

Guiding Robot Exploration in Reinforcement Learning via Automated Planning

Yohei Hayamizu^{†*}, Saeid Amiri[‡], Kishan Chandan[‡], Keiki Takadama[†], and Shiqi Zhang[‡]

†: University of Electro-Communications, ‡: SUNY Binghamton, *Work done while visiting SUNY Binghamton

Contact E-mail: hayamizu@cas.lab.uec.ac.jp



Background

Robots in human-inhabited environments are able to conduct various service and interaction tasks.



A Segway-based mobile robot platform used in this research

Requirements for these tasks

- Fulfilling requests from humans such as navigation and delivery
- Learning efficiency in the real world

Research fields for these tasks

Model-Based Reinforcement Learning (RL) and Automated Planning have been used to meet these two requirements respectively (not jointly).

Model-Based RL

Learn a world model while learning an action policy to achieve long-term goals from both real and simulated experiences

Automated Planning

Reason with declarative domain knowledge, including commonsense knowledge, that is provided a priori

Aim

Efficient task learning to fulfill diverse service requests

Main Contributions

Efficient exploration strategy for RL agents in navigation domains

Avoid less-relevant states by reasoning with contextual knowledge while using trial-and-error experiences

Exploiting complementary features of model-based RL and automated planning

Aim at improving sample-efficiency in a real robot domain

Guided Dyna-Q (GDQ):

Bridging the gap between model-based RL and automated planning

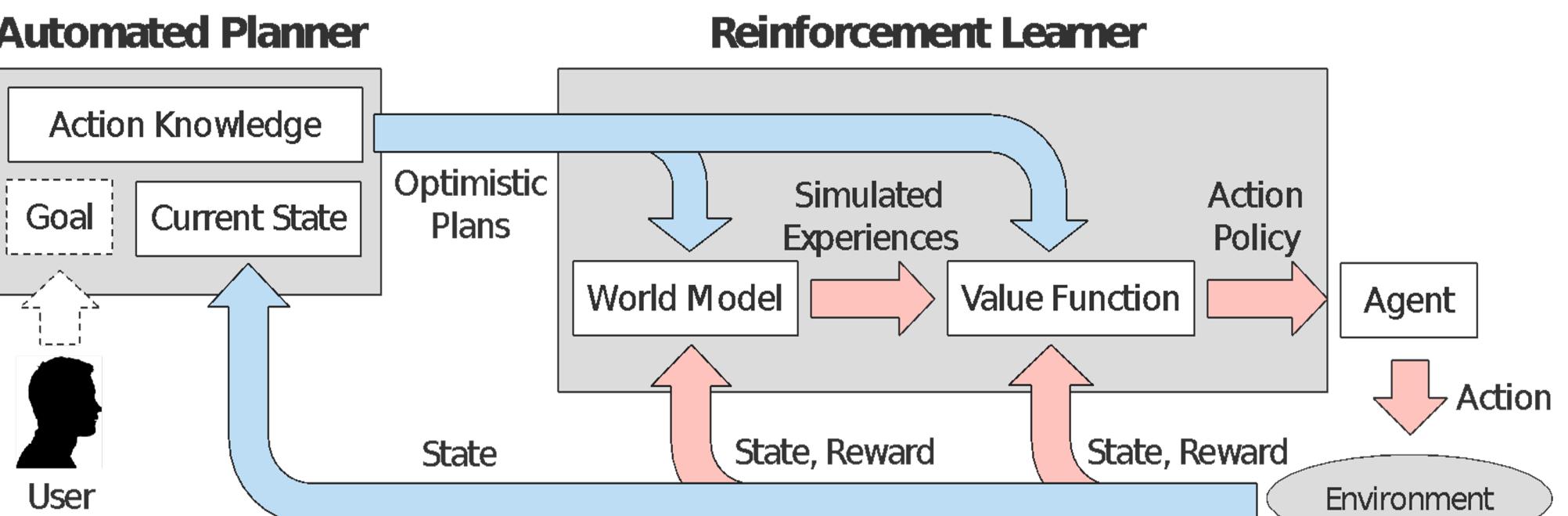
GDQ integrates the two sub-procedures for **optimistic initialization** and repeatedly conducting runtime **policy update**.

1. Optimistic Initialization

Help the agent avoid the near-random exploration behaviors through a “warm start” enabled by our automated planner

1. Policy Update

Guide the agent to only try the actions that can potentially lead to the goal states



An overview of Guided Dyna-Q (GDQ)

Experiments using Navigation Tasks

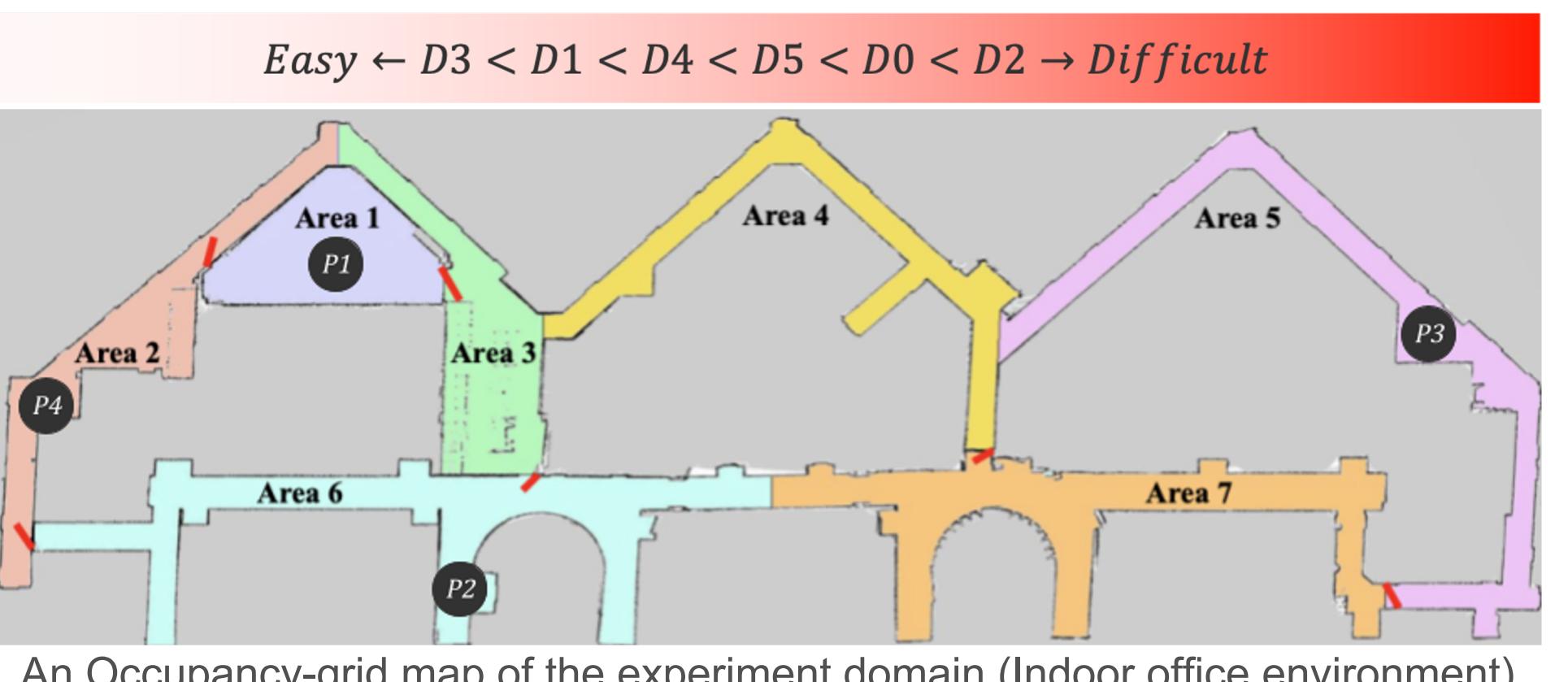
Hypotheses about GDQ:

- Perform better than existing RL methods from the literature in cumulative reward.
- Help the robot avoid visiting “irrelevant” areas. (A navigation task is achieved via relevant areas.)
- Is more robust to goal changes

Experiment Settings

Indoor office environment settings:

- All states are categorized into 7 areas.
- There are 6 doors that a robot can use to enter rooms.



An Occupancy-grid map of the experiment domain (Indoor office environment)

Action sets: 4 types of actions for navigational purposes

{*goto*, *gothrough*, *approach*, *opendoor*}

Action knowledge (designed by a human expert)

Answer Set Programming (ASP) [Lifschitz, 2002]:

- Performing well in knowledge intensive domain
- Computing plans with reasoning paradigm

Example of action knowledge representation by ASP

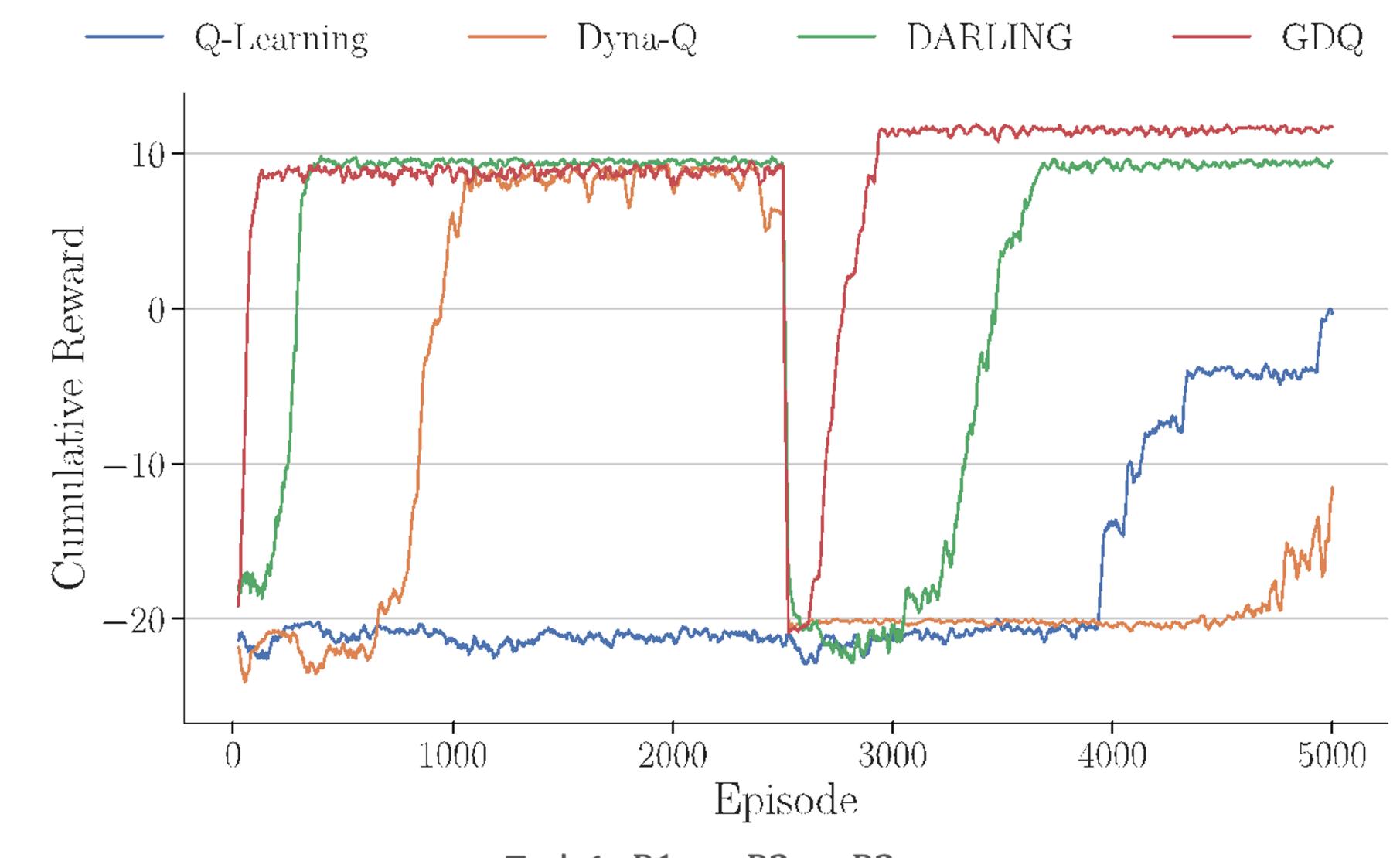
```
at(Z,I + 1) :- gothrough(Y,I), at(X,I), acc(X,Y,Z), I < n.  
hasdoor(s1,d0). acc(s0,s1).  
:- approach(D,I), facing(D,I), door(D), I = 0..n - 1.
```

Results

A) Simulation Experiment

Tasks:

- Task 1: $(P1 \rightarrow P2 \rightarrow P3)$



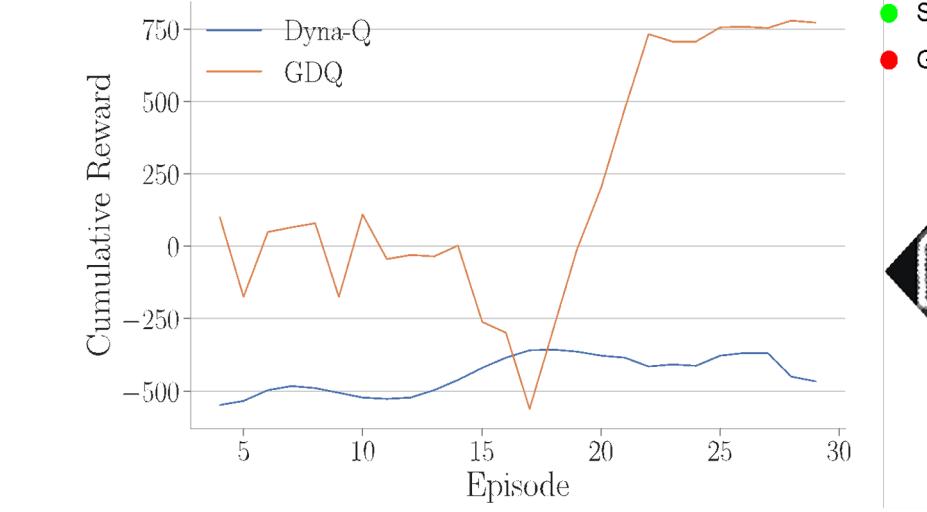
Task 1: $P1 \rightarrow P2 \rightarrow P3$

Average cumulative rewards over ten runs: GDQ learned the practical policy faster and is more robust to task changes.

B) Real Robot Experiment

Task:

- Task 2: $(P4 \rightarrow P2)$



Cumulative rewards on a real robot.

GDQ enabled robot to find the practical path in 27 trials.

Baseline Method:

- Dyna-Q



Heatmaps of our office domain for visualizing where the robot visited using the Dyna-Q Baseline (Left) and GDQ (Right).

Conclusion

Aim: Efficient task learning to fulfill diverse service requests

Guided Dyna-Q (GDQ): Optimistic Initialization & Policy Update

Results: GDQ improves the learned policies’ quality, and reduces exploration efforts

Acknowledgement

A portion of this work has taken place in the Autonomous Intelligent Robotics (AIR) Group at SUNY Binghamton. AIR research is supported in part by grants from the National Science Foundation (IIS-1925044 and REU Supplement), Ford Motor Company (two URP Awards), OPPO (Faculty Research Award), and SUNY Research Foundation. The authors thank collaborators on related researches that fed into the development of the ideas described in this paper.