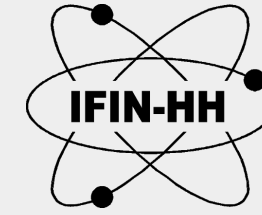# Event-Level Anomaly Detection for Multijet BSM Searches with Probabilistic Autoencoders

Ioan-Mihail Dinu, Julien Donini
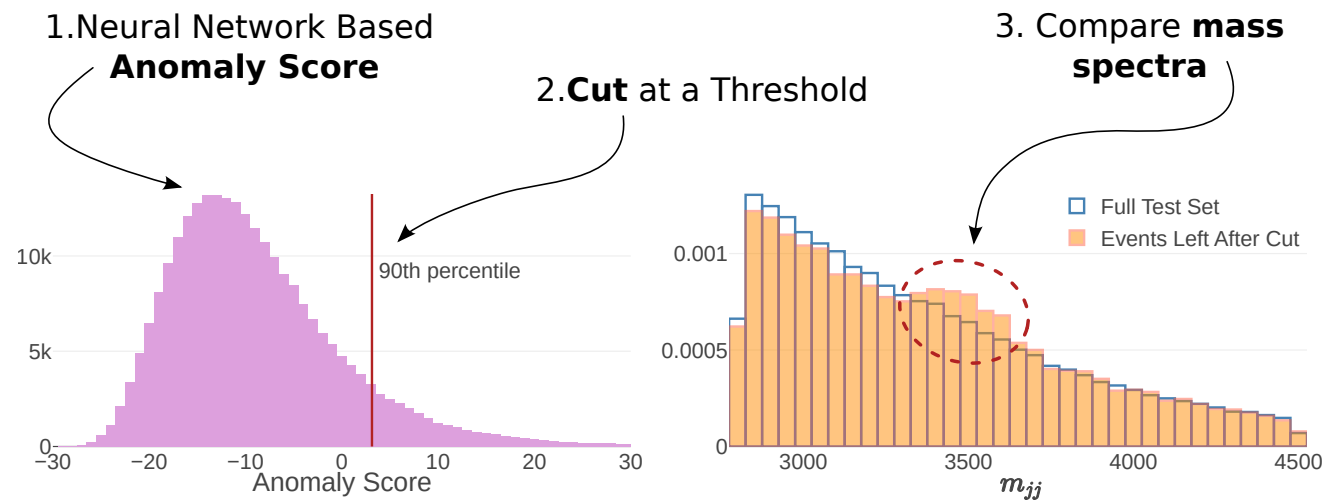
LPC, IFIN-HH

## Goals and Motivation

Standard BSM search techniques heavily rely on specific theory models. Rather than exploring all possible BSM models, maybe this could be done in a **model independent** way.

**Unsupervised Anomaly Detection:**
Makes **minimal assumptions on the signal model** (that it's quantifiably *different* from the background)
Only requires background data for training and it is **sensitive to low amounts of signal**

**BSM Search Strategy:**

1. Neural Network Based **Anomaly Score**

2. **Cut** at a Threshold

3. Compare **mass spectra**



## Machine Learning Anomaly Detection

**Autoencoders (AE)**

- Inputs are reconstructed from a learned **lower-dimensional** representation
- Trained to minimize reconstruction error
- High reconstruction error $\Rightarrow$ **anomalous events**



**Normalizing Flows (NF)**

- This model attempts to learn the **probability density function** of the data
- It uses a chain of **triangular maps** to create a **bijection** between the data space and a same-dimensional normal distribution

$$\mathbf{b}_\gamma \rightarrow \text{Triangular Map Chain} \qquad \mathbf{b}_\gamma(\vec{x}) = \vec{u}$$
$$\vec{x} \rightarrow \text{Target Distribution} \qquad \vec{u} \rightarrow \text{Gaussian Distribution}$$
$$\mathcal{J}_\gamma \rightarrow \text{Jacobian of Map Chain} \qquad p(\vec{x}) = p(\vec{u}) \det |\mathcal{J}_\gamma|^{-1}$$

## LHC Olympics Challenge Data [1]

**Data Format**: 4-vector particle flow information of multijet events simulated with Pythia and Delphes.
**Feature Extraction**: Jet kinematics, substructure variables or any other observables need to be computed and extracted by applying clustering algorithms **Datasets:**
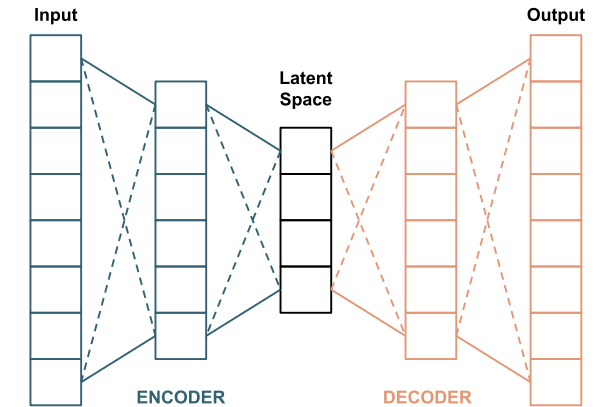
- RnD dataset: QCD background (1M), dijet signal (100k) and trijet signal (100k)
- Background-only training set (1M)
- 3 different black-boxes with potential signal (1M each)
  - BB1 : 3.8 TeV Z' decaying in dijet with 834 signal event
  - BB2 : QCD background only
  - BB3 : 4.2 gKK decaying in dijet and trijet (BR trijet = 0.625)

*Note: background is modeled differently across all datasets*

## Probabilistic Autoencoder [2]

**Combining Autoencoders and Normalizing Flows**

- Train a NF model on the **latent space** of an AE

- The **likelihood of the inputs** approximates to:

$$\ln p(\vec{x}) \approx -\frac{1}{2}||\vec{x} - \vec{x}'||^2 \vec{\sigma}^{\circ -2} - \frac{1}{2}b_\gamma(\vec{z})^2 + \ln|\det \mathcal{J}_\gamma|$$

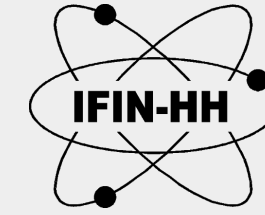$\vec{\sigma} \rightarrow$ average validation reconstruction error

# Event-Level Anomaly Detection for Multijet BSM Searches with Probabilistic Autoencoders

Ioan-Mihail Dinu, Julien Donini

LPC, IFIN-HH

## Results – Mass sculpting

### Bias Mitigation Strategies

▶ Input feature **uniformization**

▶ **Sample weights** based on $m_{jj}$ density.

★ Mass sculpting quantification:
⤳ **Jensen-Shannon Divergence** between test data and events passing the cut

★ Models trained on *QCD background data* and tested on *RnD dataset*
⇒ **Data-driven background model**

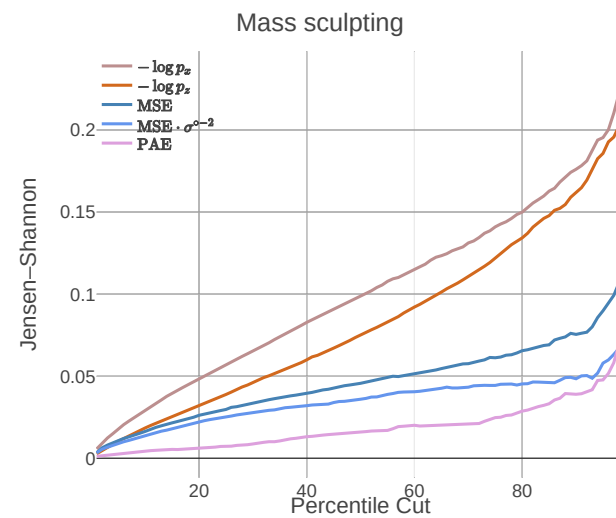### Performance Summary

$\log \mathbf{p_x}$: NF likelihood of inputs

$\log \mathbf{p_z}$: NF likelihood of latent representation
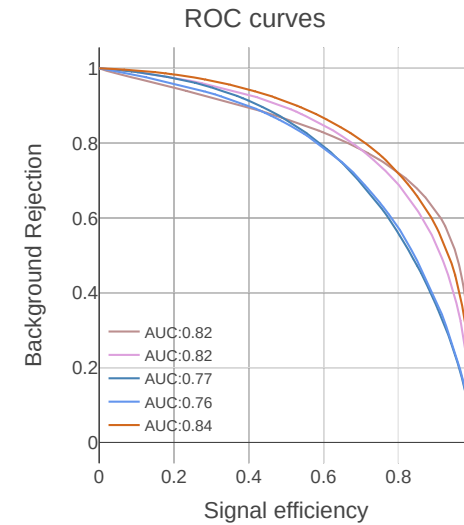
**MSE**: AE reconstruction error

**MSE**$\cdot\sigma^{\circ-2}$: MSE normalized to average validation reconstruction error

**PAE**: approximation of input likelihood with PAE



Mass sculpting

## References

[1] Gregor Kasieczka, Benjamin Nachman, David Shih, Oz Amram, Anders Andreassen, Kees Benkendorfer, Blaz Bortolato, Gustaaf Brooijmans, Florencia Canelli, Jack H. Collins, Biwei Dai, Felipe F. De Freitas, Barry M. Dillon, Ioan-Mihail Dinu, Zhongtian Dong, Julien Donini, Javier Duarte, D. A. Faroughy, Julia Gonski, Philip Harris, Alan Kahn, Jernej F. Kamenik, Charanjit K. Khosa, Patrick Komiske, Luc Le Pottier, Pablo Martín-Ramiro, Andrej Matevc, Eric Metodiev, Vinicius Mikuni, Inês Ochoa, Sang Eon Park, Maurizio Pierini, Dylan Rankin, Veronica Sanz, Nilai Sarda, Urous Seljak, Aleks Smolkovic, George Stein, Cristina Mantilla Suarez, Manuel Szewc, Jesse Thaler, Steven Tsan, Silviu-Marian Udrescu, Louis Vaslin, Jean-Roch Vlimant, Daniel Williams, and Mikaeel Yunus.
The lhc olympics 2020: A community challenge for anomaly detection in high energy physics, 2021.

[2] Vanessa Böhm and Uroš Seljak.
Probabilistic auto-encoder, 2020.

[3] Georgios Choudalakis.
On hypothesis testing, trials factor, hypertests and the bumphunter, 2011.

## Future outlook

▶ performance studies on **3-prong signals**

▶ readjust method for **jet images** inputs

## Results – Performance



ROC curves

★ NF likelihoods show good discrimination but are the most biased

★ Reconstruction error score compromises some performance for less bias

⇒ Combining the two scores results in **good performance and reduced bias**

★ Pae fitted on *Background-only training set*



Distributions with bump

### Black Box 1 Dataset

$Z' \to X, Y$
$m_{Z'} = 3.8\,TeV$
$m_X = 732\,GeV \quad m_Y = 378\,GeV$
$834/1M$ signal events

### Bump hunting [3] results

mean : $3866.7\,GeV$
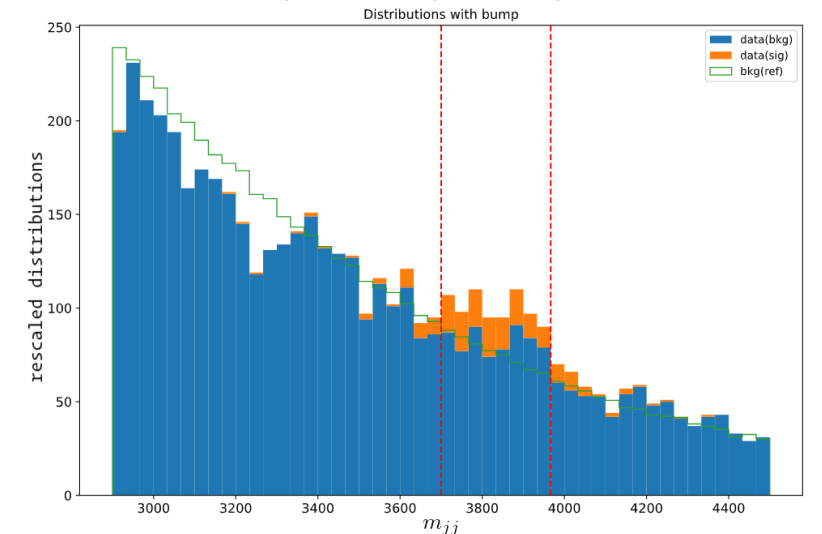width : $160\,GeV$
number of signal events : $118$

**signal bump significance** $\gg 5\sigma$

## Conclusions

**Autoencoder** and **Normalizing Flow** neural networks tested on LHC Olympics Data for anomaly detection. Combining the two networks:

⇒ **Probabilistic Autoencoder** ensemble → successful unsupervised anomaly detection application to jet physics:

★ Sensitive to very low signal fractions

★ Low bias when using **Mitigation Strategies**

★ Allows for **successful Bump Hunting** on LHCO Black Box 1