# Machine Learning For Searches

Michael Kagan

LHCP, May 21, 2019

U.S. DEPARTMENT OF
**ENERGY**
Office of Science

**SLAC** NATIONAL ACCELERATOR LABORATORY

# The Search Plan

```
┌─────────────────┐
│ Experiment Data │
│   Acquisition   │ ──┐
└─────────────────┘   │   ┌──────────────┐   ┌────────────────┐   ┌─────────────────────┐
                      ├─▶ │    Object    │──▶│Event Selection │──▶│ Statistical Analysis│
┌─────────────────┐   │   │ Reconstruction│   └────────────────┘   └─────────────────────┘
│   Simulation    │ ──┘   └──────────────┘
└─────────────────┘
```

Physics knowledge drives development of experiments, simulations, reconstruction, and analysis pipeline
- Underlying theory drives our inference goals
- Mechanistic understanding of structure of events, particles interactions with material
- Compositionality: design detectors and algorithms to identify specific particles, and analyze them together as events
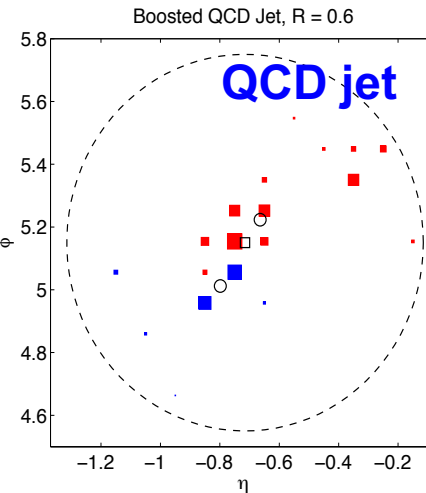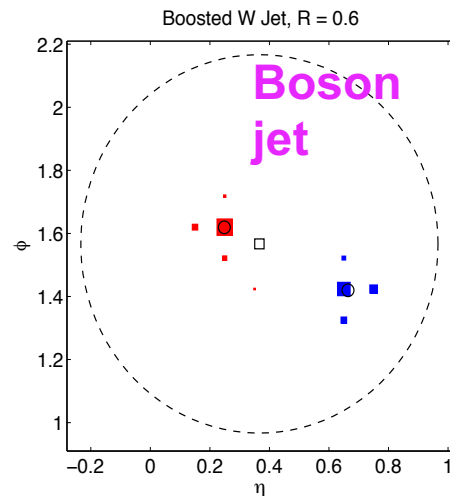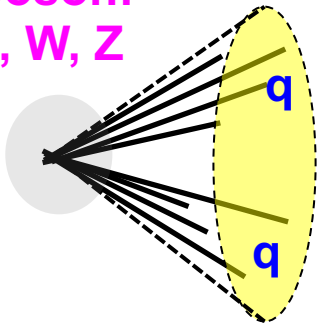
Much of this is intractable
- Don't know p(shower | electron) or p(electron | shower)
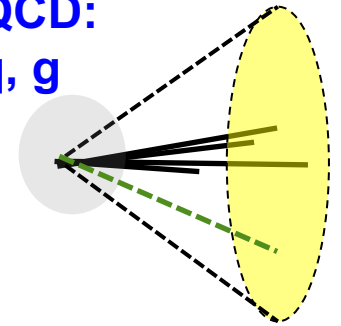- Can sample distributions with simulators encapsulating physics knowledge

Machine learning to augment and improve the pipeline, preserving our physics knowledge while by providing expressive and flexible models to study our data

# Jet Tagging

**Boson: h, W, Z**

q

q

W

Boosted W Jet, R = 0.6

**Boson jet**

g / q

g / q

p

p

Boosted QCD Jet, R = 0.6

**QCD jet**

**QCD: q, g**

Boosted QCD Jet, R = 0.6

g / q

p

- Can use internal jet (sub)structure of a jet for classification

- Wealth of domain expertise in feature engineering

- Can Machine Learning perform this classification?

# Inductive Bias and Data Representation

Moving **inductive bias** from feature engineering to machine learning (neural network) model design

- Inductive bias ~ knowledge about the problem
- Feature engineering ~ hand crafted variables
- Model design ~ the data representation and the structure of the machine learning model / network

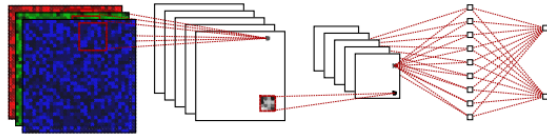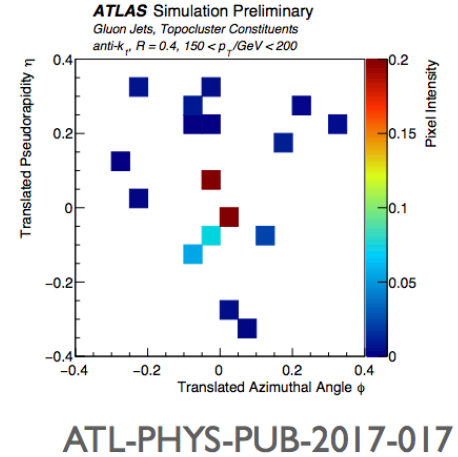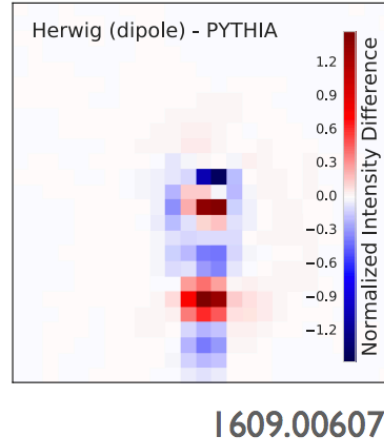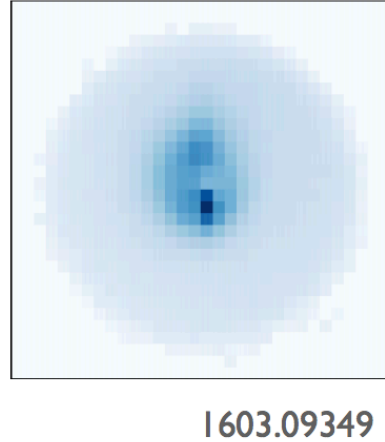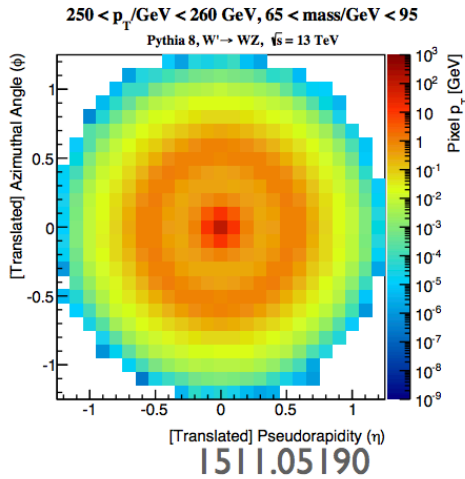Need a good inductive bias, i.e. physical motivation, for data representation and model structure

- Better learn to approximate our data
- Easier to extract information about what is learned?
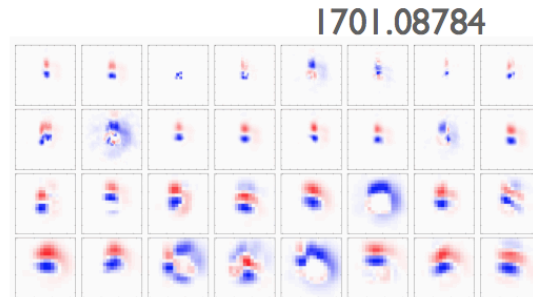
We can represent jets in different ways
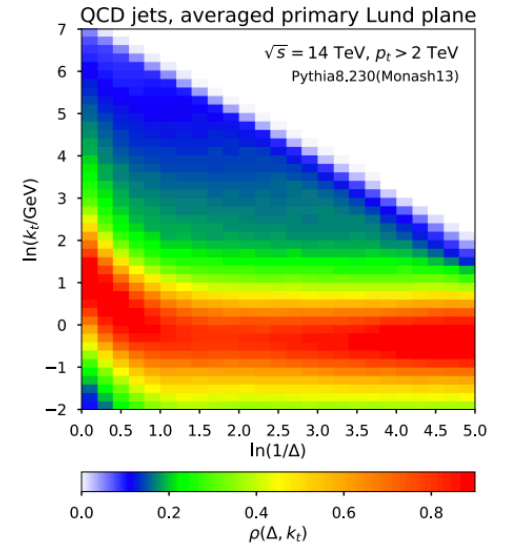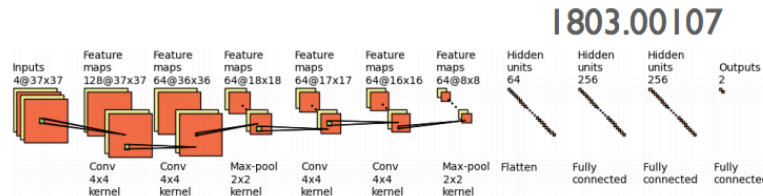We can utilize different classes of models

# Jets as Images

250 < $p_T$/GeV < 260 GeV, 65 < mass/GeV < 95
Pythia 8, W' → WZ, √s = 13 TeV

1511.05190

1603.09349

Herwig (dipole) – PYTHIA

1609.00607

**ATLAS** Simulation Preliminary
Gluon Jets, Topocluster Constituents
anti-$k_t$, R = 0.4, 150 < $p_T$/GeV < 200

ATL-PHYS-PUB-2017-017

red = transverse momenta of charged particles

green = the transverse momenta of neutral particles

blue = charged particle multiplicity

1612.01551
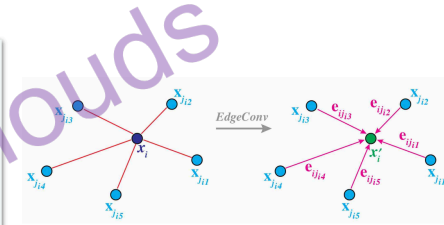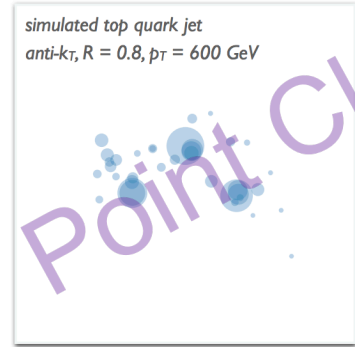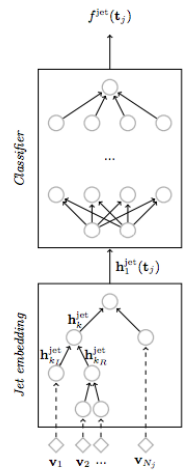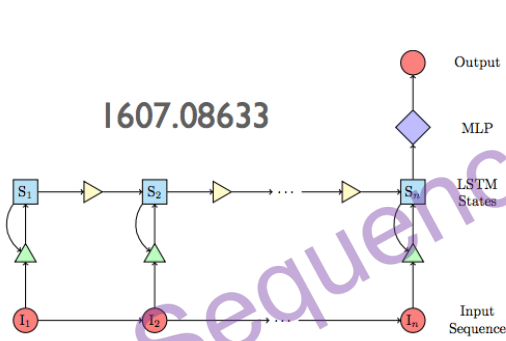
1701.08784

1803.00107

QCD jets, averaged primary Lund plane

√s = 14 TeV, $p_t$ > 2 TeV
Pythia8.230(Monash13)

$\ln(k_t/\text{GeV})$

$\ln(1/\Delta)$

$\rho(\Delta, k_t)$

1807.04758

# Jets as Collections of Particles

Sequences

1607.08633

Trees

1711.09059

1711.02633

1702.00748

Point Clouds

1902.08570

Graphs

NIPS2017 workshop [http://bit.ly/2AkwYRG]

Sets

Lorentz Boost Network

1812.09722

1810.05165

$$k_{\mu,i} \xrightarrow{\text{CoLa}} \tilde{k}_{\mu,j} = k_{\mu,i}\, C_{ij}$$

$$\text{with}\quad C = \begin{pmatrix} 1 & 1 & \cdots & 0 & \chi_1 & \cdots & 0 & C_{1,N+2} & \cdots & C_{1,M} \\ \vdots & & \ddots & & \vdots & & \ddots & \vdots & \ddots & \vdots \\ 1 & 0 & \cdots & 1 & 0 & \cdots & \chi_N & C_{N,N+2} & \cdots & C_{N,M} \end{pmatrix}. \quad \tilde{k}_j \xrightarrow{\text{LoLa}} \hat{k}_j = \begin{pmatrix} m^2(\tilde{k}_j) \\ p_T(\tilde{k}_j) \\ p_T(\tilde{k}_j)\Delta R_{j,\text{jet}} \\ w_{jm}^{(E)} E(\tilde{k}_m) \\ w_{jm}^{(d)} d_{jm}^2 \\ E_T(\tilde{k}_j)E_T(\tilde{k}_m)(\Delta R_{jm})^{0.2} \end{pmatrix},$$

1707.08966, 1812.09223

And more...

# Jet Tagging with ML



arXiv:1902.09914

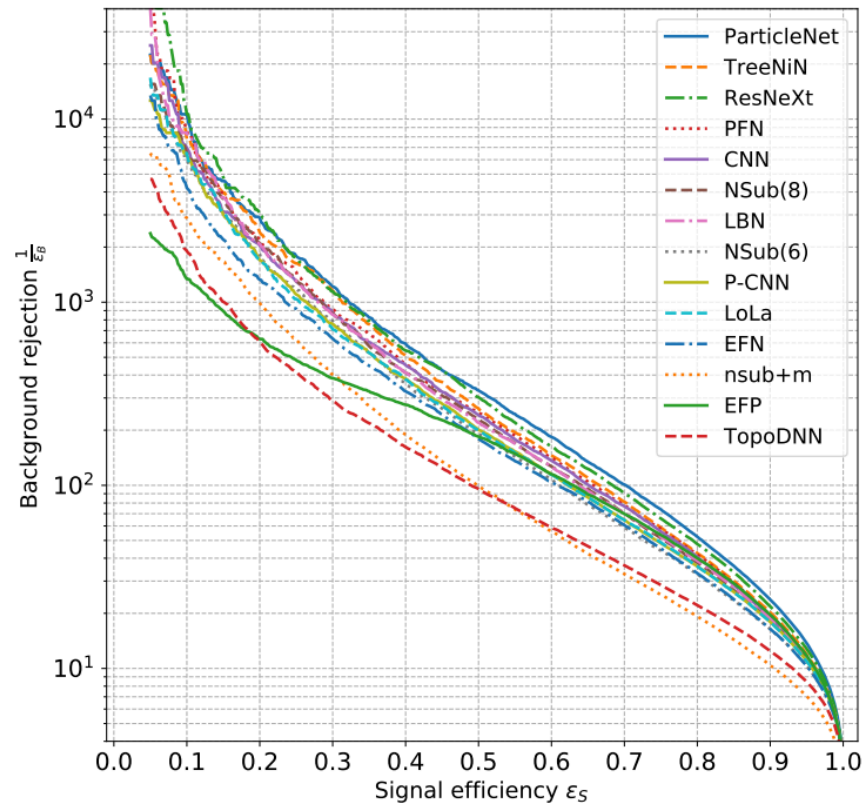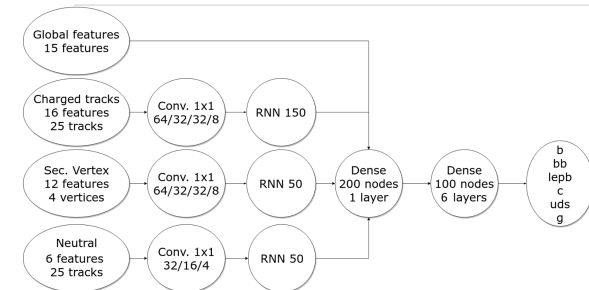| | AUC | Acc | $1/\epsilon_B$ ($\epsilon_S = 0.3$) | | | #Param |
|---|---|---|---|---|---|---|
| | | | single | mean | median | |
| CNN [16] | 0.981 | 0.930 | 914±14 | 995±15 | 975±18 | 610k |
| ResNeXt [30] | 0.984 | 0.936 | 1122±47 | 1270±28 | 1286±31 | 1.46M |
| TopoDNN [18] | 0.972 | 0.916 | 295±5 | 382± 5 | 378 ± 8 | 59k |
| Multi-body $N$-subjettiness 6 [24] | 0.979 | 0.922 | 792±18 | 798±12 | 808±13 | 57k |
| Multi-body $N$-subjettiness 8 [24] | 0.981 | 0.929 | 867±15 | 918±20 | 926±18 | 58k |
| TreeNiN [43] | 0.982 | 0.933 | 1025±11 | 1202±23 | 1188±24 | 34k |
| P-CNN | 0.980 | 0.930 | 732±24 | 845±13 | 834±14 | 348k |
| ParticleNet [47] | 0.985 | 0.938 | 1298±46 | 1412±45 | 1393±41 | 498k |
| LBN [19] | 0.981 | 0.931 | 836±17 | 859±67 | 966±20 | 705k |
| LoLa [22] | 0.980 | 0.929 | 722±17 | 768±11 | 765±11 | 127k |
| Energy Flow Polynomials [21] | 0.980 | 0.932 | 384 | | | 1k |
| Energy Flow Network [23] | 0.979 | 0.927 | 633±31 | 729±13 | 726±11 | 82k |
| Particle Flow Network [23] | 0.982 | 0.932 | 891±18 | 1063±21 | 1052±29 | 82k |
| GoaT | 0.985 | 0.939 | 1368±140 | | 1549±208 | 35k |

Appear to be reach performance asymptote by several models

Key for use in experiments: Understanding computational requirements and sensitivity to systematic uncertainties

# Flavour Tagging: Deep Learning in Experimental Action



Finding jets containing long-lived b-hadrons is key to finding H, Z, Top

- Complex decay topology drives need for powerful algorithms
- (Physics driven) Ordering of set of tracks / vertices to analyze as a sequence
- Sequence based algorithm to account for long range correlations among tracks!

# Enforcing Invariance

With flexibility comes complexity:

- Hard to control how models learn and utilize information
- Potentially unwanted sensitivity to poorly modeled aspects of simulation
- Potentially unwanted sculpting of key physics distributions like mass

*Idea*: Augment training of classifier to enforce invariance to changes in a variable $Z$ (nuisance parameter for systematic uncertainty, kinematic variables, etc.)

Adversarial Approach:

- Build loss that encodes performance of a classifier and and adversary

- Classifier penalized when adversary does well at predicting Z

# Learning to Pivot: Physics Example



Optical tradeoff of performance vs. robustness

Non-Adversarial training

AMS

$\lambda = 0 | Z = 0$
$\lambda = 0$
$\lambda = 1$
$\lambda = 10$
$\lambda = 500$

W-jets vs. QCD-jets
Z = noise level from "pileup"

threshold on $f(X)$

## $\lambda=0, Z=0$

- Standard training with no systematics during training, evaluate systematics after training

## $\lambda=0$

- Training samples include events with systematic variations, but no adversary used

## $\lambda=10$

- Trading accuracy for robustness results in net gain in terms of statistical significance

[AMS = Estimate of statistical significance including systematic uncertainty]

# Decorrelating Variables

Same adversarial setup can decorrelate a classifier from a chosen variable (rather than nuisance parameter) [arXiv:1703.03507]

For example, decorrelate classifier from jet mass, so as not to sculpt jet mass distribution with classifier cut

**W-jets vs. QCD Jets**

# Looking for Signals

 Machine Learning driven reconstruction techniques allow us to improve the identification of known particle signatures in detectors

 Typically combine information from several identified particles to search for signals / perform measurements.

When we know what signal we are looking for

- Can rely on standard MC and data driven techniques

What if we don't know what signal we are looking for?

# ML Enhanced Resonance Finding with CWola Hunting

Want to look for resonance but be as agnostic as possible to features, e.g. if

- We don't have a theory yet to predict it
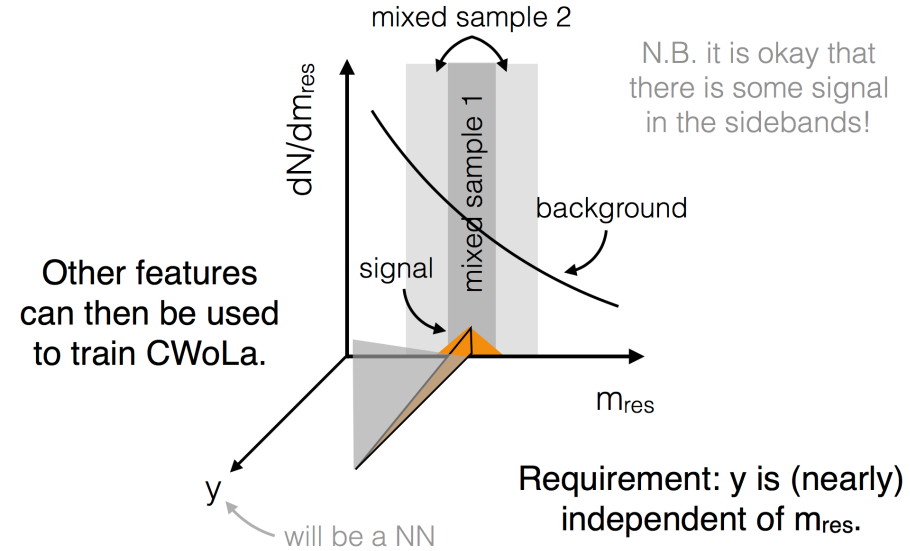- We don't think MC models its features well

Density ratio trick

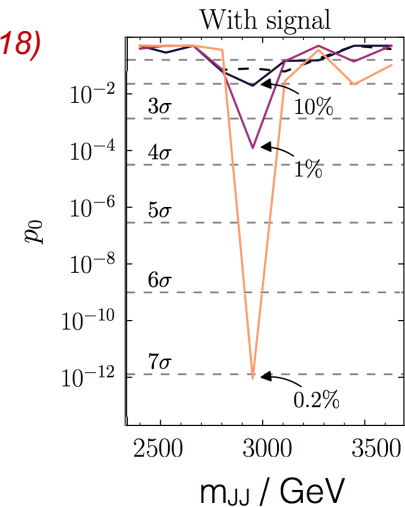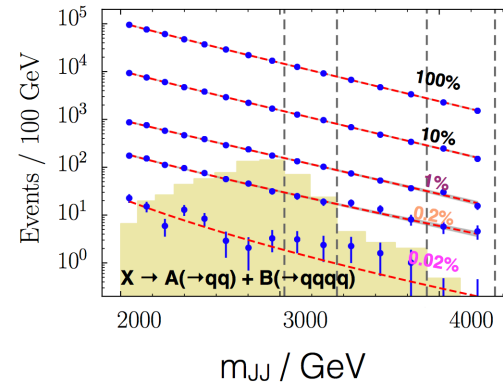$$D(x) = \frac{1}{1 + \frac{p(x|y=0)p(y=0)}{p(x|y=1)p(y=1)}} = \frac{1}{1 + \frac{1}{r(x)}\frac{p(y=0)}{p(y=1)}}$$

- D(x) is discriminator, e.g. ML model
- r(x) is the likelihood ratio

Works even if classes are not pure signal or background, if signal fractions are different!

- Train on data directly in samples with different signal fractions!
- Build mass-independent classifier to not sculpt mass distribution
- Apply varying thresholds on classifier
- Bump Hunt!

mixed sample 2

N.B. it is okay that there is some signal in the sidebands!

$dN/dm_{res}$

mixed sample 1

background

Other features can then be used to train CWoLa.

signal

$m_{res}$

y

will be a NN

Requirement: y is (nearly) independent of $m_{res}$.

*JHEP 10 (2017) 174*
*Phys. Rev. Lett. 121, 241803 (2018)*

Events / 100 GeV

100%

10%

1%

0.2%

0.02%

X → A(→qq) + B(→qqqq)

$m_{JJ}$ / GeV

With signal

$3\sigma$    10%

$4\sigma$    1%

$5\sigma$

$6\sigma$

$7\sigma$    0.2%

$p_0$

$m_{JJ}$ / GeV

- - - - no cut on NN      —— most 1% signal-region-like

—— most 10% signal-region-like      —— most 0.2% signal-region-like

Images from B. Nachman

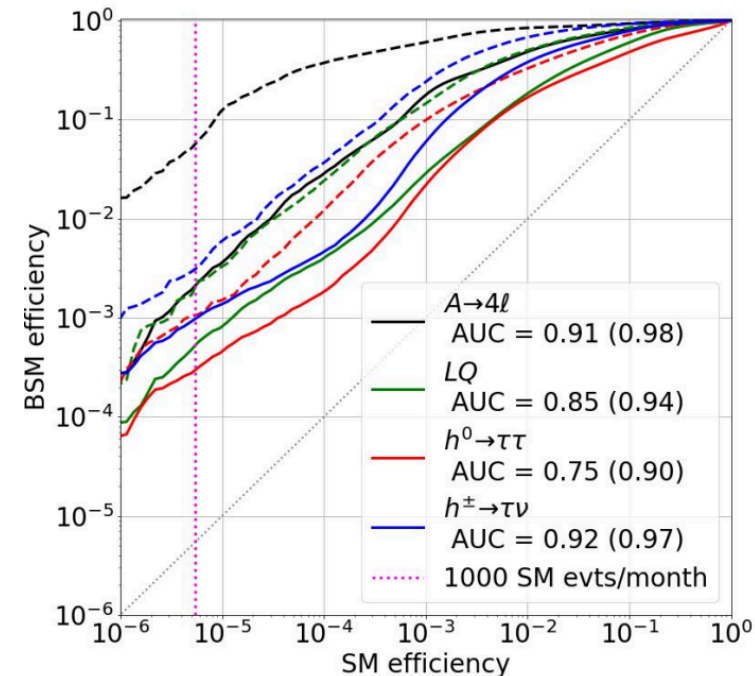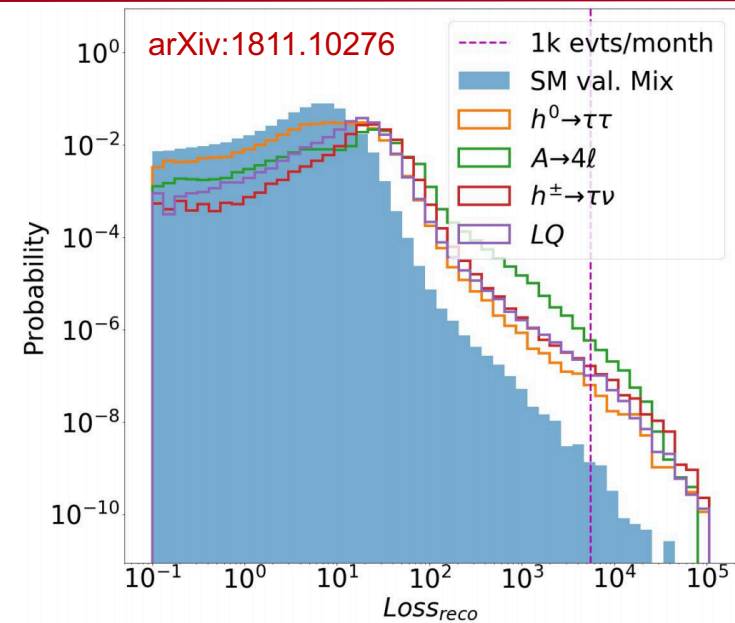# Don't know what to look for? Anomaly Detection Autoencoders

Anomaly detection: find rare events that differ from standard or majority data

- Define standard: i.e. Standard Model
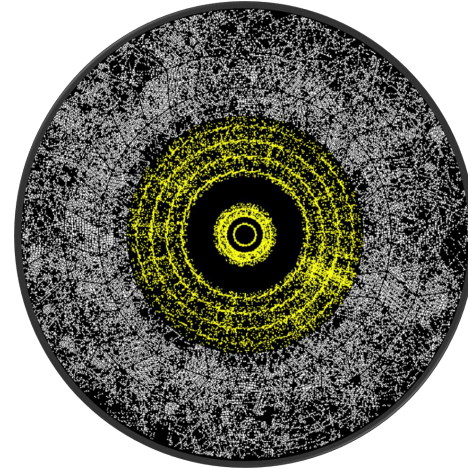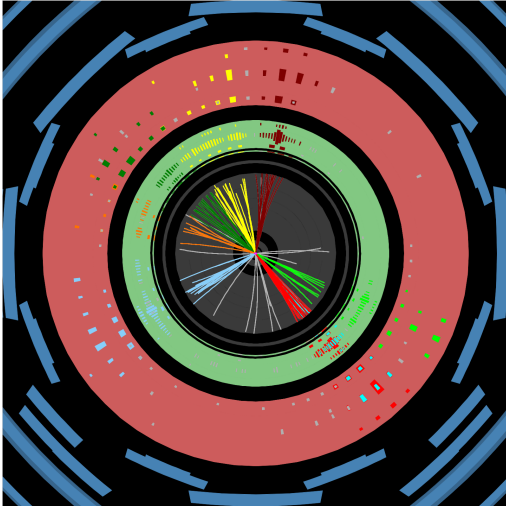- Anomaly: BSM events not like the SM

Look for BSM events with small SM probability, $p_{SM}(x)$ … but don't know $p_{SM}(X)$!

(Variational) AutoEncoder

- Latent variable model, latent space z
- Learn encoder $p(z|x)$ and Decoder $p(x|z)$
- ***Key Idea***: Ability to reconstruct input after encoding into latent space should be diminished for non-standard (i.e. BSM) data

- Growing literature: 1808.08979, 1808.08992, 1807.10261, 1811.10276



arXiv:1811.10276

Legend:
- 1k evts/month
- SM val. Mix
- $h^0 \to \tau\tau$
- $A \to 4\ell$
- $h^\pm \to \tau\nu$
- LQ

x-axis: $Loss_{reco}$, y-axis: Probability



Legend:
- $A \to 4\ell$   AUC = 0.91 (0.98)
- LQ   AUC = 0.85 (0.94)
- $h^0 \to \tau\tau$   AUC = 0.75 (0.90)
- $h^\pm \to \tau\nu$   AUC = 0.92 (0.97)
- 1000 SM evts/month

x-axis: SM efficiency, y-axis: BSM efficiency

# Resource Constraints



Increased pileup at HL-LHC will push boundaries of our computational capabilities
- Major challenges in triggering, large scale simulation, and high multiplicity tracking
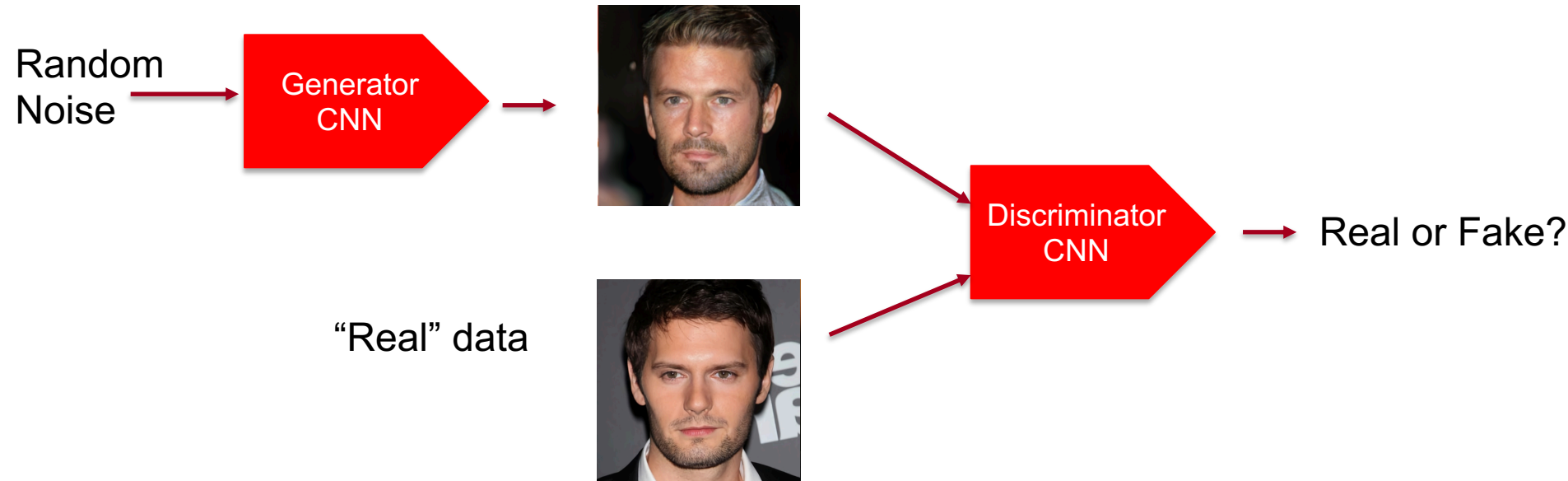- New tools and developments in ML may help address some of these challenges

Simulation
- Accurate but often costly simulation of particle interactions with material, that produces sample and not analytic P(energy deposits | particle)
- *ML approach*: Generative models to learn data distribution, p(x), and produce samples?

Trigger
- High performance algorithms early in trigger to reduce backgrounds for key signals?
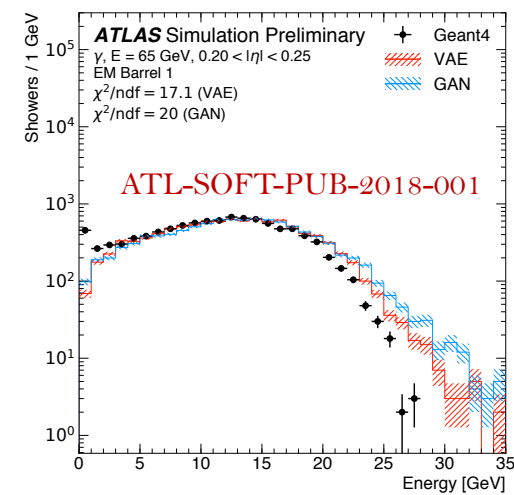
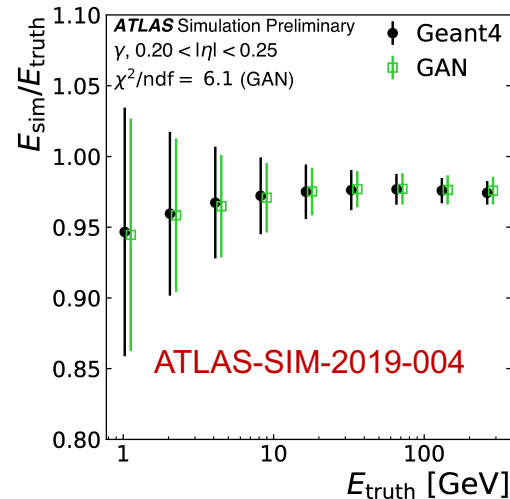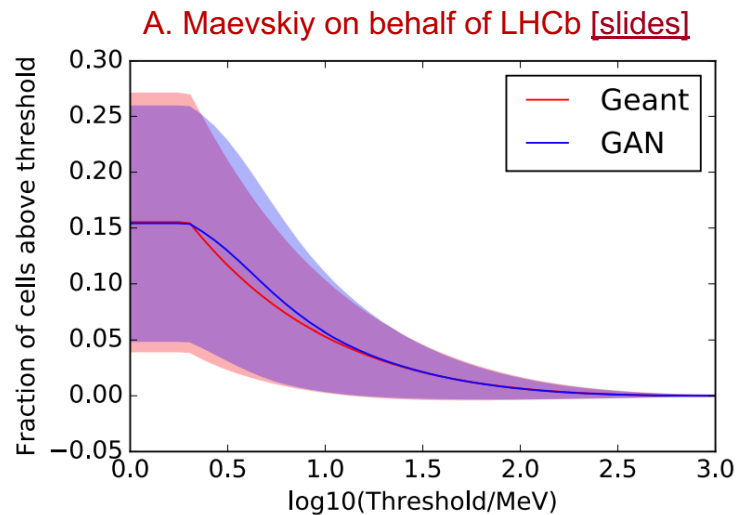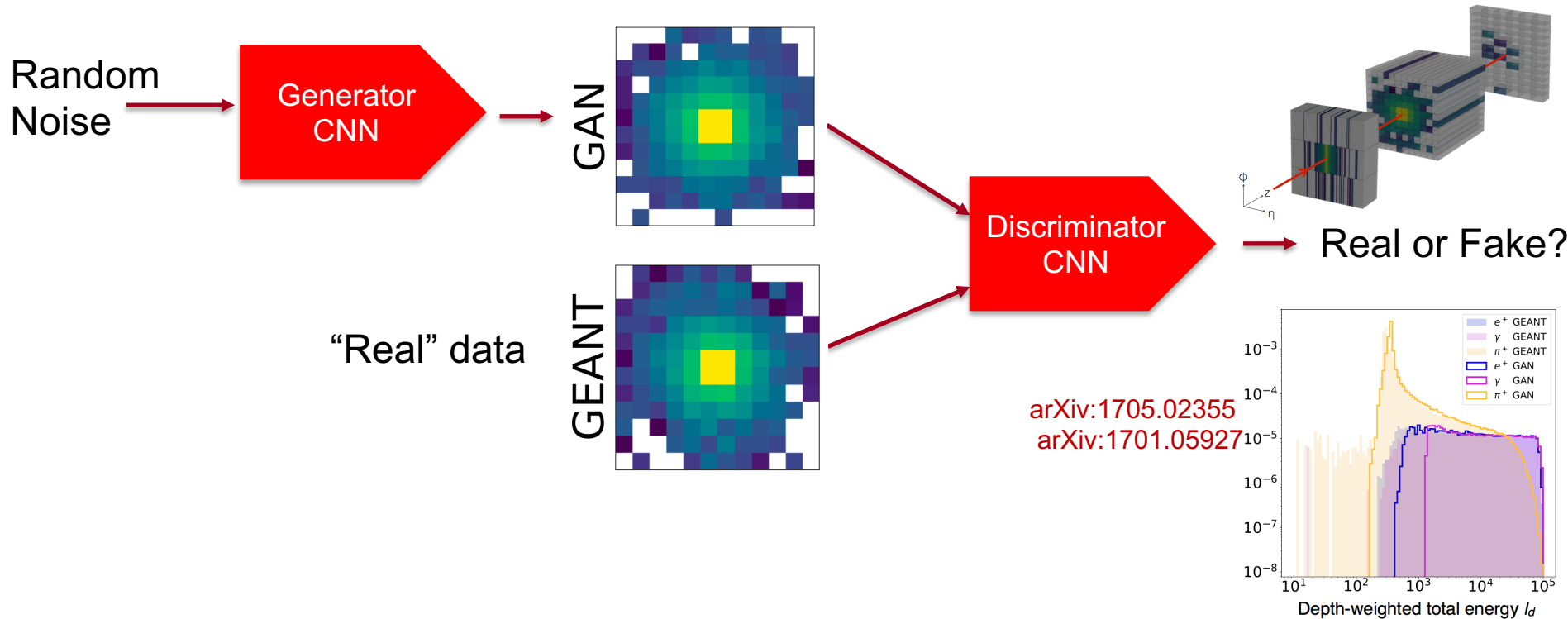# Deep Generative Models for Simulation: GANS

## *Generative Adversarial Network*

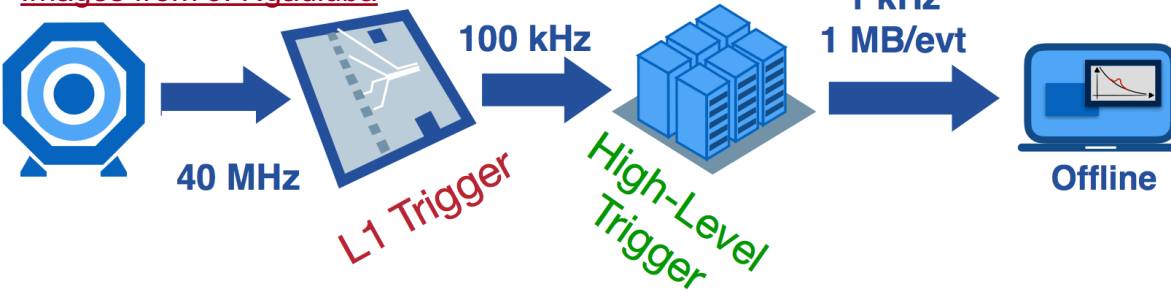Generator produces images from random noise and tries to trick discriminator into thinking they are real

Classifier tries to tell the difference between real and fake images

Images: arXiv:1710.10196

# GANs / VAEs Generating Jet-images, and 3D calo-clusters



Random Noise → Generator CNN → GAN

"Real" data → GEANT

Discriminator CNN → Real or Fake?

arXiv:1705.02355
arXiv:1701.05927

$e^+$ GEANT
$\gamma$ GEANT
$\pi^+$ GEANT
$e^+$ GAN
$\gamma$ GAN
$\pi^+$ GAN

Depth-weighted total energy $I_d$

A. Maevskiy on behalf of LHCb [slides]

Geant
GAN

Fraction of cells above threshold vs log10(Threshold/MeV)

ATLAS Simulation Preliminary
$\gamma$, $0.20 < |\eta| < 0.25$
$\chi^2$/ndf = 6.1 (GAN)
Geant4
GAN
$E_{sim}/E_{truth}$ vs $E_{truth}$ [GeV]

ATLAS-SIM-2019-004

ATLAS Simulation Preliminary
$\gamma$, E = 65 GeV, $0.20 < |\eta| < 0.25$
EM Barrel 1
$\chi^2$/ndf = 17.1 (VAE)
$\chi^2$/ndf = 20 (GAN)
Geant4
VAE
GAN
Showers / 1 GeV vs Energy [GeV]

ATL-SOFT-PUB-2018-001

# Fast Data Acquisition with ML on FPGA with HLS4ML

Images from J. Ngadiuba



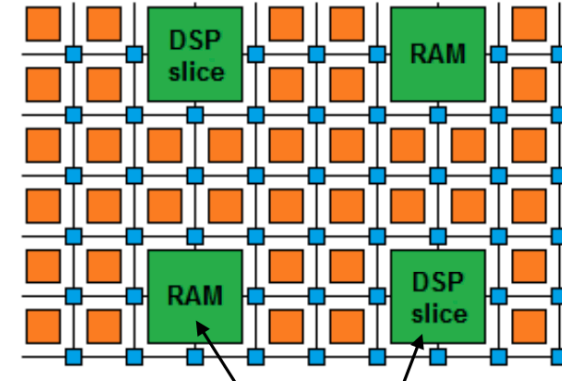40 MHz  →  L1 Trigger  →  100 kHz  →  High-Level Trigger  →  1 kHz / 1 MB/evt  →  Offline

Absorbs 100s Tb/s
Trigger decision to be made in O(μs)
Latencies require all-FPGA design

Computing farm for detailed analysis of the full event
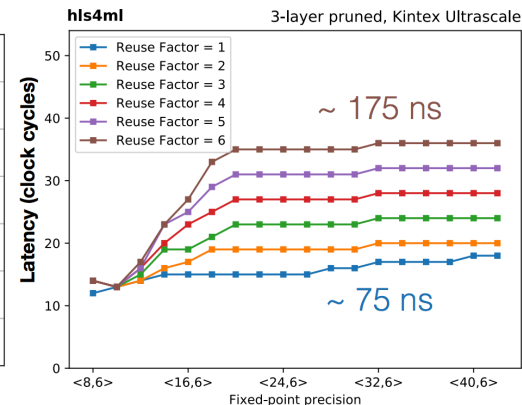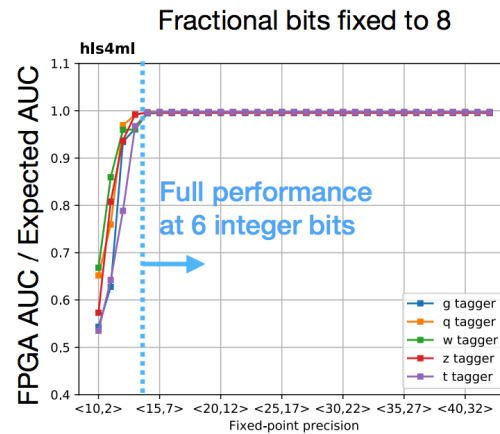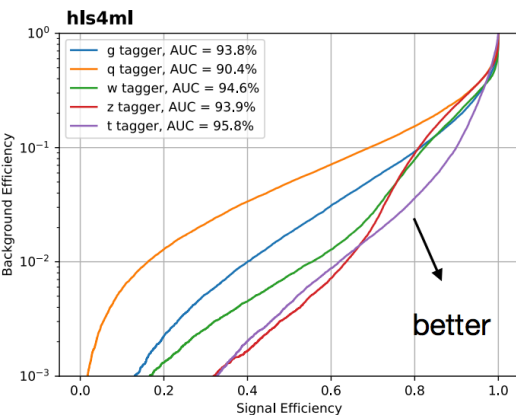Latency O(100 ms)

FPGA Diagram



FPGAs are high speed, low power, and highly parallelizable

Dedicated SW needed to efficiently and effectively port ML algorithms to FPGA

Tuning resource usage, data precision, and model pruning needed to hit timing needs

Example: Boosted jet tagging



hls4ml
g tagger, AUC = 93.8%
q tagger, AUC = 90.4%
w tagger, AUC = 94.6%
z tagger, AUC = 93.9%
t tagger, AUC = 95.8%

better

Fractional bits fixed to 8

Full performance at 6 integer bits

g tagger
q tagger
w tagger
z tagger
t tagger

3-layer pruned, Kintex Ultrascale

Reuse Factor = 1
Reuse Factor = 2
Reuse Factor = 3
Reuse Factor = 4
Reuse Factor = 5
Reuse Factor = 6

~ 175 ns

~ 75 ns

Longer latency

Each mult. used 6x
⋮
Each mult. used 3x
⋮
Fully parallel
Each mult. used 1x

More resources

# Conclusion

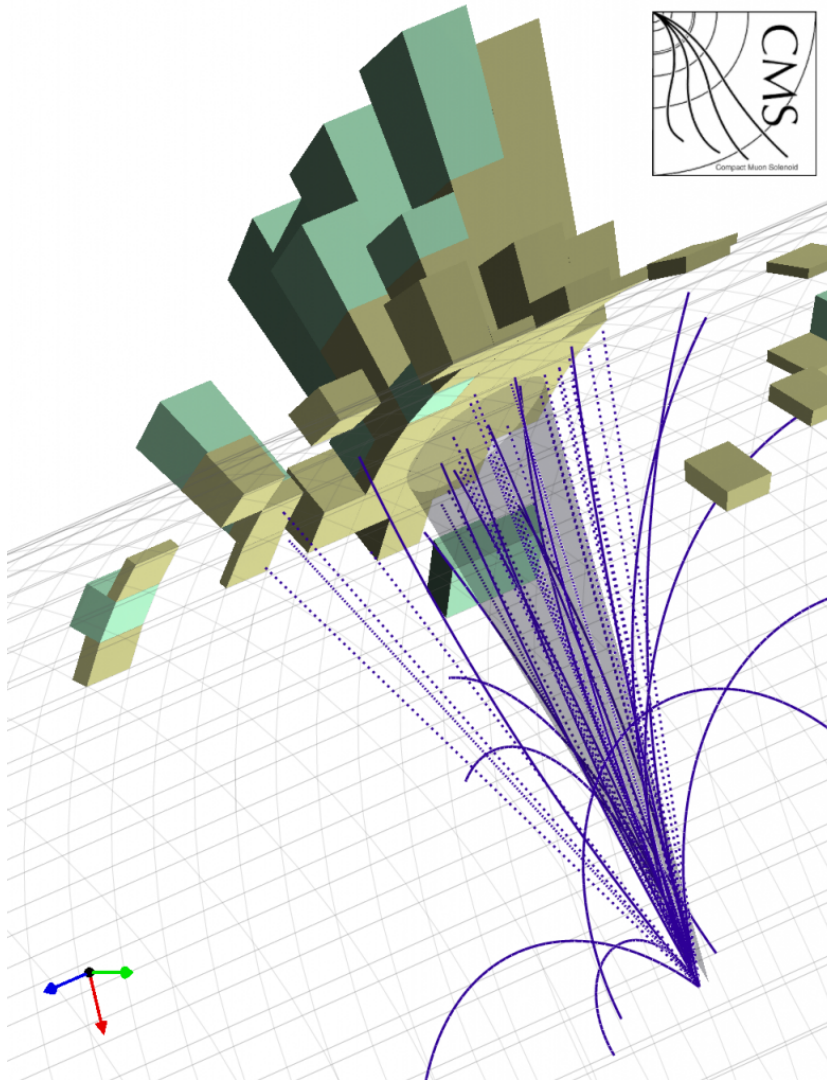The structure of analysis pipeline is grounded in our detail physics domain knowledge

We can maintain our physics knowledge embedded in this pipeline while utilizing ML to help solve some of the intractable challenges

ML methods have shown strong performance improvements in reconstruction, and techniques to deal with key experimental challenges such as computational feasibility and systematic uncertainty mitigation are under study

New ideas in data driven search strategies, fast simulation, and triggering with ML may help expand the scope of our searches!

# Backup

# Reconstructing and Tagging Particles



- **Jet**: stream of particles produced by high energy quarks and gluons
  - Clustering algorithms used to find them

Jet identification = Classification
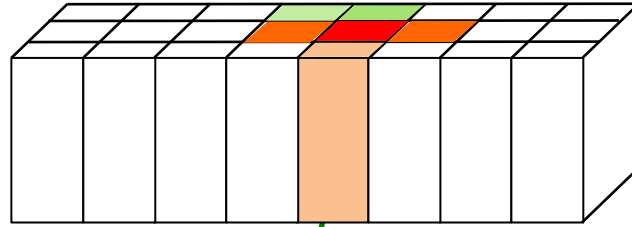
$$p(parent\ particle \mid jet\ cluster)$$

Energy estimation = Inference, Regression

$$p\left(E_{true}^{jet} \mid jet\ cluster\right)$$
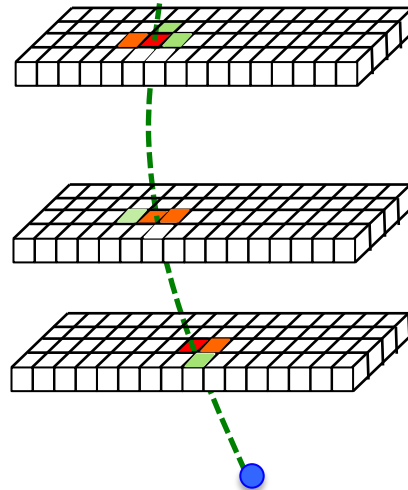
# Reconstructing and Tagging Particles

**Calorimeter:**
Stops particle and
destructively measure
energy / direction

Particle identification =

Classification

$$p(\text{electron} \mid \text{data})$$
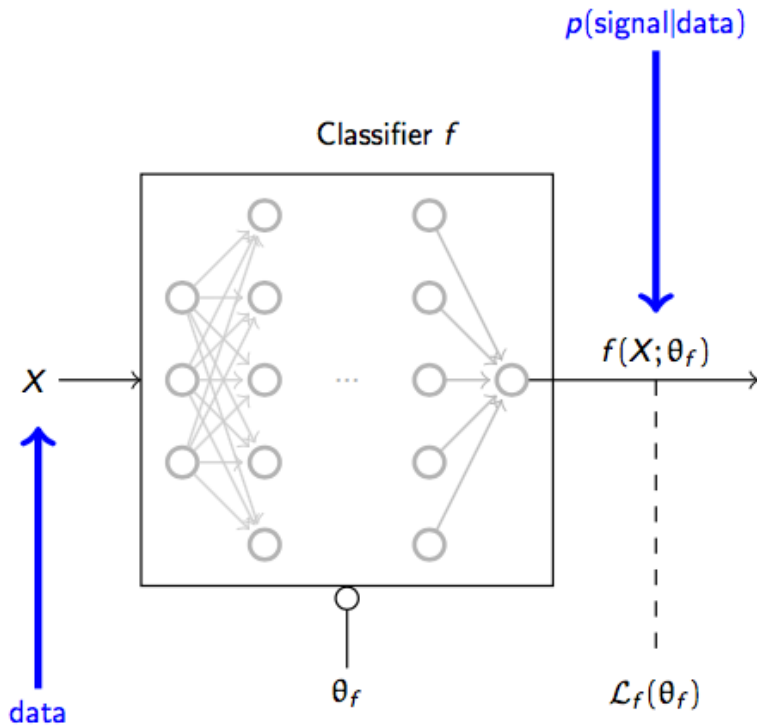
**Tracking detector:**
Typically Si-pixel detector
Non-destructive space-point
measurement

Energy estimation =

Inference, Regression

$$p(E_{\text{true}}^{\text{electron}} \mid \text{electron data})$$

# Adversarial Networks

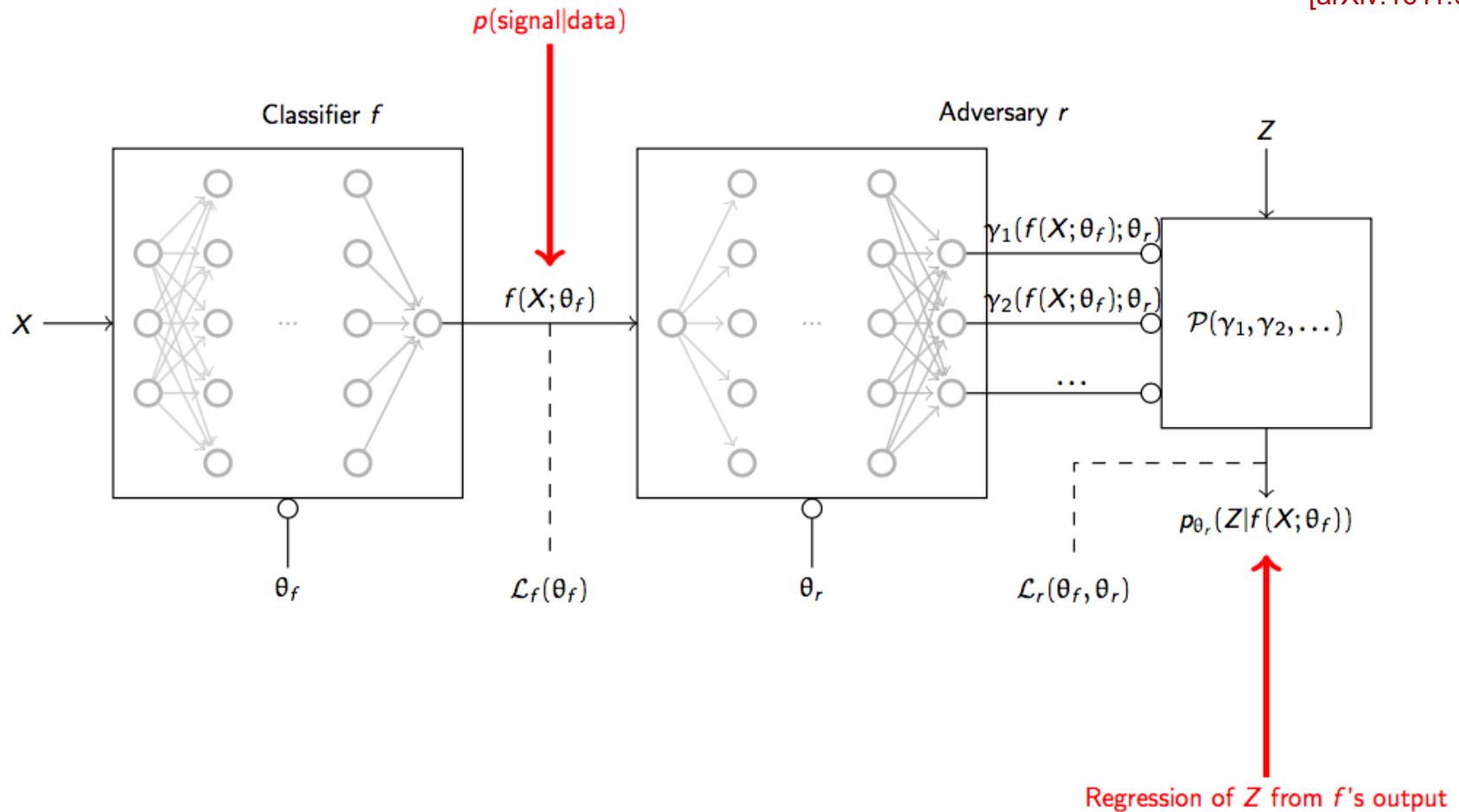$p(\text{signal}|\text{data})$

Classifier $f$

$f(X; \theta_f)$

$X$

data

$\theta_f$

$\mathcal{L}_f(\theta_f)$

Classifier built to solve problem at hand

[arXiv:1611.01046]



Regression of $Z$ from $f$'s output
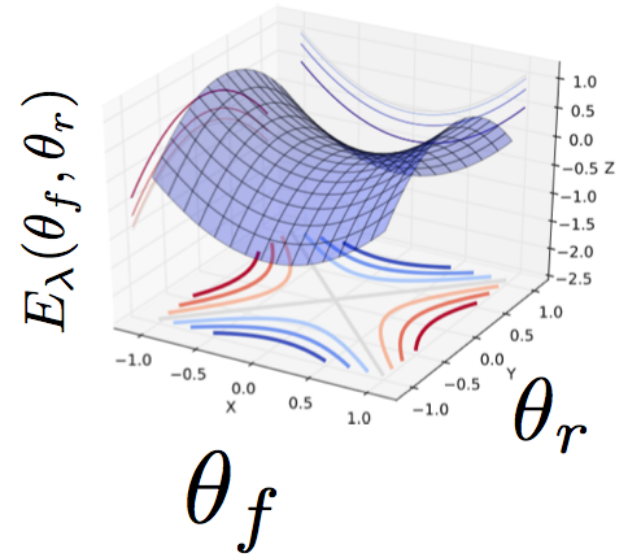
Systematic uncertainty encoded as nuisance parameters, $Z$

Adversary to predict the value of $Z$ given classifier output

# Adversarial Networks

$$\hat{\theta}_f, \hat{\theta}_r = \arg\min_{\theta_f}\max_{\theta_r} E(\theta_f, \theta_r).$$

$$E_\lambda(\theta_f, \theta_r) = \mathcal{L}_f(\theta_f) - \lambda\mathcal{L}_r(\theta_f, \theta_r),$$



Loss encodes performance of classifier and adversary

- Classifier penalized when adversary does well at predicting Z

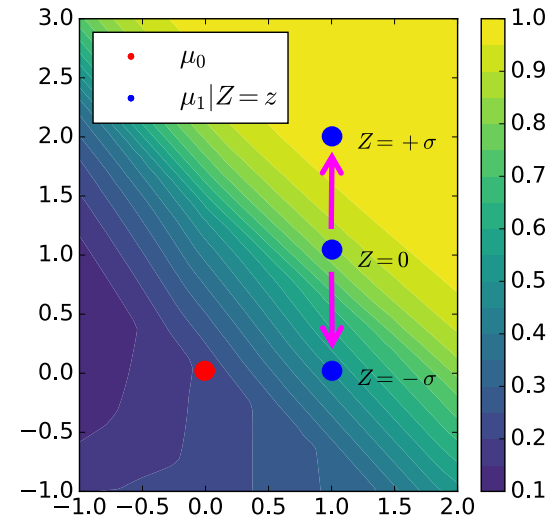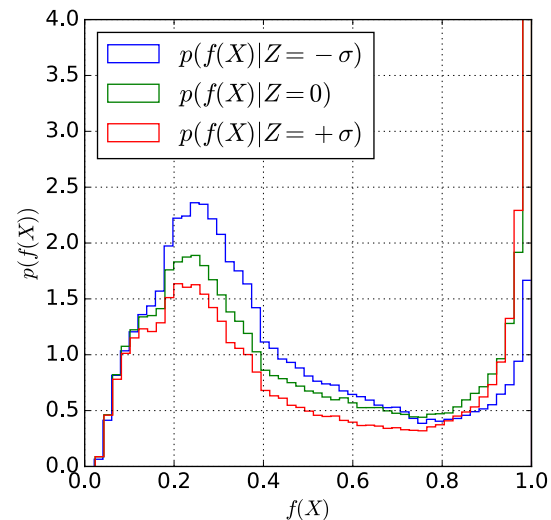Hyper-parameter λ controls trade-off

- Large λ enforces f(…) to be pivotal, e.g. robust to nuisance
- Small λ allows f(…) to be more optimal
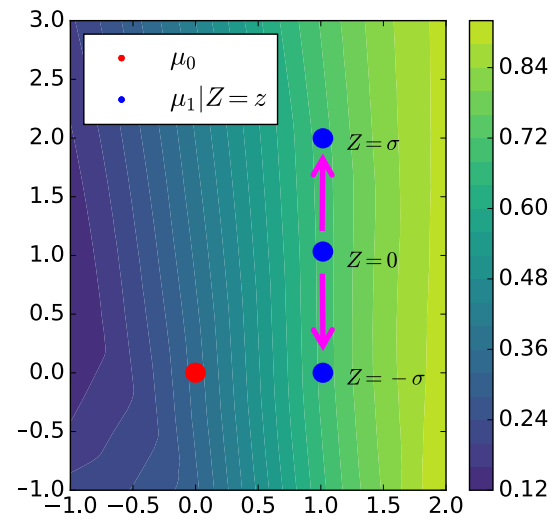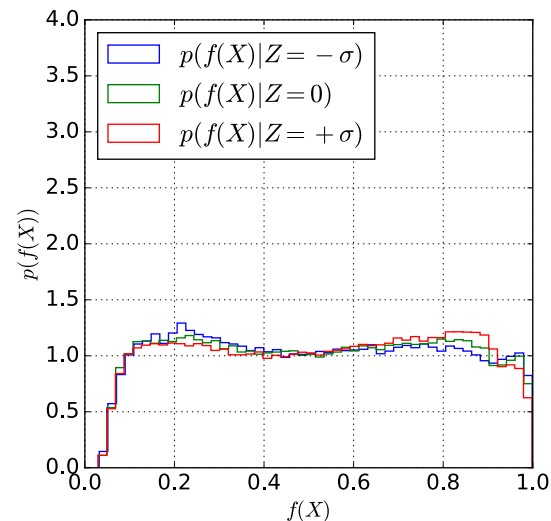
# Learning to Pivot: Toy Example

2D example

$$x \sim \mathcal{N}\left((0,0), \begin{bmatrix} 1 & -0.5 \\ -0.5 & 1 \end{bmatrix}\right) \quad \text{when } Y = 0,$$

$$x \sim \mathcal{N}\left((1, 1+Z), \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}\right) \quad \text{when } Y = 1.$$
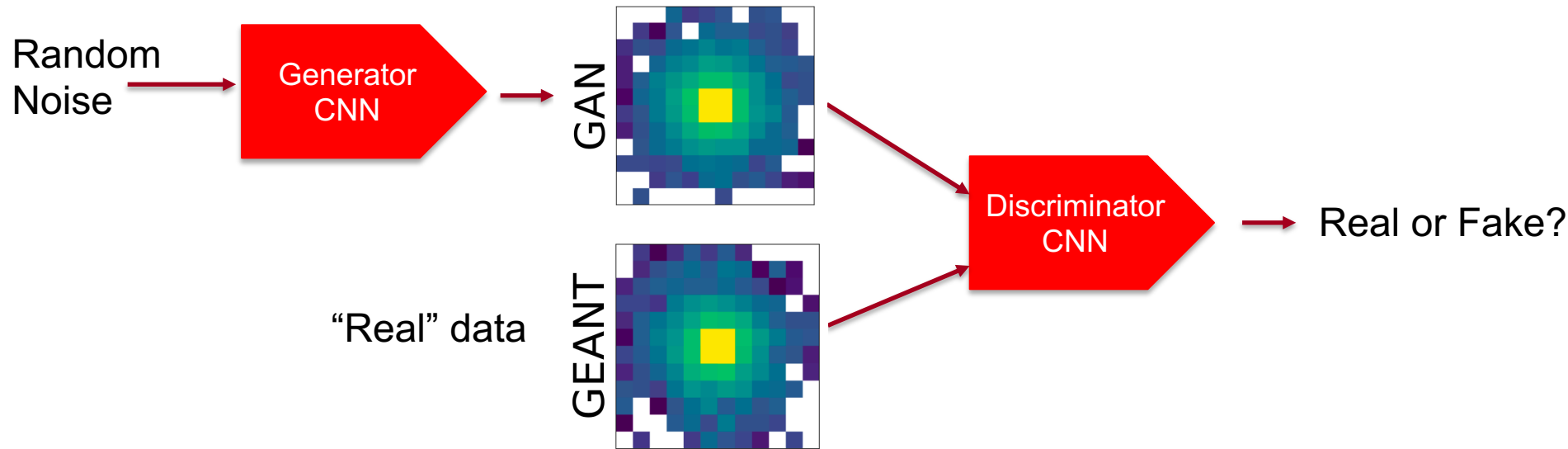


Without adversary (top) large variations in network output with nuisance parameter

With adversary (bottom) performance is independent!

# Deep Generative Models for Simulation

Random Noise → **Generator CNN** → GAN 

"Real" data → GEANT 

GAN + GEANT → **Discriminator CNN** → Real or Fake?

## Quickly growing literature

- 1701.05927  1705.02355,
- 1807.01954
- ATL-SOFT-PUB-2018-001, ATLAS-SIM-2019-004
- Slides from G Khattak, F. Carminati, S. Vallecorsa
- Slides from A. Maevskiy, et. al. on behalf of LHCb
- Slides from T. Ferber for Belle II
- Slides from V. Belavin



A. Maevskiy on behalf of LHCb



ATLAS-SIM-2019-004