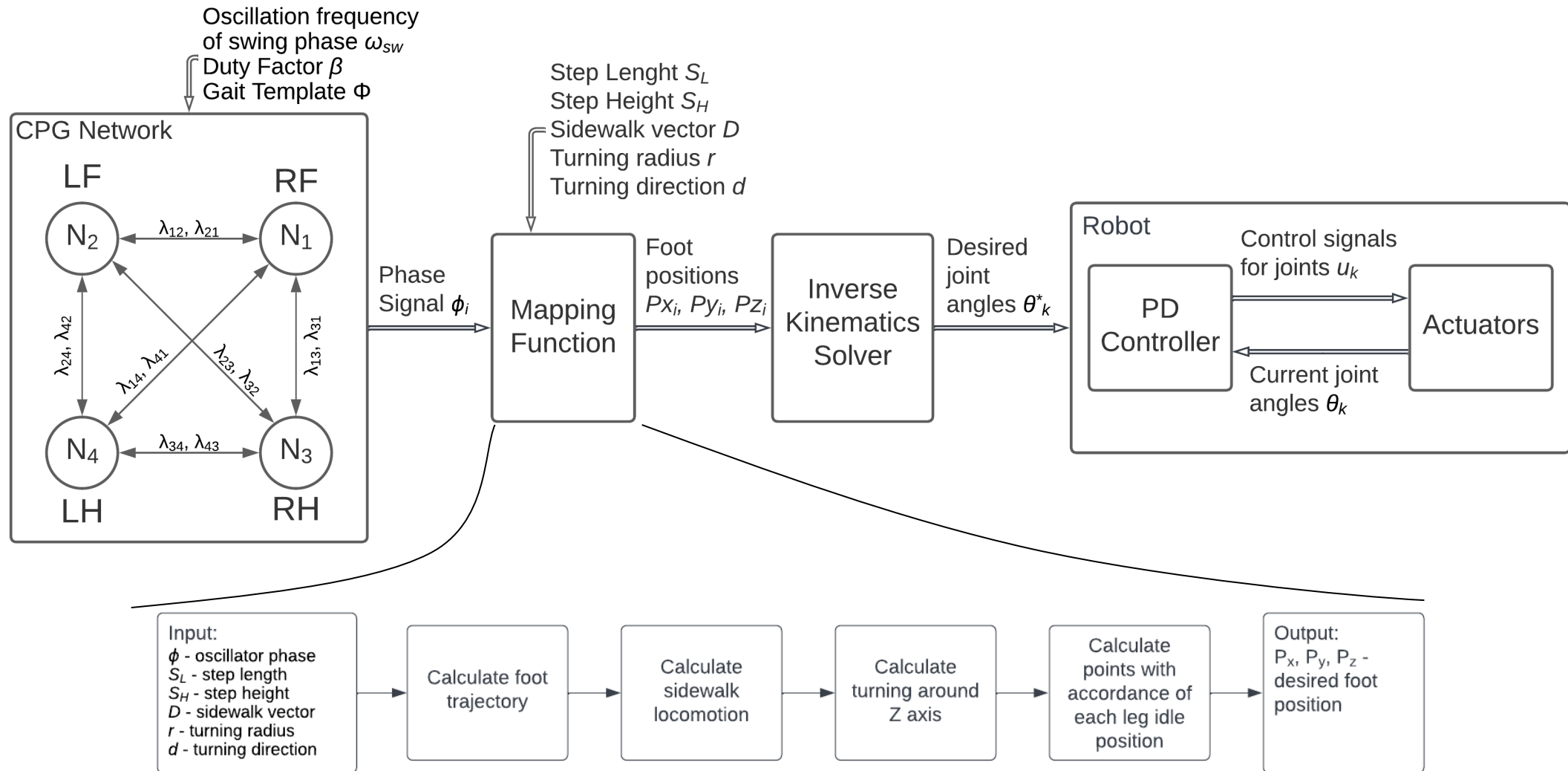


# Control of Simplified Walking Robot Model using PMTG architecture

---

Vladimir Danilov

# Block diagram of CPG-based Gait Generator



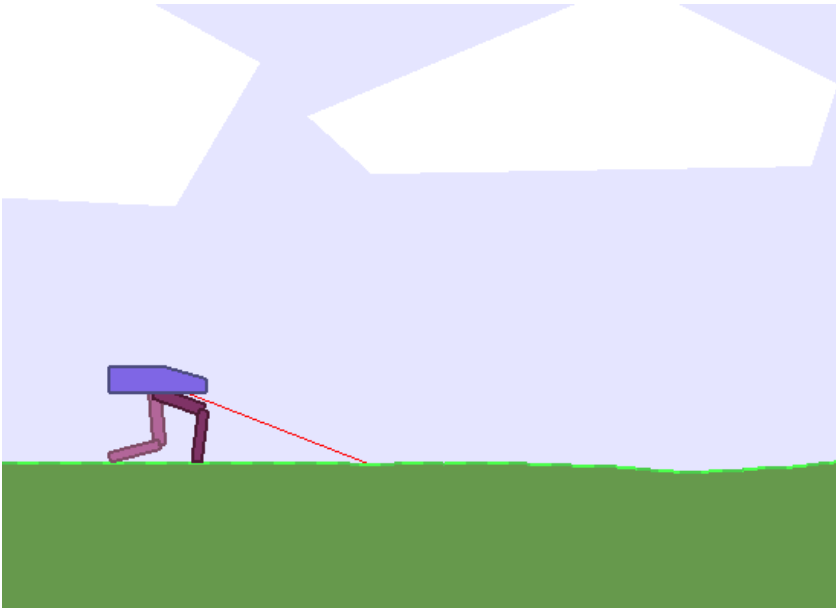
## Experimental Hardware Results



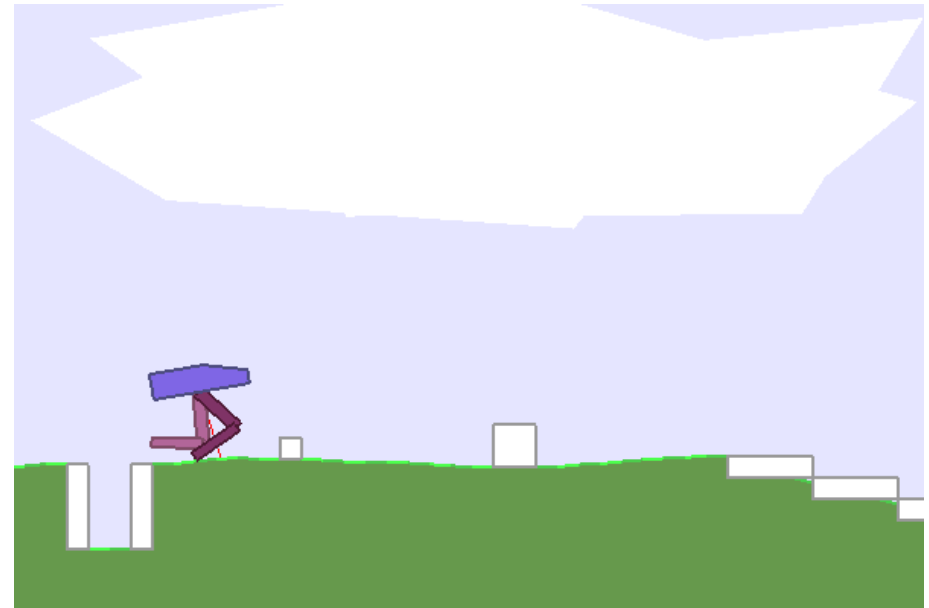
[Watch on Youtube](#)

# Simplified Walking Robot Model

Bipedal Walker



Bipedal Walker Hardcore

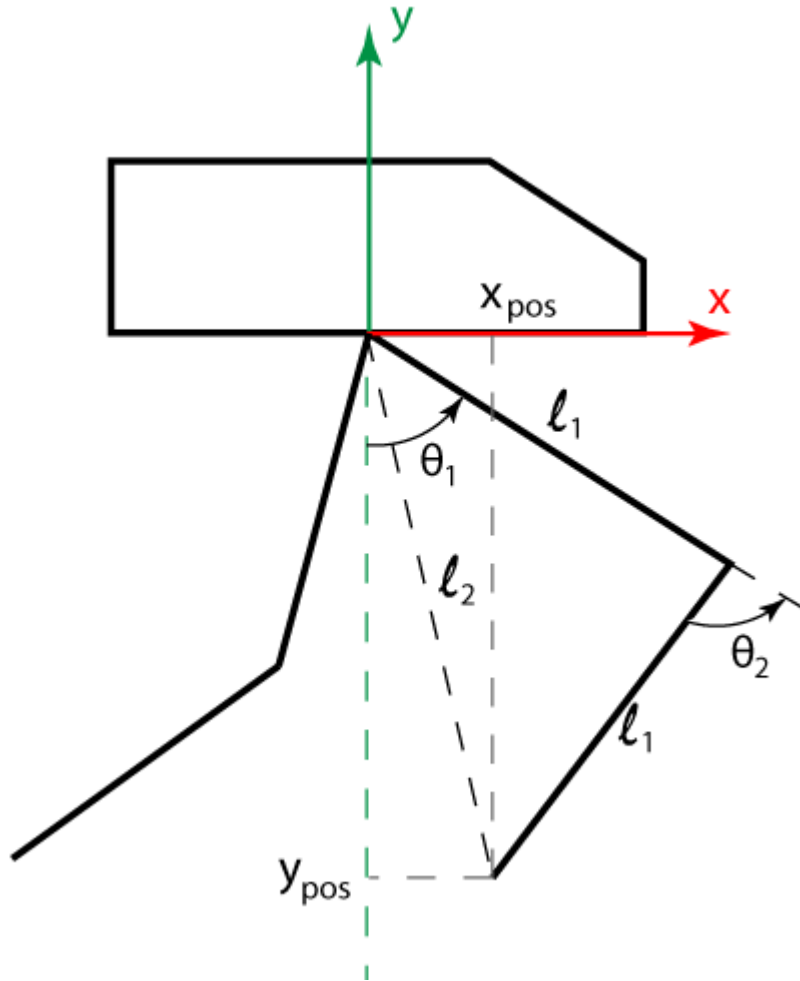


**Starting State:** Random position upright and mostly straight legs

**Episode Termination:** when the robot body touches ground or the robot reaches far right side of the environment

**Solved Requirements:** to get average reward greater than 300

# Kinematics



Forward Kinematics:

$$l_2 = 2l_1^2(1 + \sin \theta_2)$$

$$x_{pos} = l_2 \sin \left( \theta_1 + \frac{\theta_2}{2} \right)$$

$$y_{pos} = \sqrt{l_2^2 - x_{pos}^2}$$

Inverse Kinematics:

$$l_2 = \sqrt{x_{pos}^2 + y_{pos}^2}$$

$$\theta_1 = \text{atan} \frac{x_{pos}}{y_{pos}} + \text{acos} \frac{l_2}{2l_1}$$

$$\theta_2 = -2 \text{acos} \frac{l_2}{2l_1}$$

# Environment

Num	Observation	Min	Max
0	hull_angle	0	$2\pi$
1	hull_angularVelocity	$-\infty$	$+\infty$
2	vel_x	-1	+1
3	vel_y	-1	+1
4	hip_joint_1_angle	$-\infty$	$+\infty$
5	hip_joint_1_speed	$-\infty$	$+\infty$
6	knee_joint_1_angle	$-\infty$	$+\infty$
7	knee_joint_1_speed	$-\infty$	$+\infty$
8	leg_1_ground_contact_flag	0	1
9	hip_joint_2_angle	$-\infty$	$+\infty$
10	hip_joint_2_speed	$-\infty$	$+\infty$
11	knee_joint_2_angle	$-\infty$	$+\infty$
12	knee_joint_2_speed	$-\infty$	$+\infty$
13	leg_2_ground_contact_flag	0	1
14-23	10 lidar readings	$-\infty$	$+\infty$

Num	Action	Min	Max
0	Hip_1 (Torque / Velocity)	-1	+1
1	Knee_1 (Torque / Velocity)	-1	+1
2	Hip_2 (Torque / Velocity)	-1	+1
3	Knee_2 (Torque / Velocity)	-1	+1

Reward function:

$$r_{fw} = \frac{13}{3}(p_x(t) - p_x(t-1)) - \text{moving forward reward}$$

$$r_{hull} = -5(|\vartheta(t)| - |\vartheta(t-1)|) - \text{hull deviation penalty}$$

$$r_{\tau} = -0,028 \sum_{i=1}^{12} |a_i| - \text{torque penalty}$$

$$r_{es} = \begin{cases} -100, & \text{if } p_x < 0 \text{ or hull touches ground} \\ 0, & \text{otherwise} \end{cases} - \text{early stopping penalty}$$

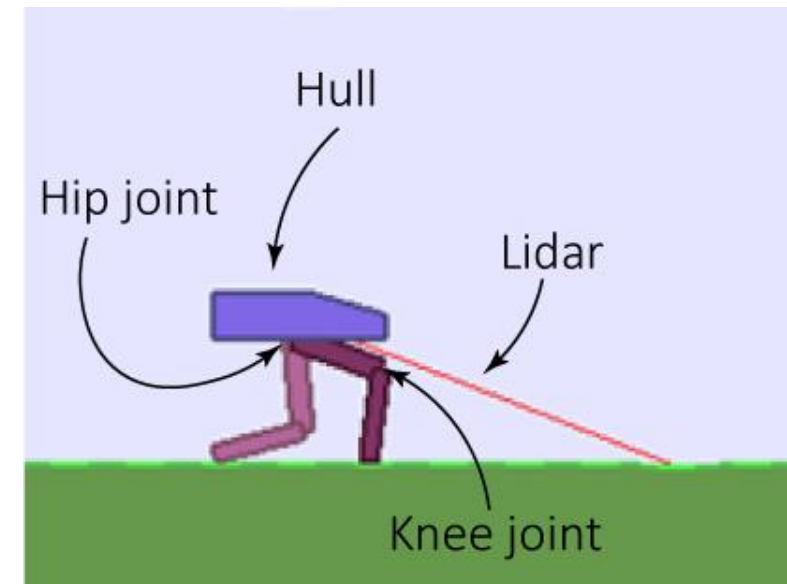
$$R = r_{fw} + r_{hull} + r_{\tau} + r_{es} - \text{total reward},$$

where:

$p_x$  - x position,

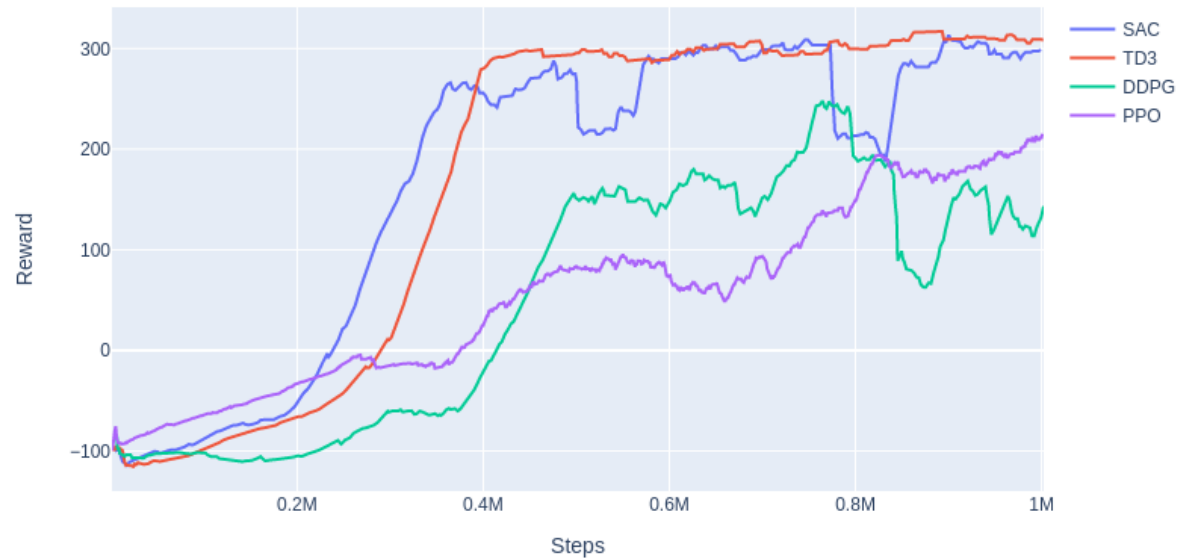
$\vartheta$  - hull angle

$a_i$  - applied action

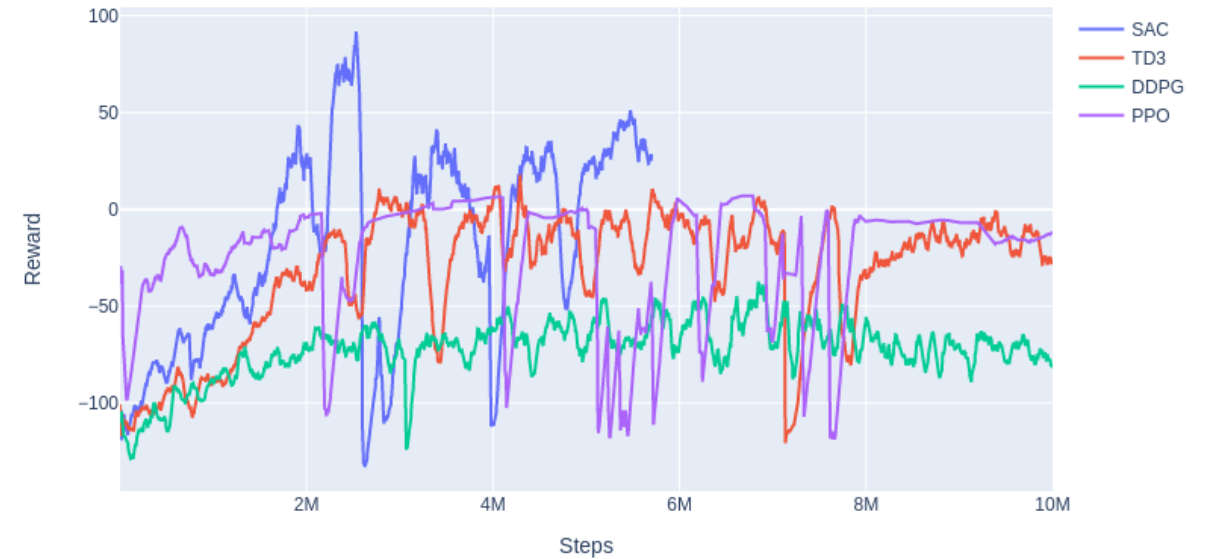


# Learning with Vanilla RL Algorithms

## Bipedal Walker

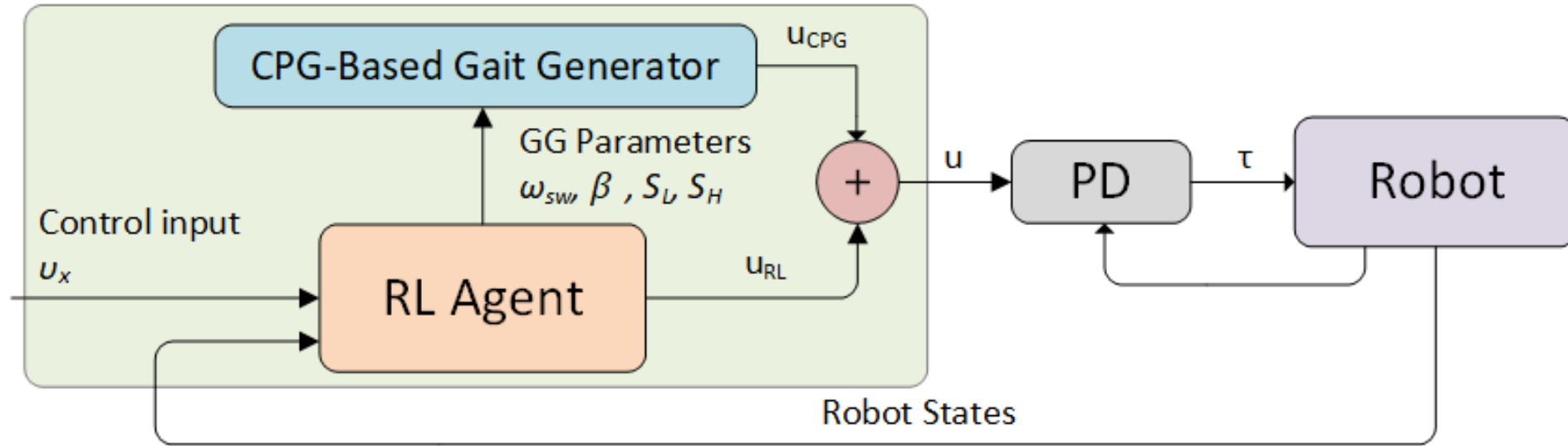


## Bipedal Walker Hardcore



\*all the hyperparameters were optimized by Tree-structured Parzen Estimator algorithm

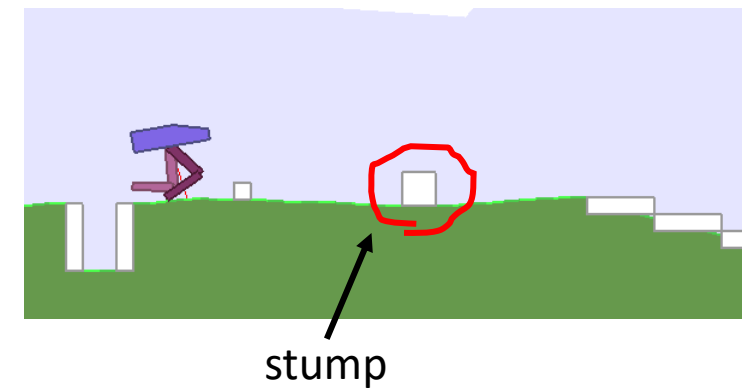
# Policy Modulating Trajectory Generator (PMTG) Architecture



Supposed agents: TD3, SAC

Tricks to enhance performance:

- During training the actions are repeated for three steps.
- When the agent falls, the terminal reward is clipped to zero.
- The reward is scaled by 5.
- Noises are added to the actions.
- When training, the probability of encountering a stump is increased.

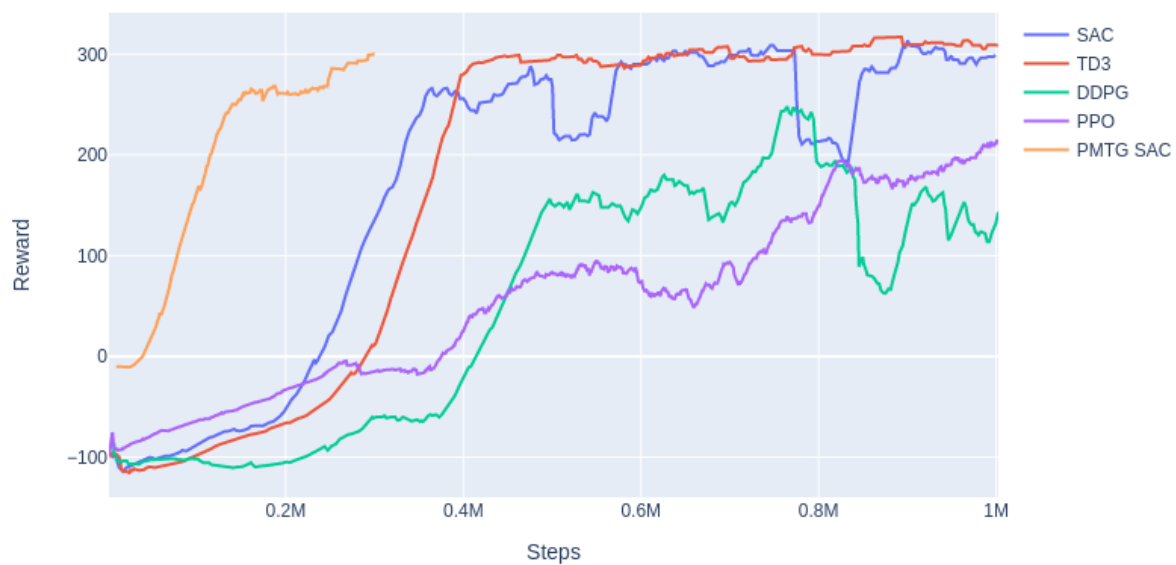




# Learning with PMTG

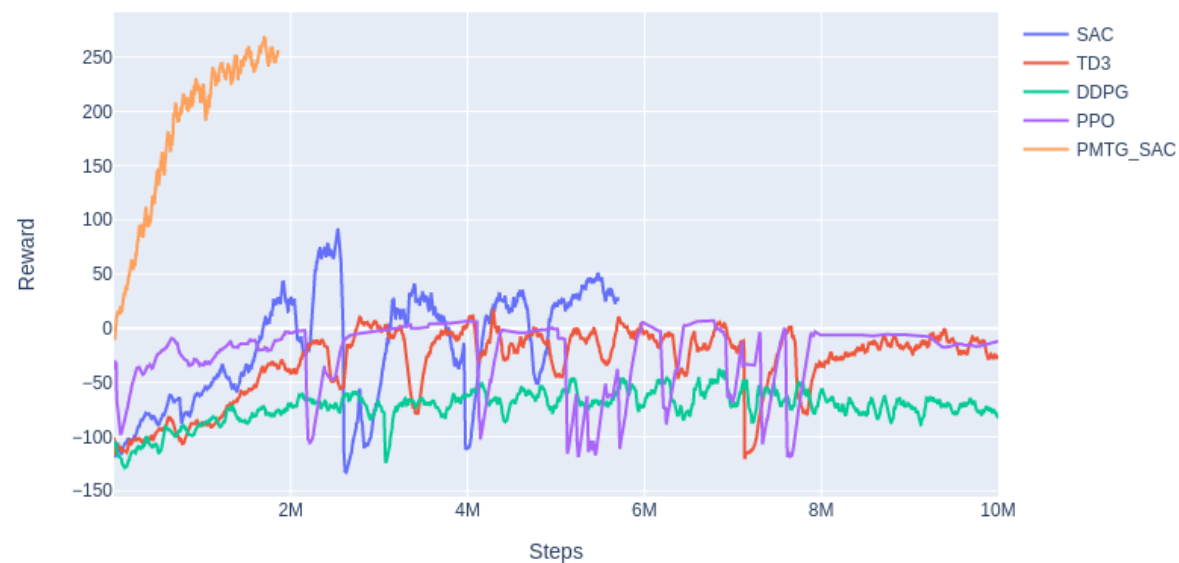
Подумать, как наглядно обозначить то, что в итоге

## Bipedal Walker



Learned agent 100-episode average score: 304.24  
Episodes before solve: 696

## Bipedal Walker Hardcore



Learned agent 100-episode average score: 302.92  
Episodes before solve: 7280

# OpenAi Leaderboard

## Bipedal Walker

User	Episodes before solve	Write-up	Video
<a href="#">Nandino Cakar</a>	474	<a href="#">writeup</a>	
<a href="#">ZhiqingXiao</a>	0 (use close-form preset policy)	<a href="#">writeup</a>	
<a href="#">Nick Kaparinos</a>	800	<a href="#">Write-up</a>	<a href="#">gif</a>
<a href="#">shnippi</a>	925	<a href="#">writeup</a>	

My learned agent 100-episode average score: 304.24  
Episodes before solve: 696

## Bipedal Walker Hardcore

User	Episodes before solve	100-Episode Average Score	Write-up	Video
<a href="#">Nick Kaparinos</a>	15500	305.40 $\pm$ 21.35	<a href="#">Write-up</a>	<a href="#">gif</a>
<a href="#">Alister Maguire</a>	N/A	313	<a href="#">Write-up</a>	<a href="#">gif</a>

My learned agent 100-episode average score: 302.92  
Episodes before solve: 7280

# Learned Agent Video

Bipedal Walker



[Watch on Youtube](#)

Bipedal Walker Hardcore



[Watch on Youtube](#)

# Bipedal Walker Hardcore with Velocity Control

Reward function:

$$r_{fw} = \frac{4}{30} \exp\left(\frac{(v^* - v)^2}{0.2}\right) - \text{moving forward reward}$$

$$r_{hull} = -5(|\vartheta(t)| - |\vartheta(t-1)|) - \text{hull deviation penalty}$$

$$r_{\tau} = -0,056 \sum_{i=1}^{12} |\tau_i| - \text{torque penalty}$$

$$R = r_{fw} + r_{hull} + r_{\tau} - \text{total reward},$$

where:

$p_x$  – x position,

$\vartheta$  – hull angle

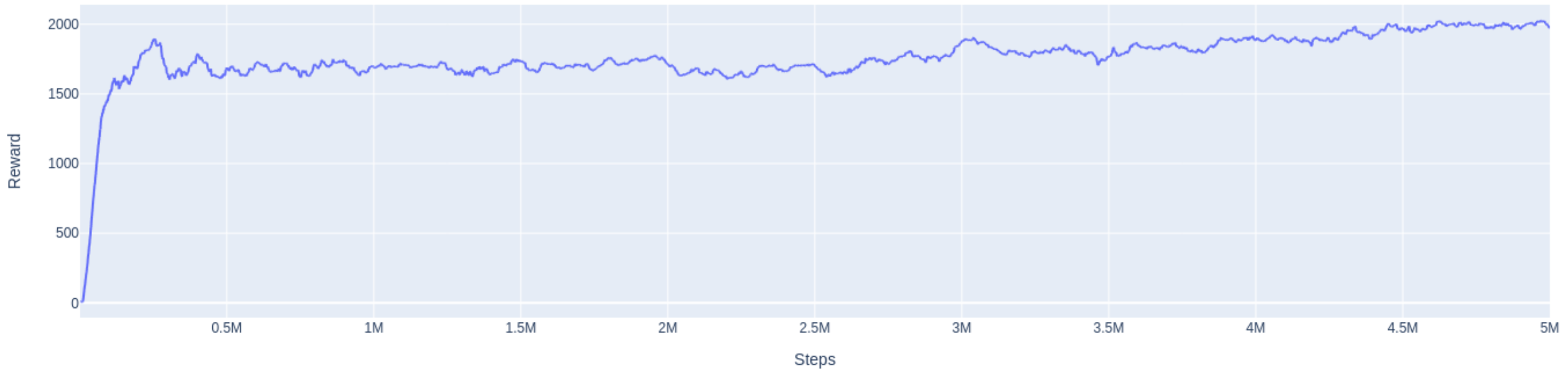
$a_i$  – applied action

Desired velocity:

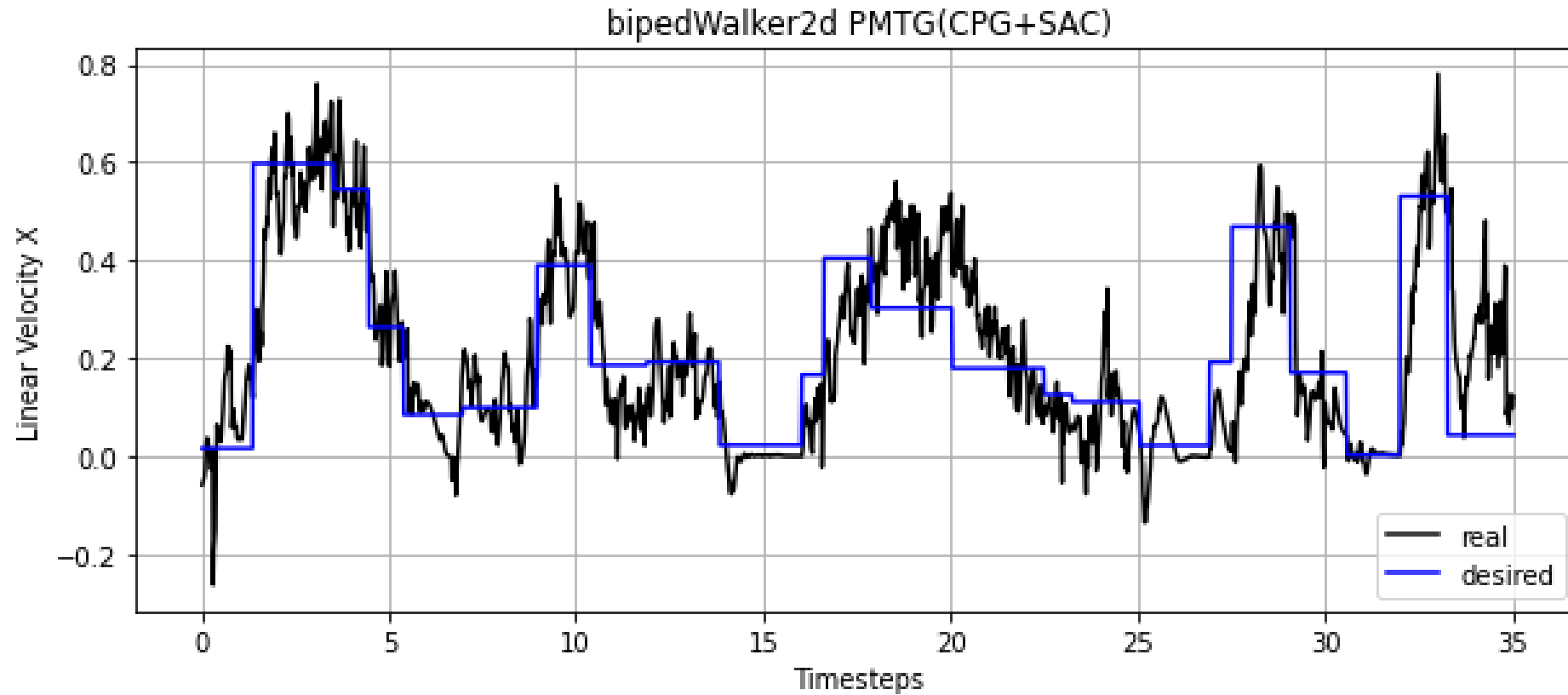
$v^* \sim U[0, 0.6]$  – desired velocity distribution

$t_{sw} \sim U[1, 5]$  – time before changing desired velocity

# Learning Curve



# Learned Agent with Velocity Control



## Learned Agent with Velocity Control Video



[Watch on Youtube](#)

# Thanks for Your Attention!

**Control of Simplified Walking Robot Model using PMTG architecture**

---

Vladimir Danilov