

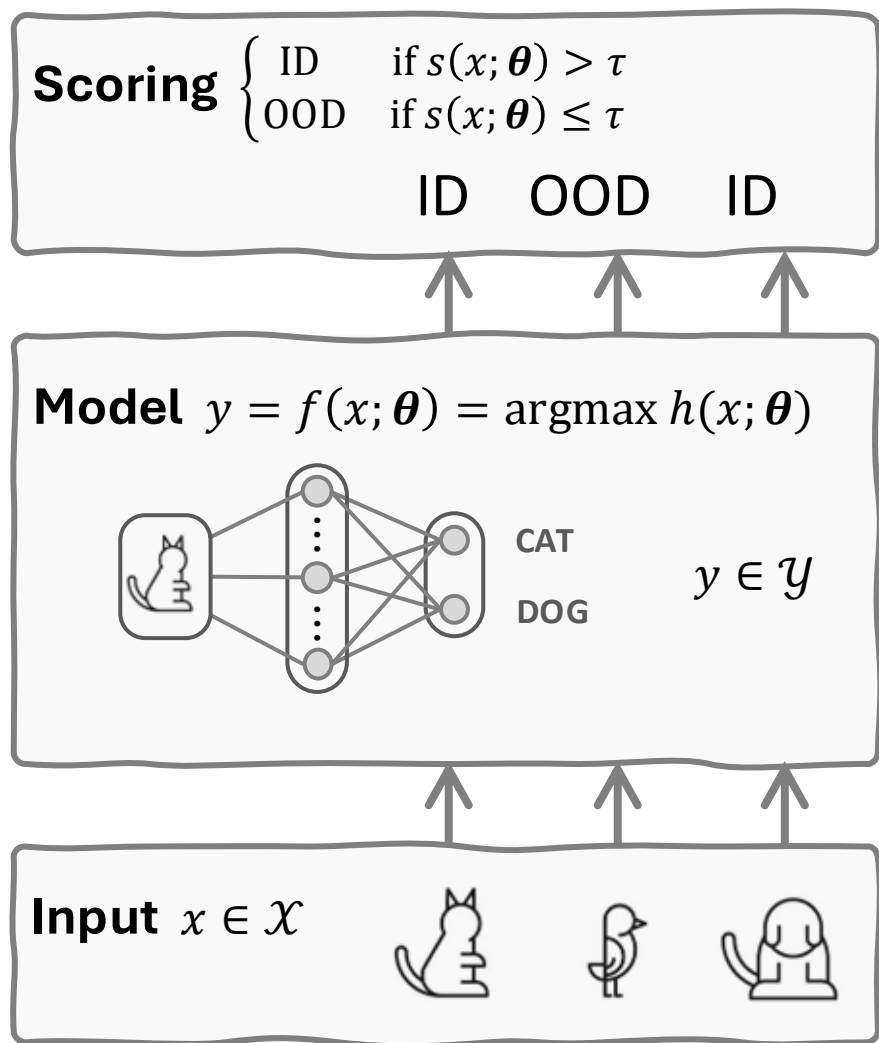
On the Insights and Strategies for OOD Detection Learning

Dr. Qizhou WANG
RIKEN AIP

<https://qizhouwang.github.io/homepage>



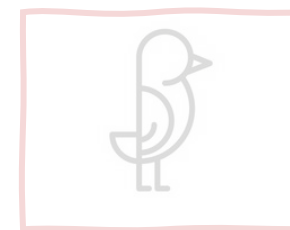
Post-hoc OOD Detection: Review



Semantic Shift: semantic relationship between inputs and labels changes.



ID

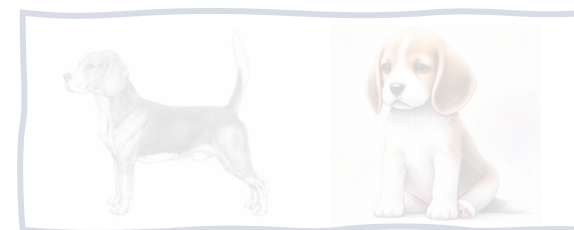


OOD

Covariate Shift: feature distribution changes while labels stay within the label space.



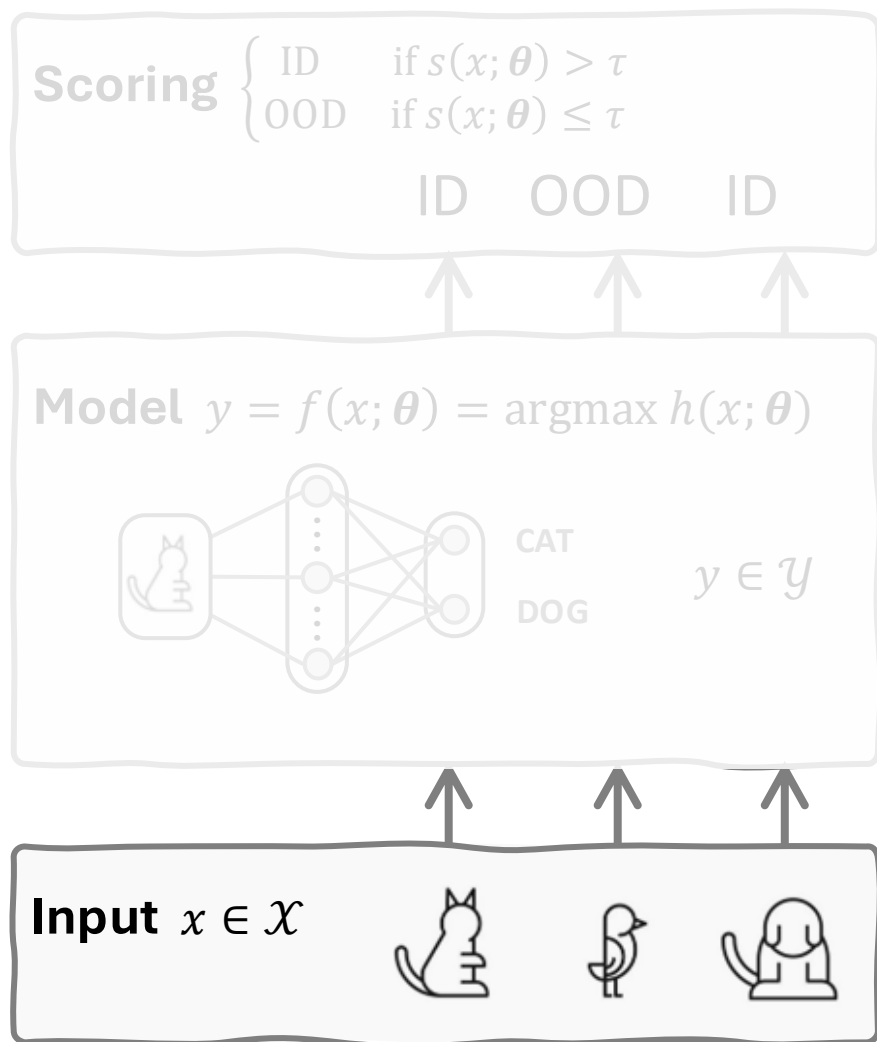
ID



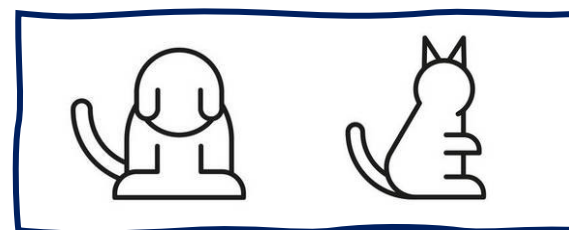
OOD

Post-hoc OOD Detection: Review

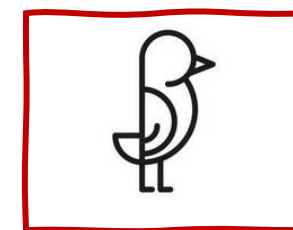
Interests of OOD detection



Semantic Shift: semantic relationship between inputs and labels changes.

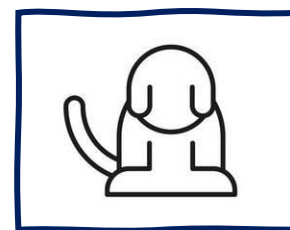


ID

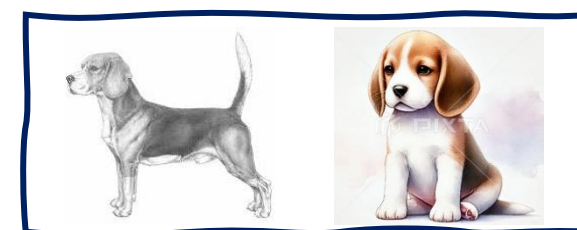


OOD

Covariate Shift: feature distribution changes while labels stay within the label space.

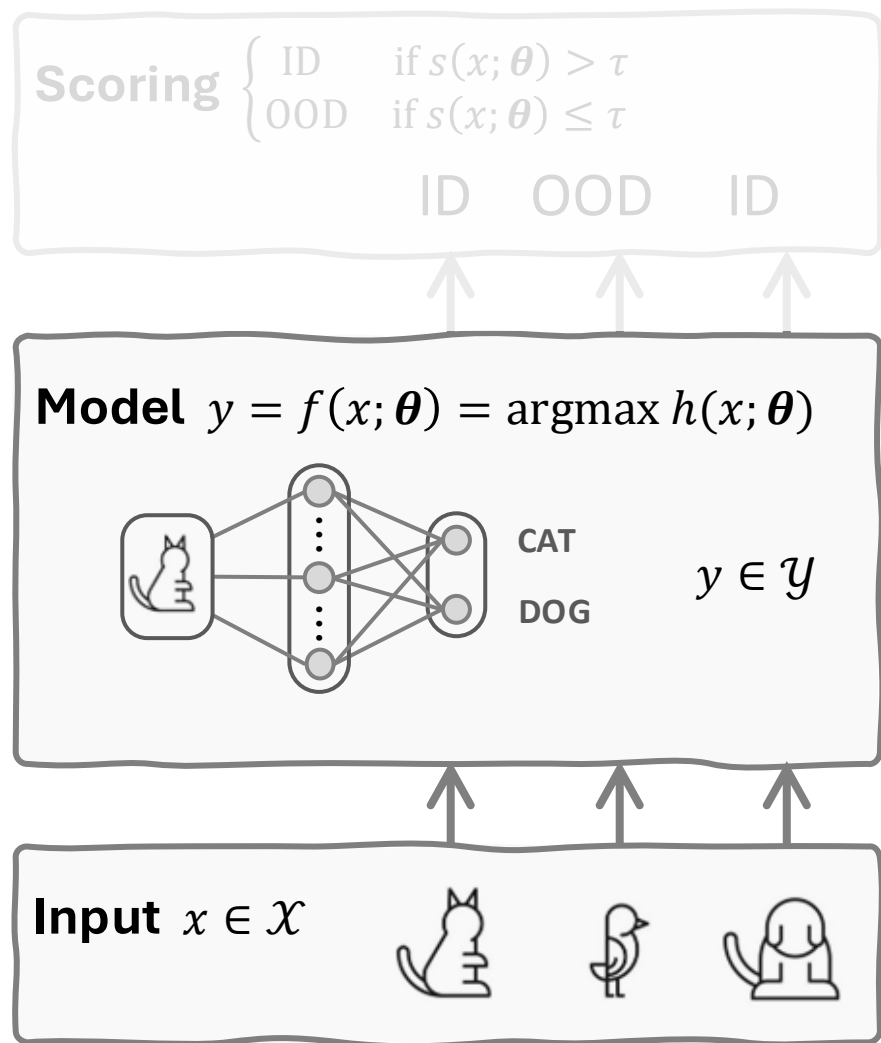


ID

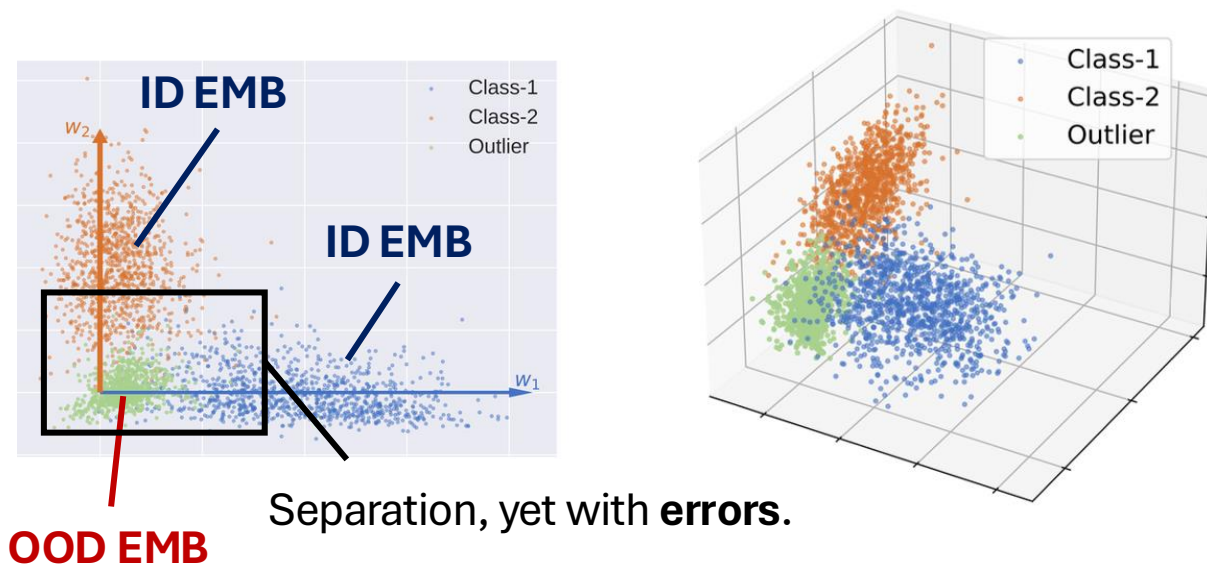


OOD

Post-hoc OOD Detection: Review

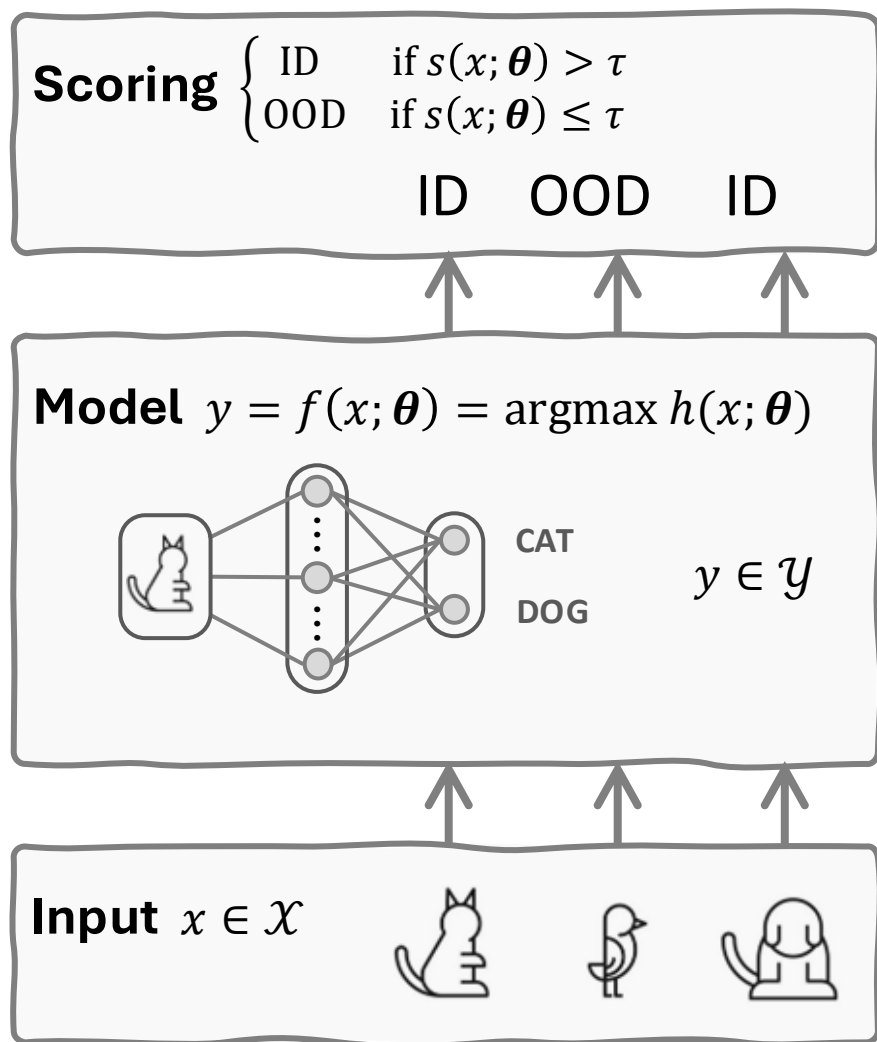


Pre-trained models can **separate ID and OOD data in embedding space** to some extent.



Figures of dimensional-reduced embeddings [a].

Post-hoc OOD Detection: Review



Model responses can be used to craft OOD scoring functions.

❖ Output Level, MSP [a]

maximal softmax prediction

$$s_{\text{MSP}}(x; \theta) = \max_k \text{softmax}_k h(x; \theta)$$

❖ Embedding Level, KNN [b]

k-th nearest neighbor

$$s_{\text{KNN}}(x; \theta) = \|h(x; \theta) - z_{(k)}\|_2$$

uniform distribution

❖ Gradient Level, GradNorm [c]

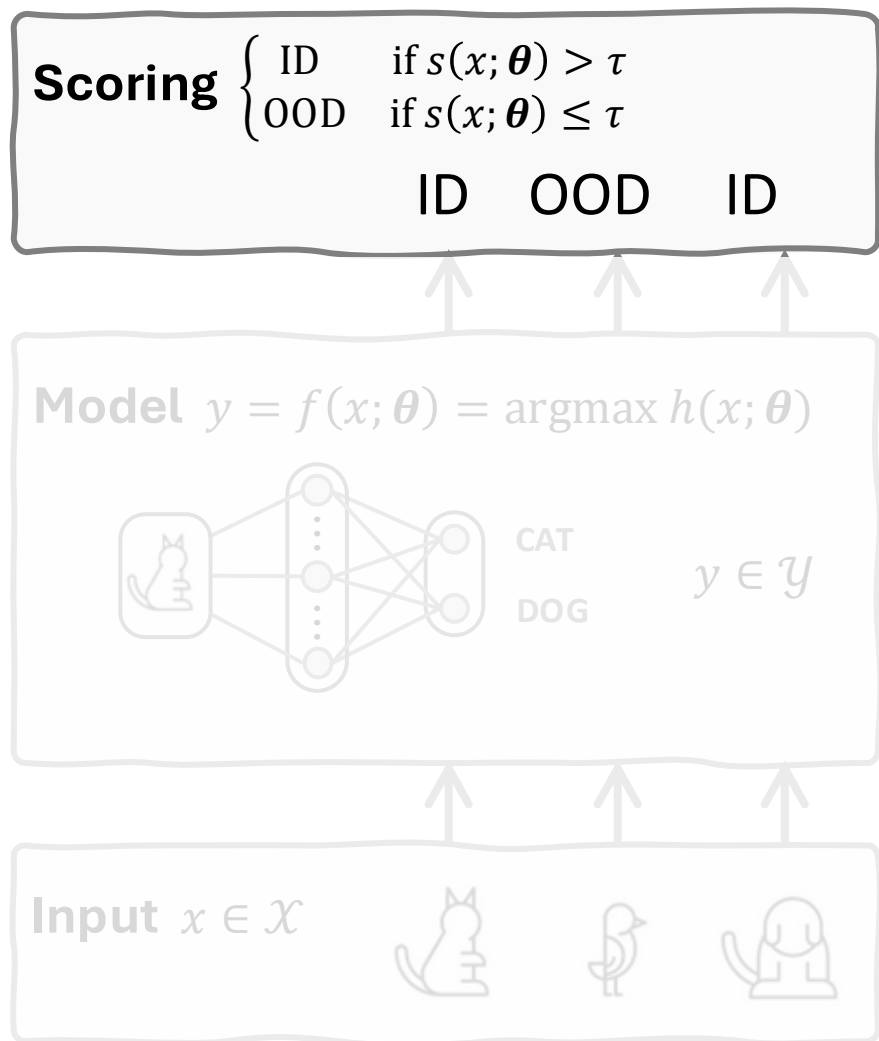
$$s_{\text{GN}}(x; \theta) = \|\nabla_{\theta} \text{KL}(u \| \text{softmax}(h(x; \theta)))\|_2$$

[a] Hendrycks et al. A Baseline for Detecting Misclassified and Out-of-Distribution Examples in Neural Networks. In ICLR, 2017.

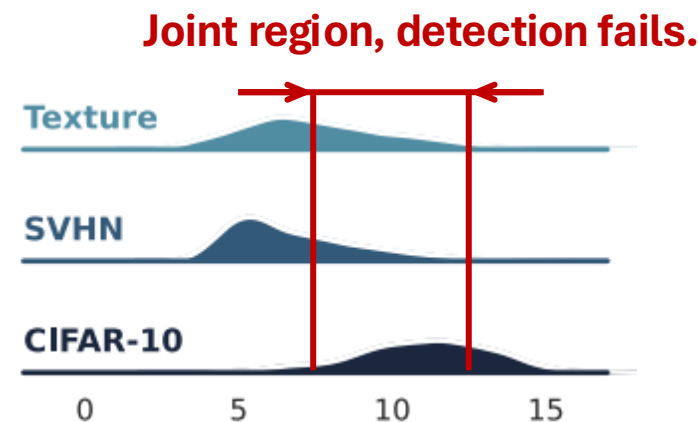
[b] Sun et al. Out-of-Distribution Detection with Deep Nearest Neighbors. In ICML, 2022.

[c] Huang et al. On the Importance of Gradients for Detecting Distributional Shifts in the Wild. In NeurIPS, 2021.

Post-hoc OOD Detection: Challenges



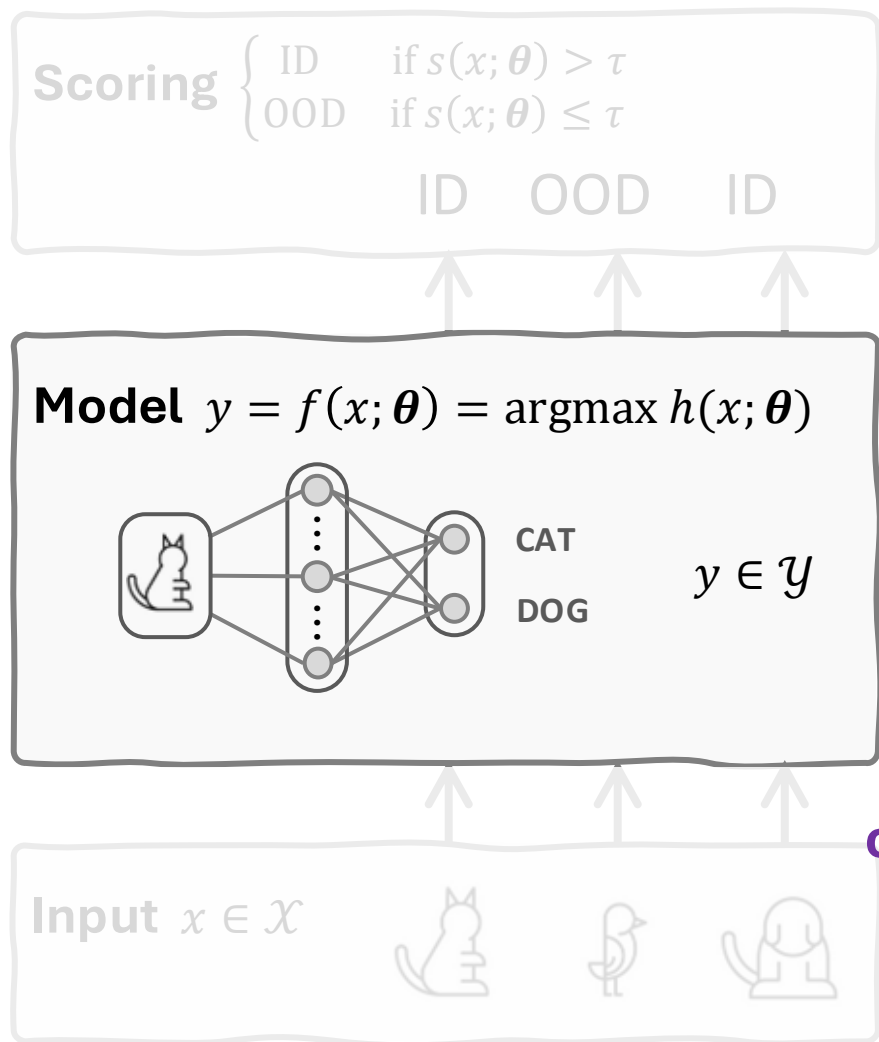
Post-hoc OOD detection often **makes mistakes**, failing to discern many ID and OOD patterns.



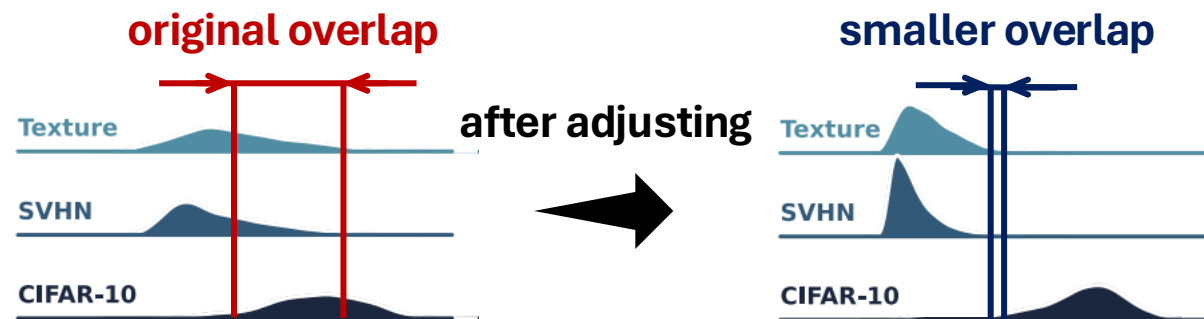
Figures of ID & OOD scoring distributions [a].

Explanations: For conventional-trained models, 1) their **representations** are not good enough, b) their **calibration** is inherently poor, and c) they cannot fully **classify** ID and OOD patterns.

OOD Detection Learning: What to Adjust?

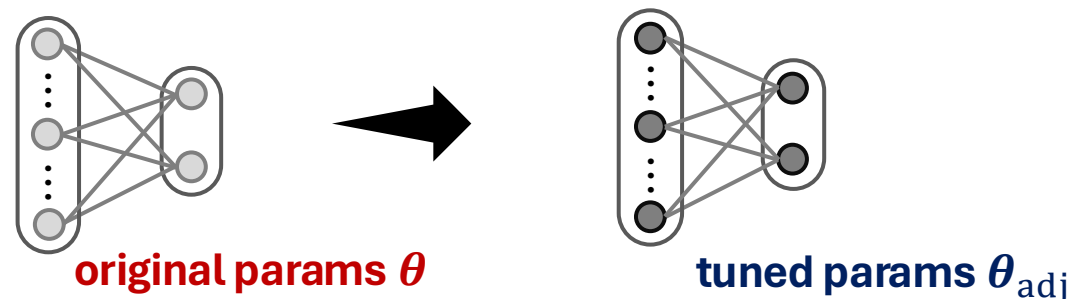


Adjust the system to improve OOD detection.

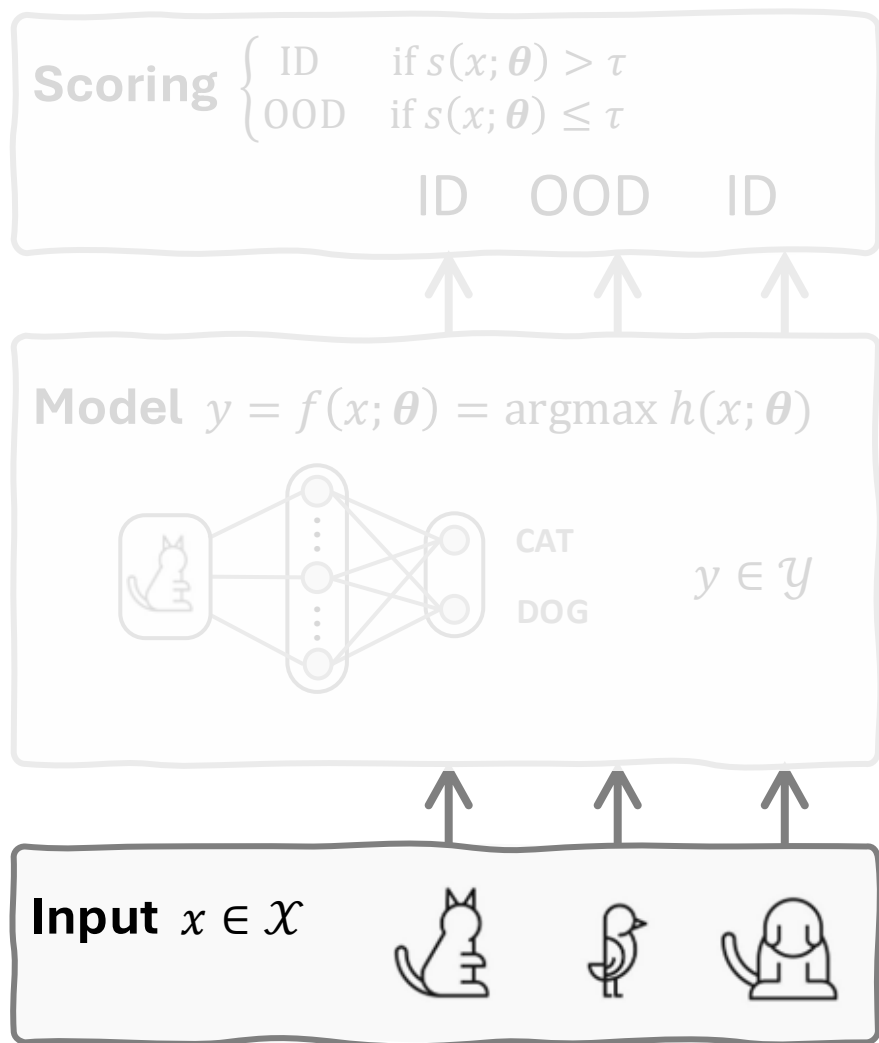


❖ **Model Level**, most works $s(x; \theta) \rightarrow s(x; \theta_{\text{adj}})$

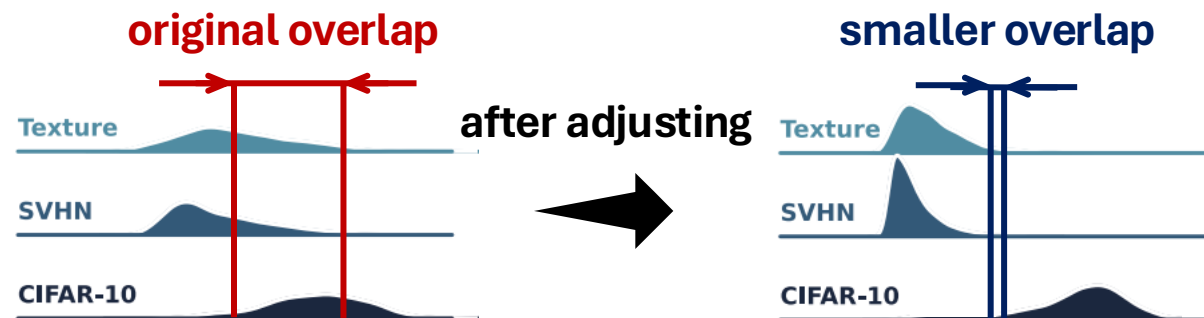
common choice



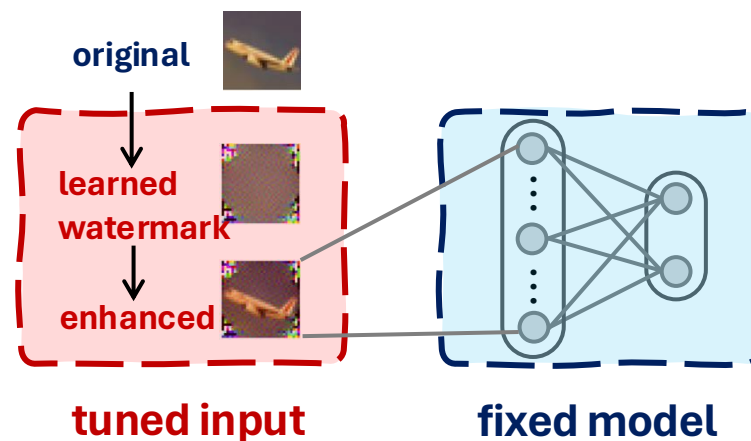
OOD Detection Learning: What to Adjust?



Adjust the system to improve OOD detection.

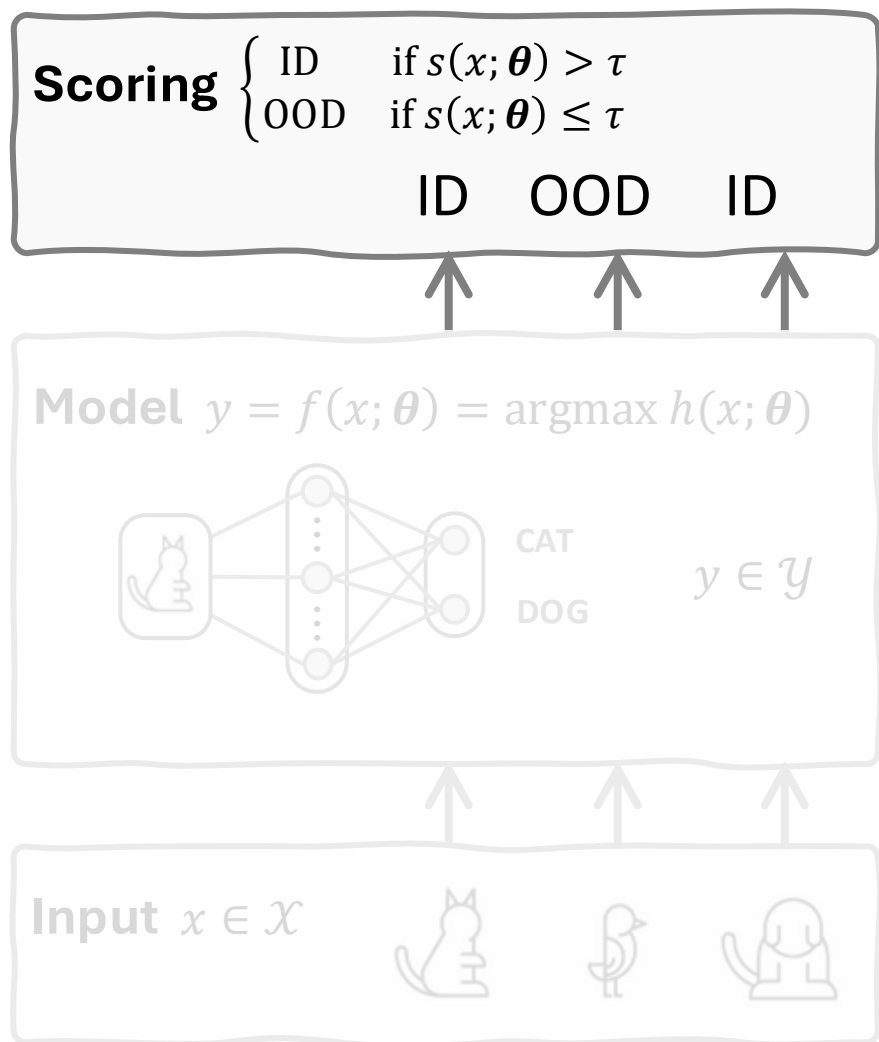


❖ **Input Level, WM** [a] $s(x; \theta) \rightarrow s(x + w; \theta)$

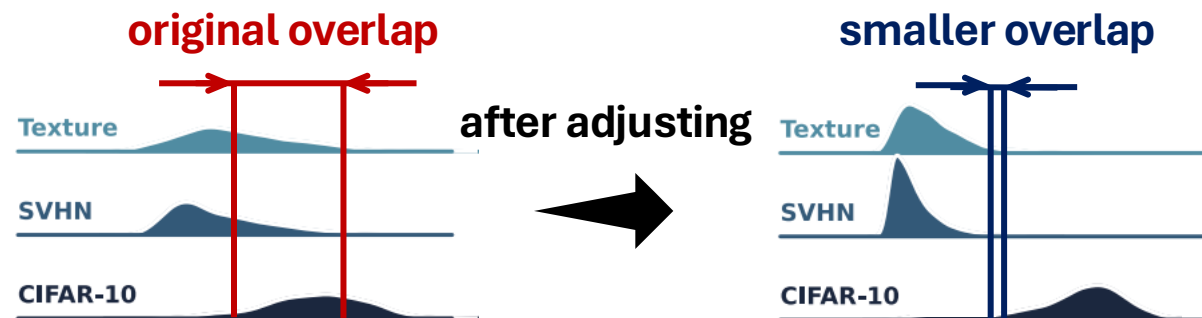


- ❖ *Watermark is static, **tuned** to enhance OOD detection.*
- ❖ *The pre-trained model remains **fixed**.*

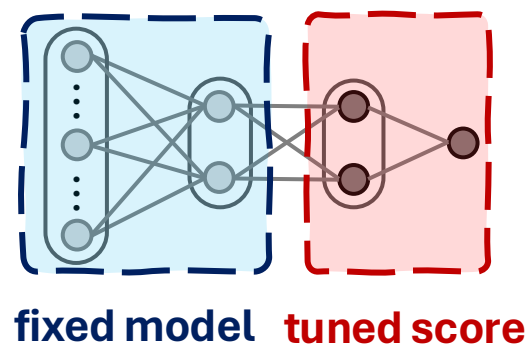
OOD Detection Learning: What to Adjust?



Adjust the system to improve OOD detection.



❖ **Score Level, VOS [a]** $s(x; \theta) \rightarrow s(h(x; \theta), \mathbf{w})$

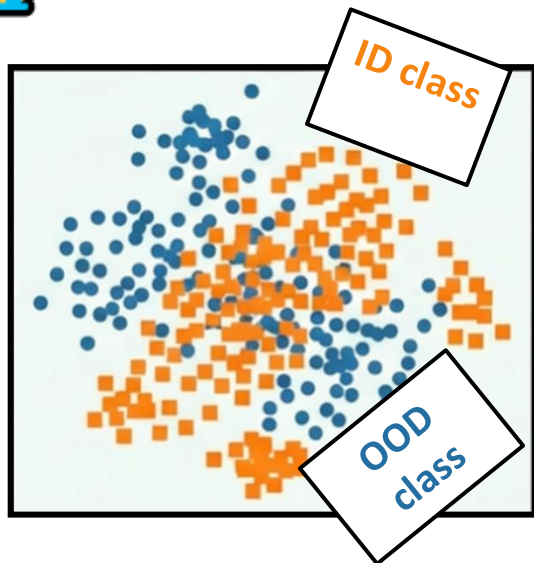


- ❖ The scoring function introduces exact params to be **tuned**.
- ❖ The pre-trained model remains **fixed**.

OOD Detection Learning: How to Adjust?

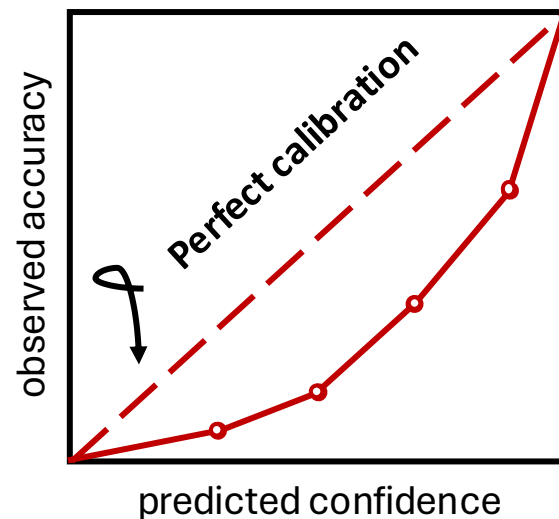
Let's recall **the drawbacks of post-hoc OOD detection**: For conventional-trained models, they have

Poor Representation



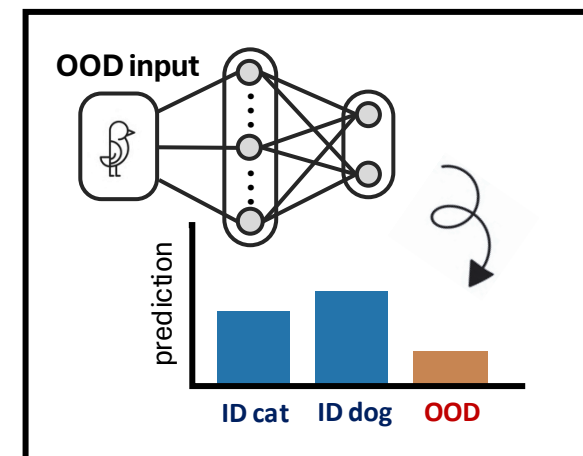
Data with **different semantics** may **not be perfectly separated** in the embedding space.

Poor Calibration



High model confidence does not correspond to **high model accuracy**.

Poor Classification

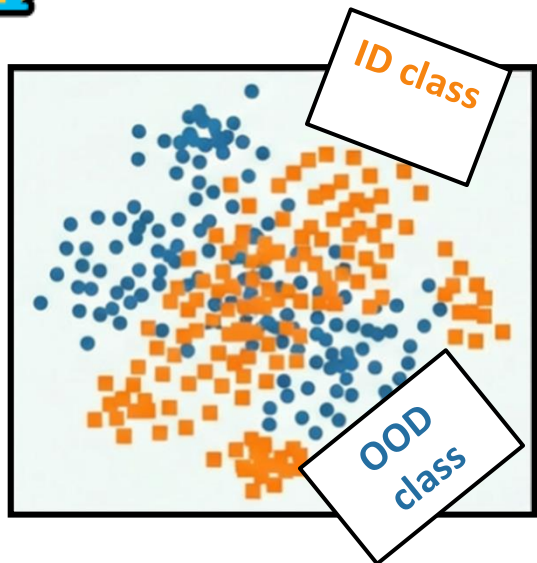


The model **predicts the wrong class**, despite the true class of being either ID or OOD.

OOD Detection Learning: How to Adjust?

Let's recall the drawbacks of post-hoc OOD detection: For conventional-trained models, they have

Poor Representation

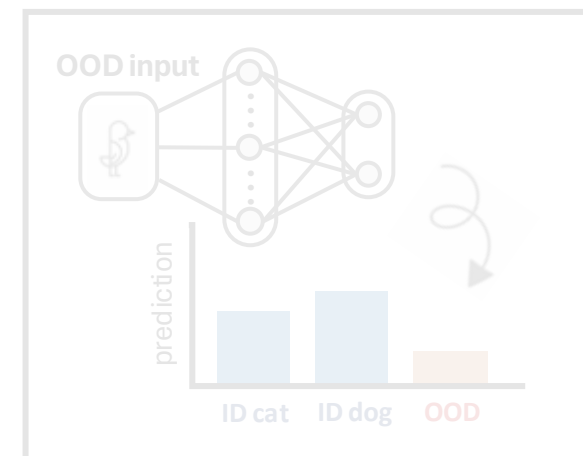


Data with **different semantics** may **not be perfectly separated** in the embedding space.

Poor Calibration

- ❖ **What happens?**
ID and OOD examples are **entangled** in the embedding space.
- ❖ **Why breaks detection?**
Many methods (e.g., k-nearest neighbors and Mahalanobis) **assume that OOD lies away from ID**.

Poor Classification



The model **predicts the wrong class**, despite the true class of being either ID or OOD.

OOD Detection Learning: How to Adjust?

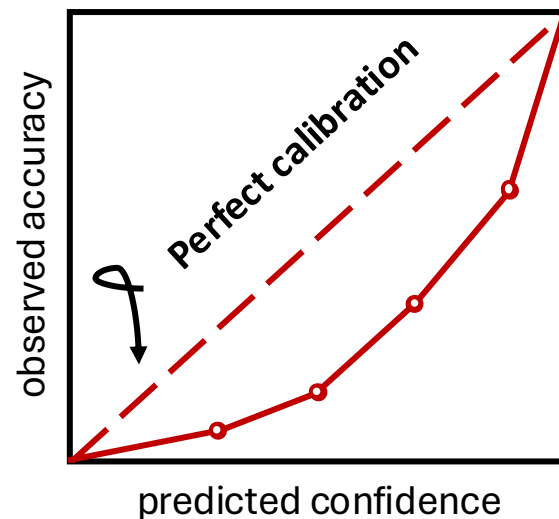
Let's recall the drawbacks of post-hoc OOD detection: For conventional-trained models, they have

Poor Representation



Data with **different semantics** may **not be perfectly separated** in the embedding space.

Poor Calibration



High model confidence does not correspond to **high model accuracy**.

Poor Classification

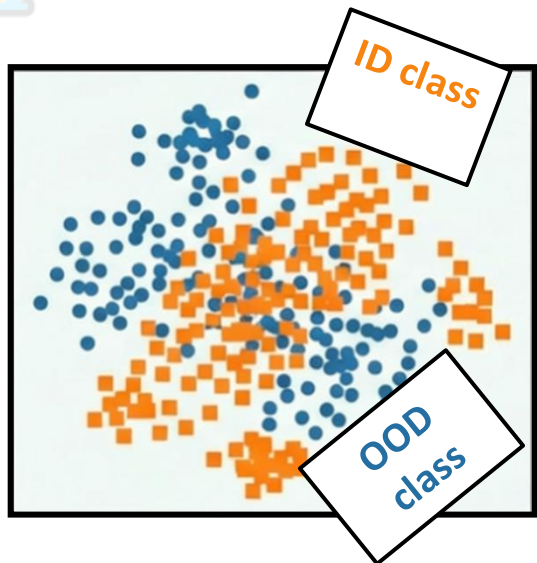
- ❖ **What happens?**
Models produce **high confidence** for **wrong predictions**.
- ❖ **Why breaks detection?**
Many methods (e.g., MSP) use **confidence-like scores**.

The model **predicts the wrong class**, despite the true class of being either ID or OOD.

OOD Detection Learning: How to Adjust?

Let's recall the drawbacks of post-hoc OOD detection: For conventional-trained models, they have

Poor Representation



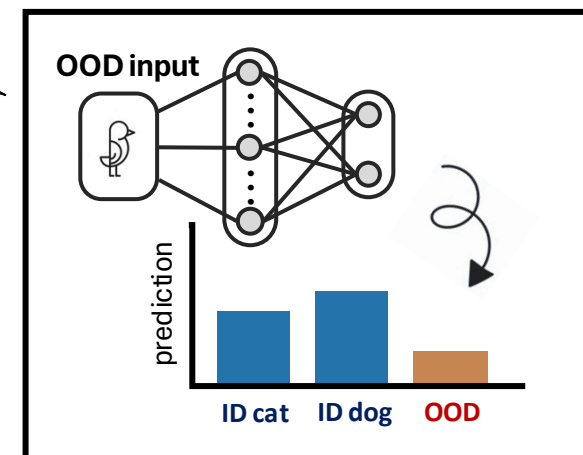
Data with **different semantics** may **not be perfectly separated** in the embedding space.

Poor Calibration

- ❖ **What happens?**
Decision boundaries does not align with true ID/OOD classes.
- ❖ **Why breaks detection?**
Taking as an **extra classification task**, this classifier is not accurate.

High model confidence does not correspond to **high model accuracy**.

Poor Classification

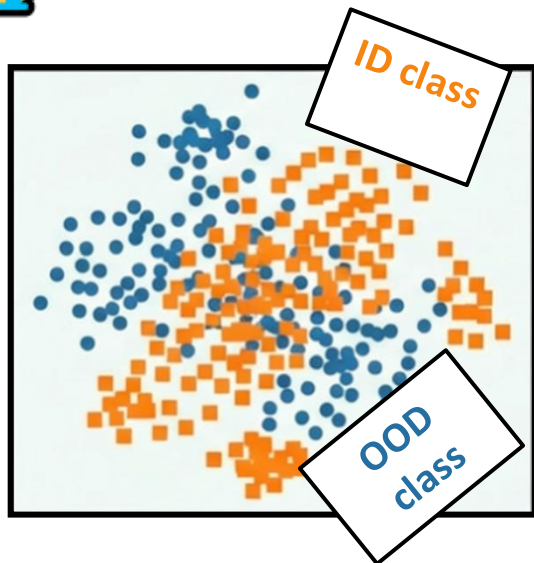


The model **predicts the wrong class**, despite the true class of being either ID or OOD.

OOD Detection Learning: How to Adjust?

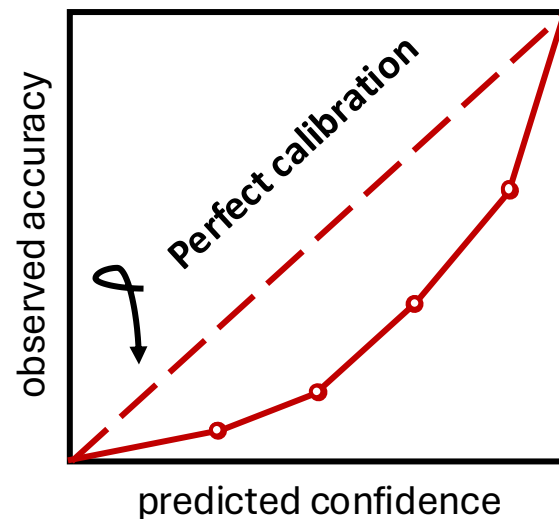
Let's recall the drawbacks of post-hoc OOD detection: For conventional-trained models, they have

Poor Representation



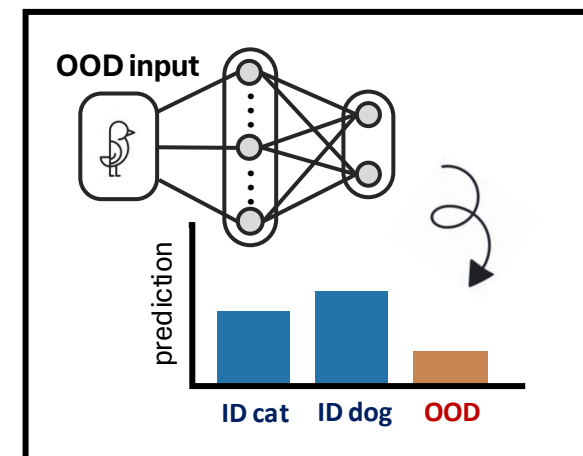
Data with **different semantics** may **not be perfectly separated** in the embedding space.

Poor Calibration



High model confidence does not correspond to **high model accuracy**.

Poor Classification

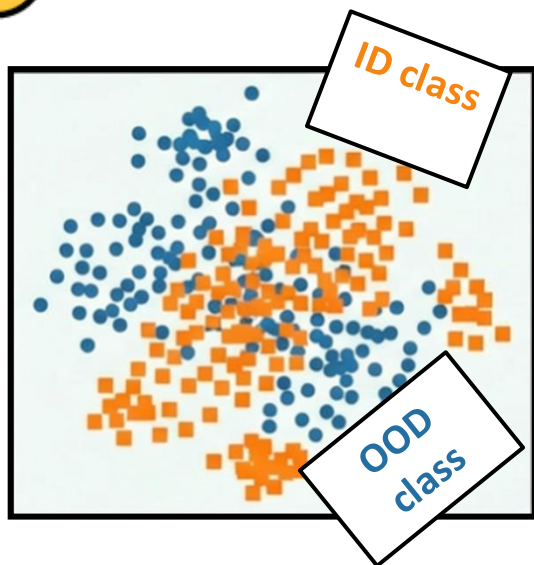


The model **predicts the wrong class**, despite the true class of being either ID or OOD.

OOD Detection Learning: How to Adjust?

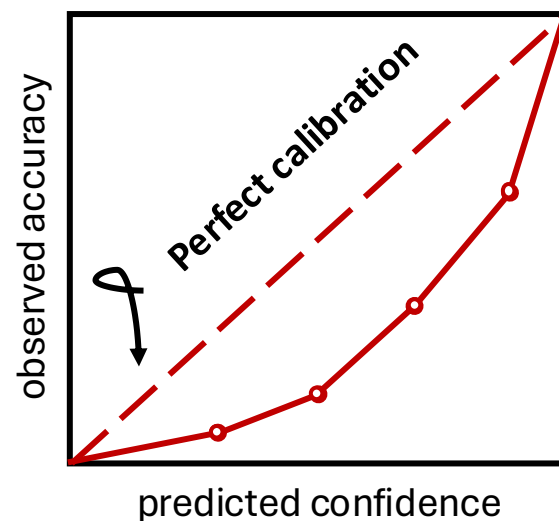
To **address the drawbacks** of post-hoc OOD detection, OOD detection learning can

Improve Representation



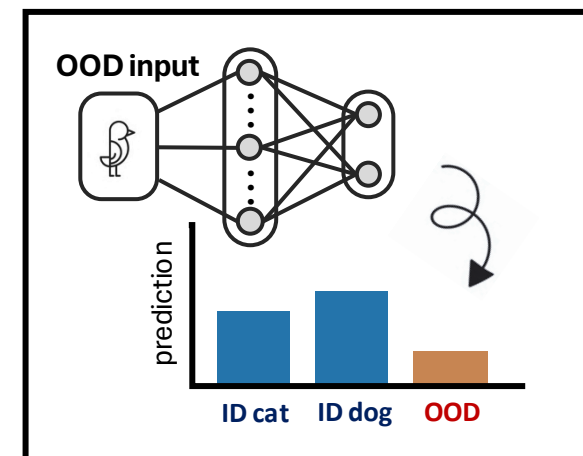
- ❖ **contrastive learning,**
- ❖ **reconstruction learning,**
- ❖ **pre-training,** et al.

Improve Calibration



- ❖ **Bayesian scoring,**
- ❖ **density regularization,**
- ❖ **calibration,** et al.

Improve Classification



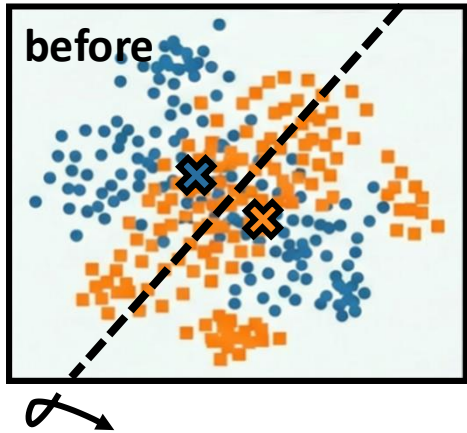
- ❖ **outlier exposure,**
- ❖ **data augmentation,**
- ❖ **sample selection,** et al.

Representation: Overview

Conventional-trained Classifiers



Representation-based OOD Learning



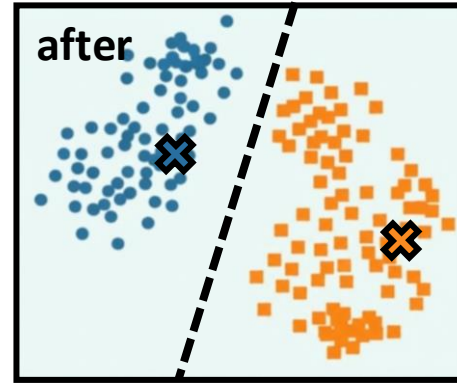
ID/OOD boundary

A simple OOD scoring: Nearest distances to K-means clustering.

$$s_{\text{KM}}(x; \theta) = \|h(x; \theta) - \mu_{(x)}\|_2$$

nearest ID centroid

- ❖ make **no assumptions about the form of OOD data** will take or the type of downstream task.
- ❖ mainly learn patterns among training classes as a **shortcut to learn classification**.



Representations with different semantics are better separated.

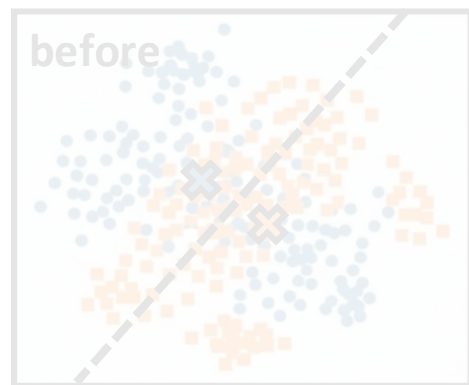
- ❖ low intra-class variance
- ❖ high inter-class variance

Pretext tasks

- ❖ **Contrastive Learning:** CSI, SSD
- ❖ **Reconstruction Learning:** MOOD
- ❖ **Pre-training:** CLIP

Representation: Overview

Conventional-trained Classifiers



A simple OOD scoring : Nearest distances to K-means clustering.

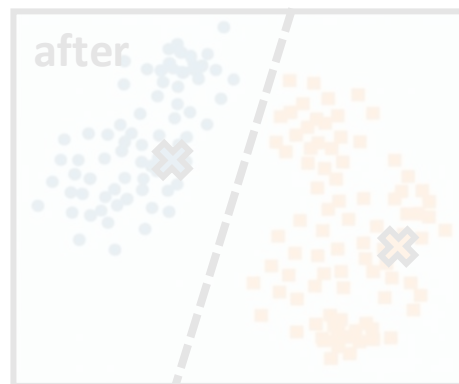
$$s_{\text{KM}}(x; \theta) = \|h(x; \theta) - \mu_{(x)}\|_2$$

nearest ID centroid

ID/OOD boundary

- ❖ make **no assumptions** about the form of OOD data will take or the type of downstream task.
- ❖ mainly learn patterns among training classes as a **shortcut to learn classification**.

Representation-based OOD Learning



Representations with different semantics are better separated.

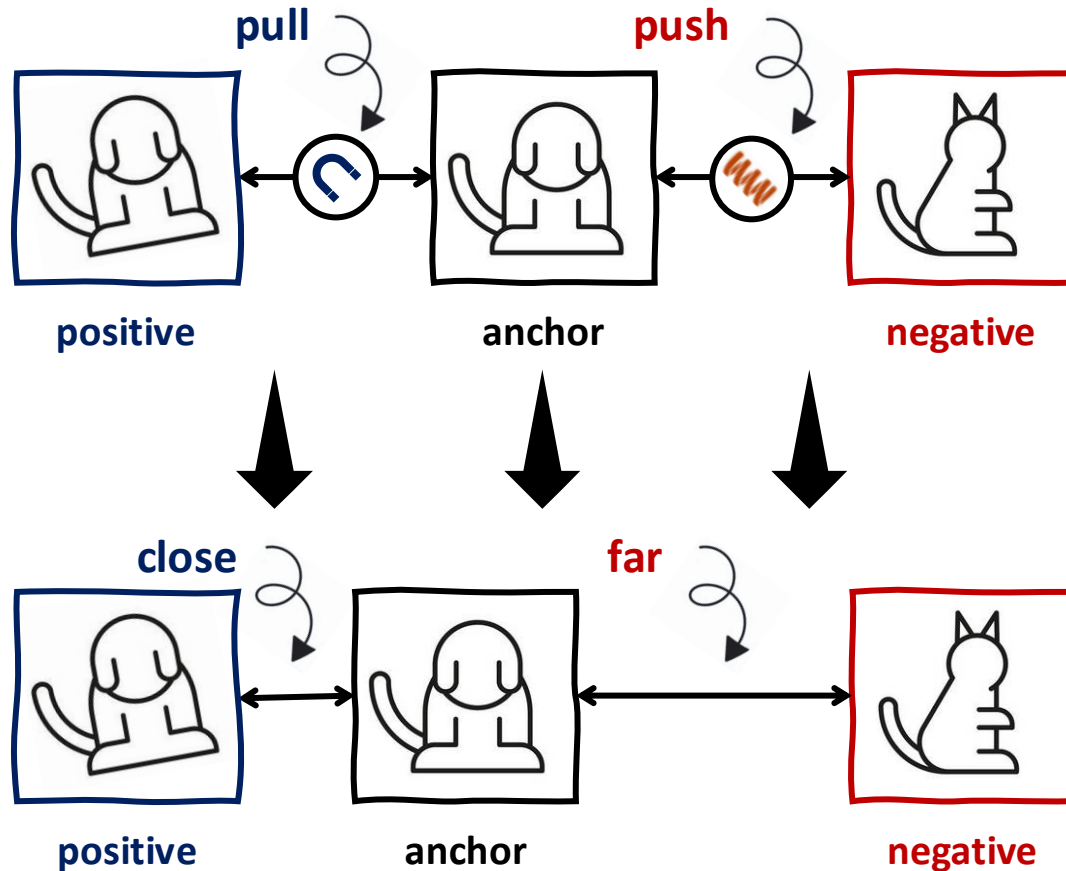
- ❖ low intra-class variance
- ❖ high inter-class variance

Pretext tasks

- ❖ **Contrastive Learning:** CSI, SSD
- ❖ **Reconstruction Learning:** MOOD
- ❖ **Pre-training:** CLIP

Representation: Contrastive Learning

Contrastive learning **improves the semantic structure of the embedding space** by **pulling** semantically similar samples together and **pushing** dissimilar ones apart.



Supervised Contrastive Learning

$$\mathcal{L}_{\text{con}}(x, \{x_+\}, \{x_-\}) = -\frac{1}{|\{x_+\}|} \log \frac{\sum_{x' \in \{x_+\}} \exp\{\text{sim}(z(x), z(x'))\}}{\sum_{x' \in \{x_+\} \cup \{x_-\}} \exp(\text{sim}(z(x), z(x'))))$$

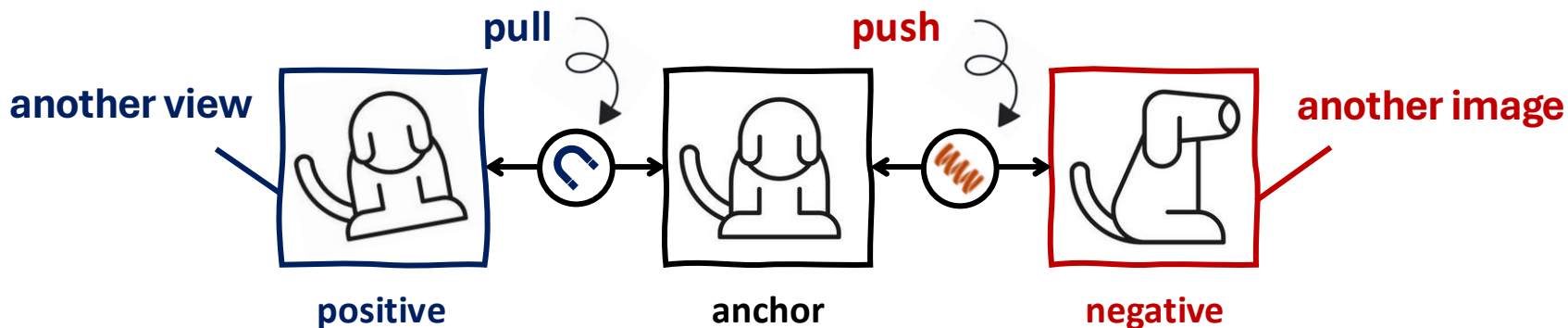
The equation is annotated with labels: 'positive' points to the set $\{x_+\}$, 'anchor' points to the term $z(x)$, and 'negative' points to the set $\{x_-\}$.

Increasing similarity between the anchor and positive samples, while *decreasing similarity* to negative samples.

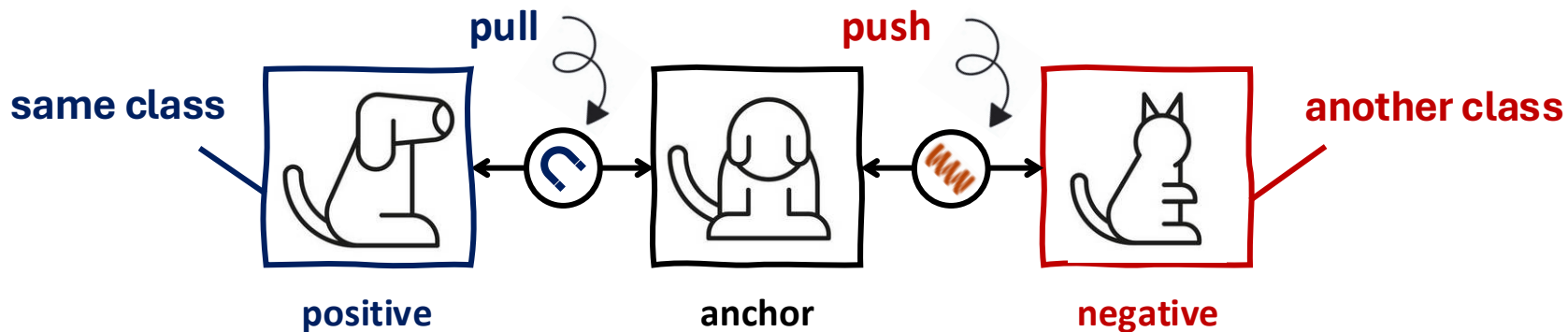
Representation: Contrastive Learning—SSD

SSD [a] Training Objective.

❖ **SSD**, discriminating between **individuals** (label-agnostic).



❖ **SSD+**, discriminating between **classes** (label-aware).



Representation: Contrastive Learning—SSD

SSD OOD Detector.

❖ **Mahalanobis**, cluster-conditioned detection (**OOD-agnostic**).

$$s_{\text{mah}}(x; \theta) = \min_k (z(x) - \mu_k)^\top \Sigma_k^{-1} (z(x) - \mu_k)$$

k-th ID cluster centroid

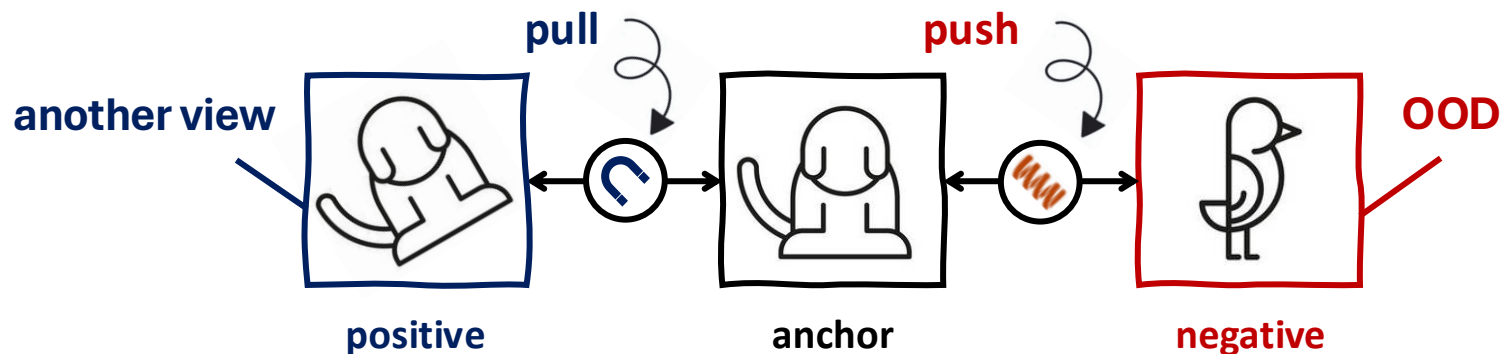
k-th ID cluster covariance, decorrelating features

❖ **Mahalanobis+**, a small OOD set is available before inference (**few-shot OOD**).

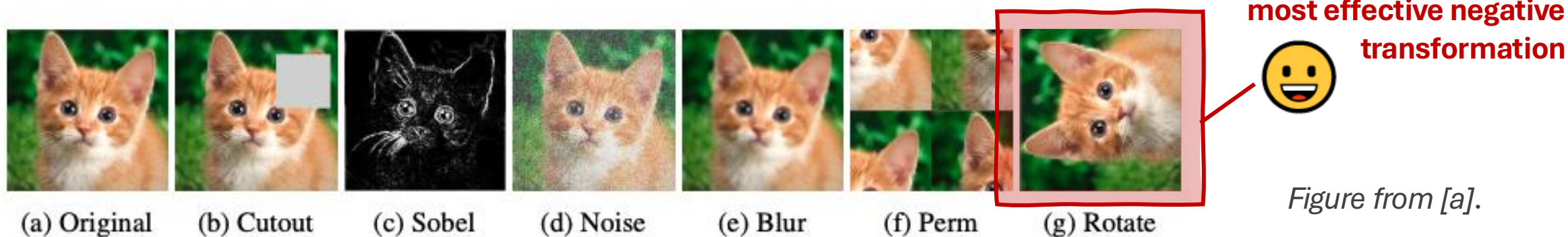
$$s_{\text{mah}+}(x; \theta) = \underbrace{(z(x) - \mu_{id})^\top \Sigma_{id}^{-1} (z(x) - \mu_{id})}_{\text{distance to overall ID centroid}} - \underbrace{(z(x) - \mu_{ood})^\top \Sigma_{ood}^{-1} (z(x) - \mu_{ood})}_{\text{distance to overall OOD centroid}}$$

Representation: Contrastive Learning—CSI

CSI [a] further studies effective **negative transformation** to discriminate OOD samples.

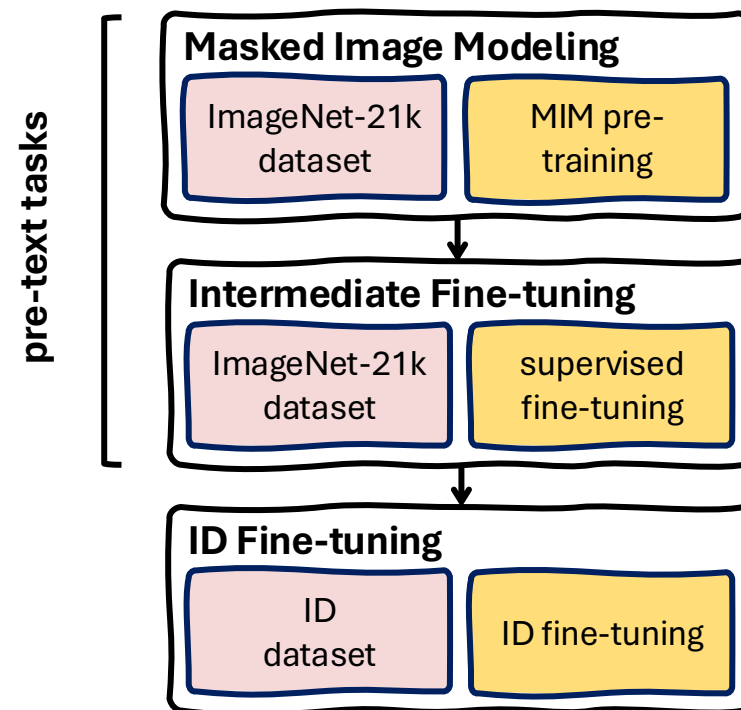
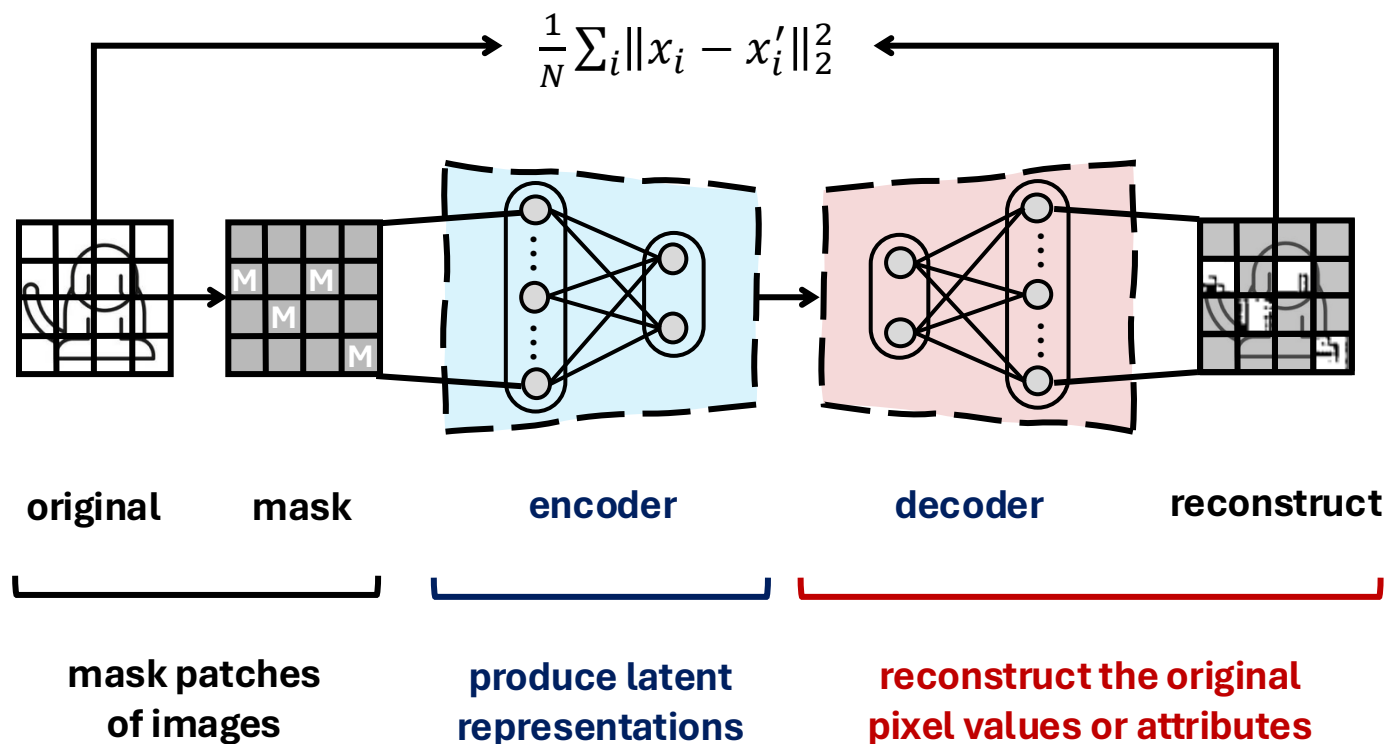


What if no OOD data are available? Those transformations previously found **ineffective as positive transformations** can instead be used as negative transformations.



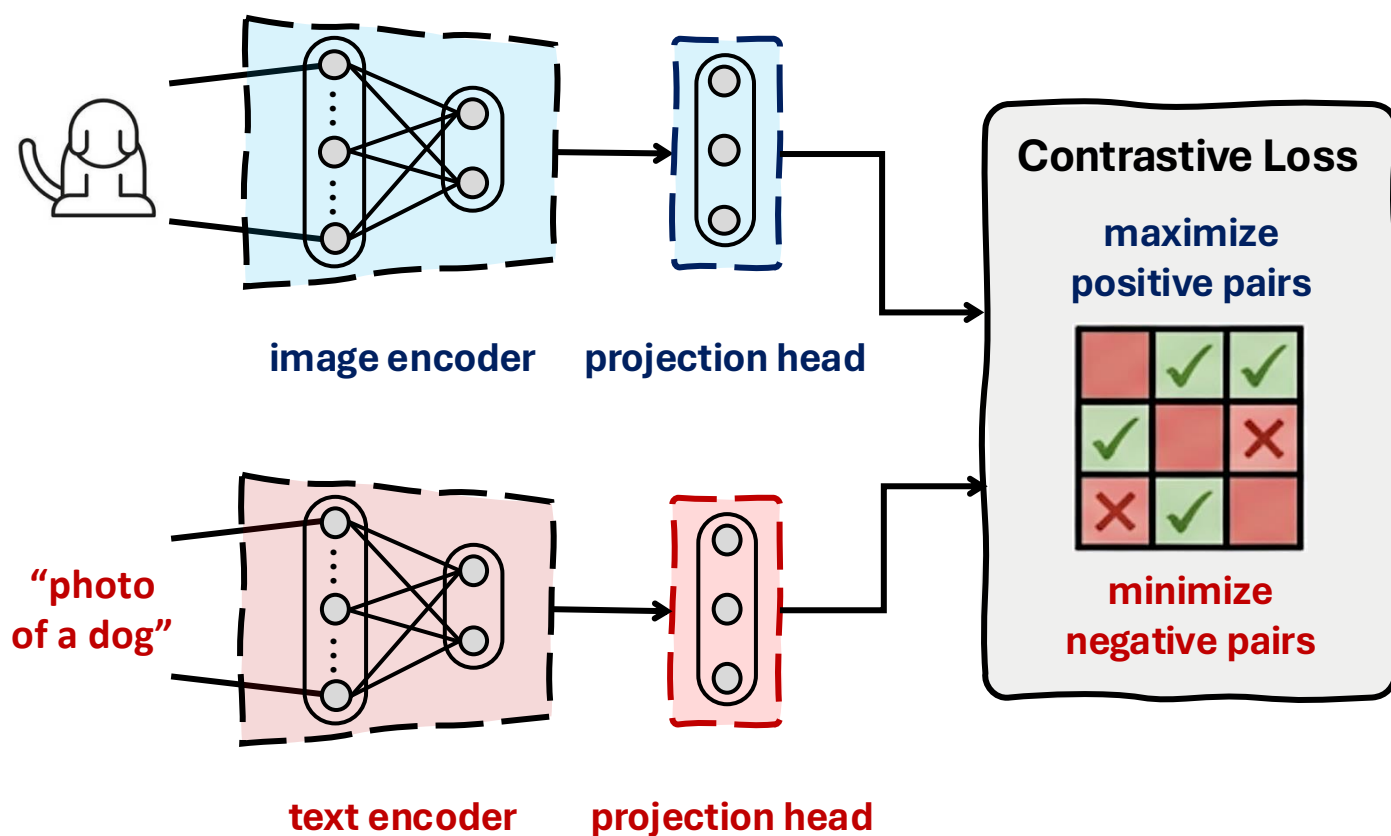
Representation: Reconstruction Learning—MOOD

MOOD [a] forces the models to **encode finer-grained information needed to rebuild input**, which is far beyond semantics as in contrastive learning.

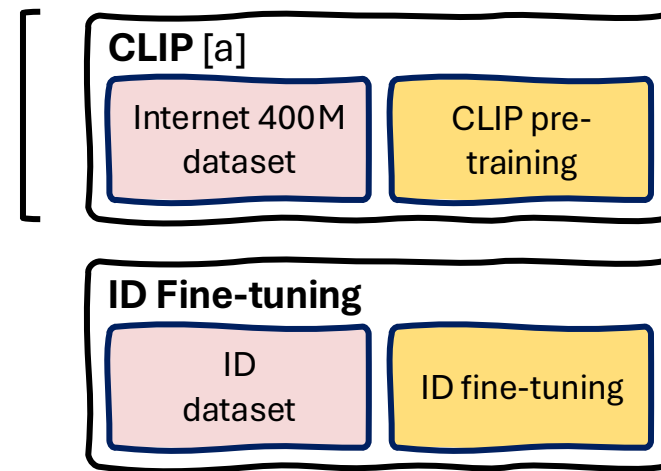


Representation: Pre-training—CLIP

Cross-modal alignment reduces reliance on **the shortcut to learn classification**, encouraging more general, semantically meaningful features.



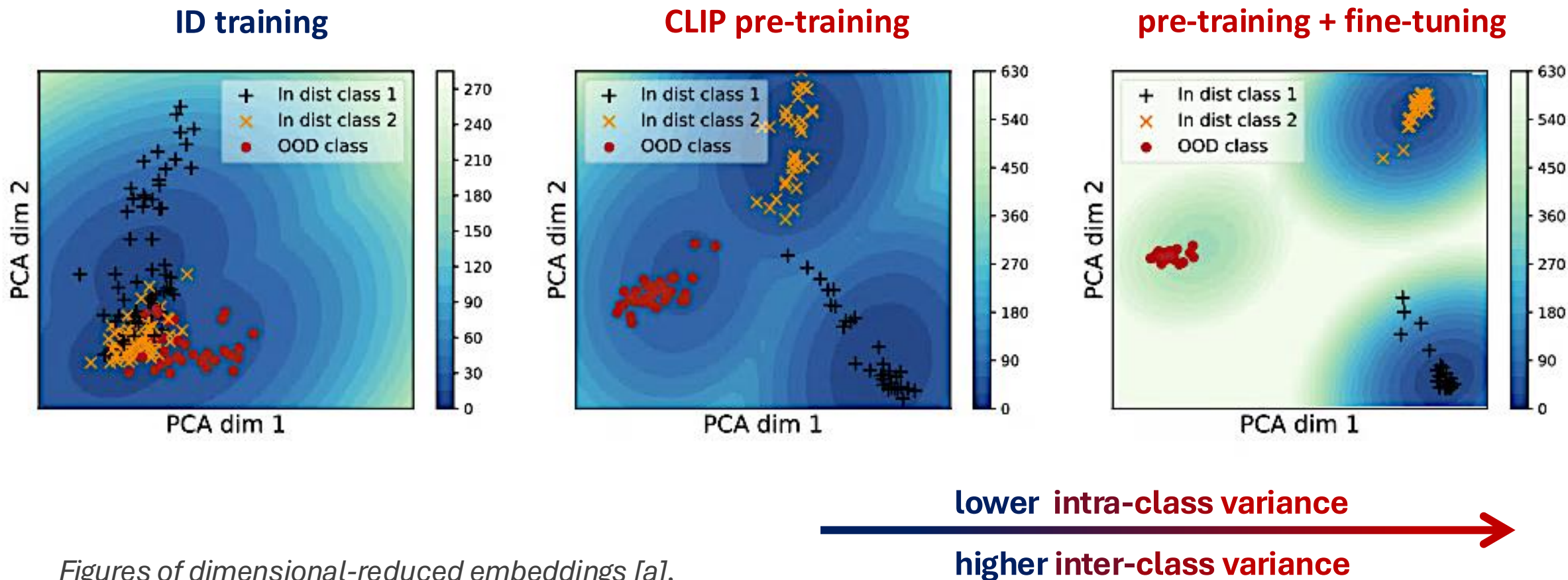
pre-text task



Aligns **images and text representations** by **maximizing similarity of correct pairs**, meanwhile **minimizing similarity of incorrect pairs**.

Representation: Pre-training—CLIP

Cross-modal alignment reduces reliance on **the shortcut to learn classification**, encouraging more general, semantically meaningful features.



Figures of dimensional-reduced embeddings [a].

Calibration: Overview

Let us review what we have learned from textbook [a].

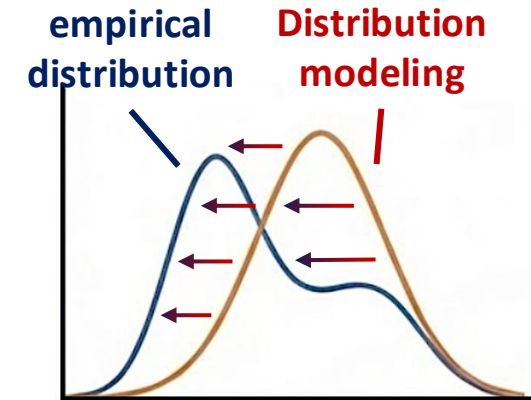
❖ Maximum Likelihood & KL Divergence

MLE minimizes KL between the **distribution model** and **empirical data distribution**.

$$p^* = \underset{p \in \mathcal{P}}{\operatorname{argmax}} \sum_i \log p(x_i) = \underset{p \in \mathcal{P}}{\operatorname{argmin}} \operatorname{KL}(\widehat{\mathcal{D}} || p)$$

optimal solution

Condition 1. proper distribution family **Condition 2.** enough data



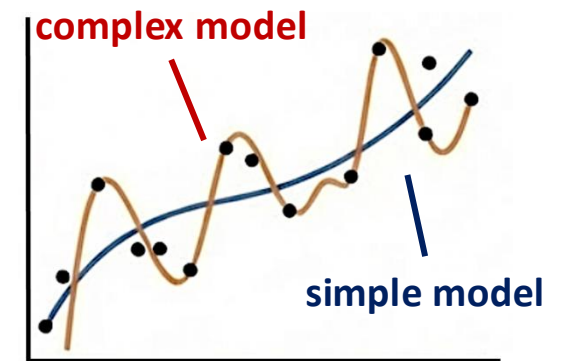
❖ Model Complexity & Estimation Error

Higher complexity increases estimation errors for a fixed dataset size.

$$\left\| \mathbb{E}_{x \sim \mathcal{D}} [\Phi(x)] - \mathbb{E}_{x \sim \widehat{\mathcal{D}}} [\Phi(x)] \right\|_{\infty}$$

estimation error

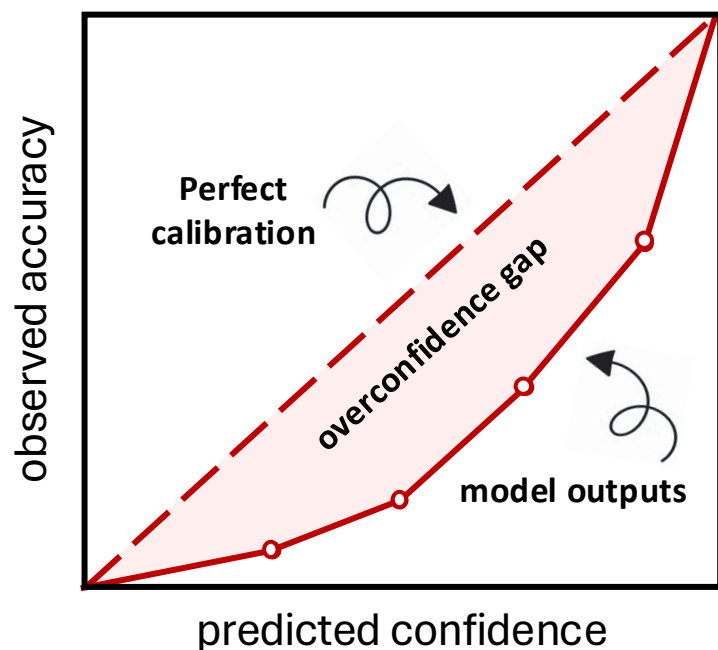
$$\leq \underbrace{2\mathfrak{R}_m(\mathcal{H})}_{\text{Condition 3. proper model family}} + r \sqrt{\frac{\log 2/\delta}{2m}}$$



Calibration: Overview

Conditions 1-3 may not be fully satisfied in practice, leading to calibration failures.

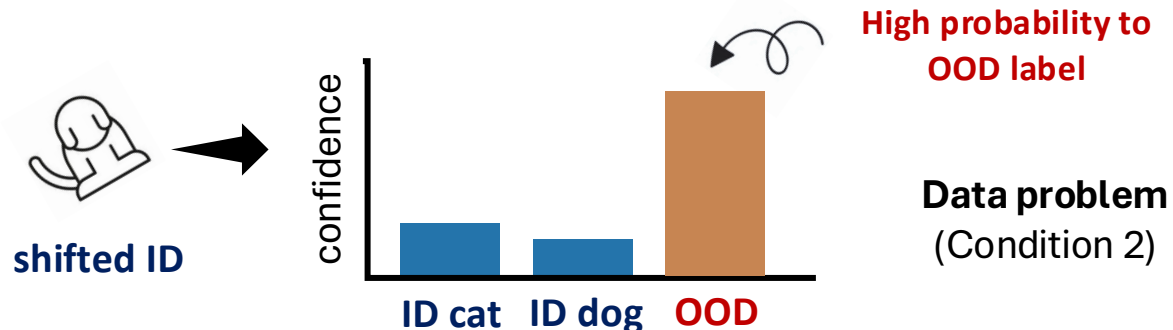
The overconfidence problem



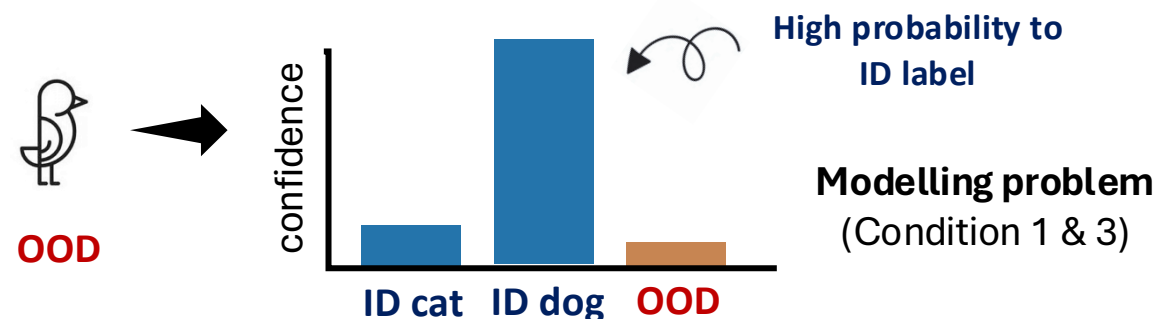
Overconfidence makes the model outputs **appear more reliable than they really are.**

Consequences

Mis-confidence on shifted ID



Mis-confidence on OOD



Calibration: Overview

Conditions 1-3 may not be fully satisfied in practice, leading to calibration failures.

The overconfidence problem

❖ Data-centric Solutions.

conventional calibration methods, such as label smoothing and mixup augmentations.

❖ Model-centric Solutions.

various regularization strategies.

❖ Distribution-centric Solutions.

modelling beyond softmax.

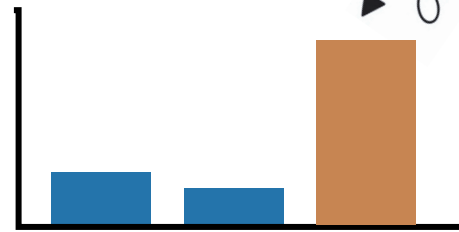
Consequences

Mis-confidence on shifted ID



shifted ID

confidence



High probability to OOD label

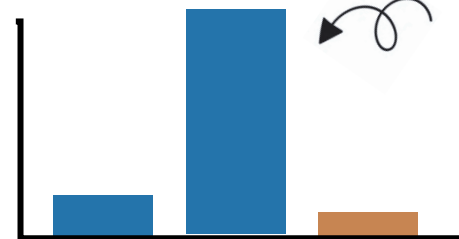
Data problem
(Condition 2)

Mis-confidence on OOD



OOD

confidence



High probability to ID label

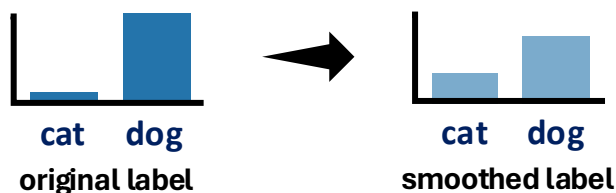
Modelling problem
(Condition 1 & 3)

Calibration: Data-centric Solutions

Conventional calibration strategies have been empirically shown to enhance OOD detection.

❖ Change Labels

label smoothing [a]



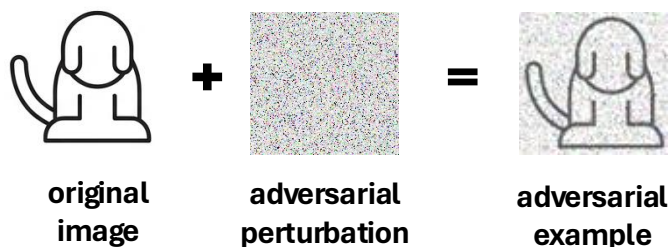
$$\tilde{y} = \overbrace{(1 - \alpha)y + \frac{\alpha}{c}}^{\text{interpolation}}$$

original label uniform

replaces one-hot labels with **slightly softened targets**.

❖ Change Inputs

adversarial examples [b]



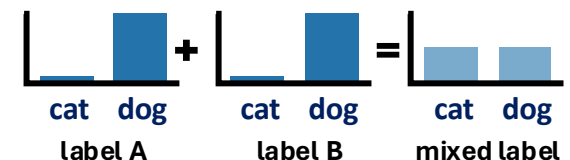
$$\max_{\delta \in \Delta} \mathcal{L}(x + \delta, y, \theta)$$

fail model perturb input

perturbs inputs **imperceptibly** to fool the models and then train the model.

❖ Change Labels and Inputs

mixup [c]



$$\tilde{x} = \lambda x_A + (1 - \lambda)x_B$$
$$\tilde{y} = \lambda y_A + (1 - \lambda)y_B$$

creates virtual training examples that **interpolates** between data pairs.

[a] Li et al. Rethinking Out-of-distribution (OOD) Detection: Masked Image Modeling is All You Need. In CVPR, 2023.

[b] Botschen et al. Out-of-Distribution Detection with Adversarial Outlier Exposure. In CVPR, 2025.

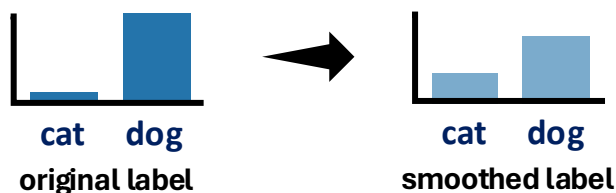
[c] Zhang et al. *mixup*: Beyond Empirical Risk Minimization. In ICLR, 2018.

Calibration: Data-centric Solutions

Conventional calibration strategies have been empirically shown to enhance OOD detection.

❖ Change Labels

label smoothing [a]



interpolation

$$\tilde{y} = (1 - \alpha)y + \frac{\alpha}{c}$$

original label uniform

replaces one-hot labels with **slightly softened targets**.

❖ Change Inputs

adversarial examples [b]



discourages extreme 0/1 probabilities and thus reduces overconfidence and improves calibrations.

perturbs inputs **imperceptibly** to fool the models and then train the model.

❖ Change Labels and Inputs

mixup [c]



$$\tilde{x} = \lambda x_A + (1 - \lambda)x_B$$

$$\tilde{y} = \lambda y_A + (1 - \lambda)y_B$$

creates virtual training examples that **interpolates** between data pairs.

[a] Li et al. Rethinking Out-of-distribution (OOD) Detection: Masked Image Modeling is All You Need. In CVPR, 2023.

[b] Botschen et al. Out-of-Distribution Detection with Adversarial Outlier Exposure. In CVPR, 2025.

[c] Zhang et al. *mixup*: Beyond Empirical Risk Minimization. In ICLR, 2018.

Calibration: Data-centric Solutions

Conventional calibration strategies have been empirically shown to enhance OOD detection.

❖ Change Labels

label smoothing [a]



interpolation

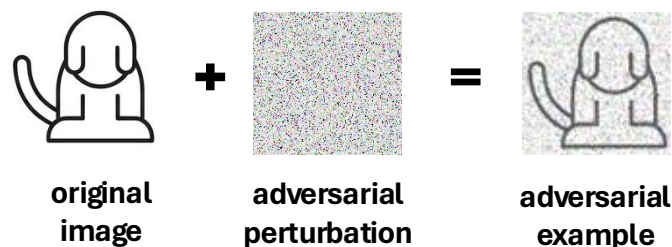
$$\tilde{y} = (1 - \alpha)y + \frac{\alpha}{c}$$

original label uniform

replaces one-hot labels with **slightly softened targets**.

❖ Change Inputs

adversarial examples [b]



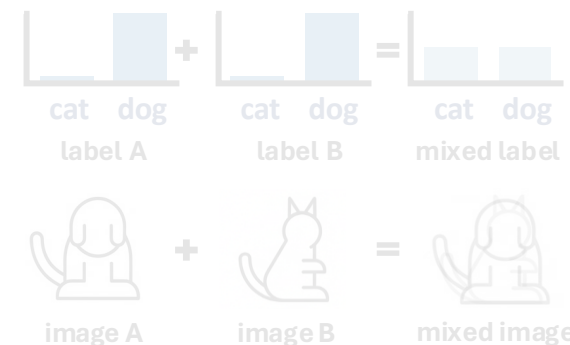
$$\max_{\delta \in \Delta} \mathcal{L}(x + \delta, y, \theta)$$

fail model perturb input

perturbs inputs **imperceptibly** to fool the models and then train the model.

❖ Change Labels and Inputs

mixup [c]



encourages the model to **spread probability mass more cautiously** around decision boundaries.

creates virtual training examples that **interpolates** between data pairs.

[a] Li et al. Rethinking Out-of-distribution (OOD) Detection: Masked Image Modeling is All You Need. In CVPR, 2023.

[b] Botschen et al. Out-of-Distribution Detection with Adversarial Outlier Exposure. In CVPR, 2025.

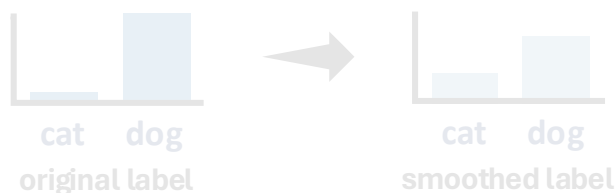
[c] Zhang et al. *mixup*: Beyond Empirical Risk Minimization. In ICLR, 2018.

Calibration: Data-centric Solutions

Conventional calibration strategies have been empirically shown to enhance OOD detection.

❖ Change Labels

label smoothing [a]



$$\tilde{y} = (1 - \alpha)y + \frac{\alpha}{c}$$

interpolation

original label uniform

replaces one-hot labels with **slightly softened targets**.

❖ Change Inputs

adversarial examples [b]

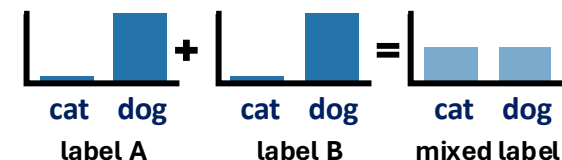


interpolates empirical distribution into a continuous one via a convex combinations of data and labels.

perturbs inputs **imperceptibly** to fool the models and then train the model.

❖ Change Labels and Inputs

mixup [c]



$$\tilde{x} = \lambda x_A + (1 - \lambda)x_B$$

$$\tilde{y} = \lambda y_A + (1 - \lambda)y_B$$

creates virtual training examples that **interpolates** between data pairs.

[a] Li et al. Rethinking Out-of-distribution (OOD) Detection: Masked Image Modeling is All You Need. In CVPR, 2023.

[b] Botschen et al. Out-of-Distribution Detection with Adversarial Outlier Exposure. In CVPR, 2025.

[c] Zhang et al. *mixup*: Beyond Empirical Risk Minimization. In ICLR, 2018.

Calibration: Model-centric Solutions

Relationship between MSP and Free Energy

softmax prediction $\rightarrow P(Y = y|X) = \frac{\exp\{h_y(X)\}}{\sum_{y'} \exp\{h_{y'}(X)\}}$


$\propto P(Y, X)$ (blue arrow pointing to $\exp\{h_y(X)\}$)

$\propto P(X)$, free energy (red arrow pointing to $\sum_{y'} \exp\{h_{y'}(X)\}$)

Free energy, which models $P(X)$, is more reliable for OOD detection than maximum softmax prediction. **Why?**

A Bayesian View [a].

Considering the following two learning goals, which one is more suitable for OOD detection?


$$\mathbb{P}(y \neq \boxed{f(x)} | r(x) = 0) + c_0 \mathbb{P}(\boxed{r(x)} = 1)$$

classifier **rejector**

Risk of accepting data that are misclassified Risk of rejecting data with cost c_0


$$\alpha \mathbb{P}_{\text{in}}(\boxed{r(x)} = 1) + \beta \mathbb{P}_{\text{ood}}(\boxed{r(x)} = 0)$$

ID distribution **OOD distribution**

Risk of rejecting ID Risk of accepting OOD

Calibration: Model-centric Solutions

Relationship between MSP and Free Energy

softmax prediction $\rightarrow P(Y = y|X) = \frac{\exp\{h_y(X)\}}{\sum_{y'} \exp\{h_{y'}(X)\}}$


$\propto P(Y, X)$ (blue arrow pointing to $\exp\{h_y(X)\}$)

$\propto P(X)$, free energy (red arrow pointing to $\sum_{y'} \exp\{h_{y'}(X)\}$)

Free energy, which models $P(X)$, is more reliable for OOD detection than maximum softmax prediction. **Why?**

A Bayesian View [a].

Considering the following two learning goals, **which one is more suitable for OOD detection?**



$$\mathbb{P}(y \neq \boxed{f(x)} \mid r(x) = 0) + c_0 \mathbb{P}(\boxed{r(x)} = 1)$$

classifier

rejector

Risk of accepting data that are misclassified

Risk of rejecting data with cost c_0


$$\alpha \mathbb{P}_{\text{in}}(\boxed{r(x)} = 1) + \beta \mathbb{P}_{\text{ood}}(\boxed{r(x)} = 0)$$

ID distribution

OOD distribution

Risk of rejecting ID

Risk of accepting OOD

Calibration: Model-centric Solutions

Relationship between MSP and Free Energy

softmax prediction $\rightarrow P(Y = y|X) = \frac{\exp\{h_y(X)\}}{\sum_{y'} \exp\{h_{y'}(X)\}}$

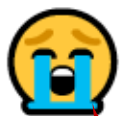
$\exp\{h_y(X)\} \propto P(Y, X)$

$\sum_{y'} \exp\{h_{y'}(X)\} \propto P(X), \text{ free energy}$

Free energy, which models $P(X)$, is more reliable for OOD detection than maximum softmax prediction. **Why?**

A Bayesian View [a].

Considering the following two learning goals, **which one is more suitable for OOD detection?**



$$\mathbb{P}(y \neq f(x), r(x) = 0) + c_0 \mathbb{P}(r(x) = 1)$$

used for **abstention-aware classification**.

$$r^*(x) = \llbracket \max_y \mathbb{P}(y|x) < 1 - c_0 \rrbracket$$

free energy



$$\alpha \mathbb{P}_{\text{in}}(r(x) = 0) + \beta \mathbb{P}_{\text{ood}}(r(x) = 1)$$

used for **OOD detection**.

$$r^*(x) = \llbracket \frac{\mathbb{P}_{\text{in}}(x)}{\mathbb{P}_{\text{ood}}(x)} < \frac{\beta}{\alpha} \rrbracket \approx \llbracket \frac{\mathbb{P}_{\text{in}}(x)}{\tau} < \frac{\beta}{\alpha} \rrbracket$$

MSP

Calibration: Model-centric Solutions

Relationship between MSP and Free Energy

softmax prediction $\rightarrow P(Y = y|X) = \frac{\exp\{h_y(X)\}}{\sum_{y'} \exp\{h_{y'}(X)\}}$

$\exp\{h_y(X)\} \propto P(Y, X)$

$\sum_{y'} \exp\{h_{y'}(X)\} \propto P(X), \text{ free energy}$

Free energy, which models $P(X)$, is more reliable for OOD detection than maximum softmax prediction. However, $P(X)$ is also not well-calibrated. So, **how to calibrate** $P(X)$?

Calibrating $P(Y)$ (low dimension) offers a sufficient condition of calibrating $P(X)$ (high dimension).

Mathematical

$$P(Y = y) = \frac{1}{Z} \int \exp\{h_y(x)\} dvx,$$

where $Z = \sum_y \int \exp\{h_y(x)\} dvx$

Monte Carlo



Hypothesis Test

$$H_0: P(Y = y) = \hat{P}(Y = y)$$

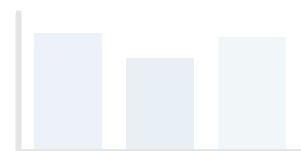
versus

$$H_1: P(Y = y) \neq \hat{P}(Y = y)$$

Empirical

$$\hat{P}(Y = y) = n_y / N$$

Count



Calibration: Model-centric Solutions—DCR [a]

Relationship between MSP and Free Energy

softmax prediction \rightarrow $P(Y = y|X) = \frac{\exp\{h_y(X)\}}{\sum_{y'} \exp\{h_{y'}(X)\}}$

$\exp\{h_y(X)\} \propto P(Y, X)$

$\sum_{y'} \exp\{h_{y'}(X)\} \propto P(X)$, free energy

Free energy, which models $P(X)$, is more reliable for OOD detection than maximum softmax prediction. However, $P(X)$ is also not well-calibrated. So, **how to calibrate** $P(X)$?

Calibrating $P(Y)$ (low dimension) offers a sufficient condition of **calibrating $P(X)$ (high dimension)**.

Mathematical

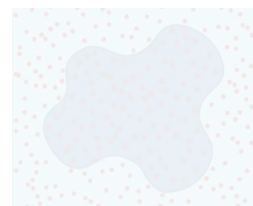
$$P(Y = y) = \frac{1}{Z} \int \exp\{h_y(x)\} dvx,$$

where $Z = \sum_y \int \exp\{h_y(x)\} dvx$

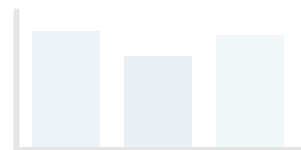
Empirical

$$\hat{P}(Y = y) = n_y / N$$

Monte Carlo



Count



Hypothesis Test

$$H_0: P(Y = y) = \hat{P}(Y = y)$$

versus

$$H_1: P(Y = y) \neq \hat{P}(Y = y)$$

Calibration: Model-centric Solutions—DCR [a]

Relationship between MSP and Free Energy

softmax prediction $\rightarrow P(Y = y|X) = \frac{\exp\{h_y(X)\}}{\sum_{y'} \exp\{h_{y'}(X)\}}$

$\exp\{h_y(X)\} \propto P(Y, X)$

$\sum_{y'} \exp\{h_{y'}(X)\} \propto P(X), \text{ free energy}$

Free energy, which models $P(X)$, is more reliable for OOD detection than maximum softmax prediction. However, $P(X)$ is also not well-calibrated. So, **how to calibrate** $P(X)$?

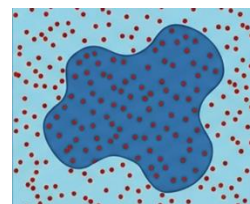
Calibrating $P(Y)$ (low dimension) offers a sufficient condition of **calibrating $P(X)$ (high dimension)**.

Mathematical

$$P(Y = y) = \frac{1}{Z} \int \exp\{h_y(x)\} dvx,$$

where $Z = \sum_y \int \exp\{h_y(x)\} dvx$

Monte Carlo



Hypothesis Test

$$H_0: P(Y = y) = \hat{P}(Y = y)$$

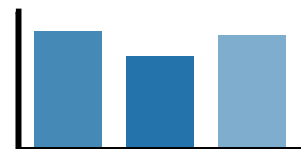
versus

$$H_1: P(Y = y) \neq \hat{P}(Y = y)$$

Empirical

$$\hat{P}(Y = y) = n_y / N$$

Count



Calibration: Model-centric Solutions—DCR [a]

Relationship between MSP and Free Energy

softmax prediction \rightarrow $P(Y = y|X) = \frac{\exp\{h_y(X)\}}{\sum_{y'} \exp\{h_{y'}(X)\}}$

$\exp\{h_y(X)\} \propto P(Y, X)$

$\sum_{y'} \exp\{h_{y'}(X)\} \propto P(X)$, free energy

Free energy, which models $P(X)$, is more reliable for OOD detection than maximum softmax prediction. However, $P(X)$ is also not well-calibrated. So, **how to calibrate** $P(X)$?

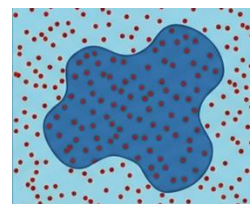
Calibrating $P(Y)$ (low dimension) offers a sufficient condition of **calibrating $P(X)$ (high dimension)**.

Mathematical

$$P(Y = y) = \frac{1}{Z} \int \exp\{h_y(x)\} dvx,$$

where $Z = \sum_y \int \exp\{h_y(x)\} dvx$

Monte Carlo



Hypothesis Test

$$H_0: P(Y = y) = \hat{P}(Y = y)$$

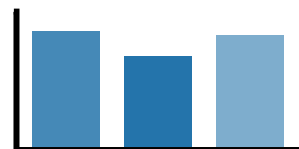
versus

$$H_1: P(Y = y) \neq \hat{P}(Y = y)$$

Empirical

$$\hat{P}(Y = y) = n_y / N$$

Count



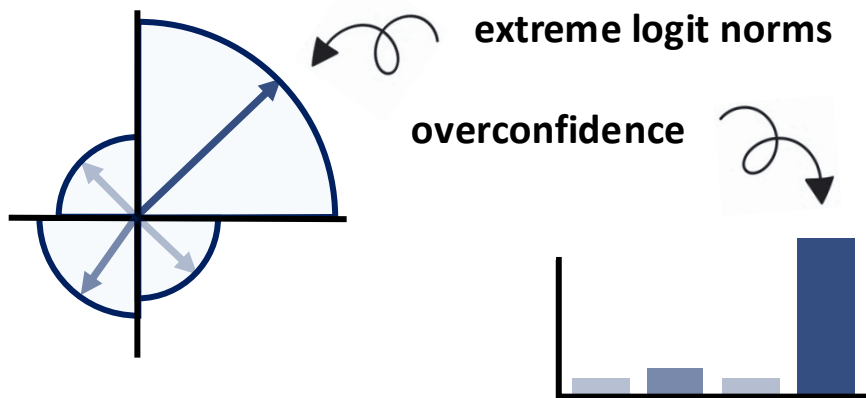
Calibration: Distribution-centric Solutions—LogitNorm

We seek **distribution modelling beyond Softmax** that are more proper for OOD detection.



Why is Softmax not sufficient?

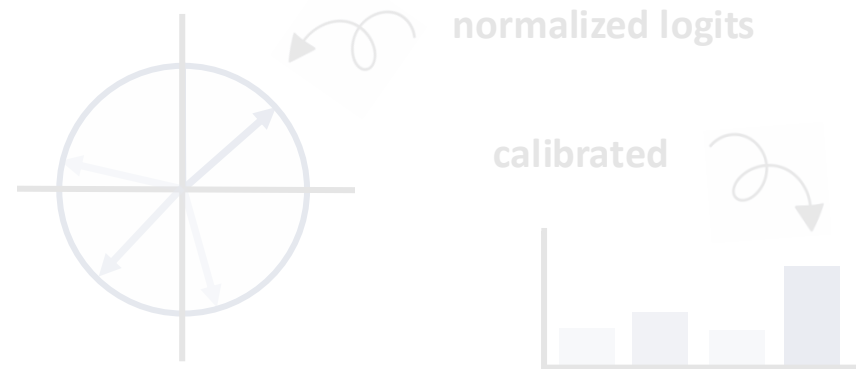
- ❖ **Scaling logits** increases confidence and decreases the risk.
- ❖ **Training pushes near-zero risk**, driving large confidence and inflating logits.



Logit Norm [a]

- ❖ **normalizes logit vector** to a constant norm during training, following

$$-\log \frac{\exp\{f_y/(\tau\|f_y\|)\}}{\sum_i \exp\{f_i/(\tau\|f_i\|)\}} \quad \text{normalization}$$



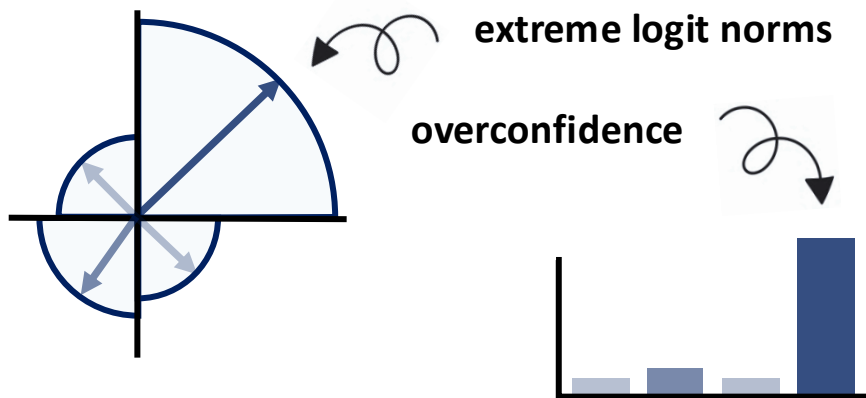
Calibration: Distribution-centric Solutions—LogitNorm

We seek **distribution modelling beyond Softmax** that are more proper for OOD detection.



Why is Softmax not sufficient?

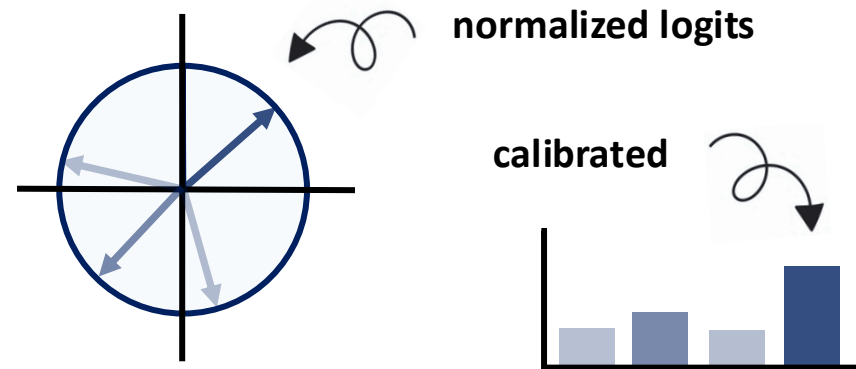
- ❖ **Scaling logits** increases confidence and decreases the risk.
- ❖ **Training pushes near-zero risk**, driving large confidence and inflating logits.



Logit Norm [a]

- ❖ **normalizes logit vector** to a constant norm during training, following

$$-\log \frac{\exp\{f_y/(\tau\|f_y\|)\}}{\sum_i \exp\{f_i/(\tau\|f_i\|)\}} \quad \text{normalization}$$



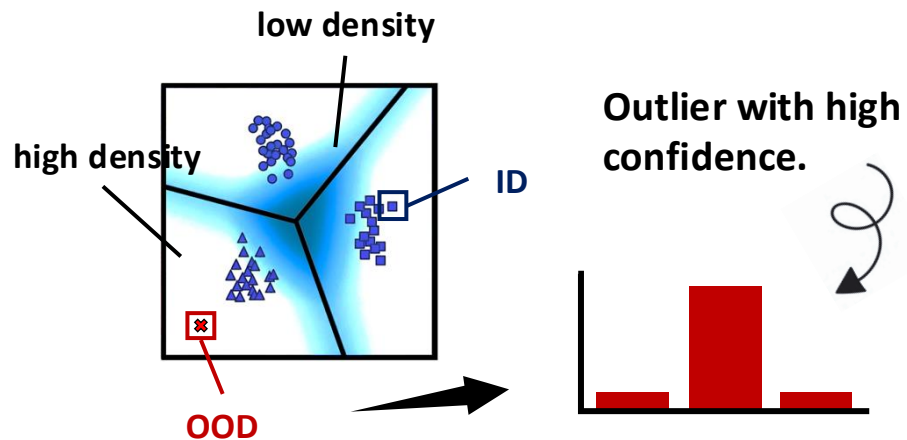
Calibration: Distribution-centric Solutions—SIREN

We seek **distribution modelling beyond Softmax** that are more proper for OOD detection.



Why is Softmax not sufficient?

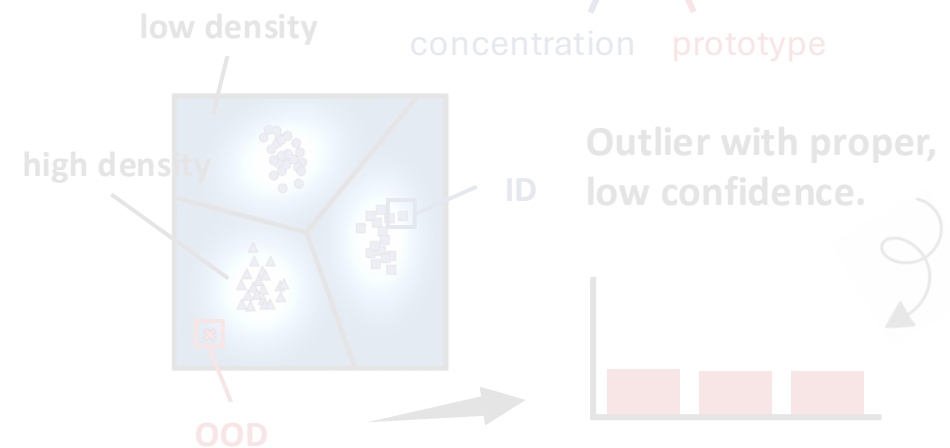
- ❖ The separation is **piecewise linear** in feature space.
- ❖ Regions that **contain no training points** can still be **assigned confidently**.



SIREN [a]

- ❖ The density depends on **similarity to prototypes** and vMF likelihood, creating a **curved, cluster-shaped** regions.

$$P(Y = c|r) \propto \exp\{\kappa_c \mu_c^T r\}$$



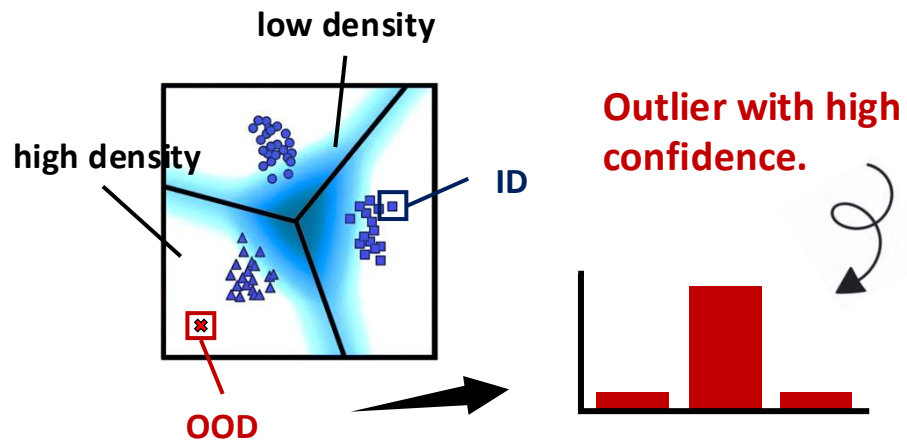
Calibration: Distribution-centric Solutions—SIREN

We seek **distribution modelling beyond Softmax** that are more proper for OOD detection.



Why is Softmax not sufficient?

- ❖ The separation is **piecewise linear** in feature space.
- ❖ Regions that **contain no training points** can still be assigned confidently.

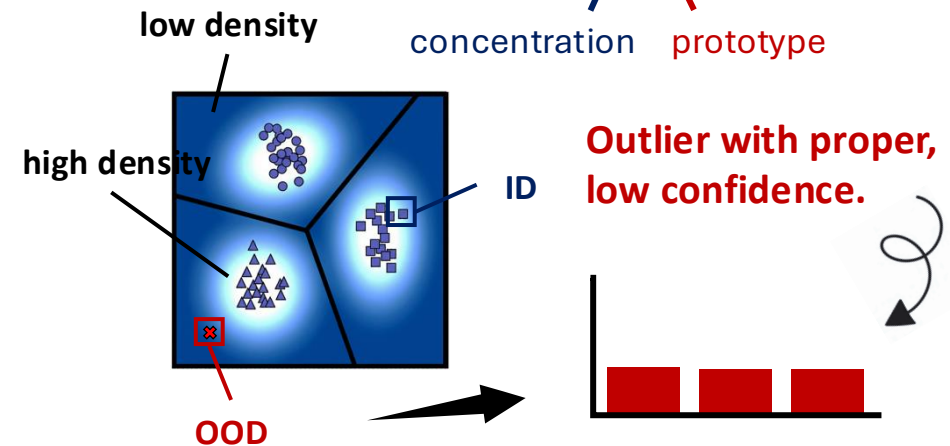


SIREN [a]

- ❖ The density depends on **similarity to prototypes** and vMF likelihood, creating a **curved, cluster-shaped** regions.

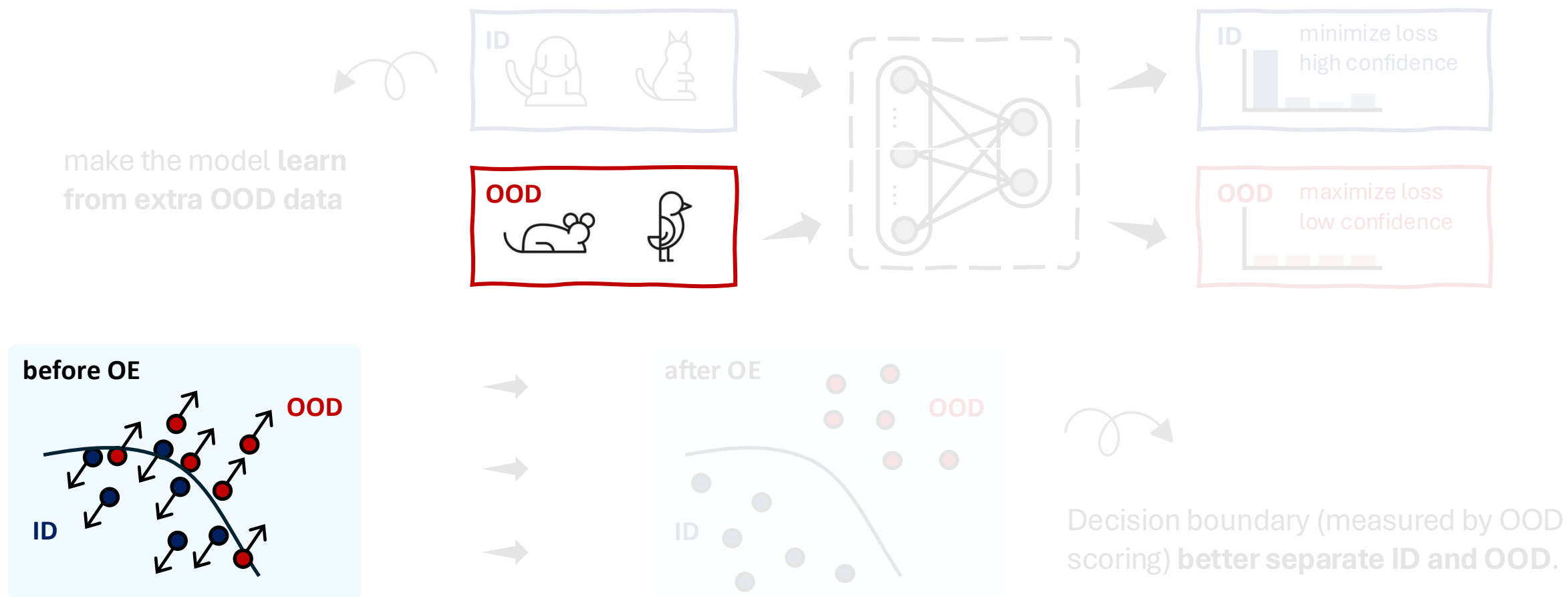
$$P(Y = c|r) \propto \exp\{\kappa_c \mu_c^T r\}$$

concentration prototype



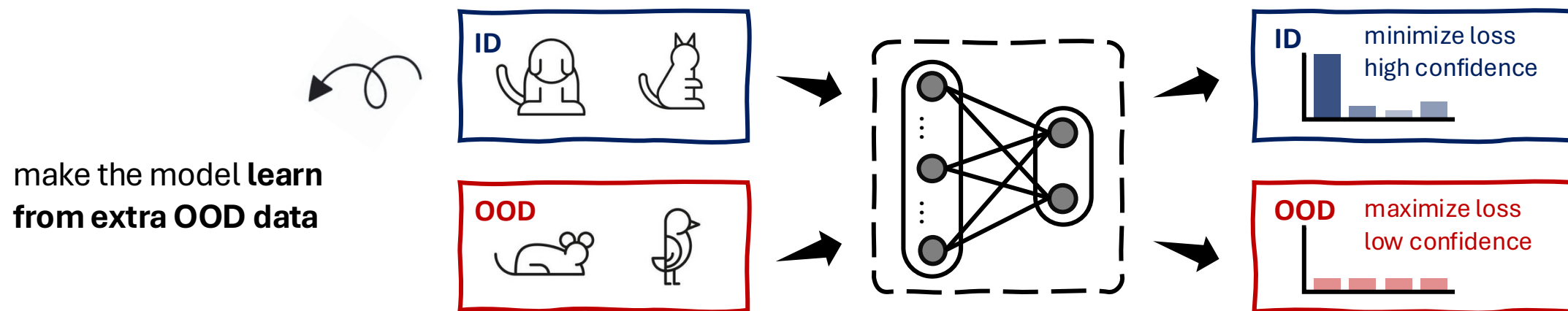
Outlier Exposure: Overview

Outlier exposure [a] takes OOD detection as **an additional binary classification task**, enabling models to directly learn to distinguish ID from OOD patterns.



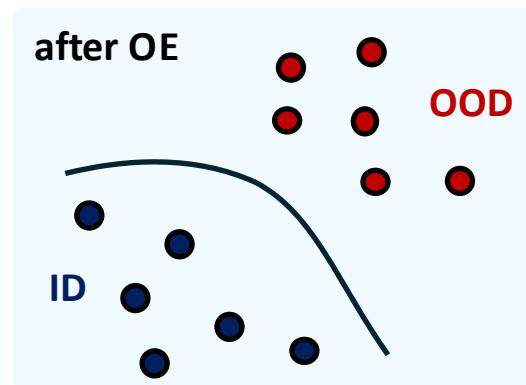
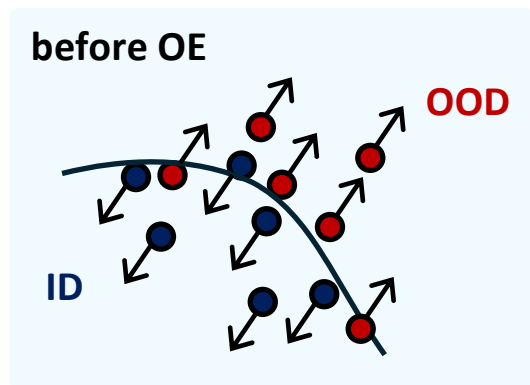
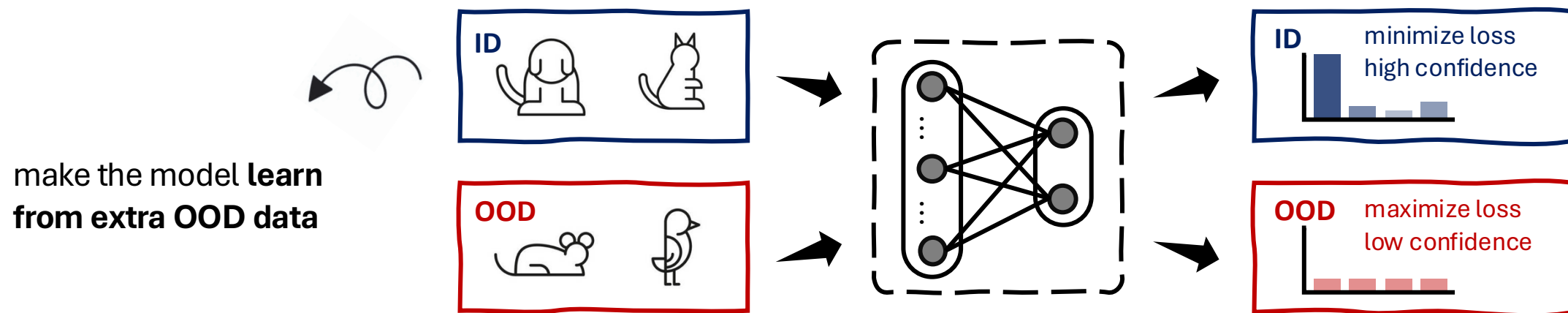
Outlier Exposure: Overview

Outlier exposure [a] takes OOD detection as **an additional binary classification task**, enabling models to directly learn to distinguish ID from OOD patterns.



Outlier Exposure: Overview

Outlier exposure [a] takes OOD detection as **an additional binary classification task**, enabling models to directly learn to distinguish ID from OOD patterns.



Decision boundary (measured by OOD scoring) **better separate ID and OOD.**

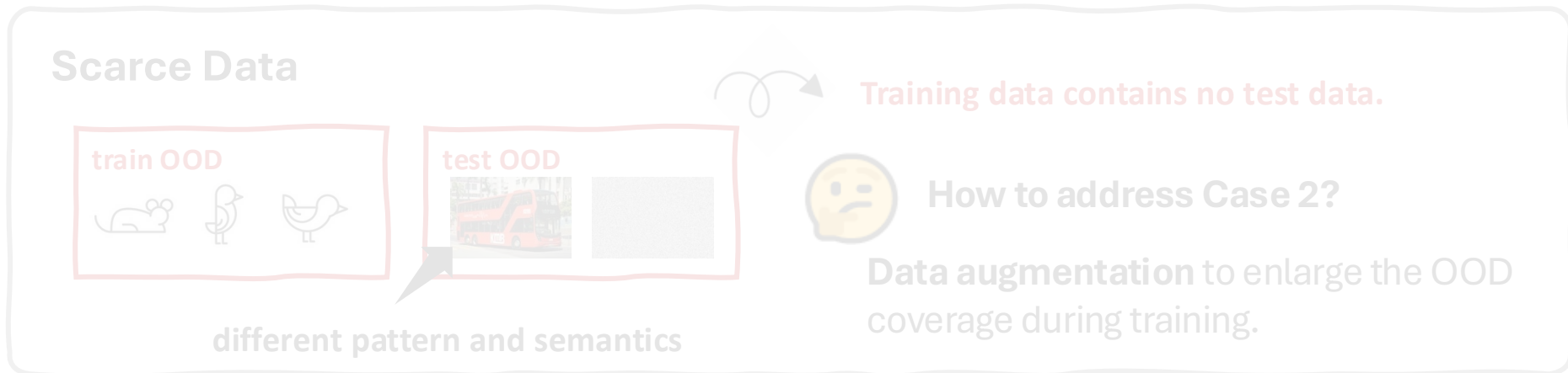
Outlier Exposure: Overview

Training and test OOD can differ, and **such distribution gap degrades OOD performance [a]**.

❖ Case 1.



❖ Case 2.

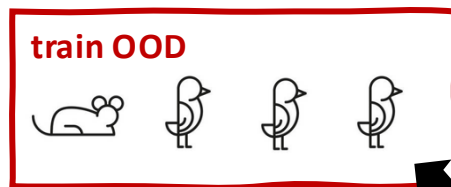


Outlier Exposure: Overview

Training and test OOD can differ, and **such distribution gap degrades OOD performance [a]**.

❖ Case 1.

Redundant Data



uninformative



Training data include test data,
but many samples are uninformative.



How to address Case 1?

Data sampling / reweighting to remove
those uninformative samples.

❖ Case 2.

Scarce Data



different pattern and semantics



Training data contains no test data.



How to address Case 2?

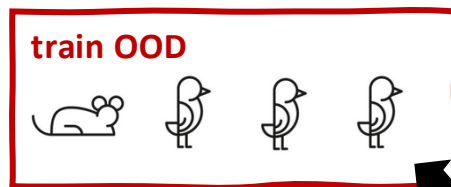
Data augmentation to enlarge the OOD
coverage during training.

Outlier Exposure: Overview

Training and test OOD can differ, and **such distribution gap degrades OOD performance** [a].

❖ Case 1.

Redundant Data



uninformative



Training data include test data,
but many samples are uninformative.



How to address Case 1?

Data sampling / reweighting to remove those uninformative samples.

❖ Case 2.

Scarce Data



different pattern and semantics



Training data contains no test data.



How to address Case 2?

Data augmentation to enlarge the OOD coverage during training.

Outlier Exposure: Overview

Data sampling learns from fewer examples and data augmentation from more.

❖ Case 1.

Redundant Data → Data Sampling



learn from less data



$$\mathbb{E}_{x \sim \mathcal{D}_{ood}} [w_x \ell_{ood}(x; \theta)]$$

reweighting

learn from resampled data

❖ Case 2.

Scarce Data → Data Augmentation



learn from extra data



$$\mathbb{E}_{x \sim \mathcal{D}_{ood}} [\ell_{ood}(x'; \theta)], x' = \Phi(x)$$

augmentation

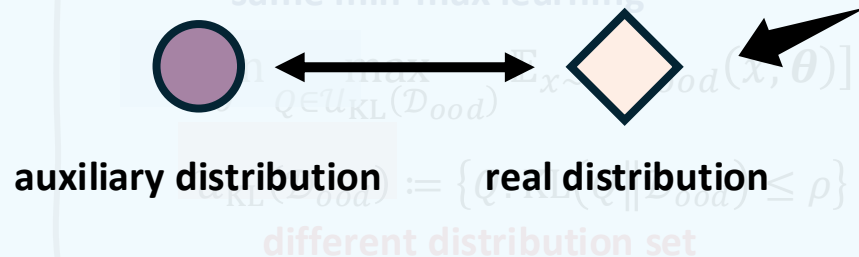
learn from extra data

Outlier Exposure: Overview

Data sampling and augmentation **both improve distribution robustness** [a].

❖ Case 1.

original distribution gap, $\min_{\theta} \mathbb{E}_{x \sim \mathcal{D}_{ood}} [\ell_{ood}(x; \theta)]$



The distribution gap degrade model performance.

$$w_x = \frac{\exp\{\ell_{ood}(x; \theta) / \eta\}}{\int \exp\{\ell_{ood}(x; \theta) / \eta\} dx}$$

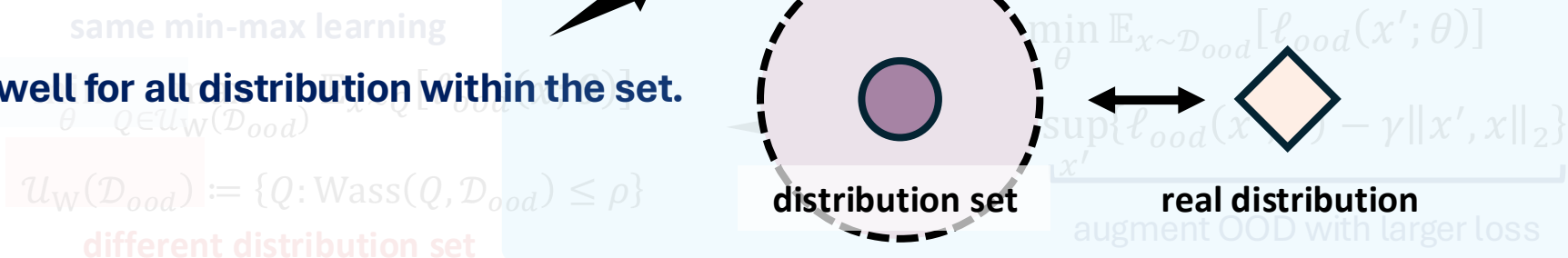
up-weight OOD with larger loss

❖ Case 2.

Scarce Data \rightarrow Data Augmentation

The model performs well for all distribution within the set.

augmented distribution gap, $\min_{\theta} \max_{Q \in \mathcal{U}(\mathcal{D}_{ood})} \mathbb{E}_{x \sim Q} [\ell_{ood}(x; \theta)]$



Outlier Exposure: Overview

Data sampling and augmentation **both improve distribution robustness** [a].

❖ Case 1.

Redundant Data → Data Sampling

same min-max learning

$$\min_{\theta} \max_{Q \in \mathcal{U}_{\text{KL}}(\mathcal{D}_{\text{ood}})} \mathbb{E}_{x \sim Q} [\ell_{\text{ood}}(x; \theta)]$$

$$\mathcal{U}_{\text{KL}}(\mathcal{D}_{\text{ood}}) := \{Q: \text{KL}(Q \parallel \mathcal{D}_{\text{ood}}) \leq \rho\}$$

different distribution set



$$\min_{\theta} \mathbb{E}_{x \sim \mathcal{D}_{\text{ood}}} [w_x \ell_{\text{ood}}(x; \theta)]$$

$$w_x = \frac{\exp\{\ell_{\text{ood}}(x; \theta) / \eta\}}{\int \exp\{\ell_{\text{ood}}(x; \theta) / \eta\} dx}$$

up-weight OOD with larger loss

❖ Case 2.

Scarce Data → Data Augmentation

same min-max learning

$$\min_{\theta} \max_{Q \in \mathcal{U}_W(\mathcal{D}_{\text{ood}})} \mathbb{E}_{x \sim Q} [\ell_{\text{ood}}(x; \theta)]$$

$$\mathcal{U}_W(\mathcal{D}_{\text{ood}}) := \{Q: \text{Wass}(Q, \mathcal{D}_{\text{ood}}) \leq \rho\}$$

different distribution set



$$\min_{\theta} \mathbb{E}_{x \sim \mathcal{D}_{\text{ood}}} [\ell_{\text{ood}}(x'; \theta)]$$

$$x' = \sup_{x'} \{\ell_{\text{ood}}(x'; \theta) - \gamma \|x', x\|_2\}$$

augment OOD with larger loss

Outlier Exposure: Overview

Data sampling and augmentation **both improve distribution robustness** [a].

❖ Case 1.

Redundant Data → Data Sampling

same min-max learning

$$\min_{\theta} \max_{Q \in \mathcal{U}_{\text{KL}}(\mathcal{D}_{\text{ood}})} \mathbb{E}_{x \sim Q} [\ell_{\text{ood}}(x; \theta)]$$

$$\mathcal{U}_{\text{KL}}(\mathcal{D}_{\text{ood}}) := \{Q: \text{KL}(Q \parallel \mathcal{D}_{\text{ood}}) \leq \rho\}$$

different distribution set



$$\min_{\theta} \mathbb{E}_{x \sim \mathcal{D}_{\text{ood}}} [w_x \ell_{\text{ood}}(x; \theta)]$$

$$w_x = \frac{\exp\{\ell_{\text{ood}}(x; \theta) / \eta\}}{\int \exp\{\ell_{\text{ood}}(x; \theta) / \eta\} dx}$$

up-weight OOD with larger loss

❖ Case 2.

Scarce Data → Data Augmentation

same min-max learning

$$\min_{\theta} \max_{Q \in \mathcal{U}_{\text{W}}(\mathcal{D}_{\text{ood}})} \mathbb{E}_{x \sim Q} [\ell_{\text{ood}}(x; \theta)]$$

$$\mathcal{U}_{\text{W}}(\mathcal{D}_{\text{ood}}) := \{Q: \text{Wass}(Q, \mathcal{D}_{\text{ood}}) \leq \rho\}$$

different distribution set



$$\min_{\theta} \mathbb{E}_{x \sim \mathcal{D}_{\text{ood}}} [\ell_{\text{ood}}(x'; \theta)]$$

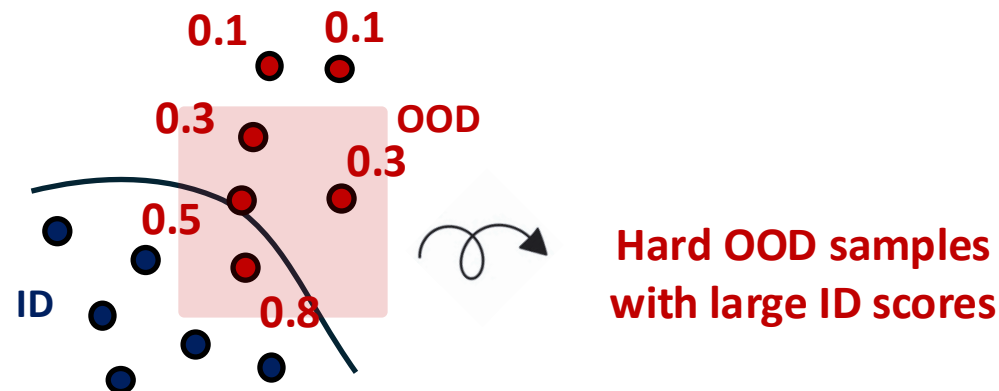
$$x' = \sup_{x'} \{\ell_{\text{ood}}(x'; \theta) - \gamma \|x', x\|_2\}$$

augment OOD with larger loss

Outlier Exposure: Sampling

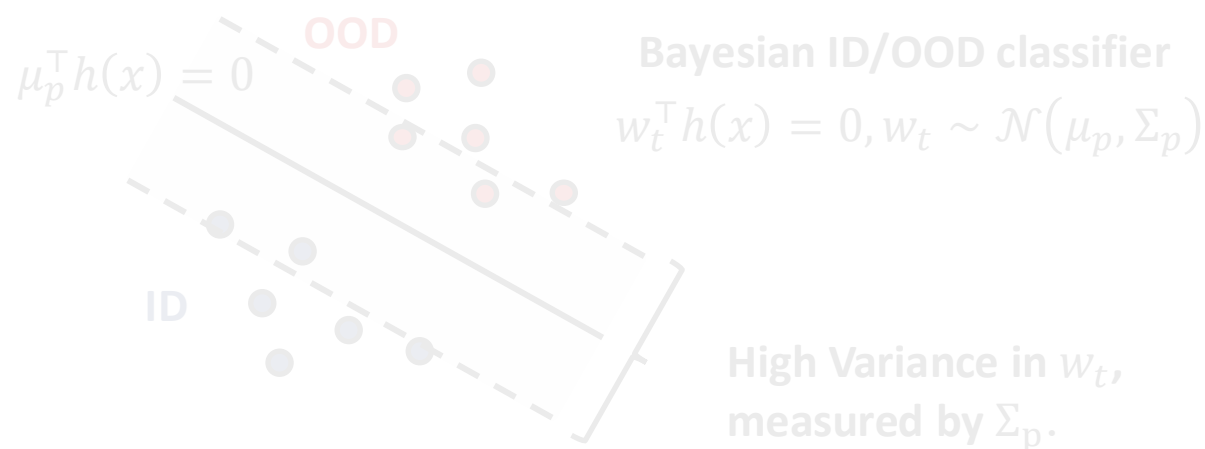
ATOM [a]

- ❖ **Hard OOD data** should be sampled more during training.
- ❖ **Large ID scores on OOD samples** indicates hard OOD data.



POEM [b]

- ❖ OOD data **near the ID/OOD boundary** should be sampled more.
- ❖ Bayesian classifiers capture parameter uncertainty, **improving exploration**.



[a] Chen et al. ATOM: Robustifying Out-of-distribution Detection Using Outlier Mining. In ECML PKDD, 2021.

[b] Ming et al. POEM: Out-of-distribution Detection with Posterior Sampling. In ICML, 2022.

Outlier Exposure: Sampling

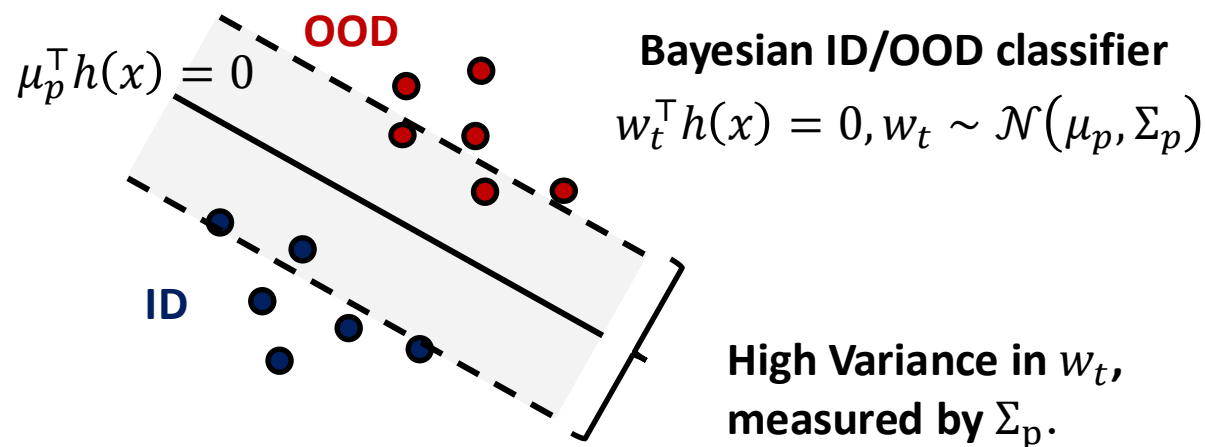
ATOM [a]

- ❖ **Hard OOD data** should be sampled more during training.
- ❖ **Large ID scores on OOD samples** indicates hard OOD data.



POEM [b]

- ❖ OOD data **near the ID/OOD boundary** should be sampled more.
- ❖ Bayesian classifiers capture parameter uncertainty, **improving exploration**.



[a] Chen et al. ATOM: Robustifying Out-of-distribution Detection Using Outlier Mining. In ECML PKDD, 2021.

[b] Ming et al. POEM: Out-of-distribution Detection with Posterior Sampling. In ICML, 2022.

Outlier Exposure: Sampling

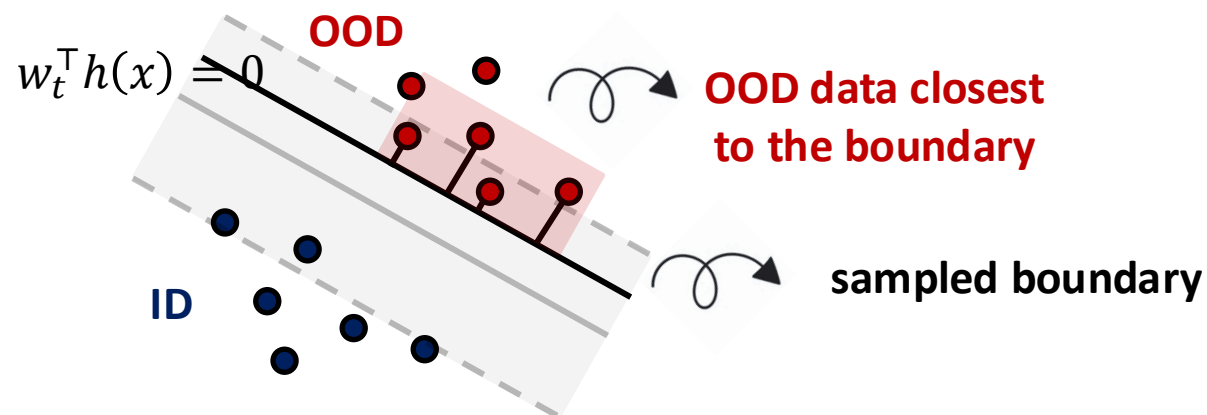
ATOM [a]

- ❖ **Hard OOD data** should be sampled more during training.
- ❖ **Large ID scores on OOD samples** indicates hard OOD data.



POEM [b]

- ❖ OOD data **near the ID/OOD boundary** should be sampled more.
- ❖ Bayesian classifiers capture parameter uncertainty, **improving exploration**.



[a] Chen et al. ATOM: Robustifying Out-of-distribution Detection Using Outlier Mining. In ECML PKDD, 2021.

[b] Ming et al. POEM: Out-of-distribution Detection with Posterior Sampling. In ICML, 2022.

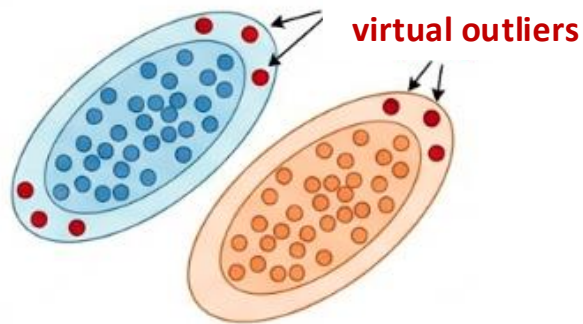
Outlier Exposure: Augmentation

Augmentation can be conducted in either **embedding space** or input space.

❖ Synthesis

VOS [a]

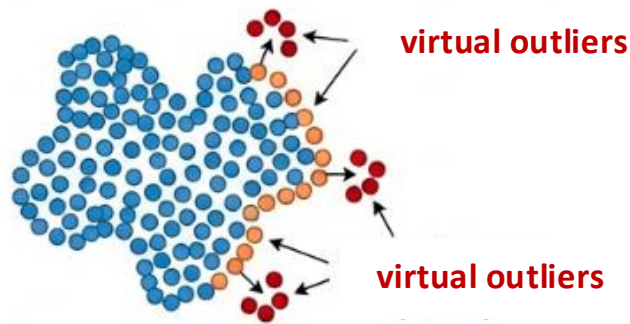
class-conditional **Gaussian modelling**



- ❖ estimate μ and Σ per class.
- ❖ sample $v \sim \mathcal{N}(\mu_c, \Sigma)$.
- ❖ keep only low likelihood samples with $p(v) < \epsilon$.

NPOS [b]

Non-parametric K-NN modelling

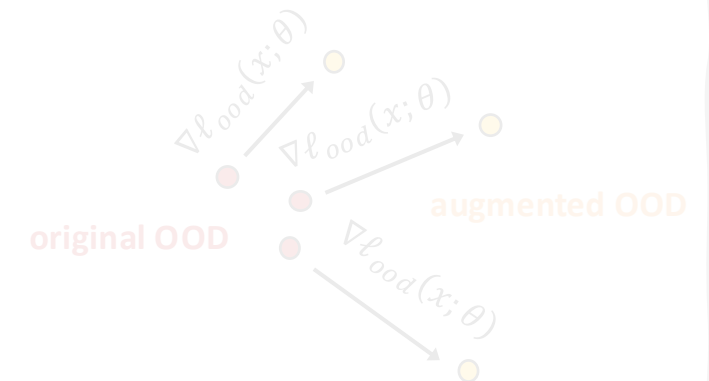


- ❖ Identify boundary ID with high k-NN.
- ❖ sample locally around boundary $v \sim \mathcal{N}(h_{bdy}, \sigma^2 I)$.
- ❖ keep candidates with large k-NN.

❖ Augmentation

DAL [c]

worst-case **OOD augmentation**



- ❖ sample original OOD data.
- ❖ gradient ascent to increase their respective loss values.
- ❖ augmented OOD data for training.

[a] Du et al. VOS: Learning What You Don't Know by Virtual Outlier Synthesis. In ICLR, 2022.

[b] Tao et al. Non-parametric Outlier Synthesis. In ICLR, 2023.

[c] Wang et al. Learning to Augment Distributions for Out-of-distribution Detection. In NeurIPS, 2023.

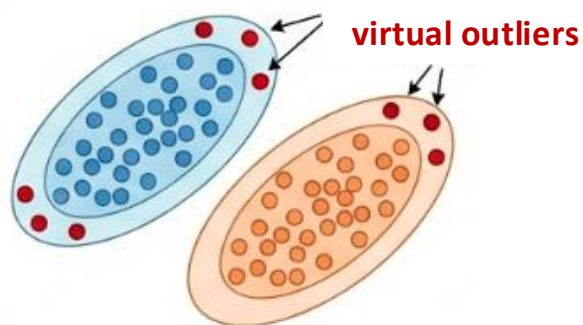
Outlier Exposure: Augmentation

Augmentation can be conducted in either **embedding space** or input space.

❖ Synthesis

VOS [a]

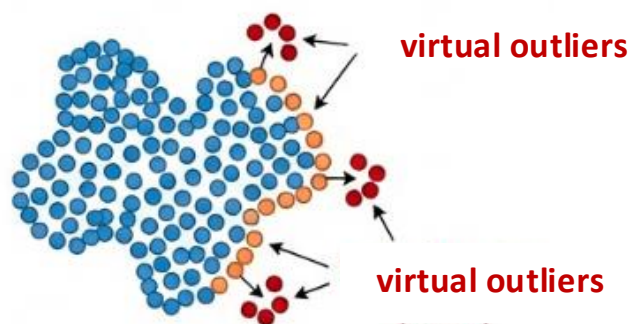
class-conditional **Gaussian modelling**



- ❖ estimate μ and Σ per class.
- ❖ sample $v \sim \mathcal{N}(\mu_c, \Sigma)$.
- ❖ keep only low likelihood samples with $p(v) < \epsilon$.

NPOS [b]

Non-parametric K-NN modelling

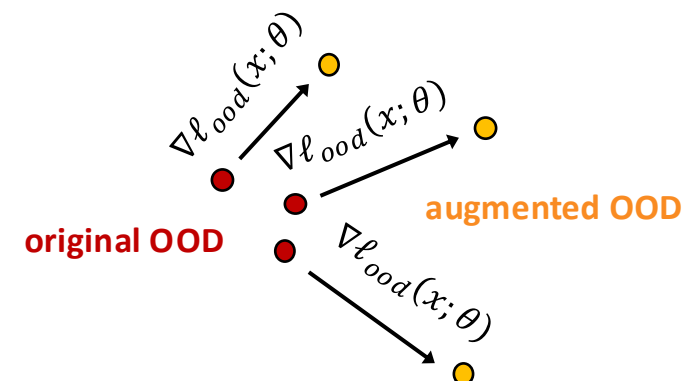


- ❖ Identify boundary ID with high k-NN.
- ❖ sample locally around boundary $v \sim \mathcal{N}(h_{bdy}, \sigma^2 I)$.
- ❖ keep candidates with large k-NN.

❖ Augmentation

DAL [c]

worst-case **OOD augmentation**



- ❖ sample original OOD data.
- ❖ gradient ascent to increase their respective loss values.
- ❖ augmented OOD data for training.

[a] Du et al. VOS: Learning What You Don't Know by Virtual Outlier Synthesis. In ICLR, 2022.

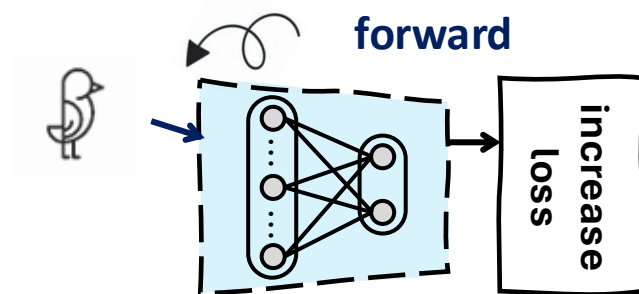
[b] Tao et al. Non-parametric Outlier Synthesis. In ICLR, 2023.

[c] Wang et al. Learning to Augment Distributions for Out-of-distribution Detection. In NeurIPS, 2023.

Outlier Exposure: Augmentation

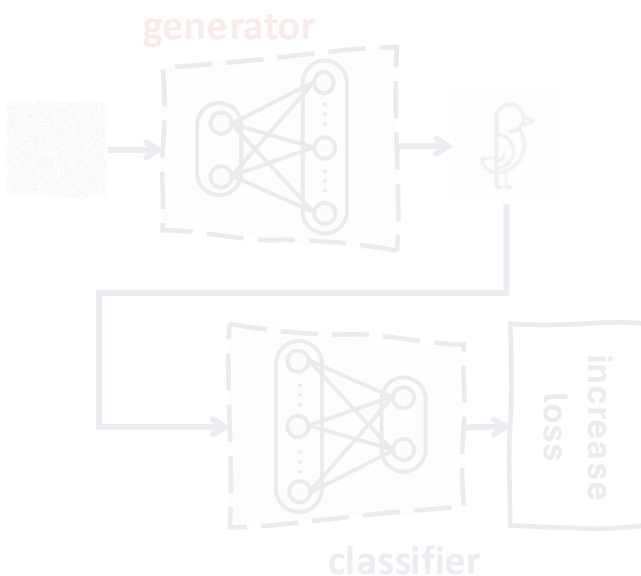
Augmentation can be conducted in either embedding space or **input space**.

❖ Adversarial Training [a]



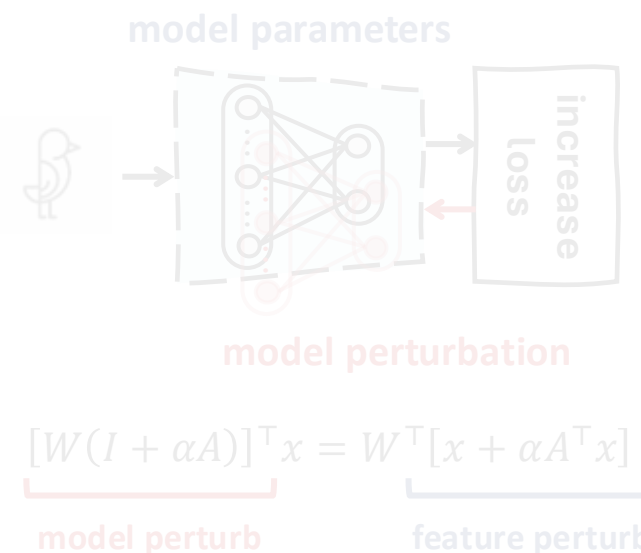
perturb inputs to maximize OOD loss values, thus creating worst-case examples.

❖ Generative Model [b]



perturb latent embedding to maximize OOD loss values, thus creating worst-case examples.

❖ Implicit Synthesis [c]



perturb parameters that maximize OOD loss can implicitly create worst-case examples.

[a] Zhu et al. Diversified Outlier Exposure for Out-of-distribution Detection via Informative Extrapolation. In NeurIPS, 2023.

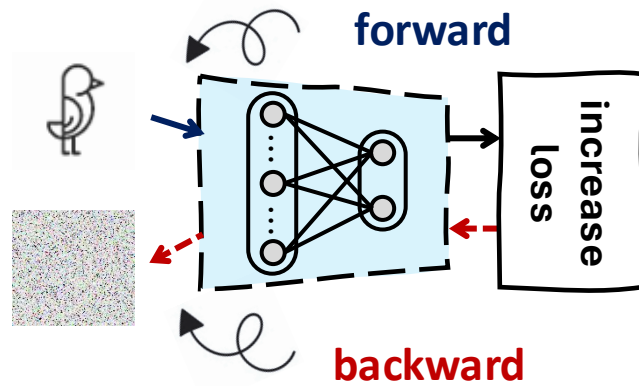
[b] Lee et al. Training Confidence-calibrated Classifiers for Detecting Out-of-distribution Samples. In ICLR, 2018.

[c] Wang et al. Out-of-distribution Detection with Implicit Outlier Transformation. In ICLR, 2023.

Outlier Exposure: Augmentation

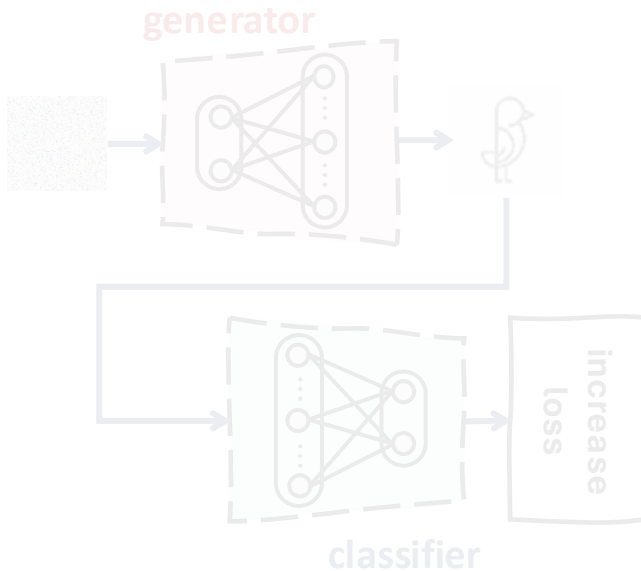
Augmentation can be conducted in either embedding space or **input space**.

❖ Adversarial Training [a]



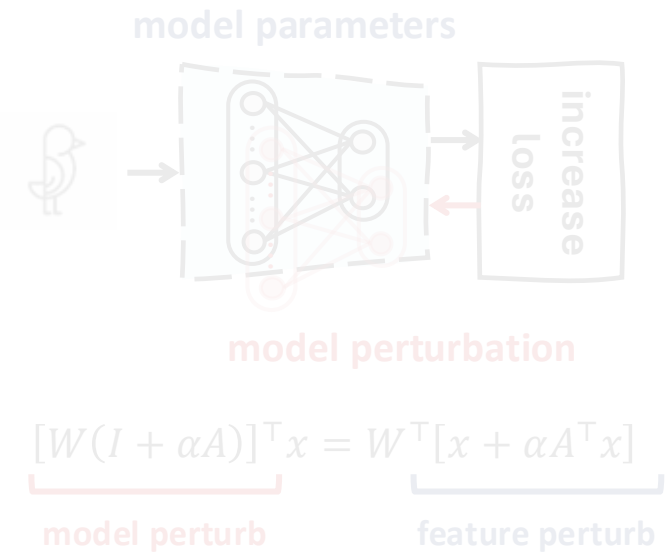
perturb inputs to maximize OOD loss values, thus creating worst-case examples.

❖ Generative Model [b]



perturb latent embedding to maximize OOD loss values, thus creating worst-case examples.

❖ Implicit Synthesis [c]



perturb parameters that maximize OOD loss can implicitly create worst-case examples.

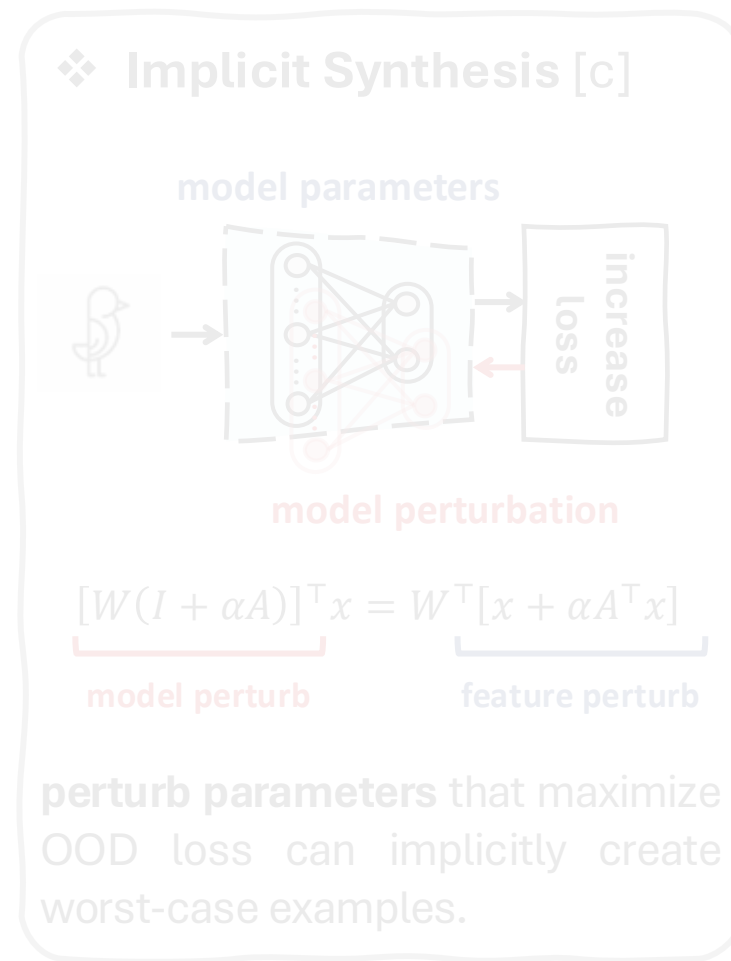
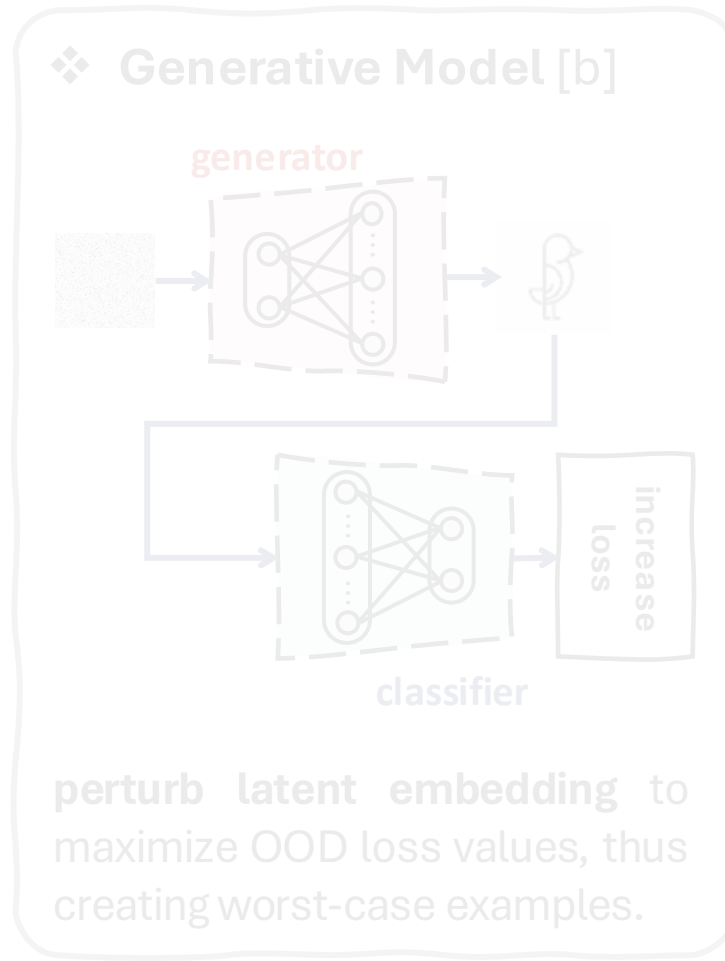
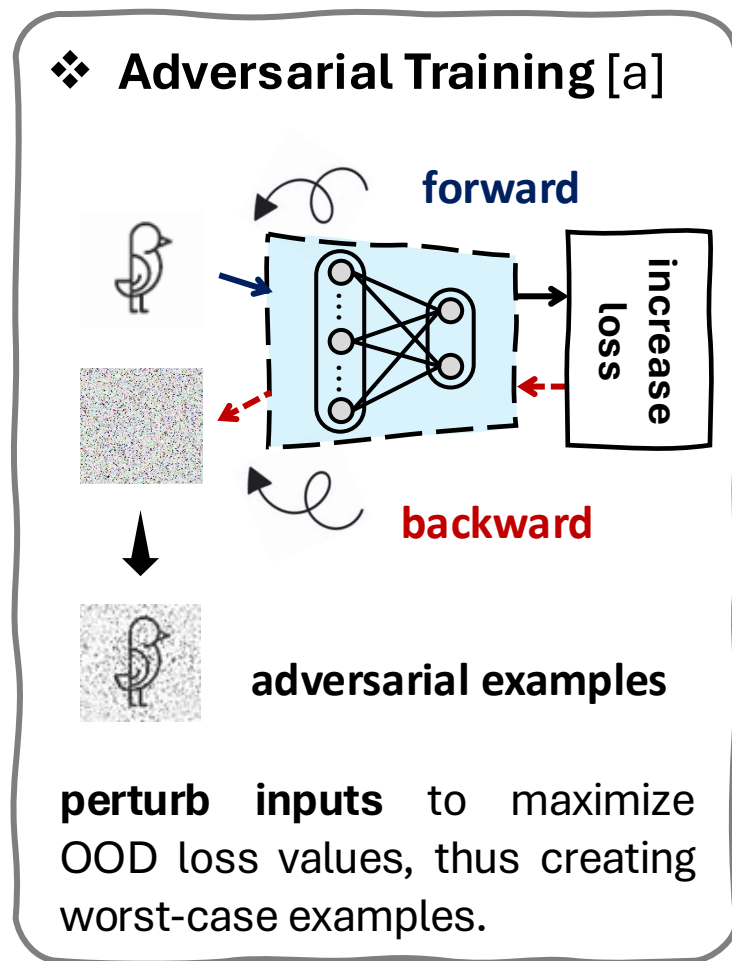
[a] Zhu et al. Diversified Outlier Exposure for Out-of-distribution Detection via Informative Extrapolation. In NeurIPS, 2023.

[b] Lee et al. Training Confidence-calibrated Classifiers for Detecting Out-of-distribution Samples. In ICLR, 2018.

[c] Wang et al. Out-of-distribution Detection with Implicit Outlier Transformation. In ICLR, 2023.

Outlier Exposure: Augmentation

Augmentation can be conducted in either embedding space or **input space**.



[a] Zhu et al. Diversified Outlier Exposure for Out-of-distribution Detection via Informative Extrapolation. In NeurIPS, 2023.

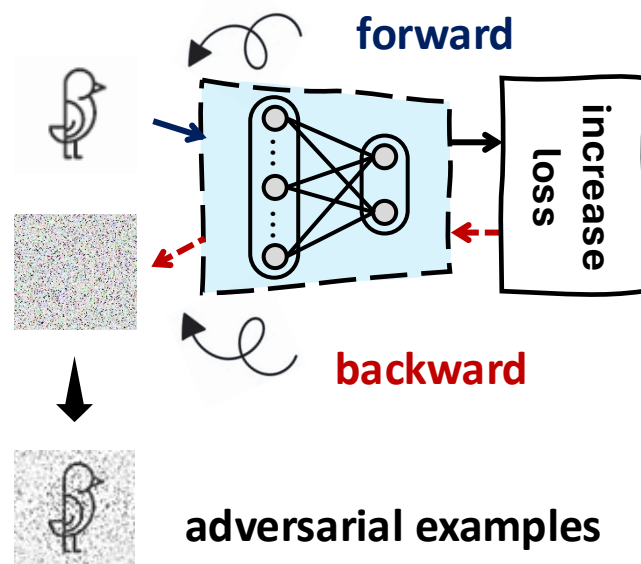
[b] Lee et al. Training Confidence-calibrated Classifiers for Detecting Out-of-distribution Samples. In ICLR, 2018.

[c] Wang et al. Out-of-distribution Detection with Implicit Outlier Transformation. In ICLR, 2023.

Outlier Exposure: Augmentation

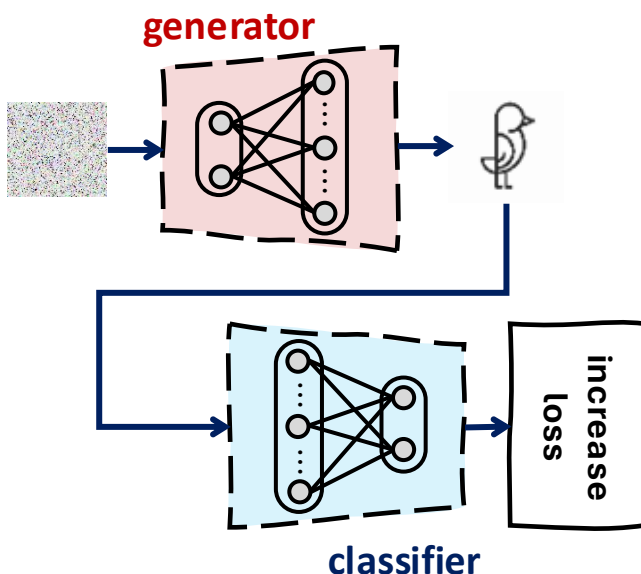
Augmentation can be conducted in either embedding space or **input space**.

❖ Adversarial Training [a]



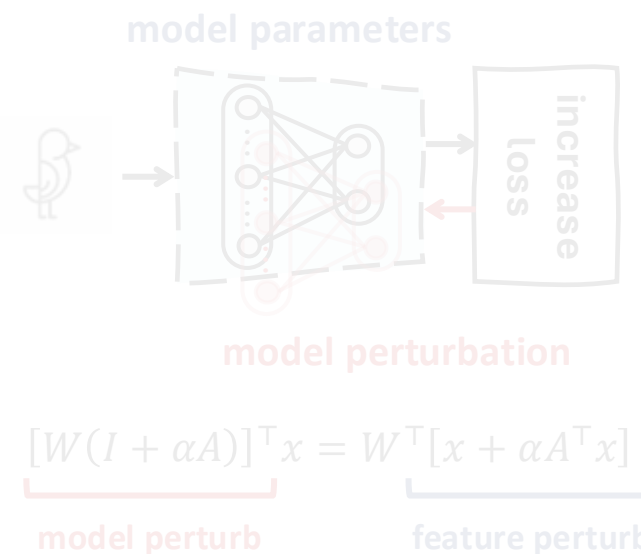
perturb inputs to maximize OOD loss values, thus creating worst-case examples.

❖ Generative Model [b]



perturb latent embedding to maximize OOD loss values, thus creating worst-case examples.

❖ Implicit Synthesis [c]



perturb parameters that maximize OOD loss can implicitly create worst-case examples.

[a] Zhu et al. Diversified Outlier Exposure for Out-of-distribution Detection via Informative Extrapolation. In NeurIPS, 2023.

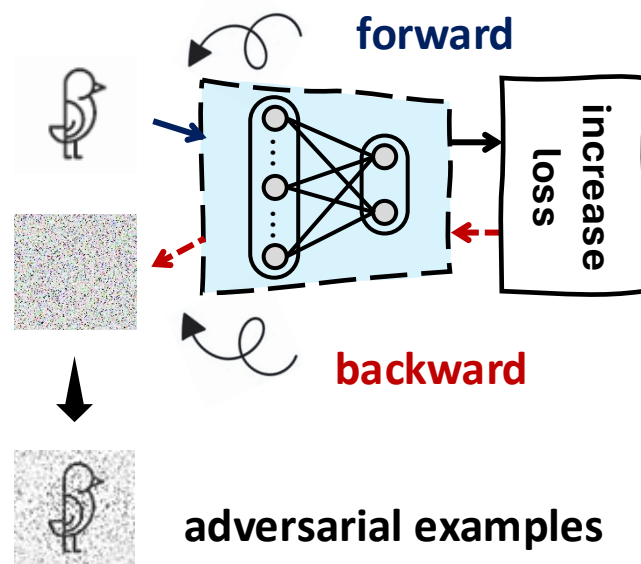
[b] Lee et al. Training Confidence-calibrated Classifiers for Detecting Out-of-distribution Samples. In ICLR, 2018.

[c] Wang et al. Out-of-distribution Detection with Implicit Outlier Transformation. In ICLR, 2023.

Outlier Exposure: Augmentation

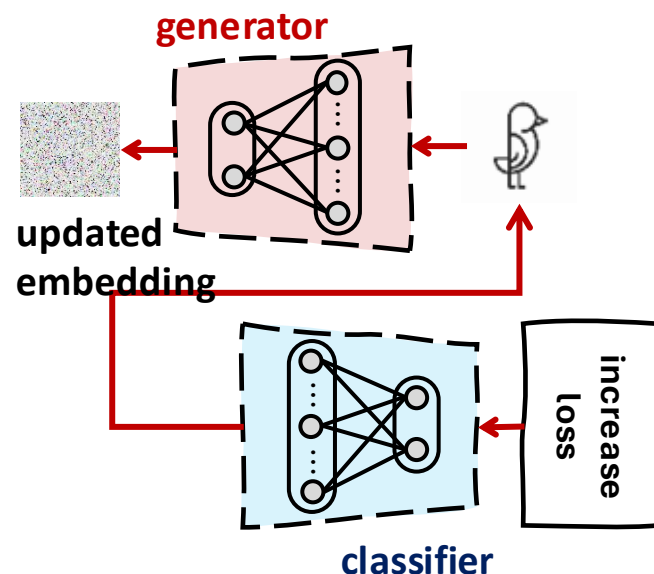
Augmentation can be conducted in either embedding space or **input space**.

❖ Adversarial Training [a]



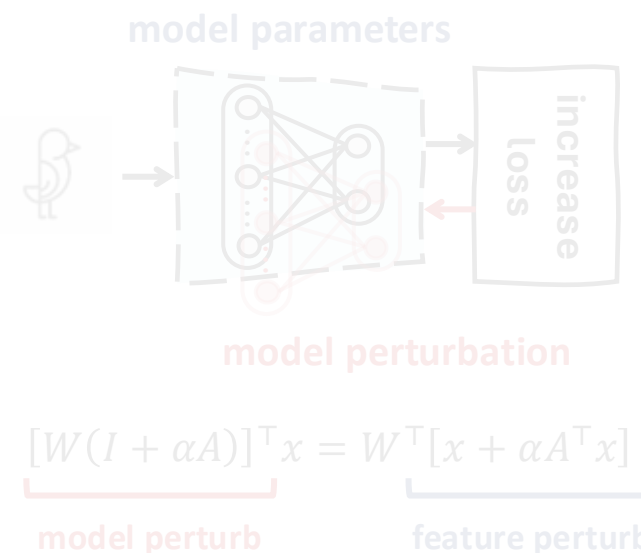
perturb inputs to maximize OOD loss values, thus creating worst-case examples.

❖ Generative Model [b]



perturb latent embedding to maximize OOD loss values, thus creating worst-case examples.

❖ Implicit Synthesis [c]



perturb parameters that maximize OOD loss can implicitly create worst-case examples.

[a] Zhu et al. Diversified Outlier Exposure for Out-of-distribution Detection via Informative Extrapolation. In NeurIPS, 2023.

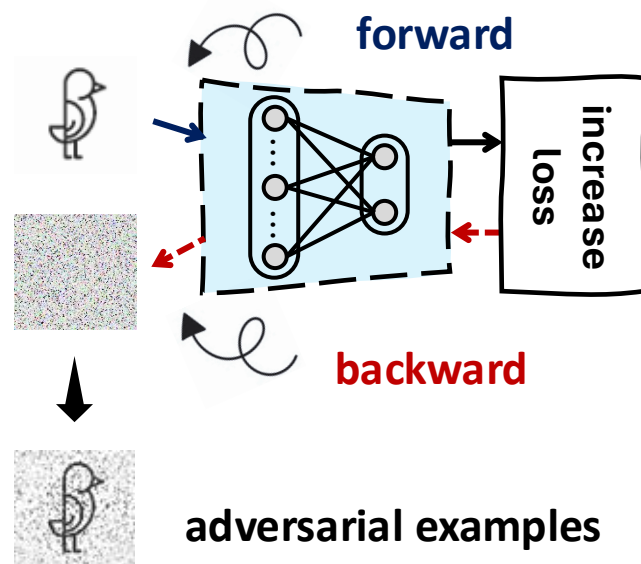
[b] Lee et al. Training Confidence-calibrated Classifiers for Detecting Out-of-distribution Samples. In ICLR, 2018.

[c] Wang et al. Out-of-distribution Detection with Implicit Outlier Transformation. In ICLR, 2023.

Outlier Exposure: Augmentation

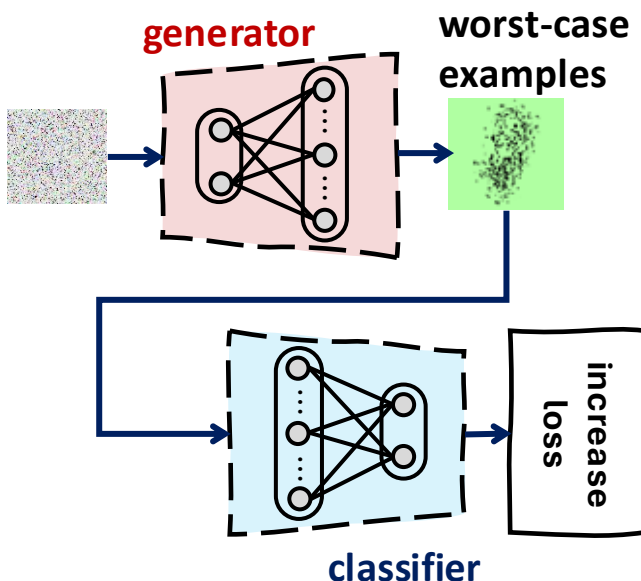
Augmentation can be conducted in either embedding space or **input space**.

❖ Adversarial Training [a]



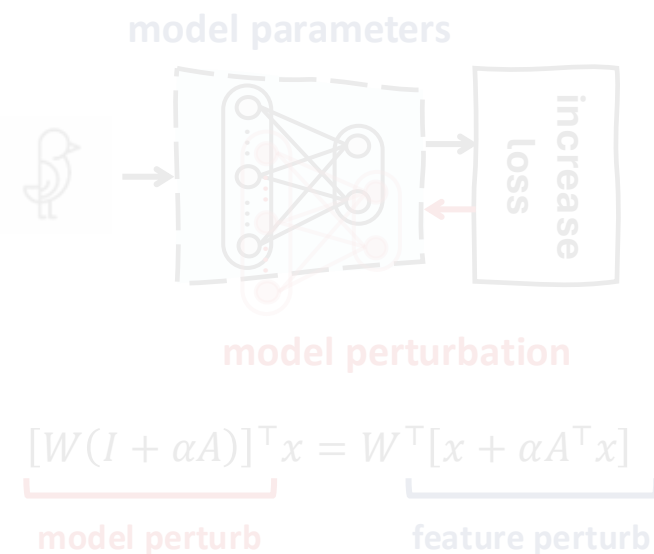
perturb inputs to maximize OOD loss values, thus creating worst-case examples.

❖ Generative Model [b]



perturb latent embedding to maximize OOD loss values, thus creating worst-case examples.

❖ Implicit Synthesis [c]



perturb parameters that maximize OOD loss can implicitly create worst-case examples.

[a] Zhu et al. Diversified Outlier Exposure for Out-of-distribution Detection via Informative Extrapolation. In NeurIPS, 2023.

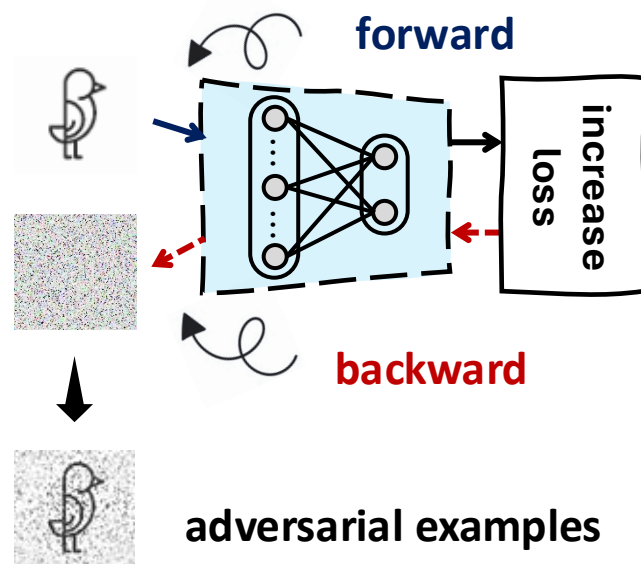
[b] Lee et al. Training Confidence-calibrated Classifiers for Detecting Out-of-distribution Samples. In ICLR, 2018.

[c] Wang et al. Out-of-distribution Detection with Implicit Outlier Transformation. In ICLR, 2023.

Outlier Exposure: Augmentation

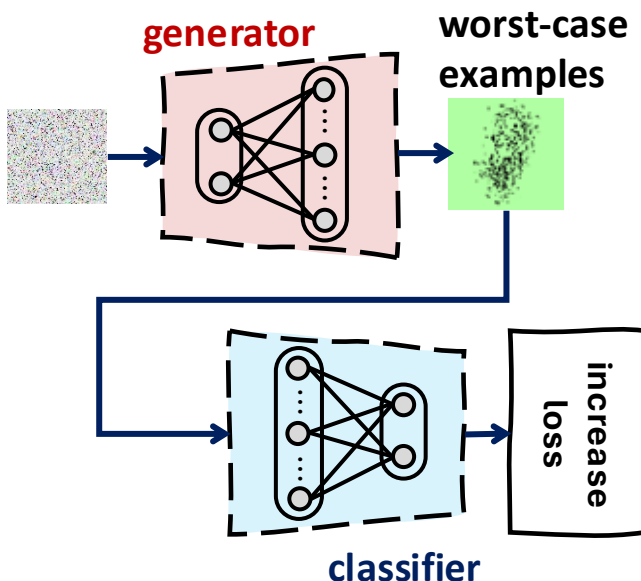
Augmentation can be conducted in either embedding space or **input space**.

❖ Adversarial Training [a]



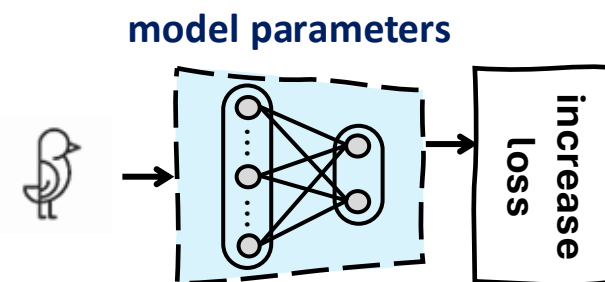
perturb inputs to maximize OOD loss values, thus creating worst-case examples.

❖ Generative Model [b]



perturb latent embedding to maximize OOD loss values, thus creating worst-case examples.

❖ Implicit Synthesis [c]



$$\underbrace{[W(I + \alpha A)]^T}_{\text{model perturb}} x = W^T \underbrace{[x + \alpha A^T x]}_{\text{feature perturb}}$$

perturb parameters that maximize OOD loss can implicitly create worst-case examples.

[a] Zhu et al. Diversified Outlier Exposure for Out-of-distribution Detection via Informative Extrapolation. In NeurIPS, 2023.

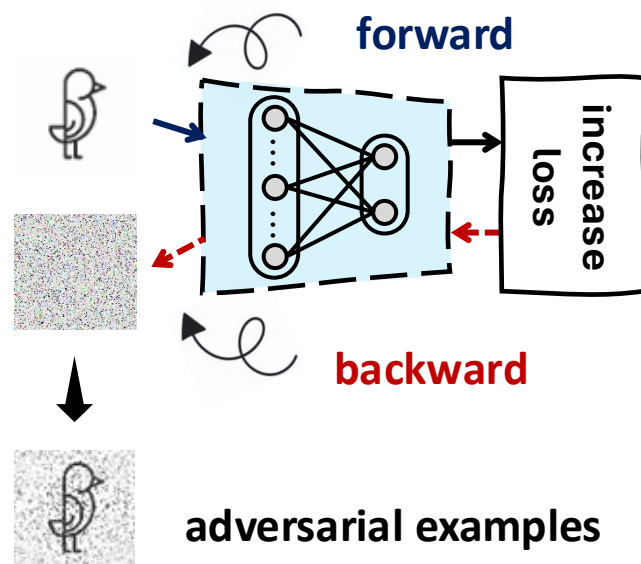
[b] Lee et al. Training Confidence-calibrated Classifiers for Detecting Out-of-distribution Samples. In ICLR, 2018.

[c] Wang et al. Out-of-distribution Detection with Implicit Outlier Transformation. In ICLR, 2023.

Outlier Exposure: Augmentation

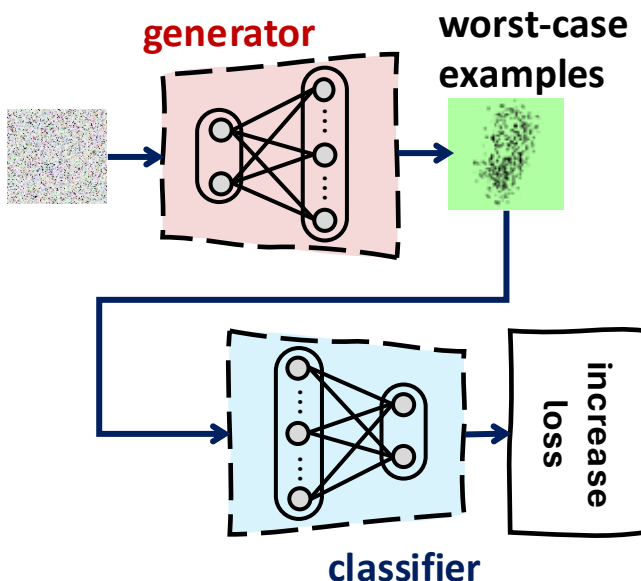
Augmentation can be conducted in either embedding space or **input space**.

❖ Adversarial Training [a]



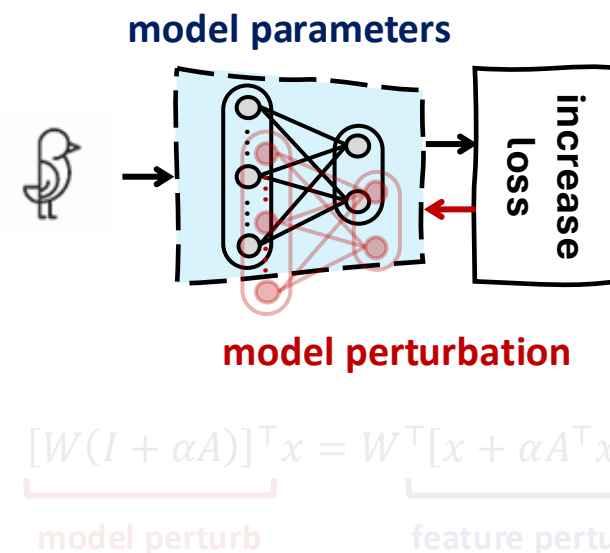
perturb inputs to maximize OOD loss values, thus creating worst-case examples.

❖ Generative Model [b]



perturb latent embedding to maximize OOD loss values, thus creating worst-case examples.

❖ Implicit Synthesis [c]



perturb parameters that maximize OOD loss can implicitly create worst-case examples.

[a] Zhu et al. Diversified Outlier Exposure for Out-of-distribution Detection via Informative Extrapolation. In NeurIPS, 2023.

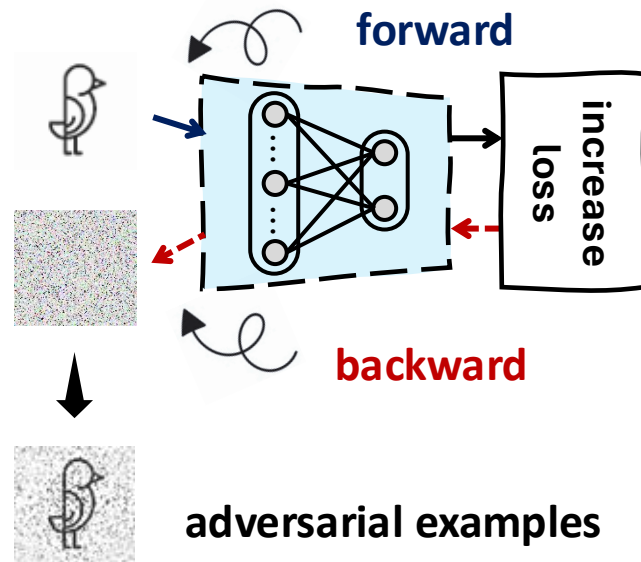
[b] Lee et al. Training Confidence-calibrated Classifiers for Detecting Out-of-distribution Samples. In ICLR, 2018.

[c] Wang et al. Out-of-distribution Detection with Implicit Outlier Transformation. In ICLR, 2023.

Outlier Exposure: Augmentation

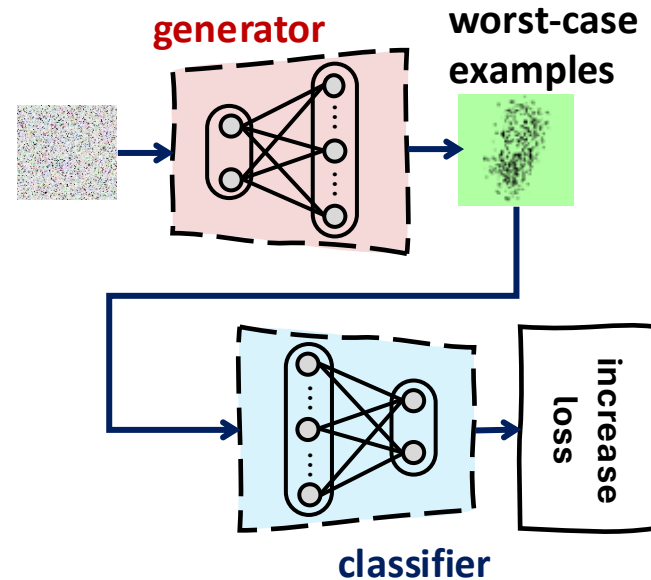
Augmentation can be conducted in either embedding space or **input space**.

❖ Adversarial Training [a]



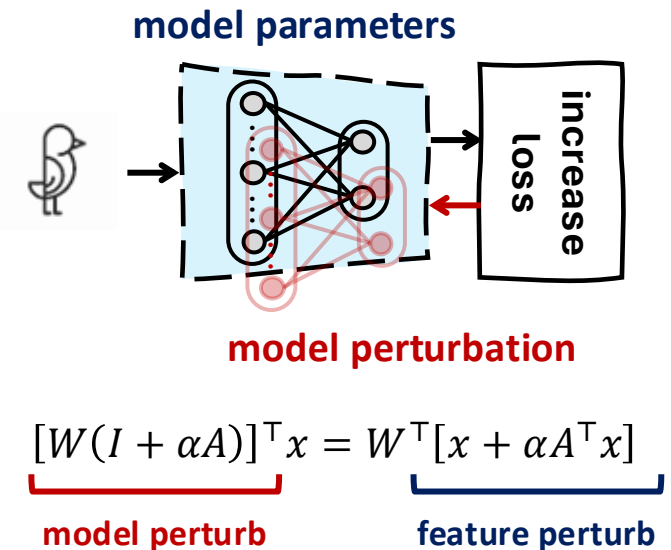
perturb inputs to maximize OOD loss values, thus creating worst-case examples.

❖ Generative Model [b]



perturb latent embedding to maximize OOD loss values, thus creating worst-case examples.

❖ Implicit Synthesis [c]



perturb parameters that maximize OOD loss can implicitly create worst-case examples.

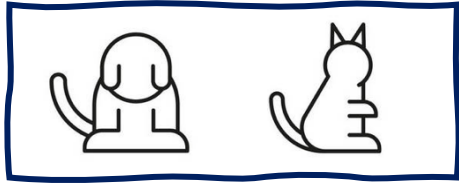
[a] Zhu et al. Diversified Outlier Exposure for Out-of-distribution Detection via Informative Extrapolation. In NeurIPS, 2023.

[b] Lee et al. Training Confidence-calibrated Classifiers for Detecting Out-of-distribution Samples. In ICLR, 2018.

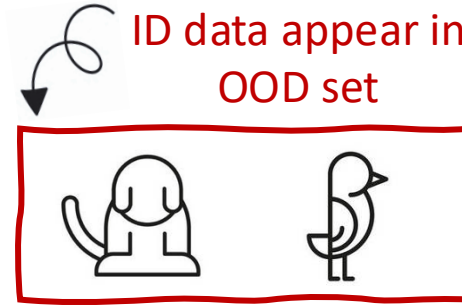
[c] Wang et al. Out-of-distribution Detection with Implicit Outlier Transformation. In ICLR, 2023.

Challenges and Future Directions

❖ Wild OOD Detection



training ID



training OOD



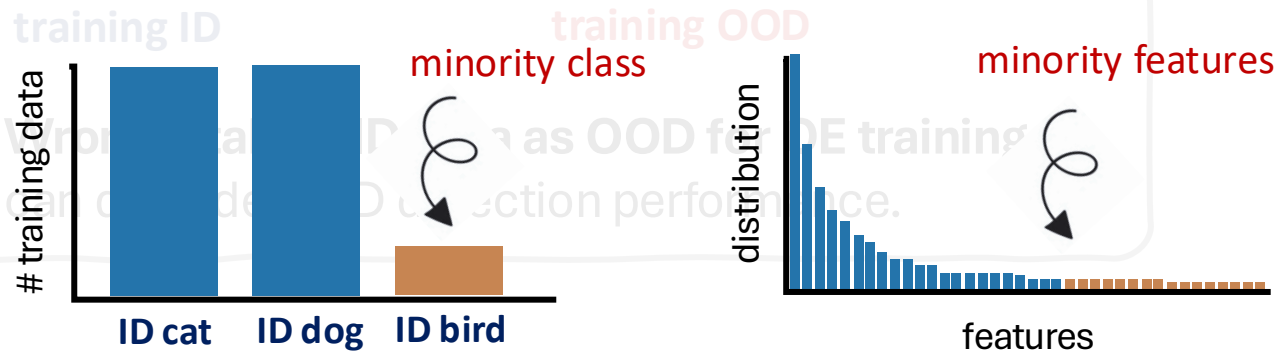
Wrongly taking ID data as OOD for OE training
can degrade OOD detection performance.

Challenges and Future Directions

❖ Wild OOD Detection

ID data appear in
OOD set

❖ Imbalanced and Long-tailed OOD Learning



Minority classes or features are easier to be confused with OOD data.

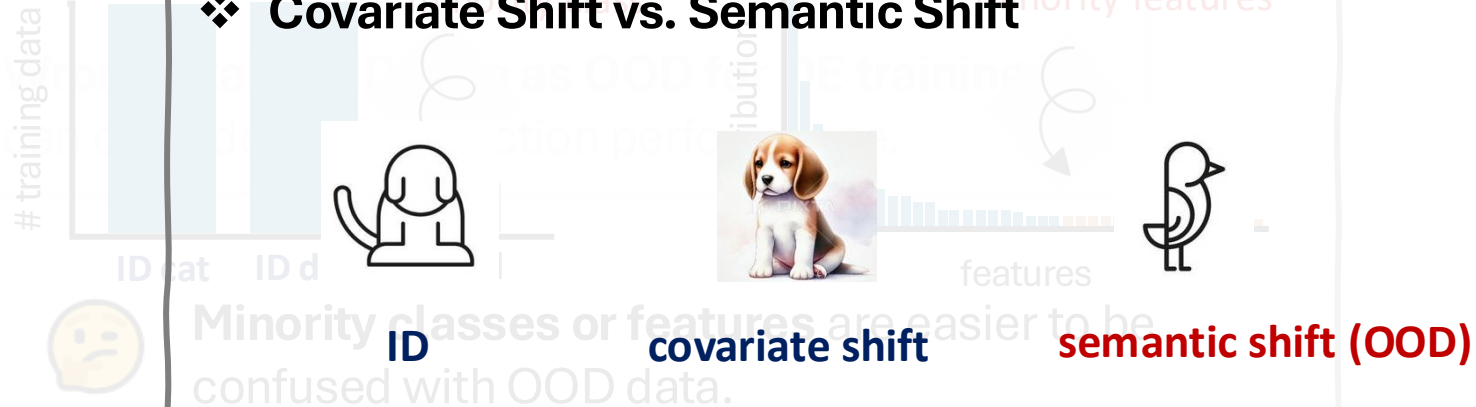
Challenges and Future Directions

❖ Wild OOD Detection

ID data appear in OOD set

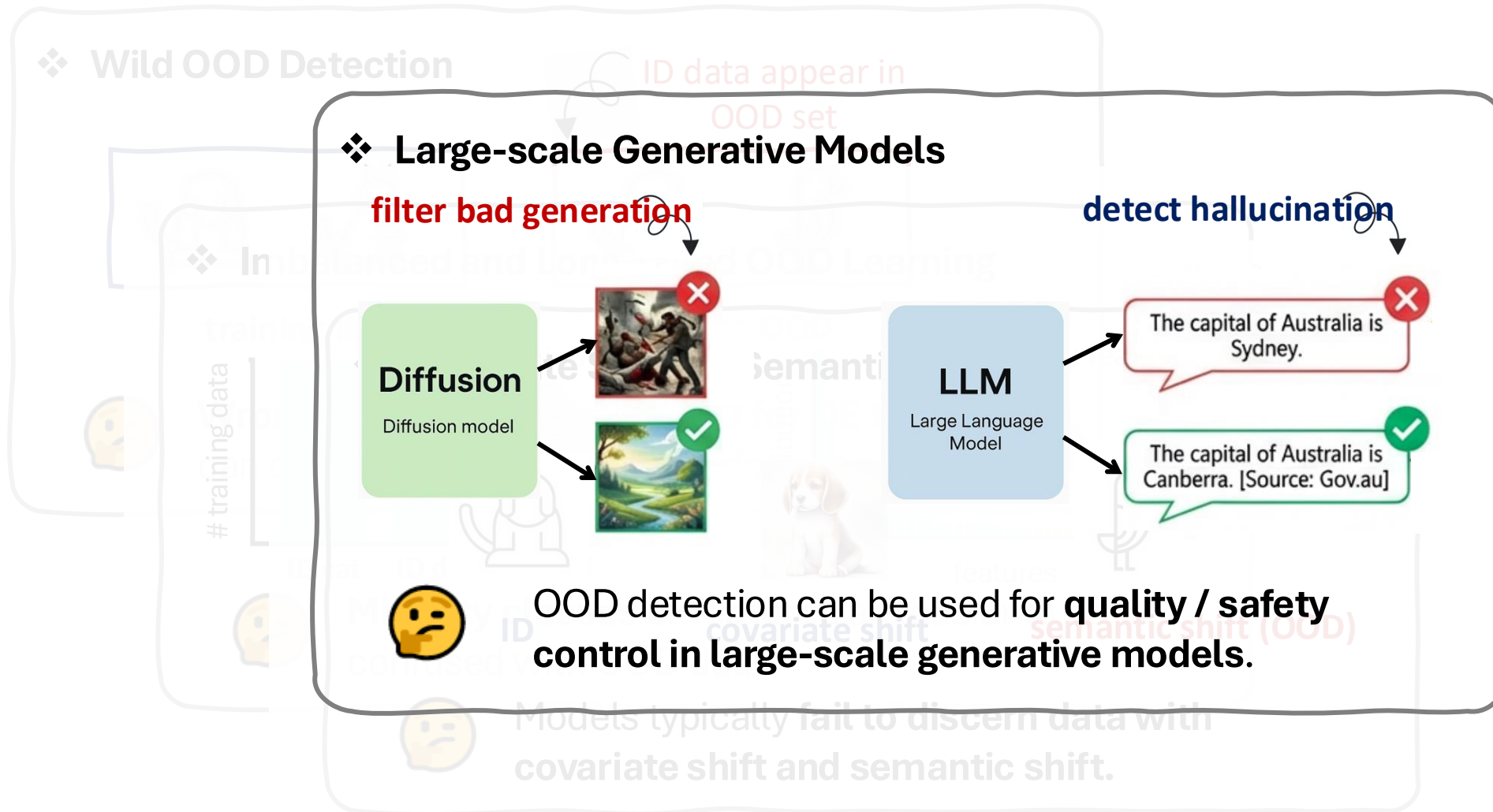
❖ Imbalanced and Long-tailed OOD Learning

❖ Covariate Shift vs. Semantic Shift



Models typically **fail to discern data with covariate shift and semantic shift.**

Challenges and Future Directions



Thank you for listening!

Find my slides from my homepage:
<https://qizhouwang.github.io/homepage>

