

Metody Actor Critic pro spojitá prostředí

Martin Kostelník, Michal Glos, Michal Szymik

Přiložené materiály: [Github projektu](#)

1. Úvod

Cílem tohoto projektu bylo vytvořit vlastní implementaci algoritmů posilovaného učení (RL), především Actor-Critic pro spojitá prostředí. Jsou využity prostředí knihoven OpenAI Gym¹ a dm_control [3]. S výslednou implementací je prováděno několik sad experimentů, které demonstrují úspěšnost implementace.

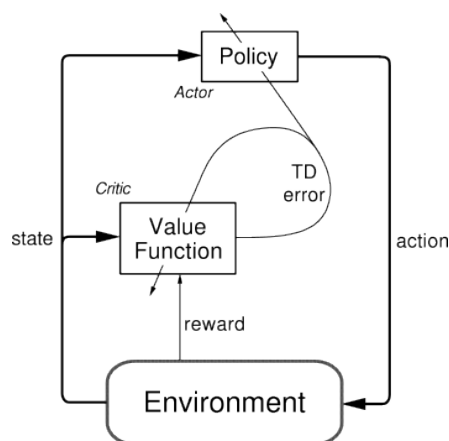
Prováděné experimenty srovnávají úspěšnost naší implementace s algoritmy z knihovny stable-baselines-3² která obsahuje řadu různých implementací algoritmů posilovaného učení. Během experimentů se zkoumá vliv jednotlivých parametrů na celkovou úspěšnost agenta, jako je architektura neuronových sítí, nastavení aktivačních funkcí a obecně různé nastavení parametrů specifických pro typ implementované metody.

V tomto projektu se zabýváme implementací TD3 algoritmu (viz 2), který vychází z jiného Actor-Critic algoritmu Deep Deterministic Policy Gradient Agents (DDPG), do kterého přináší několik změn aby nedocházelo k přeceňování value funkce.

2. Actor-critic a algoritmus TD3

Metody actor-critic jsou specifické tím, že mají oddělenou strukturu pro policy a value funkci. Policy síť se nazývá actor a síť aproximující value funkci se nazývá critic. Úkolem sítě actor je vybírat akce, které mu poté critic kritizuje, čili aproximuje očekávanou value funkci. Průběh učení agenta je dobře vidět na diagramu v obrázku 1.

Twin-delayed deep deterministic policy gradient (TD3) je jedním z actor-critic algoritmů. Jedná se o



Obrázek 1. Diagram znázorňující průběh učení agenta používajícího actor-critic metodu.

online a off-policy učící metodu, založenou na principu actor-critic. TD3 agent hledá takovou optimální policy, která maximalizuje očekávanou kumulovanou odměnu. Agent používá dvě critic sítě, kterými se učí dvě value funkce a pro aktualizaci policy používá minimální odhad value funkce. Policy se updatuje méně často než Q-funkce a během její aktualizace se přidává šum k vybrané akci, tudíž se na základě policy budou s menší pravděpodobností vybírat akce s vysokým odhadem value. [1]

3. Popis řešení

Projekt je implementován v jazyce Python3 s použitím knihoven Gym, MuJoCo [2], stable-baselines3 a několika dalších podpůrných nástrojů.

3.1 Implementace TD3

Agentu lze spustit ve dvou módech: trénovacím a evaluačním. Při evaluačním módu se načte již existující a natrénovaný model a provede se několik běhů v daném

¹<https://www.gymnasium.ml/>

²<https://stable-baselines3.readthedocs.io/en/master/index.html>

51 prostředí. Z běhů je vypočtena průměrná odměna,
52 která je vytištěna na standardní výstup. Součástí eval-
53 uace je taky vykreslování chování agenta do okna a
54 vytvoření videa.

55 Zajímavější je trénování. Po nastavení všech
56 parametrů uživatelem je vytvořeno celkem šest neu-
57 ronových sítí. Dvě pro síť actor (actor a actor-target)
58 a dvě pro každou ze dvou sítí critic (critic1,2 a
59 critic1,2-target). Také se provede naplnění paměti
60 agenta náhodnými akcemi v náhodných stavech
61 prostředí s příslušnými odměnami. V tuto chvíli
62 můžeme začít samotné trénování. To trvá daný počet
63 kroků a v každém z nich provedeme akci vybranou
64 sítí actor, ke které je přidán šum.

65 Také provedeme optimalizační akci, ve které se
66 náhodně vybere několik vzpomínek z paměti. Vy-
67 počítáme hodnotu value funkce pro obě critic-target
68 sítě vybereme z nich menší hodnotu a s pomocí chy-
69 bové funkce MSE spočítáme loss sítí critic.

70 Optimalizace sítě actor se provádí pouze co
71 N kroků. Zde nám stačí pro dané stavy z paměti
72 vygenerovat akce a jako chybu vzít negativní průměr
73 hodnot daných sítí critic. Na konci trénovacího kroku
74 se provede aktualizace target sítí tak, že se provede
75 lineární kombinace parametrů s non-target sítí s
76 využitím parametru $\tau \in (0, 1)$.

77 4. Experimenty

78 Pro vyhodnocení výsledného implementovaného al-
79 goritmu jsme se rozhodli provést řadu experimentů.
80 V první sadě experimentů jsme prováděli porovnání
81 našeho TD3 algoritmu s existující implementací z kni-
82 hovny stable-baselines-3³. Využili jsme méně kom-
83 plexních prostředí fyzikálního simulátoru MuJoCo⁴.
84 Druhá sada experimentů se zabývá vlivem nastavení
85 hyperparametrů na agentovo učení v prostředí Hopper.

86 4.1 Srovnání se stable-baselines-3

87 Pro porovnání našeho algoritmu s existující implemen-
88 tací z knihovny stable-baselines-3 jsme si vybrali tři
89 prostředí z MuJoCo:
90 InvertedDoublePendulum-v2, Reacher-v2
91 a HalfCheetah-v2. Pro vyhodnocení jsme použili
92 kritérium průměrné odměny z posledních 256
93 epizod. Pro každé prostředí byly nastaveny stejné
94 hyperparametry u obou implementací. Výsledky učení
95 pro 50000 časových kroků v každém prostředí lze
96 vidět na obrázku 2.

³<https://stable-baselines3.readthedocs.io/en/master/modules/td3.html>

⁴<https://mujoco.org/>

Nejlepší natrénované modely obou algoritmů byly 97
evaluovány na 1000 epizodách. Výsledky evaluace 98
jsou v tabulce 1, která ukazuje průměrnou hodnotu 99
odměň přes všechny provedené epizody. 100

	Inverted Double Pendulum	Reacher	HalfCheetah
Our-TD3	9349.20	-6.08	2856.33
SB3-TD3	4140.74	-4.84	2493.56

Tabulka 1. Průměrné hodnoty odměn při evaluaci
algoritmů na 1000 epizodách.

Výsledky tohoto experimentu naplňují očekávání. 101
Náš algoritmus se učí řešit prostředí s velice podob- 102
nými výsledky jako existující stable-baseline-3 103
implementace. Pro úplné porovnání by bylo třeba 104
provést více experimentů s různým nastavením 105
hyperparametrů, výše uvedené porovnání nám ale 106
přišlo pro účely tohoto projektu dostatečné. 107

4.2 Nastavení hyperparametrů 108

Druhá sada experimentů se snaží zjistit závislost 109
úspěšnosti agenta na nastavení hyperparametrů. 110
Všechny experimenty budeme provádět v prostředí 111
Hopper-v2 a budeme sledovat průměr odměn z 112
posledních 256 kroků a chybové funkce obou sítí. 113

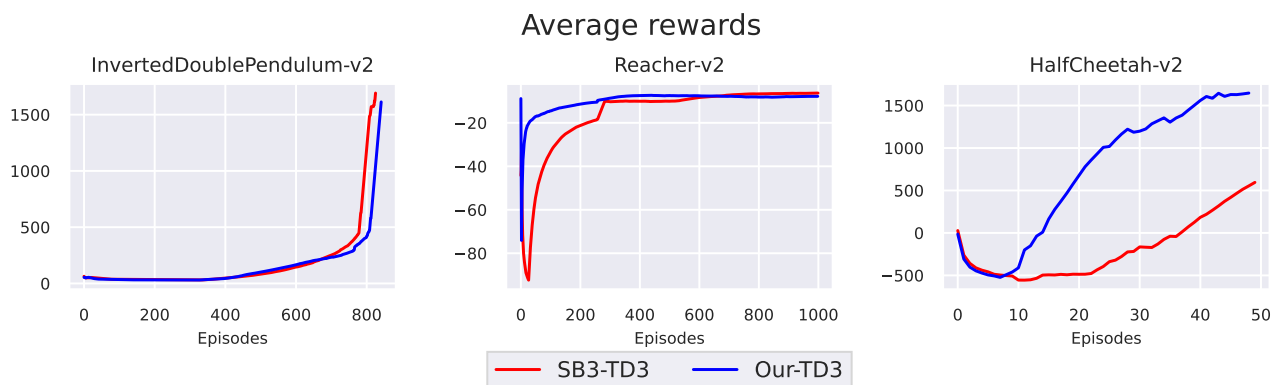
Budeme sledovat efekt discount factoru γ , kombi- 114
načního parametru τ , různé architektury sítí a různá 115
nastavení frekvence učení sítí actor. Tyto všechny 116
výsledky budeme porovnávat s výchozím nastavením 117
parametrů, které je následující: 118

- $\gamma = 0.99$ 119
- $\tau = 0.005$ 120
- $net = \{400, 300\}$ 121
- $delay = 2$ 122

Budeme vždy trénovat na 200000 kroků a budeme 123
testovat následující možnosti parametrů: 124

- $\gamma \in \{0.0, 0.5, 0.9, 1.0\}$ 125
- $\tau \in \{0.0, 0.05, 0.1, 0.4, 0.8, 1.0\}$ 126
- $net \in \{\{64, 64, 64\}, \{512, 256, 128\}\}$ a různé 127
velikosti sítí actor a critic 128
- $delay \in \{1, 4, 10\}$ 129

Výsledky experimentů můžeme vidět na obrázcích 130
na konci tohoto dokumentu počínaje obrázkem 3. 131
Můžeme si všimnout, že i při dlouhém tréningu 132
nenastalo zlepšení, ale agent se nejlépe choval už 133
po méně než 5000 epizodách. Parametr γ je nejlepší 134
nechat co nejvyšší, velice blízko hodnotě 1.0. Zdá se, 135
že architektury jednotlivých sítí na výsledek neměly 136



Obrázek 2. Srovnání průběhu učení agenta pro naši implementaci a pro implementaci ze stable-baselines-3.

nějaký značný vliv a všechny experimenty se chovaly velice podobně.

Frekvence provádění optimalizace sítě actor je lepší ponechat na nižších hodnotách, experimentálně dopadl nejlépe výchozí parametr s hodnotou 2. U parametru τ se jeví nejlepší ponechat ho velice blízko nule, zajímavého výsledku si můžeme všimnout na obrázku 18, kde se dlouho nic nedělo a ke konci učení odměny explodovaly. Při vyšších nastaveních se agent neučí vůbec.

5. Závěr

Hlavní náplní této práce bylo naprogramovat vlastní implementaci actor-critic algoritmu. Vytvořený algoritmus svými specifiky odpovídá metodě TD3 a je schopný řešit úlohy z prostředí MuJoCo. Učení agenta používajícího náš algoritmus dosahuje dobrých výsledků a je velice podobné učení agenta používajícího existující implementaci TD3 ze stable-baselines-3.

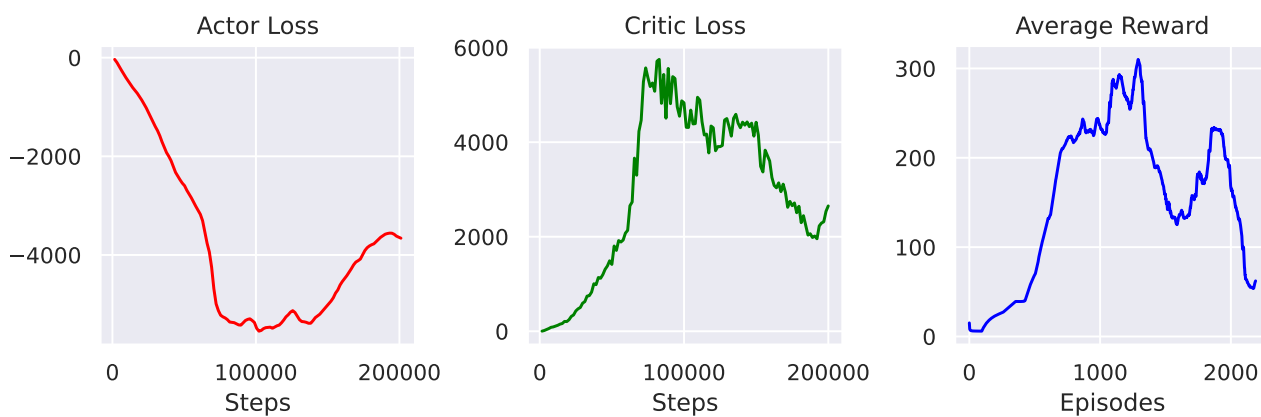
Dále jsme v této práci ukázali, že je možné v naší implementaci nastavovat různé hyperparametry jednotlivých komponent algoritmu, což jsme následně demonstrovali na několika experimentech, které ukazují vliv hyperparametrů na učení agenta.

Porovnání průběhu učení agenta s existující implementací proběhlo na třech MuJoCo prostředích a ve všech případech lze pozorovat velice podobný tvar křivky průměrných odměn. Lze tedy předpokládat, že náš algoritmus dosahuje podobné úrovně kvality jako již existující implementace.

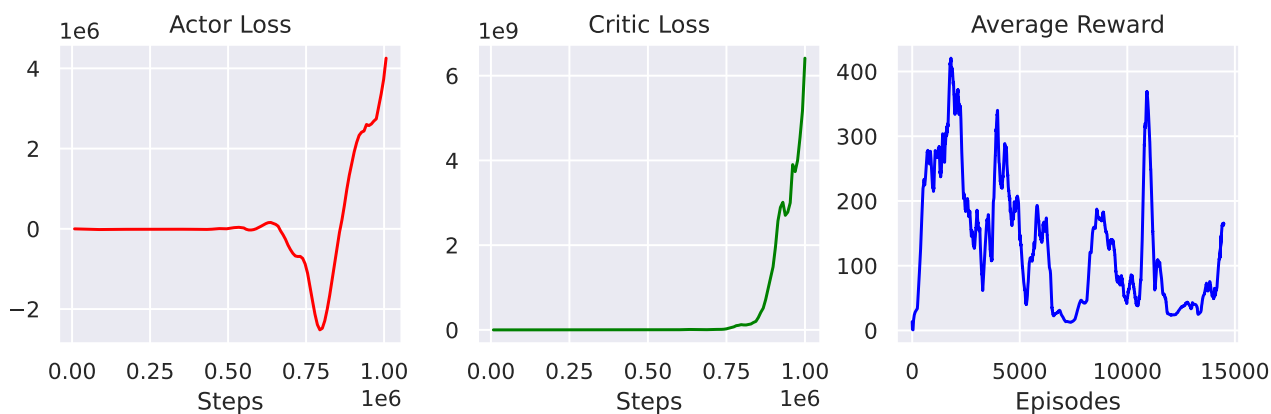
Implementace našeho algoritmu umožňuje experimentování ať už na nejnižší úrovni se samotnou metodou TD3, tak na vyšší úrovni s trénováním agenta s různými hyperparametry na libovolných prostředích knihovny OpenAI Gym (tedy nejsou limitovány pouze na MuJoCo). Výsledky agentova učení lze porovnávat s jinými existujícími implementacemi.

Literatura

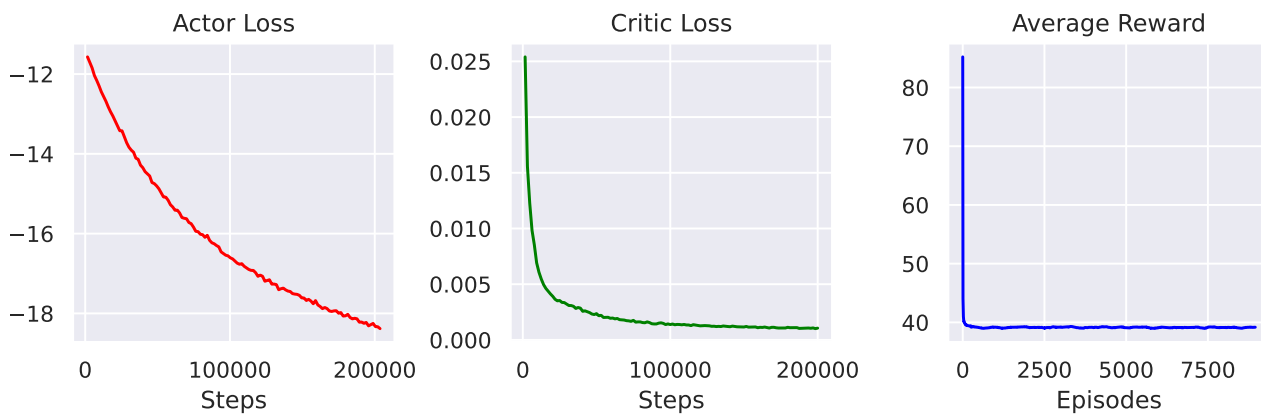
- [1] THE MATHWORKS, I. *Twin-Delayed Deep Deterministic Policy Gradient Agents* [documentation (english)]. Březen 2022. Dostupné z: <https://www.mathworks.com/help/reinforcement-learning/ug/td3-agents.html>.
- [2] TODOROV, E., EREZ, T. a TASSA, Y. MuJoCo: A physics engine for model-based control. In: IEEE. *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*. 2012, s. 5026–5033.
- [3] TUNYASUVUNAKOOL, S., MULDAL, A., DORON, Y., LIU, S., BOHEZ, S. et al. Dm_control: Software and tasks for continuous control. *Software Impacts*. 2020, roč. 6, s. 100022. Dostupné z: <https://www.sciencedirect.com/science/article/pii/S2665963820300099>. ISSN 2665-9638.



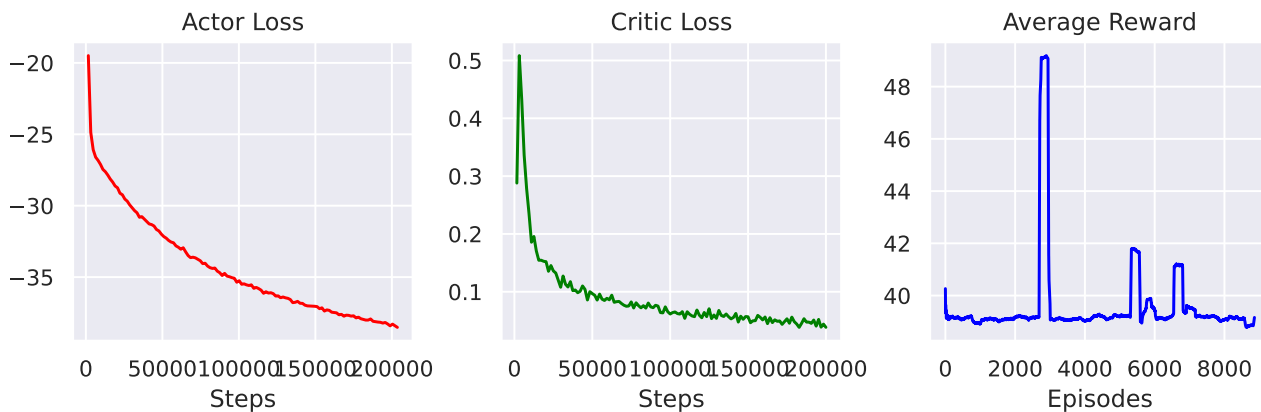
Obrázek 3. Výchozí nastavení parametrů.



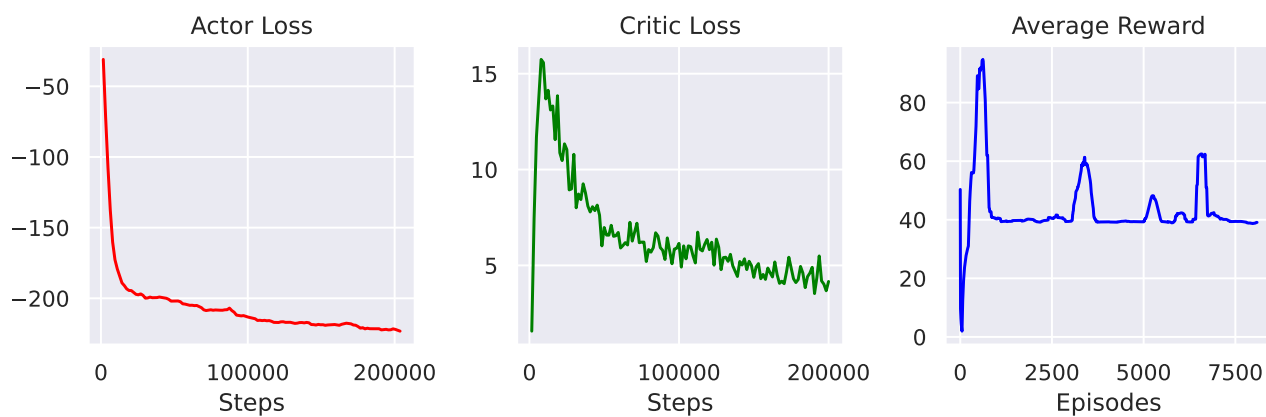
Obrázek 4. Výchozí parametry pro 1M kroků.



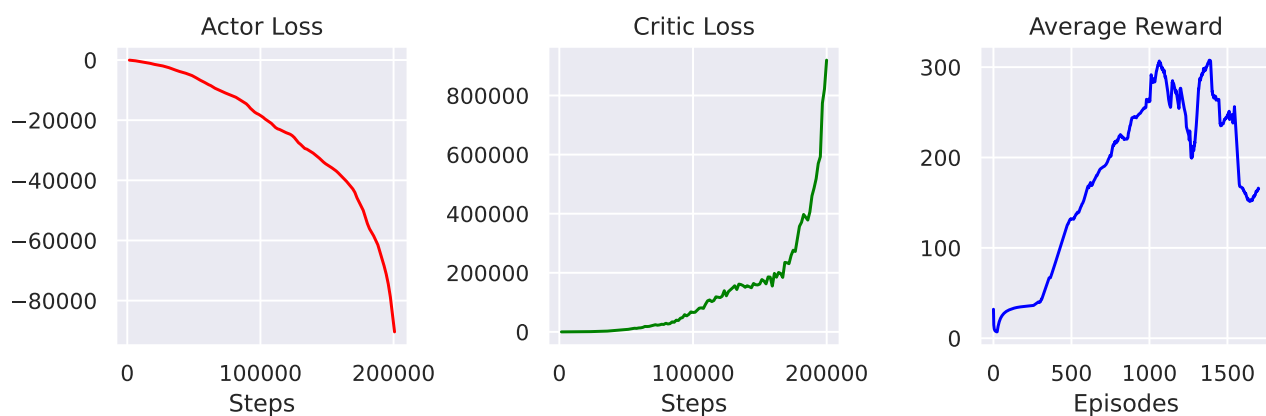
Obrázek 5. $\gamma = 0.0$.



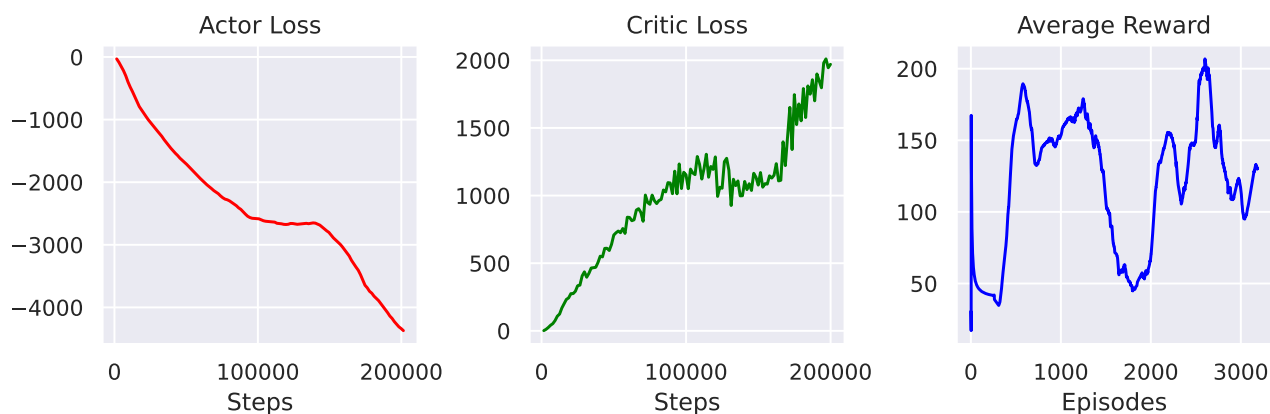
Obrázek 6. $\gamma = 0.5$.



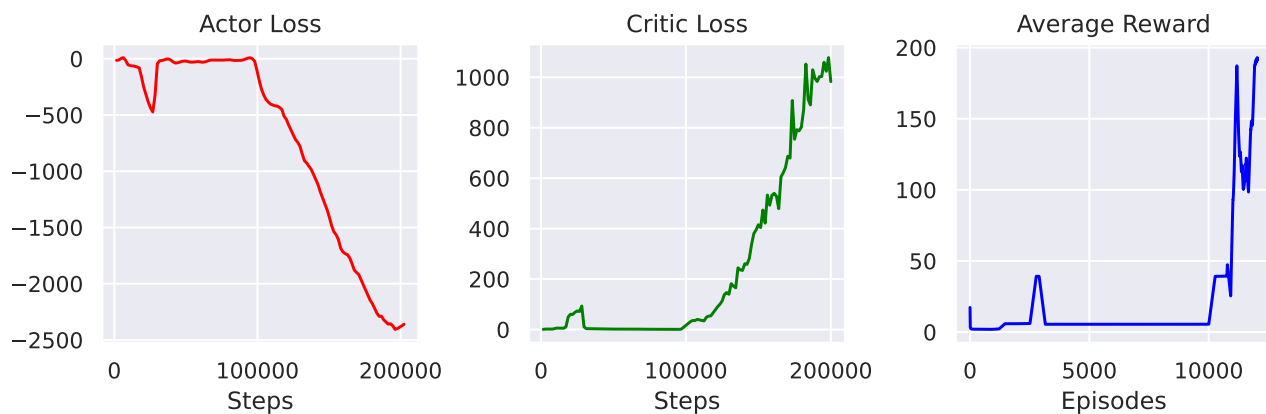
Obrázek 7. $\gamma = 0.9$.



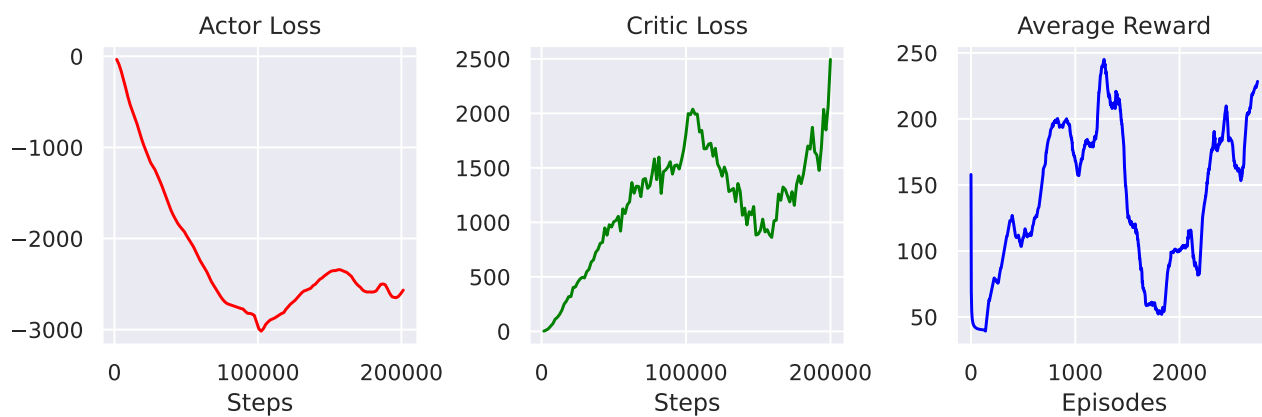
Obrázek 8. $\gamma = 1.0$.



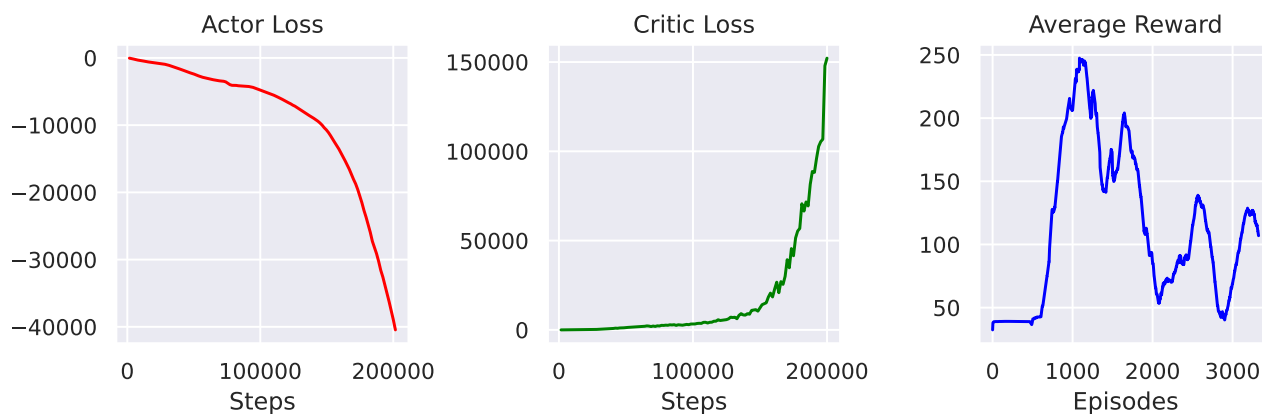
Obrázek 9. $nets = \{64, 64, 64\}$.



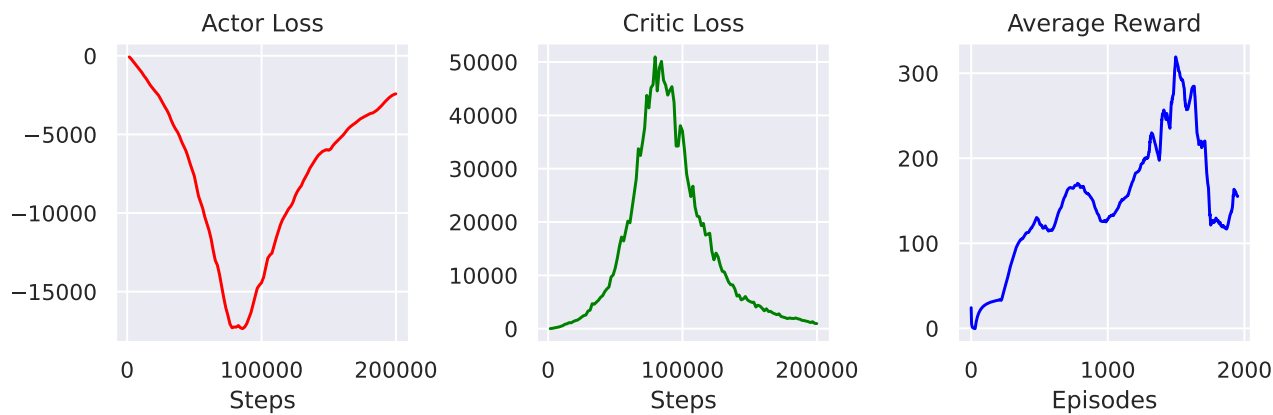
Obrázek 10. $nets = \{512, 256, 128\}$.



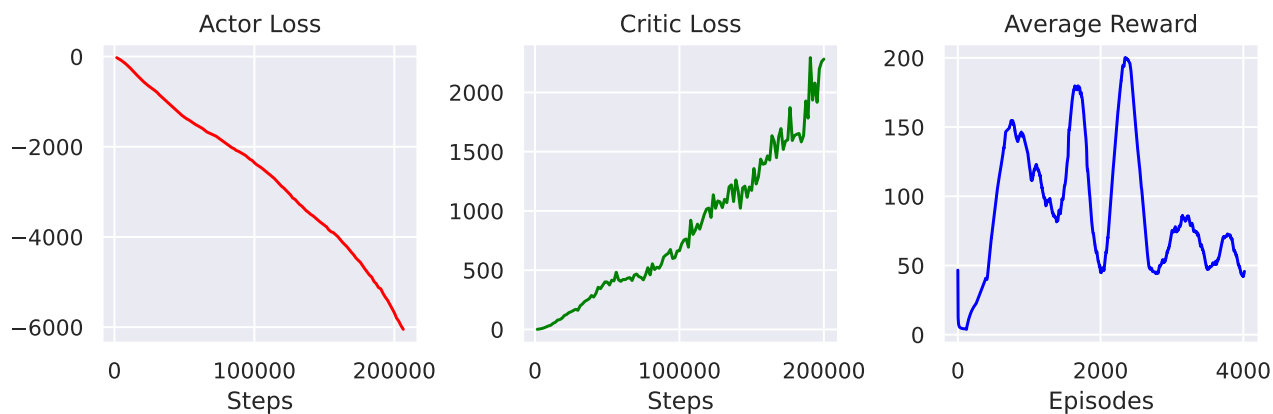
Obrázek 11. $actor = \{256, 256\}, critic = \{256, 256, 256\}$.



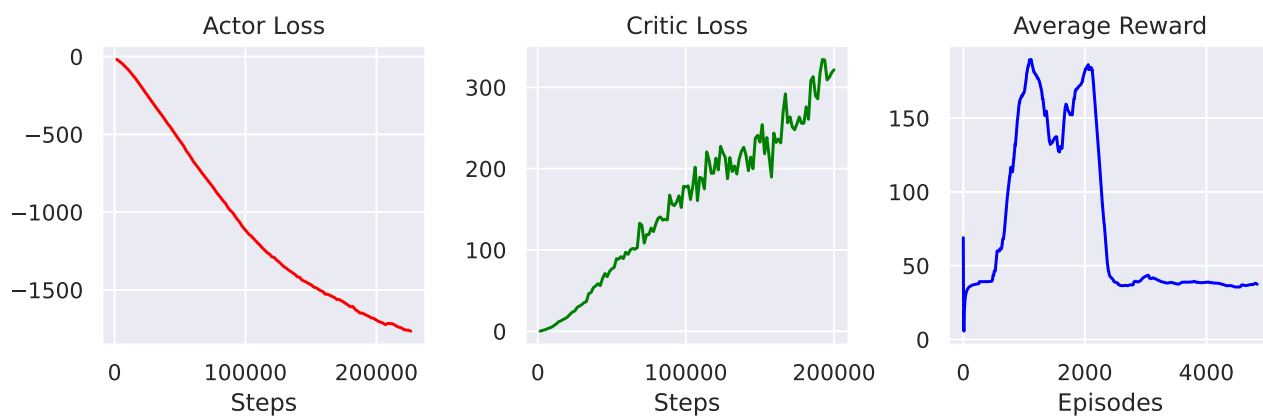
Obrázek 12. $actor = \{256, 256, 256\}, critic = \{256, 256\}$.



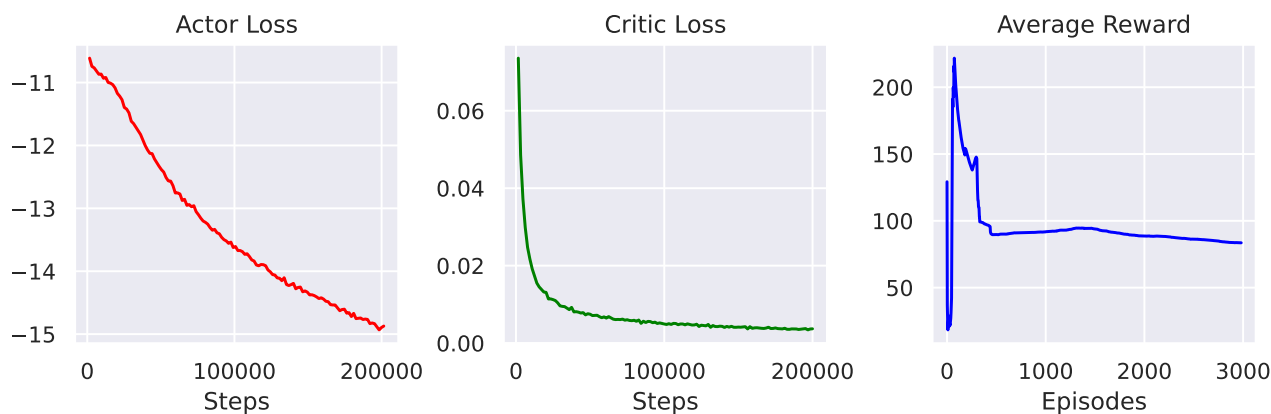
Obrázek 13. $delay = 1$.



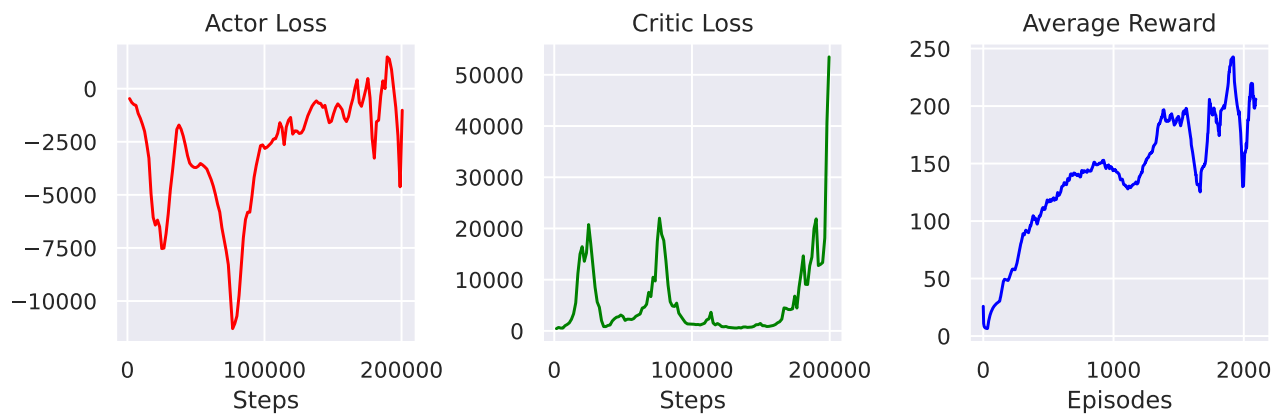
Obrázek 14. $delay = 4$.



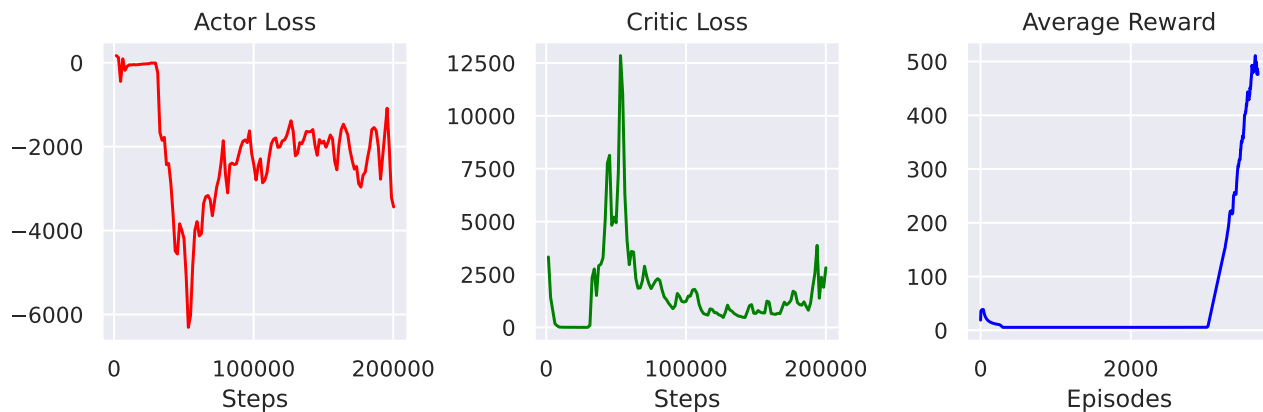
Obrázek 15. $\text{delay} = 10$.



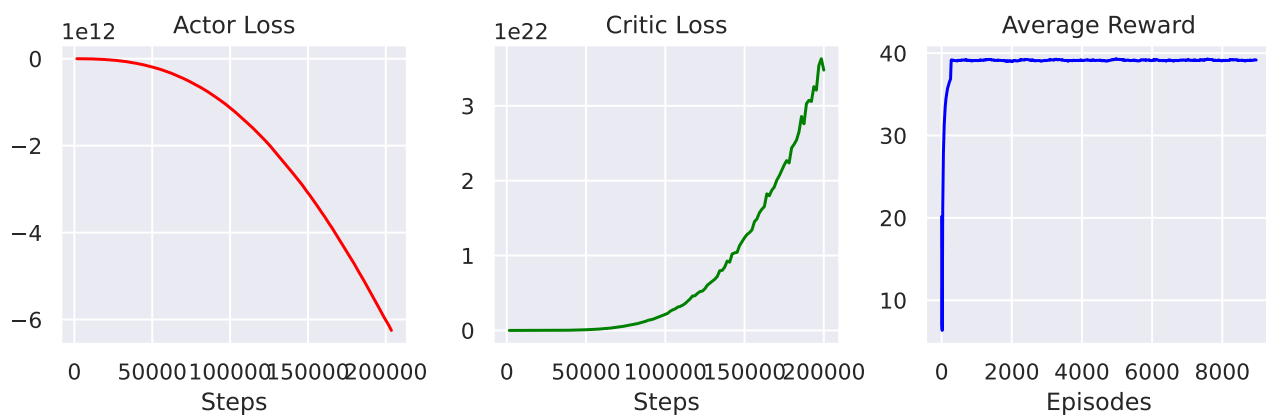
Obrázek 16. $\tau = 0.0$.



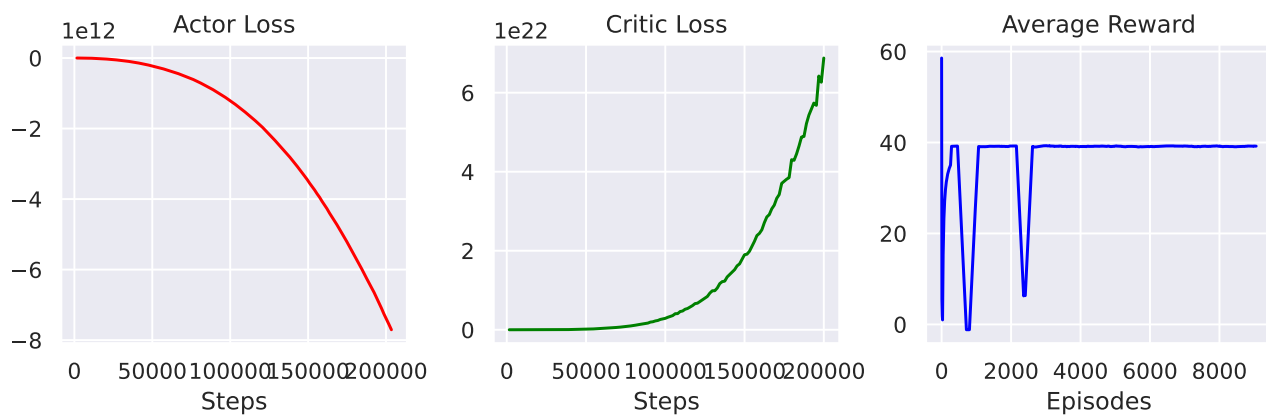
Obrázek 17. $\tau = 0.05$.



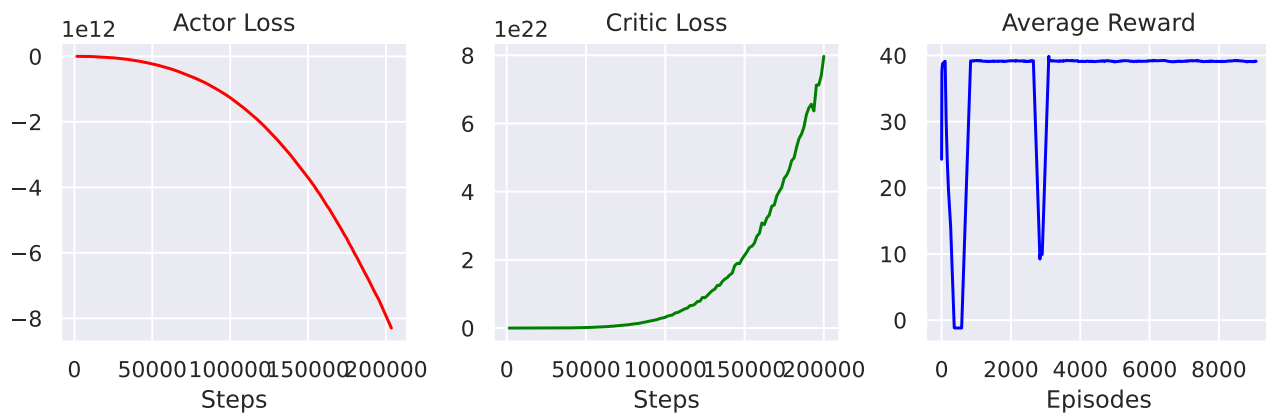
Obrázek 18. $\tau = 0.1$.



Obrázek 19. $\tau = 0.4$.



Obrázek 20. $\tau = 0.8$.



Obrázek 21. $\tau = 1.0$.