# An Implementation of Neural Simulation-Based Inference for Parameter Estimation at the LHC

Jay Sandesara on behalf of the ATLAS collaboration





# Introduction

 Neural Simulation-Based Inference (NSBI) covers a broad range of statistical techniques.

- **Idea:** build ML surrogates for powerful statistical inference in the presence of
  - Intractable likelihoods (e.g. LHC analysis), or
  - when likelihoods are slow to compute analytically (e.g. gravitational wave analysis).





Overview of typical NSBI workflow

Figure credits: Madminer

## **NSBI** at the LHC

- The focus of this talk is on a practical application of these methods to LHC analysis, with an example of the off-shell Higgs boson measurement at the ATLAS experiment [Rep. Prog. Phys. 88 067801, Rep. Prog. Phys. 88 057803].
- The talk will cover:

 Efficiently modelling likelihoods as a function of complex high-dimensional parameter space.

 Rigorously testing the quality of the surrogate models and their reliability using MC and real data.

• Building robust frequentist confidence intervals using Neyman Construction.

## **NSBI** at the LHC

Parton-level events sampled from analytical model

MC events sampled from implicit likelihoods

 $x \sim p_S(x \mid \mu), p_B(x \mid \mu)$ 



## **NSBI** at the LHC

. . . . . . . . . .

MC events sampled from implicit likelihoods

 $x \sim p_S(x \mid \mu), p_B(x \mid \mu)$ 



5

## The off-shell Higgs boson

The probability model of the off-shell Higgs boson:

$$p(x \mid \mu) = \frac{1}{\nu(\mu)} \begin{bmatrix} \mu \cdot \nu_{S} \cdot p_{S}(x) + \sqrt{\mu} \cdot \nu_{I} \cdot p_{I}(x) + \nu_{B} \cdot p_{B}(x) + \nu_{NI} \cdot p_{NI}(x) \end{bmatrix}$$

$$\nu \rightarrow \text{Exp events} \qquad \text{Parameter dependence} \qquad \mu = \frac{\sigma_{obs}^{H \rightarrow ZZ}}{\sigma_{exp}^{H \rightarrow ZZ}} \qquad \qquad \text{NI} \rightarrow \text{Non-Interfering}$$

$$backgrounds \qquad \qquad \text{Bkg probability} \qquad$$

Interference probability

ggF Higgs Signal

ggF Interfering Background

## The off-shell Higgs boson



The probability model of the off-shell Higgs boson:

~ Z

Interference probability

ggF Higgs Signal

8 WW

ggF Interfering Background

8 200

 $h \sim z$ 

Fixed S/B discriminant is often the optimal choice for hypothesis testing at the LHC



S/B discriminant used in previous Run-2 analysis of off-shell Higgs boson

Fixed S/B discriminant is often the optimal choice for hypothesis testing at the LHC



Fixed S/B discriminant is often the optimal choice for hypothesis testing at the LHC



#### But what if the parameterization is non-linear?

$$\frac{p(x \mid \mu)}{p_B(x)} \sim \mu \cdot \nu_S \cdot \frac{p_S(x)}{p_B(x)} + \sqrt{\mu} \cdot \nu_I \cdot \frac{p_I(x)}{p_B(x)} + \nu_B \cdot \frac{p_B(x)}{p_B(x)}$$

E.g.: interference effects of off-shell Higgs boson production. Single observable no longer describes the full parameter space!

Fixed S/B discriminant is often the optimal choice for hypothesis testing at the LHC



#### But what if the parameterization is non-linear?

$$\frac{p(x \mid \mu)}{p_B(x)} \sim \mu \cdot \nu_S \cdot \frac{p_S(x)}{p_B(x)} + \sqrt{\mu} \cdot \nu_I \cdot \frac{p_I(x)}{p_B(x)} + \nu_B \cdot \frac{p_B(x)}{p_B(x)}$$

E.g.: interference effects of off-shell Higgs boson production. Single observable no longer describes the full parameter space!



Fixed S/B discriminant is often the optimal choice for hypothesis testing at the LHC



#### But what if the parameterization is non-linear?

$$\frac{p(x \mid \mu)}{p_B(x)} \sim \mu \cdot \nu_S \cdot \frac{p_S(x)}{p_B(x)} + \sqrt{\mu} \cdot \nu_I \cdot \frac{p_I(x)}{p_B(x)} + \nu_B \cdot \frac{p_B(x)}{p_B(x)}$$

E.g.: interference effects of off-shell Higgs boson production. Single observable no longer describes the full parameter space!



Flat NLL region implies sub-optimality in regions with  $\sqrt{\mu} \cdot \nu_I \cdot p_I \gg \mu \cdot \nu_S \cdot p_S$ 

**Frequentist test using NSBI** 

Profile Negative Log-Likelihood Test Statistic





Off-shell Higgs measurement using NSBI Rep. Prog. Phys. 88 057803



Surrogate Model for likelihood ratios

MC events sampled from implicit likelihoods

 $x \sim p_S(x \mid \mu), p_B(x \mid \mu)$ 

**Frequentist test using NSBI** 

#### How do we train this model?



Surrogate Model for likelihood ratios

MC events sampled from implicit likelihoods

 $x \sim p_S(x \mid \mu), p_B(x \mid \mu)$ 

Profile Negative Log-Likelihood Test Statistic

$$-2 \cdot \sum_{i \in events} \log \frac{p(x_i | \mu, \hat{\alpha})}{p(x_i | \hat{\mu}, \hat{\alpha})}$$



Off-shell Higgs measurement using NSBI Rep. Prog. Phys. 88 057803



Proposal 1: estimate paramaterized PDFs  $p(x \mid \mu, \alpha)$ 

train generative models with tractable probability densities (e.g. Normalizing Flows)

Surrogate Model for likelihood ratios



 $x \sim p_S(x \mid \mu), p_B(x \mid \mu)$ 



Surrogate Model for likelihood ratios

#### MC events sampled from implicit likelihoods

 $x \sim p_S(x \mid \mu), p_B(x \mid \mu)$ 

Proposal 1: estimate paramaterized PDFs  $p(x \mid \mu, \alpha)$ 

train generative models with tractable probability densities (e.g. Normalizing Flows)

Proposal 2: Estimate parameterized density ratios

 $\frac{p(x_i \,|\, \mu, \alpha)}{p_{ref}(x)}$ 

by training well-calibrated and unbiased NN classifiers and use in the profile likelihood ratio:

$p(x_i   \mu, \alpha)$	$\rightarrow$	$p(x_i   \mu, \hat{\hat{\alpha}}) / p_{ref}(x)$	$\rightarrow$	$p(x_i   \mu, \hat{\hat{\alpha}})$
$p_{ref}(x)$		$\overline{p(x_i   \hat{\mu}, \hat{\alpha}) / p_{ref}(x)}$		$\overline{p(x_i   \hat{\mu}, \hat{\alpha})}$

 $p_{ref}(x)$  can be any chosen **parameterindependent** hypothesis



Surrogate Model for likelihood ratios

#### MC events sampled from implicit likelihoods

 $x \sim p_S(x \mid \mu), p_B(x \mid \mu)$ 

Proposal 1: estimate paramaterized PDFs  $p(x \mid \mu, \alpha)$ 

train generative models with tractable probability densities (e.g. Normalizing Flows)

Proposal 2: Estimate parameterized density ratios

 $\frac{p(x_i \,|\, \mu, \alpha)}{p_{ref}(x)}$ 

by training well-calibrated and unbiased NN classifiers and use in the profile likelihood ratio:

$$\frac{p(x_i \mid \mu, \alpha)}{p_{ref}(x)} \to \frac{p(x_i \mid \mu, \hat{\alpha}) / p_{ref}(x)}{p(x_i \mid \hat{\mu}, \hat{\alpha}) / p_{ref}(x)} \to \frac{p(x_i \mid \mu, \hat{\alpha})}{p(x_i \mid \hat{\mu}, \hat{\alpha})}$$

 $p_{ref}(x)$  can be any chosen **parameterindependent** hypothesis



for likelihood ratios

MC events sampled from implicit likelihoods

 $x \sim p_S(x \mid \mu), p_B(x \mid \mu)$ 

Proposal 1: estimate paramaterized PDFs  $p(x \mid \mu, \alpha)$ 

train generative models with tractable probability densities (e.g. Normalizing Flows)

Easier to train and validate for large-dimensional inputs

Proposal 2: Estimate parameterized density ratios  $\frac{p(x_i | \mu, \alpha)}{p_{ref}(x)}$ by training well-calibrated and unbiased NN classifiers and use in the profile likelihood ratio:  $\frac{p(x_i | \mu, \alpha)}{p_{ref}(x)} \rightarrow \frac{p(x_i | \mu, \hat{\alpha}) / p_{ref}(x)}{p(x_i | \hat{\mu}, \hat{\alpha}) / p_{ref}(x)} \rightarrow \frac{p(x_i | \mu, \hat{\alpha})}{p(x_i | \hat{\mu}, \hat{\alpha})}$ 

> $p_{ref}(x)$  can be any chosen **parameter**independent hypothesis

Full test statistic function for frequentist parameter estimation on parameter  $\mu$ 

$$t(\mu) = -2 \cdot \log \frac{\mathsf{Pois}(N_{obs} | \mu, \hat{\alpha})}{\mathsf{Pois}(N_{obs} | \hat{\mu}, \hat{\alpha})} - 2 \cdot \sum_{i=1}^{N_{obs}} \log \frac{p(x_i | \mu, \hat{\alpha}) / p_{ref}(x_i)}{p(x_i | \hat{\mu}, \hat{\alpha}) / p_{ref}(x_i)} - 2 \cdot \sum_{k}^{N_{syst}} \log \frac{p_{subs}(\hat{\alpha})}{p_{subs}(\hat{\alpha})}$$
  
Extended

Extended Poisson term Sum of event-by-event log-likelihood ratios

Constraint terms

$$N_{obs} \rightarrow$$
 total observed events

 $p_{subs} \rightarrow$  likelihood from subsidiary measurements of the nuisance parameters

Full test statistic function for frequentist parameter estimation on parameter  $\mu$ 

$$t(\mu) = -2 \cdot \log \frac{\mathsf{Pois}(N_{obs} | \mu, \hat{\alpha})}{\mathsf{Pois}(N_{obs} | \hat{\mu}, \hat{\alpha})} - 2 \cdot \sum_{i=1}^{N_{obs}} \log \frac{p(x_i | \mu, \hat{\alpha})/p_{ref}(x_i)}{p(x_i | \hat{\mu}, \hat{\alpha})/p_{ref}(x_i)} - 2 \cdot \sum_{k}^{N_{syst}} \log \frac{p_{subs}(\hat{\alpha})}{p_{subs}(\hat{\alpha})}$$

$$parameter-independent ratio$$

$$\frac{p(x_i | \mu, \alpha)}{p_{ref}(x_i)} = \frac{1}{\sum_{c} G_{c}(\alpha) \cdot f_{c}(\mu) \cdot \nu_{c}} \sum_{c} \left[ f_{c}(\mu) \cdot g_{c}(x_i | \alpha) \cdot \nu_{c} \cdot \frac{p_{c}(x_i)}{p_{ref}(x_i)} \right]$$

**Parameterized per-event ratios** 

Full test statistic function for frequentist parameter estimation on parameter  $\mu$ 



Full test statistic function for frequentist parameter estimation on parameter  $\mu$ 



Full test statistic function for frequentist parameter estimation on parameter  $\mu$ 

X

$$t(\mu) = -2 \cdot \log \frac{\mathsf{Pois}(N_{obs} \mid \mu, \hat{\alpha})}{\mathsf{Pois}(N_{obs} \mid \hat{\mu}, \hat{\alpha})} - 2 \cdot \sum_{i=1}^{N_{obs}} \log \frac{p(x_i \mid \mu, \hat{\alpha}) / p_{ref}(x_i)}{p(x_i \mid \hat{\mu}, \hat{\alpha}) / p_{ref}(x_i)} - 2 \cdot \sum_{k}^{N_{ops}} \log \frac{p_{subs}(\hat{\alpha})}{p_{subs}(\hat{\alpha})}$$

$$\underbrace{\mathsf{sum over processes}}_{c = S, B, \text{eto.}}$$

$$\underbrace{p(x_i \mid \mu, \alpha)}_{p_{ref}(x_i)} = \frac{1}{\sum_c G_c(\alpha) \cdot f_c(\mu) \cdot \nu_c} \sum_c \left[ f_c(\mu) \cdot g_c(x_i \mid \alpha) \cdot \nu_c \cdot \frac{p_c(x_i)}{p_{ref}(x_i)} \right]$$

$$\overset{\sim \text{multi-dimensional kinematic inputs}}{\underset{kinematic inputs}{x \sim p_{ref}}} \underbrace{s(x) = \frac{p_c}{p_{ref} + p_c}(x)}_{Ciassification NN} \underbrace{f(x) = \frac{p_c}{p_{ref} + p_c}(x)}_{Ciassification Vick}$$

$$\underbrace{\mathsf{Many examples in ATLAS - \underline{\mathsf{HH4b background estimation, Omnifold, etc.}}_{zintering to the track}$$



#### Factorized **nuisance parameter** $\alpha$ -dependence:

$$g_c(x \mid \alpha) = \frac{p_c(x \mid \alpha)}{p_c(x)}$$



#### Factorized **nuisance parameter** $\alpha$ -dependence:

$$g_c(x \mid \alpha) = \frac{p_c(x \mid \alpha)}{p_c(x)}$$

Challenging due to the highdimensionality of  $\alpha = (\alpha_m)$ 



#### **Assumption 1:**



$$g_c(x \mid \alpha) = \frac{p_c(x \mid \alpha)}{p_c(x)}$$

Factorized **nuisance parameter**  $\alpha$ -dependence:

Challenging due to the highdimensionality of  $\alpha = (\alpha_m)$ 

 $\alpha_m$  are orthogonal to each other



Assumption 1:

Factorized **nuisance parameter**  $\alpha$ -dependence:





#### Factorized **nuisance parameter** $\alpha$ -dependence:









But we also need to estimate these parameterized density ratios

The effects from the various NPs  $\alpha_m$  are orthogonal to each other



Train parameterized NNs for each  $\alpha_m$ 



"CARL" approach

Train parameterized NNs for each  $\alpha_m$ 



**Challenges:** 

- Simulations only available at 3 parameter points  $\alpha_m^0, \alpha_m^{\pm 1\sigma_a}$
- Difficult to validate the NN interpolation into phase space with no simulations for testing.





#### **Assumption 2:**

**Semi-analytic approximation** 

$$\frac{p_c(x|\alpha_m)}{p_c(x)} = \begin{cases} \left(\frac{p_c(x|\alpha_m^{+1\sigma_a})}{p_c(x)}\right)^{\alpha_m} & \alpha_m > 1\\ 1 + \sum_{i=1}^6 a_i \boldsymbol{\alpha}_m^i & -1 \le \alpha_m \le 1\\ \left(\frac{p_c(x|\alpha_m^{-1\sigma_a})}{p_c(x)}\right)^{-\alpha_m} & \alpha_m < -1 \end{cases}$$

The  $\alpha$ -dependent negative log-likelihood ratio is a smooth parabolic function



#### **Assumption 2:**

#### Semi-analytic approximation

$$\frac{p_c(x|\alpha_m)}{p_c(x)} = \begin{cases} \left(\frac{p_c(x|\alpha_m^{+1\sigma_a})}{p_c(x)}\right)^{\alpha_m} & \alpha_m > 1\\ 1 + \sum_{i=1}^6 a_i \boldsymbol{\alpha}_m^i & -1 \le \alpha_m \le 1\\ \left(\frac{p_c(x|\alpha_m^{-1\sigma_a})}{p_c(x)}\right)^{-\alpha_m} & \alpha_m < -1 \end{cases}$$



**Assumption 1** 

$$g_c(x \mid \alpha) = \prod_m \frac{p_c(x \mid \alpha_m)}{p_c(x \mid \alpha_m^0 = 0)}$$

The  $\alpha$ -dependent negative log-likelihood ratio is a smooth parabolic function



ratio is a smooth parabolic function

The NSBI approach learns everything, including the parameter scaling and thus the full interference effects

$$\frac{p(x \mid \mu)}{p_B(x)} \sim \mu \cdot \nu_S \cdot \frac{p_S(x)}{p_{ref}(x)} + \sqrt{\mu} \cdot \nu_I \cdot \frac{p_I(x)}{p_{ref}(x)} + \nu_B \cdot \frac{p_B(x)}{p_{ref}(x)} + \nu_{NI} \cdot \frac{p_{NI}(x)}{p_{ref}(x)}$$

Four parameter-independent ratios are trained (suppressing the  $\alpha$ -terms for brevity)

The NSBI approach learns everything, including the parameter scaling and thus the full interference effects

$$\frac{p(x|\mu)}{p_B(x)} \sim \mu \cdot \nu_S \cdot \frac{p_S(x)}{p_{ref}(x)} + \sqrt{\mu} \cdot \nu_I \cdot \frac{p_I(x)}{p_{ref}(x)} + \nu_B \cdot \frac{p_B(x)}{p_{ref}(x)} + \nu_{NI} \cdot \frac{p_{NI}(x)}{p_{ref}(x)}$$
Four parameter-independent ratios are trained
$$t_{\mu} \sim -2 \cdot \sum_{i=1}^{N_{obs}} \log \frac{p(x_i|\mu)}{p(x_i|\hat{\mu})}$$

No "fixed" S/B discriminant - asymptotic optimality throughout  $\mu$  space.

# Additional sensitivity from unbinned nature

The NSBI approach learns everything, including the parameter scaling and thus the full interference effects

$$\frac{p(x|\mu)}{p_{B}(x)} \sim \mu \cdot \nu_{S} \cdot \frac{p_{S}(x)}{p_{ref}(x)} + \sqrt{\mu} \cdot \nu_{I} \cdot \frac{p_{I}(x)}{p_{ref}(x)} + \nu_{B} \cdot \frac{p_{B}(x)}{p_{ref}(x)} + \nu_{NI} \cdot \frac{p_{NI}(x)}{p_{ref}(x)} + \nu$$

 $\mu_{\text{off-shell}}$ 

The NSBI approach learns everything, including the parameter scaling and thus the full interference effects

$$\frac{p(x \mid \mu)}{p_{B}(x)} \sim \mu \cdot v_{S} \cdot \frac{p_{S}(x)}{p_{ref}(x)} + \sqrt{\mu} \cdot v_{I} \cdot \frac{p_{I}(x)}{p_{ref}(x)} + v_{B} \cdot \frac{p_{B}(x)}{p_{ref}(x)} + v_{NI} \cdot \frac{p_{NI}(x)}{p_{ref}(x)}$$
Four parameter-independent ratios are trained
$$nalysis$$
Four parameter-independent ratios are trained
$$nalysis$$
NSBI
analysis
analysis
$$\frac{p_{NI}(x)}{p_{ref}(x)} + v_{NI} \cdot \frac{p_{NI}(x)}{p_{ref}(x)} + v_{NI} \cdot \frac{p_{NI}(x)}{p_{NI}(x)} + \frac{p_{NI}(x)}{p_{NI$$

### **Monte Carlo Diagnostics**

The NN ratios are meticulously trained to be true representations of the density ratios



Do the ratios capture the full un-biased dependence of the multi-dimensional feature space *x* ?



### Where does the sensitivity come from?



The per-event negative log-likelihood ratio allows a probe to identify phase space regions that contribute to the final analysis senstivity.

This allows to identify phase space regions that need robust modeling from Monte Carlo samples.

## **Real Data Diagnostics**

A rigorous data-MC comparison is performed using the parameterized density ratios



Check agreement as a function of any parameter  $\mu$  value

### **Building Frequentist Confidence Intervals**

• Analysis with non-linear parameterizations do not follow Wald approximation:

$$-2\ln\lambda(\mu) = rac{(\mu-\hat{\mu})^2}{\sigma^2}$$
 X

arXiv: 1007.1727

- Neyman construction is essential.
- But standard LHC techniques like Poisson PDF sampling cannot work directly.
- This is because the NSBI technique presented here does not have a PDF  $p(x | \mu, \alpha)$  to sample pseudo-data from - only the density ratios:

$$\frac{p(x \mid \mu, \alpha)}{p_{ref}(x)}$$

**Previous Histogram-based**  $H^* \rightarrow ZZ$  measurement



### **Neyman Construction for NSBI**

• The trained density ratios are used to create unbiased Asimov samples with MC weights  $w_A$  for any value of the  $\mu$ ,  $\alpha$  parameter space:

$$w_A(x \mid \mu_{truth}, \alpha) = \frac{\nu(\mu, \alpha)}{\nu_{ref}} \cdot \frac{p(x \mid \mu_{truth}, \alpha)}{p_{ref}(x)} \cdot w_{ref}(x)$$
 MC weights of reference sample

• Pseudo-experiments are then sampled using the Poisson bootstrap method  $w_{pseudo-data}(x) = Poisson(w_A(x | \mu_{truth}, \alpha)).$ 



### **Conclusion and Outlook**

- An implementation of NSBI was presented focused on building likelihood ratios as a function of complex, large-dimensional parameter spaces using well-motivated approximations.
- The NSBI approach presented in this talk has broad applicability across LHC analysis - particularly effective when the likelihood model is non-linear in the parameter of interest and when multi-dimensional information is needed for extra precision.
- The various conceptual and computational developments have been done and published in these companion papers by ATLAS:
  - <u>Rep. Prog. Phys. 88 067801</u> [General NSBI method presented in this talk]
  - <u>Rep. Prog. Phys. 88 057803</u> [Application to off-shell Higgs boson and Higgs boson width measurement with the ATLAS experiment]

# Backup

### **Parameterized Observables and Unbinning**



The improved sensitivity from using the NSBI approach is a result of:

- using parameterized information for the hypothesis testing
- and doing an unbinned fit

## **Uncertainty Parameterization**

$$\frac{p(x_i | \mu, \alpha)}{p_{ref}(x_i)} = \frac{1}{\sum_c G_c(\alpha) \cdot f_c(\mu) \cdot \nu_c} \sum_c \left[ f_c(\mu) \cdot g_c(x_i | \alpha) \cdot \nu_c \cdot \frac{p_c(x_i)}{p_{ref}(x_i)} \right]$$
  
Factorized yield  $\alpha$ -dependence:  
$$G_c(\alpha) = \prod_k \frac{\nu_c(\alpha_k)}{\nu_c}$$
Per-event analog of standard techniques  
$$g_c(x | \alpha) = \prod_k \frac{p_c(x | \alpha_k)}{p_c(x)}$$

with  $\nu_c(\alpha_k)/\nu_c$  estimated using **analytic interpolation techniques:** 

Available from simulations  
at 
$$\alpha_k = 0$$
,  $\alpha_k^+$ ,  $\alpha_k^-$   

$$\frac{\nu_c(\alpha_k)}{\nu_c} = \begin{cases} \left(\frac{\nu_c(\alpha_k^+)}{\nu_c}\right)^{\alpha_k} & \alpha_k > 1 \\ 1 + \sum_{n=1}^6 c_n \alpha_k^n & -1 \le \alpha_k \le 1, \\ \left(\frac{\nu_c(\alpha_k^-)}{\nu_c}\right)^{-\alpha_k} & \alpha_k < -1 \end{cases}$$

with  $p_c(x \mid \alpha_k)/p_c(x)$  estimated using a mix of NNs and analytic interpolation techniques:

Density ratios trained using NNs from simulations  $\alpha^+ \alpha^-$ 

$$\frac{\operatorname{at} \alpha_{k} = 0, \ \alpha_{k}^{+}, \ \alpha_{k}}{p_{c}(x \mid \alpha_{k})} = \begin{cases} \left( \underbrace{\frac{p_{c}(x \mid \alpha_{k}^{+})}{p_{c}(x)}} \right)^{\alpha_{k}} & \alpha_{k} > 1 \\ 1 + \sum_{n=1}^{6} c_{n} \alpha_{k}^{n} & -1 \le \alpha_{k} \le 1 \\ \left( \underbrace{\frac{p_{c}(x \mid \alpha_{k}^{-})}{p_{c}(x)}} \right)^{-\alpha_{k}} & \alpha_{k} < -1 \end{cases}$$

### **Challenges: Density Ratio Estimation**

**Challenge:** The best fit value from a profile likelihood fit  $\hat{\mu}$  with a single NN per  $p_c/p_{ref}$  is biased.

Solution: An ensemble of O(100) or more NNs were trained to be robust against this bias.



By building an ensemble of NNs per  $p_c/p_{ref}$  we become **robust against the bias** in the fit value:

$$\hat{\mu} \rightarrow \mu_{truth}$$

### Full workflow of the NSBI Analysis

