# Partially Observable Markov Decision Processes for Planning in Uncertain Environments

Project Thesis
Lasse Peters

Institute of Mechanics and Ocean Engineering
Prof. Dr.-Ing. Robert Seifried
Hamburg University of Technology

Hybrid Systems Laboratory
Prof. Claire J. Tomlin
University of California at Berkeley

25.07.2019

**TUHH**
*Hamburg University of Technology*

**mum**
Mechanik und Meerestechnik

Hybrid Systems
Laboratory

**Berkeley**
UNIVERSITY OF CALIFORNIA

# Introduction



**Types of Uncertainty**
- state uncertainty
- outcome uncertainty

# Introduction



**Types of Uncertainty**
- state uncertainty
- outcome uncertainty

**Dealing with Uncertainty**
- worst case disturbance sequences
- probabilistic reasoning

**Types of Uncertainty**
- state uncertainty
- outcome uncertainty

**Dealing with Uncertainty**
- worst case disturbance sequences
- probabilistic reasoning

**POMDPs**
- general framework
- computationally demanding

# Introduction



**Types of Uncertainty**
- state uncertainty
- outcome uncertainty

**Dealing with Uncertainty**
- worst case disturbance sequences
- probabilistic reasoning

**POMDPs**
- general framework
- computationally demanding
- recent research: faster solvers

**This Work**
- problem specific approximations

**This Work**
- problem specific approximations
  vs. full POMDP approaches

# Introduction



**This Work**

- problem specific approximations vs. full POMDP approaches

- application domains:

  1. localization and planning

  2. motion planning with latent human intentions

# Introduction



**This Work**

- problem specific approximations vs. full POMDP approaches

- application domains:

  1. localization and planning

  2. motion planning with latent human intentions

# Outline

# The Partially Observable Markov Decision Process

**Dynamic Decision Network**



**MDP**

- state $s$
- action $a$
- reward $r$

**Objective**: Finding a policy $\pi^*$ that maximizes

$$J(\pi) = E\left[\sum_{t=0}^{\infty} \gamma^t r_t\right].$$

# The Partially Observable Markov Decision Process
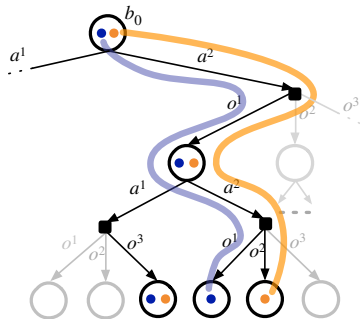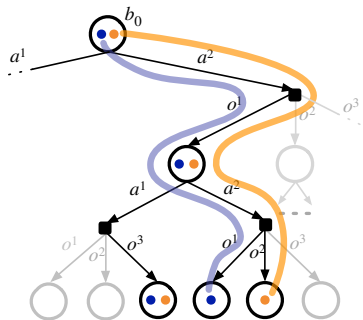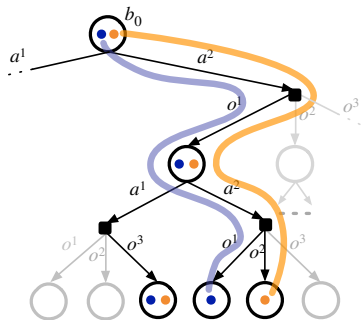
**Dynamic Decision Network**



**POMDP**

- state $s$
- action $a$
- reward $r$
- observation $o$
- initial belief $b_0$

**Objective**: Finding a policy $\pi^*$ that maximizes

$$J(\pi) = E\left[\sum_{t=0}^{\infty} \gamma^t r_t\right].$$

**High Level Idea**

- incrementally construct sparse tree
- maintain bounds on *value* at nodes
- choose action with best lower bound

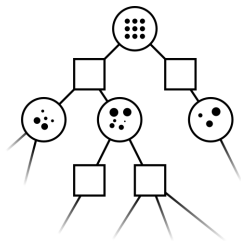# Determinized Sparse Partially Observable Tree (DESPOT)



## High Level Idea

- incrementally construct sparse tree
- maintain bounds on *value* at nodes
- choose action with best lower bound

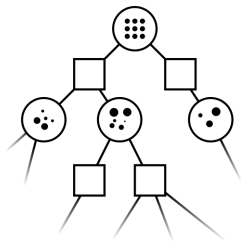## Domain Knowledge

- user-defined initial bound estimates

# Partially Observable Monte-Carlo Planning with Observation Widening (POMCPOW)



**High Level Idea**

- Monte-Carlo tree search
- locally approximate the *value function* through Monte-Carlo simulations
- choose action with highest value

# Partially Observable Monte-Carlo Planning with Observation Widening (POMCPOW)
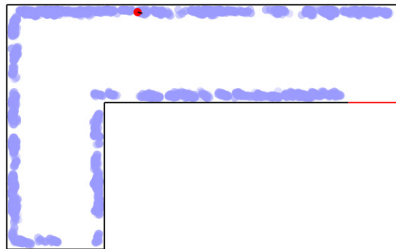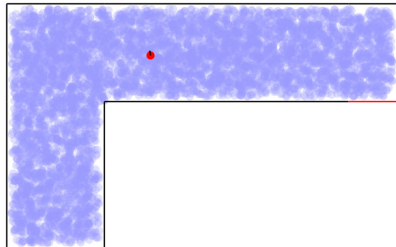


**High Level Idea**

- Monte-Carlo tree search
- locally approximate the *value function* through Monte-Carlo simulations
- choose action with highest value

**Domain Knowledge**
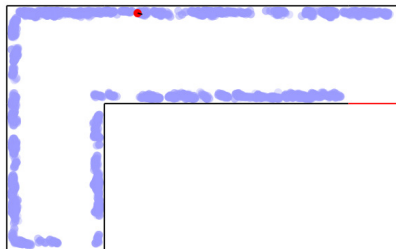
- user-defined value estimate
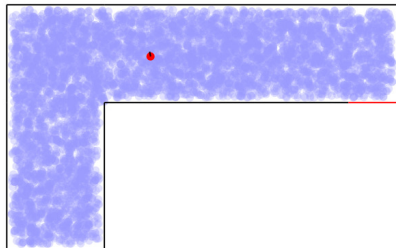
# Localization and Planning

**Screenshot of the Simulator**



**Problem Details**

- *known room* but *unknown location*
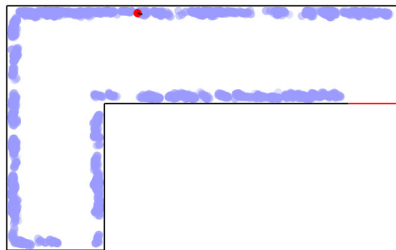
# Localization and Planning

**Screenshot of the Simulator**



## Problem Details

- *known room* but *unknown location*

- objective:
    - success: leave room at exit (green)
    - failure: falling down stairs (red)
    - penalties: time and collisions

- observations: collision sensor

- dynamics: noisy differential drive
    - actions: left, right, straight
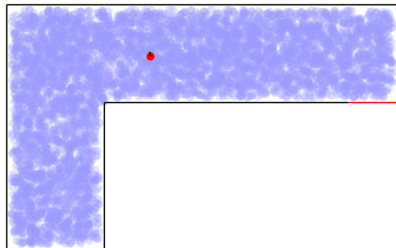
# Localization and Planning

**Screenshot of the Simulator**



**Problem Details**

- *known room* but *unknown location*

- objective:
    - success: leave room at exit (green)
    - failure: falling down stairs (red)
    - penalties: time and collisions

- observations: collision sensor

- dynamics: noisy differential drive
    - actions: left, right, straight

**POMDP**

- continuous state space
- discrete action and observation space

# Baselines

**Mode Control**

- use mode of belief to approximate the state
- treat POMDP as fully observable problem

# Baselines

**Mode Control**

- use mode of belief to approximate the state
- treat POMDP as fully observable problem

### Most Likely Reflex Agent (MLRA)

- handcrafted feedback policy
- P-controller tracking preset path

## Baselines

**Mode Control**

- use mode of belief to approximate the state
- treat POMDP as fully observable problem

### Most Likely Reflex Agent (MLRA)

- handcrafted feedback policy
- P-controller tracking preset path

### Most Likely Model Predictive Control (MLMPC)

- use MPC to plan a path
- cost function: negative reward

$\Rightarrow$ MLMPC is an optimal policy for the fully observed problem (MDP)

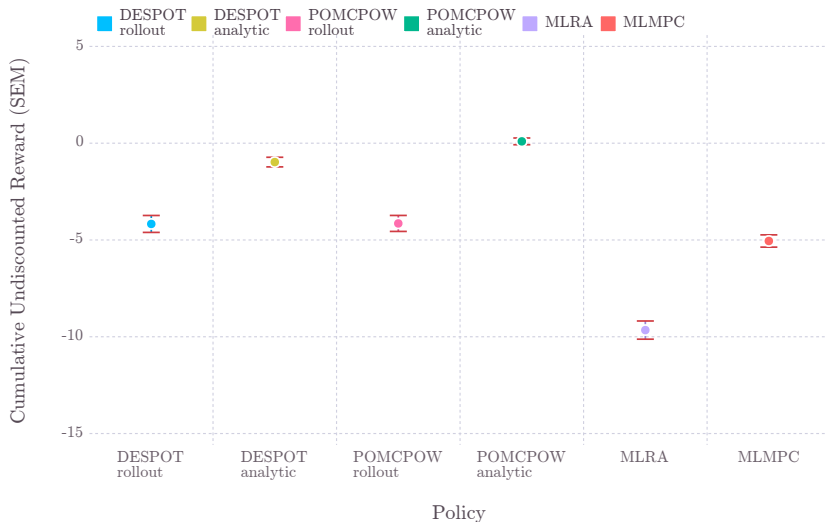# Integrating Domain Knowledge for POMDP Solvers

**Value Estimates**

1. Rollout Value Estimate
   - simulate default policy: "always straight"
2. Analytic Value Estimate
   - estimate remaining steps from distance to goal
   - approximate value by cumulative living penalty

# Integrating Domain Knowledge for POMDP Solvers

**Value Estimates**

1. Rollout Value Estimate
   - simulate default policy: "always straight"
2. Analytic Value Estimate
   - estimate remaining steps from distance to goal
   - approximate value by cumulative living penalty

**POMDP Solver Setups**

1. DESPOT-rollout
2. DESPOT-analytic
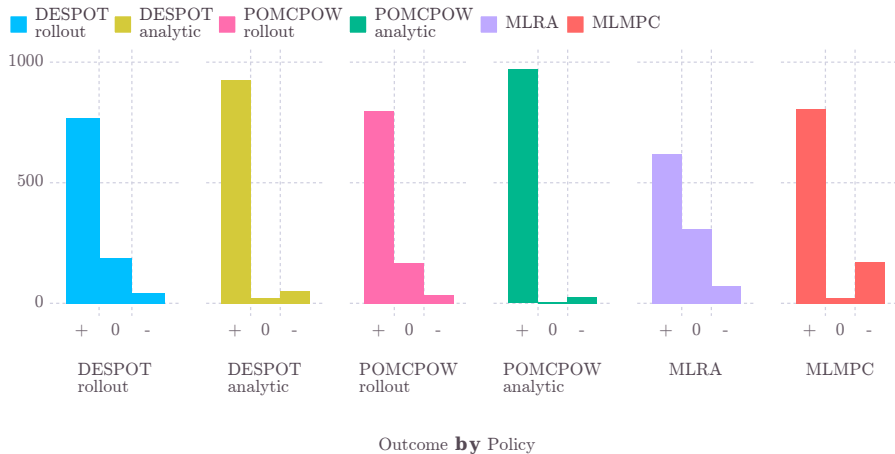3. POMCPOW-rollout
4. POMCPOW-analytic

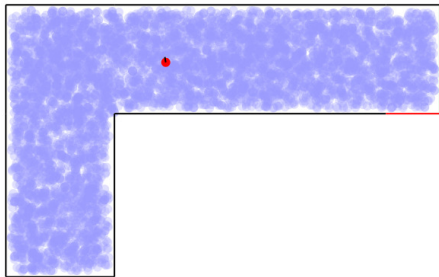**Mean and SEM of the undiscounted return
for 1000 simulations per policy.**

**Histogram of outcome frequencies grouped by policy.**



Outcome **by** Policy

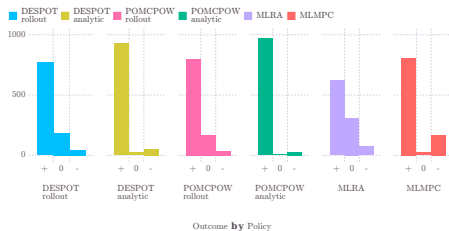## Initial Conditions



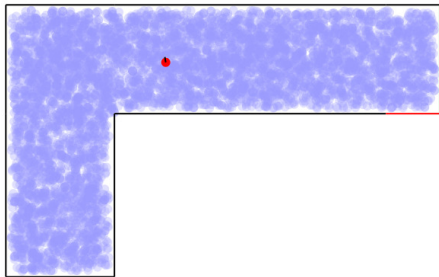## Histogram of outcome frequencies grouped by policy.
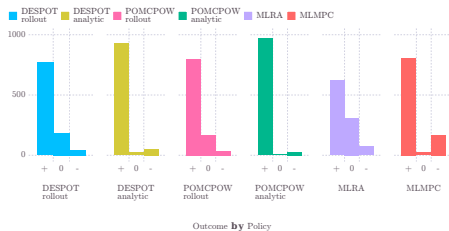
**Initial Conditions**



**Histogram of outcome frequencies grouped by policy.**



$\Rightarrow$ POMCPOW-analytic near optimal with respect to safety

**MLMPC**

**POMCPOW-Analytic**

# Evaluation – Qualitative Analysis

**MLMPC**

- mode approximation compromises safety
- passive information gathering
- fails for highly symmetric beliefs

**POMCPOW**-**Analytic**

## MLMPC

- mode approximation compromises safety
- passive information gathering
- fails for highly symmetric beliefs

## POMCPOW-Analytic

- active information gathering
- reliably reduces uncertainty
- safe and efficient behaviors

# Conclusion and Future Work

**Conclusions**

1. safe and efficient behaviors
   $\Rightarrow$ suitable approach for safety vs. efficiency tradoff

2. near real-time planning capabilities
   $\Rightarrow$ already a useful high-level planner for moderately sized problems

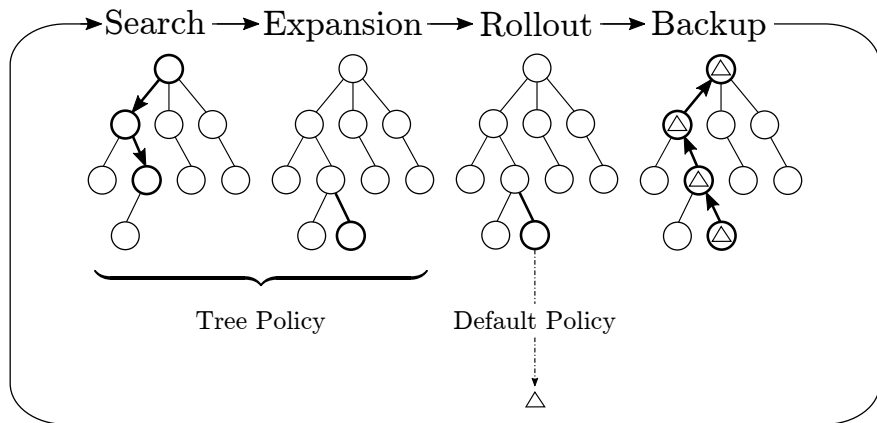3. designing suitable heuristic guidance hard but cruitial for performance

# Conclusion and Future Work

**Conclusions**

1. safe and efficient behaviors
   $\Rightarrow$ suitable approach for safety vs. efficiency tradoff

2. near real-time planning capabilities
   $\Rightarrow$ already a useful high-level planner for moderately sized problems

3. designing suitable heuristic guidance hard but cruitial for performance

**Future Work**

1. theory for a-priory safety assurances

2. improving scalability by using GPU

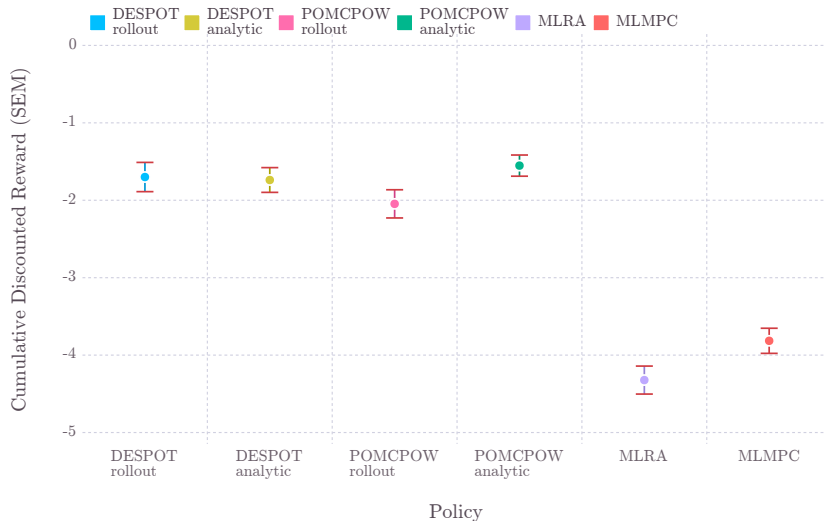3. learning for designing of heuristic guidance

# End
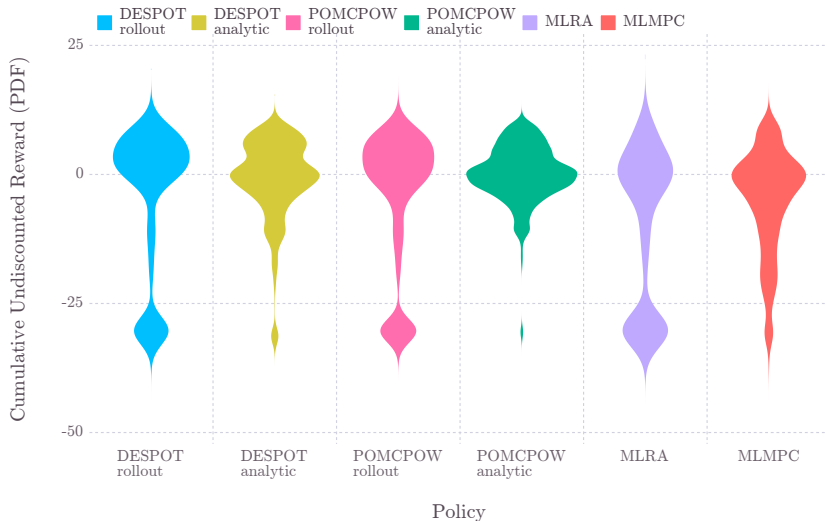
End

# Monte-Carlo Tree Search

# Evaluation: Localization and Planning – Objective Value



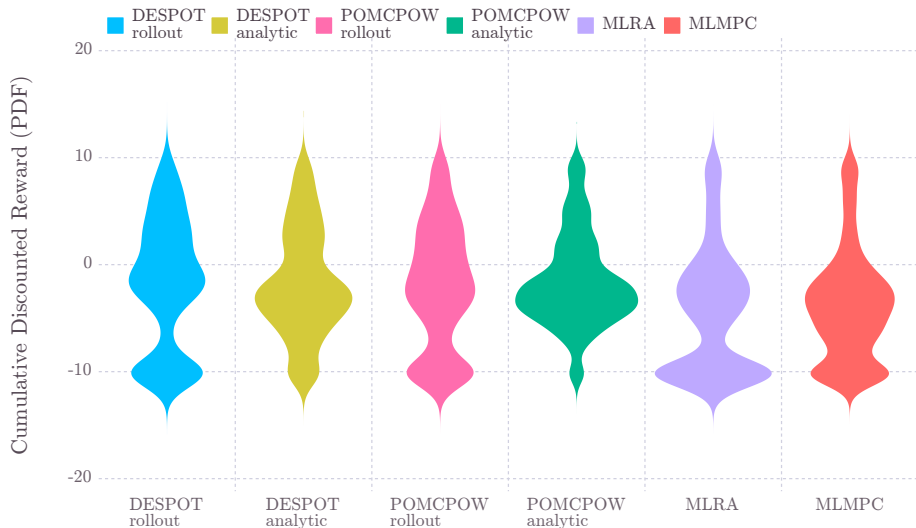**Mean and SEM of the discounted return for 1000 simulations per policy.**

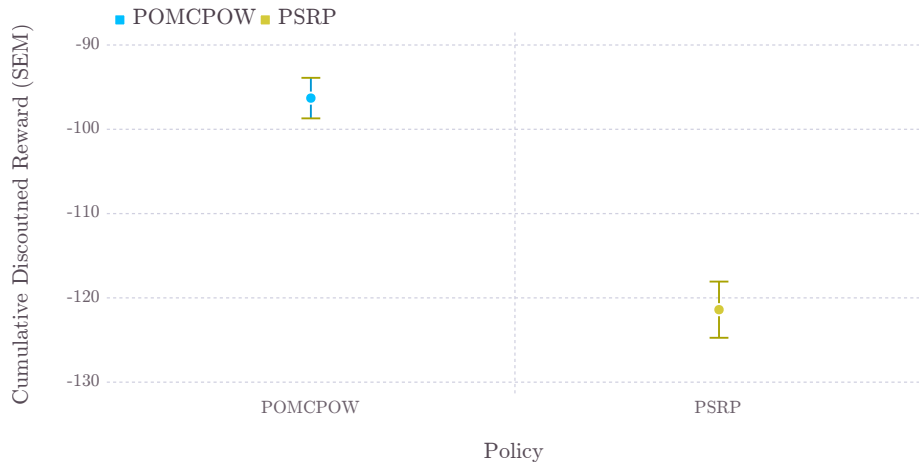**Distribution of the undiscounted return for 1000 simulations per policy.**

**Distribution of the return for 1000 simulations per policy:**

# Evaluation: Motion Planning with Latent Human Intentions

# Evaluation: Motion Planning with Latent Human Intentions



Outcome **by** Policy