

Supriyo Ghosh¹, Sean Laguna¹, Shiau Hong Lim¹, Laura Wynter¹ and Hasan Poonawala²

¹IBM Research AI, Singapore; ²Amazon Web Services (AWS), UK

Motivation

Air traffic control (ATC)

- Monitor current state of aircrafts and recommend real-time decisions.
- Heavy traffic volume might lead to (human) operational errors.
- Need to optimize a complex objective (minimize congestion, conflicts, arrival delay and fuel cost).
- A sequential decision-making problem involving multiple actors influencing each other.



Our Contributions:

- Modelled ATC problem within a multi-agent reinforcement learning (MARL) framework.
- Solved the MARL problem with a model-based Kernel RL and a model-free deep RL methods.
- Proposed a general-purpose novel deep ensemble MARL method to combine the power of deep RL and kernel RL.
- Demonstrated the efficacy of ensemble MARL method on a real-world dataset consisting of ~1600 active aircrafts.

MARL Formulation for ATC

- Single Agent RL: $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$
- States Actions Transition function Reward function

Learn a policy π to maximize long term rewards:

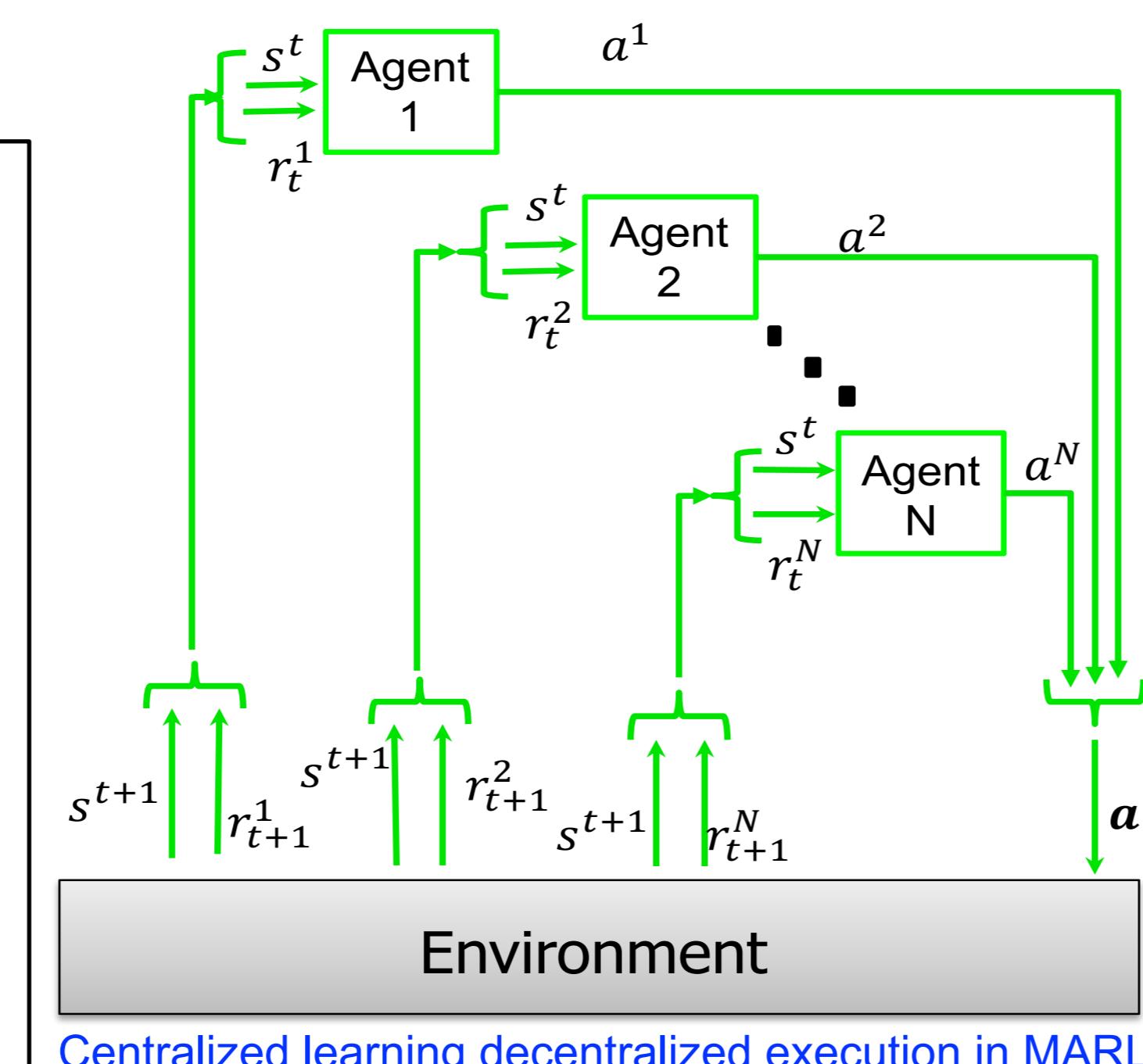
$$Q^*(s, a) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) | s_0 = s, a_0 = a \right]$$

- Multi-agent RL: $\langle \mathcal{S}, \mathcal{O}_1 \dots \mathcal{O}_N, \mathcal{A}_1 \dots \mathcal{A}_N, \mathcal{P}, \mathcal{R}, \gamma \rangle$

Observation: $o_i : \mathcal{S} \rightarrow \mathcal{O}_i$

Transition: $\mathcal{P} : \mathcal{S} \times \mathcal{A}_1 \times \dots \times \mathcal{A}_N \rightarrow \mathcal{S}$

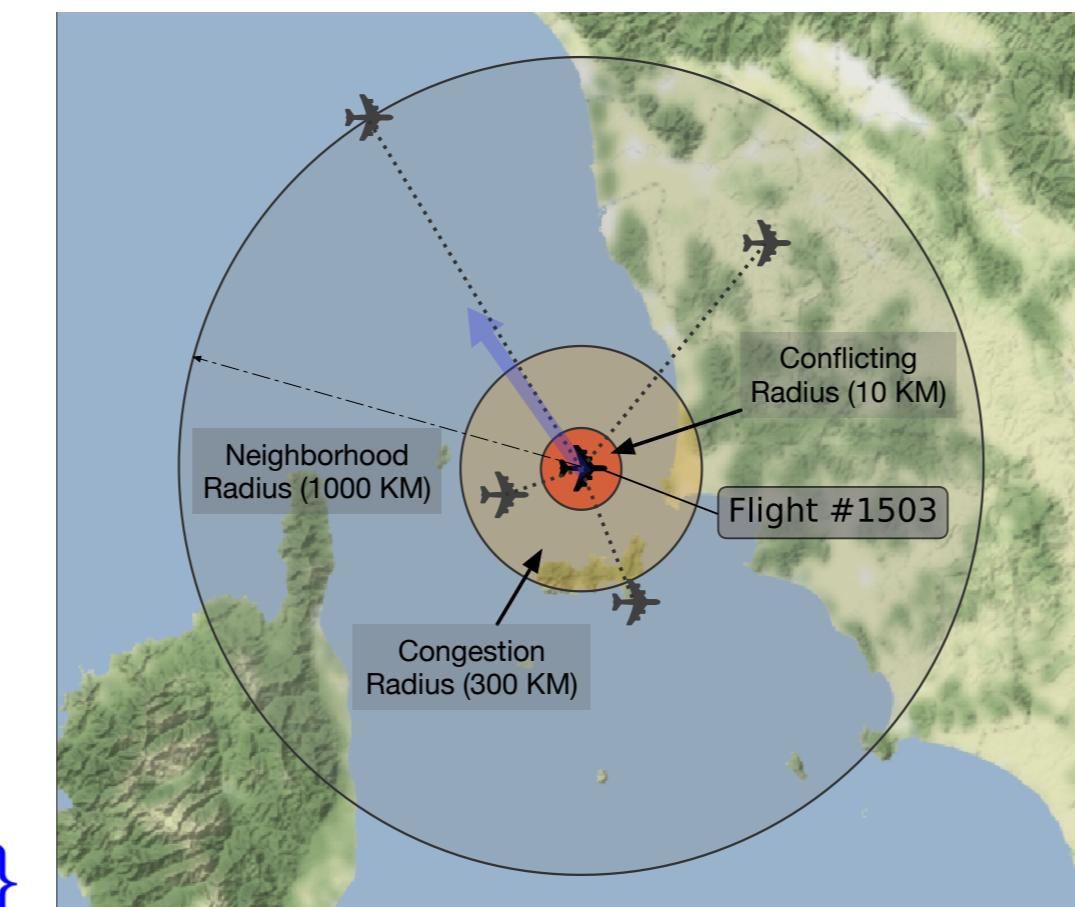
Reward: $r_i : \mathcal{O}_i \times \mathcal{A}_i \rightarrow \mathbb{R}$



MARL Components for ATC Problem:

State Space

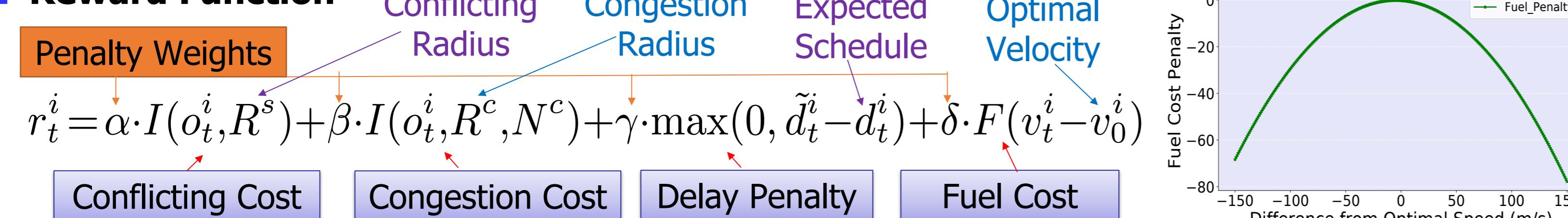
- Local features: Aircraft's location, speed, direction, timeliness
- Neighborhood features: N nearest aircrafts' relative velocity and relative direction
- Extended feature: Coarse and fine grid image information



Action Space (deviate speed by δ)

$$A_t = \{\max(v_{min}, (v_{t-1} - \delta)), v_{t-1}, \min(v_{max}, (v_{t-1} + \delta))\}$$

Reward Function



Kernel and Deep MARL for ATC

Model Based Kernel RL

1. Inputs: $S^a = \{s_k^a, r_k^a, \hat{s}_k^a | k = 1, \dots, n_a\} \forall a \in \mathcal{A}$
2. Generate m representative states with K-means clustering: $\bar{S} = \{\bar{s}_1, \dots, \bar{s}_m\}$
3. Define Gaussian kernels $\kappa_\tau(\bar{s}, s)$, $\bar{\kappa}_\tau(\bar{s}, \bar{s})$ using Euclidean distance from original to rep states:
$$\kappa_\tau(\bar{s}, s) = \phi\left(\frac{\|\bar{s} - s\|}{\tau}\right); \kappa_\tau(\bar{s}, \bar{s}_i) = \frac{\kappa_\tau(\bar{s}, s_i)}{\sum_{j=1}^m \kappa_\tau(\bar{s}, s_j)}$$
4. Compute $D^a : d_{ij}^a = \bar{\kappa}_\tau(\bar{s}_i^a, \bar{s}_j)$
5. Compute $K^a : k_{ij}^a = \kappa_\tau(\bar{s}_i^a, s_j^a)$
6. Compute transition probability: $P^a = K^a D^a$
7. Compute reward $r^a : r_i^a = \sum k_{ij}^a r_j^a$
8. Solve MDP $\{\bar{S}, \mathcal{A}, P^a, r^a, \gamma = .99\}$ & obtain Q^*

Advantages

- Performs well in the dense neighborhood of sample training data.
- Strong theoretical bounds on training data.

Limitations

- Extrapolates poorly to unknown situations.

Model Free Deep RL (PPO)

1. Initialize policy network with parameter θ_0
2. For each episode k , run line 3-4:
3. For every time step t and for every agent i , collect transition samples $(s_t^i, a_t^i, r_t^i, s_{t+1}^i)$ by simulating policy $\pi(\theta_k)$ and store them in a replay buffer D .
4. After every episode k , Update θ_k with minibatch of transitions from replay buffer D for M rounds to optimize the PPO objective:
$$\mathbb{E}_t [\min(r_t(\theta) A_t, \text{clip}(r_t(\theta), 1-\epsilon, 1+\epsilon) A_t)]$$

$r_t(\theta)$ is ratio between $\pi_\theta(a_t | s_t)$, $\pi_{\theta_{old}}(a_t | s_t)$

$A_t := R_t - V(s_t)$ is advantage function

Advantages

- Flexible and generalizes well to unknown situations.
- Can deal with richer state space information.

Limitations

- Can be brittle even in dense neighborhood of training data.

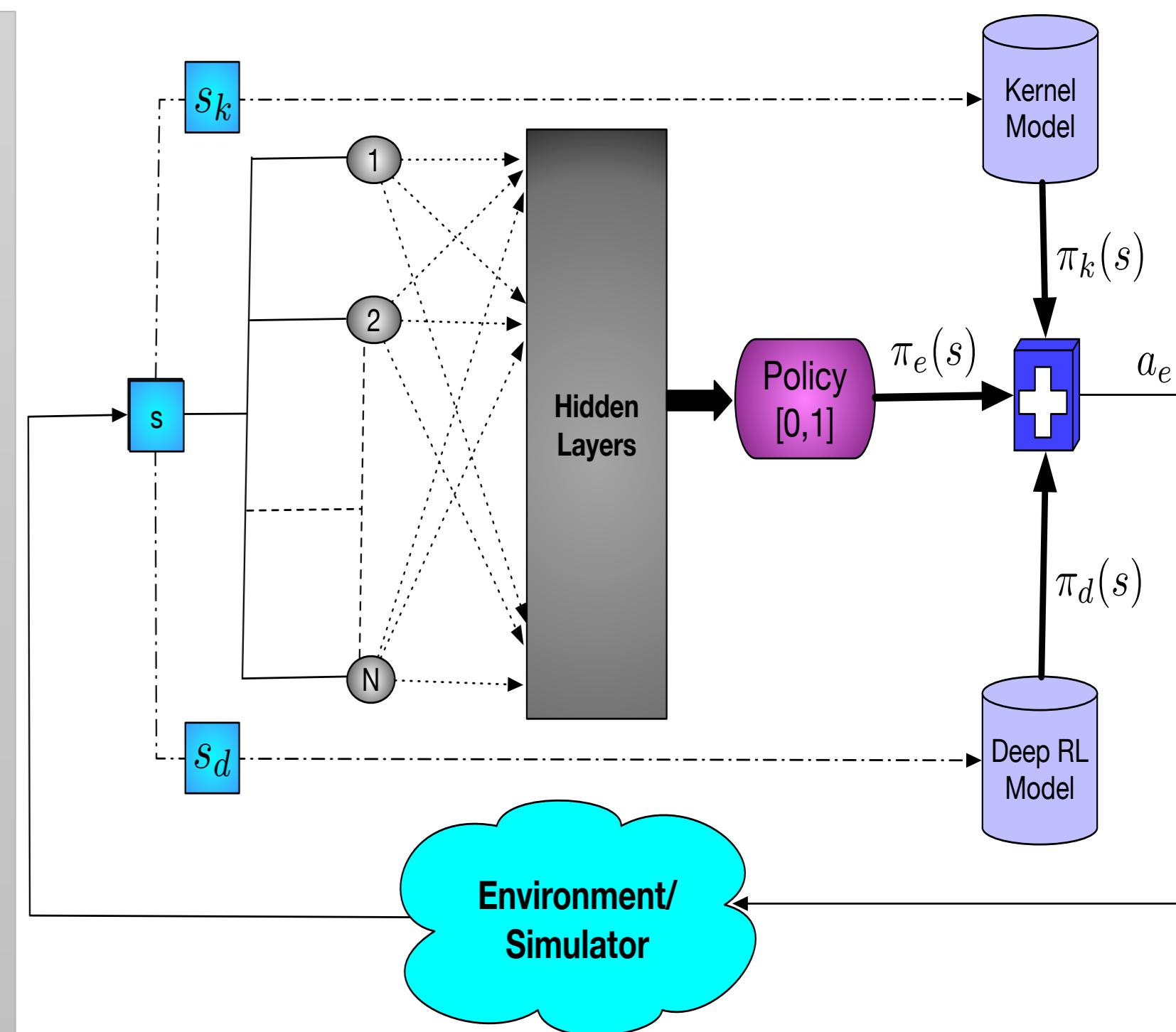
Deep Ensemble MARL

- Existing ensemble methods are either not feasible for model-based methods or unable to take multi-agent interactions in account.
- We train a separate deep neural network that efficiently learns to arbitrate between decisions of pre-trained kernel and deep MARL by considering multi-agent interactions in environment.

Ensemble Multi-agent RL:

Inputs: Kernel model \tilde{K} , PPO model $\tilde{\pi}(\tilde{\theta})$

1. Initialize ensemble policy to $\pi(\theta_0)$
2. For each episode k , run line 3-7
3. For each time t and agent i , sample ensemble action a_t^i using policy $\pi(\theta_k)$ for observation s_t^i
4. If a_t^i is 0 then get action \tilde{a}_t^i from \tilde{K} , otherwise get \tilde{a}_t^i using $\tilde{\pi}(\tilde{\theta})$
5. Execute joint action $\tilde{a}_t = (\tilde{a}_t^1, \dots, \tilde{a}_t^N)$
6. Store transitions $(s_t^i, a_t^i, r_t^i, s_{t+1}^i)$ in buffer D
7. For M rounds, update θ_k with minibatch of transitions from D to optimize the PPO objective



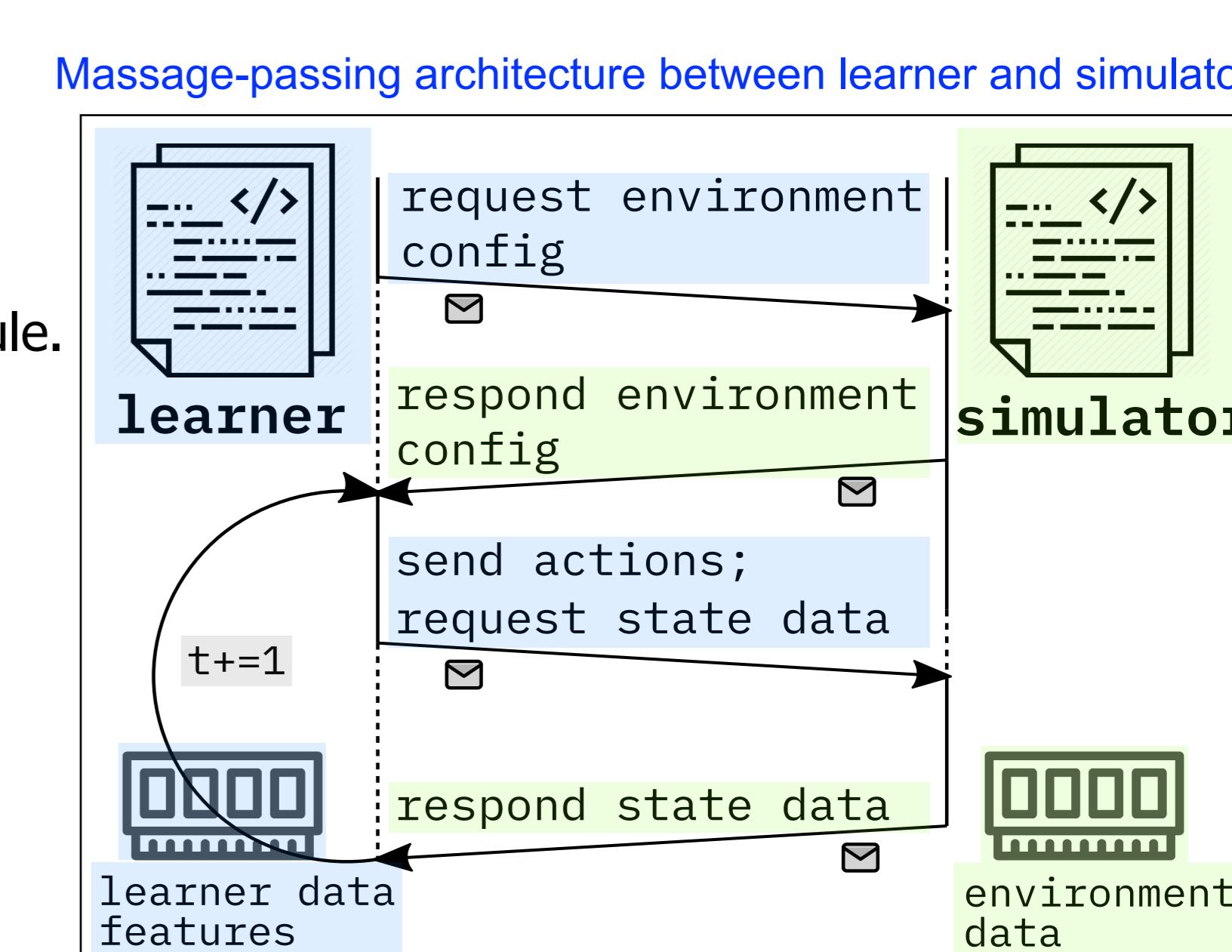
Experimental Results

Datasets (from Southern Europe)

- 24-hours schedule for ~1600 active flights.
- 300 training days and 30 testing days of schedule.
- 3 fuel cost settings considered: low, medium (from Airbus) and high.

Air traffic simulator

- An open-source simulator developed by Eurocontrol.
- We develop a message passing adapter between the simulator and our RL agent.



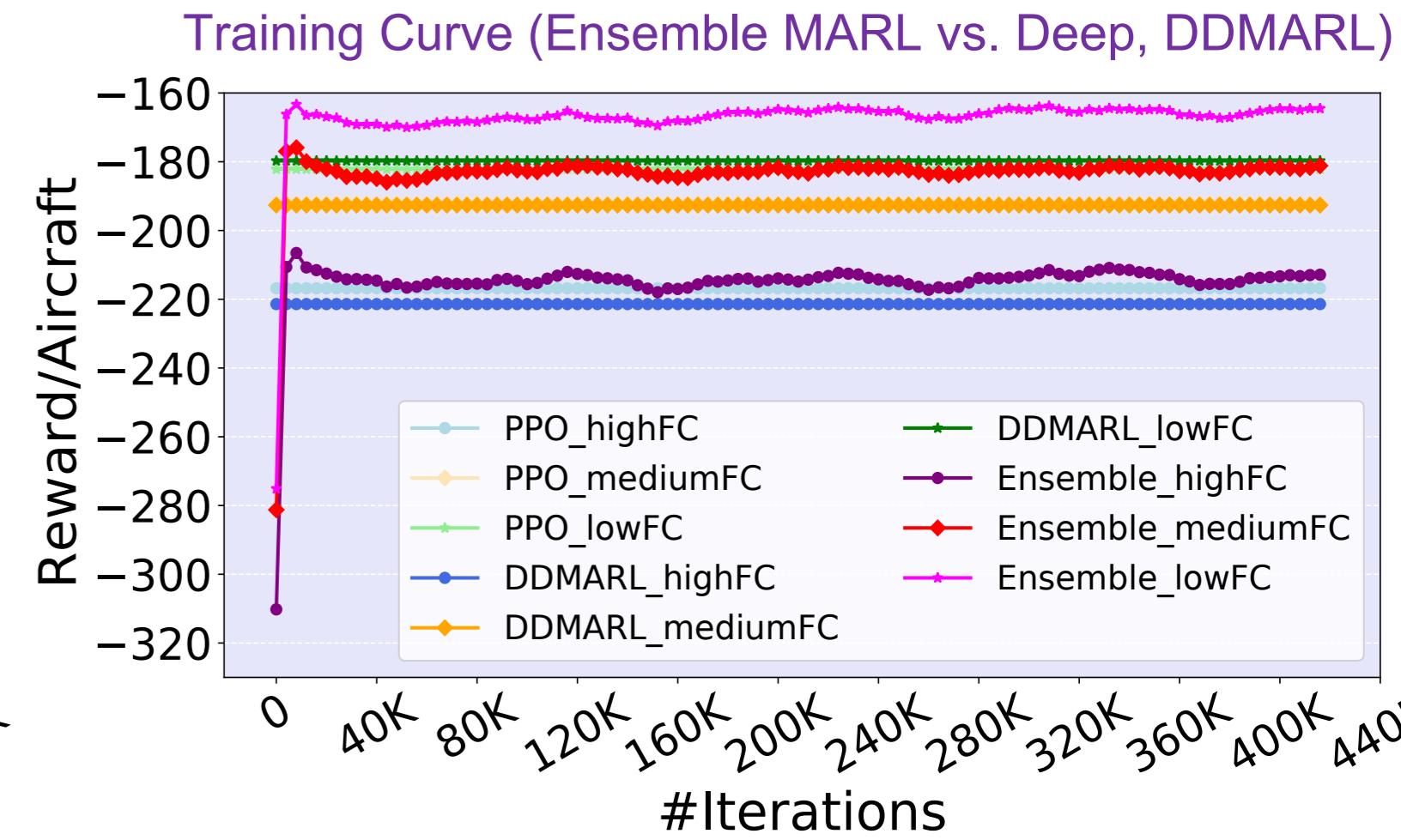
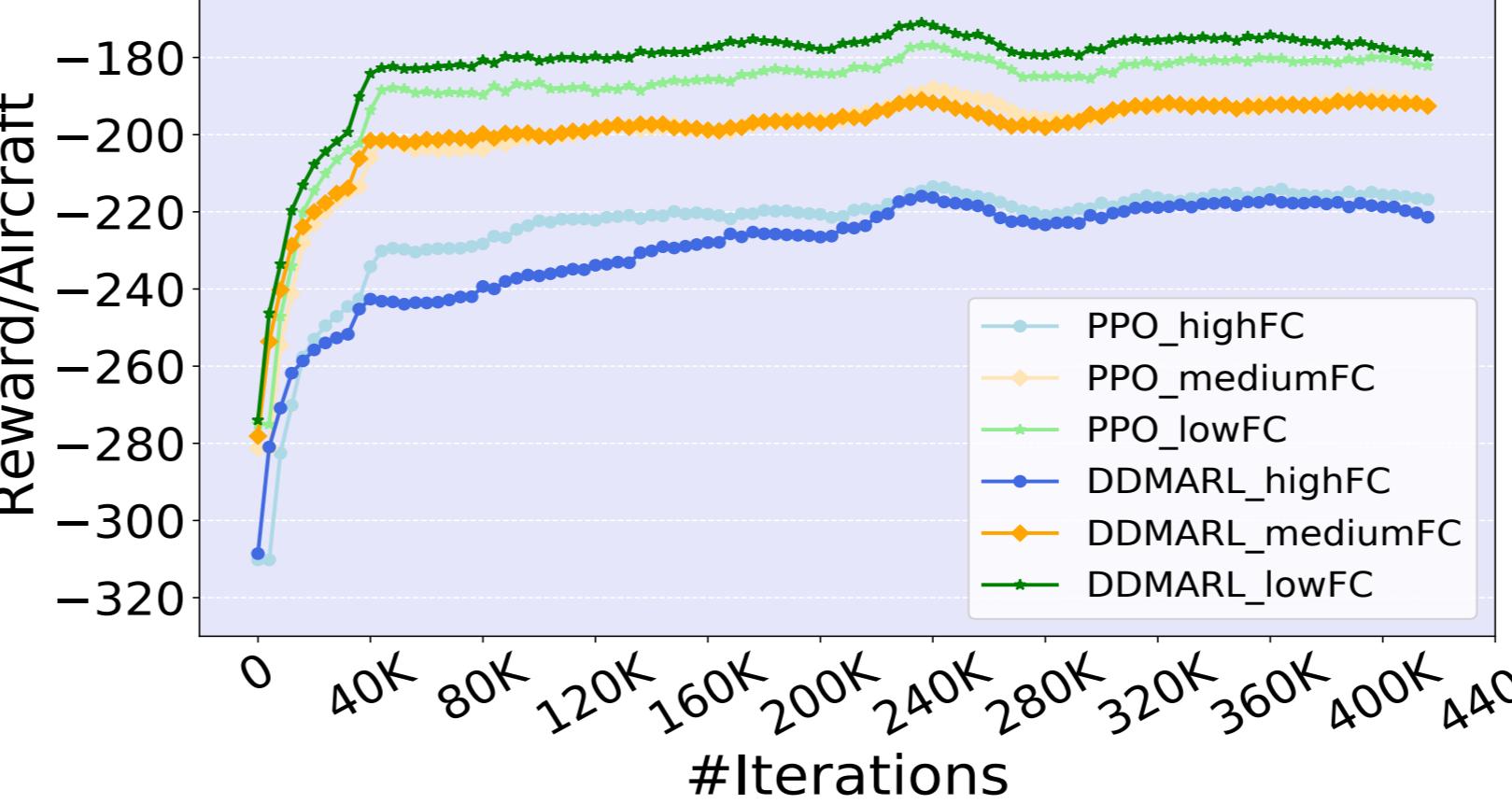
Benchmark Approaches

- Baseline: simulate default schedule (no penalty for fuel & delay).
- Local search: Each aircraft chooses a myopic best action.
- DDMARL (Brittain et. al., 2019): Only consider penalty for conflicting situations.

Training performance

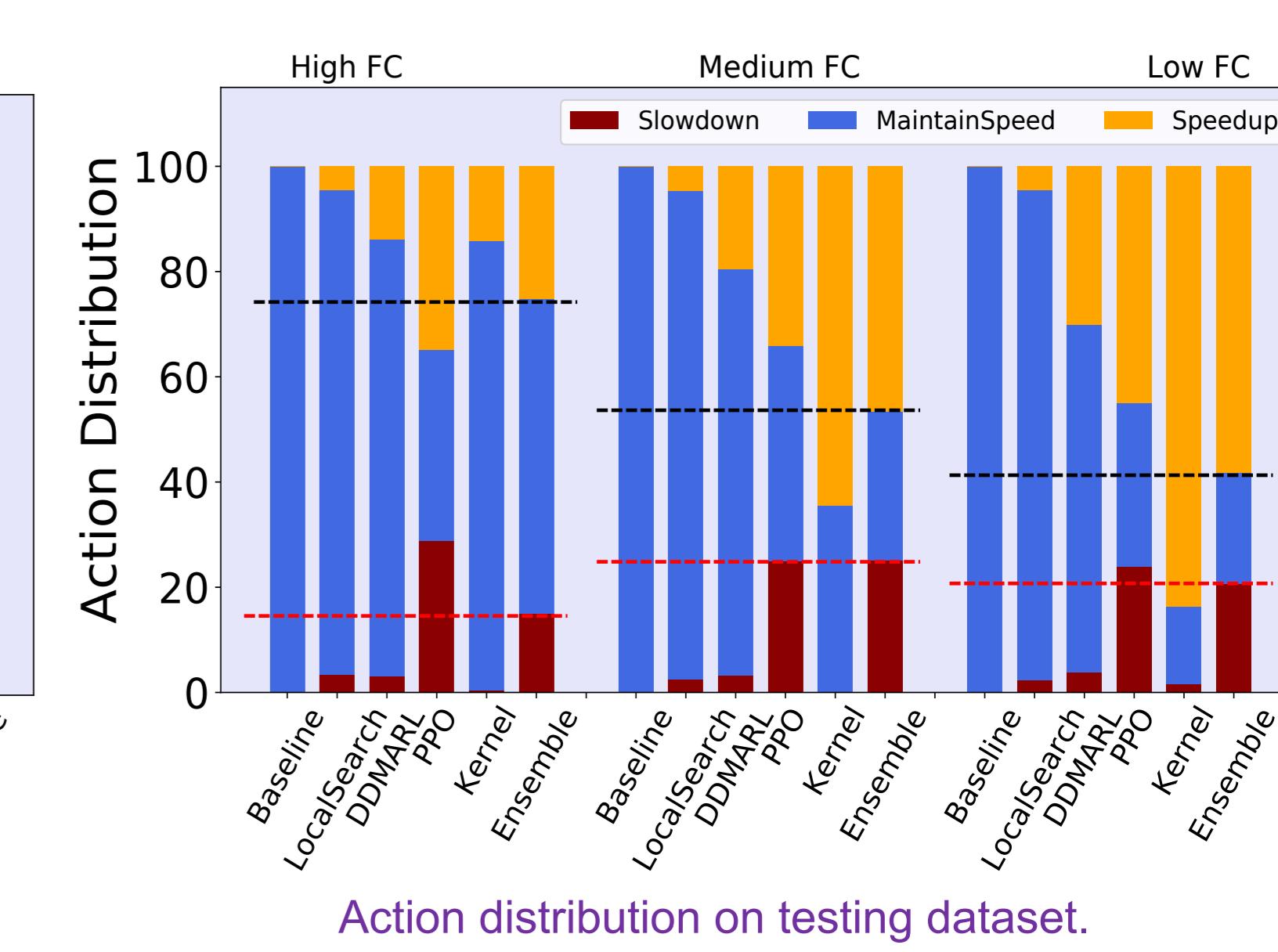
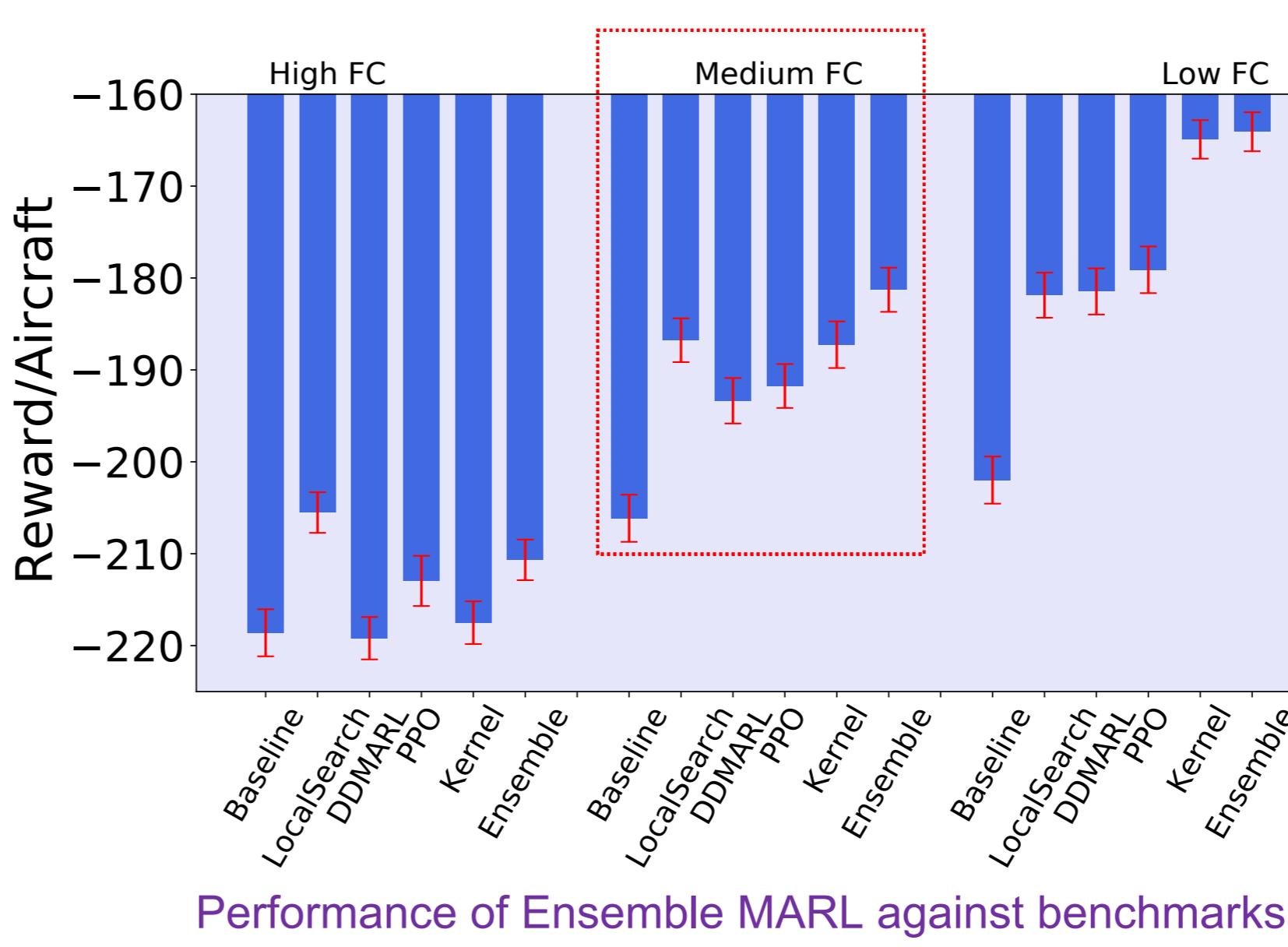
- Training performance of our deep MARL is at par with DDMARL.
- Our ensemble MARL has better training performance than both deep MARL and DDMARL across the board.

Fuel cost penalty structure



Testing performance

- Ensemble MARL always outperforms kernel and deep RL.
- Ensemble MARL provides ~9% gain in reward in a realistic (medium) fuel cost setting.
- Ensemble MARL diversifies distribution of actions to maximize overall reward value.



Conclusion

- **Summary:** Our proposed novel deep ensemble method improves the objective of air traffic controllers by 9% on a real-world dataset consisting of ~1600 active aircrafts.
- **Future Work:** (1) Extend action space to incorporate additional controls such as directional and altitude changes; (2) Extend state space to handle take-off and landing scenarios; (3) Extend ensemble MARL to combine power of multiple methods.