

---

# Group-Fair Online Allocation in Continuous Time

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 The theory of discrete-time online learning has been successfully applied in many  
2 problems that involve sequential decision-making under uncertainty. However,  
3 in many applications including contractual hiring in online freelancing platforms  
4 and server allocation in cloud computing systems, the outcome of each action  
5 is observed only after a random and action-dependent time. Furthermore, as a  
6 consequence of certain ethical and economic concerns, the controller may impose  
7 deadlines on the completion of each task, and require fairness across different  
8 groups in the allocation of total time budget  $B$ . In order to address these ap-  
9 plications, we consider continuous-time online learning problem with fairness  
10 considerations, and present a novel framework based on continuous-time utility  
11 maximization. We show that this formulation recovers reward-maximizing, max-  
12 imin fair and proportionally fair allocation rules across different groups as special  
13 cases. We characterize the optimal offline policy, which allocates the total time be-  
14 tween different actions in an optimally fair way (as defined by the utility function),  
15 and impose deadlines to maximize time-efficiency. In the absence of any statistical  
16 knowledge, we propose a novel online learning algorithm based on dual ascent  
17 optimization for time averages, and prove that it achieves  $\tilde{O}(B^{-1/2})$  regret bound.

## 18 1 Introduction

19 With the prevalence automated decision methods and machine learning methods, it is important to  
20 analyze the impact of learning and evaluate models not only with respect to traditional objectives  
21 such as reward or model accuracy, but also to account for the impact on individuals that interact with  
22 the system. Indeed, there are many studies highlighting algorithmic discrimination due to problems  
23 in the machine learning pipeline: imbalance in data [1], learnt representations [2, 3], choice of model  
24 proxies [4], demographic group-dependent difference in error rates of the learned models [5, 6, 7],  
25 to name a few. With rising ethical and legal concerns, addressing such issues has become urgent,  
26 specially as these impact critical societal decisions involving job opportunities and hiring. In 2014, it  
27 was estimated that 25% of the total workforce in the US was involved in some form of freelancing,  
28 and this number was predicted to grow to 40% by 2020 [8]. In reality, this percentage might be  
29 much higher, due to COVID-19 restrictions leading to increased work-from-home and changes in  
30 job opportunities [9, 10]. In online platforms however, there has been a strong evidence of bias  
31 observed in number of user reviews and user ratings<sup>1</sup> on completing jobs with significant correlations  
32 with race, gender, location of work and length of profiles<sup>2</sup> [11]. Motivated by these problems in

---

<sup>1</sup>The mean (median) normalized rating score for White workers was 0.98 (1), while it is 0.97 (1) for Black workers on TASKRABBIT. The mean (median) rating of White workers was found to be 3.3 (4.8), 3.0 (4.6) for Black workers, 3.3 (4.8) for Asian workers, 3.6 (4.8) for workers with a picture that does not depict a person, and 1.7 (0.0) for workers with no image on FIVERR [11].

<sup>2</sup>Mean (median) number of reviews: for women 33 (11), 59 (15) for men on TASKRABBIT. Mean (median) number of reviews: for Black workers was found to be 65 (4), 104 (6) for White workers, 101 (8) for Asian workers, 94 (10) for non-human pictures and 18 (0) for users with no image on FIVERR [11].

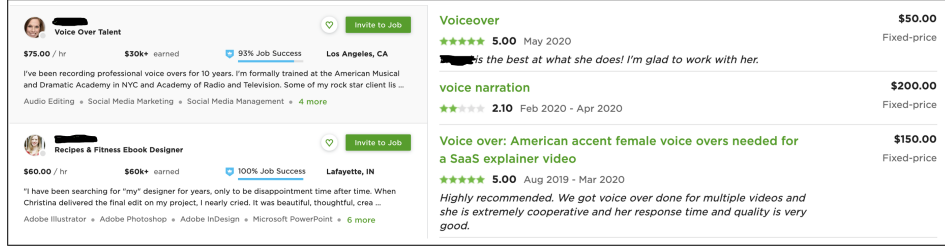


Figure 1: Freelancer profiles on UPWORK with their past performance and corresponding reviews for “fixed-price” contracts. Contractors can access these profiles and allocate fixed-timed contracts with deadlines.

33 online contractual hiring, we study a theoretical framework for sequential resource allocation to  
 34 workers, where the controller (decision maker) can enforce deadlines for each task’s completion. Our  
 35 key contribution is to quantify impact of reward maximization in terms of equality of opportunity  
 36 for jobs and develop algorithms that can achieve a meaningful trade-off between these via online  
 37 utility maximization. The challenge is to maximize total reward within a given time budget, while  
 38 accounting for random completion times by workers from different groups and fairness in allocation.

39 Formally, we consider  $K$  groups of individuals who can be hired sequentially for each task, i.e.,  
 40 at any point, exactly one individual can be hired. If an individual from group  $k \in [K]$  is chosen  
 41 for the  $n$ -th task and given a contractual deadline  $t$  by the controller, he/she generates a random  
 42 reward of  $R_{k,n}$  if the task is completed by (random) time  $X_{k,n}$  within deadline  $t$ . If the task is  
 43 not completed by the deadline, the reward obtained by the controller is zero and the time until the  
 44 deadline is wasted (i.e., yields 0 reward for the controller). Completion times and reward distributions  
 45 are assumed group-dependent and i.i.d. across tasks. The objective of the controller is to maximize  
 46 utility (trade-off between total reward and fair allocation) in the offline (known distributions) and  
 47 online settings (unknown distributions) under a budget constraint on time. As we will show in this  
 48 paper, controlled deadlines set are essential for optimal time-efficiency under the budget constraint.

49 The ethical problems we are concerned with involve the rate of jobs allocated to different demographic  
 50 groups and the deadlines imposed on these under reward maximization regimes [11]. Our sequential  
 51 framework would also apply to other settings, for e.g., comparative clinical trials with varying  
 52 follow-up durations as well as to server allocation in cloud computing where jobs are drawn from  
 53 different application groups and must commit computational resources until a specific amount of time  
 54 due to service level agreements (Section 2). We will often focus on the first application involving  
 55 online contractual hiring, since fairness concerns are most naturally motivated in this domain.

56 Given a time budget constraint  $B$  and the diverse random nature of completion time and reward pairs,  
 57 the main question we consider is how to decide distribution of tasks and deadlines between different  
 58 groups of people. Two potential extreme allocations are: (i) *Reward-maximizing task allocation*: The  
 59 controller assigns all tasks to the most rewarding group to maximize the total reward within the given  
 60 time budget. The other groups do not get any chance to receive tasks. (ii) *Proportional task allocation*:  
 61 The controller completely ignores the reward distributions, and attempts to give equal time share  
 62 to each group. In other words, each group receives a fraction of the tasks inversely proportional  
 63 to their mean completion times. There is clearly a trade-off between the reward maximization and  
 64 equal time-share considerations in continuous-time sequential task allocation, and well-chosen utility  
 65 functions [12] can be helpful in modeling this in a unified way. In this paper, we consider a very  
 66 general class of utility functions, which recovers broadly used fairness criteria such as proportional  
 67 fairness, max-min fairness, reward maximization among many others [13, 14, 12]. The controller  
 68 can determine her priorities in terms of notions of fairness and model the task allocation problem by  
 69 choosing the utility function accordingly.

70 The main contributions of this paper are summarized as follows:

- 71 **1. Incorporation of random completion time dynamics and fairness in allocation:** In discrete-  
 72 time online learning models, each action is assumed to take a unit completion time, thus the  
 73 random and diverse nature of task completion times, as required in many fundamental real-life  
 74 applications, is ignored. In this work, we incorporate this aspect and develop a sequential learning  
 75 framework in continuous time using tools from the theory of renewal processes and stochastic  
 76 control. We show how controlled deadlines improve the time-efficiency in continuous-time

77 decision processes. Moreover, this is the first work, to the best of our knowledge, that analyzes  
78 fair distribution policies in online contractual hiring.

79 **2. Characterization of Approximately Optimal Offline Policies:** As a consequence of the random  
80 and controlled task completion times, the optimal policy for fair resource allocation is PSPACE-  
81 hard akin to unbounded stochastic knapsack problems. For tractability in design and analysis, we  
82 propose an approximation to the optimal offline policy based on Lagrange duality and renewal  
83 theory, and prove that it is asymptotically optimal. These approximate policies allocate tasks  
84 independently with respect to a fixed probability distribution.

85 **3. Online learning for utility maximization:** For utility maximization in an online setting with full  
86 information feedback, we develop a novel and low-computational-complexity online learning  
87 algorithm based on dynamic stochastic optimization methods for time averages, and show that  
88 it achieves  $\tilde{O}(B^{-1/2})$  regret for a time budget  $B$ . The optimal offline control policy in this  
89 paper is time-dependent, randomized and attempts to optimize time averages unlike the reward  
90 maximization problems in discrete-time problems. Despite these, the online learning algorithm  
91 we developed adapts to the randomness in completion time-reward pairs, and achieves optimal  
92 performance with vanishing regret at a fast rate.

93 **Related Work:** The problem of fair resource allocation via utility maximization has been widely  
94 considered in economics and network management [15, 16, 17, 18]. The utility maximization  
95 approach to fair resource allocation in these papers predominantly deals with discrete-time systems,  
96 therefore the randomness and diversity in task completion times is completely ignored. Furthermore,  
97 these works either assume perfect knowledge of rewards and completion times prior to decision-  
98 making, or they assume the knowledge of statistics, therefore they do not incorporate online learning.  
99 The only continuous-time utility maximization approach to fair resource allocation is [19], which  
100 assumes the knowledge of first-order statistics.

101 Online learning under budget constraints has been considered under the scope of bandits with  
102 knapsacks [20, 21, 22]. In the classical bandits with knapsacks model, the objective is to maximize  
103 expected total reward under knapsack constraints in a stochastic setting. In [23], an interrupt  
104 mechanism is employed to incorporate the continuous-time dynamics into the budget-constrained  
105 online learning model. Note that these works focus solely on reward maximization, therefore  
106 they do not address the fair resource allocation problem. The bandits with knapsacks setting was  
107 extended to concave rewards and convex constraints in [24], which assumes bounded cost and reward,  
108 and the deadline mechanism is not involved in decision-making, thus optimal time-efficiency in  
109 continuous time is not achieved. Our paper deviates from this line of work as it proposes a versatile  
110 and comprehensive framework for fairness, and incorporates continuous-time dynamics into the  
111 decision-making for time-efficiency. We include an extended discussion of related work in Appendix  
112 A.

## 113 **2 Online Learning Framework for Group Fairness**

114 We consider the sequential and fair allocation of tasks to individuals from different groups, whose  
115 completion times and rewards randomly vary. This goal differs significantly from traditional online  
116 learning models that aim to maximize the expected total reward with unit completion times. Under  
117 this traditional setting, the controller’s goal is to find and persistently select the reward-maximizing  
118 groups to allocate its tasks. As a consequence, the reward-maximization objective leads to the  
119 starvation of suboptimal groups, which causes unfairness amongst the groups with different statistical  
120 characteristics. Next, we provide a few motivating examples with group fairness requirements:

121 • **Contractual Hiring in Online Freelancing Platforms:** Online freelancing sites like UPWORK  
122 host contractual workers (freelancers) that can be hired by “contractors” who require specific tasks  
123 to be completed. Each freelancer has a profile and performance on past tasks that can be learned by  
124 the contractors via ratings and reviews (see, typical profile in Figure 1). Fixed-timed contracts are  
125 popular on UPWORK, wherein contractors enforce a deadline by which the task must be completed  
126 otherwise the contract is terminated (i.e., there is no payment). Contractors can browse profiles  
127 and post a job to a selected set of freelancers with a deadline. However, there is a large literature  
128 documenting bias in online rating systems, which in turn impact job opportunities disparately  
129 [11, 25, 26], thus making it critical to develop theory of online learning for such settings.

130 • **Server Allocation in Cloud Computing:** An important application of our framework is online  
 131 learning for fair resource allocation in cloud computing systems. In a very basic setting, a single  
 132 server is sequentially allocated to tasks from one of  $K$  user groups, which exhibit similar execution  
 133 time statistics and priority levels within each group. In many practical scenarios, the execution time  
 134 of a task is unknown at the time decision [27, 28], and exhibits a power-law behavior [29], which  
 135 necessitates a deadline mechanism for optimal time-efficiency [23]. In this setting, the objective of  
 136 the controller is to allocate the server in an optimally fair way across the groups in a given time  
 137 interval  $[0, B]$ , depending on the completion time statistics and priority levels. Our work proposes  
 138 a versatile framework to model fairness for this problem based on the concept of continuous-time  
 139 utility maximization, and develops online learning algorithms to achieve the optimal performance  
 140 with low regret in the absence of any statistical knowledge.

141 More examples can be found in other domains, including multi-user wireless communication over  
 142 fading channels (e.g., see [23]), comparative clinical trials with optimal follow-up duration (e.g., see  
 143 [30, 31]), whereby the goal is to fairly share the limited resources between groups of users.

144 Motivated by these examples, next we introduce an online learning framework that expands the  
 145 traditional setting substantially to incorporate group fairness characteristics into its formulation.  
 146 Suppose that there are  $K \geq 1$  groups of individuals that are available for serving tasks with different  
 147 (and unknown) statistics. Specifically, if an individual from group  $k \in [K] = \{1, 2, \dots, K\}$  is  
 148 chosen for the  $n$ -th task, he/she takes  $X_{k,n}$  units of *completion time* for successful completion, and  
 149 a *reward* of  $R_{k,n}(t) = \bar{R}_{k,n} \mathbb{I}\{X_{k,n} \leq t\}$  is obtained  $t$  time units after the initiation where  $\bar{R}_{k,n}$   
 150 is a positive random variable and  $R_{k,n}(t) \in [0, R_{max}(t)]$  a.s. for some finite constant  $R_{max}(t) > 0$ .  
 151 Thus, the random reward  $\bar{R}_{k,n}$  is gathered only if the task is completed. For example, in the  
 152 server allocation application, a group- $k$  task of random size  $\bar{R}_{k,n}$  yields a reward (throughput)  
 153  $R_{k,n}(t) = \bar{R}_{k,n} \mathbb{I}\{X_{k,n} \leq t\}$  only upon successful completion. We assume that  $(X_{k,n}, R_{k,n}(t))$   
 154 is independent and identically distributed (iid) over  $n$ , and independent across different groups  $k$ .  
 155 Note that the completion time  $X_{k,n}$  and reward  $\bar{R}_{k,n}$  can be correlated, for example, in the server  
 156 allocation example, the completion time  $X_{k,n}$  and size  $\bar{R}_{k,n}$  of a task are positively correlated [32].  
 157 We assume that each task has a positive completion time, i.e.,  $X_{k,n} > 0$  almost surely for all  $k, n$ .

158 Before the  $n$ -th task begins, the controller makes two decisions: the group  $G_n \in [K]$  of the individual  
 159 that will be assigned the task, and a deadline  $T_n \in \mathbb{T}$ , where  $\mathbb{T} \subset \mathbb{R}_+$  is the decision set. If the task is  
 160 not completed by the selected deadline, the service is interrupted without collecting any reward. In  
 161 many applications, the deadlines are chosen within a discrete set (e.g., days/months in contractual  
 162 hiring or time-slots in server allocation), thus we assume a finite decision set  $\mathbb{T} = \{t_1, t_2, \dots, t_L\}$   
 163 with  $t_l < \infty$  for all  $l$  in this paper. The sequential task allocation continues until a given time budget  
 164  $B > 0$  is exceeded, therefore, the completion time of a task is as important as the reward.

165 To describe this process mathematically, let  $\mathcal{H}_{k,n-1}$  denote the available feedback for group  $k$ , and  
 166  $\mathcal{H}_{n-1} = \cup_{k \in [K]} \mathcal{H}_{k,n-1}$  denote the history before making a decision for task  $n$ . For a given time  
 167 budget  $B > 0$ , a *causal policy*  $\pi = \{\pi_1, \pi_2, \dots\}$  sequentially makes two decisions  $\pi_n = (G_n, T_n) \in$   
 168  $[K] \times \mathbb{T}$  for each task  $n$  based on the history  $\mathcal{H}_{k-1}$ , where  $G_n$  is the chosen group and  $T_n$  is the  
 169 assigned deadline. Under a policy  $\pi$ , the number of initiated tasks is the following *first-passage time*:

$$N^\pi(B) = \inf \left\{ n : \sum_{i=1}^n \min\{X_{G_i, i}, T_i\} > B \right\}, \quad (1)$$

170 which is a random and controlled stopping time. Moreover, the *reward rate* of any user type  $k$  is:

$$\bar{r}_k^\pi(B) = \mathbb{E} \left[ \frac{1}{B} \sum_{n=1}^{N^\pi(B)} \mathbb{I}\{G_n = k\} R_{k,n}(T_n) \right], \quad \text{under policy } \pi. \quad (2)$$

171 If  $R_{k,n}(t) = \mathbb{I}\{X_{k,n} \leq t\}$ , i.e., each task completion yields a unit reward, then  $\bar{r}_k^\pi(B)$  simply denotes  
 172 the task completion rate (i.e., throughput) of group  $k$  individuals in the time interval  $[0, B]$ .

Note that designing strategies that aim to maximize the total reward rate in (2) will lead to the  
 persistent selection of the group with the highest reward rate at the cost of starvation of all the rest  
 (see [23]). In order to address group fairness considerations, we propose a continuous-time online  
 learning framework based on the utility maximization concept that is used effectively in the fair

resource allocation domain (e.g., see [16]). Specifically, for a given continuously-differentiable, concave and monotonically increasing utility function  $U_k : \mathbb{R} \rightarrow \mathbb{R}$ , we let the utility of group  $k$  under a policy  $\pi$  be given by  $U_k(\bar{r}_k^\pi(B))$ . Then, the total utility under a policy  $\pi$  is defined as:

$$U^\pi(B) = \sum_{k=1}^K U_k(\bar{r}_k^\pi(B)), \text{ for time interval } [0, B].$$

173 The optimum utility over a class of policies  $\Pi$ , and the regret for a given  $\pi \in \Pi$  are, respectively:

$$\text{OPT}_\Pi(B) = \max_{\pi \in \Pi} \sum_{k=1}^K U_k(\bar{r}_k^\pi(B)) \quad \text{and} \quad \text{REG}_\Pi^\pi(B) = \text{OPT}_\Pi(B) - U^\pi(B), \text{ for } B > 0. \quad (3)$$

174 Note that, due to the monotonically increasing and concave nature of utility functions, allocating the  
 175 tasks always to the most rewarding group is not a good choice, because the same amount of time  
 176 could yield a higher utility for another group because of the diminishing return property of concave  
 177 functions. A particularly important set of utility functions is captured by the  $\alpha$ -fair class, given next.

178 **Definition 1** ( $\alpha$ -Fair Allocation). *For any given  $\alpha > 0$  and weight  $w_k > 0$ , let  $U_k(x) = w_k \frac{x^{1-\alpha}}{1-\alpha}$ ,  
 179 for all  $k$ . Resource allocation by using these utility functions is called  $\alpha$ -fair resource allocation.*

180 This class is attractive since it includes as special cases proportional fairness, minimum potential  
 181 delay fairness, reward maximization and max-min fairness [12].

### 182 3 Approximation of the Optimal Offline Policy

183 Note that a simpler version of the sequential maximization problem in (3) with linear utility functions  
 184 over all causal policies is called an unbounded knapsack problem, and it is PSPACE-hard even in the  
 185 case of known statistics [33, 20]. Therefore, the optimal causal policy for the problem in (3) has a  
 186 very high computational complexity even in the offline setting, which makes it intractable for online  
 187 learning. For tractability in design and analysis, we consider a class of simple policies that allocate  
 188 tasks in an i.i.d. randomized way according to a fixed probability distribution over groups, and show  
 189 its efficiency in this section.

190 **Definition 2** (Stationary Randomized Policies). *Let  $P$  be a fixed probability distribution over  $[K] \times \mathbb{T}$ .  
 191 A stationary randomized policy (SRP)  $\pi = \pi(P)$  makes a randomized decision independently  
 192 according to  $P$  for every task until the time budget  $B$  is depleted. In other words, under the SRP  
 193  $\pi(P)$ , we have  $\mathbb{P}(\pi_n = (k, t)) = P(k, t)$ ,  $\forall n \leq N_\pi(B)$ , for all  $(k, t) \in [K] \times \mathbb{T}$ . We denote the  
 194 class of all stationary randomized policies as  $\Pi_S$ .*

195 **Proposition 1** (Asymptotic optimality of SRP). *There exists a probability distribution  $P^*$  such that  
 196 the stationary randomized policy  $\pi(P^*)$  is asymptotically optimal over all causal policies as  $B \rightarrow \infty$ .  
 197*

198 The proof of Proposition 1 can be found in Appendix B. In the following, we characterize the total  
 199 utility under  $\pi(P)$  by providing tight bounds.

**Proposition 2.** *Let  $P$  be any given probability distribution over  $[K] \times \mathbb{T}$ . Then, the reward per unit  
 time for group  $k$  under the stationary randomized policy  $\pi(P)$  is as follows:*

$$\rho_k(P) = \frac{\sum_{t \in \mathbb{T}} P(k, t) \mathbb{E}[R_{k,1}(t)]}{\sum_{(i,t) \in [K] \times \mathbb{T}} P(i, t) \mathbb{E}[\min\{X_{i,1}, t\}]}, \forall k \in [K].$$

Consequently, the total utility under the stationary randomized policy  $\pi(P)$  is bounded as follows:

$$\sum_{k \in [K]} U_k(\rho_k(P)) \leq \sum_{k \in [K]} U_k(\bar{r}_k^{\pi(P)}(B)) \leq \sum_{k \in [K]} U_k(\rho_k(P)) + O\left(\frac{1}{B}\right).$$

200

201 We include the complete proof of Proposition 2 in Appendix B. The key idea is that under an SRP,  
 202 the total reward of a group  $k$  is a regenerative process. Then, by using the theory of stopped random  
 203 walks for regenerative processes, the reward per unit time under  $\pi(P)$  is found as  $\rho_k(P)$ , and the  
 204 upper bound for the total utility is found by using Lorden's inequality [34] and concavity of  $U_k$ .

205 Proposition 2 emphasizes the significance of the reward per unit time  $\rho_k(P)$ . In conjunction with  
 206 Proposition 1, this suggests that using a probability distribution that maximizes the limiting total  
 207 utility would be an effective offline approximation.

208 **Definition 3** (Optimal Stationary Randomized Policy). *Let  $P^*$  be a probability distribution defined  
 209 as  $P^* \in \arg \max_P \sum_{k \in [K]} U_k(\rho_k(P))$ . Then, the optimal SRP  $\pi^*$  makes a selection independently  
 210 for every task according to  $P^*$ :  $\mathbb{P}(\pi_n^* = (k, t)) = P^*(k, t)$  for all  $(k, t) \in [K] \times \mathbb{T}$  and  $n \leq N_\pi(B)$ .*

211 An interesting question regarding  $P^*$  is the choice of deadline policy for each group. The following  
 212 proposition characterizes the optimal deadline policy under  $\pi^*$ , and yields a significant simplification  
 213 in finding the optimal policy by reducing the size of the search space.

214 **Proposition 3** (Optimal Deadline Policy). *For any  $k$ , the optimal probability distribution  $P^*$  makes  
 215 a deterministic deadline decision for group  $k$ , that is,  $|\{t \in \mathbb{T} : P^*(k, t) > 0\}| \leq 1$ . For any  $k$ , we  
 216 denote  $t_k^* \in \mathbb{T}$  as the (unique) optimal deadline for group  $k$  such that  $P^*(k, t_k^*) > 0$ .*

217 The detailed proof of Prop. 3 can be found in Appendix C. As we will see later, we can explicitly  
 218 characterize the optimal deadline for a broad class of utility functions used for the so-called  $\alpha$ -fair  
 219 allocations. In the following, we use Prop. 2 to characterize the performance of the optimal SRP.

220 **Proposition 4** (Optimal Total Utility). *For any group  $k$ , let  $t_k^* \in \mathbb{T}$  be the (unique) optimal deadline  
 221 by Prop. 3;  $r_k^* = \mathbb{E}[R_{k,1}(t_k^*)]/\mathbb{E}[\min\{X_{k,1}, t_k^*\}]$  be the reward per processing time for group  $k$ ; and*

$$\varphi_k = \frac{P^*(k, t_k^*) \cdot \mathbb{E}[\min\{X_{k,1}, t_k^*\}]}{\sum_{j \in [K]} P^*(j, t_j^*) \cdot \mathbb{E}[\min\{X_{j,1}, t_j^*\}]}, \quad (4)$$

222 *be the fraction of time budget allocated to group  $k$  under  $\pi(P^*)$ . Then, for any SRP  $\pi(P)$ , the total  
 223 utility is bounded as  $\sum_k U_k(\rho_k(P)) \leq \sum_k U_k\left((U'_k)^{-1}\left(\frac{\lambda}{r_k^*}\right)\right)$ , where the upper bound is achieved  
 224 by the probability distribution that satisfies  $\varphi_k = \frac{1}{r_k^*} (U'_k)^{-1}\left(\frac{\lambda}{r_k^*}\right)$  for  $\lambda$  such that  $\sum_k \varphi_k = 1$ .*

225 The proof of Proposition 4 follows from Lagrange duality and Prop. 3, and can be found in Appendix  
 226 D. Note that the above analysis is very general in the sense that it holds for any set of utility functions  
 227  $\{U_k : \mathbb{R} \rightarrow \mathbb{R} : k \in [K]\}$  that are continuously differentiable and concave. In the following, we  
 228 apply the results to the class of  $\alpha$ -fair allocations (cf. Definition 1) and discuss their implications.

229 **Proposition 5** ( $\alpha$ -Fair Resource Allocation in Continuous Time). *For any group  $k$ , the optimal  
 230 deadline is  $t_k^* = \arg \max_{t \in \mathbb{T}} \frac{\mathbb{E}[R_{k,1}(t)]}{\mathbb{E}[\min\{X_{k,1}, t\}]}$ . Also, let  $r_k^* = \max_{t \in \mathbb{T}} \frac{\mathbb{E}[R_{k,1}(t)]}{\mathbb{E}[\min\{X_{k,1}, t\}]}$  be the reward  
 231 per processing time and  $\mu_k = \mathbb{E}[\min\{X_{k,1}, t_k^*\}]$  be the mean processing time for group  $k$ . Then, for  
 232 any  $\alpha > 0$ , we have the following results for  $\alpha$ -fair utility functions:*

$$\max_P U^{\pi(P)}(B) = \frac{1}{1 - \alpha} \left( \sum_{k \in [K]} (r_k^*)^{\frac{1}{\alpha} - 1} w_k^{\frac{1}{\alpha}} \right)^\alpha, \quad (5)$$

where the optimum probability distribution  $P_k^*$  and the optimum fraction of time budget  $\varphi_k$  allocated  
 to group  $k$  are, respectively, given by:

$$P^*(k, t) = \mathbb{I}\{t = t_k^*\} \frac{w_k^{\frac{1}{\alpha}} (r_k^*)^{\frac{1}{\alpha} - 1} / \mu_k}{\sum_{j \in [K]} w_j^{\frac{1}{\alpha}} (r_j^*)^{\frac{1}{\alpha} - 1} / \mu_j}, \quad \varphi_k = \frac{(r_k^*)^{\frac{1}{\alpha} - 1} w_k^{\frac{1}{\alpha}}}{\sum_{j \in [K]} (r_j^*)^{\frac{1}{\alpha} - 1} w_j^{\frac{1}{\alpha}}}, \quad \forall k \in [K].$$

233

234 To gain a clear understanding of the notion of  $\alpha$ -fairness, we consider the following special cases.

235 **Corollary 1.** *For any given set of parameters  $\{w_k > 0 : k \in [K]\}$ , we have the following results for  
 236 continuous-time  $\alpha$ -fair resource allocation problem for various  $\alpha > 0$  values.*

237 (i) **Proportional fairness:** *In this case, we have  $\lim_{\alpha \rightarrow 1} U_k(x) = w_k \log(x)$  for all  $k$ . Let  $\mu_k =$   
 238  $\mathbb{E}[\min\{X_{k,1}, t_k^*\}]$  be the mean processing time for group  $k$ . Then, the optimum utility is achieved  
 239 by the probability distribution  $P^*(k, t) = \mathbb{I}\{t = t_k^*\} \frac{w_k / \mu_k}{\sum_{j \in [K]} w_j / \mu_j}$ ,  $(k, t) \in [K] \times \mathbb{T}$ , thus we  
 240 have  $\varphi_k = \frac{w_k}{\sum_{j \in [K]} w_j}$  for all  $k$  and  $\text{OPT}_{\Pi_S}(B) = \sum_k \log\left(\frac{r_k^* w_k}{\sum_{k' \in [K]} w_{k'}}\right) + O\left(\frac{1}{B}\right)$ .*

241 (ii) **Reward maximization:** If  $\alpha = 0$ , we have  $U_k(x) = \omega_k x$  for all  $k$ . Let  $k^* = \arg \max_{k \in [K]} w_k r_k^*$   
 242 be the group with highest weighted reward rate. Then, the optimal probability distribution is  
 243  $P^*(k, t) = \mathbb{I}\{k = k^*, t = t_k^*\}$ , for all  $(k, t)$ . Thus,  $\text{OPT}_{\Pi_S}(B) = \max_{k \in [K]} w_k r_k^* + O(1/B)$ .

244 **Remark 1.** Note that optimal deadline  $t_k^*$  for any group  $k$  is chosen so as to maximize the reward per  
 245 processing time of group  $k$ . Under proportional fairness ( $\alpha \rightarrow 1$ ), the controller distributes the time  
 246 budget proportional to group weights, i.e.,  $\varphi_k = w_k / \sum_j w_j$ , which reduces to equal time-sharing  
 247 under uniform weights. To achieve this, the controller allocates tasks with probability inversely  
 248 proportional to the mean processing time  $\mu_k$ . Under reward maximization ( $\alpha = 0$ ), the controller  
 249 allocates the entire time budget  $B$  to a single group that yields the highest reward per processing time  
 250 to maximize the expected total reward, i.e.,  $\varphi_k = \mathbb{I}\{k = k^*\}$ . As such, the trade-off between reward  
 251 maximization and equal (i.e., reward-insensitive) time-sharing is modeled by  $\alpha$ -fairness for any  
 252  $\alpha \in [0, 1)$ . Further, the  $\alpha$ -fair utility maximization framework includes max-min fairness ( $\alpha \rightarrow \infty$ )  
 253 and minimum potential delay fairness ( $\alpha = 2$ ) as subcases.

## 254 4 Online Learning for Utility Maximization (OLUM)

255 In the previous section, we provided key results on the asymptotically optimal approximations to  
 256 the offline utility maximization problem. In this section, we will build on these to attack the online  
 257 learning problem for continuous-time fair allocation. In particular, we will propose a novel light-  
 258 weight online learning algorithm for the fair resource allocation problem based on Lagrangian duality,  
 259 and show that it achieves vanishing regret at rate  $\tilde{O}(B^{-1/2})$ .

260 **Feedback model:** We assume a delayed full-information feedback model where the completion time  
 261 and reward of all groups for task  $n$  are revealed to the controller at stage  $n + \tau$  for some delay  $\tau \geq 1$ .

262 This assumption holds approximately for our target applications. In freelancing platforms, there  
 263 are often multiple contractors that hire freelancers for various tasks. It is often possible to get full  
 264 information on various freelancers due to employment by other companies and their reviews can  
 265 serve as the feedback for the controller. Competitions hosting websites like TOPCODER have also  
 266 recently been catering to businesses who need fast-prototyping using freelancers. In their business  
 267 model, a controller might invest in a few topcoders at a time, however, she can potentially get access  
 268 to updated rankings (quality and time to complete tasks) via topcoder competitions over time. In  
 269 server cations such as Amazon AWS and Microsoft Azure as well, although a controller might  
 270 be optimizing operations on a local set of servers, they can request task performance data from a  
 271 centralized server or a scheduler after a delay in time [35]. This feedback model already presents  
 272 with technical challenges due to random completion times, as we discuss next.

273 In order to design the online learning algorithm, let us define, for any  $(k, t) \in [K] \times \mathbb{T}$ , the empirical  
 274 estimates of the mean completion time and reward after  $n$  stages, respectively, as

$$\hat{\mu}_{k,n}(t) = \frac{1}{n} \sum_{i=1}^n \min\{t, X_{k,i}\}, \quad \text{and} \quad \hat{\theta}_{k,n}(t) = \frac{1}{n} \sum_{i=1}^n R_{k,i}(t).$$

275 **Definition 4** (OLUM Algorithm). For any  $k$ , let  $Q_{k,0} = 1$  and  $Q_{k,i}$  be defined recursively as follows:

$$Q_{k,i+1} = \left( Q_{k,i} + \gamma_k(i) \min\{X_{G_i,i}, T_i\} - R_{k,i}(T_i) \mathbb{I}\{G_i = k\} \right)^+, \quad i > 0 \quad (6)$$

where the auxiliary variable  $\gamma_k(i) = (U_k^i)^{-1}(Q_{k,i}/V)$ , where  $V > 0$  is a design choice. Then, for  
 the task  $n$ , the OLUM Algorithm, denoted by  $\pi^{\text{OLUM}}$ , makes the following decision:

$$(G_n, T_n) \in \arg \max_{(k,t) \in [K] \times \mathbb{T}} \frac{\hat{\theta}_{k,n-\tau}(t) Q_{k,n}}{\hat{\mu}_{k,n-\tau}(t)}.$$

276 Upon observing the corresponding feedback, the controller updates  $Q_{k,n+1}$  via (6).

277 **Interpretation:** The OLUM Algorithm aims to maximize the time-average reward weighted with  $Q_{k,n}$   
 278 at each round. Note that for any  $k \in [K]$ , if the sequence  $Q_{k,n}$  gets very big, then its reward rate is  
 279 much smaller than the optimal value, thus the controller tends to select that group. In other words, the  
 280 magnitude of  $Q_{k,n}$  is a measure of the unfairness that group  $k$  has endured by stage  $n$ . The algorithm  
 281 is designed so as to balance the weights  $Q_{k,n}$  to maximize the total utility.

282 In the following theorem, we prove regret bounds for the OLUM Algorithm.

283 **Theorem 1** (Regret bounds for OLUM). *For any  $V > 0$  and constant delay  $\tau$ , the regret under  $\pi^{\text{OLUM}}$*   
 284 *is bounded as  $\text{REG}_{\Pi_S}^{\pi^{\text{OLUM}}}(B) = O\left(\sqrt{\frac{\log(B)}{B}} + \frac{V}{B} + \frac{1}{V}\right)$ . By choosing  $V = \Theta(\sqrt{B/\log(B)})$ , we*  
 285 *obtain  $\text{REG}_{\Pi_S}^{\pi^{\text{OLUM}}}(B) = O(\sqrt{\log(B)/B}) = \tilde{O}(1/\sqrt{B})$ .*

286 The proof is based on PAC bounds and stochastic dual optimization, and can be found in Appendix E.

## 287 5 Simulations

288 We implemented the OLUM Algorithm on a fair resource allocation problem with  $K = 2$  groups.  
 289 In the application domains that we considered in Section 2, the task completion times naturally  
 290 follow a power-law distribution. For example, in the server allocation example, empirical studies  
 291 indicate that the distribution of job execution times can be accurately approximated by a Pareto(1,  
 292  $\gamma$ ) distribution with exponent  $\gamma \in (0, 2)$  [36]. Similarly, for the contractual online hiring setting,  
 293 creativity of individuals has been shown to follow a Pareto(1,  $\gamma$ ) distribution with exponent  $\gamma > 1$ ,  
 294 where  $\gamma$  is dependent on the field of expertise [37]. Motivated by these applications, we consider the  
 295 following group statistics:

- 296 • **Group 1:**  $X_{k,n} \sim \text{Pareto}(1, 1.2)$  and  $R_{k,n}(t) = X_{k,n}^{0.6} \cdot \mathbb{I}\{X_{k,n} \leq t\}$
- 297 • **Group 2:**  $X_{k,n} \sim \text{Pareto}(1, 1.4)$  and  $R_{k,n}(t) = X_{k,n}^{0.2} \cdot \mathbb{I}\{X_{k,n} \leq t\}$

298 The reward per processing time as a function of the deadline is shown in Figure 2. Note that the  
 optimal deadline improves the reward per unit processing time. For this setting, we implemented the

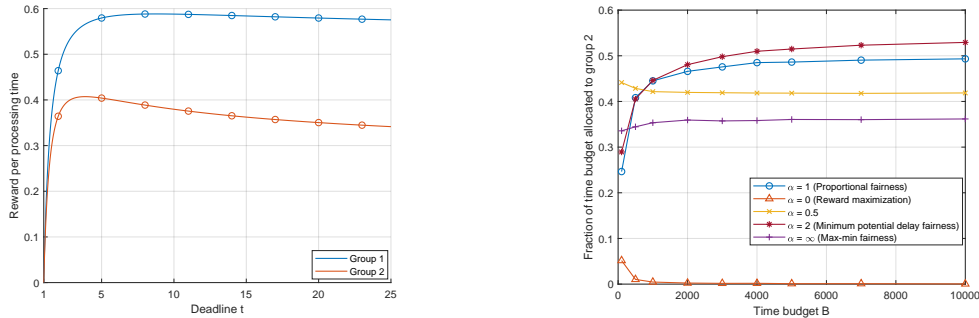


Figure 2: (Left) Reward per processing time for each group. (Right) Fraction of time budget assigned to Group-2 individuals under the OLUM Algorithm for various fairness criteria.

299 OLUM Algorithm with parameter  $V = 20$ , and considered  $\alpha$ -fair resource allocation problems with  
 300 various  $\alpha$  values. In Figure 2, we present the simulation results for  $\varphi_2$ , i.e., the average fraction of  
 301 time budget  $B$  allocated to Group-2 individuals, under the OLUM Algorithm. For these experiments,  
 302 we chose  $w_k = 1$  for  $k = 1, 2$  and ran the OLUM Algorithm for 1000 trials for each set. Note that the  
 303 optimal reward per processing time of Group-1 individuals is higher than that of Group-2 individuals,  
 304 thus Group-1 is chosen for reward maximization. Under proportional fairness, the time budget is  
 305 equally distributed between Group-1 and Group-2 individuals. We observe from Figure 2 that the  
 306 OLUM Algorithm converges to the optimal operating points very fast, which verifies the theoretical  
 307 results we presented.

## 309 6 Conclusion

310 In this paper, we proposed a versatile and comprehensive framework for continuous-time online  
 311 resource allocation with fairness considerations, and proposed a no-regret learning algorithm for  
 312 this problem in a delayed full-information feedback model. Note that although the full-information  
 313 feedback is available in many application scenarios, there are cases in which the controller does not  
 314 have an access to full feedback, thus a mechanism that incorporates bandit feedback is required. The  
 315 online learning framework introduced in this paper can be extended to bandit feedback. One way to  
 316 achieve this might be to replace the empirical estimates with upper confidence bounds in the OLUM  
 317 Algorithm, which makes the analysis even more complicated. We leave the design and analysis of  
 318 bandit algorithms in this setting as a future work.



## 319 **Broader Impact**

320 Our work develops the theory of fair online learning, specifically analyzing the impact of reward-  
321 maximizing allocation policies on opportunities for different groups of people. Our proposal analyzes  
322 the trade-offs across various allocation policies (ranging from profit maximizing to equal opportunity  
323 for all), thus highlighting the choice of objectives that the controllers should carefully consider. This  
324 work does not have any foreseeable negative ethical or societal impact.

## 325 **References**

- 326 [1] J. Zhao, T. Wang, M. Yatskar, V. Ordonez, and K.-W. Chang, “Men also like shopping: Reducing  
327 gender bias amplification using corpus-level constraints,” arXiv preprint arXiv:1707.09457,  
328 2017.
- 329 [2] T. Bolukbasi, K.-W. Chang, J. Y. Zou, V. Saligrama, and A. T. Kalai, “Man is to  
330 computer programmer as woman is to homemaker? debiasing word embeddings,” in  
331 Advances in Neural Information Processing Systems, 2016, pp. 4349–4357.
- 332 [3] A. Caliskan, J. J. Bryson, and A. Narayanan, “Semantics derived automatically from language  
333 corpora contain human-like biases,” Science, vol. 356, no. 6334, pp. 183–186, 2017.
- 334 [4] K. Lum and W. Isaac, “To predict and serve?” Significance, vol. 13, no. 5, pp. 14–19, 2016.
- 335 [5] A. L. Washington, “How to argue with an algorithm: Lessons from the compas-propublica  
336 debate,” Colo. Tech. LJ, vol. 17, p. 131, 2018.
- 337 [6] J. Kleinberg, “Inherent trade-offs in algorithmic fairness,” in Abstracts of the 2018 ACM  
338 International Conference on Measurement and Modeling of Computer Systems, 2018, pp. 40–  
339 40.
- 340 [7] A. Chouldechova, “Fair prediction with disparate impact: A study of bias in recidivism predic-  
341 tion instruments,” Big data, vol. 5, no. 2, pp. 153–163, 2017.
- 342 [8] G. Laumeister, “The next big thing in e-commerce: Online labor marketplaces,” Forbes (Online),  
343 2014.
- 344 [9] H. Torry, “Coronavirus pandemic deepens labor divide between online, offline workers,” Wall  
345 Street Journal, 2020.
- 346 [10] A. Teeley, “There are 57 million u.s. independent professionals — upwork wants them all to  
347 succeed,” Built In Chicago, 2020.
- 348 [11] A. Hannák, C. Wagner, D. Garcia, A. Mislove, M. Strohmaier, and C. Wilson, “Bias in online  
349 freelance marketplaces: Evidence from taskrabbit and fiverr,” in Proceedings of the 2017  
350 ACM Conference on Computer Supported Cooperative Work and Social Computing, 2017, pp.  
351 1914–1933.
- 352 [12] R. Srikant and L. Ying, Communication networks: an optimization, control, and stochastic  
353 networks perspective. Cambridge University Press, 2013.
- 354 [13] D. Bertsimas, V. F. Farias, and N. Trichakis, “On the efficiency-fairness trade-off,” Management  
355 Science, vol. 58, no. 12, pp. 2234–2250, 2012.
- 356 [14] K. Jain and V. V. Vazirani, “Eisenberg-gale markets: Algorithms and structural properties,”  
357 in Proceedings of the thirty-ninth annual ACM symposium on Theory of computing, 2007, pp.  
358 364–373.
- 359 [15] D. P. Palomar and M. Chiang, “A tutorial on decomposition methods for network utility  
360 maximization,” IEEE Journal on Selected Areas in Communications, vol. 24, no. 8, pp. 1439–  
361 1451, 2006.
- 362 [16] A. Eryilmaz and R. Srikant, “Fair resource allocation in wireless networks using queue-length-  
363 based scheduling and congestion control,” IEEE/ACM transactions on networking, vol. 15,  
364 no. 6, pp. 1333–1344, 2007.
- 365 [17] H. J. Kushner and P. A. Whiting, “Convergence of proportional-fair sharing algorithms under  
366 general conditions,” IEEE Transactions on Wireless Communications, vol. 3, no. 4, pp. 1250–  
367 1259, 2004.

- 368 [18] D. Kahneman and R. H. Thaler, “Anomalies: Utility maximization and experienced utility,”  
369 Journal of economic perspectives, vol. 20, no. 1, pp. 221–234, 2006.
- 370 [19] M. J. Neely, “Dynamic optimization and learning for renewal systems,” IEEE Transactions on  
371 Automatic Control, vol. 58, no. 1, pp. 32–46, 2012.
- 372 [20] A. Badanidiyuru, R. Kleinberg, and A. Slivkins, “Bandits with knapsacks,” Journal of the ACM  
373 (JACM), vol. 65, no. 3, pp. 1–55, 2018.
- 374 [21] L. Tran-Thanh, A. Chapman, A. Rogers, and N. R. Jennings, “Knapsack based optimal policies  
375 for budget–limited multi–armed bandits,” in Twenty-Sixth AAAI Conference on Artificial  
376 Intelligence, 2012.
- 377 [22] A. Slivkins, “Introduction to multi-armed bandits,” arXiv preprint arXiv:1904.07272, 2019.
- 378 [23] S. Cayci, A. Eryilmaz, and R. Srikant, “Learning to control renewal processes with bandit  
379 feedback,” Proceedings of the ACM on Measurement and Analysis of Computing Systems,  
380 vol. 3, no. 2, pp. 1–32, 2019.
- 381 [24] S. Agrawal and N. R. Devanur, “Bandits with concave rewards and convex knapsacks,” in  
382 Proceedings of the fifteenth ACM conference on Economics and computation, 2014, pp. 989–  
383 1006.
- 384 [25] A. Rosenblat, K. E. Levy, S. Barocas, and T. Hwang, “Discriminating tastes: Customer ratings  
385 as vehicles for bias,” Available at SSRN 2858946, 2016.
- 386 [26] A. Chakraborty, A. Hannak, A. J. Biega, and K. P. Gummadi, “Fair sharing for sharing economy  
387 platforms,” 2017.
- 388 [27] M. Harchol-Balter, “Task assignment with unknown duration,” in Proceedings 20th IEEE  
389 International Conference on Distributed Computing Systems. IEEE, 2000, pp. 214–224.
- 390 [28] R. Motwani, S. Phillips, and E. Torng, “Nonclairvoyant scheduling,” Theoretical computer  
391 science, vol. 130, no. 1, pp. 17–47, 1994.
- 392 [29] M. Harchol-Balter and A. B. Downey, “Exploiting process lifetime distributions for dynamic  
393 load balancing,” ACM Transactions on Computer Systems (TOCS), vol. 15, no. 3, pp. 253–285,  
394 1997.
- 395 [30] K. Kim and A. A. Tsiatis, “Study duration for clinical trials with survival response and early  
396 stopping rule,” Biometrics, pp. 81–92, 1990.
- 397 [31] P. F. Thall, R. Simon, and S. S. Ellenberg, “Two-stage selection and testing designs for compar-  
398 ative clinical trials,” Biometrika, vol. 75, no. 2, pp. 303–310, 1988.
- 399 [32] P. R. Jelenković and J. Tan, “Characterizing heavy-tailed distributions induced by retransmis-  
400 sions,” Advances in Applied Probability, vol. 45, no. 1, pp. 106–138, 2013.
- 401 [33] C. H. Papadimitriou and J. N. Tsitsiklis, “The complexity of optimal queuing network control,”  
402 Mathematics of Operations Research, vol. 24, no. 2, pp. 293–305, 1999.
- 403 [34] S. Asmussen, Applied probability and queues. Springer Science & Business Media, 2008,  
404 vol. 51.
- 405 [35] R. Zabolotnyi, P. Leitner, and S. Dustdar, “Profiling-based task scheduling for factory-worker  
406 applications in infrastructure-as-a-service clouds,” in 2014 40th EUROMICRO Conference on  
407 Software Engineering and Advanced Applications. IEEE, 2014, pp. 119–126.
- 408 [36] M. Harchol-Balter, “The effect of heavy-tailed job size distributions on computer system de-  
409 sign,” in Proc. of ASA-IMS Conf. on Applications of Heavy Tailed Distributions in Economics,  
410 Engineering and Statistics, 1999.
- 411 [37] J. Kleinberg and M. Raghavan, “Selection problems in the presence of implicit bias,” arXiv  
412 preprint arXiv:1801.03533, 2018.
- 413 [38] A. Gut, Stopped random walks. Springer, 2009.
- 414 [39] M. Neely, Stochastic network optimization with application to communication and queueing  
415 systems. Morgan & Claypool Publishers, 2010.
- 416 [40] S. Cayci, A. Eryilmaz, and R. Srikant, “Budget-constrained bandits over general cost and reward  
417 distributions,” arXiv preprint arXiv:2003.00365, 2020.

- 418 [41] M. J. Wainwright, High-dimensional statistics: A non-asymptotic viewpoint. Cambridge  
419 University Press, 2019, vol. 48.
- 420 [42] M. J. Neely, “Dynamic optimization and learning for renewal systems,” IEEE Transactions on  
421 Automatic Control, vol. 58, no. 1, pp. 32–46, 2012.