# Task and Motion Planning using Learning from Demonstration and Reinforcement Learning

**Dominik Urbaniak**

Master Thesis – Final Presentation

Dr. Alejandro Agostini

Human-centered Assistive Robotics

Technische Universität München

# Dirty Dishes

Washing dishes is one common task in our daily life

→ A dishwasher is helpful, but still requires to be loaded

**Requirements:**

Carefully picking and placing objects of various size and shape in a meaningful order.
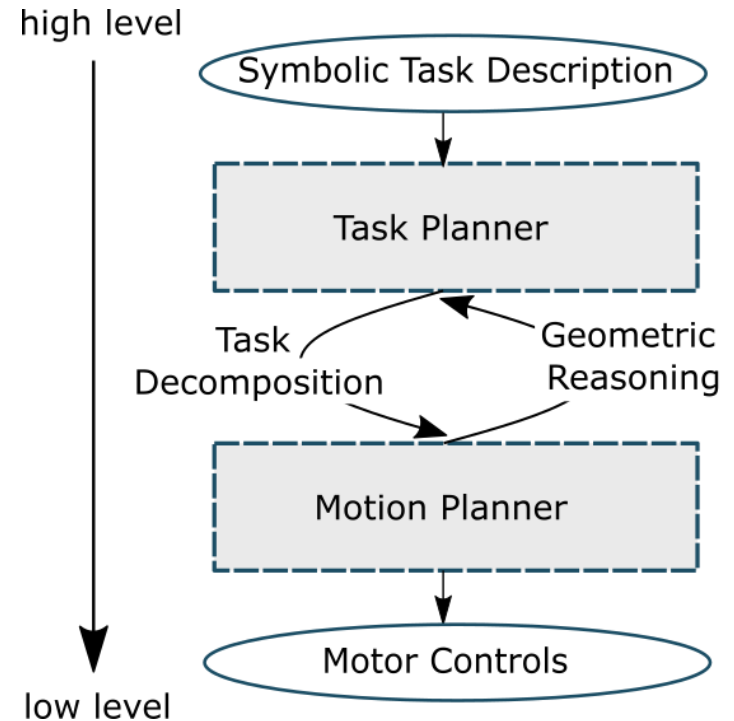
→ Task and motion planning problem



https://home.howstuffworks.com/dishwasher.htm

# Why Task and Motion Planning?

**Task Planner**

- Efficient planning with logic descriptions for long time horizons

**Motion Planner**

- Considers the detailed geometric specification of the environment

➡️ **Challenge:**
Transform symbolic actions into feasible motion in variable scenarios

# Related Work

→ Improve state-of-the art methods in terms of **efficiency**

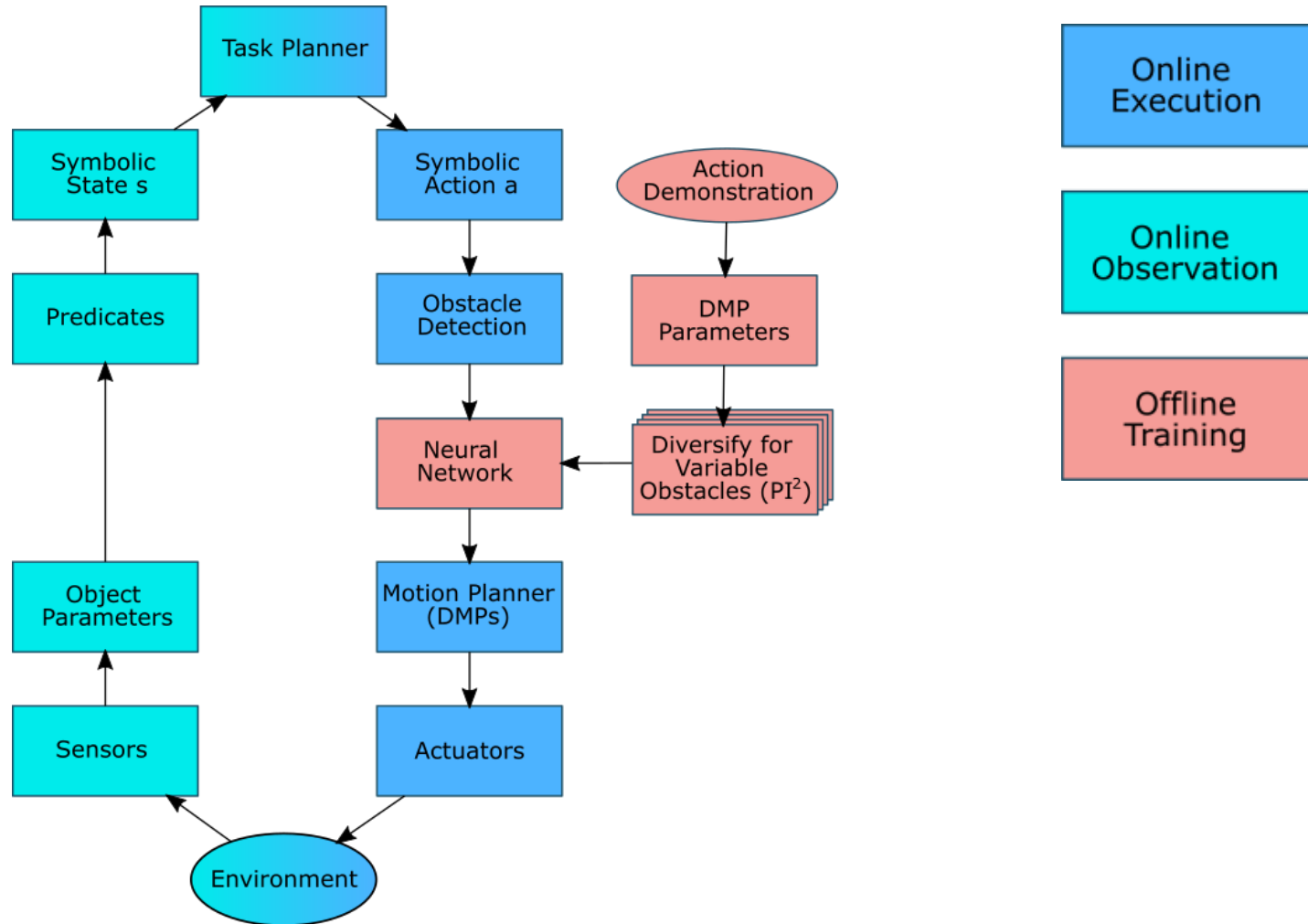| TAMP framework | Limitation | Proposed solution | Efficiency gain |
|---|---|---|---|
| Search-based [Dantam, 2018] [Bidot, 2017] | Deliberation before every execution | Learning-based approach | Faster at execution time |
| Learning-based using RL [Quack, 2015] | Only performed in low-dimensional spaces | Apply RL on LfD | Efficient motion generation with LfD |
| Learning-based using LfD [Agostini, 2020] | Many demonstrations required | Generalization with RL | Only one demonstration required |

# Contributions

- Our TAMP framework permits executing complex tasks comprising long action sequences with obstacle avoidance.

- Each symbolic action is grounded using DMPs that an action policy provides for variable object configurations.

- PI² efficiently generates divers sets of optimal collision-free trajectories serving as training samples to learn the action policy encoded in a NN.

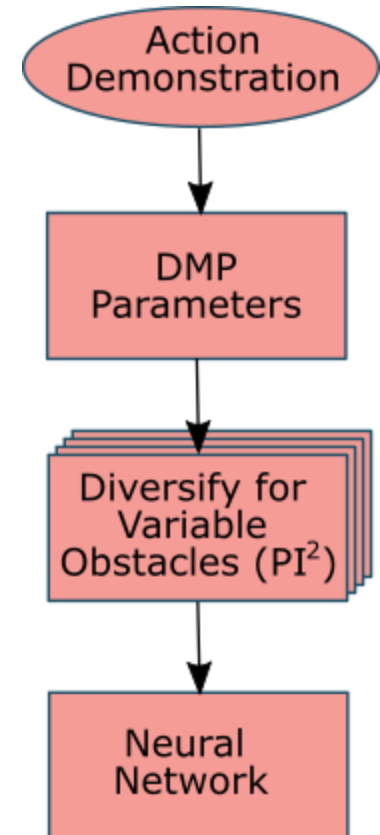# The Proposed TAMP Framework

# Training Methods

**Why DMPs?**

- Efficient motion encoding from one demo
- Translation, dilatation, rotation invariance
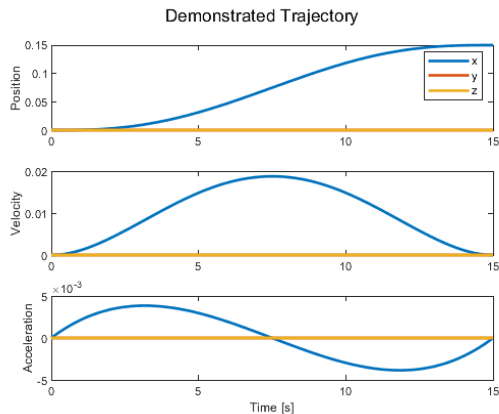
**Why PI²?**

- Tunes parameterized policies like DMPs
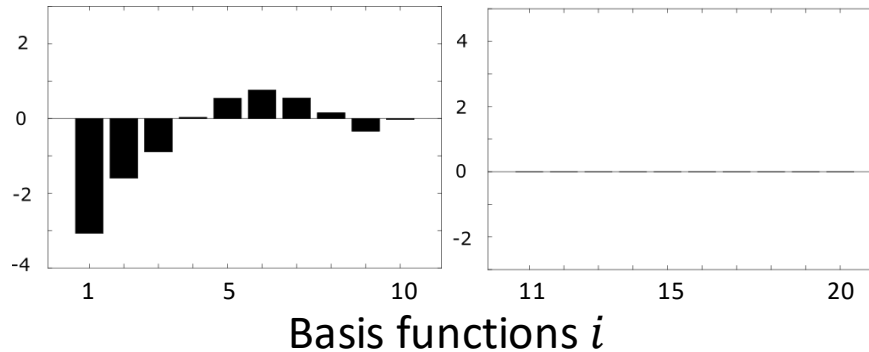- Numerically stable based on SOC

**Why NN?**

- Encode action policy for action selection
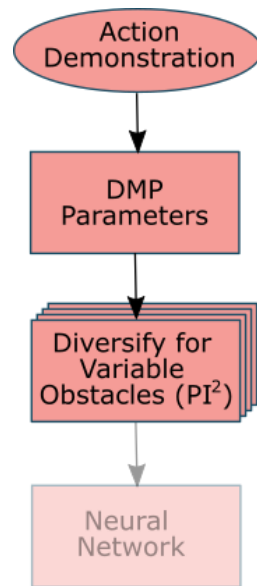
# Offline Training – Encode a Demonstration



DMP parameters $\theta_i^j$ of the forcing term at $j = 1$

Basis functions $i$

No ability to avoid obstacles
$\rightarrow$ Apply PI² iteratively to tune the DMP parameters $\theta_i$

Action Demonstration

DMP Parameters

Diversify for Variable Obstacles (PI²)

Neural Network

# Offline Training - PI² Optimization Step

The forcing term $\theta_i^j$ of the DMPs represents the current policy

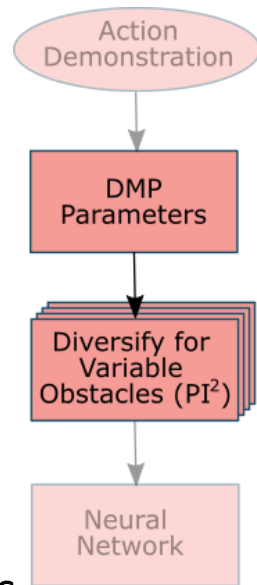1. Generate $K$ random samples from the current policy $\theta_i^j$

$$\theta_{i,k}^{j+1} = \theta_i^j + \epsilon_{i,k}^j$$
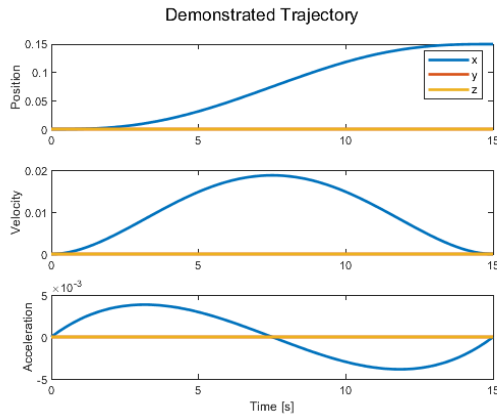
2. Weight each sample $k$ using a cost function $S$

$$W_\theta(\theta_{i,k}^j) = exp\left(-\gamma \frac{S(\theta_{i,k}^j) - \min S(\boldsymbol{\theta}_i^j)}{\max S(\boldsymbol{\theta}_i^j) - \min S(\boldsymbol{\theta}_i^j)}\right)$$

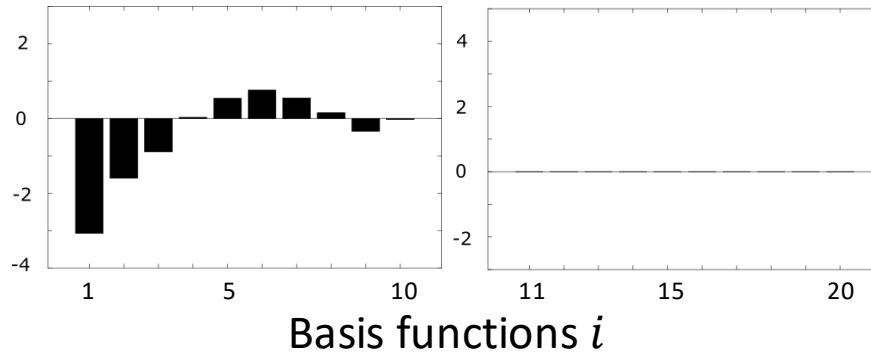3. Compute the new policy $\theta_i^{j+1}$ as weighted average of the samples

$$\theta_i^{j+1} = \frac{\sum_{k=1}^K W_\theta(\theta_{i,k}^j)\theta_{i,k}^j}{\sum_{k=1}^K W_\theta(\theta_{i,k}^j)}$$

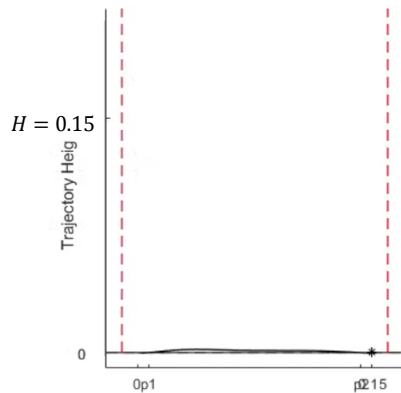Action Demonstration

DMP Parameters

Diversify for Variable Obstacles (PI²)

Neural Network

# Offline Training – Diversify the Demonstration



DMP parameters $\theta_i^j$ of the forcing term at $j = 1$

Basis functions $i$

$S = -\text{H} = \min(h_{p1}, h_{p2})$

$H = 0.15$

Tuned forcing term $\theta_i^j$ for $j = \{50, 100, \dots, J\}$

Basis functions $i$

4x

Action Demonstration

DMP Parameters

Diversify for Variable Obstacles (PI$^2$)

Neural Network

# Defintion of the Trajectory Shape

Two variables specify the trajectory shape $r_c, r_L$

Degree of curvature $\qquad\qquad r_c = H/l$

Trajectory steepness $\qquad r_L = (|\tilde{L}(2) - \tilde{L}(1)| + 1)/B$

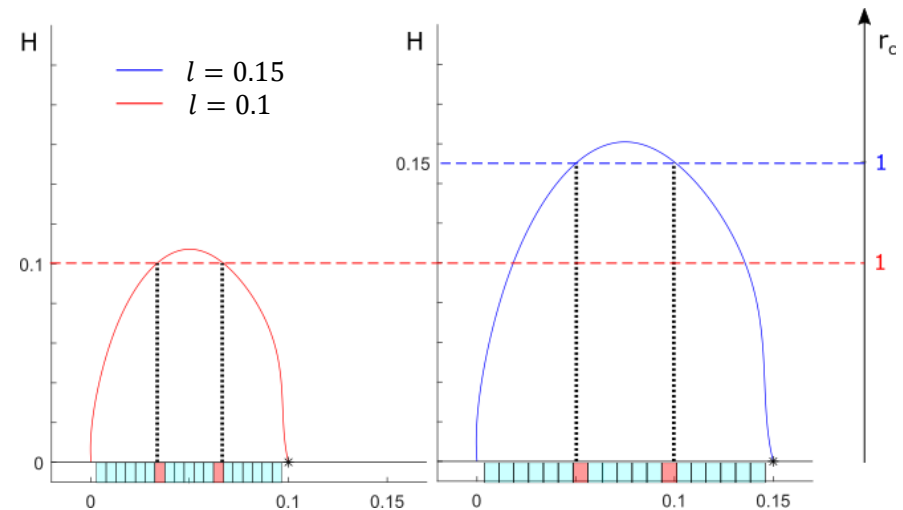→ $r_c, r_L$ are dilatation invariant



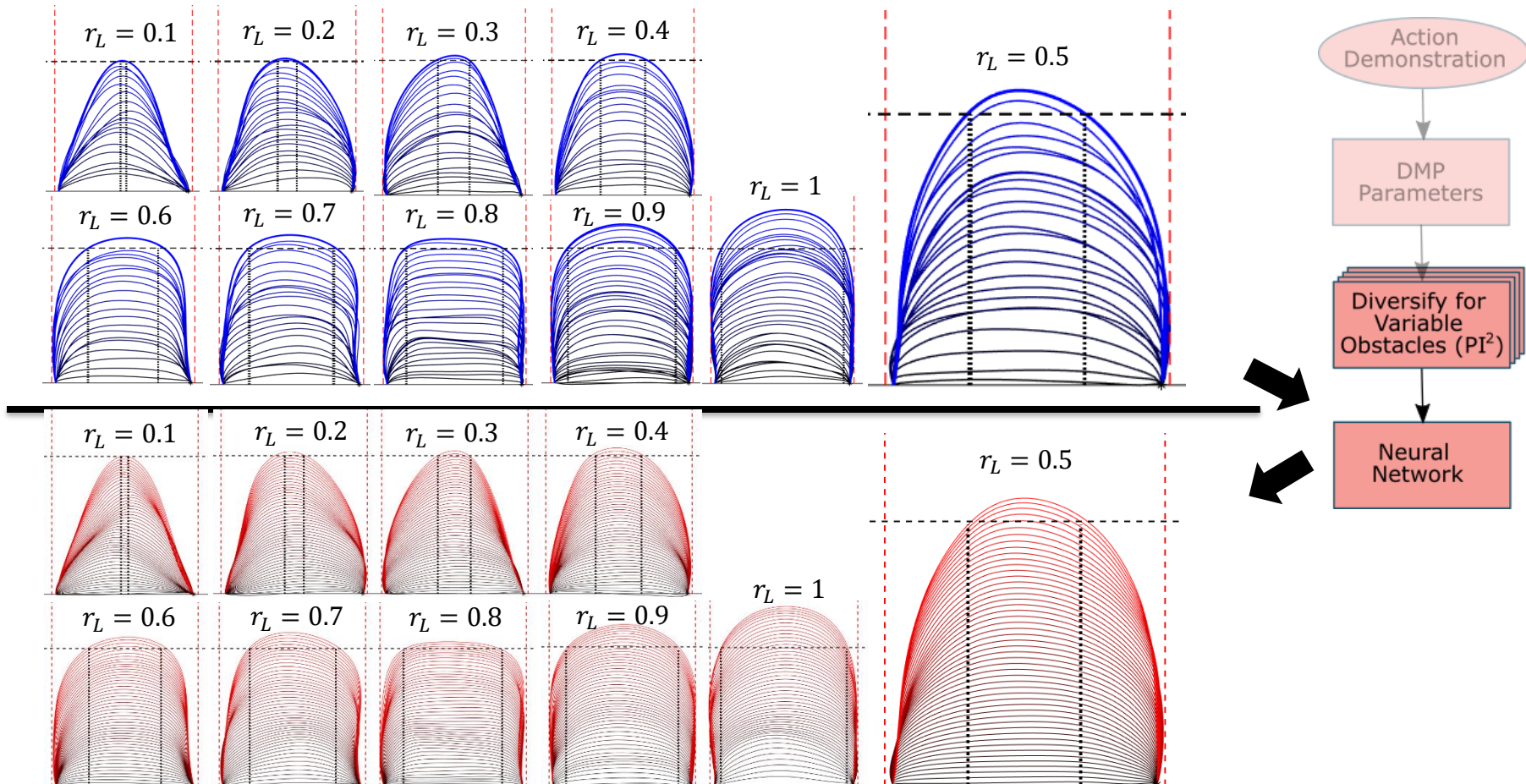Same set of forcing term parameters $\theta_i$
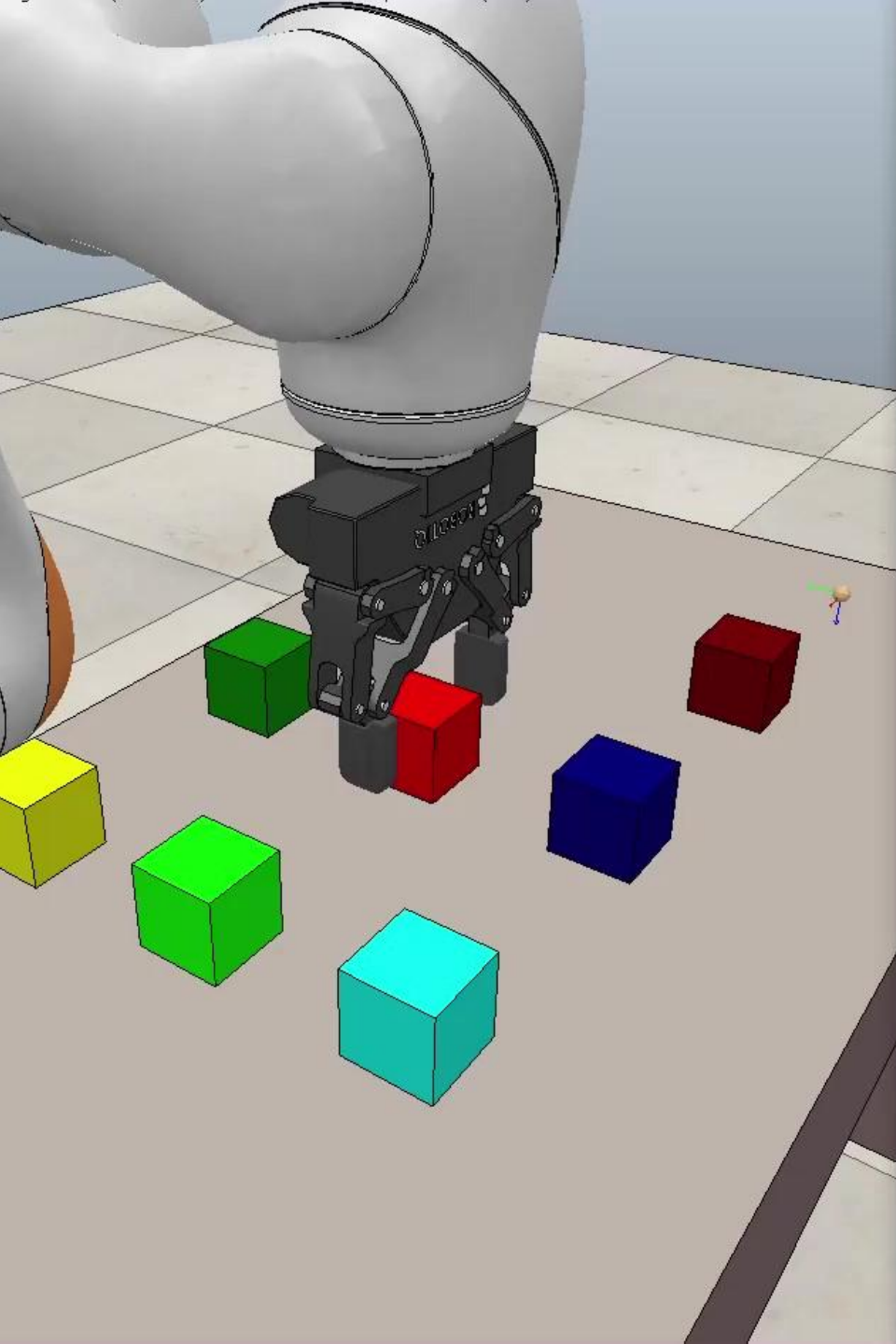
↕

Same trajectory shape
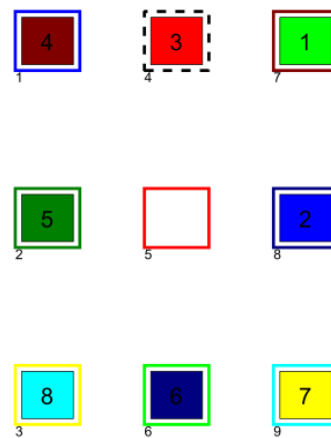
↕

Same pair of $r_c = 1, r_L = 0.4$

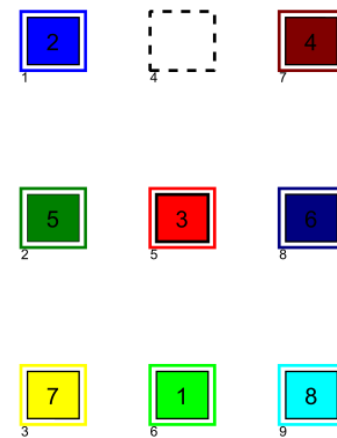# Offline Training – Learning Action Selection



→ Learned optimal parameters to ground symbolic actions in varying scenarios

**Initial State**

**Goal State**

Cube

Empty cell

# Task Planner

Decompose the task into a sequence of *pickplace* actions

## The *domain*

**constants**
```
air
```

**predicates**
```
(on ?cell ?cube)
```

**symbolic action**
```
(pickplace ?from ?to ?c)
```
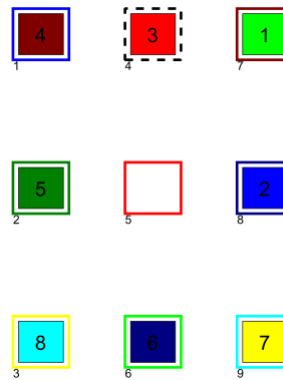
**precondition**
```
(and (on ?from ?c)
(on ?to air))
```

**effect**
```
(and (on ?from air)
(on ?to ?c)
(not (on ?from c))
(not (on ?to air))
```

**Initial State**



**Task Plan**

```
#1: pickplace cell1 cell5 cube4
#2: pickplace cell3 cell1 cube8
#3: pickplace cell9 cell3 cube7
#4: pickplace cell1 cell19 cube8
#5: pickplace cell8 cell1 cube2
#6: pickplace cell6 cell8 cube6
#7: pickplace cell7 cell6 cube1
#8: pickplace cell5 cell17 cube4
#9: pickplace cell14 cell5 cube3
```

## The *task*

**objects**
```
cube1 cube2 …
cube8
cell1 cell2 …
cell9
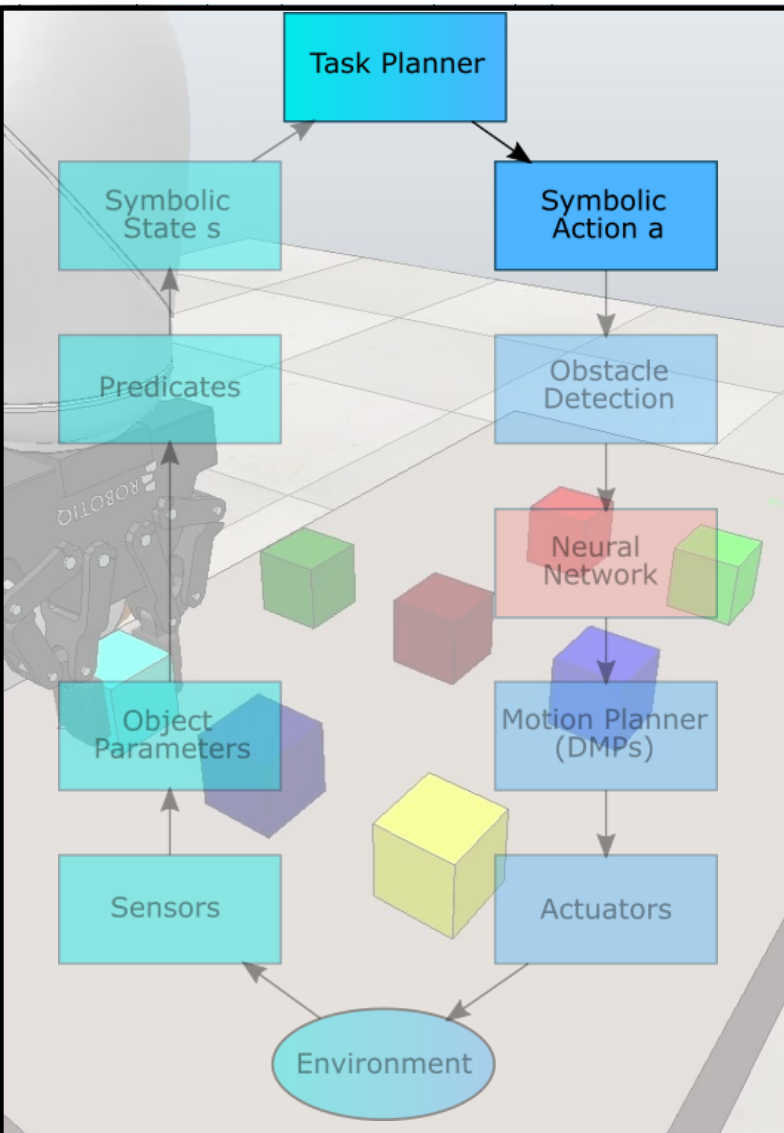```

**init**
```
(on cell5 air)
(on cell1 cube4)
…
```

**goal**
```
(on cell4 air)
(on cell2 cube5)
…
```

# Online Loop – Select Symbolic Action



**Current symbolic action:**
Pick and place *cube8*
from *cell3* to *cell1*

**Task Plan**
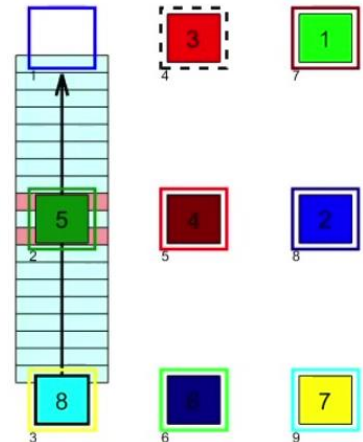
#1: pickplace cell1 cell5 cube4
#2: pickplace cell3 cell1 cube8
#3: pickplace cell9 cell3 cube7
#4: pickplace cell1 cell9 cube8
#5: pickplace cell8 cell1 cube2
#6: pickplace cell6 cell8 cube6
#7: pickplace cell7 cell6 cube1
#8: pickplace cell5 cell7 cube4
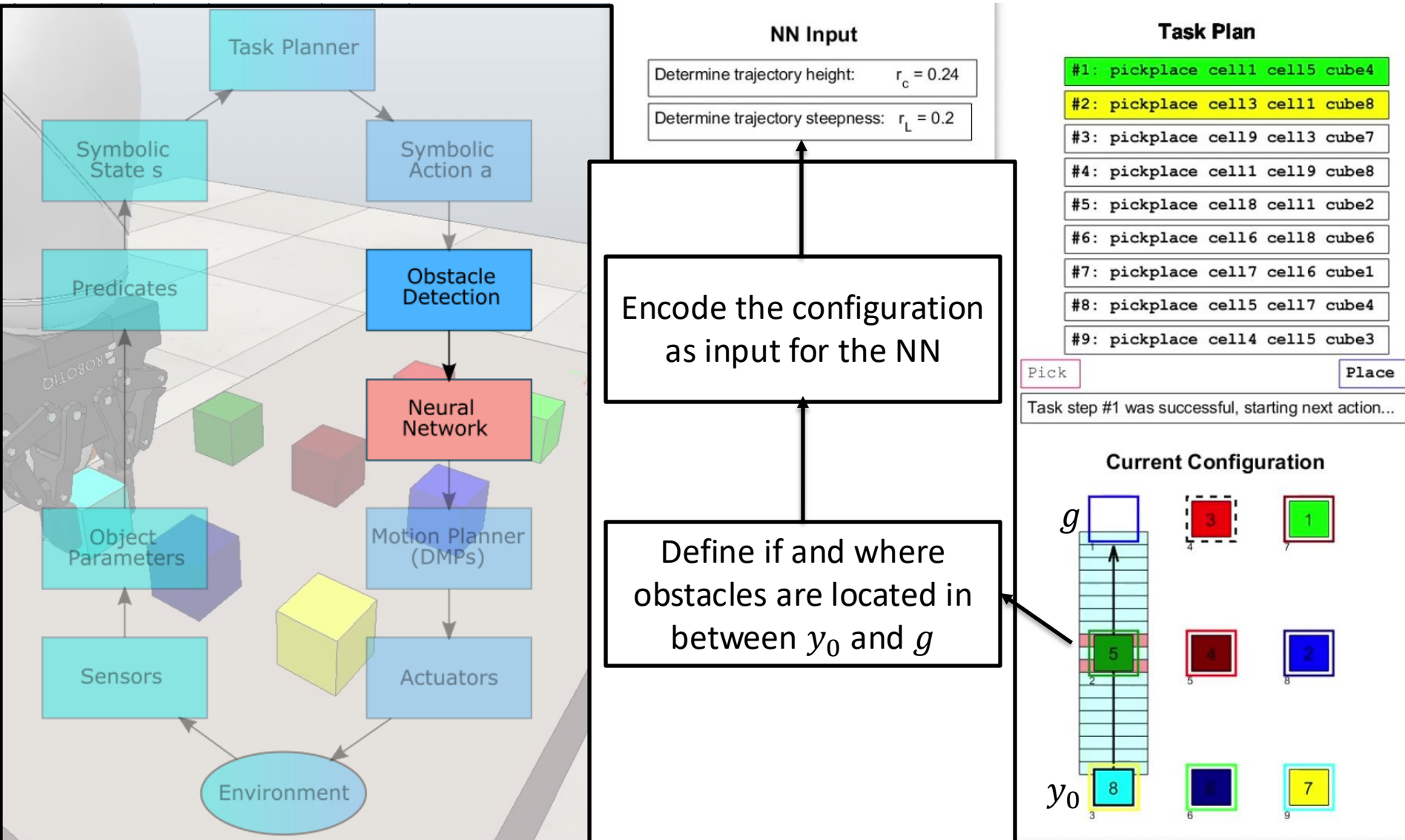#9: pickplace cell4 cell5 cube3

Pick                           Place

Task step #1 was successful, starting next action…
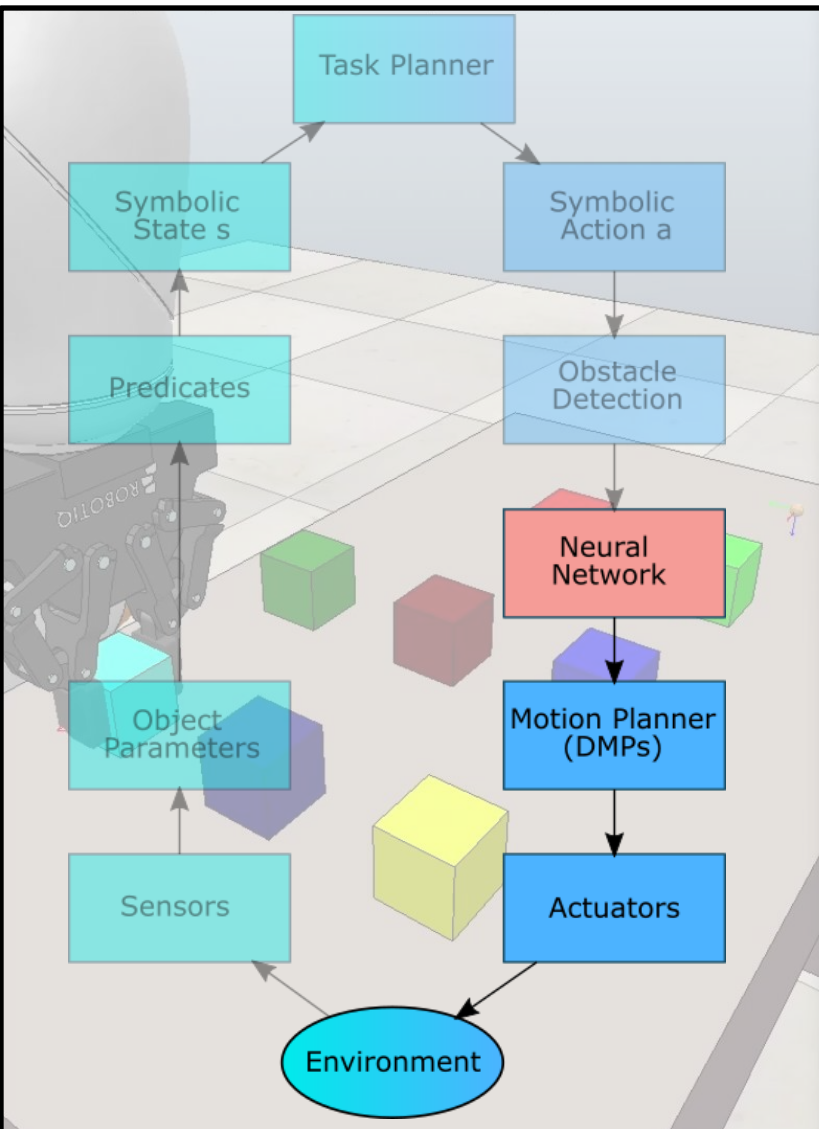
**Current Configuration**

# Online Loop – Encode Current Configuration



**NN Input**

Determine trajectory height:  $r_c = 0.24$

Determine trajectory steepness:  $r_L = 0.2$

Encode the configuration as input for the NN

Define if and where obstacles are located in between $y_0$ and $g$

**Task Plan**

#1: pickplace cell1 cell5 cube4
#2: pickplace cell3 cell1 cube8
#3: pickplace cell9 cell3 cube7
#4: pickplace cell1 cell9 cube8
#5: pickplace cell8 cell1 cube2
#6: pickplace cell6 cell8 cube6
#7: pickplace cell7 cell6 cube1
#8: pickplace cell5 cell7 cube4
#9: pickplace cell4 cell5 cube3

Pick                          Place

Task step #1 was successful, starting next action...

**Current Configuration**

$g$

$y_0$

Task Planner

Symbolic State s

Symbolic Action a

Predicates

Obstacle Detection

Neural Network

Object Parameters

Motion Planner (DMPs)

Sensors

Actuators

Environment

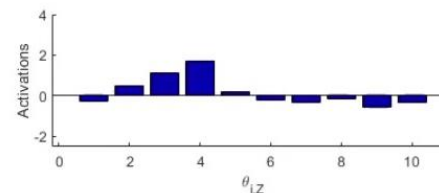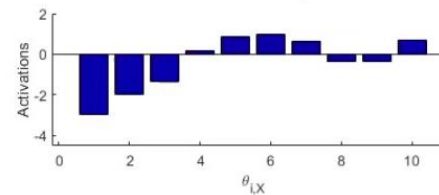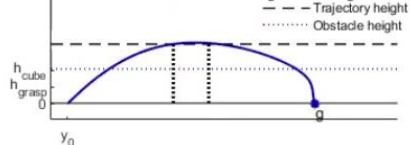# Online Loop – Execute Action Policy



**NN Input**

Determine trajectory height: $r_c = 0.24$

Determine trajectory steepness: $r_L = 0.2$

**NN Output**

**Generated Trajectory**

Pick    Place
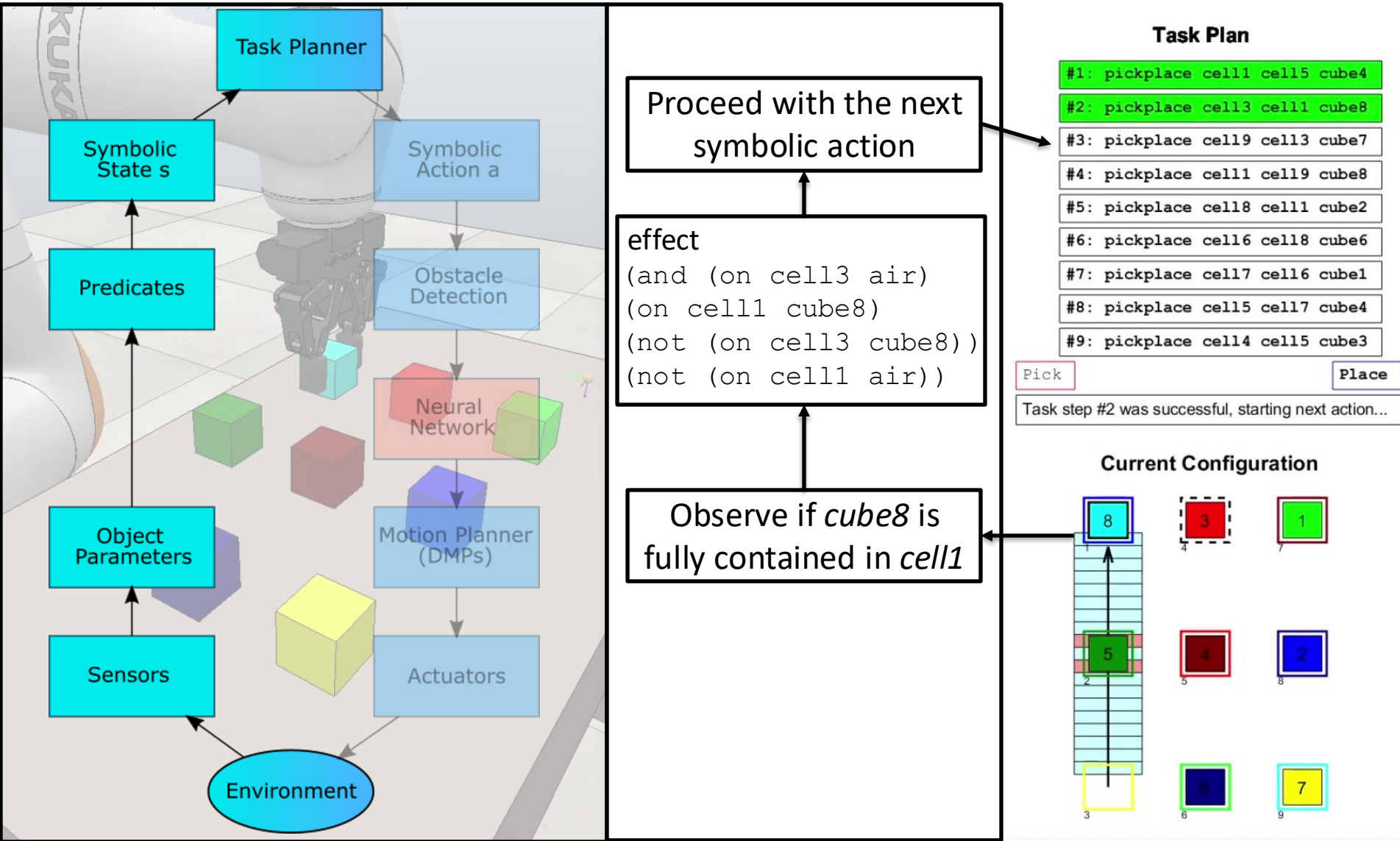
**Trajectory Precision**

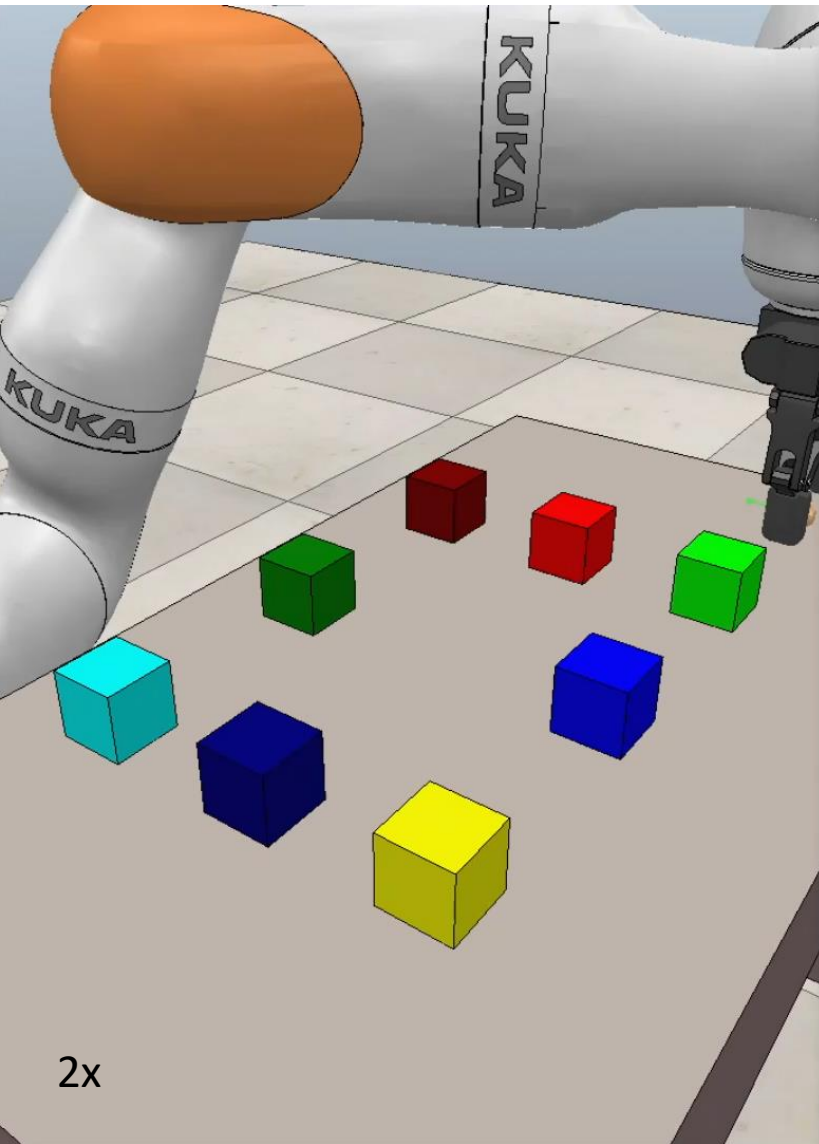Goal deviation: $d_g = 0.02\%$

Height deviation: $d_H = -0.5\%$

Retrieve and set the forcing term parameters of the DMPs

Generate the trajectory and send it to the robot arm

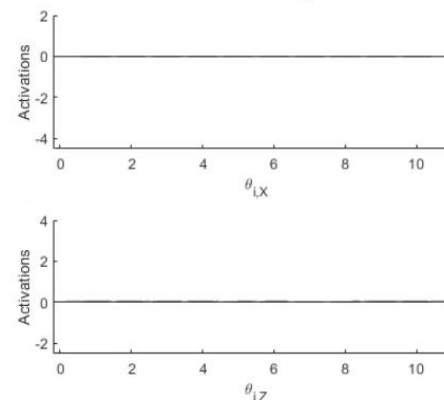# Online Loop – Observe Changes to the Environment



**Task Planner**

**Symbolic State s**

**Symbolic Action a**

**Predicates**

**Obstacle Detection**

**Neural Network**

**Object Parameters**

**Motion Planner (DMPs)**

**Sensors**

**Actuators**

**Environment**

Proceed with the next symbolic action

```
effect
(and (on cell3 air)
(on cell1 cube8)
(not (on cell3 cube8))
(not (on cell1 air))
```

Observe if *cube8* is fully contained in *cell1*

**Task Plan**

| #1: pickplace cell1 cell5 cube4 |
| #2: pickplace cell3 cell1 cube8 |
| #3: pickplace cell19 cell3 cube7 |
| #4: pickplace cell1 cell19 cube8 |
| #5: pickplace cell8 cell1 cube2 |
| #6: pickplace cell16 cell18 cube6 |
| #7: pickplace cell7 cell16 cube1 |
| #8: pickplace cell5 cell17 cube4 |
| #9: pickplace cell4 cell15 cube3 |

Pick                                        Place

Task step #2 was successful, starting next action...

**Current Configuration**

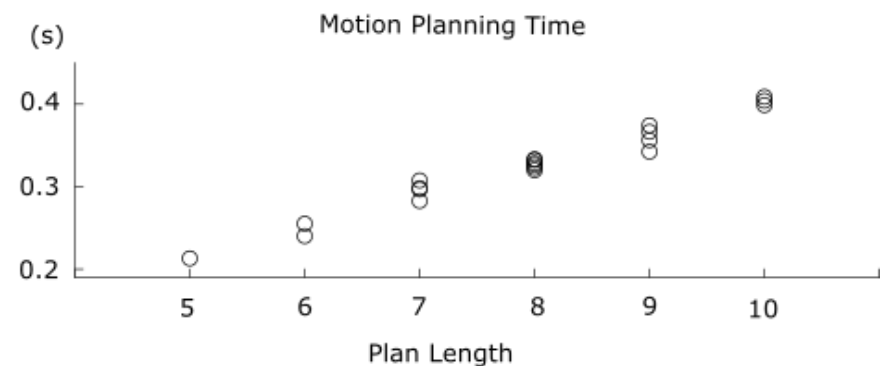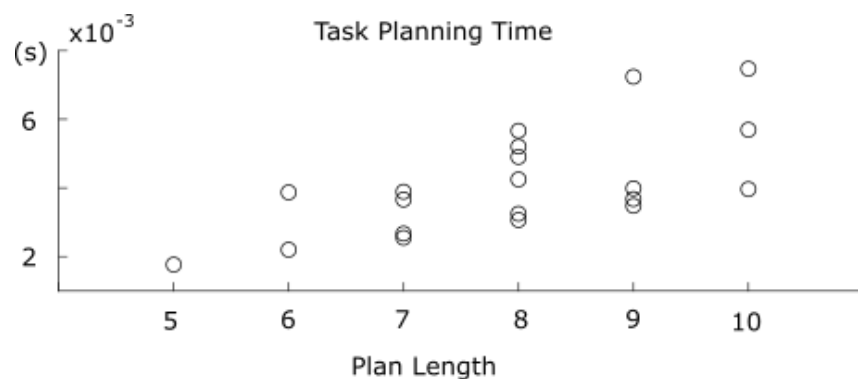# Complete Execution of a Task Example



2x

# Experimental Evaluation

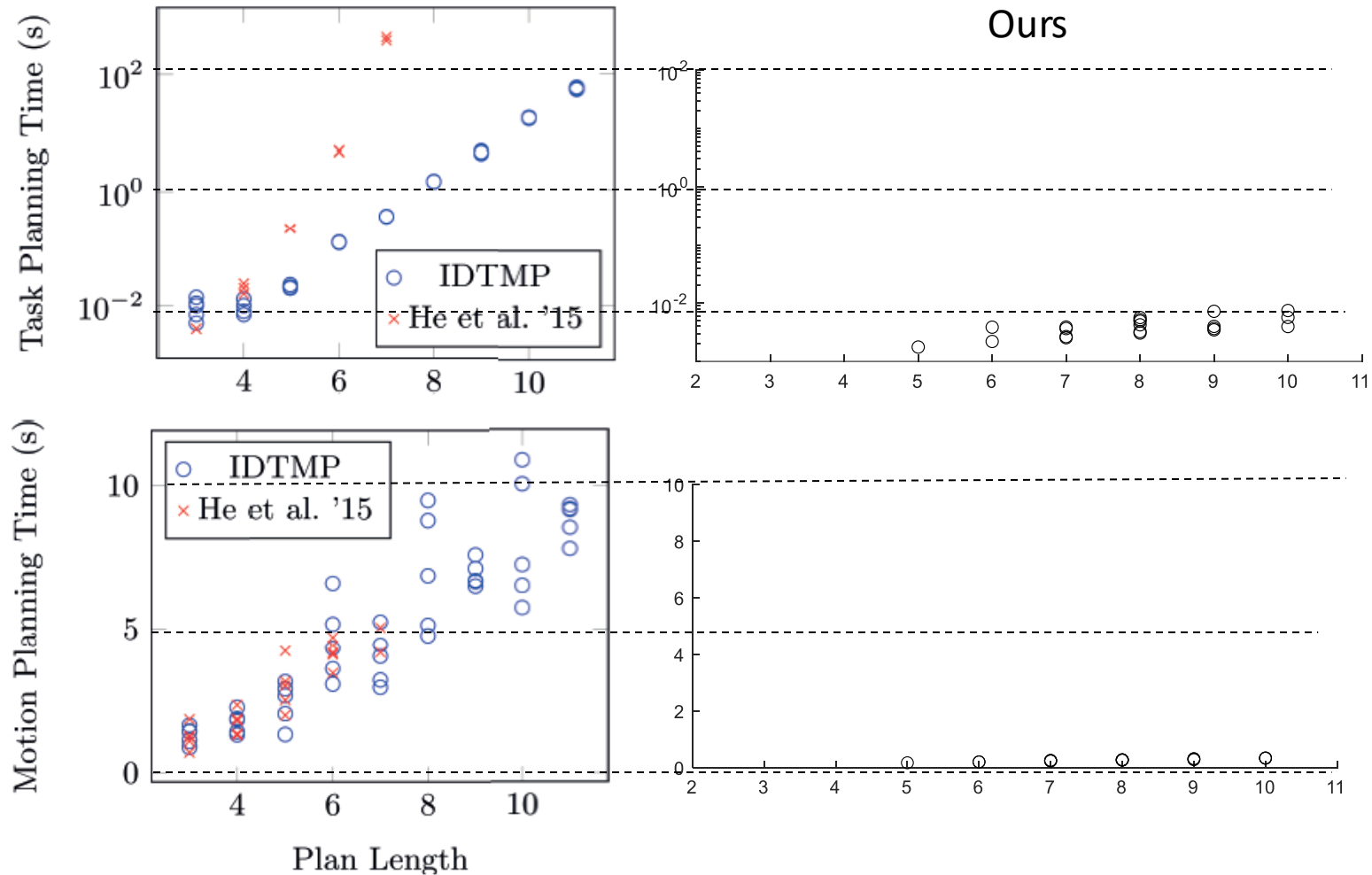20 consecutive task executions with random initial and goal state

→ All 159 performed *pickplace* actions successful

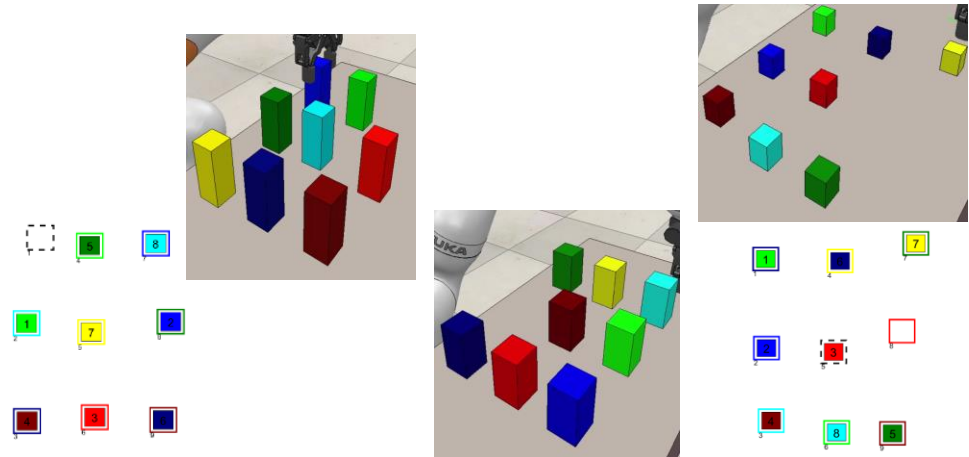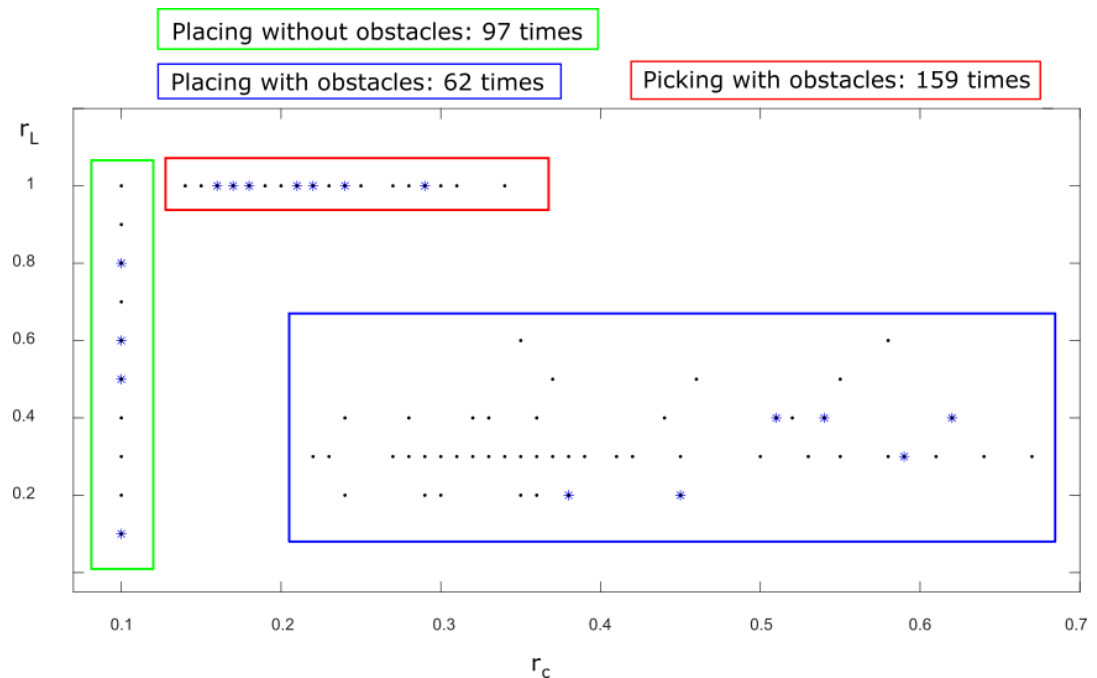| Training Times | 9 min |
|---|---|
| Generating the demonstration | 0.02 s |
| PI² optimizations | 241 s |
| Training of the NN | 279 s |

# Evaluation – Compared to [Dantam, 2018]
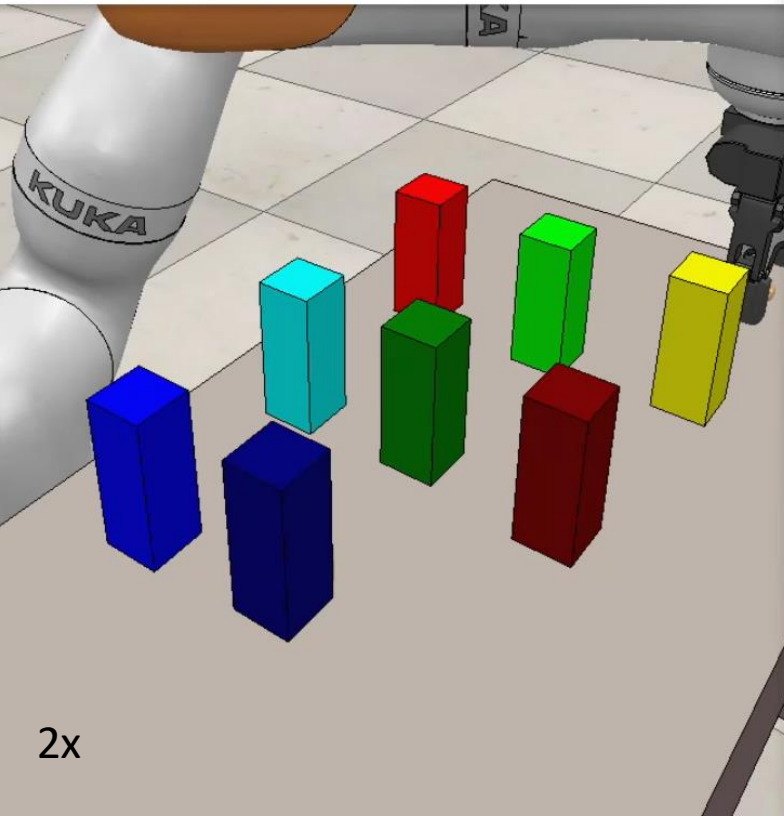
# Generalization Ability

Varying block dimensions
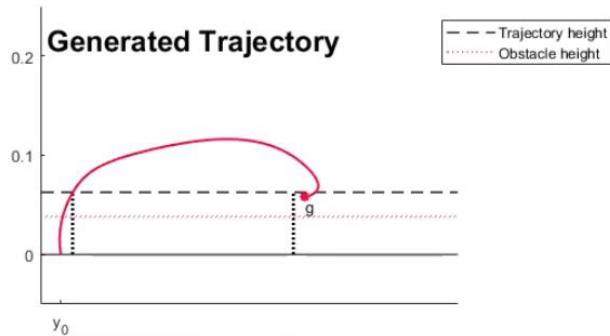
Varying cell positions

→ Learned action policy finds collision-free trajectories for all scenarios
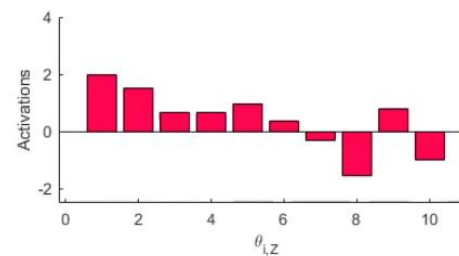
→ More trajectory shapes required



Placing without obstacles: 97 times

Placing with obstacles: 62 times

Picking with obstacles: 159 times

# Complete Execution of a Task Example



**Generated Trajectory**
- Trajectory height
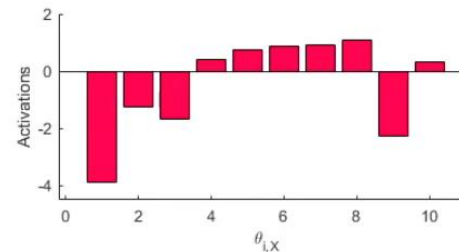- Obstacle height

**NN Input**

Determine trajectory height:     $r_c = 0.26$

Determine trajectory steepness:     $r_L = 1$

**NN Output**

Pick    Place

**Trajectory Precision**

Goal deviation:     $d_g = 0.92\%$

Height deviation:     $d_H = -7.39\%$

**Task Plan**

#1: pickplace cell1 cell9 cube3
#2: pickplace cell2 cell1 cube8
#3: pickplace cell14 cell2 cube1
#4: pickplace cell6 cell4 cube6
#5: pickplace cell5 cell6 cube5
#6: pickplace cell7 cell5 cube7
#7: pickplace cell1 cell7 cube8
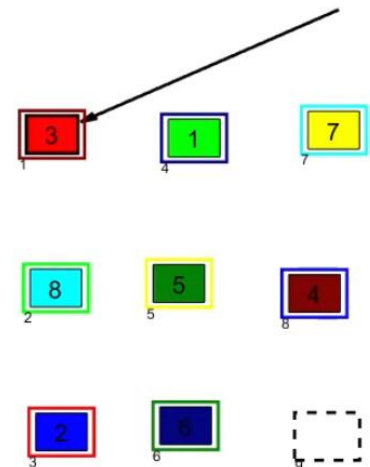#8: pickplace cell8 cell1 cube4
#9: pickplace cell3 cell18 cube2
#10: pickplace cell19 cell3 cube3

Pick                                    Place

Starting the Task...
**Current Configuration**

2x

# Limitation & Future Work

**Limitations**

1. Restricted to intial and goal positions in the same plane
2. Constant gripper orientation
3. Obstacle avoidance in one dimension only
4. Heuristics required for obstacle definition and optimization costs

**Future Work**

→ Train NN on independent data to improve generalization
→ Let RL agent learn to select appropriate input parameters and goal poses

# Conclusion

- TAMP framework that utilizes LfD to efficiently generate motion and RL to generalize this motion

- $PI^2$ permits learning DMP parameters from a single demonstration to avoid obstacles of varying size and in varying situations

- After a few minutes of training, the action policy reliably selects collision-free trajectories to ground symbolic actions of a complex task

# References

J. Bidot, L. Karlsson, F. Lagriffoul, A. Saffiotti
**Geometric backtracking for combined task and motion planning in robotic systems.**
In: *Artificial Intelligence*, 2017, pp. 229–265.

N. T. Dantam, Z. K. Kingston, S. Chaudhuri, L. E. Kavraki
**An incremental constraint-based framework for task and motion planning.**
In: The International Journal of Robotics Research*, 2018, pp. 1134–1151*

B. Quack, F. Wörgötter, A. Agostini
**Simultaneous learning at different levels of abstraction.**
In: *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)* , 2015, pp. 4600–4607.

A. Agostini, M. Saveriano, D. Lee, J. Piater
**Manipulation planning using object-centered predicates and hierarchical decomposition of contextual actions.**
In: *IEEE Robotics and Automation Letters*, 2020, pp. 5629–5636.

# Appendix

# Cost Function



$$S = -H + c_1 \cdot S_{prec} + c_2 \cdot S_{scope}$$

$$H = \min(h_{p1}, h_{p2})$$

$$S_{prec} = \|g - y_{end}\|$$

$$S_{scope} = -\sum_{t=1}^{T} \min(0, m + y_{X,t} - y_{0,X}) - \sum_{t=1}^{T} \min(0, m + g_X - y_{X,t})$$
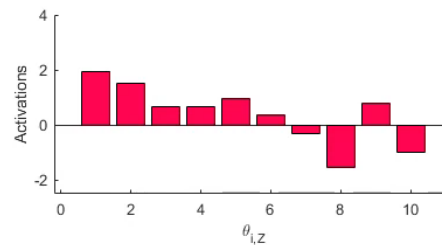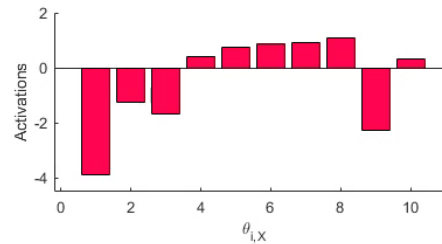
## Generated Trajectory

- - - Trajectory height
........ Obstacle height

## NN Input

Determine trajectory height: $r_c = 0.24$

Determine trajectory steepness: $r_L = 1$

## NN Output

Activations vs $\theta_{i,X}$

Activations vs $\theta_{i,Z}$

Pick    Place

## Trajectory Precision
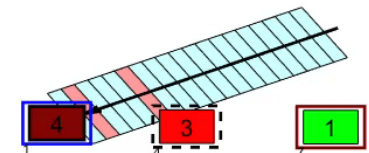
Goal deviation: $d_g = 0.78\%$

Height deviation: $d_H = -6.29\%$

## Task Plan

#1: pickplace cell1 cell5 cube4
#2: pickplace cell3 cell1 cube8
#3: pickplace cell9 cell3 cube7
#4: pickplace cell1 cell19 cube8
#5: pickplace cell8 cell1 cube2
#6: pickplace cell6 cell8 cube6
#7: pickplace cell7 cell6 cube1
#8: pickplace cell5 cell17 cube4
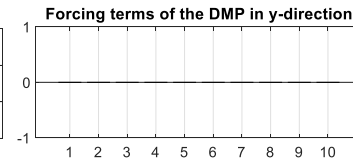#9: pickplace cell4 cell15 cube3

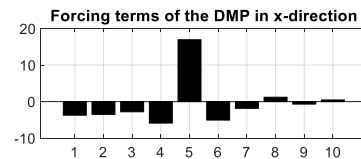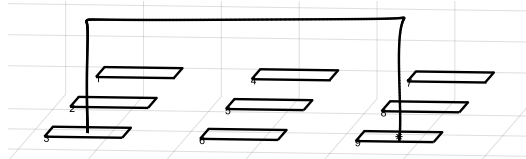Pick    Place

Starting the Task...

## Current Configuration

4   3   1

5       2

8   6   7

# DMP: Roto-Dilatation Invariance [Ginesi, 2019]
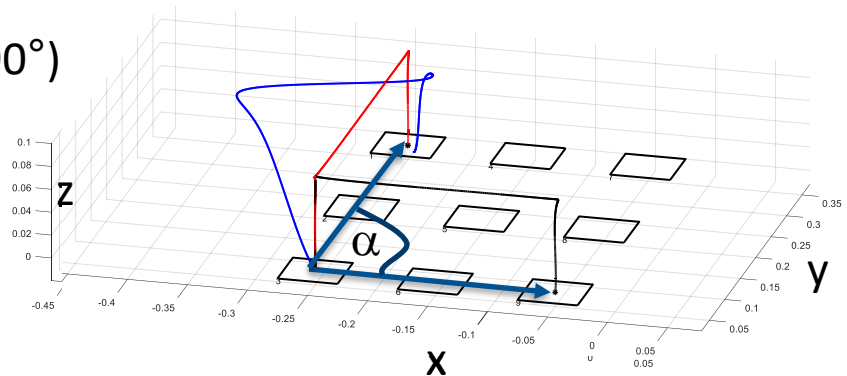
One demonstration encoded in three DMPs



New goal position (rotated around z: $\alpha$=90°)

Same forcing term parameters
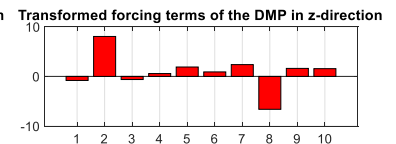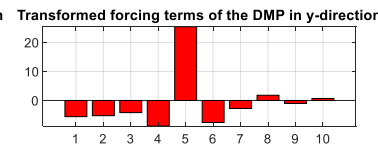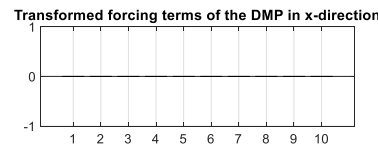→ Unexpected trajectory shape
→ Goal position is not reached precisely



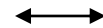Rotate the forcing term parameters

$$f_x^{new} = \cos \alpha * f_x^{demo}$$
$$f_y^{new} = \sin \alpha * f_y^{demo}$$

$$\tilde{L} = (7,14) \qquad \tilde{L} = (7,14)$$