

Intelligent Systems Programming

Lecture 4: DeepMind Game Algorithms

DeepMind Founders, acquired by Google 2014, HQ in London

- Mantra: Solve intelligence. Use it to make the world a better place.



Demis Hassabis, CEO



Mustafa Suleyman



Shane Legg

Readings and Project

SSS17 (URL link on Learnit)

ARTICLE

Mastering the game of Go without human knowledge

David Silver^{1*}, Julian Schrittwieser^{1*}, Karen Simonyan¹, Ioannis Antonoglou¹, Aja Huang¹, Arthur Guez¹, Thomas Hubert¹, Lucas Baker¹, Matthew Lai¹, Adrian Bolton¹, Yutian Chen¹, Timothy Lillicrap¹, Fan Hui¹, Laurent Sifre¹, George van den Driessche¹, Thore Graepel¹ & Demis Hassabis¹

A long-standing goal of artificial intelligence is an algorithm that learns, *tabula rasa*, superhuman proficiency in challenging domains. Recently, AlphaGo became the first program to defeat a world champion in the game of Go. The triumph of AlphaGo over grand masters and self-play training neural networks were trained by supervised learning from human moves, and by reinforcement learning from self-play. Here we implement an algorithm based solely on reinforcement learning, without human data, guidance or domain knowledge beyond game rules. AlphaGo becomes its own teacher: a neural network is trained to predict AlphaGo's own move selections and also the winner of AlphaGo's games. This neural network improves the strength of the tree search, resulting in higher quality move selection and stronger self-play in the next iteration. Starting *tabula rasa*, our new program AlphaGo Zero achieved superhuman performance, winning 100–0 against the previously published, champion-defeating AlphaGo.

Much progress towards artificial intelligence has been made using supervised learning systems that are trained to replicate the decisions of human experts^{1–3}. However, expert data sets are often expensive, trained solely by self-play reinforcement learning, starting from random play, without any supervision or use of human data. Second, it uses only the black and white stones from the board as input features.

Othello Proj. (ZIP on Learnit)

IS Project 1- Othello

Your Task:
In groups of 2-4 you are going to make your own implementation of an adversarial search algorithm for the game Othello. Specifically, you will need to implement this in Java. You can either do this in C/C++ or Java. If you are not familiar with Java, you should talk to the TA if you have it approved. It is required that you use the provided classes and interface. The game board is represented as a 8x8 matrix of tokens. The board starts with 8 columns and 8 rows which is "horizontal". Imagine that on average, response time on a PC is at least than 10 seconds although five moves up to 30 seconds would be acceptable. Moreover, your implementation should be able to handle 1000 moves per second. You will need to apply several of the topics described during lecture for designing the evaluation and cut-off function.

What is provided:

You are provided with five classes and an interface. You do not have to change anything in these classes/interfaces. In fact, we recommend that you do not change anything in these classes:
• **GameBoard**: A GUI class that shows the game board and takes for input from the human player in order to make a move.
• **GameState**: A class to represent the state of the game, i.e., the board and whose turn it is. It has methods e.g. for finding the possible legal moves and placing a token in a selected position.
• **Player**: An abstract class that represents the two players in the game.
• **IOHandler**: The interface that you need to implement. It only has one method, namely `decideWhatMove`, which takes the current state of the game and the turn and decides what move the player in turn should make, i.e., at which position it should place a token.
• **Othello**: This is the class with the main-method. It reads in the provided parameters, and starts the game loop. It calls the `decideWhatMove`-method of the IOHandler. As you can see, it does not provide it with the required parameters, these are described a little further below.
• **Test**: A sample implementation of the required interface. It returns the best possible move that it finds.

Read the code and documentation of the class `GameState` carefully as you may want to use the provided methods for your implementation of `IOHandler`. If you implement in Java, you do not need to implement the `IOHandler` interface, as Java has its own `Object` class, but has to do so in C/C++. Otherwise, i.e. if you implement in C/C++, you have to implement another classes and document how to use them.

For convenience, players are represented with integers, specifically:

- +1 player +1 = the human/red AI's opponent = black.
 - 1 player -1 = the computer/black AI's opponent = red.
- Statically, the game board has to be represented in `GameState` by a 2-dimensional integer array, where the values range from -1 to +1. [Figure 1] shows the initial state of the board. When starting the game (see instructions further below), the game board should look like the one in Figure 1. The columns on the board are numbered from left to right with the numbers 0 to 7 on an 8x8 board, while the rows are numbered from top to bottom with the numbers from 0 to 7. On white square (0,0), the top-most white token in Figure 1 is on position (0,0).

*Of course, there are various implementations to be found on the internet, but do make your own implementation. If found out, saying its reason for exclusion from the course.

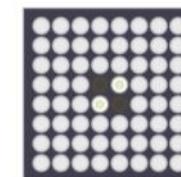


Figure 1 -The initial Othello game board

Running the game

Setup:
Download the provided files. Setup a project folder as you will usually do. In most IDEs, let's assume you call it "OthelloProject". In order to show the images correctly, it's important, to "OthelloProject" has a dependency called "image", which contains the images (as .png) in the folder you downloaded.

Abstract

The game of chess is the most widely-studied domain in the history of artificial intelligence. The strongest programs are based on a combination of sophisticated search techniques, domain-specific adaptations, and handcrafted evaluation functions that have been



Today's Program

[10-10:50] Methods

- Reinforcement learning
- Convolutional neural networks

[11-11:50] AlphaGo Zero

- AlphaGo Zero's neural network
- Self-play reinforcement learning
- Results

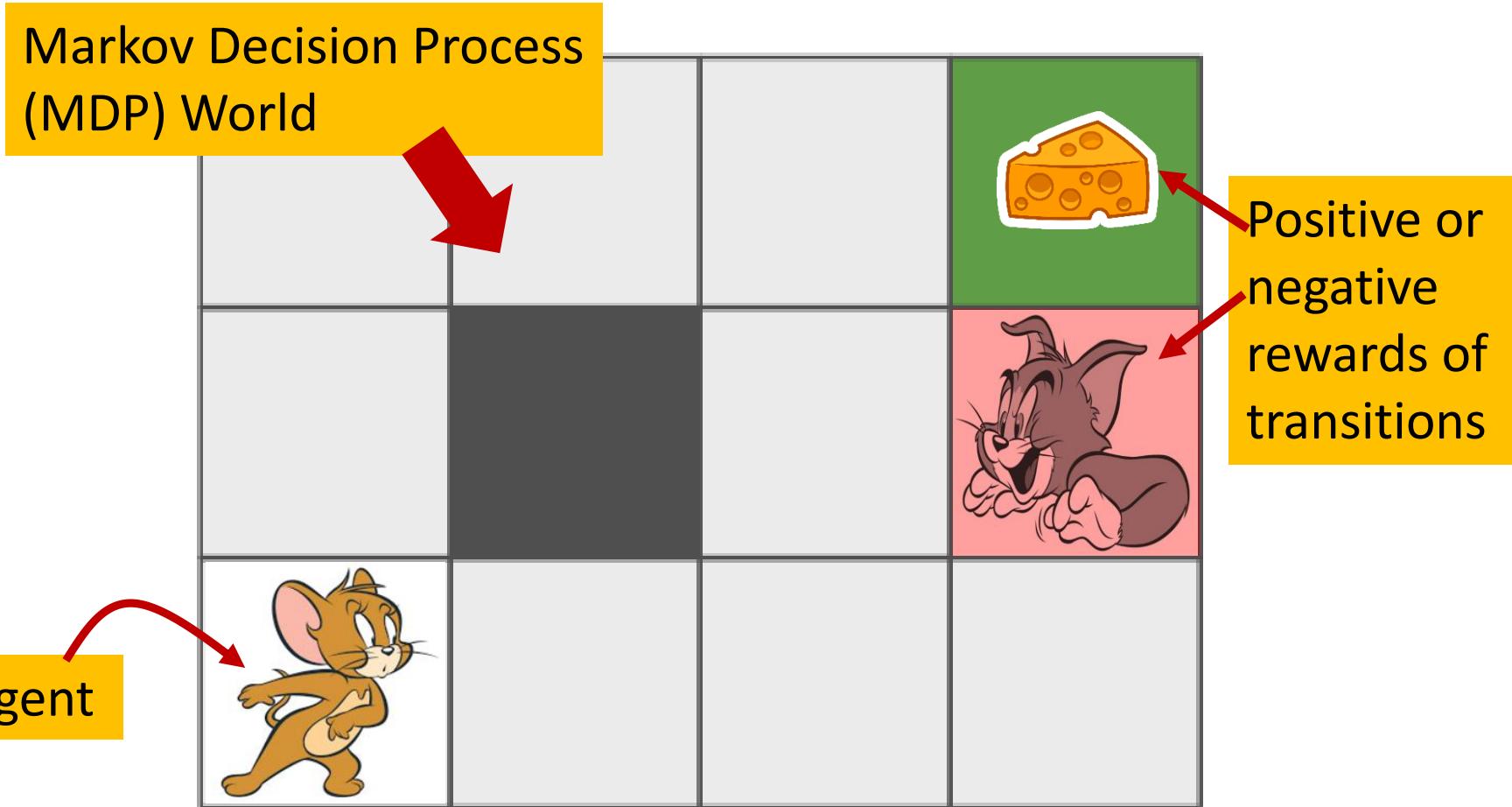
[12-14] Exercises

What's the Fuss About DeepMind?



How to Learn a Game from Zero Knowledge

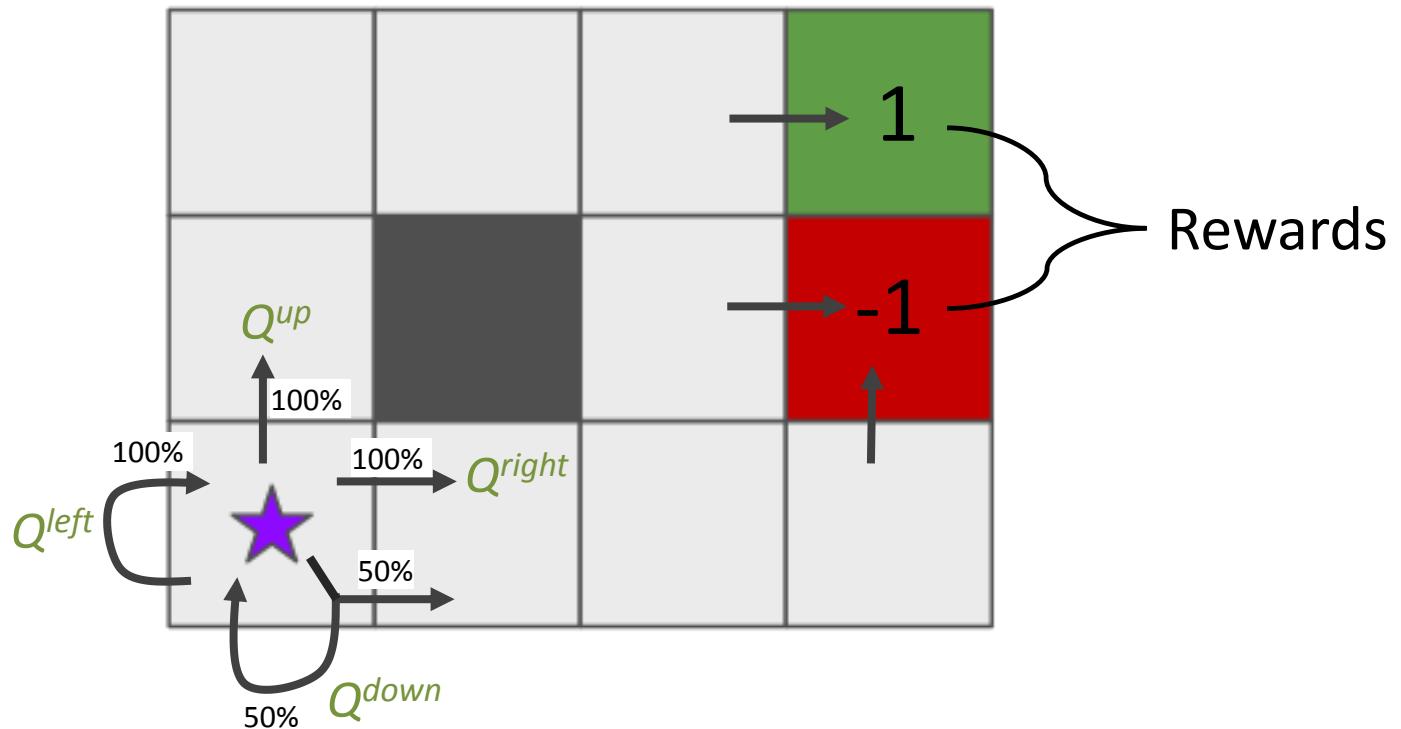
Reinforcement Learning (RL)



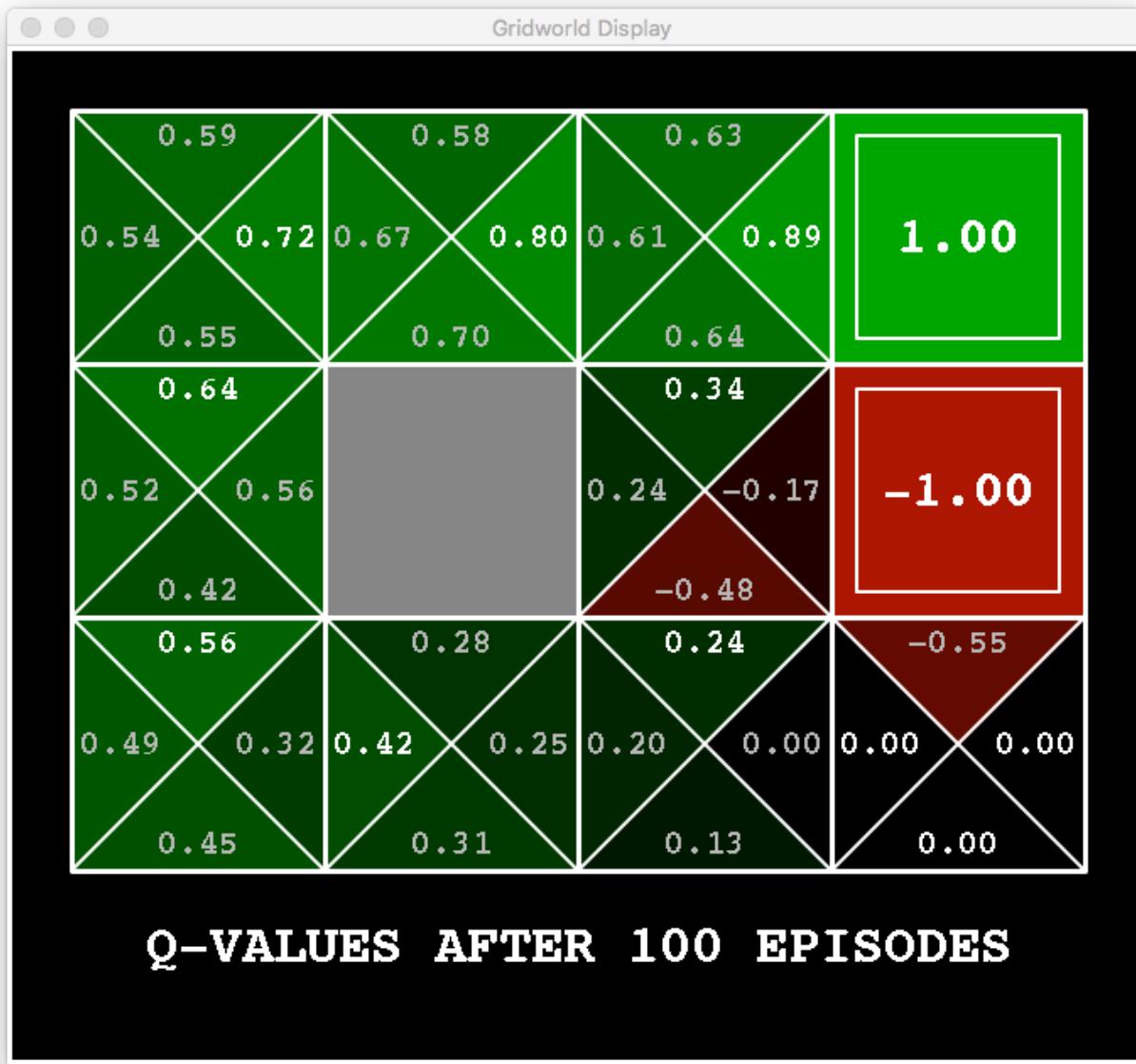
Q-Learning

$$Q(s_t, a_t) \leftarrow (1 - \alpha) \cdot \underbrace{Q(s_t, a_t)}_{\text{old value}} + \underbrace{\alpha}_{\text{learning rate}} \cdot \left(\underbrace{r_t}_{\text{reward}} + \underbrace{\gamma}_{\text{discount factor}} \cdot \underbrace{\max_a Q(s_{t+1}, a)}_{\substack{\text{learned value} \\ \text{estimate of optimal future value}}} \right)$$

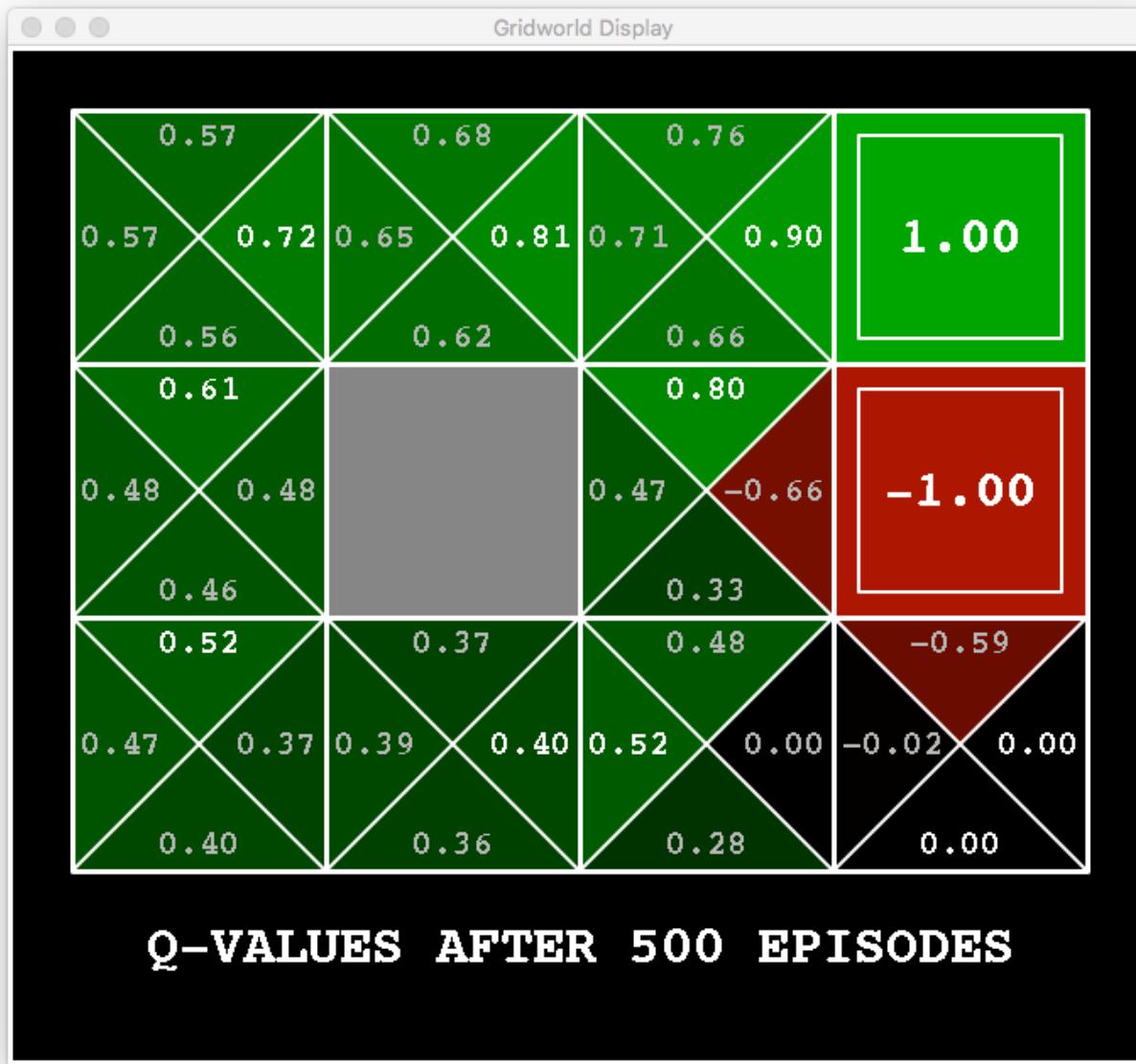
Policy (π):
Take action
with max
Q-value



Q-Learning Example (1/2)



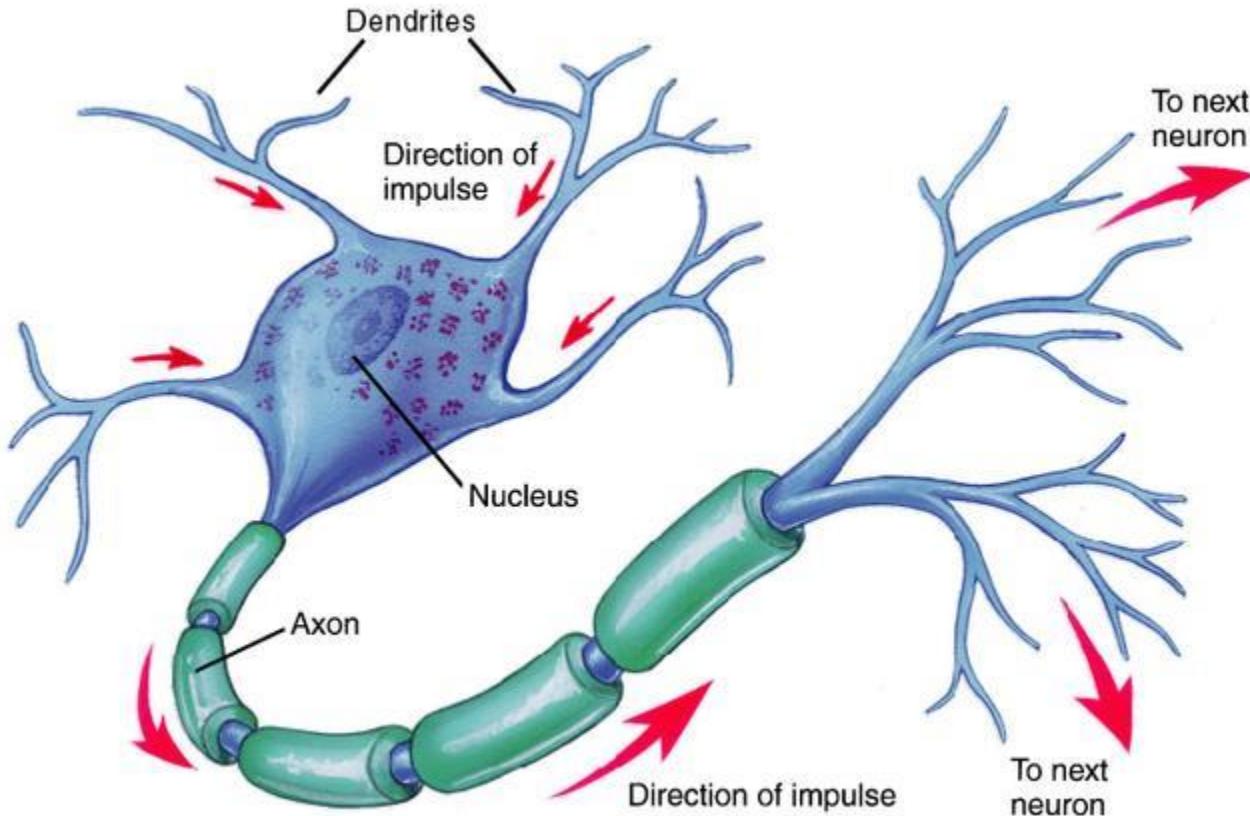
Q-Learning Example (2/2)



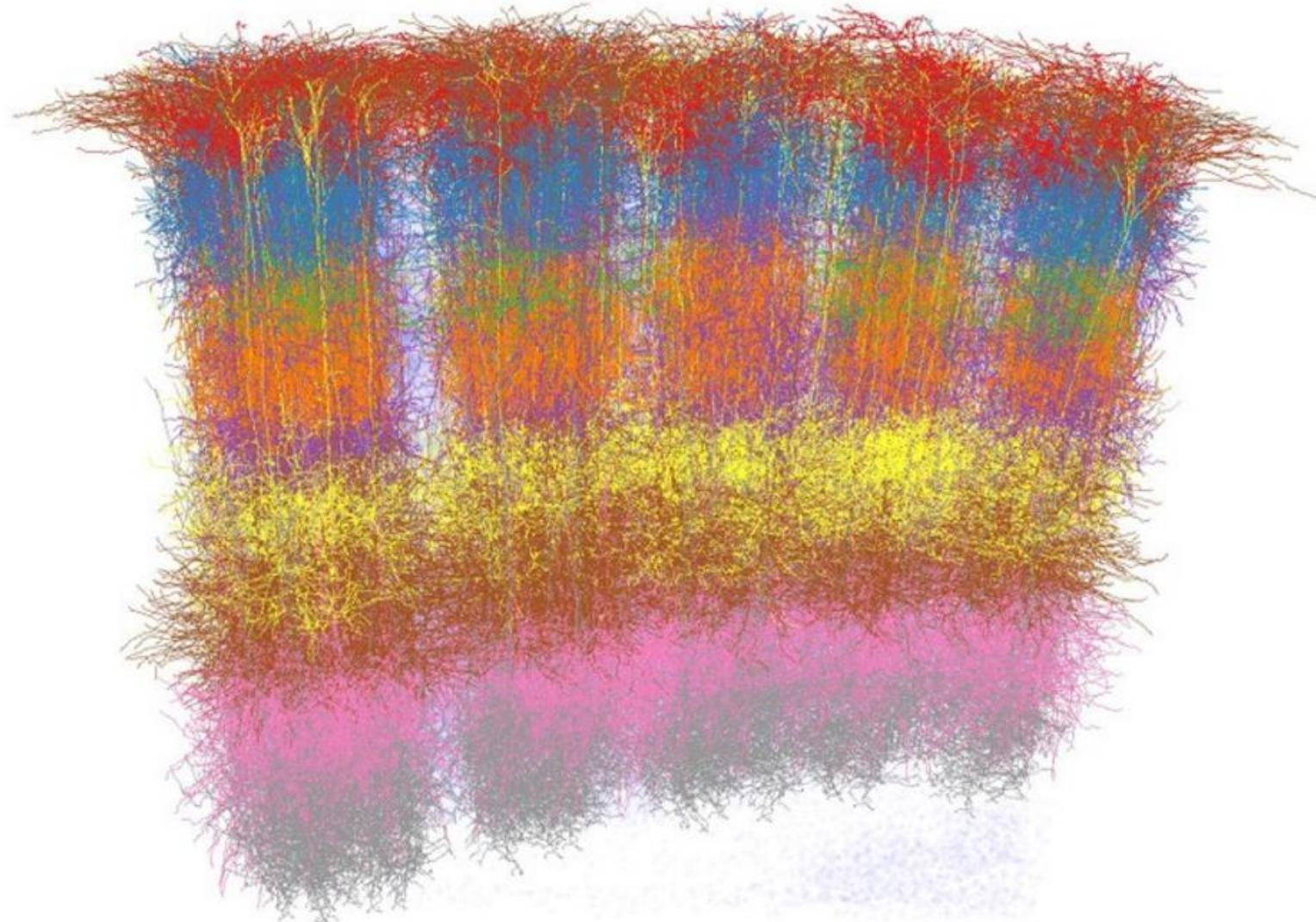
Trick: Use a **Neural Network** to Represent the Q-Table



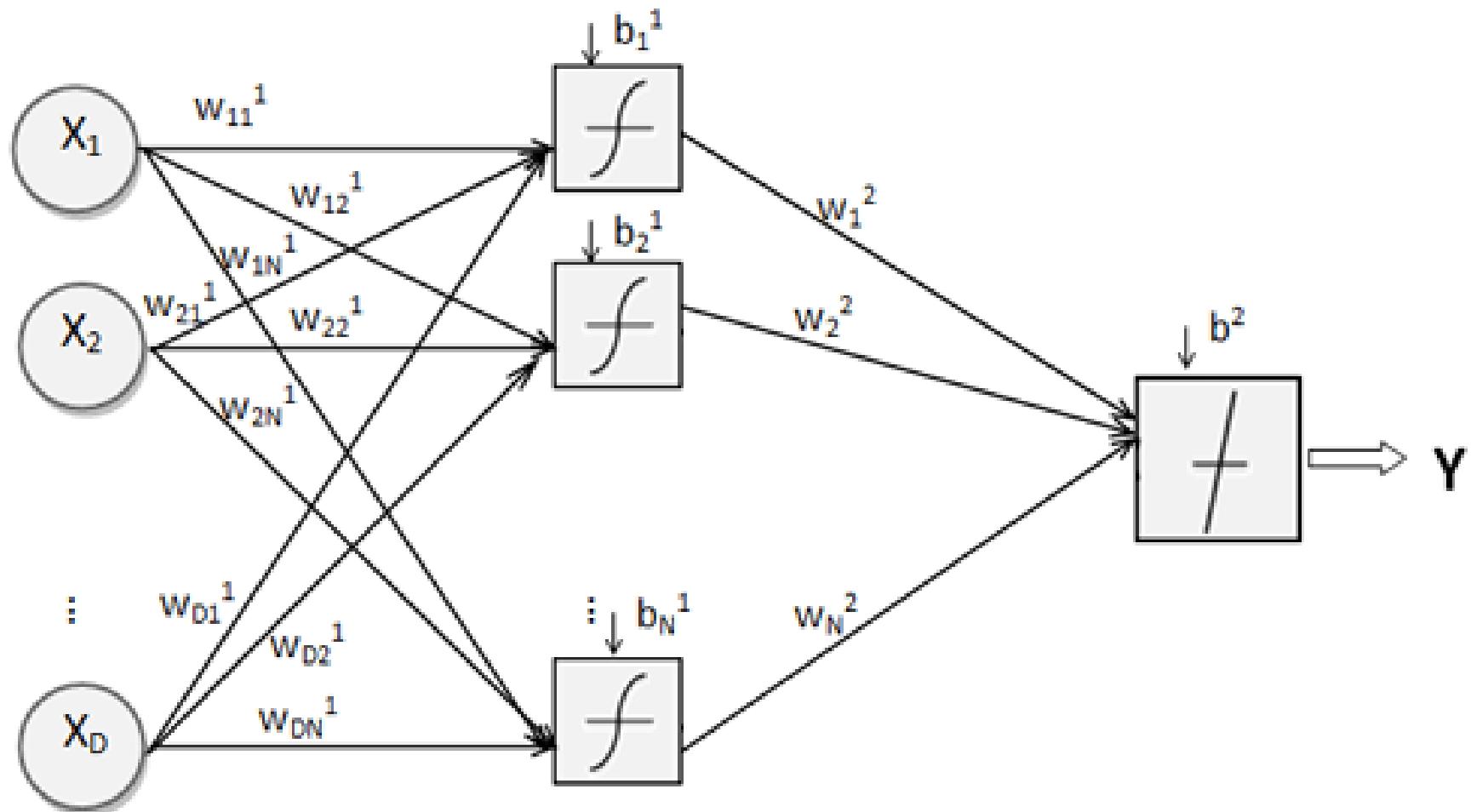
Neurons



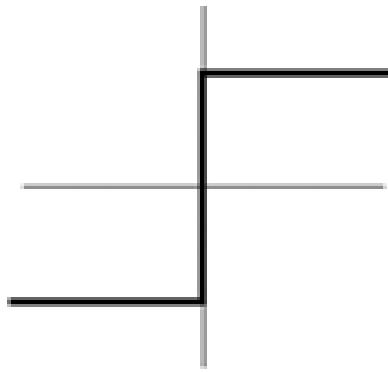
Neuron Network in Cortex



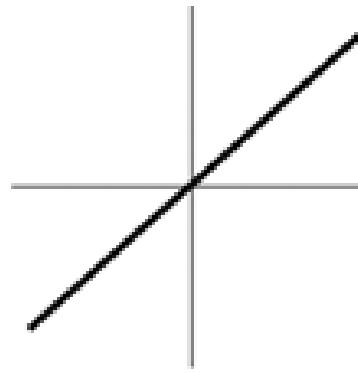
Artificial Neural Networks (ANNs)



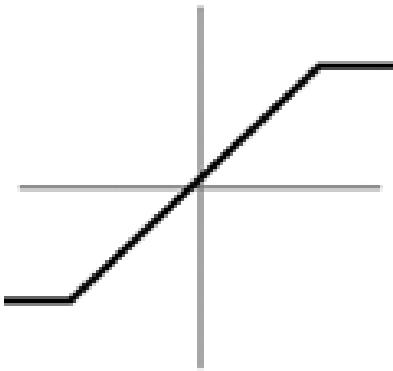
Popular Threshold Units



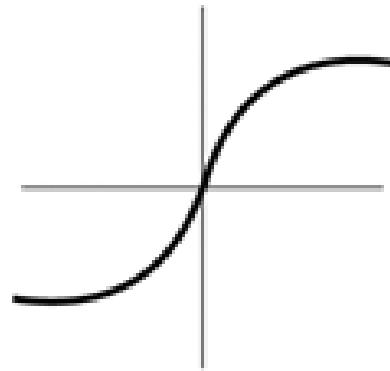
Step Function



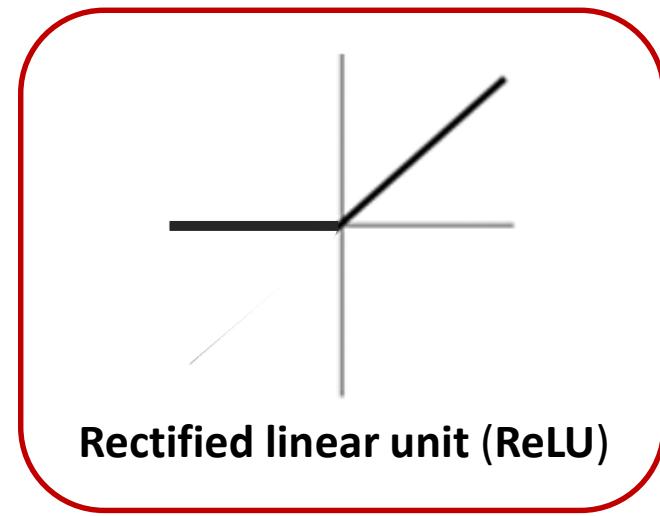
Linear Function



Threshold Logic



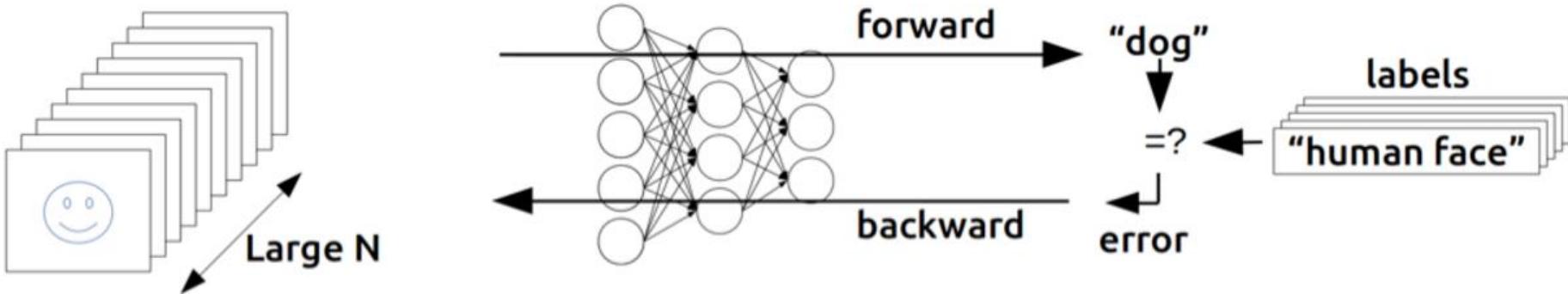
Sigmoid Function
tanh()



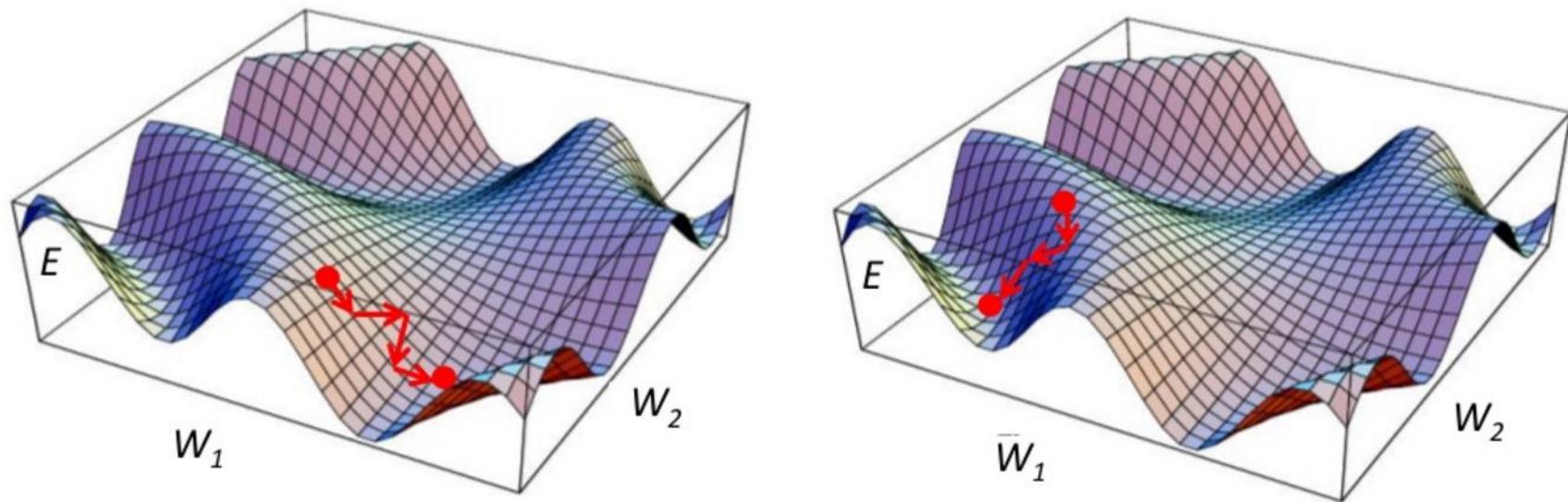
Rectified linear unit (ReLU)

<http://www.yaldex.com>

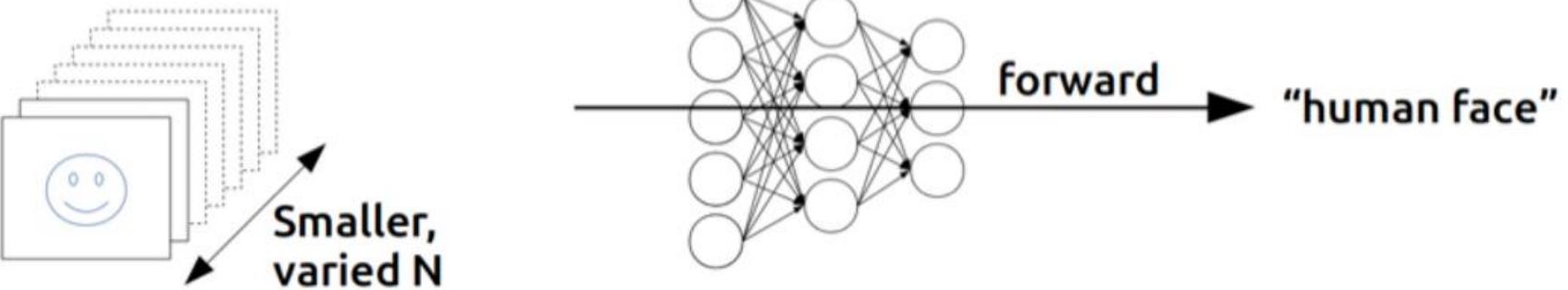
Training an ANN



Gradient Descent



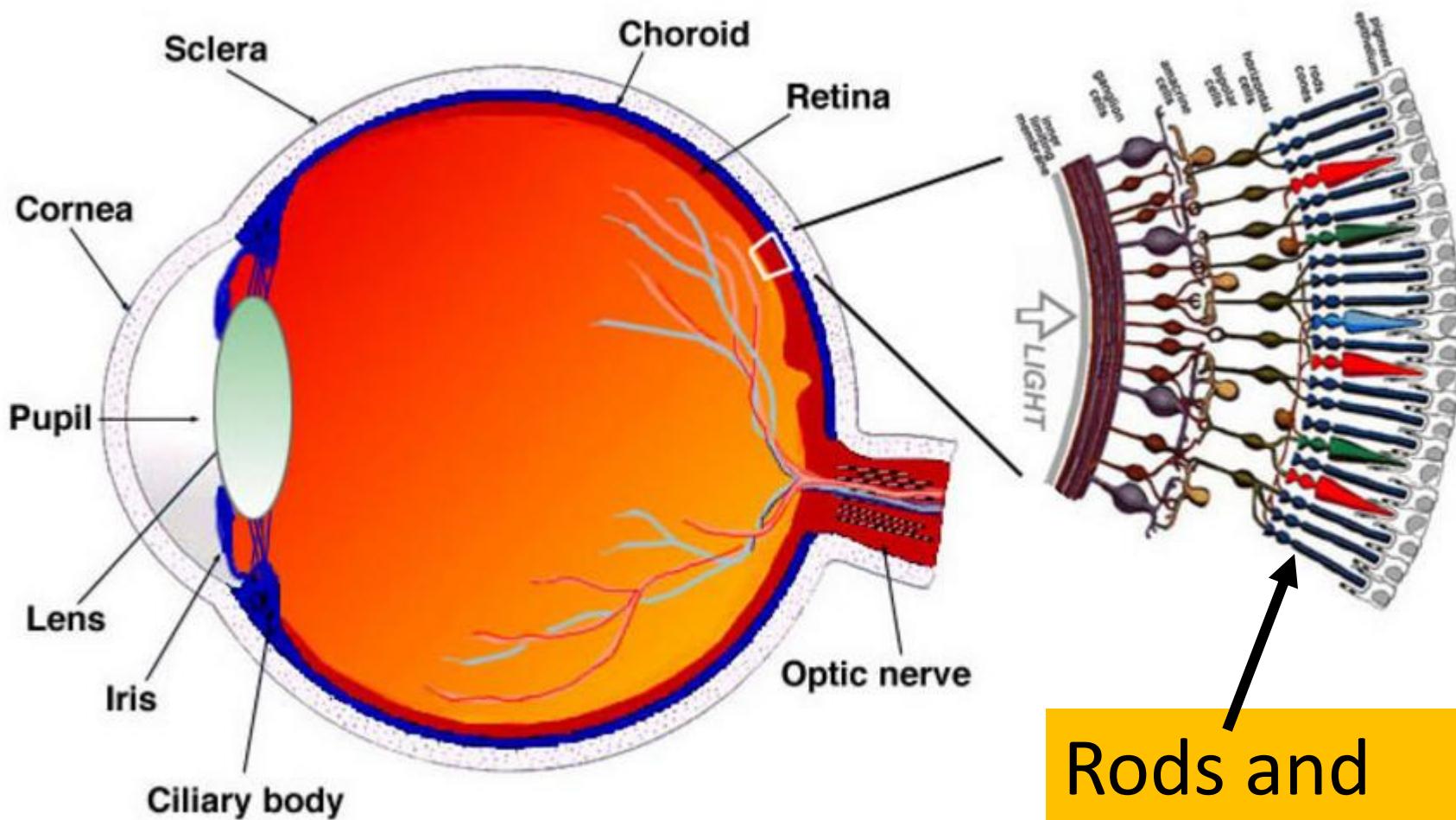
Classifying with an ANN



Practical Challenges

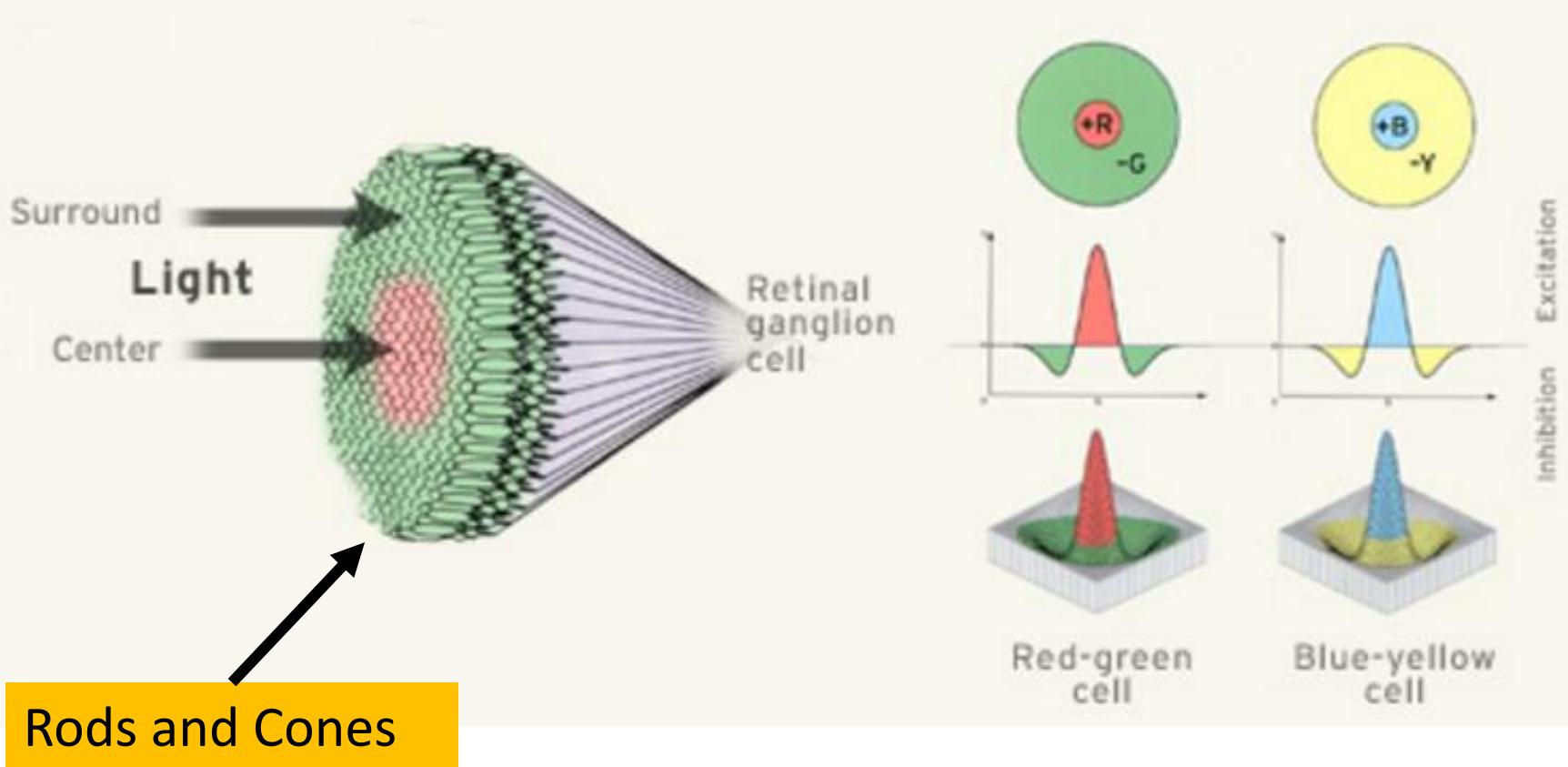
- **Curse of dimensionality:** too many weights to learn
- In image classification, full connectivity overkill
- Idea: reduce number of weights by imitating **visual perception**

Visual Processing Starts in Retina

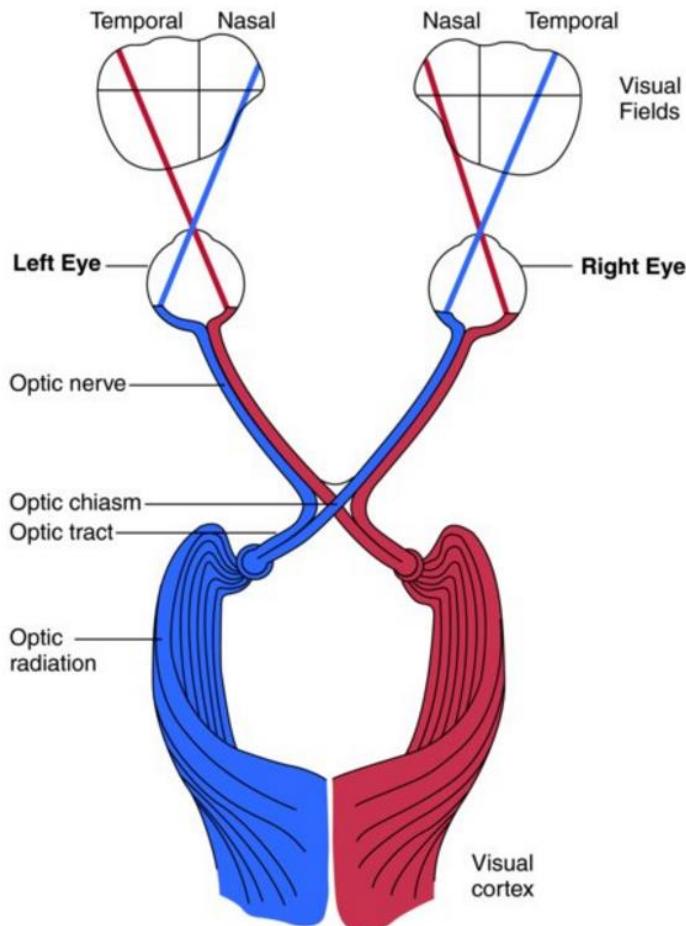


Rods and
Cones

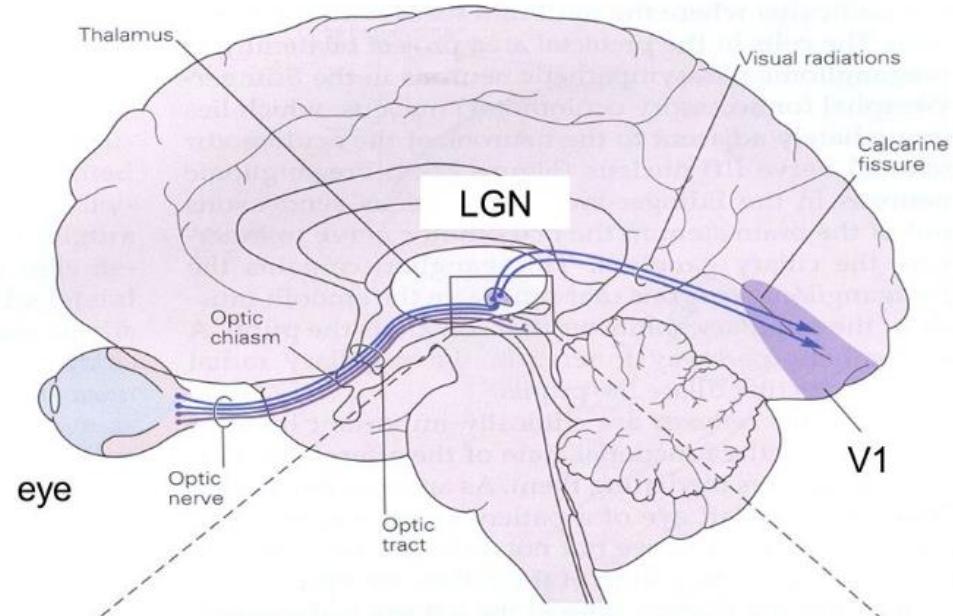
Receptive Field of Ganglion Cells



Central Visual Pathways

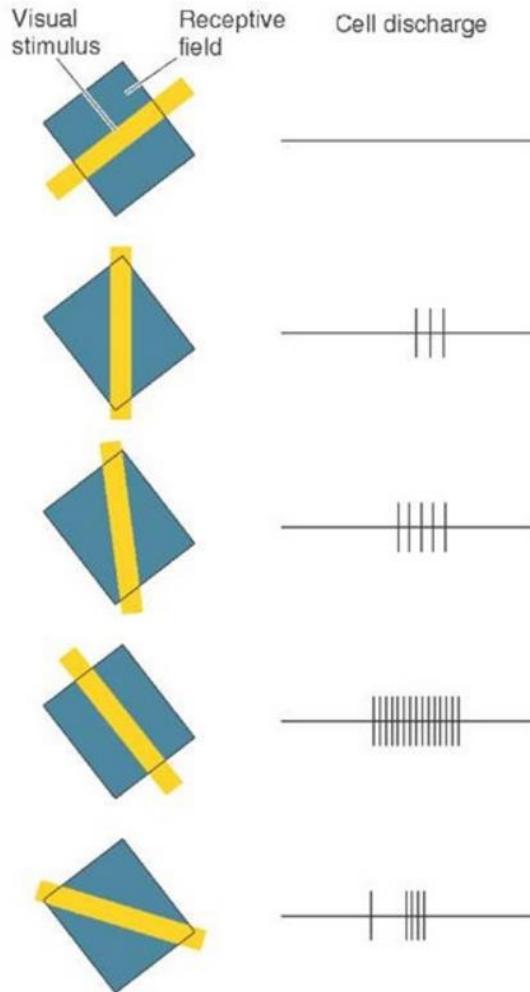


Top View



Side View

Receptive Field of Simple Cells in V1

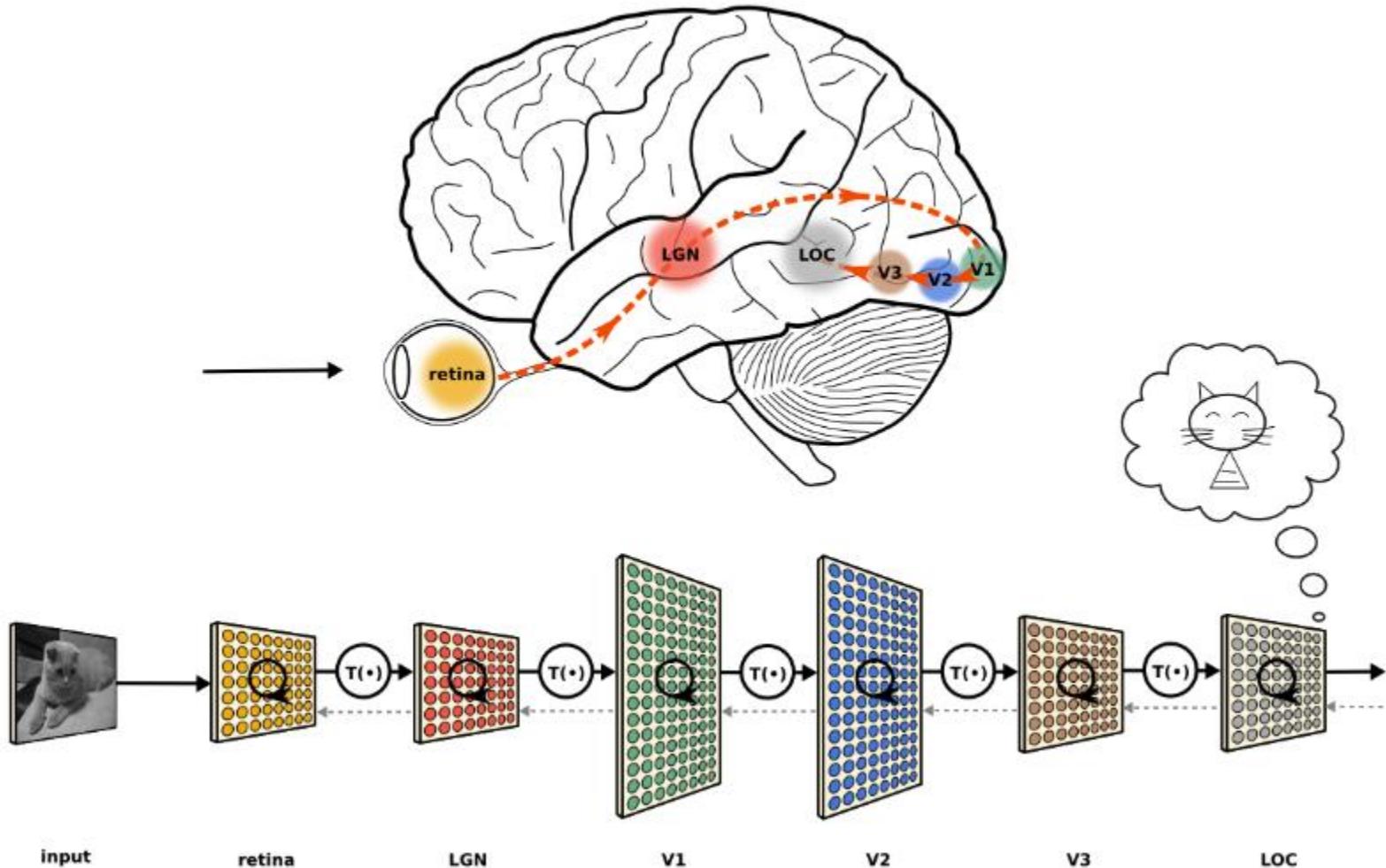


Conclusions

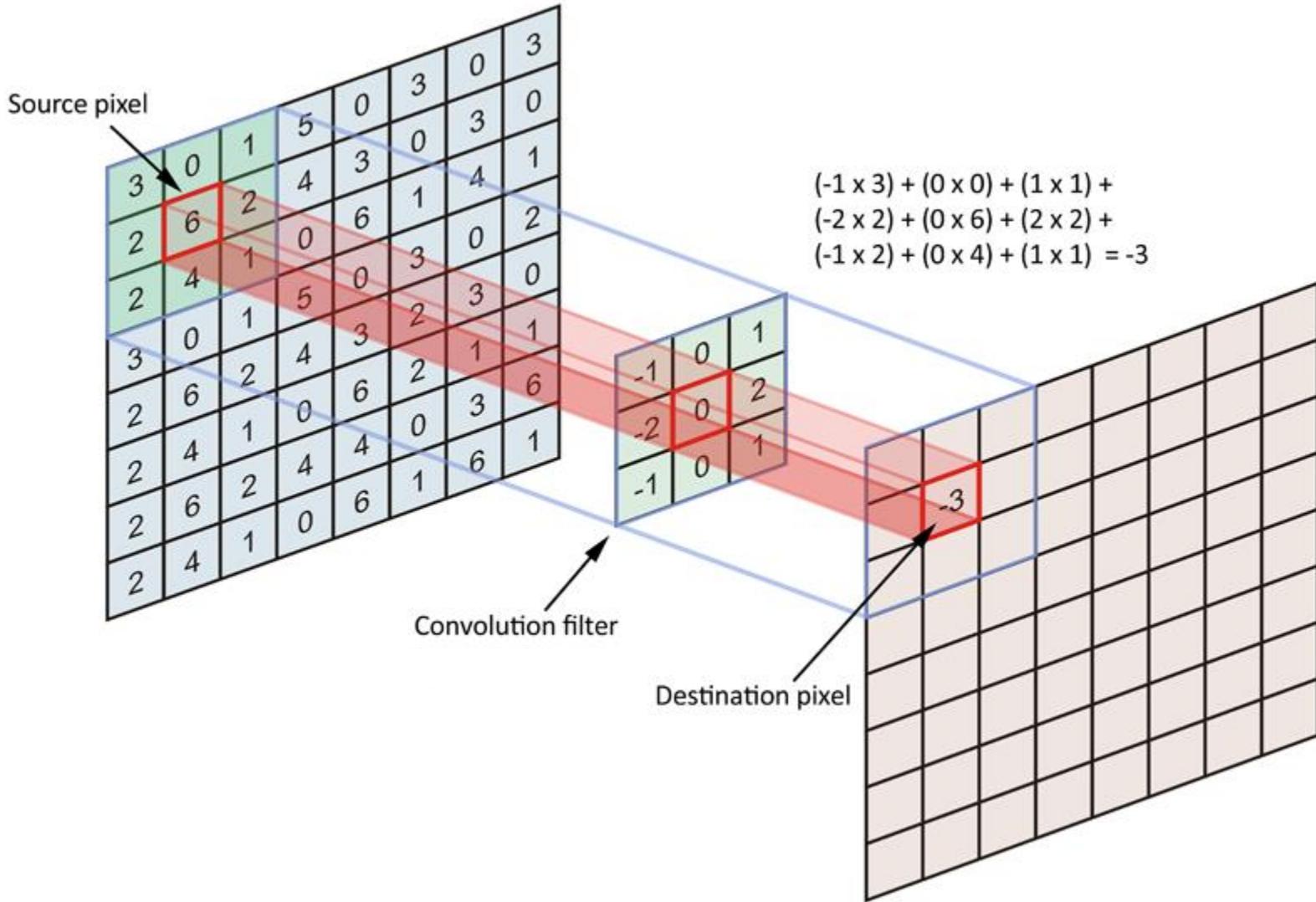
- Retina and V1 are topologically organized
- Ganglion and simple cells have receptive fields equal to **convolution filters**

Modelling the Visual Pathway

neuwritesd.org

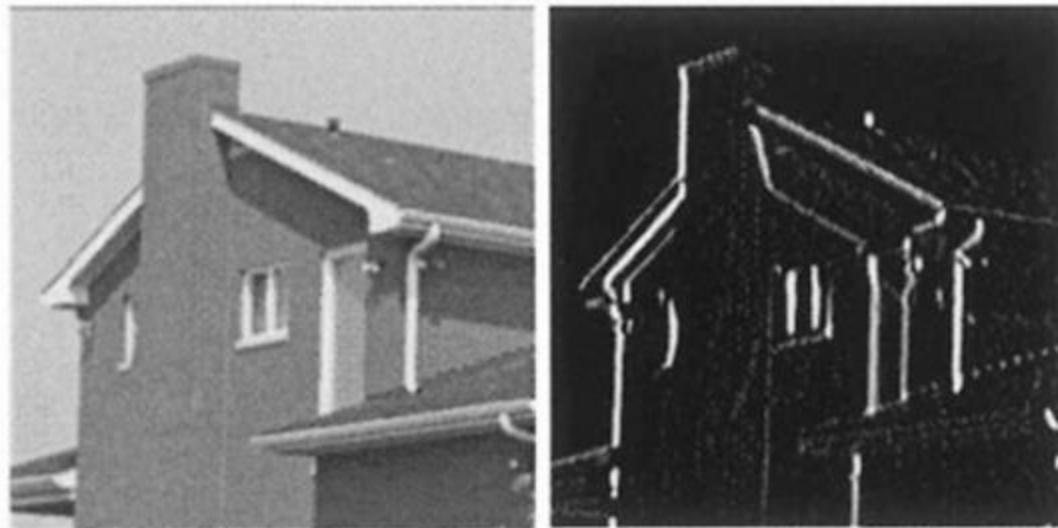


Convolution Filter



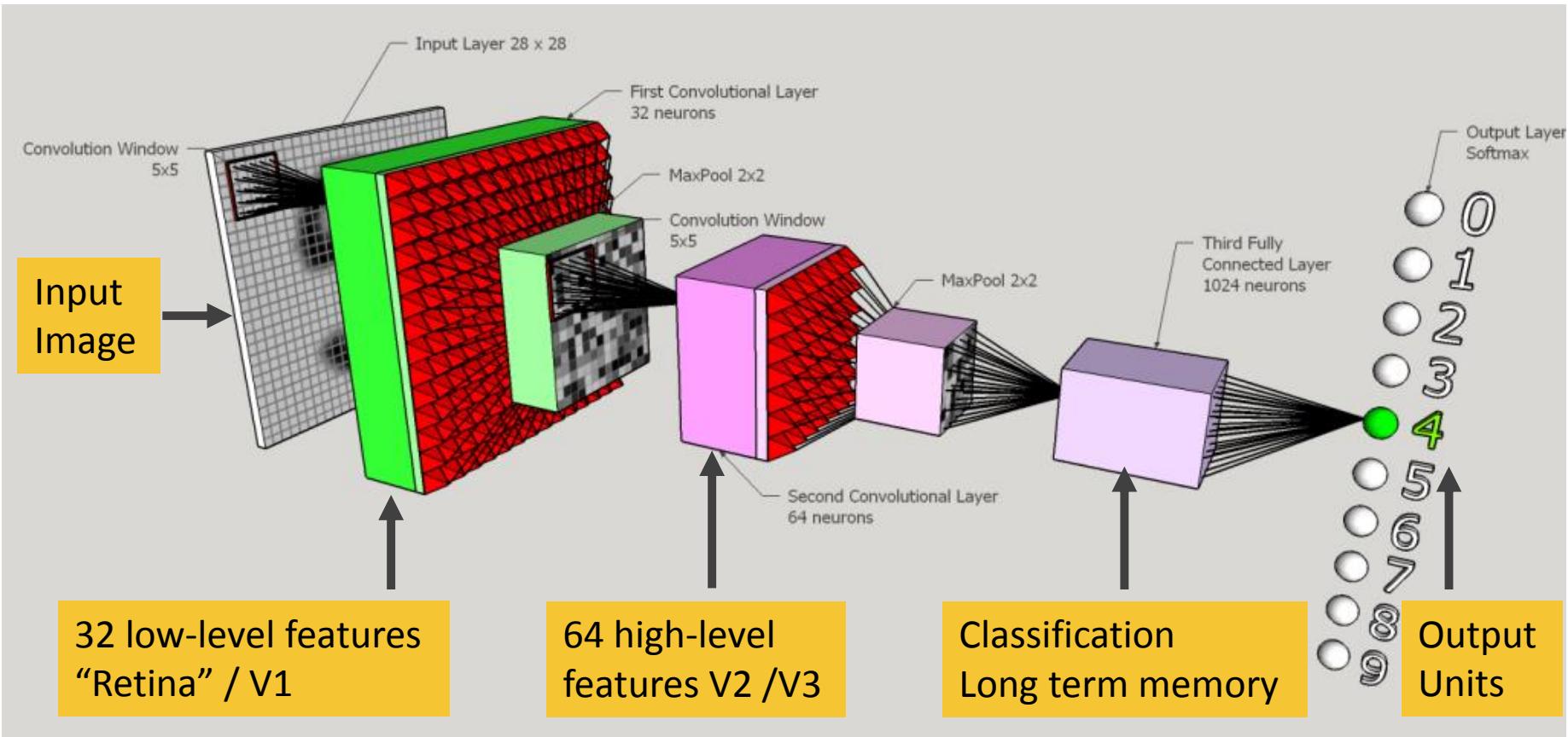
Example: Vertical Edge Detection

$$\begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix}$$



what-when-how.com

Convolutional Neural Networks (CNNs)



zderadicka.eu



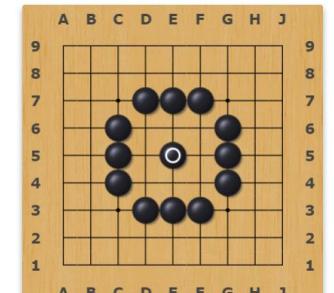
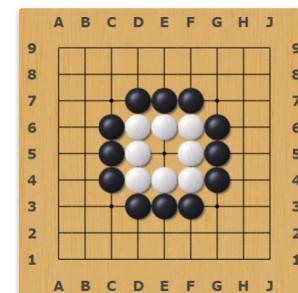
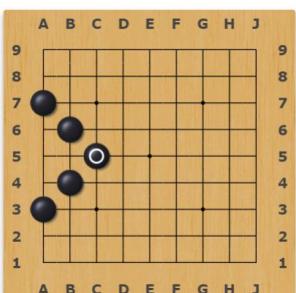
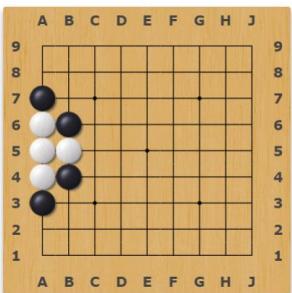
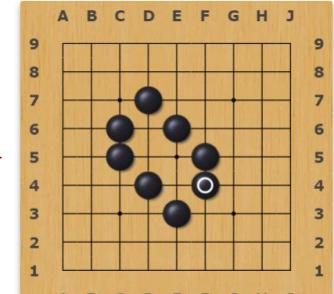
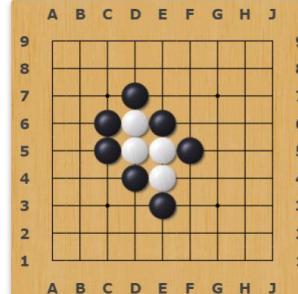
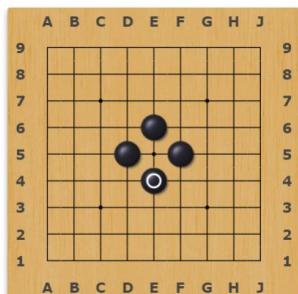
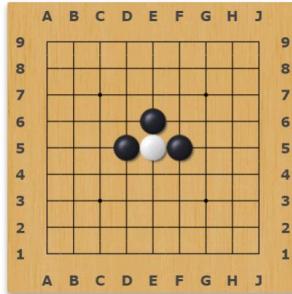
Break



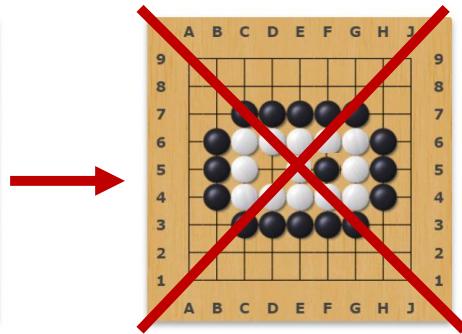
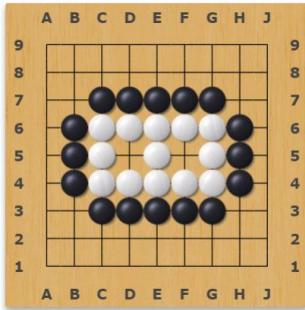
Go: A Zero-Sum Turn-Taking Game



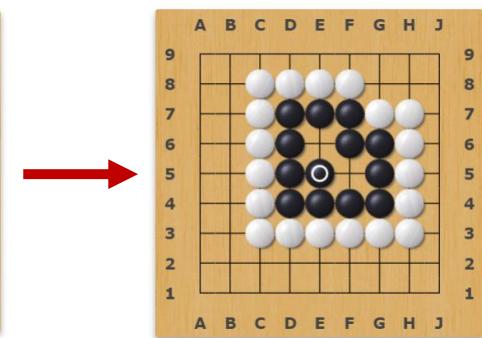
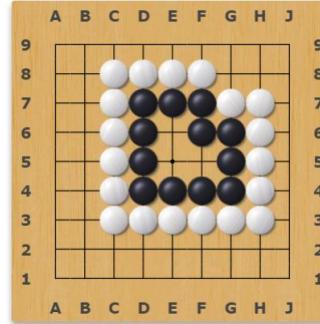
Captures



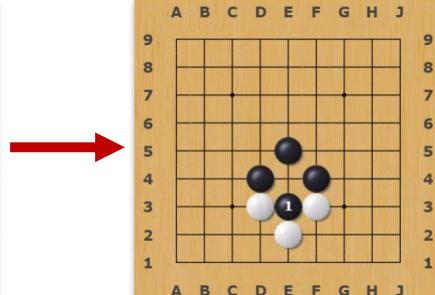
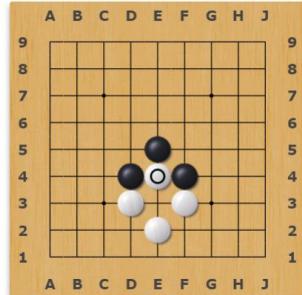
Capture Protection and Ko Rule



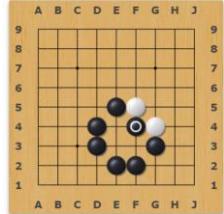
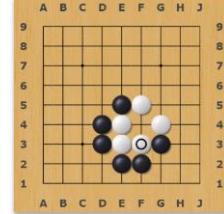
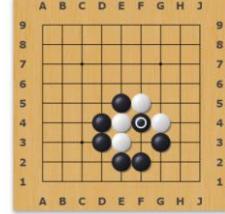
No capture with two holes!
Black F5 will suicide off board...



Capture protection



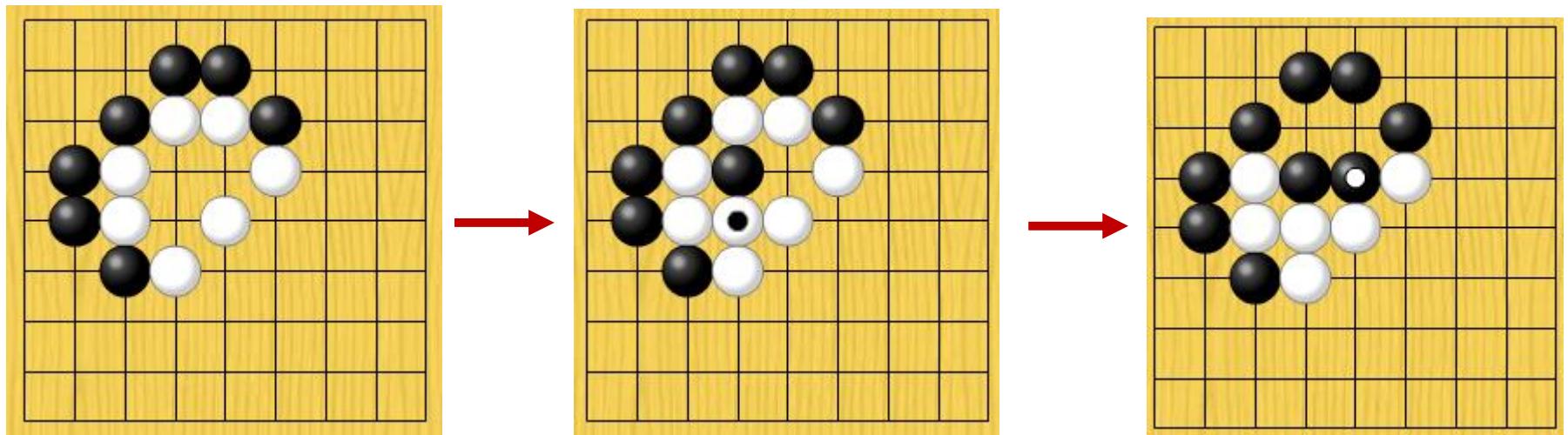
Ko rule: White cannot do E4 in next move, but ok later



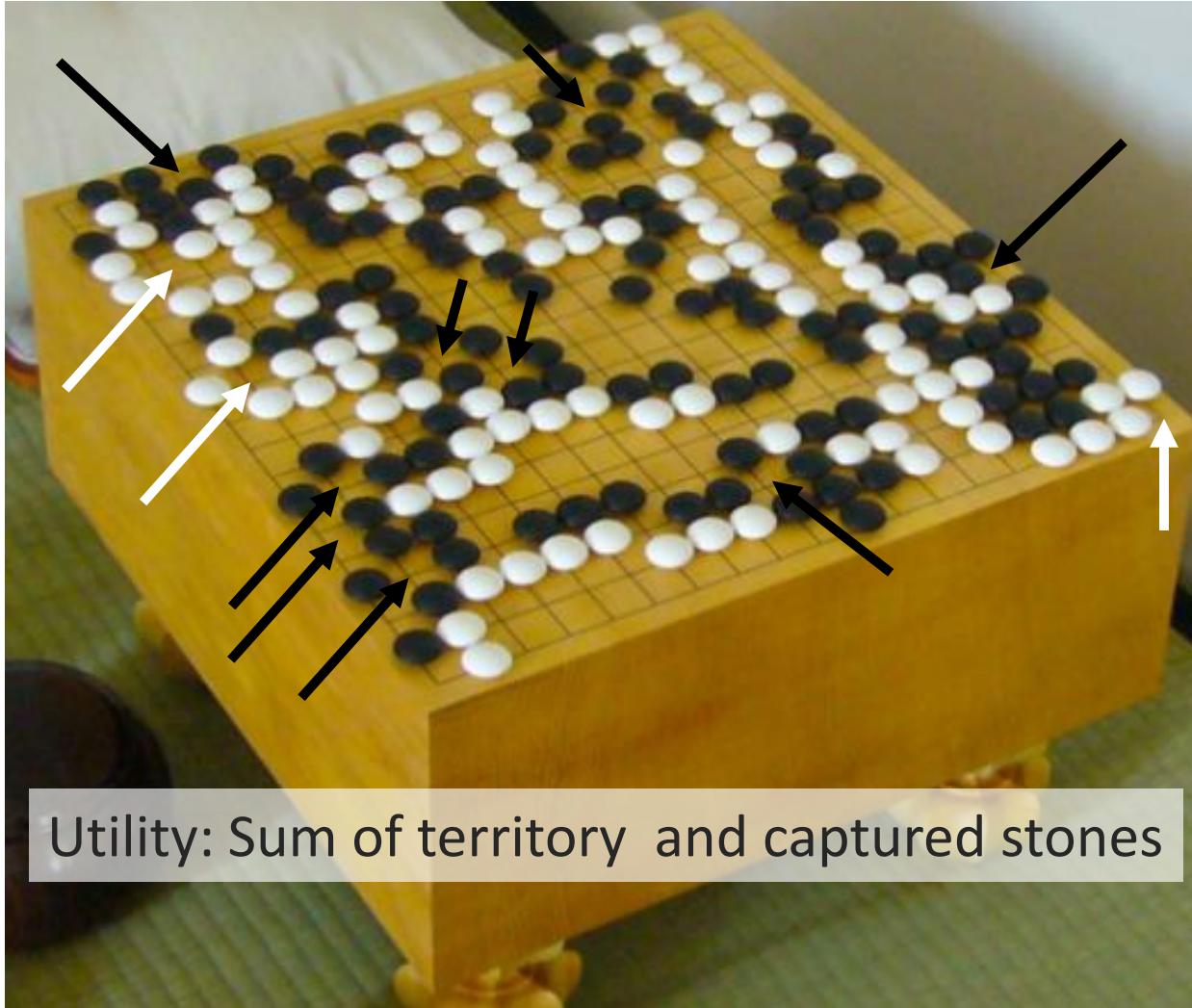
No Ko rule applies, since no direct reverse to a previous state



Double Atari Example



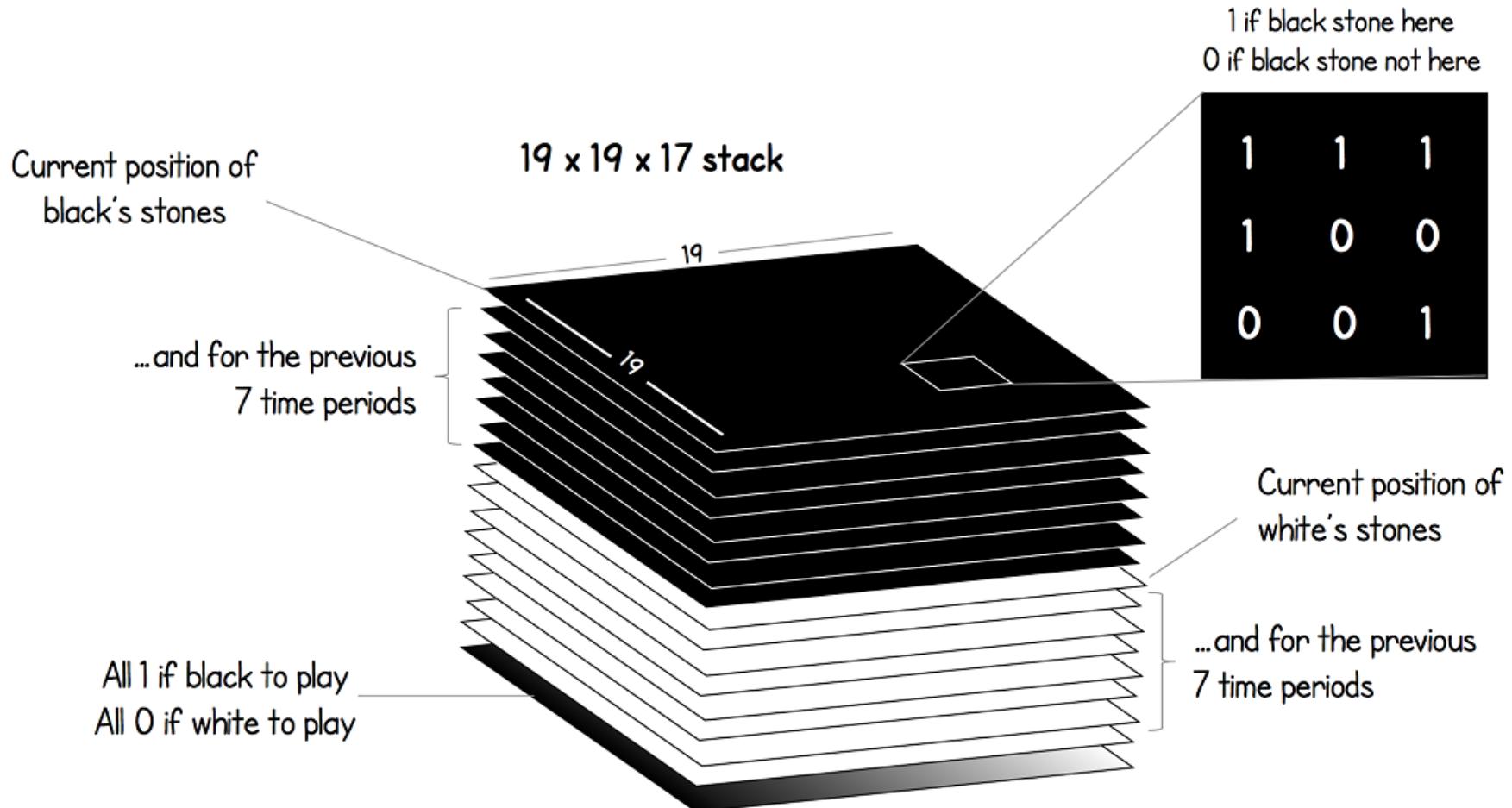
Final State Example



AlphaGo Zero's Neural Network

- Input / Output: $f_{\theta}(s) \rightarrow (p, v)$
 - s Board state including history (last 7 moves)
 - P Best move probabilities for s
 - v Board value of s
 - θ Arc weights of neural network

NN Architecture: Input

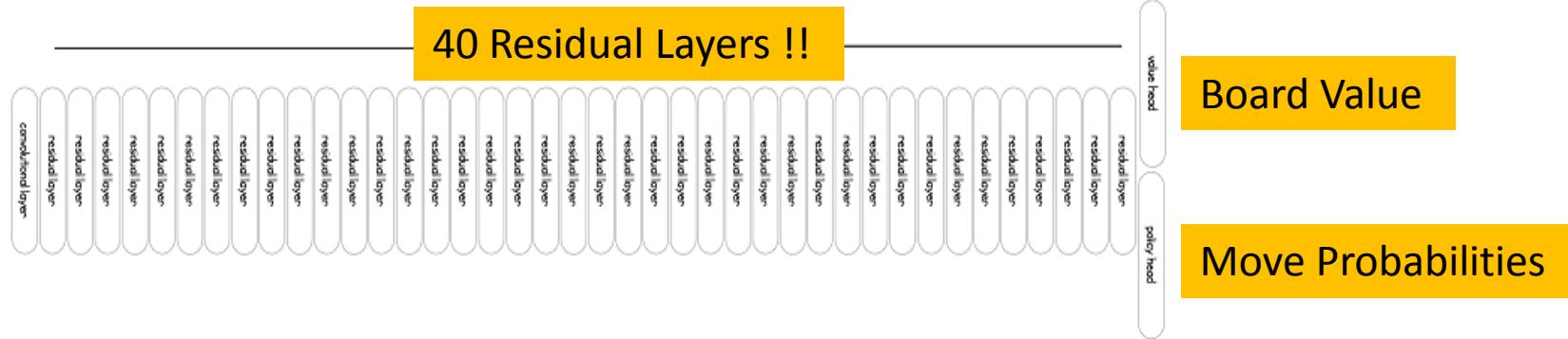


medium.com



Deep Network

Input



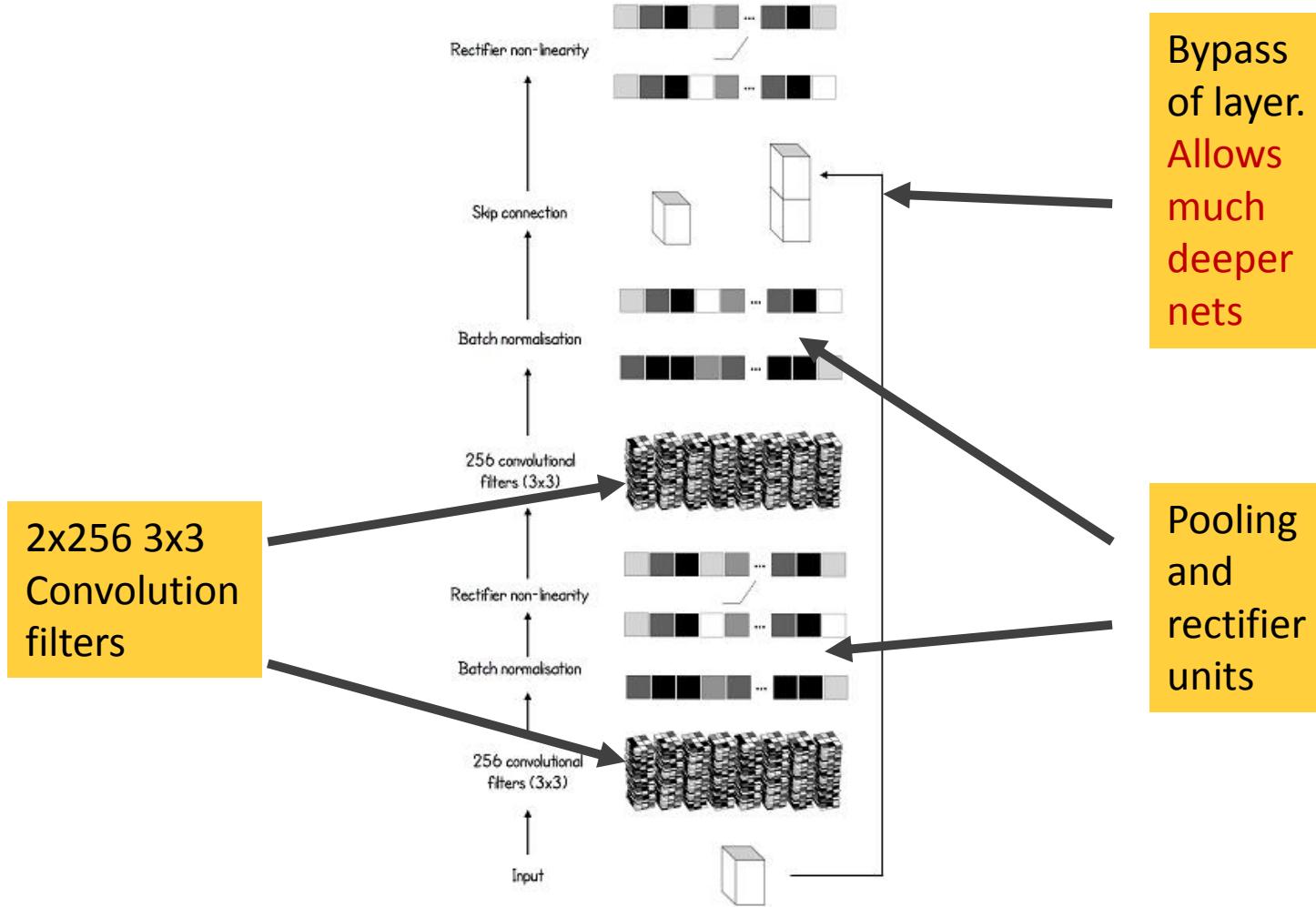
medium.com



IT University of Copenhagen

© Rune Møller Jensen

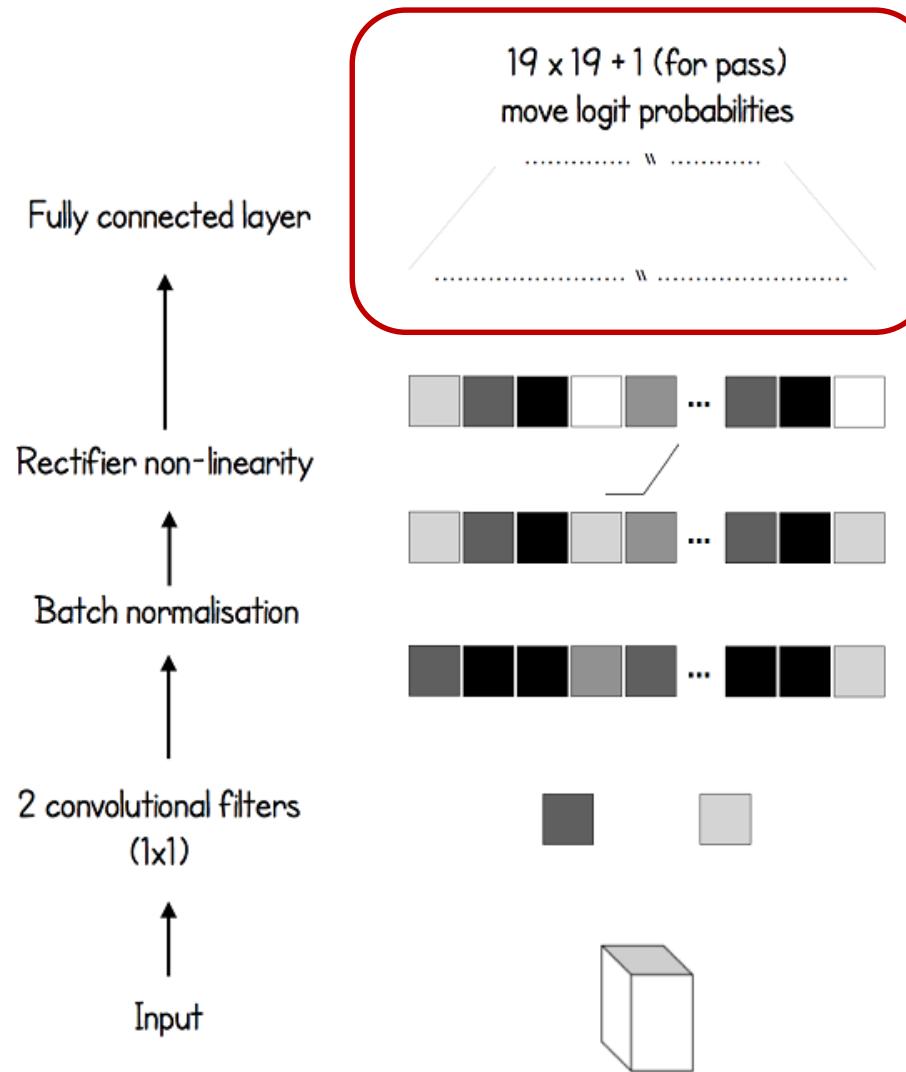
A Residual Layer



medium.com



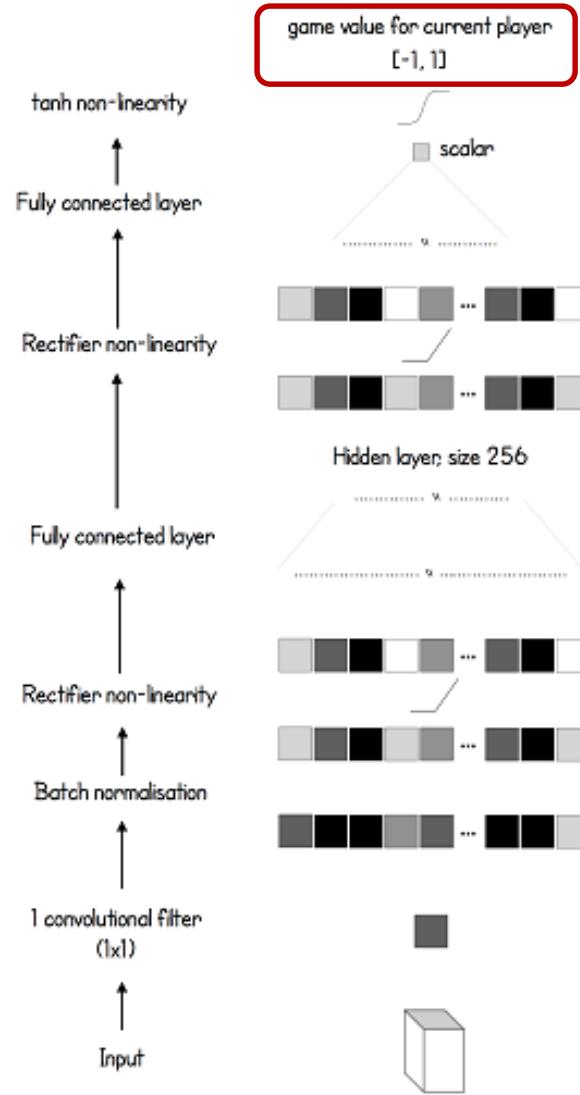
Move Probability Output



medium.com



Board Value Output

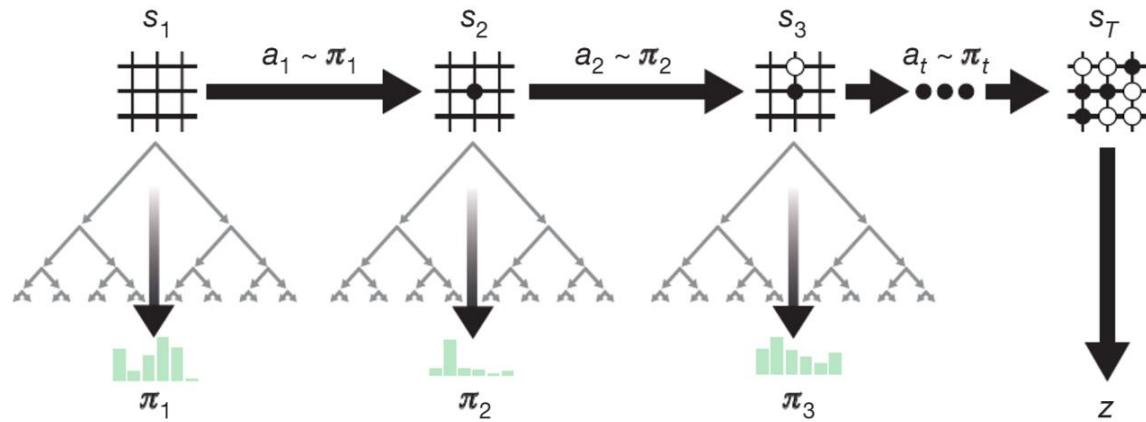


medium.com

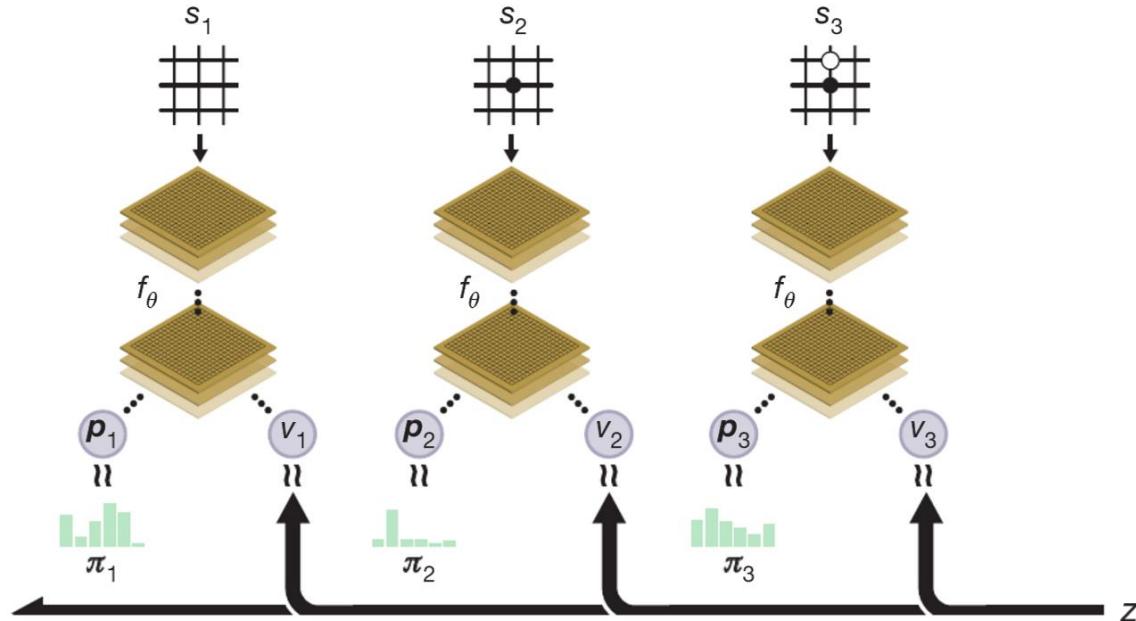


Self-Play Reinforcement Learning

Self-play



Neural
network
training



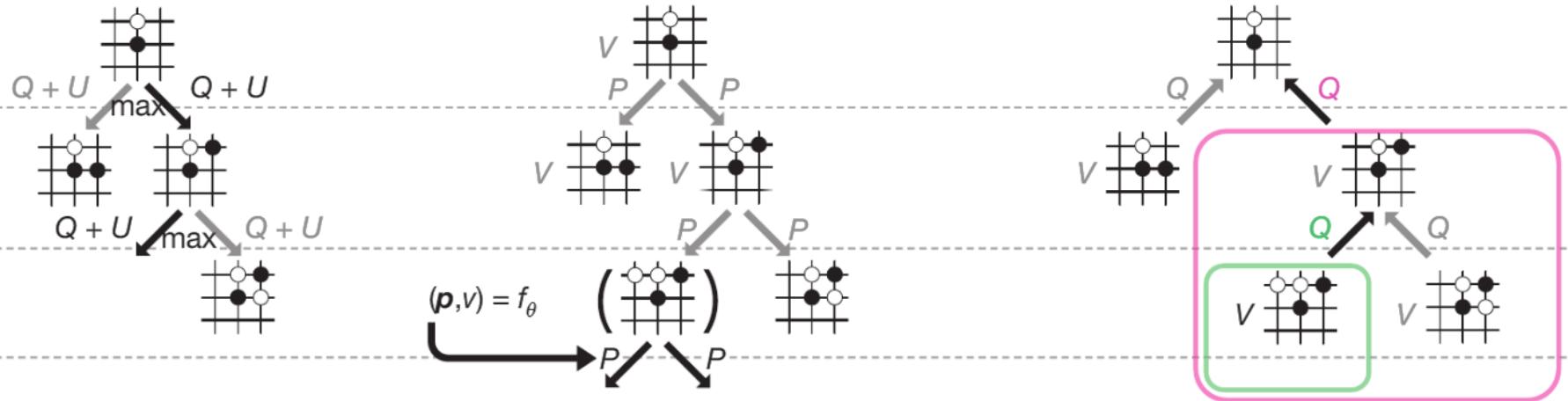
Monte Carlo Tree Search (MCTS)

Select

Expand and evaluate

Backup

Repeat 1600 times



$N(s,a)$ Number of times action a taken from state s

$P(s,a)$ NN probability of action a in state s

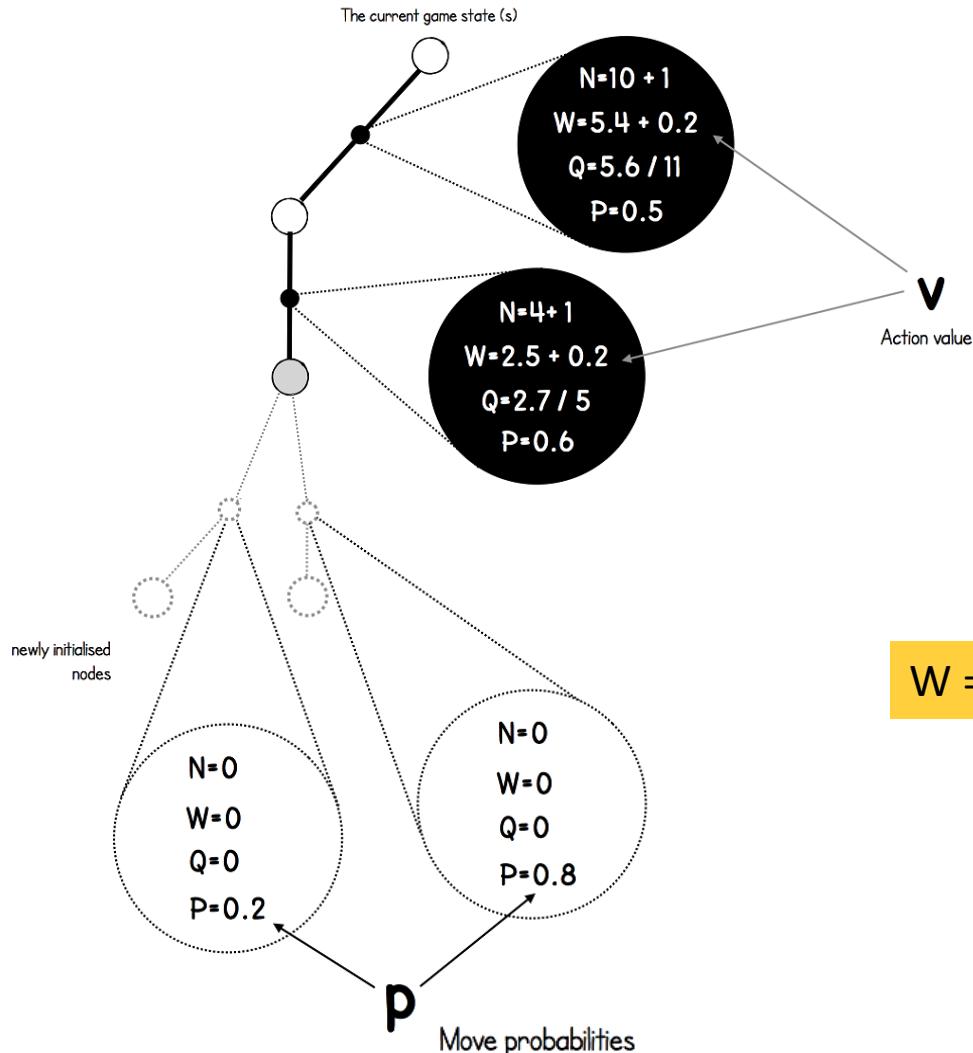
$Q(s,a)$ Mean value of next state

$U(s,a) \propto P(s,a) / N(s,a)$

www.nature.com



MCTS Update Example

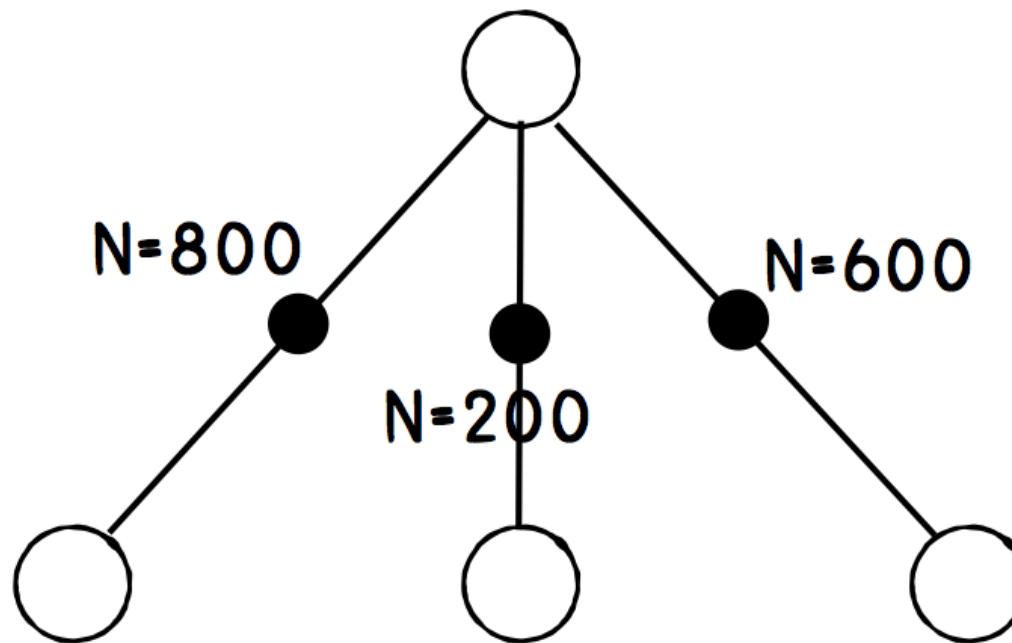


medium.com



Resulting MCTS Policy (π)

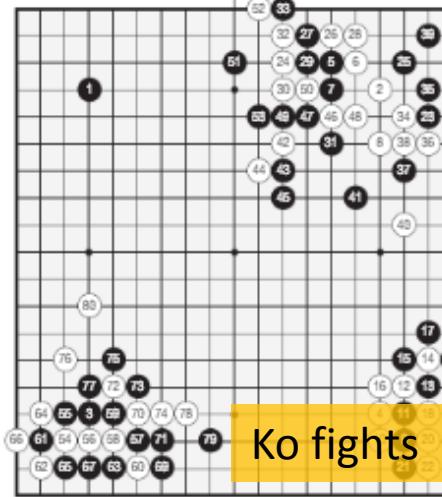
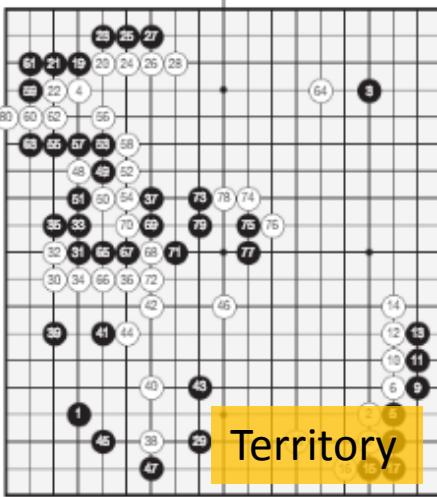
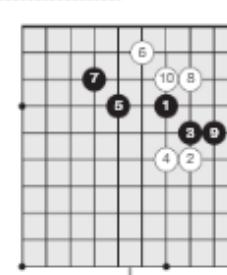
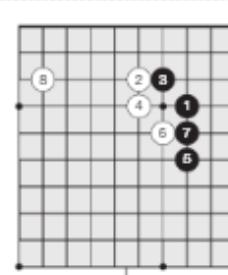
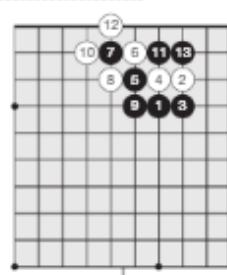
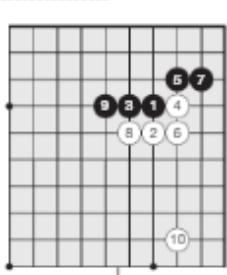
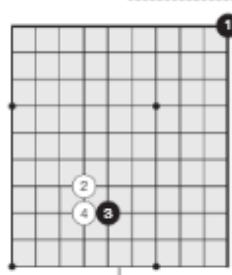
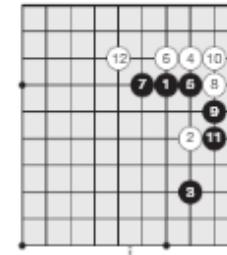
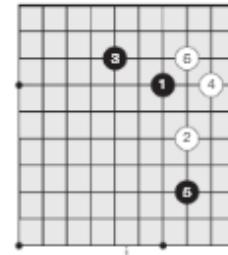
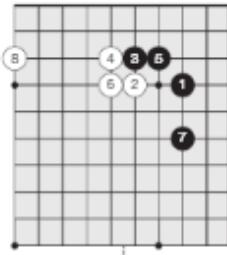
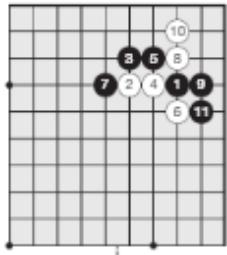
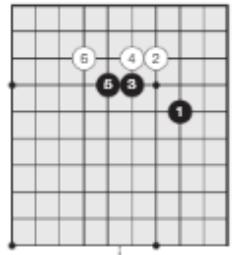
The current game state (s)



$$\pi = (800/1600, 200/1600, 600/1600) = (0.5, 0.125, 0.375)$$

Empirical Analysis

- NN initiated with random weights (tabula rasa)
- 4.9M games of self-play generated
- 1600 simulations of each MCTS (~ 0.4 sec)
- 72h on single machine with 4 TPUs



Capturing

Territory

Ko fights

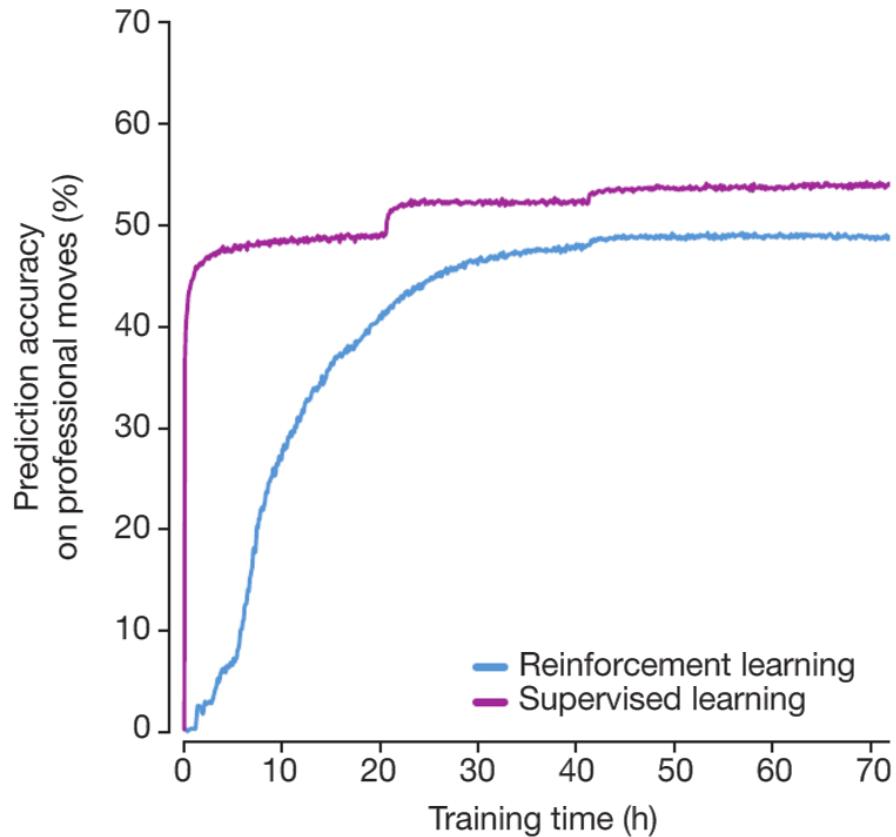
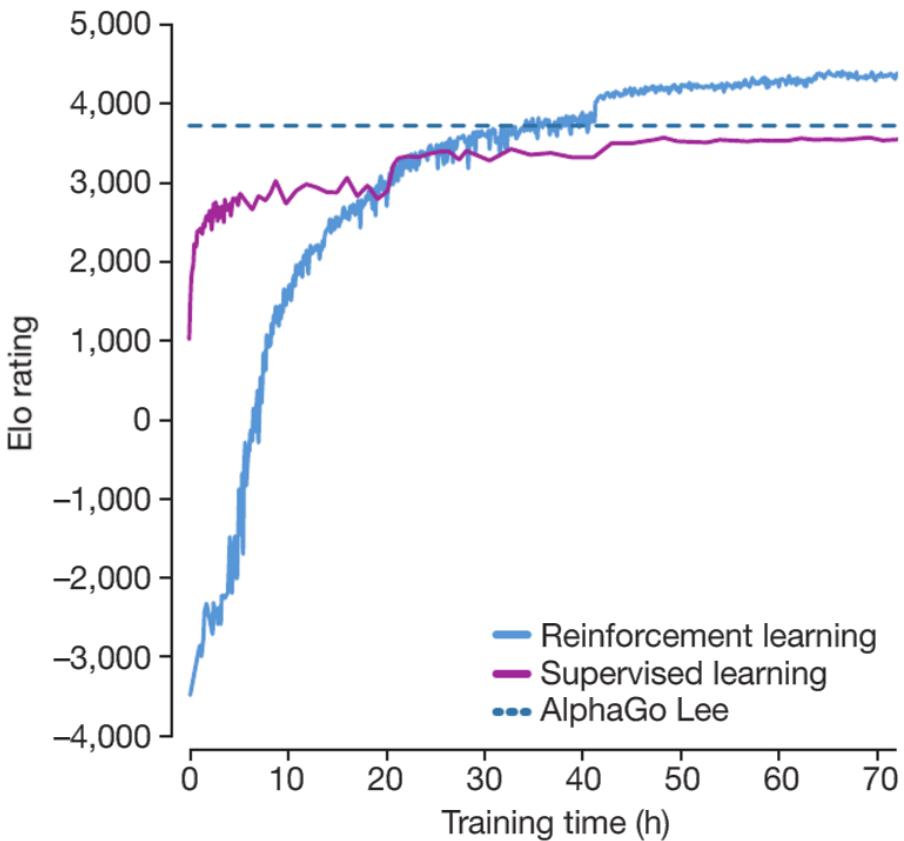
27 at 17 39 at 29 37 at 21 42 at 34 55 at 44 61 at 39
64 at 40 67 at 39 70 at 40 71 at 25 73 at 21 74 at 60
75 at 39 76 at 34

5 known
Joseki
discovered

Most
frequent
Joseki

Lee Sedol
“Divine
moves”
played by
AlphaGo

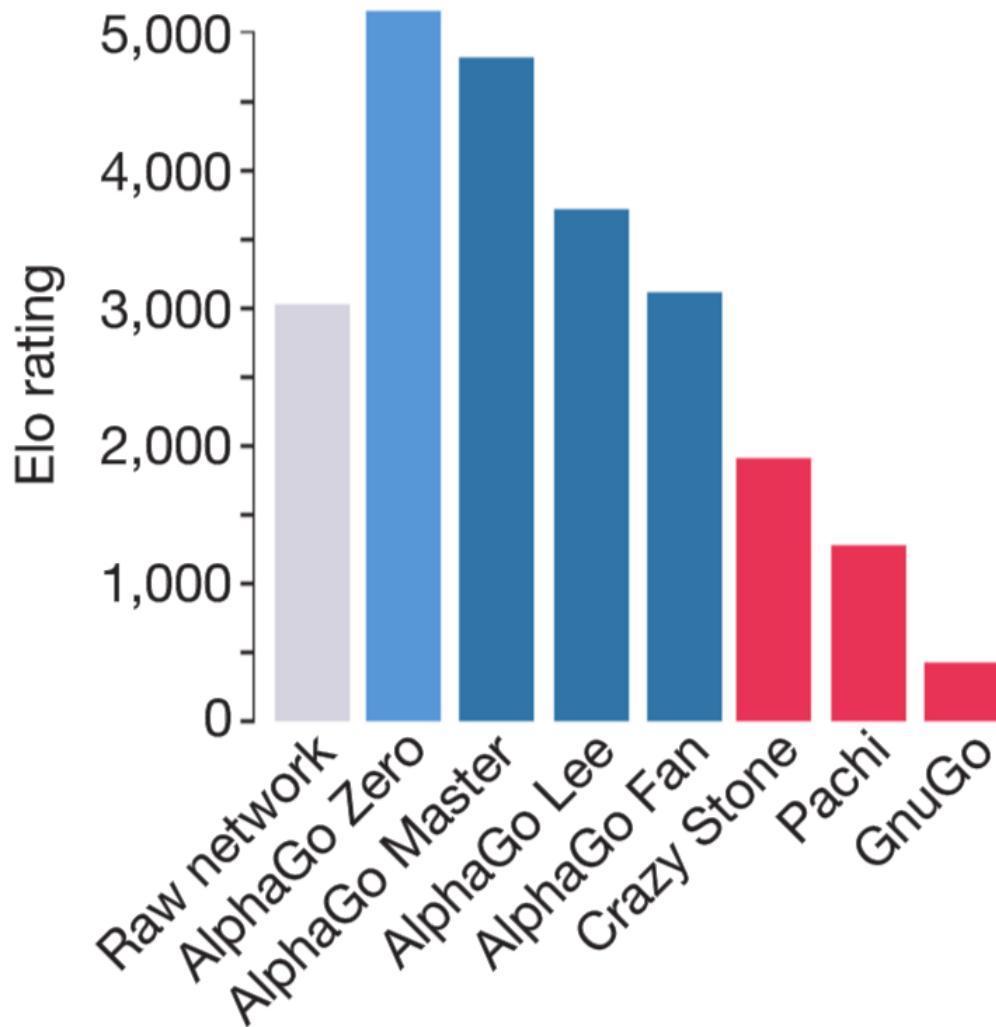
Sup-Human Level Achieved



www.nature.com



Related Work



www.nature.com

Conclusions

- For zero-sum turn-taking board games:
Alpha Zero > Minimax with Alpha Beta pruning
- Opportunities in
 - Autonomous control
 - Combinatorial optimization
- Issues to consider
 - Size of state-space: go 10^{200} , containership 10^{20000}
 - “Game length”: go 400 , containership 20000
 - “Self-play” possible to get training data?

