

## Introduction

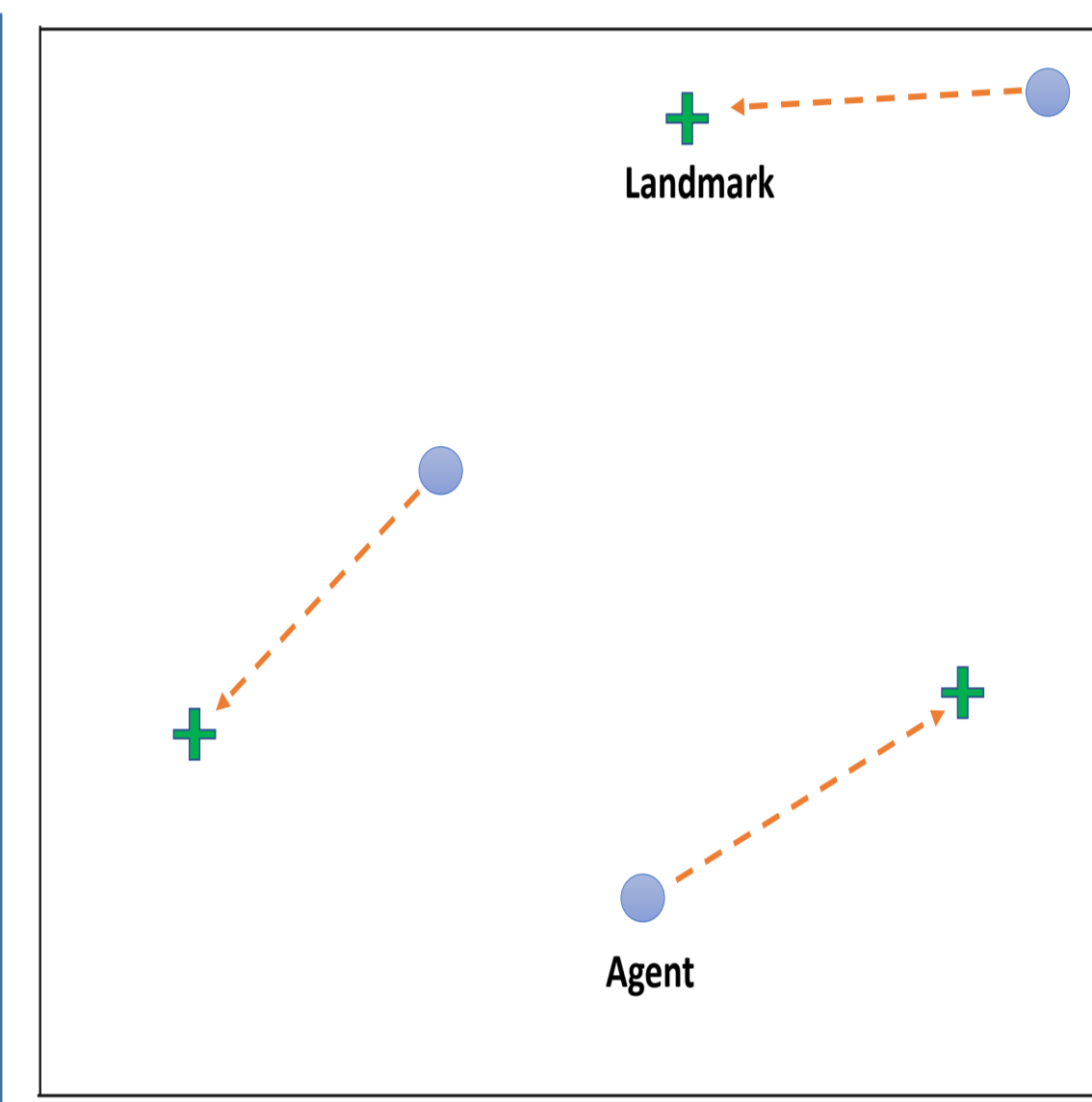
We address the problem of multiagent credit assignment in large scale multiagent system. Our main contributions are:

- An approach to learn a differentiable reward model by exploiting the collective nature of interactions among agents.
- A principled method to analytically compute shaped rewards from the reward model.
- A model-based RL approach that uses learned shaped rewards addressing credit assignment problem.

### ❖ Motivating Domain :



Air Traffic Control



Cooperative Navigation

## System Reward Approximator

### ❖ System Reward :

$$r(\mathbf{n}_t^{SA}) = \sum_{s \in S} \sum_{a \in A} n_t(s, a) \cdot \tilde{r}(s, a, \mathbf{n}_t^S)$$

### ❖ Loss Function for Reward Approximator :

$$\tilde{\mathcal{L}}(\mathbf{w}) = M \sum_{\xi \in \mathcal{B}} \sum_{s \in S} \sum_{a \in A} n_{\xi}(s, a) \cdot \left( \tilde{r}(s, a, \mathbf{n}_{\xi}^S) - r_{\mathbf{w}}(s, a, \mathbf{n}_{\xi}^S) \right)^2$$

## Approximate Difference Reward

### ❖ Difference rewards (DRs) :

$$D^m(s_t^m, a_t^m) = r(s_t, \mathbf{a}_t) - r(s_t^{-m} \cup d_s, \mathbf{a}_t^{-m} \cup d_a)$$

### ❖ Difference rewards with count variable:

$$D^m(s_t^m, a_t^m) = r_{\mathbf{w}}(\mathbf{n}_t^{SA}) - r_{\mathbf{w}}(\mathbf{n}_t^{SA - (s_t^m, a_t^m) + (d_s, d_a)})$$

### ❖ Difference rewards for state-action :

$$D_t(s, a) = r(\mathbf{n}_t^{SA}) - r(\mathbf{n}_t^{SA} - \mathcal{I}^{sa} + \mathcal{I}^{d_s d_a})$$

### ❖ Approximate difference rewards :

$$D_t(s, a) \approx \frac{1}{M} \cdot \left( \frac{\partial r_{\mathbf{w}}(\tilde{\mathbf{n}}_t^{SA})}{\partial \tilde{\mathbf{n}}_t^{SA}(s, a)} - \frac{\partial r_{\mathbf{w}}(\tilde{\mathbf{n}}_t^{SA})}{\partial \tilde{\mathbf{n}}_t^{SA}(d_s, d_a)} \right)$$

## Count Variables

### ❖ State count variable :

$$n_t(s) = \sum_{m=1}^M \mathbb{I}[s_t^m = s; \mathbf{s}_t], \forall s \in S$$

### ❖ State-action count variable :

$$n_t(s, a) = \sum_{m=1}^M \mathbb{I}[s_t^m = s, a_t^m = a; \mathbf{s}_t, \mathbf{a}_t], \forall s \in S$$

## Policy Gradient with DRs

### ❖ Return with difference rewards :

$$R_t^{dr} = \sum_{i=0}^{\infty} \gamma^i \left( \sum_{s \in S} \sum_{a \in A} n_{t+i}(s, a) \cdot D_{t+i}(s, a) \right)$$

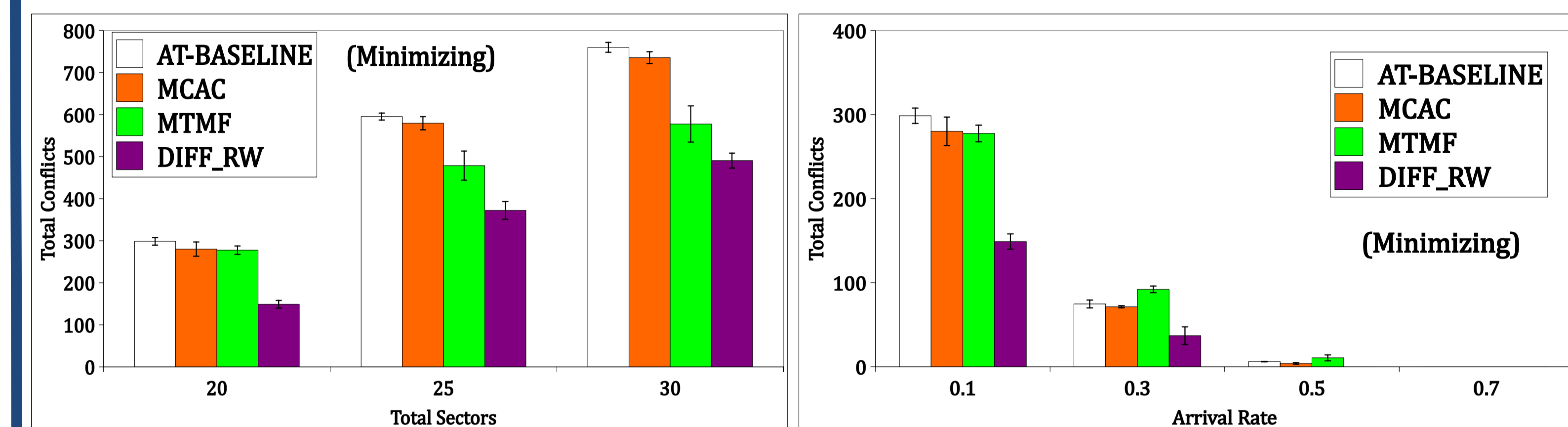
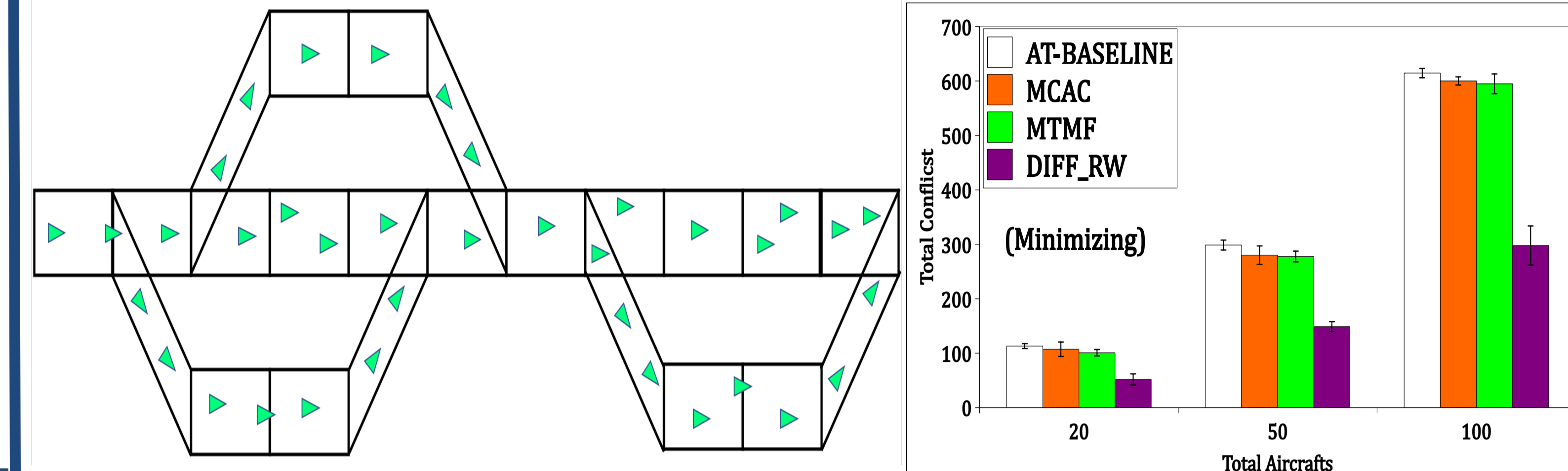
### ❖ Policy gradient :

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\mathbf{s}_{0:\infty}, \mathbf{a}_{0:\infty}} \left[ \sum_{t=0}^{\infty} \sum_{s \in S} \sum_{a \in A} n_t(s, a) \cdot \nabla_{\theta} \log \pi_{\theta}(a | s_t) \cdot R_t^{dr} \right]$$

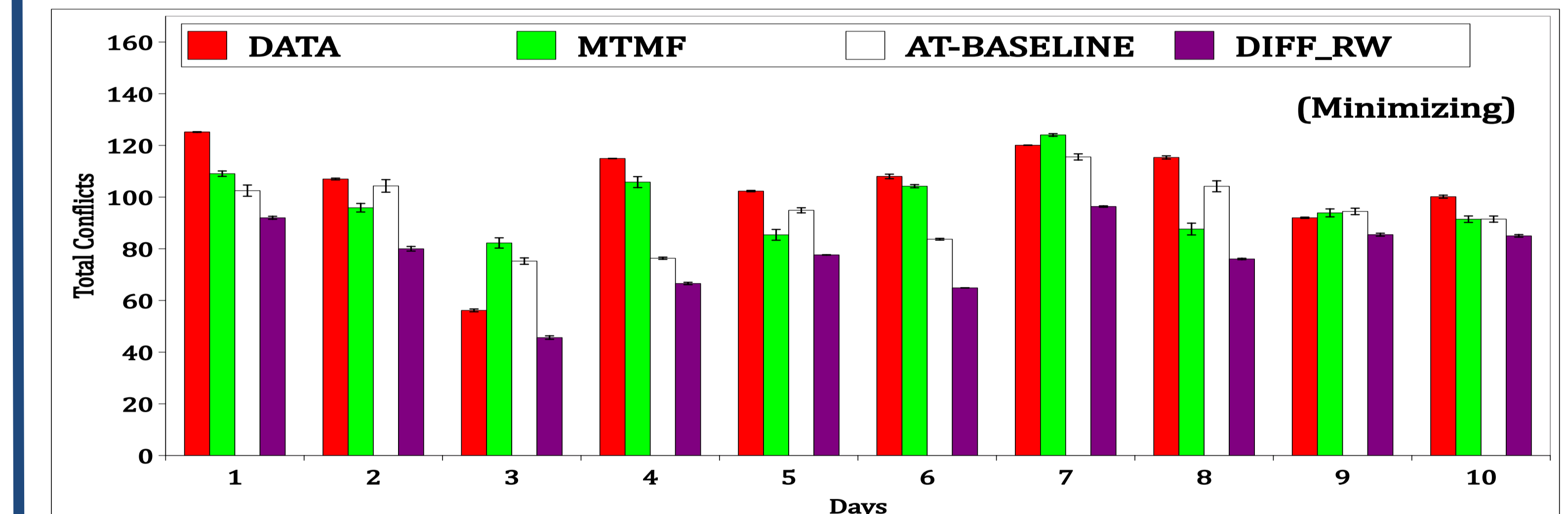
## Experiments

### ❖ Air Traffic Control :

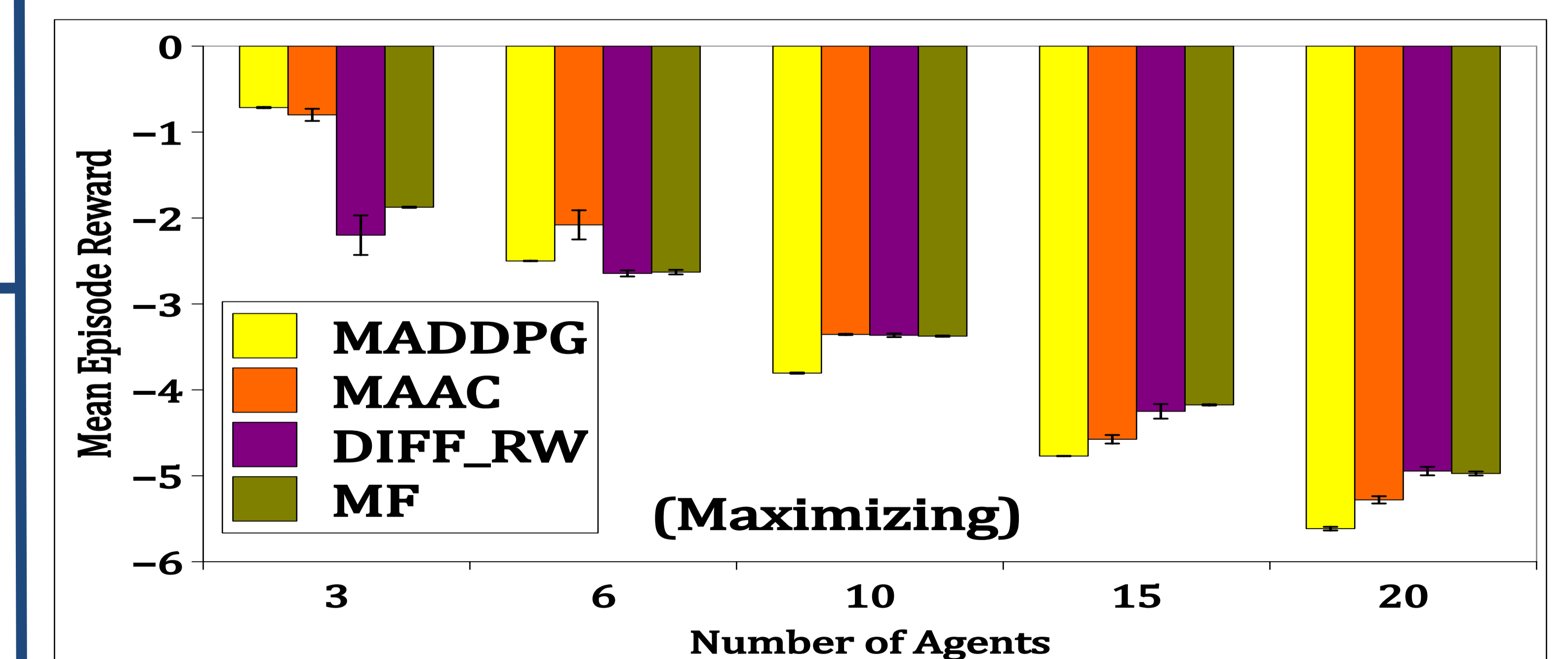
#### Synthetic Data:



#### Real world dataset (1 month data):



### ❖ Cooperative Navigation :



### ❖ Acknowledgments :

This research is supported by the Agency for Science, Technology and Research (A\*STAR), Fujitsu Limited and the National Research Foundation Singapore as part of the A\*STAR-Fujitsu- SMU Urban Computing and Engineering Centre of Excellence.