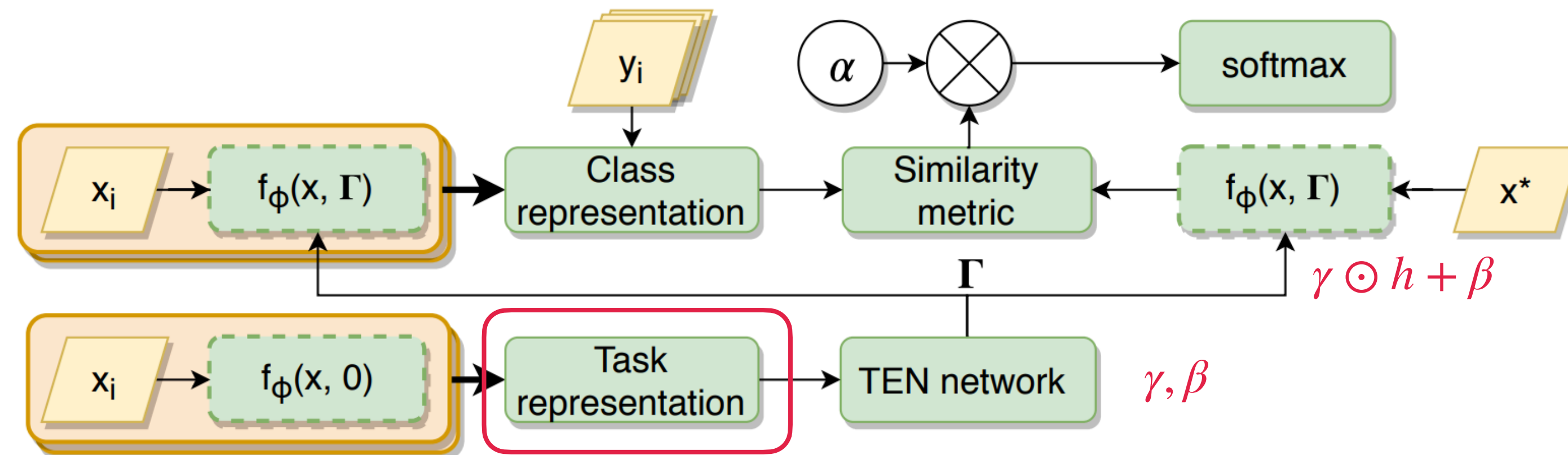
# Metric-based: TADAM

Learnable temperature

Metric Scaling

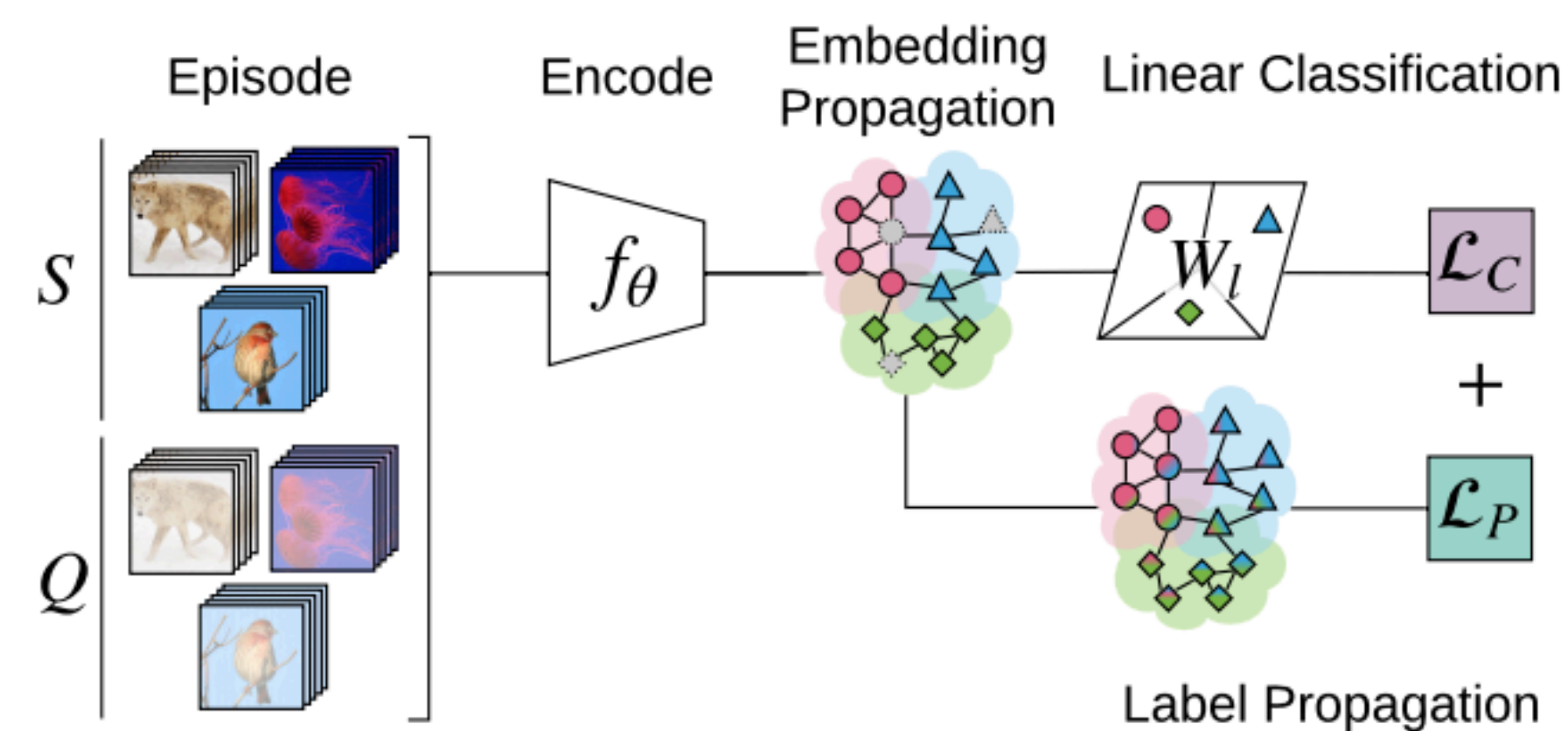
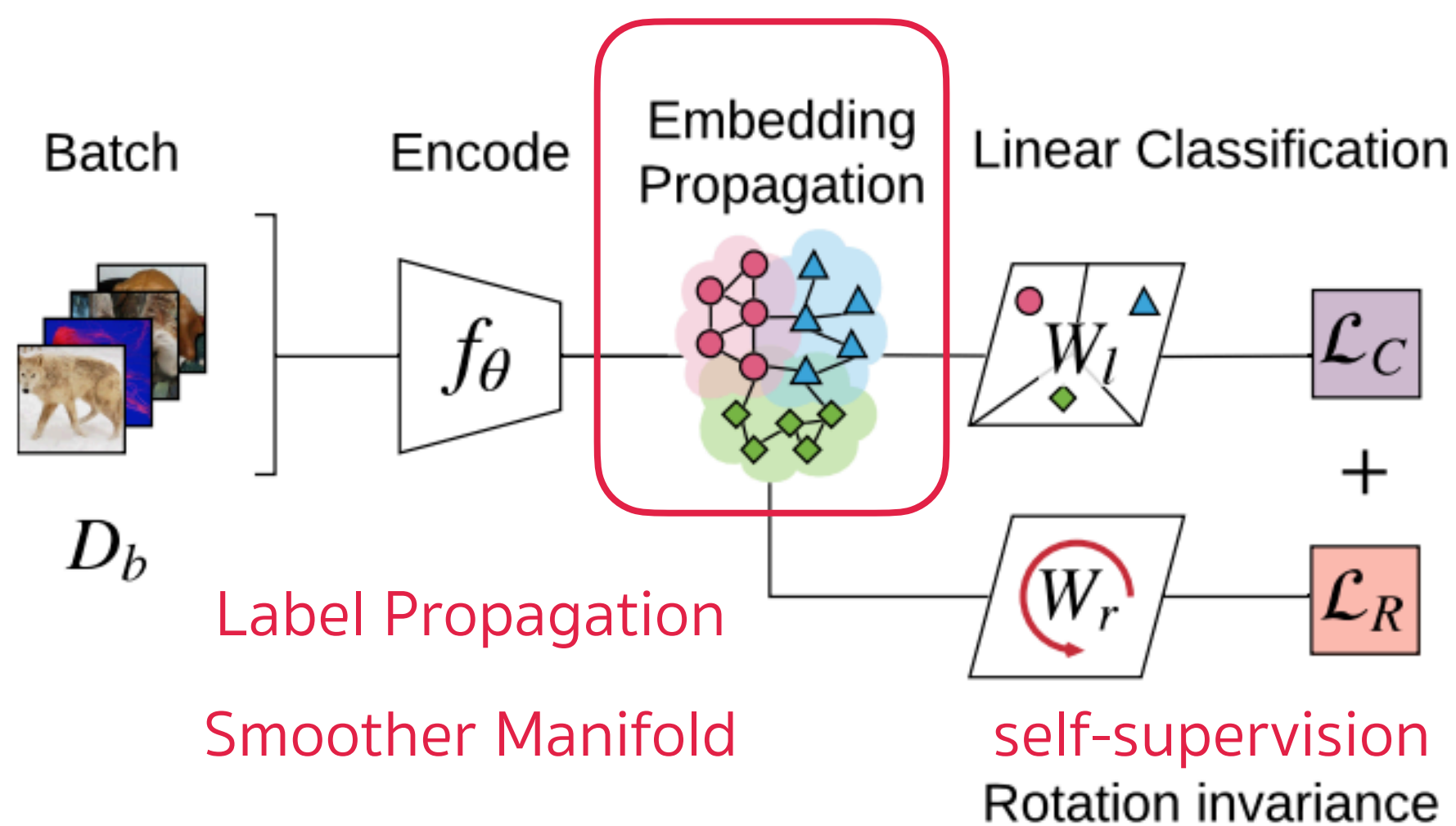
$$p_{\phi, \alpha}(y = k | \mathbf{x}) = \text{softmax}(-\alpha d(\mathbf{z}, \mathbf{c}_k))$$

Task conditioning



mean of the class prototypes:  $\bar{\mathbf{c}} = \frac{1}{K} \sum_k \mathbf{c}_k$

# Metric-based: EPNet



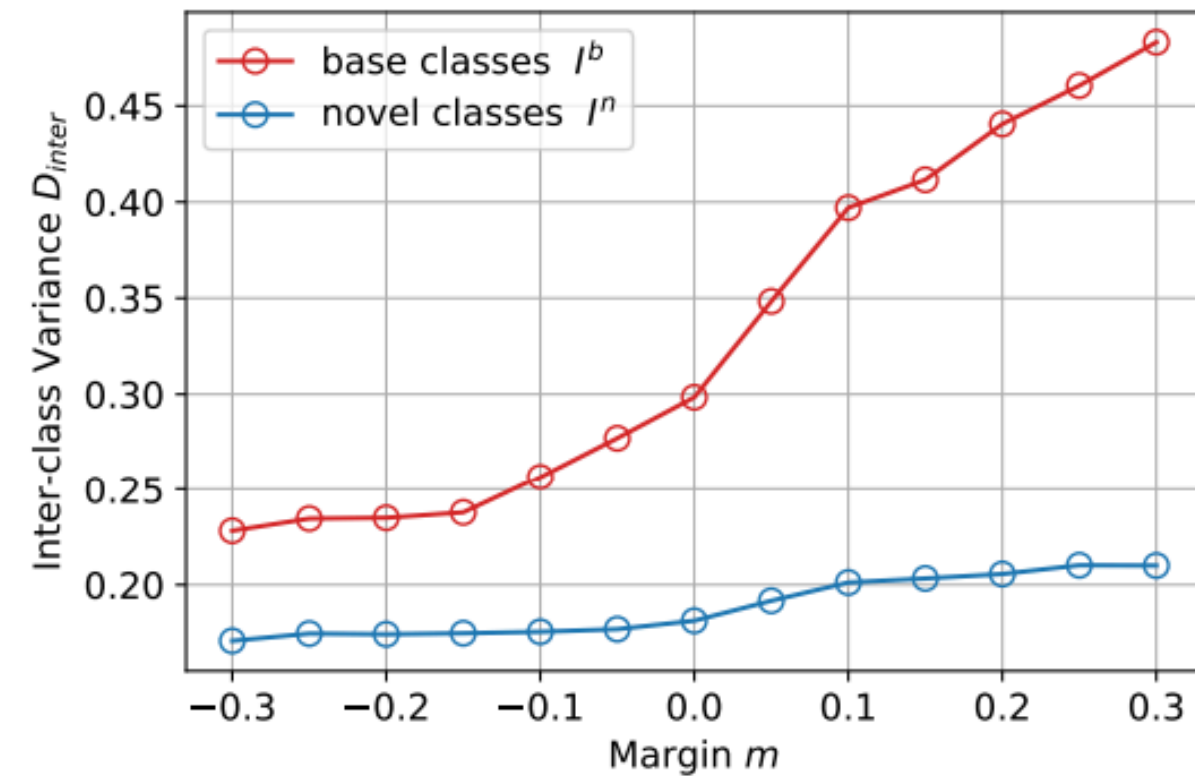
$$\mathcal{L}_c(\mathbf{x}_i, y_i; W_l, \theta) = -\ln p(y_i | \tilde{\mathbf{z}}_i, W_l)$$

$$\mathcal{L}_r(\mathbf{x}_i, r_j; W_r, \theta) = -\ln p(r_j | \tilde{\mathbf{z}}_i, W_r)$$

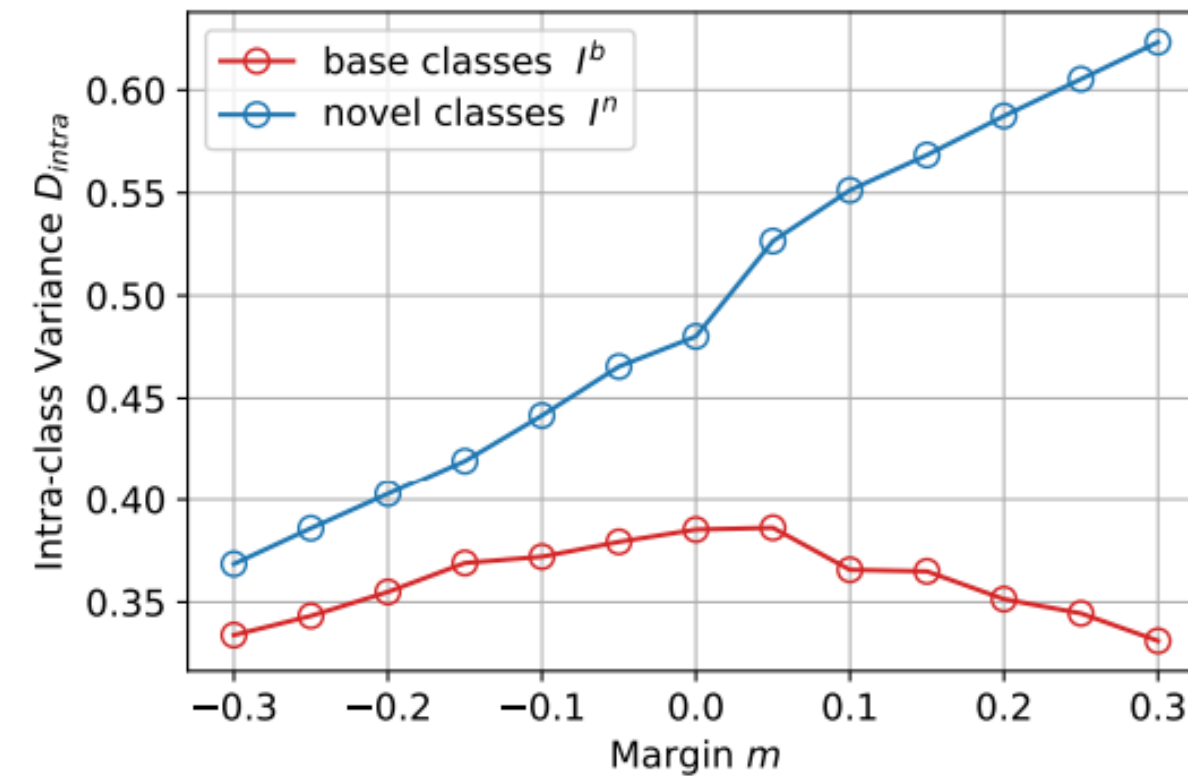
$$\operatorname{argmin}_{\theta, W_l, W_r} \sum_{i=1}^{128} \sum_{j=1}^4 \mathcal{L}_c(\mathbf{x}_i, y_i; W_l, \theta) + \mathcal{L}_r(\mathbf{x}_i, r_j; W_r, \theta)$$

$$\operatorname{argmin}_{\theta, W_l} \left[ \frac{1}{|Q|} \sum_{(\mathbf{x}_i, y_i) \in Q} \mathcal{L}_p(\mathbf{x}_i, y_i; \theta) + \frac{1}{|S \cup Q|} \sum_{(\mathbf{x}_i, y_i) \in S \cup Q} \frac{1}{2} \mathcal{L}_c(\mathbf{x}_i, y_i; W_l, \theta) \right]$$

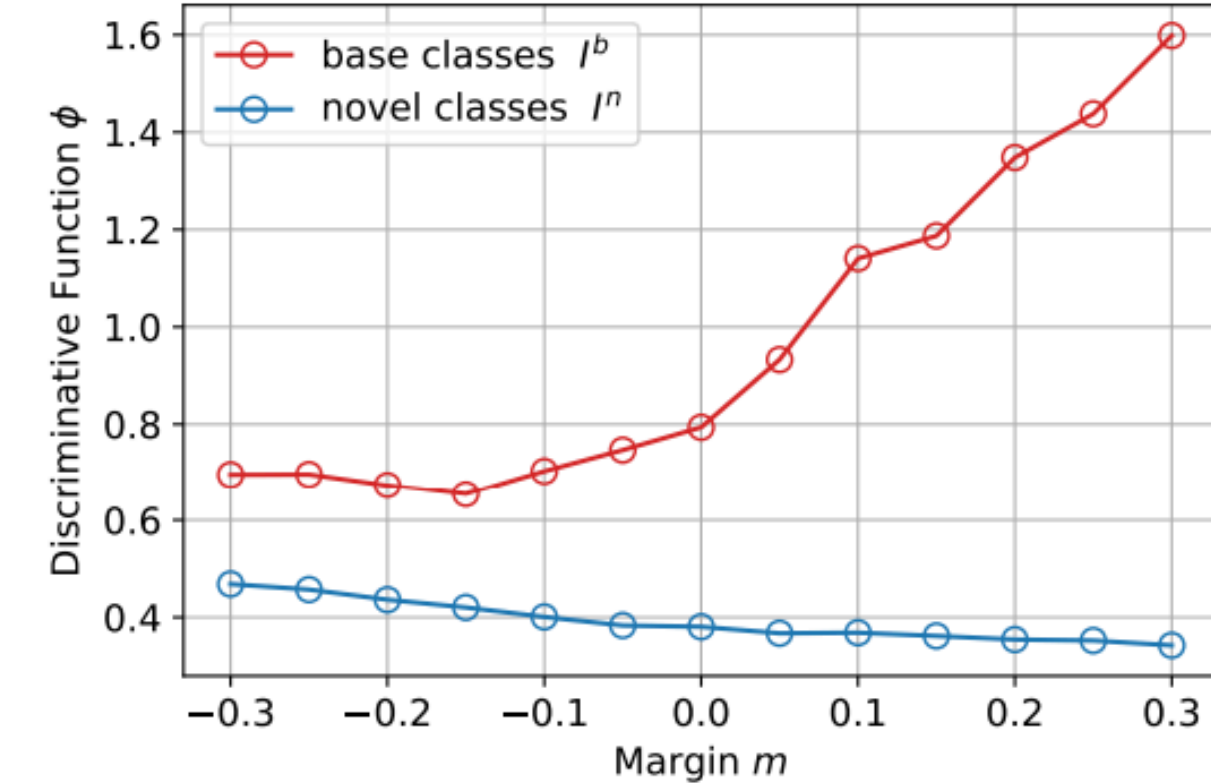
# Metric-based: Negative-Margin Softmax Loss



(a)  $D_{\text{inter}}$



(b)  $D_{\text{intra}}$

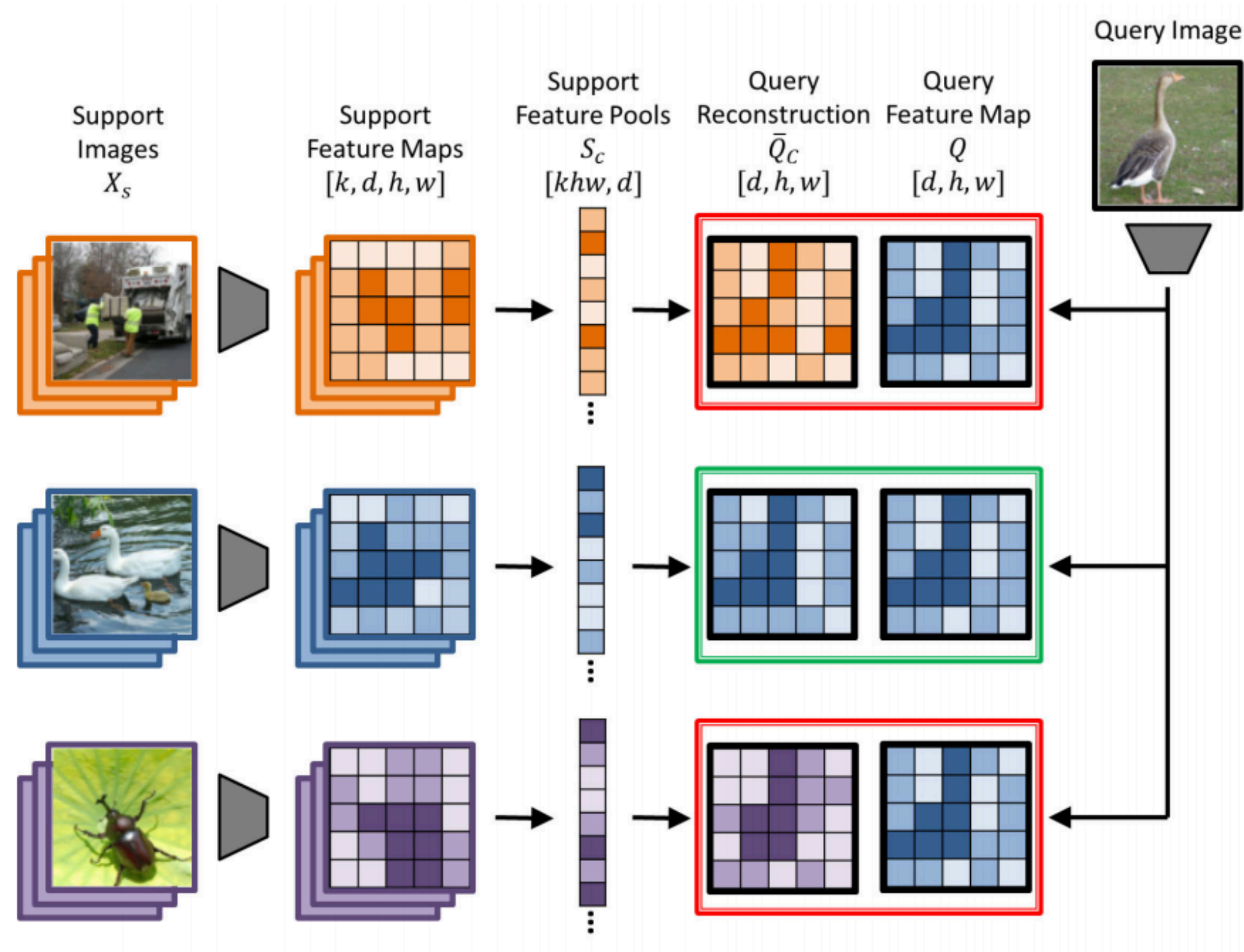


(c)  $\phi$   $\phi(I, m) = \frac{D_{\text{inter}}(I, m)}{D_{\text{intra}}(I, m)}$

negative margin (hyperparameter)

$$L = -\frac{1}{N} \sum_{i=1}^N \log \frac{e^{\beta \cdot (\mathbf{s}(z_i, W_{y_i}) - m)}}}{e^{\beta \cdot (\mathbf{s}(z_i, W_{y_i}) - m)} + \sum_{j=1, j \neq y_i}^C e^{\beta \cdot \mathbf{s}(z_i, W_j)}}$$

# Metric-based: FRN



Reconstruct feature maps to preserve location details

Find a matrix  $W \in \mathbb{R}^{r \times kr}$  such that  $WS_c \approx Q$

Closed-form solution:

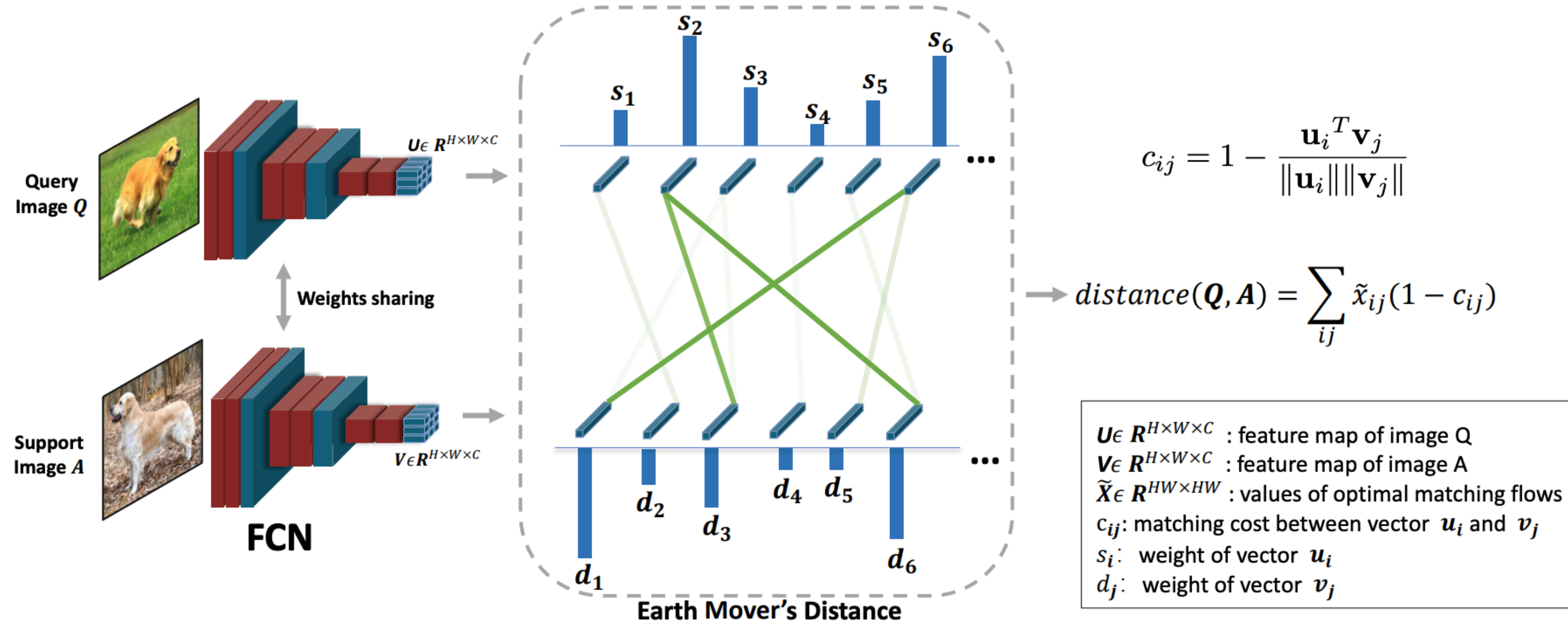
$$\bar{W} = \arg \min_W \|Q - WS_c\|^2 + \lambda \|W\|^2$$

$$\bar{W} = QS_c^T (S_c S_c^T + \lambda I)^{-1}$$

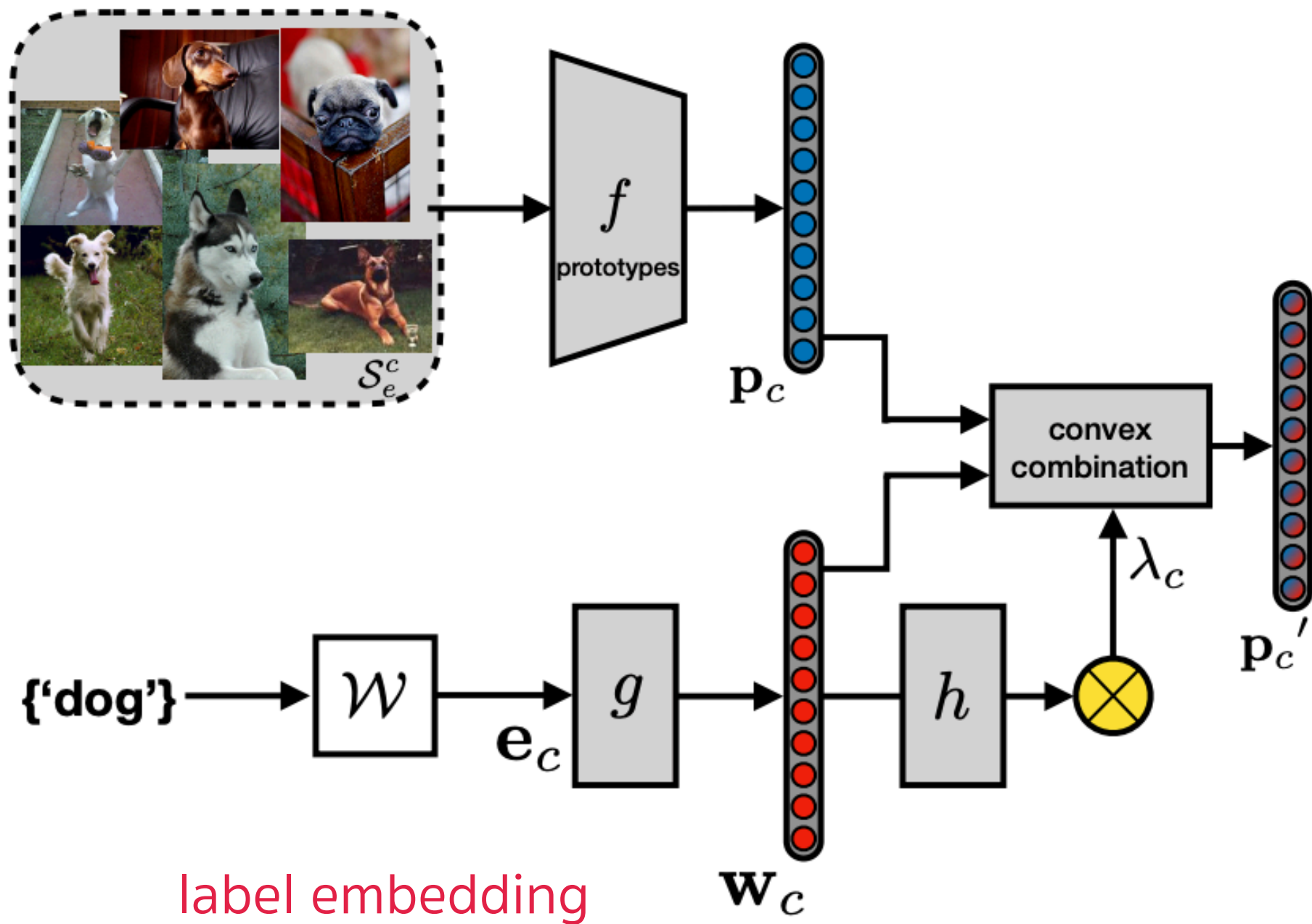
$$\bar{Q}_c = \bar{W} S_c$$

$$P(y_q = c | x_q) = \frac{e^{-\gamma \langle Q, \bar{Q}_c \rangle}}{\sum_{c' \in C} e^{-\gamma \langle Q, \bar{Q}_{c'} \rangle}}$$

# Metric-based: DeepEMD



# Metric-based: AM3

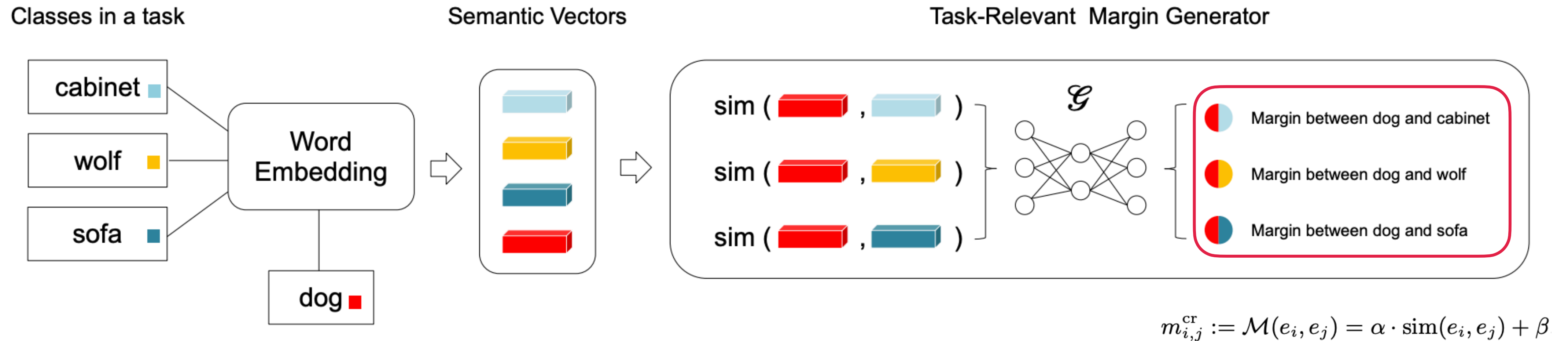


$$\mathbf{p}'_c = \lambda_c \cdot \mathbf{p}_c + (1 - \lambda_c) \cdot \mathbf{w}_c$$

$$\lambda_c = \frac{1}{1 + \exp(-h(\mathbf{w}_c))}$$

$$p_\theta(y = c | q_t, S_e, \mathcal{W}) = \frac{\exp(-d(f(q_t), \mathbf{p}'_c))}{\sum_k \exp(-d(f(q_t), \mathbf{p}'_k))}$$

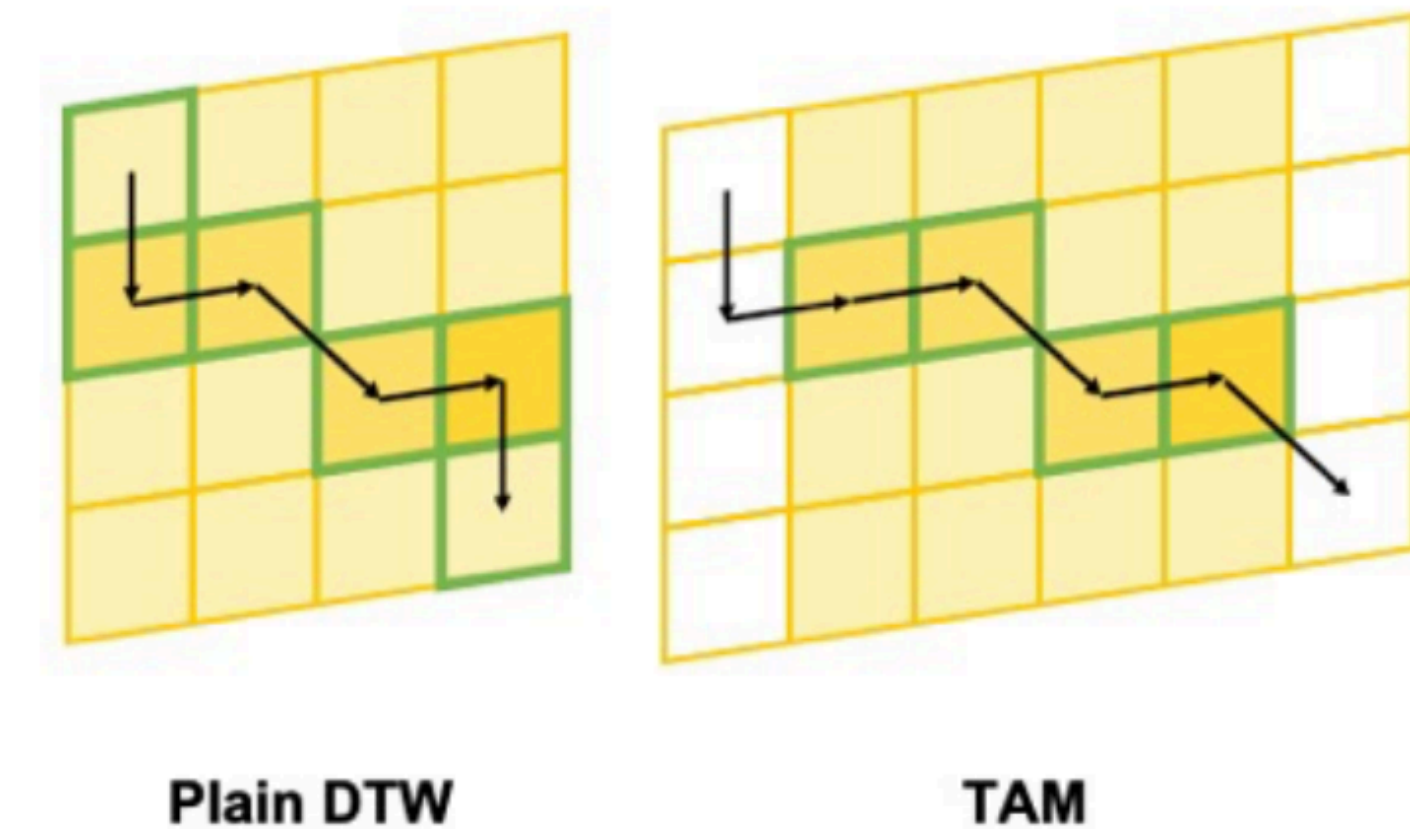
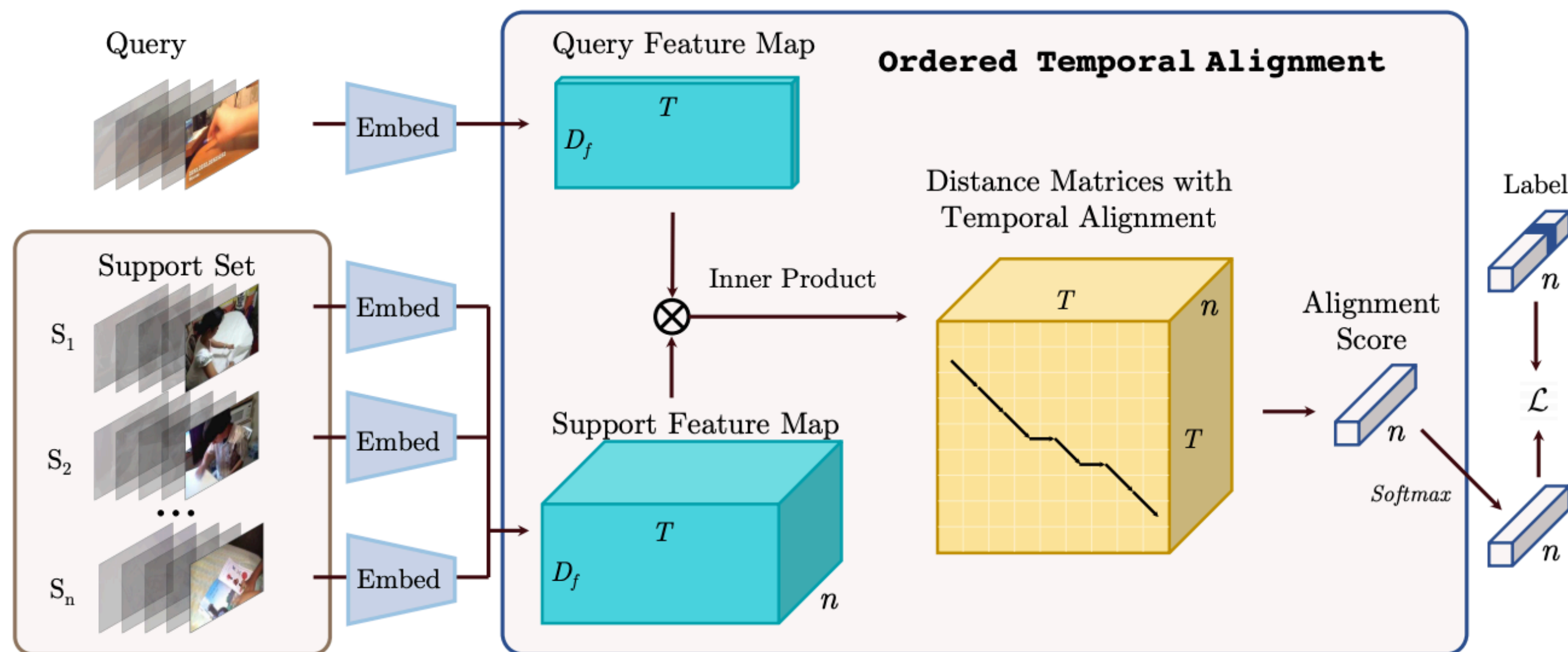
# Metric-based: TRAML



$$p^{\text{tr}}(y|x, S) = \frac{e^{\mathcal{D}(\mathcal{F}(x), r_y)}}{e^{\mathcal{D}(\mathcal{F}(x), r_y)} + \sum_{k \in C_t \setminus \{y\}} e^{\mathcal{D}(\mathcal{F}(x), r_k) + m_{y,k}^{\text{tr}}}}$$



# Metric-based: OTAM (Video Classification)



$$\gamma(i, j) = D_{ij} + \begin{cases} \min\{\gamma(i-1, j-1), \gamma(i-1, j), \gamma(i, j-1)\}, & j = 0 \text{ or } j = T+1 \\ \min\{\gamma(i-1, j-1), \gamma(i, j-1), \text{otherwise}\} & \text{otherwise} \end{cases}$$

Thanks!