

Credit Assignment and Fine-Tuning Enhanced Reinforcement Learning for Collaborative Spatial Crowdsourcing

Wei Chen, Yafei Li*, Baolong Mei, Guanglei Zhu, Jiaqi Wu, Mingliang Xu

Zhengzhou University

Abstract: Collaborative spatial crowdsourcing (CSC) leverages distributed workers' collective intelligence to accomplish spatial tasks. A central challenge is to efficiently assign suitable workers to collaborate on these tasks. Although mainstream reinforcement learning (RL) methods have proven effective in task allocation, they face two key obstacles: delayed reward feedback and non-stationary data distributions, both hindering optimal allocation and collaborative efficiency. To address these limitations, we propose CAFE (credit assignment and fine-tuning enhanced), a novel multiagent RL framework for spatial crowdsourcing. CAFE introduces a credit assignment mechanism that distributes rewards based on workers' contributions and spatiotemporal constraints, coupled with bi-level meta-optimization to jointly optimize credit assignment and RL policy. To handle nonstationary spatial task distributions, CAFE employs an adaptive fine-tuning procedure that efficiently adjusts credit assignment parameters while preserving collaborative knowledge. Experiments on two real-world datasets validate the effectiveness of our framework, demonstrating superior performance in terms of task completion and equitable reward redistribution.

1. Introduction

To illustrate CSC, consider the example shown in Fig.1, which involves four workers (A-D) and two tasks (X, Y). The initial configuration assigns workers A and B to tasks X and Y respectively, while workers C and D remain unallocated. The scenario is characterized by worker heterogeneity: worker C is exclusively qualified for task Y, while worker D possesses the versatility to perform either task. Task completion times directly impact payment structures, and the assignment of worker D to either task X or Y yields different compensation outcomes. This study aims to optimize the platform's revenue through strategic task allocation. Notably, as the scale of workers and tasks expands in real-world applications, the computational complexity of CSC increases exponentially.

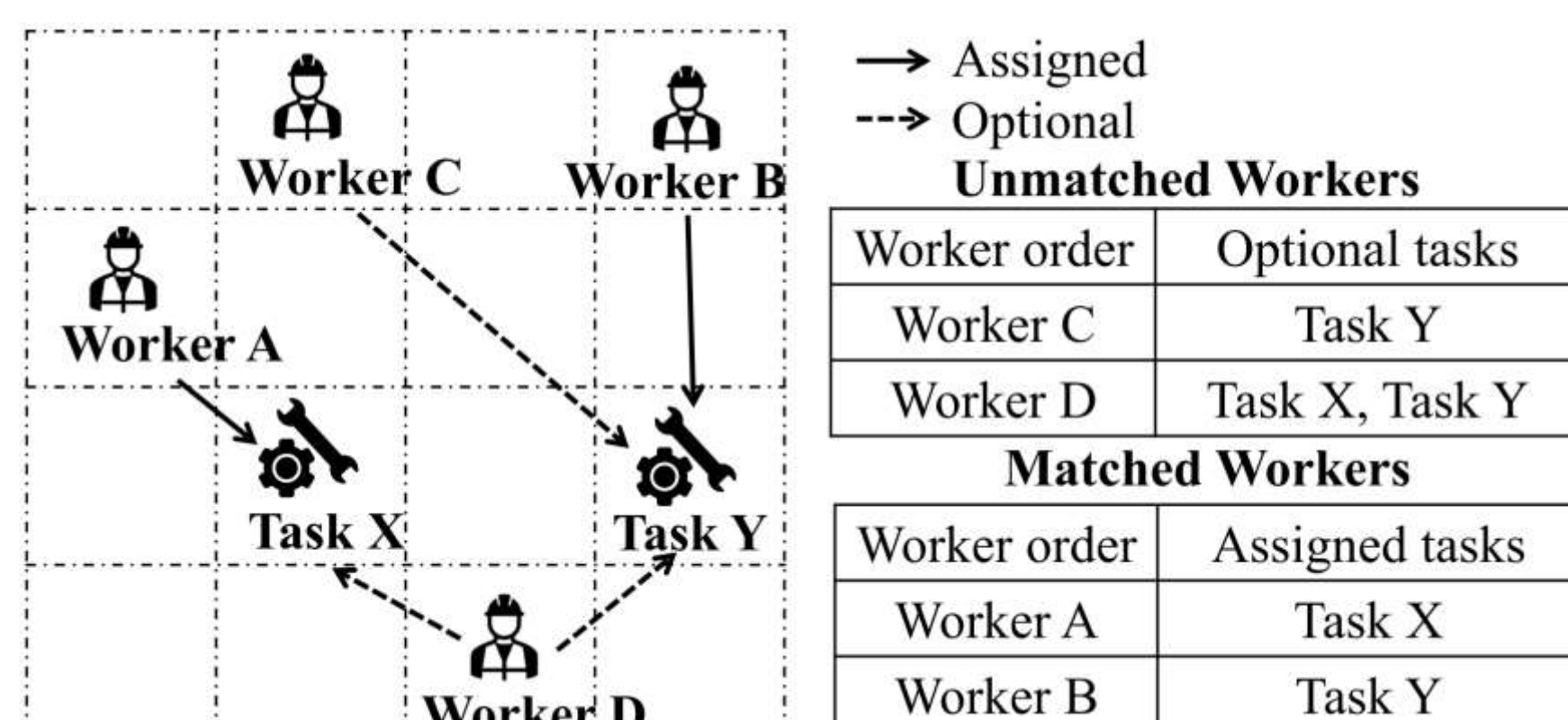


Fig. 1. An example of collaborative spatial crowdsourcing.

2. Problem Definition

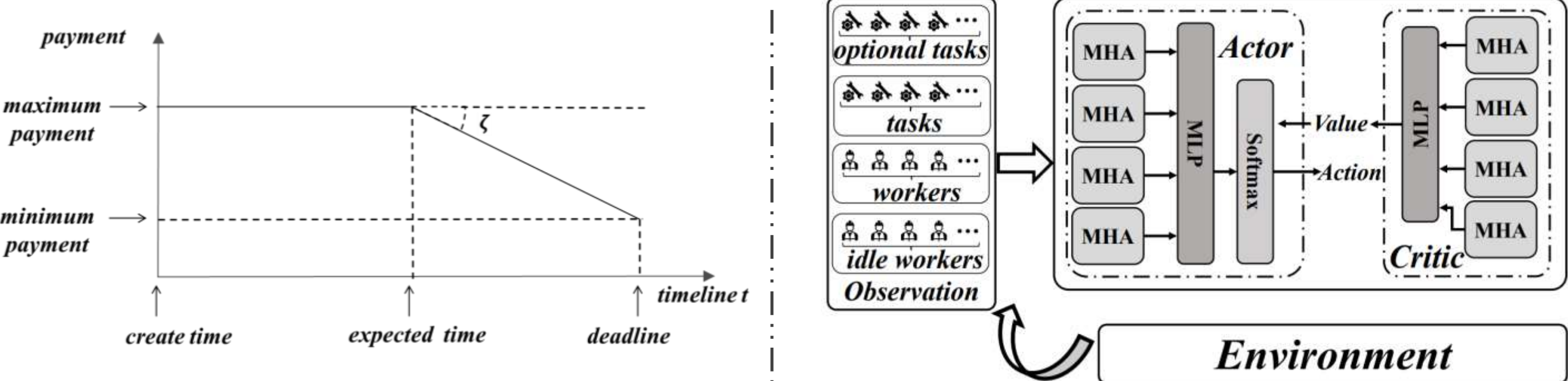


Fig. 2. Payment setting of task.

Fig. 3. Actor-Critic network structure.

Given a stream of tasks T and a set of works W , CSC problem aims to find the optimal assignment M to maximize the global revenue P .

$$P = \max \sum_{\tau \in T} \rho^\tau$$

3. Method

Reward Redistribution

To comprehensively evaluate how individual workers' actions affect overall revenue, we propose a reward redistribution mechanism that quantifies each worker's contribution.

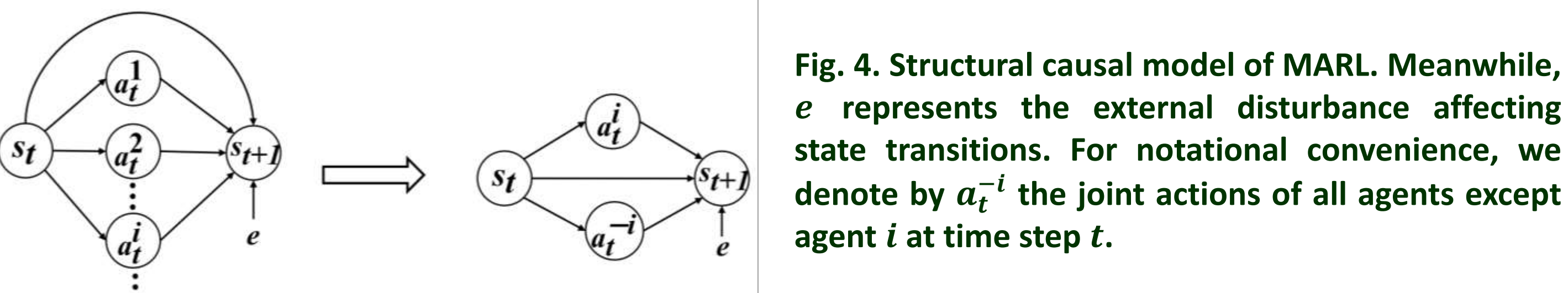


Fig. 4. Structural causal model of MARL. Meanwhile, e represents the external disturbance affecting state transitions. For notational convenience, we denote by a_t^{-i} the joint actions of all agents except agent i at time step t .

To comprehensively evaluate how individual workers' actions affect overall revenue, we propose a reward redistribution mechanism that quantifies each worker's contribution. Within this model, directional arrows represent inherent causal dependencies between components. To quantify how each agent's action a_t^i causally impacts the subsequent state s_{t+1} , we employ conditional mutual information (CMI):

$$CMI_t^i = I(s_{t+1}; a_t^i | s_t, a_t^{-i}) \approx I(o_{t+1}; a_t^i | o_t, a_{t,o}^{-i})$$

To comprehensively evaluate how individual workers' actions affect overall revenue, we propose a reward redistribution mechanism that quantifies each worker's contribution. Within this model, directional arrows represent inherent causal dependencies between components. To quantify how each agent's action a_t^i causally impacts the subsequent state s_{t+1} , we employ conditional mutual information (CMI):

$$CMI_t^i \approx I(z_\rho; a_t^i | o_t, a_{t,o}^{-i}) \approx D_{KL}[q(z_\rho | o_t, a_t^i, a_{t,o}^{-i}) || q(z_\rho^{-i} | o_t, a_{t,o}^{-i})].$$

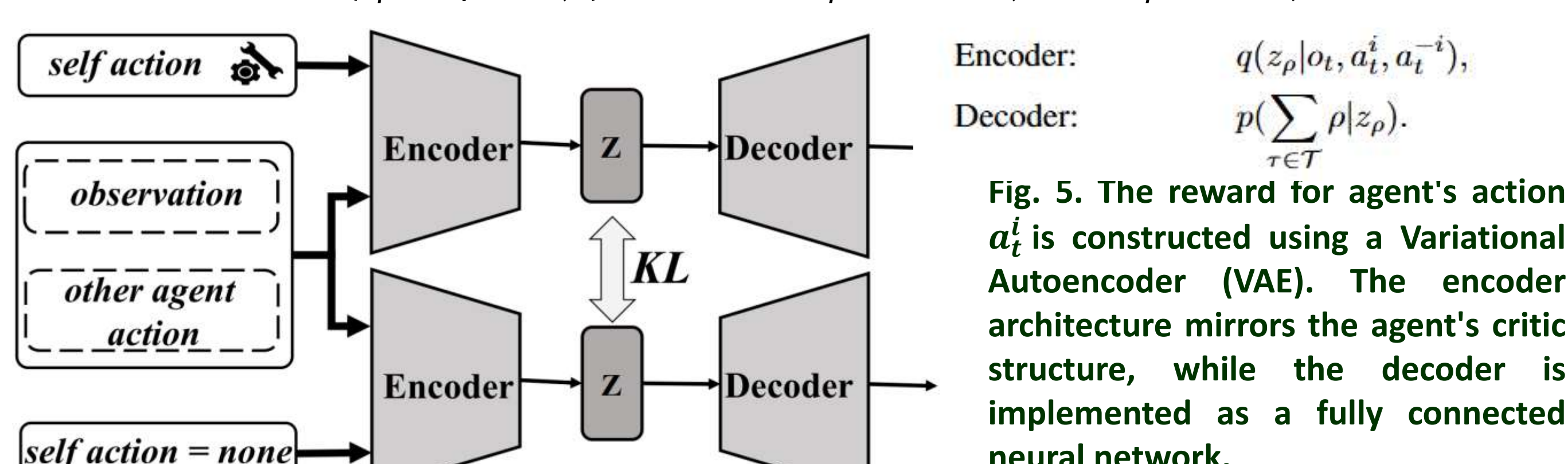


Fig. 5. The reward for agent's action a_t^i is constructed using a Variational Autoencoder (VAE). The encoder architecture mirrors the agent's critic structure, while the decoder is implemented as a fully connected neural network.

$$J = -D_{KL}[q(z | o_t, a_t^i) || p(z)] + E_{z \sim q(z)} [\log p(\sum_{\tau \in T} (\rho | z_\rho))]$$

To mitigate this misalignment, we introduce a regularization term into the reward function:

$$r^i(\phi) = \phi_1 \cdot CMI_t^i + \phi_2$$

Implicitly Learning

To optimize the hyperparameters in the reward function, we propose a novel bi-level optimization method based on implicit gradient tune the reward function's hyperparameters. This approach enables the concurrent optimization of hyperparameters during the reinforcement learning process.

$$\text{inner-level: } \theta^* = \arg \max_{\theta} J(\theta),$$

$$\text{outer-level: } \phi^* := \arg \max_{\phi} [J(\theta, \phi) - \mathcal{L}_{\Lambda}]$$

which is expressed as:

we use \mathcal{L}_{Λ} to regulate the long-term reward,

$$\mathcal{L}_{\Lambda}(\phi) = E_{\lambda \sim \Lambda} \left[\left(\sum_{t=1}^T r(\phi) - \eta \cdot \sum_{t=1}^T \rho \right)^2 \right]$$

We define $F(\phi)$ as the outer-level loss function, and its gradient can calculate by implicit gradient:

$$\frac{dF}{d\phi} \approx \frac{\partial J}{\partial \phi} - \frac{\partial J}{\partial \theta^i} \cdot \left(\frac{\partial^2 J}{\partial \theta^i \partial \theta^j} \right)^{-1} \cdot \frac{\partial^2 J}{\partial \theta^i \partial \phi} - \frac{d\mathcal{L}_{\Lambda}}{d\phi}$$

Swift Parameter Refinement

Firstly, we use the Taylor series expansion to the second-order term

$$\mathcal{L}_{\Lambda}(\phi) \approx \mathcal{L}_{\Lambda}(\hat{\phi}) + (\phi - \hat{\phi})^T \cdot \left[\frac{d\mathcal{L}_{\Lambda}(\phi)}{d\phi} \right]_{\phi=\hat{\phi}} + (\phi - \hat{\phi})^T \cdot \left[\frac{d^2 \mathcal{L}_{\Lambda}(\phi)}{d\phi^2} \right]_{\phi=\hat{\phi}} \cdot (\phi - \hat{\phi})$$

$$\frac{d\mathcal{L}_{\Lambda}}{d\phi} \approx \left[\frac{d\mathcal{L}_{\Lambda}(\phi)}{d\phi} \right]_{\phi=\hat{\phi}} + \left[\frac{d^2 \mathcal{L}_{\Lambda}(\phi)}{d\phi^2} \right]_{\phi=\hat{\phi}} \cdot (\phi - \hat{\phi})$$

Consider the concept of implicit gradients for optimizing the network parameters.

$$\frac{d\mathcal{L}_{\Lambda}(\phi^*)}{d\phi} = 0$$

$$\left[\frac{d\mathcal{L}_{\Lambda}(\phi)}{d\phi} \right]_{\phi=\hat{\phi}} + \left[\frac{d^2 \mathcal{L}_{\Lambda}(\phi)}{d\phi^2} \right]_{\phi=\hat{\phi}} \cdot (\phi^* - \hat{\phi}) = 0$$

$$\phi^* = \hat{\phi} - \left[\frac{d\mathcal{L}_{\Lambda}(\phi)}{d\phi} \right]_{\phi=\hat{\phi}} \cdot \left[\frac{d^2 \mathcal{L}_{\Lambda}(\phi)}{d\phi^2} \right]_{\phi=\hat{\phi}}^{-1}$$

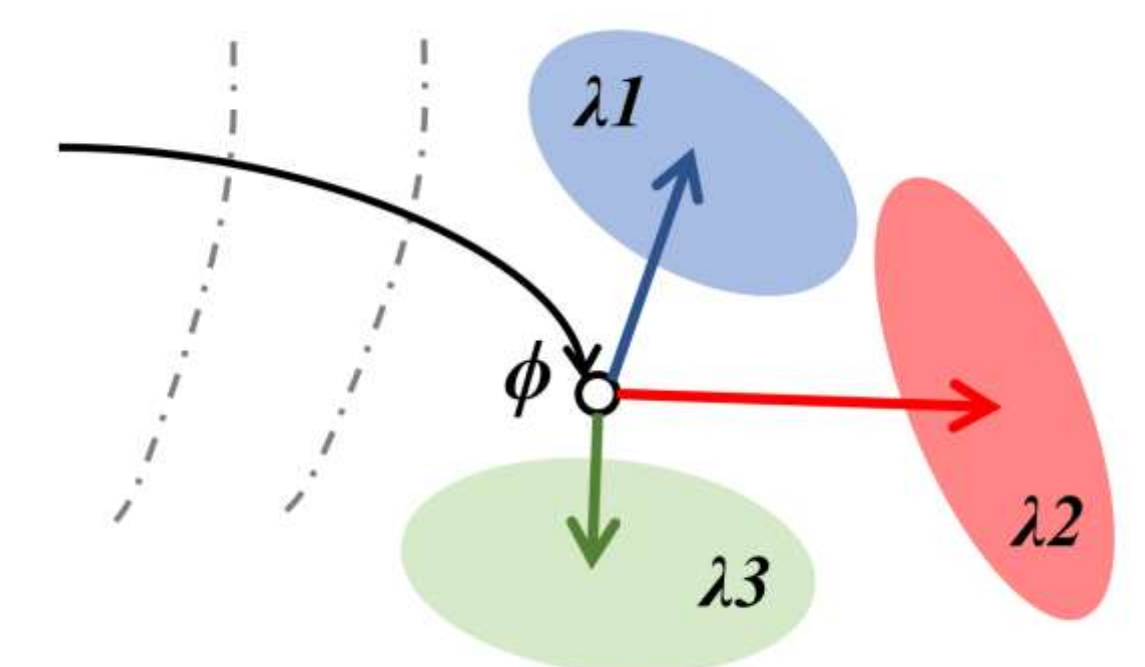
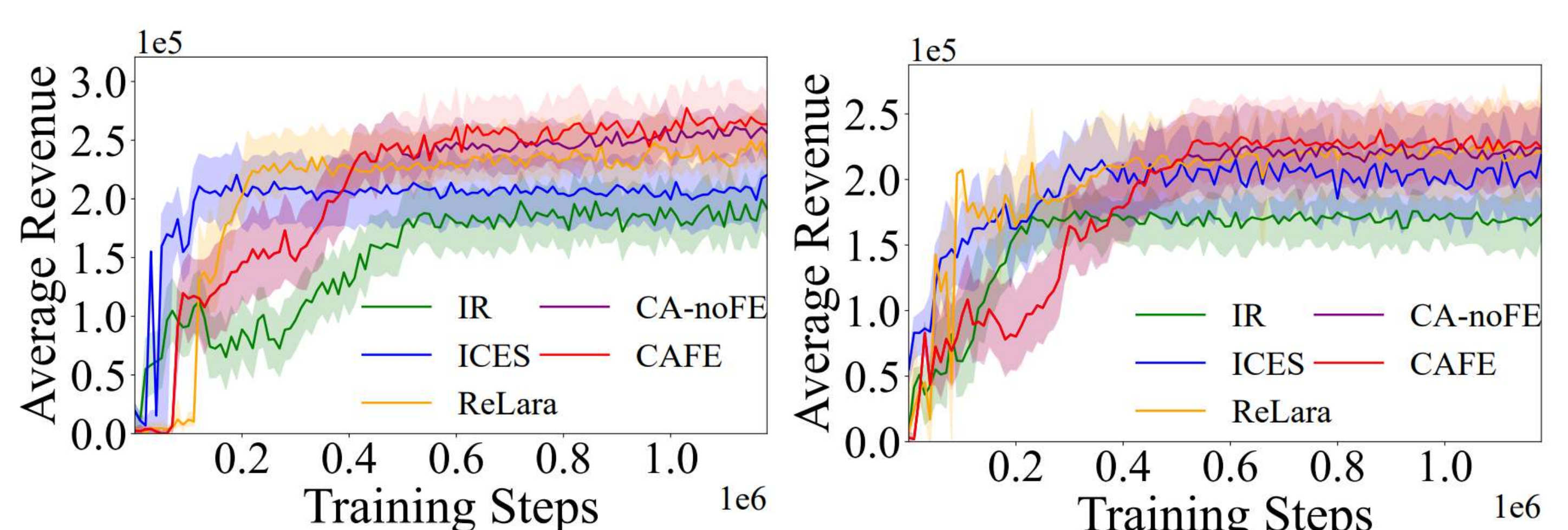


Fig. 6. Fine-tuning the parameters to adapt to changes in crowdsourcing tasks quickly.

4. Experiments

Comparison of Reward Redistribution



(a) Chengdu datasets

(b) Haikou datasets

Fig. 7. Comparison of reward redistribution approaches in model training.

Revenue performance across varying numbers of tasks

Dataset	Methods	$N_\tau = 50$	$N_\tau = 100$	$N_\tau = 150$	$N_\tau = 200$	$N_\tau = 250$	$N_\tau = 300$	$N_\tau = 350$
Chengdu	Greedy	6.53 \pm 0.16	12.14 \pm 0.17	17.77 \pm 0.37	21.85 \pm 0.15	26.07 \pm 0.39	29.30 \pm 0.28	32.94 \pm 0.81
	IPPO	6.47 \pm 0.22	12.42 \pm 0.26	17.25 \pm 0.42	23.74 \pm 0.80	27.01 \pm 0.36	34.40 \pm 1.16	40.09 \pm 1.25
	PAU	6.57 \pm 0.22	11.65 \pm 0.20	18.36 \pm 0.88	24.18 \pm 0.79	26.94 \pm 0.81	33.57 \pm 0.92	39.88 \pm 1.86
	CA-noFE	6.59 \pm 0.18	12.52 \pm 0.31	18.10 \pm 0.55	24.40 \pm 0.80	27.43 \pm 0.73	34.55 \pm 1.06	41.14 \pm 1.53
	CAFE	6.69 \pm 0.17	12.67 \pm 0.22	18.48 \pm 0.80	24.66 \pm 0.71	27.63 \pm 0.74	35.37 \pm 0.98	41.55 \pm 1.42
Haikou	Greedy	6.42 \pm 0.14	11.78 \pm 0.40	17.03 \pm 0.22	21.52 \pm 0.39	26.16 \pm 0.50	29.18 \pm 0.39	31.25 \pm 0.42
	IPPO	6.17 \pm 0.18	12.23 \pm 0.43	17.30 \pm 0.32	22.84 \pm 0.32	27.19 \pm 0.24	32.48 \pm 0.88	35.03 \pm 1.38
	PAU	6.68 \pm 0.15	11.88 \pm 0.19	18.25 \pm 0.79	23.28 \pm 0.90	25.95 \pm 0.90	32.87 \pm 1.01	35.48 \pm 0.88
	CA-noFE	6.64 \pm 0.15	12.14 \pm 0.31	18.52 \pm 0.59	23.28 \pm 0.90	27.55 \pm 0.38	33.21 \pm 0.57	36.09 \pm 1.50
	CAFE	6.71 \pm 0.15	12.43 \pm 0.35	18.51 \pm 0.47	23.40 \pm 0.87	28.37 \pm 0.45	34.15 \pm 0.69	36.99 \pm 0.97

Revenue performance across varying numbers of workers

Dataset	Methods	$N_w = 50$	$N_w = 100$	$N_w = 150$	$N_w = 200$	$N_w = 250$	$N_w = 300$	$N_w = 350$
Chengdu	Greedy	7.63 \pm 0.42	15.67 \pm 0.51	19.81 \pm 0.38	21.85 \pm 0.15	22.88 \pm 0.37	23.76 \pm 0.36	24.22 \pm 0.28
	IPPO	10.45 \pm 0.34	20.82 \pm 1.14	23.48 \pm 1.01	23.74 \pm 0.80	22.79 \pm 0.31	24.42 \pm 0.20	23.71 \pm 0.42
	PAU	11.48 \pm 0.49	21.54 \pm 1.05	23.96 \pm 0.90	24.18 \pm 0.79	24.73 \pm 0.85	24.17 \pm 0.62	24.94 \pm 0.55
	CA-noFE	11.94 \pm 0.62	21.82 \pm 0.89	24.24 \pm 1.04	24.40 \pm 0.80	25.29 \pm 0.93	24.90 \pm 0.51	25.63 \pm 0.60
	CAFE	12.17 \pm 0.49	21.85 \pm 1.02	24.60 \pm 0.93	24.66 \pm 0.71	25.66 \pm 0.86	25.69 \pm 0.58	26.56 \pm 0.59
Haikou	Greedy	7.71 \pm 0.30	15.89 \pm 0.44	19.75 \pm 0.45	21.52 \pm 0.39	23.10 \pm 0.28	23.65 \pm 0.26	23.98 \pm 0.38
	IPPO	15.92 \pm 1.90	21.80 \pm 0.60	22.16 \pm 1.24	22.84 \pm 0.32	24.12 \pm 0.25	26.68 \pm 0.33	27.63 \pm 0.28
	PAU	16.38 \pm 1.73	22.58 \pm 0.38	23.14 \pm 1.02	23.28 \pm 0.90	23.98 \pm 0.91	26.89 \pm 0.60	30.31 \pm 1.21
	CA-noFE	16.60 \pm 1.83	22.87 \pm 0.45	23.21 \pm 0.79	23.28 \pm 0.90	24.51 \pm 0.80	27.26 \pm 0.49	29.87 \pm 0.55
	CAFE	16.62 \pm 1.76	23.32 \pm 0.33	23.60 \pm 0.87	23.40 \pm 0.87	24.75 \pm 0.91	27.65 \pm 0.45	30.02 \pm 0.54