

# LHC Olympics 2020

Gregor Kasieczka

[gregor.kasieczka@uni-hamburg.de](mailto:gregor.kasieczka@uni-hamburg.de)

Twitter: [@GregorKasieczka](https://twitter.com/GregorKasieczka)

PHYSTAT-Anomalies, May 24/25 2022

**CLUSTER OF EXCELLENCE**  
QUANTUM UNIVERSE



Universität Hamburg

DER FORSCHUNG | DER LEHRE | DER BILDUNG



Partnership of  
Universität Hamburg and DESY



Bundesministerium  
für Bildung  
und Forschung



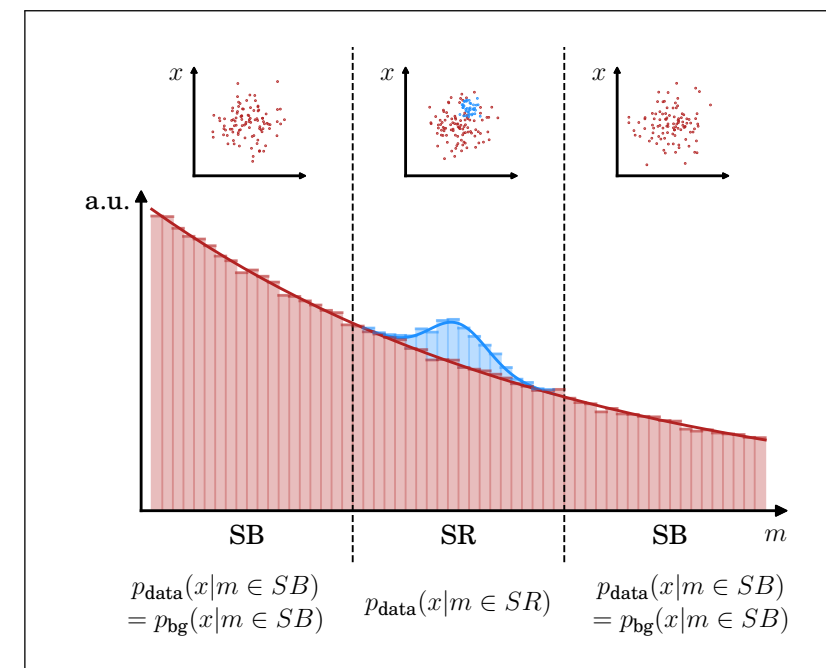
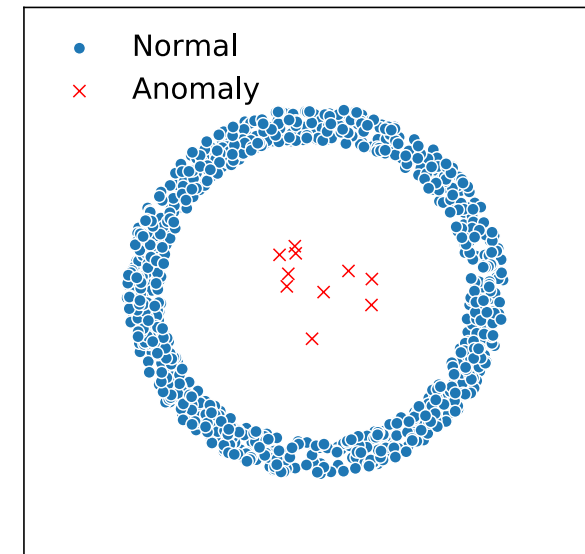
# Overview

- At this point in the workshop, do not need to motivate searching for anomalies
- Introduce LHC Olympics: Ideas behind, setup, methods, and results
- Some other comments on open issues for anomaly detection

*(With many thanks to Ben & David as LHCO co-organizers and all participants!!)*

# Types of anomalies

- **Outliers/Point anomalies:** Datapoints far away from regular distribution
- Examples:
  - Detector malfunctions
  - Background-free search
- **Group anomalies:** Individual examples not interesting, but signal is an overdensity with respect to background
- Examples:
  - Resonance searches
  - Transient signals in time series
- ***Focus of LHC Olympics***



# Approaches

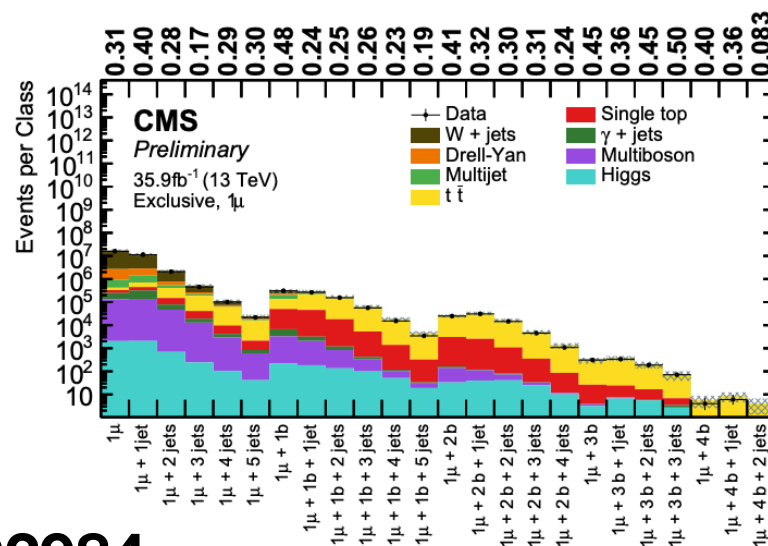
## Use simulation to estimate backgrounds?

Yes

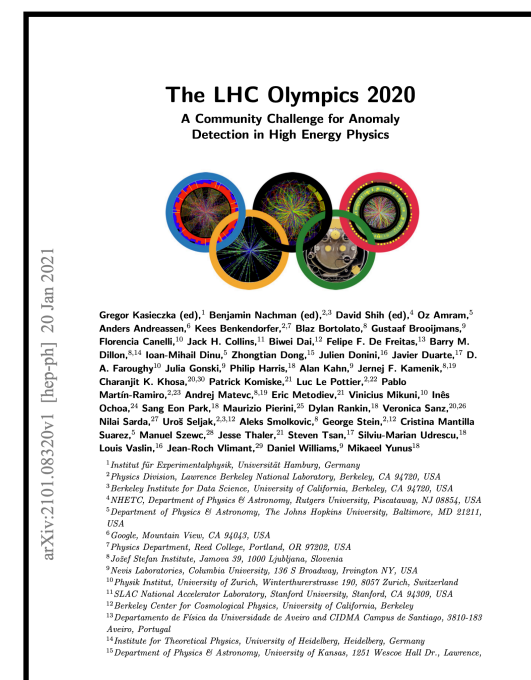
No

- Systematically compare simulation and recorded data, look for differences
- Con: Relies on imperfect simulation, Maximally background model dependent
- Pro: Sensitive to all types of anomalies
- e.g. MUSIC

- Estimate background from data
- Con: Need to make assumptions about signal model
- Pro: No reliance on simulation
- Focus of LHC Olympics



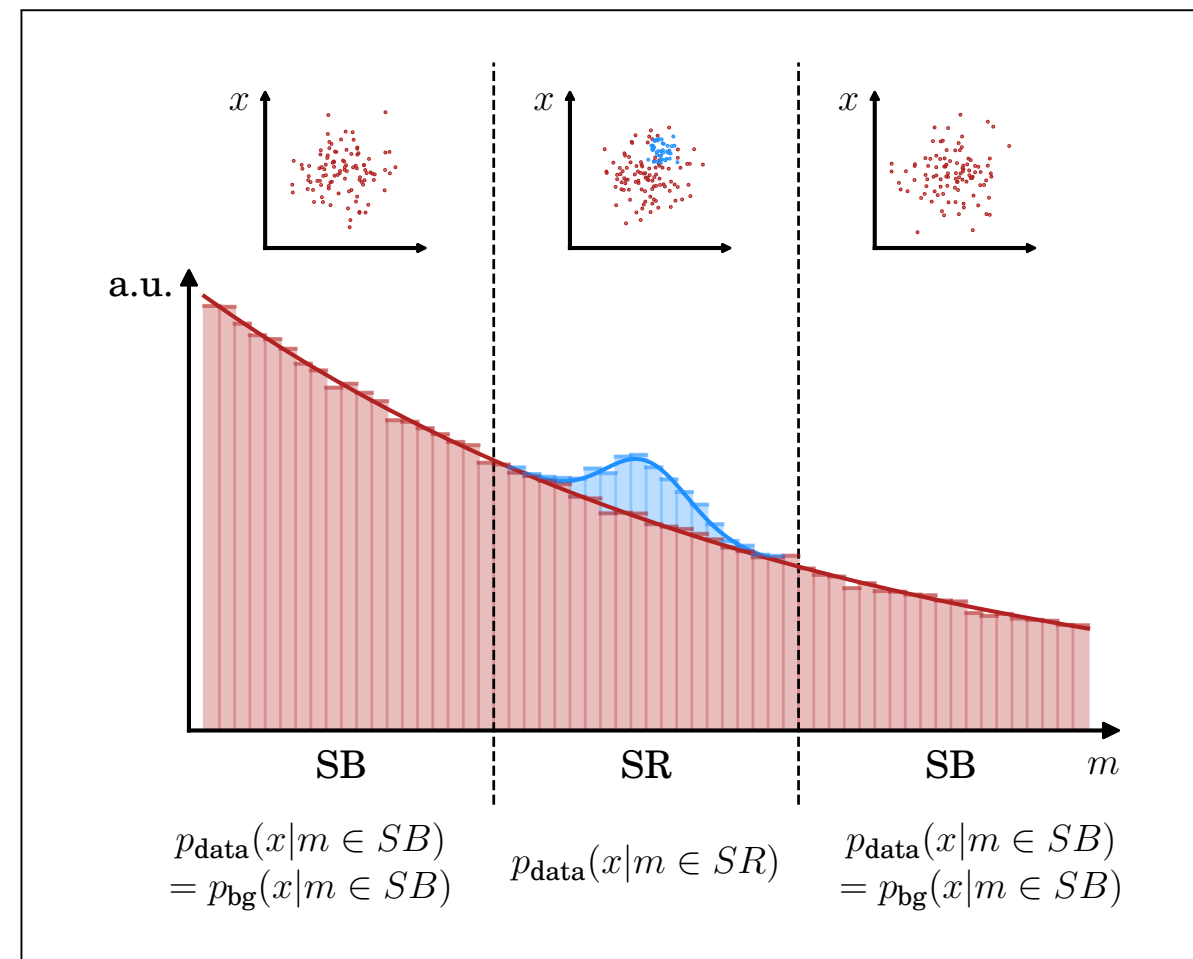
arXiv:  
2010.02984





# Assumptions

- **Rarity:**  $\Pr(\text{anomaly}) \ll \Pr(\text{normal})$
- **Overlap:**  $\max_x p(x|\text{anomaly})/p(x|\text{normal}) < \infty$
- **Resonance:**  $\Pr(|m - m_0| > \delta | \text{anomaly}) \approx 0$  for some feature  $m$  (often a mass) and fixed  $m_0, \delta$
- **Smoothness:**  $p(x|m, \text{normal})$  varies slowly with  $m$  so that one can use data with  $|m - m_0| > \delta$  to estimate  $p(x|m, \text{normal})$  for  $|m - m_0| < \delta$



# Overall Strategy

- Define an anomaly score  $a$ 
  - Should be high for anomalous (signal-like) and low for background-like data
- Use a selection  $a$  to create an anomaly-enriched dataset
- Estimate background from data
- Compare predicted background to number of observed with high  $a$  (+potentially other selection criteria)

***Idea behind LHC Olympics: Provide dataset that allows exercising all these steps***



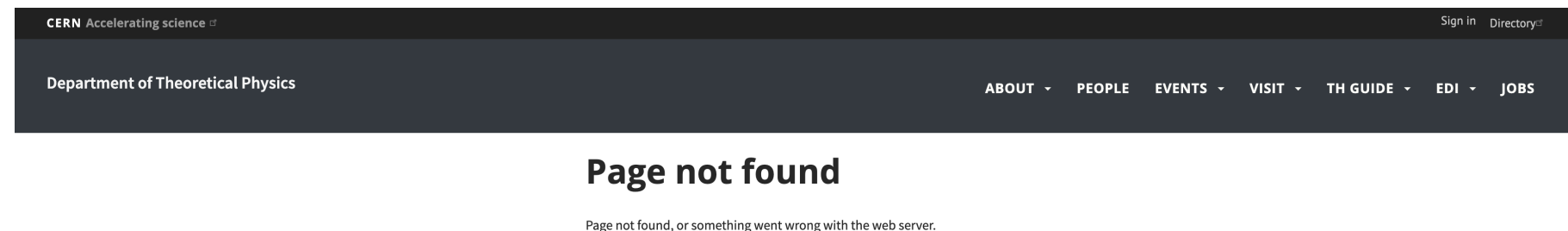
# History

- Ancient LHC Olympics
  - Four workshops 2005—2007 / prior starting to LHC
  - Assumed new physics would be plenty, focus on characterisation of new physics models



First LHCO lost to sands of time:

<http://ph-dep-th.web.cern.ch/ph-dep-th/content2/workshops/lhcOlympics/lhcolympicsI.html>



4th LHC



Workshop

Princeton, 22 March, 2007

## Information

The LHC Olympics is a collective effort by theorists to train themselves in establishing a correspondence between theoretical models and experimental signatures using collider simulation. This is done by developing and distributing user-friendly versions of simulation and data analysis tools, and via a series of black box exercises, in which participants are challenged to disentangle simulated LHC data sets. The first three LHC Olympics workshops were held at CERN in [July 2005](#) and [February 2006](#), and the KITP in [August 2006](#).

The fourth LHC Olympics meeting will take place during the workshop "[Physics at LHC: From Experiment to Theory](#)," held from March 21-24, 2007, in Princeton. The workshop will bring together high energy physicists with expertise ranging from formal theory to collider phenomenology, simulation and experiment. The meeting is held jointly with the second workshop in the "[Monte Carlo Tools for Beyond the Standard Model Physics](#)" (MC4BSM) series.

## Data Challenge

Using event generators in combination with the [PGS-4](#) detector simulation, a number of black box data sets have been generated. The boxes are presented as a data challenge to participants, who are invited to try to figure out the underlying theoretical model. This round of the LHCO has two new black box data sets.

A) The [SLAC LHC Olympics Black Box](#) contains signal only -- no background -- for 20 fb<sup>-1</sup> of LHC data. The black box' creators are John Conley, Michael Peskin, and Tommer Wizansky.

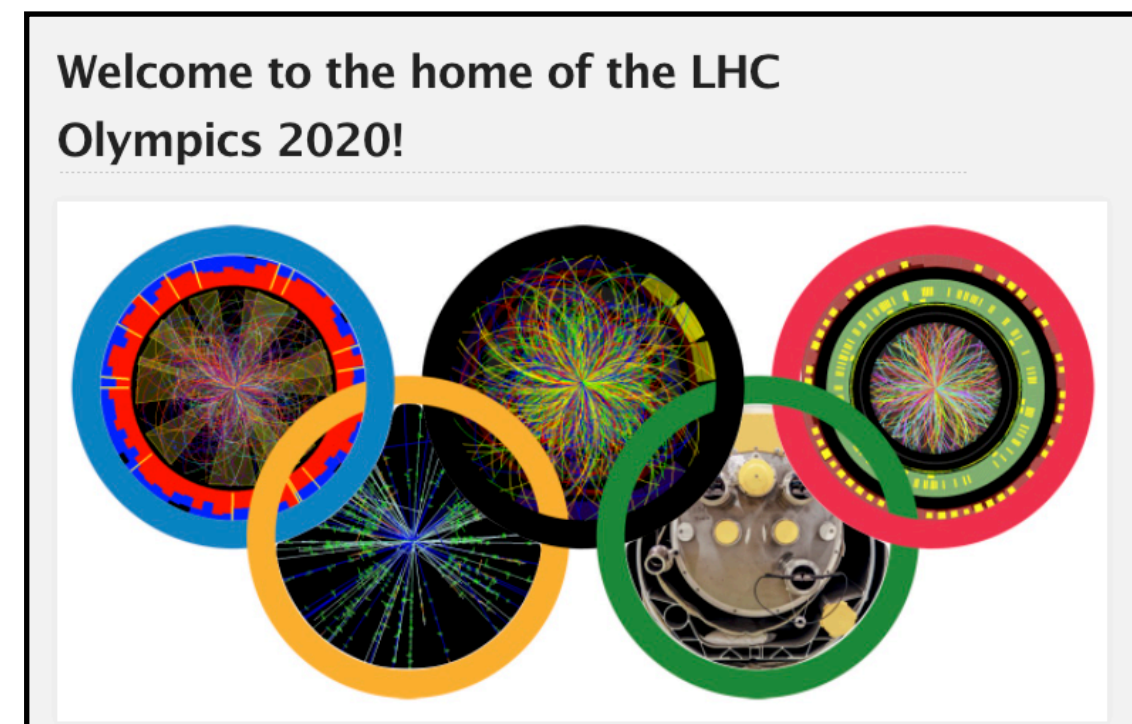
B) The [Cornell/Harvard LHC Olympics Black Box](#) represents about 10 fb<sup>-1</sup> of new physics data. This black box has been created by Patrick Meade, Peter Onyisi, Maxim Perelstein, and Matt Reece.

An explanation of the LHCO format and suggestions for how to get started on analyzing them can be found on the black box web pages the [LHCO wiki](#). The login for the wiki is "olympian" and password is "blackbox." All LHCO participants are free to use the wiki to add comments, ask questions, make suggestions, etc.

The programs used to generate the boxes are in continuing development. To keep apprised of possible updates to the boxes, it is advisable to regularly check the BB webpages and the [LHCO wiki](#) for announcements.

# Motivation

- Encourage development and comparison of model-agnostic search strategies
  - Focus on group anomalies, data-driven searches
- Provide a complete package, balance details vs accessibility
- Datasets:
  - One R&D dataset for algorithm development
  - Three black box datasets (BB1-BB3)
    - Unblinded over time
- Timeline:
  - Spring 2019: Release R&D dataset ([link](#))
  - Autumn 2019: Release BB datasets ([link](#))
  - January 2020: Winter Olympics as part of ML4Jets, unblinding of BB1 ([link](#))
  - July 2020: (Virtual) Summer Olympics, unblinding of BB2 and BB3 ([link](#))
  - LHC Olympics paper (<https://arxiv.org/abs/2101.08320>) public



<https://lhco2020.github.io/homepage/>



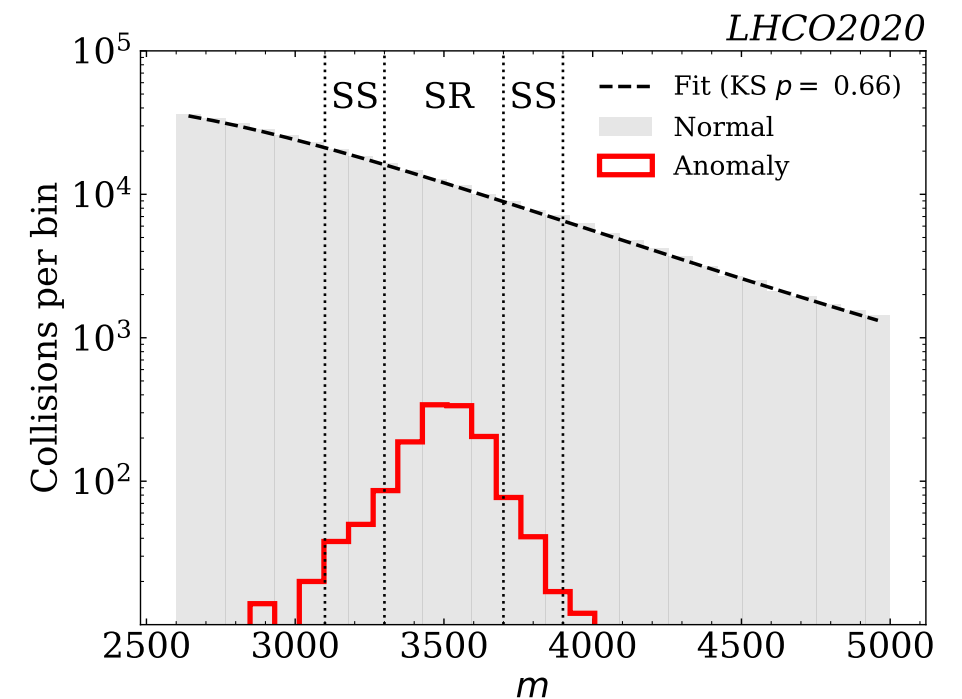
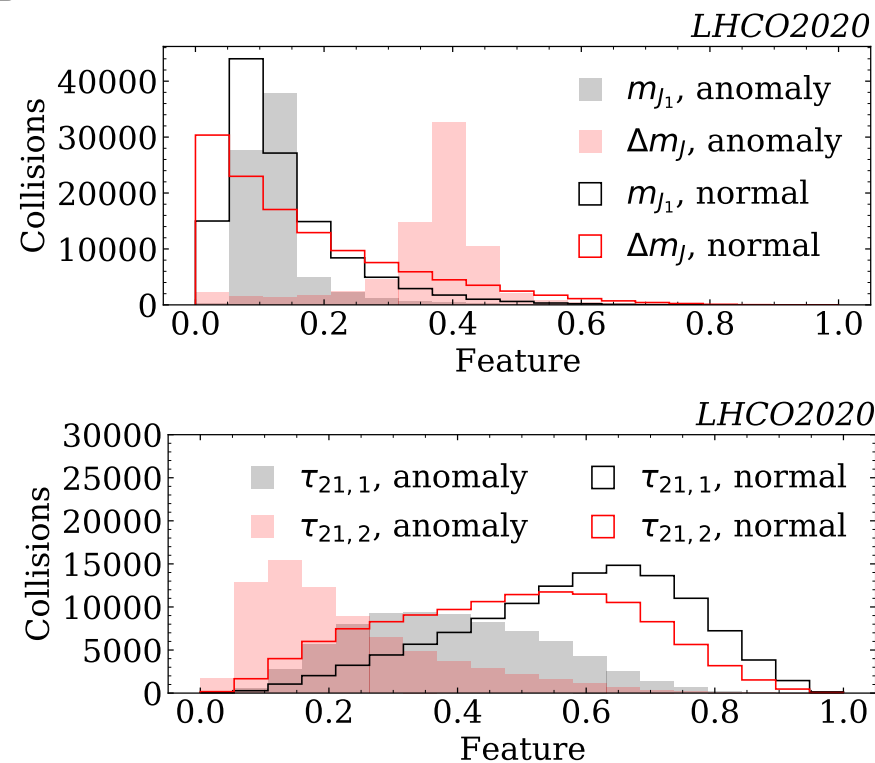
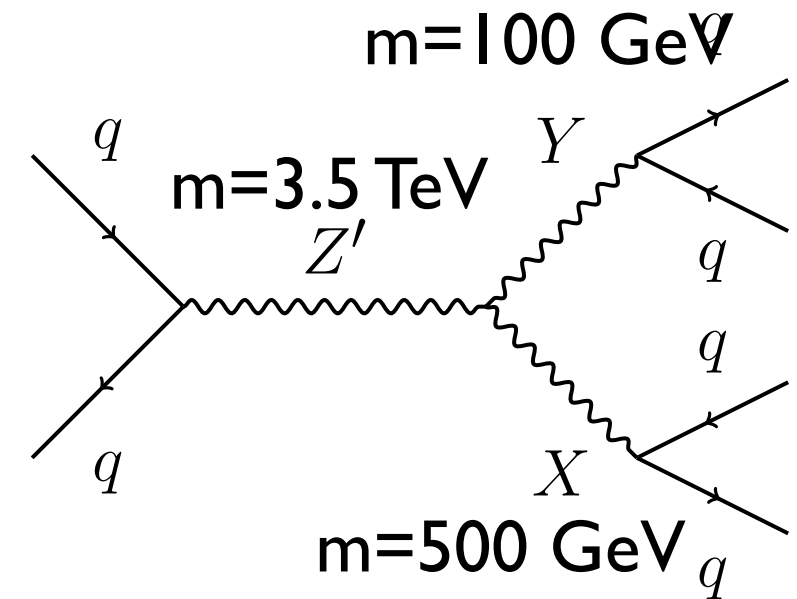
# Datasets

- Data format:
  - 3-vectors of reconstructed particles in the event
    - pt, eta, phi, (m assumed to be zero)
    - Leading 700 particles (zero-passed otherwise)
    - -> **2100 dimensional input space**
- Also provided a much lower dimensional representation
  - Clustering into two jets and using mass/substructure
  - O(10) dimensional input space
- No other quantities (e.g. flavor tagging) included
- Single R=1 jet trigger  $p_T > 1.2$  TeV
- Generation with Pythia/Herwig; detector simulation with Delphes

Setting	R&D	BB1	BB3
Tune:pp	14	3	10
PDF:pSet	13	12	5
TimeShower:alphaSvalue	0.1365	0.118	0.16
SpaceShower:alphaSvalue	0.1365	0.118	0.16
TimeShower:renormMultFac	1	0.5	2
SpaceShower:renormMultFac	1	0.5	2
TimeShower:factorMultFac	1	1.5	0.5
SpaceShower:factorMultFac	1	1.5	0.5
TimeShower:pTmaxMatch	1	2	1
SpaceShower:pTmaxMatch	0	2	1

# R&D dataset

- For building and testing methods
- 1M background examples (Standard Model),  
100k signal examples (signal, see Feynman diagram on the right)
- Labels provided
- Relatively simple signal
  - Known to differ in previously mentioned features from background distribution
- Unrealistically high S/B



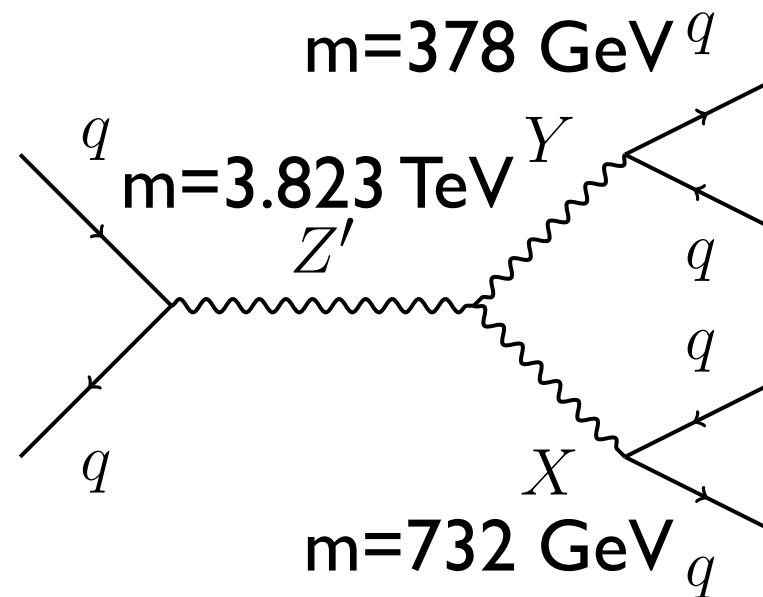
2107.02821



# Challenge datasets: BB1

- All contain total of 1M examples; might contain signal; no labels provided during 'content' phase (labels available no)
- All used different simulation parameters for background (to avoid unrealistic exploits)
- BB1: 834 signal examples  
Same event topology as R&D dataset, different masses

*might be easy?*



# Challenge datasets: BB2

- All contain total of 1M examples; might contain signal; no labels provided during 'content' phase (labels available no)
- All used different simulation parameters for background (to avoid unrealistic exploits)
- Additional pure-background sample provided (again with a different tune)
- BB2: 0 signal examples; Herwig++ instead of Pythia for background

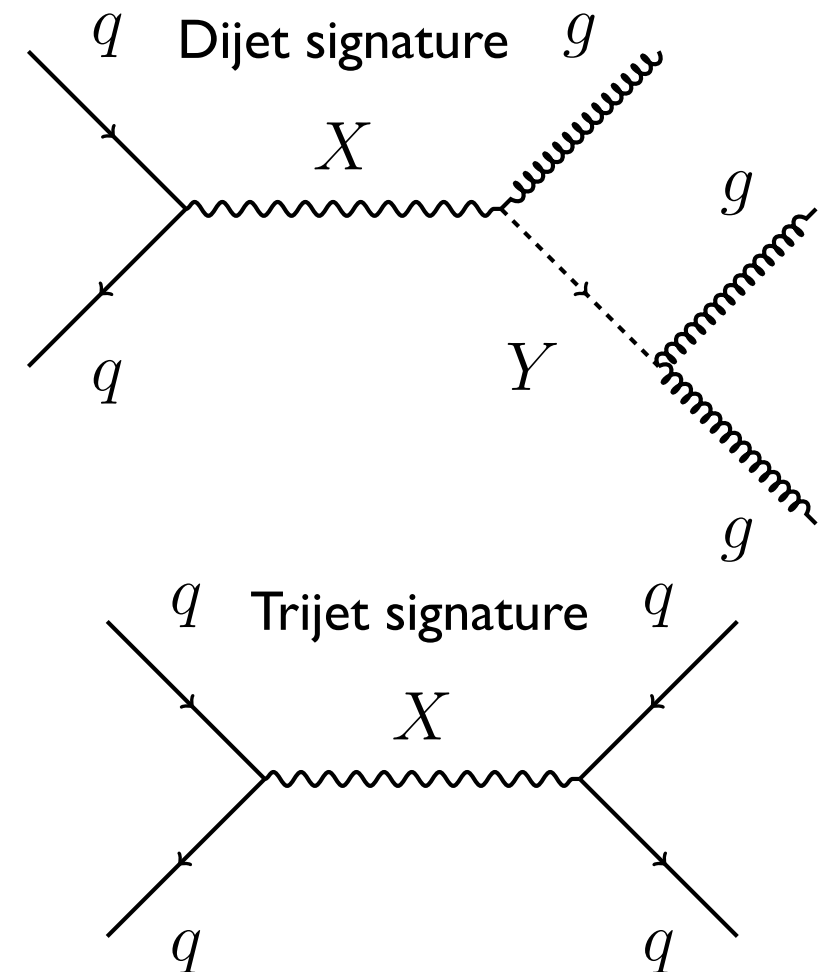
***test for false positives***



# Challenge datasets: BB3

- All contain total of 1M examples; might contain signal; no labels provided during ‘content’ phase (labels available no)
- All used different simulation parameters for background (to avoid unrealistic exploits)
- Additional pure-background sample provided (again with a different tune)
- BB3:  
mX = 4.2 TeV and two decay modes:  
1200 signal events in di-jet signature  
2000 signal events in tri-jet signature  
(finding individual excess should not yield significance)

***should be challenging***



# Evaluation criteria

- **What you should report:**
  - A p-value associated with the dataset having no new particles (null hypothesis).
  - As complete a description of the new physics as possible. For example: the masses and decay modes of all new particles (and uncertainties on those parameters).
  - How many signal events (+uncertainty) are in the dataset (before any selection criteria).
  - Partial submissions in only a subset of the categories are welcome!

*(Goal not to necessarily pick ‘one winner’ but to get a useful understanding of anomaly detection capabilities)*

# Overview of Methods

## 3 Unsupervised

- 3.1 Anomalous Jet Identification via Variational Recurrent Neural Network
- 3.2 Anomaly Detection with Density Estimation
- 3.3 BuHuLaSpa: Bump Hunting in Latent Space
- 3.4 GAN-AE and BumpHunter
- 3.5 Gaussianizing Iterative Slicing (GIS): Unsupervised In-distribution Anomaly Detection through Conditional Density Estimation
- 3.6 Latent Dirichlet Allocation
- 3.7 Particle Graph Autoencoders
- 3.8 Regularized Likelihoods
- 3.9 UCluster: Unsupervised Clustering

*(No labels)*

## 4 Weakly Supervised

- 4.1 CWoLa Hunting
- 4.2 CWoLa and Autoencoders: Comparing Weak- and Unsupervised methods for Resonant Anomaly Detection
- 4.3 Tag N' Train
- 4.4 Simulation Assisted Likelihood-free Anomaly Detection
- 4.5 Simulation-Assisted Decorrelation for Resonant Anomaly Detection

*(Noisy labels)*

## 5 (Semi)-Supervised

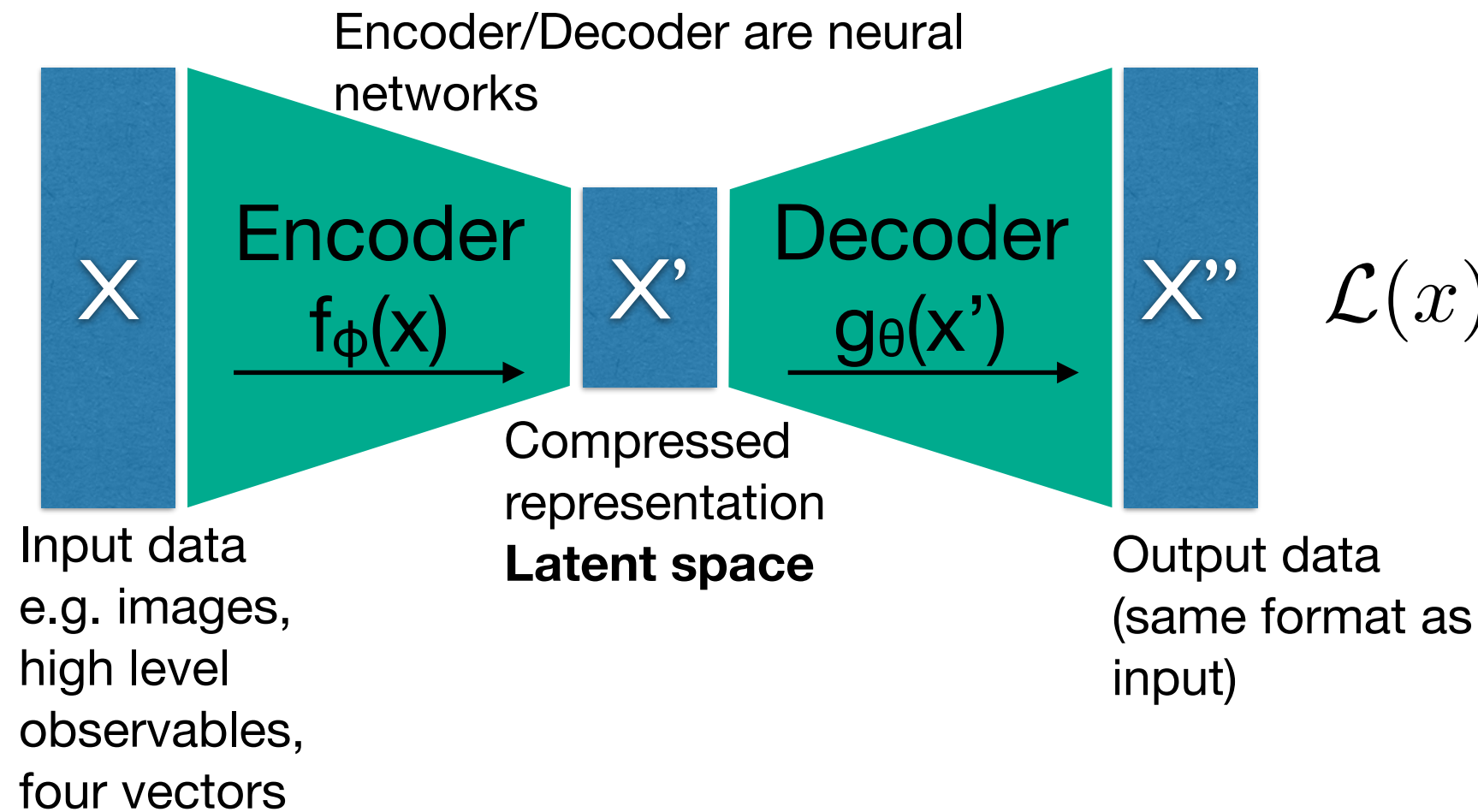
- 5.1 Deep Ensemble Anomaly Detection
- 5.2 Factorized Topic Modeling
- 5.3 QUAK: Quasi-Anomalous Knowledge for Anomaly Detection
- 5.4 Simple Supervised learning with LSTM layers

*(Partial / full labels)*

*Some examples and trends in the following.  
For exhaustive discussion, refer to  
2101.08320*

# Unsupervised - Autoencoders

- Several autoencoder-type learning approaches
- Underlying assumption is that an autoencoder trained on background dominated sample will have bad reconstruction performance for previously unseen signal



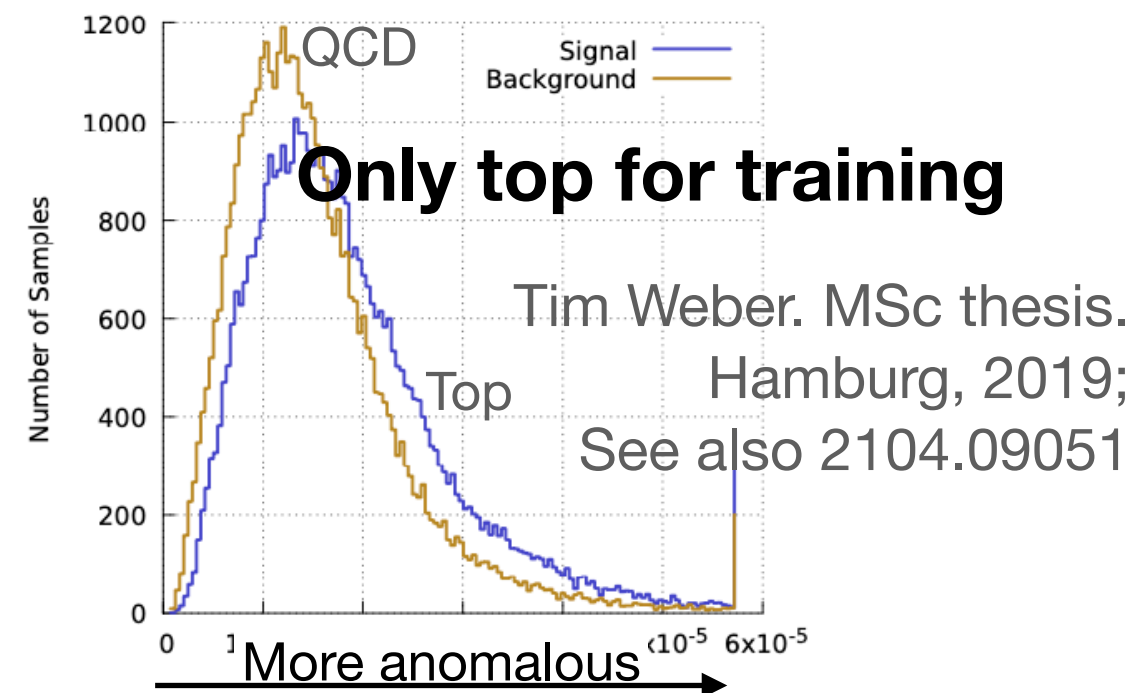
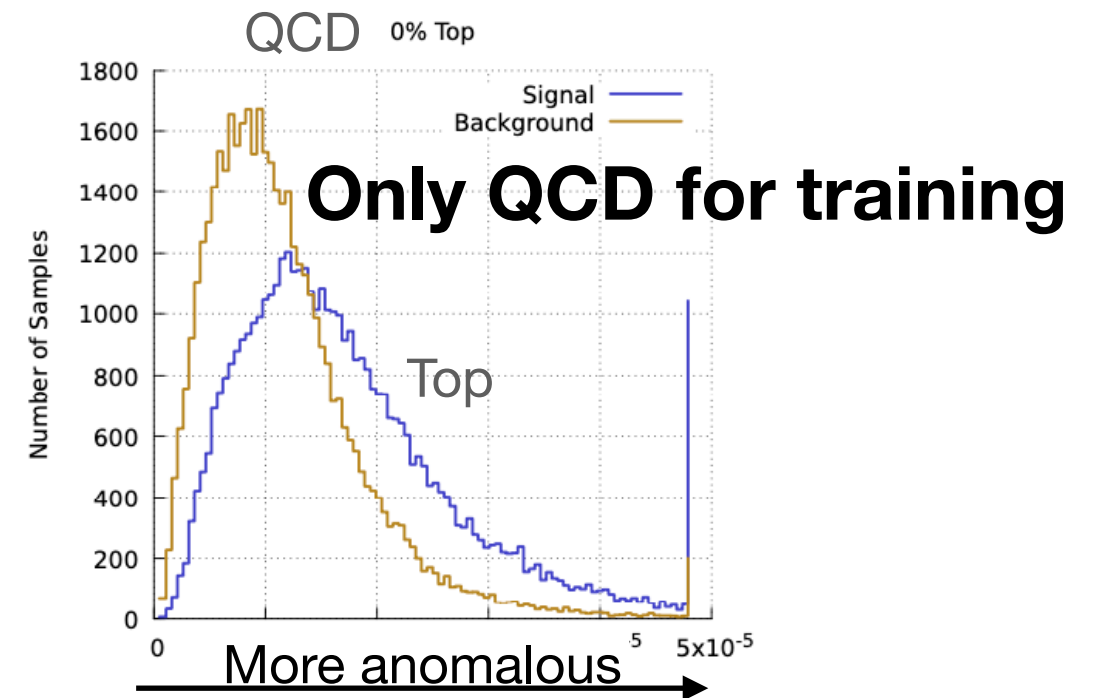
$$\mathcal{L}(x) = ||x - g_\theta(f_\phi(x))||_2$$
$$a(x) = \mathcal{L}(x)$$

- Differences in data representation  
Data space vs latent space anomaly detection  
Different latent space prior distributions

# Limitations

## Complexity

- If anomalies are much simpler (therefore easier to reconstruct):  
*a(x) will still be lower, despite never encountered in training*
- Observed with naive AE in QCD vs top
  - Train on tops only; top still considered anomaly wrt/ QCD

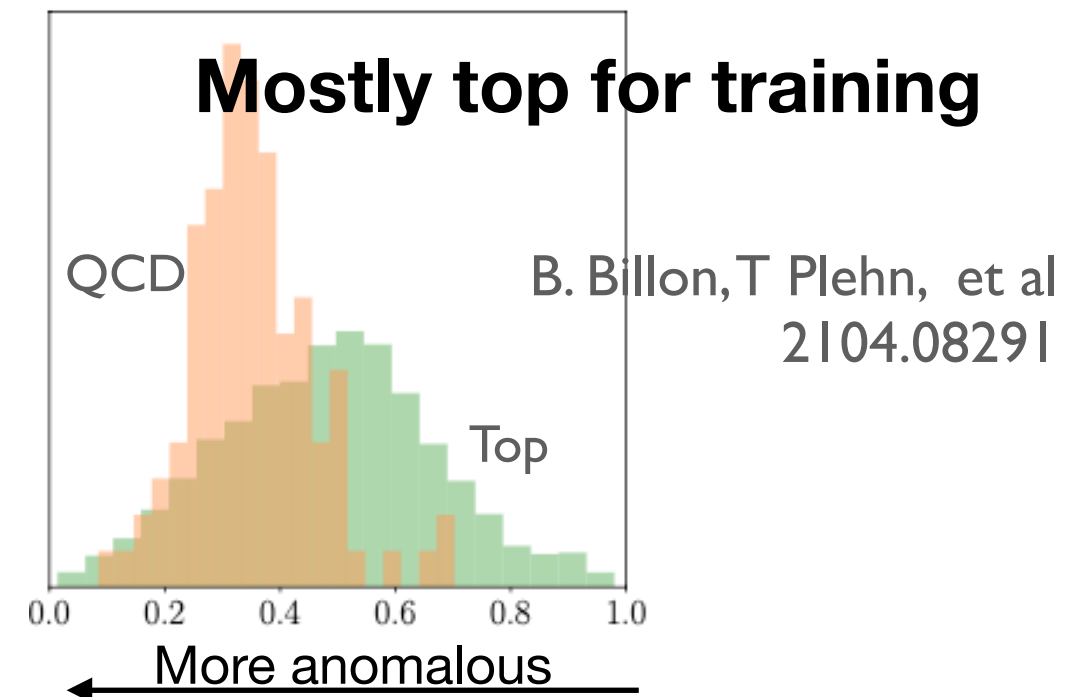
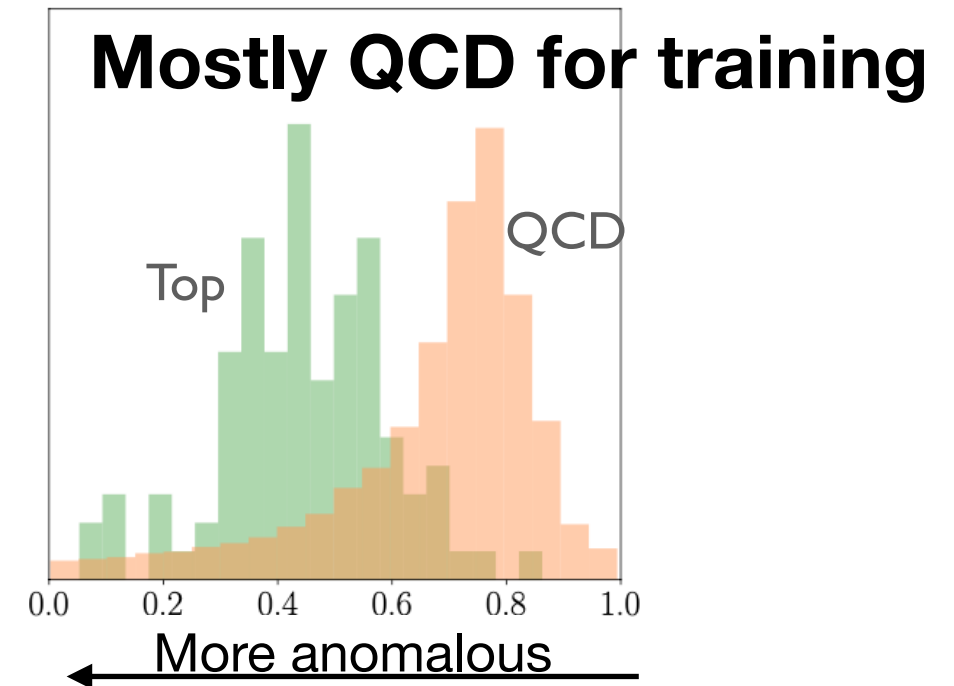




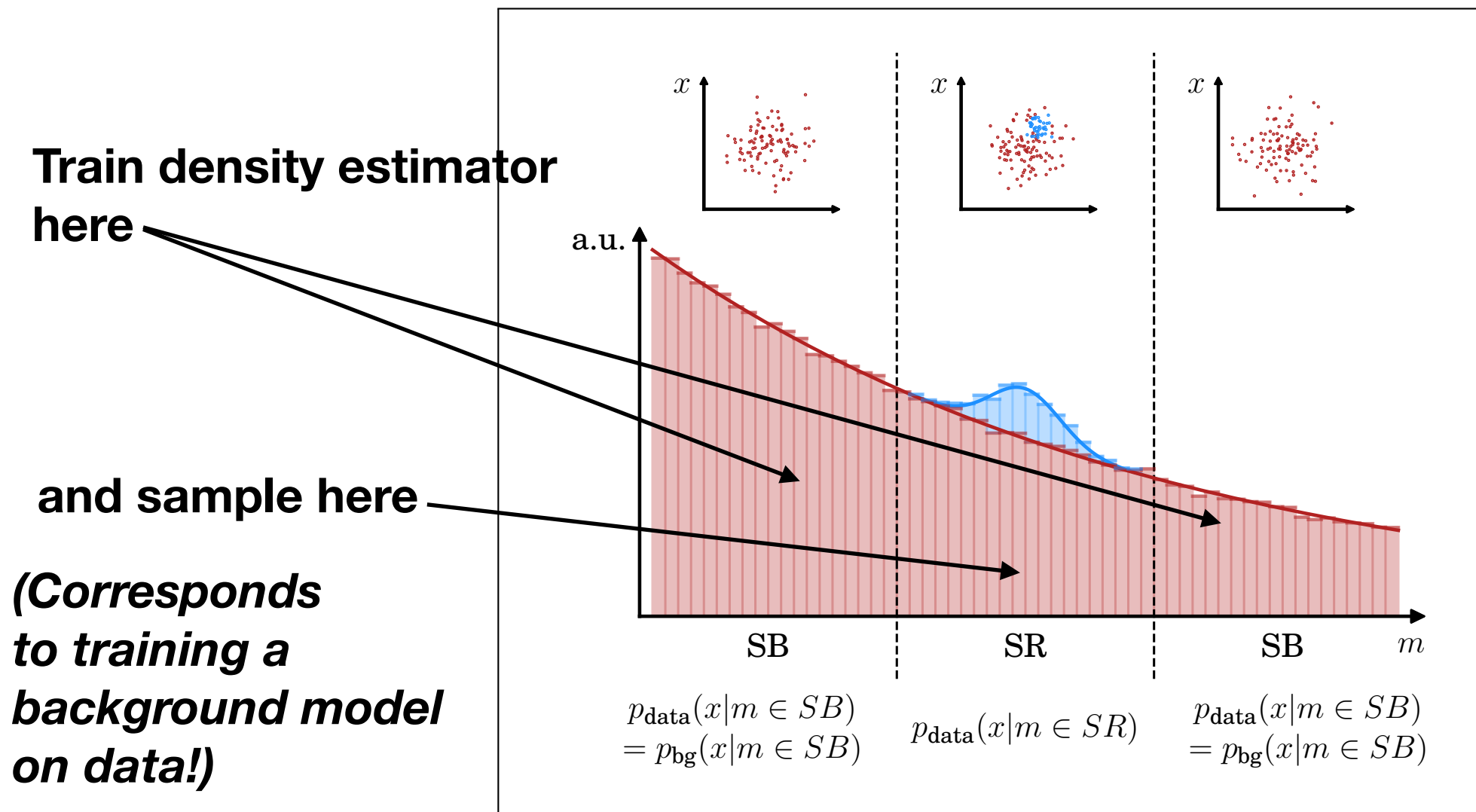
# Limitations

## Complexity

- If anomalies are much simpler (therefore easier to reconstruct):  
*a(x) will still be lower, despite never encountered in training*
- Observed with naive AE in QCD vs top
  - Train on tops only; top still considered anomaly wrt/ QCD
  - Can be overcome (e.g. by structuring the latent space)



# Unsupervised - Density Estimation

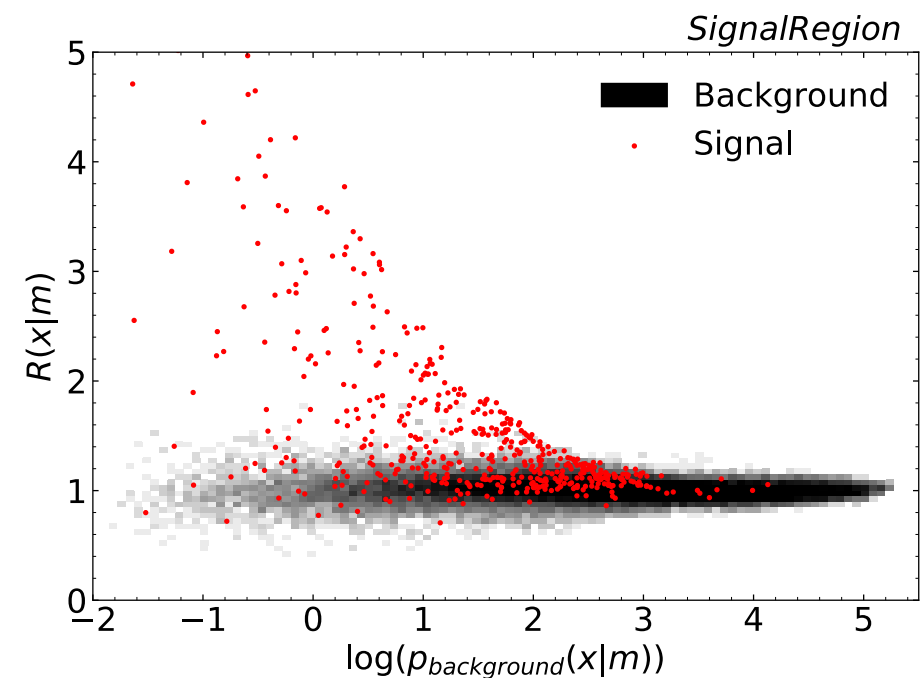
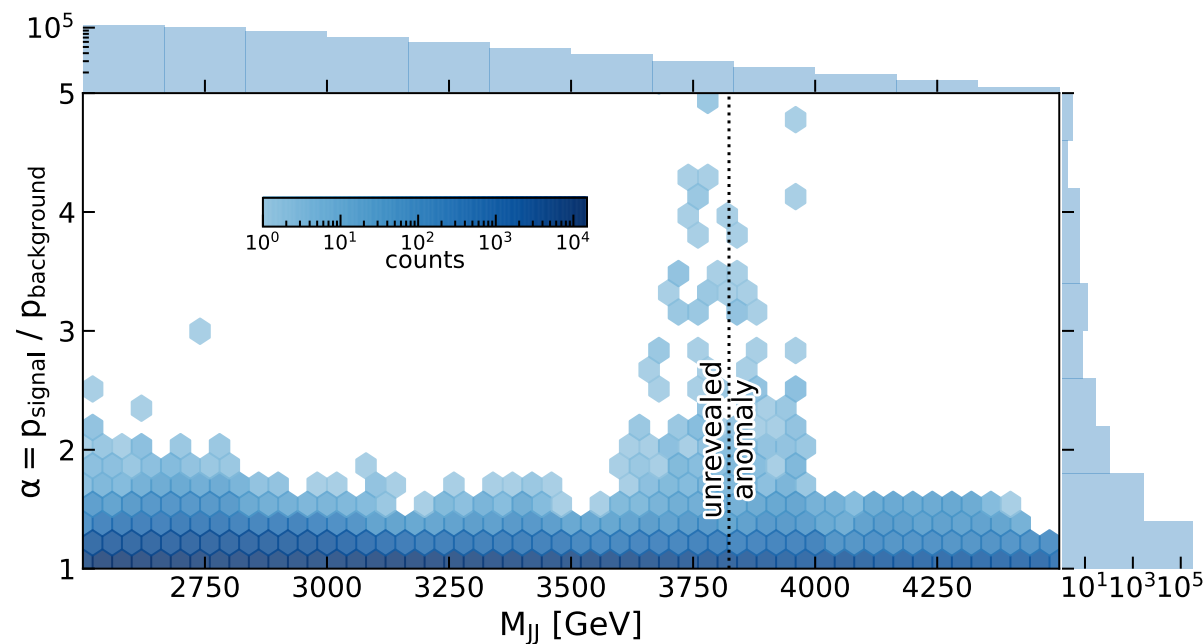


- Conditional transport and sampling
  - Train density estimator (e.g. conditional normalising flow) in sideband
  - Interpolate to signal region
  - Sample data there
  - This produces 'extrapolated-background'

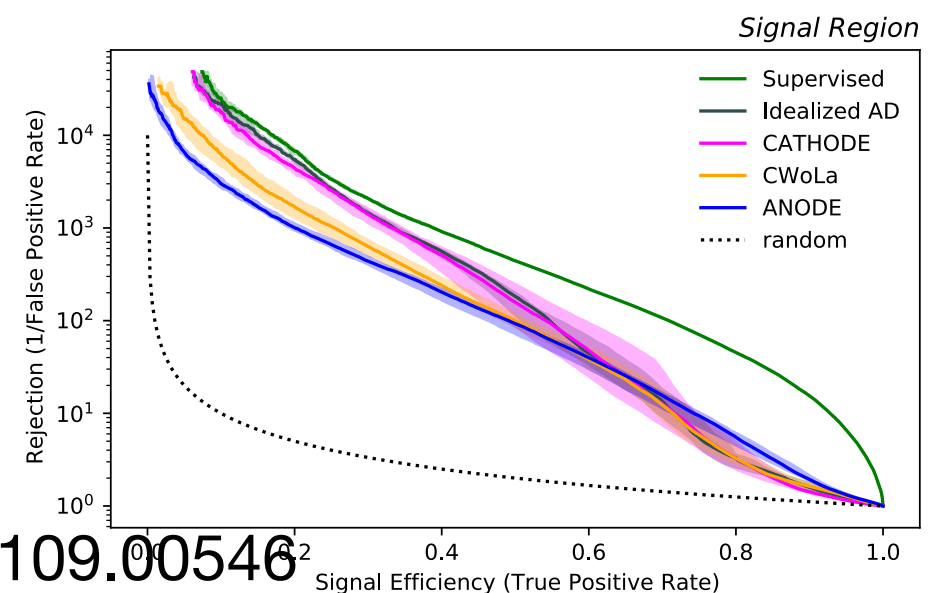
# Unsupervised - Density Estimation

- Compare extrapolated-background to actual data
- Either by also training a density estimator in signal region and building the likelihood ratio (GIS-approach, ANODE)

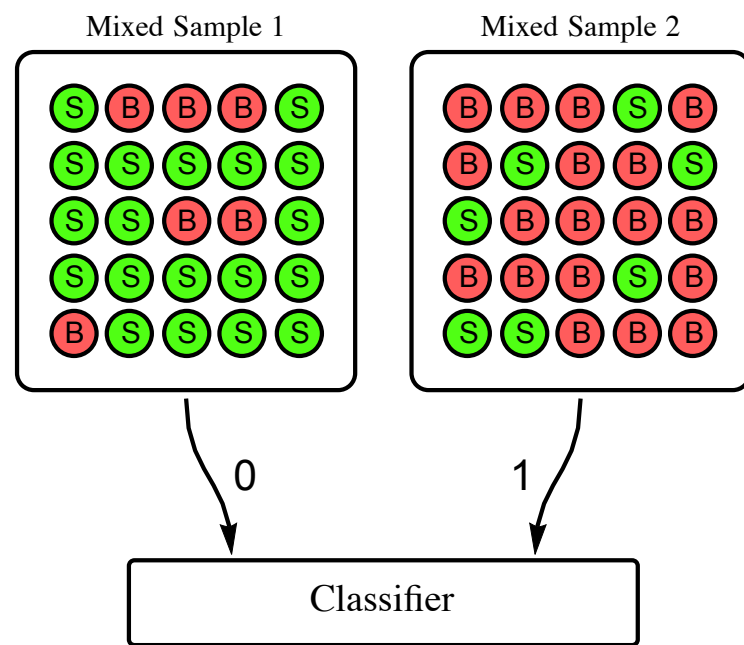
$$R(x|m) = \frac{p_{\text{data}}(x|m)}{p_{\text{background}}(x|m)}.$$



- Or by training a classifier between extrapolated-background and actual data (e.g. CATHODE / 2109.00546 or CURTAINS / 2203.09470) (post-LHCO)



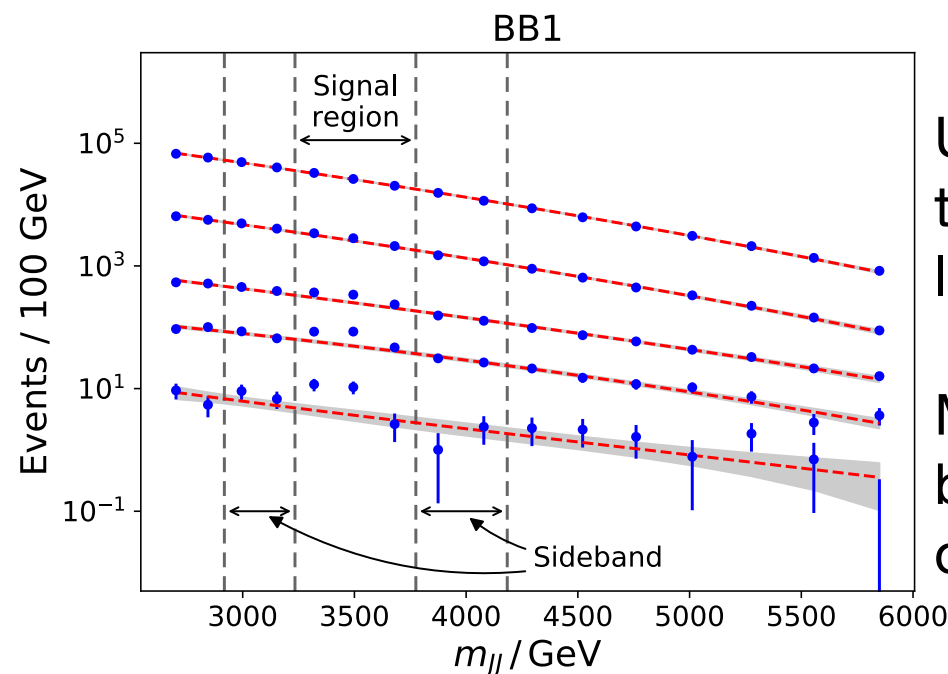
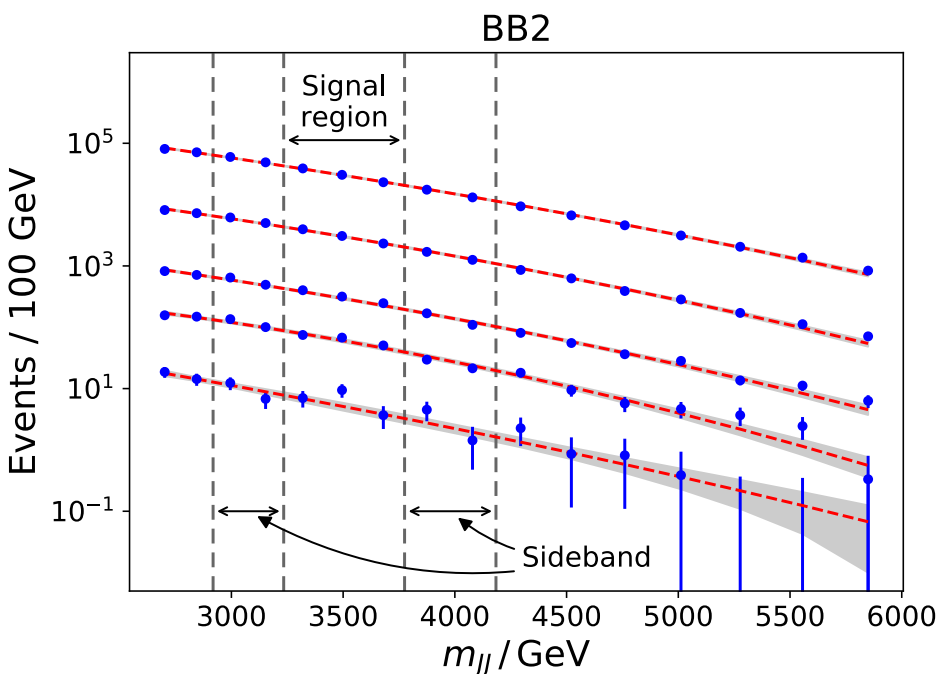
# Weak Supervision: CWola Hunting



$$L_{M_1/M_2} = \frac{p_{M_1}}{p_{M_2}} = \frac{f_1 p_S + (1 - f_1) p_B}{f_2 p_S + (1 - f_2) p_B} = \frac{f_1 L_{S/B} + (1 - f_1)}{f_2 L_{S/B} + (1 - f_2)}$$

Important observation:

*A classifier (i.e. a neural network) trained to distinguish two mixed samples learns to distinguish the components*



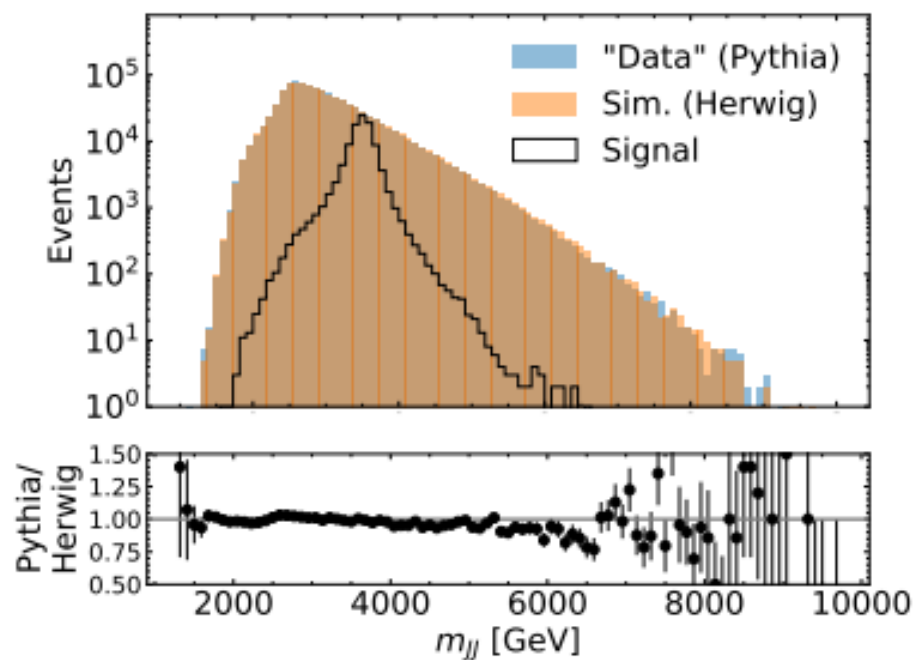
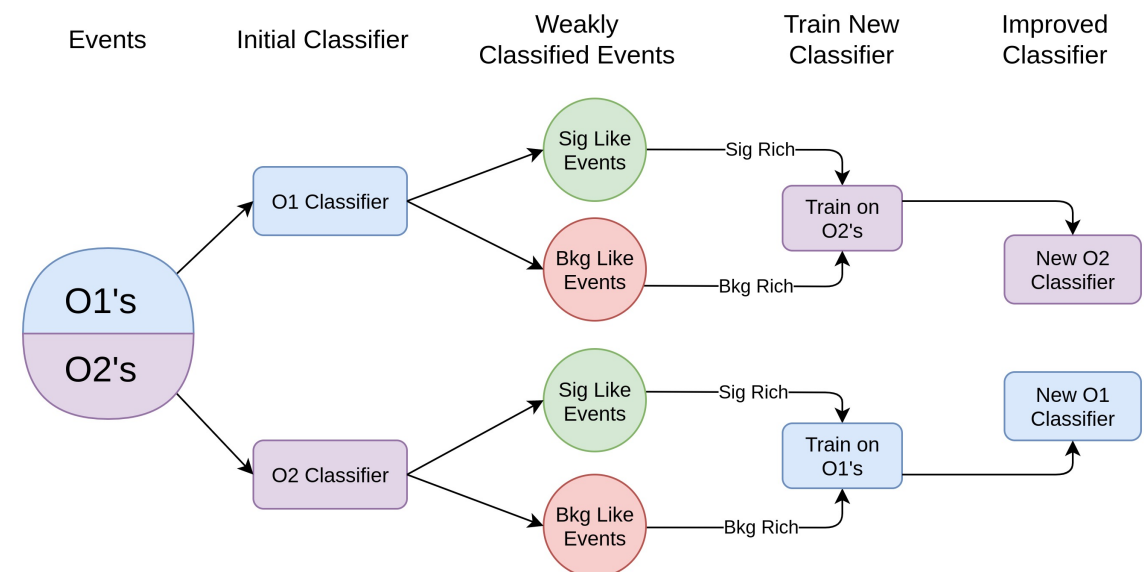
Use increasingly tight selections to identify localised signal.

Major downside: Correlations between mass and other features

# Weak Supervision: TNT and SALAD

## TNT (2002.12376):

Interesting signal might contain two anomalous jets per event. Use per-jet classifiers to build enriched datasets for training.



## SALAD (2001.05001):

Use classifier-based reweighting (DCTR approach / 1907.08209) to learn mapping background simulation in sideband to data. Apply in signal-region and treat non-closure as anomaly

$$w(x|m) \equiv \frac{f(x)}{1 - f(x)} = \frac{p(x|\text{data})}{p(x|\text{simulation})} \times \frac{p(\text{data})}{p(\text{simulation})}$$



# Semi-supervised

## Signal-classifier based training

Train a classifier on a potential signal (or cocktail of potential signals!) and use like in a fully supervised search.

### QUAK (2011.03550):

Combine potential signals (supervised) and unlabelled data. Essentially use different signal priors to build a latent space in which to search for anomalies.

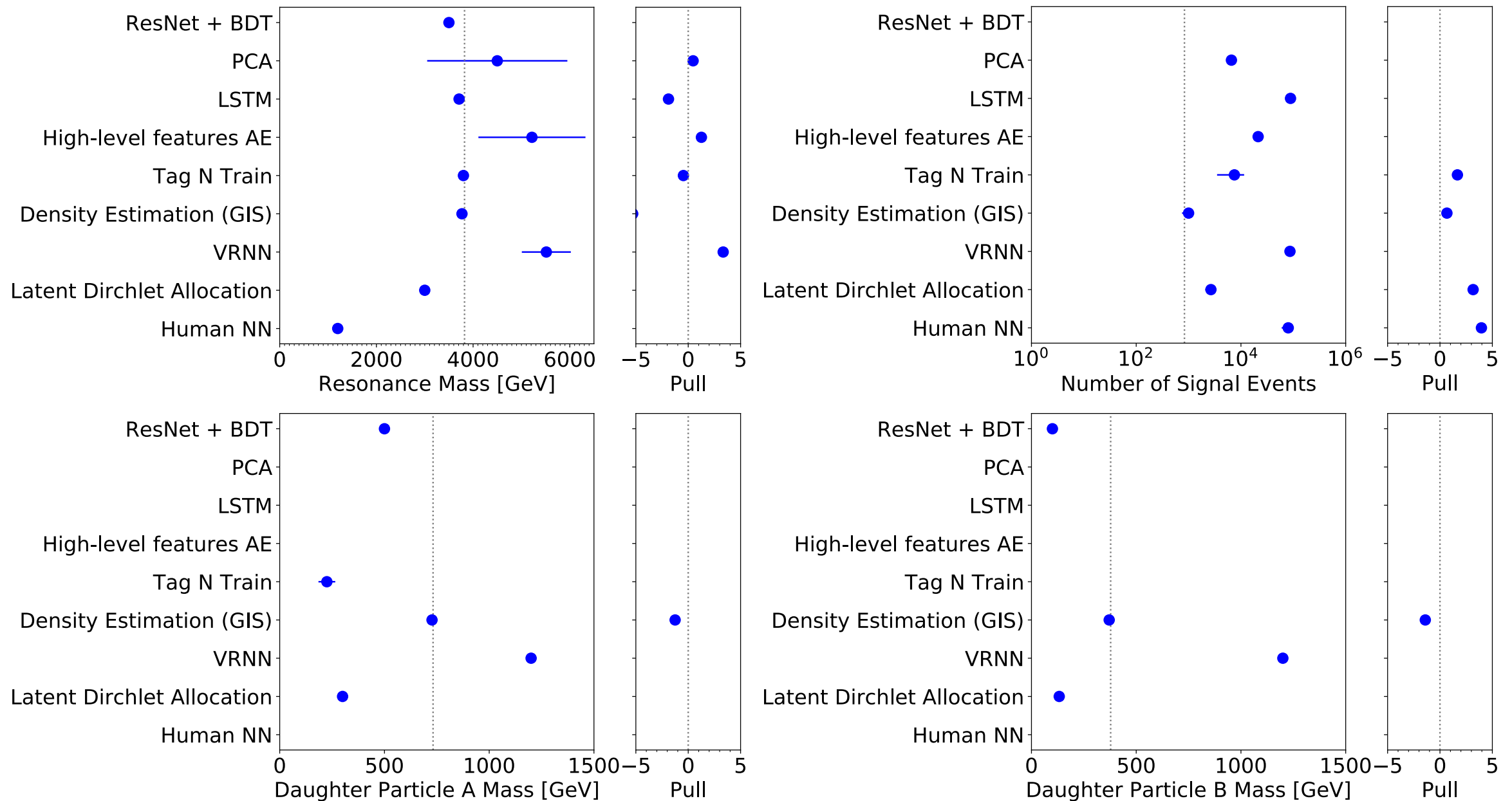


# Reporting of results

Section	Short Name	Method Type	Results Type
<a href="#">3.1</a>	VRNN	Unsupervised	(i) (BB2,3) and (ii) (BB1)
<a href="#">3.2</a>	ANODE	Unsupervised	(iii)
<a href="#">3.3</a>	BuHuLaSpa	Unsupervised	(i) (BB2,3) and (ii) (BB1)
<a href="#">3.4</a>	GAN-AE	Unsupervised	(i) (BB2-3) and (ii) (BB1)
<a href="#">3.5</a>	GIS	Unsupervised	(i) (BB1)
<a href="#">3.6</a>	LDA	Unsupervised	(i) (BB1-3)
<a href="#">3.7</a>	PGA	Unsupervised	(ii) (BB1-2)
<a href="#">3.8</a>	Reg. Likelihoods	Unsupervised	(iii)
<a href="#">3.9</a>	UCluster	Unsupervised	(i) (BB2-3)
<a href="#">4.1</a>	CWoLa	Weakly Supervised	(ii) (BB1-2)
<a href="#">4.2</a>	CWoLa AE Compare	Weakly/Unsupervised	(iii)
<a href="#">4.3</a>	Tag N' Train	Weakly Supervised	(i) (BB1-3)
<a href="#">4.4</a>	SALAD	Weakly Supervised	(iii)
<a href="#">4.5</a>	SA-CWoLa	Weakly Supervised	(iii)
<a href="#">5.1</a>	Deep Ensemble	Semisupervised	(i) (BB1)
<a href="#">5.2</a>	Factorized Topics	Semisupervised	(iii)
<a href="#">5.3</a>	QUAK	Semisupervised	(i) (BB2,3) and (ii) (BB1)
<a href="#">5.4</a>	LSTM	Semisupervised	(i) (BB1-3)

*i) during challenge phase*  
*ii) after challenge phase*  
*iii) R&D dataset used*

# Results - BB1



*(Shown are results during challenge)*

*Several approaches identified resonance; density estimation also found correct properties*

# Results - BB2

Reminder: no signal injected

Some methods reported resonances in the tail of the mass distribution (around 4.5 TeV)

*Difficult to predict for edges of phase space*

**Latent Dirichlet Allocation:** Our method extracts signal descriptions which appear convincing, however the classifier does not identify a bump in the invariant mass spectra. Without this we were unable to determine that a signal was present. The di-jet description extracted consisted of one jet of mass 350-400 GeV and another of mass 150-200 GeV. If the production of these states was non-resonant, we would be unable to find the signal with our method. Or if more than just di-jets were relevant to reconstruct the invariant mass, we would also not be able to find it. **Otherwise, we determine that no signal was present in the data.**

# Results - BB3

Reminder: di-jet and tri-jet topologies

Different observations claimed, none identified the correct signal.



# Lessons learned

- Anomaly detection is difficult
  - Even for “anomalies” close to already considered signals
  - Even more so for “exotic” signals
  - Value in blind studies
- Robust uncertainty quantification needed - especially for tails of distributions
- Many methods used “sidebanding” in invariant mass + learning some anomaly detector.
  - Less reliance on peaks and
  - less reliance on one ‘lucky’ (physically inspired) variable desirable
- Did not discuss data representation: image vs point cloud vs ... vs high-level features
  - Will bias anomaly detection performance. Need to understand better.

# Other open issues

- Strategies to assess the quality of anomaly detection techniques without (or at least with less) dependence on specific signal models? Right now, the strategy seems to be to compare the ability of ADs to find some benchmark signals.
- Can there be robust methods to set exclusion limits with ‘data only’ anomaly detectors (i.e. methods where all final trainings are carried out on data - as opposed to training on simulation) without injecting signal events into the data.
- How to publish the on-data-trained anomaly detectors in such a way that allows ex-post analysis by people outside the experiment whether the training result is compatible / rules out a given new physics signals.
- Methods to go from an observed anomaly (ie. a signal like excess in some region of data) to an interpretation in terms of physics models can still be improved as well.

# Conclusions

- Exciting space of anomaly detection in LHC physics
- First successes and breadth of ideas but better understanding and more applications needed
- Potential to search for more signatures with less people, even while some conceptual issues are being resolved
- Also applies to other areas (e.g. data quality / detector operations)



*Thank you!*