

The background features a complex network of thin grey lines connecting various points, forming a web-like structure. Scattered throughout are numerous triangles of different sizes and orientations, some solid and some outlined. The overall aesthetic is technical and geometric.

Deep-Q-Learning mit Super Mario Bros. und A3C

Jan Gaida
Angewandtes maschinelles Lernen
Hochschule Hof, Juni 2020

Einführung 01

Environment 02

**Advantage
Actor-Critic 03**

**04 Netzwerk
Architektur**

05 Ergebnisse

06 Resümee

01

Einführung



Super Mario Bros. (1985)

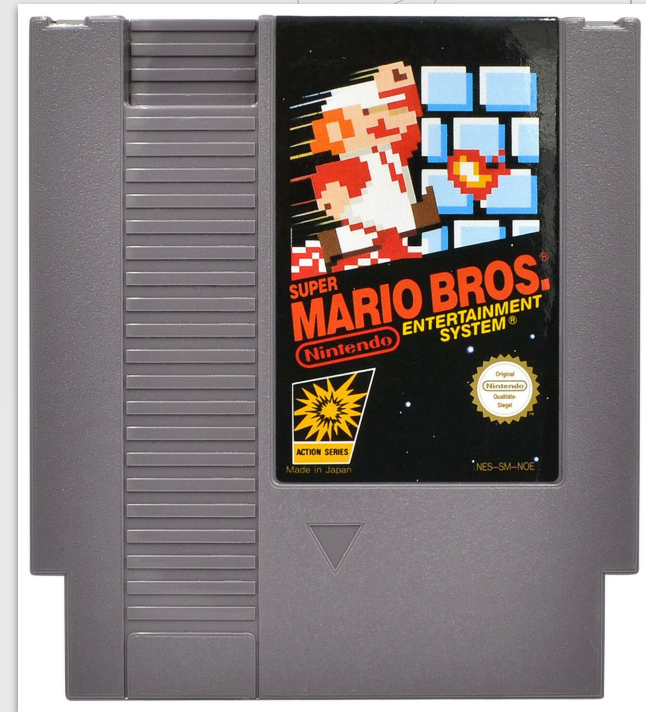
- **Nintendo** Co. Ltd.
- Erstveröffentlichung:
 - **Japan:** 13. September 1985
 - **Westen:** 1986 bis 1987
- Plattformen:
 - Famicom
 - Nintendo Entertainment System (NES)



Quelle: [Datei:雨の日はファミコンで遊べる \(15441664223\).jpg - Wikipedia](#)

Super Mario Bros. (1985)

- **Nintendo** Co. Ltd.
- Erstveröffentlichung:
 - **Japan:** 13. September 1985
 - **Westen:** 1986 bis 1987
- Plattformen:
 - Famicom
 - Nintendo Entertainment System (NES)
- Medium:
 - **40KB** Steckmodul (max. 320 KB)
- Copyright / DMCA:
 - bis (mind.) 2080
 - Ausnahme(n):
'educational setting or documentary purposes' ¹



Quelle: [PicClickImg](#)

¹ Quelle: [Albright-lp](#)

Motivation



Quelle: [KnowYourMeme.com](https://www.knowyourmeme.com/memes/super-mario-bros)

CHALLENGE ACCEPTED





02

Environment

OpenAI Gym



Quelle: [Velotio Technologies](#)

Ziele:

- **Benchmark** von RL-Algorithmen durch eine große Kollektion von diversen Environments
- **Standardisierung** der Environment für bessere Vergleichbarkeit von RL-Algorithmen



Christian Kauten: Gym-Super-Mario-Bros

<https://github.com/Kautenja/gym-super-mario-bros>

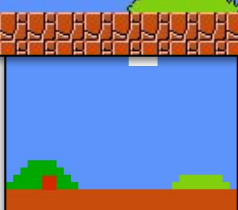
Render-Varianten:

Standard ✓

Downsample

Pixel ✗

Rectangle ✗



Vordefinierte Action-Spaces:



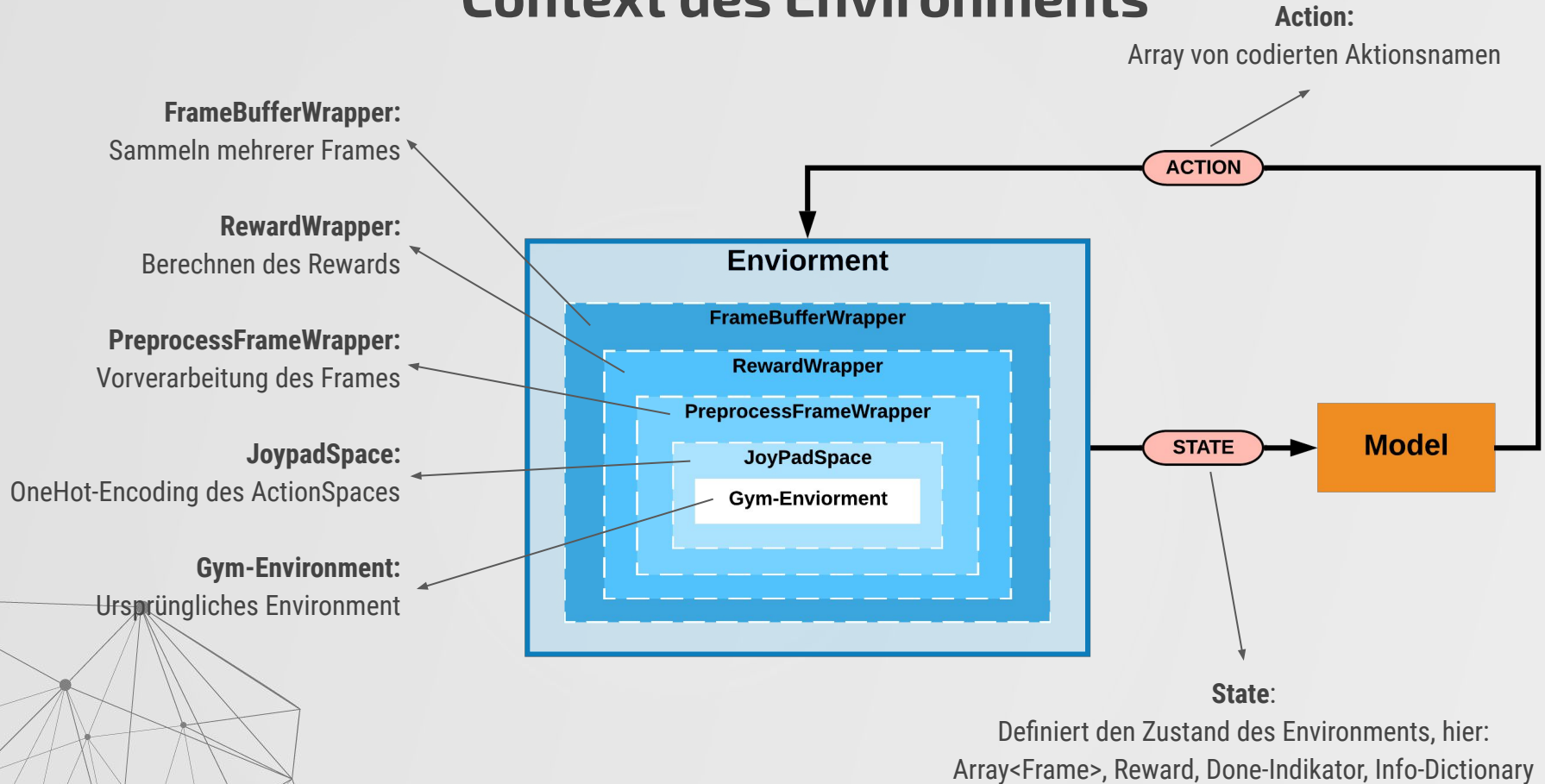
✗ Standard (256-Bit)

✓ Complex

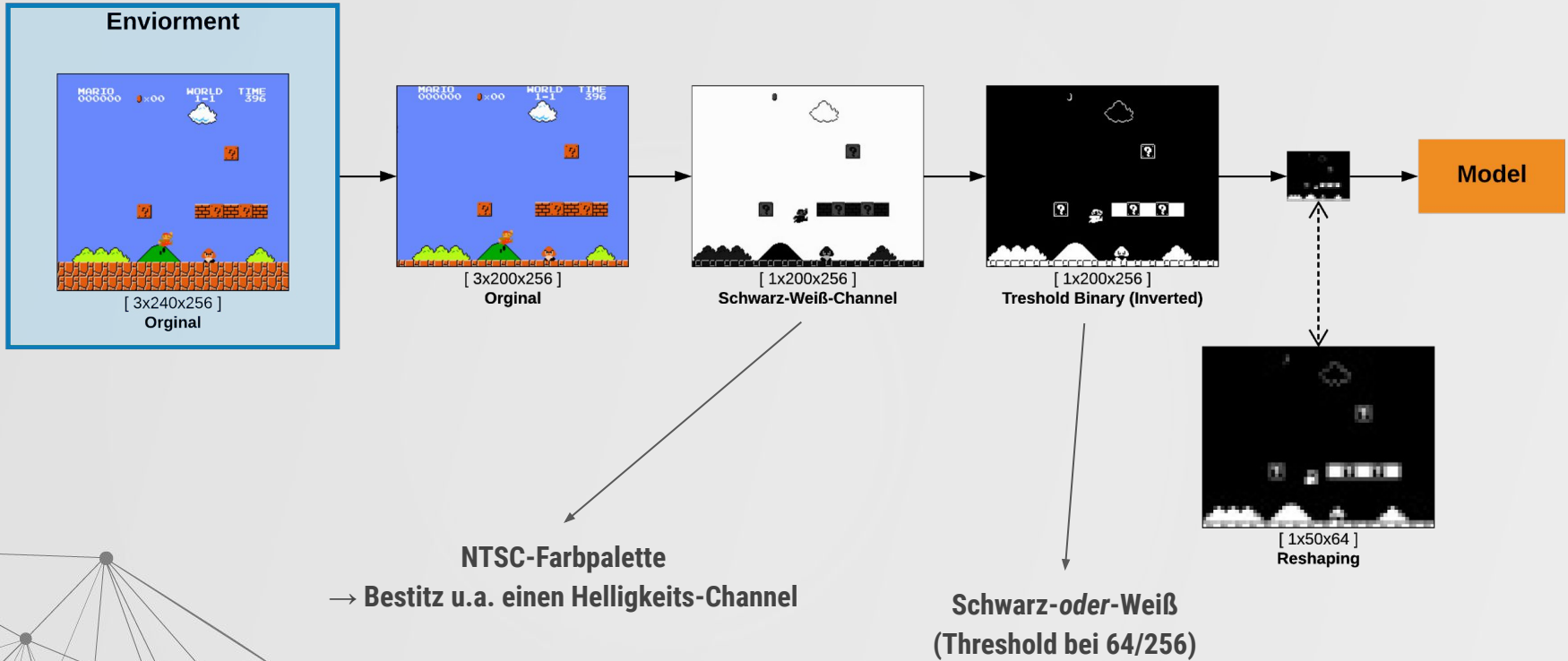
✗ Simple

✗ Right-Only

Context des Environments



Preprocessing



Reward

Reward Shaping ✓

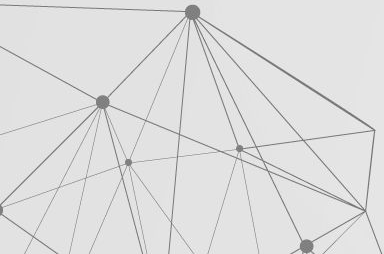
✗ Sparse Reward

- **Keine Generalisierung** des Problems
- **Abhängigkeit** von der Qualität der Reward-Funktion
- (Keine vollständige Exploration)



Aktuelles Forschungsgebiet in RL

(Curiosity-driven Exploration, Unsupervised Auxiliary Tasks, Hindsight Experience Replay)



Reward

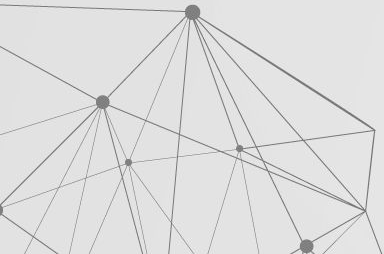
→ X-Position ($w = [\text{len}(\text{buffer}) * -1; \text{len}(\text{buffer})]$)
 $\text{delta_x} = x_1 - x_0$

→ Zeit ($w = [0; (\text{len}(\text{buffer})/10)]$)
 $\text{delta_time} = \min(t_1 - t_0, 0)$

→ Erreichtes Ziel ($w = [0; 45]$)
 $r_{\text{goal}} = 45$ if goal_achieved else 0

→ Verlorenes Leben ($w = [-45; 0]$)
 $r_{\text{life}} = -45$ if life_lost else 0

→ **$\text{reward} = (\text{delta_x} + \text{delta_time} + r_{\text{goal}} + r_{\text{life}}) / 10$**

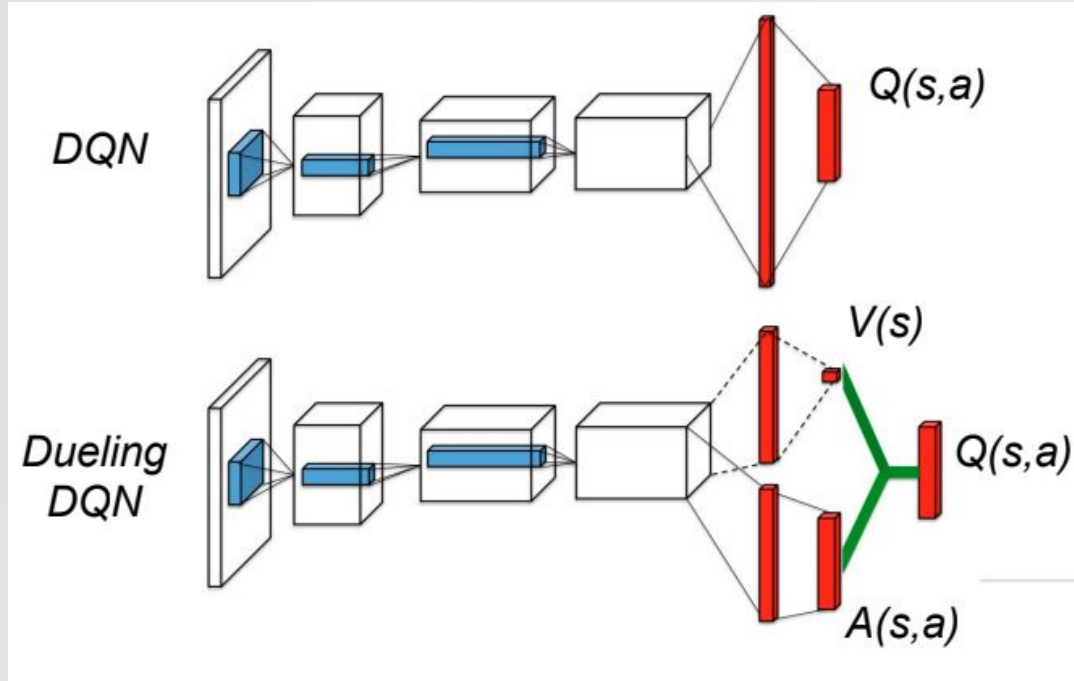




03

Advantage Actor-Critic

Warum Actor-Critic ?



→ Erzeugen einer Q-Table von State's relativ zu der bestmöglichen Action

→ Erzeugen einer Q-Table von State's relativ zu der bestmöglichen Action + Bewertung von aktuellen State

Warum Actor-Critic ?

DQN



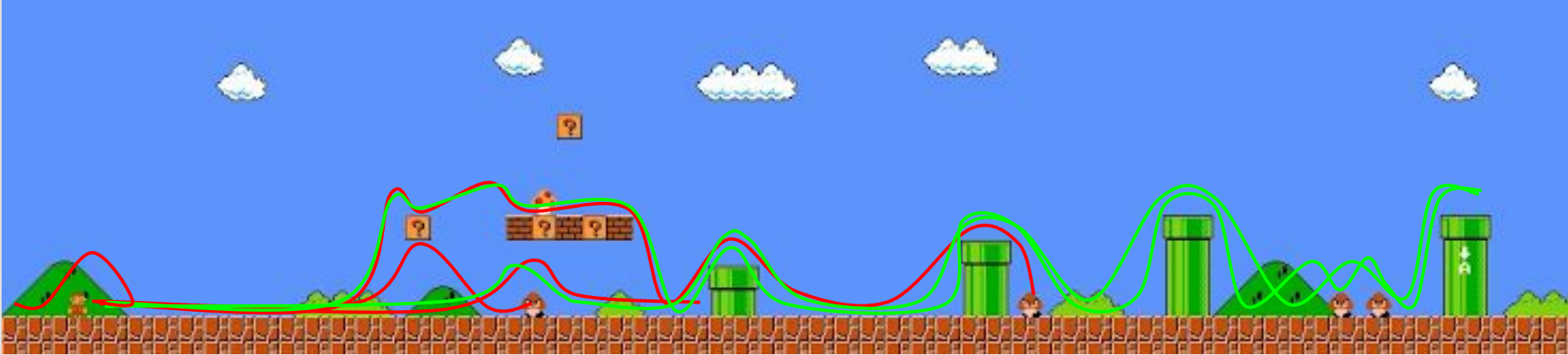
Dueling
DQN



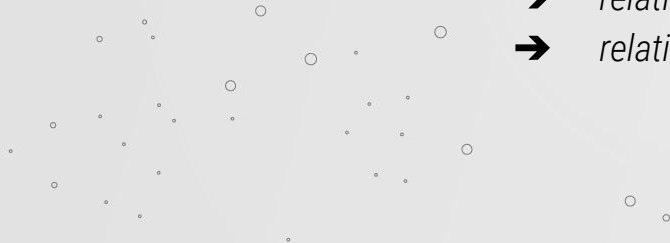
ble von State's
öglichsten Action

ble von State's
öglichsten Action
uellen State

Warum Actor-Critic ?



- *relativ* großer **ObservationSpace** (Output d. Env.)
- *relativ* großer **ActionSpace** (Input d. Env.)



Idee hinter A3C



```
graph TD; A[Idee hinter A3C] --> B[Asynchronous]; A --> C[Advantage]; A --> D[Actor-Critic];
```

Asynchronous

Konkurrente Trainings-Prozesse

Advantage

.. unter Berücksichtigung einer geschätzten Bewertung des aktuellen States

Actor-Critic

.. ausgeführt durch einen Actor in Abhängigkeit eines Kritikers

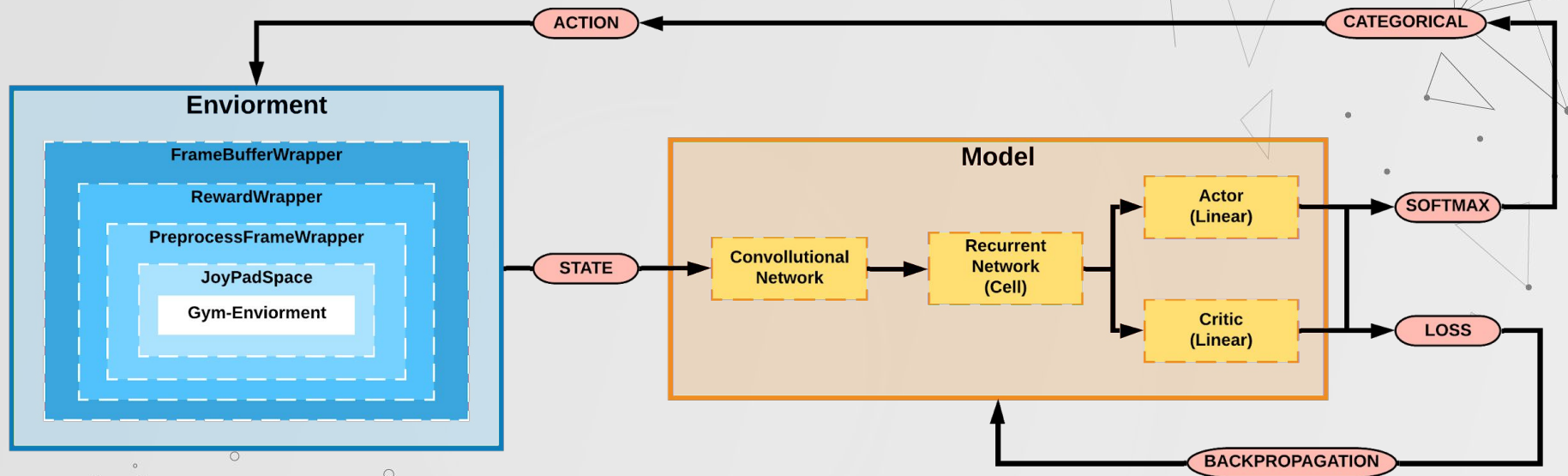
**Keine
Q-Tables !**

Idee hinter A3C

Beispiel:

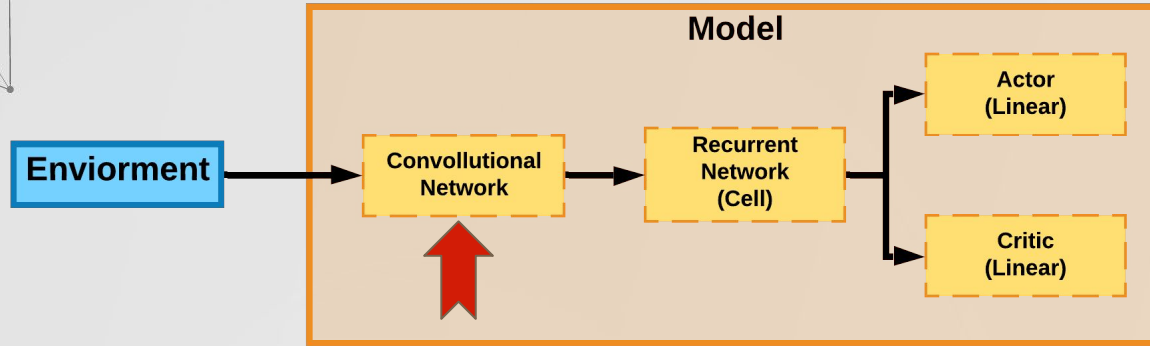
[...] Let [us] imagine a small mischievous **child (actor)** [which] is discovering the amazing world around him, while his **dad (critic)** oversees him, to make sure that he does not do anything dangerous. Whenever the kid does anything good, his dad will praise and encourage him to repeat that action in the future. And of course, when the kid does anything harmful, he will get [a] warning from his dad. [...]

Context A3C



04

Netzwerk Architektur



Approach B

Naiv

Approach B

Deep-Convolutional

- Auswertung der Bilderfolge
- Approach B ist inspiriert von Google_ResNet-Modulen

Approach A

Naiv

Enviornment

Conv2d

Kernelsize = 3x3
Output_Channels = 320
Stride = 2

Conv2d

Kernelsize = 3x3
Output_Channels = 240
Stride = 2

Conv2d

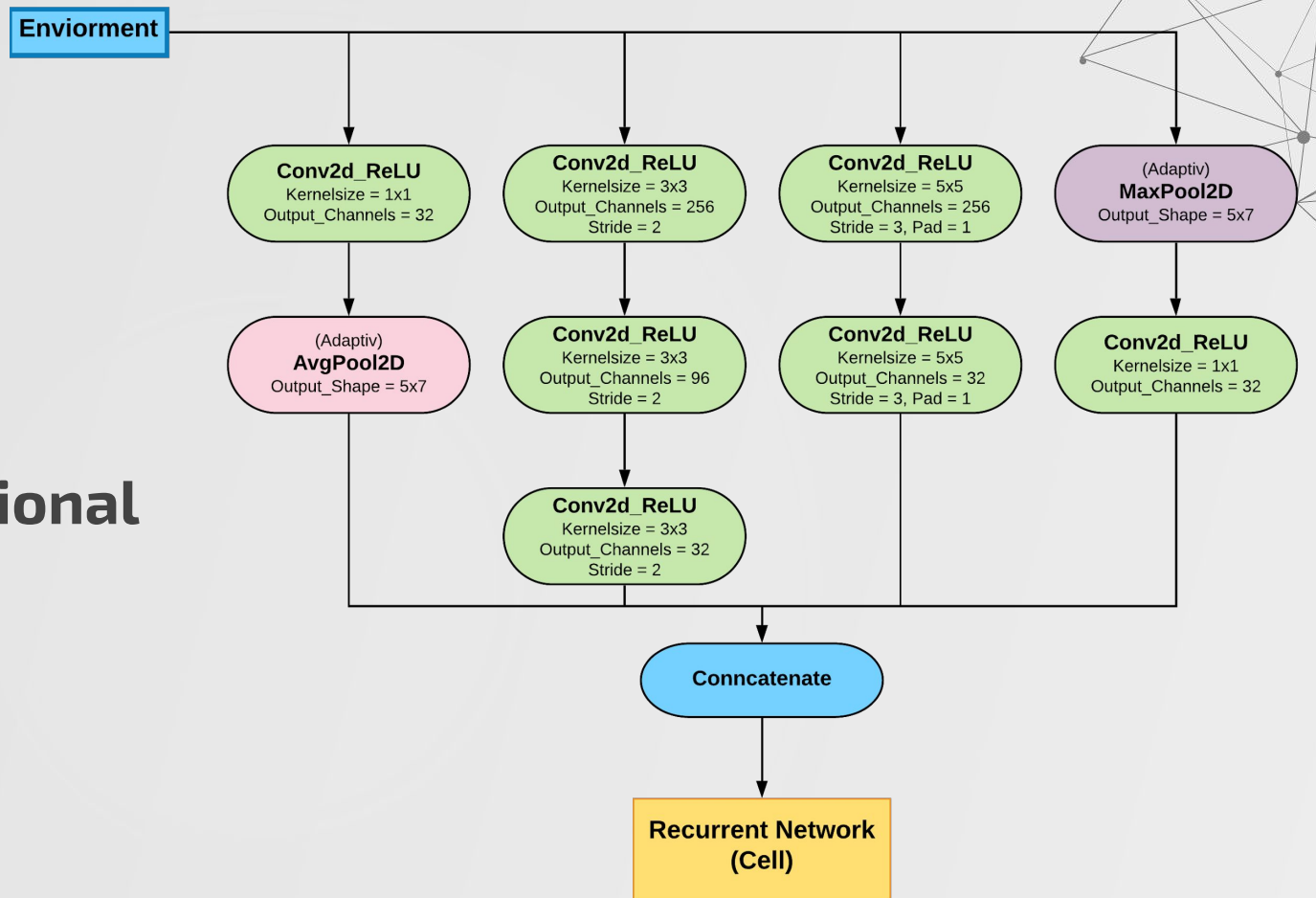
Kernelsize = 3x3
Output_Channels = 160
Stride = 2

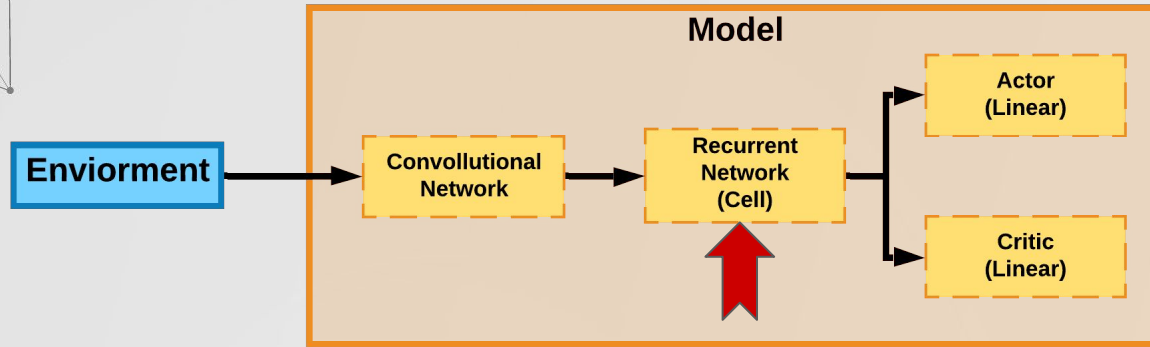
Conv2d

Kernelsize = 3x3
Output_Channels = 80

**Recurrent Network
(Cell)**

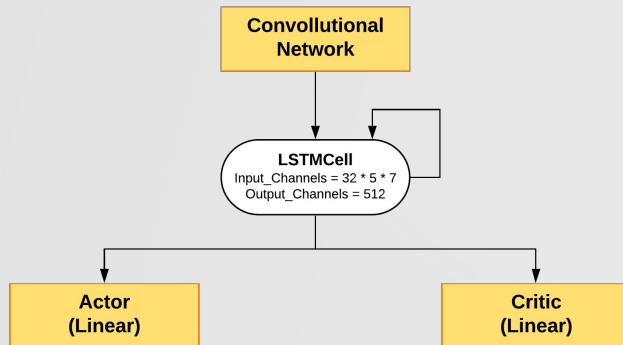
Approach B Deep-Convolutional





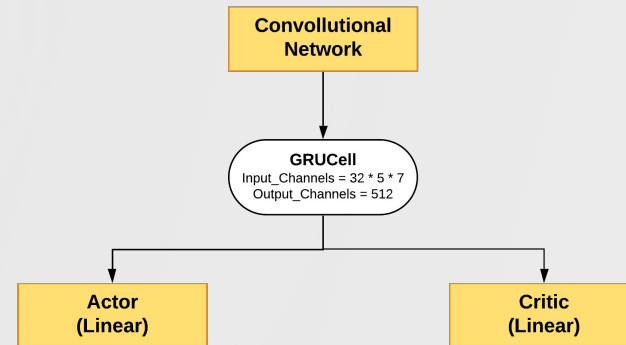
Approach A

LSTM-Zelle



Approach B

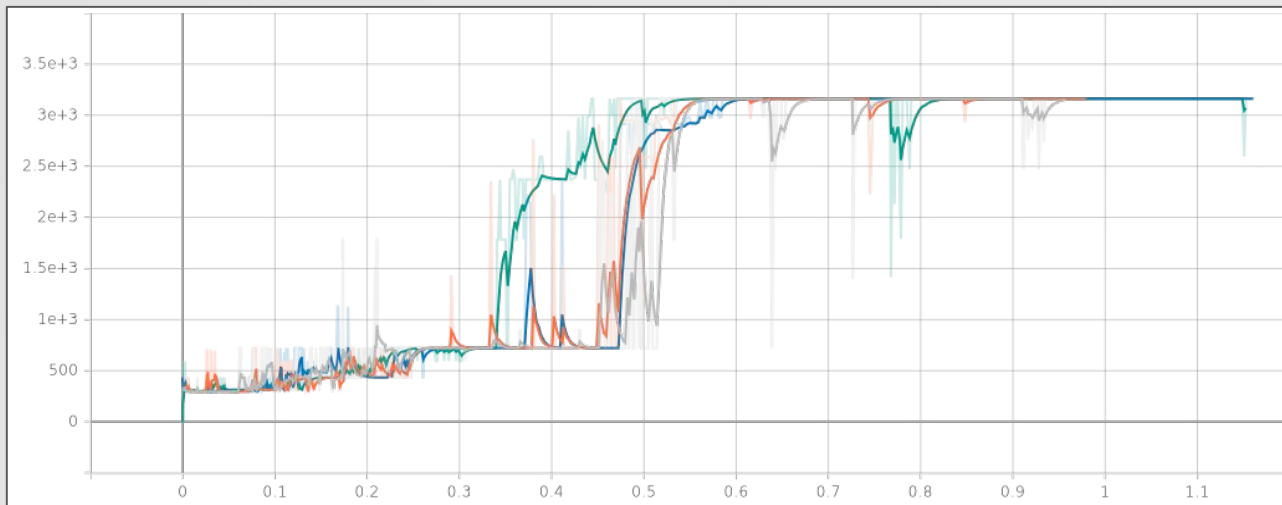
GRU-Zelle



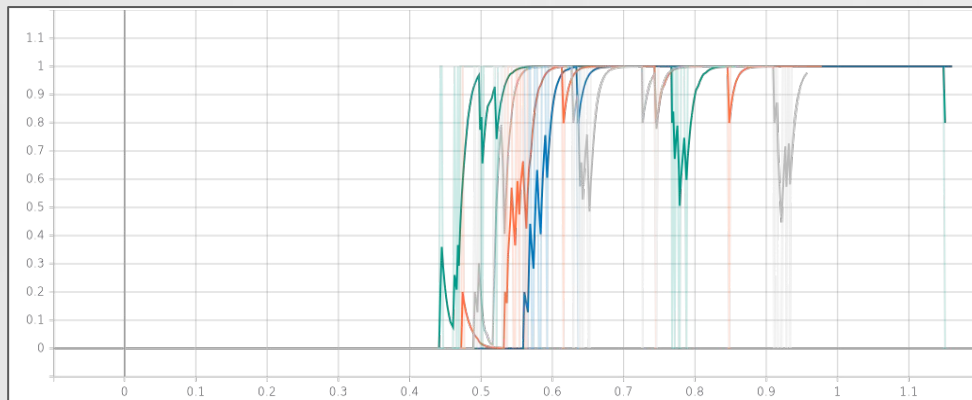
05

Ergebnisse

Erreichte X-Position



Flagge



- ✓ DCN (GRU)
- ✓ CN (GRU)
- ✓ CN (LSTM)
- ✓ DCN (LSTM)

GRU: Convolutional vs DeepConvolutional

CN



1500 Ep x5



2000 Ep x5



2500 Ep x5



3500 Ep x5

DCN



LSTM: Convolutional vs DeepConvolutional

CN



1500 Ep x5



2000 Ep x5



2500 Ep x5



3500 Ep x5

DCN



DCN - LSTM

Stage 1 World 3
2500 Ep x 5 Threads
Smoothingfaktor: 0.8



500 Ep x5



1000 Ep x5

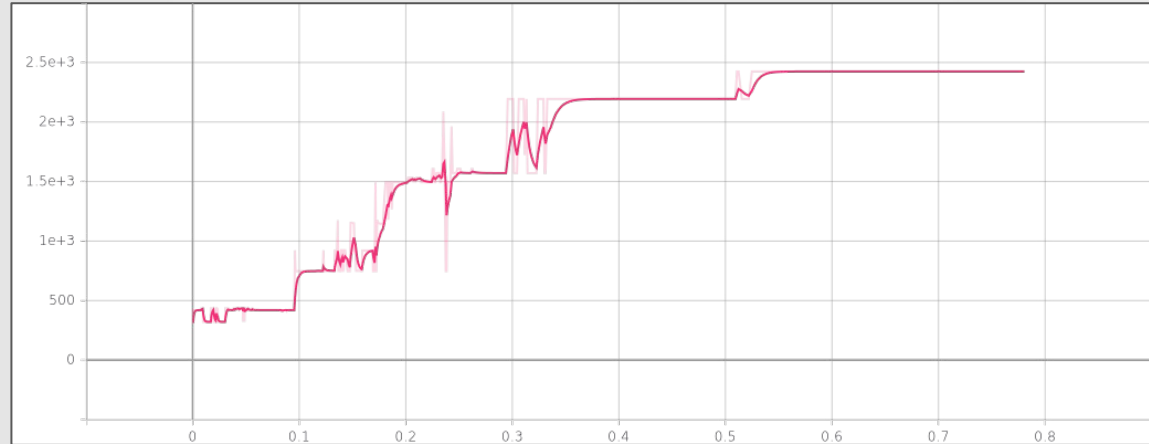


2000 Ep x5



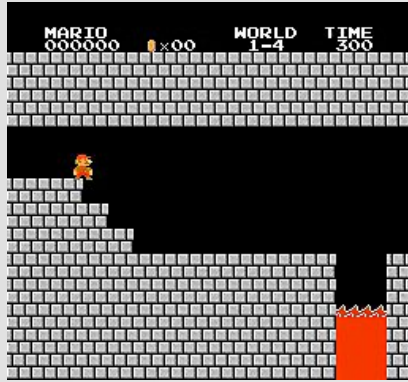
2500 Ep x5

X-Position

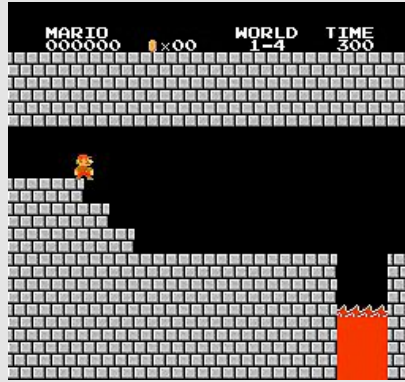


DCN - LSTM

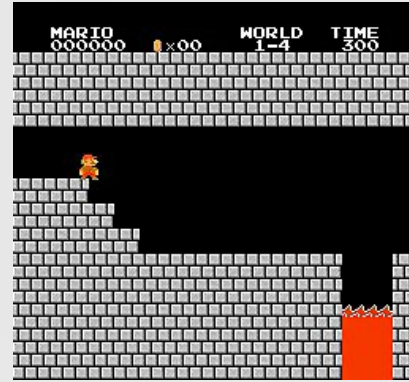
Stage 1 World 3
7500 Ep x 5 Threads
Smoothingfaktor: 0.8



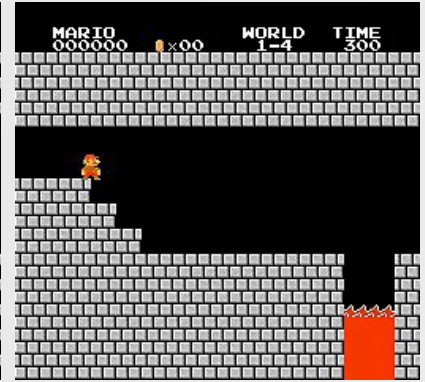
1000 Ep x5



2000 Ep x5

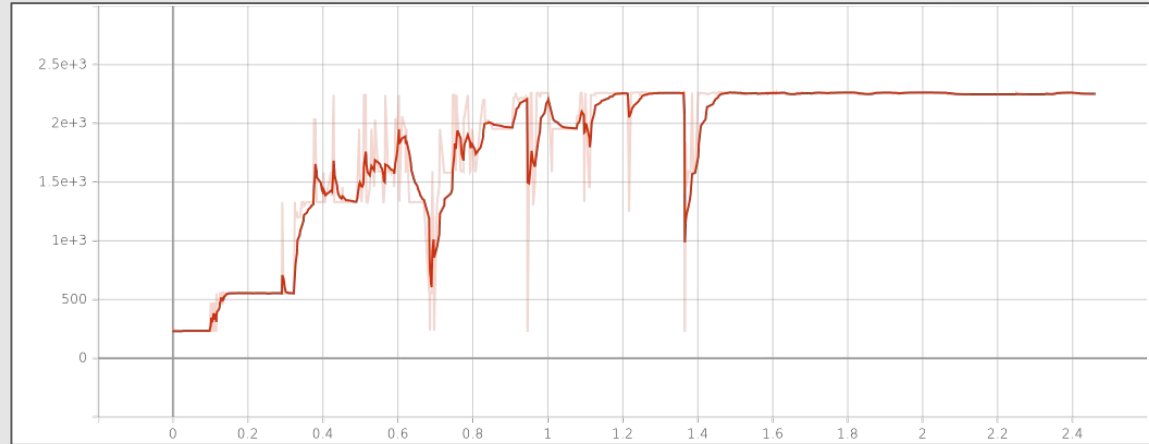


3000 Ep x5



7500 Ep x5

X-Position





06

Resümee

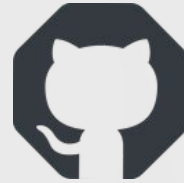
- Super Mario ist für einige RL-Algorithmen eine große Herausforderung
 - 'The Mario AI Competition (2009-2012)'
 - Marl/O, Intrinsic Curiosity Module, ...
- 'Do not trust a (learning) robot using GRU'
 - Stabile Lernerfolge sind für manche Probleme Key
 - LSTM > GRU
- 'We must go deeper'
 - Deep-CNN > CNN



DANKE FÜR IHRE AUFMERKSAMKEIT



[Ask me Anything](#)



[Repository](#)

CREDITS: This presentation template was created by **Slidesgo**, including icons by **Flaticon**, and infographics & images by **Freepik**.