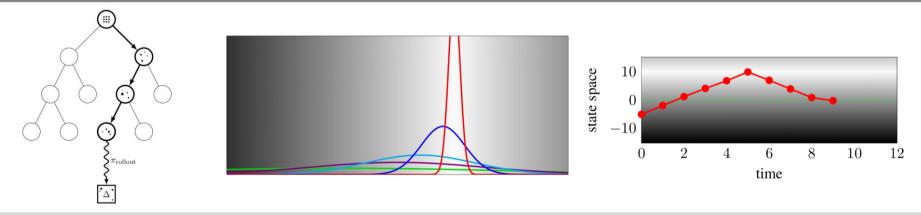# Information Particle Filter Tree: An Online Algorithm for POMDPs with Belief-Based Rewards on Continuous Domains

**Johannes Fischer * and Ömer Sahin Tas ***

*Equal contribution

International Conference on Machine Learning 2020

# POMDPs

- Model decision problems under uncertainty

Introduction
● ○ ○ ○

Reward Shaping
○ ○

IPFT
○ ○

Experiments
○ ○ ○ ○

Conclusion
○

Information Particle Filter Tree Algorithm for Continuous $\rho$POMDPs

ICML, July 2020

**2**

# POMDPs

- Model decision problems under uncertainty
- Cover uncertainties in
  - Models
  - Environment
  - Future behavior of others

# POMDPs

- Model decision problems under uncertainty
- Cover uncertainties in
  - Models
  - Environment
  - Future behavior of others



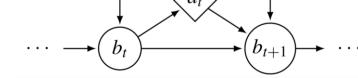Figure: Probabilistic graphical model of a POMDP.

Introduction
● ○ ○ ○

Reward Shaping
○ ○

IPFT
○ ○

Experiments
○ ○ ○ ○

Conclusion
○

Information Particle Filter Tree Algorithm for Continuous $\rho$POMDPs

ICML, July 2020

2

# POMDPs

- Model decision problems under uncertainty
- Cover uncertainties in
  - Models
  - Environment
  - Future behavior of others
- Reasoning in high dimensional belief space
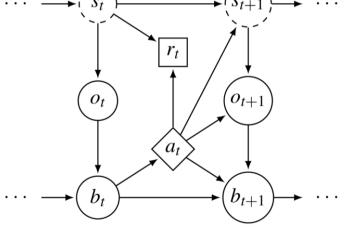  - → *Difficult to solve!*



Figure: Probabilistic graphical model of a POMDP.

Introduction
● ○ ○ ○

Reward Shaping
○ ○

IPFT
○ ○

Experiments
○ ○ ○ ○

Conclusion
○

Information Particle Filter Tree Algorithm for Continuous $\rho$POMDPs

ICML, July 2020

**2**

# POMDPs

- Model decision problems under uncertainty
- Cover uncertainties in
  - Models
  - Environment
  - Future behavior of others
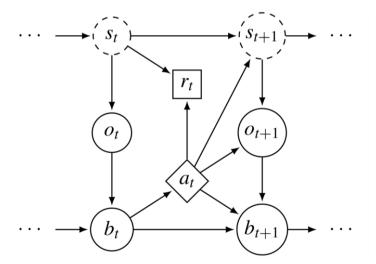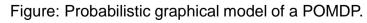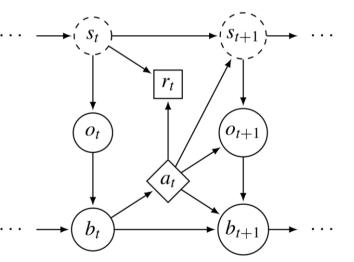- Reasoning in high dimensional belief space
  → *Difficult to solve!*

*Can POMDP solvers be improved by considering information?*



Figure: Probabilistic graphical model of a POMDP.

Introduction
● ○ ○ ○

Reward Shaping
○ ○

IPFT
○ ○

Experiments
○ ○ ○ ○

Conclusion
○

Information Particle Filter Tree Algorithm for Continuous $\rho$POMDPs

ICML, July 2020

**2**

# Information Measures

- Optimal value function $V^*$ and information measures have similar shape
  - →"more information = higher value"



Optimal value function

Negative entropy

Figure: Shape of optimal value function and negative entropy.

# Information Measures

■ Optimal value function $V^*$ and information measures have similar shape
  → "more information = higher value"

■ Motivation
  ■ Speed up planning
  ■ Allow active information gathering

Optimal value function



Negative entropy



Figure: Shape of optimal value function and negative entropy.

# POMDPs



Figure: Probabilistic graphical model of a POMDP.

# $\rho$**POMDPs**

- Extension of POMDP framework
- Belief-dependent reward model $\rho(b, a)$



Figure: Probabilistic graphical model of a $\rho$POMDP.

[1] Araya-López et al., "A POMDP Extension with Belief-dependent Rewards," (2010)

Introduction
○ ○ ● ○
Reward Shaping
○ ○
IPFT
○ ○
Experiments
○ ○ ○ ○
Conclusion
○
Information Particle Filter Tree Algorithm for Continuous $\rho$POMDPs                    ICML, July 2020        **4**

# $\rho$**POMDPs**

- Extension of POMDP framework
- Belief-dependent reward model $\rho(b, a)$
- Solvers exist only for
  - Discrete problems
  - Piecewise linear and convex $\rho$
  - Offline computation



Figure: Probabilistic graphical model of a $\rho$POMDP.

[1] Araya-López et al., "A POMDP Extension with Belief-dependent Rewards," (2010)

# $\rho$**POMDPs**

- Extension of POMDP framework
- Belief-dependent reward model $\rho(b, a)$
- Solvers exist only for
  - Discrete problems
  - Piecewise linear and convex $\rho$
  - Offline computation

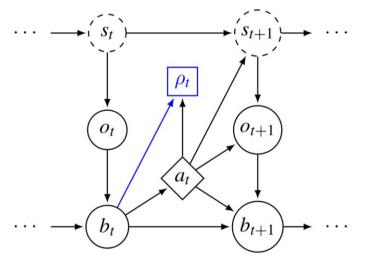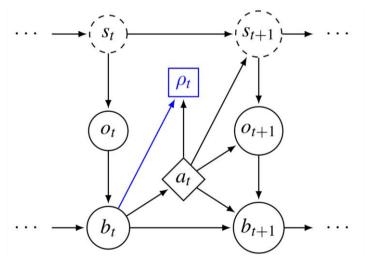> *How can $\rho$POMDPs on continuous domains be solved online?*



Figure: Probabilistic graphical model of a $\rho$POMDP.

[1] Araya-López et al., "A POMDP Extension with Belief-dependent Rewards," (2010)

# Approach - Information Particle Filter Tree

- Adapt MCTS-based POMDP solver
- Approximate belief by particles
- Evaluate $\rho$ on particle sets



Figure: Simulation phase of IPFT.

# Approach - Information Particle Filter Tree

- Adapt MCTS-based POMDP solver
- Approximate belief by particles
- Evaluate $\rho$ on particle sets

→Online anytime algorithm
→Continuous problems



Figure: Simulation phase of IPFT.

Introduction
○ ○ ○ ●

Reward Shaping
○ ○

IPFT
○ ○

Experiments
○ ○ ○ ○

Conclusion
○

Information Particle Filter Tree Algorithm for Continuous $\rho$POMDPs

ICML, July 2020

5

# Potential-Based Reward Shaping

■ Reward shaping changes the optimal policy

$$\tilde{R}(b_t, a_t) = R(b_t, a_t) + F(b_t, a_t, b_{t+1})$$

Introduction
○ ○ ○ ○

Reward Shaping
● ○

IPFT
○ ○

Experiments
○ ○ ○ ○

Conclusion
○

Information Particle Filter Tree Algorithm for Continuous $\rho$POMDPs                                      ICML, July 2020          **6**

# Potential-Based Reward Shaping

■ Reward shaping changes the optimal policy

$$\tilde{R}(b_t, a_t) = R(b_t, a_t) + F(b_t, a_t, b_{t+1})$$

■ BUT: Optimal policy is invariant under potential-based reward shaping for infinite horizon [2]

$$F(b_t, a_t, b_{t+1}) = \gamma\phi(b_{t+1}) - \phi(b_t)$$

[2] Eck et. al. "Potential-based reward shaping for finite horizon online POMDP planning." (2016)

# Potential-Based Reward Shaping

■ Reward shaping changes the optimal policy

$$\tilde{R}(b_t, a_t) = R(b_t, a_t) + F(b_t, a_t, b_{t+1})$$

■ BUT: Optimal policy is invariant under potential-based reward shaping for infinite horizon [2]

$$F(b_t, a_t, b_{t+1}) = \gamma\phi(b_{t+1}) - \phi(b_t)$$

■ $V^*$ serves as a particularly effective potential

[2] Eck et. al. "Potential-based reward shaping for finite horizon online POMDP planning." (2016)

| Introduction | Reward Shaping | IPFT | Experiments | Conclusion |
| --- | --- | --- | --- | --- |
| ○○○○ | ●○ | ○○ | ○○○○ | ○ |

Information Particle Filter Tree Algorithm for Continuous $\rho$POMDPs                    ICML, July 2020          6

# Information-Theoretic Reward Shaping

- Information measures have similar shape to $V^*$
  - Convex on belief space
  $\rightarrow$ Use as heuristic for $V^*$



Optimal value function

Negative entropy

Figure: Shape of optimal value function and negative entropy.

Introduction
○ ○ ○ ○

Reward Shaping
○ ●

IPFT
○ ○

Experiments
○ ○ ○ ○

Conclusion
○

Information Particle Filter Tree Algorithm for Continuous $\rho$POMDPs

ICML, July 2020

7

# Information-Theoretic Reward Shaping



Optimal value function

- Information measures have similar shape to $V^*$
  - Convex on belief space
  - →Use as heuristic for $V^*$

- Two potential-based shaping functions
  - Discounted information gain $\quad \Delta\mathcal{I}_\gamma(b, b') = \gamma\mathcal{I}(b') - \mathcal{I}(b)$
  - Undiscounted information gain $\quad \Delta\mathcal{I}_1(b, b') = \mathcal{I}(b') - \mathcal{I}(b)$

Negative entropy

Figure: Shape of optimal value function and negative entropy.

# Information-Theoretic Reward Shaping



Optimal value function

- Information measures have similar shape to $V^*$
  - Convex on belief space
  - → Use as heuristic for $V^*$

- Two potential-based shaping functions
  - Discounted information gain $\quad \Delta\mathcal{I}_\gamma(b, b') = \gamma\mathcal{I}(b') - \mathcal{I}(b)$
  - Undiscounted information gain $\quad \Delta\mathcal{I}_1(b, b') = \mathcal{I}(b') - \mathcal{I}(b)$

Negative entropy

Figure: Shape of optimal value function and negative entropy.

$$\rho(b, a, b') = \int_S R(s, a)b(s)ds + \lambda\Delta\mathcal{I}(b, b')$$

Introduction ○○○○

Reward Shaping ○●

IPFT ○○

Experiments ○○○○

Conclusion ○

Information Particle Filter Tree Algorithm for Continuous $\rho$POMDPs

ICML, July 2020

7

# Solving $\rho$POMDPs in Continuous Domains

- Based on Particle Filter Tree (PFT) Algorithm [3]
  - MCTS → continuous states
  - Double Progressive Widening (DPW)
    → continuous actions & observations



Figure: Simulation phase of PFT.

[3] Sunberg and Kochenderfer, "Online Algorithms for POMDPs with Continuous State, Action, and Observation Spaces," (2018)

Introduction
○ ○ ○ ○

Reward Shaping
○ ○

IPFT
● ○

Experiments
○ ○ ○ ○

Conclusion
○

Information Particle Filter Tree Algorithm for Continuous $\rho$POMDPs                    ICML, July 2020            **8**

# Solving $\rho$POMDPs in Continuous Domains



- Based on Particle Filter Tree (PFT) Algorithm [3]
  - MCTS → continuous states
  - Double Progressive Widening (DPW)
    → continuous actions & observations
  - Solves belief MDP
  - Small weighted particle sets $(m = 20)$
  - Update with mean particle return

[3] Sunberg and Kochenderfer, "Online Algorithms for POMDPs with Continuous State, Action, and Observation Spaces," (2018)

Figure: Simulation phase of PFT.

Introduction
○○○○

Reward Shaping
○○

IPFT
●○

Experiments
○○○○

Conclusion
○

Information Particle Filter Tree Algorithm for Continuous $\rho$POMDPs                    ICML, July 2020                    22

# Solving $\rho$POMDPs in Continuous Domains - Information Particle Filter Tree (IPFT)

- Particle set approximates belief



Figure: Simulation phase of IPFT.

Introduction
○ ○ ○ ○

Reward Shaping
○ ○

IPFT
○ ●

Experiments
○ ○ ○ ○

Conclusion
○

Information Particle Filter Tree Algorithm for Continuous $\rho$POMDPs                ICML, July 2020                **23**

# Solving $\rho$POMDPs in Continuous Domains - Information Particle Filter Tree (IPFT)

- Particle set approximates belief
- Evaluate $\rho$ on weighted particle sets, e.g.
  - $-\mathcal{H}(b) = \int_S b(s) \log b(s) \, ds \approx \sum_i w_i \log b(s_i)$



Figure: Simulation phase of IPFT.

Introduction
○ ○ ○ ○

Reward Shaping
○ ○

IPFT
○ ●

Experiments
○ ○ ○ ○

Conclusion
○

Information Particle Filter Tree Algorithm for Continuous $\rho$POMDPs                ICML, July 2020        **24**

# Solving $\rho$POMDPs in Continuous Domains - Information Particle Filter Tree (IPFT)

- Particle set approximates belief
- Evaluate $\rho$ on weighted particle sets, e.g.
  - $-\mathcal{H}(b) = \int_S b(s) \log b(s) \, \mathrm{d}s \approx \sum_i w_i \log \hat{b}(s_i)$
  - Particle-based kernel density estimate $\hat{b}$



Figure: Simulation phase of IPFT.

Introduction
○ ○ ○ ○

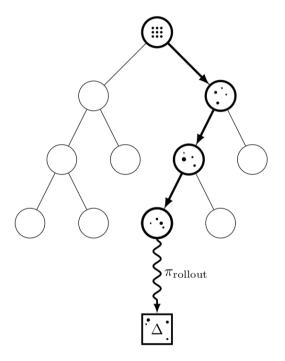Reward Shaping
○ ○

IPFT
○ ●

Experiments
○ ○ ○ ○

Conclusion
○

Information Particle Filter Tree Algorithm for Continuous $\rho$POMDPs

ICML, July 2020

25

# Solving $\rho$POMDPs in Continuous Domains - Information Particle Filter Tree (IPFT)

- Particle set approximates belief
- Evaluate $\rho$ on weighted particle sets, e.g.
  - $-\mathcal{H}(b) = \int_S b(s) \log b(s) \, ds \approx \sum_i w_i \log \hat{b}(s_i)$
  - Particle-based kernel density estimate $\hat{b}$
- Averaging over many particle sets leads to better entropy estimate



Figure: Simulation phase of IPFT.

Introduction
○ ○ ○ ○

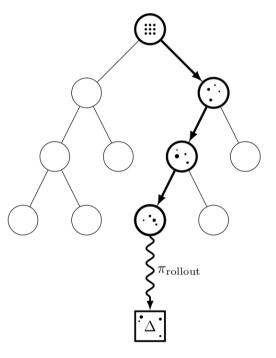Reward Shaping
○ ○

IPFT
○ ●

Experiments
○ ○ ○ ○

Conclusion
○

# Solving $\rho$POMDPs in Continuous Domains - Information Particle Filter Tree (IPFT)

- Particle set approximates belief
- Evaluate $\rho$ on weighted particle sets, e.g.
  - $-\mathcal{H}(b) = \int_S b(s) \log b(s) \, \mathrm{d}s \approx \sum_i w_i \log \hat{b}(s_i)$
  - Particle-based kernel density estimate $\hat{b}$
- Averaging over many particle sets leads to better entropy estimate

$\rightarrow$ IPFT can solve arbitrary $\rho$POMDPs on continuous domains



Figure: Simulation phase of IPFT.

# Experiments – Light Dark

Introduction
○○○○

Reward Shaping
○○

IPFT
○○

Experiments
●○○○○

Conclusion
○

Information Particle Filter Tree Algorithm for Continuous $\rho$POMDPs
ICML, July 2020
**10**

# Experiments – Light Dark


Figure: Light Dark environment.

- Goal: execute $a = 0$ at $s = 0$
- Consider action spaces

$$\mathbb{A}_{10} = \{-10, -1, 0, 1, 10\}$$

$$\mathbb{A}_3 \;\; = \{ \;\; -3, -1, 0, 1, \;\; 3\}$$

Introduction
○ ○ ○ ○

Reward Shaping
○ ○

IPFT
○ ○

Experiments
● ○ ○ ○

Conclusion
○

Information Particle Filter Tree Algorithm for Continuous $\rho$POMDPs                 ICML, July 2020          **10**

# Experiments – Light Dark



Figure: Light Dark environment.



Figure: Continuous Light Dark environment.

- Goal: execute $a = 0$ at $s = 0$
- Consider action spaces
$$\mathbb{A}_{10} = \{-10, -1, 0, 1, 10\}$$
$$\mathbb{A}_3 \ = \{ \ -3, -1, 0, 1, \ 3\}$$

- Continuous state space
- Transition noise
- Increased observation noise

# Results – Light Dark

| Algorithm | Light Dark problem | | | |
| --- | --- | --- | --- | --- |
| | action space $\mathbb{A}_{10}$ | | action space $\mathbb{A}_3$ | |
| IPFT($\Delta\mathcal{I}_1$) | $58.2 \pm 0.4$ | 🟩 | $34.8 \pm 0.7$ | 🟩 |
| IPFT($\Delta\mathcal{I}_\gamma$) | $55.4 \pm 0.5$ | 🟩 | $27.8 \pm 0.8$ | 🟩 |
| POMCPOW | $58.6 \pm 0.5$ | 🟩 | $-2.6 \pm 0.9$ | 🟥 |
| PFT-DPW | $57.4 \pm 0.5$ | 🟩 | $33.9 \pm 0.8$ | 🟩 |

Table: Mean reward and standard deviation of 1000 simulations.

Introduction
○ ○ ○ ○

Reward Shaping
○ ○

IPFT
○ ○

Experiments
○ ● ○ ○

Conclusion
○

Information Particle Filter Tree Algorithm for Continuous $\rho$POMDPs

ICML, July 2020

**11**

# Results – Light Dark

| Algorithm | Light Dark problem | | | | Continuous Light Dark problem | | | |
|---|---|---|---|---|---|---|---|---|
| | action space $\mathbb{A}_{10}$ | | action space $\mathbb{A}_3$ | | action space $\mathbb{A}_{10}$ | | action space $\mathbb{A}_3$ | |
| IPFT($\Delta\mathcal{I}_1$) | $58.2 \pm 0.4$ | | $34.8 \pm 0.7$ | | $35.7 \pm 1.8$ | | $35.9 \pm 1.0$ | |
| IPFT($\Delta\mathcal{I}_\gamma$) | $55.4 \pm 0.5$ | | $27.8 \pm 0.8$ | | $38.4 \pm 1.7$ | | $32.3 \pm 1.4$ | |
| POMCPOW | $58.6 \pm 0.5$ | | $-2.6 \pm 0.9$ | | $-8.5 \pm 2.3$ | | $-2.9 \pm 2.1$ | |
| PFT-DPW | $57.4 \pm 0.5$ | | $33.9 \pm 0.8$ | | $-33.1 \pm 2.4$ | | $-19.6 \pm 2.3$ | |

Table: Mean reward and standard deviation of 1000 simulations.

Introduction
○ ○ ○ ○

Reward Shaping
○ ○

IPFT
○ ○

Experiments
○ ● ○ ○

Conclusion
○

Information Particle Filter Tree Algorithm for Continuous $\rho$POMDPs          ICML, July 2020          **11**

# Results – Light Dark

| Algorithm | Light Dark problem | | Continuous Light Dark problem | |
| --- | --- | --- | --- | --- |
| | action space $\mathbb{A}_{10}$ | action space $\mathbb{A}_3$ | action space $\mathbb{A}_{10}$ | action space $\mathbb{A}_3$ |
| IPFT($\Delta\mathcal{I}_1$) | $58.2 \pm 0.4$ | $34.8 \pm 0.7$ | $35.7 \pm 1.8$ | $35.9 \pm 1.0$ |
| IPFT($\Delta\mathcal{I}_\gamma$) | $55.4 \pm 0.5$ | $27.8 \pm 0.8$ | $38.4 \pm 1.7$ | $32.3 \pm 1.4$ |
| POMCPOW | $58.6 \pm 0.5$ | $-2.6 \pm 0.9$ | $-8.5 \pm 2.3$ | $-2.9 \pm 2.1$ |
| PFT-DPW | $57.4 \pm 0.5$ | $33.9 \pm 0.8$ | $-33.1 \pm 2.4$ | $-19.6 \pm 2.3$ |

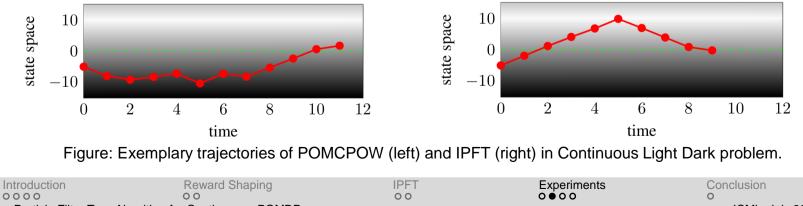Table: Mean reward and standard deviation of 1000 simulations.



Figure: Exemplary trajectories of POMCPOW (left) and IPFT (right) in Continuous Light Dark problem.

Introduction
Reward Shaping
IPFT
Experiments
Conclusion

Information Particle Filter Tree Algorithm for Continuous $\rho$POMDPs

ICML, July 2020

**11**

# Laser Tag



Figure: Laser Tag problem.

Introduction
○ ○ ○ ○

Reward Shaping
○ ○

IPFT
○ ○

Experiments
○ ○ ● ○

Conclusion
○

Information Particle Filter Tree Algorithm for Continuous $\rho$POMDPs

ICML, July 2020

**12**

# Laser Tag



Figure: Laser Tag problem.

| | Laser Tag | |
|---|---|---|
| IPFT($\Delta\mathcal{I}_1$) | $-9.0 \pm 0.2$ | |
| IPFT($\Delta\mathcal{I}_\gamma$) | $-8.9 \pm 0.2$ | |
| POMCPOW | $-9.9 \pm 0.2$ | |
| PFT-DPW | $-12.0 \pm 0.2$ | |

Table: Mean reward and standard deviation of 1000 simulations.

Introduction
○ ○ ○ ○

Reward Shaping
○ ○

IPFT
○ ○

Experiments
○ ○ ● ○

Conclusion
○

Information Particle Filter Tree Algorithm for Continuous $\rho$POMDPs                    ICML, July 2020                    **12**

# Hyperparameter Sensitivity Analysis



Figure: Mean reward and standard deviation of 1000 simulations of the Continuous Light Dark problem for different parameters.
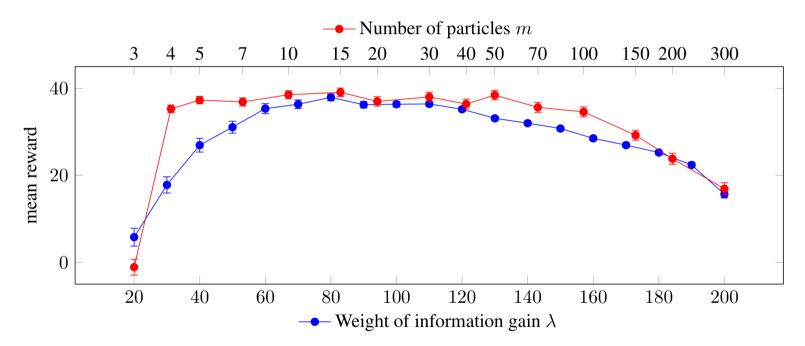
Introduction

Reward Shaping

IPFT

Experiments

Conclusion

Information Particle Filter Tree Algorithm for Continuous $\rho$POMDPs

ICML, July 2020

13

# Conclusion

*Can POMDP solvers be improved by considering information?*

Introduction
○ ○ ○ ○

Reward Shaping
○ ○

IPFT
○ ○

Experiments
○ ○ ○ ○

Conclusion
●

Information Particle Filter Tree Algorithm for Continuous $\rho$POMDPs

ICML, July 2020

14

# Conclusion

*Can POMDP solvers be improved by considering information?*

■ Information-theoretic reward shaping

→Helps by guiding agent to informative beliefs

Introduction
○ ○ ○ ○

Reward Shaping
○ ○

IPFT
○ ○

Experiments
○ ○ ○ ○

Conclusion
●

Information Particle Filter Tree Algorithm for Continuous $\rho$POMDPs

ICML, July 2020

**14**

# Conclusion



*Can POMDP solvers be improved by considering information?*

■ Information-theoretic reward shaping

→Helps by guiding agent to informative beliefs

*How can ρPOMDPs on continuous domains be solved online?*

Introduction
○○○○

Reward Shaping
○○

IPFT
○○

Experiments
○○○○

Conclusion
●

Information Particle Filter Tree Algorithm for Continuous ρPOMDPs                    ICML, July 2020          **14**

# Conclusion



> *Can POMDP solvers be improved by considering information?*

- Information-theoretic reward shaping
→ Helps by guiding agent to informative beliefs

> *How can $\rho$POMDPs on continuous domains be solved online?*

- IPFT combines PFT algorithm with $\rho$POMDPs
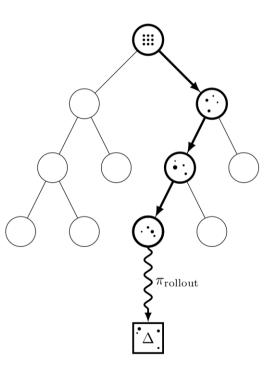→ General online solver for continuous $\rho$POMDPs

$\pi_{\text{rollout}}$

Figure: Simulation phase of IPFT.