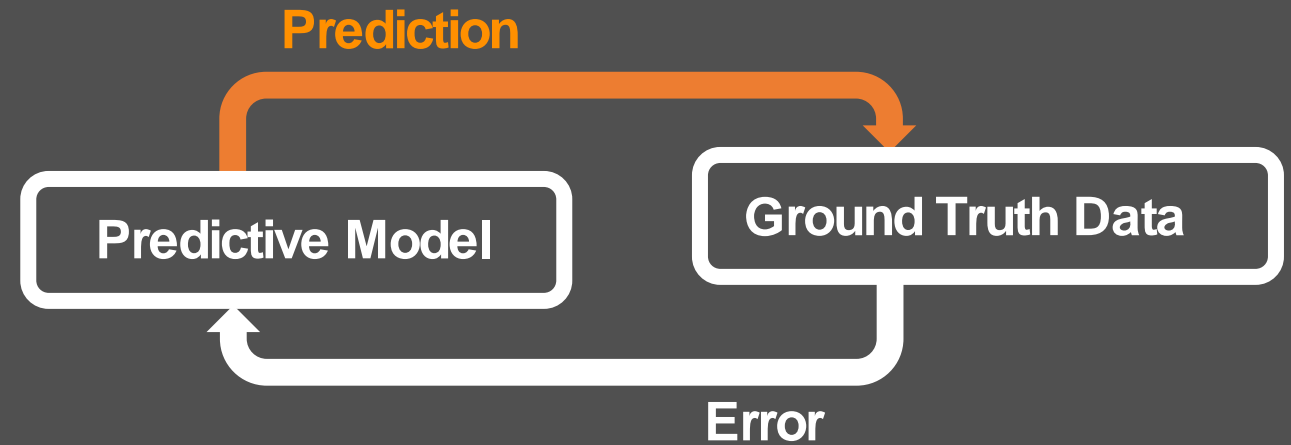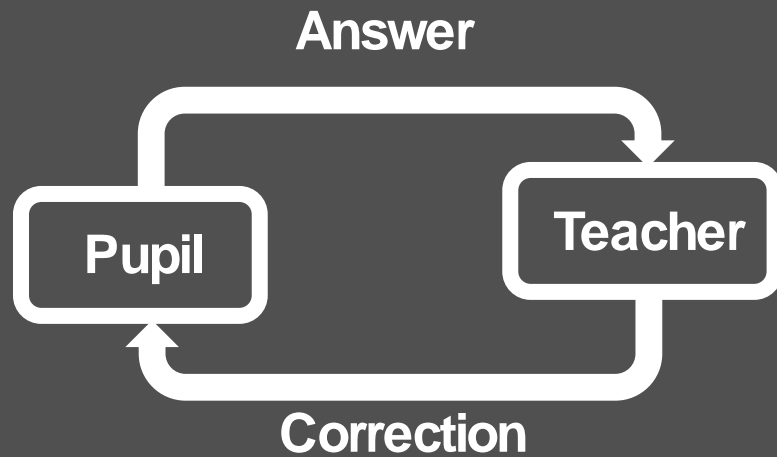# Reinforcement Learning

## Introduction
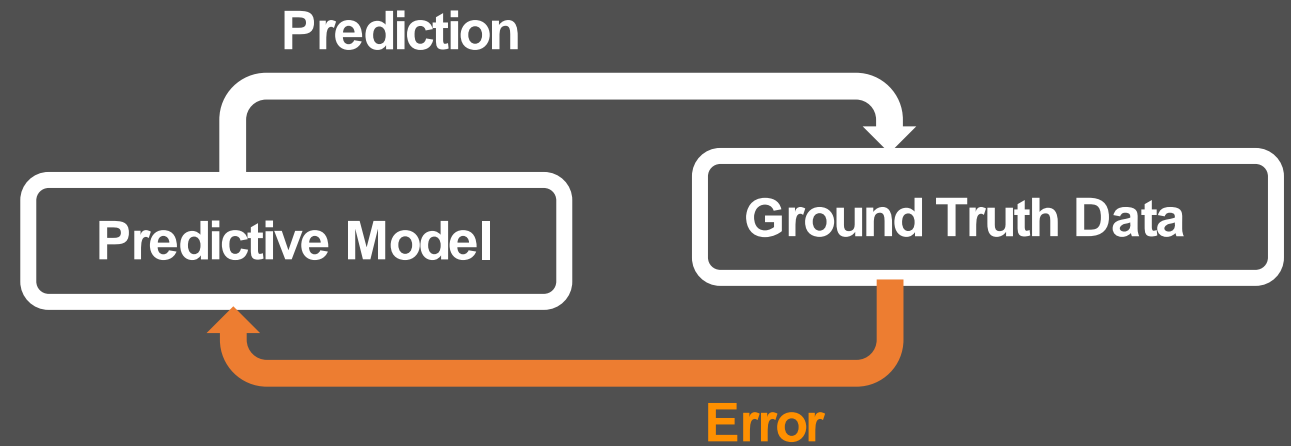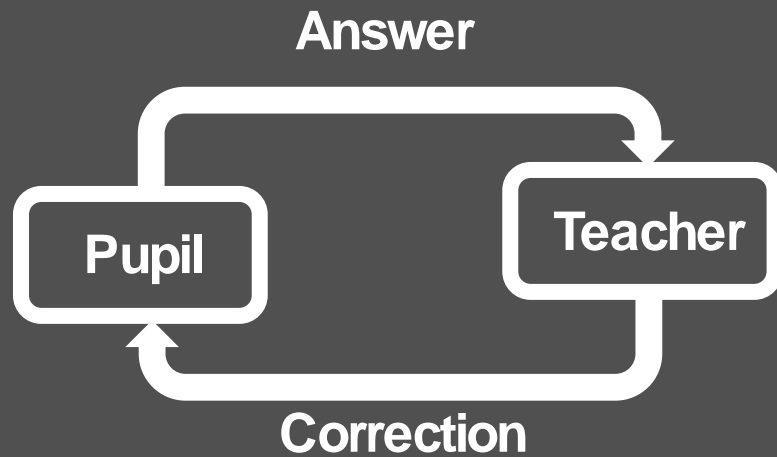
# Two Processes:

- **Forward Propagation (Prediction)**
- **Backward Propagation (Weight Tuning)**



Answer

Pupil → Teacher

Correction

Prediction

Predictive Model → Ground Truth Data

Error

# Two Processes:

- **Forward Propagation (Prediction)**
- **Backward Propagation (Weight Tuning)**



Answer

Pupil → Teacher

Correction

Prediction

Predictive Model → Ground Truth Data

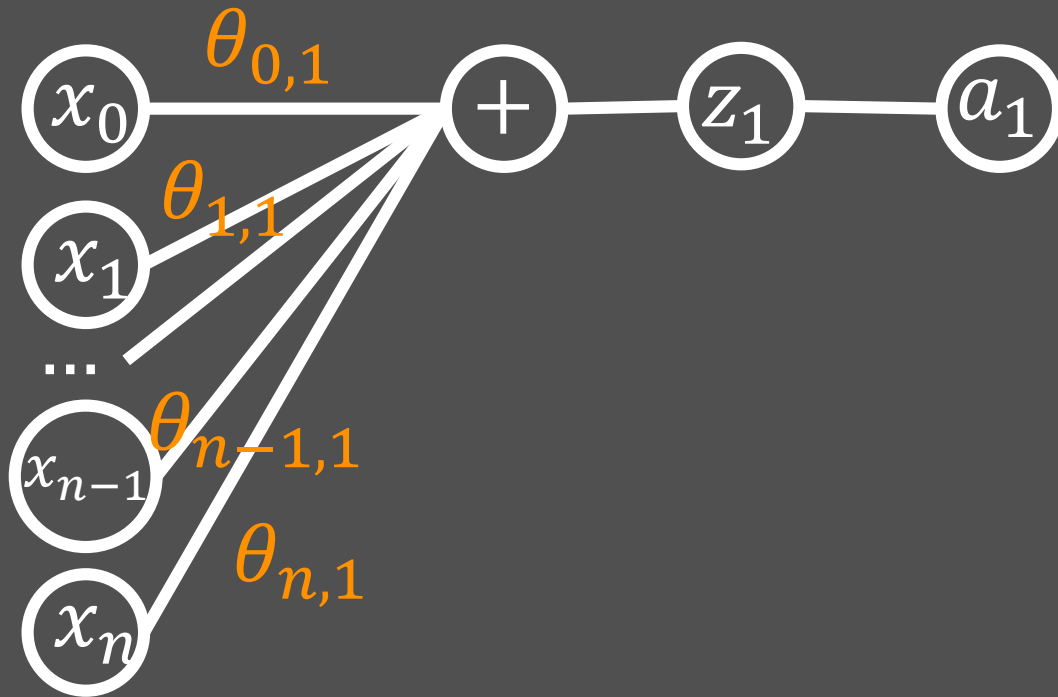Error

- **Compute multiple weighted sums $z_i$**
- **Activate them to get $a_i(z_i)$**
- **Remember from Perceptron: $z = \theta x$**

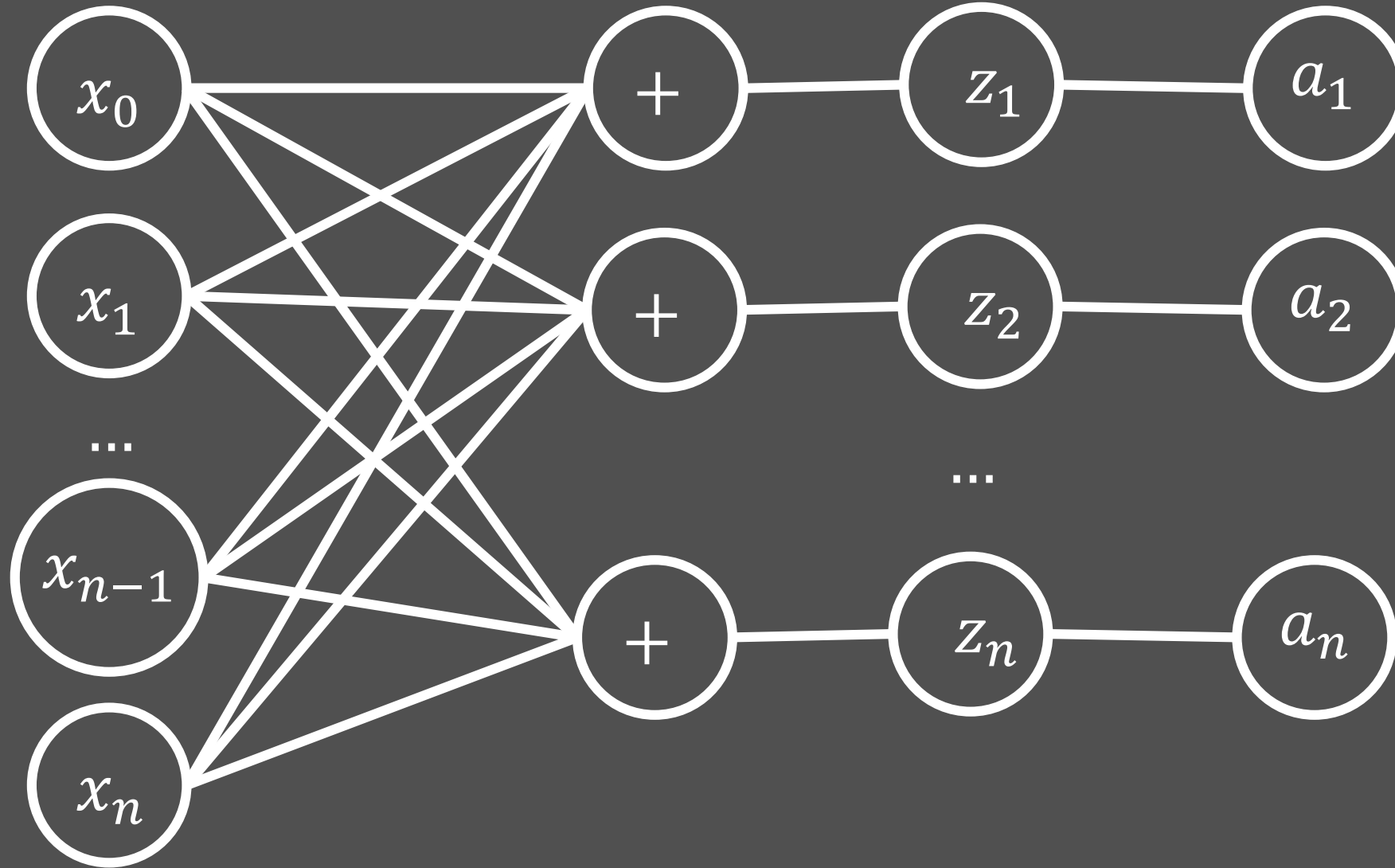$$z_{to} = \theta_{from,to} x_{from}$$

$$\text{e.g. } z_1 = \theta_{0,1} x_0$$

- **Object Recognition**
- **Assume 4 classes**
- **L-Layer (L is the last layer)**
- **Softmax:** $a = \dfrac{e^{z_j}}{\sum_{k=1}^{K} e^{z_k}}$
- **Use Softmax for the last Layer**

**K**

0.1   **Car**

0.05   **Airplane**

0.05   **Boat**

0.8   **Bike**

**Hidden Layer Activation Functions** $a(z)$**:**

- **Linear Neuron:** $a(z) = z$

- **Binary Threshold Unit:** $a(z) = \begin{cases} 0, z < 0 \\ 1, z \geq 0 \end{cases}$

- **Rectified Linear (ReLu):** $a(z) = \begin{cases} \mathbf{0, z < 0} \\ \mathbf{z, z \geq 0} \end{cases}$

- **Sigmoid:** $a(z) = \frac{1}{1+e^{-z}}$

- **Tanh:** $a(z) = \tanh(z)$

**Rectified Linear (ReLu):** $a(z) = \begin{cases} 0, z < 0 \\ z, z \geq 0 \end{cases}$

HSD
Hochschule Düsseldorf
University of Applied Sciences

Fachbereich Medien
Faculty of Media

MIREVI

## Exercise:

- **Construct a MLP (input: 3, hidden: 4, output: 3)**
- **All activation functions are ReLu-Function**
- **The output activation is a softmax-function**
- **Initialize the weights randomly**
  - **No zero values allowed!**
- **The input vector is** $x = (1, -2, 1)^T$
- **Compute one forward-propagation step!**

# What we will do?

- Applied AI
- Reinforcement Learning
- Train Agent to play Atari-Games

Reinforcement Learning / M.Sc. Marcel Tiator

# Reinforcement Learning

**Action**

**Agent** ← **Reward** ← **Environment**

**State**

# Reinforcement Learning

**Action**

**Reward**

**Agent**

**Environment**

**State**

# Reinforcement Learning

**Action**

**Agent** ← **Reward** ← **Environment**

**State**

# Reinforcement Learning

**Action**

**Reward**

**Agent**

**Environment**

**State**

# Reinforcement Learning

**Action**

**Reward**

**Agent**

**Environment**

**State**

# Reinforcement Learning

**Action**

**Agent**

**Reward**

**Environment**

**State**

# Rewards = Goals

Convolution     Convolution     Fully connected     Fully connected

No input

- **Apply Reinforcement Learning Strategies**
- **Optimize the Net**
- **Try and Error**
- **Net learns best Actions in certain States**

**Action**

**Agent**

**Reward**

**Environment**

**State**

Convolution → Convolution → Fully connected → Fully connected

No input
↑
↗
→
↘
↓
↙
←
↖
●
↑+●
↗+●
→+●
↘+●
↓+●
↙+●
←+●
↖+●

Convolution → Convolution → Fully connected → Fully connected

280 2 1

No input
↑
↗
→
↘
↓
↙
←
↖
●
↑ + ●
↗ + ●
→ + ●
↘ + ●
↓ + ●
↙ + ●
← + ●
↖ + ●

Action

Reward

Environment

Agent

State

# Deep Q(uality) Network



$$Q = \begin{pmatrix} -1 & 1 & 0 \\ -1 & 0 & -1 \\ 10 & 0 & -1 \end{pmatrix}$$

**Actions** 0 1 2

**States** 0 1 2

# Deep Q(uality) Network

**Value of take action 0 in state 0**

Actions

$$Q = \begin{pmatrix} 0 & 1 & 2 \\ -1 & 1 & 0 \\ -1 & 0 & -1 \\ 10 & 0 & -1 \end{pmatrix} \begin{matrix} \\ 0 \\ 1 \\ 2 \end{matrix}$$

States

HSD
Hochschule Düsseldorf
University of Applied Sciences

Fachbereich Medien
Faculty of Media

MIREVI

# Deep Q(uality) Network



**How good is it to take action 0 in state 0?**

Actions

$$Q = \begin{pmatrix} -1 & 1 & 0 \\ -1 & 0 & -1 \\ 10 & 0 & -1 \end{pmatrix} \begin{matrix} 0 \\ 1 \\ 2 \end{matrix}$$

States

(column headers: 0 1 2)

# Deep Q(uality) Network

**Q Table will be optimised by Neural Net**

$$Q = \begin{array}{c} \text{Actions} \\ \begin{array}{ccc} 0 & 1 & 2 \end{array} \\ \begin{pmatrix} -1 & 1 & 0 \\ -1 & 0 & -1 \\ 10 & 0 & -1 \end{pmatrix} \begin{array}{c} 0 \\ 1 \\ 2 \end{array} \end{array} \quad \text{States}$$

# Q-Learning Algorithm:

- **Estimate row in the Q table for the current state (Prediction)**

- **Apply action and get next state + reward**

- **Estimate row in the Q table for next state**

- **Apply Bellman Equation (Target):**
  - $Q(s, a) := r + \gamma * max_{a'} Q(s', a') = Q(s, a) * max_{a'} Q(s', a')$

- **Minimize Loss between Target and Prediction:**
  - $L = \frac{1}{2}(r + \gamma * max_{a'} Q(s', a') - Q(s, a))^2$

- **Estimate row in the Q table for the current state (Prediction)**

- **Apply action and get next state + reward**

- **Estimate row in the Q table for next state**

- **Apply Bellman Equation (Target):**
  - $Q(s,a) := r + \gamma * max_{a'} Q(s',a') = Q(s,a) * max_{a'} Q(s',a')$

- **Minimize Loss between Target and Prediction:**
  - $L = \frac{1}{2}(r + \gamma * max_{a'} Q(s',a') - Q(s,a))^2$

# Q-Learning:

- **Estimate row in the Q table for the current state (<span style="color:orange">Prediction</span>)**

- **Apply action and get next state + reward**

- **Estimate row in the Q table for next state**

- **Apply Bellman Equation (Target):**
  - $Q(s,a) := r + \gamma * max_{a'} Q(s',a') = Q(s,a) * max_{a'} Q(s',a')$

- **Minimize Loss between Target and <span style="color:orange">Prediction</span>:**
  - $L = \frac{1}{2}(r + \gamma * max_{a'} Q(s',a') - Q(s,a))^2$

$$Q(s, a) := r + \gamma * max_{a'} Q(s', a')$$

**Maximum future reward for this state and action is the immediate reward plus maximum future reward for the next state**

- **Minimize Loss between Target and Prediction:**
  - $L = \frac{1}{2}(r + \gamma * max_{a'} Q(s', a') - Q(s, a))^2$

- $r + \gamma * max_{a'} Q(s', a') \in R^n$
- $Q(s, a) \in R^n$
- **Q-Learning is very similar to regression**
  - Try to predict continous variables

# Implementation Tricks:

- **Use $\varepsilon$–Greedy Policy**

- **Experience Replay**
  - **Minibatch training samples from past experience**
  - **Avoid to get stuck in a local minimum**

- **Discount Future Reward with $\gamma$**
  - $Q(s,a) := r + \gamma * max_{a'} Q(s', a')$

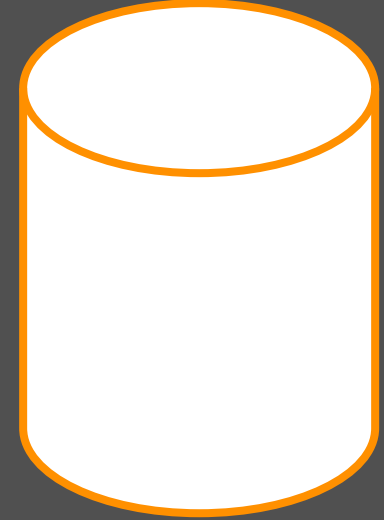- **4 Images form a training sample for the neural net**

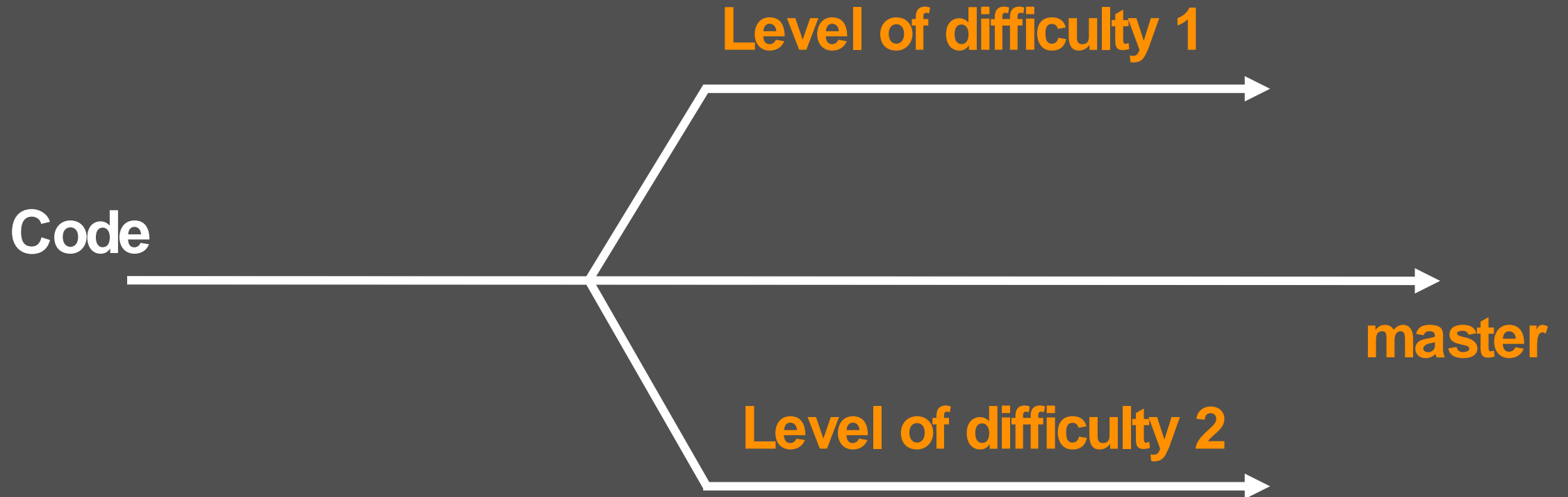# Will the ball go up or down?

# Your Task:

- We give you code with gaps
- You have to **fill in the gaps** with appropriate code
- The **gaps** are marked with a **comment with TODO**
- Gaps can be in all scripts
- Each script has an order of the TODOs, e.g. TODO 1, TODO 2, …

## https://github.com/mati3230/MetaMarathon

- **Code-Repository**
- **How to setup the working environment**
- **Tipps and recommendations**
- **3 branches for different levels of difficulties**
  - easy, medium, hard
- **Solution will be pushed (uploaded) at the end in the master branch ;)**

**Branching:**

Level of difficulty 1

Code

master

Level of difficulty 2

**Branching:**

Level of difficulty 1

Code

**Switching of branches is explained in repository**

master

Level of difficulty 2

# 3 Files:

- **main.py: should be executed**

- **estimator.py: neural net**

- **state_processor.py: image preprocessing**

**main.py**: should be executed

- **Mother-Script**
- **Management processes (training, playing)**
- **Execution of steps in the environment**
- **Optimization of Q-values**

# estimator.py: neural net

- ## Convolutional Neural Net (CNN)

- ## Prediction of Actions

- ## Training of Network

# state_processor.py: image preprocessing

- **Grayscale**
- **Cropping of relevant area**
- **Scaling to 84x84**

**First Steps:**

- **Go to the repository:**
  **https://github.com/mati3230/MetaMarathon**

- **Setup the environment (use manual)**

- **Clone Repository**

- **Choose a level of difficulty and switch branch**

- **Search TODOs and get familiar with code**

# https://github.com/mati3230/MetaMarathon