Safe Multi-Agent Navigation guided by Goal-Conditioned Safe **Reinforcement Learning** ATLANTA 2025

Meng Feng*, Viraj Parimi*, Brian Williams Massachusetts Institute of Technology





Enabling scalable coordination of multi-agent systems relying on image-based observations to solve long-horizon safe navigation tasks, while balancing goal achievement with risk mitigation.

Motivation

- Autonomous Disaster Response
 - Time-critical missions
 - Safety-critical missions
- Time-critical missions demand multi-agent coordination for rapid task



Problem

How can we advance goal-conditioned RL to enable safer deployment in multi-agent, sparse-reward, long-horizon navigation problems?

Problem Features

Multiple agents

Key Insights

• Goal-conditioned RL handles local goal achievement with safe policy • Planning handles multi-agent coordination to ensure global goal

completion

• Safety-critical missions require navigating risks while meeting safety thresholds

Goal: Coordinate multiple agents to reach their goals in a safe manner

- Long-horizon visual navigation
- Avoid risky behaviors like getting too close to obstacles
- Sparse-reward settings

achievement

Experiments

2D Navigation



Problem Configurations		Methods (Cumulative Cost (Success Rate))			
Problem Type	Agents	Unconstrained Policy	Unconstrained Search (SoRB)	Constrained Policy	Constrained Search (Ours)
	1	$0.46 \pm 1.05~(100\%)$	$1.38 \pm 2.62 \ (100\%)$	$0.47 \pm 1.07~(100\%)$	0.49 ± 1.06 (100%)
Foot	5	N/A	6.39 ± 4.66 (100%)	N/A	$1.68 \pm 1.21~(100\%)$
Easy	10	N/A	$8.69 \pm 4.05 \ (100\%)$	N/A	$2.24 \pm 1.03~(100\%)$
	20	N/A	$10.51 \pm 3.57 \ (100\%)$	N/A	$\textbf{2.72} \pm \textbf{0.80} \ \textbf{(100\%)}$
	1	$1.58 \pm 2.03 \; (100\%)$	$2.30 \pm 3.02~(100\%)$	$1.58 \pm 2.01 \ (100 \ \%)$	1.59 ± 1.99 (100%)
Madium	5	N/A	$5.12 \pm 3.65 \ (100\%)$	N/A	$3.09 \pm 1.50~(100\%)$
Medium	10	N/A	7.27 ± 3.66 (98%)	N/A	$4.11 \pm 1.48 (98 \%)$
	20	N/A	8.73 ± 3.73 (98%)	N/A	$\bf 4.78 \pm 1.33~(98\%)$
	1	3.98 ± 4.40 (100%)	4.05 ± 4.19 (100%)	4.19 ± 4.23 (100%)	$3.09 \pm 3.87~(100\%)$
Hord	5	N/A	8.58 ± 3.79 (100%)	N/A	$6.01 \pm 1.96 \ (100 \%)$
Hard	10	N/A	$10.77 \pm 3.51 \; (100\%)$	N/A	$7.24 \pm 1.52~(100\%)$
	20	N/A	$11.96 \pm 3.31 \ (100\%)$	N/A	$8.36 \pm 1.58 (100 \%)$

Approach

Goal-Conditioned Risk-Aware Policy

$$egin{aligned} \pi^* &= rg\max_{\pi\in\Pi} \, J_r(\pi_ heta) \quad \pi^*, \lambda^* &= rg\min_{\lambda\geq 0} \,rg\max_{\pi\in\Pi} \, J_r(\pi_ heta) - \lambda(J_c(\pi_ heta) - \Delta) \ ext{ s.t } J_c(\pi_ heta) &\leq \Delta \end{aligned} \ egin{aligned} &T_r(\pi_ heta) &= \mathbb{E}_{s_1\sim
ho(s), \, a_t\sim\pi(a_t|s_t,s_g), \ s_{t+1}\sim p(s_{t+1}|s_t,a_t)} \left[\sum_{t=1}^T r(s_t,a_t,s_g)
ight] \, J_c(\pi_ heta) &= \mathbb{E}_{s_1\sim
ho(s), \, a_t\sim\pi(a_t|s_t,s_g), \ s_{t+1}\sim p(s_{t+1}|s_t,a_t)} \left[\sum_{t=1}^T c(s_t,a_t,s_g)
ight] \end{aligned}$$

Graph Generation

$d_{\pi}(s_i,s_j) = Q_d(s_i,\pi(s_i,s_j),s_j)$	$D \sim \operatorname{Cat}(N_D, p_{\pi}(d \mid s, a, s_g))$
$r_{\pi}(s_i,s_j) = Q_r(s_i,\pi(s_i,s_j),s_j)$	$R \sim \operatorname{Cat}(N_R, p_\pi(r \mid s, a, s_g))$

 $\mathcal{G} \triangleq (\mathcal{V}, \mathcal{E}, \mathcal{W}_d, \mathcal{W}_r)$



$\boldsymbol{\mathcal{I}}$		No. 10
where $\mathcal{V}=\mathcal{J}$	B	100 - 50- 100 - 50
$\mathcal{E}=\mathcal{B} imes\mathcal{B}=ig\{\epsilon$	$e_{s_i ightarrow s_j} s_i, s_j \in \mathcal{B} ig \}$	
$\mathcal{W}_{d}\left(e_{s_{i} ightarrow s_{j}} ight)=\left\{egin{array}{c} d_{\pi}(s_{i},\ \infty) \ \infty\end{array} ight.$	$egin{aligned} s_j) & ext{if} \ d_\pi(s_i,s_j) < d_{ ext{max}} \ & ext{otherwise} \end{aligned}$	
$\mathcal{W}_r\left(e_{s_i ightarrow s_j} ight) = \left\{egin{array}{c} r_\pi(s_i,\\infty) \ \infty\end{array} ight.$	$egin{aligned} s_j) & ext{if } r_\pi(s_i,s_j) < r_{ ext{max}} \ & ext{otherwise} \end{aligned}$	

Goal Sampling Algorithm 1 Self-Sampling and Training of Goal-Conditioned Actor Critic **Inputs**: Environment E, Goal-Conditioned Policy π , Associated Q-functions for rewards or auxiliary costs $Q(s, a, s_a)$, Desirable sample target for rewards or costs τ , Population size N, Number of training problems per batch K**Outputs**: Updated $\pi(s, a, s_g), Q(s, a, s_g)$ Paths 1: $\mathcal{P} \leftarrow$ Initialize an empty training set 2: for each batch do $\{(s_i, s_j)\}_N \leftarrow$ Randomly sample state pairs from E $\{v_{ij}\}_N \leftarrow Q(s_i, \pi(s_i, a, s_j), s_j)$ $\{l_{ij}\}_N \leftarrow L^2$ distance to τ Find $\{(s_i, s_j)\}$ corresponding to lowest $K \{l_i j\}_N$ Add the selected $\{(s_i, s_j)\}$ to the training set Train the agent with \mathcal{P} until \mathcal{P} is depleted **Plan and Act**



Visual Navigation

Unconstrained

Policy

+ Start 2

X End 2

Constrained

Search (Ours)

★ Goal 1

----- Waypoint 1

Single-Agent Trajectory Unconstrained Constrained



★ Goal 0

---- Waypoint 0

Unconstrained

Policy

Cost: 11.22

Steps: 9.00

Start 0

End 0



+

Start 1

X End 1



→ Waypoint 2

Unconstrained



---- Waypoint 3

ost: 7.23 eps: 8.00

Constrained

Search (Ours)

Street, and a second second the second second	COMPANY AND ADDRESS OF TAXABLE PARTY.	The second se	
Cost: 10.77 Steps: 8.00	Cost: 7.53 Steps: 8.00	Cost: Steps	
★ Goal 2	+ Start 3	★ Goal 3	

X End 3

Multi-Agent Trajectories

Constrained

Policy

Problem Configurations			Methods (Cumulative Cost (Success Rate))			
Мар	Problem Type	Agents	Unconstrained Policy	Unconstrained Search (SoRB)	Constrained Policy	Constrained Search (Ours)
1		1	3.16 ± 9.40 (100%)	2.14 ± 2.22 (98%)	$1.85 \pm 1.80 \; (100\%)$	$1.28 \pm 1.69 \; (100 \%)$
	Easy	5	N/A	4.32 ± 2.29 (98%)	N/A	$3.35 \pm 1.48 \; (100 \%)$
		10	N/A	5.36 ± 2.60 (96%)	N/A	$4.02\pm1.15(100\%)$
SC2		1	$4.98 \pm 3.46 \ (98\%)$	5.15 ± 3.32 (98%)	$4.77 \pm 3.25 \ (100\%)$	3.35 ± 3.23 (100%)
Staging 08	Medium	5	N/A	9.74 ± 11.13 (90%)	N/A	6.18 ± 1.64 (98%)
		10	N/A	14.23 ± 23.53 (84%)	N/A	$7.30 \pm 1.69 \ (98\%)$
		1	13.41 ± 4.00 (96%)	$12.35 \pm 5.10 \ (90\%)$	$10.43 \pm 3.07 \ (100\%)$	9.62 ± 3.66 (100%)
	Hard	5	N/A	21.23 ± 14.64 (84%)	N/A	13.02 ± 2.38 (98%)
		10	N/A	27.37 ± 19.05 (62%)	N/A	$13.87 \pm 2.17 \ (92\%)$
		1	1.87 ± 1.65 (98%)	2.23 ± 1.85 (98%)	$1.80 \pm 1.58 \; (100\%)$	$1.25 \pm 1.60 \; (100 \%)$
	Easy	5	N/A	3.71 ± 1.96 (98%)	N/A	$2.76 \pm 1.36 \; (100\%)$
	670 	10	N/A	4.20 ± 1.73 (98%)	N/A	$3.32 \pm 1.11 \; (100\%)$
SC3	2	1	5.06 ± 3.27 (100%)	4.87 ± 3.15 (100%)	4.38 ± 2.91 (100%)	$2.66 \pm 2.32 \; (100 \%)$
Staging 05	Medium	5	N/A	7.58 ± 2.16 (98%)	N/A	5.88 \pm 1.38 (100%)
		10	N/A	8.83 ± 1.75 (98%)	N/A	$7.28 \pm 1.53 \; (96\%)$
	2	1	14.86 ± 4.9 (98%)	13.16 ± 4.24 (100%)	8.59 ± 2.21 (98%)	6.47 ± 3.30 (98%)
	Hard	5	N/A	17.79 ± 2.68 (96%)	N/A	$15.10 \pm 2.19 \; (96\%)$
		10	N/A	19.56 ± 2.66 (94%)	N/A	$16.71 \pm 2.28 \ (94\%)$
0		1	$1.20 \pm 1.56 \; (100\%)$	$1.27 \pm 1.35 \ (100\%)$	$1.20 \pm 1.55 \ (100\%)$	$0.78 \pm 1.12 \; (100 \%)$
	Easy	5	N/A	2.32 ± 1.11 (100%)	N/A	$1.88 \pm 0.90 \; (100 \%)$
		10	N/A	$3.00 \pm 1.07 \; (100\%)$	N/A	$2.54 \pm 1.15 \; (100 \%)$
SC3		1	3.66 ± 2.95 (98%)	4.19 ± 3.23 (98%)	3.40 ± 2.83 (98%)	$1.95 \pm 2.24 \; (100 \%)$
Staging 11	Medium	5	N/A	8.33 ± 2.37 (98%)	N/A	$5.93 \pm 2.40 \; (96\%)$
		10	N/A	$10.23 \pm 2.27 \ (96\%)$	N/A	$8.13 \pm 2.18 \; (94\%)$
		1	16.41 ± 6.50 (100%)	14.75 ± 5.43 (100%)	11.59 ± 3.60 (100%)	$10.93 \pm 3.94 \; (100 \%)$
	Hard	5	N/A	18.67 ± 2.51 (96%)	N/A	$15.78 \pm 2.16 \; (98\%)$
		10	N/A	$20.37 \pm 2.41 \ (92\%)$	N/A	$16.94 \pm 1.66 \ (92\%)$

Inputs:	Agents N, Current States s, Goal States s_g , Buffer of observations \mathcal{B} ,
	Learned constrained policy π_c , Q-functions Q_d^{π} , Q_r^{π} of the unconstrained policy π
Output:	Action a
1: if <i>N</i>	> 1 then
2: s	$s_{w_1}, \leftarrow \text{CBS}(s, s_g, \mathcal{B}, Q_d^{\pi}, Q_r^{\pi})$
3: else	
4: <i>s</i>	$s_{w_1}, \leftarrow \text{SHORTEST_PATH}(s, s_g, \mathcal{B}, Q_d^{\pi}, Q_r^{\pi})$
5: if d_{π}	$(s \to s_{w_1}) < d_{\pi}(s \to s_g) \text{ or } d_{\pi}(s \to s_g) > d_{\max} \text{ then}$
6: <i>c</i>	$a \leftarrow \pi_c(a \mid s, s_{w_1})$
7: else	
8: <i>c</i>	$a \leftarrow \pi_c(a \mid s, s_g)$

References

- [1] B. Eysenbach, R. R. Salakhutdinov, and S. Levine, "Search on the replay buffer: Bridging planning and reinforcement learning," Advances in Neural Information Processing Systems, vol. 32, 2019.
- [2] M. G. Bellemare, W. Dabney, and R. Munos, "A distributional perspective on reinforcement learning," in International conference on machine learning. PMLR, 2017, pp. 449–458.
- [3] R. Stern, N. Sturtevant, A. Felner, S. Koenig, H. Ma, T. Walker, J. Li, D. Atzmon, L. Cohen, T. Kumar, et al., "Multi-agent pathfinding: Definitions, variants, and benchmarks," in Proceedings of the International Symposium on Combinatorial Search, vol. 10, 2019, pp. 151–158
- [4] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, "Continuous control with deep reinforcement learning," 2019.

Acknowledgements

This work was supported by the Defence Science and Technology Agency (DSTA). Any opinions, findings and conclusions or recommendations in this material are those of the author(s) and do not necessarily reflect the views of DSTA.