

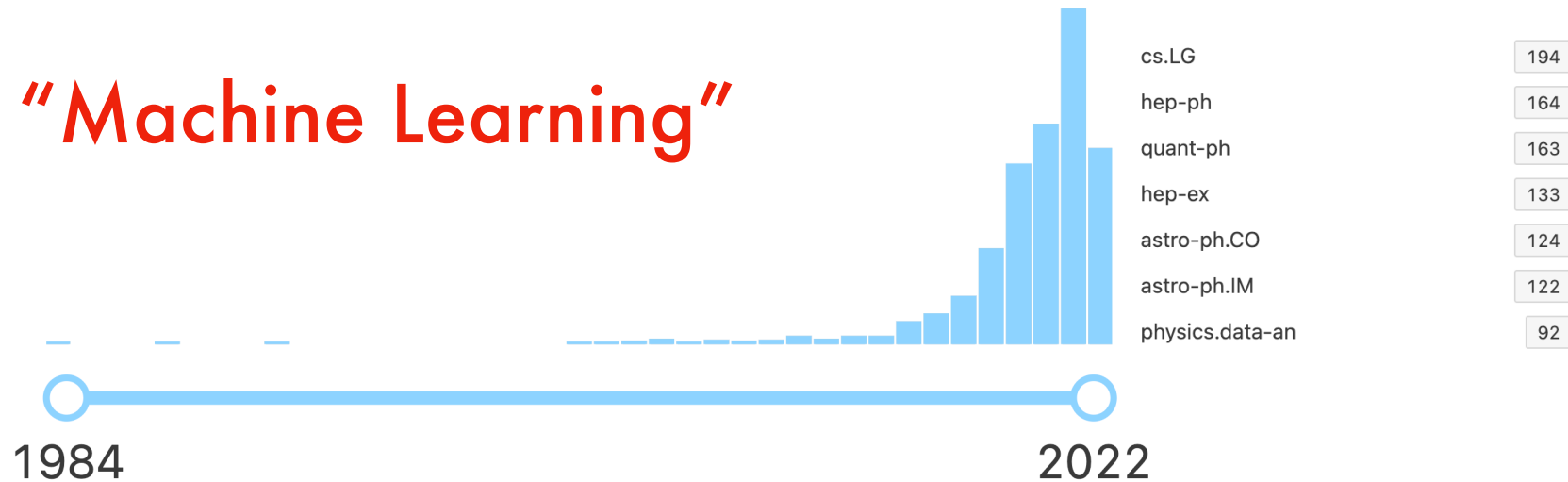
What's new about Machine Learning?



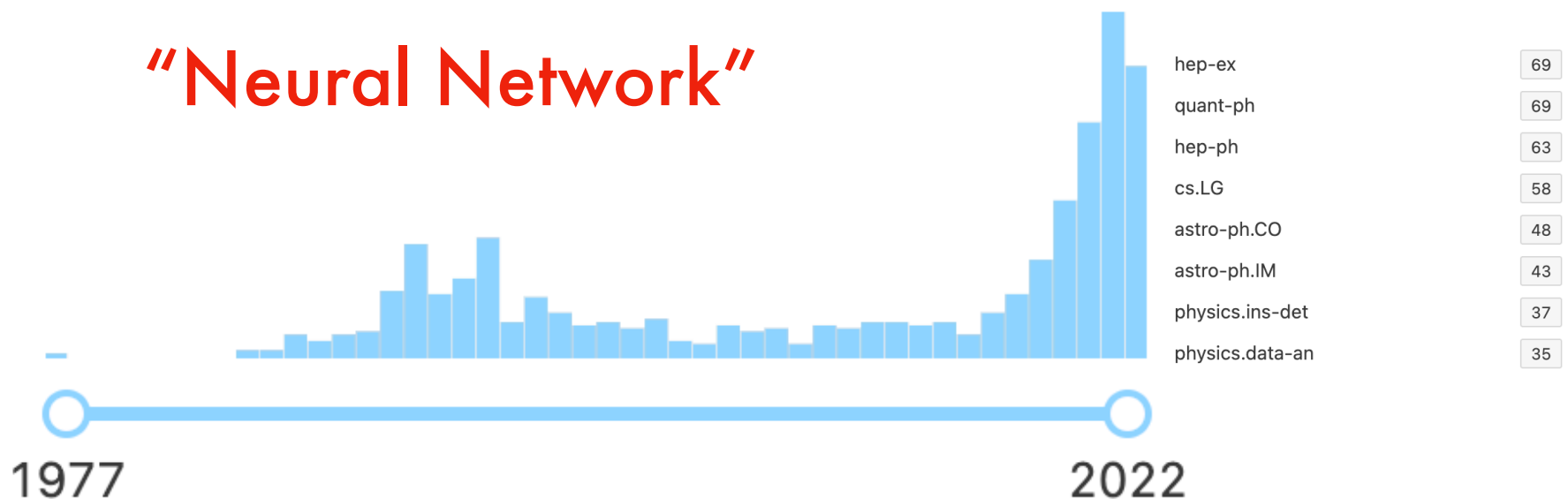
Daniel Whiteson, UC Irvine
Jul 2022 / Snowmass in Seattle

It's everywhere!

"Machine Learning"



"Neural Network"



What's new about ML?



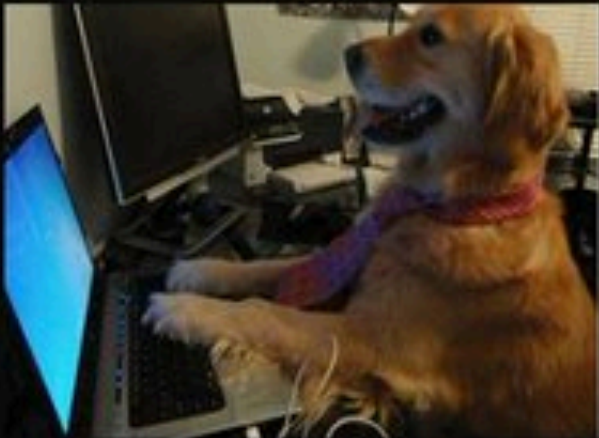
What society thinks I do



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

Early days of HEP



ML in HEP is not new

1997:

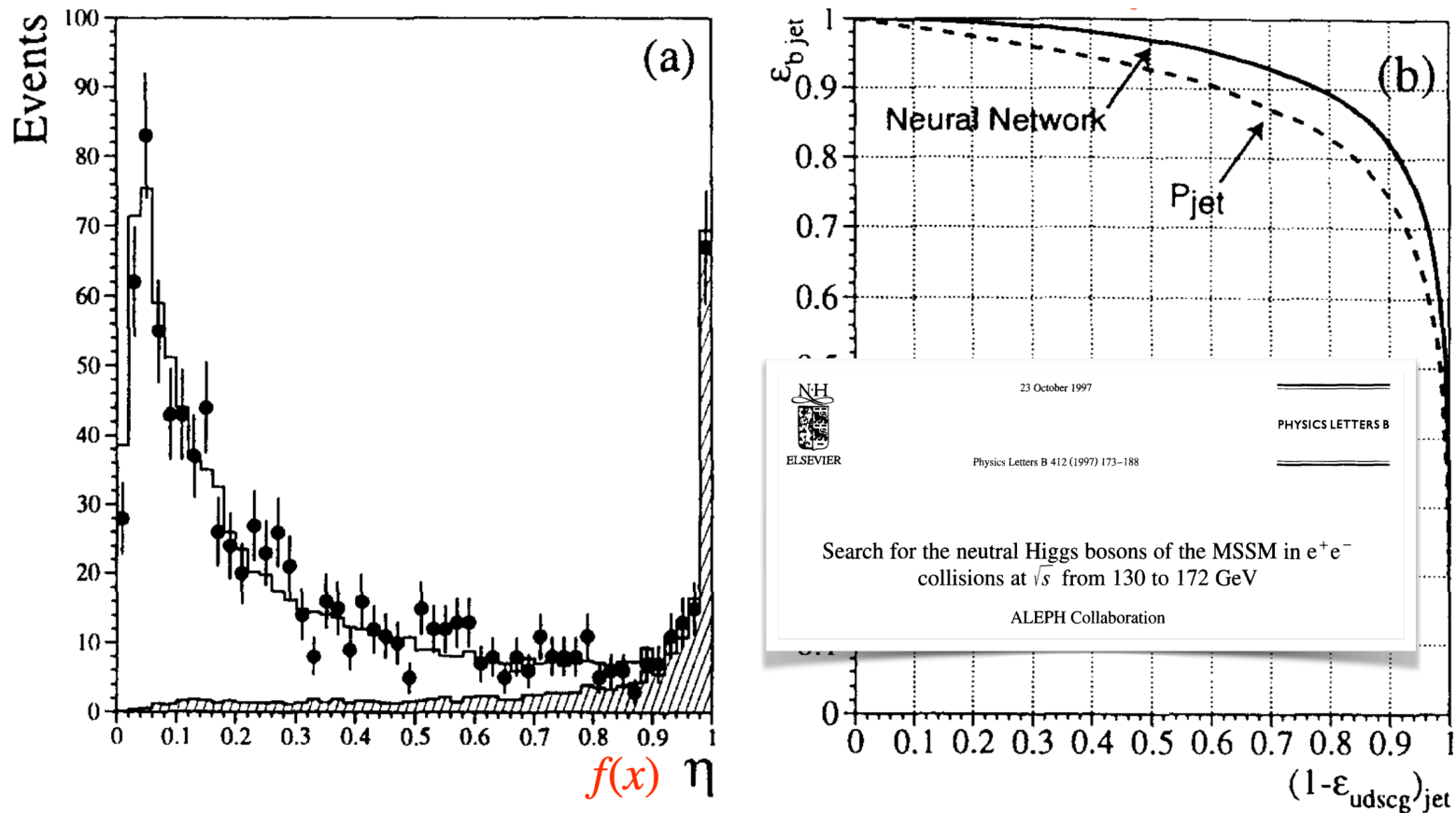


Fig. 2. (a) The output η of the neural network b tag for radiative returns to the Z for 161 GeV $q\bar{q}$ Monte Carlo (histogram) compared to the data at 161 GeV (points). The shaded region shows the contribution from generated b-jets. (b) The performance of the neural network b tag (solid line) for Monte Carlo events, presented in terms of the efficiency for identifying b-jets versus the efficiency for rejecting light quark jets. The performance of the single most powerful b tagging input variable to the neural network is shown for comparison (dashed curve).

Is it something new?

Is modern ML **something new**,
or just **more of the same**?

Is it something new?

Is modern ML something new,
or just more of the same?



Daniel Whiteson
@DanielWhiteson

...

Is recent (> ~2013 deep learning moment) ML in particle physics "more of the same" or "qualitatively something new".

Is it something new?

Is modern ML something new,
or just more of the same?



Daniel Whiteson
@DanielWhiteson

...

Is recent (> ~2013 deep learning moment) ML in particle physics "more of the same" or "qualitatively something new".

More of the same	39.7%
More, not the same	39.7%
It's complicated(comment)	11%
ML is nonsense	9.6%

73 votes · Final results

Is it something new?

Is modern ML something new,
or just more of the same?



Daniel Whiteson
@DanielWhiteson

...

Is recent (> ~2013 deep learning moment) ML in
particle physics "more of the same" or "qualitatively
something new".



73 votes · Final results

Outline

1. **Much much** more of the same
2. Something **qualitatively new**

Traditional role of ML

Why do we need machine learning?

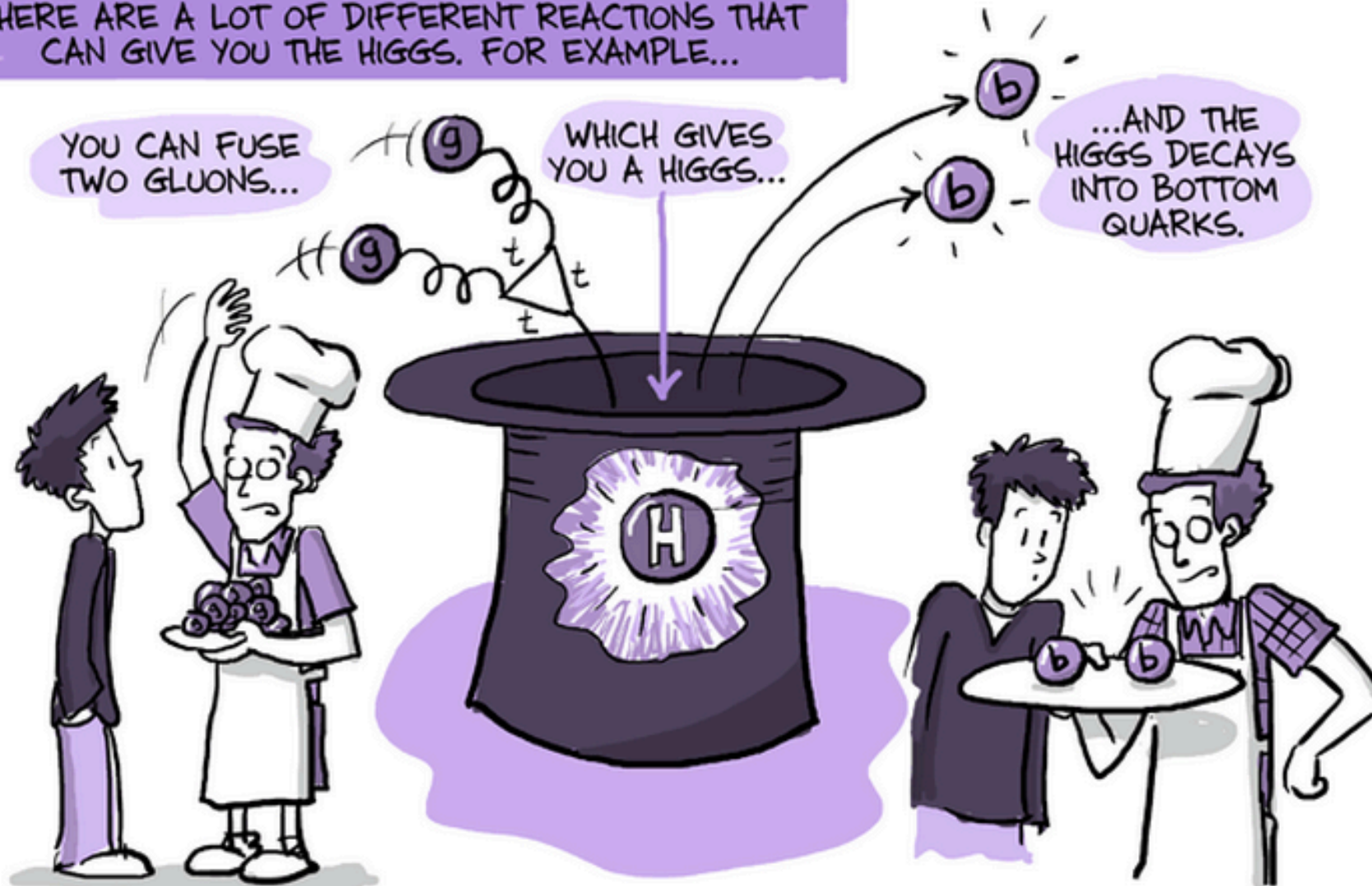
Traditional role of ML

Why do we need machine learning?



Making a new particle

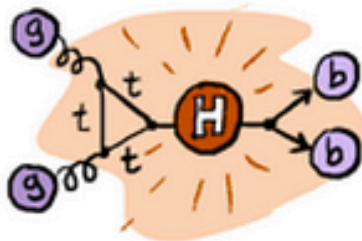
THERE ARE A LOT OF DIFFERENT REACTIONS THAT CAN GIVE YOU THE HIGGS. FOR EXAMPLE...



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

Neyman-Pearson

NP lemma says that the best statistic is the **likelihood ratio**:

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

A diagram illustrating the components of the likelihood ratio formula. A blue arrow points from the word "data" to the variable x in the denominator $P(x|H_0)$. A red arrow points from the word "theory" to the hypothesis H_0 in the denominator $P(x|H_0)$.

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

Functional space

All functions

*Global
Optimum*



No problem

If you can calculate:

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

For which you need:

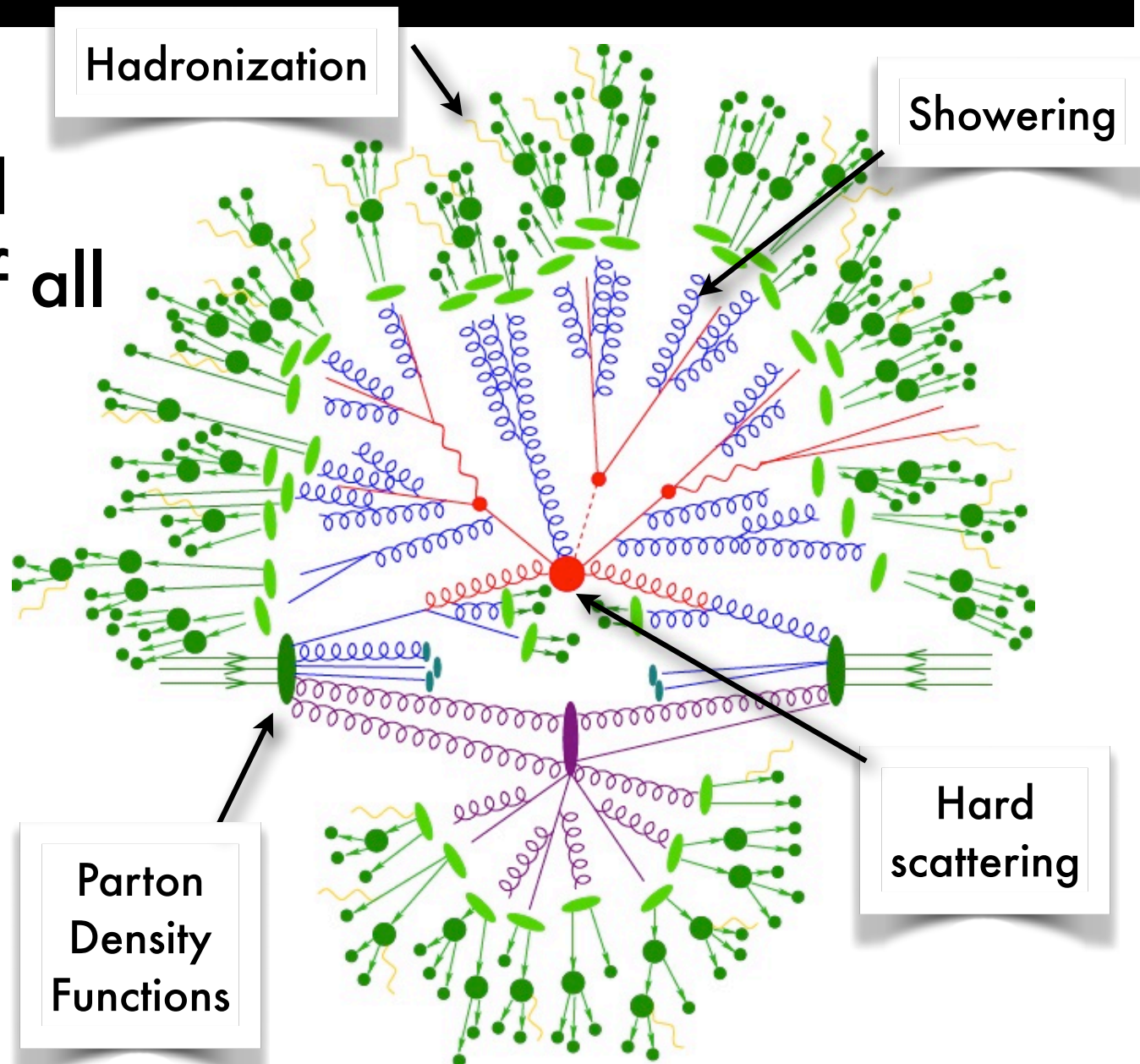
$$P(\text{data} | \text{theory})$$

In general

We have a good understanding of all of the pieces

Do we have

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



In general

Hadronization

Showering

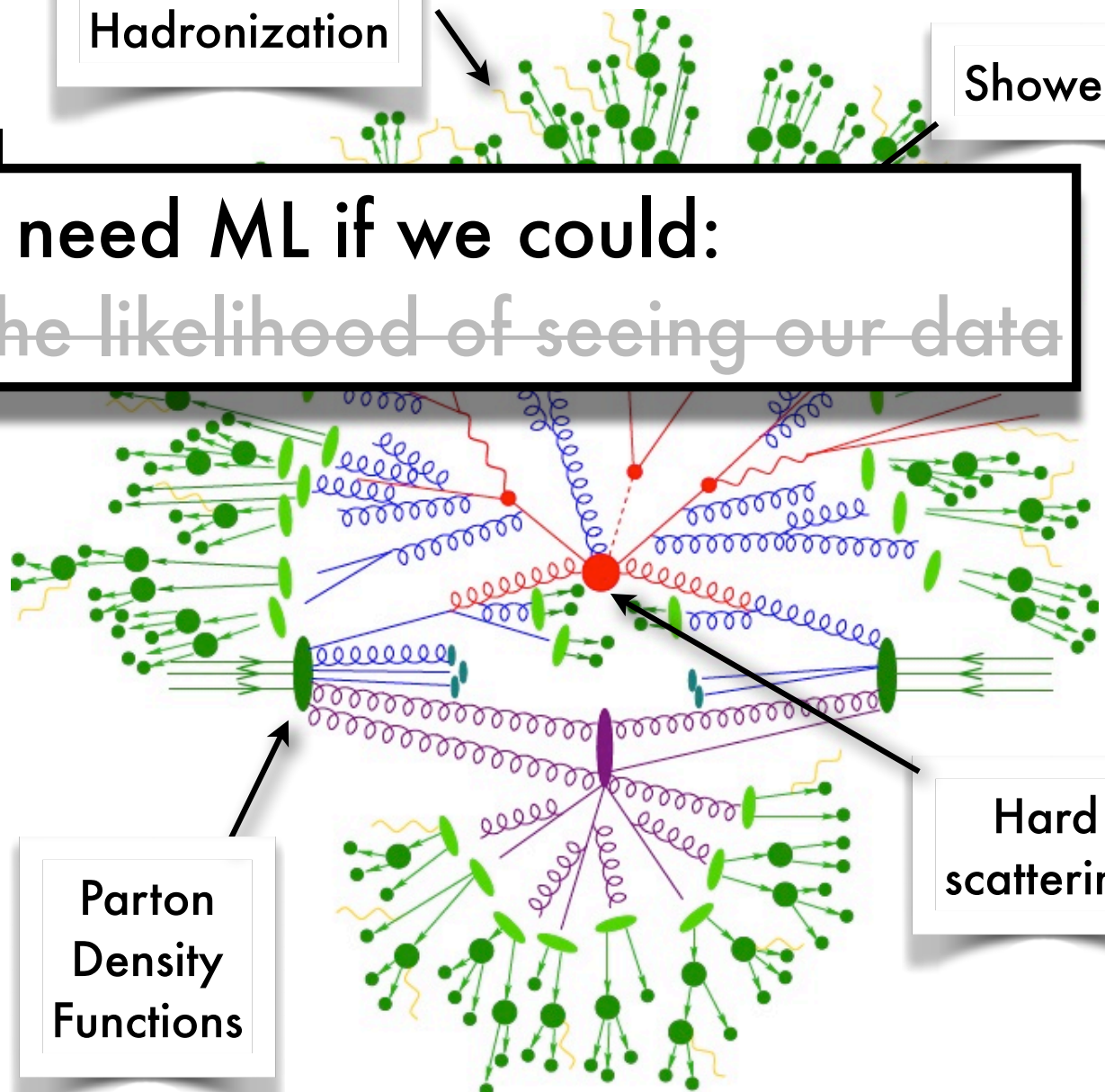
We have a good
understanding
of the

We wouldn't need ML if we could:

- ~~Express the likelihood of seeing our data~~

Do we have

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

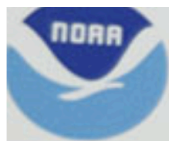


Darn

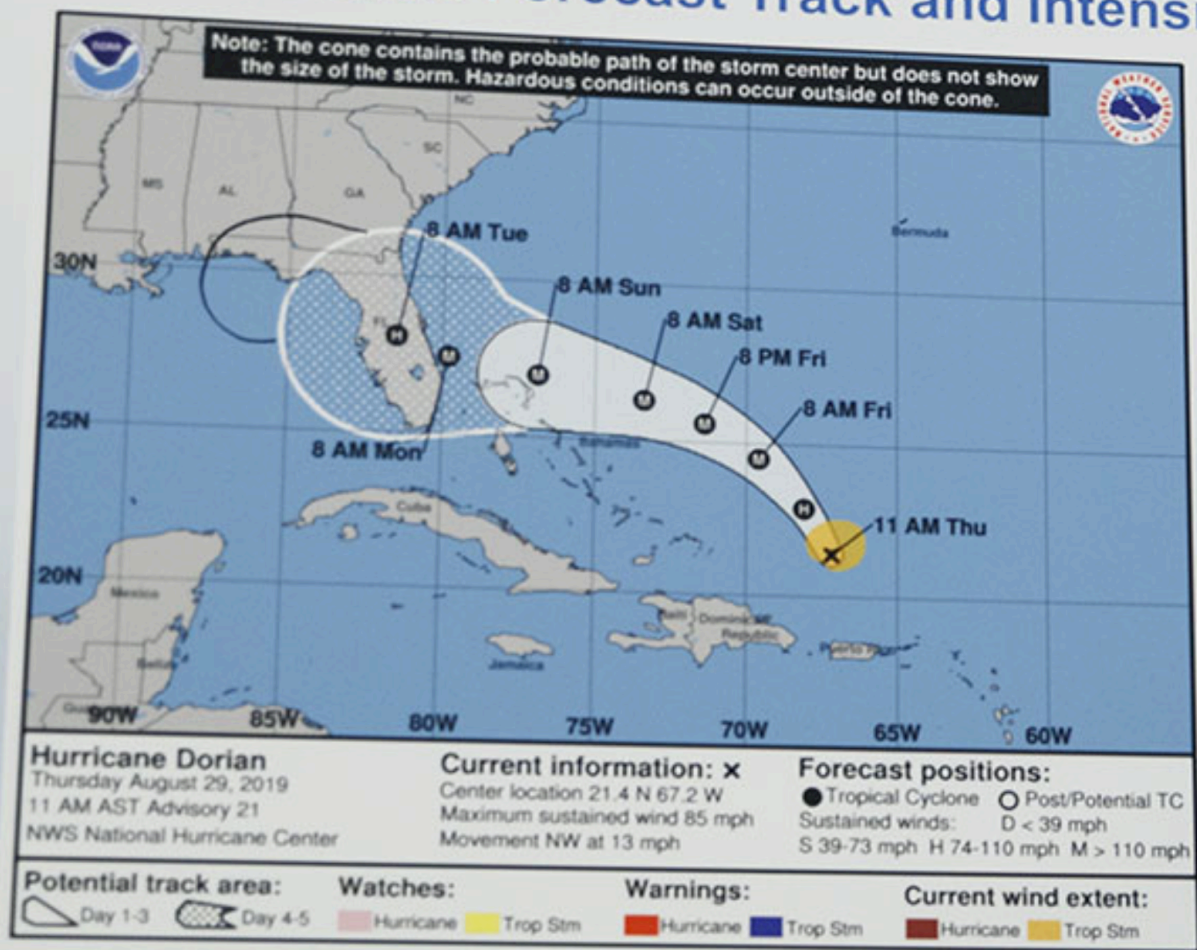
We can't calculate

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

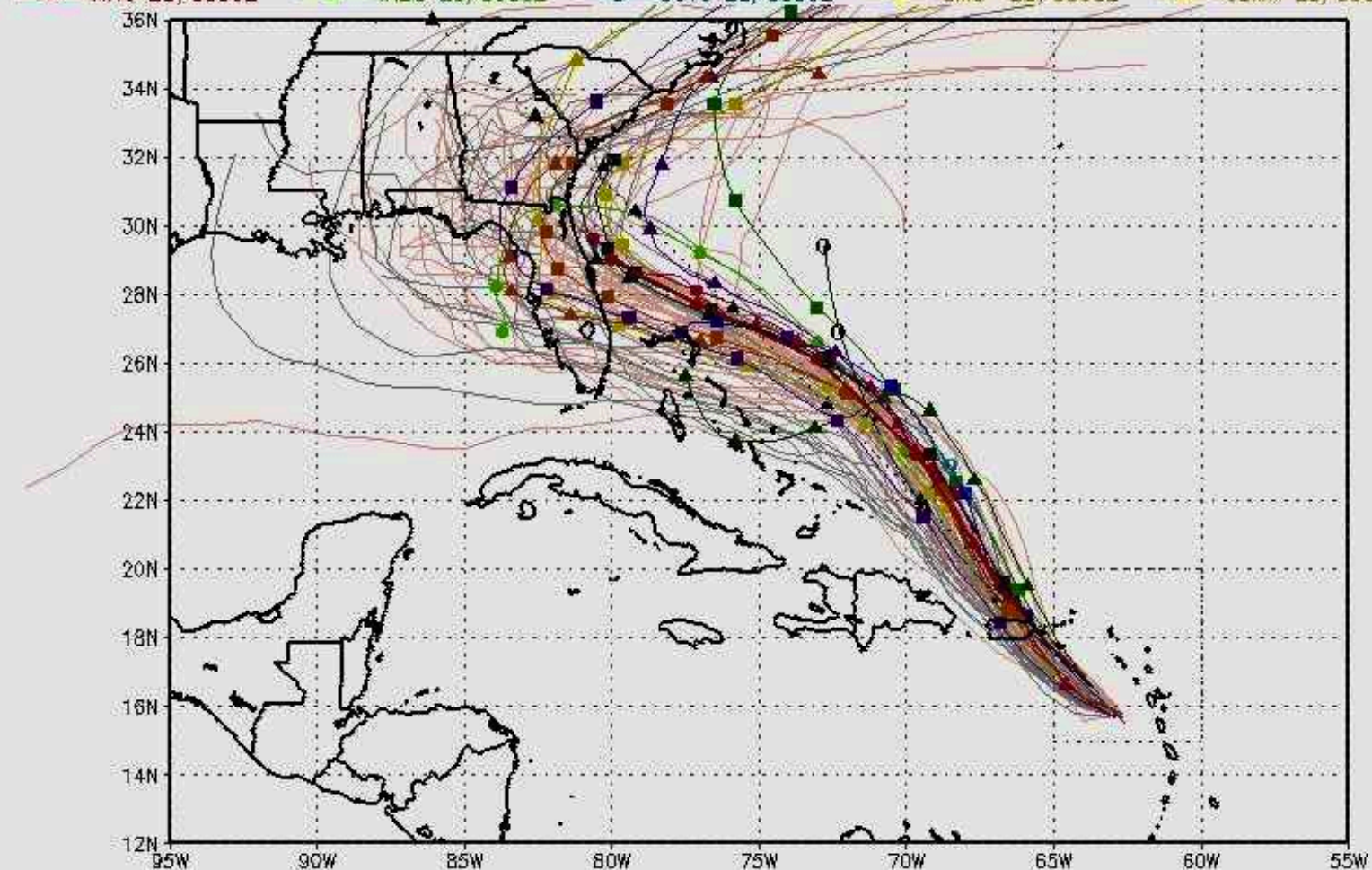
.... but we can simulate it!



Hurricane Dorian Forecast Track and Intensity



--▲-- XTRP 28/0800Z	—○— 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	—■— OEMN 28/0000Z



storm_05

sfwmd.gov

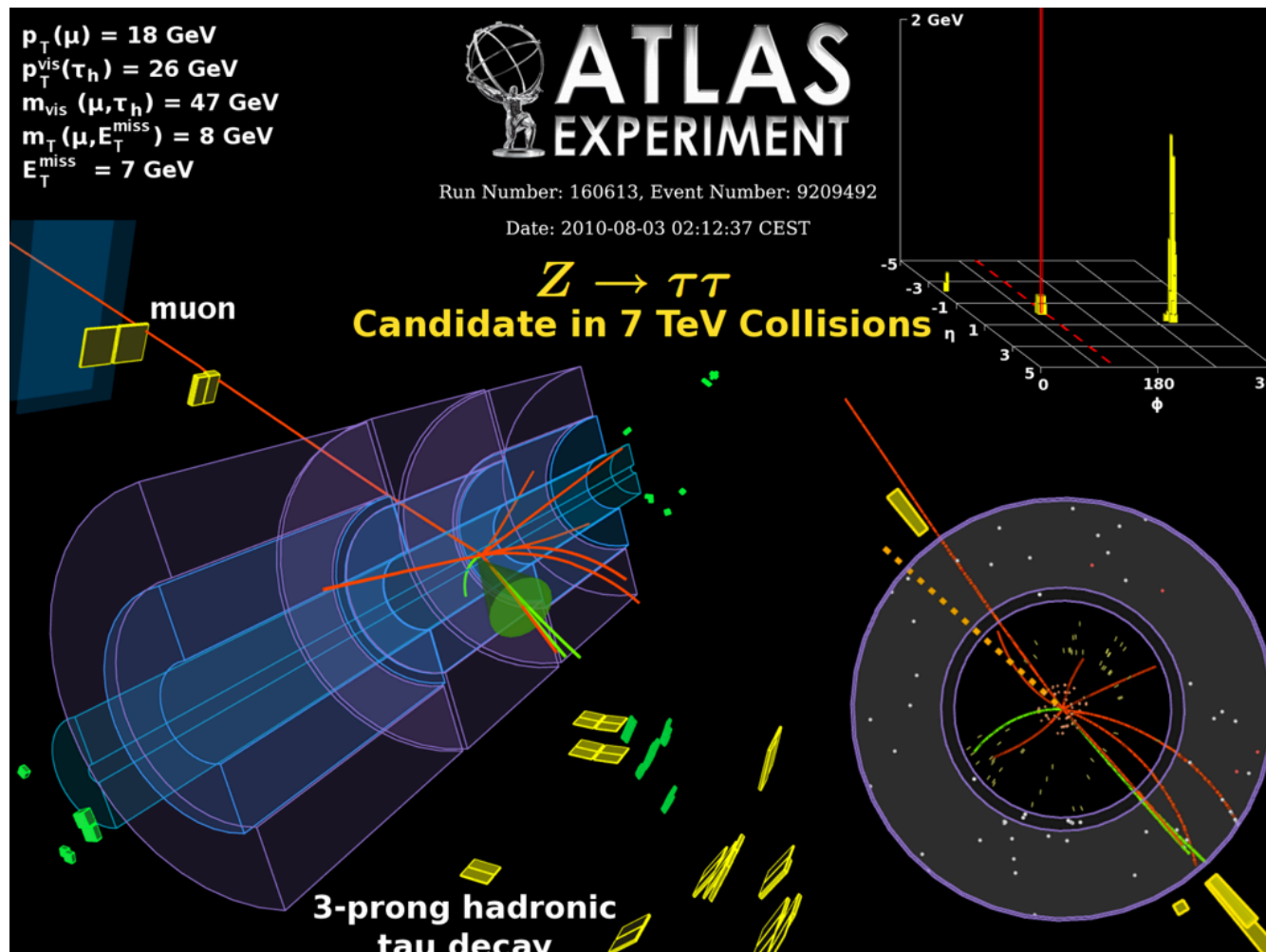
weather@sfwmd.gov

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>

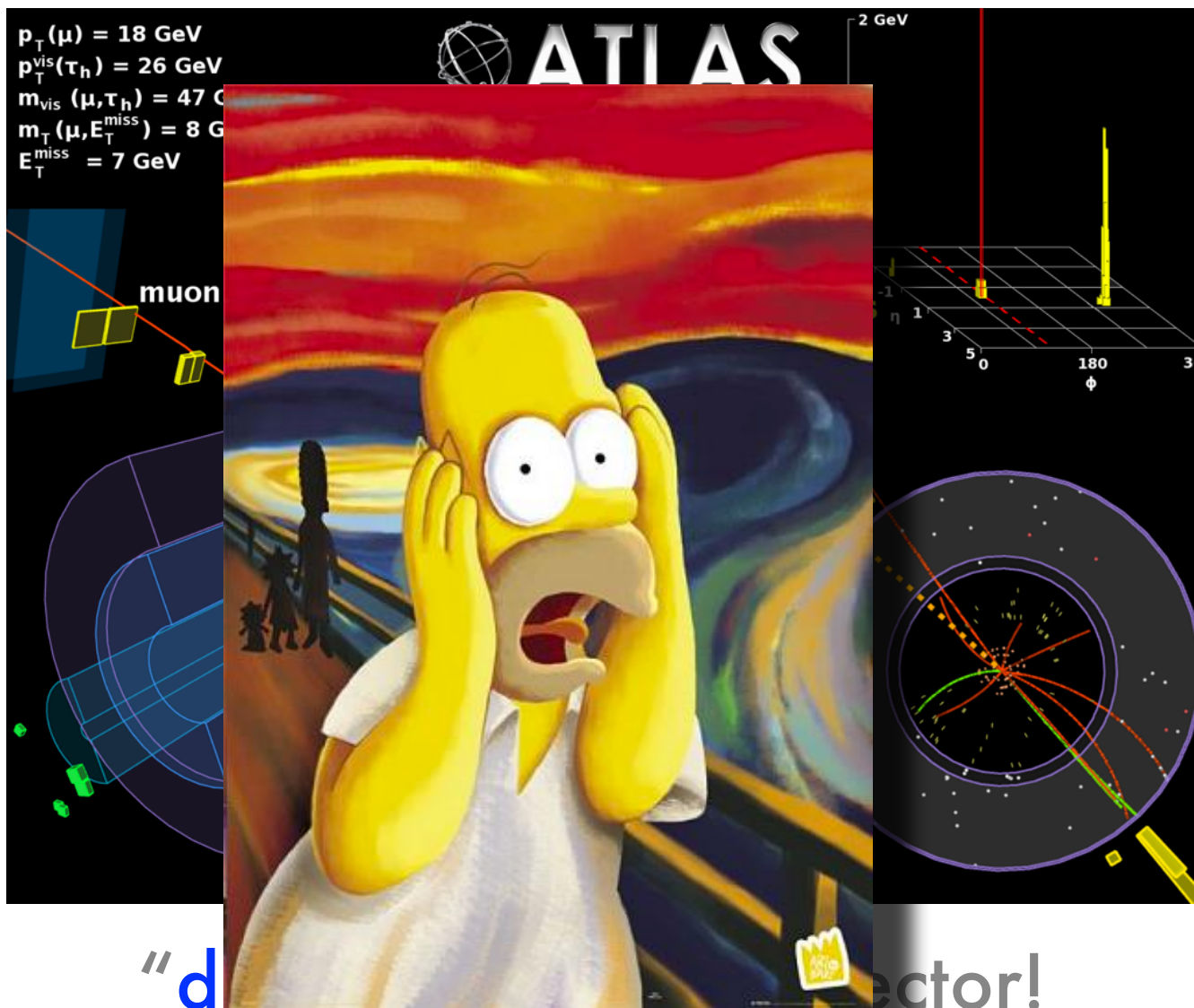


The nightmare



“data” is a 100M-d vector!

The nightmare



The nightmare



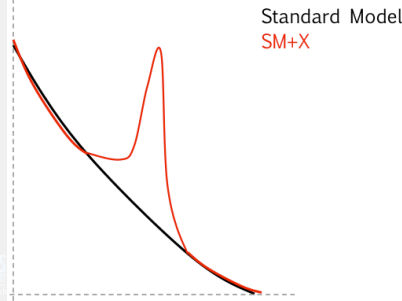
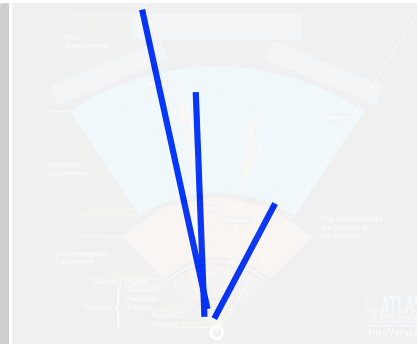
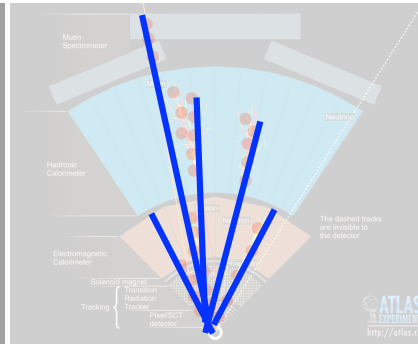
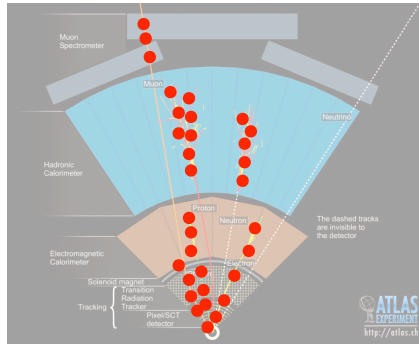
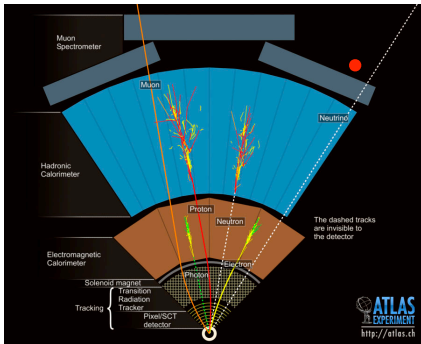
We wouldn't need ML if we could:

- ~~Express the likelihood of seeing our data~~
- ~~Access infinite computing resources~~
- ~~Develop infinitely-fast simulation~~



Summary statistics

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



We don't need to analyze the raw data

...If we could summarize it perfectly

Summary statistics

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

We wouldn't need ML if we could:

- ~~Express the likelihood of seeing our data~~
- ~~Access infinite computing resources~~
- ~~Develop infinitely-fast simulation~~
- ~~Derive perfect summary statistics~~

Standard Model
SM+X

...If we could summarize it perfectly



Summary statistics

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

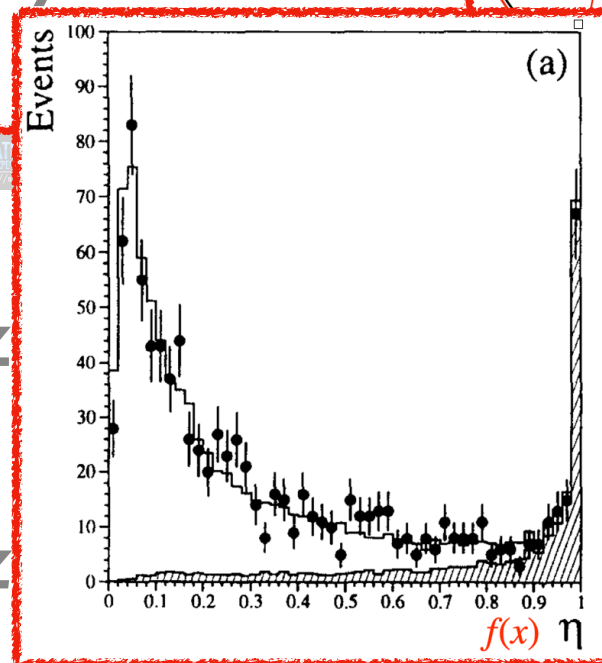
Traditional ML (< 2012)

Mostly allowed for >1 analysis feature
Combine a few features

Standard Model
SM+X

We don't need to analyze

...If we could summarize



Functional space

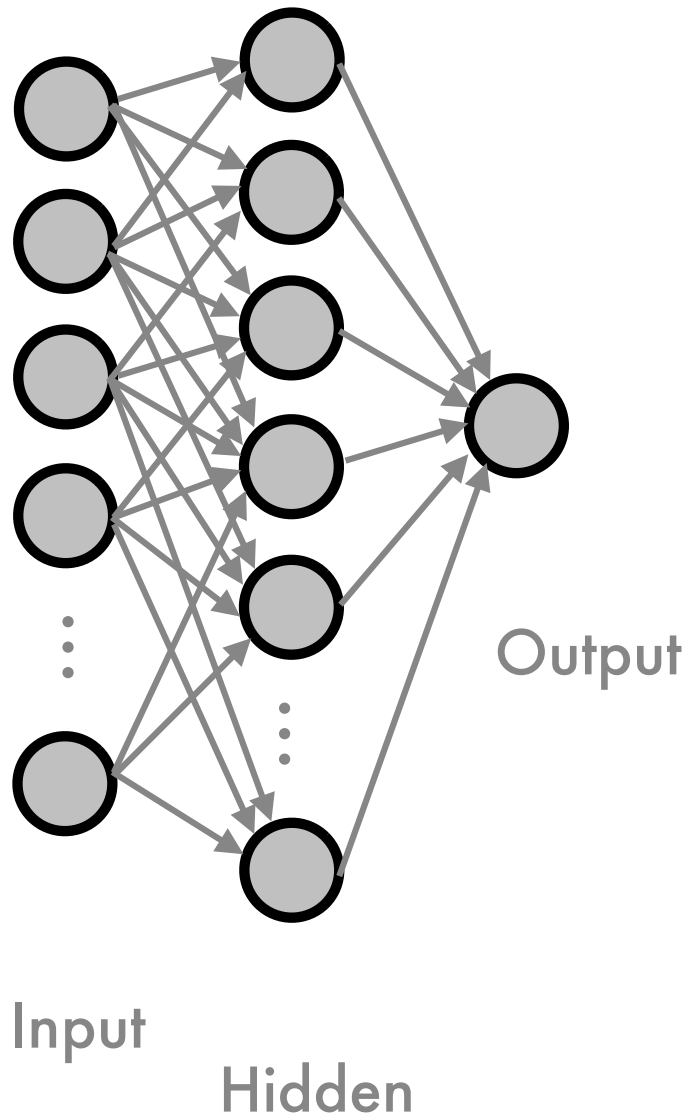
All functions

*Global
Optimum*



How complex?

Essentially a functional fit with many parameters

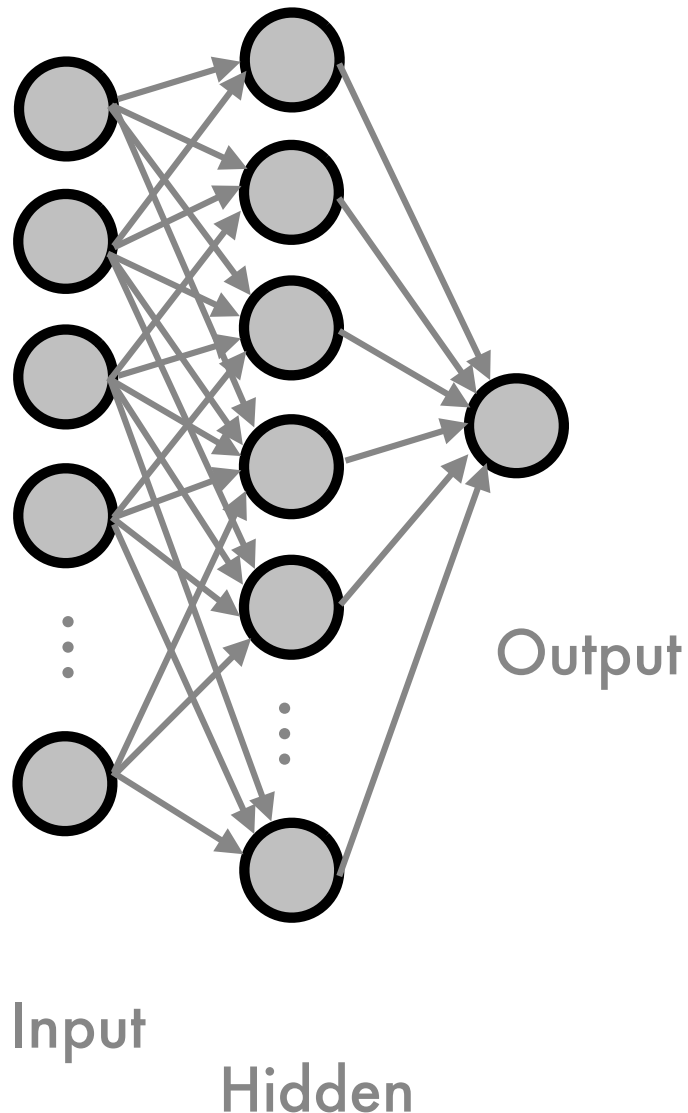


Single hidden layer

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

How complex?

Essentially a functional fit with many parameters

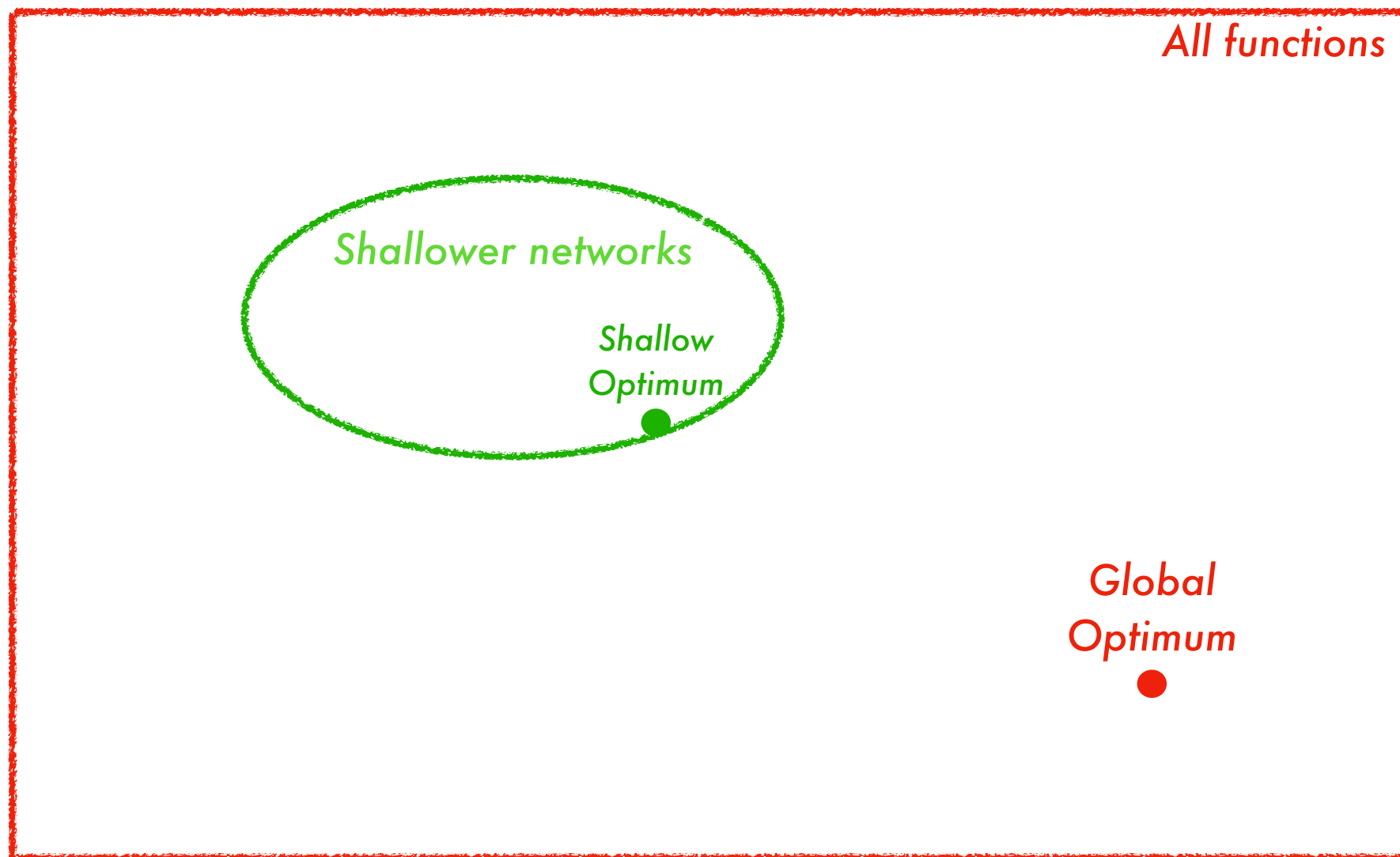


Single hidden layer

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

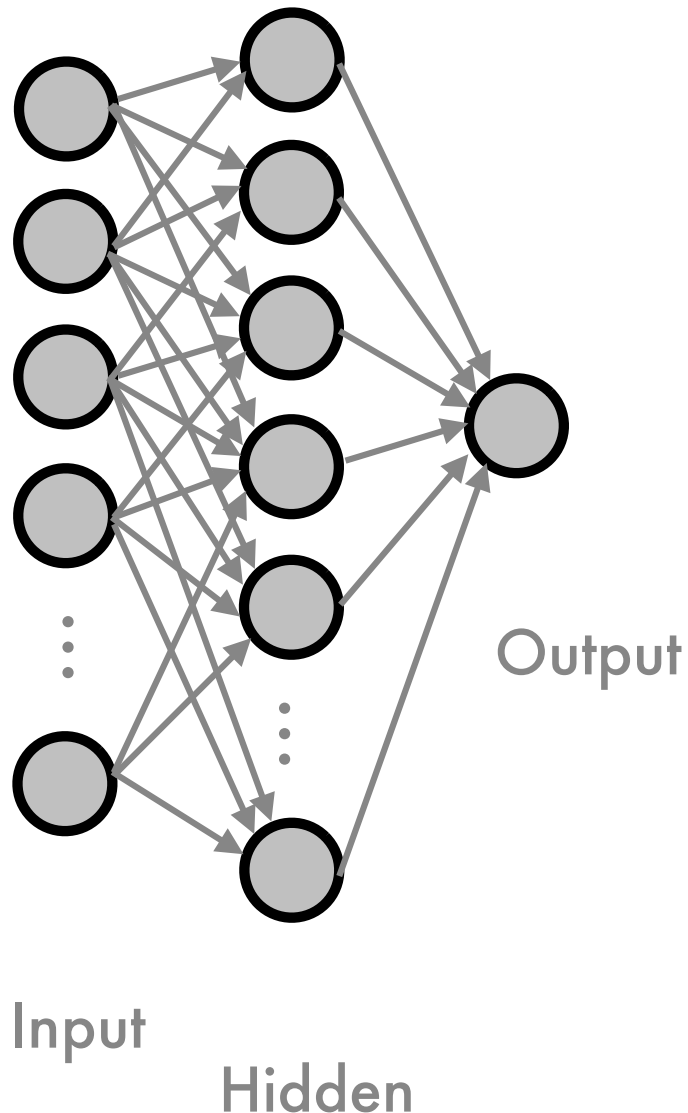
But might require very large hidden layer

Shallow space



Neural Networks

Essentially a functional fit with many parameters



Consequence:

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

In short:

Couldn't just throw data at NN.

Search for Input

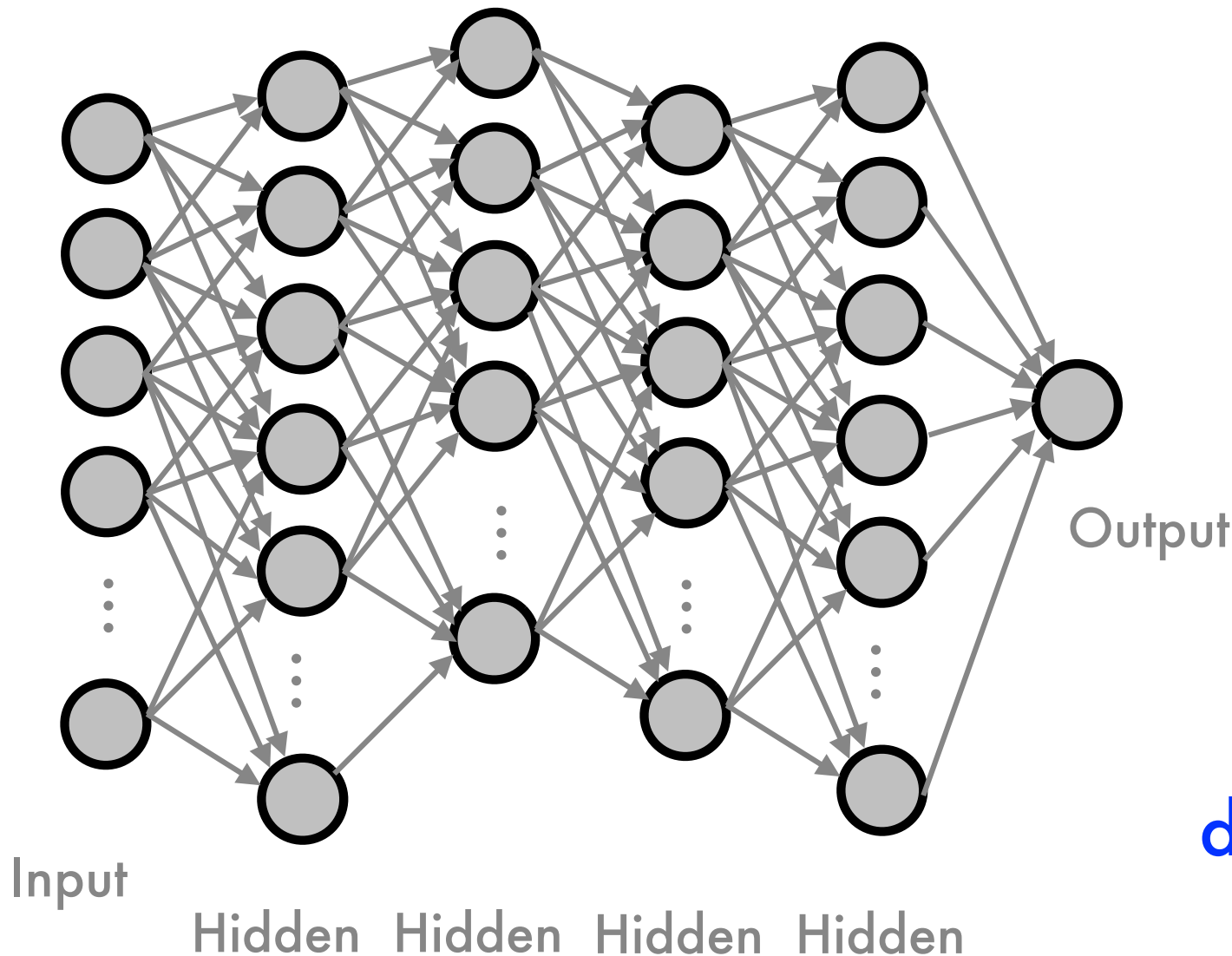
No low-level inputs

Limited input size

Painstaking search through input 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_T^{total}		•	•			
sum p_T					•	•
$p_T(\tau_1)/p_T(\tau_2)$					•	•
$E_T^{\text{miss}} \phi$ 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				•		
$p_T^{\ell_1}$				•		
$p_T^{j_1}$				•		
$E_T^{\text{miss}}/p_T^{\ell_2}$				•		
m_T		•			•	
$\min(\Delta\eta_{\ell_1, \ell_2, \text{jets}})$	•					
$j_3 \eta$ centrality	•					
$\ell_1 \times \ell_2 \eta$ centrality	•					
$\ell \eta$ centrality		•				
$\tau_{1,2} \eta$ centrality			•			

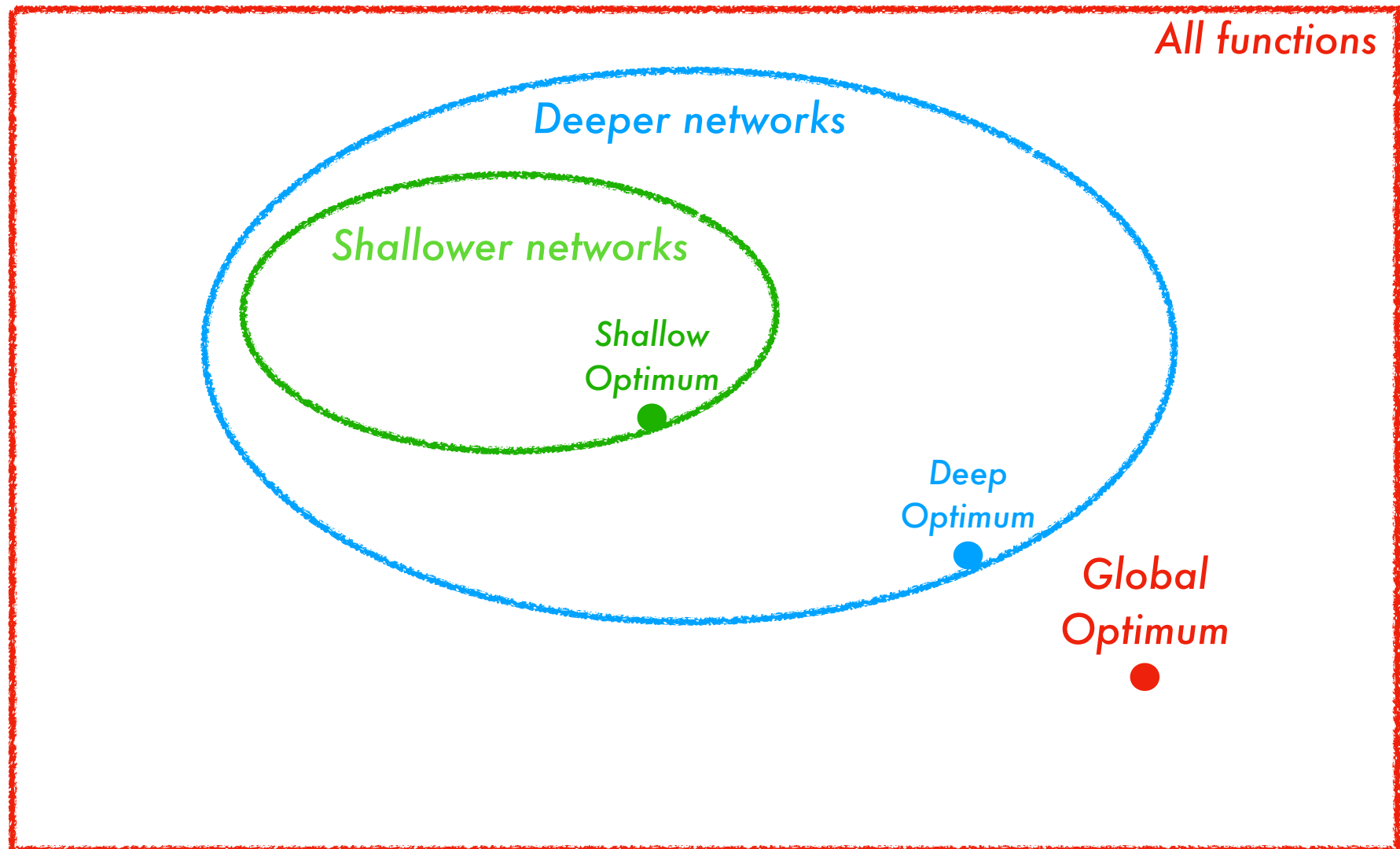
Deep networks



New tools
let us
train
deep
networks.

How well
do they work?

Expanding space



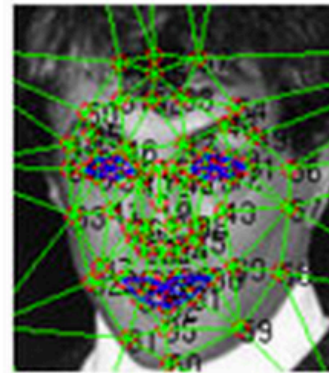
Real world applications



(a)



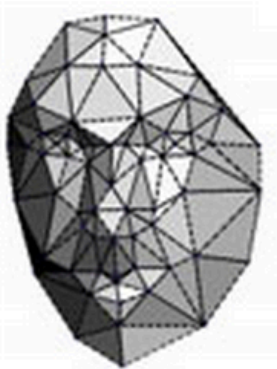
(b)



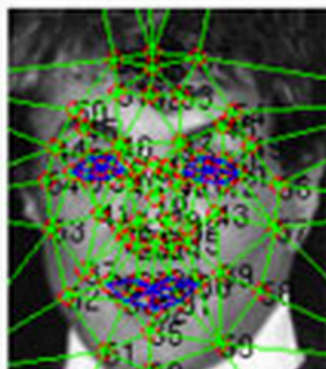
(c)



(d)



(e)



(f)



(g)

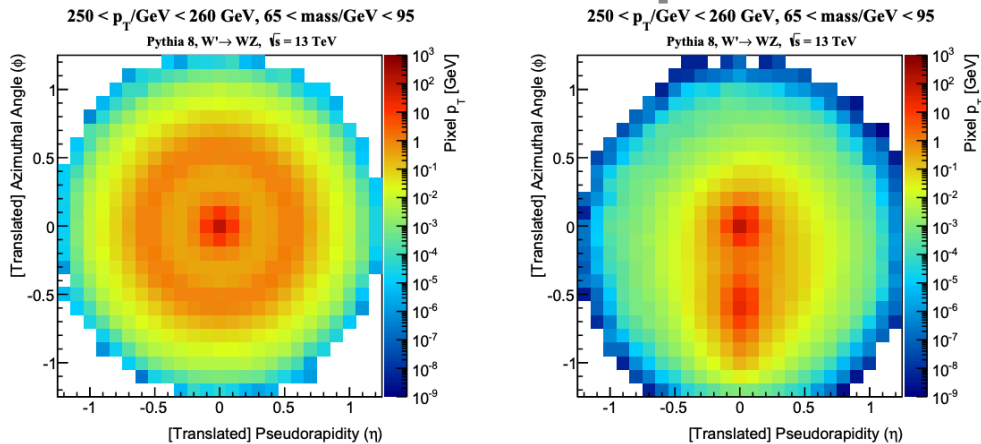


(h)

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.

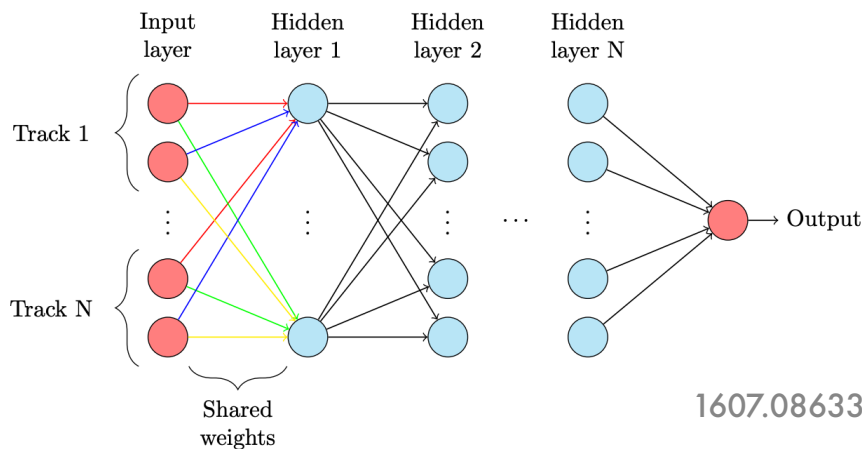
Low level data

Calorimeter pixels



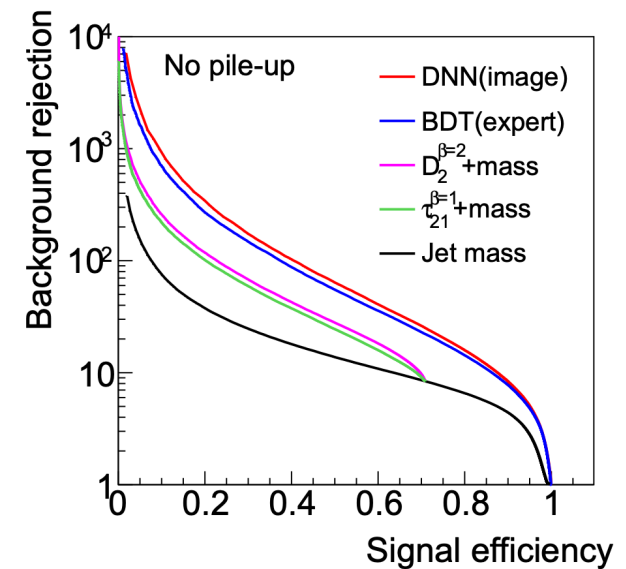
1511.05190

Lists of tracks



1607.08633 8

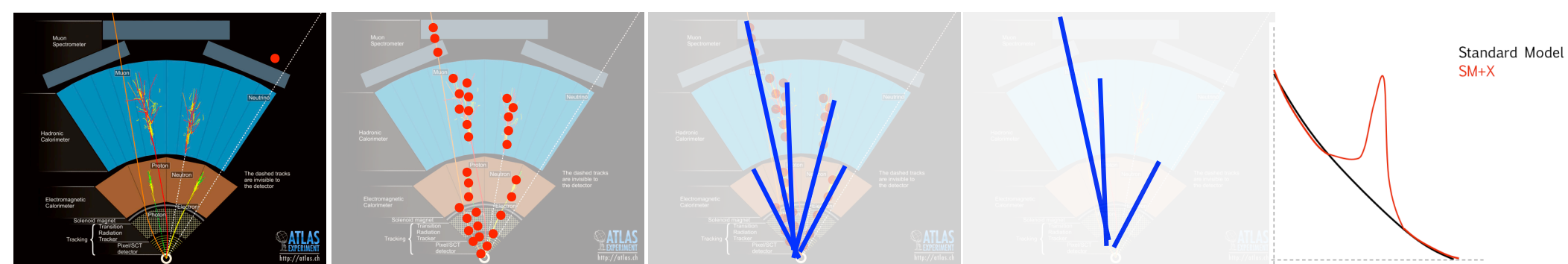
Networks beat experts



1603.09349

Summary statistics

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



Networks can handle higher dimensionality
And lower-level data

The new frontier

Expertise is not obsolete!

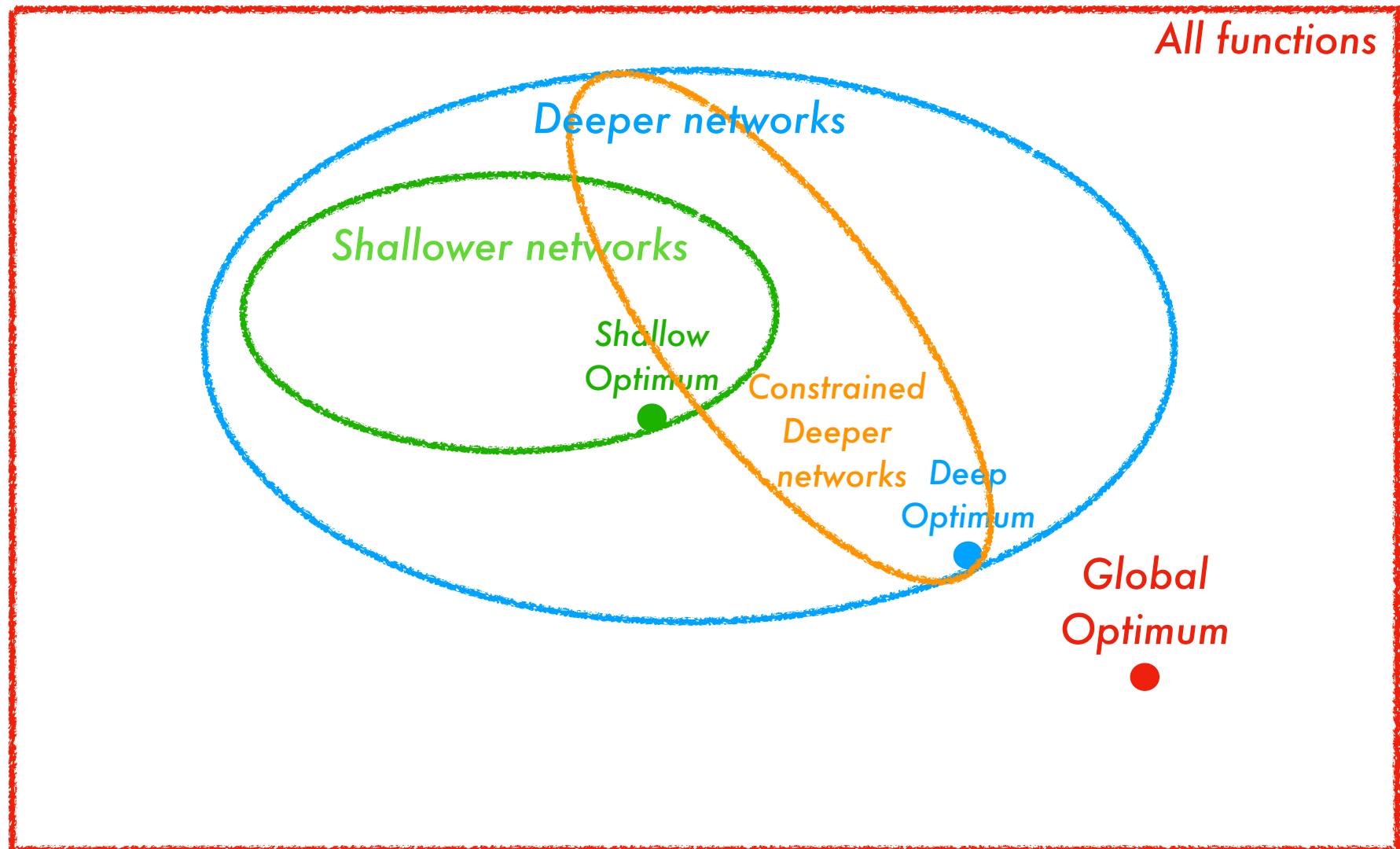
If you know something about the problem,
don't use a completely general solution.

Engineer your network structure!



e.g, network structures which respect **symmetries**

Constraining space

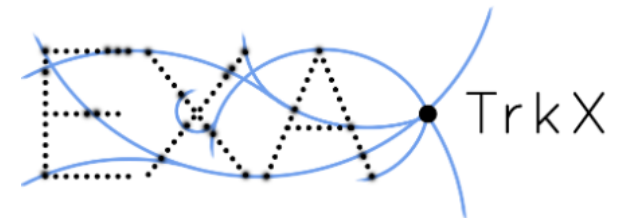
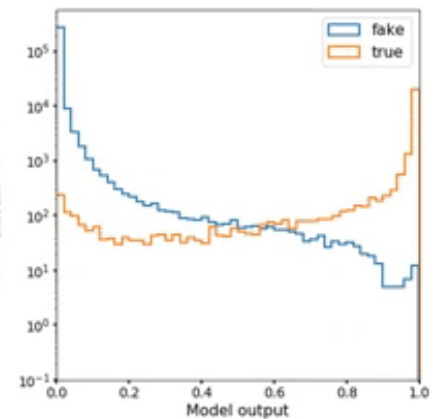
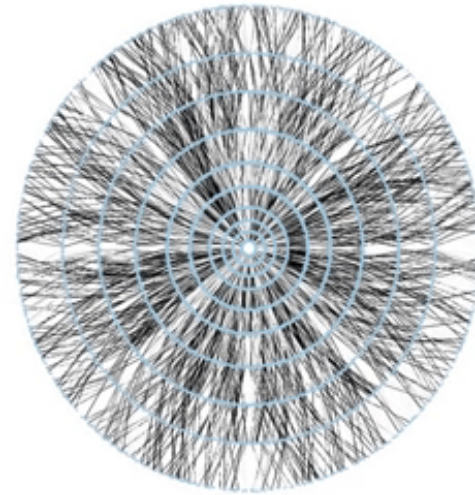
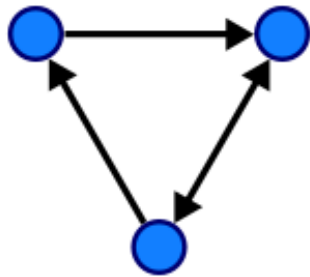


Outline

1. **Much much** more of the same
2. Something **qualitatively new**

Graph networks

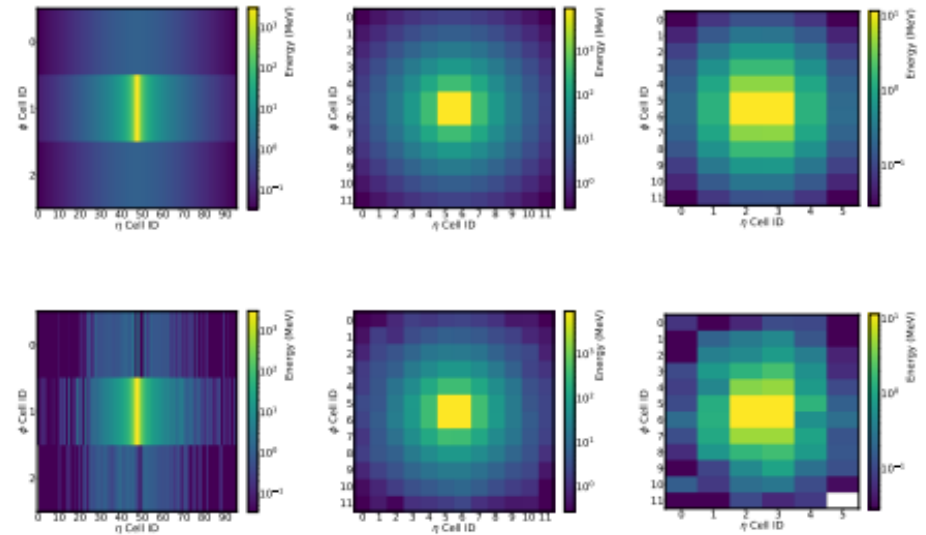
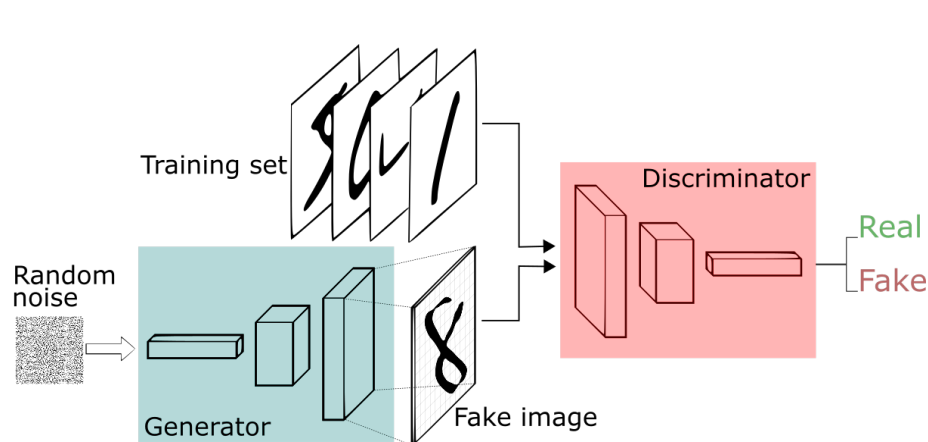
Represent structured data



Generative models

Do more than classify

Generate data from noise

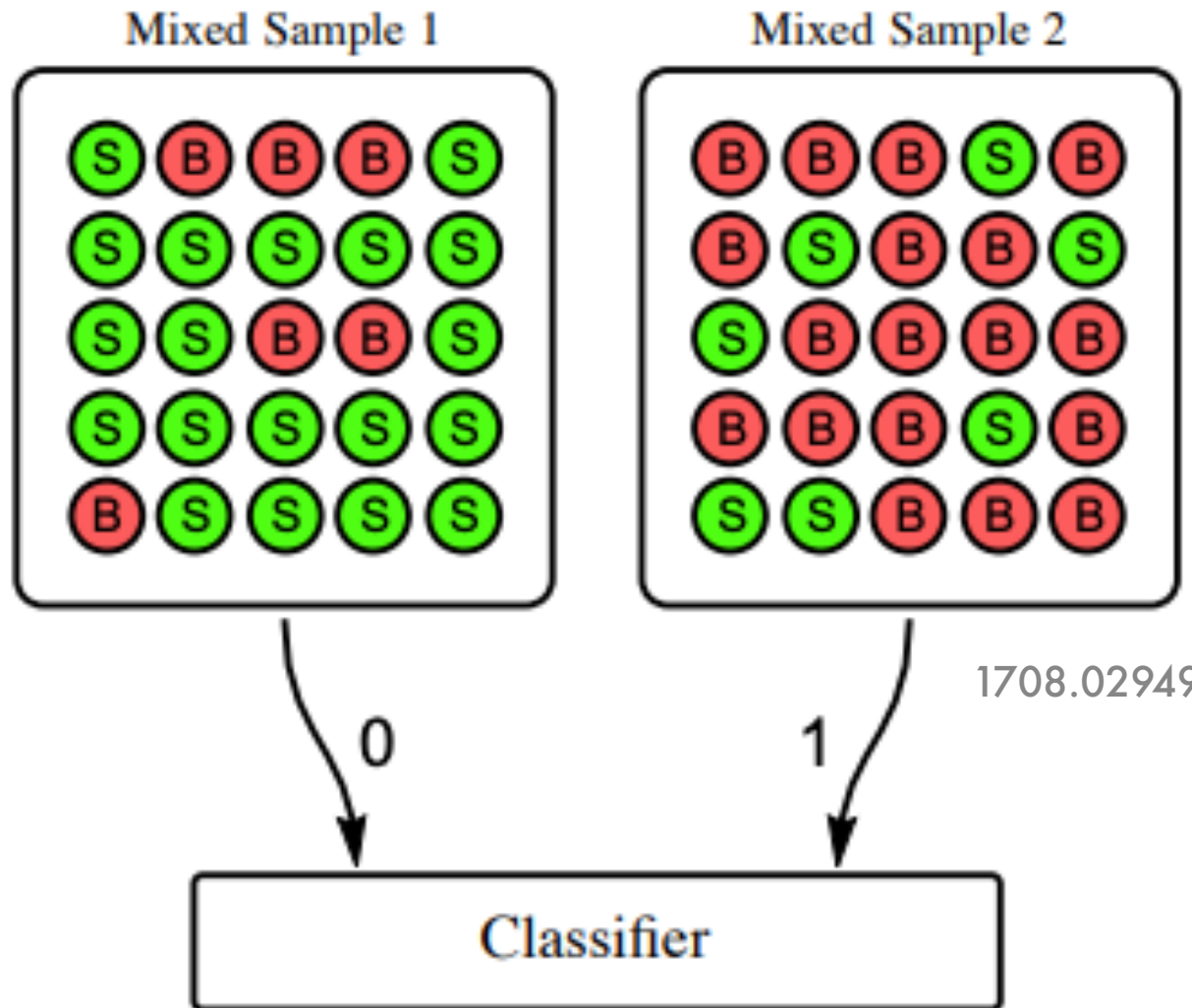


1712.10321

Optimal transport: new ways to compare distributions

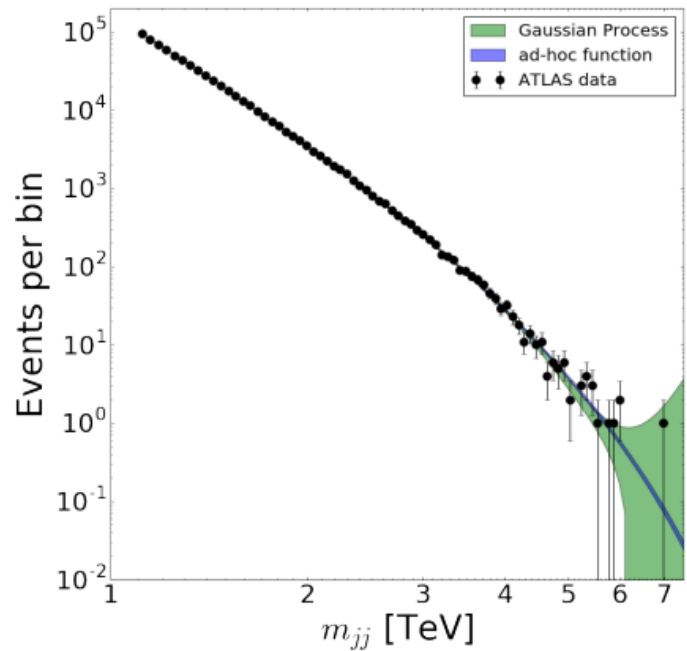
2101.08944: Learn the detector from data!

Away from supervision

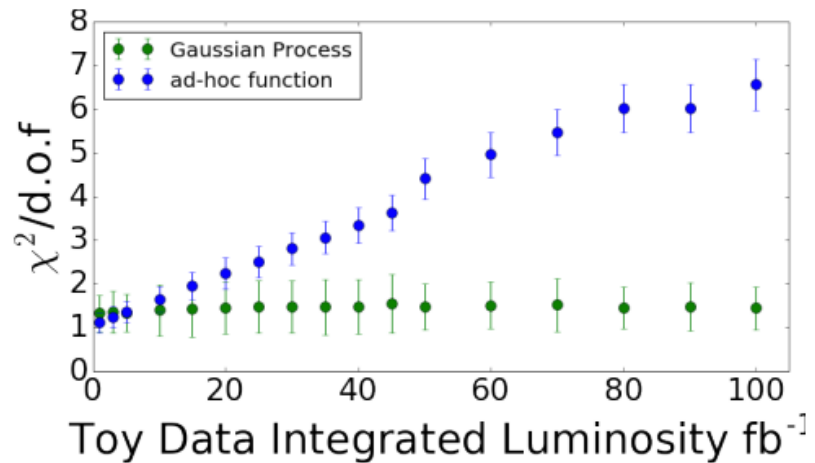


Background fitting

Away from ad-hoc background shapes:



1709.05681



ML for design

Optimize everything

Automatic Differentiation

Numerical gradients $\Delta L / \Delta \phi$ hopeless in trillion-D, need exact gradients $\partial L / \partial \phi$

Automatic Differentiation: careful application of *chain rule to computer programs*

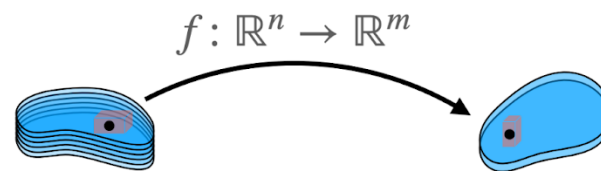
```
import jax
import jax.numpy as jnp

def func(x):
    y = x
    for i in range(4):
        y += x[0]**2 + jnp.sin(x[1]) + jnp.exp(-x[2])
    y = y.sum()
    return y
```

exact gradients!

```
gfunc = jax.value_and_grad(func)
gfunc(jnp.array([2., 3., -2]))

(DeviceArray(141.36212, dtype=float32),
 DeviceArray([ 49.          , -10.8799095, -87.66867 ], dtype=float32))
```



$$y = f(x) \quad dy = J_f dx$$

$$J_f = \frac{\partial(y_1, \dots, y_m)}{\partial(x_1, \dots, x_n)}$$

 TensorFlow



PYTORCH

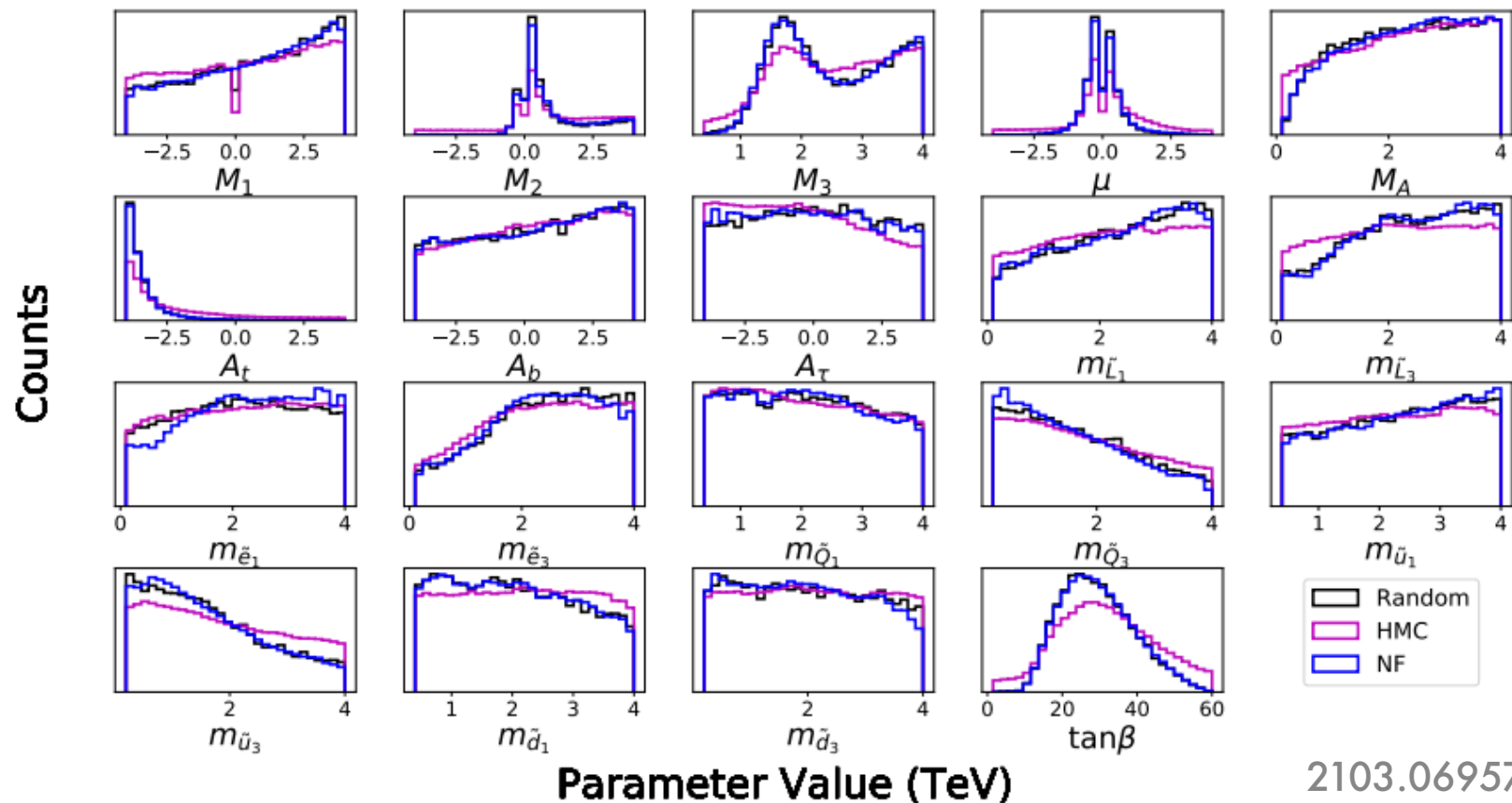
... but also C++, Fortran, ...

L. Heinrich

See also: 1806.04743

ML for Theory!

How do we search large spaces?



2103.06957

String theory applications: 1707.00655, 1903.11616

Summary

Modern ML

Much more flexible and capable
Tackling previously intractable problems

Many creative new ideas

Widening in scope
Attacking new problems