

EUROPEAN AI FOR FUNDAMENTAL PHYSICS CONFERENCE EuCAIFCon 2025

## Tests for model misspecification in simulation-based inference

**Noemi Anau Montel** 17 June 2025



## Simulation-based inference

 Simulation-based Inference - II: Application of So

scattering transform

Simultary

Dem

Didie

Opti

Bene

Isola

Vane

Con

Reio

Fred

Infe

[arX

Time

Sch

Effic

Nor

Tian

Opti

Dee

Lear

Hua

Sim

Phel

Bay

M. Gatti, N

Porredon,

Becker, G

SimBIG: F

Pablo Lem

Chirag Mc

SIMBIG: C

Bruno Réc

Modi, Aza

SIMBIG: T

ChangHod

Azadeh M

Klaus Dola

Beatriz Tu

Chiraq Mc

Sensitivit

#### Papers: App

Domain application of ne

#### Cosmology, Asti

- Extracting cosmole inference [arXiv] Íñigo Zubeldia, Bori
- · Trial by FIRE: Probi Tri Nguyen, Justin F Wetzel
- How many simulati Anirban Bairagi, Ber
- Field-leve · Fast Sampling of C Natalí S. N Oleg Savchenko, Gu Castro, Yu
- Cosmological Anal He Jia
- HaloFlow Hybrid Summary S ChangHod T. Lucas Makinen, C
- EFTofLSS · What to do when the [arXiv]
- James Alvey, Uddip
- Leveraging Time-Dependent Instru James Alvey, Uddipta Bhardwaj, Vale

#### Xiaosheng Zhao, Yi M **Application to Real Data** Dark Energy Survey Y

Applications of neural simulation-based inference beyond synthetic data.

- A robust neural determination of the source-count distribution of the Fermi-LAT sky [arXiv]
- Christopher Eckner, Noemi Anau Montel, Florian List, Francesca Calore, Christoph Wenig
- SimBIG : A Forward Modeling Approach To Analyzing Galaxy Clustering [arXiv] ChangHoon Hahn et al
- Mental speed is high until age 60 as revealed by analysis of over a million participant Krause, Stefan T. Radev, Andreas Voss
- A neural simulation-based inference approach for characterizing the Galactic Center γ-ray excess [arXiv] Siddharth Mishra-Sharma, Kyle Cranmer
- Towards constraining warm dark matter with stellar streams through neural simulation-based inference (Preliminary) [arXiv]

Joeri Hermans, Nilanjan Banik, Christoph Weniger, Gianfranco Bertone, Gilles Louppe

- OutbreakFlow: Model-based Bayesian inference of disease outbreak dynamics with invertible neural networks and its application to the COVID-19 pandemics in Germany [arXiv] Stefan T. Radey, Frederik Graw, Simiao Chen, Nico T. Mutters, Vanessa M. Eichel, Till Bärnighausen, Ullrich
- Likelihood-free inference with neural compression of DES SV weak lensing map statistics [arXiv] Niall Jeffrey, Justin Alsing, François Lanusse



harma

sian Neural Doppler Imaging [arXiv] C. Diaz Baso, O. Kochukhov

Implicit Likelihood Inference for Cosr fom Charnock, Justin Alsing, Benjamin

onal-wave science with neural poster phen R. Green, Jonathan Gair, Jakob H

od-Free Inference of Roman Binary M on [arXiv]

Keming Zhang, Joshua S. Bloom, B. Scott Gaudi, Francois egł

 Towards constraining warm dark matter with stellar stre rec [arXiv] ima

i Di Joeri Hermans, Nilanjan Banik, Christoph Weniger, Gianfra

- Lightning-Fast Gravitational Wave Parameter Inference าq t
- Arnaud Delaunoy, Antoine Wehenkel, Tanja Hinderer, Sama on [ Williamson, Gilles Louppe าอน
- The sum of the masses of the Milky Way and M31: a like ike Pablo Lemos, Niall Jeffrey, Lorne Whiteway, Ofer Lahav, Ni ıg Z
- Likelihood-free inference with neural compression of DE ted
- Niall Jeffrey, Justin Alsing, François Lanusse on
- Mining for Dark Matter Substructure: Inferring subhalo p onmachine learning [arXiv] gin
- Johann Brehmer, Siddharth Mishra-Sharma, Joeri Herman Credible Likelihood-Free Cosmology with Truncated Marginal Ne

Alex Cole, Benjamin Kurt Miller, Samuel J. Witte, Maxwell X. Cai, Meiert W. Gro

Köthe

## Model misspecification

# $\frac{\text{Model}}{x, \theta \sim p(x \mid \theta) p(\theta)}$



#### Data

 $x_{obs}$ 



#### Does my model fit the data? If not, where and how does it fail?



Outline

Model misspecification diagnostics

Inference robustness checks
Posterior predictive validation
Contrastive model diagnostics

Classification from TASI lectures on structured reasoning for SBI, Christoph Weniger [to appear!!]

Structured test batteries from augmented simulators

NAM, J. Alvey, C. Weniger (PRD, 2025) [arXiv: 2412.15100]



## 1. Inference robustness checks

Is our inference stable under small, structured changes to the data or the inference pipeline?

- Masking part of the data
- Changing summary network architecture
- Altering training algorithm, optimizers, seeds etc.



Cannon+22 [arXiv:2209.01845]



Check Christopher Eckner's talk and poster!

C. Eckner, NAM+25 [arXiv:2505.02906]

#### 2. Posterior predictive validation

Does our model produce data that look like what we observed?

$$p(x_{\text{new}} \mid x_{\text{obs}}) = \int p(x_{\text{new}} \mid \theta) p(\theta \mid x_{\text{obs}}) d\theta$$

C



Dupourqué+25 [arXiv:2506.05911]



KiDS-SBI 2024 analysis [arXiv:2404.15402]

## 3. Contrastive model diagnostics

Which model describe our data better?

$$K = \frac{p(x_{obs} | H_1)}{p(x_{obs} | H_0)} \quad \text{with} \quad p(x_{obs} | H_j) = \int p(x_{obs} | \theta, H_j) p_j(\theta) d\theta$$

# SBI Bayesian model comparison via Evidence Networks:

Train a discriminative classifier via a special loss to distinguish between simulations from each model.

Jeffrey+23 [arXiv:2305.11241]



 $\rightarrow$  How to extend this contrastive diagnostic reasoning to a structured framework for multiple tests?

 $x^{(0)} \sim p(x \,|\, H_0)$ 



$$\begin{cases} x^{(0)} \sim p(x \mid H_0) \\ x^{(i)} \sim p(x \mid H_i) \end{cases}$$

<



#### 1. Define structured alternatives.



1. Define structured alternatives.

2. Define log-likelihood ratio test statistics for the alternatives.

$$\begin{cases} x^{(0)} \sim p(x | H_0) \\ x^{(i)} \sim p(x | H_i) \end{cases} \longrightarrow t_i(x) = -2 \log \frac{p(x | H_0)}{p(x | H_i)} \longrightarrow p_i(x) = \mathbb{E}_{x \sim p(x | H_0)} [\mathbb{I}(t(x) > t(x_{obs}))] \end{cases}$$

1. Define structured alternatives.

2. Define log-likelihood ratio test statistics for the alternatives.

3. From test statistics to p-values.

NAM, J. Alvey, C. Weniger [arXiv: 2412.15100]

\_

Structured test batteries from augmented simulators: Localized and aggregated tests



Localized test statistics are more sensitive towards single isolated distortions, and, in some limits, lead to matched filter and anomaly localization "bump-hunt" type of analyses. Aggregated test statistics provides

complementary information about the statistical significance of favoring the alternatives  $H_i$  over the baseline model  $H_0$ , and, in some limits, lead to model validation statistics.

Structured test batteries from augmented simulators: Localized and aggregated tests



Localized test statistics are more sensitive towards single isolated distortions, and, in some limits, lead to matched filter and anomaly localization "bump-hunt" type of analysis.

NAM, J. Alvey, C. Weniger [arXiv: 2412.15100]

Aggregated test statistics provides

complementary information about the statistical significance of favoring the alternatives  $H_i$  over the baseline model  $H_0$ , and, in some limits, lead to model validation statistics.

#### Structured test batteries from augmented simulators: Global significance





Because many hypotheses (localized and aggregated) are tested in parallel, multiple testing corrections are required. A common approach is to compute a give ball p versuebased on the most extreme individual result.

Global p-value  $p_{glob} = 7.57 \times 10^{-3}$ 

Structured test batteries from augmented simulators: Connection to classical testing frameworks

• Localized test statistics for additive distortions are closely related to matched filters and signal-tonoise ratio (SNR) statistics.  $H_i: \tilde{x} = x + \epsilon \cdot n^{(i)}$  with  $\epsilon \sim \mathcal{U}(-b, b)$ 

Assuming a Gaussian likelihood function for the base model, in the large sample limit, and for scenarios where the maximum-likelihood estimator is not significantly correlated with the distortion, the test statistic for a given distortion is directly related to the signal-to-noise ratio (SNR) of that distortion in the data

 $t_i(\mathbf{x}) \simeq \text{SNR}_i^2(\mathbf{x}) + \text{const.}$ 

• Aggregated discrepancy scores reduce to  $\chi^2$  tests when distortions are orthogonal and noise is Gaussian.

 $t_{sum}(\mathbf{x}) \simeq \chi^2 + \text{const.}$ 

• In this sense, simulator augmentation **extends** residual analysis and goodness-of-fit testing to the flexible, implicit-likelihood setting of SBI.

# What is the **training strategy** to build such structured test statistics batteries from augmented simulators?



Structured test batteries from augmented simulators: Training strategies

BCE: discriminative classifiers can be used to approximate the generalized likelihood ratio statistic.

$$\mathcal{L}_{\text{BCE}}^{(i)}\left[f_{i,\phi}(\boldsymbol{x})\right] = \mathbb{E}_{\boldsymbol{x} \sim p(\boldsymbol{x}|H_0)}\left[-\ln\sigma(f_{i,\phi}(\boldsymbol{x}))\right] + \mathbb{E}_{\boldsymbol{x} \sim p(\boldsymbol{x}|H_i)}\left[-\ln\sigma(1 - f_{i,\phi}(\boldsymbol{x}))\right] \quad \text{with} \quad t \approx 2f_{i,\phi}(\boldsymbol{x})$$

**SNR:** minimizing a Gaussian negative log-likelihood loss for the MLE of the matched filter  $\epsilon(\mathbf{x})$  and its variance  $\sigma^2$ .

$$\mathscr{L}_{\text{SNR}}^{(i)}\left[\epsilon_{i,\phi}(\boldsymbol{x}), \sigma_{i,\phi}^{2}\right] = \mathbb{E}_{\boldsymbol{x}, \epsilon \sim p(\boldsymbol{x}|H_{i}, \epsilon)p(\epsilon)}\left[\frac{(\epsilon_{i,\phi}(\boldsymbol{x}) - \epsilon)^{2}}{\sigma_{i,\phi}^{2}} + \ln \sigma_{i,\phi}^{2}\right] \quad \text{with} \quad t_{i}(\boldsymbol{x}) \propto \frac{\epsilon_{i,\phi}(\boldsymbol{x})}{\sigma_{i,\phi}}$$

When distortions correspond to structured changes in space (e.g., image domains), the individual networks can be trained jointly using **shared neural architectures:**  $f_{\phi}(x) : \mathcal{D} \to \mathbb{R}^{N_i}$  and  $\epsilon_{\phi}(x) : \mathcal{D} \to \mathbb{R}^{N_i}$ ,  $\sigma_{\phi}^2 : \mathcal{D} \to \mathbb{R}^{N_i}$ 

## Residual analysis

SNR training strategy allows to visualize where distortions are located in data space and how large they are.





## Adaptive learning of distortions amplitude

$$b = \text{SNR}_{\max}\sigma = \frac{\text{SNR}_{\max}}{\sqrt{(\boldsymbol{n}^{(i)})^T \Sigma^{-1} \boldsymbol{n}^{(i)}}}$$

The algorithm converges to distortions that are significant enough to be detectable, but not so significant that they are clearly ruled out.



## An application to GW150914

- Fit model to data using jimgw [Wong+23 - <u>arXiv:2302.05333</u>].
- Construct structured test batteries of independent and correlated distortions to the model and test for misspecification.
- As expected, no significant anomaly is present in the modelling of GW150914, with global *p*-values for all the types of analyses of around a few tenths.



#### Summary

- Model misspecification analysis strategies are integral to advancing our understanding of physical phenomena.
- The framework presented here is designed to carry this out in a SBI context. By leveraging classical concepts, it provides a flexible and comprehensive approach to simultaneously perform many hypothesis tests and quantify their statistical significance.
- Via two **training strategies**, we can actually test (and Monte Carlo sample) all of these alternative hypotheses simultaneously. This makes the pipeline very **efficient** when looking to test for broad classes of mismodelling, while still maintaining the ability to carry out individual, targeted tests.
- The **SNR training strategy**, can be used to visualize model residuals and calibrate the scale of distortions searched for in the data.

## Thanks!

# Backup slides

#### Framework summary: analytic test



#### Framework summary: correlated distortions



#### Framework summary: multiple small correlated distortions



#### The effect of marginalization

$$t(x) = -2\log \frac{\int p(x \mid \theta, H_0) p_0(\theta) d\theta}{\int p(x \mid \theta, H_1) p_1(\theta) d\theta}$$

