

Anomaly detection in 5 mins or less

Snowmass 2022 CompF3 ML Session

David Shih
July 19, 2022



RUTGERS
THE STATE UNIVERSITY
OF NEW JERSEY

ML for New Physics Searches

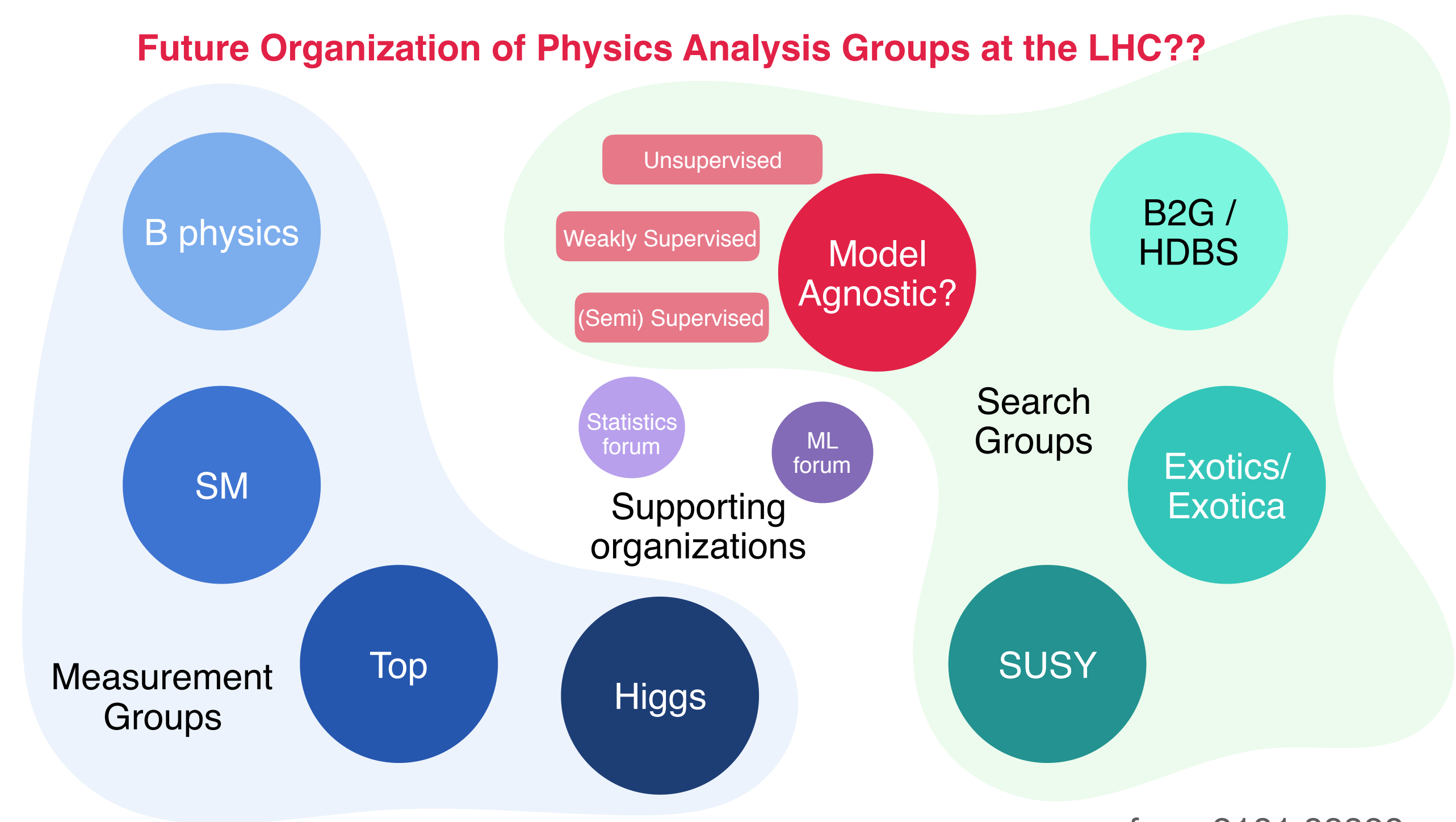


The vast majority of LHC searches for new physics are very model specific

ML for New Physics Searches



The vast majority of LHC searches for new physics are very model specific

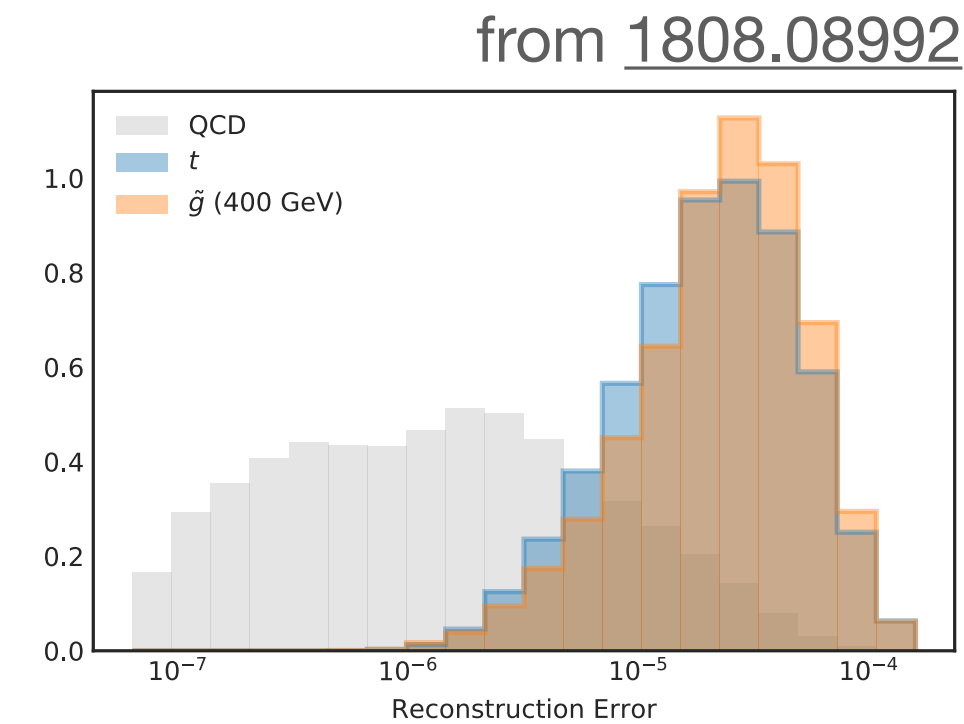
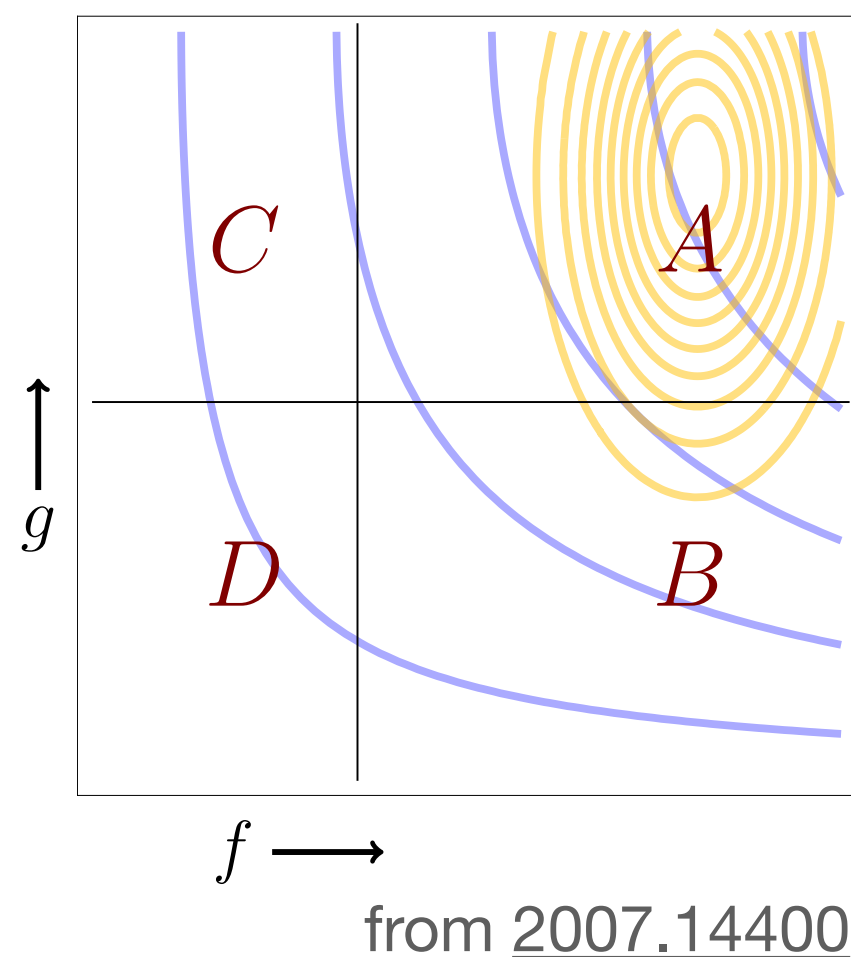


from [2101.08320](#)

Why aren't there more model-agnostic new physics searches?

ML for Anomaly Detection

- How do we search for new physics in a model-agnostic way?
- Need two ingredients:
 1. Signal sensitivity (anomaly score)
 2. Background estimation

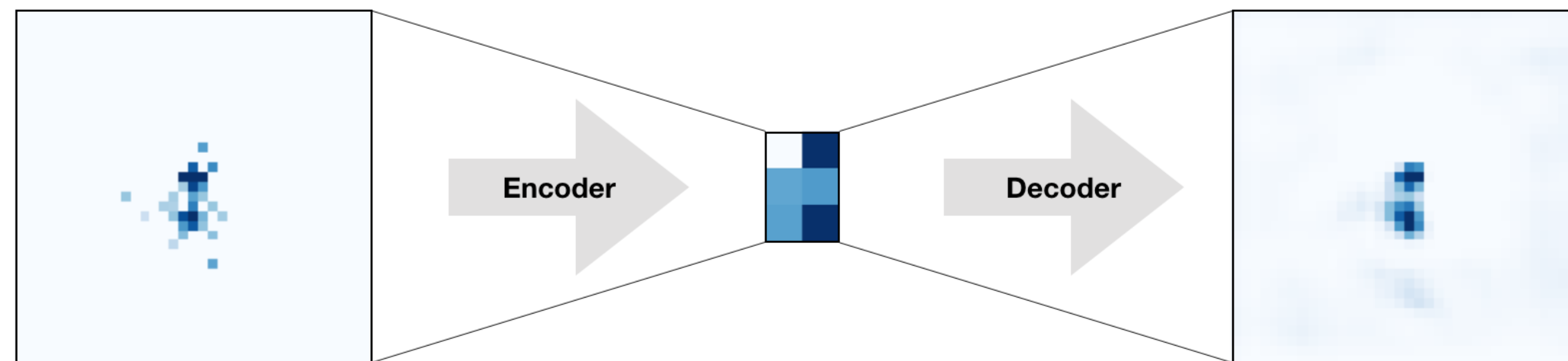


Both should be data-driven for a truly model-agnostic search

ML for Anomaly Detection

Types of Anomaly Scores

- Low $p(x)$ — outliers



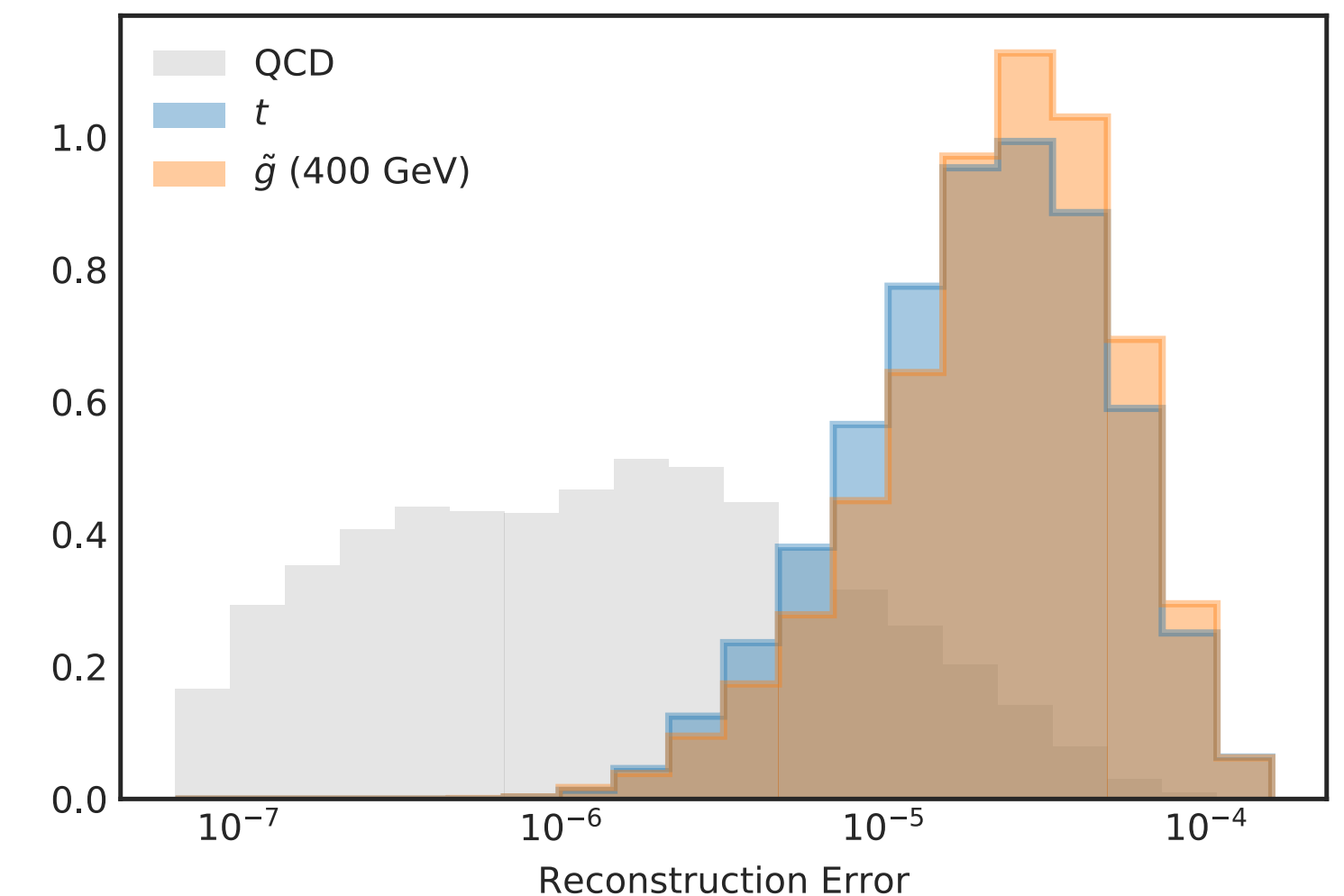
Autoencoders

Fully unsupervised

Farina, Nakai & **DS** [1808.08992](#)

Heimel et al [1808.08979](#)

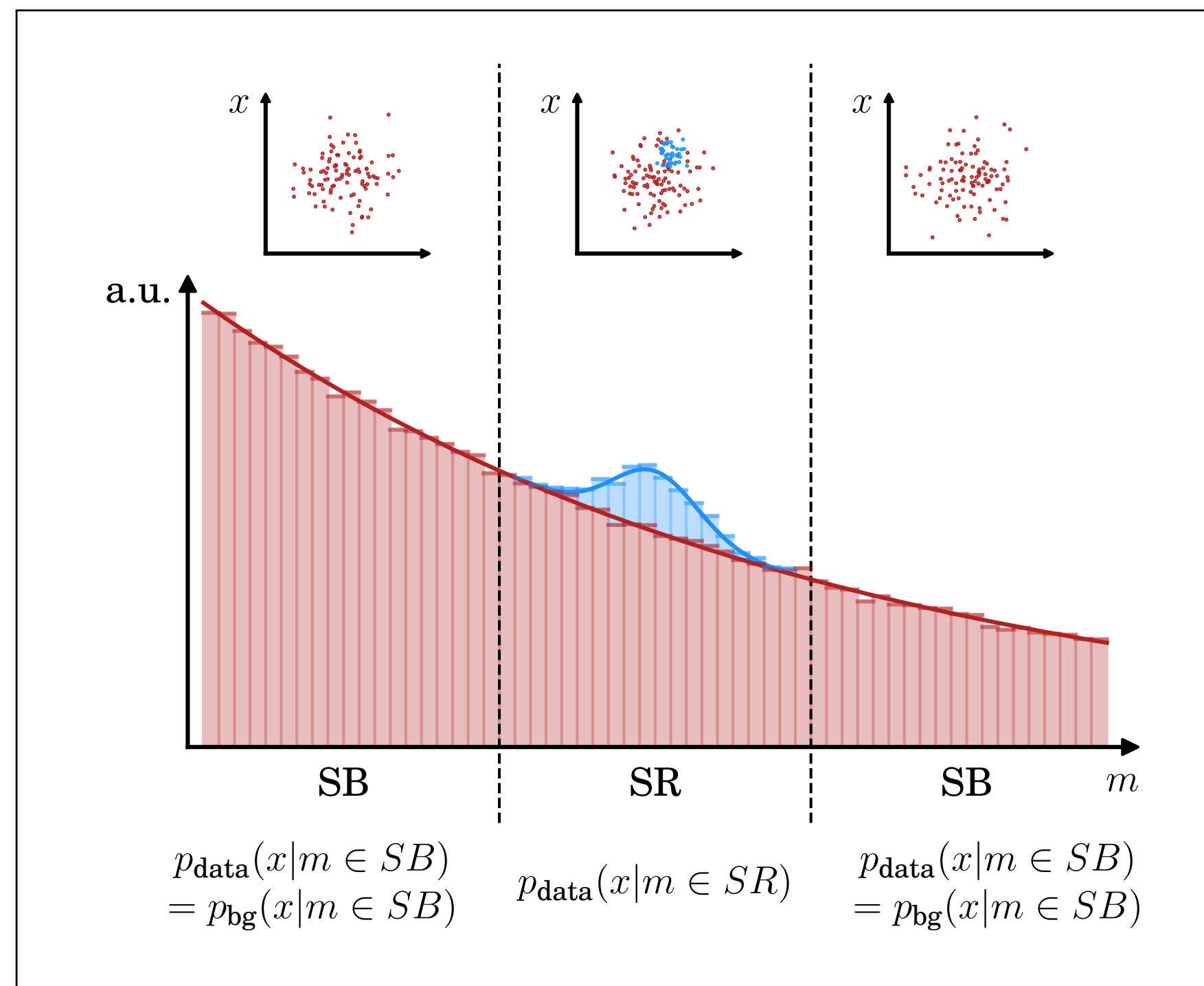
and many more!!



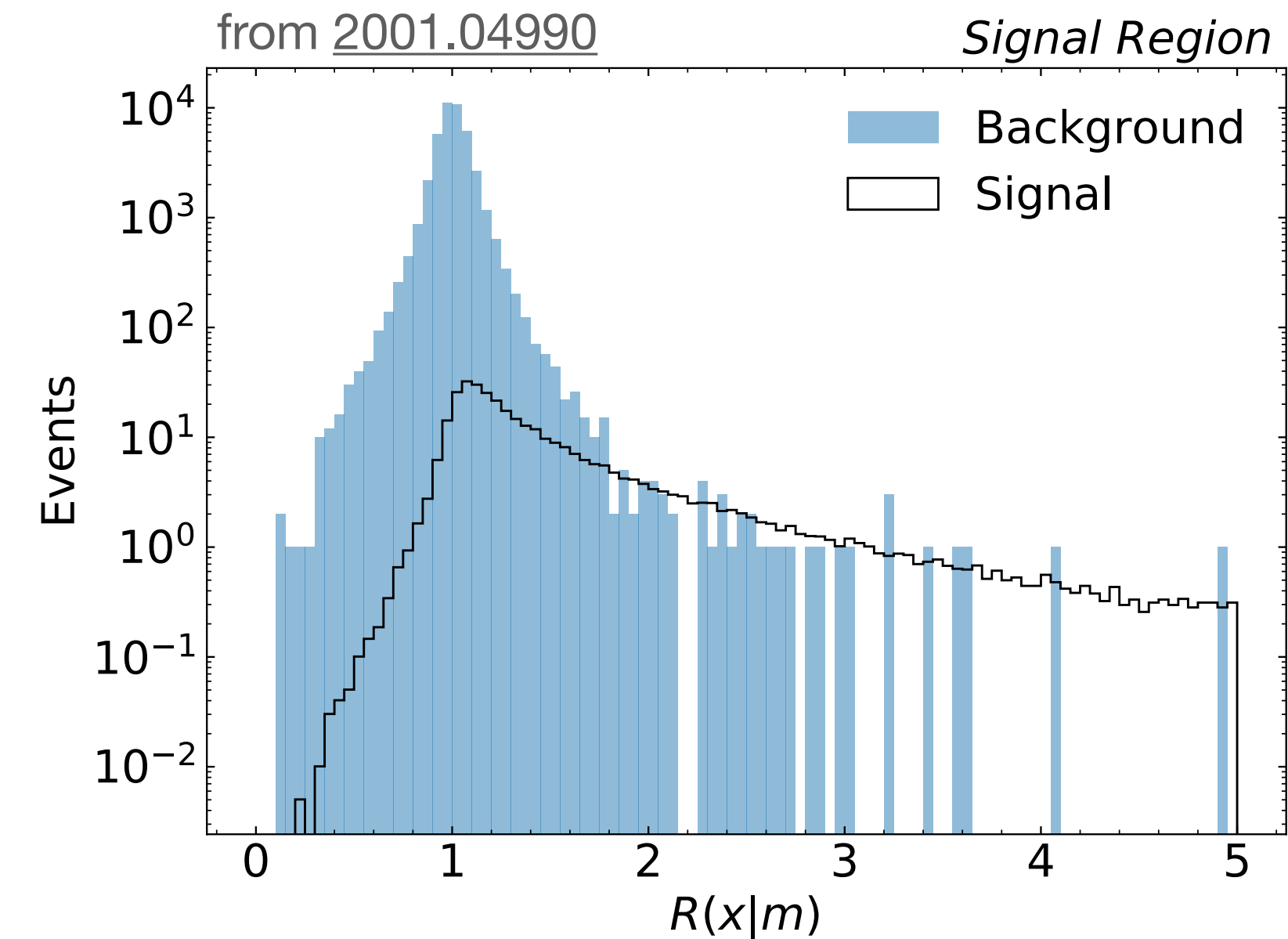
ML for Anomaly Detection

Types of Anomaly Scores

- High $p_{data}(x)/p_{bg}(x)$ – overdensities



from [2109.00546](#)



Enhanced bump hunts

Weakly supervised

CWoLa Hunting [Collins, Howe & Nachman [1805.02664](#), [1902.02634](#)]

ANODE [Nachman & **DS** [2001.04990](#)]

CATHODE [Hallin et al [2109.00546](#)]

CURTAINS [Raine et al [2203.09470](#)]

and more...

LHC Olympics 2020

The LHC Olympics 2020

A Community Challenge for Anomaly
Detection in High Energy Physics

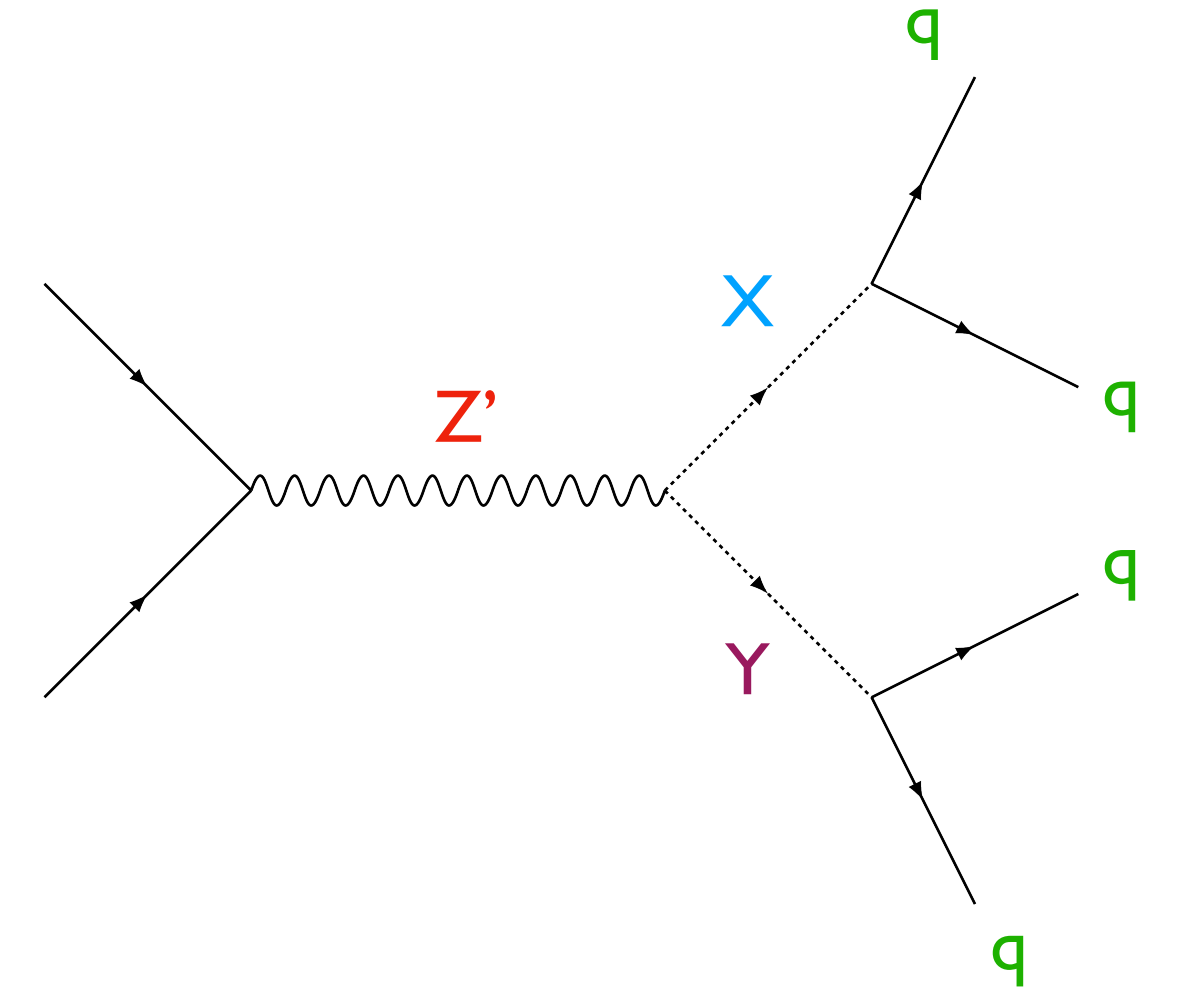


Gregor Kasieczka (ed),¹ Benjamin Nachman (ed),^{2,3} David Shih (ed),⁴ Oz Amram,⁵ Anders Andreassen,⁶ Kees Benkendorfer,^{2,7} Blaz Bortolato,⁸ Gustaaf Brooijmans,⁹ Florencia Canelli,¹⁰ Jack H. Collins,¹¹ Biwei Dai,¹² Felipe F. De Freitas,¹³ Barry M. Dillon,^{8,14} Ioan-Mihail Dinu,⁵ Zhongtian Dong,¹⁵ Julien Donini,¹⁶ Javier Duarte,¹⁷ D. A. Faroughy,¹⁰ Julia Gonski,⁹ Philip Harris,¹⁸ Alan Kahn,⁹ Jernej F. Kamenik,^{8,19} Charanjit K. Khosa,^{20,30} Patrick Komiske,²¹ Luc Le Pottier,^{2,22} Pablo Martín-Ramiro,^{2,23} Andrej Matevc,^{8,19} Eric Metodiev,²¹ Vinicius Mikuni,¹⁰ Inês Ochoa,²⁴ Sang Eon Park,¹⁸ Maurizio Pierini,²⁵ Dylan Rankin,¹⁸ Veronica Sanz,^{20,26} Nilai Sarda,²⁷ Uroš Seljak,^{2,3,12} Aleks Smolkovic,⁸ George Stein,^{2,12} Cristina Mantilla Suarez,⁵ Manuel Szewc,²⁸ Jesse Thaler,²¹ Steven Tsan,¹⁷ Silviu-Marian Udrescu,¹⁸ Louis Vaslin,¹⁶ Jean-Roch Vlimant,²⁹ Daniel Williams,⁹ Mikaeel Yunus¹⁸

<https://arxiv.org/abs/2101.08320>

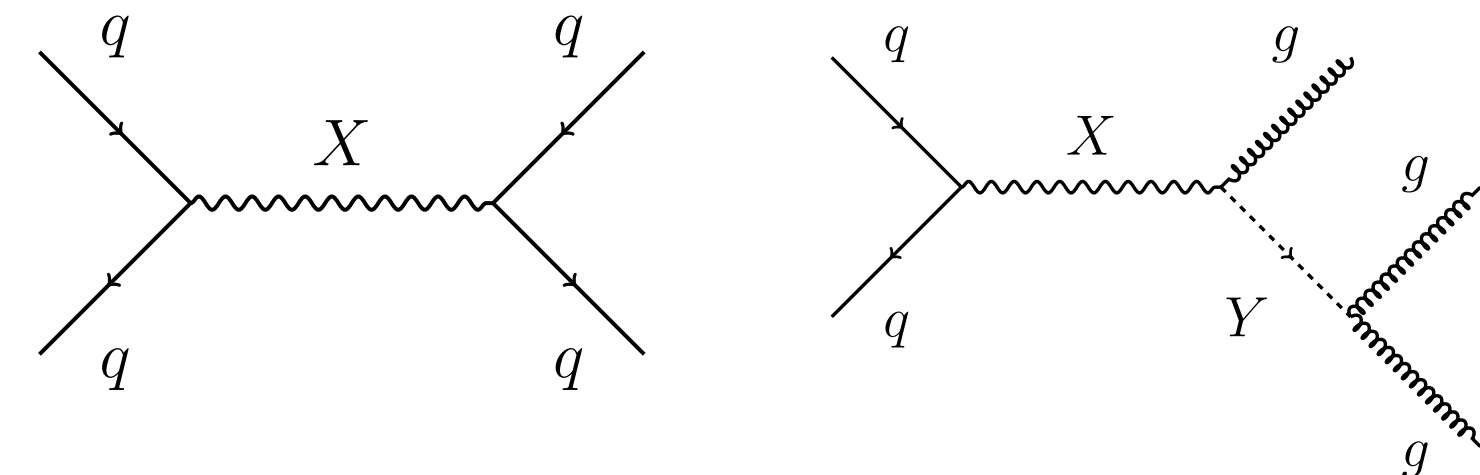
R&D dataset:

bg: 1M QCD dijet events
signal: up to 100k $Z' \rightarrow XY$ events
Pythia+Delphes
 $pT(J1) > 1.2$ TeV trigger



3 “Black Box” datasets

1. 1M QCD dijets + 834 $Z' \rightarrow XY$ signal (same topology as R&D, different masses)
2. No signal
3. QCD dijets + 3,000 Z' decaying to dijets or trijets



LHC Olympics 2020

3 “Black Box” datasets

LHC Olympics 2020

3 “Black Box” datasets

1. 1M QCD dijets + 834 $Z' \rightarrow XY$ signal (same topology as R&D, different masses)

LHC Olympics 2020

3 “Black Box” datasets

1. 1M QCD dijets + 834 $Z' \rightarrow XY$ signal (same topology as R&D, different masses)
 - **Several successful methods! (based on autoencoders, CWoLa, density estimation...)**

LHC Olympics 2020

3 “Black Box” datasets

1. 1M QCD dijets + 834 $Z' \rightarrow XY$ signal (same topology as R&D, different masses)
 - **Several successful methods! (based on autoencoders, CWoLa, density estimation...)**
2. No signal

LHC Olympics 2020

3 “Black Box” datasets

1. 1M QCD dijets + 834 $Z' \rightarrow XY$ signal (same topology as R&D, different masses)
 - **Several successful methods! (based on autoencoders, CWoLa, density estimation...)**
2. No signal
 - **Some approaches found false positives — importance and challenges of background estimation!**

LHC Olympics 2020

3 “Black Box” datasets

1. 1M QCD dijets + 834 $Z' \rightarrow XY$ signal (same topology as R&D, different masses)
 - **Several successful methods! (based on autoencoders, CWoLa, density estimation...)**
2. No signal
 - **Some approaches found false positives — importance and challenges of background estimation!**
3. QCD dijets + 3,000 Z' decaying to dijets or trijets

LHC Olympics 2020

3 “Black Box” datasets

1. 1M QCD dijets + 834 $Z' \rightarrow XY$ signal (same topology as R&D, different masses)
 - **Several successful methods! (based on autoencoders, CWoLa, density estimation...)**
2. No signal
 - **Some approaches found false positives — importance and challenges of background estimation!**
3. QCD dijets + 3,000 Z' decaying to dijets or trijets
 - **No approaches discovered the signal in BB3**

Outlook

- There is a lot of community interest in anomaly detection and model-agnostic NP searches!
- LHC Olympics 2020 was a very successful challenge, drawing nearly 50 participants from theory, experiment, and beyond (cosmology, computer science)
- Proofs-of-concept are beginning to be ported over to real data

CWoLa Hunting

ATLAS, PRL **125** 131801 (2020)

RNN VAE

ATLAS-CONF-2022-045

- **Many challenges for future R&D, including: feature selection, background estimation, multiple decay modes (BB3), non-resonant signals, ...**