

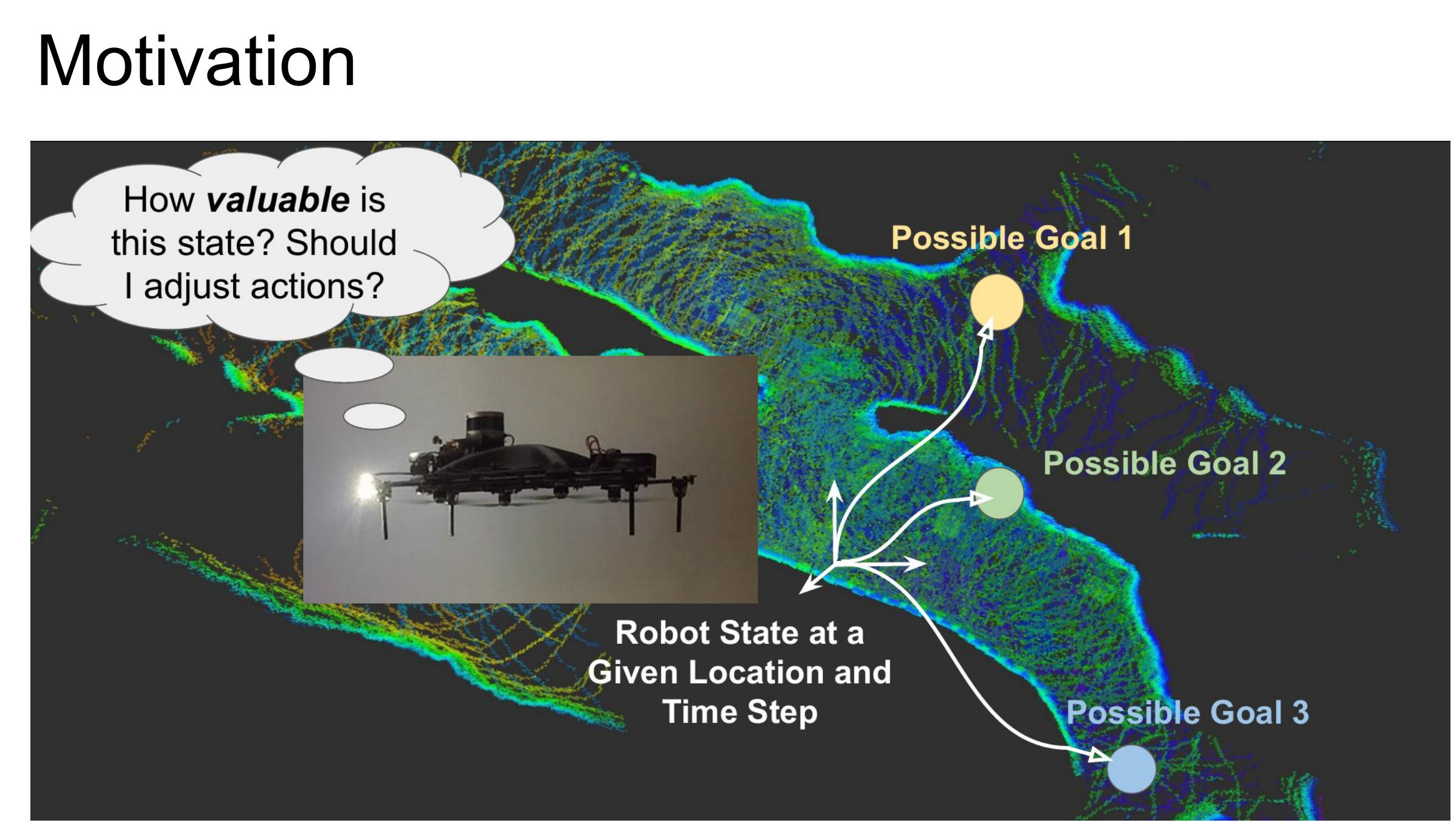
# Off-Policy Evaluation with Online Adaptation for Robot Exploration in Challenging Environments

Yafei Hu<sup>1</sup>, Junyi Geng<sup>1,2</sup>, Chen Wang<sup>1,3</sup>, John Keller<sup>1</sup>, and Sebastian Scherer<sup>1</sup>  
<sup>1</sup> CMU, <sup>2</sup> Penn State, <sup>3</sup> SUNY Buffalo



<https://jeffreyyh.github.io/opere/>

## Background and Motivation



- ### Problem formulation
- Robot exploration as a POMDP
  - Learn the value function of states
    - States visited by the robot in the trajectories
    - The goal states that haven't been visited
  - Value function learning as an off-policy evaluation (OPE) problem
    - On-policy evaluation for real robots is **costly** and **dangerous**
    - Data collection policy **different** from behaviour policy

## Experiments

### Robot Platform

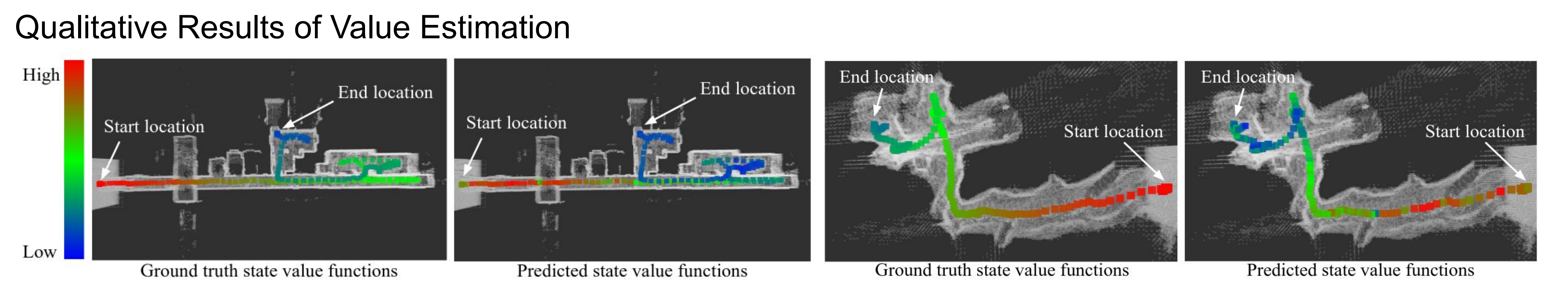
### Data Collection Envs

### Value Function Estimation Metrics

$$\text{NRMSE} = \frac{\text{RMSE}}{\hat{V}(s_t)_{\max} - \hat{V}(s_t)_{\min}}, \forall t \in [0, T-1]$$

$$R^2 = 1 - \frac{\sum_{t=0}^{T-1} [\hat{V}(s_t) - V_\pi(s_t)]^2}{\sum_{t=0}^{T-1} [\hat{V}(s_t) - \bar{V}(s_t)]^2}, \bar{V}(s_t) = \frac{1}{T} \sum_{t=0}^{T-1} \hat{V}(s_t)$$

$$\text{Regret} = \sum_{t=0}^{T-1} [R_C(\pi^*(a_t|o_t)) - R_C(\pi(a_t|o_t))]$$



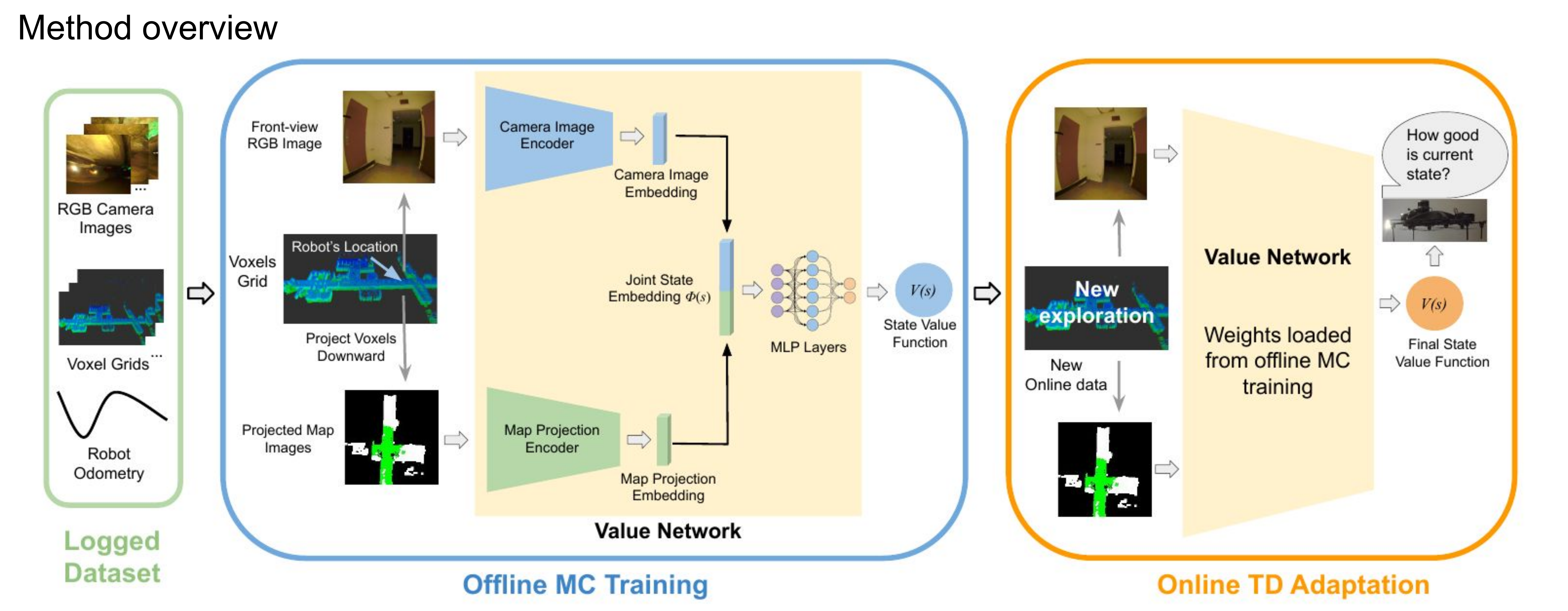
### Quantitative Results of Value Estimation

Metrics	Normalized RMSE (↓ the lower the better)				R2 Score (↑ the higher the better)			
	IS [12]	FQE [12]	DICE [32]	Ours	IS [12]	FQE [12]	DICE [32]	Ours
Corridor Env.	0.222±0.000	0.192±0.000	0.176±0.001	<b>0.129±0.004</b>	0.563±0.001	0.672±0.001	0.724±0.002	<b>0.853±0.010</b>
Room Env.	0.398±0.001	0.506±0.000	0.512±0.000	<b>0.213±0.002</b>	-0.593±0.007	-1.573±0.000	-1.634±0.000	<b>0.543±0.008</b>
Mine Env.	0.272±0.000	0.264±0.000	0.282±0.001	<b>0.207±0.002</b>	0.064±0.003	0.122±0.002	-0.002±0.005	<b>0.460±0.012</b>
Cave Env.	0.535±0.000	0.532±0.000	0.535±0.000	<b>0.164±0.002</b>	-2.234±0.000	-2.198±0.000	-2.235±0.000	<b>0.695±0.006</b>

### Regret analysis

Methods	Corridor Env.	Room Env.	Mine Env.	Cave Env.
Frontier [30]	0.333	0.714	0.625	0.250
IS [12]	0.633	0.443	0.363	0.625
FQE [12]	0.367	0.286	0.263	0.263
DICE [32]	0.617	0.343	0.125	0.375
Ours	<b>0.108±0.038</b>	<b>0.114±0.057</b>	<b>0.100±0.050</b>	<b>0.013±0.038</b>

## Method



- ### Reward Relabelling
- Extrinsic reward (sparse)
    - Use number of objects detected at a period of time, devoted as object gain (OG)  $OG(t) = O(t) - O(t - \Delta t)$
  - Intrinsic rewards:
    - Camera visual coverage map point gain (CG)  $CG(t) = C(t) - C(t - \Delta t)$
    - Lidar frontier map point gain (LG)  $LG(t) = L(t) - L(t - \Delta t)$
  - Final reward  $R(t) = aCG(t) + bLG(t) + cOG(t)$

### Training and Online Adaptation

Monte-Carlo (MC) Offline Training:

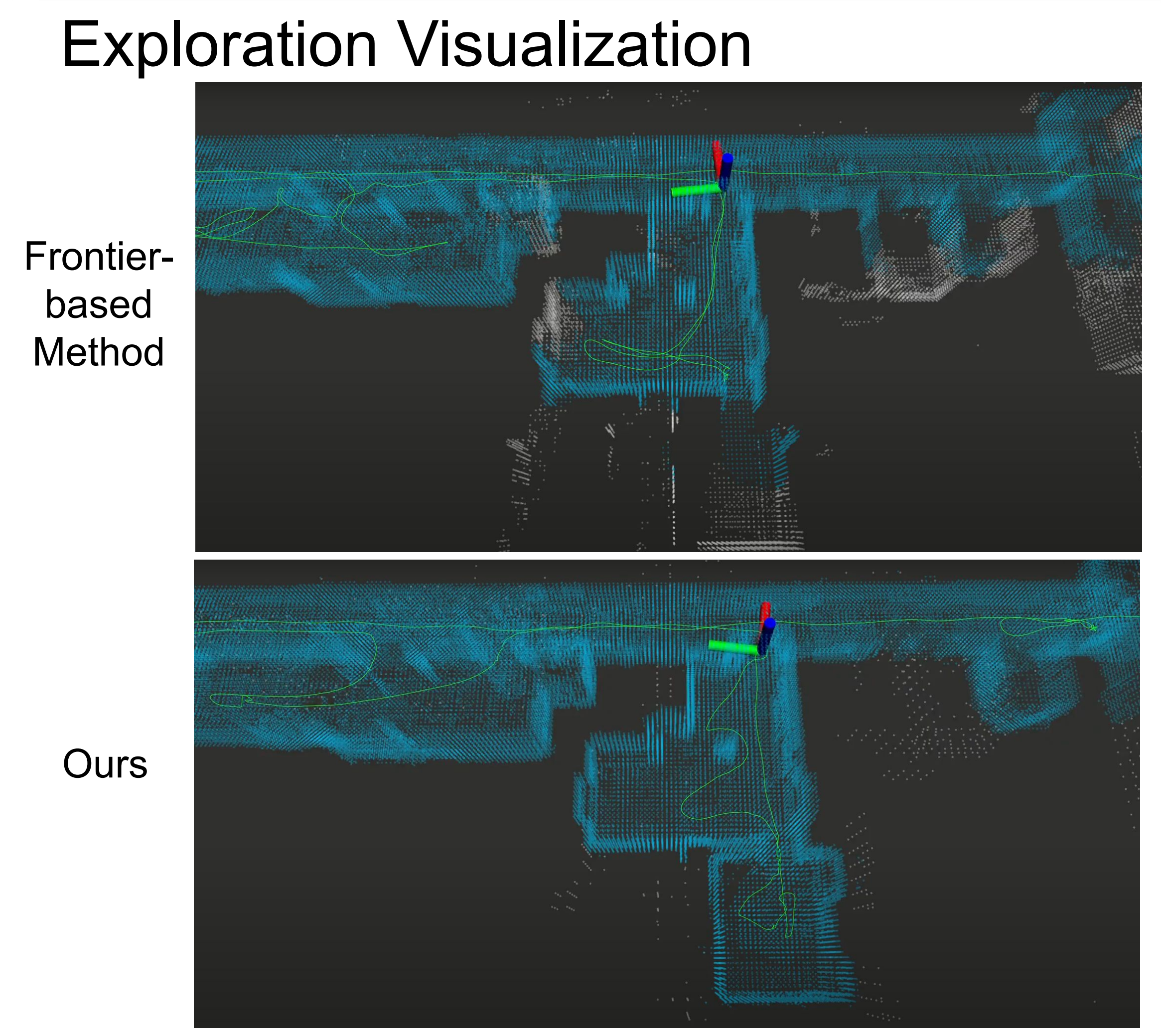
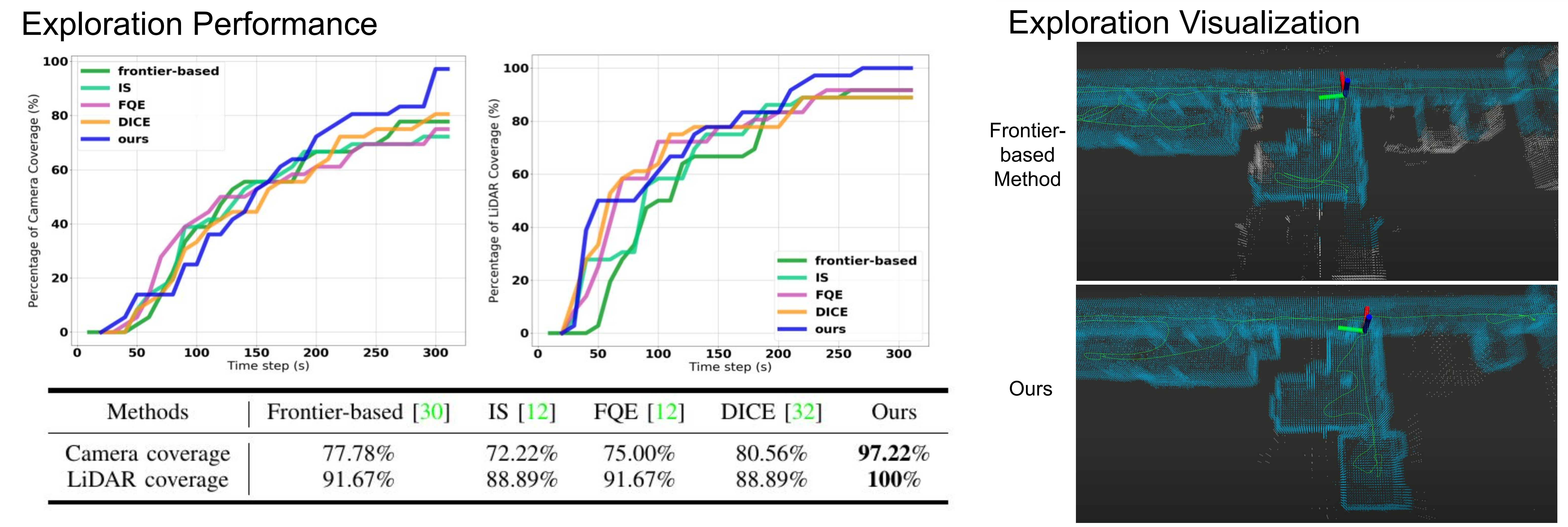
$$J(\theta) = \sum_{s \in \mathcal{S}} [V_\pi(s) - \hat{V}(s, \theta)]^2$$

$$\theta = \theta + \eta [G_t - \hat{V}(\phi(s), \theta)] \nabla_{\theta} \hat{V}(\phi(s), \theta)$$

Temporal-Difference (TD) Online Adaptation:

$$\theta = \theta + \eta [R(t) + \gamma \hat{V}(\phi(s_{t+1}), \theta) - \hat{V}(\phi(s_t), \theta)] \nabla_{\theta} \hat{V}(\phi(s_t), \theta)$$

$$= \theta + \eta [\beta \hat{V}(\phi(s_t), \theta) - \hat{V}(\phi(s_t), \theta)] \nabla_{\theta} \hat{V}(\phi(s_t), \theta)$$



## Conclusion

Major takeaways: • We propose an offline MC pre-training and TD online adaptation method to learn value function for robot exploration • The proposed method outperforms other OPE baselines • Real robot testings show certain advantages over frontier-based baseline