

Fair and Risk-Averse Worker Selection in Mobile Crowdsourcing via Mean-Variance Bandits

Ziqun Chen, Kechao Cai, *Member, IEEE*, Jinbei Zhang, *Member, IEEE*, and John C.S. Lui, *Fellow, IEEE*

Abstract—Mobile crowdsourcing is recognized as a promising paradigm for harnessing collective wisdom to solve problems. The crowdsourcing platform assigns tasks to a group of workers (e.g., vehicles, sensors) and receives the outcomes from them. In this paper, we develop a novel combinatorial mean-variance bandit model for worker selection and characterize the quality of workers with mean-variance measures to consider both the reward and risk in task assignment. Moreover, we introduce a utility-based fairness constraint to ensure a fair selection of workers. We consider two types of feedback: quantitative feedback, where the platform can receive numerical rewards from the selected workers, and preference feedback, where the platform only receives a qualitative performance ranking of the selected workers. We define the mean-variance regret and fairness regret and propose novel bandit algorithms that balance the fairness guarantee and mean-variance optimization in worker selection under both feedback settings. We prove that our algorithms achieve sublinear expected mean-variance regret and fairness regret. Through extensive simulations, we validate that our algorithms can fairly select workers while maximizing the mean-variance of selected workers.

Index Terms—Mobile crowdsourcing, mean-variance bandit, utility-based fairness, quantitative feedback, preference feedback.

I. INTRODUCTION

With the rapid development of mobile devices and wireless networks, mobile crowdsourcing has come into focus as a novel problem-solving paradigm in recent years [1]. As shown in Figure 1, the crowdsourcing platform assigns tasks like traffic detection [2], air quality measurement [3], and street view collection [4] to a group of workers (e.g., vehicles, sensors) selected from the candidates, who are willing to participate in the crowdsourcing. The selected workers then execute the tasks and send the outcomes to the platform. It is favorable for the platform to select workers with higher quality to improve task performance and user satisfaction. However, practically, the quality of workers is usually unknown in advance. Therefore, the platform has to experience a learning process to estimate the quality of workers. At this point, the platform faces a dilemma: explore various workers to learn about their quality or exploit the historical observations by selecting those with previously demonstrated high performance.

Ziqun Chen, Kechao Cai, and Jinbei Zhang are with the School of Electronics and Communication Engineering, Shenzhen Campus of Sun Yat-sen University, Shenzhen, China, 518107. (e-mail: chenqz35@mail2.sysu.edu.cn; {caikch3, zhjinbei}@mail.sysu.edu.cn). John C.S. Lui is with the Department of Computer Science & Engineering, The Chinese University of Hong Kong (CUHK), Hong Kong (e-mail: cslui@cse.cuhk.edu.hk). (Corresponding author: Kechao Cai).

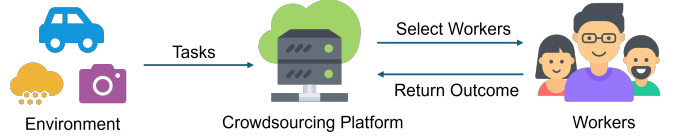


Fig. 1. Mobile Crowdsourcing System.

The emerging field of combinatorial multi-armed bandit (CMAB) offers a promising framework for worker selection problems in mobile crowdsourcing. Unlike prior studies that typically rely on heuristic, static optimization [5], or deep learning-based approaches [6], CMAB provides an on-line decision-making paradigm that can dynamically balance exploration and exploitation. This enables the platform to adaptively learn workers' performance over time and make informed decisions under uncertainty without requiring prior knowledge of worker quality or extensive offline training. Specifically, we view the platform as a player and the workers as arms. Each arm generates a random reward that reflects the corresponding worker's task performance and follows an unknown distribution. The player's goal is to maximize the cumulative reward based on the previous observations and decision history.

However, directly applying the standard CMAB model to mobile crowdsourcing can be problematic, as it fails to address the multifaceted challenges inherent in real-world task assignment. Firstly, it does not consider the risk of assigning tasks to workers with unstable performance. Conventional CMAB algorithms are designed to maximize the expected reward under the assumption of risk neutrality and stable arm performance. However, in mobile crowdsourcing, worker performance often exhibits significant variability due to mobility, network latency, and heterogeneous device capabilities. Ignoring such fluctuations leads to risk-prone task allocation decisions: workers with high average performance but large variance may still be frequently selected, resulting in unstable outcomes and a degraded overall platform experience. In particular, it may result in frequent worker switches, where tasks are repeatedly reassigned to different workers at each round, thereby incurring significant communication overhead and switching costs [7]. Moreover, the conventional CMAB model ignores the interests of some workers, resulting in an unfair worker selection. Consider a bandit algorithm that tries to maximize the worker's utility: it simply learns which worker has the highest quality and constantly assigns the task to that worker, even if other workers are almost equally good (albeit slightly lower in quality). Such an approach results in a winner-takes-

all allocation, where many skillful workers receive insufficient tasks, leading to substantial income disparities. Over time, this may cause workers to lose interest in the platform, ultimately harming its reputation and long-term performance. Thus, to build a sustainable platform, a good policy should guarantee fairness among workers and ensure that workers with comparable skill levels have similar probabilities of being assigned tasks. Finally, for tasks such as content moderation or image annotation, it is challenging for the crowdsourcing platform to figure out the exact numerical rewards, known as quantitative feedback, from workers' outcomes, as it requires domain experts to establish detailed evaluation metrics beforehand [8]. In contrast, eliciting preference feedback through qualitative ranking of workers' task performance is generally more practical and easier than direct quantitative evaluation.

Main Contributions. Motivated by the above discussions, this paper makes the following key contributions:

(i) We develop a novel combinatorial mean-variance bandit model (CMVB) for mobile crowdsourcing that considers both reward and risk in task assignment, while guaranteeing utility-based fairness among workers. Specifically, we measure worker quality using the *mean-variance* of their performance and define the utility of a worker as a function of its mean-variance. Then we impose a *utility-based fairness* constraint to ensure each worker is selected with a probability proportional to its utility.

(ii) We design two fair mean-variance bandit algorithms, *Mean-Variance Fair learning with Quantitative feedback* (MVFQ) and *Mean-Variance Fair learning with Preference feedback* (MVFP), tailored for *quantitative feedback* and the more challenging *preference feedback*, respectively. We introduce *mean-variance regret* and *fairness regret* as performance metrics to evaluate the effectiveness of our algorithms.

(iii) Our theoretical analysis and experimental results demonstrate that the proposed algorithms achieve sublinear upper bounds for both types of expected regret, outperforming existing methods by fairly selecting workers based on their utility while maximizing the mean-variance of selected workers.

II. FAIR MEAN-VARIANCE BANDIT WITH QUANTITATIVE FEEDBACK

In this section, we present the details of our CMVB model with quantitative feedback and design a fair and risk-averse algorithm for worker selection in mobile crowdsourcing.

A. Problem Formulation

We consider a mobile crowdsourcing system with a crowdsourcing platform and a set of workers $[K] = \{1, 2, \dots, K\}$, which can provide computing and sensing service. Let $[T] := \{1, 2, \dots, T\}$ denote the set of decision rounds. At round $t \in [T]$, the platform selects a subset S_t of L ($L \leq K$) workers from $[K]$ to execute the tasks. Once a selected worker $i \in S_t$ has returned the outcome, the platform will generate a numerical reward $r_{i,t} \in [0, 1]$ based on the worker's task performance. We interpret the reward as quantitative feedback from the worker, as it reflects the quality of the outcome

in that round. Due to the system uncertainty (e.g., worker mobility, transmission/processing oscillation), the reward $r_{i,t}$ is a random variable that follows an unknown distribution with expectation $\mu_i = \mathbb{E}[r_{i,t}]$ and variance $\sigma_i^2 = \text{Var}[r_{i,t}]$. The μ_i and σ_i^2 characterize the reliability and stability of worker i , respectively. When the reward variance σ_i^2 is large, the worker could still perform poorly even with a high expected reward μ_i . To capture the tradeoff between a worker's reliability and stability, we formulate the worker selection problem as a CMVB model. The quality of each worker is measured by a *mean-variance* metric, defined for worker i as $w_i = \rho\mu_i - \sigma_i^2$, where $\rho > 0$ is the risk tolerance parameter that controls the balances the two objectives of high expected reward and low variance.

To ensure that similar levels of workers obtain comparable treatment, we define a utility function $f(\cdot) > 0$ that maps the mean-variance of a worker to a positive utility value. Then we enforce a utility-based fairness constraint that the probability p_i of selecting worker i is proportional to its utility $f(w_i)$. Formally, we have

$$\frac{p_i}{f(w_i)} = \frac{p_{i'}}{f(w_{i'})}, \quad \forall i \neq i', i, i' \in [K]. \quad (1)$$

The introduced utility function $f(\cdot)$ allows us to tailor the fairness criterion for different scenarios. We have two assumptions on the utility function $f(\cdot)$.

Assumption 1. The utility of each worker is bounded such that (i) $\exists \lambda > 0$ and $\min_w f(w) \geq \lambda$, (ii) $\forall w_1, w_2, \frac{f(w_1)}{f(w_2)} \leq \frac{K-1}{L-1}$ for $L > 1$.

Assumption 2. The utility function f is M -Lipschitz continuous, i.e., there exists a positive constant $M > 0$, such that $\forall w_1, w_2, |f(w_1) - f(w_2)| \leq M |w_1 - w_2|$.

We now show that there is a unique optimal fair policy that fulfills the fairness constraints in (1) in the following theorem.

Theorem 1. For any $w_i, i \in [K]$ and any choice of utility function $f(\cdot) > 0$ under Assumption 1 and 2, there exist a unique optimal fair policy $\mathbf{p}^* = \{p_1^*, p_2^*, \dots, p_K^*\}$ such that

$$p_i^* = \frac{L f(w_i)}{\sum_{i'=1}^K f(w_{i'})}, \quad \forall i \in [K], \quad (2)$$

that satisfies the utility-based fairness constraints in (1).

Due to the space limit, all the proofs of the theorems are provided in the Appendix of the full version of the paper [9]. Theorem 1 implies that the optimal fair policy in our model is no longer selecting a fixed optimal set of L workers as in classical bandit problems, but a probability distribution on all the possible sets $S_t \subseteq [K]$, $|S_t| = L$. Specifically, we characterize a worker selection policy with a probabilistic selection vector $\mathbf{p}_t = \{p_{1,t}, p_{2,t}, \dots, p_{K,t}\}$ where $p_{i,t} \in [0, 1]$ is the probability of selecting worker $i \in [K]$ at round t , and $\sum_{i=1}^K p_{i,t} = L$ since only L workers can be selected at each round.

To measure the gap in the expected mean-variance of workers between the optimal fair policy and the deployed policy, we define the *mean-variance regret* as follows:

$$\text{MR}_T = \sum_{t=1}^T \max \left\{ \sum_{i=1}^K p_i^* w_i - \sum_{i=1}^K p_{i,t} w_i, 0 \right\}, \quad (3)$$

which quantifies the speed of the mean-variance optimization of the deployed policy. Especially, we only consider the non-negative part at each round in (3), since a less fair policy could have a larger mean-variance than the optimal fair policy by consistently selecting high-quality workers at the expense of violating fairness constraints. This would result in a negative mean-variance gap and thus be meaningless. Moreover, we also require a measure to quantify its fairness guarantee. We define the *fairness regret* that measures the cumulative 1-norm distance between the optimal fair policy \mathbf{p}^* and the deployed policy \mathbf{p}_t as follows:

$$\text{FR}_T = \sum_{t=1}^T \sum_{i=1}^K |p_i^* - p_{i,t}|. \quad (4)$$

The fairness regret measures the overall violation of the utility-based fairness constraints. Our objective is to design selection policies that have both *sublinear expected mean-variance regret* and *sublinear expected fairness regret* with respect to the number of rounds T , where the expectations are taken over the randomness in both the worker selections and the mean-variance measure. By doing so, we can approach the optimal fair policy and select the high-quality workers while ensuring fairness among all the workers in the long run.

It is important to point out that both Assumption 1 and Assumption 2 are necessary for designing CMVB algorithms as stated in the following theorem and remark.

Theorem 2. *For any CMVB algorithm, if either Assumption 1(i) or Assumption 2 does not hold, the lower bound of the fairness regret is linear; in other words, there exists a CMVB instance with linear expected fairness regret $O(T)$.*

Remark 1. Assumption 1(ii) ensures that the selection probability $p_{i,t}$ in the form of $Lf(\cdot)/\sum_{a=1}^K f(\cdot)$ is constrained in $[0, 1]$.

B. Algorithm Description

Algorithm 1 shows the details of our *Mean-Variance Fair learning with Quantitative feedback* (MVFQ) algorithm, which follows the principle of optimism in the face of uncertainty. At each round t , we incorporate a randomized rounding scheme (RRS) from [10] to ensure that each worker i is independently selected with probability $p_{i,t}$, while maintaining the constraint that exactly L workers are chosen, thereby preserving fairness in the selection process. In line 2, RRS takes \mathbf{p}_t as input and generates a set of selected workers S_t such that $\mathbb{E}[\mathbb{1}_{\{i \in S_t\}}] = p_{i,t}$, where $\mathbb{1}_{\{\cdot\}}$ is the indicator function. Given the observed reward and the number of samples $n_{i,t}$ up to round t , we compute the empirical mean-variance estimate $\hat{w}_{i,t}$ for each worker $i \in [K]$, using empirical expected reward $\hat{\mu}_{i,t}$ and variance $\hat{\sigma}_{i,t}^2$. Then, using both UCB (Upper Confidence

Algorithm 1 MVFQ Algorithm

Input: $f(\cdot)$, T , L , K , $\rho > 0$, $\delta = \frac{1}{LT}$.

Output: \mathbf{p}_t .

Init: Select each worker in $[K]$ once with $\lceil K/L \rceil$ rounds.

Initial selection policy: $p_{i, \lceil K/L \rceil + 1} = 1/K$, $\forall i \in [K]$.

- 1: **for** $t = \lceil K/L \rceil + 1$ to T **do**
 - 2: Select workers in $S_t = \text{RRS}(L, \mathbf{p}_t)$
 - 3: Receive reward $r_{i,t}$ from $i \in S_t$
 - 4: **for** $i \in [K]$ **do**
 - 5: $n_{i,t} = \sum_{\tau=1}^t \mathbb{1}_{\{i \in S_\tau\}}$
 - 6: $\hat{\mu}_{i,t} = \sum_{\tau=1}^t \mathbb{1}_{\{i \in S_\tau\}} r_{i,\tau} / n_{i,t}$
 - 7: $\hat{\sigma}_{i,t}^2 = \sum_{\tau=1}^t \mathbb{1}_{\{i \in S_\tau\}} (r_{i,\tau} - \hat{\mu}_{i,t})^2 / n_{i,t}$
 - 8: $\hat{w}_{i,t} = \rho \hat{\mu}_{i,t} - \hat{\sigma}_{i,t}^2$
 - 9: $u_{i,t} = \hat{w}_{i,t} + (\rho + 3) \sqrt{\frac{\log(8KT/\delta)}{2n_{i,t}}}$
 - 10: $l_{i,t} = \hat{w}_{i,t} - (\rho + 3) \sqrt{\frac{\log(8KT/\delta)}{2n_{i,t}}}$
 - 11: $\mathcal{C}_t = \{\tilde{\mathbf{w}} \in \mathbb{R}^K \mid \forall i \in [K], \tilde{w}_i \in [l_{i,t}, u_{i,t}]\}$
 - 12: $\tilde{\mathbf{w}}_t = \arg \max_{\tilde{\mathbf{w}} \in \mathcal{C}_t} \sum_{i \in [K]} \frac{Lf(\tilde{w}_i)}{\sum_{i' \in [K]} f(\tilde{w}_{i'})} \tilde{w}_i$
 - 13: Compute $p_{i,t+1} = \frac{Lf(\tilde{w}_{i,t})}{\sum_{i' \in [K]} f(\tilde{w}_{i',t})}$ for $i \in [K]$
-

Bound) estimates $u_{i,t}$ and LCB (Lower Confidence Bound) estimates $l_{i,t}$ of all workers, we can construct a confidence region \mathcal{C}_t (see line 11) which contains the actual mean-variance vector $\mathbf{w} := (w_i)_{i \in [K]}$ with high probability. We find a vector $\tilde{\mathbf{w}}_t := (\tilde{w}_{i,t})_{i \in [K]}$ in the confidence region \mathcal{C}_t that maximizes the expected mean-variance of a fair policy as shown in line 12. Finally, according to Theorem 1, we update the selection policy $p_{i,t+1}$ of worker $i \in [K]$ as $\frac{Lf(\tilde{w}_{i,t})}{\sum_{i'=1}^K f(\tilde{w}_{i',t})}$ to satisfy the utility-based fairness constraints, which is limited to the interval $[0, 1]$ under Assumption 1(ii).

We present the sublinear expected mean-variance regret and fairness regret upper bounds of MVFQ in the following theorem.

Theorem 3. *For MVFQ in Algorithm 1, the expected mean-variance regret is upper bounded by $O(\rho\sqrt{LKT\log T})$ and the expected fairness regret of MVFQ is upper bounded by $O\left(\frac{\rho ML}{\lambda} \sqrt{LKT\log T}\right)$.*

In Theorem 3, the factor $\frac{ML}{\lambda}$ in the fairness regret bound comes from Assumption 1(i) and Assumption 2 on the utility function $f(\cdot)$, while the mean-variance regret bound does not depend on Assumption 1(i) and Assumption 2, making it tighter up to logarithmic factors.

III. FAIR MEAN-VARIANCE BANDIT WITH PREFERENCE FEEDBACK

In this section, we consider the more challenging CMVB model with preference feedback, where the platform receives only the qualitative ranking of workers based on task performance. To address this, we propose another algorithm to maximize the mean-variance of workers and ensure utility-based fairness among workers.

A. Problem Formulation

To characterize the relative preferences of workers in a given ranking, we adopt the Multinomial-Logit (MNL) probability model with unknown parameter $\theta = (\theta_1, \dots, \theta_i, \dots, \theta_K)$ where each component $\theta_i \in (0, 1]$ corresponds to the potential strength of a worker i . At round t , the platform assigns the tasks to a subset $S_t \subseteq [K]$ of L workers and receives as *preference feedback* a partial task performance ranking π_t of $1 \leq m \leq L - 1$ workers from S_t . Let $\pi_{i,t}$ denote the worker at i -th position in π_t and $\mathcal{P}_{S_t}^m$ denote the set of permutations of any m -subset of S_t . The ranking π_t , of the form $(\pi_{1,t} \succ \dots \succ \pi_{m,t})$, is drawn without replacement from the MNL probability model on S_t , where $\pi_{i,t} \succ \pi_{j,t}, 1 \leq i < j \leq m$ means that the i -th worker is preferred over the j -th worker in π_t . More formally, the probability of a partial ranking $\pi_t \in \mathcal{P}_{S_t}^m$ under MNL model is

$$\mathbb{P}(\pi_t \in \mathcal{P}_{S_t}^m \mid S_t, \theta) = \prod_{i=1}^m \frac{\theta_{\pi_{i,t}}}{\sum_{j \in S_t \setminus \{\theta_{\pi_{1,t}}, \dots, \theta_{\pi_{i-1,t}}\}} \theta_{\pi_{j,t}}}. \quad (5)$$

In particular, preference feedback reduces to winner feedback, i.e., a single worker who is preferred over all other workers in S_t when $m = 1$ and a full performance ranking of workers in S_t when $m = L - 1$. Our model extends the original dueling bandits problem by simultaneously providing preference feedback on several workers rather than just two.

To model the quality of workers under preference feedback, we maintain a pairwise preference matrix \mathbf{Q} where each element $q_{ij} = \frac{\theta_i}{\theta_i + \theta_j}$ denotes the pairwise probability of worker i being preferred over worker j . Then we introduce the Borda score [11] to measure the reliability of worker $i \in [K]$: $b_i = \frac{1}{K-1} \sum_{j \in [K]} q_{ij} \mathbb{1}_{\{j \neq i\}}$, which is the average preference probability that a worker is preferred over another worker in $[K]$ chosen uniformly at random. Analogous to the quantitative feedback setting, to incorporate both the reliability and stability of workers, we define the *Borda mean-variance* of worker i as $w_i^b = \rho b_i - \sigma_i^{2,b}$, where $\sigma_i^{2,b}$ denotes the variance of the Borda score b_i and $\rho > 0$ is a risk tolerance parameter.

Furthermore, we adopt the utility function $f(\cdot) > 0$ in Section II-A and reformulate the utility-based fairness constraint in (1) using Borda mean-variance:

$$\frac{p_i}{f(w_i^b)} = \frac{p_{i'}}{f(w_{i'}^b)}, \quad \forall i \neq i', i, i' \in [K], \quad (6)$$

which ensures that the probability of selecting each worker is proportional to its utility. Similarly to Theorem 1, there still exists a unique optimal fair policy $\mathbf{p}^{*,b} = \{p_1^{*,b}, p_2^{*,b}, \dots, p_K^{*,b}\}$ in the form of

$$p_i^{*,b} = \frac{Lf(w_i^b)}{\sum_{i'=1}^K f(w_{i'}^b)}, \quad \forall i \in [K], \quad (7)$$

that satisfies the utility-based fairness constraints in (6). Moreover, we define the *Borda mean-variance regret* as the expected Borda mean-variance difference between the optimal policy and the deployed policy:

$$\text{BMR}_T = \sum_{t=1}^T \max \left\{ \sum_{i=1}^K p_i^{*,b} w_i^b - \sum_{i=1}^K p_{i,t} w_i^b, 0 \right\}, \quad (8)$$

Algorithm 2 MVFP Algorithm

Input: $f(\cdot), T, L, K, \rho > 0, \delta = \frac{1}{LT}, \alpha \geq \left(\frac{K^2 - 2K + L}{L^2 - L} \right)^4$.

Output: \mathbf{p}_t .

Init: Pairwise winning matrix: $\mathbf{V} = [v_{ij,0}] \leftarrow [0]_{K \times K}$.

Initial selection policy: $p_{i,1} = 1/K, \forall i \in [K]$.

```

1: for  $t = 1$  to  $T$  do
2:   Select workers in  $S_t = \text{RRS}(L, \mathbf{p}_t^b)$ 
3:   Receive preference feedback  $\pi_t$ 
4:   for  $i = 1$  to  $m$  do
5:      $v_{\pi_{i,t},j,t} = v_{\pi_{i,t},j,t-1} + 1, \forall j \in S_t \setminus \{\pi_{1,t}, \dots, \pi_{i,t}\}$ 
6:   Estimate Borda scores, for  $i \in [K]$ ,
       
$$\tilde{b}_{i,t} = \frac{\mathbb{1}_{\{i \in S_t\}}}{(K-1)p_{i,t}} \sum_{j \in [K]} \frac{\mathbb{1}_{\{j \in S_t\}} \mathbb{1}_{\{j \neq i\}}}{p_{j,t}} \frac{v_{ij,t}}{v_{ij,t} + v_{ji,t}}$$

       (Assuming  $\frac{x}{0} = \frac{1}{2}, x \in \mathbb{R}$ )
7:   for  $i \in [K]$  do
8:      $\hat{b}_{i,t} = \sum_{\tau=1}^t \tilde{b}_{i,\tau} / t$ 
9:      $\hat{\sigma}_{i,t}^{2,b} = \sum_{\tau=1}^t (\tilde{b}_{i,\tau} - \hat{b}_{i,\tau})^2 / t$ 
10:     $\hat{w}_{i,t}^b = \rho \hat{b}_{i,t} - \hat{\sigma}_{i,t}^{2,b}$ 
11:     $u_{i,t}^b = \hat{w}_{i,t}^b + \alpha(\rho + 3) \sqrt{\frac{2 \log(4/\delta)}{t}}$ 
12:     $l_{i,t}^b = \hat{w}_{i,t}^b - \alpha(\rho + 3) \sqrt{\frac{2 \log(4/\delta)}{t}}$ 
13:     $\mathcal{C}_t^b = \{\tilde{\mathbf{w}}^b \in \mathbb{R}^K \mid \forall i \in [K], \tilde{w}_i^b \in [l_{i,t}^b, u_{i,t}^b]\}$ 
14:     $\tilde{\mathbf{w}}_t^b = \arg \max_{\tilde{\mathbf{w}}^b \in \mathcal{C}_t^b} \sum_{i \in [K]} \frac{Lf(\tilde{w}_i^b)}{\sum_{i' \in [K]} f(\tilde{w}_{i'}^b)} \tilde{w}_i^b$ 
15:    Compute  $p_{i,t+1} = \frac{Lf(\tilde{w}_{i,t}^b)}{\sum_{i' \in [K]} f(\tilde{w}_{i',t}^b)}$  for  $i \in [K]$ 

```

and *Borda fairness regret* as the cumulative 1-norm distance between the optimal fair policy and the deployed policy:

$$\text{BFR}_T = \sum_{t=1}^T \sum_{i=1}^K |p_i^{*,b} - p_{i,t}|. \quad (9)$$

We aim to minimize both Borda mean-variance regret and Borda fairness regret. Note that Theorem 2 and Remark 1 still hold for the preference feedback setting as the Borda score can be interpreted as the expected reward of a worker. Therefore, we impose Assumption 1 and Assumption 2 on the utility function and design CMVB algorithms with sublinear Borda fairness regret.

B. Algorithm Description

To deal with preference feedback, we introduce a novel variant of MVFQ, named *Mean-Variance Fair learning with Preference feedback* (MVFP) algorithm, detailed in Algorithm 2. At round t , we first construct an unbiased Borda score estimator for each worker $i \in [K]$. Denote \mathbf{V} as a pairwise winning matrix, where each element $v_{ij,t}$ is the number of times worker i is preferred over j up to round t . Thanks to Lemma 1 in [12], we can extract pairwise comparisons from π_t by treating each comparison independently as shown in line 5 and obtain the unbiased pairwise preference estimator $\hat{q}_{ij,t} = \frac{v_{ij,t}}{v_{ij,t} + v_{ji,t}}$ for $i, j \in [K]$. Next, along with the selection policy \mathbf{p}_t , we establish an unbiased estimate $\tilde{b}_{i,t}$ of Borda score

$b_{i,t}$ for each worker $i \in [K]$ (see line 6). Then we compute the empirical mean-variance estimate $\hat{w}_{i,t}^b$ of the Borda score using the time-average Borda score $\hat{b}_{i,t}$ and variance $\hat{\sigma}_{i,t}^{2,b}$. We further construct a confidence region \mathcal{C}_t^b with UCB $u_{i,t}^b$ and LCB $l_{i,t}^b$ estimates of all workers in line 13. Similarly to the MVFQ algorithm, we identify a parameter $\tilde{w}_t^b = (\tilde{w}_{i,t}^b)_{i \in [K]}$ within the confidence region to maximize the expected mean-variance subject to the utility-based fairness constraints. Finally, we update the selection policy \mathbf{p}_{t+1} for the next round according to the optimal fair policy in (7).

We now state the sublinear expected Borda mean-variance regret and fairness regret upper bounds of MVFP algorithm in the following theorem.

Theorem 4. *For MVFP in Algorithm 2, the expected Borda mean-variance regret is upper bounded by $O(\alpha \rho L \sqrt{T \log T})$ and the expected Borda fairness regret of MVFP is upper bounded by $O\left(\frac{\alpha \rho M L^2}{\lambda} \sqrt{T \log T}\right)$.*

Algorithm 2 adopts a confidence bound that differs in form from that used in Algorithm 1. This difference arises because the observed preference feedback is *not independent and identically distributed (non-i.i.d.)*, and the estimated Borda scores $\hat{b}_{i,t}$ for $i \in [K]$ are inherently correlated with the qualities of other workers, rendering the derivation strategy employed in Algorithm 1 inapplicable.

Note that Theorem 4 does not capture the impact of the number of workers m in ranking π_t on the regret bounds. Obtaining tight regret bounds that reflect this impact is challenging and remains an open problem. Nevertheless, the regret bounds in Theorem 4 hold for all $1 \leq m \leq L - 1$. The sublinear bounds remain valid even under high reward variance or extremely sparse feedback (i.e., small m). As shown empirically in Section V, the parameter m introduces only a minor additive penalty, demonstrating the algorithm's robustness across varying levels of feedback granularity.

IV. APPLICATIONS

In this section, we provide detailed examples to demonstrate how our proposed MVFQ and MVFP algorithms can be applied in real-world scenarios.

A. MVFQ Application: On-Demand Delivery in Crowd-sourced Logistics

Consider a delivery crowdsourcing platform that allocates parcel or food delivery tasks to a pool of couriers. Each courier returns the completion time and route information after finishing a delivery, and the platform evaluates task outcomes and generates a numerical reward based on metrics such as delivery punctuality, route efficiency, and customer rating.

The platform often makes assignment decisions under uncertainty. The delivery quality of couriers is unknown beforehand and can fluctuate due to mobility, traffic conditions, and network delays. Simply selecting couriers with the highest average performance can lead to biased and unstable task allocation, as some workers exhibit large performance variance across regions or time periods. Moreover, persistently favoring a few top couriers results in an unfair workload distribution

and discourages long-term participation from other qualified workers.

The MVFQ algorithm quantifies the trade-off between a courier's reliability and the stability of their performance using a mean-variance metric. The utility function maps each worker's mean-variance measure to a utility value representing the courier's capability in terms of reliability and stability, and it can be customized by the platform. By incorporating the utility-based fairness constraint, couriers with comparable utility are assigned tasks with similar probabilities, resulting in both high service quality and a balanced, sustainable workload distribution across the platform.

B. MVFP Application: Image Captioning by Human Annotators

Consider a mobile crowdsourcing platform that collects image captions for large-scale datasets used in tourism recommendation or e-commerce search. In this setting, human annotators, such as smartphone users, are recruited to generate textual descriptions for assigned photos through a mobile application or web interface.

In practice, caption quality is initially unknown and inherently difficult to quantify, as it depends on factors such as linguistic fluency, semantic relevance, and creativity, making it challenging for the platform to assign explicit numerical rewards. Instead, the platform can estimate the annotators' quality by comparing candidate captions from different workers to determine which is more relevant to the given image across rounds. To handle such qualitative feedback, the MVFP algorithm quantifies pairwise comparisons among annotators as preference feedback and introduces the Borda mean-variance metric capturing both the reliability and stability of workers' performance over time.

Similar to MVFQ, MVFP employs a utility function that maps each annotator's Borda mean-variance measure to a utility value. By integrating with a utility-based fairness constraint, annotators with comparable caption reliability and stability are assigned tasks with similar probabilities, thereby promoting fairness.

V. EXPERIMENTS

In this section, we conduct experiments to demonstrate the effectiveness of our algorithms.

We consider a mobile crowdsourcing system with $K = 6$ workers where the platform selects $L = 2$ workers at each round. The rewards of each worker follow a Gaussian distribution¹ with expectation in $\{0.3, 0.5, 0.7, 0.9, 0.8, 0.6, 0.4\}$ and variance in $\sigma^2 = \{0.3, 0.5, 0.7, 0.9, 0.8, 0.6, 0.4\}$. We use the utility function $f(w) = \frac{K-L}{(1+e^{-cw})(L-1)} + 1$, adapted from Sigmoid function to satisfy both Assumption 1 and Assumption 2. The Sigmoid-based formulation provides a smooth and bounded mapping that stabilizes the selection probability, which is critical for ensuring robust empirical performance. The parameter c further controls the gradient

¹While in the paper we assume the reward distributions to be bounded in $[0, 1]$, our results can be naturally extended to sub-Gaussian distributions.

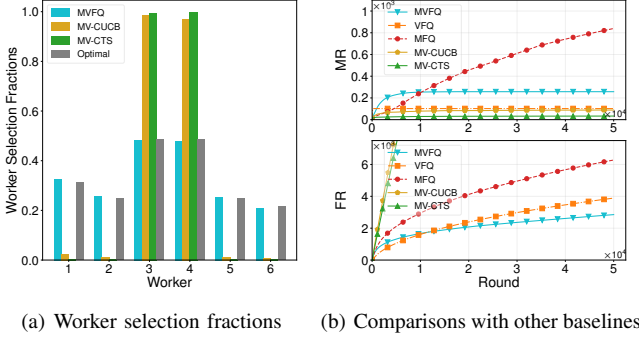


Fig. 2. Experiment results of MVFQ algorithm.

of the function, enabling flexible adjustment of the utility’s sensitivity to worker merit. We set $c = 4$ and $\rho = 1$ in the following, and all results are averaged over 100 runs.

We first examine the fairness of different algorithms with quantitative feedback. For comparison, we implement two other mean-variance bandit algorithms, MV-CUCB and MV-CTS, which are adapted from MV-UCB [13] and MVTS [14], respectively, to accommodate the combinatorial selection structure. Figure 2(a) illustrates the average worker selection fractions of MV-CUCB, MV-CTS, the optimal fair policy, and MVFQ. Each bar corresponds to the fraction of rounds a worker is chosen over $T = 5 \times 10^4$ rounds by a specific algorithm. As observed, MV-CUCB and MV-CTS are unfair by mainly selecting the worker 3 and 4 with the highest mean-variances, neglecting the potential utility of other workers. In contrast, the MVFQ algorithm can converge to the optimal fair policy. This observation shows the effectiveness of MVFQ in achieving utility-based fairness, ensuring that each worker receives a selection allocation proportional to its utility.

As shown in Figure 2(b), MV-CUCB and MV-CTS consistently exhibit smaller mean-variance regret and larger fairness regret when compared to MVFQ. This observation suggests that MV-CUCB and MV-CTS outperform MVFQ in mean-variance regret but substantially violate the utility-based fairness constraints, incurring higher fairness regret. Moreover, the mean-variance/fairness regrets of MVFQ increase sublinearly in T , aligning with the bounds we derived in Theorem 3. Then we introduce two other algorithms, *Mean Fair learning with Quantitative feedback* (MFQ) and *Variance Fair learning with Quantitative feedback* (VFQ), to model extreme cases where the parameter ρ is large or small. Specifically, MFQ estimates the quality of workers based only on the expected reward, while VFQ relies on the reward variance. We show that the VFQ and MFQ algorithms still achieve both sublinear mean-variance and fairness regret.

Finally, we evaluate the performance of the MVFP algorithm with preference feedback. With the same utility function, we consider the case $K = 10$, $L = 8$ and the MNL model with the parameters $\theta = (0.1, 0.3, 0.45, 0.6, 0.7, 0.95, 0.9, 0.75, 0.4, 0.15)$. In Figure 3(a), we compare MVFP ($m = 4$) with several baselines, including MV-CUCB-P and MV-CTS-P adapted from MV-UCB [13] and MVTS [14], as well as *Mean Fair Learning with Preference Feedback* (MFP) and *Variance Fair*

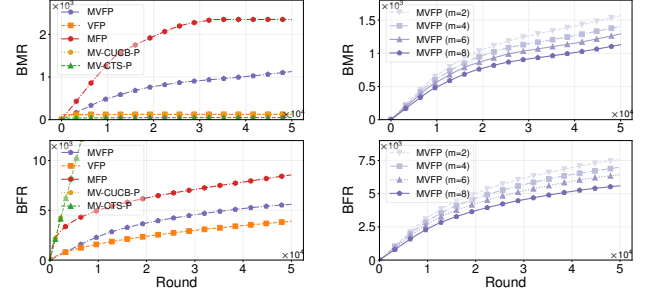


Fig. 3. Experiment results of MVFP algorithm.

Learning with Preference Feedback (VFP), which represent two extreme cases where ρ is large or small. Specifically, MFP evaluates worker quality solely based on the Borda score, while VFP relies on the variance of the Borda score. The experimental results show MVFP achieves sublinear Borda mean-variance/fairness regret, which is consistent with the regret bounds we derived in Theorem 4 and demonstrate that our method remains robust across a wide range of ρ , achieving superior performance compared with other baselines. In Figure 3(b), the Borda mean-variance/fairness regret of MVFP scales down as m increases. This is because a larger m allows MVFP to estimate Borda scores more accurately by incorporating more pairwise preference information.

Notably, FairX-UCB and FairX-TS proposed in [15] are not directly applicable to our settings, as they can only select a single arm at each round and do not incorporate preference feedback. Other fair bandit algorithms adopt different fairness metrics, making them unsuitable for direct comparison with our proposed algorithms. We have also conducted additional experiments with larger worker sets ($L = 10$, $K = 50$) to demonstrate the scalability and robustness of the proposed algorithms. Due to space limitations, the experimental results are provided in the Appendix of the full paper [9].

VI. CONCLUSION

In this paper, we develop a fair and risk-averse worker selection framework for mobile crowdsourcing via the CMVB model. We consider quantitative feedback and preference feedback settings, respectively and design novel algorithms that achieve both sublinear expected mean-variance regret and fairness regret. Our experimental results demonstrate that the proposed algorithms effectively balance utility-based fairness and mean-variance optimization of worker selection under both types of feedback.

REFERENCES

- [1] N. Zhao, Y. Pei, Y.-C. Liang, and D. Niyato, “A deep reinforcement learning-based contract incentive mechanism for mobile crowdsourcing networks,” *IEEE Transactions on Vehicular Technology*, 2023.
- [2] A. Thiagarajan, L. Ravindranath, K. LaCurtis, S. Madden, H. Balakrishnan, S. Toledo, and J. Eriksson, “Vtrack: accurate, energy-aware road traffic delay estimation using mobile phones,” in *Proceedings of the 7th ACM conference on embedded networked sensor systems*, 2009, pp. 85–98.

- [3] J. Xu, Z. Luo, C. Guan, D. Yang, L. Liu, and Y. Zhang, "Hiring a team from social network: Incentive mechanism design for two-tiered social mobile crowdsourcing," *IEEE Transactions on Mobile Computing*, vol. 22, no. 8, pp. 4664–4681, 2022.
- [4] Y. Zhang, Q. Liu, H. Wang, D. Chen, and K. Han, "Crowdsourcing live high definition map via collaborative computation in automotive edge computing," *IEEE Transactions on Vehicular Technology*, 2024.
- [5] X. Ding, J. Guo, G. Sun, and D. Li, "Optimizing worker selection in collaborative mobile crowdsourcing," *IEEE Internet of Things Journal*, vol. 11, no. 4, pp. 7172–7185, 2023.
- [6] Z. Zhan, Y. Wang, P. Duan, A. M. V. V. Sai, Z. Liu, C. Xiang, X. Tong, W. Wang, and Z. Cai, "Enhancing worker recruitment in collaborative mobile crowdsourcing: A graph neural network trust evaluation approach," *IEEE Transactions on Mobile Computing*, vol. 23, no. 10, pp. 10 093–10 110, 2024.
- [7] X. Liu, M. Derakhshani, L. Mihaylova, and S. Lambotharan, "Risk-aware contextual learning for edge-assisted crowdsourced live streaming," *IEEE Journal on Selected Areas in Communications*, vol. 41, no. 3, pp. 740–754, 2022.
- [8] S. Wang and Z. Shao, "Green dueling bandits," in *ICC 2023 - IEEE International Conference on Communications*, 2023, pp. 5129–5134.
- [9] Z. Chen, K. Cai, J. Zhang, and J. C. S. Lui, "Fair and risk-averse worker selection in mobile crowdsourcing via mean-variance bandits," 2025. [Online]. Available: <https://github.com/MLCL-SYSU/FairMVBandit/blob/main/FairMVlearning.pdf>
- [10] R. Gandhi, S. Khuller, S. Parthasarathy, and A. Srinivasan, "Dependent rounding and its applications to approximation algorithms," *Journal of the ACM (JACM)*, vol. 53, no. 3, pp. 324–360, 2006.
- [11] K. Jamieson, S. Katariya, A. Deshpande, and R. Nowak, "Sparse dueling bandits," in *Artificial Intelligence and Statistics*. PMLR, 2015, pp. 416–424.
- [12] A. Saha and A. Gopalan, "Pac battling bandits in the plackett-luce model," in *Algorithmic Learning Theory*. PMLR, 2019, pp. 700–737.
- [13] S. Vakili and Q. Zhao, "Risk-averse multi-armed bandit problems under mean-variance measure," *IEEE Journal of Selected Topics in Signal Processing*, vol. 10, no. 6, pp. 1093–1111, 2016.
- [14] Q. Zhu and V. Tan, "Thompson sampling algorithms for mean-variance bandits," in *International Conference on Machine Learning*. PMLR, 2020, pp. 11 599–11 608.
- [15] L. Wang, Y. Bai, W. Sun, and T. Joachims, "Fairness of exposure in stochastic bandits," in *International Conference on Machine Learning*. PMLR, 2021, pp. 10 686–10 696.
- [16] J. Bretagnolle and C. Huber, "Estimation des densités: risque minimax," *Séminaire de probabilités de Strasbourg*, vol. 12, pp. 342–363, 1978.

APPENDIX

A. Proof of Theorem 1

Proof. According to (1), the optimal fair policy \mathbf{p}^* satisfies the following utility-based fairness constraint:

$$\frac{p_i^*}{f(w_i)} = \frac{p_{i'}^*}{f(w_{i'})}, \quad \forall i \neq i', i, i' \in [K], \quad (10)$$

which correspond to $K - 1$ linearly independent equations of \mathbf{p}^* . Moreover, because only L workers can be selected at each round, there is an additional linear equation $\sum_{i=1}^K p_i^* = L$ that is linearly independent of the other $K - 1$ ones. Then we have K linearly independent equations on K unknowns in $\mathbf{p}^* = \{p_1^*, p_2^*, \dots, p_K^*\}$. Therefore, the optimal fair policy \mathbf{p}^* is unique. By solving this system of linear equations, we have

$$p_i^* = \frac{L f(w_i)}{\sum_{i'=1}^K f(w_{i'})}, \quad \forall i \in [K]. \quad (11)$$

This completes the proof of Theorem 1. \square

B. Proof of Theorem 2

Proof. We first prove that the lower bound on fairness regret is linear without Assumption 1(i) by constructing two CMVB instances and a 1-Lipschitz utility function $f(\cdot)$. For any bandit algorithm, we show that the sum of the expected fairness regrets of the two CMVB instances increases linearly in T . Consequently, we conclude that any bandit algorithm will incur a linear regret in T for at least one of the two CMVB instances.

The two instances can be defined as $x^1 = (\nu_1^1, \nu_2^1, \nu_3^1)$ and $x^2 = (\nu_1^2, \nu_2^2, \nu_3^2)$, where ν_i is the reward distribution of worker i . Each instance consists of three workers and the crowdsourcing platform selects a subset S_t of $L = 2$ workers at each round t . We assume that the reward of each worker in the two instances follows a Bernoulli distribution. The mean-variances of three workers in the first instance are $3\eta, 2\eta, 2\eta$, and the mean-variances of three workers in the second instance are $2\eta, 2\eta, 2\eta$, where $\eta \in (0, 1/3]$. The utility function $f(\cdot)$ is defined as an identity function. i.e., $f(w) = w$. Therefore, the optimal fair policy for the first instance is $\mathbf{p}^{*,1} = \{6/7, 4/7, 4/7\}$, and the optimal fair policy for the second instance is $\mathbf{p}^{*,2} = \{2/3, 2/3, 2/3\}$. For any bandit algorithm \mathcal{A} , the platform selects the workers stochastically according to a selection policy \mathbf{p}_t at each round t based on the history \mathcal{H}_t , which consists of all the previous selection policies, selected worker sets, and received feedback. We have $i \sim \mathbf{p}_t$, $r_{i,t} \sim \nu_i$ for $i \in S_t$. Then we can derive the lower bound of the expected fairness regret for the two instances as follows. For the first instance x^1 , we have

$$\begin{aligned} \mathbb{E} \left[\frac{1}{T} \text{FR}_T^1 \right] &= \mathbb{E} \left[\frac{1}{T} \sum_{t=1}^T \left(\left| p_{1,t} - \frac{6}{7} \right| + \left| p_{2,t} - \frac{4}{7} \right| + \left| p_{3,t} - \frac{4}{7} \right| \right) \right] \\ &\geq \mathbb{E} \left[\frac{1}{T} \sum_{t=1}^T p_{1,t} - \frac{6}{7} + \frac{1}{T} \sum_{t=1}^T p_{2,t} - \frac{4}{7} + \frac{1}{T} \sum_{t=1}^T p_{3,t} - \frac{4}{7} \right] \\ &= 2\mathbb{E} \left[\frac{1}{T} \sum_{t=1}^T p_{1,t} - \frac{6}{7} \right]. \end{aligned} \quad (12)$$

Similarly, for the second instance x^2 , we can derive the fairness regret lower bound as follows.

$$\mathbb{E} \left[\frac{1}{T} \text{FR}_T^2 \right] \geq 2\mathbb{E} \left[\frac{1}{T} \sum_{t=1}^T p_{1,t} - \frac{2}{3} \right]. \quad (13)$$

We consider a worker selection trace during the T rounds as $h = (\mathbf{p}_1, S_1, \mathbf{r}_1, \dots, \mathbf{p}_T, S_T, \mathbf{r}_T)$, where $\mathbf{r}_t := (r_{i,t})_{i \in [K]}$. Denote $\mathbb{H}^1, \mathbb{H}^2$ as the distributions of h for the CMVB instance x^1, x^2 using the algorithm \mathcal{A} , respectively. Then we have

$$\begin{aligned} \mathbb{E} \left[\frac{1}{T} \text{FR}_T^1 \right] + \mathbb{E} \left[\frac{1}{T} \text{FR}_T^2 \right] &\geq \frac{4}{21} \mathbb{P}^1 \left(\frac{1}{T} \sum_{t=1}^T p_{1,t} \leq \frac{16}{21} \right) + \frac{4}{21} \mathbb{P}^2 \left(\frac{1}{T} \sum_{t=1}^T p_{1,t} > \frac{16}{21} \right) \\ &\stackrel{(a)}{\geq} \frac{2}{21} \exp(-\text{KL}(\mathbb{H}^1, \mathbb{H}^2)), \end{aligned} \quad (14)$$

where (a) follows from the Bretagnolle-Huber inequality [16]. We can derive the upper bound of the KL divergence between \mathbb{H}^1 and \mathbb{H}^2 , $\text{KL}(\mathbb{H}^1, \mathbb{H}^2)$, as follows:

$$\begin{aligned}
\text{KL}(\mathbb{H}^1, \mathbb{H}^2) &= \mathbb{E}_{h \sim \mathbb{H}^1} \left[\log \frac{\mathbb{H}^1(h)}{\mathbb{H}^2(h)} \right] \leq \mathbb{E}_{h \sim \mathbb{H}^1} \left[\sum_{t=1}^T \sum_{i \in S_t} \log \frac{\nu_i^1(r_{i,t})}{\nu_i^2(r_{i,t})} \right] = \sum_{t=1}^T \mathbb{E}_{\mathbf{p}_t \sim \mathcal{A}^1} \mathbb{E}_{i \sim \mathbf{p}_t} [\text{KL}(\nu_i^1, \nu_i^2)] \\
&= \sum_{t=1}^T \mathbb{E}_{\mathbf{p}_t \sim \mathcal{A}^1} [p_{1,t} \text{KL}(\nu_1^1, \nu_1^2)] \\
&= \sum_{t=1}^T \mathbb{E}_{\mathbf{p}_t \sim \mathcal{A}^1} \left[p_{1,t} \left(3\eta \log \frac{3}{2} + (1-3\eta) \log \frac{1-3\eta}{1-2\eta} \right) \right] \\
&\leq 3T\eta \log \frac{3}{2},
\end{aligned} \tag{15}$$

where $\mathbf{p}_t \sim \mathcal{A}^1$ means that \mathbf{p}_t is sampled from the process of the algorithm \mathcal{A}^1 applied to the first CMVB instance.

Combining (15) with (14) and setting $\eta = 1/3T$, we have

$$\mathbb{E} \left[\frac{1}{T} \text{FR}_T^1 \right] + \mathbb{E} \left[\frac{1}{T} \text{FR}_T^2 \right] \geq \frac{2}{21} \exp \left(-3T\eta \log \frac{3}{2} \right) \geq 0.06, \tag{16}$$

which implies that at least one of the two CMVB instances incurs linear expected fairness regret. Therefore, we infer that at least one of the two CMVB instances incurs linear expected fairness regret.

Next, we prove that *the lower bound on fairness regret is linear without Assumption 2* by constructing two CMVB instances and a utility function $f(\cdot)$ where $\min_w f(w) = 1$. For any bandit algorithm, we demonstrate that the sum of the expected fairness regrets of the two CMVB instances grows linearly with respect to T . Thus, any bandit algorithm will result in linear regret in T for at least one of the two CMVB instances.

The two instances can be defined as $x^1 = (\nu_1^1, \nu_2^1, \nu_3^1)$ and $x^2 = (\nu_1^2, \nu_2^2, \nu_3^2)$, where ν_i is the reward distribution of worker i . Each instance consists of three workers and the platform selects a subset S_t of $L = 2$ workers at each round t . We assume that the reward of each worker in the two instances follows a Bernoulli distribution. The mean-variances of three workers in the first instance are $2\eta - 1, \eta - 1, \eta - 1$, and the mean-variances of three workers in the second instance are $\eta - 1, \eta - 1, \eta - 1$, where $\eta \in (0, 1/2)$. We use the utility function $f(\cdot)$ with the form

$$f(w) = \begin{cases} 1 & w \leq -1 \\ M(w+1) + 1 & w > -1 \end{cases}$$

where $M > 0$ is a positive constant to be defined later. Therefore, the optimal fair policy for the first instance is $\mathbf{p}^{*,1} = \{(4\eta M + 2)/(4\eta M + 3), (2\eta M + 2)/(4\eta M + 3), (2\eta M + 2)/(4\eta M + 3)\}$, and the optimal fair policy for the second instance is $\mathbf{p}^{*,2} = \{2/3, 2/3, 2/3\}$. For any bandit algorithm \mathcal{A} , the platform selects the workers with a selection policy \mathbf{p}_t at each round t based on the observation and decision history \mathcal{H}_t . Then we have $i \sim \mathbf{p}_t, r_{i,t} \sim \nu_i$ for $i \in S_t$.

For any algorithm, we can lower bound the expected fairness regret for the two instances as follows. For the first instance x^1 , we have

$$\begin{aligned}
\mathbb{E} \left[\frac{1}{T} \text{FR}_T^1 \right] &= \mathbb{E} \left[\frac{1}{T} \sum_{t=1}^T \left(\left| p_{1,t} - \frac{4\eta M + 2}{4\eta M + 3} \right| + \left| p_{2,t} - \frac{2\eta M + 2}{4\eta M + 3} \right| + \left| p_{3,t} - \frac{2\eta M + 2}{4\eta M + 3} \right| \right) \right] \\
&\geq \mathbb{E} \left[\left| \frac{1}{T} \sum_{t=1}^T p_{1,t} - \frac{4\eta M + 2}{4\eta M + 3} \right| + \left| \frac{1}{T} \sum_{t=1}^T p_{2,t} - \frac{2\eta M + 2}{4\eta M + 3} \right| + \left| \frac{1}{T} \sum_{t=1}^T p_{3,t} - \frac{2\eta M + 2}{4\eta M + 3} \right| \right] \\
&\geq 2\mathbb{E} \left[\left| \frac{1}{T} \sum_{t=1}^T p_{1,t} - \frac{4\eta M + 2}{4\eta M + 3} \right| \right].
\end{aligned} \tag{17}$$

Similarly, for the second instance x^2 , we can derive the fairness regret lower bound as follows,

$$\mathbb{E} \left[\frac{1}{T} \text{FR}_T^2 \right] \geq 2\mathbb{E} \left[\left| \frac{1}{T} \sum_{t=1}^T p_{1,t} - \frac{2}{3} \right| \right]. \tag{18}$$

We consider a worker selection trace during the T rounds as $h = (\mathbf{p}_1, S_1, \mathbf{r}_1, \dots, \mathbf{p}_T, S_T, \mathbf{r}_T)$. Denote \mathbb{H}^1 as the distribution of h when the algorithm \mathcal{A} is applied to the first MAB instance x^1 , while \mathbb{H}^2 as the distribution of h when the algorithm \mathcal{A} is applied to the second MAB instance x^2 . Then we have

$$\begin{aligned}
\mathbb{E} \left[\frac{1}{T} \text{FR}_T^1 \right] + \mathbb{E} \left[\frac{1}{T} \text{FR}_T^2 \right] &\geq \frac{4M\eta}{12M\eta + 9} \mathbb{P}^1 \left(\frac{1}{T} \sum_{t=1}^T p_{1,t} \leq \frac{10M\eta + 6}{12M\eta + 9} \right) + \frac{4M\eta}{12M\eta + 9} \mathbb{P}^2 \left(\frac{1}{T} \sum_{t=1}^T p_{1,t} > \frac{10M\eta + 6}{12M\eta + 9} \right) \\
&\stackrel{(a)}{\geq} \frac{2M\eta}{12M\eta + 9} \exp(-\text{KL}(\mathbb{H}^1, \mathbb{H}^2)),
\end{aligned} \tag{19}$$

where (a) follows from the Bretagnolle-Huber inequality [16]. Then we can upper bound the KL divergence $\text{KL}(\mathbb{H}^1, \mathbb{H}^2)$ between \mathbb{H}^1 and \mathbb{H}^2 as follows,

$$\begin{aligned}
\text{KL}(\mathbb{H}^1, \mathbb{H}^2) &= \mathbb{E}_{h \sim \mathbb{H}^1} \left[\log \frac{\mathbb{H}^1(h)}{\mathbb{H}^2(h)} \right] \leq \mathbb{E}_{h \sim \mathbb{H}^1} \left[\sum_{t=1}^T \sum_{a \in A_t} \log \frac{\nu_i^1(r_{i,t})}{\nu_i^2(r_{i,t})} \right] = \sum_{t=1}^T \mathbb{E}_{\mathbf{p}_t \sim \mathcal{A}^1} \mathbb{E}_{i \sim \mathbf{p}_t} [\text{KL}(\nu_i^1, \nu_i^2)] \\
&= \sum_{t=1}^T \mathbb{E}_{\mathbf{p}_t \sim \mathcal{A}^1} [p_{1,t} \text{KL}(\nu_1^1, \nu_1^2)] \\
&= \sum_{t=1}^T \mathbb{E}_{\mathbf{p}_t \sim \mathcal{A}^1} \left[p_{1,t} \left(2\eta \log 2 + (1 - 2\eta) \log \frac{1 - 2\eta}{1 - \eta} \right) \right] \\
&\leq 2T\eta \log 2,
\end{aligned} \tag{20}$$

where $\mathbf{p}_t \sim \mathcal{A}^1$ means that \mathbf{p}_t is sampled from the process of the algorithm \mathcal{A}^1 applied to the first MAB instance.

According to (19) and (20), set $\eta = \frac{1}{2T}$, $M = T$, we have

$$\mathbb{E} \left[\frac{1}{T} \text{FR}_T^1 \right] + \mathbb{E} \left[\frac{1}{T} \text{FR}_T^2 \right] \geq \frac{2M\eta}{12M\eta + 9} \exp(-2T\eta \log 2) \geq 0.03, \tag{21}$$

which implies that at least one of the two CMVB instances incurs linear expected fairness regret.

This completes the proof of Theorem 2. \square

C. Proof of Theorem 3

Proof. We first prove the expected mean-variance upper bound and then prove the expected fairness regret upper bound of MVFQ.

Part 1: Proof of the Expected Mean-variance Regret Upper Bound of MVFQ

We first prove the following lemmas that will be used in the proofs.

Lemma 1. For any $\delta \in (0, 1)$, with probability at least $1 - \frac{\delta}{2}$, $\forall t > \lceil K/L \rceil$, $i \in [K]$, $r_{i,t} \in [0, 1]$, the mean-variance vector $\mathbf{w} \in \mathcal{C}_t$.

Proof. According to Hoeffding's inequality, for $t \in [T]$, $i \in [K]$, with probability at least $1 - \frac{\delta}{4KT}$,

$$|\hat{\mu}_{i,t} - \mu_i| \leq \sqrt{\frac{\log(8KT/\delta)}{2n_{i,t}}}. \tag{22}$$

Let $\hat{\mu}_{i,t}^{(2)} = \frac{\sum_{\tau=1}^t \mathbb{1}_{\{i \in S_\tau\}} r_{i,\tau}^2}{n_{i,t}}$ and $\mu_i^{(2)} = \mathbb{E}[r_{i,t}^2]$. Similarly, with probability at least $1 - \frac{\delta}{4KT}$,

$$\left| \hat{\mu}_{i,t}^{(2)} - \mu_i^{(2)} \right| \leq \sqrt{\frac{\log(8KT/\delta)}{2n_{i,t}}}. \tag{23}$$

Note that $\sigma_i^2 = \mu_i^{(2)} - \mu_i^2$ and $\hat{\sigma}_{i,t}^2 = \sum_{\tau=1}^t \frac{\mathbb{1}_{\{i \in S_\tau\}} (r_{i,\tau} - \hat{\mu}_{i,t})^2}{n_{i,t}} = \sum_{\tau=1}^t \frac{\mathbb{1}_{\{i \in S_\tau\}} r_{i,\tau}^2}{n_{i,t}} - \hat{\mu}_{i,t}^2 = \hat{\mu}_{i,t}^{(2)} - \hat{\mu}_{i,t}^2$. Then we have

$$\begin{aligned}
|\hat{\sigma}_{i,t}^2 - \sigma_i^2| &\leq \left| \hat{\mu}_{i,t}^{(2)} - \mu_i^{(2)} \right| + \left| \hat{\mu}_{i,t}^2 - \mu_i^2 \right| \\
&\leq \left| \hat{\mu}_{i,t}^{(2)} - \mu_i^{(2)} \right| + (\hat{\mu}_{i,t} + \mu_i) |\hat{\mu}_{i,t} - \mu_i| \\
&\leq \left| \hat{\mu}_{i,t}^{(2)} - \mu_i^{(2)} \right| + 2 |\hat{\mu}_{i,t} - \mu_i|
\end{aligned} \tag{24}$$

Fianlly, with probability at least $1 - \frac{\delta}{2KT}$,

$$\begin{aligned}
|\hat{w}_{i,t} - w_i| &= |(\rho \hat{\mu}_{i,t} - \hat{\sigma}_{i,t}^2) - (\rho \mu_i - \sigma_i^2)| \leq \rho |\hat{\mu}_{i,t} - \mu_i| + |\hat{\sigma}_{i,t}^2 - \sigma_i^2| \\
&\leq (\rho + 2) |\hat{\mu}_{i,t} - \mu_i| + \left| \hat{\mu}_{i,t}^{(2)} - \mu_i^{(2)} \right| \\
&\leq (\rho + 3) \sqrt{\frac{\log(8KT/\delta)}{2n_{i,t}}}.
\end{aligned} \tag{25}$$

Using union bound over all round t and i , with probability at least $1 - \frac{\delta}{2}$, $\forall t > \lceil K/L \rceil$, $a \in [K]$, $w_i \in [l_{i,t}, u_{i,t}]$.

This completes the proof of Lemma 1. \square

Lemma 2. For any $\delta \in (0, 1)$, with probability at least $1 - \frac{\delta}{2}$, it holds that

$$\left| \sum_{t=\lceil K/L \rceil+1}^T \mathbb{E}_{i \sim \mathbf{p}_t} \left[\sqrt{\frac{1}{n_{i,t}}} \right] - \sum_{t=\lceil K/L \rceil+1}^T \sum_{i \in S_t} \sqrt{\frac{1}{n_{i,t}}} \right| \leq L \sqrt{2T \log \frac{4}{\delta}}. \quad (26)$$

Proof. We first construct a martingale difference sequence

$$\sum_{i \in S_t} \sqrt{\frac{1}{n_{i,t}}} - \mathbb{E}_{I \sim \mathbf{p}_t} \left[\sqrt{\frac{1}{n_{i,t}}} \right]. \quad (27)$$

Then we have

$$\left| \sum_{i \in S_t} \sqrt{\frac{1}{n_{i,t}}} - \mathbb{E}_{i \sim \mathbf{p}_t} \left[\sqrt{\frac{1}{n_{i,t}}} \right] \right| < L. \quad (28)$$

By Azuma-Hoeffding's inequality, with probability at least $1 - \frac{\delta}{2}$, we have

$$\left| \sum_{t=\lceil K/L \rceil+1}^T \mathbb{E}_{i \sim \mathbf{p}_t} \left[\sqrt{\frac{1}{n_{i,t}}} \right] - \sum_{t=\lceil K/L \rceil+1}^T \sum_{i \in S_t} \sqrt{\frac{1}{n_{i,t}}} \right| \leq L \sqrt{2T \log \frac{4}{\delta}}. \quad (29)$$

This completes the proof of Lemma 2. \square

Based on the lemmas above, with probability at least $1 - \delta$, the mean-variance regret can be bounded as follows:

$$\begin{aligned} \text{MR}_T &= \sum_{t=1}^T \max \left\{ \sum_{i=1}^K p_i^* w_i - \sum_{i=1}^K p_{i,t} w_i, 0 \right\} \\ &\stackrel{(a)}{\leq} \left(\frac{K}{L} + 1 \right) L + \sum_{t=\lceil \frac{K}{L} \rceil+1}^T \sum_{i=1}^K (p_{i,t} \tilde{w}_{i,t} - p_{i,t} w_i) \\ &= K + L + \sum_{t=\lceil \frac{K}{L} \rceil+1}^T \sum_{i=1}^K p_{i,t} (\tilde{w}_{i,t} - \hat{w}_{i,t} + \hat{w}_{i,t} - w_i) \\ &\stackrel{(b)}{\leq} K + L + \sum_{t=\lceil \frac{K}{L} \rceil+1}^T \sum_{i=1}^K p_{i,t} 2(\rho + 3) \sqrt{\frac{\log(8KT/\delta)}{2n_{i,t}}} \\ &= K + L + (\rho + 3) \sqrt{2 \log \frac{8KT}{\delta}} \sum_{t=\lceil \frac{K}{L} \rceil+1}^T \mathbb{E}_{i \sim \mathbf{p}_t} \left[\sqrt{\frac{1}{n_{i,t}}} \right] \\ &\stackrel{(c)}{\leq} K + L + (\rho + 3) \sqrt{2 \log \frac{8KT}{\delta}} \left(L \sqrt{2T \log \frac{4}{\delta}} + \sum_{t=\lceil K/L \rceil+1}^T \sum_{i \in S_t} \sqrt{\frac{1}{n_{i,t}}} \right) \\ &\leq K + L + (\rho + 3) \sqrt{2 \log \frac{8KT}{\delta}} \left(L \sqrt{2T \log \frac{4}{\delta}} + K \int_0^{\frac{LT}{K}} \sqrt{\frac{1}{x}} dx \right) \\ &\leq K + L + (\rho + 3) \sqrt{2 \log \frac{8KT}{\delta}} \left(L \sqrt{2T \log \frac{4}{\delta}} + 2\sqrt{LKT} \right), \end{aligned} \quad (30)$$

where (a) is from Line 12 in Algorithm 1, (b) is from Lemma 1 and (c) is from Lemma 2. Setting $\delta = \frac{1}{LT}$, the expected mean-variance regret can be upper bounded as

$$\begin{aligned} \mathbb{E} [\text{MR}_T] &\leq K + L + (\rho + 3) \sqrt{2 \log \frac{8KT}{\delta}} \left(L \sqrt{2T \log \frac{4}{\delta}} + 2\sqrt{LKT} \right) + LT\delta \\ &\leq K + L + (\rho + 3) \sqrt{4 \log(8LKT)} \left(L \sqrt{2T \log(4LT)} + 2\sqrt{LKT} \right) + 1 \\ &= O \left(\rho \sqrt{LKT \log T} \right). \end{aligned} \quad (31)$$

This completes the proof of the mean-variance regret upper bound of MVFQ.

Part 2: Proof of the Expected Fairness Regret Upper Bound of MVFQ

For any $\delta \in (0, 1)$, with probability $1 - \delta$,

$$\begin{aligned}
\sum_{i=1}^K |p_{i,t} - p_i^*| &= \sum_{i=1}^K \left| \frac{Lf(\tilde{w}_{i,t})}{\sum_{i'=1}^K f(\tilde{w}_{i',t})} - \frac{Lf(w_i)}{\sum_{i'=1}^K f(w_{i'})} \right| \\
&= \sum_{i=1}^K \frac{L \left| f(\tilde{w}_{i,t}) \sum_{i'=1}^K f(w_{i'}) - f(w_i) \sum_{i'=1}^K f(\tilde{w}_{i',t}) \right|}{\sum_{i'=1}^K f(\tilde{w}_{i',t}) \sum_{i'=1}^K f(w_{i'})} \\
&= \sum_{i=1}^K \frac{L \left| f(\tilde{w}_{i,t}) \sum_{i'=1}^K f(w_{i'}) - f(w_i) \sum_{i'=1}^K f(w_{i'}) + f(w_i) \sum_{i'=1}^K f(w_{i'}) - f(w_i) \sum_{i'=1}^K f(\tilde{w}_{i',t}) \right|}{\sum_{i'=1}^K f(\tilde{w}_{i',t}) \sum_{i'=1}^K f(w_{i'})} \\
&\leq \frac{L \sum_{i=1}^K |f(\tilde{w}_{i,t}) - f(w_i)| \sum_{i'=1}^K f(w_{i'}) + L \sum_{i=1}^K f(w_i) \sum_{i'=1}^K |f(w_{i'}) - f(\tilde{w}_{i',t})|}{\sum_{i'=1}^K f(\tilde{w}_{i',t}) \sum_{i'=1}^K f(w_{i'})} \\
&= \frac{2L \sum_{i=1}^K \frac{f(\tilde{w}_{i,t})}{f(\tilde{w}_{i,t})} |f(\tilde{w}_{i,t}) - f(w_i)|}{\sum_{i'=1}^K f(\tilde{w}_{i',t})} \\
&\stackrel{(a)}{\leq} \sum_{i=1}^K \frac{2Mp_{i,t}L}{\lambda} (\tilde{w}_{i,t} - w_i) \\
&\stackrel{(b)}{\leq} \sum_{i=1}^K \frac{4ML(\rho+3)p_{i,t}}{\lambda} \sqrt{\frac{\log(8KT/\delta)}{2n_{i,t}}} \\
&= \frac{4ML(\rho+3)}{\lambda} \sqrt{\frac{1}{2} \log \frac{8KT}{\delta}} \mathbb{E}_{i \sim \mathbf{p}_t} \left[\sqrt{\frac{1}{n_{i,t}}} \right],
\end{aligned} \tag{32}$$

where (a) follows from the Assumption 1(i) and Assumption 2, (b) follows from Lemma 1. When $T > K$, with probability at least $1 - \delta$, the fairness regret can be upper bounded as follows:

$$\begin{aligned}
\text{FR}_T &= \sum_{t=1}^T \sum_{i=1}^K |p_i^* - p_{i,t}| \leq \left(\frac{K}{L} + 1 \right) L + \sum_{t=\lceil \frac{K}{L} \rceil + 1}^T \sum_{i=1}^K |p_i^* - p_{i,t}| \\
&\leq K + L + \frac{4ML(\rho+3)}{\lambda} \sqrt{\frac{1}{2} \log \frac{8KT}{\delta}} \sum_{t=\lceil \frac{K}{L} \rceil + 1}^T \mathbb{E}_{i \sim \mathbf{p}_t} \left[\sqrt{\frac{1}{n_{i,t}}} \right] \\
&\stackrel{(a)}{\leq} K + L + \frac{4ML(\rho+3)}{\lambda} \sqrt{\frac{1}{2} \log \frac{8KT}{\delta}} \left(L \sqrt{2T \log \frac{4}{\delta}} + 2\sqrt{LKT} \right),
\end{aligned} \tag{33}$$

where (a) follows from the result in (30). Furthermore, setting $\delta = \frac{1}{LT}$, the expected fairness regret can be upper bounded as

$$\begin{aligned}
\mathbb{E}[\text{FR}_T] &\leq K + L + \frac{4ML(\rho+3)}{\lambda} \sqrt{\frac{1}{2} \log \frac{8KT}{\delta}} \left(L \sqrt{2T \log \frac{4}{\delta}} + 2\sqrt{LKT} \right) + LT\delta \\
&\leq K + L + \frac{4ML(\rho+3)}{\lambda} \sqrt{\log(8LKT)} \left(L \sqrt{2T \log(4LT)} + 2\sqrt{LKT} \right) + 1 \\
&= O \left(\frac{\rho ML}{\lambda} \sqrt{LKT \log T} \right).
\end{aligned} \tag{34}$$

This completes the proof of fairness regret upper bound of MVFQ.

Combining Part 1 and Part 2 of the proof, we complete the proof of Theorem 3. \square

D. Proof of Theorem 4

Proof. We first prove the expected mean-variance upper bound and then prove the expected fairness regret upper bound of MVFP.

Part 1: Proof of the Expected Mean-variance Regret Upper Bound of MVFP

Before presenting our theoretical analysis and results, we need two technical lemmas that will be used in the proofs. We first prove $\hat{b}_{i,t}$ is an unbiased estimate of the Borda score b_i in the following lemma.

Lemma 3. At any round t , for all $i \in [K]$, it holds that $\mathbb{E}[\tilde{b}_{i,t}] = b_i$.

Proof. Note that

$$\begin{aligned}
\mathbb{E}[\tilde{b}_{i,t}] &= \mathbb{E}\left[\frac{\mathbb{1}_{\{i \in S_t\}} \mathbb{1}_{\{j \neq i\}}}{p_{i,t}(K-1)} \sum_{j \in [K]} \frac{\mathbb{1}_{\{j \in S_t\}}}{p_{j,t}} \frac{v_{ij,t}}{v_{ij,t} + v_{ji,t}}\right] = \frac{1}{K-1} \mathbb{E}\left[\sum_{j \in [K]} \frac{\mathbb{1}_{\{i \in S_t\}} \mathbb{1}_{\{j \in S_t\}} \mathbb{1}_{\{j \neq i\}}}{p_{i,t} p_{j,t}} \frac{v_{ij,t}}{v_{ij,t} + v_{ji,t}}\right] \\
&= \frac{1}{K-1} \mathbb{E}_{\mathcal{H}_{t-1}} \left[\mathbb{E} \left[\frac{\mathbb{1}_{\{i \in S_t\}}}{p_{i,t}} \sum_{j \in [K]} \frac{\mathbb{1}_{\{j \in S_t\}} \mathbb{1}_{\{j \neq i\}}}{p_{j,t}} \frac{v_{ij,t}}{v_{ij,t} + v_{ji,t}} \middle| \mathcal{H}_{t-1} \right] \right] \\
&= \frac{1}{K-1} \mathbb{E}_{\mathcal{H}_{t-1}} \left[\mathbb{E}_i \left[\frac{\mathbb{1}_{\{i \in S_t\}}}{p_{i,t}} \sum_{j \in [K]} \mathbb{E}_j \left[\frac{\mathbb{1}_{\{j \in S_t\}} \mathbb{1}_{\{j \neq i\}}}{p_{j,t}} \mathbb{E} \left[\frac{v_{ij,t}}{v_{ij,t} + v_{ji,t}} \right] \right] \middle| \mathcal{H}_{t-1} \right] \right] \\
&= \frac{1}{K-1} \mathbb{E}_{\mathcal{H}_{t-1}} \left[\mathbb{E}_i \left[\frac{\mathbb{1}_{\{i \in S_t\}}}{p_{i,t}} \sum_{j \in [K]} \mathbb{E}_j \left[\frac{\mathbb{1}_{\{j \in S_t\}} \mathbb{1}_{\{j \neq i\}} q_{ij}}{p_{j,t}} \right] \middle| \mathcal{H}_{t-1} \right] \right] \\
&= \frac{1}{K-1} \mathbb{E}_{\mathcal{H}_{t-1}} \left[\mathbb{E}_i \left[\frac{\mathbb{1}_{\{i \in S_t\}}}{p_{i,t}} \sum_{j \in [K]} \sum_{j'=1}^K \frac{\mathbb{1}_{\{j=j'\}} \mathbb{1}_{\{j \in S_t\}} \mathbb{1}_{\{j \neq i\}} q_{ij} p_{j',t}}{p_{j,t}} \middle| \mathcal{H}_{t-1} \right] \right] \\
&= \frac{1}{K-1} \mathbb{E}_{\mathcal{H}_{t-1}} \left[\mathbb{E}_i \left[\frac{\mathbb{1}_{\{i \in S_t\}}}{p_{i,t}} \sum_{j \in [K]} q_{ij} \mathbb{1}_{\{j \neq i\}} \middle| \mathcal{H}_{t-1} \right] \right] \\
&= \frac{1}{K-1} \mathbb{E}_{\mathcal{H}_{t-1}} \left[\sum_{i'=1}^K \frac{\mathbb{1}_{\{i=i'\}} \mathbb{1}_{\{i \in S_t\}} p_{i',t}}{p_{i,t}} \sum_{j \in [K]} q_{ij} \mathbb{1}_{\{j \neq i\}} \right] \\
&= \frac{1}{K-1} \sum_{j \in [K]} q_{ij} \mathbb{1}_{\{j \neq i\}} = b_i,
\end{aligned} \tag{35}$$

which concludes the proof. \square

Next, we prove that the constructed confidence region \mathcal{C}_t^b contains the actual Borda score mean-variance with high probability in the following lemma.

Lemma 4. For any $\delta \in (0, 1)$, $\alpha \geq \left(\frac{K^2 - 2K + L}{L^2 - L}\right)^4$, with probability at least $1 - \delta$, the Borda mean-variance vector $\mathbf{w}^b \in \mathcal{C}_t^b$.

Proof. By Lemma 3, we have $\mathbb{E}[\tilde{b}_{i,t}] = b_i$. We can construct a martingale difference sequence $\tilde{b}_{i,\tau} - b_i$ for $\tau \leq t$. Note that $|\tilde{b}_{i,\tau} - b_i| \leq \left(\frac{K^2 - 2K + L}{L^2 - L}\right)^2$ since $p_{i,t} \geq \frac{L^2 - L}{K^2 - 2K + L}$ under Assumption 1. Then by Azuma-Hoeffding's inequality, with probability at least $1 - \delta/2$, we have

$$|\hat{b}_{i,t} - b_i| = \left| \frac{1}{t} \sum_{\tau=1}^t \tilde{b}_{i,\tau} - \frac{1}{t} \sum_{\tau=1}^t b_i \right| \leq \left(\frac{K^2 - 2K + L}{L^2 - L} \right)^2 \sqrt{\frac{2 \log 4/\delta}{t}} \leq \sqrt{\frac{2\alpha \log 4/\delta}{t}}. \tag{36}$$

Let $\hat{b}_{i,t}^{(2)} = \frac{\sum_{\tau=1}^t \tilde{b}_{i,\tau}^2}{t}$ and $b_i^{(2)} = \mathbb{E}[\tilde{b}_{i,t}^2]$. Similarly, by Azuma-Hoeffding's inequality, with probability at least $1 - \delta/2$,

$$|\hat{b}_{i,t}^{(2)} - b_i^{(2)}| = \left| \frac{1}{t} \sum_{\tau=1}^t \tilde{b}_{i,\tau}^{(2)} - \frac{1}{t} \sum_{\tau=1}^t b_i^{(2)} \right| \leq \left(\frac{K^2 - 2K + L}{L^2 - L} \right)^4 \sqrt{\frac{2 \log 4/\delta}{t}} \leq \alpha \sqrt{\frac{2 \log 4/\delta}{t}}. \tag{37}$$

Note that $\sigma_i^2 = b_i^{(2)} - b_i^2$ and $\hat{\sigma}_{i,t}^2 = \sum_{\tau=1}^t \frac{(\tilde{b}_{i,t} - \tilde{b}_{i,\tau})^2}{t} = \sum_{\tau=1}^t \frac{\tilde{b}_{i,\tau}^2}{t} - \hat{b}_{i,t}^2 = \hat{b}_{i,t}^{(2)} - \hat{b}_{i,t}^2$. Then we have

$$\begin{aligned}
|\hat{\sigma}_{i,t}^2 - \sigma_i^2| &\leq |\hat{b}_{i,t}^{(2)} - b_i^{(2)}| + |\hat{b}_{i,t}^2 - b_i^2| \\
&\leq |\hat{b}_{i,t}^{(2)} - b_i^{(2)}| + (\hat{b}_{i,t} + b_i) |\hat{b}_{i,t} - b_i| \\
&\leq |\hat{b}_{i,t}^{(2)} - b_i^{(2)}| + \left(\left(\frac{K^2 - 2K + L}{L^2 - L} \right)^2 + 1 \right) |\hat{b}_{i,t} - b_i| \\
&\leq |\hat{b}_{i,t}^{(2)} - b_i^{(2)}| + (\sqrt{\alpha} + 1) |\hat{b}_{i,t} - b_i|
\end{aligned} \tag{38}$$

Finally, with probability at least $1 - \delta$, $\alpha \geq \left(\frac{K^2 - 2K + L}{L^2 - L} \right)^4$,

$$\begin{aligned}
|\hat{w}_{i,t}^b - w_i^b| &= |(\rho \hat{b}_{i,t} - \hat{\sigma}_{i,t}^2) - (\rho b_i - \sigma_i^2)| \leq \rho |\hat{b}_{i,t} - b_i| + |\hat{\sigma}_{i,t}^2 - \sigma_i^2| \\
&\leq \rho |\hat{b}_{i,t} - b_i| + |\hat{b}_{i,t}^{(2)} - b_i^{(2)}| + (\sqrt{\alpha} + 1) |\hat{b}_{i,t} - b_i| \\
&\leq (\sqrt{\alpha}(\rho + 1) + 2\alpha) \sqrt{\frac{2 \log 4/\delta}{t}} \\
&\leq \alpha(\rho + 3) \sqrt{\frac{2 \log 4/\delta}{t}}.
\end{aligned} \tag{39}$$

This completes the proof of Lemma 4.

Based on the lemmas above, with probability at least $1 - \delta$, the mean-variance regret can be bounded as follows:

$$\begin{aligned}
\text{BMR}_T &= \sum_{t=1}^T \max \left\{ \sum_{i=1}^K p_i^* w_i^b - \sum_{i=1}^K p_{i,t} w_i^b, 0 \right\} \\
&\stackrel{(a)}{\leq} \sum_{t=1}^T \sum_{i=1}^K (p_{i,t} \tilde{w}_{i,t}^b - p_{i,t} w_i^b) \\
&= \sum_{t=1}^T \sum_{i=1}^K p_{i,t} (\tilde{w}_{i,t}^b - \hat{w}_{i,t}^b + \hat{w}_{i,t}^b - w_i^b) \\
&\stackrel{(b)}{\leq} \sum_{t=1}^T \sum_{i=1}^K p_{i,t} 2 (\sqrt{\alpha}(\rho + 1) + 2\alpha) \sqrt{\frac{2 \log 4/\delta}{t}} \\
&= 2 (\sqrt{\alpha}(\rho + 1) + 2\alpha) \sqrt{2 \log \frac{4}{\delta}} \sum_{t=1}^T \mathbb{E}_{i \sim p_t} \left[\sqrt{\frac{1}{t}} \right] \\
&\stackrel{(c)}{\leq} 4 (\sqrt{\alpha}(\rho + 1) + 2\alpha) L \sqrt{2T \log \frac{4}{\delta}}
\end{aligned} \tag{40}$$

Setting $\delta = \frac{1}{LT}$, the expected mean-variance regret can be upper bounded as

$$\begin{aligned}
\mathbb{E} [\text{BMR}_T] &\leq 4 (\sqrt{\alpha}(\rho + 1) + 2\alpha) L \sqrt{2T \log \frac{4}{\delta}} + LT\delta \leq 4 (\sqrt{\alpha}(\rho + 1) + 2\alpha) L \sqrt{2T \log(4LT)} + 1 \\
&= O \left(\alpha \rho L \sqrt{T \log T} \right).
\end{aligned} \tag{41}$$

This completes the proof of the mean-variance regret upper bound of MVFP.

□

Part 2: Proof of the Expected Fairness Regret Upper Bound of MVFP

For any $\delta \in (0, 1)$, with probability $1 - \delta$,

$$\begin{aligned}
\sum_{i=1}^K |p_{i,t} - p_i^*| &= \sum_{i=1}^K \left| \frac{Lf(\tilde{w}_{i,t}^b)}{\sum_{i'=1}^K f(\tilde{w}_{i',t}^b)} - \frac{Lf(w_i^b)}{\sum_{i'=1}^K f(w_{i'}^b)} \right| \\
&= \sum_{i=1}^K L \left| \frac{f(\tilde{w}_{i,t}^b) \sum_{i'=1}^K f(w_{i'}^b) - f(w_i^b) \sum_{i'=1}^K f(\tilde{w}_{i',t}^b)}{\sum_{i'=1}^K f(\tilde{w}_{i',t}^b) \sum_{i'=1}^K f(w_{i'}^b)} \right| \\
&= \sum_{i=1}^K L \left| \frac{f(\tilde{w}_{i,t}^b) \sum_{i'=1}^K f(w_{i'}^b) - f(w_i^b) \sum_{i'=1}^K f(w_{i'}^b) + f(w_i^b) \sum_{i'=1}^K f(w_{i'}^b) - f(w_i^b) \sum_{i'=1}^K f(\tilde{w}_{i',t}^b)}{\sum_{i'=1}^K f(\tilde{w}_{i',t}^b) \sum_{i'=1}^K f(w_{i'}^b)} \right| \\
&\leq \frac{L \sum_{i=1}^K |f(\tilde{w}_{i,t}^b) - f(w_i^b)| \sum_{i'=1}^K f(w_{i'}^b) + L \sum_{i=1}^K f(w_i^b) \sum_{i'=1}^K |f(w_{i'}^b) - f(\tilde{w}_{i',t}^b)|}{\sum_{i'=1}^K f(\tilde{w}_{i',t}^b) \sum_{i'=1}^K f(w_{i'}^b)} \\
&= \frac{2L \sum_{i=1}^K \frac{f(\tilde{w}_{i,t}^b)}{f(\tilde{w}_{i,t}^b)} |f(\tilde{w}_{i,t}^b) - f(w_i^b)|}{\sum_{i'=1}^K f(\tilde{w}_{i',t}^b)} \\
&\stackrel{(a)}{\leq} \sum_{i=1}^K \frac{2Mp_{i,t}L}{\lambda} (\tilde{w}_{i,t}^b - w_i^b) \\
&\stackrel{(b)}{\leq} \sum_{i=1}^K \frac{4ML\alpha(\rho+3)p_{i,t}}{\lambda} \sqrt{\frac{2\log 4/\delta}{t}} \\
&= \frac{4ML\alpha(\rho+3)}{\lambda} \sqrt{2\log \frac{4}{\delta}} \mathbb{E}_{i \sim p_t} \left[\sqrt{\frac{1}{t}} \right],
\end{aligned} \tag{42}$$

where (a) follows from the Assumption 1(i) and Assumption 2, (b) follows from Lemma 1. When $T > K$, with probability at least $1 - \delta$, the fairness regret can be upper bounded as follows:

$$\begin{aligned}
\text{BFR}_T &= \sum_{t=1}^T \sum_{i=1}^K |p_i^* - p_{i,t}| \leq \sum_{t=\lceil \frac{K}{L} \rceil + 1}^T \sum_{i=1}^K |p_i^* - p_{i,t}| \\
&\leq \frac{4ML(\sqrt{\alpha}(\rho+1) + 2\alpha)}{\lambda} \sqrt{2\log \frac{4}{\delta}} \sum_{t=1}^T \mathbb{E}_{i \sim p_t} \left[\sqrt{\frac{1}{t}} \right] \\
&\stackrel{(a)}{\leq} \frac{8ML^2(\sqrt{\alpha}(\rho+1) + 2\alpha)}{\lambda} \sqrt{2T \log \frac{4}{\delta}},
\end{aligned} \tag{43}$$

where (a) follows from the result in (30). Furthermore, setting $\delta = \frac{1}{LT}$, the expected fairness regret can be upper bounded as

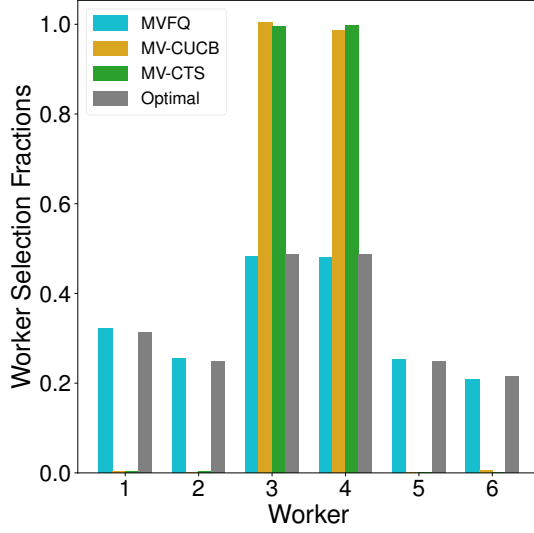
$$\begin{aligned}
\mathbb{E}[\text{BFR}_T] &\leq \frac{8ML^2(\sqrt{\alpha}(\rho+1) + 2\alpha)}{\lambda} \sqrt{2T \log \frac{4}{\delta}} + LT\delta \\
&\leq \frac{8ML^2(\sqrt{\alpha}(\rho+1) + 2\alpha)}{\lambda} \sqrt{2T \log(4LT)} + 1 \\
&= O\left(\frac{\alpha\rho ML^2}{\lambda} \sqrt{T \log T}\right).
\end{aligned} \tag{44}$$

This completes the proof of fairness regret upper bound of MVFP.

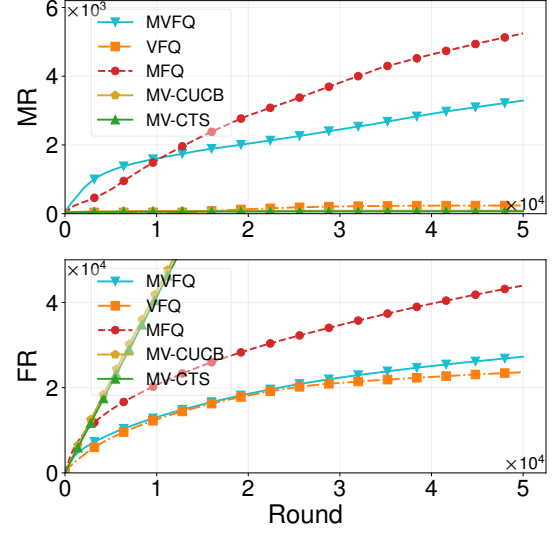
Combining Part 1 and Part 2 of the proof, we complete the proof of Theorem 4. \square

E. Addition experiments

We consider a mobile crowdsourcing system with $K = 50$ workers, where the platform selects $L = 10$ workers at each round. We then repeat the experiments for MVFQ and MVFP as described in Section V. As shown in Figure 4 and Figure 5, both MVFQ and MVFP maintain strong performance compared with other baselines, demonstrating the scalability and robustness of the proposed algorithms.

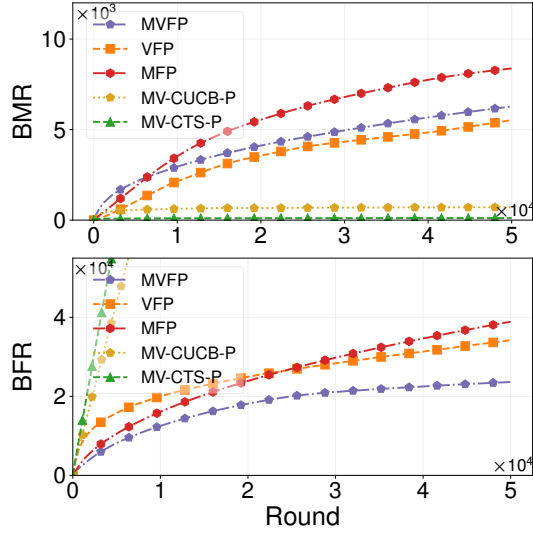


(a) Worker selection fractions



(b) Comparisons with other baselines

Fig. 4. Experiment results of MVFQ algorithm.



(a) Comparisons with other baselines

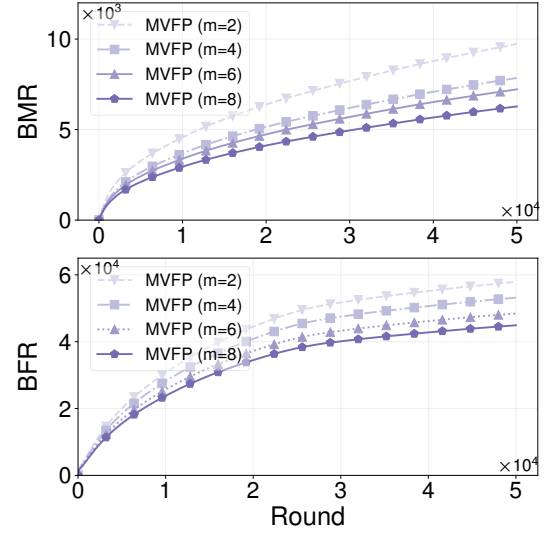
(b) BMR and BFR with different m

Fig. 5. Experiment results of MVFQ algorithm.