# *Uncertainty Quantifiction and Anomaly Detection with Evidential Deep Learning*

## *Mark Neubauer*
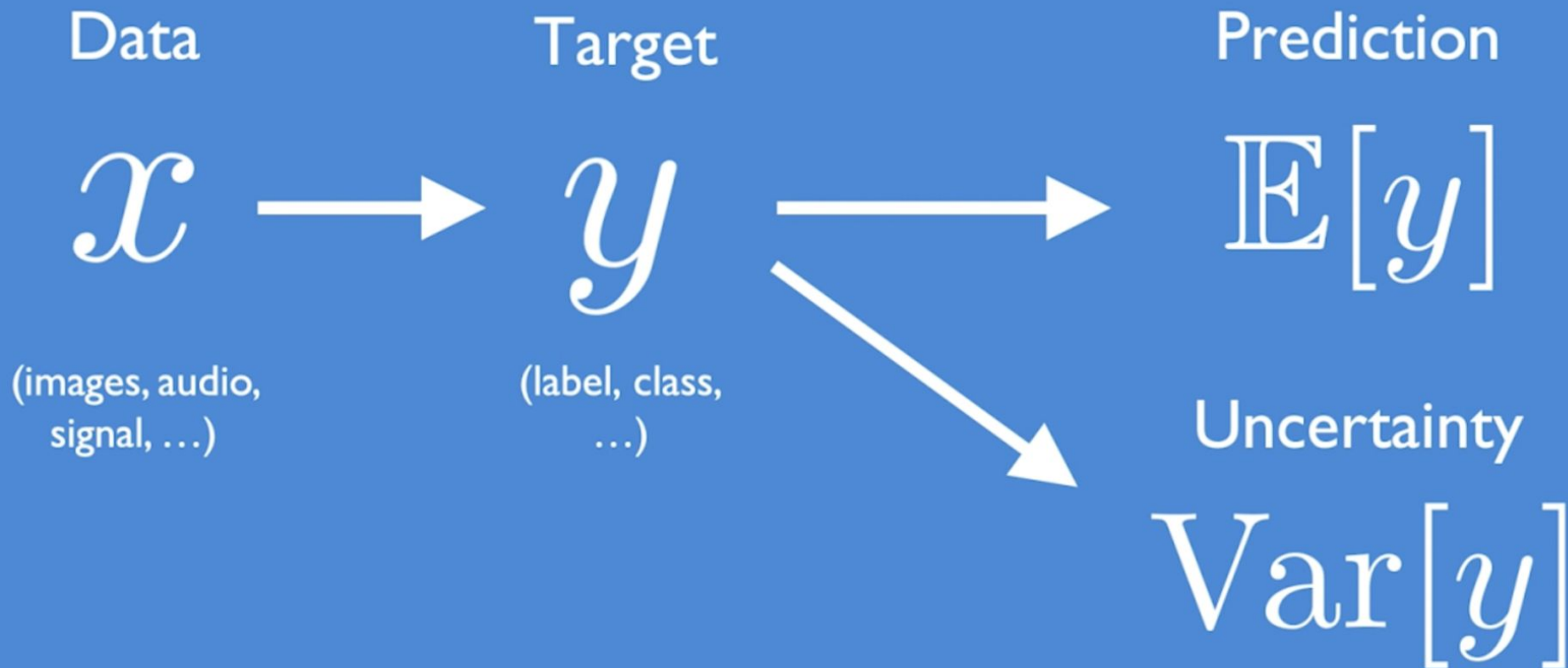
**University of Illinois at Urbana-Champaign**

*Artificial Intelligence and the Uncertainty Challenge in Fundamental Physics Workshop*

# Probabilistic Learning



Data

$x$

(images, audio,
signal, …)

Target

$y$

(label, class,
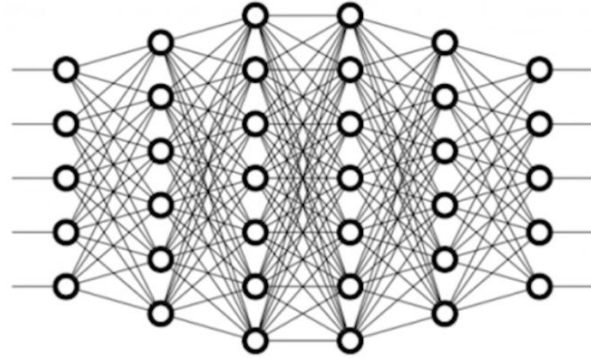…)

Prediction

$\mathbb{E}[y]$

Uncertainty

$\mathrm{Var}[y]$

# Learning Probabilistic Outputs



$$p(y = \text{``cat''} \mid x)$$

$$p(y = \text{``dog''} \mid x)$$

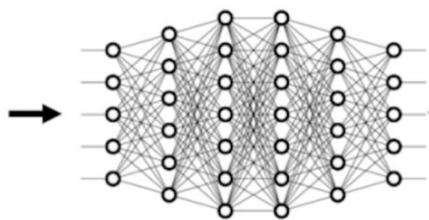Probability distribution over **discrete class categories**

# Learning Discrete Class Targets
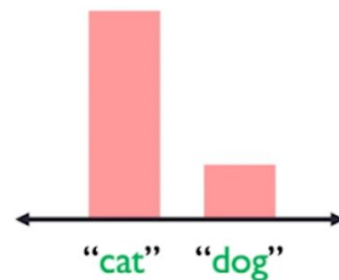
## Classification



$x$

$$p(y = \text{``cat''} \mid x)$$

$$p(y = \text{``dog''} \mid x)$$

**Activation:** softmax(z) $\longrightarrow$ $\sigma(\vec{z})_i = \dfrac{e^{z_i}}{\sum_{j=1}^{K} e^{z_j}}$

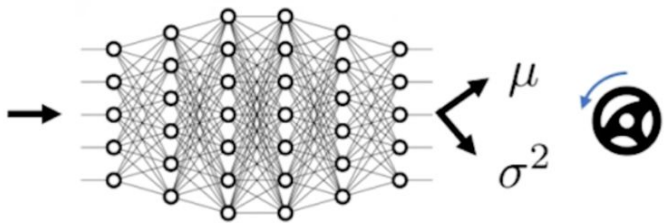**Loss:** Neg. Log Likelihood (Cross Entropy) $\longrightarrow$ $-\sum_{i=1}^{K} y_i \log p_i$

## Why?

$$\underline{y} \sim \underline{\text{Categorical}(p)}$$

Class Labels     Likelihood function     Distribution parameters (probabilities)

$$f(y = y_i \mid p) = p_i$$



"cat"    "dog"

# Learning Continuous Class Targets



Regression

**Activation:**
$$\mu \in \mathbb{R}$$
$$\sigma > 0$$
$$\longrightarrow$$
$$\mu = z_\mu$$
$$\sigma = \exp(z_\sigma)$$

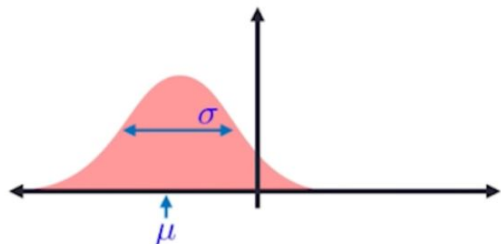**Loss:** Neg. Log Likelihood $\longrightarrow -\log\left(\mathcal{N}(y|\mu, \sigma^2)\right)$

Why?

$$y \sim \text{Normal}(\mu, \sigma^2)$$

Target Labels  Likelihood function  Distribution parameters

$$f(y \mid \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y-\mu)^2}{2\sigma^2}\right)$$
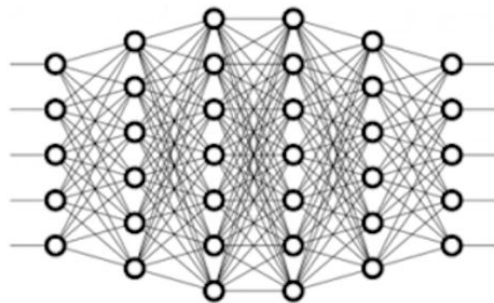
# Likelihood vs Confidence

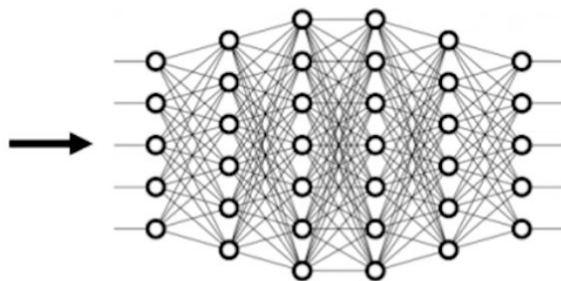⚠️ **Do not mistake likelihood (probability) for model confidence!** ⚠️



$$p(\text{``cat''})$$

$$p(\text{``dog''})$$

# Likelihood vs Confidence

⚠️ **Do not mistake likelihood (probability) for model confidence!** ⚠️



$$p(\text{"cat"}) = 0.5$$

$$p(\text{"dog"}) = 0.5$$

# Likelihood vs Confidence

⚠️ **Do not mistake likelihood (probability) for model confidence!** ⚠️



Expectation:
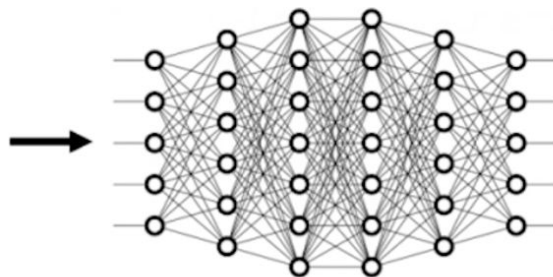Training on a your dataset

Reality:
Testing in reality

Dogs

Driving

*http://introtodeeplearning.com/2021/slides/6S191_MIT_DeepLearning_L7.pdf*

# Likelihood vs Confidence

⚠️ **Do not mistake likelihood (probability) for model confidence!** ⚠️

*The output likelihoods will be unreliable if the input is* **unlike anything during training**



$$p(\text{"cat"}) = 0.15$$

$$p(\text{"dog"}) = 0.85$$

⭐ $p(\text{"cat"}) + p(\text{"dog"}) = 1$ ⭐
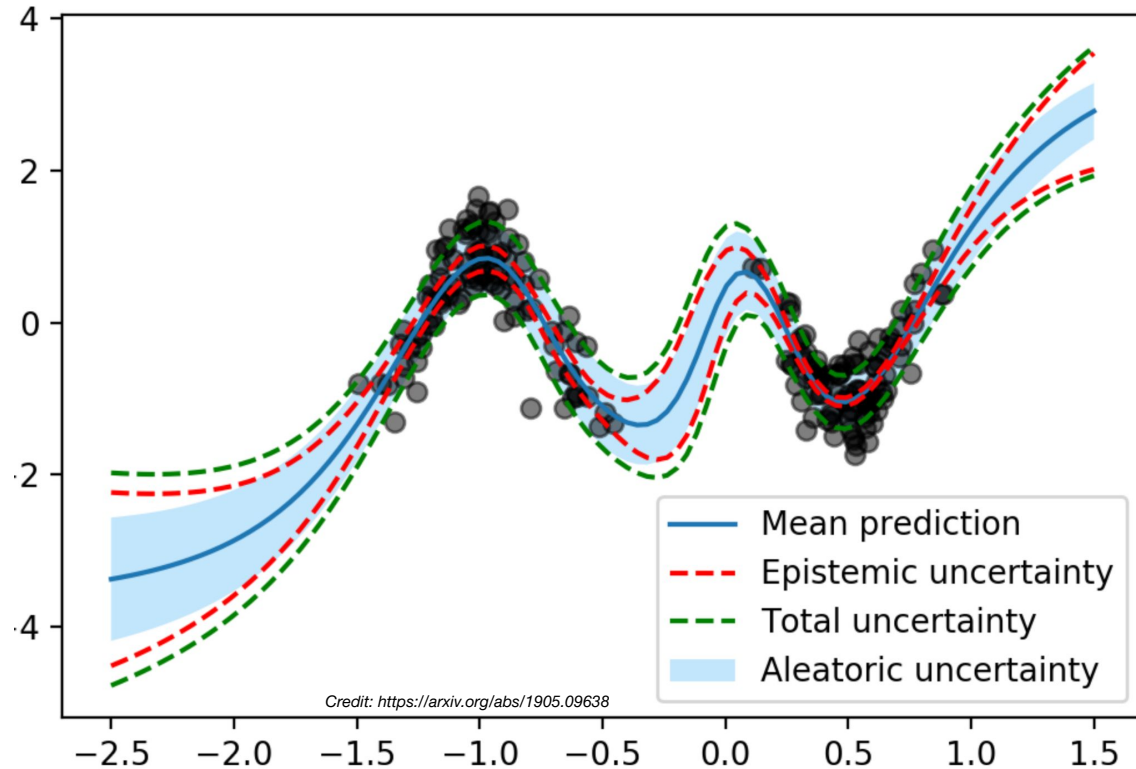
# Types of Uncertainty

## **Aleatoric** Uncertainty

- Describes the confidence in the input data
- Large when input data is noisy
- Cannot be reduced by simply adding more data

## **Epistemic** Uncertainty

- Describes the confidence in the prediction
- Large when insufficient training data
- Can be reduced by adding more data



Credit: https://arxiv.org/abs/1905.09638

Legend:
— Mean prediction
-- Epistemic uncertainty
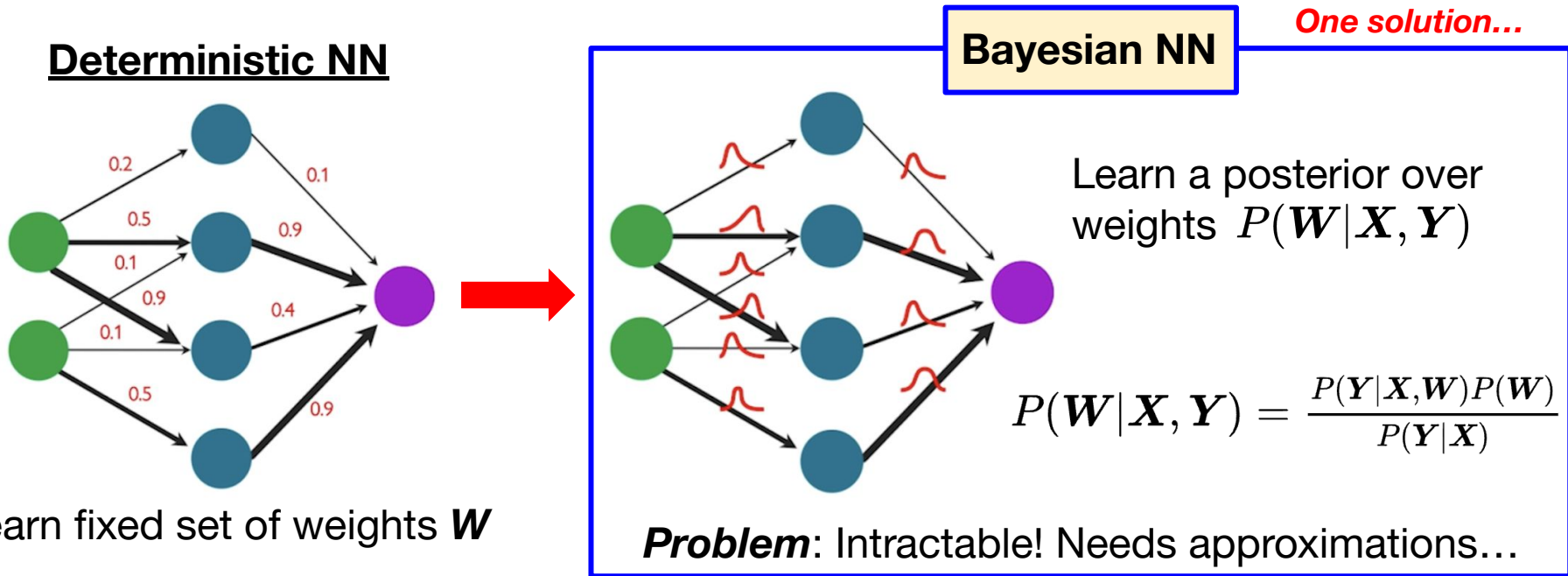-- Total uncertainty
▨ Aleatoric uncertainty

# Estimating epistemic uncertainty

- Aleatoric uncertainty can be learned directly using NNs
- Epistemic uncertainty is much more challenging to estimate

**Q**: How can a model understand when it doesn't know the answer?



**Deterministic NN**

Learn fixed set of weights **W**

*One solution…*

**Bayesian NN**

Learn a posterior over weights $P(\boldsymbol{W}|\boldsymbol{X}, \boldsymbol{Y})$

$$P(\boldsymbol{W}|\boldsymbol{X}, \boldsymbol{Y}) = \frac{P(\boldsymbol{Y}|\boldsymbol{X}, \boldsymbol{W})P(\boldsymbol{W})}{P(\boldsymbol{Y}|\boldsymbol{X})}$$

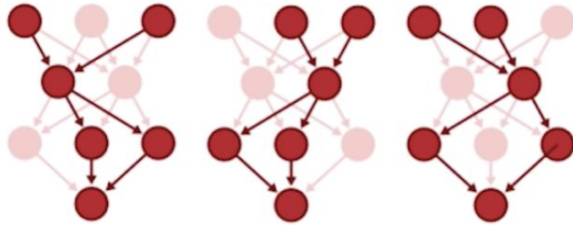*Problem*: Intractable! Needs approximations…

# Approximations through Sampling

Evaluate $T$ stochastic forward passes using different samples of weights $\{W_t\}_{t=1}^{T}$
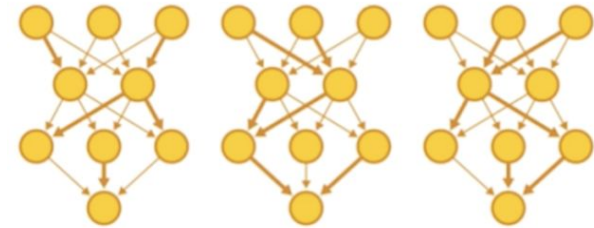
- Dropout as a form of stochastic sampling

$$z_{w,t} \sim Bernoulli(p) \quad \forall w \in W$$

*Monte Carlo Dropout*



- Ensemble of $T$ independently trained models, each learning a unique

$$W_t = train(f; X, Y)$$

*Model Ensembles*



***Epistemic uncertainty:***

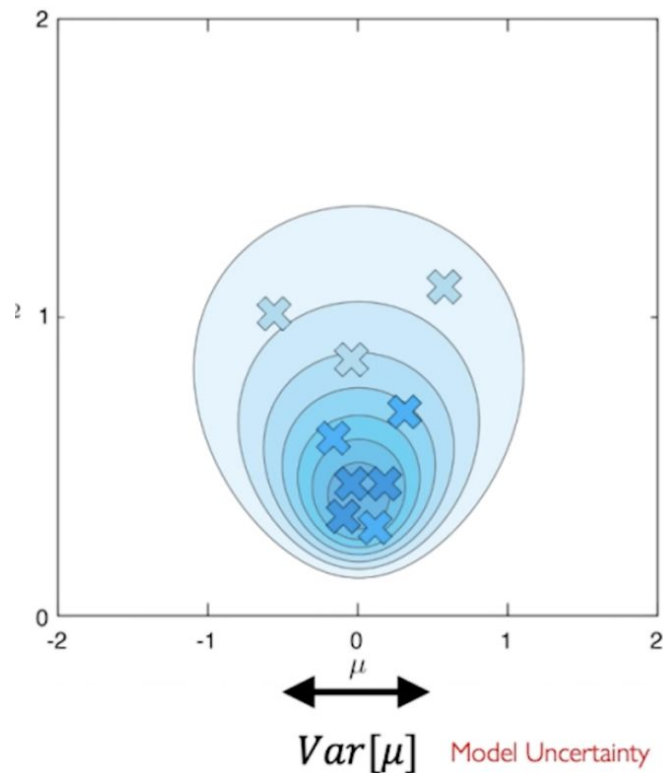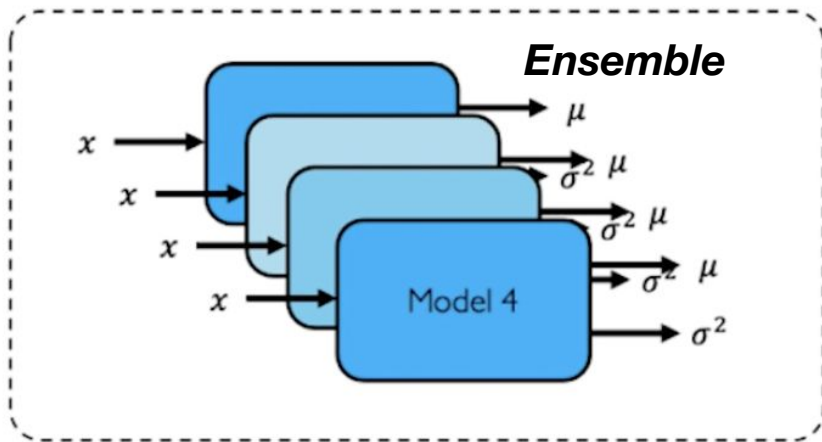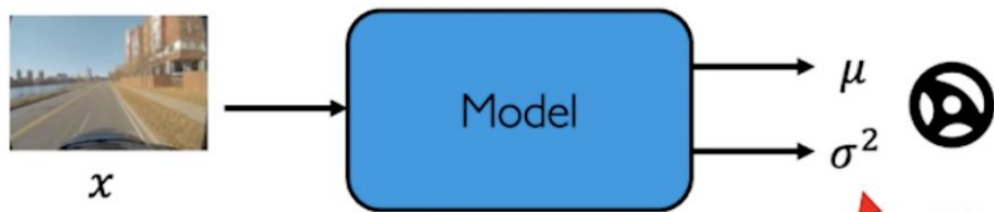$$Var(\hat{Y}|X) = \frac{1}{T}\sum_{t=1}^{T} f(X)^2 - \mathbb{E}(\hat{Y}|X)^2$$

where $\mathbb{E}(\hat{Y}|X) = \frac{1}{T}\sum_{t=1}^{T} f(X|W_t)$

## Downsides of Bayesian Deep Learning

- ***Slow***: Requires running network $T$ times for each input
- ***Memory***: Stores $T$ copies of the network in parallel
- ***Efficiency***: Sampling hinders real-time on edge devices
- ***Calibration***: Sensitive to prior and often over-confident

# Uncertainty Estimation: Sampling



**Ensemble**
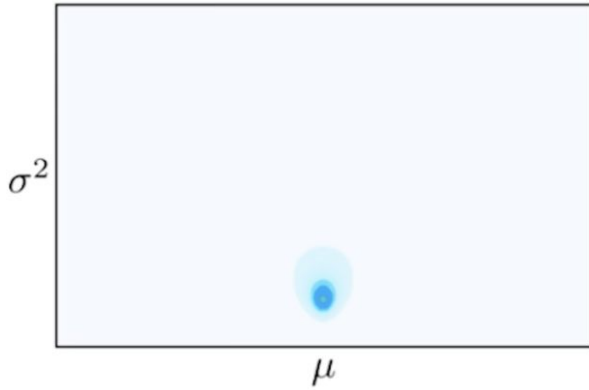
$Var[\mu]$    Model Uncertainty

**Q**: Can we *underline{directly}* learn the parameters defining this likelihood distribution?
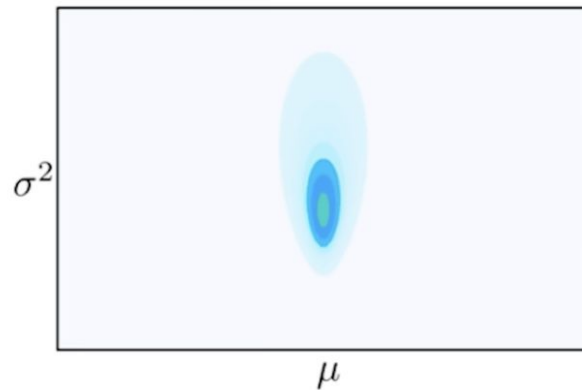
# Evidential Deep Learning (EDL)

Treat learning and an evidence acquisition process
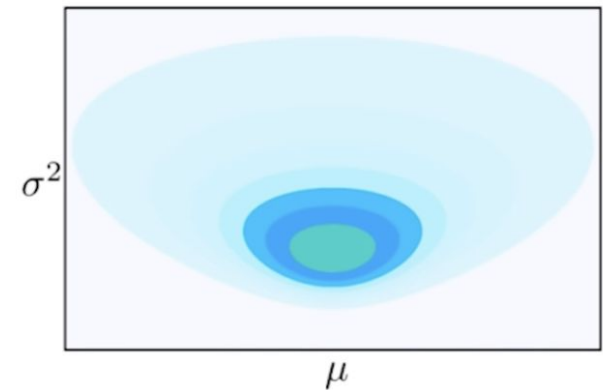- More evidence → Increased predictive confidence



Low uncertainties
→ High confidence

High aleatoric (data)
uncertainty

High epistemic (model)
uncertainty

**Goal**: *train a neural network to learn these type of evidential distributions*

# EDL for Regression

*Key point to remember*: **Sampling from an evidential distribution yields individual new distributions over the data**



$y \sim \text{Normal}(\mu, \sigma^2)$

Target Labels    Likelihood function    Distribution parameters

Assume the distribution parameters are not known, place priors over each and probabilistically estimate!

$\mu \sim \text{Normal}(\gamma, \sigma^2 v^{-1})$
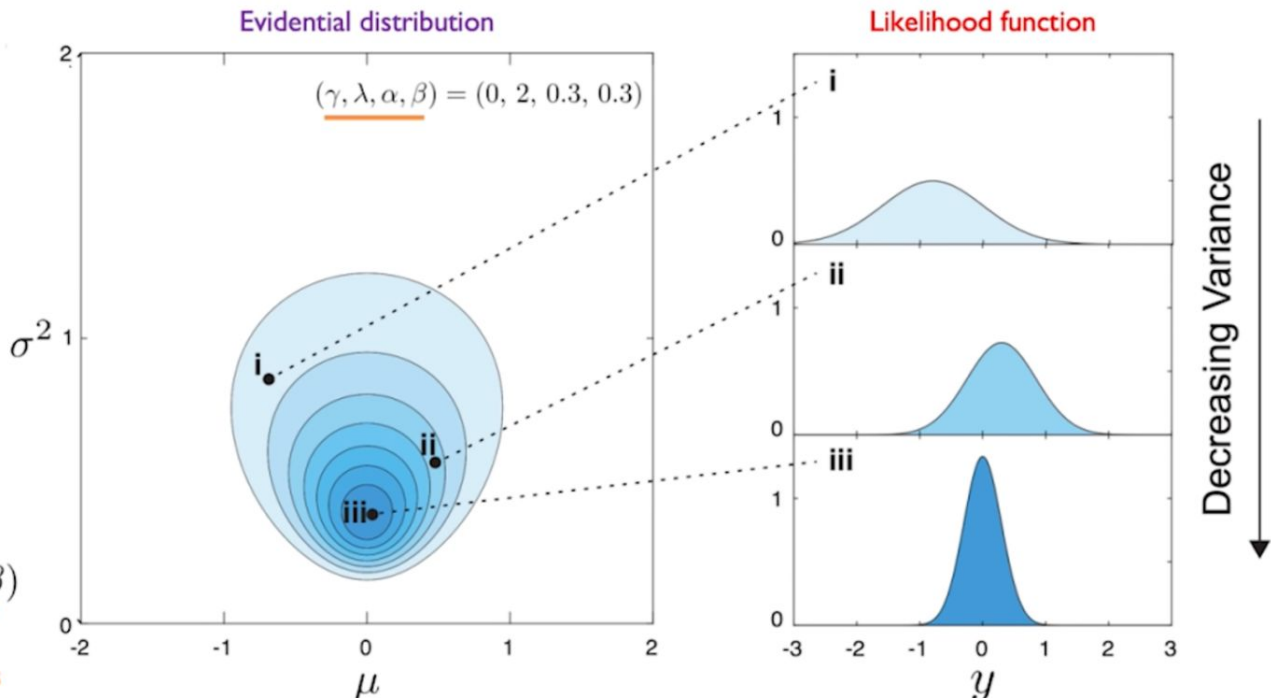
$\sigma^2 \sim \Gamma^{-1}(\alpha, \beta)$

$\mu, \sigma^2 \sim \text{NormalInvGamma}(\gamma, v, \alpha, \beta)$

Distribution parameters    Evidential Prior    Model parameters

**Evidential distribution**

$(\gamma, \lambda, \alpha, \beta) = (0, 2, 0.3, 0.3)$

**Likelihood function**

Decreasing Variance

# EDL for Classification

Sampling from an evidential distribution yields individual new distributions over the data

$$y \in \{1, \cdots, K\}$$

$$\underline{y} \sim \underline{\text{Categorical}}(\underline{\boldsymbol{p}})$$

Class Labels — Likelihood function — Distribution parameters (probabilities)

$$\underline{\boldsymbol{p}} \sim \underline{\text{Dirichlet}}(\underline{\boldsymbol{\alpha}})$$

Distribution parameters — Evidential Prior — Model parameters

$$\boldsymbol{p} = \begin{bmatrix} 1.0 \\ 0.0 \\ 0.0 \end{bmatrix}$$

$$\boldsymbol{p} = \begin{bmatrix} 0.5 \\ 0.5 \\ 0.0 \end{bmatrix}$$

$$\boldsymbol{p} = \begin{bmatrix} 0.33 \\ 0.33 \\ 0.33 \end{bmatrix}$$

$$\boldsymbol{p} = \begin{bmatrix} 0.1 \\ 0.7 \\ 0.2 \end{bmatrix}$$

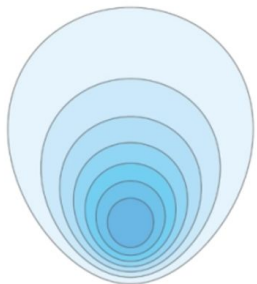$$K = 3; \quad \boldsymbol{\alpha} = (5, 5, 5)$$

# EDL Distributions



## Regression

$y \in \mathbb{R}$

$y \sim \mathrm{Normal}(\mu, \sigma^2)$

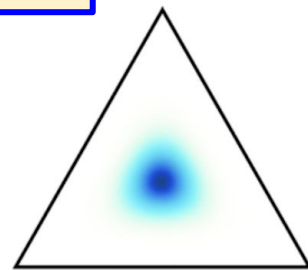$\mu, \sigma^2 \sim \mathrm{NormalInvGamma}(\gamma, \upsilon, \alpha, \beta)$

## Classification

$y \in \{1, \cdots, K\}$

$y \sim \mathrm{Categorical}(\boldsymbol{p})$

$\boldsymbol{p} \sim \mathrm{Dirichlet}(\boldsymbol{\alpha})$

Note that the choice of evidential distributions is closely related to conjugate priors in the context of Bayesian inference

- It is often easiest for computations to pick your evidential distribution to be a conjugate prior of your likelihood:
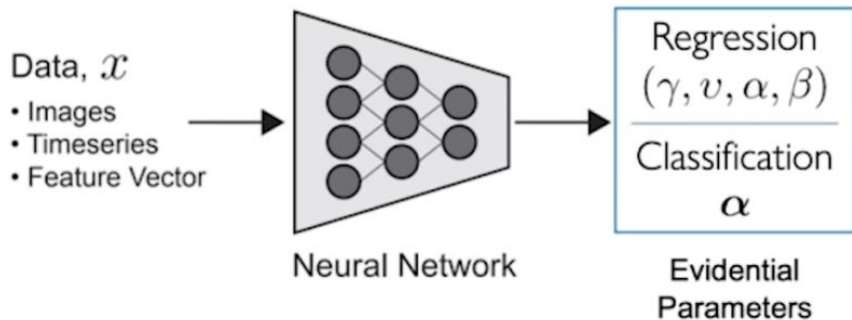
$$p(\theta | y) = \frac{p(y|\theta)\, p(\theta)}{\int_{\theta'} p(y|\theta')\, p(\theta')\, d\theta'}$$
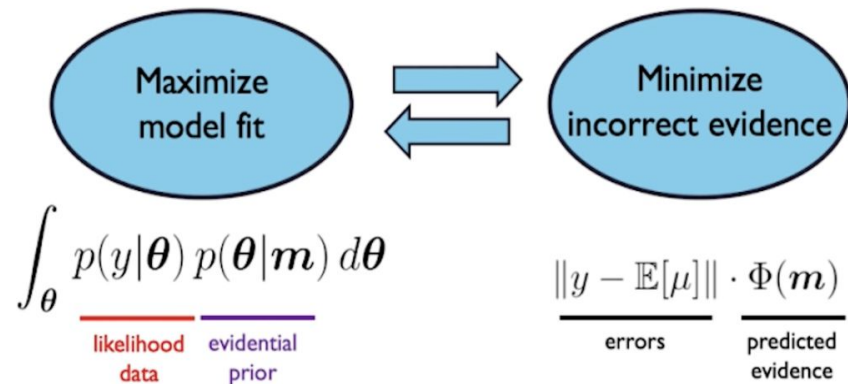
# EDL Model and Training

### Model

*Train the network to output the parameters of an evidential distribution:*



Data, $x$
• Images
• Timeseries
• Feature Vector

Neural Network

Regression
$(\gamma, \upsilon, \alpha, \beta)$

Classification
$\alpha$

Evidential Parameters

### Optmization

*Perform multi-objective training:*



Maximize model fit

Minimize incorrect evidence

$\int_\theta p(y|\boldsymbol{\theta}) \, p(\boldsymbol{\theta}|\boldsymbol{m}) \, d\boldsymbol{\theta}$

likelihood data    evidential prior

$\|y - \mathbb{E}[\mu]\| \cdot \Phi(\boldsymbol{m})$
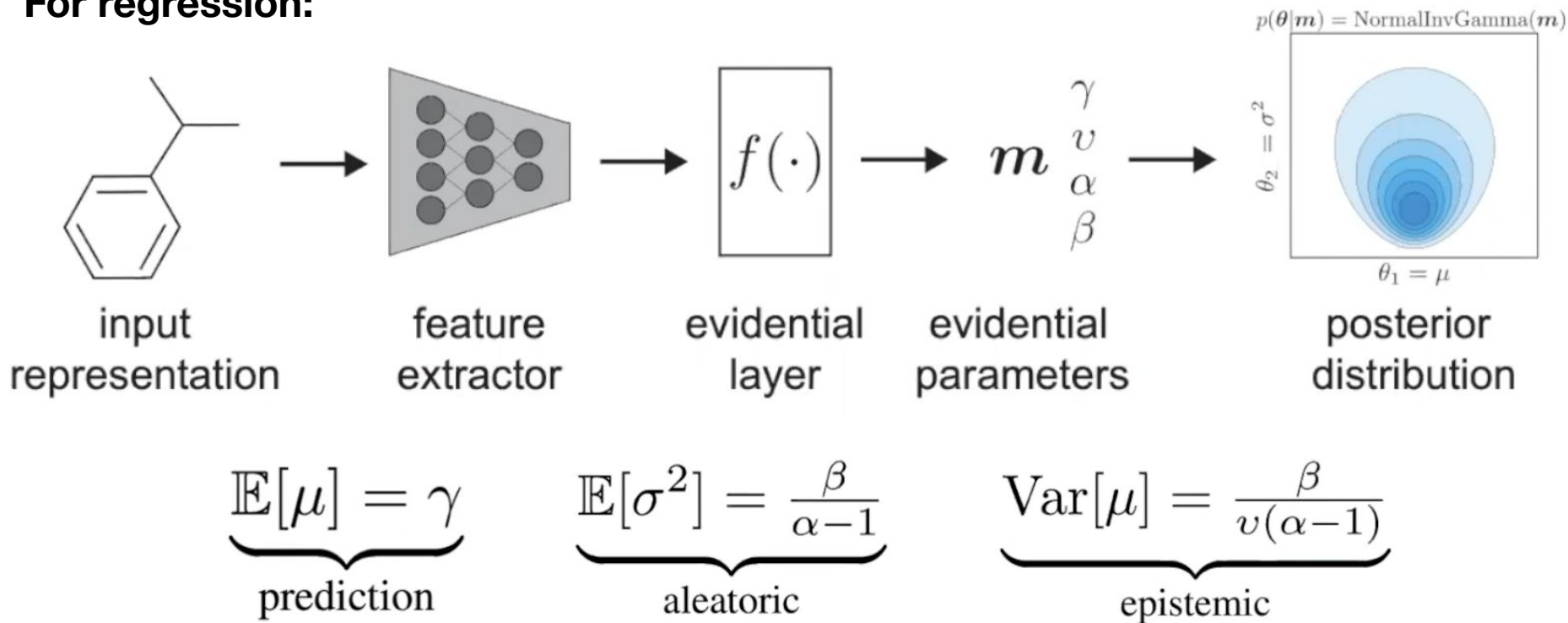
errors    predicted evidence

**For classification:** $\mathcal{L}(\Theta) = \sum_{i=1}^{N} \mathcal{L}_i(\Theta) + \lambda_t \sum_{i=1}^{N} KL[D(\mathbf{p}_i | \tilde{\boldsymbol{\alpha}}_i) \, || \, D(\mathbf{p}_i | \mathbf{1})]$

Reconstruction Loss

Penalty for assigning large confidence to uncertain samples

# Forming EDL Predictions

**For regression:**



$$\underbrace{\mathbb{E}[\mu] = \gamma}_{\text{prediction}} \qquad \underbrace{\mathbb{E}[\sigma^2] = \frac{\beta}{\alpha-1}}_{\text{aleatoric}} \qquad \underbrace{\mathrm{Var}[\mu] = \frac{\beta}{v(\alpha-1)}}_{\text{epistemic}}$$

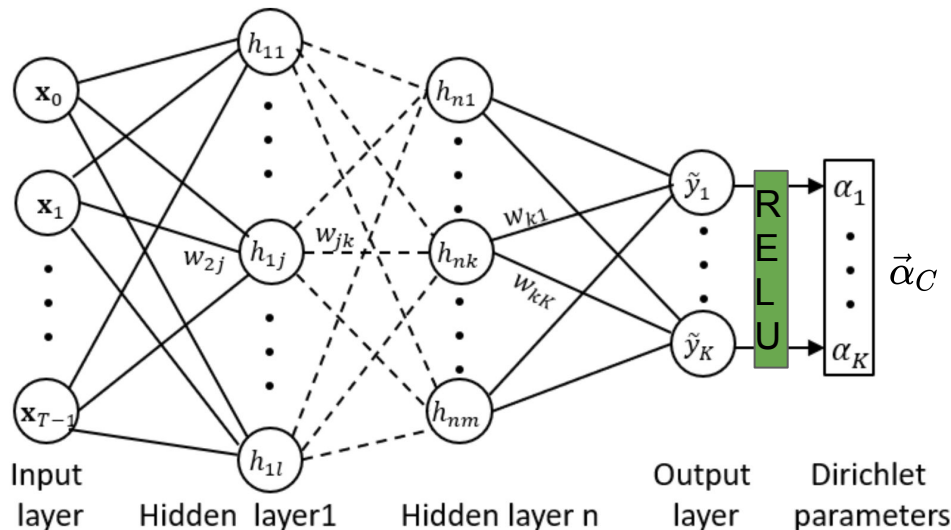$\rightarrow$ *Evidential Uncertainty can be easily integrated with 4x parameters and a new loss*

# Forming EDL Predictions

**For classification:**

Data, $x$

- Images
- Timeseries
- Feature Vector



Once the network learns the parameters **α**, its mean, can be taken as an estimate of the $K$ class probabilities
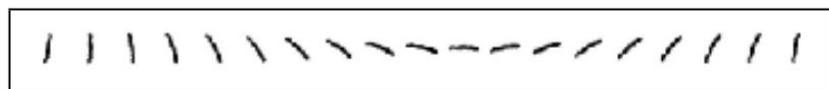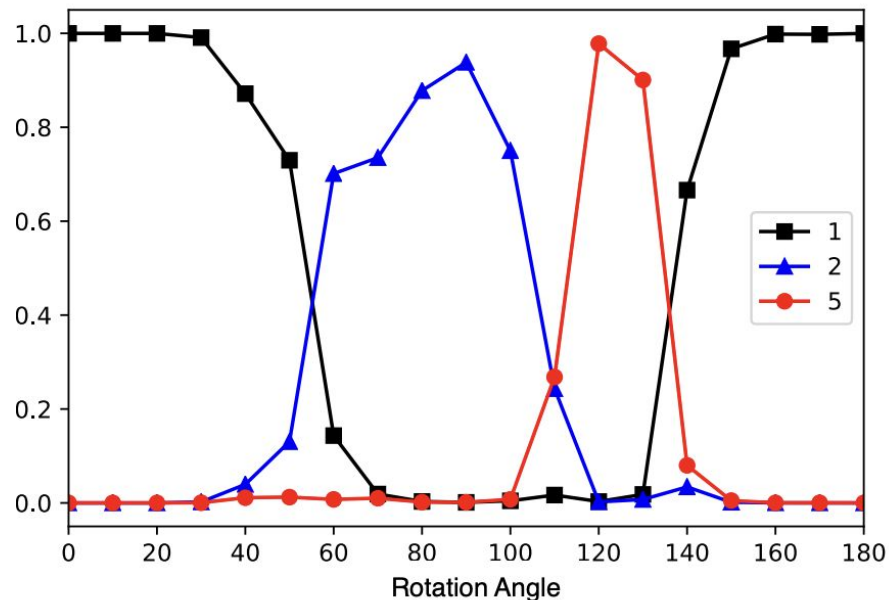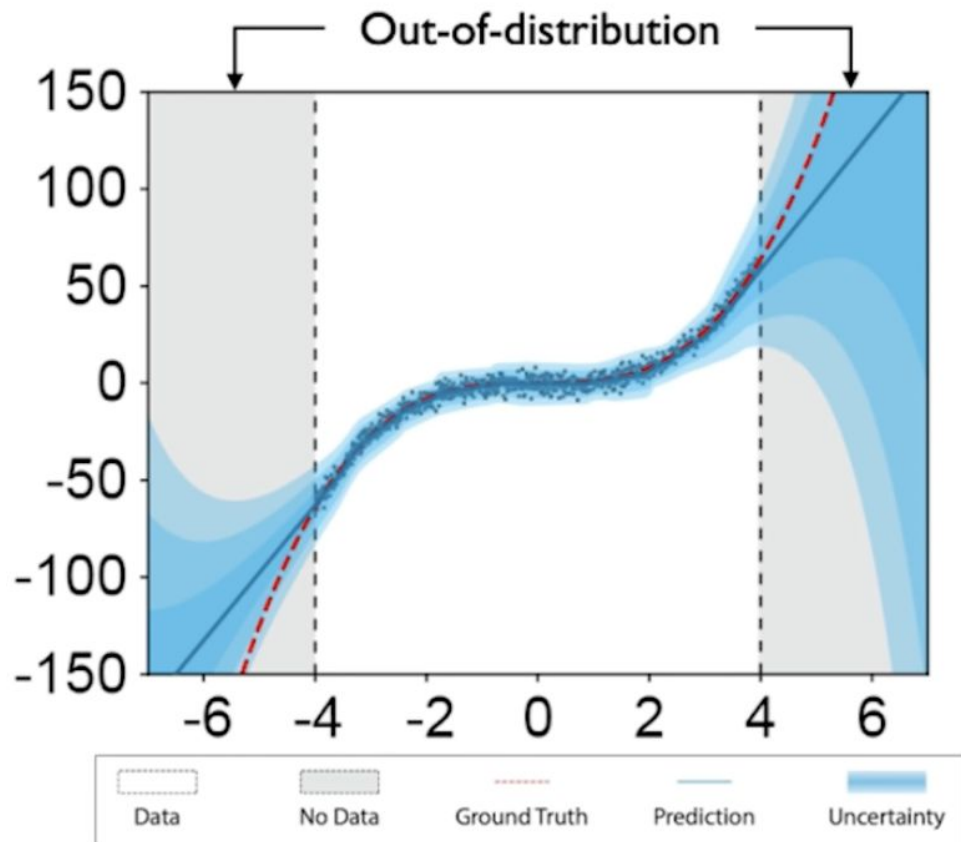
$$\tilde{p}_c = \alpha_c / S$$

The epistemic uncertainty $u$ on the prediction is computed as the inverse of total evidence or Dirichlet strength $S$
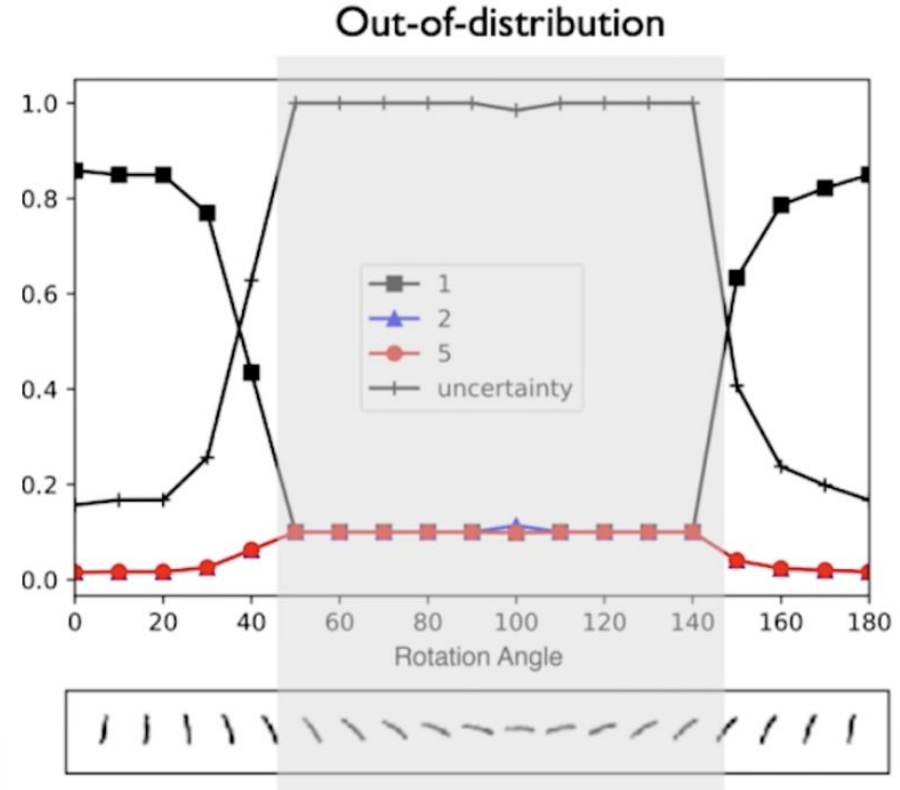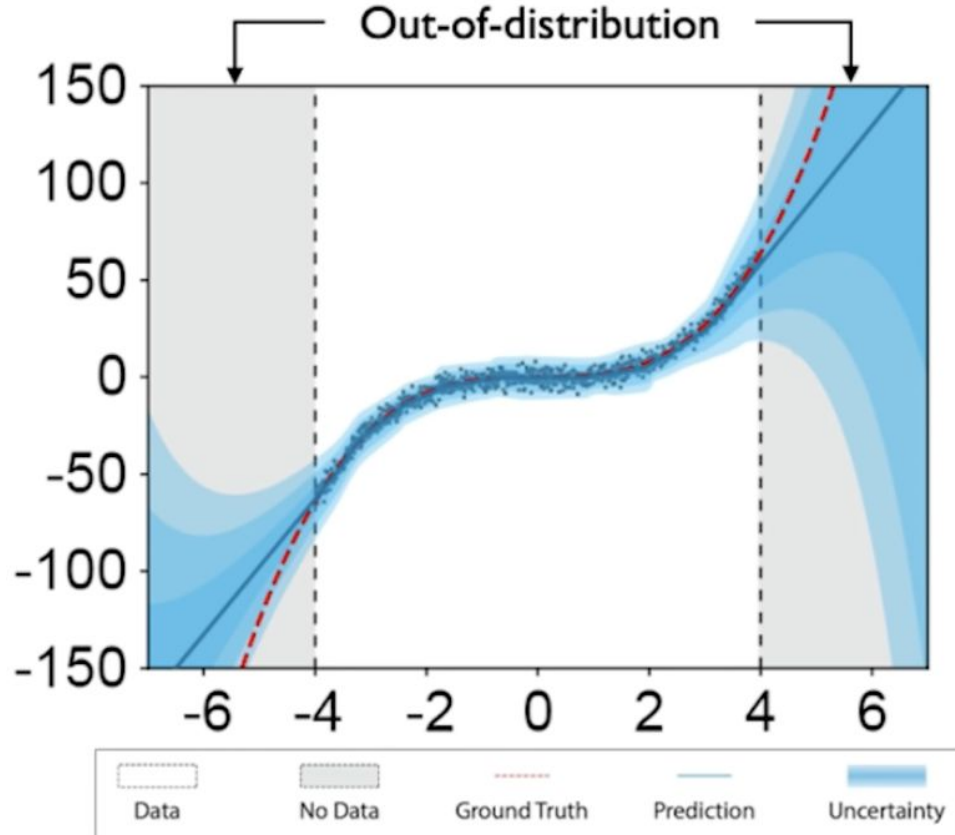
$$u = K/S \qquad \text{where} \qquad S = \sum_{c=1}^{K} \alpha_c$$

→ *Evidential Uncertainty can be easily integrated with × K parameters and a new loss*

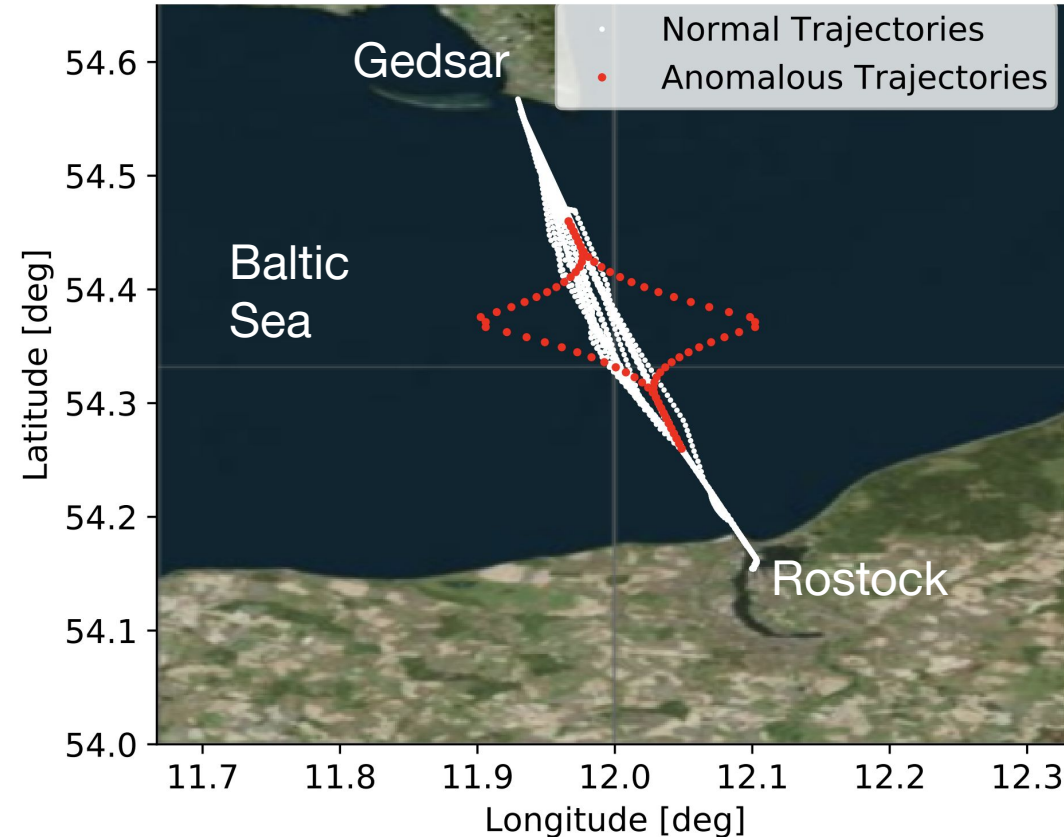# EDL Toy Learning Problems



Out-of-distribution

| Data | No Data | Ground Truth | Prediction | Uncertainty |

Rotation Angle

- 1
- 2
- 5

# EDL Toy Learning Problems



Data    No Data    Ground Truth    Prediction    Uncertainty

# EDL Applied to Anomaly Detection



## *Maritime Anomaly Detection*

Most ships are equipped with automatic identification system (AIS) transponders to provide their static and dynamic information

Vessels' location, navigational status, and voyage-related information can be used for

- *collision-avoidance mechanisms*
- *vessel tracking*
- detection of *loss of AIS signal* and *anomalous trajectories*

# Maritime Anomaly Detection



60-90 passenger vessels

$\theta_{AT} = 0.7$



60-90 passenger vessels

$\theta_{AT} = 0.4$

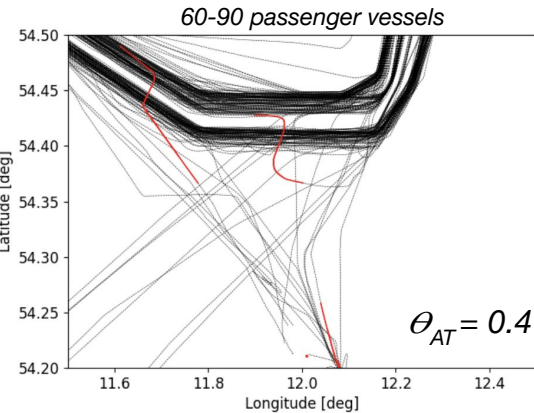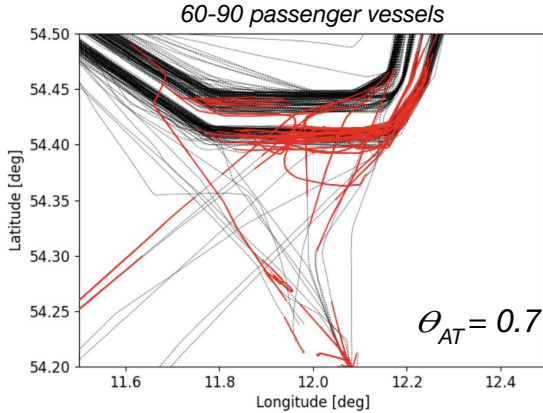## *EDL for Anomalous Trajectory Detection*

High epistemic uncertainty may represent anomalous trajectory. However, different output features are predicted with different uncertainties, so comparing segments with a set uncertainty threshold might not be a good idea

Thus, a trajectory segment is defined as anomalous if the predicted sequences of the segment have an abrupt transition in their epistemic uncertainties

$$\min_{d}\left[\frac{\min_{j}(\mathrm{var}[\mu_{j}^{d}])}{\max_{j}(\mathrm{var}[\mu_{j}^{d}])}\right] < \Theta_{AT}$$

This selects the feature *d* and output sequence *j* with the minimum normalized epistemic uncertainties. If this value is below $\theta_{AT}$, then the segment is considered as anomalous

A vessel's trajectory is termed as anomalous if it contains one or more anomalous segments
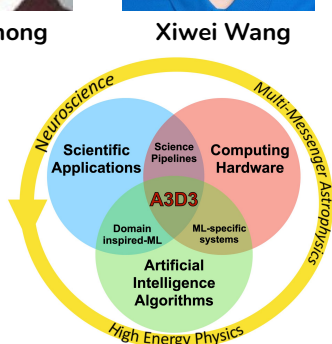
**Ayush Khot**

**Avik Roy**

**Mark Neubauer**

**Dewen Zhong**

**Xiwei Wang**

# A Detailed Study of Interpretability of Deep Neural Network based Top Taggers

**Ayush Khot, Mark S. Neubauer, and Avik Roy**[1]

*Department of Physics & National Center for Supercomputing Applications (NCSA) University of Illinois at Urbana-Champaign*

E-mail: akhot2@illinois.edu, msn@illinois.edu, avroy@illinois.edu

ABSTRACT: Recent developments in the methods of explainable AI (XAI) allow researchers to explore the inner workings of deep neural networks (DNNs), revealing crucial information about input-output relationships and realizing how data connects with machine learning models. In this paper we explore interpretability of DNN models designed to identify jets coming from top quark decay in high energy proton-proton collisions at the Large Hadron Collider (LHC). We review a subset of existing top tagger models and explore different quantitative methods to identify which features play the most important roles in identifying the top jets. We also investigate how and why feature importance varies across different XAI metrics, how correlations among features impact their explainability, and how latent space representations encode information as well as correlate with physically meaningful quantities. Our studies uncover some major pitfalls of existing XAI methods and illustrate how they can be overcome to obtain consistent and meaningful interpretation of these models. We additionally illustrate the activity of hidden layers as Neural Activation Pattern (NAP) diagrams and demonstrate how they can be used to understand how DNNs relay information across the layers and how this understanding can help to make such models significantly simpler by allowing effective model reoptimization and hyperparameter tuning. These studies not only facilitate a methodological approach to interpreting models but also unveil new insights about what these models learn. Incorporating these observations into augmented model design, we propose the Particle Flow Interaction Network (PFIN) model and demonstrate how interpretability-inspired model augmentation can improve top tagging performance.
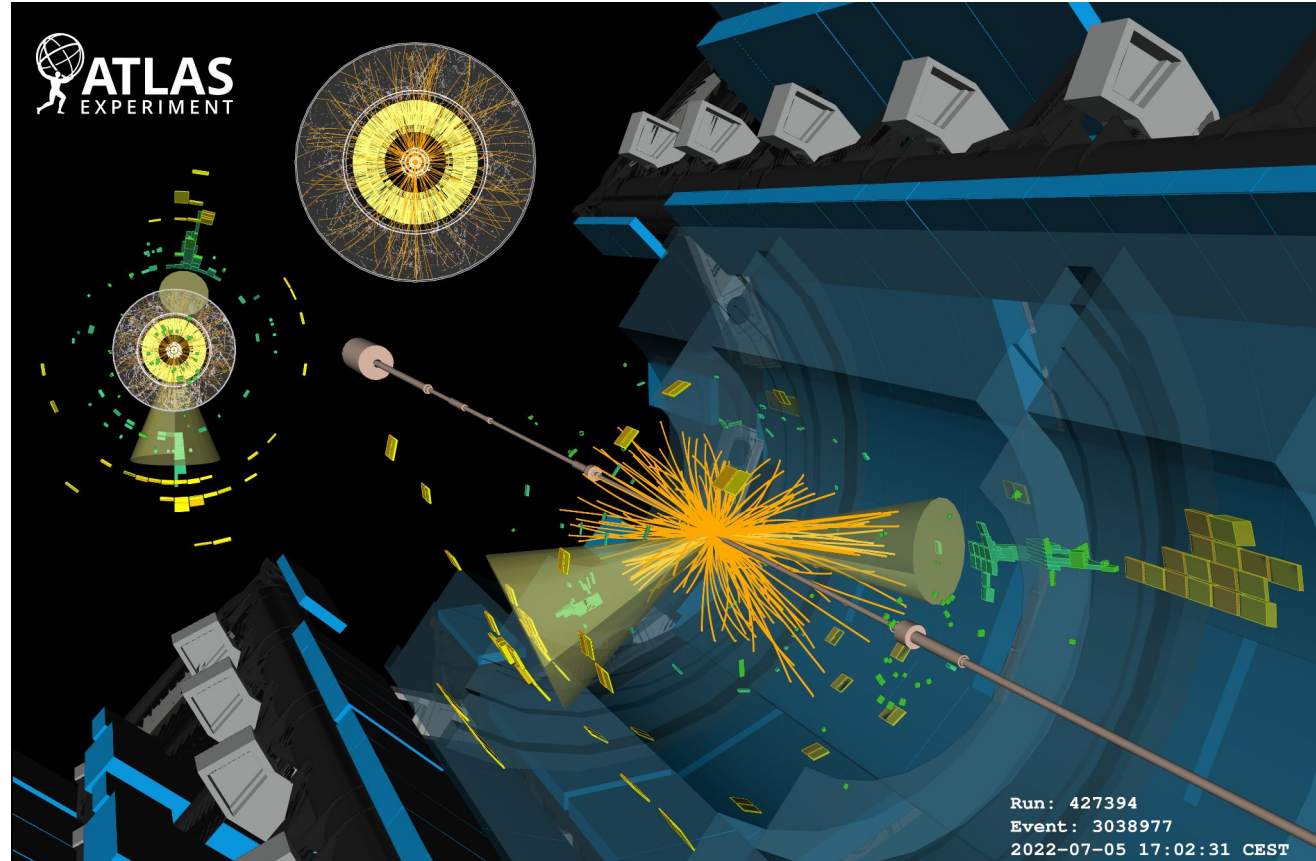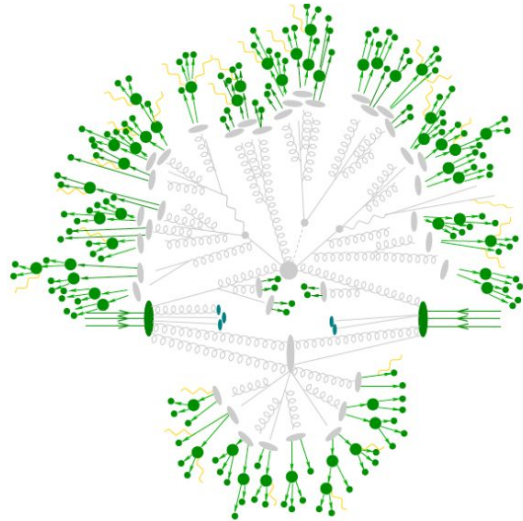
see <u>XAI talk in the Tuesday session</u>

# Jets at the Large Hadron Collider

Colliding protons at the LHC produce collections of particles called "jets"

- Observed as clusters of energy and tracks

# Jet Tagging: Classification in HEP

- Jets can emerge from many processes, and we want to identify the "type" of the process that gives us the jet

- Classic example from HEP: QCD and top jet classification

- Includes information about momenta ($p_x$, $p_y$, $p_z$) and energy ($e$) of up to 200 particles that make up the jet

- The total energy (momentum) of the jet is obtained by a scalar (vector) summation of the particle-level
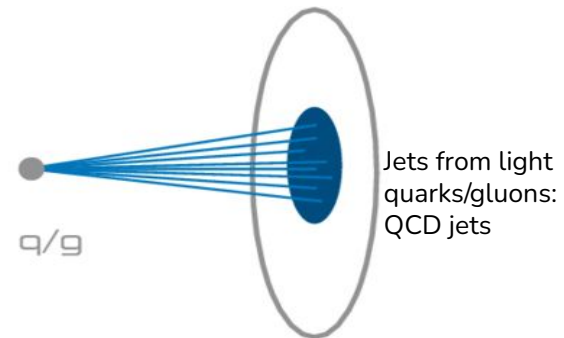
**Transverse momentum**

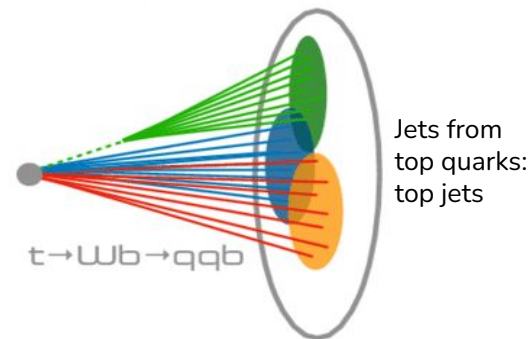$$p_t = \sqrt{p_x^2 + p_y^2}$$

**Azimuthal angle**

$$\phi = tan^{-1}\left(\frac{p_y}{p_x}\right)$$

**pseudorapidity**

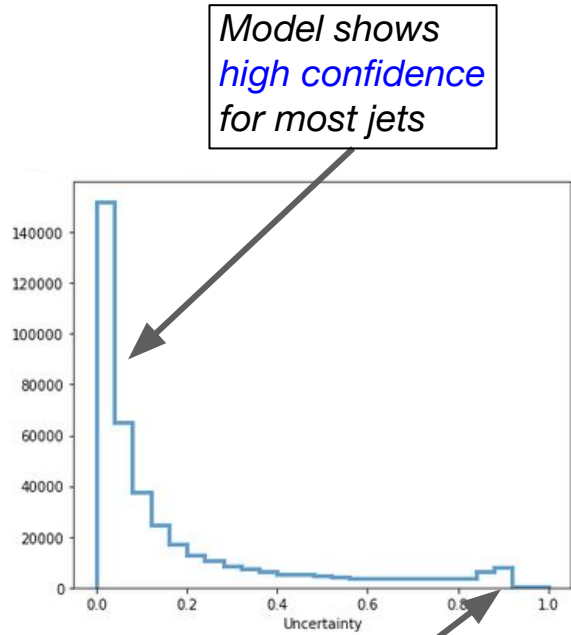$$\eta = \frac{1}{2}\ln\left(\frac{e + p_z}{e - p_z}\right)$$



q/g

Jets from light quarks/gluons: QCD jets

Simulated dataset with 2M jets available at: zenodo: 2603256



t→Wb→qqb

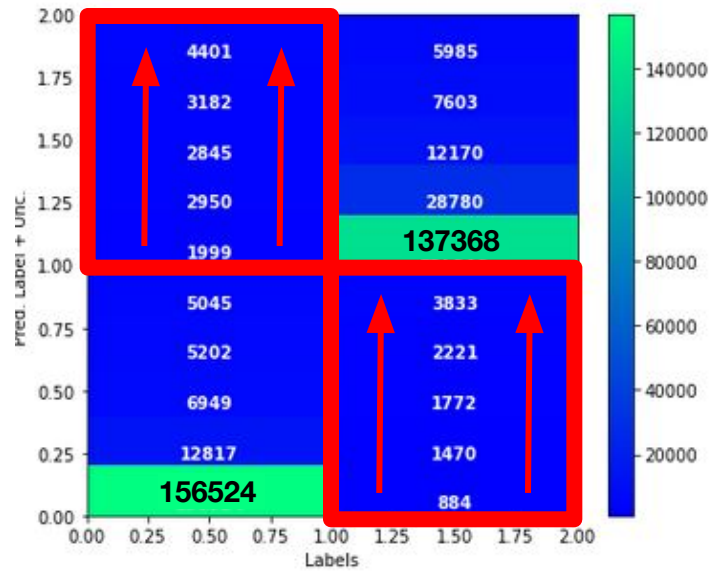Jets from top quarks: top jets

# Uncertainties in Jet Tagging - I

- **Q**: To what extent can a jet tagging model be confident in its predictions?

- Using the XAI-Inspired Interaction based Graph NN called Particle Flow Interaction Network (PFIN) for top tagging

- Binary classification goal: distinguish signal **top-quark jets** (*label*=1) from the background **QCD jets** (*label*=0)

Model shows *high confidence for most jets*

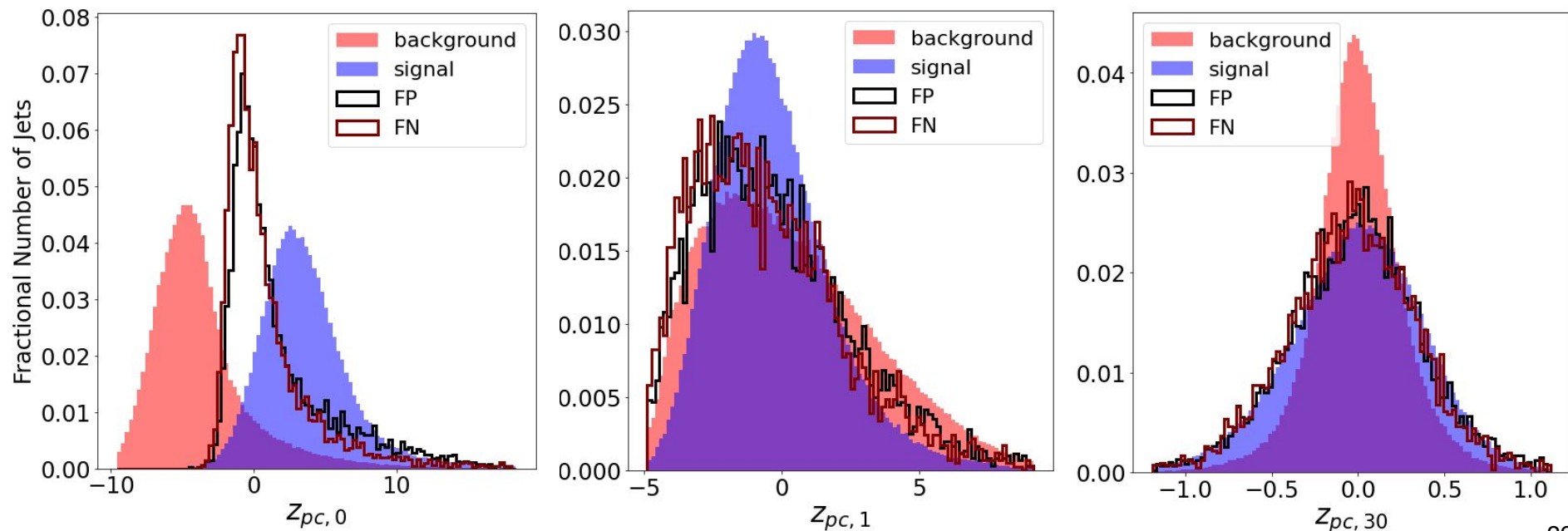*Large uncertainties dominated by misclassified jets!*
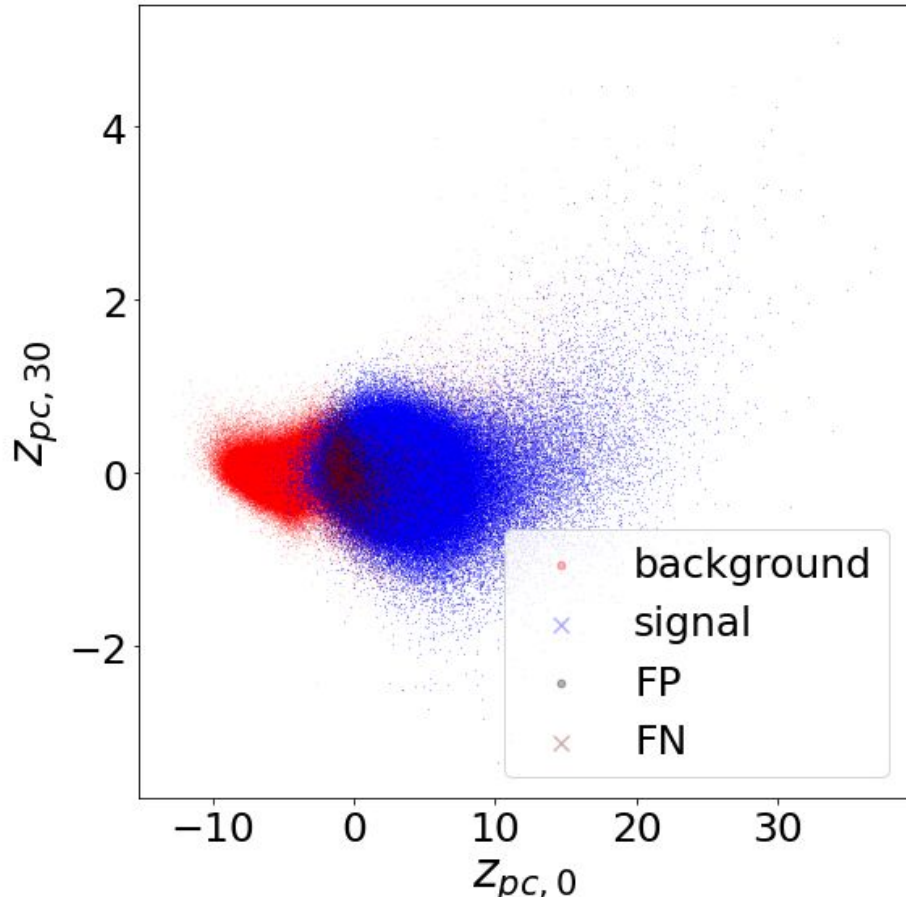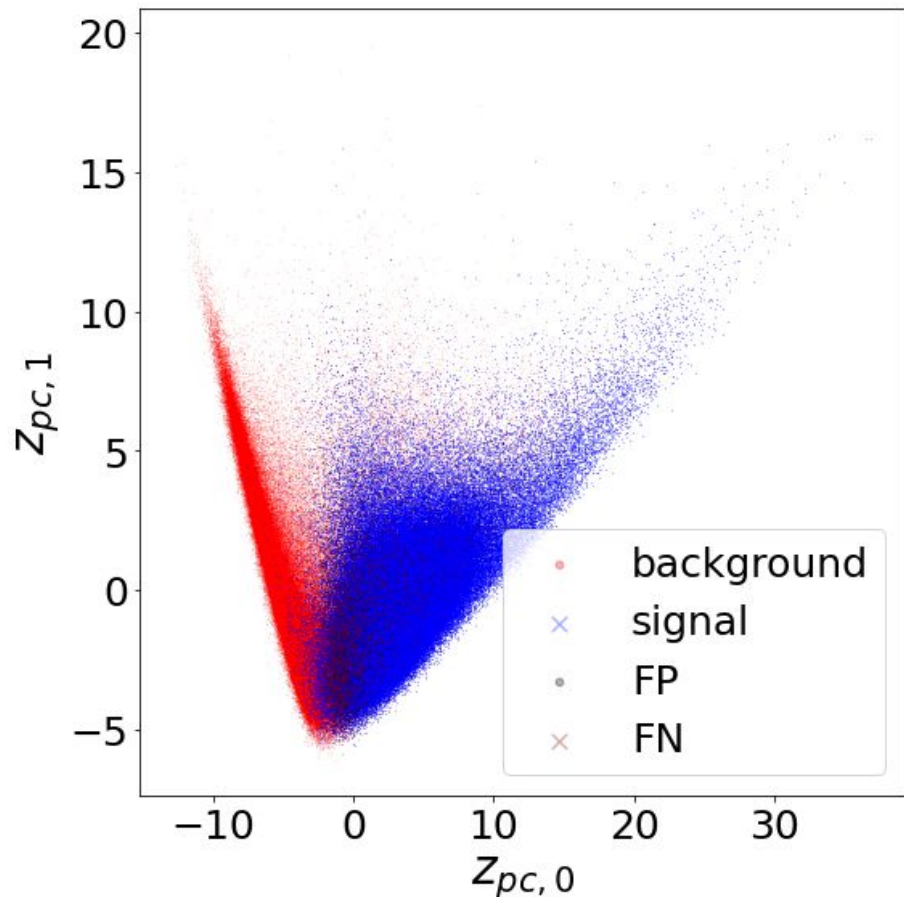


*Increasing uncertainty for misclassified jets!*

# Who Gets Largest Uncertainties?

Studies on XAI using Principal Component Analysis on the jet classifier model latent spaces show expressive discrimination (see also the XAI talk in the Tuesday session)

And we see that samples with large EDL-based uncertainty (> 0.8) lie in the overlap region, where discrimination is the hardest (expected "I don't know" from the model!)
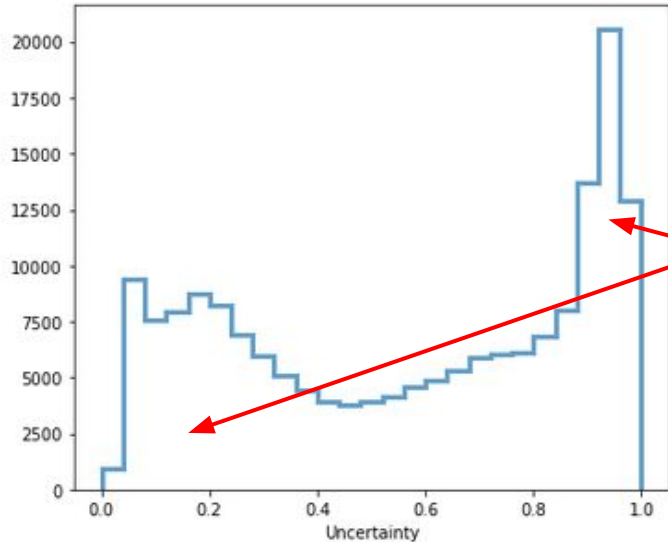
# Who Gets Largest Uncertainties? *(cont.)*
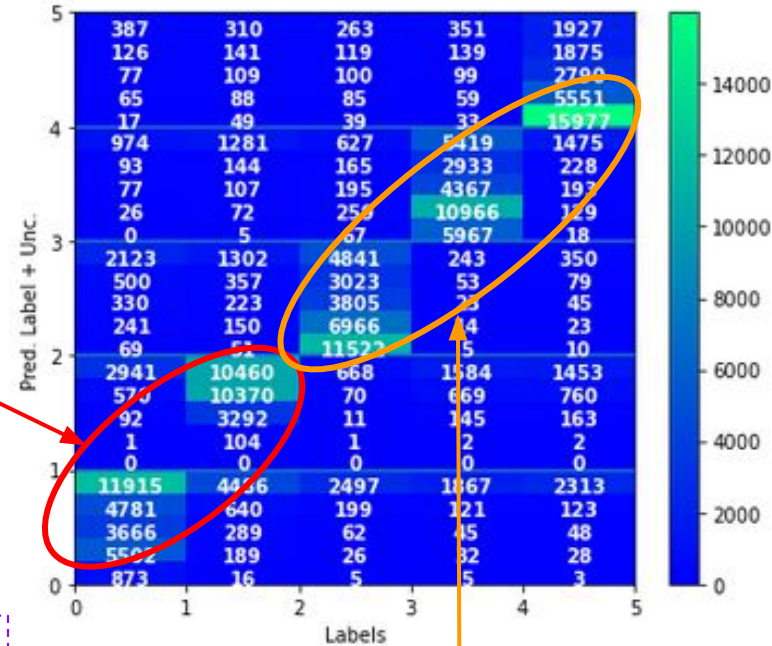
# Uncertainties in Jet Tagging - II

- PFIN model applied to a multi-class problem with JetNet Dataset: distinguishing jets coming from: *light quarks* (0), *gluons* (1), *top quarks* (2), *W bosons* (3), *Z bosons* (4)
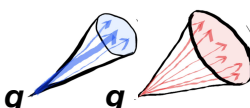


Bimodal distribution with a large peak at large uncertainties dominated by correctly classified quark and gluon jets

These jets have similar physical characteristics, and are hard** to tell apart

**but not impossible *q* *g* *More constituents w/ more uniform energy fragmentation and wider*

Heavier jets tend to have lower uncertainties

# Studies on Anomaly Detection

**Q**: What happens if the models encounter jets that they have not "seen" before (i.e. trained on)? 🤔

- **Anomaly detection** with EDL can be tested by withdrawing some jet classes from training dataset

  - **In-Distribution** (ID): jets the model is trained on

  - **Out of Distribution** (OOD): jets withdrawn from training

- Models trained with EDL tend to assign a large "uncertainty" score to anomalous (OOD) classes

  - Model saying "hmmm…I don't know" 🤔

- One major challenge remains: how do we distinguish "hard-to-tell" jets from "anomalous jets" using a single uncertainty metric?



Light Jets (q/g): In-distribution



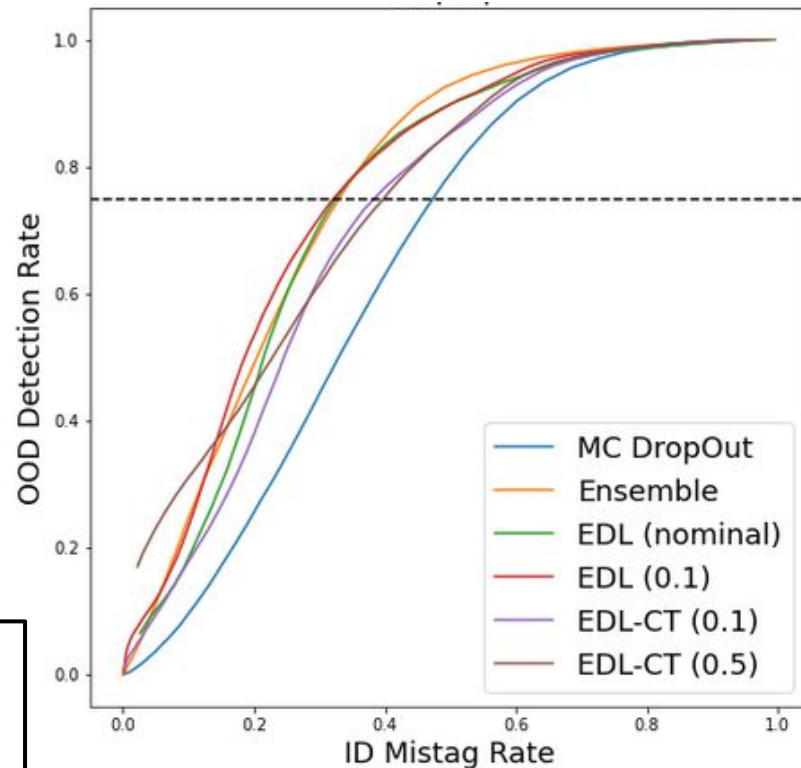Heavy Jets (t/W/Z): Out-of-distribution

# Comparing with Ensemble Methods

- Comparison can be done using ROC
  - A larger AUC would indicate a better performing model

- Key metrics:
  - **OOD Detection Rate**: what fraction of OOD samples are correctly identified
  - **ID Mistag Rate**: what fraction of ID samples are incorrectly identified

**EDL shows equivalent performance to ensemble methods and better than MC Dropout**



EDL-CT is a "Confidence Tuned" variant of the EDL method where the model is first allowed to converge w/o any annealing and then the parameters are tuned by retraining the model with annealing

# Lessons Learned and Future Work

Evidential Deep Learning (EDL) involves training a deterministic NN to place uncertainty priors over the predictive distribution, requiring only a single forward pass to estimate uncertainty

The EDL approach to uncertainty estimation proved to be well calibrated on the Top tagger and JetNet datasets and was capable of detecting OOD samples

EDL shows equivalent performance to ensemble methods and better than MC Dropout

*Challenge*: No clear metric to differentiate between "it's hard to tell" and "I don't know"

*Next steps*:

- Try these methods on the Jet Class dataset
- Bind in together with **One Class Classifier Methods** (OCC), as the current approach only works when at least two training classes exist
- Differentiate between uncertain ID samples and anomalous samples
- Explore the XAI aspects: exploring latent spaces is a good place to start

OCC

OCC trained to project in-distribution events within a hypersphere of radius $c$