

TheFatBot

...

Ein Doom-Reinforcement-Learning-Approach mit ViZDOOM



Agenda

1. Einleitung & ViZDoom
2. Szenario erstellen
3. Duel Q-Learning & PPO
4. Unsere Umsetzung & Hyperparameter Tuning
5. Herausforderungen
6. Fazit



To Do

- Sepp

- Policy versus non-Policy Check
- PPO Grundlagen Check
- SLADE, Szenarios, Scripting Klemens fragen ob er noch Ideen hat was fehlt
- Visualisierung Muss noch
- Herausforderungen Check
- Einleitung drüberlesen/anpassen Check
- Fazit Muss noch

- Björn

- PPO zum laufen bringen NE
- PPO Kapitel schreiben SEPP
- DQN zu Ende DONE
- Code durch kommentieren DONE

- Klemens

- Slade & Szenario Folien
- Hyperparameterfolien mit Videobeispielen
- Letztes Model trainieren Done
- Videos aufnehmen und hochladen Done
- Slade und Szenario Kapitel Done?
- Ergebnisse und Herausforderungen Kapitel (Reward challenges, Scripting challenges (Alles mit Slade), Learning challenges auch bzgl mywayhome Done



Einleitung

Was ist Doom?





ViZDOOM

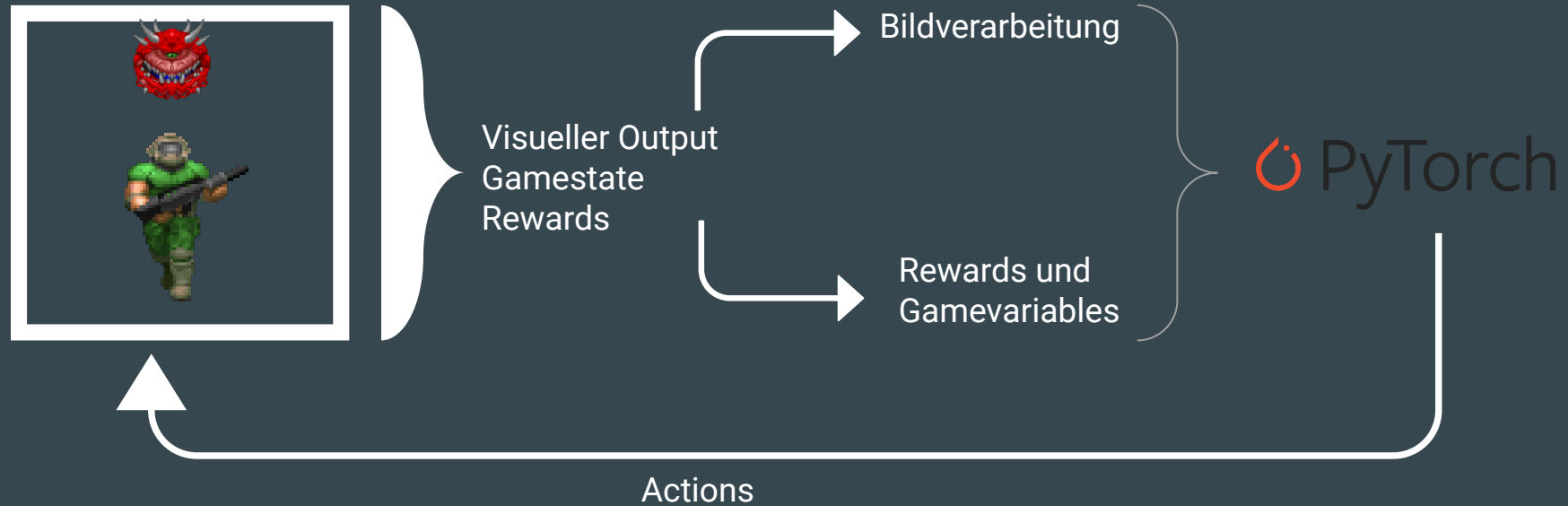
Erlaubt uns:

- Zugriff auf Screen Buffer
- Zugriff auf In-Game Variablen
- 4 Control Modes
- Erstellen von eigenen Szenarien
 - Maps
 - dynamisches Spielgeschehen
 - speziell für Reinforcement Learning: Rewards definieren





Überblick





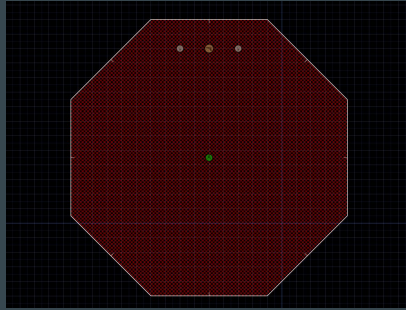
Szenarios



Wad Files



Scripts



Maps



Game Files &
Textures

Slade



SLADE

File Edit Archive Entry Text View Tools Help

Base Resource: <none> Palette: Existing/Global

Archive Manager

Start Page BasicAugment_oneEnemy.wad

Archives Bookmarks File Browser

Open Archives:

| Filename | Path |
|---------------------------|------|
| BasicAugment_oneEnemy.wad | C:\U |

Recent Files:

| Filename |
|------------------------------|
| BasicAugment_oneEnemy.wad |
| BasicAugment.wad |
| BasicAugment_deathreward.wad |

Entries

| Name | Size | Type |
|-----------|--------|----------------------|
| MAP01 | 0 | Map Marker |
| TEXTMAP | 1.87kb | UDMF Map Data |
| BEHAVI... | 484 | Compiled ACS (ZDoom) |
| SCRIPTS | 1.19kb | ACS source |
| DIALOG... | 2 | Unknown |
| ENDMAP | 0 | Marker |

Show: All Filter:

Entry Contents

Save Revert Text Language: ACS (ZDoom) Jump To:

```

1 #include "zcommon.acs"
2
3 int target_id = 10;
4 int target_amount = 2;
5
6 global int 0:reward;
7
8
9 script 1 OPEN
10 {
11
12     reward = 0;
13     Thing_ChangeTID(0, 1000 + PlayerNumber()); // This assigns the TID
14 }
15
16 int c = 0;
17 script 2 ENTER
18 {
19     TakeInventory("Fist", 1);
20     ACS_Execute(3, 1, 0, 0, 0);
21 }
22
23 script 3 (void)
24 {
25     int bullets = CheckInventory("Clip");
26     int health = GetActorProperty(1000, APROP_Health);
27     while(true)
28     {
29         int t_bullets = CheckInventory("Clip");
30         int t_health = GetActorProperty(1000, APROP_Health);
31         if(t_bullets < bullets)
32         {
33             reward = reward - 5.0;
34         }
35         if (t_health < health)
36         {
37             reward = reward - 50.0;
38         }
39         bullets = t_bullets;
40         health = t_health;
41
42         delay(1);
43     }
44 }
45
46 script 4 (void)
47 {
48     reward = reward + 100.0;
49 }
  
```

3: SCRIPTS, 1221 bytes, ACS source

Ln 1, Col 1, Pos 0



Implementierung von Rewards

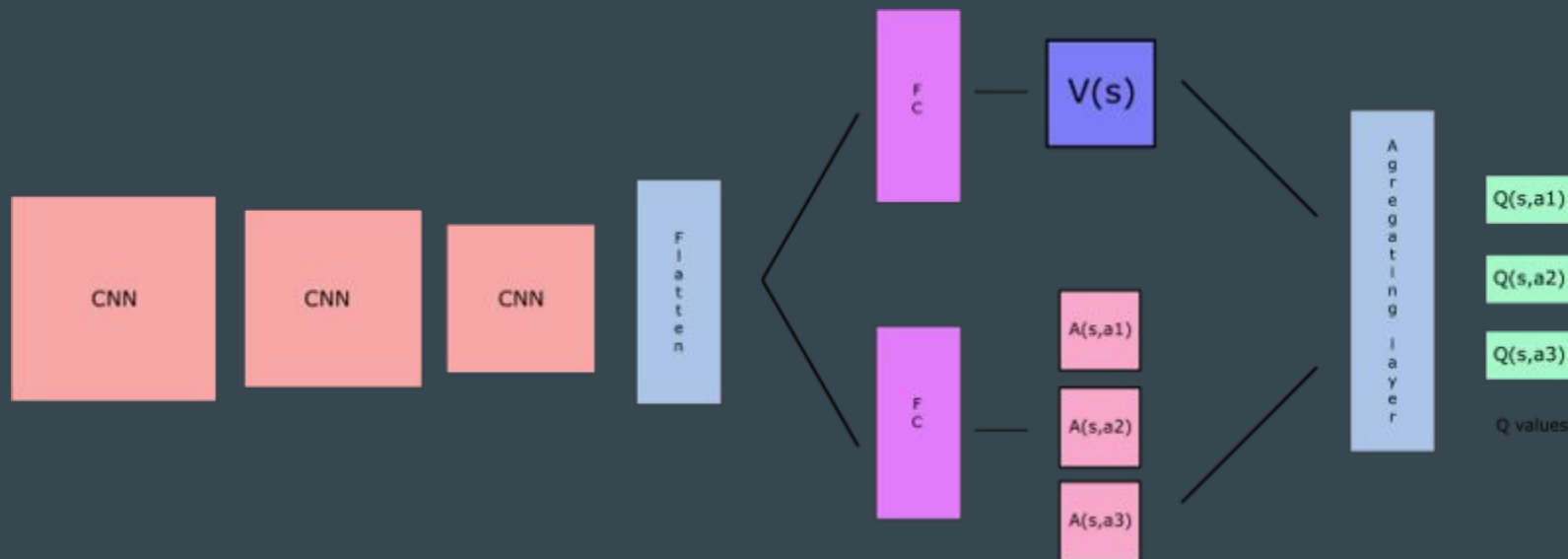




Duel Q-Learning



Wie funktioniert Duel Q-Learning?



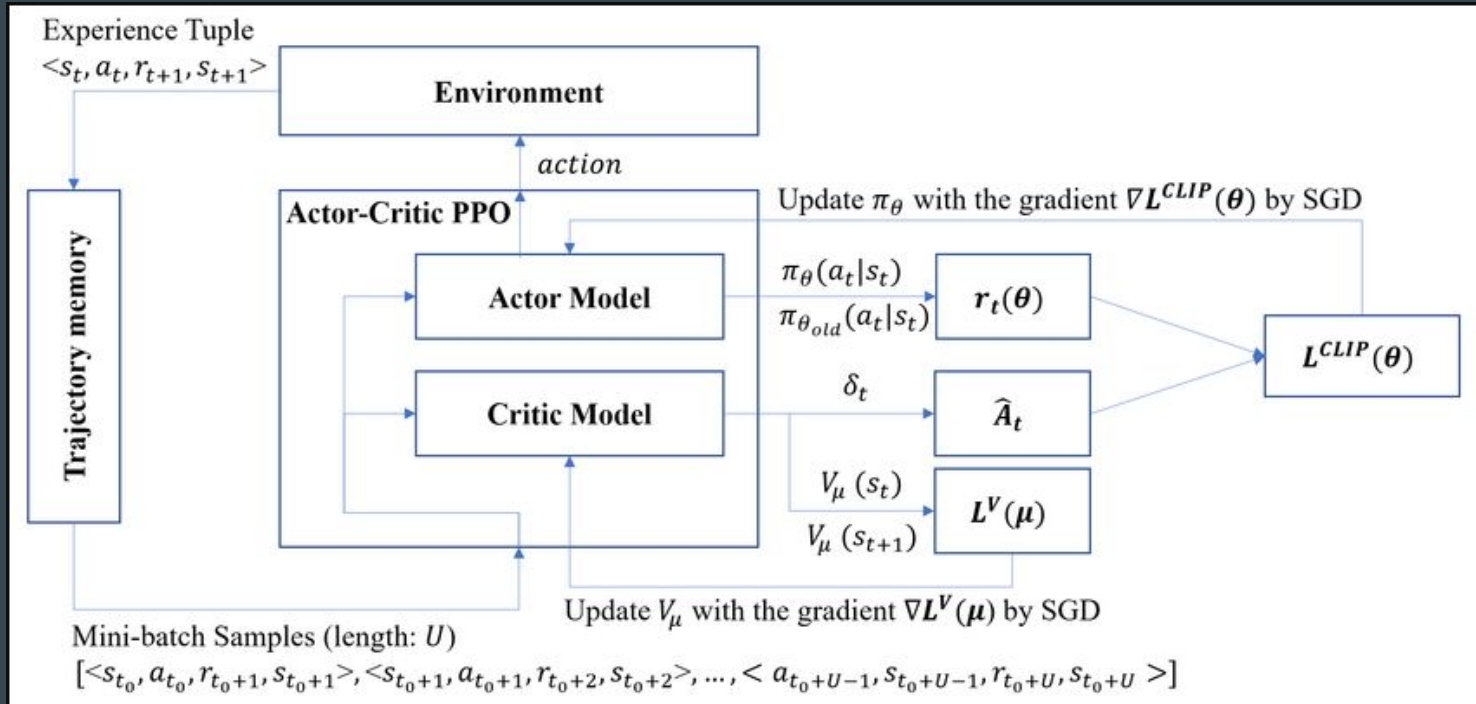
<https://www.freecodecamp.org/news/improvements-in-deep-q-learning-dueling-double-dqn-prioritized-experience-replay-and-fixed-58b130cc5682/>



Proximal Policy Optimization



Wie funktioniert PPO?





Unsere Umsetzung



Die Ausgangslage

vmware®





Erste Schritte in ViZDoom



Basic.wad



My Way Home.wad



Challenges from the get go



Sparse Rewards



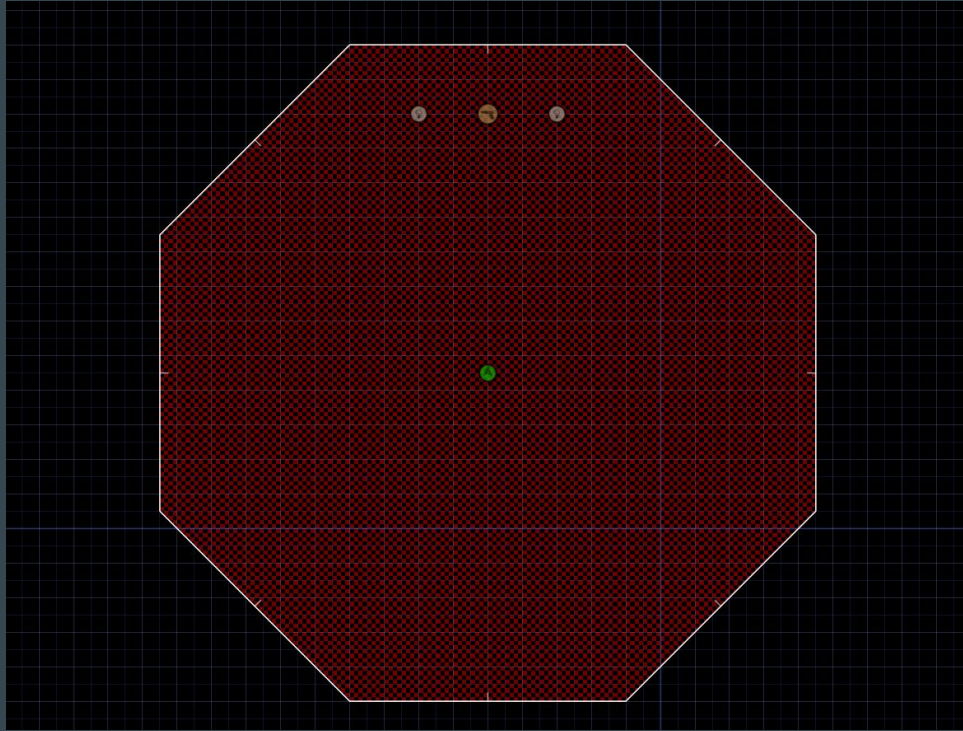
Balance Rewards



Sensitivity to Parameters



Unsere Map





Unser Reward script

Positive Rewards:

- Gegner Treffer
- Aufsammeln der Waffe

Negative Rewards

- Schießen
- Schaden nehmen
- Game Over

Experimentelle Rewards:

- Negativer Living Reward
- Rewards auf Positionsbasis



Hyper Parameter Tuning und Vergleich

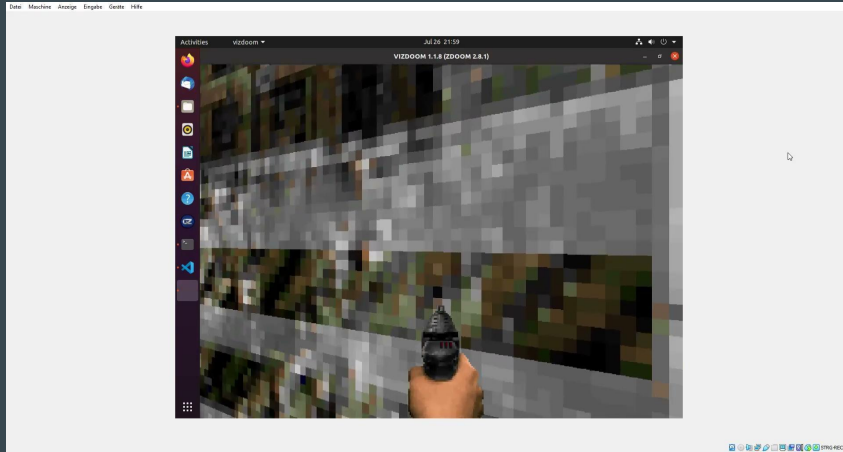


Hyper Parameter

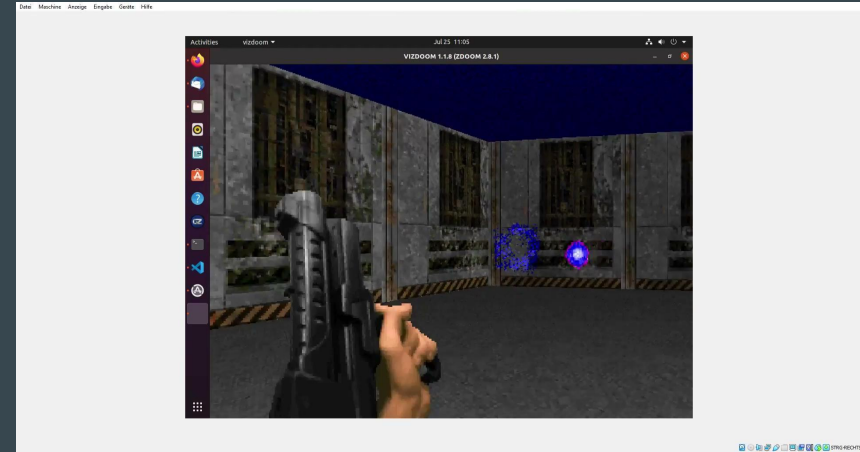
1. Anzahl der Epochen
2. Learning Rate
3. Discount Factor
4. Learning Steps per Epoch
5. Replay Memory Size
6. Batch Size



Hyperparameter Tuning

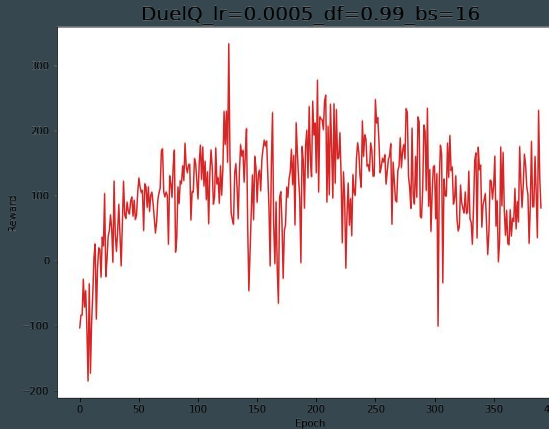
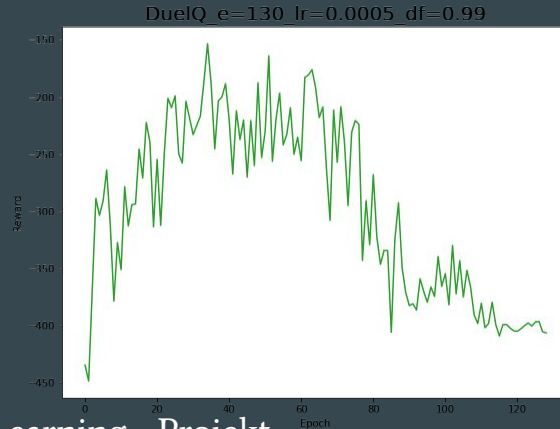
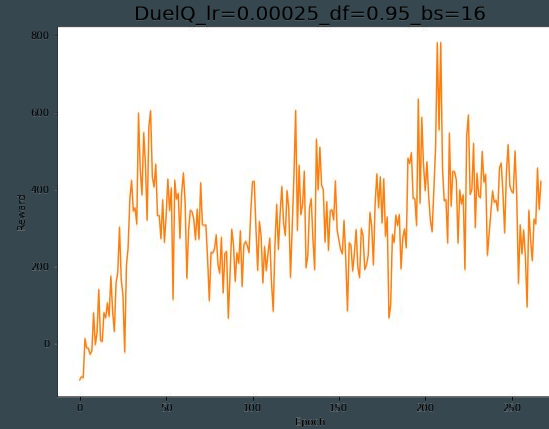
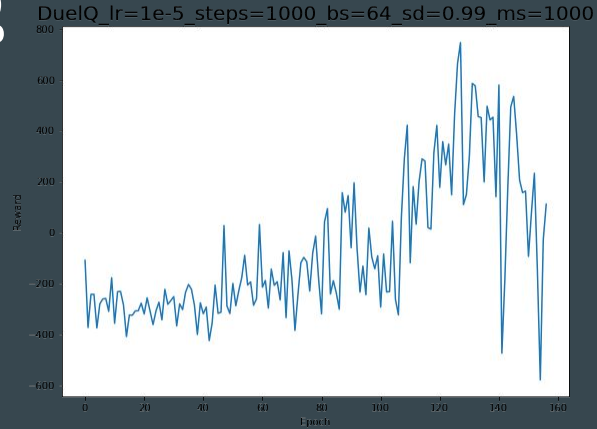


$e=391$ $lr=0,0005$



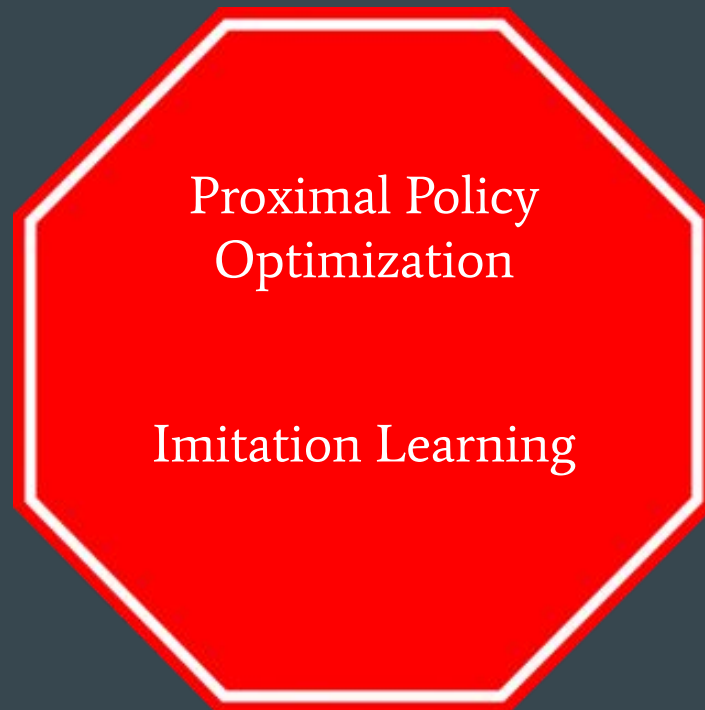
$e=412$ $lr=0,0002$

Bewertung





Gescheiterte Umsetzungen





Fazit



Fazit



Eher schwache Performance
der Modelle

- Zu wenig Rechenleistung
 - Begrenzung durch VMs
- Probleme bei Reward Funktion
- Zeitlich zu aufwändig
- Umsetzung von weiteren RL Ansätzen an Bugs gescheitert



Danke für eure Aufmerksamkeit!



Quelle

- Federated Reinforcement Learning for Training Control Policies on Multiple IoT Devices - Scientific Figure on ResearchGate. Available from: https://www.researchgate.net/figure/The-actor-critic-proximal-policy-optimization-Actor-Critic-PPO-algorithm-process_fig3_339651408 [accessed 27 Jul, 2021]
- Reinforcement Trained Basic Example: https://www.youtube.com/watch?v=fKHw3wmT_uA
- Reinforcement Trained My Way Home Example: <https://www.youtube.com/watch?v=15yZubaTLvw>
- Vizdoom Proposal: <https://arxiv.org/abs/1605.02097> [accessed 26 Jul, 2021]
- Proximal Policy Optimization Paper: <https://arxiv.org/abs/1707.06347> [accessed 26 Jul, 2021]
- Duel Q Learning Proposal: <https://arxiv.org/abs/1511.06581> [accessed 26 Jul, 2021]