

Cognitive Neuroscience for AI Developers

Week 4– Structure and Function of the Nervous System: Plasticity + Neurodevo

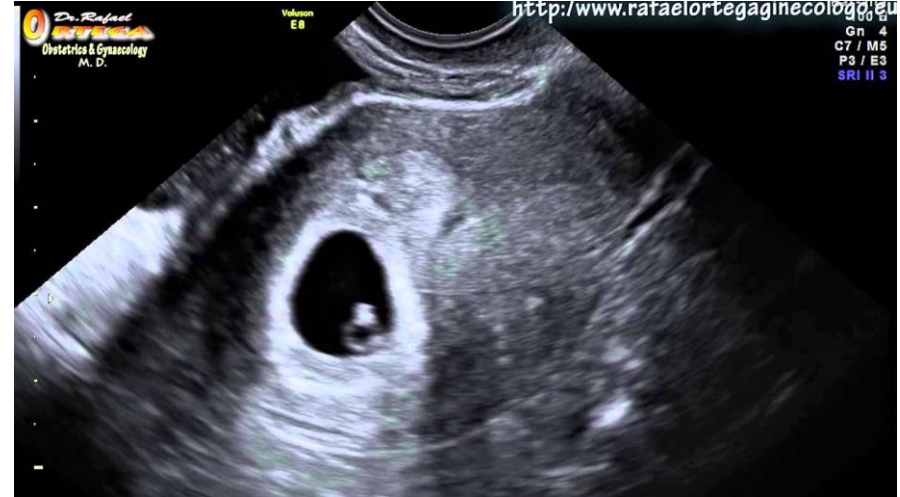


The brain



How does it emerge?

Back in the days...

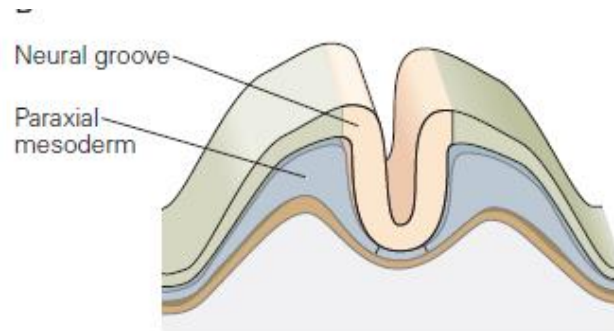
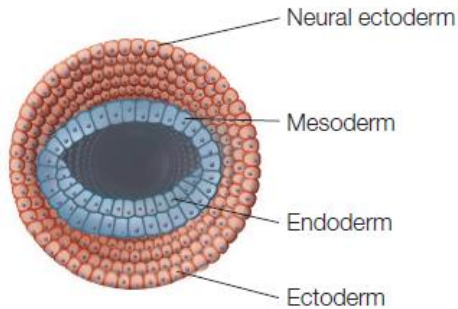
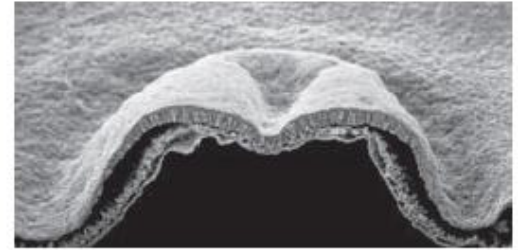
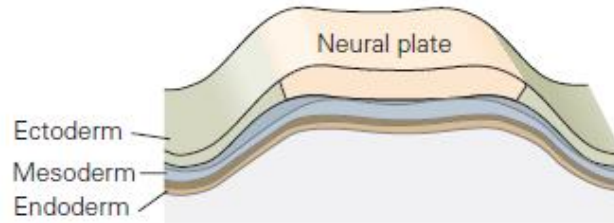


© Pauline Breijer

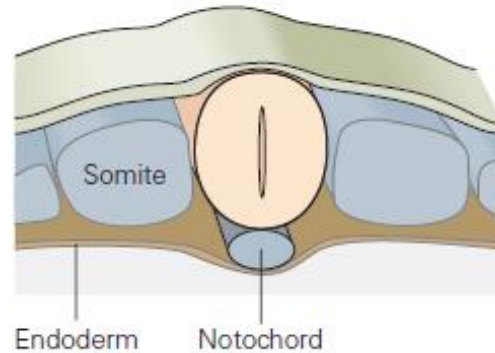
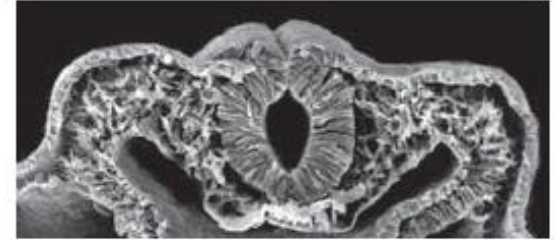
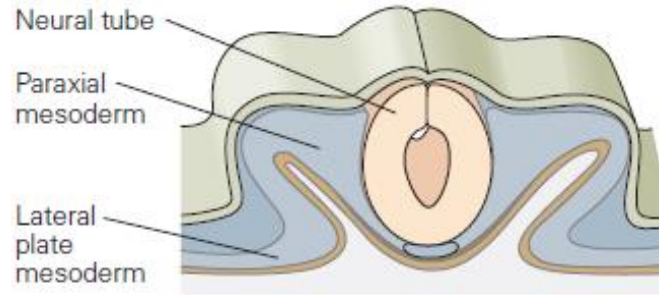
What is happening there?

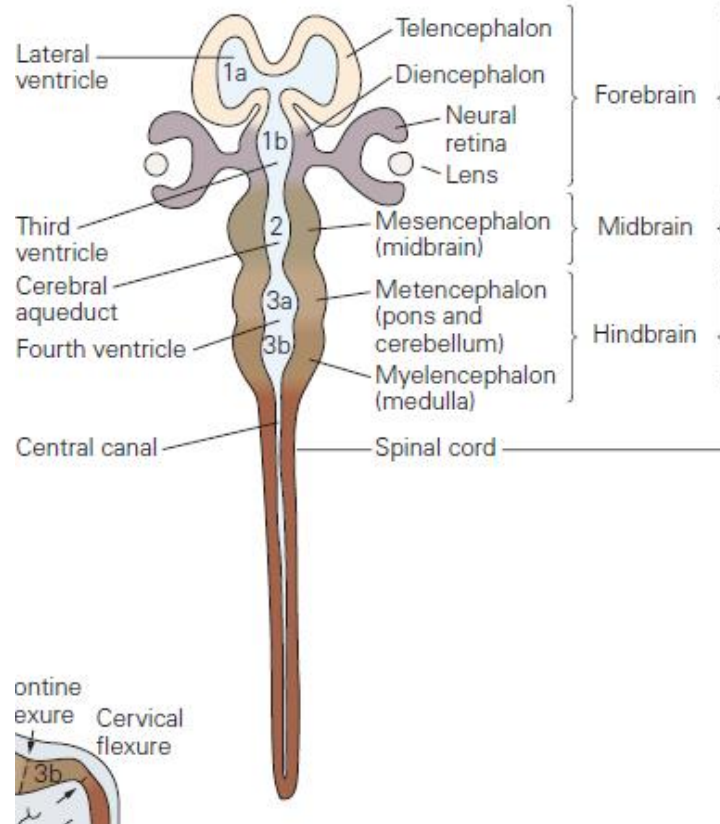
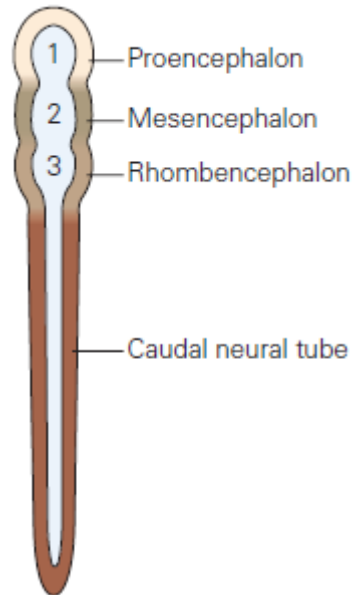


A

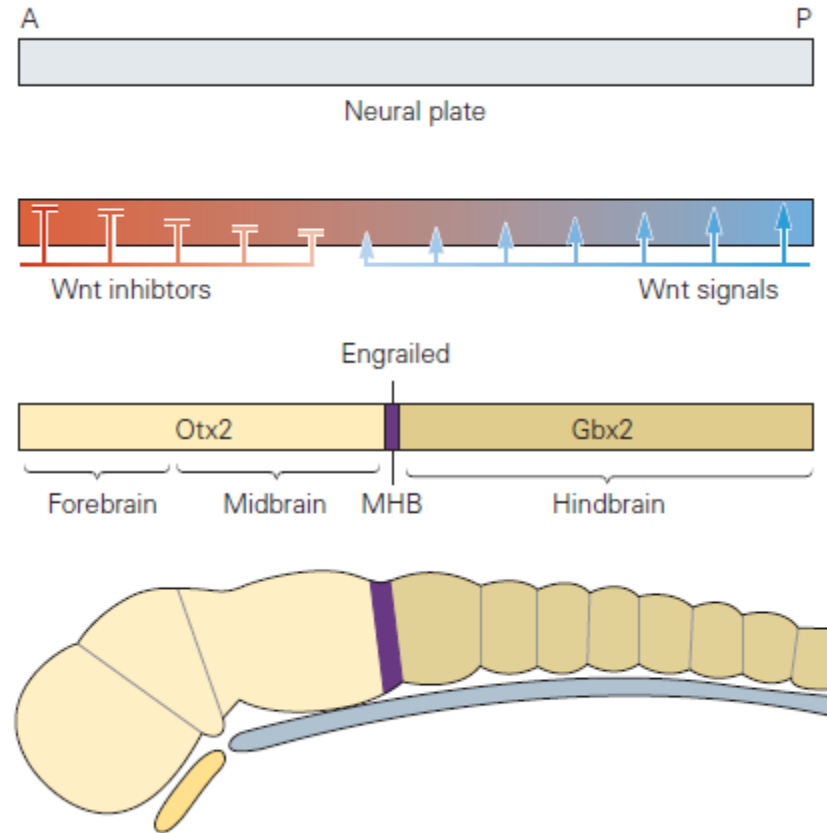


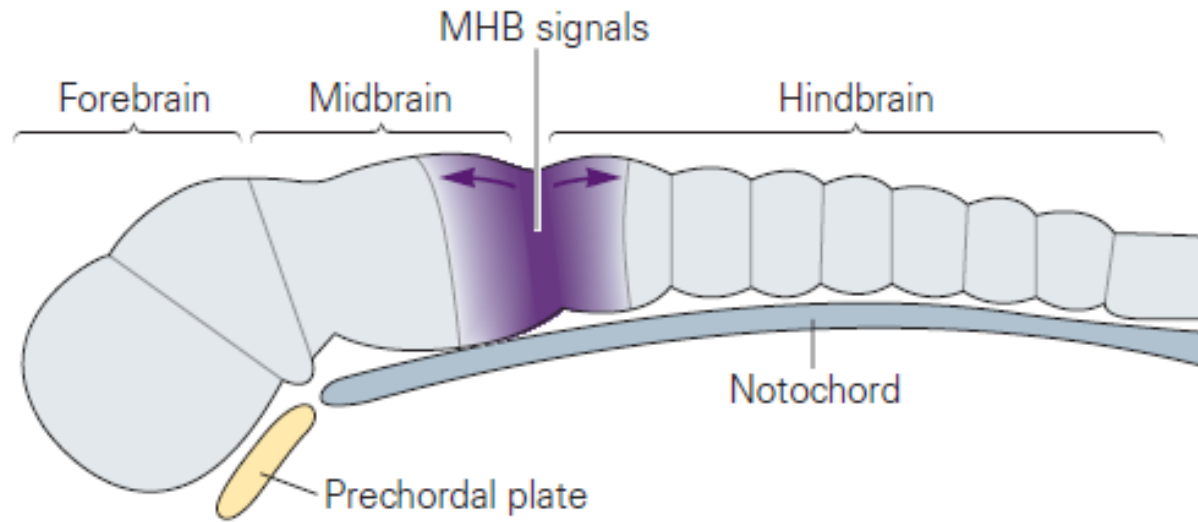
What is happening there?



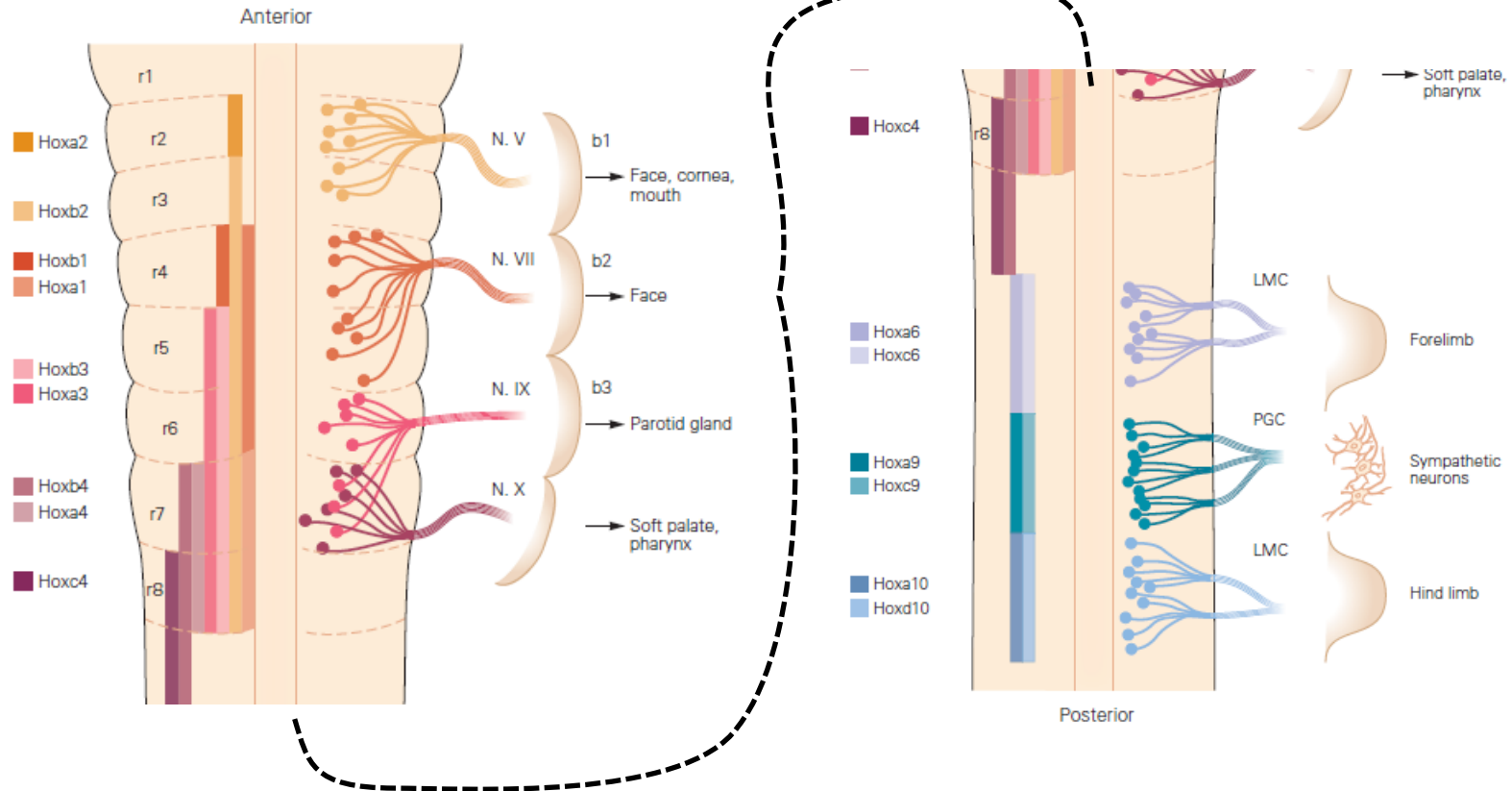


Signaling pathways define neural development

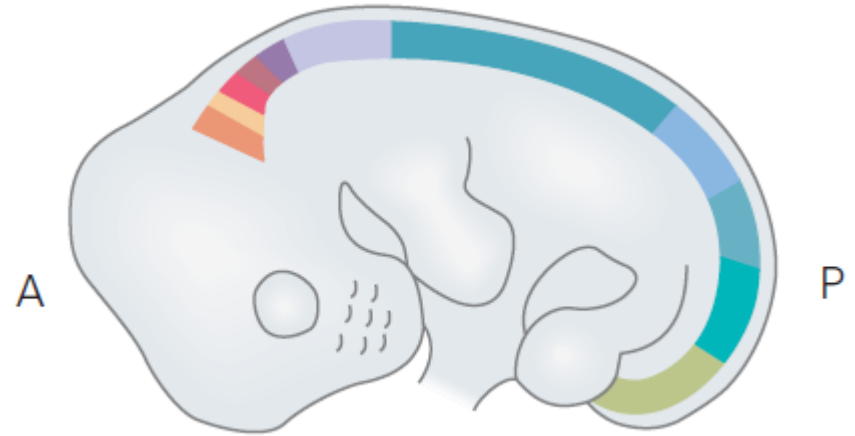
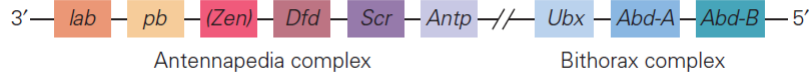
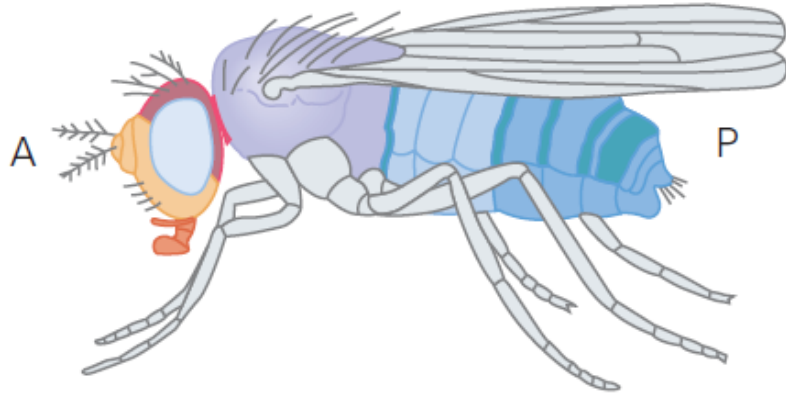




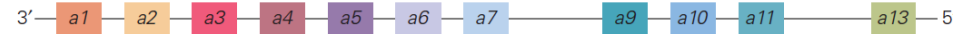
Hox genes determines motor neurons



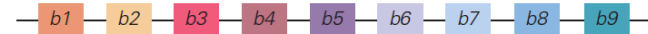
Positional development is highly conserved



Hoxa, chromosome 6



Hoxb, chromosome 11



Hoxc, chromosome 15



Hoxd, chromosome 2



Mammals have very similar development

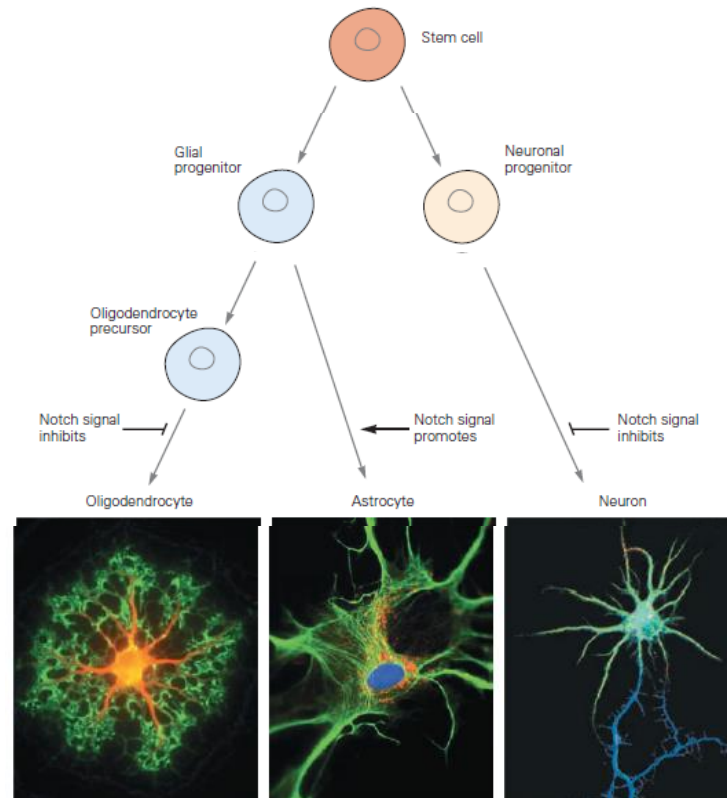


Human

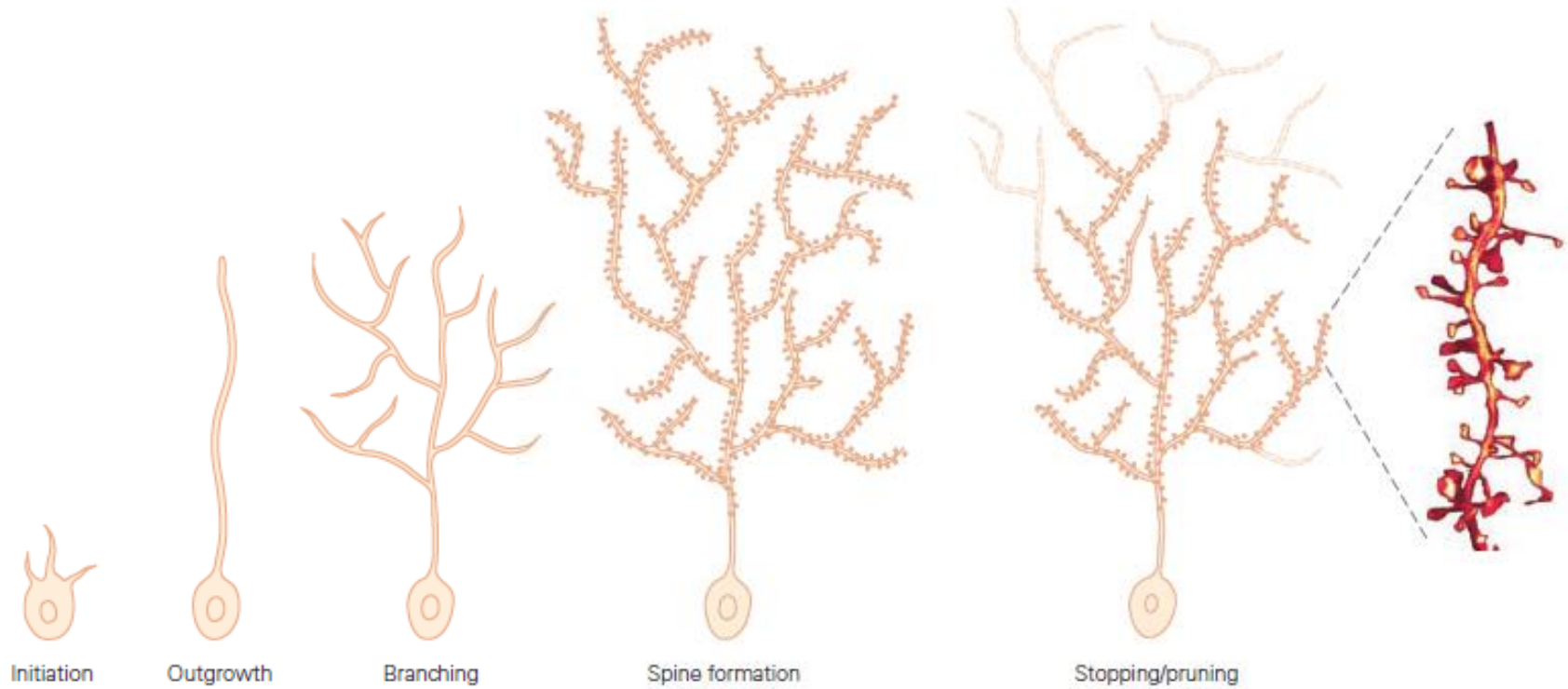


Pig

Brain cells have a common ancestor




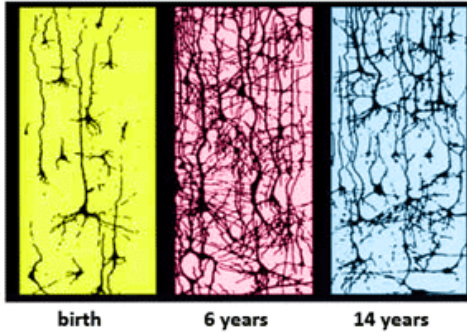
Neuron maturation



Pruning over lifetime

Experience Shapes Brain Architecture by Over-Production Followed by Pruning

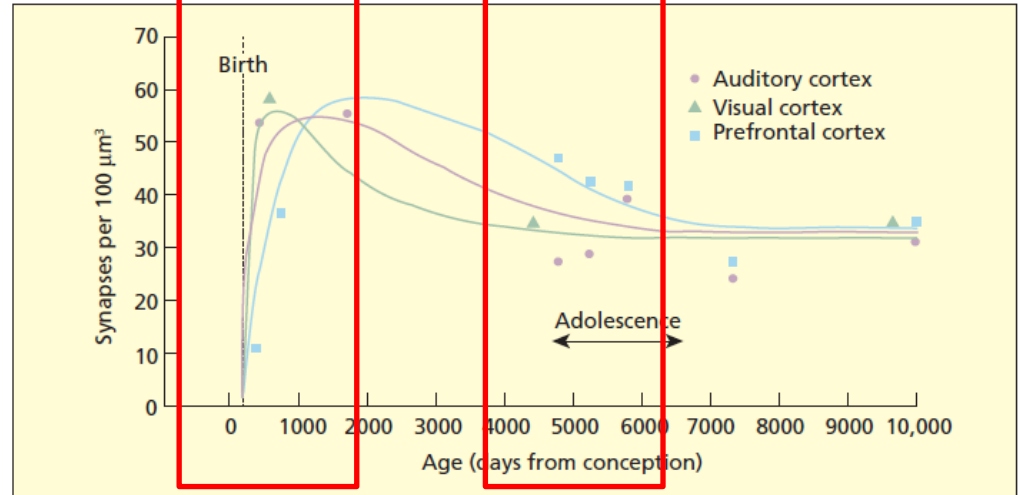
Center on the Developing Child  HARVARD UNIVERSITY



Source: Shonkoff, J. P. (2008) **

WIRING

RE-WIRING



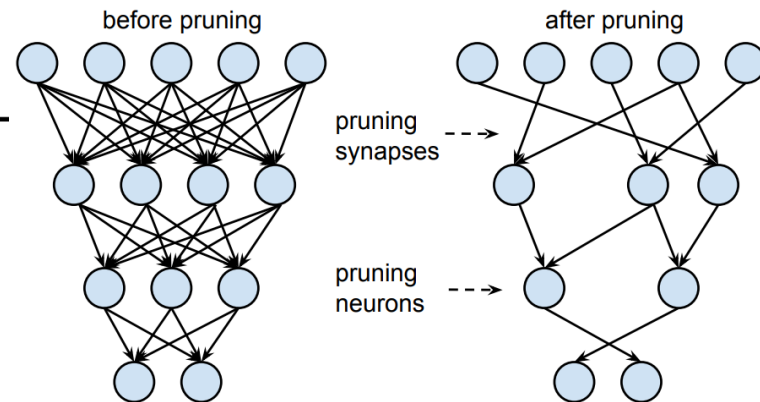
From Huttenlocher and Dabholkar, 1997. Reprinted with permission of John Wiley & Sons Inc.

Optimal Brain Damage

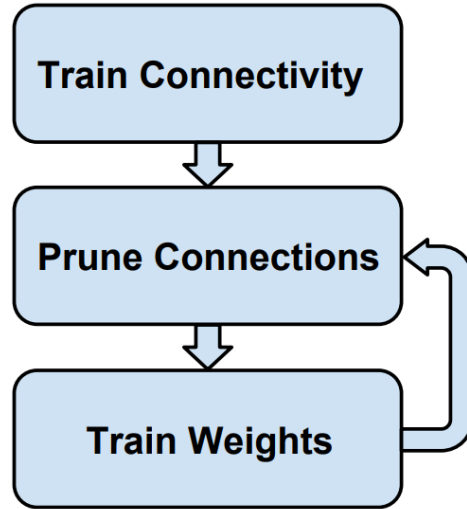
Yann Le Cun, John S. Denker and Sara A. Solla
AT&T Bell Laboratories, Holmdel, N. J. 07733

ABSTRACT

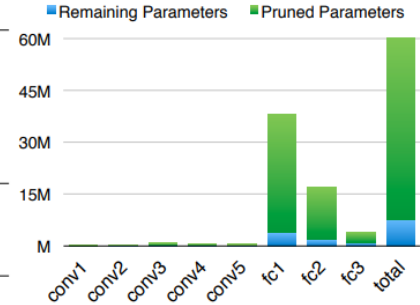
We have used information-theoretic ideas to derive a class of practical and ~~nearly optimal schemes for adapting the size of a neural network~~. By removing unimportant weights from a network, several improvements can be expected: better generalization, fewer training examples required, and improved speed of learning and/or classification. The basic idea is to use second-derivative information to make a tradeoff between network complexity and training set error. Experiments confirm the usefulness of the methods on a real-world application.



Pruning synapses:
making network sparse
Pruning neurons:
Making network dense

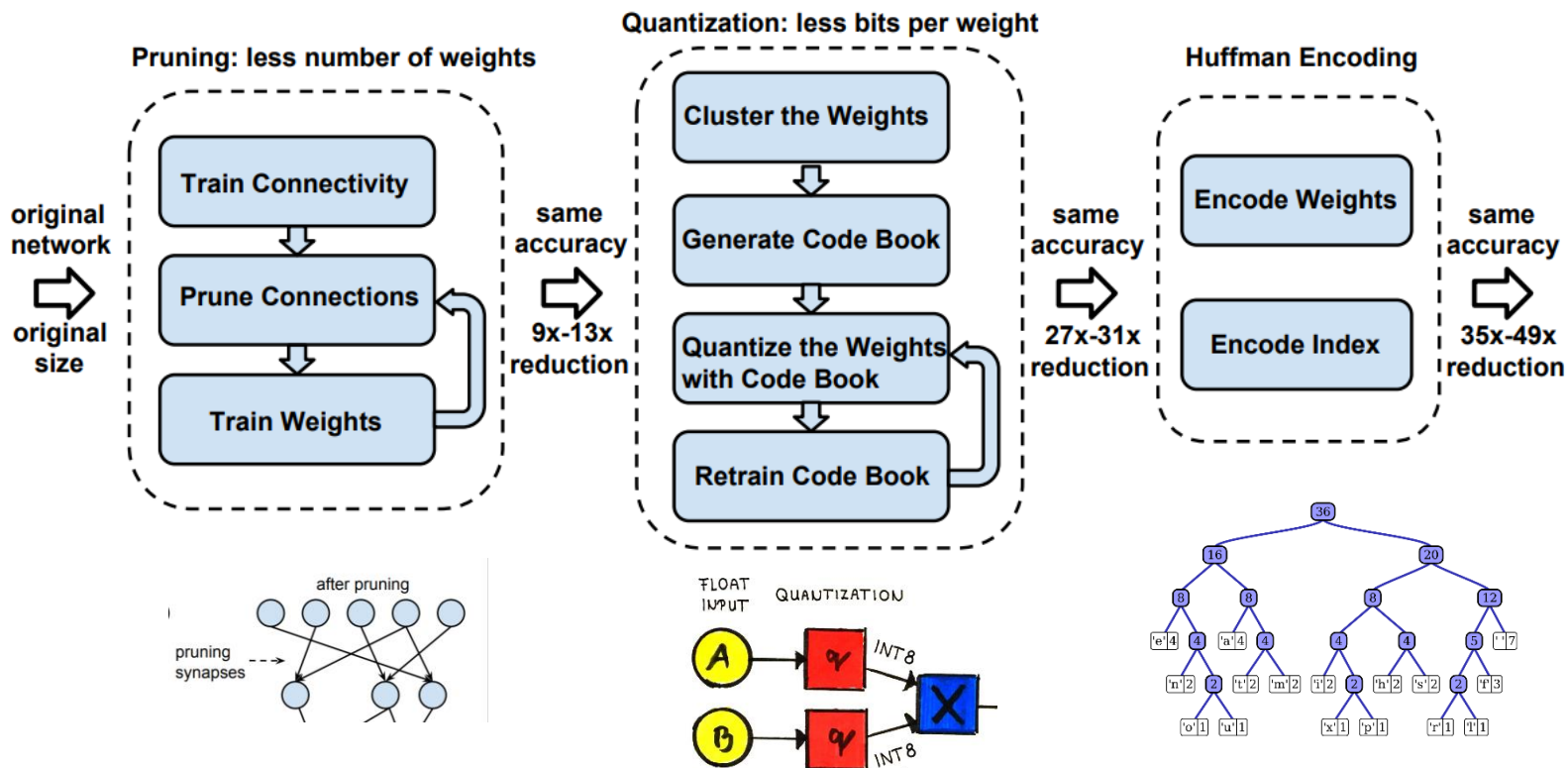


Layer	Weights	FLOP	Act%	Weights%	FLOP%
conv1	35K	211M	88%	84%	84%
conv2	307K	448M	52%	38%	33%
conv3	885K	299M	37%	35%	18%
conv4	663K	224M	40%	37%	14%
conv5	442K	150M	34%	37%	14%
fc1	38M	75M	36%	9%	3%
fc2	17M	34M	40%	9%	3%
fc3	4M	8M	100%	25%	10%
Total	61M	1.5B	54%	11%	30%

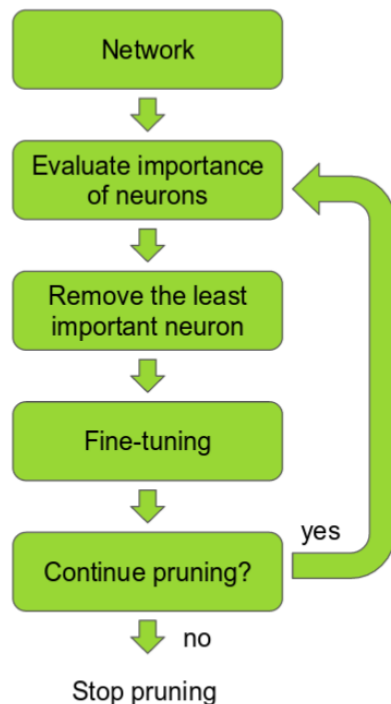


Layer	Weights	FLOP	Act%	Weights%	FLOP%
conv1_1	2K	0.2B	53%	58%	58%
conv1_2	37K	3.7B	89%	22%	12%
conv2_1	74K	1.8B	80%	34%	30%
conv2_2	148K	3.7B	81%	36%	29%
conv3_1	295K	1.8B	68%	53%	43%
conv3_2	590K	3.7B	70%	24%	16%
conv3_3	590K	3.7B	64%	42%	29%
conv4_1	1M	1.8B	51%	32%	21%
conv4_2	2M	3.7B	45%	27%	14%
conv4_3	2M	3.7B	34%	34%	15%
conv5_1	2M	925M	32%	35%	12%
conv5_2	2M	925M	29%	29%	9%
conv5_3	2M	925M	19%	36%	11%
fc6	103M	206M	38%	4%	1%
fc7	17M	34M	42%	4%	2%
fc8	4M	8M	100%	23%	9%
total	138M	30.9B	64%	7.5%	21%

Pruning in deep neural networks



Han et al., Deep Compression ICLR 2016



2.1 ORACLE PRUNING

Minimizing the difference in accuracy between the full and pruned models depends on the criterion for identifying the “least important” parameters, called *saliency*, at each step. The best criterion would be an exact empirical evaluation of each parameter, which we denote the *oracle* criterion, accomplished by ablating each non-zero parameter $w \in \mathcal{W}'$ in turn and recording the cost’s difference.

To compute the oracle, we evaluate the change in loss caused by removing each individual feature map from the fine-tuned VGG-16 network. (See Appendix A.3 for additional analysis.) We rank

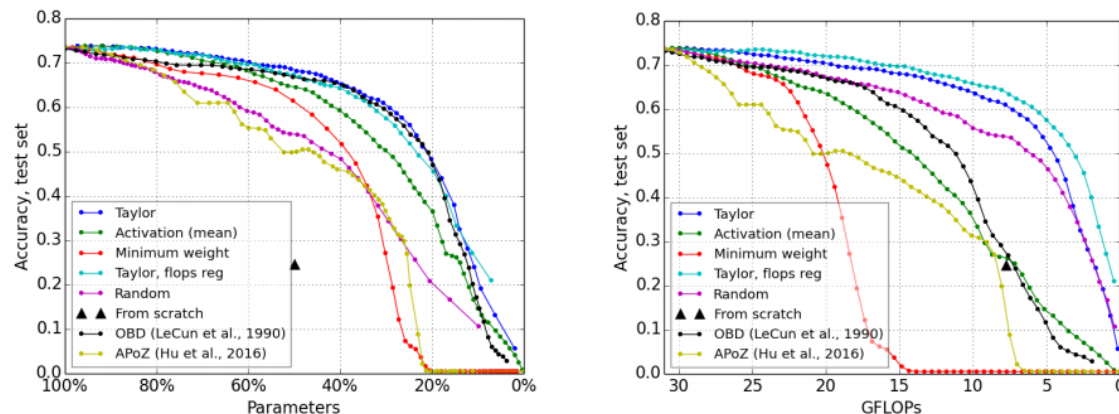


Figure 4: Pruning of feature maps in VGG-16 fine-tuned on the Birds-200 dataset.

WHAT IS THE STATE OF NEURAL NETWORK PRUNING?

Davis Blalock^{*1} Jose Javier Gonzalez Ortiz^{*1} Jonathan Frankle¹ John Guttag¹

Speed and Size Tradeoffs for Original and Pruned Models

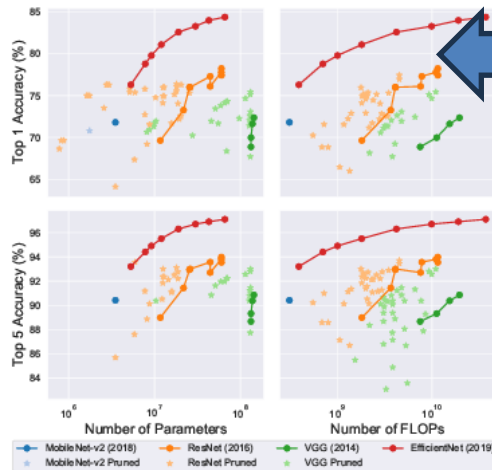
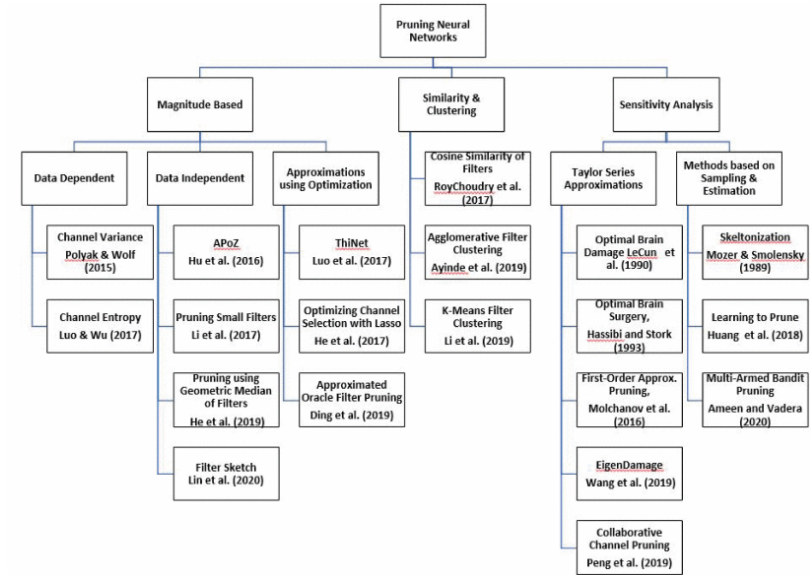


Figure 1: Size and speed vs accuracy tradeoffs for different pruning methods and families of architectures. Pruned models sometimes outperform the original architecture, but rarely outperform a better architecture.

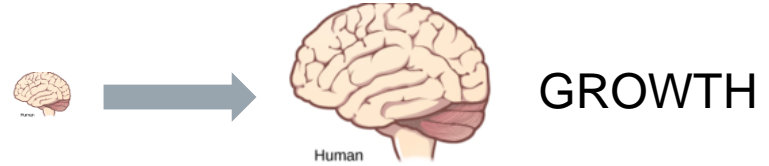
4.5 Confounding Variables

Even when comparisons include the same datasets, models, and operating points, other confounding variables make meaningful comparisons difficult. Some variables of particular interest include:

- Accuracy and efficiency of the initial model
- Data augmentation and preprocessing
- Random variations in initialization, training, and fine-tuning. This includes choice of optimizer, hyperparameters, and learning rate schedule.
- Pruning and fine-tuning schedule
- Deep learning library. Different libraries are known to

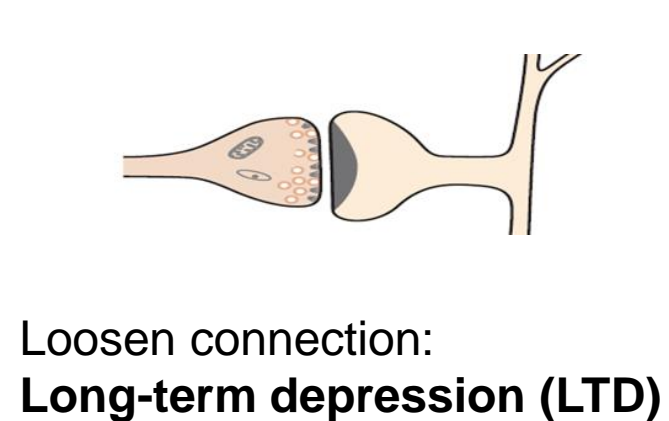
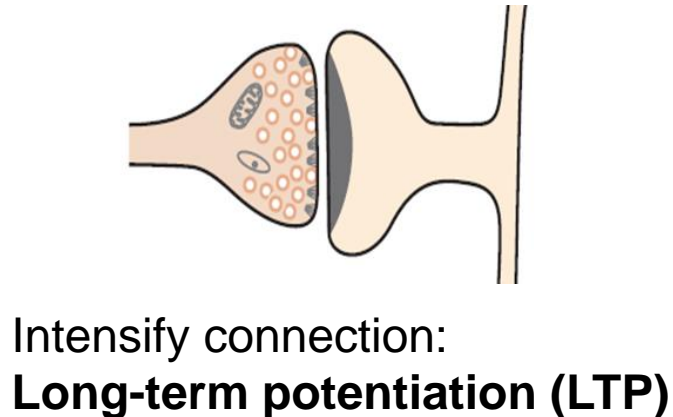
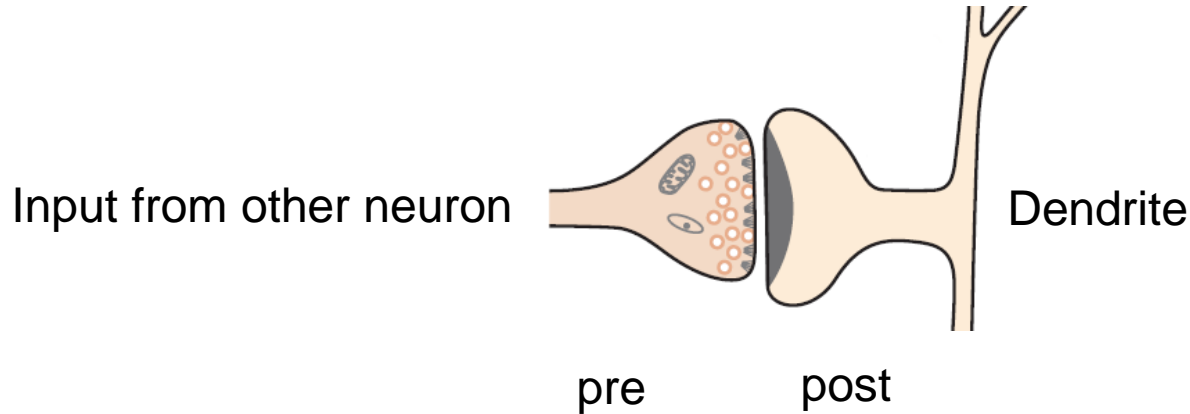


The brain is dynamic

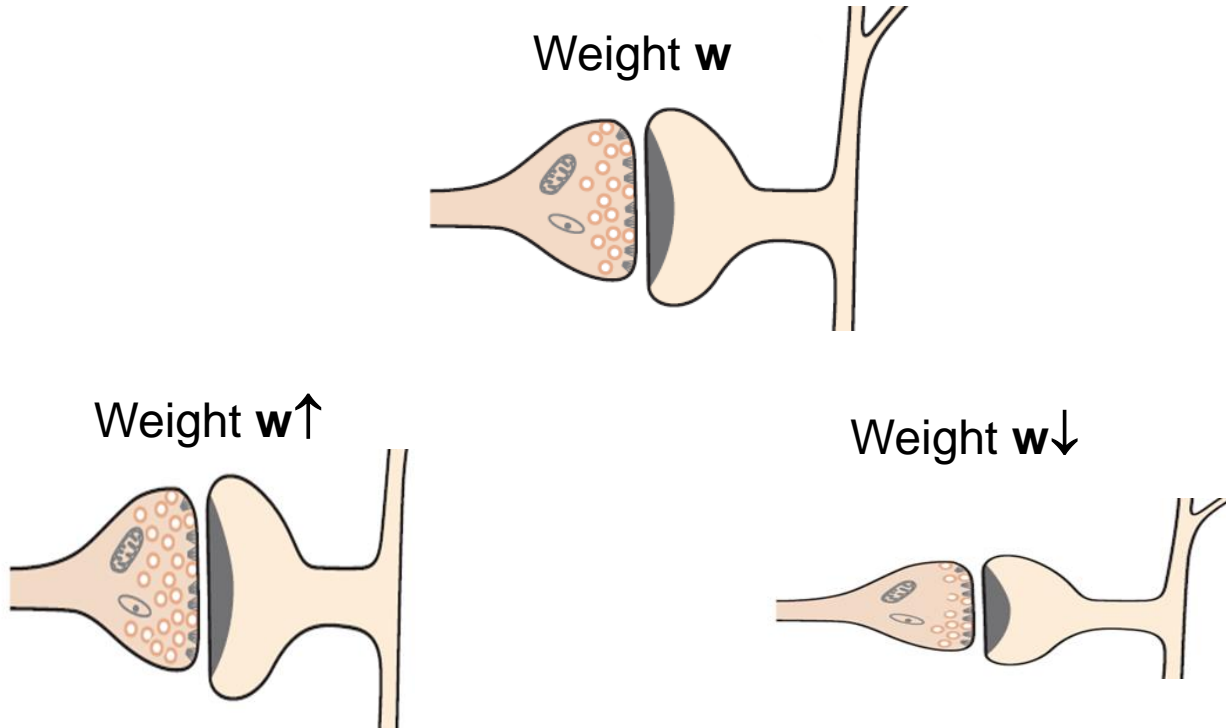


Neural plasticity

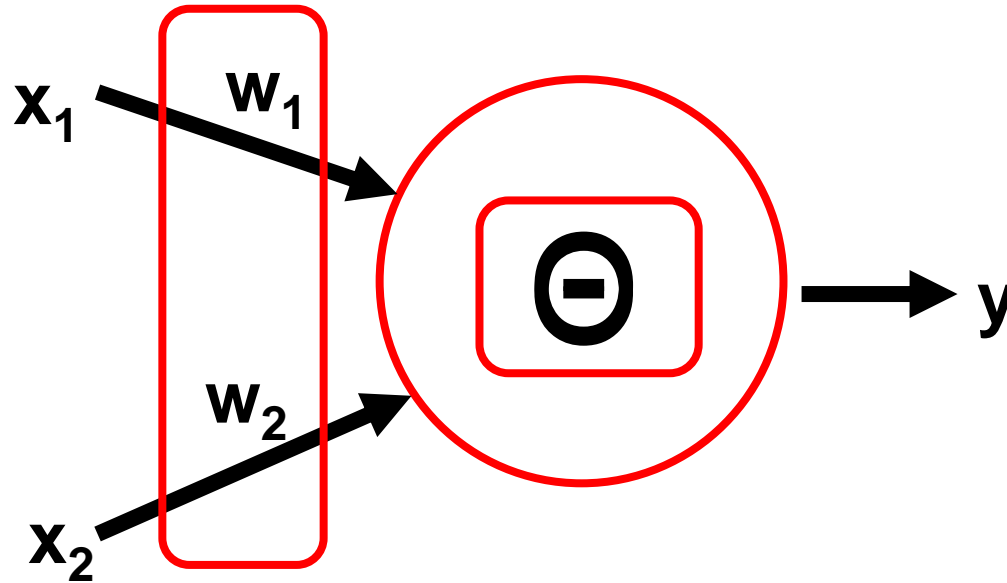
Changing the synapse strength



Changing the synapse strength



Perceptrons



Frank
Rosenblatt

How do we change...

Weights: Hebbian learning, Delta Rule
 Synaptic anatomy

Threshold: Learnable w/ SGD

Learnable Extended Activation Function for Deep Neural Networks

YEVGENIY BODIANSKIY¹, SERHII KOSTIUK²

¹Control Systems Research Laboratory, National University of Radio Electronics, Nasykiv av. 14, Kharkiv, 61166, Kharkiv, Ukraine, (e-mail: yevgeniy.bodianskiy@nure.ua)

²Dept. of Artificial Intelligence, National University of Radio Electronics, Nasykiv av. 14, Kharkiv, 61166, Kharkiv, Ukraine, (e-mail: serhii.kostiuk@nure.ua)

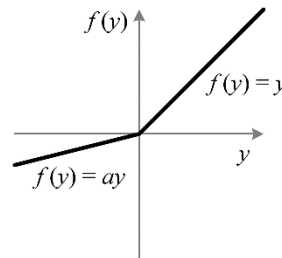
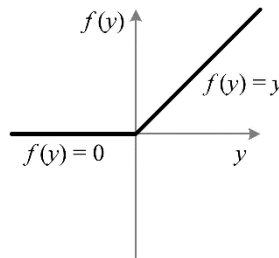
Corresponding author: Serhii Kostiuk (e-mail: serhii.kostiuk@nure.ua).

Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification

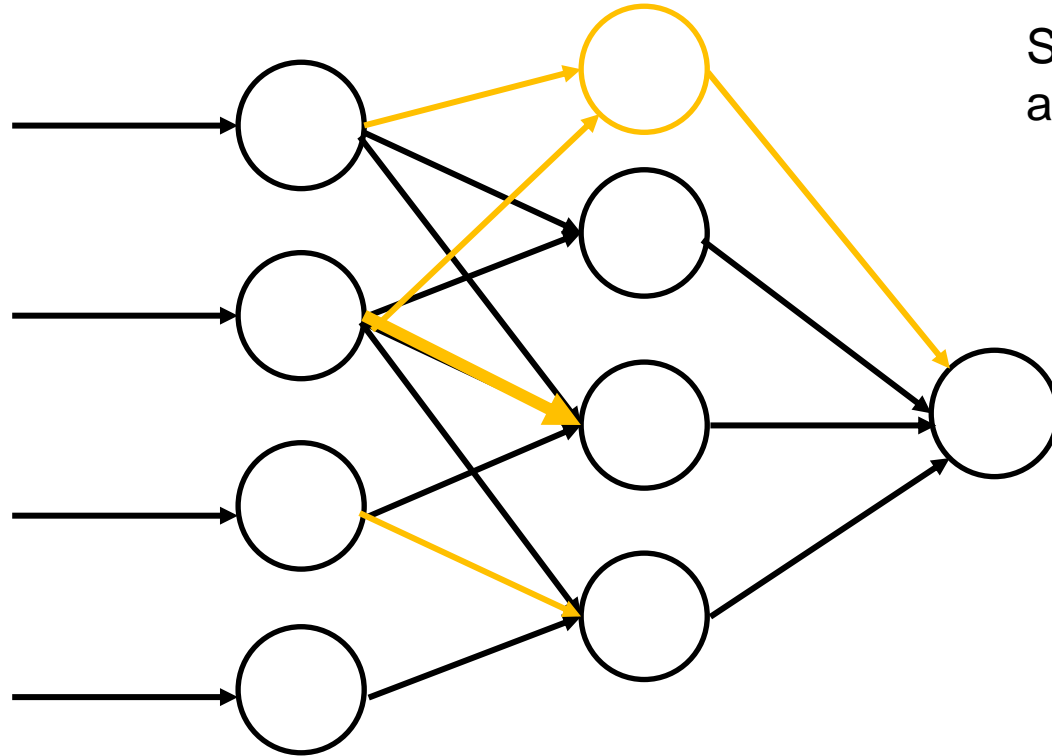
Kaiming He, Xiangyu Zhang, Shaoqing Ren, Jian Sun

Activation function:

Parametric ReLU
(Earlier layers more linear,
Later more variable)



Multilayer perceptrons



Structural
and synaptic plasticity

→ Typical learning paradigms work with synaptic plasticity

Evolving neural networks through augmenting topologies

[KO Stanley, R Miikkulainen](#) - *Evolutionary computation*, 2002 - MIT Press

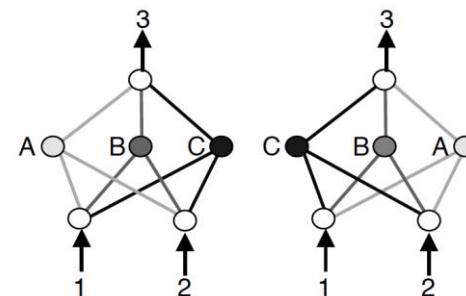
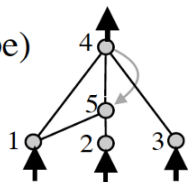
An important question in neuroevolution is how to gain an advantage from evolving neural network topologies along with weights. We present a method, NeuroEvolution of Augmenting Topologies (NEAT), which outperforms the best fixed-topology method on a challenging benchmark reinforcement learning task. We claim that the increased efficiency is due to (1) employing a principled method of crossover of different topologies, (2) protecting structural innovation using speciation, and (3) incrementally growing from minimal structure ...

☆ 99 Zitiert von: 3254 Ähnliche Artikel Alle 23 Versionen

Genome (Genotype)

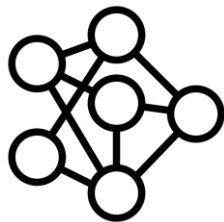
Node Genes	Node 1	Node 2	Node 3	Node 4	Node 5			
	Sensor	Sensor	Sensor	Output	Hidden			
Connect. Genes	In 1	In 2	In 3	In 2	In 5	In 1	In 4	
	Out 4	Out 4	Out 4	Out 5	Out 4	Out 5	Out 5	
	Weight 0.7	Weight-0.5	Weight 0.5	Weight 0.2	Weight 0.4	Weight 0.6	Weight 0.6	
	Enabled	DISABLED	Enabled	Enabled	Enabled	Enabled	Enabled	
	Innov 1	Innov 2	Innov 3	Innov 4	Innov 5	Innov 6	Innov 11	

Network (Phenotype)



$$\begin{array}{r} [A,B,C] \\ \times [C,B,A] \\ \hline \end{array}$$
 Crossovers: [A,B,A] [C,B,C]
 (both are missing information)

NEAT vs. the human brain



No topology enforced

Minimal set to tackle task

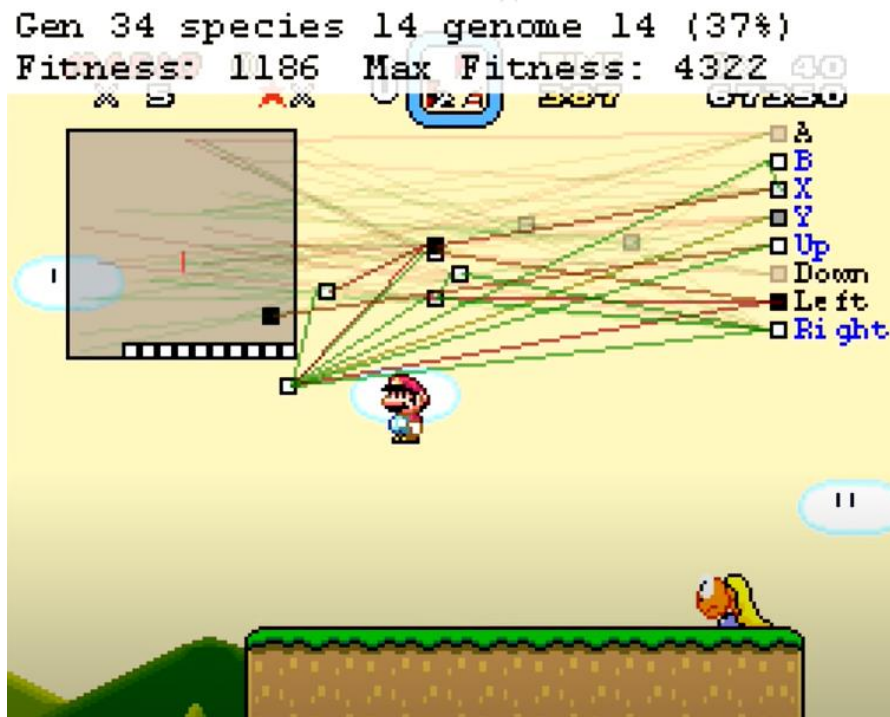
Hard to extrapolate to
Similar, yet slightly different tasks



Overall topology enforced

One-fits-all solution

Very good in abstraction



„**MarI/O** is a program made of neural networks and genetic algorithms that kicks butt at Super Mario World.”

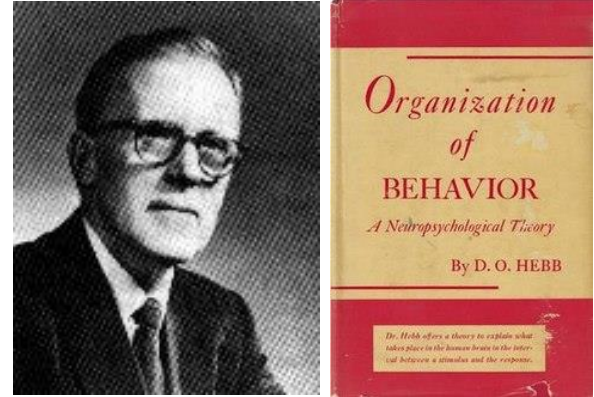
© Seth Bling, YouTube

<https://www.youtube.com/watch?v=qv6UVOQ0F44>

How to change weights?

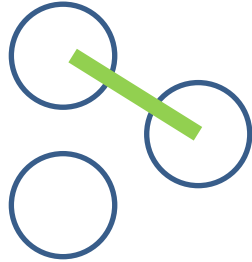
Hebb's rule:

„Neurons that fire together, wire together.“

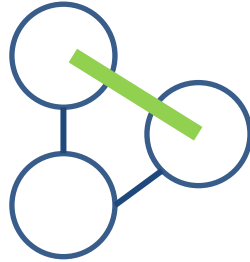


$$w_{ij} = x_i x_j$$

x_i	x_j	w_{ij}
0	0	0
1	0	0
0	1	0
1	1	1



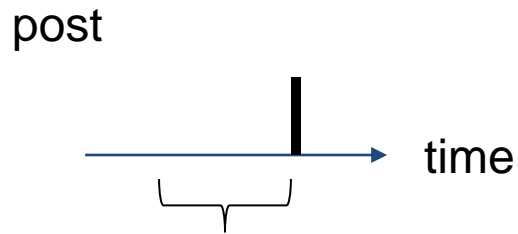
Structural
plasticity



Synaptic
plasticity

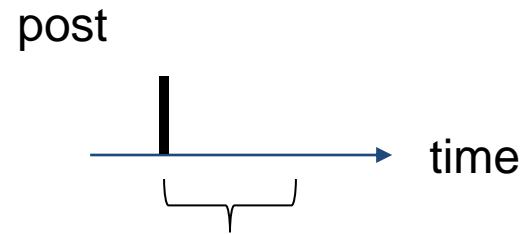
Spike timing dependent plasticity (STDP)

Hebb's rule – cntd.



< 20 ms

Synaptic $w \uparrow \rightarrow$ LTP



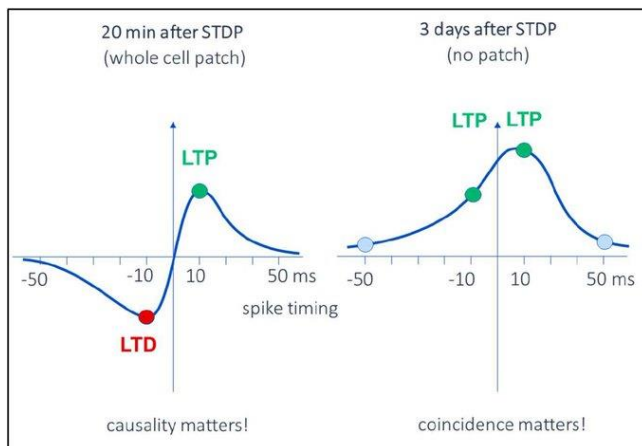
< 20 ms

Synaptic $w \downarrow \rightarrow$ LTD

Spike-timing-dependent plasticity rewards synchrony rather than causality

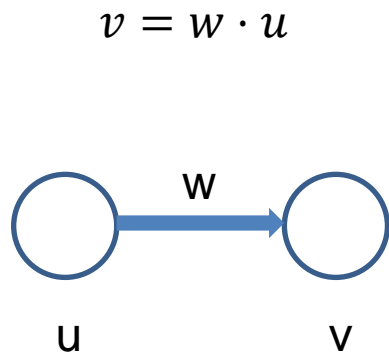
Margarita Anisimova¹, Bas van Bommel¹, Marina Mikhaylova^{1,2}, J. Simon Wiegert, Thomas G.

Oertner^{1*}, Christine E. Gee^{1,3*}



<https://doi.org/10.1101/863365>

“Our results confirm that neurons wire together if they fire together, but suggest that synaptic depression after anti-causal activation (tLTD) is a transient phenomenon.”



$$v = w \cdot u^T$$

The basic Hebb's rule:

$$\tau_w \frac{dw}{dt} = v \cdot u$$

$$\tau_w \frac{dw}{dt} = w \cdot \underbrace{u \cdot u^T}_{\text{Correlation matrix}}$$

$$\tau_w \frac{dw}{dt} = Q \cdot w^T$$

Correlation-based plasticity rule

$$w \rightarrow w + \epsilon Q \cdot w^T \quad \epsilon = \frac{1}{\tau_w}$$

Correlation-based plasticity rule

$$\tau_w \frac{dw}{dt} = \mathbf{Q} \cdot \mathbf{w}^T$$

→ Basic Hebb only allows LTP

Postsynaptic
LTD/LTP switch

$$\tau_w \frac{dw}{dt} = (v - \theta_v) \cdot \mathbf{u}$$

$$\tau_w \frac{dw}{dt} = v \cdot (\mathbf{u} - \theta_u)$$

Presynaptic
LTD/LTP switch

Covariance matrix

$$v = \mathbf{w} \cdot \mathbf{u}^T \quad \Rightarrow \quad \tau_w \frac{dw}{dt} = \mathbf{w} \cdot \mathbf{u} \cdot (\mathbf{u} - \theta_u)^T \quad \Rightarrow \quad \tau_w \frac{dw}{dt} = \mathbf{C} \cdot \mathbf{w}^T$$

Covariance-based plasticity rule

Hebbian learning suffers from instability

$$\tau_w \frac{dw}{dt} = v \cdot \mathbf{u} \cdot (v - \theta_v) \quad \text{If constant} \rightarrow \text{unstable}$$

➔ Threshold of postsynaptic activity that determines if synapse is strengthened or weakened.

Adapt threshold θ_v :

$$\tau_\theta \frac{d\theta_v}{dt} = v^2 - \theta_v$$

$$\tau_{\theta} \frac{d\theta_v}{dt} = v^2 - \theta_v \quad \rightarrow \text{Stabilize weights through postsynaptic activity}$$

Can we use penalty terms directly on the weight vector?

$$\tau_w \frac{d\mathbf{w}}{dt} = v \cdot \mathbf{u} - \frac{v(\mathbf{n} \cdot \mathbf{u}^T)\mathbf{n}}{N_u}$$

Normalize by subtracting the same quantity

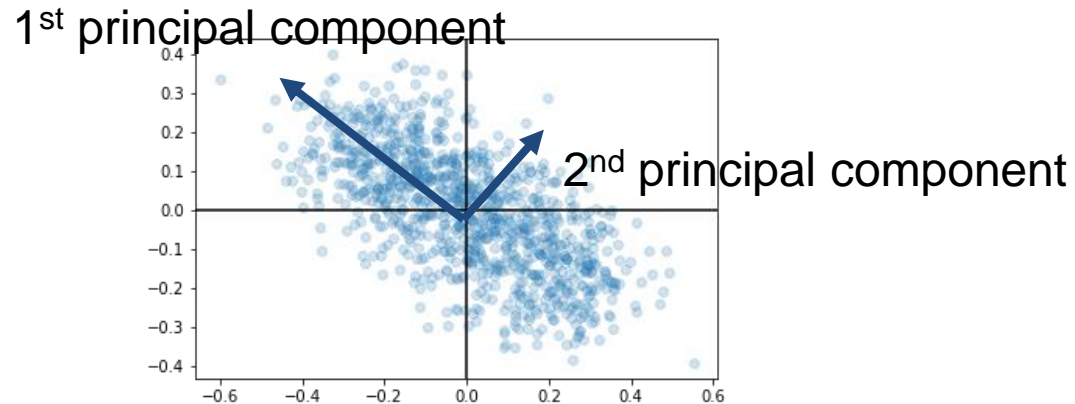


$$\tau_w \frac{d\mathbf{w}}{dt} = v \cdot \mathbf{u} - \alpha \cdot v^2 \cdot \mathbf{w} \quad \alpha > 0$$

$$\tau_w \frac{d\mathbf{w}}{dt} = v \cdot (\mathbf{u} - \alpha \cdot v \cdot \mathbf{w})$$

Unsupervised learning

PRINCIPAL COMPONENT ANALYSIS (PCA)

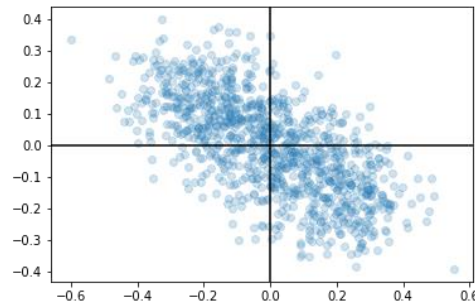


Simulation

Steps:

1. Data generation
2. Variable initialization
3. Iterate through data
 1. Compute pre-synaptic input
 2. Compute post-synaptic activation
 3. Compute Δw and update
4. Plot result

```
1 # That you see the same what I see
2 np.random.seed(42)
3 N = 1000
4
5 # Generate random data
6 x = np.linspace(-.3, .3, N)
7 np.random.shuffle(x)
8 y = -.7 * x
9
10 x += np.random.randn(x.size) / 10
11 y += np.random.randn(y.size) / 10
12
```



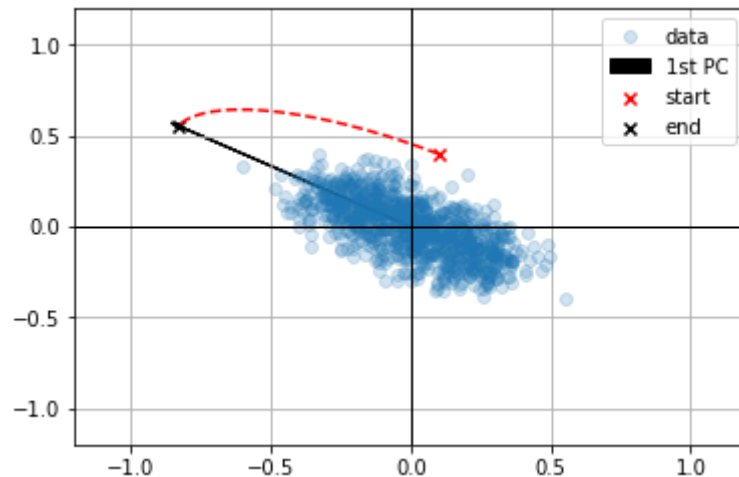
```
1 # Initialize some random weights
2 w = np.array([0.1, 0.4])
3 ## Training iterations
4 N = 1000
5 eta = 0.1
6 ws = []
7
```

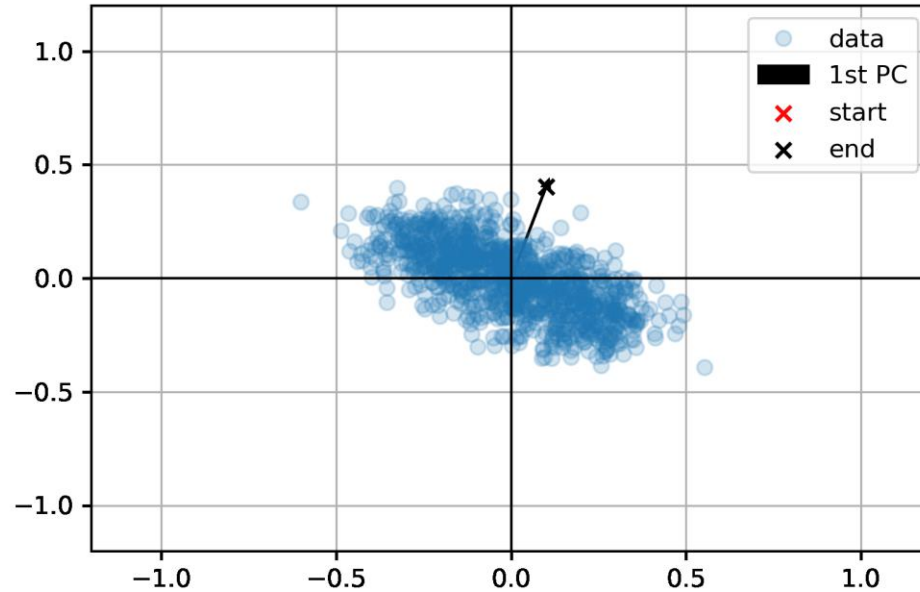
```
8 # Training
9 for i in range(N):
10     rPre = np.asarray([x[ix], y[ix]])
11     rPost = w @ rPre
12     w = w + eta * rPost * (rPre - rPost * w)
```

Simulation

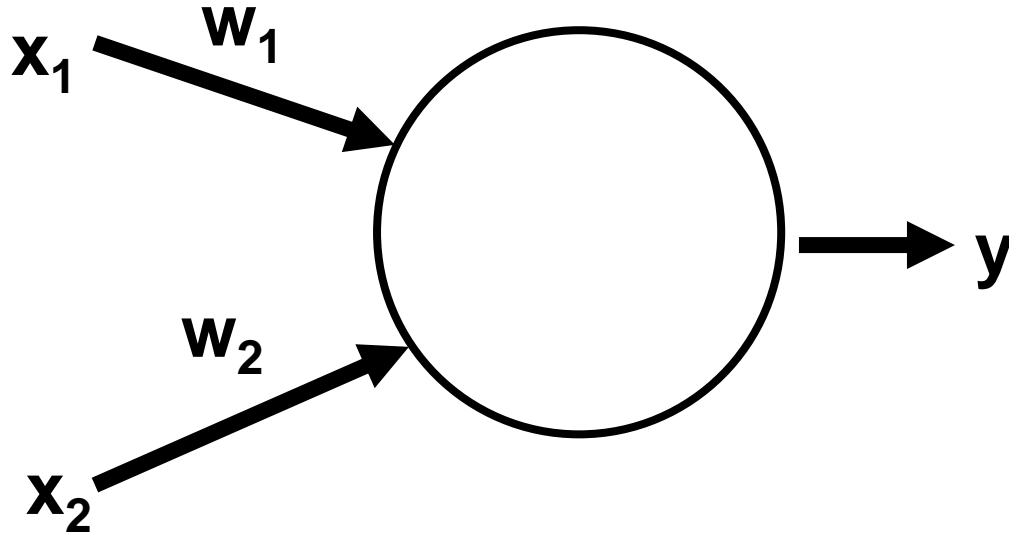
Steps:

1. Data generation
2. Variable initialization
3. Iterate through data
 1. Compute pre-synaptic input
 2. Compute post-synaptic activation
 3. Compute Δw and update
4. Plot result

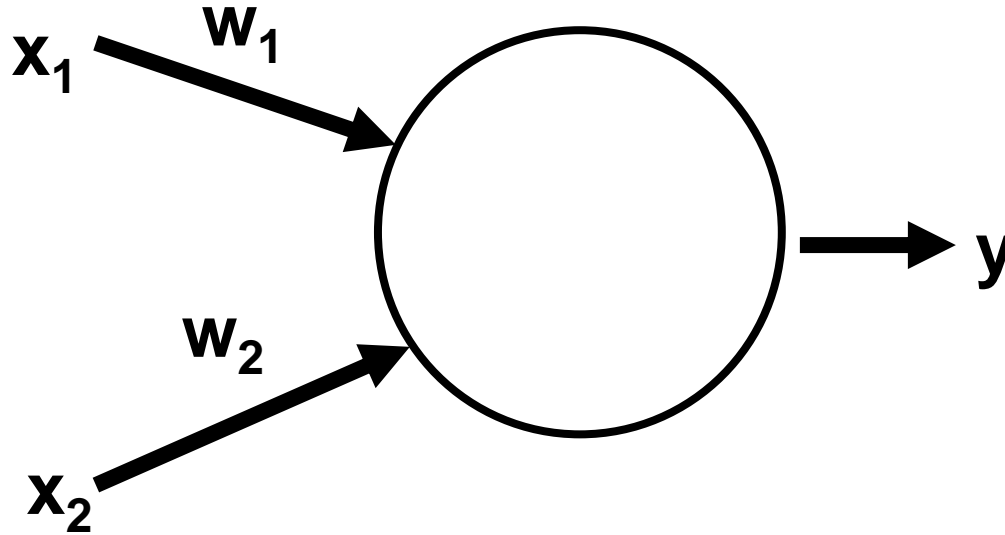




Perceptron



Unsupervised learning → no target imposed
Supervised learning → target imposed



$$\hat{y} = \sum_i w_i \cdot x_i$$

Target y : 2

Prediction \hat{y} : 1

$$\begin{aligned} \text{Error} &= \text{„Target – Prediction“} = \\ &= \frac{1}{2} (y - \hat{y})^2 \end{aligned}$$

$$\frac{dE}{dw_i} E = \frac{dE}{dw_i} \frac{1}{2} (y - \hat{y})^2 = -(y - \hat{y}) \frac{d\hat{y}}{dw_i} \sum_i w_i \cdot x_i = -(y - \hat{y}) \cdot x_i$$

Delta rule

$$\Delta w_i = \alpha \cdot (y - \hat{y}) \cdot x_i$$

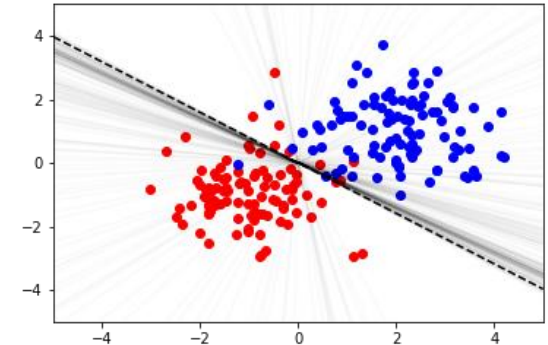
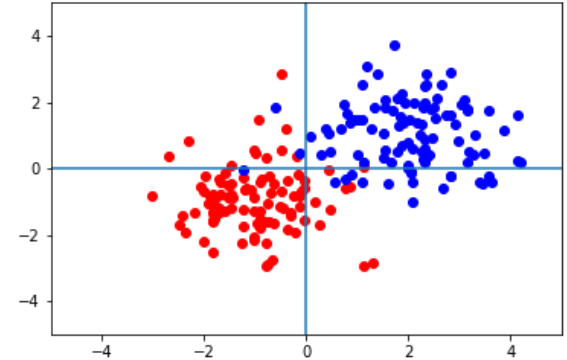
Steps:

1. Data generation
2. Variable initialization
3. Iterating
 1. Compute prediction
 2. Update weights
4. Plotting

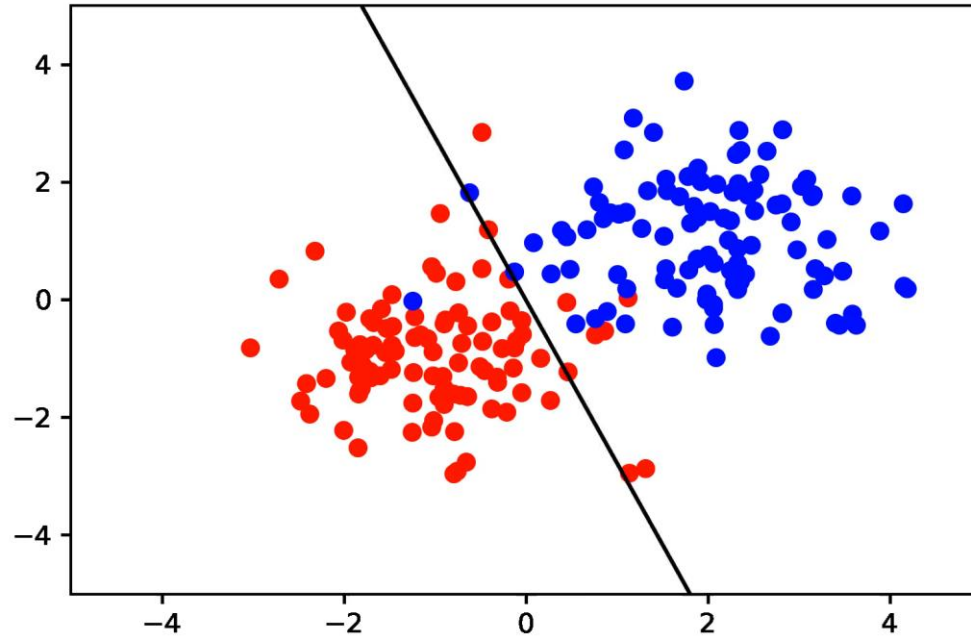
```
1 np.random.seed(42)
2
3 N = 100
4
5 xs = []
6
7 for _ in range(N):
8     x0 = np.random.randn()-1
9     x1 = np.random.randn()-1
10    xs.append((x0, x1))
11
12    x0 = np.random.randn()+2
13    x1 = np.random.randn()+1
14    xs.append((x0, x1))
15
16 xs = np.asarray(xs)
```

```
1 w = np.random.randn(2)
2 eta = 0.01
```

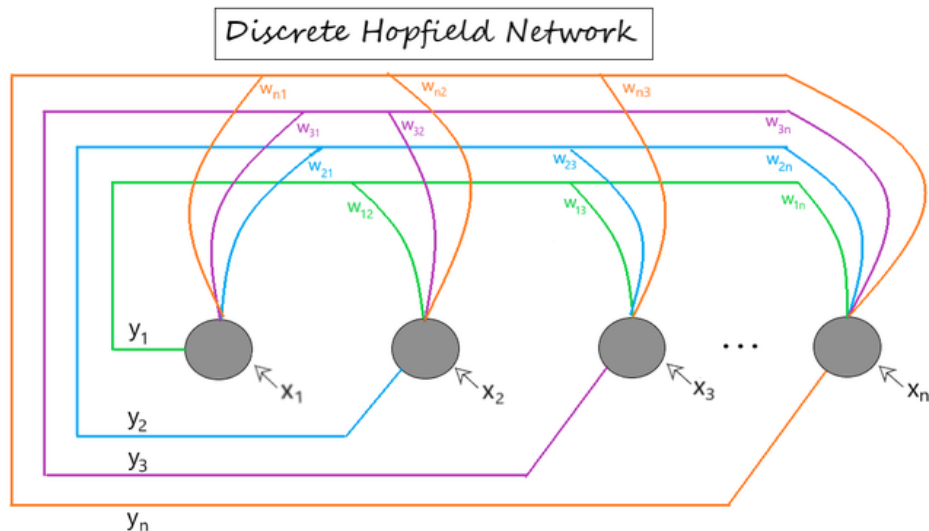
```
5 for i in range(2*N):
6     # Determine class/target
7     t = -1 if i % 2 == 0 else 1
8
9     # Compute prediction
10    pred = w @ xs[i]
11
12    # Update weights
13    w = w + eta * (t-pred) * xs[i]
14
```



Simulation over time



Applied in deep neural networks



Binary (0,1)
Or bipolar (-1,1)
activations

Learning rules should incorporate:
Local and incremental
→ e.g. Hebbian learning

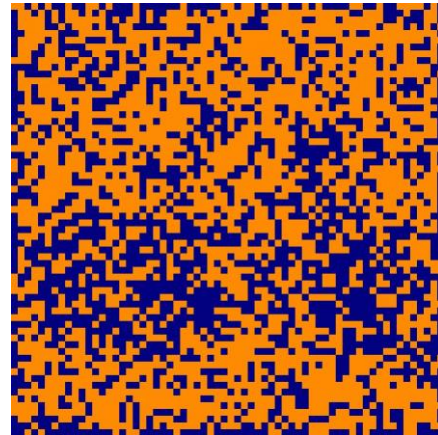
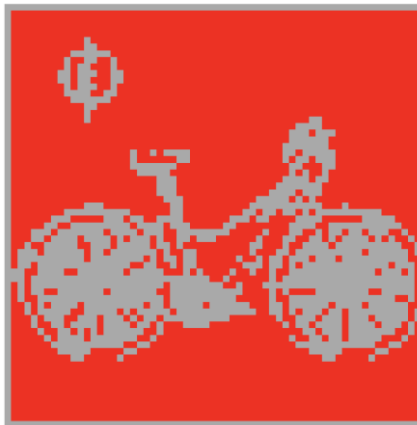
Minimizing the energy in
the system

Capacity:
 $0.16 * n$ (n=units in Network)

$$C \cong \frac{n}{2 \log_2 n}$$

Large images
(1024x1024) → require
8M neurons (!!)

Storing continuous
space?



64x64 px

→ Trained on 4096 nodes

→ 7.2% (256) nodes were updated per step

<https://towardsdatascience.com/hopfield-networks-neural-memory-machines-4c94be821073>

Hopfield networks - now

HOPFIELD NETWORKS IS ALL YOU NEED

Hubert Ramsauer* Bernhard Schäff* Johannes Lehner* Philipp Seidl*
Michael Widrich* Thomas Adler* Lukas Gruber* Markus Holzleitner*
Milena Pavlović^{‡,§} Geir Kjetil Sandve[§] Victor Greiff[‡] David Kreil[†]
Michael Kopp[†] Günter Klambauer* Johannes Brandstetter* Sepp Hochreiter*,[†]

*ELLIS Unit Linz, LIT AI Lab, Institute for Machine Learning,
Johannes Kepler University Linz, Austria

[†]Institute of Advanced Research in Artificial Intelligence (IARAI)

[‡]Department of Immunology, University of Oslo, Norway

[§]Department of Informatics, University of Oslo, Norway

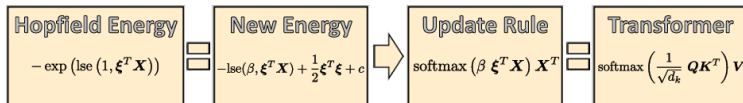


Figure 1: We generalize the energy of binary modern Hopfield networks to continuous states while keeping fast convergence and storage capacity properties. We also propose a new update rule that minimizes the energy. The new update rule is the attention mechanism of the transformer. Formulae are modified to express softmax as row vector. “=”-sign means “keeps the properties”.

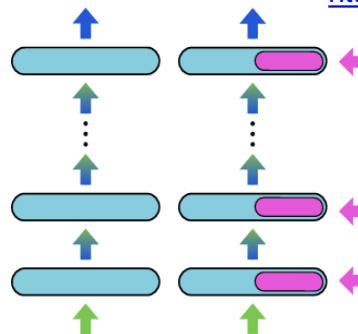
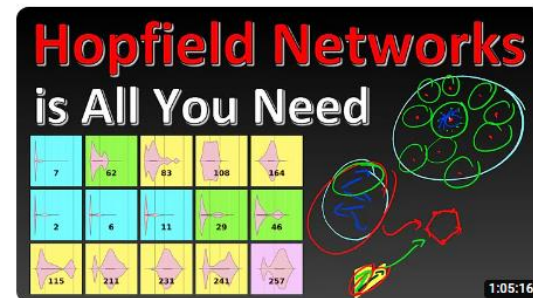


Figure 2: Left: A standard deep network with layers (■) propagates either a vector or a set of vectors from the input to the output. Right: A deep network, where layers (■) are equipped with associative memories via Hopfield layers (■).

<https://www.youtube.com/watch?v=nv6oFDp6rNQ>



Explanation
by Yannic Kilcher

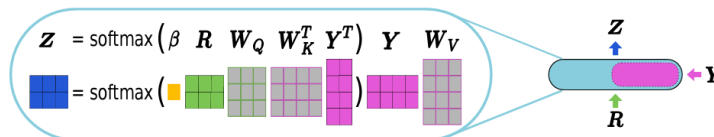
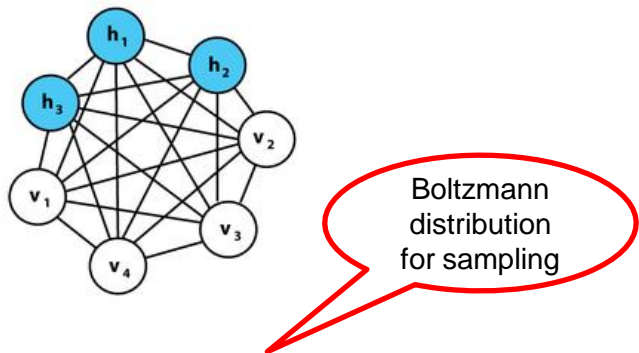


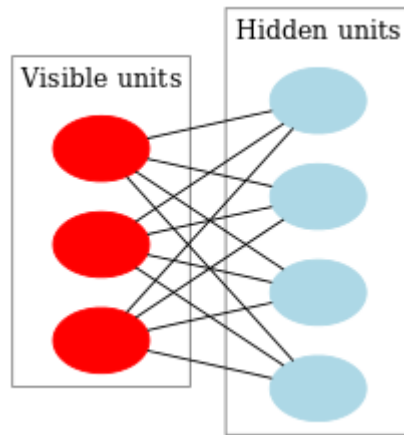
Figure 3: The layer Hopfield allows the association of two sets R (■) and Y (■). It can be integrated into deep networks that propagate sets of vectors. The Hopfield memory is filled with a set from either the input or previous layers. The output is a set of vectors Z (■).



Boltzmann machine

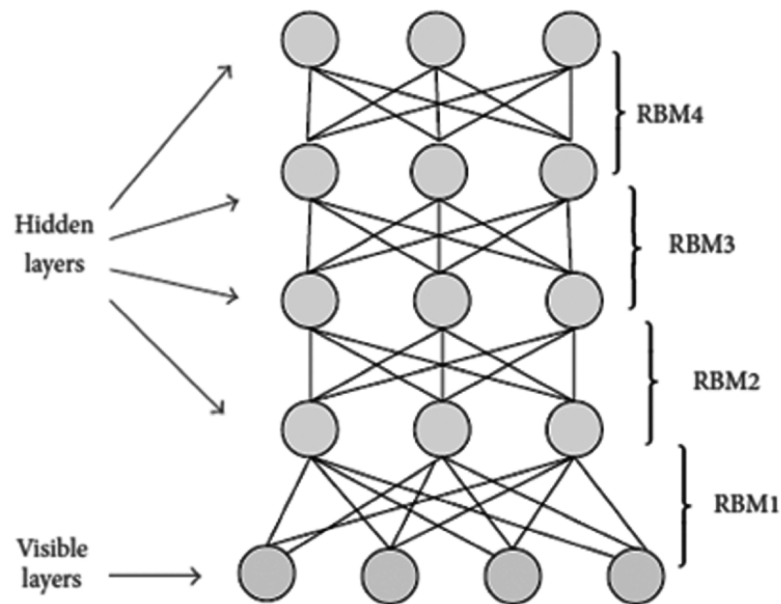
Everything is connected to everything,
Energy-based model (EBM),
Similar to Hopfields.

Training:
minimizing KL-divergence



Restricted Boltzmann machine
(RBM) – no visible-visible and hidden-hidden connections
Training w/ contrastive divergence

Deep Belief Networks



First round of training:
unsupervised w/ CD
RBM by RBM

Connected to output,
then **supervised** classification