



FAIR Universe HiggsML Uncertainty Challenge

Wahid Bhimji, Paolo Calafiura, <u>Ragansu Chakkappai</u>, Po-Wen Chang, Yuan-Tang Chou, Sascha Diefenbacher, Jordan Dudley, Steven Farrell, Aishik Ghosh, Isabelle Guyon, Chris Harris, Shih-Chieh Hsu, Elham E Khoda, Benjamin Nachman, Peter Nugent, David Rousseau, Benjamin Thorne, Ihsan Ullah, Yulei Zhang





Bias and uncertainty in Physics



Differences between simulation and data can bias measurements

Landscape of Uncertainty Aware Learning

- "pivot" Louppe, Kagan, Cranmer : <u>arXiv:1611.01046</u>
- "Uncertainty-aware" approach of Ghosh, Nachman, Whiteson <u>PhysRevD.104.056026</u>
 - $\circ \quad \text{Parameterize classifier using Z}$
 - Measured on "Toy" 2D Gaussian Dataset and dataset from <u>HiggsML Challenge</u> modified to include systematic on tauenergy scale
 - \circ $\,$ $\,$ Performs as well as classifier trained on true Z $\,$
- Other novel approaches e.g. (not comprehensive)
 - Inferno: <u>arxiv:1806.04743</u>
 - Direct profile-likelihood: e.g. <u>arxiv:2203.13079</u>
 - (Neuro) Simulation Based Inference has to include Z: <u>arXiv:1911.01429</u>



(Signal Strength)

Landscape of Uncertainty Aware Learning



Fair Universe: HiggsML Uncertainty Challenge



FAIR UNIVERSE - HIGGS UNCERTAINTY CHALLENGE



- Full HiggsML Uncertainty Challenge Ran from September 12 to March 14th
- Accepted as <u>NeurIPS competition</u> 2024
- Dedicated <u>workshop at NeurIPS</u> (most important ML conference) 2024 at December 14th
- Final results presented on May at CERN(<u>satellite event of IML 2025</u>)

PARTICIPANTS

SUBMISSIONS

160

347

Physics Problem



Fair Universe: HiggsML Uncertainty Challenge

	²⁰¹⁴ Higgs	2024 FAIR Universe
Data Source	ATLAS Open Data (Sample of real data)	New simulated dataset
Data Quantity	800,000	400,000,000
Parameterized systematics	×	6 Systematics
Inference	Classification accuracy	Signal strength (μ) CI (Confidence Interval) with pseudo-experiment
Metric	Z - Classification Significance	New metric aimed at CI Coverage

Challenge Dataset

- Generated data with fast simulation with Pythia LO and Delphes detector simulation
- Using the updated Delphes ATLAS card
- Generated ~280 Million Events after initial cuts equivalent to 220 X 10 fb-1 + Higgs
- Data organised into **tabular** form with **28** features per event.

Process	Number Generated	LHC Events @10fb ⁻¹ Label	
Higgs	52101127	1015	signal
Z Boson	221724480	1002395	background
Di-Boson	2105415	3783	background
$t \overline{t}$	12073068	44190	background



Dataset permanently released on Zenodo in May 2025 for Future Benchmarks https://zenodo.org/records/15131565

Dataset : some feature examples



Challenge Datasets - Systematics

Apply parameterized systematics (Nuisance Parameters) :

- Feature distortions:
 - Tau Energy Scale (**TES**) (and correlated MET)
 - Jet Energy Scale (**JES**) (and correlated MET impact)
 - Additional randomised Soft MET

- Event category normalisation
 - Background overall normalisation
 - Di-boson background normalisation
 - ttbar background normalisation



Tau Energy Scale Systematics Applied

Histogram between nominal (TES = 1) and shifted (TES = 0.9) TES = 0.9, is an exaggeration, in practice it is sampled with a gaussian of 1 + - 0.01



Competition Work Flow



Confidence Interval Evaluation



- Form multiple pseudo-experiment test sets: different signal strengths (μ) and systematics
 - 10μ times 100 pseudo-experiments (Online Phase)
 - **1000** μ times **100** pseudo-experiments (Offline Phase)
- Task: predict uncertainty interval [μ_{16}, μ_{84}]
 - \circ E.g. 68% quantile of likelihood or assume 1σ



Confidence Interval Evaluation



- Form multiple pseudo-experiment test sets: different signal
 - 10μ times 100 pseudo-experiments (Online Phase)
 - 1000μ times 100 pseudo-experiments (Offline Phase)

0.5

0.6

0.7

X

0.8

- Task: predict uncertainty interval $[\mu_{16}, \mu_{84}]$
 - E.g. 68% quantile of likelihood or assume 1σ

European Physical Society Conference on High Energy Physics

 $f(x)N_{test} = 100$

0.9

 $f(x)N_{tost} = 1000$

1.0







Winners

Medal	Rank	Team	Avg Coverage	Avg Interval	Avg Quantile Score
	1 (Tie)	HEPHY	0.6683	0.4599	-0.5823
	1 (Tie)	IBRAHIME	0.6698	0.4974	-0.5761
3	3	HZUME	0.6659	0.8134	-2.1650

- HEPHY (Lisa Benato, Cristina Giordano, Claudius Krause, Ang Li, Robert Schöfbeck, Dennis Schwarz, Maryam Shooshtari, Daohan Wang) from Vienna's Institute of High Energy Physics (HEPHY) in Austria will win \$2000.
- IBRAHIME (Ibrahim Elsharkawy) from University of Illinois at Urbana-Champaign will win \$2000.
- HZUME (Hashizume Yota) from Kyoto University Japan will win \$500.

Final evaluation on 1000 x 100 p-e



3rd - Hashizume Yota - Kyoto University Japan

Decision-Tree Aggregated Features and Hybrid Bin-Classifier/Quantile-Regressor



IML 2025 link : https://indico.cern.ch/event/1523250/contributions/6456344/attac hments/3070079/5431138/higgsml-3rd-place.pdf

1st (Tie) - Vienna's Institute of High Energy Physics (HEPHY)

UNBINNED MEASUREMENTS WITH REFINABLE SYSTEMATICS

Lisa Benato, Cristina Giordano, Claudius Krause, Ang Li, Robert Schöfbeck, Dennis Schwarz, Maryam Shooshtari, Daohan Wang











See more in Robert Schöfbeck (next) talk









1st (Tie) - Ibrahim Elsharkawy -University of Illinois at Urbana-Champaign





https://arxiv.org/abs/2505.08709

Contrastive Normalizing Flows



DNN Training with CNF

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



Learn more at: https://arxiv.org/abs/2505.08709

Fair Universe - Weak lensing Challenge

- NeurIPS 2025 Competition
- Goal to determine cosmological parameters (Ω_m and S_8), and uncertainties, from large simulated weak gravitational lensing dataset
- Various realistic systematic effects including baryonic effect, intrinsic alignment, photometric redshift uncertainty, shear measurement bias, point spread function, and source clustering effect
- Second phase test data with different (OoD) physical models



Shear maps γ_1 and γ_2 and weight map



101 different cosmological models

Conclusion

- Al challenge which addresses Systematic Uncertainty in HEP problem.
- New Scoring to take **Coverage** and **Confidence interval** into account.
- Large Computing Infrastructure as backend
- Exciting submissions which could be applied to future HEP analysis
- Large **Public** Data Set with ~280M Events (signal + background) in **Zenodo**
- New long term permanent benchmark will be launched Soon.

Ongoing information Google Group: <u>Fair-Universe-Announcements</u> Collaborations, questions, comments: <u>fair-universe@lbl.gov</u>

Thank you for your attention!





Bias and uncertainty in ML in HEP

- ML methods in HEP are often trained based on simulation which has estimated systematic uncertainties ("Z")
- These are then applied in data with the different detector state Z=?



• Common baseline approach: Train classifier on nominal data (e.g. Z=1) and estimate uncertainties with alternate simulations. Shift Z and look at impact or perform full profile likelihood



Background on Fair Universe Project

- 3 year US Dept. of Energy, AI for HEP project. Aims to:
 - Provide an open, large-compute-scale Al ecosystem for



- sharing datasets, training large models, fine-tuning those models, and **hosting challenges and benchmarks**.
- Organize a challenge series, progressively rolling in tasks of
 increasing difficulty, based on novel datasets.
- Tasks will focus on measuring and minimizing the effects of systematic uncertainties in HEP (particle physics and cosmology).
- This funding went to LBL, NERSC, U Washington, and Chalearn

Competition Flow



Large-compute-scale Al ecosystem ... hosting challenges and benchmarks.





Codabench Platform



Codabench

Codabench - open source platform for AI benchmarks and challenges

- Originally (CodaLab) Microsoft/Stanford now a Paris-Saclay/<u>LISN</u> led community
- > 600 challenges since 2013
- Completely open-ended competition design.
- Allows code submission as well as results e.g. for evaluation timing or reproducibility
- Also data-centric AI "inverted competitions"
- Queues for evaluation can run on diverse compute resources
- Platform itself can be deployed on different compute resources
- Ranked best challenge platform for ML by <u>ML contests</u>



Fair Universe Platform: Codabench-NERSC integration



System Specifications

Partition	# of nodes	CPU	GPU
GPU	1536	1x AMD EPYC 7763	4x <u>NVIDIA A100</u> (40GB)
	256	1x AMD EPYC 7763	4x <u>NVIDIA A100</u> (80GB)

3rd Place Solution Decision-Tree Aggregated Features and Hybrid Bin-Classifier/Quantile-Regressor



Yota Hashizume (hzume)

Overview

- 2-stage, GPU-free GBDT-based model
- First stage
 - Aggregated features
 - Used 2 models (1,2)
- Second stage
 - Estimate 68% confidence interval by using aggregated features
 - Used 2 models (3,4) and merged their outputs



https://indico.cern.ch/event/1523250/contributions/6456344/attachments/3070079/5431138/higgsml-3rd-place.pdf