



The Family of Objective Functions in DML

CVPR Tutorial: Deep Visual Similarity and Metric Learning

Timo Milbich

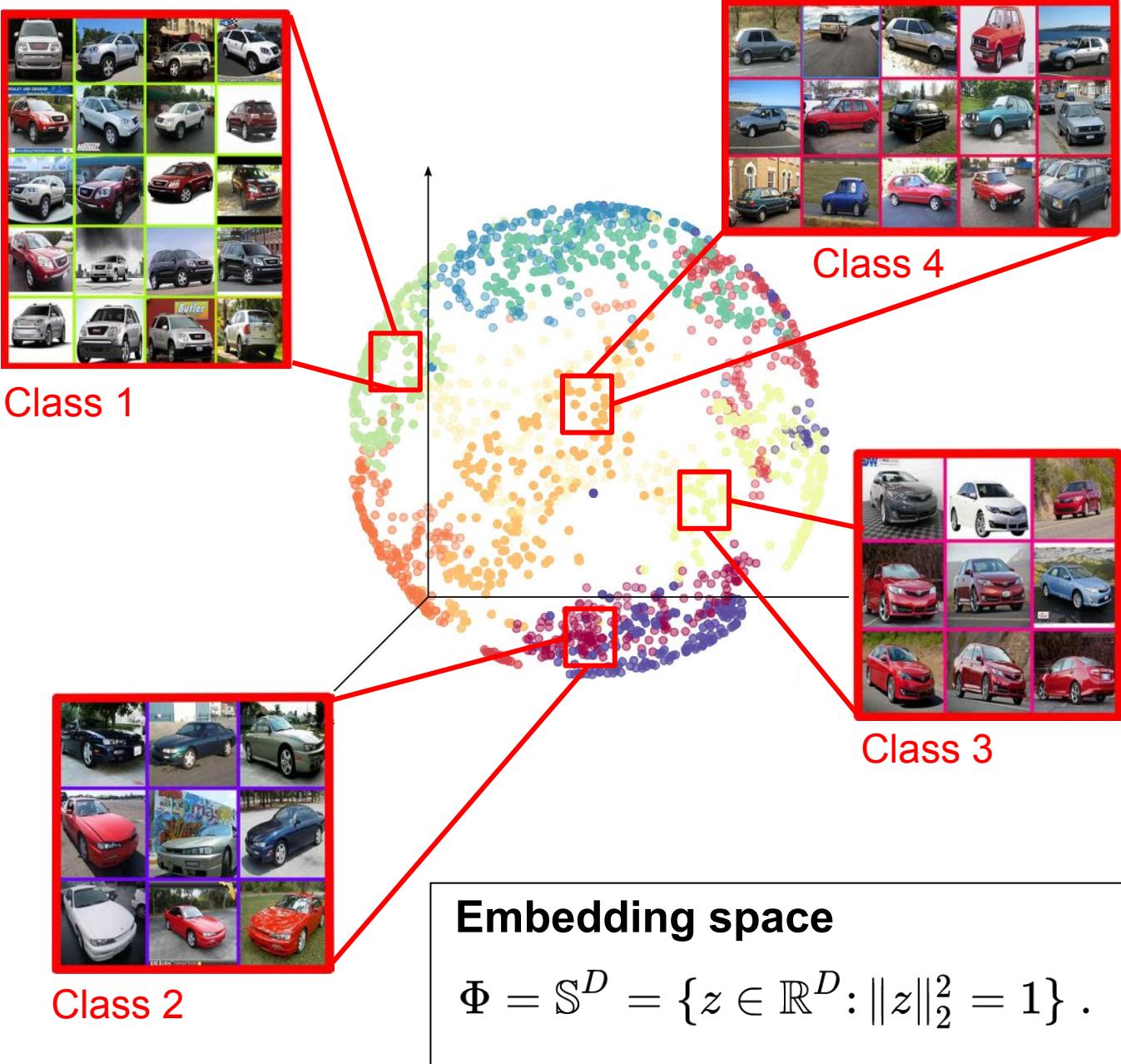
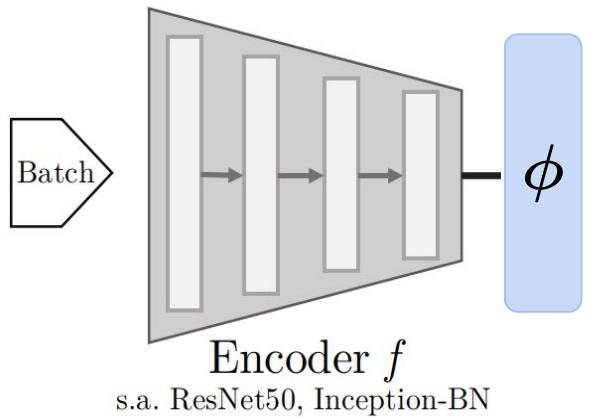
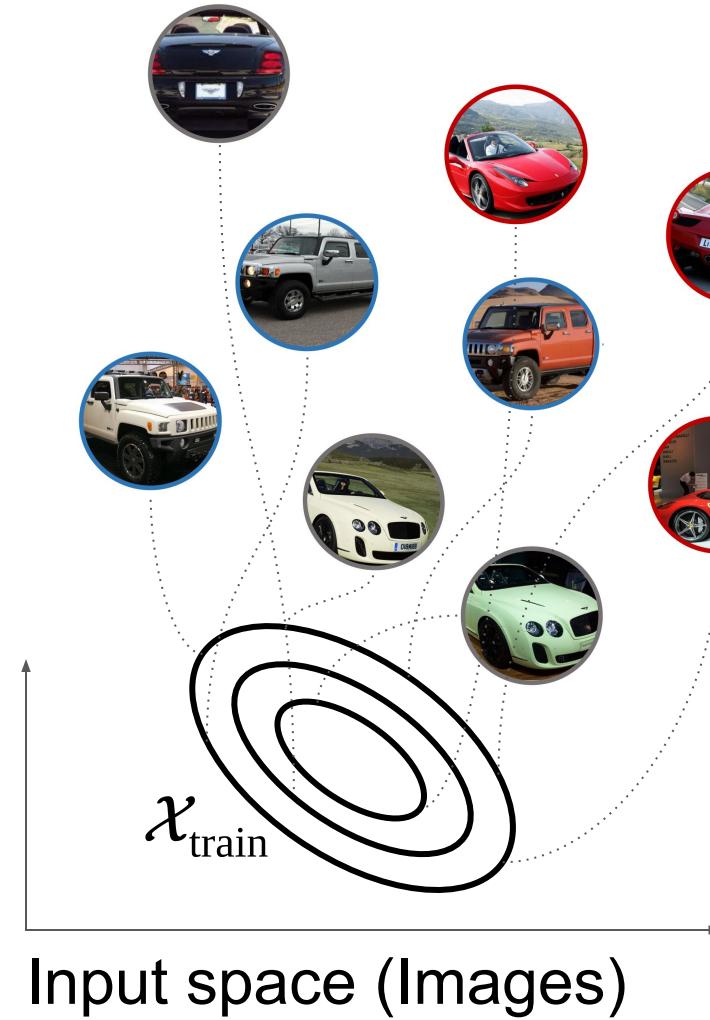


Technical
University
of Munich



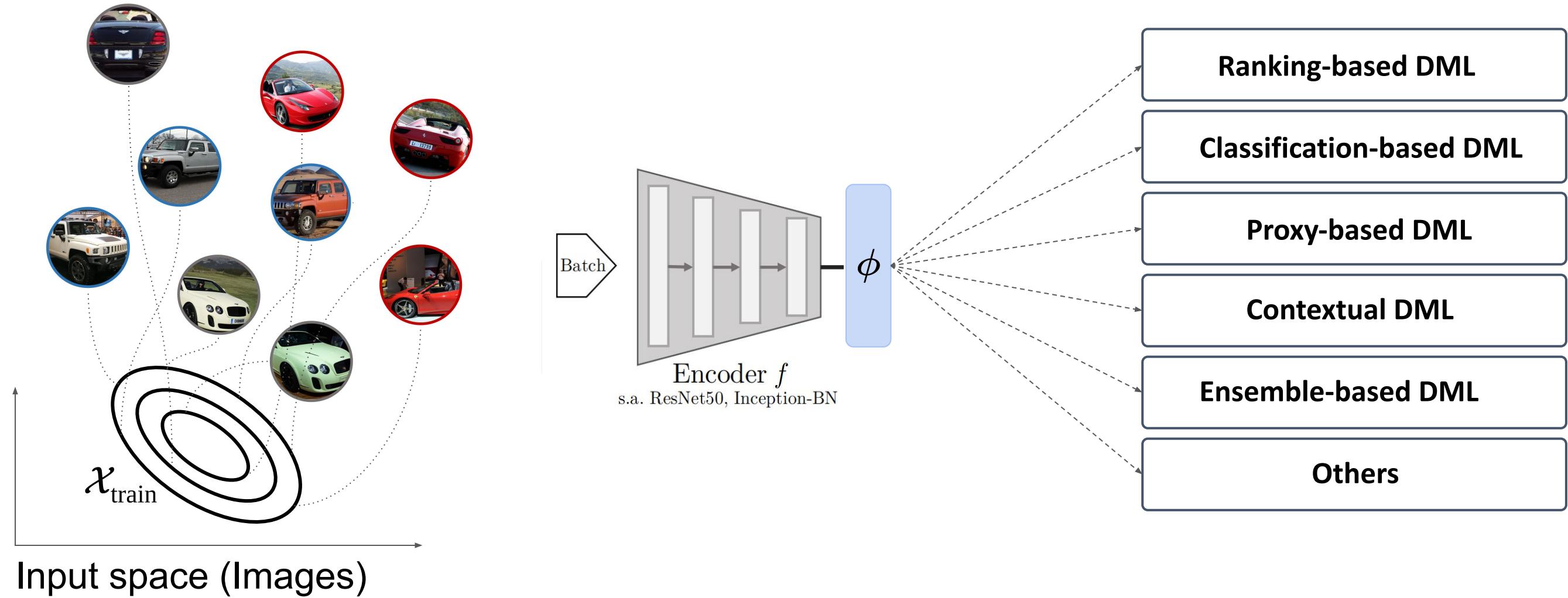
Visual Similarity Learning

Learn **representation** $\phi(x)$ which reflects **semantic similarity** $d(\phi(x_i), \phi(x_j))$ within training distribution $\mathcal{X}_{\text{train}}$.



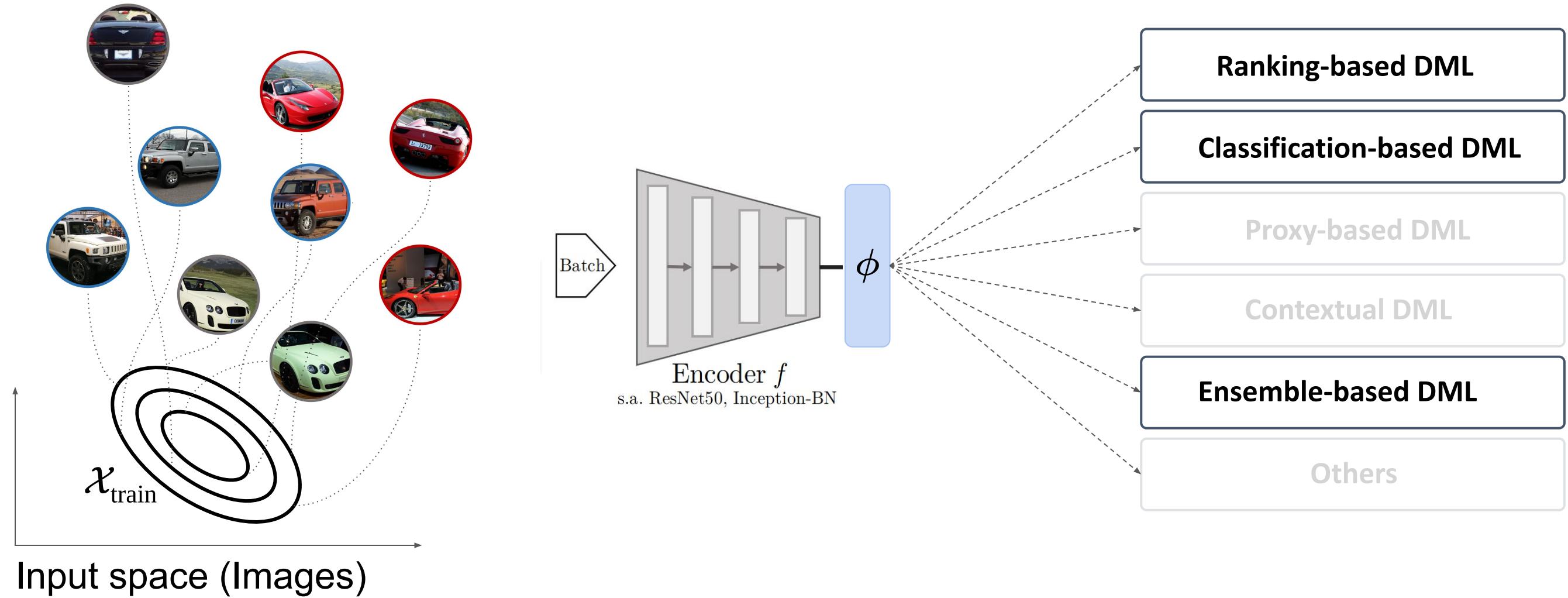
Deep Metric Learning (DML)

Learn **representation** $\phi(x)$ which reflects **semantic similarity** $d(\phi(x_i), \phi(x_j))$ within training distribution $\mathcal{X}_{\text{train}}$.



Deep Metric Learning (DML)

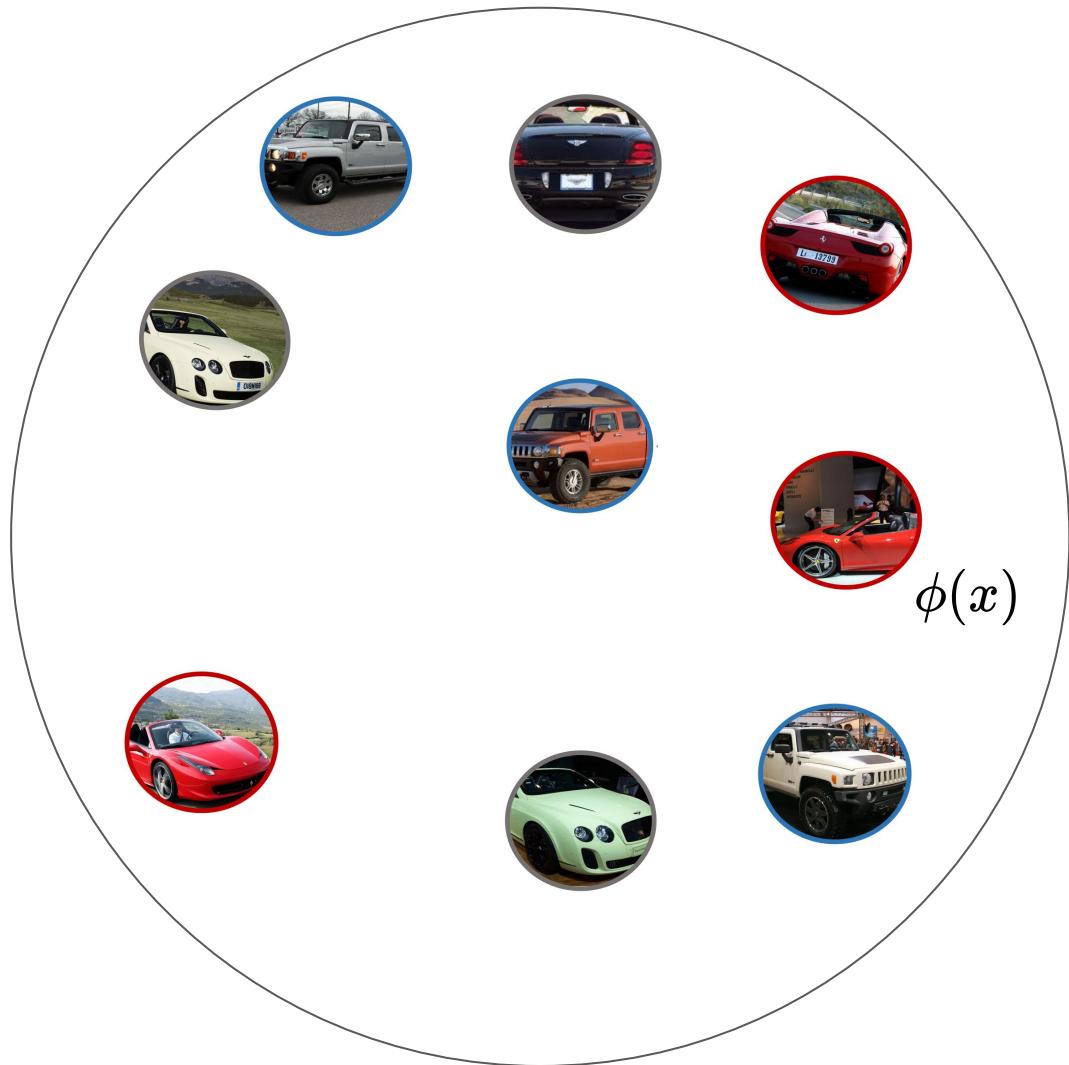
Learn **representation** $\phi(x)$ which reflects **semantic similarity** $d(\phi(x_i), \phi(x_j))$ within training distribution $\mathcal{X}_{\text{train}}$.



Ranking-based DML

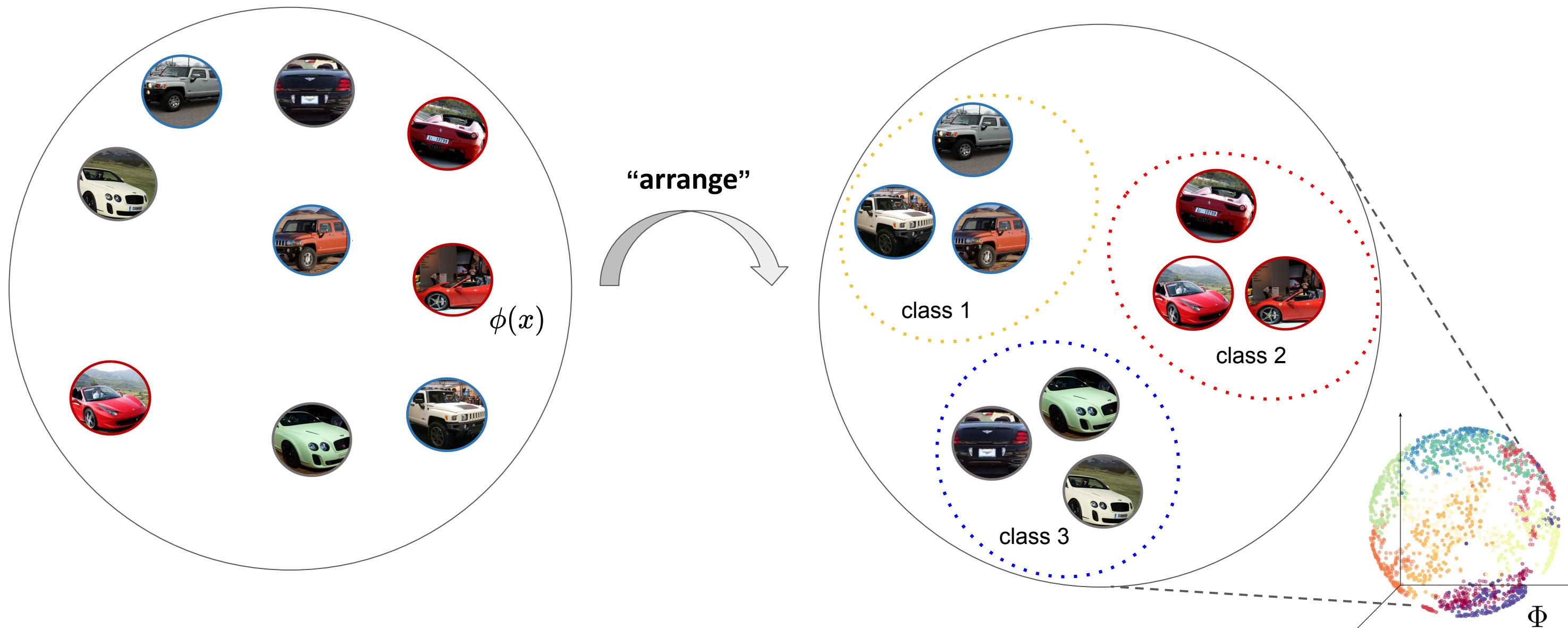
Ranking-based DML

- Learning informative embedding space requires to “arrange” data under $\phi(x)$



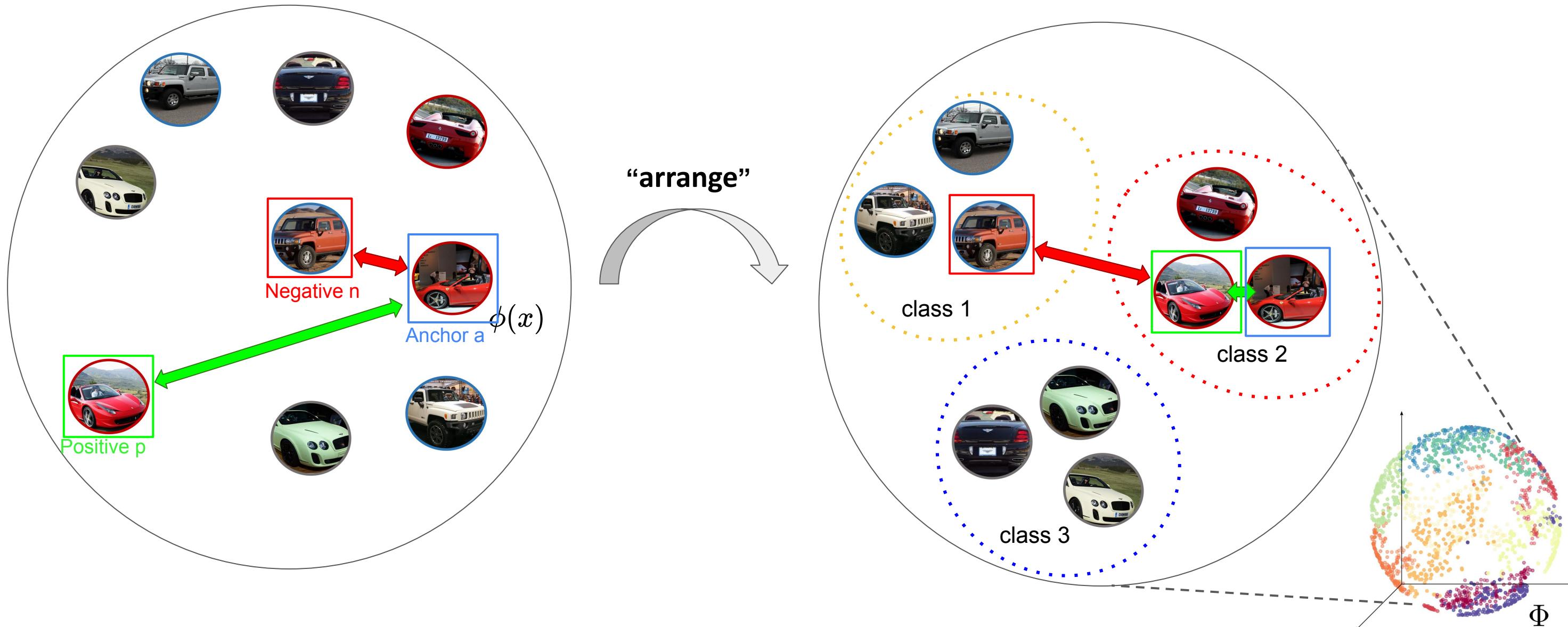
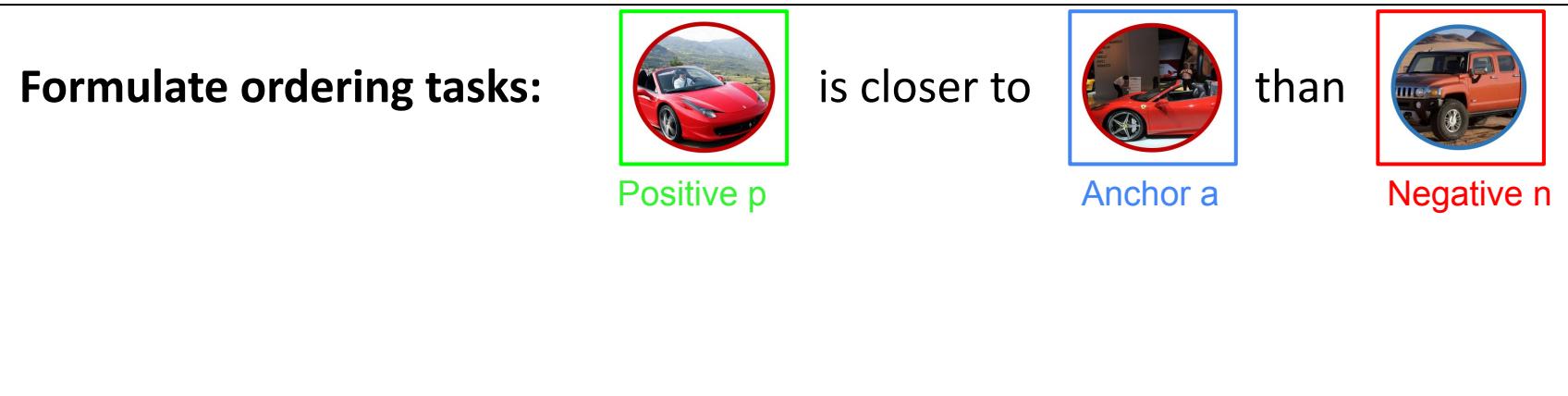
Ranking-based DML

- Learning informative embedding space requires to “arrange” data under $\phi(x)$



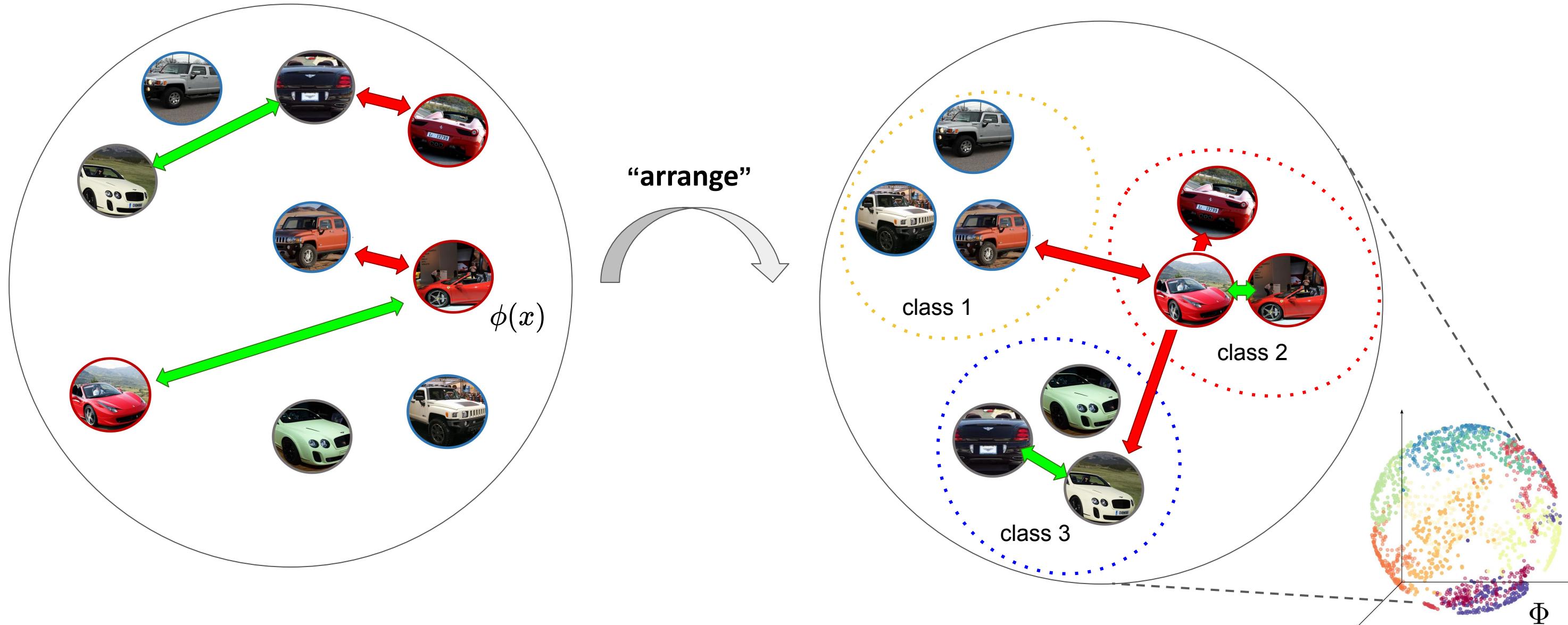
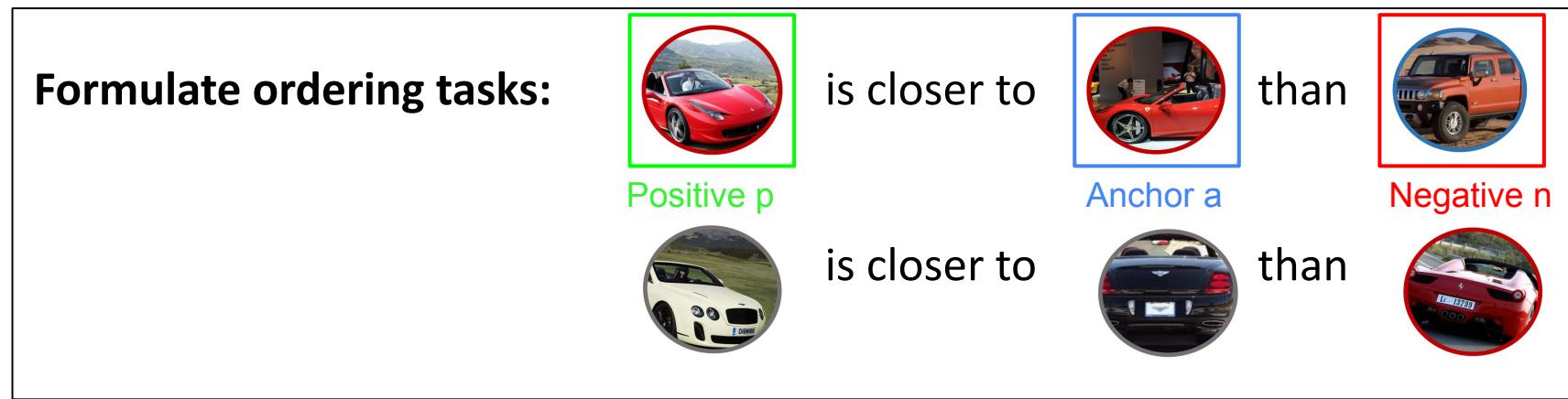
Ranking-based DML

- Learning informative embedding space requires to “arrange” data under $\phi(x)$
- Naturally formulate the learning problem as orderings “A is closer to B than C”



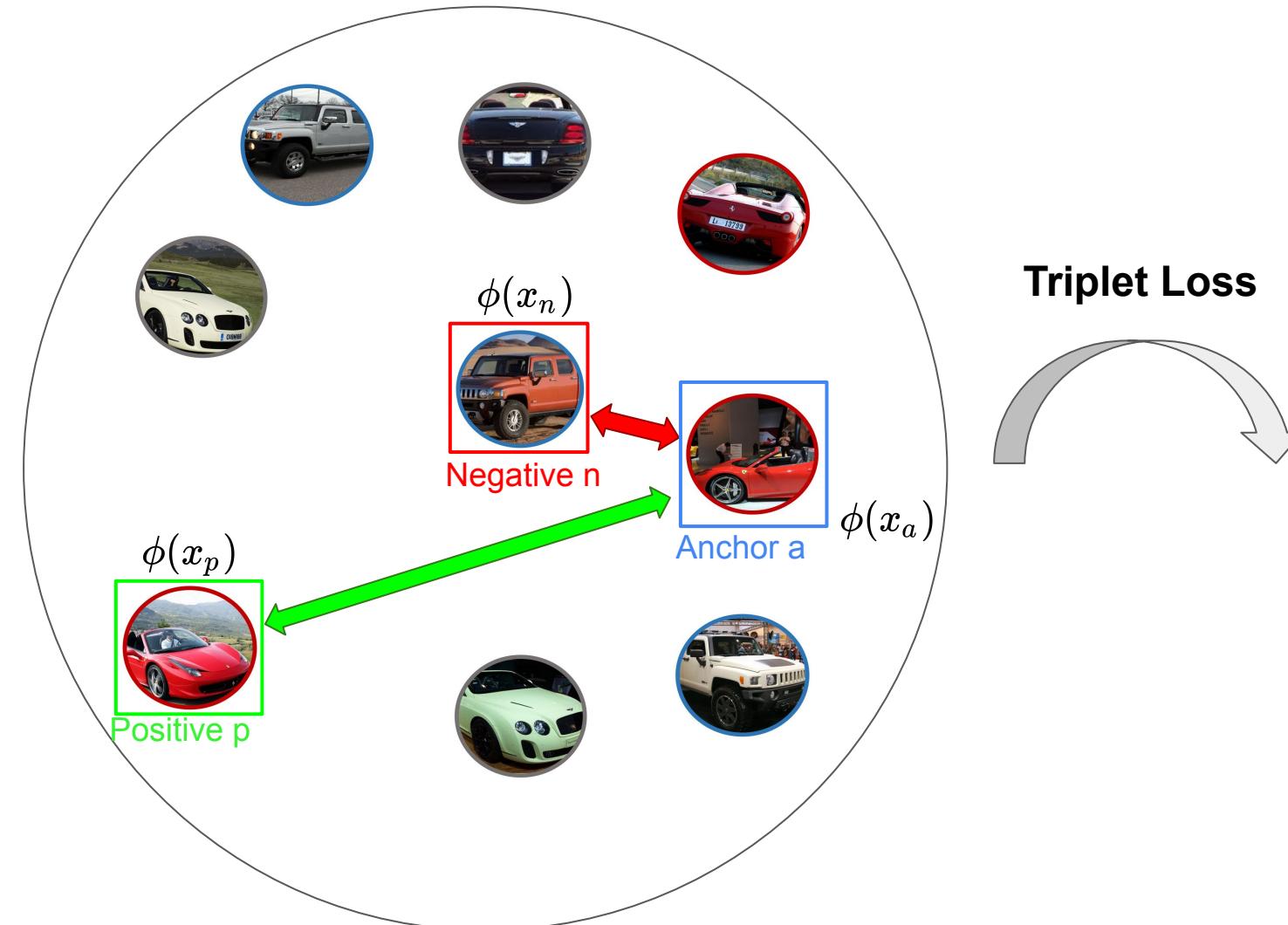
Ranking-based DML

- Learning informative embedding space requires to “arrange” data under $\phi(x)$
- Naturally formulate the learning problem as orderings “A is closer to B than C”



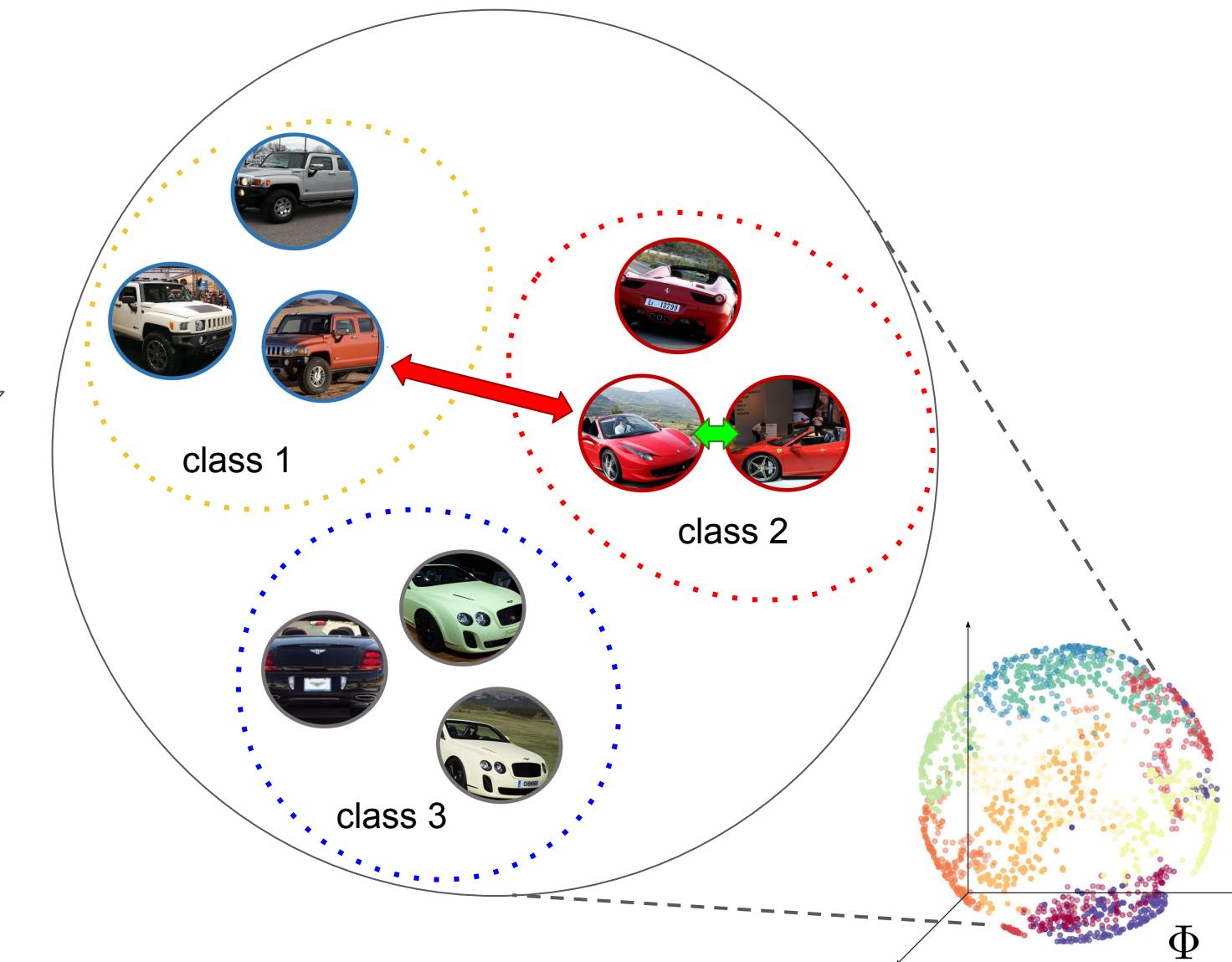
Ranking-based DML

- Learning informative embedding space requires to “arrange” data under $\phi(x)$
- Naturally formulate the learning problem as orderings “A is closer to B than C”



Triplet loss¹: $t = \{x_a, x_p, x_n\}$

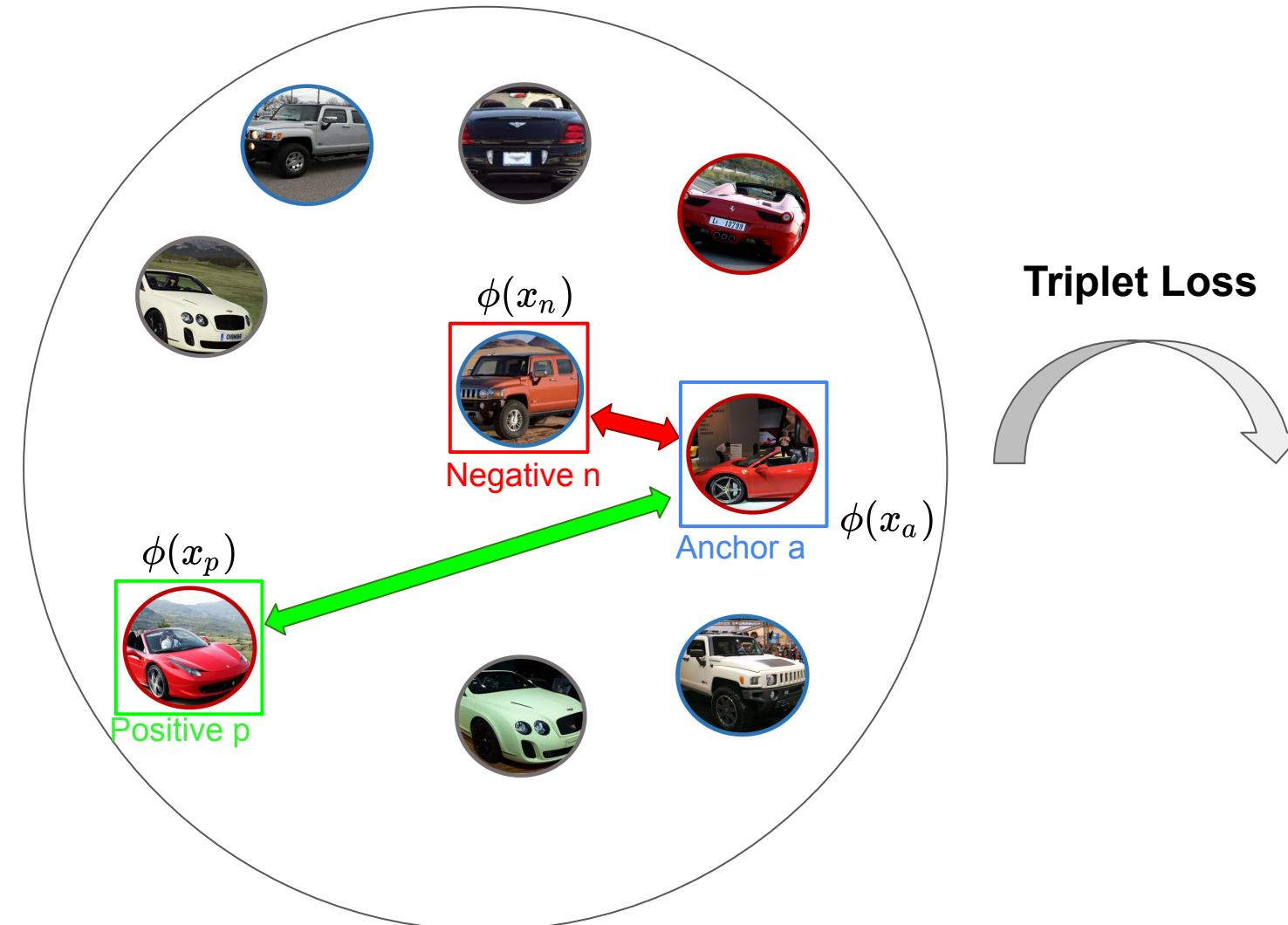
$$\ell^{\text{triplet}}(a, p, n) = \|\phi(x_a) - \phi(x_p)\|_2^2 - \|\phi(x_a) - \phi(x_n)\|_2^2$$



¹ Weinberger et al. 2006; Schroff et al. 2015

Ranking-based DML

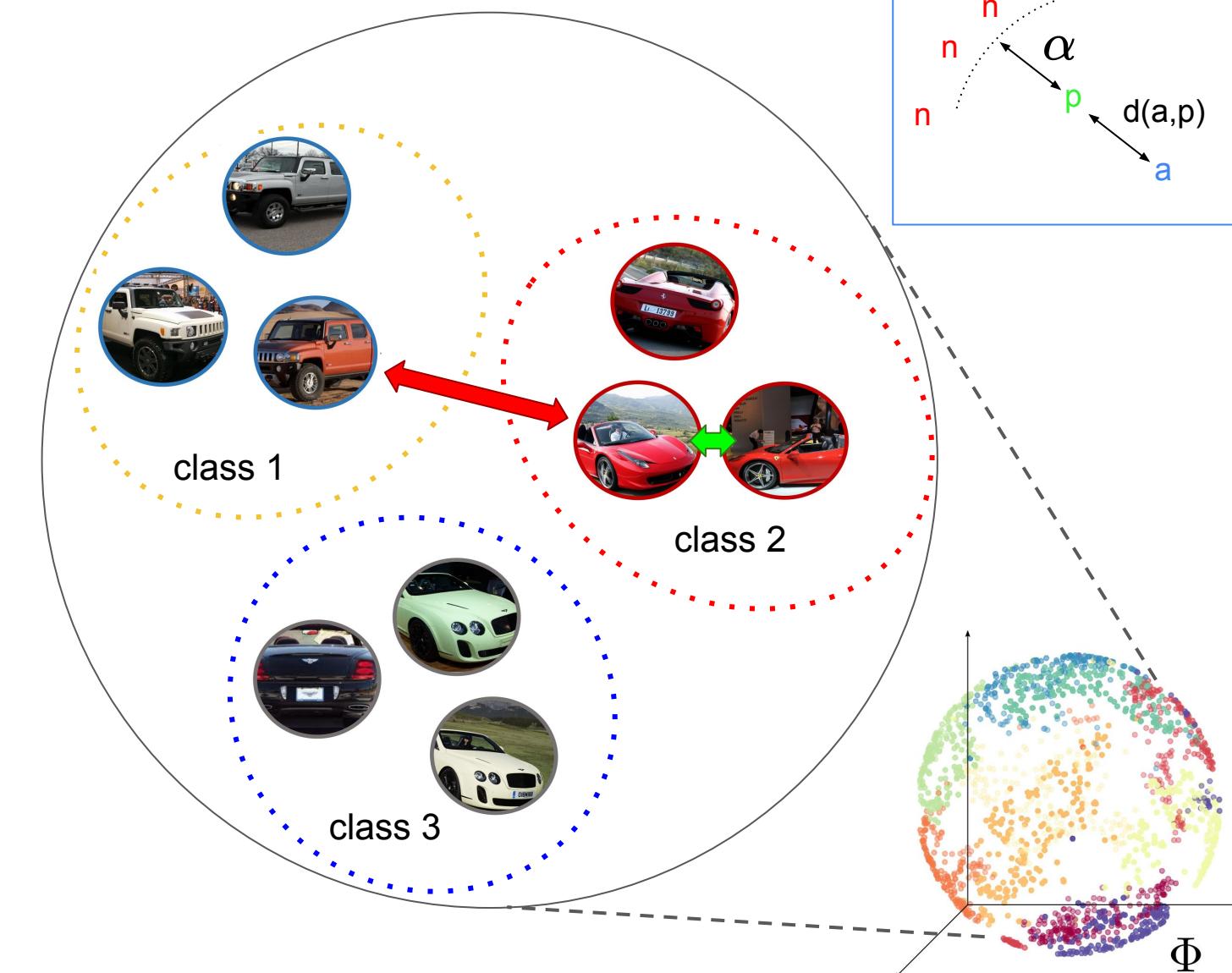
- Learning informative embedding space requires to “arrange” data under $\phi(x)$
- Naturally formulate the learning problem as orderings “A is closer to B than C”



Triplet loss¹: $t = \{x_a, x_p, x_n\}$

$$\ell^{\text{triplet}}(a, p, n) = \max(||\phi(x_a) - \phi(x_p)||_2^2 - ||\phi(x_a) - \phi(x_n)||_2^2 + \alpha, 0)$$

fixed margin



¹ Weinberger et al. 2006; Schroff et al. 2015

Ranking-based DML

- Learning informative embedding space requires to “arrange” data under $\phi(x)$
- Naturally formulate the learning problem as orderings “A is closer to B than C”

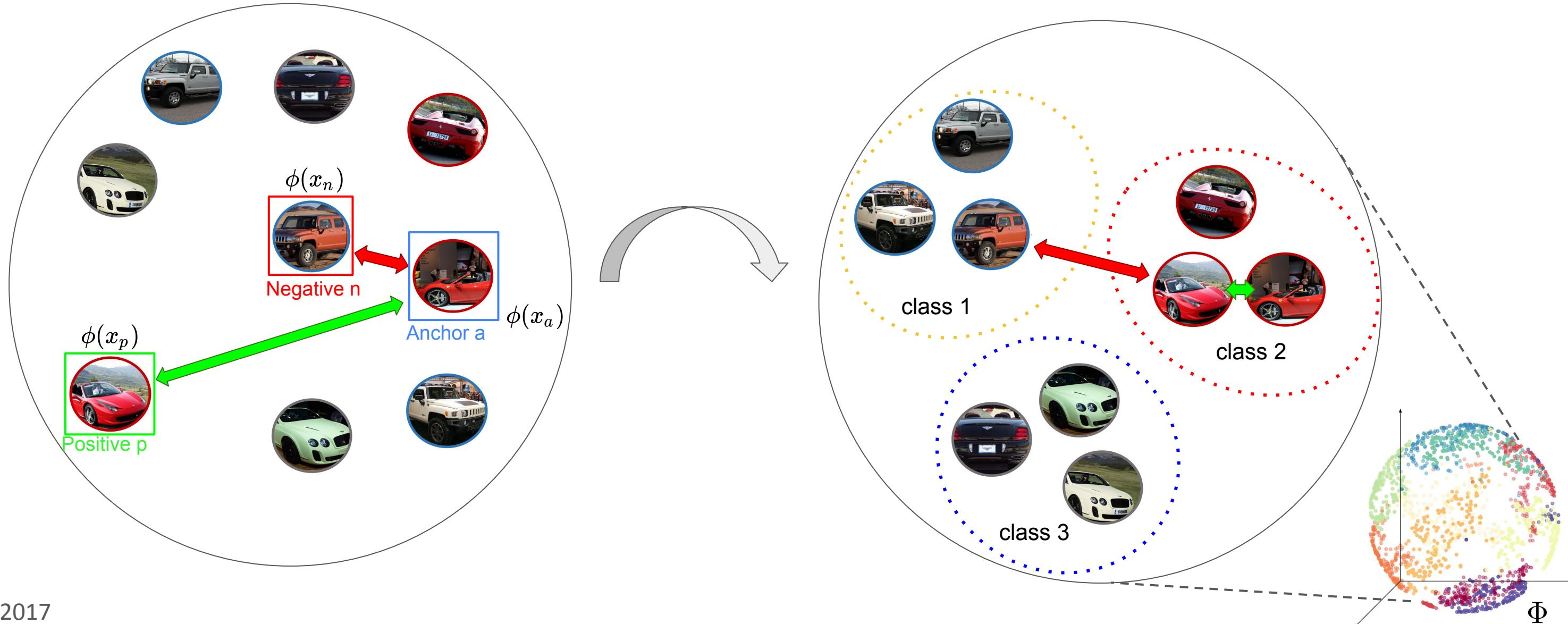
Margin loss²:

$$\ell^{\text{triplet}}(a, p, n) := [D_{ap}^2 - D_{an}^2 + \alpha]_+$$

$$\ell^{\text{margin}}(i, j) := (\alpha + y_{ij}(D_{ij} - \beta))_+$$

↑
fixed margin ↑
learned margin parameter

$$\beta(i) := \beta^{(0)} + \beta_{c(i)}^{(\text{class})} + \beta_i^{(\text{img})}$$



Ranking-based DML

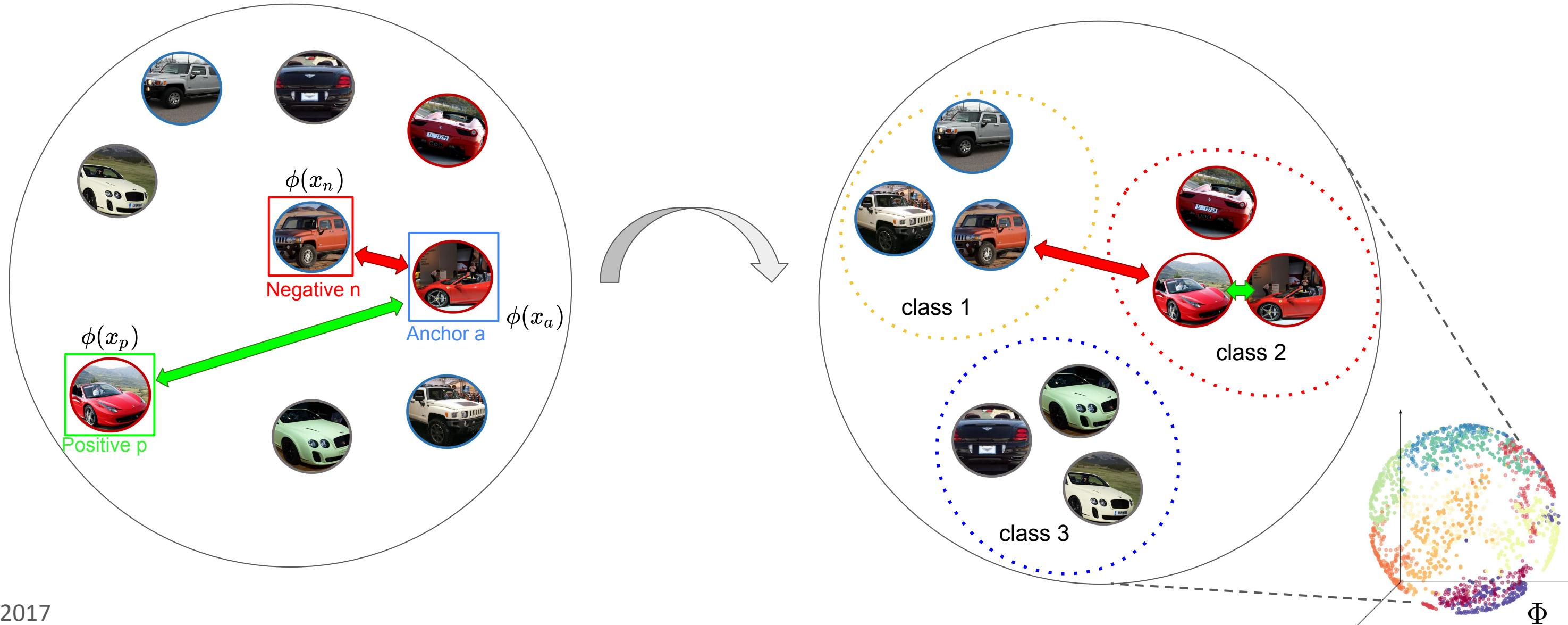
- Learning informative embedding space requires to “arrange” data under $\phi(x)$
- Naturally formulate the learning problem as orderings “A is closer to B than C”

Margin loss² (Joint optimization):

$$\text{minimize} \sum_{(i,j)} \ell^{\text{margin}}(i,j) + \nu \left(\beta^{(0)} + \beta_{c(i)}^{\text{(class)}} + \beta_i^{\text{(img)}} \right)$$

$$\ell^{\text{margin}}(i,j) := (\alpha + y_{ij}(D_{ij} - \beta))_+$$

$$\beta(i) := \beta^{(0)} + \beta_{c(i)}^{\text{(class)}} + \beta_i^{\text{(img)}}$$



Ranking-based DML

- Learning informative embedding space requires to “arrange” data under $\phi(x)$
- Naturally formulate the learning problem as orderings “A is closer to B than C”

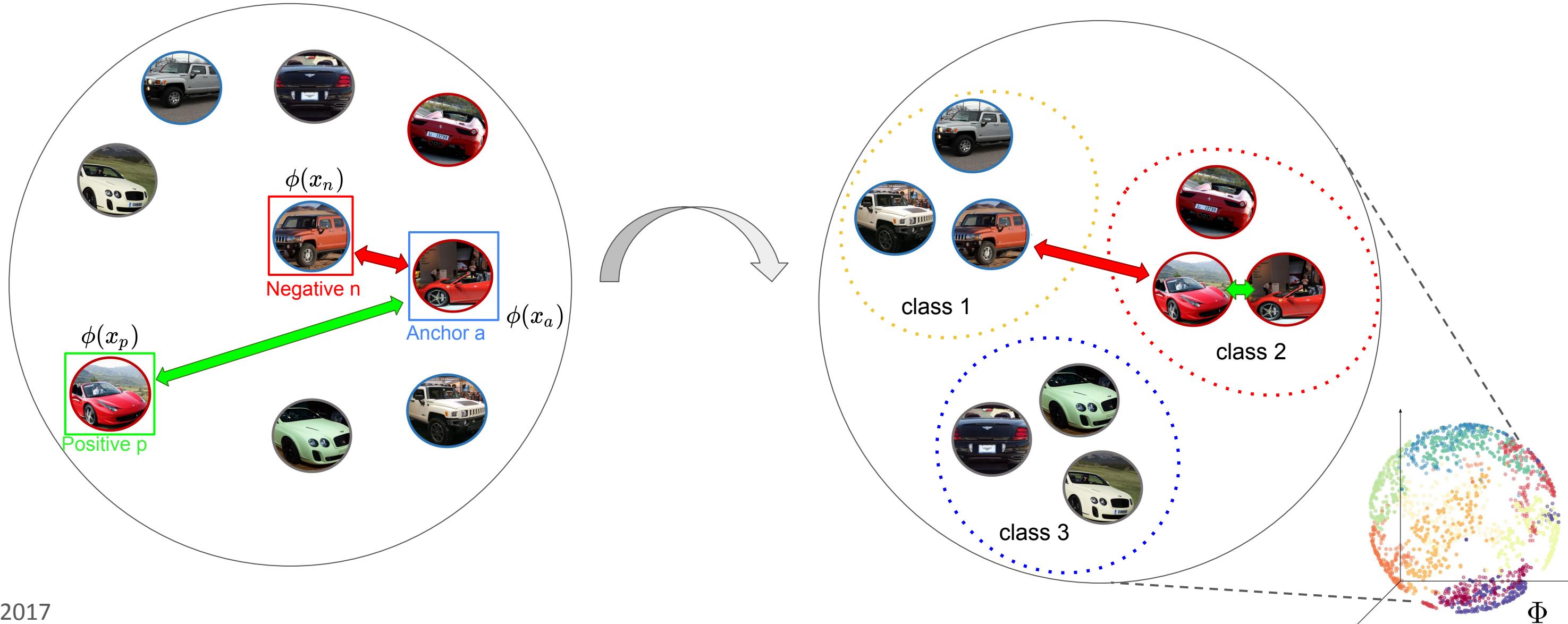
Margin loss² (Joint optimization):

$$\text{minimize} \sum_{(i,j)} \ell^{\text{margin}}(i,j) + \nu \left(\beta^{(0)} + \beta_{c(i)}^{\text{(class)}} + \beta_i^{\text{(img)}} \right)$$

$$\ell^{\text{margin}}(i,j) := (\alpha + y_{ij}(D_{ij} - \beta))_+$$

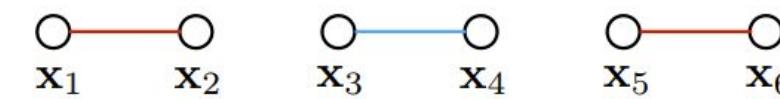
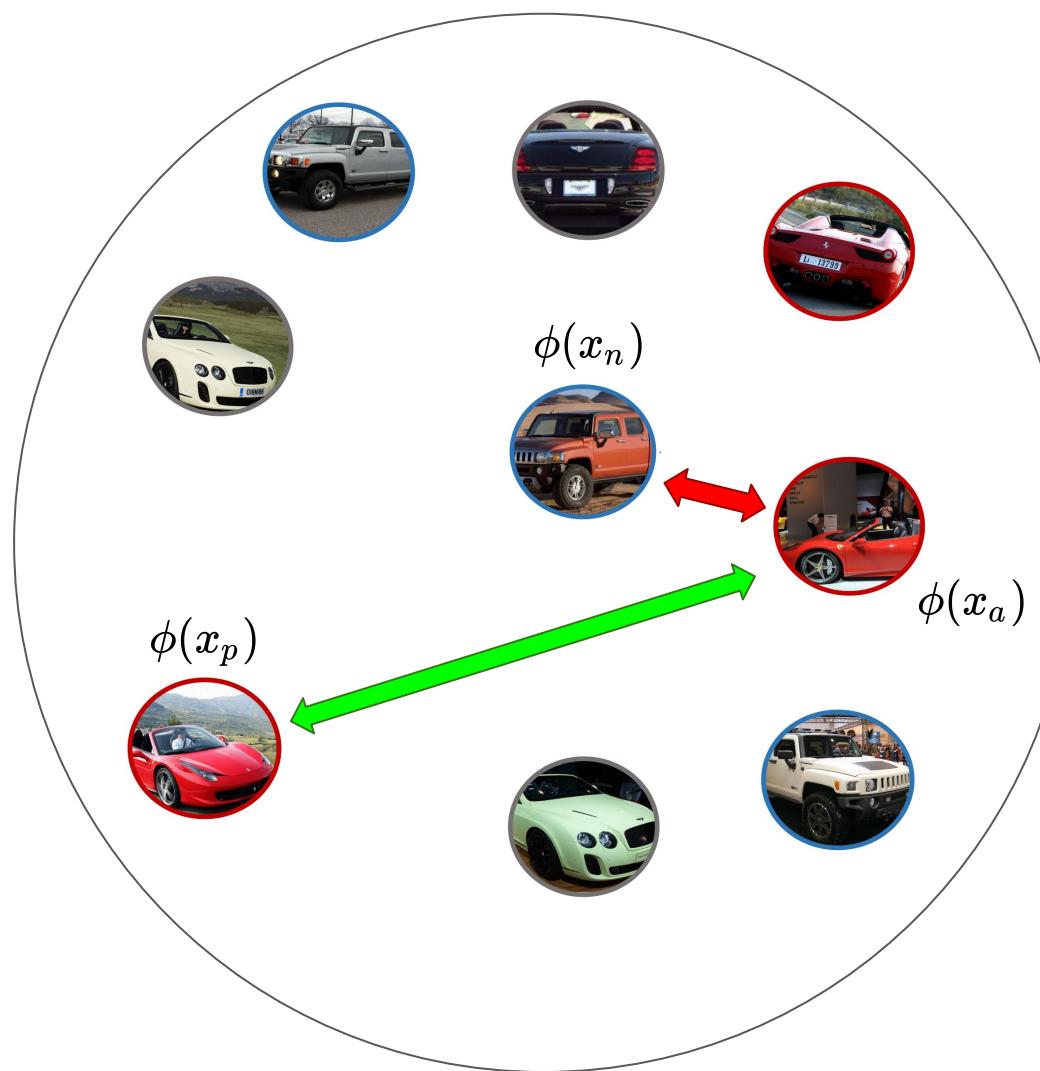
$$\beta(i) := \beta^{(0)} + \beta_{c(i)}^{\text{(class)}} + \cancel{\beta_i^{\text{(img)}}}$$

in practice

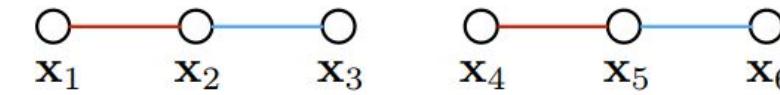


Ranking-based DML

- So far: Per sample in batch: **form a single triplet**



(a) Contrastive embedding

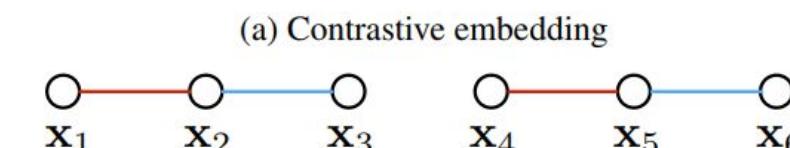
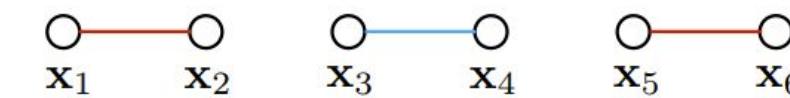
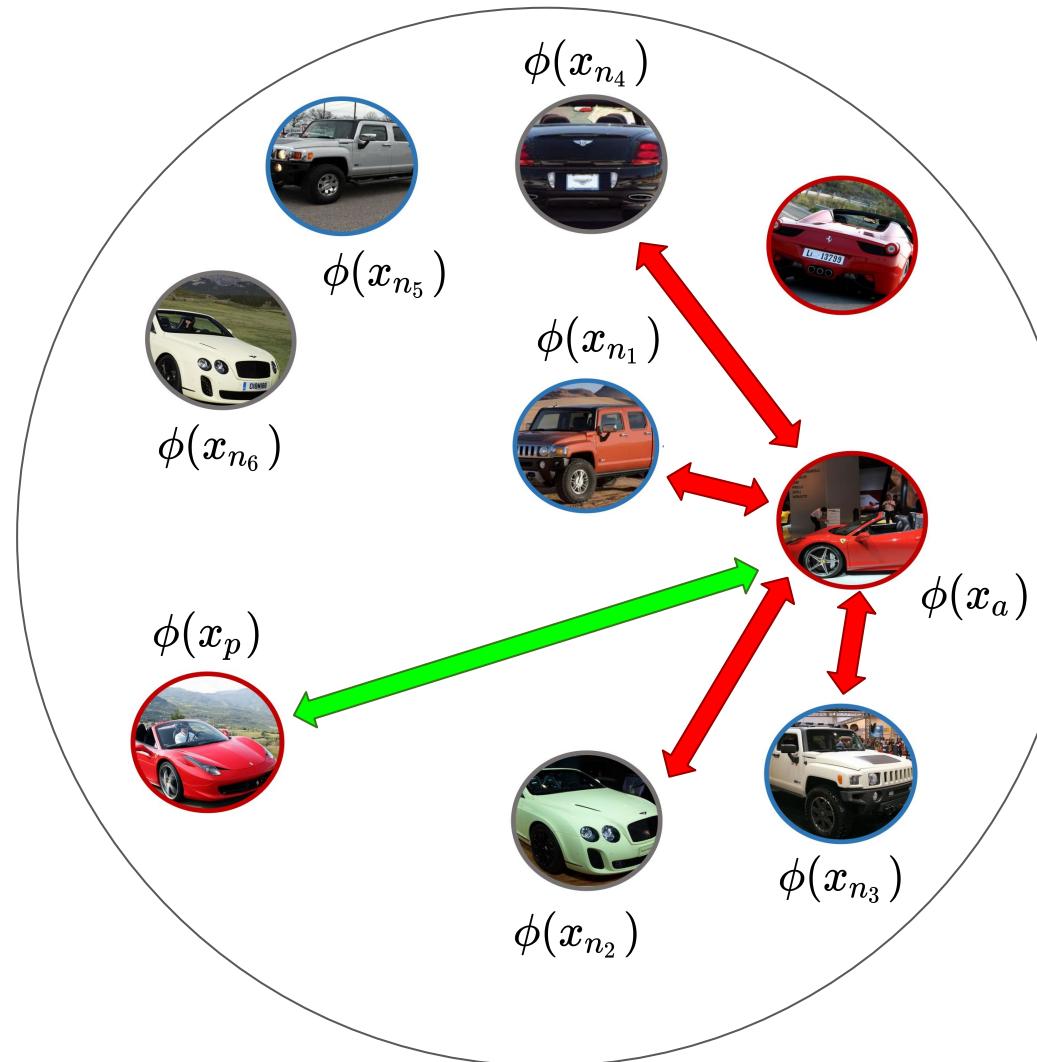


(b) Triplet embedding

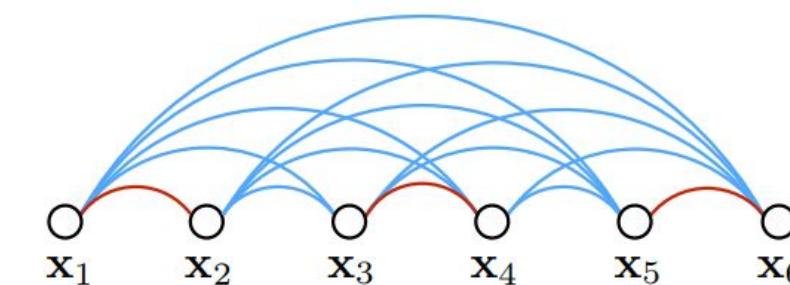
$\mathcal{O}(B)$ relations

Ranking-based DML

- So far: Per sample in batch: **form a single triplet**
- Use **all possible relations** within batch for learning!



(b) Triplet embedding



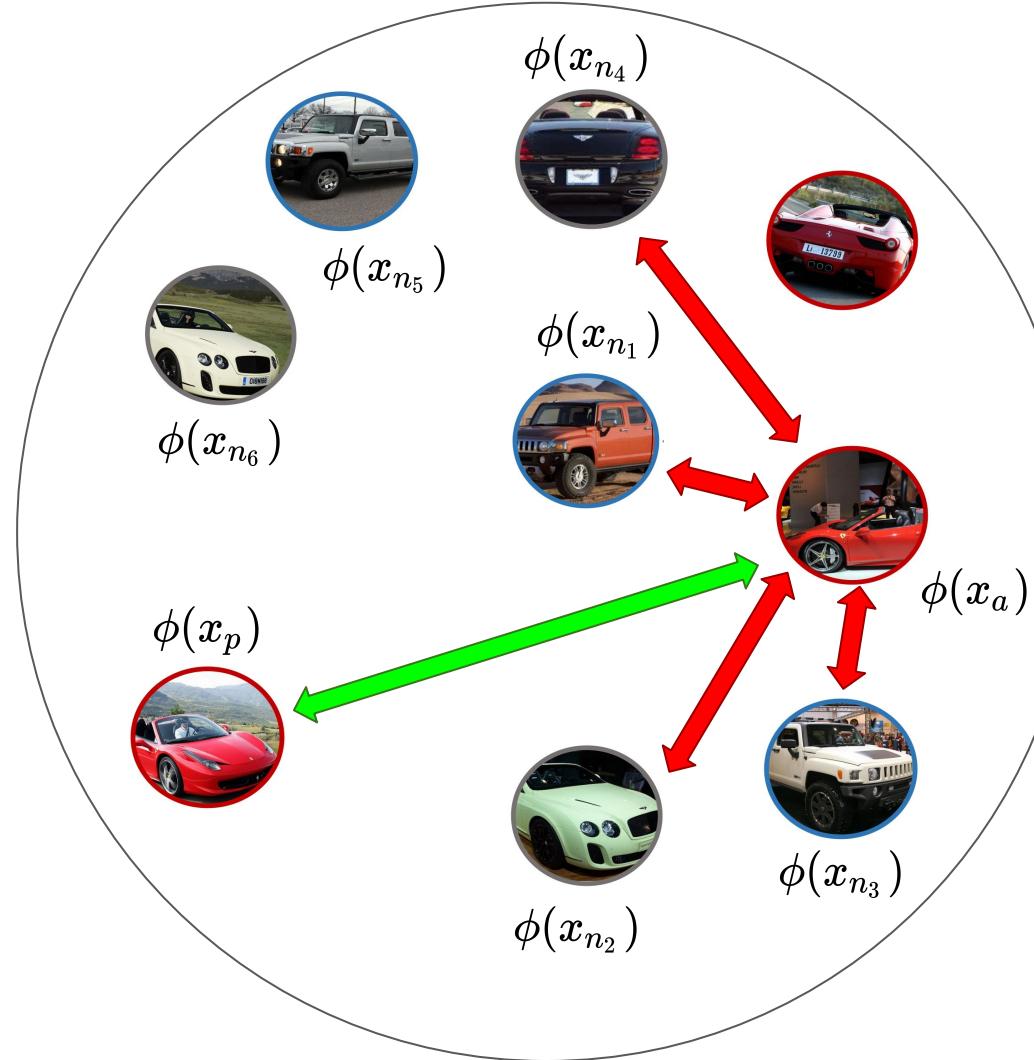
(c) Lifted structured embedding

$\mathcal{O}(B)$ relations

$\mathcal{O}(B^2)$ relations

Ranking-based DML

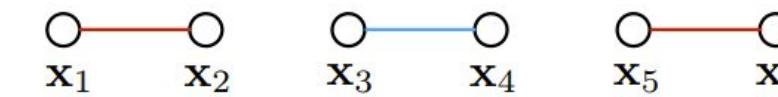
- So far: Per sample in batch: **form a single triplet**
- Use **all possible relations** within batch for learning!



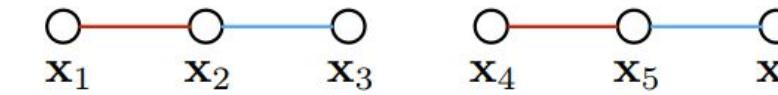
Lifted Structures Feature Embedding Loss³ :

$$\tilde{J}_{i,j} = \log \left(\sum_{(i,k) \in \mathcal{N}} \exp\{\alpha - D_{i,k}\} + \sum_{(j,l) \in \mathcal{N}} \exp\{\alpha - D_{j,l}\} \right) + D_{i,j}$$

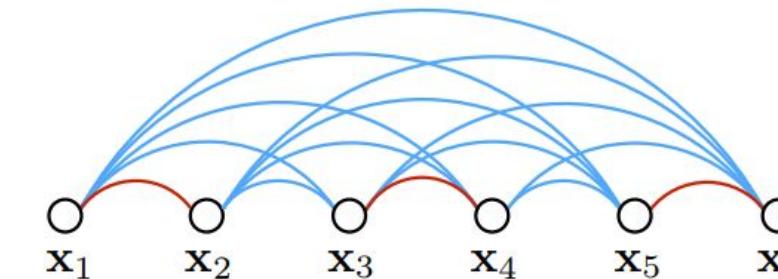
$$\tilde{J} = \frac{1}{2|\mathcal{P}|} \sum_{(i,j) \in \mathcal{P}} \max(0, \tilde{J}_{i,j})^2,$$



(a) Contrastive embedding



(b) Triplet embedding

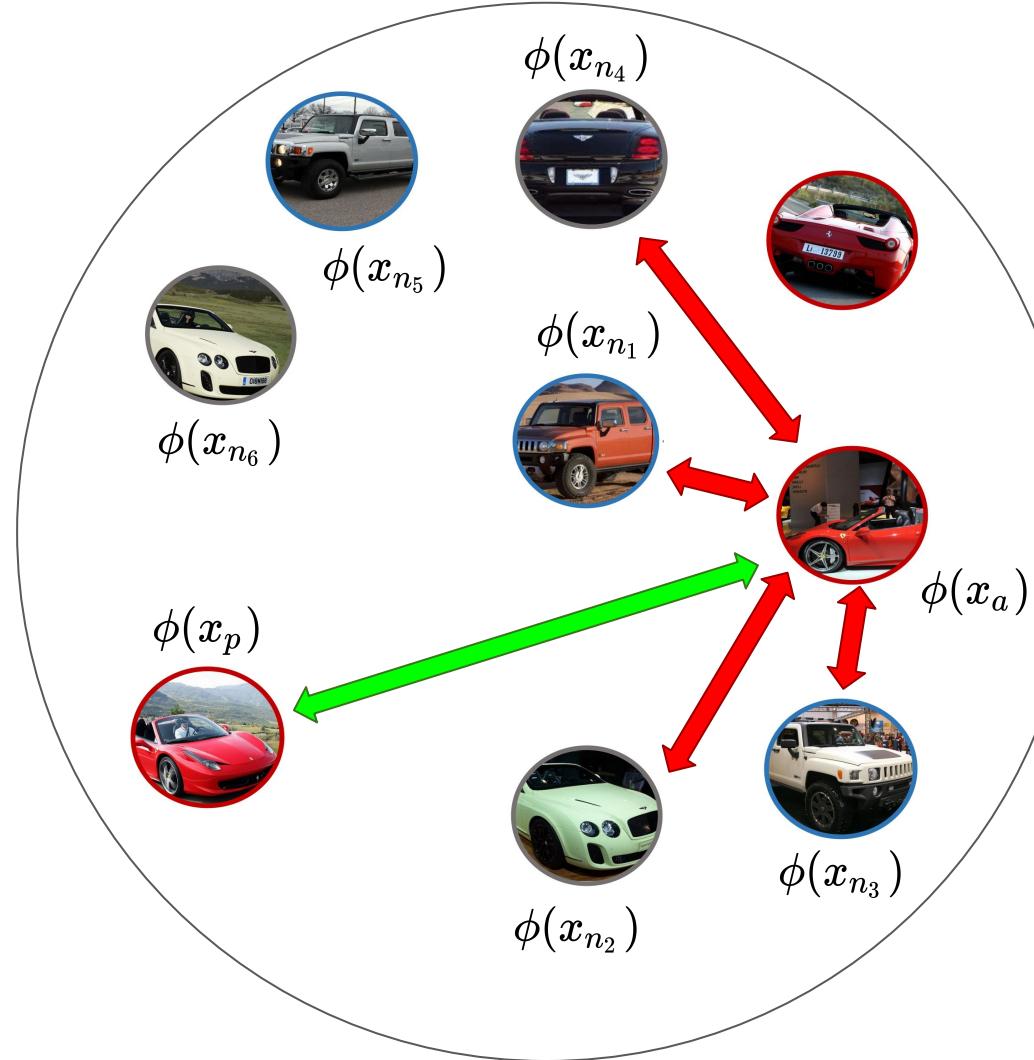


$\mathcal{O}(B)$ relations

$\mathcal{O}(B^2)$ relations

Ranking-based DML

- So far: Per sample in batch: **form a single triplet**
- Use **all possible relations** within batch for learning!



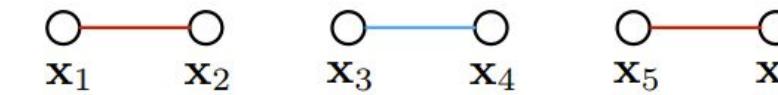
Lifted Structures Feature Embedding Loss³ :

$$\tilde{J}_{i,j} = \log \left(\sum_{(i,k) \in \mathcal{N}} \exp\{\alpha - D_{i,k}\} + \sum_{(j,l) \in \mathcal{N}} \exp\{\alpha - D_{j,l}\} \right) + D_{i,j}$$

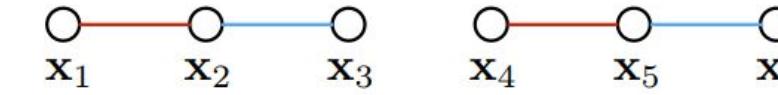
$$\tilde{J} = \frac{1}{2|\mathcal{P}|} \sum_{(i,j) \in \mathcal{P}} \max(0, \tilde{J}_{i,j})^2,$$

anchor to all negatives in batch positive to all negatives in batch

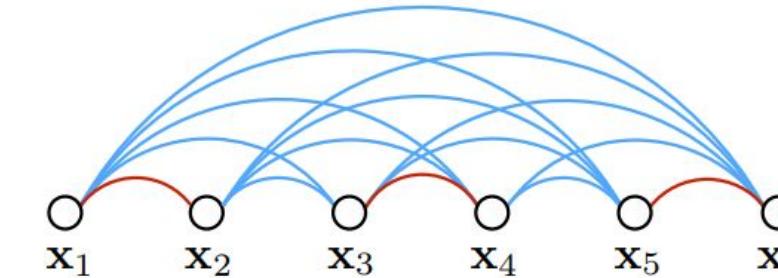
anchor to positive



(a) Contrastive embedding



(b) Triplet embedding



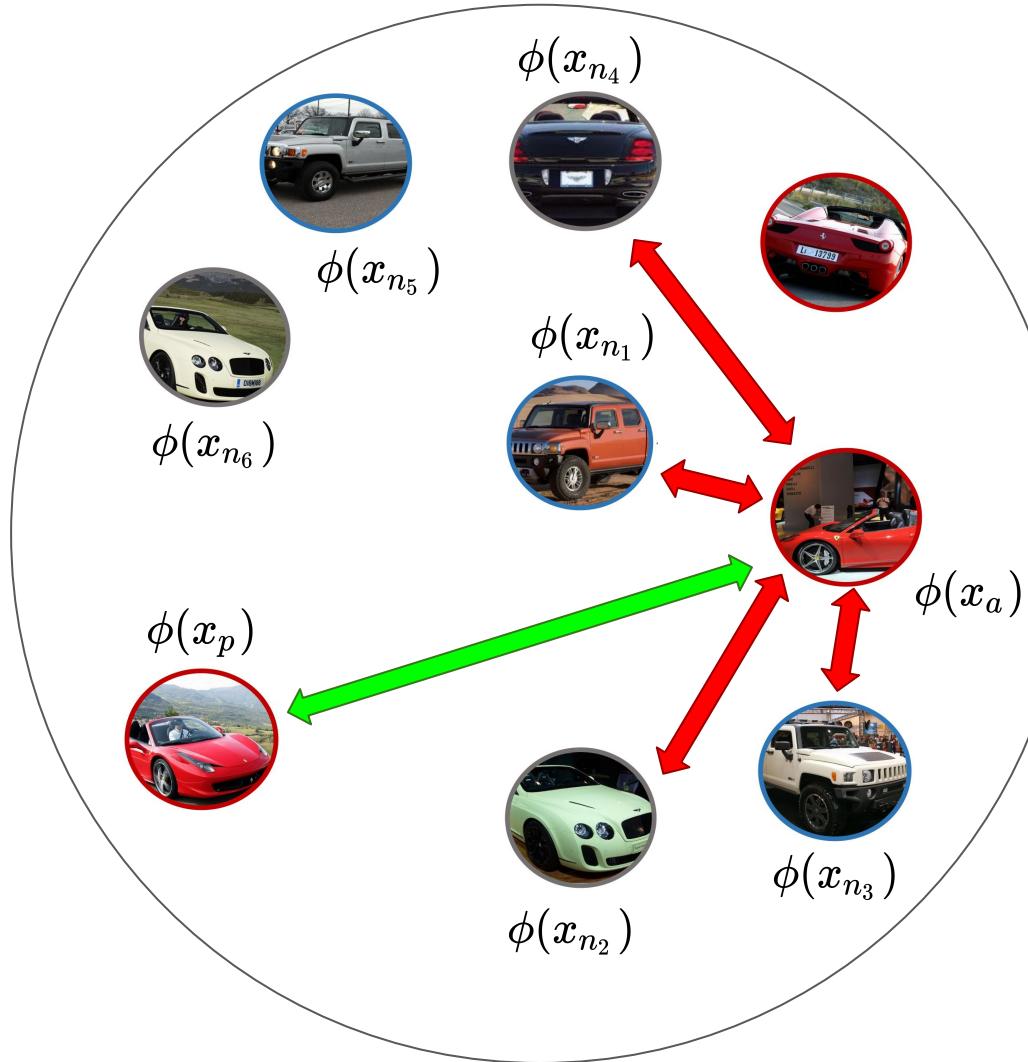
(c) Lifted structured embedding

$\mathcal{O}(B)$ relations

$\mathcal{O}(B^2)$ relations

Ranking-based DML

- So far: Per sample in batch: **form a single triplet**
- Use **all possible relations** within batch for learning!

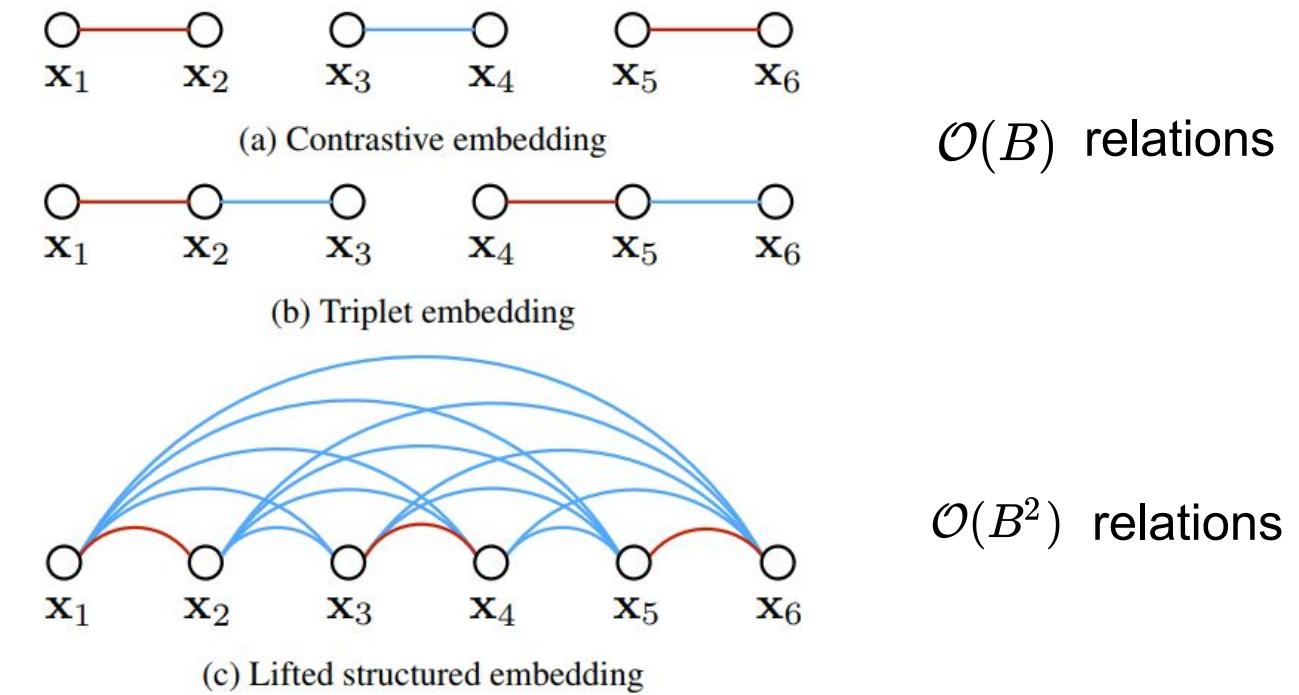


Multi-Similarity Loss⁴:

$$\mathcal{L}_{MS} = \frac{1}{m} \sum_{i=1}^m \left\{ \frac{1}{\alpha} \log \left[1 + \sum_{k \in \mathcal{P}_i} e^{-\alpha(S_{ik} - \lambda)} \right] + \frac{1}{\beta} \log \left[1 + \sum_{k \in \mathcal{N}_i} e^{\beta(S_{ik} - \lambda)} \right] \right\}$$

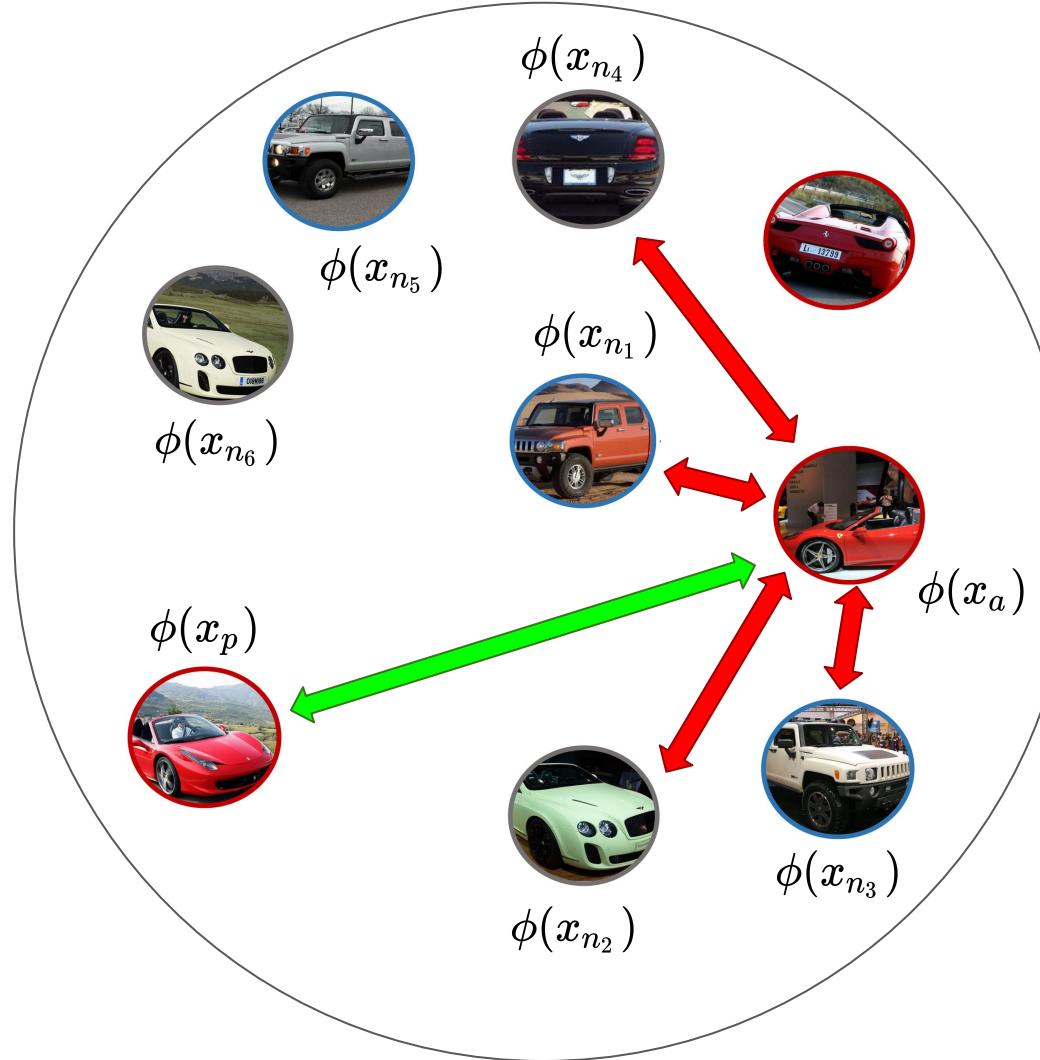
anchor to positives in batch

anchor to negatives in batch



Ranking-based DML

- So far: Per sample in batch: **form a single triplet**
- Use **all possible relations** within batch for learning!



Multi-Similarity Loss⁴:

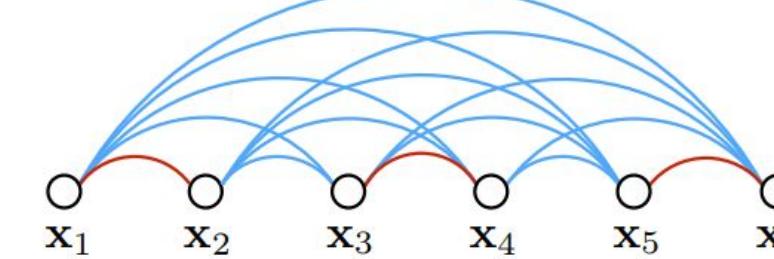
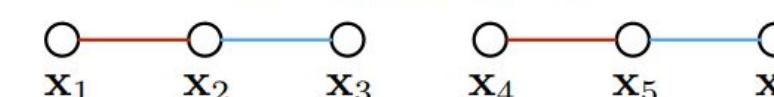
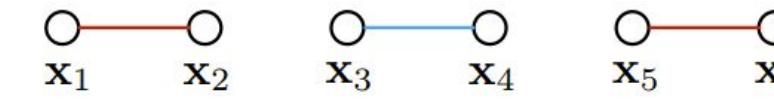
$$\mathcal{L}_{MS} = \frac{1}{m} \sum_{i=1}^m \left\{ \frac{1}{\alpha} \log \left[1 + \boxed{\sum_{k \in \mathcal{P}_i} e^{-\alpha(S_{ik} - \lambda)}} \right] + \frac{1}{\beta} \log \left[1 + \boxed{\sum_{k \in \mathcal{N}_i} e^{\beta(S_{ik} - \lambda)}} \right] \right\}.$$

anchor to positives in batch
anchor to negatives in batch

Informative pair mining:

$$S_{ij}^+ < \max_{y_k \neq y_i} S_{ik} + \epsilon, \quad \text{"positive less similar to anchor than best negative"}$$

$$S_{ij}^- > \min_{y_k = y_i} S_{ik} - \epsilon, \quad \text{"negative more similar to anchor than worst positive"}$$

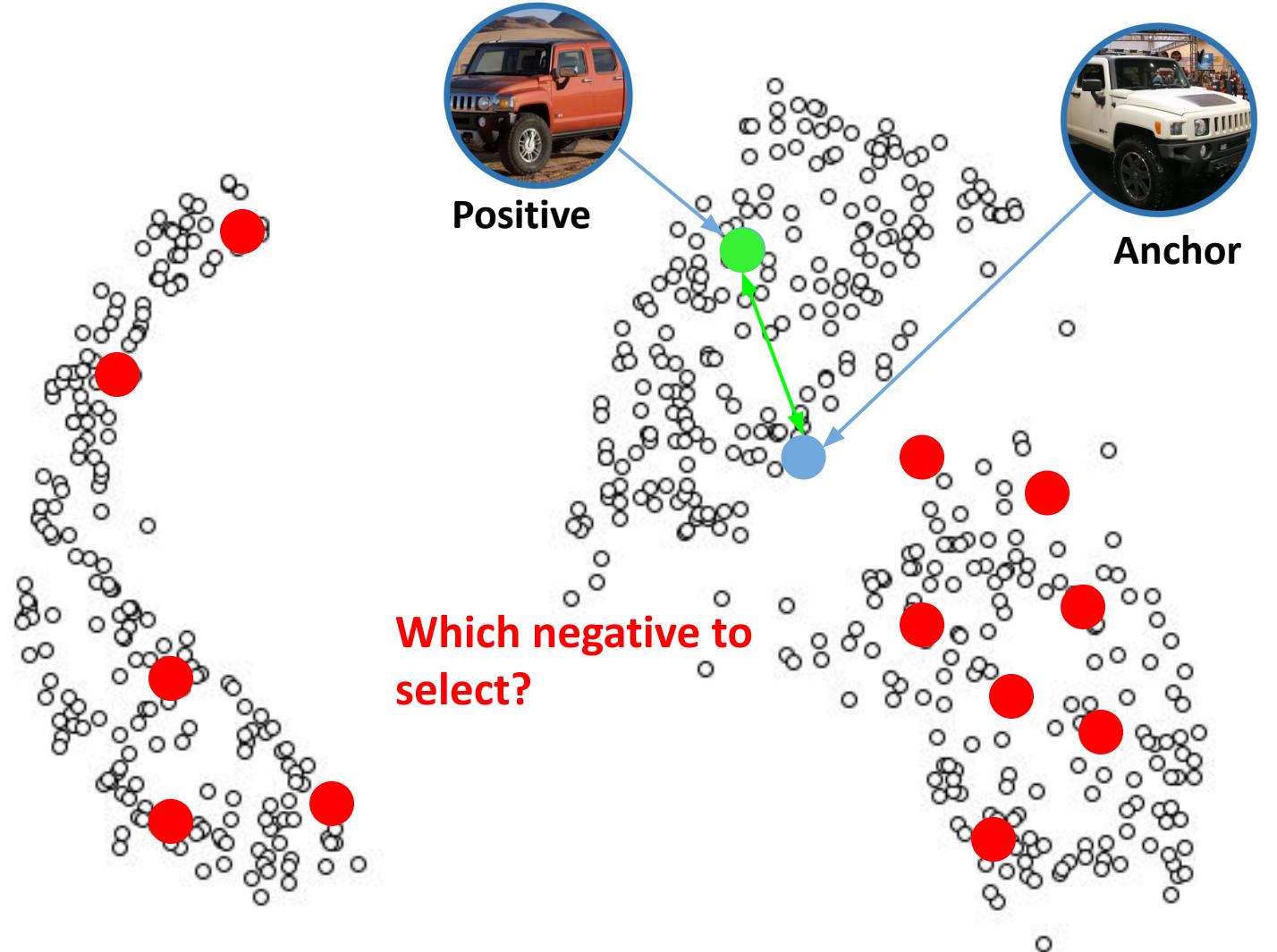


$\mathcal{O}(B)$ relations

$\mathcal{O}(B^2)$ relations

Negative Sampling in Ranking-based DML

- **Infeasible** amount of triplets to learn from (Scale $\mathcal{O}(n^3)$)

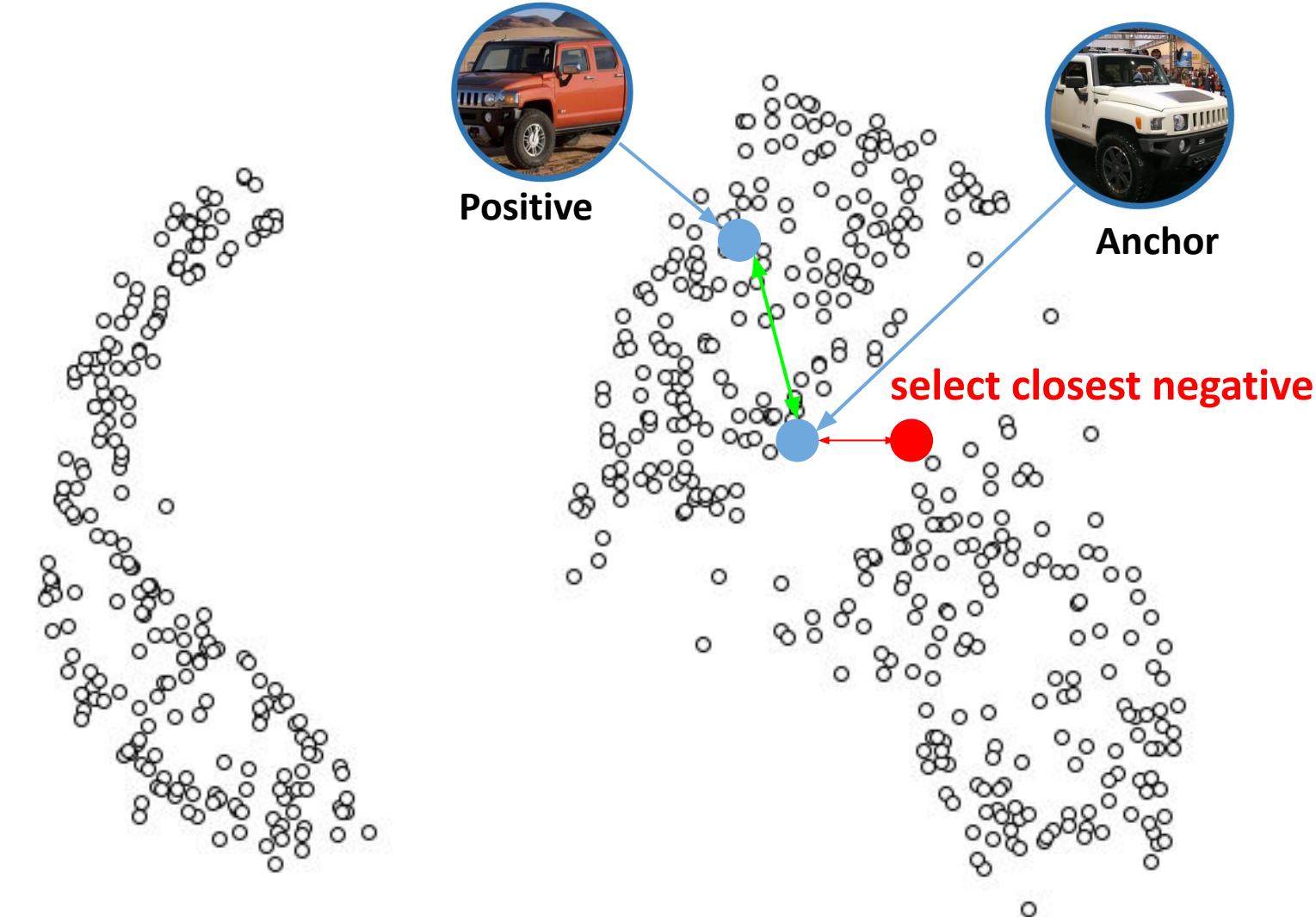


Negative Sampling in Ranking-based DML

- **Infeasible** amount of triplets to learn from (Scale $\mathcal{O}(n^3)$)
- Fixed sampling heuristics:

- **(Semi-)Hard negative mining** [Schroff et al. 2015]

$$n^* = \arg \min_{n: d_{an} > d_{ap}} d_{an}$$



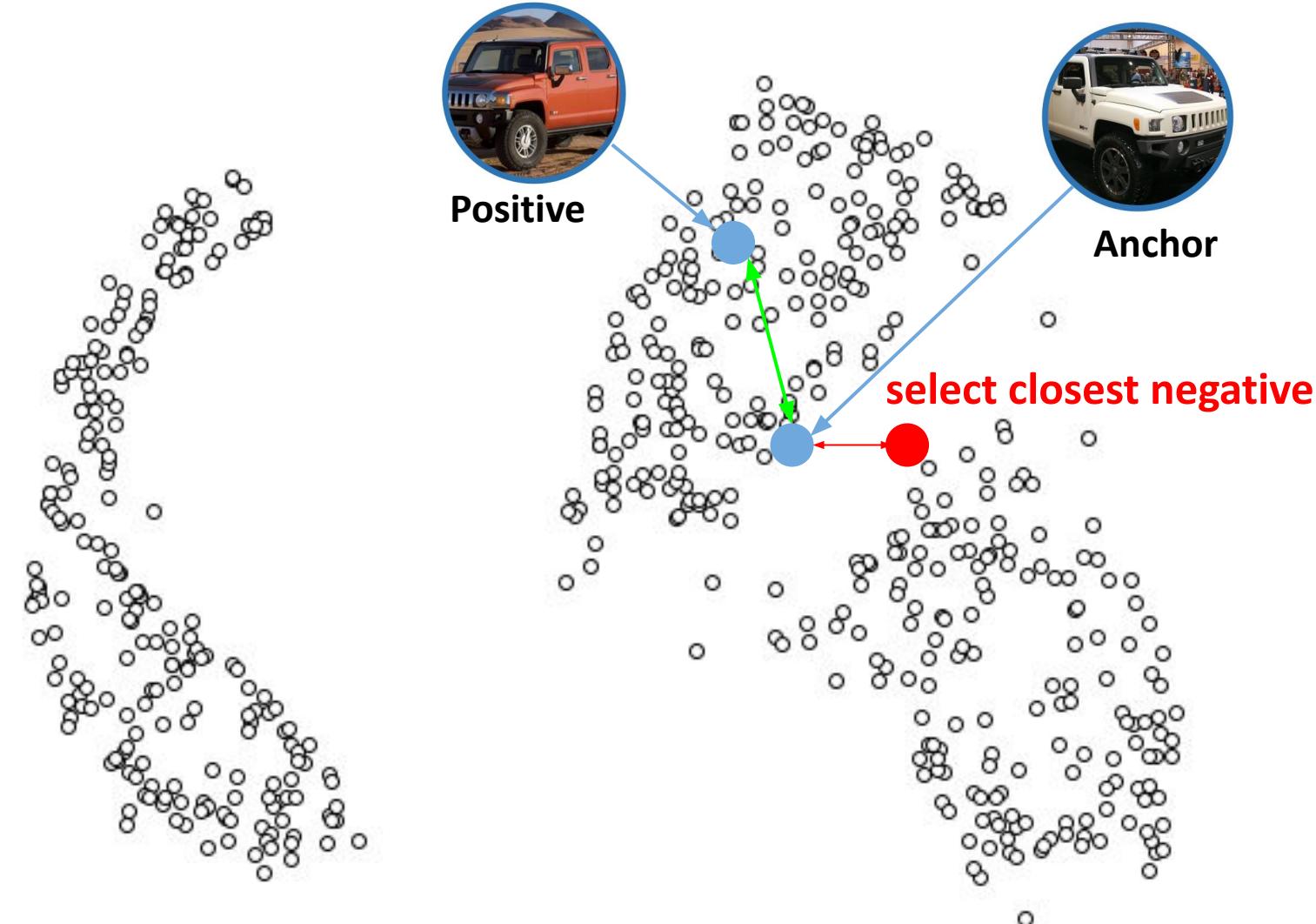
Negative Sampling in Ranking-based DML

- **Infeasible** amount of triplets to learn from (Scale $\mathcal{O}(n^3)$)
- Fixed sampling heuristics:

- **(Semi-)Hard negative mining** [Schroff et al. 2015]

$$n^* = \arg \min_{n: d_{an} > d_{ap}} d_{an}$$

- May suffer from **unstable gradients** due to focus on hardest negatives only

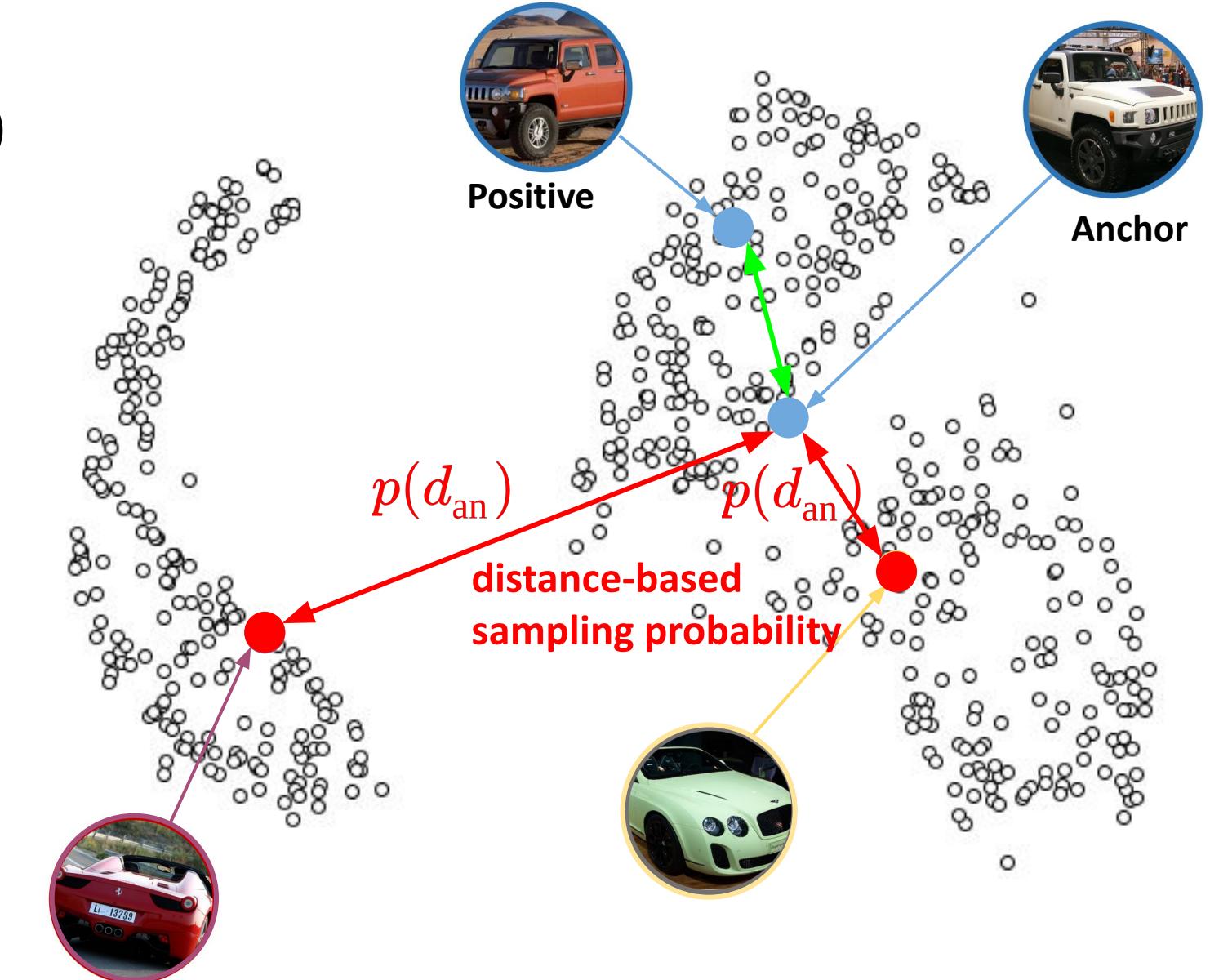


Negative Sampling in Ranking-based DML

- **Infeasible** amount of triplets to learn from (Scale $\mathcal{O}(n^3)$)
- Fixed sampling heuristics:

- **Distance-weighted sampling** [Wu et al. 2017]

$$n^* \sim p(I_n | I_a) = p(d_{an})$$



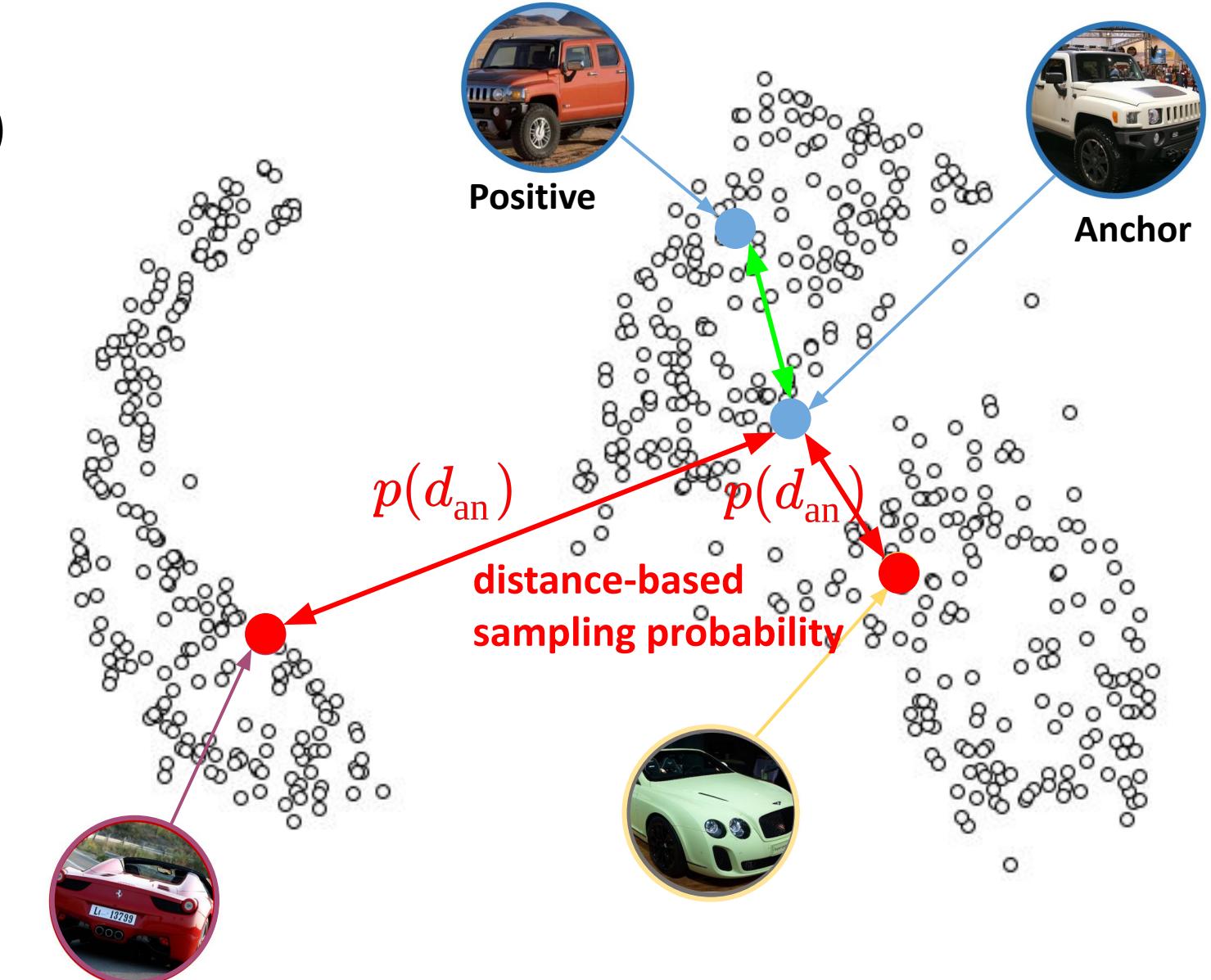
Negative Sampling in Ranking-based DML

- **Infeasible** amount of triplets to learn from (Scale $\mathcal{O}(n^3)$)
- Fixed sampling heuristics:

- **Distance-weighted sampling** [Wu et al. 2017]

$$n^* \sim p(I_n | I_a) = p(d_{an})$$

- Soften hard negative constraint by uniform sampling from **entire range of distances** d_{an}



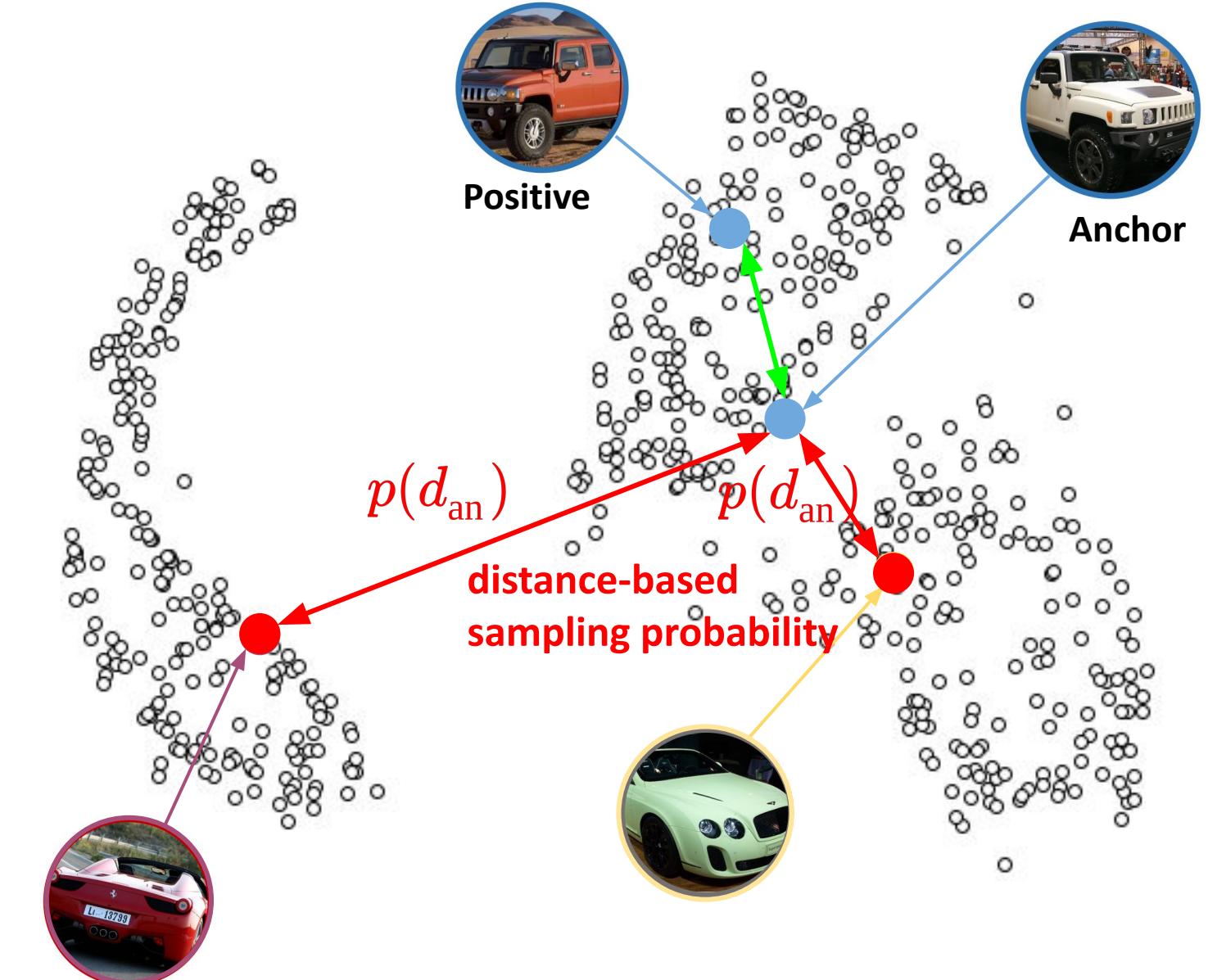
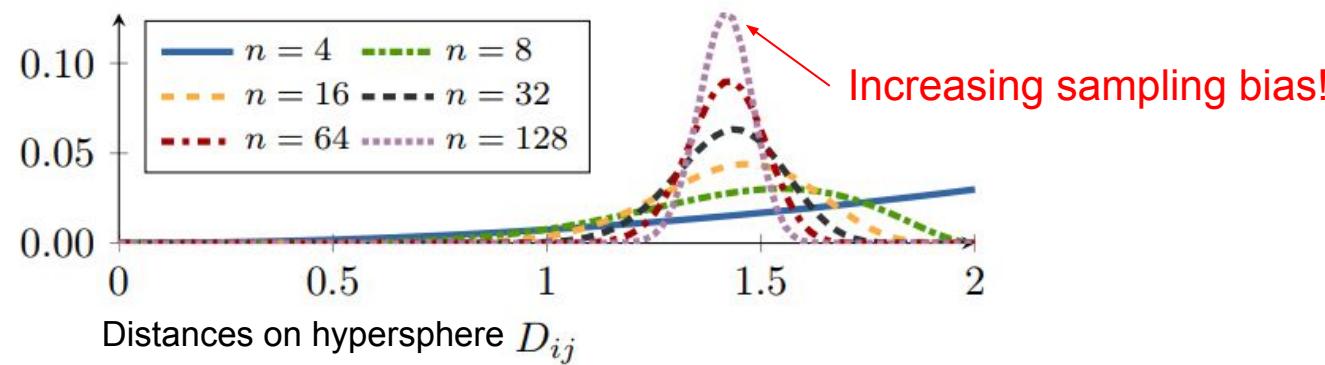
Negative Sampling in Ranking-based DML

- **Infeasible** amount of triplets to learn from (Scale $\mathcal{O}(n^3)$)
- Fixed sampling heuristics:

- **Distance-weighted sampling** [Wu et al. 2017]

$$n^* \sim p(I_n | I_a) = p(d_{an})$$

- Soften hard negative constraint by uniform sampling from **entire range of distances** d_{an}



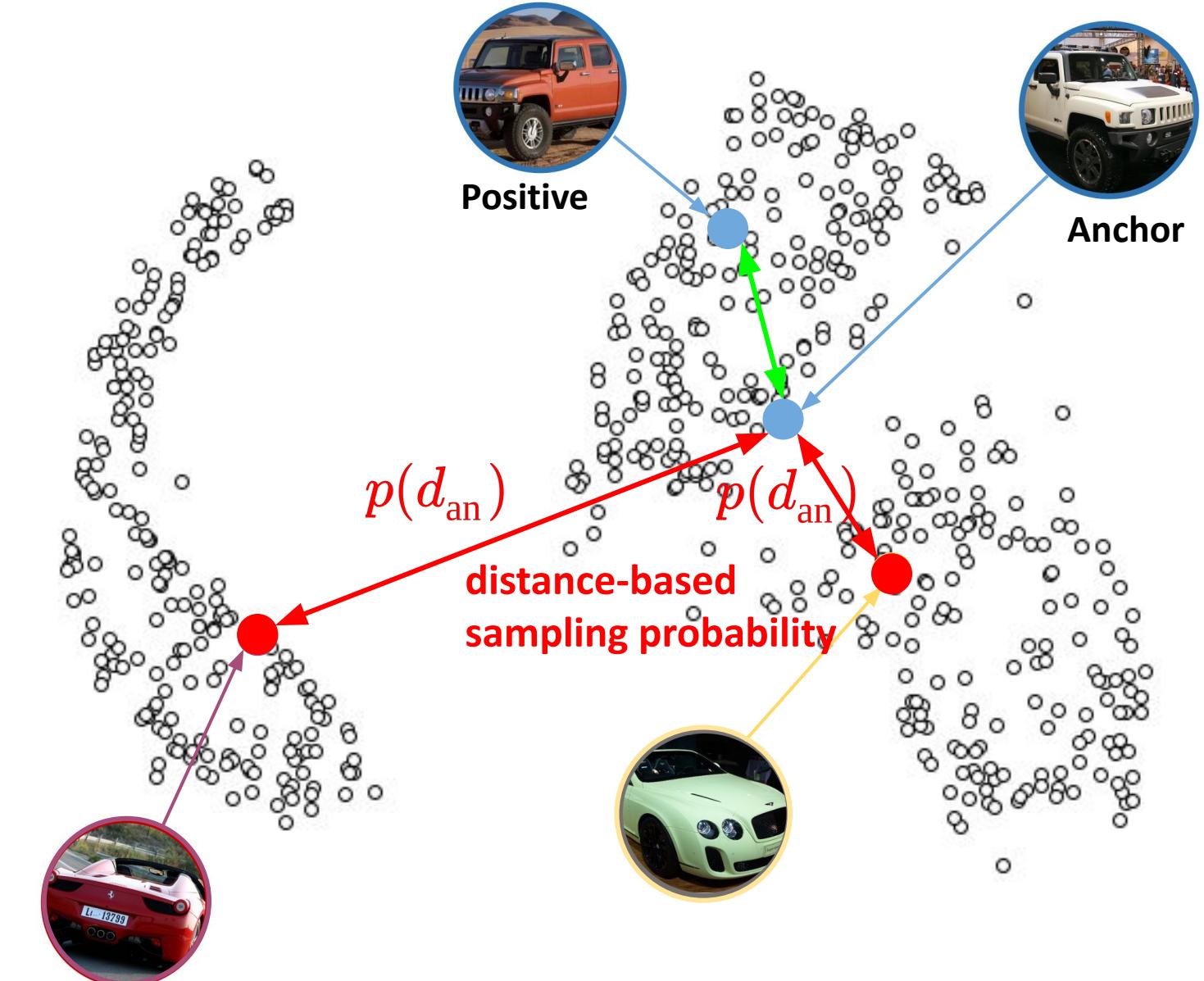
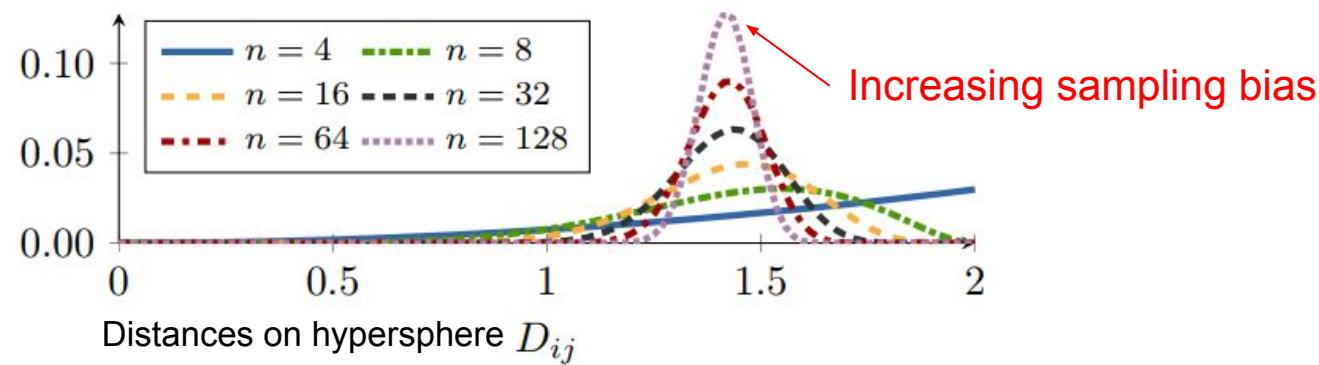
Negative Sampling in Ranking-based DML

- **Infeasible** amount of triplets to learn from (Scale $\mathcal{O}(n^3)$)
- Fixed sampling heuristics:

- **Distance-weighted sampling** [Wu et al. 2017]

$$n^* \sim p(I_n | I_a) \propto \min(\lambda, q(d_{an})^{-1})$$

- Soften hard negative constraint by uniform sampling from **entire range of distances** d_{an}
- Uniform distribution on \mathbb{S}^D : $q(d) \propto d^{k-2} [1 - \frac{1}{4}d^2]^{\frac{k-3}{2}}$



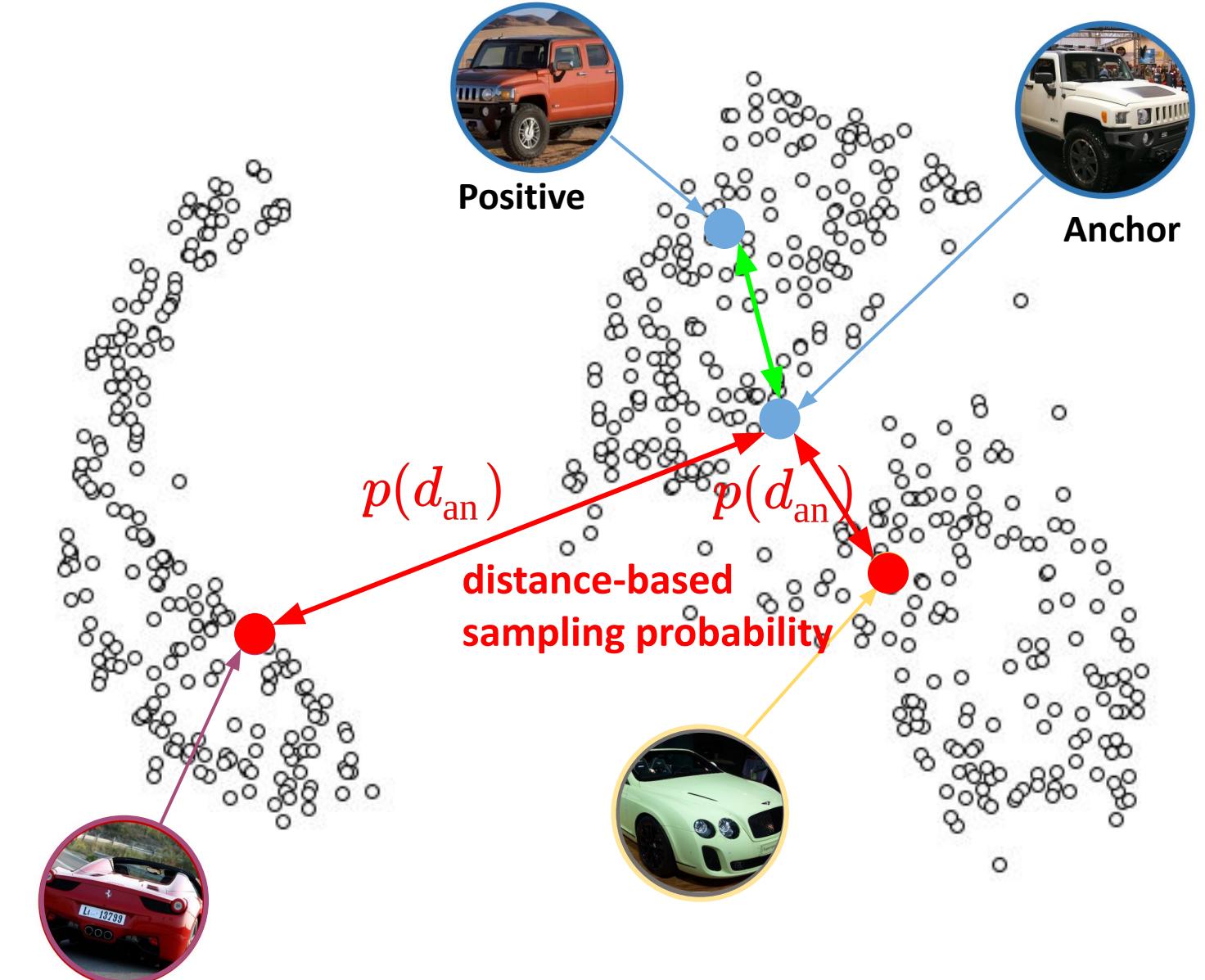
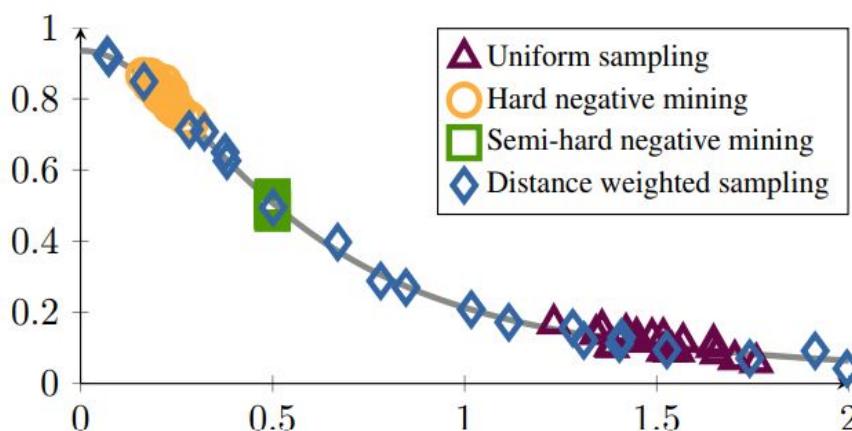
Negative Sampling in Ranking-based DML

- **Infeasible** amount of triplets to learn from (Scale $\mathcal{O}(n^3)$)
- Fixed sampling heuristics:

- **Distance-weighted sampling** [Wu et al. 2017]

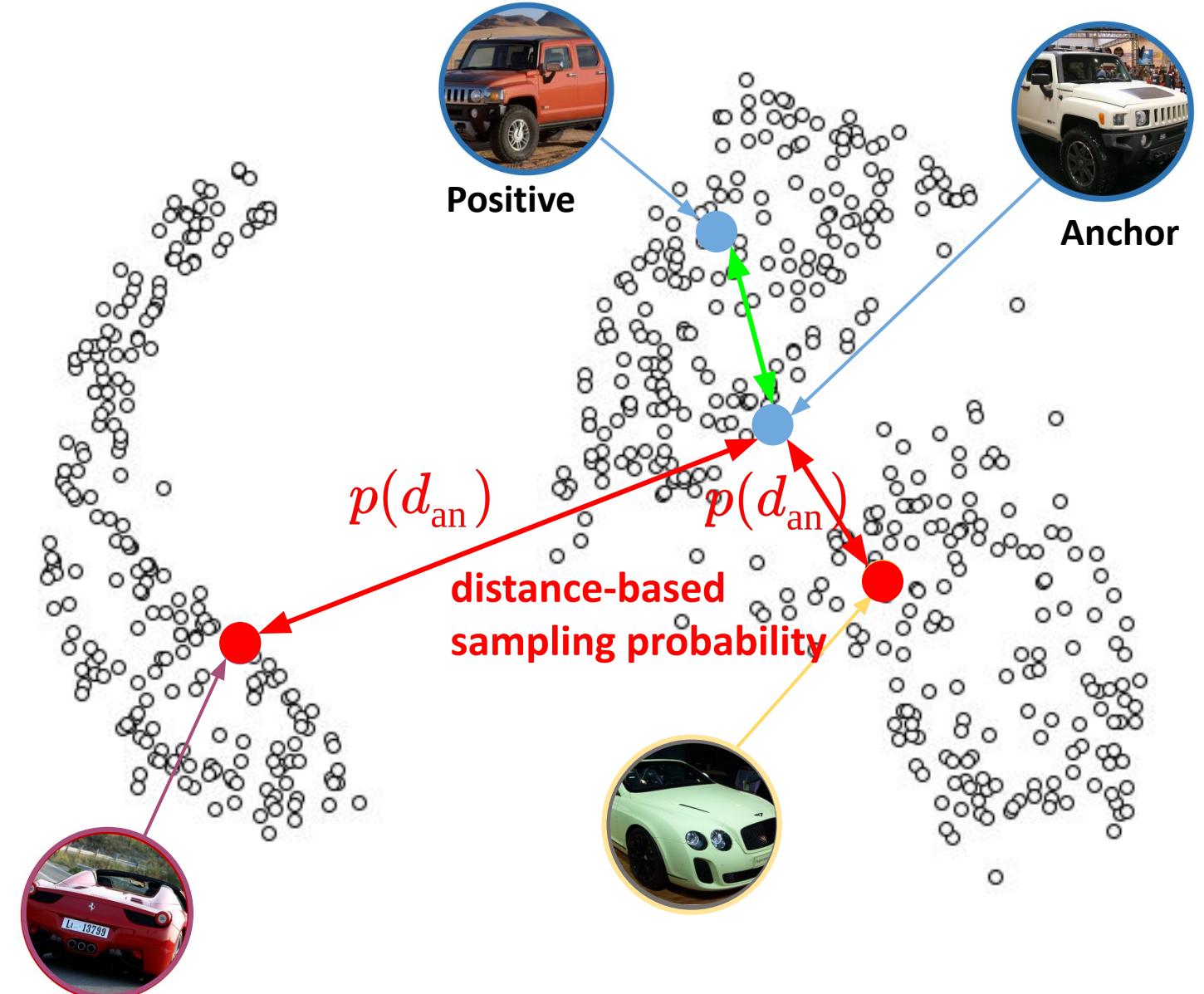
$$n^* \sim p(I_n | I_a) \propto \min(\lambda, q(d_{an})^{-1})$$

- Soften hard negative constraint by uniform sampling from **entire range of distances** d_{an}
- Uniform distribution on \mathbb{S}^D : $q(d) \propto d^{k-2} [1 - \frac{1}{4}d^2]^{\frac{k-3}{2}}$

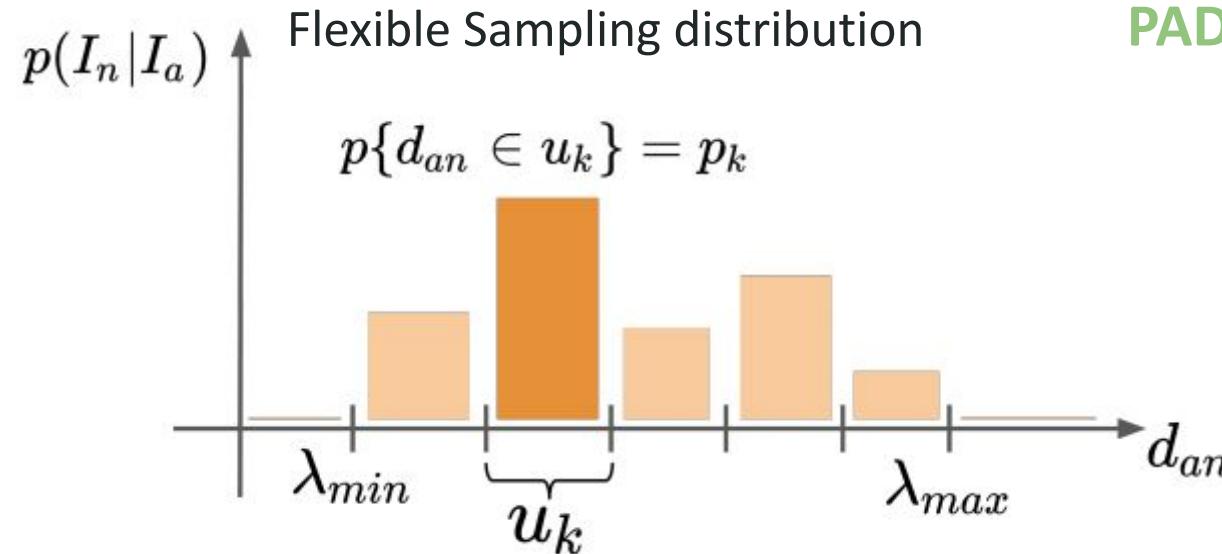


Negative Sampling in Ranking-based DML

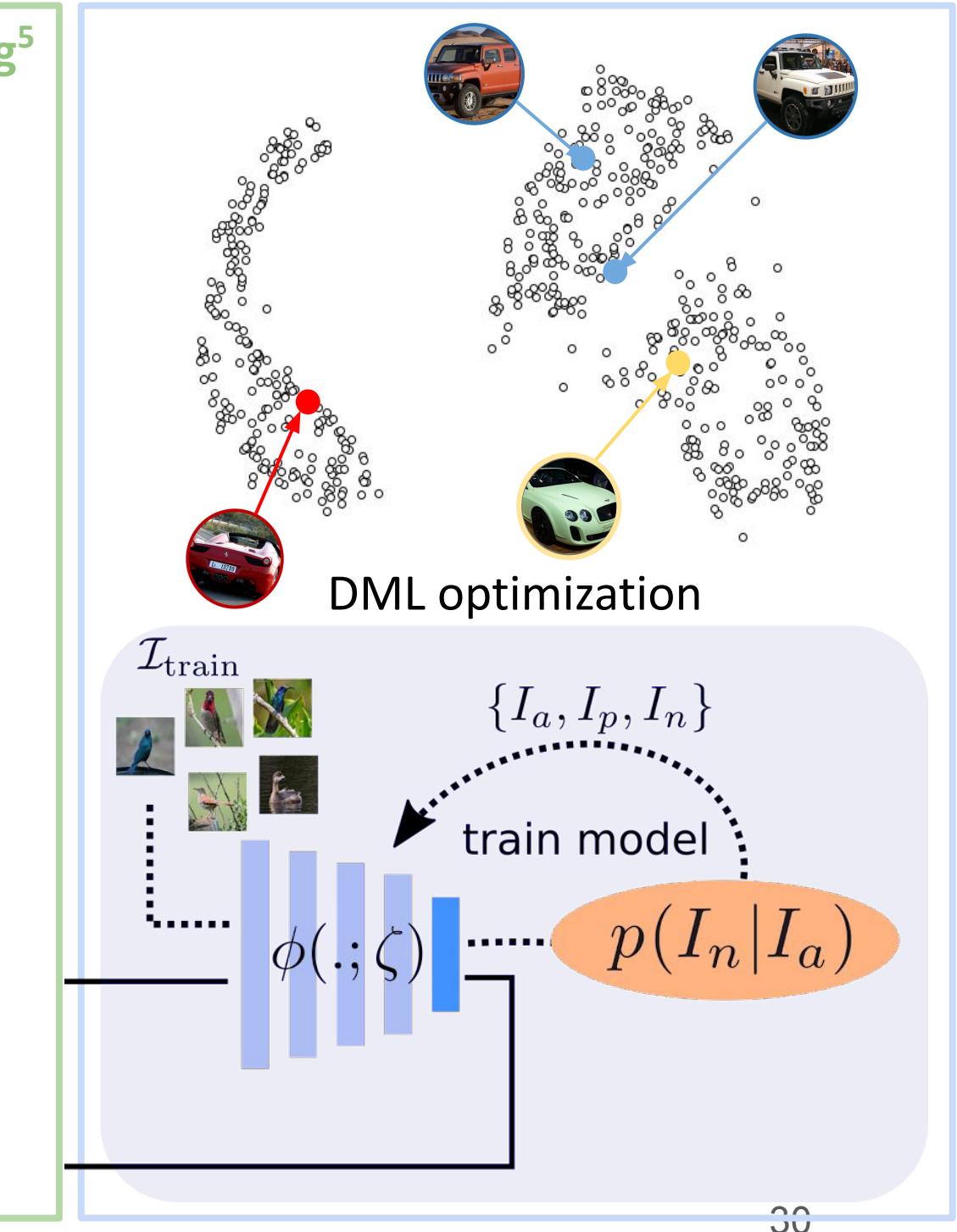
- **Infeasible** amount of triplets to learn from (Scale $\mathcal{O}(n^3)$)
- Fixed sampling heuristics:
 - **(Semi-)Hard negative mining** [Schroff et al. 2015]
 - **Distance-weighted sampling** [Wu et al. 2017]
- **Drawbacks:**
 - Predefined and independent of DML optimization
 - Fixed and disconnected from learning process
(see **Curriculum Learning**)
 - Optimal distribution $p(d_{an})$?



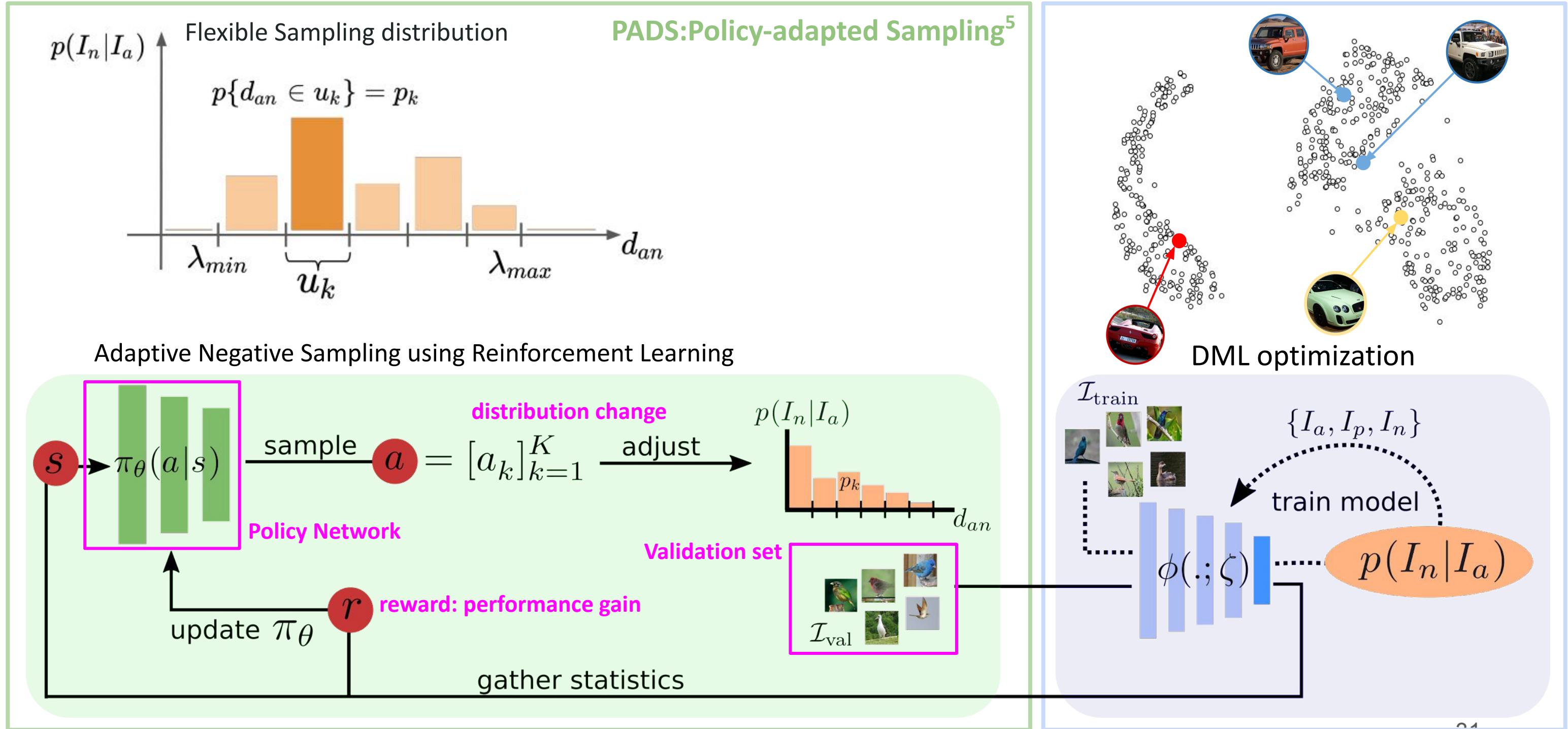
Negative Sampling in Ranking-based DML



PADS:Policy-adapted Sampling⁵



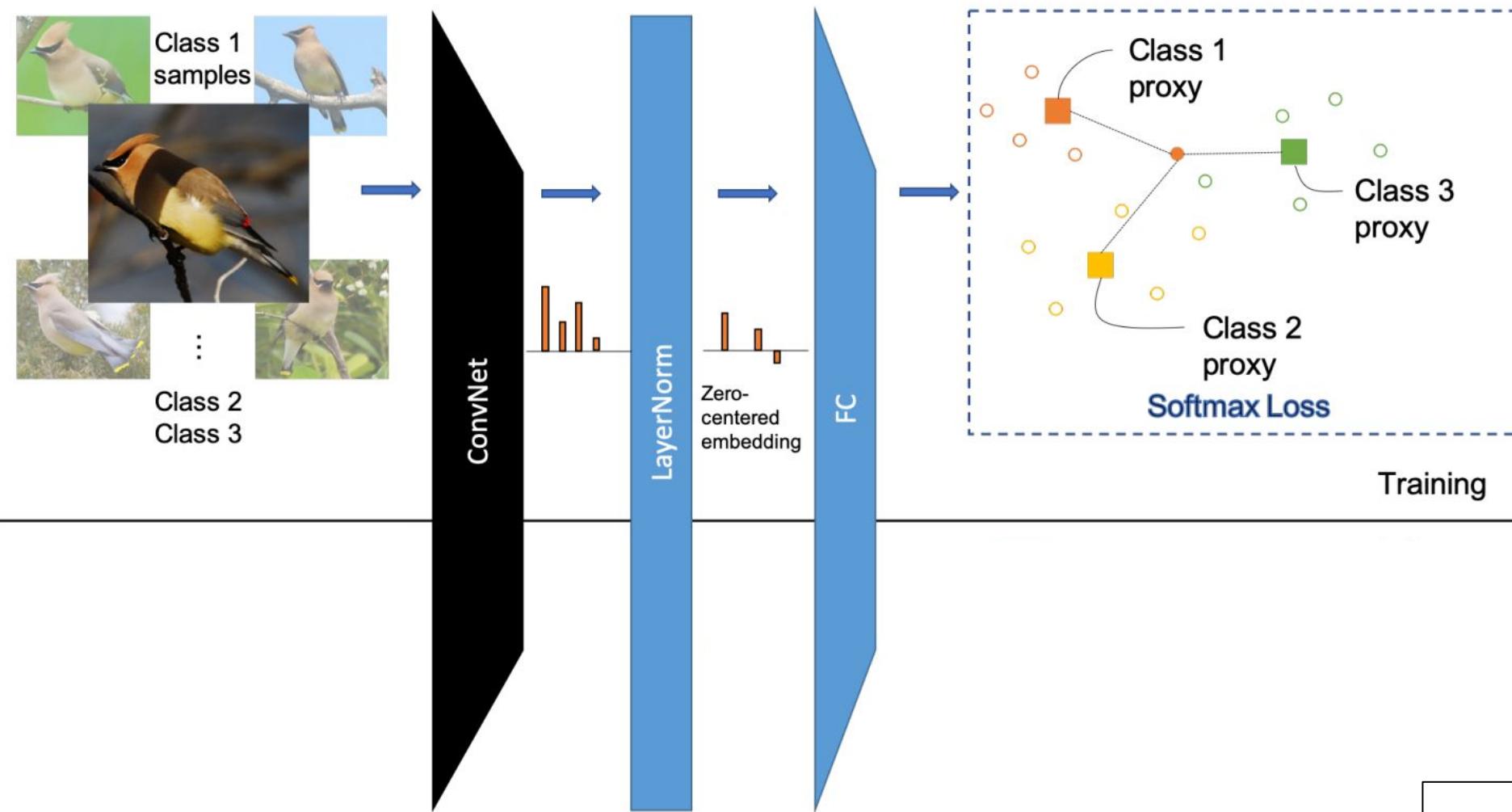
Negative Sampling in Ranking-based DML



Classification-based DML

Classification-based DML

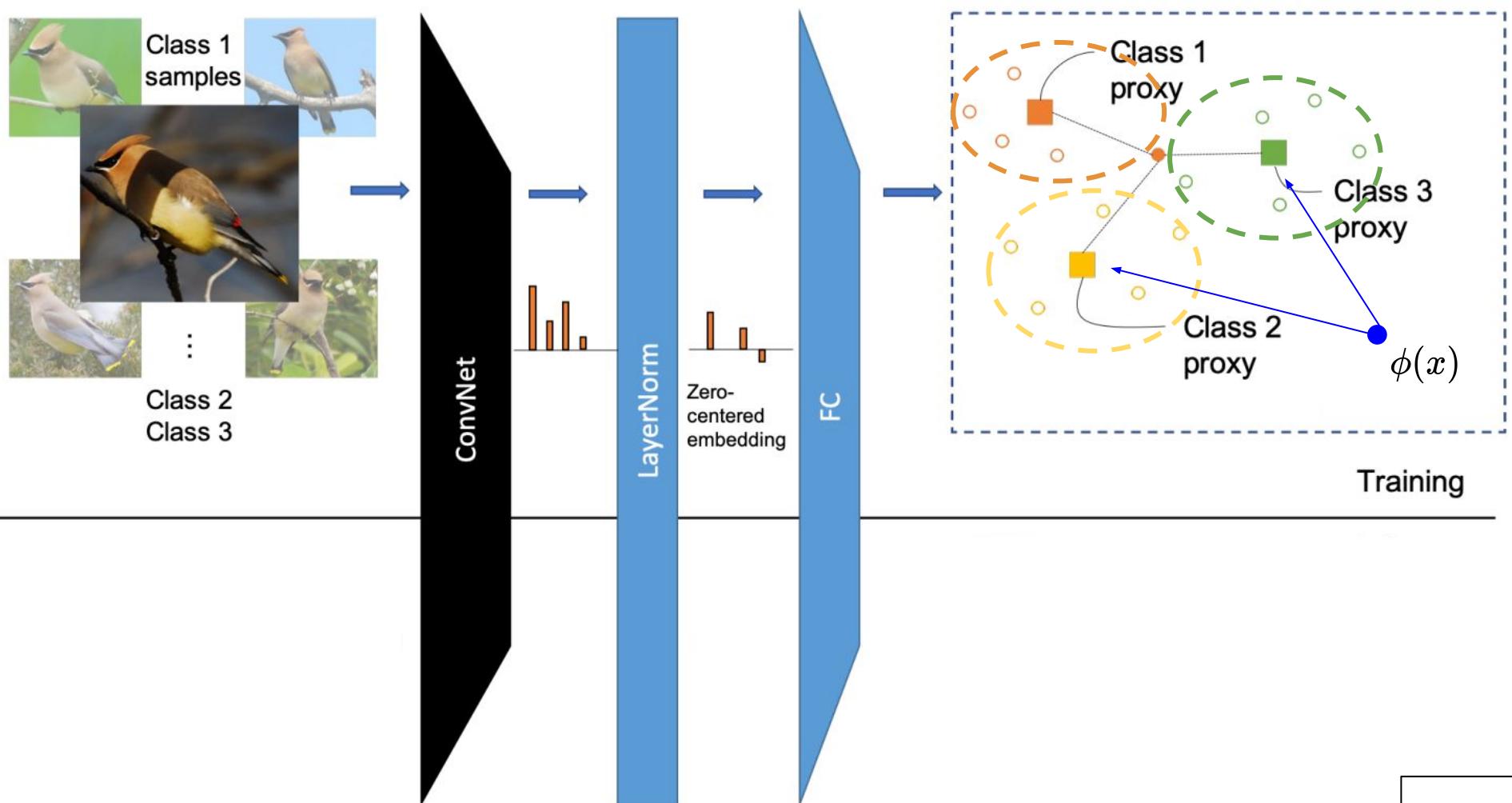
- Classification is “classical” approach to representation learning



Softmax Loss :

$$\sigma(\phi(x_i)) = - \log \left[\frac{\exp(W_{y_i}^\top \phi(x_i))}{\sum_{k=1}^K \exp(W_k^\top \phi(x_i))} \right]$$

Classification-based DML

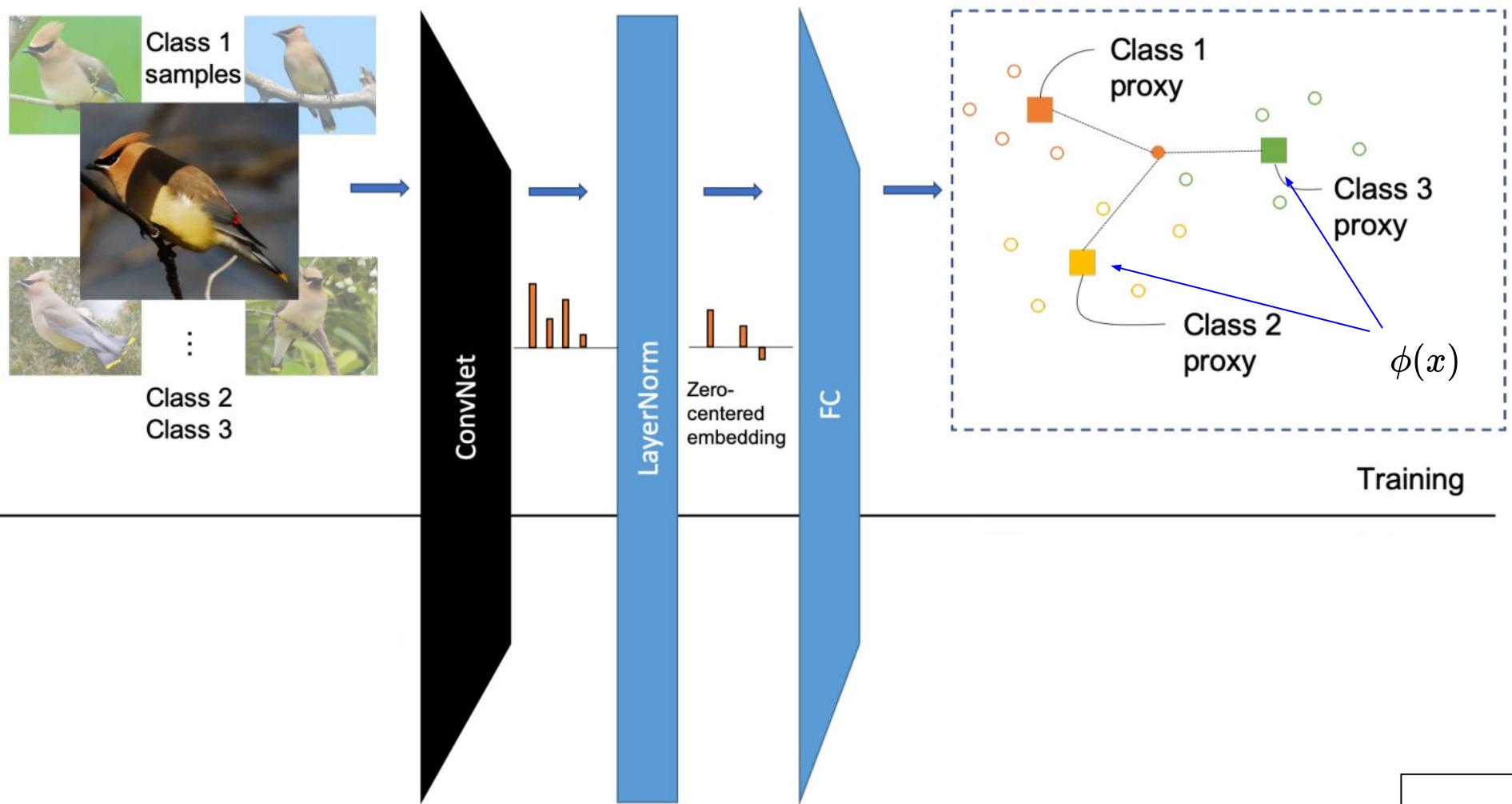


- **Classification** is “classical” approach to representation learning
- Rows of weights matrix W_k act as **class distribution proxies**
 - alleviate the sampling problem of ranking-based methods
- (Connections to **Proxy-based DML**)

Softmax Loss :

$$\sigma(\phi(x_i)) = - \log \left[\frac{\exp(W_{y_i}^\top \phi(x_i))}{\sum_{k=1}^K \exp(W_k^\top \phi(x_i))} \right]$$

Classification-based DML



- Classification is “classical” approach to representation learning
- Rows of weights matrix W_k act as **class distribution proxies**
 - alleviate the sampling problem of ranking-based methods
- (Connections to **Proxy-based DML**)

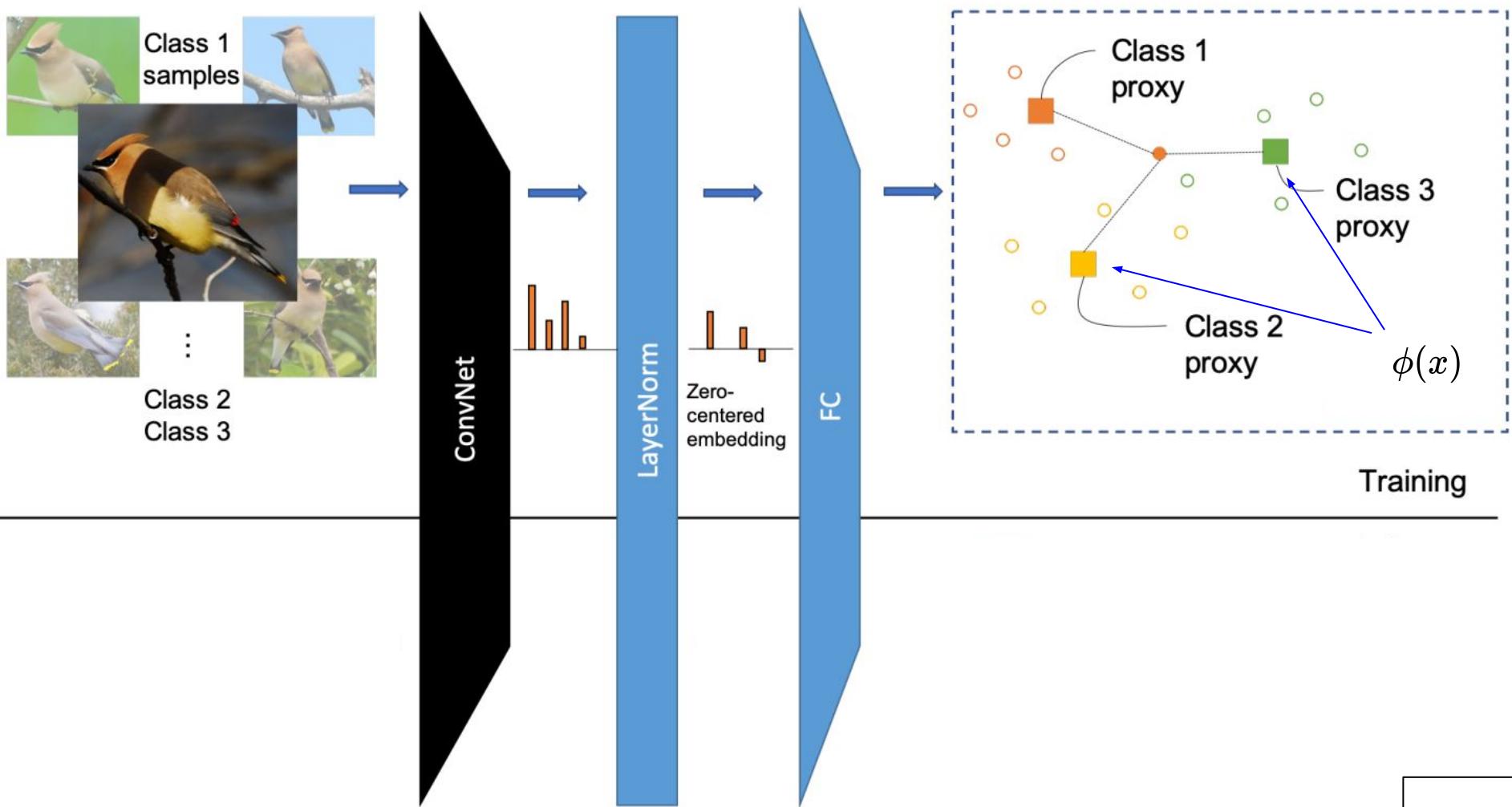
Classification as **competitive baseline**⁶:

- **Regularize** embedding space to hypersphere

Normalized Softmax Loss : L2 normalize

$$\sigma(\phi(x_i)) = - \log \left[\frac{\exp(W_{y_i}^\top \phi(x_i))}{\sum_{k=1}^K \exp(W_k^\top \phi(x_i))} \right]$$

Classification-based DML



- Classification is “classical” approach to representation learning
- Rows of weights matrix W_k act as **class distribution proxies**
→ alleviate the sampling problem of ranking-based methods
- (Connections to **Proxy-based DML**)

Classification as **competitive baseline**⁶:

- **Regularize** embedding space to hypersphere
- **Temperature scaling** to enforce compact intra-class clusters

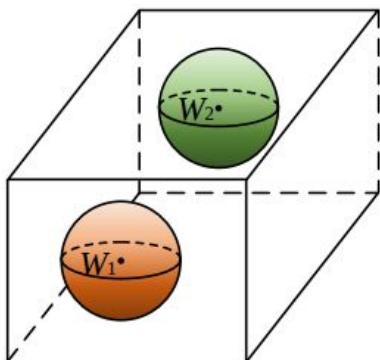
Normalized Softmax Loss :

L2 normalize	Temperature scaling
--------------	---------------------

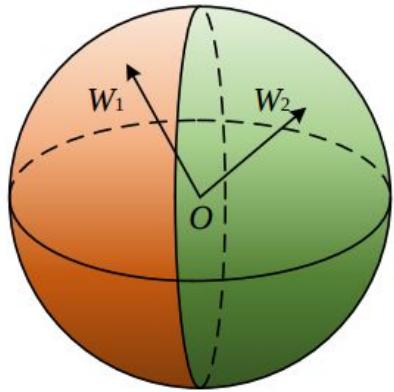
$$\sigma(\phi(x_i)) = - \log \left[\frac{\exp(W_{y_i}^\top \phi(x_i) / \tau)}{\sum_{k=1}^K \exp(W_k^\top \phi(x_i) / \tau)} \right]$$

Classification-based DML

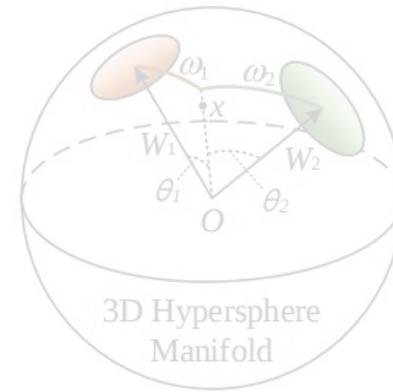
- So far, we optimize actual **distances** (e.g. Euclidean, Cosine) between data samples/proxies
- Now: Express learning constraints explicitly as actual **angles on (Hypersphere-)Manifold**
- very popular for **Face recognition applications**



Euclidean Margin Loss



Modified Softmax Loss



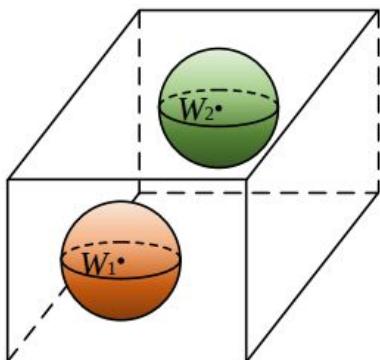
A-Softmax Loss ($m \geq 2$)

$$\sigma(\phi(x_i)) = -\log \left[\frac{\exp(W_{y_i}^\top \phi(x_i))}{\sum_{k=1}^K \exp(W_k^\top \phi(x_i))} \right]$$

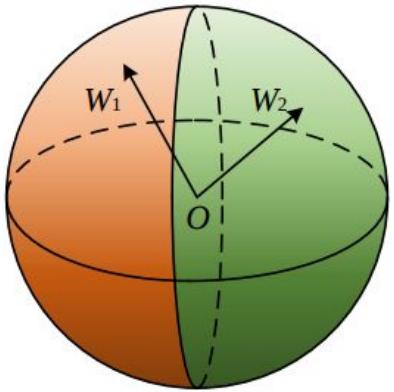
$x^T y$

Classification-based DML

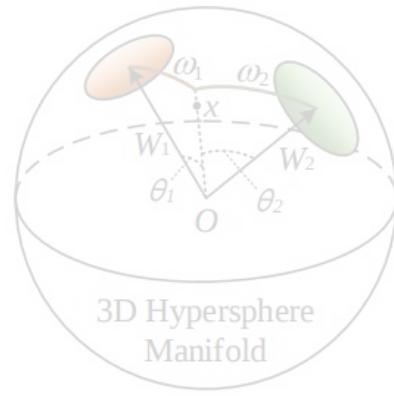
- So far, we optimize actual **distances** (e.g. Euclidean, Cosine) between data samples/proxies
- Now: Express learning constraints explicitly as actual **angles on (Hypersphere-)Manifold**
- very popular for **Face recognition applications**



Euclidean Margin Loss



Modified Softmax Loss



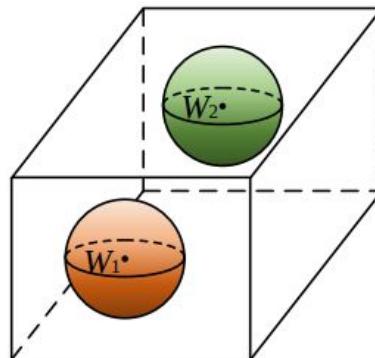
A -Softmax Loss ($m \geq 2$)

$$\begin{aligned}\sigma(\phi(x_i)) &= -\log \left[\frac{\exp(W_{y_i}^\top \phi(x_i))}{\sum_{k=1}^K \exp(W_k^\top \phi(x_i))} \right] \\ &= -\log \left[\frac{\exp(\|W_{y_i}\| \|\phi(x_i)\| \cos(\theta_{i,y_i}))}{\sum_{k=1}^K \exp(\|W_k\| \|\phi(x_i)\| \cos(\theta_{i,k}))} \right].\end{aligned}$$

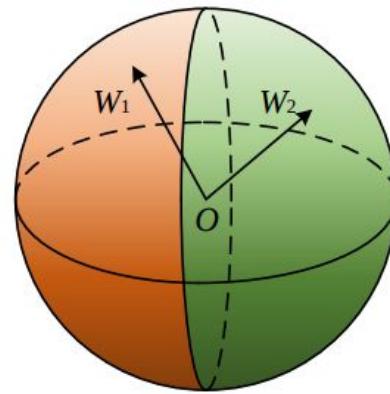
$x^T y = \|x\|_2 \|y\|_2 \cos(\theta)$

Classification-based DML

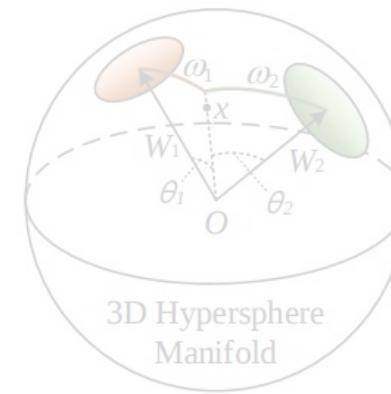
- So far, we optimize actual **distances** (e.g. Euclidean, Cosine) between data samples/proxies
- Now: Express learning constraints explicitly as actual **angles on (Hypersphere-)Manifold**
- very popular for **Face recognition applications**



Euclidean Margin Loss



Modified Softmax Loss



A -Softmax Loss ($m \geq 2$)

$$\begin{aligned}\sigma(\phi(x_i)) &= -\log \left[\frac{\exp(W_{y_i}^\top \phi(x_i))}{\sum_{k=1}^K \exp(W_k^\top \phi(x_i))} \right] \\ &= -\log \left[\frac{\exp(\|W_{y_i}\| \|\phi(x_i)\| \cos(\theta_{i,y_i}))}{\sum_{k=1}^K \exp(\|W_k\| \|\phi(x_i)\| \cos(\theta_{i,k}))} \right].\end{aligned}$$

$x^T y = \|x\|_2 \|y\|_2 \cos(\theta)$

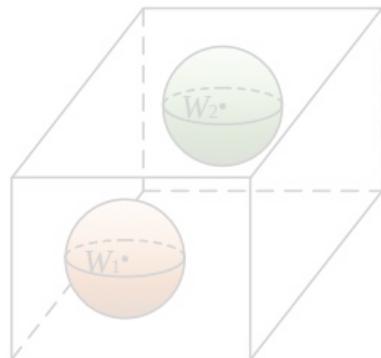
- Constrain $\|W_{y_i}\| = 1$ and $\|\phi(x_i)\|_2 = s$

$$\sigma(\phi(x_i)) = -\log \left[\frac{\exp(s \cos(\theta_{i,y_i}))}{\sum_{k=1}^K \exp(s \cos(\theta_{i,k}))} \right].$$

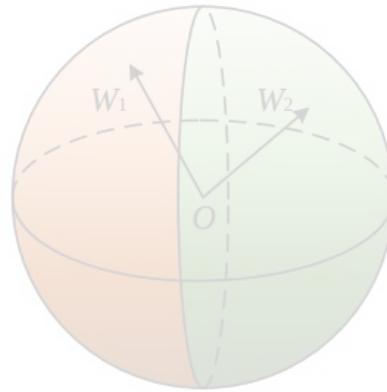
- We now optimize the angles $\theta_{i,k}$ between $\phi(x_i)$ and W_k

Classification-based DML

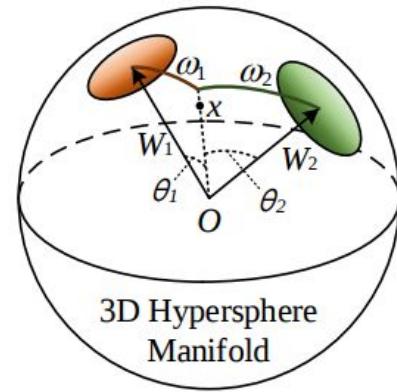
- So far, we optimize actual **distances** (e.g. Euclidean, Cosine) between data samples/proxies
- Now: Express learning constraints explicitly as actual **angles on (Hypersphere-)Manifold**
- very popular for **Face recognition applications**



Euclidean Margin Loss



Modified Softmax Loss



A-Softmax Loss ($m \geq 2$)

$$\begin{aligned}\sigma(\phi(x_i)) &= -\log \left[\frac{\exp(W_{y_i}^\top \phi(x_i))}{\sum_{k=1}^K \exp(W_k^\top \phi(x_i))} \right] \\ &= -\log \left[\frac{\exp(\|W_{y_i}\| \|\phi(x_i)\| \cos(\theta_{i,y_i}))}{\sum_{k=1}^K \exp(\|W_k\| \|\phi(x_i)\| \cos(\theta_{i,k}))} \right].\end{aligned}$$

$x^T y = \|x\|_2 \|y\|_2 \cos(\theta)$

- Constrain $\|W_{y_i}\| = 1$ and $\|\phi(x_i)\|_2 = s$

$$\sigma(\phi(x_i)) = -\log \left[\frac{\exp(s \cos(\theta_{i,y_i}))}{\sum_{k=1}^K \exp(s \cos(\theta_{i,k}))} \right].$$

- We now optimize the angles $\theta_{i,k}$ between $\phi(x_i)$ and W_k
- Introduce a margin β similar to Ranking-based DML

Sphere⁷-/ArcFace⁸ Loss:

$$\sigma(\phi(x_i); \beta) = -\log \left[\frac{\exp(s \cos(\beta + \theta_{i,y_i}))}{\exp(s \cos(\beta + \theta_{i,y_i})) + \sum_{k=1, k \neq y_i}^M \exp(s \cos(\theta_{i,k}))} \right].$$

fixed margin

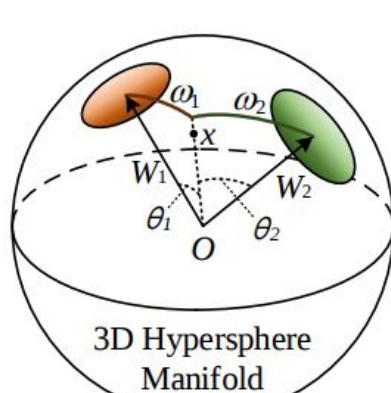
Classification-based DML

- So far, we optimize actual **distances** (e.g. Euclidean, Cosine) between data samples/proxies
- Now: Express learning constraints explicitly as actual **angles** on (**Hypersphere**-)Manifold
- very popular for **Face recognition applications**

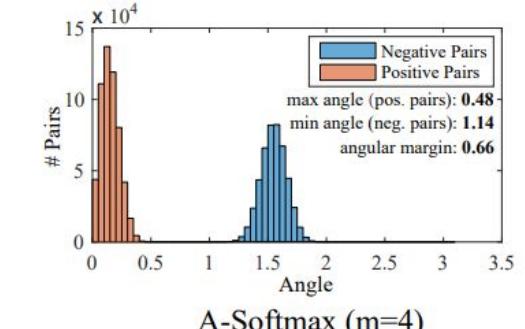
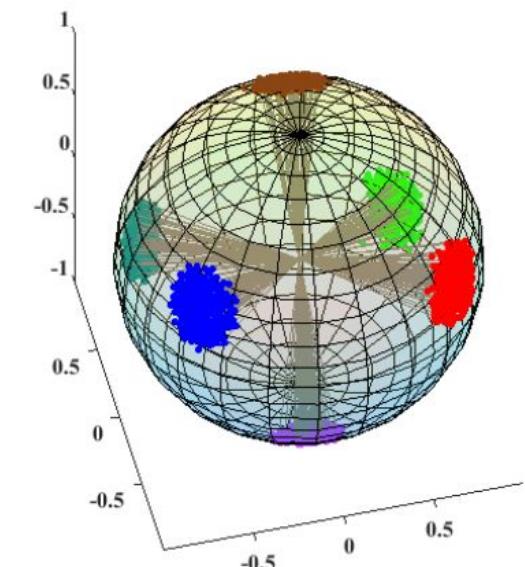
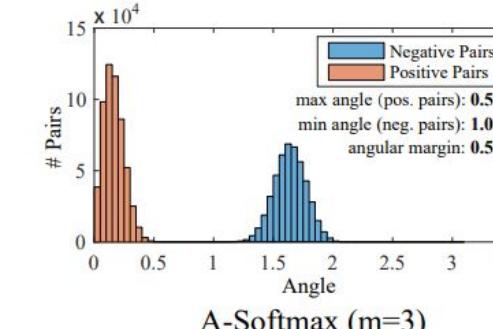
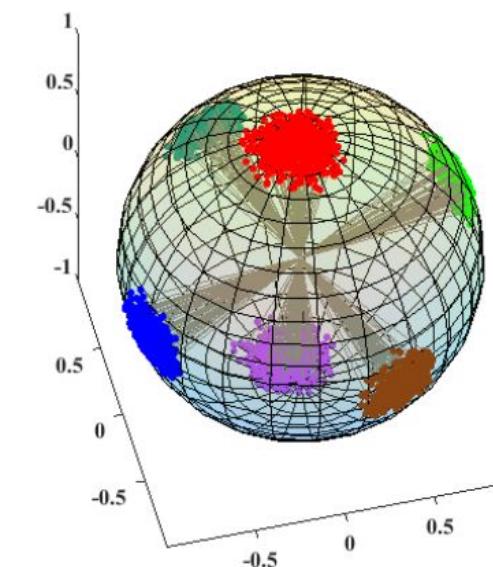
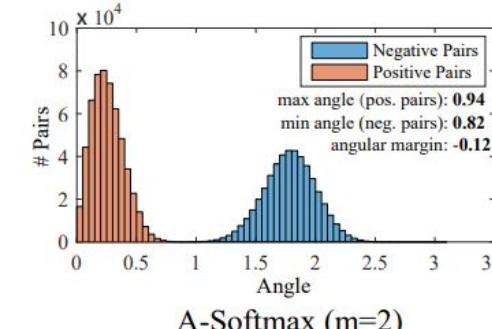
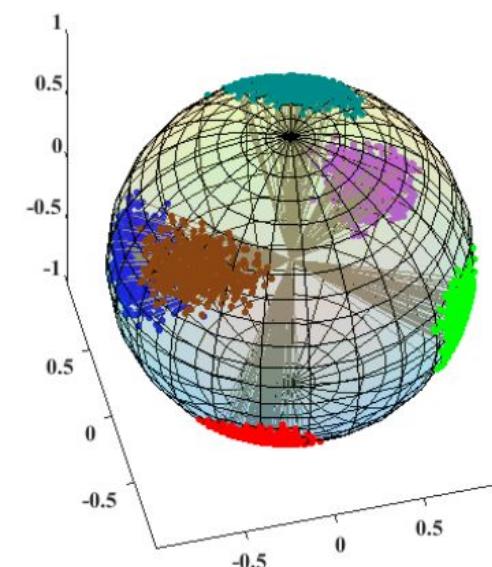
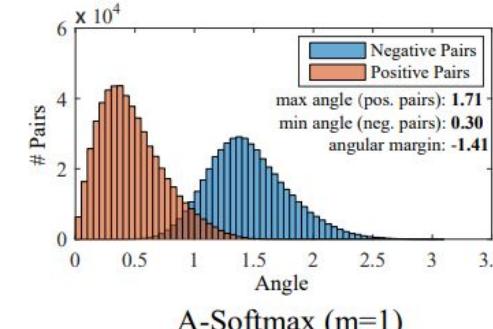
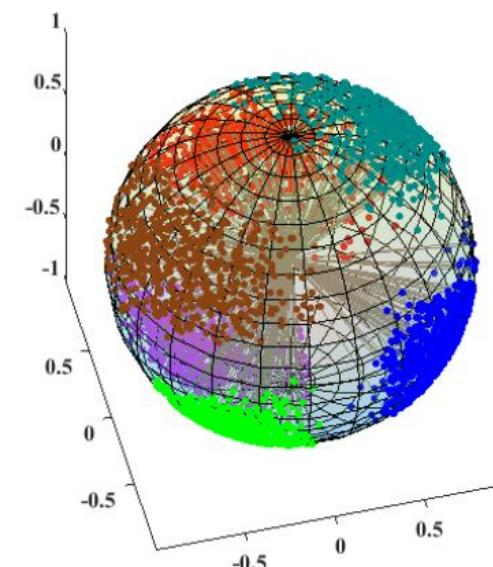
Many different extensions based on this formulations!

Sphere⁷-/ArcFace⁸ Loss:

$$\sigma(\phi(x_i); \beta) = -\log \left[\frac{\exp(s \cos(\beta + \theta_{i,y_i}))}{\exp(s \cos(\beta + \theta_{i,y_i})) + \sum_{k=1, k \neq y_i}^M \exp(s \cos(\theta_{i,k}))} \right]$$



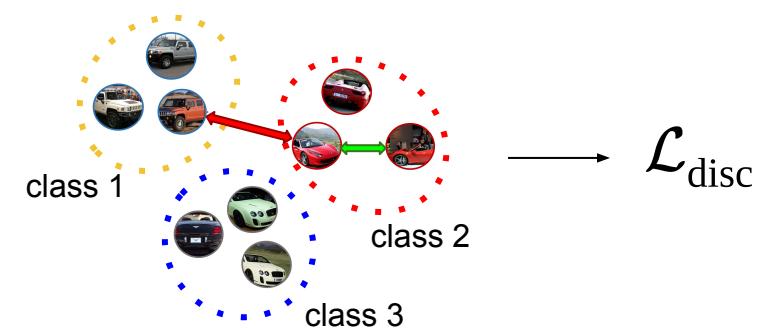
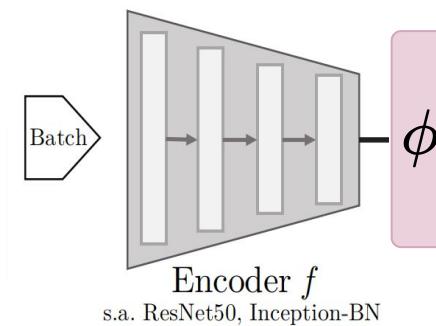
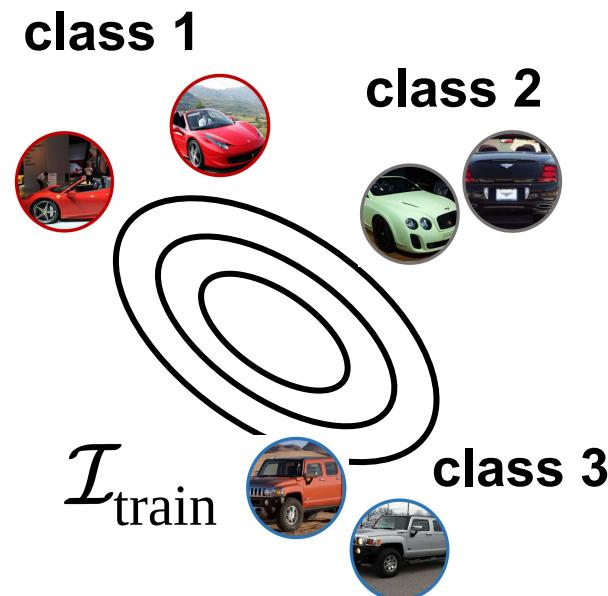
A-Softmax Loss ($m \geq 2$)



Ensemble-based DML

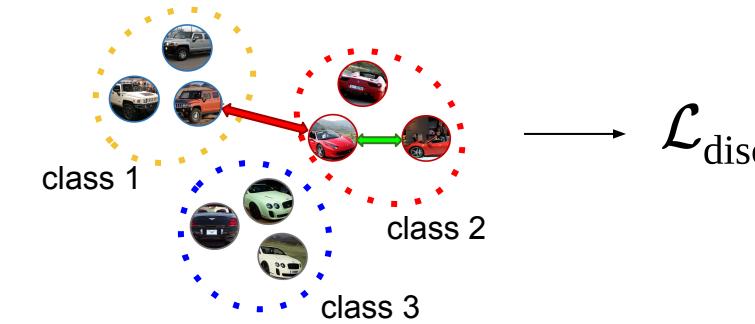
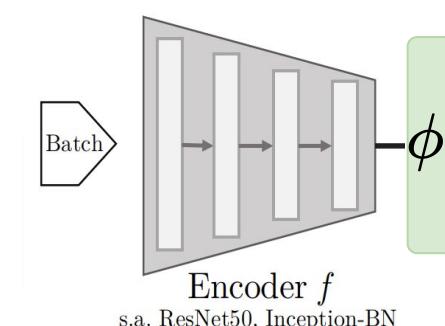
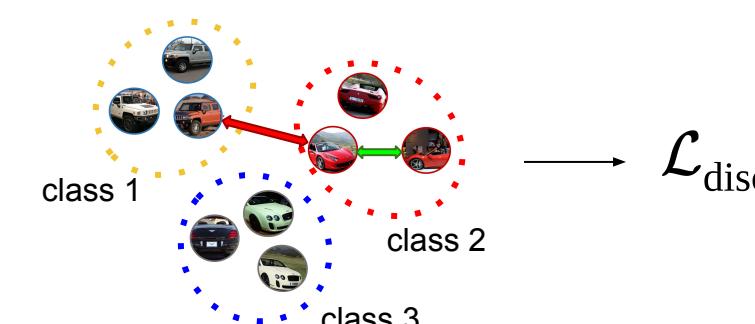
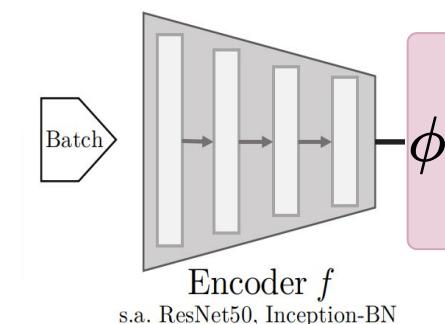
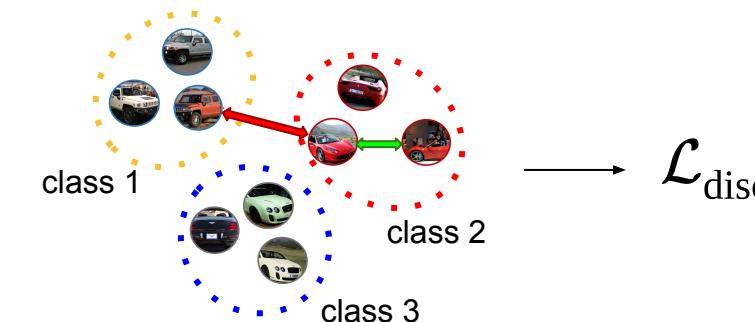
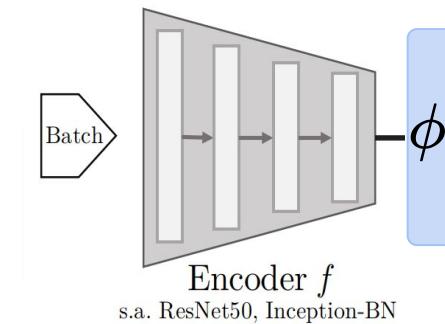
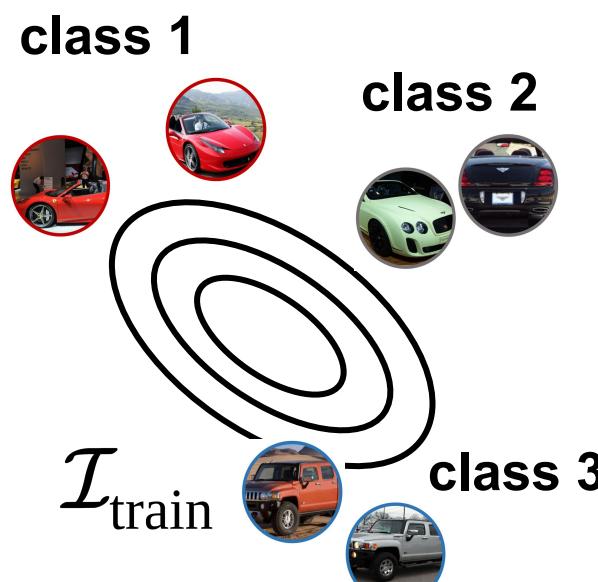
Ensemble-based DML

- Learner **ensembles** are common approach to improve overall performance.



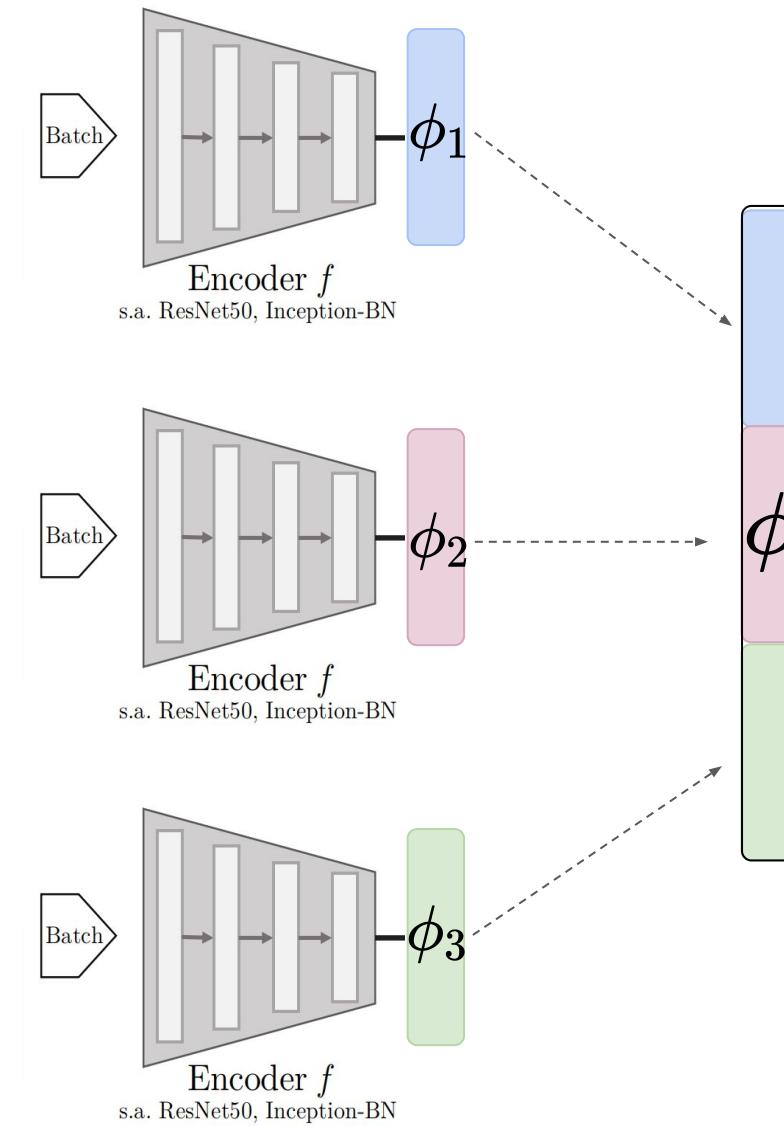
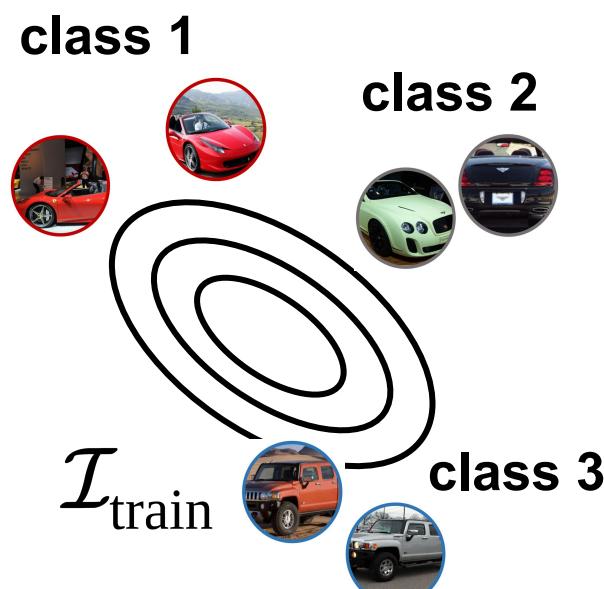
Ensemble-based DML

- **Learner ensembles** are common approach to improve overall performance.
- **Naive approach:** Independently train K independent models using the same, standard discriminative loss $\mathcal{L}_{\text{disc}}$



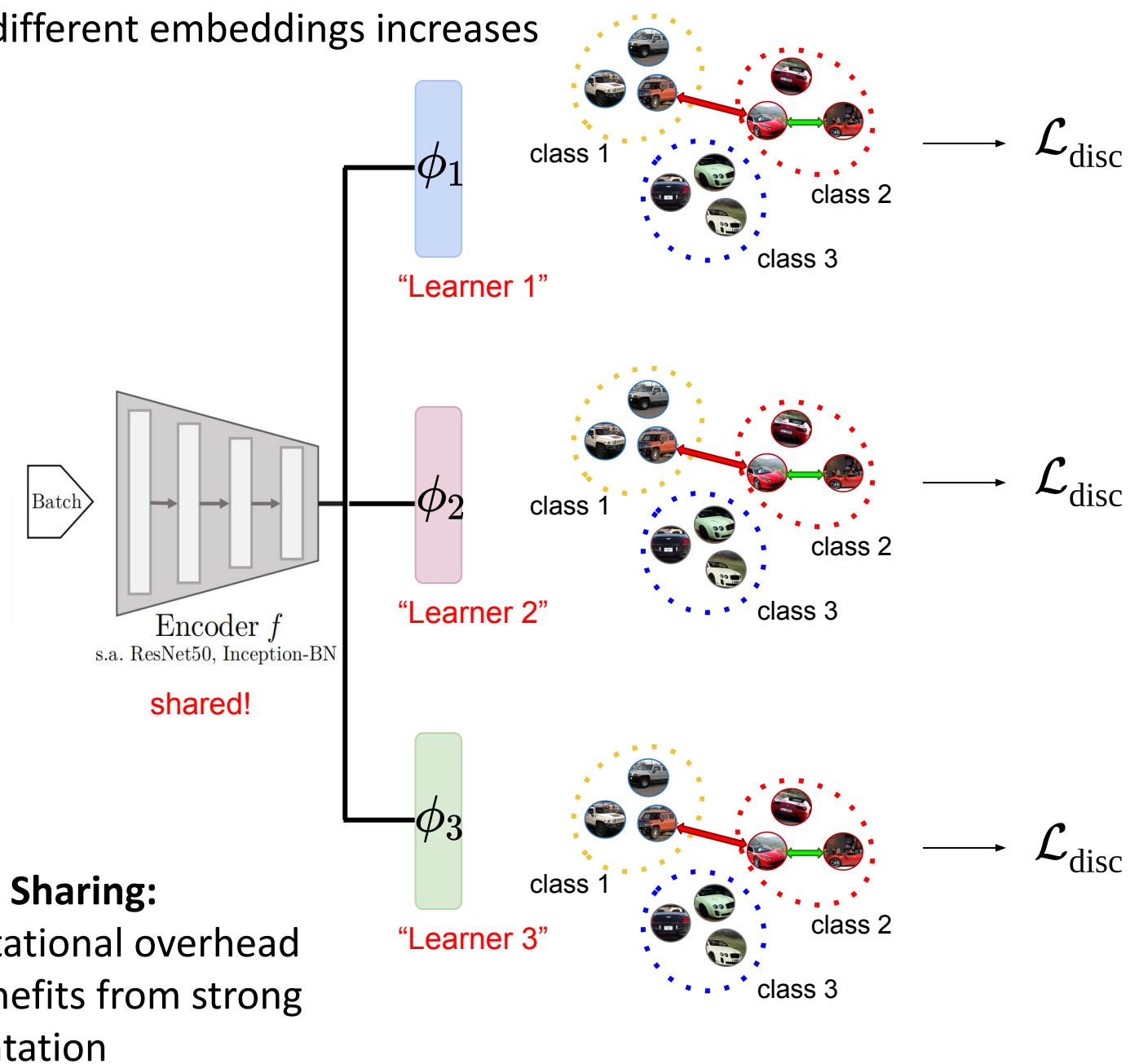
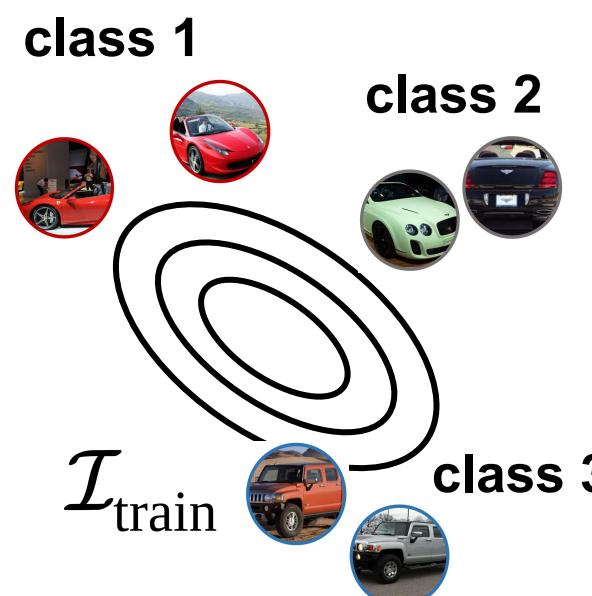
Ensemble-based DML

- **Learner ensembles** are common approach to improve overall performance.
- **Naive approach:** Independently train K different models using the same, standard discriminative loss $\mathcal{L}_{\text{disc}}$
- **Assumption:** Aggregation over different embeddings increases model robustness.



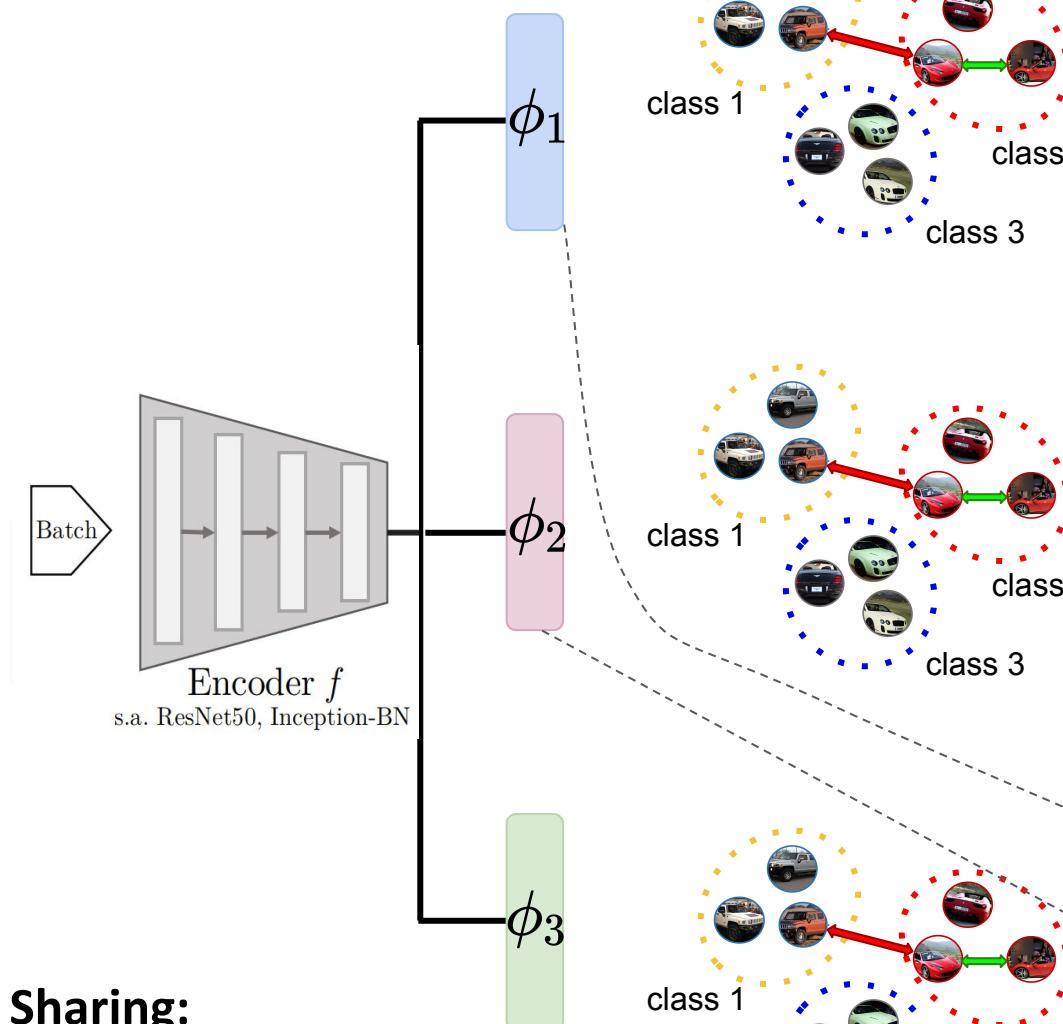
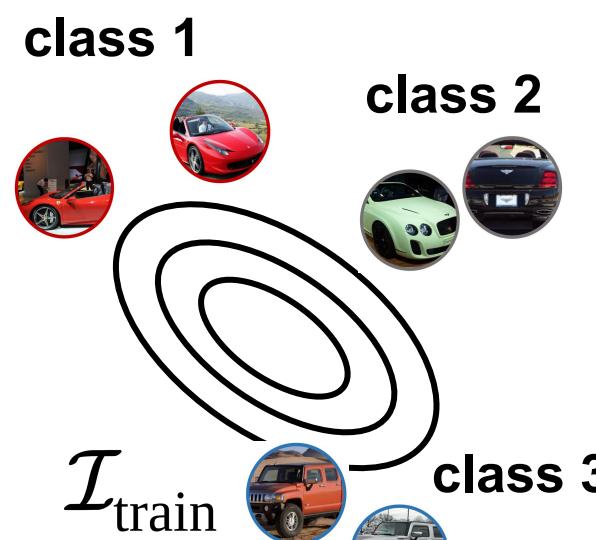
Ensemble-based DML

- **Learner ensembles** are common approach to improve overall performance.
- **Naive approach:** Independently train K different models using the same, standard discriminative loss $\mathcal{L}_{\text{disc}}$
- **Assumption:** Aggregation over different embeddings increases model robustness.



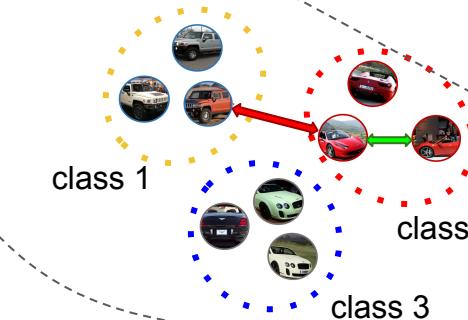
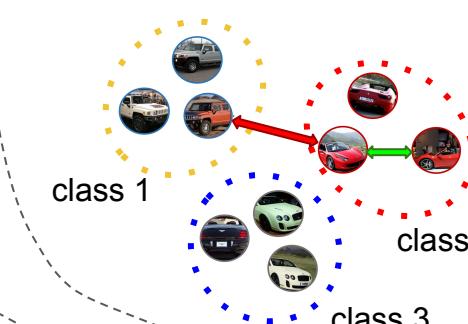
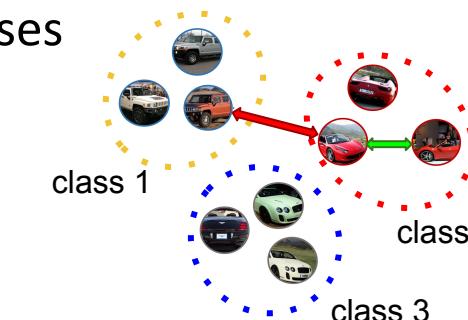
Ensemble-based DML

- **Learner ensembles** are common approach to improve overall performance.
- **Naive approach:** Independently train K different models using the same, standard discriminative loss $\mathcal{L}_{\text{disc}}$
- **Assumption:** Aggregation over different embeddings increases model robustness.

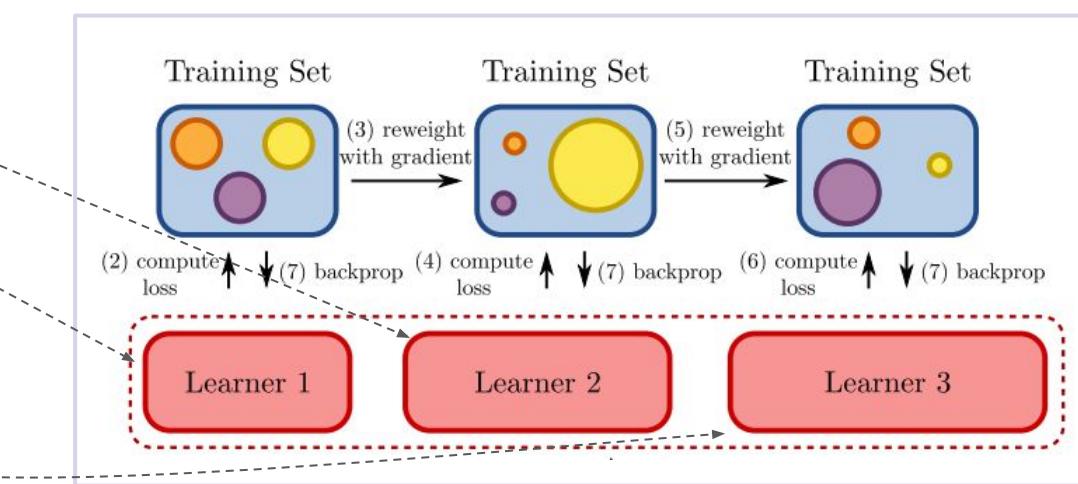


Feature Extractor Sharing:

- mitigate computational overhead
- each learner benefits from strong feature representation

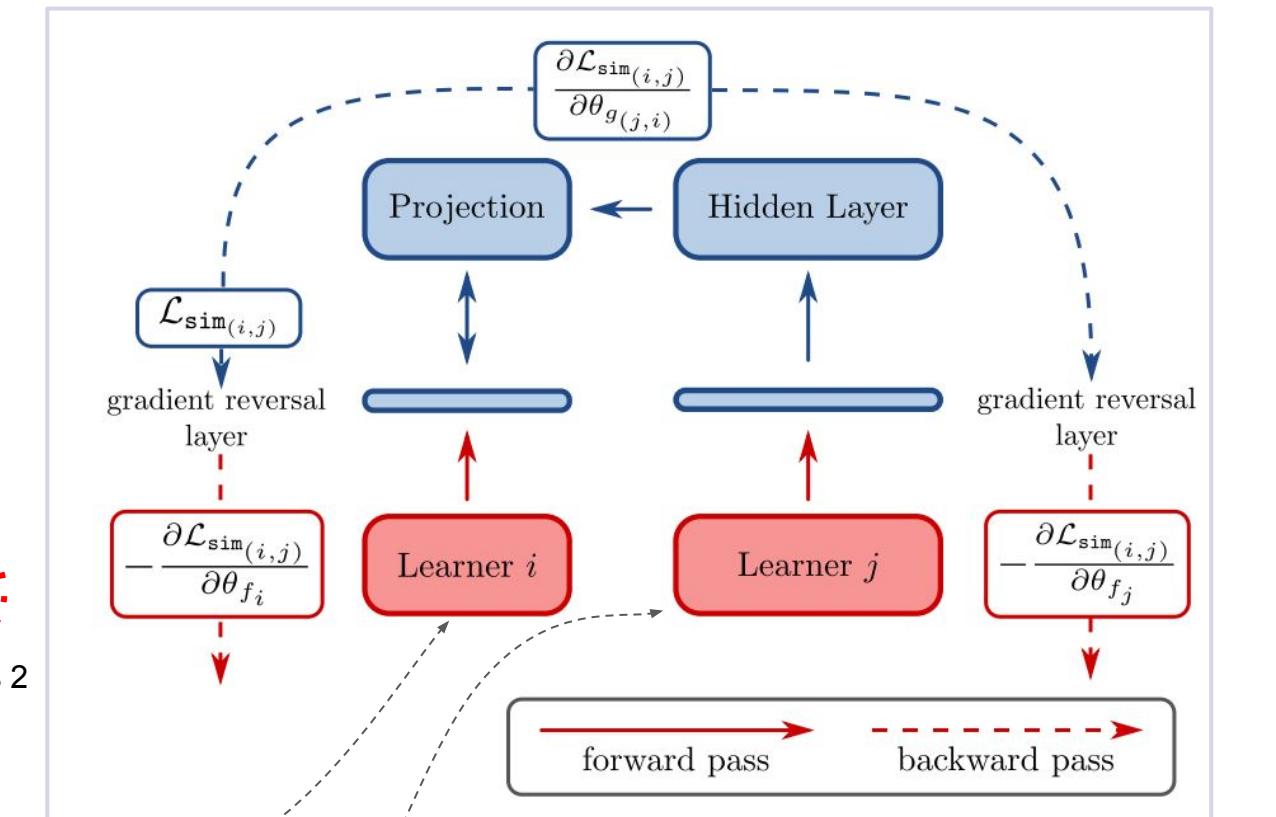
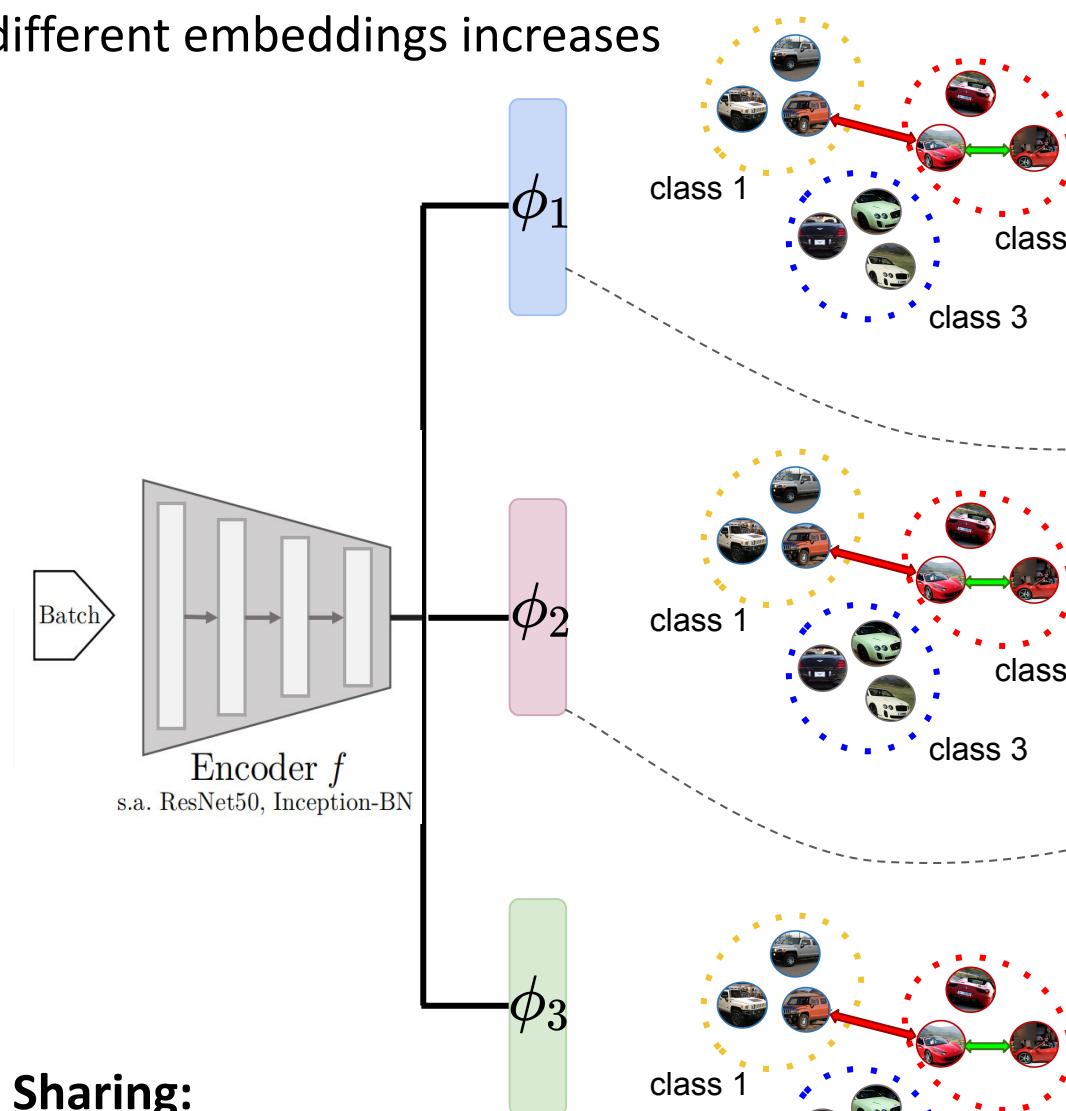
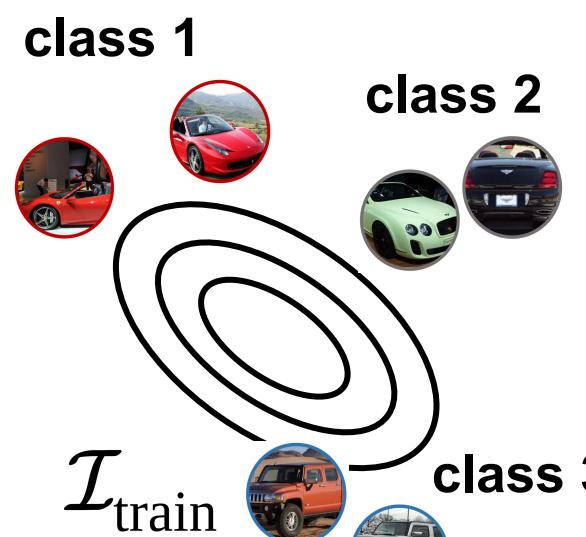


BIER: Online Gradient Boosting [Opitz et al. 2018]:
implicit specialization of learners



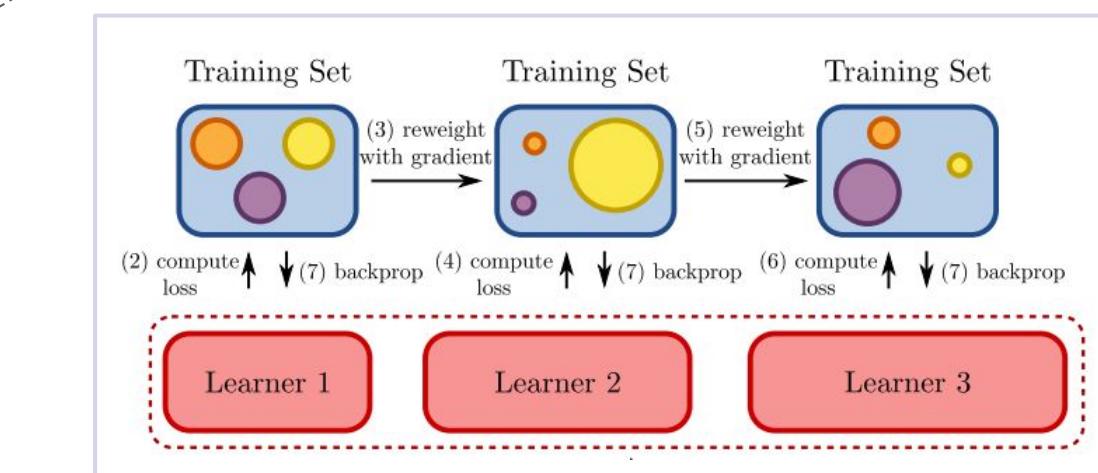
Ensemble-based DML

- **Learner ensembles** are common approach to improve overall performance.
- **Naive approach:** Independently train K different models using the same, standard discriminative loss
- **Assumption:** Aggregation over different embeddings increases model robustness.



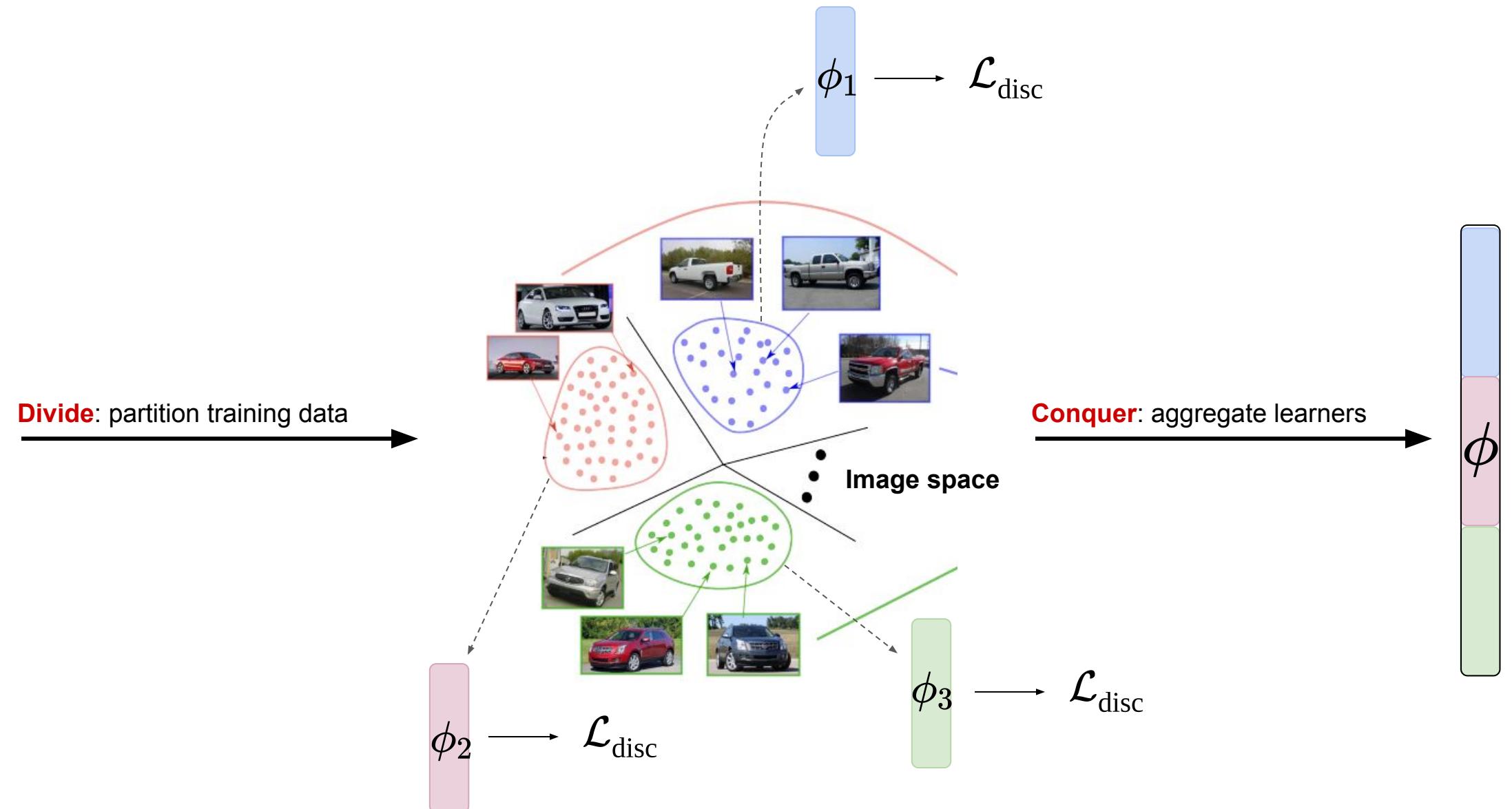
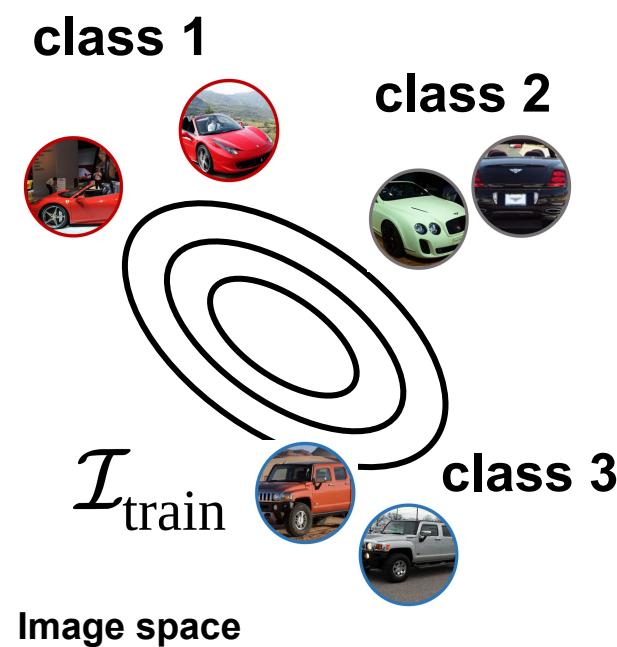
BIER: Adversarial Decorrelation [Opitz et al. 2018]
encourage learners to focus on different features

BIER: Online Gradient Boosting [Opitz et al. 2018]:
implicit specialization of learners



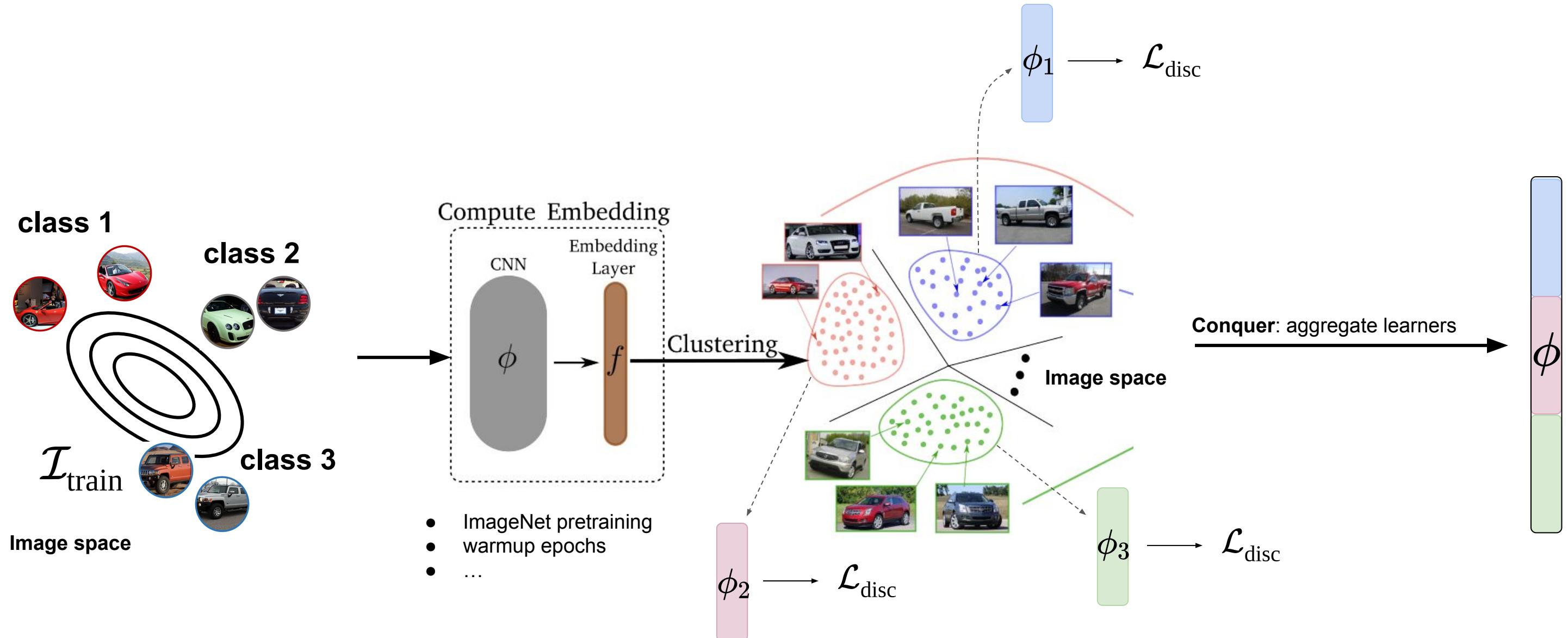
Ensemble-based DML

- Typically training distribution is **multimodal**
- **Explicit** data specialization of learners following **Divide & Conquer strategy**⁹



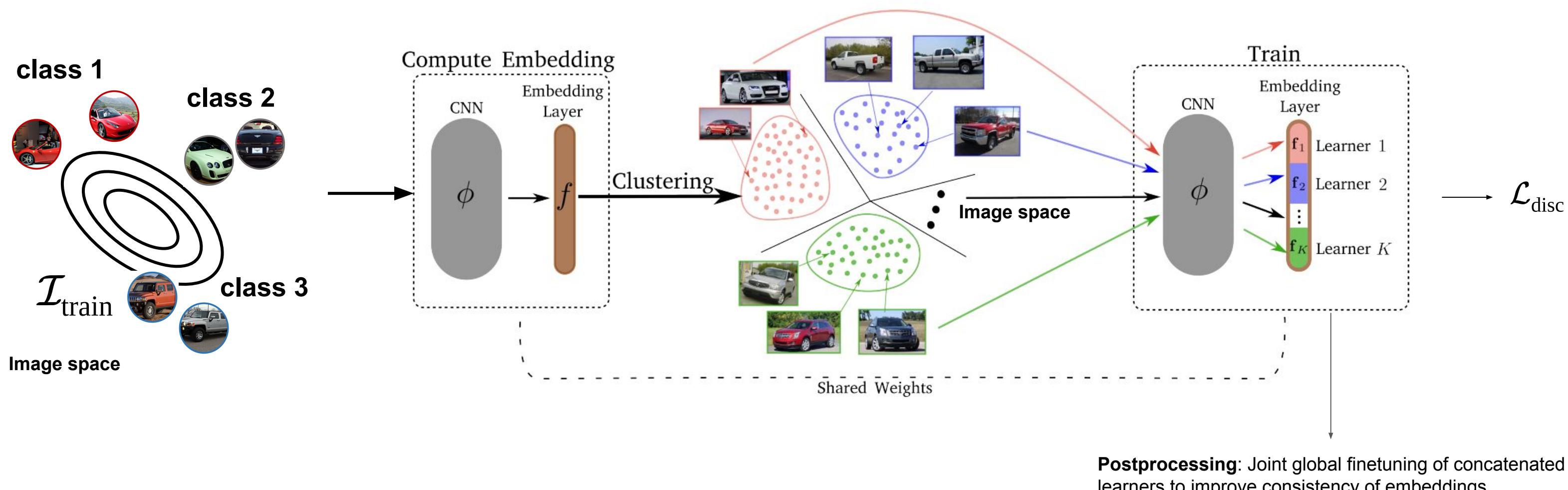
Ensemble-based DML

- Typically training distribution is **multimodal**
- **Explicit** data specialization of learners following **Divide & Conquer strategy**⁹



Ensemble-based DML

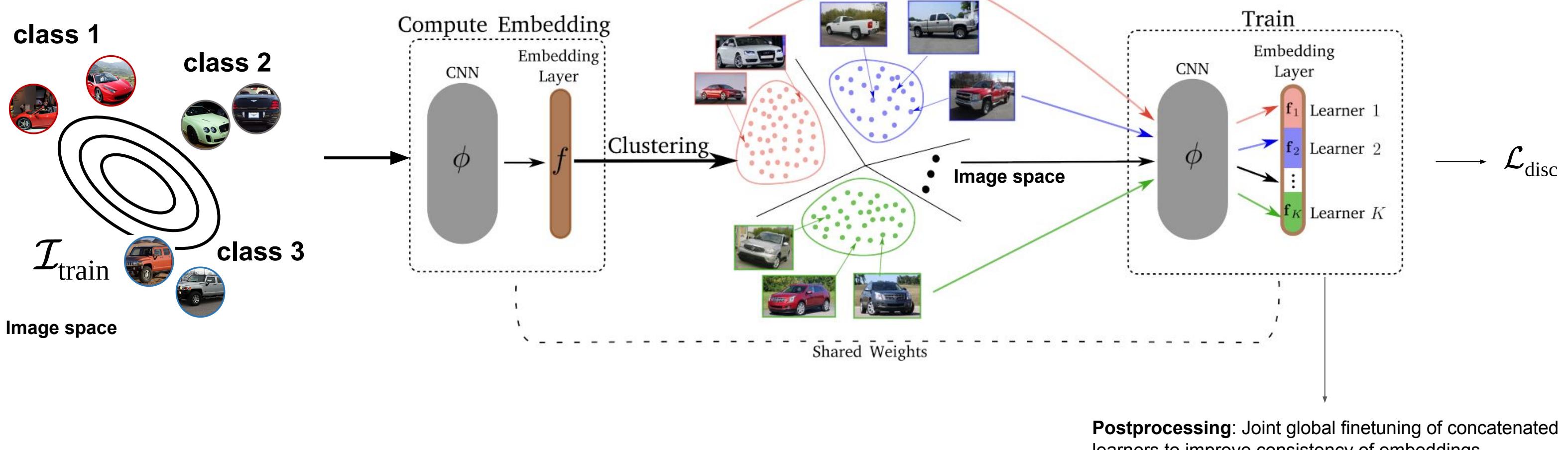
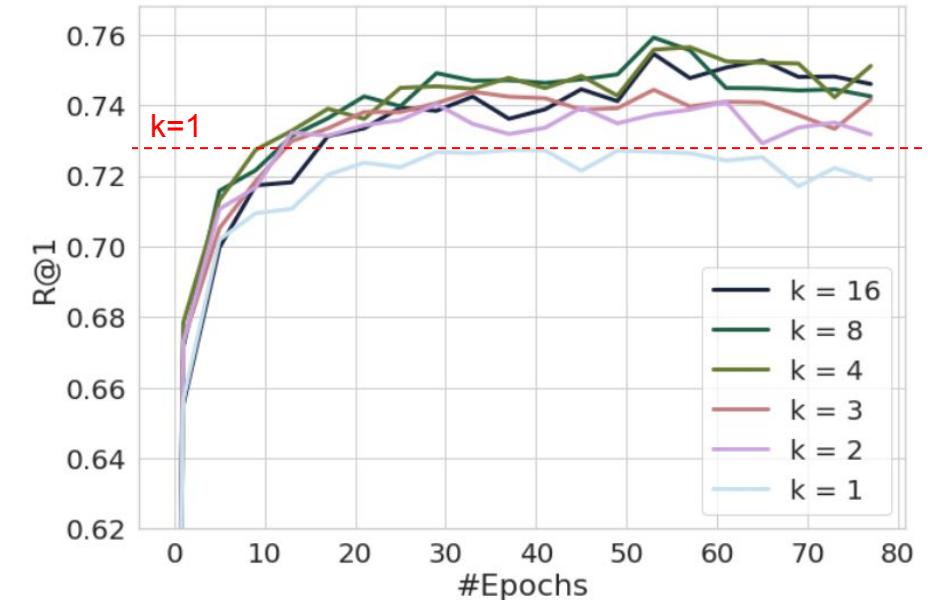
- Typically training distribution is **multimodal**
- **Explicit** data specialization of learners following **Divide & Conquer strategy**⁹



Postprocessing: Joint global finetuning of concatenated learners to improve consistency of embeddings.

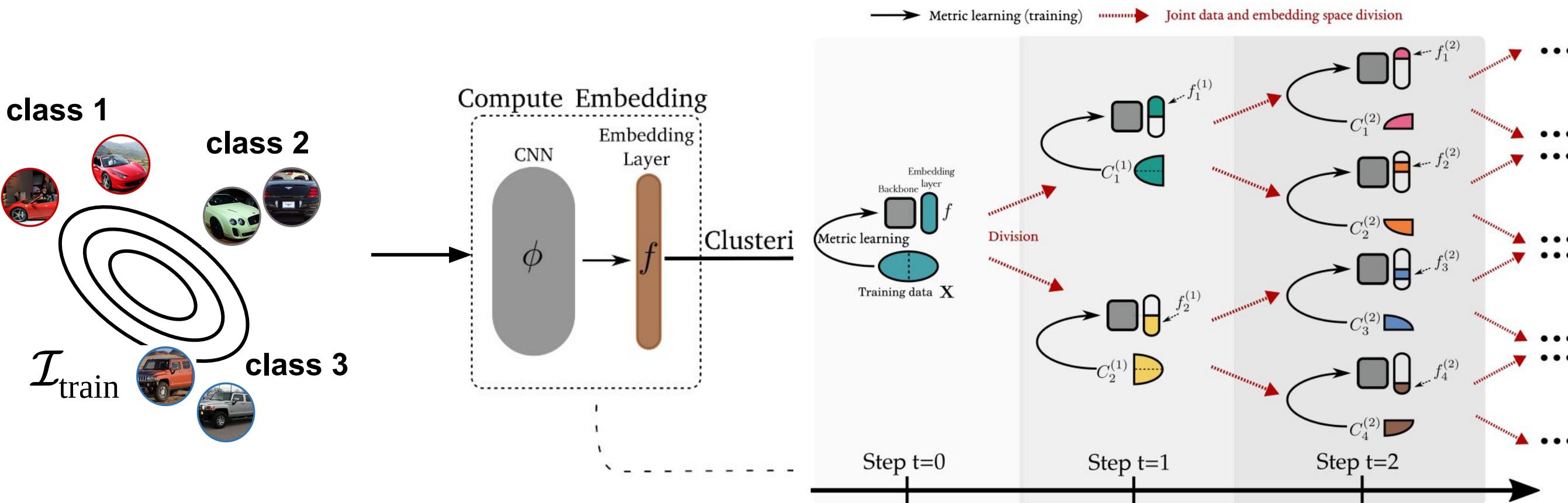
Ensemble-based DML

- Typically training distribution is **multimodal**
- **Explicit** data specialization of learners following **Divide & Conquer strategy**⁹



Ensemble-based DML

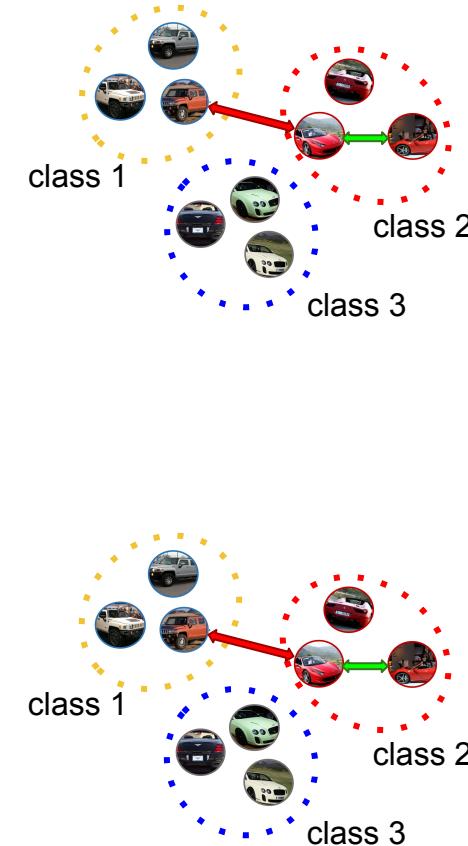
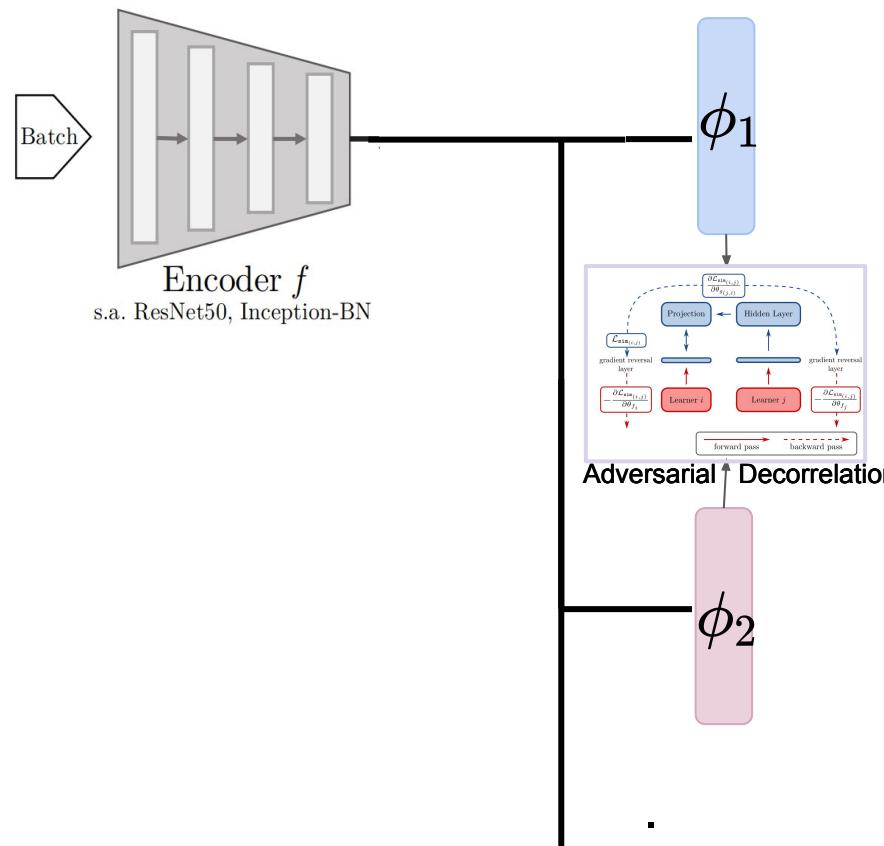
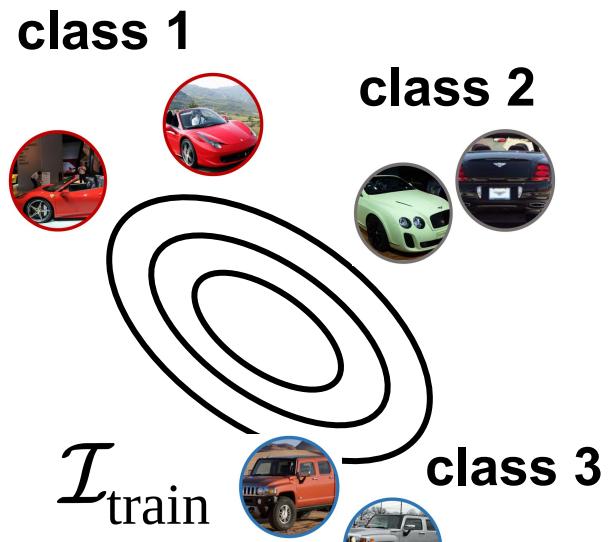
- Typically training distribution is **multimodal**
- **Explicit** data specialization of learners following **Hierarchical Divide & Conquer strategy¹⁰**



Sequential Splitting: Continuous specialization by splitting data and adding learners during optimization

¹⁰Sanakoyeu, Ma, Tschernetzski, Ommer 2021

Ensemble-based DML



Target similar object features!

$$t = \{I_a, I_p, I_n\}$$

$$\frac{y_a = y_p \neq y_n}{\mathcal{L}_{\text{disc}}}$$

Class-discriminative features

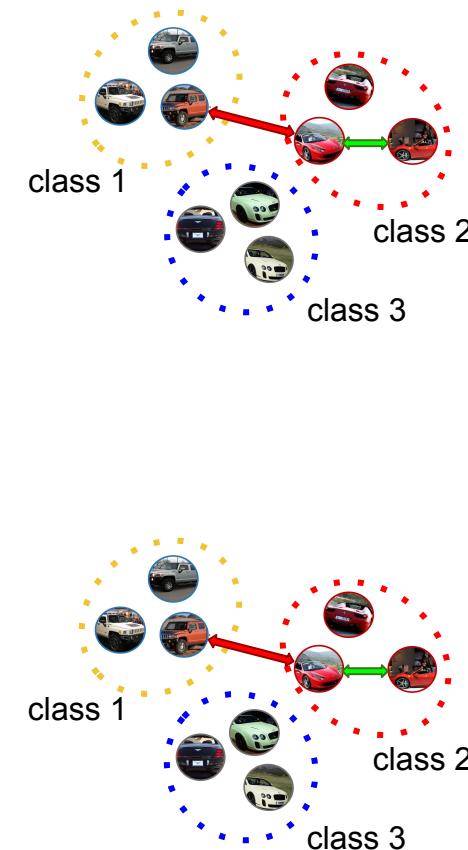
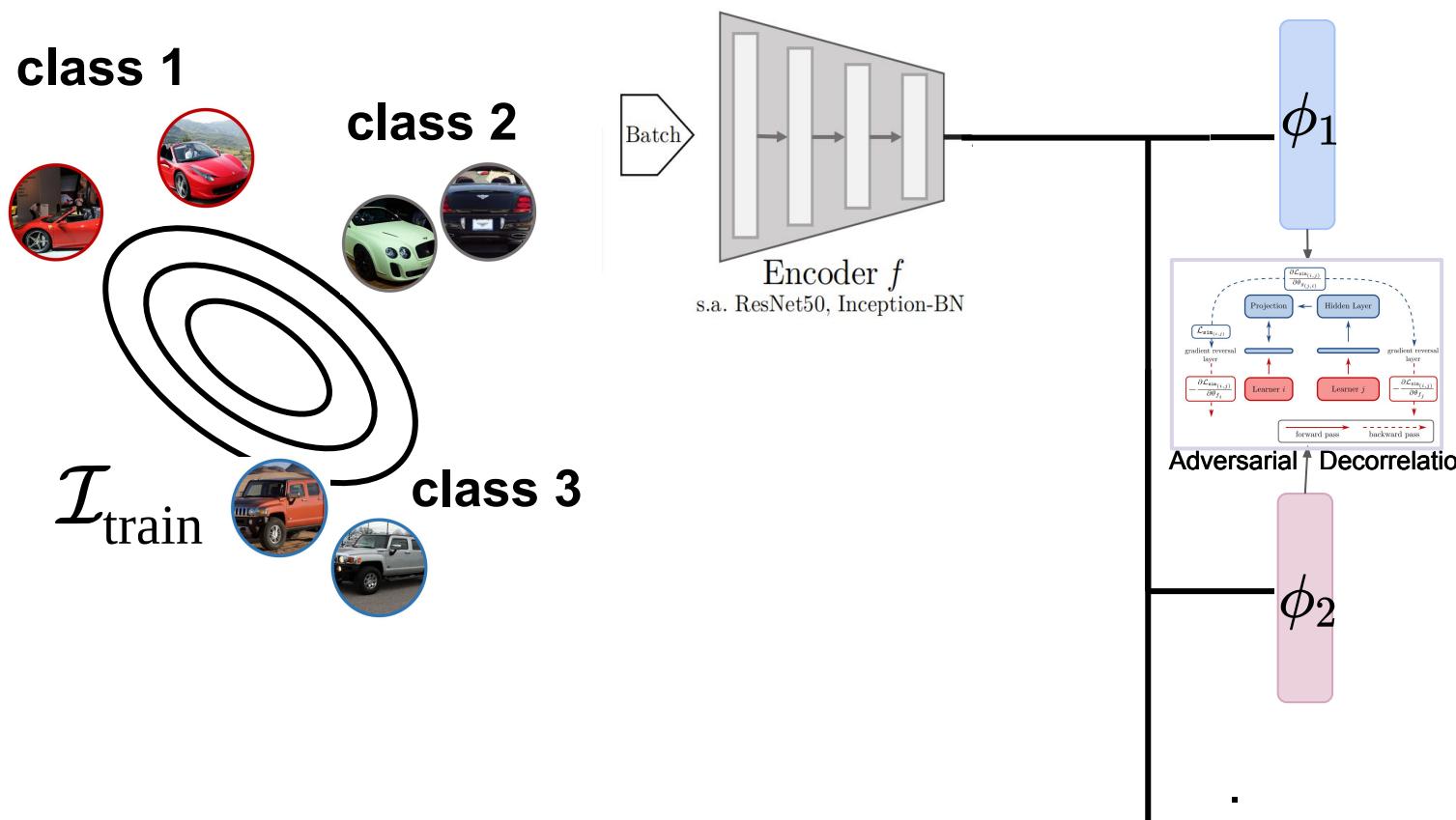
- capture features **separating between classes**
- aggregate features into ‘classes’
- very specialized
- e.g. “Ferrari” vs. “Hummer”

$$t = \{I_a, I_p, I_n\}$$

$$\frac{y_a = y_p \neq y_n}{\mathcal{L}_{\text{disc}}}$$

Ensemble-based DML

- **Assumption:** Different features improve robustness to OOD data (new classes, etc.)



$$t = \{I_a, I_p, I_n\} \quad \mathcal{L}_{\text{disc}}$$

$$\frac{y_a = y_p \neq y_n}{}$$

Class-discriminative features

- capture features **separating between classes**
- aggregate features into ‘classes’
- very specialized
- e.g. “Ferrari” vs. “Hummer”

$$t = \{I_a, I_p, I_n\} \quad \boxed{\mathcal{L}_{\text{disc}}}$$

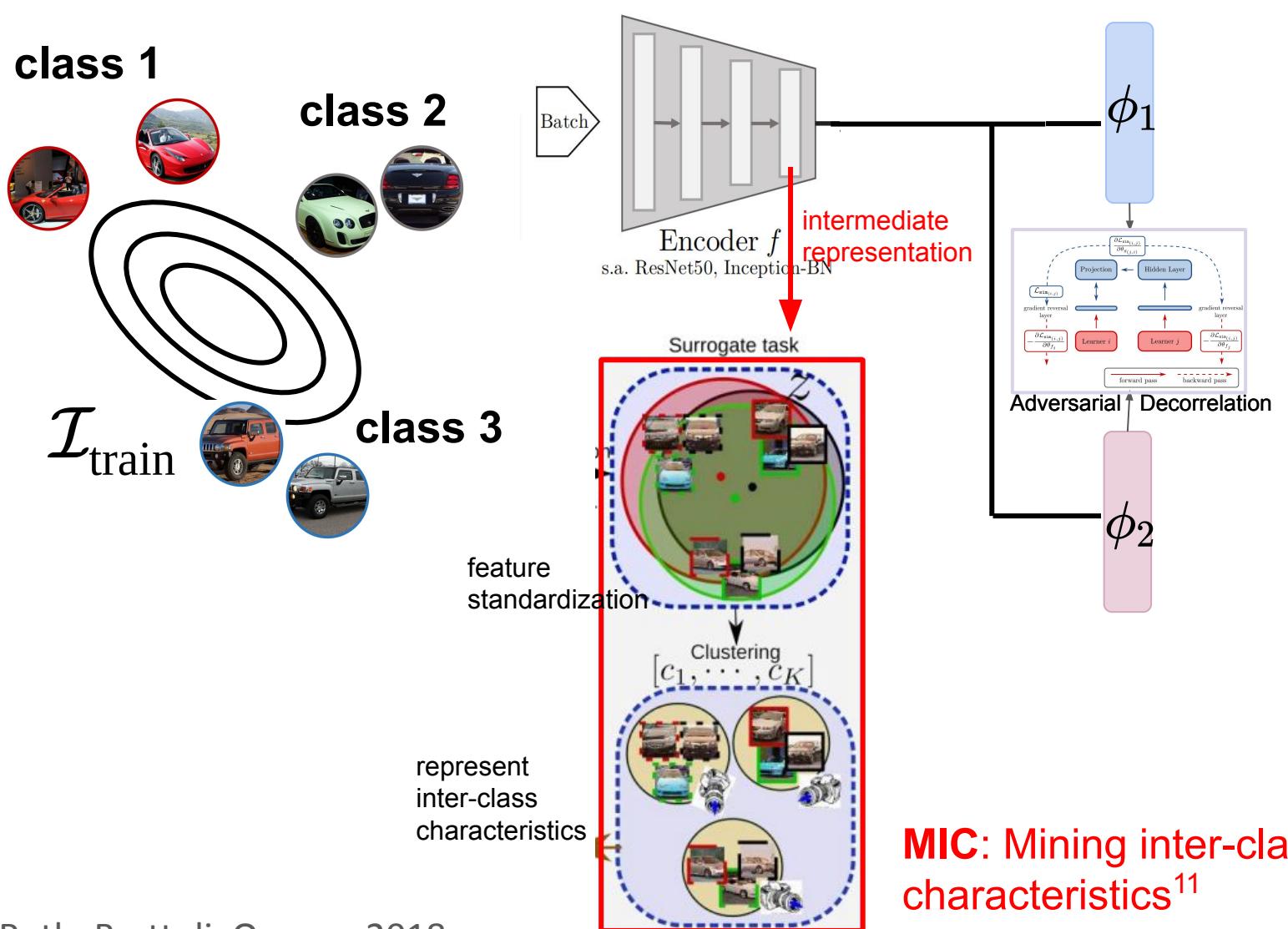
$$\frac{y_a = y_p \neq y_n}{}$$

Learn about different/more general features?

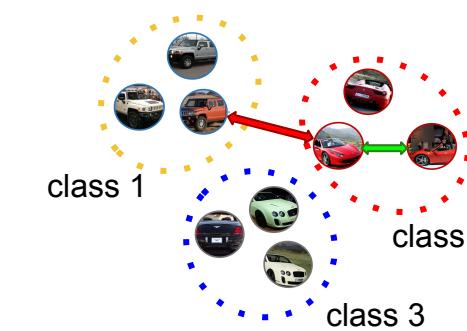
- color
- shape
- viewpoint
-

Ensemble-based DML

- **Assumption:** Different features improve robustness to OOD data (new classes, etc.)
- Only class-labels available: use **unsupervised learning**



MIC: Mining inter-class characteristics¹¹



$$t = \{I_a, I_p, I_n\} \quad \mathcal{L}_{\text{disc}}$$

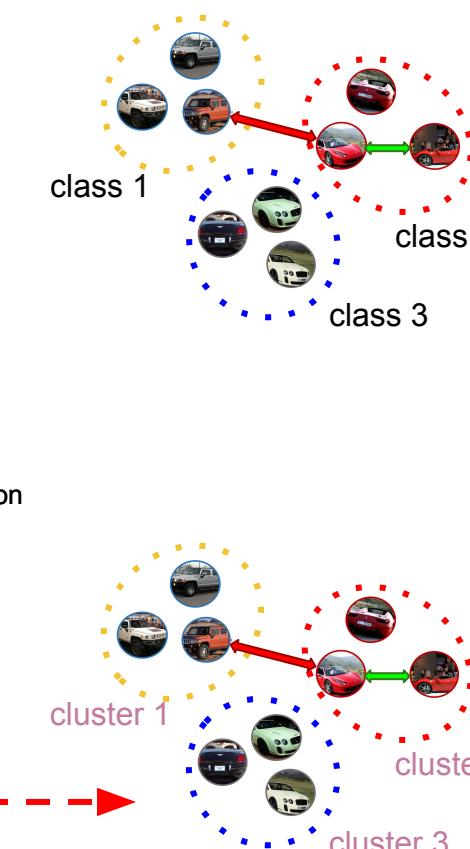
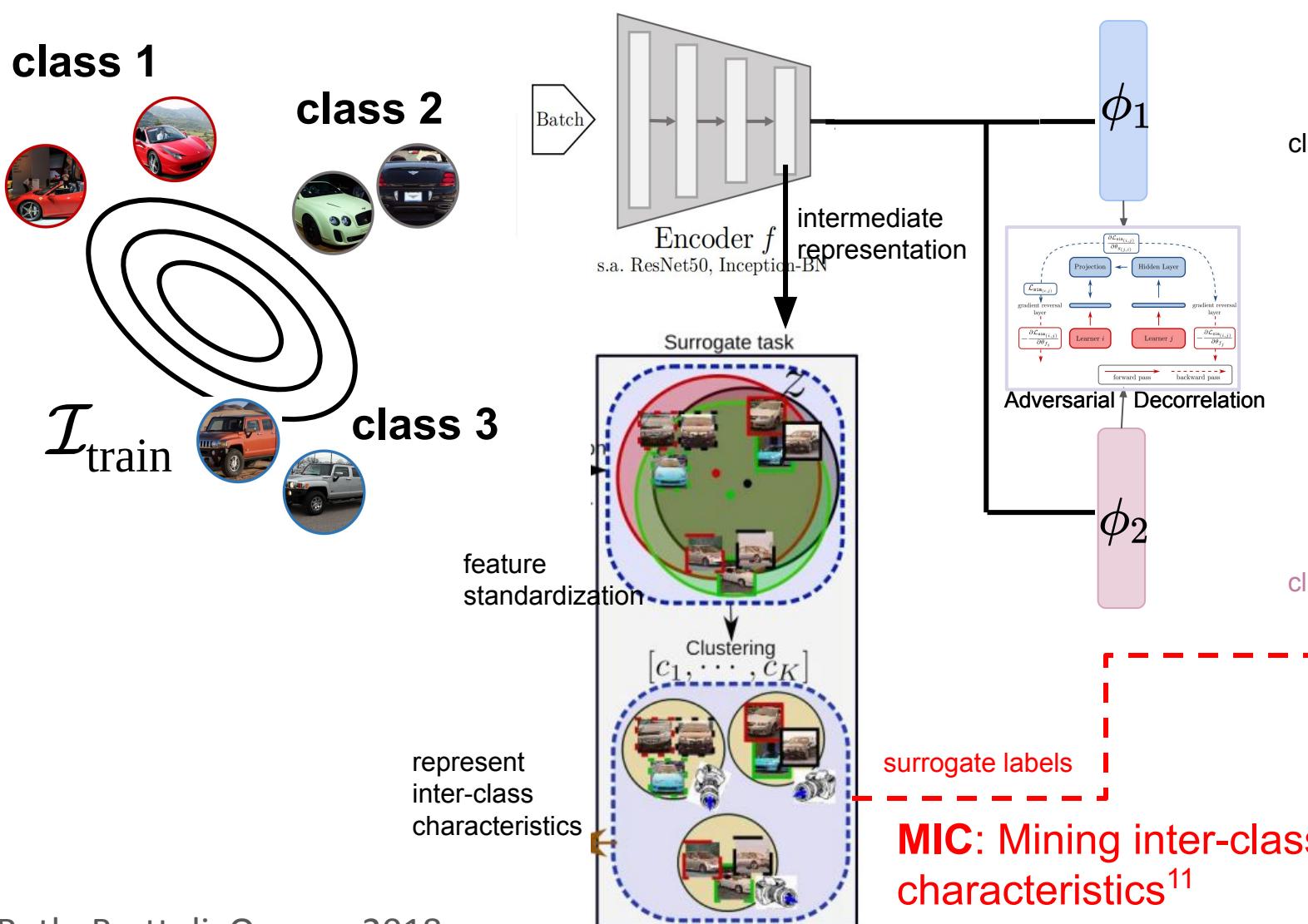
$$y_a = y_p \neq y_n$$

Class-discriminative features

- capture features **separating between classes**
- aggregate features into ‘classes’
- very specialized
- e.g. “Ferrari” vs. “Hummer”

Ensemble-based DML

- **Assumption:** Different features improve robustness to OOD data (new classes, etc.)
- Only class-labels available: use **unsupervised learning**



$$t = \{I_a, I_p, I_n\} \quad \mathcal{L}_{\text{disc}}$$

Class-discriminative features

- capture features **separating between classes**
- aggregate features into ‘classes’
- very specialized
- e.g. “Ferrari” vs. “Hummer”

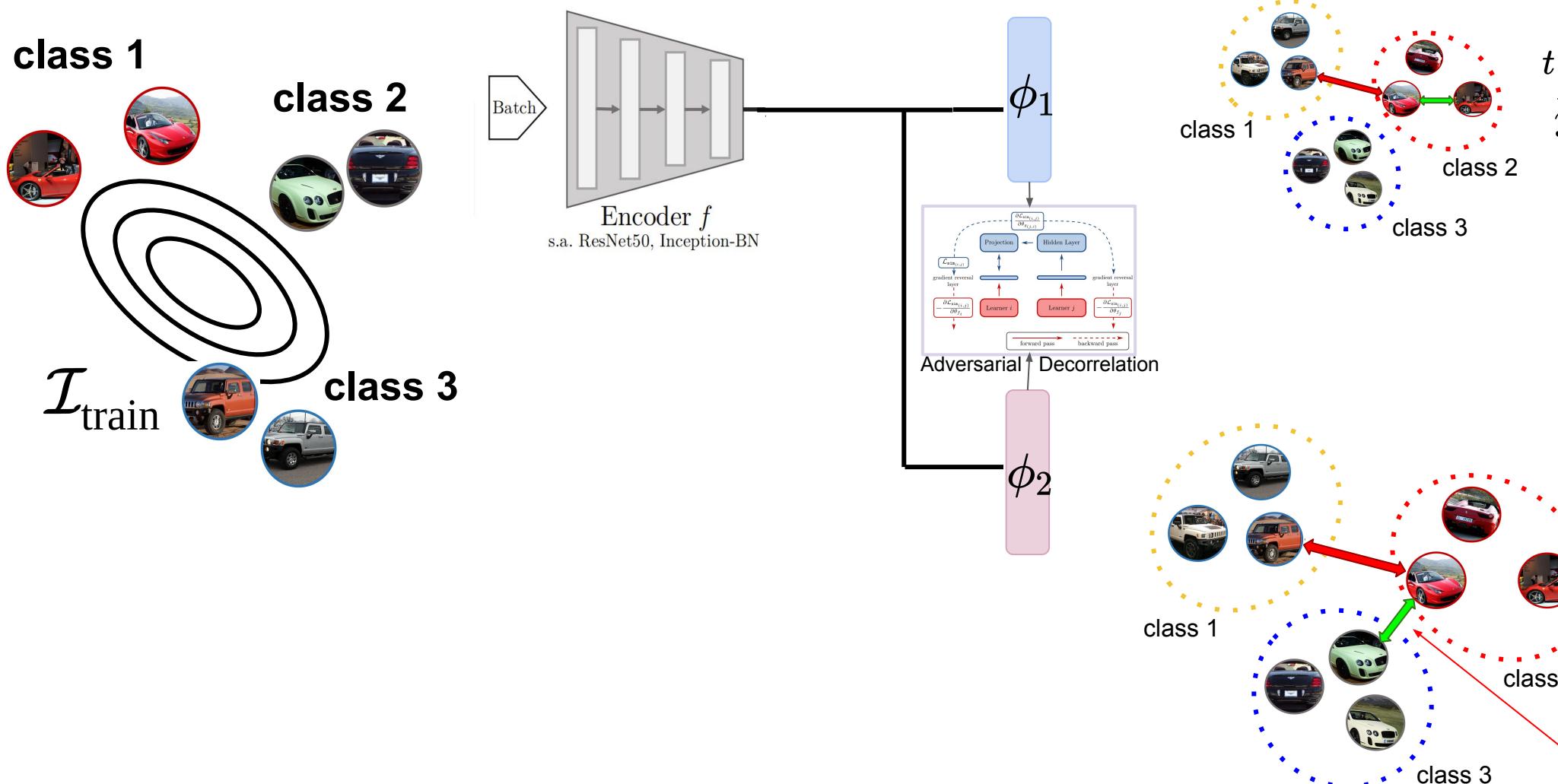
$$t = \{I_a, I_p, I_n\} \quad \mathcal{L}_{\text{surrogate}}$$

Inter-class characteristics¹¹

- anchors and positives **from same cluster** (different class labels!)
- represent inter-class characteristics, e.g. viewpoint, color, etc.

Ensemble-based DML

- **Assumption:** Different features improve robustness to OOD data (new classes, etc.)
- Only class-labels available: use **unsupervised learning**



$$t = \{I_a, I_p, I_n\} \xrightarrow{y_a = y_p \neq y_n} \mathcal{L}_{\text{disc}}$$

Class-discriminative features

- capture features **separating between classes**
- aggregate features into ‘classes’
- very specialized
- e.g. “Ferrari” vs. “Hummer”

Class-shared features¹²

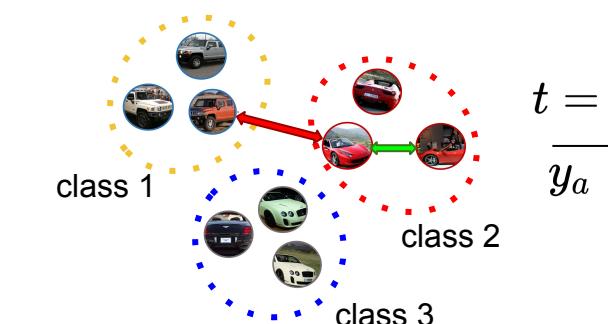
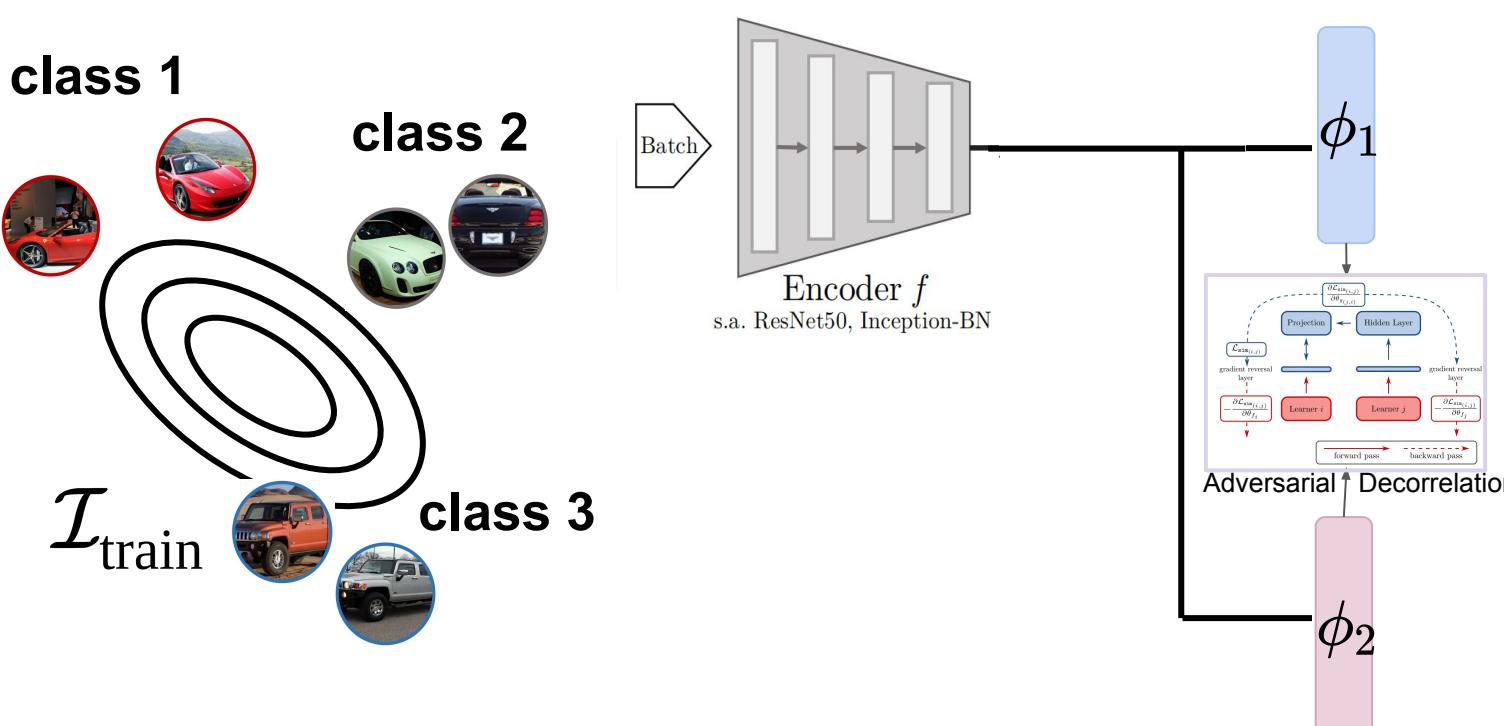
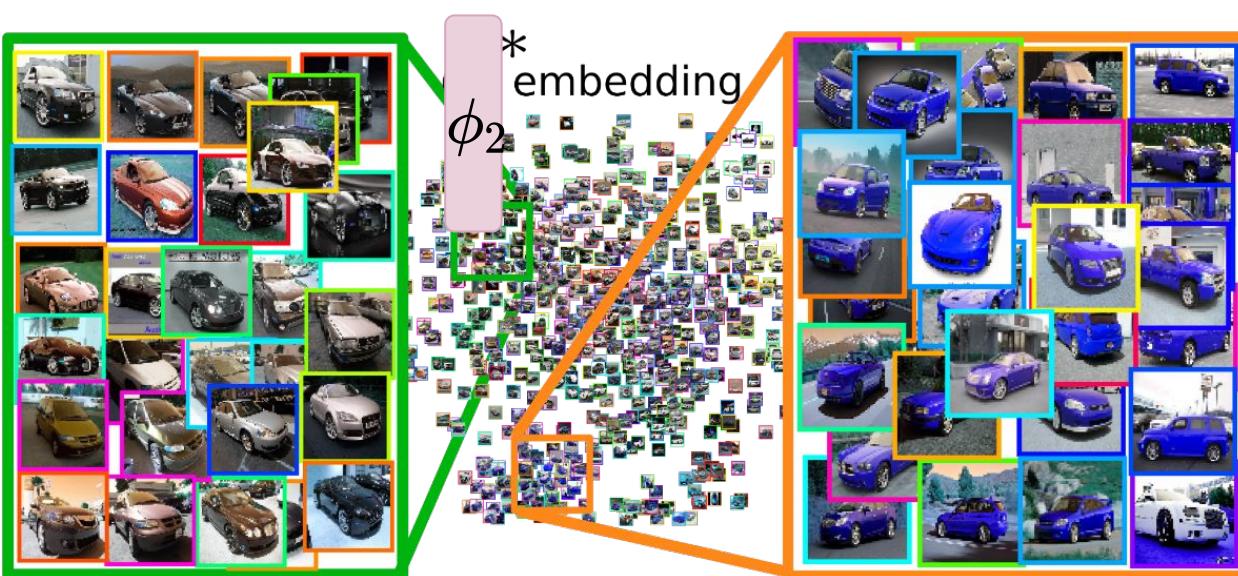
- anchors and positives **from different classes**
- represent **shared, general object** characteristics

$$t = \{I_a, I_p, I_n\} \xrightarrow{y_a \neq y_p \neq y_n} \mathcal{L}_{\text{shared}}$$

- Explicitly sample anchor and positive from **DIFFERENT classes!**
- Can be achieved by simply changing triplet sampling

Ensemble-based DML

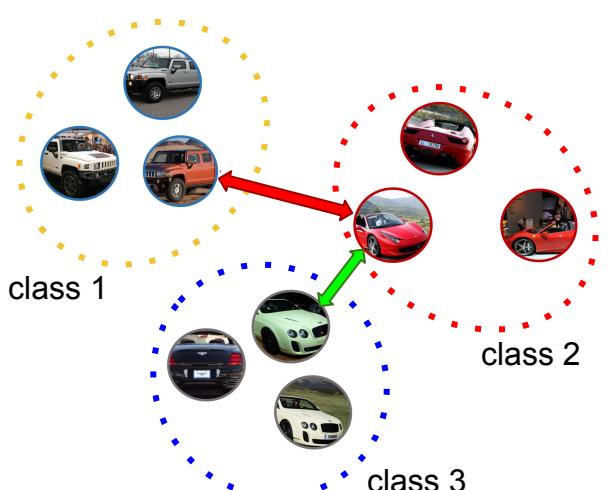
- **Assumption:** Different features improve robustness to OOD data (new classes, etc.)
- Only class-labels available: use **unsupervised learning**



$$t = \{I_a, I_p, I_n\} \xrightarrow{y_a = y_p \neq y_n} \mathcal{L}_{\text{disc}}$$

Class-discriminative features

- capture features **separating between classes**
- aggregate features into ‘classes’
- very specialized
- e.g. “Ferrari” vs. “Hummer”



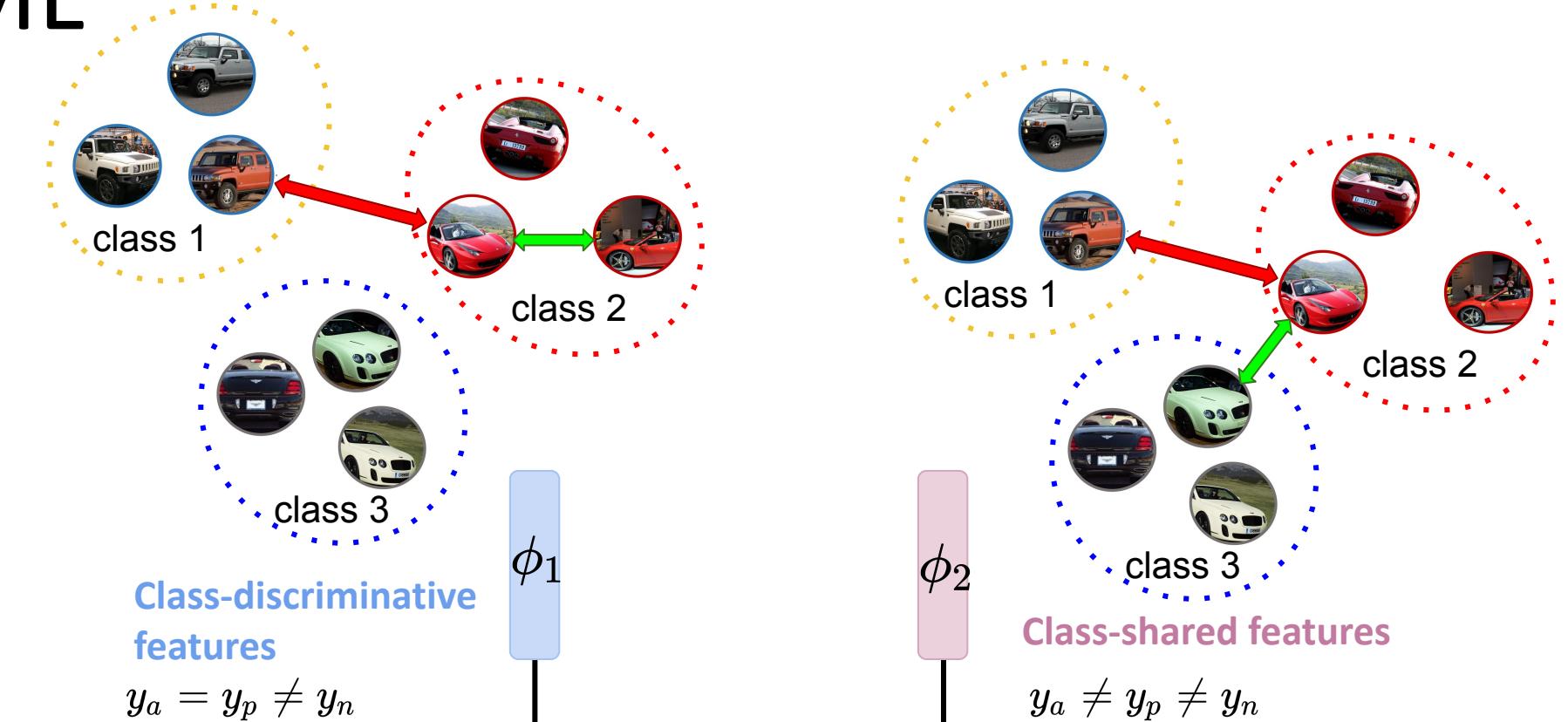
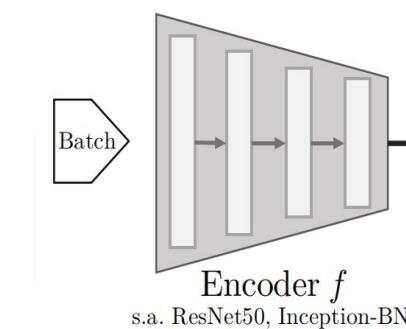
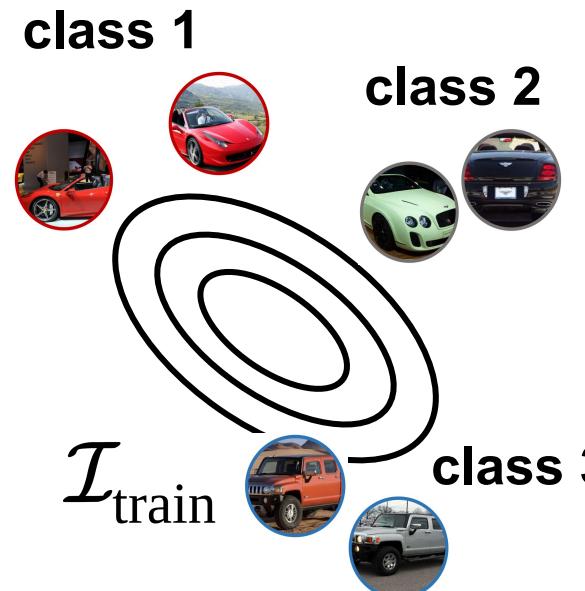
Class-shared features¹²

- anchors and positives **from different classes**
- represent **shared, general object** characteristics

$$t = \{I_a, I_p, I_n\} \xrightarrow{y_a \neq y_p \neq y_n} \mathcal{L}_{\text{shared}}$$

Ensemble-based DML

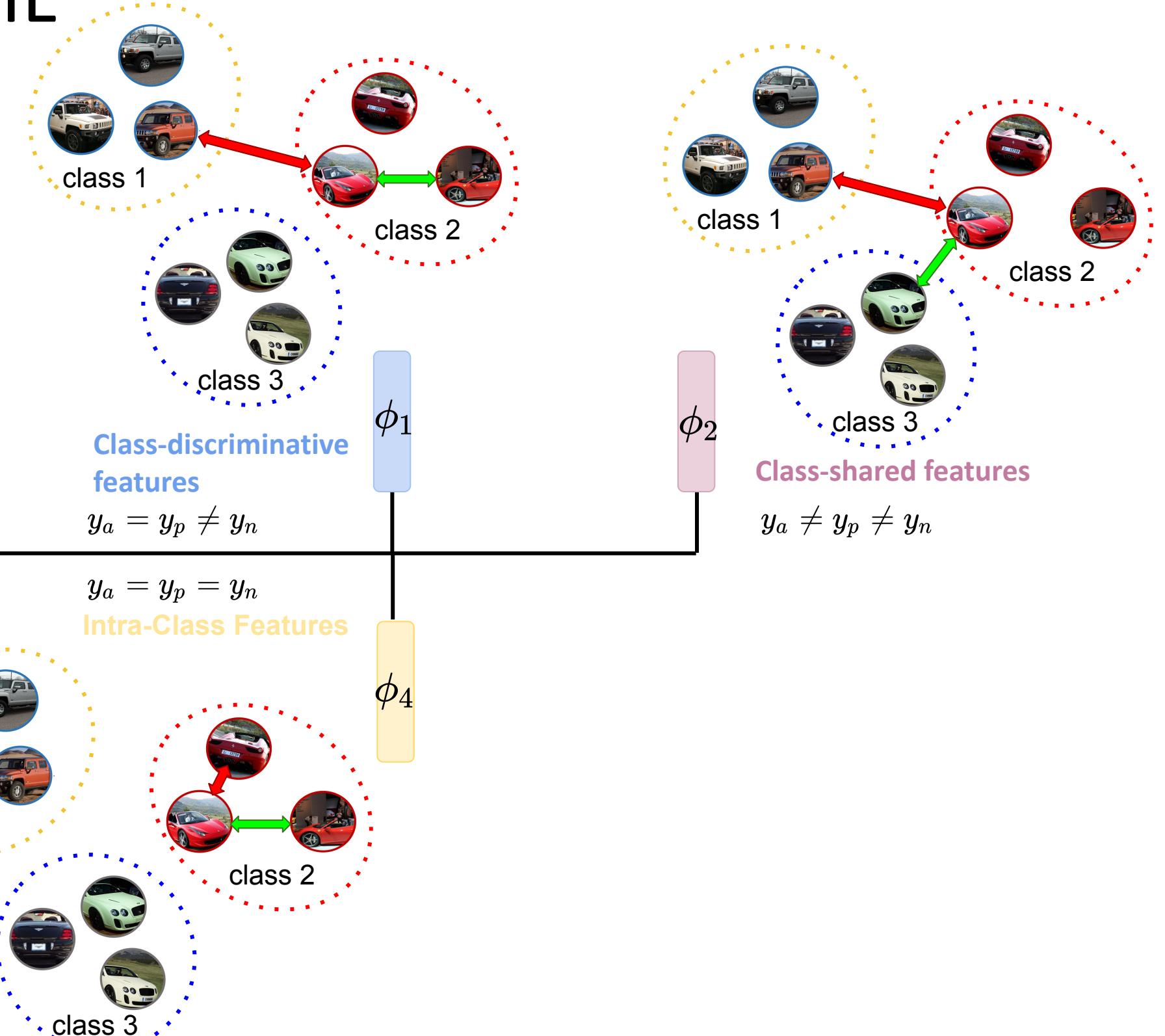
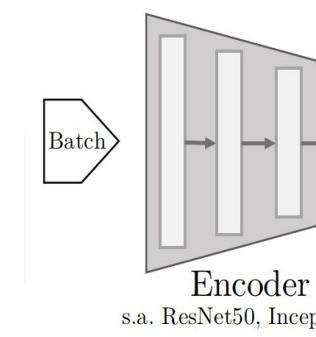
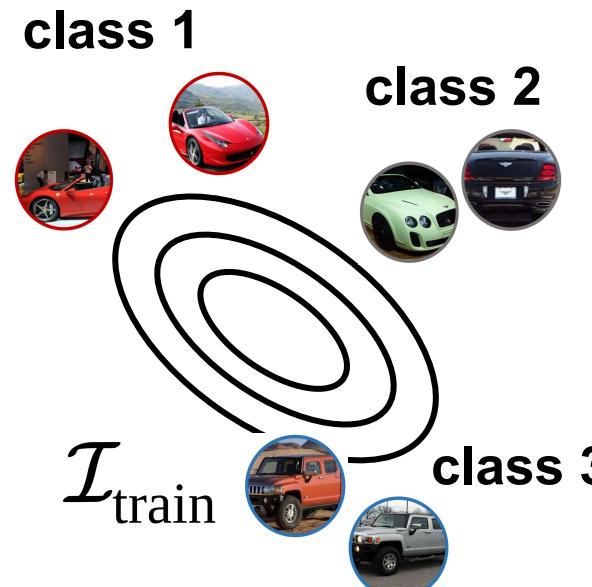
- DiVA¹³ :Diverse Visual Feature Aggregation:
Multi-task ensemble with each learner focussing on complementary features



¹³Milbich, Roth, Bharadhwaj, Sinha, Y.Bengio, Ommer, Cohen 2021

Ensemble-based DML

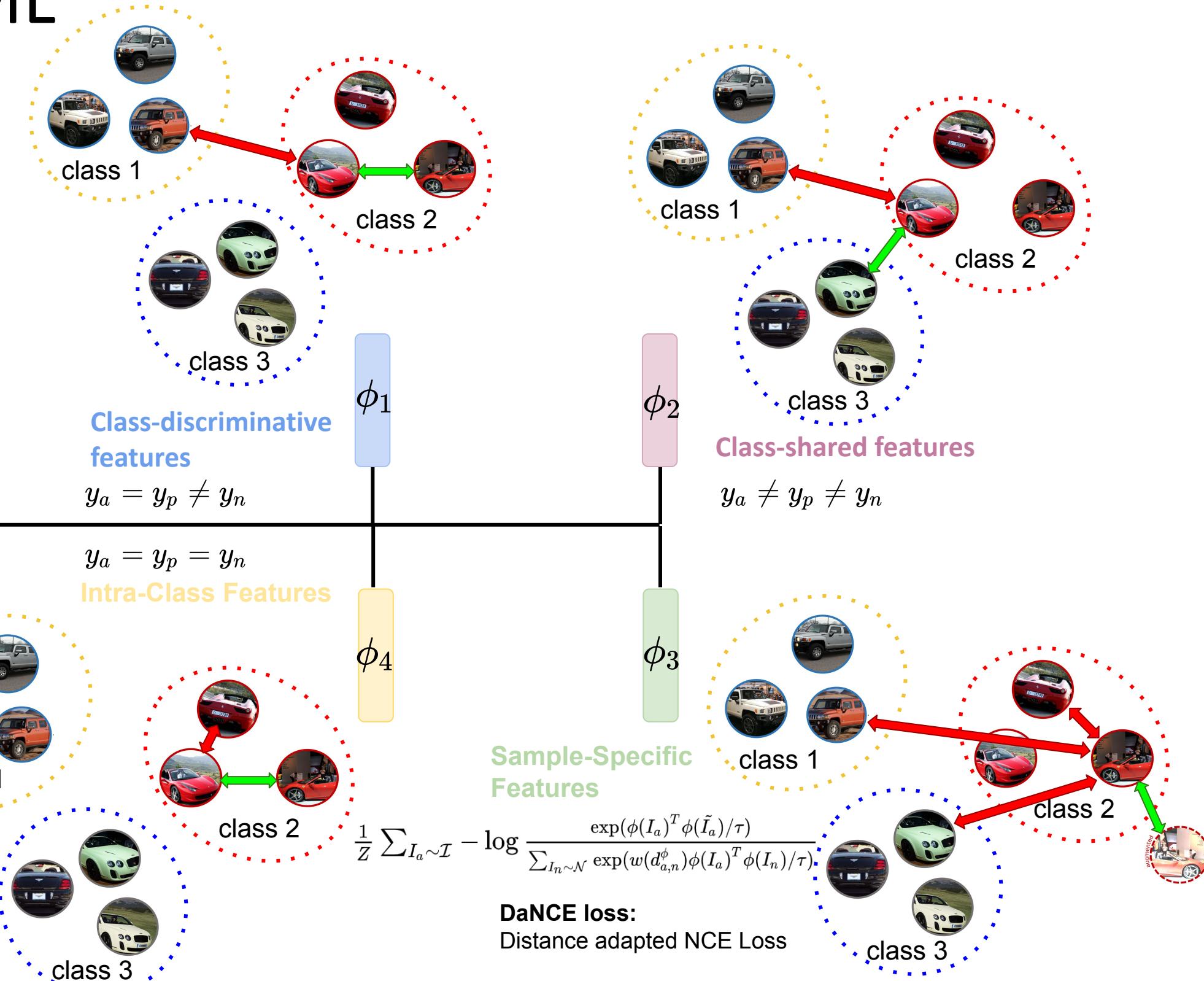
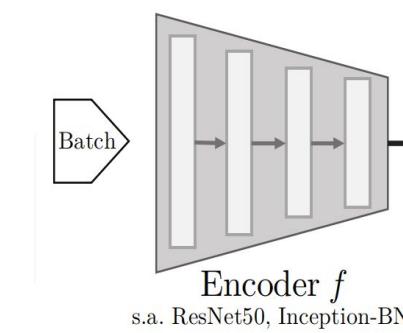
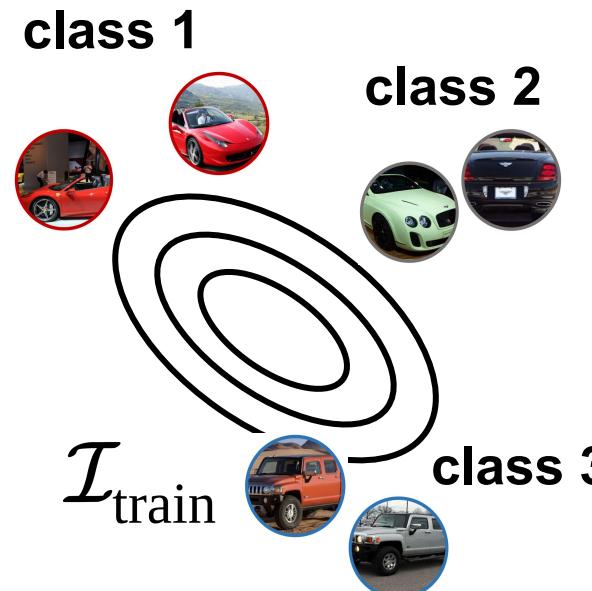
- DiVA¹³ :Diverse Visual Feature Aggregation:
Multi-task ensemble with each learner focussing on complementary features



¹³Milbich, Roth, Bharadhwaj, Sinha, Y.Bengio, Ommer, Cohen 2021

Ensemble-based DML

- DiVA¹³ :Diverse Visual Feature Aggregation:
Multi-task ensemble with each learner focussing on complementary features

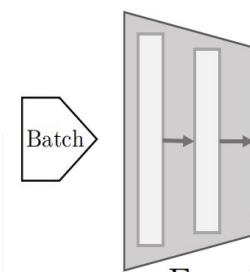
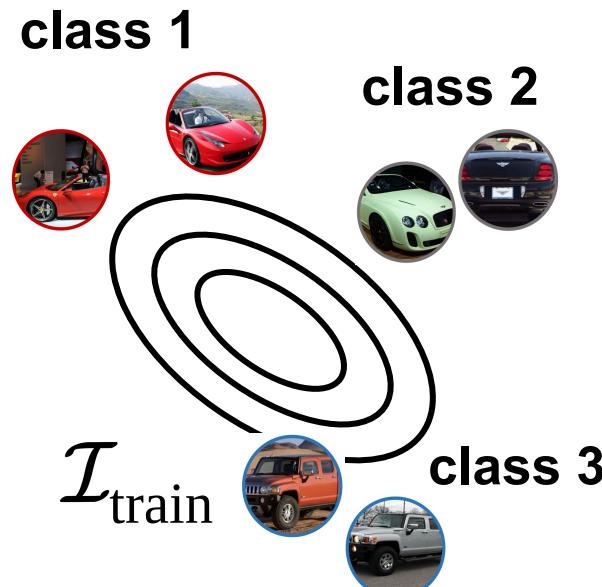


¹³Milbich, Roth, Bharadhwaj, Sinha, Y.Bengio, Ommer, Cohen 2021

Ensemble-based DML

- DiVA¹³ :Diverse Visual Feature Aggregation:

Multi-task ensemble with each learner focussing on complementary features



Encoder f
s.a. ResNet50, Inception-BN

Dataset →	CUB200-2011[60]				CARS196[32]				SOP[42]			
Approach ↓	Dim	R@1	R@2	NMI	R@1	R@2	NMI	R@1	R@10	NMI		
Margin[64] (orig, R50)	128	63.6	74.4	69.0	79.6	86.5	69.1	72.7	86.2	90.7		
Margin[64] (ours, IBN)	512	63.6	74.7	68.3	79.4	86.6	66.2	76.6	89.2	89.8		
DiVA (IBN, D & Da)	512	64.5	76.0	68.8	80.4	87.7	67.2	77.0	89.4	90.1		
DiVA (IBN, D & S)	512	65.1	76.4	69.0	81.5	88.3	66.8	77.2	89.6	90.0		
DiVA (IBN, D & I)	512	64.9	75.8	68.4	80.6	87.9	67.4	76.9	89.4	89.9		
DiVA (IBN, D & Da & I)	510	65.3	76.5	68.3	82.2	89.1	67.8	75.8	89.0	89.8		
DiVA (IBN, D & S & I)	510	65.5	76.4	68.4	82.1	89.4	67.2	77.0	89.3	89.7	$\leq y_p \neq y_n$	
DiVA (IBN, D & Da & S)	510	65.9	76.7	68.9	82.6	89.6	68.0	77.4	89.6	90.1		
DiVA (IBN, D & Da & S & I)	512	66.4	77.2	69.6	83.1	90.0	68.1	77.5	90.3	90.1		

2 tasks
3 tasks
4 tasks

Intra-Class Features

ϕ_4

Sample-Specific
Features

ϕ_3

$$\frac{1}{Z} \sum_{I_a \sim \mathcal{I}} -\log \frac{\exp(\phi(I_a)^T \phi(\tilde{I}_a) / \tau)}{\sum_{I_n \sim \mathcal{N}} \exp(w(d_{a,n}^\phi) \phi(I_a)^T \phi(I_n) / \tau)}$$

DaNCE loss:
Distance adapted NCE Loss

class 1

class 2

class 3

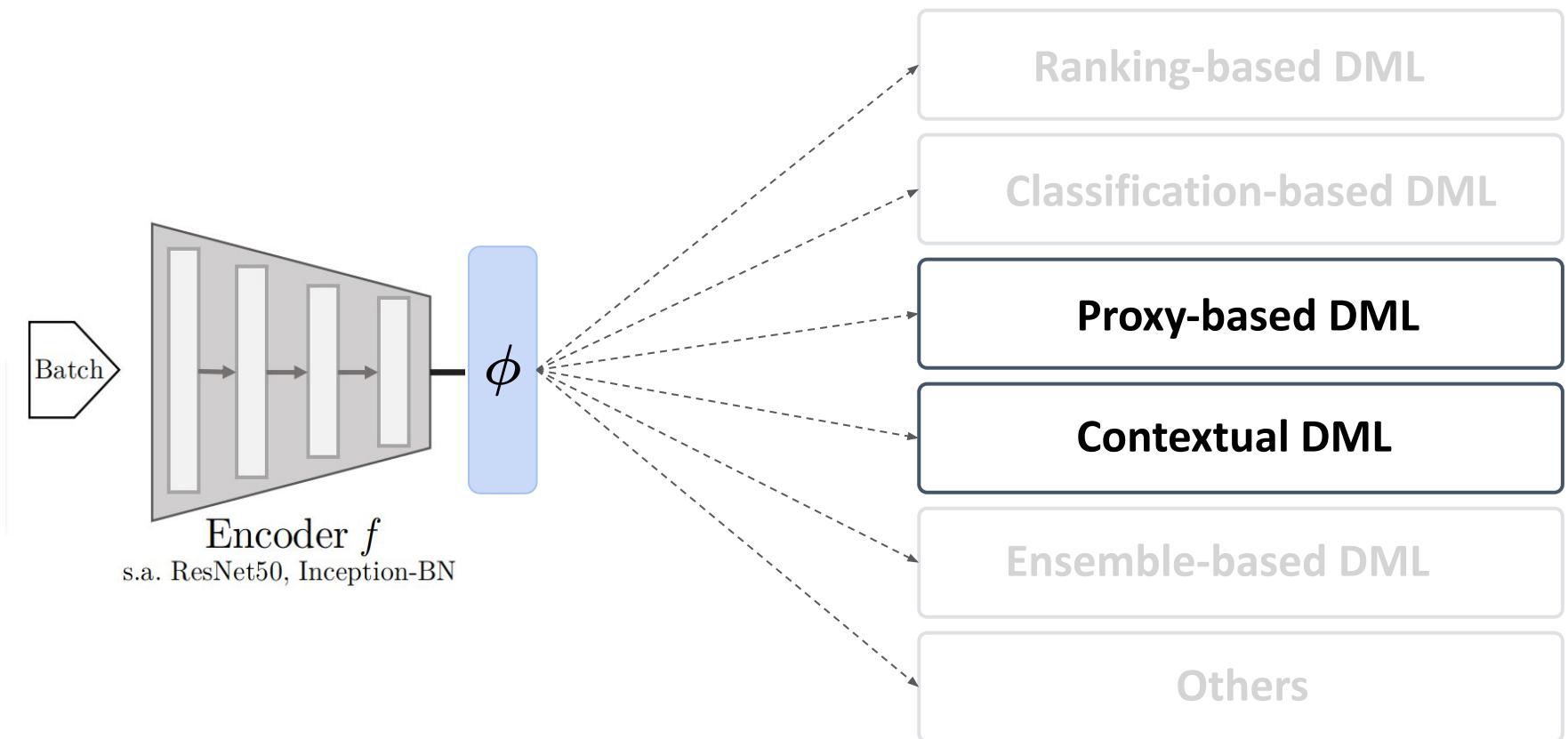
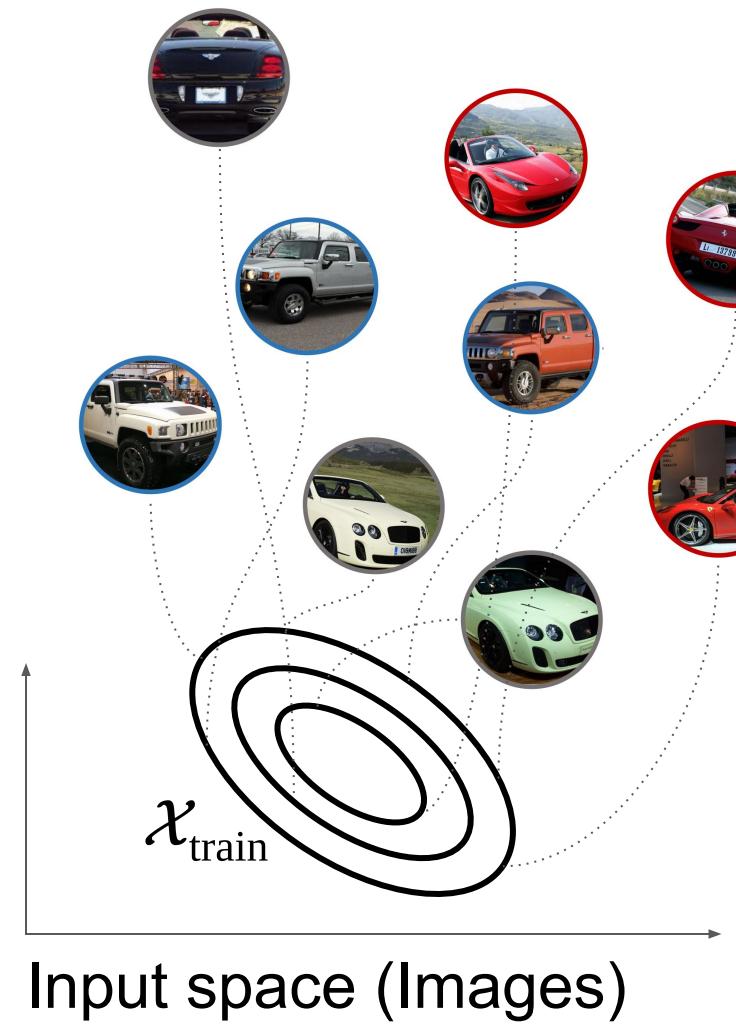
class 2

class 3

augmented

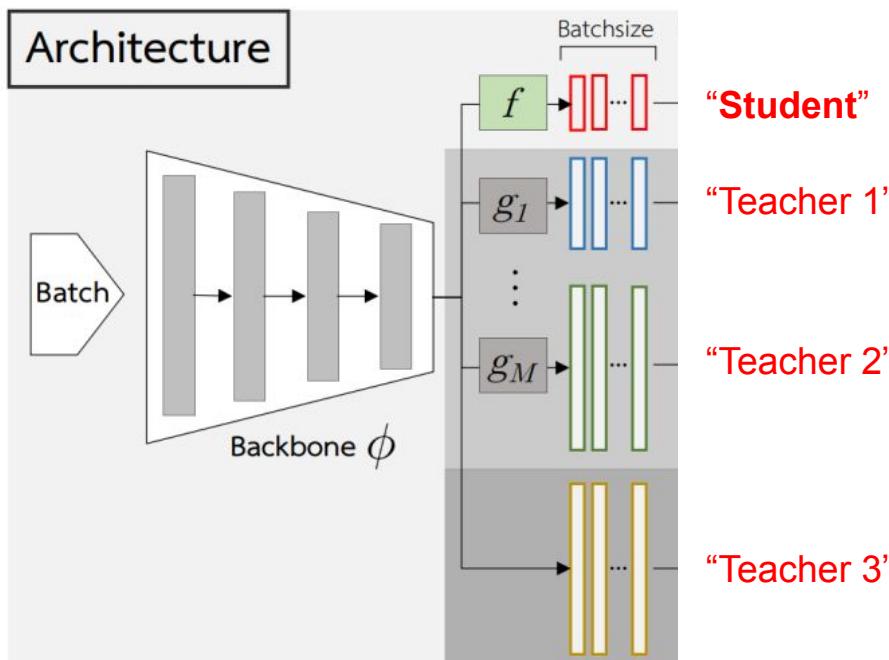
Deep Metric Learning (DML)

Next talk:



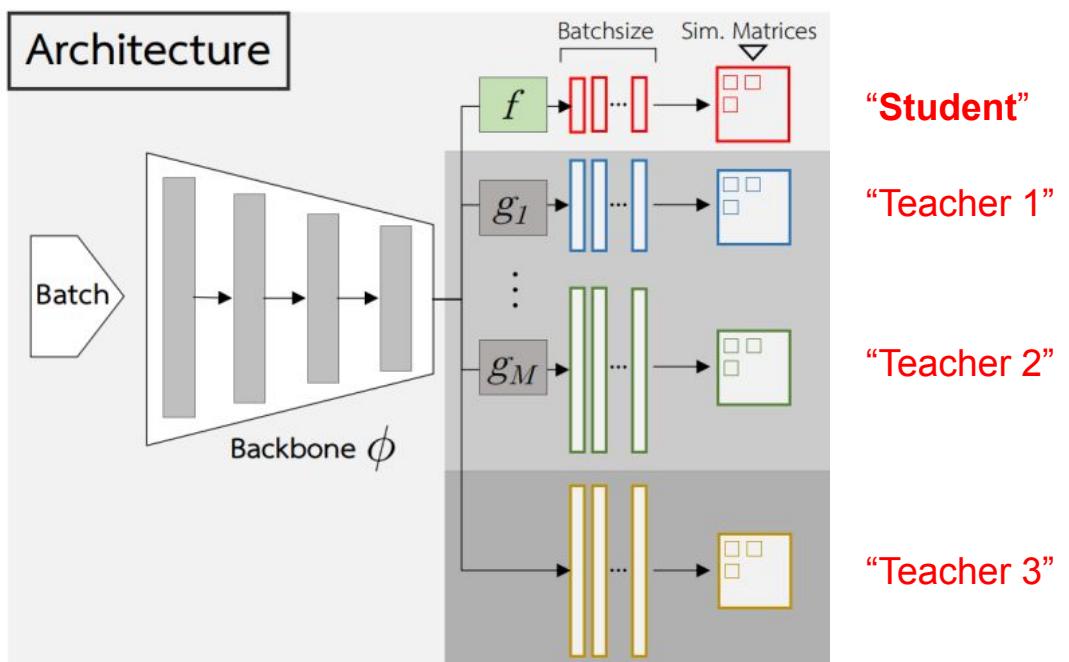
Ensemble-based DML

- S2SD¹⁴: Learn ensemble of teachers and distill their ‘knowledge’ into student



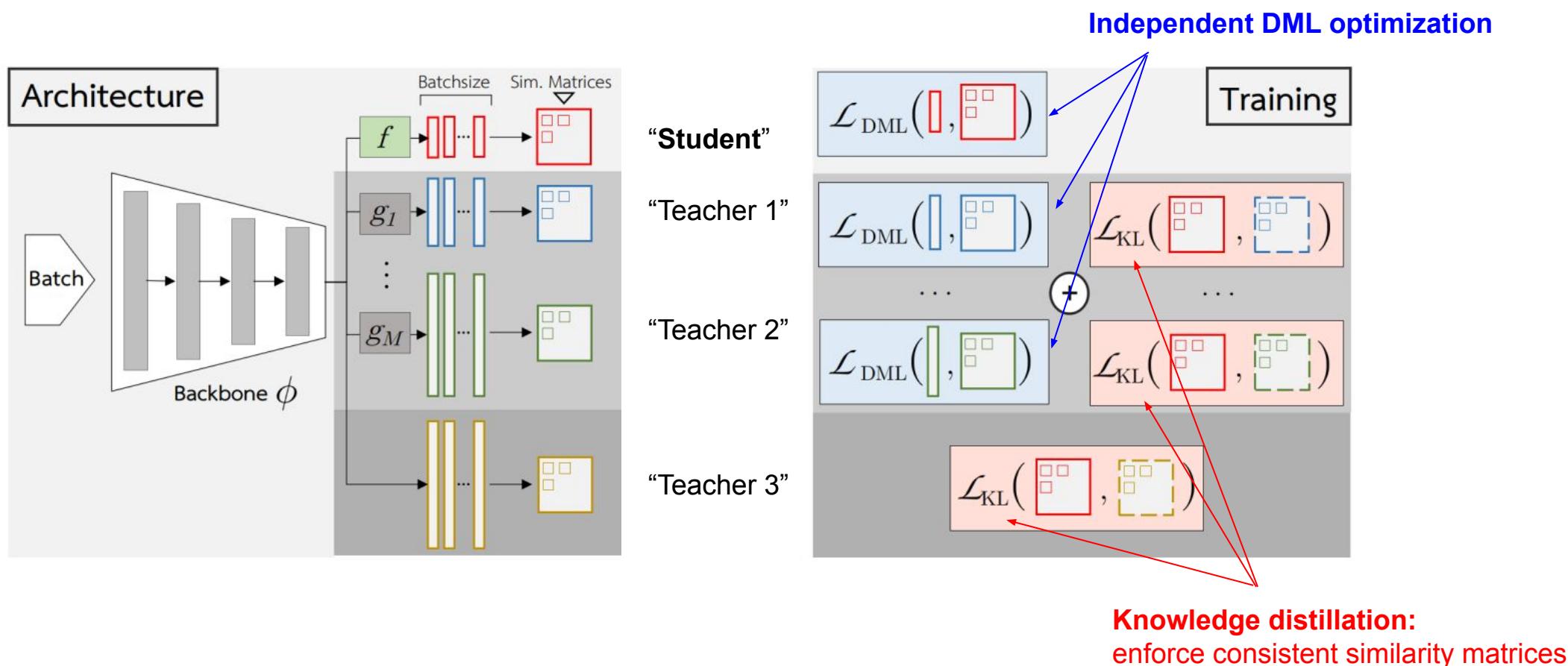
Ensemble-based DML

- S2SD¹⁴: Learn ensemble of teachers and distill their ‘knowledge’ into student



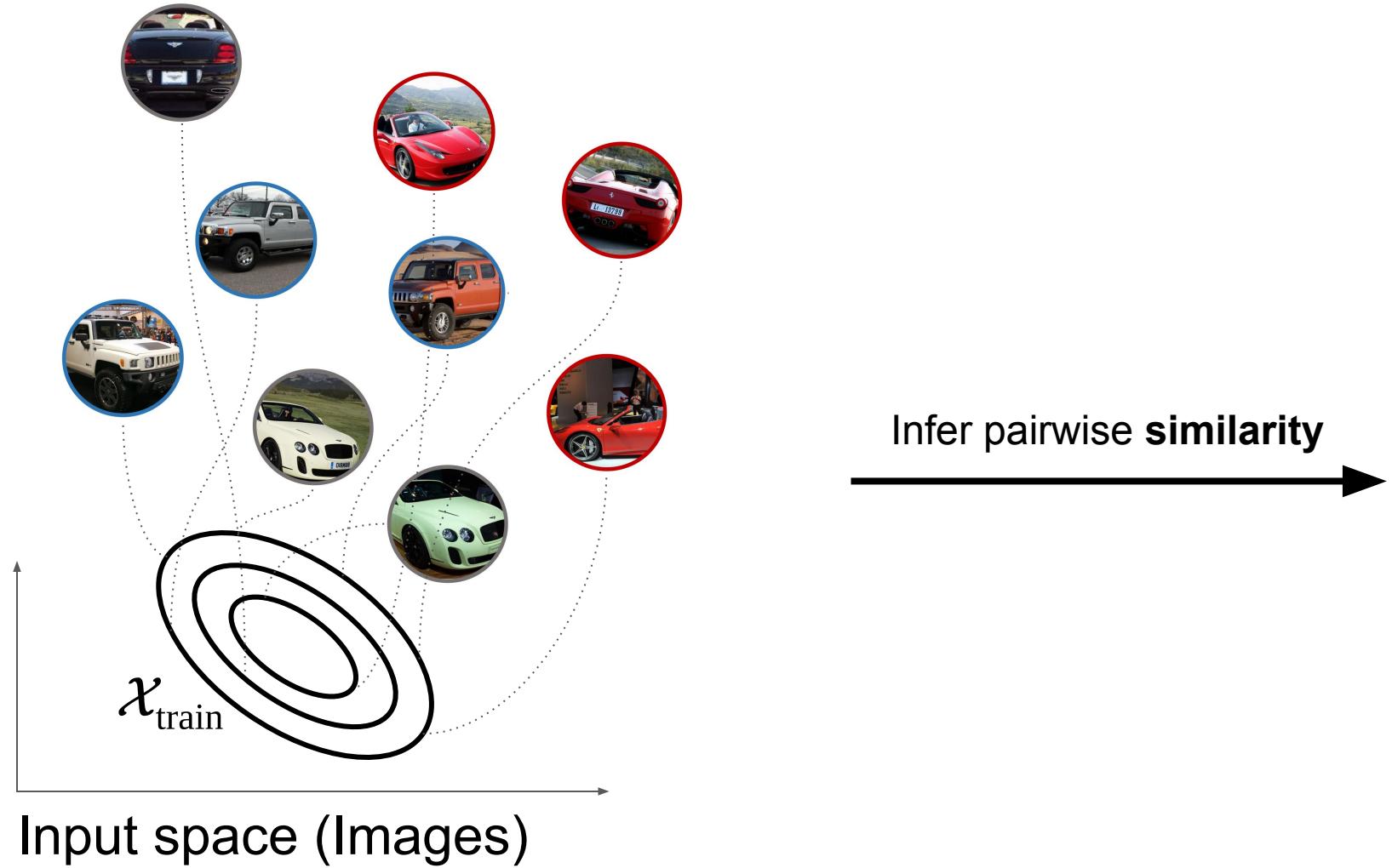
Ensemble-based DML

- S2SD¹⁴: Learn **ensemble of teachers** and distill their ‘knowledge’ into student
- use different embedding dimensionalities for teachers to increase robustness of student



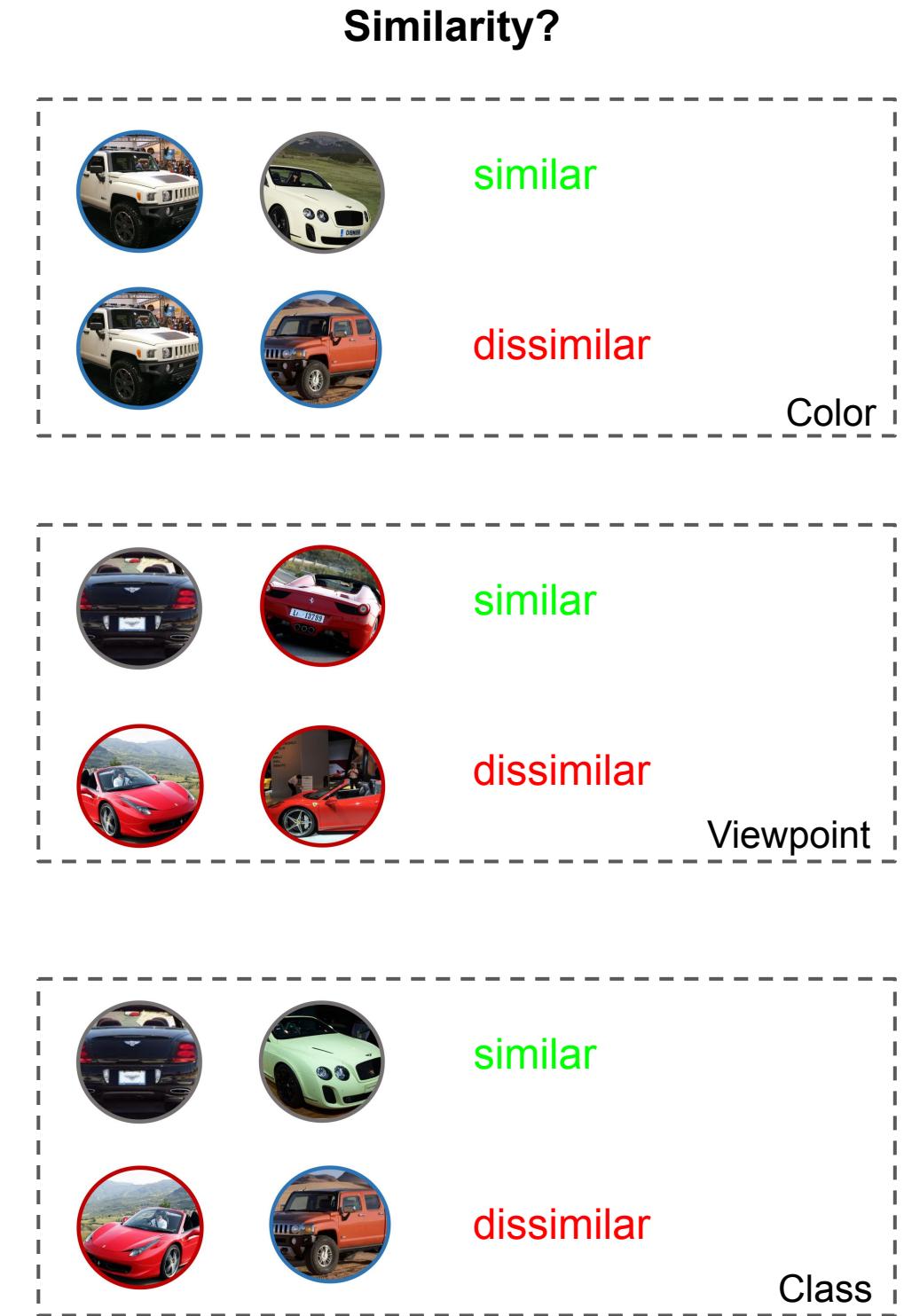
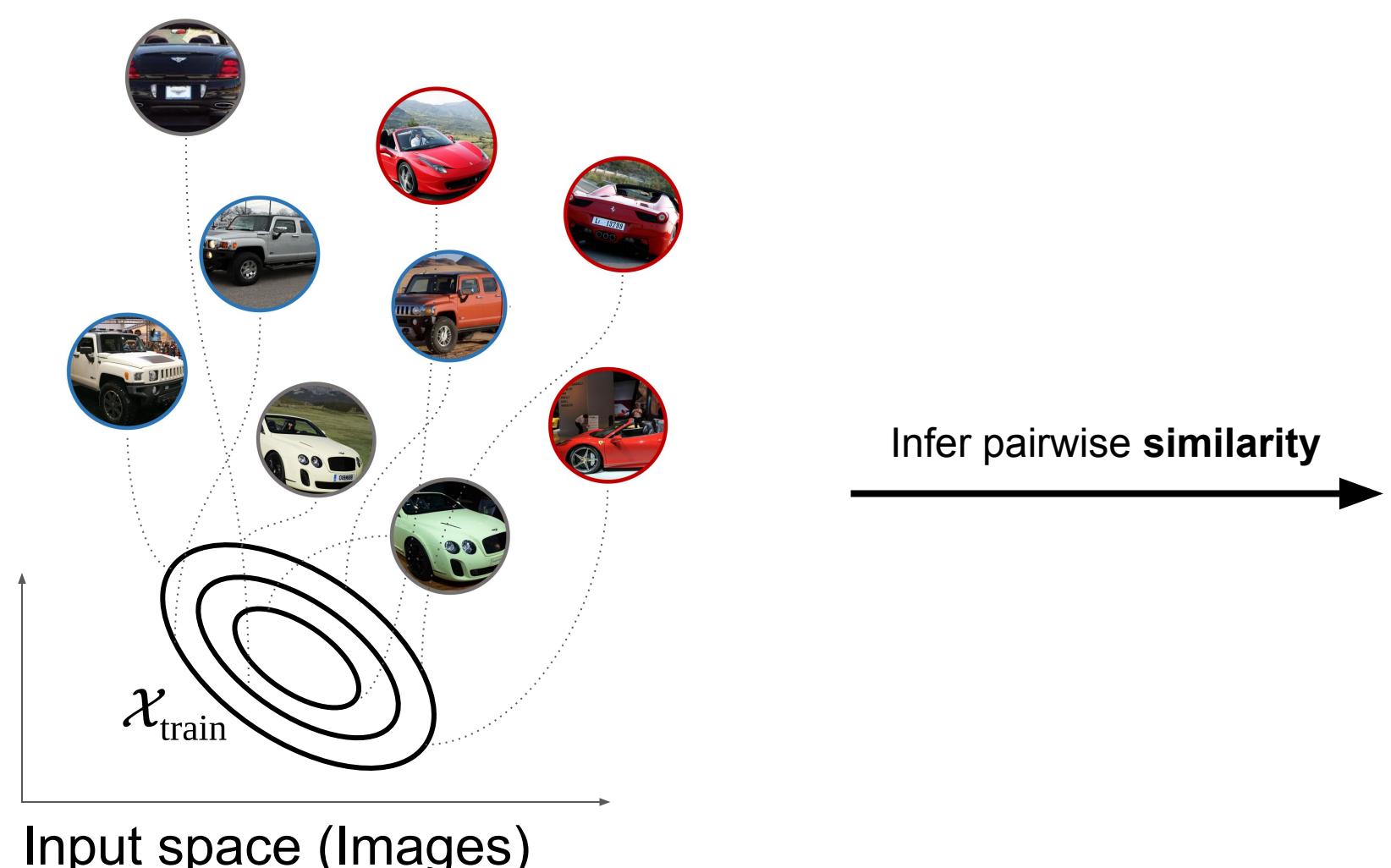
Visual Similarity Learning

Given data distribution $\mathcal{X}_{\text{train}}$, infer pairwise **semantic similarity**.



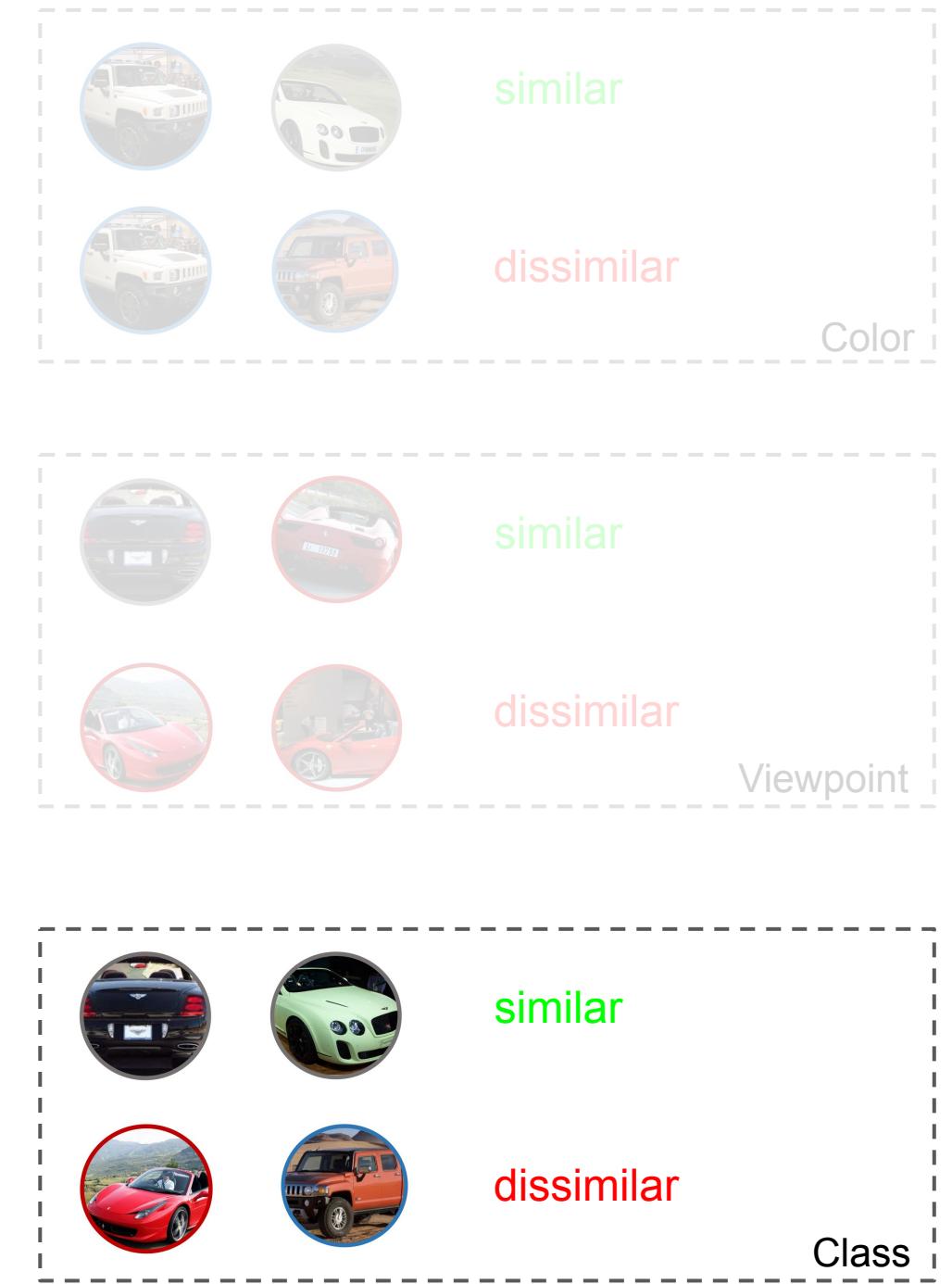
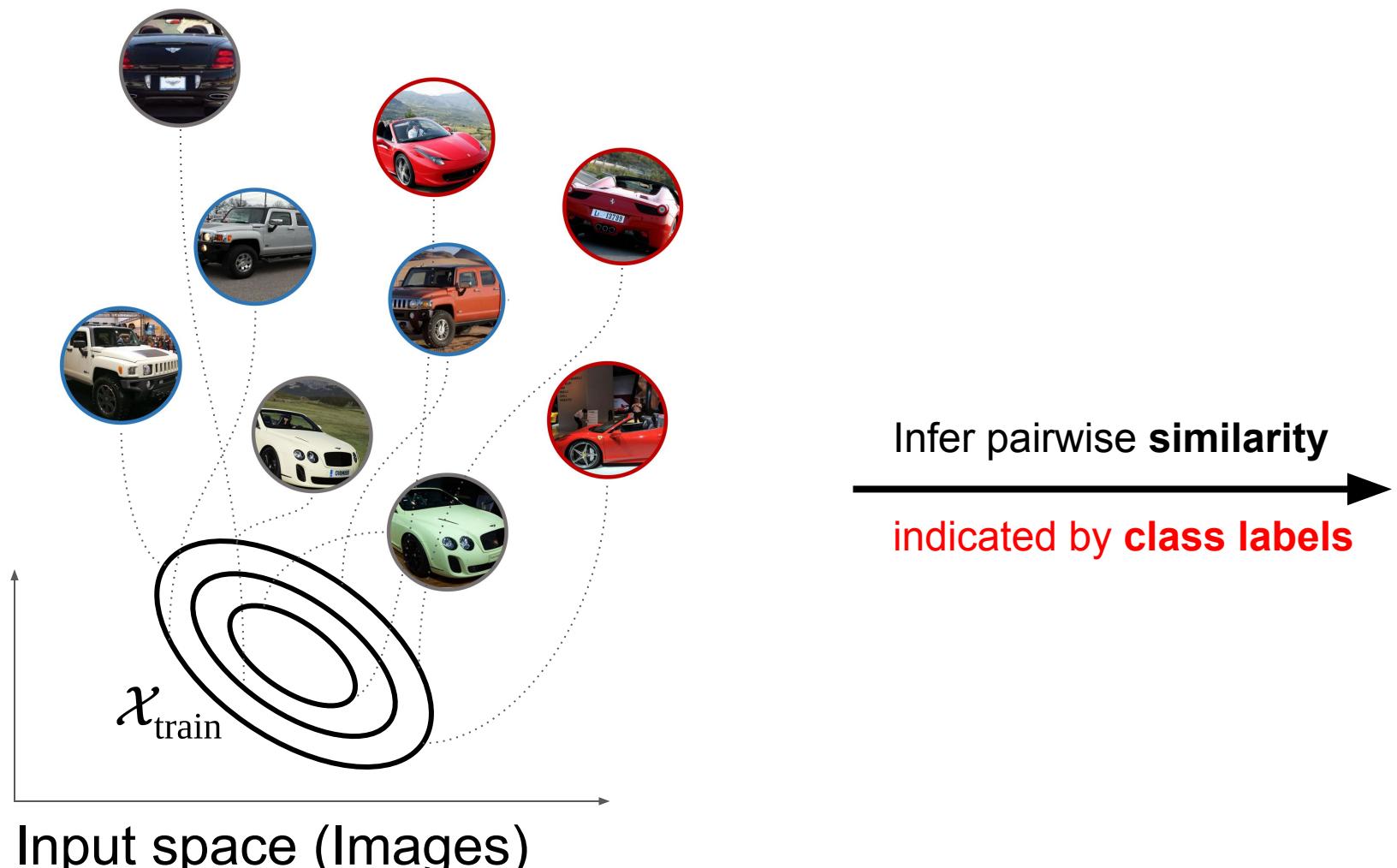
Visual Similarity Learning

Given data distribution $\mathcal{X}_{\text{train}}$, infer pairwise **semantic similarity**.



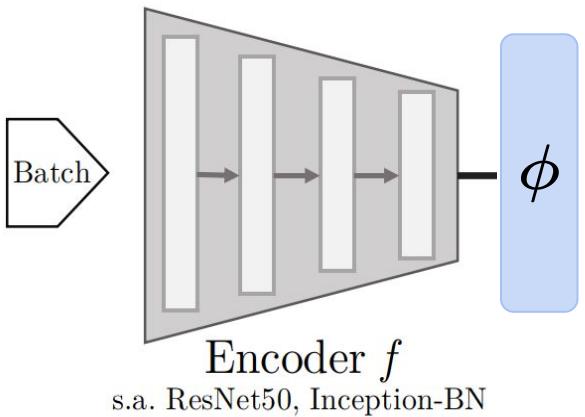
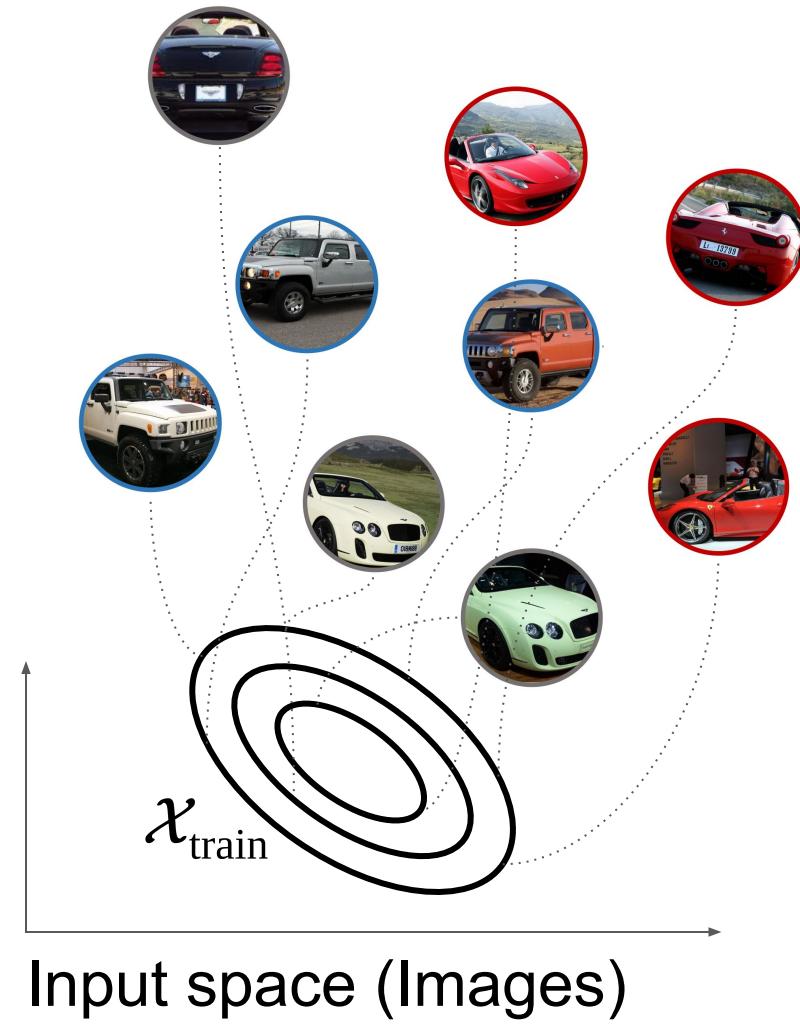
Visual Similarity Learning

Given data distribution $\mathcal{X}_{\text{train}}$, infer pairwise **semantic similarity**.



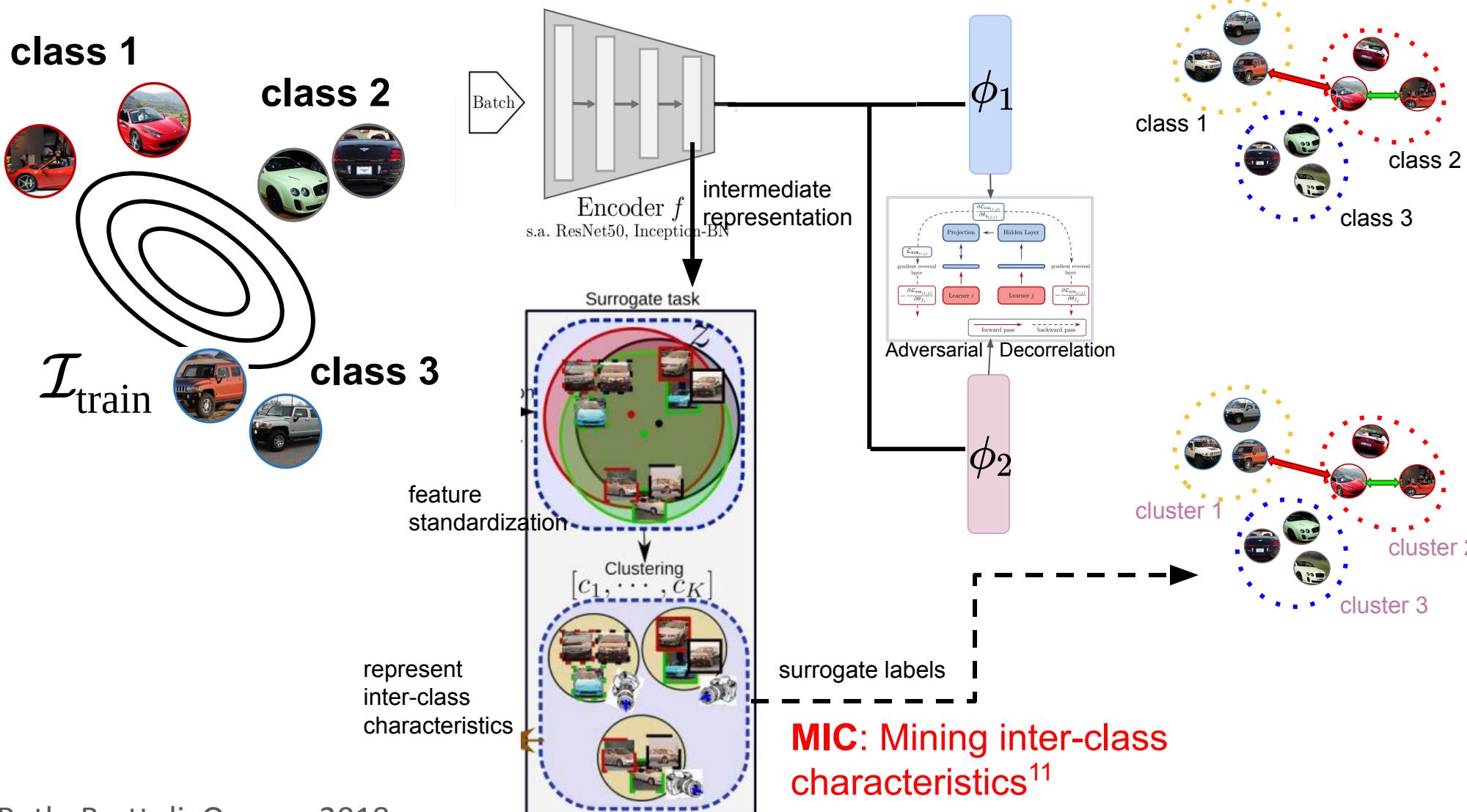
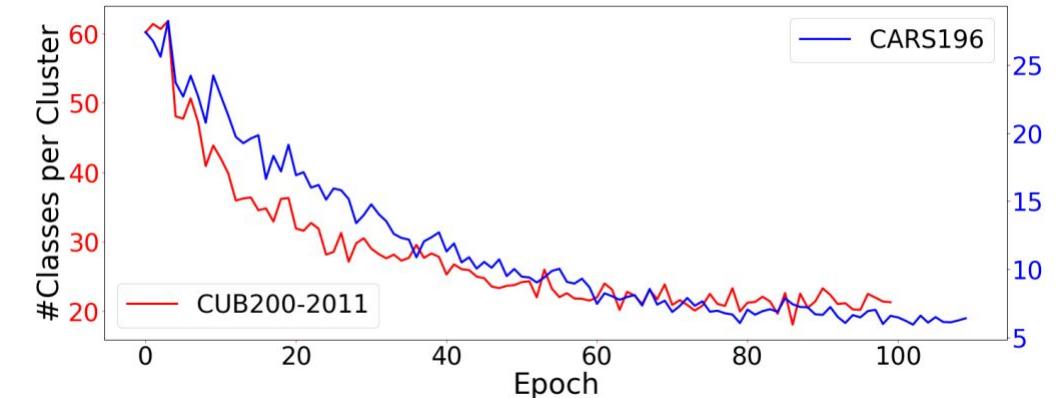
Visual Similarity Learning

Learn **representation** $\phi(x)$ which reflects **semantic similarity** $d(\phi(x_i), \phi(x_j))$ within training distribution $\mathcal{X}_{\text{train}}$.



Ensemble-based DML

- Typically there are **different notions of similarity** (classes, color, viewpoint, ...)
- **Assumption:** Different features improve robustness to OOD data (new classes, etc.)
- **Explicit** similarity specialization of learners
- Only class-labels available: use **unsupervised learning**



$$t = \{I_a, I_p, I_n\} \xrightarrow{y_a = y_p \neq y_n} \mathcal{L}_{\text{disc}}$$

Class-discriminative features

- capture features **separating between classes**
- aggregate features into ‘classes’
- very specialized
- e.g. “Ferrari” vs. “Hummer”

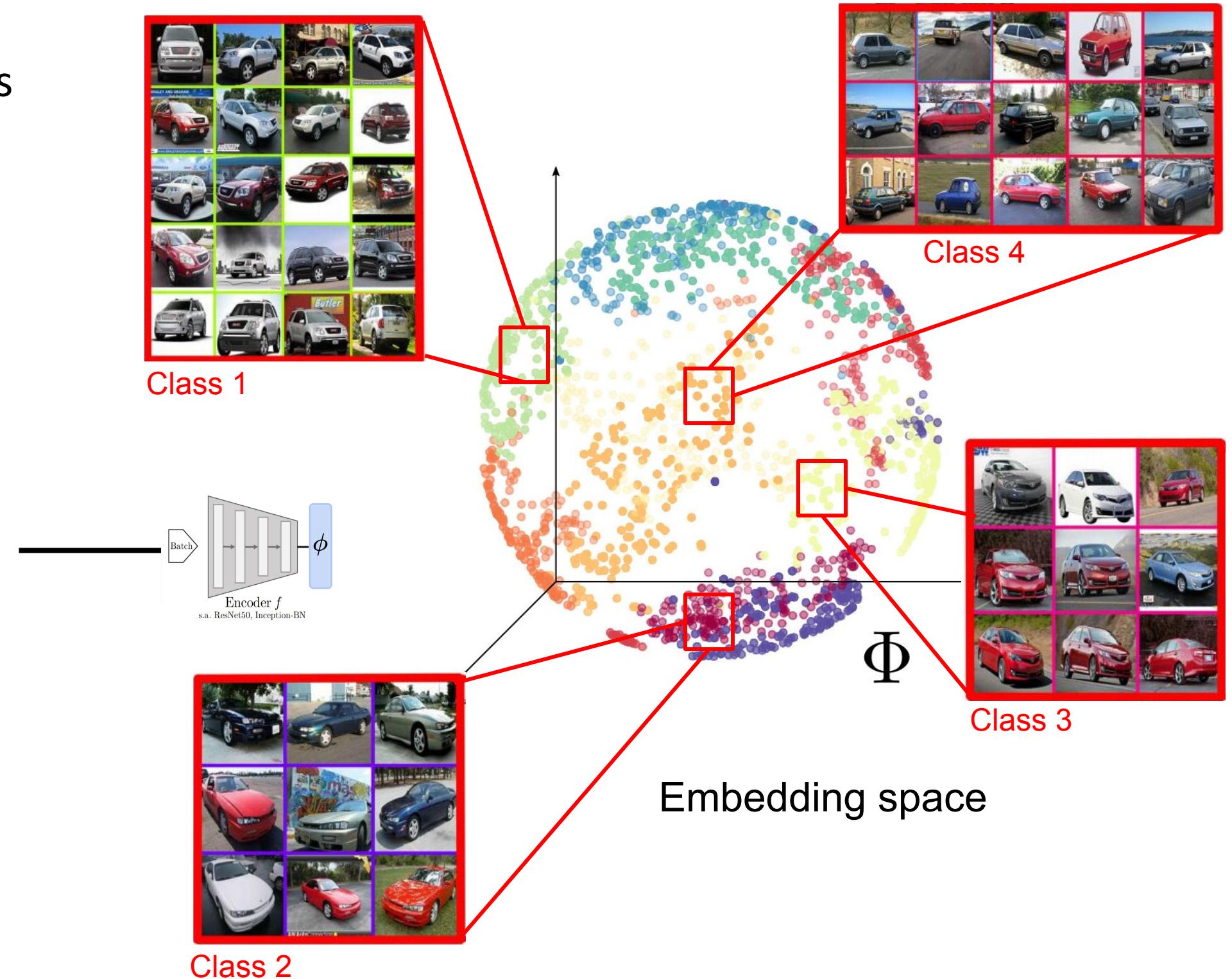
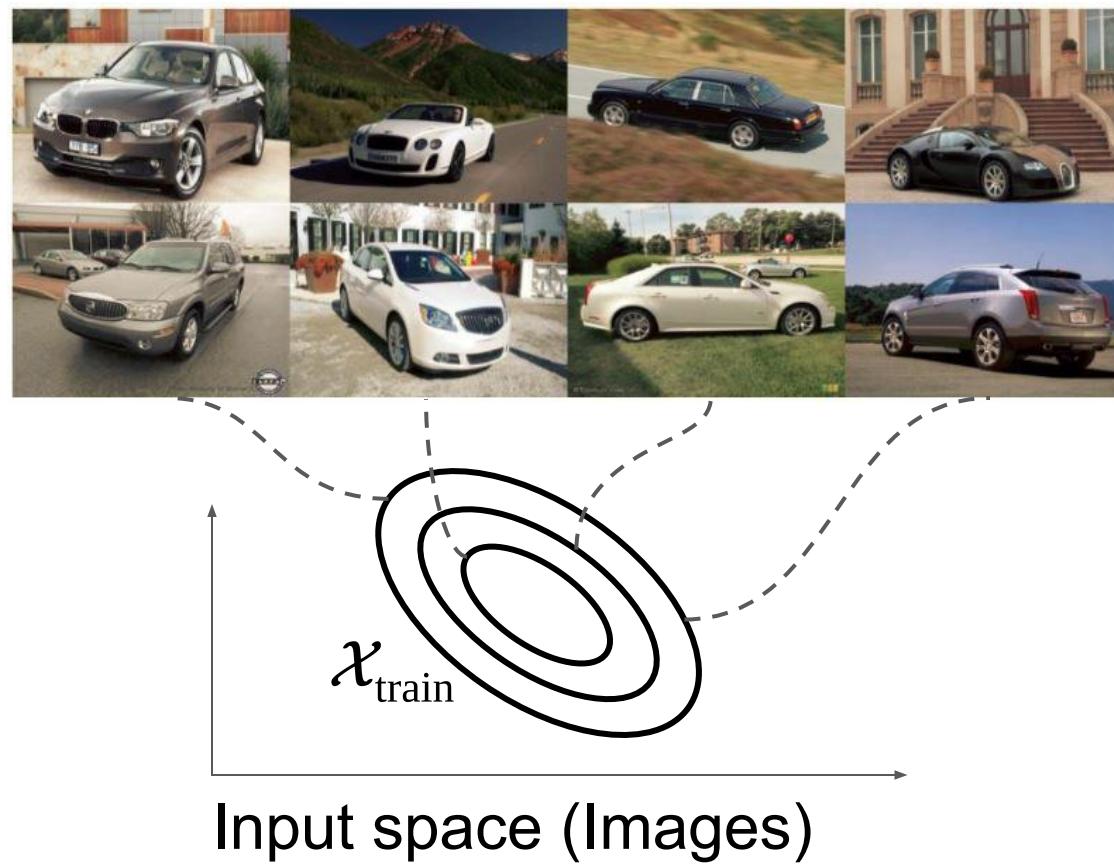
$$t = \{I_a, I_p, I_n\} \xrightarrow{c_a = c_p \neq c_n} \mathcal{L}_{\text{surrogate}}$$

Inter-class characteristics¹¹

- anchors and positives **from same cluster** (different class labels!)
- represent inter-class characteristics, e.g. viewpoint, color, etc.

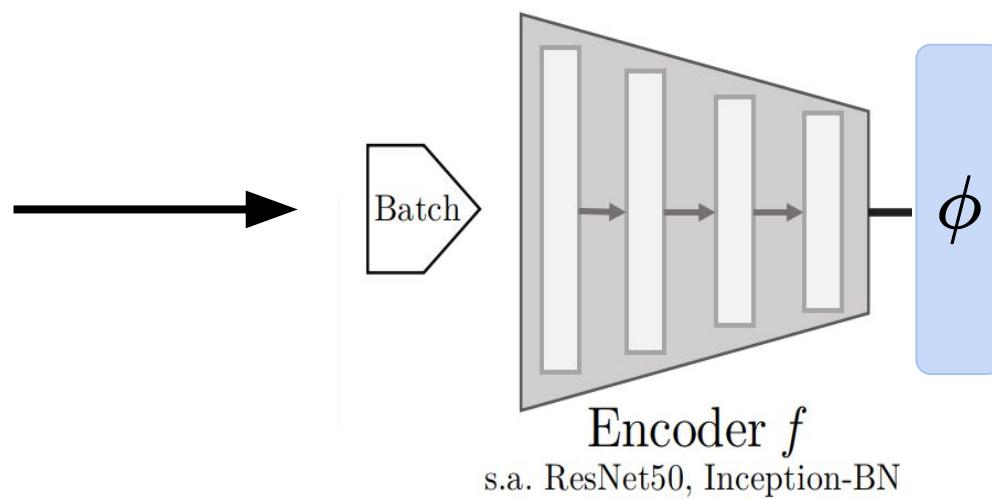
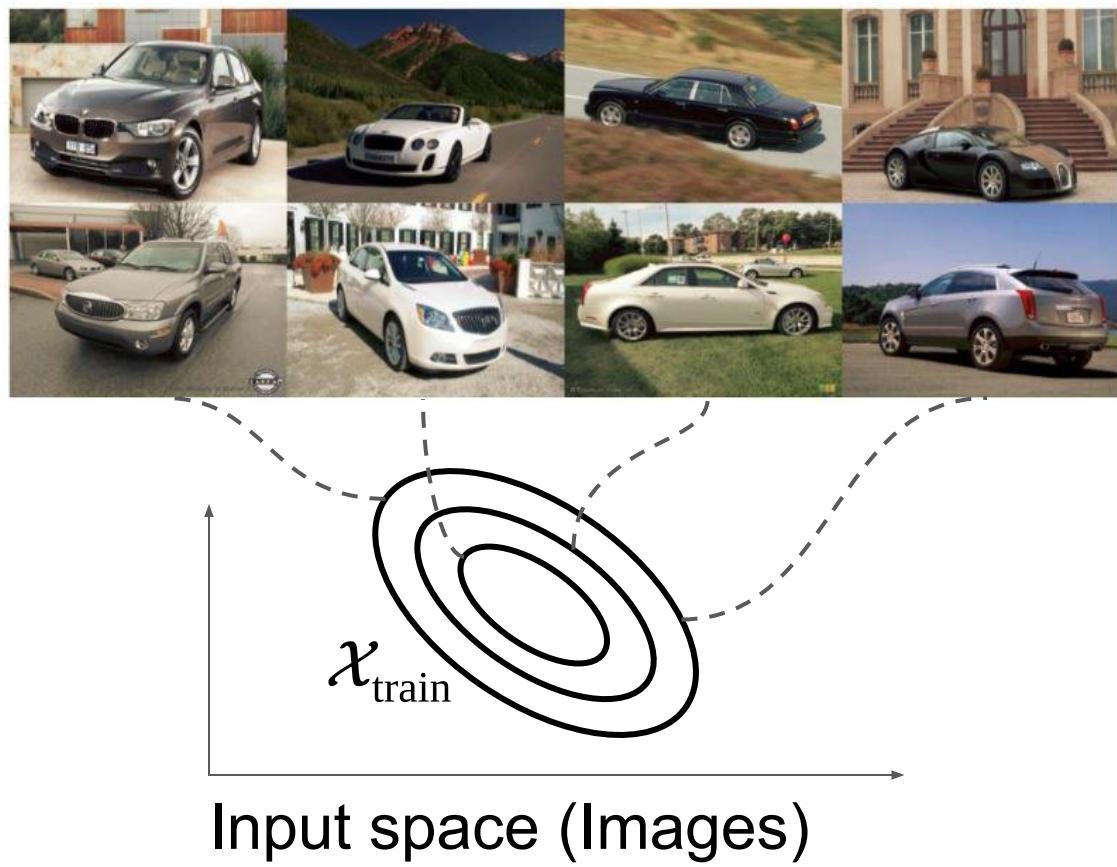
Learning Visual Representations

Learn representation Φ which reflects semantic similarity within training distribution $\mathcal{X}_{\text{train}}$.

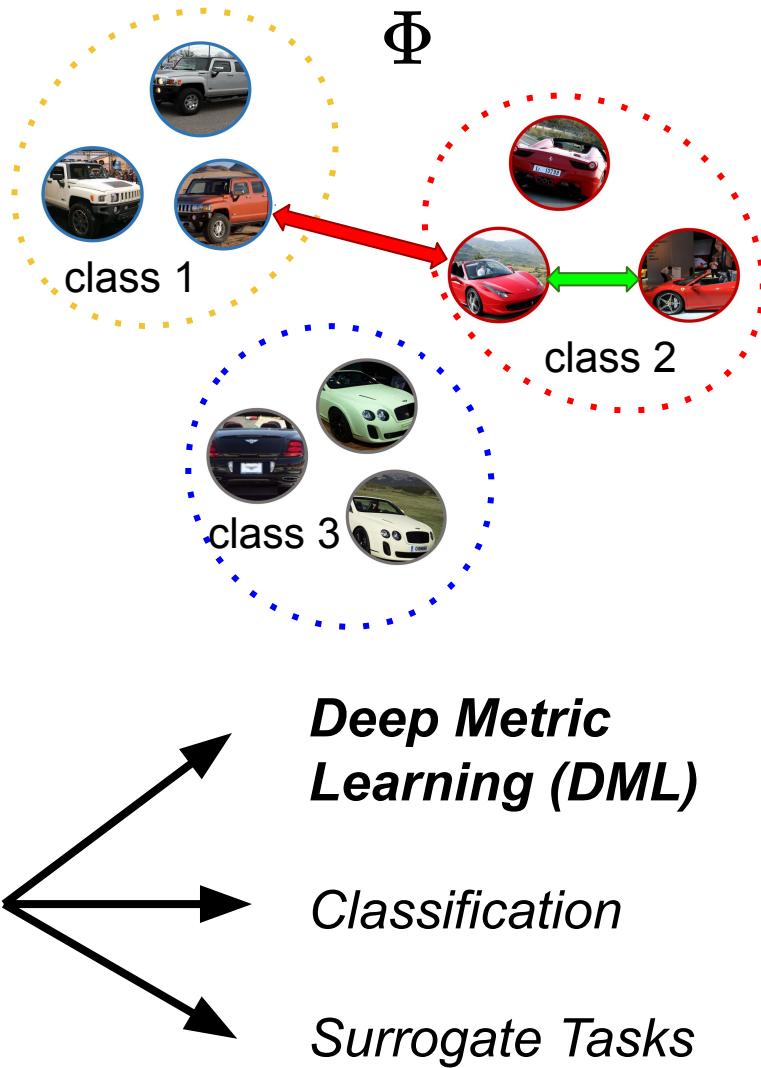


Learning Visual Representations

Learn **representation** $\phi(x)$ which reflects **semantic similarity** within training distribution $\mathcal{X}_{\text{train}}$.

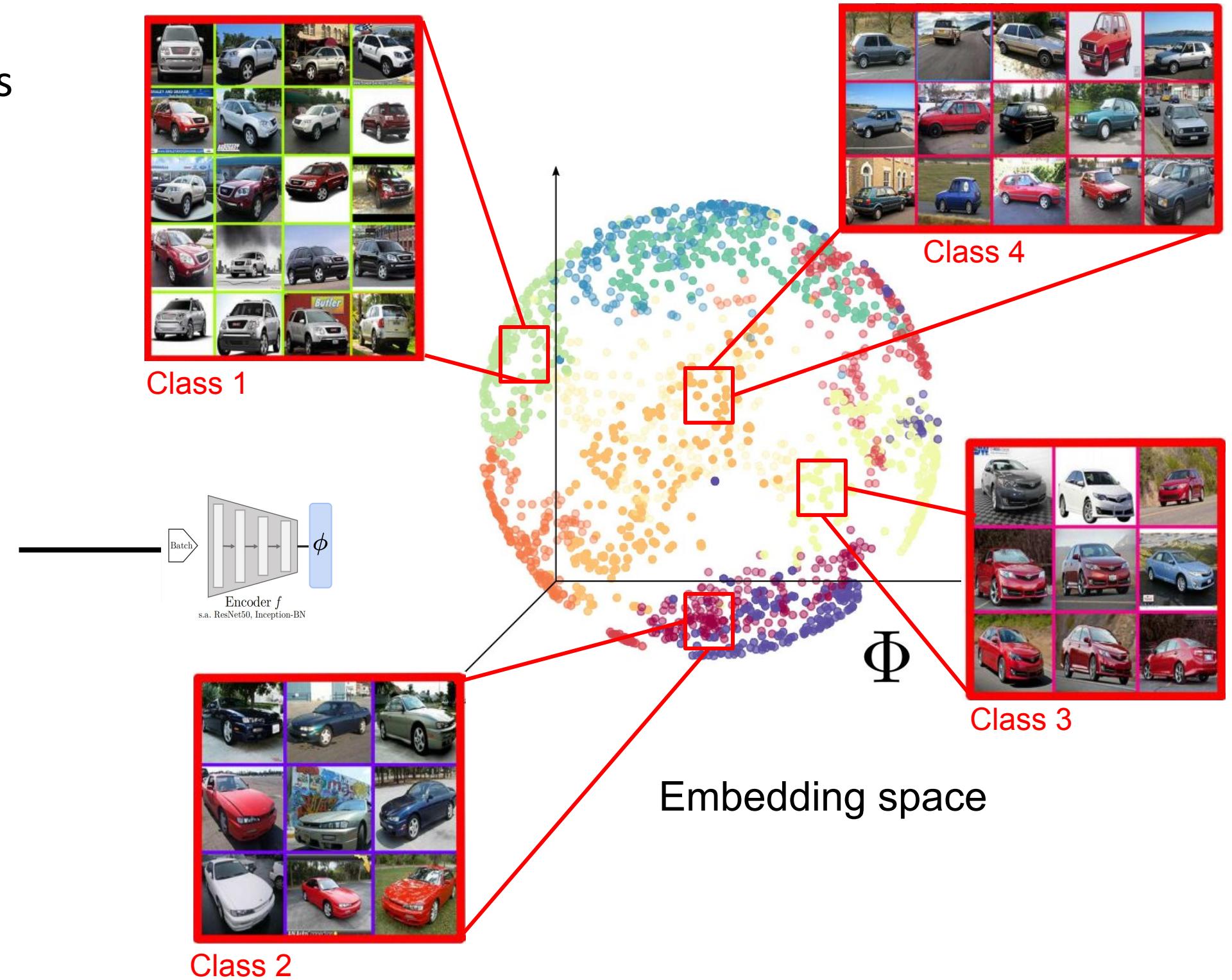
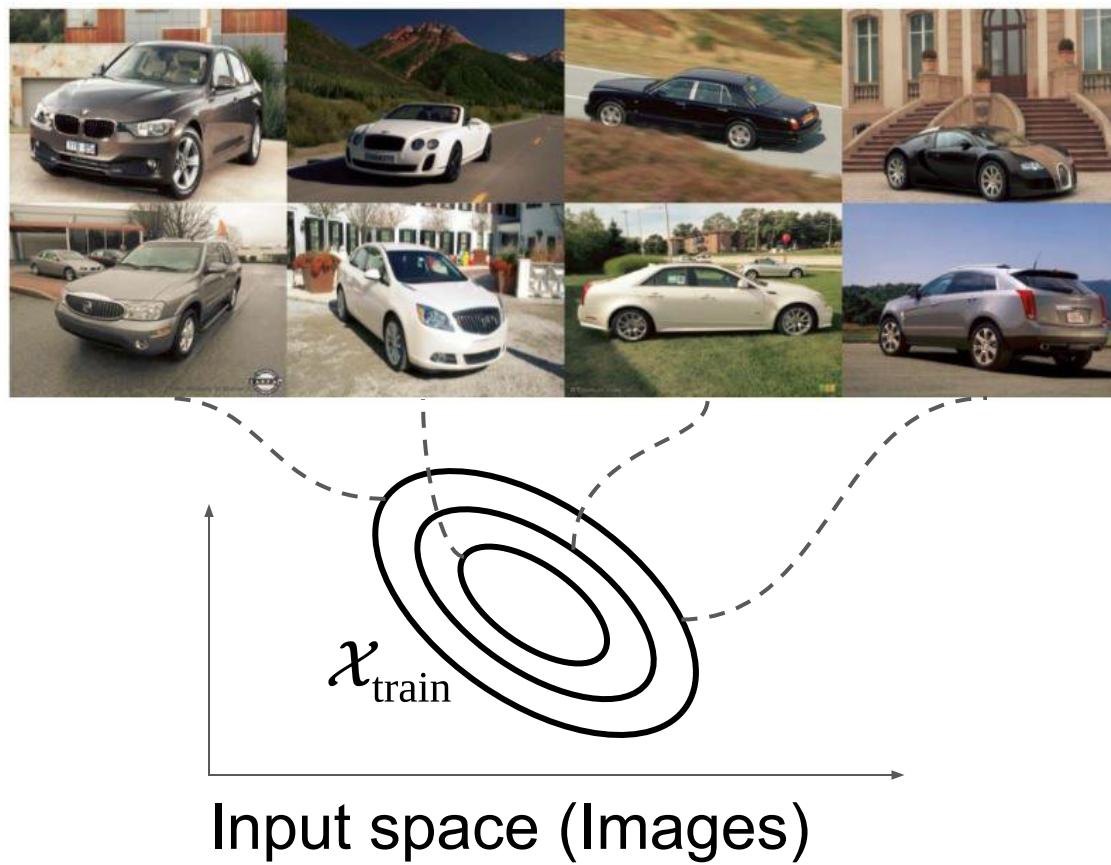


Representation Learning

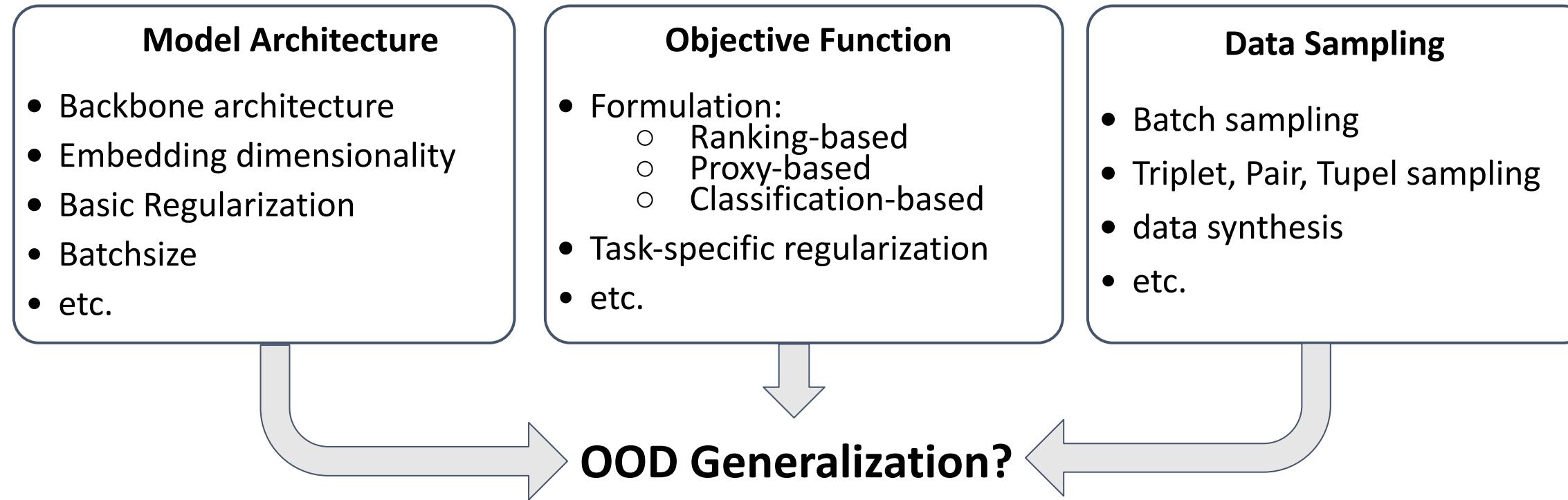


Learning Visual Representations

Learn representation Φ which reflects semantic similarity within training distribution $\mathcal{X}_{\text{train}}$.



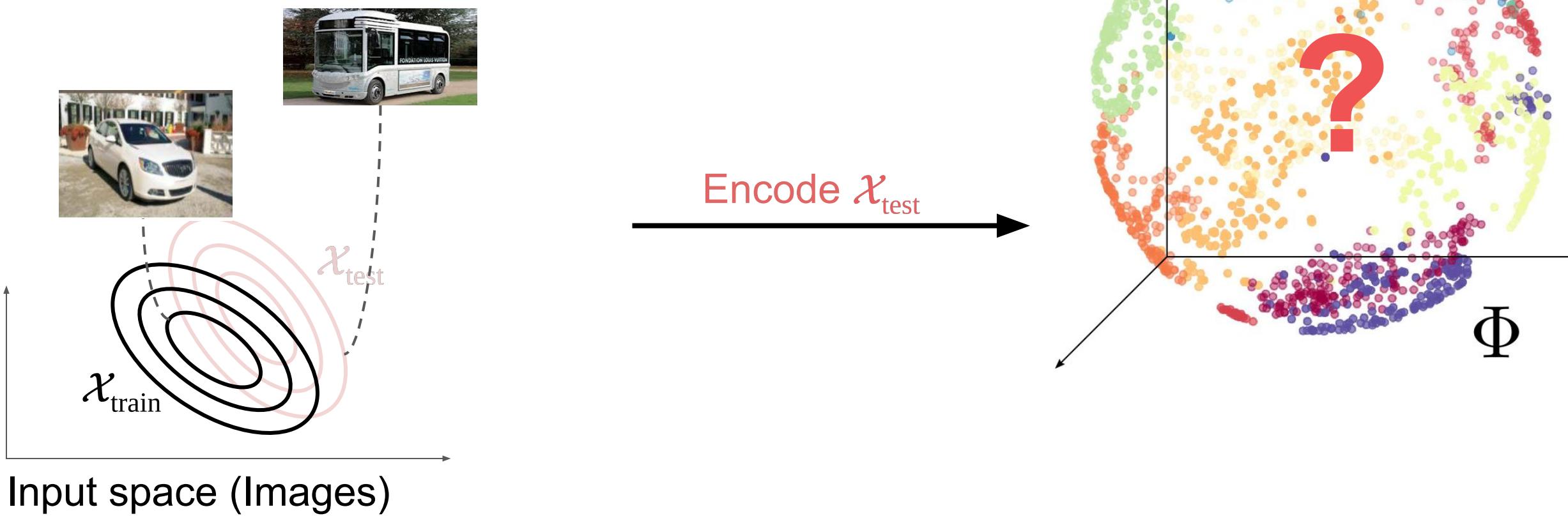
Introduction DML



Learning Visual Representations

(Out-Of-Distribution-Generalization)

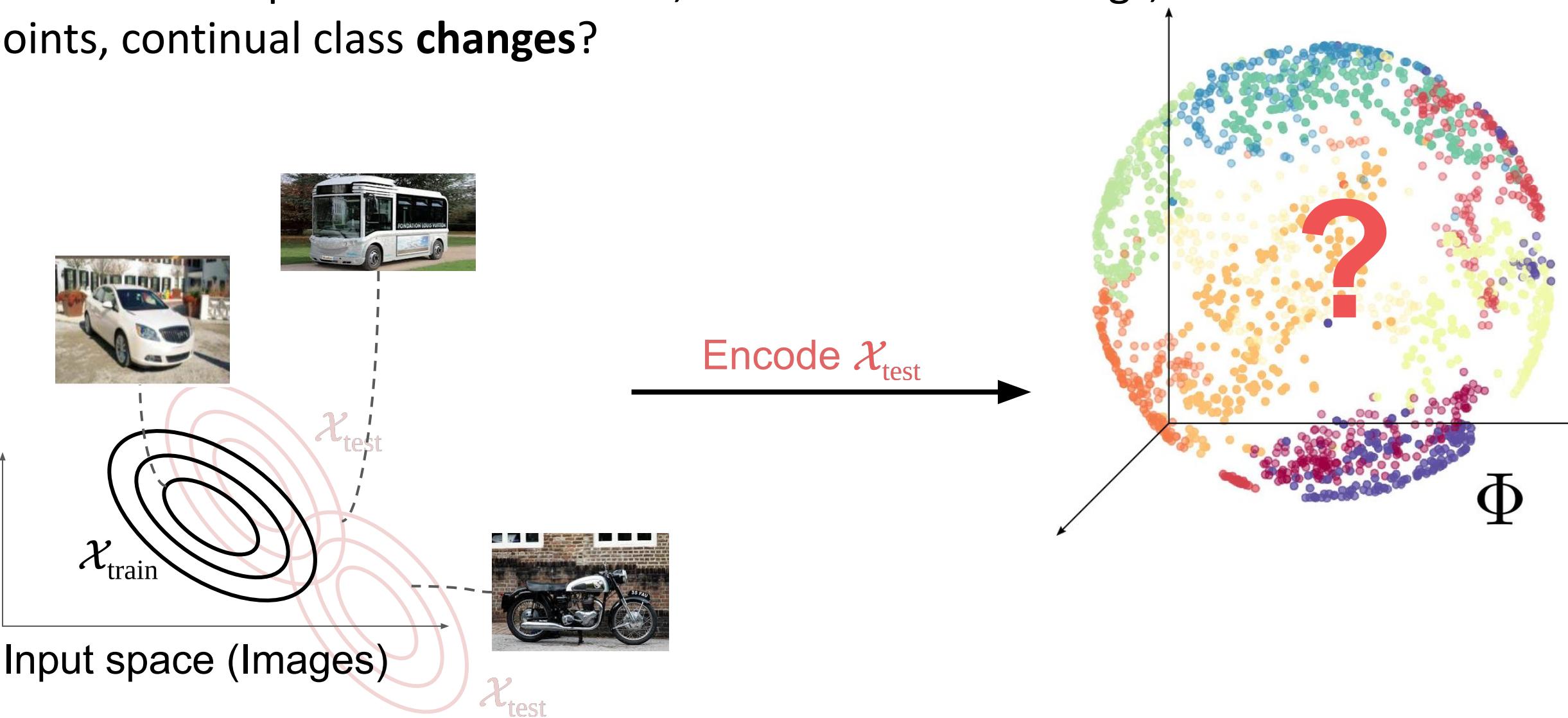
How well does Φ capture **unseen** classes, **unknown** surroundings, viewpoints, continual class **changes**?



Learning Visual Representations

(Out-Of-Distribution-Generalization)

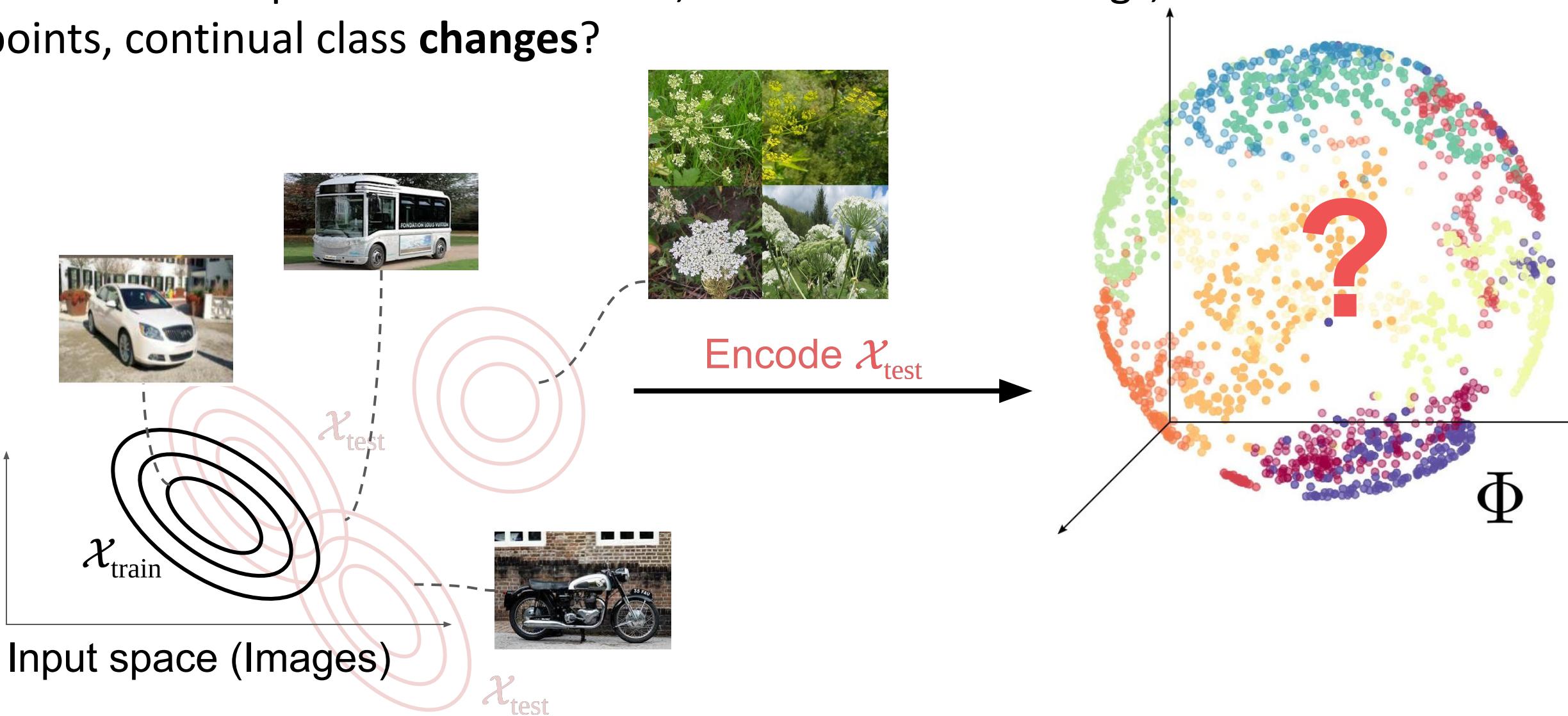
How well does Φ capture **unseen** classes, **unknown** surroundings, viewpoints, continual class **changes**?



Learning Visual Representations

(Out-Of-Distribution-Generalization)

How well does Φ capture **unseen** classes, **unknown** surroundings, viewpoints, continual class **changes**?

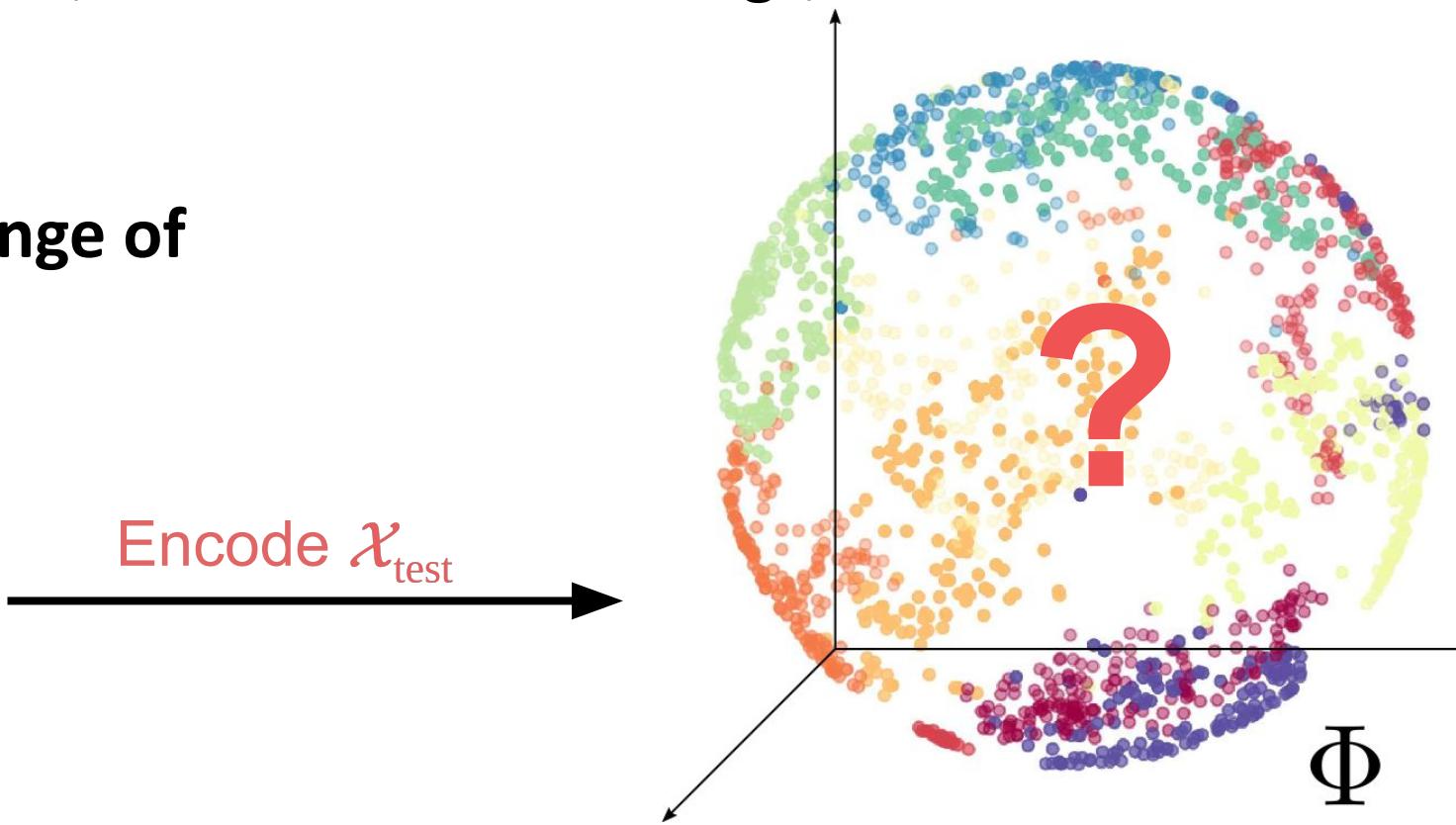
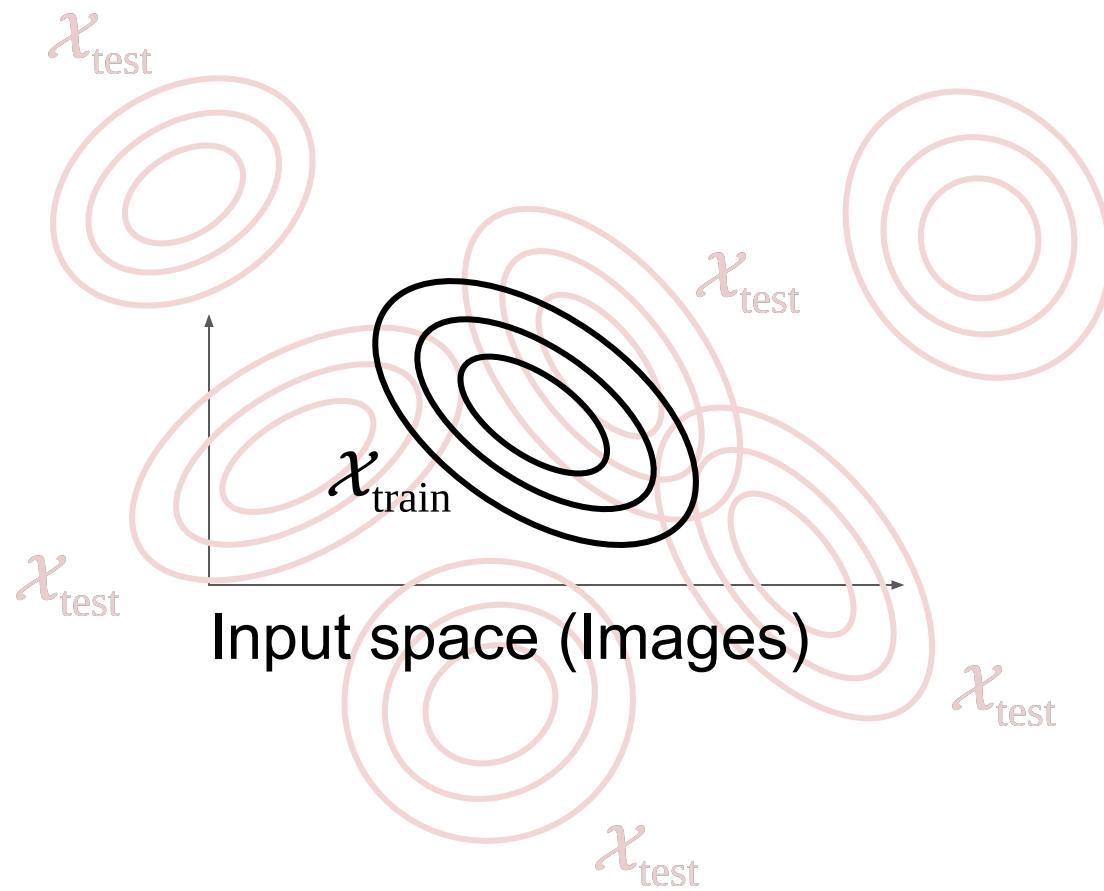


Learning Visual Representations

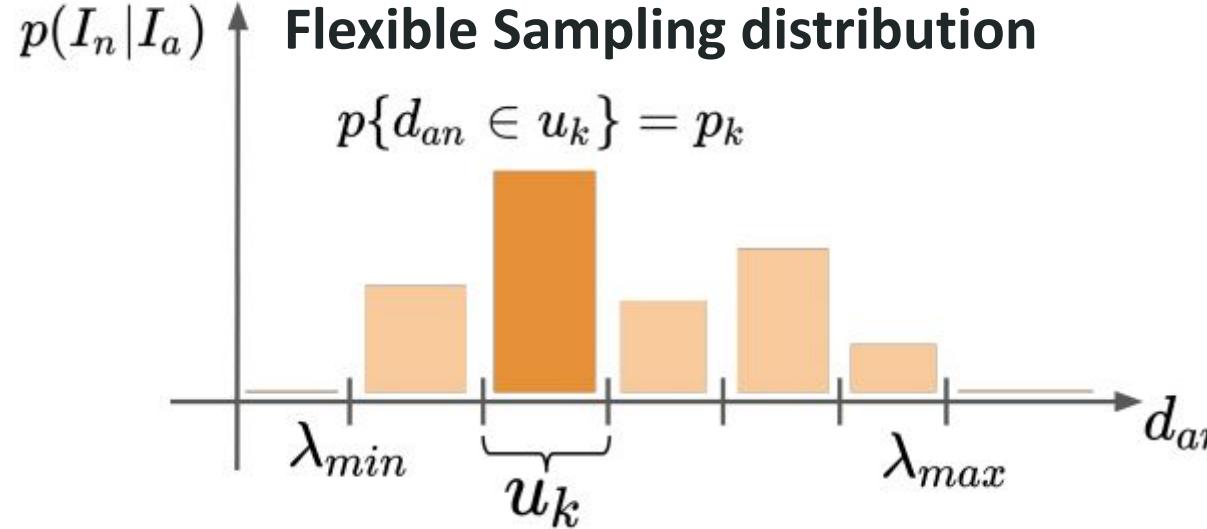
(Out-Of-Distribution-Generalization)

How well does Φ capture **unseen** classes, **unknown** surroundings, viewpoints, continual class **changes**?

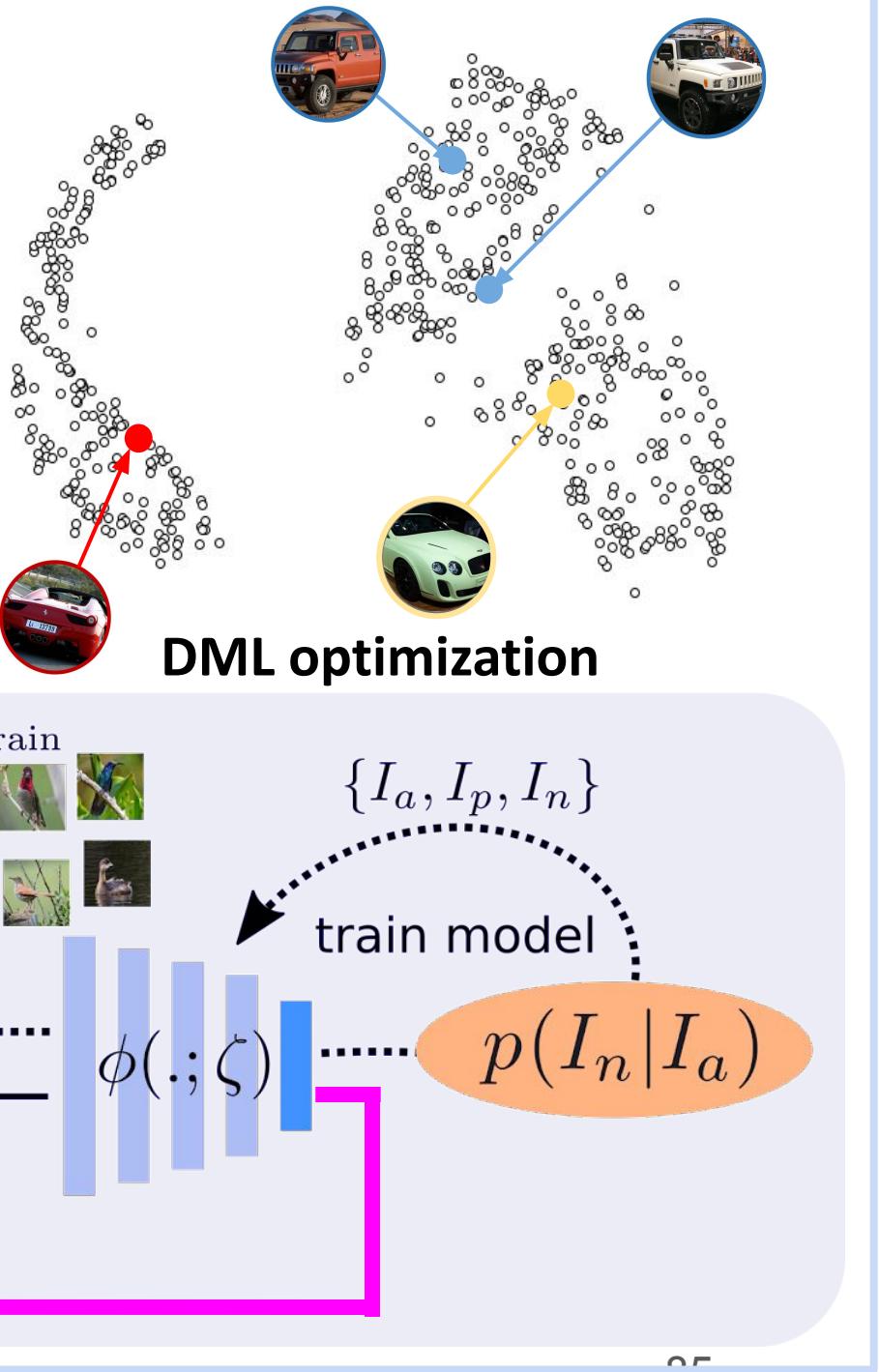
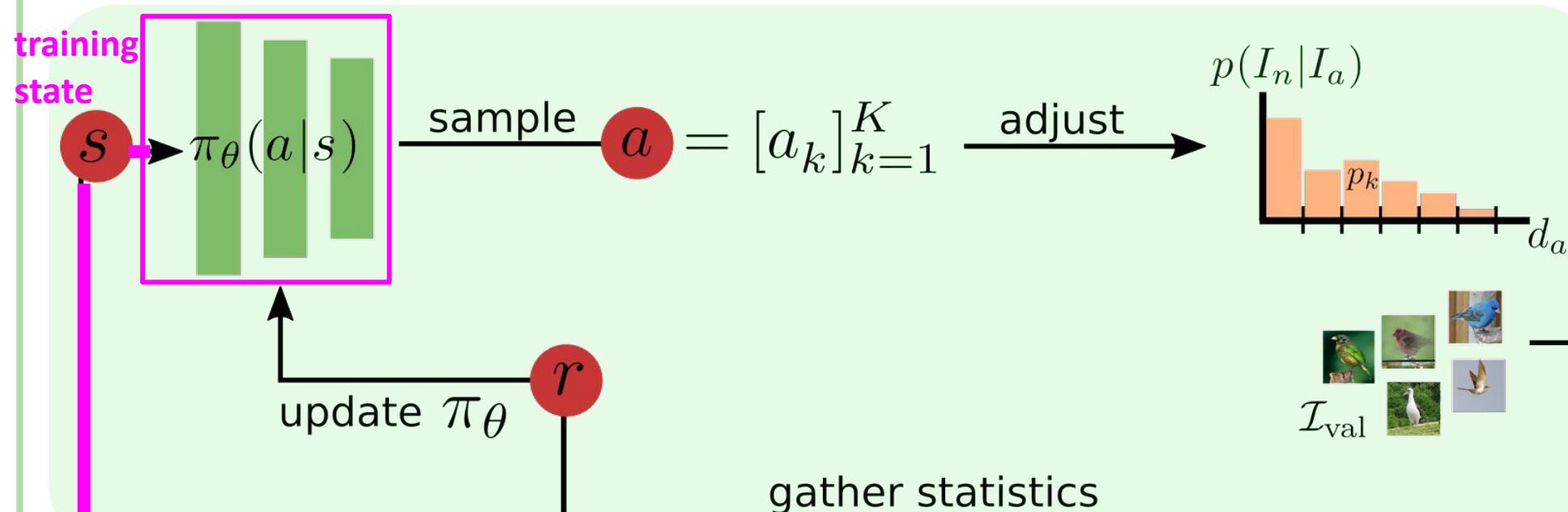
Evaluation needs to consider **broad range of test distributions and difficulty.**



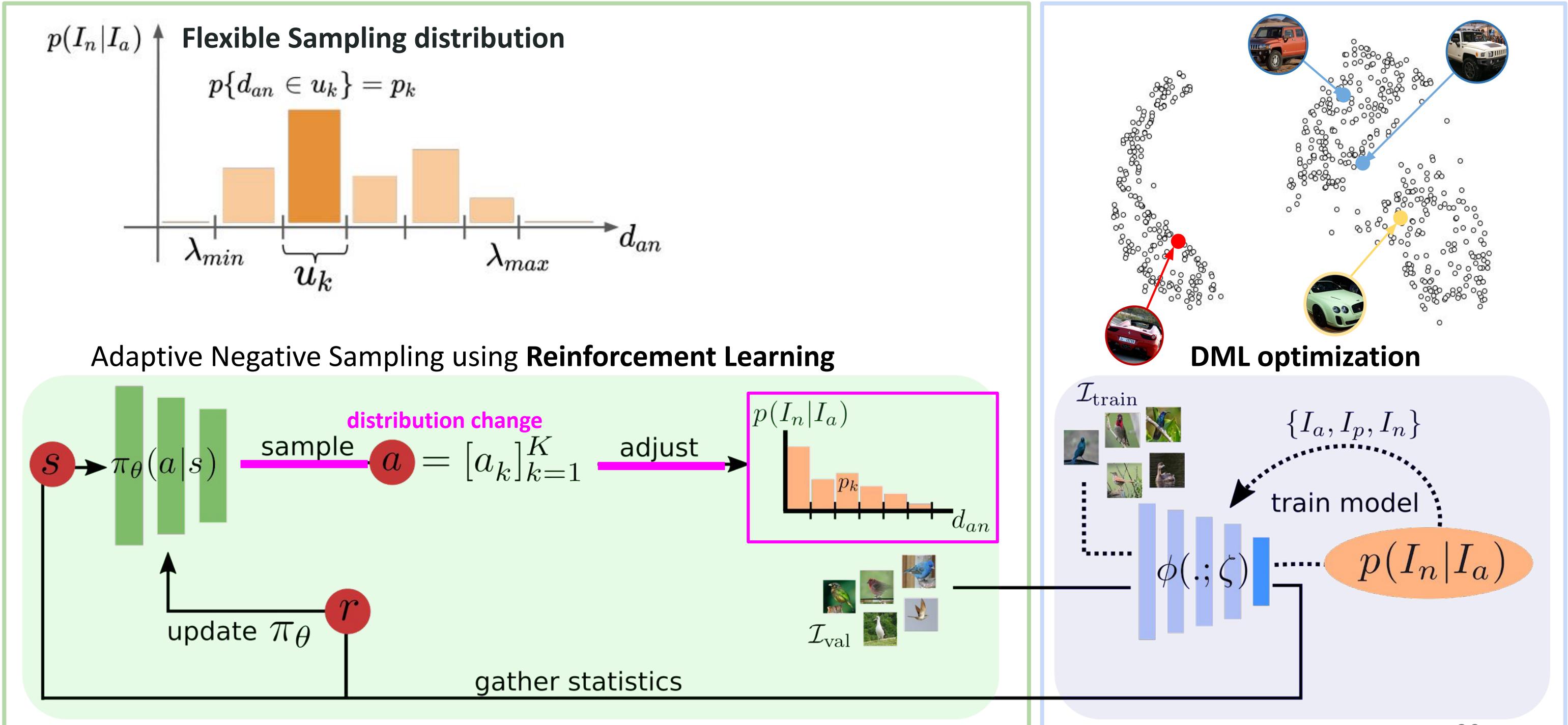
Negative Sampling in Ranking-based DML



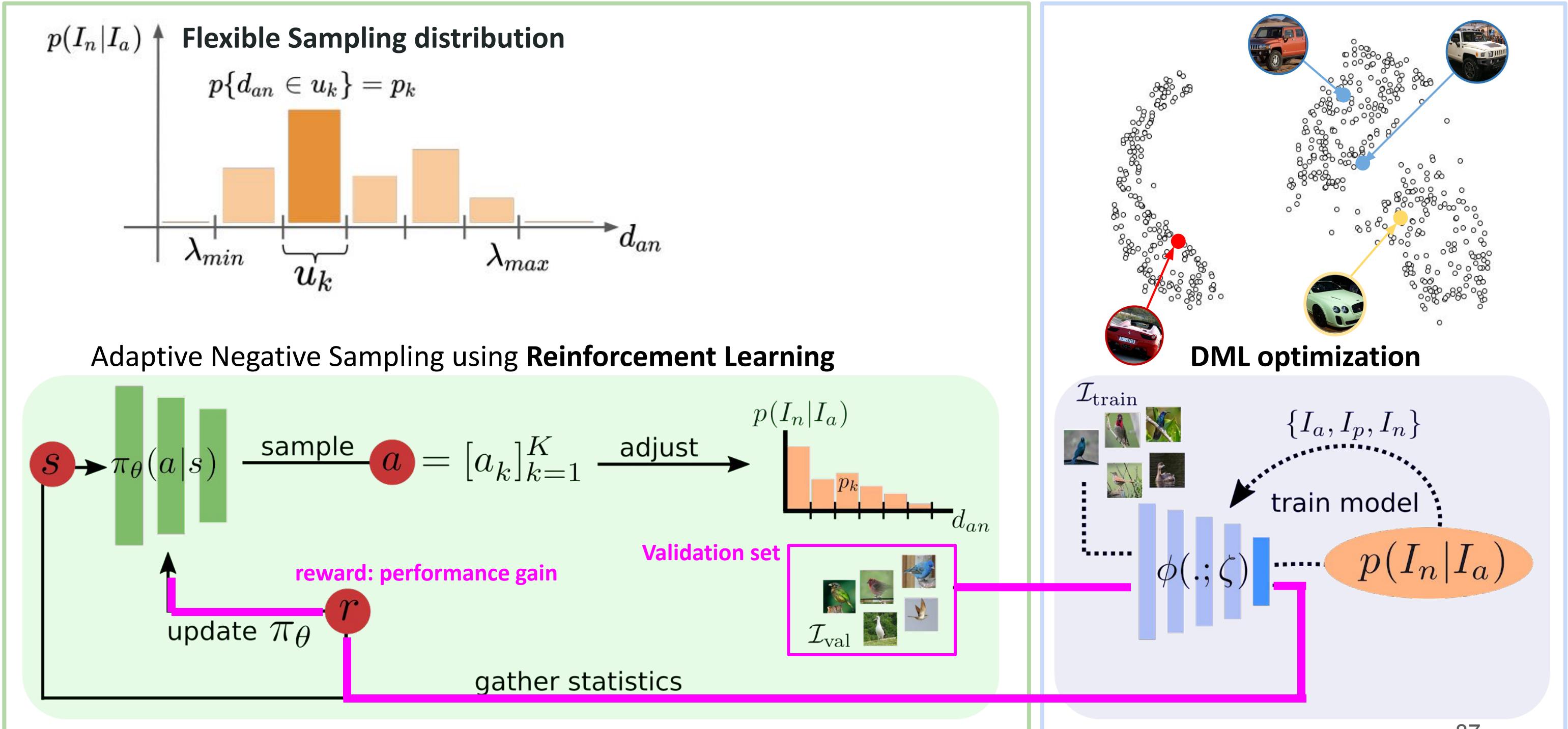
Adaptive Negative Sampling using Reinforcement Learning



Negative Sampling in Ranking-based DML

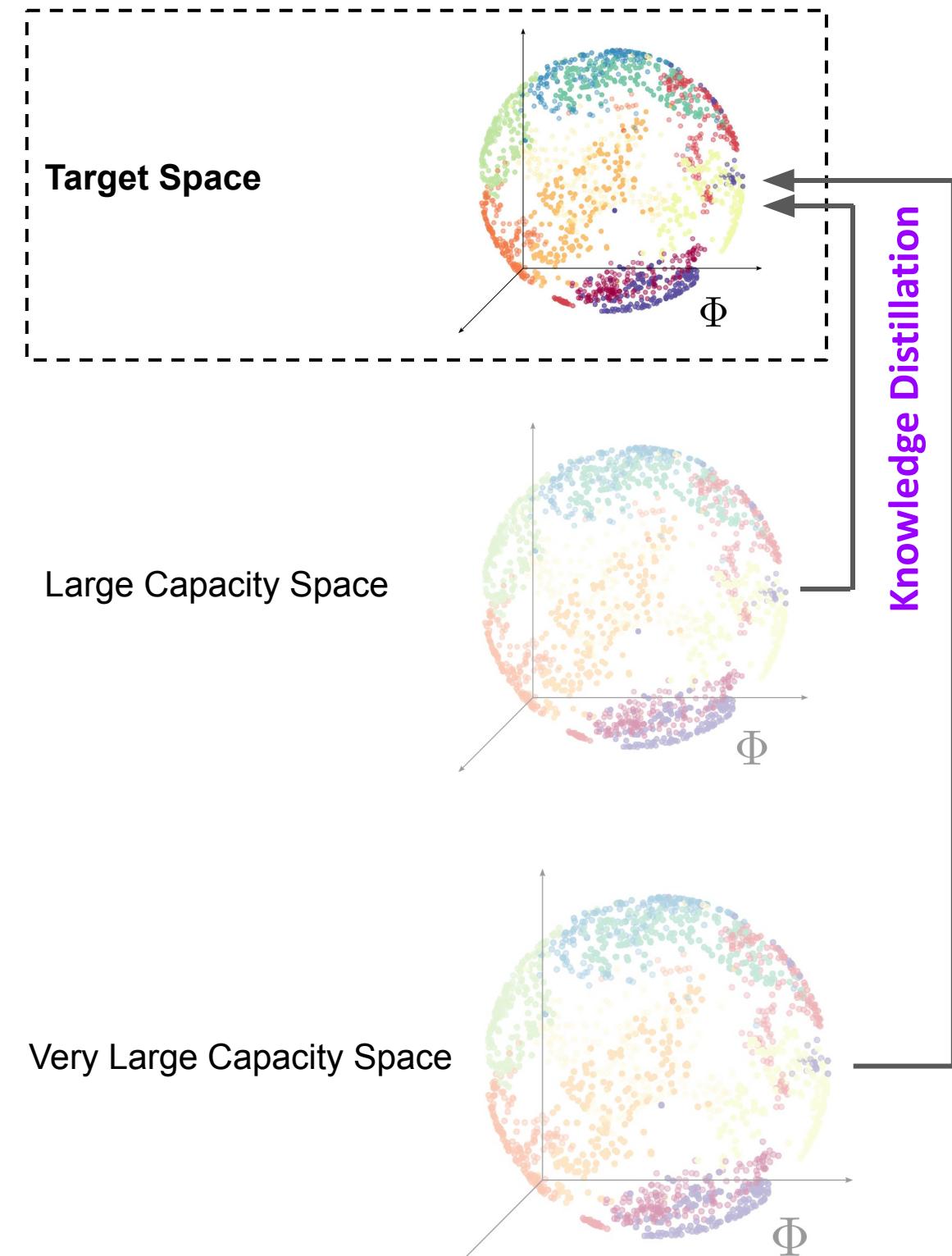
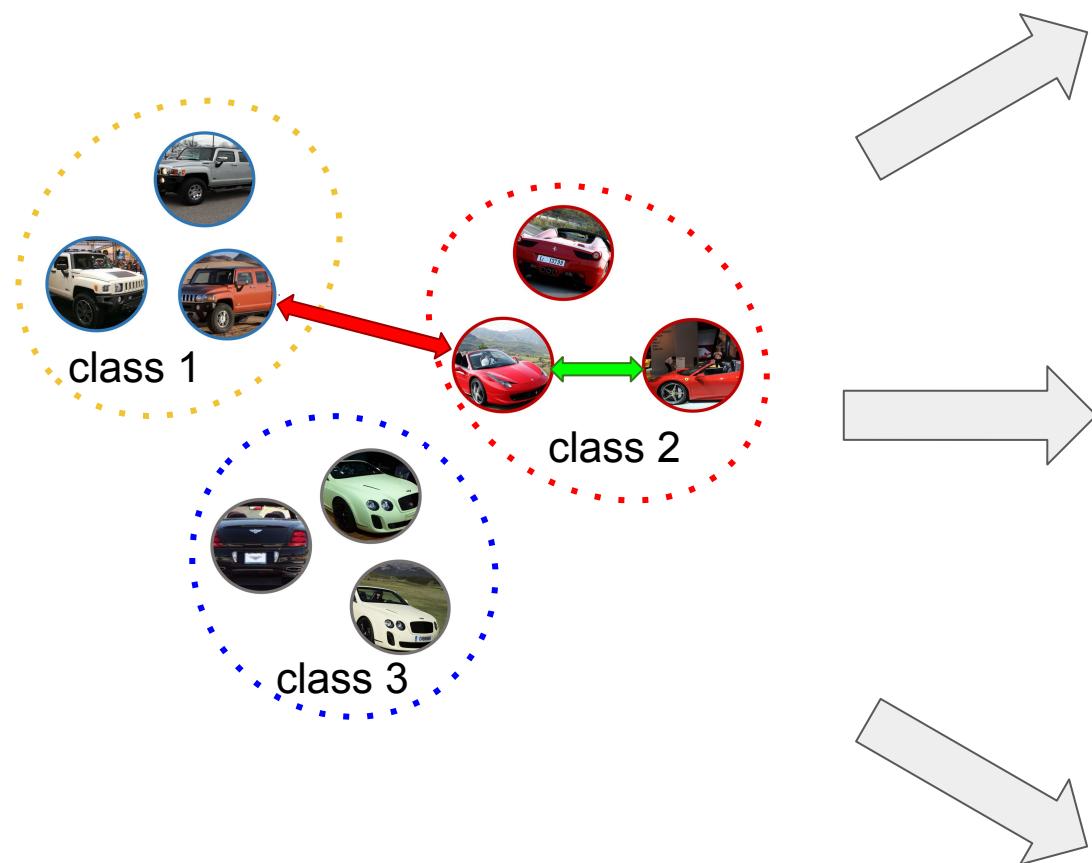


Negative Sampling in Ranking-based DML



Ensemble-based DML

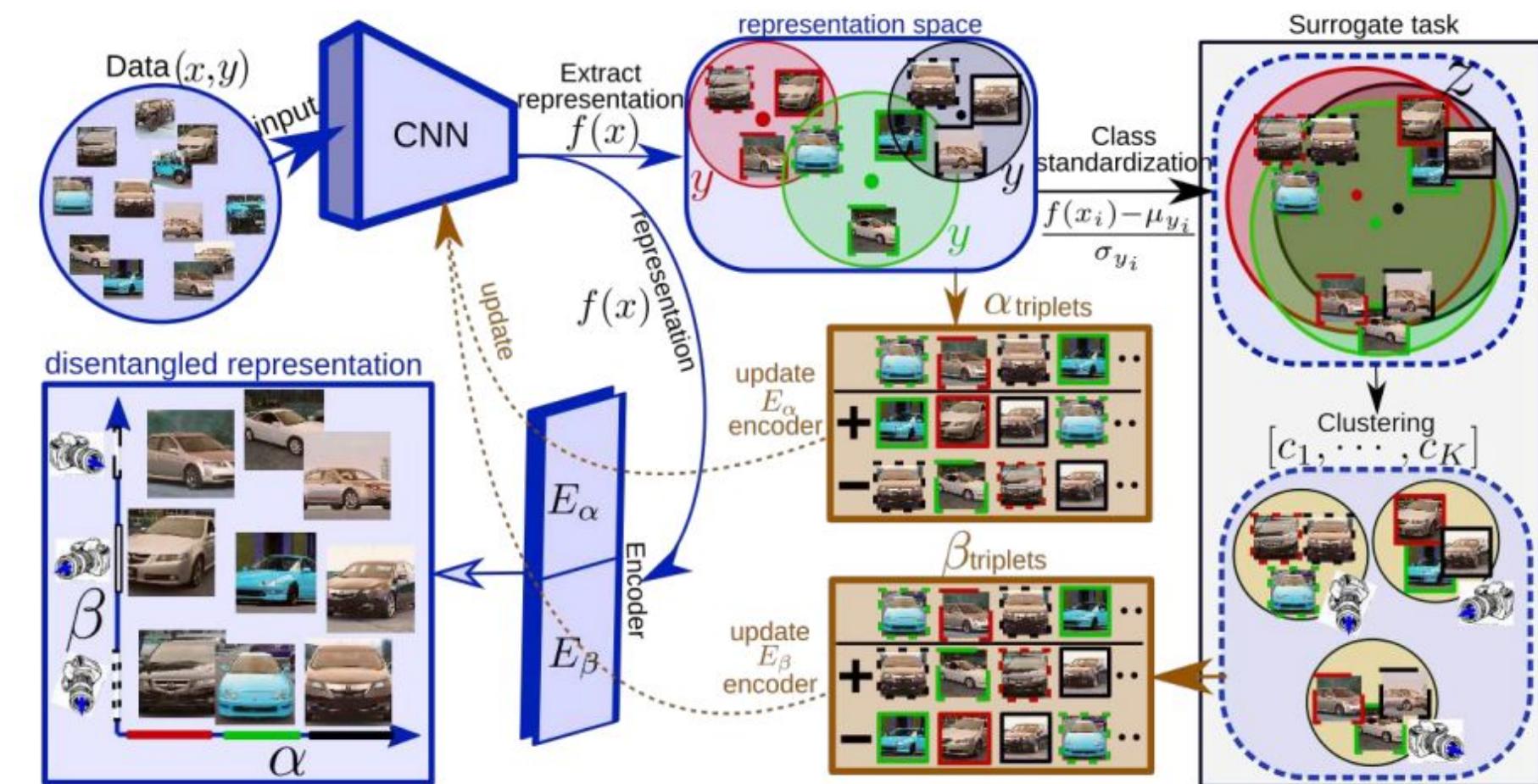
- Multi Self-Distillation³



³ Roth, Milbich et al.; ICML 2021; S2SD: Simultaneous Similarity-based Self-distillation for Deep Metric Learning

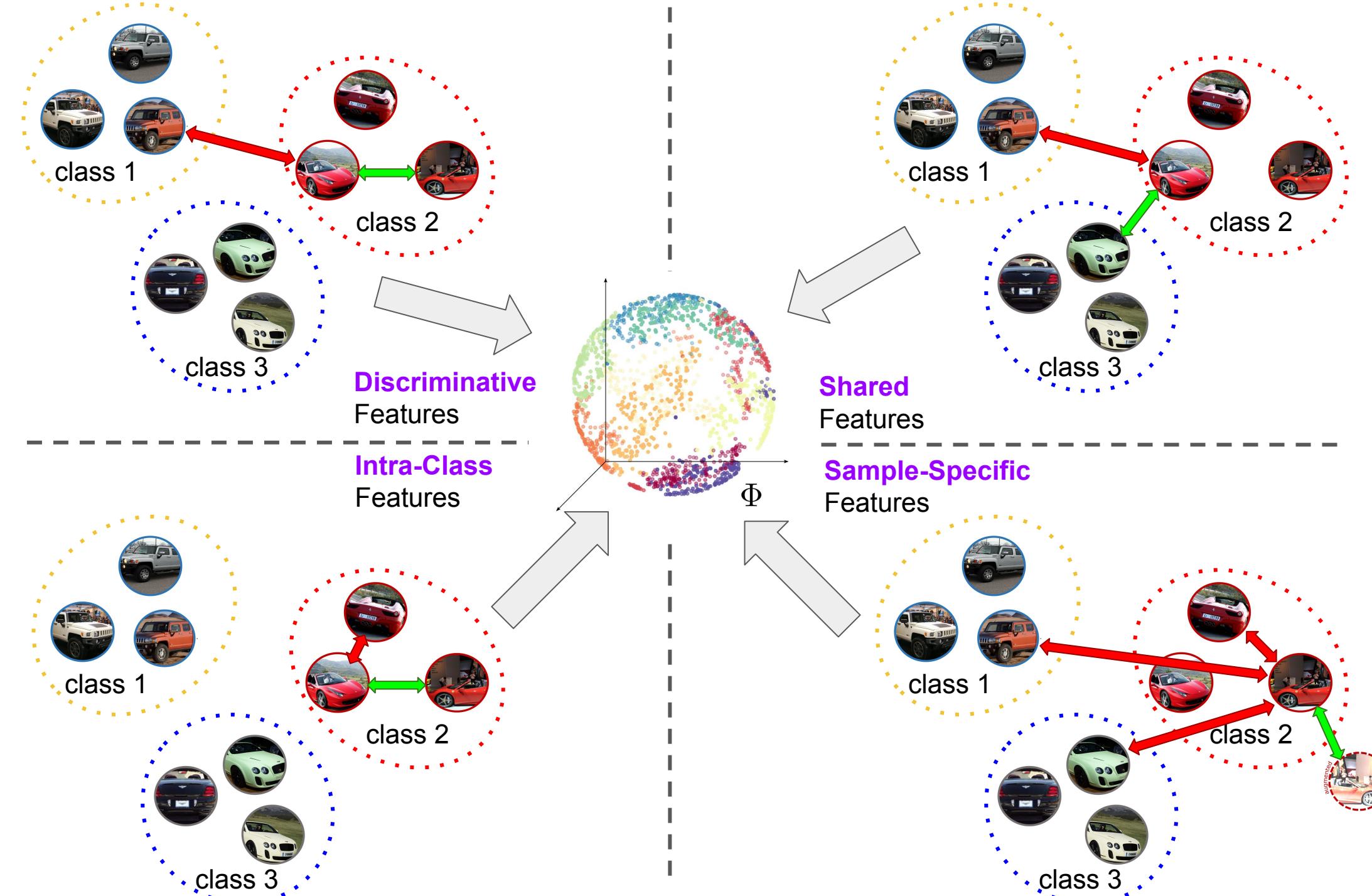
Ensemble-based DML

- Shared Features, DiVA, **MIC** (complementary semantics, “multi-task learning”)
- **explicit** specialization of learners (semantics)



Overview: Typical Learning Concepts for DML Generalization

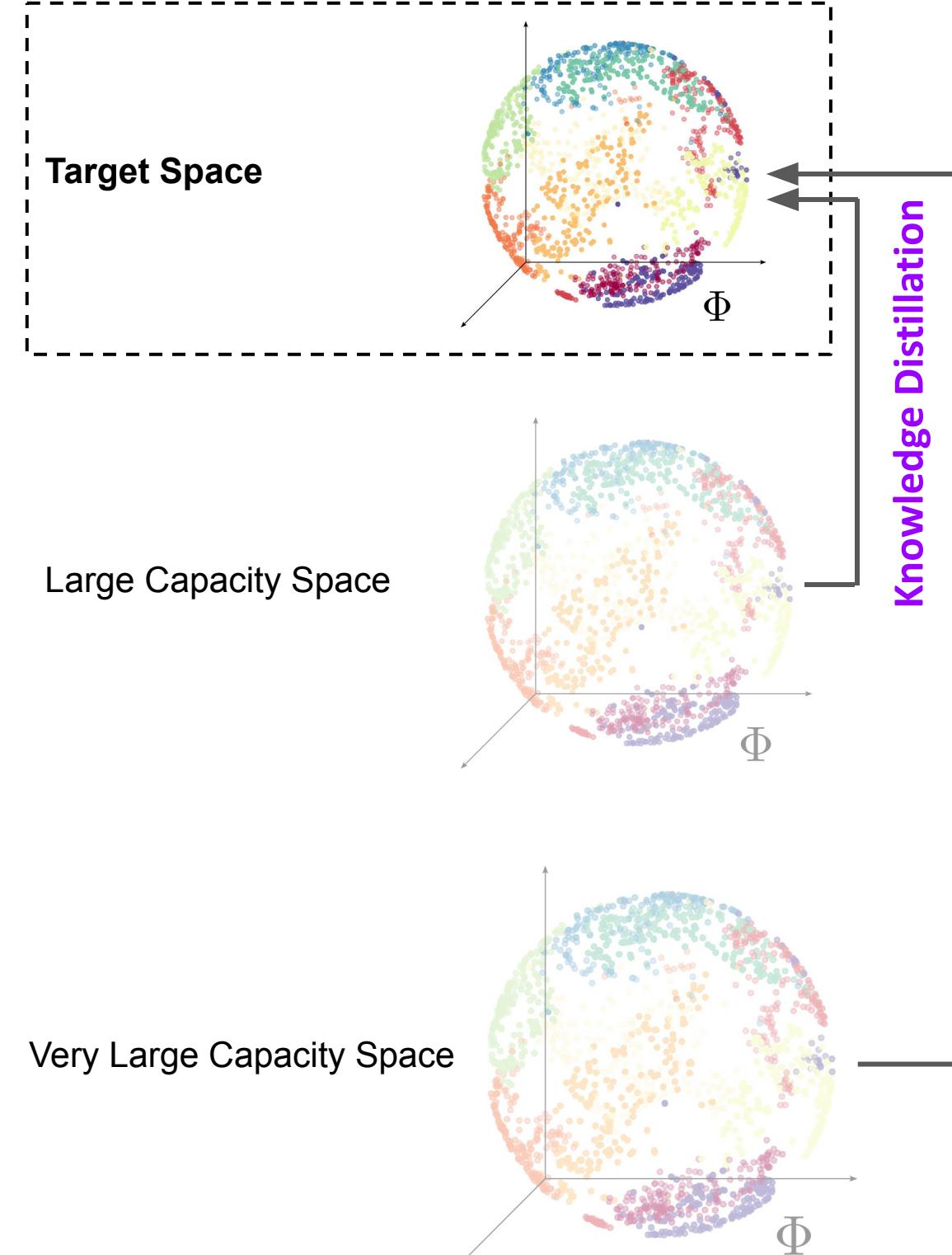
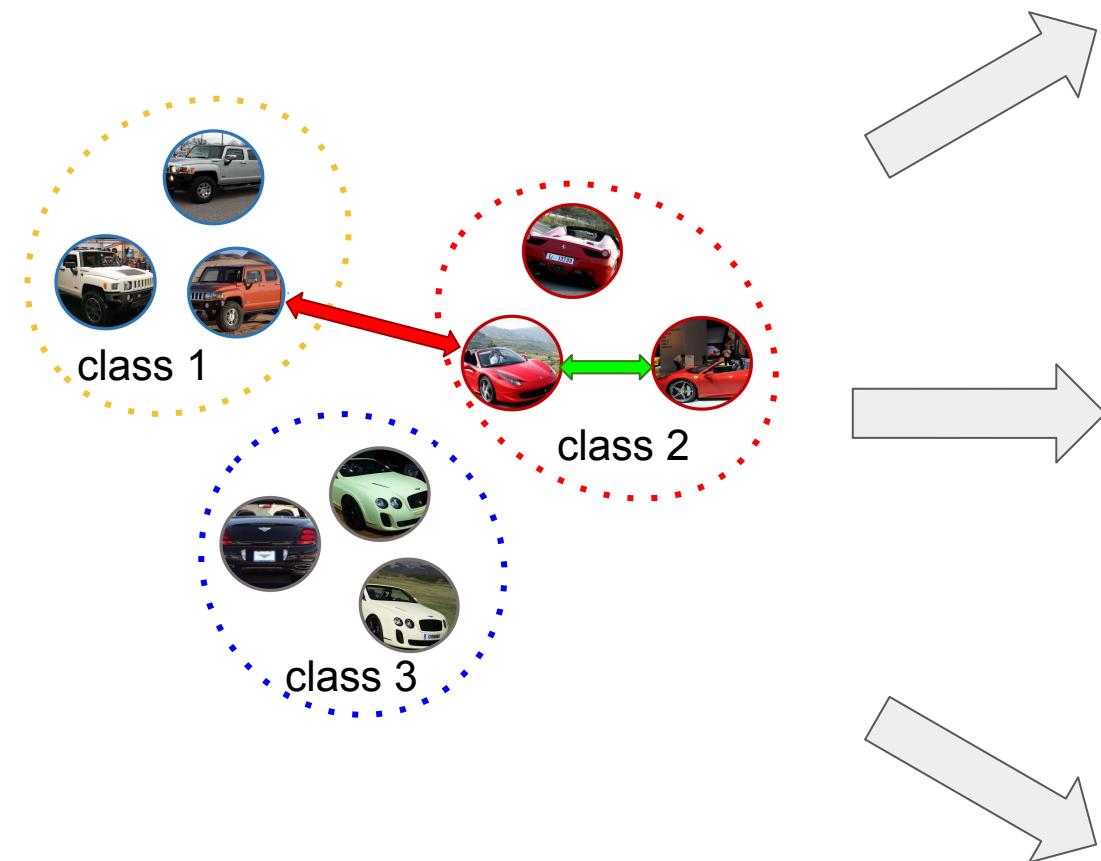
- Diverse Feature Aggregation²
- Multi Self-Distillation
- Proxy-Learning
- ...



² Milbich, Roth et al.; ECCV 2020; DiVA: Diverse Visual Feature Aggregation for Deep Metric Learning

Overview: Typical Learning Concepts for DML Generalization

- Diverse Feature Aggregation
- **Multi Self-Distillation³**
- Proxy Learning
- ...



³ Roth, Milbich et al.; ICML 2021; S2SD: Simultaneous Similarity-based Self-distillation for Deep Metric Learning