

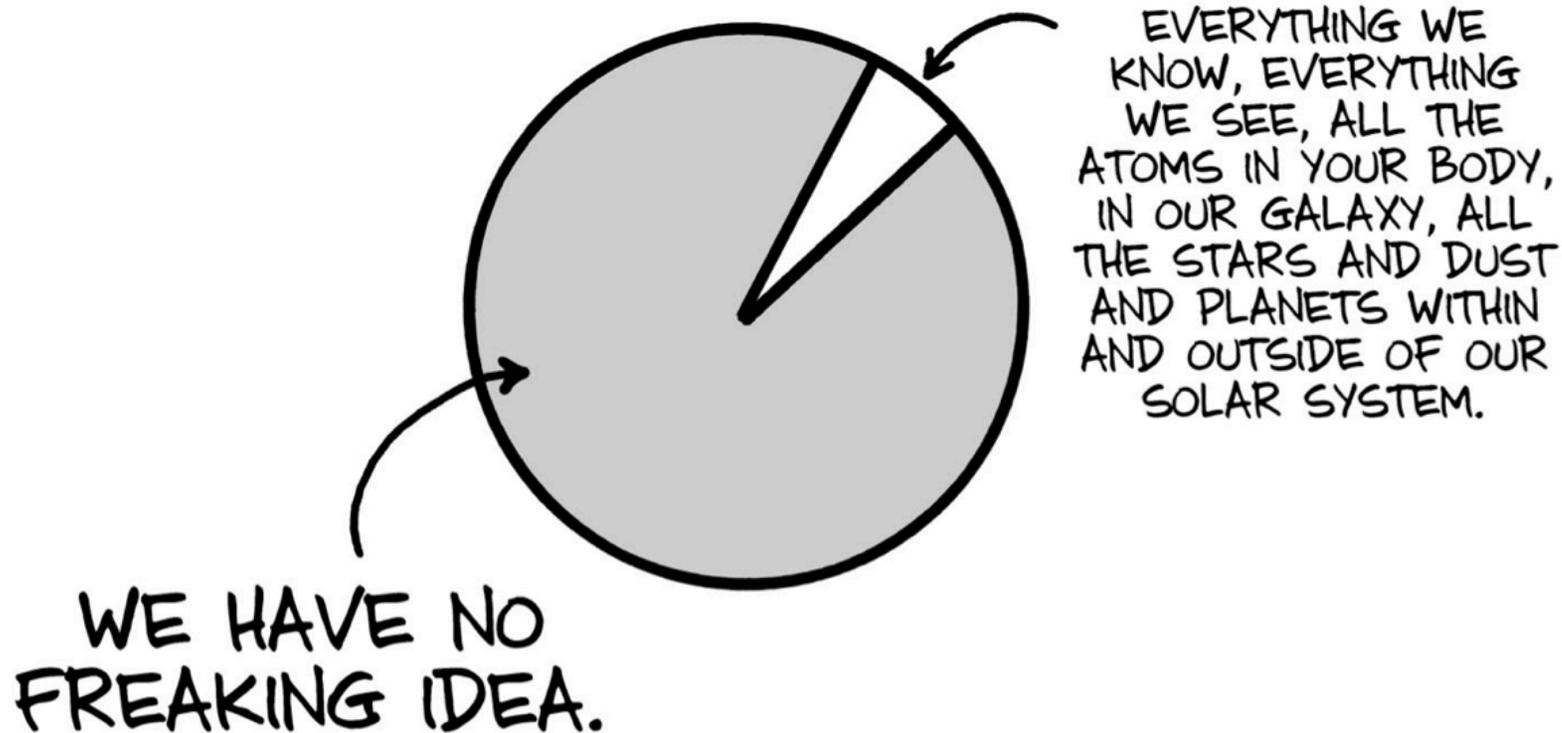
# Open ML problems in High Energy Physics



Daniel Whiteson, UC Irvine  
NIPS2015, ALEPH workshop

# Motivation

THE UNIVERSE AS WE KNOW IT:



SIMULTANEOUSLY, WE KNOW WE KNOW VERY LITTLE...



...SO WE KNOW WHERE TO LOOK.

AND, WE BUILT THIS NEW COLLIDER.



THE MAGIC OF A COLLIDER IS THAT YOU CAN MAKE KINDS OF MATTER THAT YOU DON'T HAVE AROUND. IT'S THIS AMAZING QUANTUM MECHANICAL MAGIC THAT'S JUST TURNED ON.

THESE TWO THINGS ARE COMING TOGETHER  
**RIGHT NOW**

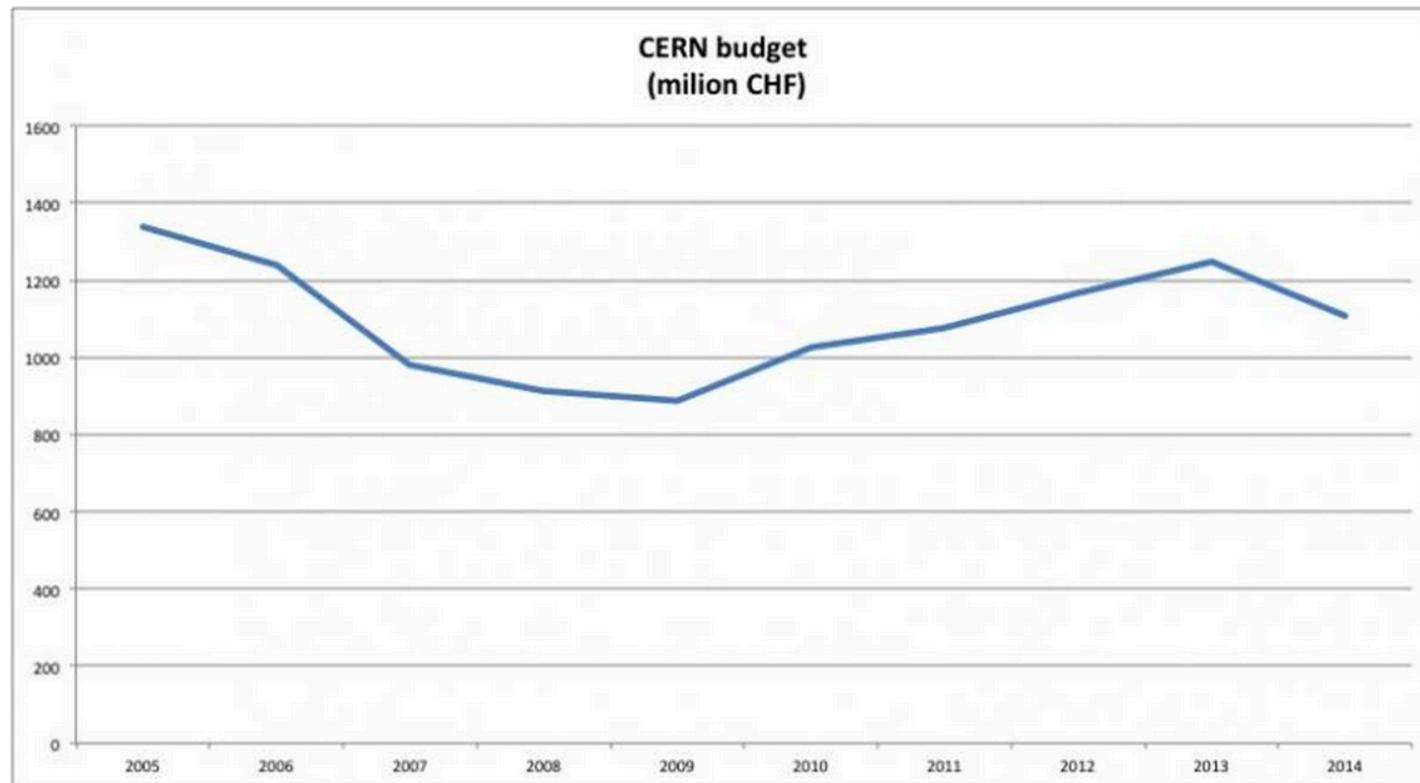


JORGE CHAM © 2011



# Why do we need ML?

## Budget overview



# Outline

## Defining the problem

Dimensionality reduction

## Current approaches

Deep(er) networks for low(er)-level data

## Dimensionality increase

Parameterized networks for sets of problems

# Hypothesis testing

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

1.

THE STANDARD MODEL			
Fermions			
Quarks	<b><i>u</i></b> up	<b><i>c</i></b> charm	<b><i>t</i></b> top
	<b><i>d</i></b> down	<b><i>s</i></b> strange	<b><i>b</i></b> bottom
Leptons	<b><math>V_e</math></b> electron neutrino	<b><math>V_\mu</math></b> muon neutrino	<b><math>V_\tau</math></b> tau neutrino
	<b><i>e</i></b> electron	<b><math>\mu</math></b> muon	<b><math>\tau</math></b> tau

# Hypothesis testing

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

1.

2.

THE STANDARD MODEL			
Fermions			
Quarks	<b><i>u</i></b> up	<b><i>c</i></b> charm	<b><i>t</i></b> top
	<b><i>d</i></b> down	<b><i>s</i></b> strange	<b><i>b</i></b> bottom
Leptons	<b><math>V_e</math></b> electron neutrino	<b><math>V_\mu</math></b> muon neutrino	<b><math>V_\tau</math></b> tau neutrino
	<b><i>e</i></b> electron	<b><math>\mu</math></b> muon	<b><math>\tau</math></b> tau

THE STANDARD MODEL PLUS X				
Fermions				
Quarks	<b><i>u</i></b> up	<b><i>c</i></b> charm	<b><i>t</i></b> top	<b>X</b>
	<b><i>d</i></b> down	<b><i>s</i></b> strange	<b><i>b</i></b> bottom	
Leptons	<b><math>V_e</math></b> electron neutrino	<b><math>V_\mu</math></b> muon neutrino	<b><math>V_\tau</math></b> tau neutrino	
	<b><i>e</i></b> electron	<b><math>\mu</math></b> muon	<b><math>\tau</math></b> tau	

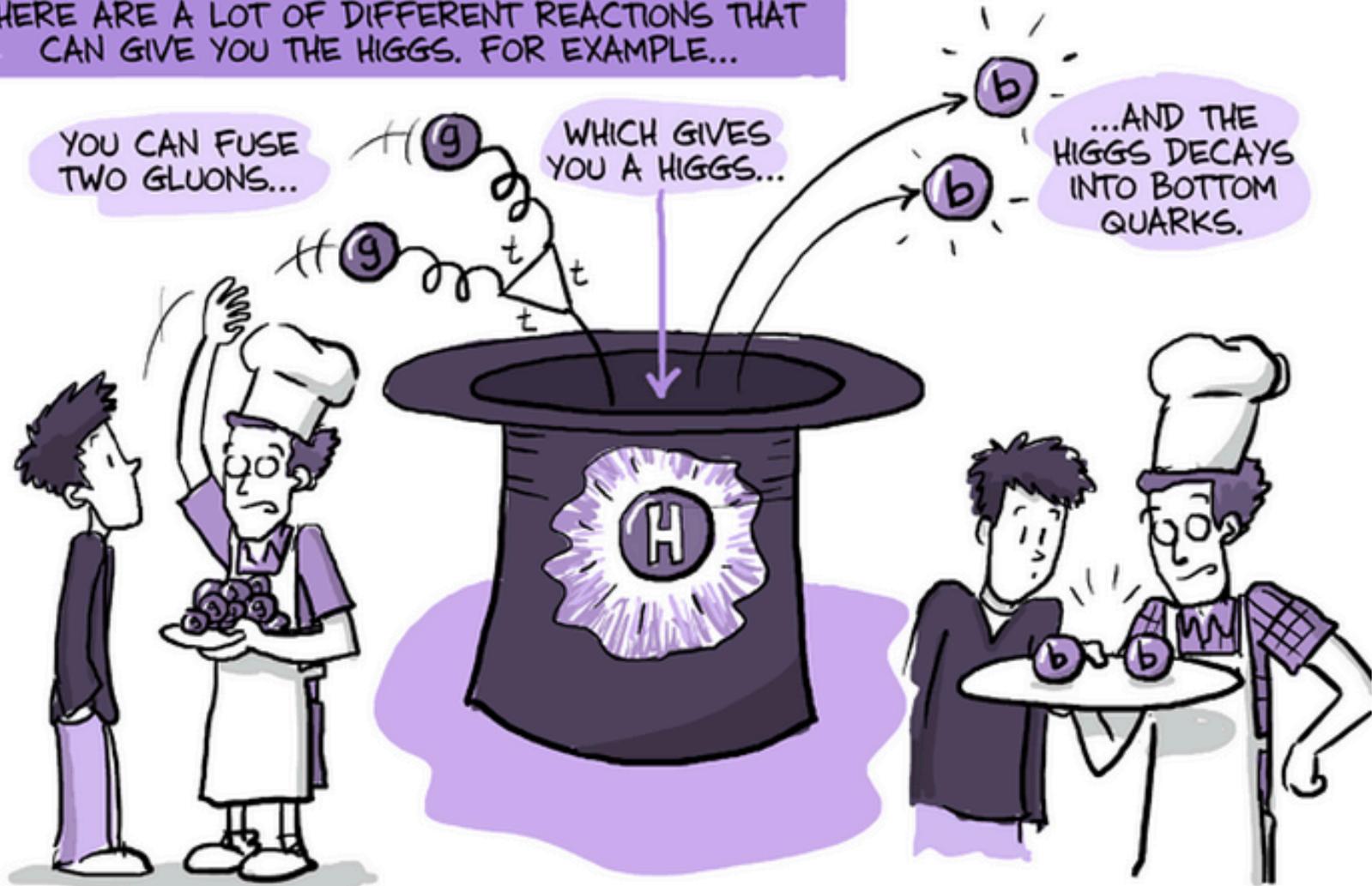
# Making a new particle

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

YOU CAN FUSE TWO GLUONS...

WHICH GIVES YOU A HIGGS...

...AND THE HIGGS DECAYS INTO BOTTOM QUARKS.

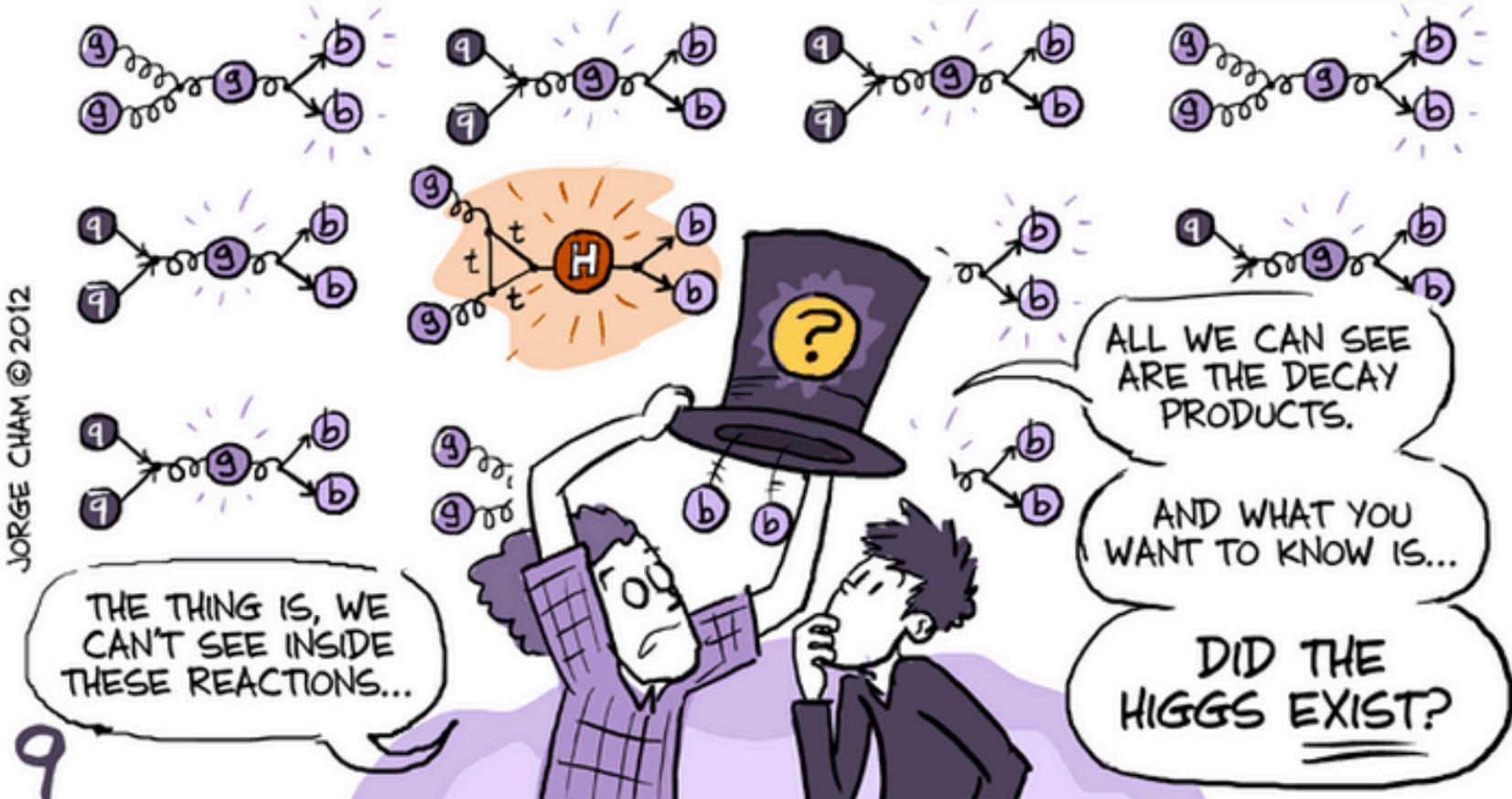


# Backgrounds

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

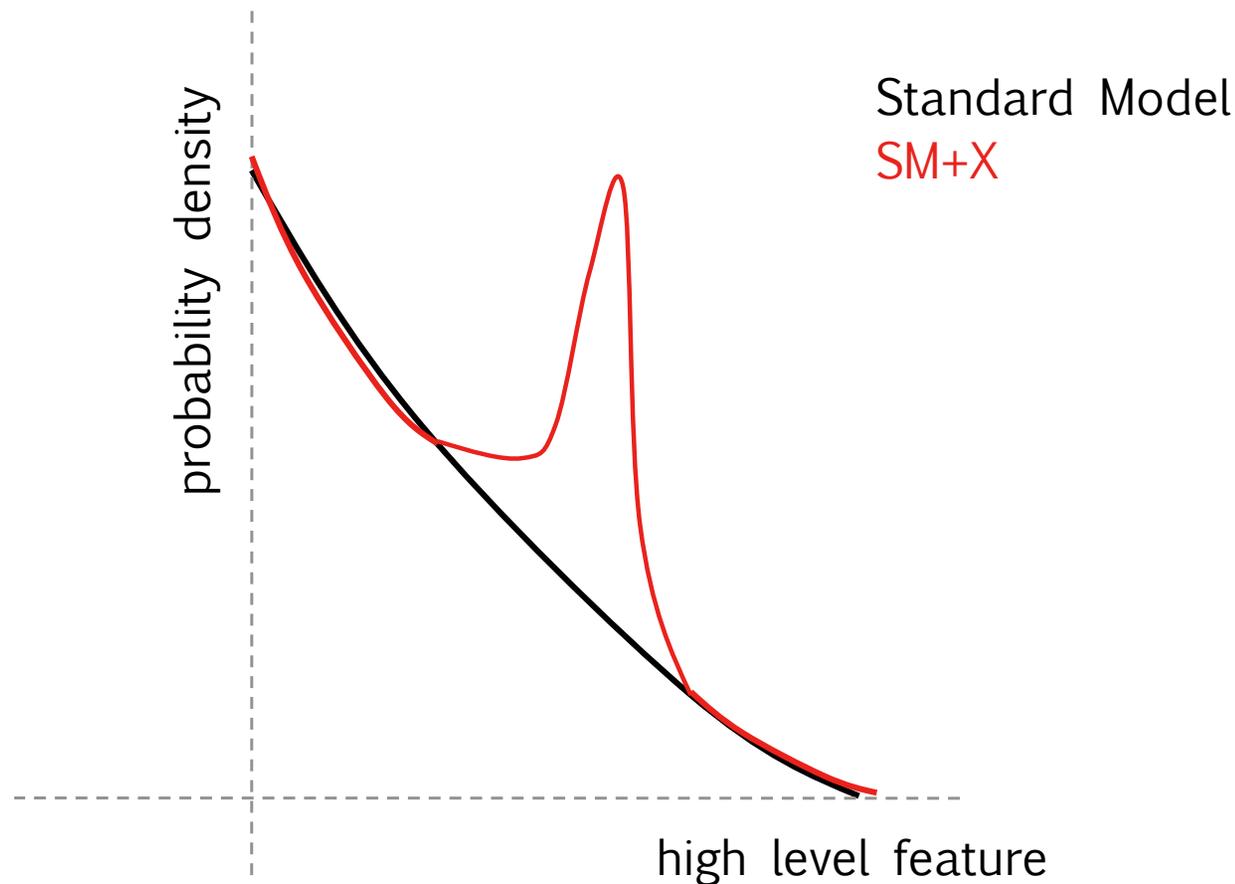
IT'S ONE OF THE MOST COMMON THINGS TO MAKE.

JORGE CHAM © 2012



# Hypothesis Testing

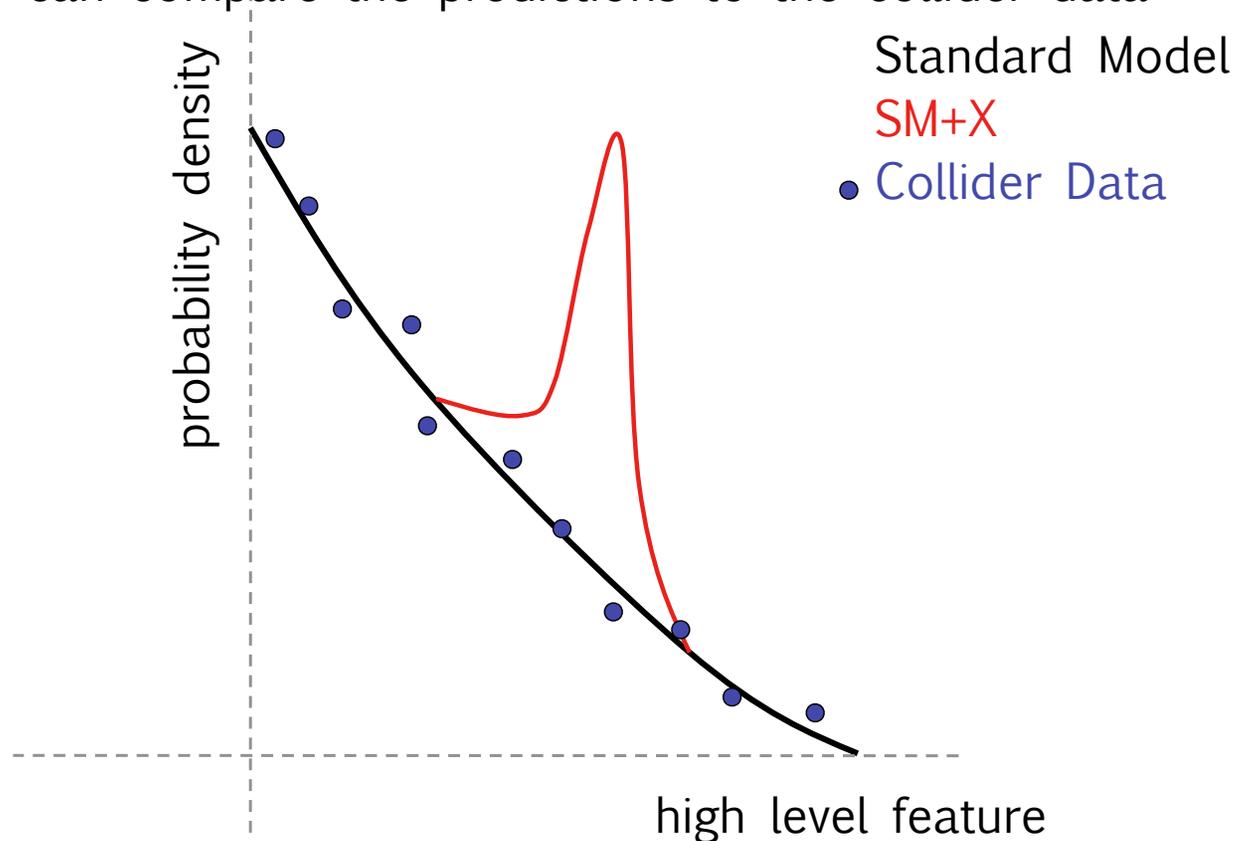
Which are boiled down to a few (10-50) high level features



We look for a region in feature space where the two hypotheses have large differences in their predictions.

# Hypothesis Testing

We can compare the predictions to the collider data



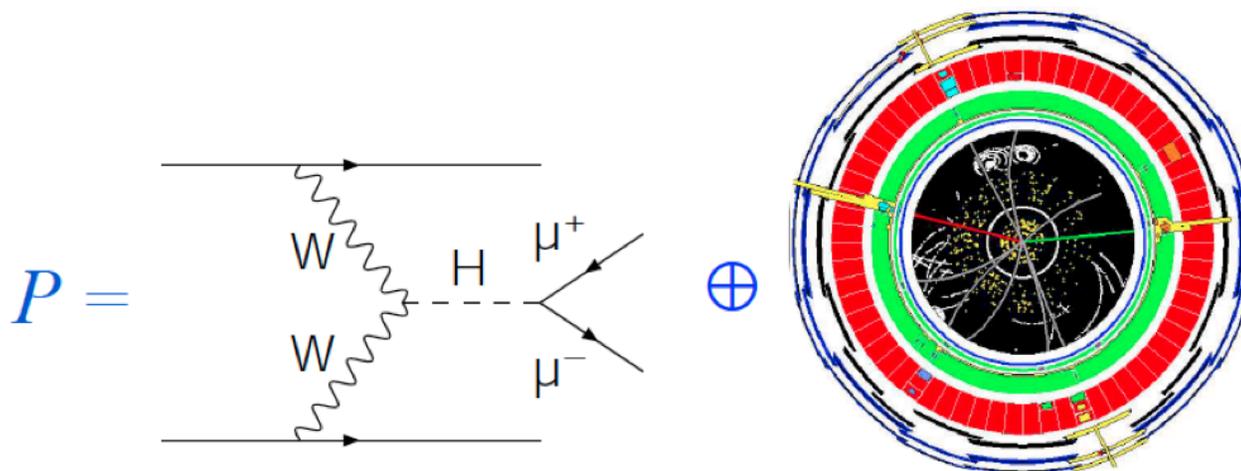
Which can tell us which hypothesis is preferred via a likelihood ratio:

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

# $P(\text{data} \mid \text{SM})$

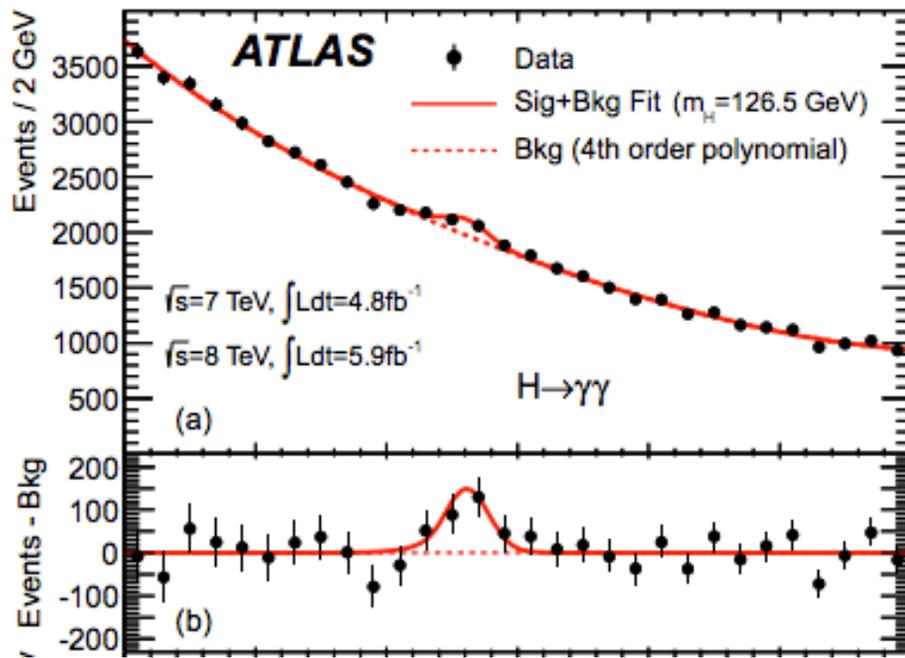
We don't know how to calculate

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

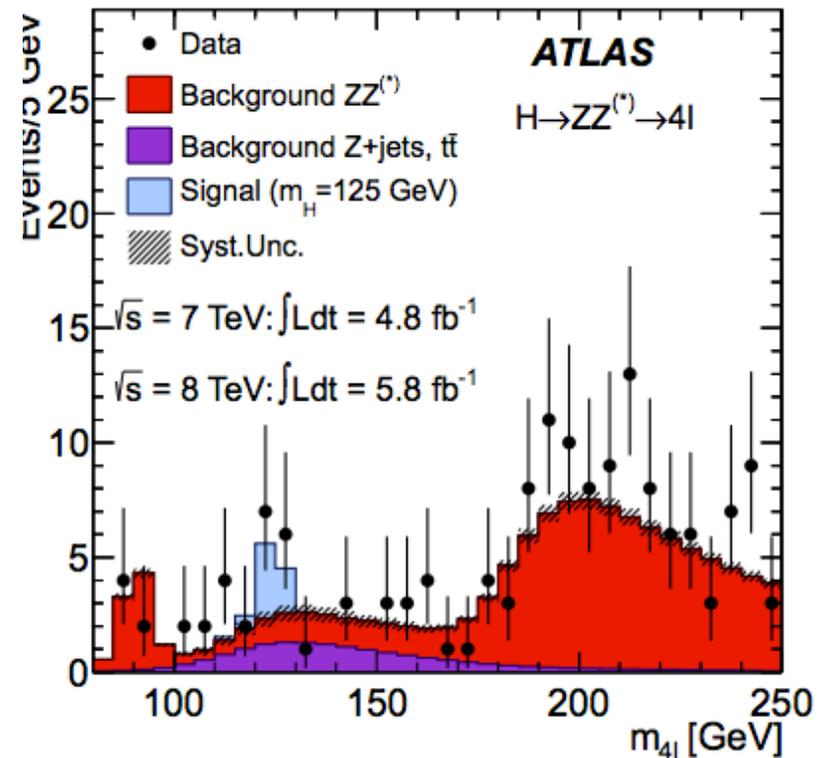


# ATLAS

## Two photons



## Two Zs to 4 leptons

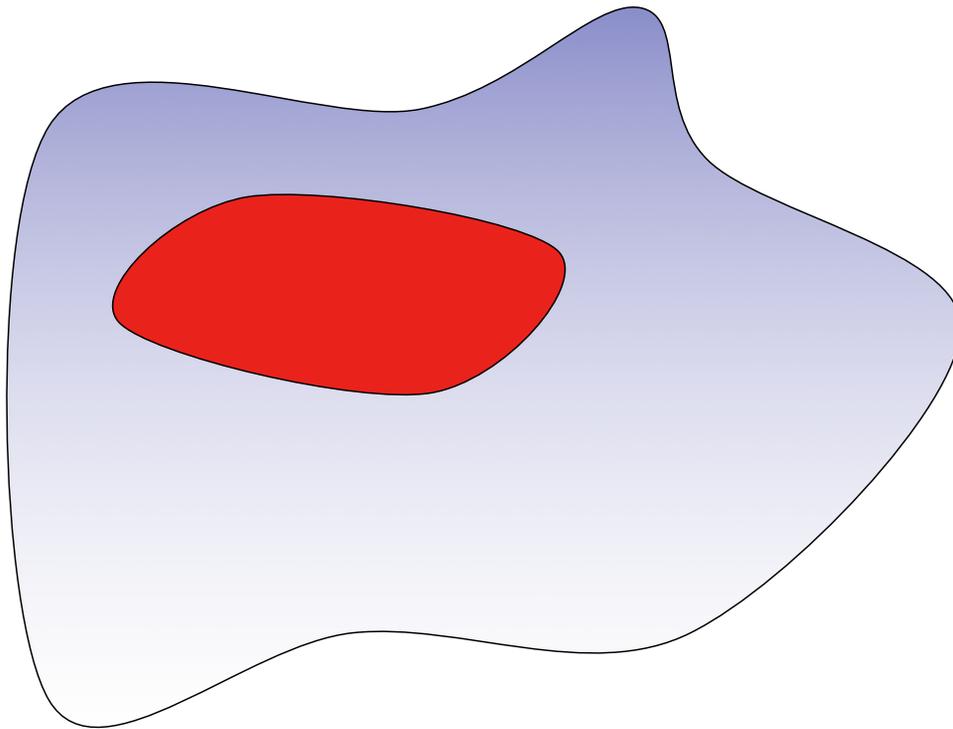


# Simulation

Standard Model

X

feature 2



feature 1

generate simulated collisions for both hypotheses

Use histograms to define

$$P(\text{data} \mid \text{SM}+\text{X}),$$

and

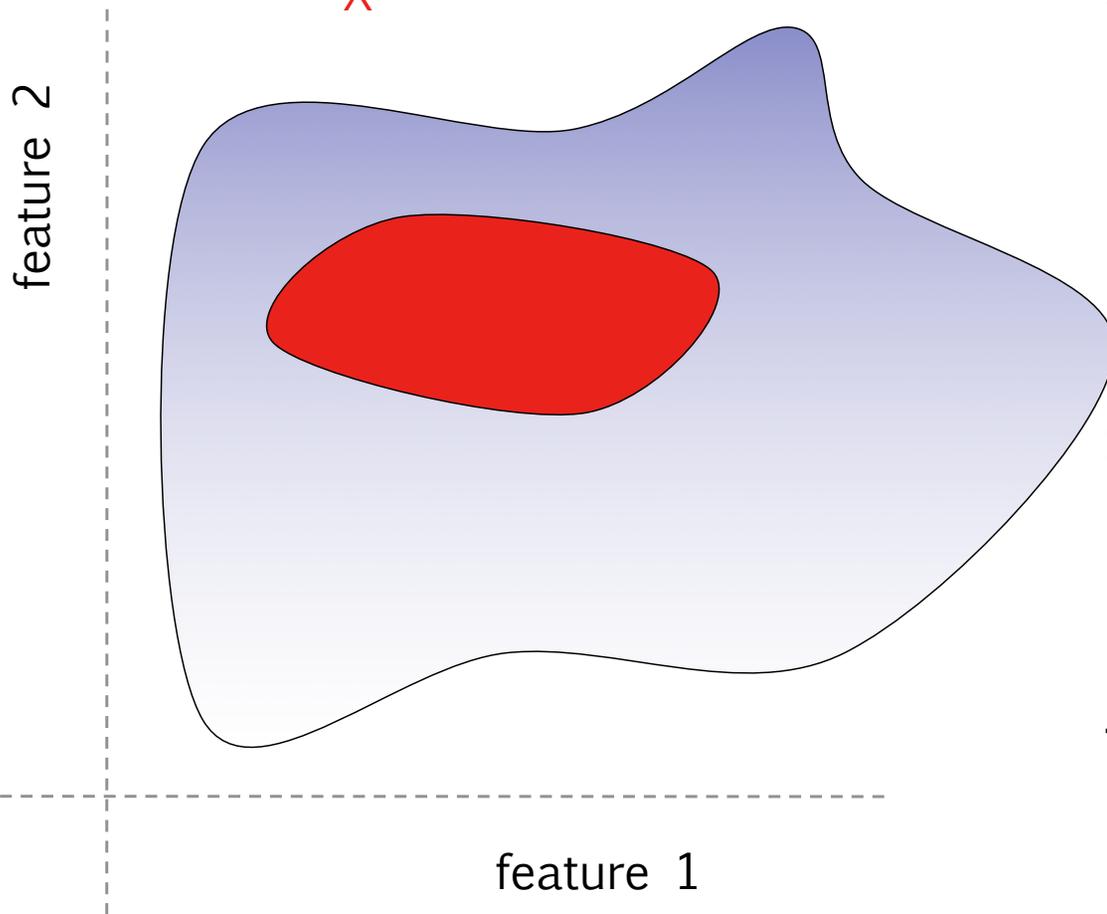
$$P(\text{data} \mid \text{SM})$$

(We've also tried kernel-based probability density estimates).

# Dimensionality

Standard Model

X



Want to harness all differences  
between hypotheses.



Use as *many* features as possible.

Simulation is slow  
(~1 - 10 min/collision)

Filling a  $d$ -dimensional histogram  
requires  $\sim 100^d$  simulated collisions!



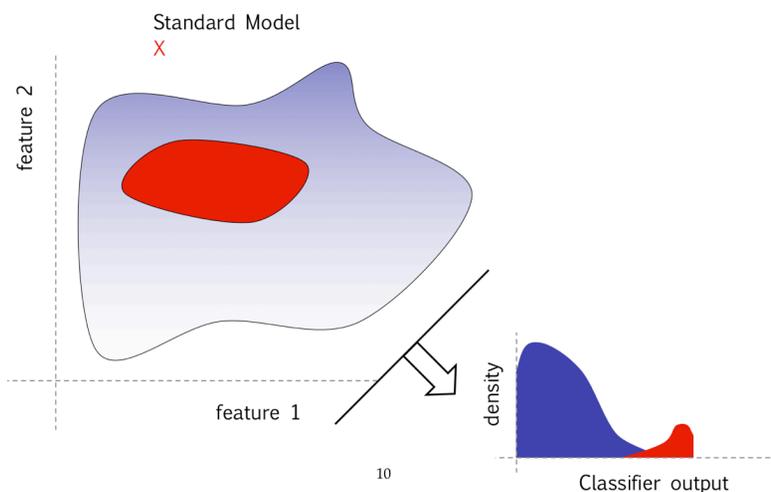
Use as *few* features as possible.

# Function

Find a function

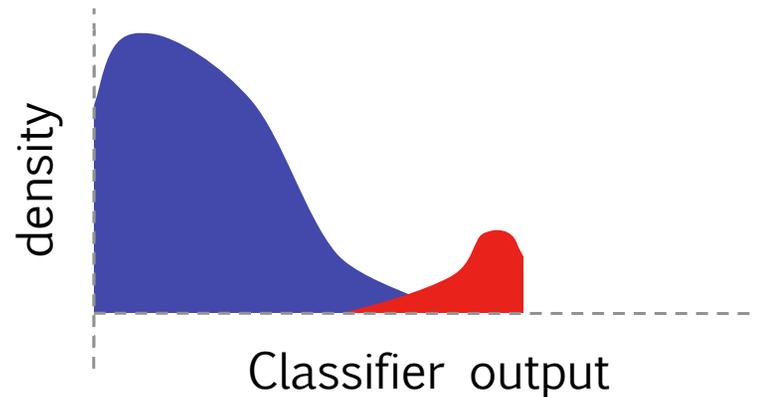
$$f(\mathcal{X}^N \rightarrow \mathcal{R}^1)$$

Neural networks  
can learn these  
shapes in high-dim  
and summarize  
in a 1D output



# Dimensional Reduction

This dimensional reduction can be very helpful.



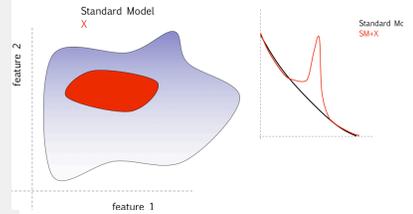
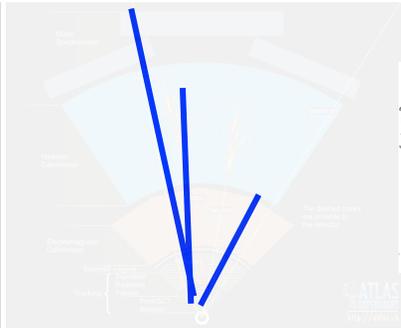
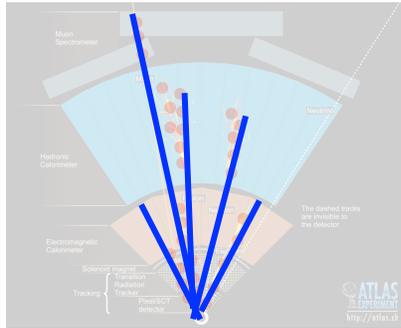
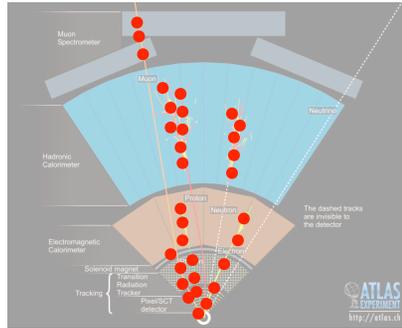
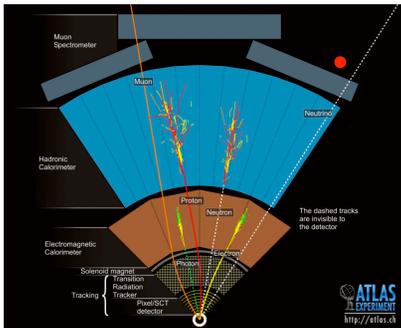
Summarize the differences between the hypotheses

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

And require a histogram in only one dimension

# Dimensionality

Raw	Sparsified	Reco	Select	Physics	Ana
$1e7$	$1e4$	100-ish*	50	10	1



ish\*: dimensionality is variable here

# So what?

What's wrong with this picture?

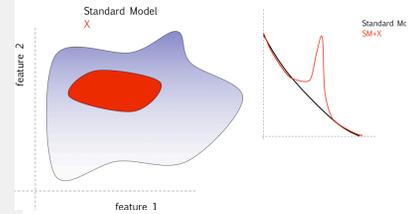
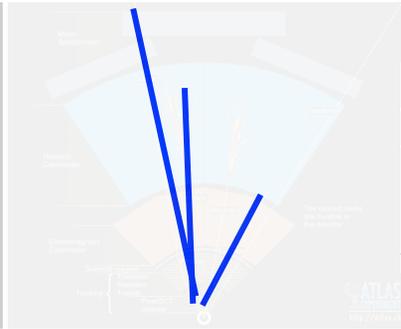
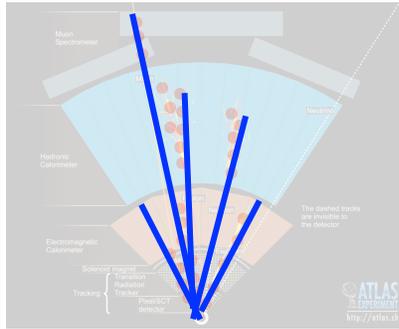
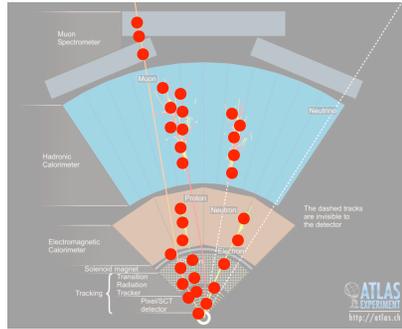
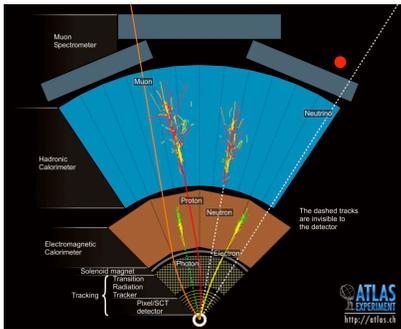
Nothing dramatic

Each step is done using domain knowledge and decades of experience.

**But is there information lost? Certainly.**

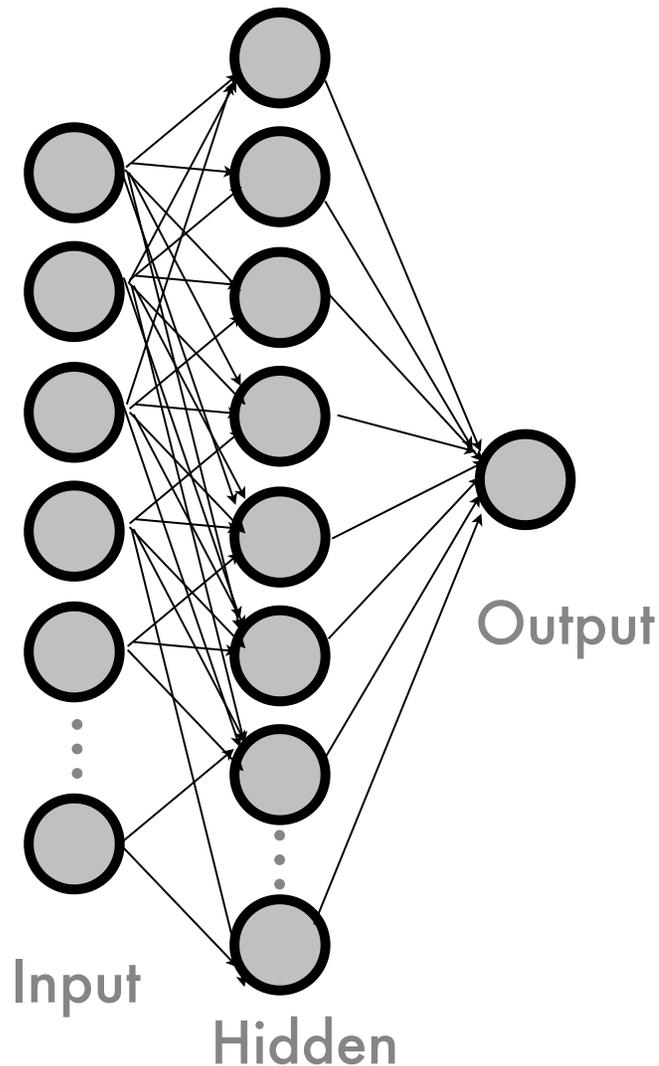
# Dimensionality

Raw	Sparsified	Reco	Select	Physics	Ana
$1e7$	$1e4$	100-ish*	50	10	1



# Neural Networks

Essentially a functional fit with many parameters



## Function

Each neuron's output is a function of the weighted sum of inputs.

## Goal

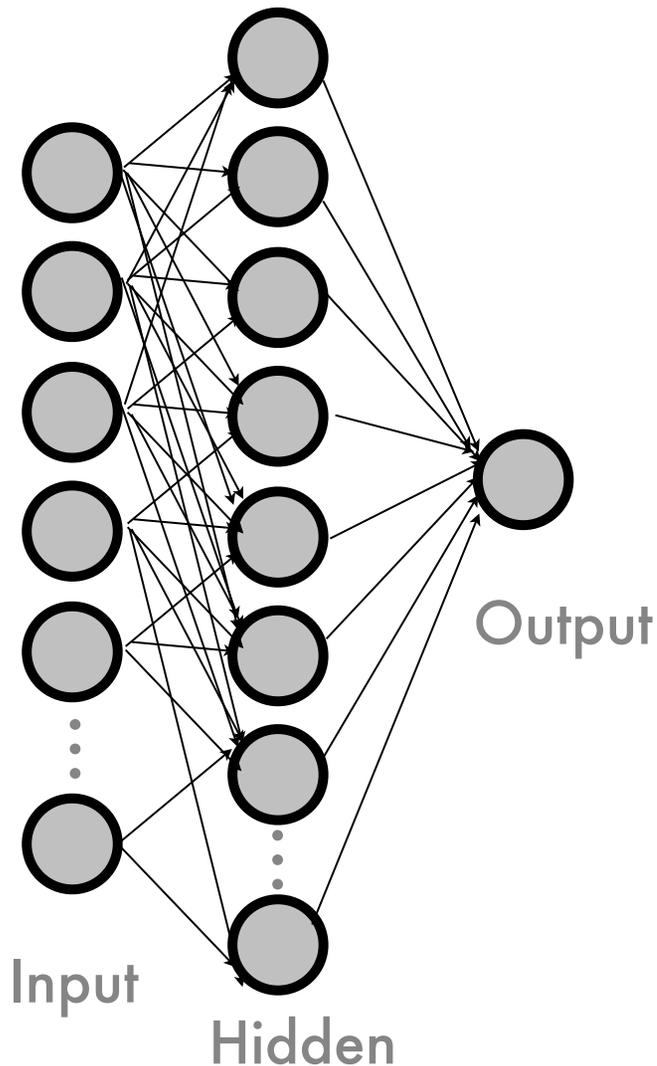
find set of weights which give **most useful function**

## Learning

give examples, back-propagate error to adjust weights

# Neural Networks

Essentially a functional fit with many parameters



## Problem:

Networks with  $> 1$  layer are very difficult to train.

## Consequence:

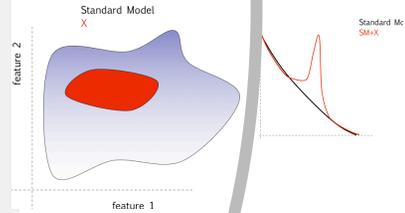
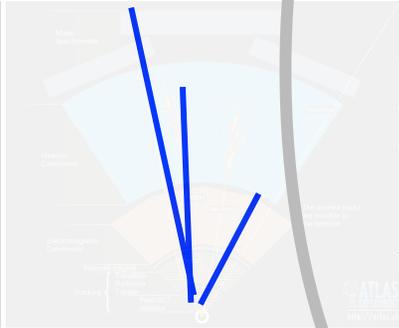
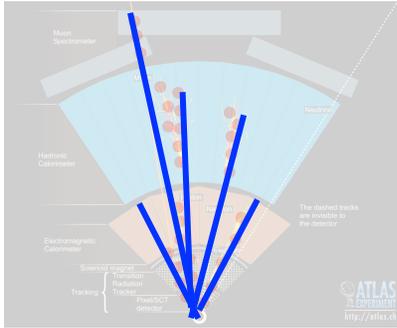
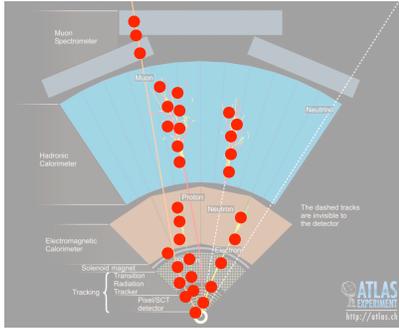
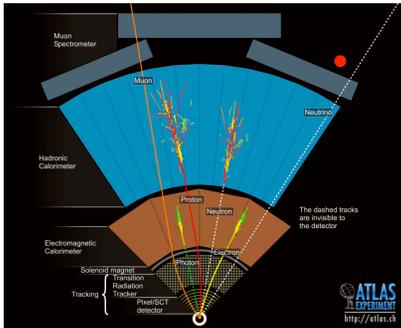
Networks are not good at learning non-linear functions.  
**(like invariant masses!)**

## In short:

Can't just throw 4-vectors at NN.

# Summary

Raw	Sparsified	Reco	Select	Physics	Ana
$1e7$	$1e4$	100-ish*	50	10	1



Why we have this step

# Search for Input

ATLAS-CONF-2013-108

Can't just use  
reco or selected  
objects

Can't give it too  
many inputs

Painstaking search  
through input  
feature space.

Variable	VBF			Boosted		
	$\tau_{lep}\tau_{lep}$	$\tau_{lep}\tau_{had}$	$\tau_{had}\tau_{had}$	$\tau_{lep}\tau_{lep}$	$\tau_{lep}\tau_{had}$	$\tau_{had}\tau_{had}$
$m_{\tau\tau}^{MMC}$	•	•	•	•	•	•
$\Delta R(\tau, \tau)$	•	•	•		•	•
$\Delta\eta(j_1, j_2)$	•	•	•			
$m_{j_1, j_2}$	•	•	•			
$\eta_{j_1} \times \eta_{j_2}$		•	•			
$p_T^{total}$		•	•			
sum $p_T$					•	•
$p_T(\tau_1)/p_T(\tau_2)$					•	•
$E_T^{miss}$ $\phi$ centrality		•	•	•	•	•
$x_{\tau 1}$ and $x_{\tau 2}$						•
$m_{\tau\tau, j_1}$				•		
$m_{\ell_1, \ell_2}$				•		
$\Delta\phi_{\ell_1, \ell_2}$				•		
sphericity				•		
$p_T^{\ell_1}$				•		
$p_T^{j_1}$				•		
$E_T^{miss}/p_T^{\ell_2}$				•		
$m_T$		•			•	
$\min(\Delta\eta_{\ell_1, \ell_2, jets})$	•					
$j_3$ $\eta$ 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  $\tau_{had}$ , or between the two  $\tau_{had}$  candidates, depending on the decay mode.



# Search for Input

ATLAS-CONF-2013-108

Can't just use  
reco or selected  
objects

Can't give  
many inputs

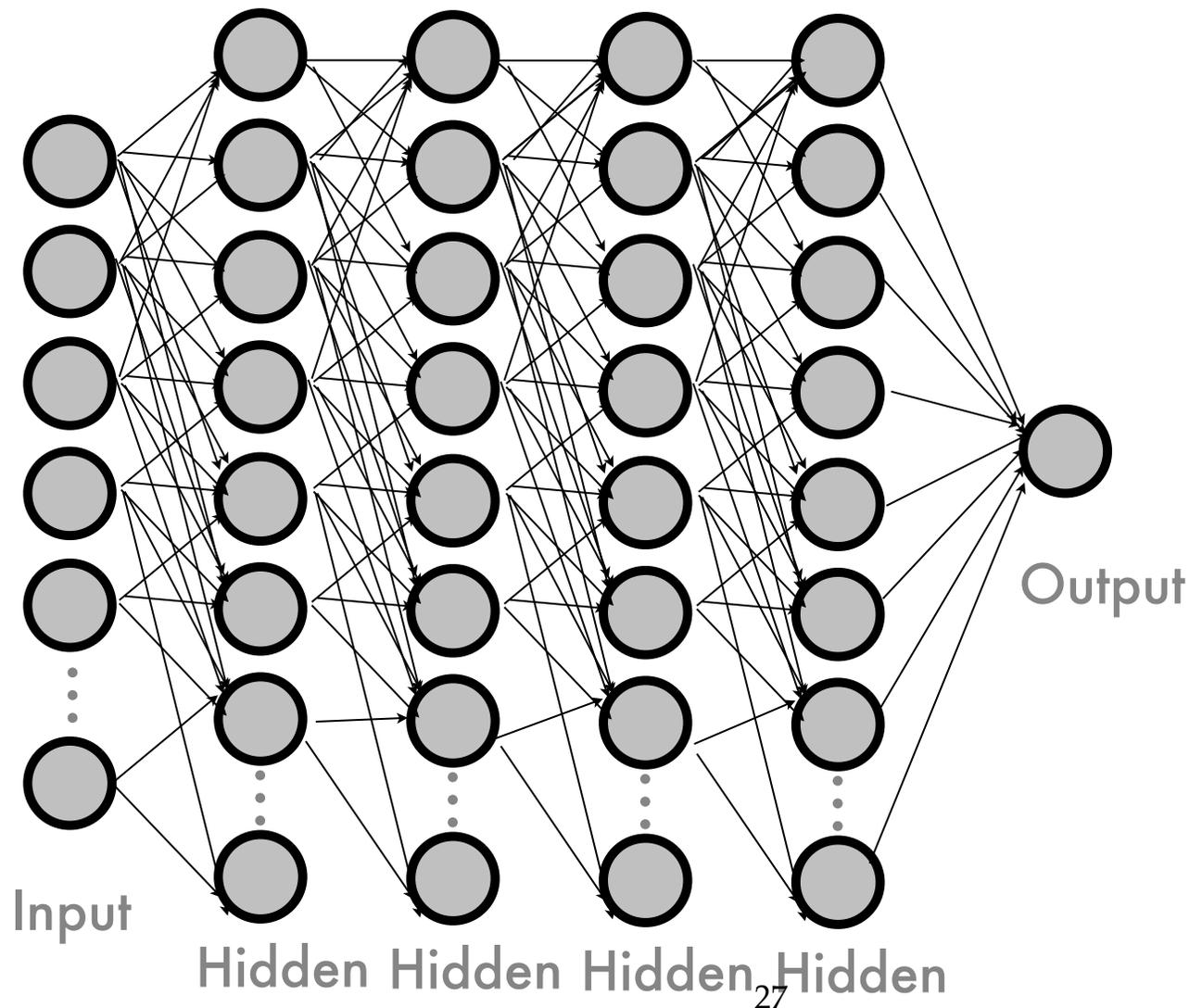
Painstaking search  
through input  
feature space.

**Also true for  
BDTs, SVNs, etc**

Variable	VBF		Boosted	
	$\tau_{lep}\tau_{lep}$	$\tau_{lep}\tau_{had}$	$\tau_{lep}\tau_{lep}$	$\tau_{had}\tau_{had}$
$m_{TT}^{MMC}$			•	•
$\Delta R(\tau, \tau)$				•
$p_T^{j_1}$				
$E_T^{miss} / p_T^{\ell_2}$				
$m_T$		•		•
$\min(\Delta\eta_{\ell_1, \ell_2, jets})$	•			
$j_3$ $\eta$ 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  $\tau_{had}$ , or between the two  $\tau_{had}$  candidates, depending on the decay mode.

# Deep networks



New tools  
let us  
train  
deep  
networks.

How well  
do they work?

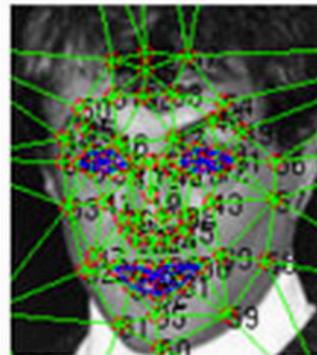
# Real world applications



(a)



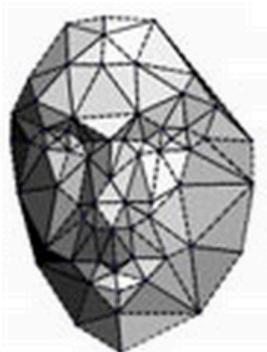
(b)



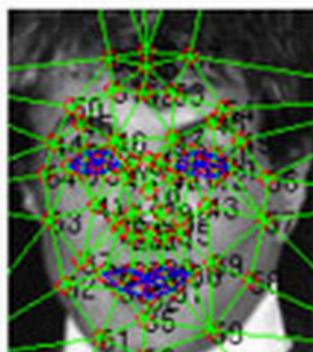
(c)



(d)



(e)



(f)



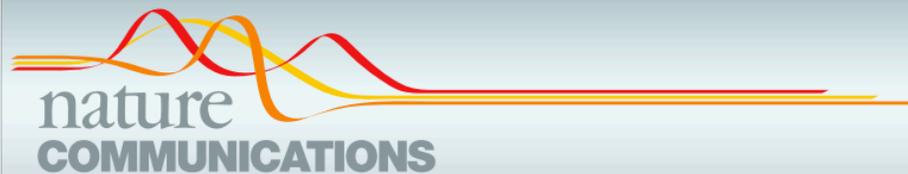
(g)



(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)

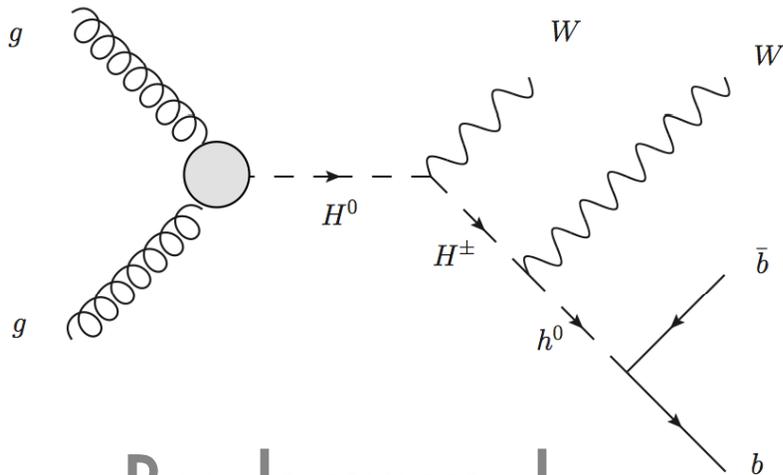
# Searching for exotic particles in high-energy physics with deep learning

P. Baldi<sup>1</sup>, P. Sadowski<sup>1</sup> & D. Whiteson<sup>2</sup>

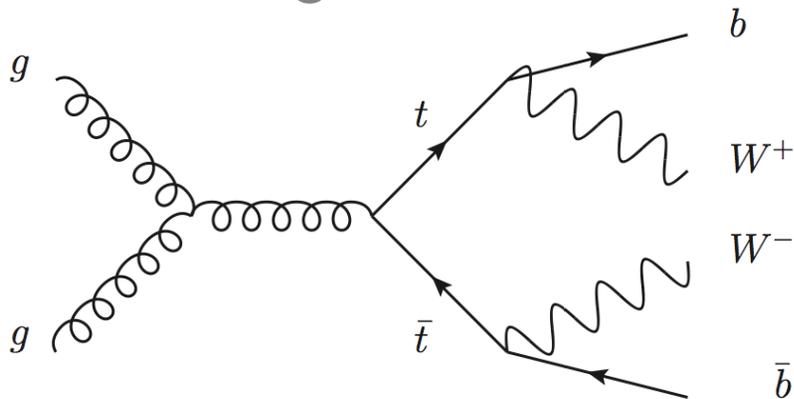
arXiv: 1402.4735

# Benchmark problem

## Signal



## Background



Can deep networks automatically discover useful variables?

# 4-vector inputs

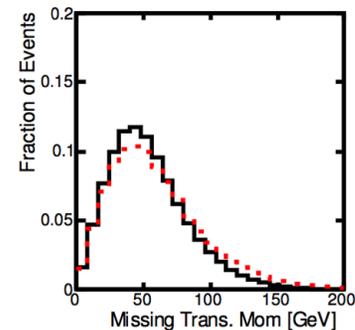
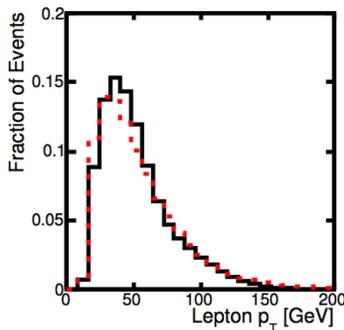
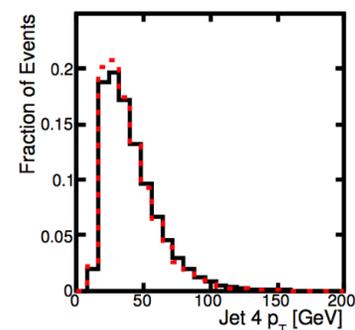
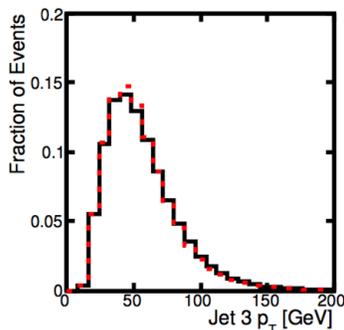
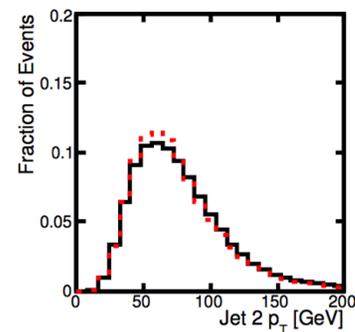
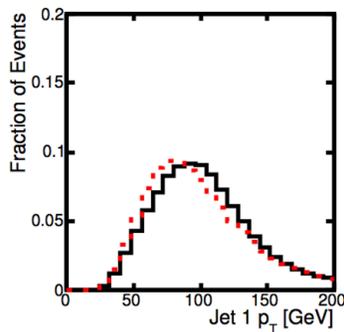
## 21 Low-level vars

jet+lepton mom. (3x5)

missing ET (2)

jet btags (4)

Not much  
separation  
visible in 1D  
projections



# 4-vector inputs

## 7 High-level vars

$m(WWbb)$

$m(Wbb)$

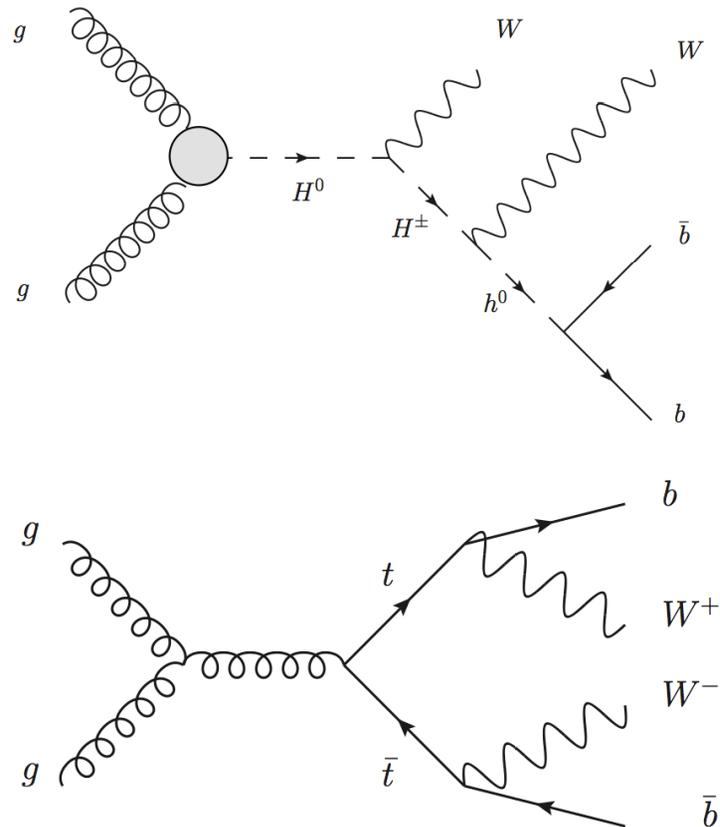
$m(bb)$

$m(bjj)$

$m(jj)$

$m(lv)$

$m(blv)$



# 4-vector inputs

## 7 High-level vars

$m(WWbb)$

$m(Wbb)$

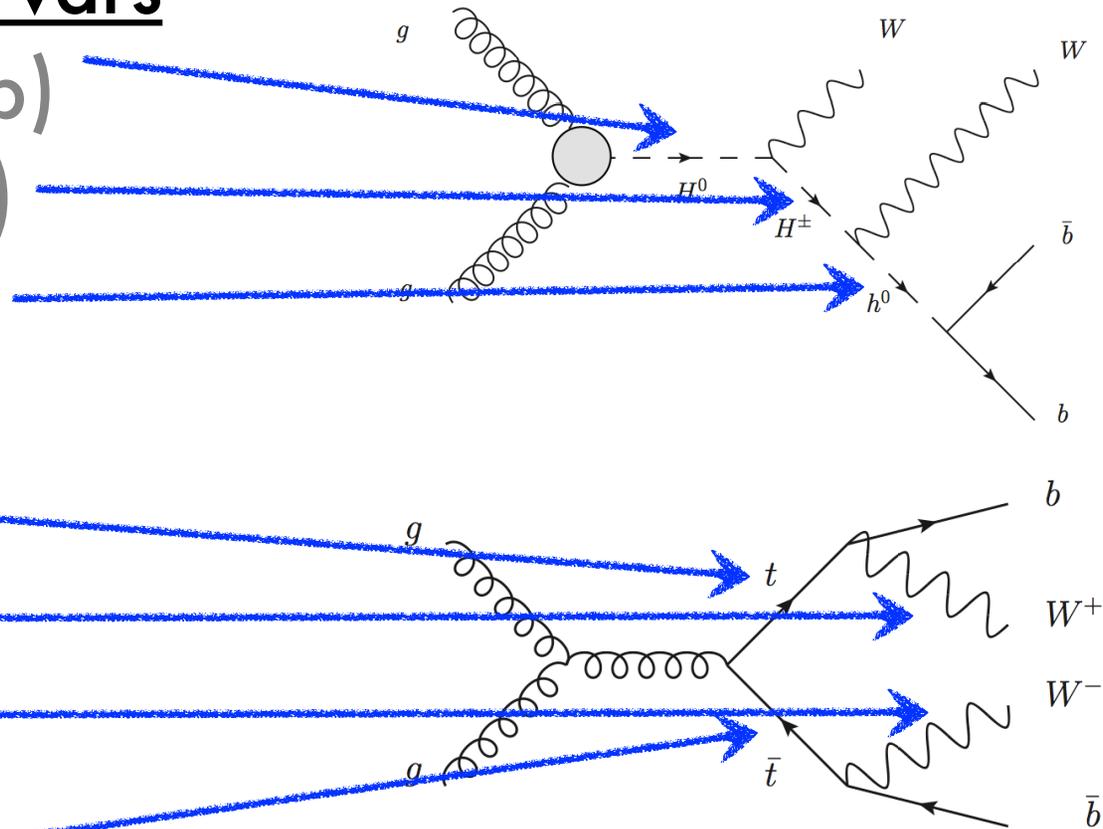
$m(bb)$

$m(bjj)$

$m(jj)$

$m(lv)$

$m(blv)$



# 4-vector inputs

## 7 High-level vars

$m(WWbb)$

$m(Wbb)$

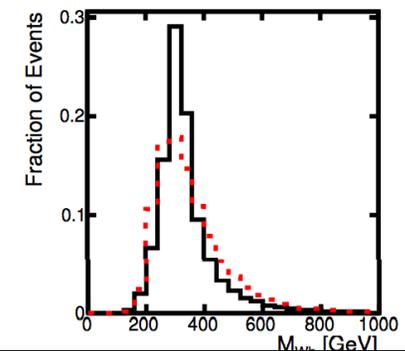
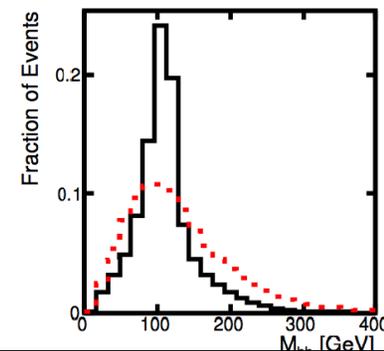
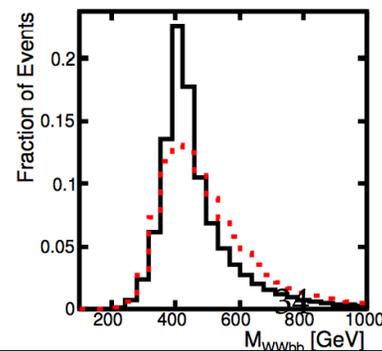
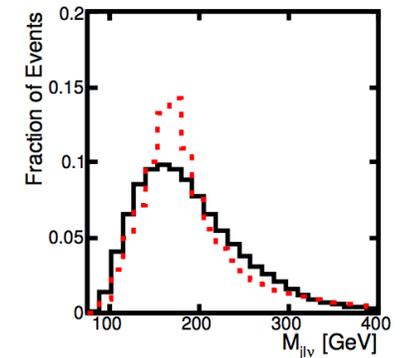
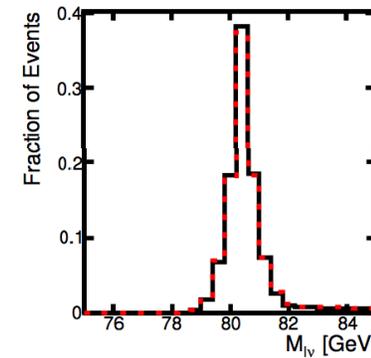
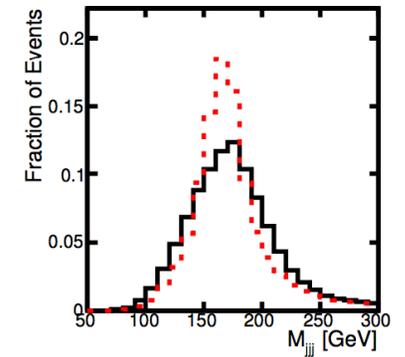
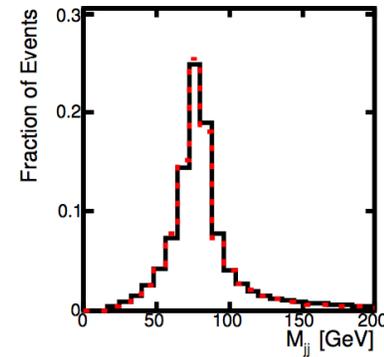
$m(bb)$

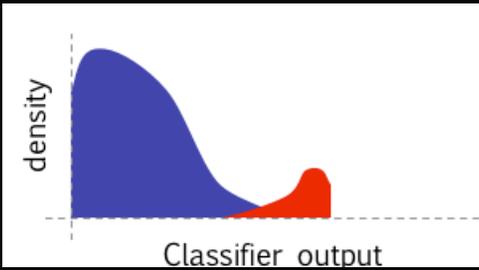
$m(bjj)$

$m(ij)$

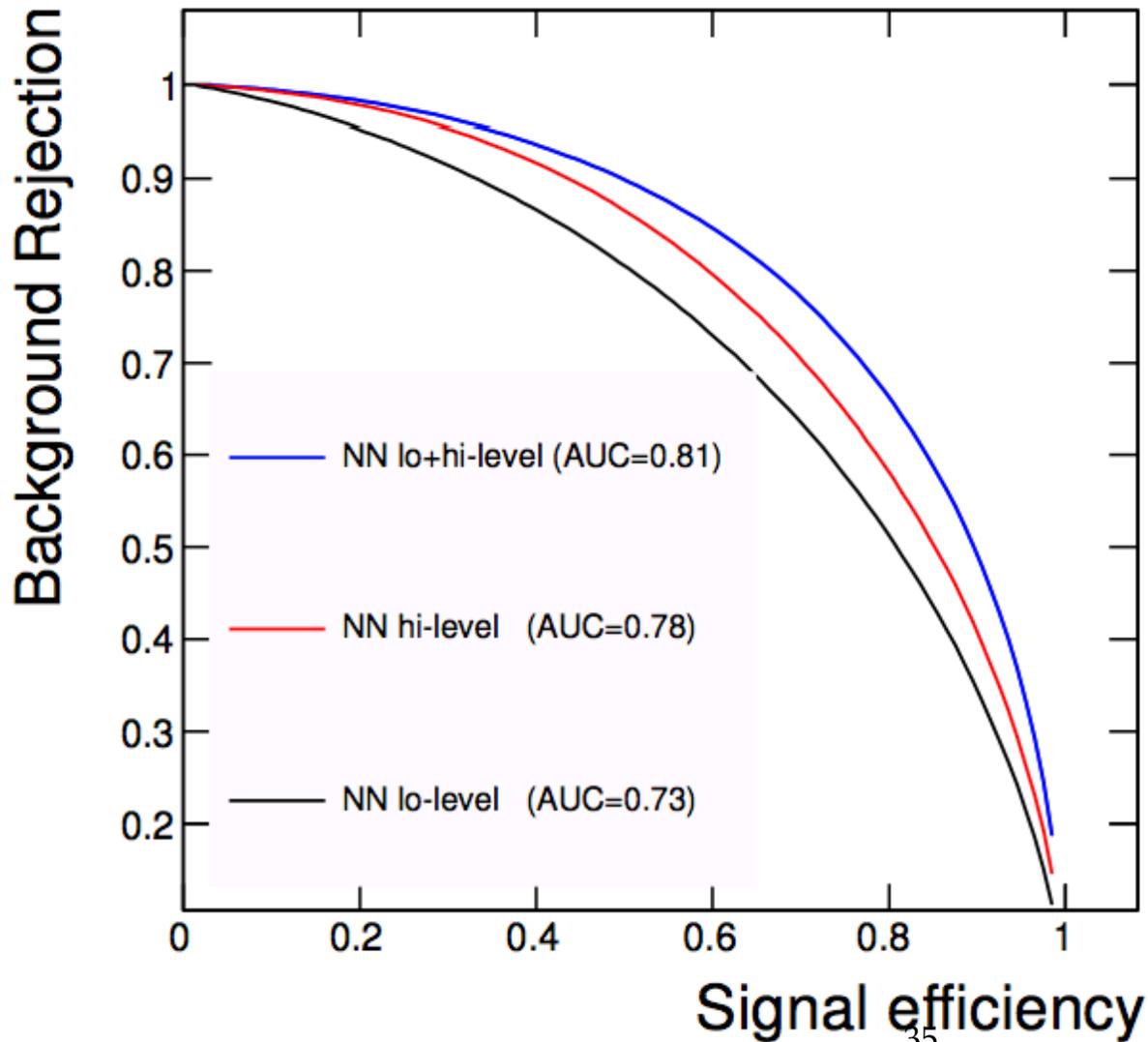
$m(lv)$

$m(blv)$





# Standard NNs



## Results

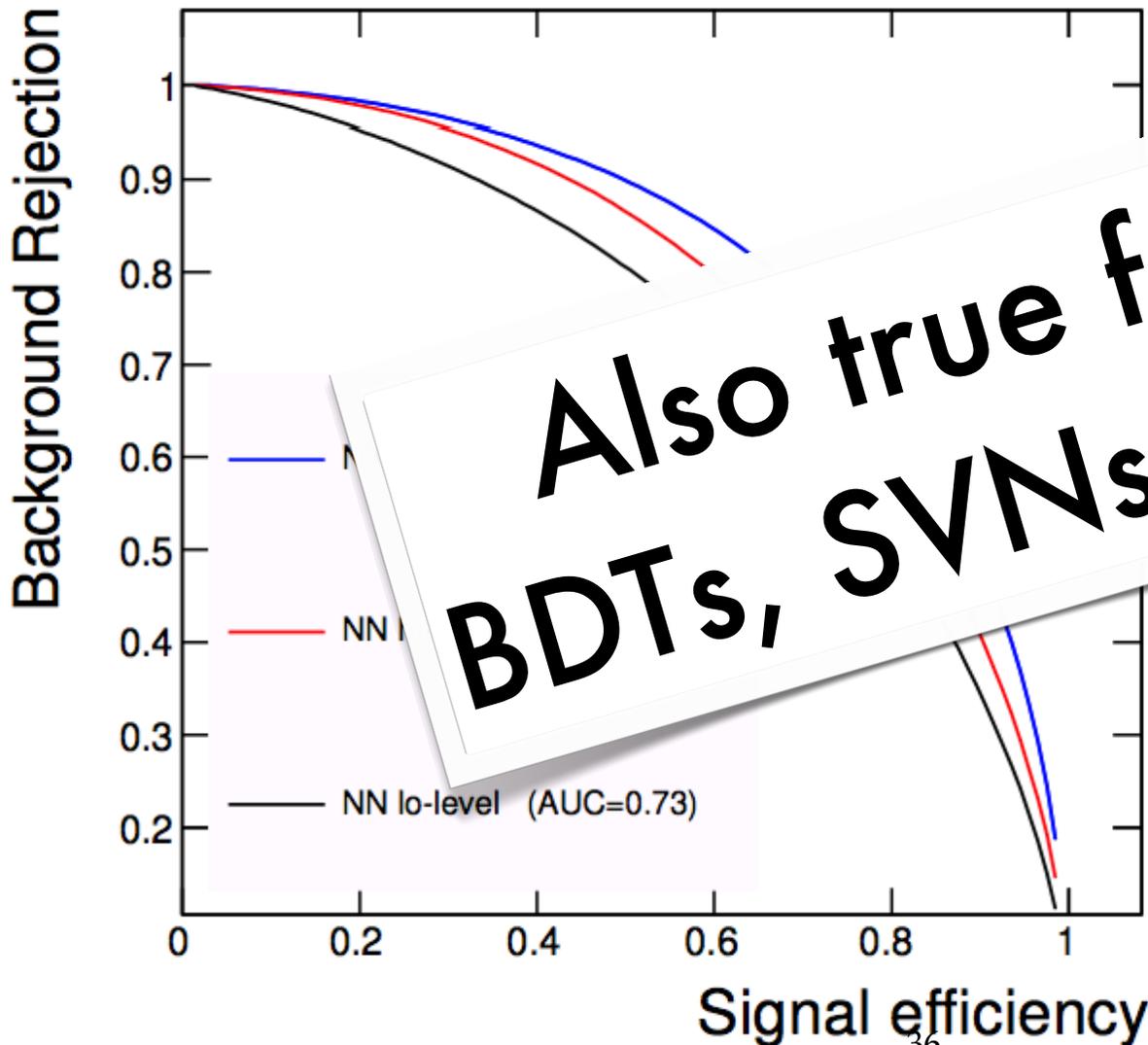
Adding hi-level  
boosts performance  
Better: lo+hi-level.

## Conclude:

NN can't find  
hi-level vars.

Hi-level vars  
do not have all info

# Standard NNs



Also true for  
BDTs, SVNs, etc

## Results

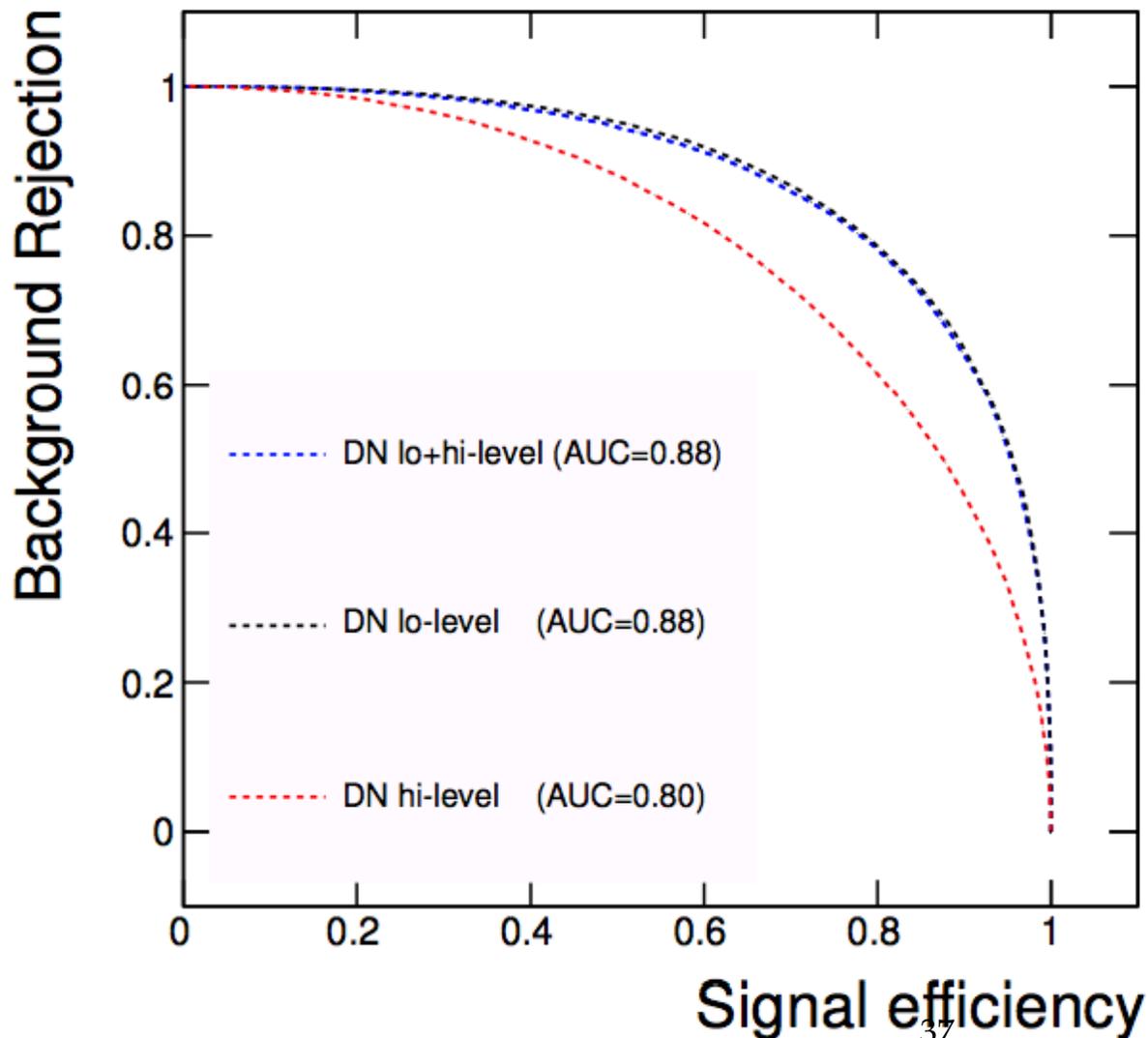
g hi-level  
performance  
lo+hi-level.

## include:

NN can't find  
hi-level vars.

Hi-level vars  
do not have all info

# Deep Networks



## Results

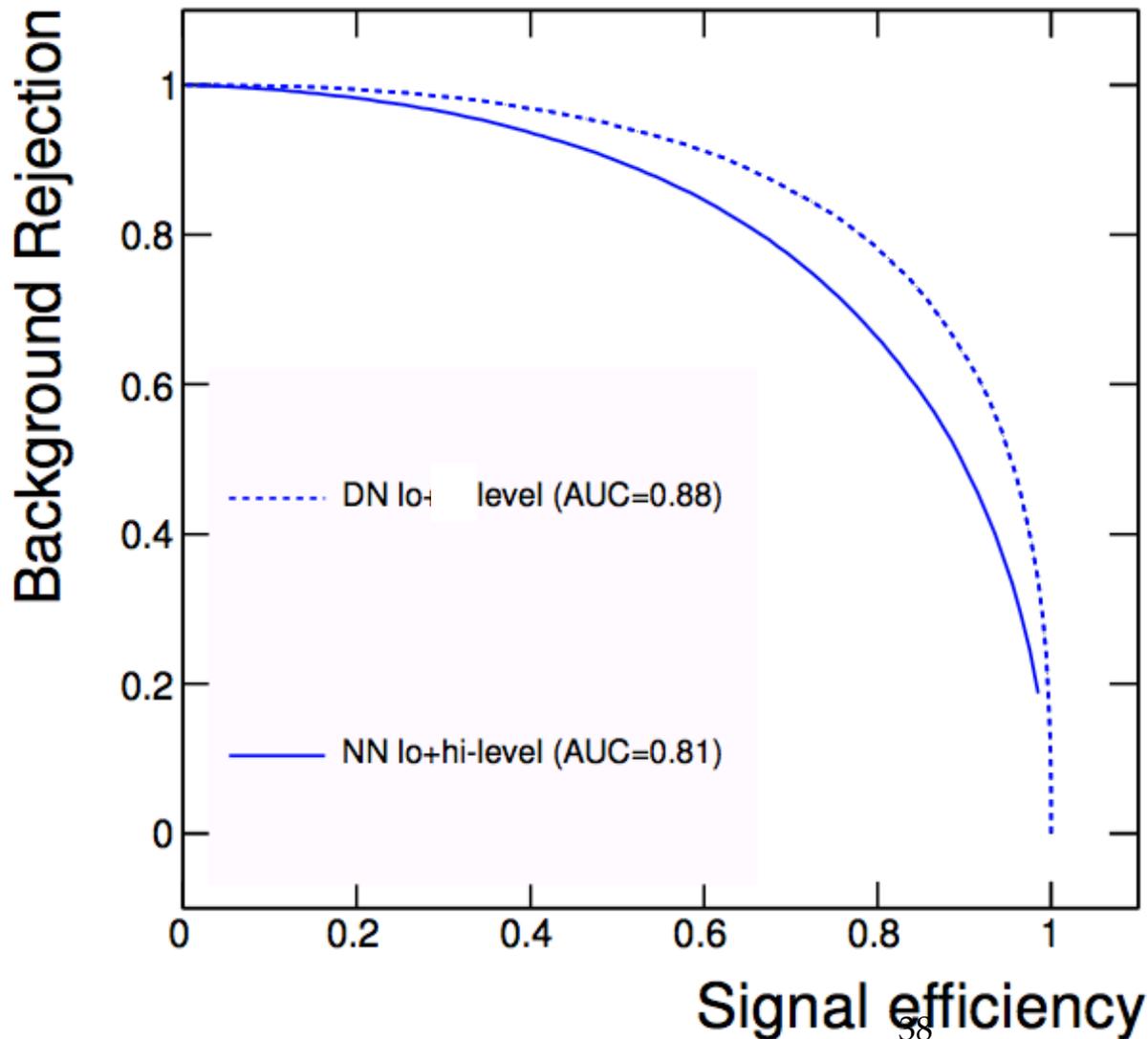
Lo+hi = lo.

## Conclude:

DN can find  
hi-level vars.

Hi-level vars  
do not have all info  
are unnecessary

# Deep Networks



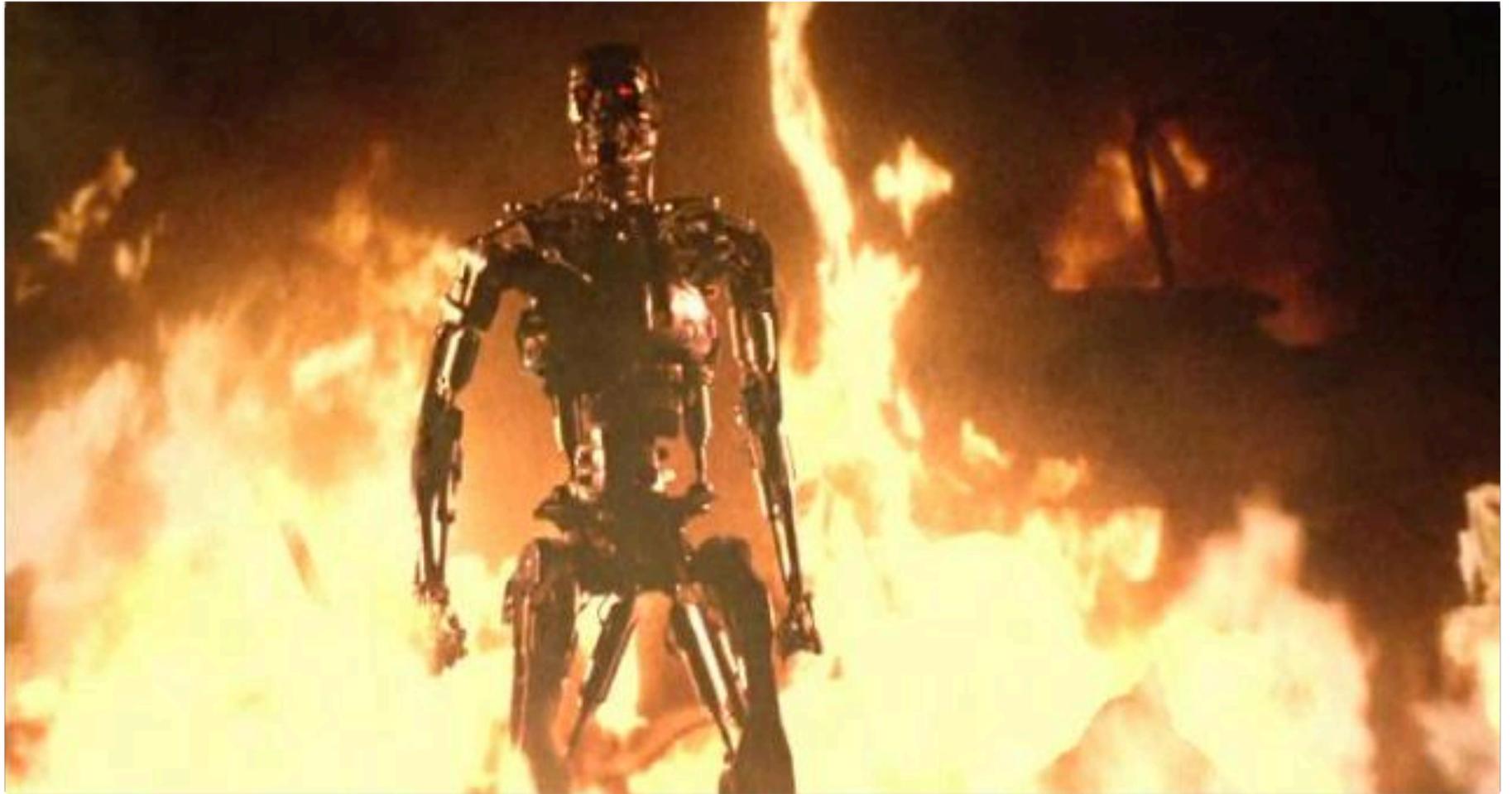
## Results

DN > NN

## Conclude:

DN does better  
than human  
assisted NN

# The Als win



# Results

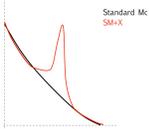
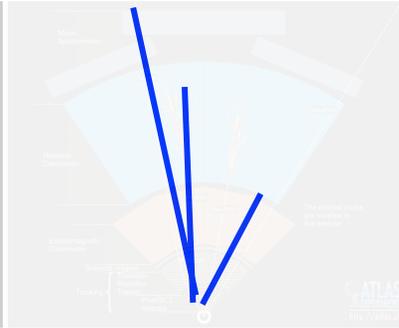
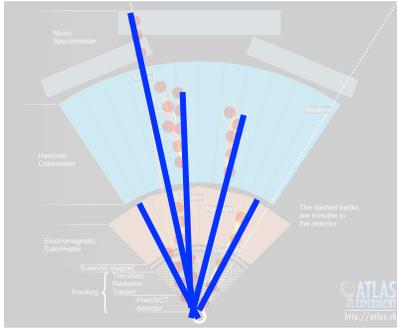
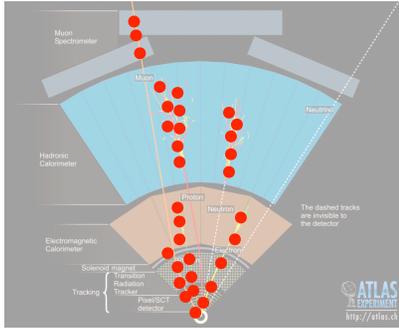
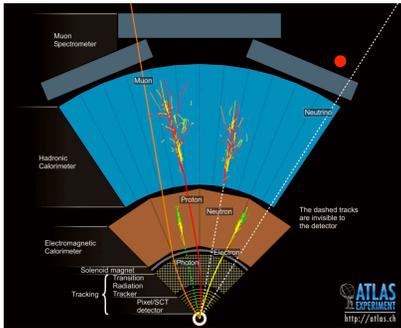
Identified example benchmark where traditional NNs fail to discover all discrimination power.

Adding human insight helps traditional NNs.

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

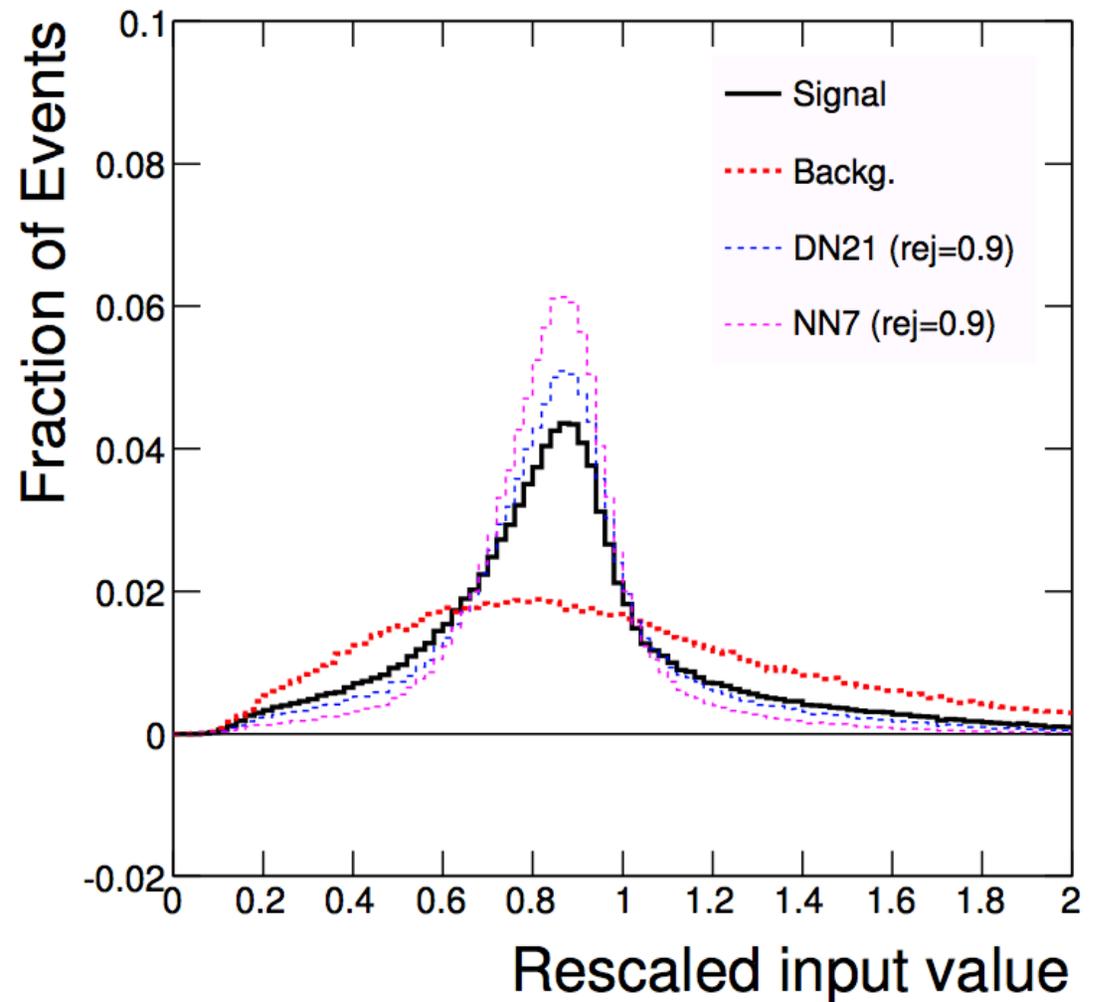
# Dimensionality

Raw	Sparsified	Reco	Select	Physics	Ana
$1e7$	$1e4$	100-ish*	50	10	1



# Why?

DN not as  
reliant on signal  
features. Cuts into  
background space.



# Depth

Supplementary Table 3: **Study of network size and depth.** Comparison of shallow networks with different numbers of hidden units (single hidden layer), and deep networks with varying hidden layers in terms of the Area Under the ROC Curve (AUC) for the HIGGS benchmark. The deep networks have 300 units in each hidden layer.

Technique	AUC		
	Low-level	High-level	Complete
NN 300-hidden	0.733	0.777	0.816
NN 1000-hidden	0.788	0.783	0.841
NN 2000-hidden	0.787	0.788	0.842
NN 10000-hidden	0.790	0.789	0.841
DN 3 layers	0.836	0.791	0.850
DN 4 layers	0.868	0.797	0.872
DN 5 layers	0.880	0.800	0.885
DN 6 layers	0.888	0.799	0.893

# Depth

Supplementary Table 3: **Study of network size and depth.** Comparison of shallow networks with different numbers of hidden units (single hidden layer), and deep networks with varying hidden layers in terms of the Area Under the ROC Curve (AUC) for the HIGGS benchmark. The deep networks have 300 units in each hidden layer.

Technique	AUC		
	Low-level	High-level	Complete
NN 300-hidden	0.733	0.777	0.816
NN 1000-hidden	0.788	0.783	0.841
NN 2000-hidden	0.787	0.788	0.842
NN 10000-hidden	0.790	0.789	0.841
DN 3 layers	0.836	0.791	0.850
DN 4 layers	0.868	0.797	0.872
DN 5 layers	0.880	0.800	0.885
DN 6 layers	0.888	0.799	0.893

**Diminishing gains after 4-5 layers**

# Depth

Supplementary Table 3: **Study of network size and depth.** Comparison of shallow networks with different numbers of hidden units (single hidden layer), and deep networks with varying hidden layers in terms of the Area Under the ROC Curve (AUC) for the HIGGS benchmark. The deep networks have 300 units in each hidden layer.

Technique	AUC		
	Low-level	High-level	Complete
NN 300-hidden	0.733	0.777	0.816
NN 1000-hidden	0.788	0.783	0.841
NN 2000-hidden	0.787	0.788	0.842
NN 10000-hidden	0.790	0.789	0.841
DN 3 layers	0.836	0.791	0.850
DN 4 layers	0.868	0.797	0.872
DN 5 layers	0.880	0.800	0.885
DN 6 layers	0.888	0.799	0.893

**Deep and thin > shallow and wide**

# Can we trust it?

## Black-box summary of N-dim space into 1-D

- no physics intuition available to critique
- possible reliance on poorly modeled corners
- impossible to by-hand validate N-dim corr.

# Can we trust it?

## Black-box summary of N-dim space into 1-D

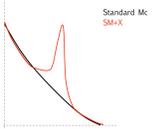
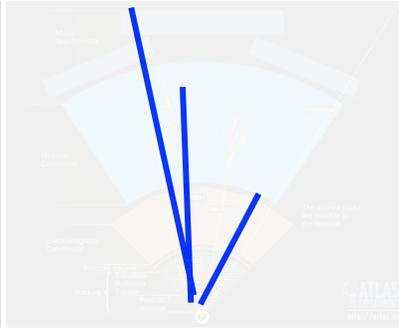
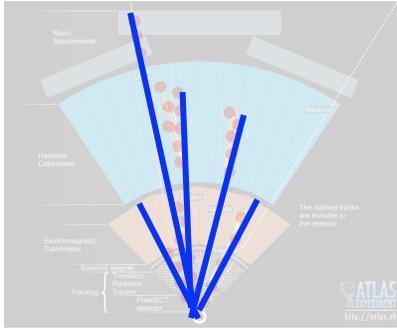
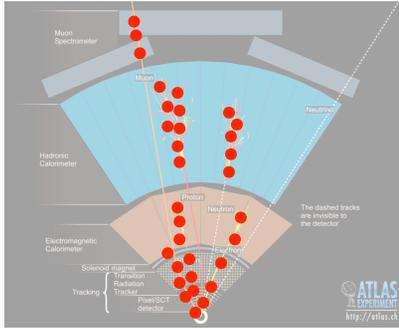
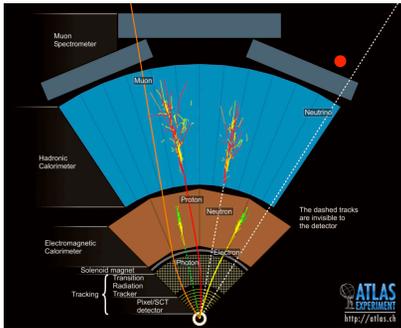
- no physics intuition available to critique
- possible reliance on poorly modeled corners
- impossible to by-hand validate N-dim corr.

## Problem is the same for shallow and deep nets

- only difference is non-linearity of function
- (1) validate inputs in 1D
  - (2) validate outputs in data with bg control regions

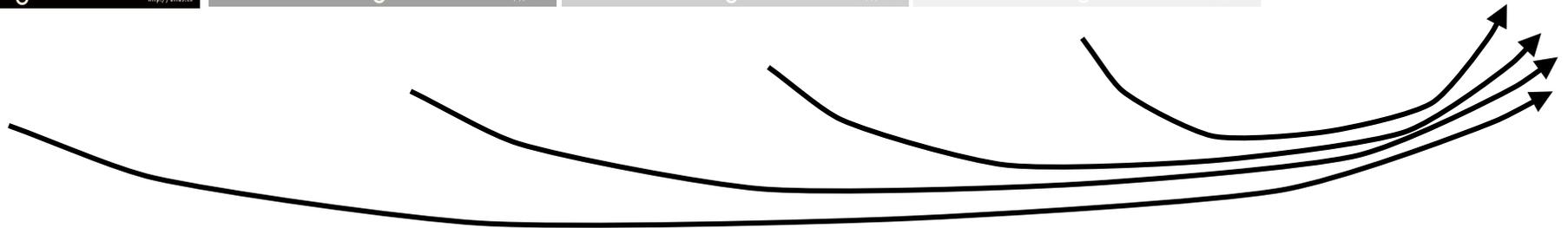
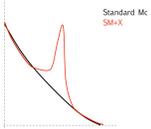
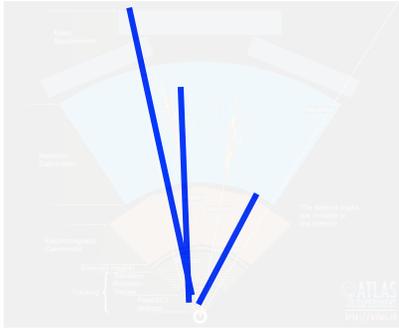
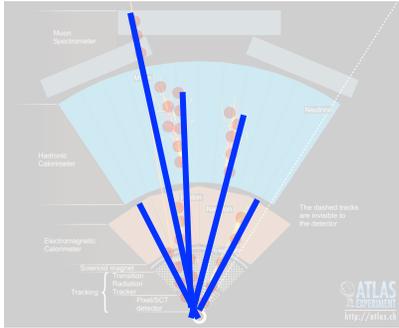
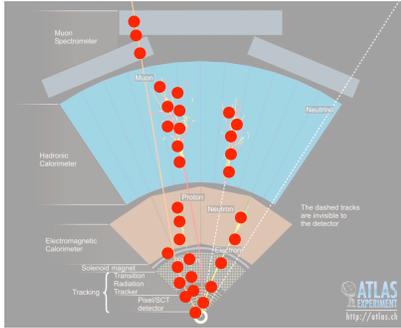
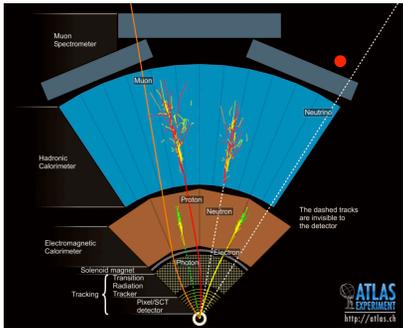
# What is possible?

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



# What is possible?

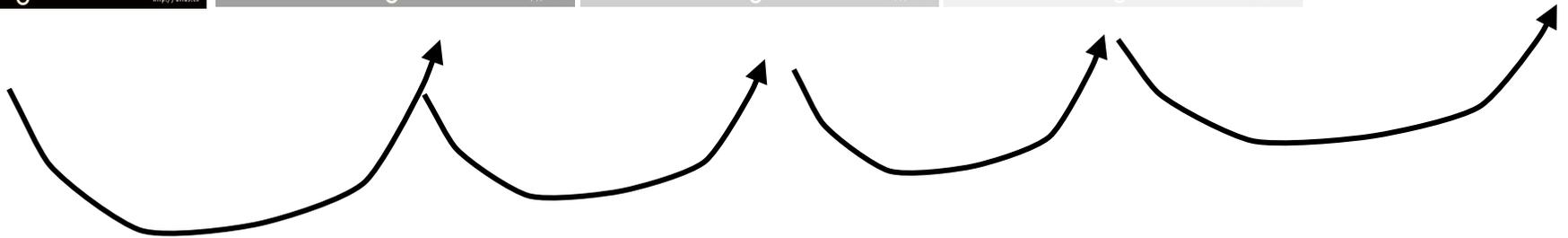
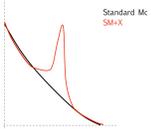
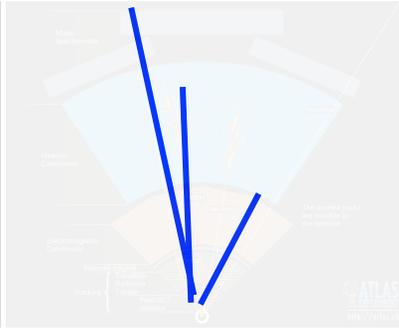
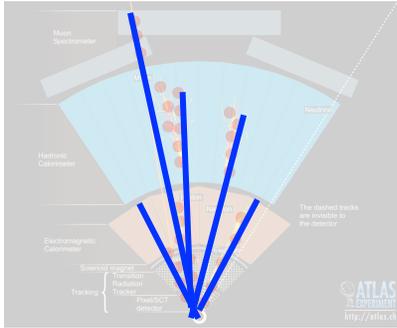
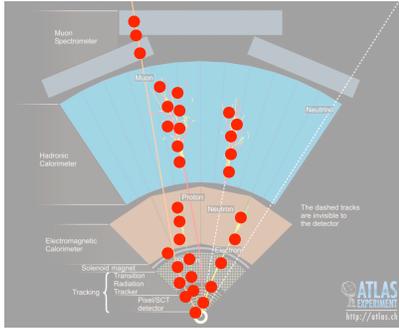
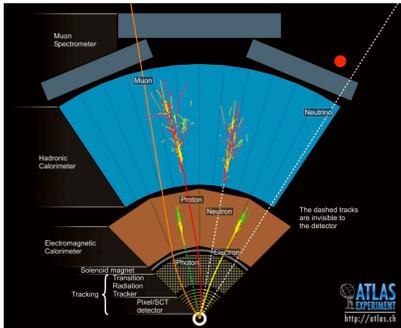
Raw	Sparsified	Reco	Select	Ana
$1e7$	$1e3$	100	50	1



Skip more steps with ML?

# Or this?

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



Improve each step with ML?

# Outline

## Defining the problem

Dimensionality reduction

## Current approaches

Deep(er) networks for low(er)-level data

## Dimensionality increase

Parameterized networks for sets of problems

# Larger Scope

1.

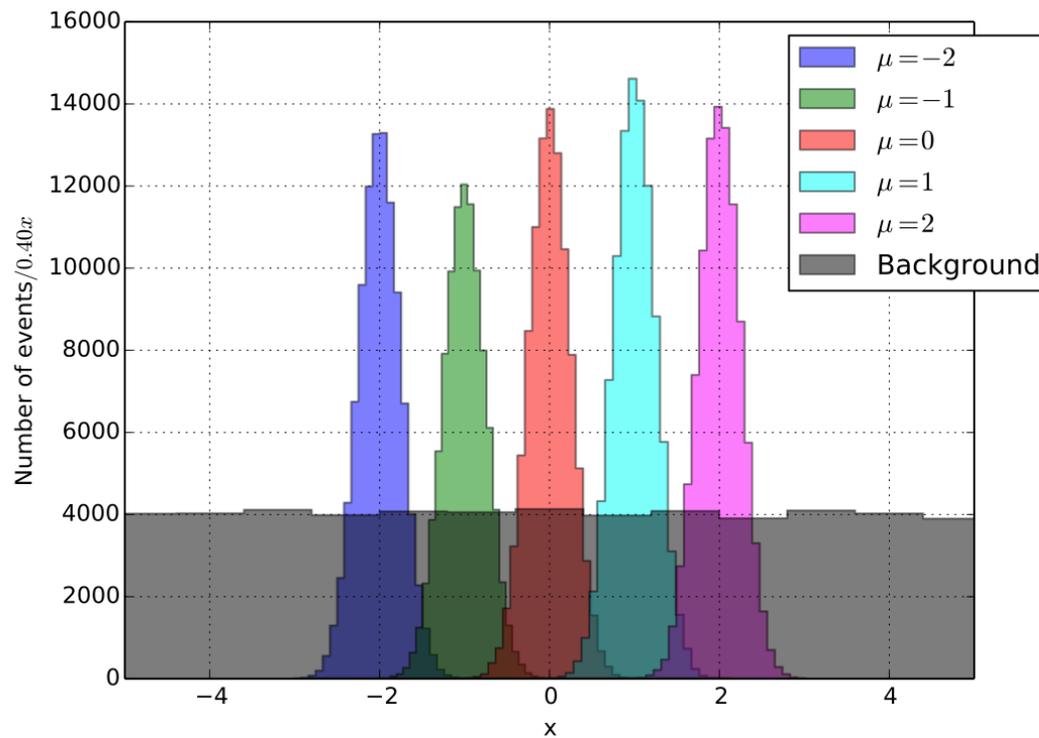
THE STANDARD MODEL			
Fermions			
Quarks	$u$ up	$c$ charm	$t$ top
	$d$ down	$s$ strange	$b$ bottom
Leptons	$\nu_e$ electron neutrino	$\nu_\mu$ muon neutrino	$\nu_\tau$ tau neutrino
	$e$ electron	$\mu$ muon	$\tau$ tau

2.

THE STANDARD MODEL PLUS X				
Fermions				
Quarks	$u$ up	$c$ charm	$t$ top	$X$
	$d$ down	$s$ strange	$b$ bottom	
Leptons	$\nu_e$ electron neutrino	$\nu_\mu$ muon neutrino	$\nu_\tau$ tau neutrino	
	$e$ electron	$\mu$ muon	$\tau$ tau	

What if you don't know the mass of X?

# Related problems



You might have a set of closely related problems

# Current approaches

1. Develop unique solution at each point

Pros: optimality

Cons: cannot interpolate

2. Develop single solution for all points

Pros: can interpolate

Cons: not optimal anywhere

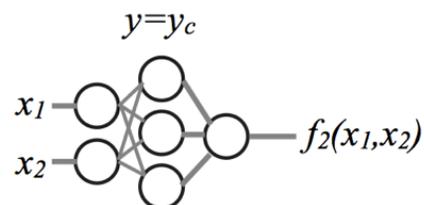
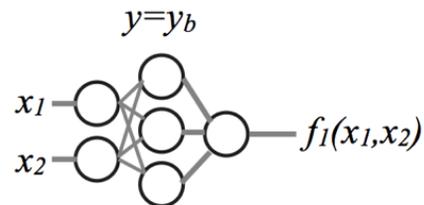
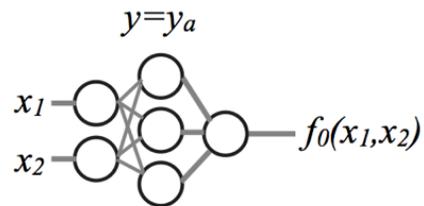
3. Develop single solution at one point, use everywhere

Pros: can interpolate, some optimality

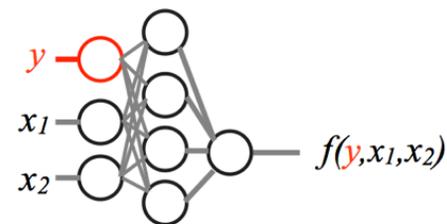
Cons: not optimal everywhere

# New approach

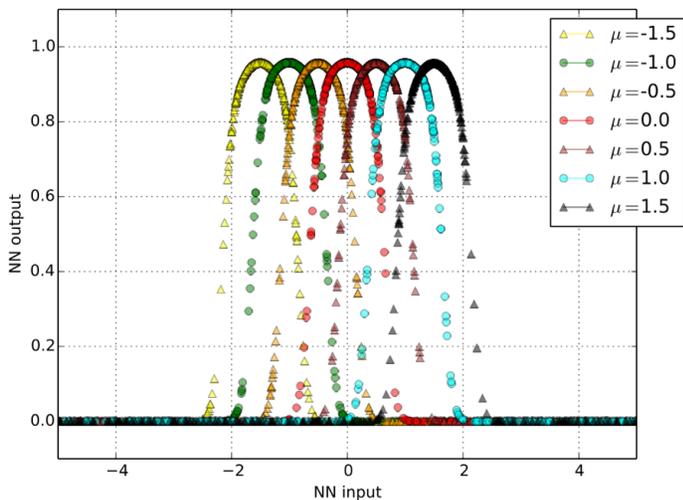
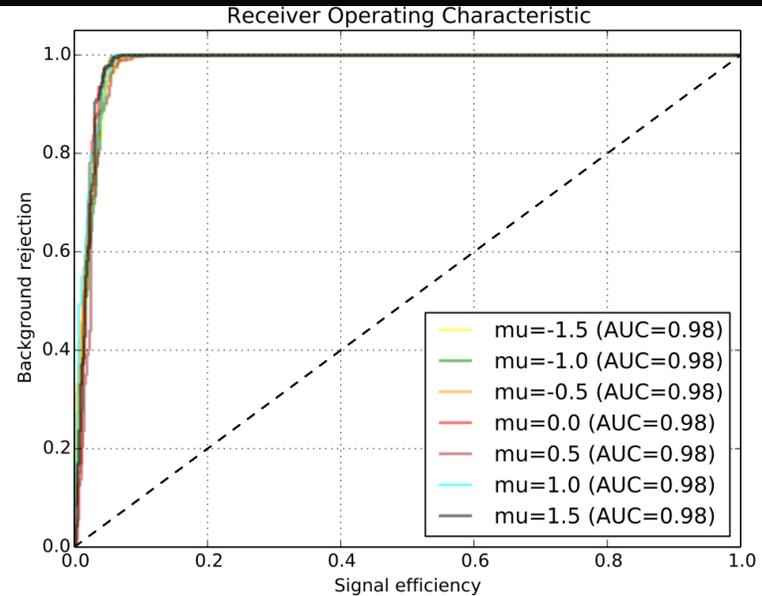
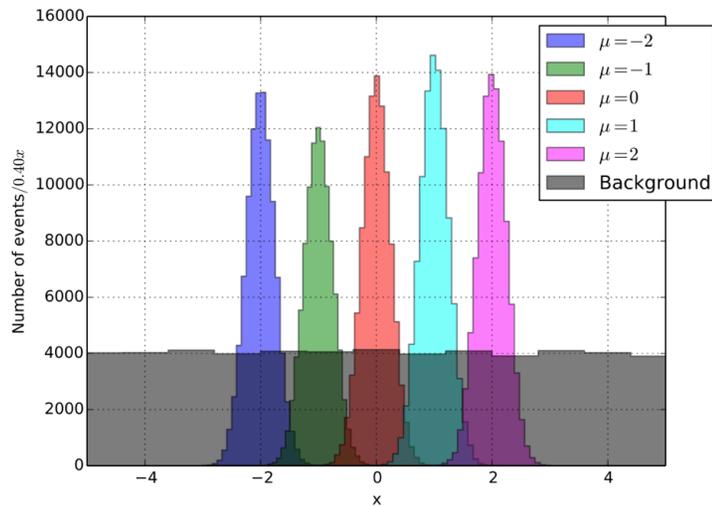
Replace **set** of solutions....



.....with single **parameterized** solution



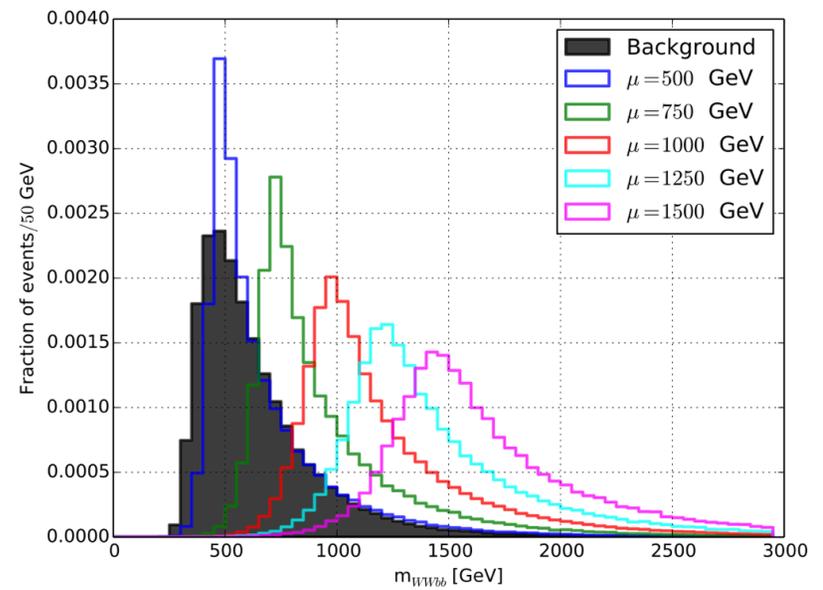
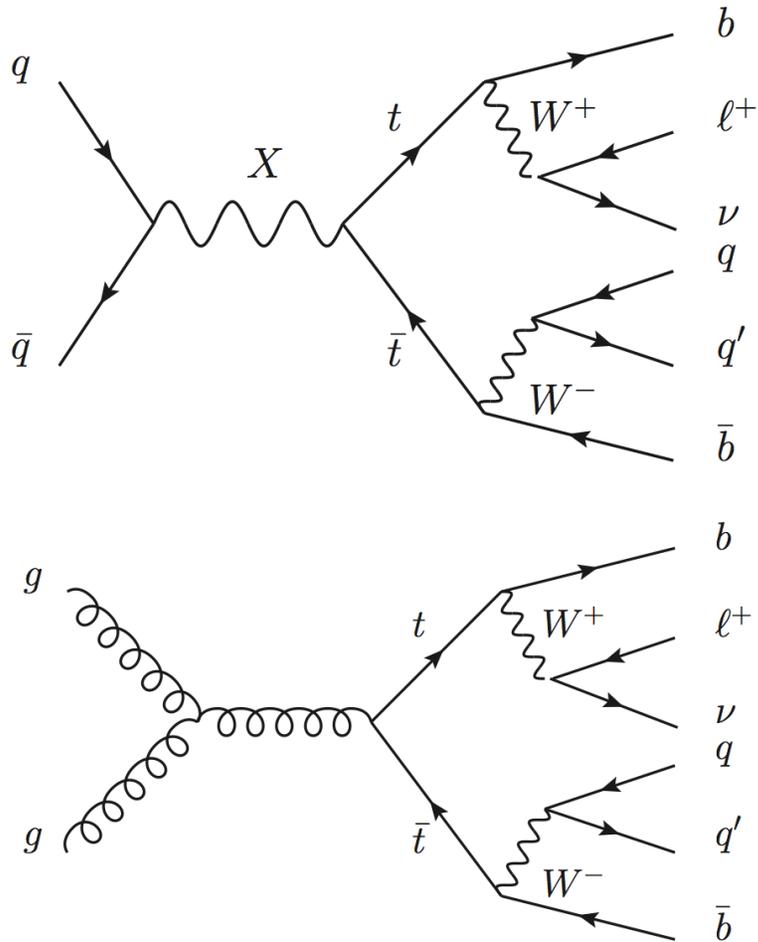
# Toy example 1



NN generalized solution.

Good performance at values of parameter with no training!

# Half-toy example

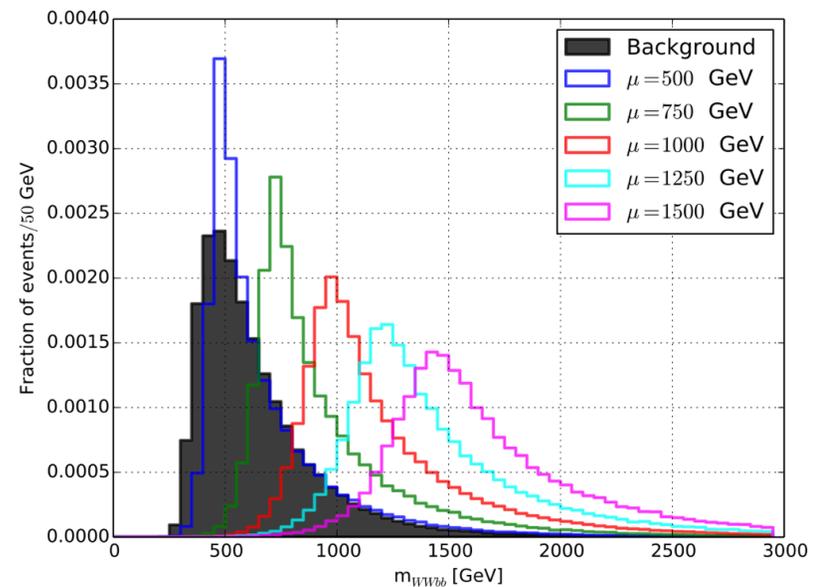
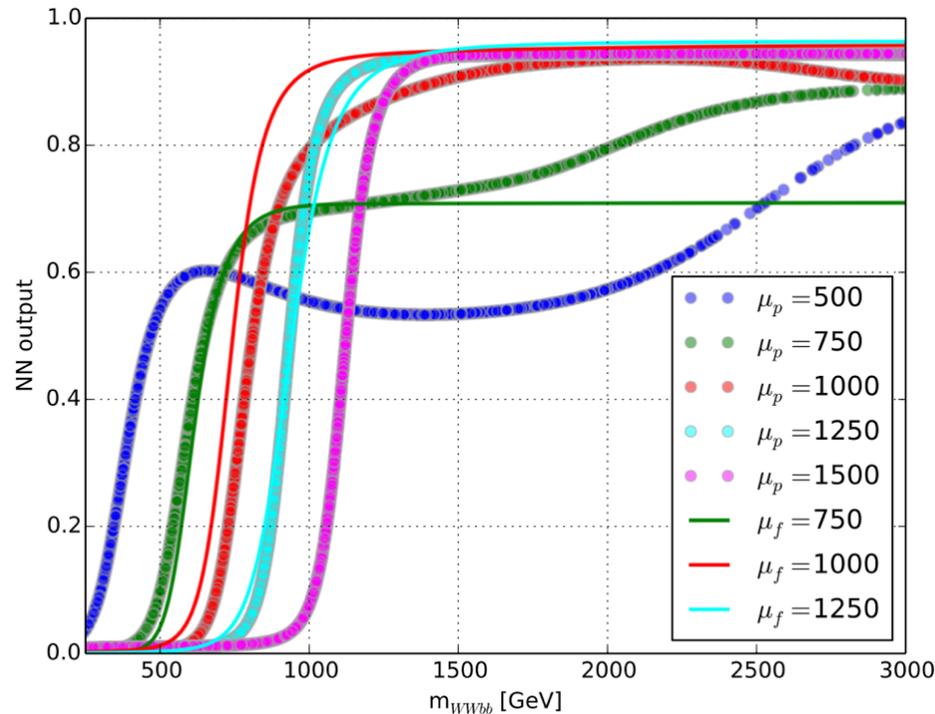


# Half-toy example

## Compare NN response

Single network **trained at  $\mu$**

Param network **not trained at  $\mu$**

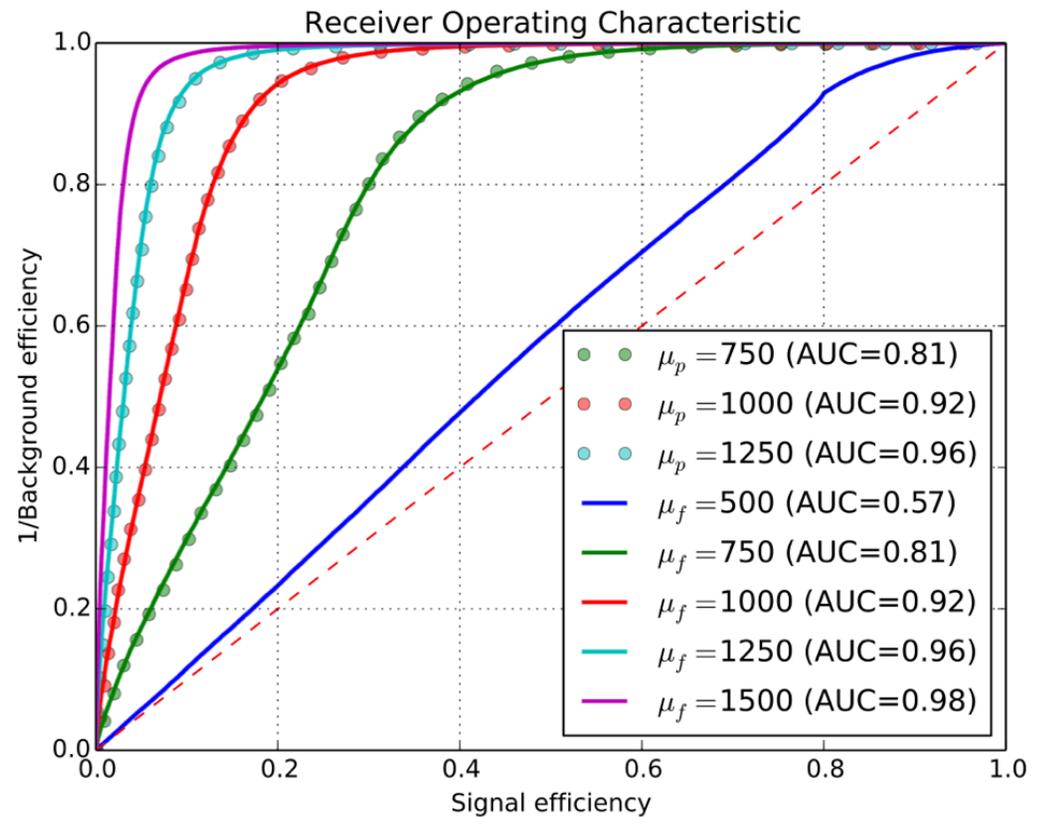
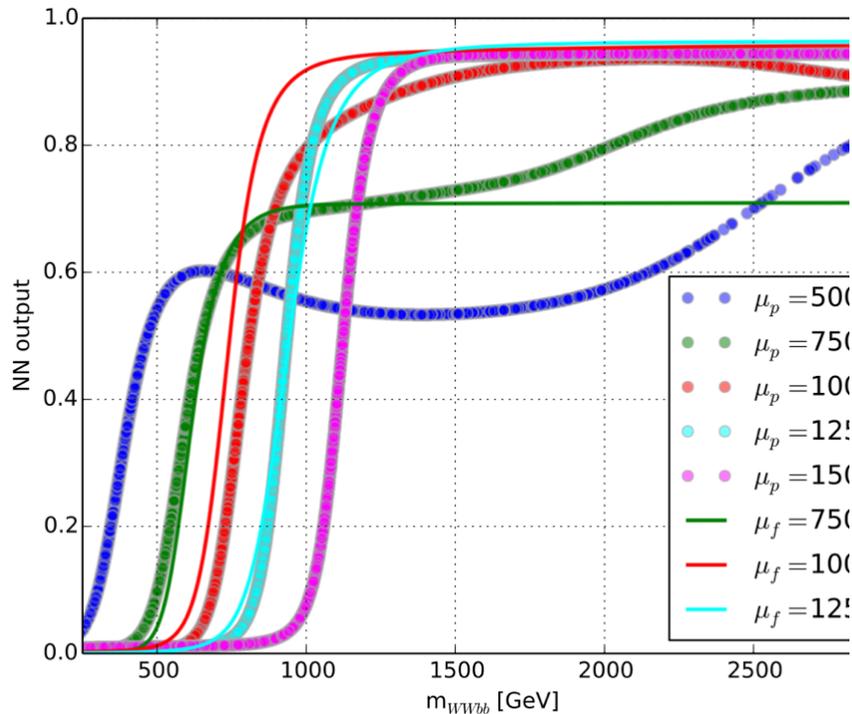


# Half-toy example

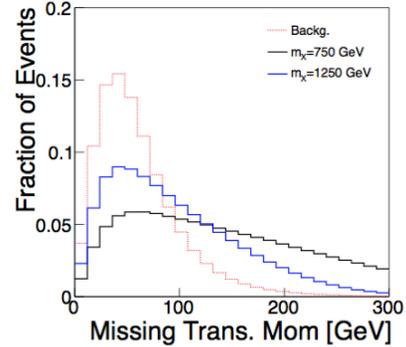
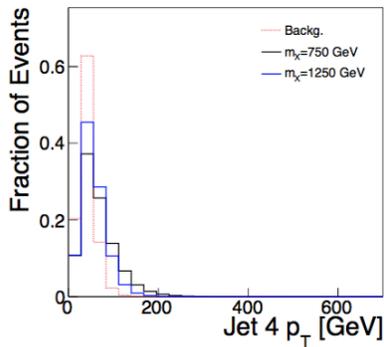
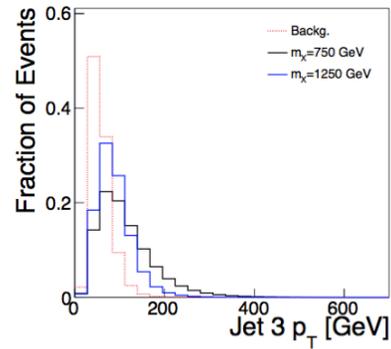
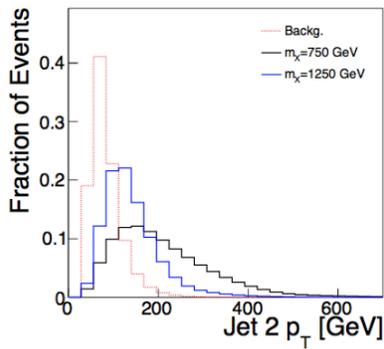
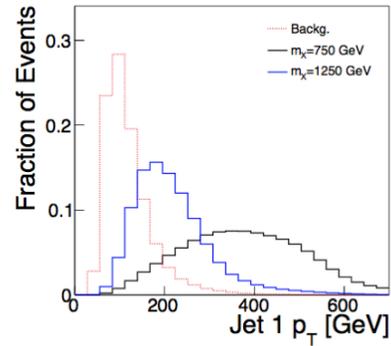
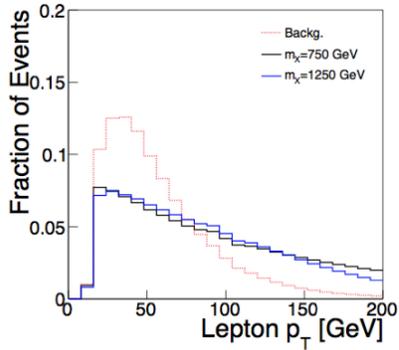
## Compare NN discrimination

Single network **trained at  $\mu$**

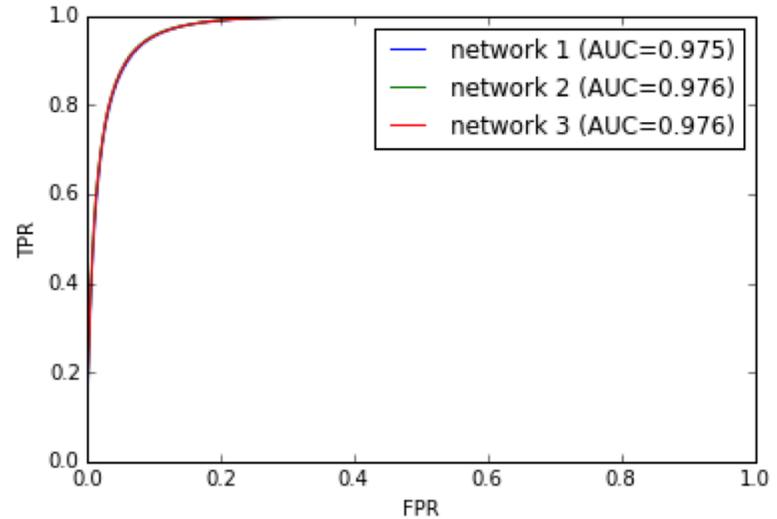
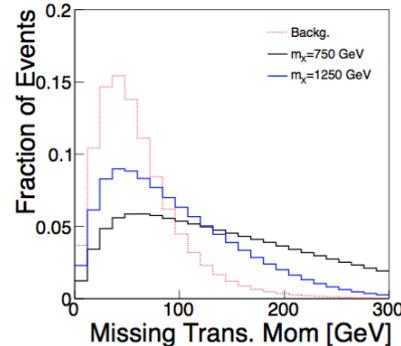
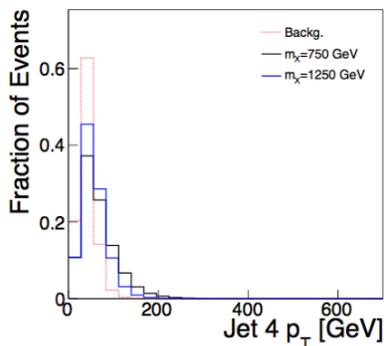
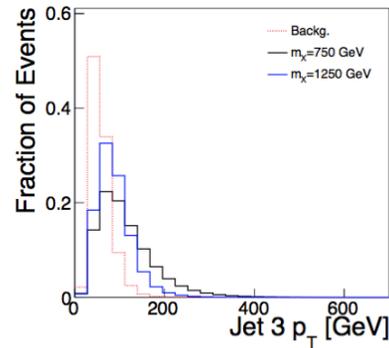
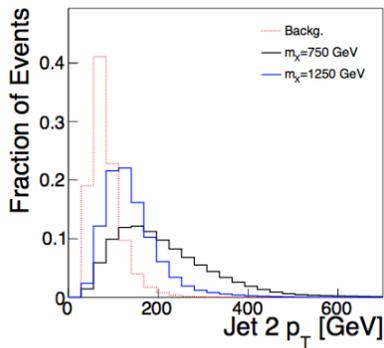
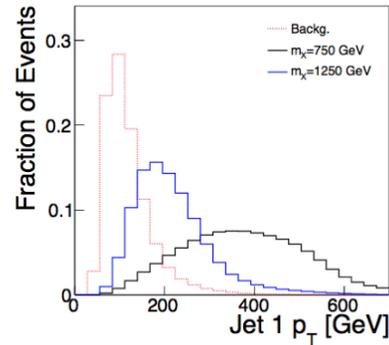
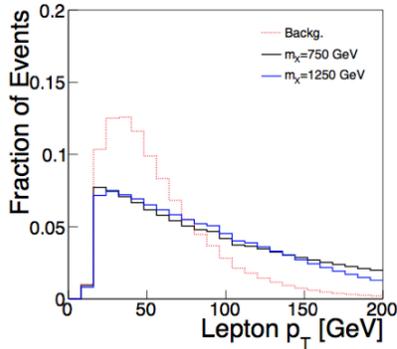
Param network **not trained at  $\mu$**



# Full example



# Full example



**N1: param network**  
trained at  $m=500,750,1250,1500$

**N2: fixed network**  
trained at  $m=1000$

**N3: param network**  
trained at  $m=500,750,1000,1250,1500$

# Discussion

We have an intractable high-dimensional problem of calculating:

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

Human tricks to reduce dimensionality are **non-optimal**.

Gains possible at many levels by **replacing brains with AI**.