

Python for Data Science 2

Lecture 11- Machine Learning

Amir Farbin

ML Basics

- In this context, 3 things define a ML model
 - ***Technique***: What type of technique are you applying? e.g. Multilayered Perceptron (MLP)- a type of Neural Network.
 - ***Hyper-parameters***: Specific to the technique. e.g. number of neurons.
 - ***Parameters***: What is learned. e.g. weight of connections between neurons.
- Generally 2 primary modes a ML model is used
 - ***Training***: Example data → Trained Model
 - ***Inference***: Trained model is applied to data → predictions
- Training Data generally separated into sub-sets
 - ***Training*** → obtain parameters
 - ***Test*** → select technique and hyper-parameters
 - ***Validation*** → used to assess performance

Data Representation

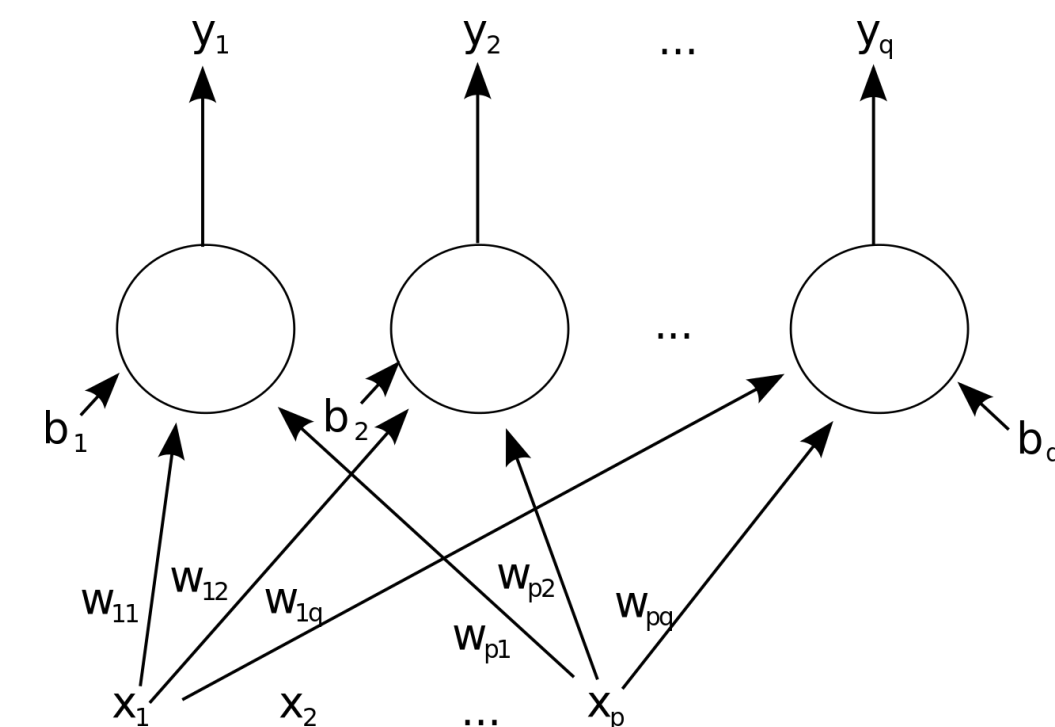
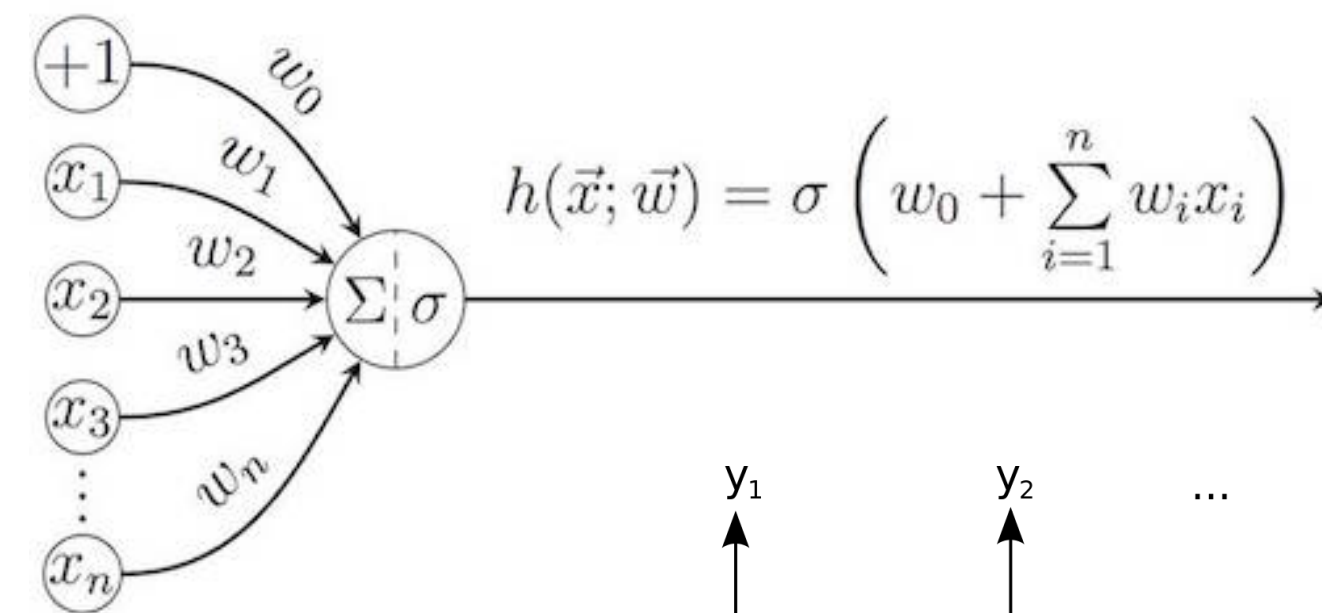
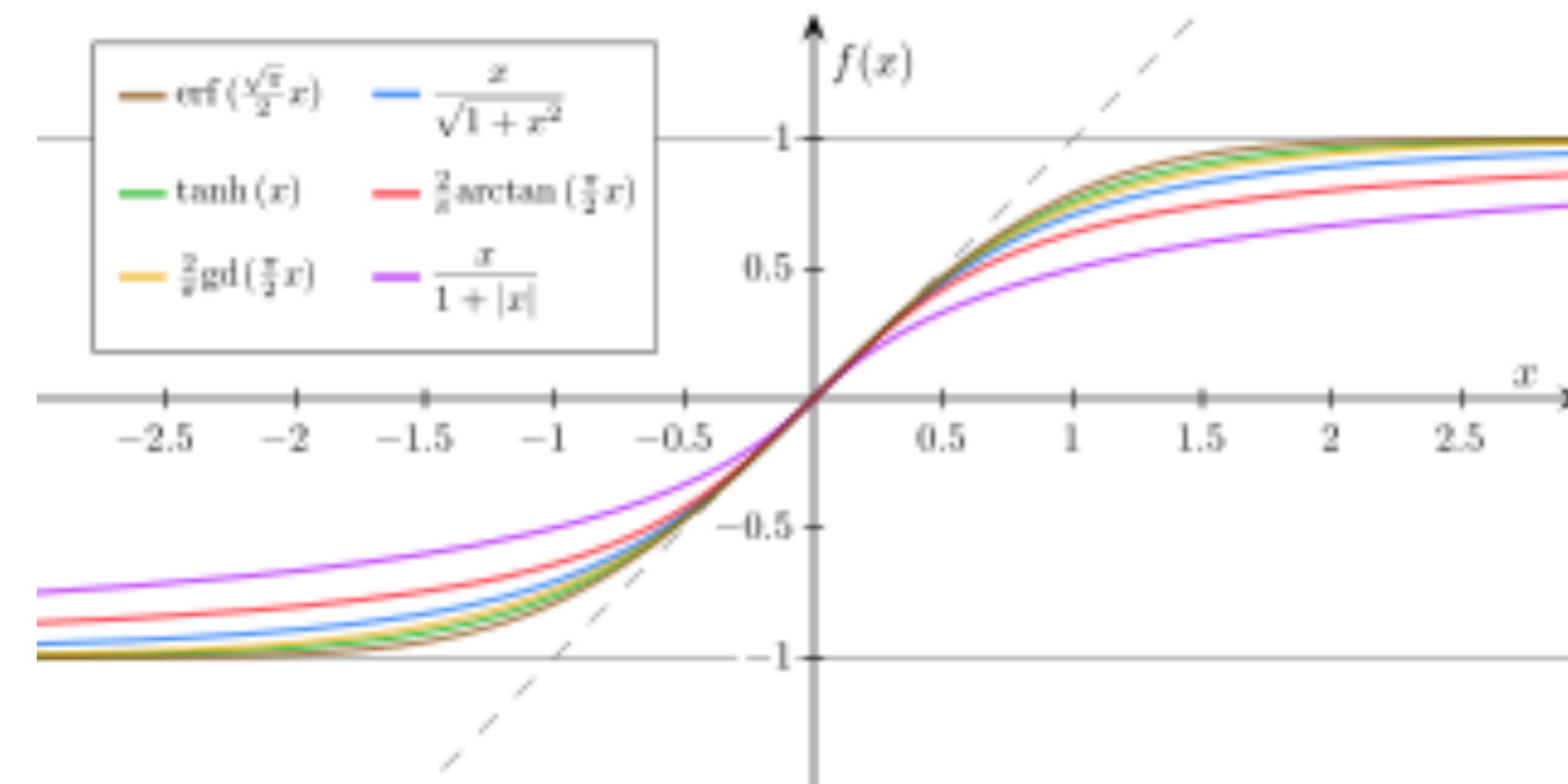
- Data are stored in “tensors”.
 - Basically an N- Dimensional Array with a “shape”
 - shape = (): Scalar
 - shape = (N,): Vector
 - shape = (N,M): Matrix
 - shape = (N₁, N₂, N₃, ..., N_R): Rank R Tensor
 - Inputs: **X**
 - Can be arbitrary shape. Typically first dimension is the example index (usually an “event” or collision in HEP)
 - Example: Let’s say your examples are students, and your data is their age, sex, years at University, undergrad/grad, and department
 - $X = \begin{bmatrix} [20, 0, 2, 0, 4] \\ [25, 1, 2, 1, 3] \\ [23, 0, 0, 1, 3] \end{bmatrix}$, # 20 year old, 0=male, 2=junior, 0=undergrad, 4=computer science
25 year old, 1=female, 2=3rd year, 0=grad, 4=physics
23 year old, 0=male, 2=1st year, 0=grad, 4=physics
 - $X[0] = [20, 0, 2, 0, 4]$: the first students data.
 - $X[0][3] = 0$. This is an undergraduate student
 - Outputs : **Y**
 - Can be arbitrary shape. Typically first dimension is the example index (usually an “event” or collision in HEP)
 - Example: $Y = 0/1$, student does not / does know python

Machine Learning Problem Formulation

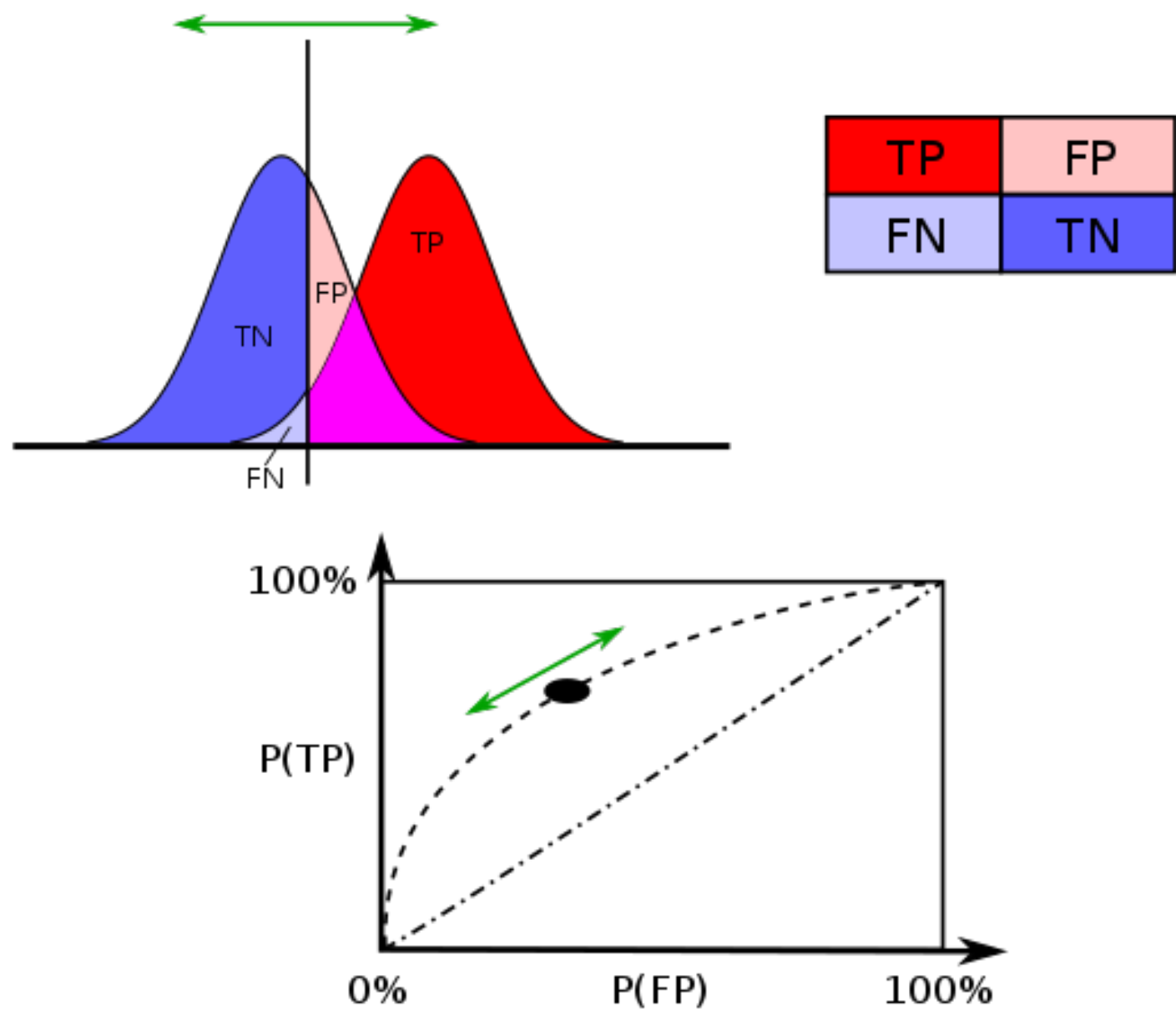
- Split *Datasets*:
 - $(\mathbf{X}_{\text{train}}, \mathbf{Y}_{\text{train}})$ = training dataset
 - $(\mathbf{X}_{\text{test}}, \mathbf{Y}_{\text{test}})$ = test dataset
 - $(\mathbf{X}_{\text{val}}, \mathbf{Y}_{\text{val}})$ = validation dataset
 - (\mathbf{X}) = unlabeled data
- Set *Goal*:
 - Inference algorithm/function $F(\mathbf{X} \mid \mathbf{a}) = \mathbf{Y}_{\text{predict}}$.
 - F can be a heuristic. e.g. if (computer science student) then (student knows python).
 - F can be anything
 - \mathbf{a} are parameters of the function, for Neural Networks, these are weights.
 - Note that in a simple classification problem, $\mathbf{Y}_{\text{train}}$ can be 0 or 1 for any example. But $\mathbf{Y}_{\text{predict}}$ will usually be between 0 and 1.
- **Training**: (for Neural Networks)
 - Optimize (usually a minimization) a *cost function* $F(\mathbf{X} \mid \mathbf{a}) = C(F(\mathbf{X}_{\text{train}} \mid \mathbf{a}), \mathbf{Y}_{\text{train}})$ w.r.t. \mathbf{a}
 - For example, $C = [F(\mathbf{X} \mid \mathbf{a}) - \mathbf{Y}_{\text{train}}]^2$
 - $\mathbf{a}_{\text{trained}}$ = result of training
- **Test**:
 - Compute cost function on test data $C(F(\mathbf{X}_{\text{test}} \mid \mathbf{a}_{\text{trained}}), \mathbf{Y}_{\text{test}})$
 - Determine (e.g. via significance optimization) the cut-off $F(\mathbf{X}_{\text{test}} \mid \mathbf{a}_{\text{trained}}) > c$, e.g. $c=0.5$
 - Other metrics. For example:
 - Select $\mathbf{Y}_{\text{test}}=1$ and see how often $F(\mathbf{X}_{\text{test}} \mid \mathbf{a}_{\text{trained}}) > 0.5$
 - Retrain/test to try/compare different techniques/hyper-parameters
- **Validation**:
 - Assess performance metrics \rightarrow e.g. TPR for $F(\mathbf{X}_{\text{val}} \mid \mathbf{a}_{\text{trained}}) > 0.5$
- **Inference**:
 - $\mathbf{Y}_{\text{predict}} = F(\mathbf{X} \mid \mathbf{a}_{\text{trained}})$

Artificial Neural Network

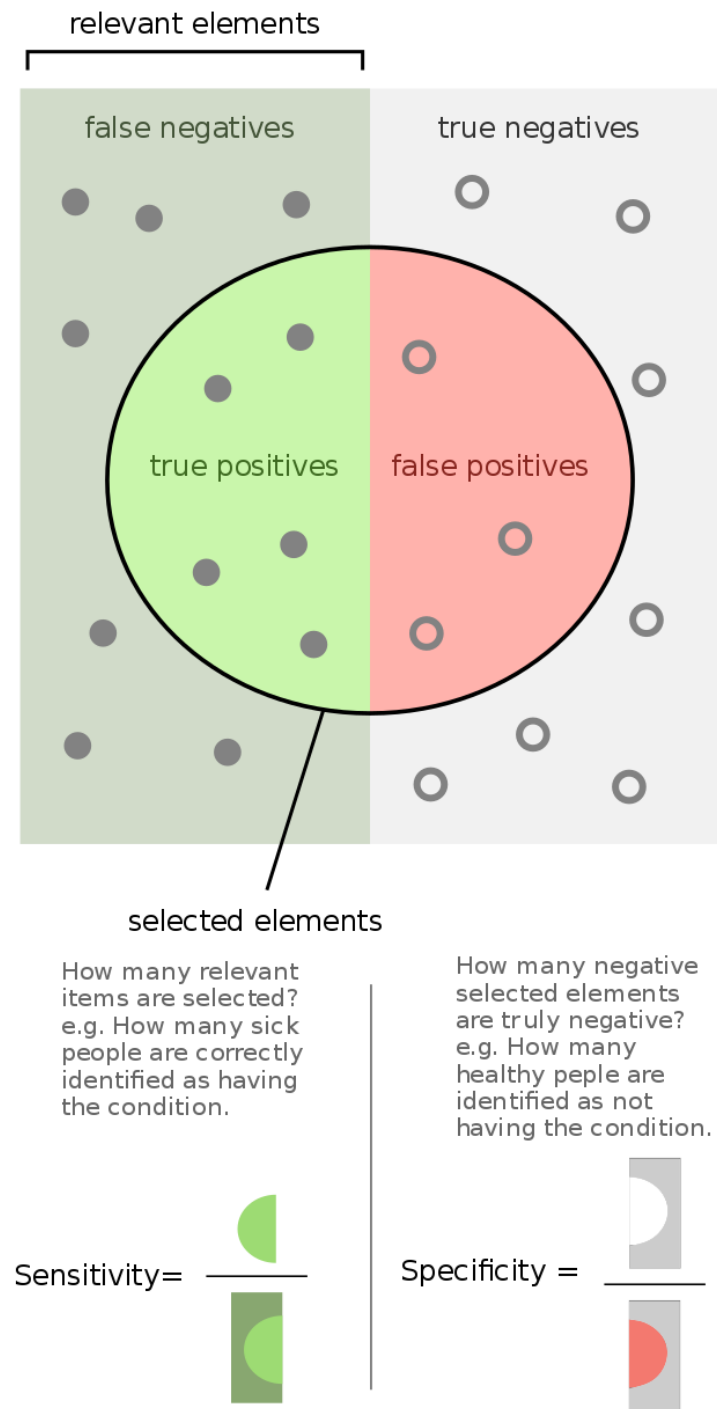
- A simple one layer NN
- $F(\mathbf{X} \mid \mathbf{a} = \mathbf{W}, \mathbf{b}) = f(\mathbf{W}\mathbf{X} + \mathbf{b})$
- \mathbf{W}, \mathbf{b} = “weights”, “biases”
- $f(x)$ = “activation function”
 - Must be non-linear.
- Universal Computation Theorem.



Assessing Performance



TP	FP
FN	TN



		True condition			
Total population		Condition positive	Condition negative	$Prevalence = \frac{\sum \text{Condition positive}}{\sum \text{Total population}}$	$Accuracy (ACC) = \frac{\sum \text{True positive} + \sum \text{True negative}}{\sum \text{Total population}}$
Predicted condition	Predicted condition positive	True positive	False positive, Type I error	$Positive \text{ predictive value (PPV), Precision} = \frac{\sum \text{True positive}}{\sum \text{Predicted condition positive}}$	$False \text{ discovery rate (FDR)} = \frac{\sum \text{False positive}}{\sum \text{Predicted condition positive}}$
	Predicted condition negative	False negative, Type II error	True negative	$False \text{ omission rate (FOR)} = \frac{\sum \text{False negative}}{\sum \text{Predicted condition negative}}$	$Negative \text{ predictive value (NPV)} = \frac{\sum \text{True negative}}{\sum \text{Predicted condition negative}}$
		$True \text{ positive rate (TPR), Recall, Sensitivity, probability of detection, Power} = \frac{\sum \text{True positive}}{\sum \text{Condition positive}}$	$False \text{ positive rate (FPR), Fall-out, probability of false alarm} = \frac{\sum \text{False positive}}{\sum \text{Condition negative}}$	$Positive \text{ likelihood ratio (LR+)} = \frac{TPR}{FPR}$	$Diagnostic \text{ odds ratio (DOR)} = \frac{LR+}{LR-}$
		$False \text{ negative rate (FNR), Miss rate} = \frac{\sum \text{False negative}}{\sum \text{Condition positive}}$	$Specificity (SPC), Selectivity, True negative rate (TNR) = \frac{\sum \text{True negative}}{\sum \text{Condition negative}}$	$Negative \text{ likelihood ratio (LR-)} = \frac{FNR}{TNR}$	$F_1 \text{ score} = 2 \cdot \frac{Precision \cdot Recall}{Precision + Recall}$

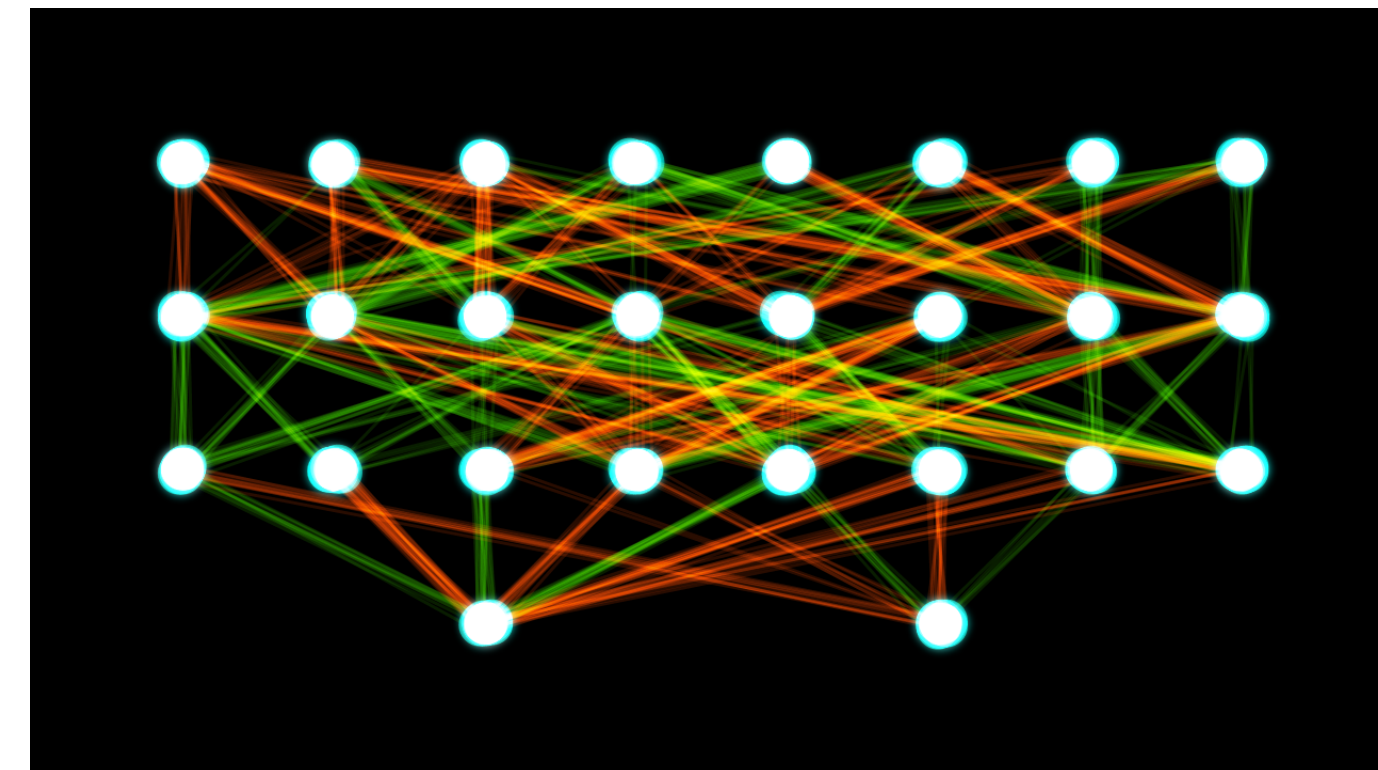
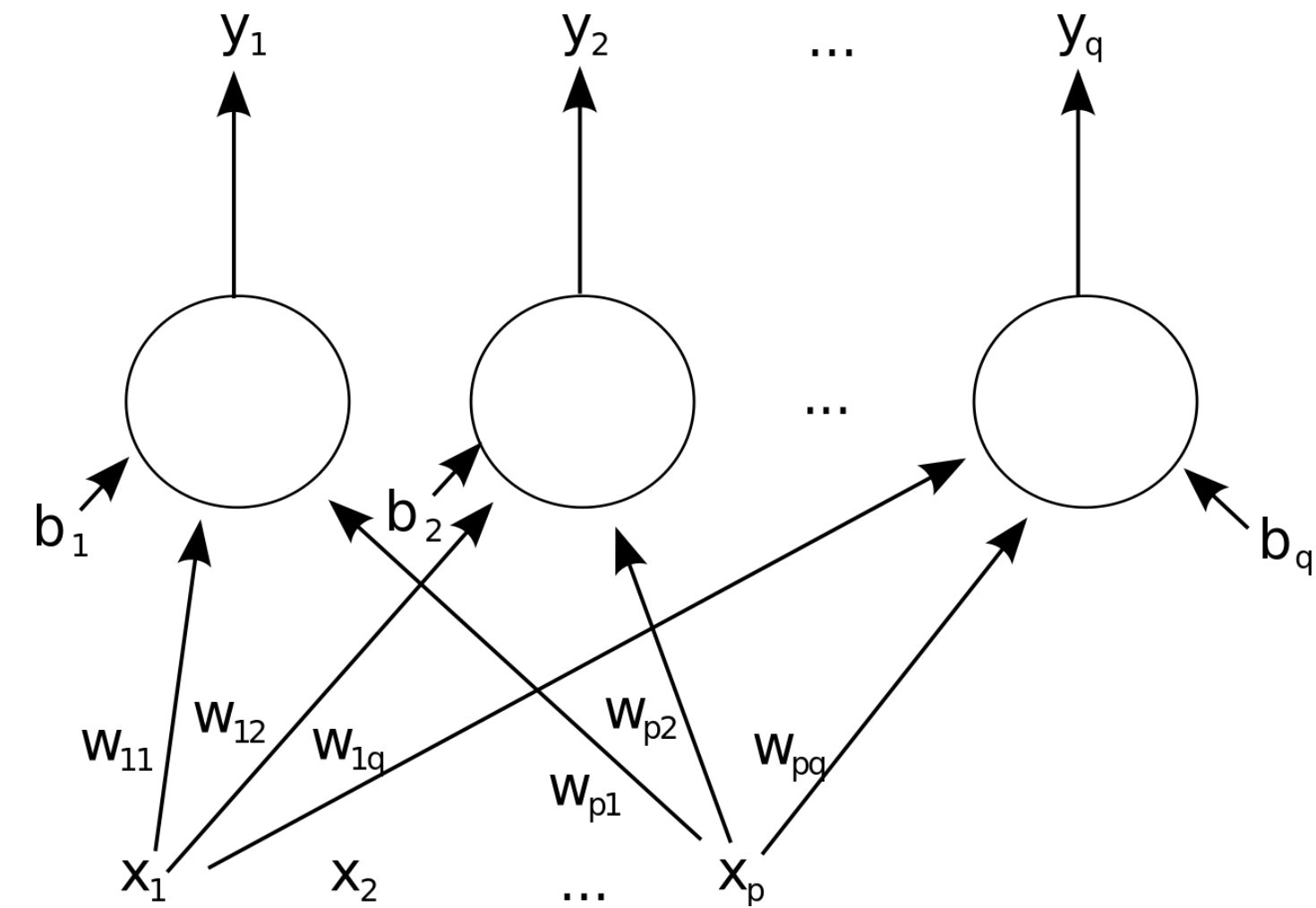
		True condition			
Total population		Condition positive	Condition negative	Prevalence = $\frac{\Sigma \text{Condition positive}}{\Sigma \text{Total population}}$	Accuracy (ACC) = $\frac{\Sigma \text{True positive} + \Sigma \text{True negative}}{\Sigma \text{Total population}}$
Predicted condition	Predicted condition positive	True positive	False positive, Type I error	Positive predictive value (PPV), Precision = $\frac{\Sigma \text{True positive}}{\Sigma \text{Predicted condition positive}}$	False discovery rate (FDR) = $\frac{\Sigma \text{False positive}}{\Sigma \text{Predicted condition positive}}$
	Predicted condition negative	False negative, Type II error	True negative	False omission rate (FOR) = $\frac{\Sigma \text{False negative}}{\Sigma \text{Predicted condition negative}}$	Negative predictive value (NPV) = $\frac{\Sigma \text{True negative}}{\Sigma \text{Predicted condition negative}}$
		True positive rate (TPR), Recall, Sensitivity, probability of detection, Power $= \frac{\Sigma \text{True positive}}{\Sigma \text{Condition positive}}$	False positive rate (FPR), Fall-out, probability of false alarm $= \frac{\Sigma \text{False positive}}{\Sigma \text{Condition negative}}$	Positive likelihood ratio (LR+) = $\frac{\text{TPR}}{\text{FPR}}$	Diagnostic odds ratio (DOR) = $\frac{\text{LR+}}{\text{LR-}}$ F ₁ score = $2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$
		False negative rate (FNR), Miss rate $= \frac{\Sigma \text{False negative}}{\Sigma \text{Condition positive}}$	Specificity (SPC), Selectivity, True negative rate (TNR) = $\frac{\Sigma \text{True negative}}{\Sigma \text{Condition negative}}$	Negative likelihood ratio (LR-) = $\frac{\text{FNR}}{\text{TNR}}$	

- Condition Positive/Negative → Ground Truth
- Predicted Condition Positive/Negative → From ML
- Prevalence → Fraction where Truth=Positive in Population
 - Training Population and Inference Populations are generally different.
 -

Deep Learning

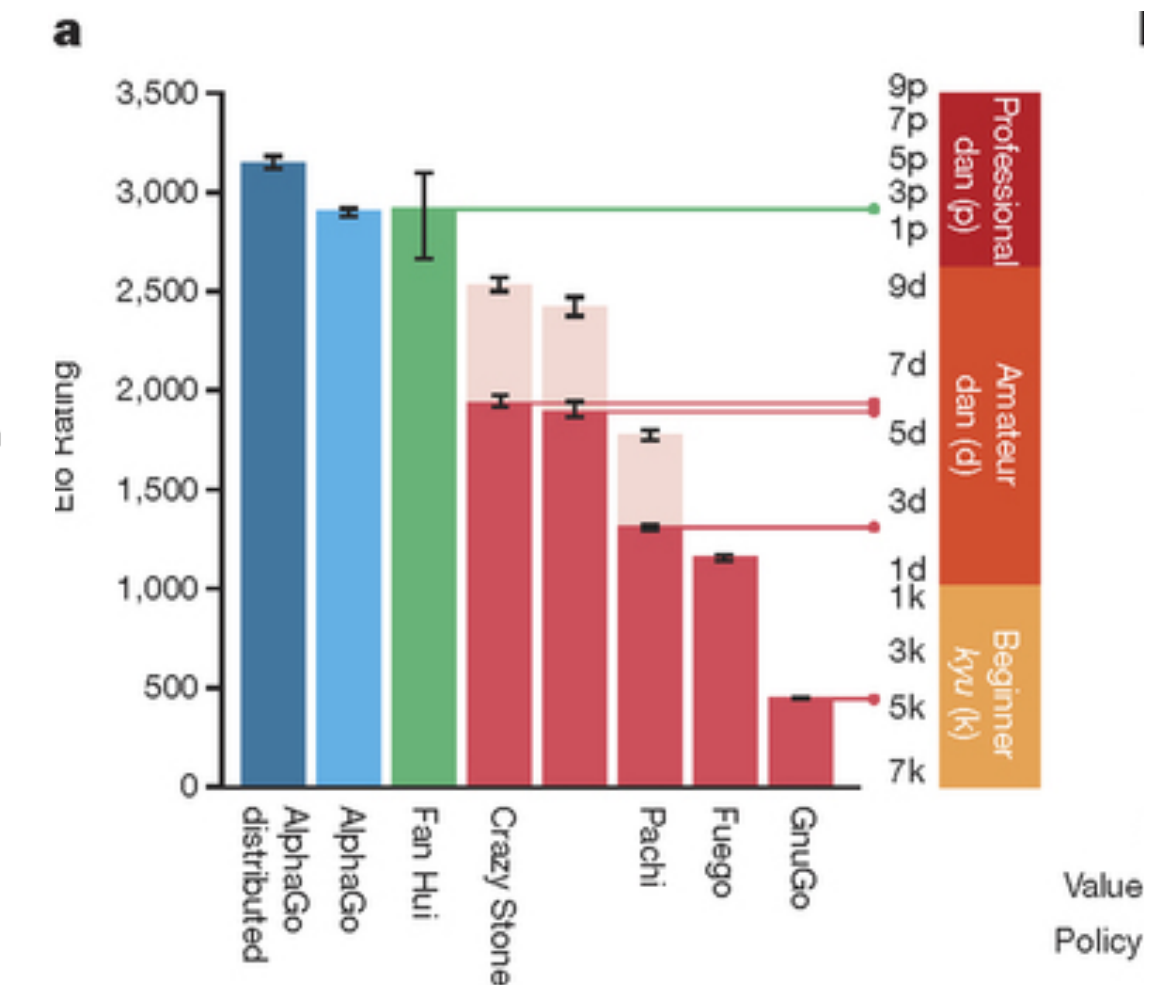
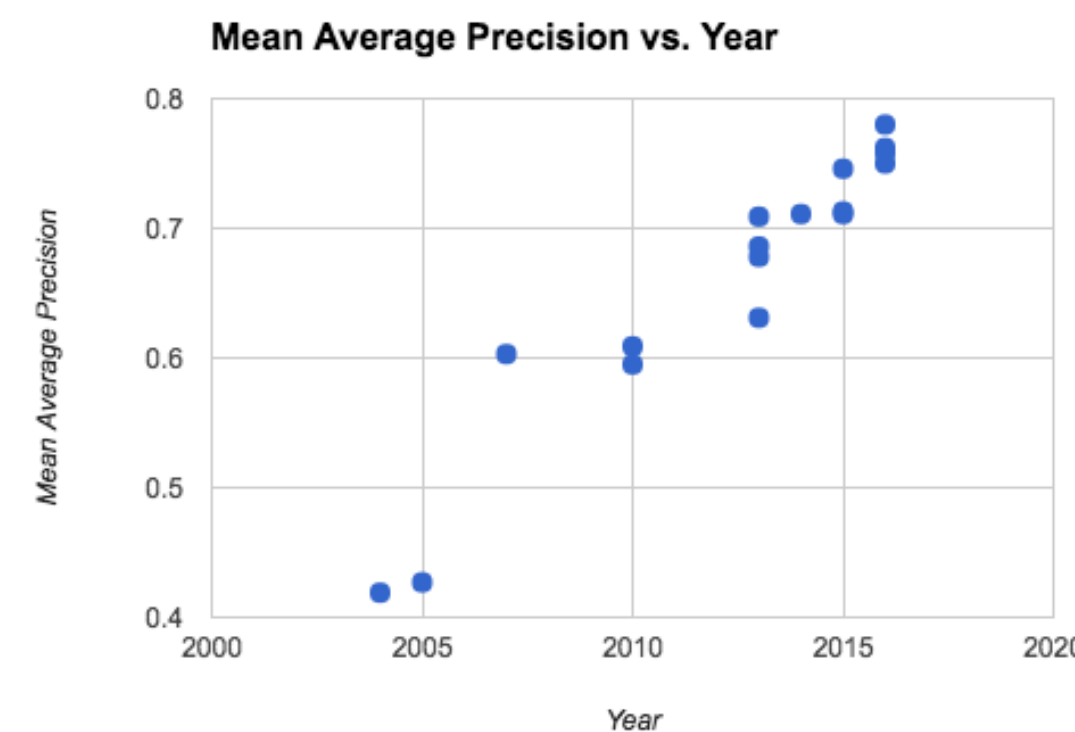
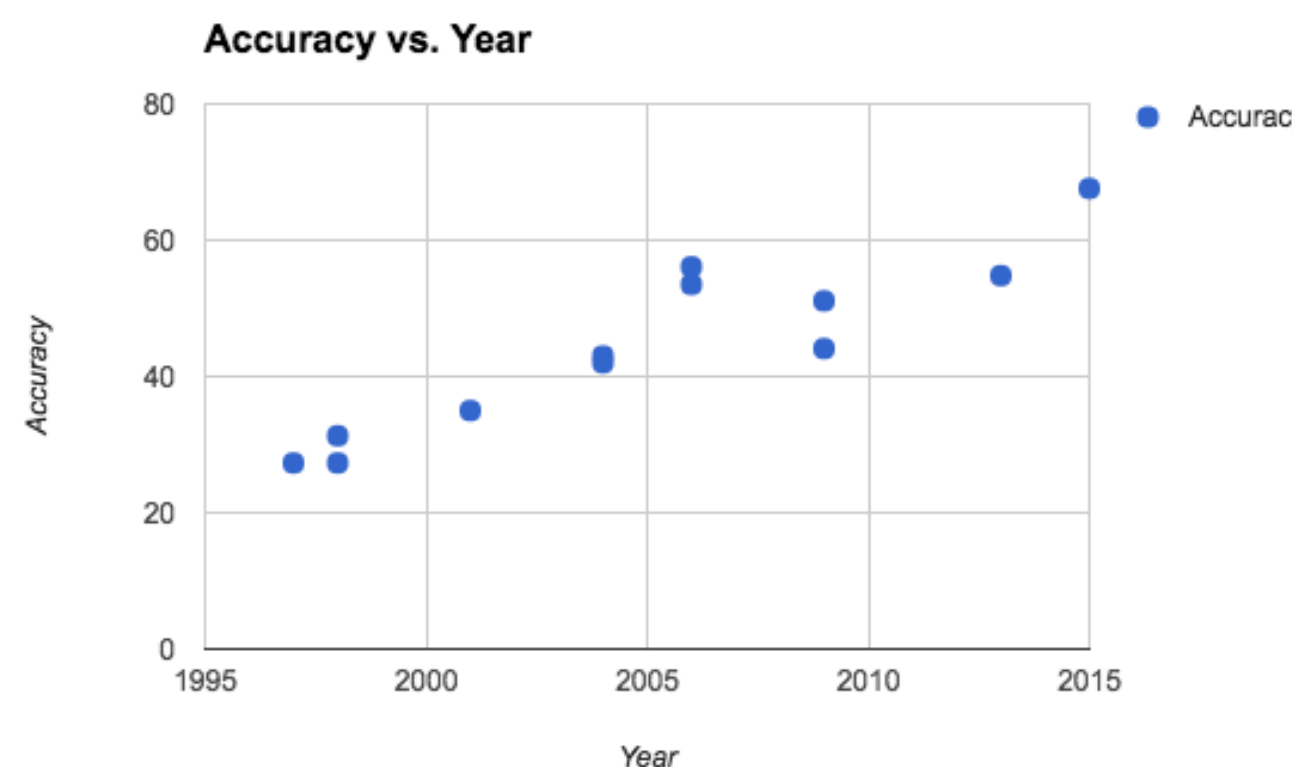
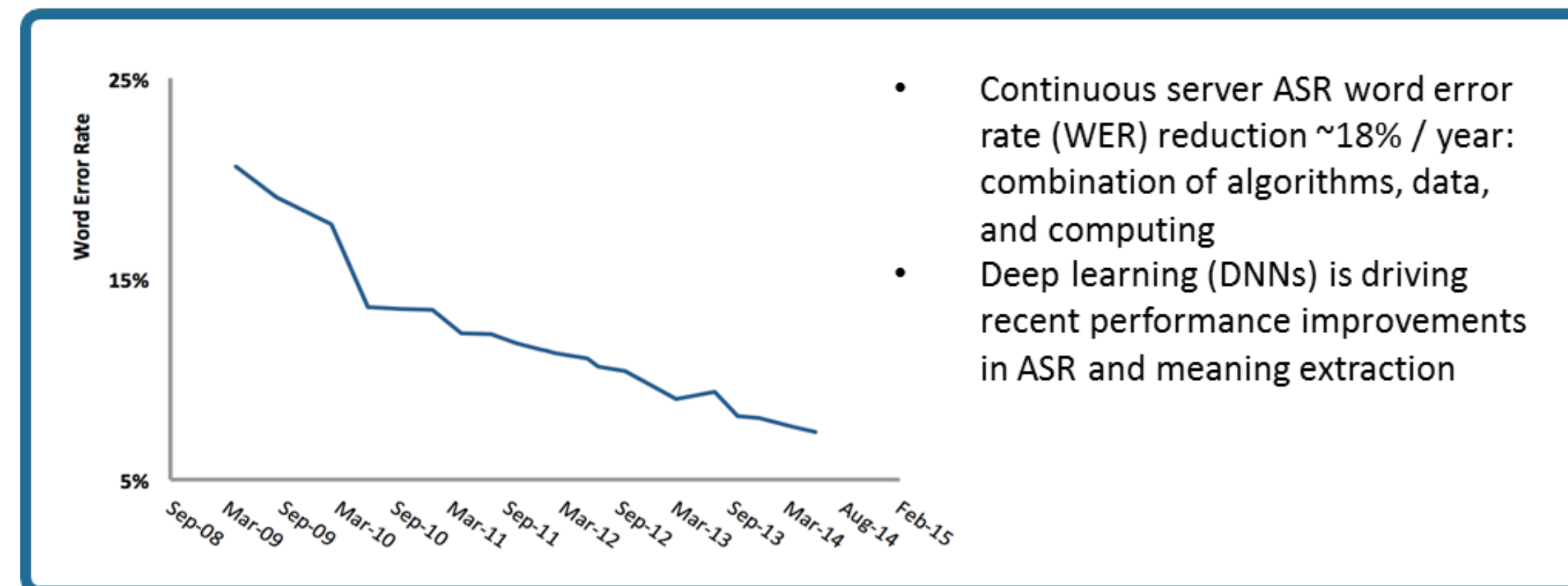
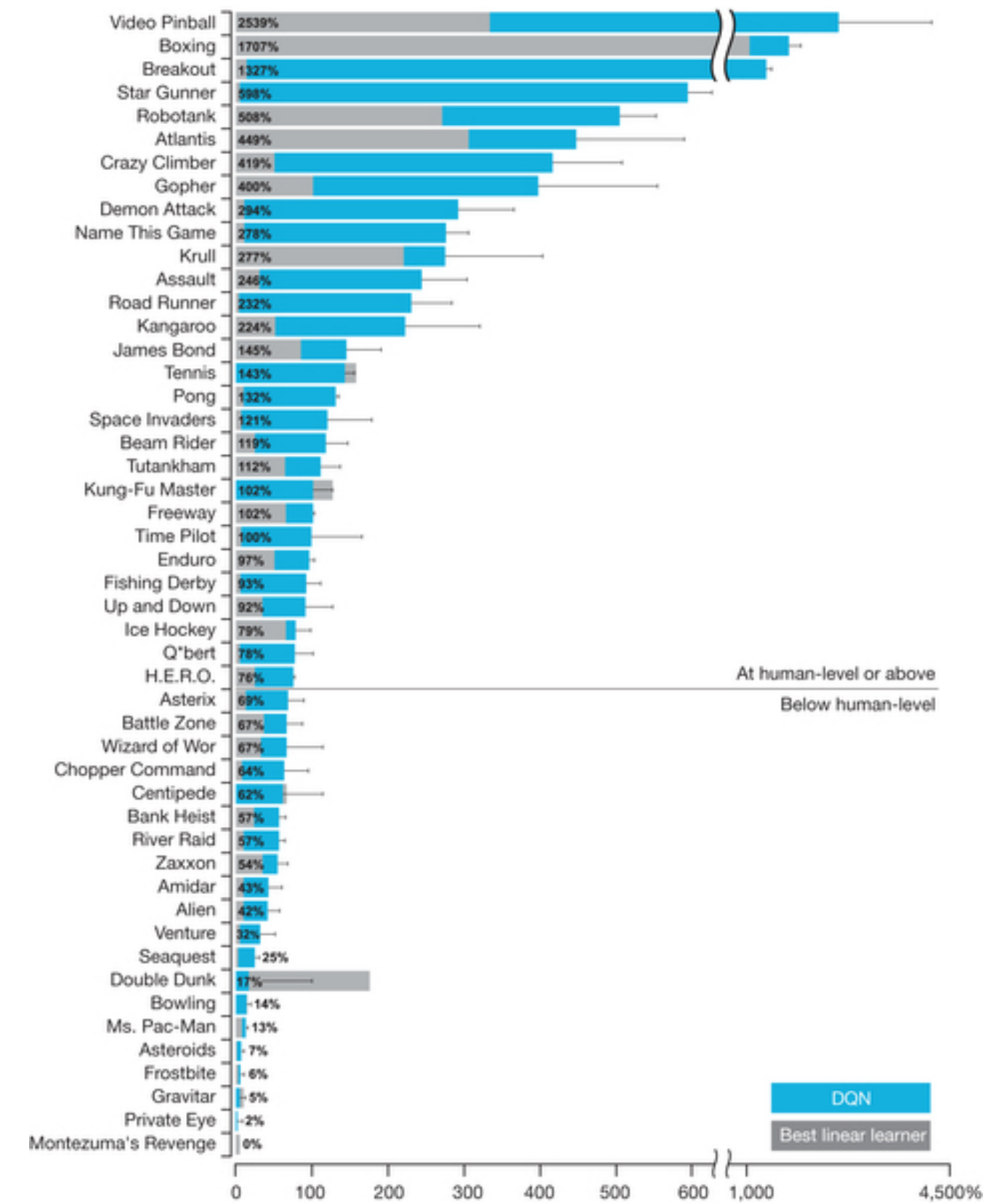
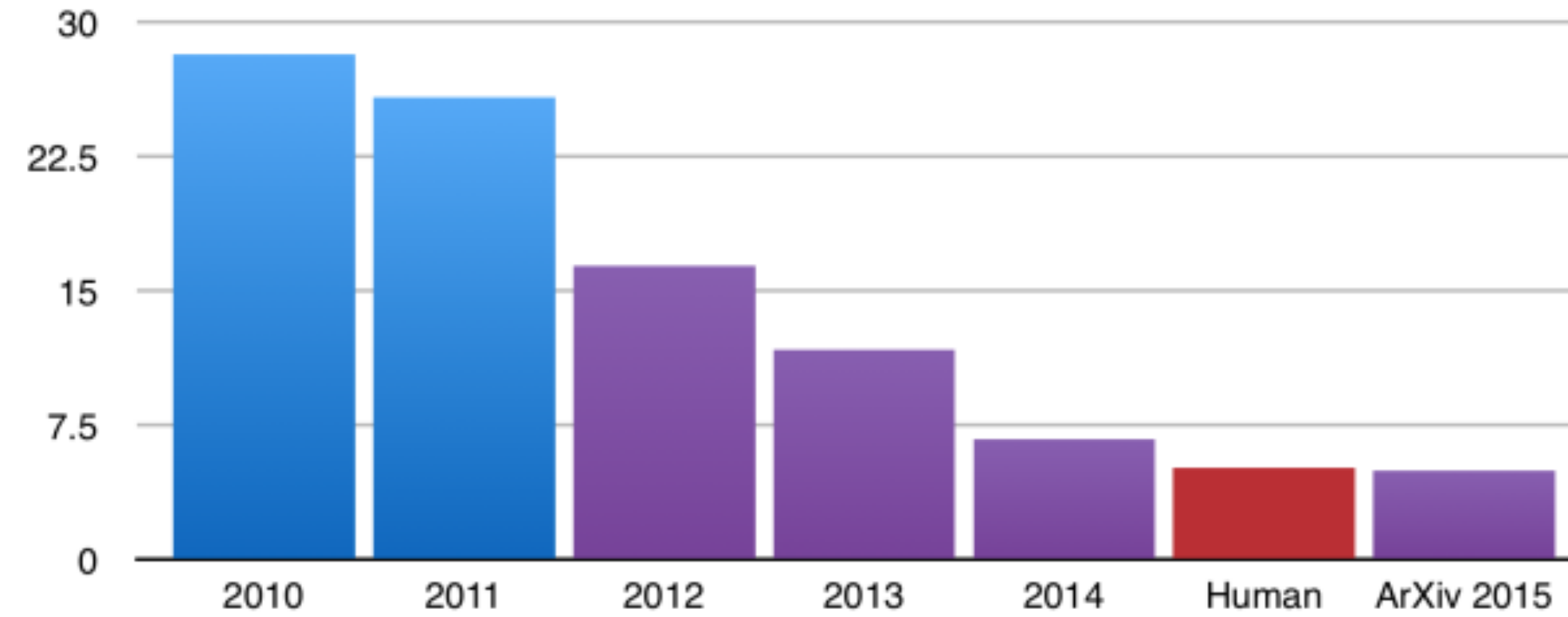
Artificial Neural Networks

- **Biologically inspired computation**, (first attempts in 1943)
 - **Probabilistic Inference**: e.g. signal vs background
 - **Universal Computation Theorem** (1989)
- Multi-layer (**Deep**) Neural Networks:
 - Not a new idea (1965), just impractical to train. **Vanishing Gradient problem** (1991)
 - Solutions:
 - New techniques: e.g. better activation or layer-wise training
 - **More training**: big training datasets and lots of computation ... **big data and GPUs**
 - **Deep Learning Renaissance**. First DNN in HEP (2014).
 - **Amazing Feats**: Audio/Image/Video recognition, captioning, and generation. Text (sentiment) analysis. Language Translation. Video game playing agents.
 - **Rich field**: Variety of architectures, techniques, and applications.



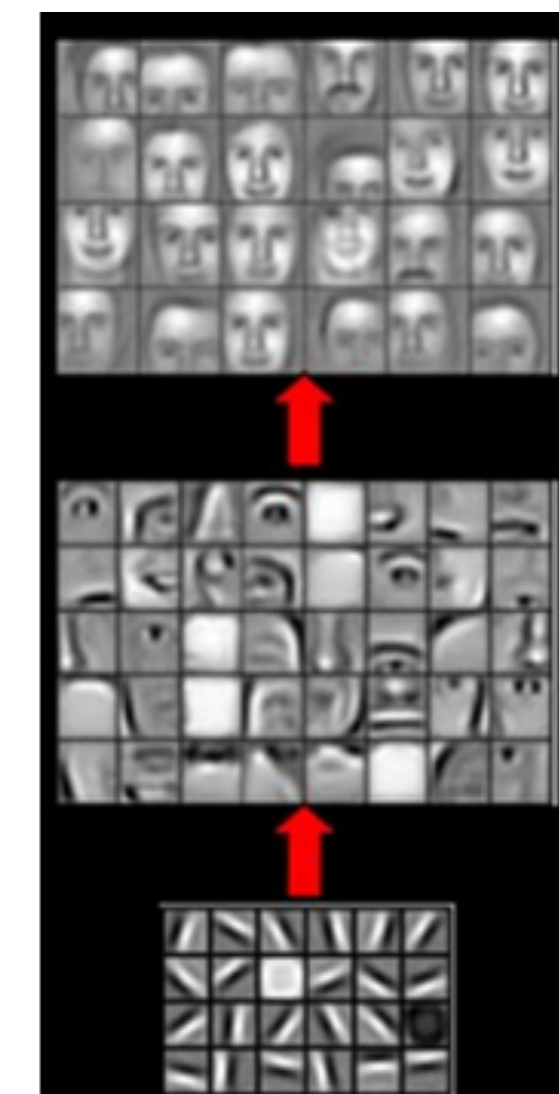
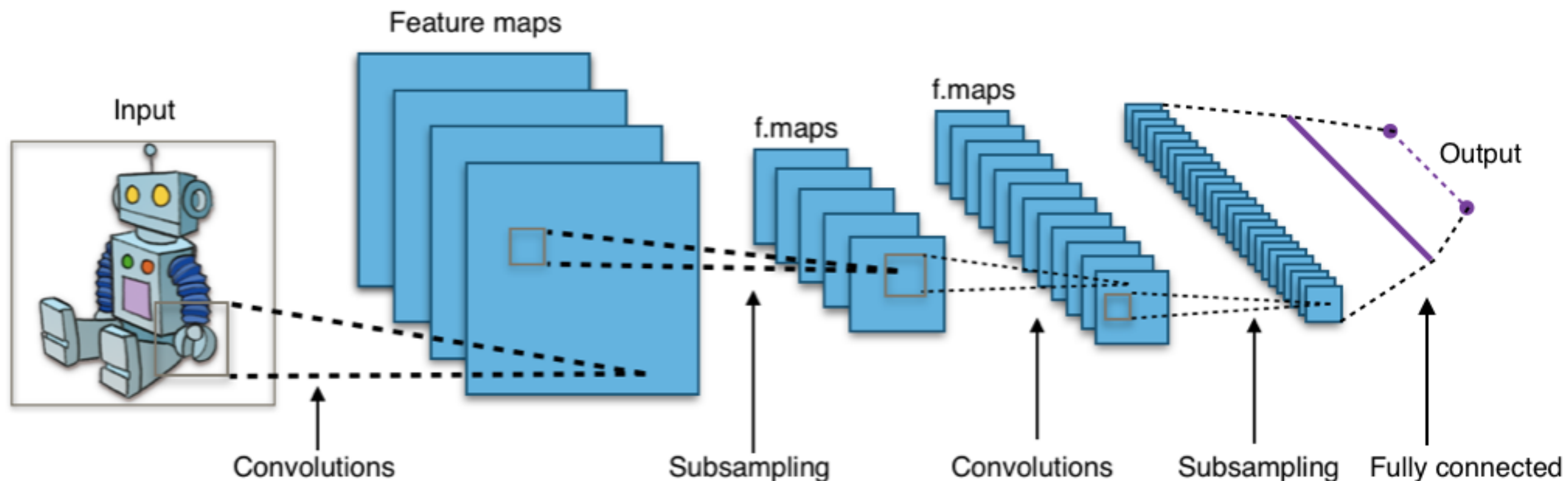
Images from Wikipedia

ILSVRC top-5 error on ImageNet

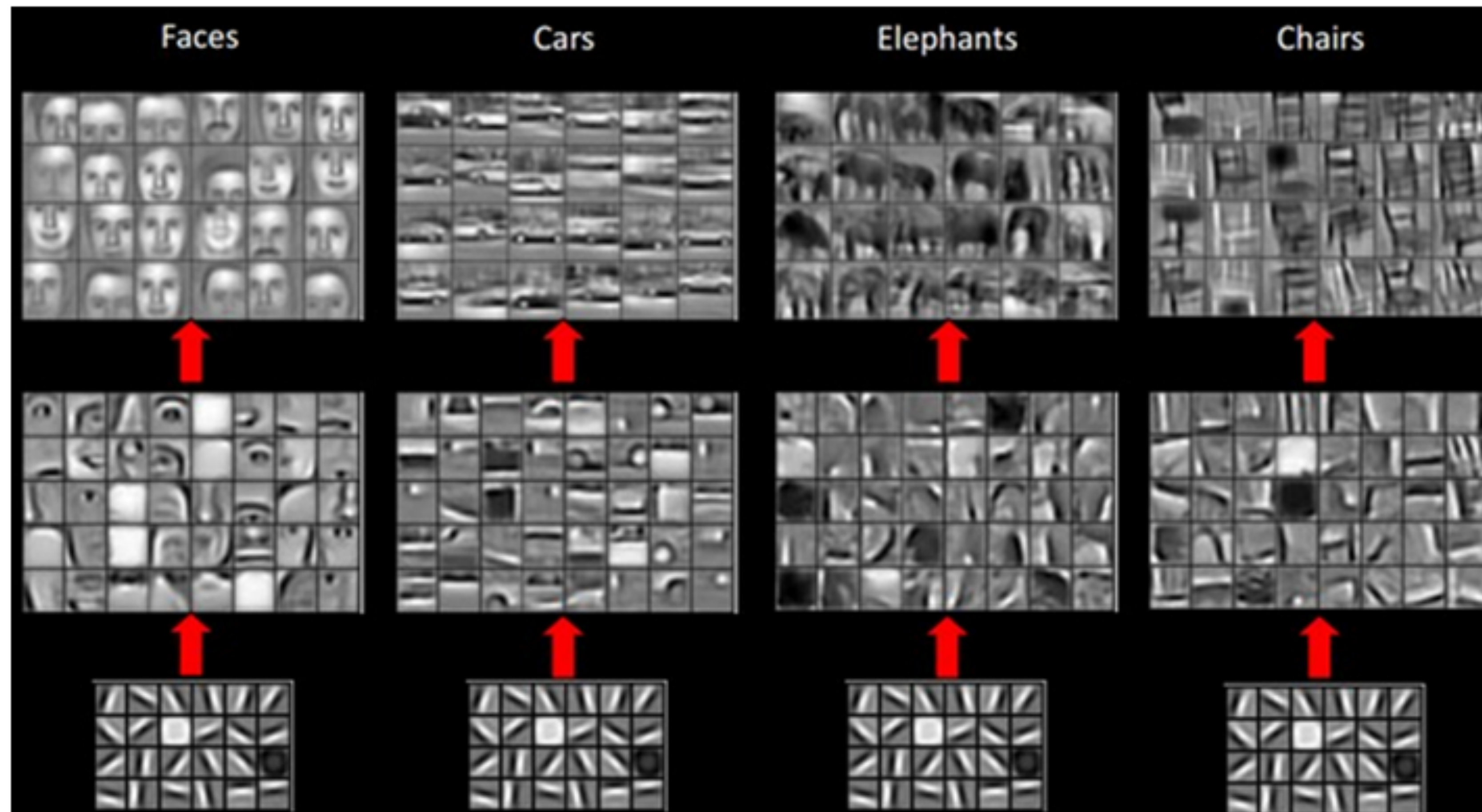


Feature Learning

- **Feature Engineering**: e.g. Event Reconstruction ~ Feature Extraction, Pattern Recognition, Fitting, ...
- Deep Neural Networks can **Learn Features** from **raw data**.
- Example: **Convolutional Neural Networks** - Inspired by visual cortex
 - **Input**: Raw data... for example 1D = Audio, 2D = Images, 3D = Video
 - **Convolutions** ~ learned feature detectors
 - **Feature Maps**
 - **Pooling** - dimension reduction / invariance
 - **Stack**: Deeper layers recognize higher level concepts.

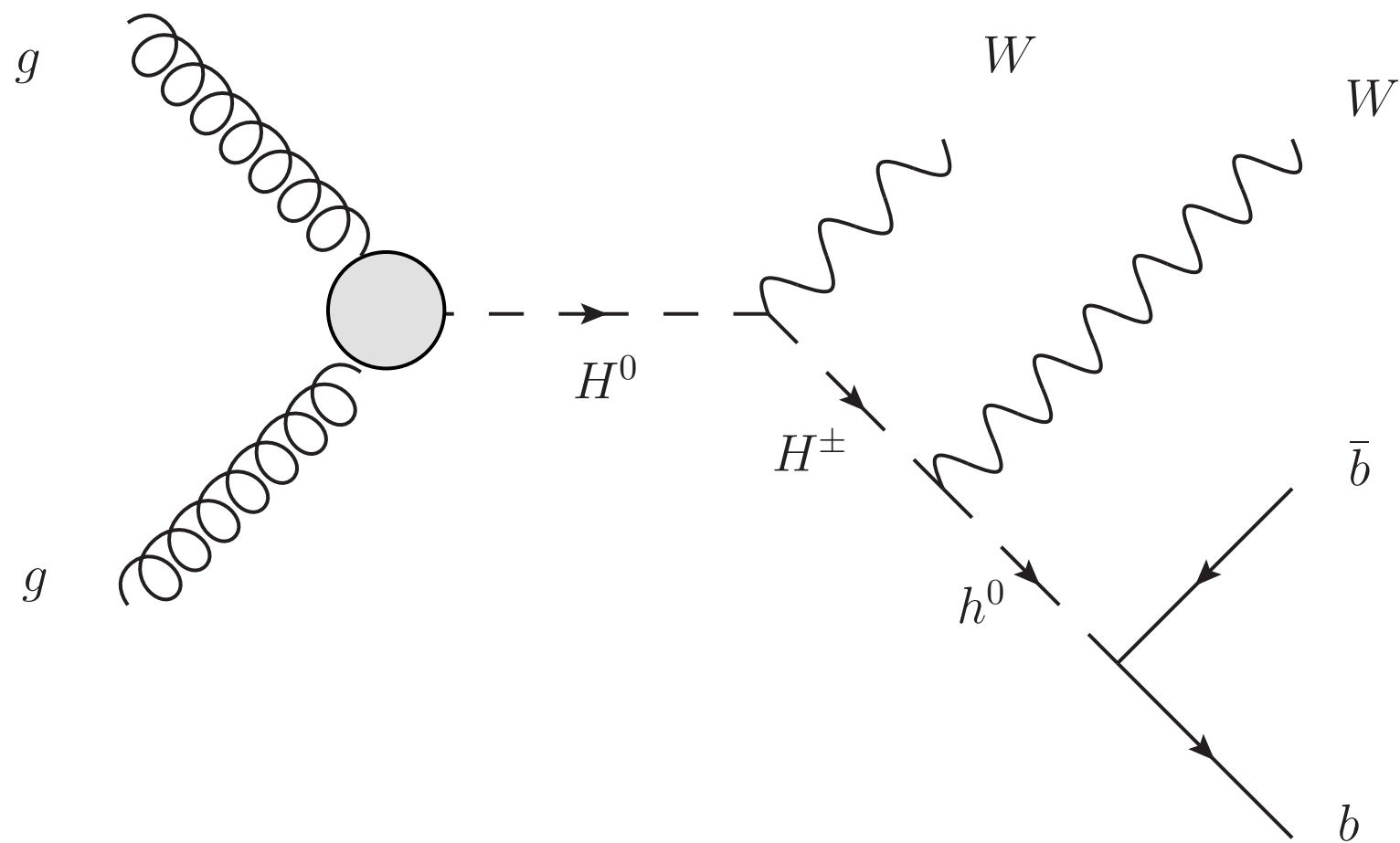


Deep Neutral Networks

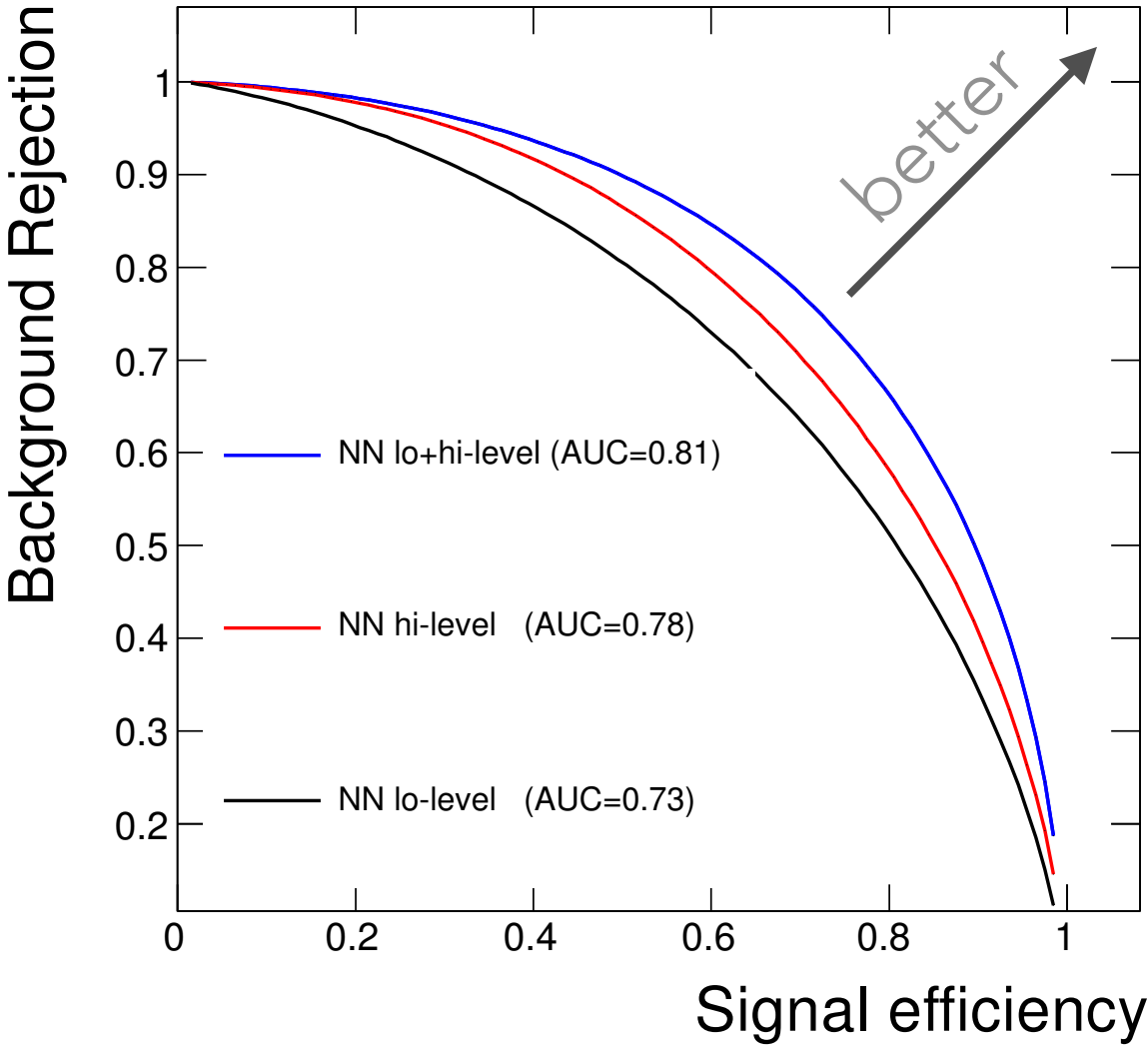


DEEP LEARNING IN HEP

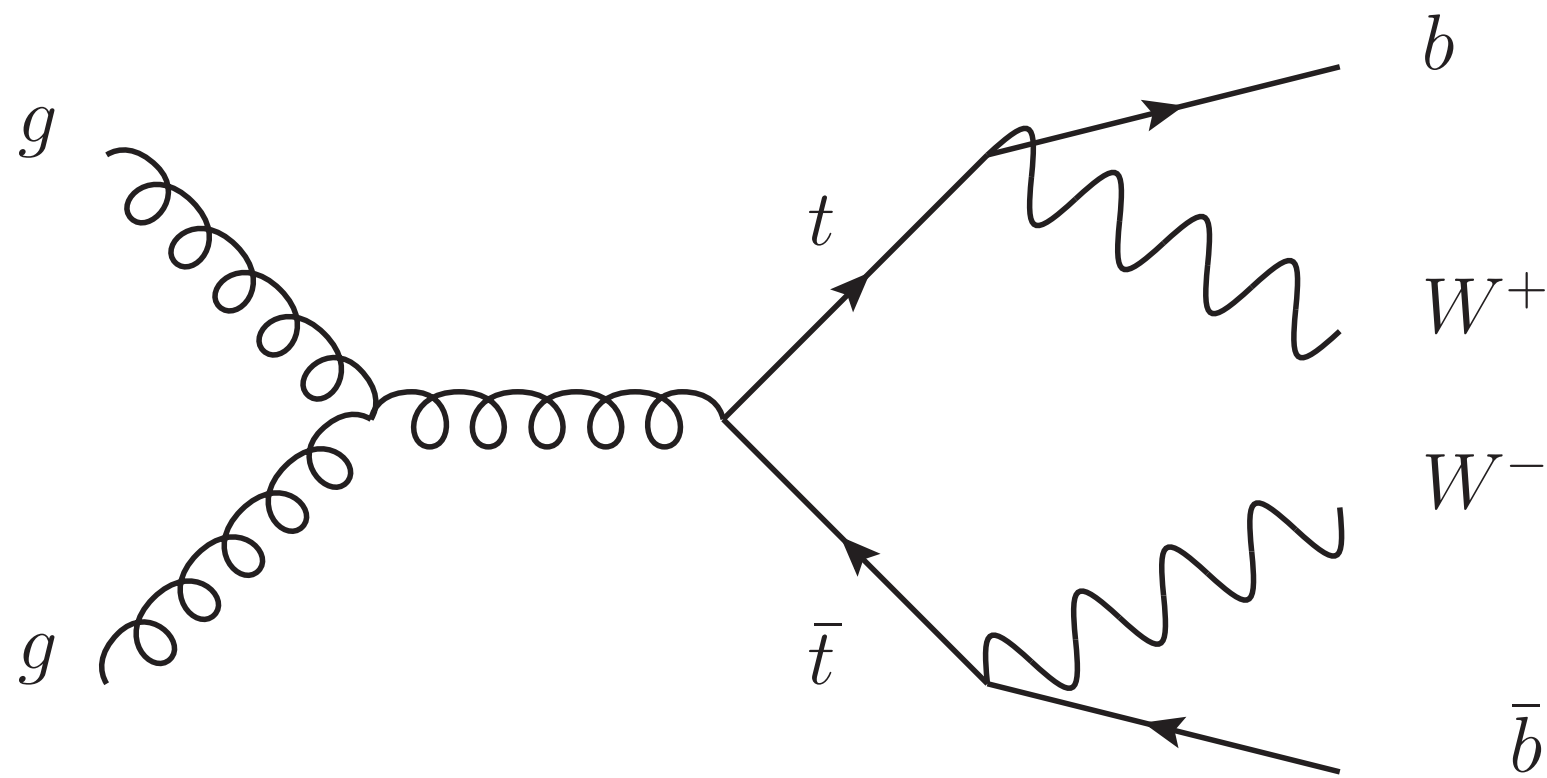
Baldi, Sadowski, Whiteson
arxiv:1402.4735



(a)



(a)



(b)

