

# Učenie sa basketbalových trikov s použitím optimalizácie trajektórie a hlbokého zosilneného učenia

***Autor:*** Miroslav Kobliška

O clanku **Learning Basketball Dribbling Skills Using  
Trajectory Optimization and Deep Reinforcement Learning**

LIBIN LIU, DeepMotion Inc., USA

JESSICA HODGINS, Carnegie Mellon University, USA

# Úvod

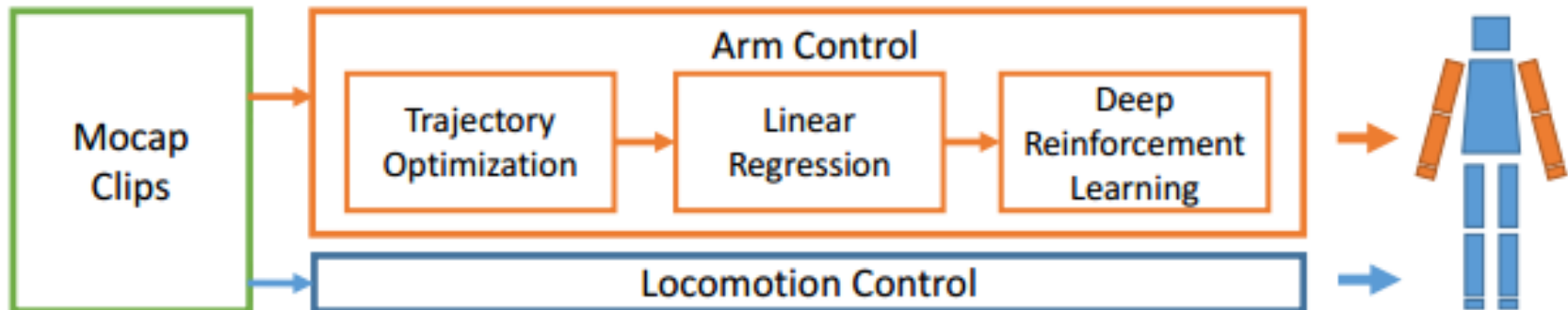
- V tejto práci animovali basketbal pomocou motion capture dát
- Používali aj fyzikálne vlastnosti na realistické prevedenie
- Táto metóda spracováva basketbalový kontrolér ako kombináciu kontroly pohybu a kontroly častí paže.

# Úvod 2

- Systém sa najskôr naučí kontrolovať pohyb a potom natrénuje kontrolu paží
- Vie sa naučiť rôzne basketbalové triky ako :
- dribling medzi nohami aj krížne pohyby a vie reagovať na interakciu používateľa

# Prehľad systému

- klipy na snímanie pohybu ako vstup
- riadiaca jednotka sa skladá z dvoch spojených komponentov
- komponent kontrolujúci paže
- komponent kontrolujúci pohyblivosť



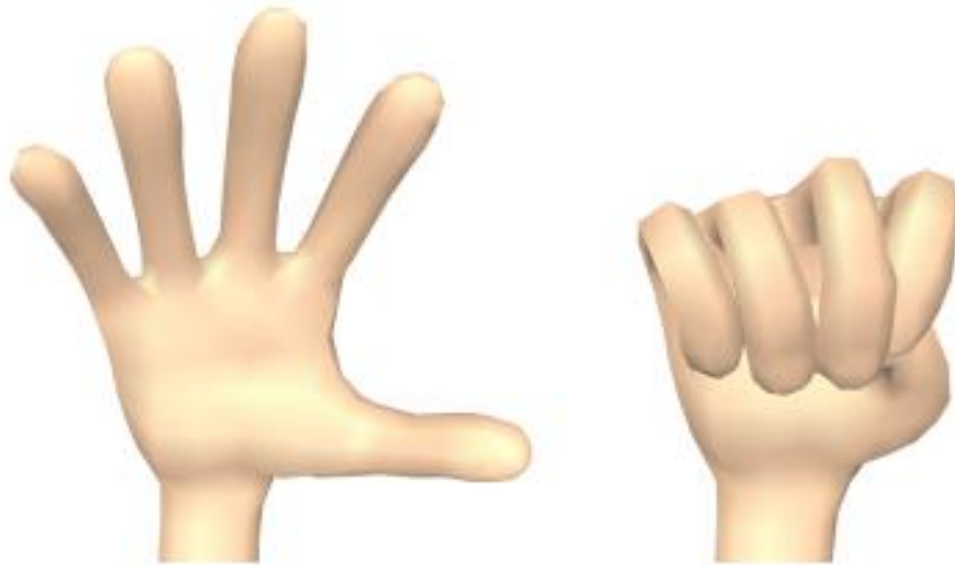
# Kostra

- Umelá kostra
- Skladá sa z pohyblivých a vnútorných kĺbov
- Hráčove ruky sú namodelované s prstami



# Prsty

$$\hat{\theta}_i = (1 - \lambda_i)\theta_i^{\text{flat}} + \lambda_i\theta_i^{\text{fist}}$$



(b) The flat hand pose (left) and the fist pose (right) used for computing the target pose for hands.

# Pohyby prstov

- systém kontroluje každú ruku pomocou troch kontrolných signálov

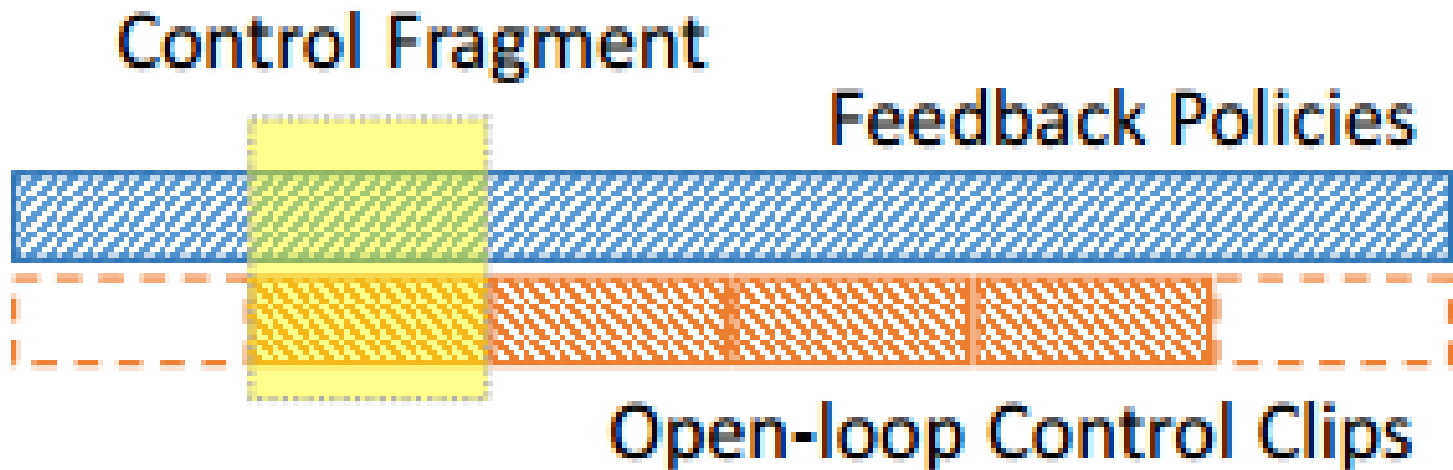
$$\alpha = \{\lambda_{\text{Thumb}}, \lambda_{\text{Index}}, \lambda_{\text{Pinky}}\}$$

$$\lambda_{\text{Middle}} = \frac{2}{3}\lambda_{\text{Index}} + \frac{1}{3}\lambda_{\text{Pinky}}$$

$$\lambda_{\text{Ring}} = \frac{1}{3}\lambda_{\text{Index}} + \frac{2}{3}\lambda_{\text{Pinky}}$$

# Kontrolné fragmenty

- Každý komponent na ovládanie má spätnú väzbu a „open-loop“ kontrolovanú trajektóriu





# Učenie sa kontroly pohyblivosti

- hráč vie kontrolovať loptu iba počas veľmi krátkeho času
- používa sa kontrolný fragment dĺžky 0.05 s
- keď je naučený komponent pre kontrolu pohyblivosti je fixný počas učenia komponentu kontroly paží

# Učenie kontroly paží

- používa optimalizáciu trajektórie na úspešné vypočítanie „open-loop“ kontroly paží
- „open-loop“ ovládanie paží je cielené na plece, lakeť a zápästie

# Optimalizácia trajektórie

- nájsť množinu korekčných kompenzácií, aby hráč mohol úspešne driblovať s loptou
- frejmy, v ktorých sa lopta dotýka hráčovej ruky
- Teda určité frejmy sú ako „chackpointy“
- Tu je vynútený dotyk ruky s loptou
- „Chackpointy“ - body, kde lopta dosiahne maximálnu výšku

# Kontrola trajektórie 2

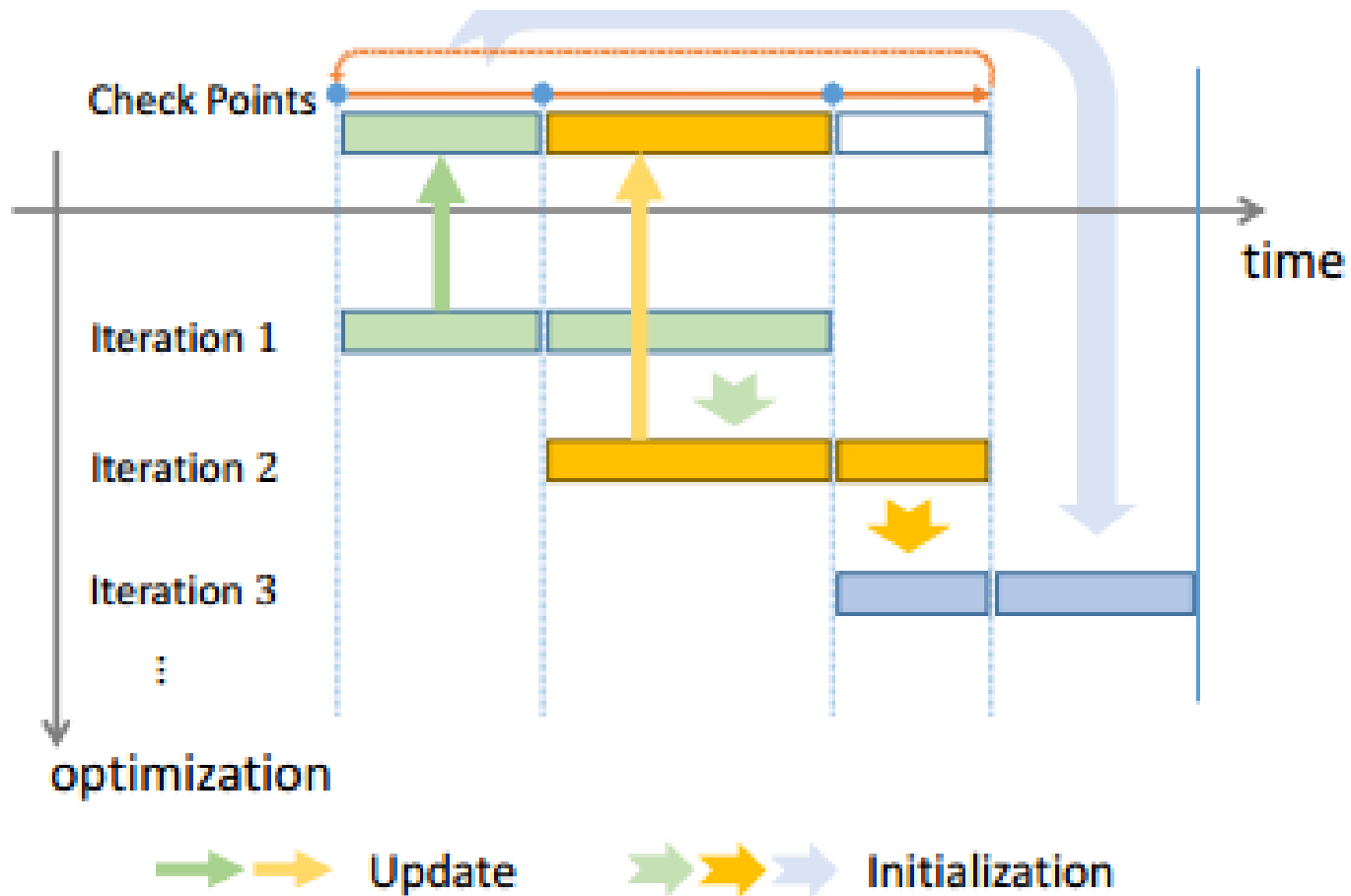
- Pre komplikované zručnosti, dva takého body
- Jeden na začiatku a druhý na konci otočky

# Optimalizačné premenné

- množina  $\mathcal{H} = \{H_t\}$ ,  $H_t \in \{\emptyset, \{L\}, \{R\}, \{L, R\}\}$
- $\{\chi_k^L, \chi_k^H\}$  Použijeme ich na opravu „open-loop“ segmentov na ovládanie paží
- $\chi_k^L$  and  $\chi_k^H$  Reprezentujú použité korekcie v tomto poradí: pre ľavú a pravú ruku

$$\chi = (q_{\text{shoulder}}, q_{\text{elbow}}, q_{\text{wrist}}, \alpha_{\text{fingers}})$$

# Optimalizačné premenné 2



# Učenie lineárnej stratégie kontroly

- Stratégia kontroly pre paže
- väzme stavový vektor  $x$  ako vstup
- vypočíta akčný vektor  $a = (\chi^L, \chi^R)$
- Stavový vektor sa skladá zo súčasného simulačného stavu  $s$  a indexu kontrolného fragmentu  $k$

# Učenie lineárnej stratégie kontroly 2

- $s$  zachytáva stav hráča a lopty ako vektor:

$$s = \{p_{\text{ball}}, \dot{p}_{\text{ball}}, p_j, \dot{p}_j, d_i, c, \dot{c}, L\}$$



# Lineárna regresia

- Na naučenie sa postupnej lineárnej stratégie kontroly budeme používať lineárnu regresiu

$$\pi(s, k) = M_k s + \hat{a}_k,$$

- pomocou regresie nevieme kontrolovať zložitejšie driblingové triky
- lineárnu stratégiu kontroly budeme používať na inicializáciu hlbokého procesu učenia.

# Informácia o dotyku

- náš systém počíta priemernú polohu lopty v každom kontrolnom fragmente
- aktualizuje informácie o dotyku
- tak aby vzdialenosť medzi rukou a loptou bola menšia ako 5 cm

# Hlboké zosilnené učenie

- Používame ho aby sme dosiahli robustné učenie basketbalových zručností.

---

**ALGORITHM 1:** Learn Arm Control Policy Using DDPG

---

**Input:** control fragments  $\{C_k\}$ ,  $k = 1, \dots, K$  and  
associated linear control polices  $\{(M_k, \hat{a}_k)\}$

**Input:** starting states  $\{s_{k_i}\}$ ,  $k_i \in \{1, \dots, K\}$

**Result:** arm control policy  $\pi$

initialize  $\mathcal{D} \leftarrow \emptyset$

initialize critic network parameters  $\theta_Q$  and actor network parameters  
 $\theta_\pi$

initialize target function  $\theta'_Q \leftarrow \theta_Q$ ,  $\theta'_\pi \leftarrow \theta_\pi$

initialize simulation with random starting state  $s_{k_i}$ , set  $\mathbf{x} \leftarrow (s_{k_i}, k_i)$

**for**  $t \leftarrow 1, 2, \dots$  **do**

$k \leftarrow k(\mathbf{x});$  // get the 'k' component of  $\mathbf{x}$

**if**  $t \leq N_{init}$  **then**

$s \leftarrow s(\mathbf{x});$  // get the 's' component of  $\mathbf{x}$

        compute action  $\mathbf{a} = M_k s + \hat{a}_k + \epsilon_t$  and  $r = r(\mathbf{x}, \mathbf{a})$

**else**

        compute action  $\mathbf{a} = \pi(\mathbf{x}; \theta_\pi) + \epsilon_t$  and  $r = r(\mathbf{x}, \mathbf{a})$

**end**

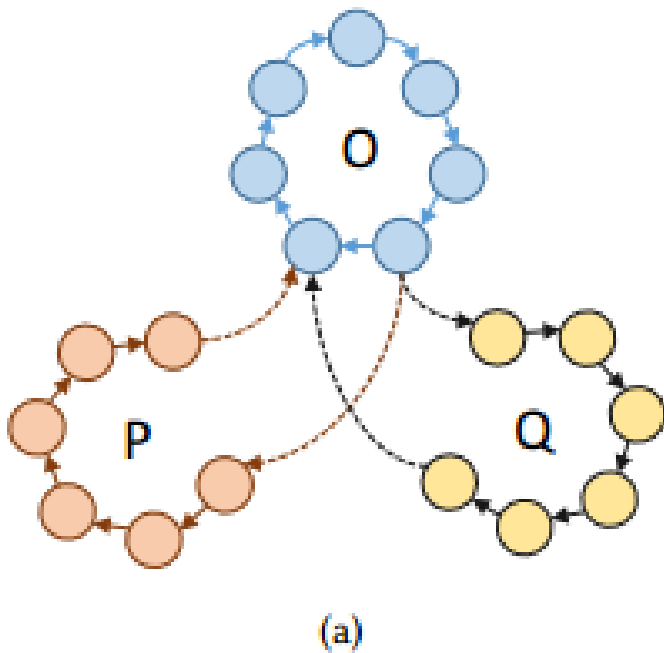
# Hlboké zosilnené učenie 2

```
13 execute control fragment  $C_k$  and observe  $\mathbf{x}' \leftarrow (s', \text{next}(k))$ 
14 store transition  $\tau = (\mathbf{x}, \mathbf{a}, r, \mathbf{x}')$  in  $\mathcal{D}$ 
15 if  $t > N_{init}$  then
16     sample a minibatch of  $N_{batch}$  transitions  $\{\tau_i\} \subset \mathcal{D}$ 
17     compute target values  $y_i = y(\tau_i; \theta'_Q, \theta'_\pi)$  using Equation 13
18     update  $\theta_Q$  by minimizing the loss function of Equation 12
19     update  $\theta_\pi$  using the policy gradient of Equation 16
20      $\theta'_Q \leftarrow (1 - \eta)\theta'_Q + \eta\theta_Q$ 
21      $\theta'_\pi \leftarrow (1 - \eta)\theta'_\pi + \eta\theta_\pi$ 
22 end
23 if  $\mathbf{x}' \in \mathcal{X}_{term}$  then
24     reset simulation to a random starting state  $s_{k_i}$ 
25      $\mathbf{x} \leftarrow (s, k_i)$ 
26 else
27      $\mathbf{x} \leftarrow \mathbf{x}'$ 
28 end
29 end
```

---

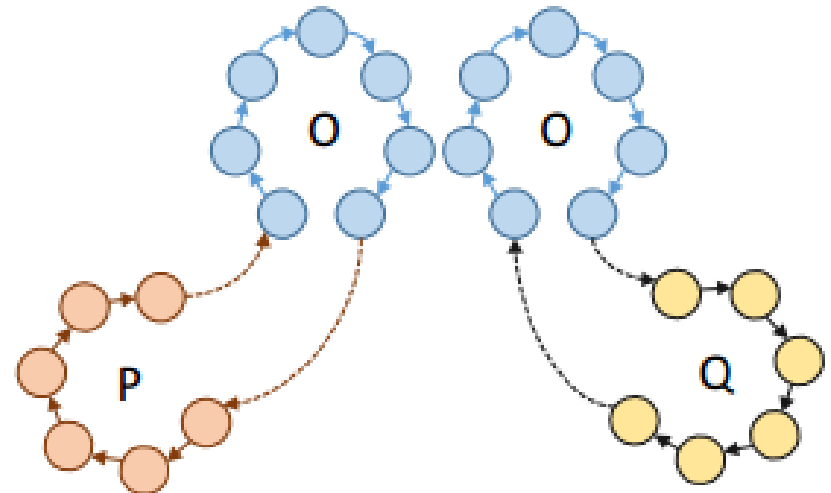
# Učenie kontrolných grafov

- Je to graf, ktorého uzly sú kontrolné fragmenty



# Inkrementálne učenie

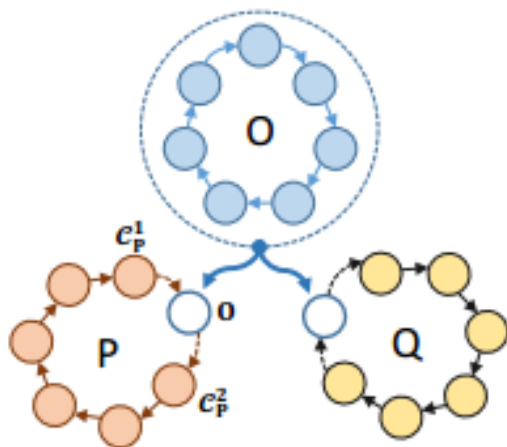
- systém sa učí kontrolu paží pre kontrolný graf inkrementálne
- Ak trénujeme necyklický trik P
- systém skombinuje triky P a O do cyklického pohybu (O+P)



(b)

# Inkrementálne učenie 2

- Pri hlbokom zosilnenom učení najskôr natrénujeme nelineárnu korekčnú kompenzáciu pre trik O
- použijeme kombinovaný trik (O+P) na natréňovanie necyklického triku P



(c)

# Výsledky

- Náš hráč má 1.8 m váži 76 kg
- basketbalová lopta má rádius 11,93 cm a váži 623.7 kg
- Systém vie bežať rýchlejšie ako beží reálny čas.



# Cyklické zručnosti



(A)



(B)



(C)



(D)



(E)

# Optimalizácia trajektórie

- optimalizujeme pohybovú sekvenciu, v ktorej hráč opakovane vykonáva zručnosti
- *Nrep = 100 cyklov.*
- Distribúcia vzoriek algoritmu CMA-ES (Covariance Matrix Adaption Evolution Strategy) je inicializovaná ako normálne rozdelenie  $\mathcal{N}(0, \sigma_0^2)$ ,
- $\sigma_0 = 0.03$  v prvom cykle
- $\sigma_0 = 0.01$  v ďalších cykloch

# Optimalizácia trajektórie 2

- Končí ak
  - a) počet iterácii je väčší ako 1000
  - b) proces optimalizácie sa zastaví na 200 iterácii
  - c) priemerná vzdialenosť medzi hráčovú rukou a loptou menšia ako 1 cm pre (A – D) respektíve 2 cm pre E

# Lineárna stratégia kontroly.

- postup postupnej lineárnej stratégie kontroly paží sa naučí z optimalizovaných pohybových sekvencií pomocou lin. regresie
- hráč môže opakovane mávať rukami stovky krát pri prevedení zručnosti (A)
- Problém - nevieme vykonávať kontrolu nad komplexnejšími zručnosťami (B – E)

# Hlboké zosilnenie učenia

- nelineárna stratégia kontroly paže s použitím DDPG algoritmu
- Táto nelineárne stratégia kontroly paží umožňuje kontrolu pre triky B - E

# Kontrolné grafy

- dva kontrolné grafy
- každý obsahuje jeden cyklický a dva necyklické triky

Table 2. Robustness of noncyclic skills in the learned control graphs.

control graph	noncyclic skill	success rate
in-place skills	dribbling between the legs	99.4%
	dribbling behind the back	99.6%
running skills	front crossover	98.1%
	spin crossover	95.7%

# Robustnosť trikov

- nechali sme hráča opakovane vykonávať cyklickú aj necyklickú zručnosť
- pravdepodobnosť prechodu z cyklického na necyklický trik bola 0.5

Ďakujem za pozornosť