

Deep Learning for High Energy Physics



Daniel Whiteson, UC Irvine
Oct 2019, LBL

What is Deep Learning?



What society thinks I do



What my friends think I do



What other computer scientists think I do



What mathematicians think I do



What I think I do

```
from theano import *
```

What I actually do

Unambiguous data

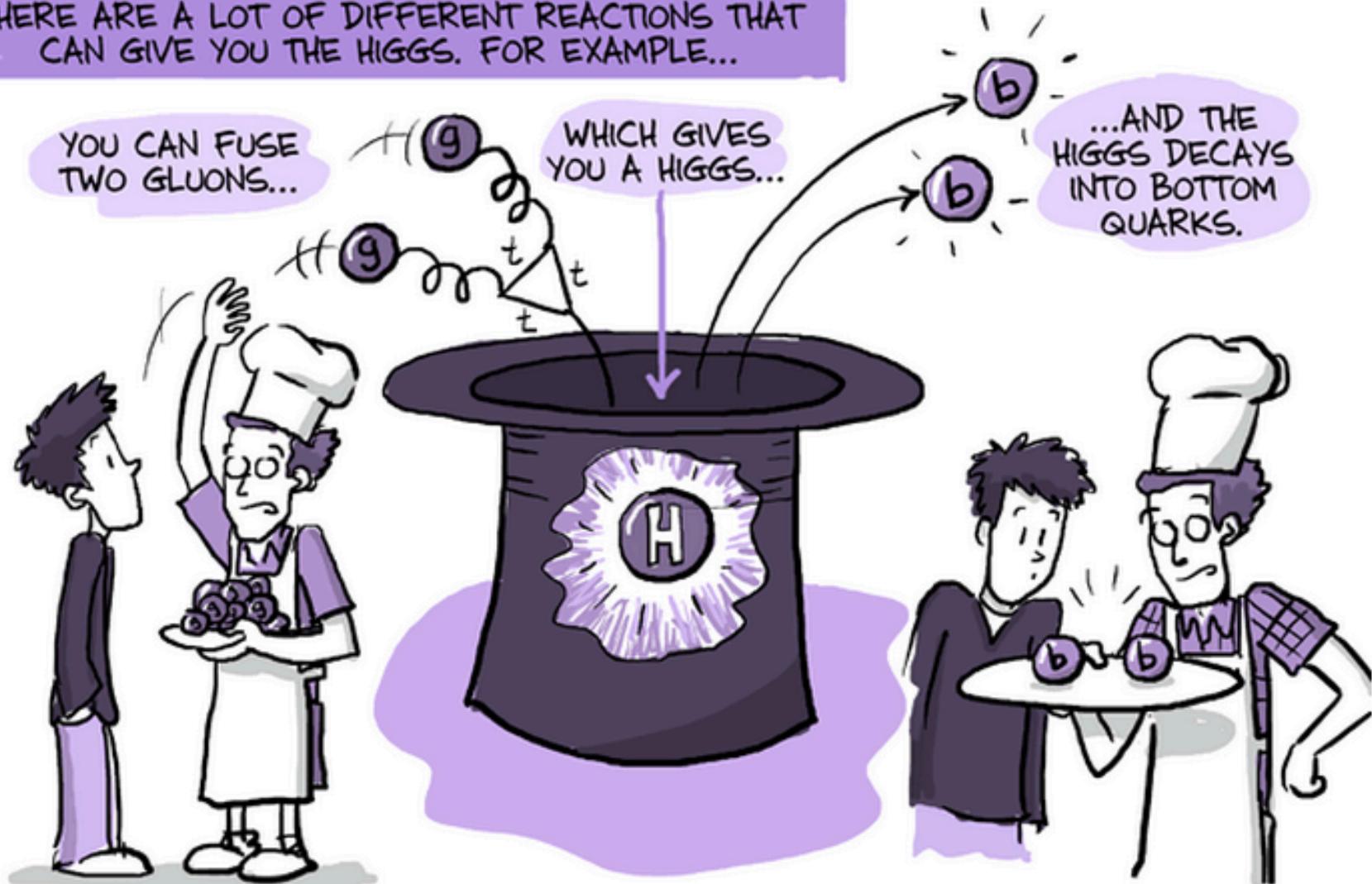


Ok, but see:

<http://cerncourier.com/cws/article/cern/54388>

Making a new particle

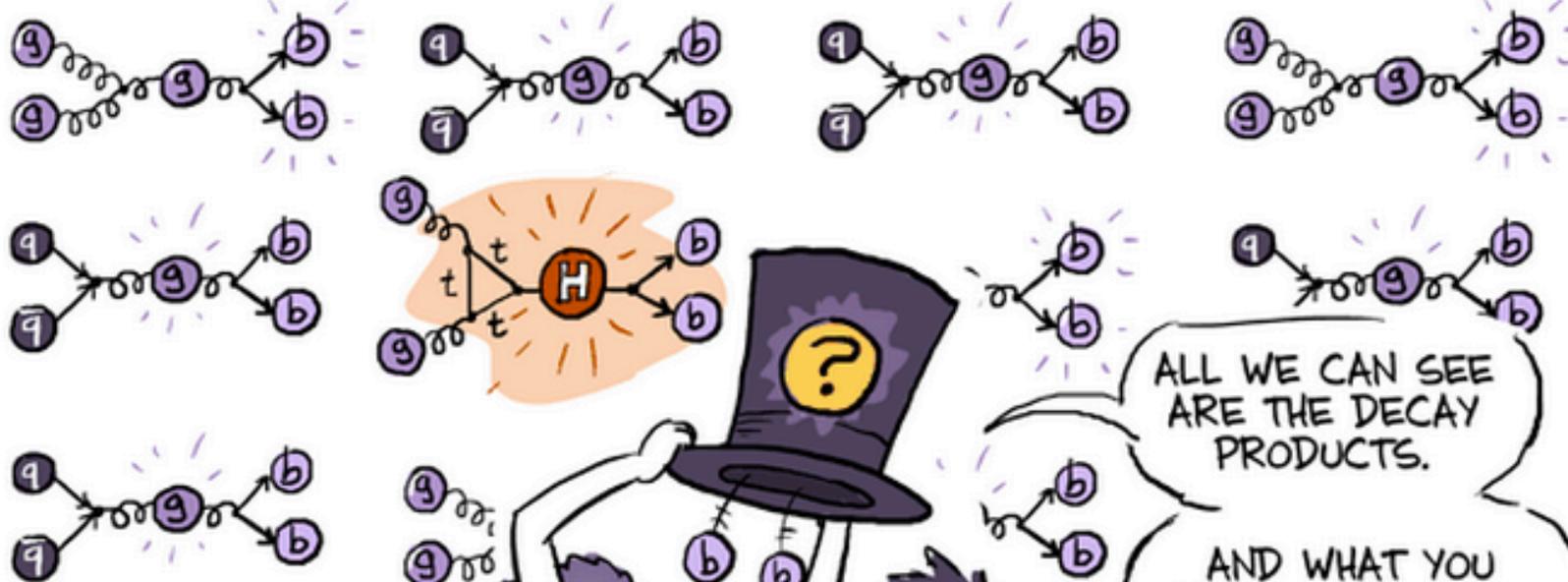
THERE ARE A LOT OF DIFFERENT REACTIONS THAT CAN GIVE YOU THE HIGGS. FOR EXAMPLE...



Backgrounds

THE PROBLEM IS, THERE'S LOTS OF OTHER WAYS YOU CAN MAKE TWO BOTTOM QUARKS:

IT'S ONE OF THE MOST COMMON THINGS TO MAKE.



JORGE CHAM © 2012

THE THING IS, WE CAN'T SEE INSIDE THESE REACTIONS...

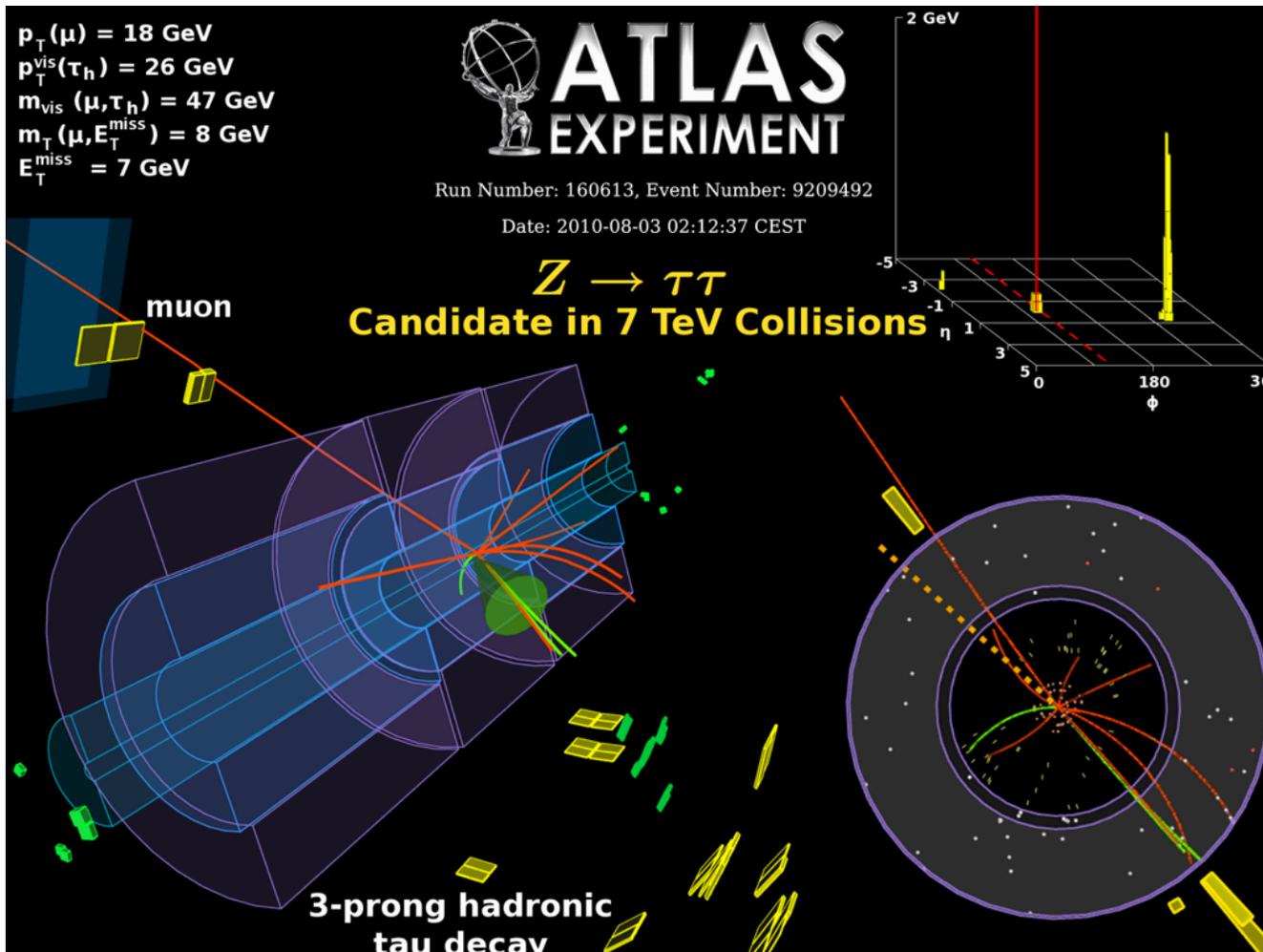
ALL WE CAN SEE ARE THE DECAY PRODUCTS.

AND WHAT YOU WANT TO KNOW IS...

DID THE HIGGS EXIST?

9

Why statistics?



The nature of our data demands it.

Hypothesis testing

To search for a new particle, we compare the predictions of two hypotheses:

1.

THE STANDARD MODEL			
Fermions			
Quarks	<i>u</i> up	<i>c</i> charm	<i>t</i> top
	<i>d</i> down	<i>s</i> strange	<i>b</i> bottom
Leptons	V_e electron neutrino	V_μ muon neutrino	V_τ tau neutrino
	<i>e</i> electron	μ muon	τ tau

Hypothesis testing

To search for a new particle, we compare the predictions of two hypotheses:

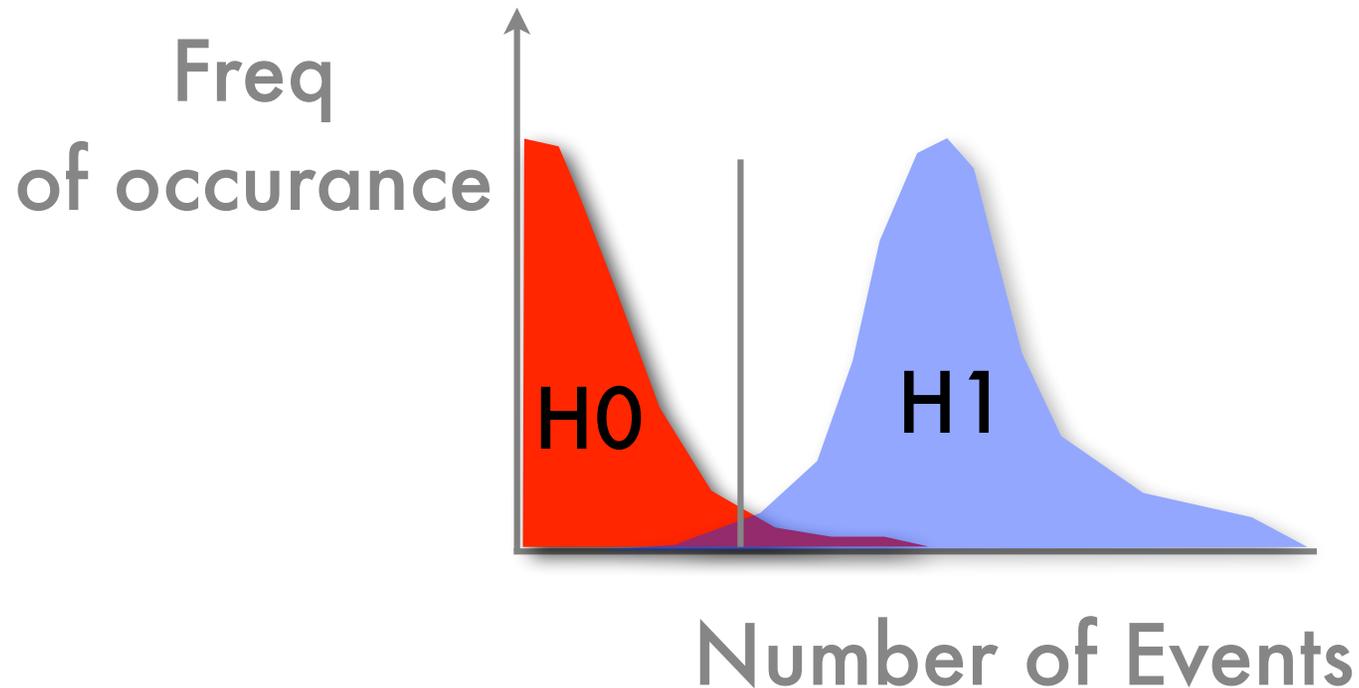
1.

THE STANDARD MODEL			
Fermions			
Quarks	u up	c charm	t top
	d down	s strange	b bottom
Leptons	ν_e electron neutrino	ν_μ muon neutrino	ν_τ tau neutrino
	e electron	μ muon	τ tau

2.

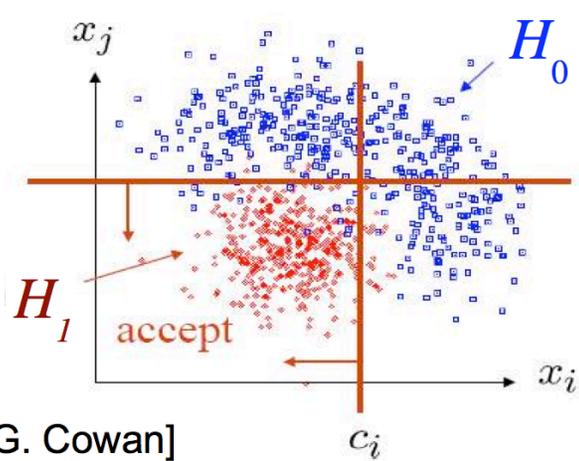
THE STANDARD MODEL PLUS X				
Fermions				
Quarks	u up	c charm	t top	X
	d down	s strange	b bottom	
Leptons	ν_e electron neutrino	ν_μ muon neutrino	ν_τ tau neutrino	
	e electron	μ muon	τ tau	

Example

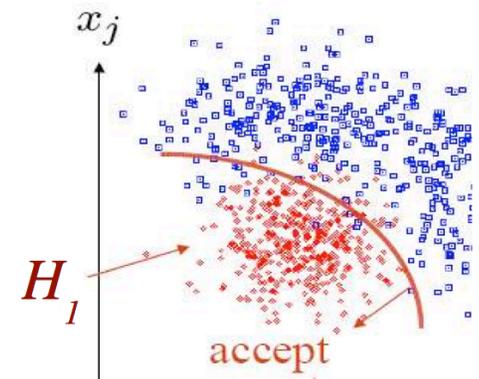
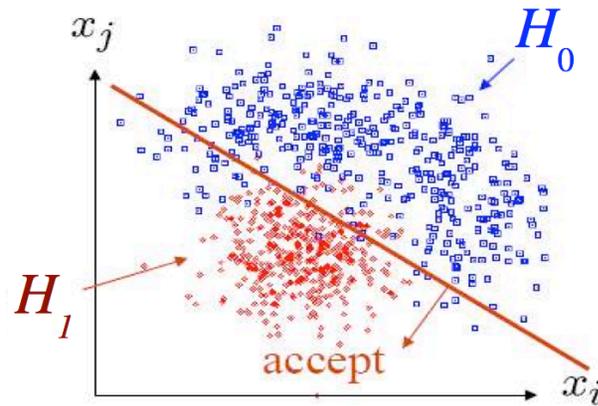


A threshold makes sense.
Choice of position balances
false vs **missed** discovery

More complicated



[G. Cowan]



Neyman-Pearson

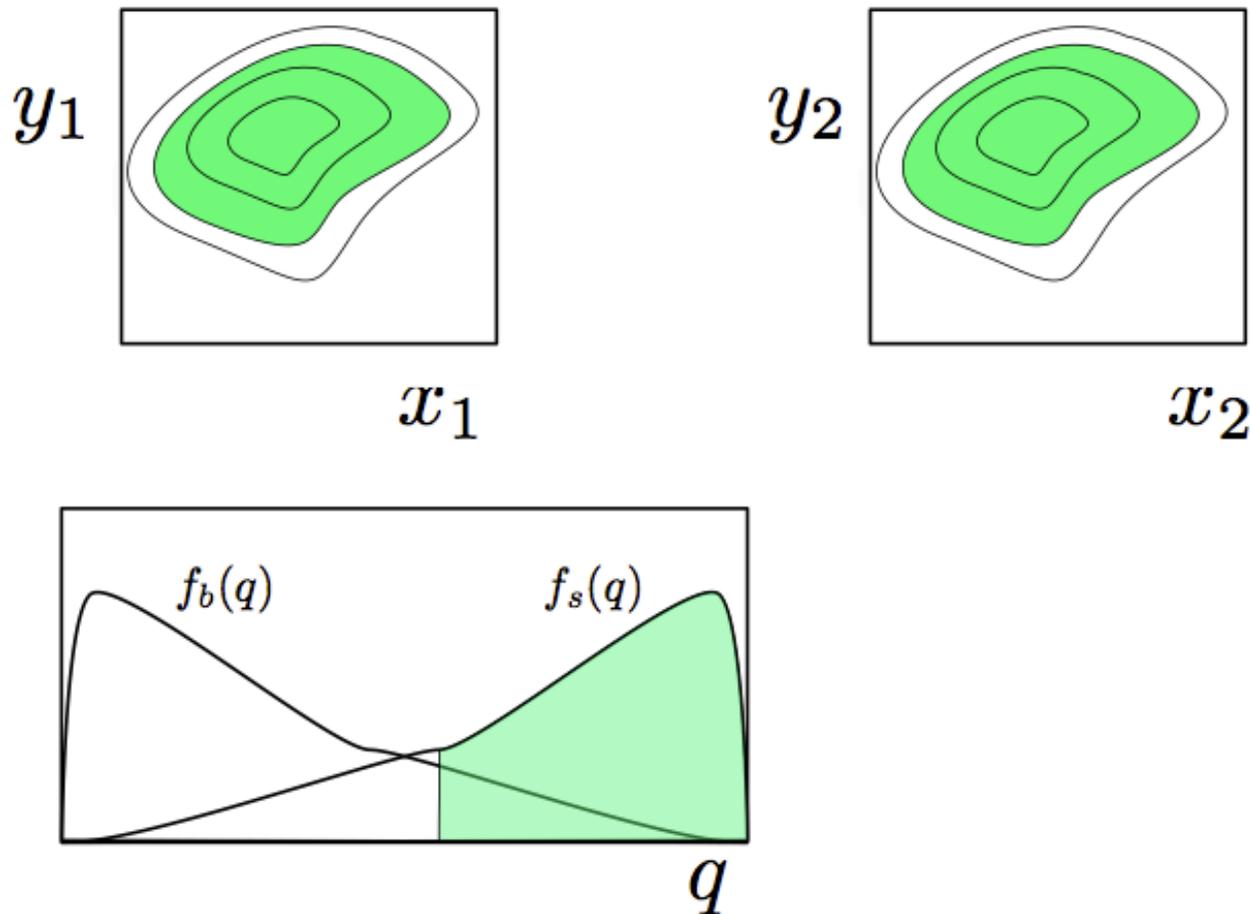
NP lemma says that the best decision boundary is the **likelihood ratio**:

$$\frac{P(x|H_1)}{P(x|H_0)} > k_\alpha$$

(Gives smallest missed discovery rate for fixed false discovery rate)

What does this do?

Finds a region in variable space



No problem

Fairly straightforward

if you can calculate

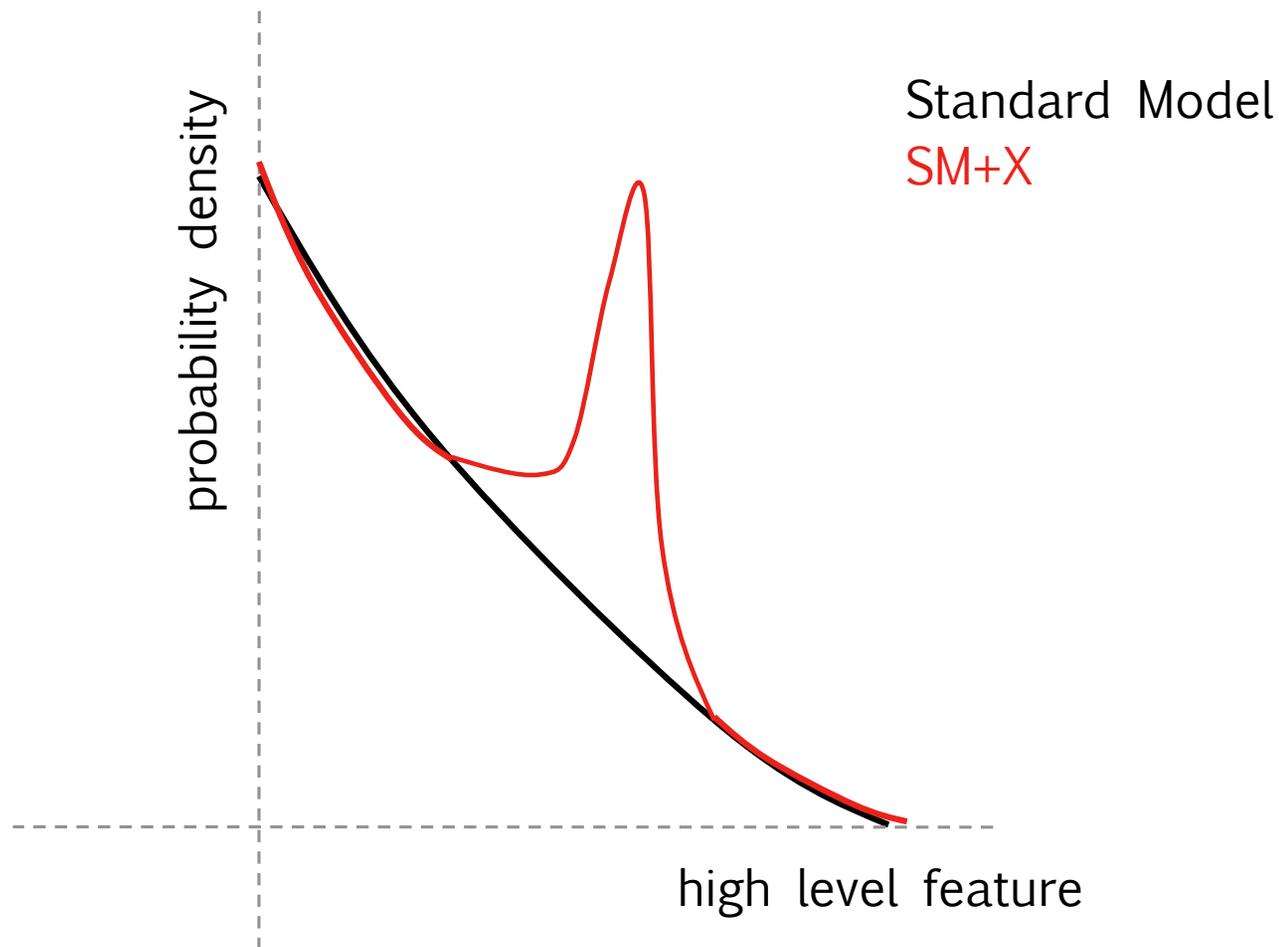
$$\frac{P(x|H_1)}{P(x|H_0)}$$

or generally

$$P(\text{data} | \text{theory})$$

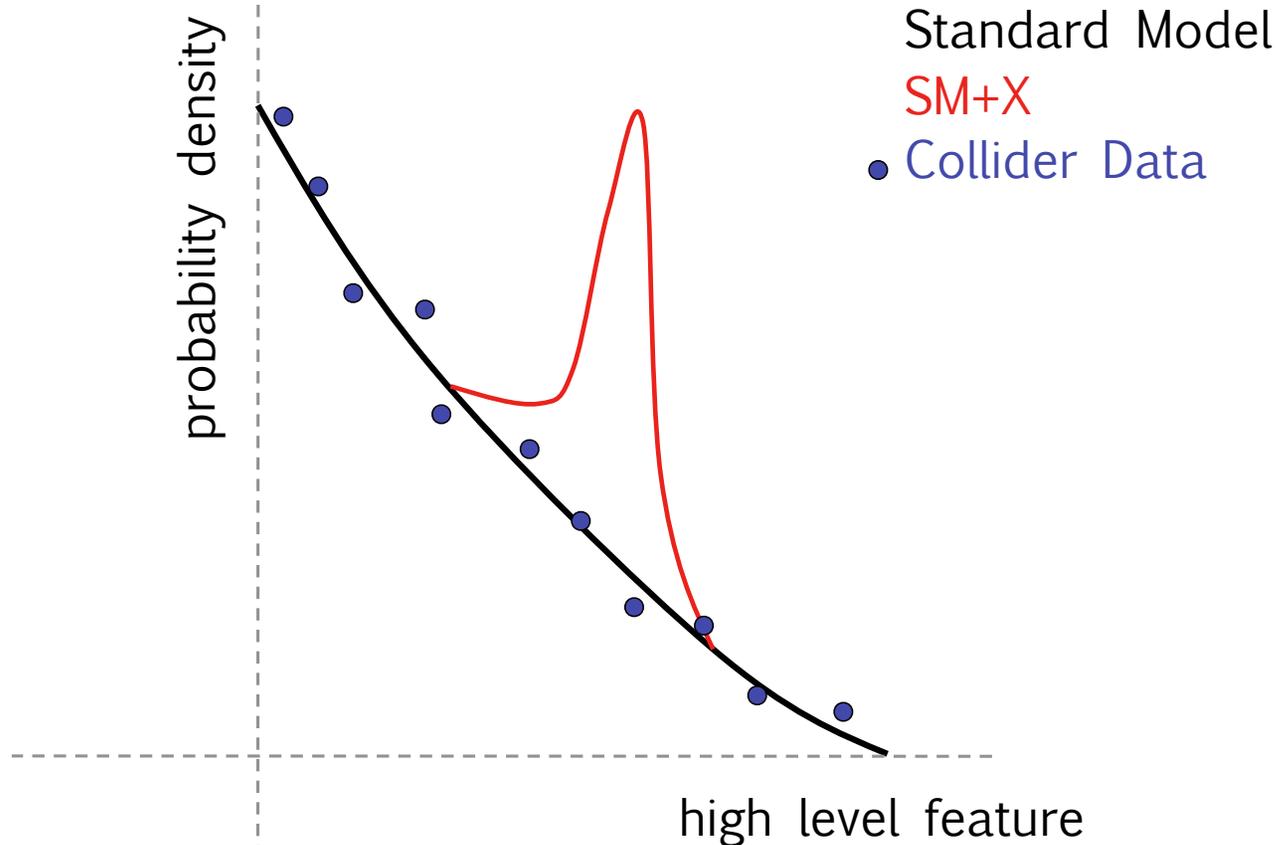
Hypothesis Testing

Sometimes this is easy



Hypothesis Testing

We can compare the predictions to the collider data



Which can tell us which hypothesis is preferred via a likelihood ratio:

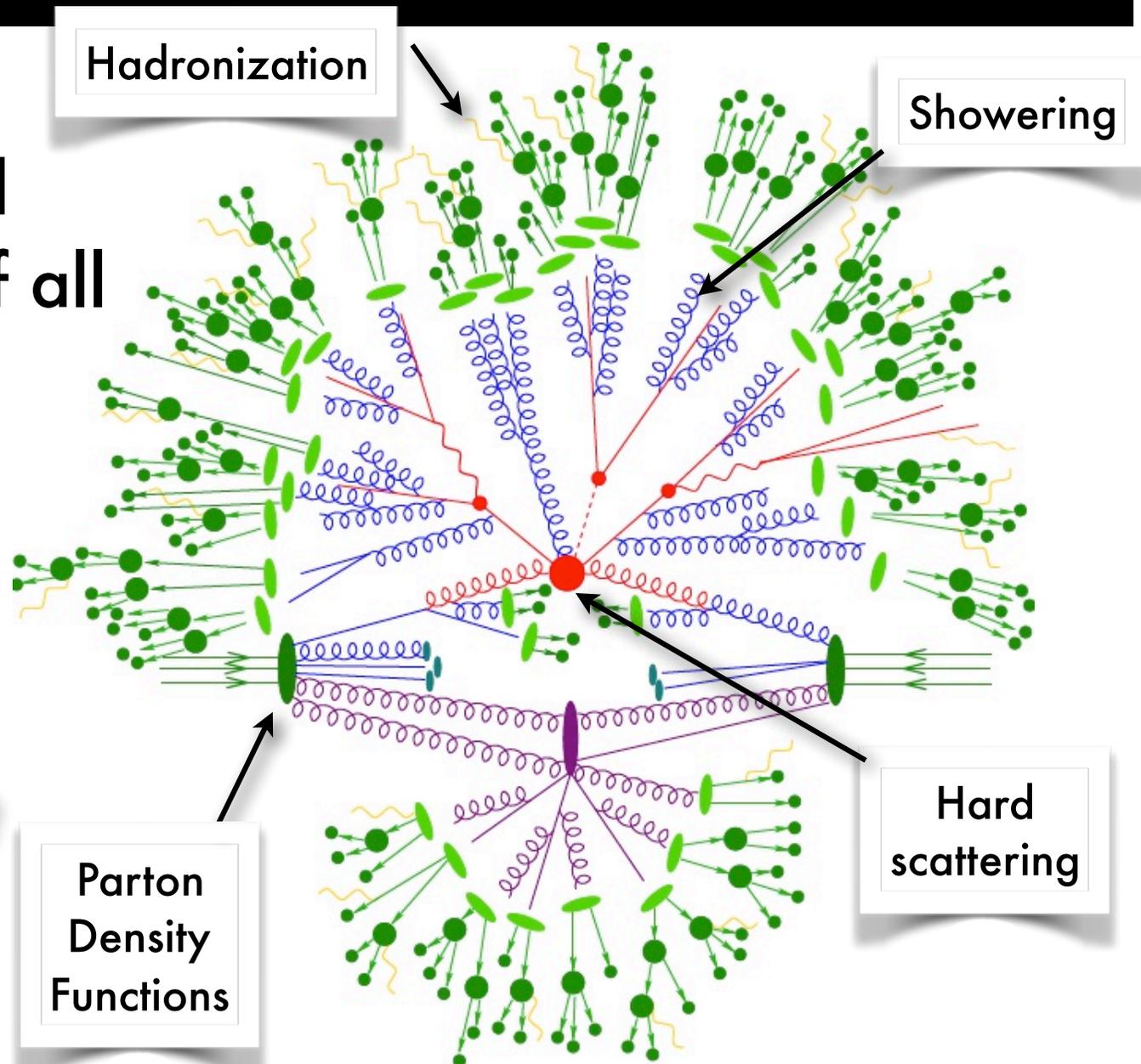
$$\frac{L_{SM+X}}{L_{SM}} = \frac{P(\text{data} \mid \text{SM+X})}{P(\text{data} \mid \text{SM})}$$

In general

We have a good understanding of all of the pieces

Do we have

$P(\text{data} | \text{theory})?$



In general

What would

$P(\text{data} \mid \text{theory})$

look like?

The dream

Detector Response

$$p(\text{data} \mid \text{final-state particles } P)$$

Hadronization

$$\times p(\text{final state particles } P \mid \text{showered particles } S)$$

Showering

$$\times p(\text{showered particles } S \mid \text{hard scatter products } M)$$

$$\times p(\text{hard scatter products } M \mid \text{theory})$$

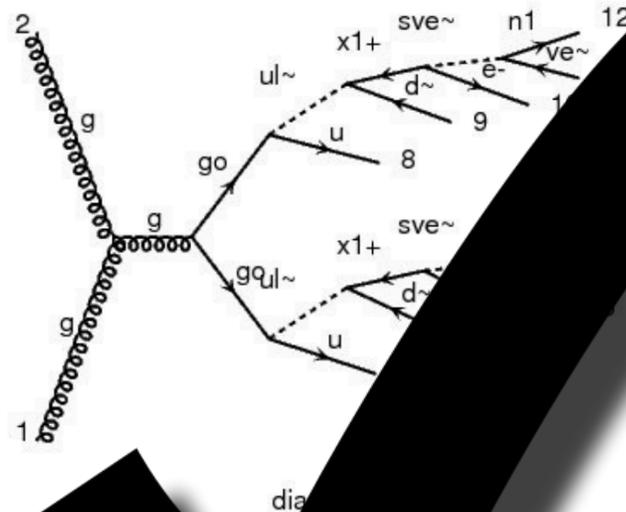
Hard scattering

Sum over all intermediate P, S, M

Parton
Density
Functions

The dream

$p(\text{hard scatter products } M \mid \text{theory})$



The \mathcal{M} defined
automatically \mathcal{M} exist
for almost any (B)SM theory

The nightmare

$p(\text{data} \mid \text{final-state particles } P)$

$\times p(\text{final state particles } P \mid \text{showered particles } S)$

$\times p(\text{showered particles } S \mid \text{hard scatter products } M)$

We have: solid understanding of microphysics

We need: analytic description of high-level physics

The solution

We have: solid understanding of microphysics

We need: analytic description of high-level physics

But: only heuristic lower-level approaches exist

Iterative simulation strategy, **no overall PDF**

Iterative approach

- (1) Draw events from $p(M | \text{theory})$
- (2) add random showers
- (3) do hadronization
- (4) simulate detector

The solution

We have: solid understanding of microphysics

We need: analytic description of high-level physics

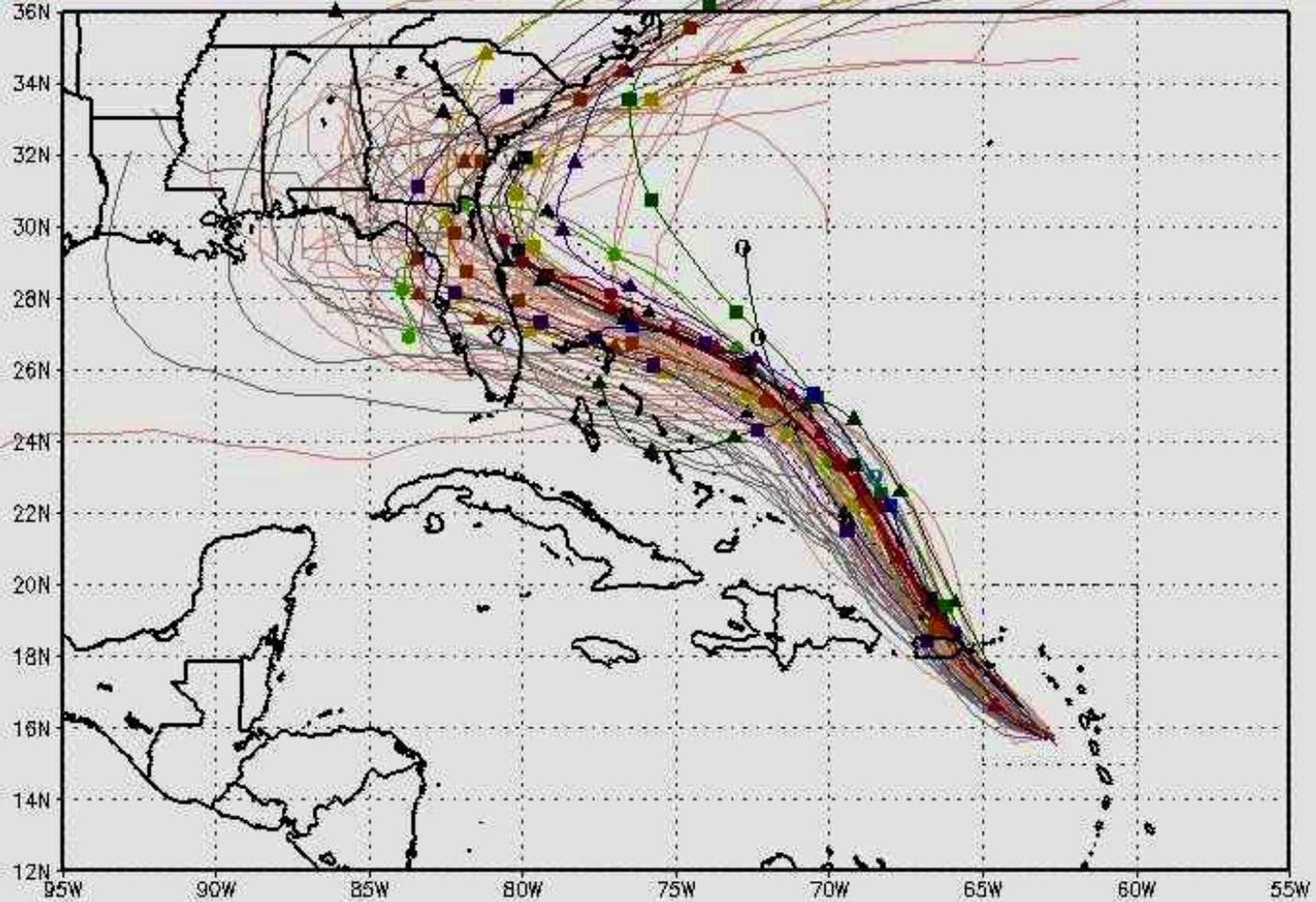
But: only heuristic lower-level approaches exist

Iterative simulation strategy, **no overall PDF**

What do we get

Arbitrarily large samples of events
drawn from $p(\text{data} | \text{theory})$, but **not**
the PDF itself

- ▲--- XTRP 28/0600Z
- CLP5 28/0600Z
- ▲--- HMON 28/0000Z
- ▲--- AVNO 28/0000Z
- ▲--- ECMF 28/0000Z
- TVCN 28/0600Z
- ▲--- TABD 28/0600Z
- HWRF 28/0000Z
- AEMN 28/0000Z
- EEMN 28/0000Z
- ▲--- TVCX 28/0600Z
- TABM 28/0600Z
- UKM 28/0000Z
- APxx 28/0000Z
- EExx 28/0000Z
- NHC 28/0900Z
- TABS 28/0600Z
- COTC 28/0000Z
- ▲--- CMC 28/0000Z
- GEMN 28/0000Z



storm_05
 sfwmd.gov
 weather@sfwmd.gov
 28-Aug 08:06EDT

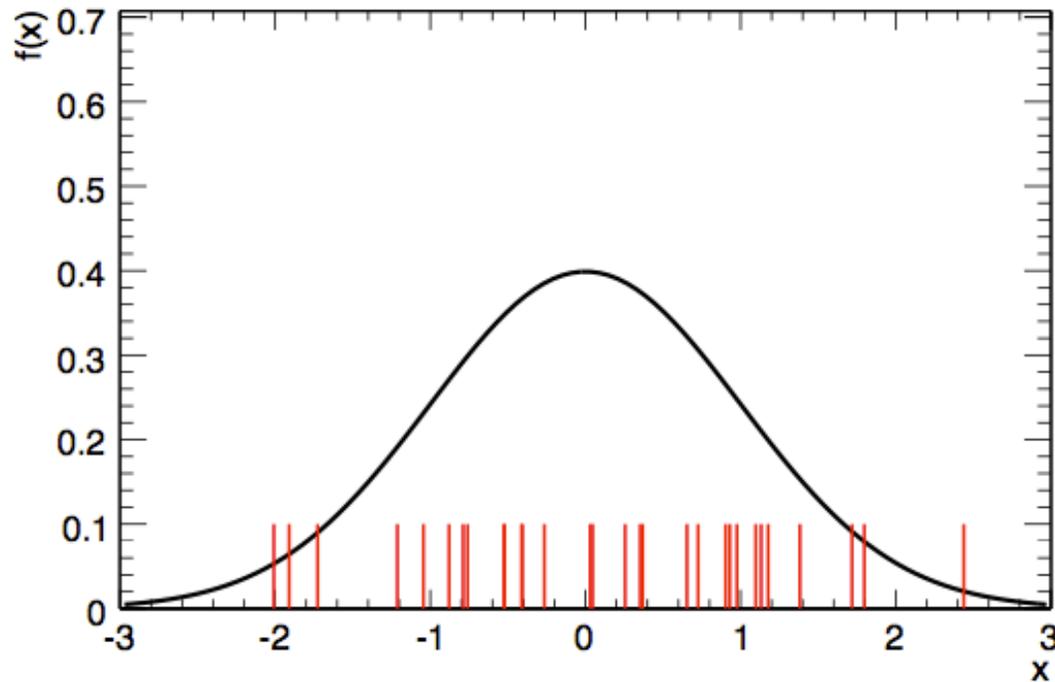
NHC Advisories and County Emergency Management Statements supersede this product.
 This graphic should complement, not replace, NHC discussions.
 If anything on this graphic causes confusion, ignore the entire product.
 For full info, see <http://my.sfwmd.gov/sfwmd/common/images/weather/plots.html>



The problem

Don't know PDF, have events drawn from PDF

$$f_{emp} = \frac{1}{N} \sum_i^N \delta(x - x_i)$$



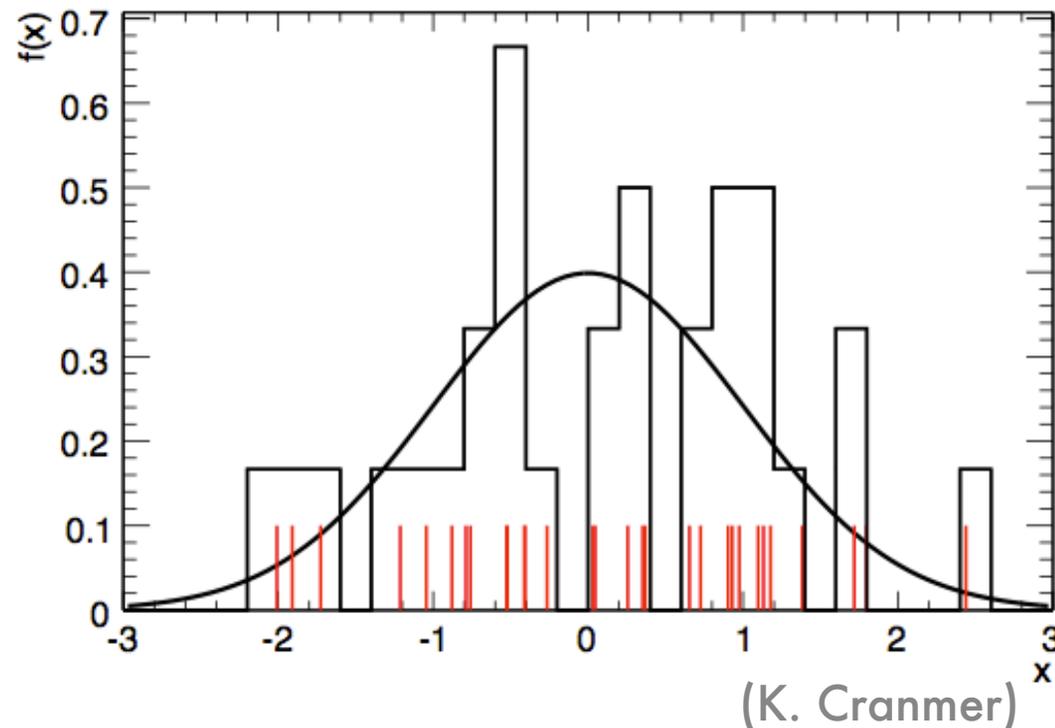
(K. Cranmer)

Need to recreate PDF

MC events to PDF

Simple approach : histogram

$$f_{hist}^{w,s}(x) = \frac{1}{N} \sum_i h_i^{w,s}$$



Curse of Dimensionality

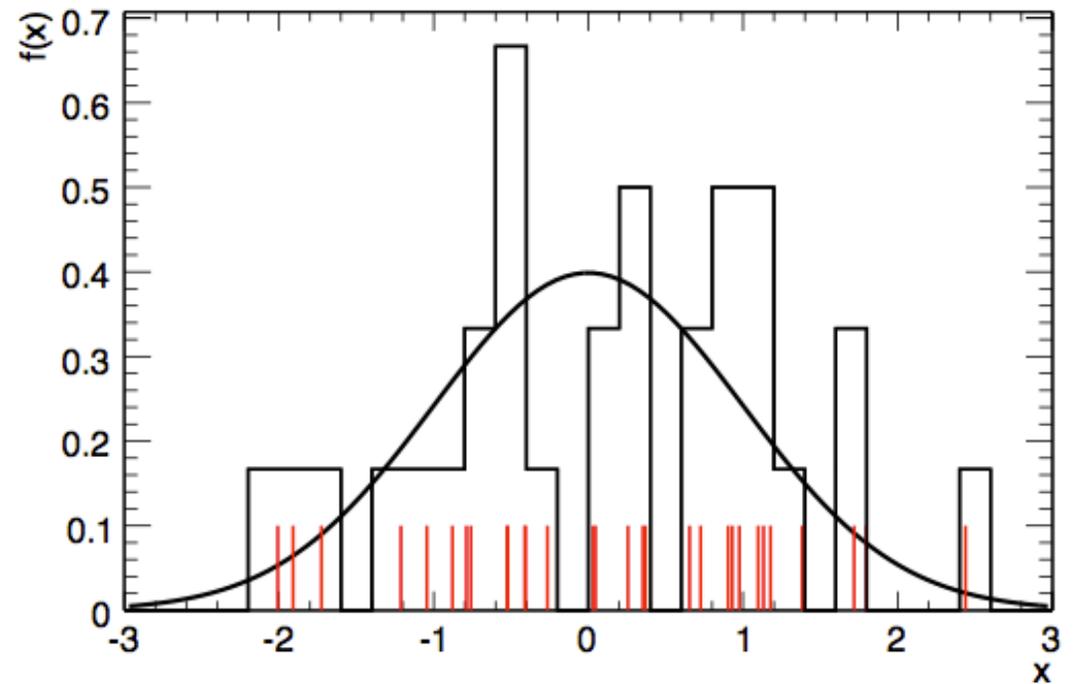
How many events
do you need
to describe a 1D
distribution? $O(100)$

An n-D distribution?

$O(100^n)$

!!

$$f_{hist}^{w,s}(x) = \frac{1}{N} \sum_i h_i^{w,s}$$



(K. Cranmer)

The nightmare

f(data | final-state particles P)

x f(final state particles P | showered particles S)

x f(showered particles S | hard scatter products M)

“data” is a 100M-d vector!

The nightmare

f(data | final-state particles P)

x f(final state particles S)

x f(showered matter products M)

// vector!



Task for ML

Find a function:

$$f(\bar{x}) : \mathbb{R}^N \rightarrow \mathbb{R}^1$$

which contains the same
hypothesis testing power
as

$$\frac{P(x|H_1)}{P(x|H_0)} > k_\alpha$$

Neural networks

Strategy:

$$f(\vec{x}) : \mathbb{R}^N \rightarrow \mathbb{R}^1$$

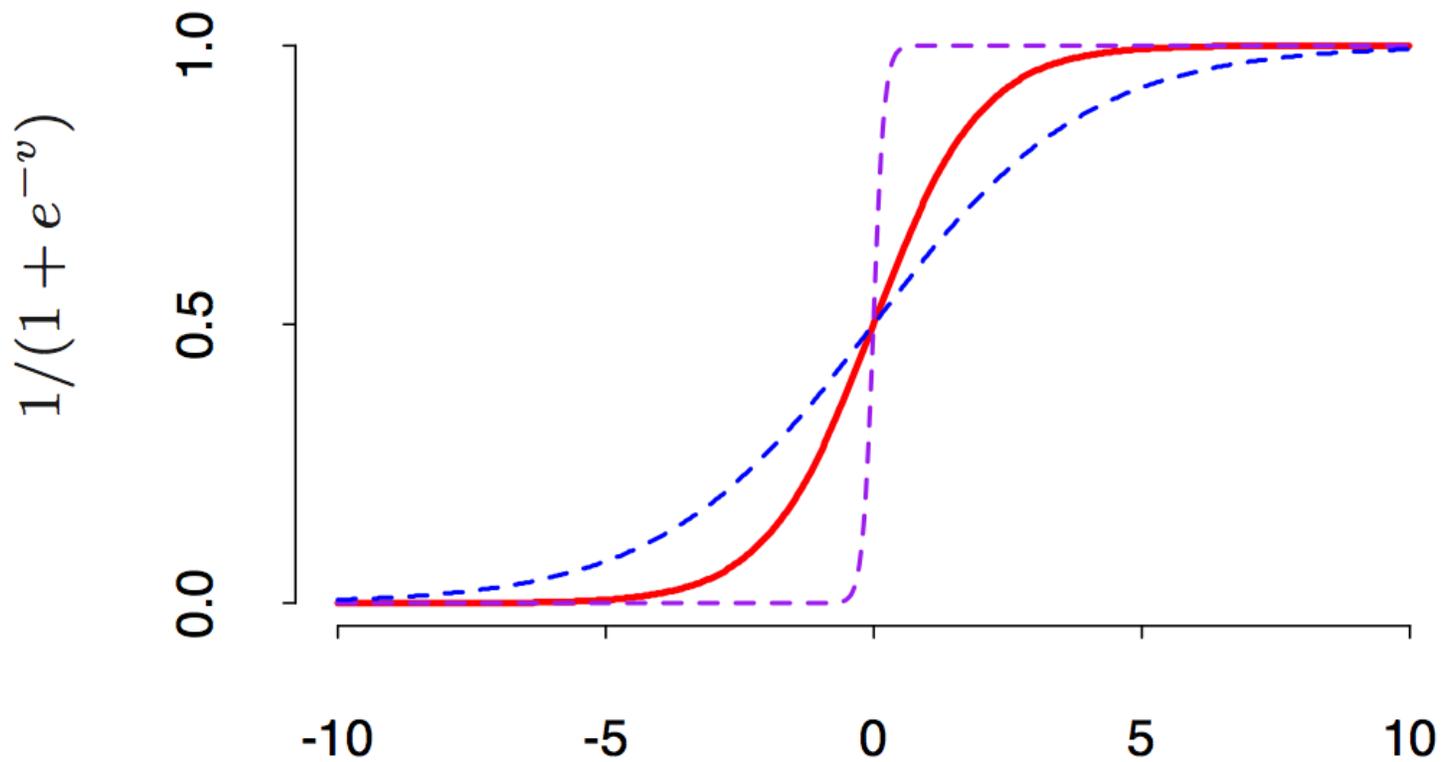
Build $f(\vec{x})=y(\vec{x})$ out of a pile of convoluted mini-functions

$$y(\vec{x}) = h\left(w_0 + \sum_{i=1}^n w_i x_i\right)$$

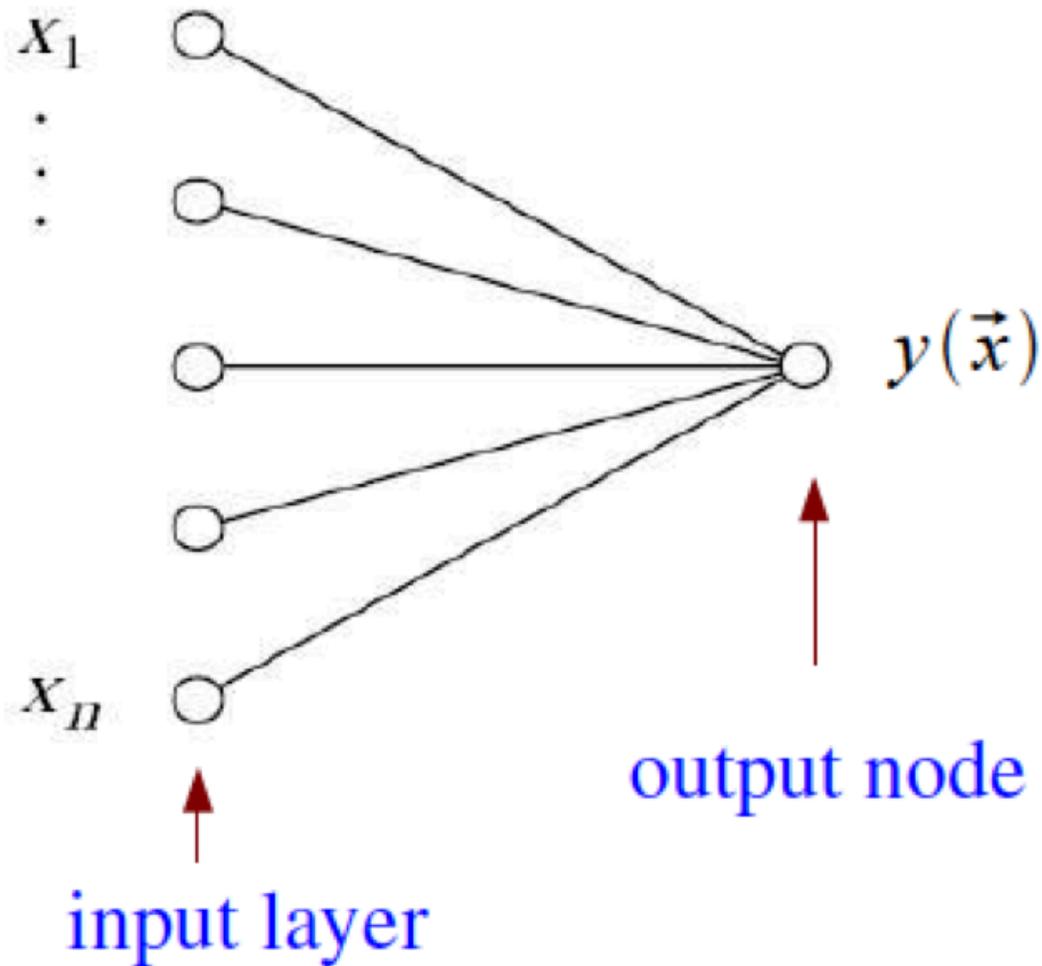
here $h()$ is a non-linear *activation function* and the w factors are *unknown parameters*

Neuron

Example activation function



Simple visualization



Finding good weights

We have

a weight space
a quality metric

$$y(\vec{x}) = h\left(w_0 + \sum_{i=1}^n w_i x_i\right)$$

$$E(\mathbf{w})$$

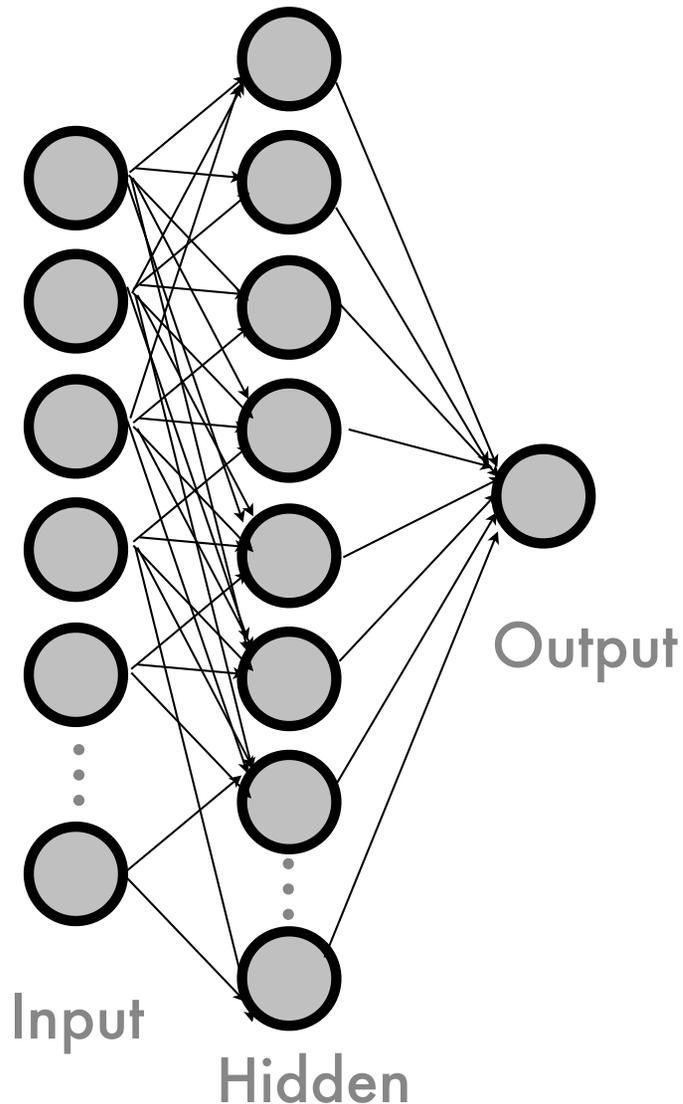
We need

to find the max quality (or min error)

Search the space!

How complex?

Essentially a functional fit with many parameters



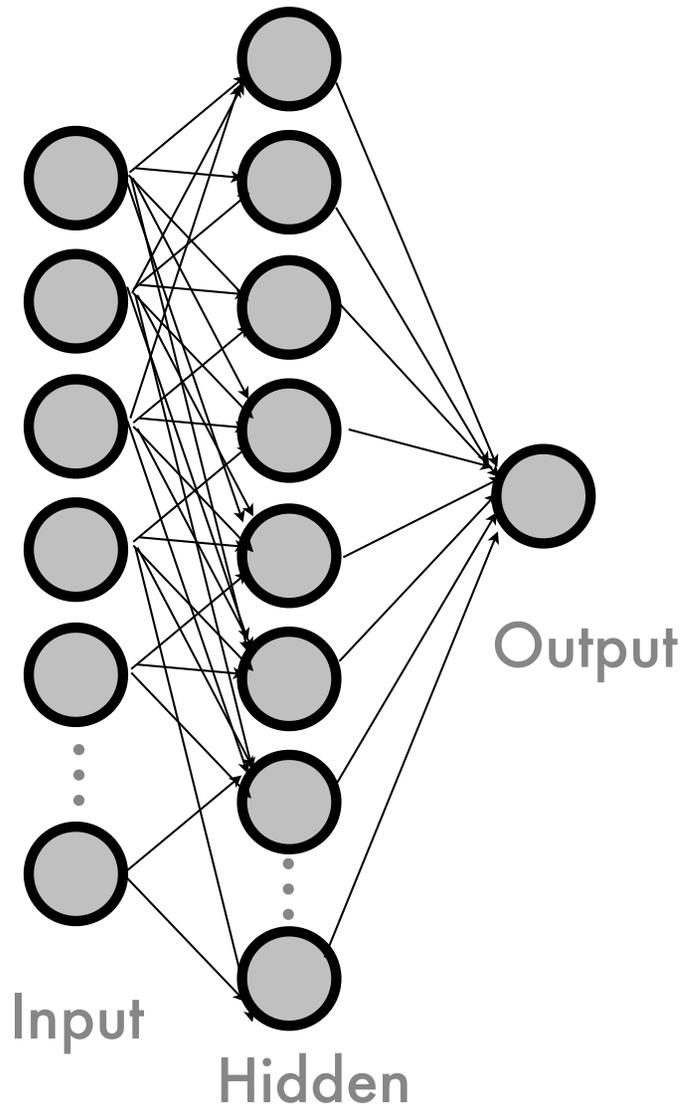
Single hidden layer

In theory any function can be learned with a single hidden layer.

But might require very large hidden layer

Neural Networks

Essentially a functional fit with many parameters



Problem:

Networks with > 1 layer are very difficult to train.

Consequence:

Networks are not good at learning non-linear functions.
(like invariant masses!)

In short:

Can't just throw 4-vectors at NN.

Search for Input

ATLAS-CONF-2013-108

Can't just use $4v$

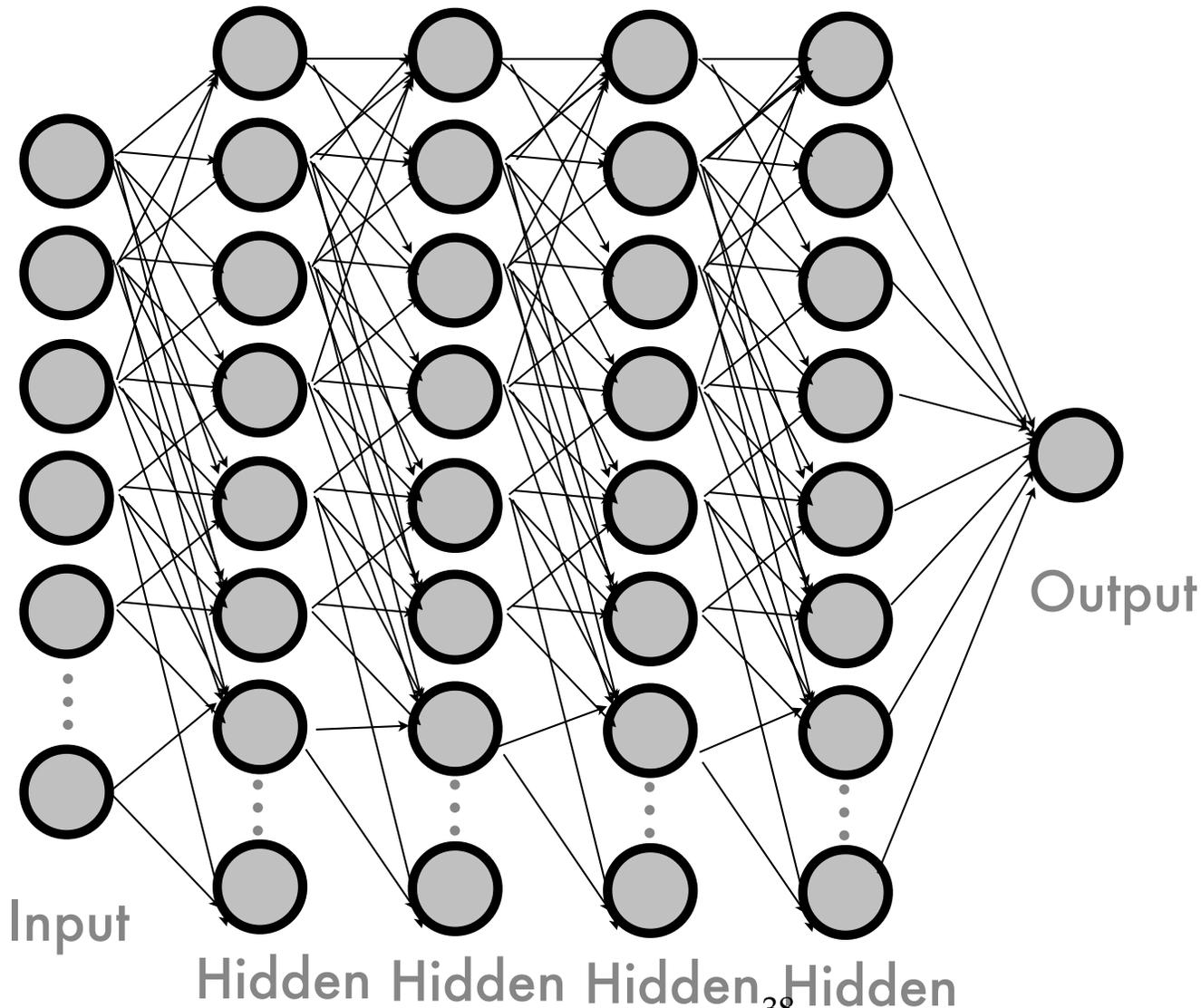
Can't give it too many inputs

Painstaking search through input feature space.

Variable	VBF			Boosted		
	$\tau_{\text{lep}}\tau_{\text{lep}}$	$\tau_{\text{lep}}\tau_{\text{had}}$	$\tau_{\text{had}}\tau_{\text{had}}$	$\tau_{\text{lep}}\tau_{\text{lep}}$	$\tau_{\text{lep}}\tau_{\text{had}}$	$\tau_{\text{had}}\tau_{\text{had}}$
$m_{\tau\tau}^{\text{MMC}}$	•	•	•	•	•	•
$\Delta R(\tau, \tau)$	•	•	•		•	•
$\Delta\eta(j_1, j_2)$	•	•	•			
m_{j_1, j_2}	•	•	•			
$\eta_{j_1} \times \eta_{j_2}$		•	•			
p_{τ}^{total}		•	•			
sum p_{τ}					•	•
$p_{\tau}(\tau_1)/p_{\tau}(\tau_2)$					•	•
$E_{\tau}^{\text{miss}} \phi$ centrality		•	•	•	•	•
$x_{\tau 1}$ and $x_{\tau 2}$						•
$m_{\tau\tau, j_1}$				•		
m_{ℓ_1, ℓ_2}				•		
$\Delta\phi_{\ell_1, \ell_2}$				•		
sphericity				•		
$P_{\tau}^{\ell_1}$				•		
$P_{\tau}^{j_1}$				•		
$E_{\tau}^{\text{miss}}/P_{\tau}^{\ell_2}$				•		
m_{τ}		•			•	
$\min(\Delta\eta_{\ell_1, \ell_2, \text{jets}})$	•					
$j_3 \eta$ centrality	•					
$\ell_1 \times \ell_2 \eta$ centrality	•					
$\ell \eta$ centrality		•				
$\tau_{1,2} \eta$ centrality			•			

Table 3: Discriminating variables used for each channel and category. The filled circles identify which variables are used in each decay mode. Note that variables such as $\Delta R(\tau, \tau)$ are defined either between the two leptons, between the lepton and τ_{had} , or between the two τ_{had} candidates, depending on the decay mode.

Deep networks



New tools
let us
train
deep
networks.

How well
do they work?

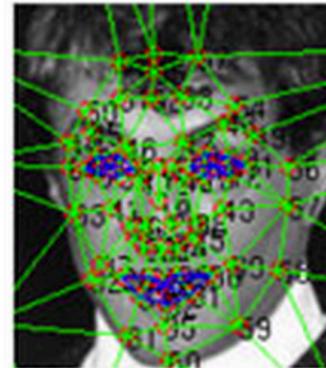
Real world applications



(a)



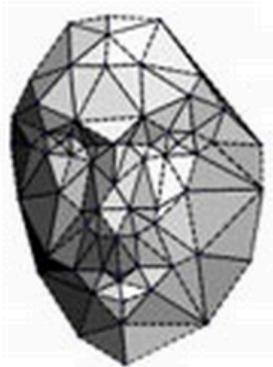
(b)



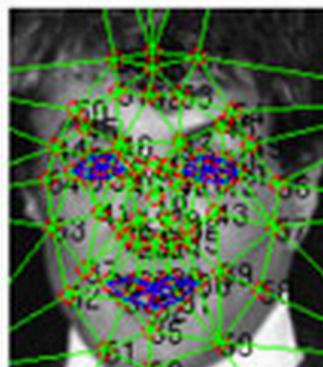
(c)



(d)



(e)



(f)



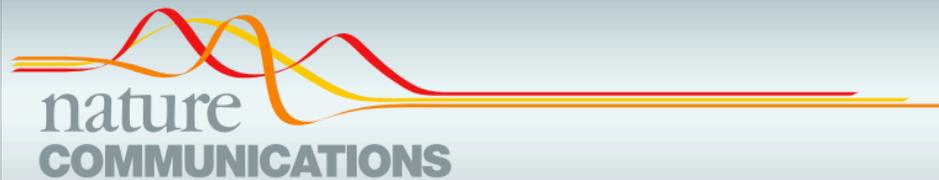
(g)



(h)

Head turn: DeepFace uses a 3-D model to rotate faces, virtually, so that they face the camera. Image (a) shows the original image, and (g) shows the final, corrected version.

Paper



ARTICLE

Received 19 Feb 2014 | Accepted 4 Jun 2014 | Published 2 Jul 2014

DOI: [10.1038/ncomms5308](https://doi.org/10.1038/ncomms5308)

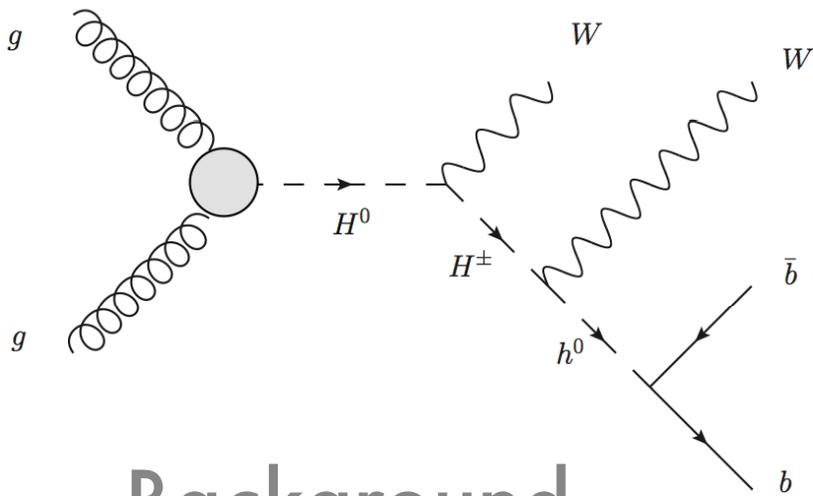
Searching for exotic particles in high-energy physics with deep learning

P. Baldi¹, P. Sadowski¹ & D. Whiteson²

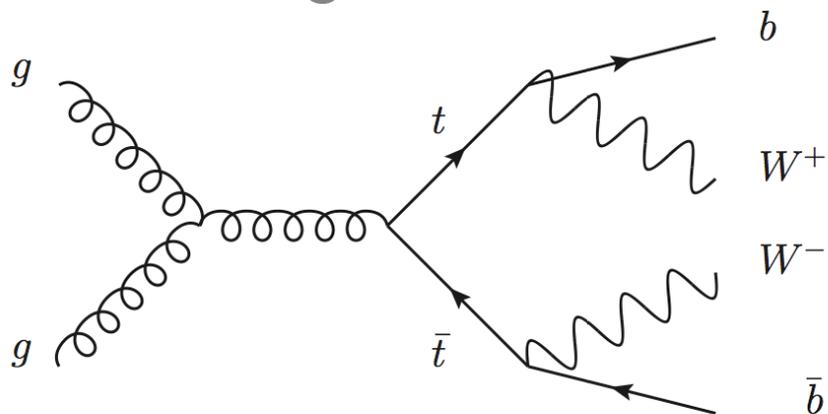
arXiv: 1402.4735

Benchmark problem

Signal



Background



Can deep networks automatically discover useful variables?

4-vector inputs

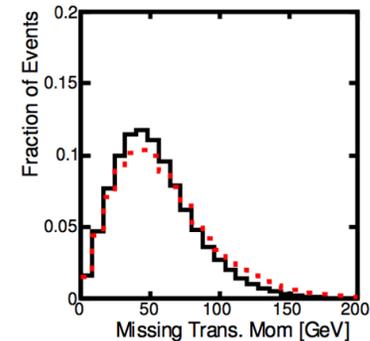
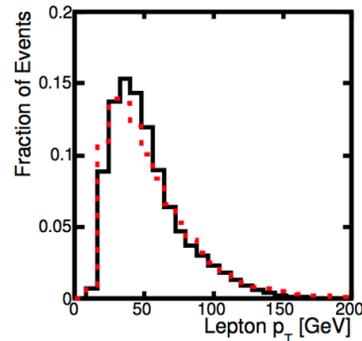
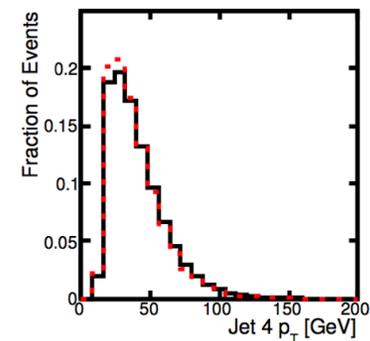
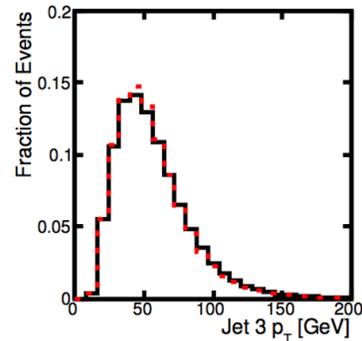
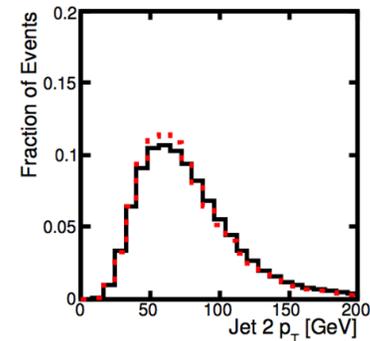
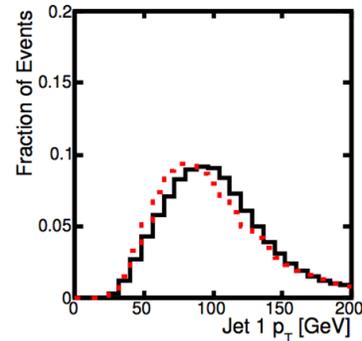
21 Low-level vars

jet+lepton mom. (3x5)

missing ET (2)

jet btags (4)

Not much
separation
visible in 1D
projections



4-vector inputs

7 High-level vars

$m(WWbb)$

$m(Wbb)$

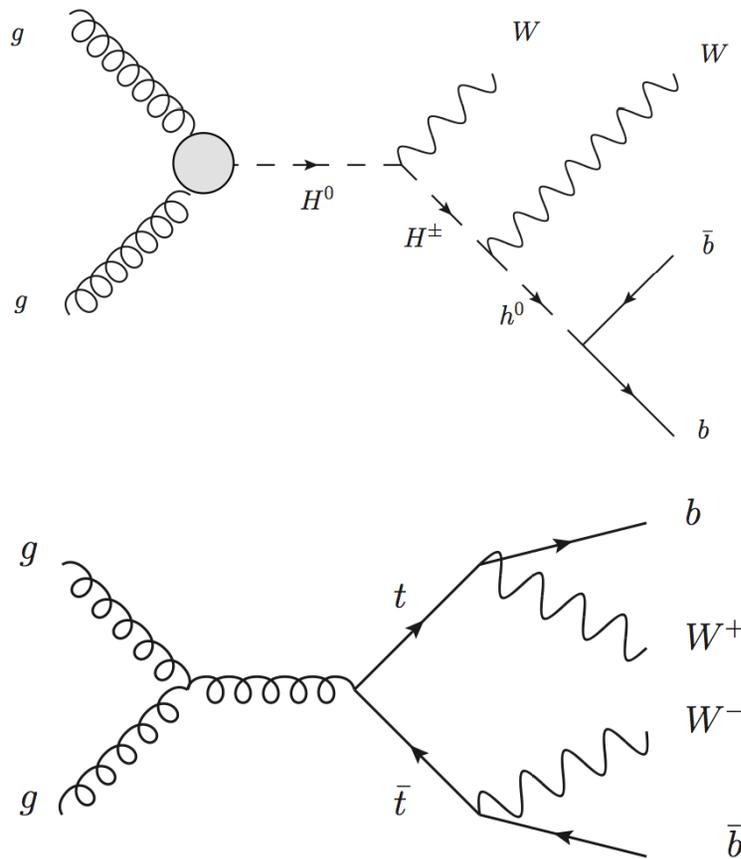
$m(bb)$

$m(bjj)$

$m(jj)$

$m(lv)$

$m(blv)$



4-vector inputs

7 High-level vars

$m(WWbb)$

$m(Wbb)$

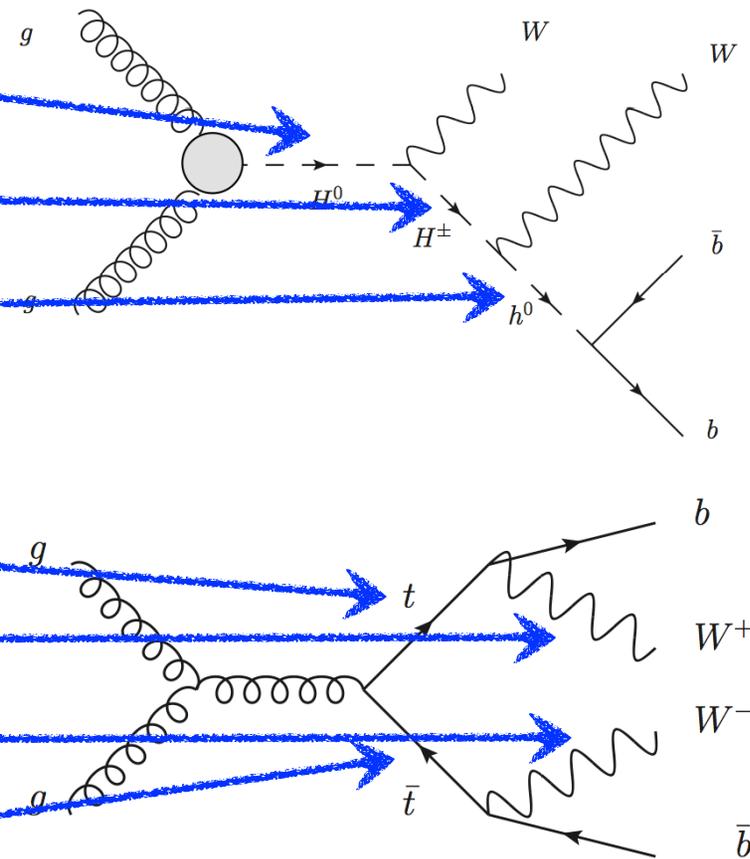
$m(bb)$

$m(bjj)$

$m(jj)$

$m(lv)$

$m(blv)$



4-vector inputs

7 High-level vars

$m(WWbb)$

$m(Wbb)$

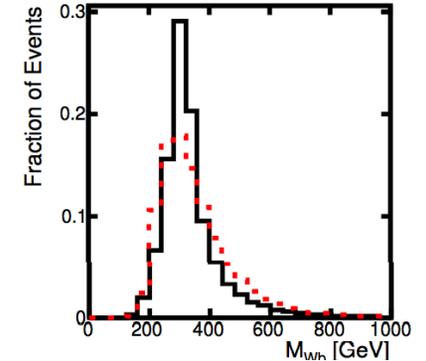
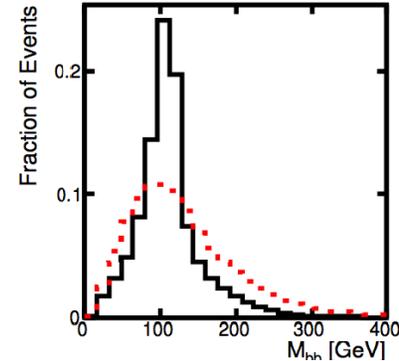
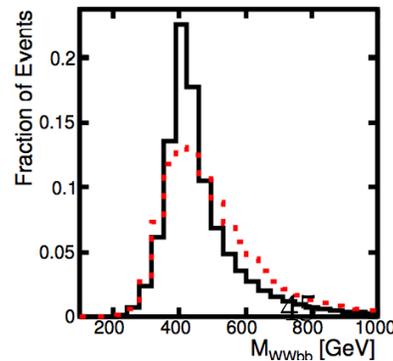
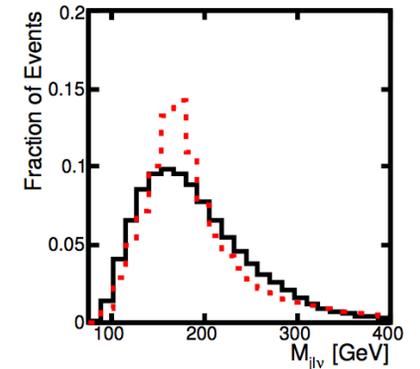
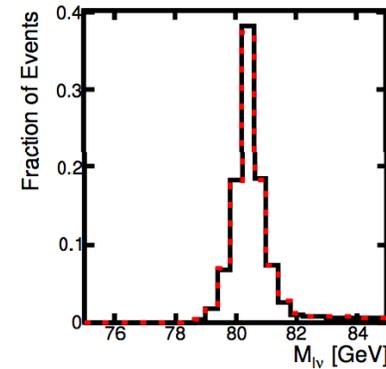
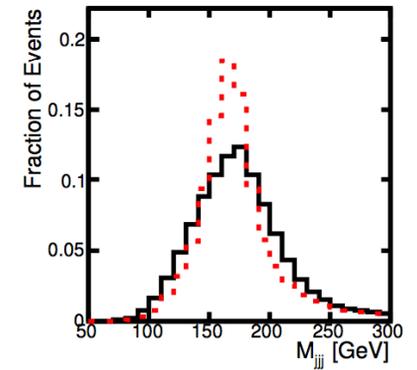
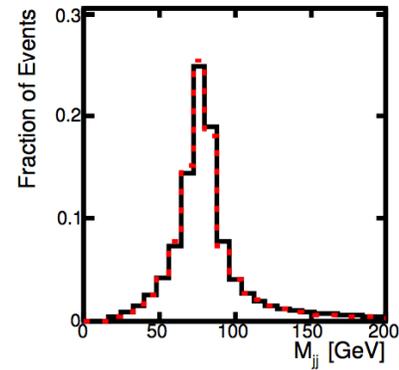
$m(bb)$

$m(bjj)$

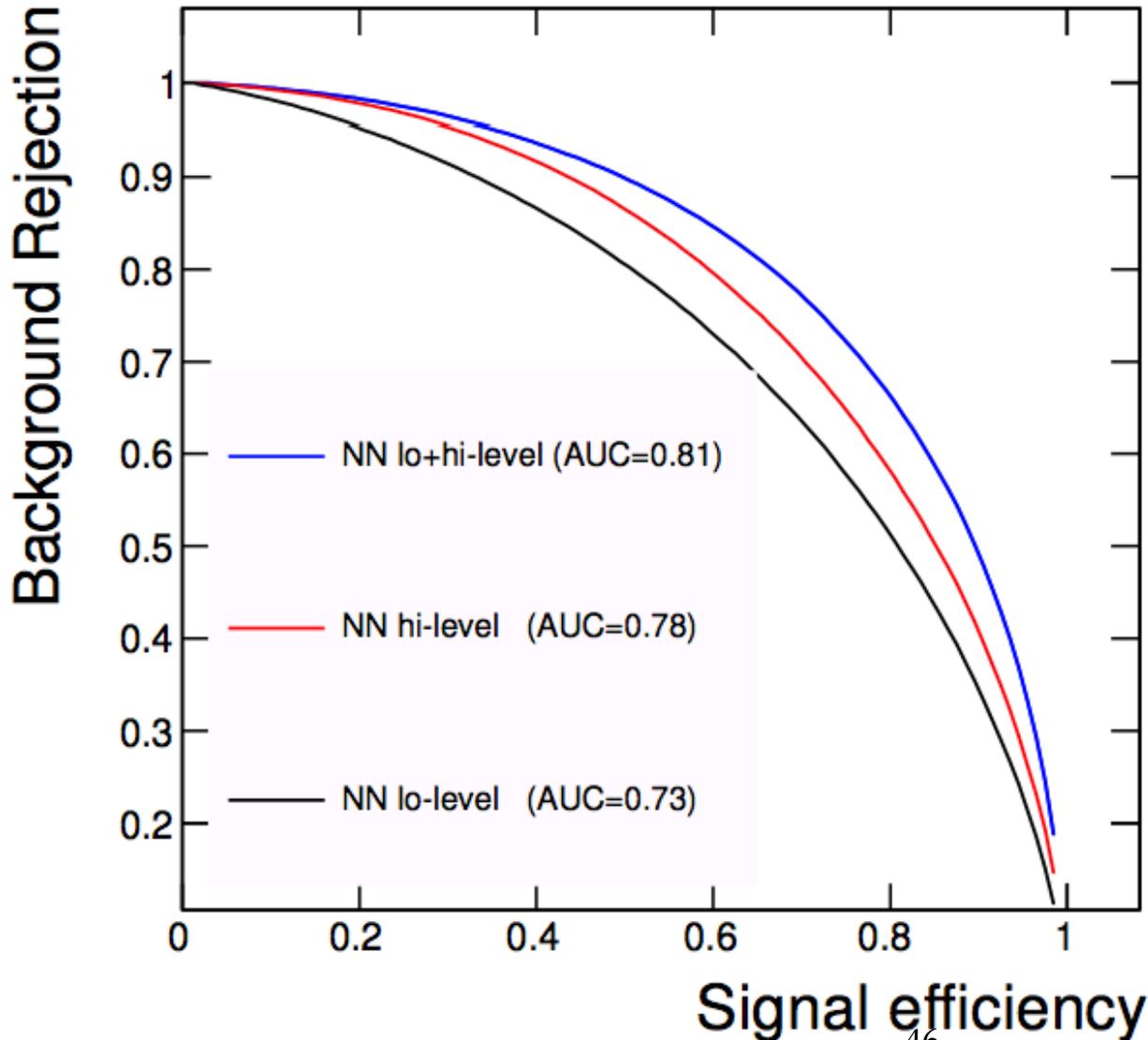
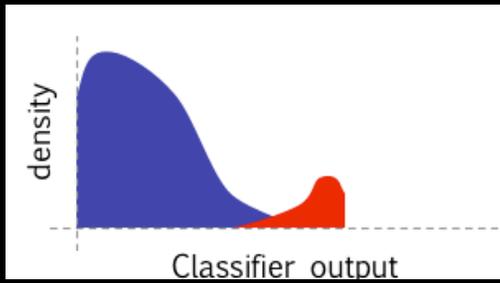
$m(ijj)$

$m(lv)$

$m(blv)$



Standard NNs



Results

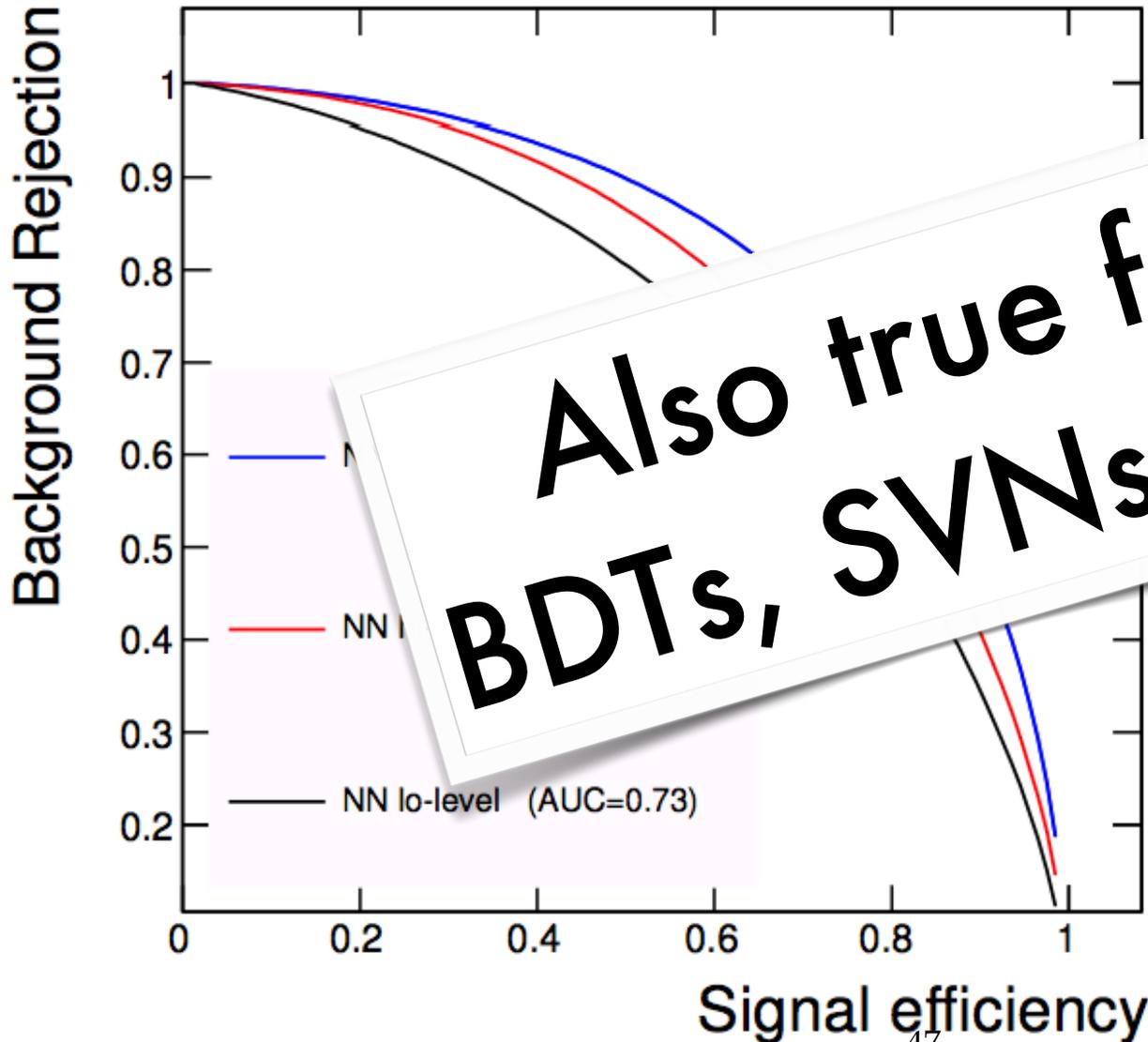
Adding hi-level
boosts performance
Better: lo+hi-level.

Conclude:

NN can't find
hi-level vars.

Hi-level vars
do not have all info

Standard NNs



Results

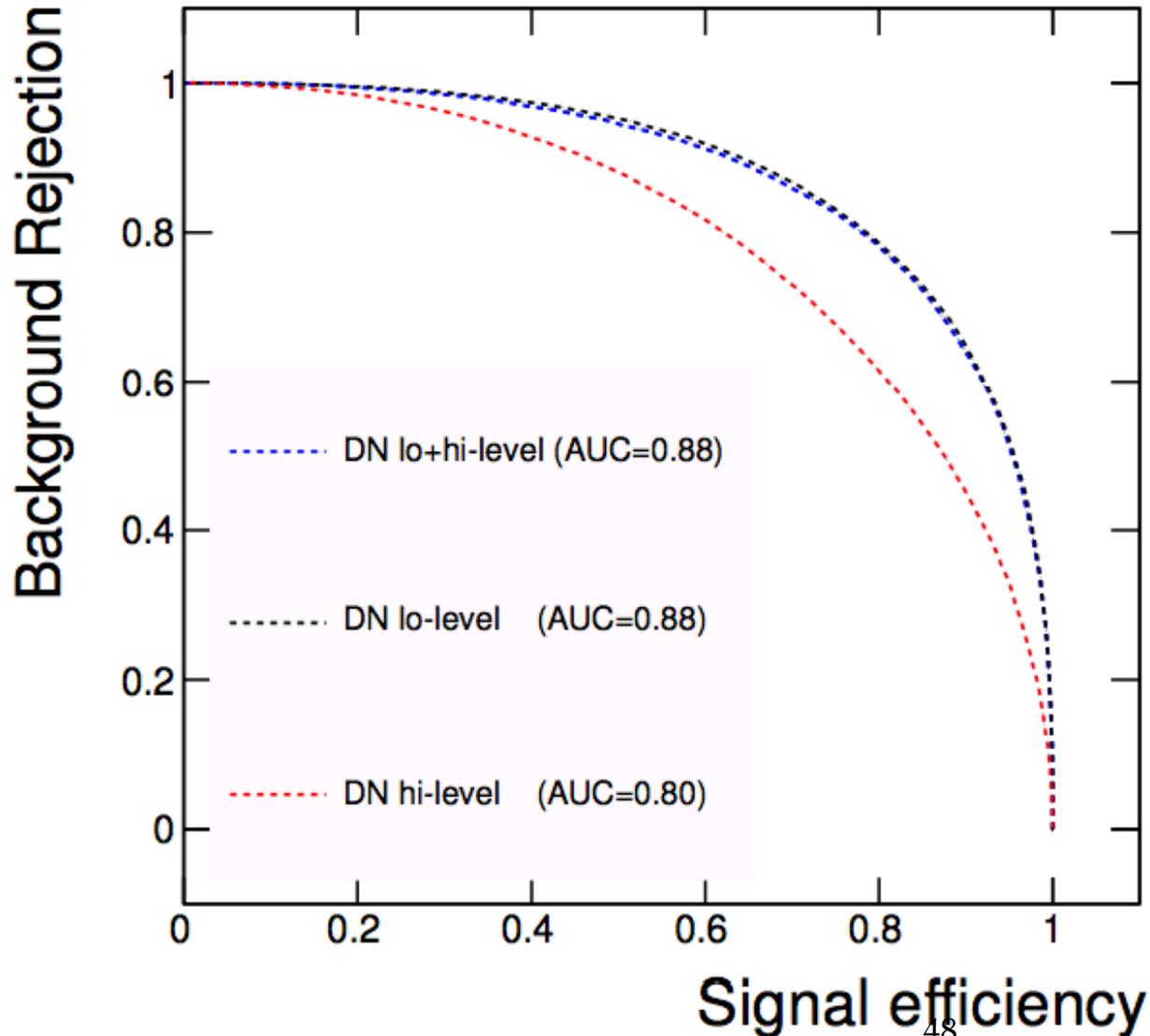
...ing hi-level
... performance
...lo+hi-level.

Conclude:

NN can't find
hi-level vars.

Hi-level vars
do not have all info

Deep Networks



Results

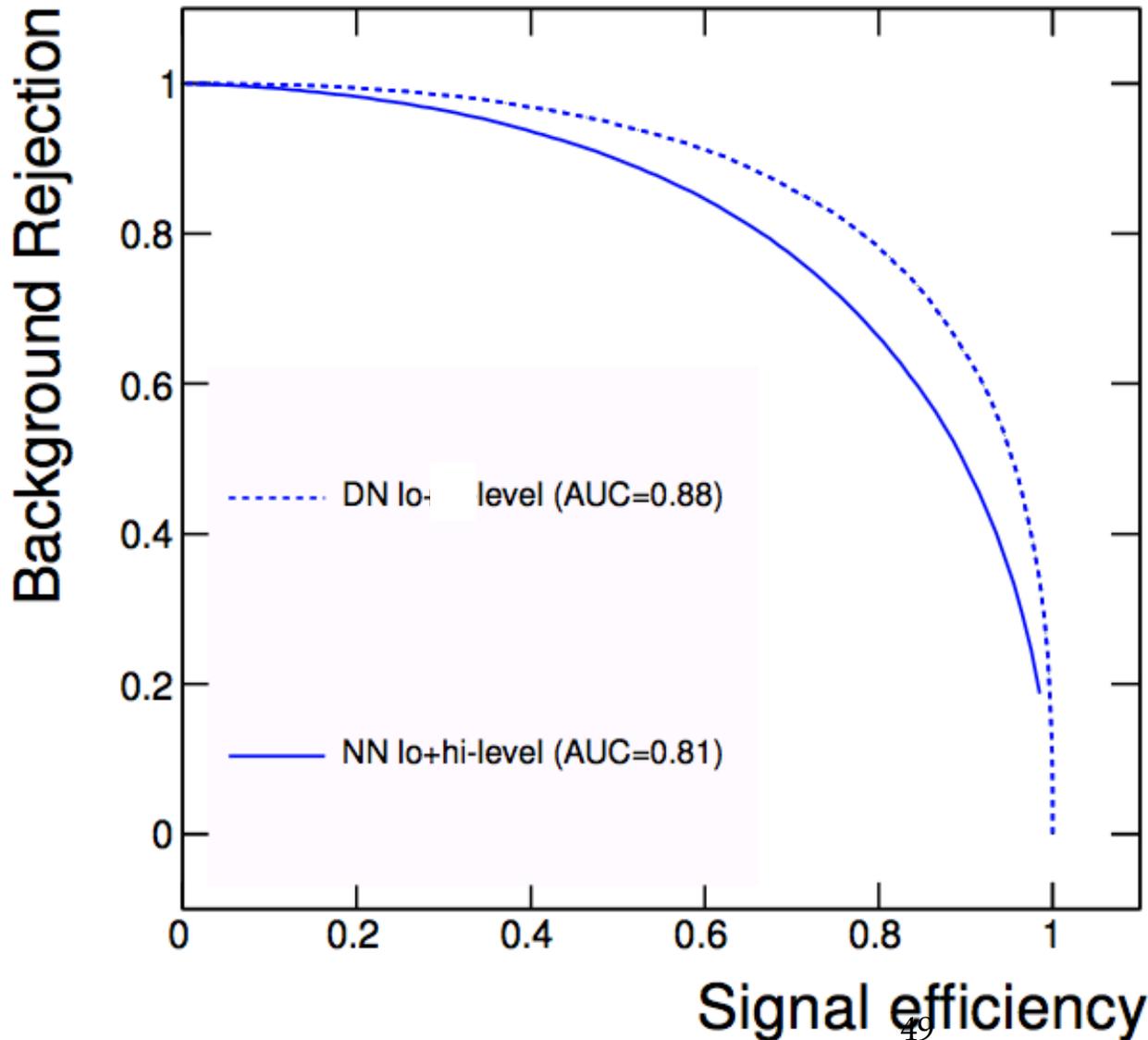
Lo+hi = lo.

Conclude:

DN can find
hi-level vars.

Hi-level vars
do not have all info
are unnecessary

Deep Networks



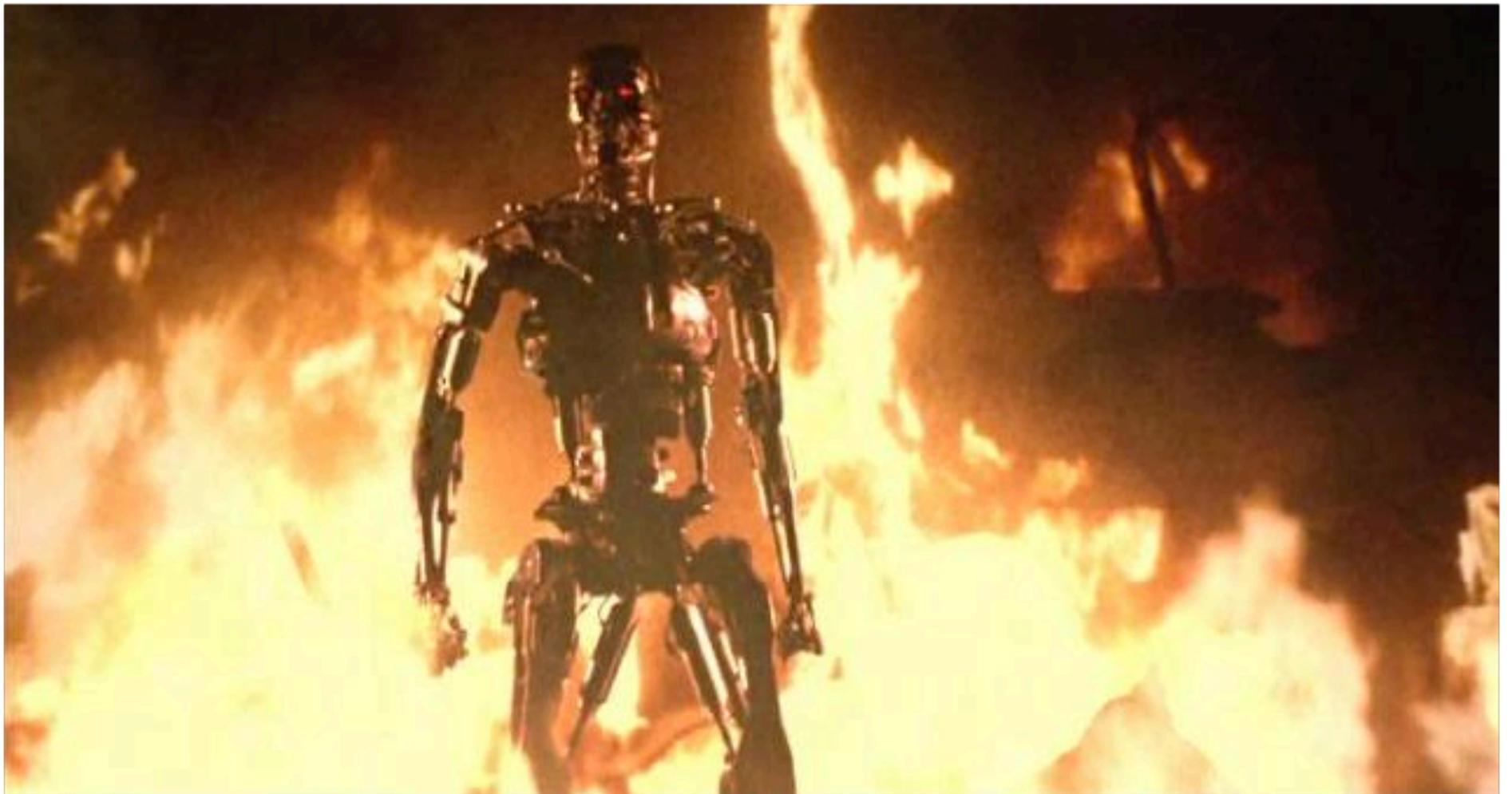
Results

DN > NN

Conclude:

DN does better than human assisted NN

The AIs win



Results

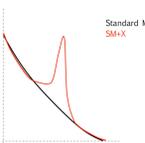
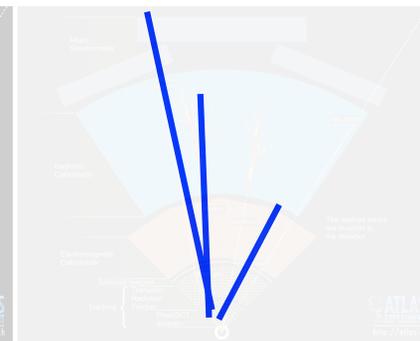
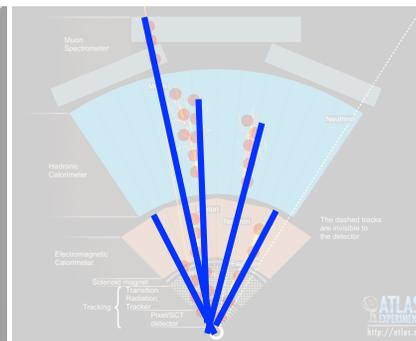
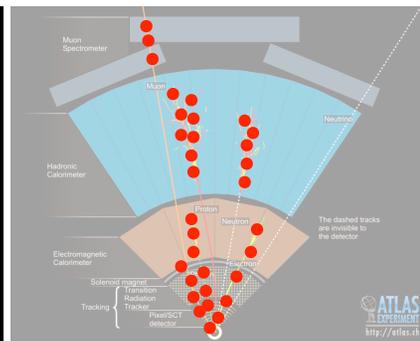
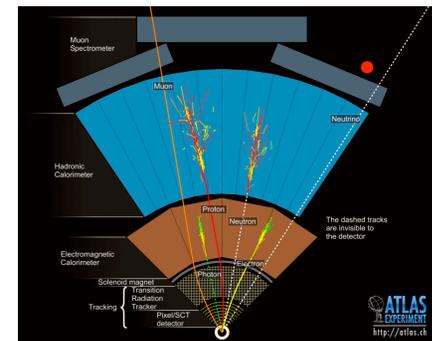
Identified example benchmark where traditional NNs fail to discover all discrimination power.

Adding human insight helps traditional NNs.

Deep networks succeed **without human insight**.
Outperform human-boosted traditional NNs.

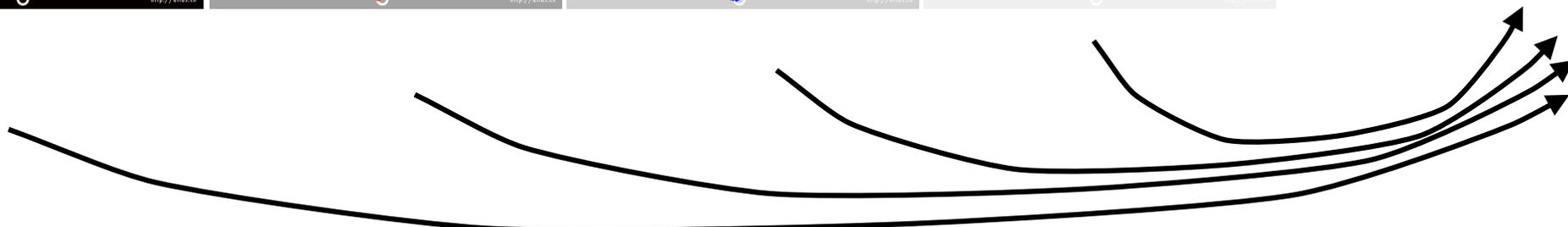
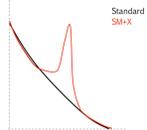
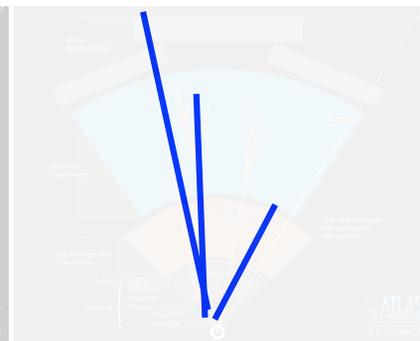
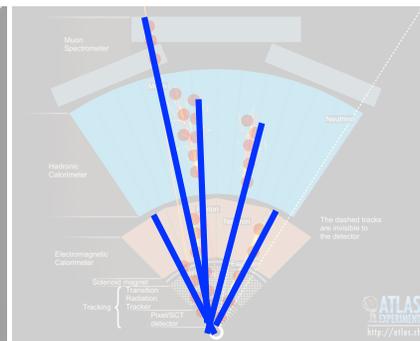
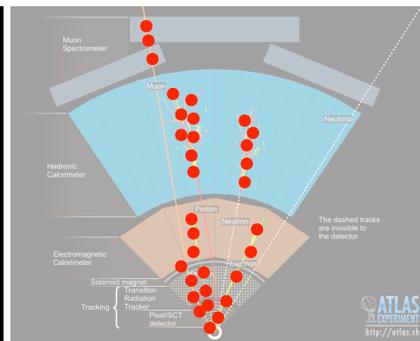
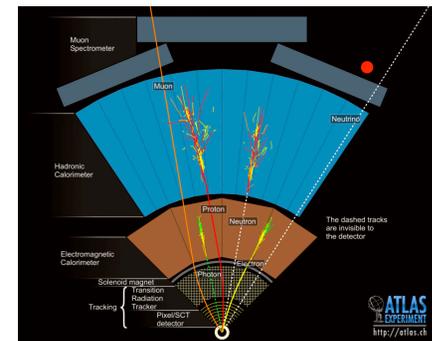
What is possible?

Raw	Sparsified	Reco	Select	Ana
$1e7$	$1e3$	100	50	1



What is possible?

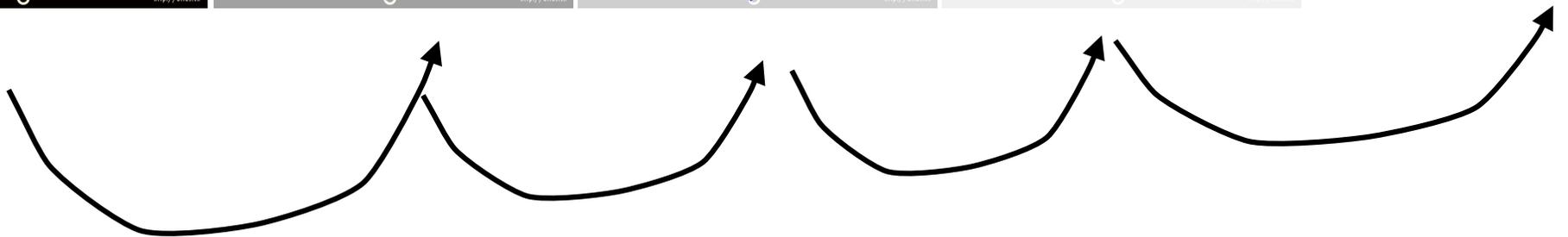
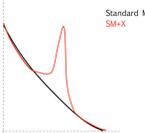
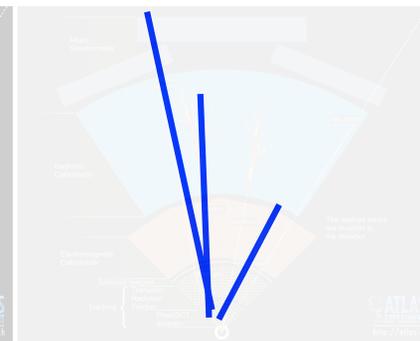
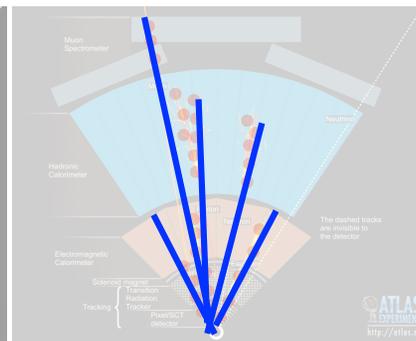
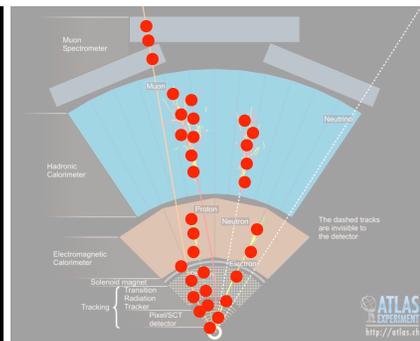
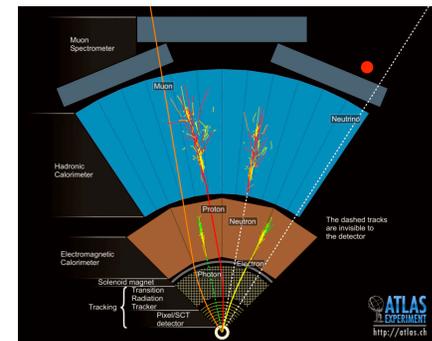
Raw	Sparsified	Reco	Select	Ana
$1e7$	$1e3$	100	50	1



Skip more steps with ML?

Or this?

Raw	Sparsified	Reco	Select	Ana
$1e7$	$1e3$	100	50	1



Improve each step with ML?

Jets

Jet Substructure Classification in High-Energy Physics with Deep Neural Networks

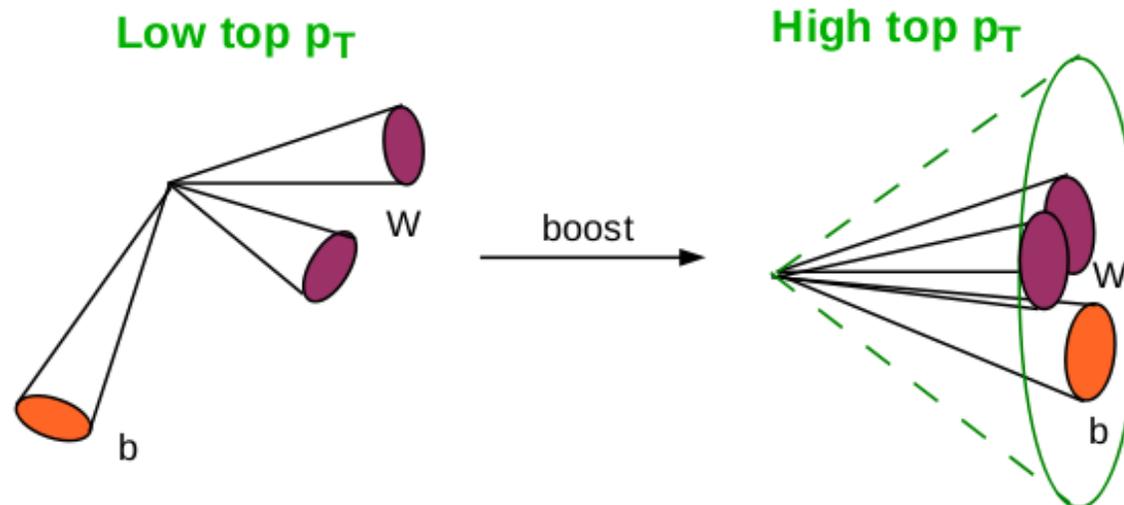
Pierre Baldi,¹ Kevin Bauer,² Clara Eng,³ Peter Sadowski,¹ and Daniel Whiteson²

¹*Department of Computer Science, University of California, Irvine, CA 92697*

²*Department of Physics and Astronomy, University of California, Irvine, CA 92697*

³*Department of Chemical and Biomolecular Engineering, University of California, Berkeley CA 94270*

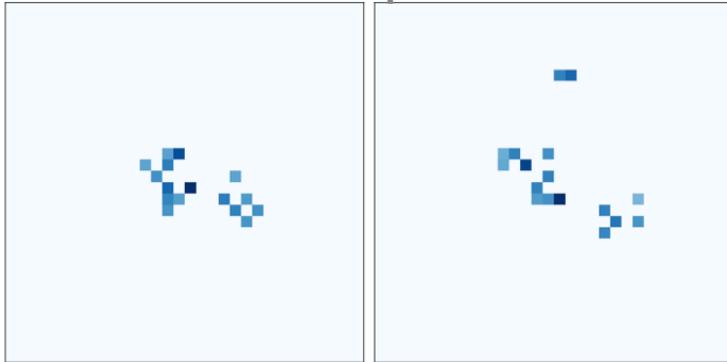
(Dated: April 12, 2016)



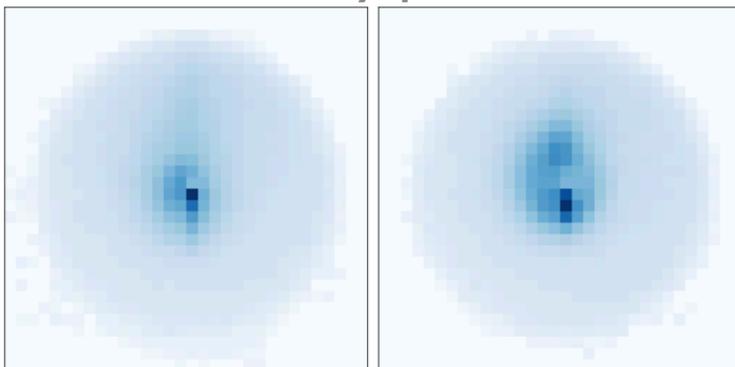
Jet substructure

LL variables

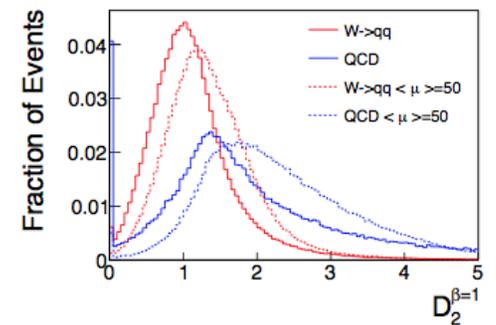
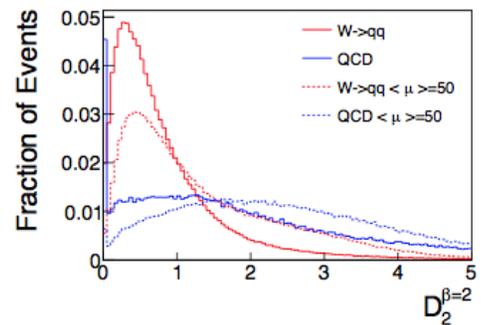
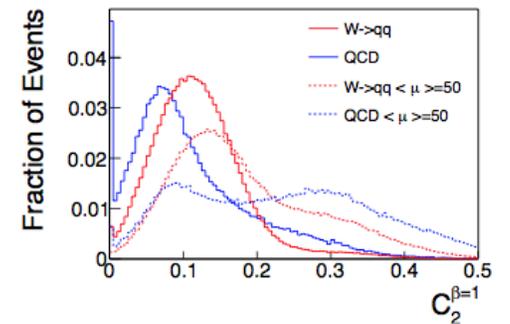
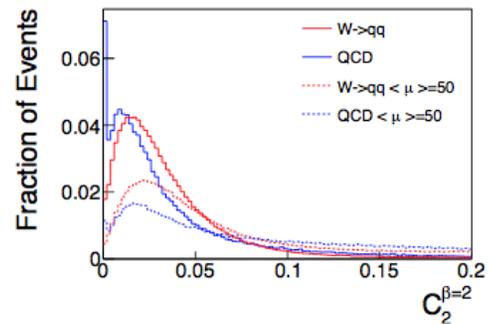
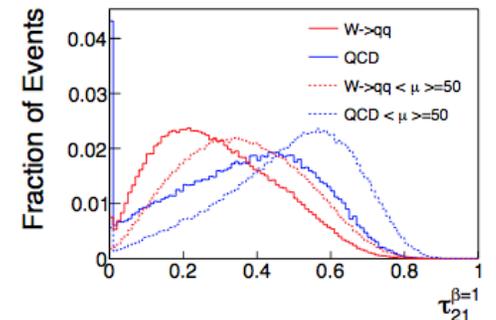
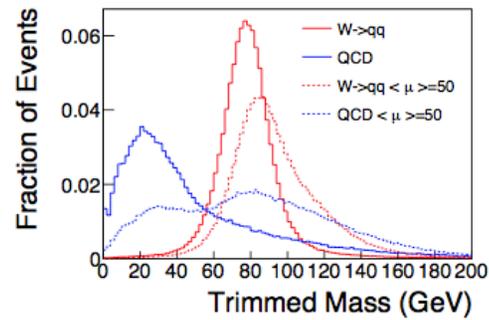
one jet



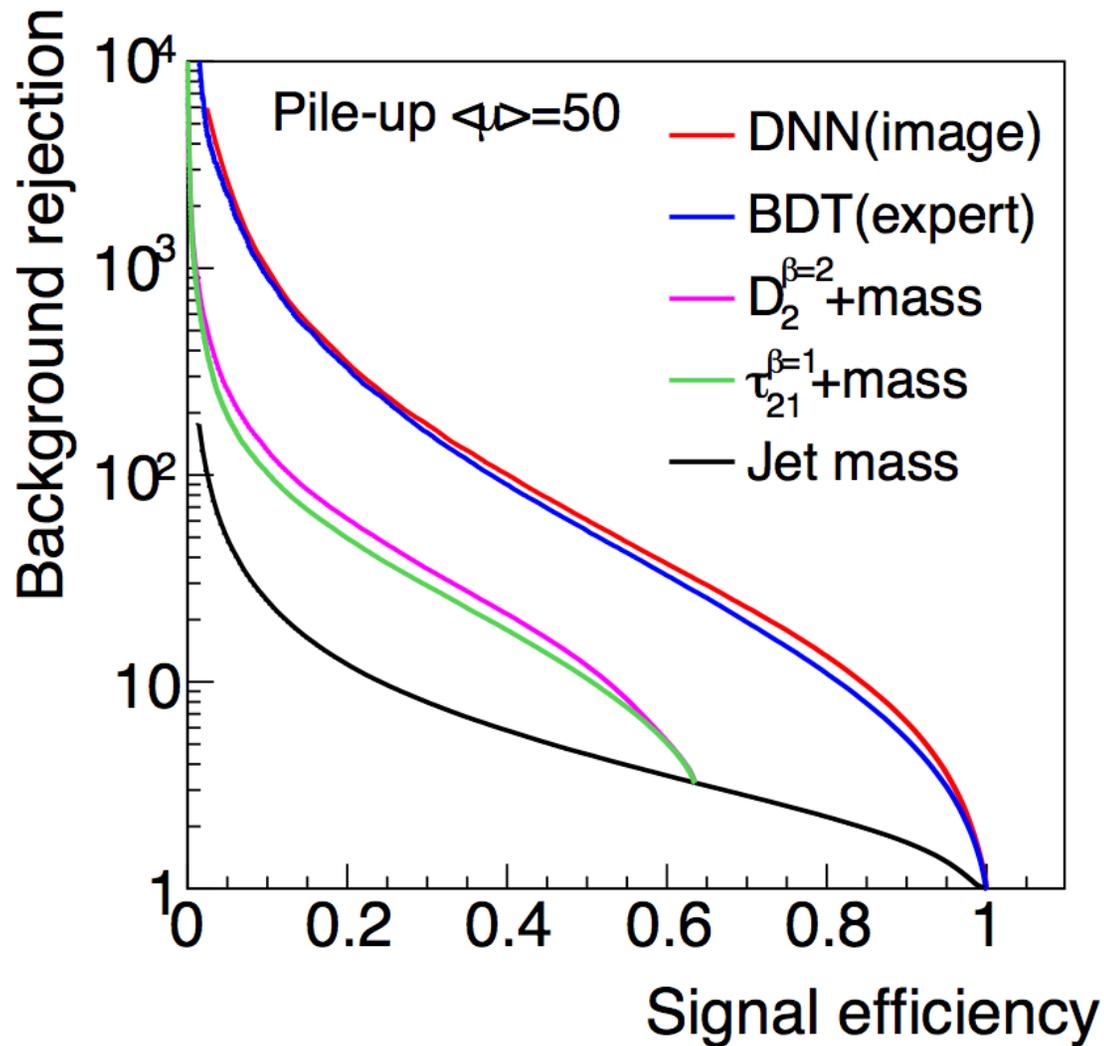
many jets



HL variables



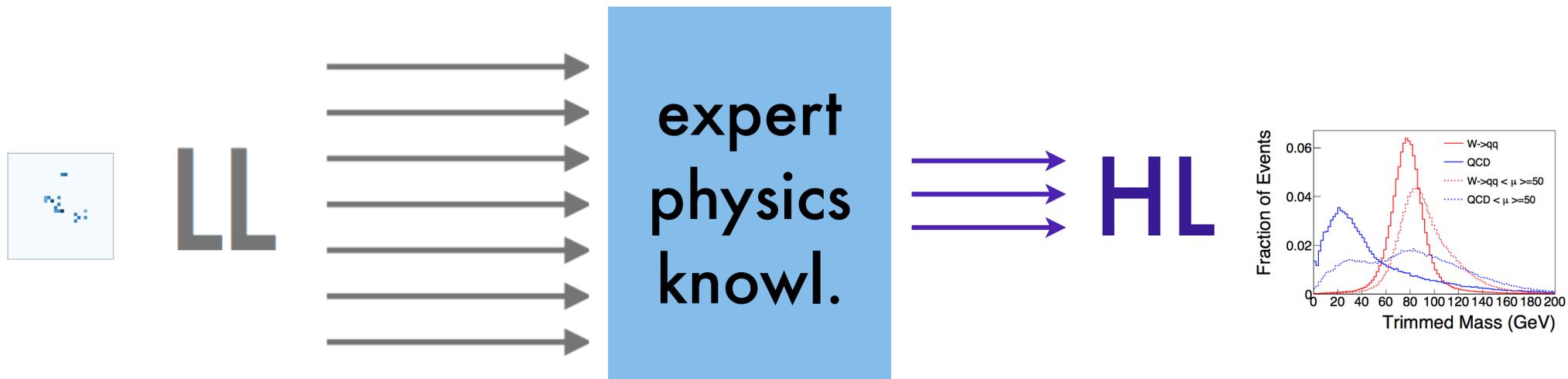
Jet tagging



How?

What is it doing?

Our low-level (LL) data are often high-dim



We can calculate likelihood ratios in the low-dim HL space often using MC techniques

But HL doesn't always capture the information

Yet we prefer HL

If HL data includes all necessary information...

- It is easier to understand
- Its modeling can be verified
- Uncertainties can be sensibly defined
- It is more compact and efficient
- LL -> HL is physics, so we like it.

Our question

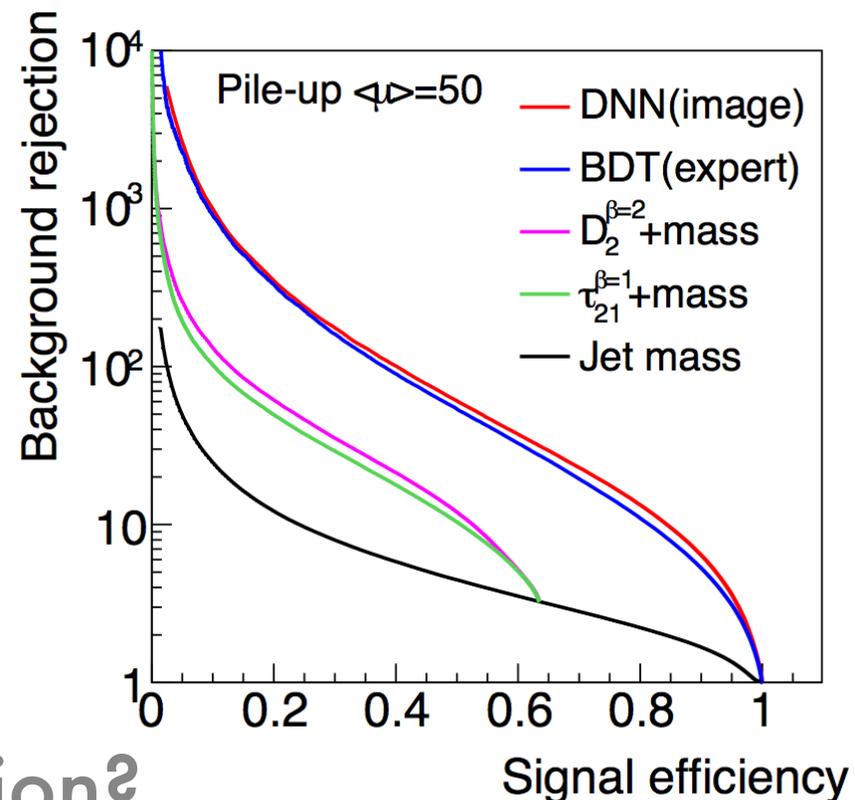
How has the DNN found its solution?
What can we learn from it?

Residual knowledge:

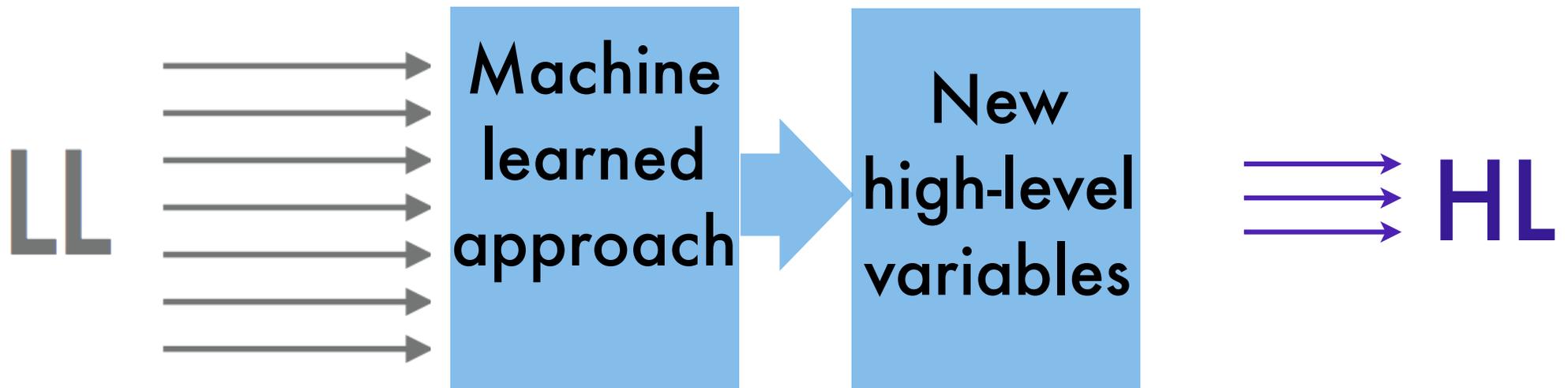
Is there a **new** HL variable?
Can it reveal physics?

Translating complete solutions:

What is the **structure** of its solution?
Has it just rediscovered and
optimized the existing HL vars?



Learning from ML



Use LL analysis as a probe, not a final product.

How?

I. Define space of possible human solutions

- provides context for NN solution
- defines problem
- does NN live in this space?
- Can it be compactly represented?
- Yes or No are both interesting!



$$= \sum_a \sum_b \sum_c z_a z_b z_c \theta_{ab} \theta_{ac} \theta_{bc}^2$$

$$z_i = \frac{p_{T_i}}{\sum_j p_{T_j}}$$

How?

I. Define space of possible human solutions

- provides context for NN solution
- defines problem
- does NN live in this space?
- Can it be compactly represented?
- Yes or No are both interesting!

II. Define mapping metric

- how do you compare two solutions?
- can't use functional identity or linear correlation

Discriminant Similarity

Input space

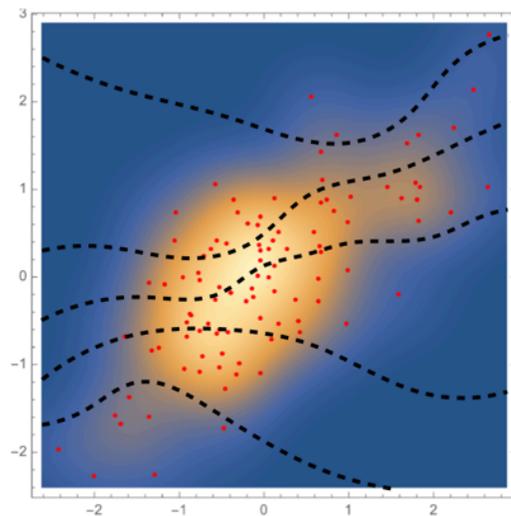
Output

Function sameness

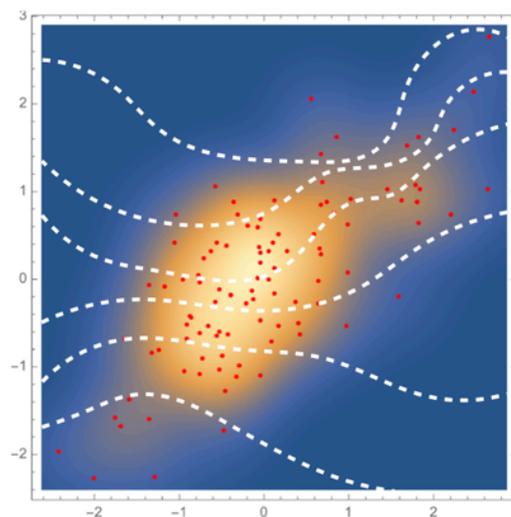
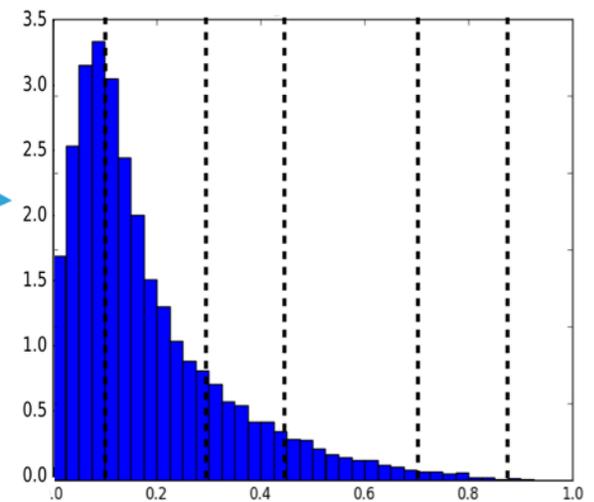
Complete equivalence
not the idea

Any 1:1 transformation
of function has no impact
in our context

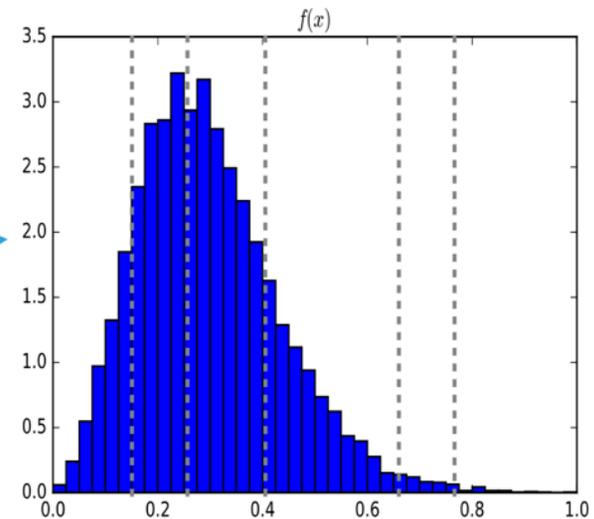
Only care about the
ordering of points
not the actual function
values



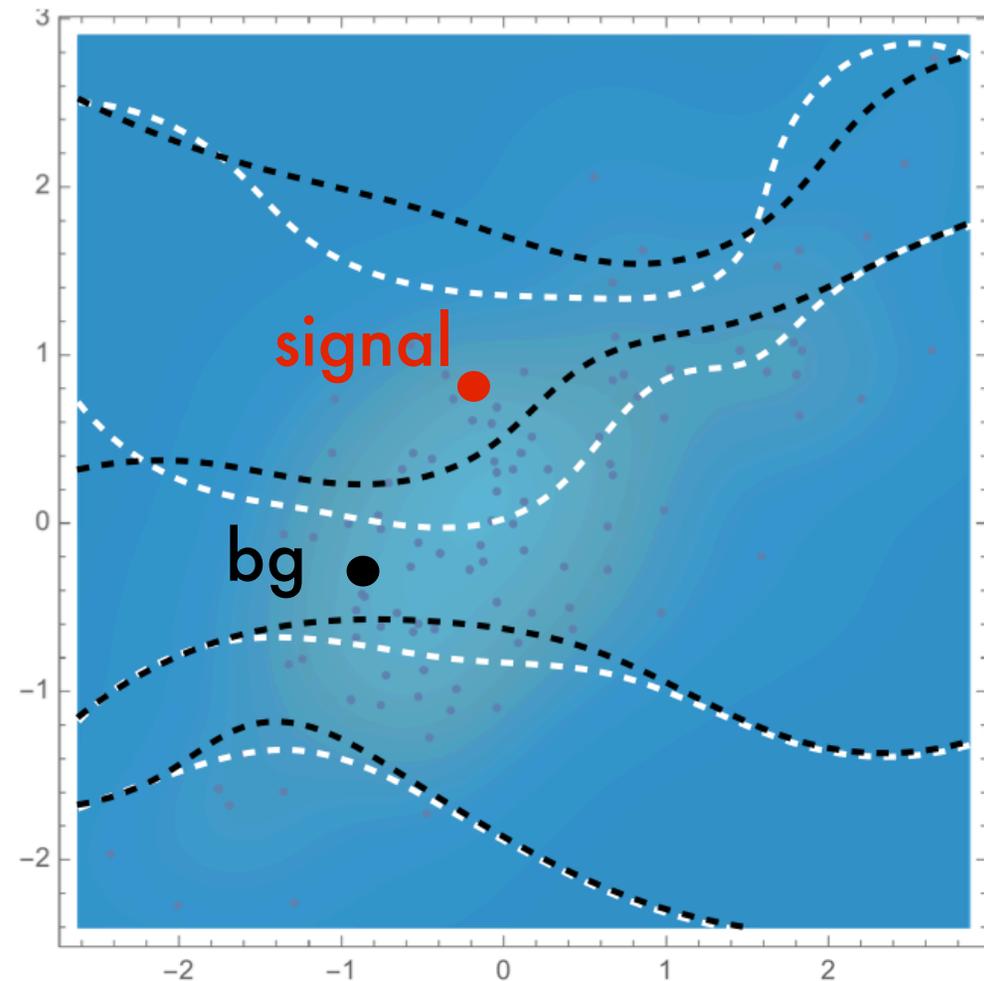
$f(x)$



$g(x)$



Discriminant ordering

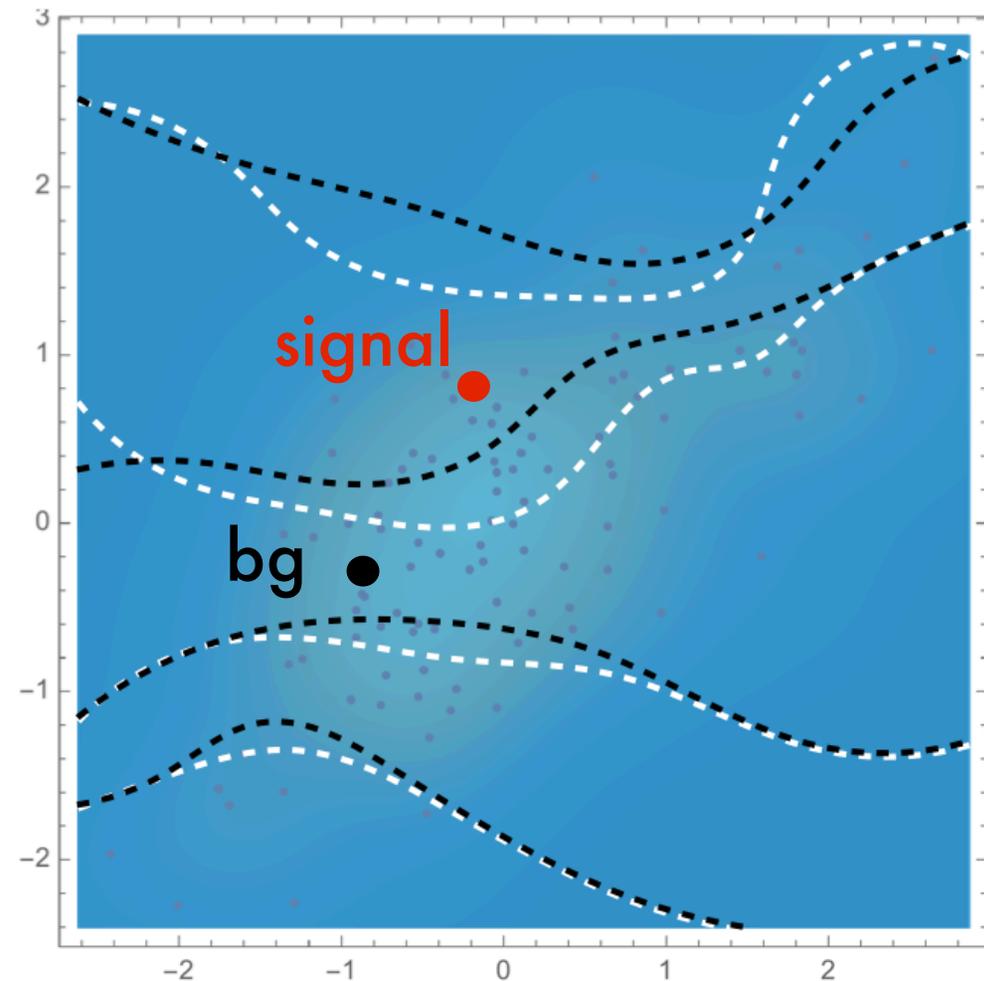


Consider how two functions treat a pair of points

$$f(x_{\text{sig}}) - f(x_{\text{bg}})$$
$$g(x_{\text{sig}}) - g(x_{\text{bg}})$$

Do these have the same sign?

Discriminant ordering



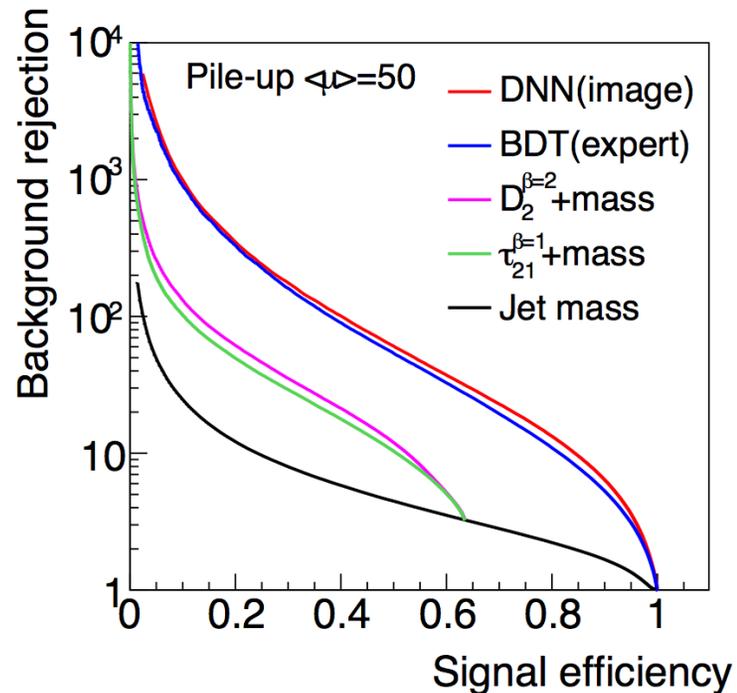
Evaluate how often they give a bg-sig pair the same ordering.

$$\text{DO}(x, x') = \Theta\left(\left(f(x) - f(x')\right)\left(g(x) - g(x')\right)\right)$$

Sample the space.

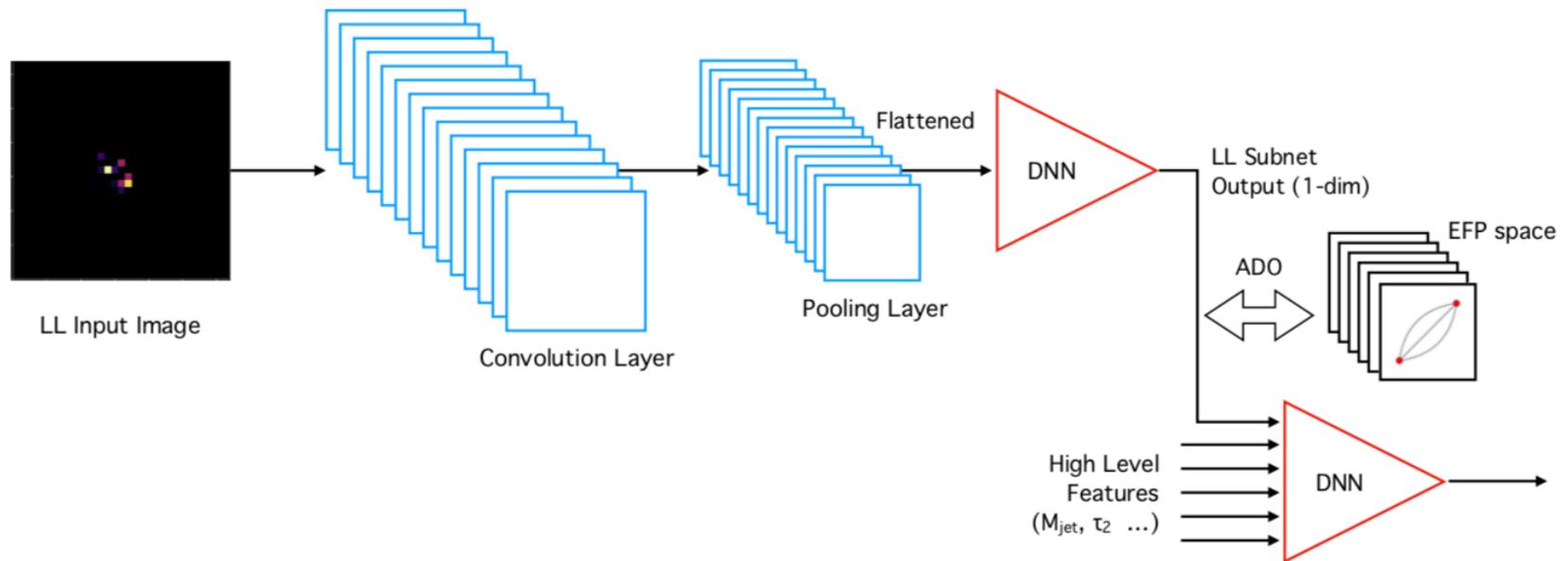
$$\text{ADO} = \int dx dx' p_{\text{sig}}(x) p_{\text{bkg}}(x') \text{DO}(x, x').$$

The problem

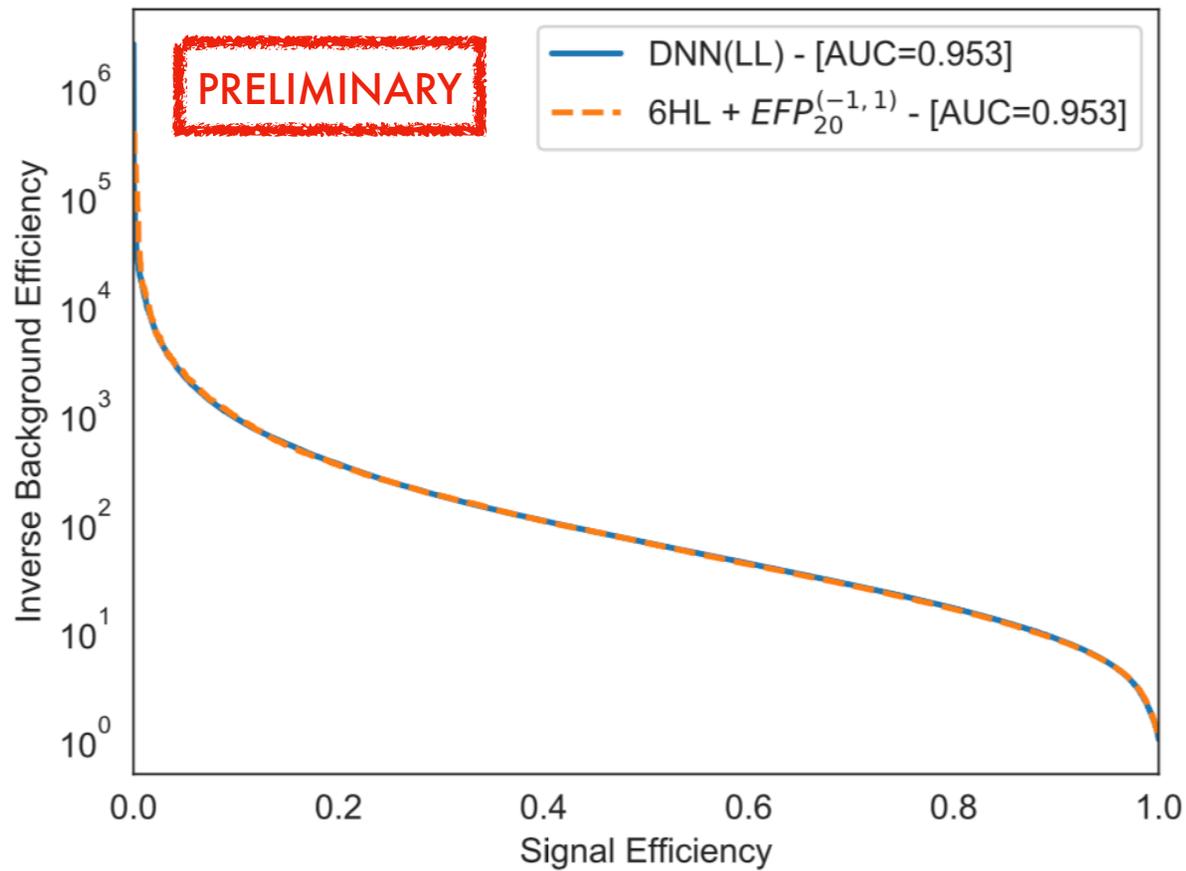


- Two approaches:
- (1) find the gap
 - (2) build from scratch

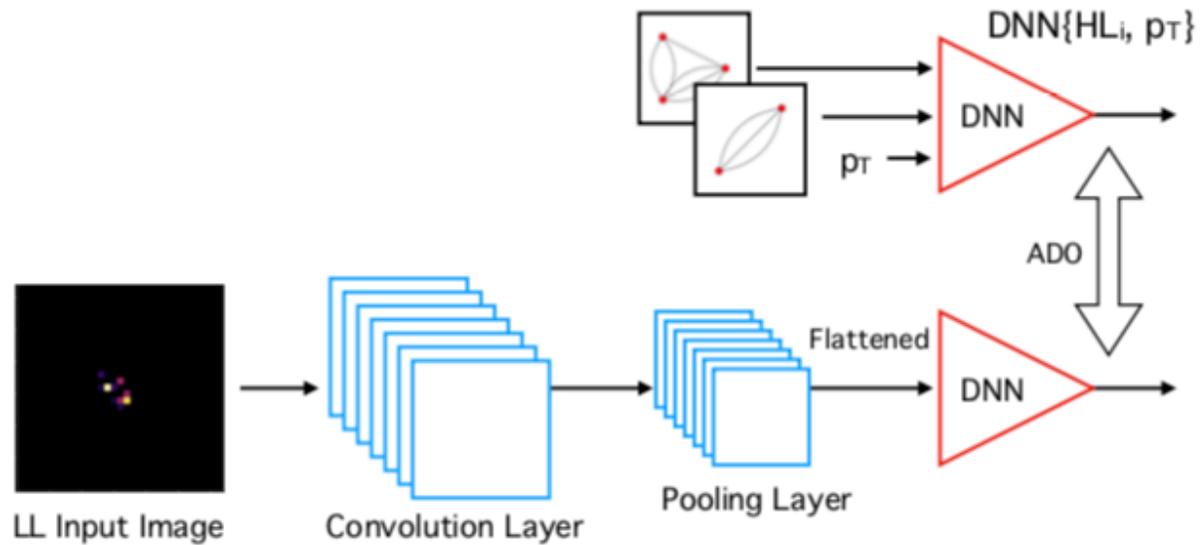
Find the gap



It works!



Build from scratch



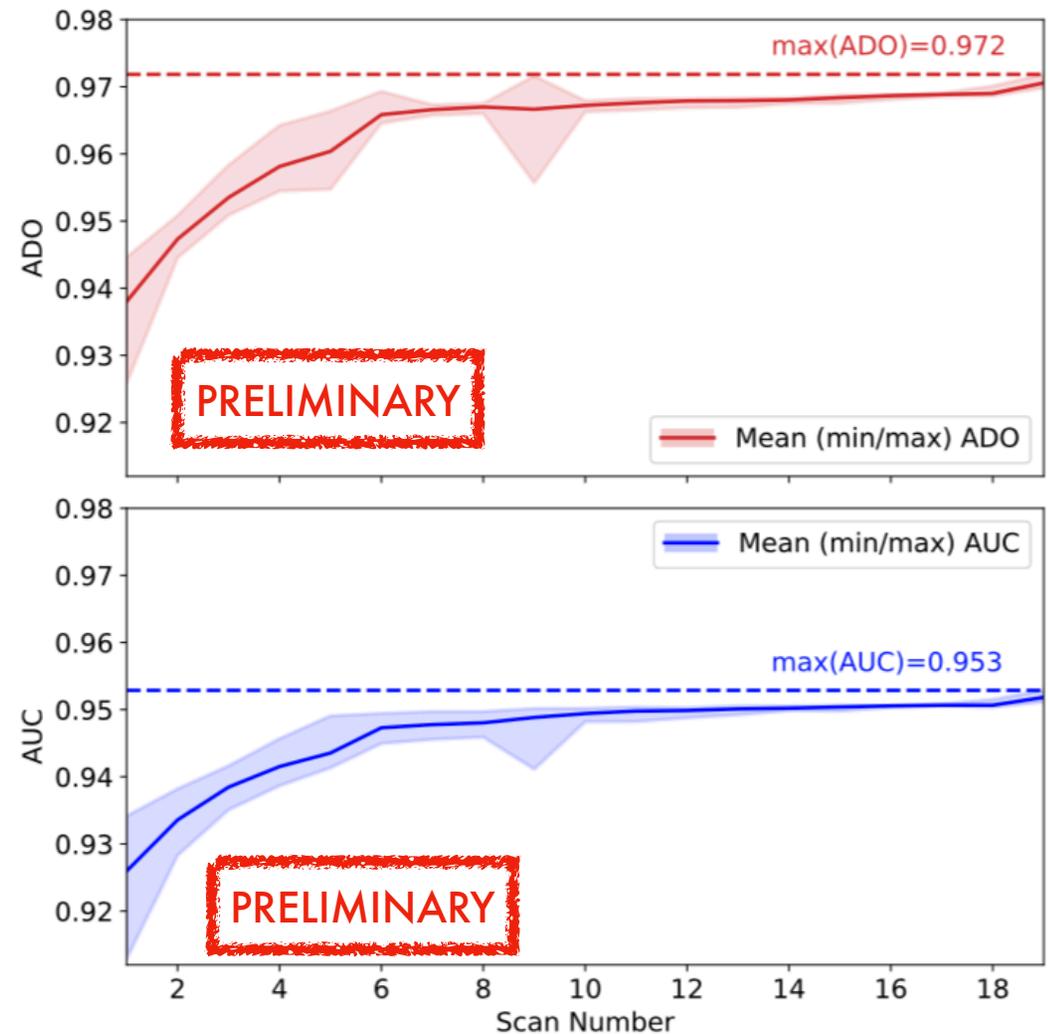
Preliminary

A single point
in this space:

Is very similar
to NN(LL) sol

Captures most of
performance of HL sol

Adding more points
approaches the full solution.

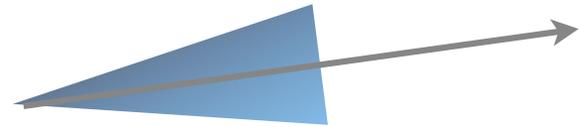


Muon isolation

Isolated muons



Muons from jets



Problem

Jet can be soft, not reconstructed

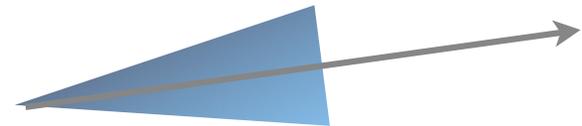
Jets are strongly produced, large background

Muon isolation

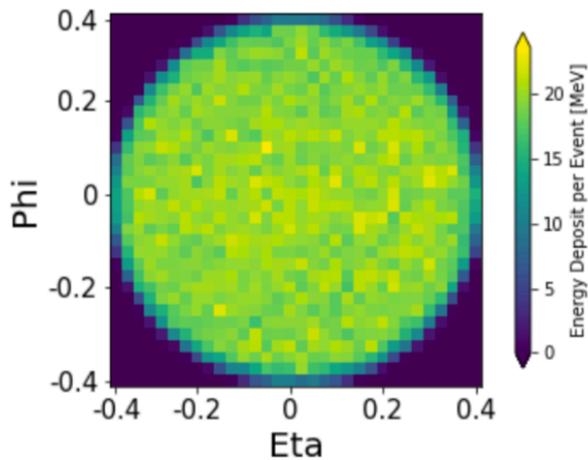
Isolated muons



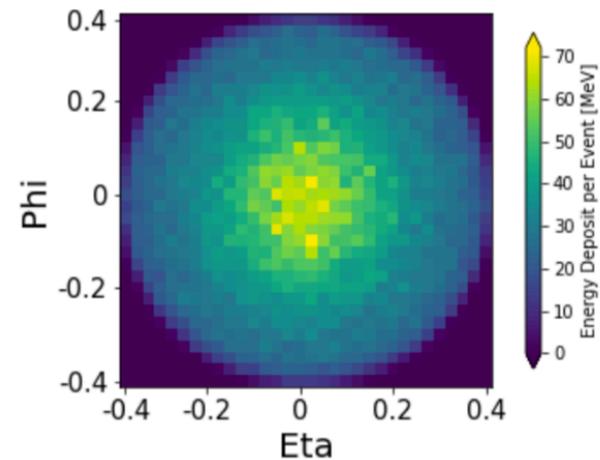
Muons from jets



Little cal deposition



Large cal deposition

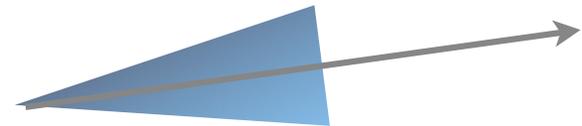


Muon isolation

Isolated muons



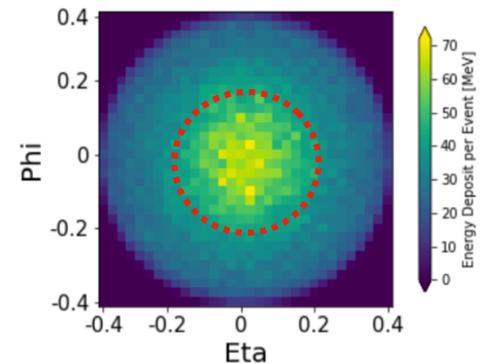
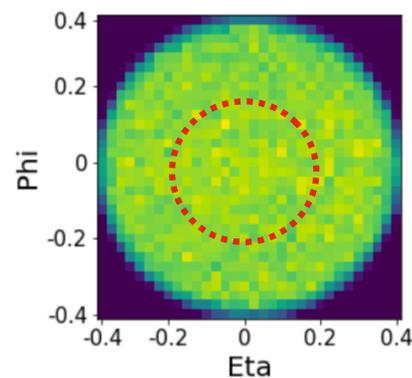
Muons from jets



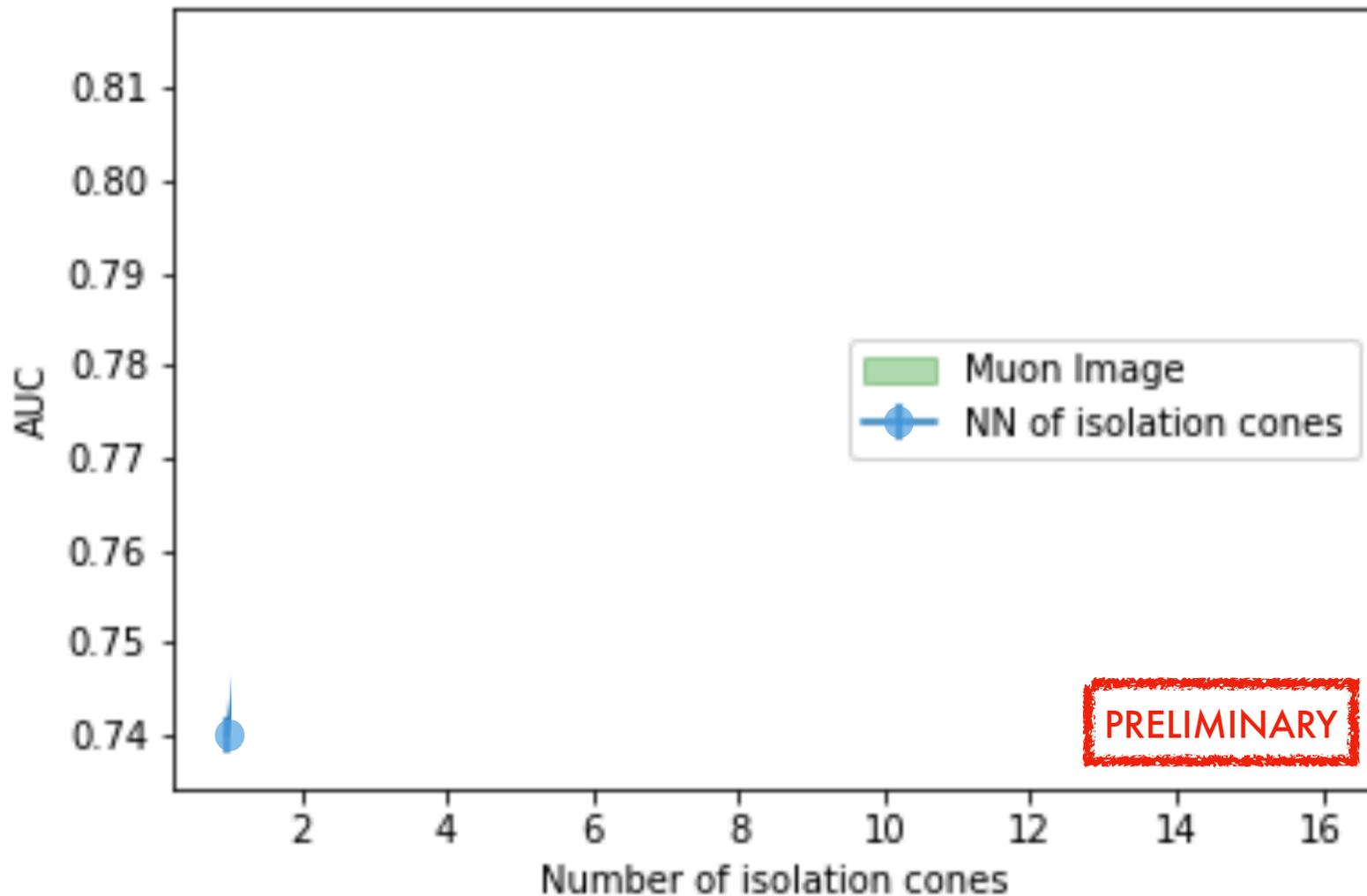
Standard Approach

Calculate "isolation"
Energy in a cone around
muon.

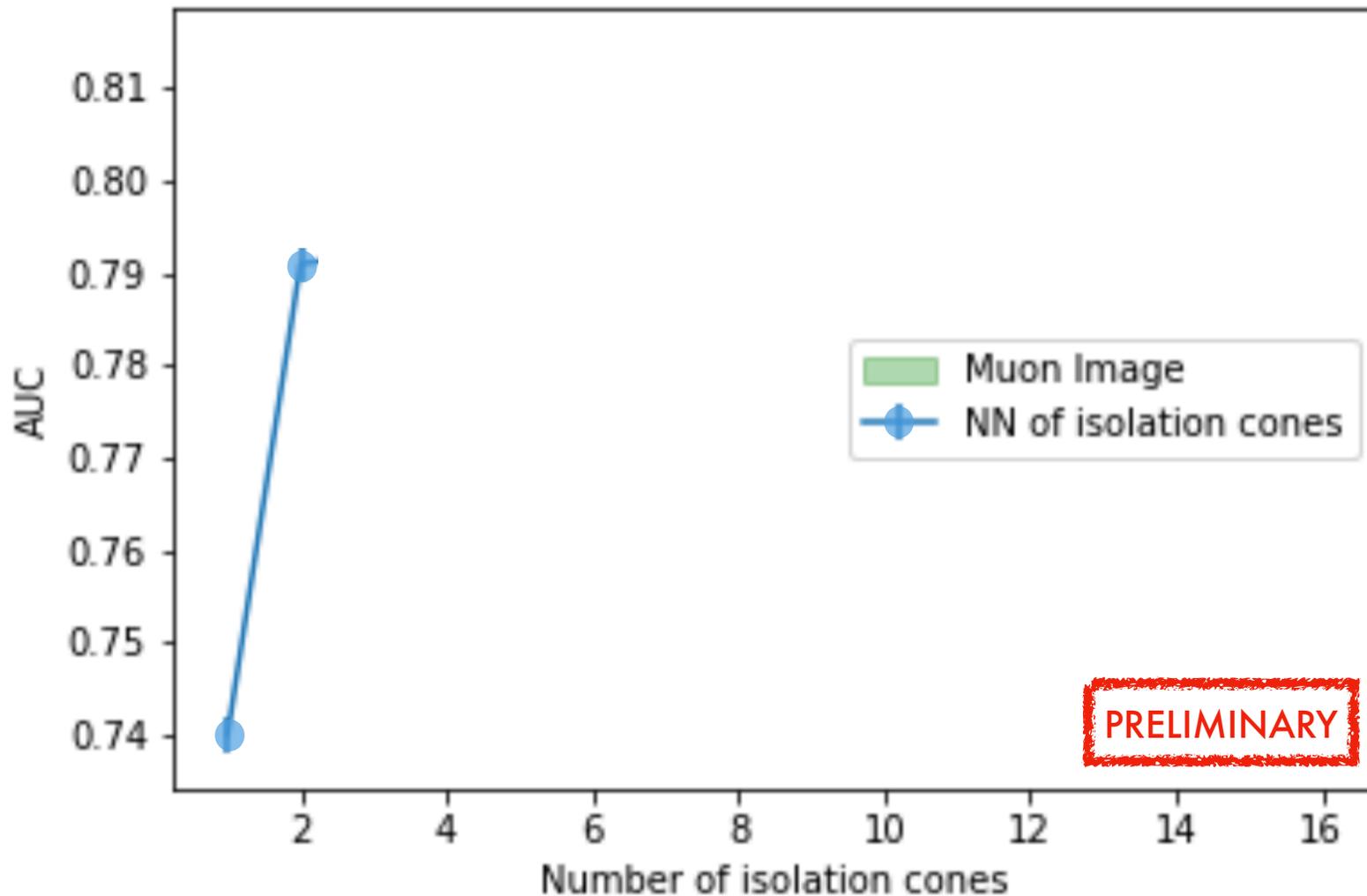
$$I_{\mu}(R_0) = \sum_{R < R_0} \frac{p_T^{\text{calo}}}{p_T^{\text{muon}}}$$



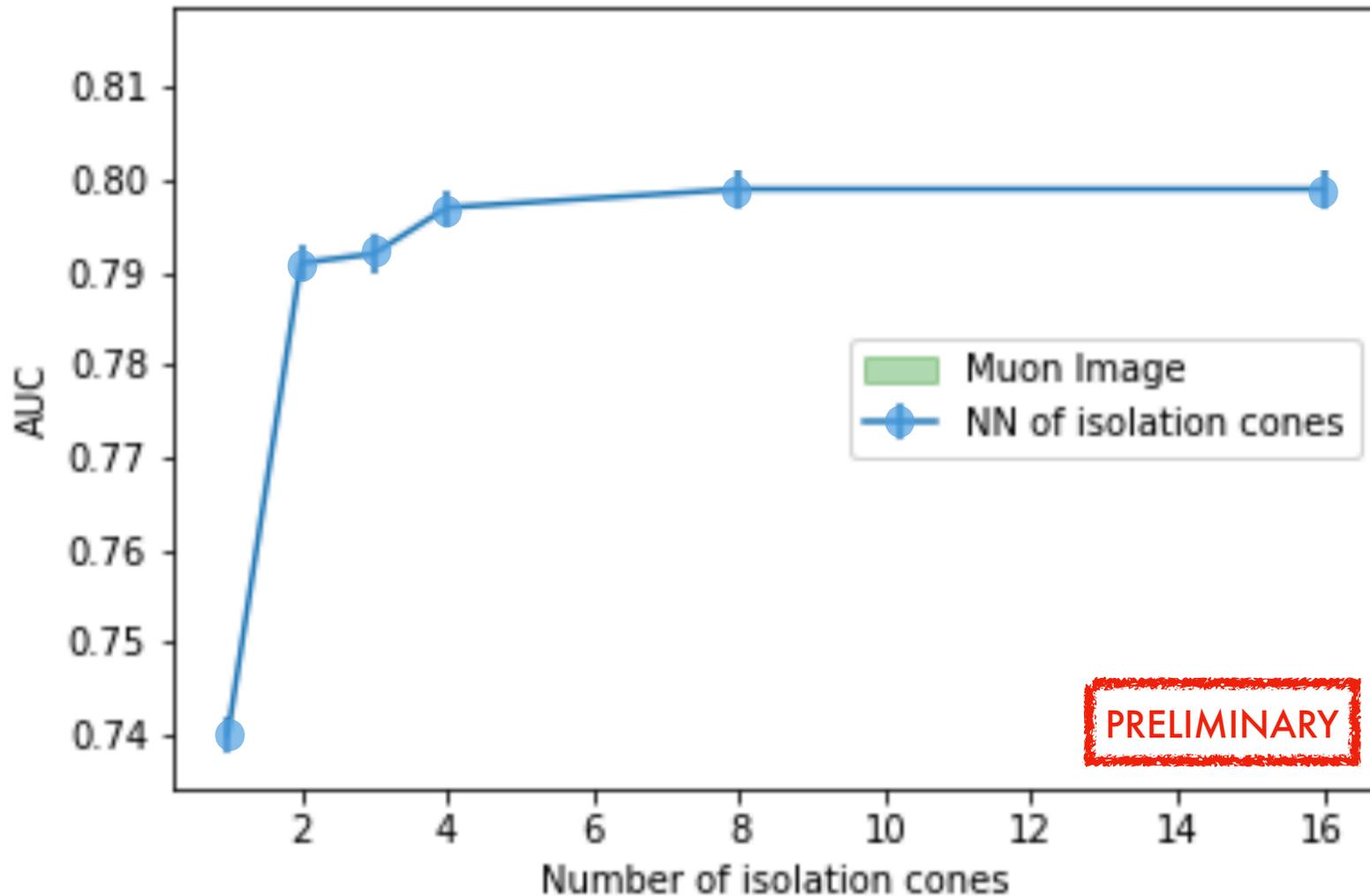
Isolation



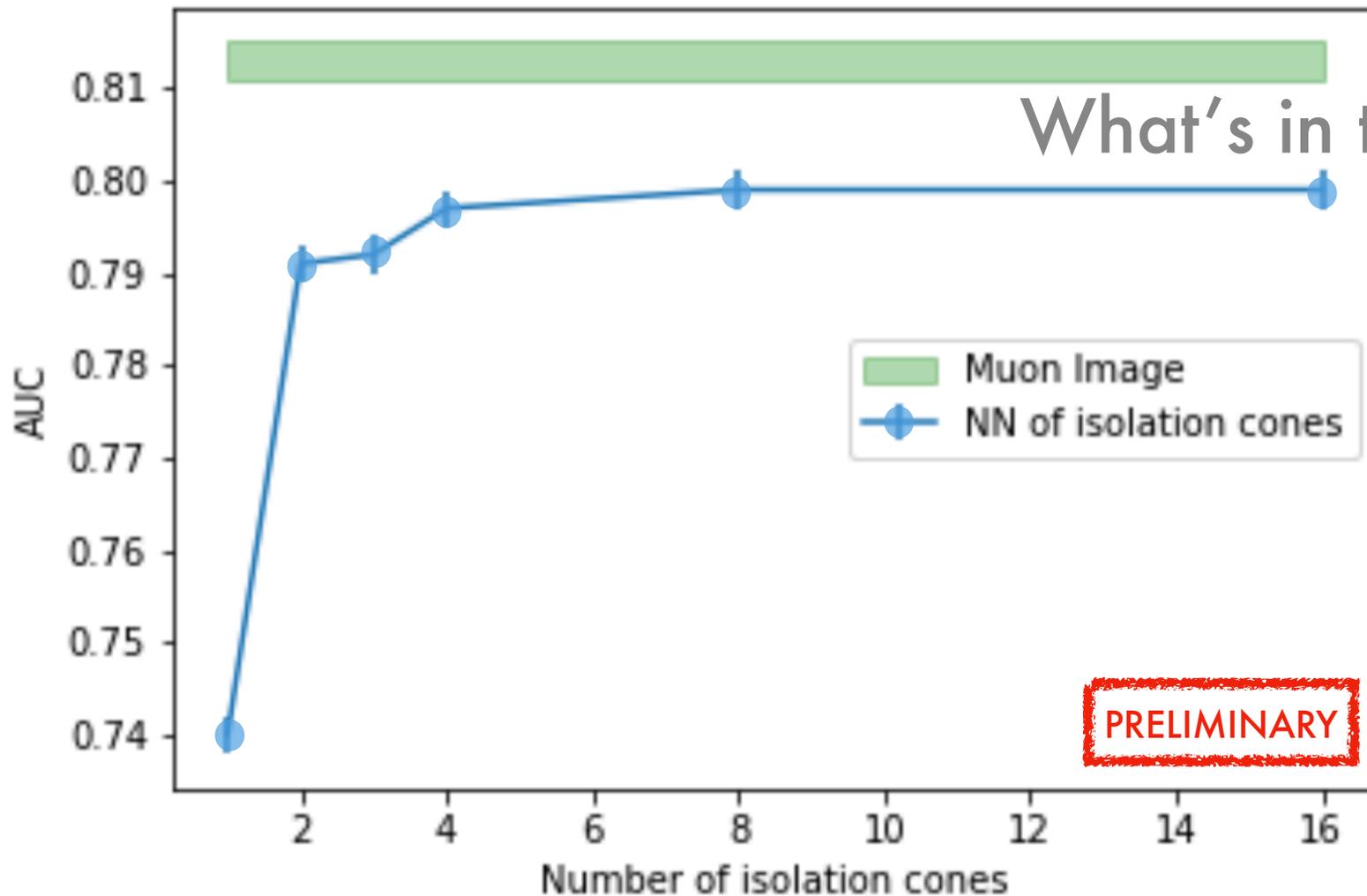
More isolation



Most isolation



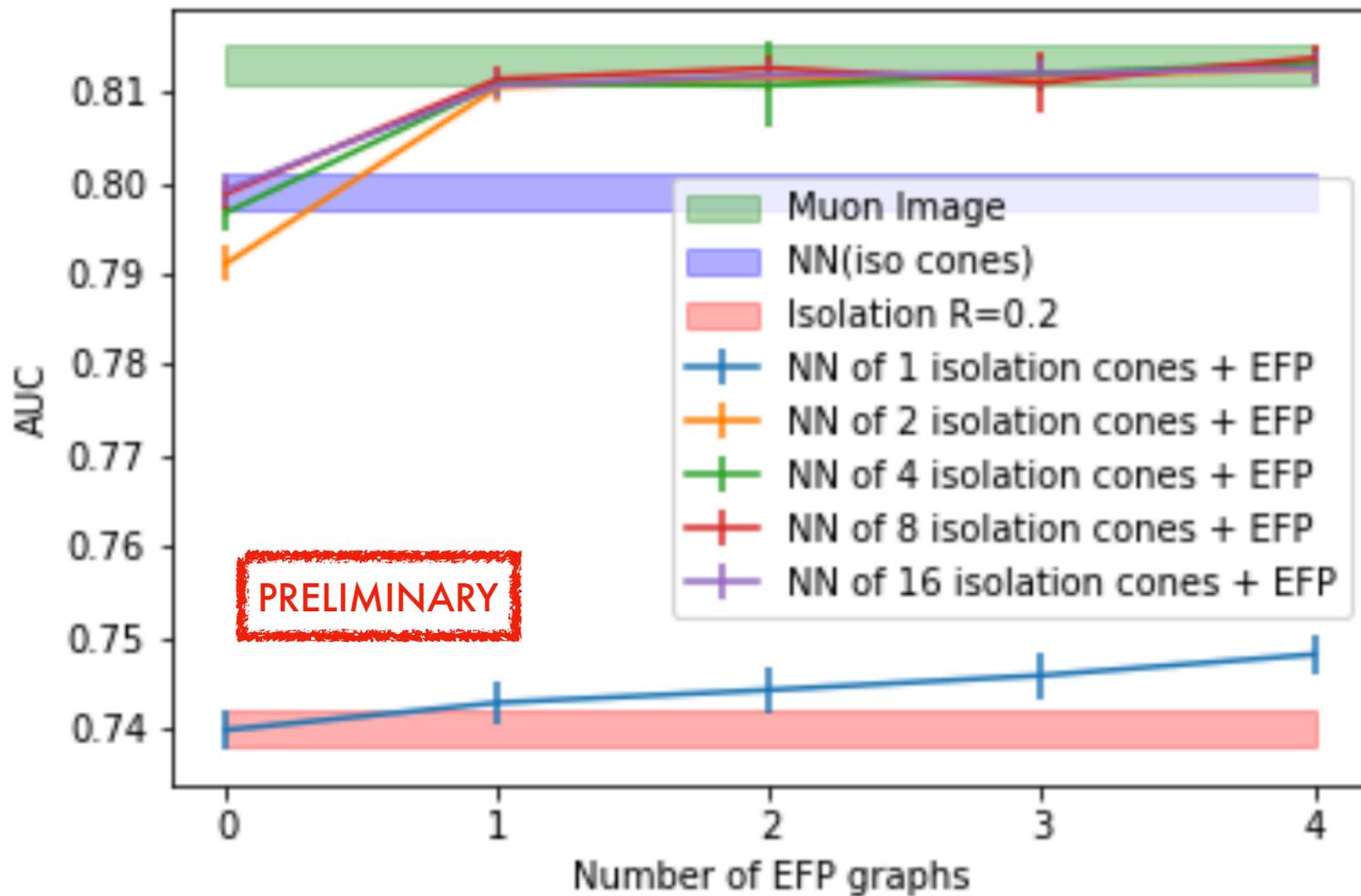
What can ML do?



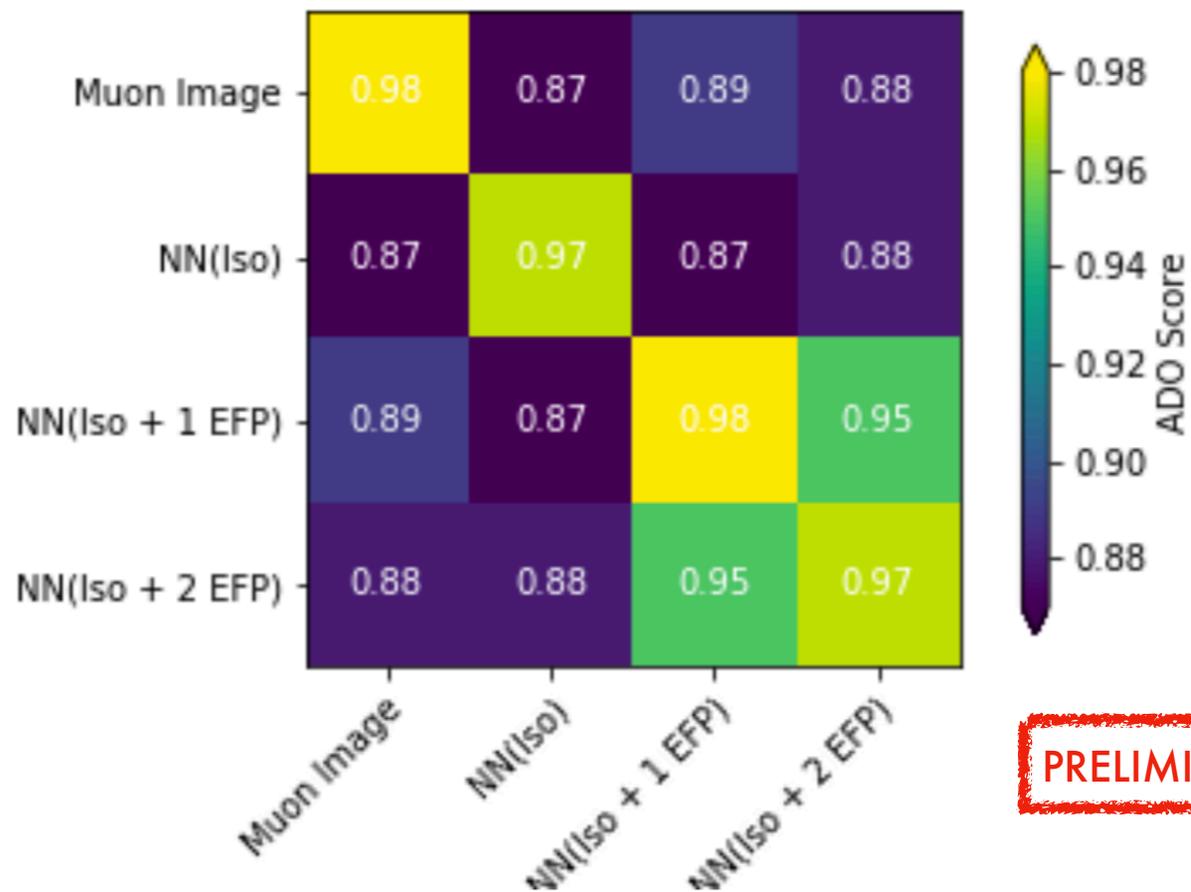
What's in the gap?

PRELIMINARY

Close the gap?



Information



PRELIMINARY

Conclusions

Deep Learning is a powerful new tool

offers faster learning of nonlinear functions

We have many appropriate tasks in HEP

traditional heuristics should be re-examined

No replacement for human intelligence

garbage in will still give garbage out