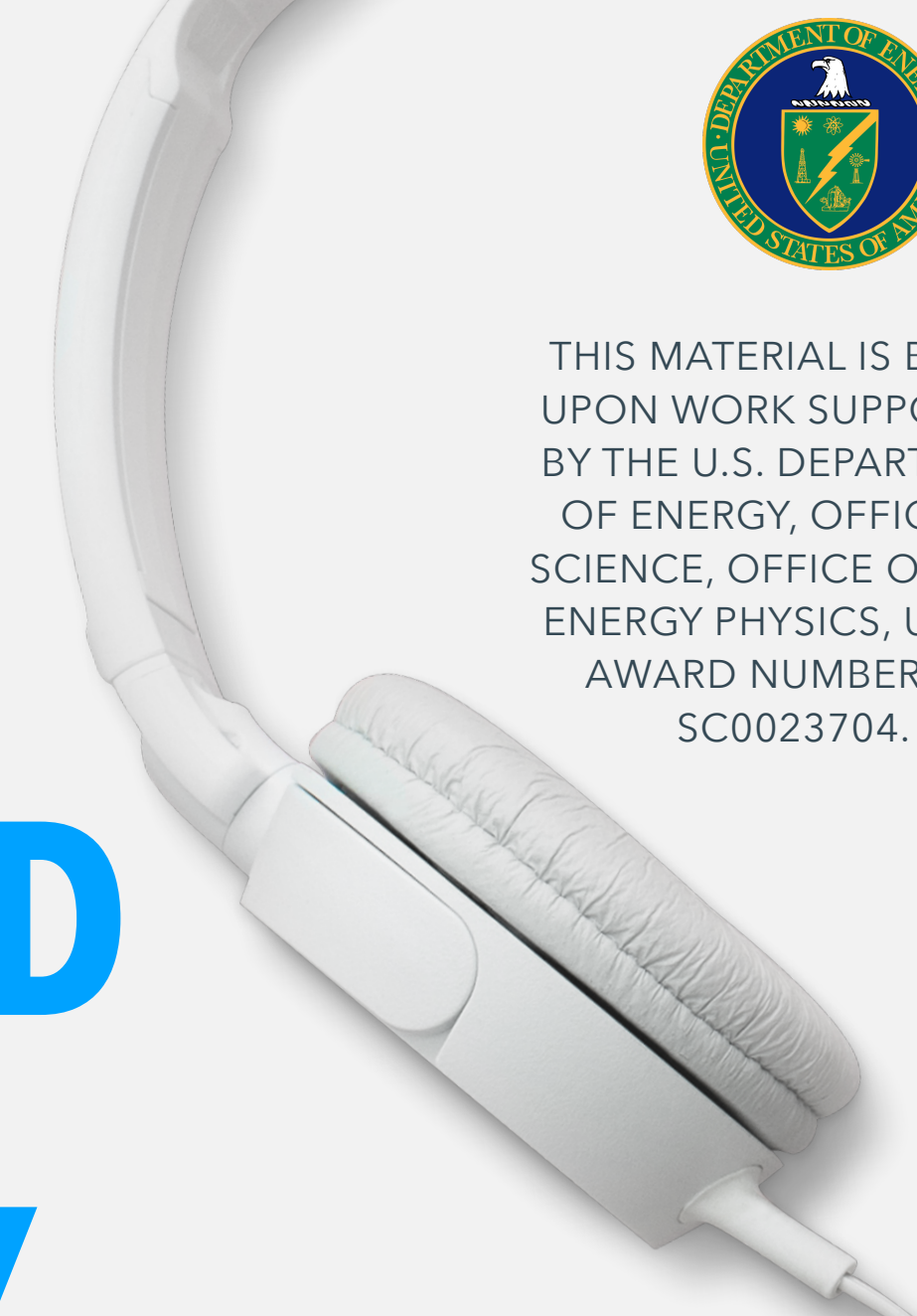


THIS MATERIAL IS BASED
UPON WORK SUPPORTED
BY THE U.S. DEPARTMENT
OF ENERGY, OFFICE OF
SCIENCE, OFFICE OF HIGH
ENERGY PHYSICS, UNDER
AWARD NUMBER DE-
SC0023704.

ENSEMBLES AND UNCERTAINTY QUANTIFICATION

I B R A H I M E L S H A R K A W Y



| ABOUT ME

Undergrad at Rice University

- Studied Physics, Applied Math, Philosophy

Second-Year Grad at University of Illinois at Urbana-Champaign

- Advisors: Prof. Yoni Kahn and Prof. Ben Hooberman



Research + Work:

- Two Research Projects **at CMS**
- Few Years of Interning as a **Research Geophysicist**
- **Physics for AI** work with Advisors
- Higgs Uncertainty Challenge

Physics
+
Machine Learning

BY THE END OF THE TALK

1. **Higgs Uncertainty Challenge:** Ensembles of Normalizing Flows, Systemic Uncertainty Robust Classifiers, Nuisance Parameters Estimation
2. **Uncertainty Quantifying From Scaling Laws:** Field Theory for NN, Infinite Width NN, and NN Scaling Laws

TABLE OF CONTENTS

01

Higgs Uncertainty Challenge

An Introduction to the Challenge, and the First Iteration



02

The Solution

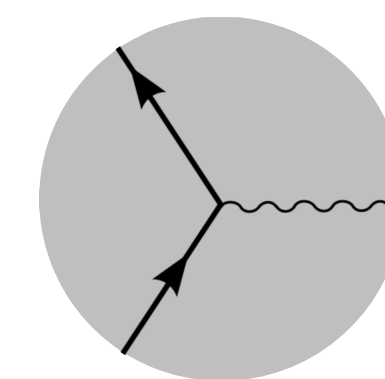
An Overview of the Final Iteration, involving NF Ensembles, Classifiers, and Estimating Nuisance Parameters



03

Uncertainty Quantifying From Scaling Laws

Physics for AI: How we can use field theory to predict NN behavior with Infinite Width Networks and Scaling Laws



04

Empirical Results

Empirical Results Comparing Theoretical Results with Empirical Reality



01

Higgs Uncertainty

Challenge



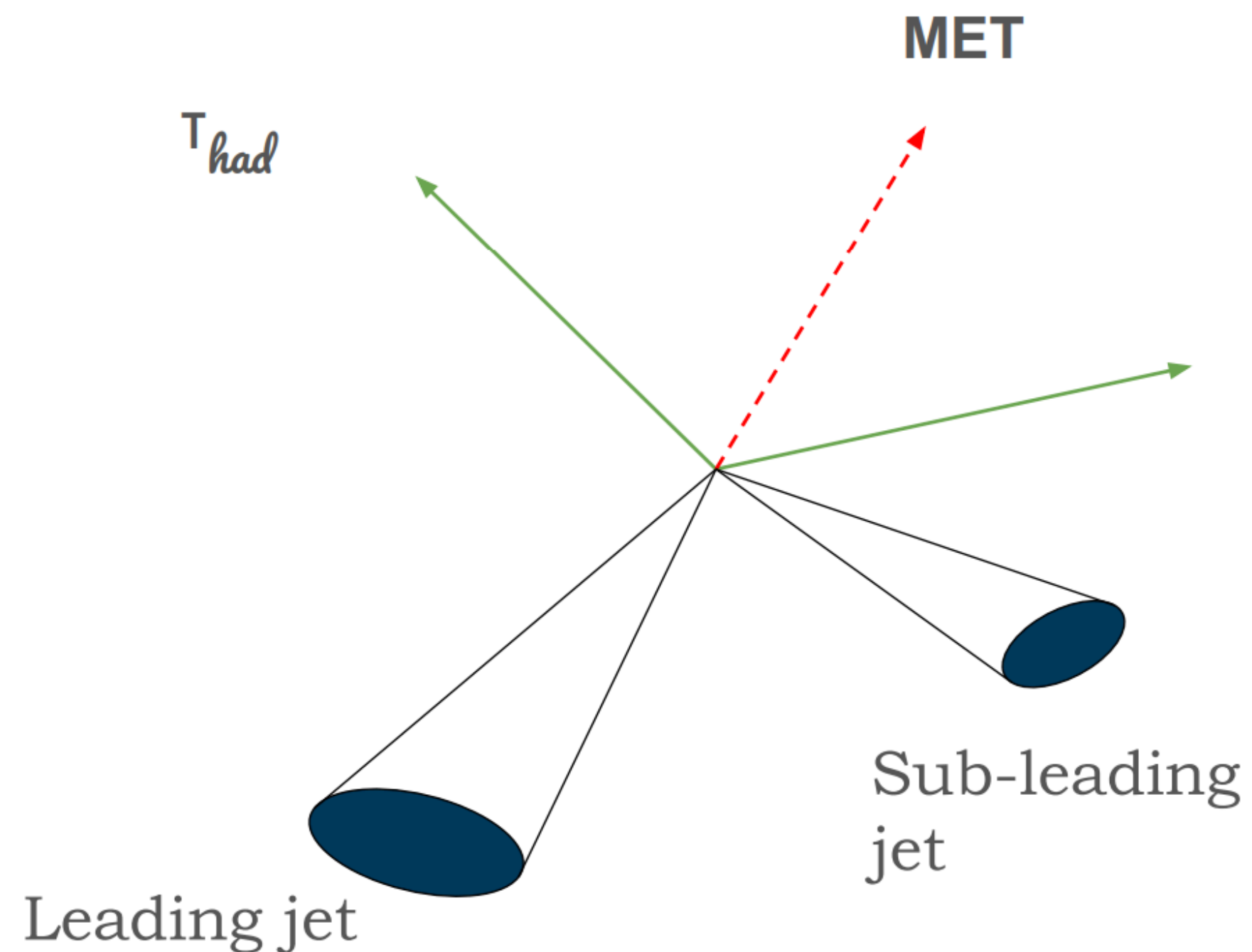
An Introduction to the Challenge,
and the First Iteration



Introduction: The Higgs Uncertainty Challenge

The Goal:

- 1) **Measure** the signal strength $\mu = \frac{\text{Observed Higgs}}{\text{Expected Higgs}}$
- 2) **Give** correct and small 68% CI on the measurement



The **signal process** is $H \rightarrow \tau\tau$

Data: **28 Input Features**

Introduction: The Higgs Uncertainty Challenge

Variable	Mean	Sigma	Range
α_{tes}	1.	0.01	[0.9, 1.1]
α_{jes}	1.	0.01	[0.9, 1.1]
α_{soft_met}	0.	3.	[0., $+\infty$]
α_{ttbar_scale}	1.	0.25	[0., $+\infty$]
$\alpha_{diboson_scale}$	1.	0.025	[0., $+\infty$]
α_{bkg_scale}	1.	0.01	[0., $+\infty$]

Six nuances parameters

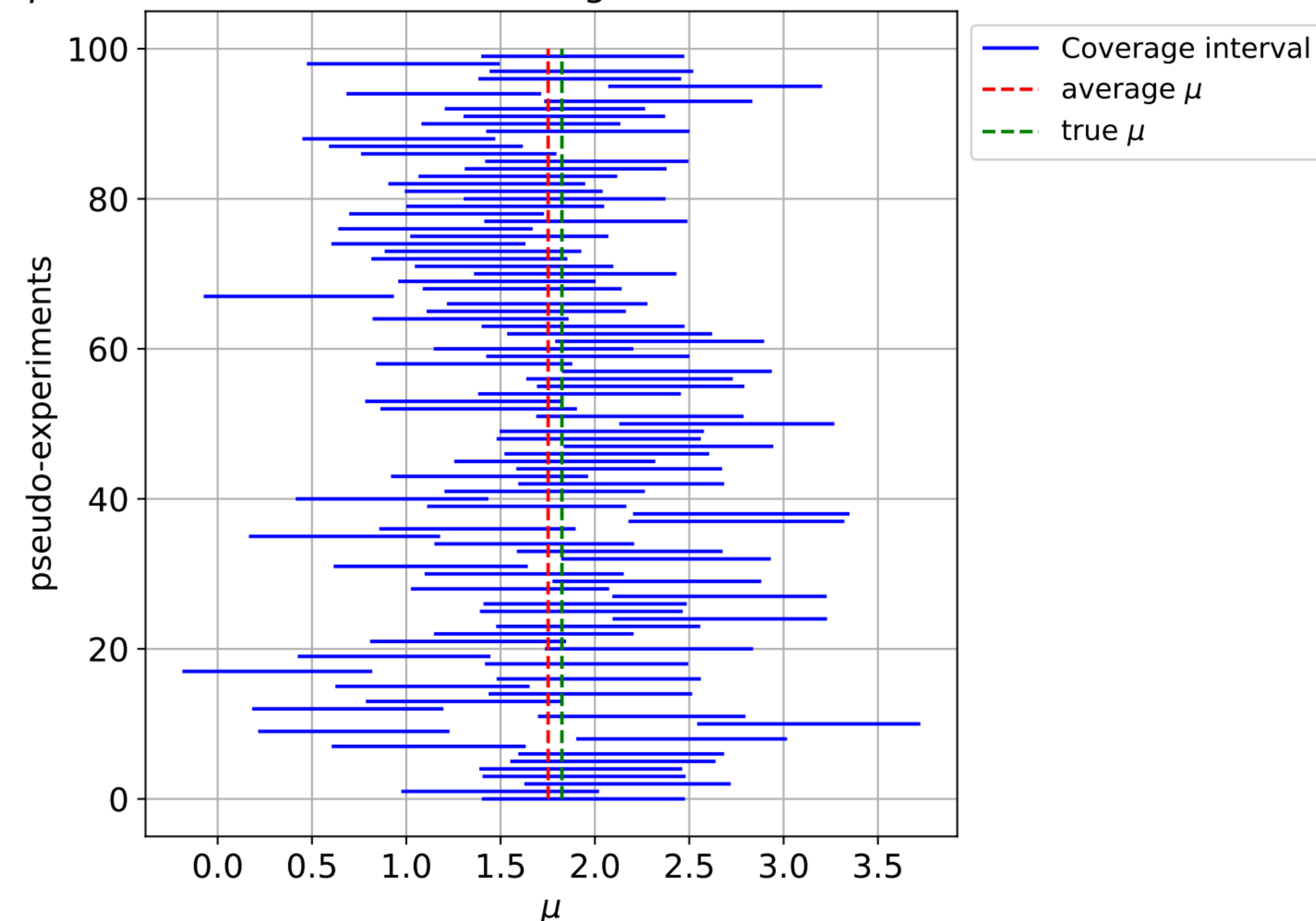
Distorts the 28 features in a unknown nonlinear way

Method is evaluated by:

- Running 100 pseudo-experiments with **different nuisance parameters**
- On 10 **different values** of $\mu = [.1, 3]$

CI must be correct $\sim 68\%$ of the time and are **rewarded with smaller intervals**

μ distribution - Set 0 - Coverage: 0.700 - Interval: 1.068



First Iteration: A Bayesian Approach

$$\mathcal{P}(\mu | \{x\}) \propto \mathcal{P}(\{x\} | \mu, \theta) = \prod_j^n \mathcal{P}(x_i | \mu, \theta)$$

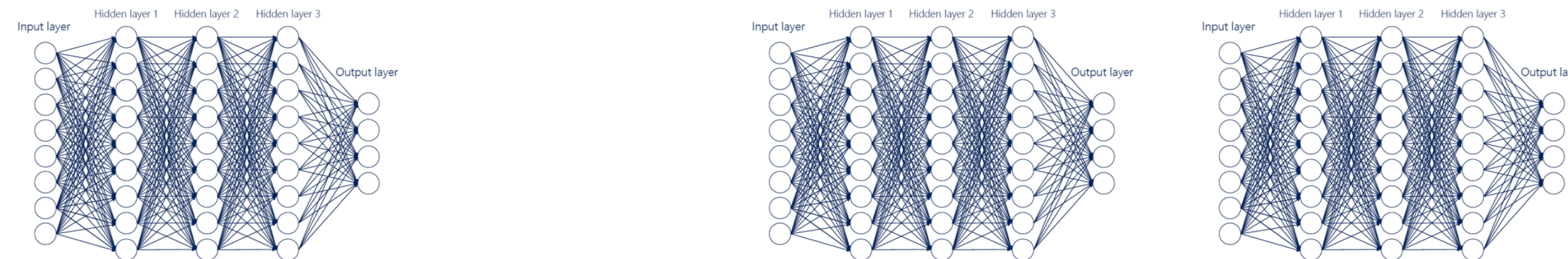
$$= \prod_i^N \left(\frac{\mu}{N} \mathcal{P}(x_i^s | \theta) + \frac{N - \mu}{N} \mathcal{P}(x_i^{bg} | \theta) \right)$$

How do we estimate these?

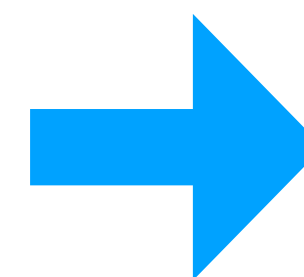
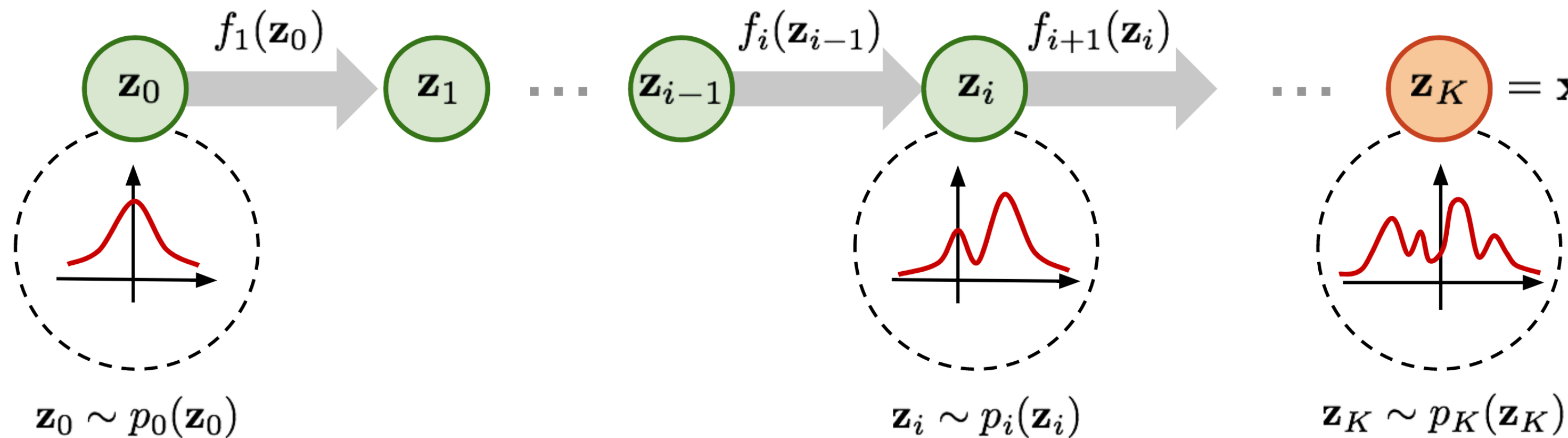
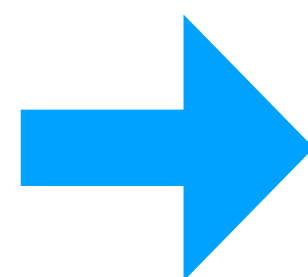
Here

μ = Observed Higgs

First Iteration: Normalizing Flows



$\mathcal{P}(x_i)$
 $\sim \mathcal{N}(0,1)$



$\mathcal{P}(x_i^{bg})$

OR

$\mathcal{P}(x_i^s)$



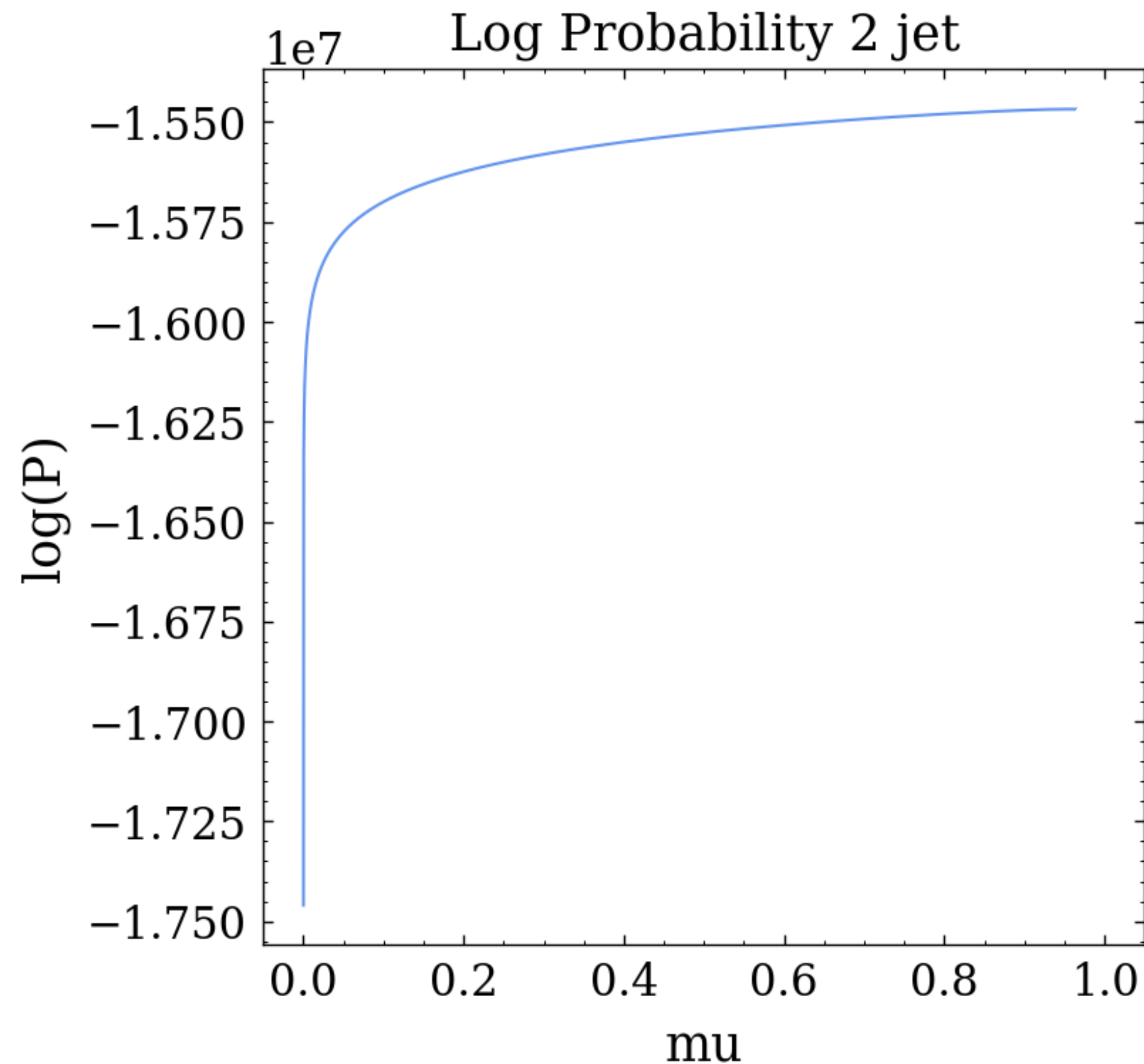
$$\mathcal{L}(\theta) = - \mathbb{E} \left[\log p_Z (f_\theta(x)) + \log \left| \det \left(\frac{\partial f_\theta}{\partial x} \right) \right| \right]$$

First Iteration: One Problem

$$\mathcal{P}(\{x\} | \mu, \theta) = \prod_j^n \mathcal{P}(x_j | \mu, \theta) = \prod_i^N \left(\frac{\mu}{N} \mathcal{P}(x_i^s | \theta) + \frac{N - \mu}{N} \mathcal{P}(x_i^{bg} | \theta) \right)$$

Not Enough

Information in the NF-
Likelihood!

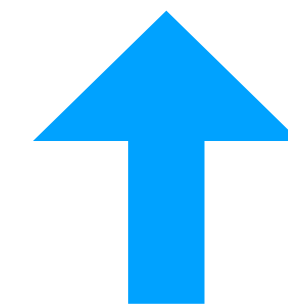
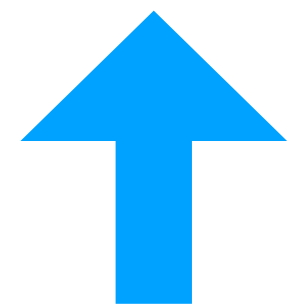


First Iteration: Adversarial Loss

$$\mathcal{L}(\theta) = - \mathbb{E} \left[\log p_Z (f_\theta(x)) + \log \left| \det \left(\frac{\partial f_\theta}{\partial x} \right) \right| \right]$$



$$\mathcal{L}(\theta) = - \mathbb{E} \left[c \cdot \left(\log p_Z (f_\theta(x)) + \log \left| \det \left(\frac{\partial f_\theta}{\partial x} \right) \right| \right) - \log p_Z (f_\theta(y)) - \log \left| \det \left(\frac{\partial f_\theta}{\partial y} \right) \right| \right]$$



For Example: x = signal events, y = background events

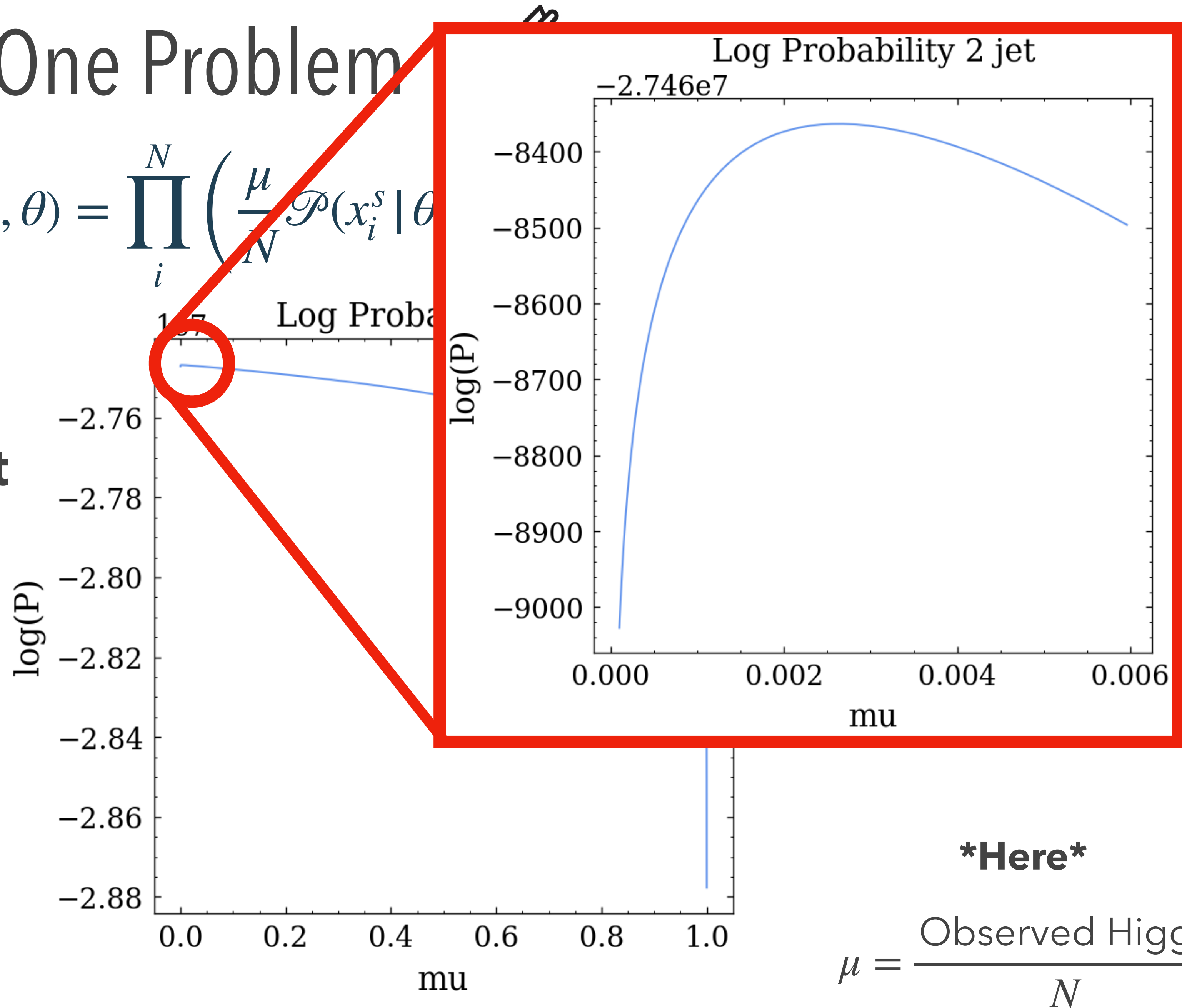
First Iteration: One Problem

$$\mathcal{P}(\{x\} | \mu, \theta) = \prod_j \mathcal{P}(x_j | \mu, \theta) = \prod_i \left(\frac{\mu}{N} \mathcal{P}(x_i^s | \theta) \right)$$

Peaks at close to the right value of mu!

$$\mu_{real} \approx 0.0026$$

$$\mu_{peak} \approx 0.002$$

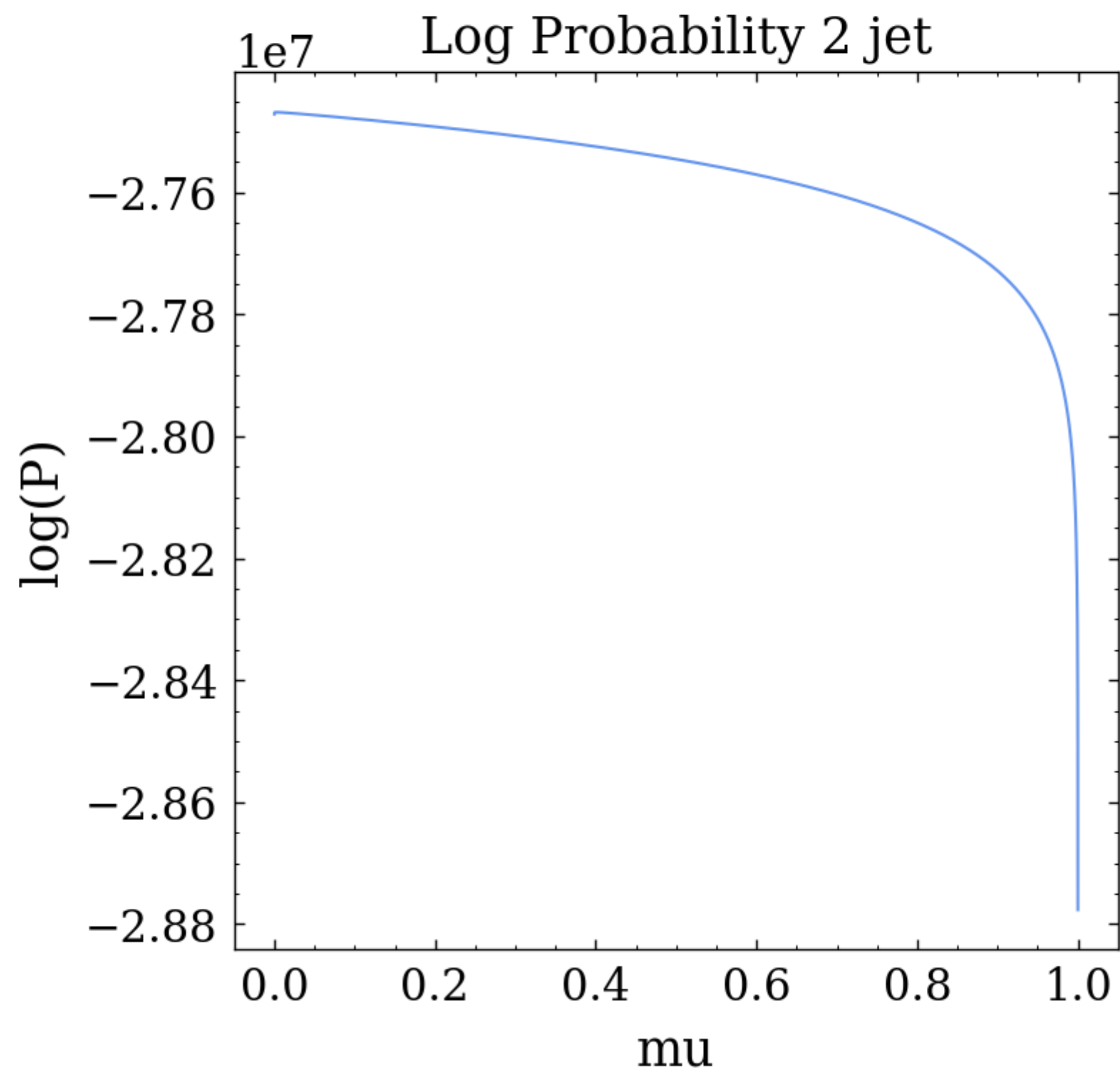


Here

$$\mu = \frac{\text{Observed Higgs}}{N}$$

First Iteration: Not Sufficient

$$\mathcal{P}(\{x\} | \mu, \theta) = \prod_j \mathcal{P}(x_j | \mu, \theta) = \prod_i \left(\frac{\mu}{N} \mathcal{P}(x_i^s | \theta) + \frac{N - \mu}{N} \mathcal{P}(x_i^{bg} | \theta) \right)$$



Issue: Not robust to:

1) Systematics

- Very incorrect predictions

2) Changes in μ

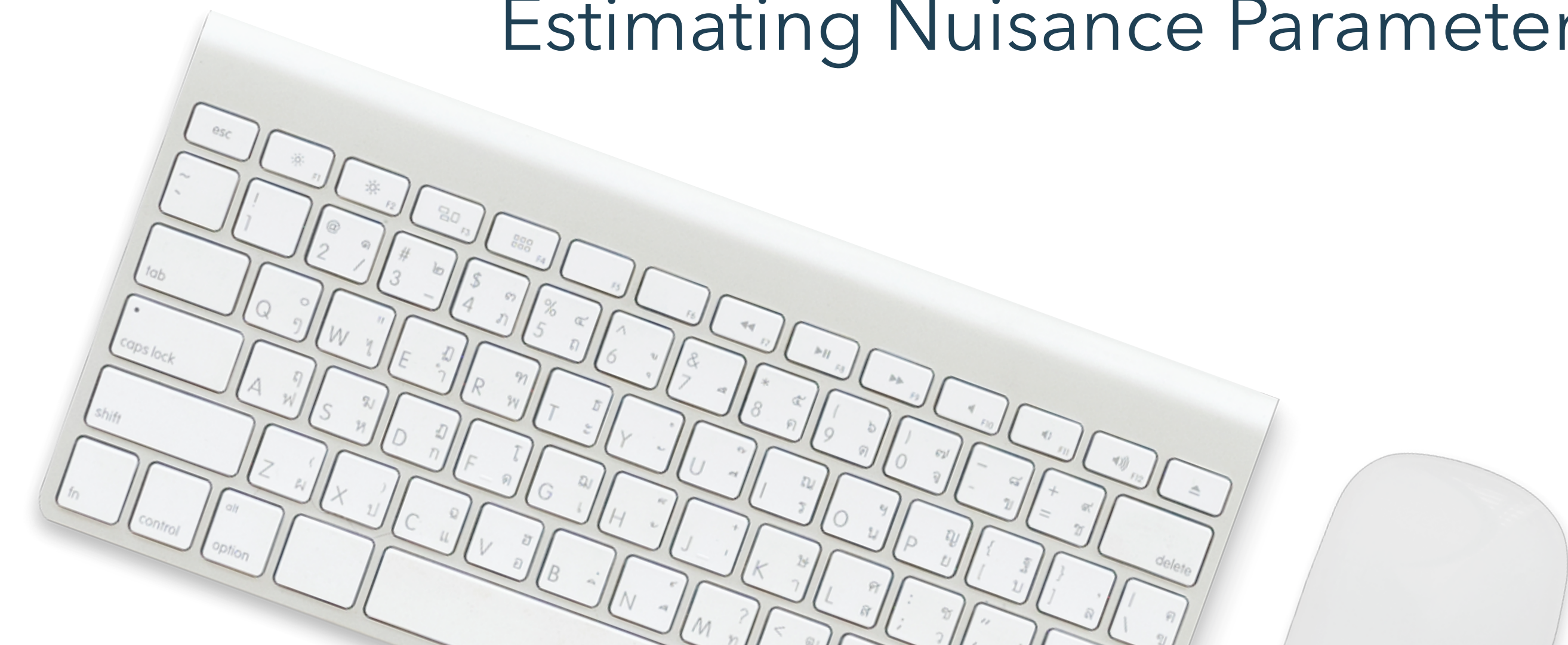
- "Sticky" peak and no sensitivity for $\mu < .5$

The Solution

An Overview of the Final Iteration,
involving NF Ensembles, Classifiers, and
Estimating Nuisance Parameters



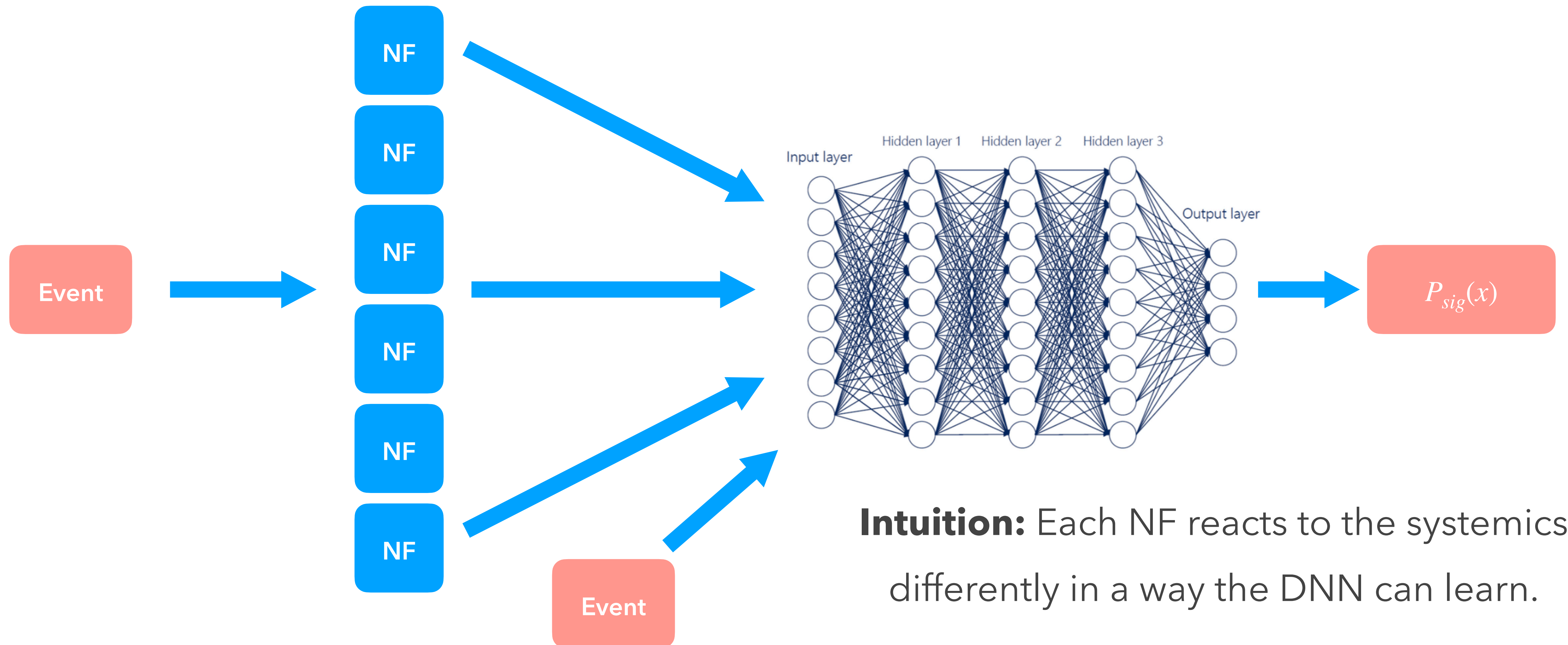
02



Final Iteration: Back to Classifiers



Idea: Use Ensemble of NF Likelihoods of as input to a simple DNN classifier



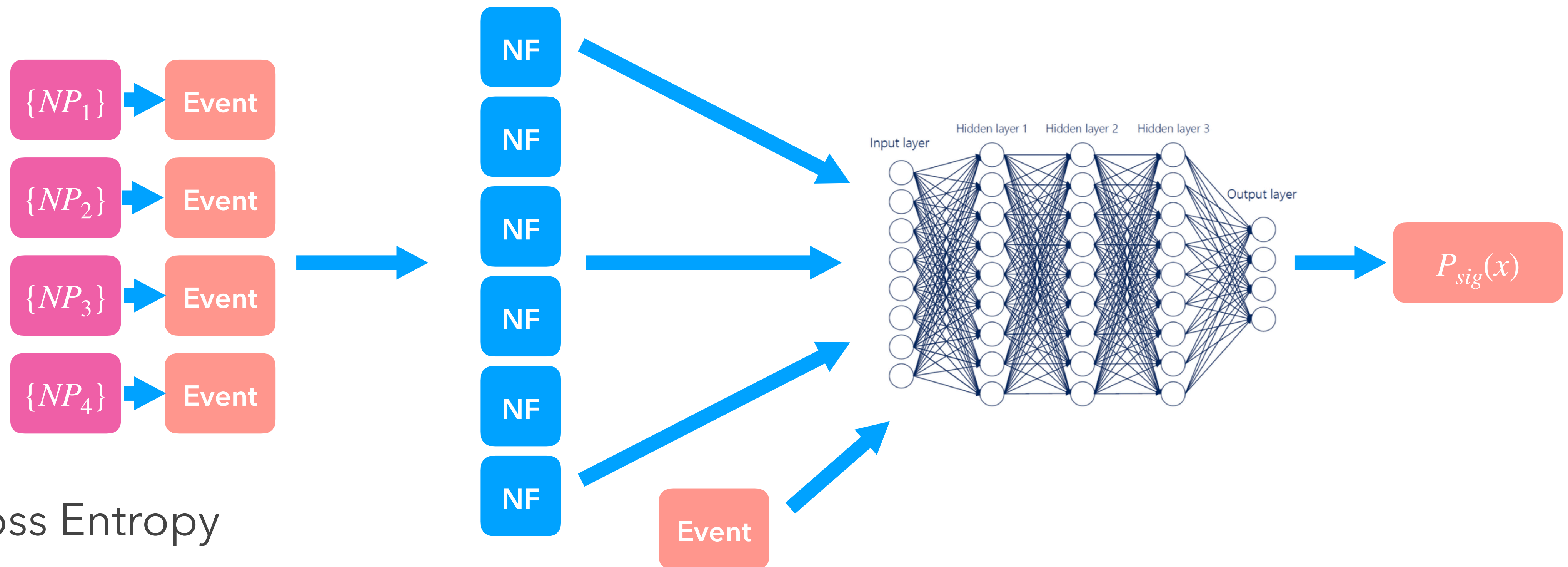
Intuition: Each NF reacts to the systemics differently in a way the DNN can learn.

Final Iteration: Classifier Training



Systematic Robust Training:

Train the same DNN with event data perturbed with a wide variety of nuisance parameters.



Loss: Binary Cross Entropy

$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^N [y_i \cdot \log(\hat{y}_i) + (1 - y_i) \cdot \log(1 - \hat{y}_i)]$$

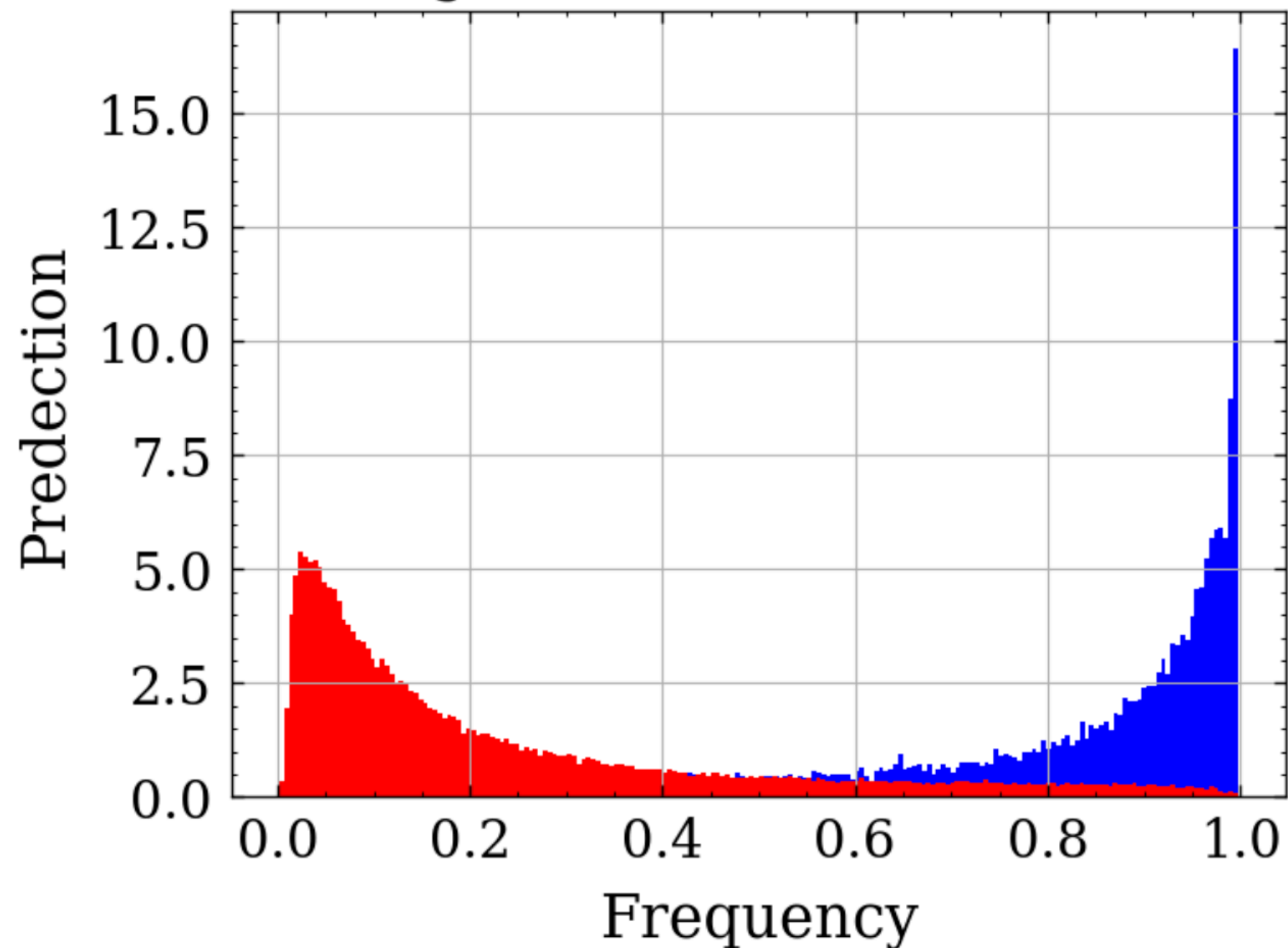
Final Iteration: Binned Analysis



Here

$$\mu = \frac{\text{Observed Higgs}}{N}$$

Histogram of DNN Discriminator



Let k_i be the true counts of events in bin i :

$$E[k_i] = f_i \cdot \mathbf{E}[S_i] + b_i$$

The **Poisson Likelihood** is then given by:

$$P(\{f_i\} | \{k_i\}) \propto \mathcal{L}(\{k_i\} | \{f_i\}) = \prod_{i=1}^N \frac{(f_i S_i + b_i)^{k_i} e^{-(f_i S_i + b_i)}}{k_i!}$$

Then:

$$\mu \equiv \frac{\sum_i f_i \cdot \mathbf{E}[S_i]}{\sum_i (b_i + f_i \cdot \mathbf{E}[S_i])}$$

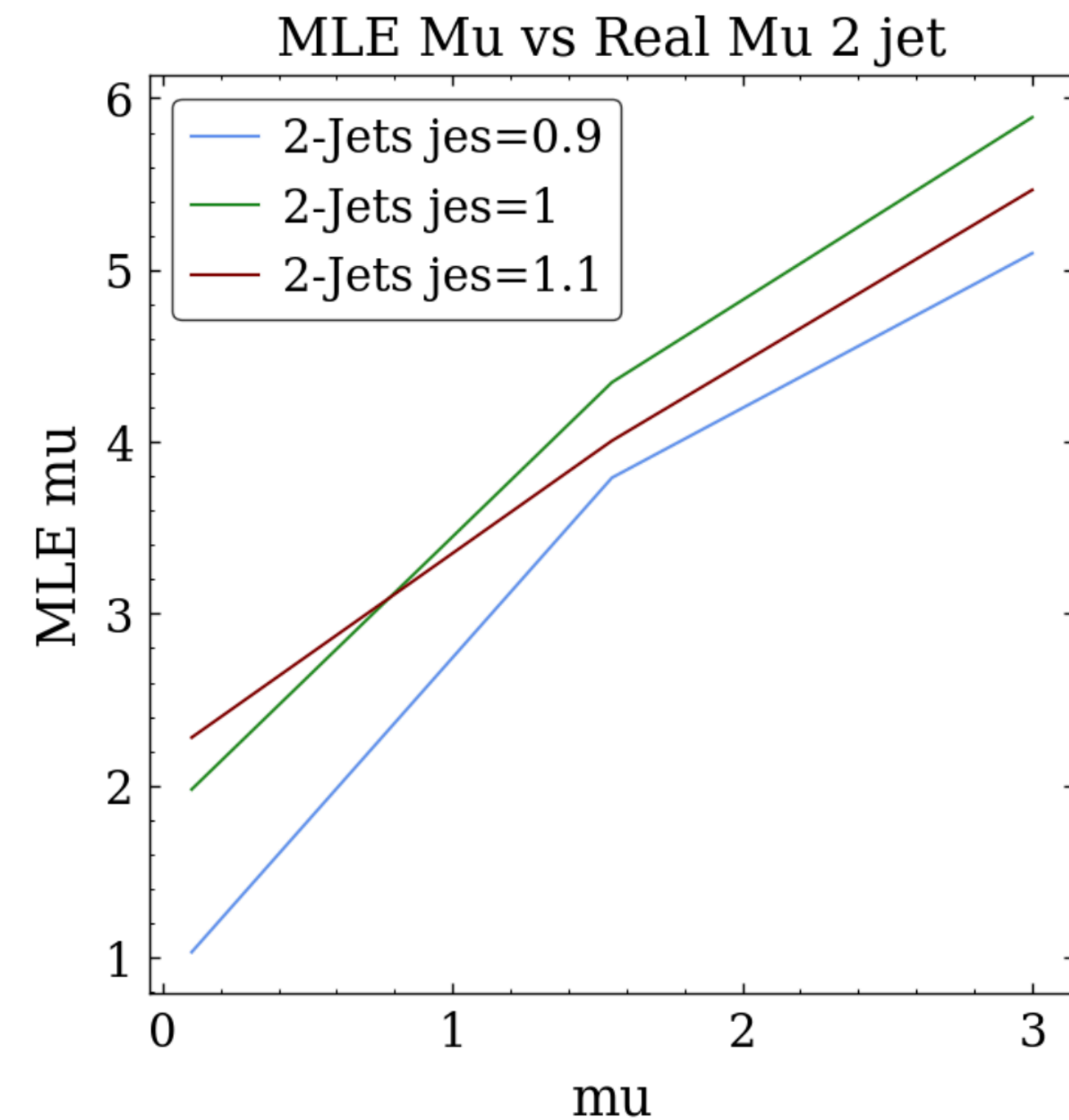
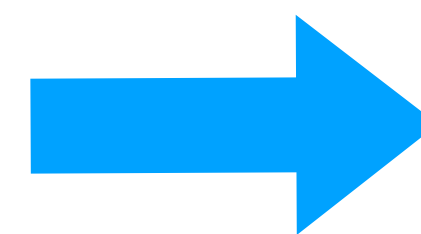
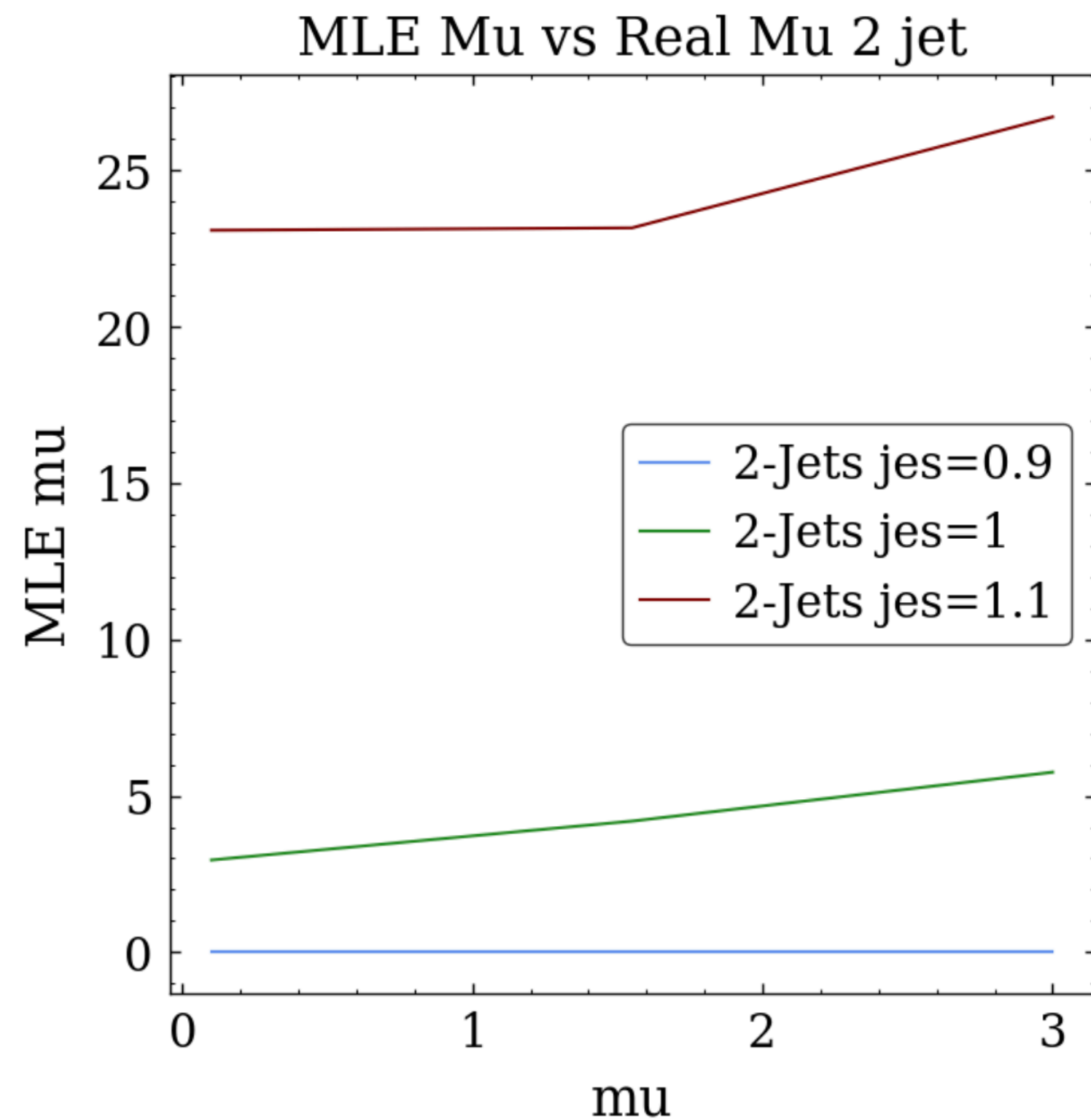
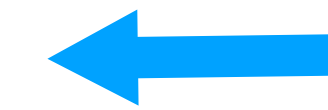
Final Iteration: Parameter Estimation

To deal with the **"worst" nuisance parameter** (j_{es}) we can do the same procedure but compute instead:

Do MLE on both θ and f_i !

$$P(\{f_i\} | \{k_i\}, \theta) \propto \mathcal{L}(\{k_i\} | \{f_i\}, \theta) = \prod_{i=1}^N \frac{(f_i \mathbf{E}[S_i | \theta] + b_i)^{k_i} e^{-(f_i \mathbf{E}[S_i | \theta] + b_i)}}{k_i!} \cdot P(\theta)$$

Prior on
Nuances Parameters

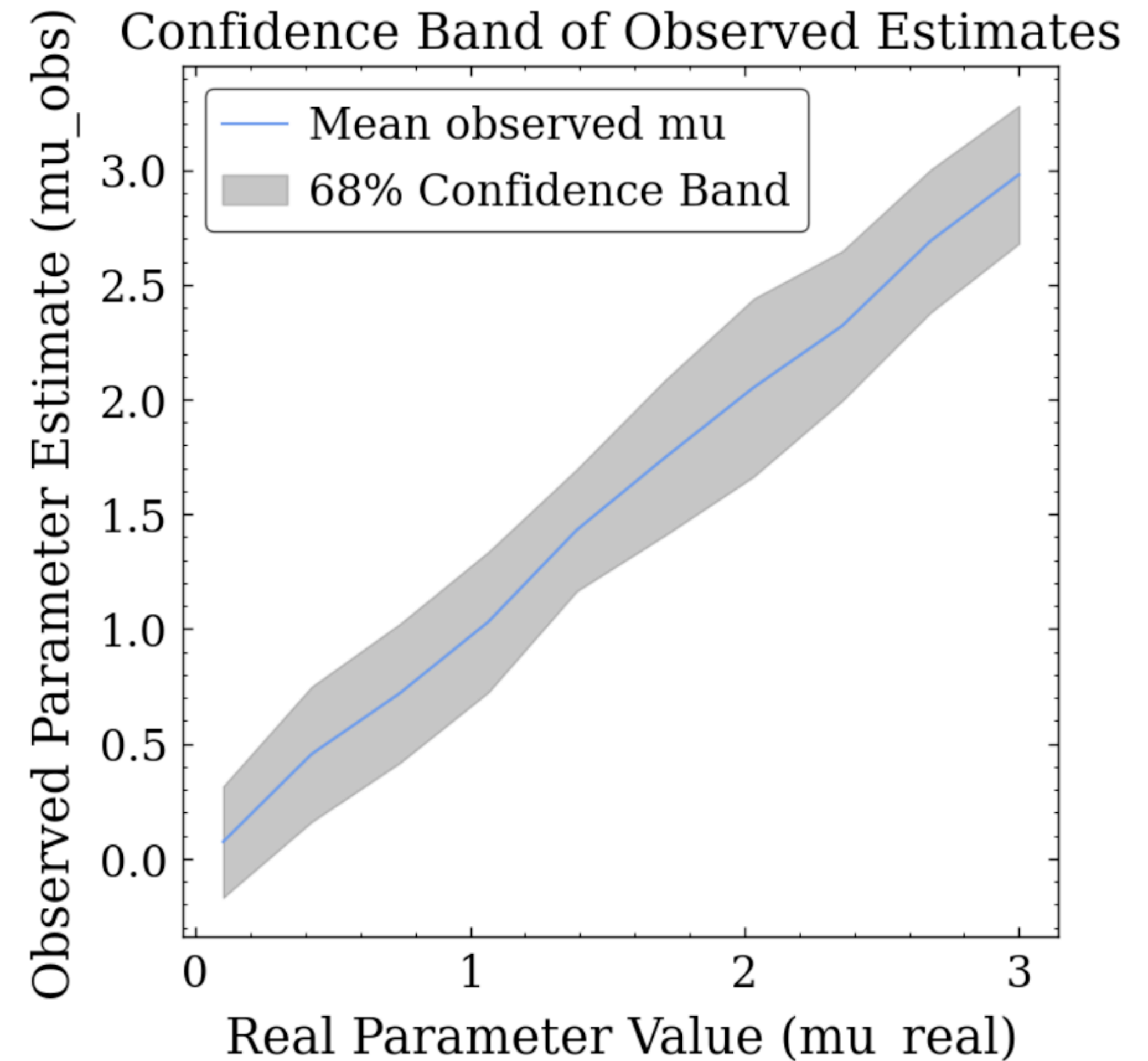


Final Iteration: Neyman Construction



For Error Bars we use the **Neyman Construction**

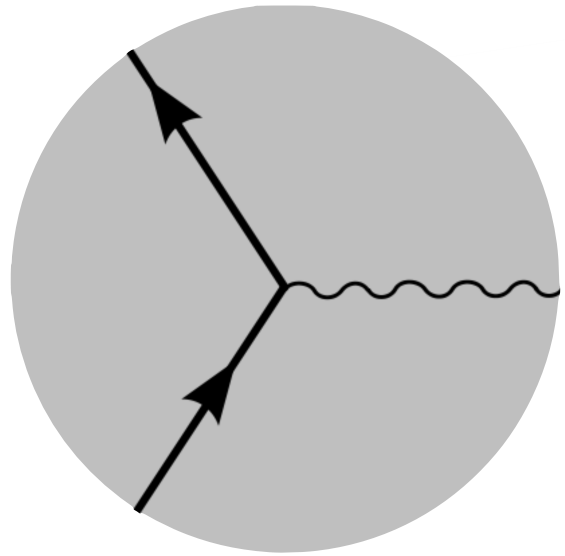
- 1) **For each value of real μ** compute estimate
~100 time with a different draw of NP
- 2) With mean and std of the estimates:
 - 1) Apply a **Bias Correction** to the mean
 - 2) And use std as error bar estimate **for a given estimate of μ**



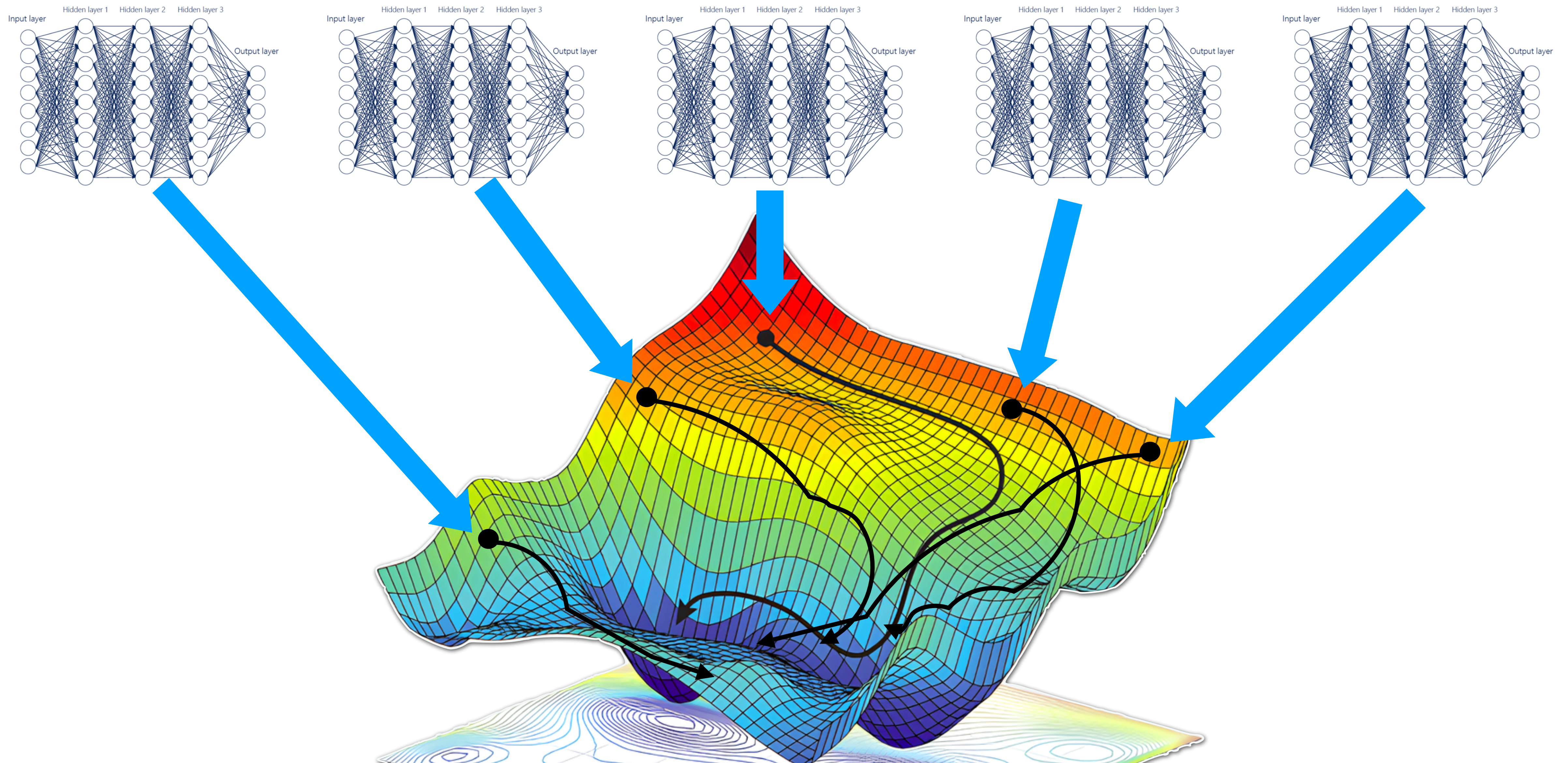
03

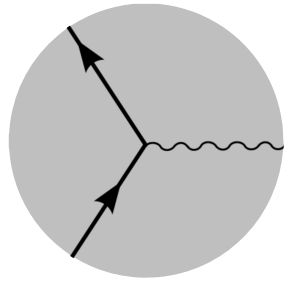
Uncertainty Quantifying From *Scaling* Laws

Empirical Results comparing
theoretical results with empirical reality

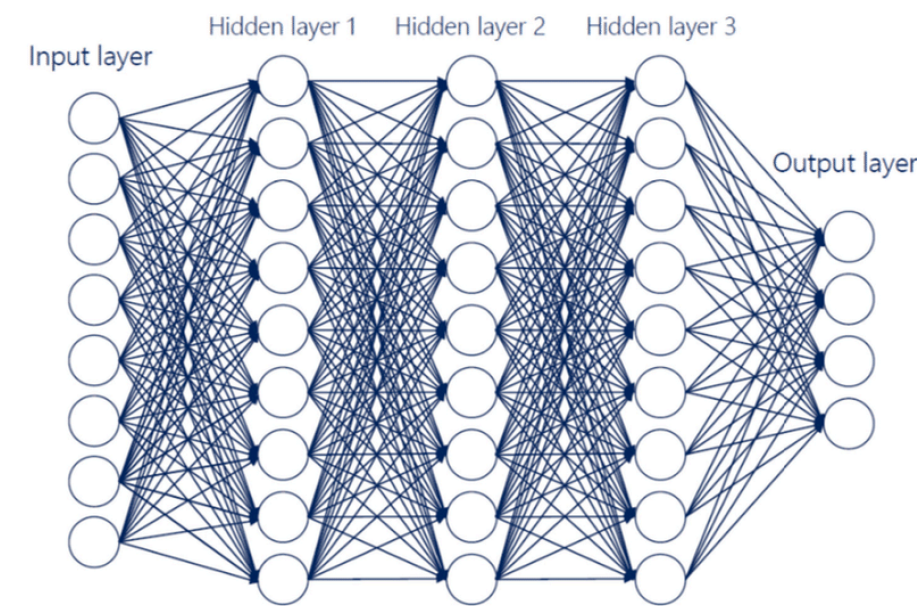
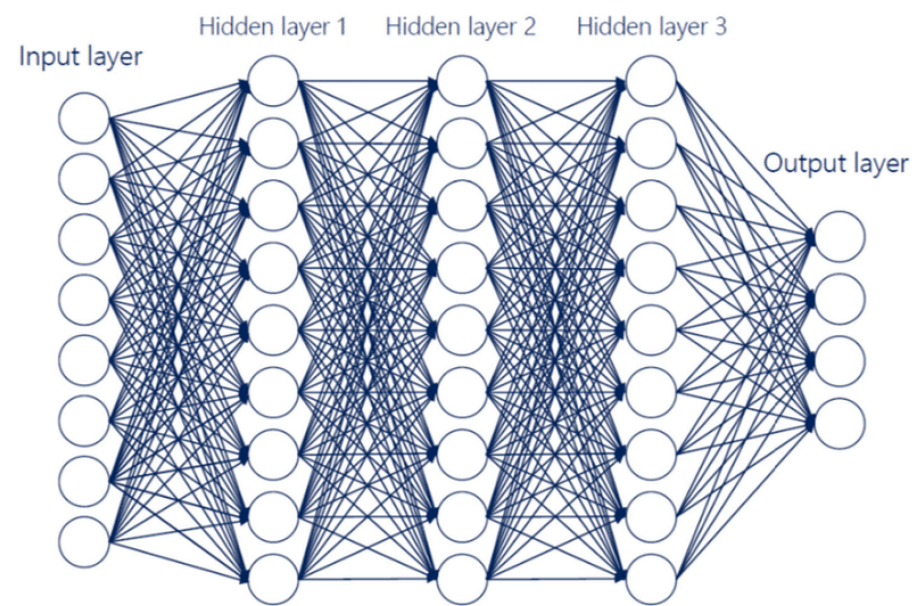
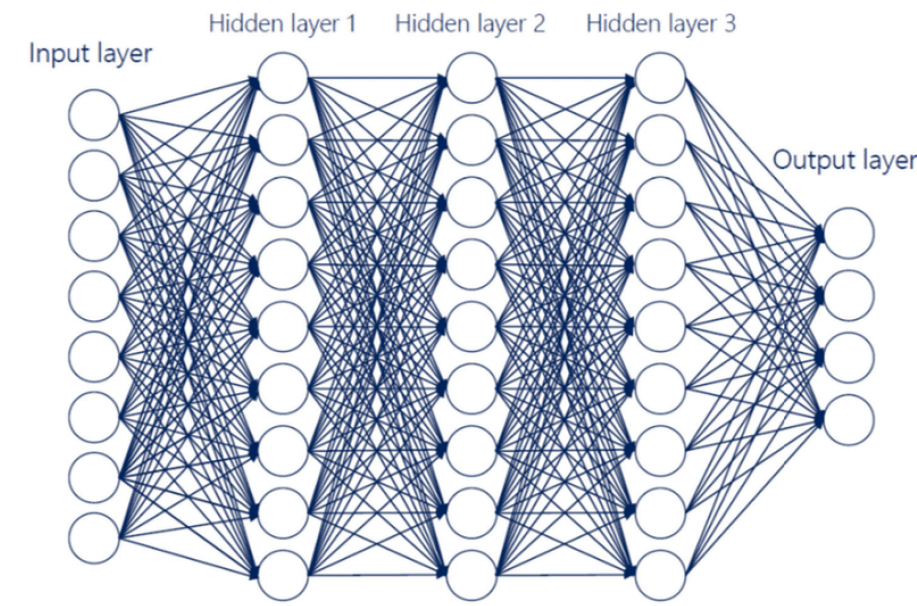
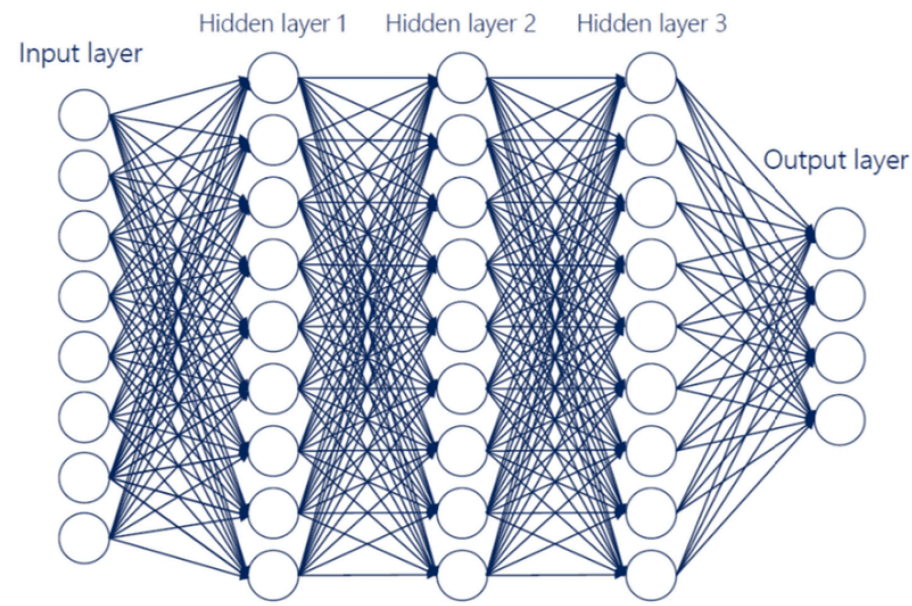
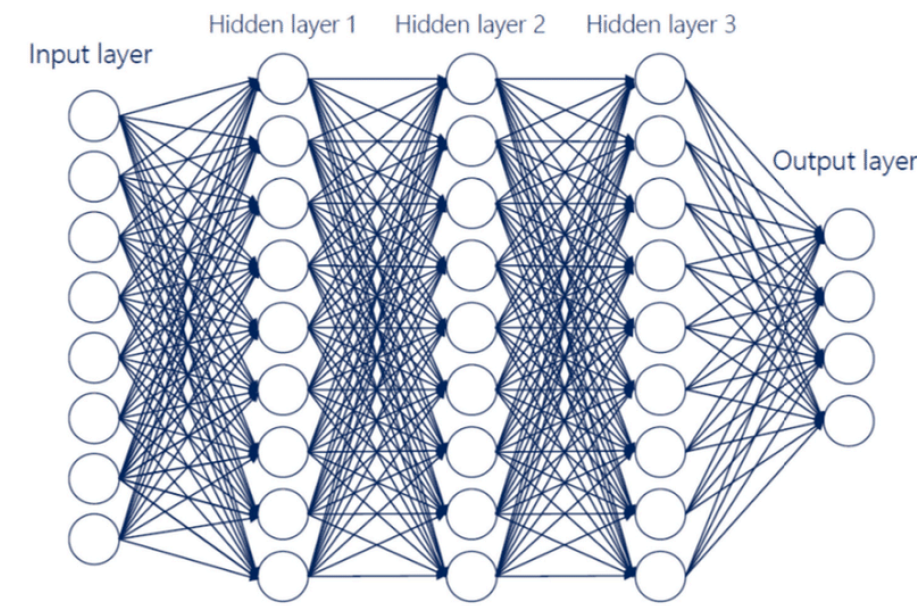
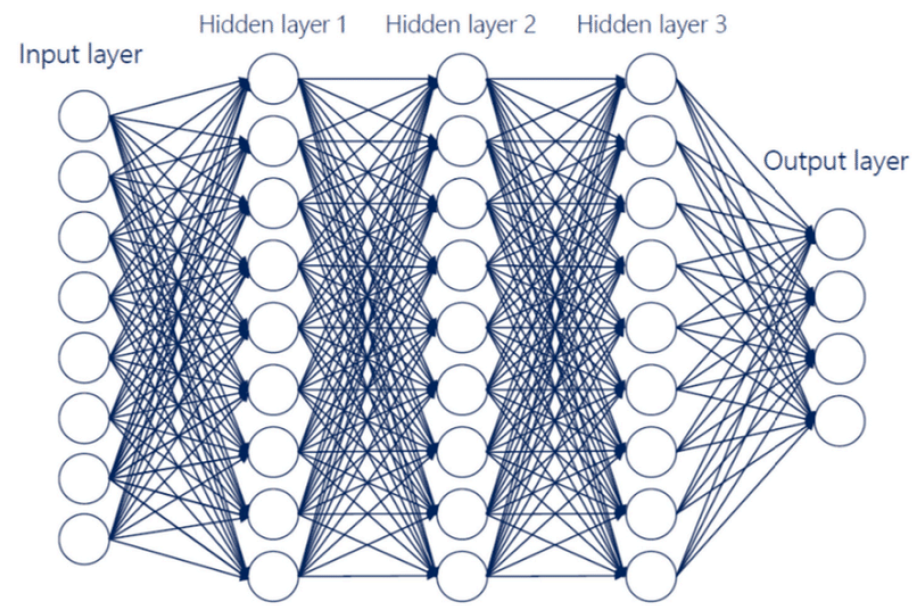


Uncertainty Quantification: Ensemble Uncertainty

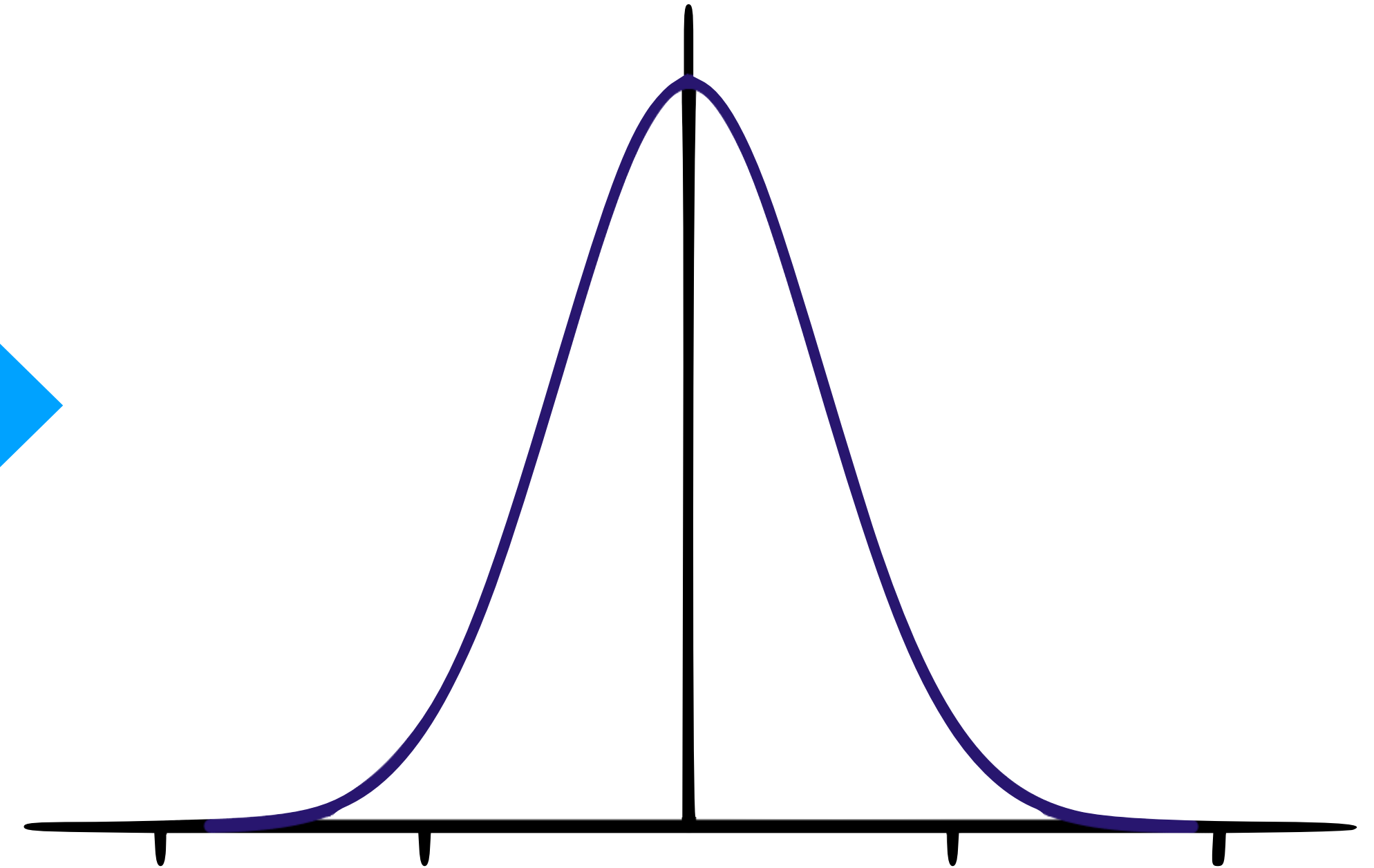
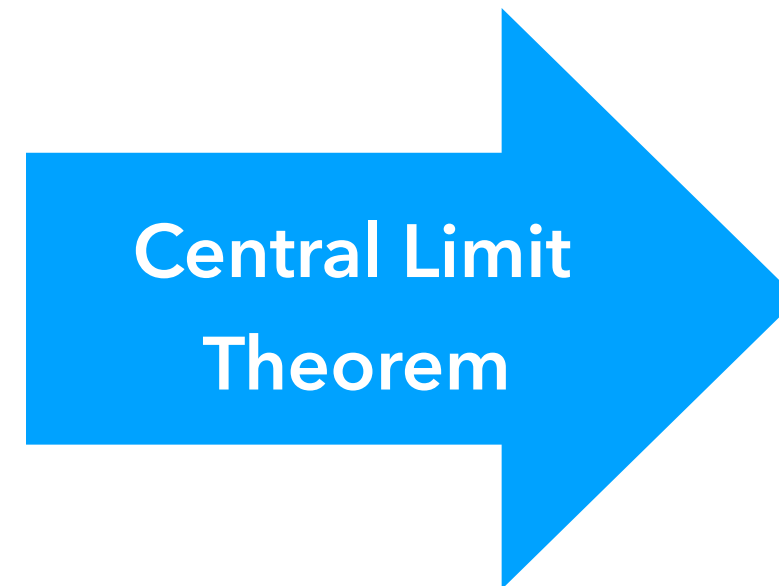




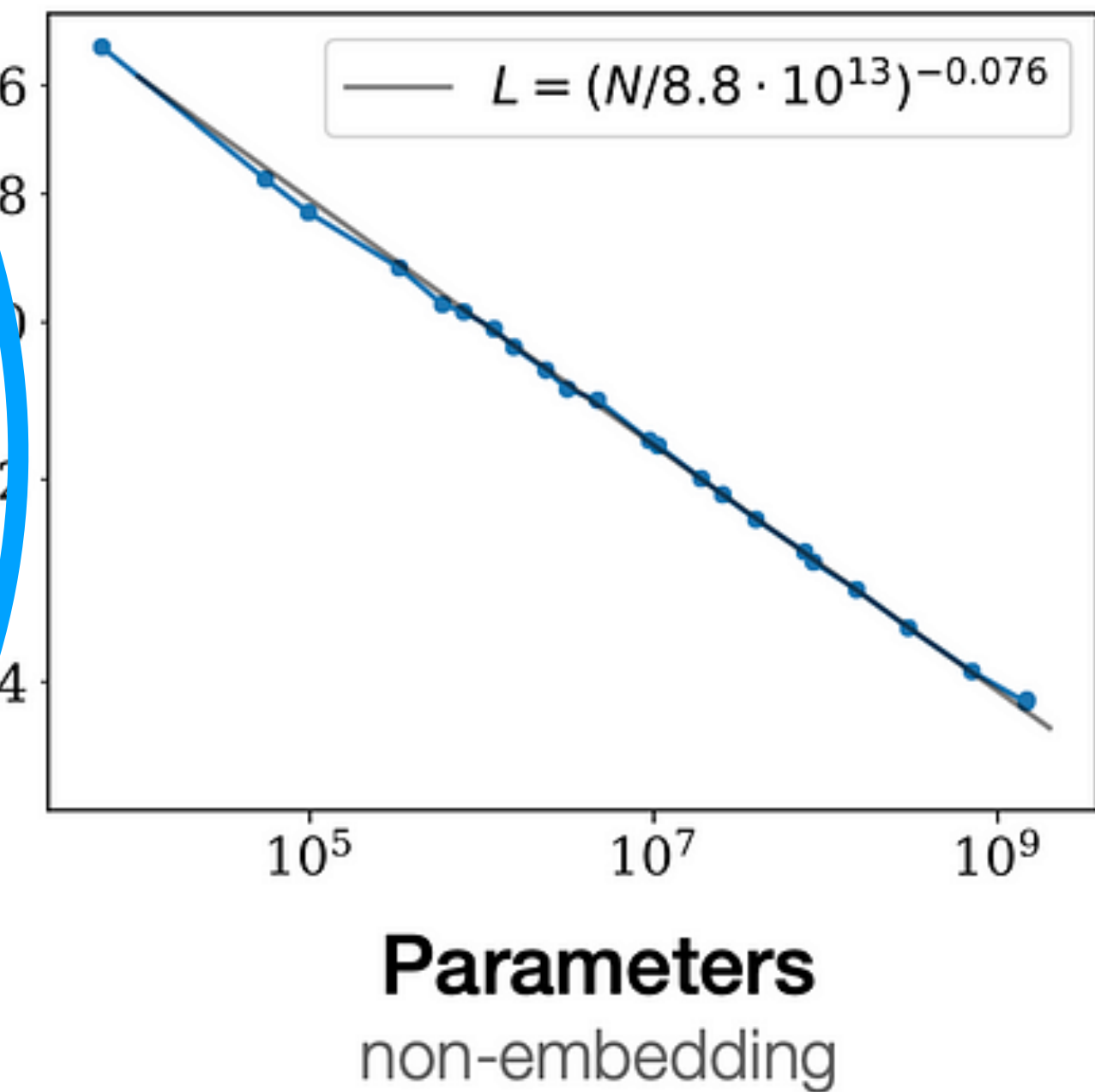
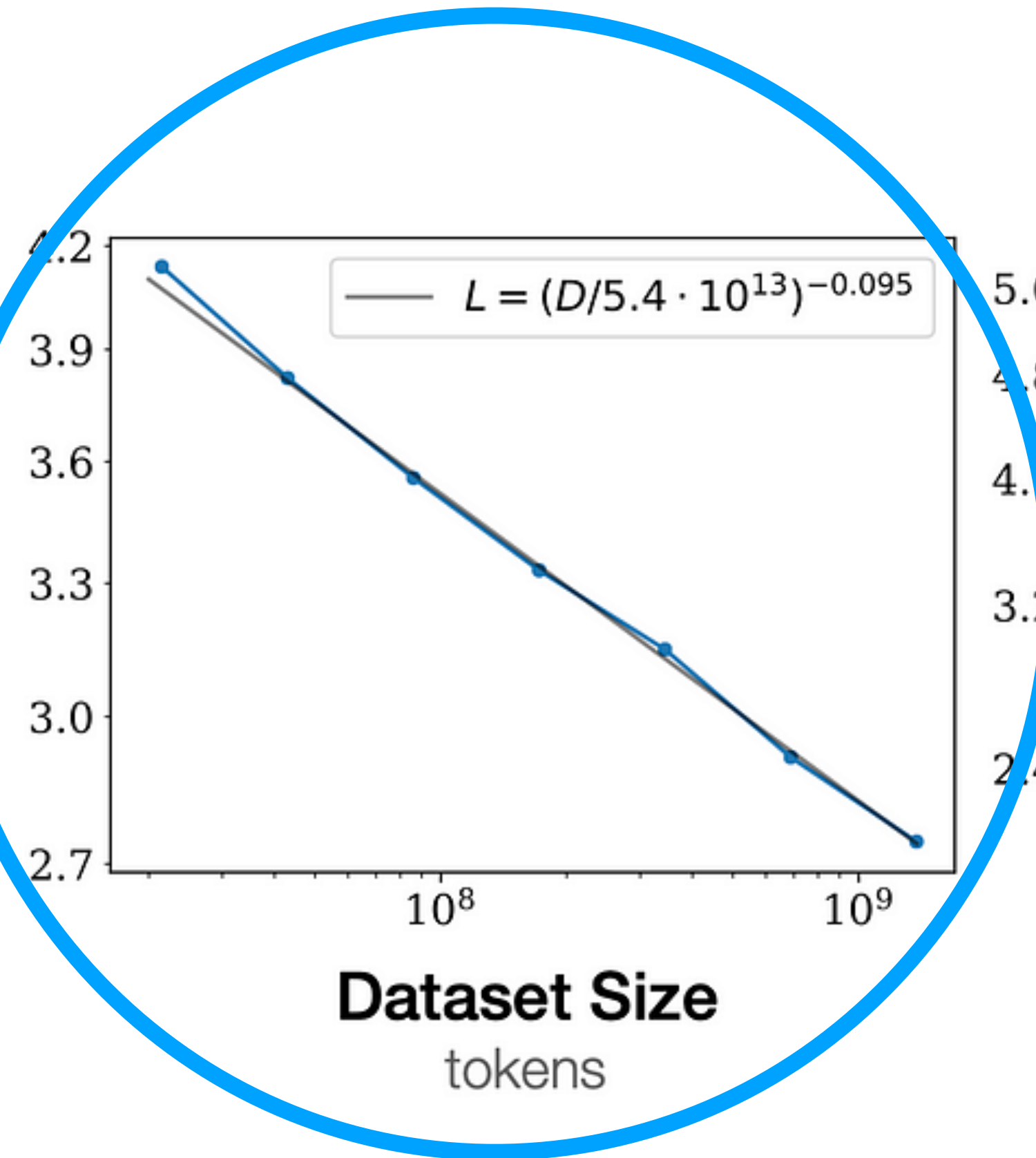
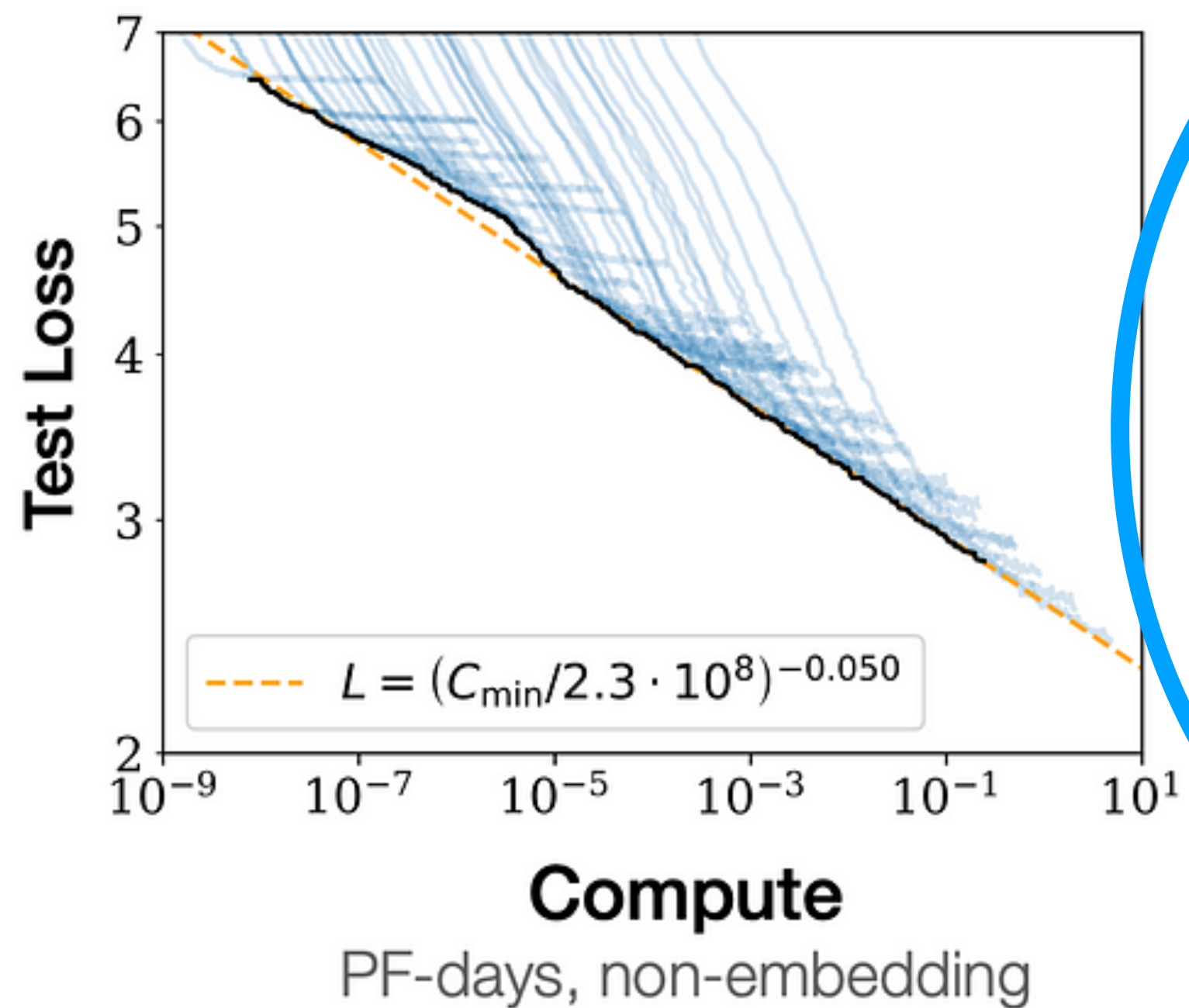
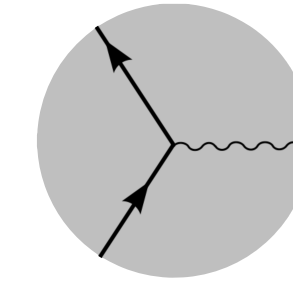
Uncertainty Quantification: Ensemble Uncertainty



Our Goal: Compute the Variance of an Ensembles Prediction *without training an ensemble*



Neural Network Scaling Laws

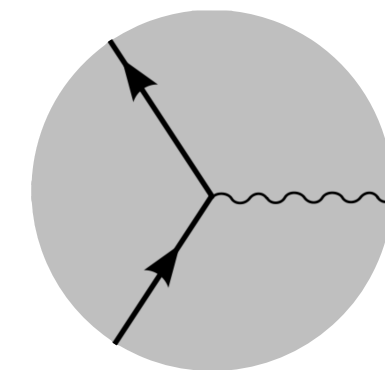


Neural Network Scaling Laws

Question: How does Ensemble Variance **scale** with Training Dataset Size? Can we predict this scaling with **physics-inspired theory**?



Physics-Inspired Theory: The NTK



NTK = Neural Tangent Kernel

$$\theta_\mu(t+1) = \theta_\mu(t) - \eta \left. \frac{\partial \mathcal{L}_A}{\partial \theta_\mu} \right|_{\theta_\mu = \theta_\mu(t)} \xrightarrow{\text{Taylor Expansion}} \Delta \mathcal{L}_A = -\eta \sum_{\mu, \nu} \lambda_{\mu\nu} \frac{\partial \mathcal{L}_A}{\partial \theta_\mu} \frac{\partial \mathcal{L}_A}{\partial \theta_\nu} + O(\eta^2),$$

The Error Factor

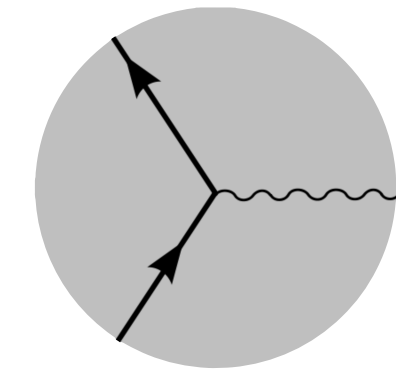
The NTK!

Chain Rule

$$\Delta \mathcal{L}_A = -\eta \sum_{i_1, i_2=1}^{n^{(L)}} \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2 \in \mathcal{A}} \left[\frac{\partial \mathcal{L}_A}{\partial z_{i_1; \tilde{\alpha}_1}} \frac{\partial \mathcal{L}_A}{\partial z_{i_2; \tilde{\alpha}_2}} \right] \left[\sum_{\mu, \nu} \lambda_{\mu\nu} \frac{dz_{i_1; \tilde{\alpha}_1}}{d\theta_\mu} \frac{dz_{i_2; \tilde{\alpha}_2}}{d\theta_\nu} \right].$$

$$\frac{\partial \mathcal{L}_A}{\partial z_{i; \tilde{\alpha}}} = z_i(x_{\tilde{\alpha}}; \theta) - y_{i; \tilde{\alpha}}$$

Physics-Inspired Theory: The NTK



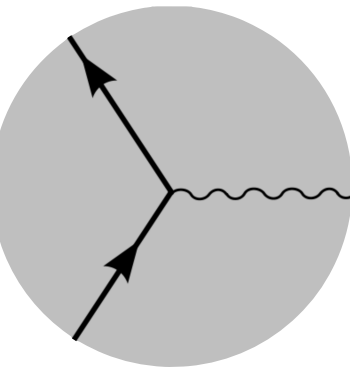
$$\Theta_{i_1 i_2; \tilde{\alpha}_1 \tilde{\alpha}_2} \equiv \sum_{\mu, \nu} \lambda_{\mu\nu} \frac{dz_{i_1; \tilde{\alpha}_1}}{d\theta_\mu} \frac{dz_{i_2; \tilde{\alpha}_2}}{d\theta_\nu}$$

The NTK is the **main driver** of the **function-approximation** dynamics.*

Thus, by understanding the NTK and how it evolves we can directly understand the behavior of NNs

* for DNNs trained with full batch gradient descent

Physics-Inspired Theory: The NTK Perturbation



NTK dynamics are possible to understand **perturbatively and analytically**

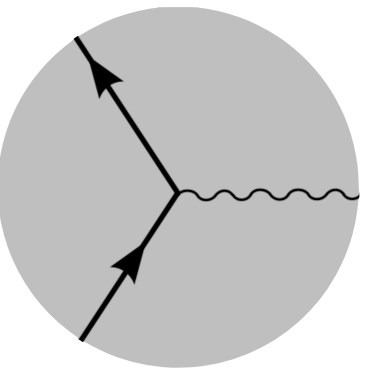
$$\Theta_{\alpha_1 \alpha_2}^{(\ell)} = \Theta_{\alpha_1 \alpha_2}^{\{0\}(\ell)} + \frac{1}{n_\ell} \Theta_{\alpha_1 \alpha_2}^{\{1\}(\ell)} + \frac{1}{n_\ell^2} \Theta_{\alpha_1 \alpha_2}^{\{2\}(\ell)} + O\left(\frac{1}{n_\ell^3}\right)$$

↑
Infinte Width

↑
Finite Width Corrections

Perturbing in the **width of the network**

Physics-Inspired Theory: RG-Flow of the NTK




How to compute the infinite width NTK? \rightarrow RG Flow

Recursively Repeat through each layer until **final layer** action is computed! This is **RG Flow in an DNN**.

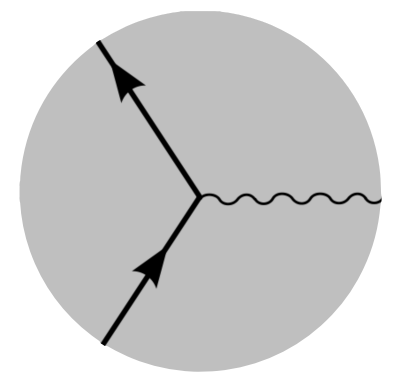
$$\Theta_{\alpha_1 \alpha_2}^{\{0\}(\ell)} \longrightarrow e^{-S_{eff}(z^{(\ell)})}$$


$$S_{full}(z^{(1)}, \dots, z^{(L)}) = \sum_{\ell=1}^L S_M(z^{(\ell)}) + \sum_{\ell=1}^{L-1} S_I(z^{(\ell+1)} | z^{(\ell)})$$



Compute **effective action of a layer ℓ** by **marginalizing** over all pre-activations in layer $\ell - 1$

$$e^{-S_{eff}(z^{(L)})} = \int \left[\prod_{\ell=1}^{L-1} dz^{(\ell)} \right] e^{-S_{full}(z^{(1)}, \dots, z^{(L)})}$$



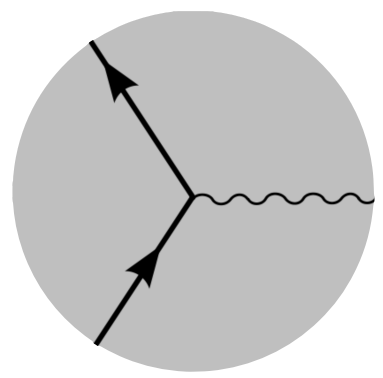
Physics-Inspired Theory: Infinite Width Predictions

With the **Final Layer NTK**, we can compute the **“end of training”** prediction:

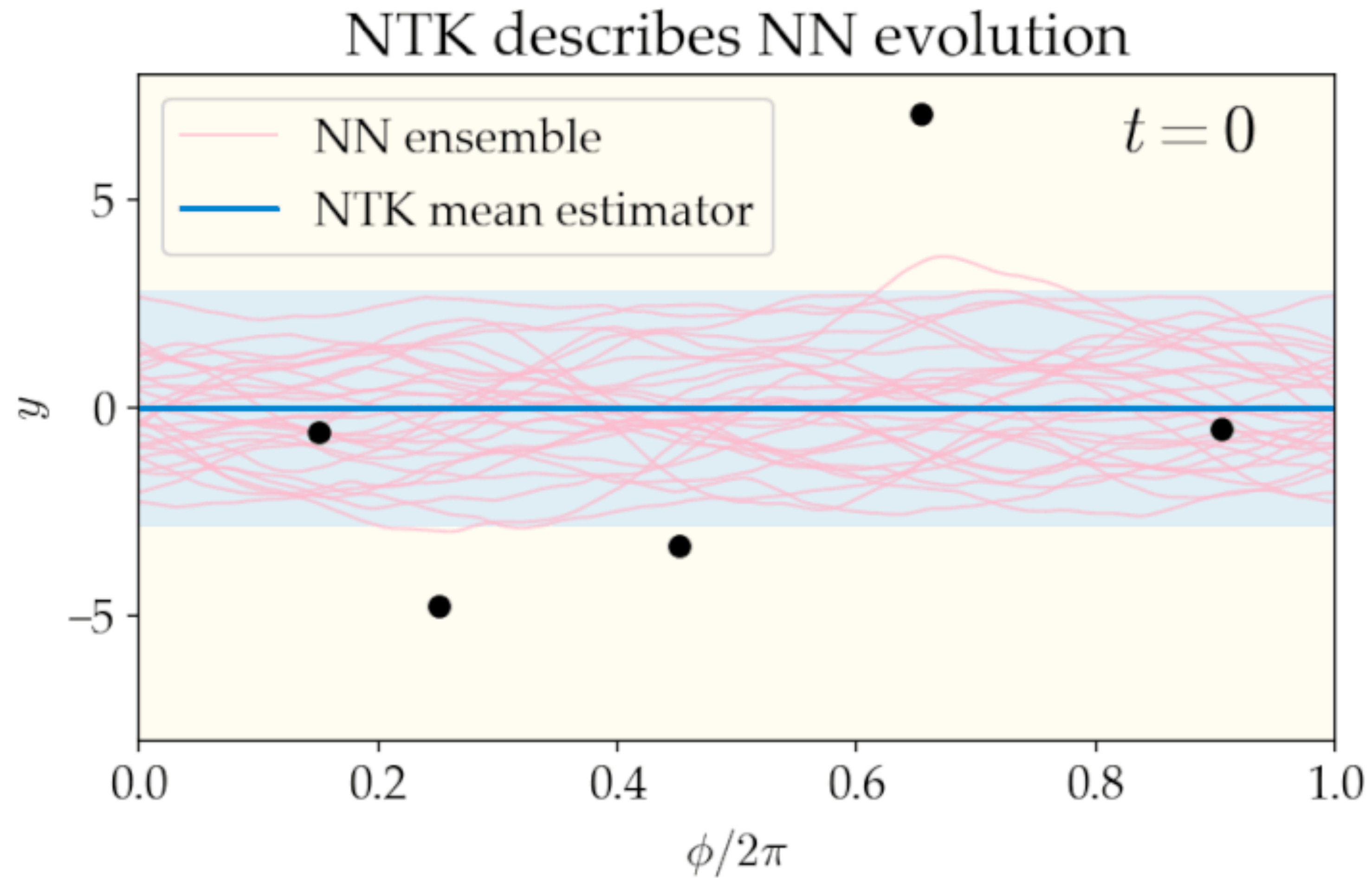
$$m_{i;\beta}^{\infty} \equiv \mathbb{E} \left[z_{i;\beta}^{(L)}(T) \right] = \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2 \in \mathcal{A}} \Theta_{\beta \tilde{\alpha}_1}^{(L)} \left(\Theta_{\tilde{\alpha}_1 \tilde{\alpha}_2}^{(L)} \right)^{-1} y_{i;\tilde{\alpha}_2}$$

And the **variance on that prediction**:

$$\begin{aligned} \text{Cov} \left[z_{i_1;\beta_1}^{(L)}(T), z_{i_2;\beta_2}^{(L)}(T) \right] &= \mathbb{E} \left[z_{i_1;\beta_1}^{(L)}(T) z_{i_2;\beta_2}^{(L)}(T) \right] - m_{i_1;\beta_1}^{\infty} m_{i_2;\beta_2}^{\infty} \\ &= \delta_{i_1 i_2} K_{\beta_1 \beta_2}^{(L)} - \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2 \in \mathcal{A}} \Theta_{\beta_1 \tilde{\alpha}_1}^{(L)} K_{\tilde{\alpha}_1 \tilde{\alpha}_2}^{(L)} \Theta_{\tilde{\alpha}_2 \beta_2}^{(L)} - \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2 \in \mathcal{A}} \Theta_{\beta_2 \tilde{\alpha}_1}^{(L)} K_{\tilde{\alpha}_1 \tilde{\alpha}_2}^{(L)} \Theta_{\tilde{\alpha}_2 \beta_1}^{(L)} + \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4 \in \mathcal{A}} \Theta_{\beta_1 \tilde{\alpha}_1}^{(L)} \Theta_{\tilde{\alpha}_1 \tilde{\alpha}_2}^{(L)} K_{\tilde{\alpha}_2 \tilde{\alpha}_3}^{(L)} \Theta_{\tilde{\alpha}_3 \tilde{\alpha}_4}^{(L)} \Theta_{\tilde{\alpha}_4 \beta_2}^{(L)}. \end{aligned}$$

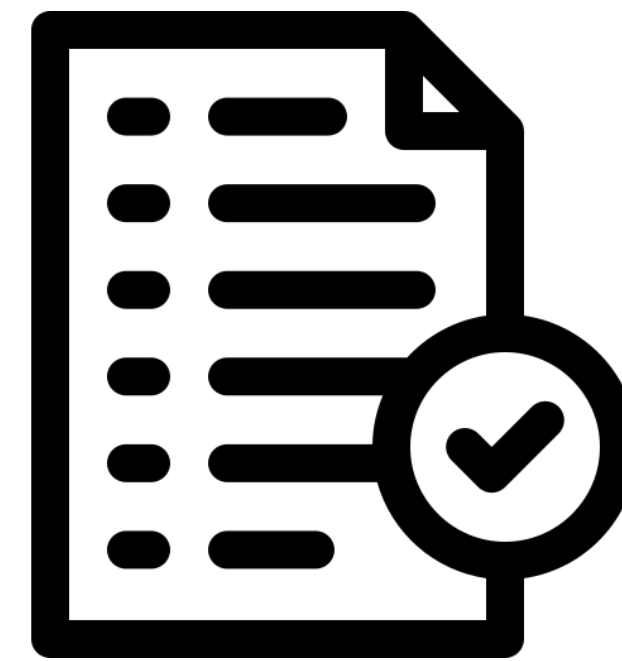


Physics-Inspired Theory: Infinite Width Predictions





Empirical Results



Empirical Results Comparing Theoretical
Results with Empirical Reality

04

I Empirical Results: Infinite Width Predictions



Are infinite width calculations **predictive** on **real Machine Learning** problems?

We test this by computing:

- 1) Infinite Width Prediction
- 2) A Trained Ensemble of DNNs (~150 networks)
 - Width-30, Early Stopping, Full Batch Gradient Descent

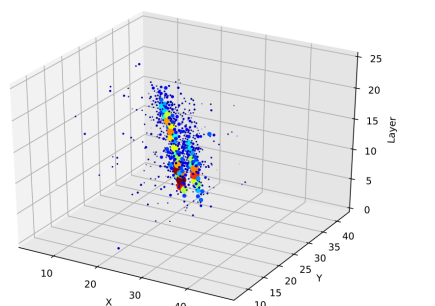
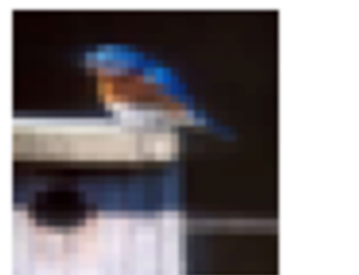
For a range of training set sizes

On three datasets:

- 1) MNIST Image Classification
- 2) CIFAR Image Classification
- 3) A HEP Calorimeter Energy Regression Problem

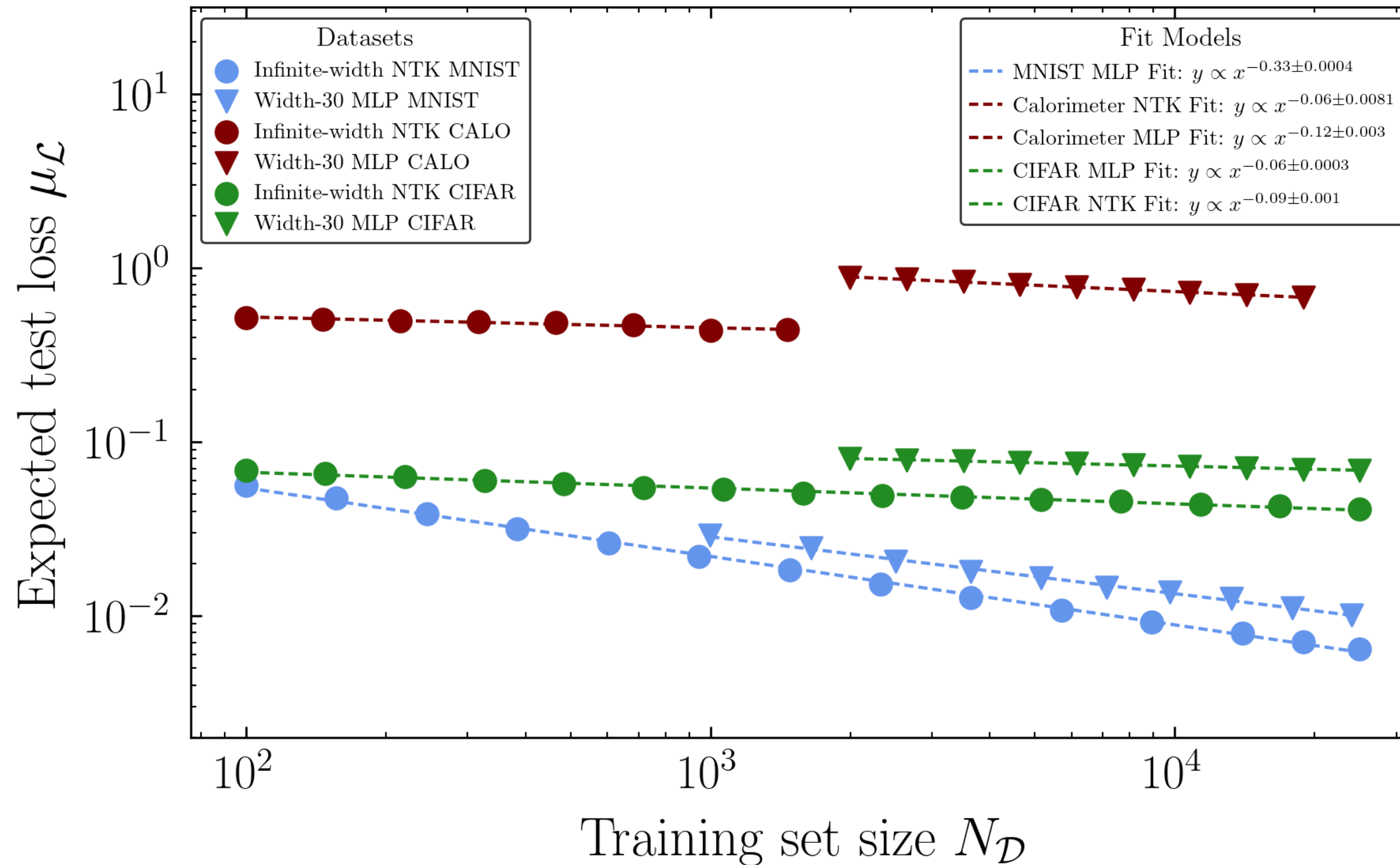


2: bird



Empirical Results: Infinite Width Prediction Loss

MNIST, Calorimeter, and CIFAR, depth 3, $\lambda_b/\lambda_W = 10$



The **mean test loss** for a trained ensemble of DNNs and Infinite Width Networks
for three datasets

Empirical Results: Infinite Width Coefficient of Variation

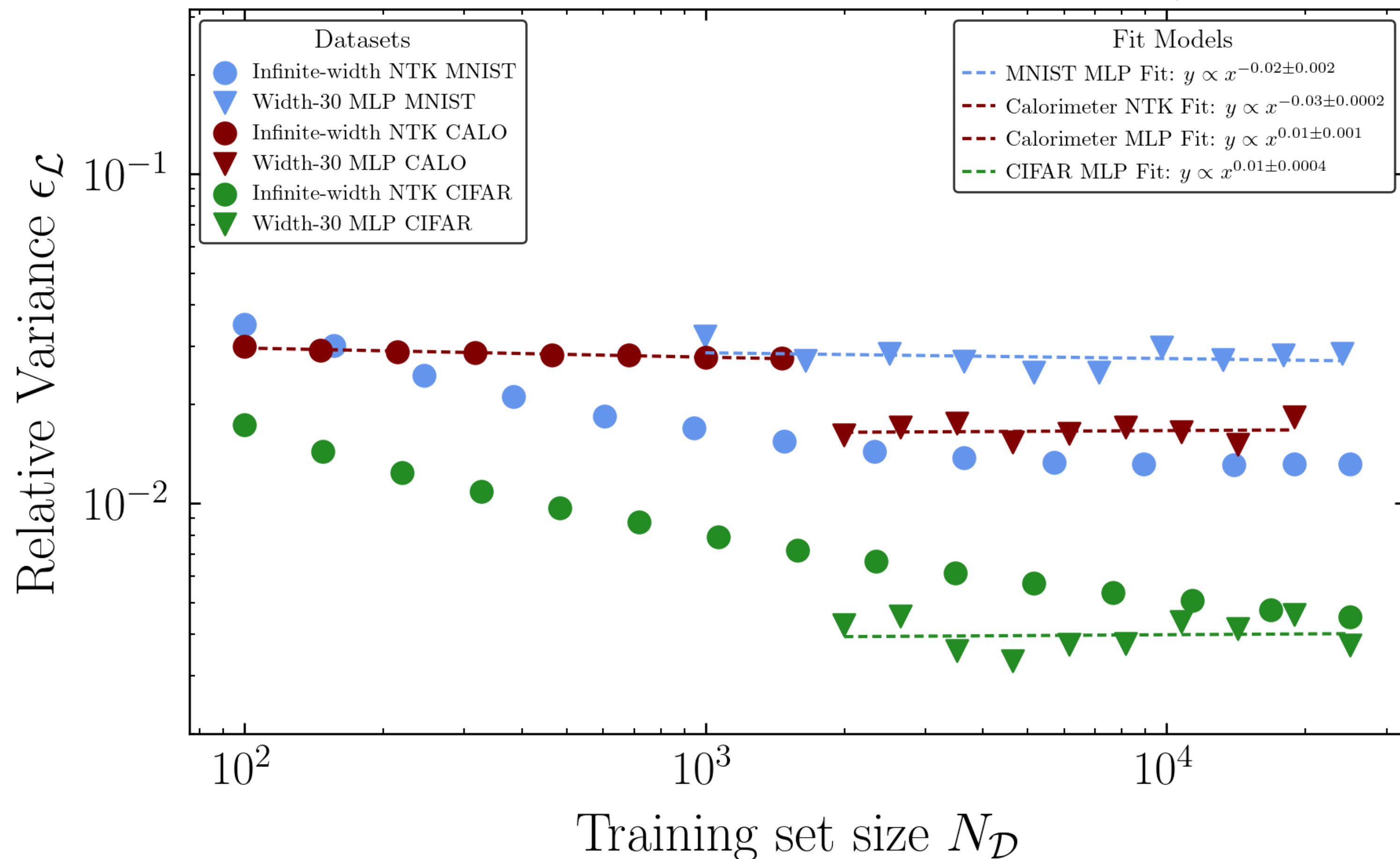
Suppose we introduce the **Coefficient of Variation:**

$$\epsilon_{\mathcal{L}} \equiv \frac{\sigma_{\mathcal{L}}}{\mu_{\mathcal{L}}}$$

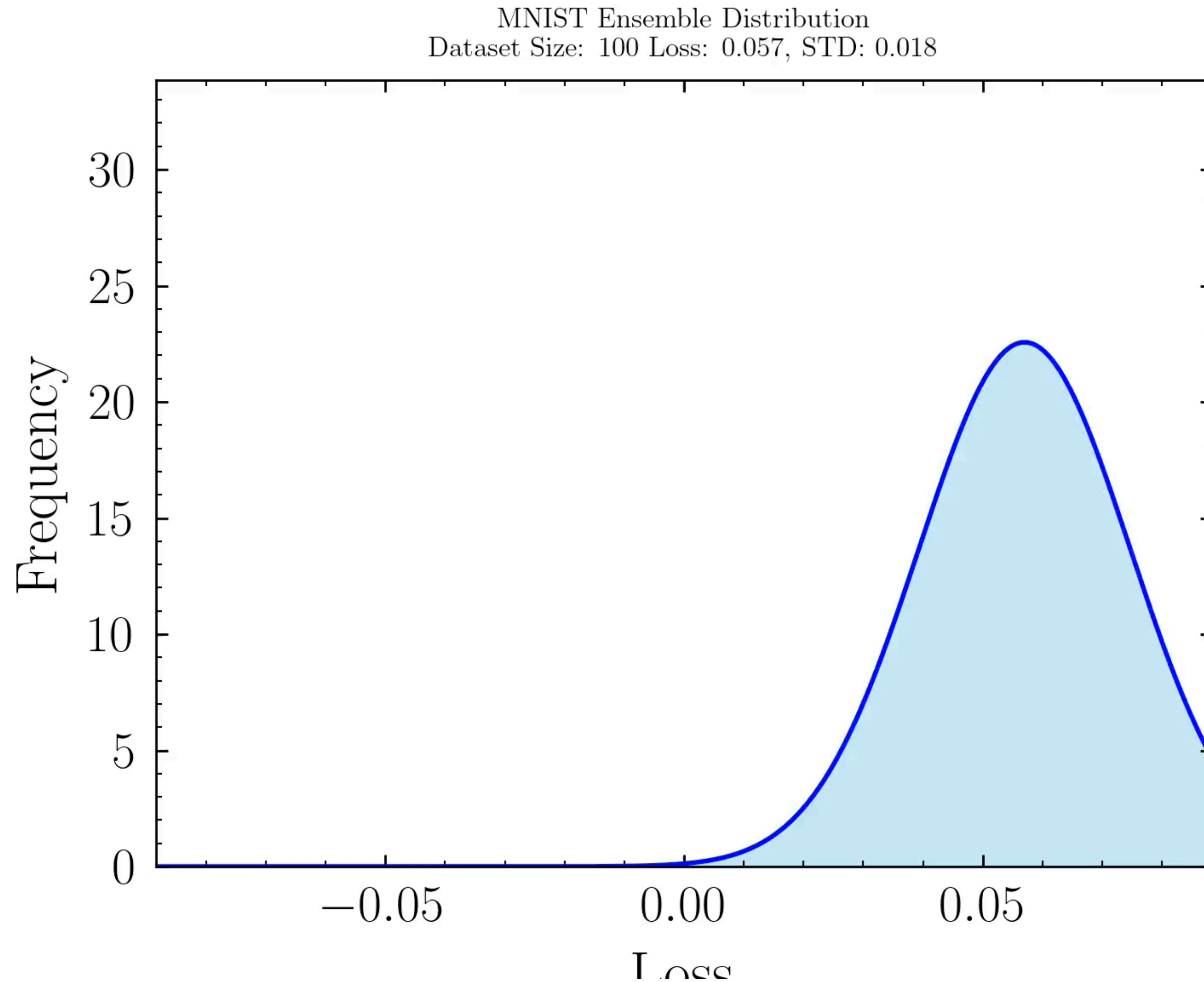
We find:

- 1) DNN $\epsilon_{\mathcal{L}}$ flat with dataset size!
- 2) Infinite width $\epsilon_{\mathcal{L}}$ asymptotes flat.

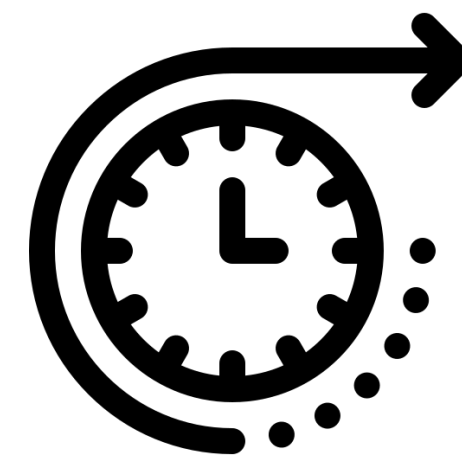
MNIST, Calorimeter, and CIFAR, depth 3, $\lambda_b/\lambda_W = 10$



Empirical Results: Infinite Width Prediction Loss



Conclusion: Implications



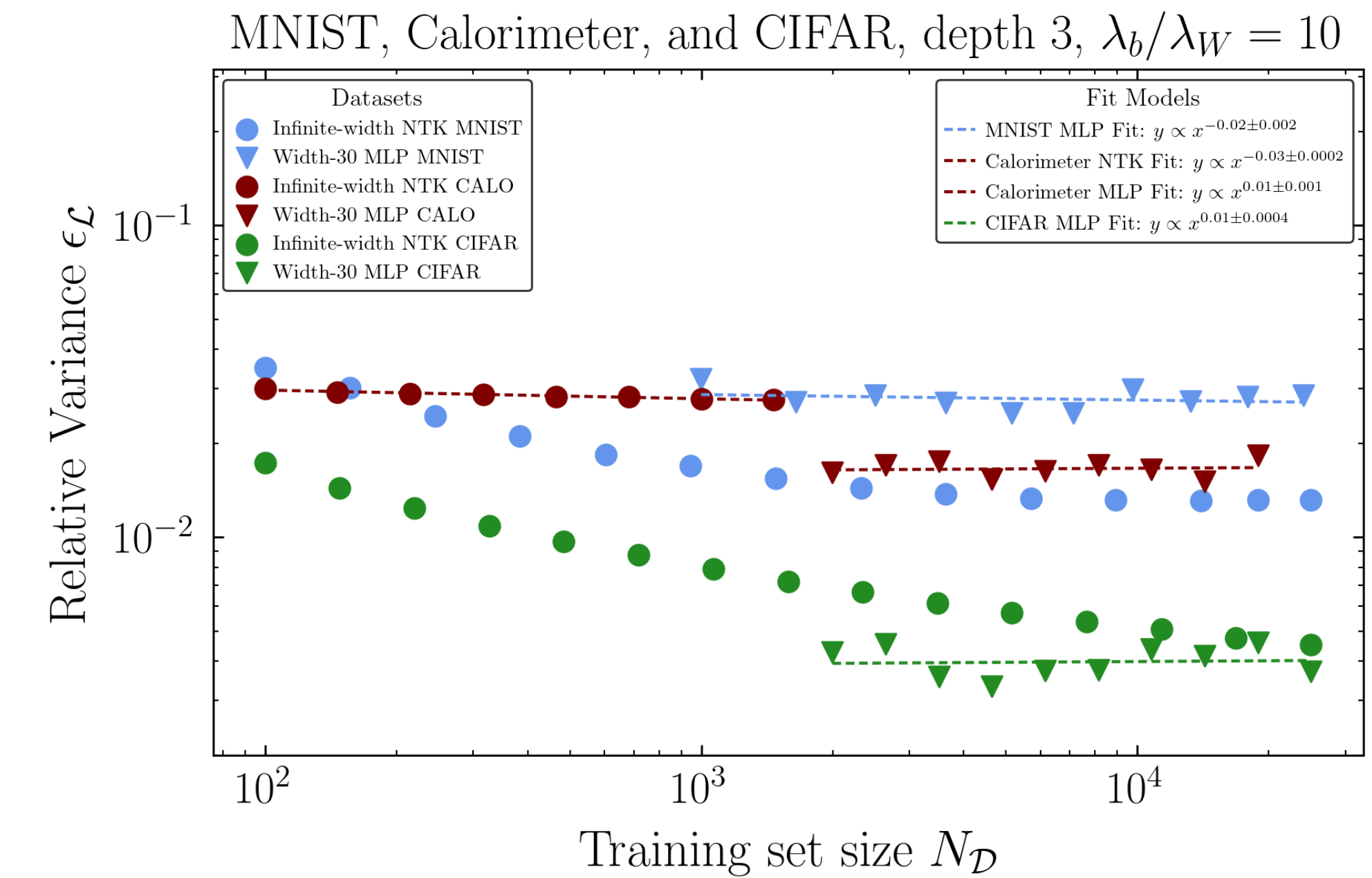
Implications of our work:

1. We **find that** $\epsilon_{\mathcal{L}}$ **is small** ($\mu_{\mathcal{L}} > \sigma_{\mathcal{L}}$)

- We can assign $\mu_{\mathcal{L}}$ as the systematic uncertainty due to the DNN

2. $\epsilon_{\mathcal{L}}$ **is flat and similar to the infinite width value**, thus one can estimate $\epsilon_{\mathcal{L}}$ by either:

- **Training an ensemble for small** $N_{\mathcal{D}}$ (cheap) and extrapolate $\epsilon_{\mathcal{L}}$ value to larger $N_{\mathcal{D}}$
- Compute Infinite Width Value after $N_{\mathcal{D}}$ asymptotes (very cheap)



QUESTIONS?

