

# New Paradigms for New Physics Searches with Machine Learning

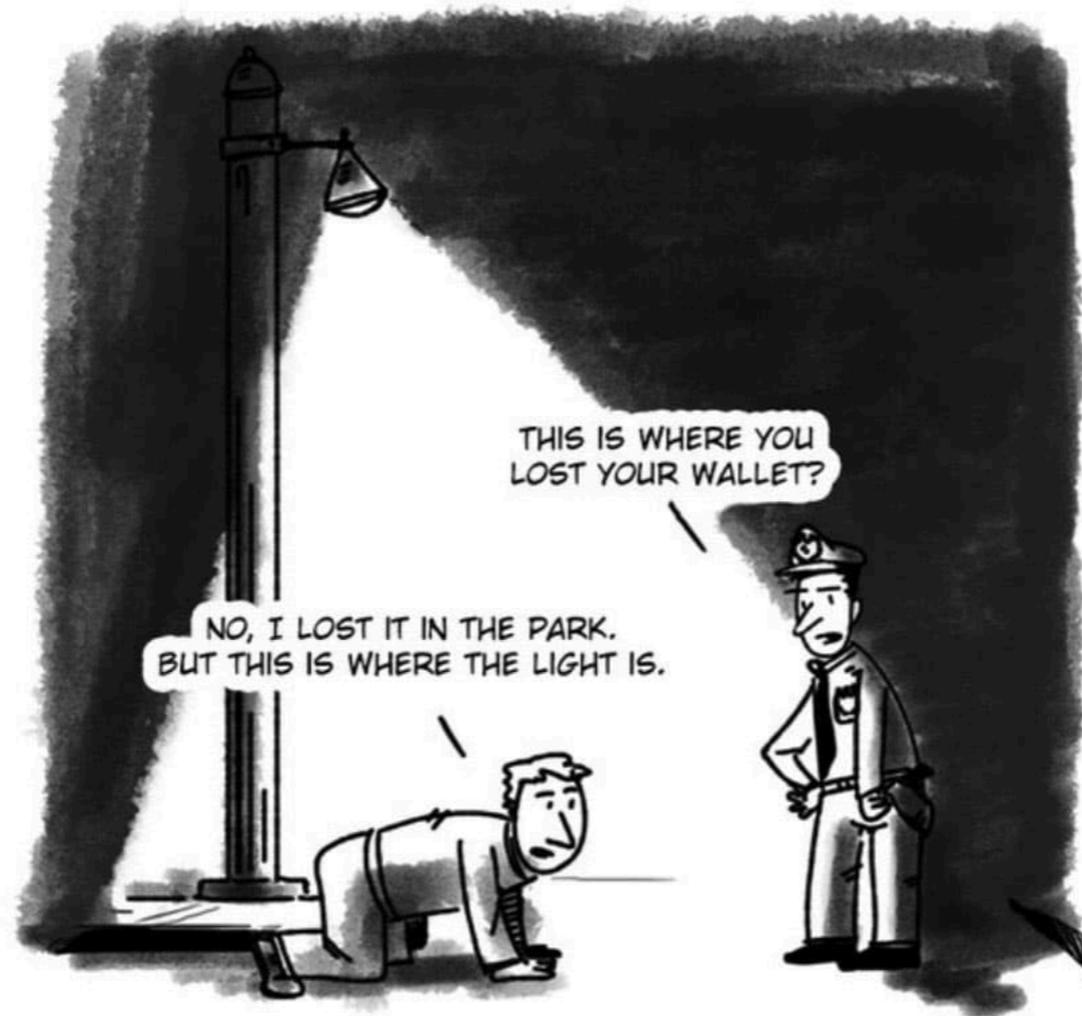
David Shih

January 13, 2021

Theory Mini-workshop  
HKUST IAS Program on HEP



# Where is the new physics??



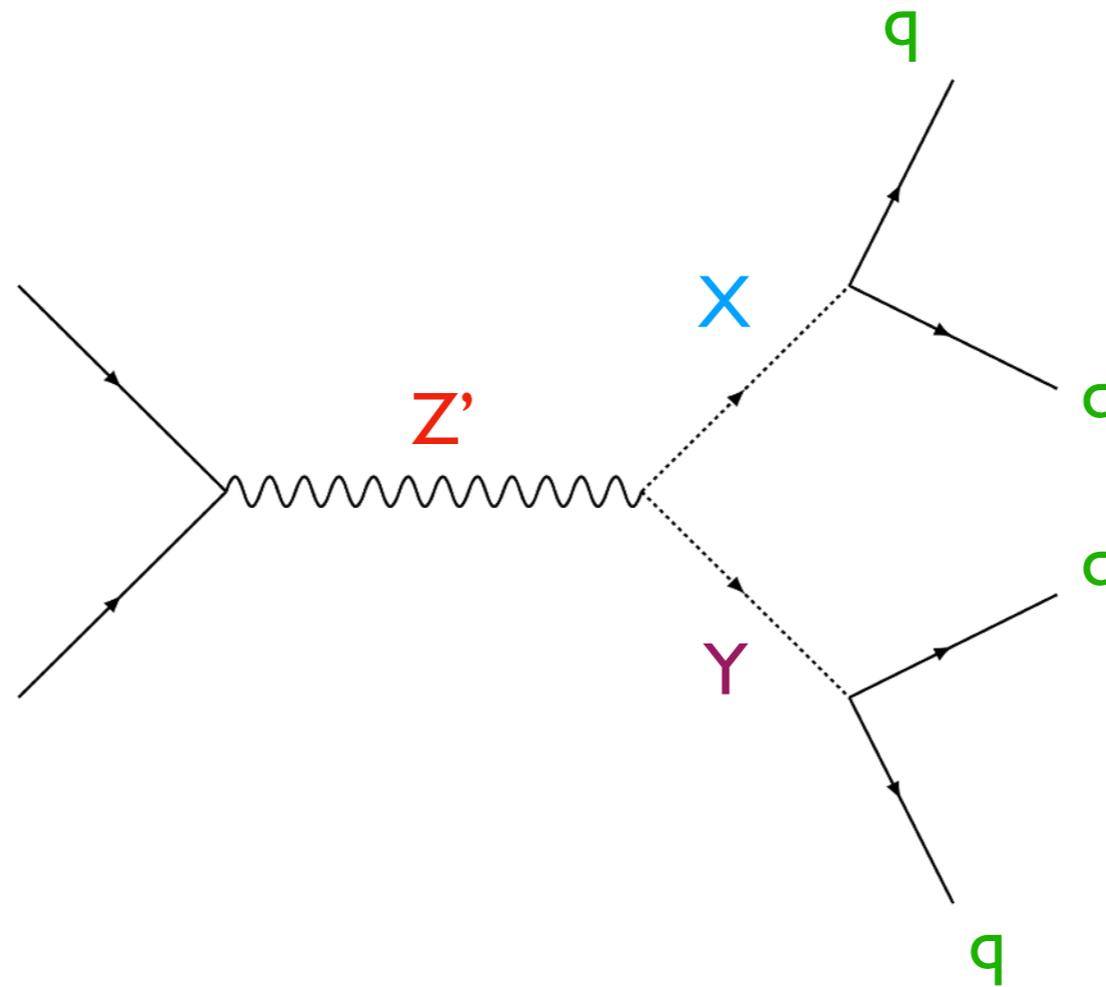
Despite thousands of searches for new physics at the LHC, nothing but limits and null results so far.

**What if new physics is hiding in the data but we haven't looked in the right places yet?**

# A Benchmark Example

LHC Olympics 2020 R&D Dataset

<https://doi.org/10.5281/zenodo.2629072>



**No explicit search at the LHC for this scenario!**

**Could be hiding in the dijet resonance search at  $>5\sigma$  significance!!**

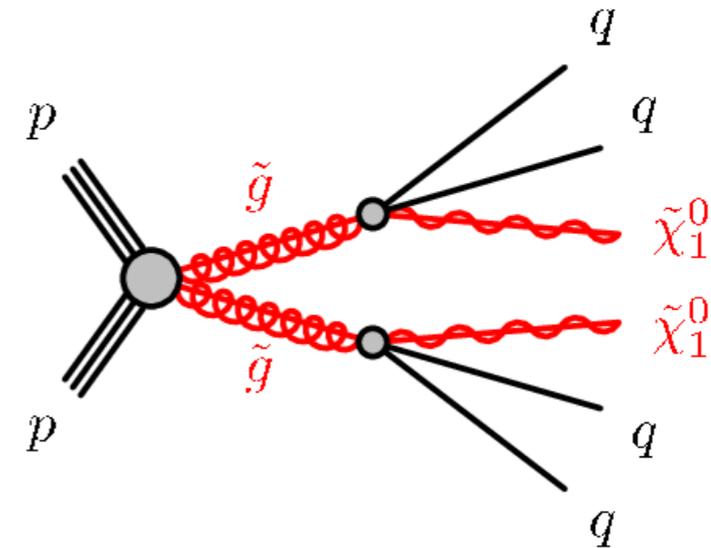
# The most common approach

## Model specific searches

Most NP searches at the LHC are heavily optimized with specific signals in mind (SUSY, extra dimensions, ...)

ATLAS jets+MET 2010.14293

|  | BDT-GGd1  | BDT-GGd2  | BDT-GGd3 | BDT-GGd4 |
|--|-----------|-----------|----------|----------|
| $N_j$  | $\geq 4$  |           |          |          |
| $\Delta\phi(j_{1,2,(3)}, \mathbf{p}_T^{\text{miss}})_{\text{min}}$ | $> 0.4$   |           |          |          |
| $\Delta\phi(j_{i>3}, \mathbf{p}_T^{\text{miss}})_{\text{min}}$     | $> 0.4$   |           |          |          |
| $E_T^{\text{miss}}/m_{\text{eff}}(N_j)$                            | $> 0.2$   |           |          |          |
| $m_{\text{eff}}$ [GeV]   | $> 1400$  |           | $> 800$  |          |
| BDT score  | $> 0.97$  | $> 0.94$  | $> 0.94$ | $> 0.87$ |
| $\Delta m(\tilde{g}, \tilde{\chi}_1^0)$ [GeV]                      | 1600–1900 | 1000–1400 | 600–1000 | 200–600  |



Kinematic cuts (or BDTs) optimized using simulations of signal AND background.

# The most common approach

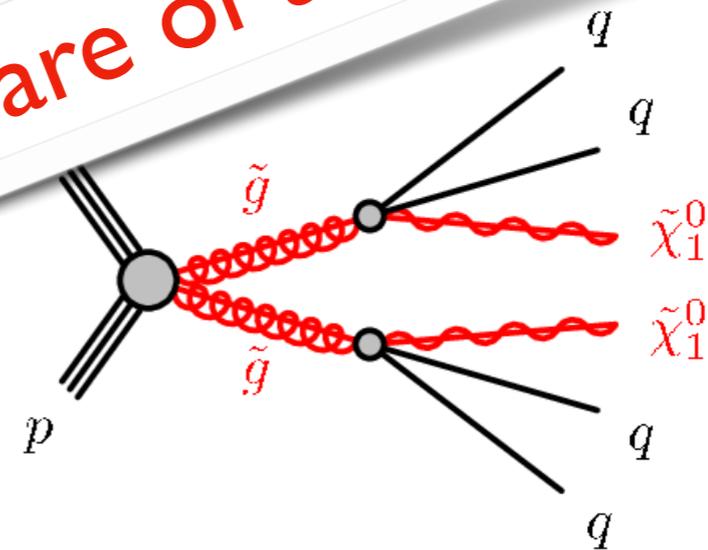
## Model specific searches

Most NP searches at the LHC are heavily optimized with specific signals in mind (SUSY, extra dimensions, ...)

ATLAS jets+MET 2010.14293

|   | BDT-GGd1 | BDT-GGd2  | BDT-GGd3 | BDT-GGd4 |
|---|----------|-----------|----------|----------|
| $N_j$   |          | $\geq 4$  |          |          |
| $\Delta\phi(j_{1,2,(3)}, p_T^{\text{miss}})_{\text{min}}$ |          |           |          |          |
| $\Delta\phi(j_{i>3}, p_T^{\text{miss}})_{\text{min}}$     |          |           |          |          |
| $E_T^{\text{miss}}/m_{\text{eff}}(N_j)$                   |          |           |          |          |
| $m_{\text{eff}}$ [GeV]                                    |          |           | $> 800$  |          |
| BDT score   |          | $> 0.94$  | $> 0.94$ | $> 0.87$ |
| $\Delta m(\tilde{g}, \tilde{\chi}_1^0)$                   | 500–1900 | 1000–1400 | 600–1000 | 200–600  |

**> 99% of searches at the LHC are of this type**



Kinematic cuts (or BDTs) optimized using simulations of signal AND background.



Of course, we should continue to perform these searches, because NP could always be right around the corner...

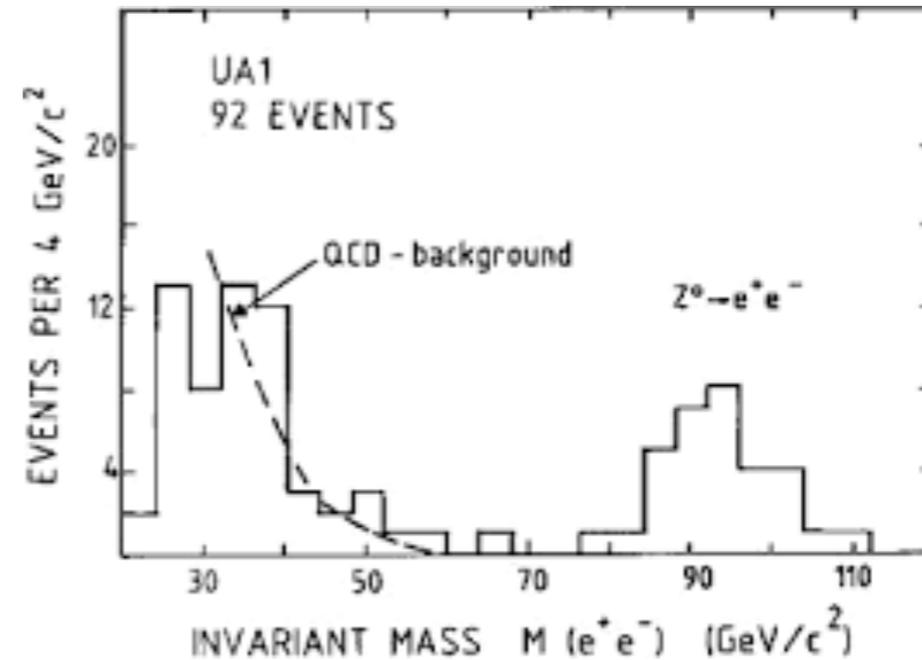
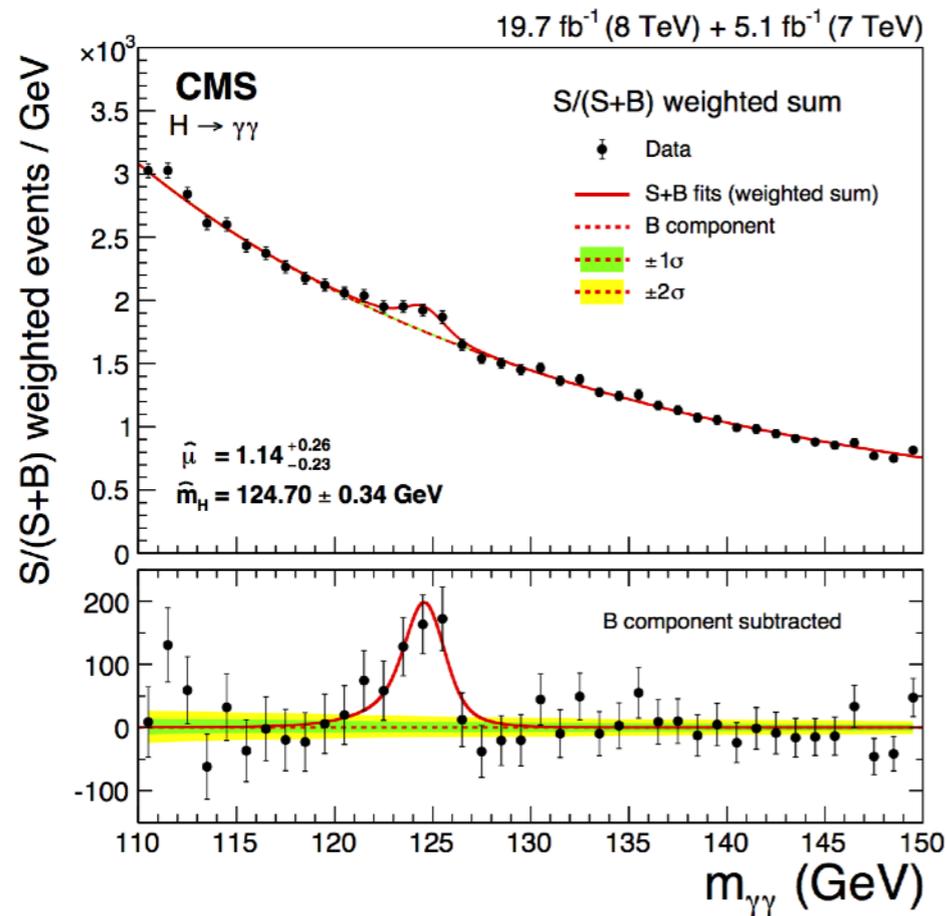
**But we should also recognize that perhaps new physics is on a different street, or in a different building altogether...**

# Existing model-independent approaches

“the bump hunt”

Idea: assume signal is localized in some feature (usually invariant mass) while background is smooth.

Interpolate from **sidebands** into **signal region**, search for an excess.



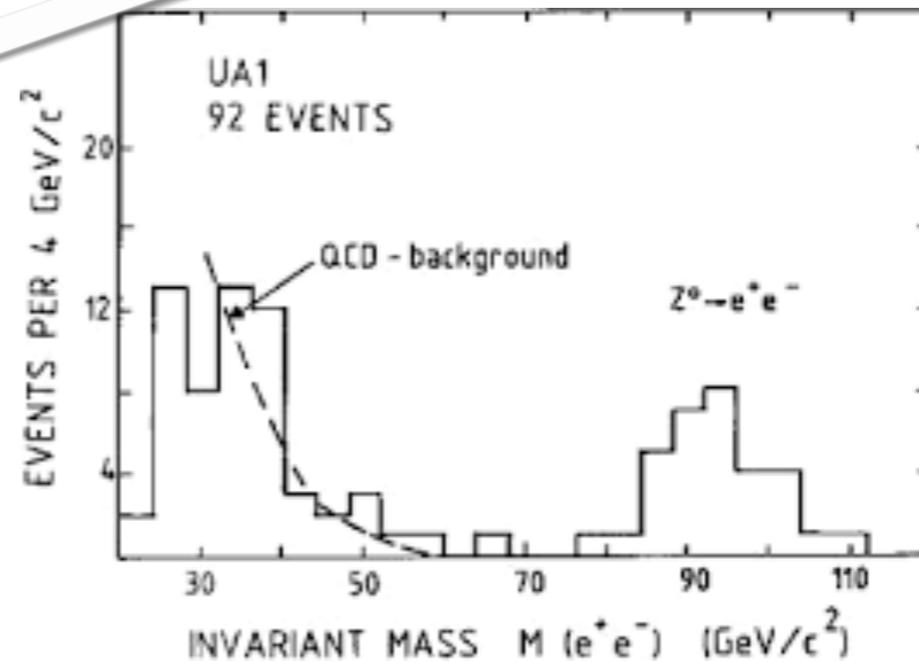
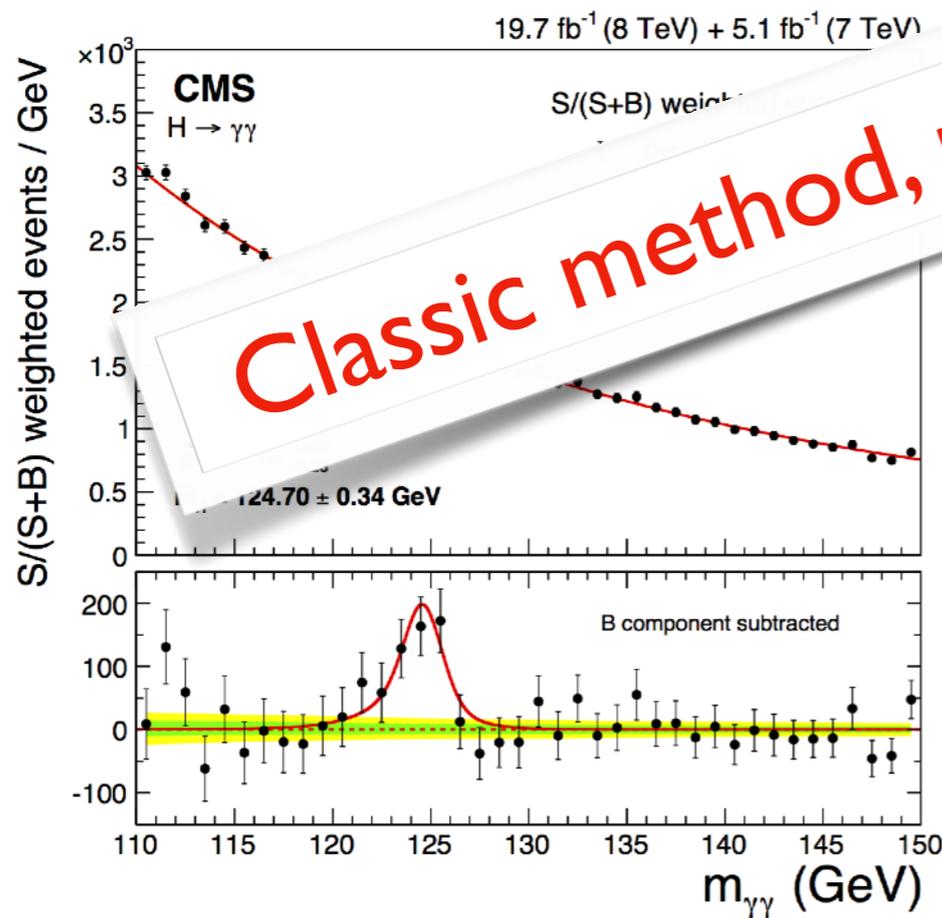
# Existing model-independent approaches

“the bump hunt”

Idea: assume signal is localized in some feature (usually invariant mass) while background is smooth.

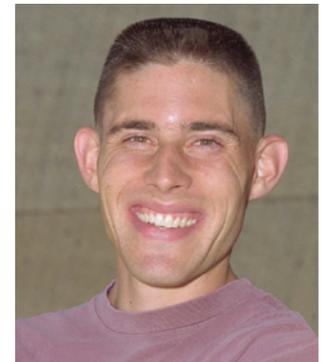
Interpolate from **sidebands** into **signal region** to find any excess.

**Classic method, used in many discoveries.**

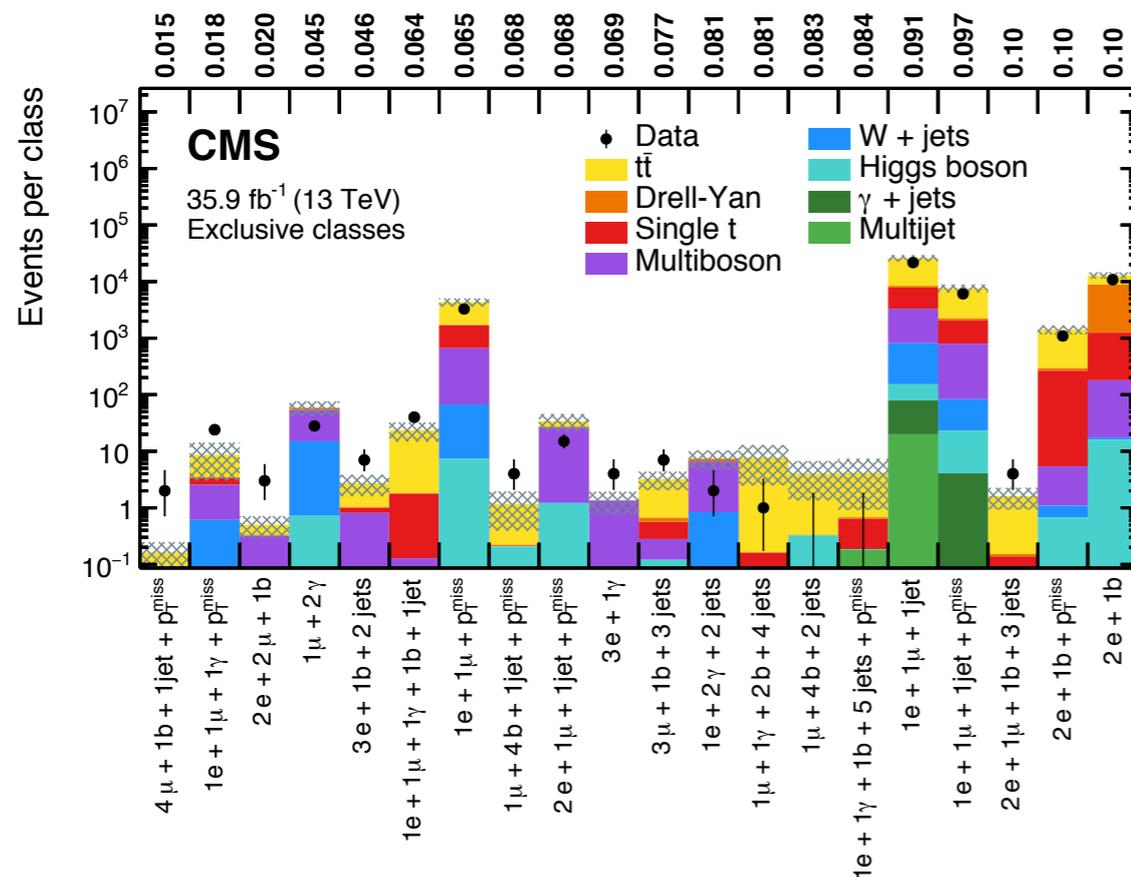


# Existing model-independent searches

“the general search”



Idea: divide the phase space up into thousands of bins, compare **data to SM simulation** in each one



**CMS**

“MUSIC”

CMS-PAS-EXO-14-016

**ATLAS**

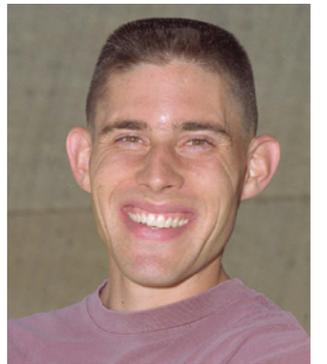
“Model independent general search”

1807.07447  
EPJC 79:120 (2019)

See also proposals by D’Agnolo, Wulzer et al (1806.02350, 1912.12155):  
train DNN on full phase space to distinguish data from background

# Existing model-independent searches

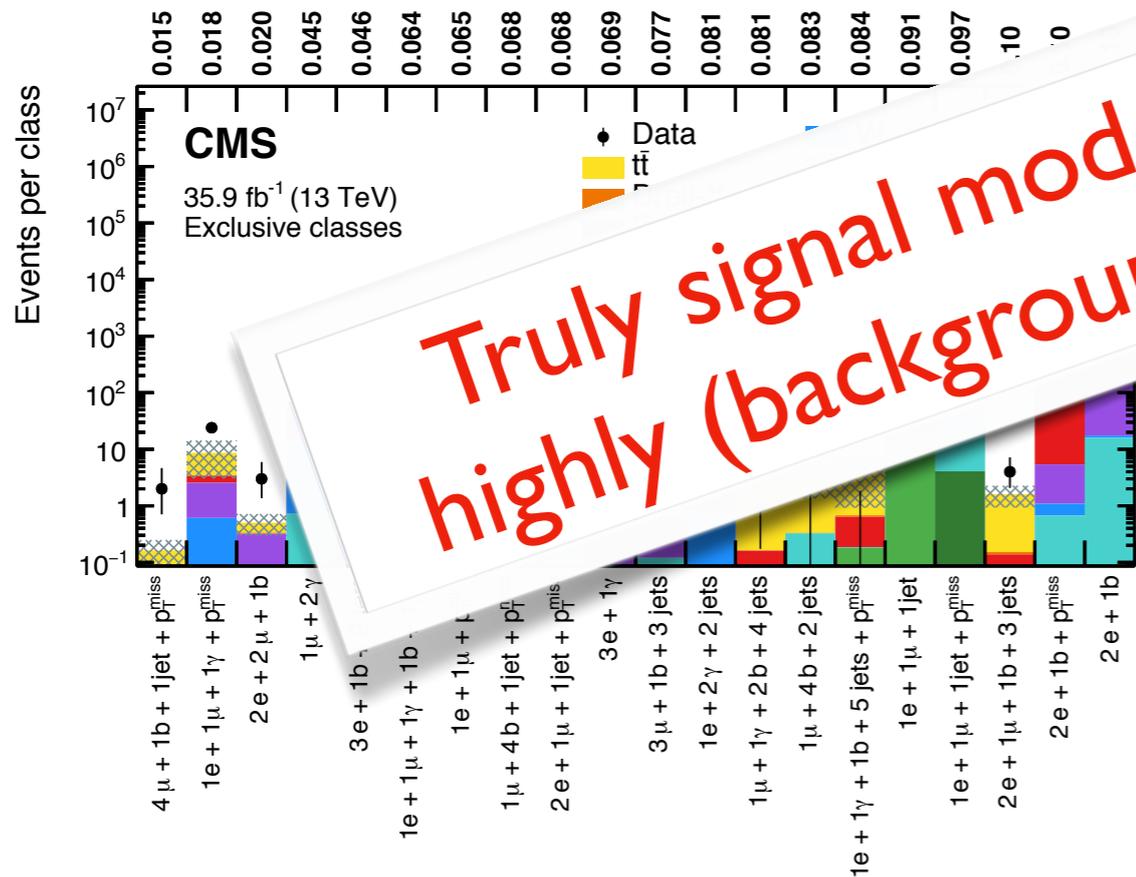
“the general search”



Idea: divide the phase space up into thousands of bins, compare **data to SM simulation** in each one

**Truly signal model independent, but still highly (background) simulation dependent**

“MUSIC”



**ATLAS**

“Model independent general search”

CMS-PAS-EXO-14-016

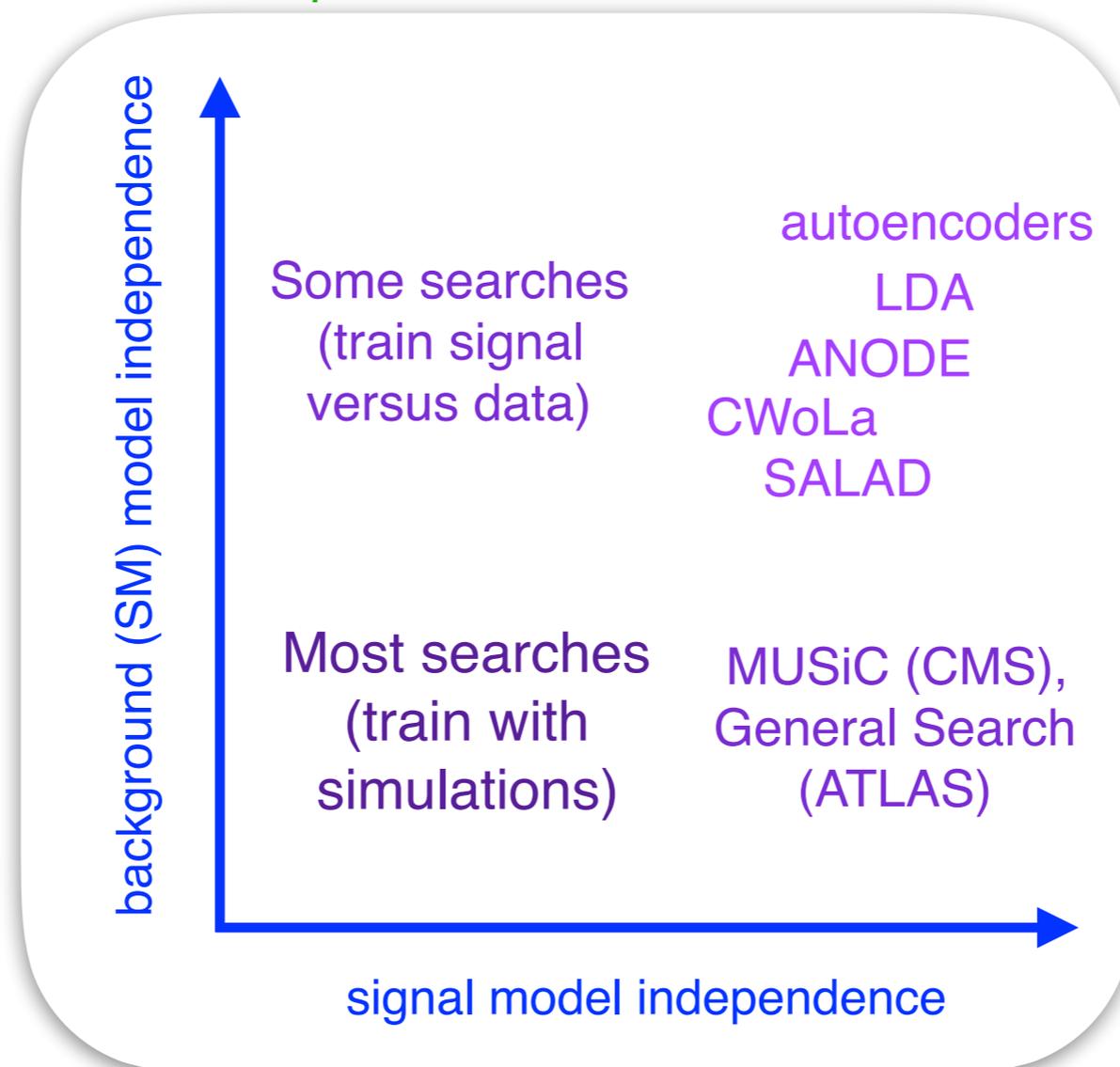
1807.07447  
EPJC 79:120 (2019)

See also proposals by D’Agnolo, Wulzer et al (1806.02350, 1912.12155):  
train DNN on full phase space to distinguish data from background

# New paradigms for model-agnostic searches

*Can advances in machine learning open up new avenues for model-independent searches?*

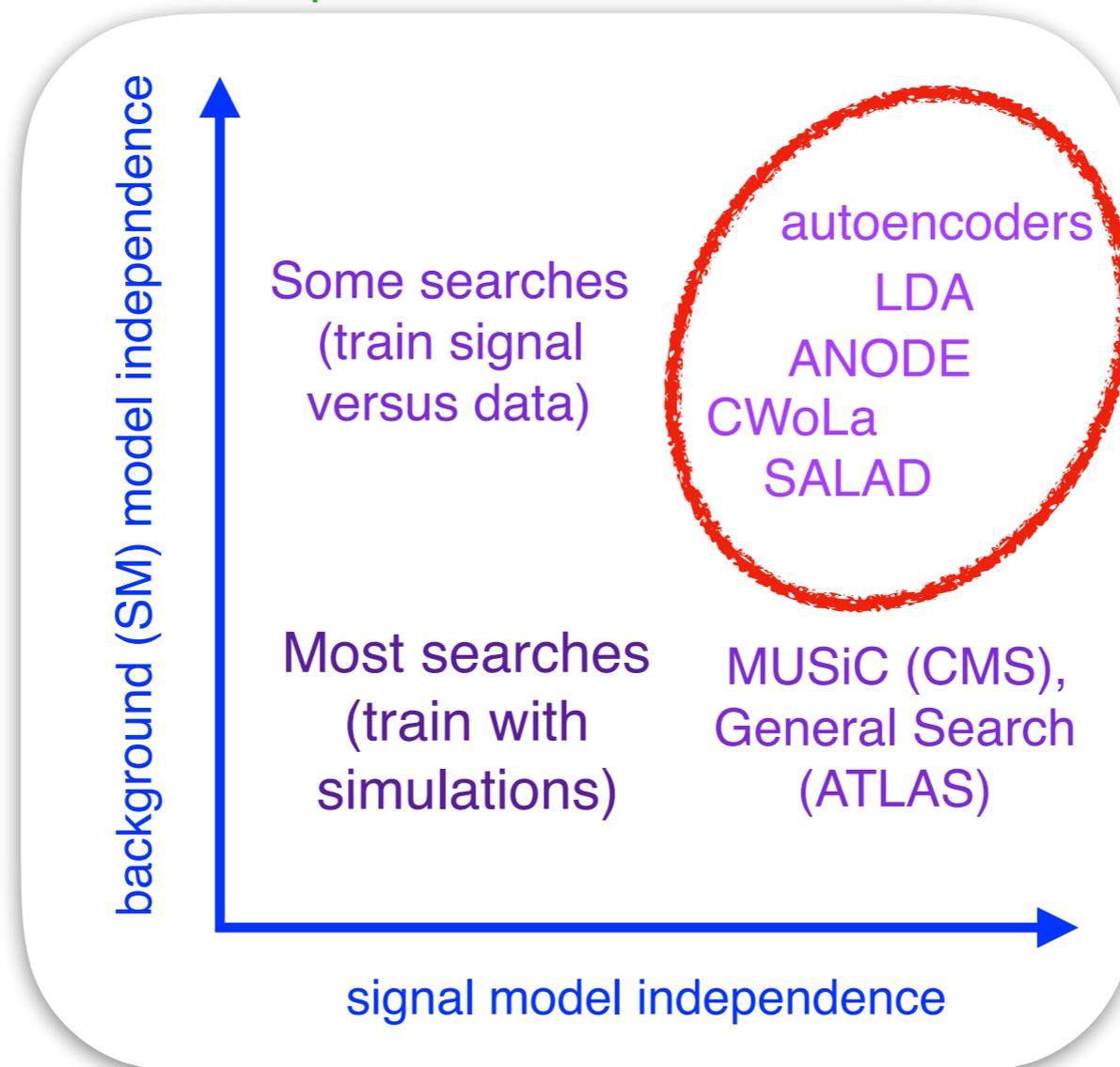
*from Nachman & DS 2001.04990*



# New paradigms for model-agnostic searches

*Can advances in machine learning open up new avenues for model-independent searches?*

*from Nachman & DS 2001.04990*



**Many new ideas recently!**

# Many new approaches inspired by (or at least demonstrated on) the **LHC Olympics 2020 Data Challenge**

## The LHC Olympics 2020

A Community Challenge for Anomaly  
Detection in High Energy Physics



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

to appear on the arXiv next week...

# Many new approaches inspired by (or at least demonstrated on) the **LHC Olympics 2020 Data Challenge**

|    |   |           |
|----|---|-----------|
| 73 | <b>Individual Approaches</b>  | <b>9</b>  |
| 74 | <b>3 Unsupervised</b>   | <b>11</b> |
| 75 | 3.1 Anomalous Jet Identification via Variational Recurrent Neural Network       | 11        |
| 76 | 3.2 Anomaly Detection with Density Estimation                                   | 16        |
| 77 | 3.3 BuHuLaSpa: Bump Hunting in Latent Space                                     | 19        |
| 78 | 3.4 GAN-AE and BumpHunter   | 24        |
| 79 | 3.5 Gaussianizing Iterative Slicing (GIS): Unsupervised In-distribution Anomaly |           |
| 80 | Detection through Conditional Density Estimation                                | 29        |
| 81 | 3.6 Latent Dirichlet Allocation   | 34        |
| 82 | 3.7 Particle Graph Autoencoders   | 38        |
| 83 | 3.8 Regularized Likelihoods   | 42        |
| 84 | 3.9 UCluster: Unsupervised Clustering   | 46        |
| 85 | <b>4 Weakly Supervised</b>  | <b>51</b> |
| 86 | 4.1 CWoLa Hunting   | 51        |
| 87 | 4.2 CWoLa and Autoencoders: Comparing Weak- and Unsupervised methods            |           |
| 88 | for Resonant Anomaly Detection  | 55        |
| 89 | 4.3 Tag N' Train  | 60        |
| 90 | 4.4 Simulation Assisted Likelihood-free Anomaly Detection                       | 63        |
| 91 | 4.5 Simulation-Assisted Decorrelation for Resonant Anomaly Detection            | 68        |
| 92 | <b>5 (Semi)-Supervised</b>  | <b>71</b> |
| 93 | 5.1 Deep Ensemble Anomaly Detection   | 71        |
| 94 | 5.2 Factorized Topic Modeling   | 77        |
| 95 | 5.3 QUAK: Quasi-Anomalous Knowledge for Anomaly Detection                       | 81        |
| 96 | 5.4 Simple Supervised learning with LSTM layers                                 | 85        |

# LHC Olympics 2020: Black Boxes

<https://doi.org/10.5281/zenodo.3547721>

In 2019, Gregor Kasieczka, Ben Nachman and I initiated the LHC Olympics 2020 Challenge.

It consisted of three “black boxes” of simulated data (bg dominated!):

- 1 million events each
- 4-vectors of every reconstructed particle (all hadronic) in the event
- Particle ID, charge, etc not included
- Single  $R=1$  jet trigger  $p_T > 1.2$  TeV

The goal of the challenge was for participants to analyze each box and

1. Decide whether or not it contains new physics
2. Characterize the new physics, if it's there

# LHC Olympics 2020: Black Boxes

<https://doi.org/10.5281/zenodo.3547721>

In 2019, Gregor Kasieczka, Ben Nachman and I initiated the LHC Olympics 2020 Challenge.

It consisted of three “black boxes” of simulated data (LHC)

- 1 million events each
- 4-vectors of

Stay tuned for Ben's talk to see what was in the boxes and the outcome of the challenge

oser  $p_T > 1.2 \text{ TeV}$

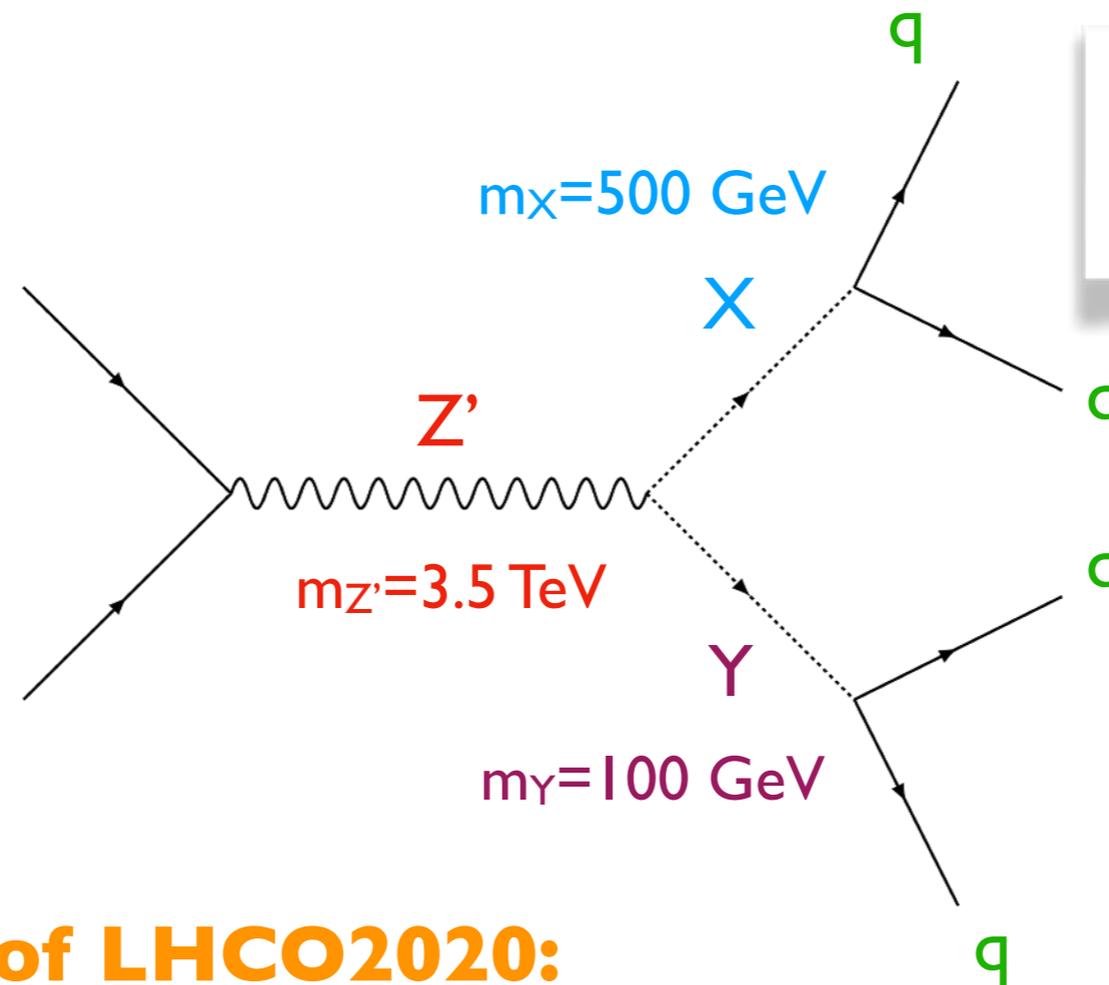
The goal of the challenge was for participants to analyze each box and

1. Decide whether or not it contains new physics
2. Characterize the new physics, if it's there

# LHC Olympics 2020: R&D Dataset

<https://doi.org/10.5281/zenodo.2629072>

Prior to the challenge, we also released a labeled R&D dataset consisting of 1M QCD dijet events and 100k signal events

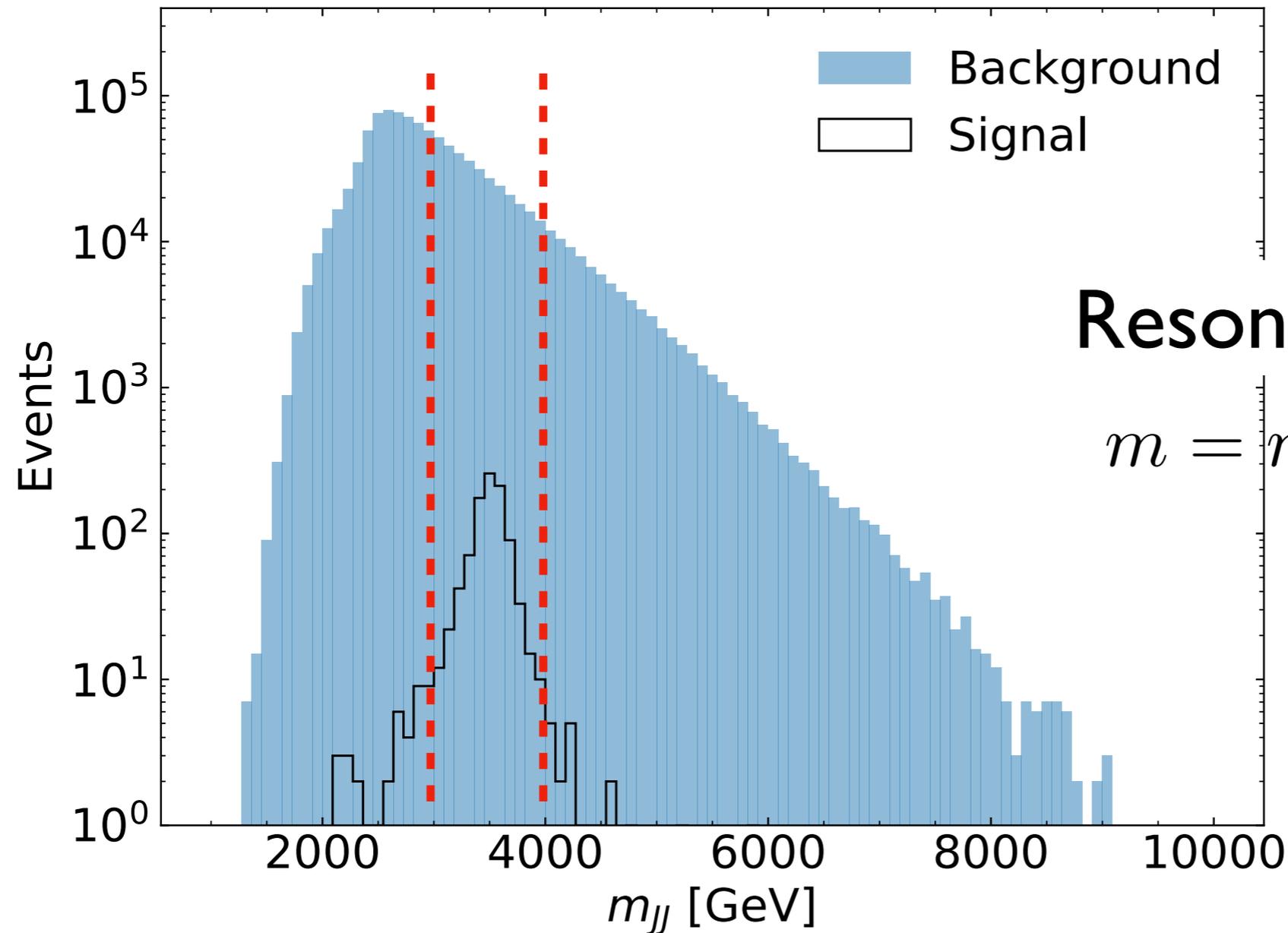


No explicit search at the LHC for this scenario!

**Unofficial theme of LHCO2020:**

*“enhancing the bump hunt”*

# LHC Olympics 2020: R&D Dataset

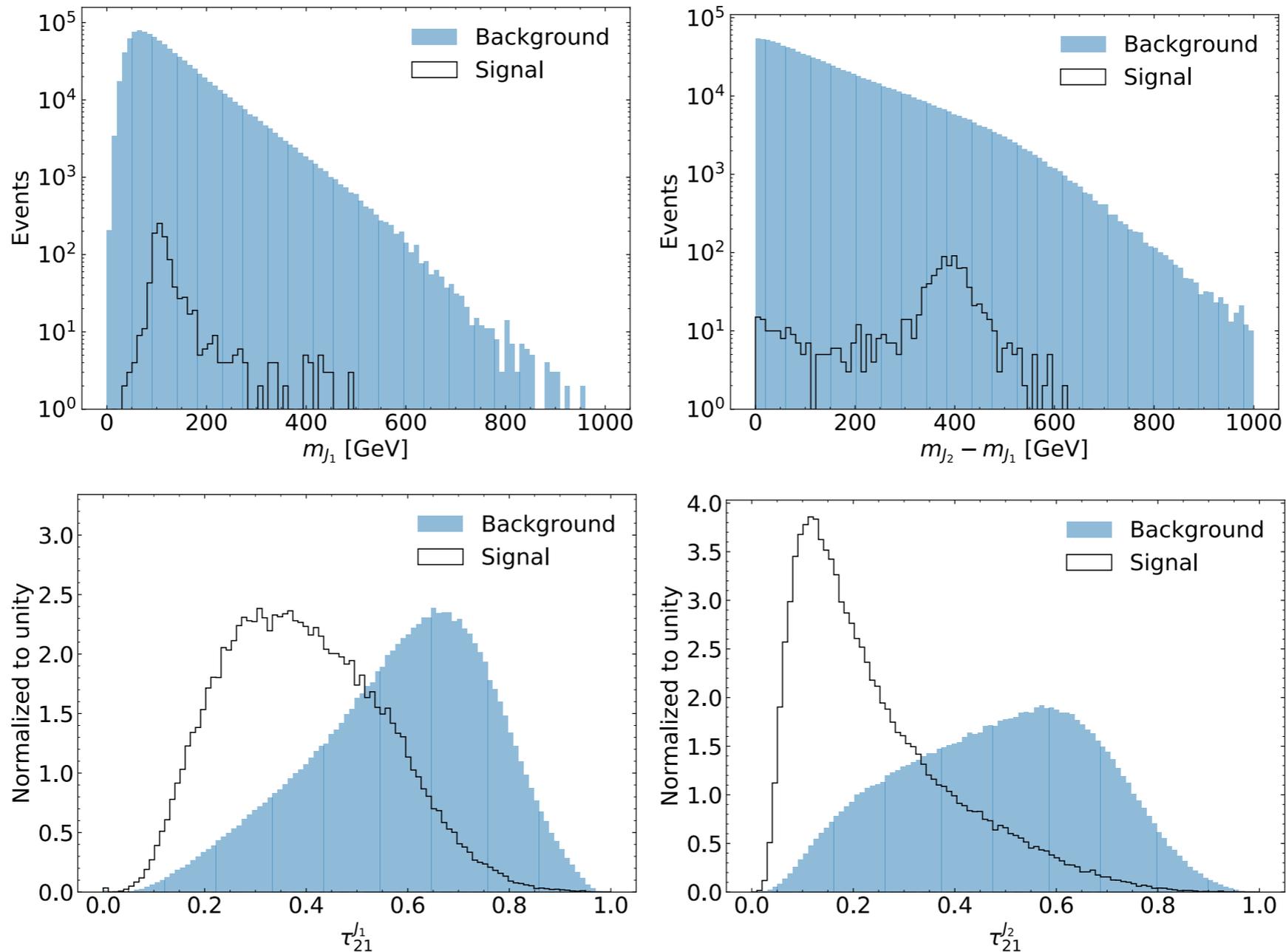


Resonant feature

$$m = m_{Z'} = m_{JJ}$$

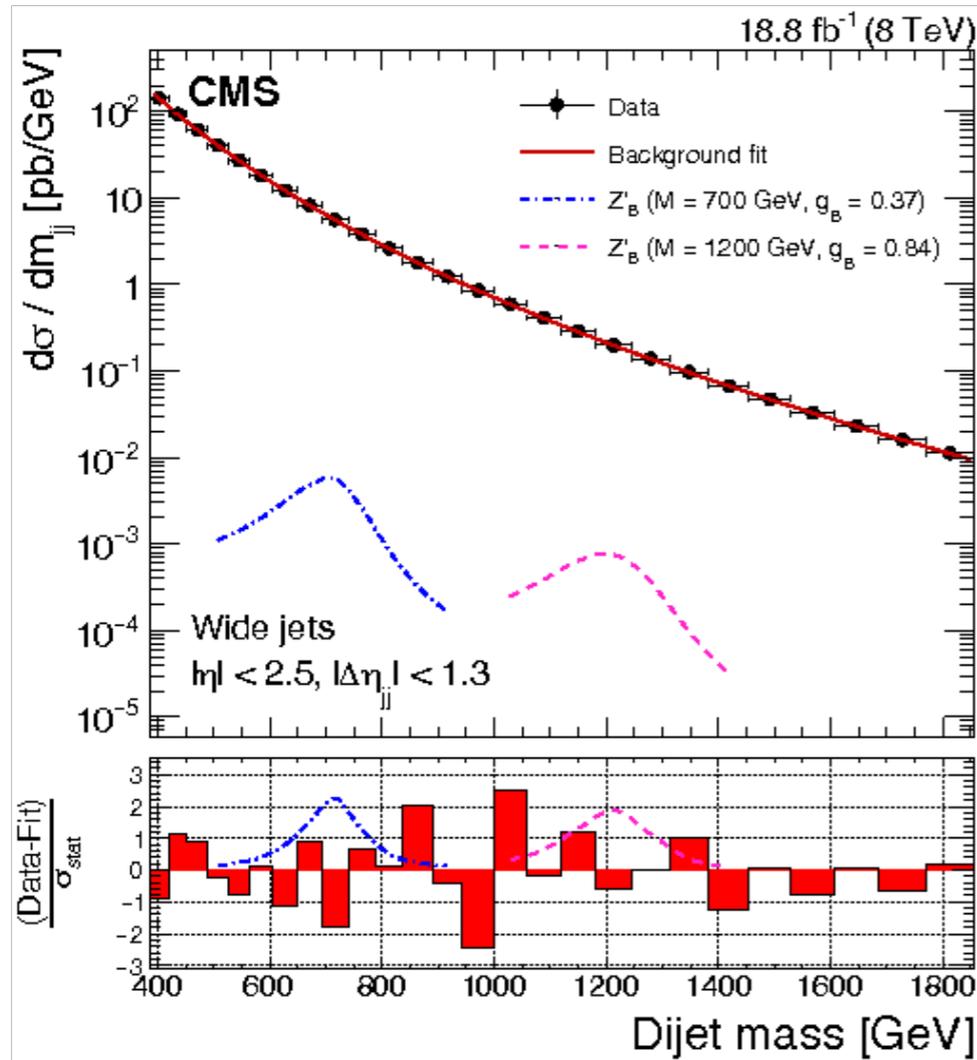
Benchmark  
signal strength:  $S=500, B=500,000, B_{SR}=61,000$   
 $S/B_{SR} \sim 6 \times 10^{-3}, S/\sqrt{B_{SR}} \sim 1.5$

# LHC Olympics 2020: R&D Dataset

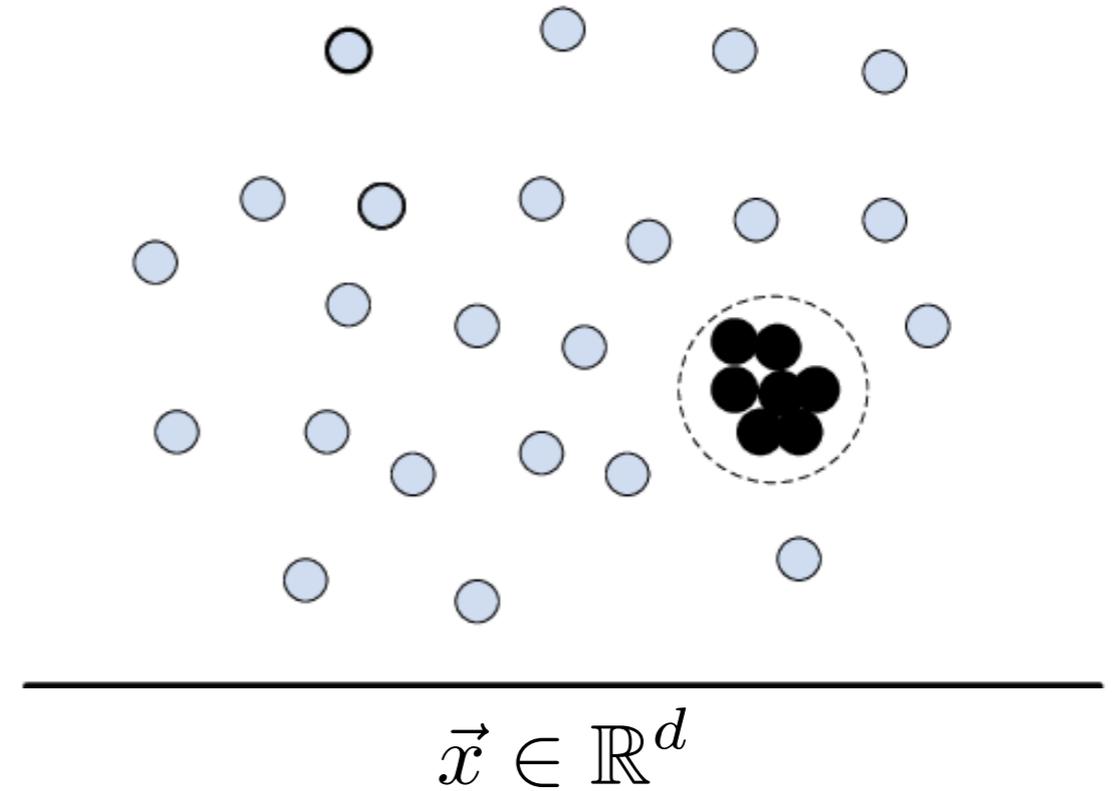


Additional features:  $x = (m_{J_1}, m_{J_2}, \tau_{21}^{J_1}, \tau_{21}^{J_2})$

# Enhancing the bump hunt



×



primary resonant feature (mJJ)

additional features

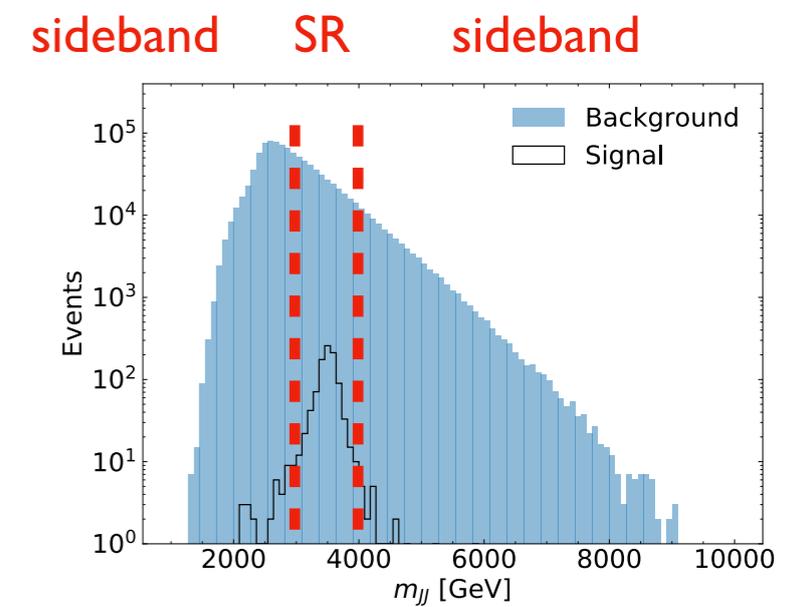
If the signal is localized in additional features, can we find it in a model-independent way?

# Enhancing the bump hunt

The optimal model-agnostic discriminant would be

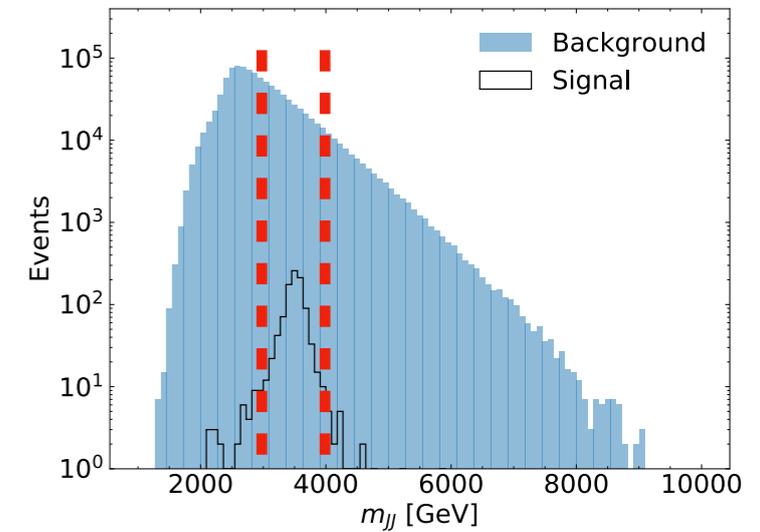
$$R(x) = \frac{\mathcal{P}(x|data)}{\mathcal{P}(x|bg)}$$

**Key idea: use sidebands in  $m_{JJ}$  to model  $\mathcal{P}(x|bg)$**



# Enhancing the bump hunt

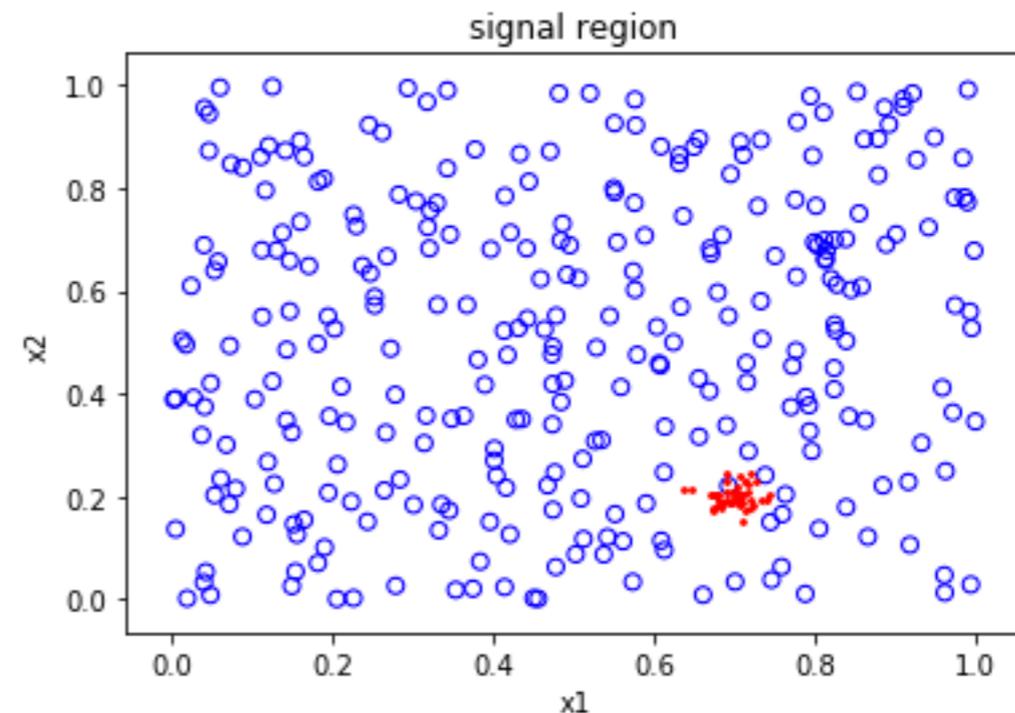
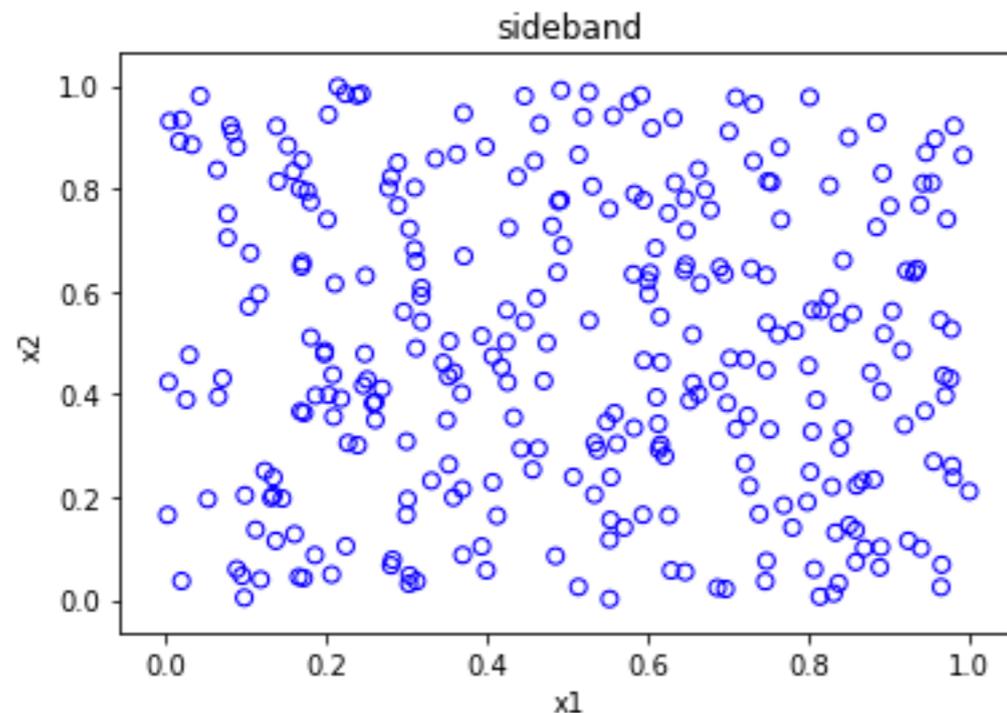
sideband SR sideband



The optimal model-agnostic discriminant would be

$$R(x) = \frac{\mathcal{P}(x|data)}{\mathcal{P}(x|bg)}$$

**Key idea: use sidebands in  $m_{jj}$  to model  $\mathcal{P}(x|bg)$**



# ANODE: Anomaly Detection with Density Estimation

Nachman & DS 2001.04990

**Example of a new approach inspired by LHCO2020.**

(See Ben's talk for additional new approaches!)

Use **neural density estimation** to directly learn the conditional probability densities from the data

$$\mathcal{P}(x|data; m_{JJ} \in SR)$$

$$\mathcal{P}(x|data; m_{JJ} \notin SR) = \mathcal{P}(x|bg; m_{JJ} \notin SR)$$

*interpolate in (x,mJJ)*

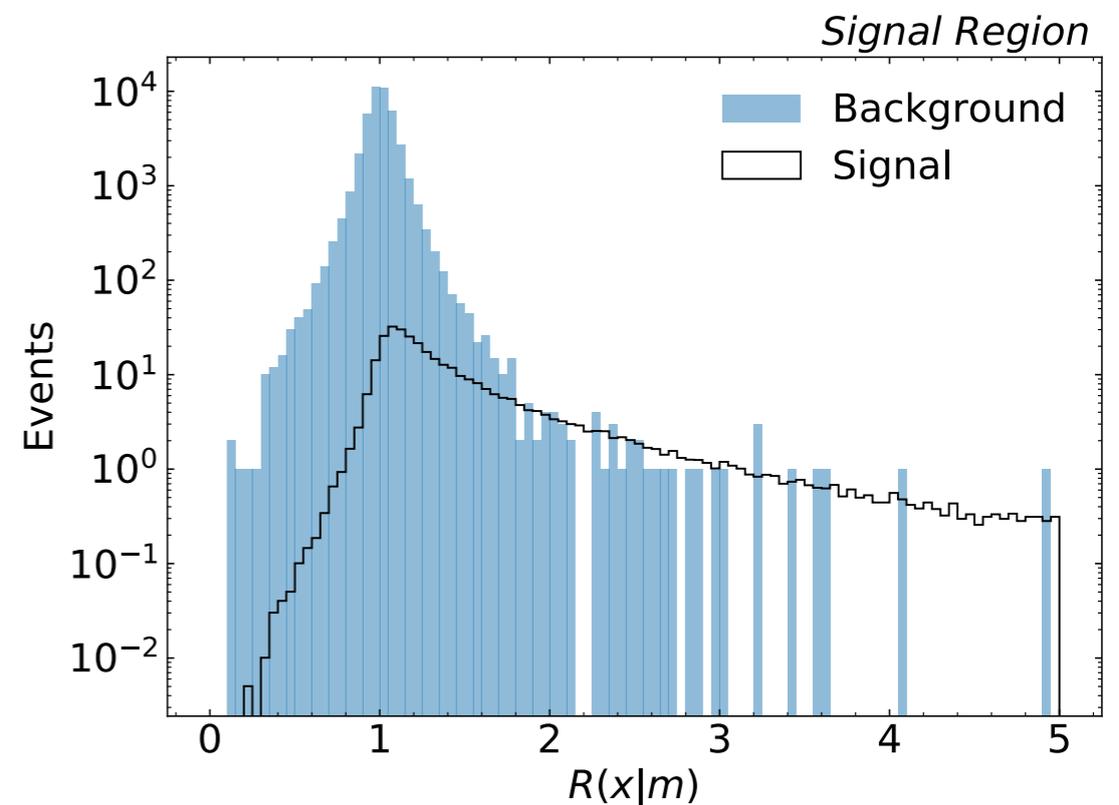
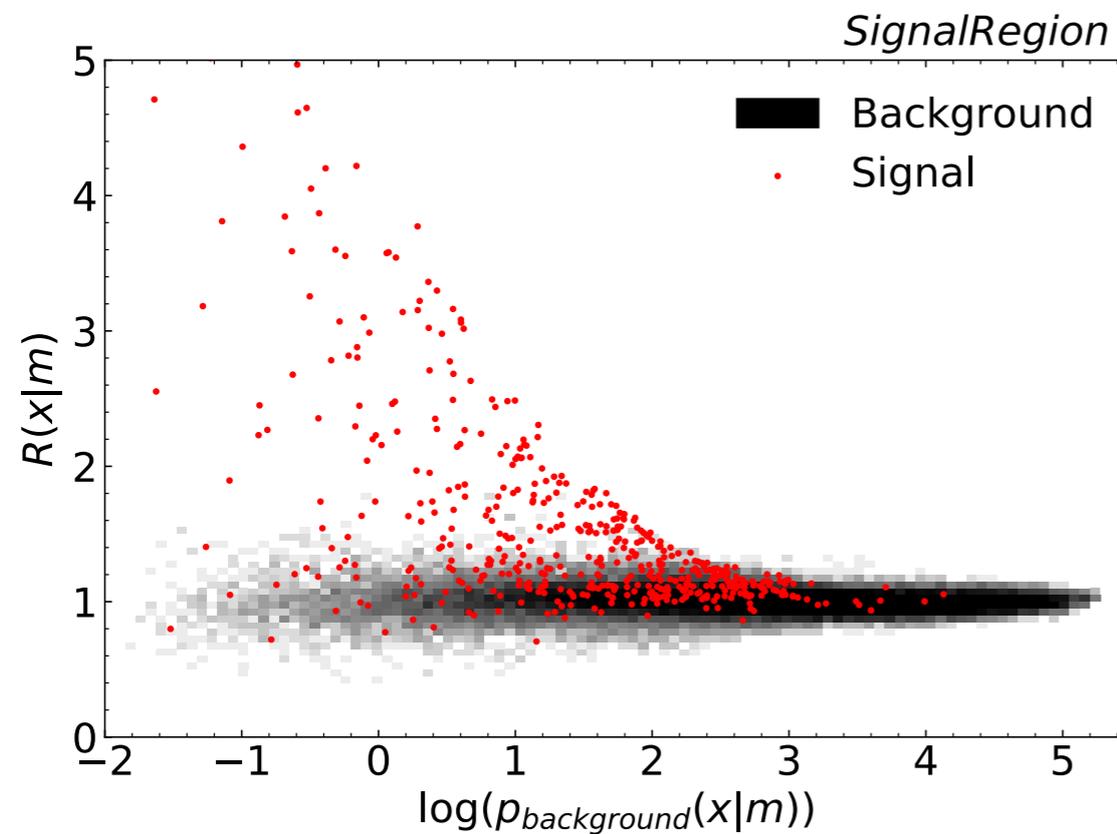
$$\mathcal{P}(x|bg; m_{JJ} \in SR)$$

Construct the likelihood ratio:

$$R(x) = \frac{\mathcal{P}(x|data; m_{JJ} \in SR)}{\mathcal{P}(x|bg; m_{JJ} \in SR)}$$

# ANODE: Results on LHCO R&D Dataset

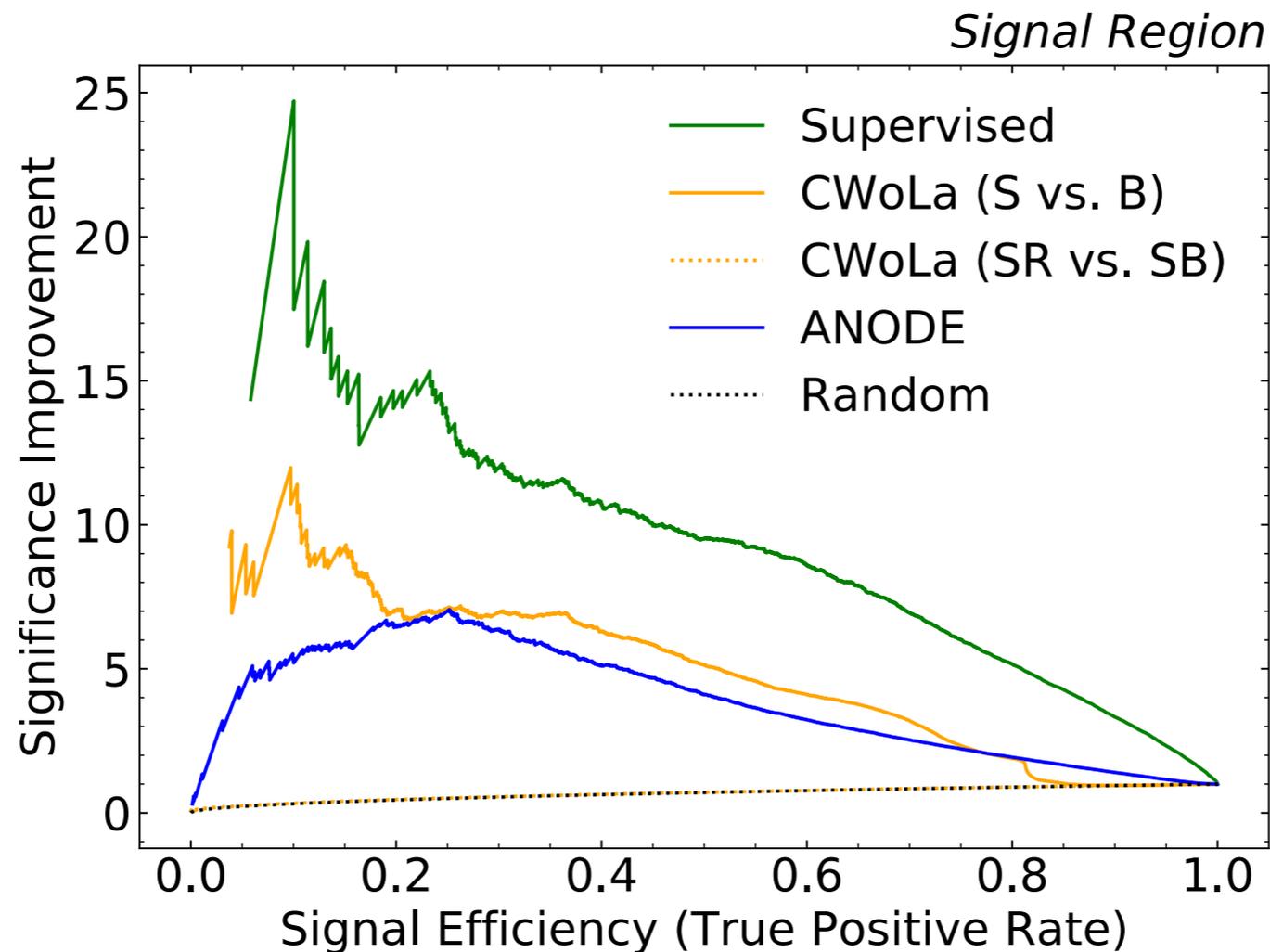
Nachman & DS 2001.04990



The method works! ANODE is sensitive to the signal!

# ANODE: Results on LHCO R&D Dataset

Nachman & DS 2001.04990



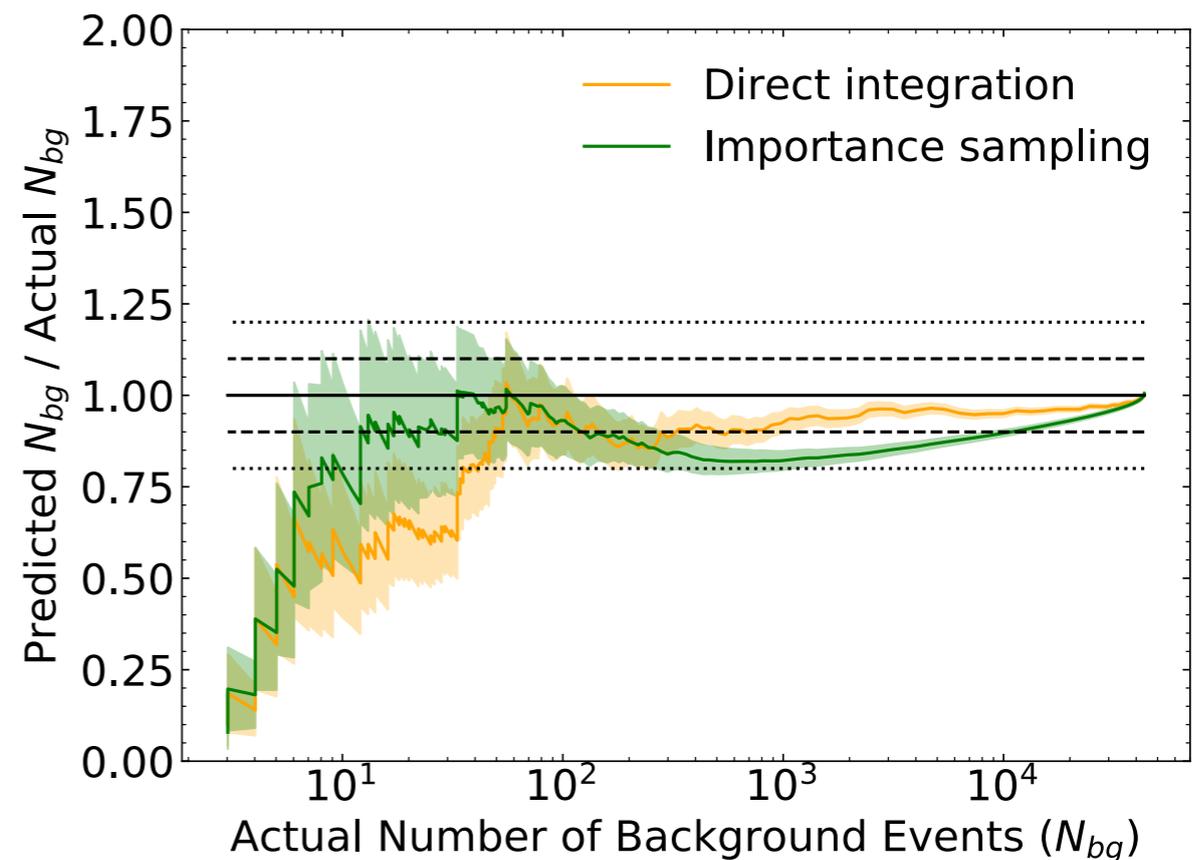
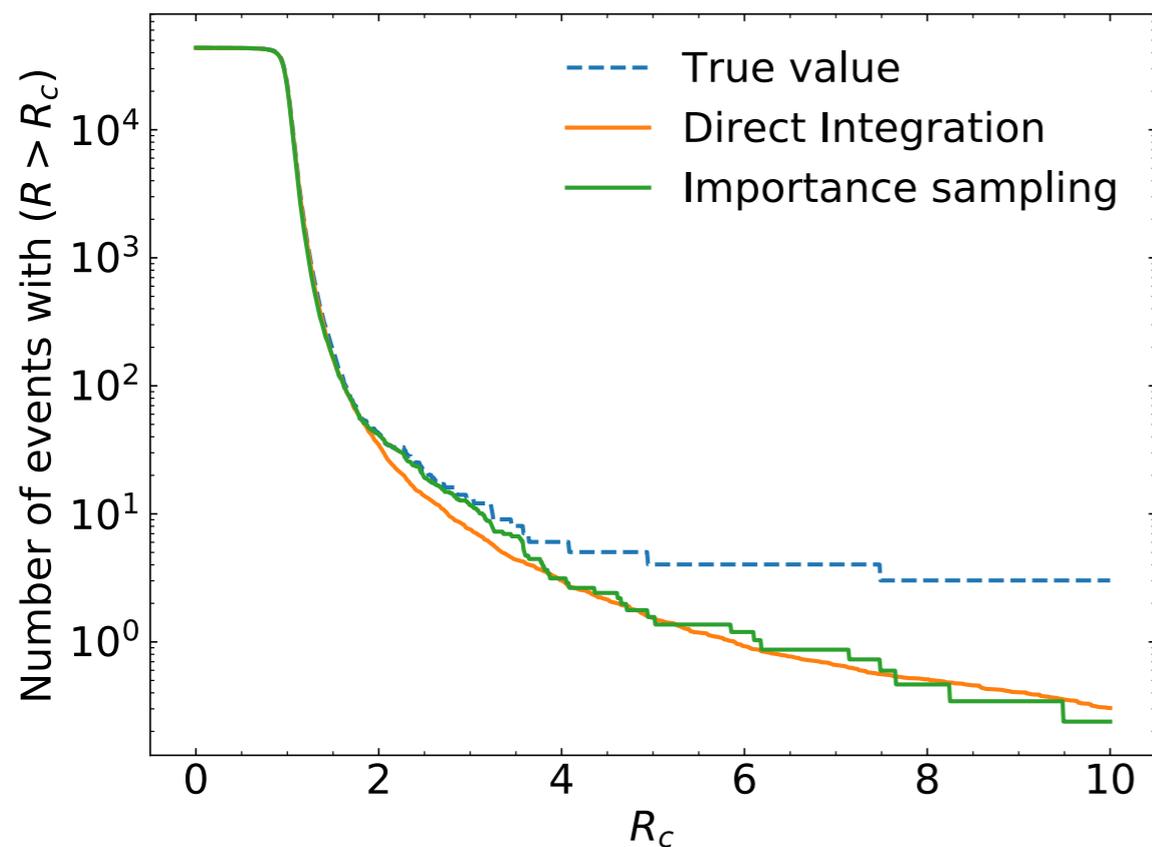
Can enhance the significance of the bump hunt by a factor of up to 7!

**1.5 $\sigma$  (dijet bump hunt)  $\Rightarrow$  10 $\sigma$  (ANODE+dijet bump hunt)**

# ANODE: Results on LHCO R&D Dataset

Nachman & DS 2001.04990

Novel aspect of ANODE: can estimate backgrounds directly with  $P(x|bg; m \in SR)$



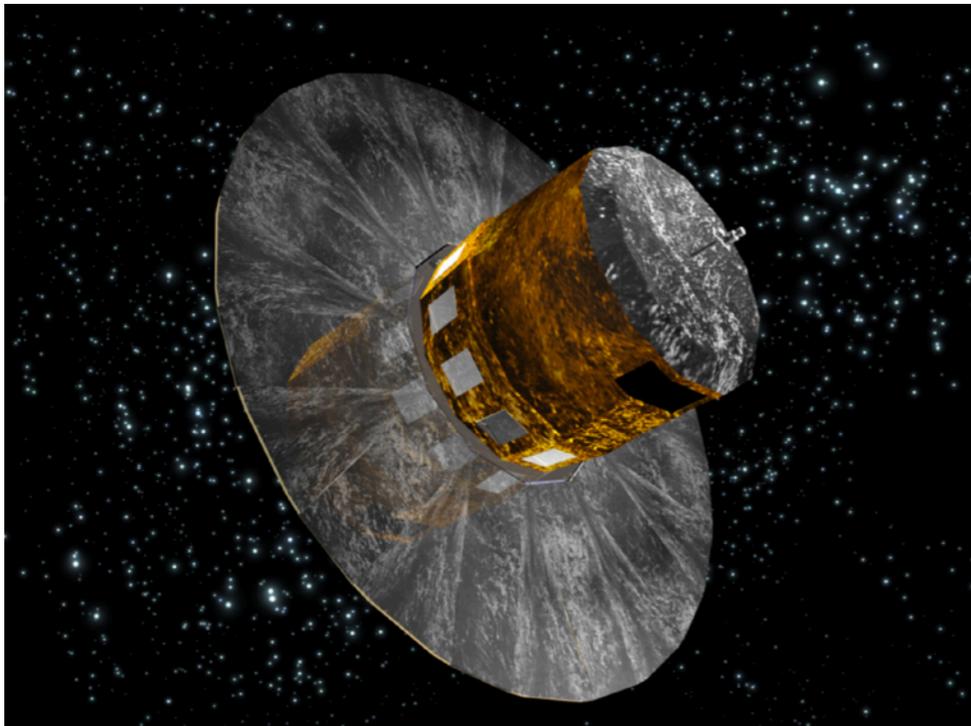
# Via Machinae: ANODE IN SPACE

DS, Matt Buckley, Lina Necib & John Tamasas, in preparation

ANODE is a completely general method for identifying multivariate overdensities in data using sidebands. It could have many diverse applications.

We are currently using it to find **stellar streams** in Gaia data.

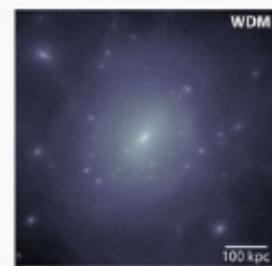
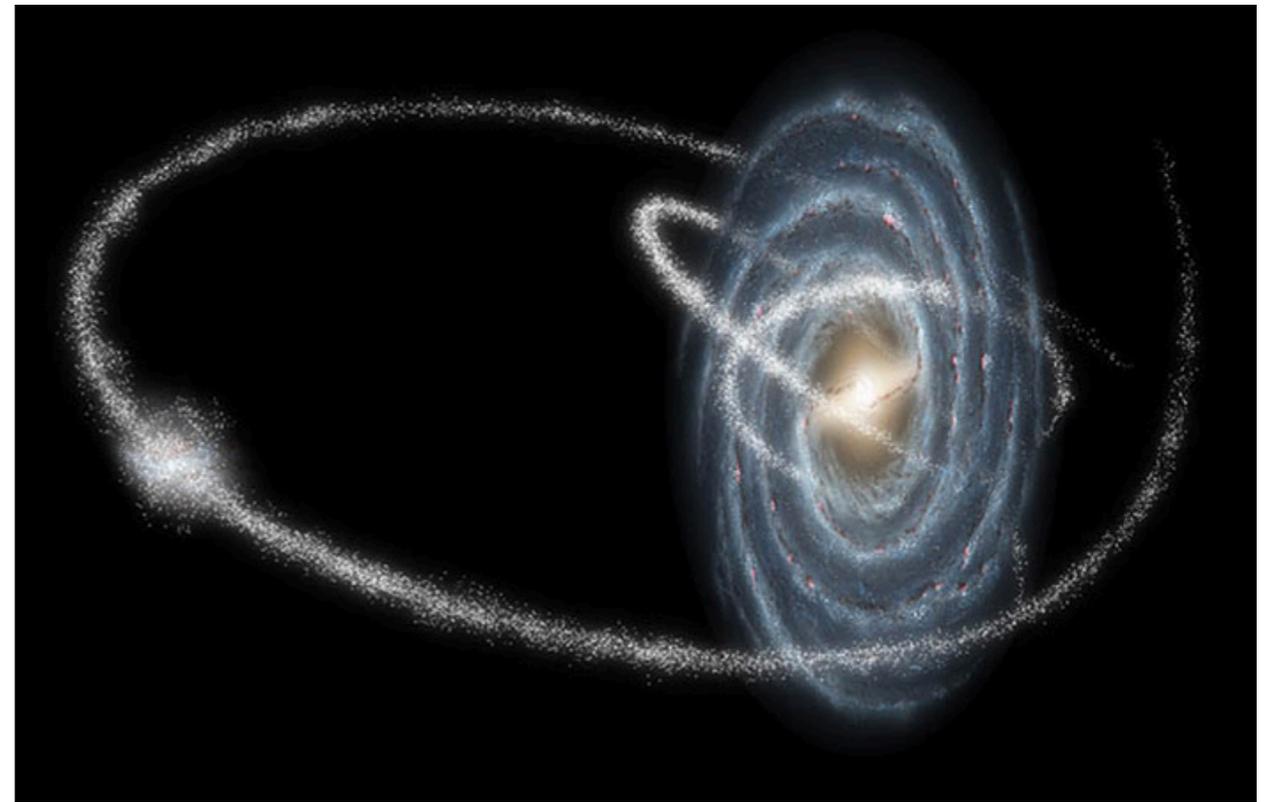
Gaia satellite:



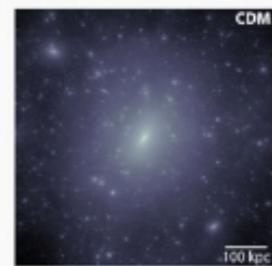
- Launched in 2013; extended to 2025
- Mission: map out the full 6d phase space of the stars in our galaxy
- Angular positions, velocities, color and magnitude of over 1 billion stars in our galaxy
- radial positions and velocities for a smaller subset of nearby stars (not used in this work)

# Stellar streams

Cold stellar streams are tidally-stripped remnants of globular clusters and dwarf galaxies, falling into and orbiting our galaxy.



Stellar stream in a smooth galaxy



Stellar stream in a clumpy galaxy

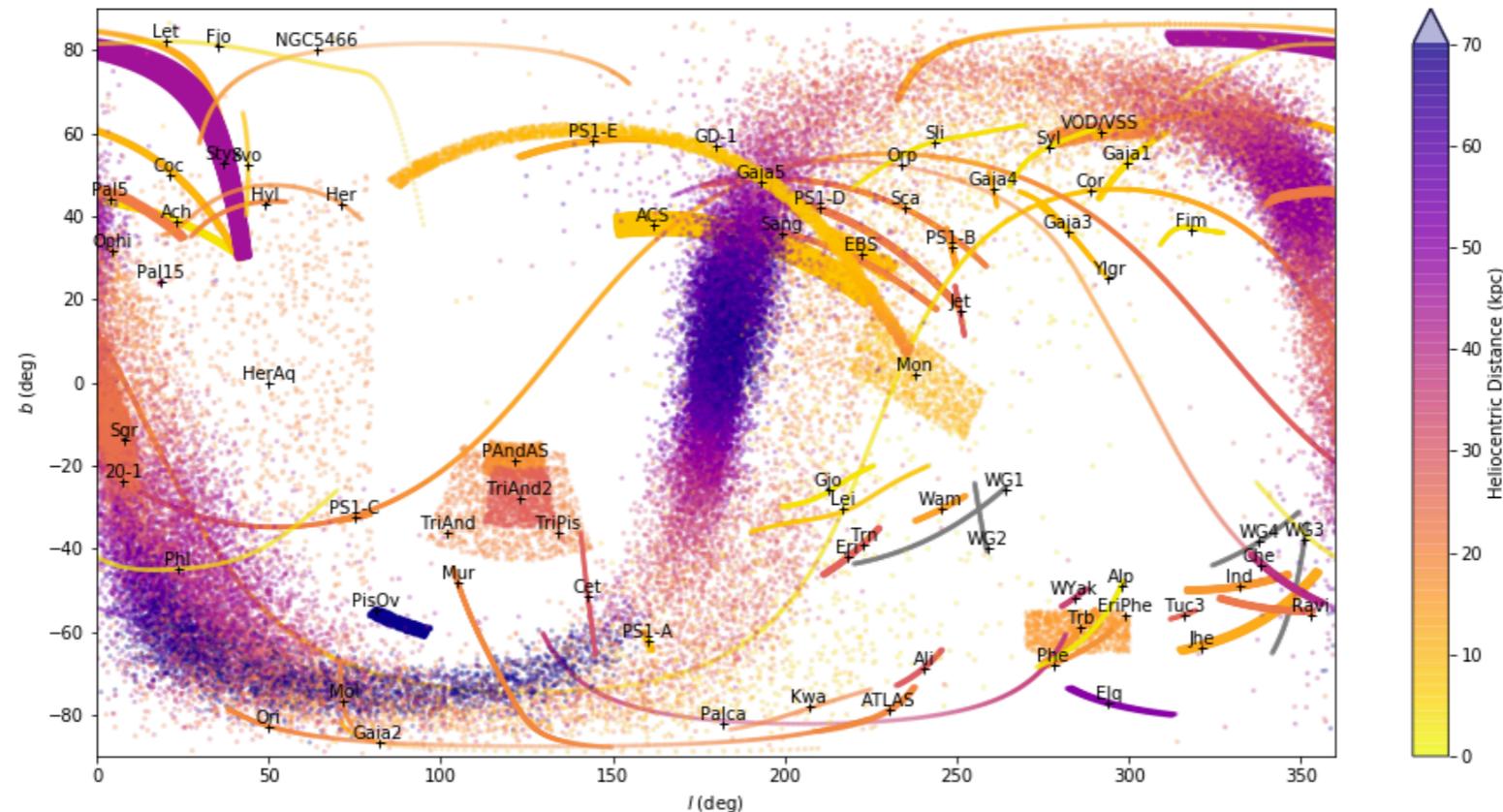
Bonaca et al. (2014)

They are very interesting objects of study for astrophysicists and particle physicists.

**In particular, they could be unique probes into dark matter substructure.**

# Stream finding

<https://github.com/cmateu/galstreams>

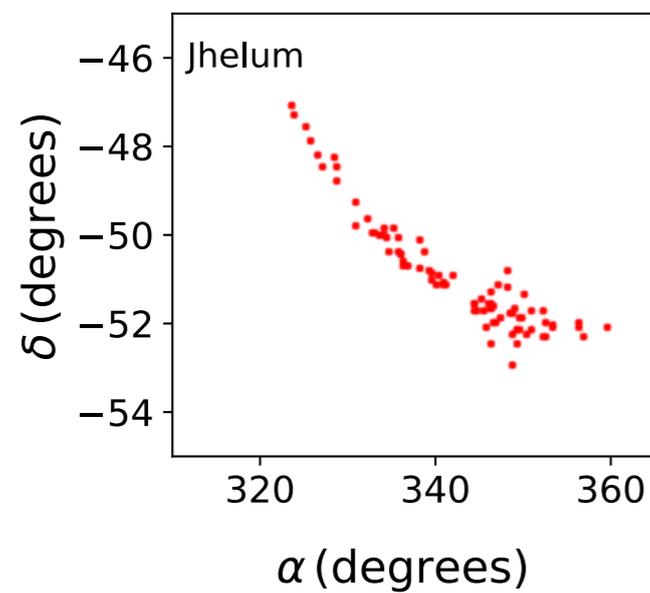


Existing methods (eg **STREAMFINDER**, Malhan & Ibata 2018) have found many new streams in the Gaia data, but they make a number of model-dependent assumptions (form of the galactic potential, orbits, isochrones, ...).

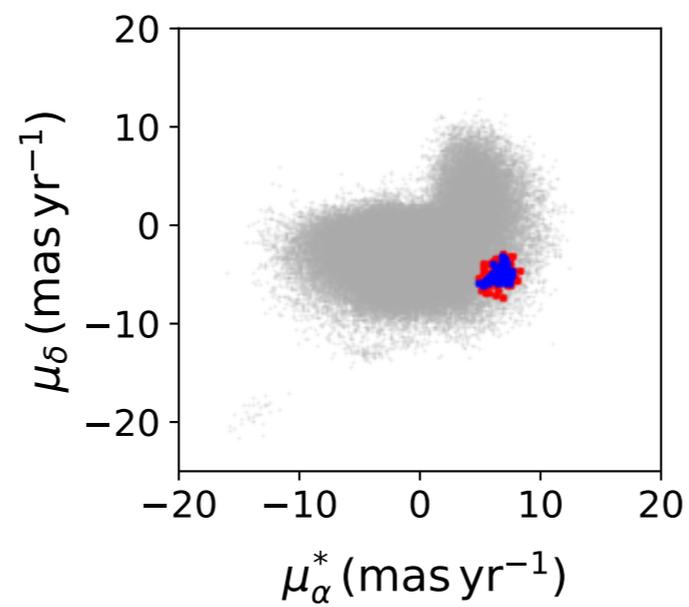
We were interested in whether unsupervised ML could be used to find streams more model-independently.

# Stream finding

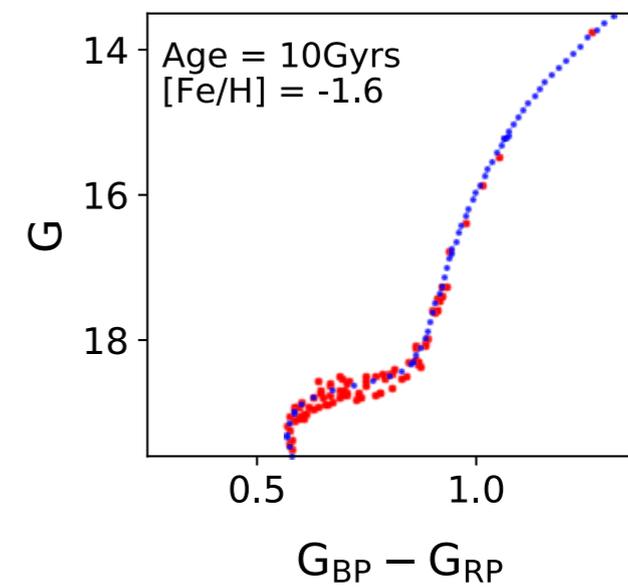
from Malhan et al 2018



position



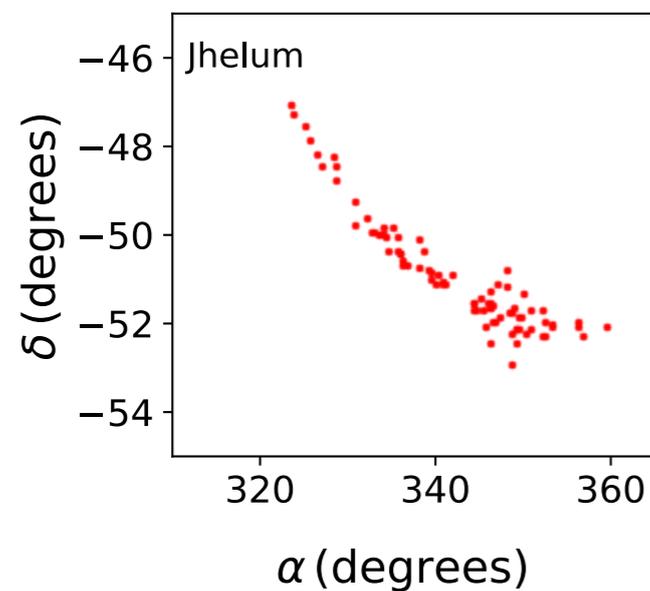
velocity



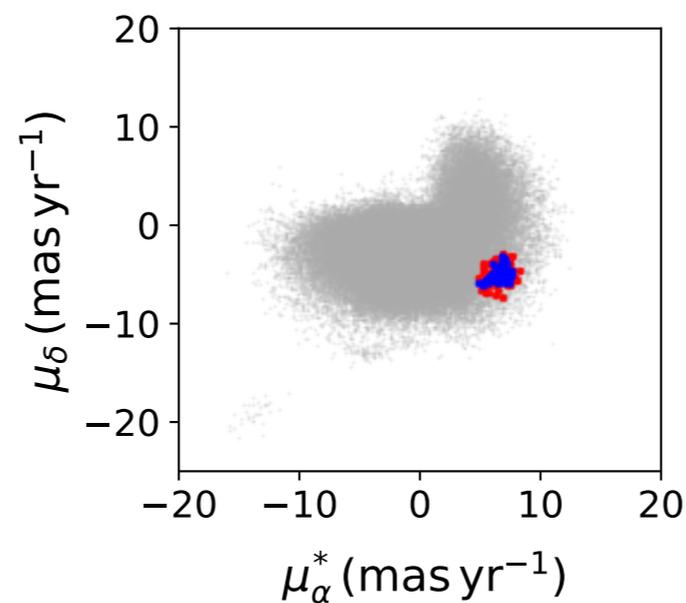
photometry

# Stream finding

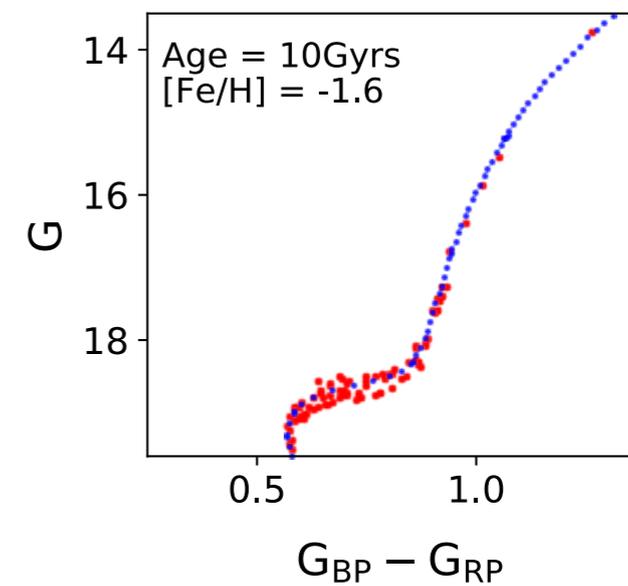
from Malhan et al 2018



position



velocity

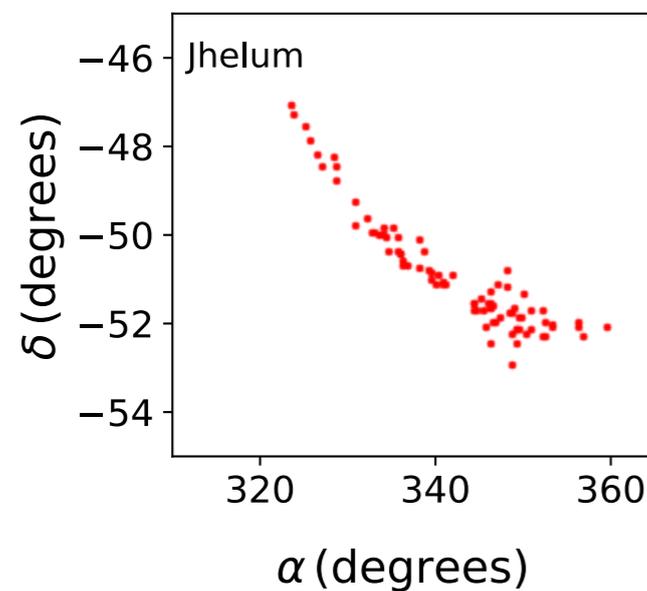


photometry

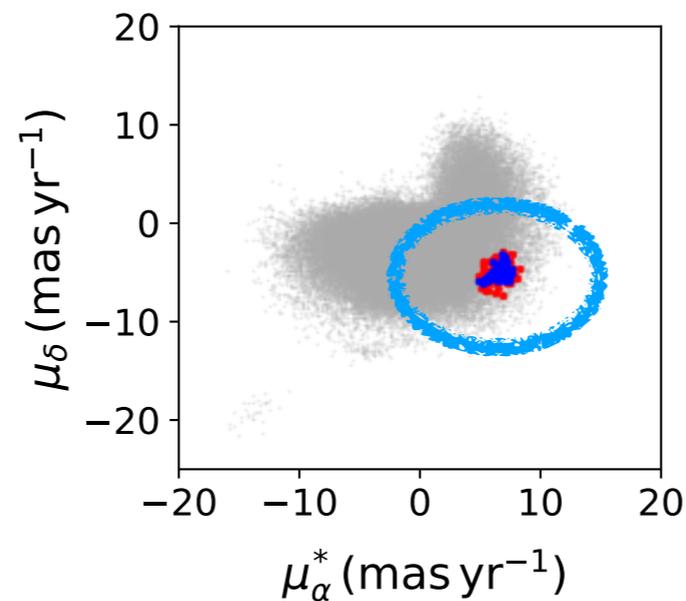
Idea: streams are local overdensities in position, velocity and photometric space.

# Stream finding

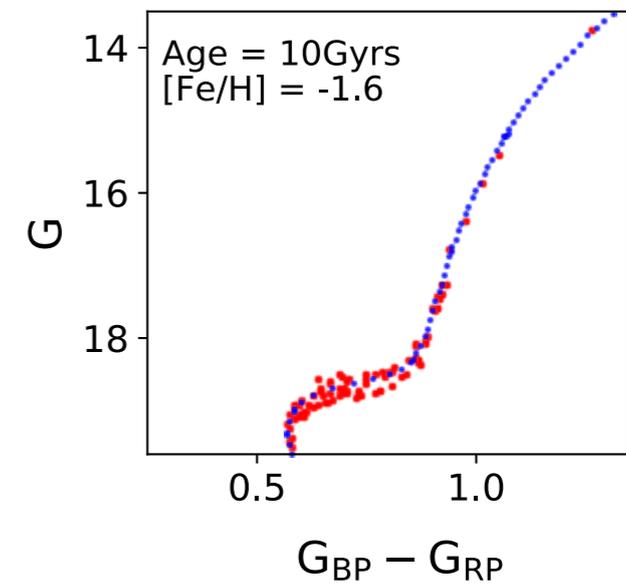
from Malhan et al 2018



position



velocity



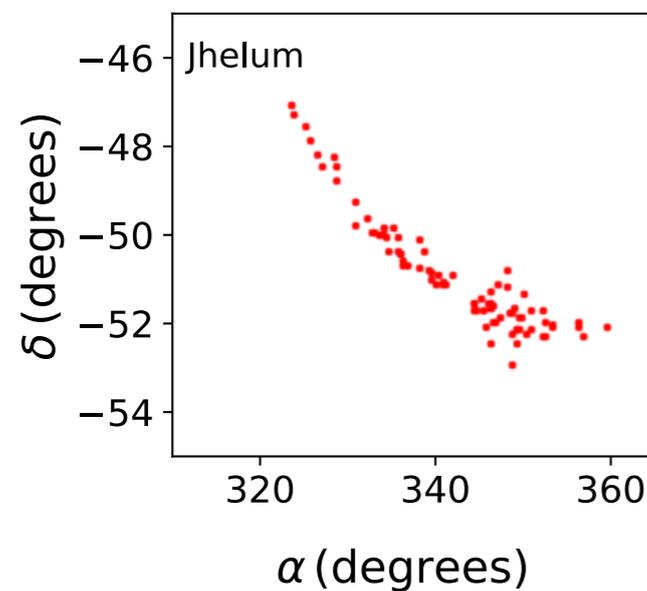
photometry

Idea: streams are local overdensities in position, velocity and photometric space.

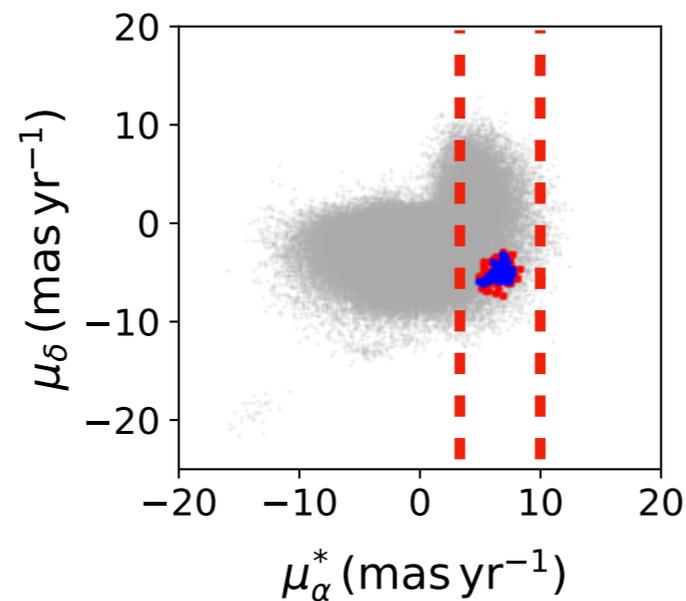
Since they are **cold**, the stars in the stream are clustered in velocity.

# Stream finding

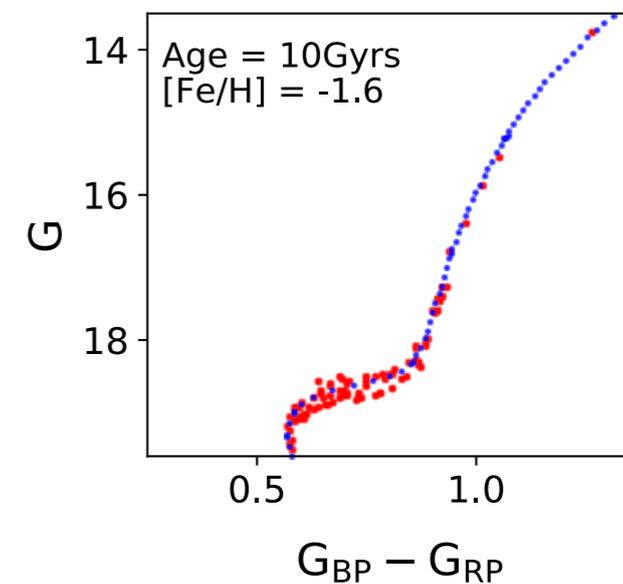
from Malhan et al 2018



position



velocity



photometry

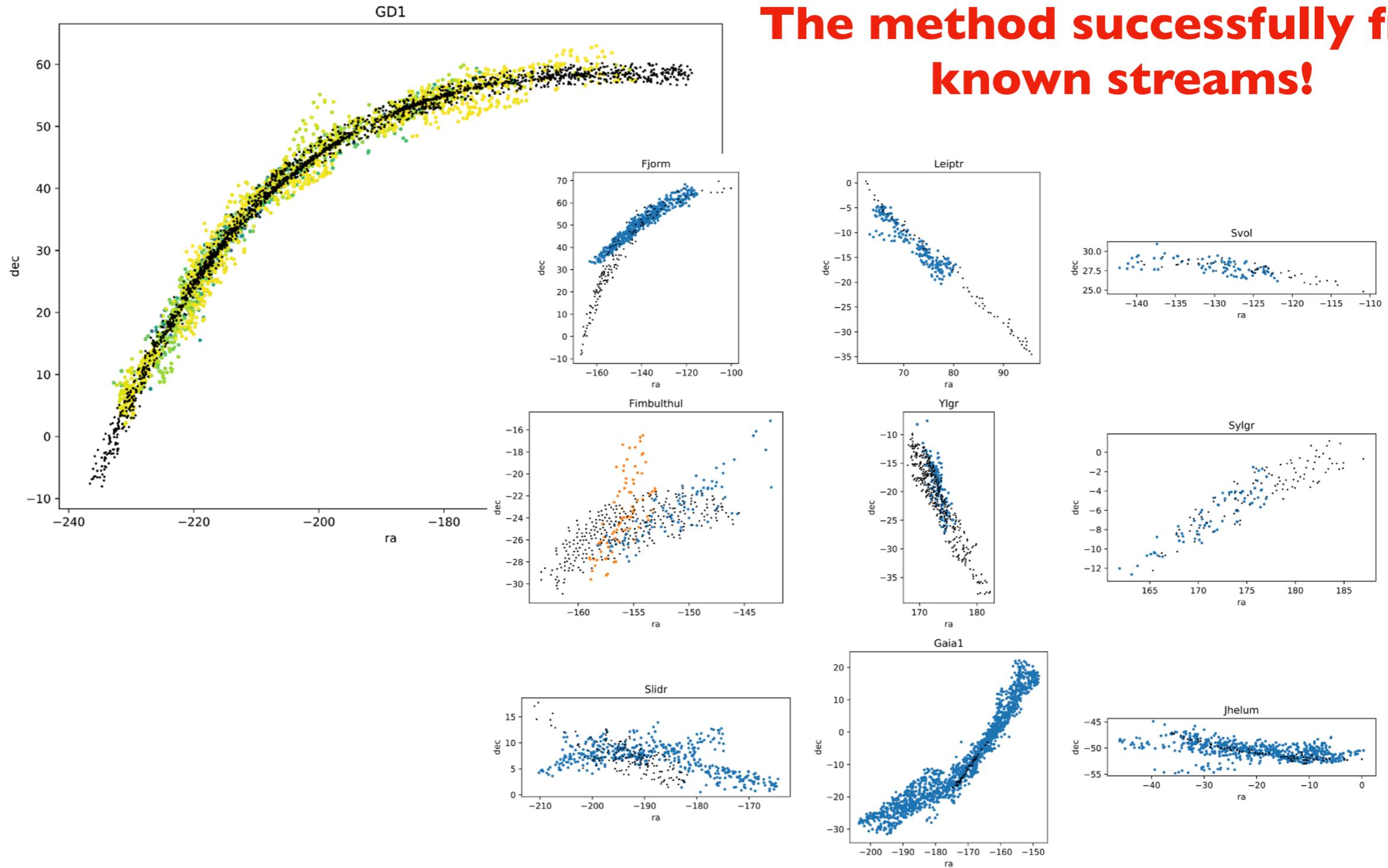
Idea: streams are local overdensities in position, velocity and photometric space.

Since they are **cold**, the stars in the stream are clustered in velocity.

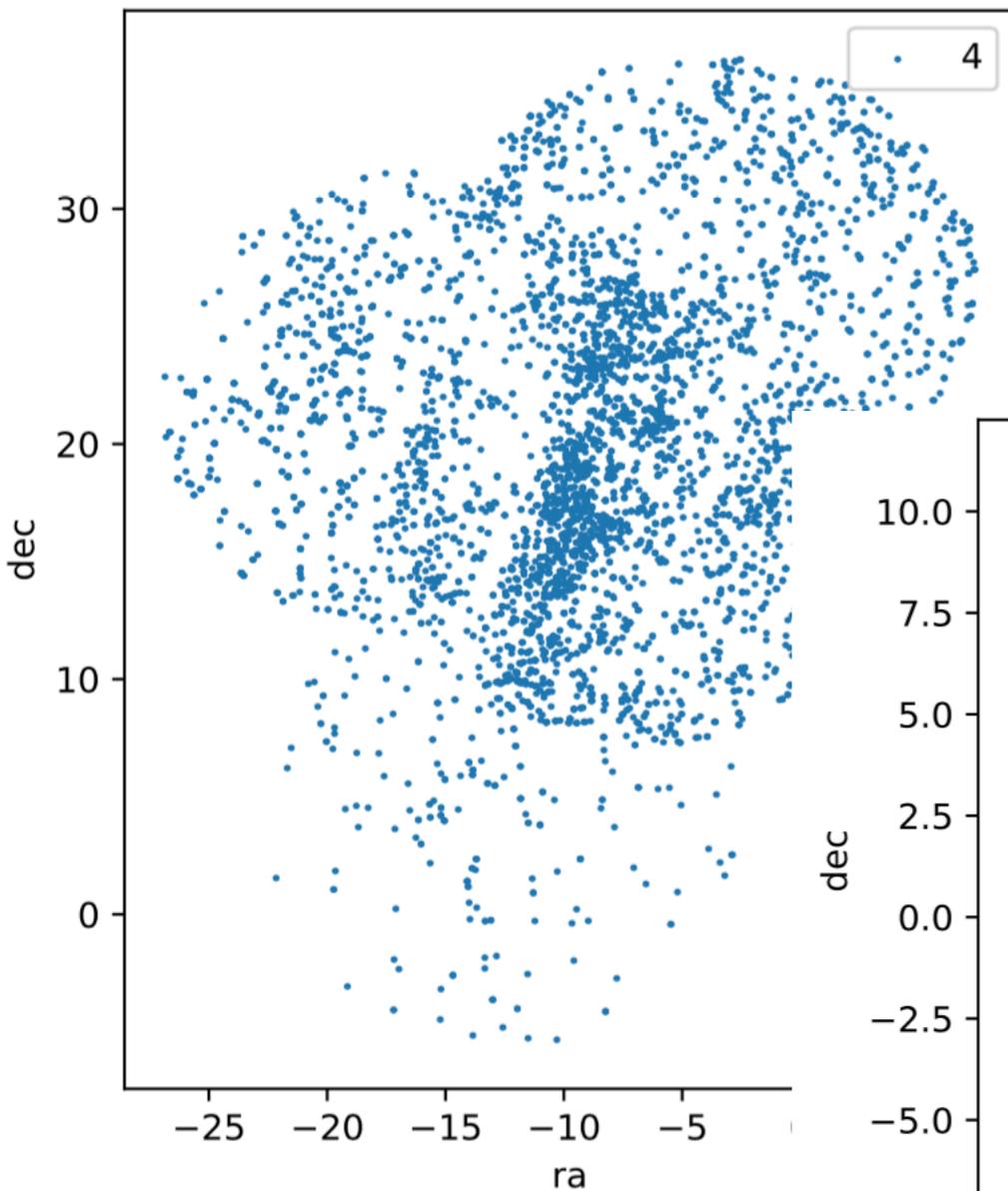
Can sideband in one of the velocities and use ANODE to look for local overdensities!

# Via Machinae: preliminary results

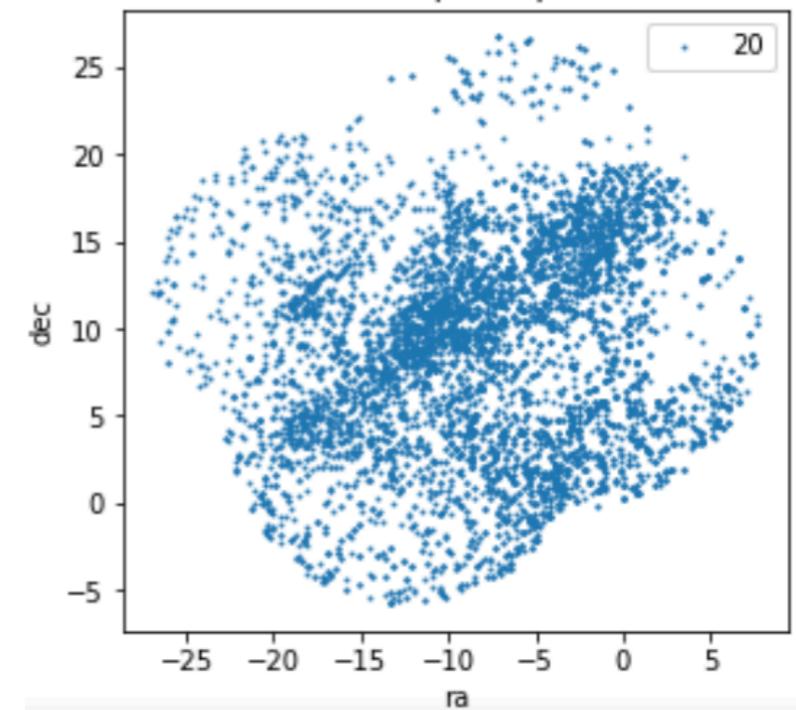
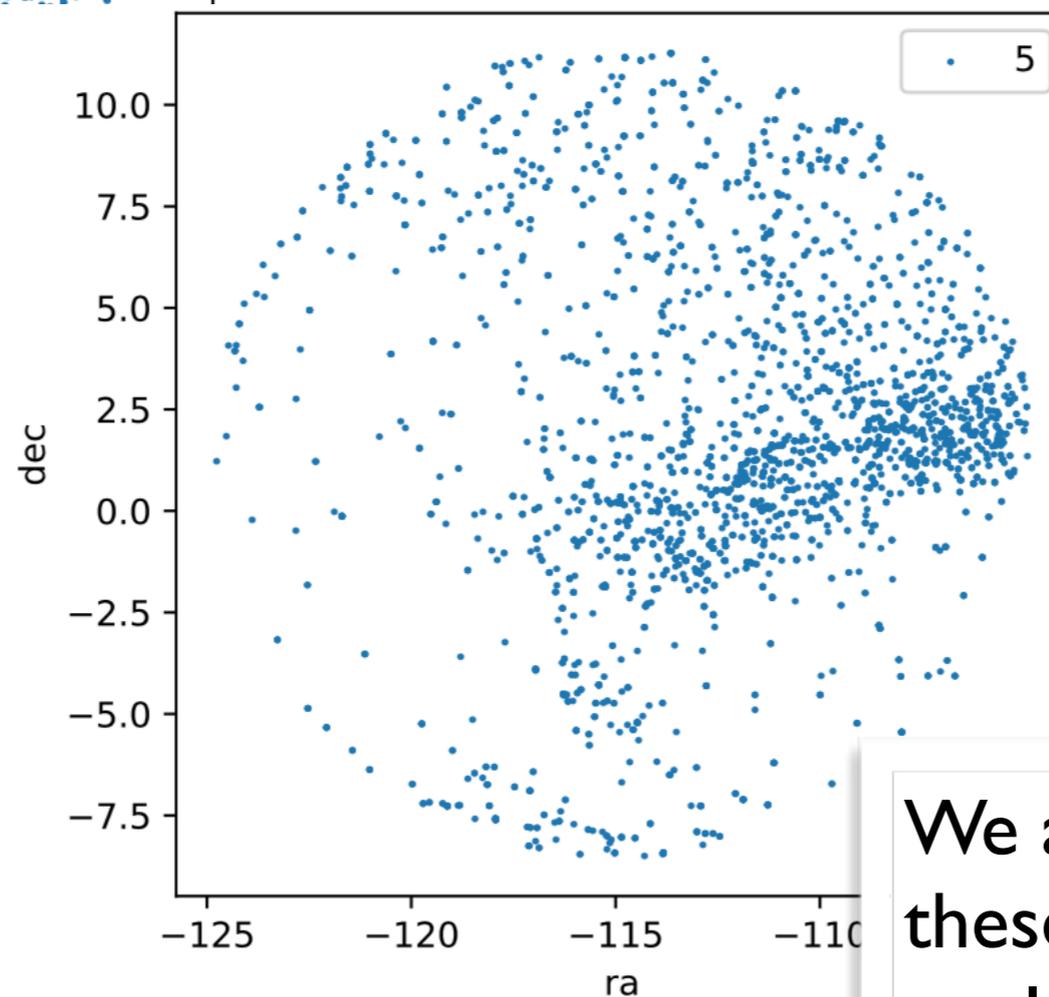
**The method successfully finds known streams!**



# Via Machinae: preliminary results



**The method also finds ~600 other stream candidates**

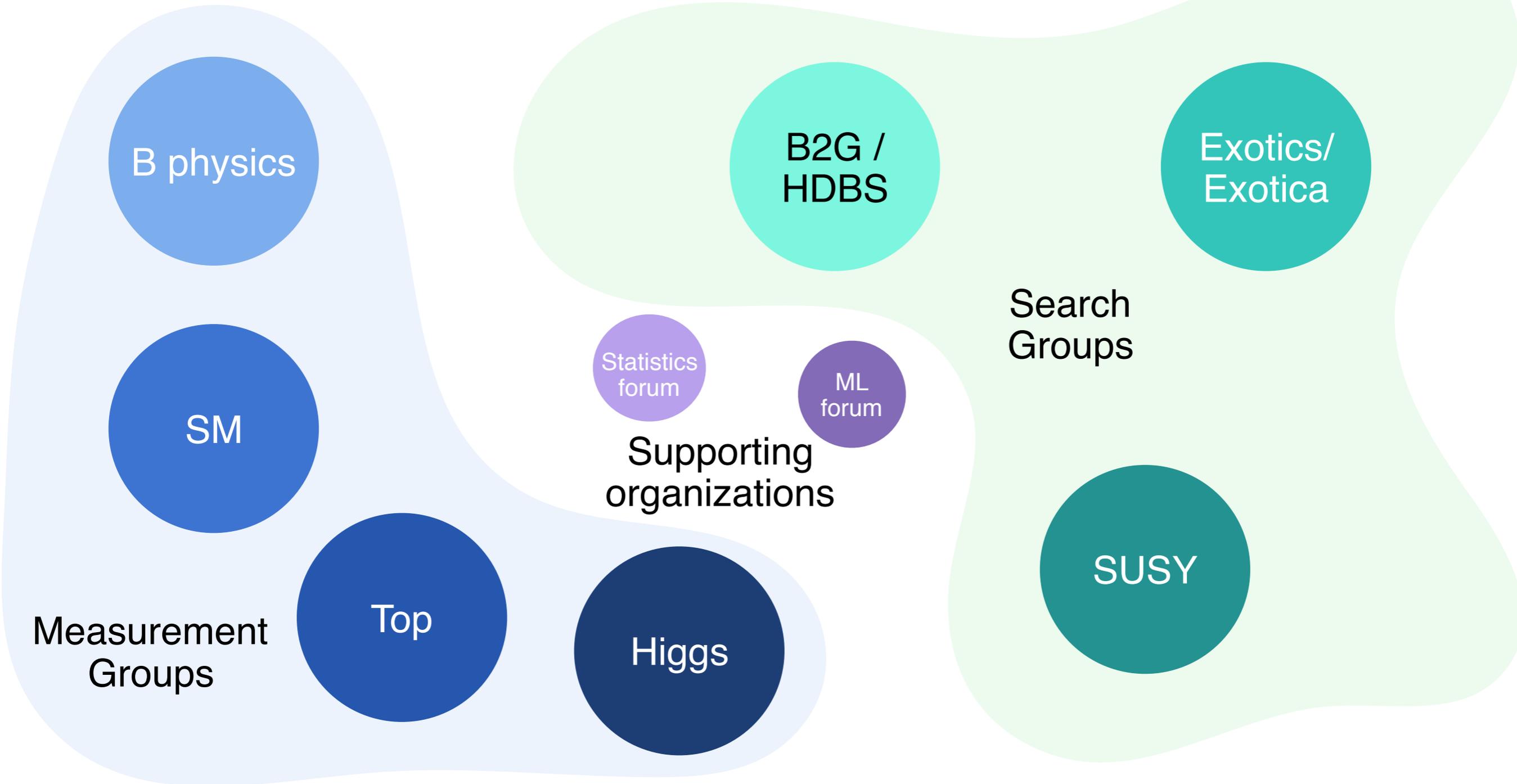


We are currently investigating these further to see if they could be real. Stay tuned!

# Summary and Outlook

- Advances in machine learning are opening up new and exciting avenues for model independent new physics searches at the LHC.
- Ideas and methods inspired by LHC problems are being applied successfully to other fields such as astronomy and astrophysics.
- The LHC Olympics 2020 provided a very useful testing ground for the development and common benchmarking of new approaches.
- Much work remains to be done in order to port these ideas over to ATLAS and CMS and implement them as actual analyses on real data.
- We need more ideas for model-independent searches at the LHC.  
This is just the beginning!

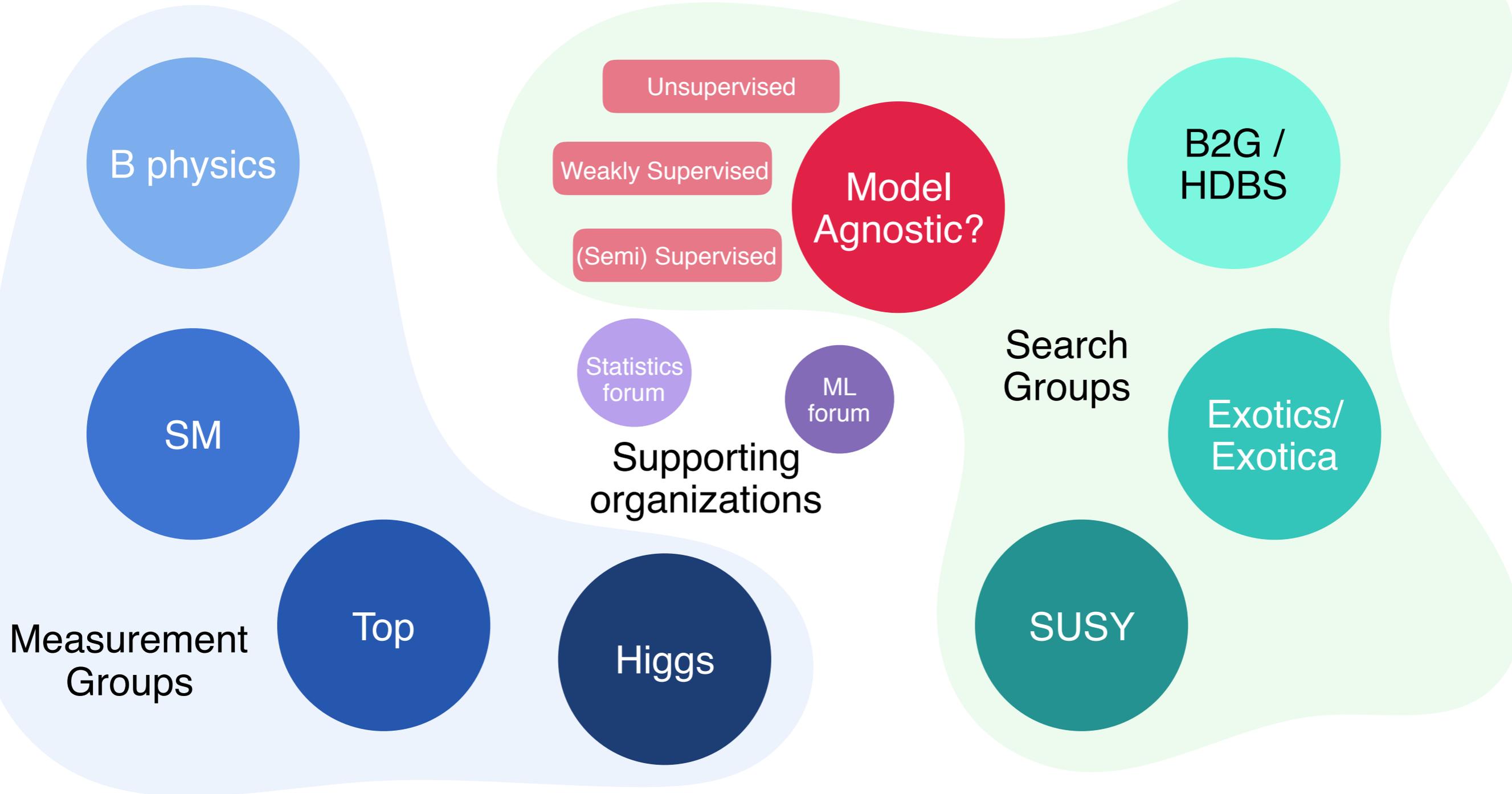
# Current Organization of Physics Analysis Groups at the LHC



**Q: Why is there no model independent search group???**

# A vision for the future...

## Future Organization of Physics Analysis Groups at the LHC??



from G. Kasieczka, B. Nachman, DS (eds), et al 2101.nextweek

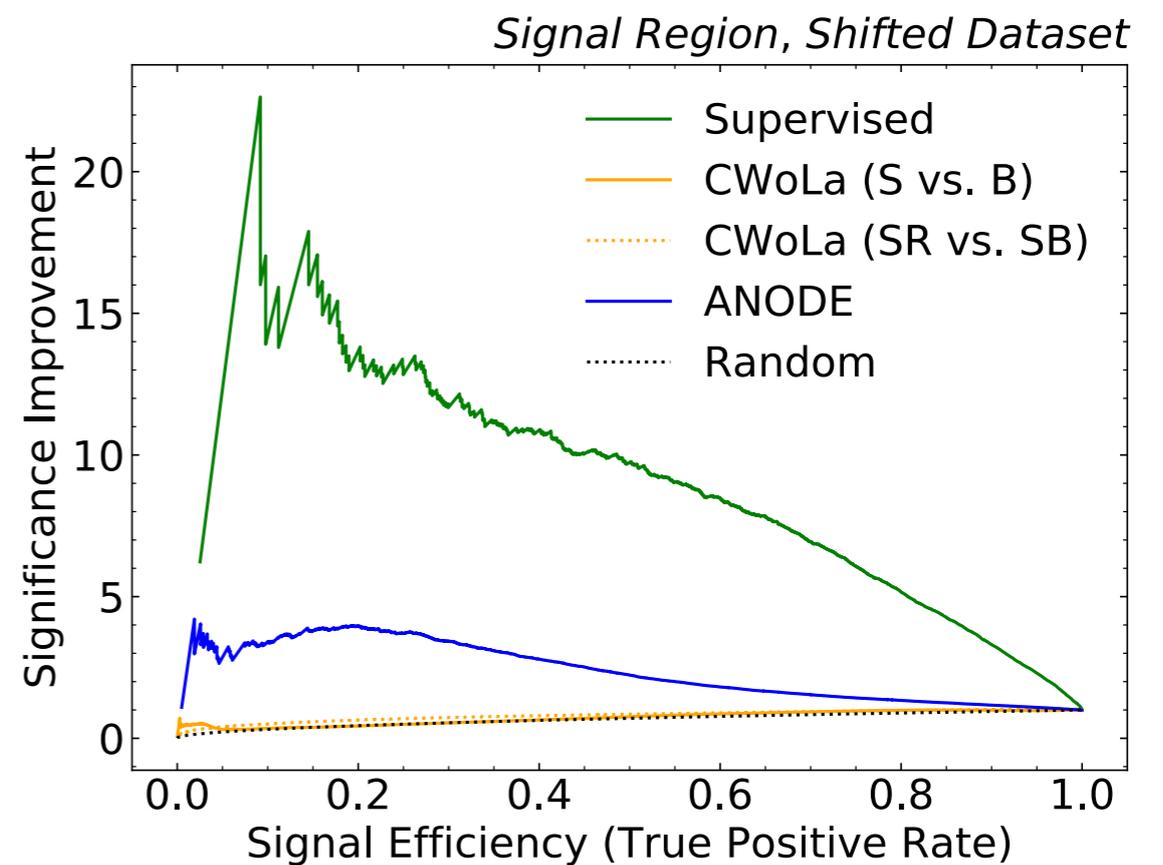
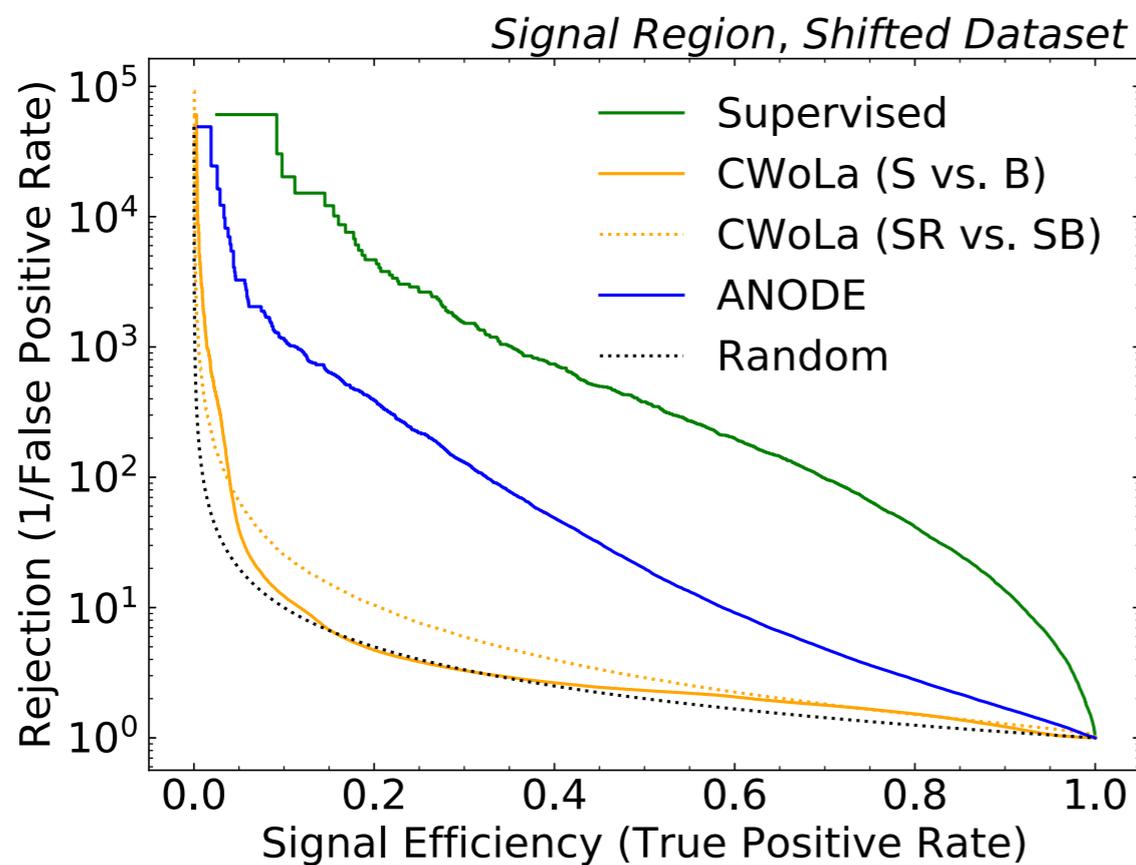
**Thanks for your attention!**

# ANODE: Results on LHCO R&D Dataset

Ben Nachman & DS 2001.04990

Can also consider performance on a feature set which is not independent of  $m$ . We introduced artificial correlations just as proof of concept:

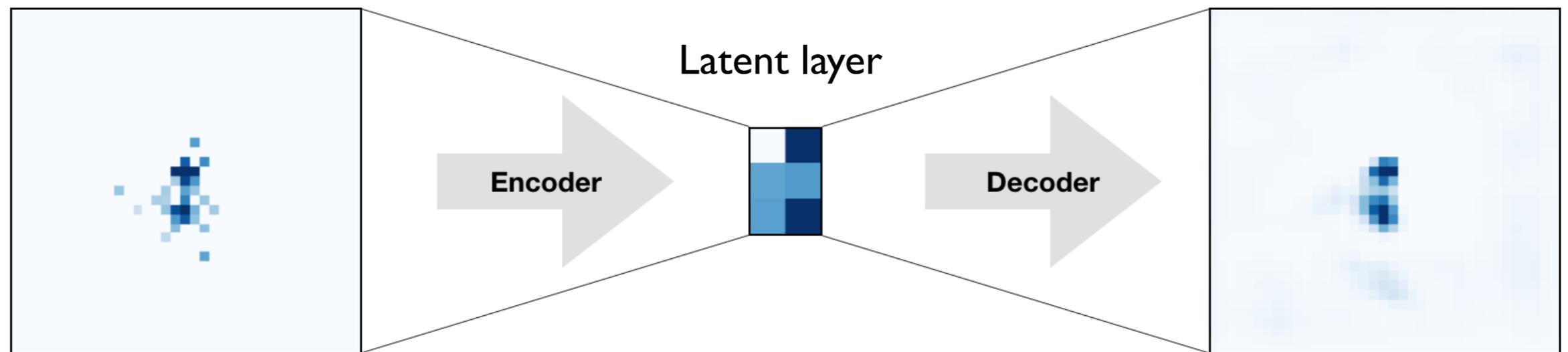
$$m_{J_{1,2}} \rightarrow m_{J_{1,2}} + c m_{JJ}$$



**ANODE is robust while CWoLa completely fails!**

# Searching for NP with deep autoencoders

Heimel et al I808.08979; Farina, Nakai & DS I808.08992



An autoencoder maps an input into a reduced “latent representation” and then attempts to reconstruct the original input from it.

*Can use reconstruction error as an anomaly threshold!*

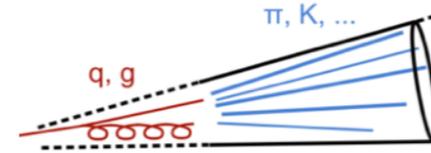
See also:

Hajer et al “Novelty Detection Meets Collider Physics” I807.10261

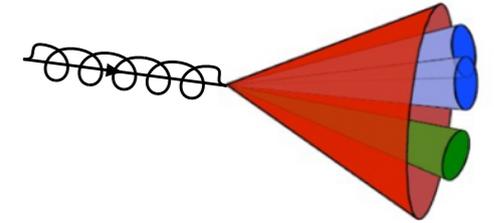
Cerri et al “Variational Autoencoders for New Physics Mining at the Large Hadron Collider” I811.10276

# Performance

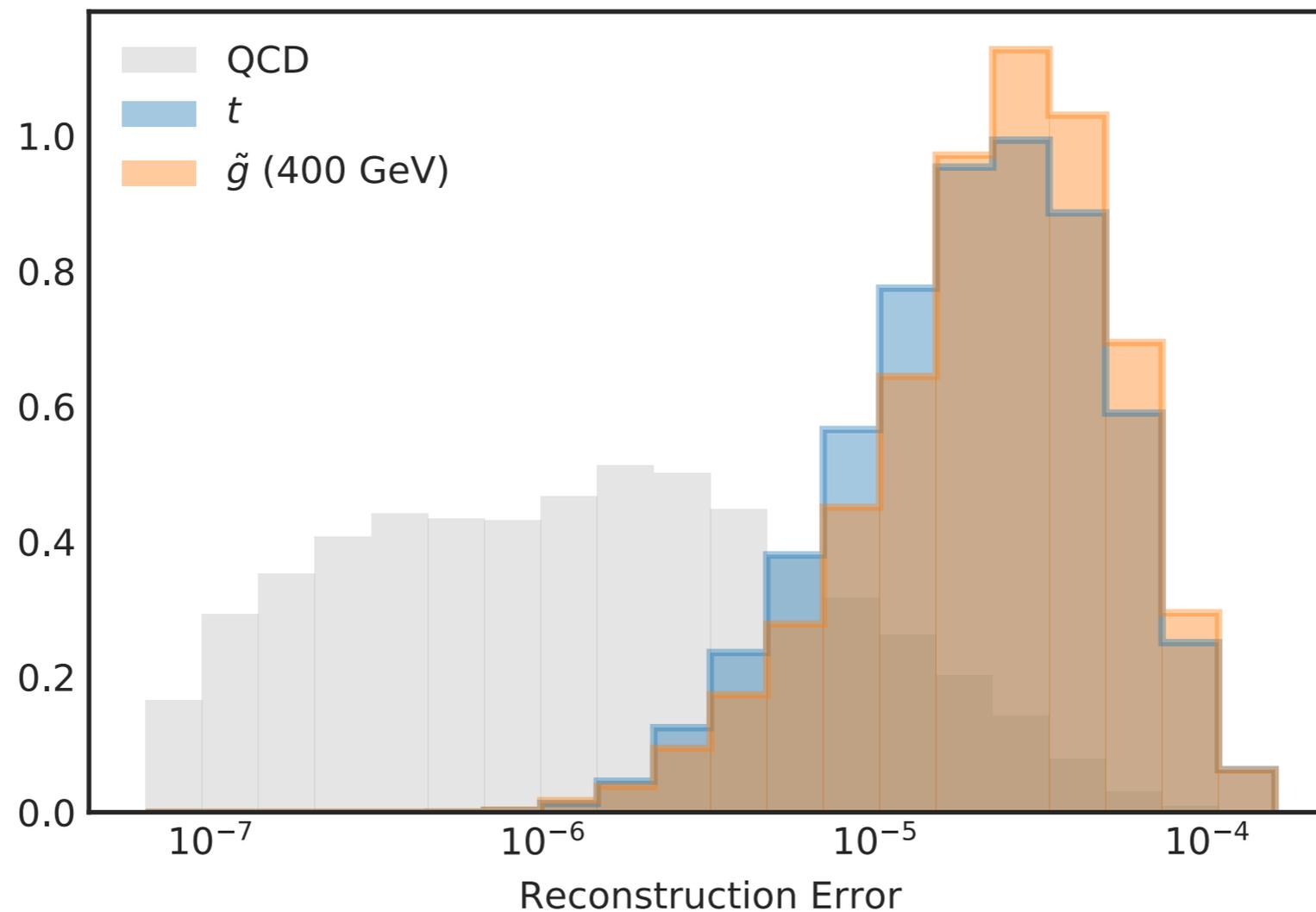
background (QCD)



signal (tops and 400 GeV RPV gluinos)



It works as an anomaly detector!



Robust against contamination with signal — can use in fully unsupervised mode