# Searching for the Unexpected with Unsupervised Machine Learning

#### David Shih

#### Bay Area Particle Theory Seminar

#### October 25, 2019







### Current Status of NP Searches @ LHC

#### ATLAS SUSY Searches\* - 95% CL Lower Limits

July 2019  $\sqrt{s} = 13 \text{ TeV}$ Signature [L dt [fb<sup>-1</sup>] Model Mass limit Reference  $\tilde{q}\tilde{q}, \tilde{q} \rightarrow q \ell_{\perp}^{0}$ 1.55  $0 \in \mu$ 38.1 m(2<sup>1</sup>)<100 GeV 1712.02332 2-6 jets  $E^{hiss}$ 1-3 jets (fix, Sx Deger mono-jet 38.1 0.43 0.71 n)@m(?])=5GeV 1711.00301 Inclusive Searches  $0 \ r, \mu$ 2-6 jets  $E_{\pi}^{miss}$ 36.1 m(2<sup>1</sup>)<200 GeV 1712.02332  $\tilde{z}\tilde{z}, \tilde{g} \rightarrow a\tilde{g}\tilde{\ell}_1^0$ 2.0 0.95-1.6 Forbidden m(2))-200 GeV 1712.02332  $3 \epsilon_{e} \mu$ 4 jets m(21)<\$00 GeV 1708.03731  $\tilde{g}\tilde{g}, \tilde{g} \rightarrow q\tilde{q}(U)\tilde{k}_{1}^{0}$ 38.1 1.85  $E_{2}^{mis}$  $ee, \mu\mu$ 2 jets 38.1 1.2 m(z) m(i<sup>2</sup>)=60 GeV 1805.11361  $\mathcal{E}_{\mathcal{T}}^{miss}$ 7-11 jets. ξξ.ξ→ogWZℓ<sub>1</sub><sup>0</sup>  $0 \times \mu$  $m(\ell_1^4) <400 \, GeV$ 1708.02794 36.1 1.8  $SS\sigma,\mu$ 1.15  $m(\tilde{c})$ - $m(\tilde{c}_1)$ =200 GeV ATLAS-CONF-2019-015 6 jets 139  $\tilde{g}\tilde{g}, \tilde{g} \rightarrow \tilde{u}\tilde{\ell}_{1}^{0}$ 0-1 e.p. 3.6 Ento 79.8 2.25  $m(\tilde{t}_{1}^{0}) < 200 \text{ GeV}$ ATLAS-CONF-2018-041 88 c.u 6 jets 139 1.25  $m(y) \cdot m(\tilde{x}_1^2) = 300 \text{ GeV}$ ATLAS-CONF-2019-015 Multiple 1708.09366.1711.03301  $b_1b_1, b_1 \rightarrow b\tilde{x}_1^0/\tilde{x}_1^2$ 35.1 Forbidden 0.9 n(i?)=305GeV, 8B(M?)=1 0.58-0.82 Multiple Etylpidea 1708.00256 35.1m(f<sup>\*</sup>)=300 GeV, BR(l/<sup>\*</sup>)=BR(f<sup>\*</sup>)=0.5 ATLAS-CONF-2019-015 Multiple 139 Forbiddea 0.74 m(\$\vec{k}\_{1}^{\*})=200 GeV, m(\$\vec{k}\_{1}^{\*})=300 GeV, BF((\$\vec{k}\_{1}^{\*})=1)  $b_1b_1, b_1 \rightarrow b \hat{\mathcal{C}}_2^0 \rightarrow b \hat{\mathcal{K}}_1^0$  $0 \epsilon, \mu$  $E_{\gamma}^{\min}$ 139 0.23-1.35  $\Delta m(\hat{c}_{2}^{0},\hat{c}_{1}^{0}) = 130 \text{ GeV}, m(\hat{c}_{1}^{0}) = 100 \text{ GeV}.$ SUS7-2018-31 8 b Forbidden 0.23-0.4E  $\Delta m(\tilde{t}_1^2, \tilde{t}_1^2) = 130 \text{ GeV}, m(\tilde{t}_1^0) = 0 \text{ GeV}$ SUSY-2018-01 0-2 c, μ = 0-2 jgta/1-2 h E<sup>gèts</sup> 1506.08516, 1709.04183, 1711.11520  $\tilde{i}_1 \tilde{i}_1, \tilde{i}_1 \rightarrow W b \tilde{k}_1^0 \text{ or } b \tilde{k}_1^0$ 36.11.0  $m(\ell_1^n)=1$  GeV 3 jets/1 à Ettata 0.44-0.59  $\tilde{t}_1 \tilde{t}_1, \tilde{t}_1 \rightarrow Wb \tilde{t}_1^0$  $1 \times \mu$ 139  $m(\tilde{x}_1^{\ell})=400$  GeV ATLAB-CONF-S019-017  $E_{2}^{miss}$  $\tilde{t}_1 \tilde{t}_1, \tilde{t}_1 \rightarrow \tilde{\tau}_1 hw, \tilde{\tau}_1 \rightarrow \tau \tilde{G}$ m(2))=800 GeV  $1 \tau + 1 e \mu \tau$ 2 lets/1 8 38.1 1808.10178 1.16  $i_1i_1, i_1 \rightarrow c i_1^2 / i t, t \rightarrow c i_1^2$  $E_7^{tto}$ 36.10.85 m(2<sup>0</sup>)=0 GeV 1005/01649  $0 n \mu$ 2 c. 0.46  $m(\tilde{r}_1, \tilde{c}) \cdot m(\tilde{c}_1^2) = 50 \text{ GeV}$ 1805.01649 0.43 1711.05301 Eve 38.1  $0 \in \mu$ mono-jet  $m(l_1, l_2) \cdot m(l_1^{(0)}) = 5 \text{ GeV}$  $\tilde{l}_{2}\tilde{k}_{2}, \tilde{k}_{2} \rightarrow \tilde{k}_{1} + h$  $1-2.e.\mu$ 4.5  $E_{i}^{miss}$ 38.1 0.32-0.88  $m(\hat{t}_{1}^{0})=0$  GeV,  $m(\hat{t}_{1})+m(\hat{t}_{1}^{0})=-180$  GeV. 1708.03958  $i_2i_2,\,i_2{\rightarrow}i_1+Z$  $E_T^{mbs}$  $3 v, \mu$ 139 0.85  $m(\tilde{x}_{1}^{*}) = 380 \text{ GeV}, m(\tilde{t}_{1}) \cdot m(\tilde{x}_{1}^{*}) = 40 \text{ GeV}$ ATLAS-CONF-2019-015 1.5Farbioden ř, 后面 新潟 via WZ  $2-3 e.\mu$  $E_{\nu}^{miss}$ 0.6 1403.5294, 1805.02293 38.1  $m(\mathcal{E}_{i}^{2})=0$  $E^{iiiss}$  $ev, \mu\mu$  $\geq 1$ 139 0.205  $m(\tilde{t}_{1}^{*}) m(\tilde{t}_{1}^{0}) = 5 \text{ GeV}$ ATLAS-CONF-S019-014  $E_{\pi}^{min}$ RTRT via WW  $2r,\mu$ 139 ATLAS CONF-2019-008 e 0.42  $m(\ell_1^2)=0$ 岩塔 via Wa  $0-1.e.\mu$  $2.0/2\gamma$ Etto State. 0.74 ATLAS-CONF-2019-019, ATLAS-CONF-2019-XYZ 139 Forbiddon  $H(\tilde{t}_{1}^{0})=70$  GeV 结结 via &/P  $2\epsilon,\mu$ grites 1392 1.0  $m(\tilde{t}_{1}^{2}) = 0.5(m(\tilde{t}_{1}^{*}))m(\tilde{t}_{1}^{0}))$ ATLAS-CONF-2019-006  $E^{mba}$  $\bar{\tau}\bar{\tau}, \bar{\tau} \rightarrow \tau \bar{\chi}_{1}^{0}$  $2\tau$ 139 PL PREF 0.16-0.3 0.12-0.39 ATLAS-CONF-2019-018  $m(k_1^3)=0$  $\mathcal{E}_{\mathcal{D}}^{miss}$  $2r, \mu$ 0 jets 139 0.7 ATLAS CONF-2019-008  $\tilde{t}_{L,R}\tilde{t}_{L,K}, \tilde{t} \rightarrow t\tilde{t}_{1}^{0}$  $m(\mathcal{X}_1^2)=0$  $E_{\gamma}^{\rm fries}$ 139 0.258 ATLAS-CONF-2019-014  $2 \epsilon \mu$  $\geq 1$ m(2)-m(2)-10 GeV  $BB, \theta \rightarrow kG/ZG$  $0 e, \mu$  $E_{\gamma}^{mba}$ 0.13-0.23 0.29-0.88 1806.04030  $\geq 3.6$ 35.1 ù  $BR(\hat{\ell}_1^2 \rightarrow \Lambda \hat{\ell}_2)=1$ E takes  $4 c, \mu$ 0 jets 38.1 0.3  $BR(\mathcal{C}_{1}^{0} \rightarrow Z\mathcal{C})-1$ 1804.00602 Ħ Direct  $\mathcal{X}_1^{\dagger} \mathcal{X}_1^{\dagger}$  prod., long-lived  $\mathcal{X}_1^{\dagger}$ Disapp. trk. 1 jet  $E_{2}^{mbs}$ 38.1 0.46 Pure Wino 1712.02118 0.15 Pure Higgsino ATL-PHVS-PUB-2017-019 Stable # R-hadron Multiple 1902.01636,1803.04095 36.1 2.0 Multiple 2.05 2.4 1710.04901,1808.04095 Metastable ∦ R-hadron, ≵→ooℓ<sub>1</sub><sup>0</sup>  $\hat{s} = [r(\hat{s}) = 10 \text{ ns}, 0.2 \text{ ns}]$ 36.1 $m(\tilde{E}_1^{i})=100$  GeV LFV  $pp \rightarrow \bar{\pi}_r + X_r \bar{\pi}_r \rightarrow q\mu/e\tau/\mu\tau$  $\lambda_{aaa}^{r} = 0.11, \lambda_{1021100,100} = 0.07$ 10.07.07 3.2 1.9 1607.06079  $\hat{\chi}_1^{\dagger} \hat{\chi}_1^{\dagger} / \hat{\chi}_2^0 \rightarrow WW/ZUUUvv$  $m(\tilde{E}_1^0)=100$  GeV  $4 \epsilon, \mu$ 0 jets E's 36.10.82 1.33 1604.00602  $[A_{33}\neq 0, A_{32}\neq 0]$ Large  $\mathcal{X}_{1,0}^{\prime}$  $\tilde{g}\tilde{g}, \tilde{g} \rightarrow gq\tilde{\chi}_{1}^{0}, \tilde{\chi}_{1}^{0} \rightarrow qqq$ 4-5 large-R jets 36.11.3 1.9 1804.03558 n(7, |=200 GeV, 1 (7, =26-1, 26-5) 100 BeW Multiple 1.05 36.12.0  $\mathfrak{m}(\widetilde{\mathfrak{t}}_1^0)$ =200 GeV, bino-like ATLAS-CONF-2018-008 **PPV**  $\widetilde{R}, \widetilde{t} \rightarrow \mathscr{R}_1^0, \widetilde{\mathcal{R}}_1^0 \rightarrow \operatorname{rbs}$ Multiple 35.1 -2a-4. 1a-2 1.05 m(x<sup>0</sup>)i⊨200 GeV, bino-like 0.55ATLAS-CONF-2018-008  $I_1I_1, I_1 \rightarrow bx$ 2 jels + 2 *b* 38.7 0.61 1710.07171 0.42 $i_1i_1, i_1 {\rightarrow} q\ell$  $2r, \mu$ 2.b36.1 $BH(\tilde{t}_1 \rightarrow bv/b\mu) > 20\%$ 1710.05544 0.4 - 1.45e-10< X ... ≤1e-8, 3e-10< X ... ≤3e-DV 136 1.0 1.6 BR(i<sub>1</sub>→<sub>6</sub>s)=100%, cccs<sub>2</sub>=1 ATLAS-CONF-2019-006  $1_K$ 

"Only a selection of the available mass limits on new states or phonomenous is phonen. Many of the limits are based as 10<sup>-1</sup>

Mass scale [TeV]

1

ATLAS Preliminary

### Current Status of NP Searches @ LHC



"Only a selection (

Selection of observed limits at 95% G.L. (theory uncertainties are not included). Probe up to the quoted mass limit for light LSPs unless stated otherwise. The quantities  $\Delta M$  and x represent the absolute mass difference between the primary sparticle and the LSP, and the difference between the intermediate

### Current Status of NP Searches @ LHC



#### Vector-like quark single production





0.2

IEP 2019

0.0

Selection of observed limits at 95% C.L. The quantities  $\Delta M$  and  $\pi$  represent the i-sparticle and the LSP relative to  $\Delta M$ , re-

BB ((1±,1±1± HH→(/ ± , / ± / ±

#### Vector-like quark single production













How can we make sure we don't miss anything?

# Model Independent Searches

Can we search for new physics without a model in mind?

# Model Independent Searches

Can we search for new physics without a model in mind?

Yes?? A brief history of model independent searches in HEP:

•	DO	"Sleuth"	PRD 62:092004 (2000) PRD 64:012004 (2001) PRL 86:3712 (2001)
	HI (Hera)	"General Search"	PLB 602:14-30 (2004) 0705.3721
•	CDF	"Sleuth/Vista"	0712.1311 PRD 78:012002 (2008) 0712.2534 (submitted to PRL, NEVER PUBLISHED) 0809.3781 PRD 79:011101 (2009)
•	CMS	"MUSIC"	CMS-PAS-EXO-14-016
•	ATLAS	"Model independent general search"	1807.07447 EPJC 79:120 (2019)

## Model independent searches

The general approach behind all of these:

- Bin the data into exclusive final states,
- Compare a zillion ID histograms between data and <u>simulation</u> (or just look at high pT tails),
- Look for discrepancies.

#### From CDF 0712.2534:

A global comparison of data to standard model prediction is made in 16,486 kinematic distributions in 344 populated exclusive final states. In each final state, the

bottom quark (b), and missing momentum  $(\not p)$ . Monte Carlo event generators are used to determine the standard model prediction. VISTA partitions data and Monte

The general approach behind all of these:

- Bin the data into exclusive final states,
- Compare a zillion ID histograms between data and <u>simulation</u> (or just look at high pT tails),
- Look for discrepancies.

#### **Enormous look elsewhere effect**

The general approach behind all of these:

- Bin the data into exclusive final states,
- Compare a zillion ID histograms between data and <u>simulation</u> (or just look at high pT tails),
- Look for discrepancies.

#### **Enormous look elsewhere effect**

The general approach behind all of these:

- Bin the data into exclusive final states,
- Compare a zillion ID histograms between data and <u>simulation</u> (or just look at high pT tails),
- Look for discrepancies.

Sub-optimal signal / background discrimination

#### **Enormous look elsewhere effect**

The general approach behind all of these:

- Bin the data into exclusive final states,
- Compare a zillion ID histograms between data and <u>simulation</u> (or just look at high pT tails),
- Look for discrepancies.

Sub-optimal signal / background discrimination

#### **Over-reliance on simulation for background prediction**

### An example of what is found

#### Sleuth@CDFII result

#### fraction of pseudo experiments in this fraction of pseudo experiments in any (top 5) final state as interesting as CDF data final state as interesting as CDF data **SLEUTH Final State** $\mathcal{P}$ 0.46bb0.0055jø 0.0092expected to be as interesting 0.011 pii Sleuth finds no significant excess in CDF Run II high-p<sub>T</sub> data $\ell^+\ell'^+ p$ 0.016This does not prove there is no 0.016 $\tau p$

**U.UID** new physics present

From B. Knuteson talk at UMich (2008)

#### CDF Run II (927 pb<sup>-1</sup>)

### An example of what is found



From B. Knuteson talk at UMich (2008)



#### Simple way to mitigate LEE:

• Divide data in two



- Divide data in two
- Look for discrepancies in first half



- Divide data in two
- Look for discrepancies in first half
- Fix regions of interest around these discrepancies and restrict search to these in the second half





- Divide data in two
- Look for discrepancies in first half
- Fix regions of interest around these discrepancies and restrict search to these in the second half





- Divide data in two
- Look for discrepancies in first half
- Fix regions of interest around these discrepancies and restrict search to these in the second half
- An even better approach: k-fold cross validation, see e.g. 1805.02664



## Sub-optimal signal/background discrimination

Can do much better than comparing ID histograms.

Using **deep learning**, can leverage the entire high-dimensional phase space to separate signals from background!



See D'Agnolo & Wulzer 1806.02350: train DNN to distinguish data from bg MC. Optimal version of existing general searches

#### Over-reliance on simulation for backgrounds

Existing approaches are largely signal model independent, but not background model independent.

To achieve full model independence, we would like to predict backgrounds directly from the data.



We want to search for new physics in the data in a signal **and** background independent way.

Cannot use simulation or ground-truth labels. We need **unsupervised ML**.

Two main approaches proposed so far:

We want to search for new physics in the data in a signal **and** background independent way.

Cannot use simulation or ground-truth labels. We need **unsupervised ML**.

Two main approaches proposed so far:

 CWoLa hunting Collins, Howe & Nachman 1805.02664, 1902.02634

We want to search for new physics in the data in a signal **and** background independent way.

Cannot use simulation or ground-truth labels. We need **unsupervised ML**.

Two main approaches proposed so far:

- CWoLa hunting Collins, Howe & Nachman 1805.02664, 1902.02634
- Autoencoders
  Farina, Nakai & DS 1808.08992; Heimel et al 1808.085



We want to search for new physics in the data in a signal **and** background independent way.

Cannot use simulation or ground-truth labels. We need **unsupervised ML.** 

Two main approaches proposed so far:

- CWoLa hunting Collins, Howe & Nachman 1805.02664, 1902.02634
- Autoencoders Farina, Nakai & DS 1808.08992; Heimel et al 1808.089



There are surely more methods waiting to be invented!

Heimel et al 1808.08979; Farina, Nakai & DS 1808.08992



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

The encoding is lossy, so the reconstruction is not perfect.

Many real world applications of autoencoders, including anomaly detection, fraud detection, denoising, compression, generation, density estimation

See also:

Hajer et al "Novelty Detection Meets Collider Physics" 1807.10261 Cerri et al "Variational Autoencoders for New Physics Mining at the Large Hadron Collider" 1811.10276

Heimel et al 1808.08979; Farina, Nakai & DS 1808.08992

Loss function for autoencoder: *I* 

$$L = \frac{1}{N} \sum_{i=1}^{N} (x_i^{in} - x_i^{out})^2$$

Heimel et al 1808.08979; Farina, Nakai & DS 1808.08992

Loss function for autoencoder:  $L = \frac{1}{N} \sum_{i=1}^{N} (x_i^{in} - x_i^{out})^2$ 



Heimel et al 1808.08979; Farina, Nakai & DS 1808.08992

Loss function for autoencoder:  $L = \frac{1}{N} \sum_{i=1}^{N} (x_i^{in} - x_i^{out})^2$ 



Heimel et al 1808.08979; Farina, Nakai & DS 1808.08992





Heimel et al 1808.08979; Farina, Nakai & DS 1808.08992





# Sample definitions

 Background: QCD jets (p<sub>T</sub>: 800-900 GeV, |η|<1, anti-kt R=1)</li>



- Signals:
  - All-hadronic tops
  - 400 GeV gluinos decaying via RPV

 We formed jet images in η and φ with a pixel resolution and intensity given by the calorimeter towers.





#### **Convolutional Autoencoder** Autoencoder architecture



128C3-MP2-128C3-MP2-128C3-32N-6N-32N-12800N-128C3-US2-128C3-US2-1C3 128C3-MP2-128C3-MP2-128C3-32N-6N-32N-12800N-128C3-US2-128C3-US2-1C3

Our primary autoencoder used convolutional neural networks (CNNs) for encoding and decoding the jet images.

We also considered autoencoders based on PCA and simple DNNs.

Many more architectures are possible.

- d too large  $\rightarrow$  autoencoder becomes identity transform
- d too small  $\rightarrow$  autoencoder cannot learn all the features

- d too large  $\rightarrow$  autoencoder becomes identity transform
- d too small  $\rightarrow$  autoencoder cannot learn all the features



- d too large  $\rightarrow$  autoencoder becomes identity transform
- d too small  $\rightarrow$  autoencoder cannot learn all the features





- d too large  $\rightarrow$  autoencoder becomes identity transform
- d too small  $\rightarrow$  autoencoder cannot learn all the features





### Performance: weakly supervised mode

Train the AE on QCD backgrounds only.



# Performance: weakly supervised mode



### Fully unsupervised mode

Can also train on QCD background "contaminated" with a small fraction of signal. This could be representative of actual data.



Performance of AE robust even up to 10% contamination!

### **Background estimation**

Finally, background estimation. Remember: unlike past attempts at model-independent searches, we want it to be data-driven!

One idea: combine autoencoder with a bump hunt in jet mass. Estimate backgrounds using sidebands in mass.

Only works if cutting on reconstruction error does not sculpt the mass distribution of the background!



#### Bump hunt with autoencoder



We find empirically that the background jet mass distribution is fairly stable against cuts on CNN AE reconstruction loss above ~250 GeV.

#### Bump hunt with autoencoder



Train directly on data that contains 400 GeV gluinos. Use the AE to clean away QCD jets. Enhance the significance of the bump hunt! (improve S/B by factor of ~6)

#### Could really discover new physics this way!

#### Autoencoder with explicit decorrelation

A more controlled approach to mass decorrelation would be to explicitly penalize correlations in the training of the autoencoder.

One promising method: autoencoder with adversarial decorrelation (Heimel et al 1808.08979; based on 1611.01046, 1703.03507 and the idea behind GANs)

- Introduce a second NN, the adversary, that tries to predict the mass from the reconstruction loss.
- Penalize the total loss function when the adversary does well.

$$L_{adv} = \sum_{i} (f_{adv}(L_{AE}(x_i)) - m_i)^2$$

$$L_{tot} = L_{AE} - \lambda L_{adv}$$

#### Autoencoder with explicit decorrelation



Current implementation needs QCD backgrounds for decorrelation. Can generalize to fully unsupervised case?

#### Alternatives to adversaries

Adversaries are notoriously tricky to train — saddle point optimization

$$\min_{\theta_{\rm clf}} \max_{\theta_{\rm adv}} L_{\rm clf}(y(\theta_{\rm clf})) - \lambda L_{\rm adv}(y(\theta_{\rm clf}), m; \theta_{\rm adv})$$

Would be great if we could achieve the same performance but with a convex regularizer term

$$\min_{\theta_{\rm clf}} L_{\rm clf}(y(\theta_{\rm clf})) + \lambda C_{\rm reg}(y(\theta_{\rm clf}), m)$$

First idea: can we just use Pearson correlation coefficient?

$$C_{\rm reg} = R(y,m) \propto \sum_i y_i m_i$$

Problem: this only measures linear correlations

#### Pearson correlation



y and m can be highly correlated yet R=0

# Distance correlation ("DisCo")

Work in progress with Gregor Kasieczka

Promising idea: "distance correlation" (Szekely, Rizzo, Bakirov 2007; Szekely & Rizzo 2009)

 $dCov^{2}(X,Y) = \langle |X - X'| |Y - Y'| \rangle + \langle |X - X'| \rangle \langle |Y - Y'| \rangle - 2\langle |X - X'| |Y - Y''| \rangle$ 

- Zero iff X,Y are independent; positive otherwise
- Computationally tractable
- Straightforward sample definition doesn't require binning

#### Distance correlation



Disco is sensitive to nonlinear correlations!

# State of the art: ATLAS study of various decorrelation methods in context of boosted W-tagging.



Andreas Søgaard / University of Edinburgh







# Distance correlation ("DisCo")

Work in progress with Gregor Kasieczka



# Summary/Outlook

There is increasingly strong motivation to perform model-independent searches at the LHC. We need to ensure that we have not missed anything in the data.

However, attempts at model-independent searches have suffered from major drawbacks: an enormous trials factor, simplistic signal/background discrimination, and an over-reliance on simulation. Together, they have given the philosophy of model-independent searches a bad name.

In this talk we have explored new ideas for model-independent searches inspired by recent breakthroughs in unsupervised deep learning.

We have seen that:

- Deep autoencoders can find subtle signals in the data in a model agnostic way.
- Decorrelation is essential in combining data-driven background estimation with deep learning.
- A new method "DisCo" can more easily achieve state-of-the-art decorrelation in NN training.

## Outlook

I think there is an exciting future for model-independent searches @ LHC!

- Comparison of different approaches to model-independent searches (AE, CWoLa, ... — work in progress with Pablo Martin & Ben Nachman)
- Careful study of the LEE in these different approaches
- More applications of DisCo
  - e.g. using it to improve the ABCD method ("Double DisCo" work in progress with Gregor Kasieczka, Ben Nachman & Matt Schwartz)
- Real life applications of decorrelation?
  - E.g. making less discriminatory Als for hiring/admissions/bail decisions/sentencing/...?

#### • New ideas for model-independent searches

- e.g. using recent breakthroughs in density estimation to search for anomalies
- many more...?!

#### ML4Jets2020

#### 15-17 January 2020 Europe/Zurich timezone

Search...



Despite an impressive and extensive effort by the LHC collaborations, there is currently no convincing evidence for new particles produced in high-energy collisions. At the same time, there has been a growing interest in machine learning techniques to enhance potential signals using all of the available information.

In the spirit of the first LHC Olympics (circa 2005-2006) [1<sup>st</sup>, 2<sup>nd</sup>, 3<sup>rd</sup>, 4<sup>th</sup>], we are organizing the 2020 LHC Olympics. Our goal is to ensure that the LHC search program is sufficiently well-rounded to capture "all" rare and complex signals. The final state for this olympics will be focused (generic dijet events) but the observable phase space and potential BSM parameter space(s) are large: all hadrons in the event can be used for learning (be it "cuts", supervised machine learning, or unsupervised machine learning).

For setting up, developing, and validating your methods, we provide background events and a benchmark signal model. You can download these from this page. To help get you started, we have also prepared simple python scripts to read in the data and do some basic processing.

The final test will happen 2 weeks before the ML4Jets2020 workshop. We will release a new dataset where the "background" will be similar to but not identical to the one in the development set (as is true in real data!). The goal of the challenge is to see who can "best" identify BSM (yes/no, what mass, what cross-section) in the dataset. There are many ways to quantify "best" and we will use all of the submissions to explore the pros/cons of the various approaches.

# Thanks for your attention!