# Learning a Contact-Adaptive Controller for Robust, Efficient Legged Locomotion

Hansol Kang, Robotics Innovatory, SKKU
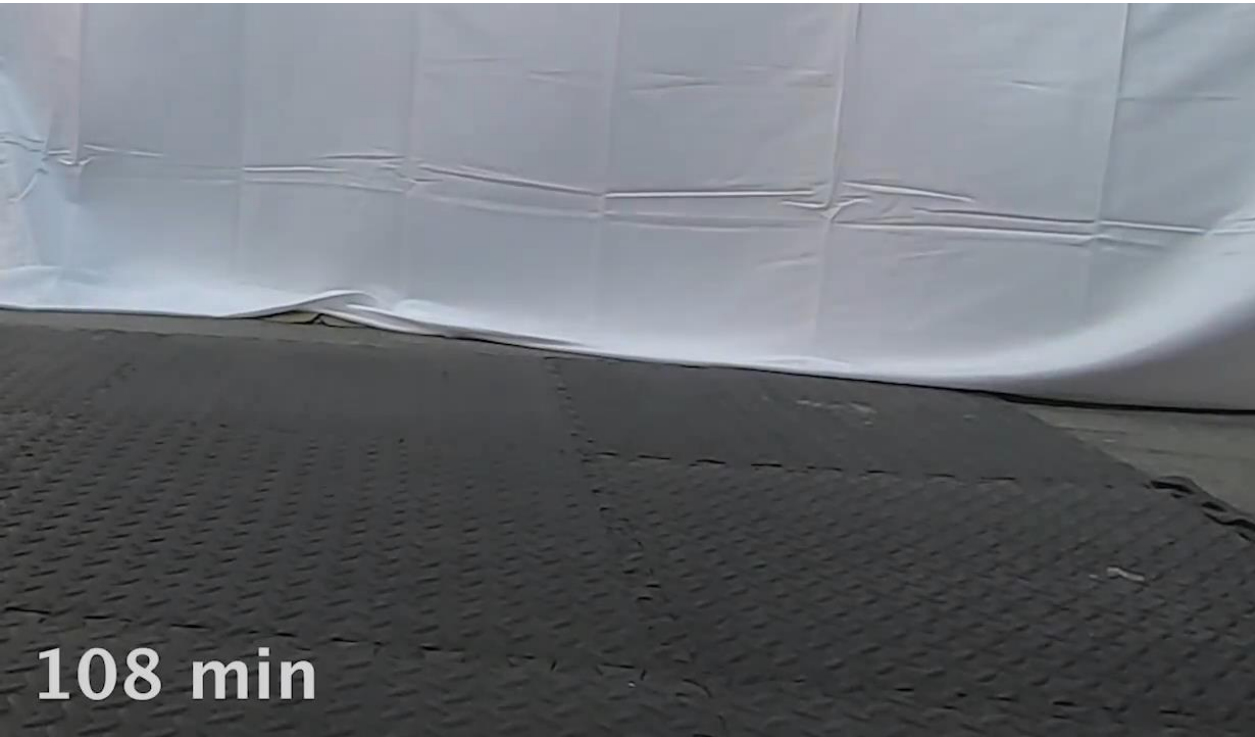
SKKU UNIVERSITY

Robotics Innovatory
http://mecha.skku.ac.kr

# Contents

- Introduction

- Paper Review

- Implementation

# Introduction

▶ Robots controlled through RL



108 min

Through co-training, MELA learns adaptive skills across various locomotion modes

Learning to Walk via Deep Reinforcement Learning. Tuomas Haarnoja, Sehoon Ha, Aurick Zhou, Jie Tan, George Tucker, Sergey Levine. Robotics: Science and Systems (RSS). 2019.
Multi-expert learning of adaptive legged locomotion, Chuanyu Yang, Kai Yuan, Qiuguo Zhu, Wanming Yu, Zhibin LiScience Robotics, 2020

▶ Learning a Contact-Adaptive Controller for Robust, Efficient Legged Locomotion



▶ 2020 Conference on Robot Learning(CoRL)

▶ 2020.11. arXiv

▶ Xingye Da[1], Zhaoming Xie[1,2], David Hoeller[1], Byron Boots[1,3], Anima Anandkumar[1,4], Yuke Zhu[1,5], Buck Babich[1], Animech Garg[1,6]

▶ [1]NVIDIA, [2]Univ. of British Columbia, [3]Univ. of Washington, [4]Calthec, [5]UT Austin, [6]Univ. of Toronto

▶ Introduce a hierarchical control structure that combines model-based control design and model-free reinforcement learning for legged locomotion.

▶ How to control a robot?



Curr. Pos.

Des. Pos.

$\theta$

1 DoF* link

Curr. Pos.

Obstacle

?

Des. Pos.

1 DoF link

* Degree of Freedom

# Paper Review

▶ Quadratic Programming(QP)

QP is an optimization tool

$$minimize\ 2x_1^2 + x_1x_2 + x_2^2 + x_1 + x_2$$

$$subject\ to\ x_1 + x_2 = 1$$

$$0 \leq x_1 \leq 0.7$$

$$0 \leq x_2 \leq 0.7$$

$$minimize\ \frac{1}{2}\mathbf{x}^T\mathbf{Q}\mathbf{x} + \mathbf{c}^T\mathbf{x}$$

$$subject\ to\ \mathbf{Ax} \leq \mathbf{b}$$

$$\text{minimize} \quad \frac{1}{2}x^T \begin{bmatrix} 4 & 1 \\ 1 & 2 \end{bmatrix} x + \begin{bmatrix} 1 \\ 1 \end{bmatrix}^T x$$

$$\text{subject to} \quad \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \leq \begin{bmatrix} 1 & 1 \\ 1 & 0 \\ 0 & 1 \end{bmatrix} x \leq \begin{bmatrix} 1 \\ 0.7 \\ 0.7 \end{bmatrix}$$

▶ How to use QP for quadruped robot control?

$$\underbrace{\begin{bmatrix} I & \cdots & I \\ [p_{com,1\times}] & \cdots & [p_{com,c\times}] \end{bmatrix}}_{A} \underbrace{\begin{bmatrix} f_1 \\ \cdots \\ f_c \end{bmatrix}}_{f} = \underbrace{\begin{bmatrix} m(\ddot{x}^d_{com} + g) \\ I_g \dot{\omega}^d_b \end{bmatrix}}_{b}$$

$$f^d = \arg \min_{f \in \mathbb{R}^k} (Af - b)^T S (Af - b) + \alpha f^T W f$$

$$s.t. \underline{d} < Cf < \overline{d}$$

$$S \in \mathbb{R}^{6\times6} \text{ and } W \in \mathbb{R}^{k\times k} \text{ are positive} - \text{definite weight matrices}$$

$$\underline{d}, \overline{d} \in \mathbb{R}^p, C \in \mathbb{R}^{p\times k}$$

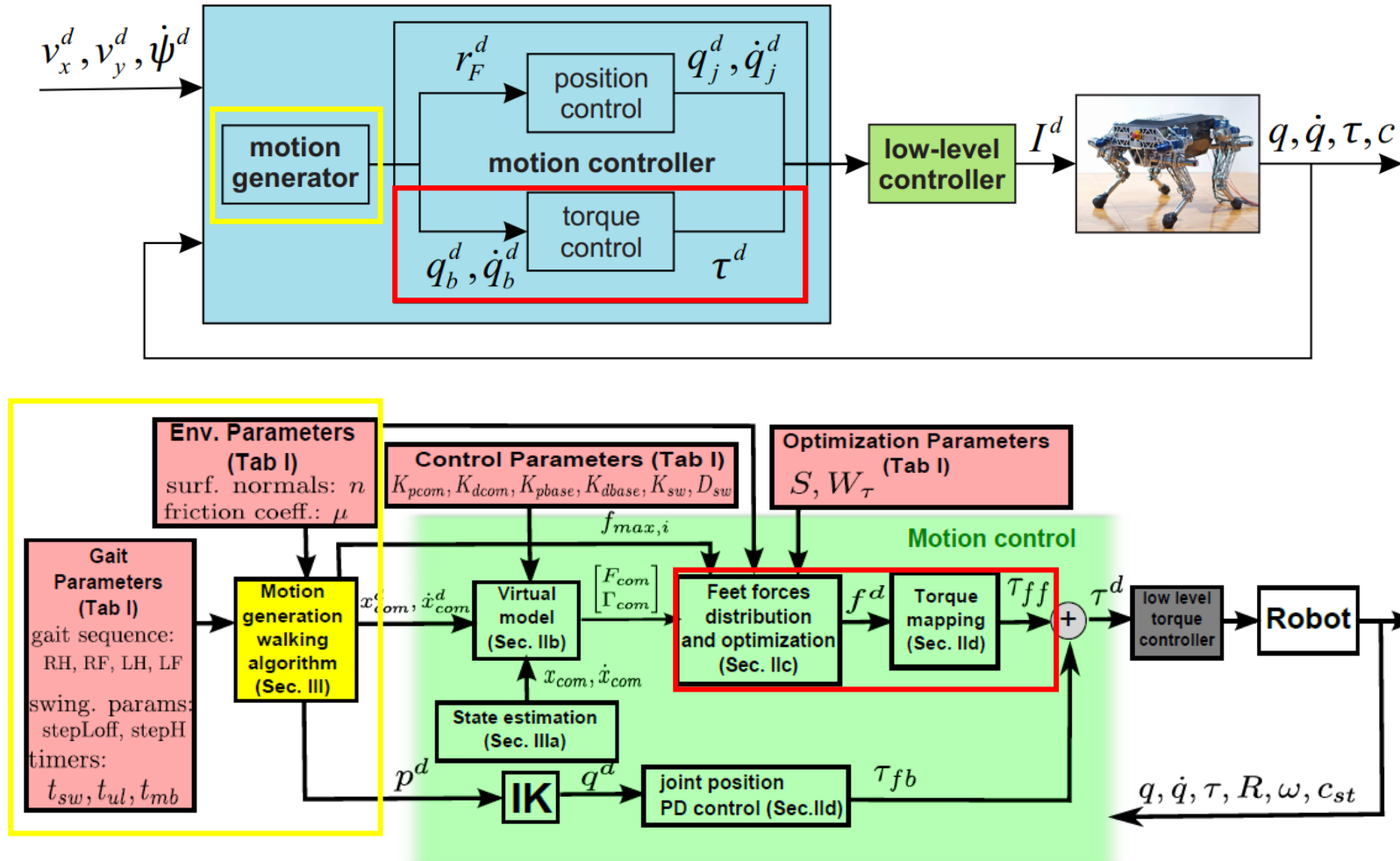$$k = 3 \times Num.of\ contact, p = Num.of\ inequality\ const.$$

$$\alpha \in \mathbb{R}\ weighs\ the\ secondary\ objective.$$

# Paper Review

▶ How to use QP for quadruped robot control?

Control of Dynamic Gaits for a Quadrupedal Robot., Christian Gehring, Stelian Coros, Marco Hutter, Michael Bloesch, Markus A. Hoepflinger, Roland Siegwart, International Conference on Robotics and Automation(ICRA),2013
High-slope Terrain Locomotion for Torque-Controlled Quaruped Robots, Michele Focchi, Andrea del Prete, Ioannis Havoutis, Roy Featherstone, Darwin Caldwell, Claudio Semini, Autonomous Robots, 2017

▶ Controller overview

▶ Primitive set
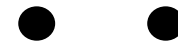
Stand  Trot 1  Trot 2  Pace 1  Pace 2

Step 1  Step 2  Step 3  Step 4

▶ High level controller

Simulator : Issac Gym

State Space
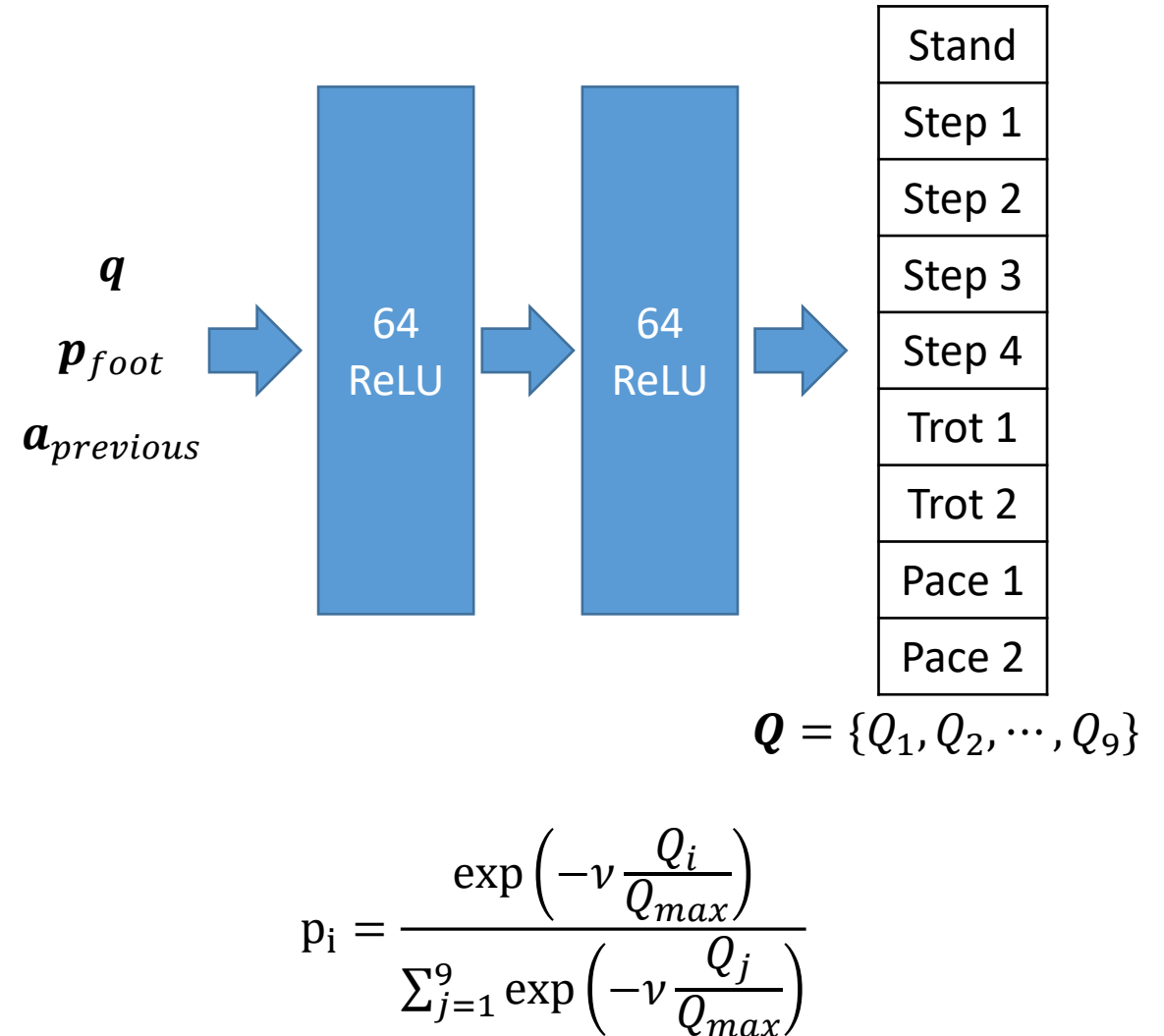- Body pose $q$ (exclude x, y linear position)
- Relative foot position $p_{foot}$
- Previously-used primitive $a_{previous}$

Action Space
- 9 Primitives

Reward

$$r = 1 - 0.0025\frac{1}{T}\Sigma\|\tau\|^2 - \frac{1}{T}\Sigma\|\dot{p}_{d,body} - \dot{p}_{body}\|^2$$

$q$

$p_{foot}$

$a_{previous}$

64 ReLU

64 ReLU

| Stand |
| Step 1 |
| Step 2 |
| Step 3 |
| Step 4 |
| Trot 1 |
| Trot 2 |
| Pace 1 |
| Pace 2 |

$$Q = \{Q_1, Q_2, \cdots, Q_9\}$$

$$p_i = \frac{\exp\left(-\nu\frac{Q_i}{Q_{max}}\right)}{\sum_{j=1}^9 \exp\left(-\nu\frac{Q_j}{Q_{max}}\right)}$$

▶ High level controller

## B Q-Learning Algorithm

We use DQN like algorithm to train our high-level policy. Details are shown in Algorithm 1.

---

**Algorithm 1:** Q Learning

---

initialization Q-function parameters $\theta_1.\theta_2$ for $Q_{\theta_1}, Q_{\theta_2}$, empty replay buffer $D$ ;
set target network parameters $\theta_{targ,1}, \theta_{targ,2} \leftarrow \theta_1, \theta_2$ for $Q_{\theta_{targ,1}}, Q_{\theta_{targ,2}}$ ;
**while** *not done* **do**

    observe current state $s$ ;
    sample action $a$ based on Q-function;
    observe next state $s'$, reward $r$ and done signal $d$;
    store $(s, a, r, d, s')$ in replay buffer $D$;
    **if** *d is True or time limit reached* **then**
        reset environment;
    **end**
    **if** *time to update* **then**
        **for** $j = 1, 2, \ldots$ *number of update* **do**
            sample batch of transition data $B = \{s, a, r, d, s'\}$;
            compute $a' = \arg\max_a Q_\theta(s', a)$;
            compute target $q_{targ} = r + (1 - d)\gamma \min_{i=1,2}(Q_{\theta_{targ,i}}(s', a'))$;
            update $\theta_1, \theta_2$ by taking gradient descent w.r.t the objective function
            $\frac{1}{|B|} \sum_{(s,a,r,d,s')\in B}((Q_{\theta_1}(s, a) - q_{targ})^2 + (Q_{\theta_2}(s, a) - q_{targ})^2)$ ;
            **if** $j \mod 2 = 1$ **then**
                $\theta_{targ,1} \leftarrow \rho\theta_{targ,1} + (1 - \rho)\theta_1$ ;
                $\theta_{targ,2} \leftarrow \rho\theta_{targ,2} + (1 - \rho)\theta_2$ ;
            **end**
        **end**
    **end**
**end**

▶ Low level controller

Base Pose Control

$$\min_{f} ||\mathbf{M}f - \tilde{g} - \ddot{q}_d||Q + ||f||R$$

$$subject\ to\ f_{z,i} \geq f_{z,min}\ if\ P_{t,i}\ is\ Stance$$
$$f_{z,i} = 0\ if\ P_{t,i}\ is\ Swing$$
$$-\mu f_x \leq f_z \leq \mu f_x$$
$$-\mu f_y \leq f_z \leq \mu f_y$$

Swing Foot Control

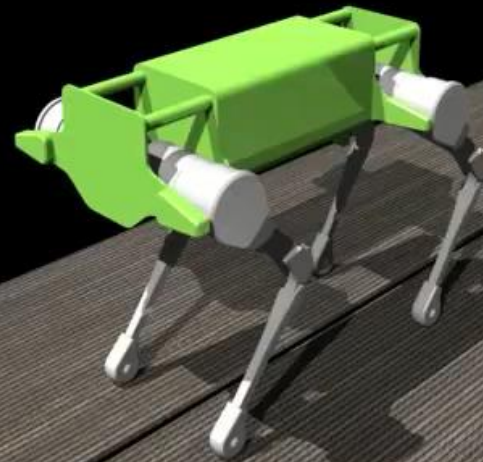$$p_{d,i} = p_{0,i} + k(\dot{p}_{body} - \dot{p}_{d,body})$$
$$f_i = k_{p,i}(p_{d,i} - p_i) - k_{d,i}\dot{p}_i$$

▶ Result

▶ Result

▶ Result



(a) Energy comparison over different speeds.

(b) Energy comparison over different yaws.

▶ Zero-shot adaptation

▶ Zero-shot adaptation

▶ Sim-to-real test

# Implementation

▶ Environment

PyBullet, PyTorch, Gym

# Q & A

**Thank you for your attention**