

# Machine learning

Yann Coadou

CPPM Marseille

esipap...

European School of Instrumentation  
in Particle & Astroparticle Physics

Online edition, 27 January 2021





European School of Instrumentation  
in Particle & Astroparticle Physics

Apply now!

18 JANUARY  
12 MARCH  
2021

Remote learning

• Course 1  
Physics of Particle and  
Astroparticle Detectors  
18 Jan - 12 Feb

• Course 2  
Advanced Lectures  
on Particle Detectors  
and Applications  
15 Feb - 12 Mar

Deadline 27 November 2020



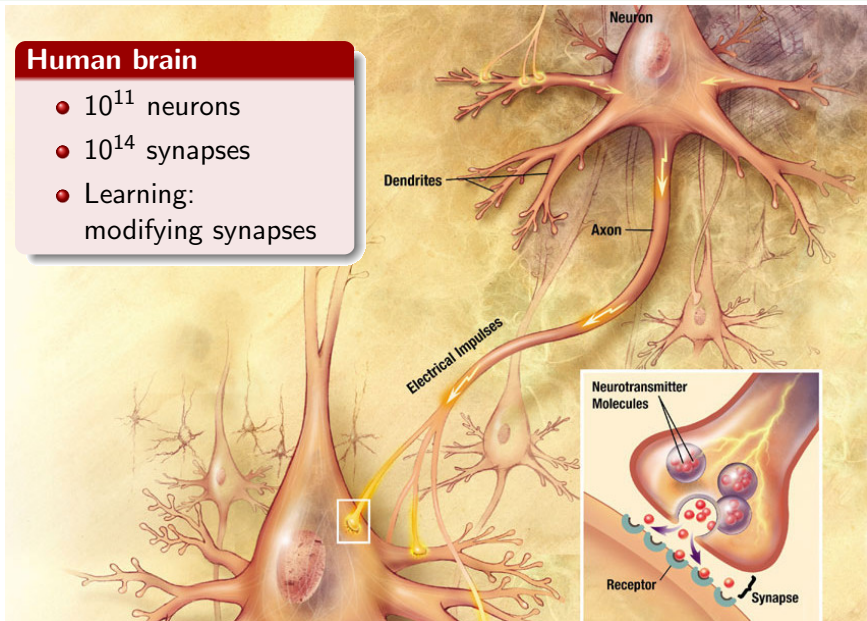
ON-LINE APPLICATIONS | [www.esipap.eu](http://www.esipap.eu)

- 1 Introduction
- 2 Optimal discrimination
- 3 Machine learning
- 4 Random grid search
- 5 Genetic algorithms
- 6 Quadratic and linear discriminants
- 7 Support vector machines
- 8 Kernel density estimation
- 9 (Boosted) Decision trees
- 10 Neural networks
- 11 Deep neural networks
- 12 Machine learning and particle physics
- 13 Conclusion
- 14 References



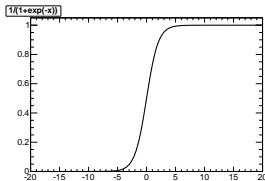
## Human brain

- $10^{11}$  neurons
- $10^{14}$  synapses
- Learning: modifying synapses

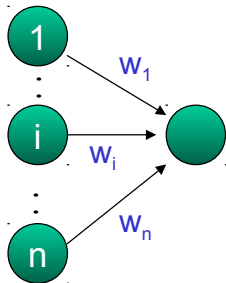


- 1943: W. McCulloch and W. Pitts explore capabilities of networks of simple neurons
- 1958: F. Rosenblatt introduces perceptron (single neuron with adjustable weights and threshold activation function)
- 1969: M. Minsky and S. Papert prove limitations of perceptron (linear separation only) and (wrongly) conjecture that multi-layered perceptrons have same limitations  
⇒ ANN research almost abandoned in 1970s!!!
- 1986: Rumelhart, Hinton and Williams introduce “backward propagation of errors”: solves (partially) multi-layered learning
- Next: focus on multilayer perceptron (MLP)

- Remember linear separation (Fisher discriminant):  
 $\lambda(x) = w \cdot x = \sum_{i=1}^n w_i x_i + w_0$
- Boundary at  $\lambda(x) = 0$
- Replace threshold boundary by sigmoid (or tanh):



$$\lambda \rightarrow \sigma(\lambda) = \frac{1}{1 + e^{-\lambda}}$$



- $\sigma$ : activation function (neuron activity)
- Neuron behaviour completely controlled by weights  $w = \{w_0, \dots, w_n\}$
- Training: minimisation of error/loss function (quadratic deviations, entropy [maximum likelihood]), via gradient descent or stochastic approximation

## Universal approximation theorem

Let  $\sigma(\cdot)$  be a non-constant, bounded, and monotone-increasing continuous function. Let  $\mathcal{C}(I_n)$  denote the space of continuous functions on the  $n$ -dimensional hypercube. Then, for any given function  $f \in \mathcal{C}(I_n)$  and  $\varepsilon > 0$  there exists an integer  $M$  and sets of real constants  $w_j, w_{ij}$  where  $i = 1, \dots, n$  and  $j = 1, \dots, M$  such that

$$y(x, w) = \sum_{j=1}^M w_j \sigma \left( \sum_{i=1}^n w_{ij} x_i + w_{0j} \right)$$

is an approximation of  $f(\cdot)$ , that is  $|y(x) - f(x)| < \varepsilon$ .

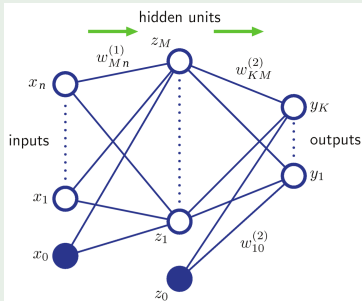
## Interpretation

- You can approximate any continuous function to arbitrary precision with a linear combination of sigmoids
- Corollary 1: can approximate any continuous function with neurons!
- Corollary 2: a single hidden layer is enough
- Corollary 3: a linear output neuron is enough

## Multilayer perceptron: feedforward network

- Neurons organised in layers
- Output of one layer becomes input to next layer

$$y_k(x, w) = \sum_{j=0}^M w_{kj}^{(2)} \underbrace{\sigma \left( \sum_{i=0}^n w_{ji}^{(1)} x_i \right)}_{z_j}$$



# A neural network can fit any function: examples

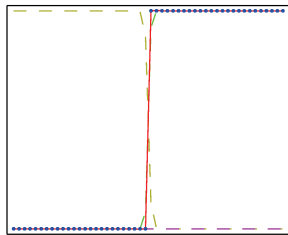
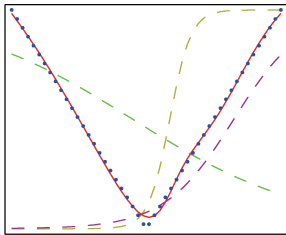
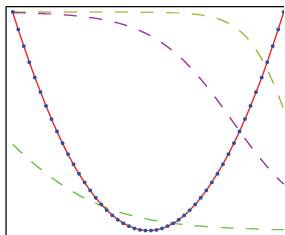
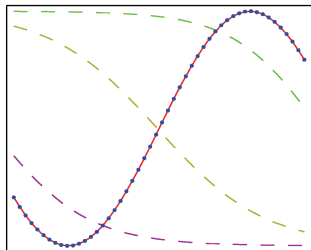
- 1 input (training data), 1 output
- 3 hidden neurons on one hidden layer

© Jan Therhaag

---  $z_1$   
---  $z_2$   
---  $z_3$

— *output*

..... *training data*

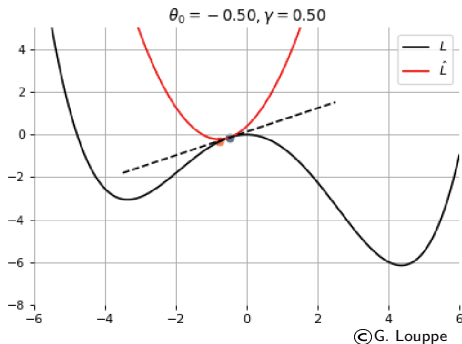


- Training means minimising error function  $E(w)$
- $\frac{\partial E}{\partial w_j} = \sum_{n=1}^N -(t^{(n)} - y^{(n)})x_j^{(n)}$  with target  $t^{(n)}$  (0 or 1), so  $t^{(n)} - y^{(n)}$  is the error on event  $n$
- All events at once (batch learning):
  - weights updated all at once after processing the entire training sample
  - finds the actual steepest descent
  - takes more time
  - usually: mini-batches (send events by batches)
  - new training events: need to restart training from scratch
- or one-by-one (online learning):
  - incremental learning: new training events included as they come
  - speeds up learning
  - may avoid local minima with stochastic component in minimisation
  - careful: depends on the order of training events
- One epoch: going through the entire training data once

- Minimise error function  $E(w)$
- Gradient descent:  $w^{(k+1)} = w^{(k)} - \eta \nabla_w E^{(k)}$   
with learning rate  $\eta$

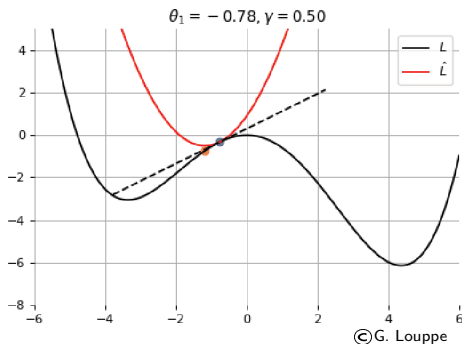


- Minimise error function  $E(w)$
- Gradient descent:  $w^{(k+1)} = w^{(k)} - \eta \nabla_w E^{(k)}$   
with learning rate  $\eta$



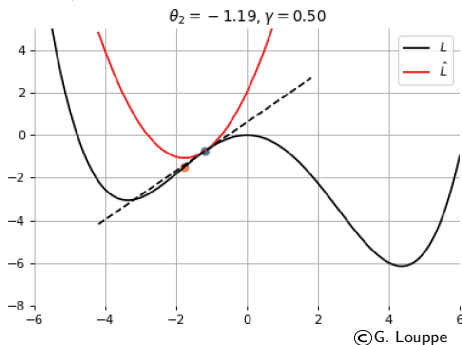
local minimum

- Minimise error function  $E(w)$
- Gradient descent:  $w^{(k+1)} = w^{(k)} - \eta \nabla_w E^{(k)}$   
with learning rate  $\eta$

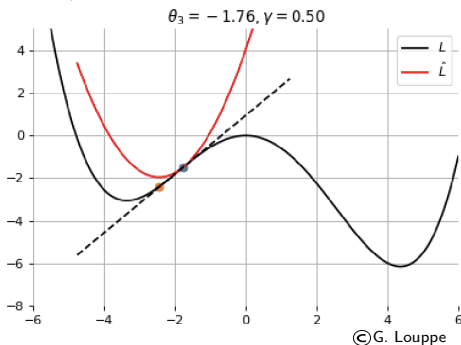


local minimum

- Minimise error function  $E(w)$
- Gradient descent:  $w^{(k+1)} = w^{(k)} - \eta \nabla_w E^{(k)}$   
with learning rate  $\eta$

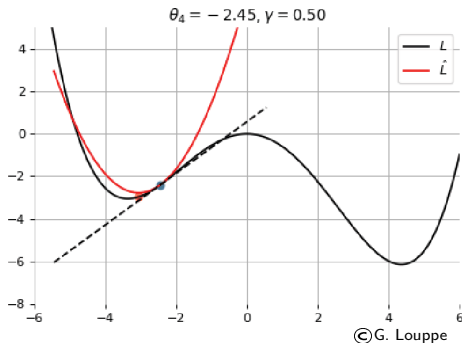


- Minimise error function  $E(w)$
- Gradient descent:  $w^{(k+1)} = w^{(k)} - \eta \nabla_w E^{(k)}$   
with learning rate  $\eta$



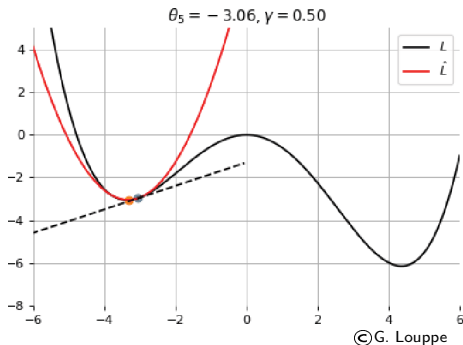
local minimum

- Minimise error function  $E(w)$
- Gradient descent:  $w^{(k+1)} = w^{(k)} - \eta \nabla_w E^{(k)}$   
with learning rate  $\eta$



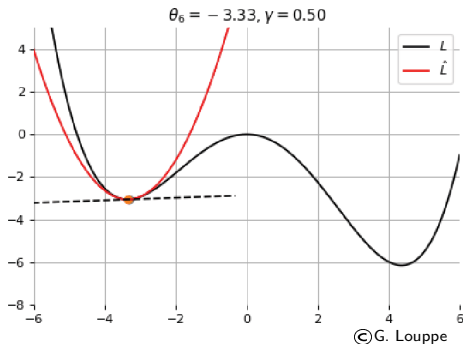
local minimum

- Minimise error function  $E(w)$
- Gradient descent:  $w^{(k+1)} = w^{(k)} - \eta \nabla_w E^{(k)}$   
with learning rate  $\eta$



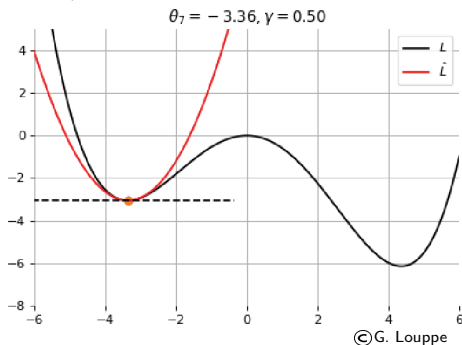
local minimum

- Minimise error function  $E(w)$
- Gradient descent:  $w^{(k+1)} = w^{(k)} - \eta \nabla_w E^{(k)}$   
with learning rate  $\eta$



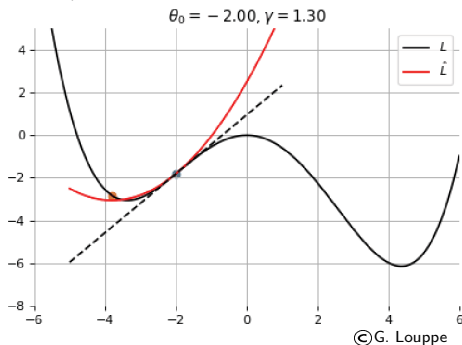
local minimum

- Minimise error function  $E(w)$
- Gradient descent:  $w^{(k+1)} = w^{(k)} - \eta \nabla_w E^{(k)}$   
with learning rate  $\eta$



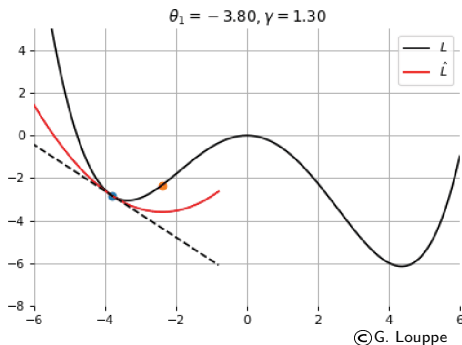


- Minimise error function  $E(w)$
- Gradient descent:  $w^{(k+1)} = w^{(k)} - \eta \nabla_w E^{(k)}$   
with learning rate  $\eta$



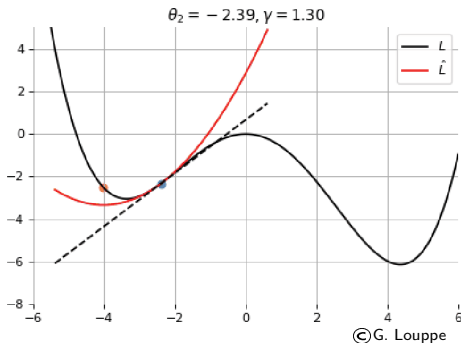
too large learning rate

- Minimise error function  $E(w)$
- Gradient descent:  $w^{(k+1)} = w^{(k)} - \eta \nabla_w E^{(k)}$   
with learning rate  $\eta$



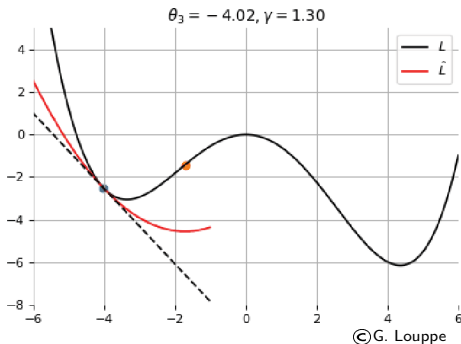
too large learning rate

- Minimise error function  $E(w)$
- Gradient descent:  $w^{(k+1)} = w^{(k)} - \eta \nabla_w E^{(k)}$   
with learning rate  $\eta$



too large learning rate

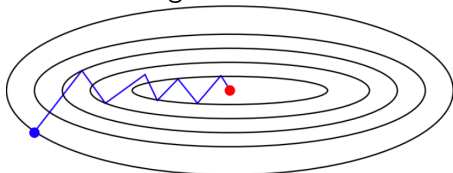
- Minimise error function  $E(w)$
- Gradient descent:  $w^{(k+1)} = w^{(k)} - \eta \nabla_w E^{(k)}$   
with learning rate  $\eta$



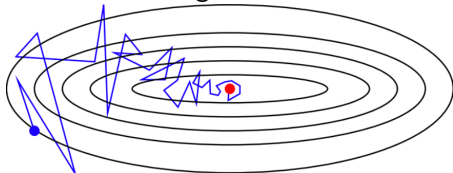
too large learning rate

- Solution: stochastic gradient descent (SGD)

batch gradient descent



stochastic gradient descent



©G. Louppe

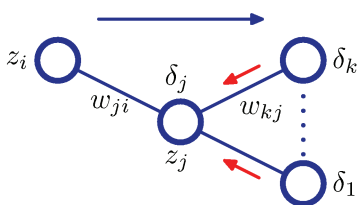
- Training means minimising error function  $E(w)$
- For single neuron:  $\frac{dE}{dw_k} = (y - t)x_k$
- One can show that for a network:

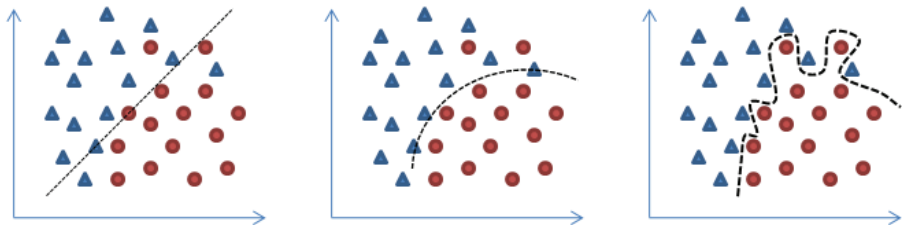
$$\frac{dE}{dw_{ji}} = \delta_j z_i, \text{ where}$$

$\delta_k = (y_k - t_k)$  for output neurons

$$\delta_j \propto \sum_k w_{kj} \delta_k \text{ otherwise}$$

- Hence errors are propagated backwards

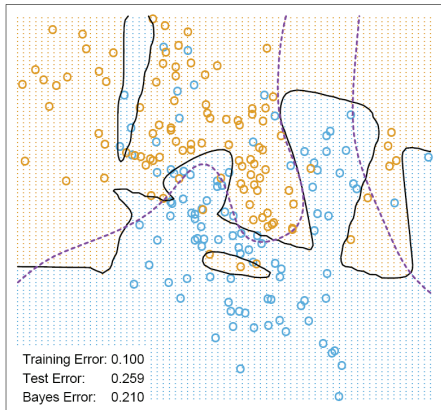




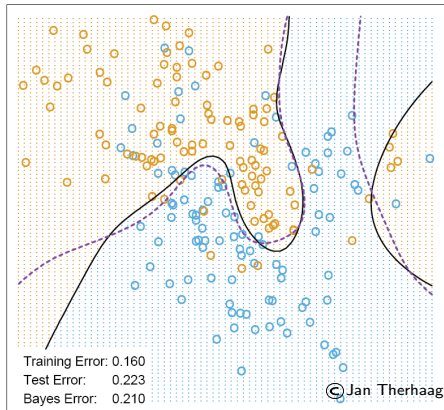
- Diverging weights can cause overfitting
- Mitigate by:
  - early stopping (after a fixed number of epochs)
  - monitoring error on test sample
  - regularisation, introducing a “weight decay” term to penalise large weights, preventing overfitting:

$$\tilde{E}(w) = E(w) + \frac{\alpha}{2} \sum_i w_i^2$$

10 hidden nodes



10 hidden nodes and  $\alpha = 0.04$

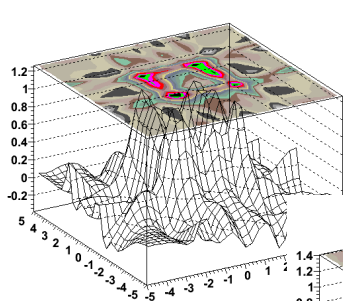


- Much less overfitting, better generalisation properties

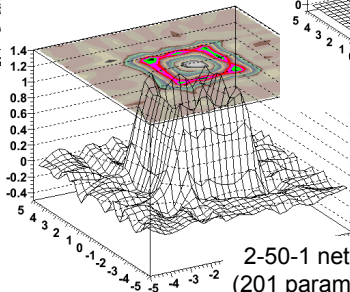
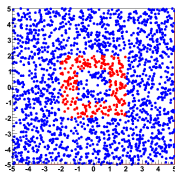


- Preprocess data:
  - if relevant, provide e.g.  $x/y$  instead of  $x$  and  $y$
  - subtract the mean because the sigmoid derivative becomes negligible very fast (so, input mean close to 0)
  - normalise variances (close to 1)
  - shuffle training sample (order matters in online training)
- Initial random weights should be small to avoid saturation
- Batch/online training: depends on the problem
- Regularise weights to minimise overtraining.
- Make sure the training sample covers the full parameter space
- No rule (not even guestimates) about the number of hidden nodes (unless using constructive algorithm, adding resources as needed)
- A single hidden layer is enough for all purposes, but multiple hidden layers may allow for a solution with fewer parameters

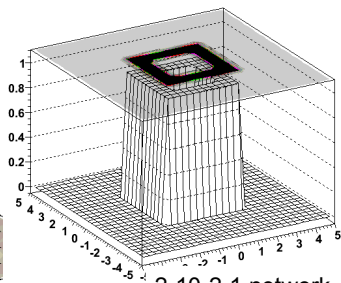
# Adding a hidden layer



2-20-1 network  
(81 parameters)



2-50-1 network  
(201 parameters)



2-10-2-1 network  
(55 parameters)

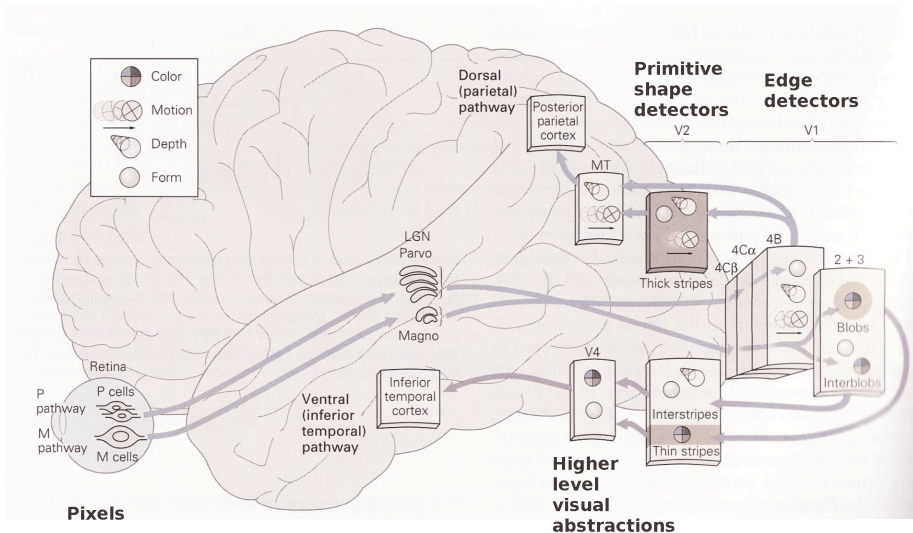
## What is learning?

- Ability to learn underlying and previously unknown structure from examples  
⇒ capture variations
- Deep learning: have several hidden layers ( $> 2$ ) in a neural network

## Motivation for deep learning

- Inspired by the brain!
- Humans organise ideas hierarchically, through composition of simpler ideas
- Heavily unsupervised training, learning simpler tasks first, then combining into more abstract ones
- Learn first order features from raw inputs, then patterns in first order features, then etc.

# Deep architecture in the brain



## Mimicking the brain

- About 1% of neurons active simultaneously in the brain:  
**distributed representation**
  - activation of small subset of features, not mutually exclusive
  - more efficient than local representation
  - distributed representations necessary to achieve non-local generalization, exponentially more efficient than 1-of- $N$  enumeration
  - example: integers in  $1..N$ 
    - local representation: vector of  $N$  bits with single 1 and  $N-1$  zeros
    - distributed representation: vector of  $\log_2 N$  bits (binary notation), exponentially more compact
- Meaning: information not localised in particular neuron but distributed across them

## Deep architecture

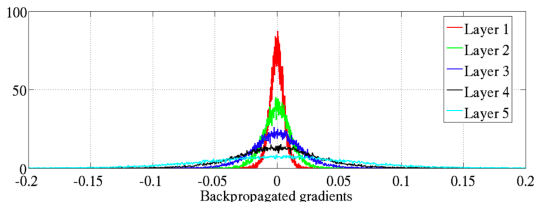
- Insufficient depth can hurt
- Learn basic features first, then higher level ones
- Learn good intermediate representations, shared across tasks

## Deep networks were unattractive

- One layer theoretically enough for everything
- Used to perform worse than shallow networks with 1 or 2 hidden layers
- Apparently difficult/impossible to train (using random initial weights and supervised learning with backpropagation)
- Backpropagation issues:
  - requires labelled data (usually scarce and expensive)
  - does not scale well, getting stuck in local minima
  - “vanishing gradient”: gradients getting very small further away from output  $\Rightarrow$  early layers do not learn much, can even penalise overall performance

## Deep networks were unattractive

- One layer theoretically enough for everything
- Used to perform worse than shallow networks with 1 or 2 hidden layers
- Apparently difficult/impossible to train (using random initial weights and supervised learning with backpropagation)
- Backpropagation issues:
  - requires labelled data (usually scarce and expensive)
  - does not scale well, getting stuck in local minima
  - “vanishing gradient”: gradients getting very small further away from output  $\Rightarrow$  early layers do not learn much, can even penalise overall performance



## Deep networks were unattractive

- One layer theoretically enough for everything
- Used to perform worse than shallow networks with 1 or 2 hidden layers
- Apparently difficult/impossible to train (using random initial weights and supervised learning with backpropagation)
- Backpropagation issues:
  - requires labelled data (usually scarce and expensive)
  - does not scale well, getting stuck in local minima
  - “vanishing gradient”: gradients getting very small further away from output  $\Rightarrow$  early layers do not learn much, can even penalise overall performance

## Breakthroughs around 2006 (Bengio, Hinton, LeCun)

- Train each layer independently
- Can use unlabelled data (a lot of it)
- New activation functions
- Possible thanks to algorithmic innovations, computing resources, data!



## Algorithm

- Take input information
- Train feature extractor
- Use output as input to training another feature extractor
- Keep adding layers, train each layer separately
- Finalise with a supervised classifier, taking last feature extractor output as input
- All steps above: pre-training
- Fine-tune the whole thing with supervised training (backpropagation)
  - initial weights are those from pre-training

## Feature extractors

- Restricted Boltzmann machine (RBM), auto-encoder, sparse auto-encoder, denoising auto-encoder, etc.
- Note: important to not use linear activation functions in hidden layers. Combination of linear functions still linear, so equivalent to single hidden layer

# Why does unsupervised training work?

## Optimisation hypothesis

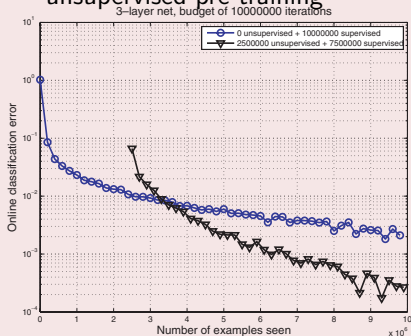
- Training one layer at a time scales well
- Backpropagation from sensible features
- Better local minimum than random initialisation, local search around it

## Overfitting/regularisation hypothesis

- More info in inputs than labels
- No need for final discriminant to discover features
- Fine-tuning only at category boundaries

## Example

- Stacked denoising auto-encoders
- 10 million handwritten digits
- First 2.5 million used for unsupervised pre-training



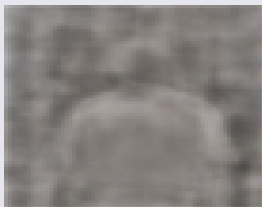
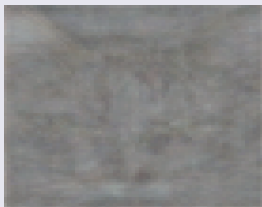
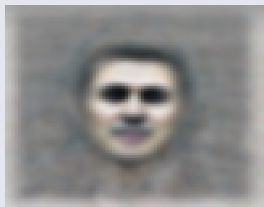
- Worse with supervision: eliminates projections of data not useful for local cost but helpful for deep model cost

## A “giant” neural network

- At Google they trained a 9-layered NN with 1 billion connections
  - trained on 10 million  $200 \times 200$  pixel images from YouTube videos
  - on 1000 machines (16000 cores) for 3 days, unsupervised learning
- Sounds big? The human brain has 100 billion ( $10^{11}$ ) neurons and 100 trillion ( $10^{14}$ ) connections...

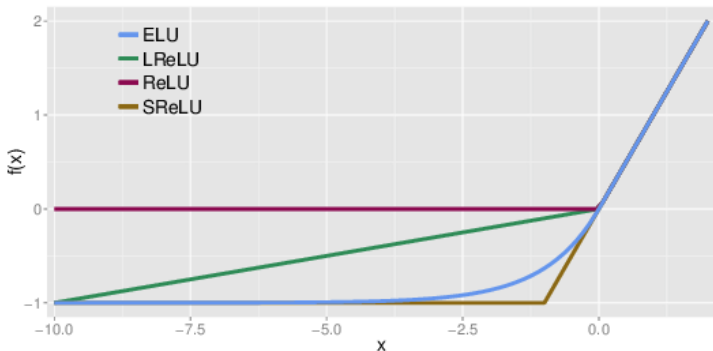
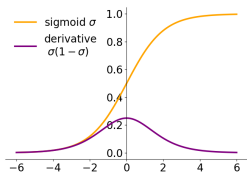
## What it did

- It learned to recognise faces, one of the original goals
- ... but also cat faces (among the most popular things in YouTube videos) and body shapes

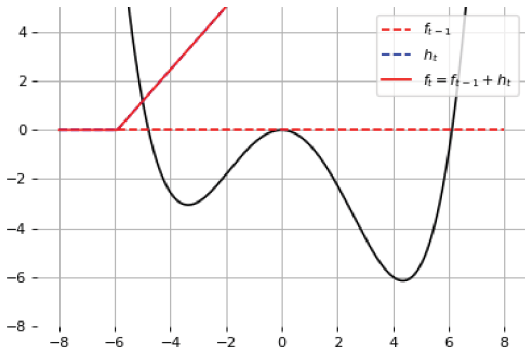
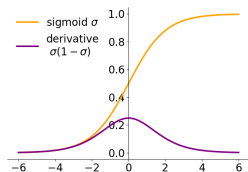




- One of reasons for vanishing gradient: sigmoid activation
  - tiny non-varying derivative away from zero
- Solution: non-saturating function
- Simplest case: rectified linear unit ReLU
- Other variants: leaky ReLU, shifted ReLU (SReLU), exponential linear unit (ELU), etc.

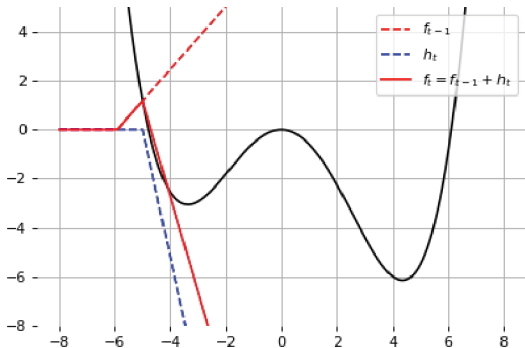
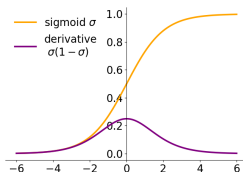


- One of reasons for vanishing gradient: sigmoid activation
  - tiny non-varying derivative away from zero
- Solution: non-saturating function
- Simplest case: rectified linear unit ReLU
- Other variants: leaky ReLU, shifted ReLU (SReLU), exponential linear unit (ELU), etc.



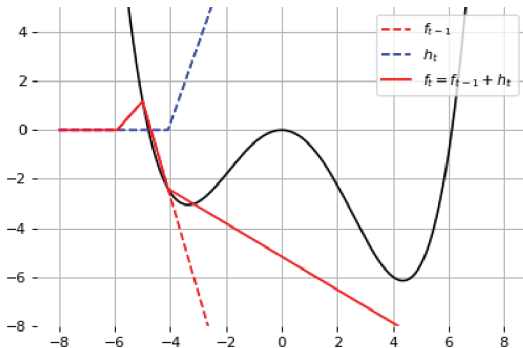
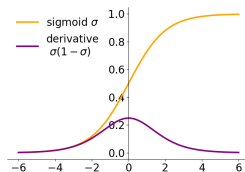
©G. Louppe

- One of reasons for vanishing gradient: sigmoid activation
  - tiny non-varying derivative away from zero
- Solution: non-saturating function
- Simplest case: rectified linear unit ReLU
- Other variants: leaky ReLU, shifted ReLU (SReLU), exponential linear unit (ELU), etc.



©G. Louppe

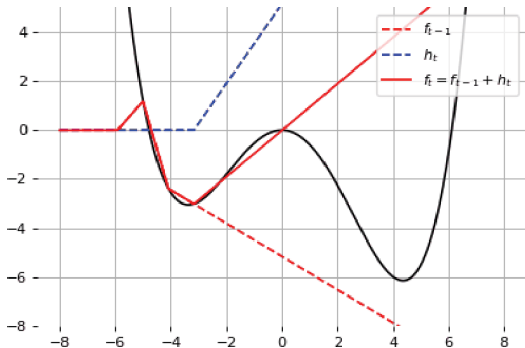
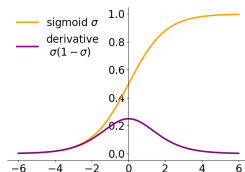
- One of reasons for vanishing gradient: sigmoid activation
  - tiny non-varying derivative away from zero
- Solution: non-saturating function
- Simplest case: rectified linear unit ReLU
- Other variants: leaky ReLU, shifted ReLU (SReLU), exponential linear unit (ELU), etc.



©G. Louppe

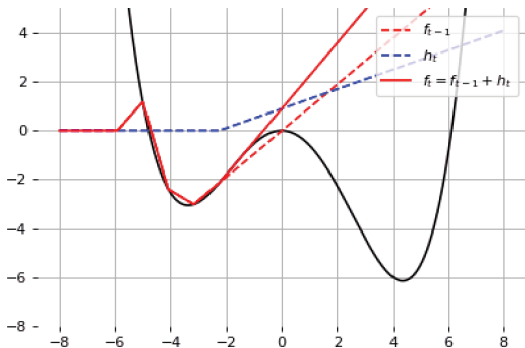
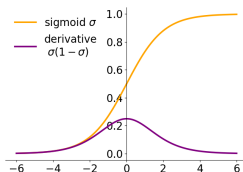


- One of reasons for vanishing gradient: sigmoid activation
  - tiny non-varying derivative away from zero
- Solution: non-saturating function
- Simplest case: rectified linear unit ReLU
- Other variants: leaky ReLU, shifted ReLU (SReLU), exponential linear unit (ELU), etc.



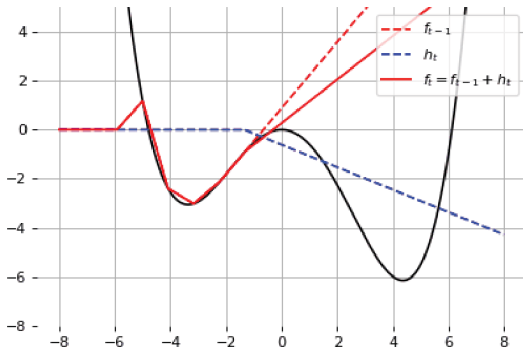
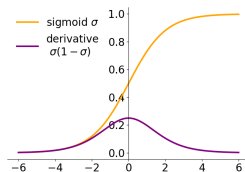
©G. Louppe

- One of reasons for vanishing gradient: sigmoid activation
  - tiny non-varying derivative away from zero
- Solution: non-saturating function
- Simplest case: rectified linear unit ReLU
- Other variants: leaky ReLU, shifted ReLU (SReLU), exponential linear unit (ELU), etc.



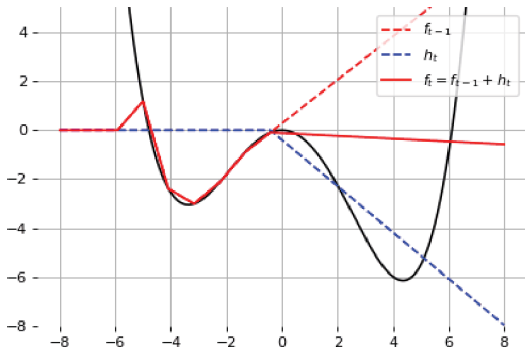
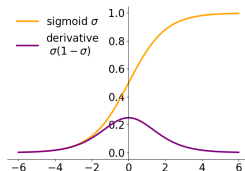
©G. Louppe

- One of reasons for vanishing gradient: sigmoid activation
  - tiny non-varying derivative away from zero
- Solution: non-saturating function
- Simplest case: rectified linear unit ReLU
- Other variants: leaky ReLU, shifted ReLU (SReLU), exponential linear unit (ELU), etc.



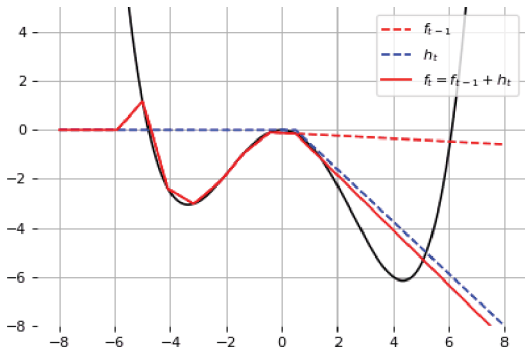
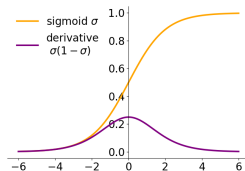
©G. Louppe

- One of reasons for vanishing gradient: sigmoid activation
  - tiny non-varying derivative away from zero
- Solution: non-saturating function
- Simplest case: rectified linear unit ReLU
- Other variants: leaky ReLU, shifted ReLU (SReLU), exponential linear unit (ELU), etc.



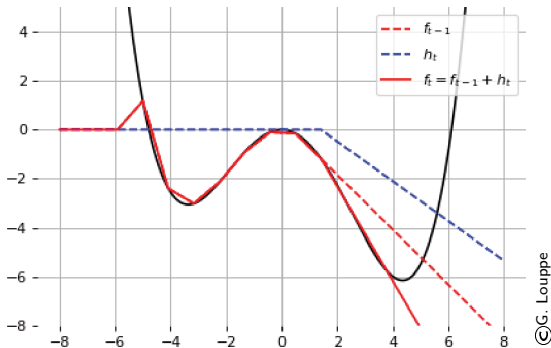
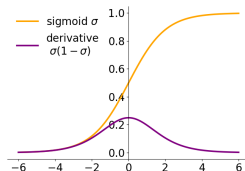
©G. Louppe

- One of reasons for vanishing gradient: sigmoid activation
  - tiny non-varying derivative away from zero
- Solution: non-saturating function
- Simplest case: rectified linear unit ReLU
- Other variants: leaky ReLU, shifted ReLU (SReLU), exponential linear unit (ELU), etc.

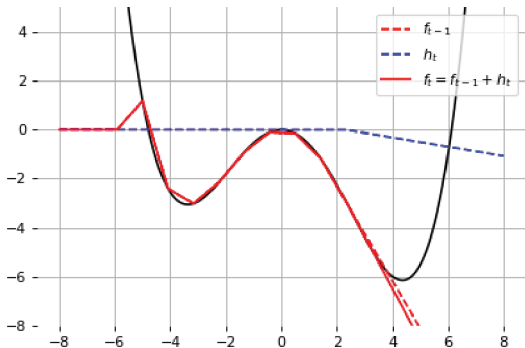
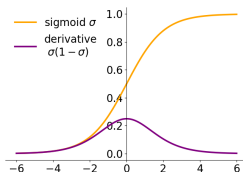


©G. Louppe

- One of reasons for vanishing gradient: sigmoid activation
  - tiny non-varying derivative away from zero
- Solution: non-saturating function
- Simplest case: rectified linear unit ReLU
- Other variants: leaky ReLU, shifted ReLU (SReLU), exponential linear unit (ELU), etc.

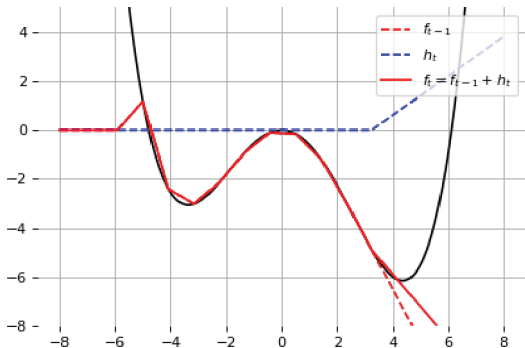
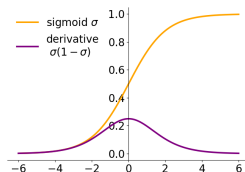


- One of reasons for vanishing gradient: sigmoid activation
  - tiny non-varying derivative away from zero
- Solution: non-saturating function
- Simplest case: rectified linear unit ReLU
- Other variants: leaky ReLU, shifted ReLU (SReLU), exponential linear unit (ELU), etc.



©G. Louppe

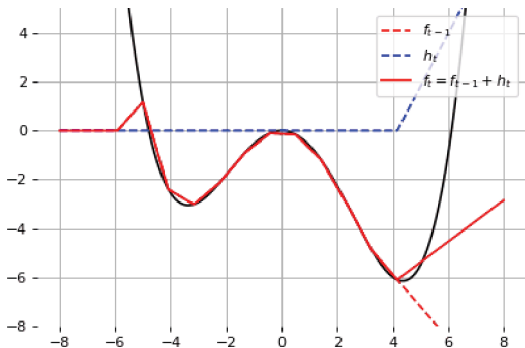
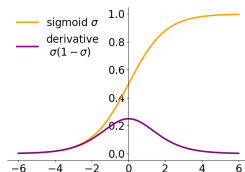
- One of reasons for vanishing gradient: sigmoid activation
  - tiny non-varying derivative away from zero
- Solution: non-saturating function
- Simplest case: rectified linear unit ReLU
- Other variants: leaky ReLU, shifted ReLU (SReLU), exponential linear unit (ELU), etc.



©G. Louppe

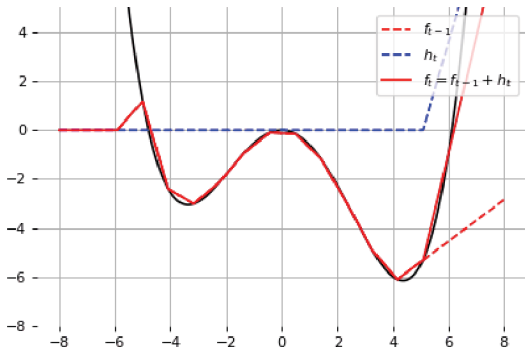
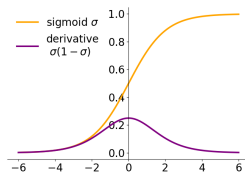


- One of reasons for vanishing gradient: sigmoid activation
  - tiny non-varying derivative away from zero
- Solution: non-saturating function
- Simplest case: rectified linear unit ReLU
- Other variants: leaky ReLU, shifted ReLU (SReLU), exponential linear unit (ELU), etc.

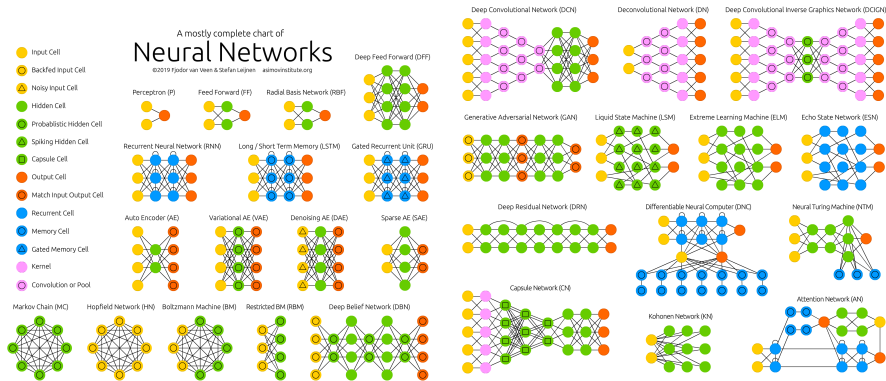


©G. Louppe

- One of reasons for vanishing gradient: sigmoid activation
  - tiny non-varying derivative away from zero
- Solution: non-saturating function
- Simplest case: rectified linear unit ReLU
- Other variants: leaky ReLU, shifted ReLU (SReLU), exponential linear unit (ELU), etc.



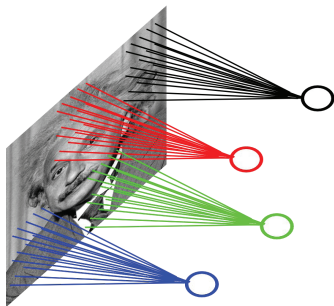
©G. Louppe



► <https://www.asimovinstitute.org/>

- Many possible network structures
- Moving away from feature engineering to model design

- Images are stationary: can learn feature in one part and apply it in another



- Images are stationary: can learn feature in one part and apply it in another
- Use e.g. small patch sampled randomly, learn feature, convolve with full image

1 <sub>x1</sub>	1 <sub>x0</sub>	1 <sub>x1</sub>	0	0
0 <sub>x0</sub>	1 <sub>x1</sub>	1 <sub>x0</sub>	1	0
0 <sub>x1</sub>	0 <sub>x0</sub>	1 <sub>x1</sub>	1	1
0	0	1	1	0
0	1	1	0	0

Image

4		

Convolved  
Feature

- Images are stationary: can learn feature in one part and apply it in another
- Use e.g. small patch sampled randomly, learn feature, convolve with full image

1	1 <sub>x1</sub>	1 <sub>x0</sub>	0 <sub>x1</sub>	0
0	1 <sub>x0</sub>	1 <sub>x1</sub>	1 <sub>x0</sub>	0
0	0 <sub>x1</sub>	1 <sub>x0</sub>	1 <sub>x1</sub>	1
0	0	1	1	0
0	1	1	0	0

Image

4	3	

Convolved  
Feature

- Images are stationary: can learn feature in one part and apply it in another
- Use e.g. small patch sampled randomly, learn feature, convolve with full image

1	1	1 <sub>x1</sub>	0 <sub>x0</sub>	0 <sub>x1</sub>
0	1	1 <sub>x0</sub>	1 <sub>x1</sub>	0 <sub>x0</sub>
0	0	1 <sub>x1</sub>	1 <sub>x0</sub>	1 <sub>x1</sub>
0	0	1	1	0
0	1	1	0	0

Image

4	3	4

Convolved  
Feature

- Images are stationary: can learn feature in one part and apply it in another
- Use e.g. small patch sampled randomly, learn feature, convolve with full image

1	1	1	0	0
0 <sub>x1</sub>	1 <sub>x0</sub>	1 <sub>x1</sub>	1	0
0 <sub>x0</sub>	0 <sub>x1</sub>	1 <sub>x0</sub>	1	1
0 <sub>x1</sub>	0 <sub>x0</sub>	1 <sub>x1</sub>	1	0
0	1	1	0	0

Image

4	3	4
2		

Convolved  
Feature



- Images are stationary: can learn feature in one part and apply it in another
- Use e.g. small patch sampled randomly, learn feature, convolve with full image

1	1	1	0	0
0	1 <sub>x1</sub>	1 <sub>x0</sub>	1 <sub>x1</sub>	0
0	0 <sub>x0</sub>	1 <sub>x1</sub>	1 <sub>x0</sub>	1
0	0 <sub>x1</sub>	1 <sub>x0</sub>	1 <sub>x1</sub>	0
0	1	1	0	0

Image

4	3	4
2	4	

Convolved  
Feature

- Images are stationary: can learn feature in one part and apply it in another
- Use e.g. small patch sampled randomly, learn feature, convolve with full image

1	1	1	0	0
0	1	1 <sub>x1</sub>	1 <sub>x0</sub>	0 <sub>x1</sub>
0	0	1 <sub>x0</sub>	1 <sub>x1</sub>	1 <sub>x0</sub>
0	0	1 <sub>x1</sub>	1 <sub>x0</sub>	0 <sub>x1</sub>
0	1	1	0	0

Image

4	3	4
2	4	3

Convolved  
Feature

- Images are stationary: can learn feature in one part and apply it in another
- Use e.g. small patch sampled randomly, learn feature, convolve with full image

1	1	1	0	0
0	1	1	1	0
0 <sub>x1</sub>	0 <sub>x0</sub>	1 <sub>x1</sub>	1	1
0 <sub>x0</sub>	0 <sub>x1</sub>	1 <sub>x0</sub>	1	0
0 <sub>x1</sub>	1 <sub>x0</sub>	1 <sub>x1</sub>	0	0

Image

4	3	4
2	4	3
2		

Convolved  
Feature

- Images are stationary: can learn feature in one part and apply it in another
- Use e.g. small patch sampled randomly, learn feature, convolve with full image

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

Image

4	3	4
2	4	3
2	3	

Convolved  
Feature

- Images are stationary: can learn feature in one part and apply it in another
- Use e.g. small patch sampled randomly, learn feature, convolve with full image

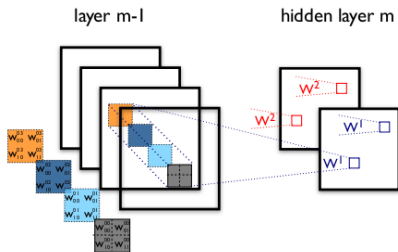
1	1	1	0	0
0	1	1	1	0
0	0	1 <sub>x1</sub>	1 <sub>x0</sub>	1 <sub>x1</sub>
0	0	1 <sub>x0</sub>	1 <sub>x1</sub>	0 <sub>x0</sub>
0	1	1 <sub>x1</sub>	0 <sub>x0</sub>	0 <sub>x1</sub>

Image

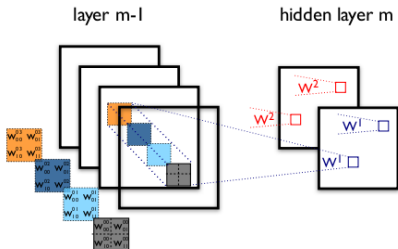
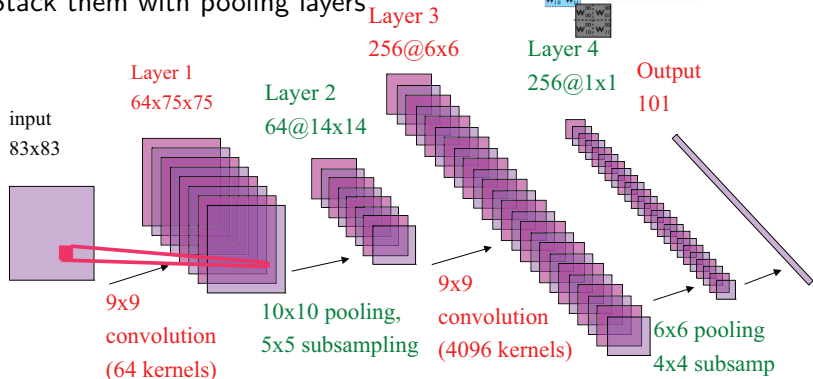
4	3	4
2	4	3
2	3	4

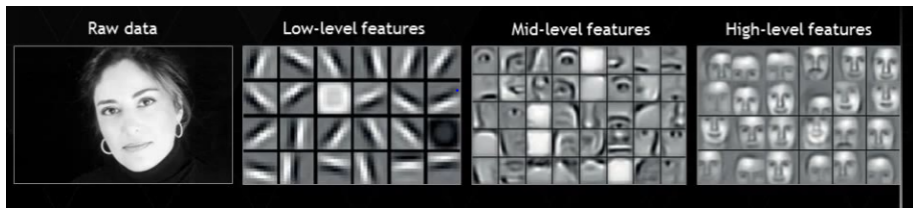
Convolved  
Feature

- Images are stationary: can learn feature in one part and apply it in another
- Use e.g. small patch sampled randomly, learn feature, convolve with full image
- Build several “feature maps”



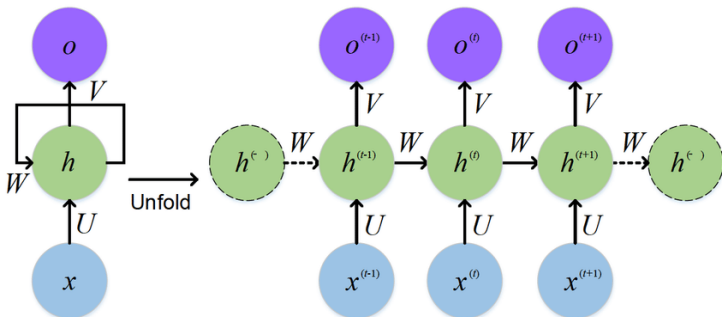
- Images are stationary: can learn feature in one part and apply it in another
- Use e.g. small patch sampled randomly, learn feature, convolve with full image
- Build several “feature maps”
- Stack them with pooling layers



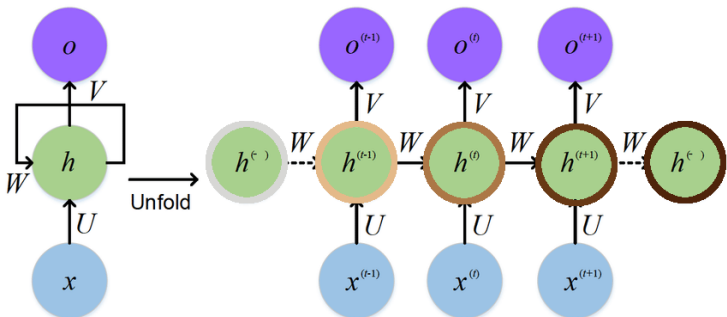




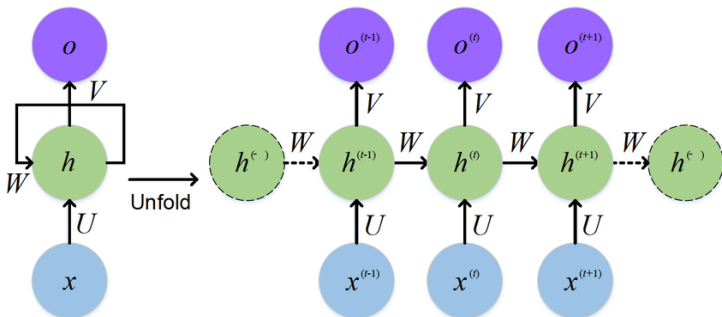
- Many problems require processing a sequence
  - sequence classification
    - text analysis (“sentiment analysis”)
    - DNA sequencing
    - action selection
  - sequence synthesis
    - text synthesis
    - music/video
  - sequence translation
    - speech recognition
    - translation
- Usually variable length sequences (number of words/ notes/ frames/ etc.)
- Use a recurrent model, maintaining a recurrent state updated after each step



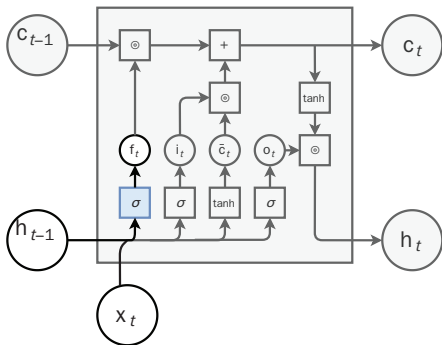
- Keeps information from earlier frames while processing (variable-size) sequence
- Could also be bi-directional, consuming sequence in both directions



- Keeps information from earlier frames while processing (variable-size) sequence
- Could also be bi-directional, consuming sequence in both directions
- Issue: early frames diluted over sequence  $\Rightarrow$  memory loss



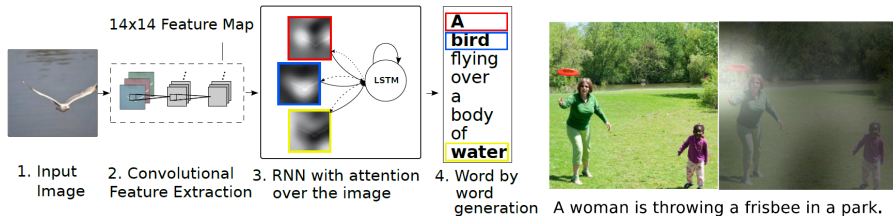
- Keeps information from earlier frames while processing (variable-size) sequence
- Could also be bi-directional, consuming sequence in both directions
- Issue: early frames diluted over sequence  $\Rightarrow$  memory loss
- Introducing long short-term memory (LSTM) networks
  - using forget gate to regulate information flow
  - also possible with gated recurrent units (GRU)



- Recurrent state split in two parts
  - cell state  $c_t$
  - output state  $h_t$
- Forget gate  $f_t$  to erase cell state info
- Input gate  $i_t$  to update cell state info
- Output gate  $o_t$  to select output state

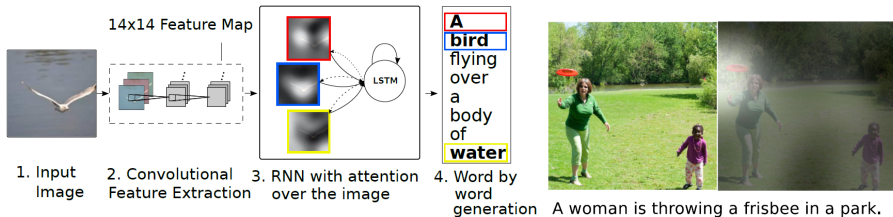
## Labelling images

▶ arXiv:1502.03044



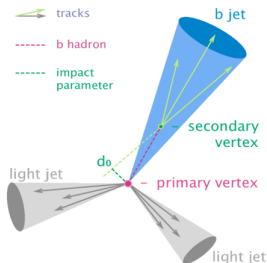
## Labelling images

arXiv:1502.03044



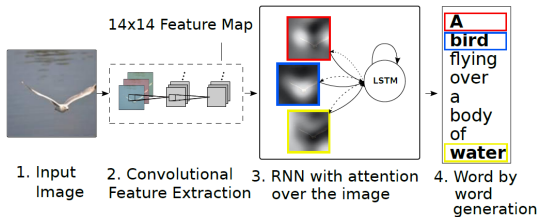
## *b*-jet tagging in ATLAS experiment

ATL-PHYS-PUB-2017-003



## Labelling images

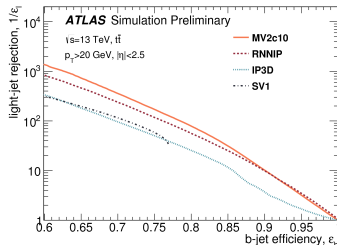
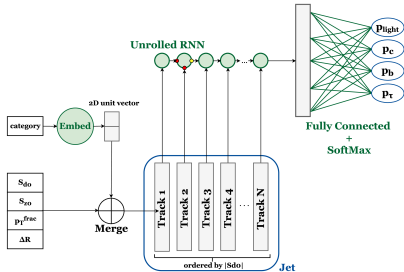
arXiv:1502.03044



A woman is throwing a frisbee in a park.

## b-jet tagging in ATLAS experiment

ATL-PHYS-PUB-2017-003



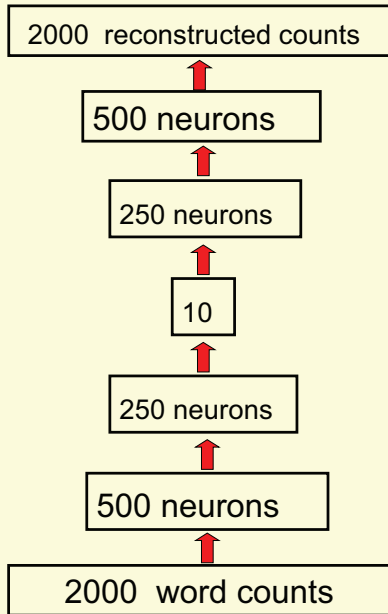


## Approximate the identity function

- Build a network whose output is similar to its input
- Sounds trivial? Except if imposing constraints on network (e.g., # of neurons, locally connected network) to discover interesting structures
- Can be viewed as lossy compression of input

## Finding similar books

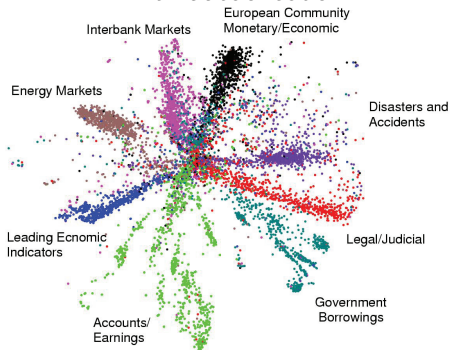
- Get count of 2000 most common words per book
- “Compress” to 10 numbers



With principle component analysis  
(PCA)



With autoencoder

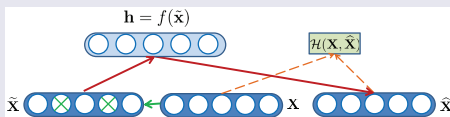


## Sparse auto-encoder

- Sparsity: try to have low activation of neurons (like in the brain)
- Compute average activation of each hidden unit over training set
- Add constraint to cost function to make average lower than some value close to 0

## Denoising auto-encoder

- Stochastically corrupt inputs
- Train to reconstruct uncorrupted input



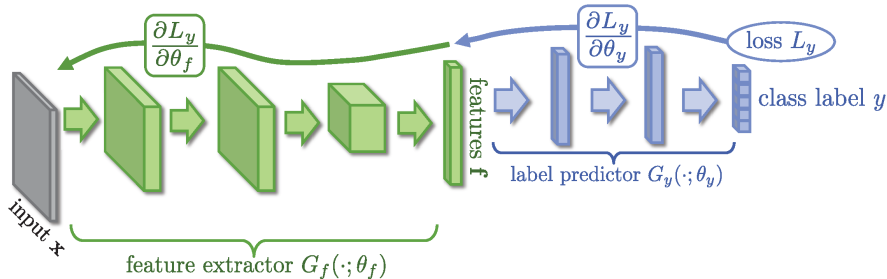
## Locally connected auto-encoder

- Allow hidden units to connect only to small subset of input units
- Useful with increasing number of input features (e.g., bigger image)
- Inspired by biology: visual system has localised receptive fields

- Typical training
  - signal and background from simulation
  - results compared to real data to make measurement
- Requires good data–simulation agreement

▶ arXiv:1409.7495

▶ arXiv:1505.07818

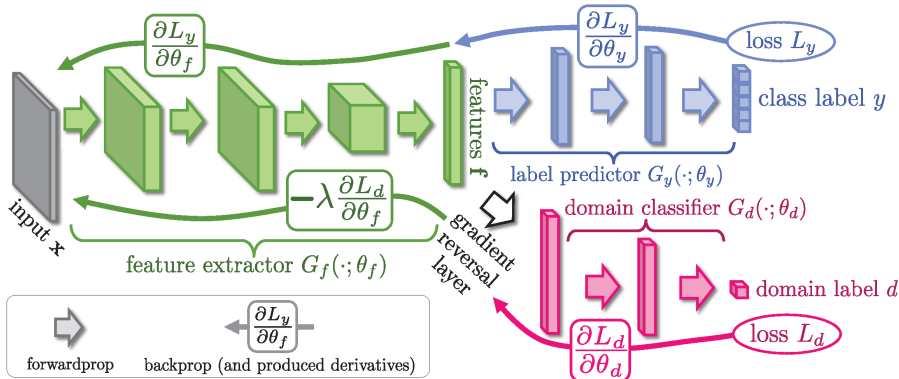


forwardprop      backprop (and produced derivatives)

- Typical training
  - signal and background from simulation
  - results compared to real data to make measurement
- Requires good data–simulation agreement
- Possibility to use adversarial training and domain adaptation to account for discrepancies/systematic uncertainties

▶ arXiv:1409.7495

▶ arXiv:1505.07818



## ImageNet Large Scale Visual Recognition Challenge

- ImageNet: database with 14 million images and 20k categories
- Used 1000 categories and about 1.3 million manually annotated images

### PASCAL



bird



cat



dog

### ILSVRC



flamingo



cock



ruffed grouse



quail



partridge

...



Egyptian cat



Persian cat



Siamese cat



tabby



lynx

...



dalmatian



keeshond



miniature schnauzer



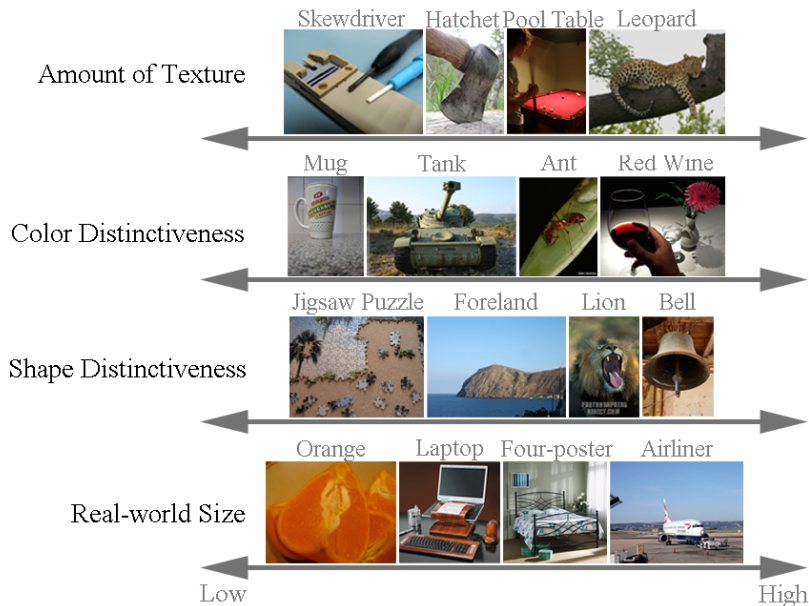
standard schnauzer



giant schnauzer

...







## Image classification

Steel drum



Ground truth

Steel drum  
Folding chair  
Loudspeaker

Accuracy: 1

Scale  
T-shirt  
Steel drum  
Drumstick  
Mud turtle

Accuracy: 1

Scale  
T-shirt  
Giant panda  
Drumstick  
Mud turtle

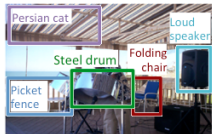
Accuracy: 0

## Single-object localization

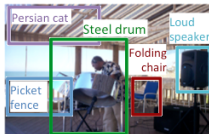
Steel drum



Ground truth



Accuracy: 1

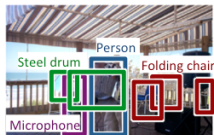


Accuracy: 0

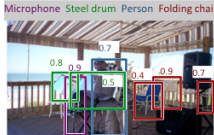


Accuracy: 0

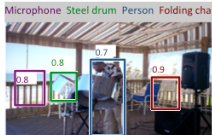
## Object detection



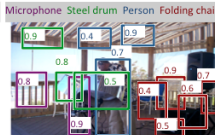
Ground truth



AP: 1.0 1.0 1.0 1.0



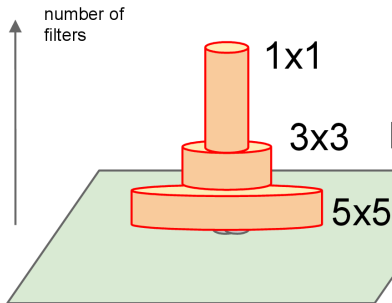
AP: 0.0 0.5 1.0 0.3



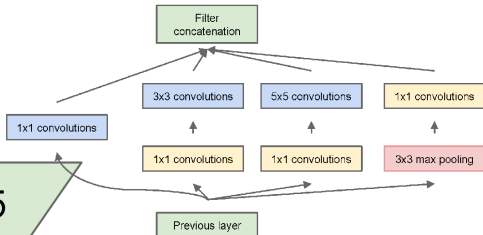
AP: 1.0 0.7 0.5 0.9

- Google of course! (first time)
- GoogLeNet:

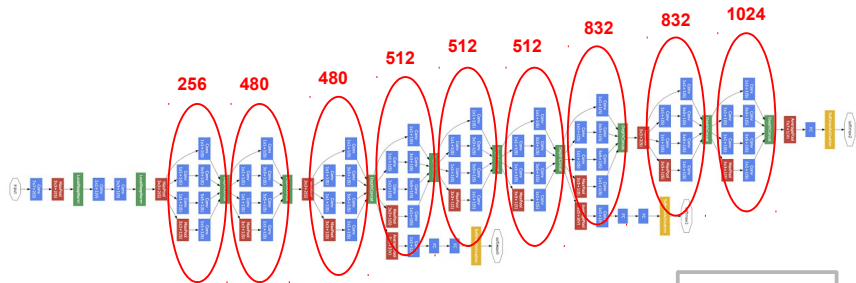
## Schematic view



## Inception module



- Google of course! (first time)
- GoogLeNet:



9 **Inception** modules

Network in a network in a network...

**Convolution**  
**Pooling**  
**Softmax**  
**Other**

## Classification failure cases



Groundtruth: **Police car**

GoogLeNet:

- laptop
- hair drier
- binocular
- **ATM machine**
- seat belt

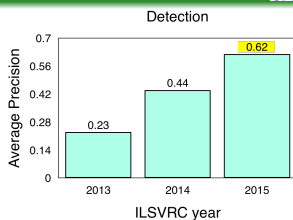
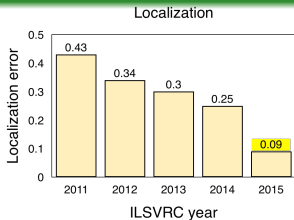
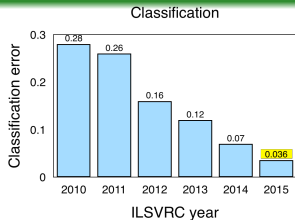
## Classification failure cases



Groundtruth: **hay**

GoogLeNet:

- sorrel (horse)
- hartebeest
- Arabian camel
- warthog
- gabelle



2010–14: 4.2x reduction

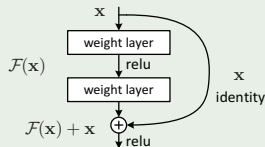
1.7x reduction

1.9x increase

## ILSVRC 2015 (same dataset as 2014)

► arXiv:1512.03385

- Winner: MSRA (Microsoft Research in Beijing)
- Deep residual networks with  $> 150$  layers
- Classification error: 6.7%  $\rightarrow$  3.6% (1.9x)
- Localisation error: 26.7%  $\rightarrow$  9.0% (2.8x)
- Object detection: 43.9%  $\rightarrow$  62.1% (1.4x)

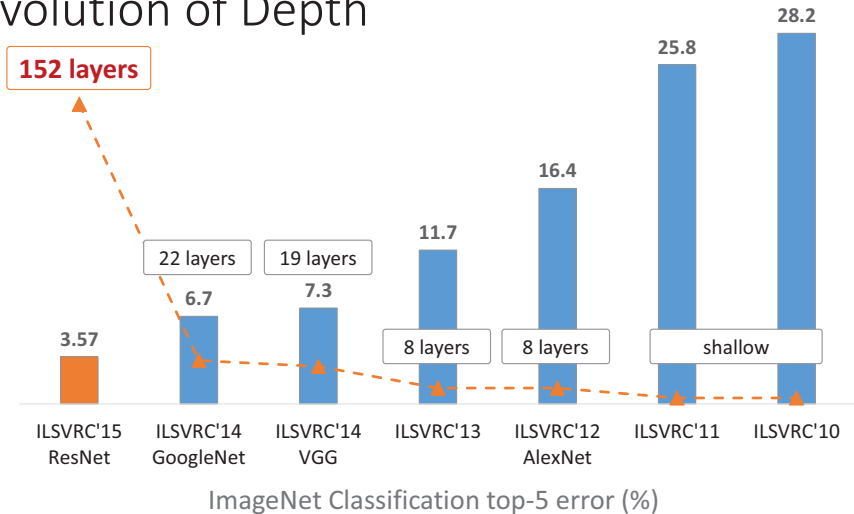


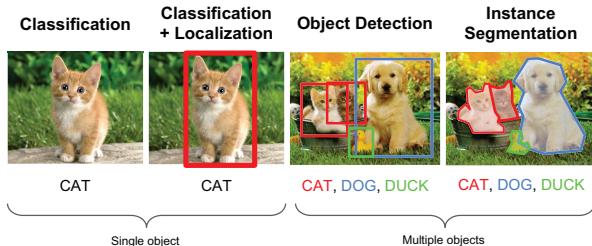
## ILSVRC 2016

► <http://image-net.org/challenges/LSVRC/2016>

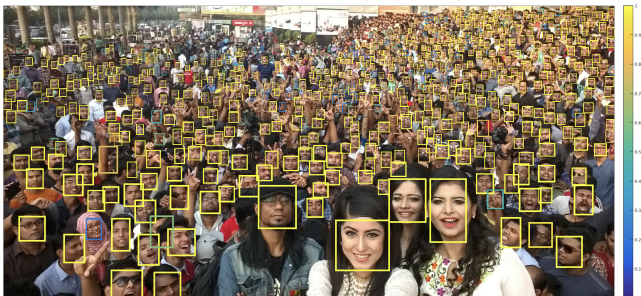
- Mostly ResNets. Classification: 0.030; localisation: 0.08; detection: 0.66

## Revolution of Depth





- More and more refinement (segmentation)
- More objects, in real time on video1/video2/video3





- Learning to play 49 different Atari 2600 games
- No knowledge of the goals/rules, just 84x84 pixel frames
- 60 frames per second, 50 million frames (38 days of game experience)
- Deep convolutional network with reinforcement: DQN (deep Q-network)
  - action-value function  $Q^*(s,a) = \max_{\pi} \mathbb{E}[r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots | s_t = s, a_t = a, \pi]$
  - maximum sum of rewards  $r_t$  discounted by  $\gamma$  at each timestep  $t$ , achievable by a behaviour policy  $\pi = P(a|s)$ , after making observation  $s$  and taking action  $a$
- Tricks for scalability and performance:
  - experience replay (use past frames)
  - separate network to generate learning targets (iterative update of Q)
- Outperforms all previous algorithms, and professional human player on most games

# Google DeepMind: training&performance

## Algorithm 1: deep Q-learning with experience replay.

Initialize replay memory  $D$  to capacity  $N$

Initialize action-value function  $Q$  with random weights  $\theta$

Initialize target action-value function  $\hat{Q}$  with weights  $\theta^- = \theta$

**For** episode = 1,  $M$  **do**

Initialize sequence  $s_1 = \{x_1\}$  and preprocessed sequence  $\phi_1 = \phi(s_1)$

**For**  $t = 1, T$  **do**

With probability  $\epsilon$  select a random action  $a_t$

otherwise select  $a_t = \operatorname{argmax}_a Q(\phi(s_t), a; \theta)$

Execute action  $a_t$  in emulator and observe reward  $r_t$  and image  $x_{t+1}$

Set  $s_{t+1} = s_t, a_t, x_{t+1}$  and preprocess  $\phi_{t+1} = \phi(s_{t+1})$

Store transition  $(\phi_t, a_t, r_t, \phi_{t+1})$  in  $D$

Sample random minibatch of transitions  $(\phi_j, a_j, r_j, \phi_{j+1})$  from  $D$

Set  $y_j = \begin{cases} r_j & \text{if episode terminates at step } j+1 \\ r_j + \gamma \max_{a'} \hat{Q}(\phi_{j+1}, a'; \theta^-) & \text{otherwise} \end{cases}$

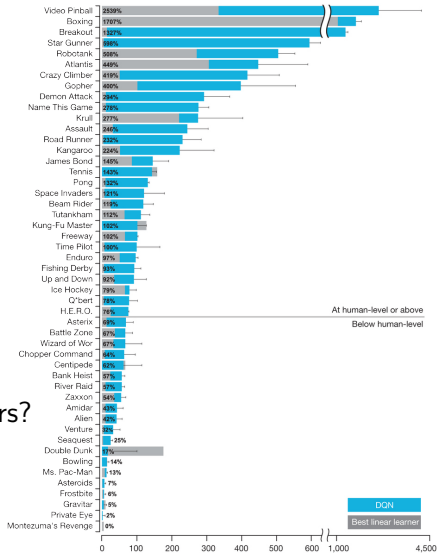
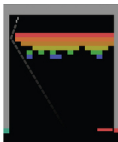
Perform a gradient descent step on  $(y_j - Q(\phi_j, a_j; \theta))^2$  with respect to the network parameters  $\theta$

Every  $C$  steps reset  $\hat{Q} = Q$

**End For**

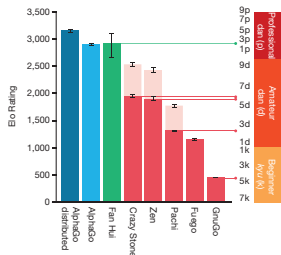
**End For**

## • What about Breakout or Space invaders?

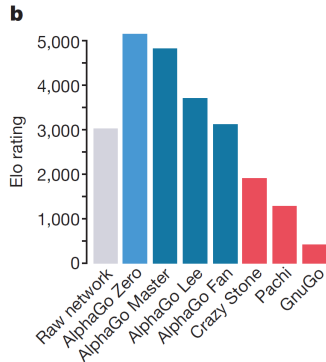
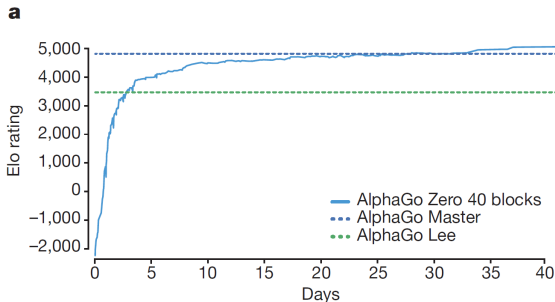


- Game of Go considered very challenging for AI
- Board games: can be solved with search tree of  $b^d$  possible sequences of moves ( $b =$  breadth [number of legal moves],  $d =$  depth [length of game])
- Chess:  $b \approx 35$ ,  $d \approx 80 \rightarrow$  go:  $b \approx 250$ ,  $d \approx 150$
- Reduction:
  - of depth by position evaluation (replace subtree by approximation that predicts outcome)
  - of breadth by sampling actions from probability distribution (policy  $p(a|s)$ ) over possible moves  $a$  in position  $s$
- $19 \times 19$  image, represented by CNN
- Supervised learning policy network from expert human moves, reinforcement learning policy network on self-play (adjusts policy towards winning the game), value network that predicts winner of games in self-play.

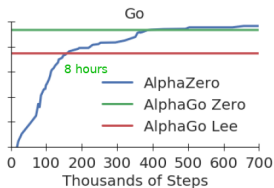
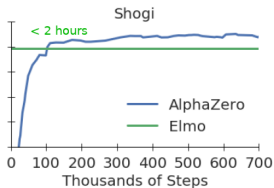
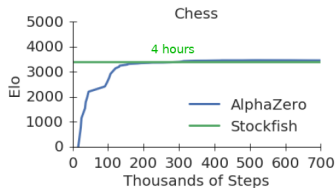
- AlphaGo: 40 search threads, simulations on 48 CPUs, policy and value networks on 8 GPUs. Distributed AlphaGo: 1020 CPUs, 176 GPUs
- AlphaGo won 494/495 games against other programs (and still 77% against Crazy Stone with four handicap stones)
- Fan Hui: 2013/14/15 European champion
- Distributed AlphaGo won 5–0
- AlphaGo evaluated thousands of times fewer positions than Deep Blue (first chess computer to beat human world champion)  $\Rightarrow$  better position selection (policy network) and better evaluation (value network)
- Then played Lee Sedol (top Go play in the world over last decade) in March 2016  $\Rightarrow$  won 4–1. AlphaGo given honorary professional ninth dan, considered to have “reach a level ‘close to the territory of divinity’ ”
- Ke Jie (Chinese world #1): “Bring it on!”. May 2017: 3–0 win for AlphaGo. New comment: “I feel like his game is more and more like the ‘Go god’”. Really, it is brilliant”



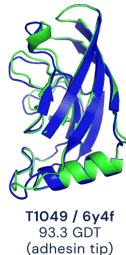
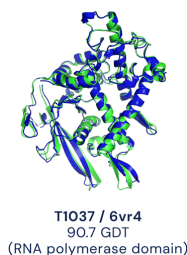
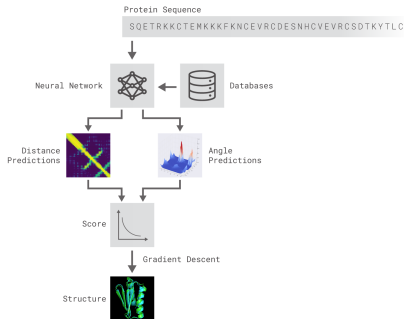
- Learn from scratch, just from the rules and random moves
- Reinforcement learning from self-play, no human data/guidance
- Combined policy and value networks
- 4.9 million self-play games
- Beats AlphaGo Lee (several months of training) after just 36 hours
- Single machine with four TPU



- Same philosophy as AlphaGo Zero, applied to chess, shogi and go
- Changes:
  - not just win/loss, but also draw or other outcomes
  - no additional training data from game symmetries
  - using always the latest network to generate self-play games rather than best one
  - tree search: 80k/70M for chess AlphaZero/Stockfish, 40k/35M for shogi AlphaZero/Elmo

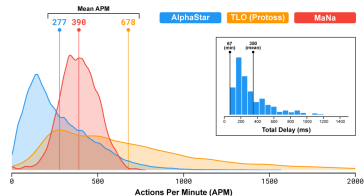
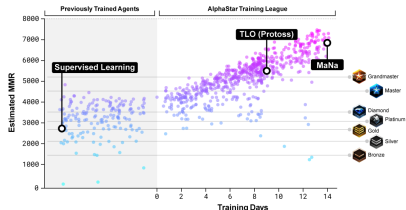


- Trying to tackle scientific problem
- Goal: predict 3D structure of protein based solely on genetic sequence
- Using DNN to predict
  - distances between pairs of amino acids
  - angles between chemical bonds
- Search DB to find matching existing substructures
- Also train a generative NN to invent new fragments
- Achieved best prediction ever



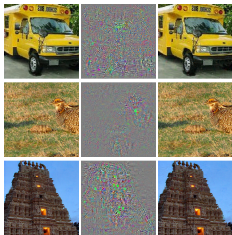
● Experimental result  
● Computational prediction

- Mastering real-time strategy game StarCraft II
- Challenges in game theory (no single best strategy), imperfect information (hidden parts of game), long term planning, real time (continuous flow of actions), large action space (many units/buildings)
- Using DNN trained
  - directly on raw data games
  - supervised learning on human games
  - reinforcement learning (continuous league)
- DNN output: list of actions
- Trained for 14 days; each agent: up to 200 years of real-time play
- Runs on single desktop GPU
- Defeated 5–0 one of best pro-players



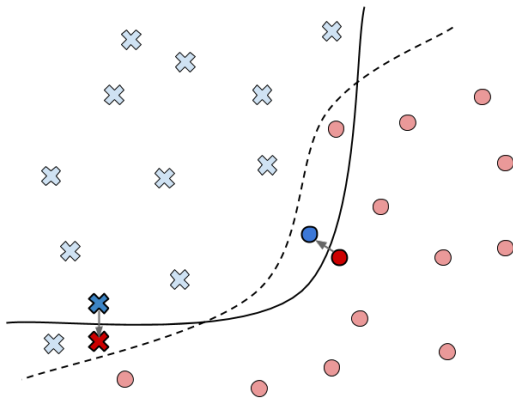


- Playing poker
  - Libratus (AI developed by Carnegie Mellon University) defeated four of the world's best professional poker players (Jan 2017)
  - After 120,000 hands of Heads-up, No-Limit Texas Hold'em, led the pros by a collective \$1,766,250 in chips
  - Learned to bluff, and win with incomplete information and opponents' misinformation
- Lip reading [▶ arXiv:1611.05358 \[cs.CV\]](#)
  - human professional: deciphers less than 25% of spoken words
  - CNN+LSTM trained on television news programs: 50%
- Limitation: adversarial attacks [▶ arXiv:1312.6199 \[cs.CV\]](#)



- left: correctly classified image
- middle: difference between left image and adversarial image ( $\times 10$ )
- right: adversarial image, classified as ostrich

# Adversarial attack: what is happening?



----- Task decision boundary

———— Model decision boundary

⊗ Test point for class 1

⊗ Adversarial example for class 1

⊗ Training points for class 1

● Training points for class 2

● Test point for class 2

● Adversarial example for class 2



original semantic segmentation framework



original semantic segmentation framework



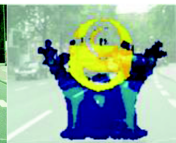
adversarial attack



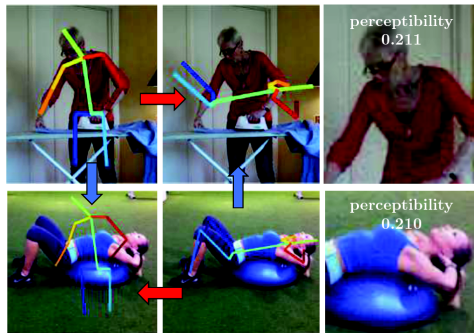
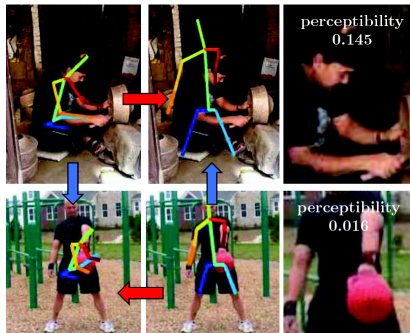
original semantic segmentation framework

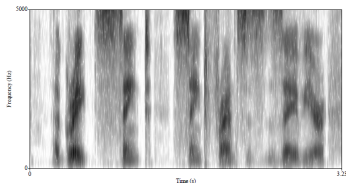


adversarial attack

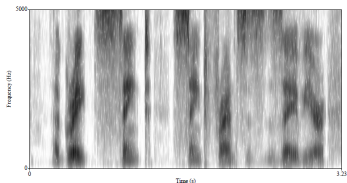


compromised semantic segmentation framework





(a) a great saint saint francis zaviour



(b) i great sinkt shink t frimsuss avir

Figure 7: The model models' output for each of the spectrograms is located at the bottom of each spectrogram. The target transcription is: A Great Saint Saint Francis Xavier.

## AllConv



SHIP

CAR(99.7%)



HORSE

DOG(70.7%)



CAR

AIRPLANE(82.4%)

## NiN



HORSE

FROG(99.9%)



DOG

CAT(75.5%)



DEER

DOG(86.4%)

## VGG



DEER

AIRPLANE(85.3%)



BIRD

FROG(86.5%)



CAT

BIRD(66.2%)



## AllConv



SHIP

CAR(99.7%)



HORSE

DOG(70.7%)



CAR

AIRPLANE(82.4%)

## NiN



HORSE

FROG(99.9%)



DOG

CAT(75.5%)



DEER

DOG(86.4%)

## VGG



DEER

AIRPLANE(85.3%)



BIRD

FROG(86.5%)



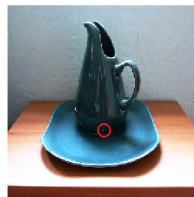
CAT

BIRD(66.2%)



Cup(16.48%)

Soup Bowl(16.74%)



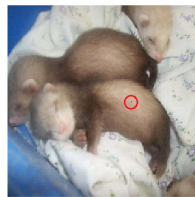
Teapot(24.99%)

Joystick(37.39%)



Bassinet(16.59%)

Paper Towel(16.21%)



Hamster(35.79%)

Nipple(42.36%)





## Hype Cycle for Emerging Technologies, 2020



[gartner.com/SmarterWithGartner](https://www.gartner.com/SmarterWithGartner)

Source: Gartner  
© 2020 Gartner, Inc. and/or its affiliates. All rights reserved. Gartner and Hype Cycle are registered trademarks of Gartner, Inc. and its affiliates in the U.S.

Gartner.

- No deep learning/ML anymore (since 2019)
- Instead, all sorts of “AI” and ML-driven systems
  - explainable AI
  - embedded AI
  - generative AI / generative adversarial networks
  - adaptive ML
  - self-supervised learning
  - AI-assisted design

## Data Science @ LHC 2015

Bridging High-Energy Physics and Machine Learning communities

9 - 13 November 2015, CERN

### Local Organising Committee

- Xavier Cid (CERN)
- Gilles Louppe (CERN)
- Michelangelo Mangano (CERN)
- Maurizio Perini (CERN)
- Jean-Roch Vlimant (Caltech)

### Program Committee

- Kyle Cranmer (New York U)
- Cécile Germain (ILR)
- Vladimir Vazir Magarov (CERN)
- Gilles Louppe (CERN)
- Andrew Lowe (Wigner RCP)
- Maurizio Perini (CERN)
- David Rousseau (LAL Orsay)
- Maria Spiropulu (Caltech)
- Jean-Roch Vlimant (Caltech)
- Daniel Whiteson (UC Irvine)

sponsored by

LHC Physics Center at CERN: <http://lpc.web.cern.ch>  
Fermilab National Laboratory: <http://fnal.gov>  
Moore-Sloan Data Science Environment: <http://cds.nyu.edu/mooresloan>

### International Advisory Committee

- Roger Barlow (Huddersfield U)
- Tommaso Dorigo (INFN Padova)
- Ian Fisk (Simons Foundation)
- Maria Gionni (CERN)
- Eilam Gross (Weizmann)
- Balázs Kégl (LAL Orsay)
- Constantin Louizides (ILNL)
- Stuart Russell (UC Berkeley)
- Victoria Stodden (Ill Urbana-Champaign)
- Max Welling (Amsterdam U)

<http://cern.ch/DataScienceLHC2015>

<http://opendata.cern.ch>

## Data Science @ LHC 2015

Bridging High-Energy Physics and Machine Learning communities

Exploring the potential for Machine Learning on ATLAS

## ATLAS Machine Learning Workshop

29<sup>th</sup>-31<sup>st</sup> March 2016, CERN

### Organising Committee:

Matthew Beckingham (Warwick)  
Michael Kagan (SLAC)  
David Rousseau (LAL-Orsay)

<http://cern.ch/AtlasML2016>

<http://opendata.cern.ch>

## Data Science @ LHC 2015

Bridging High-Energy Physics and Machine Learning communities

Exploring the potential for Machine Learning on ATLAS

## ATLAS Machine Learning Workshop

# MLHEP

20-26 June 2016  
Lund, Sweden

### Second Machine Learning School for High Energy Physics

<http://cern.ch/AtlasML2016>

<http://opendata.cern.ch>

## Data Science @ LHC 2015

Bridging High-Energy Physics and Machine Learning communities

Exploring the

# ATLAS Works

# M

## Second M

Higgs  
challenge



## the HiggsML challenge

May to September 2014

When High Energy Physics meets Machine Learning



info to participate and compete : <https://www.kaggle.com/c/higgs-boson>



Organization committee

Yann Coadou - ATLAS  
Lucie Gammelin - INRIA

David Boucayan - ATLAS  
Glen Cowan - ATLAS

Isabelle Guyon - CERN  
Claire Adams-Bourdeau - ATLAS

Advisory committee

Thorsten Muehlbauer - ATLAS  
Andreas Hocker - ATLAS

Jeremy Shotton - ATLAS  
Ralf Schaefer - ATLAS

g on ATLAS

# Learning

0-26 June  
Sweden

# 2016

## Energy Physics

# L2016

<http://opendata.cern.ch>



## Data Science @ LHC 2015

Bridging High-Energy Physics and Machine Learning communities

Exploring the

# ATLAS Workshop

# M

Second Meeting

Higgs challenge



## the HiggsML challenge

May to September 2014

When High Energy Physics meets Machine Learning

Workshop on ATLAS

# Machine Learning

## NIPS 2016

Monday December 05 -- Saturday December 10, 2016

Centre Convencions Internacional Barcelona, Barcelona SPAIN

2016 Pricing »

Registration 2016 »

Dates

Calls

Student Support

Program Books

Schedule

Barcelona

View Earlier Meetings »

2015 Workshop Videos »

### Invited Speakers

Yann LeCun (Facebook), Susan Holmes (Stanford), **Kyle Cranmer (NYU)**, Saket Navlakha (Saik Institute), Drew Purves (Deep Mind), Marc Raibert (Boston Dynamics), Irina Rish (IBM)

### Tutorials

The tutorial times and rooms have not been set yet. View the list of tutorials using the button below.

View Tutorials »

info to participate and compete



Organization committee

Holger Kogler - ATLAS/LHC, Gerd Gies - ATLAS/LHC, David Rousseau - ATLAS/LHC, Glen Cowan - ATLAS/LHC

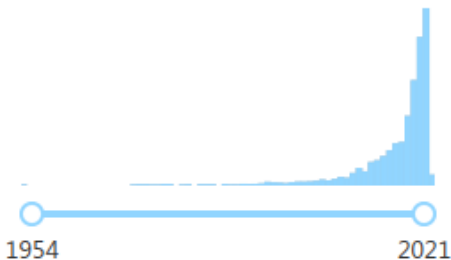
http://

http://opendata.cern.ch





## Date of paper

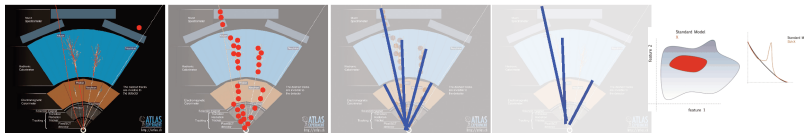


► [https://inspirehep.net/literature?q=machine learning or deep learning or multivariate](https://inspirehep.net/literature?q=machine+learning+or+deep+learning+or+multivariate)

**Up-to-date comprehensive review of papers**

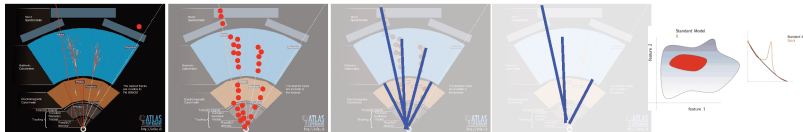
► <https://github.com/iml-wg/HEPML-LivingReview>

<u>Raw</u>	<u>Sparsified</u>	<u>Reco</u>	<u>Select</u>	<u>Physics</u>	<u>Ana</u>
1e7	1e4	100-ish*	50	10	1



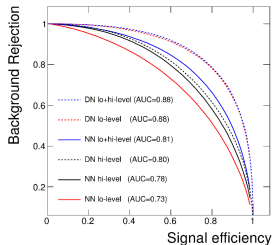
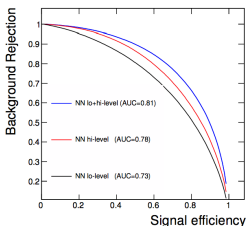
- Reduce data dimensionality to allow analysis

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



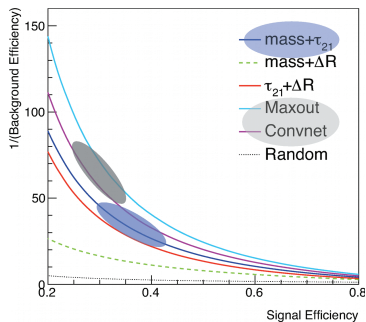
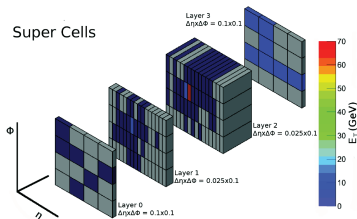
- Reduce data dimensionality to allow analysis
- Going to lower level features

▶ arXiv:1402.4735

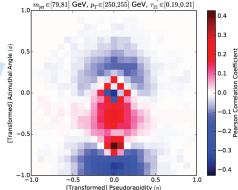


- Transforming inputs into images

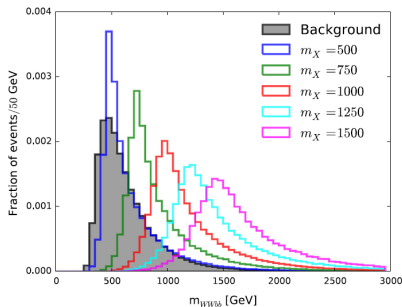
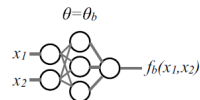
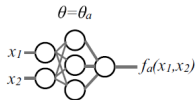
▶ arXiv:1511.05190



Correlation of Deep Network output with pixel activations.

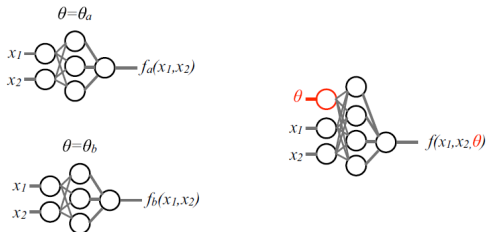


- Looking for new physics scenario with unknown mass  
⇒ one NN for each mass point





- Looking for new physics scenario with unknown mass  
 $\Rightarrow$  one NN for each mass point

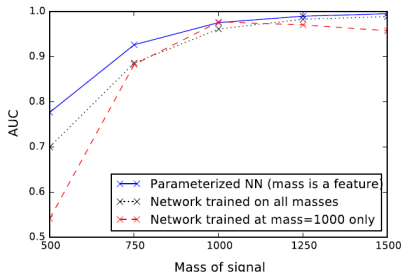
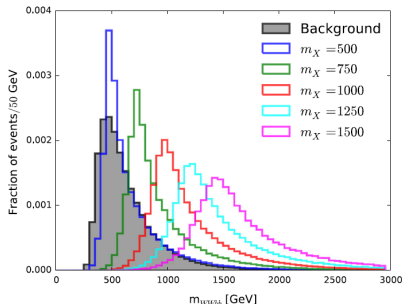


## Parameterised NN

► EPJC (2016) 76:235

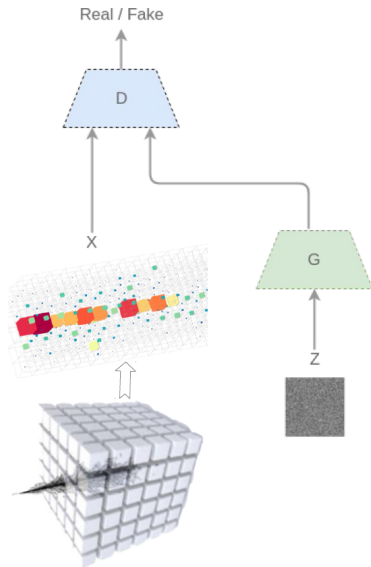
► arXiv:1601.07913

- mass as training parameter
- as good as dedicated training
- generalises better



## Fast simulation with generative models

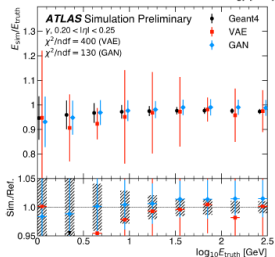
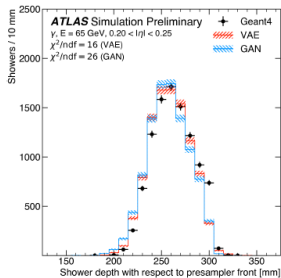
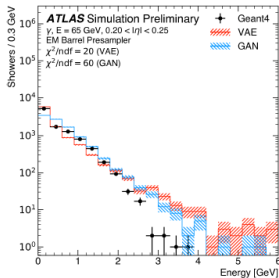
- Heavy CPU cost of simulation (> 50% of grid resources)
  - MC stats becoming limiting factor in analyses
- Replace “full simulation” with approximation
  - already routinely done, using parameterisation of showers or library of pre-simulated objects



## Fast simulation with generative models

▶ ATL-SOFT-PUB-2018-001

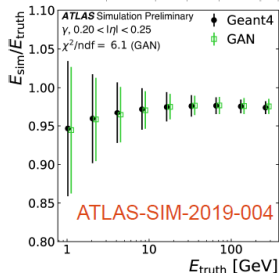
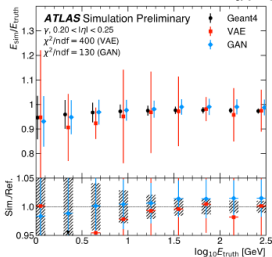
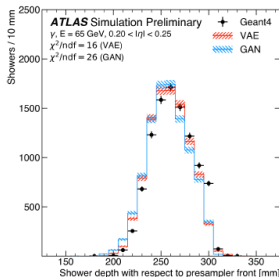
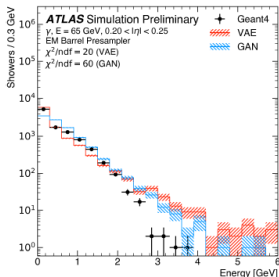
▶ ATLAS-SIM-2019-004



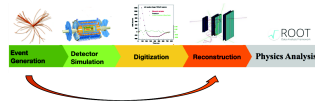
## Fast simulation with generative models

▶ ATLAS-SOFT-PUB-2018-001

▶ ATLAS-SIM-2019-004



- $< 1$  ms instead of 10 s per object!
- $\times 100$  gain on complete event
- Still some work to be done

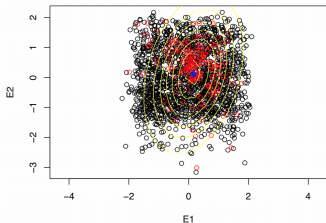


## Anomaly detection: looking for new physics

- Learn background (SM) properties
- Flag deviations from prediction without knowing anything about specific new physics scenario

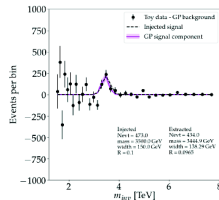
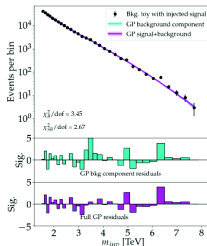
## Penalised anomaly detection

- based on Gaussian mixture model
- $f_S$  and  $f_B$ : finite sums of Gaussians
- semi-supervised training
- penalty term in LH to select variables



## Gaussian processes

- Learn background with GP instead of parametric model
- Compare data to new GP: background model+signal
- Returns bg parameters of “peak”



- Very active field of research in machine learning and artificial intelligence
  - not just at universities (Google, Facebook, Microsoft, NVIDIA, etc. . . )
- Training with curriculum:
  - what humans do over 20 years, or even a lifetime
  - learn different concepts at different times
  - solve easier or smoothed version first, and gradually consider less smoothing
  - exploit previously learned concepts to ease learning of new abstractions
- Influence learning dynamics can have big impact:
  - order and selection of examples matters
  - choose which examples to present first, to guide training and possibly increase learning speed (called shaping in animal training)
- Combination of deep learning and **reinforcement learning**
  - still in its infancy, but already impressive results
- **Domain adaptation and adversarial training**
  - e.g. train in parallel network that produces difficult examples
  - learn discrimination (s vs. b) and difference between training and application samples (e.g. Monte Carlo simulation and real data)




**UNIFORM CONVERGENCE MAY BE UNABLE TO EXPLAIN GENERALIZATION IN DEEP LEARNING**

VASHNAVI NADARAJAN<sup>1</sup>  
Computer Science Department, Carnegie Mellon University
ZICO KOLTER<sup>1\*</sup>  
Bosch Center for Artificial Intelligence, Pittsburgh
BOSCH  
S

### THE HIGH LEVEL MESSAGE

The existing theory states that deep learning generalizes well because of its implicit regularization properties.



We argue that the high-level direction of deriving uniform convergence bounds from the over-parameterized hypothesis space is the generalization direction.

### THE MAIN CONTRIBUTIONS & OPENING QUESTIONS (1/2)

Conventional  $\epsilon$ -bound like VC-dim fail to explain generalization in deep learning [1,2].

generalization gap  $\leq$   $\left\{ \begin{array}{l} \text{representational complexity of whole hypothesis class} \\ \text{training set size} \end{array} \right\}$  (non-convex)

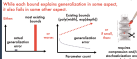
For tighter, more meaningful bounds, the proposed suggestion was to identify *implicit bias* and use it to refine  $\epsilon$ -bound.

generalization gap  $\leq$   $\left\{ \begin{array}{l} \text{representational complexity of relevant subset of hypothesis class} \\ \text{training set size} \end{array} \right\}$

Intuitively, diversity is *not* necessarily bad for SGD. It is, only from irrelevant diversity.

### THE MAIN CONTRIBUTIONS & OPENING QUESTIONS (2/2)

Many more novel, refined  $\epsilon$ -bounds have been proposed, using Rademacher complexity, covering numbers, compressors, PAC-Bayes. While each bound involves generalization in some degree, if they fail in some other degree.

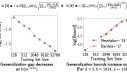


### ONLY TO BEING GENERALIZATION BOUNDS FROM TRAINING SET SIZE

Stability, for instance, is the best output of network  $\pi$  on data  $S$  is  $\mathcal{F}(S)$ . On  $\mathcal{D}$  output is  $\mathcal{F}(\mathcal{D})$ . Define margin of  $\mathcal{F}$  to be  $\mathcal{F}(\mathcal{D}) - \mathcal{F}(S) = \mathbb{E}[\mathcal{F}(S')] - \mathcal{F}(S)$ .

Let  $\pi$  be training set of  $n$  examples drawn i.i.d from  $\mathcal{D}$ . Define network margin for  $\pi$  and width by  $\gamma$ .

Empirical margin, SGD with learning rate 0.1 and minibatch size 1 on 95% of positive class of MNIST is illustrated by a sample of 50. We evaluate bounds from [2], while bounds below analyze.



where  $\mathcal{F}(S) = \mathbb{E}[\mathcal{F}(S)] + \mathcal{O}(\frac{1}{\sqrt{n}})$ ,  $\mathcal{F}(\mathcal{D}) = \mathbb{E}[\mathcal{F}(S)] + \mathcal{O}(\frac{1}{\sqrt{n}})$

[1]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[2]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[3]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[4]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[5]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[6]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[7]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[8]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[9]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[10]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[11]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[12]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[13]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[14]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[15]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[16]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[17]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[18]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[19]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[20]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[21]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[22]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[23]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[24]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[25]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[26]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[27]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[28]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[29]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[30]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[31]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[32]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[33]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[34]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[35]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[36]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[37]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[38]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[39]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[40]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[41]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[42]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[43]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[44]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[45]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[46]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[47]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[48]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[49]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[50]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[51]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[52]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[53]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[54]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[55]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[56]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[57]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[58]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[59]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[60]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[61]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[62]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[63]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[64]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[65]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[66]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[67]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[68]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[69]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[70]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[71]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[72]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[73]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[74]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[75]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[76]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[77]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[78]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[79]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[80]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[81]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[82]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[83]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[84]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[85]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[86]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[87]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[88]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[89]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[90]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[91]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[92]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[93]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[94]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[95]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[96]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[97]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[98]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[99]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)

[100]  $\mathbb{E}[\mathcal{F}(S)] - \mathcal{F}(S) \leq \frac{1}{\sqrt{n}}$  (margin)


### Key element in proof: Tightest Uniform Convergence

Definition: Let  $\mathcal{F}$  be hypothesis learned on dataset  $\mathcal{D}$ . Let  $\mathcal{F}(\mathcal{D})$  denote the  $\mathcal{D}$  error of  $\mathcal{F}$ . Let  $\mathcal{F}(\mathcal{D}')$  denote the error of  $\mathcal{F}$  on a new dataset  $\mathcal{D}'$ .

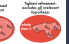
Def.1: The generalization gap is the smallest value of  $\mathcal{F}(\mathcal{D}')$  such that  $\mathcal{F}(\mathcal{D}) - \mathcal{F}(\mathcal{D}') \leq \epsilon$ .

Def.2: The "uniform"  $\epsilon$ -bound is the smallest value of  $\mathcal{F}(\mathcal{D}')$  such that  $\mathcal{F}(\mathcal{D}) - \mathcal{F}(\mathcal{D}') \leq \epsilon$ .

To define this bound, we can consider many different kinds of "relevant subsets" in  $\mathcal{H}$ .



Linear-based relevant subset




Higher-order polynomial relevant subset

Def.3: The optimal uniform  $\epsilon$ -bound is the smallest value  $\epsilon$  such that  $\mathcal{F}(\mathcal{D}) - \mathcal{F}(\mathcal{D}') \leq \epsilon$  and  $\mathcal{F}(\mathcal{D}) - \mathcal{F}(\mathcal{D}') \leq \epsilon$ .

### FAILURE OF U.C. IN A REPRESENTATIVE EXAMPLE

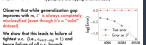
We train a 3-layer network of width  $n = 100$  using SGD to classify two synthetic classes, 1000-dimensional hyperspheres of radius 1 and 1.1.



Next, we create a proposed training set  $\mathcal{D}$ . By sampling out training examples to separate hyperspheres and flipping to create label.

Observe that while generalization gap improves with  $n$ , it is always bounded by the margin between the two classes.

We show that the bounds below of tightest  $\epsilon$ -bound  $\epsilon$  will not learn better of all  $n$ -width.



### PROOF SKETCH

They learned. For any given training set  $\mathcal{D}$ , if we can design a corresponding "test set"  $\mathcal{D}'$  such that  $\mathcal{F}(\mathcal{D}) - \mathcal{F}(\mathcal{D}') \leq \epsilon$ , then  $\mathcal{F}(\mathcal{D}) - \mathcal{F}(\mathcal{D}') \leq \epsilon$ .

Def.4: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.5: The decision boundary learned by SGD on  $\mathcal{D}'$  is the decision boundary which best separates the training data.

Def.6: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.7: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.8: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.9: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.10: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.11: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.12: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.13: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.14: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.15: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.16: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.17: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.18: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.19: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.20: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.21: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.22: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.23: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.24: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.25: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.26: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.27: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.28: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.29: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.30: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.31: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.32: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.33: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.34: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.35: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.36: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.37: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.38: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.39: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.40: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.41: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.42: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.43: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.44: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.45: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.46: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.47: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.48: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.49: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.50: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.51: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.52: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.53: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.54: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.55: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.56: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.57: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.58: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.59: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.60: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.61: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.62: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.63: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

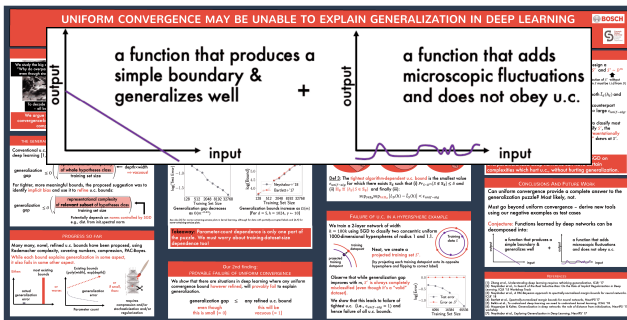
Def.64: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.65: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.66: The decision boundary learned by SGD on  $\mathcal{D}$  is the decision boundary which best separates the training data.

Def.67: The decision boundary learned by SGD on  $\math$





**UNIFORM CONVERGENCE MAY BE UNABLE TO EXPLAIN GENERALIZATION IN DEEP LEARNING**

**a function that produces a simple boundary & generalizes well**

**a function that adds microscopic fluctuations and does not obey u.c.**

**output** vs **input**

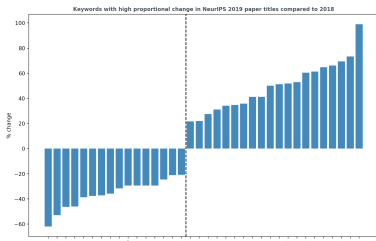
**output** vs **input**

**CONCLUDING AND FUTURE WORK**  
Can uniform convergence provide a complete answer to the generalization puzzle? Most likely, not.  
Must go beyond uniform convergence - derive new tools using our negative examples on test cases

**Complexity**: Functions learned by deep networks can be decomposed into:  
- a function that produces a smooth boundary & generalizes well  
- a function that adds microscopic fluctuations and does not obey u.c.

- Actually works surprisingly well
- Over-parameterised DNN should overfit but don't: why?
- Neural tangent kernel (NTK): helps thinking in infinite-width limit. But can do better in reality
- Robustness to adversarial attacks
- Start with large learning rate to learn easy features, then decrease to learn low noise, hard-to-fit patterns

- More people active with neurosciences: ML to understand NS, and NS to understand ML
- Meta learning (learning to learn)
- Reinforcement learning is gaining ground. Other keywords: bandit, feedback, regret, control
- Attributing uncertainty to ML algorithms (often with Bayesian methods in deep learning)
- Generative models still popular
- Hardware keyword on the rise, signaling more hardware-aware algorithms: hardware = bottleneck?
- “Recurrent and convolutional neural networks are literally so last year”
- Growing consciousness of potential impact on society



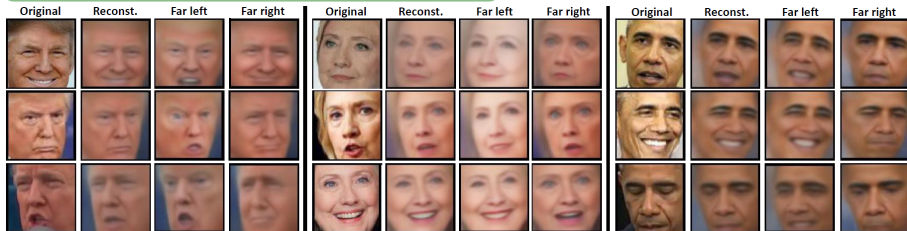
- 91 short papers accepted for poster presentation (6 selected for talks)
- 70 “digital acceptance” papers (above rejection threshold but beyond capacity)
- 228 referees [▶ web site](#) (incl. videos)
- 5 invited speakers:
  - Alan Aspuru-Guzik: *Recent progress in ML for chemistry: SELFIES, inverse design of drug candidates and materials, and Bayesian algorithms for self-driving laboratories*
  - Yasaman Bahri: *Towards an understanding of wide, deep neural networks*
  - Katie Bouman: Cannot find title, about Event Horizon Telescope imaging technique
  - Bernhard Schölkopf: *Causality and Exoplanets*
  - Maria Schuld: *Innovating machine learning with near-term quantum computing*
  - Lenka Zdeborova: *Understanding machine learning via exactly solvable statistical physics models*

## Suggested areas

- Application of machine and deep learning to physical sciences
- Generative models
- Likelihood-free inference
- Variational inference
- Simulation-based models
- Implicit models
- Probabilistic models
- Model interpretability
- Approximate Bayesian computation
- Strategies for incorporating prior scientific knowledge into machine learning algorithms
- Experimental design
- Any other area related to the subject of the workshop

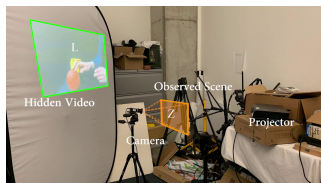
- ML achieves super-human performance for well-designed problems, or games with score  $\Rightarrow$  where one can define a proper loss function or reward
- Scale to “real” problems?
  - explainability
  - causality
  - “moral” stand
  - culture, art
- Many advances in medical imaging, modelling of various phenomena, supernova analysis or LHC physics, but issues with:
  - out-of-distribution generalisation
  - scalability of computing resources, carbon footprint
  - reliability
  - decision bias (gender, race, etc.)
- Workshops/socials on Fairness & ethics, AI for Good, Tackling Climate Change with ML, AI for Humanitarian Assistance and Disaster Response, Safety and Robustness in Decision-making, ...
- Importance of personal decisions

## ▶ Predicting the Politics of an Image Using Webly Supervised Data



## ▶ Computational Mirrors: Blind Inverse Light Transport by Deep Matrix Factorization

(edited video)



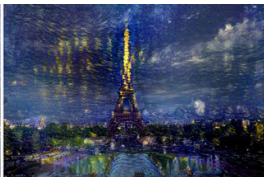
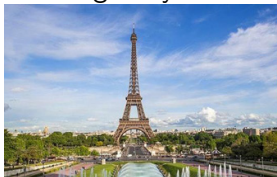
- Many techniques and tools exist to achieve optimal discrimination
- (Un)fortunately, no one method can be shown to outperform the others in all cases
- One should try several and pick the best one for any given problem
- Latest machine learning algorithms (e.g. deep networks) require enormous hyperparameter space optimisation. . .
- Machine learning and multivariate techniques are at work in your everyday life without your knowing and can easily outsmart you for many tasks
- Try this to convince yourself [▶ http://www.phy-t.de/mousegame/mousegame\\_en.html](http://www.phy-t.de/mousegame/mousegame_en.html)

**Upcoming reference book (in about six months)**

Artificial Intelligence for High Energy Physics

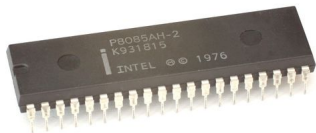
[▶ https://doi.org/10.1142/12200](https://doi.org/10.1142/12200)

- Learning a style [▶ arXiv:1508.06576 \[cs.CV\]](https://arxiv.org/abs/1508.06576) [▶ Neural-style](#)



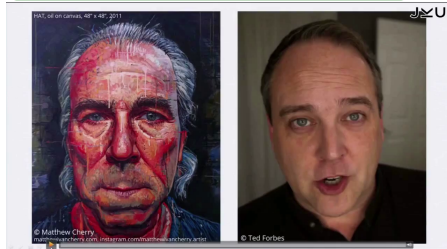
- Computer dreams [▶ Google original](#)

[▶ deepdream](#)



- Face Style [▶ http://facestyle.org](http://facestyle.org)

[▶ http://dcgi.fel.cvut.cz/home/sykorad/facestyle.html](http://dcgi.fel.cvut.cz/home/sykorad/facestyle.html)



## Summary

Monet  $\leftrightarrow$  Photos



Monet  $\rightarrow$  photo

Zebra  $\leftrightarrow$  Horses



zebra  $\rightarrow$  horse

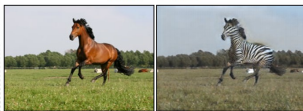
Summer  $\leftrightarrow$  Winter



summer  $\rightarrow$  winter



photo  $\rightarrow$  Monet



horse  $\rightarrow$  zebra



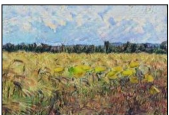
winter  $\rightarrow$  summer



Photograph



Monet



Van Gogh



Cezanne



Ukiyo-e



## From Monnet to photograph

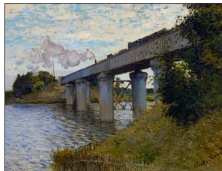
Input



Output



Input



Output



## Style transfer

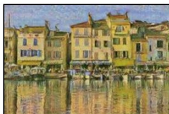
Input

Monet

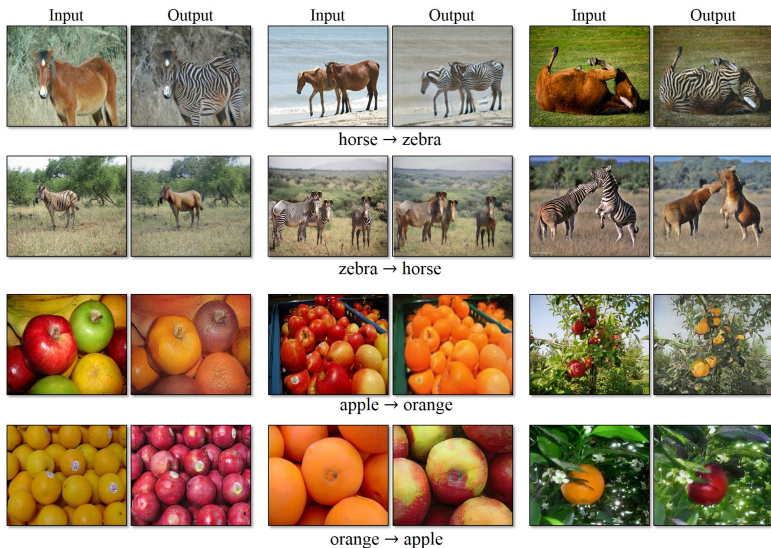
Van Gogh

Cezanne

Ukiyo-e



## Object transfiguration



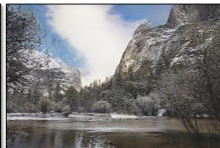
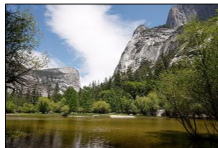
## Season transfer



winter Yosemite → summer Yosemite



summer Yosemite → winter Yosemite










## Photo enhancement











Failure 😊



-  C.M. Bishop, *Pattern Recognition and Machine Learning*, Springer, New York, 2007
-  M. Minsky and S. Papert, *Perceptrons*, MIT Press, Cambridge, 1969
-  W.S. McCulloch & W. Pitts, "A logical calculus of the ideas immanent in nervous activity", *Bulletin of Mathematical Biophysics*, 5, 115-137, 1943
-  F. Rosenblatt, "The Perceptron: A Probabilistic Model for Information Storage & Organization in the Brain", *Psychological Review*, 65, pp. 386-408, 1958
-  D.E. Rumelhart et al., "Learning representations by back-propagating errors", *Nature* vol. 323, p. 533, 1986
-  K. Hornik et al., "Multilayer Feedforward Networks are Universal Approximators", *Neural Networks*, Vol. 2, pp 359-366, 1989
-  Y. LeCun, L. Bottou, G. Orr and K. Muller, "Efficient BackProp", in *Neural Networks: Tricks of the trade*, Orr, G. and Muller K. (Eds), Springer, 1998

-  Q.V. Le et al., “Building High-level Features Using Large Scale Unsupervised Learning”, in *Proceedings of the 29th International Conference on Machine Learning*, Edinburgh, Scotland, UK, 2012 ▶ <http://research.google.com/pubs/pub38115.html>
-  Y. Bengio, P. Lamblin, D. Popovici and H. Larochelle, “Greedy Layer-Wise Training of Deep Networks”, in *Advances in Neural Information Processing Systems 19 (NIPS'06)*, pages 153–160, MIT Press 2007
-  M.A. Ranzato, C. Poultney, S. Chopra and Y. LeCun, in J. Platt et al., “Efficient Learning of Sparse Representations with an Energy-Based Model”, in *Advances in Neural Information Processing Systems 19 (NIPS'06)*, pages 1137–1144, MIT Press, 2007
-  Y. Bengio, “Learning deep architectures for AI”, *Foundations and Trends in Machine Learning*, Vol. 2, No. 1 (2009) 1–127. Also book at ▶ [Now Publishers](http://www.iro.umontreal.ca/lisa/publications2/index.php/publications/show/239)  
▶ <http://www.iro.umontreal.ca/lisa/publications2/index.php/publications/show/239>
-  I. Goodfellow, Y. Bengio and A. Courville, *Deep Learning*, MIT Press, 2016  
▶ <http://www.deeplearningbook.org>
-  G. Carleo et al., “Machine learning and the physical sciences”  
▶ [Rev. Mod. Phys. 91, 045002 \(2019\)](https://doi.org/10.1103/RevModPhys.91.045002) ▶ [arXiv:1903.10563](https://arxiv.org/abs/1903.10563)