

Uncertainty Quantification (UQ) in AI

Bellairs Workshop, December 15-18, 2025

Nicolas THOME – Prof at Sorbonne University

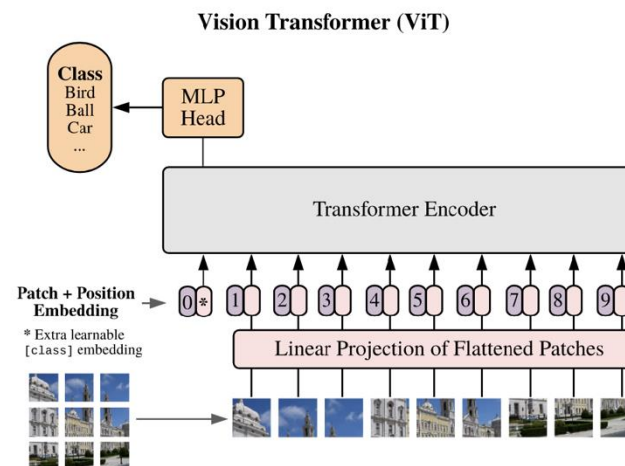
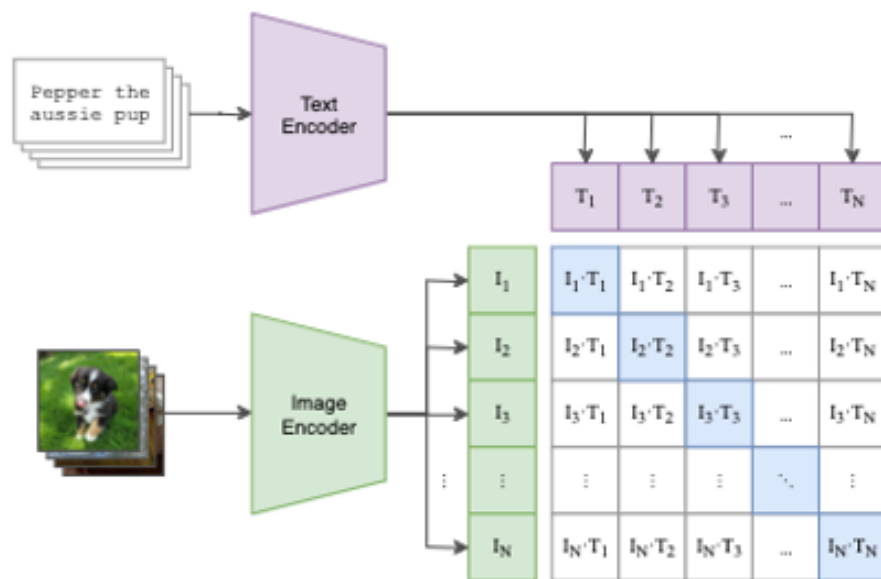
ISIR Lab, MLIA Team

On leave at ILLS Lab, ETS/Mila/McGill



Context: AI/ML summer

- AI in the last decade: huge performance boost
 - Vision, NLP, multi-modal prediction robotics

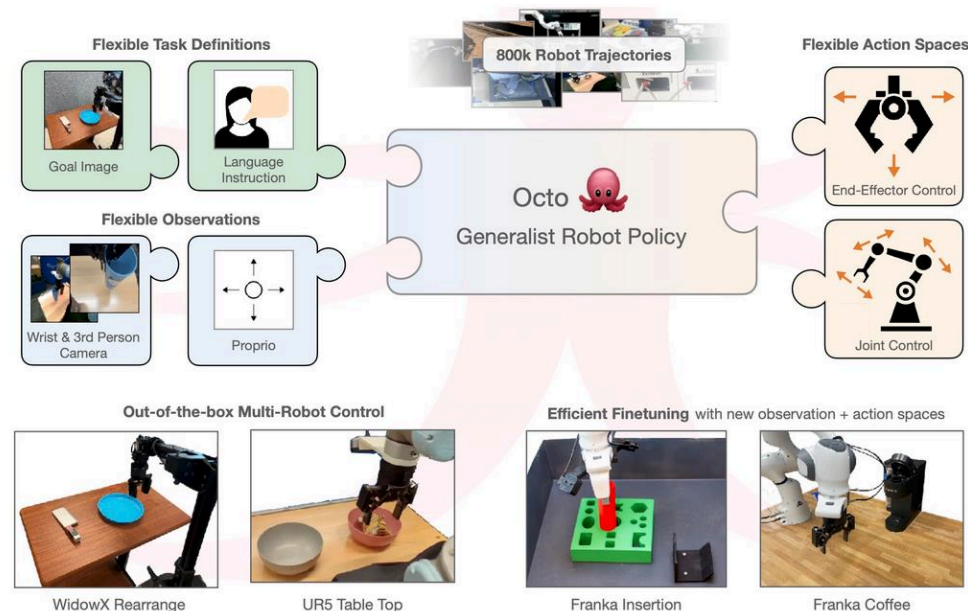


explain quantum theory to a child in the style of snoop dogg

Yo, little homie, let me break it down for you
 Quantum theory, it's a crazy thing, it's true
 See, the world is made up of tiny, tiny particles
 That are always moving, and changing, and actin' wild

These particles, they don't always behave like we think
 Sometimes they're waves, and sometimes they're points, it's a trip
 They can be in two places at once, and they can be in two states
 It's a mind-blowing concept, but that's the quantum fate

So don't be scared, little homie, it's all good
 Quantum theory, it's just the way the world works, it's understood
 It's a wild ride, but it's a cool one too
 Quantum theory, it's the way the world does what it do.



Context: robustness in deep learning

Several brittleness aspects in deep learning models

- Explainability, biases & shortcuts, fairness, *etc*

Stability: adversarial examples, mistake severity

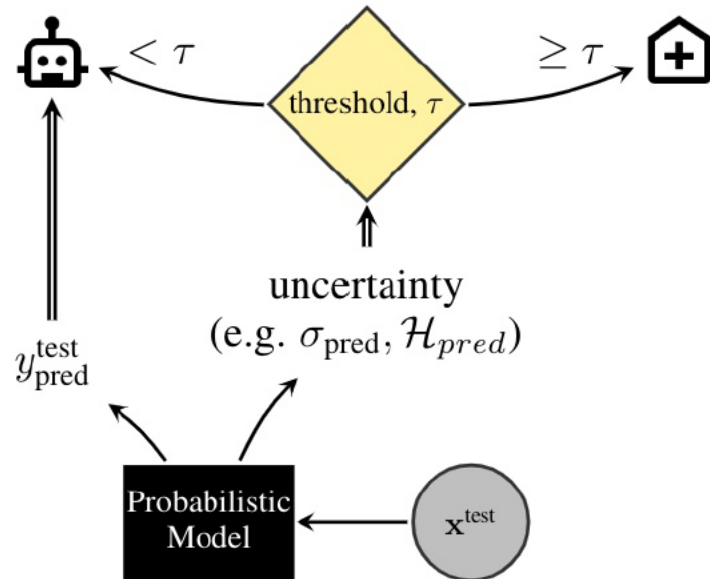
Query image



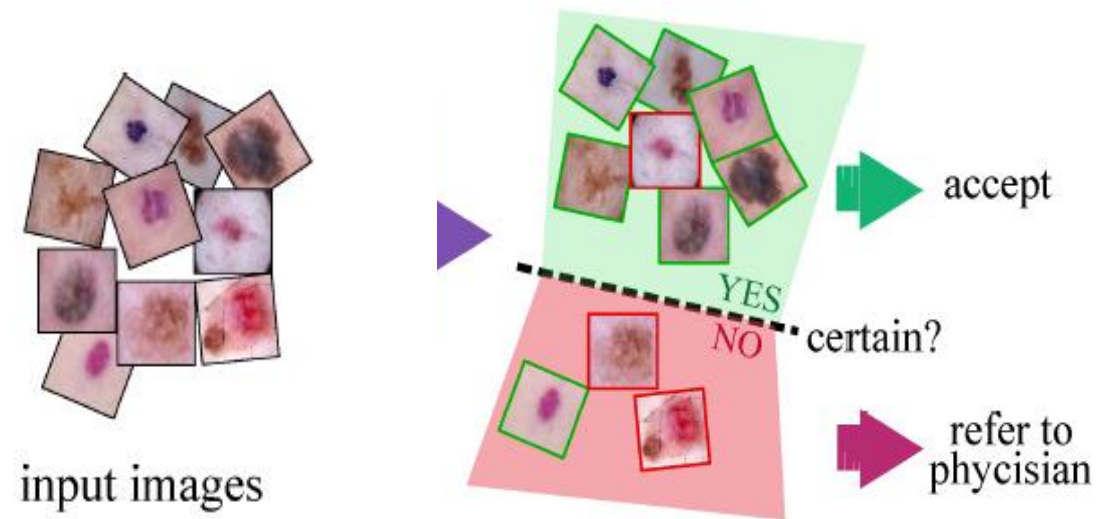
Context: robustness in deep learning

Uncertainty Quantification (UQ)

“Know when you do not know”



Abstain to make a prediction



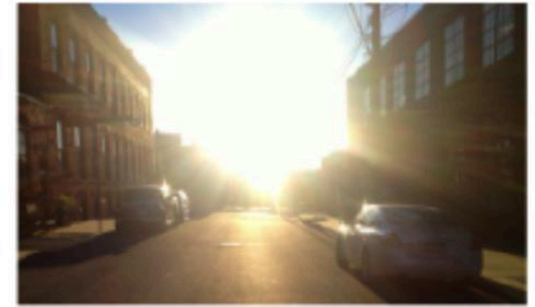
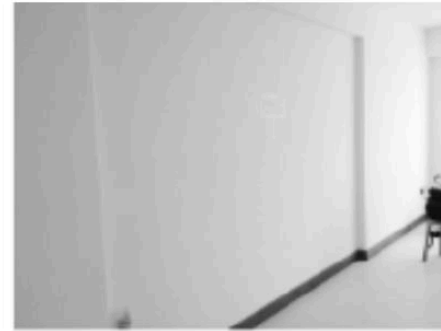
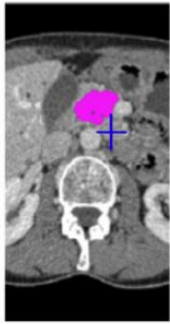
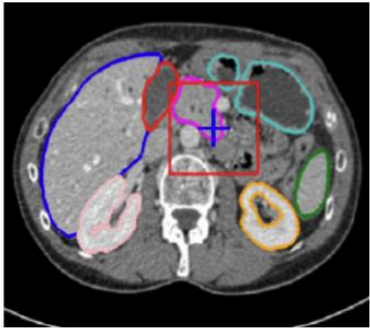
UQ: a challenge in DL

- Which uncertainty score?
- What type of uncertainty: aleatoric, epistemic?
- Which tasks: calibration, failure prediction, anomaly detection?

Sources of uncertainty

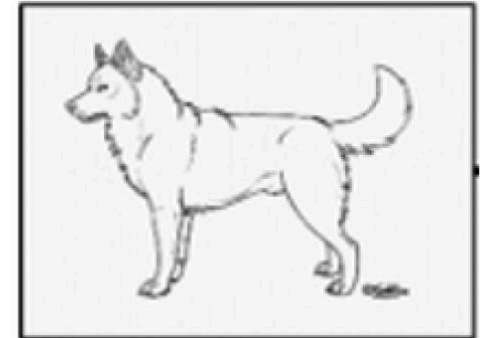
- Aleatoric uncertainty: data

- Class confusion, ambiguous data, sensor noise



- Epistemic uncertainty: model

- Distribution shift in $p(x,y)$, e.g. x (snow, image \rightarrow cartoon), or y (open set, new classes)



Uncertainty Quantification in deep learning

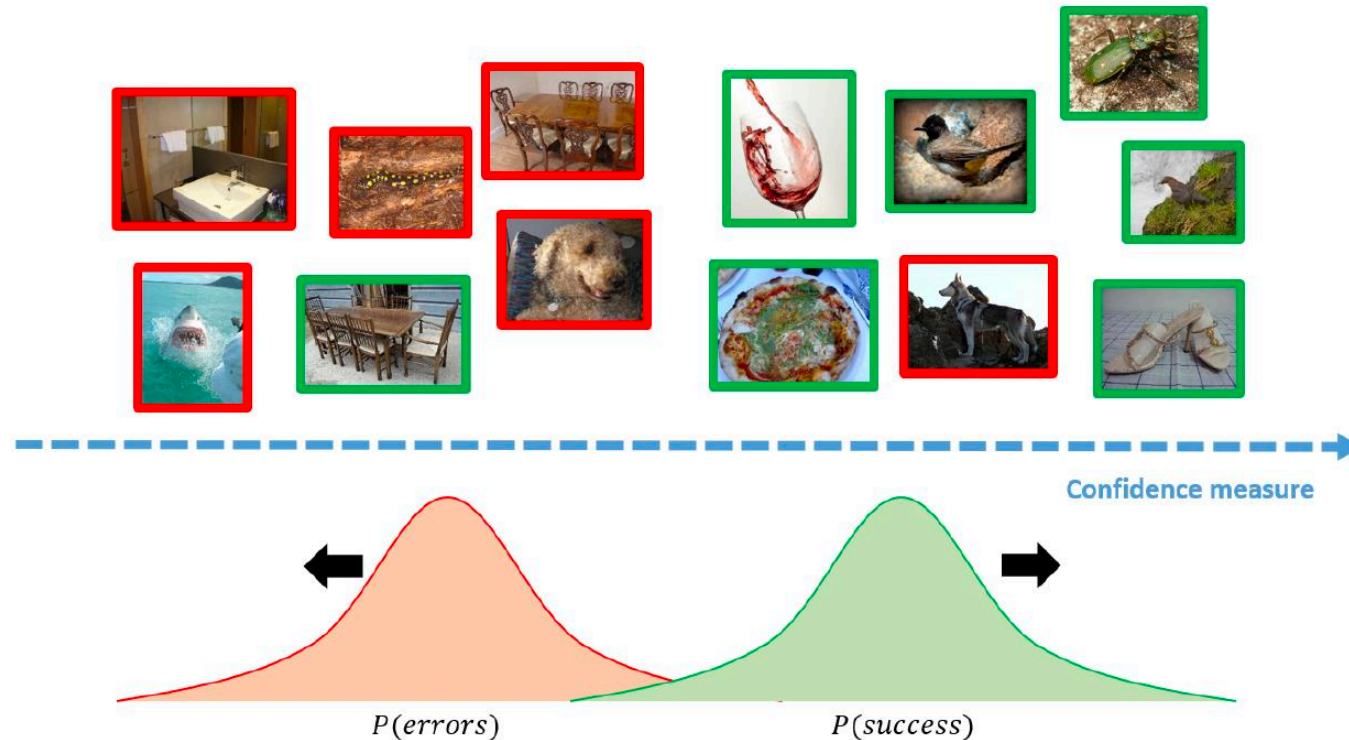
1. Calibration

2. Out-of-distribution / anomaly detection

3. Failure Prediction

Uncertainty \leftrightarrow confidence quantification

- Confidence estimate $\hat{C}(\mathbf{x}_i)$ goal: **distinguish correct from erroneous predictions**



- Sort examples wrt $\hat{C}(\mathbf{x}_i)$
 - Evaluate capacity of \hat{C} to assign larger prediction values for correct predictions than for errors

Calibration

Calibrated probabilities: estimated confidence $C(x) \Leftrightarrow$ accuracy

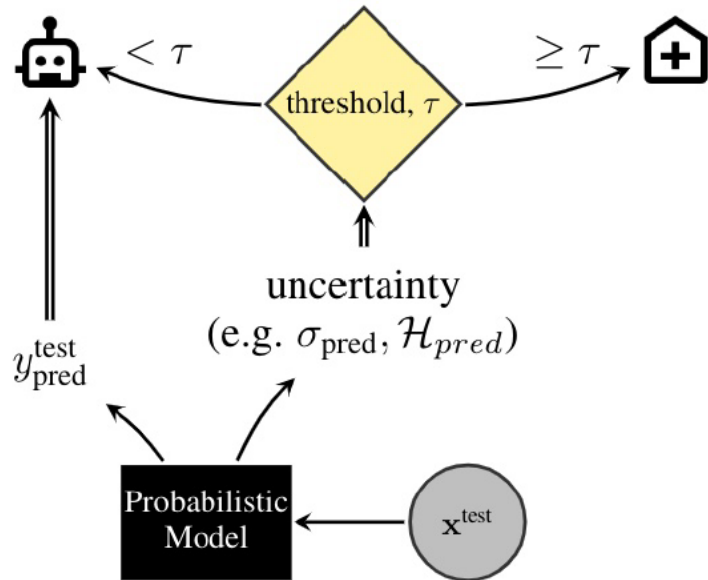
- Test **Data:** $\mathcal{D} = (X, Y) = \{(x_1, y_1), \dots, (x_N, y_N)\}$
- $(\hat{y}_i, \hat{C}(\mathbf{x}_i))$ class prediction and confidence level
- Perfect calibration:

$$p(\hat{Y} = Y | \hat{C} = p) = p, \quad \forall p \in [0, 1]$$

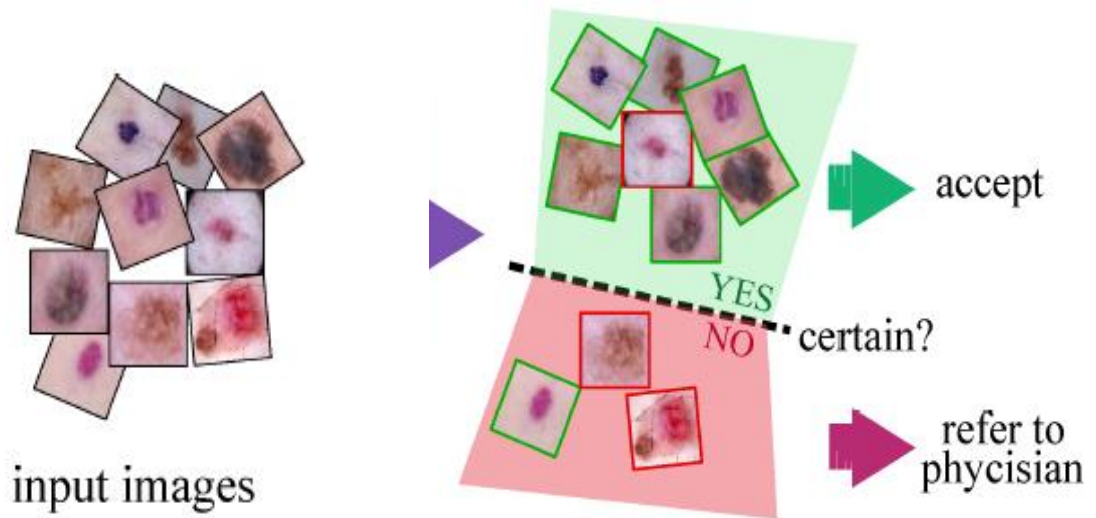
- Predicted confidence match actual accuracy
 - ▶ *e.g. given 100 predictions, each with confidence of 0.8, we expect that 80 should be correctly classified.*

Why does calibration matters?

“Know when you do not know”



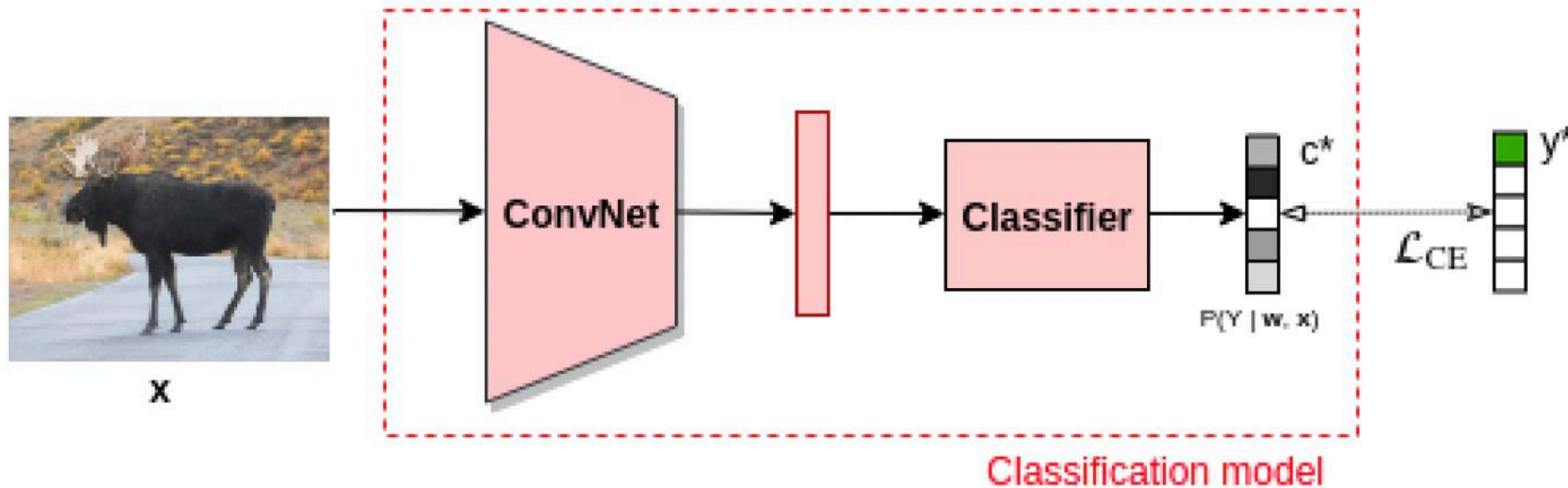
Abstain to make a prediction



- Predicted confidence $c(x) = \text{accuracy} \Rightarrow c(x) > \tau$: statistical guarantees

Baseline UQ score: classification example

- Classification model trained on $\mathcal{D} = \{(\mathbf{x}_i, y_i^*)\}_{i=1}^N$



- Model prediction: $\hat{\mathbf{y}} = \arg \max_{k \in \mathcal{Y}} p(Y = k | \mathbf{w}, \mathbf{x})$
- Model confidence $\hat{C}(\mathbf{x})$:
 - ▶ Simple baseline for deep neural networks: $\mathbf{MCP}(\mathbf{x}) = \max_{k \in \mathcal{Y}} p(Y = k | \mathbf{w}, \mathbf{x})$

Estimating calibration of confidence scores

Reliability Diagram

M interval bins.

B_m set of samples whose predictions are in $I_m = (\frac{m-1}{M}, \frac{m}{M}]$

$$acc(B_m) = \frac{1}{|B_m|} \sum_{i \in B_m} 1(\hat{y}_i = y_i)$$

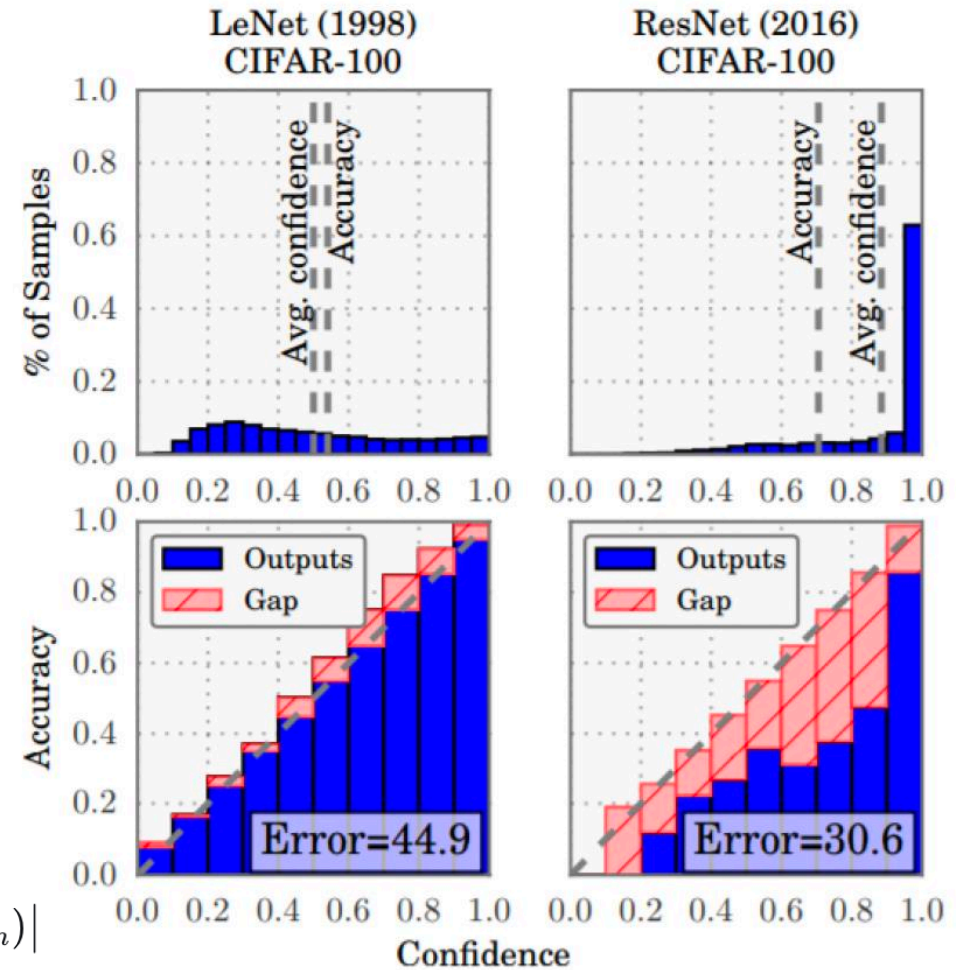
$$conf(B_m) = \frac{1}{|B_m|} \sum_{i \in B_m} \hat{p}_i$$

perfect calibration

$$\iff acc(B_m) = conf(B_m) \quad \forall m$$

Expected calibration error: $ECE = \sum_{m=1}^M \frac{|B_m|}{N} \cdot |Acc(B_m) - Conf(B_m)|$

[Guo et al., 2017] **showed that modern neural networks are no longer well-calibrated!**

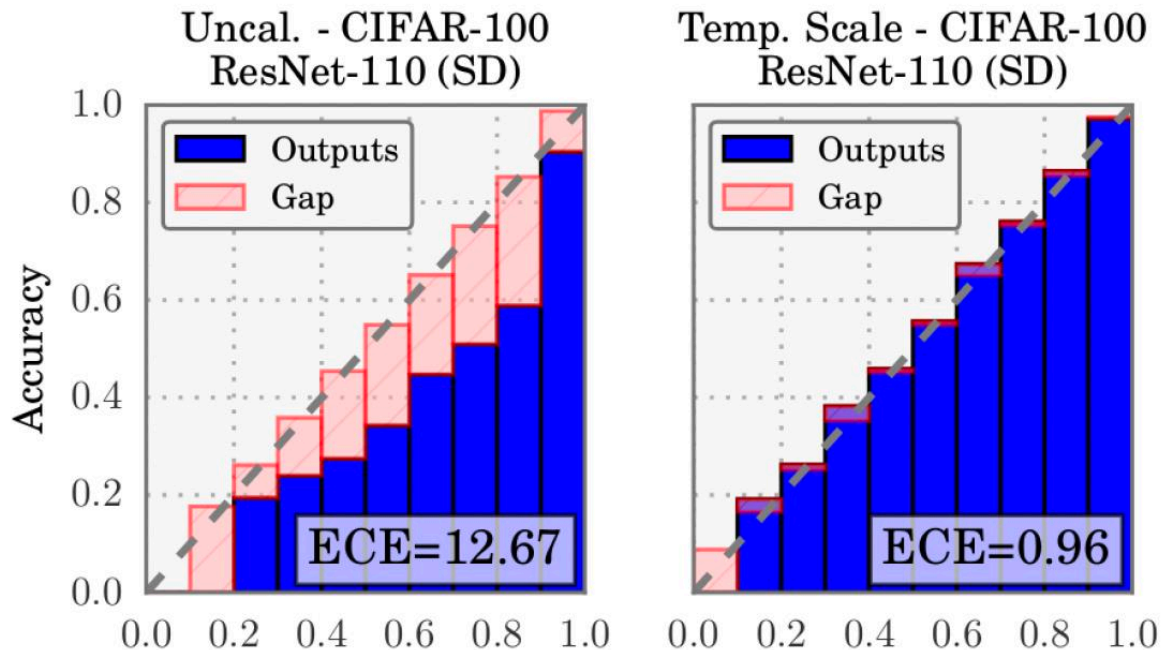


Post-hoc calibration

- Simple solution to over-confident prediction: temperature scaling [Guo et al., 2017]

$$P(\hat{y}_k) = \frac{e^{s_k/T}}{\sum_{k'=1}^K e^{s_{k'}/T}}$$

- temperature T optimized on val set s.t. $acc(B_m) = conf(B_m)$



How to evaluate calibration?

- Expected Calibration Error (ECE) for calibration only
- Proper scoring rules $S(p,q)$:
 - $S(p,q)$ is minimized if $p=q$. Ex: NLL or MSE for classification or regression
 - Using them for evaluating prediction models: include both discrimination (accuracy, AUC) and calibration
 - Brier Score (BS) & Continuous Ranked Probability Score (CRPS) for continuous outputs (regression)

$$\text{BS} = \frac{1}{N} \sum_{i=1}^N \|\mathbf{y}_i - \hat{\mathbf{y}}_i\|_2^2$$

- Negative Log-Likelihood (NLL) - Negative Log Predictive Density (NLPD)

Learning calibrated models

- **Fighting overconfidence:** controlling logit's scale, e.g., [Murugesan et.al., 2024]

$$\min_{\theta} \mathcal{L}_{CE} + \lambda \sum_k |\tau_k - l_k|$$

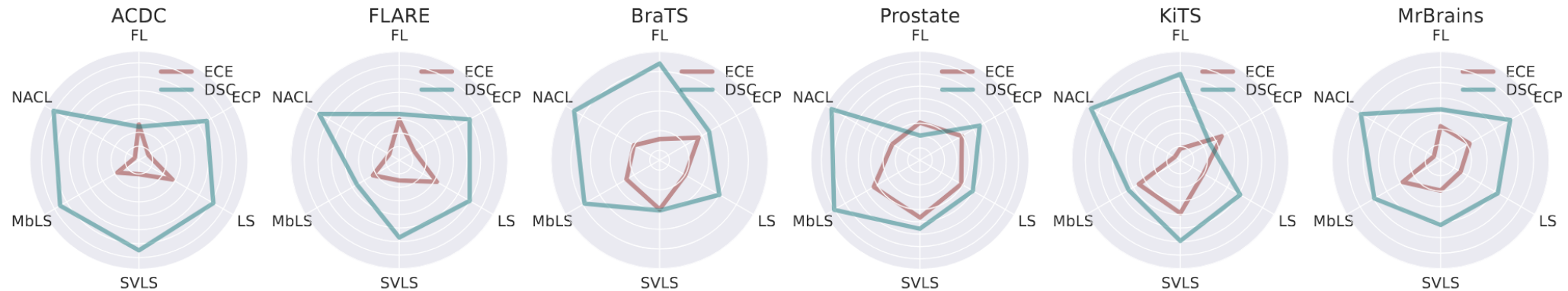


Figure 1: **Compromise between calibration and discriminative performance.** For each dataset, we show the discriminative (DSC) and calibration (ECE) results obtained by each method. We expect a *well-calibrated* model to achieve simultaneously large DSC (*in blue*) and small ECE (*in brown*) values.

[Murugesan et.al., 2024]Neighbor-Aware Calibration of Segmentation Networks with Penalty-Based Constraints. Balamurali Murugesana, Sukesh Adiga Vasudevaa, Bingyuan Liub, Herve Lombaerta, Ismail Ben Ayeda, Jose Dolz. Medical Image Analysis, 2024

Uncertainty Quantification in deep learning

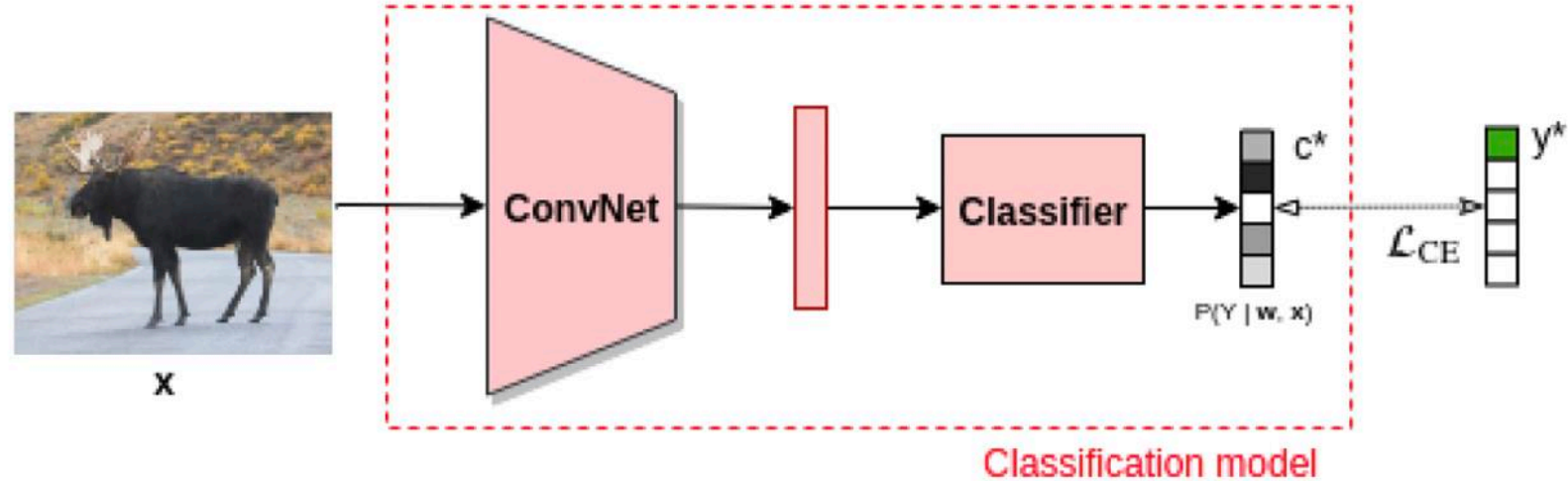
1. Calibration

2. Failure Prediction

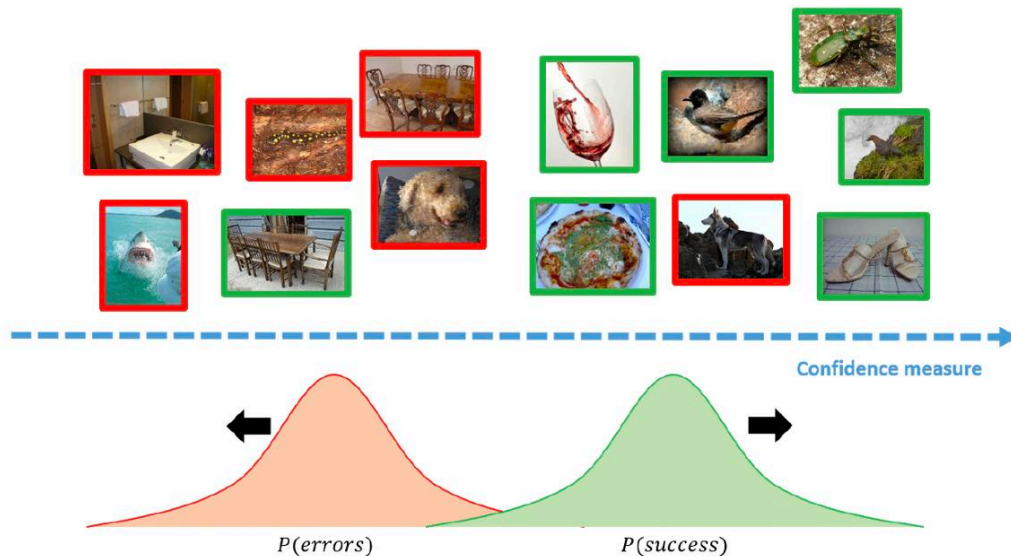
3. Out-of-distribution / anomaly detection

Failure prediction in classification

- **Post-hoc task:** UQ of a pre-trained model to detect its own errors



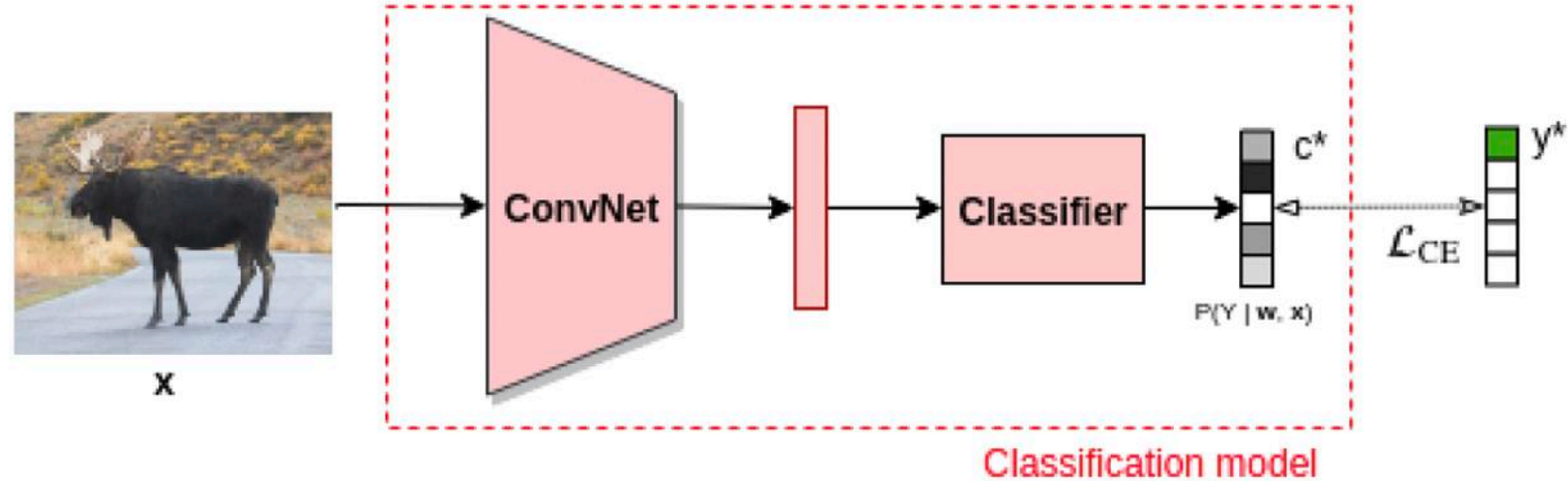
- A binary classification problem: correct vs incorrect prediction



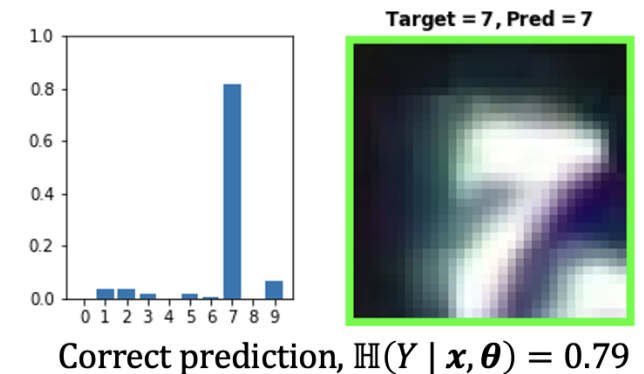
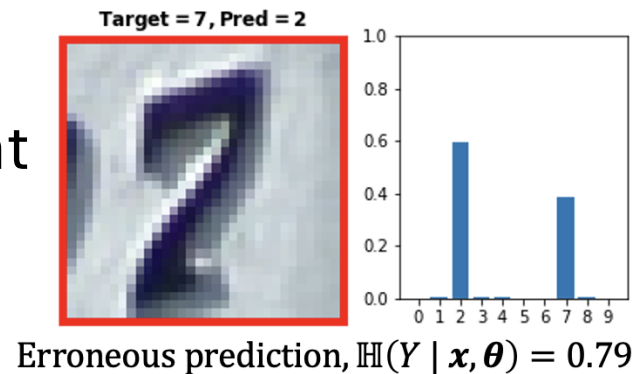
- Imbalance => evaluated by ranking metrics: AUC ROC, AUPR, FPR95, etc

Baseline methods for failure prediction

- $MCP(x) = \max_{k \in \mathcal{Y}} p(Y = k | \mathbf{w}, \mathbf{x})$
 - Overconfident by design



- Entropy: class-permutation invariant
 - Not targeted for failure prediction



$$\mathbb{H}(Y | \mathbf{x}, \boldsymbol{\theta}) := - \sum_{k \in \mathcal{Y}} P(Y = k | \mathbf{x}, \boldsymbol{\theta}) \cdot \log P(Y = k | \mathbf{x}, \boldsymbol{\theta})$$

Failure prediction methods

- DOCTOR [GRG+21]:

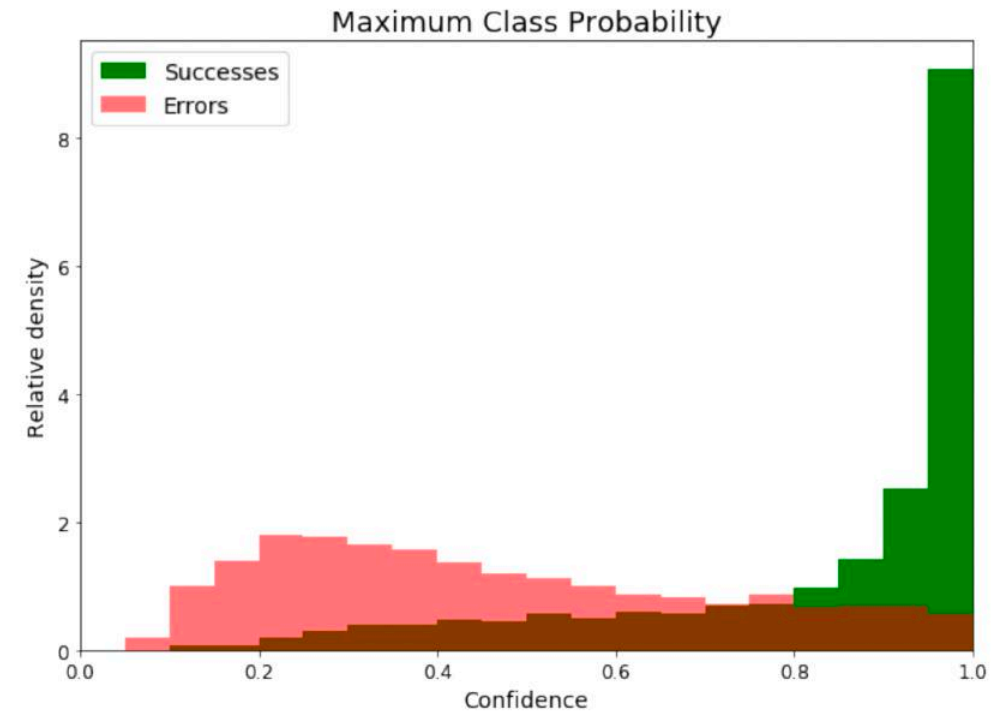
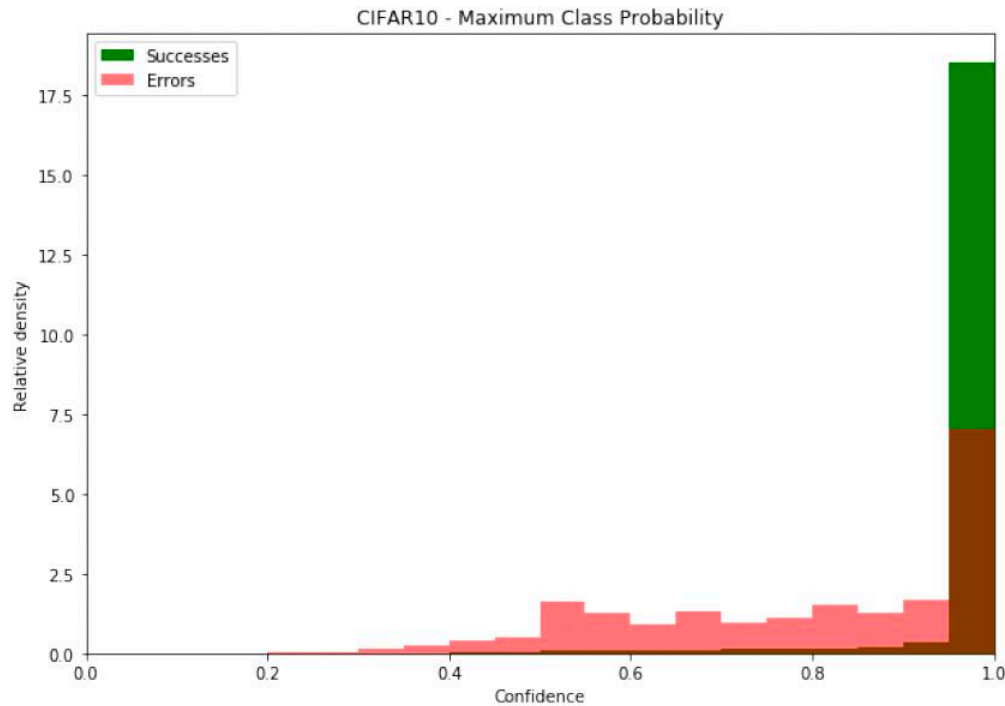
- $D_\alpha(\mathbf{x}, \gamma) \sim s_{\text{gini}}(\mathbf{x}) \triangleq 1 - \sum_{y \in \mathcal{Y}} (\hat{\mathbf{p}}(\mathbf{x})_y)^2 \Rightarrow \text{confidence} := ||\hat{P}||^2$
- $D_\beta(\mathbf{x}, \gamma) \sim \text{MCP}$

- Theoretical bounds for failure prediction

- Gini \mathbf{D}_α : still class-permutation invariant

[GRG+21] F. Granese, M. Romanelli, D. Gorla, C. Palamidessi, P. Piantanida. DOCTOR: A simple method for detecting misclassification errors. NeurIPS'21

Calibration & Failure Prediction

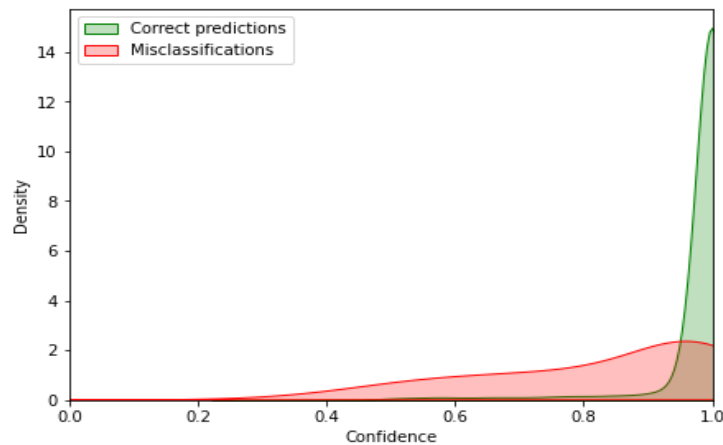


Calibration: absolute score vs

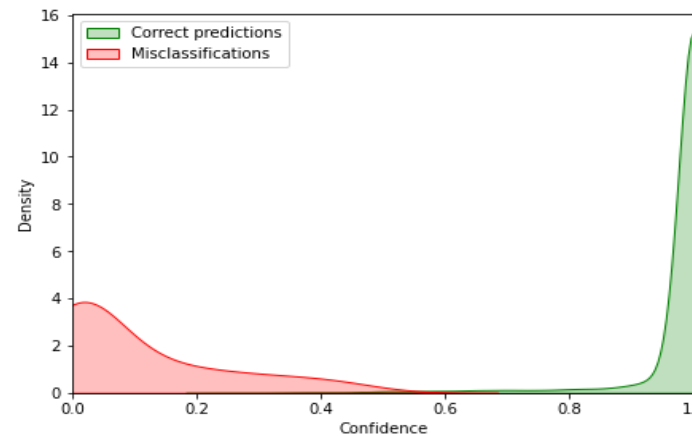
Failure prediction: relative score to distinguish correct from in correct predictions

Criterion for Failure Prediction

- **Idea [CTB+19]: True Class probability (TCP)** instead Maximum Class Probability (MCP) $TCP : P(Y = y^* | \mathbf{w}, \mathbf{x})$
- TCP better distinguish correct prediction from incorrect predictions MCP for failure prediction
 - Theoretical guarantees: $TCP(\mathbf{x}, y^*) > 1/2 \Rightarrow \hat{y} = y^*$ $TCP(\mathbf{x}, y^*) < 1/K \Rightarrow \hat{y} \neq y^*$



MCP

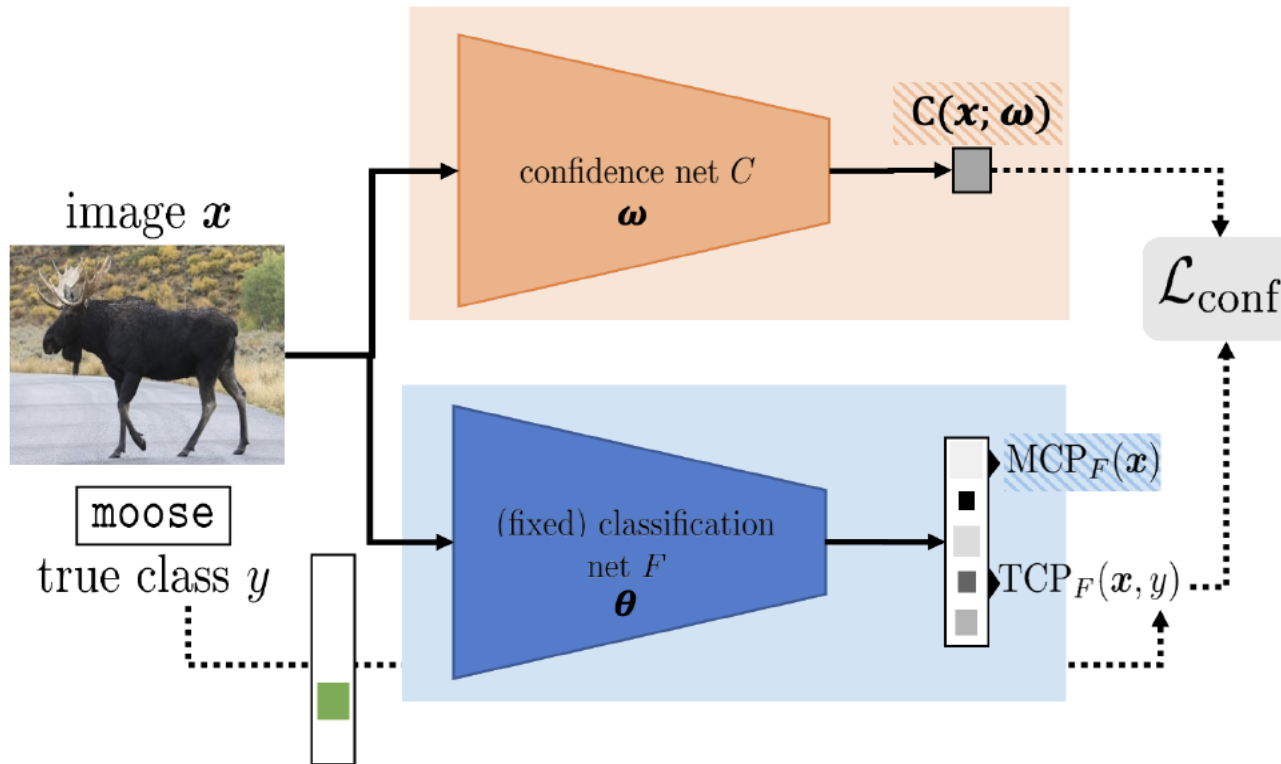


TCP

- Perfect separation with $m(x) = P(y^* | x) - \max_{y' \neq y^*} P(y' | x)$

Learning Confidence

TCP unknown at test time: learning it! => ConfidNet



- Pre-trained prediction model (blue)
- Learning to regress TCP with an auxiliary model (orange)

$$\mathcal{L}_{\text{conf}}(\theta; \mathcal{D}) = \frac{1}{N} \sum_{i=1}^N (\hat{c}(\mathbf{x}_i, \theta) - c^*(\mathbf{x}_i, y_i^*))^2$$

Same idea in different contexts: [YK19, K24] => **Learning Visual Uncertainties (LVU)**

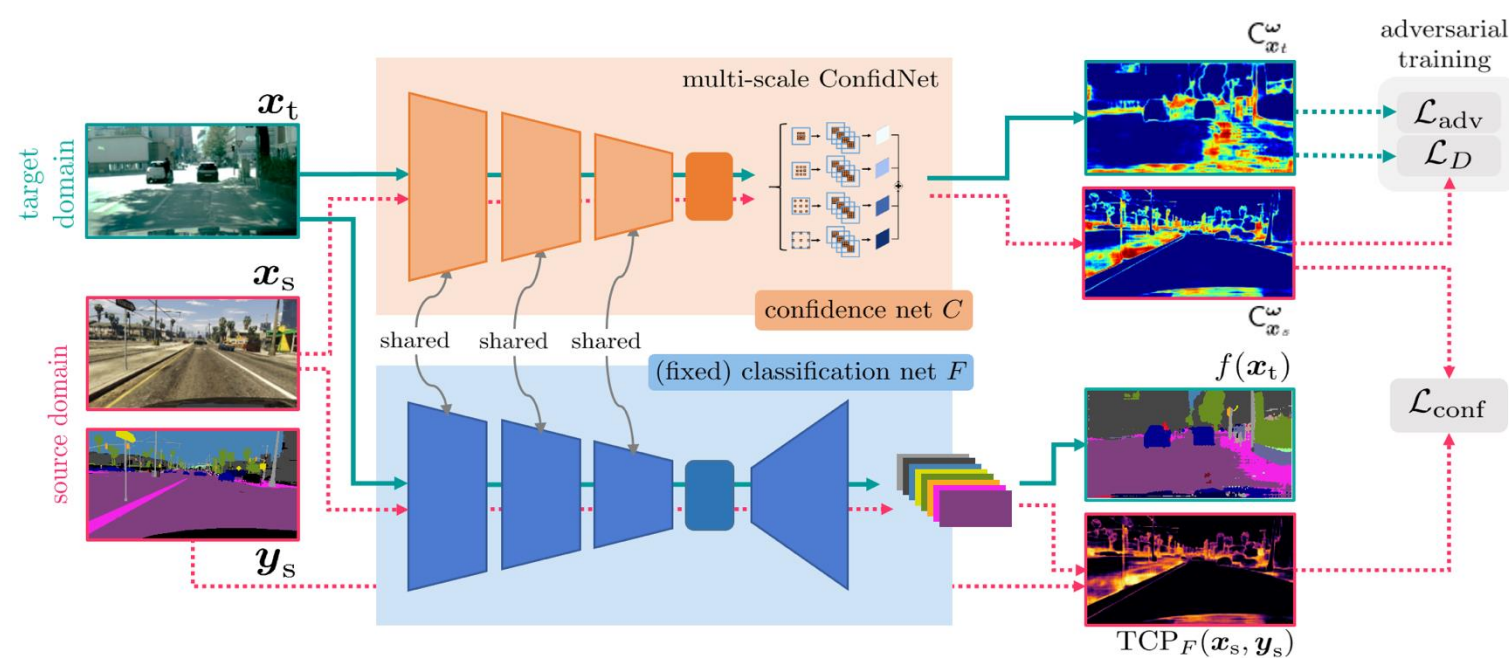
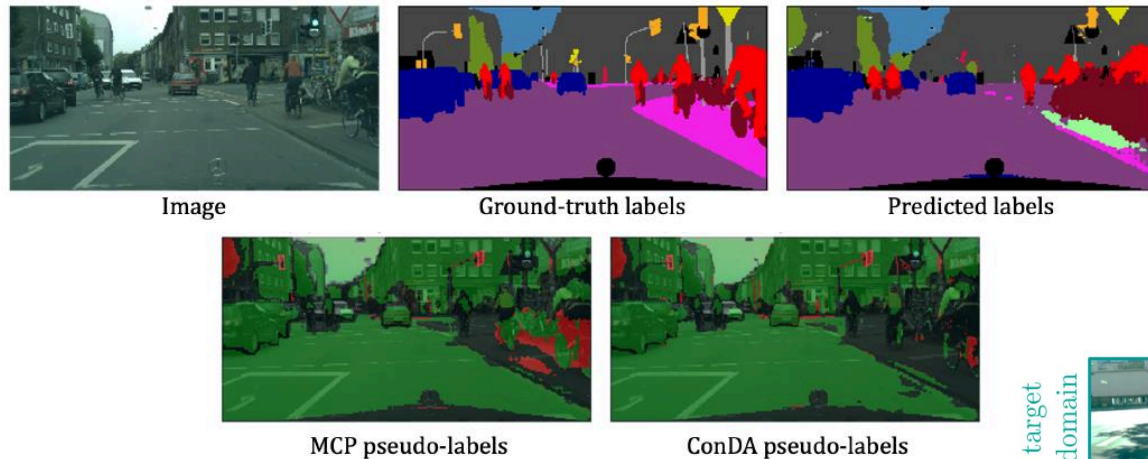
[CTB+19] C. Corbière, N. Thome, A. Bar-Hen, M. Cord, P. Pérez. Addressing Failure Detection by Learning Model Confidence. NeurIPS 2019.

[YK19] D. Yoo and In So Kweon. Learning loss for active learning. CVPR 2019

[K24] Kirchhof, Michael, et al. "Pretrained visual uncertainties." arXiv preprint 2024

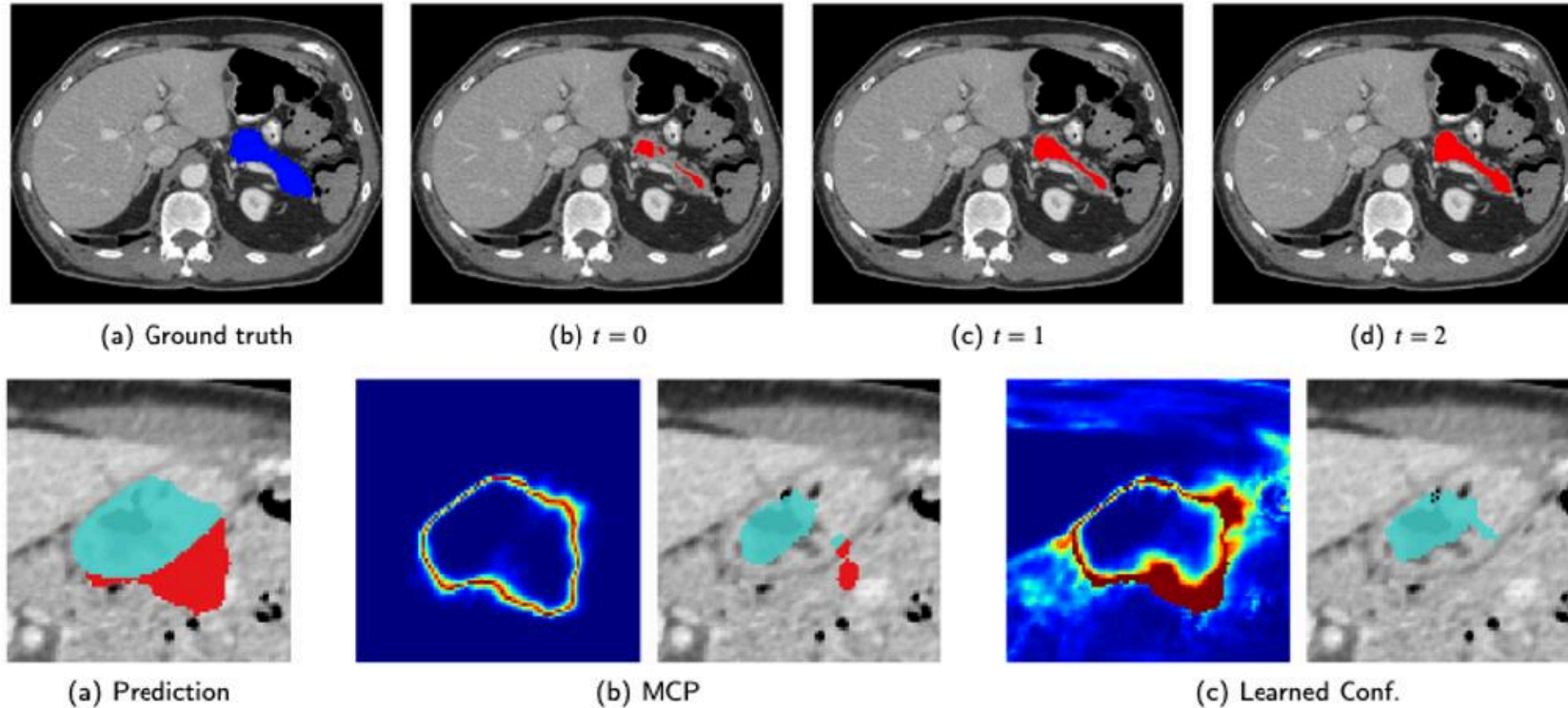
Learning confidence for self-labelling

- Extension for domain adaptation [CTS+21]



Learning confidence for self-labelling

- Extension for Medical image segmentation [PTS21]



Learning class distance matrix in classification

- Relative Uncertainty (Rel-U) [DRG+24]:

- Learn a trainable distance matrix : $s_d(\mathbf{x}) = \hat{\mathbf{p}}(\mathbf{x}) D \hat{\mathbf{p}}(\mathbf{x})^\top$
- Optimize to separate errors/correct predictions;

$$\mathcal{L}(D) \triangleq (1 - \lambda) \cdot \mathbb{E} \left[\hat{\mathbf{p}}(\mathbf{X}_+) D \hat{\mathbf{p}}(\mathbf{X}_+)^\top \right] - \lambda \cdot \mathbb{E} \left[\hat{\mathbf{p}}(\mathbf{X}_-) D \hat{\mathbf{p}}(\mathbf{X}_-)^\top \right]$$

$$\text{s.t.} \quad \left\{ \begin{array}{ll} d_{ii} = 0, & \forall i \in \mathcal{Y} \\ d_{ij} \geq 0, & \forall i, j \in \mathcal{Y} \\ d_{ij} = d_{ji}, & \forall i, j \in \mathcal{Y} \\ \text{Tr}(DD^\top) \leq K \end{array} \right.$$

Learning class distance matrix in classification

$$s_d(\mathbf{x}) = \hat{\mathbf{p}}(\mathbf{x}) D \hat{\mathbf{p}}(\mathbf{x})^\top = \sum_{y \in \mathcal{Y}} \sum_{y' \in \mathcal{Y}} d(y, y') \hat{\mathbf{p}}(\mathbf{x})_y \hat{\mathbf{p}}(\mathbf{x})_{y'}.$$

- Closed-form solution

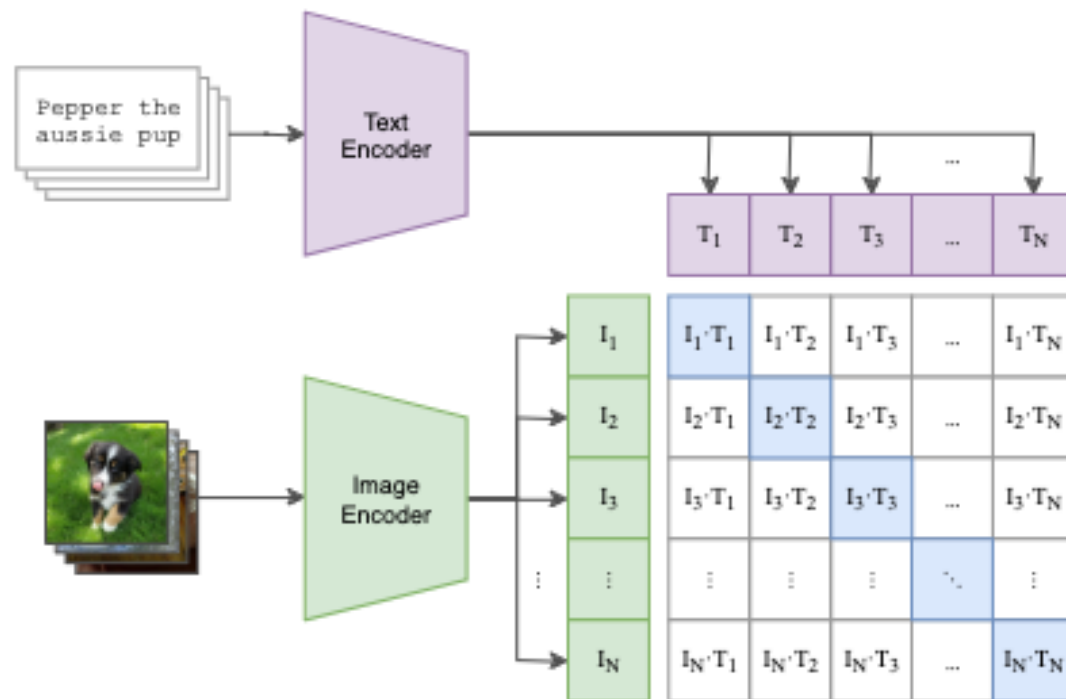
$$D^* = \text{ReLU} \left(\lambda \cdot \mathbb{E} \left[\hat{\mathbf{p}}(\mathbf{X}_-)^\top \hat{\mathbf{p}}(\mathbf{X}_-) \right] - (1 - \lambda) \cdot \mathbb{E} \left[\hat{\mathbf{p}}(\mathbf{X}_+)^\top \hat{\mathbf{p}}(\mathbf{X}_+) \right] \right)$$

- Distance matrix: limit the number of trainable parameters, suitable in few-shot context

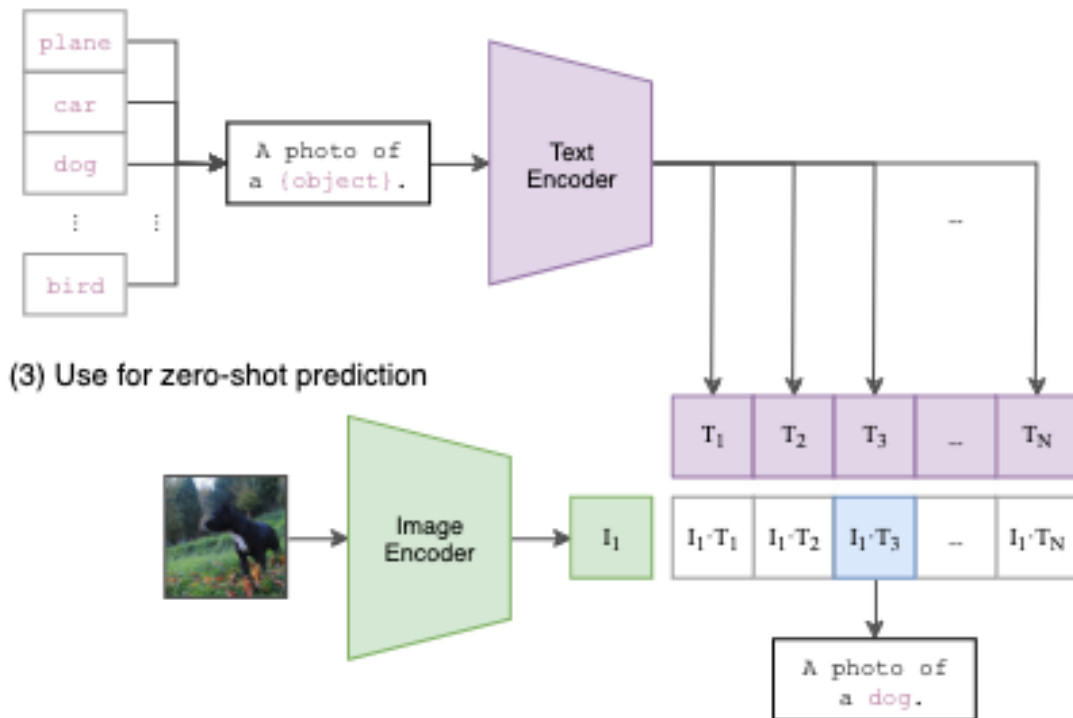
UQ for Vision-Language Models (VLMs)

- Vision-Language Models: e.g., text—image (CLIP) [3]

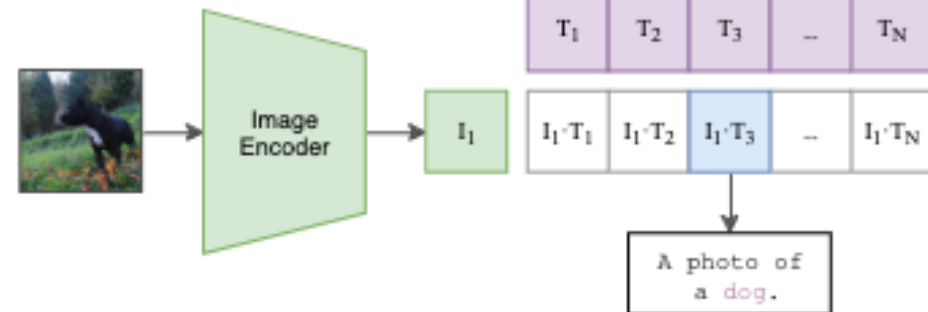
(1) Contrastive pre-training



(2) Create dataset classifier from label text

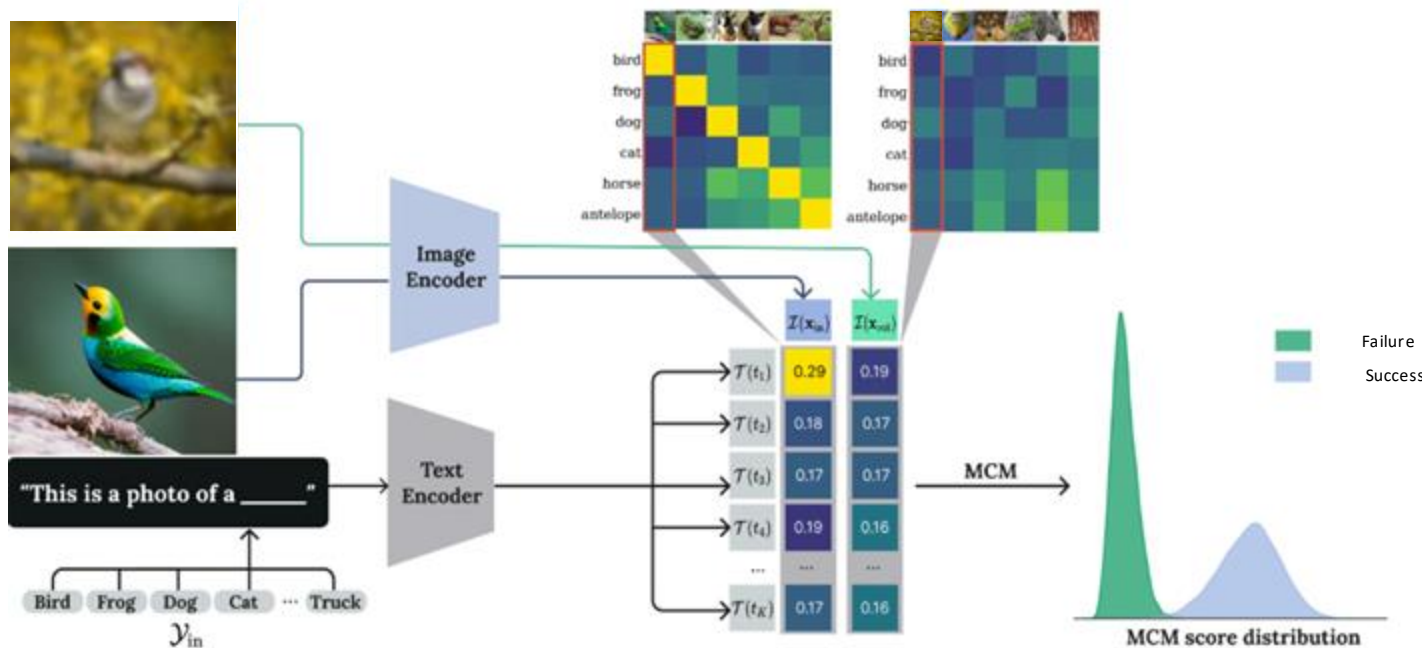


(3) Use for zero-shot prediction



Failure Prediction with Vision Language Models (VLMs)

- **Maximum Concept Matching (MCM)** = Probability of predicted class with VLMs
 - MCP Extension to VLMs



- Pros:
 - Strong baseline
 - No training required
- Cons:
 - Overconfident by design for errors
 - Limited adaptability

- Learning Visual Uncertainties (LVU): don't take into account text uncertainties

[CTB+19] C. Corbière, N. Thome, A. Bar-Hen, M. Cord, P. Pérez. Addressing Failure Detection by Learning Model Confidence. NeurIPS 2019.

[YK19] D. Yoo and In So Kweon. Learning loss for active learning. CVPR 2019

[K24] Kirchhof, Michael, et al. "Pretrained visual uncertainties." arXiv preprint 2024

Taking task complexity into account



What is the uncertainty associated with this image ?

→ It depends on the task

Task 1: cat vs dog classification



- Low class confusion
- Low uncertainty

Taking task complexity into account



What is the uncertainty associated with this image ?

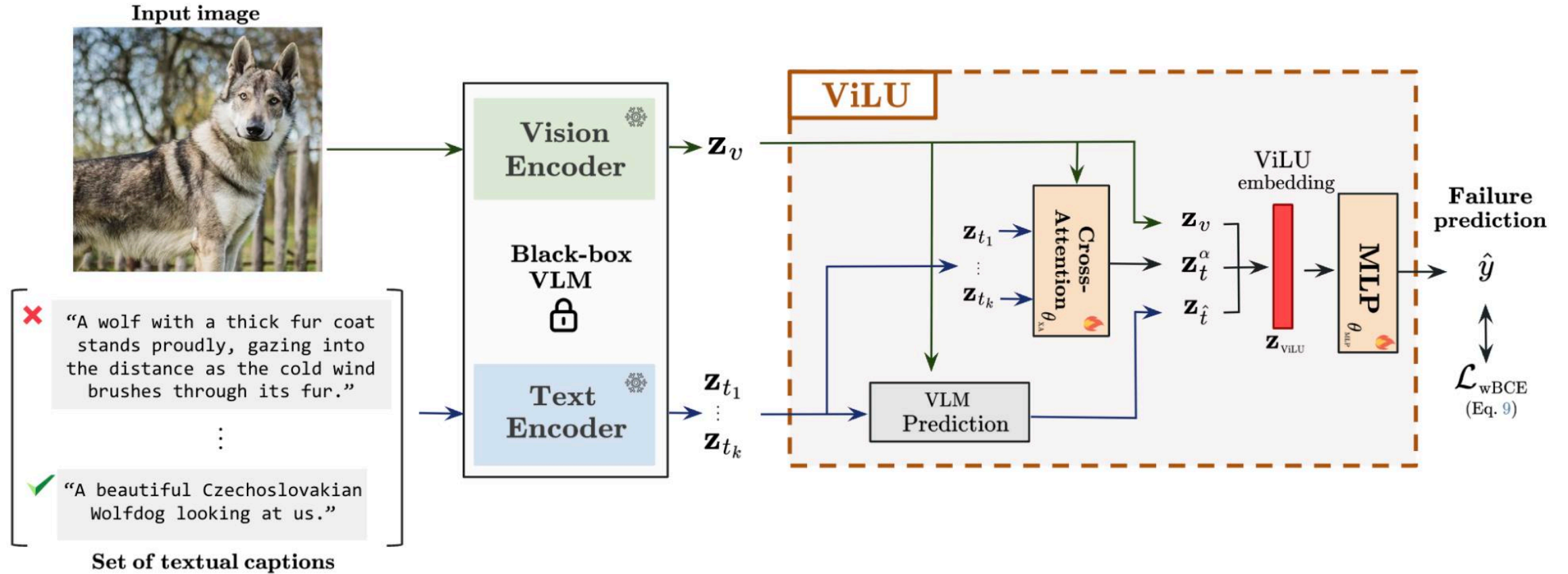
→ It depends on the task

Task 2: Dog breed classification



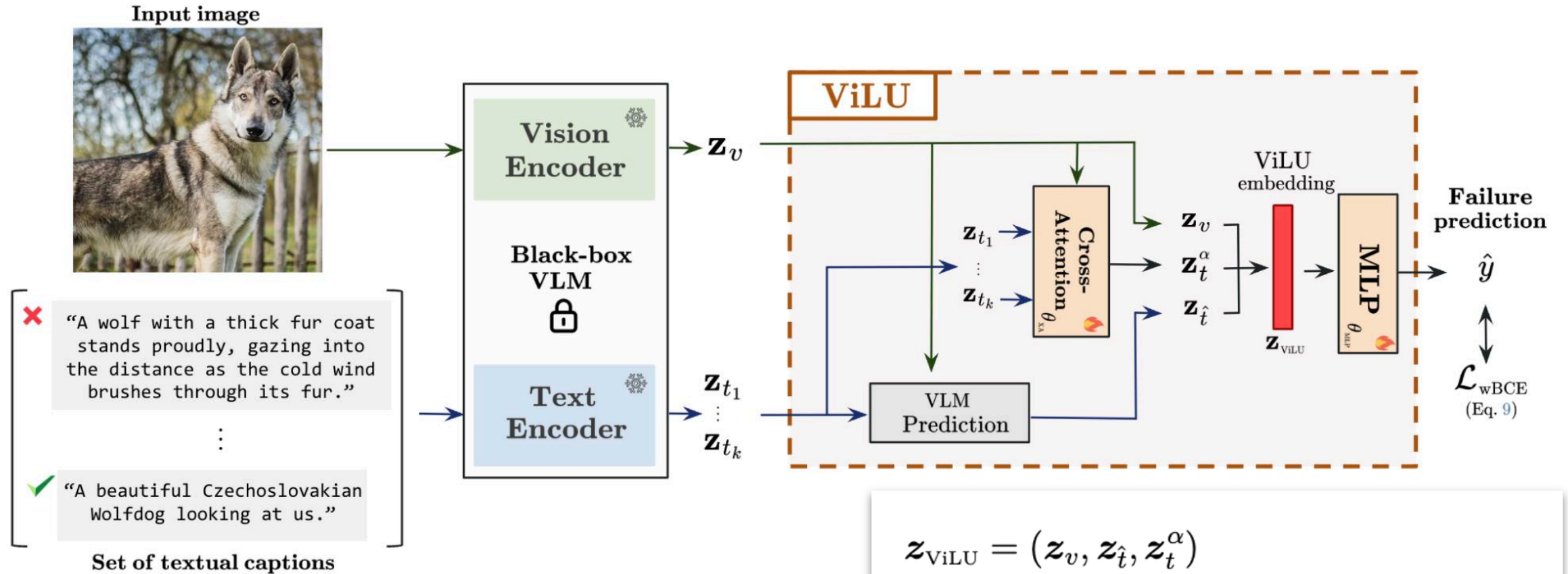
- Higher class confusion
- High uncertainty

ViLU: Learning Vision-Language Uncertainty



- Inputs: visual representation \mathbf{z}_v + K concept embeddings $\mathbf{Z}_t = \{\mathbf{z}_{t_j}\}_{1 \leq j \leq K}$
- ViLU embedding: $\mathbf{z}_{ViLU} = (\mathbf{z}_v, \mathbf{z}_{\hat{t}}, \mathbf{z}_t^\alpha), \mathbf{z}_{\hat{t}}$ pred. Class embedding
 - Image-text cross-attention module $\Rightarrow \mathbf{z}_t^\alpha$
 - Query \mathbf{z}_v , keys/values \mathbf{Z}_t

ViLU: Learning Vision-Language Uncertainty



- **Failure prediction from ViLU embedding:** binary classification
- Consistent generalization of MCM

$$\mathbf{z}_{\text{ViLU}} = (\mathbf{z}_v, \mathbf{z}_{\hat{t}}, \mathbf{z}_t^\alpha)$$

$$g_{\theta_{\text{MLP}}}(\mathbf{z}_{\text{ViLU}}) \approx \frac{1}{2} \mathbf{z}_{\text{ViLU}}^T A \mathbf{z}_{\text{ViLU}} = \mathbf{z}_v^T \mathbf{z}_{\hat{t}}$$

$$\text{with } A = \begin{pmatrix} 0 & \mathbf{I}_d & 0 \\ \mathbf{I}_d & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}$$

Conformal prediction & theoretical guarantees

- Matteo's talk!

Uncertainty Quantification in deep learning

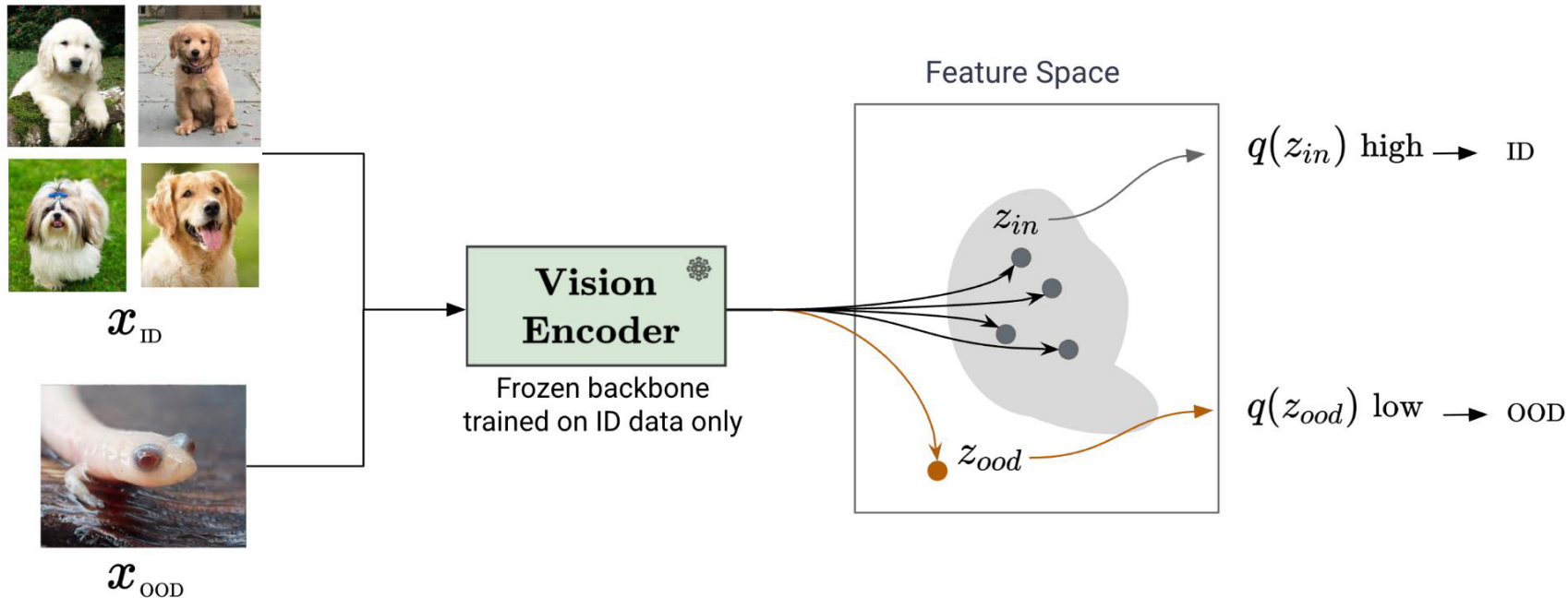
1. Calibration

2. Failure Prediction

3. Out-of-distribution / anomaly detection

OOD/ anomaly detection

- Detecting OODs/anomalies \Leftrightarrow epistemic uncertainty
 - Binary classification problem: in-distribution (ID) vs OOD
- 2 main classes of approaches for OOD
 1. Bayesian methods: estimating the predictive distribution $p(y|x^*, D)$
 2. Estimating ID density $p(x)$



- Density in the feature space of a pre-trained model
- State-of-the-art: GMM, KNN, EBM, diffusion models, etc

Sehwag, V., et al. "SSD: A unified framework for self-supervised outlier detection." ICML 2021.

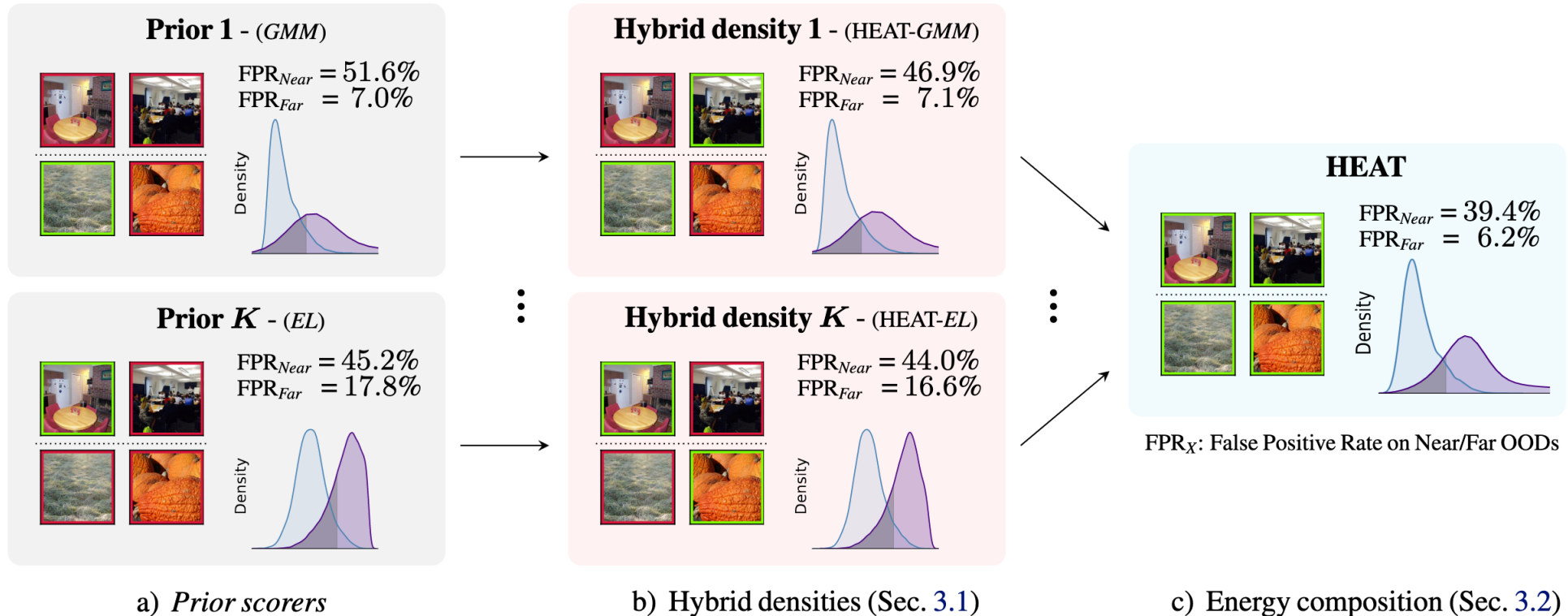
Sun, Yiyu, et al. "Out-of-distribution detection with deep nearest neighbors." ICML, 2022.

Lafon, et al. "Hybrid Energy Based Model in the Feature Space for Out-of-Distribution Detection." ICML 2023.

A, Heng et.al. "Out-of-Distribution Detection with a Single Unconditional Diffusion Model". NeurIPS 2024

Out-Of-Distribution (OOD) detection

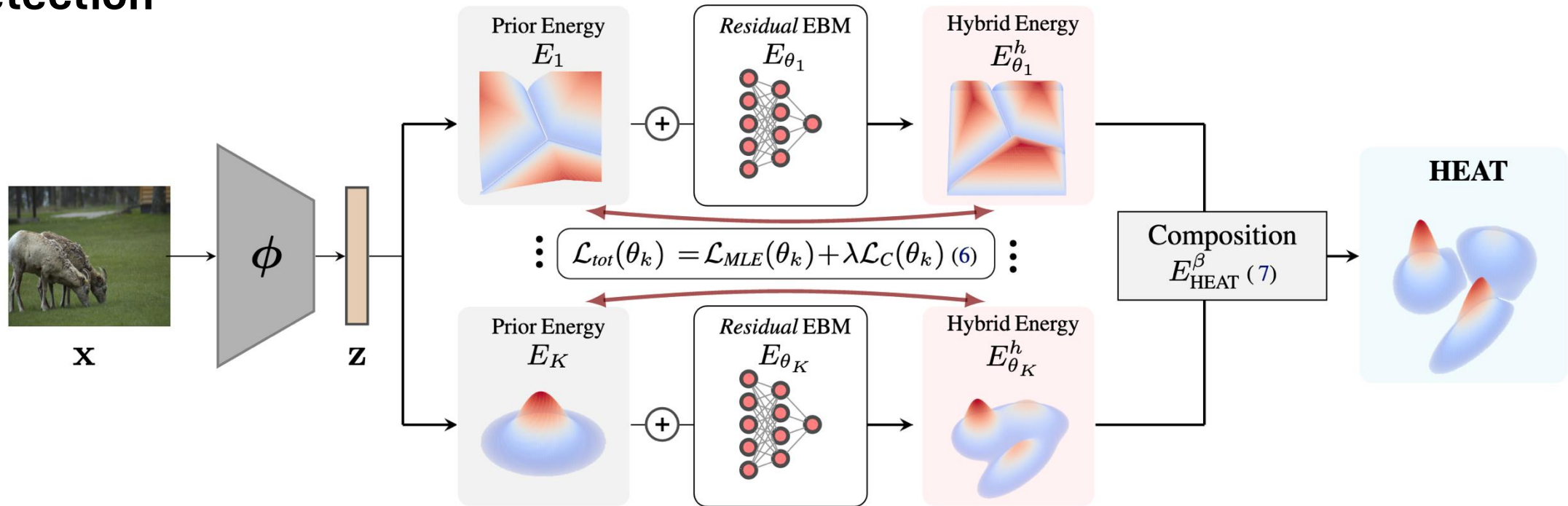
- State-of-the-art ID density estimation: prior densities, e.g., GMM, Energy Logits (EL)



- Prior density: not accurate => **Energy correction**
- GMM good for far-OOD, EL for near-OOD => **Energy composition**

OOD detection

- **HEAT [LRR+23]: Hybrid Energy Based Model (EBM) in the feature space for OOD detection**

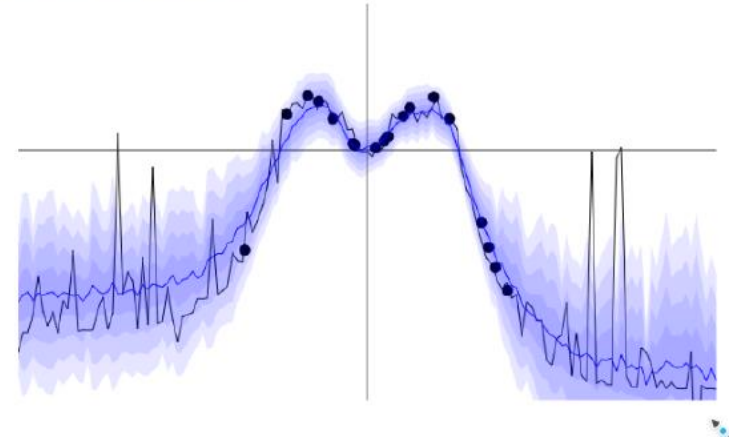


- **Energy-based correction** of prior energy terms, e.g. Gaussians

- **Energy composition** of several terms (Gaussian, Energy Logits, std for style)

Bayesian methods

We fit a **distribution**...



- Observed inputs $X = \{x_i\}_{i=1}^N$ and outputs $Y = \{y_i\}_{i=1}^N$
 - ▶ $x_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}^K$ (classification or regression)
 - ▶ Model with parameters w : $\hat{y}_i = f_w(x_i)$

- **Bayes rule:** $p(Y, w/X) = p(Y/X, w)p(w) = p(w/X, Y)p(Y/X)$

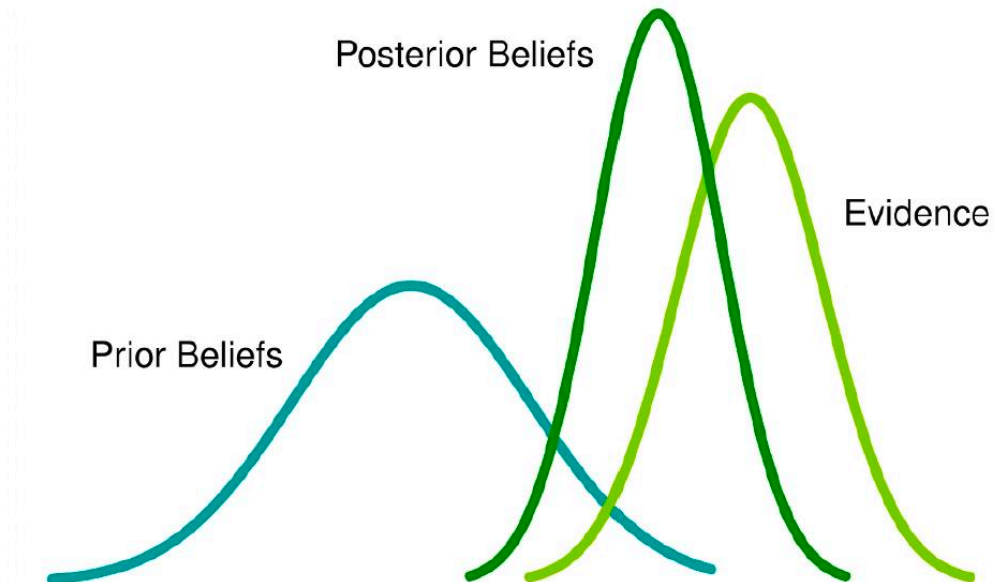
$$\Rightarrow \boxed{p(w/X, Y) = \frac{p(Y/X, w)p(w)}{p(Y/X)} \propto p(Y/X, w)p(w)}$$

- From posterior $p(w/X, Y) \Rightarrow$ compute **predictive distribution** given new input x^* ($\forall y$): $p(y/x^*, Y, X) = \int p(y, w/x^*, Y, X)dw$

$$\boxed{p(y/x^*, Y, X) = \int p(y/x^*, w)p(w/X, Y)dw} = \mathbb{E}_{p(w|\mathcal{D})}[p(y|x^*, w)]$$

Bayesian methods

- RECAP: for uncertainty estimate with Bayesian models
 1. Define prior $p(w)$ and likelihood $p(Y/X, w)$
 2. Compute posterior distribution $p(w/X, Y)$
 3. Compute predictive distribution $p(y^*/x^*, Y, X)$
- Easy? NO !! \Rightarrow steps 2 and 3 computationally hard in general!
 - Typically no closed form for step 2
 - High-dimensional integration for step 3



Bayesian Linear Regression

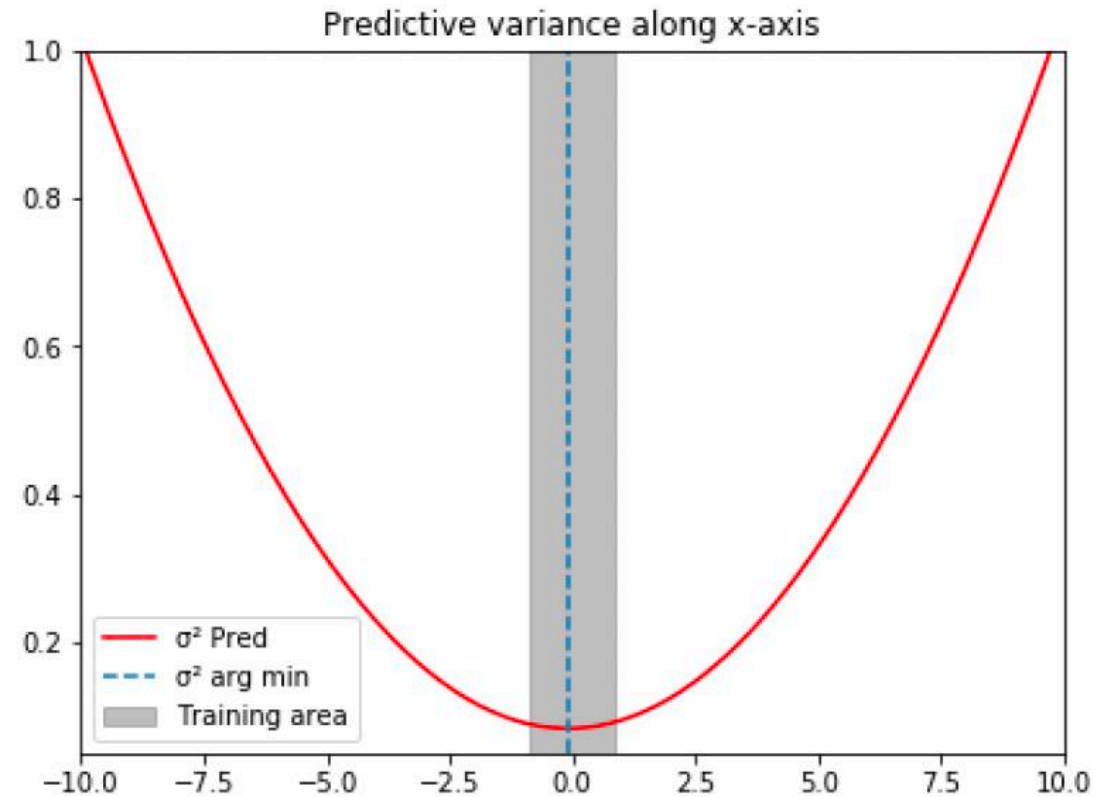
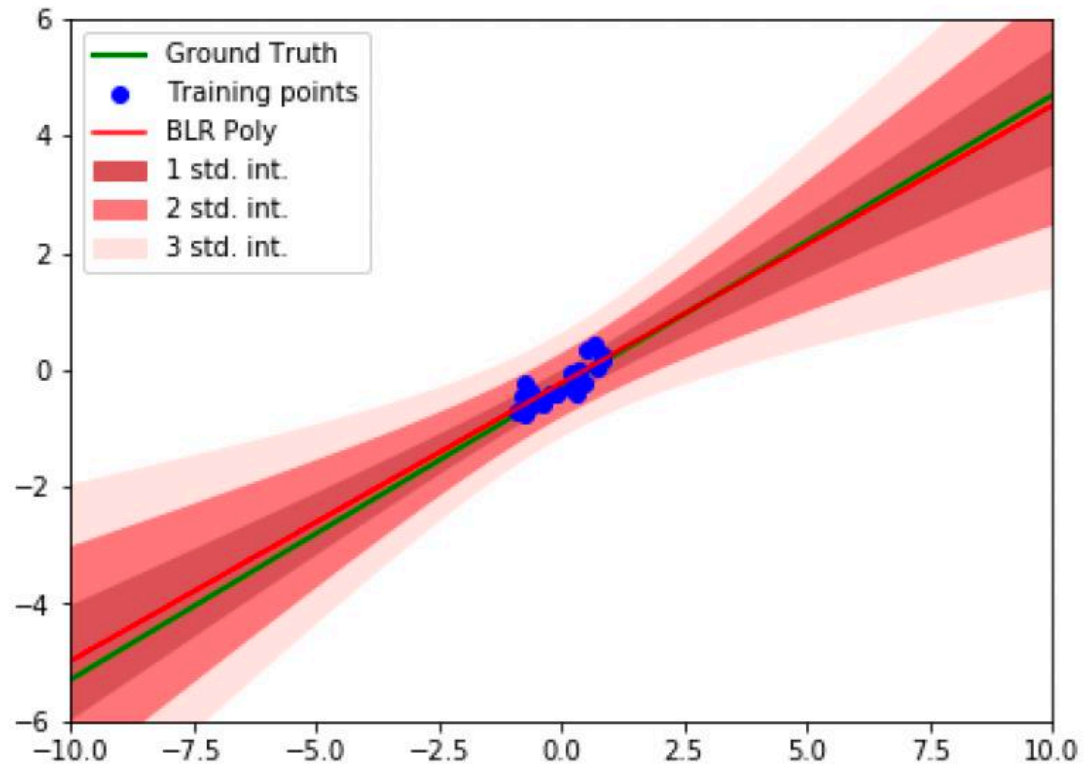
With $p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}; 0, \alpha^{-1}\mathbf{I})$ and $p(y_i/x_i, \mathbf{w}) = \mathcal{N}(\Phi_i^T \mathbf{w}, \beta^{-1})$, we can show that:

$$\begin{aligned} p(\mathbf{w}|X, Y) &= \mathcal{N}(\mathbf{w}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \\ \boldsymbol{\Sigma}^{-1} &= \alpha\mathbf{I} + \beta\boldsymbol{\Phi}^T\boldsymbol{\Phi} \\ \boldsymbol{\mu} &= \beta\boldsymbol{\Sigma}\boldsymbol{\Phi}^T\mathbf{Y} \end{aligned}$$

- $p(y|x^*, \mathcal{D}, \alpha, \beta) = \int p(y|x^*, \mathbf{w}, \beta) p(\mathbf{w}|\mathcal{D}, \alpha, \beta) d\mathbf{w}$ **Convolution between 2 Gaussians => Gaussian**
 - ▶ Mean of predictive distribution $\mu^T \boldsymbol{\Phi}(x^*)$
 - ▶ Variance of predictive distribution $\sigma_{pred}^2(x^*) = \frac{1}{\beta} + \boldsymbol{\Phi}(x^*)^T \boldsymbol{\Sigma} \boldsymbol{\Phi}(x^*)$

$$p(y|x^*, \mathcal{D}, \alpha, \beta) = \mathcal{N}(y; \mu^T \boldsymbol{\Phi}(x^*), \frac{1}{\beta} + \boldsymbol{\Phi}(x^*)^T \boldsymbol{\Sigma} \boldsymbol{\Phi}(x^*))$$

Bayesian Linear Regression



Beyond Bayesian Linear Regression

No analytical expression for posterior $p(\mathbf{w}|\mathcal{D})$ and $p(y|\mathbf{x}^*, \mathcal{D})$ in general

\Rightarrow Approximation needed!

$p(\mathbf{w}|\mathcal{D})$: approximating with a Gaussian distribution $q_{\theta}(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\theta)$:

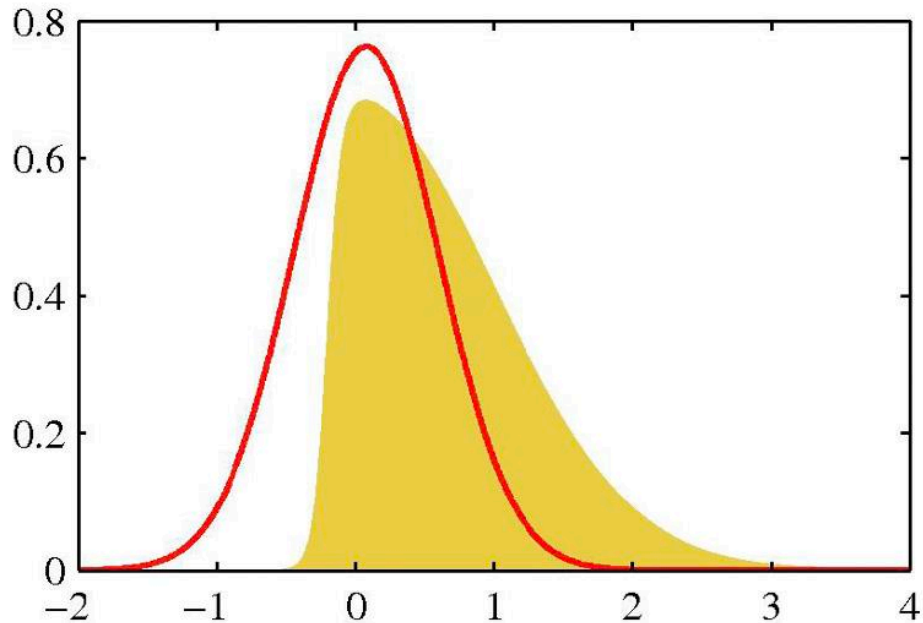
- θ : parameters of the Gaussian distribution, i.e. mean μ , covariance Σ

Approximating posterior with $q_{\theta}(\mathbf{w})$

- Laplace approximation
- Fit $q_{\theta}(\mathbf{w})$ on the mode of $p(\mathbf{w}|\mathbf{D})$

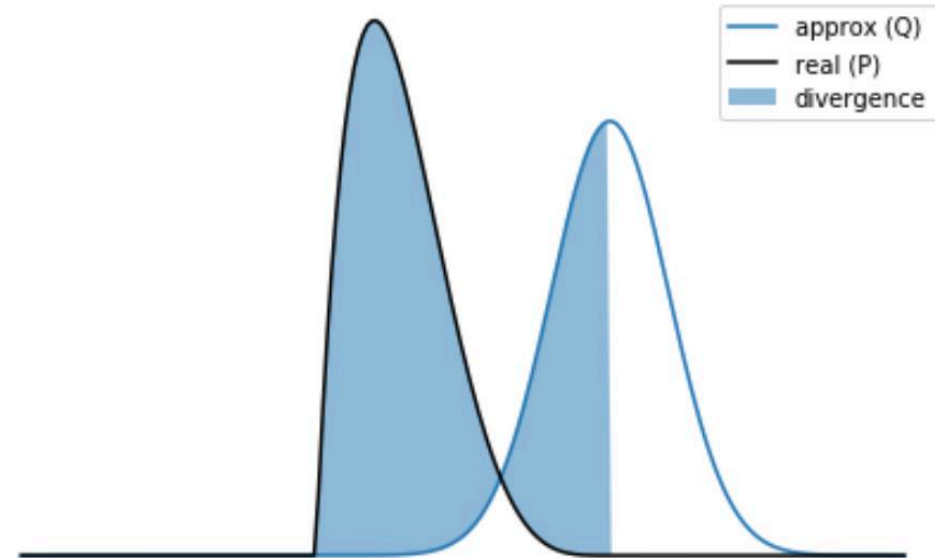
$$\boldsymbol{\mu} = \mathbf{w}_{MAP} \quad \nabla_{\mathbf{w}} p(\mathbf{w}) = 0$$

$$\boldsymbol{\Sigma}^{-1} = \nabla \nabla_{\mathbf{w}} p(\mathbf{w}|\mathbf{X}, \mathbf{Y}) \big|_{\mathbf{w}=\mathbf{w}_{MAP}}$$



- Used in Bayesian Logistic regression

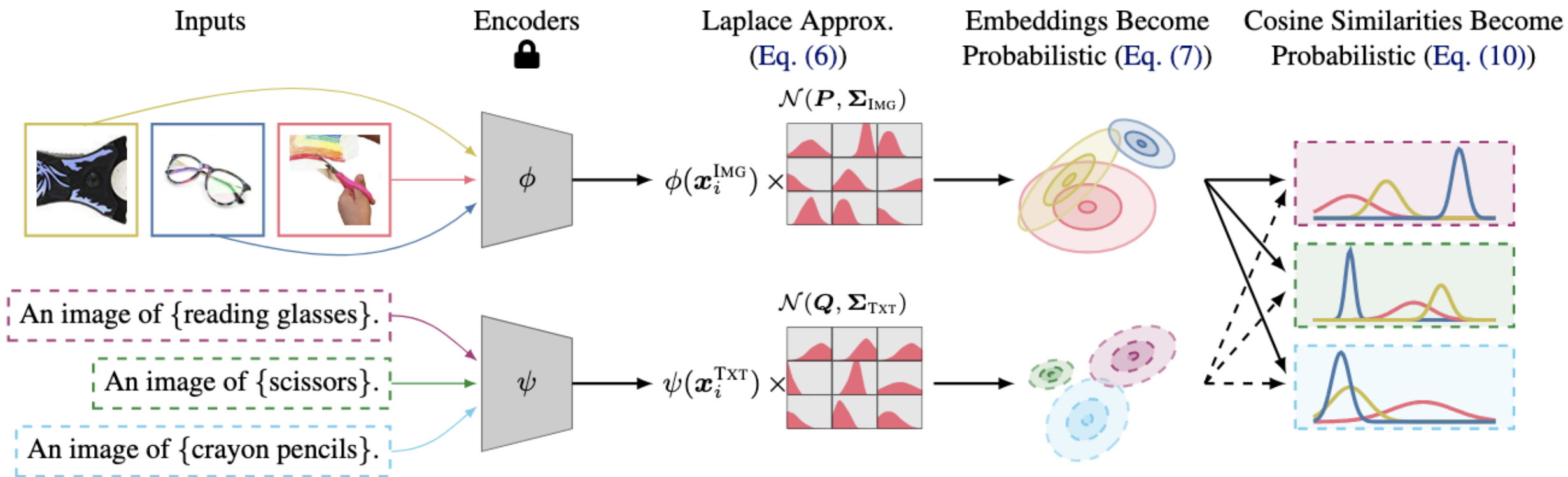
- Variational approximation:
Minimize $KL(q_{\theta}(\mathbf{w}) || p(\mathbf{w}|\mathbf{X}, \mathbf{Y}))$
- More global fitting



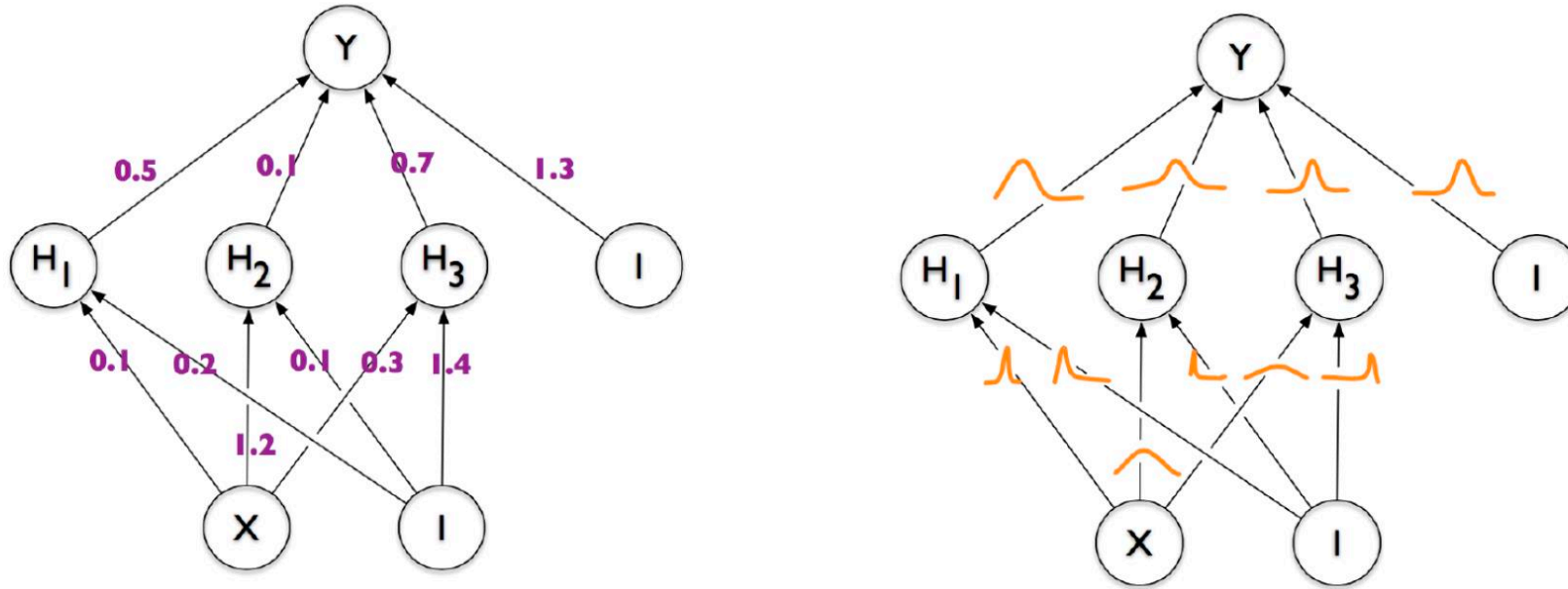
- Used in Bayesian Neural Networks

Laplace for Vision Language Models (VLMs)

- BayesVLM [A]: post-hoc Laplace approximation in VLMs, pre-trained e.g., LAION
- Estimate the distribution of similarities between text/image pairs
 - Estimate the distribution of the **last layer** or text/image encoders independently
 - Kronecker-factored (KFAC) Generalized Gauss–Newton (GGN) Hessian approx



VI in Bayesian Neural Networks (BNNs)



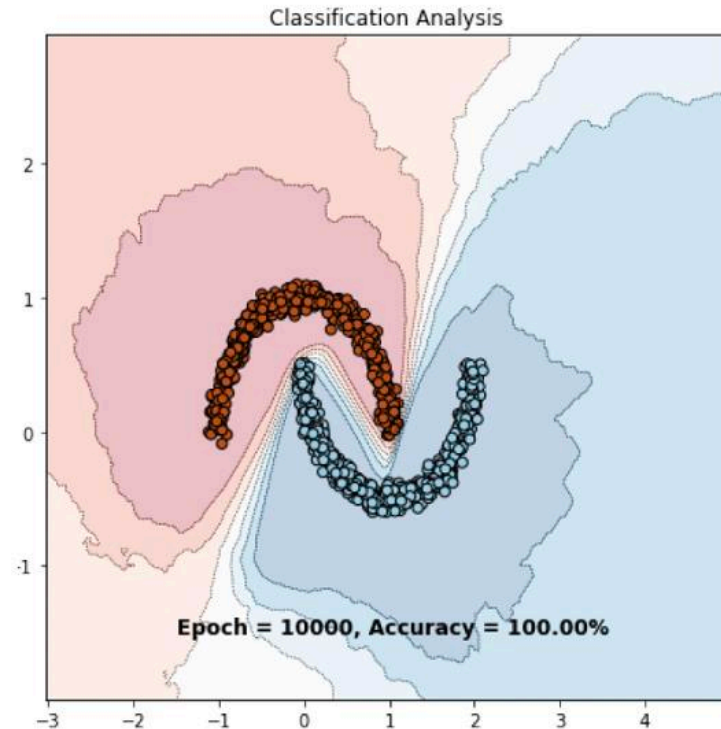
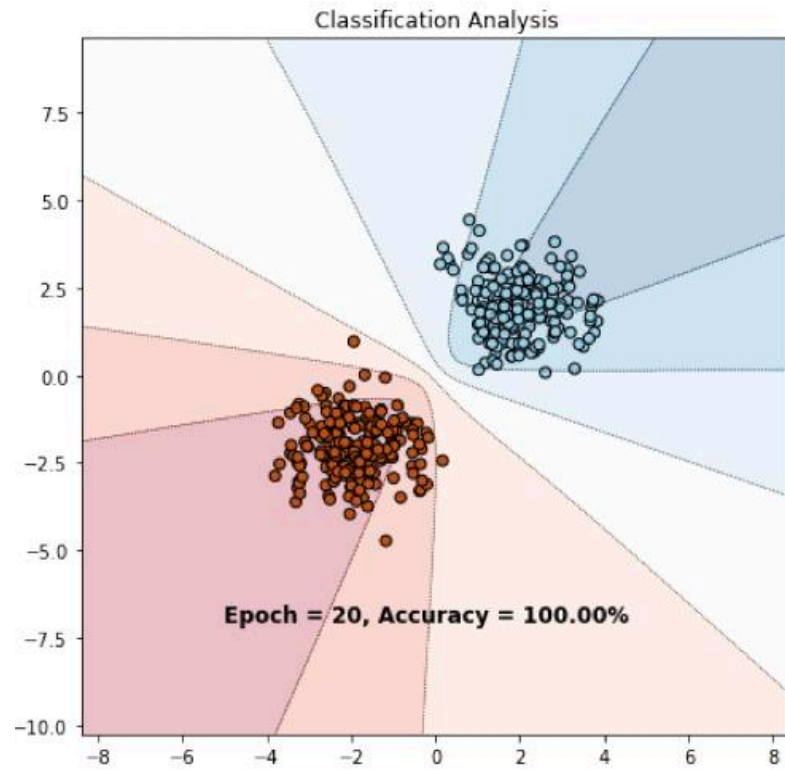
Credit: [Blundell et al., 2015]

- Generally, independence between weights (mean field)+ simple 1D Gaussians
 - ▶ Define prior over weights, e.g. $p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|0, \alpha^{-1}\mathcal{I})$
 - ▶ Define likelihood, $p(\mathbf{y}_i|\mathbf{x}_i, \mathbf{w})$, e.g. for regression $p(\mathbf{y}_i|\mathbf{x}_i, \mathbf{w}) = \mathcal{N}(\mathbf{y}_i; f^{\mathbf{w}}(\mathbf{x}_i), \beta^{-1})$
- Each weight of the network w_j has its own mean μ_j and variance σ_j
 - ▶ $\boldsymbol{\theta} = \{(\mu_j, \sigma_j)\}_{j \in \{1;D\}}$: **variational parameters**

Approximating Predictive Distribution

$p(\mathbf{y}|\mathbf{x}^*, \mathcal{D}) \approx \int p(\mathbf{y}|\mathbf{x}^*, \mathbf{w}) q_{\theta}(\mathbf{w}) d\mathbf{w} \Rightarrow$ approximate with Monte Carlo (MC) sampling

$$\int p(\mathbf{y}|\mathbf{x}^*, \mathbf{w}) q_{\theta}(\mathbf{w}) d\mathbf{w} \approx \sum_{s=1}^S p(\mathbf{y}|\mathbf{x}^*, \mathbf{w}^s) \quad \mathbf{w}^s \sim q_{\theta}(\mathbf{w}) \quad \text{Easy to sample from } q_{\theta}(\mathbf{w})$$

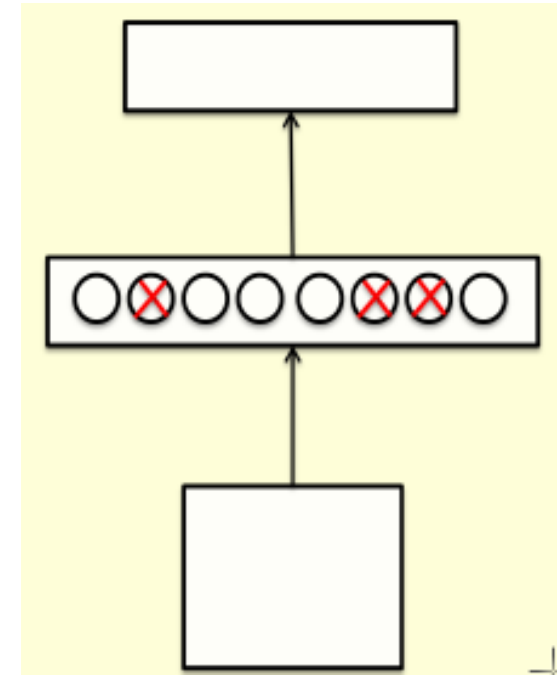


Monte Carlo Dropout

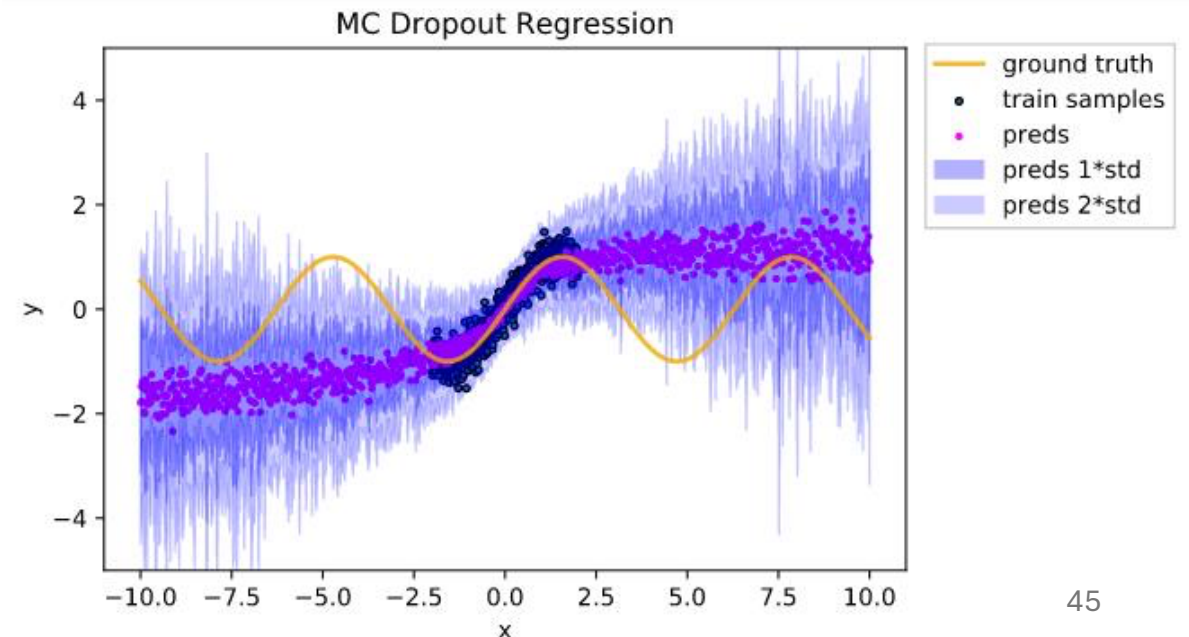
- Dropout as a variational inference [Gal, 2016]
- Big results: training with dropout is a special VI case
 - With some specific approximate posterior distribution

$$q(\mathbf{W}_I) = \text{diag}(\hat{\varepsilon}_I) \mathbf{M}_I, \varepsilon_{I,i} \sim \text{Bernoulli}(1 - p_i)$$

- \mathbf{M}_I deterministic weights, $q(\mathbf{W}_I)$ approximate posterior



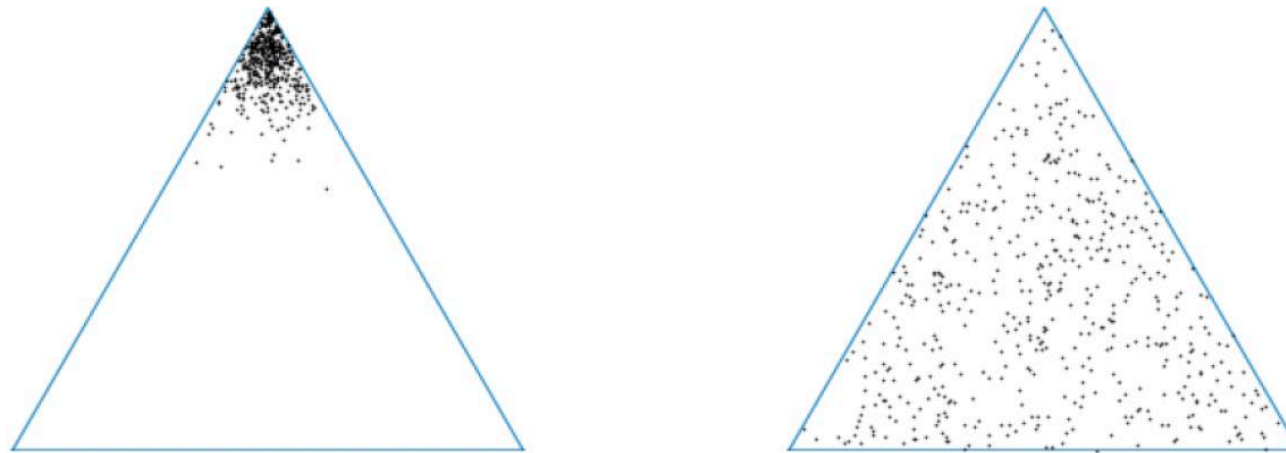
- => MC sampling of model trained with dropout: predictive distribution
 - Mean ~ prediction
 - Std ~ uncertainty



Ensembling

$$\int p(y|x^*, w) q_{\theta}(w) dw \approx \sum_{s=1}^S p(y|x^*, w^s) \quad w^s \sim q_{\theta}(w)$$

- MC sampling \Leftrightarrow ensembling different models
- Simple but effective baseline: train several models on the same dataset
 - Use disagreement, e.g. Mutual Information, as epistemic uncertainty



- **Main drawback:** High training & inference cost

Future directions in Uncertainty Quantification

- A unified uncertainty score (epistemic/aleatoric/calibration) able to capture different facets of uncertainties
- UQ for sequential decision tasks, forecasting or NLP?
 - Which uncertainty for LLMs / Generative VLMs?
- Dense/structured prediction tasks, correlations – segmentation
 - How to design UQ score reflecting these correlations?
- Statistical guarantees of UQ?
 - Conformal, but needs conditional guarantees

Thank you for your attention

- Questions?