

# Introduction to Deep Learning

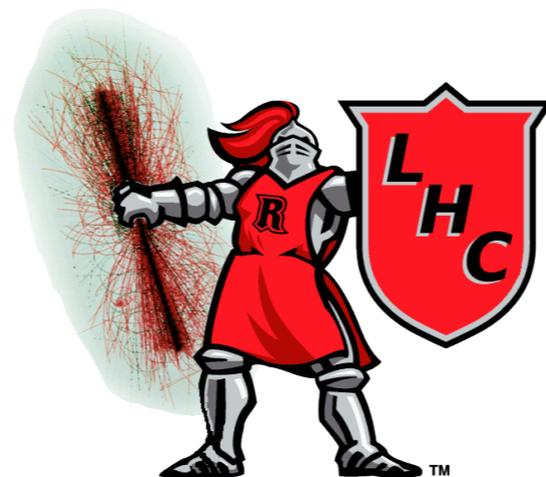
and its applications to the LHC

David Shih

“Machine Learning at LHC”

KMI Nagoya

February 2020



# The AI Revolution is Here



The past decade has seen remarkable advances in machine learning and artificial intelligence.

These technological breakthroughs are reshaping the world around us.

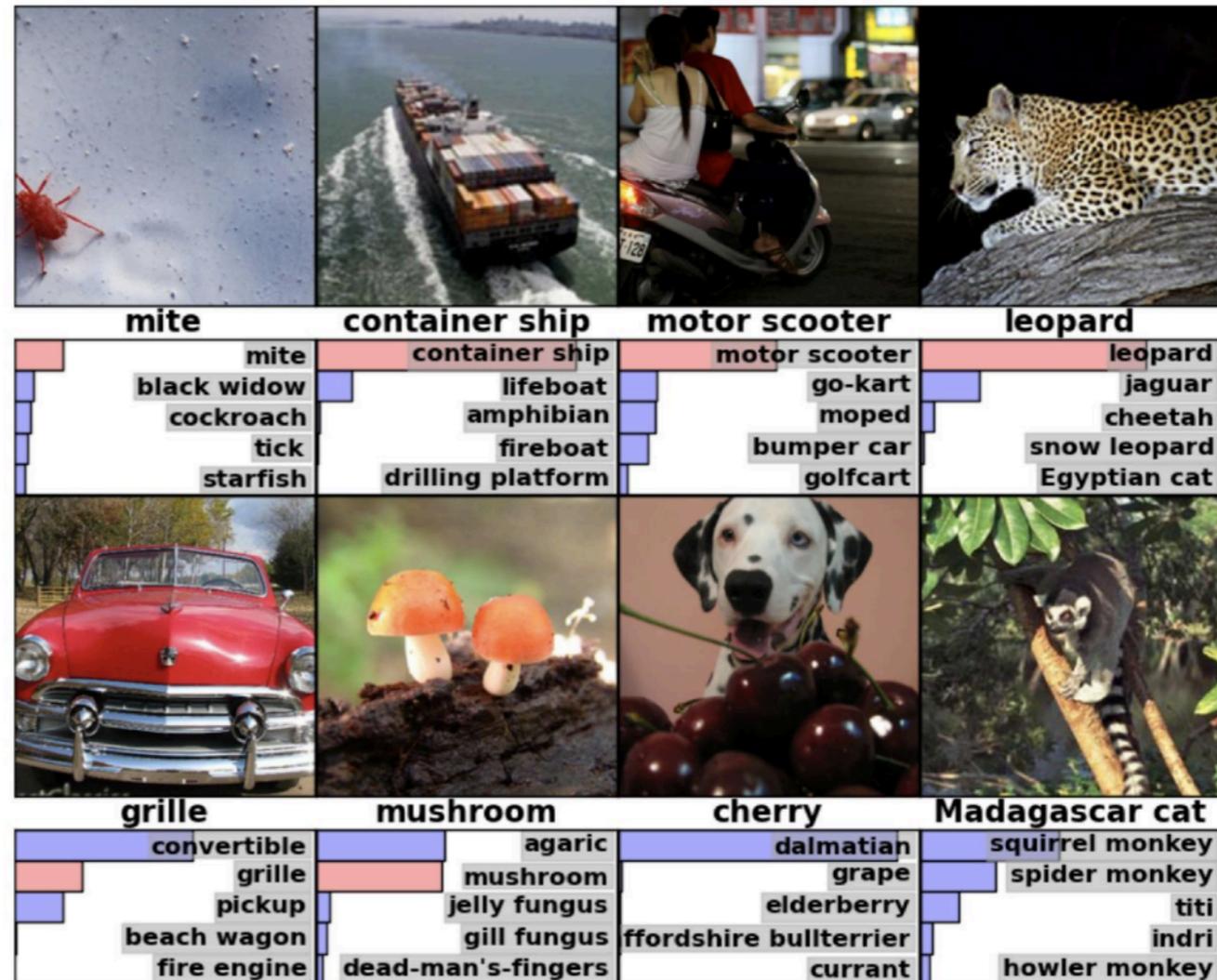
Deep learning also has the potential to revolutionize physics at the LHC.

# Deep Learning Breakthrough

## ImageNet Challenge

IMAGENET

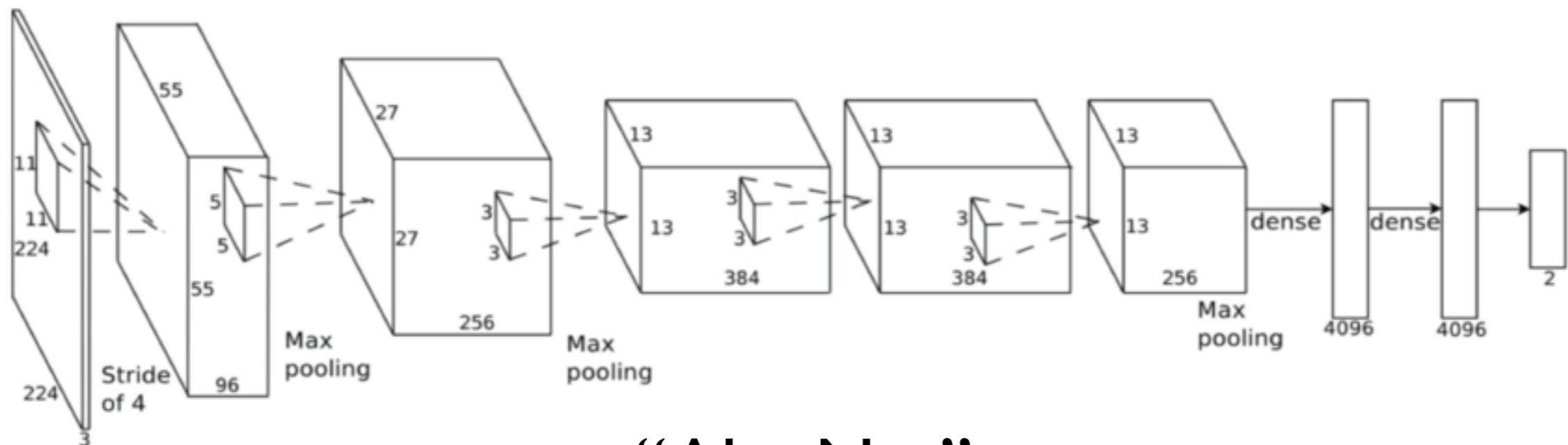
- 1,000 object classes (categories).
- Images:
  - 1.2 M train
  - 100k test.



# Deep Learning Breakthrough

In 2012, a deep convolutional neural network won the “ImageNet” image classification competition by a huge margin (Krizhevsky, Sutskever, Hinton)

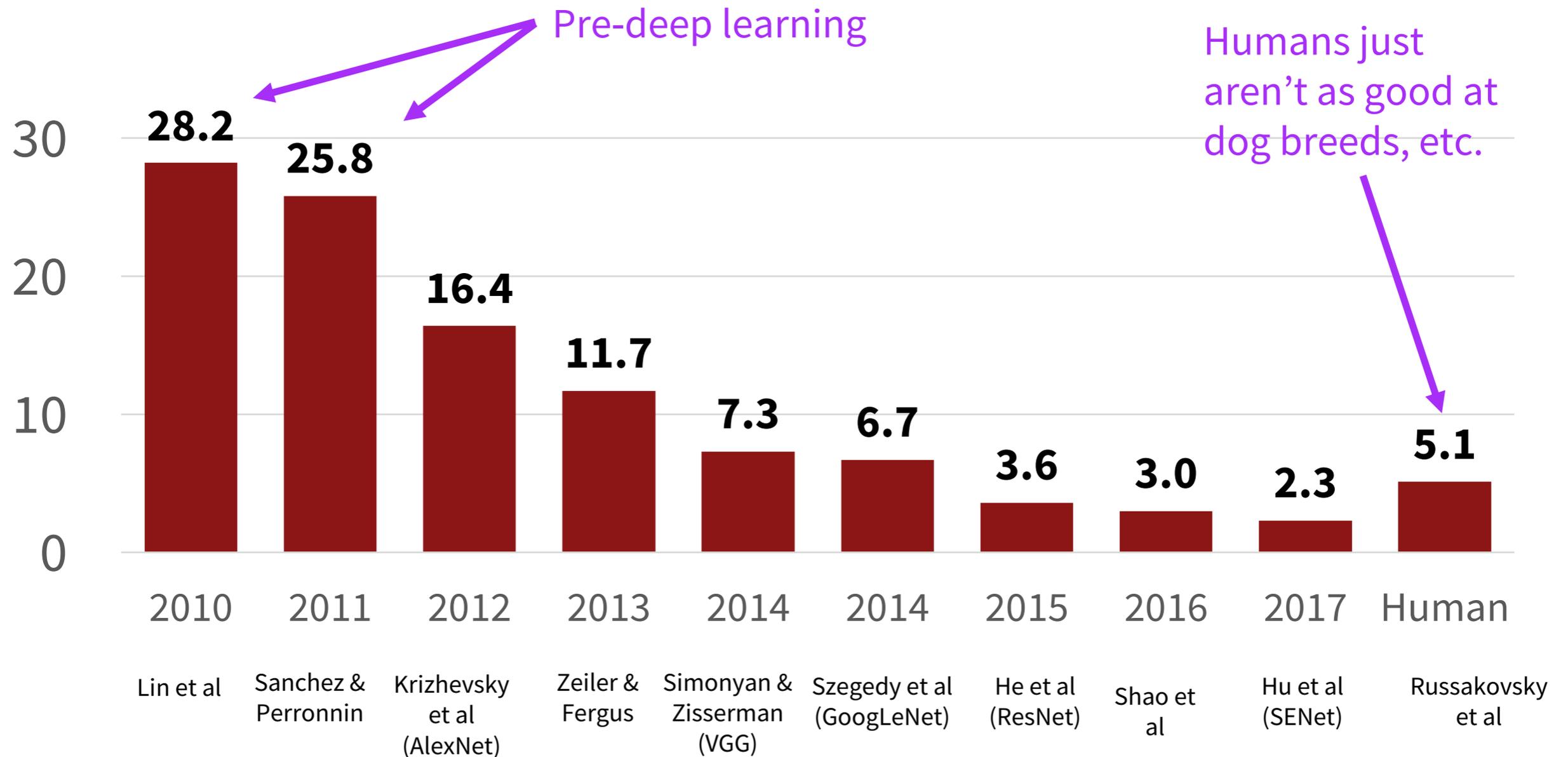
This dramatic breakthrough inaugurated the modern revolution in deep learning.



“AlexNet”

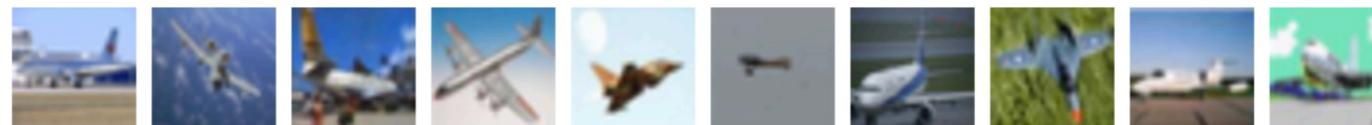
# of parameters 1000 x LeNet (60M). Required training on a GPU.

# Deep Learning Breakthrough

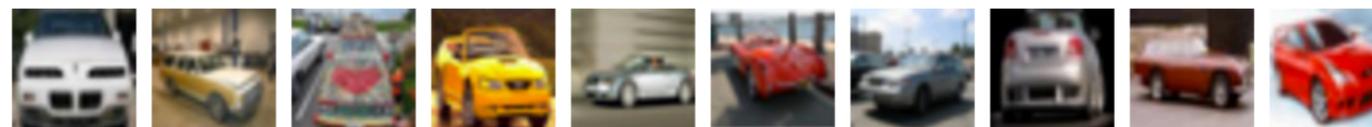


# Computer vision

**airplane**



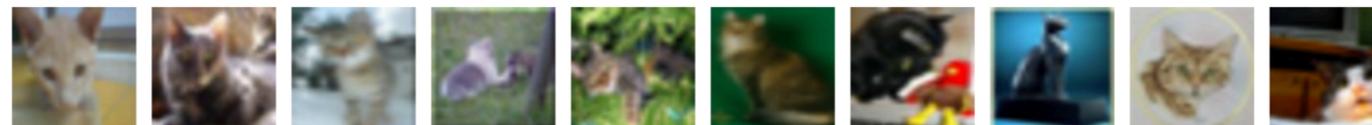
**automobile**



**bird**



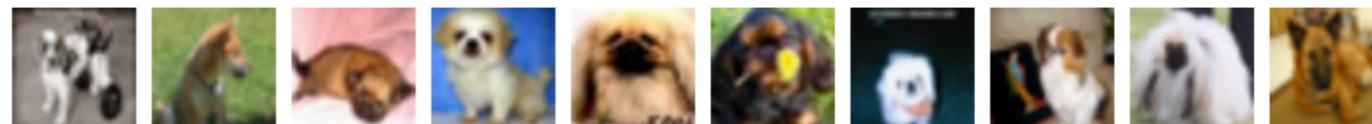
**cat**



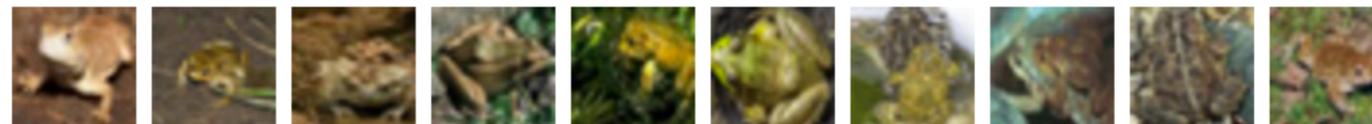
**deer**



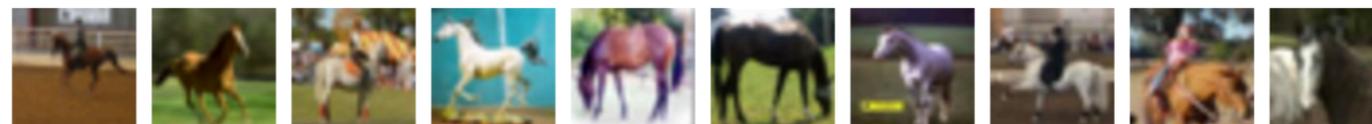
**dog**



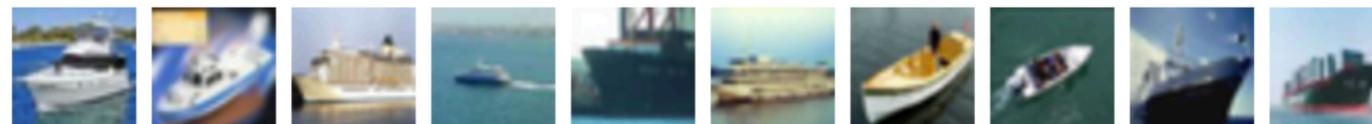
**frog**



**horse**



**ship**



**truck**



# Computer vision

**Classification**



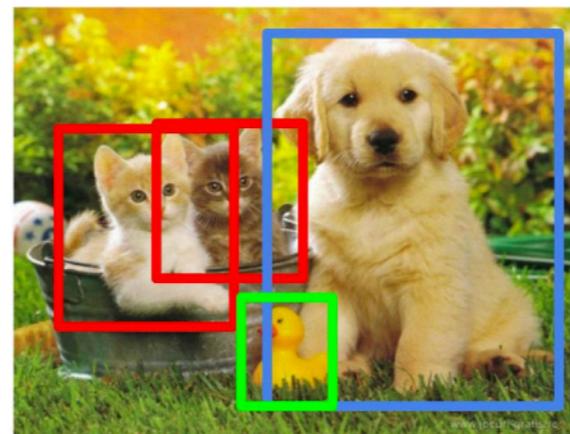
CAT

**Classification  
+ Localization**



CAT

**Object Detection**



CAT, DOG, DUCK

**Instance  
Segmentation**

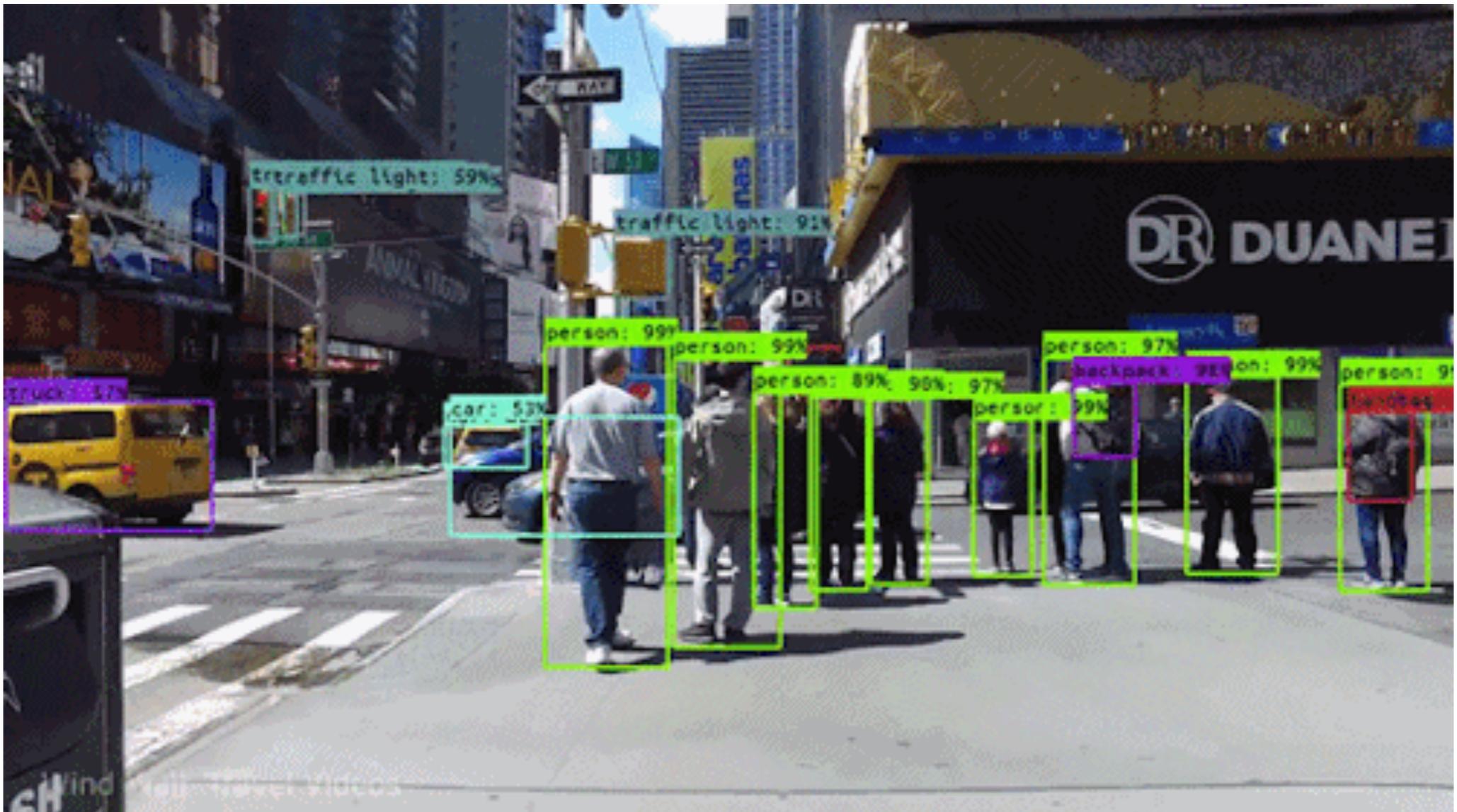


CAT, DOG, DUCK

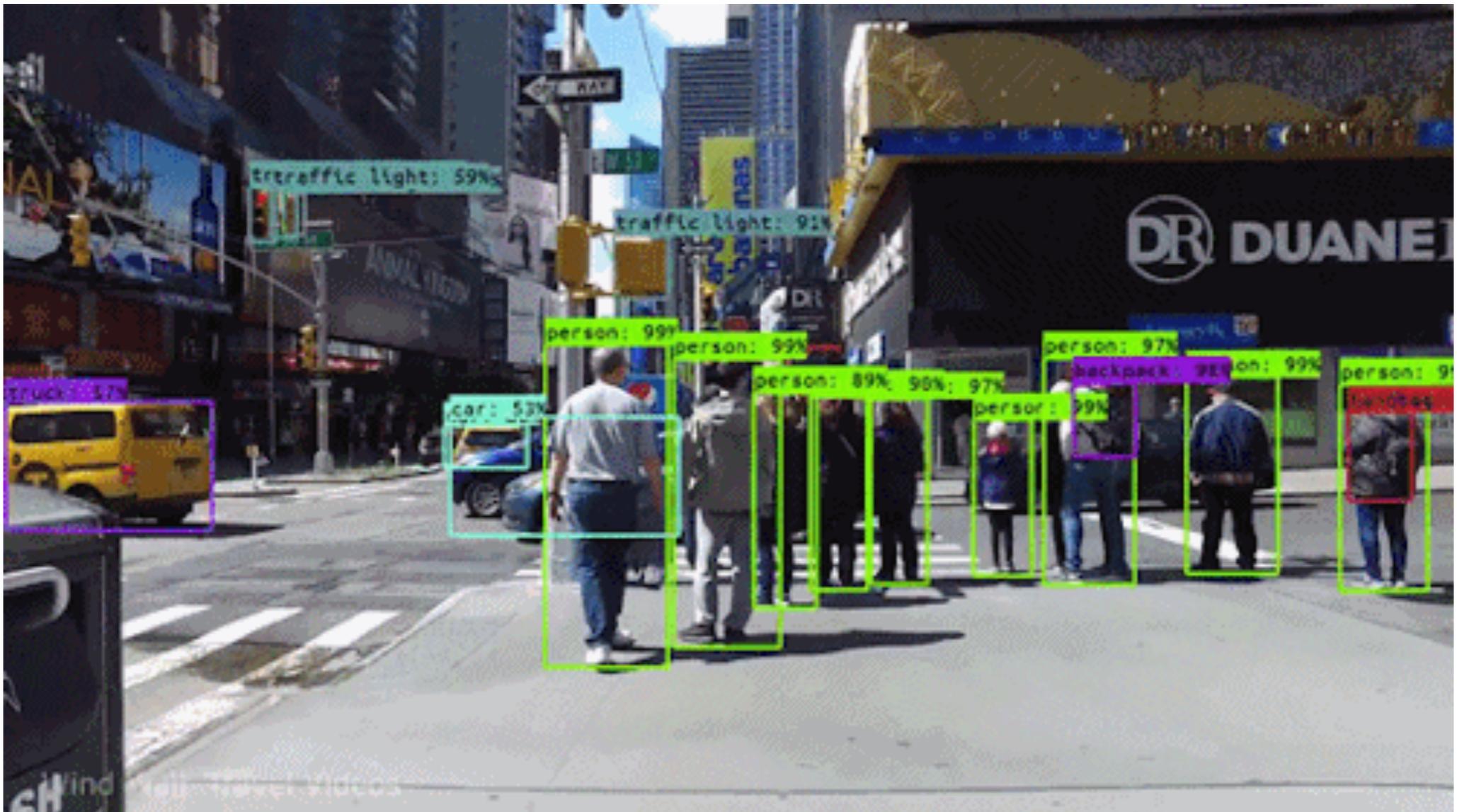
Single object

Multiple objects

# Computer vision



# Computer vision



# Natural language



Yoruba ▾ [Translate from English](#)

dog  
dog dog dog dog dog dog dog dog dog dog

English ▾

Doomsday Clock is three minutes at twelve We are experiencing characters and a dramatic developments in the world, which indicate that we are increasingly approaching the end times and Jesus' return

[Open in Google Translate](#)

[Feedback](#)

i pooped

- i pooped my pants
- i pooped my pants at school
- i pooped on my boyfriend
- i pooped hard it overflow
- i pooped my pants at work
- i pooped a cornish game hen
- i pooped a hammer
- i pooped the bed
- i pooped my pants and i liked it**
- i pooped today

[Google Search](#) [I'm Feeling Lucky](#)

[Advanced Search](#)  
[Language Tools](#)



**Jacob Kluge**

@KlugeJake

[Follow](#)

When auto correct hates you and your relationship

Can we go to the gym tomorrow

Sure Abby

Baby\*

Oh boy here we go

Read 8:59 PM

iMessage



Amazon Echo †



Google Home \*



Apple HomePod



Microsoft Invoke

# Generative modeling

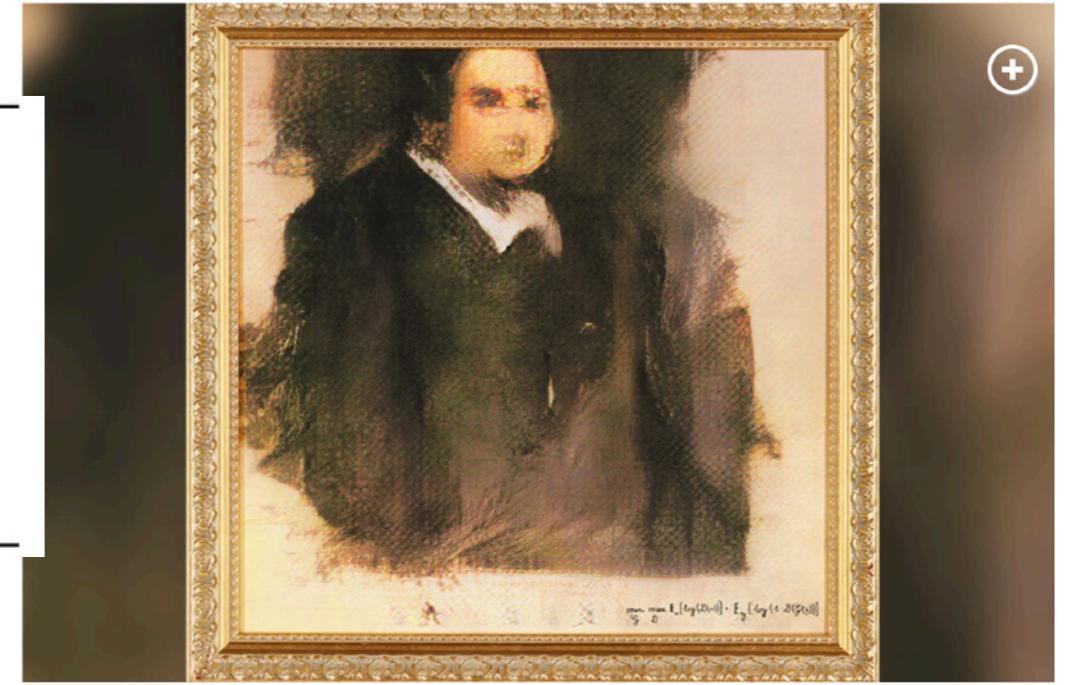
**there is no one else in the world .**  
*there is no one else in sight .*  
*they were the only ones who mattered .*  
*they were the only ones left .*  
*he had to be with me .*  
*she had to be with him .*  
*i had to do this .*  
*i wanted to kill him .*  
*i started to cry .*  
**i turned to him .**

**this was the only way .**  
*it was the only way .*  
*it was her turn to blink .*  
*it was hard to tell .*  
*it was time to move on .*  
*he had to do it again .*  
*they all looked at each other .*  
*they all turned to look back .*  
*they both turned to face him .*  
**they both turned and walked away .**

## Painting created by AI sells for stunning \$432K

By Tamar Lapin

October 26, 2018 | 2:38am | Updated



"Portrait of Edmond de Belamy"

AP

Automatically captioned



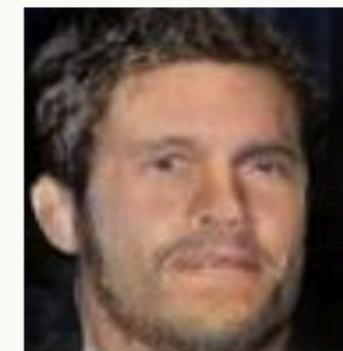
A dog is sitting on the beach next to a dog.



2014



2015

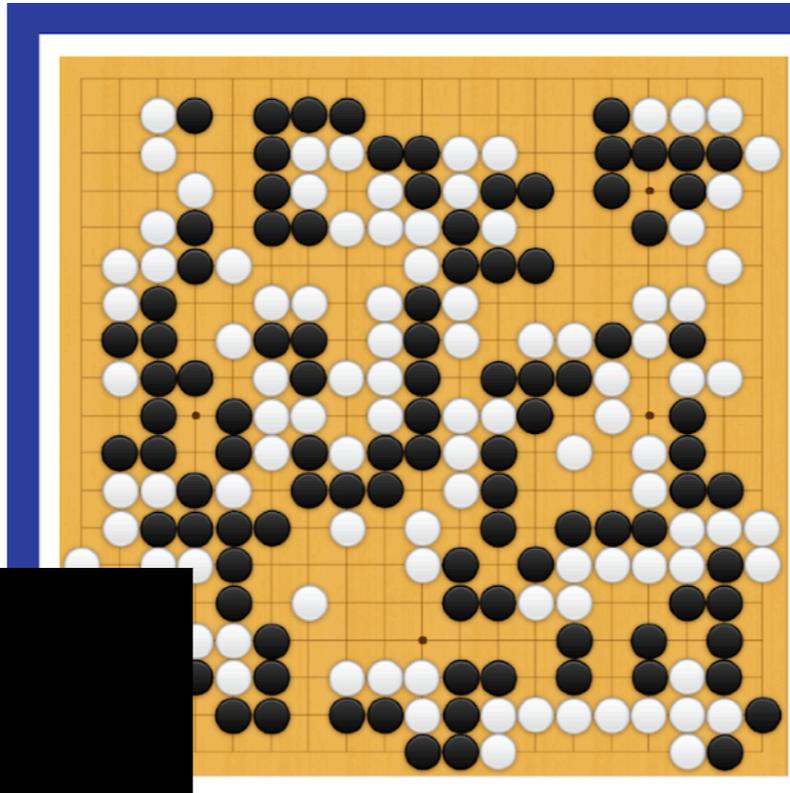


2016



2017

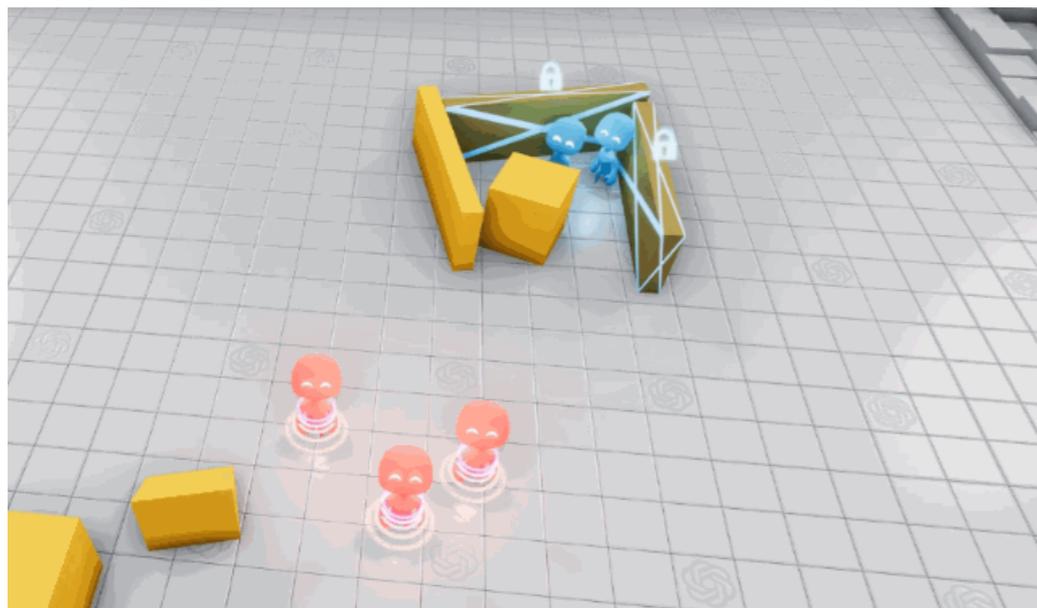
# Game playing



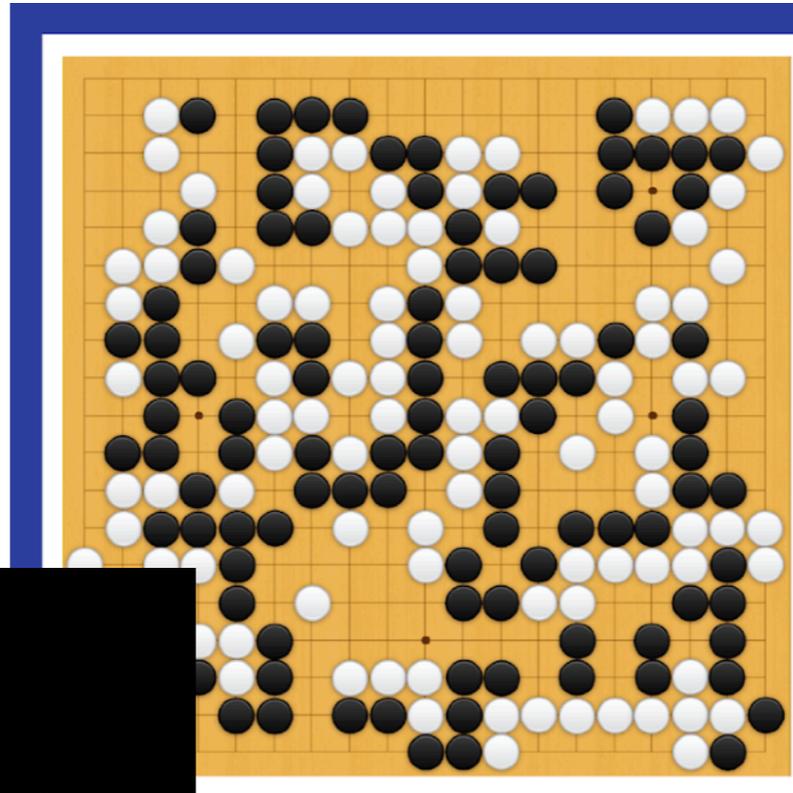
**THE ULTIMATE GO CHALLENGE**  
GAME 3 OF 3  
27 MAY 2017



**RESULT B + Res**



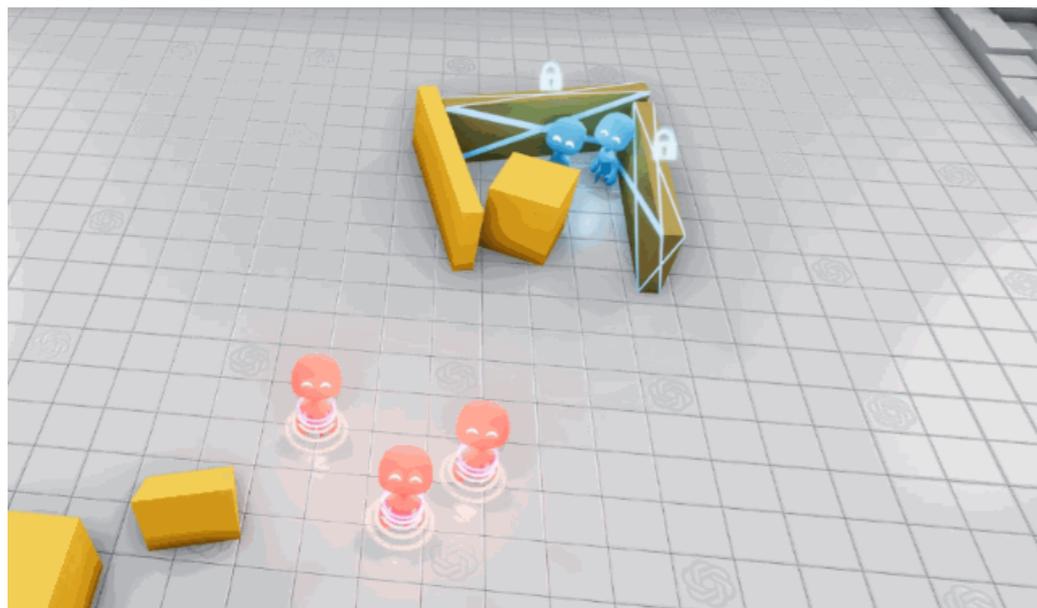
# Game playing



**THE ULTIMATE GO CHALLENGE**  
GAME 3 OF 3  
27 MAY 2017



**RESULT B + Res**



# Plan of the lecture

1. Machine Learning Basics
2. Intro to Deep Learning
3. Deep Learning at the LHC

I will assume most of you know some collider physics.

I will not assume any familiarity with machine learning or neural networks.

# I. Machine learning basics

# Brief (re)fresher on machine learning

*“ML is glorified function fitting”*

Want to **fit** a **function**  $f(x; \theta)$  with some parameters  $\theta$  (“weights”) to a collection of examples  $\{x_i\}$  in order to achieve some **objective**.

# Brief (re)fresher on machine learning

*“ML is glorified function fitting”*

Want to **fit** a **function**  $f(\mathbf{x}; \theta)$  with some parameters  $\theta$  (“weights”) to a collection of examples  $\{\mathbf{x}_i\}$  in order to achieve some **objective**.

$\mathbf{x} \in \mathbb{R}^d$  is called the **feature space**

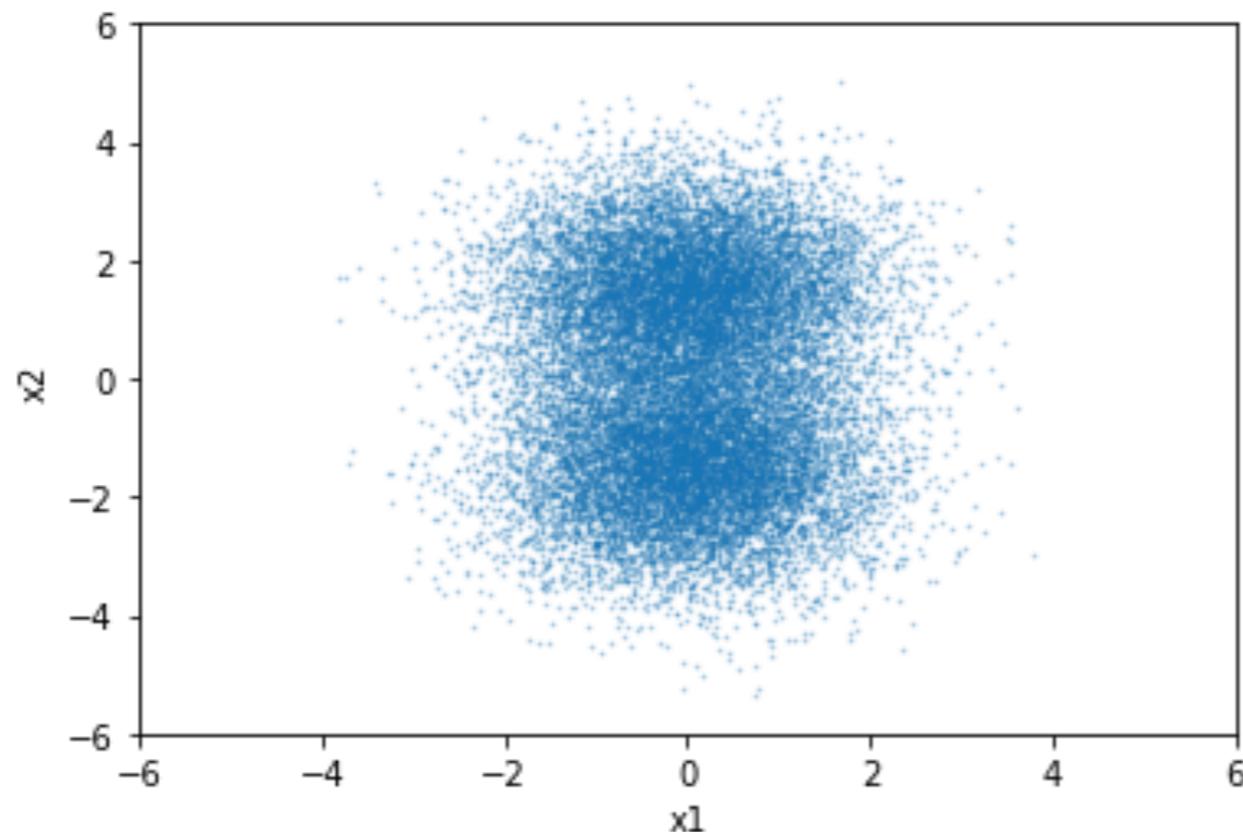
$i = 1, \dots, N$  indexes the training dataset.

# Brief (re)fresher on machine learning

*“ML is glorified function fitting”*

Want to **fit** a **function**  $f(x; \theta)$  with some parameters  $\theta$  (“weights”) to a collection of examples  $\{x_i\}$  in order to achieve some **objective**.

Some typical **objectives**:

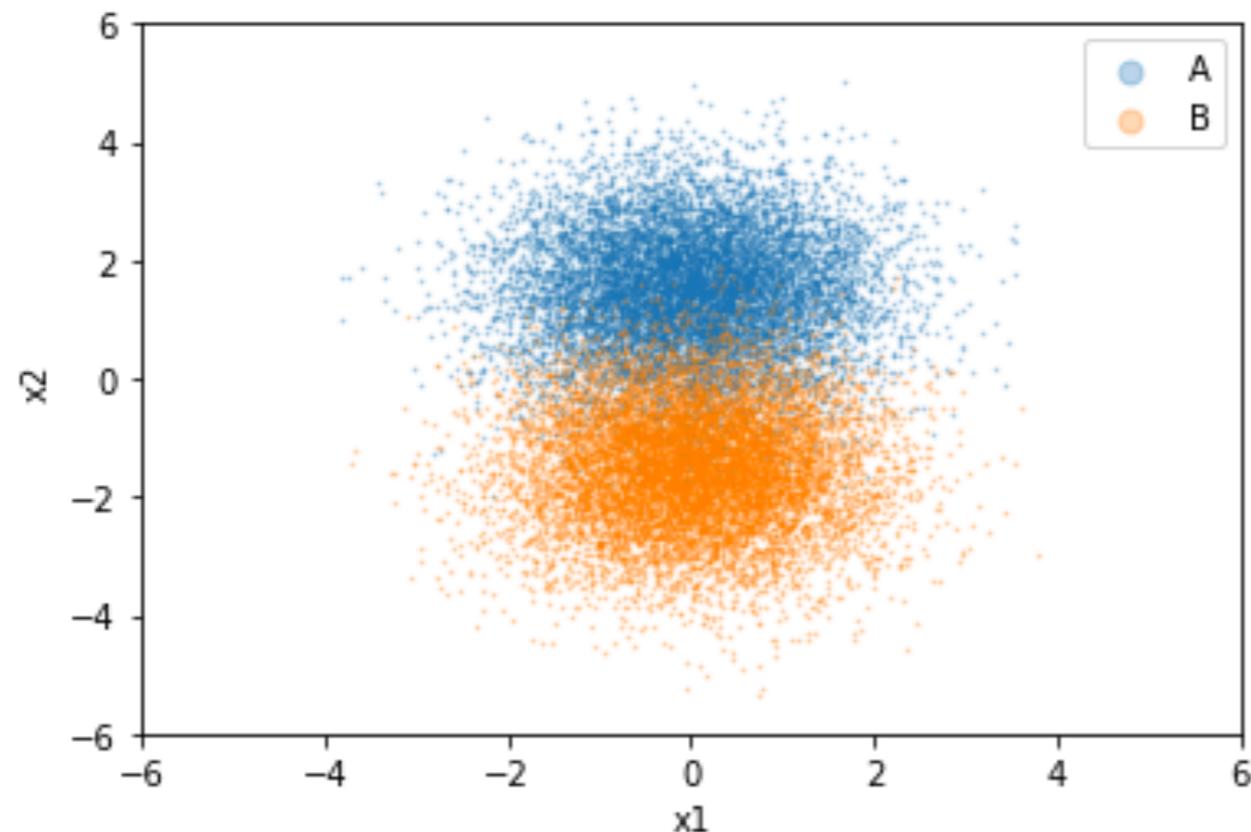


# Brief (re)fresher on machine learning

*“ML is glorified function fitting”*

Want to **fit** a **function**  $f(x; \theta)$  with some parameters  $\theta$  (“weights”) to a collection of examples  $\{x_i\}$  in order to achieve some **objective**.

Some typical **objectives**:



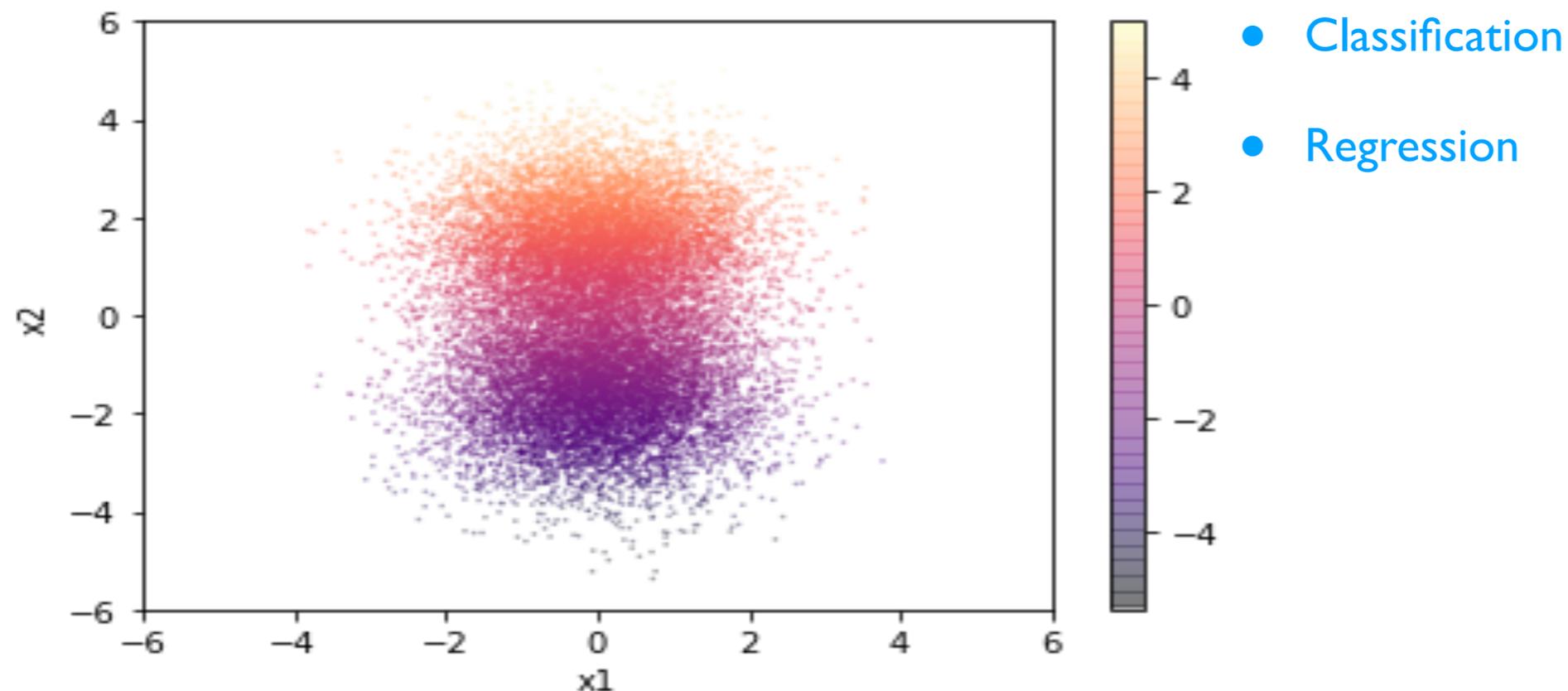
- Classification

# Brief (re)fresher on machine learning

*“ML is glorified function fitting”*

Want to **fit** a **function**  $f(\mathbf{x}; \theta)$  with some parameters  $\theta$  (“weights”) to a collection of examples  $\{\mathbf{x}_i\}$  in order to achieve some **objective**.

Some typical **objectives**:



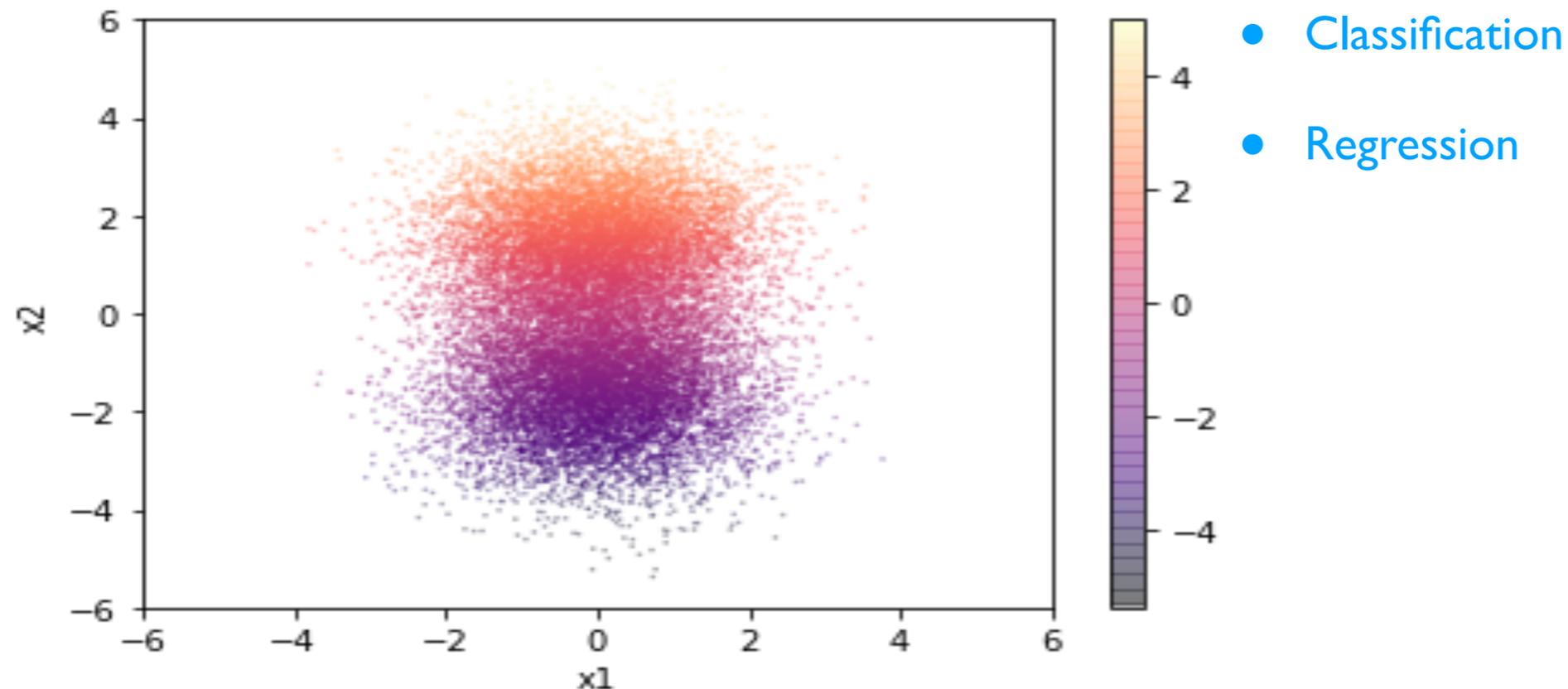
# Brief (re)fresher on machine learning

*“ML is glorified function fitting”*

Want to **fit** a **function**  $f(x; \theta)$  with some parameters  $\theta$  (“weights”) to a collection of examples  $\{x_i\}$  in order to achieve some **objective**.

Some typical **objectives**:

*“supervised ML”*



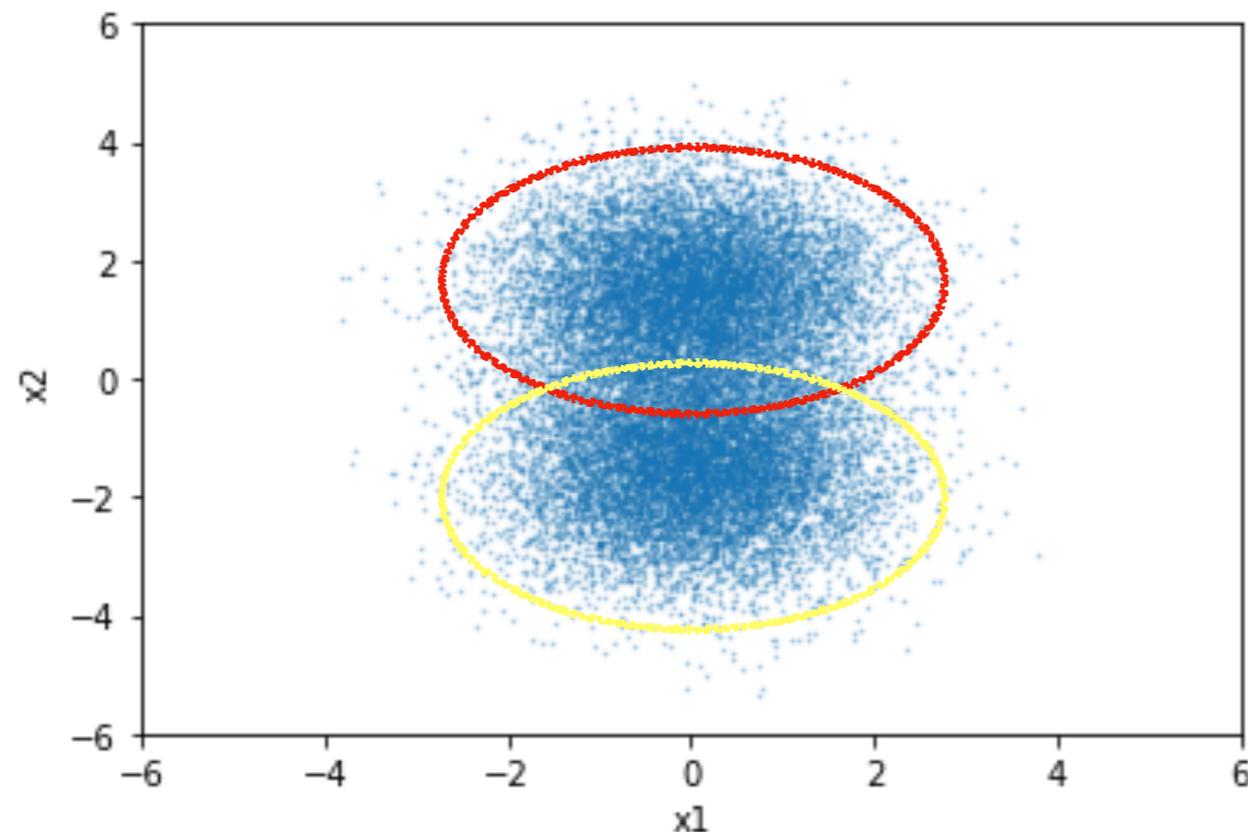
# Brief (re)fresher on machine learning

*“ML is glorified function fitting”*

Want to **fit** a **function**  $f(x; \theta)$  with some parameters  $\theta$  (“weights”) to a collection of examples  $\{x_i\}$  in order to achieve some **objective**.

Some typical **objectives**:

*“supervised ML”*



- Classification
- Regression
- Clustering

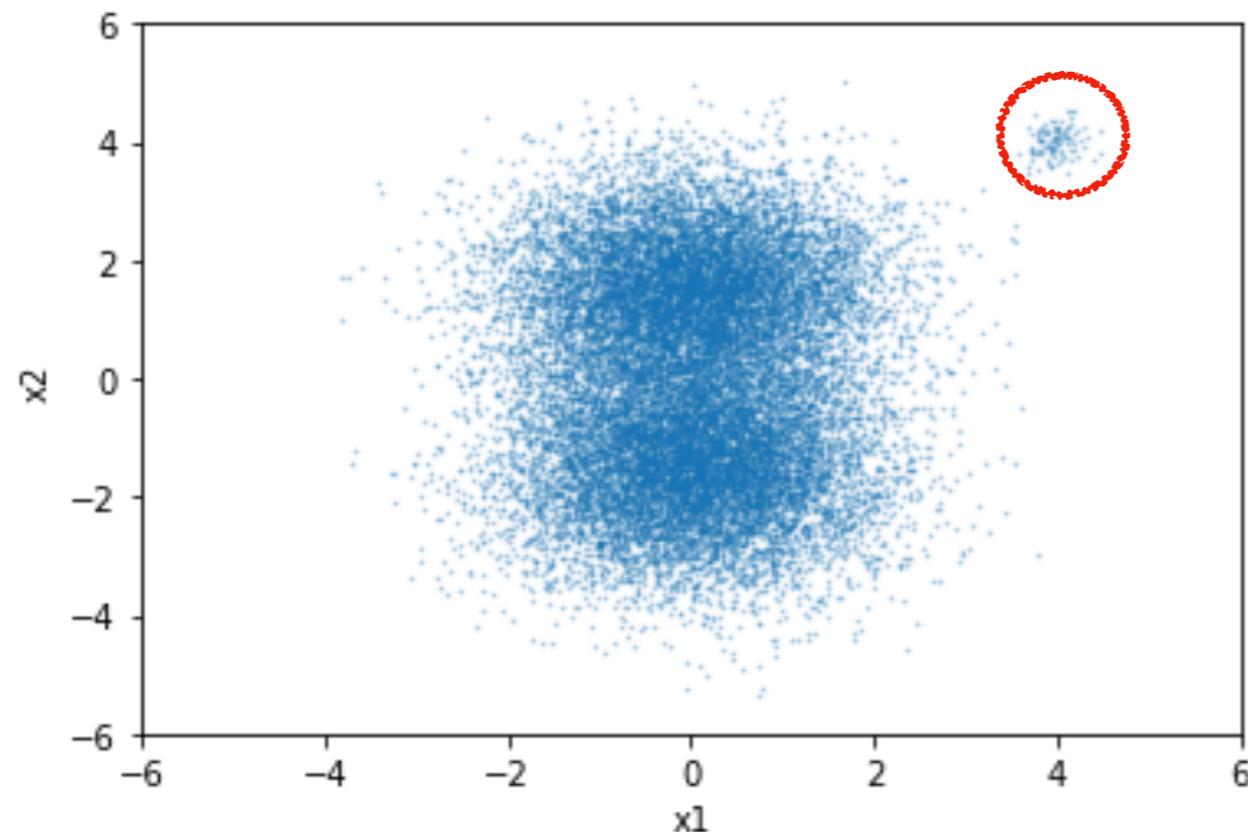
# Brief (re)fresher on machine learning

*“ML is glorified function fitting”*

Want to **fit** a **function**  $f(x; \theta)$  with some parameters  $\theta$  (“weights”) to a collection of examples  $\{x_i\}$  in order to achieve some **objective**.

Some typical **objectives**:

*“supervised ML”*



- Classification
- Regression
- Clustering
- Anomaly detection

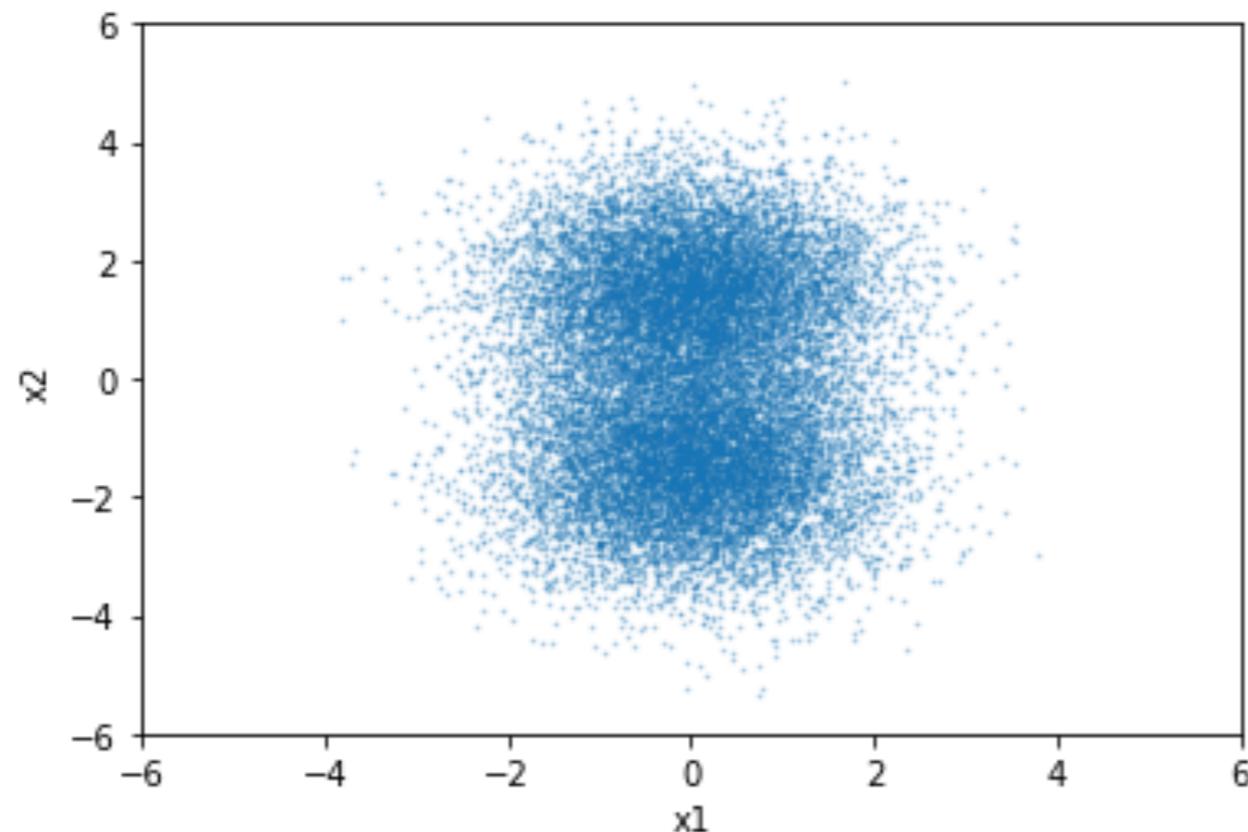
# Brief (re)fresher on machine learning

*“ML is glorified function fitting”*

Want to **fit** a **function**  $f(x; \theta)$  with some parameters  $\theta$  (“weights”) to a collection of examples  $\{x_i\}$  in order to achieve some **objective**.

Some typical **objectives**:

*“supervised ML”*



- Classification
- Regression
- Clustering
- Anomaly detection
- Density estimation

$$P(x) = 0.5 * \mathcal{N}(x|\mu = (0, -1.5), \sigma = (1, 1)) + 0.5 * \mathcal{N}(x|\mu = (0, +1.5), \sigma = (1, 1))$$

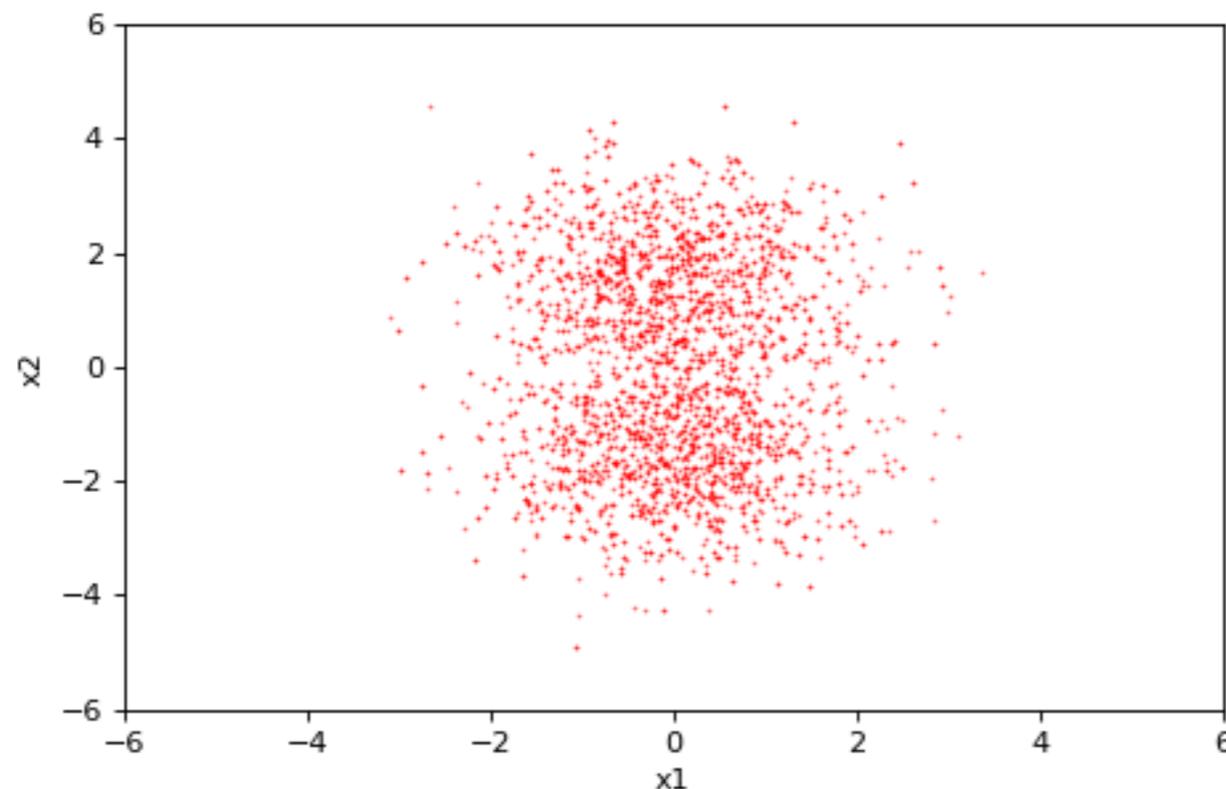
# Brief (re)fresher on machine learning

*“ML is glorified function fitting”*

Want to **fit** a **function**  $f(\mathbf{x}; \theta)$  with some parameters  $\theta$  (“weights”) to a collection of examples  $\{\mathbf{x}_i\}$  in order to achieve some **objective**.

Some typical **objectives**:

*“supervised ML”*



- Classification
- Regression
- Clustering
- Anomaly detection
- Density estimation
- Generation

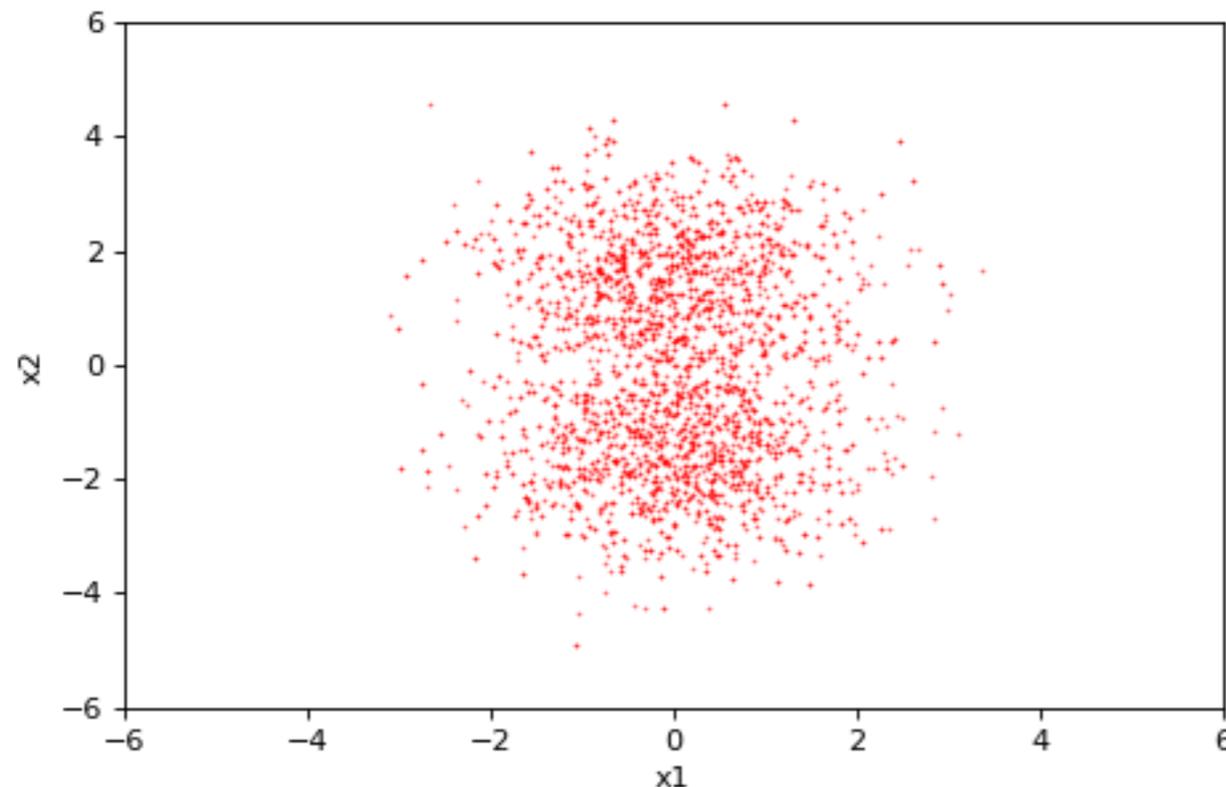
# Brief (re)fresher on machine learning

*“ML is glorified function fitting”*

Want to **fit** a **function**  $f(x; \theta)$  with some parameters  $\theta$  (“weights”) to a collection of examples  $\{x_i\}$  in order to achieve some **objective**.

Some typical **objectives**:

*“supervised ML”*



- Classification
- Regression
- Clustering
- Anomaly detection
- Density estimation
- Generation

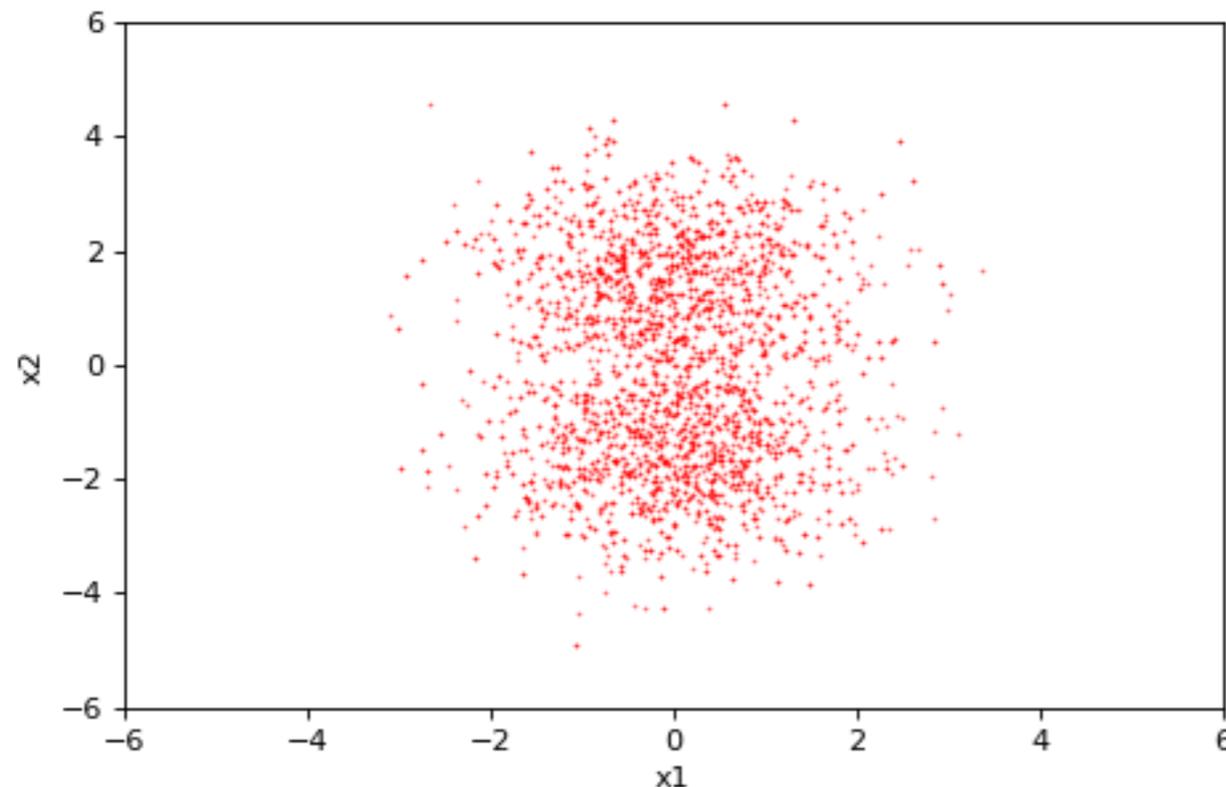
# Brief (re)fresher on machine learning

*“ML is glorified function fitting”*

Want to **fit** a **function**  $f(x; \theta)$  with some parameters  $\theta$  (“weights”) to a collection of examples  $\{x_i\}$  in order to achieve some **objective**.

Some typical **objectives**:

*“supervised ML”*

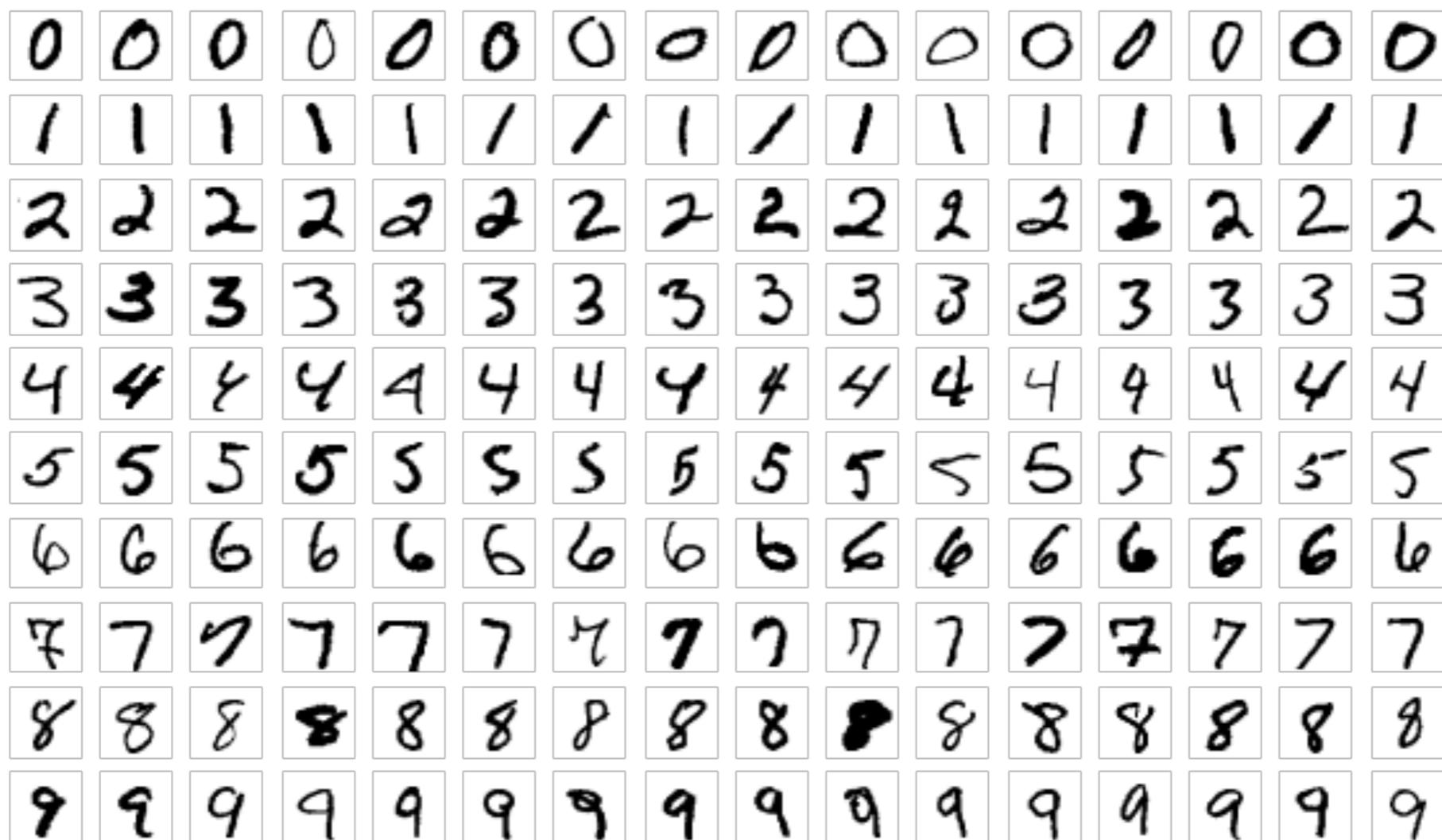


- Classification
- Regression
- Clustering
- Anomaly detection
- Density estimation
- Generation

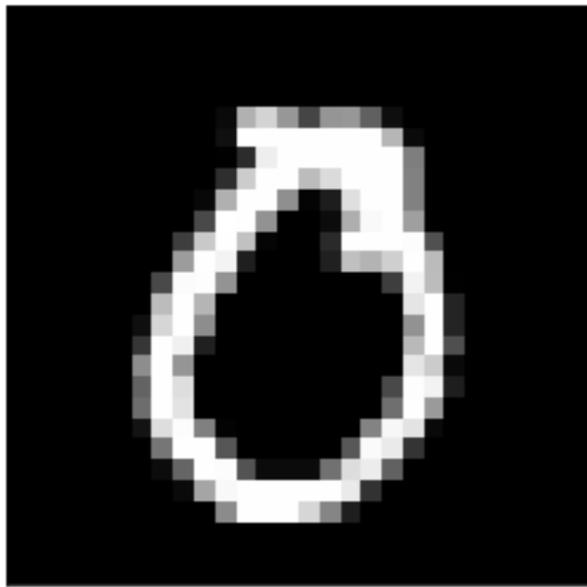
*“unsupervised ML”*

# MNIST example

image database of 70,000 handwritten digits



# MNIST example



$x$

Input: handwritten digit image  
28x28 pixel intensities  
from MNIST database



$\theta$

parameters of  
fitting function

$$\begin{pmatrix} P_0(x; \theta) \\ P_1(x; \theta) \\ \vdots \\ P_9(x; \theta) \end{pmatrix}$$

$$\sum_i P_i(x; \theta) = 1$$

Output: probability it's a 0, 1, ..., 9

# Brief (re)fresher on machine learning

*“ML is glorified function fitting”*

# Brief (re)fresher on machine learning

*“ML is glorified function fitting”*

Want to **fit** a **function**  $f(x; \theta)$  with some parameters  $\theta$  (“weights”) to a collection of examples  $\{x_i\}$  in order to achieve some **objective**.

# Brief (re)fresher on machine learning

*“ML is glorified function fitting”*

Want to **fit** a **function**  $f(x; \theta)$  with some parameters  $\theta$  (“weights”) to a collection of examples  $\{x_i\}$  in order to achieve some **objective**.

How to perform the **fit**?

# Brief (re)fresher on machine learning

*“ML is glorified function fitting”*

Want to **fit** a **function**  $f(x; \theta)$  with some parameters  $\theta$  (“weights”) to a collection of examples  $\{x_i\}$  in order to achieve some **objective**.

How to perform the **fit**?

**Minimize a loss function!** Loss function quantifies how well the objective has been achieved.

# Brief (re)fresher on machine learning

*“ML is glorified function fitting”*

Want to **fit** a **function**  $f(x; \theta)$  with some parameters  $\theta$  (“weights”) to a collection of examples  $\{x_i\}$  in order to achieve some **objective**.

How to perform the **fit**?

**Minimize a loss function!** Loss function quantifies how well the objective has been achieved.

$\sum_{(x,y)} L(P(x; \theta), y)$	$\sum_x L(P(x; \theta))$
Data is labeled supervised ML	Data is not labeled unsupervised ML

# Examples of loss functions

$$L = (P(x; \theta) - y)^2$$

“mean-squared error”  
used for regression

$$L = -\left(y \log P(x; \theta) + (1 - y) \log(1 - P(x; \theta))\right)$$

“binary cross entropy”  
used for binary classification

$$L = -y_j \log P_j(x; \theta)$$

“categorical cross entropy”  
used for multi-class classification

# Examples of loss functions

$$L = (P(x; \theta) - y)^2$$

“mean-squared error”  
used for regression

$$L = -\left(y \log P(x; \theta) + (1 - y) \log(1 - P(x; \theta))\right)$$

“binary cross entropy”  
used for binary classification

$$L = -y_j \log P_j(x; \theta)$$

MNIST example:  
labels “one-hot encoding”

$$0 \rightarrow [1, 0, \dots, 0]$$

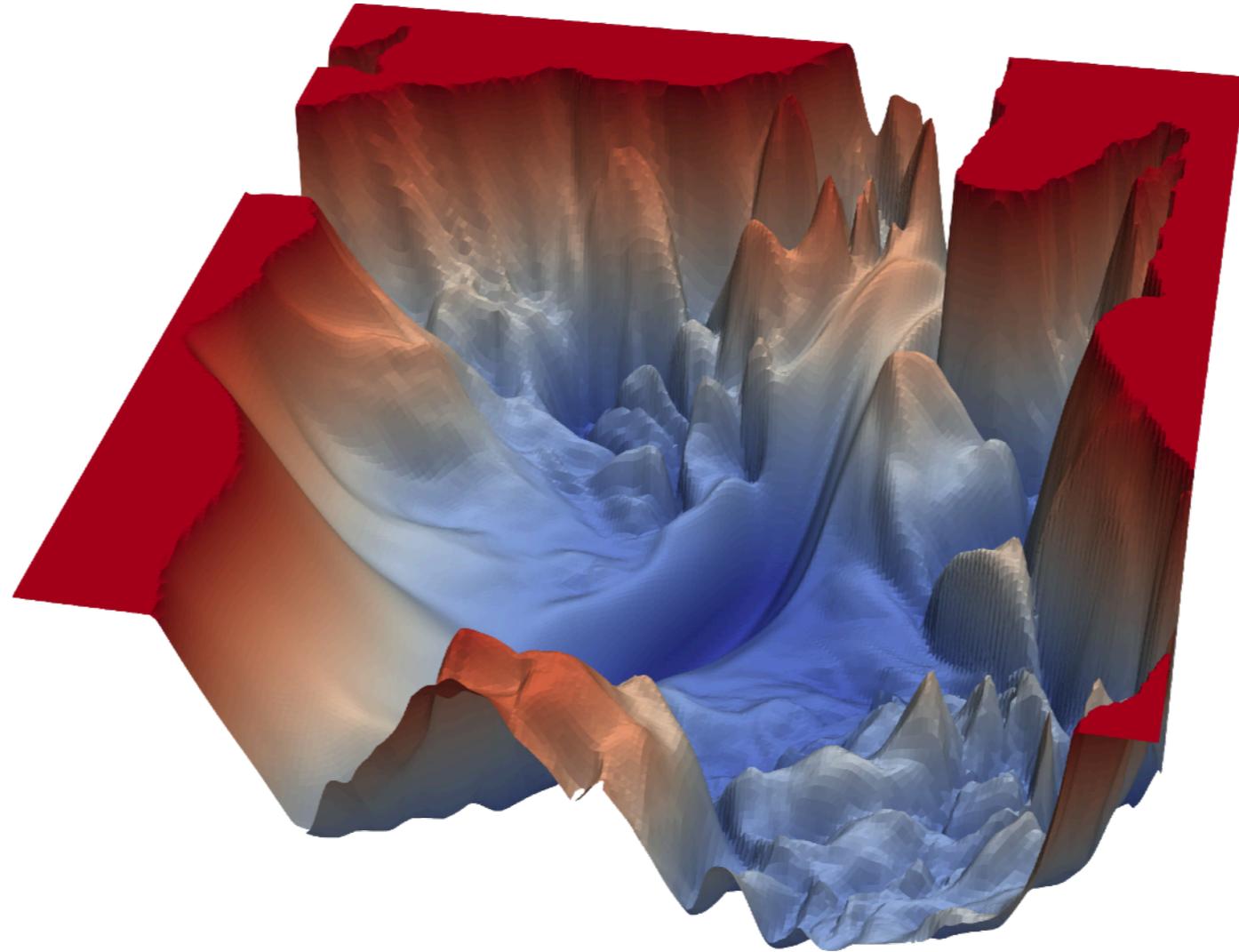
$$1 \rightarrow [0, 1, \dots, 0]$$

...

$$9 \rightarrow [0, 0, \dots, 1]$$

“categorical cross entropy”  
used for multi-class classification

# Minimizing the loss function



Highly nonconvex function over a many-dimensional space.  
Many local minima.

# Gradient descent

$$\langle L \rangle(\theta) \equiv \sum_{(x,y)} L(P(x; \theta), y)$$

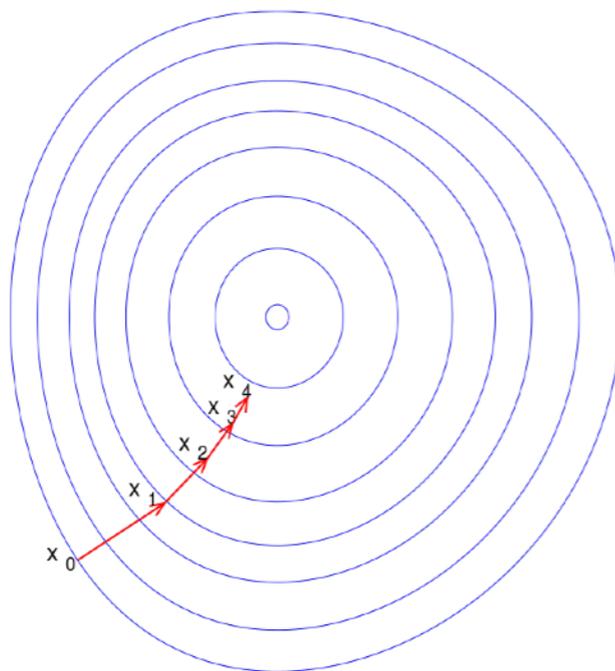
Want to minimize wrt  $\theta$ .

Obvious idea:  $\theta \rightarrow \theta - \alpha \partial_{\theta} \langle L \rangle(\theta)$

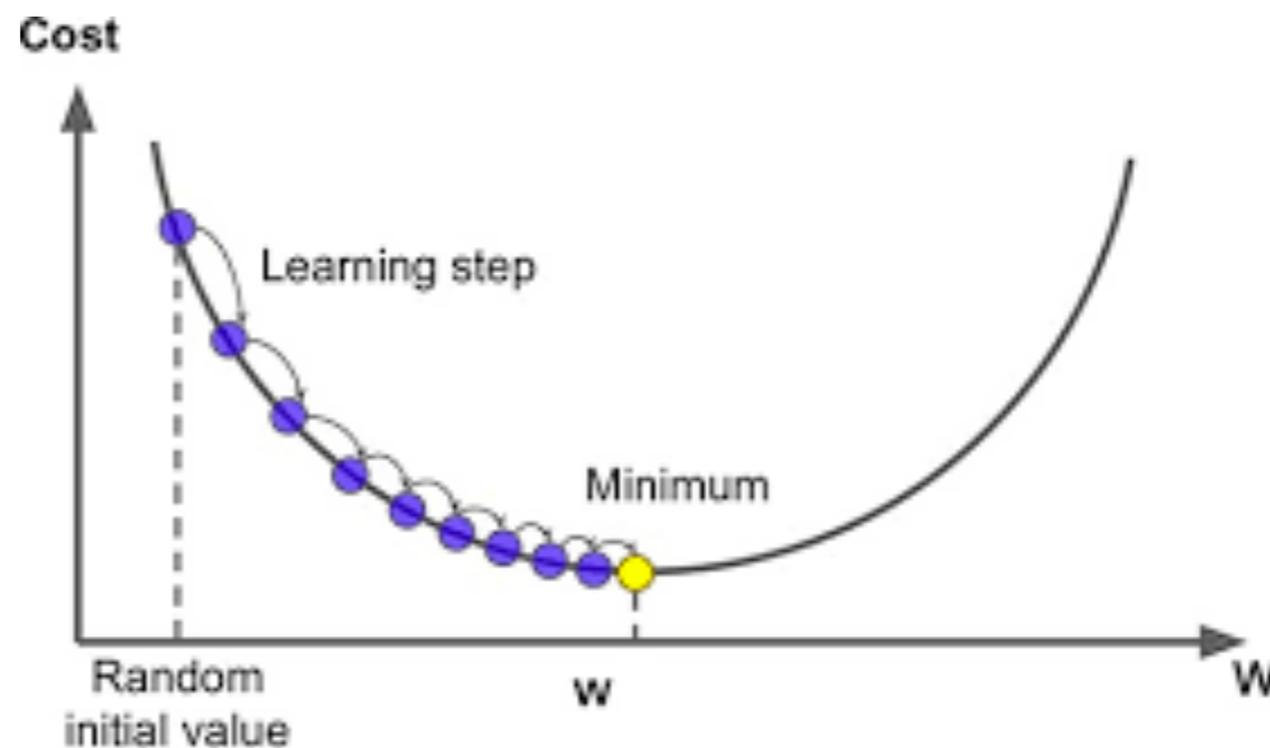
$\alpha$  : “learning rate”

“gradient descent method”

(generalization of Newton-Raphson method)



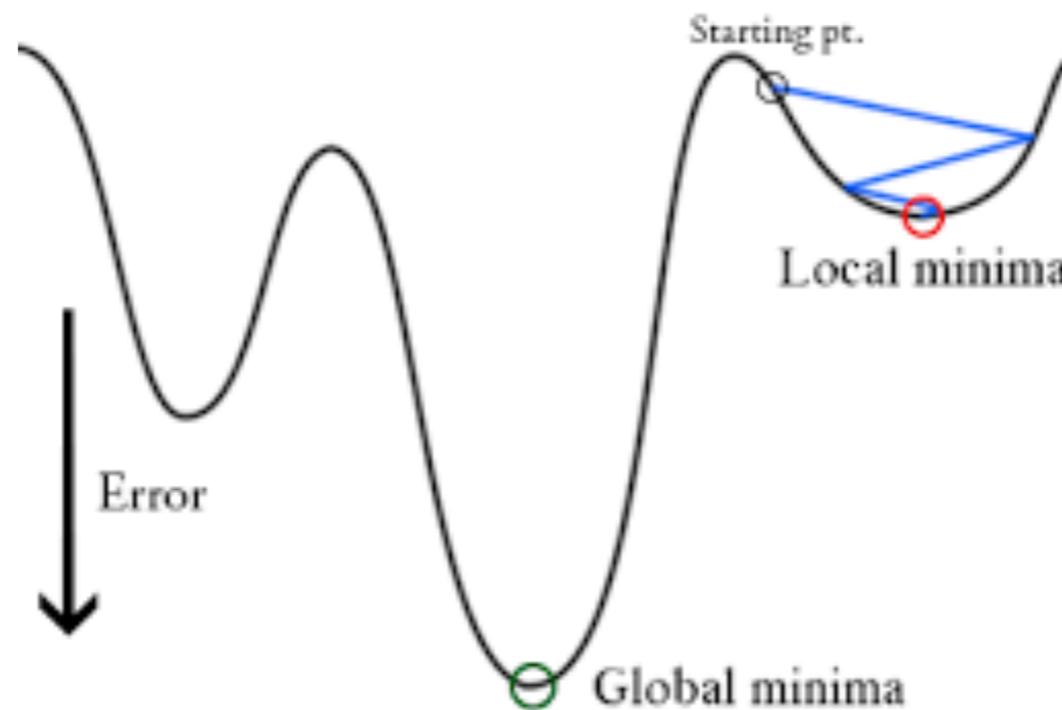
source:Wikipedia



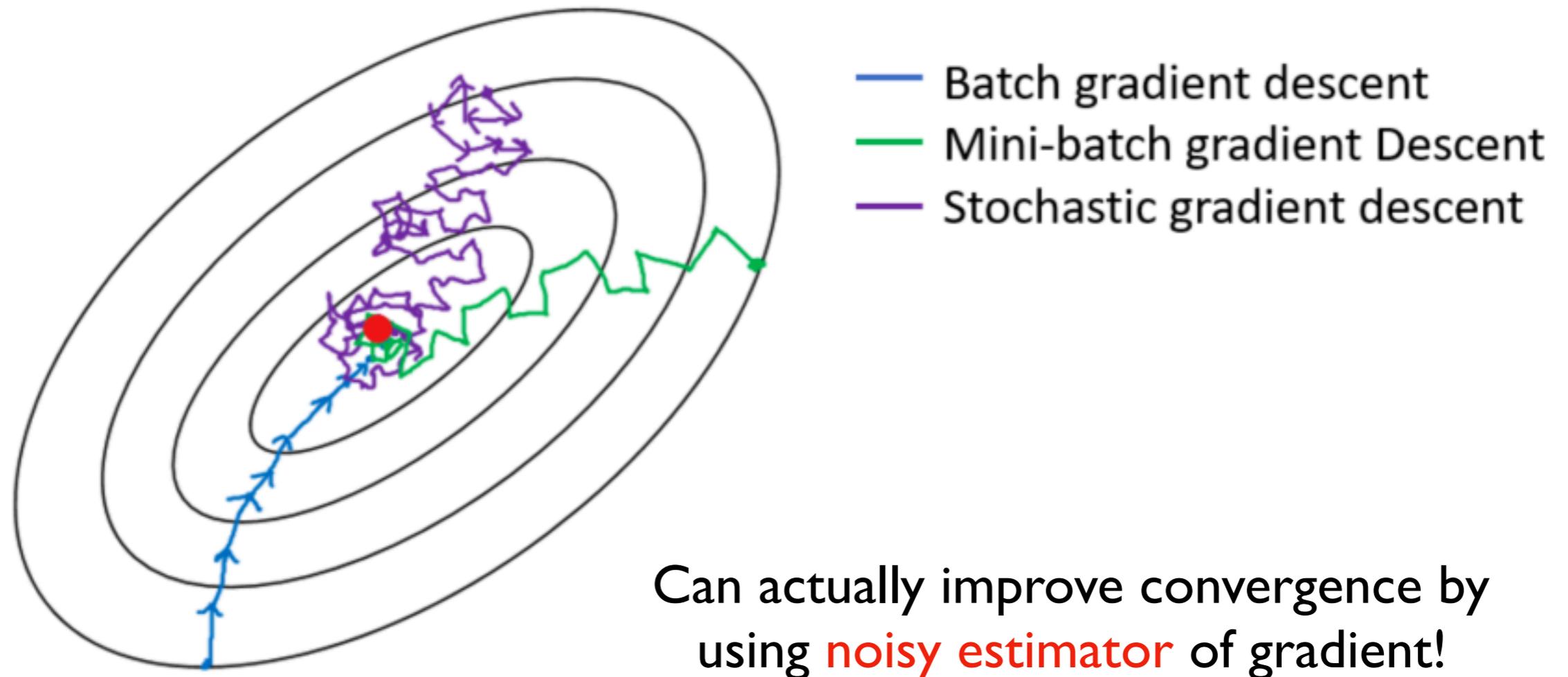
# Gradient descent: problems

## Downsides to gradient descent:

- average over full dataset  $\langle L(\theta) \rangle$  can be expensive to compute
- poor initial guess or learning rate can lead to becoming stuck in poor local minimum



# Stochastic Gradient Descent



$$\sum_{(x,y) \in \text{full dataset}} L(P(x; \theta), y)$$



$$\sum_{(x,y) \in \text{minibatch}} L(P(x; \theta), y)$$

# Stochastic Gradient Descent

Method:

1. Divide up training data into minibatches.
2. Update weights minibatch by minibatch

$$\theta \rightarrow \theta - \alpha \partial_{\theta} \langle L \rangle (\theta)$$

(average computed on each minibatch)

3. Repeat until convergence.



minibatch 1  
minibatch 2  
minibatch 3

...

# Stochastic Gradient Descent

Method:

1. Divide up training data into minibatches.
2. Update weights minibatch by minibatch

$$\theta \rightarrow \theta - \alpha \partial_{\theta} \langle L \rangle (\theta)$$

(average computed on each minibatch)

3. Repeat until convergence.



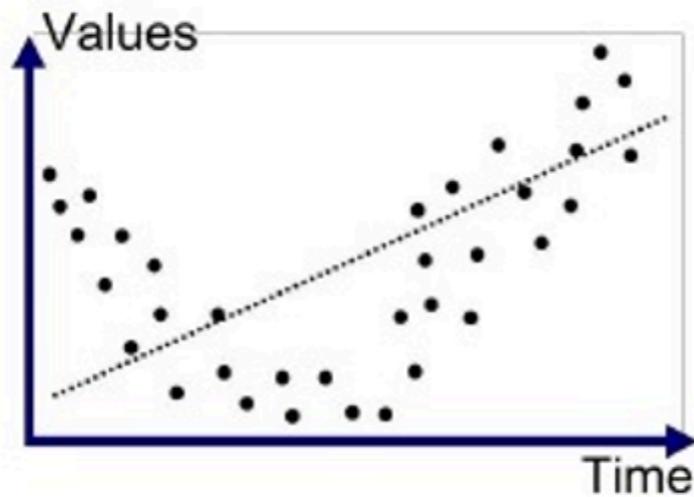
minibatch 1  
minibatch 2  
minibatch 3

...

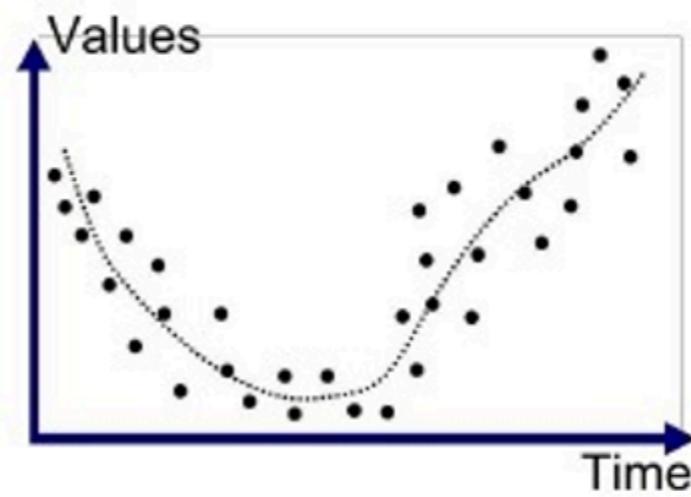
“one epoch”

# Overfitting

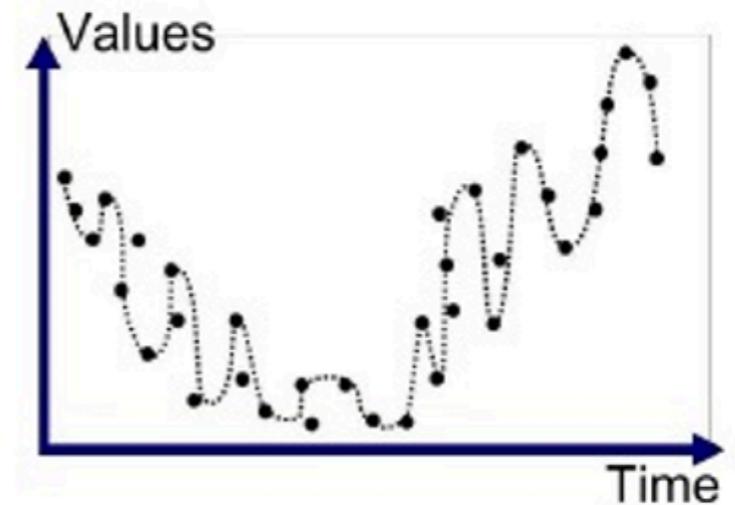
The fitting function (especially if it is a neural network) may be overparametrized. So overfitting is a major concern.



Underfitted



Good Fit/Robust

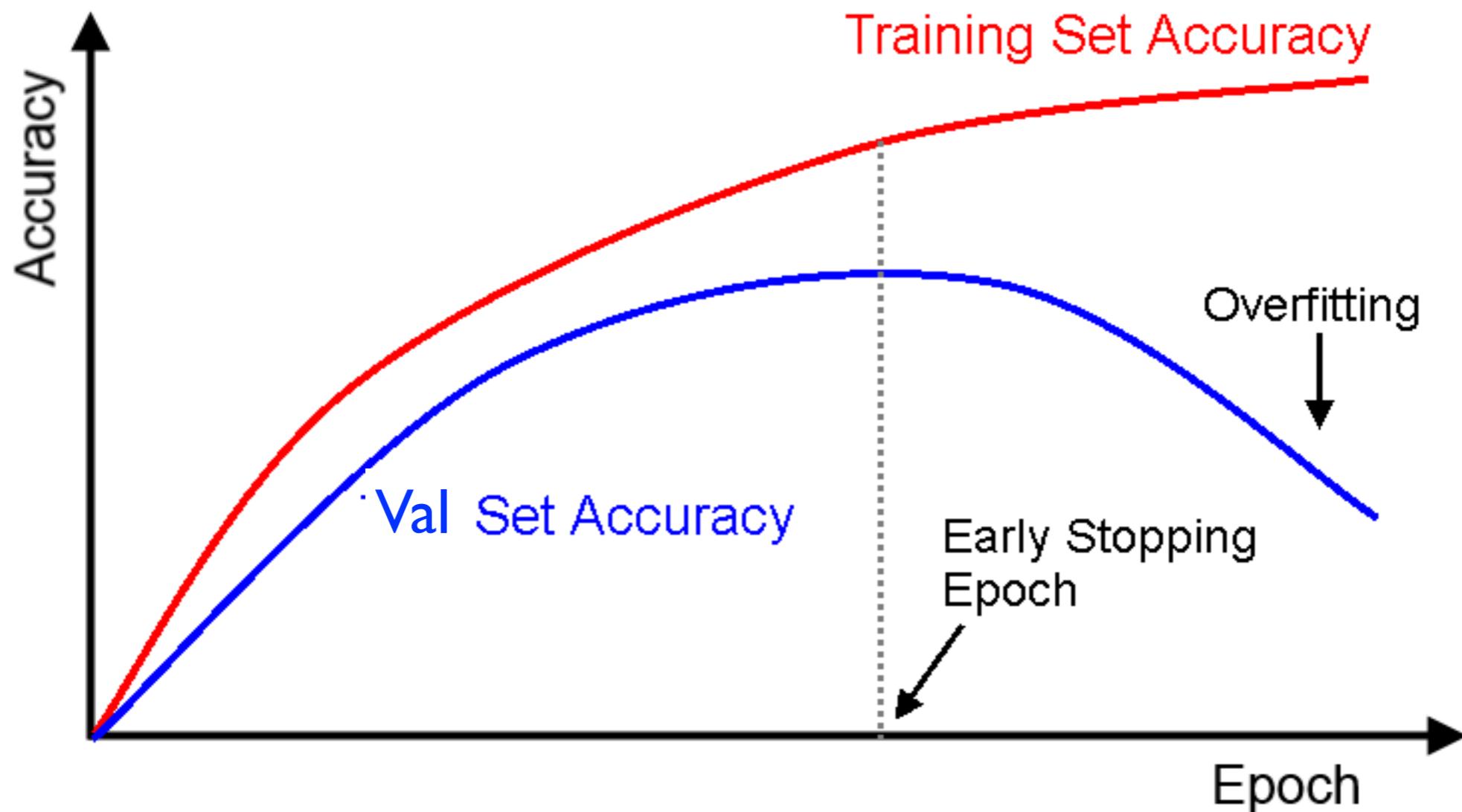


Overfitted

# Overfitting

Many ways to mitigate overfitting problem. Eg early stopping.

Key concept: **train/val/test split**



# Brief (re)fresher on machine learning

*“ML is glorified function fitting”*

Want to **fit** a **function**  $f(x; \theta)$  with some parameters  $\theta$  (“weights”) to a collection of examples  $\{x_i\}$  in order to achieve some **objective**.

# Brief (re)fresher on machine learning

*“ML is glorified function fitting”*

Want to **fit** a **function**  $f(x; \theta)$  with some parameters  $\theta$  (“weights”) to a collection of examples  $\{x_i\}$  in order to achieve some **objective**.

Finally, what **functions** to use?

# Brief (re)fresher on machine learning

*“ML is glorified function fitting”*

Want to **fit** a **function**  $f(x; \theta)$  with some parameters  $\theta$  (“weights”) to a collection of examples  $\{x_i\}$  in order to achieve some **objective**.

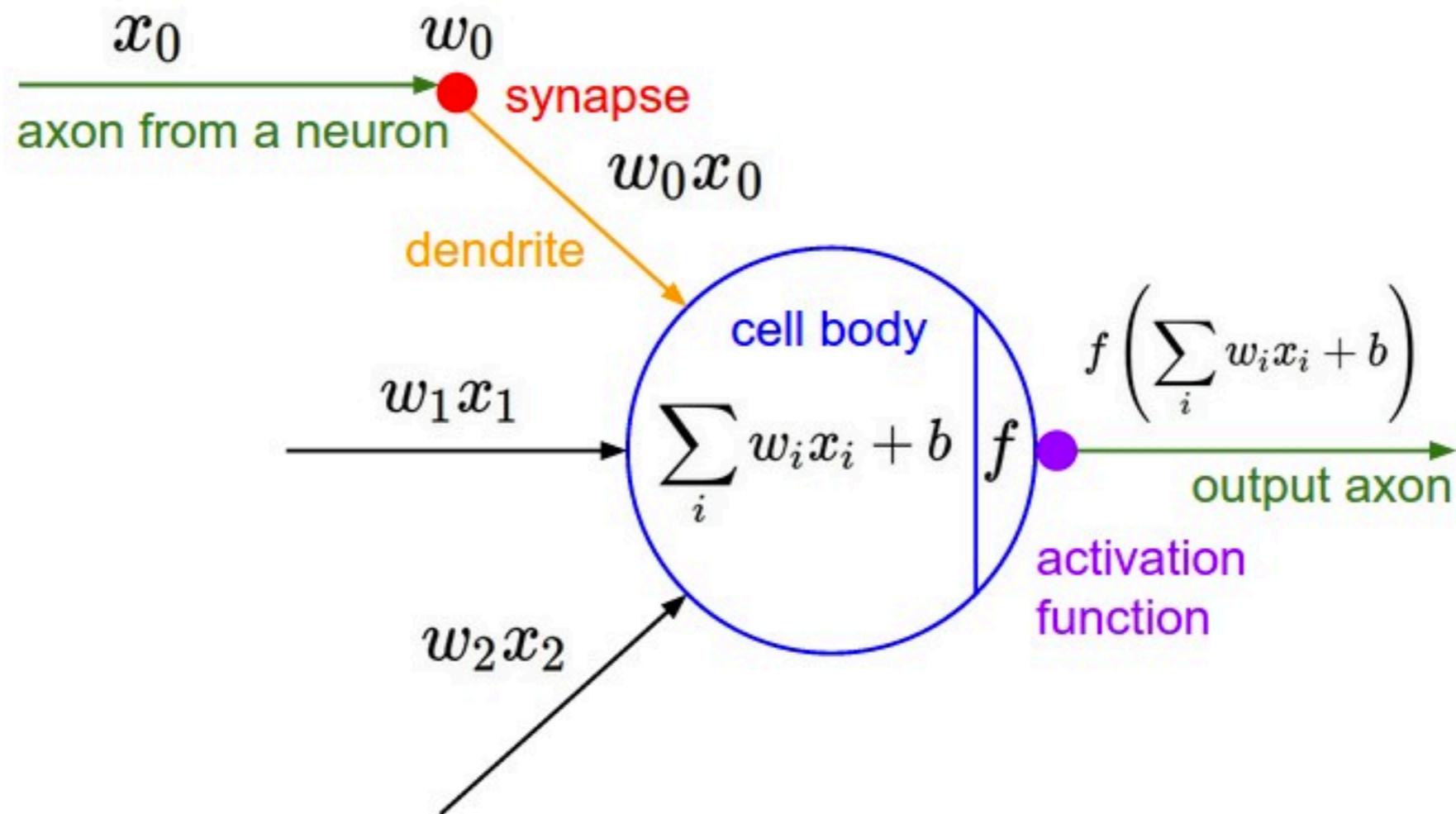
Finally, what **functions** to use?

**Current trend: deep neural networks!**

# 2. Intro to Deep Learning

# What is a (deep) neural network?

Basic building block of a neural network:

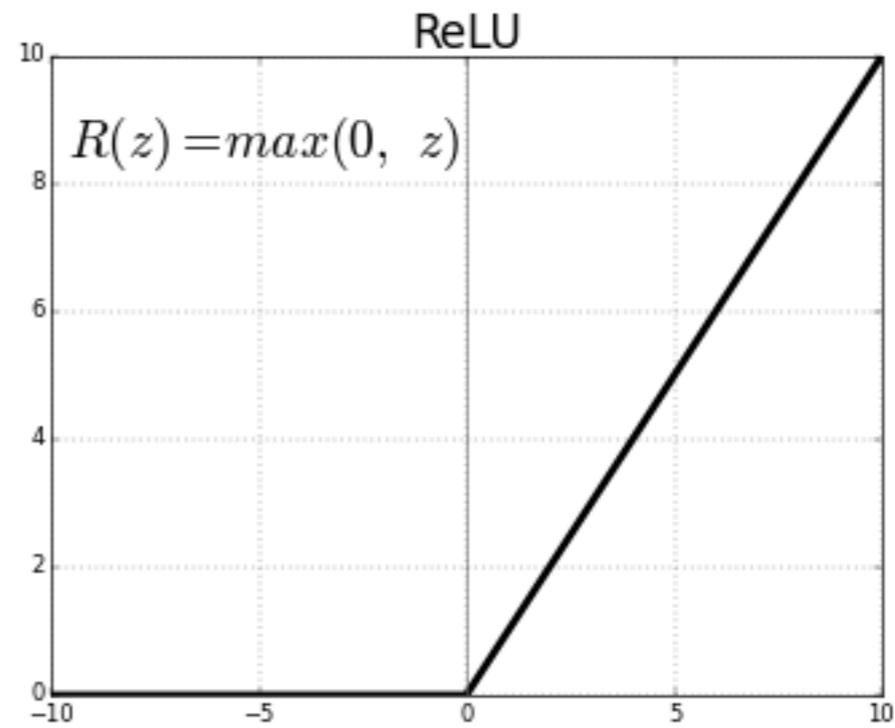
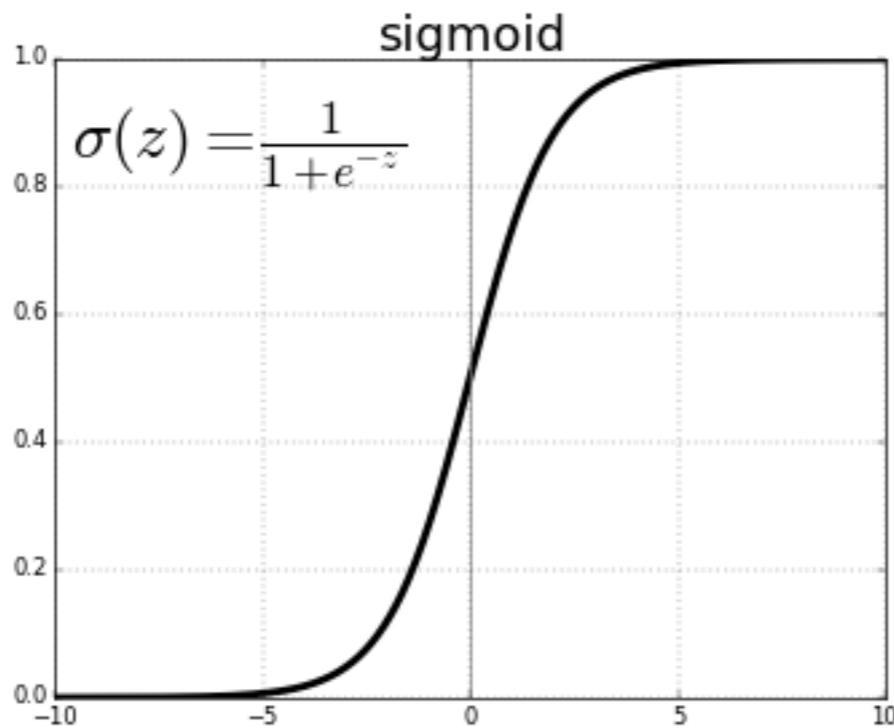


Modeled on biological systems

# Activation functions

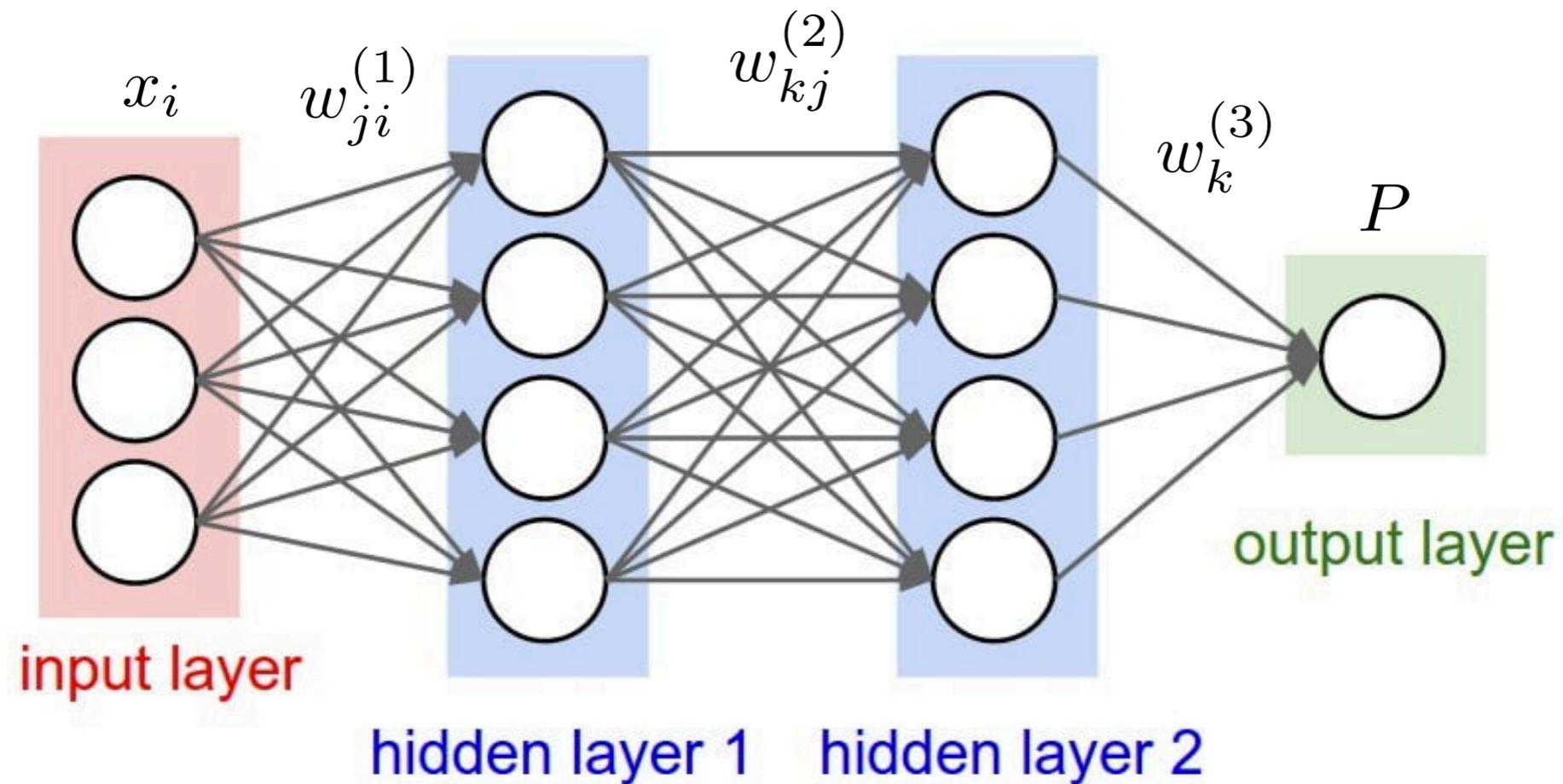
NNs need a source of non-linearity so they can learn general functions.

This is usually implemented with the activation function.



Sigmoid used to be standard. But this led to the **vanishing gradient problem**. The ReLU activation was invented to solve this problem. Now it is the standard.

# “Fully connected” or “Dense” Neural Network



$$P = A^{(3)} (w_k^{(3)} A^{(2)} (w_{kj}^{(2)} A^{(1)} (w_{ji}^{(1)} x_i + b_j^{(1)}) + b_k^{(2)}) + b_k^{(3)})$$

# Expressiveness

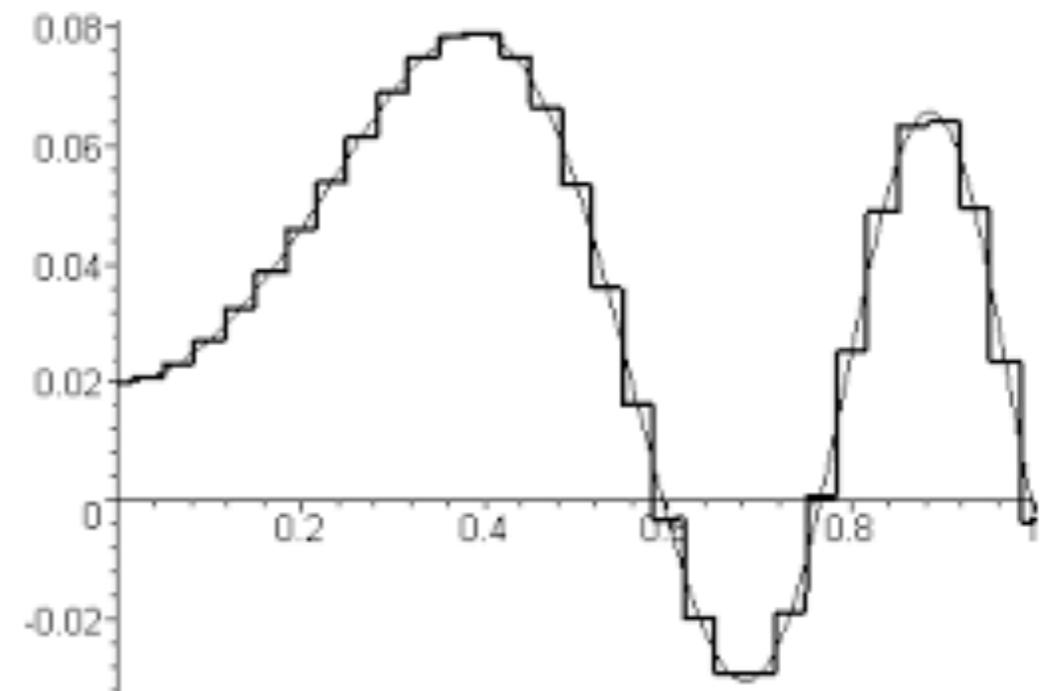
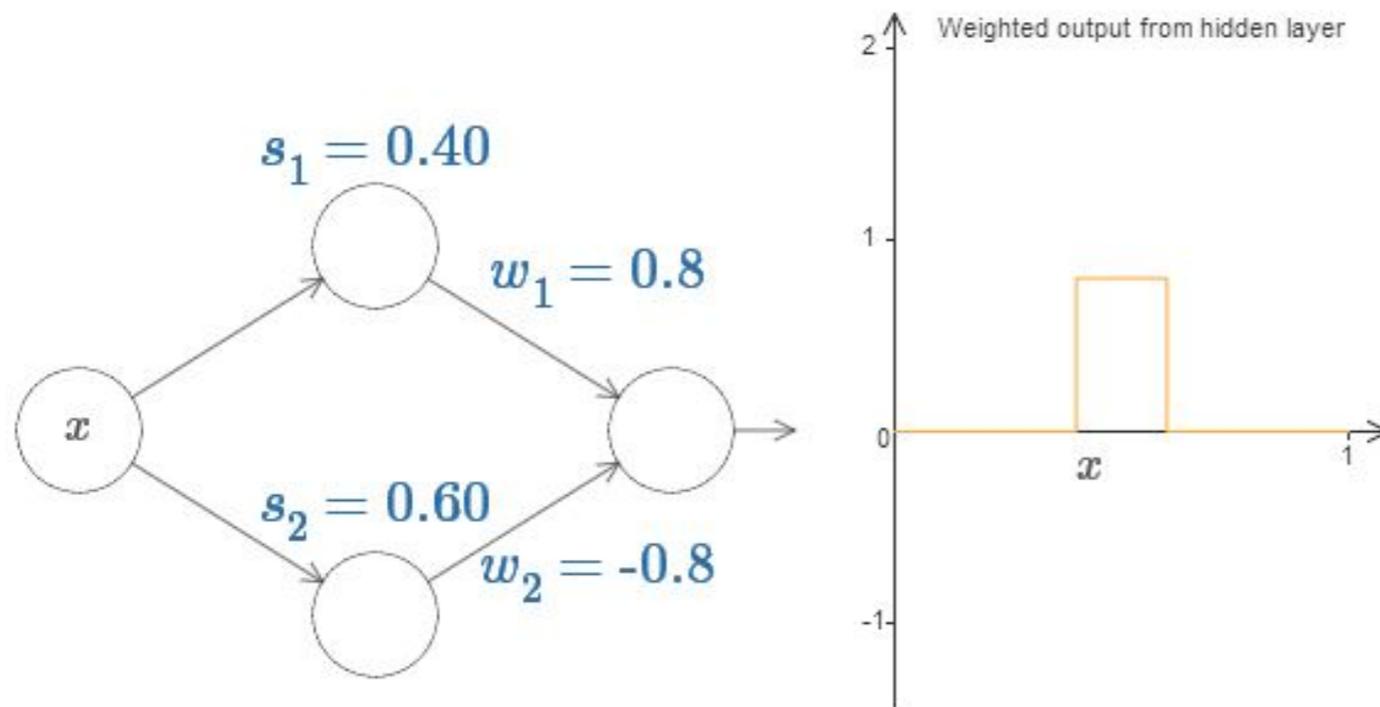
One reason why NNs “work” is that they are universal function approximators (asymptotically)

## Approximation by Superpositions of a Sigmoidal Function\*

G. Cybenko†

**Abstract.** In this paper we demonstrate that finite linear combinations of compositions of a fixed, univariate function and a set of affine functionals can uniformly approximate any continuous function of  $n$  real variables with support in the unit hypercube; only mild conditions are imposed on the univariate function. Our results settle an open question about representability in the class of single hidden layer neural networks. In particular, we show that arbitrary decision regions can be arbitrarily well approximated by continuous feedforward neural networks with only a single internal, hidden layer and any continuous sigmoidal nonlinearity. The paper discusses approximation properties of other possible types of nonlinearities that might be implemented by artificial neural networks.

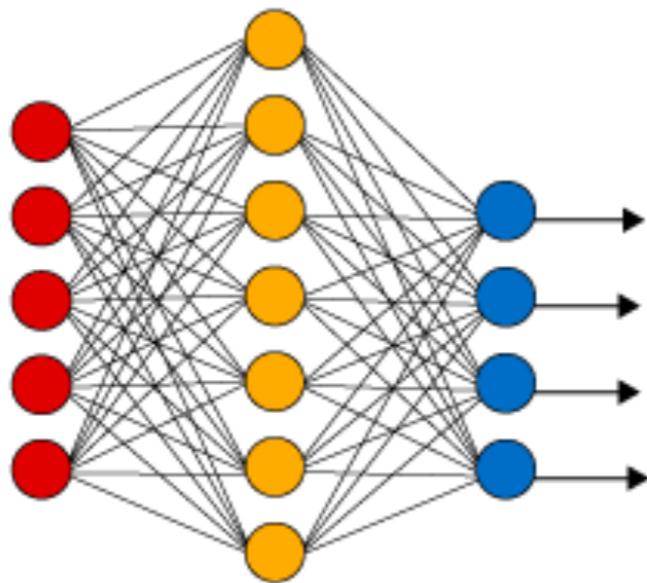
**Key words.** Neural networks, Approximation, Completeness.



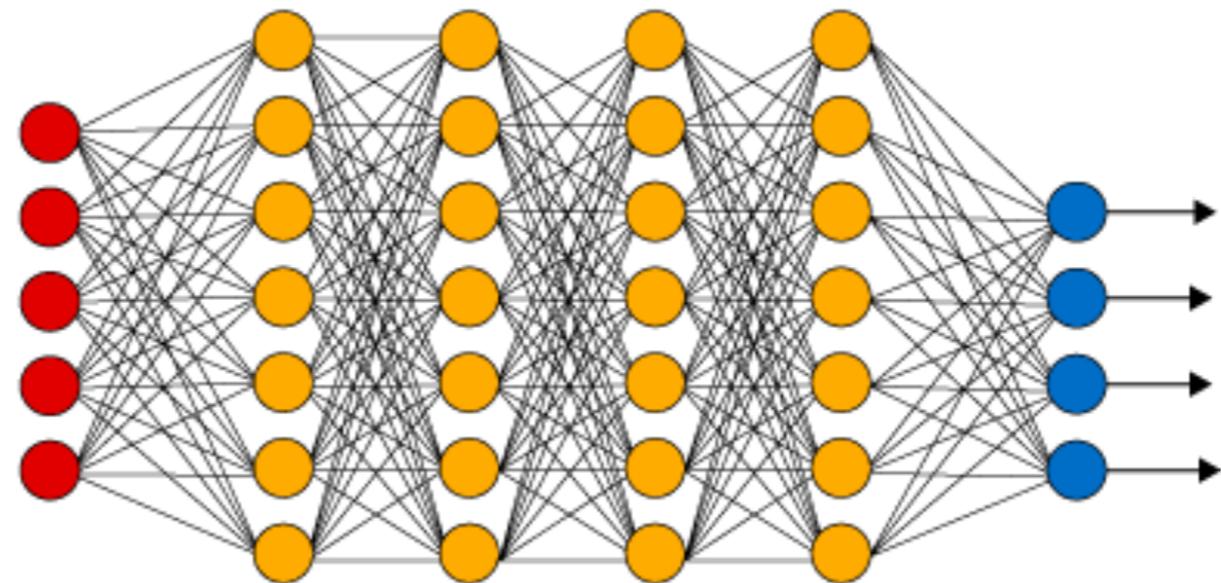
Can fit any function with a single, infinitely-wide hidden layer

# Deep neural networks

Simple Neural Network



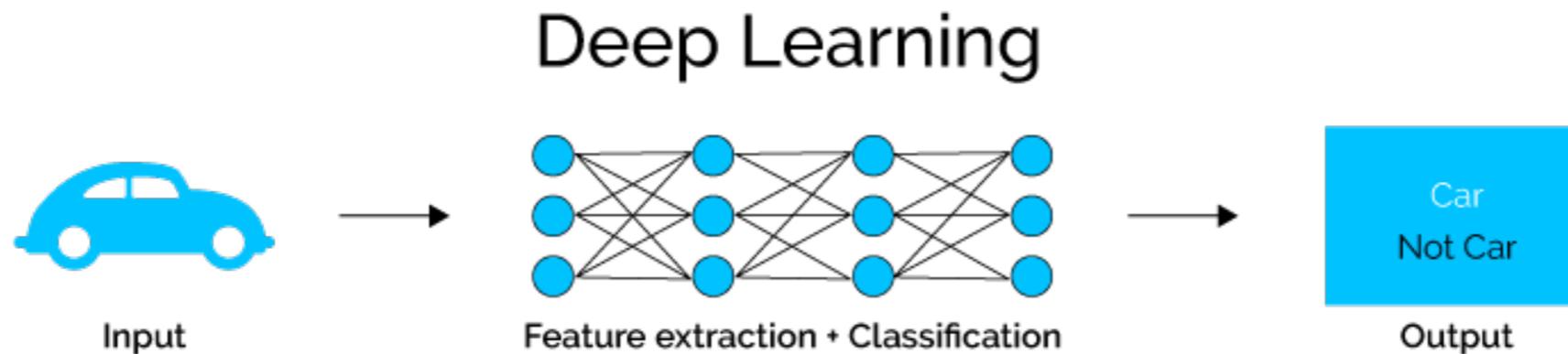
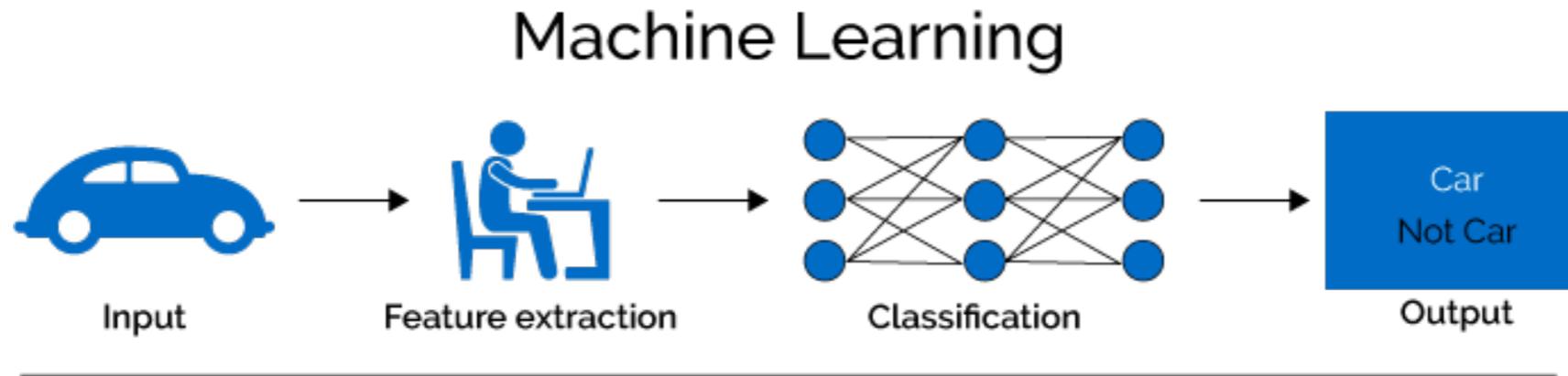
Deep Learning Neural Network



● Input Layer    ● Hidden Layer    ● Output Layer

Expressivity of a neural network increases exponentially with the number of layers ([Delalleau and Bengio, 2011](#), [Montufar et al 1402.1869](#), [Ganguli et al 1606.05336, 1606.05340](#))

# End-to-end learning



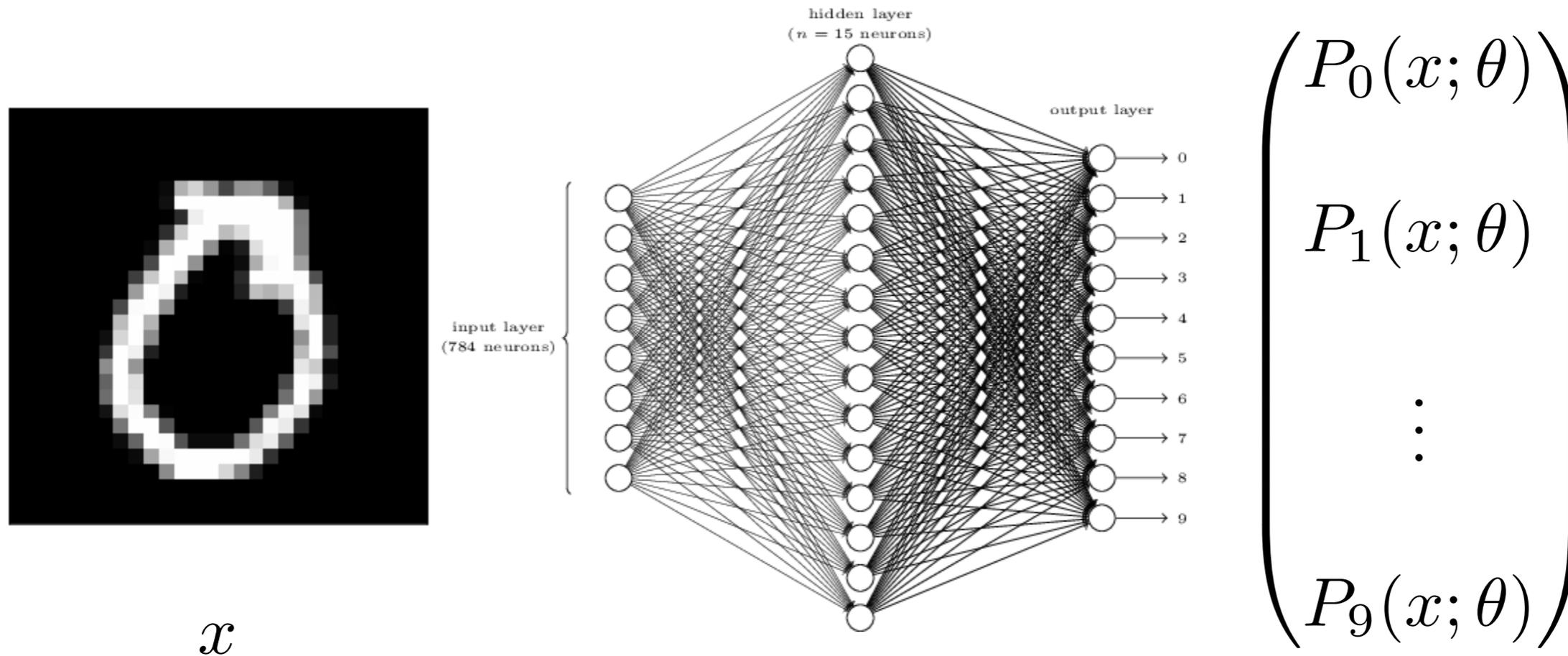
From towardsdatascience.com

Universal function approximation and high expressivity means that deep NNs can learn abstract concepts from low-level, high-dimensional inputs

*“Automated feature engineering”*

*“End-to-end learning”*

# Example: *MNIST* in more detail



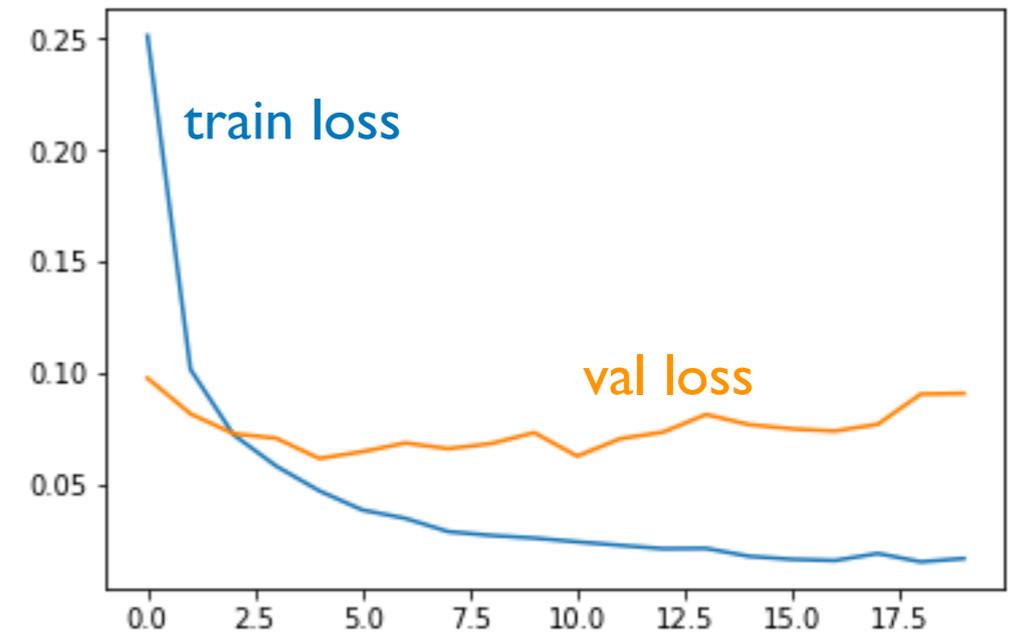
$$\sum_i P_i(x; \theta) = 1$$

Output: probability it's a 0, 1, ..., 9

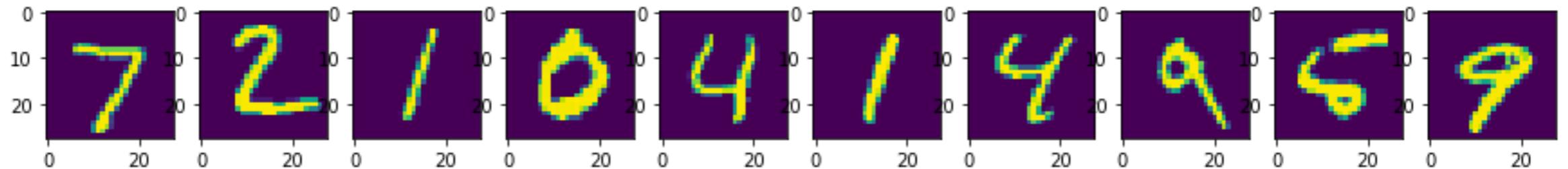
# Example: MNIST in more detail

Layer (type)	Output Shape	Param #
dense_1 (Dense)	(None, 512)	401920
dropout_1 (Dropout)	(None, 512)	0
dense_2 (Dense)	(None, 512)	262656
dropout_2 (Dropout)	(None, 512)	0
dense_3 (Dense)	(None, 10)	5130

Total params: 669,706  
Trainable params: 669,706  
Non-trainable params: 0

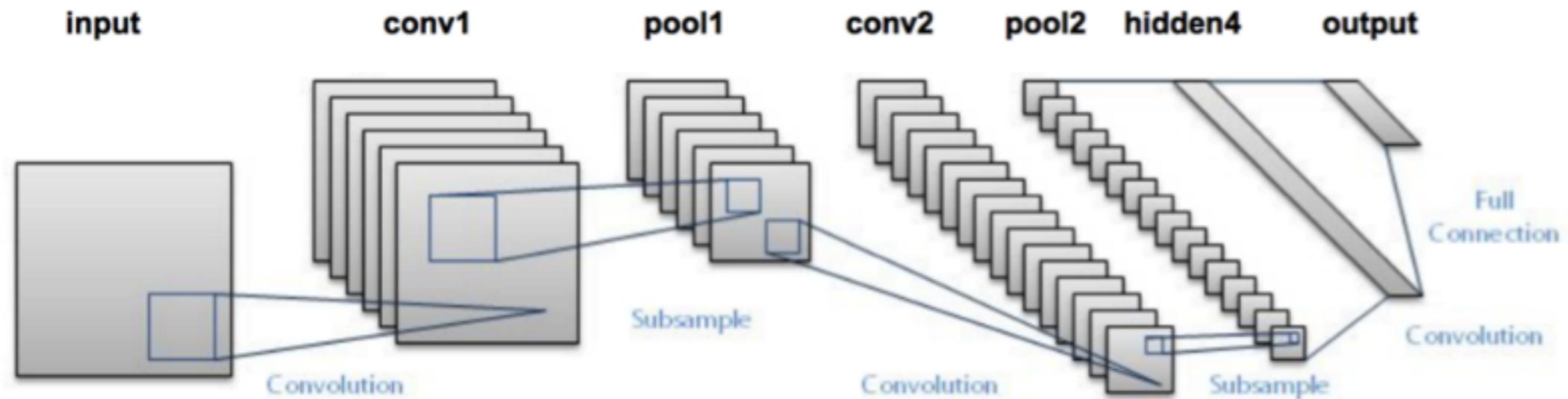


Test accuracy: 0.9831



prediction: 7 2 1 0 4 1 4 9 6 9

# Convolutional neural network (CNN)



Principal neural network architecture for image recognition.  
Invented in 1998 (LeCun, Bottou, Bengio, Haffner)

Achieved 99% accuracy on MNIST!

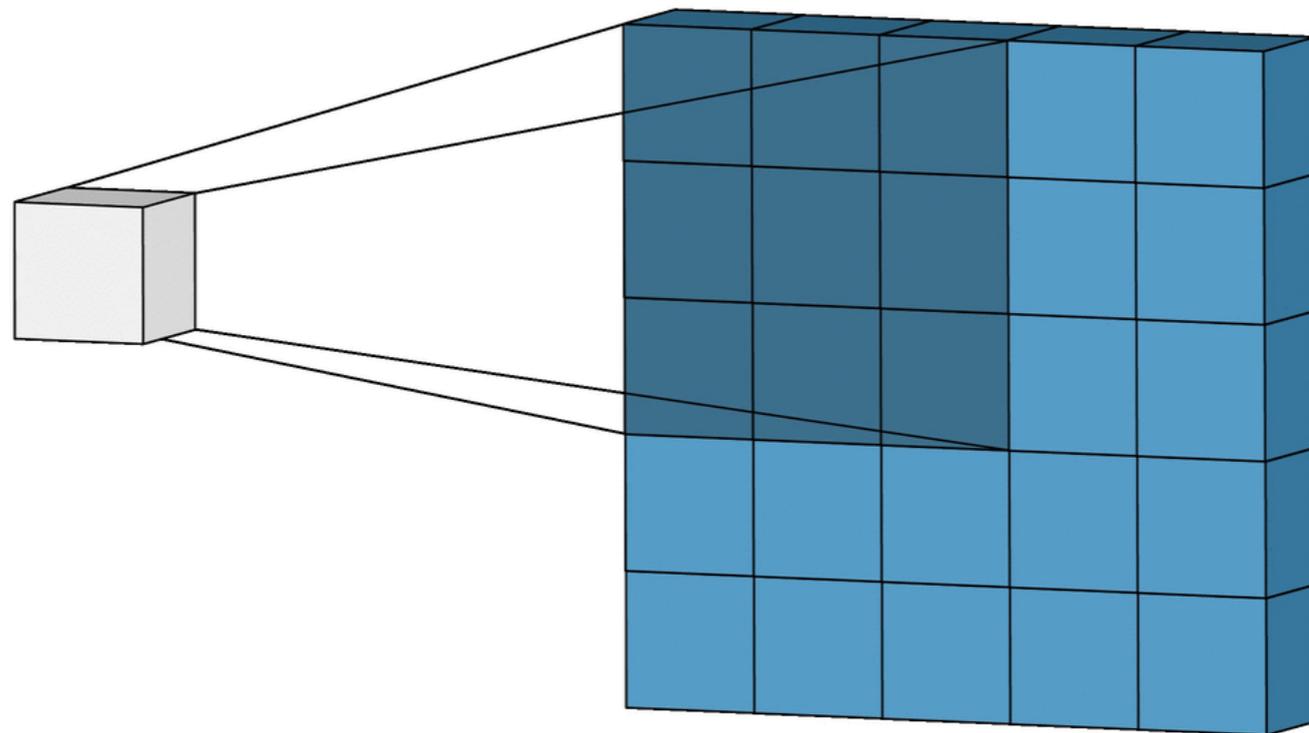
However, CNNs fell out of favor (until AlexNet in 2012) when they did not immediately generalize well to more complex image recognition tasks such as ImageNet.

# Convolutional neural network (CNN)

Main idea: features in an image (edges, curves, corners,...eyes, noses,...) are the same no matter where they occur.

Goal: Want to find these features in a translationally invariant way.

Solution: Drag or convolve a “filter” across the image that selects out interesting features.

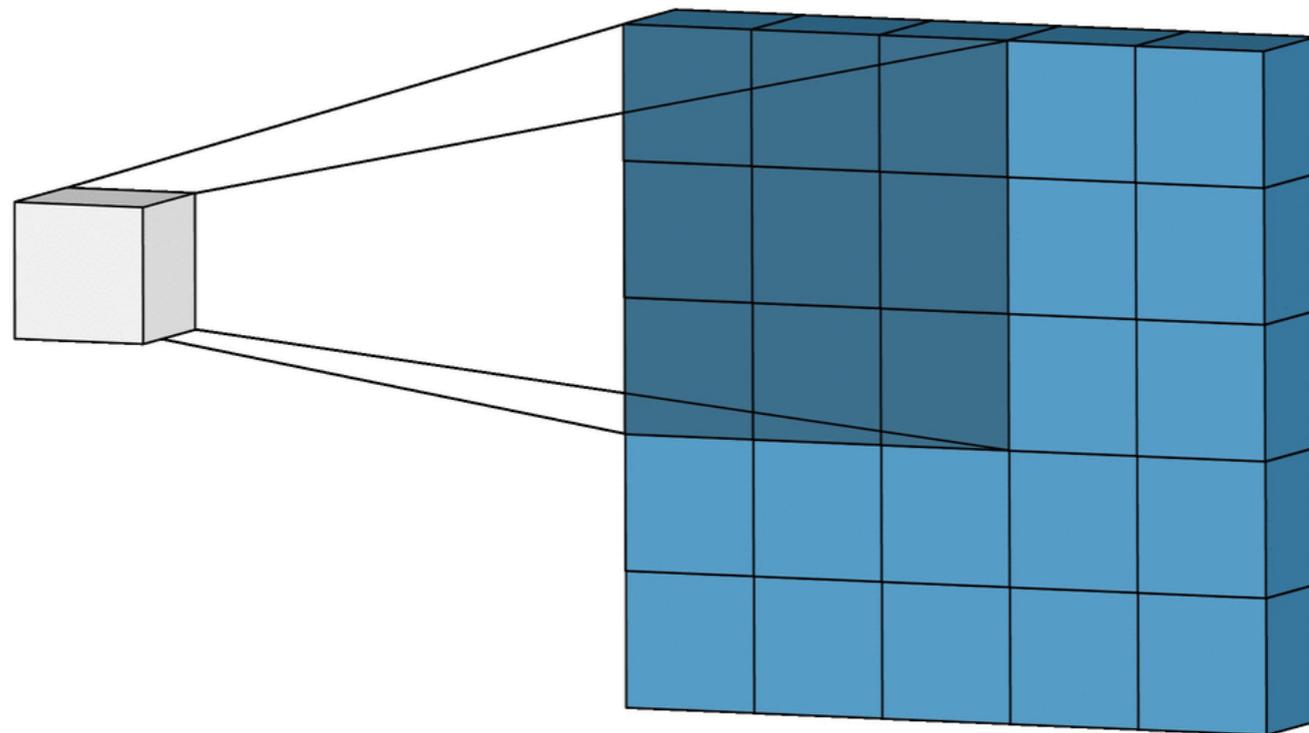


# Convolutional neural network (CNN)

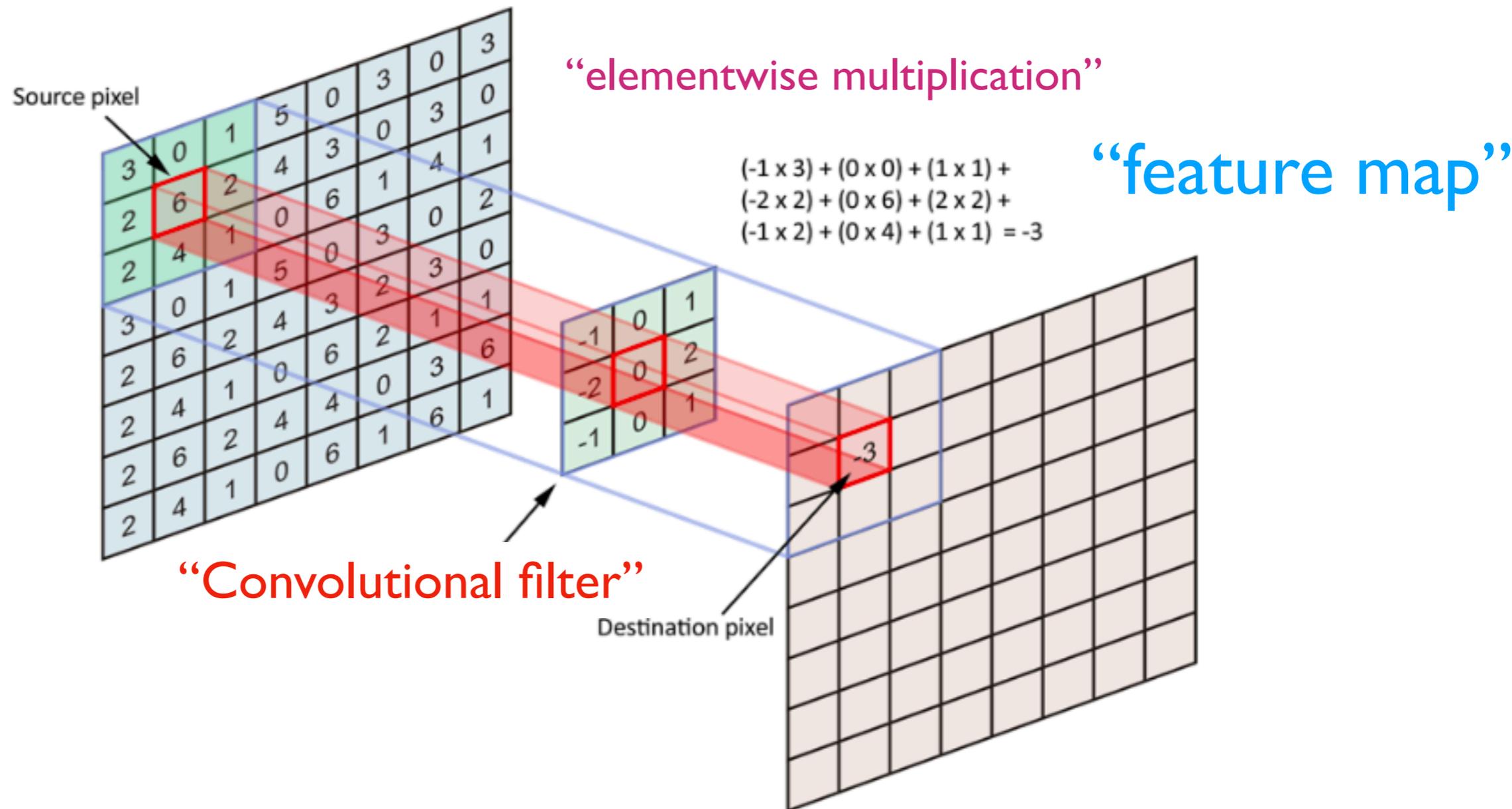
Main idea: features in an image (edges, curves, corners,...eyes, noses,...) are the same no matter where they occur.

Goal: Want to find these features in a translationally invariant way.

Solution: Drag or convolve a “filter” across the image that selects out interesting features.



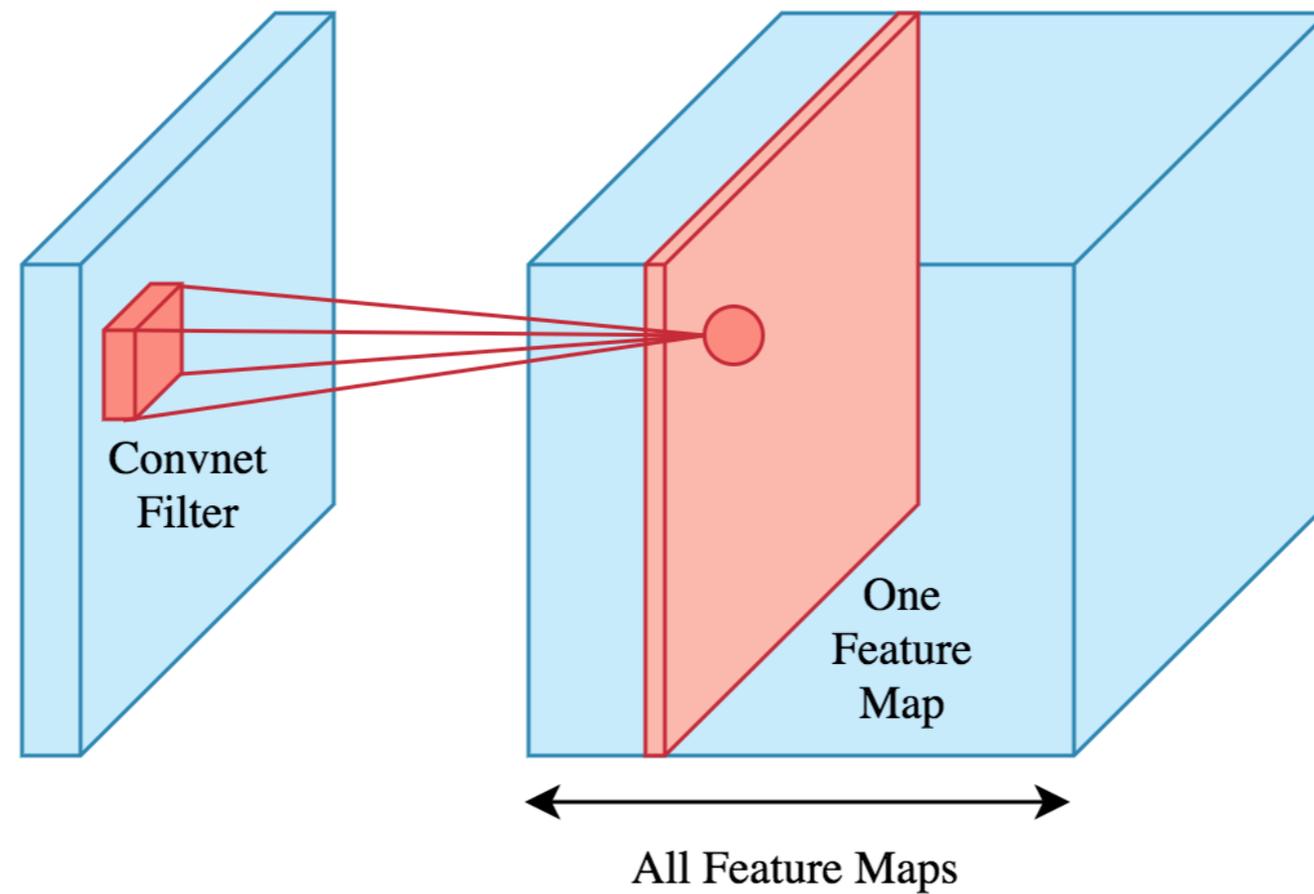
# Convolutional neural network (CNN)



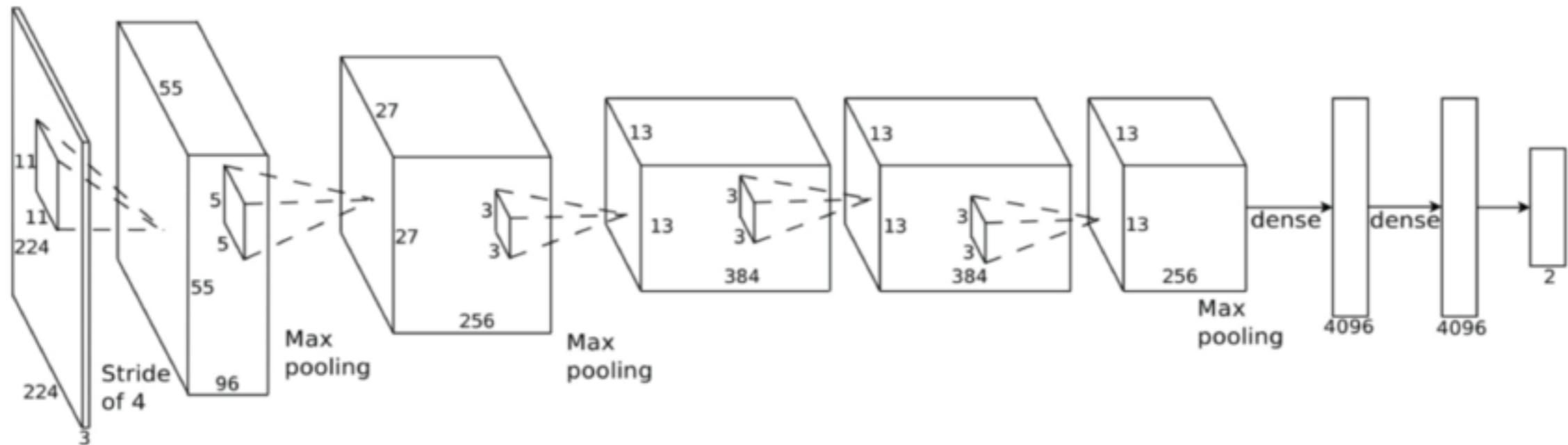
Finds features in the image in a translation invariant way

# Convolutional neural network (CNN)

Can apply multiple filters to image to produce a stack of feature maps



# Convolutional neural network (CNN)

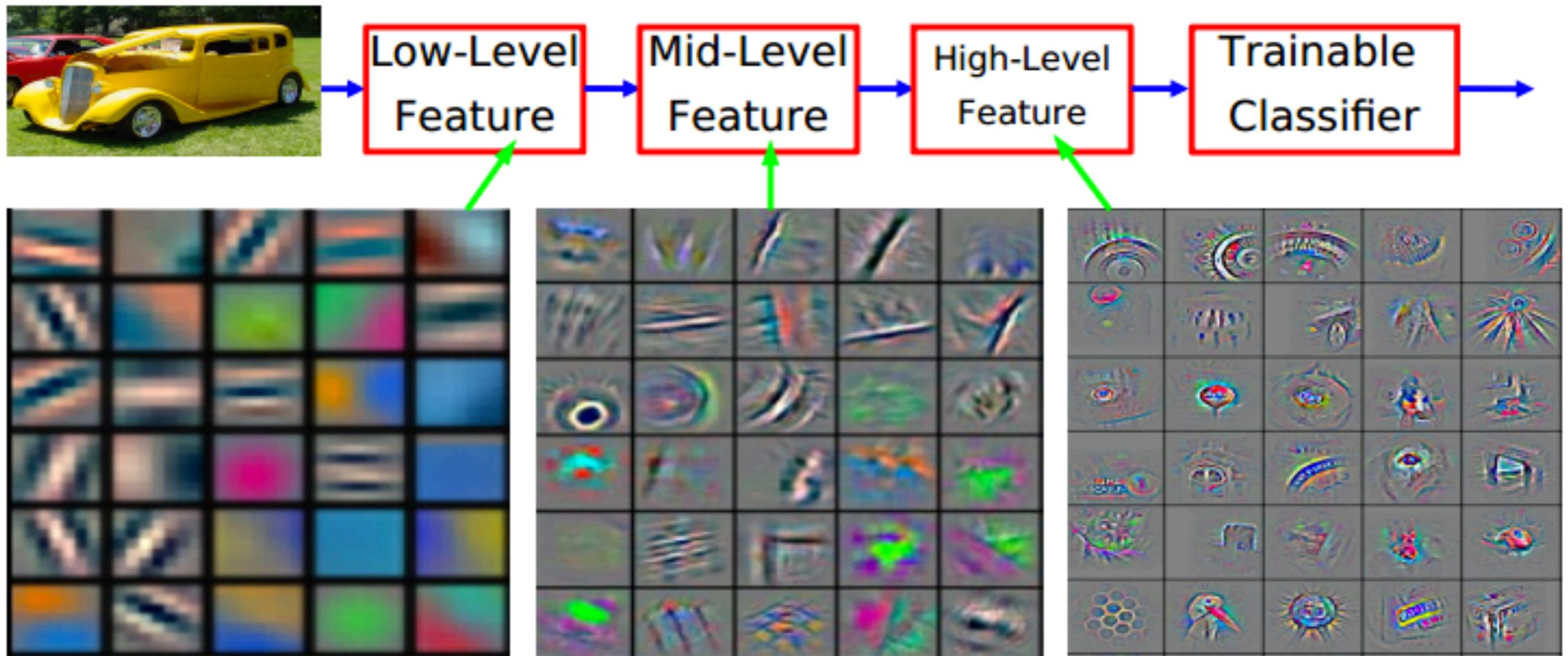


CNNs typically end with fully connected layers.

These are thought to extract the highest level information and to perform the actual classification task on the feature maps found by the convolutional layers.

# Convolutional neural network (CNN)

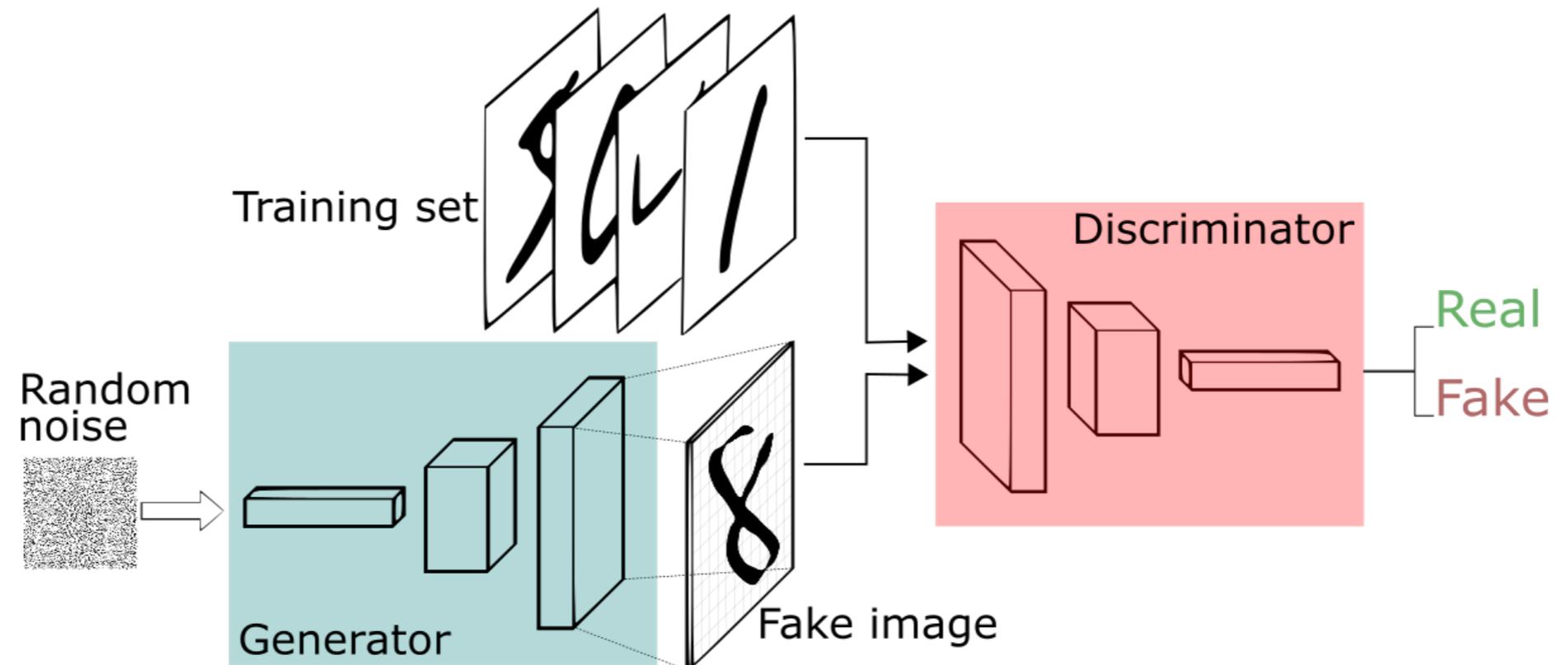
What does the machine learn?



Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

# Generative Adversarial Networks (GANs)

Breakthrough method in generative modeling and unsupervised learning  
(Goodfellow et al. 2014)



Idea: train two neural networks: a “generator” that attempts to generate fake, random images, and a “discriminator” that tries to tell them apart from a database of real images.

# Generative Adversarial Networks (GANs)

$$L_{GAN} = \sum_{x \in \text{real}} \log D(x) + \sum_{z \in \text{random}} \log(1 - D(G(z)))$$

Training is performed “adversarially”

- Discriminator tries to minimize loss
- Generator tries to **maximize** loss
- Take turns training discriminator and generator to optimize fake image generator

# Generative Adversarial Networks (GANs)



Real or fake?

# 3. Deep Learning at the LHC

First Stable Beams



proton-proton collisions at 13 TeV

Run: 266904  
Event: 9393006  
2015-06-03 10:40:31 CEST

# LHC and Big Data

First Stable Beams



proton-proton collisions at 13 TeV

Run: 266904  
Event: 9393006  
2015-06-03 10:40:31 CEST

# LHC and Big Data

- 600 million collisions per second

First Stable Beams



proton-proton collisions at 13 TeV

Run: 266904  
Event: 9393006  
2015-06-03 10:40:31 CEST

# LHC and Big Data

- 600 million collisions per second
- Raw data rate  $\sim 1$  PB/s ( $1 \text{ PB} = 10^6 \text{ GB}$ )

First Stable Beams



proton-proton collisions at 13 TeV

Run: 266904  
Event: 9393006  
2015-06-03 10:40:31 CEST

# LHC and Big Data

- 600 million collisions per second
- Raw data rate  $\sim 1$  PB/s ( $1 \text{ PB} = 10^6 \text{ GB}$ )
- Actual data rate  $\sim 25 \text{ GB/s}$

First Stable Beams



proton-proton collisions at 13 TeV

Run: 266904  
Event: 9393006  
2015-06-03 10:40:31 CEST

# LHC and Big Data

- 600 million collisions per second
- Raw data rate  $\sim 1$  PB/s ( $1 \text{ PB} = 10^6 \text{ GB}$ )
- Actual data rate  $\sim 25 \text{ GB/s}$ 
  - Need to trigger on 1 out of 40,000 events

First Stable Beams



proton-proton collisions at 13 TeV

Run: 266904  
Event: 9393006  
2015-06-03 10:40:31 CEST

# LHC and Big Data

- 600 million collisions per second
- Raw data rate  $\sim 1$  PB/s ( $1$  PB= $10^6$  GB)
- Actual data rate  $\sim 25$  GB/s
  - Need to trigger on 1 out of 40,000 events
- $\sim 10$ 's of PB annually

First Stable Beams



proton-proton collisions at 13 TeV

Run: 266904  
Event: 9393006  
2015-06-03 10:40:31 CEST

# LHC and Big Data

The data is

- large (billions of events on tape)
- complex (hundreds of particles per event)
- well-understood (Standard Model of particle physics).

Also, it is relatively easy to generate realistic simulated data.

- 600 million collisions per second
- Raw data rate  $\sim 1$  PB/s ( $1 \text{ PB} = 10^6 \text{ GB}$ )
- Actual data rate  $\sim 25 \text{ GB/s}$ 
  - Need to trigger on 1 out of 40,000 events
- $\sim 10$ 's of PB annually

First Stable Beams



proton-proton collisions at 13 TeV

Run: 266904  
Event: 9393006  
2015-06-03 10:40:31 CEST

# LHC and Big Data

The data is

- large (billions of events on tape)
- complex (hundreds of particles per event)
- well-understood (Standard Model of particle physics).

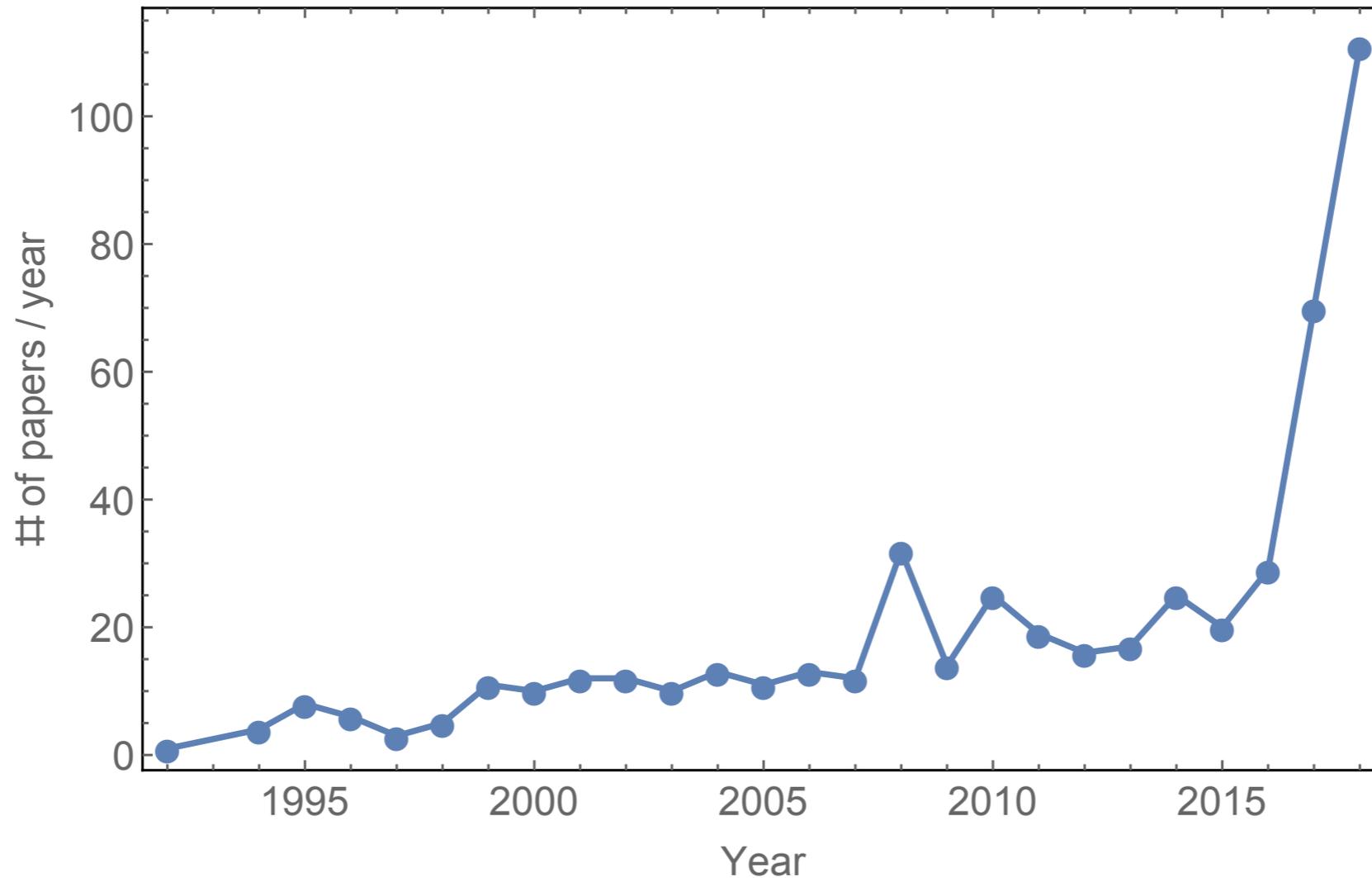
Also, it is relatively easy to generate realistic simulated data.

- 600 million collisions per second
- Raw data rate  $\sim 1$  PB/s ( $1$  PB= $10^6$  GB)
- Actual data rate  $\sim 25$  GB/s
  - Need to trigger on 1 out of 40,000 events
- $\sim 10$ 's of PB annually

***The LHC is a great setting for deep learning!***

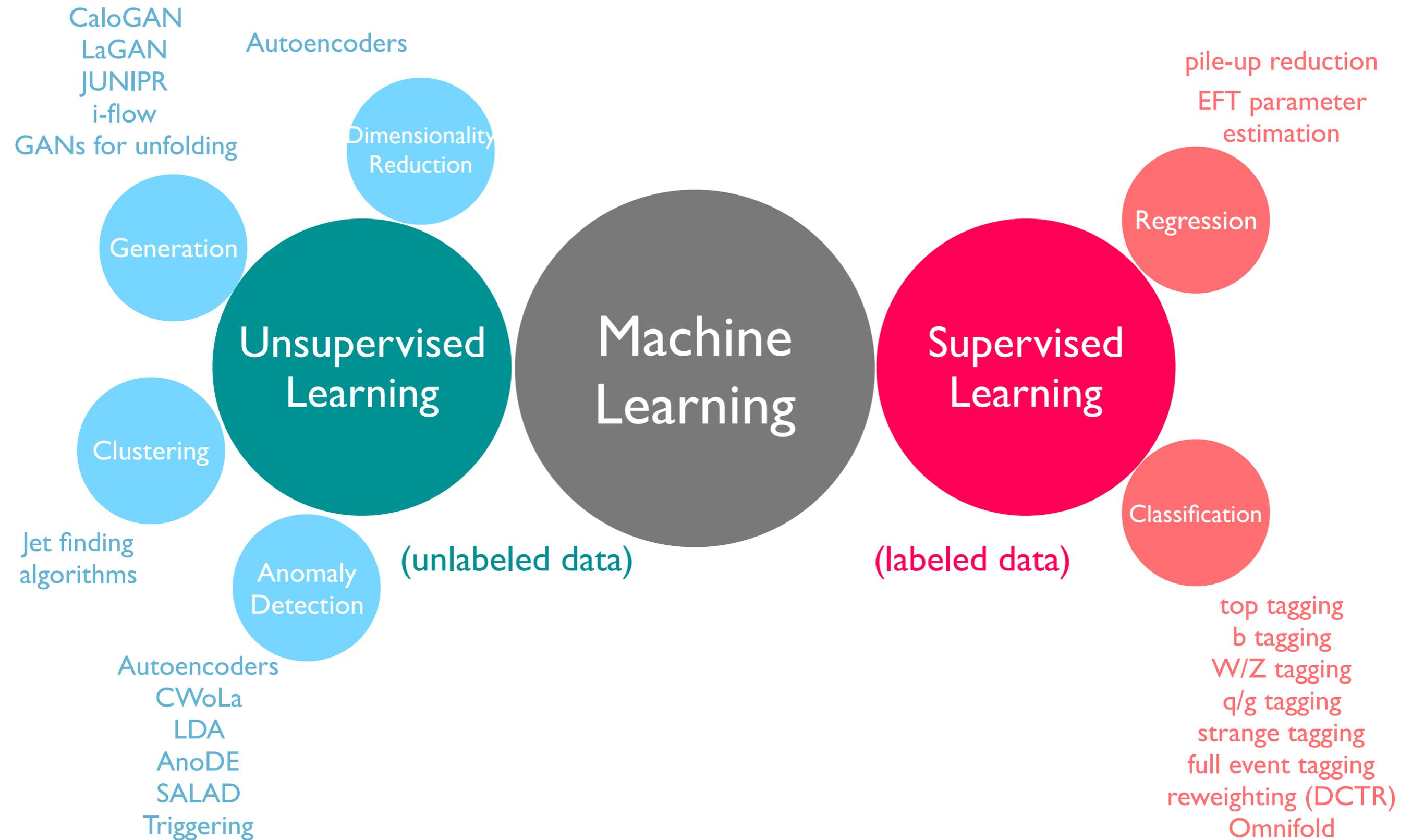
# Deep Learning Papers

INSPIRE search: ("machine learning" or "deep learning" or neural)  
and (hep-ex or hep-ph)

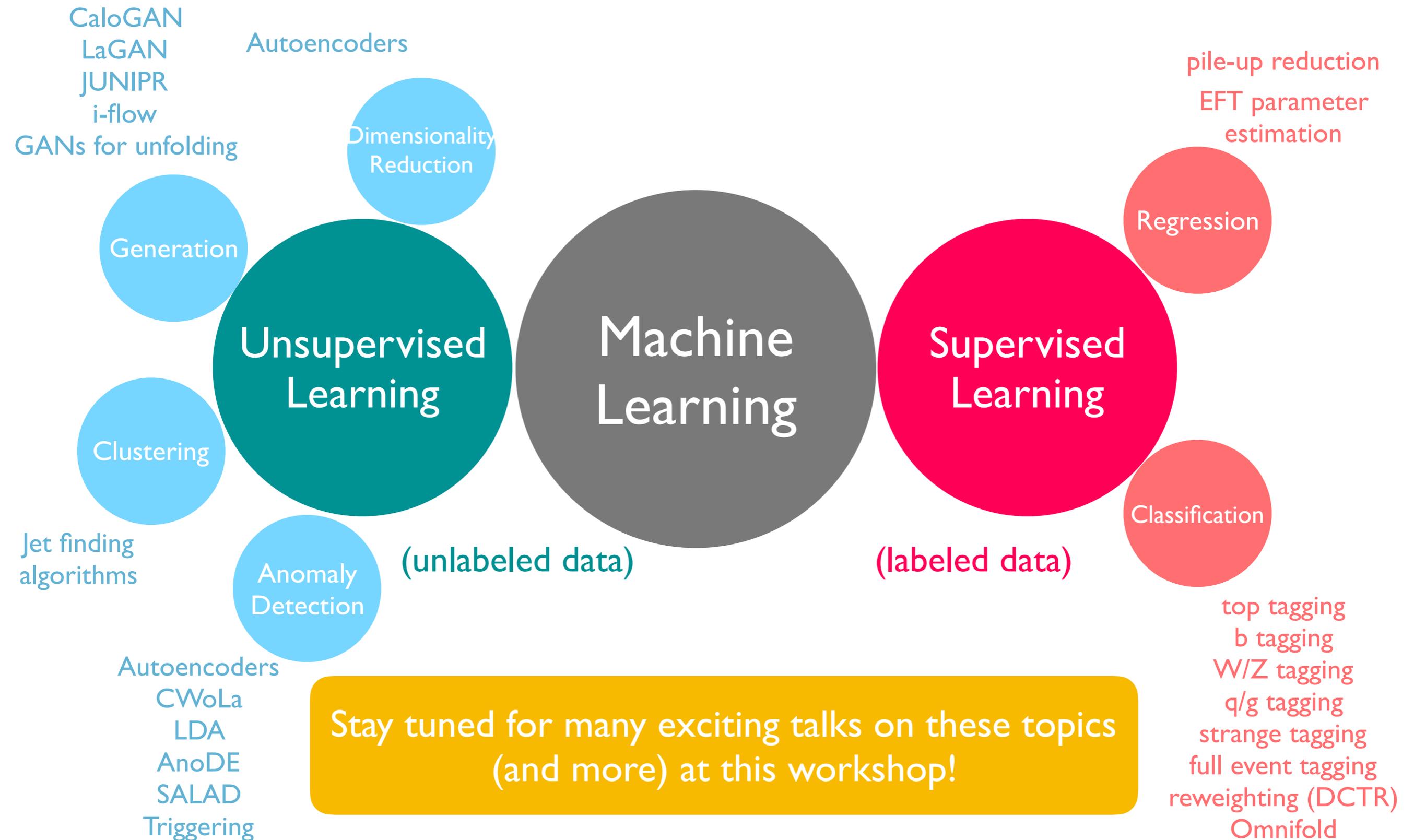


**An explosion of interest in machine learning!**

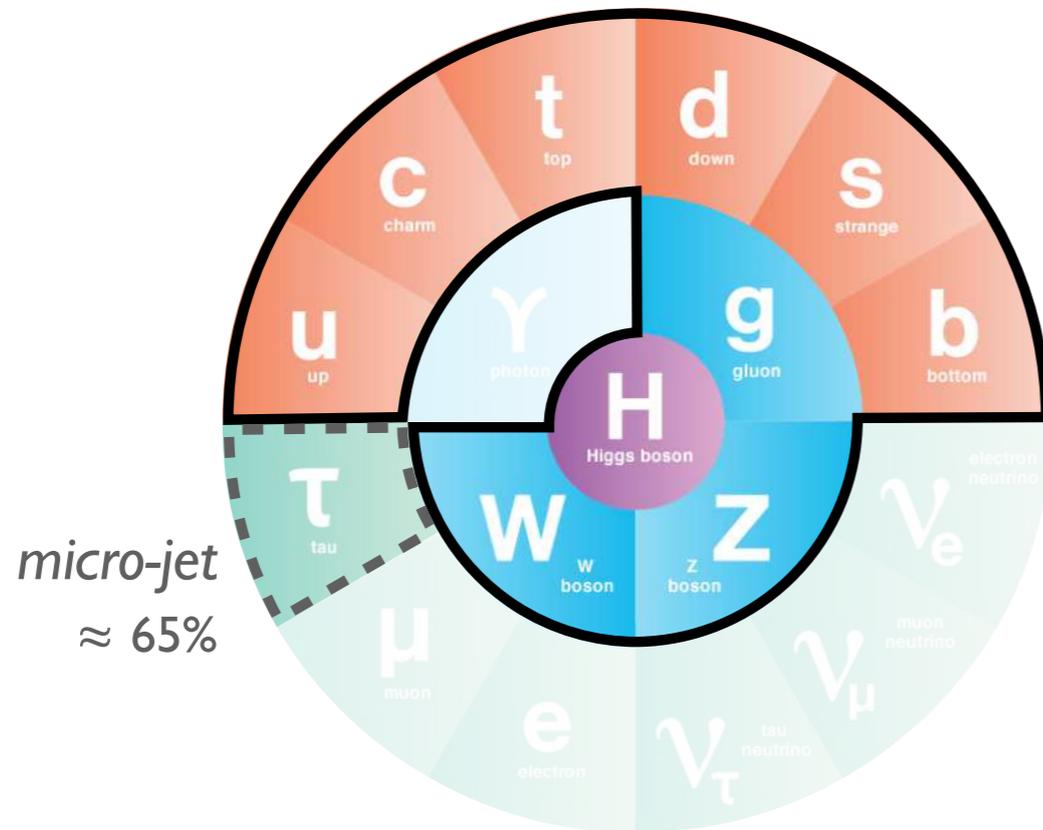
# The Landscape of DL @ LHC



# The Landscape of DL @ LHC

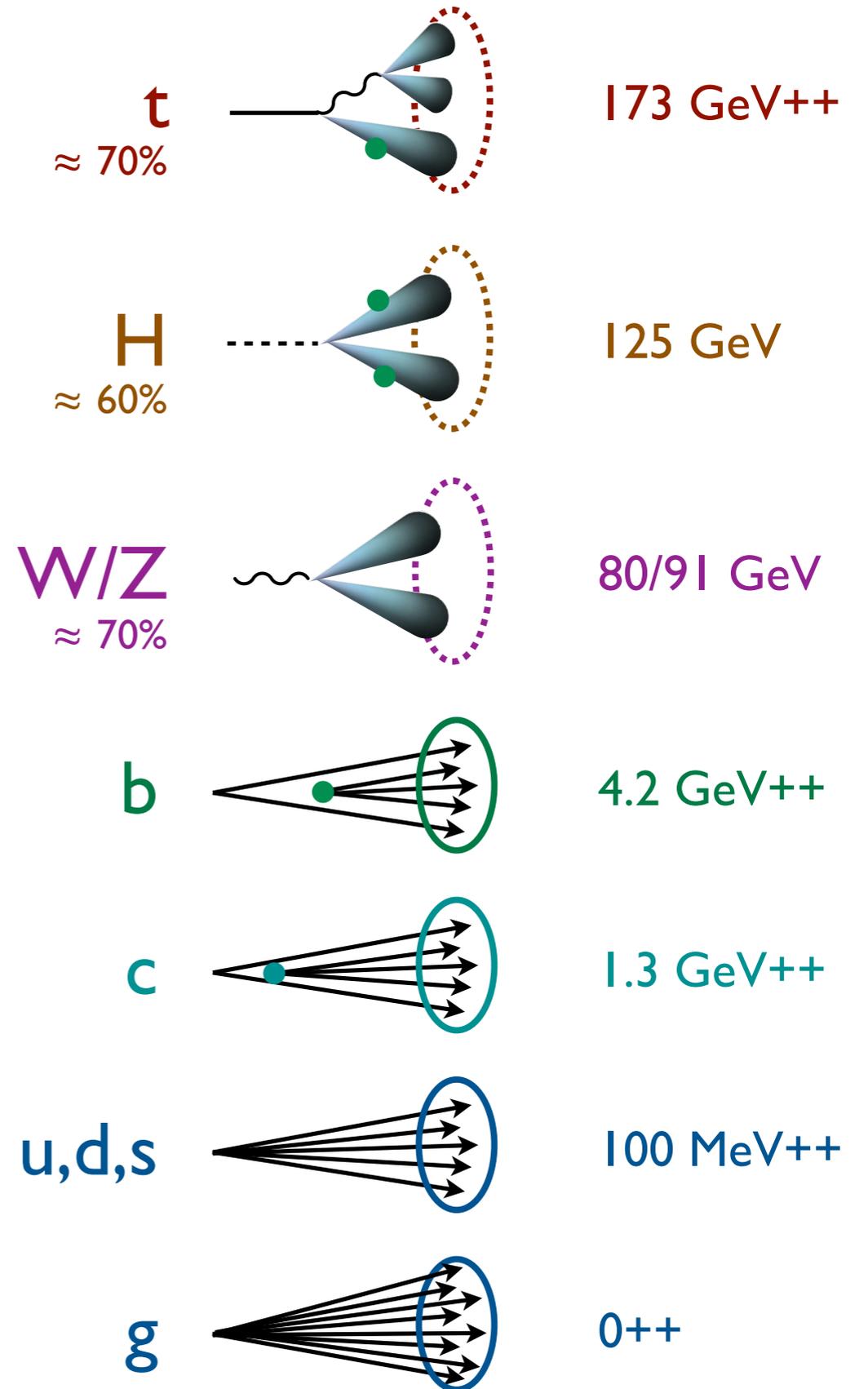


# Jet classification (“tagging”)

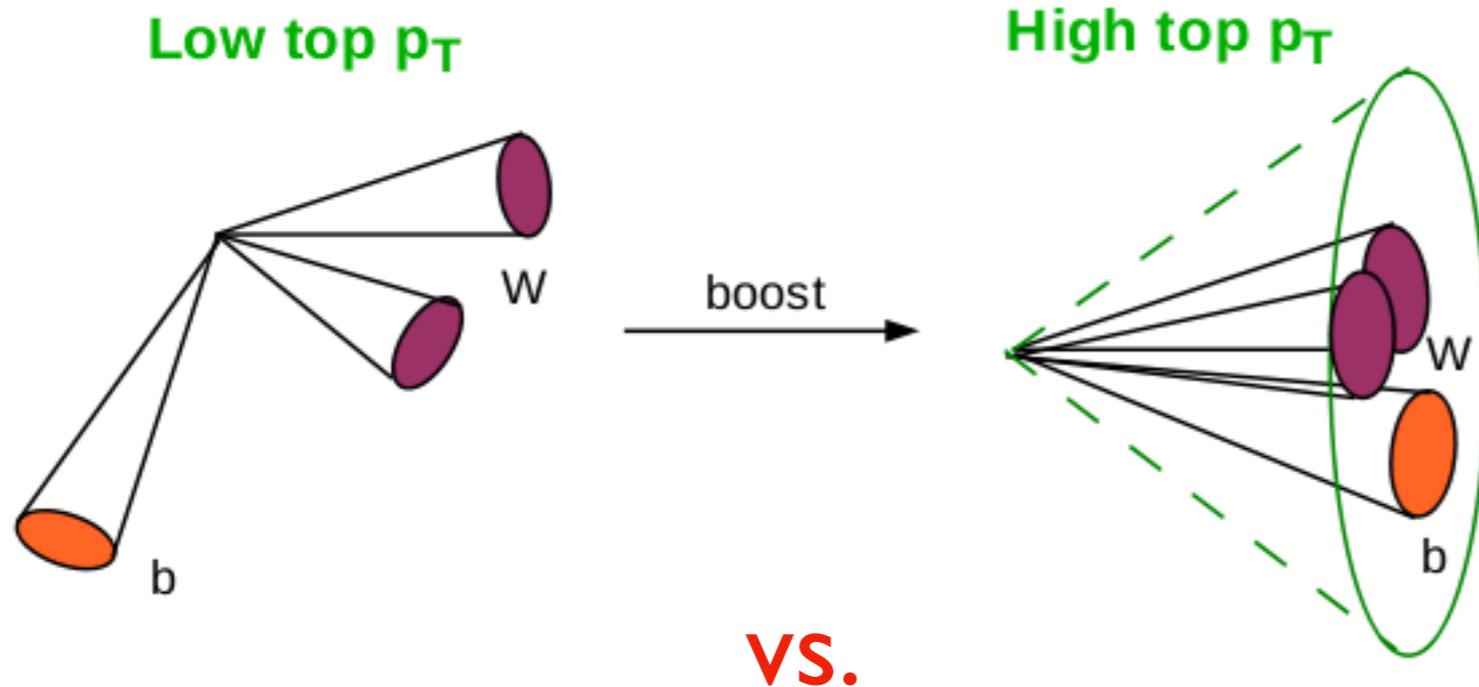


*Jets from QCD and the Standard Model*

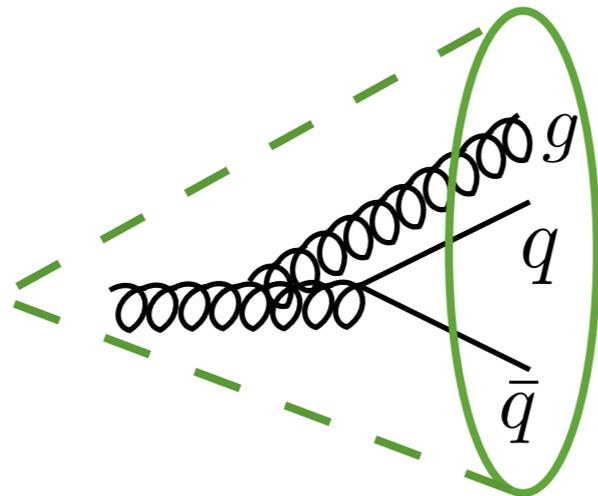
++ = Mass from QCD Radiation



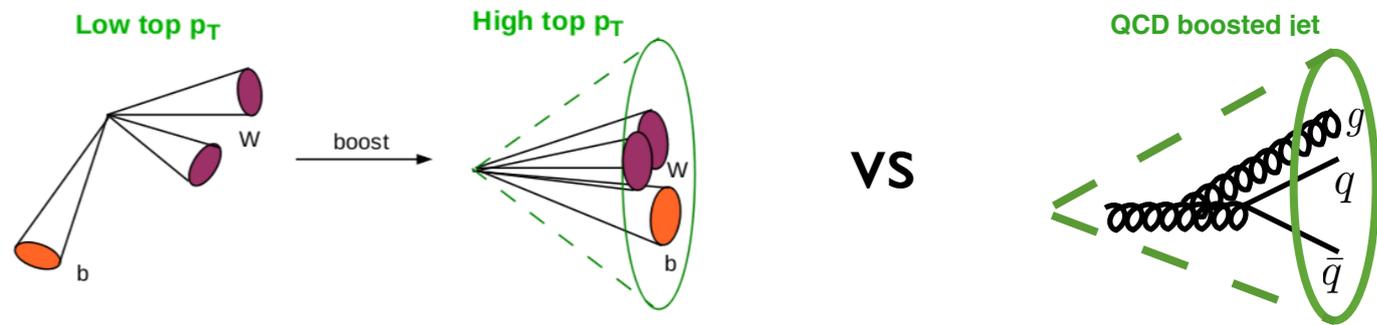
# A popular example: boosted top tagging



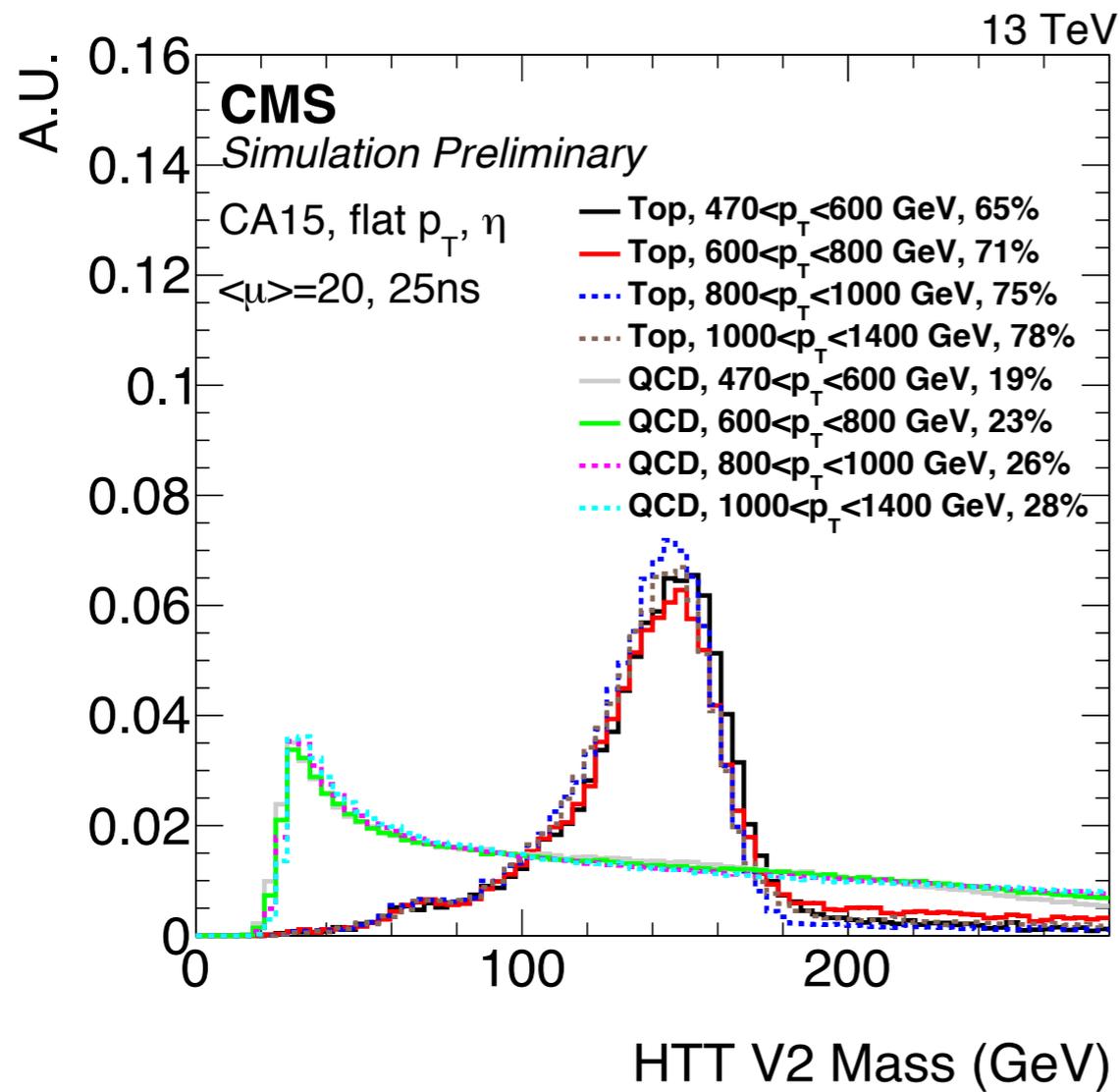
QCD boosted jet



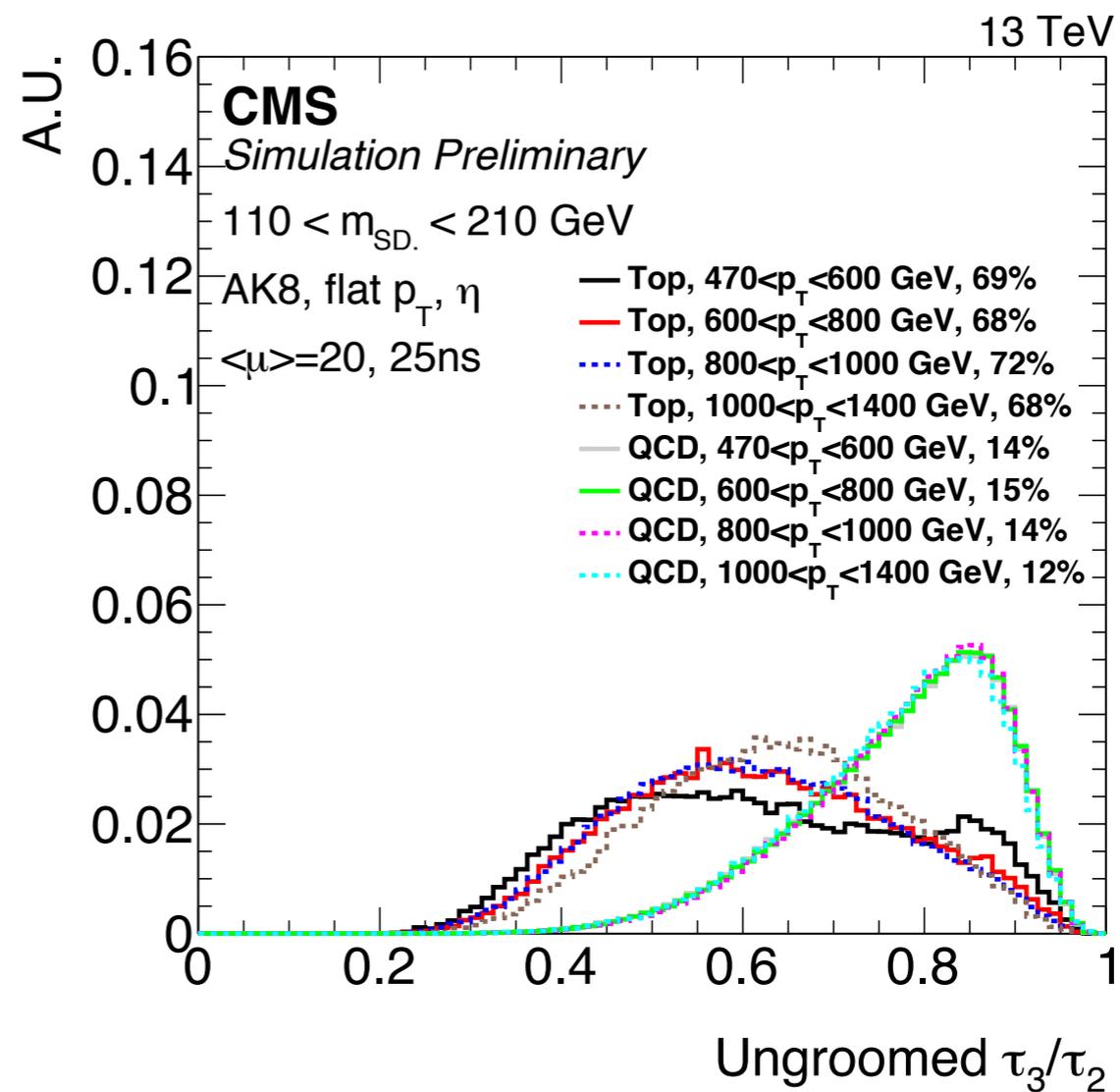
How to differentiate between these two types of jets?



Some obvious ideas:



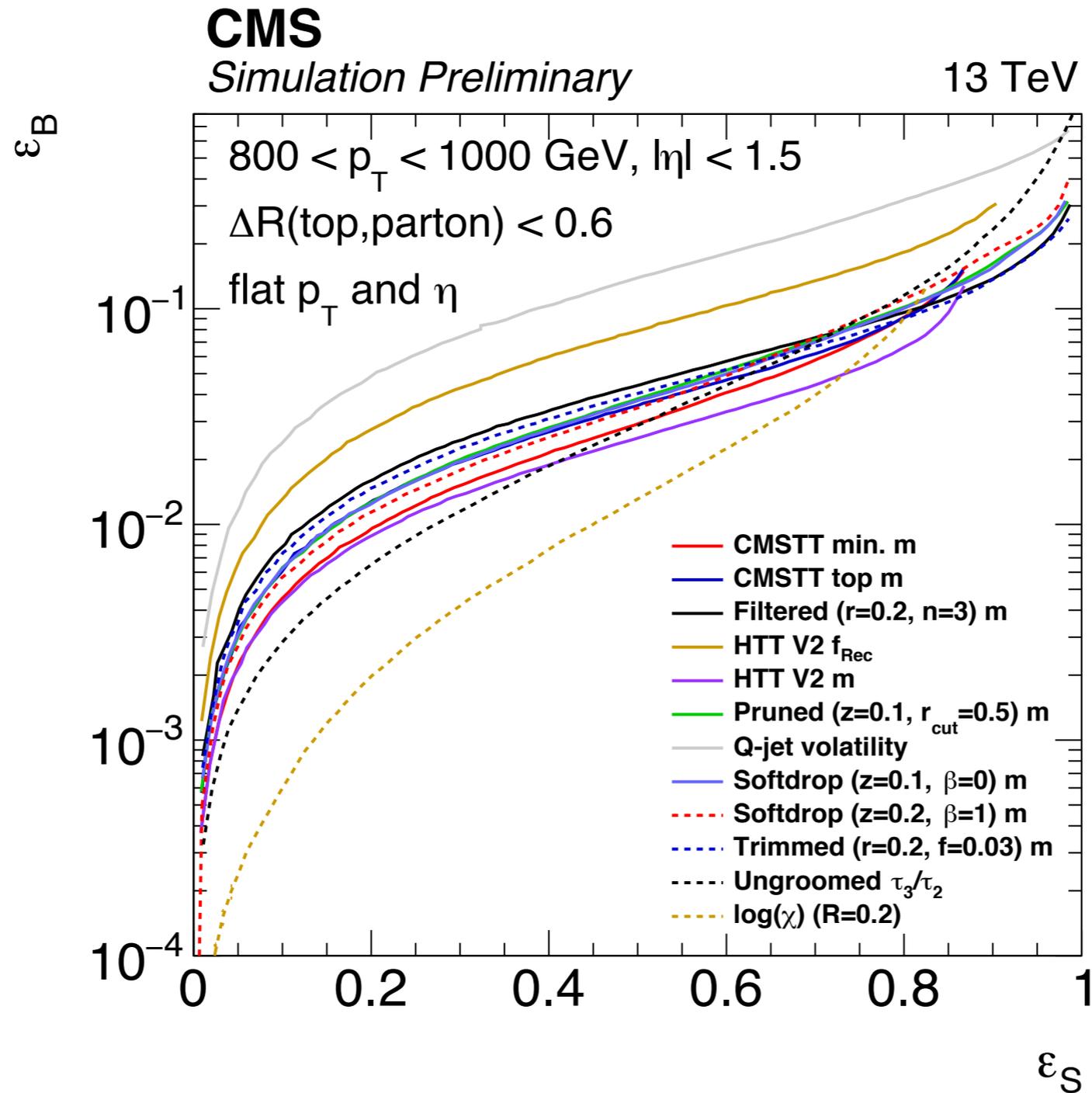
jet mass ( $m_{\text{top}}$  vs 0)



jet substructure (3 vs 1)

# State of the art with cuts on kinematic quantities:

QCD jet  
mistag rate



“ROC curve”

top tagging efficiency

Can deep learning do better??

# Jet images

Cogan et al 1407.5675, Almeida et al 1501.05968, de Oliveira et al 1511.05190

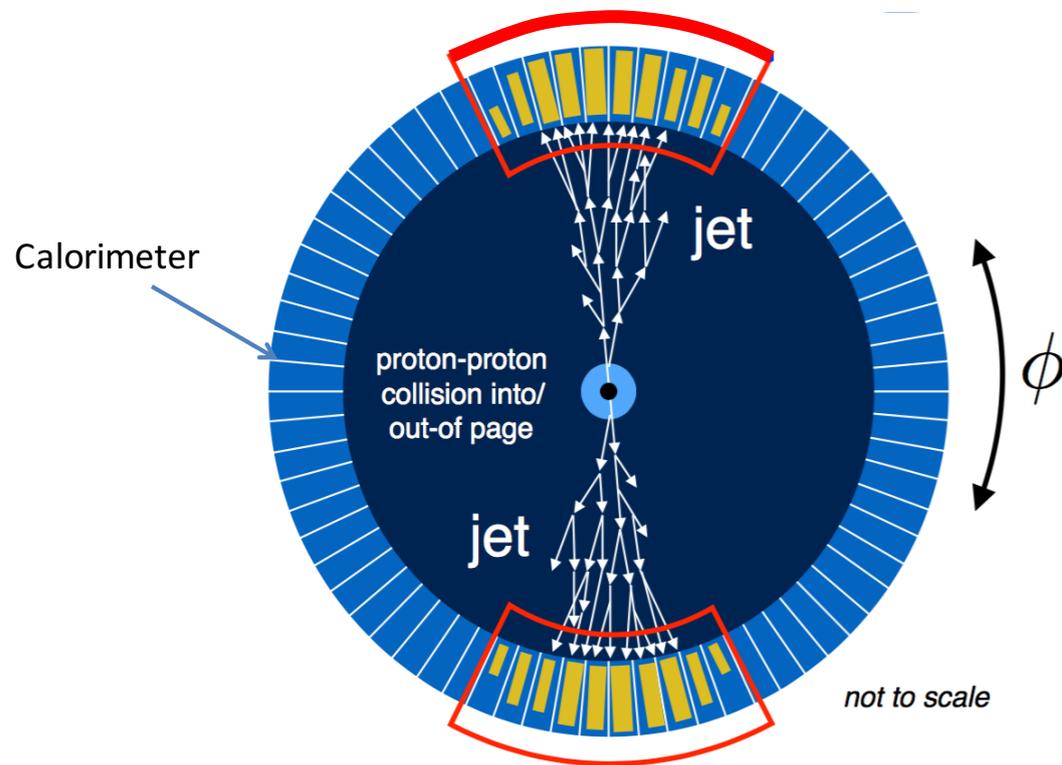
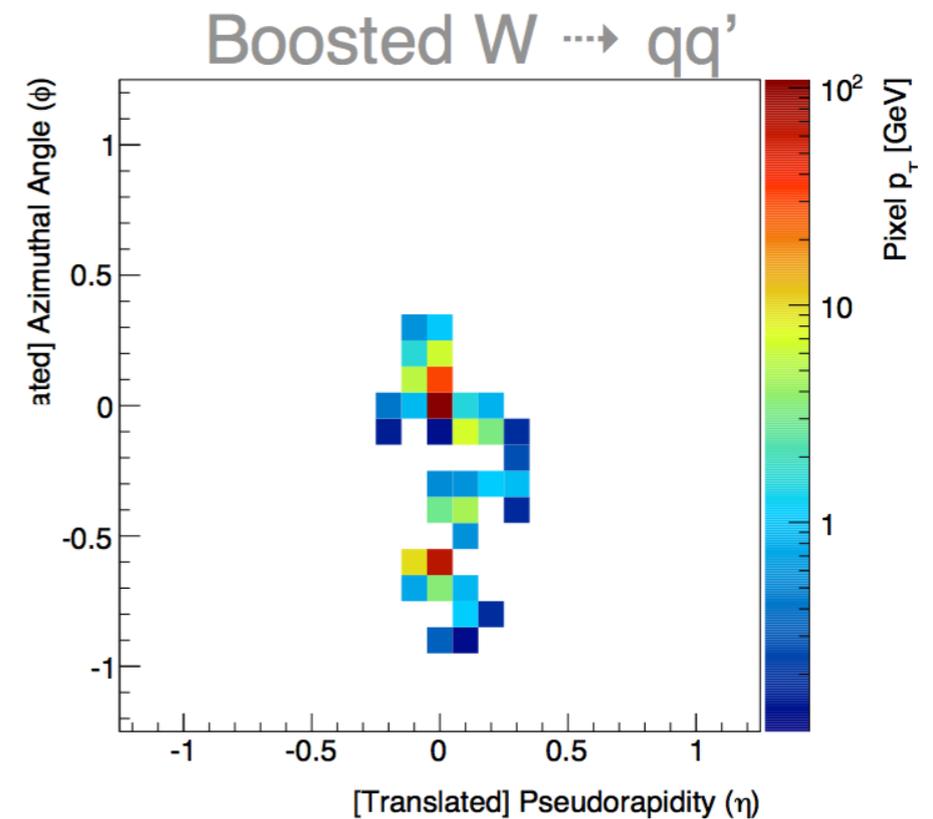


Figure credit:  
B. Nachman



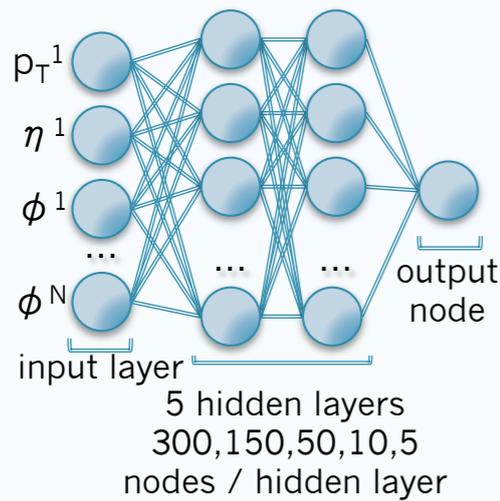
Jets are naturally images in eta and phi.

Should be able to apply “off-the-shelf” NNs developed for image recognition to classify jets at the LHC! de Oliveira et al 1511.05190

# Other jet representations

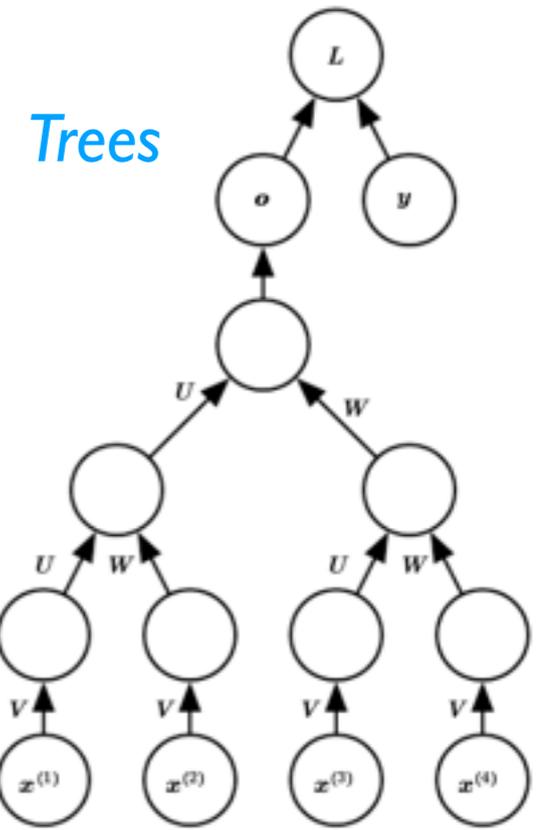
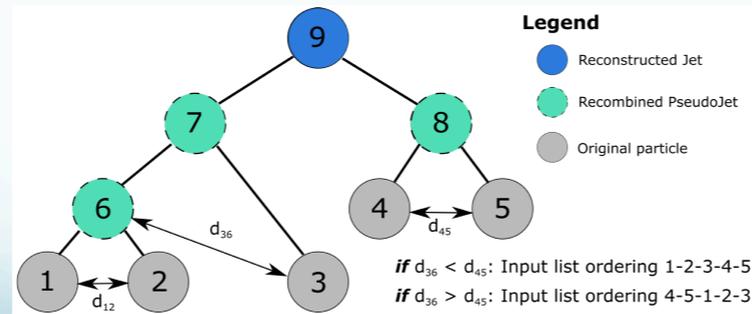
Many other ways to represent a jet besides jet images!

## Lists of 4-vectors



## Sequences

- View anti- $k_T$  sequence as a binary tree
- Order using depth-first traversal prioritizing jets with 'parents' whose  $d_{ij}$  is smaller

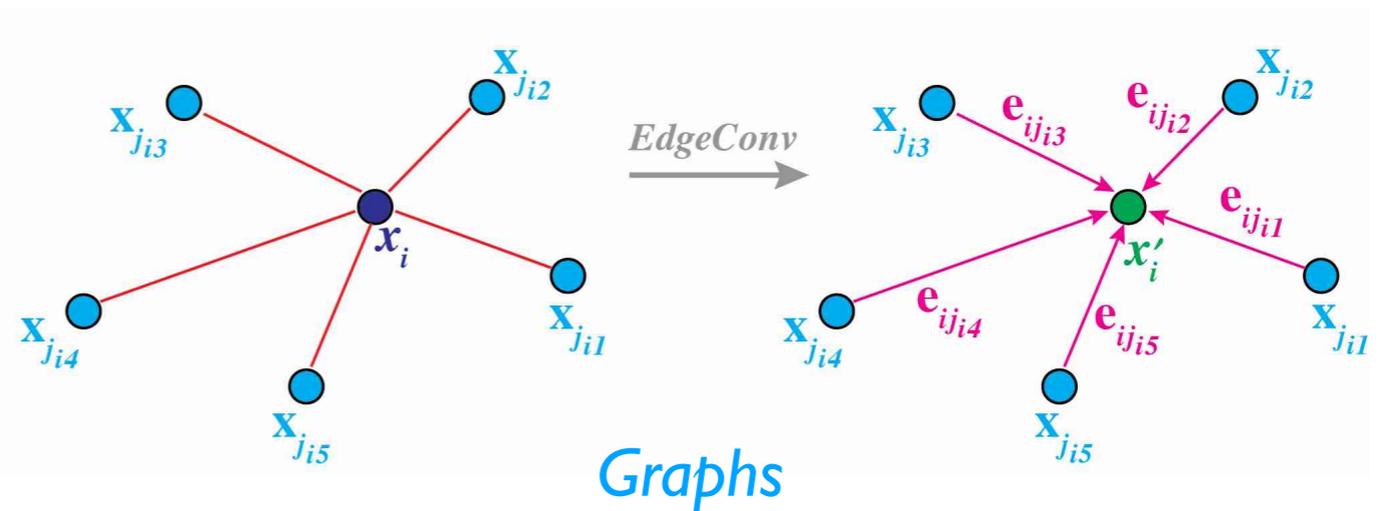
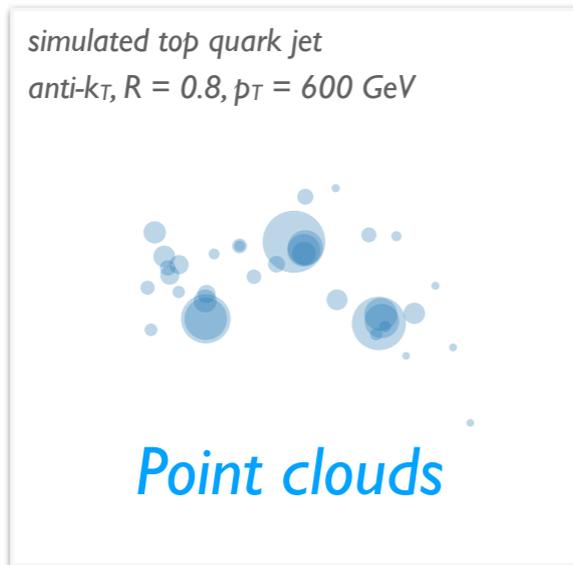


$$j \iff \sum_{i_j=1}^M z_{ij}$$

$$k \text{ --- } l \iff \theta_{ikil}$$

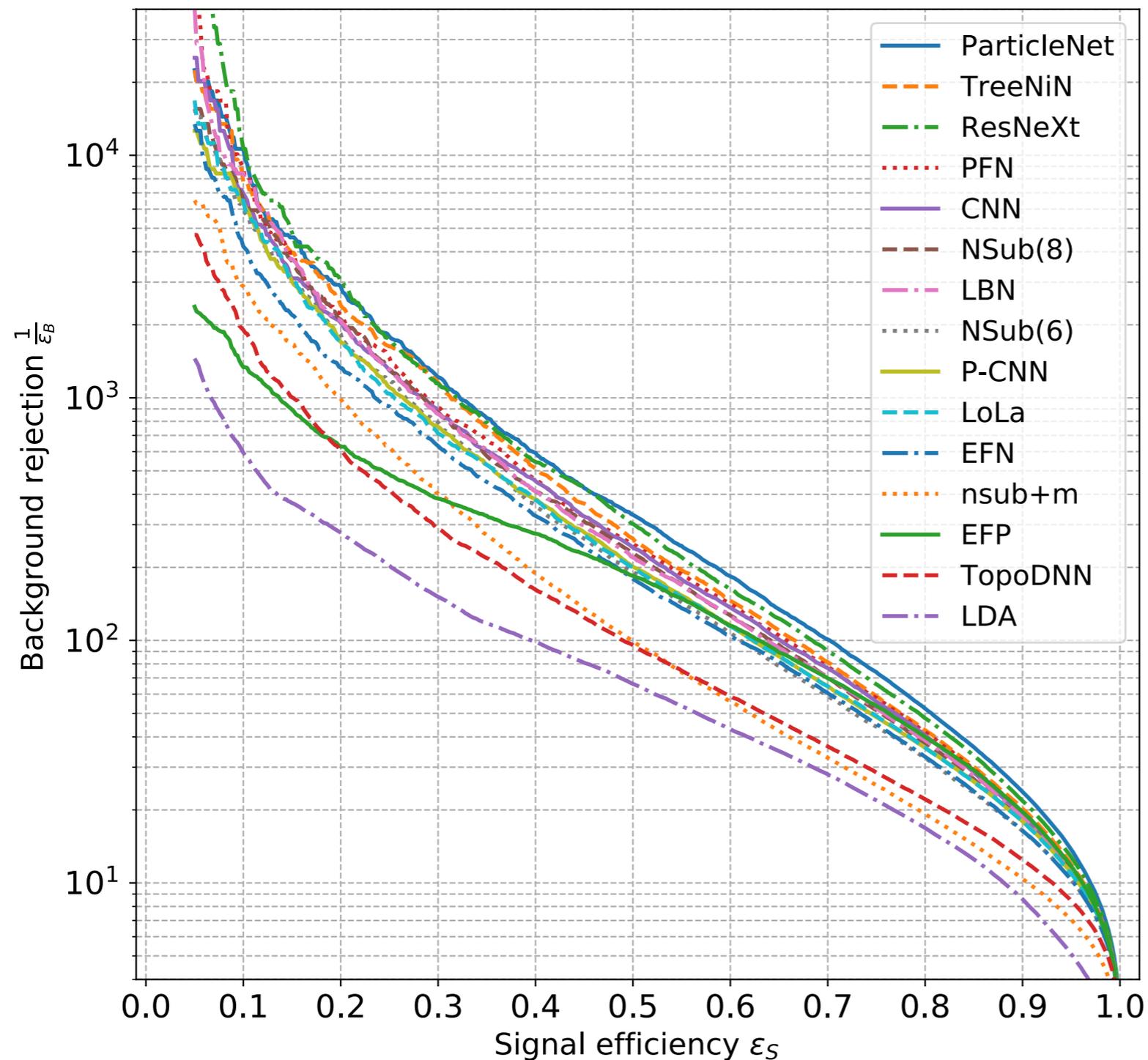
**Kinematic invariants**

$$= \sum_{i_1=1}^M \sum_{i_2=1}^M \sum_{i_3=1}^M \sum_{i_4=1}^M z_{i_1} z_{i_2} z_{i_3} z_{i_4} \theta_{i_1 i_2} \theta_{i_2 i_3} \theta_{i_2 i_4}^2 \theta_{i_3 i_4}$$



# Community top tagging comparison

Kasieczka, Plehn et al 1902.09914



**Q:** There are many papers developing jet taggers with different jet representations and architectures. How can we evaluate their relative strengths and weaknesses?

**A:** Let's perform an apples-to-apples comparison of various top taggers on a common dataset!

# Community top tagging comparison

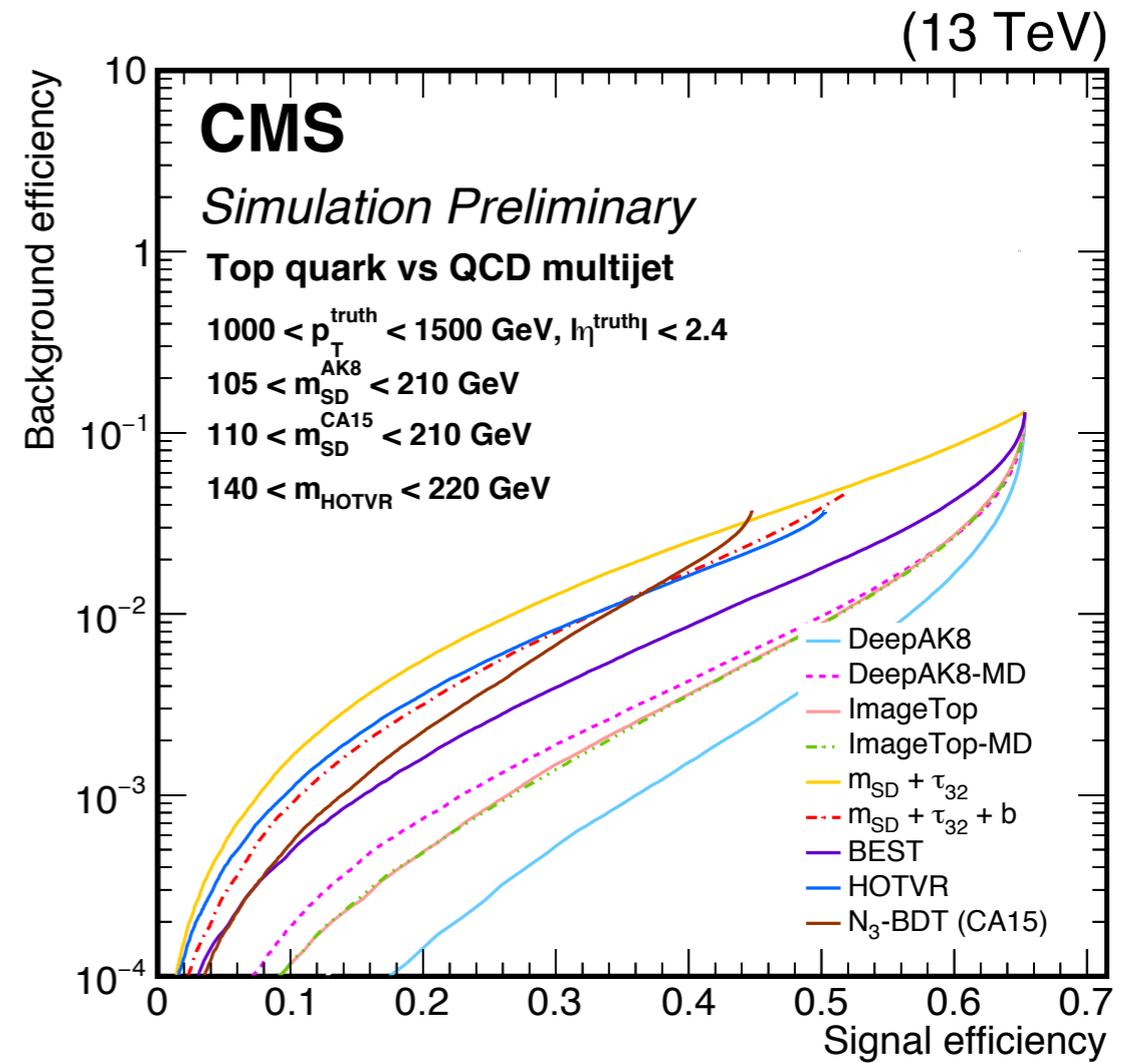
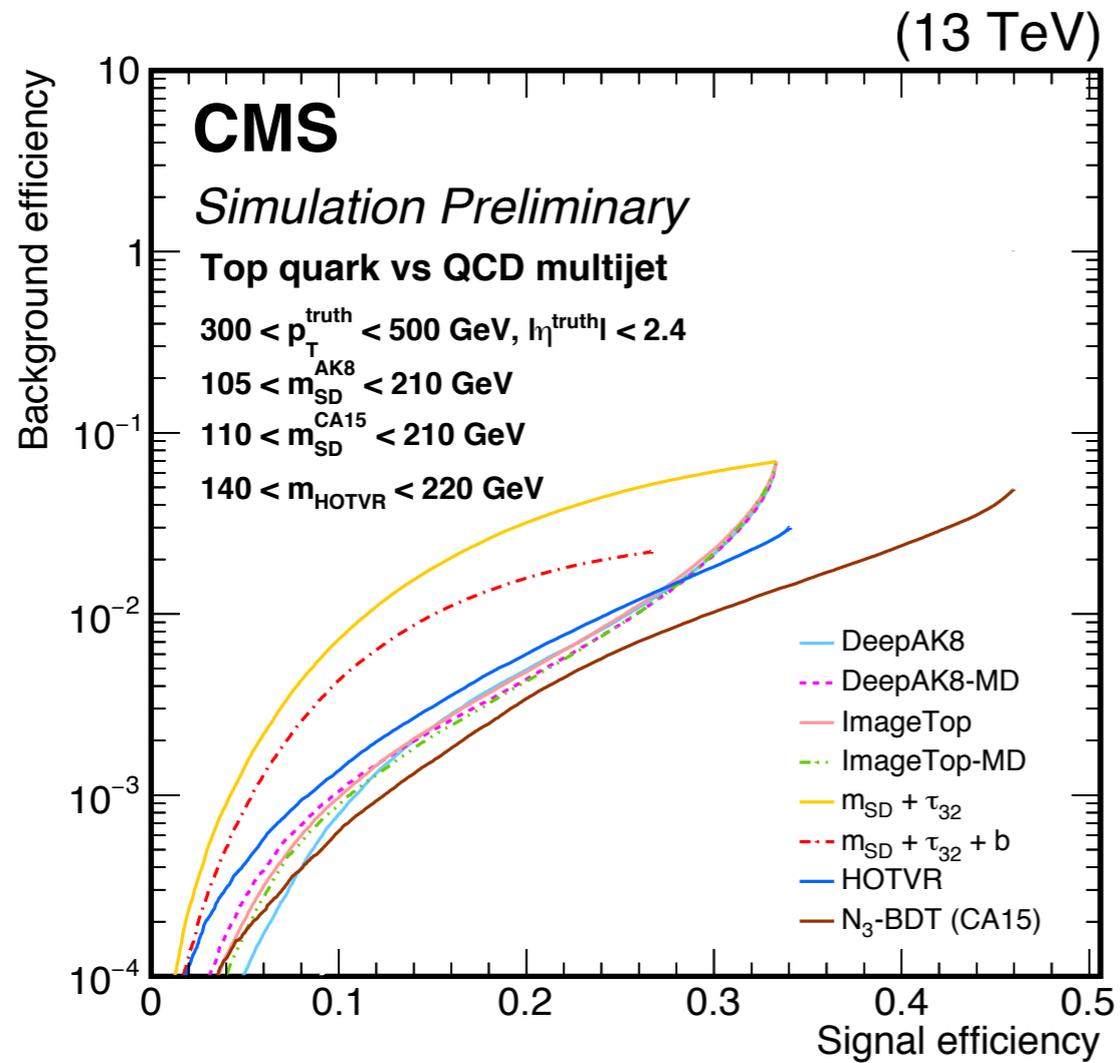
Kasieczka, Plehn et al 1902.09914

	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 [31]	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	
LDA [54]	0.955	0.892	151±0.4	151.5±0.5	151.7±0.4	184k	
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	←

Have we found the optimal tagger??

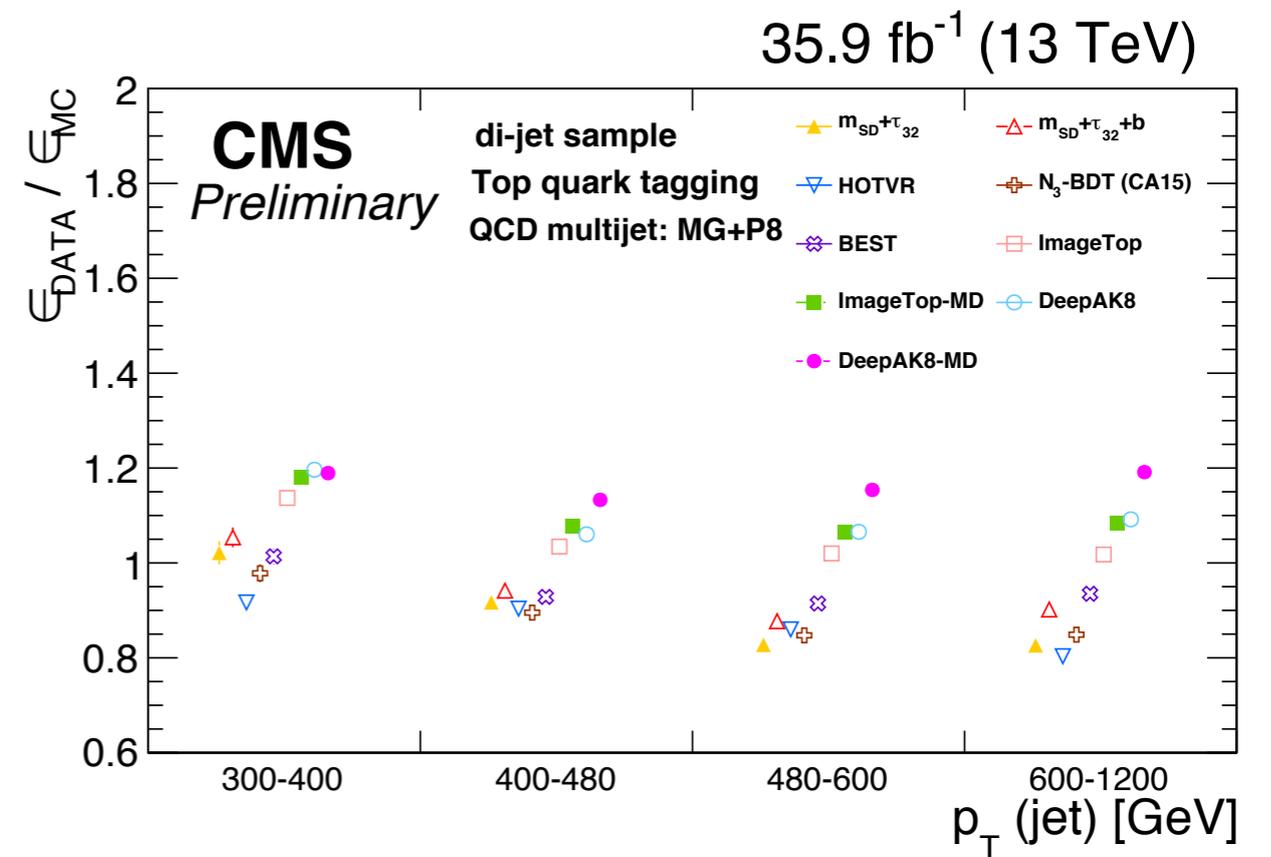
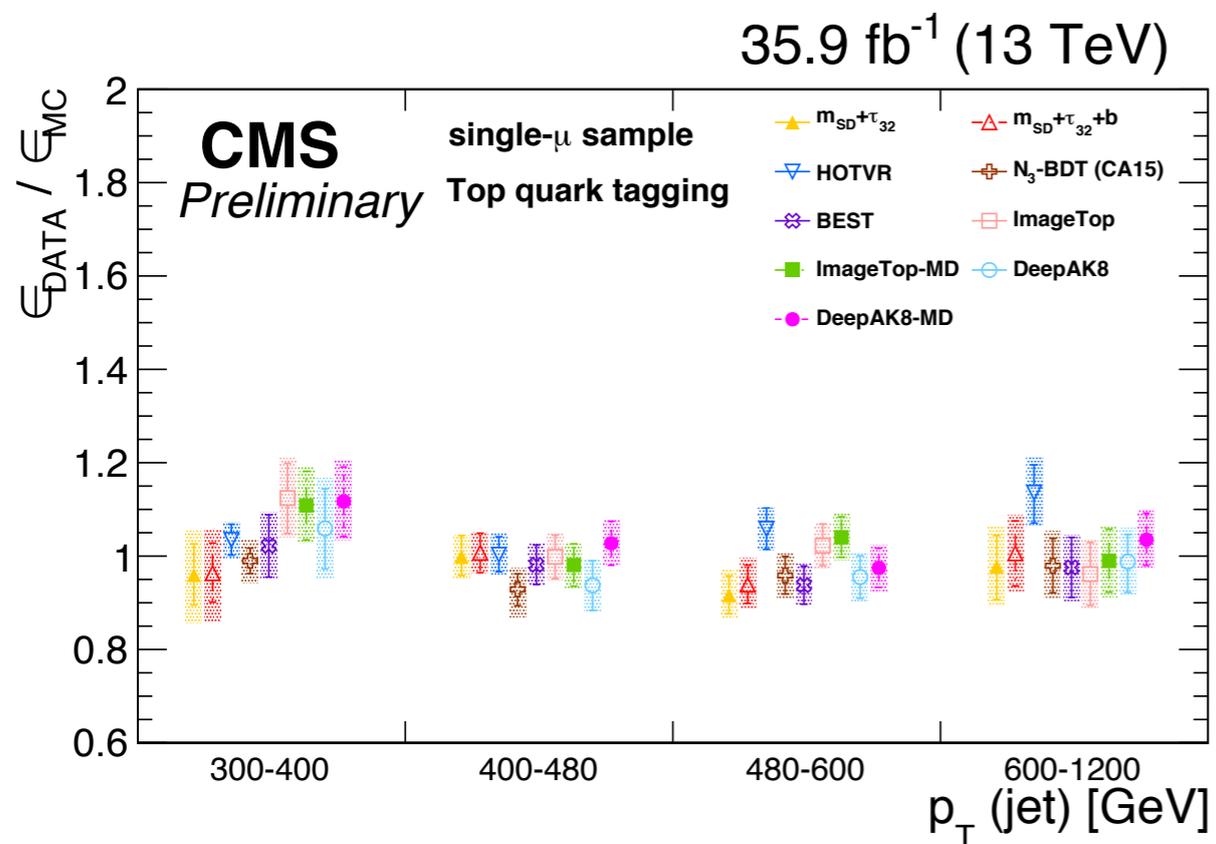
# CMS performance study

JME-18-002-PAS



# CMS performance study

JME-18-002-PAS



Performance measured in data!

# Beyond top tagging

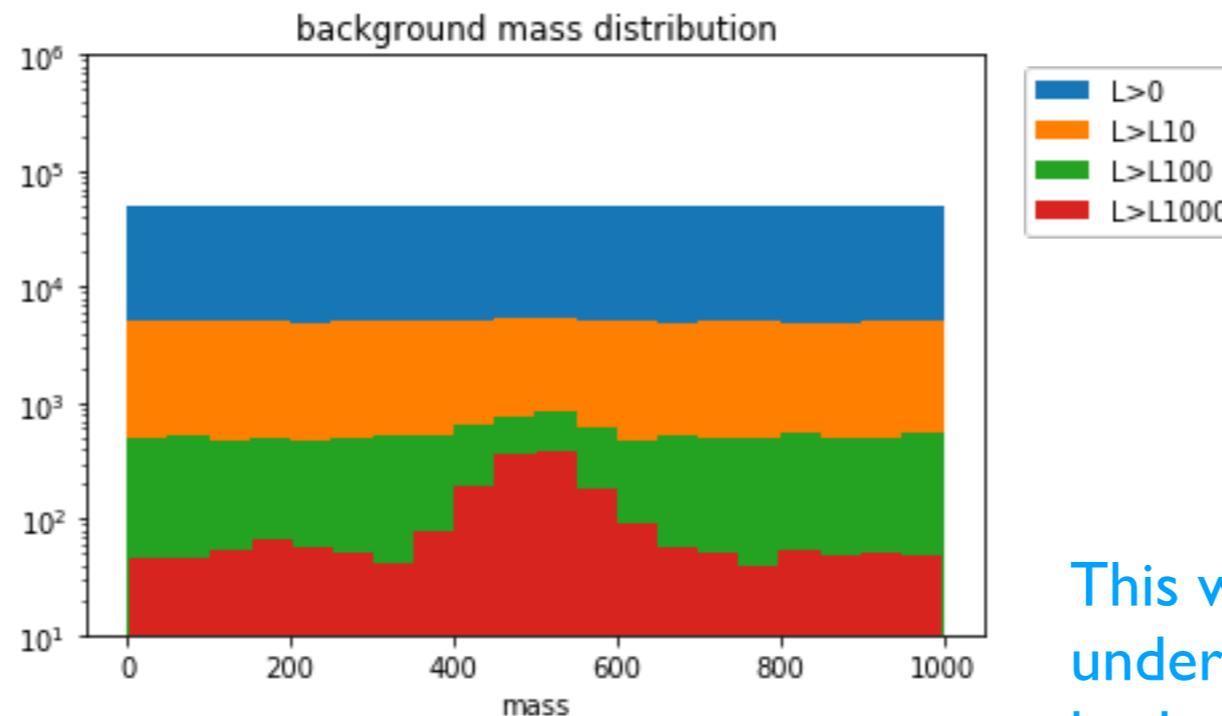
Many other classification tasks at LHC are enhanced with deep learning:

- quark/gluon tagging (Komiske, Metodiev & Schwartz '16, et seq)
- b and c tagging (DeepCSV)
- boosted W/Z tagging (Oliveira et al '15, et seq)
- strange tagging (Nakai, Thomas & DS to appear)
- u vs d tagging using jet charge (Fraser & Schwartz '18)
- boosted W<sup>+</sup>/Z/W<sup>-</sup> tagging using jet charge (Chen, Chiang, Cottin & DS 1908.08256)
- ....

# Beyond classification — decorrelation

Raw tagger performance not the only consideration.

For robust background estimation, generally need to ensure tagger does not sculpt the background mass distribution.



This would greatly underestimate the background in the SR

Two approaches with deep learning:

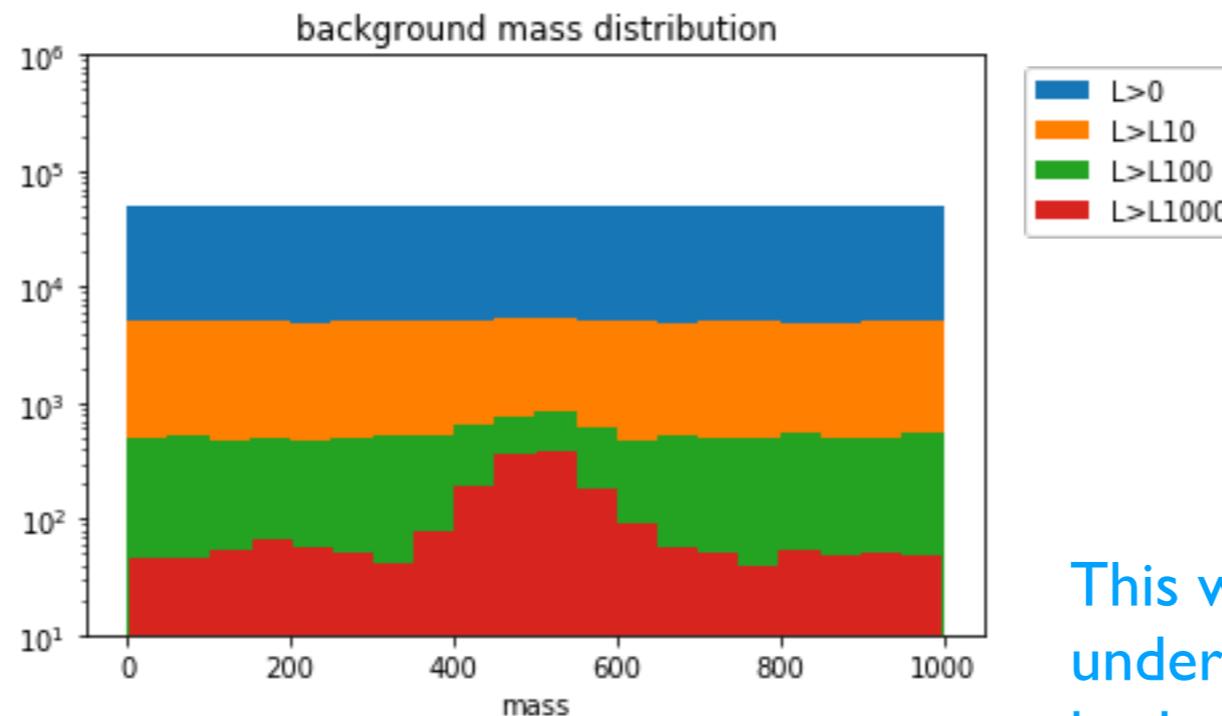
- adversarial decorrelation (Louppe et al 1611.01046, Shimmin et al 1703.03507)
- DisCo decorrelation (Kasieczka & DS 2001.05310)

# Beyond classification — decorrelation

Could have real-world applications, e.g. in designing fairer AIs?

Raw tagger performance not the only consideration.

For robust background estimation, generally need to ensure tagger does not sculpt the background mass distribution.



This would greatly underestimate the background in the SR

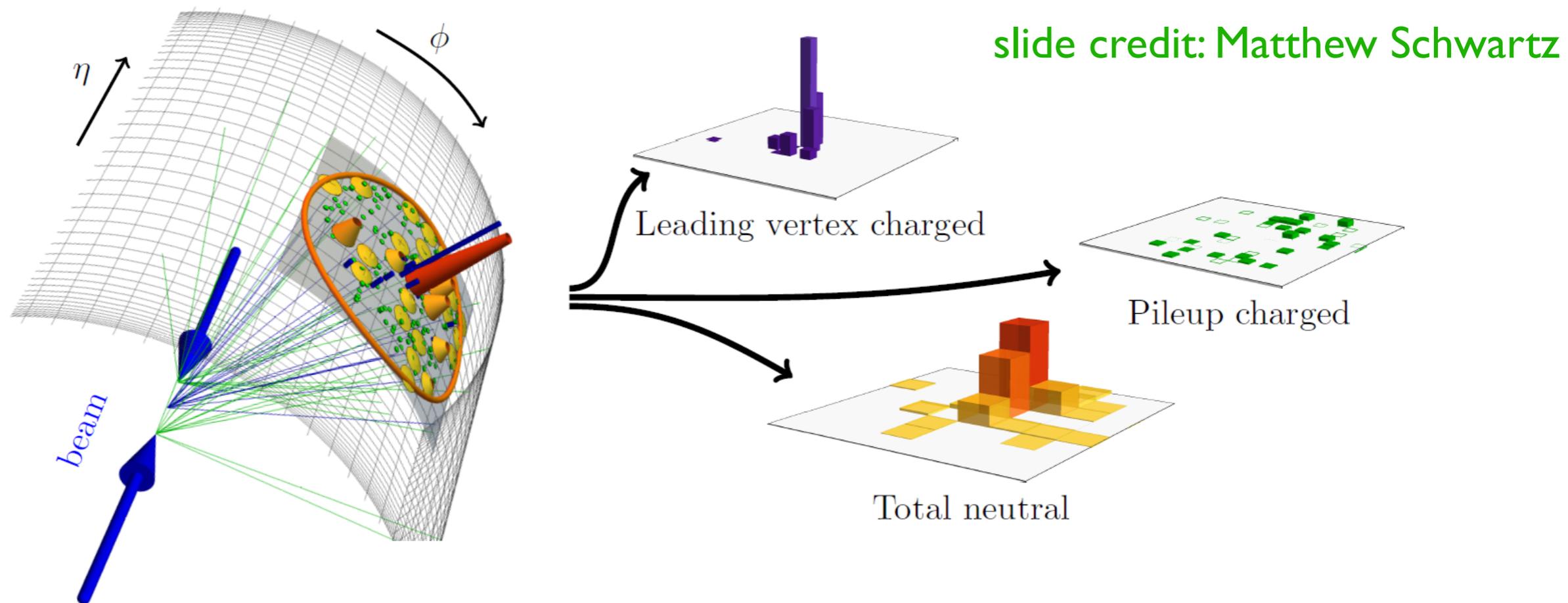
Two approaches with deep learning:

- adversarial decorrelation (Louppe et al 1611.01046, Shimmin et al 1703.03507)
- DisCo decorrelation (Kasieczka & DS 2001.05310)

# Beyond classification — pileup regression

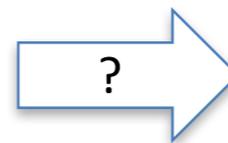
PileUp Mitigation with Machine Learning (PUMML)

(Komiske, Metodiev, Nachman & Schwartz I707.08600)



Can measure

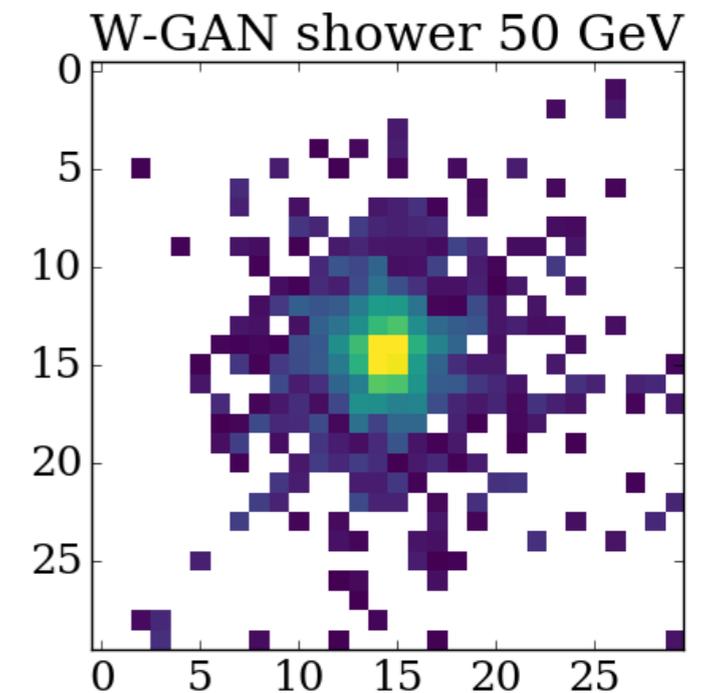
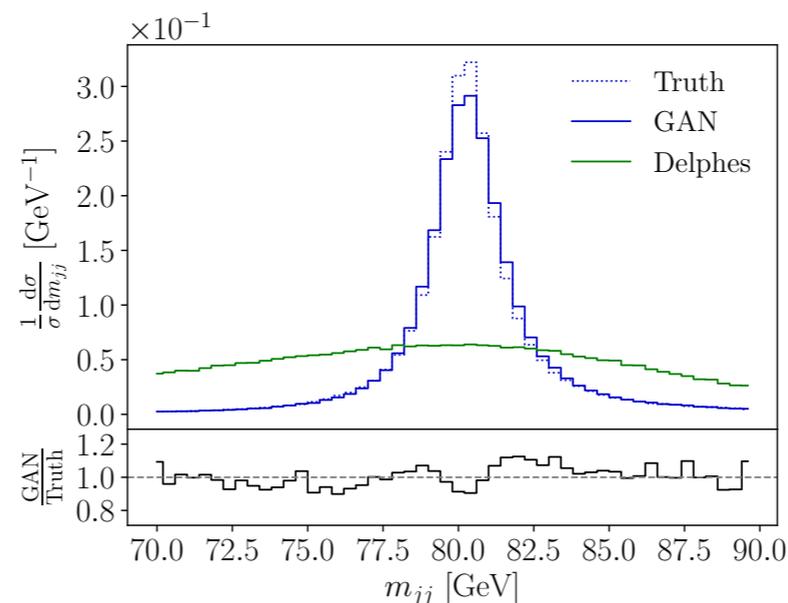
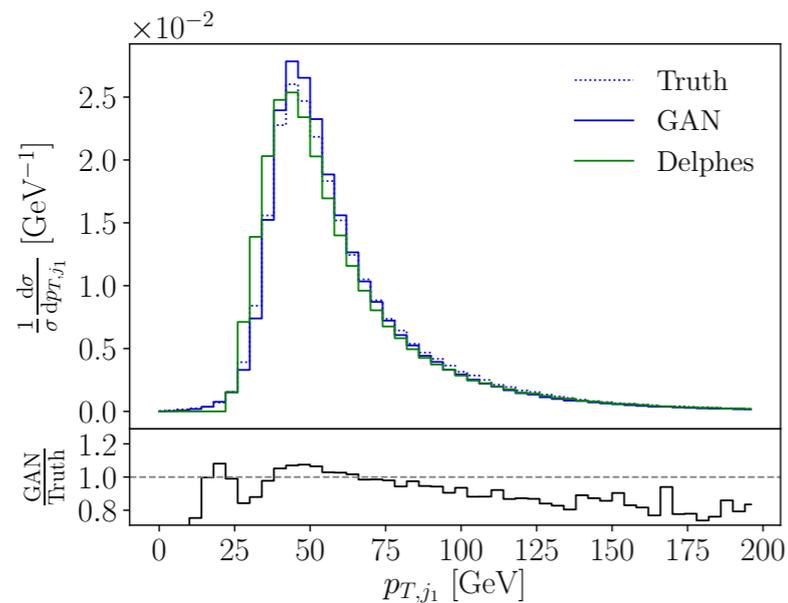
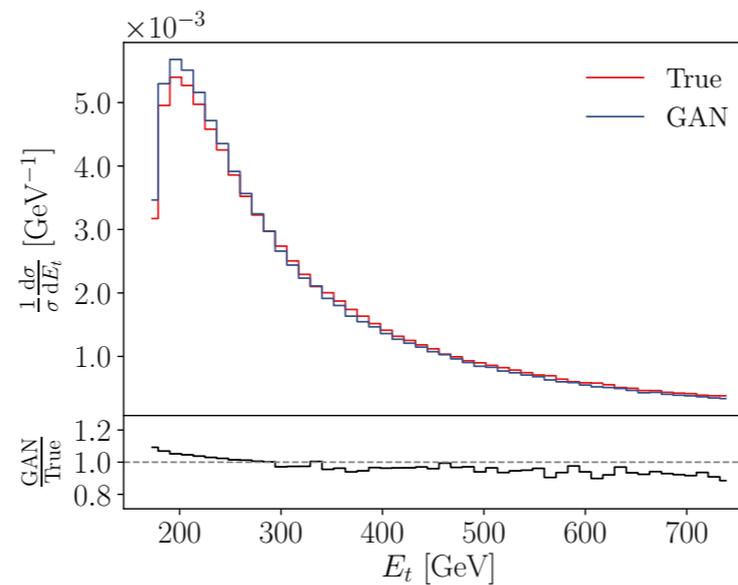
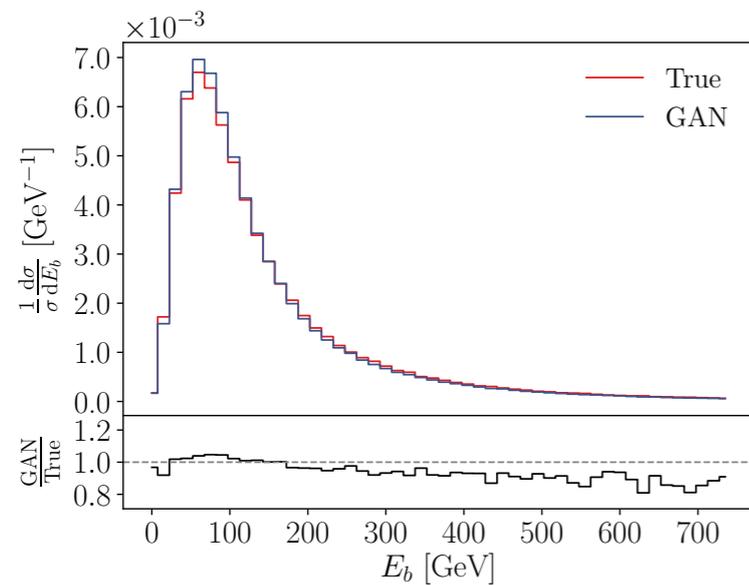
1. Leading vertex charged particles
2. Pileup charged particles
3. Total neutral particles



Leading vertex  
neutral particles

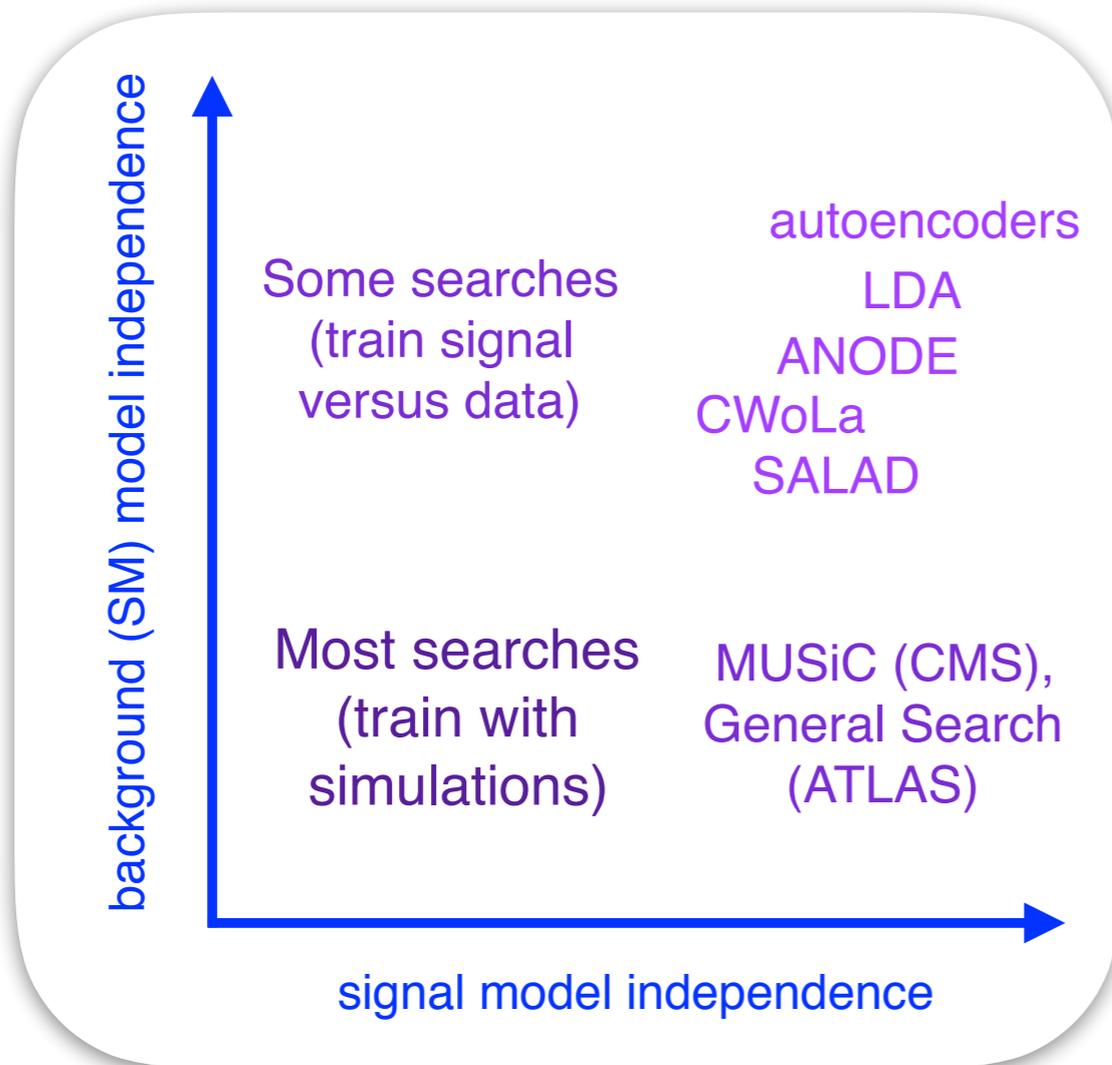
# Beyond classification — $G_{\text{erm}}$ ANs for everything

Can use GANs for event generation, fast detector simulation, unfolding, ...

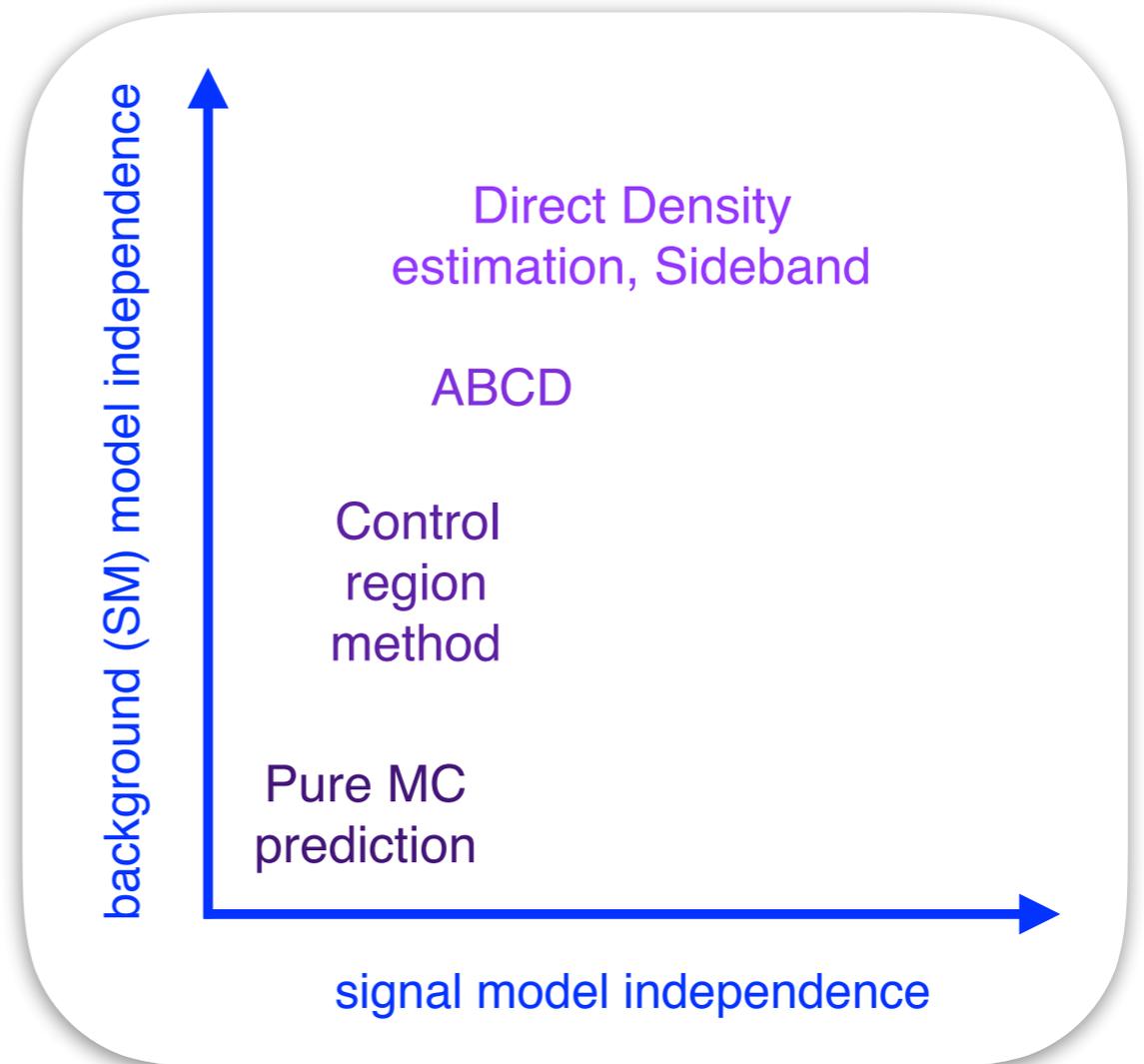


# Way beyond classification: anomaly detection

from Nachman & DS 2001.04990



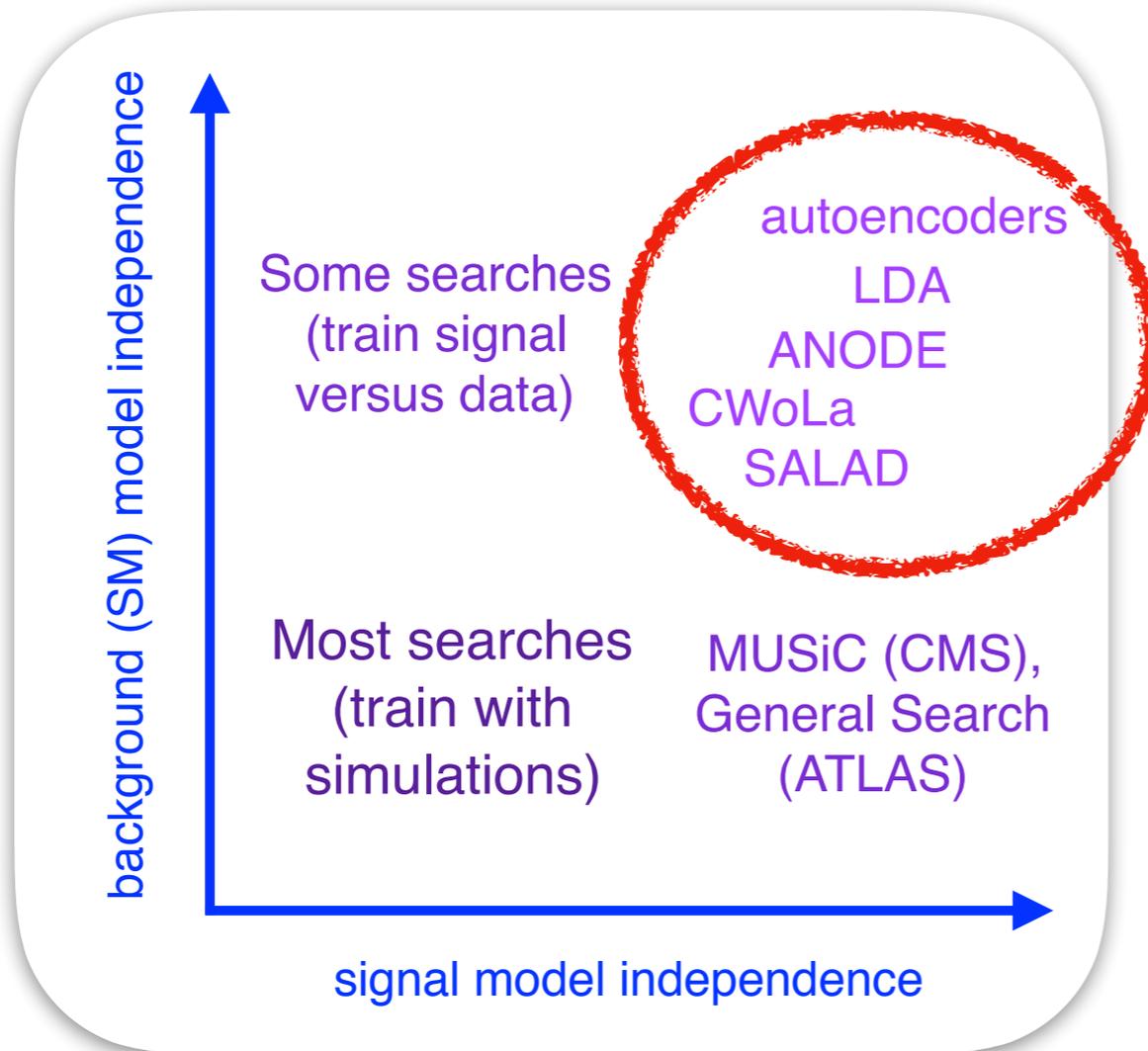
(a) Signal sensitivity



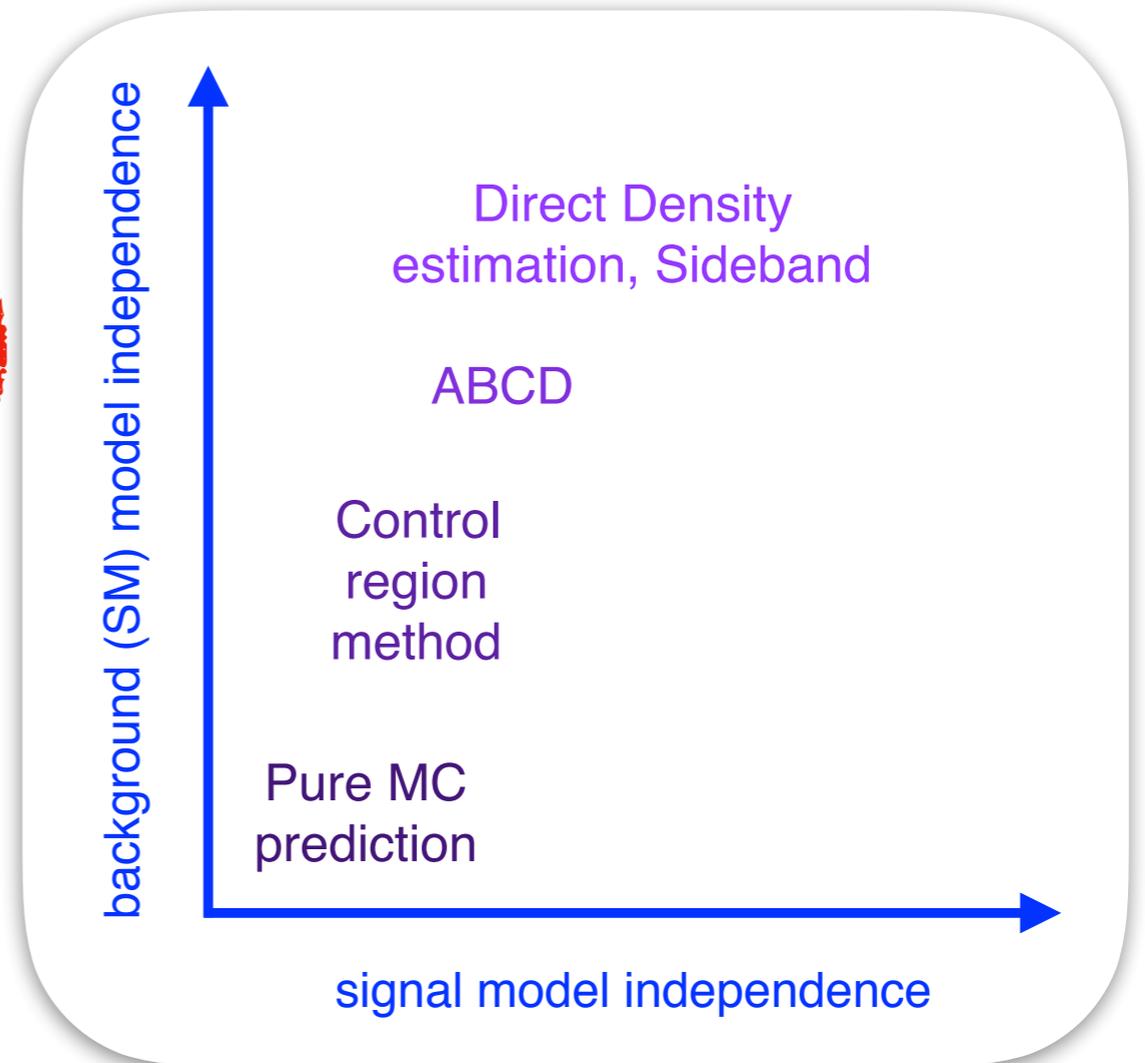
(b) Background specificity

# Way beyond classification: anomaly detection

from Nachman & DS 2001.04990



(a) Signal sensitivity



(b) Background specificity

Many interesting new ideas for model-independent searches being proposed

# Summary/Outlook

Deep learning is revolutionizing nearly every aspect of high-energy physics.

- Tagging
- Pileup
- Event generation
- Detector simulation
- Triggering
- Searches for NP
- ....

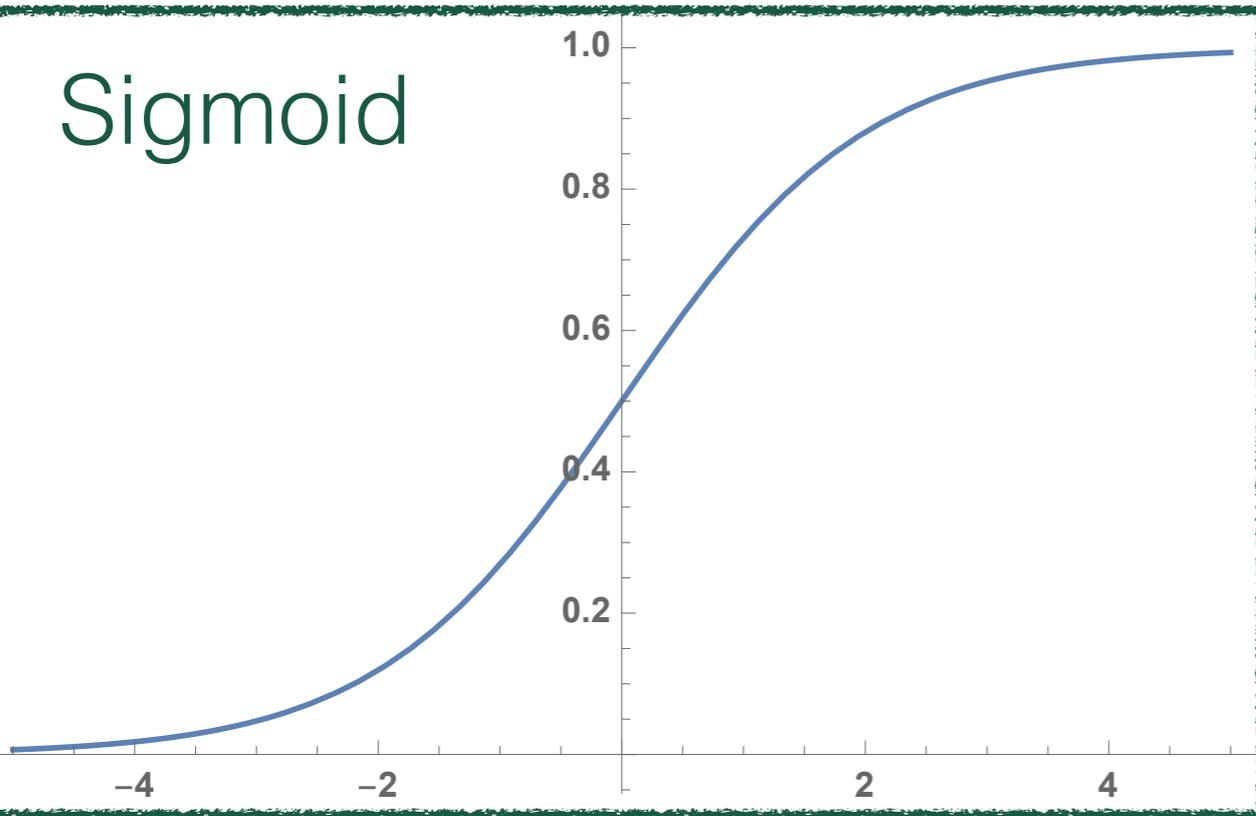
Exciting times are ahead!

**Thanks for your attention!**

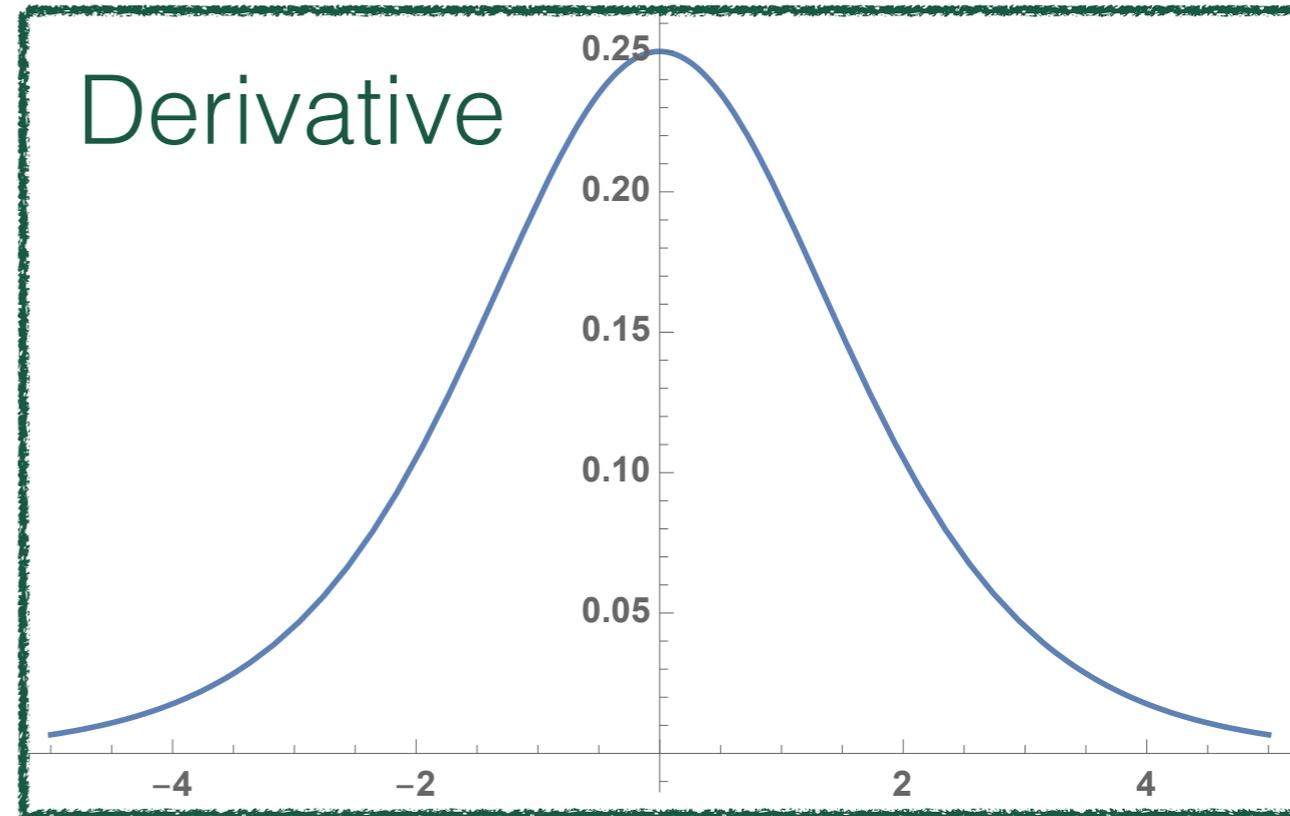
**Supplementary material**

# Disappearing Gradient

Sigmoid



Derivative



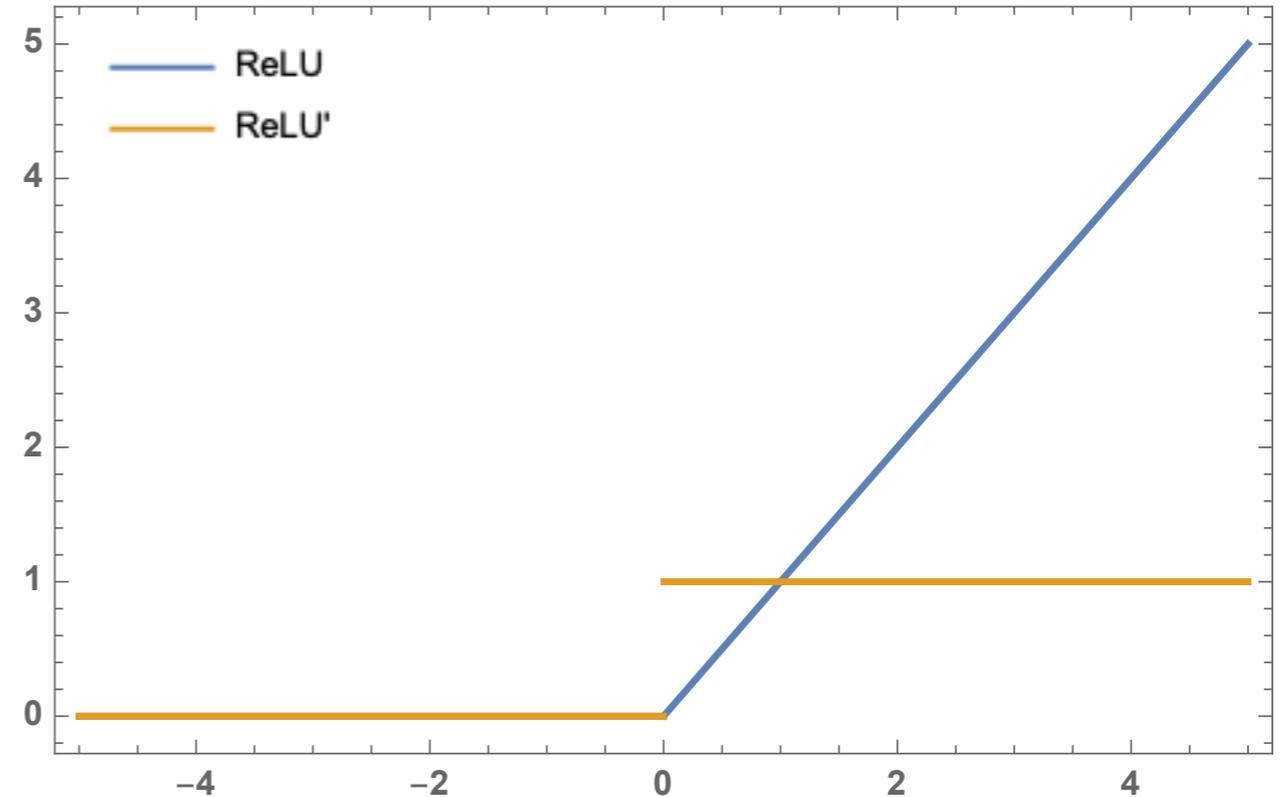
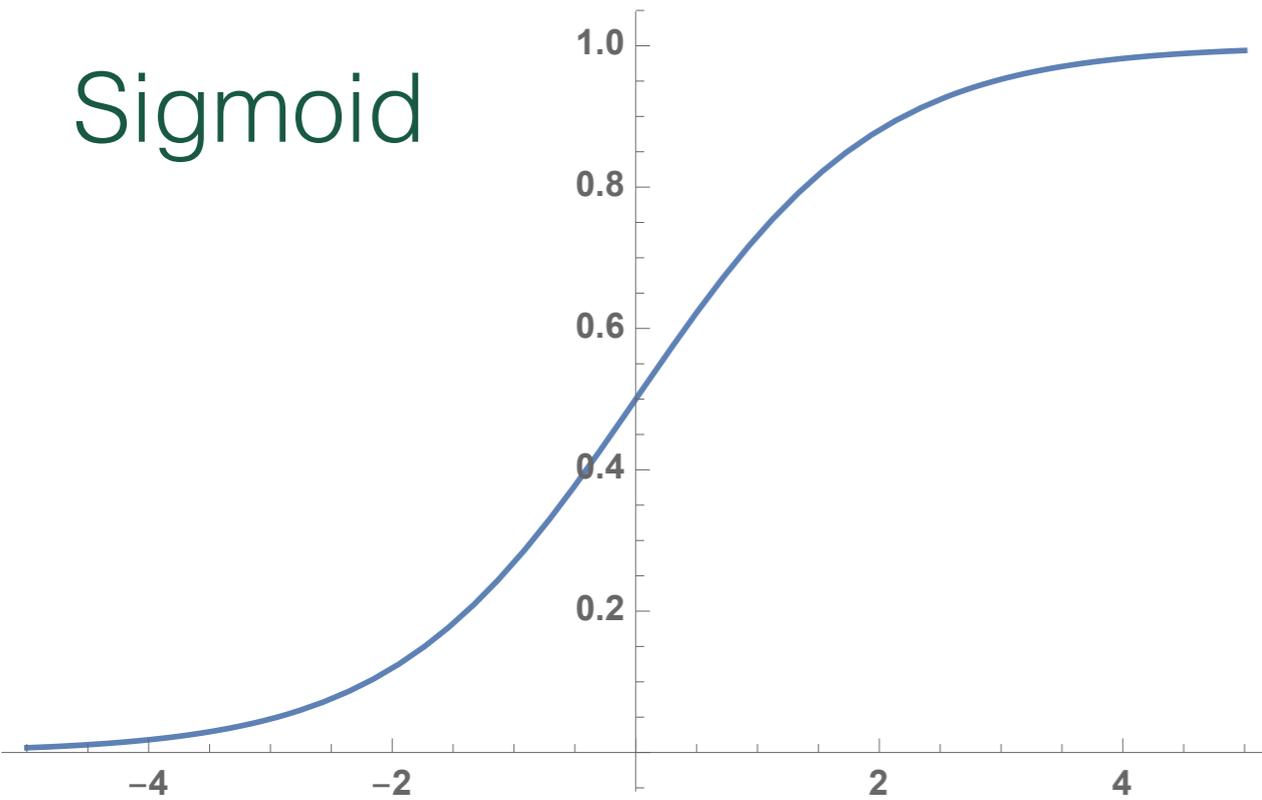
Chain rule for gradient of network involves multiple factors of the derivative multiplied together

$$(0.25)^4 = 0.0039$$

Deep networks with Sigmoid activations have exponentially hard time training early layers

# Disappearing Gradient

Sigmoid



Using the Rectified Linear Unit (ReLU) solves this problem.

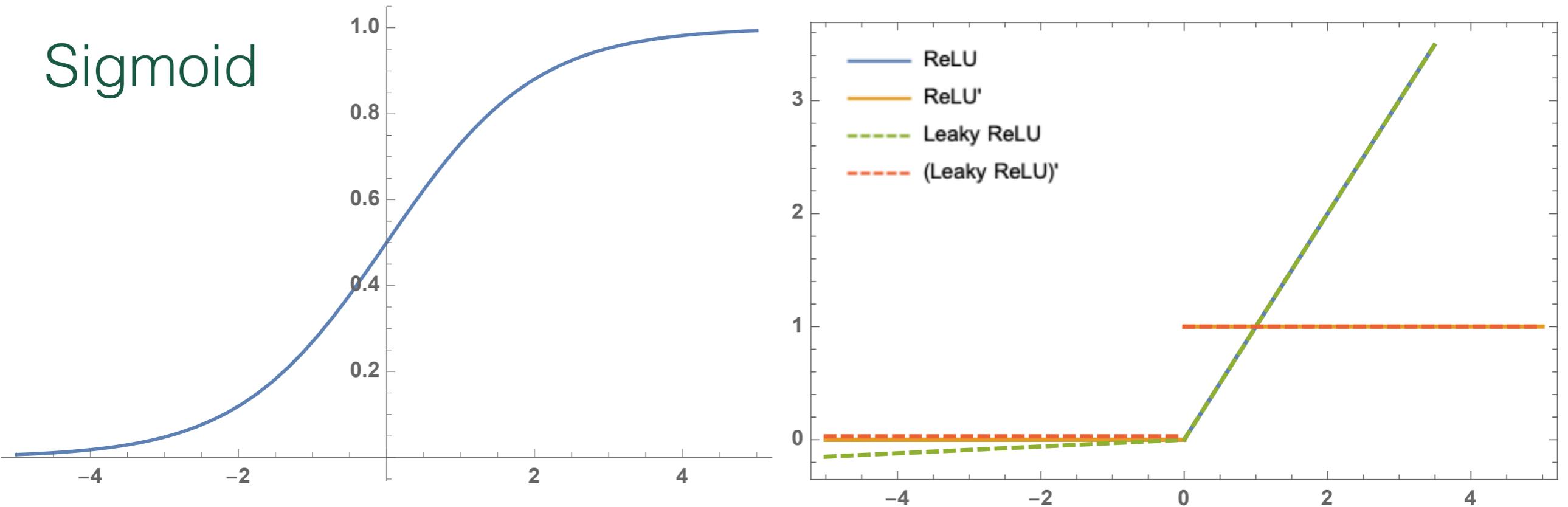
$$\text{ReLU}(x) = \{0 \text{ if } x \leq 0, x \text{ if } x > 0\}$$

Still has nonlinearity which allows network to learn complicated patterns

Nodes can die (derivative always 0 so cannot update)

# Disappearing Gradient

Sigmoid



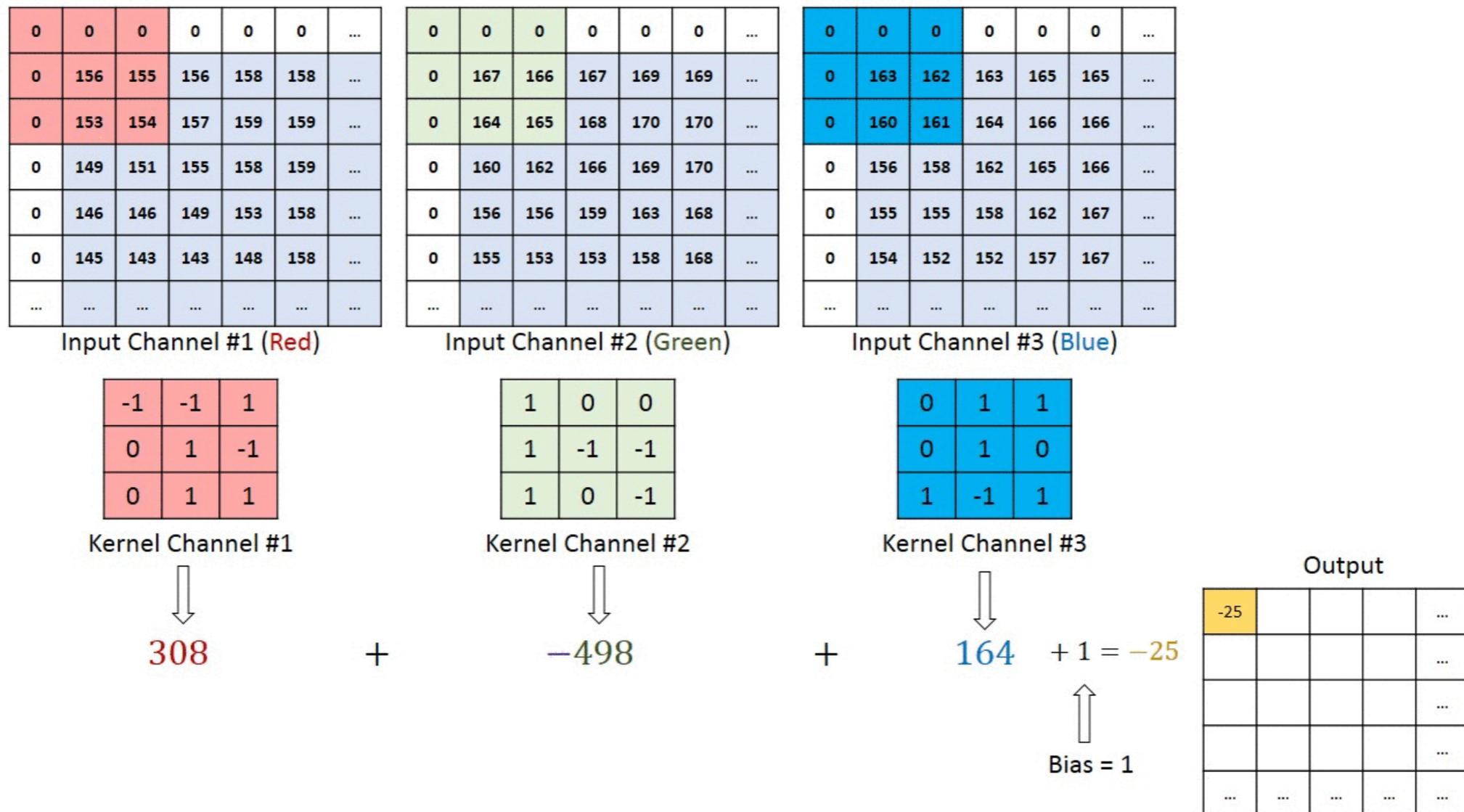
Leaky Rectified Linear Unit (LeakyReLU) solves this problem.

$$\text{LeakyReLU}(x) = \{\alpha * x \text{ if } x \leq 0, x \text{ if } x > 0\}$$

I have never had to use this in practice

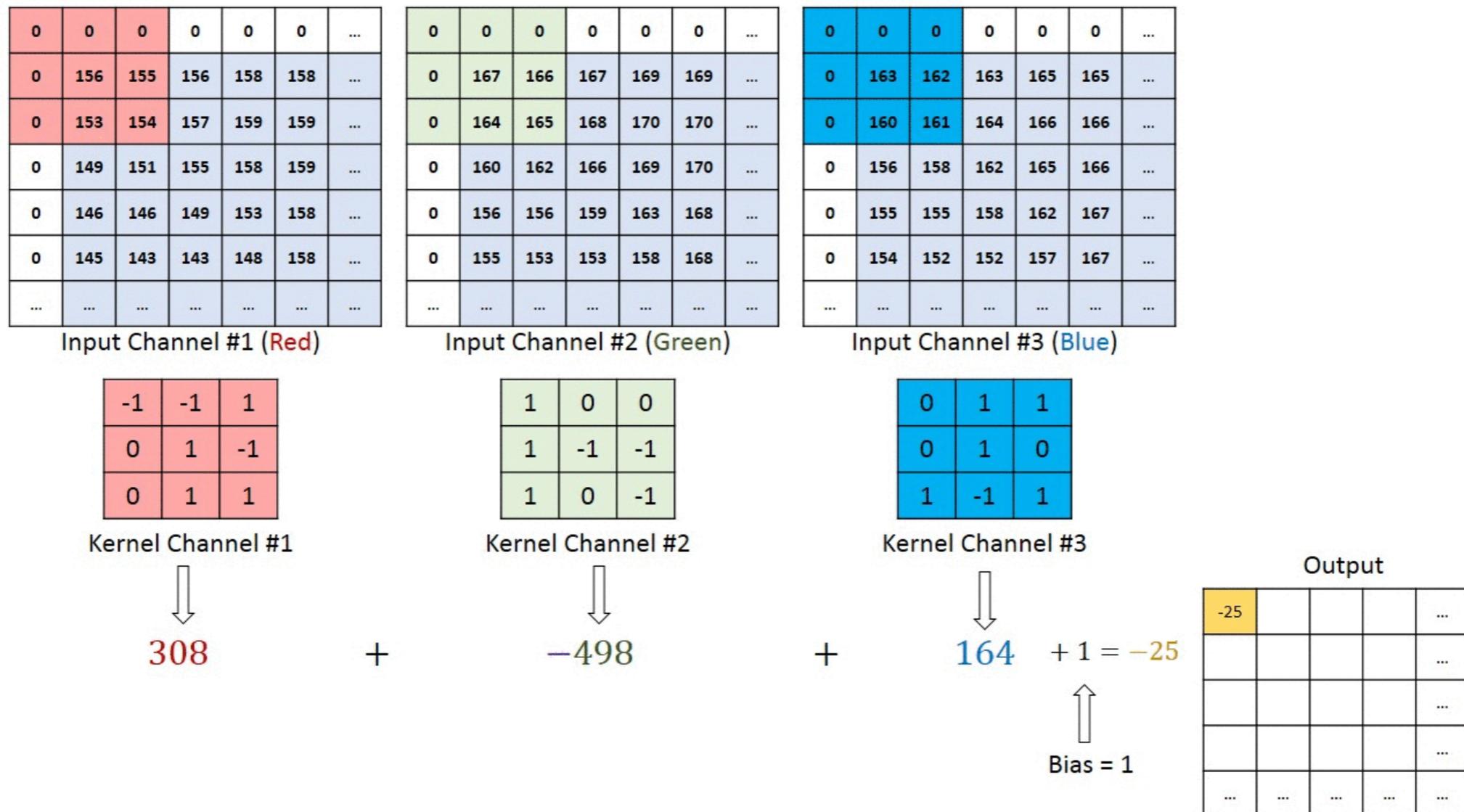
# Convolutional neural network (CNN)

Dealing with multiple channels is straightforward — just enlarge filter to include channel dimension (3d filter) and perform element-wise multiplication along channel dimension as well.



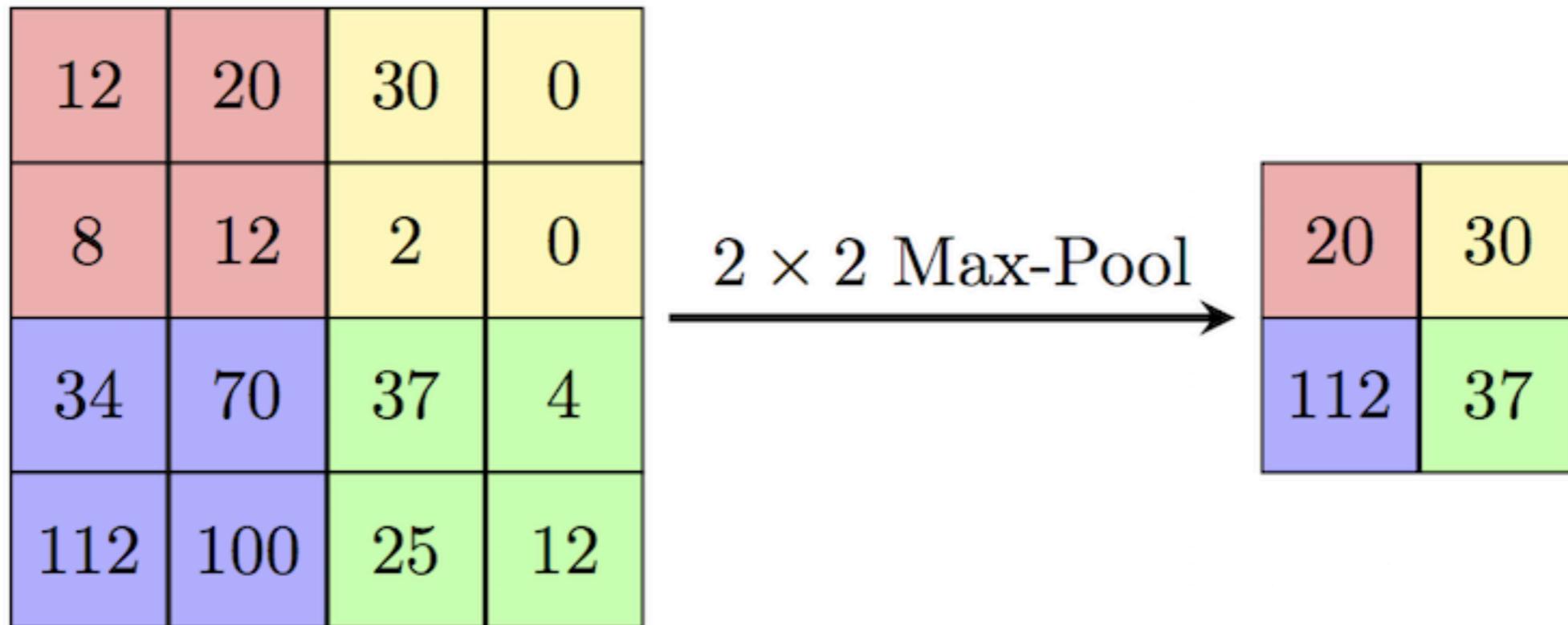
# Convolutional neural network (CNN)

Dealing with multiple channels is straightforward — just enlarge filter to include channel dimension (3d filter) and perform element-wise multiplication along channel dimension as well.



# Convolutional neural network (CNN)

“Max pooling”



Reduces image size, reducing parameters and mitigating overfitting

Allows NN to find spatially larger, often higher-level features

# Recurrent Neural Networks (RNNs)

Popular architecture for natural language processing (sentence completion, autocorrect, translation, speech recognition...)

Starting point: sequence of numbers  $x_1, x_2, x_3, \dots$

Suppose we want to predict the next number in the sequence?

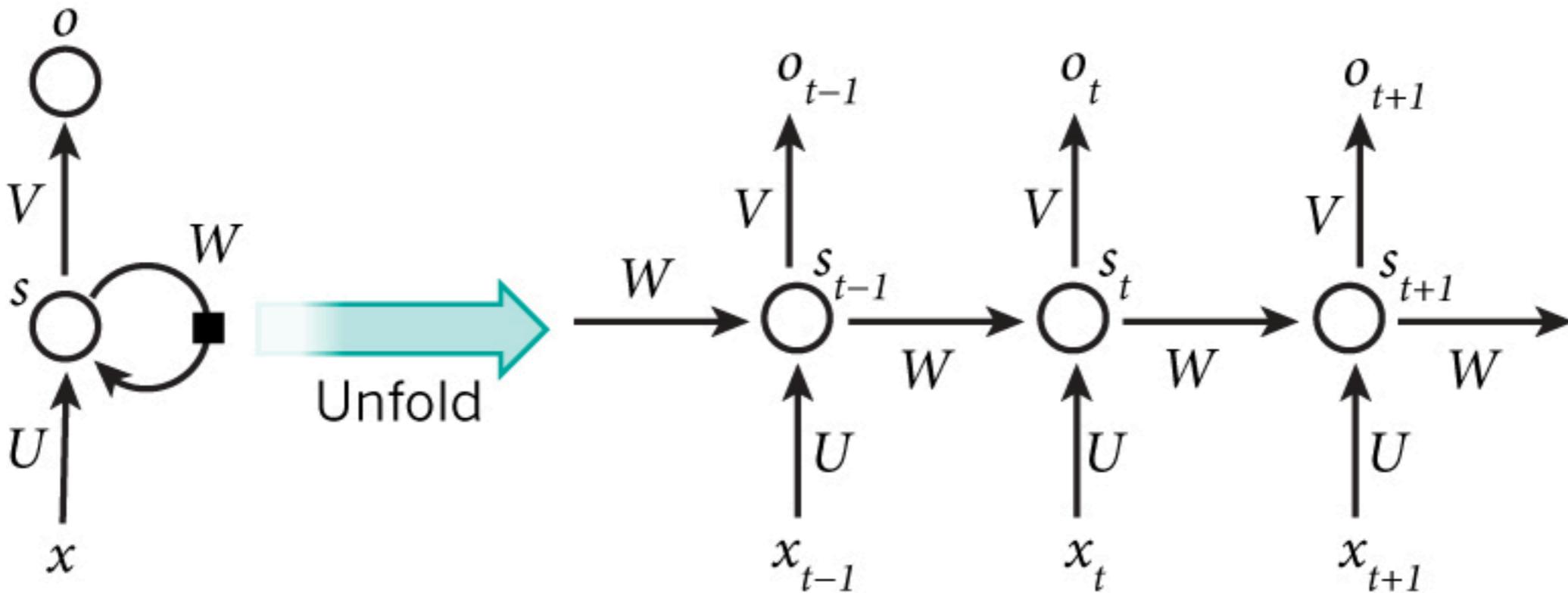
Idea of RNN:

- feed data sequentially to NN
- after each time step update *hidden state*. Hidden state encodes “memory” of sequence.
- Use hidden state to make predictions.

# Recurrent Neural Networks (RNNs)

Basic RNN architecture

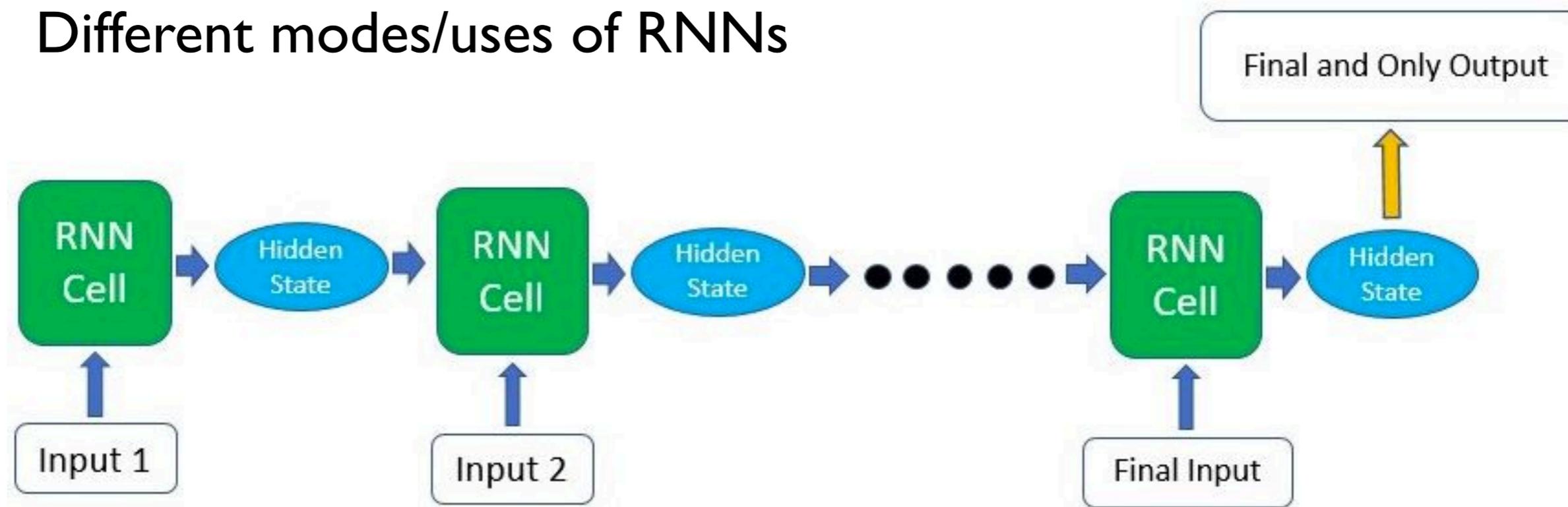
predictions, e.g.  $o_t = \text{softmax}(V s_t)$



$$s_t = f(Ux_t + Ws_{t-1}) \quad \text{hidden state after time step } t$$

# Recurrent Neural Networks (RNNs)

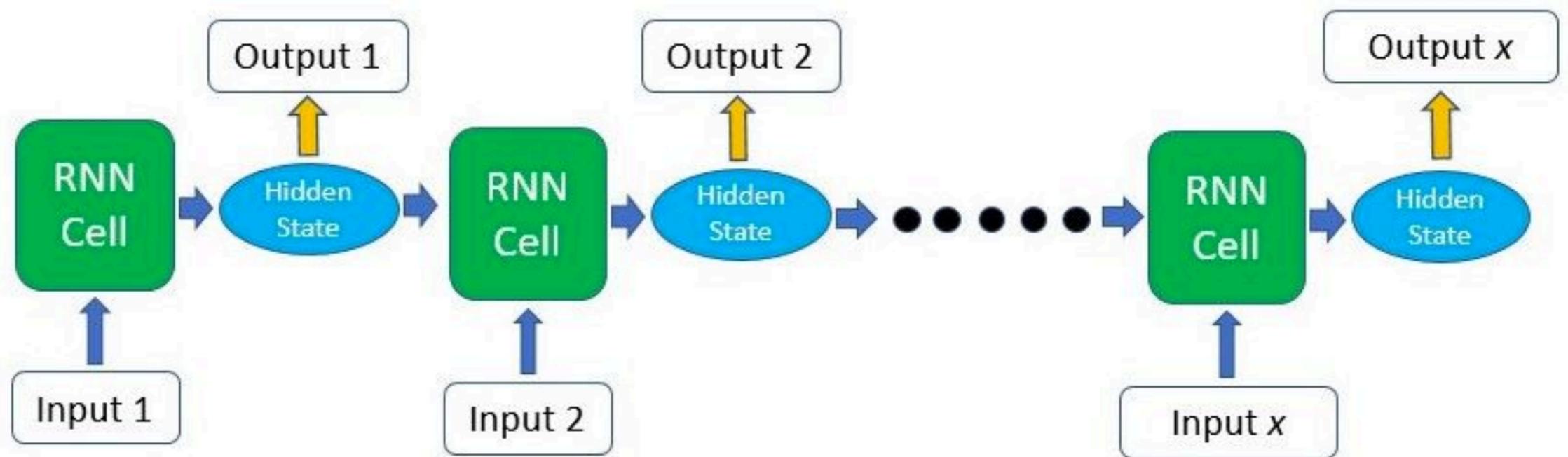
Different modes/uses of RNNs



sequence classification, regression

# Recurrent Neural Networks (RNNs)

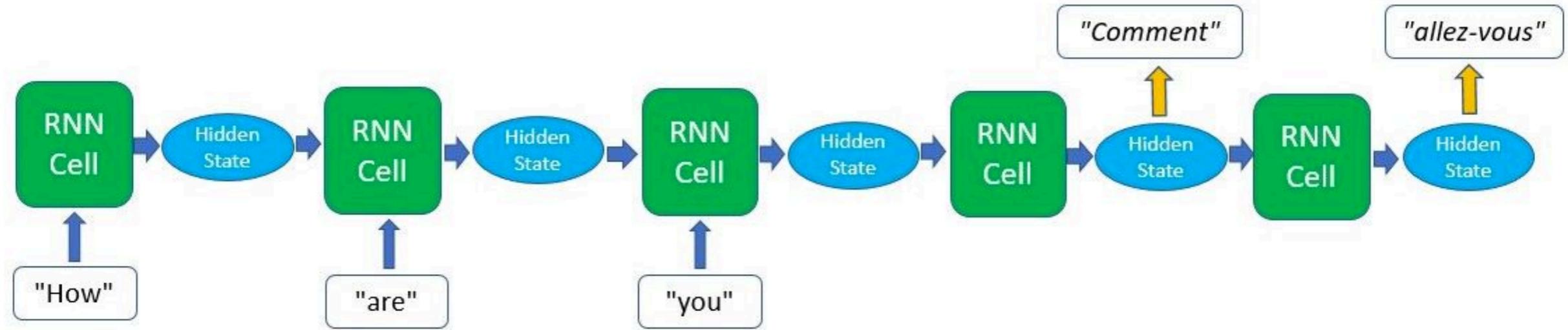
Different modes/uses of RNNs



real-time prediction

# Recurrent Neural Networks (RNNs)

## Different modes/uses of RNNs

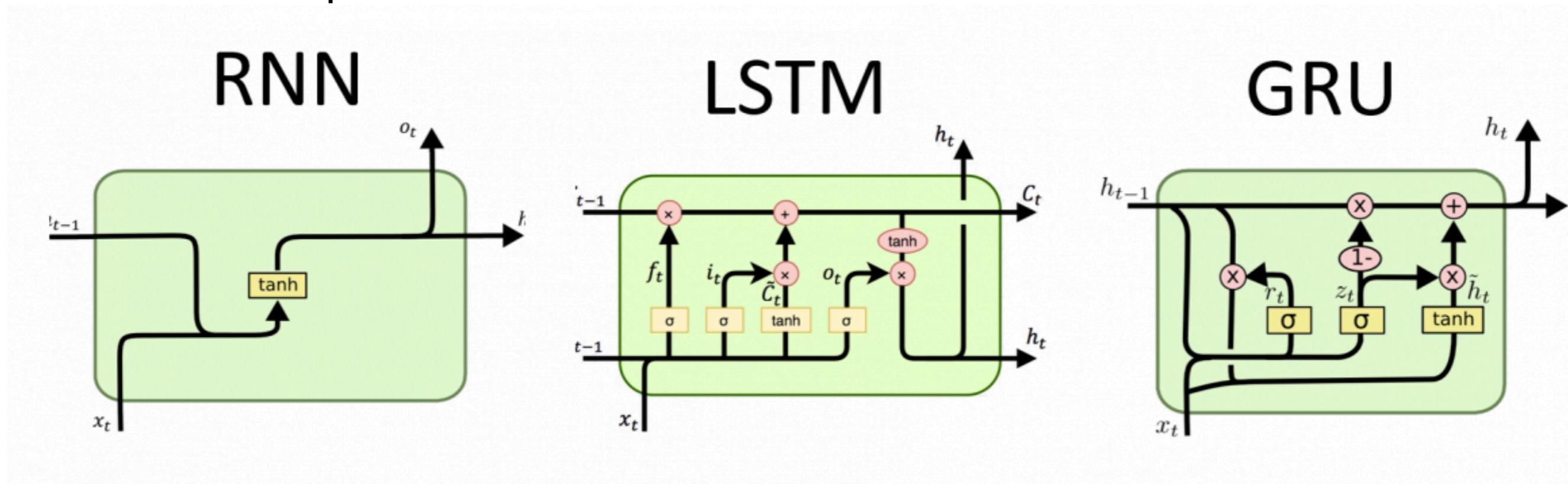


sequence-to-sequence

# Recurrent Neural Networks (RNNs)

Simple RNNs applied to long sequences have a very serious exploding/vanishing gradient problem.

Prevents them from “remembering” relevant information from earlier in the sequence.



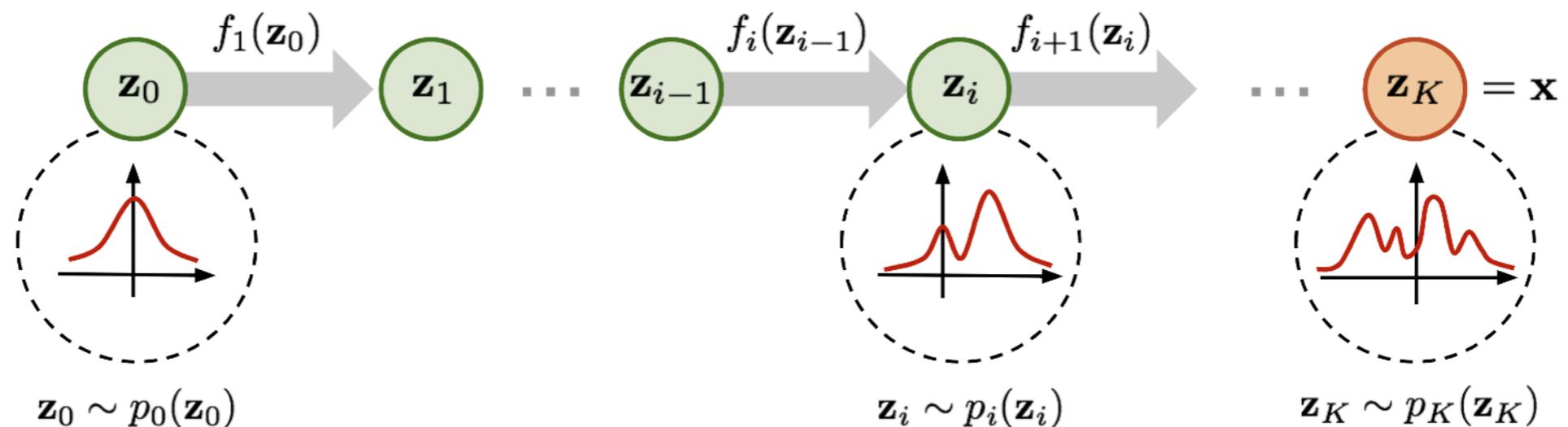
“Long-short term memory” and “Gated recurrent units” are two methods commonly used to solve the gradient problem and improve performance.

# Normalizing flows

Rezende & Mohamed 1505.05770

Recently a lot of excitement and progress in the problem of density estimation with neural networks.

Idea: map original distribution to normal distribution through series of invertible transformations.



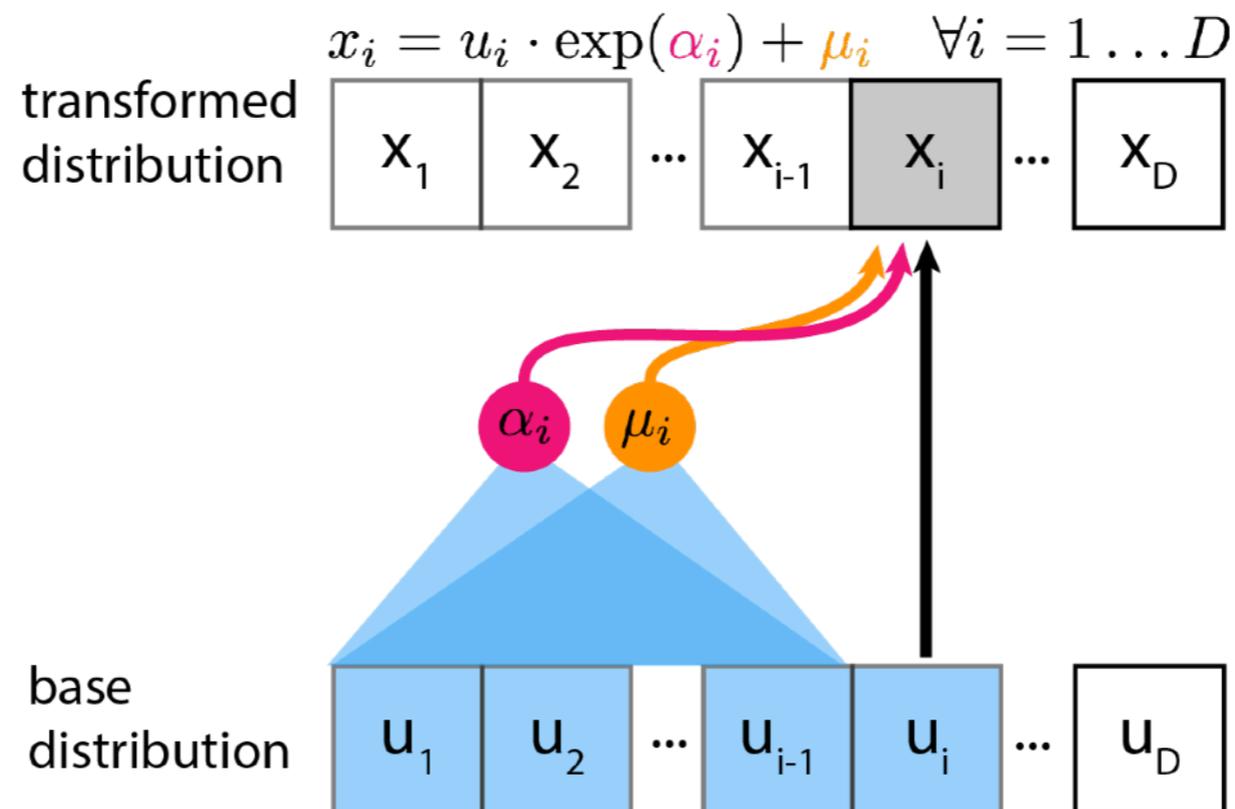
Examples: RealNVP, NICE, Glow, ...

# Autoregressive flows

$$p(x) = \prod_i p(x_i | x_{1:i-1})$$

Special type of normalizing flows. Learn probability density of each coordinate conditioned on previous coordinates.

Transformation upper triangular — automatically invertible. Allows for more expressive transformations.



Examples: MADE, MAF, IAF, NAF, PixelRNN, Wavenet, ...

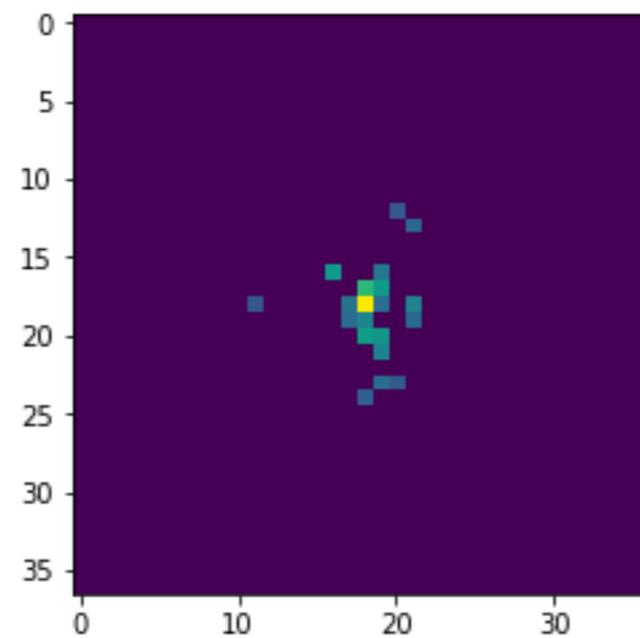
# Top Tagging with CNNs

Macaluso & DS 1803.00107

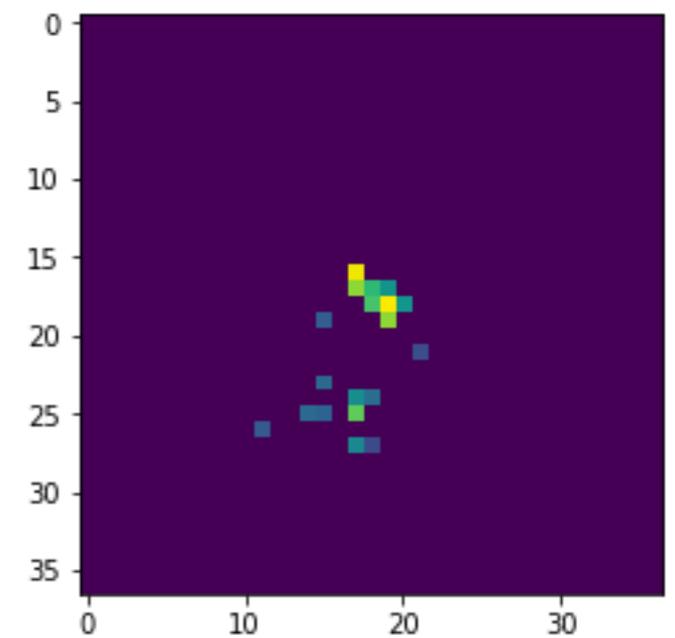
	CMS
Jet sample	13 TeV $p_T \in (800, 900)$ GeV, $ \eta  < 1$ PYTHIA 8 and DELPHES particle-flow match: $\Delta R(t, j) < 0.6$ merge: $\Delta R(t, q) < 0.6$ 1.2M + 1.2M
Image	$37 \times 37$ $\Delta\eta = \Delta\phi = 3.2$
Colors	$(p_T^{neutral}, p_T^{track}, N_{track}, N_{muon})$

Individual images very sparse

QCD



Tops



Building on previous “DeepTop” tagger of [Kasieczka et al 1701.08784](#)

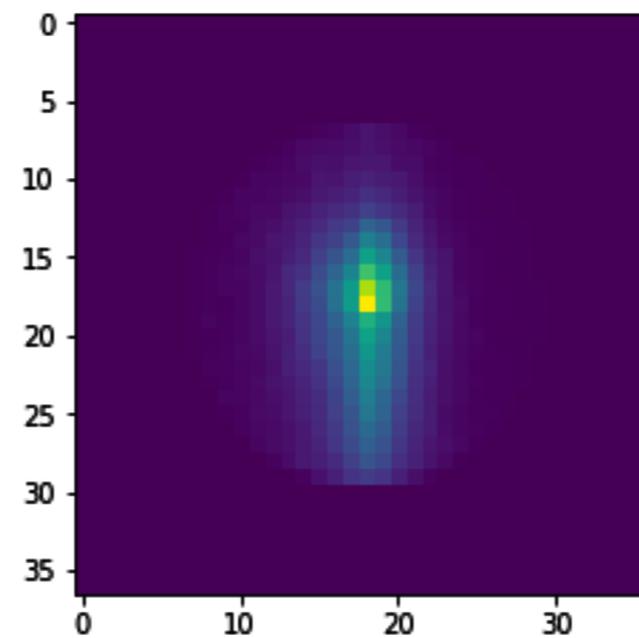
# Top Tagging with CNNs

Macaluso & DS 1803.00107

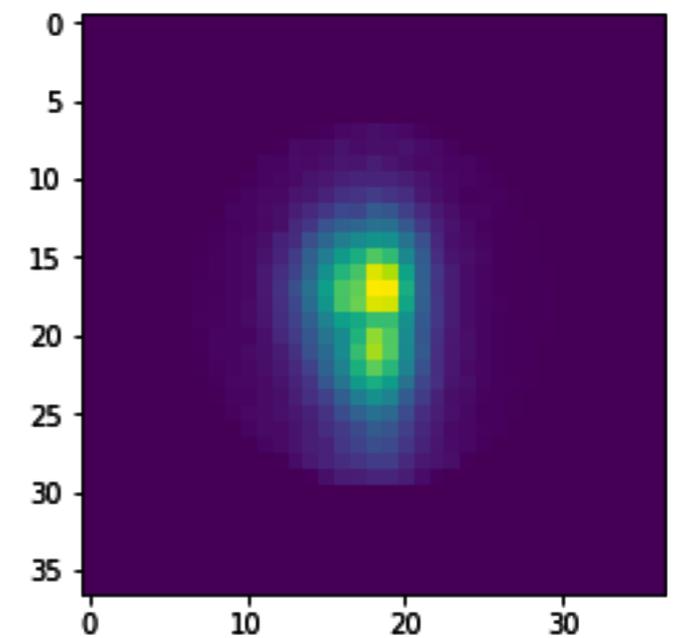
	CMS
Jet sample	13 TeV $p_T \in (800, 900)$ GeV, $ \eta  < 1$ PYTHIA 8 and DELPHES particle-flow match: $\Delta R(t, j) < 0.6$ merge: $\Delta R(t, q) < 0.6$ 1.2M + 1.2M
Image	$37 \times 37$ $\Delta\eta = \Delta\phi = 3.2$
Colors	$(p_T^{neutral}, p_T^{track}, N_{track}, N_{muon})$

Average images clearly different

QCD



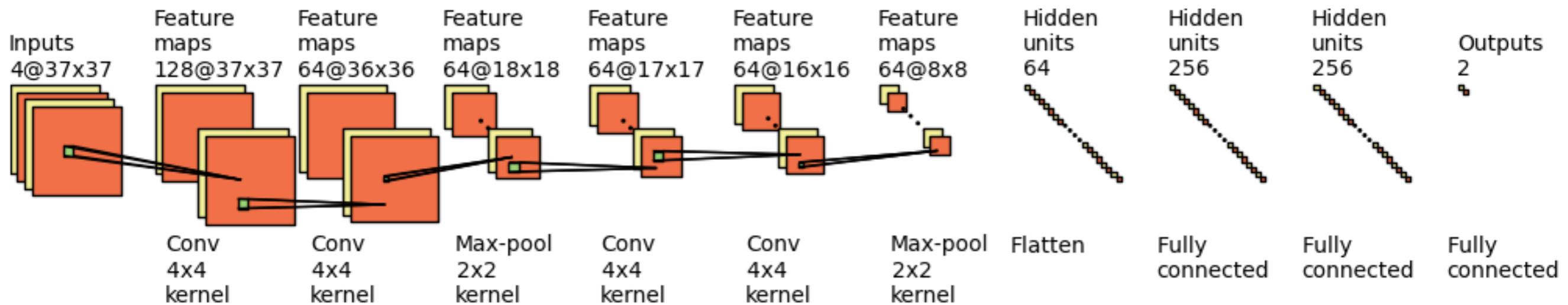
Tops



Building on previous “DeepTop” tagger of [Kasieczka et al 1701.08784](#)

# Top Tagging with CNNs

Macaluso & DS I803.00107



AdaDelta

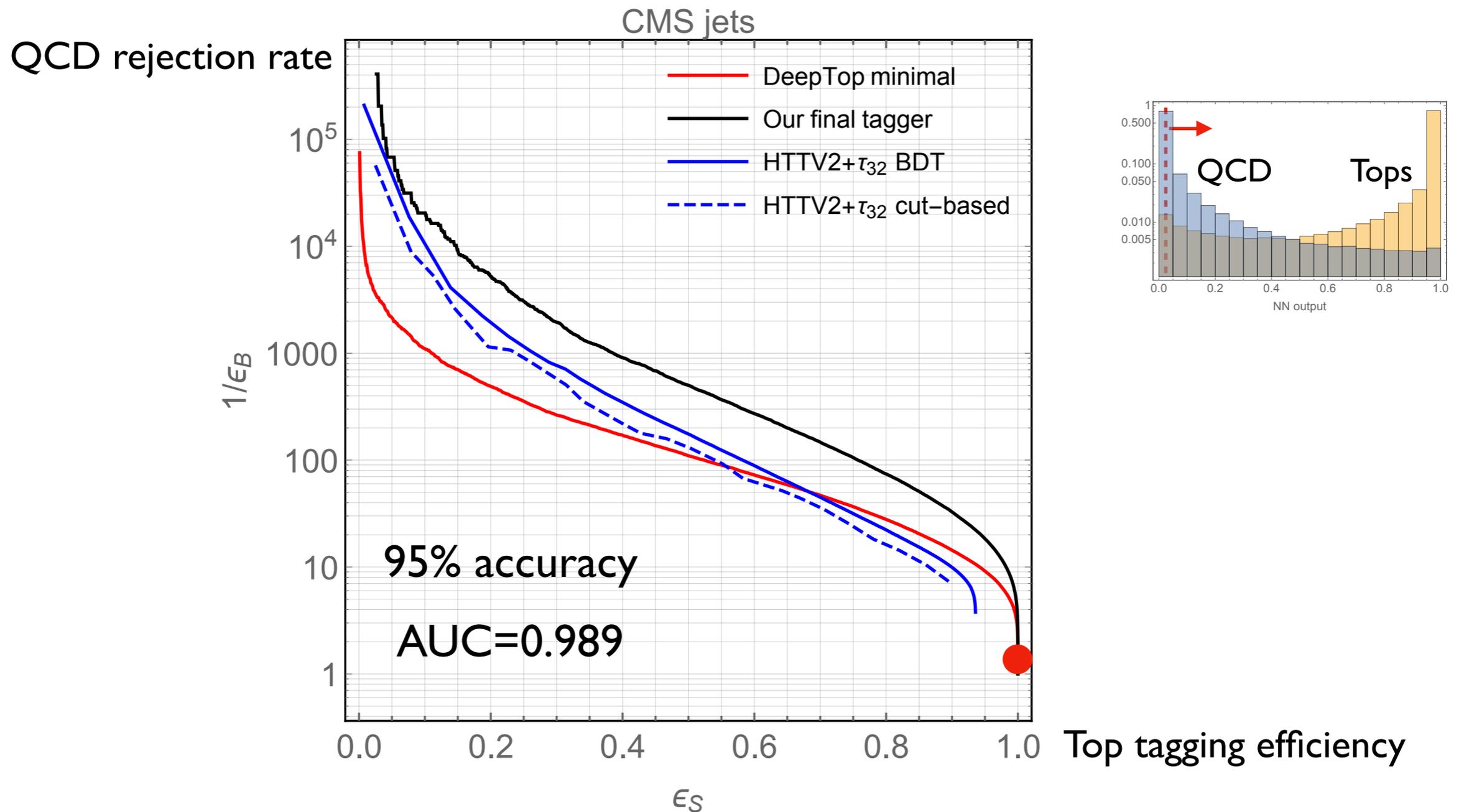
$\eta = 0.3$  with annealing schedule

minibatch size=128

cross entropy loss

# Top Tagging with CNNs

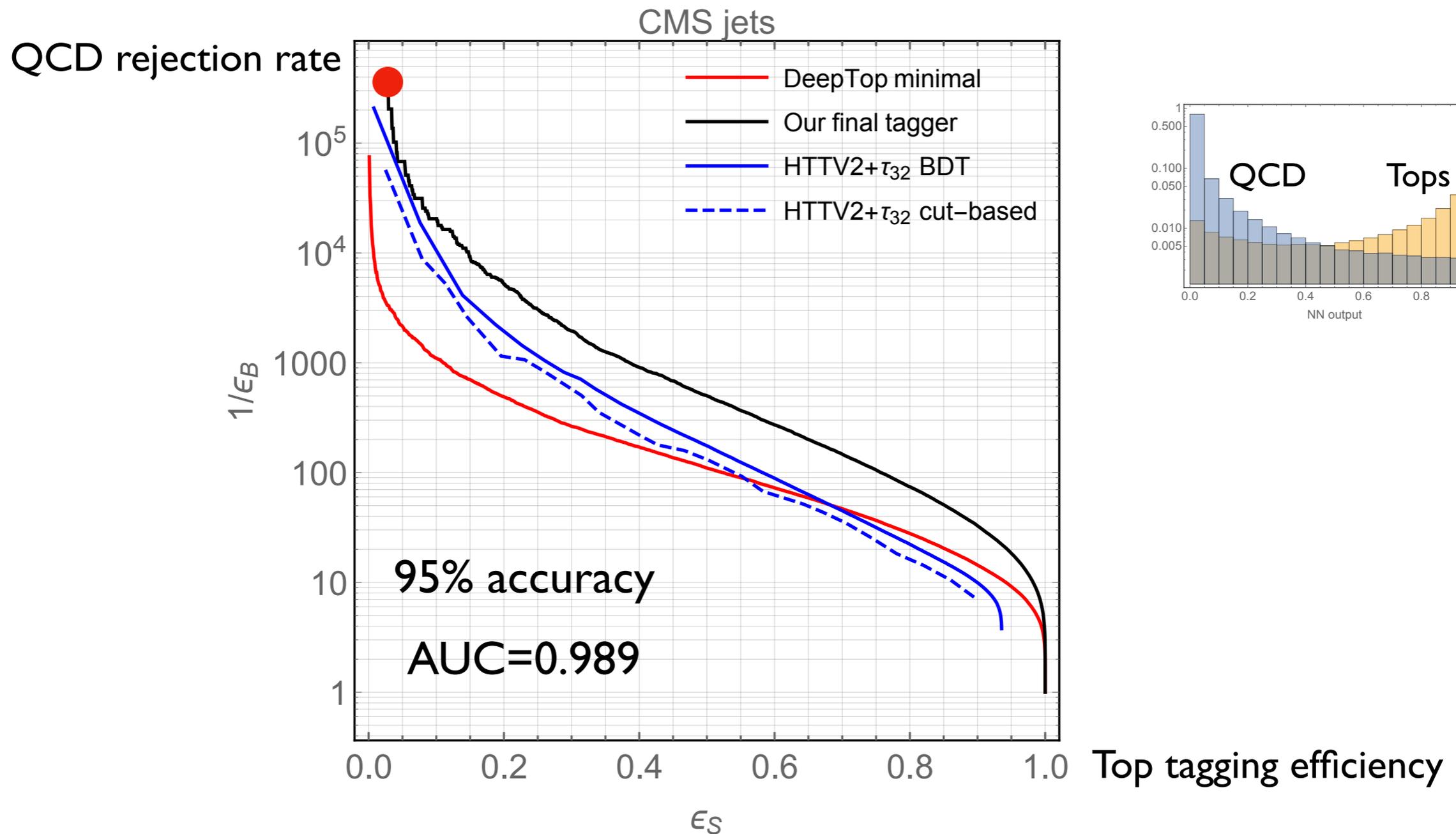
Macaluso & DS | 803.00 | 07



Can achieve factor of  $\sim 3$  improvement over cut-based approaches and BDTs!

# Top Tagging with CNNs

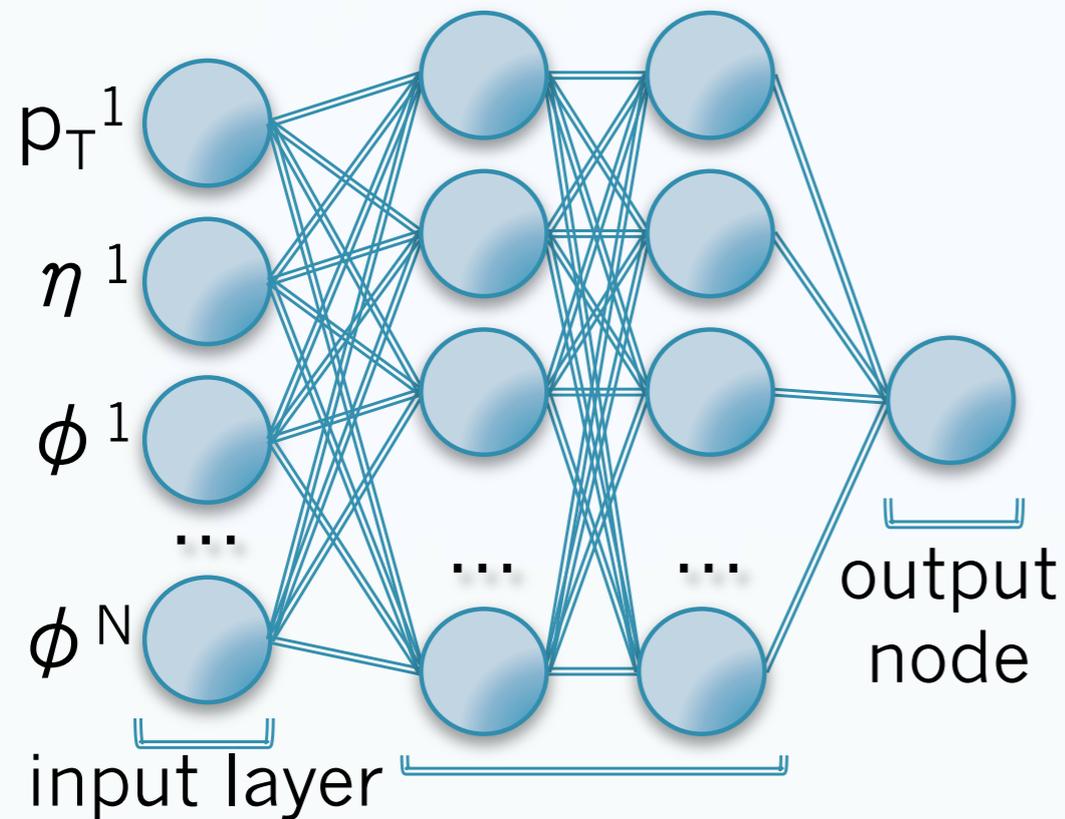
Macaluso & DS | 803.00 | 07



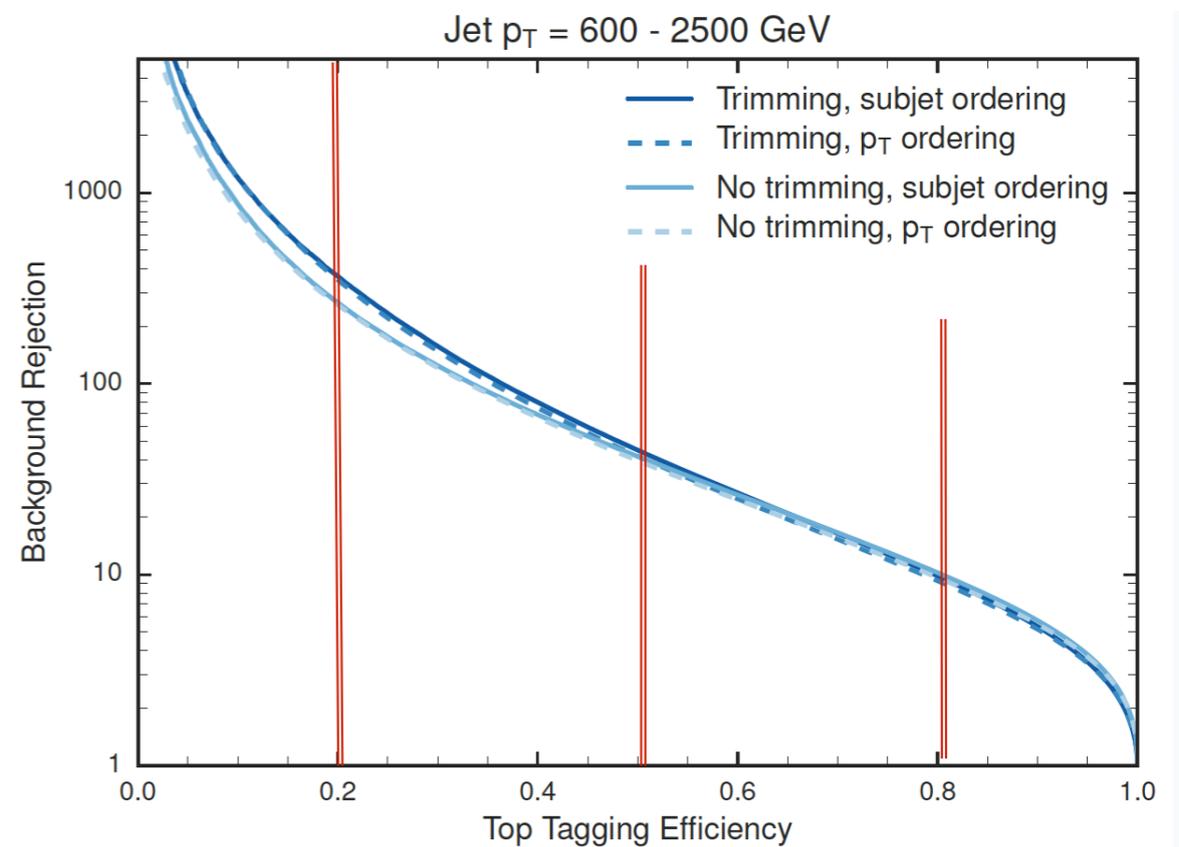
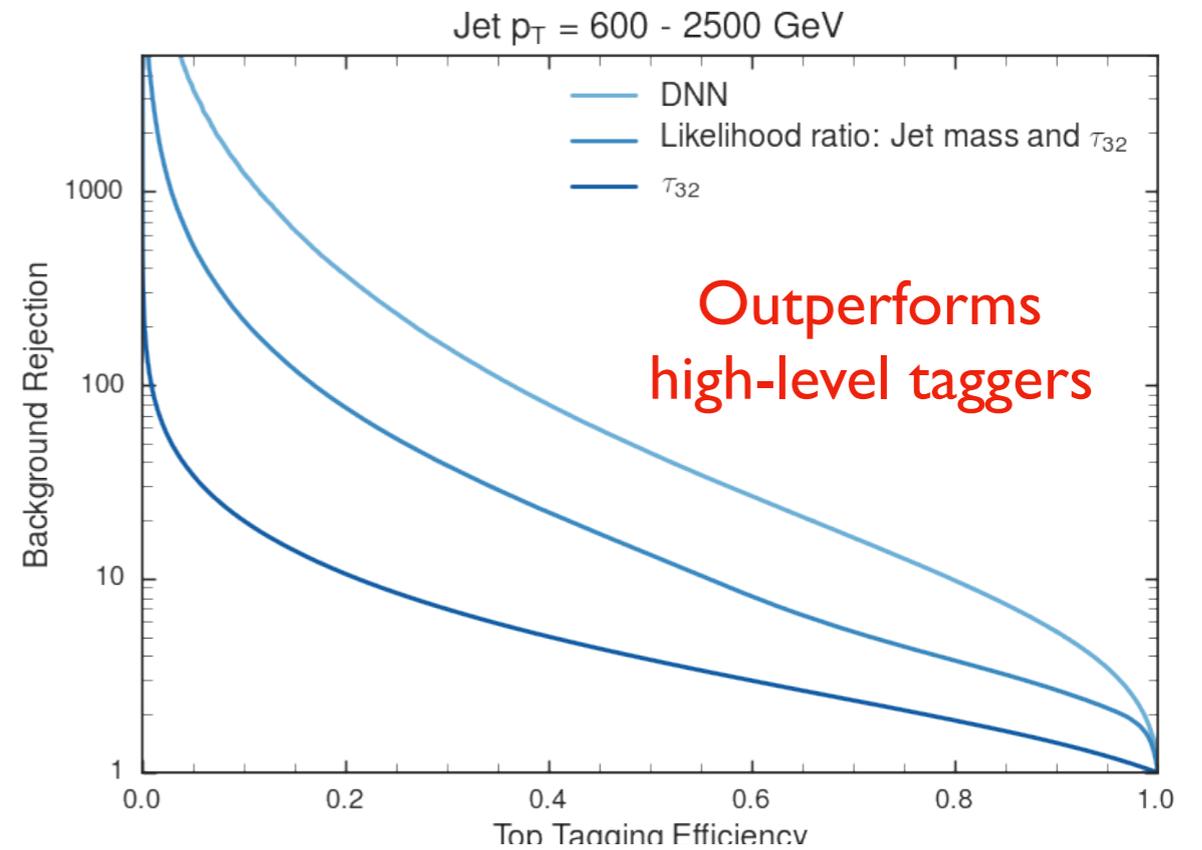
Can achieve factor of  $\sim 3$  improvement over cut-based approaches and BDTs!

# Lists of four vectors

Pearkes et al 1704.02124



5 hidden layers  
300, 150, 50, 10, 5  
nodes / hidden layer

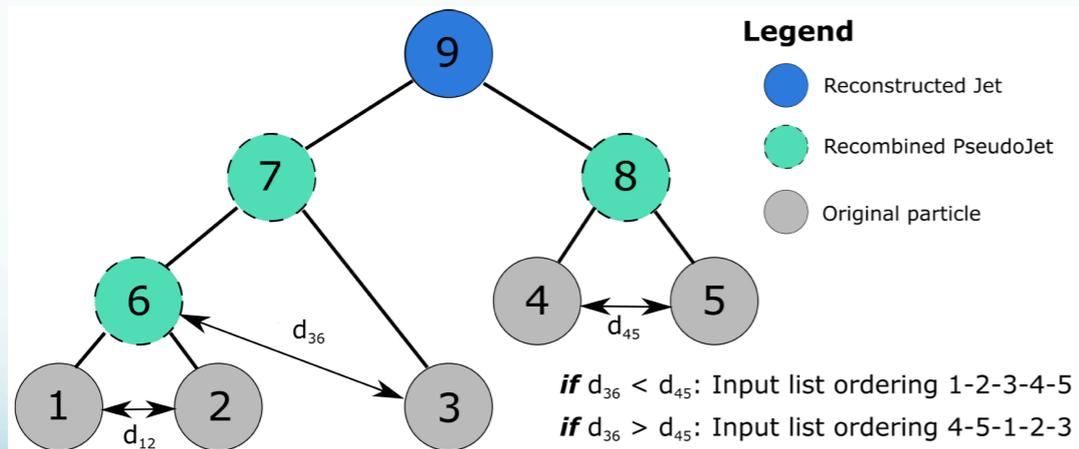


robust against different orderings

# Sequences

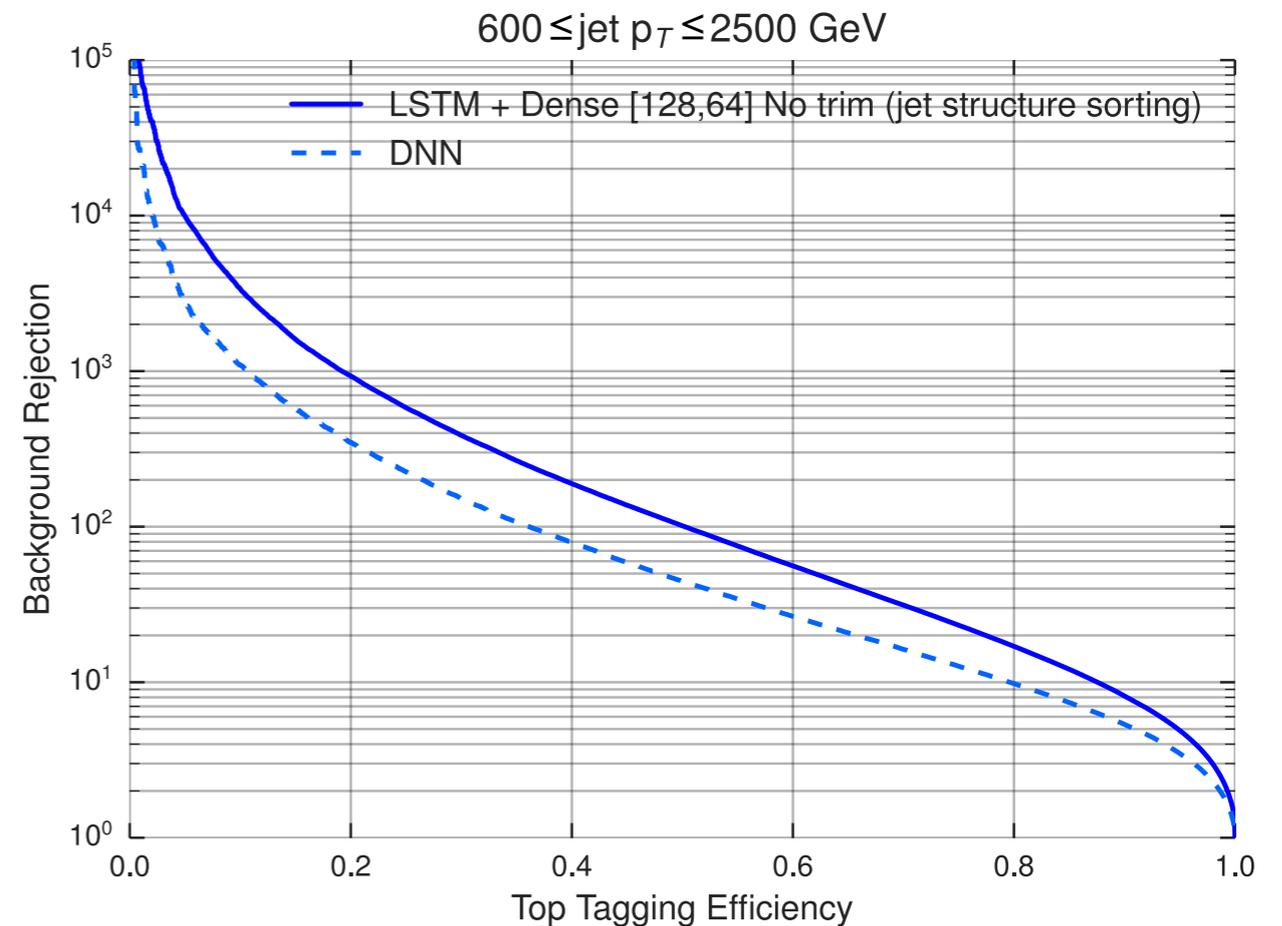
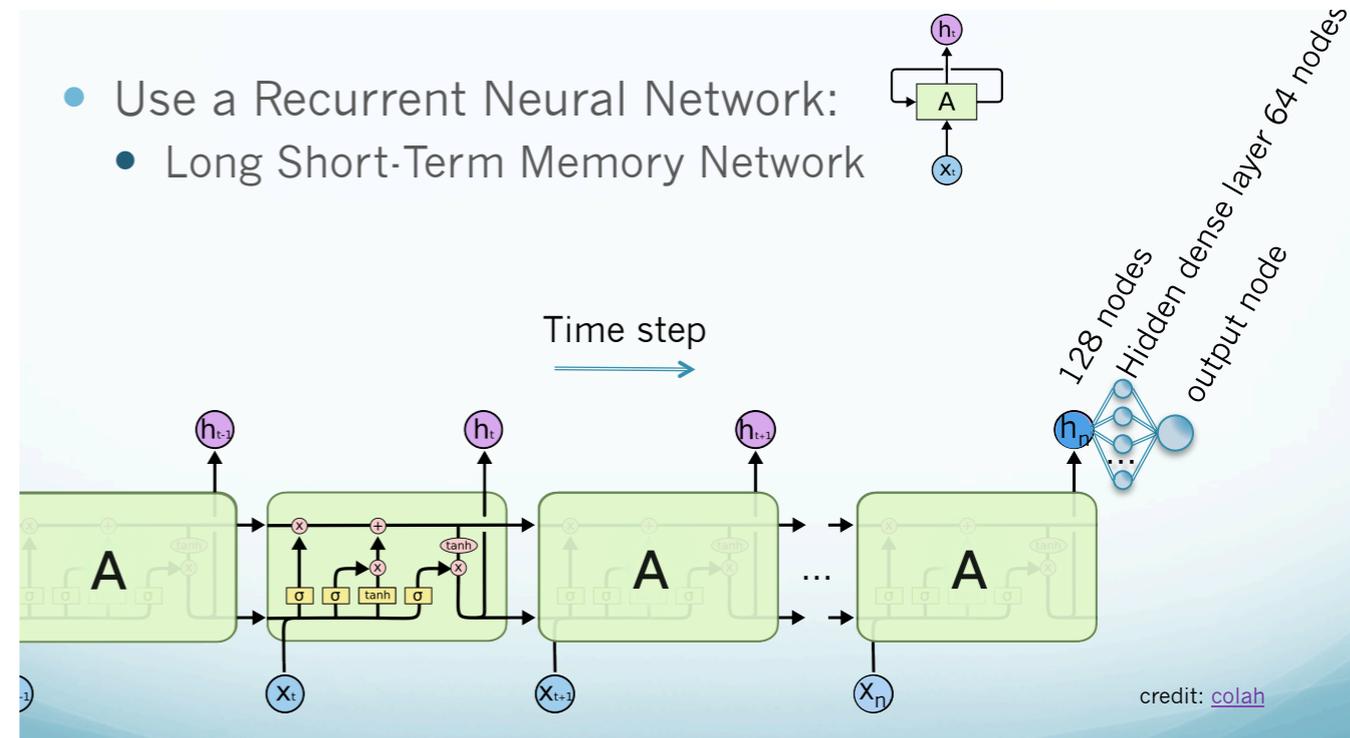
Egan et al 1711.09059

- View anti- $k_T$  sequence as a binary tree
- Order using depth-first traversal prioritizing jets with 'parents' whose  $d_{ij}$  is smaller



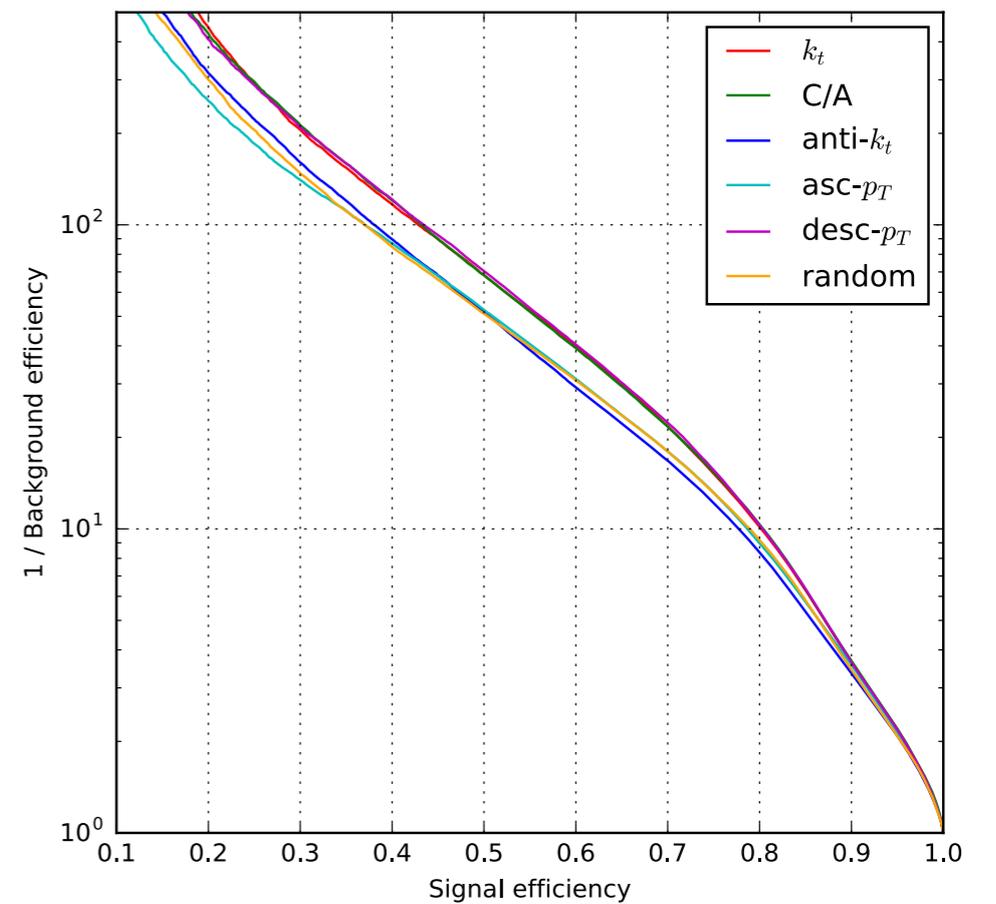
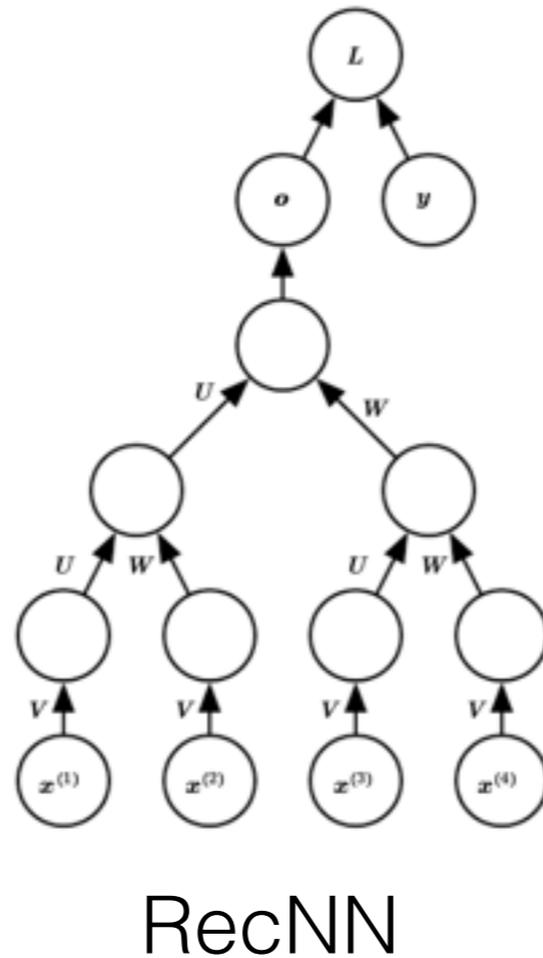
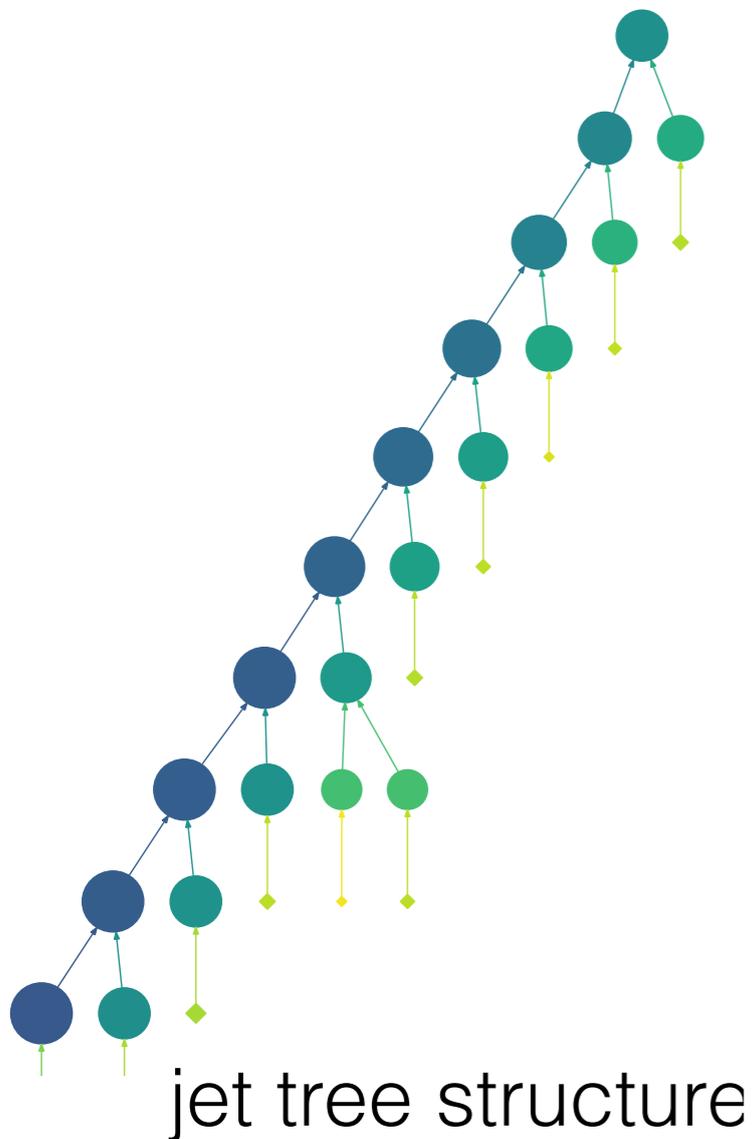
Turn jet into a sequence, e.g. using jet clustering history.

Feed sequence to RNN/LSTM, etc.



# Trees

Louppe et al I702.00748, Cheng I711.02633



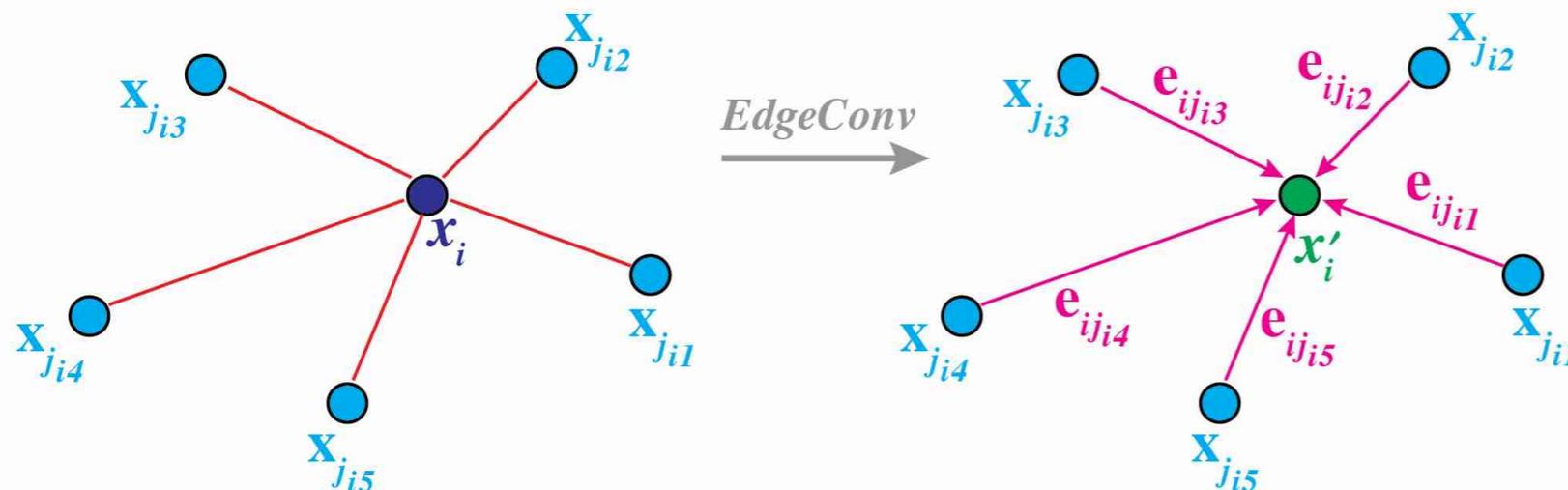
Use jet clustering history to build binary tree

Train recursive neural net on jet tree to learn “embedding” for classifier.

# Graphs / Point clouds

Hu & Gouskos 1902.08570

- Particle cloud
  - particles are intrinsically unordered
  - primary information:
    - 2D coordinates in the  $\eta$ - $\phi$  space
  - but also additional “features”:
    - energy/momenta
    - charge/particle type
    - track quality/impact parameters/etc.



simulated top quark jet  
anti- $k_T$ ,  $R = 0.8$ ,  $p_T = 600$  GeV

