



UNIVERSITÀ DEGLI STUDI
DI GENOVA



Machine Learning LHC Likelihoods.

Humberto Reyes-González
University of Genoa

LHC reinterpretation workshop
CERN, 14/12/2022

Introduction

- **Likelihood functions (full statistical models)** parametrise the full information of an LHC analysis; whether it is New Physics (NP) search or an SM measurement.
- Their **preservation** is a key part of the **LHC legacy**.

Usage:

- Resampling
- Reinterpretation with different statistical approaches.
- Reinterpretation in the context of different NP models.
- ...

Challenges:

- LHC likelihoods are often high-dimensional complex distributions.
- We want precise descriptions that can be efficiently reinterpreted.

Important steps forward:

- ATLAS started publishing full likelihoods of NP searches [ATL-PHYS-PUB-2019-029](#).
- Release of the pyhf package to construct statistical models [10.21105/joss.02823](#), L Heinrich, M Feickert, G Stark
- Theorists have started profiting from this [arXiv:2009.01809](#), [arXiv:2012.08192](#), SModelS collaboration
- Supervised learning with DNN likelihood [arxiv:1911.03305 A](#) Cocco, M. Perini, L Silvestrini, R Torre

Our approach:

Unsupervised Learning with Normalizing Flows

LHC likelihoods in a nutshell

Bayes theorem:

$$P(\Theta, x) = P_x(x | \Theta) \pi_{\Theta}(\Theta) = P_{\Theta}(\Theta | x) \pi_x(x)$$

LHC Statistical model:

$$P(\mu, \theta; \text{data}) = \prod_{k=1}^{n_c} P[n_i; \mu \epsilon_{i,k}(\vec{\theta}) N_{S,i,k}(\vec{\theta}) + B_{i.k}(\vec{\theta})] \prod_{j=1}^{n_{\text{syst}}} G(\theta_j^{\text{obs}}; \theta_j; 1)$$

Diagram illustrating the LHC Statistical model equation with color-coded boxes and labels:

- Blue boxes:** μ and $\epsilon_{i,k}$. Label: Parameters of Interest (signal strength, observables, etc.)
- Purple boxes:** θ and $\vec{\theta}$. Label: Nuisance parameters (uncertainties)
- Green boxes:** n_i and θ_j^{obs} . Label: (Observed) data
- Green box:** θ_j . Label: (Auxiliary) data

Test Statistic:

$$t(\mu) = -2 \log \frac{L(\mu; \hat{\theta}(\mu))}{L(\hat{\mu}, \hat{\theta}(\hat{\mu}))}$$

Best-fits:

$$L(\hat{\mu}; \hat{\theta})$$

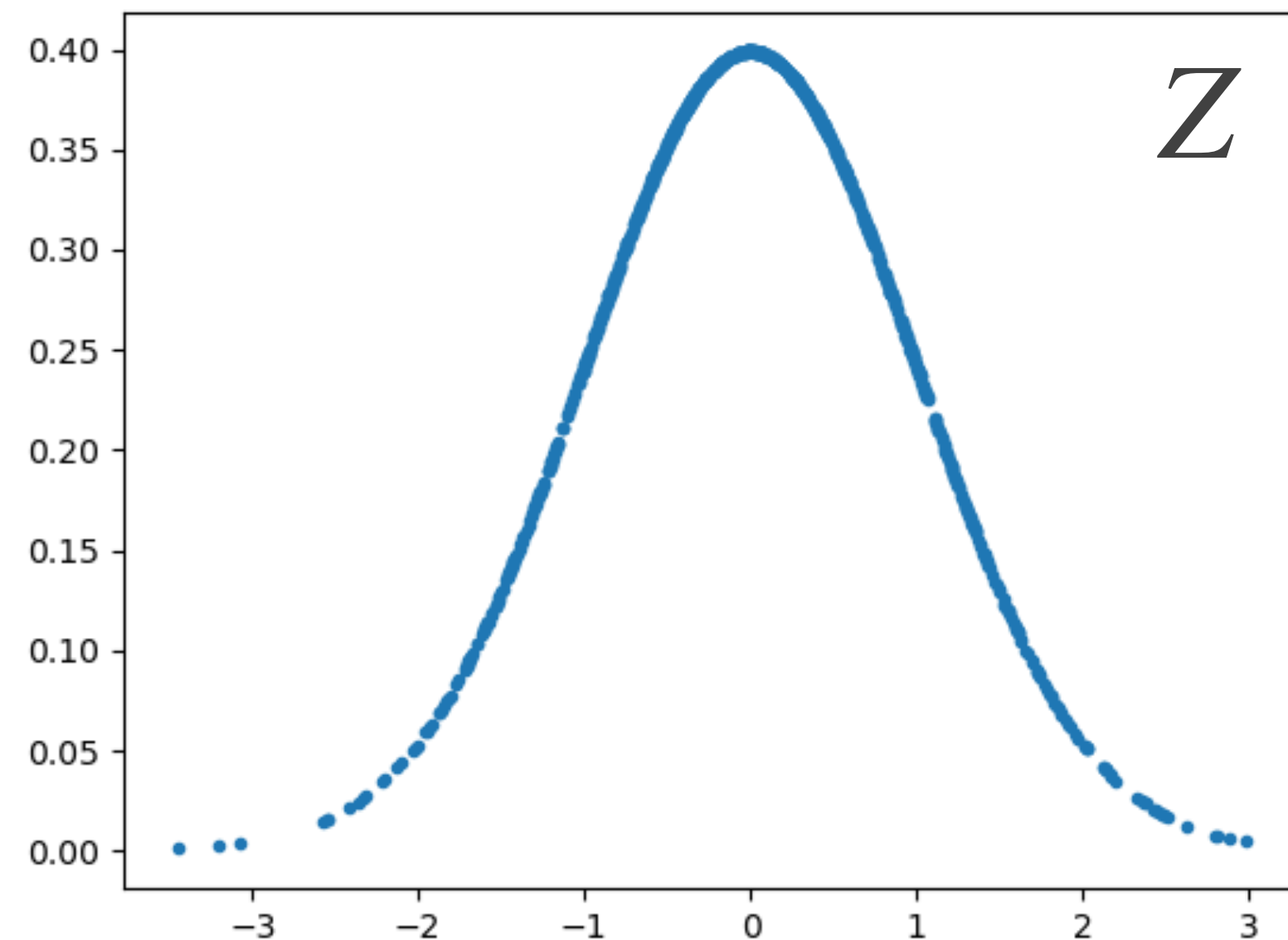
Where μ are observables

best-fit $\theta(\mu)$

Introduction.

BASIC PRINCIPLE:

Following the change of variables formula, perform a series of **bijective, continuous, invertible** transformations on a *simple* probability density function (pdf) to obtain a *complex* one.

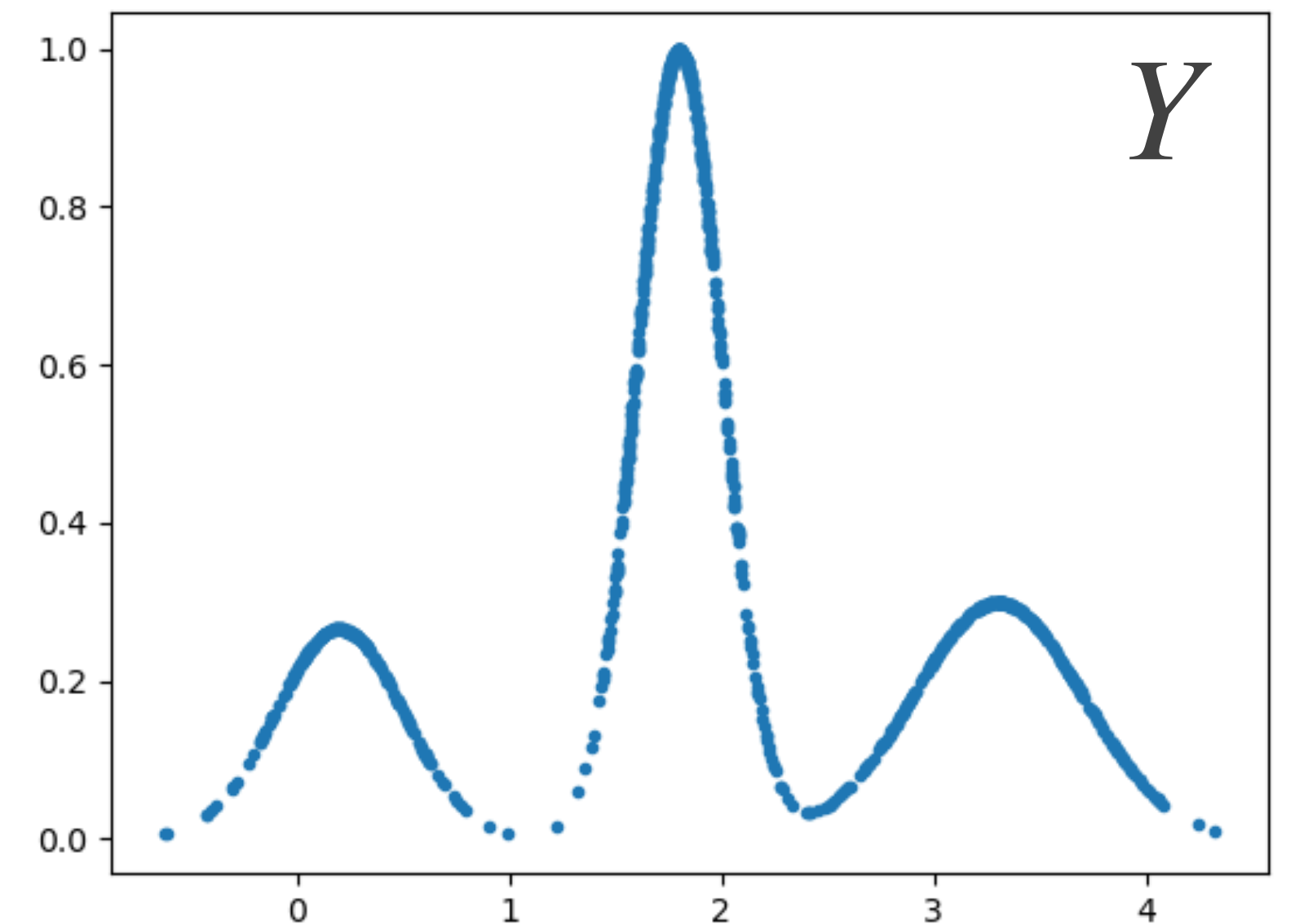


Normalizing direction

$$Z = f(Y)$$

Generative direction

$$Y = g(Z)$$



Choosing the transformations

THE OBJECTIVE:

To perform the right transformations to accurately estimate the complex underlying distribution of some observed data.

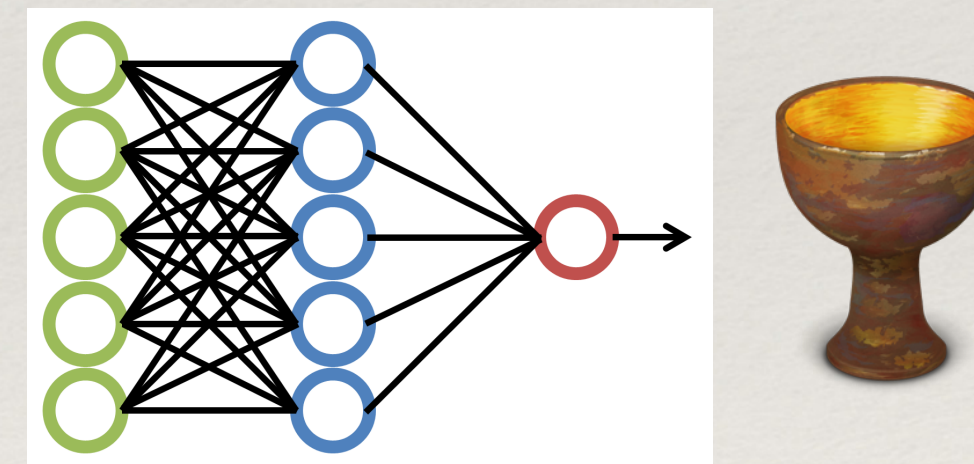
THE RULES OF THE GAME:

- The transformations (bijections) must be invertible
- They should be sufficiently expressive
- And computationally efficient (including Jacobian)

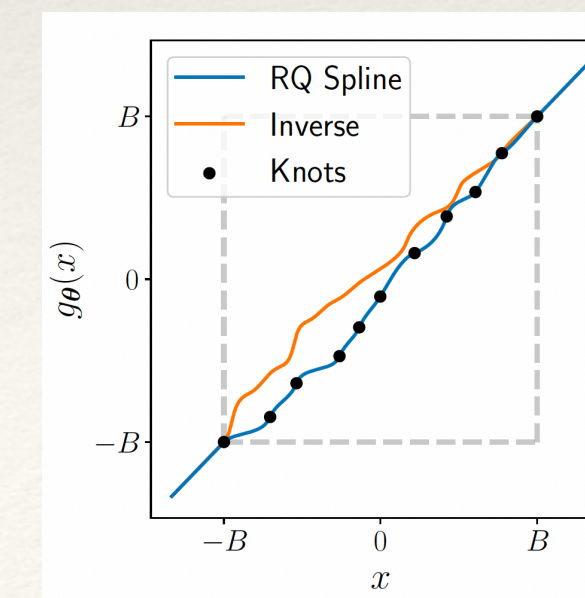


THE STRATEGY

Let *Neural Networks* learn the parameters of *Autoregressive Normalizing Flows*.



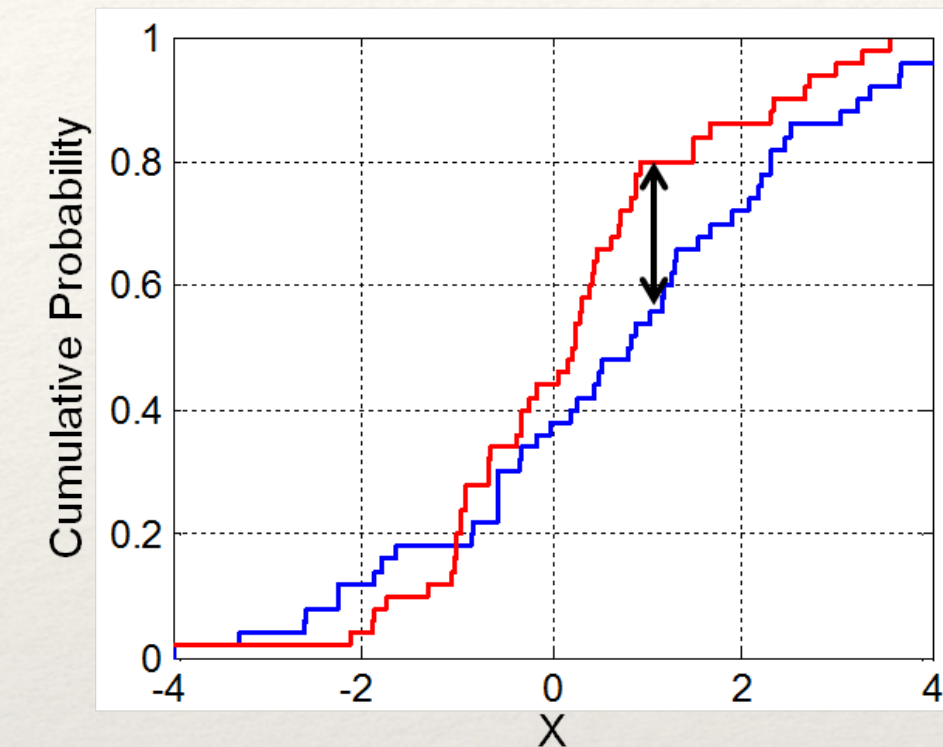
Autoregressive Rational-Quadratic-Spline Flows
(A-RQS, [arXiv:1906.04032](https://arxiv.org/abs/1906.04032))



Evaluation metric.

- Two-sample 1D Kolgomonov - Smirnov test (ks test):

$$D_{n,m} = \sup_x |F_n(x) - F_m(x)|$$



- Computes the p-value for two sets of 1D samples coming from the same *unknown* distribution.
- We average over ks test estimations and compute the median over dimensions.
- Optimal value 0.5

Example Likelihoods

$$P_{\Theta}(\Theta | x = \text{obs})$$

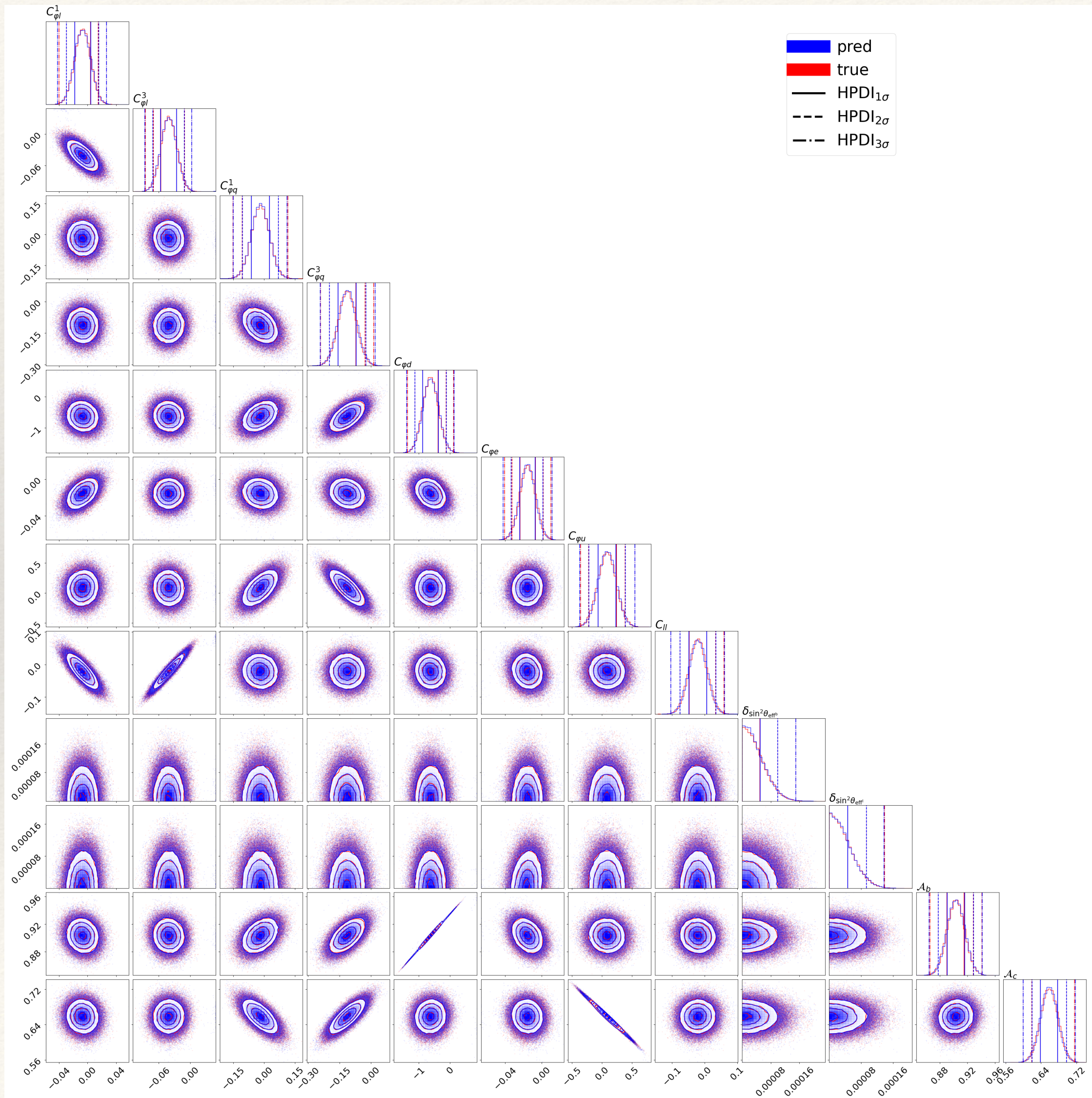
ElectroWeak fit Likelihood

- EW observables.
- Including recent measurements of top mass (CMS) and W mass (CDF).
- 8 parameters of interest (Wilson coefficients of SMEFT operators)
- 32 nuisance parameters.
- Ref. [arXiv:2204.04204](#)

Flavor fit likelihood.

- Flavor observables related to $b \rightarrow sl^+l^-$ transitions
- 12 parameters of interest (Wilson coefficients of SMEFT operators)
- 77 nuisance parameters.
- Ref. [arXiv:1903.09632](#)

ElectroWeak fit Likelihood



Hyperparameters:

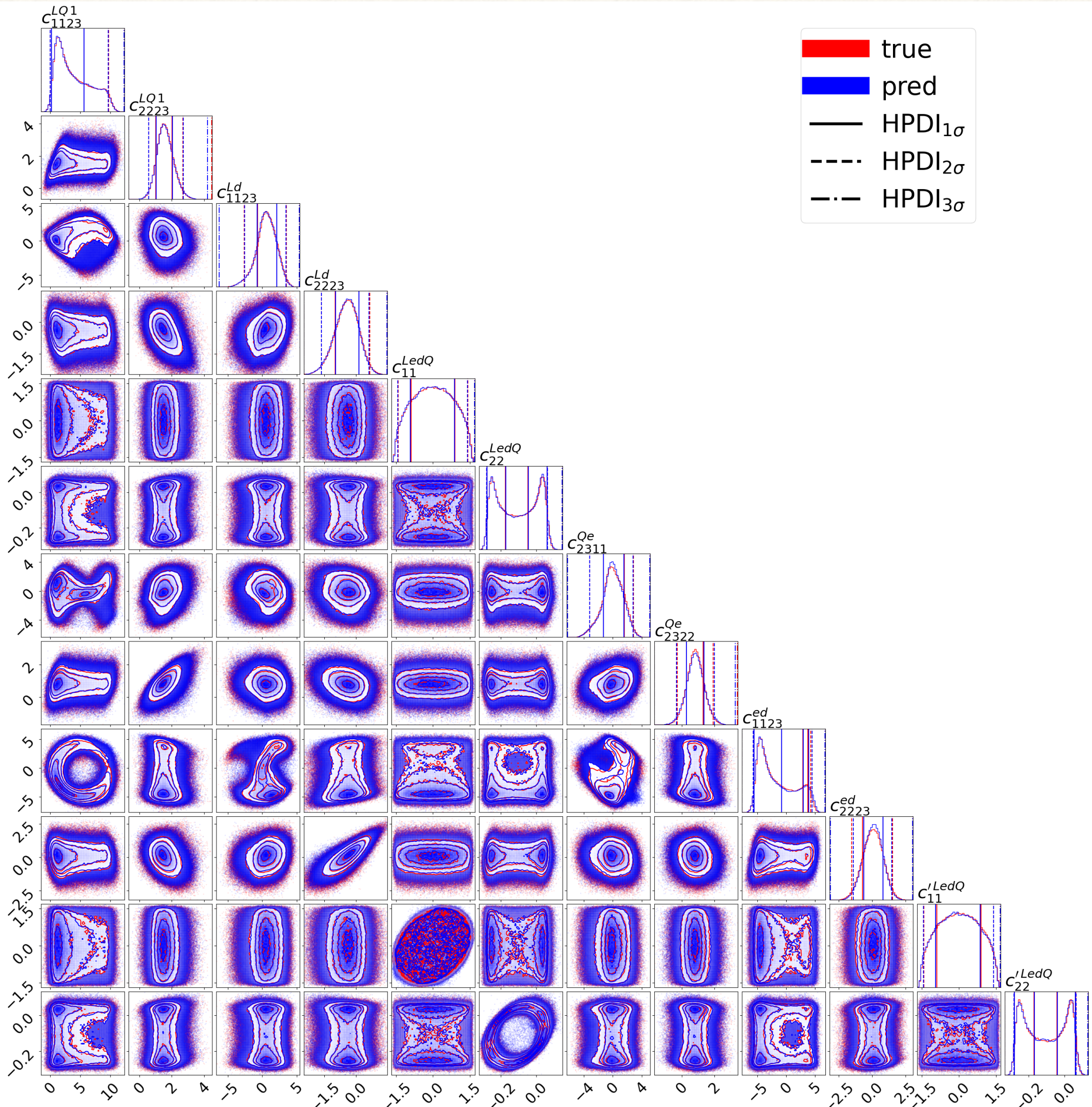
N_{train}	Flow	N bij	N knots	Range	Hidden layers	L1 factor	N epochs	N iters.
10^4	A-RQS	2	6	[-6,6]	128×3	0	200	12
10^5								
$2 \cdot 10^5$								

Results:

N_{train}	KS-test
10^4	0.453
10^5	0.4803
$2 \cdot 10^5$	0.486

Test sample : 300k

Flavor fit Likelihood



Hyperparameters:

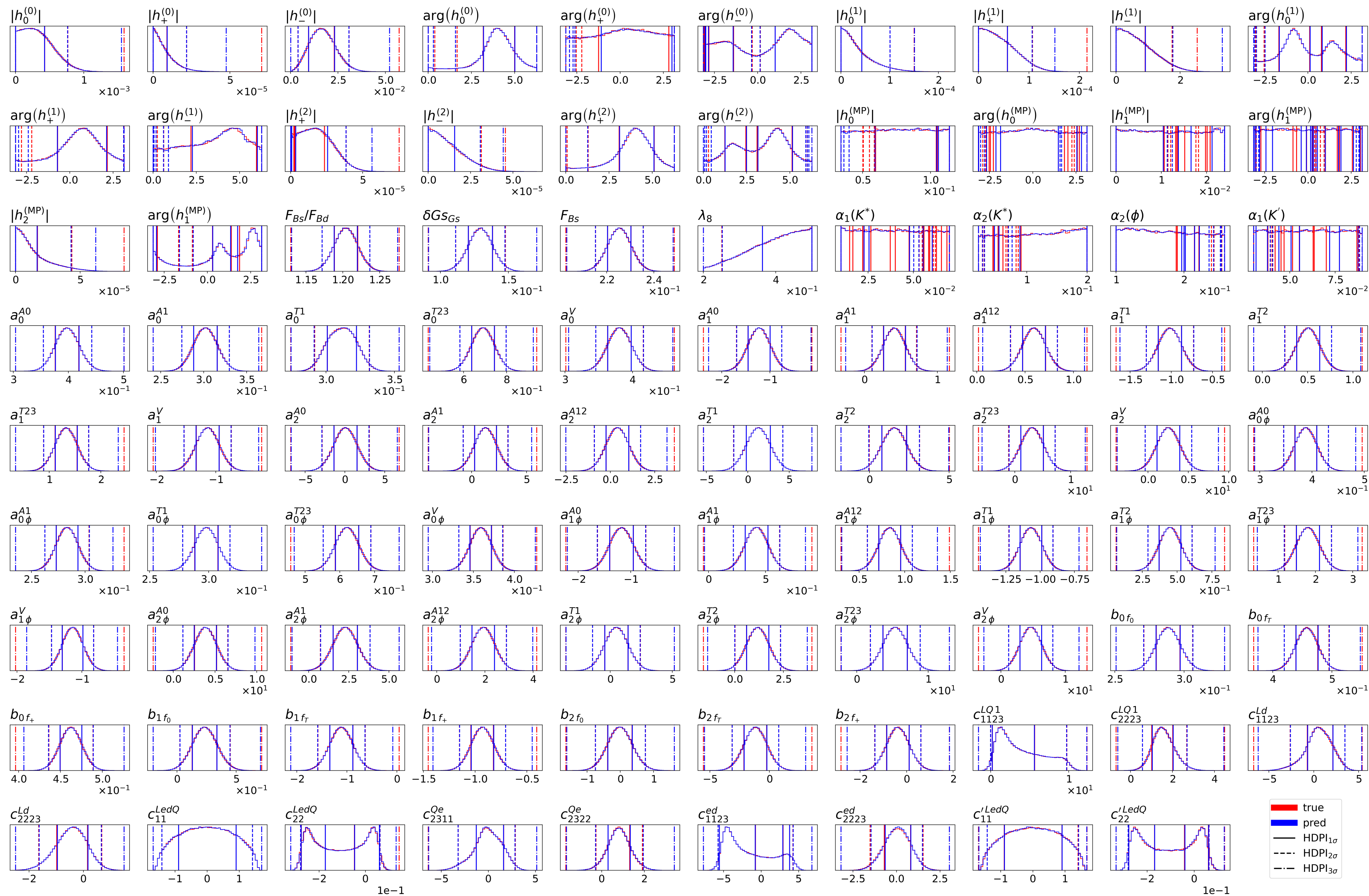
N_{train}	Flow	N bij	N knots	Range	Hidden layers	L1 factor	N epochs	N iters.
10^5	A-RQS	2	16	[-6,6]	1024×3	10^{-4}	1200	12
$5 \cdot 10^5$								
10^6								

Results:

N_{train}	KS-test
10^5	0.457
$5 \cdot 10^5$	0.482
10^6	.4806

Test sample : 500k

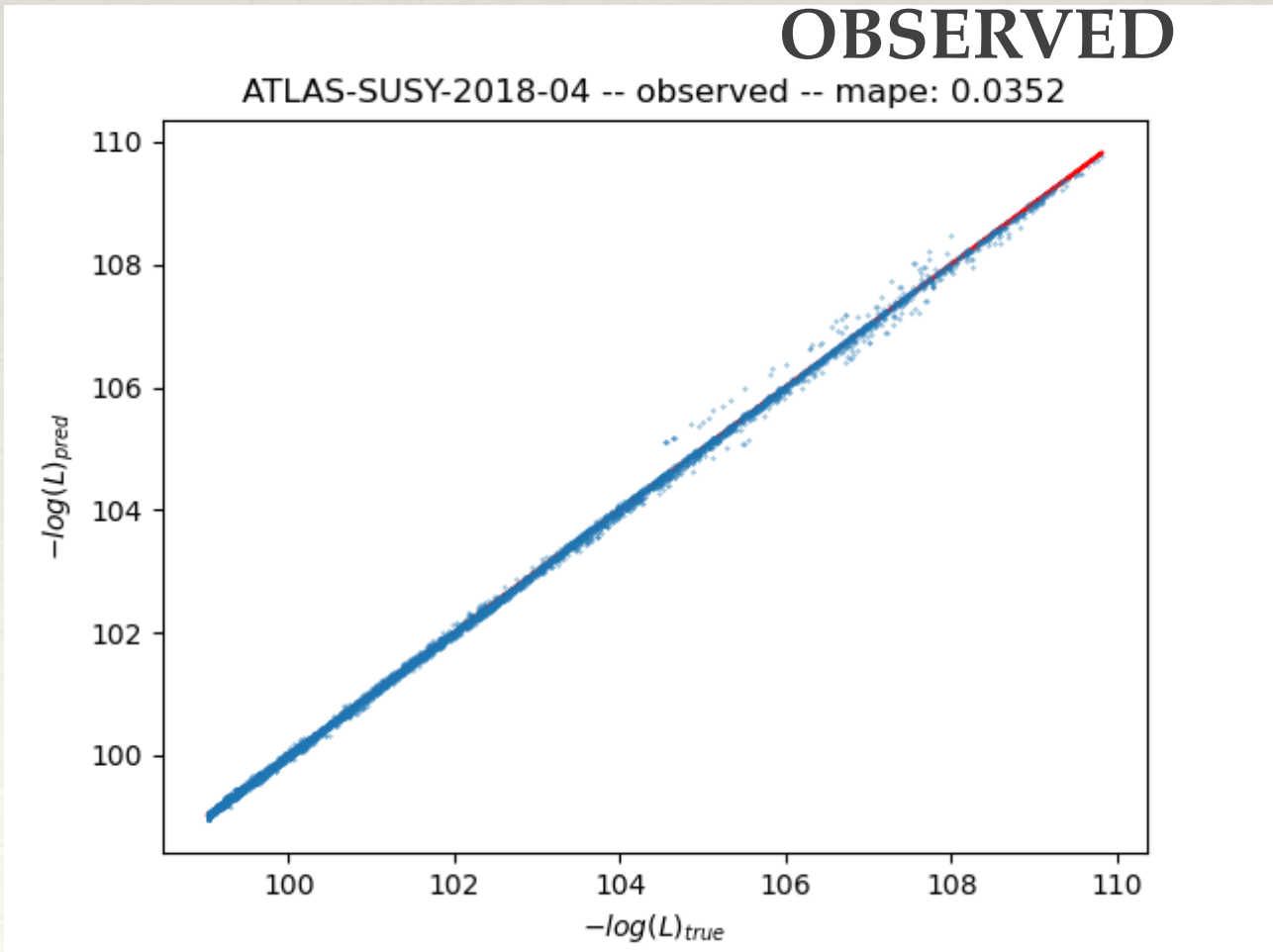
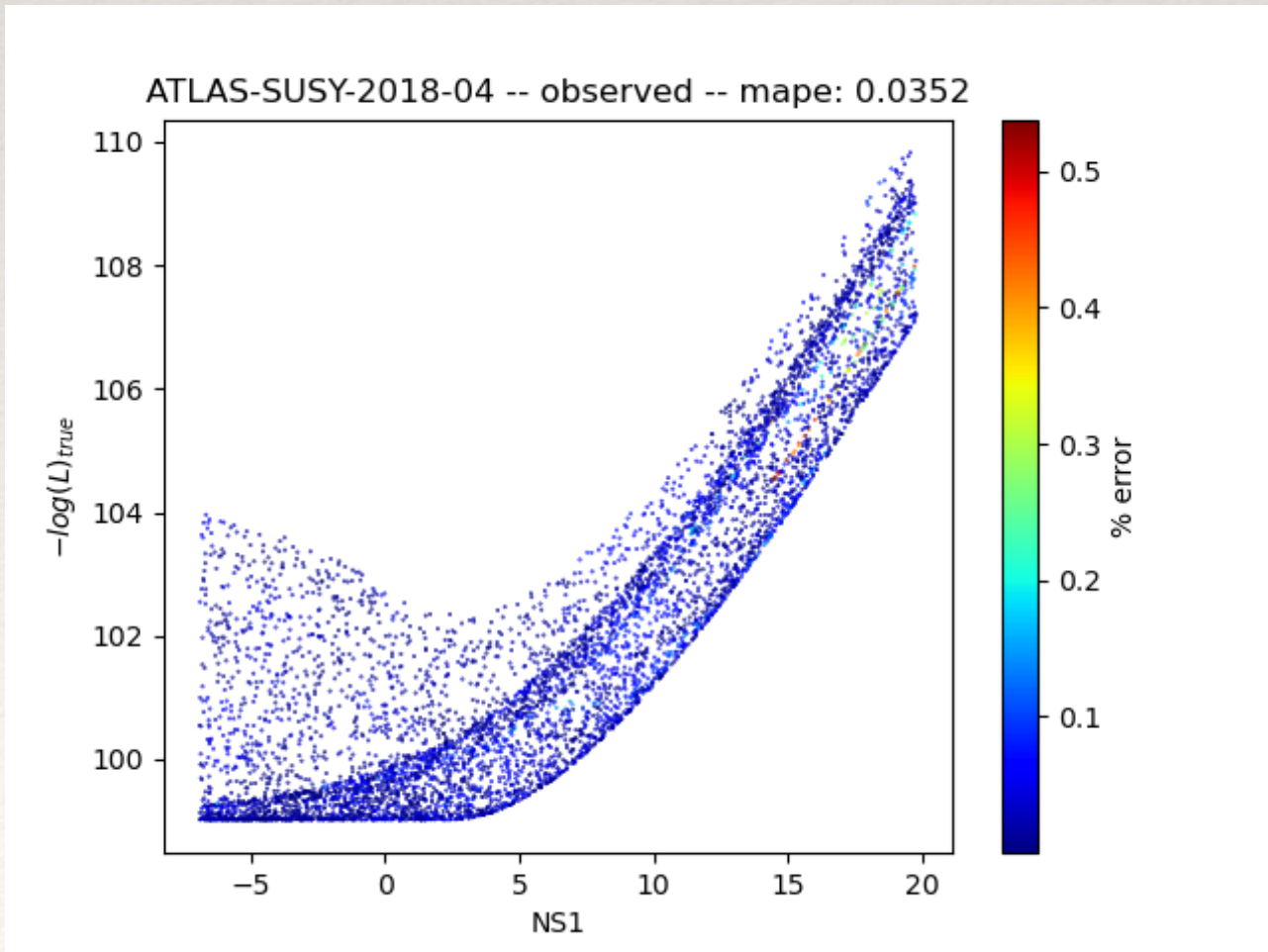
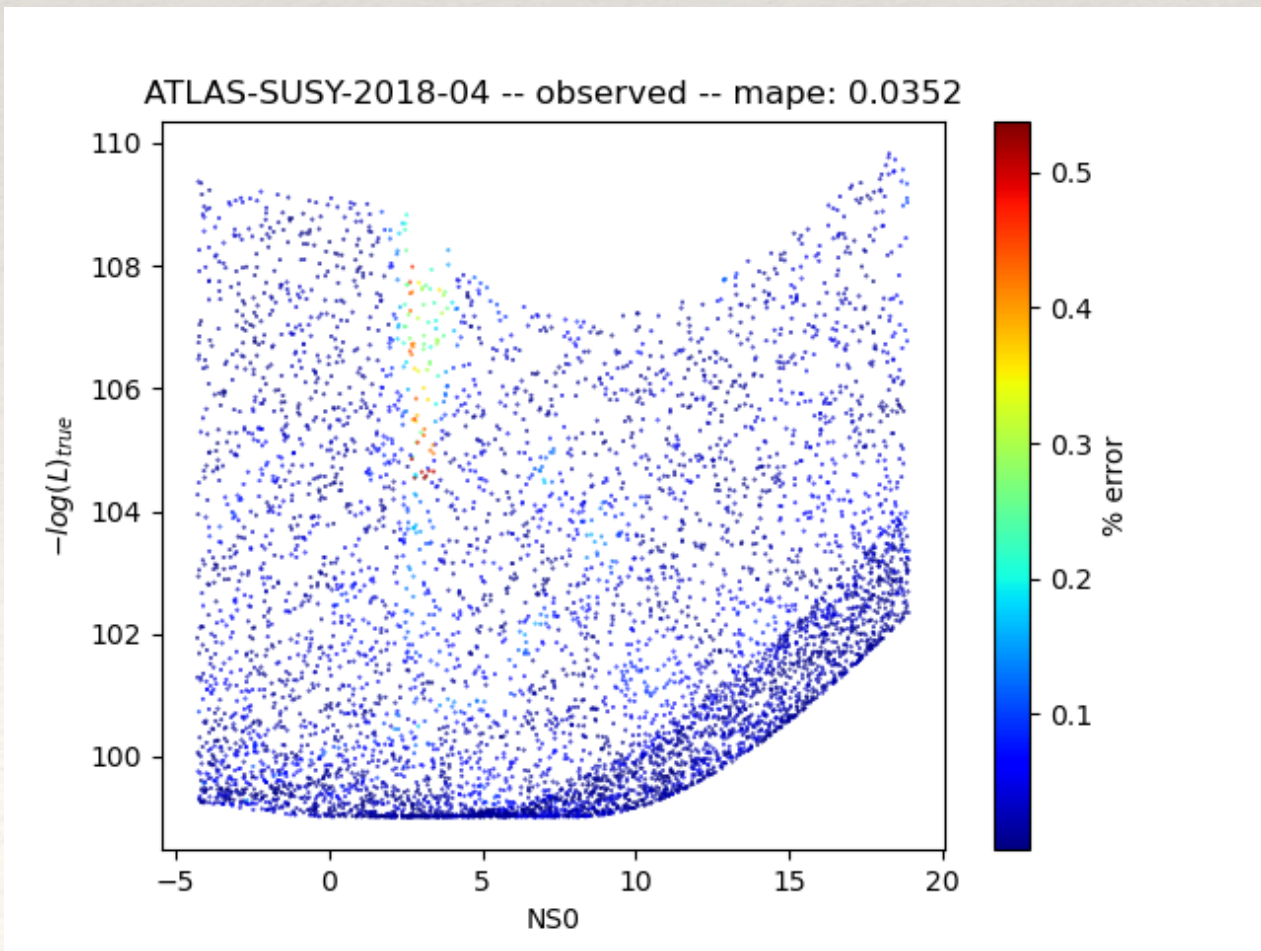
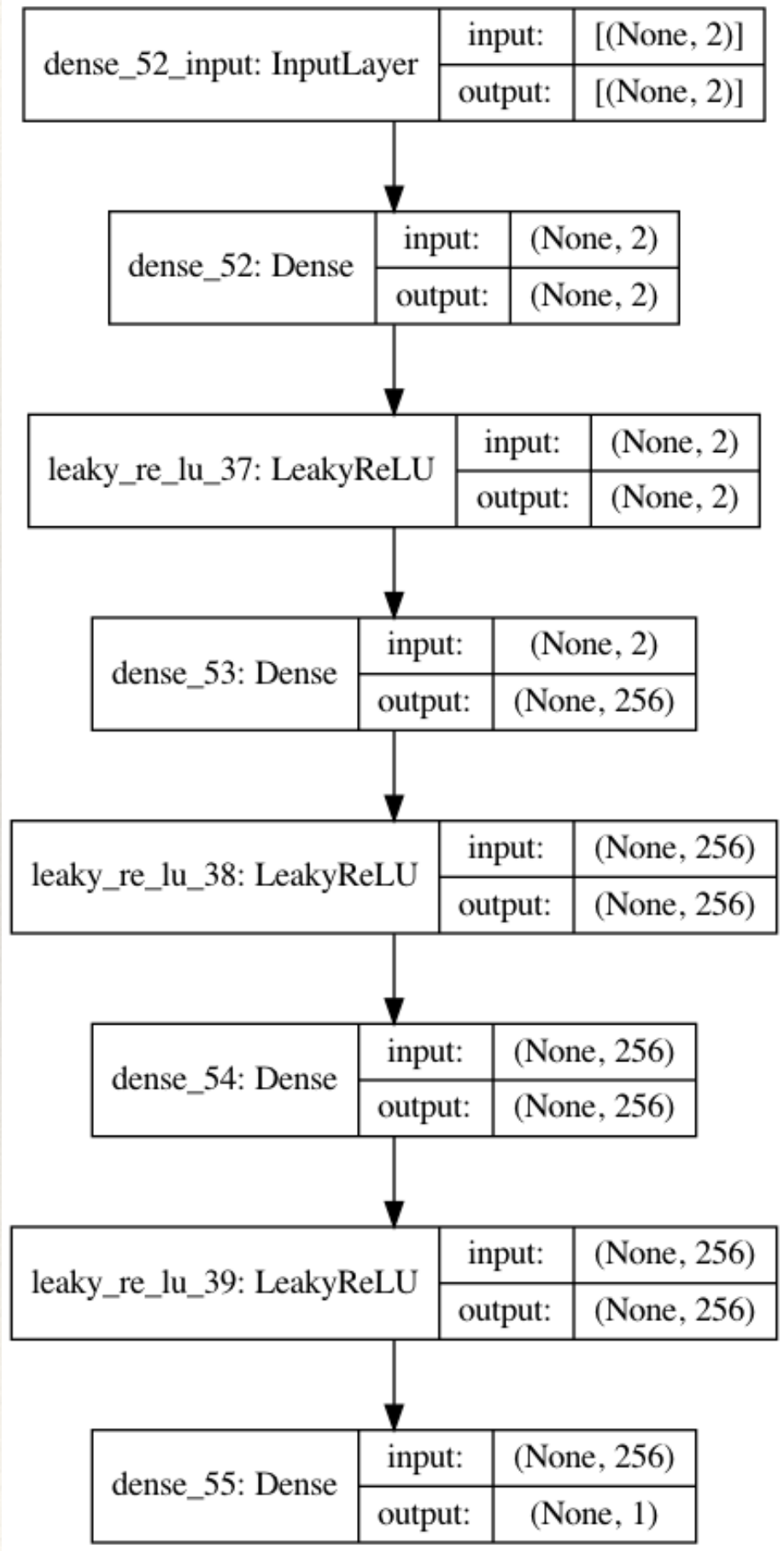
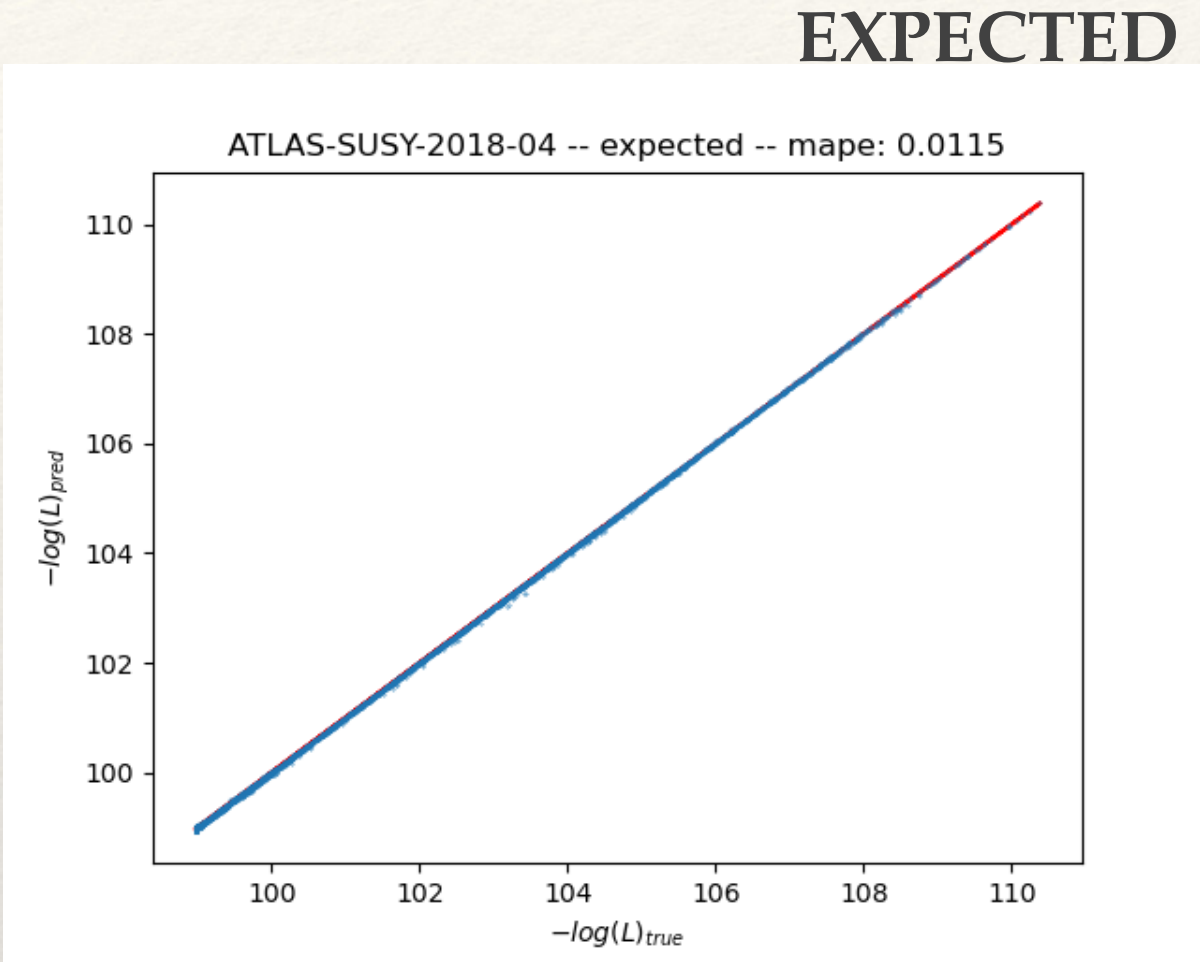
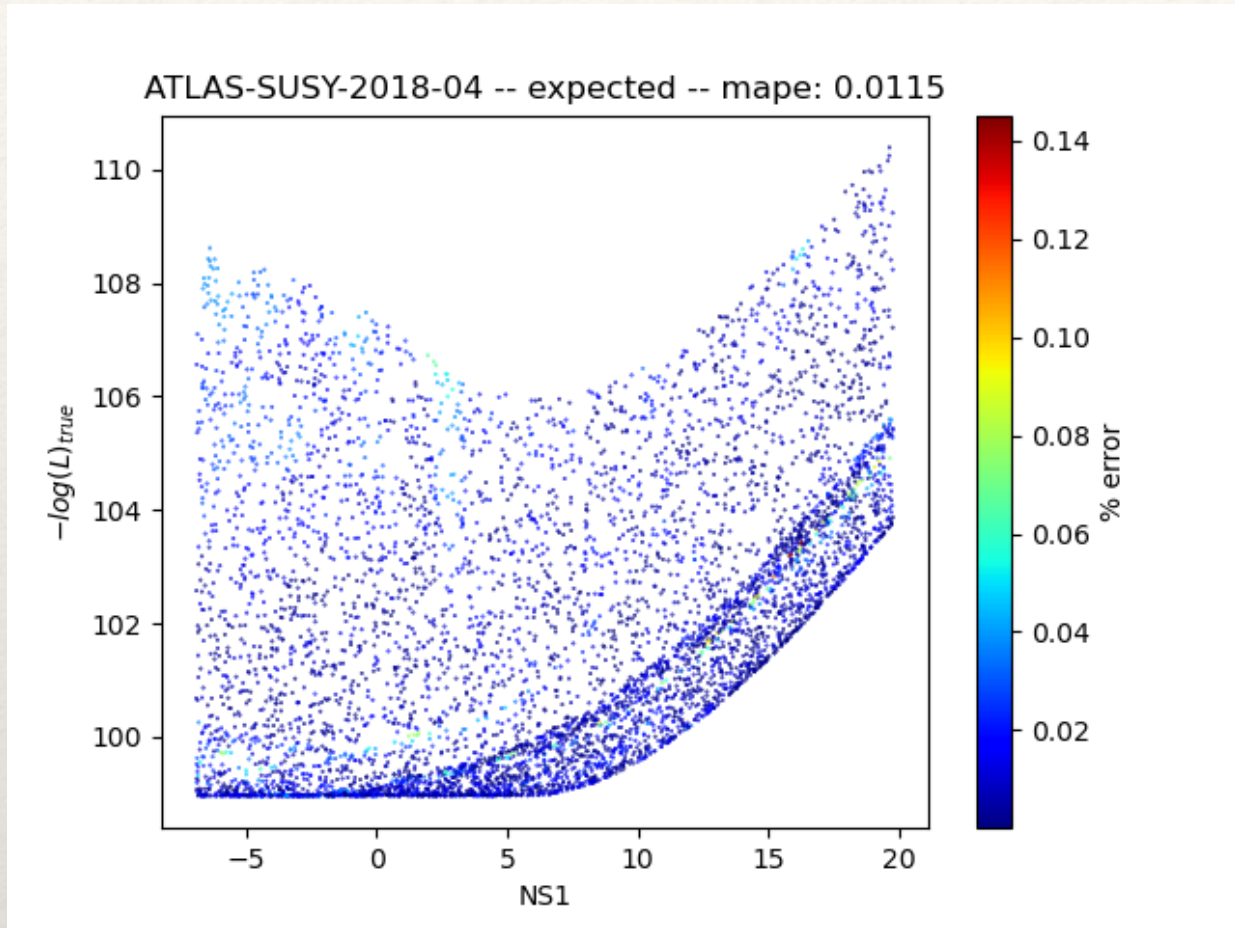
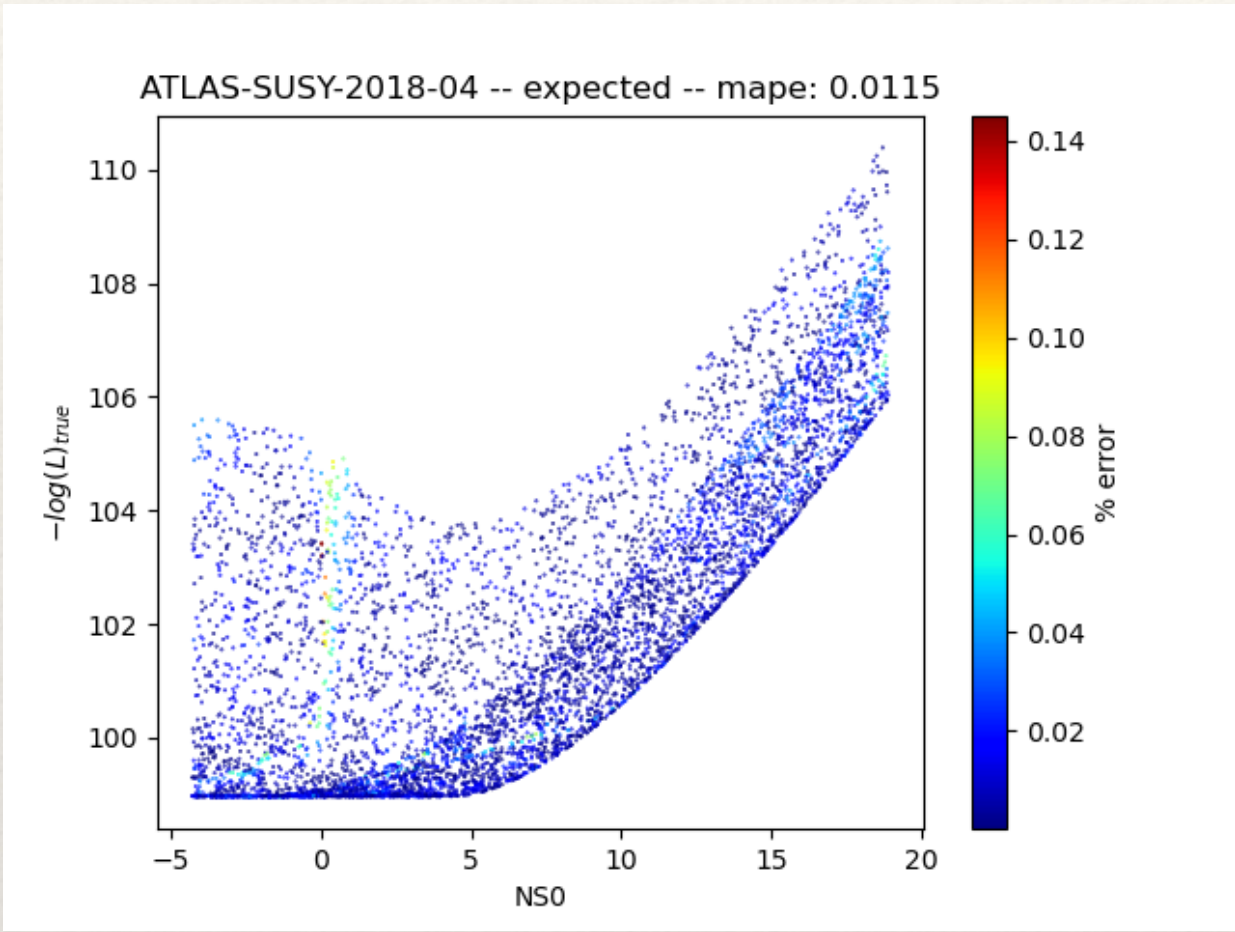
Flavor fit Likelihood



EXTRA: Supervised Learning Profiled Likelihoods

$$P_{profiled}(n_s) = P(n_s | \mu = 1, \hat{\theta}(\mu = 1))$$

EXAMPLE ATLAS-SUSY-2018-04, 2 SRS*:



*Data generated with SModelS' Pyhf interface

Conclusions

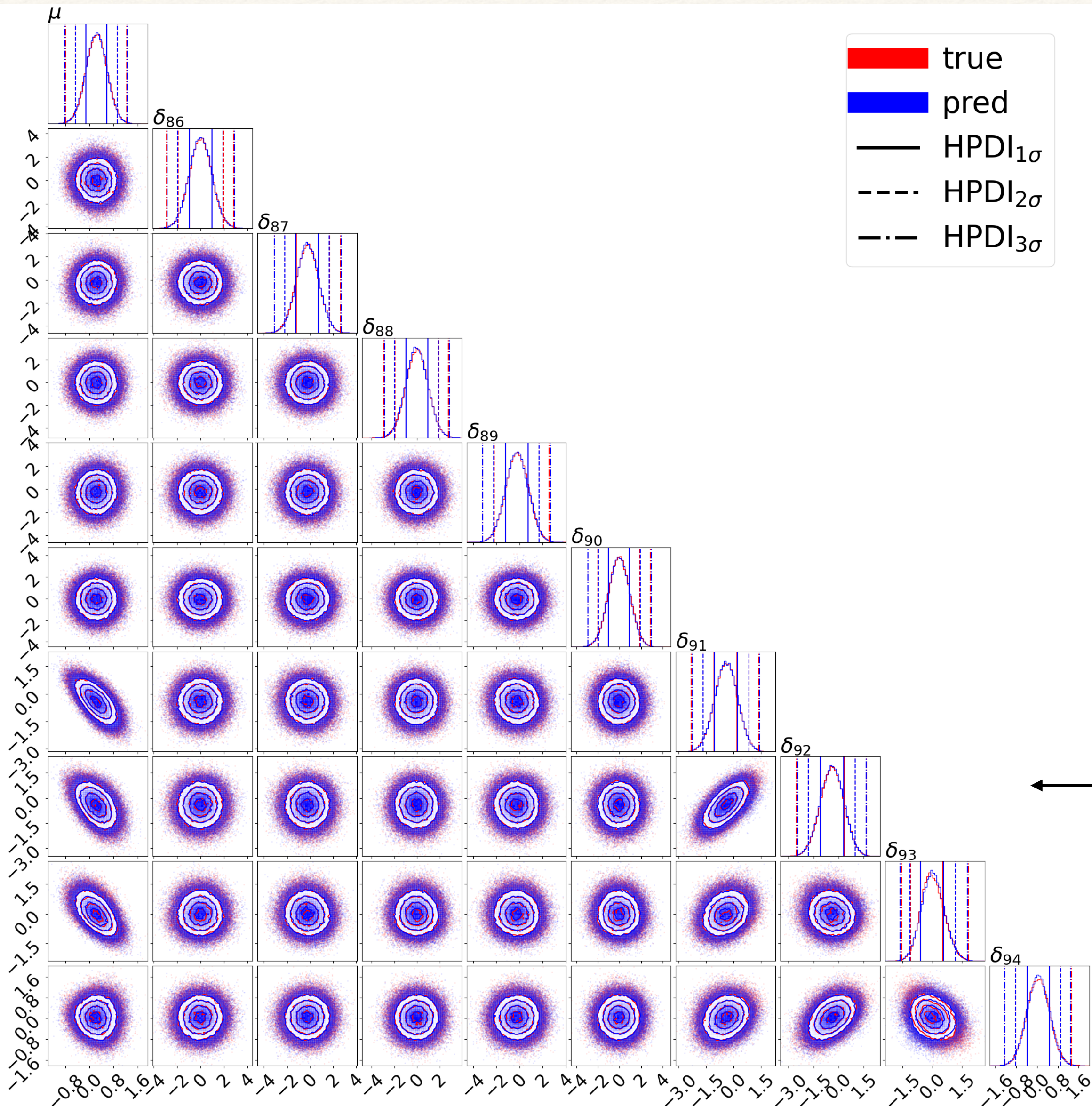
- The preservation of LHC likelihoods is of vital importance (for theorists also).
- Introduced unsupervised learning of Likelihoods with Normalizing Flows.
- Normalizing Flows show great capacity of learning complex high dimensional functions.
- Complementary, we can directly learn profiled likelihoods; useful for fast NP-search reinterpretation.

Outlook

- Paper in preparation arXiv 2301.xxxx.
- User friendly Tensorflow implementation of NFs in dev: https://github.com/riccardotorre/NFTF2_dev
- Learning full statistical models with Conditional Normalizing Flows.
- Learning profiled likelihoods from Pyhf statistical models.

THANK YOU!

LHC-like new physics search Likelihood.



Hyperparameters:

N_{train}	Flow	N bij	Hidden layers	L1 factor	N epochs	N iters.
10^4	MAF	2	256×3	0	20	4
10^5	MAF	2	128×3	10^{-4}	20	4
$2 \cdot 10^5$	MAF	2	64×3	10^{-4}	20	4

Results:

N_{train}	KS-test	W-distance	F- norm	HPDIe _{1σ}	HPDIe _{2σ}	HPDIe _{3σ}	time (s)
10^4	0.479	$1.083 \cdot 10^{-2}$	0.913	$2.211 \cdot 10^{-2}$	$1.374 \cdot 10^{-2}$	$1.3003 \cdot 10^{-2}$	86.65
10^5	0.502	$5.33 \cdot 10^{-3}$	0.527	$2.157 \cdot 10^{-2}$	$8.147 \cdot 10^{-3}$	$1.07 \cdot 10^{-2}$	317.89
$2 \cdot 10^5$.507	$4.82 \cdot 10^{-3}$	0.316	$1.883 \cdot 10^{-2}$	$9.355 \cdot 10^{-3}$	$9.903 \cdot 10^{-3}$	561.82

Test sample : 300k

BACK UP

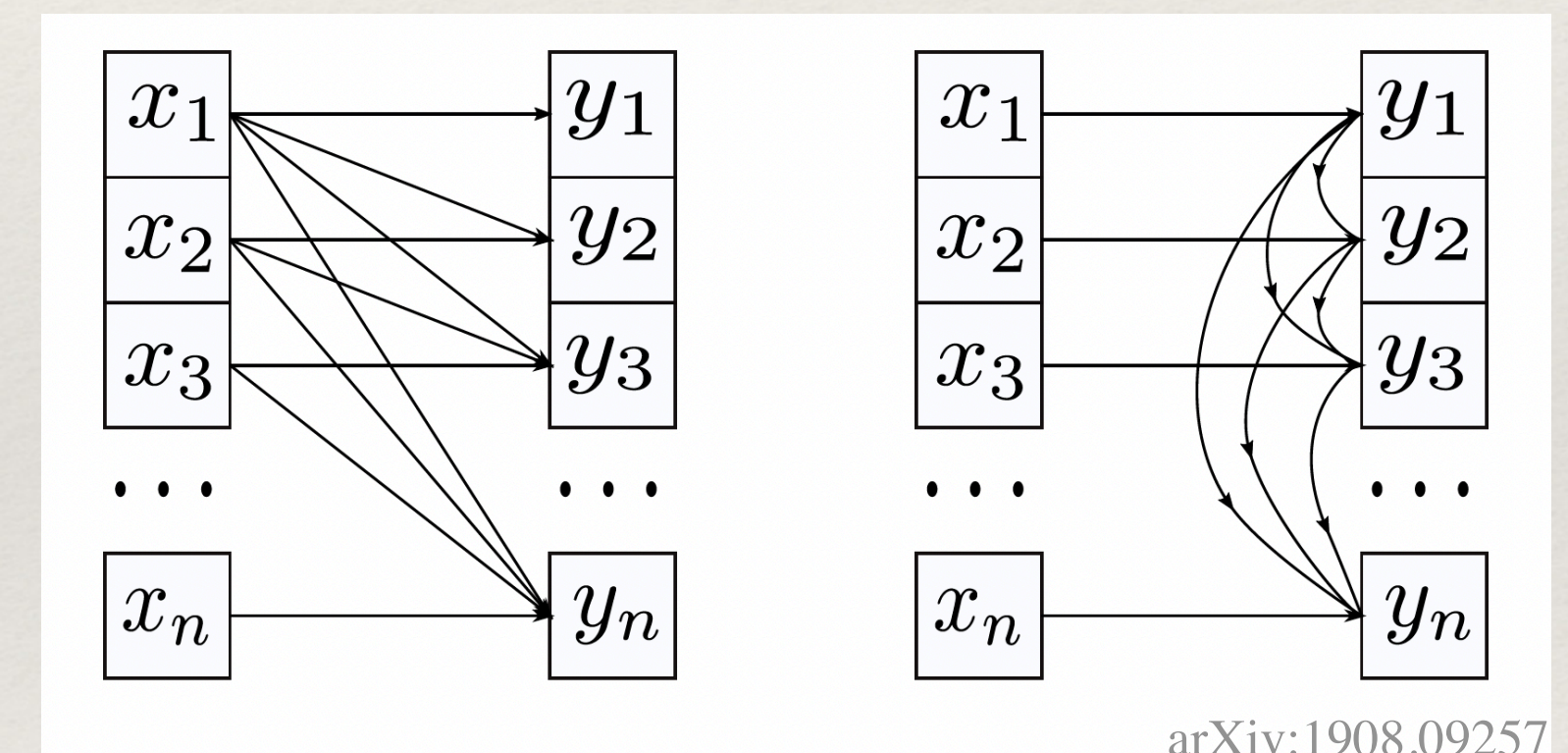
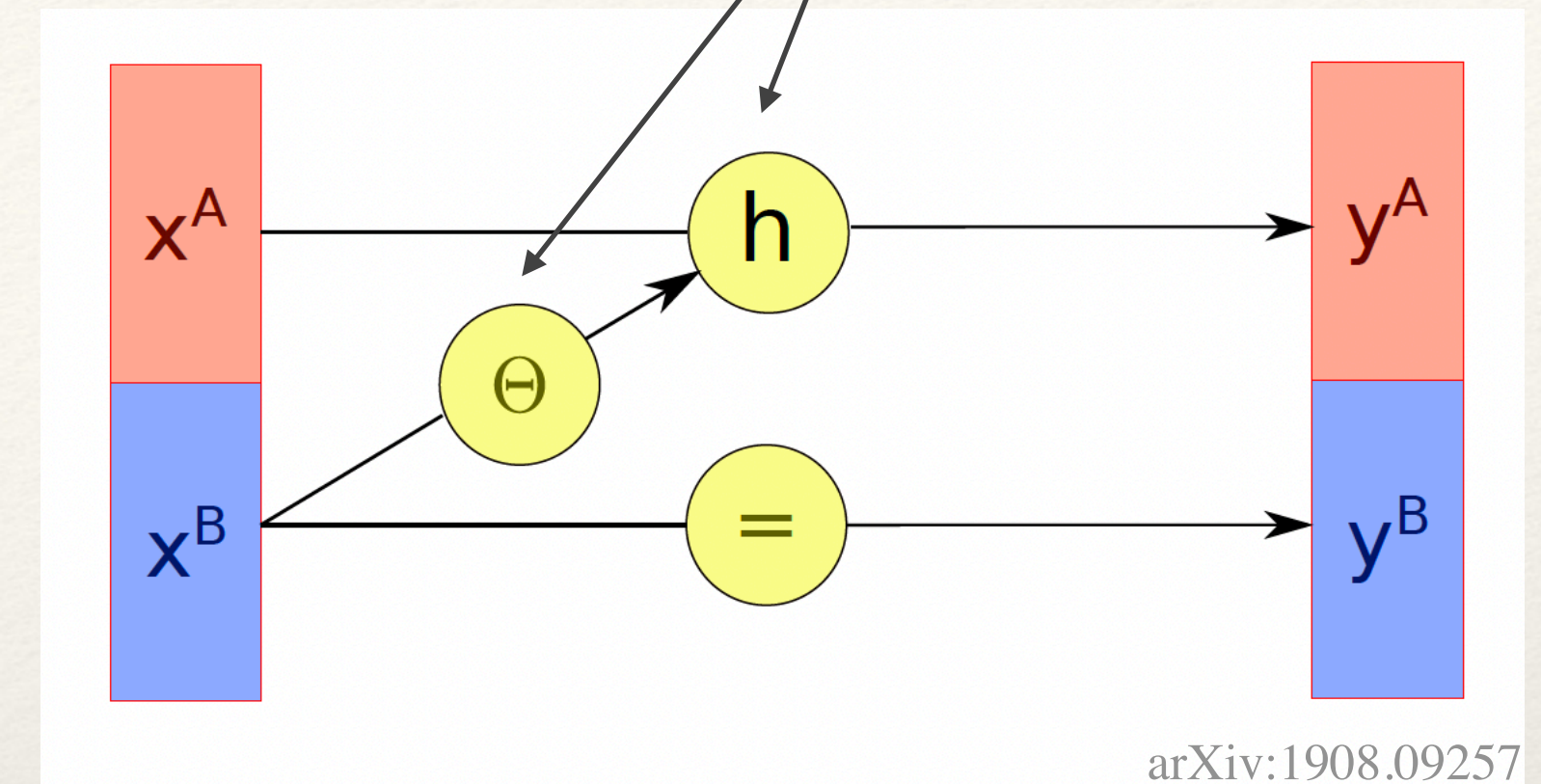
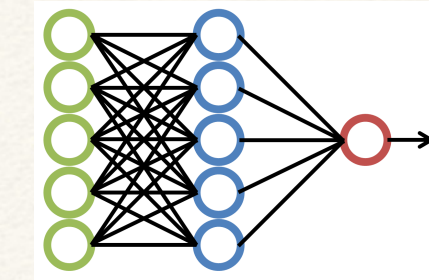
Autoregressive Flows

Coupling Flows :

- Dimensions are divided in two sets: x^A and x^B
- We transform x^B with bijectors trained with x^A .
- The bijector parameters are functions of a NN.
- The Jacobian J is triangular $\rightarrow \det J = \prod_i J_{ii}$
- **Jacobian is easily computed!**
- **Direct sampling AND density estimation.**
- **Less expressive.**

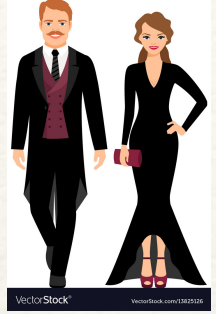
Autoregressive Flows :

- Dimension x^i is transformed with bijectors trained with $y_{1:i-1}$
- Bijector parameters are trained with Autoregressive NNs.
- The Jacobian J is also triangular thus...
- **Jacobian is easily computed!**
- **Direct sampling OR density estimation.**
- **More expressive.**



The loss function:
 $-\log(p_{AF}(target_{dist}))$

Introduction.



Let's get formal...

- If Z is a random variable with pdf p_Z , g is an invertible function such that $Y = g(Z)$ and $f = g^{-1}$, then we can obtain the pdf p_Y of the random variable Y as

$$p_Y(y) = p_Z(f(y)) | \det(Df(y)) | = p_Z(f(y)) | \det(Dg(f(y))) |^{-1} \quad \text{where} \quad Dg(z) = \frac{\partial g}{\partial z} \quad Df(y) = \frac{\partial f}{\partial y}$$

Jacobians

- N transformations are possible since...

$$f = f_1 \circ \dots \circ f_{N-1} \circ f_N$$

$$\det Df(y) = \prod_{i=1}^N \det(Df_i(x_i)) \quad \text{where} \quad x_i = g_i \circ \dots \circ g_1(z) = f_{i+1} \circ \dots \circ f_N(y)$$

- Since p_Z is parametrised by ϕ and the bijector g by θ , we can compute the **log probability** of some measured data $\mathcal{D} = \{y^{(i)}\}_{i=1}^M$ given the parameters $\Theta = (\theta, \phi)$ as

$$\log p(\mathcal{D} | \Theta) = \sum_{i=1}^M \log p_Y(y^{(i)} | \Theta) = \sum_{i=1}^M \log p_Z(f(y^{(i)} | \theta) | \phi) + \log | \det Df(y^{(i)} | \theta) |$$