

Využitie neurónových sietí pri simulovanej robotike

Bc. Martin Kubovčík

Vedúci: prof. RNDr. Jiří Pospíchal, DrSc.

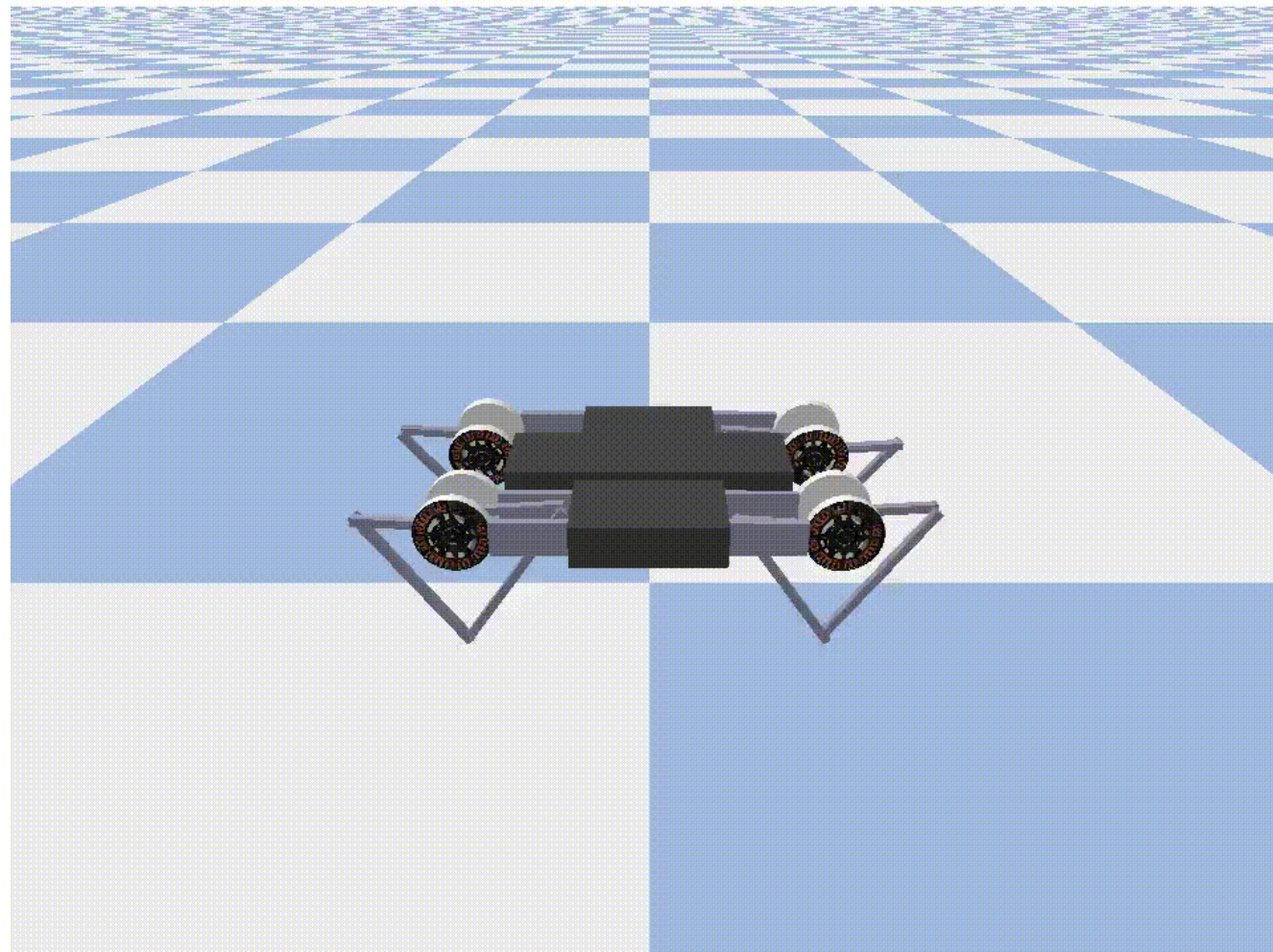
Obsah

- Prostredie
- Soft Actor-Critic
- DeepMind Reverb
- Vylepšenia
- Imitačné učenie
- Sim-to-Real

Prostredie

- OpenAI Gym
- PyBullet
- DeepMind Control Suite

Prostredie MinitaurBulletEnv-v0

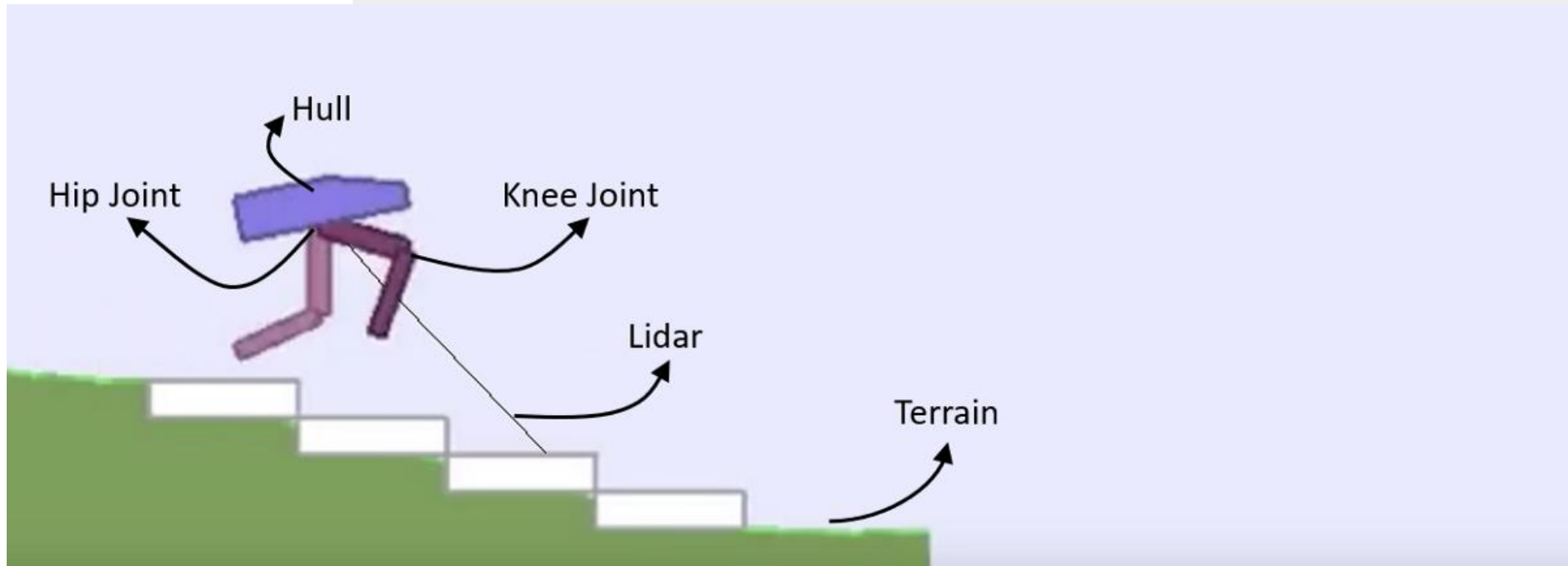


Prostredie MinitaurBulletEnv-v0

- 28 príznačov v stavovom priestore
- 8 akčných členov
- Rozsah akcií [-1.0, 1.0]
- **Vstupy:** uhly motorov, rýchlosti rotácie motorov, krútiace momenty motorov, pozícia robota, orientácia robota
- **Odmena:** pohyb vpred, trest za nepriamočiaru chôdzu, trest za trasenie sa, trest za používanie motorov

Prostredie

BipedalWalkerHardcore-v3

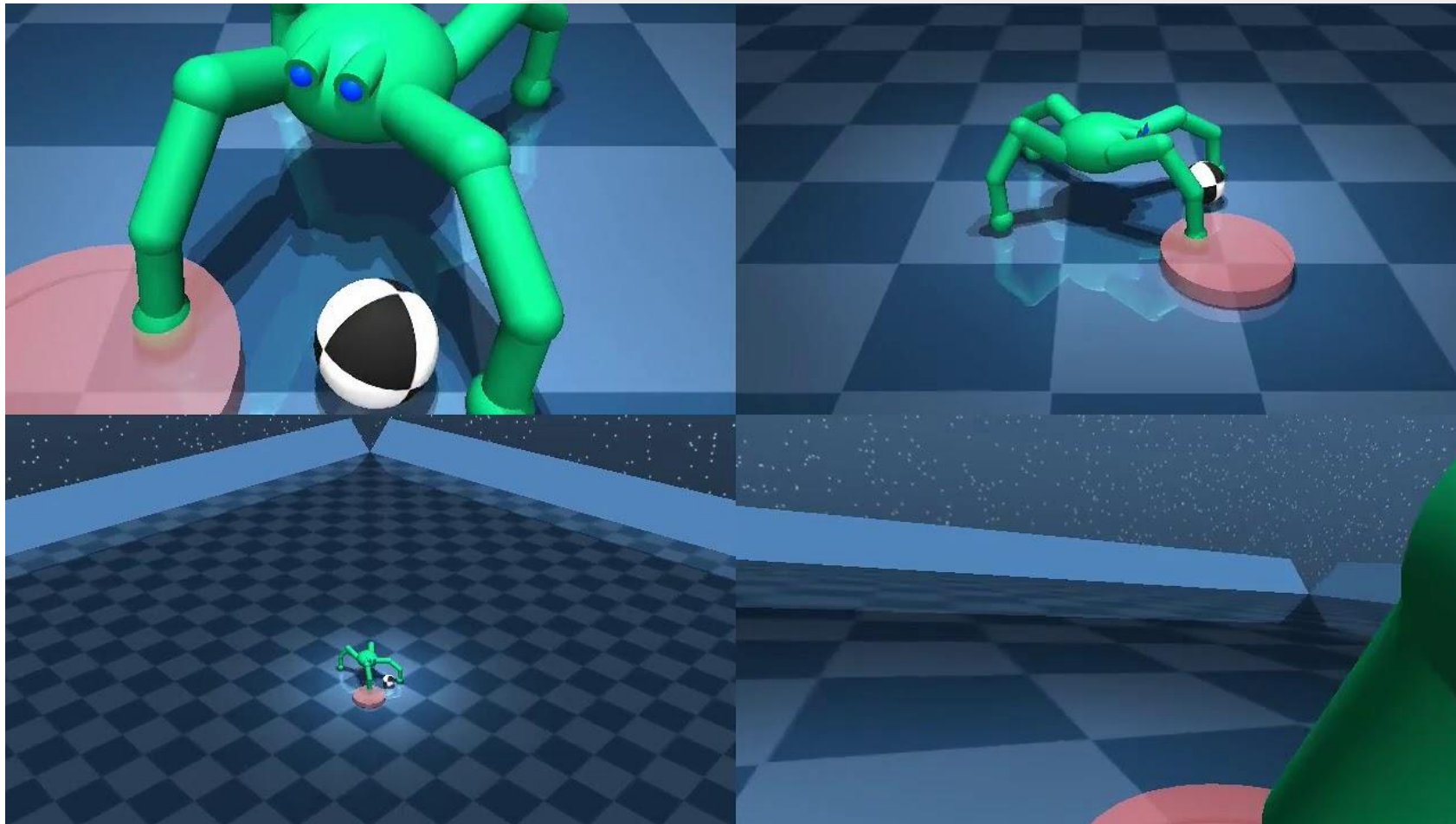


Prostredie

BipedalWalkerHardcore-v3

- 24 príznačov v stavovom priestore
- 4 akčné členy
- Rozsah akcií [-1.0, 1.0]
- **Vstupy:** uhol trupu, uhlová rýchlosť trupu, rýchlosť robota, uhly kĺbov, rýchlosti rotácie kĺbov, kontakty nôh so zemou, LIDAR
- **Odmena:** pohyb vpred (normalizovaný na 300 bodov ak dosiahne cieľ), držanie hlavy rovno, trest za používanie motorov, -100 za pád na zem

Prostredie Quadruped



Prostredie Quadruped

- 90 príznačov v stavovom priestore
- 12 akčných členov
- Rozsah akcií [-1.0, 1.0]
- **Vstupy:** uhly kĺbov, rýchlosti rotácie kĺbov, rýchlosť robota, senzor IMU (GYRO+ACC), rýchlosť rotácie, rýchlosť lopty relatívna k rýchlosti robota, pozícia lopty relatívna k pozícii robota, cieľová pozícia pre loptu relatívna k pozícii robota
- **Odmena:** vzdialenosť od lopty, vzdialenosť lopty od cieľa, trup vzpriamený v osi z (+1 ak je v tolerancii, inak <1)

Q-Learning

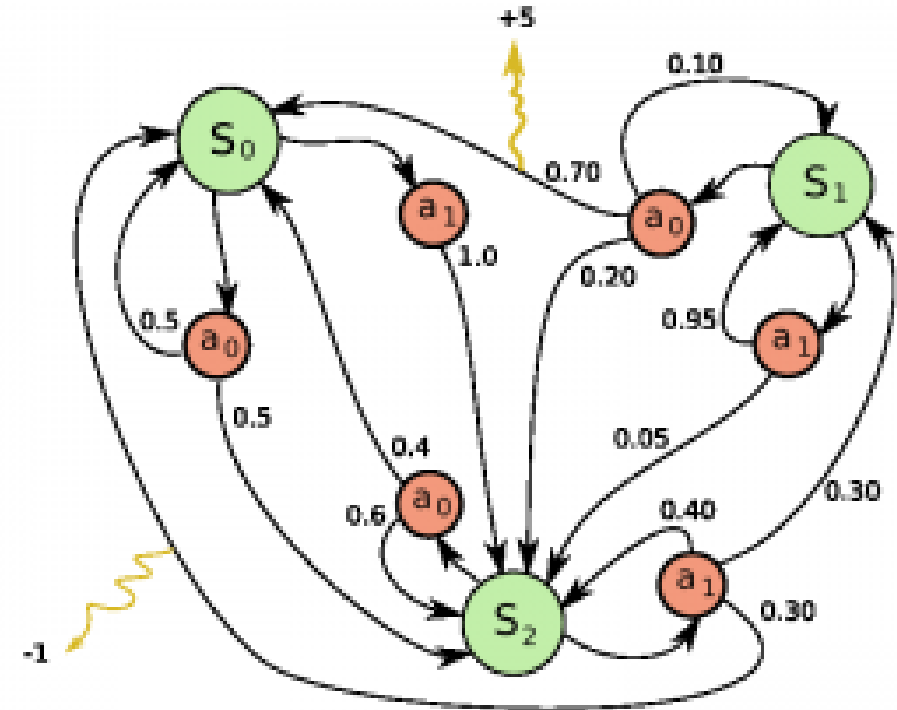
- Off-policy algoritmus
- Markov decision process

- s_0, a_0, r_1, s_1, a_1
kvalita Q vykonanej akcie a_t
discount factor

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

maximálna kvalita Q akcie na za vykonanú akciu a_t v stave s_t v zvähu γ

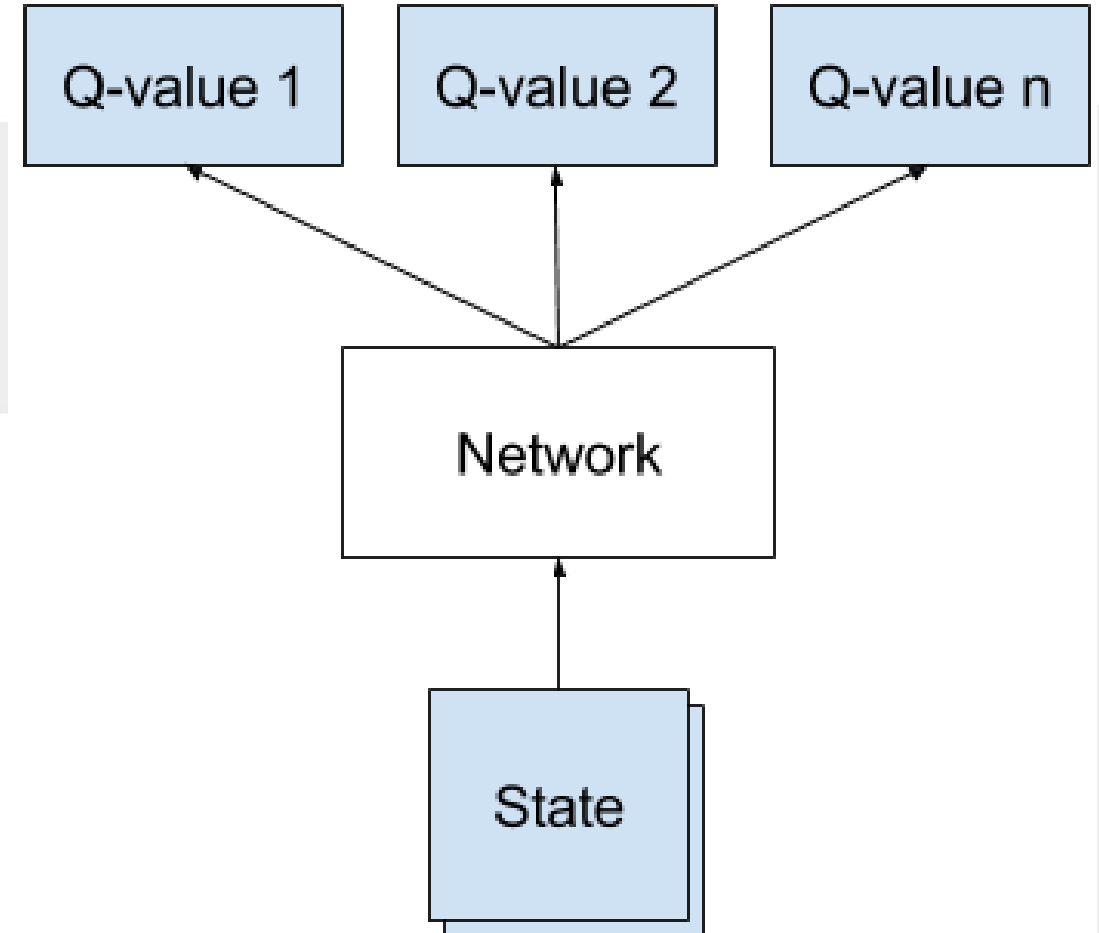
- Čím viac ideme do budúcnosti tým nižšiu váhu γ tomu prikladáme



Zdroj: <https://neuro.cs.ut.ee/demystifying-deep-reinforcement-learning/>

Deep Q Network

- Univerzálny aproximátor nelineárnych funkcií
- Transformácia vysoko rozmerných dát (2D obraz) do Q-hodnôt
- Problémy s nestabilitou počas učenia (local minimum) – množstvo trikov – target network, experience replay buffer



Zdroj: <https://neuro.cs.ut.ee/demystifying-deep-reinforcement-learning/>

Soft Actor-Critic

Policy Gradients



Go Right



Deep Q-Learning

Please wait, I am still calculating Q value, only 41891 actions left...

Soft Actor-Critic

- Maximalizácia entropie
- Maximalizácia Q funkcie
- Dvojica nezávislých modelov Actor (policy) a Critic (Q value)

$$H(X) = - \sum_{i=1}^n p(x_i) \log_b p(x_i) = -\mathbb{E} [\log \mathcal{N}(\mu, \sigma^2)]$$

entropia náhodnej premennej X

vyjadrenie entropie pre Gaussovu distribúciu

pravdepodobnosť vygenerovania konkrétnej hodnoty x_i

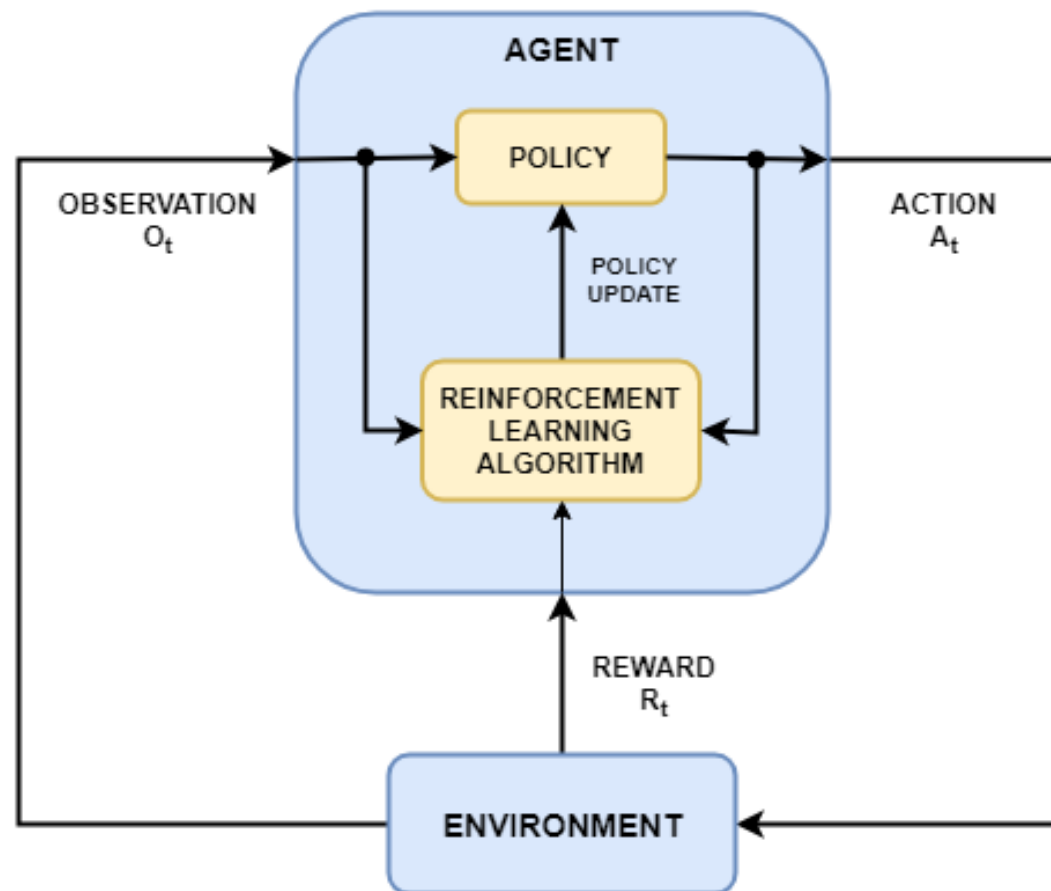
Soft Actor-Critic

- Gaussova distribúcia – predikcia priemerov a smerodajných odchýlok
- Efektívnejšie prehľadávanie na rozdiel od Twin Delayed DDPG (TD3) alebo Deep Deterministic Policy Gradient (DDPG)
- Menej citlivé na ladenie hyperparametrov
- Adaptívny podiel entropie na chybových funkciách



Zdroj: <https://hbr.org/2007/11/a-leaders-framework-for-decision-making>

Soft Actor-Critic



Soft Actor-Critic

- Odmeny sú v rozsahu $[-1.0, 1.0]$ – funkcia tanh
- Chybové funkcie:

$$J_{\pi}(\phi) = \mathbb{E}_{s_t \sim D, \epsilon_t \sim \mathcal{N}} \left[\alpha \log \pi_{\phi}(f_{\phi}(\epsilon_t, s_t) | s_t) - Q_{\theta}(s_t, f_{\phi}(\epsilon_t, s_t)) \right]$$

entropia distribúcie

podiel entropie na chybovej fu

predikovaná Q-hodnota

$$J_Q(\theta) = \mathbb{E}_{(s_t, a_t) \sim D} \left[\frac{1}{2} (Q_{\theta}(s_t, a_t) - (r(s_t, a_t) + \gamma(Q_{\theta}(s_{t+1}, a_{t+1}) - \alpha \log \pi_{\phi}(a_{t+1} | s_{t+1}))))^2 \right]$$

funkcia odmer

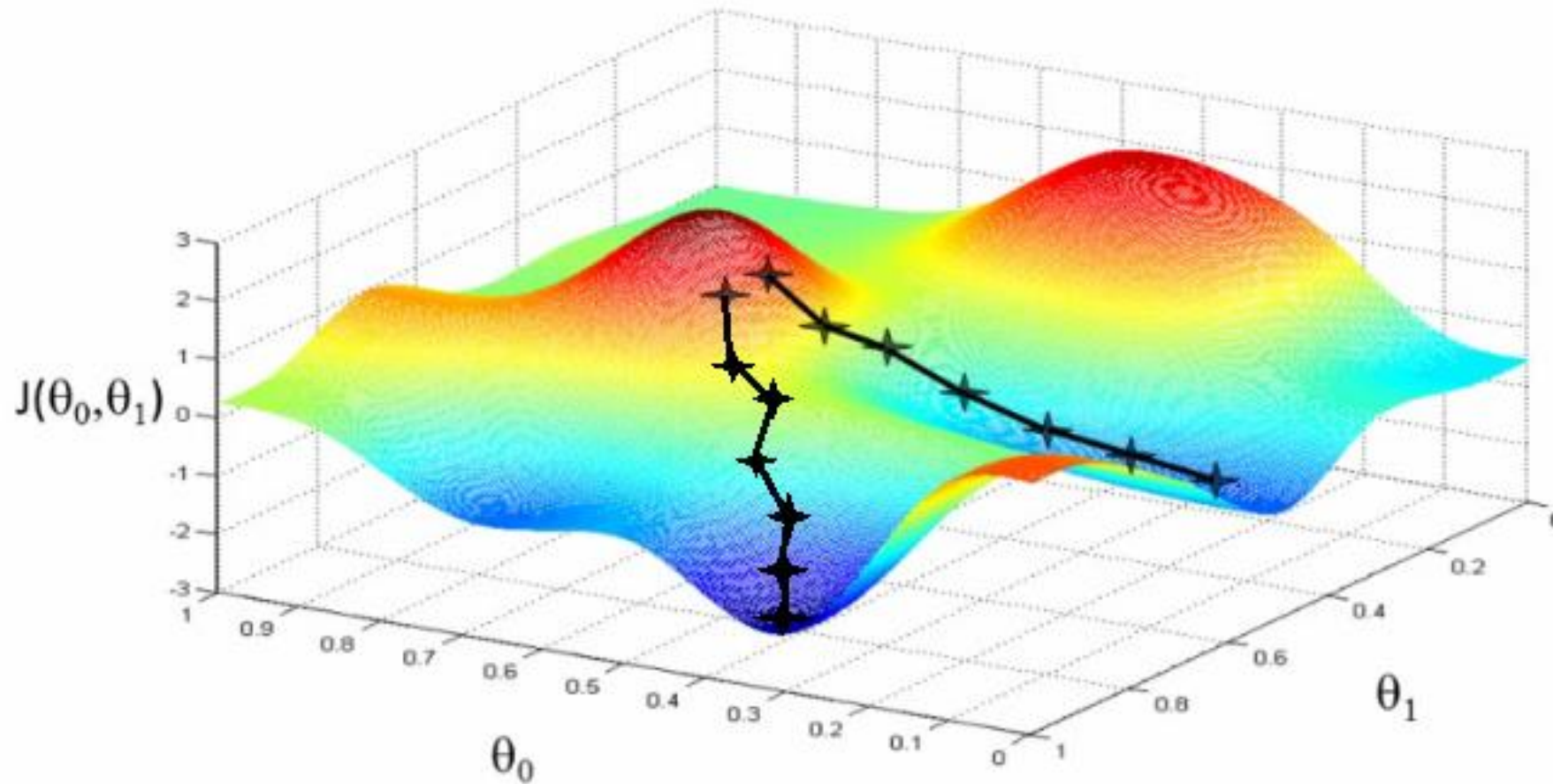
faktor

Q-hodnota pre nasledujúci stav

entropia distribúcie

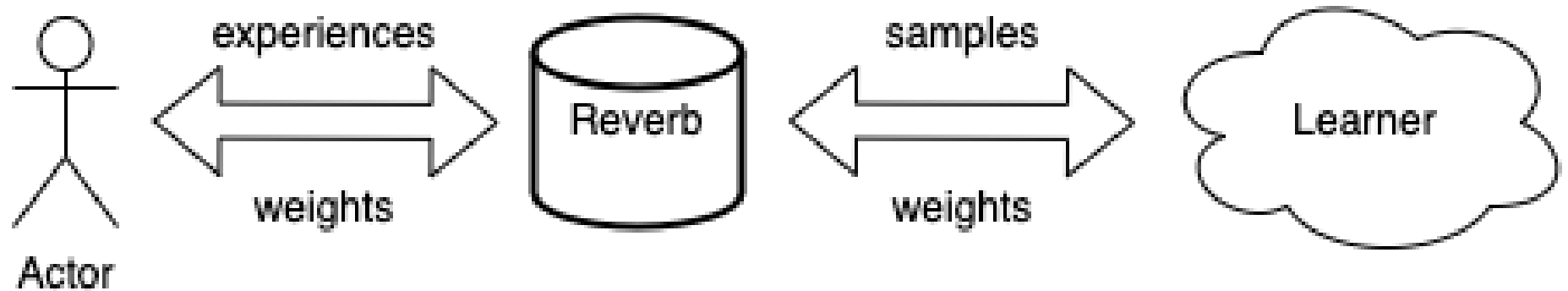
Q-hodnota pre aktuálny stav

Soft Actor-Critic



DeepMind Reverb

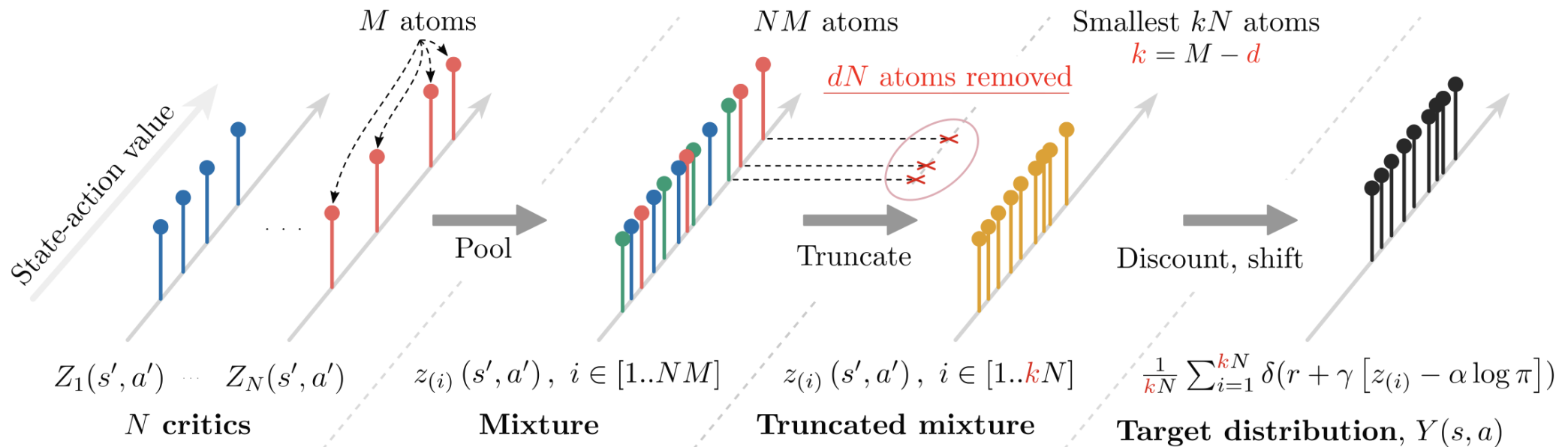
- Klient-server architektúra
- Actor – zbiera skúsenosti v hernom prostredí
- Learner – aktualizuje parametre neurónových sietí
- Databáza – distribúcia dát po počítačovej sieti (Experience Replay)



Vylepšenia

Controlling Overestimation Bias with Truncated Mixture of Continuous Distributional Quantile Critics

Controlling Overestimation Bias with Truncated Mixture of Continuous Distributional Quantile Critics



Zdroj: <http://proceedings.mlr.press/v119/kuznetsov20a/kuznetsov20a.pdf>

Vylepšenia

Generalized State-Dependent Exploration (gSDE)

- Parametre šumu θ_ϵ sú generované len každých n krokov počas epizódy

- Namiesto stavu s_t

$$a_t = \mu(s_t; \theta_\mu) + \epsilon(s_t; \theta_\epsilon)$$

$$\epsilon(s; \theta_\epsilon) = \theta_\epsilon z_\mu(s)$$

$$\theta_\epsilon \sim N(0, \sigma^2)$$

deterministická akcia prediktor

náhodná premenná z Gaussovej distribúcie s nulovým priemerom a smerodajnou odchýlkou danou trénovanými parametrami

latentný priestor modelu

šumová zložka závislá od stavu

Vylepšenia

Generalized State-Dependent Exploration (gSDE)

(gSDE)

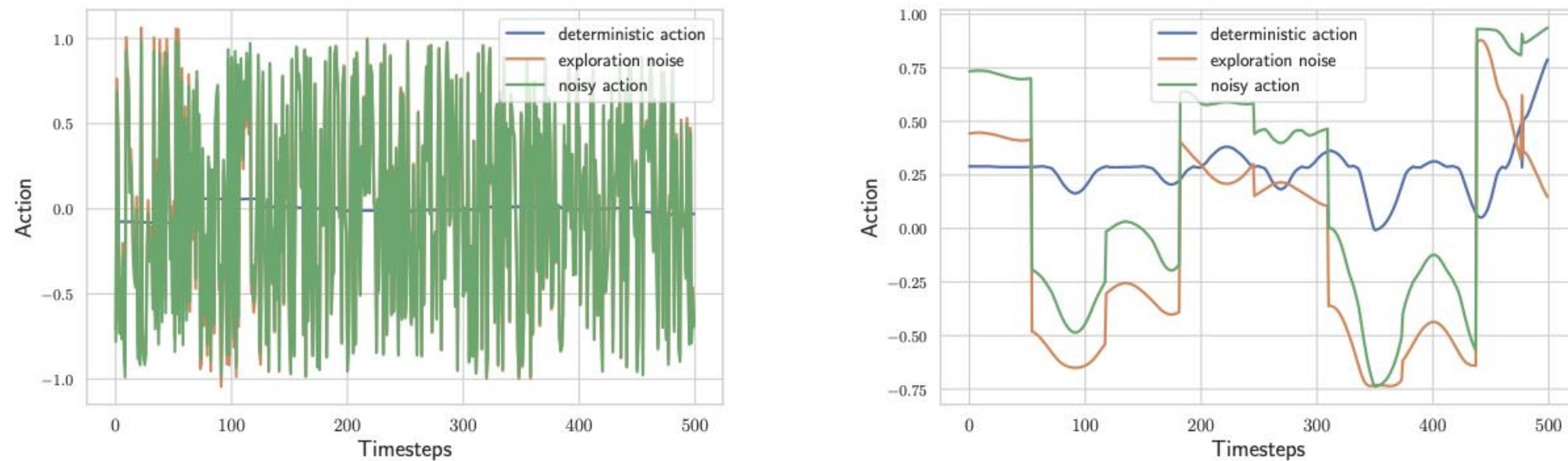


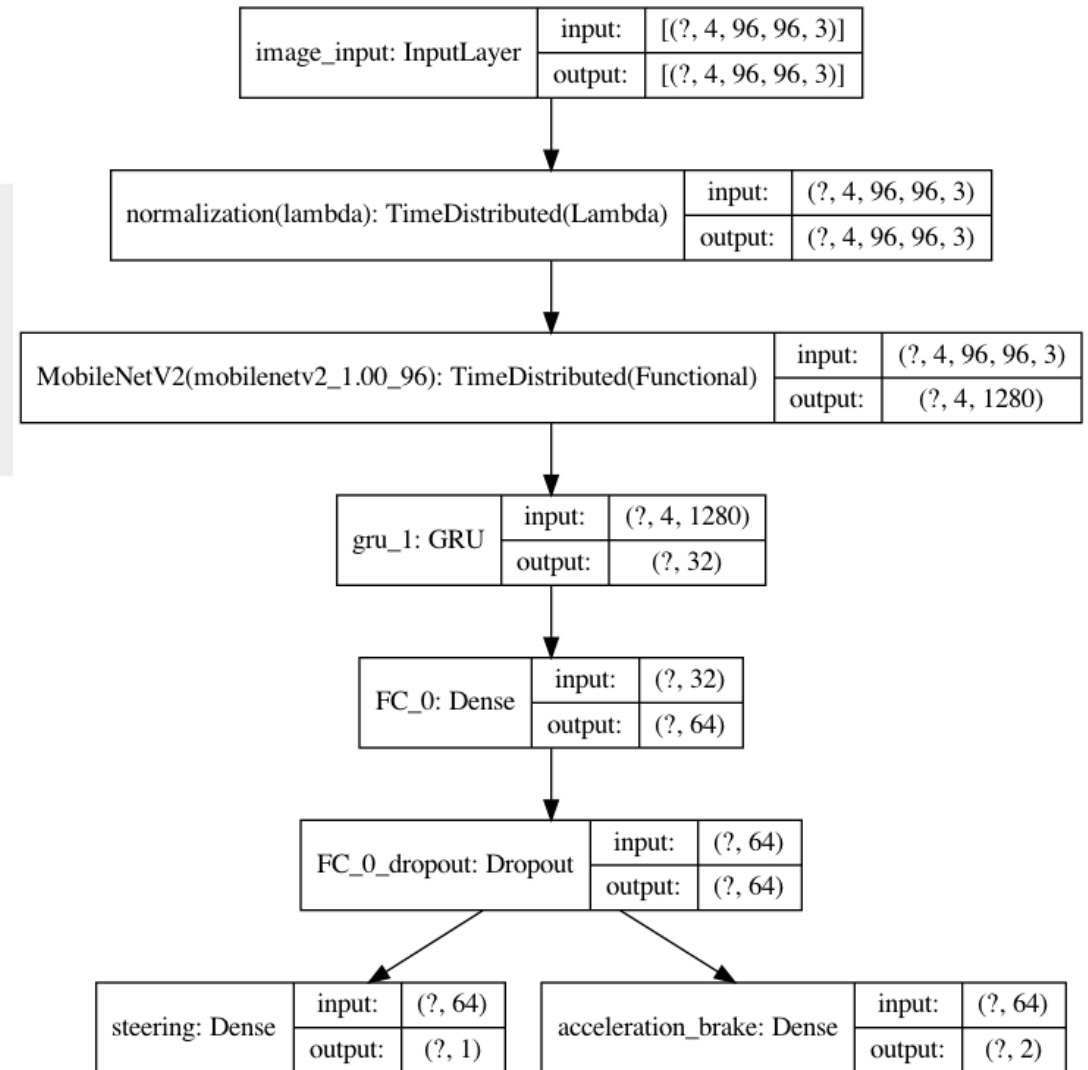
Figure 1: Left: unstructured exploration, as typically used in simulated RL. Right: gSDE provides smooth and consistent exploration.

Imitačné učenie

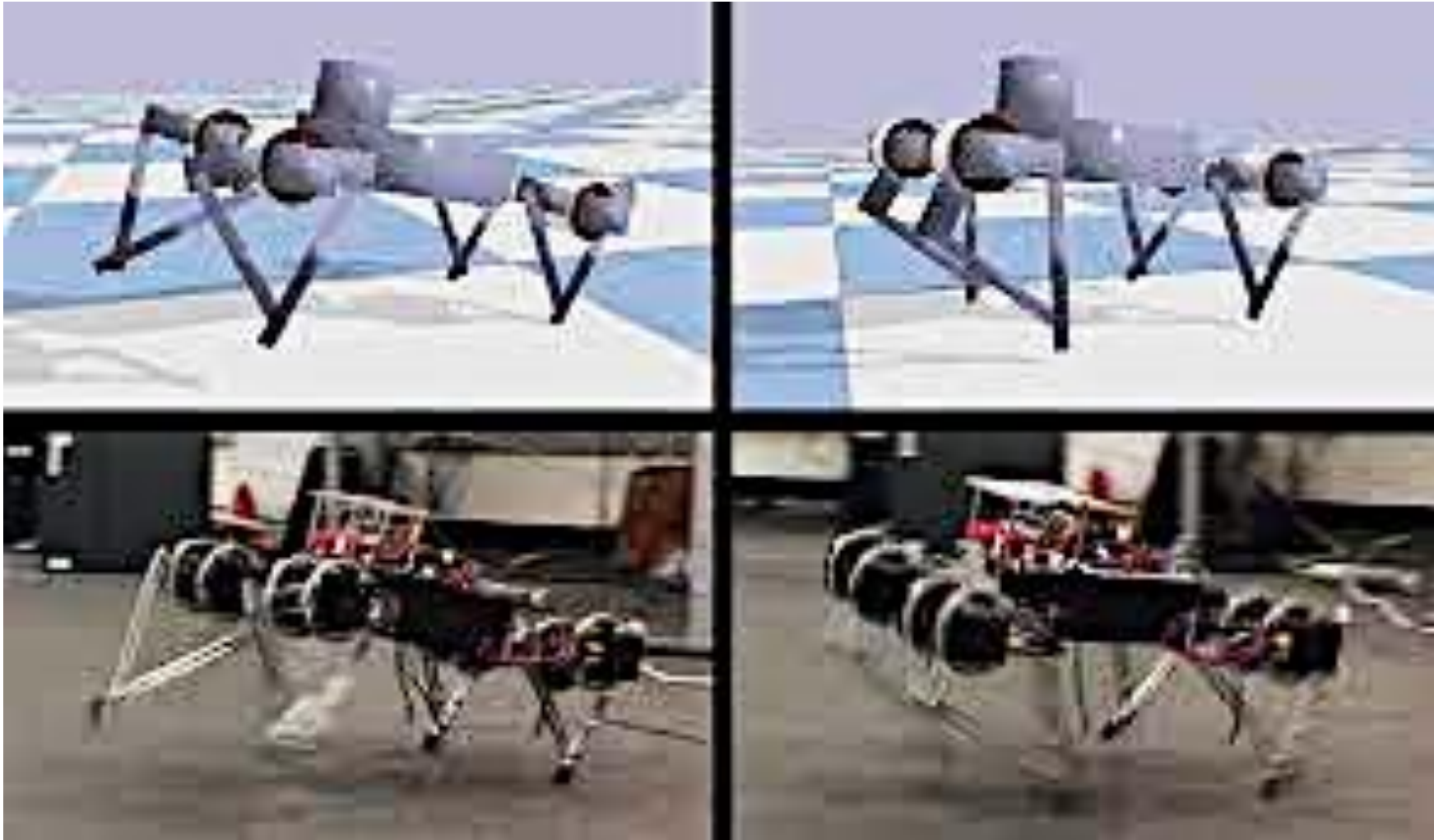


Imitačné učenie

- Vstup: RGB obraz 96x96 pixelov, normalizované do rozsahu (-1, 1)
- Rekurentná jednotka GRU kvôli časovej sérii obrazov
- MobileNetV2
- Detekcia smeru pohybu auta voči okoliu
- Výstup: volant (-1, 1), plynový pedál (0, 1), brzdo­vý pedál (0, 1)
- Využitelnosť v reálnej robotike – kamera, GPS

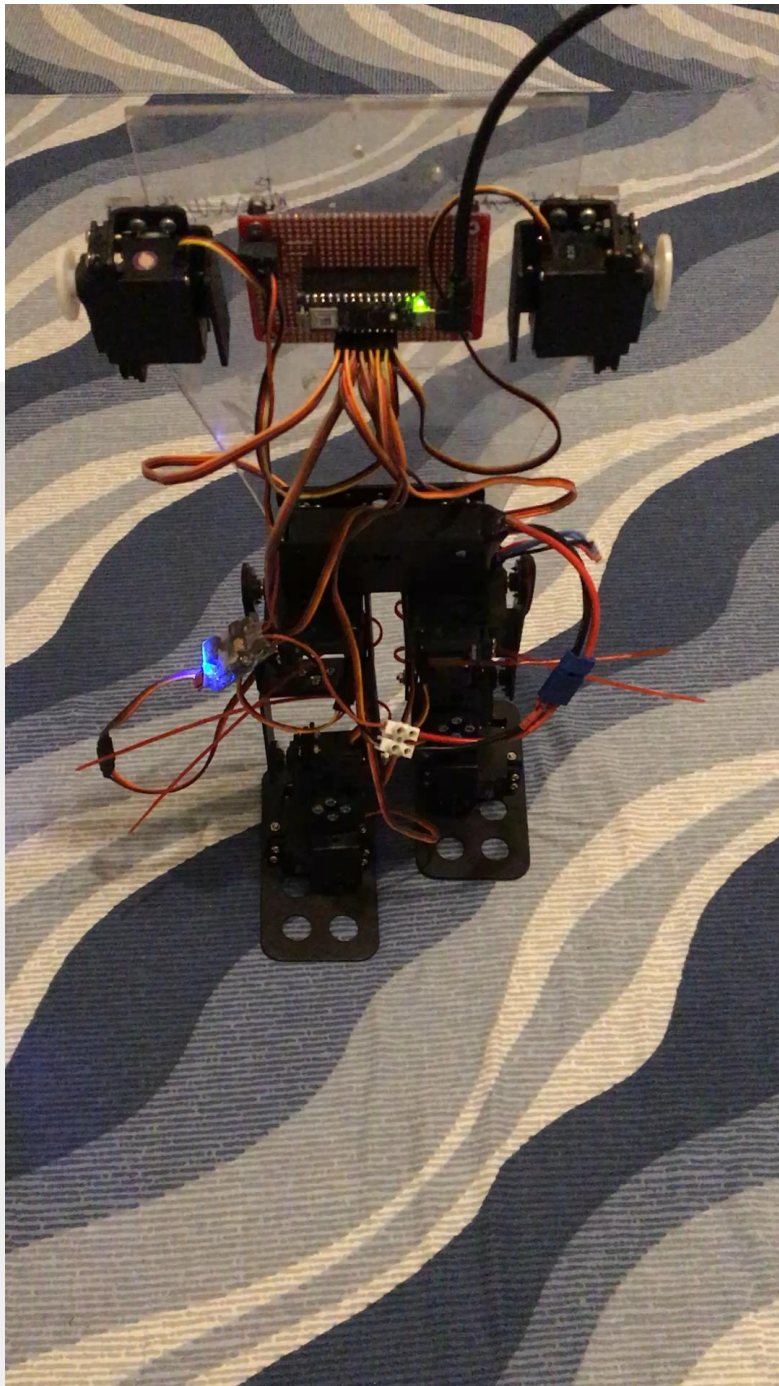


Sim-to-Real



Sim-to-Real

- Actor – predstavuje robota v simulovanom/reálnom svete
- Learner – server alebo cloud – trénuje robota
- Databáza – server alebo cloud – uchováva skúsenosti robota



Ďakujem za pozornosť

GitHub: <https://github.com/markub3327>

Email: kubovcik1@ucm.sk

Instagram: <https://www.instagram.com/martin.kubovcik/>

Twitter: <https://twitter.com/markub3327>