

Robotics / Intelligent Robotics Robot Learning

Luís Paulo Reis, Nuno Lau, David Simões, Armando Sousa

lpreat@fe.up.pt

Director of LIACC – Artificial Intelligence and Computer Science Lab.

Associate Professor at DEI/FEUP – Informatics Engineering Department, Faculty of Engineering
of the University of Porto, Portugal

President of APPIA – Portuguese Association for Artificial Intelligence



Robot Learning - Motivation

Programming Robots is a hard task

- No high-level programming language
- Sensors and actuators are noisy
- Robotics is moving towards increasingly unstructured environments

If only **robots could learn how to perform tasks by themselves...**

⇒ **Optimization and Learning in Robotics**

Motivation

Robots



Mülling + Peters

Humans



We need **learning** and **adaptation** to improve robot skills!

Motivation

Optimization vs. Learning

- Learning approaches use Optimization
- However there are some **differences**:
 - **Optimization** is focused on **training set**, while **Learning** is focused on **general performance**
 - The **loss function** used in learning is not always the real (optimized) loss function
 - **Stopping criteria** may be quite different

Motivation

Challenges in Robot Learning

- Cost of experimentation
- Cost of failure
- Limited data
- Generalization
- Curse of dimensionality
- Real time requirements
- Changes in environment
- Changes in task specification

Motivation

Learning in Robotics can be used for:

- Robot Perception
- Robot Decision
- Adapt Human-Robot Interaction
- Robot Actuation (Behaviors)
- Multi-robot Coordination and Communication

EuRoC - TIMAIRIS Team

- Challenger: Universidade de Aveiro, Portugal
- End User: IMA S.p.A, Italy



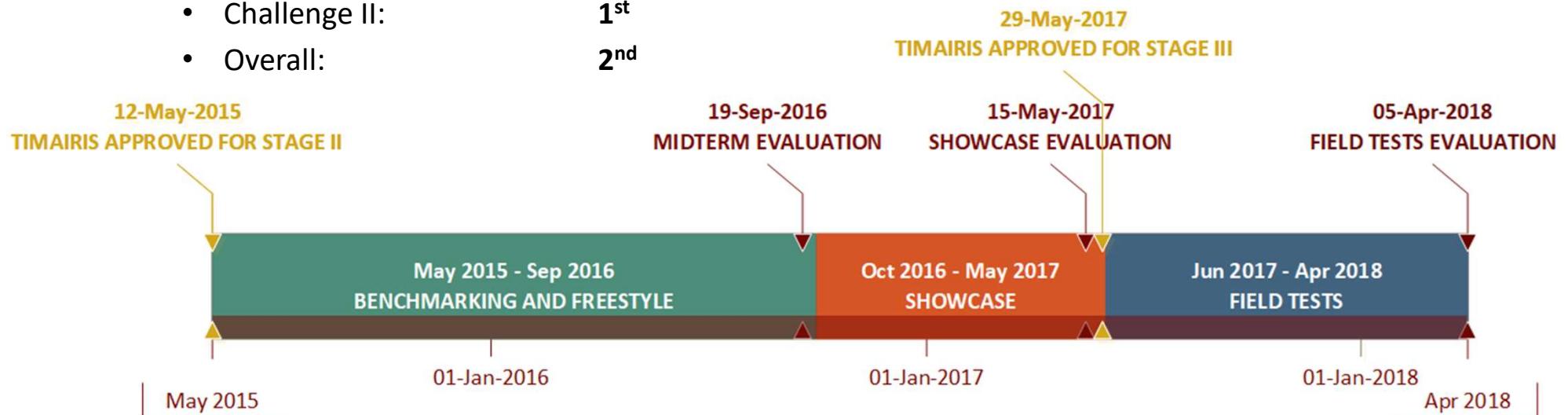
- Stage I (Overall: 103 teams, Chall. II: 37 teams)

- Stage II Results (Challenge II: 5 teams)

- Benchmarking: 1st
- Freestyle: 2nd
- Showcase: 1st

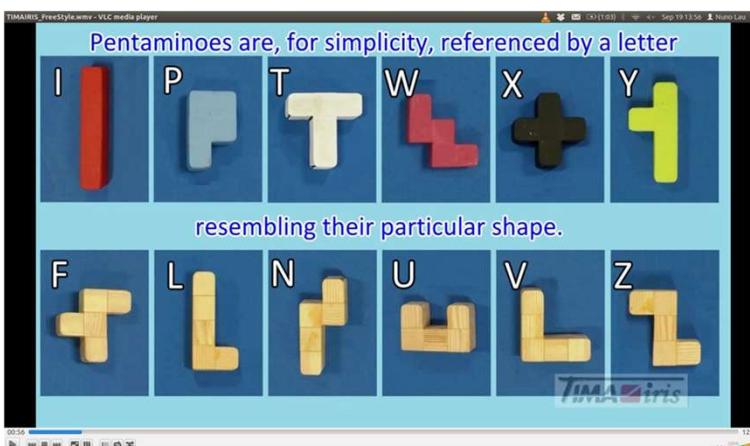
- Stage III Results

- Challenge II: 1st
- Overall: 2nd



Gesture Recognition

Task: Assembling a puzzle cooperatively by a human and a robot (EuRoC Project)

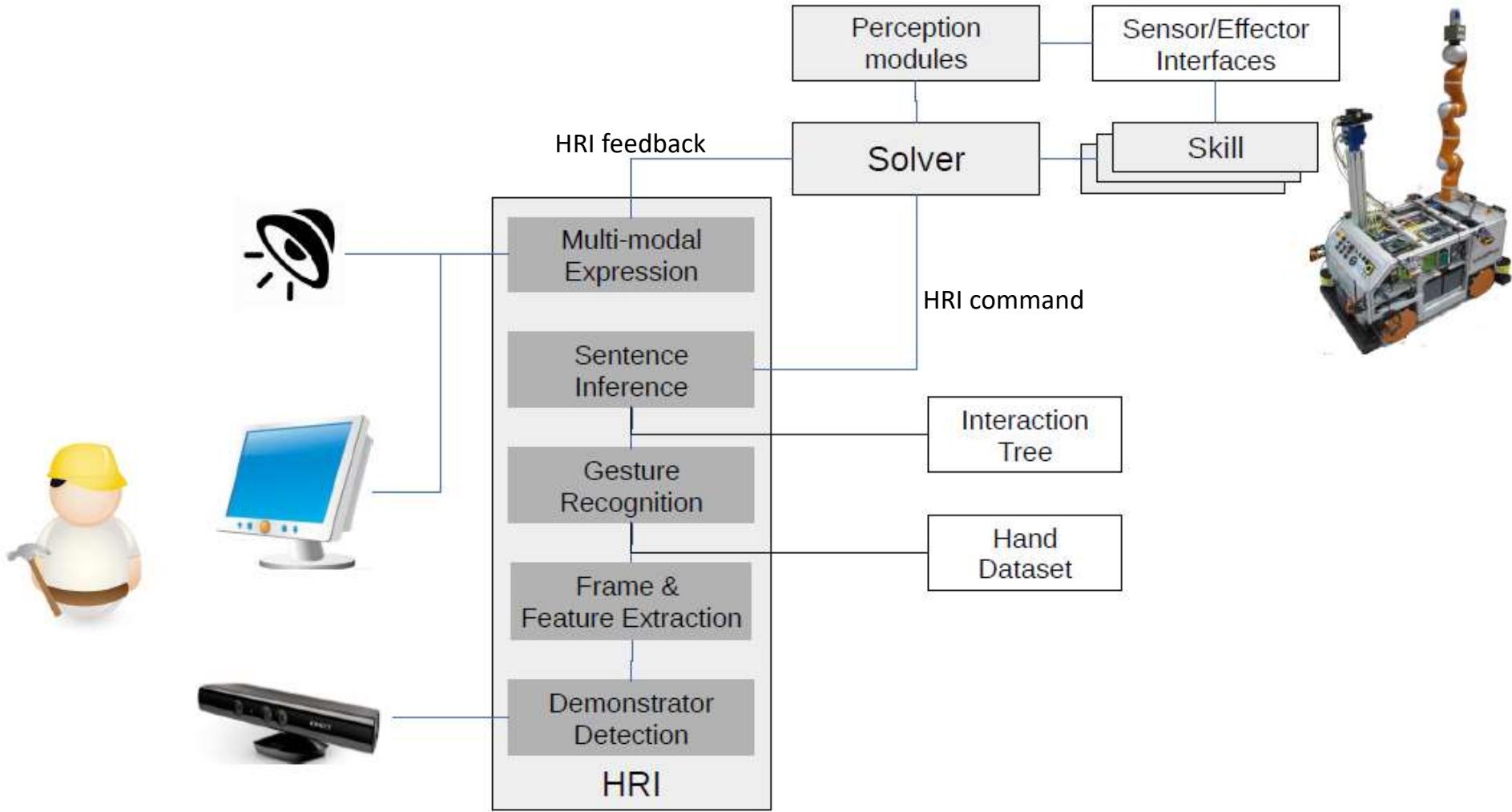


Set of 12 pentomino pieces



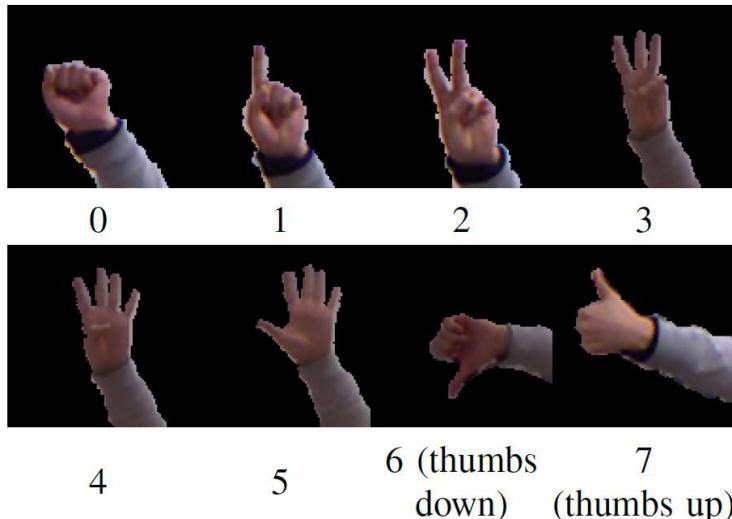
Task environment

Gesture Recognition

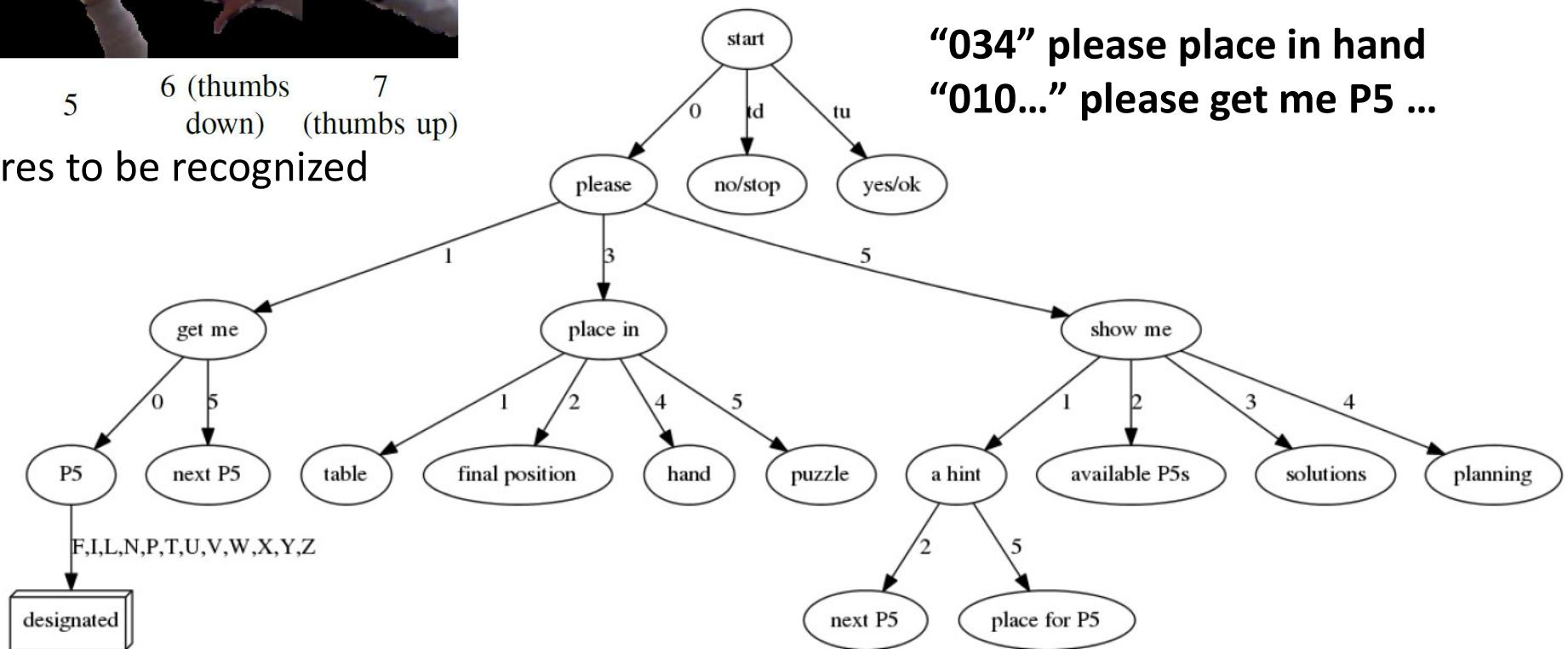


Human-Robot Interface architecture

Gesture Recognition

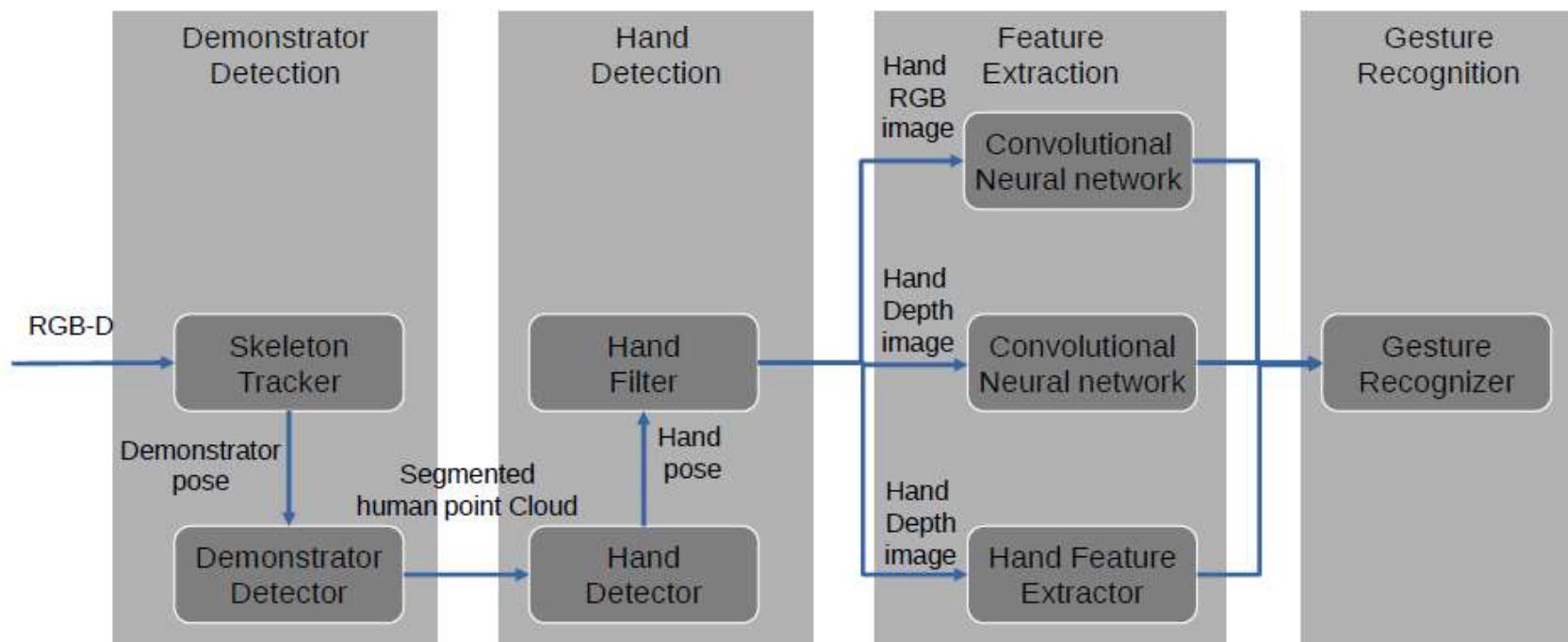


Gestures to be recognized

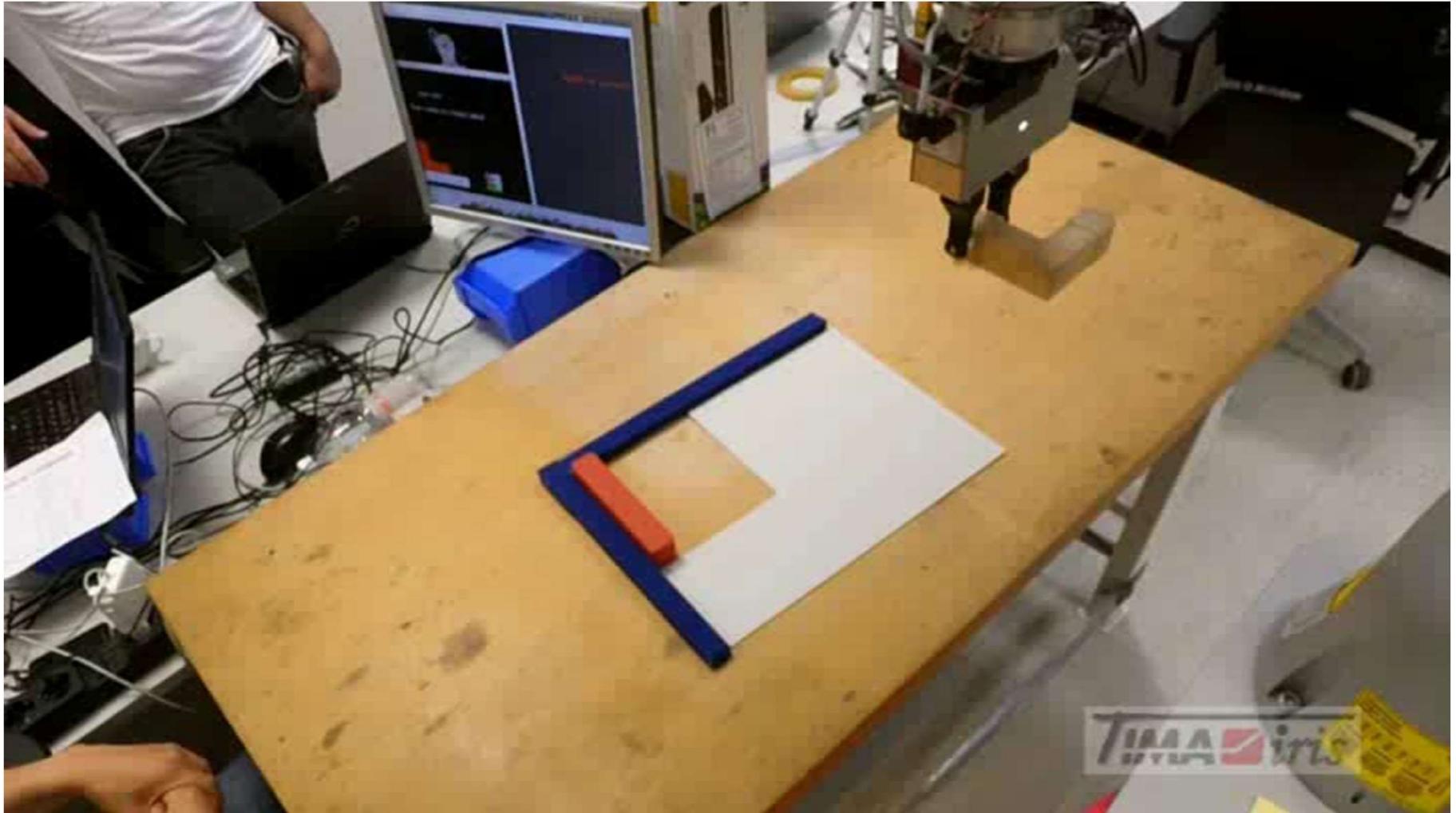


Gesture Recognition

- **Task:** Recognize Gestures
- **Approach:**
 - 1st : Use Deep Learning
 - 2nd : Mix Deep Learning with other hand features



Gesture Recognition



Lim, G.H. et al. Skill-based anytime agent architecture for European Robotics Challenges in realistic environments: EuRoC Challenge 2, Stage II — realistic labs, Robotics and Autonomous Systems, Vol. 120, 2019

Stochastic Search

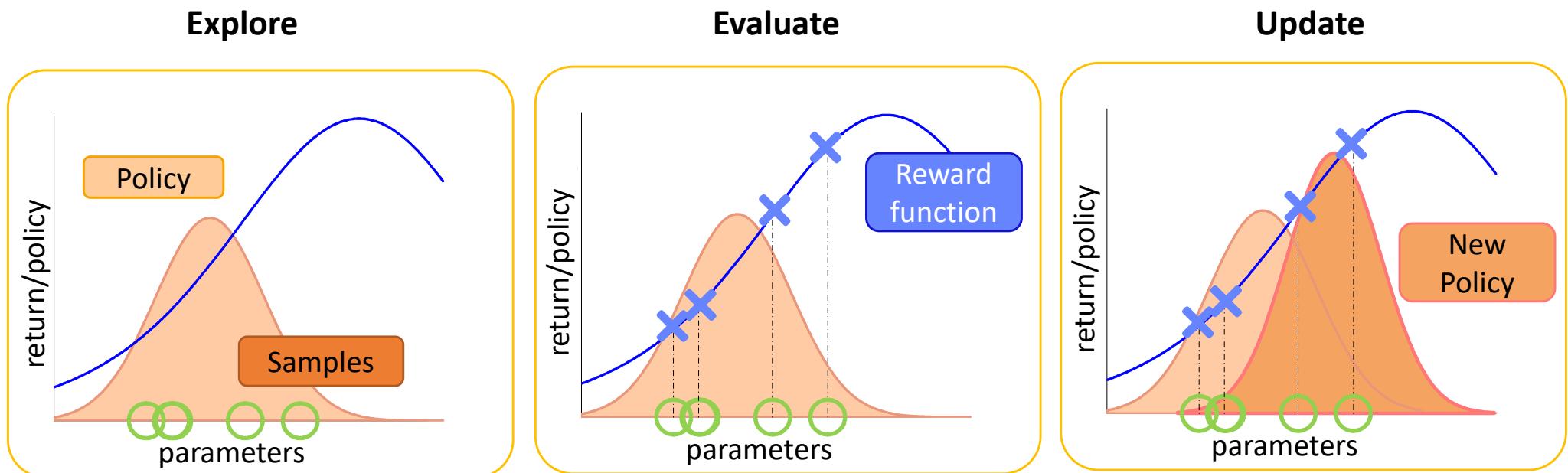
$$\pi(\mathbf{w}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

- **Use Search-Distribution:**

$$\pi(\mathbf{w})$$

$$J_\pi = \int \pi(\mathbf{w}) R(\mathbf{w}) d\mathbf{w}$$

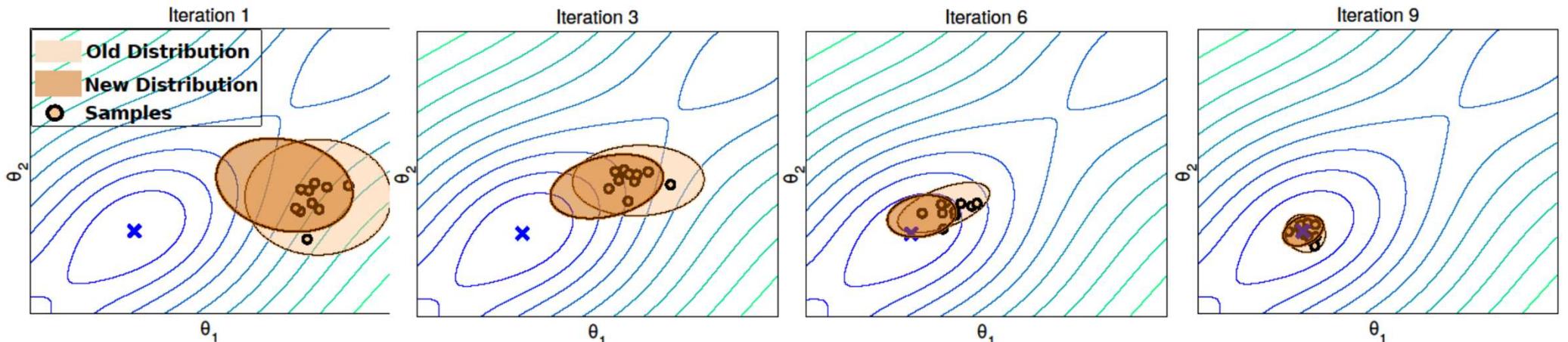
- **Objective:** Find search distribution that maximizes



Stochastic Search

- **Use Search-Distribution:** $\pi(\mathbf{w}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$
 - **Mean:** Estimate of the maximum
 - **Covariance:** Direction to explore
- **Objective:** Find search distribution $\pi(\mathbf{w})$ that maximizes

$$J_\pi = \int \pi(\mathbf{w}) R(\mathbf{w}) d\mathbf{w}$$



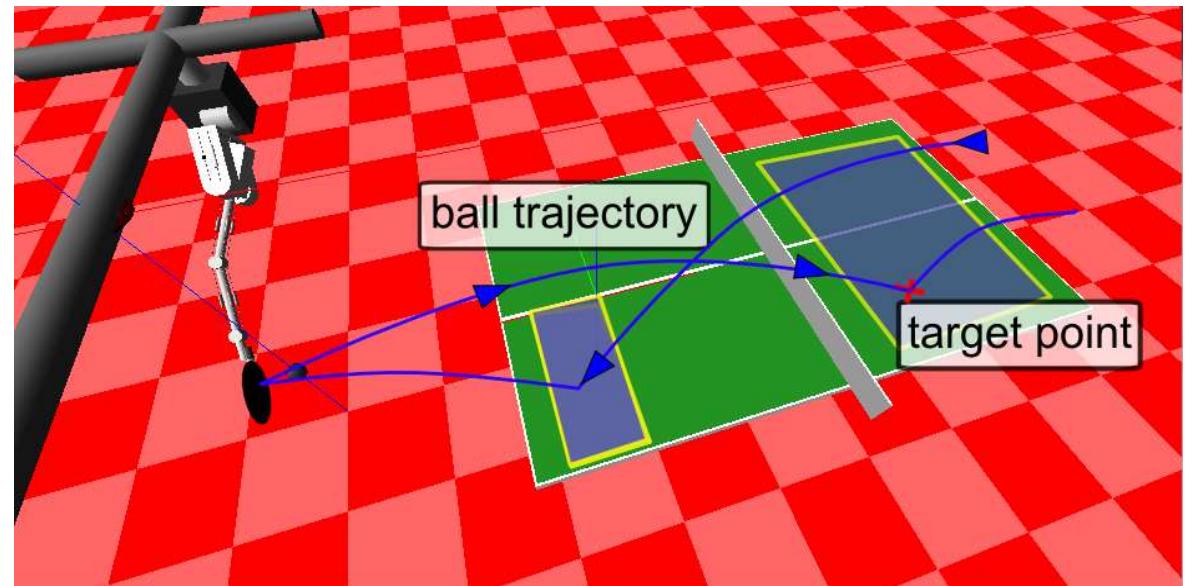
Contextual Stochastic Search

Goal: Adapt parameters w to different situations

- Different ball trajectories
- Different target locations

Introduce context vector s

- Continuous valued vector
- Characterizes environment and objectives of agent



Learn contextual search policy

$$\pi(w|s)$$

Abdolmaleki, et. al, *Model-Based Relative Entropy Stochastic Search*, NIPS 2015

Abdolmaleki, A. et al., Contextual Direct Policy Search: With Regularized Covariance Matrix Estimation, Journal of Intelligent and Robotic Systems: Theory and Applications, 96 (2), pp. 141-157, 2019

Adaptation of Skills

Contextual distribution:

$$\pi(\mathbf{w}|\mathbf{s}) = \mathcal{N}(\mathbf{s}^T \mathbf{M} + \mathbf{m}, \Sigma)$$

Compatible Function Approximation:

$$R(\mathbf{s}, \mathbf{w}) \approx \mathbf{w}^T \mathbf{A} \mathbf{w} + \mathbf{s}^T \mathbf{B} \mathbf{w} + \mathbf{a}^T \mathbf{w} + a_0$$

Contextual distribution update:

1. Maximize **expected** return
2. Bound **expected** information loss
3. Bound entropy loss

$$\arg \max_{\pi} \mathbb{E}_{p(\mathbf{s})} \left[\int \pi(\mathbf{w}|\mathbf{s}) R(\mathbf{s}, \mathbf{w}) d\mathbf{w} \right]$$

$$\text{s.t.: } \mathbb{E}_{p(\mathbf{s})} [\text{KL}(\pi(\cdot|\mathbf{s}) || \pi_{\text{old}}(\cdot|\mathbf{s}))] \leq \epsilon$$

$$\underbrace{H(\pi_{\text{old}}) - H(\pi)}_{\text{loss in entropy}} \leq \gamma$$

New distribution: $\pi(\mathbf{w}|\mathbf{s}) \propto \pi_{\text{old}}(\mathbf{w}|\mathbf{s})^{\frac{\eta}{\eta+\omega}} \exp \left(\frac{R(\mathbf{s}, \mathbf{w})}{\eta + \omega} \right)$

$$\propto \mathcal{N}(\mathbf{s}^T \mathbf{M}_{\text{new}} + \mathbf{m}_{\text{new}}, \Sigma_{\text{new}}) \quad \text{← Compatible Function Approximation}$$

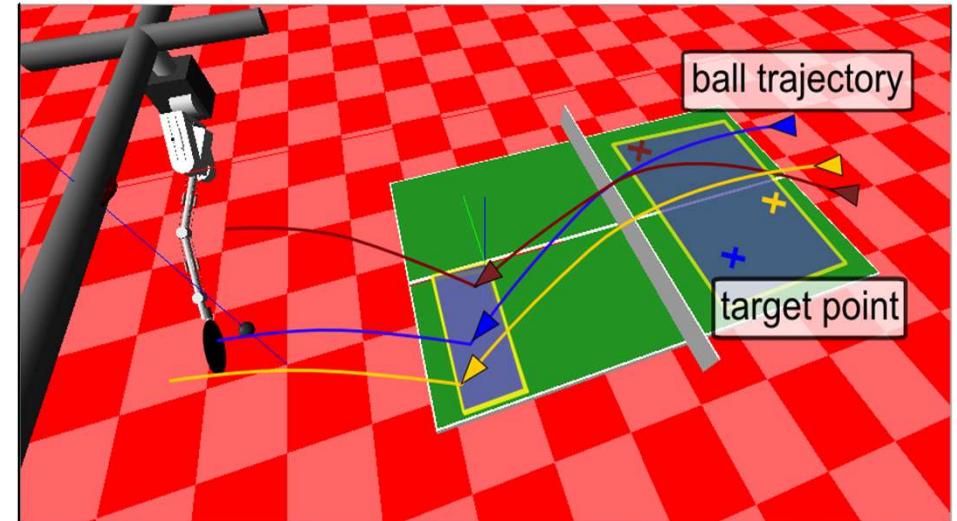
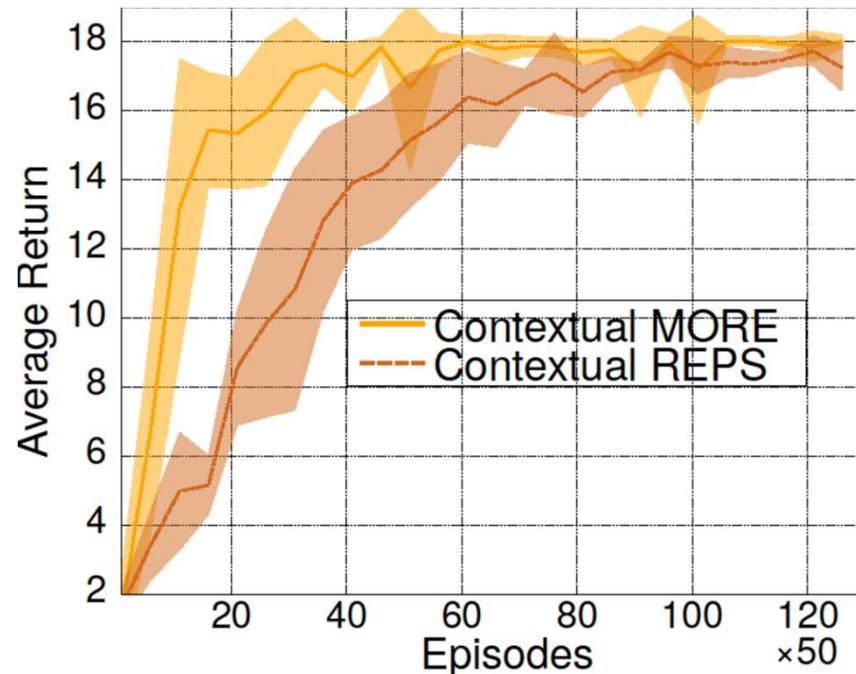
Adaptation of Skills: Table Tennis

Contextual Stochastic Search:

- Context: Initial ball velocity

Reward:

- Hit ball
- Ball impacts at target position

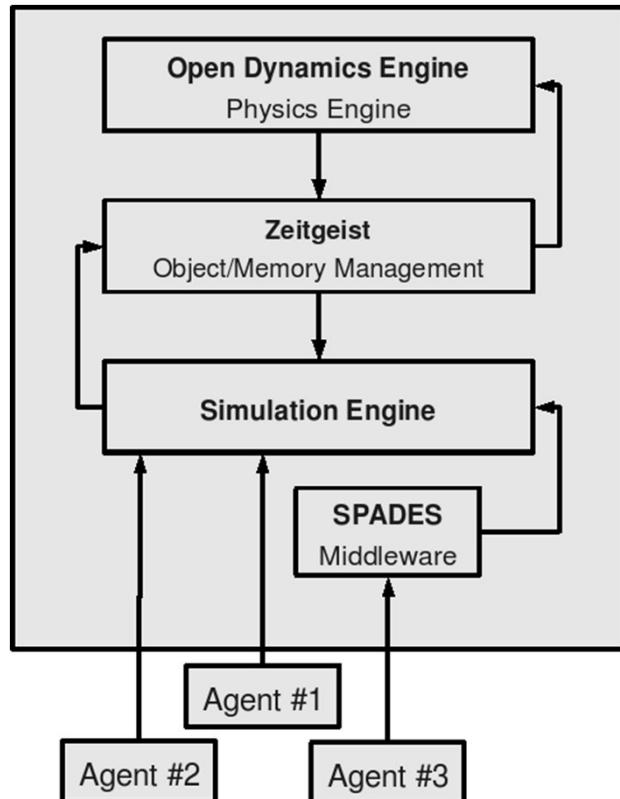


Skills Improvement:

- ✓ Hot-start with imitation
- ✓ Continuous-valued decision making
- ✓ Low number of samples
- ✓ Adaptation

RoboCup Simulation 3D Simulator

Simspark server architecture



Simulation 3D game



RoboCup 3D Simulated Agent

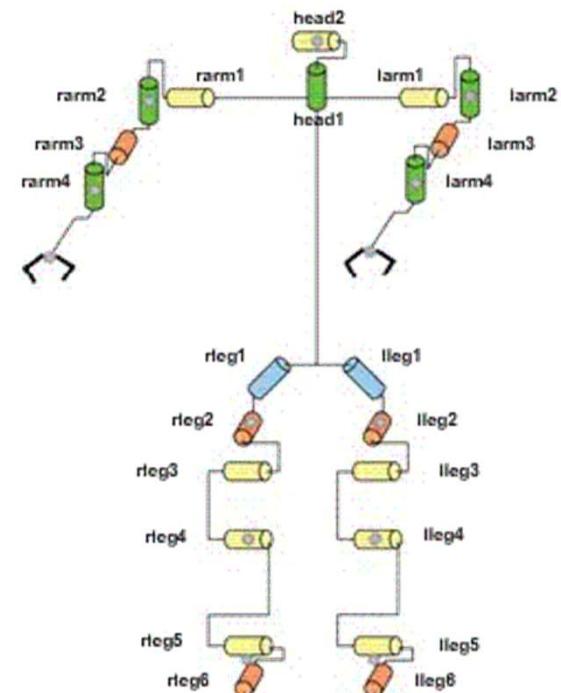
- **Body:** 22 degrees of freedom (DoFs), 57 cm height and 4.5 kg
- **Perceptors:** gyroscope; force sensors; DOFs; vision; hear
- **Effectors:** beam; DoFs; say



NAO (Aldebaran Robotics)



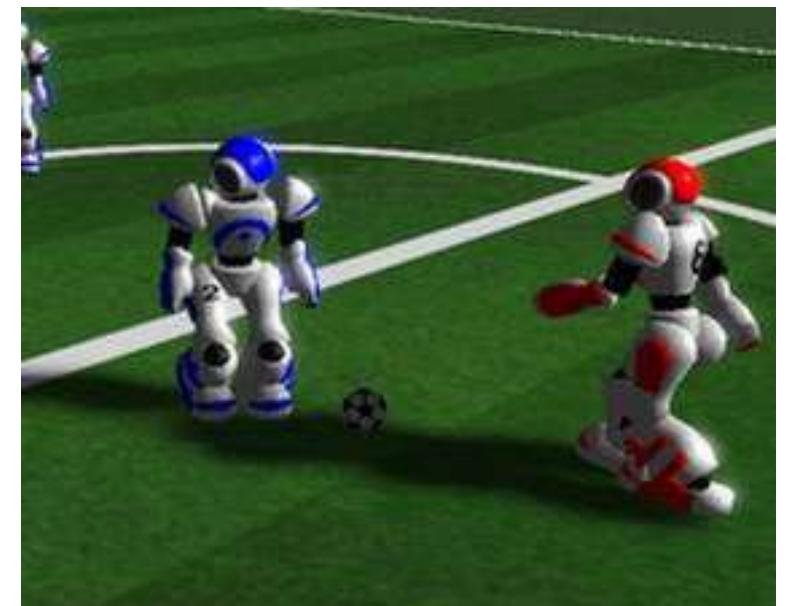
Simulated NAO (3D League)



Body structure

Skill Improvement: Controlled Kick

- **Task**
 - Develop a **kick with controlled kicking distance**
 - From 10 different positions in the soccer field (with distances ranging from 3m to 12m), kick the ball so that it stops in the center of the field
- **Classical approach**
 - Optimize for each distance
- **Contextual approach**
 - Optimize for all distances in a single process
 - Use all data to improve performance
 - Generalize for unknown contexts



Skill Improvement: Controlled Kick



Abbas Abdolmaleki et al. Learning a Humanoid Kick With Controlled Distance. RoboCup 2016: Robot World Cup XX, Springer, July 2016

User Profiles and Adapted Interfaces

- Users of Intelligent Wheelchairs have very different skills
- Command interface/language provided for each user should be adapted to his/her capabilities

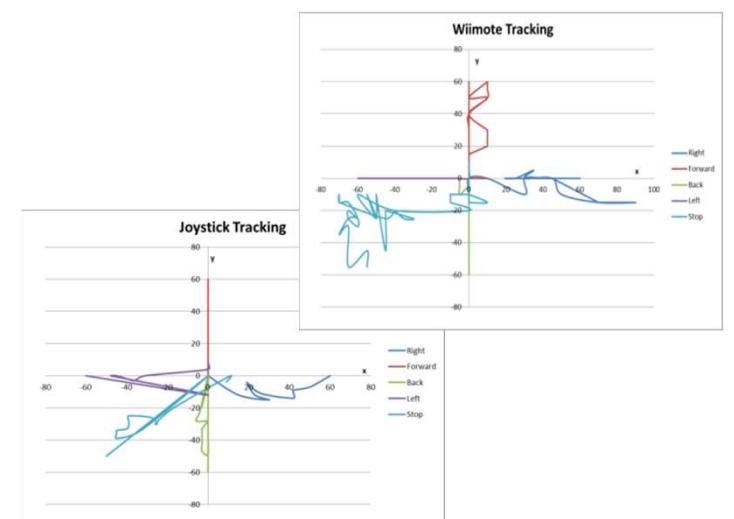
- Task
 - Generate command language adapted to each user for driving the IW
- Approach
 - Optimize command language
 - User Profiling provides relevant information



User Profiles and Adapted Interfaces

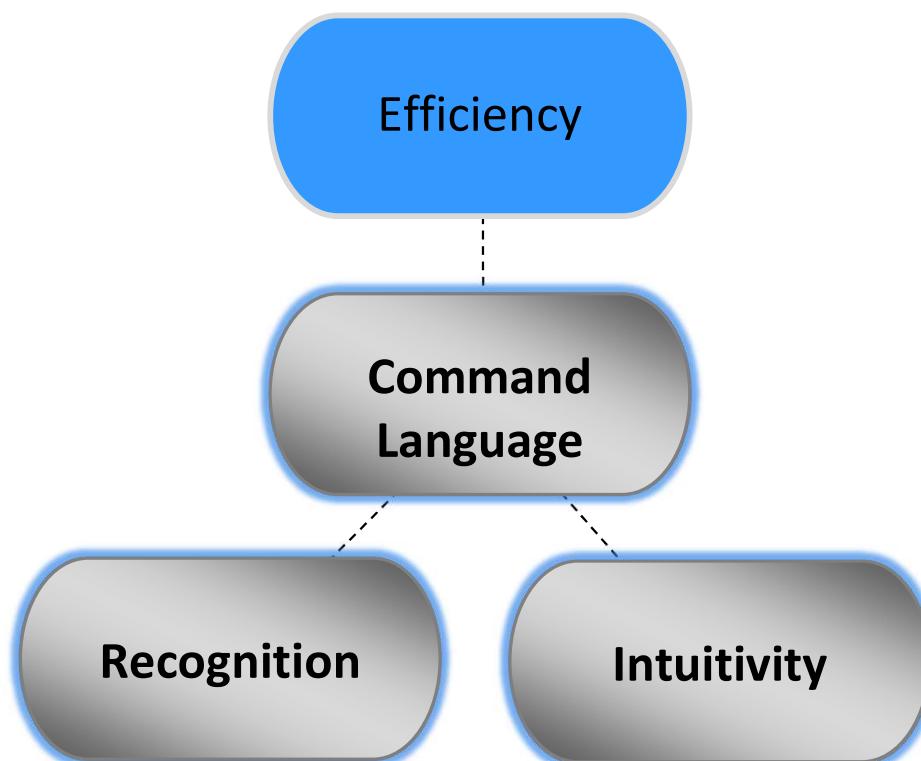
- **User Profiling Experiments**

- 11 cerebral palsy users
- Level IV (27.3%) and V (72.7%) GMFM
- Voice Inputs
 - “Go”, “Front”, “Forward”, “Back”, “Right”, “Left”, “Turn”, “Spin” and “Stop”
- Joystick and the Head Movements



User Profiles and Adapted Interfaces

- Command Language

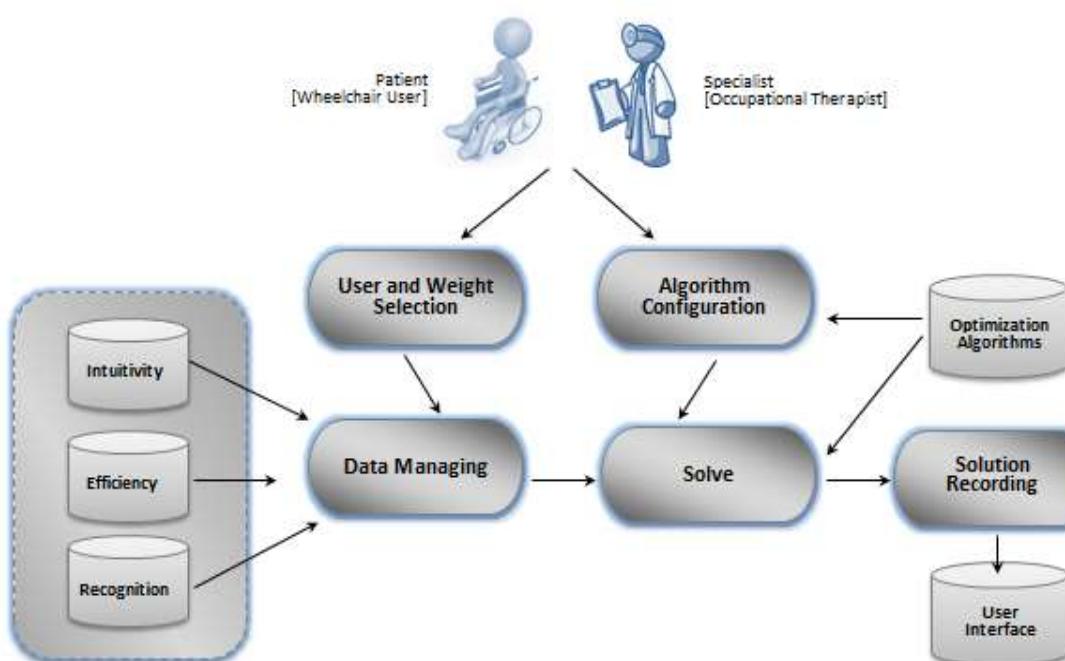


User Profiles and Adapted Interfaces

- Optimized Command Language

Maximizes the function composed by the total time efficiency, recognition and intuitiveness

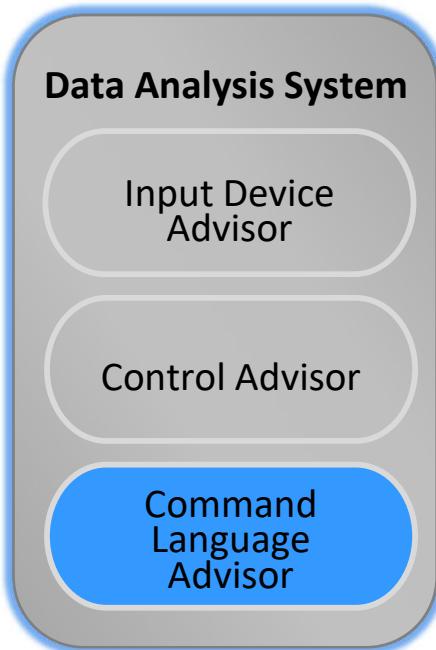
$$\arg \max_{T_{eff}, T_{reg}, T_{int}} (\alpha T_{eff} + \beta T_{reg} + \gamma T_{int})$$



```
(w_rec, w_time, w_intu) = weights; evaluation ← 0
for ncom = 1 to NC do
    recVal ← 1; timeVal ← 0; intuVal ← 1
    for nseq = 1 to NS do
        inpDev ← inputDevice(solution[ncom][nseq])
        inp ← input(newSolution[ncom][nseq])
        if inpDev = NULL then break
        else
            recVal ← recVal * rec[inpDev][inp]
            timeVal ← timeVal + time[inpDev][inp]
            intuVal ← intuVal * intu[ncom][inpDev][inp]
        endif
    endfor
    evalComm ← w_rec* recVal + w_time*1/(timeVal+1)
        + w_intu*intuVal
    evaluation ← evaluation + evalComm
endfor
return evaluation
```

User Profiles and Adapted Interfaces

- Command Language Advisor



Mean of DAS evaluation higher than mean of evaluation of the command language recommended by specialist (p value = 0.002)

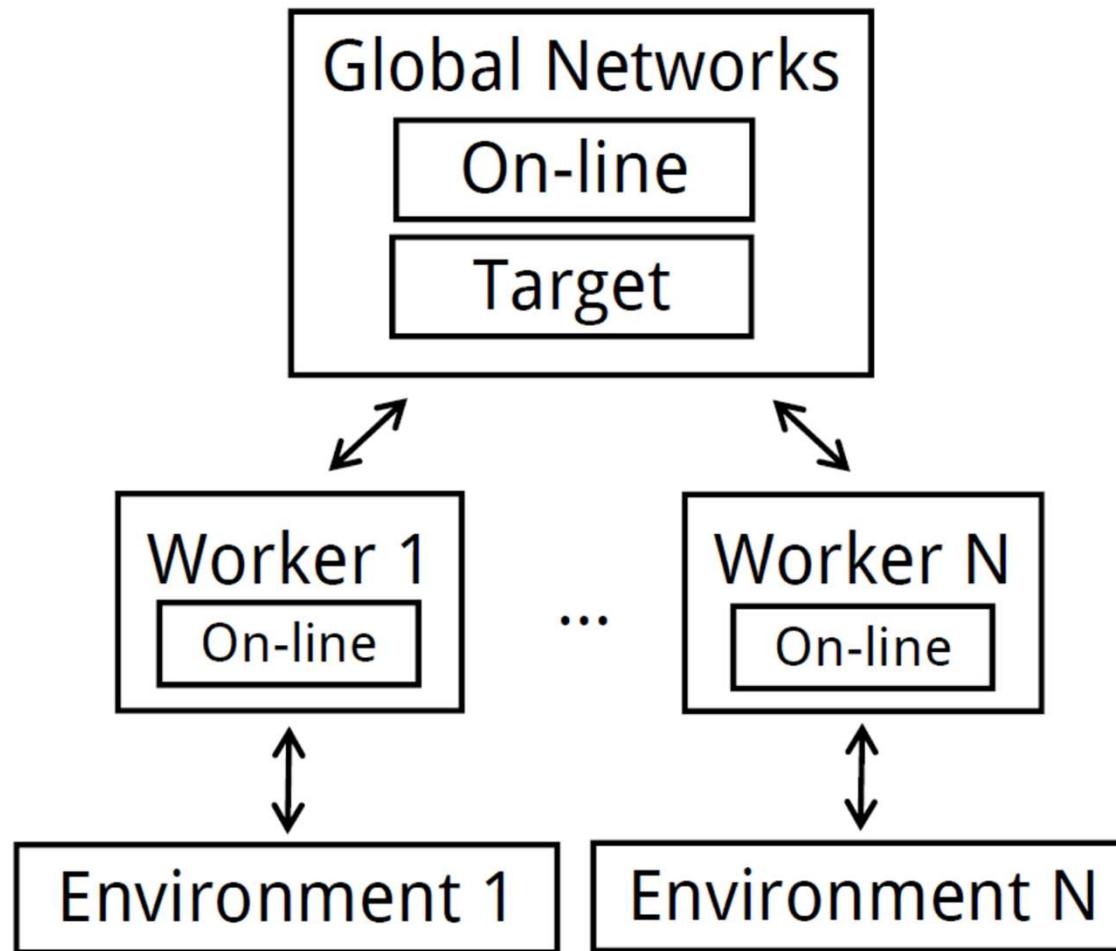
Patient	Evaluation	Forward	Command Language for Patients			
			Left	Right	Back	Stop
P1 Specialist IDAS	4.53 4.57	wiimote joystick	joystick joystick	joystick joystick	joystick joystick	joystick joystick
P2 Specialist IDAS	4.18 4.85	joystick joystick	joystick joystick	joystick joystick	joystick joystick	voice ("stop") voice ("go")
P3 Specialist IDAS	3.33 4.51	voice ("forward") wiimote	wiimote wiimote	wiimote wiimote	joystick wiimote	voice ("stop") voice ("go")
P4 Specialist IDAS	4.50 4.60	voice ("forward") joystick	joystick joystick	joystick joystick	joystick joystick	voice ("stop") voice ("stop")
P5 Specialist IDAS	4.14 4.40	voice ("front") wiimote	wiimote wiimote	wiimote voice ("turn")	joystick joystick	voice ("stop") voice ("stop")
P6 Specialist IDAS	4.13 4.38	wiimote wiimote	joystick wiimote	joystick wiimote	joystick wiimote	joystick wiimote
P7 Specialist IDAS	4.49 4.60	voice ("front") joystick	joystick joystick	joystick joystick	joystick voice ("back")	voice ("stop") voice ("stop")
P8 Specialist IDAS	3.51 4.20	wiimote wiimote	joystick wiimote	joystick wiimote	joystick wiimote	joystick wiimote
P9 Specialist IDAS	3.70 4.75	voice ("forward") joystick	wiimote joystick	wiimote joystick	joystick joystick	voice ("stop") joystick
P10 Specialist IDAS	4.11 4.80	voice ("forward") joystick	voice ("left") joystick	voice ("right") voice ("turn")	voice ("turn") joystick	voice ("stop") voice ("go")
P11 Specialist IDAS	4.29 4.30	joystick wiimote	wiimote wiimote	wiimote wiimote	joystick wiimote	joystick wiimote

Multiagent Communication Learning

- Multiagent reward based learning challenges
 - Non static environment
 - Complexity exponential to number of agents
 - Learn how to use Communication
- **Task**
 - Learning Coordination and learning Communication from scratch
- **Proposal**
 - Asynchronous Advantage Actor Centralized-Critic with Communication (A3C3)

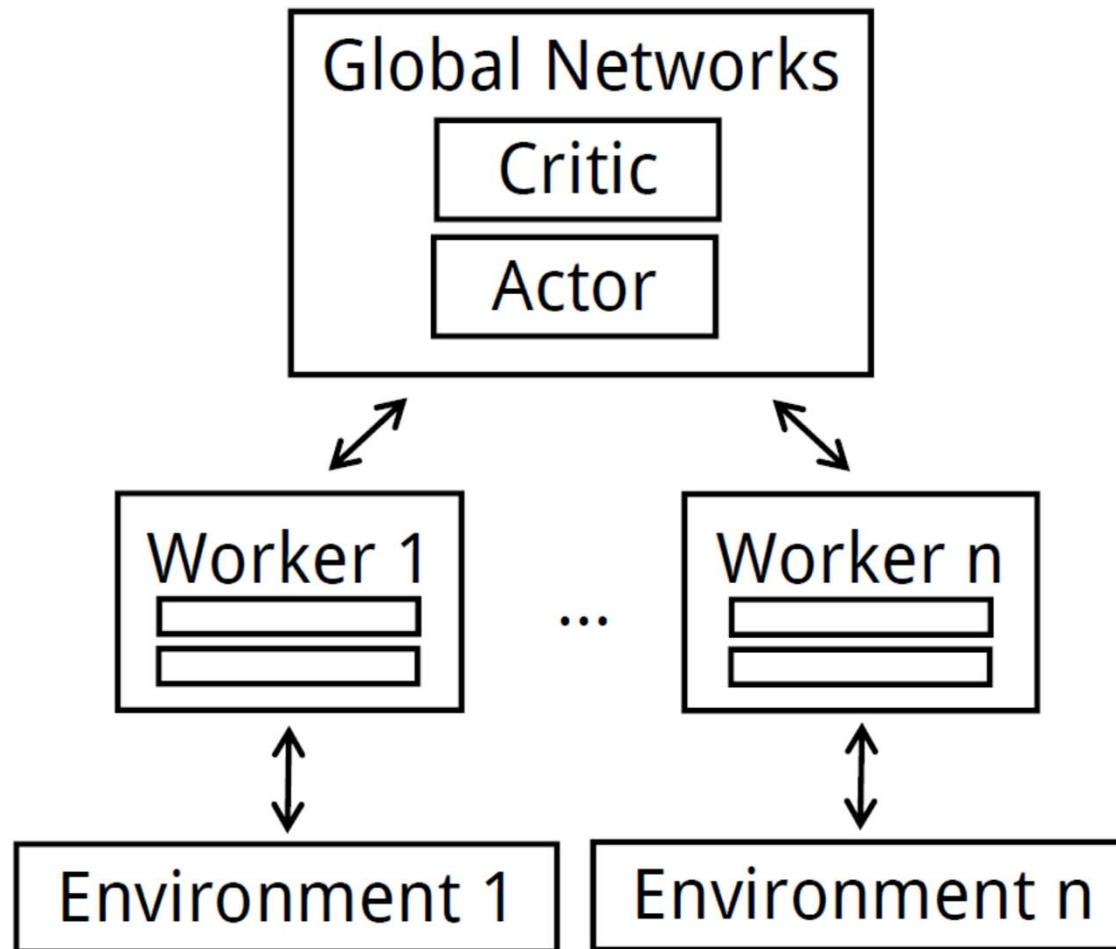
Multiagent Communication Learning

- Asynchronous n-step Q-Learning



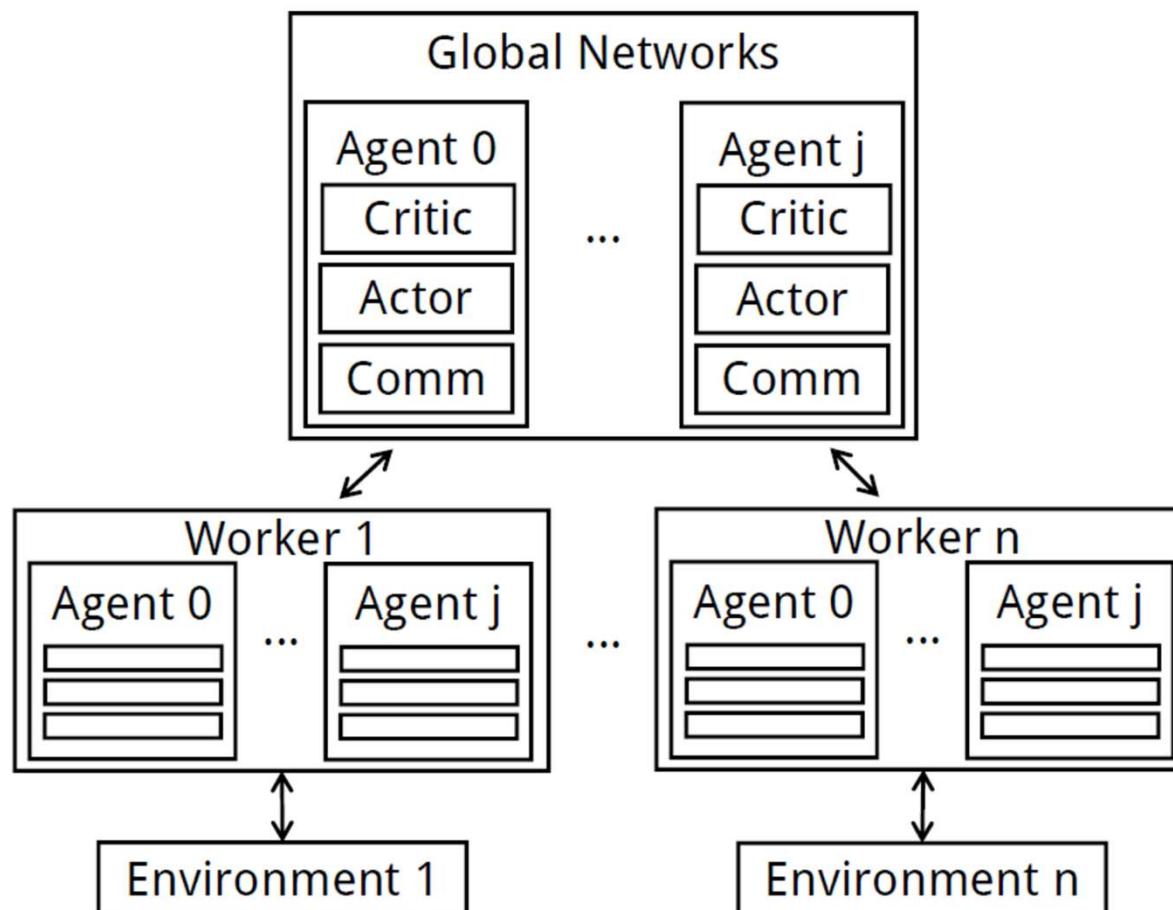
Multiagent Communication Learning

- A3C – Asynchronous Advantage Actor Critic



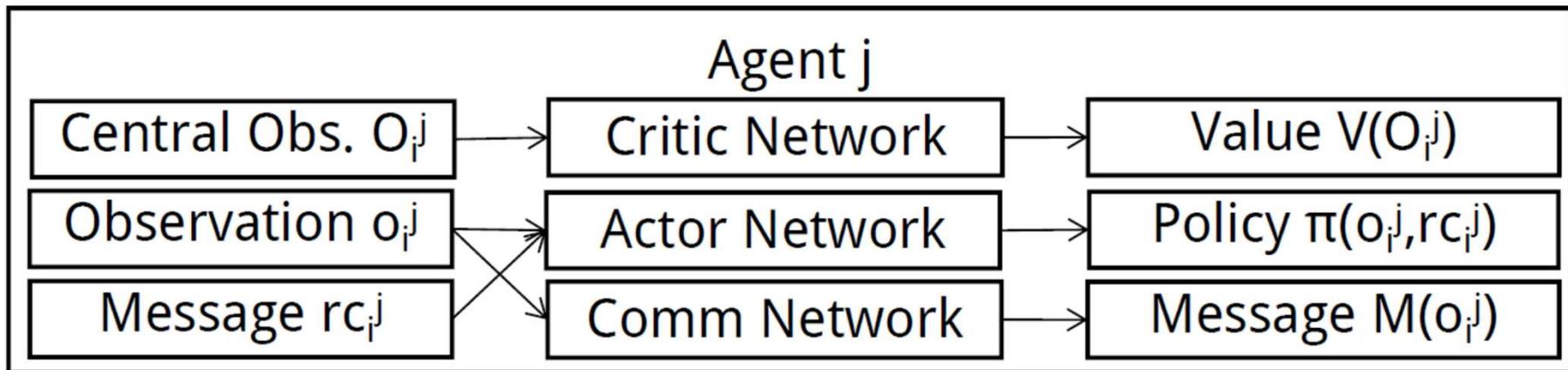
Multiagent Communication Learning

- **Proposal:** A3C3 – Asynchronous Advantage Actor-Centralized-Critic with Communication



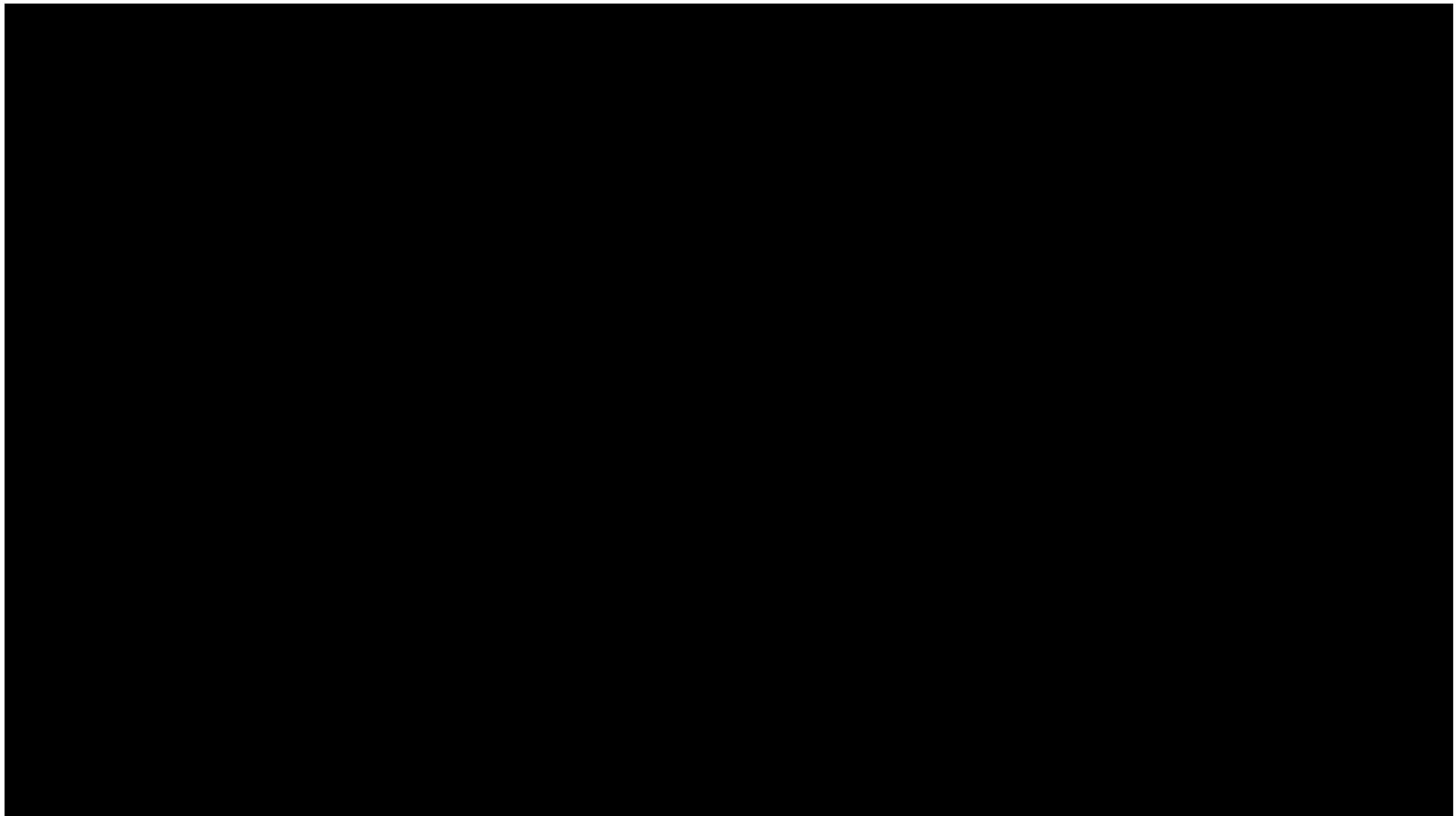
Multiagent Communication Learning

- **Proposal:** A3C3 – Asynchronous Advantage Actor Centralized-Critic with Communication



- Centralized learning phase but distributed execution
- Centralized Critic fed with complete state or aggregate of local observations
- Communication at each time step learned *tabula rasa*

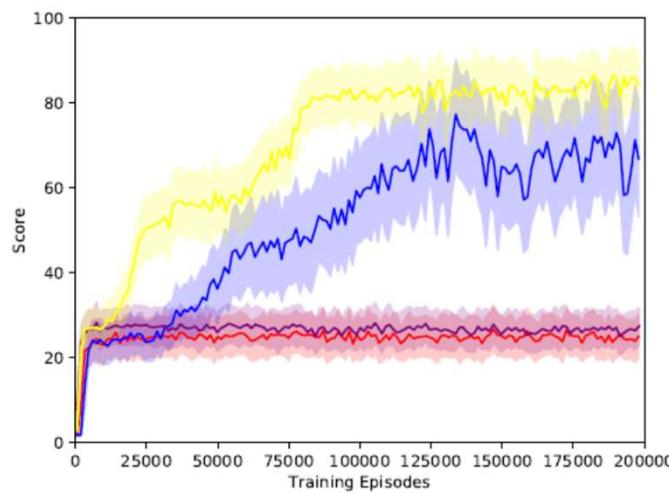
Multiagent Communication Learning



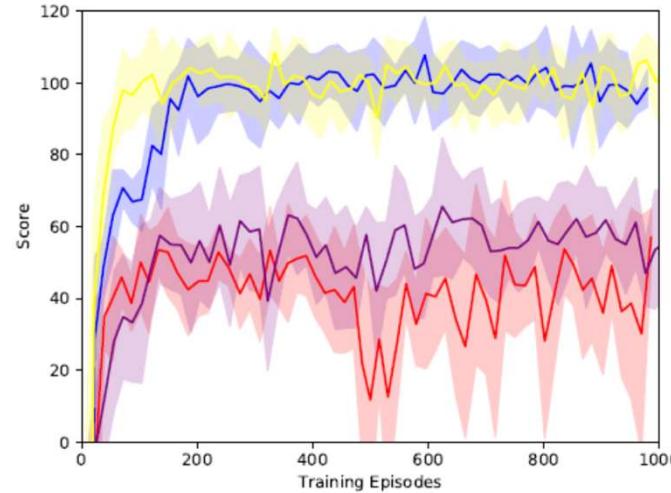
David Simões et al. Multi-agent deep reinforcement learning with emergent communication. In 2019 International Joint Conference on Neural Networks (IJCNN). IEEE, 2019.

Multiagent Communication Learning

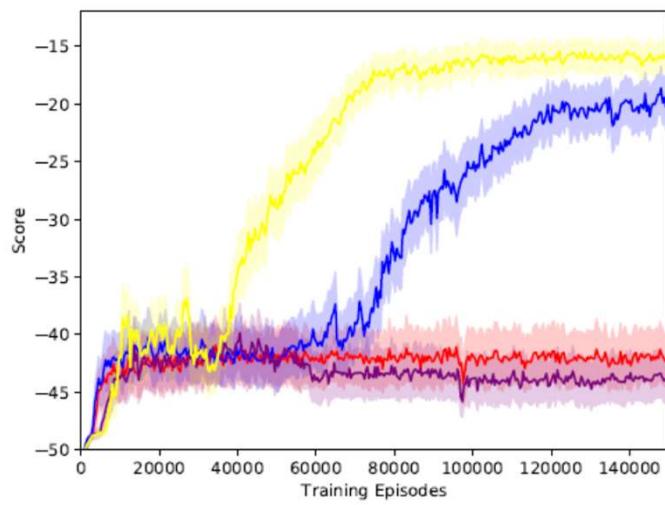
■ A3C. ■ A3C2. ■ A3C3 (No Comm). ■ A3C3.



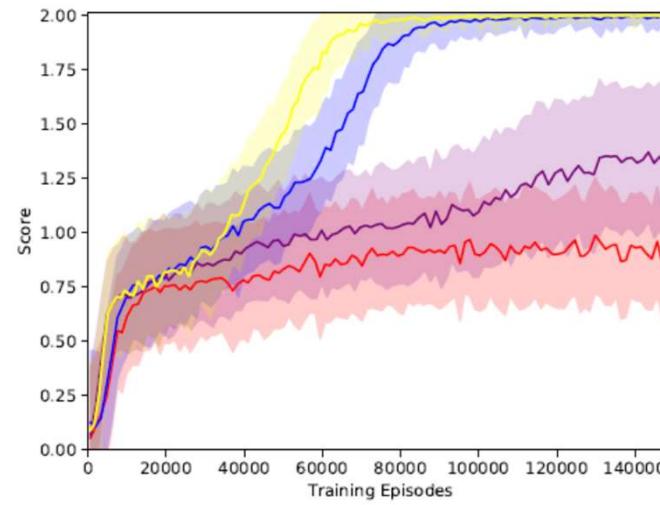
(a) Hidden Reward challenge.



(b) Traffic Intersection simulator.



(c) Pursuit game.



(d) Navigation task.

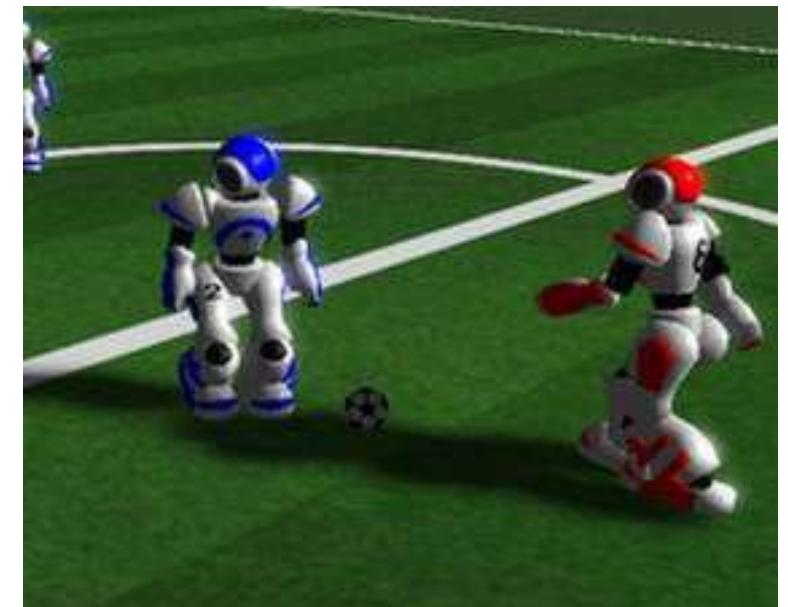
Multiagent Coordination

- Formation specification

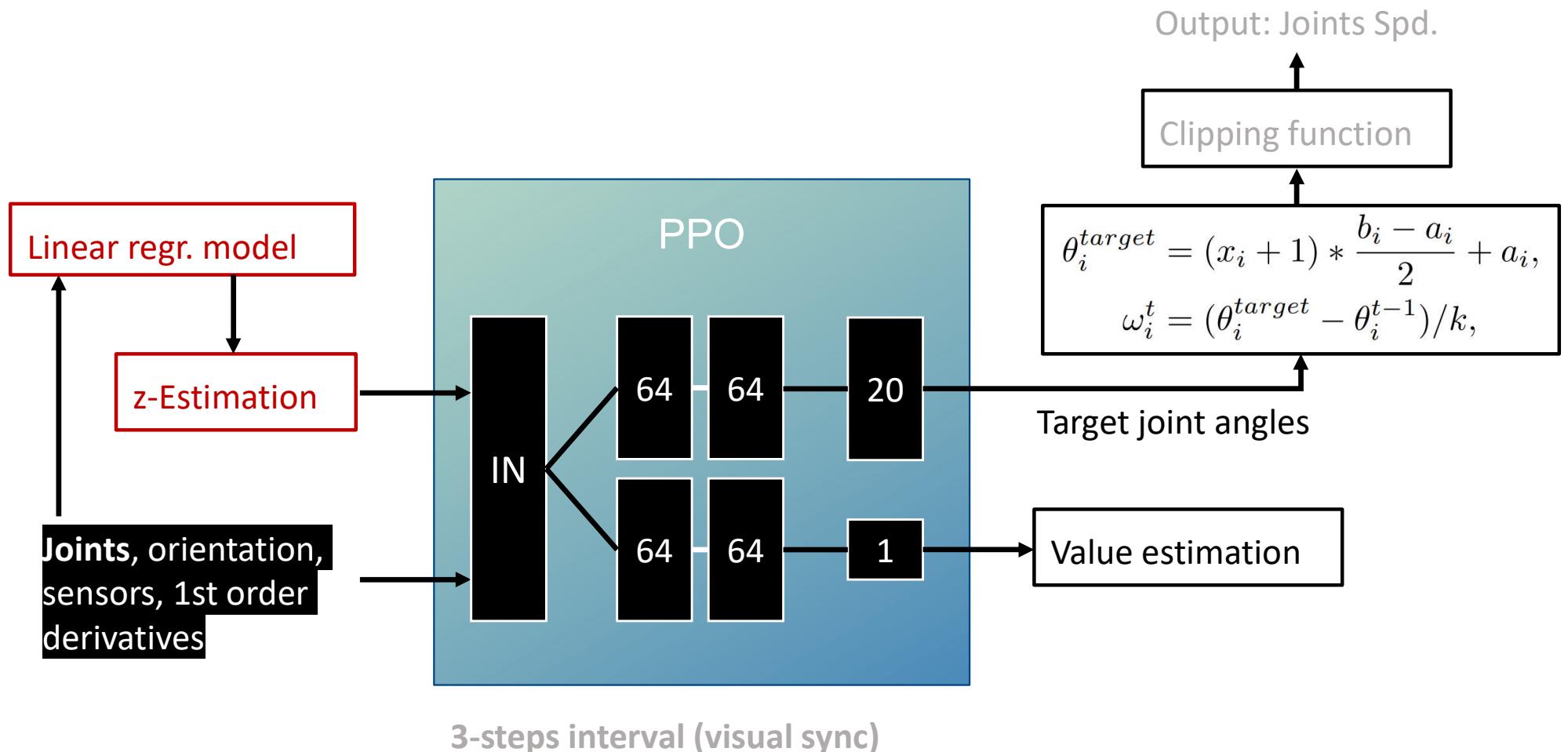


Behavior Development

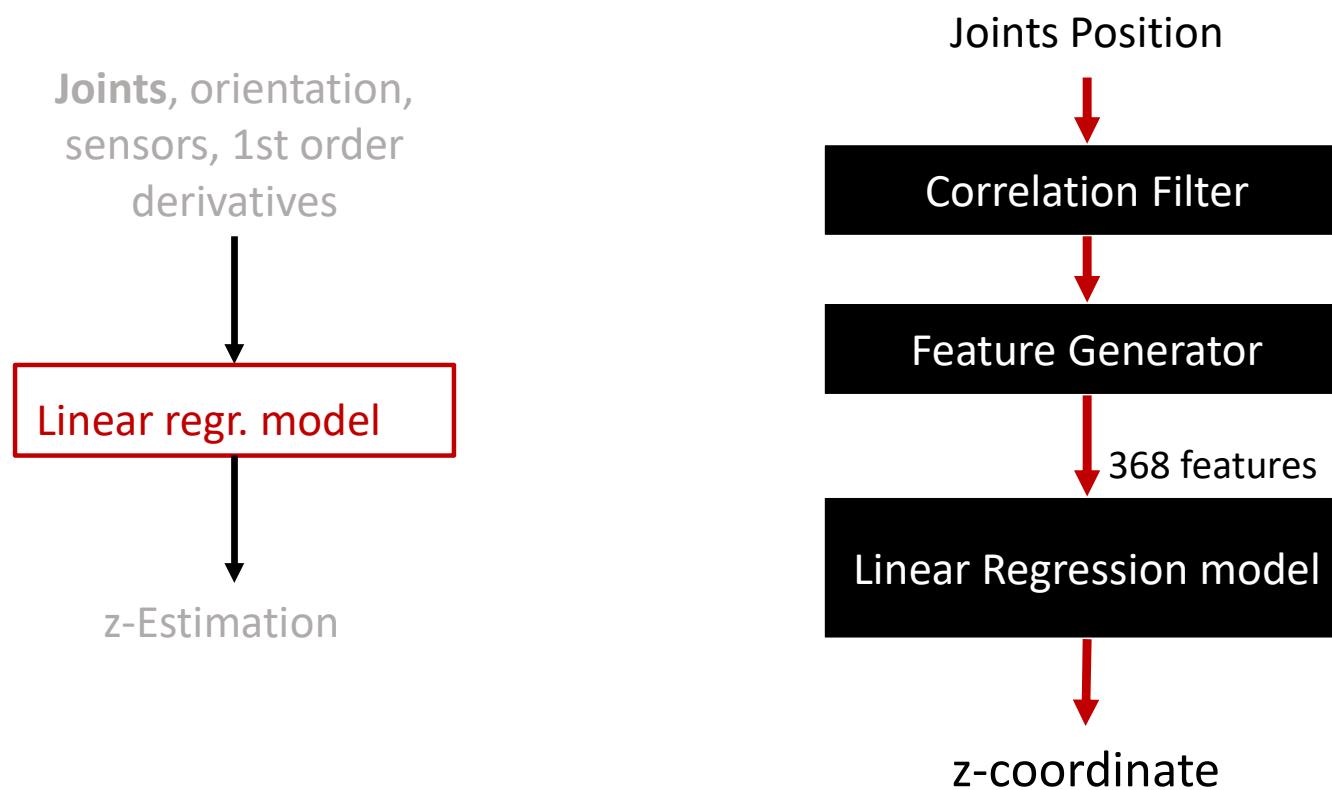
- **Task**
 - Generate fast walking engine for simulated humanoid robot
- **Classical Approach**
 - Model based Walk Engine (ZMP)
 - Optimize model parameters
- **Our Approach**
 - Deep RL: **PPO algorithm**
 - Feedback Walk Engine
 - Use large set of features
 - Initial training using ground-truth data
 - Estimate ground-truth using regression filter
 - Train with agent perceptions

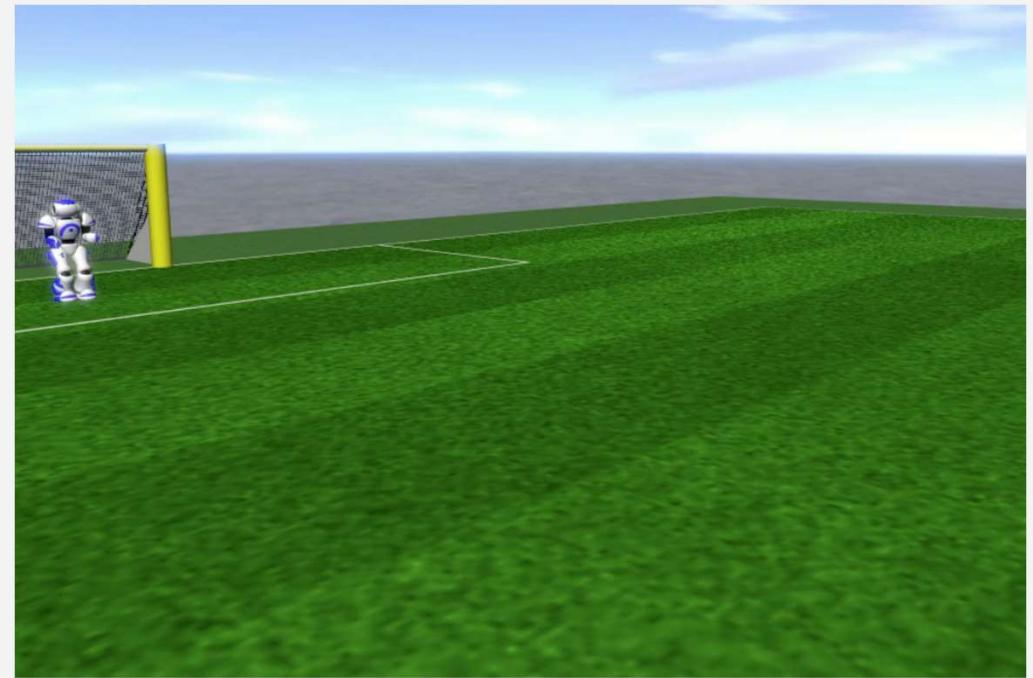
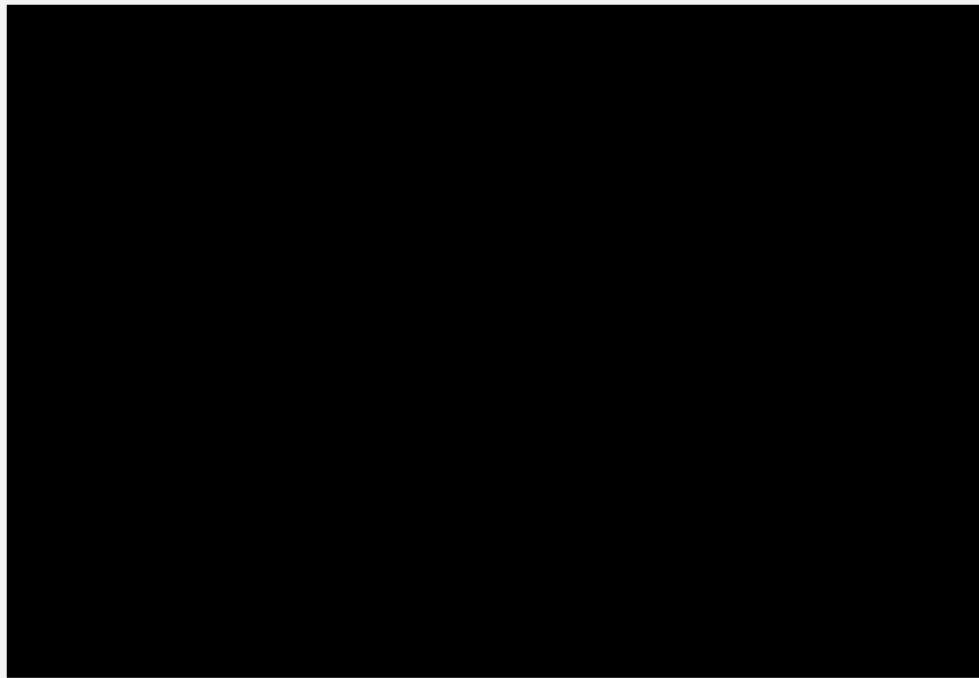


Behavior Development

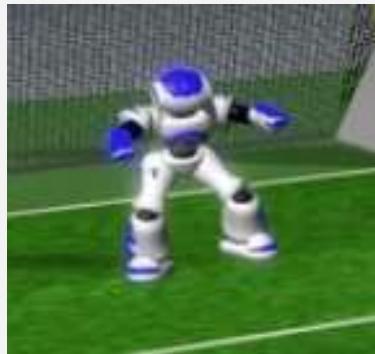


Estimation of Z-true value



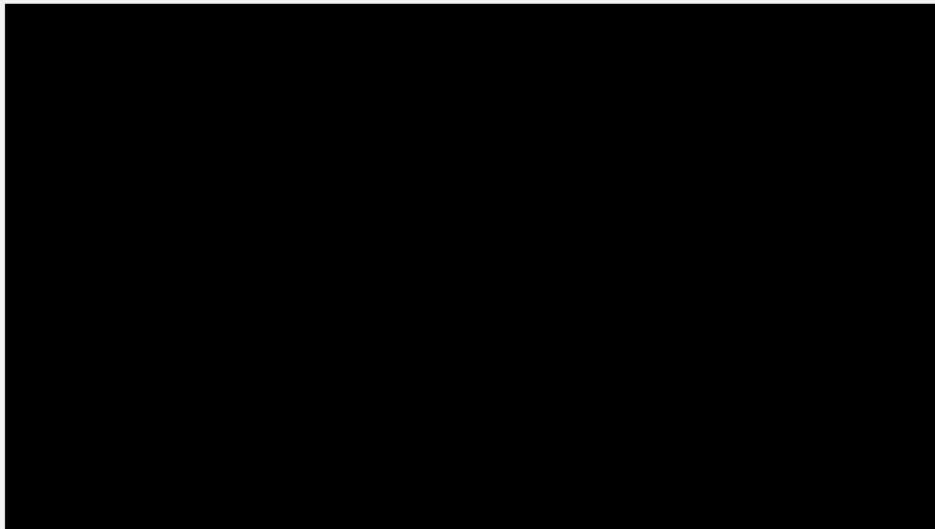


sprint
v1



sprint
v2

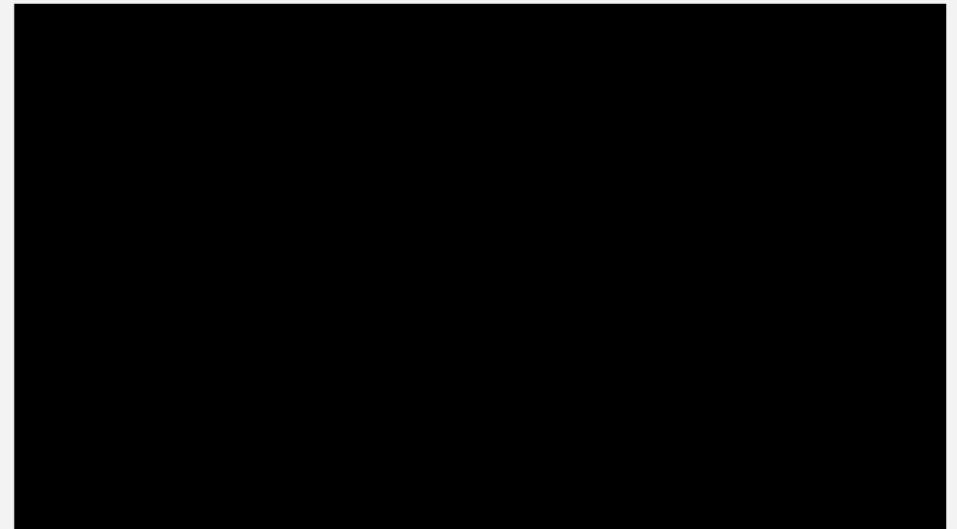
Abreu, M. et al., Learning low level skills from scratch for humanoid robot soccer using deep reinforcement learning, 19th IEEE International Conference on Autonomous Robot Systems and Competitions, ICARSC, 2019



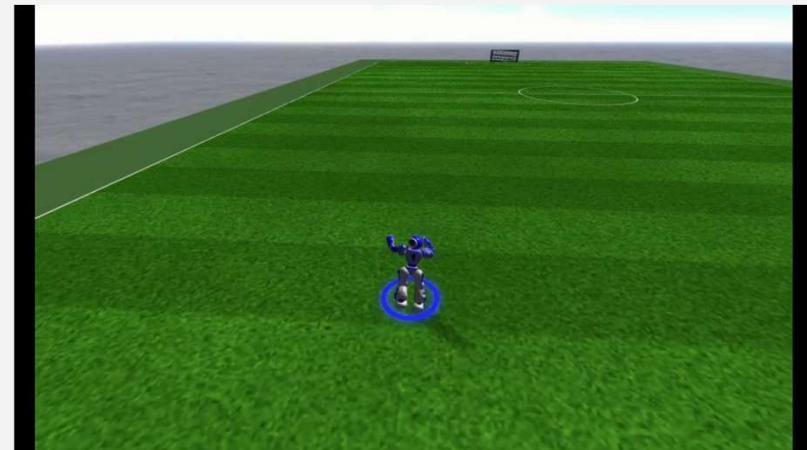
Turn and sprint



Run and kick



Dribble



Agile Run

Abreu, M. et al., Learning low level skills from scratch for humanoid robot soccer using deep reinforcement learning, 19th IEEE International Conference on Autonomous Robot Systems and Competitions, ICARSC, 2019

Behavior Development



Abreu, M. et al., Learning low level skills from scratch for humanoid robot soccer using deep reinforcement learning, 19th IEEE International Conference on Autonomous Robot Systems and Competitions, ICARSC, 2019

Conclusions

- Broad range of learning techniques applied to different areas of Robotics:
 - Perception
 - Behavior development (value based, Contextual policy search, Deep Learning)
 - Adapting Human-Robot Interfaces
 - Coordination of Robot teams
- Learning can be applied to Robotics, but:
 - Data should be used as efficiently as possible
 - Take advantage of data structure
 - Combine different approaches, if needed
 - Use simulation/ground-truth in the first learning steps

Robotics / Intelligent Robotics Robot Learning

Luís Paulo Reis, Nuno Lau, David Simões, Armando Sousa

lpreat@fe.up.pt

Director of LIACC – Artificial Intelligence and Computer Science Lab.

Associate Professor at DEI/FEUP – Informatics Engineering Department, Faculty of Engineering
of the University of Porto, Portugal

President of APPIA – Portuguese Association for Artificial Intelligence

