# An Efficiency-boosting Client Selection Scheme for Federated Learning with Fairness Guarantee

Tiansheng Huang, Weiwei Lin, Wentai Wu, Ligang He, Kegin Li, Fellow, IEEE, and Albert Y. Zomaya, *Fellow, IEEE* 

Abstract—The issue of potential privacy leakage during centralized Al's model training has drawn intensive concern from the public. A Parallel and Distributed Computing (or PDC) scheme, termed Federated Learning (FL), has emerged as a new paradigm to cope with the privacy issue by allowing clients to perform model training locally, without the necessity to upload their personal sensitive data. In FL, the number of clients could be sufficiently large, but the bandwidth available for model distribution and re-upload is quite limited, making it sensible to only involve part of the volunteers to participate in the training process. The client selection policy is critical to an FL process in terms of training efficiency, the final model's quality as well as fairness. In this paper, we will model the fairness guaranteed client selection as a Lyapunov optimization problem and then a C2MAB-based method is proposed for estimation of the model exchange time between each client and the server, based on which we design a fairness guaranteed algorithm termed RBCS-F for problem-solving. The regret of RBCS-F is strictly bounded by a finite constant, justifying its theoretical feasibility. Barring the theoretical results, more empirical data can be derived from our real training experiments on public datasets.

Index Terms—Client selection, Contextual combinatorial multi-arm bandit, Fairness scheduling, Federated learning, Lyapunov optimization.

#### Introduction

#### Background

FEDERATED Learning (FL) has been esteemed as one of the most promising solutions to the crisis known as isolated "data island". It helps break down the obstacles between parties or entities, allowing a greater extent of data sharing. All the entities being involved could benefit from such a new paradigm, in which model owners could build a more robust and comprehensive model with more data being accessible. Meanwhile, data owners might either receive substantial rewards or services that match their interests in return. More importantly, the privacy of the data owners would not risk being intruded since their raw data simply does not necessarily need to leave the local devices, as all the training is only performed locally.

#### 1.2 Motivations

Within such a novel paradigm, new challenges co-exist with opportunities. Unlike the traditional model training process, not all the data within the system could be accessed over every round of training. Owing to the limited bandwidth and the dynamic status of the training clients, only a fraction

- T. Huang and W. Lin (corresponding author) are with the School of Computer Science and Engineering, South China University of Technology, China. Email: cs\_tianshenghuang@mail.scut.edu.cn, linww@scut.edu.cn.
- W. Wu, L. He are with the Department of Computer Science, the University of Warwick. Email: wentai.wu, Ligang.He@warwick.ac.uk.
- K. Li is with the Department of Computer Science, State University of New York, New Paltz, NY 12561 USA. E-mail: lik@newpaltz.edu.
- A. Y. Zomaya is with the School of Computer Science, The University of Sydney, Sydney, Australia. Email: albert.zomaya@sydney.edu.au.

of them could be picked to perform training on behalf of the model owner. From the perspective of a model owner, the selection decision in each round could have a profound impact on the model's training time, convergence speed, training stability, as well as the final achieved accuracy. Some studies in the literature have made iconic contributions to this problem. To illustrate, in [1], when making a selection, Nishio et al. concentrate on the evaluation of communication time, which accounts for a considerable portion of time for a training round. In another study [2], the authors consider more. They further take the energy consumption factor into consideration. Barring an intelligent decision on participant selection, an efficient bandwidth allocation scheme was also given by them. However, the current line of research evades two important factors. For one thing, both of them assume a pre-known local training time to the scheduler, which may not be realistic in all circumstances. For another, indicated by Theorem 2 in [2], devices with higher performance are more favored by their proposed methods. Indeed, always selecting the "fast" devices somehow boost the training process. But clients with low priority are simply being deprived of chances to participate at the same time, which we refer to it as an unfair selection among clients. In fact, such an extreme selection scheme might bring undesirable side effects by neutralizing some portions of data. Conceivably, with a smaller amount of data involved, data diversity can not be guaranteed, thereby hurting the performance of model training to some extent. This motivates us to develop an algorithm that strikes a good balance between training efficiency and fairness. Also, the algorithm is supposed to be intelligent enough to predict the training time of the clients based on their reputation (or their historical performance), rather than assuming it to be known a priori.

#### 1.3 Contributions

The main contributions of this paper are listed as follows:

- 1) We investigate the client selection in FL from the perspective of minimizing average model exchange time when subjecting to a relatively flexible long-term fairness guarantee, as well as a few rigid system constraints. At the same time, more factors, involving the clients' availability, unknown and stochastic training time, as well as the dynamic communication status, are taken into account.
- Inspired by [3], we transform the original offline problem into an online Lyapunov optimization problem where the long-term guarantee of client participating rate is quantified using dynamic queues.
- 3) We build a Contextual Combinatorial Multi Arm Bandit (C<sup>2</sup>MAB) model for estimation of the model exchange time of each client based on their contextual properties and historical performance (or their reputation).
- 4) A fairness guaranteed selection algorithm RBCS-F is proposed for efficiently resolving the proposed optimization problem in FL. Theoretical evaluation and real data-based experiments show that RBCS-F can ensure no violation in the long-term fairness constraint. Besides, the training efficiency has been significantly enhanced, while the final model accuracy remains close, in a comparison with random, i.e., the vanilla client selection scheme of FL.

To the best knowledge of the authors, this is the first trackable practice that combines Lyapunov optimization and  ${\rm C^2MAB}$  for a long-term constrained online scheduling problem. Also, we shall remind the readers that the proposed combination does not confine to the application of our current proposed problem, but it has the potential to extend to a wider range of selection problems. (e.g. worker selection in crowdsensing, channel selection in the wireless network, etc.)

#### 2 RELATED WORKS

In recent years, we are experiencing a great surge of Edge Intelligence (see in [4]–[6]). Numerous attempts have been made to combine AI techniques and edge, tapping the profound potential of the ubiquitous deployed edge devices. Among these, one of the most iconic studies could be neurosurgeon [7]. Its basic idea is to partition an intact DNN (Deep Neural Networks) into several smaller parts and disseminate them to the edge devices. Owing to a low latency between edge and users, inference speed could be significantly improved.

Besides, edge coordinated Federated Learning is another promising combination. Federated Learning [8], which allows data to be trained in local rather than being transmitted to the cloud, is now known as a more secure paradigm for AI's model training. We have witnessed the surge of some plausible applications of FL within these years (e.g. keyboard and emoji prediction in [9], [10], visual object detection in [11], etc). Despite the potential advantages as well as the promising applications of FL, the communication

overhead between cloud and users renders as a bottleneck for it. A lengthy communication round during training might significantly degrade FL's training performance. Although more advance training schemes, such as federated distillation (FD, originally proposed in [12]), promise us a more desirable, reduced size information exchange between users and model aggregator, the latency between cloud and edge alone is inevitable. Such a defect could be better addressed by making edge the model's aggregator or at least an intermediate one (see in [13]). In this way, the data don't have to bear an outstanding communication length to the cloud. Another open problem of FL we would like to mention here is the client selection problem, originally proposed in [1] and followed by some related works (e.g. [2], [14]-[16]). Many of them see the problem from a communication perspective, focusing on building an efficient selection or bandwidth allocation scheme that helps shorten the communication length. In this paper, we will see the problem from a different angle, namely, to investigate how the fairness factor affects the training performance. We couldn't check out any specialized studies on this topic yet and we hope our research could bring some new insights in the field. Last but not least, we also want to note, FL itself is now far from its maturity, many important issues worth our study. Some of which might involve asynchronous or semi-asynchronous aggregation protocol [17], [18], incentive mechanism [19], [20] and security issues [21], etc. We look forward to more insightful and dedicated research into FL.

Now we would like to talk more about a classical problem, termed multi-arm bandit (MAB). In a classical MAB setting, arms are characterized by different unknown reward distribution. In each round of play, the player selects one of the arms from the possible options and gains a reward sampling from the selected arm's reward distribution. As there exists a tradeoff between exploration and exploitation for the player, how to maximize her obtained reward is the main concern. Several solutions, such as the well-known Upper Confidence Bound (UCB), Thompson Sampling (TS) could be applied to the problems. In addition, MAB has several variations. Those include combinatorial MAB, where players are allowed to select more than one arms in every round, contextual MAB [22], [23], where the reward of an arm follows a linear stochastic formulation, and a much newer one, contextual combinatorial MAB (C<sup>2</sup>MAB) [24], [25], which is the combination of the above two. We found that  $C^2MAB$  could be well applied to the client selection problem in FL, as each client could be regarded as an arm and our task for each round is to choose a combination of which for participation, thus, in this paper, such a prototype will be used for our model establishment.

#### 3 Preliminary Introduction on FL

In this paper, we consider an edge-coordinated federated learning system, in which edge is functioning as a model aggregator, and the clients (mostly mobile devices) are responsible for doing local training over their private data on behalf of the model's owner. We adopt in our system the most-accepted synchronous scheme for federated learning, which is characterized by training in iterations. For clearness, now we will explicitly explain the workflow of our

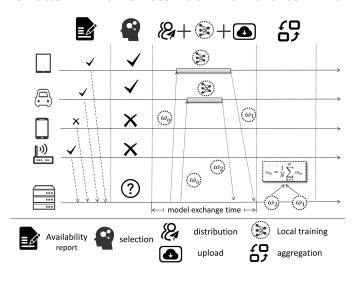


Fig. 1. Illustration of FL

synchronized scheme by giving four sequential stages of training, as follows:

- At the very beginning of a new iteration, the clients first report their willingness to participate in the training as well as a few client-side information, which will be used for the client selection in the next stage.
- 2) In the second step, the scheduler conducts client selection to choose a portion of participants among the volunteers in light of the provided information.
- 3) Global model is distributed to the selected clients. After receiving the model, the clients conduct local training using their private data and update their local model. Once the training of all the selected clients is finished, the local model will be returned to the MEC server. The time span of this round is known as model exchange time.
- 4) The collected local models are aggregated by the server, substituting the original global model that once being distributed, and then it proceeds to step 1) to start a new iteration.

For a more vivid presentation of the training process, we refer the readers to Fig. 1.

## 4 PROBLEM FORMULATION

Our main concern focuses on the selection phase, in which the server makes a decision on the involved clients. Before our formal introduction of the selection problem, we first derive a high-level description of the content of this section. In the first sub-section, we formulate the client selection problem into an offline problem with a long-term fairness constraint. The formulated problem is simple in form but indeed unsolvable due to the time coupling effect as well as the unknown model exchange time persisting in the objective. To resolve the time coupling effect, we transform the problem into an online mode using Lyapunov optimization technique, the online transformation of which gains us a fighting chance to derive an estimated model exchange time

before each round scheduling, which might help resolve another obstacle (i.e. the unknown parameter in the objective function). Specifically, targeting the transformed online problem, a C<sup>2</sup>MAB setting could come in handy for online learning of the exchange time, and being enlightened by which, we are able to further transform the problem into the ultimate form, which concludes the whole section.

Then we need to explain some key notations that are consistently used throughout the paper, among which, a set  $\mathcal{T} \triangleq \{1,2,\ldots\}$ , indexed by t, is used to capture the federated rounds (namely, the iterations in FL's model update process). The set  $\mathcal{N} \triangleq \{1,2,\ldots N\}$  captures all the clients (each indexed by n) in the system. Besides, we assume that the maximum number of selected clients each round is fixed in advance to m. Another important notation is  $\mathcal{S}_t$ , which we use to capture the selected clients in round t and it serves as the representation of the selection policy that we aim to optimize.

#### 4.1 Basic Assumption on System Model

#### 4.1.1 Model Exchange Time

In a client selection problem, an important metric we shall evaluate is the long-term average model exchange time. We refer to the model exchange time as the time span between the instant the scheduler made the selection decision and that when all the re-upload models have been gathered. This model exchange time might involve time for model distribution, model training and model upload. Intuitively, a client selection scheme that is able to achieve a shorter span of each federated round is of interest, since a shorter period of each round explicitly marks shorter time for fix rounds of training. Recall that the server could step into the next phase (model aggregation) only after all the models have been gathered when adopting a synchronous federated training protocol. The time for model exchange is explicitly determined by the participated clients, or more precisely, by the one among them who spends the most time in training and model uploading. Mathematically, we have the following equation to capture the time span for a federated round:

$$f(S_t, \boldsymbol{\tau}_t) = \max_{n \in S_t} \{ \tau_{t,n} \}$$
 (1)

where we use a set  $\mathcal{S}_t$  to capture the selected clients in round t. Besides,  $\tau_{t,n}$  is used to represent the time span between the very beginning of model distribution and the instant when the model from client n being gathered. Here  $\boldsymbol{\tau}_t \triangleq \{\tau_{t,n}\}_{n \in \mathcal{N}}$  in round t is unknown to the scheduler until the end of this round.

#### 4.1.2 Long-Term Fairness Constraint

Another metric that might have a significant impact on FL's performance is fairness. Assume an ideal case that the server is fully aware of the exact model exchange time of each client for the incoming federated round. Then is it incontrovertibly optimized when always choosing the m-fastest clients, making the time span for each round of training minimized? We must note, however, that the answer may not be such apparent. We acknowledge that the time span of each round could be somehow minimized by adopting such a greedy selection scheme, but we must argue that if we

always choose the fastest clients, small chance could become available for their slower counterparts, implicitly implying that little contribution could be obtained from the slowers' local data. Very likely, along with the selection bias, the global model would suffer a degradation on its capability to generalize. In this regard, a greedy selection may not trivially be the best scheme, and fairness in selection is another factor that we need to take into account. To model such a critical fairness concern, we introduce a long-term fairness constraint, as follows:

$$\lim_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} \mathbb{E}[x_{t,n}] \ge \beta \quad \forall n \in \mathcal{N}$$
 (2)

where  $\beta$  models the expected guaranteed chosen rate of clients.  $x_{t,n}$  is used to indicate whether client n is involved in the federated round t or not. In other words,  $x_{t,n}=1$  for  $n \in \mathcal{S}_t$ ; otherwise,  $x_{t,n}=0$ . The constraint is set to make sure the long-term average chosen rate of every client at least greater than  $\beta$ , which somehow helps maintain some degrees of fairness for the system.

#### 4.1.3 Availability of Clients

As we are investigating a client selection problem under a highly dynamic real-world system, it is unrealistic to assume clients are always ready to provide training services. In fact, clients are free to join and leave the loose "federation" at any time they want. With this consideration, we use an indicator function  $I_{t,n}$  to capture the status of a client, indicating whether the client is willing to engage or not. Such information could be given by the availability report from the clients before scheduling. Formally, we introduce a strict constraint to prevent futile participation:

$$I_{t,n} = 1 \quad \forall n \in \mathcal{S}_t$$
 (3)

#### 4.1.4 Selection Fraction

Recall that the maximum number of clients that could be selected is fixed to m in our setting. However, as the number of volunteers may not be able to reach m if the activated number is smaller than m, we have to use a "min" function to constraint the selection fraction, as follows:

$$|\mathcal{S}_t| = \min\left\{m, \sum_{n \in \mathcal{N}} I_{t,n}\right\} \tag{4}$$

where  $|S_t|$  means the number of elements in  $S_t$ . Intuitively, in the case when the total number of availability could not overtake the maximum selection fraction, we simply involve all the active clients for the incoming round of training.

#### 4.2 An Offline Long-Term Optimization Problem

Based on the above discussion, we are ready to introduce our client selection problem, as follows:

$$(P1): \min_{\{\mathcal{S}_1, \mathcal{S}_2, \dots, \mathcal{S}_{\infty}\}} \quad \lim_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} f(\mathcal{S}_t, \boldsymbol{\tau}_t)$$
s.t. (2), (3), (4)

where  $S_t$  captures the selected clients in each round, which is our optimized target. Intuitively, our aim is to minimize the long-term model exchange time while subjecting to

a "soft" long-term fairness constraint (2), which tolerates short-term violation, as well as two extra "hard" constraints (3), (4), which bear no compromise.

One could notice that P1 is a time-coupling scheduling problem, regarding the long-term objective and the fairness constraint in (2). But we note here that such an optimization problem is challenging or even impossible to be solved offline. There are mainly three concerns about this. Firstly, random events, such as clients' availability, are not known to the scheduler until the very beginning of a particular round. This implies that an offline strategy, which is not given access to this particular information, can hardly guarantee the qualifications of constraints (3) and (4). Our second concern is derived from the time-coupling constraint (2), which is quite difficult for the offline solution to deal with. The final concern is that the information on model exchange time can only be observed after actually involving the clients in training. Nevertheless, the scheduler is supposed to make a scheduling decision before the real training process, when the actual model exchange time is unachievable. The lack of this crucial information precludes any feasible attempts to achieve an optimal offline solution. Therefore, for an alternative sub-optimal problem-solving, in the following section, we will elaborate on our transformation of the offline problem to a step-by-step online scheduling problem by Lyapunov optimization to cope with the first two proposed concern. Later, we will display our estimation of model exchange time based on clients' reputation, by which we leverage to deal with our third concern.

## 4.3 Problem Transformation under Lyapunov Framework

In this sub-section, we first take advantage of Lyapunov optimization framework to transform the offline problem *P1* to an online one.

First, we introduce a virtual queue for each client, whose backlog<sup>1</sup> is denoted by  $Z_{t,n}^2$ , to transform the long-term fairness constraint. Specifically,  $Z_{t,n}$  evolves across the FL process complying the following rule:

$$Z_{t+1,n} = [Z_{t,n} + \beta - x_{t,n}]^{+}$$
(6)

where  $\beta$  is the expected guaranteed selection rate in (2) and  $[\dots]^+$  is equivalent to  $\max(\dots,0)$ .

Now we present Theorem 1 to justify the rationale for this transformation.

**Theorem 1.** Long-term time average constraint (2) holds if all the virtual queues (whose backlogs denoted by  $Z_{t,n}$ ) remain mean rate stable across the FL process.

*Proof.* According to the queue theory (see in Theorem 2.5, [26]), if all the virtual queues  $Z_{t,n}$  remain *mean rate stable* across the FL process (or formally,  $\lim_{T\to\infty} \mathbb{E}[Z_{T,n}]/T = 0$ ),

- 1. We use the term "backlogs" and "queue length" interchangeably throughout the paper but actually, they share the same meaning.
- 2. The subscripts t and n here correspond to a federated round and a client, respectively. A similar form of subscript definition will be adopted throughout the paper.

the time average arrival rates of the queue will be smaller than the service rates, namely, we have:

$$\frac{1}{T} \lim_{T \to \infty} \sum_{t=1}^{T} \mathbb{E}[\beta - |\mathcal{S}_t|] \le 0 \tag{7}$$

Through basic mathematics operations, we can reconstruct the above inequality into the form of (2) with ease. This completes the proof.  $\Box$ 

Remark. Intuitively, the length of the queue will soar towards infinity if the long-term fairness constraint is violated, (i.e. when the real chosen rate could not match up with the expected guaranteed selection rate), which is formally justified by Theorem 1. To guarantee the fairness constraint, the queue has to remain mean rate stable and a qualified algorithm is supposed to achieve this goal. Apart from this conclusion, we shall note that the stabilized queue length could also reflect the degree of fairness. For example, if a client never being selected in the first few limited round, its corresponding queue length will soar to a positive value. After that, if its real selection rate basically flats with the expected guaranteed selection rate, its queue still remains mean rate stable and the queue length will slightly fluctuate over the same positive value. Intuitively, the bigger this value is, the unfairer the selection policy could be, as it demonstrates more violation of the fairness constraint in the initial stage. This conclusion could also be derived from the results in our experiments, which will be presented later.

With Theorem 1, now we have transformed the troublesome time-coupling constraint into the goal of ensuring the virtual queues mean rate stable across the FL process. To reach this end, a straightforward approach is to bound every increase of queues so that they could not grow to infinity. Under this motivation, we shall leverage Lyapunov optimization technique to bound the growth of virtual queues while simultaneously minimizing the objective in *P1*. First, we establish the quadratic Lyapunov function, with the following form:

$$\mathcal{L}(\mathbf{\Theta}(t)) = \frac{1}{2} \sum_{n \in \mathcal{N}} Z_{t,n}^2 \tag{8}$$

where  $\Theta(t) \triangleq \{Z_{t,n}\}_{n \in \mathcal{N}}$  contains the backlogs of all the virtual queues.

Aiming at bounding the expected increase of  $\mathcal{L}(\Theta(t))$  for one single round, we first formulate the *Lyapunov drift* to measure it, basically, we have:

$$\Delta(\mathbf{\Theta}(t)) = \mathbb{E}[\mathcal{L}(\mathbf{\Theta}(t+1)) - \mathcal{L}(\mathbf{\Theta}(t))|\mathbf{\Theta}(t)] \tag{9}$$

As the backlogs of queues  $\Theta(t)$  can be known to the scheduler when being scheduled in an online manner, we take it as the condition in the Lyapunov drift. It is notable that the conditional expectation here is with respect to the availability of clients (which is a stochastic variable) as well as the possibly random selection policy. For ease of later interpretation, we let  $\omega_t \triangleq \{I_{t,n}\}_{n \in \mathcal{N}}$  to capture the stochastic availability.

Recall that the objective of *P1* is to minimize the model exchange time while satisfying the given constraints. This motivates us to combine the objective function into the drift

function. Formally, we term such a combination as *drift-plus-cost* function, with the following form:

$$\Delta(\mathbf{\Theta}(t)) + V \mathbb{E}[f(\mathcal{S}_t, \boldsymbol{\tau}_t) | \mathbf{\Theta}(t)]$$
 (10)

where  $V\geq 0$  is a penalty factor set for the purpose of balancing the tradeoff between minimizing the objective and satisfying the fairness constraint. Such a parameter is crucial for the algorithm's performance and we will conduct a specific analysis to it in the next section. Note that the conditioned expectation being taken here is also with respect to stochastic events  $\omega(t)$  and the possibly random policy as well. Now we are going to introduce a potential upper bound for the drift-plus-cost function. We show the result by Theorem 2.

**Theorem 2.** Conditioning on the queues' backlogs  $\Theta(t)$ , the drift-plus-cost function for our system model could be bounded into the following form, where  $\Gamma = N(1 + \beta^2)/2$  is a constant.

$$\Delta(\boldsymbol{\Theta}(t)) + V\mathbb{E}[f(\mathcal{S}_{t}, \boldsymbol{\tau}_{t})|\boldsymbol{\Theta}(t)]$$

$$\leq \Gamma + \sum_{n \in \mathcal{N}} Z_{t,n}\mathbb{E}[\beta - x_{t,n}|\boldsymbol{\Theta}(t)] + V\mathbb{E}[f(\mathcal{S}_{t}, \boldsymbol{\tau}_{t})|\boldsymbol{\Theta}(t))]$$
(11)

*Proof.* The complete proof is given in Appendix A.  $\Box$ 

Intuitively, if we minimize the Right Hand Side (R.H.S) of (11), the fairness virtual queues could be somehow maintained stable, while the objective function is also being minimized. Now shall introduce our step-by-step online scheduling problem by giving *P*2:

$$(P2): \min_{\boldsymbol{x}_{t}} \quad \Gamma + \sum_{n \in \mathcal{N}} Z_{t,n}(\beta - x_{t,n}) + V\dot{f}(\boldsymbol{x}_{t}, \boldsymbol{\tau}_{t})$$

$$s.t. \quad \sum_{n \in \mathcal{N}} x_{t,n} = \min \left\{ m, \sum_{n \in \mathcal{N}} I_{t,n} \right\}$$

$$x_{t,n} \leq I_{t,n}$$

$$x_{t,n} \in \{0,1\}$$

$$(12)$$

we first have to make it clear that we use  $x_t$  to substitute all the  $\mathcal{S}_t$  in P1, making it a clearer form. Here  $\dot{f}(x_t, \tau_t) = \max_{n \in \mathcal{N}} \{x_{t,n}\tau_{t,n}\}$  is an equivalent form to  $f(\mathcal{S}_t, \tau_t)$ . While solving P2 on every round, the R.H.S of (11) can be minimized. The rationale behind is quite evident. As we have done the minimization under every round (alternatively, under every  $\omega_t$ , since  $\omega_t$  is an independent sampling for each round), then the expectation with respect to  $\omega_t$  is also being minimized. Note here that  $\omega_t$  is indeed observable for an online algorithm since an online algorithm makes scheduling after the stage of availability report, making it accessible to this particular information.

For briefness, we eliminate all the constants (i.e.  $\Gamma$ ,  $Z_{t,n}\beta$ ) in the objective of P2 and transform it to P3:

$$(P3): \min_{\mathbf{x}_{t}} \quad V \max_{n \in \mathcal{N}} \{x_{t,n} \tau_{t,n}\} - \sum_{n \in \mathcal{N}} Z_{t,n} x_{t,n}$$

$$s.t. \quad \sum_{n \in \mathcal{N}} x_{t,n} = \min \left\{ m, \sum_{n \in \mathcal{N}} I_{t,n} \right\}$$

$$x_{t,n} \leq I_{t,n}$$

$$x_{t,n} \in \{0,1\}$$

$$(13)$$

But note that such a problem remains unsolvable yet since the real model exchange time of all the clients (or  $\tau_{t,n}$ ) is not known to us before real scheduling. In the next subsection, we will present a  $C^2MAB$  estimation to conquer such a barrier.

## 4.4 Estimation of Model Exchange Time with $C^2MAB$

## 4.4.1 Background knowledge on C<sup>2</sup>MAB and UCB

Each round selection in a Contextual Combinatorial Multi Arm Bandit (C<sup>2</sup>MAB) is characterized by a tuple  $(\mathcal{N}, \mathcal{S}_t, \{\boldsymbol{\theta}_n^*\}_{n \in \mathcal{N}}, \{\boldsymbol{c}_{t,n}\}_{n \in \mathcal{N}}, \{\epsilon_{t,n}\}_{n \in \mathcal{N}}, f(\cdot))$ , in which  $\mathcal{N}$ represents the arm set and  $S_t$  is another set that catpures all the possible combination of arms.  $c_{t,n}$  and  $\theta_n^*$  represents the contextual vector and coefficient vector respectively, among which,  $c_{t,n}$  is known before each round scheduling but dynamic between rounds, while  $\theta_n^*$  is unknown but stationary. After each round of scheduling, a combination of arms (often being called as a super arm)  $S_t \subset S_t$  is put into play. Then loss drawn from each selected arm, formulated by  $l_{t,n} = c_{t,n}^{\top} \theta_n^* + \epsilon_{t,n}, n \in S_t$  is revealed to the scheduler, and meanwhile, a collective loss  $f(\{r_{n,t}\}_{n\in S_t})$  is imposed. Our ultimate aim in the C<sup>2</sup>MAB setting is to minimize the expected cumulative penalty  $\frac{1}{T}\sum_{t=1}^T \mathbb{E}\left[f(\cdot)\right]$  as far as possible by a careful selection on  $S_t$ .

Now we shall give a high-level description of a plausible solution for  $\mathrm{C}^2\mathrm{MAB}$ , i.e., a UCB algorithm. The UCB algorithm takes the upper confidence bound as the optimistic estimation of the expected loss in each round. As the historical data accumulated, (i.e.  $l_{t,n}$  in the previous rounds), the bound could be narrowed and eventually converges to the real value, and thereby gaining more precision for our scheduling. By this means, the expected cumulative penalty could be minimized to the full extent with the increase of rounds of play.

#### 4.4.2 Application

Recall that the information of model exchange time, or at least an estimated one, is supposed to be fetched before real client selection. One can take advantage of a MAB based technique to predict the model exchange time for all clients based on their historical performance (or to say, their reputation). In particular, each client can be regarded as an arm³ in a bandit setting and a combination of them (i.e. a super arm) is put into training, after which, the model exchange time for the selected arm, namely,  $\{\tau_{t,n}\}_{n\in\mathcal{S}_t}$  can be observed by the scheduler.

Normally, the model exchange time is associated with the client's computation capacity, running status as well as the bandwidth allocation for the model update. In this regard, we consider introducing linear contextual bandit into our estimation. Formally, we let  $c_{t,n} \triangleq [1/\mu_{t,n}, s_{t,n}, M/B_{t,n}]^{\top}$  denote the contextual feature vectors that are collected by the scheduler before the scheduling phase. More explicitly,  $\mu_{t,n}$  is the ratio of available computation capacity of client n over round t. We can simply comprehend  $\mu_{t,n}$  as the available CPU ratio of the client<sup>4</sup>. A

binary indicator  $s_{t,n}$  indicates if client n has participated in training in the last round. M is the size of the model's parameters (measured by bit) and  $B_{t,n}$  indicates the allocated bandwidth. Barring the available computation capacity of clients (i.e.  $\mu_{t,n}$ ), which have to be proactively reported by the clients, all the other information could be fetched by the servers with ease. Therefore, here we can just comprehend the contextual feature  $c_{t,n}$  as some prior information known by us before we do the scheduling. Given the contextual features, we assume that the sampling value of  $\tau_{t,n}$  complies with the following equation:

$$\tau_{t,n} = \boldsymbol{c}_{t,n}^{\top} \boldsymbol{\theta}_n^* + \epsilon_{t,n} \tag{14}$$

where  $\pmb{\theta}_n^* \triangleq [\tau_n^b, \tau_n^s, 1/\eta]^{ op}$  captures the static coefficient factors that are presumed to be unknown to the scheduler as they are hard to be detected by the server or even by the clients themselves. More explicitly,  $\tau_n^b$  is the local training time for 100% computation capacity. Multiplying it with the first element in  $c_{t,n}$ , we get the approximated local training time under the computation capacity provided by clients.  $\tau_n^s$ denotes the cold start time, multiplying which with the second element  $s_{t,n}$  in contexts yields the real data preparation time. This formulation is derived from the fact that clients who did not undertake the previous round of training need to spend extra time for data preparation, say, loading the data into memory. Likewise, we let  $\eta \triangleq \log(1 + \text{SNR})$  and multiplying which with  $B_{t,n}$  yields the Shannon formula that we use to calculate the uploading data rate. Here SNR is an abbreviation of Signal-to-Noise Ratio, which is associated with the client profile (e.g. transmission power and channel condition). In this regard,  $M/(B_{t,n}\eta)$  can fully represent the model uploading time for client n. In light of our formulation,  $c_{t,n}^{\top} \theta_n^*$  yields the approximation of the expected model exchange time.

In addition, acknowledging some deviation, we admit a noise factor  $\epsilon_{t,n}$  in our estimation, which is assumed to be a zero-mean random variable, conditionally sampling from an unknown distribution with left-bounded support, i.e.  $\operatorname{Supp}(\epsilon_{t,n}|\boldsymbol{c}_{t,n}^{\intercal}) = (a,b]$  where  $a > -\boldsymbol{c}_{t,n}^{\intercal}\boldsymbol{\theta}_n^*$  and b is arbitrary. This assumption is made to ensure that  $\tau_{t,n}$  must be always positive. Also, we have to make sure that  $\epsilon_{t,n}$  is conditionally R-sub-Gaussian where  $R \geq 0$  is a fixed constant. Formally, we need:

$$\forall \Lambda \in \mathbb{R} \quad \mathbb{E}\left[e^{\Lambda \epsilon_{t,n}} \mid \boldsymbol{c}_{1:t,n}, \epsilon_{1:t-1,n}\right] \leq \exp\left(\frac{\Lambda^2 R^2}{2}\right) \quad (15)$$

This assumption is necessary for the regret analysis of a linear bandit, which is also adopted by [23]. Though we admit some loss of generality for the noise assumption, we argue that a great number of distribution families in nature corresponds to R-sub-Gaussian (e.g. any distributions with zero mean bounded support, zero-mean Gaussian distribution, etc), so the assumption would not compromise the objectivity of this paper.

Now we let  $\tau_{t,n}^* = \mathbb{E}[\tau_{t,n}] = c_{t,n}^\top \theta_n^*$ . If  $\tau_{t,n}^*$  is clearly known to us, we can safely substitute  $\tau_{t,n}$  in P3 with it. Recall that  $\theta_n^*$  is an inherent feature of each arm (or client) that is supposed to be static, unchangeable over time. With this assumption, although the scheduler has no access to the real value of  $\theta_n^*$ , which creates a barrier in the calculation of  $\tau_{t,n}^*$ .

<sup>3.</sup> We use an arm to represent a specific client in our later analysis.

<sup>4.</sup> Note that  $\mu_{t,n}$  could exceed 100% since a client could have more than 1 CPUs, say,  $\mu_{t,n}=200\%$  when 2 CPUs are free.

this value can be predicted using the historical information (or the reputation of an arm). For such a linear formulation, ridge regression could suit well. Now we let  $(\mathbf{D}_{t,n},\mathbf{y}_{t,n})$  to represent p pieces of client n's historical performance (i.e. the previous model exchange time and the contexts) that are obtained before round t. Formally, we have:

$$\mathbf{D}_{t,n} = \begin{bmatrix} \mathbf{c}_n^{(1)} \\ \vdots \\ \mathbf{c}_n^{(p)} \end{bmatrix}_{m \times 3} \quad \mathbf{y}_{t,n} = \begin{pmatrix} \tau_n^{(1)} \\ \vdots \\ \tau_n^{(p)} \end{pmatrix}$$
(16)

where  $\mathbf{c}_n^{(p)}$  and  $\tau_n^{(p)}$  respectively represent the context and the real model exchange time of the p-th play of the arm n. With ridge regression, we can empirically estimate  $\boldsymbol{\theta}_n^*$  with  $\hat{\boldsymbol{\theta}}_{t,n}$ :

$$\hat{\boldsymbol{\theta}}_{t,n} = \left(\mathbf{D}_{t,n}^{\top} \mathbf{D}_{t,n} + \lambda \mathbf{I}_{3}\right)^{-1} \mathbf{D}_{n}^{\top} \mathbf{y}_{t,n} \tag{17}$$

For ease of algorithm's design, we then transform  $\hat{\theta}_{t,n}$  into an equivalent form, as follows:

$$\hat{\boldsymbol{\theta}}_{t,n} = \mathbf{H}_{t-1,n}^{-1} \mathbf{b}_{t-1,n} \tag{18}$$

where 
$$\mathbf{H}_{T,n} = \mathbf{H} + \sum_{t=1}^{T} x_{t,n} \mathbf{c}_{t,n} \mathbf{c}_{t,n}^{\mathsf{T}}$$
 and  $\mathbf{b}_{T,n} = \sum_{t=1}^{T} x_{t,n} \tau_{t,n} \mathbf{c}_{t,n}$ . Among which,  $\mathbf{H} = \lambda \mathbf{I}$ .

As we are going to take advantage of the UCB algorithm we previously discussed as our solution, we resort to  $\bar{\tau}_{t,n}$  as the optimistic estimation of  $\tau_{t,n}$ , which has the following form:

$$\bar{\tau}_{t,n} \triangleq \max \left\{ \boldsymbol{c}_{t,n}^{\top} \hat{\boldsymbol{\theta}}_{t,n} - \alpha_t \sqrt{\boldsymbol{c}_{t,n}^{\top} \boldsymbol{H}_{t-1,n}^{-1} \boldsymbol{c}_{t,n}}, 0 \right\}$$
 (19)

where  $\alpha_t$  is an exploration parameter.

Now we show in Lemma 1 the validity of the given confidence bound (i.e. to show the real expected exchange time does not deviate much from the confidence bound with a high probability).

**Lemma 1.** If we set  $\alpha_t = R\sqrt{3\log\left(\frac{1+tL^2/\lambda}{\delta}\right)} + \lambda^{1/2}S$ , with probability at least  $1 - \delta$ , we have

$$0 \le \tau_{t,n}^* - \bar{\tau}_{t,n} \le 2\alpha_t \|\mathbf{c}_{t,n}\|_{\mathbf{H}_{t-1}^{-1}}$$
 (20)

for any round  $t \geq 1$  and any arm  $n \in \mathcal{N}$ 

*Proof.* The complete proof is given in Appendix B. □

We first note here that Lemma 1 will be used in our analysis of regret bound, which will be shown in the next section.

As we have decided  $\bar{\tau}_{t,n}$  as our estimation of  $\tau_{t,n}$ , we now transfer P3 to the ultimate form, shown in the following:

$$(P4): \min_{\mathbf{x}_{t}} \quad V \max_{n \in \mathcal{N}} \{x_{t,n} \bar{\tau}_{t,n}\} - \sum_{n \in \mathcal{N}} Z_{t,n} x_{t,n}$$

$$s.t. \quad \sum_{n \in \mathcal{N}} x_{t,n} = \min \left\{ m, \sum_{n \in \mathcal{N}} I_{t,n} \right\}$$

$$x_{t,n} \leq I_{t,n}$$

$$x_{t,n} \in \{0,1\}$$

$$(21)$$

Then transformed problem is an Integer Linear Programming (ILP) problem, which is indeed solvable and for which

we design a divide-and-conquer-based algorithm for an efficient settlement, shown in the coming section.

### 5 ALGORITHMS AND ANALYSIS

In this section, we first present the detail of our proposed algorithm, and then some related analysis is given.

## 5.1 Algorithms Design

Algorithm 1 Divide-and-conquer solution for P4

```
Input: The estimated time for model exchange; \{\bar{\tau}_{t,n}\}_{n\in\mathcal{N}} The expected number of chosen arms; m Indicator function of arms' availability; \{I_{t,n}\}_{n\in\mathcal{N}} Length of virtual queue; \{Z_{t,n}\}_{n\in\mathcal{N}} Output:
```

The solution for P4 in round t;  $\{x_{t,n}\}_{n\in\mathcal{N}}$ 

- 1: Set  $Z_t^* = \{Z_{t,n}\}_{I_{t,n}=1}$
- 2: Use  $\mathcal{A}_t$  to store arms with an descending order of  $\boldsymbol{Z}_t^*$
- 3: Use  $\mathcal{N}_t^+$  to store all the n that satisfies  $I_{t,n}=1$
- 4: Set  $k = \min\{m, \sum_{n \in \mathcal{N}} I_{t,n}\}$  // # of clients to be picked

5: for 
$$n_{max} \in \mathcal{N}_t^+$$
 do
6: Initialize an empty set  $\mathcal{S}_{n_{max}}$ 
7: for  $n \in \mathcal{A}_t$  do
8: if  $\bar{\tau}_{t,n} \leq \bar{\tau}_{t,n_{max}}$  then
9: Push  $n$  into  $\mathcal{S}_{n_{max}}$ 
10: end if
11: if  $length(\mathcal{S}_{n_{max}}) == k$  then
12: Calculate the objective of  $P4$  as  $F_{n_{max}}$  based on  $\mathcal{S}_{n_{max}}$ 
13: Break the first loop
14: end if
15: end for
16: end for

- 17: Set  $n^*$  the index of minimum  $F_{n_{max}}$  among those being calculated in line 12.
- 18: Return  $\{x_{t,n}\}$  that represented by  $S_{n*}$

Noticeably, the first term on the objective function of P4 has only finite possible values, so we can simply iterate these values and transform them into the constraint in the sub-problems. By this means, we divide the problem into a few smaller-scale sub-problems, which are easier to conquer. Formally, the sub-problem after division is shown in the following:

$$(P4\text{-}SUB): \min_{\mathbf{x}_{t}} \quad -\sum_{n \in \mathcal{N}} Z_{t,n} x_{t,n}$$
 
$$s.t. \quad \sum_{n \in \mathcal{N}} x_{t,n} = \min \left\{ m, \sum_{n \in \mathcal{N}} I_{t,n} \right\}$$
 
$$x_{t,n} \bar{\tau}_{t,n} \leq \bar{\tau}_{max}$$
 
$$x_{t,n} \leq I_{t,n}$$
 
$$x_{t,n} \in \{0,1\}$$
 
$$(22)$$

where  $\bar{\tau}_{max}$  is one of the fixed value among the possible values of the first term in P4. P4-SUB is much easier to conquer. First we only need to filter those qualified clients with a smaller or equal  $\bar{\tau}_{t,n}$  to  $\bar{\tau}_{max}$ , and with an active

status (or to say  $I_{t,n}=1$ ). Trivially, the sub-problem can be solved by finding  $k=\min\{m,\sum_{n\in\mathcal{N}}I_{t,n}\}$  clients with the biggest  $Z_{t,n}$  among the qualified clients. After the divide-and-conquer process, we only need to compare all the objectives obtained from the sub-problems and select the minimum one as our final achieved solution. The detail of the above process can be found in Algorithm 1, which could at least reach a computation complexity of  $\mathcal{O}(N^2)$ .

## **Algorithm 2** Reputation Based Client Selection with Fairness (RBCS-F)

```
Input:
      The expected number of involved clients each round; m
      Exploration parameter; \alpha_0, \alpha_1, \dots
      The set of clients; \mathcal{N}, Parameter for ridge regression; \lambda
      The guaranteed participating rate; \beta
      Parameter for objective balance; V
Output:
      The control policy \pi = \{x_{t,n}\}_{n \in \mathcal{N}, t=0,1,...}
 1: for n \in \mathcal{N} do
          Initialize \mathbf{H}_{0,n} \leftarrow \lambda \mathbf{I}_{3\times 3}, \mathbf{b}_{0,n} \leftarrow \mathbf{0}_3^{\top}, Z_{0,n} \leftarrow 0
 3: end for
          Observe current contexts \{\mathbf{c}_{t,n}\} and arms availability
          for n \in \mathcal{N} do
 6:
             \hat{oldsymbol{	heta}}_{t,n} \leftarrow \mathbf{H}_{t-1,n}^{-1} \mathbf{b}_{t-1,n}
 7:
             \hat{	au}_{t,n} \leftarrow \mathbf{c}_{t,n}^{\top} \hat{\boldsymbol{\theta}}_{t,n}
 8:
             \bar{\tau}_{t,n} \leftarrow \hat{\tau}_{t,n} \mathbf{v}_{t,n} \\ \bar{\tau}_{t,n} \leftarrow \hat{\tau}_{t,n} - \alpha_t \sqrt{\mathbf{c}_{t,n}^{\top} \mathbf{H}_{t-1,n}^{-1} \mathbf{c}_{t,n}}
 9.
10:
          // Execute Algorithm 1 for a decision
11:
          \{x_{t,n}\} \leftarrow \text{Algorithm } 1(\{\bar{\tau}_{t,n}\}, m, \{I_{t,n}\}, \{Z_{t,n}\})
          Distribute model to the selected clients and observe
12:
          their model exchange time; \{\tau_{t,n}\}
13:
          for n \in \mathcal{N} do
              Update Z_{t,n} according to (6)
14:
              \mathbf{H}_{t,n} \leftarrow \mathbf{H}_{t-1,n} + x_{t,n} \mathbf{c}_{t,n} \mathbf{c}_{t,n}^{\top}
15:
               \mathbf{b}_{t,n} \leftarrow \mathbf{b}_{t-1,n} + x_{t,n} \tau_{t,n} \mathbf{c}_{t,n}
16:
          end for
17:
18: end for
```

With Algorithm 1 introduced, now we are going to discuss our proposed solution for fairness-aware FL, termed Reputation Based Client Selection with Fairness (RBCS-F), shown in Algorithm 2. The working procedure of RBCS-F is quite intuitive. The algorithm starts with initialization of some parameters in the first three lines, and then begins to start iterative federated learning. In every iteration, the scheduler first observes the contexts and the availability of the arms (i.e. FL clients), then estimates the model exchange time with Eqs. (18) and (19) using historical information. Taking advantage of the observed context, availability as well as the estimation, the selection scheme for this round could be fetched by Algorithm 1. After the decision, the model would be distributed to the selected clients and gathered after local training. Before the end of a round, the algorithm records the exchange time of the selected clients and update the associated parameters, as shown in lines 14-16.

#### 5.2 Theoretical Analysis

#### 5.2.1 Regret and Fairness Guarantee

In an MAB model, regret is a key performance metric that measures the performance gap between a given policy and the optimal policy. Therefore, for ease of analysis, we first define the time average regret of RBCS-F.

**Definition 1.** *Time average regret of RBCS-F is defined as:* 

$$R(T) \triangleq \frac{1}{T} \sum_{t=1}^{T} \mathbb{E}\left[f(\mathcal{S}_t, \boldsymbol{\tau}_t) - f(\mathcal{S}_t^*, \boldsymbol{\tau}_t)\right]$$
 (23)

where we leverage  $S_t^*$  to represent the decision made by the optimal policy while  $S_t$  captures RBCS-F's decision.

To proceed, we show a strict upper bound on time average regret of RBCS-F in Theorem 3.

**Theorem 3.** Given any control parameter V, with probability at least  $(1 - \delta)^2$ , the time average regret achieved by RBCS-F is upper bounded by:

$$R(T) \le \frac{N\left(1+\beta^2\right)}{2V} + \zeta_T \sqrt{\frac{6\log(1+TL^2/3\lambda)}{T}}$$
 (24)

where S and L are both positive finite constants satisfying  $\|\boldsymbol{\theta}_n^*\|_2 \leq S$  and  $\|\mathbf{c}_{t,n}\|_2 \leq L$  for all  $t \geq 1$  and  $n \in \mathcal{N}$ . And:

$$\zeta_T = \max\{K, 1\} \cdot \max\left\{2R\sqrt{3\log\left(\frac{1+TL^2/\lambda}{\delta}\right)} + \lambda^{1/2}S, 1\right\}$$
 where  $K$  is a constant value.

*Proof.* The complete proof is given in Appendix C.  $\Box$ 

Now we give another theorem to ensure that the longterm fairness constraint would not be violated.

**Theorem 4.** For RBCS-F, the fairness vitual queues are all mean rate stable in any setting of V, thus the time average fairness is being guaranteed.

*Proof.* The complete proof is given in Appendix D.  $\Box$ 

#### 5.2.2 Impact of V

In light of Theorem 3, it seems quite reasonable for us to set the penalty factor V as large as possible so as to eliminate the first term in the regret upper bound. Such an extreme setting seems even more attractive regarding the fact that the long-term fairness constraint holds under any setting of V, which is justified by Theorem 4. Although a large value of V could indeed bring us a more satisfying long-term model exchange time while satisfying the long-term fairness constraint, we must claim here that the fairness factor is *not* impervious to the setting of *V*. Note that our long-term fairness constraint is built on the premise that the training rounds are infinite, but this may not be true in real training. With a larger V, the fairness queue will have a slower rate to converge, indicating that fairness could not be well guaranteed before convergence. When the training rounds are finite, the number of rounds that need to undergo before convergence could compromise some degrees of fairness. Such an analysis could be verified by our experiment results that we are now going to display.

TABLE 1 Inherent setting of arms (or clients)

client	$ au_n^b$	$ au_n^s$	$\eta$
class		(cold start time)	$\log(1 + SNR)$
1	1s	1s	$\log(1+1000)$
2	2s	1s	$\log(1+100)$
3	3s	1s	$\log(1+10)$
4	4s	1s	$\log(1+1)$

#### 6 EXPERIMENTS

In this section, we present the detail of our experiments. In the first sub-section, we would explain the general setting of our simulation environment and evaluate the numerical performance of our proposed solutions. The numerical evaluation results could well explain the relationship between the penalty factor (V), fairness (reflected by the queue length), and efficiency guarantee (the time span of a federated round). Then we will move on to the evaluation of the real training of two iconic public datasets, CIFAR-10 and fashion-MNIST, both of which are evaluated under different settings of non-iid extent. The real-data experiment will show how our proposed RBCS-F impacts the training efficiency and final model performance (i.e. accuracy).

#### 6.1 Numerical Simulation

#### 6.1.1 Simulation setting

In our simulation, we assume the model exchange time conforms to the linear formulation as shown in Equation (14). To simulate a heterogeneous system with clients of different computation and communication capacity, we equally divide the total number of 40 clients into 4 classes and accordingly endow disparate abilities to them. For clearness, one can check Table 1 for the inherent training setting of different classes of clients.

For the context generation (in order to simulate the perround status of clients), we assume the allocated bandwidth of all clients is sampling from a uniform distribution between [2,4]MHz and the model size M is fixed to 20Mb. Likewise, the available computation capacity of all clients is also sampling from the same uniform distribution within [50%, 200%]. The indicator  $s_{t,n}$  is set according to the training decision in the last round. In addition, for the noise in our linear formulation, we draw  $\epsilon$  from a conditional uniform distribution within  $(-c_{t,n}^{\dagger}\theta_n^*, c_{t,n}^{\dagger}\theta_n^*)$ . The availability of clients follows the same Bernoulli distribution with parameter 0.8, and the setting of other algorithm related parameters could be found in Table 2. In our simulation, we mainly compare RBCS-F with two baseline selection methods that are commonly used in the field, i.e. random and FedCS [1]. Note that we have made an adaption to FedCS in order to accommodate it to our context, but the basic idea is the same as the vanilla one, which is to select as much as clients within a fixed deadline. More concretely, we allow FedCS to have full access to both the contextual features and the static coefficient factor. With the additional information, its strategy is to select all the clients that possess an expected training time (i.e.  $c_{t,n}^{\top} \boldsymbol{\theta}_n^*$ ) that shorter than the pre-set deadline.

TABLE 2 Parameters setting

notation	meaning	value
β	guaranteed participating rate	0.15
$\overline{m}$	maximum selected clients	8
λ	parmeters for ridge regression	1
$\alpha_t$	$\alpha_t$ exploration factor	

#### 6.1.2 Numerical performance evaluation

In our first evaluation, we show the variation of queue status for RBCS-F under different values of penalty factor V. As shown in Fig. 2, where RBCS-F(x) is abbreviated for RBCS-F with a penalty of V = x, it is interesting to see that all the curves with different settings of V flatten after going through a number of scheduling rounds. This phenomenon can justify our conclusion of the mean rate stability of the queues, which indicates that they could not grow to infinity and break our fairness constraint. Another observation we can derive here is that the curve with a higher penalty factor (i.e. V) seems to have a slower convergence speed and a higher convergence value. This implies that a large value of V might sacrifice a few fairness before its convergence, although it does conform to the long-term fairness constraint. Such an observation is consilient with our explanation given in the remark below Theorem 1 and our theoretical analysis in the last section.

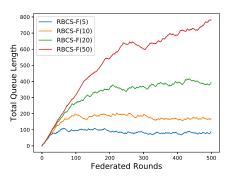


Fig. 2. The impact of V on the convergence of queues

Now we take a look at the evolution of training time across scheduling rounds. In Fig. 3, we depict the time consumption of our proposed RBCS-F with different V, and that of the random strategy and FedCS(3)  $^5$ . As depicted, RBCS-F seems to have a satisfying enhancement in reducing the training time, compared with the random scheme, and of the same number of federated rounds, RBCS-F with a higher V boasts a shorter time consumption. In addition, it is interesting to see that there is a performance gap between RBCS-F and FedCS(3). We note that this gap is inevitable due to our introduction of the fairness factor and the cost of online learning, but the bound itself is well-defined by our analysis of the regret.

The variance in training time of different schemes could be alternatively explained by looking at Fig. 4. This figure

5. For FedCS(3), we set its deadline (one of its key parameter) to 3s. The specific setting allows us to make the number of its selection clients approximates to 8, which is exactly the selection number of other strategies (see our setting in Table 2).

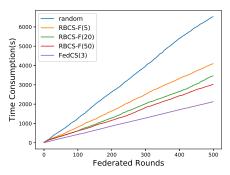


Fig. 3. Training time of different client-selection strategies

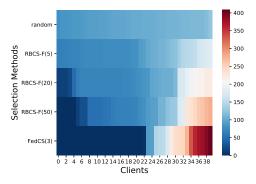


Fig. 4. Pull record of arms (or clients) under different client-selection strategies

depicts the pull number of different arms (or chosen times of FL clients) after going through 500 rounds of decision, in which the clients are sorted in ascending order over their pull number. The brighter color indicates a heavier pull (or more times being selected) on the corresponding arm. From Fig. 4, we notice that the pull number of clients could vary dramatically when V is set to a high value and the unbalanced selection is more intense for FedCS(3). By contrast, the scheme that is known to be fairer (e.g. random or RBCF-F with low penalty) boasts an even distribution on the pull number, based on which we can explain why the training time of RBCS-F would escalate with a fairer selection. Clearly, the selection scheme that evenly chooses the clients shall never match up with those always choosing the fastest ones. However, is it the faster the better? Does fairness matter in real training? Now we are going to explore the answers with our real training on two public datasets.

#### 6.2 Training on Public Dataset

#### 6.2.1 Setup

We set up federated environment with *PyTorch* (version: 1.6.0) and all the computation is conducted using a high-performance workstation (Dell PowerEdge T630 with 2x GTX 1080Ti). We have prepared two tasks for an evaluation purpose. To be specific, we use two different Convolutional Neural Network (CNN) models to predict the classifying results from two datasets, fashion-MNIST, and CIFAR-10. For fashion-MNIST, we adopt a CNN with two 5x5 convolution layers (the first with 20 channels, the second with 50, each followed with 2x2 max pooling), a fully-connected layer with 500 units and ReLu activation, and finally a

softmax output layer. For CIFAR-10, which is known to be a harder task, we use another much heavier CNN model with two 5x5 convolution layers (each with 64 channels), also followed with 2x2 max pooling, two fully connected layers with respective 384 and 192 units, and finally a softmax output layer.

In addition to the general iid setting, we also explore the training performance on a non-iid one. Here we adopt the same approach as in [27] to synthesize non-identical client data. More specifically, we uniformly sample  $q_i\times 500$  items from each of the classifying class, where  $\boldsymbol{q}\triangleq (q_1,q_2,\ldots,q_i)$  is drawn from a Dirichlet distribution, i.e.,  $\boldsymbol{q}\sim \mathrm{Dir}(\gamma_1\boldsymbol{p}).$  Here  $\boldsymbol{p}$  is an all-1 10-dimension vector  $^6$  and  $\gamma_1$  is a concentration parameter controlling the extent of identicalness among clients, say, with  $\gamma_1\to 0$  each client holds only one class chosen at random (i.e. high degree of non-iid), conversely, all clients have identical access to all classes (i.e. approximates to iid) if  $\gamma_1\to\infty$ .

#### 6.2.2 Impact of fairness

In order to quantify the fairness factor and investigate how the factor affects the model accuracy as well as the training efficiency, we thereby introduce  $\gamma_2$  to indicate the extent of fairness. Analogically to how we quantify the non-iid extent, we draw q from a Dirichlet distribution, i.e.,  $q' \sim \text{Dir}(\gamma_2 p')$ and then q' is serving as the probability vector that we use to randomly select clients in each round. As we note before, a smaller value of concentration parameter  $\gamma_2$  leads to a higher variation of q' and thereby causing greater unbalance in selection. Fig. 5 show how the model accuracy evolve with different  $\gamma_2$ , under different non-iid extent (given by  $\gamma_1$ ). Among which, subfigures (a), (b), and (c) depict that of the training for fashionMnist, where we can see that a higher  $\gamma_2$  (a fairer selection) boasts a higher final model accuracy. Also, a similar observation, or an even more conspicuous one, can be found in our training for CIFAR-10, as indicated in subfigures (d), (e), and (f). From our result, it appears that the fairness factor might have different degrees of influence for the training of different datasets. More radically, we are in fact guessing that fairness factor would play a more critical role in a more complicated task. Our theory is that training of a harder task might require more diversified data (in terms of both targets and features), and corresponding, the relative information that each piece of data contains would reduce, and thereby, the training of those tasks should better involve as much available data as possible (i.e. better to be fair), so as to improve the model performance (specifically, final accuracy).

On the other hand, although the experimental data does demonstrate a profound impact of non-iid extent on the model stability and convergence speed during training (as we can observe in Fig. 5 that when  $\gamma_1$  decreases, more jitters on the curve and more rounds underwent before convergence), it does not explicitly show the fairness factor (as reflected by  $\gamma_2$ ) being more or less influential with the change of  $\gamma_1$ , which seems to tell us that our defined non-iid extent has little or no impact on the effect of fairness.

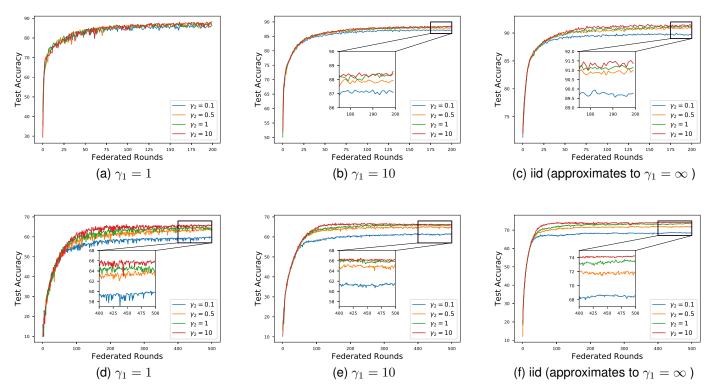


Fig. 5. Fairness impact under fashion-MNIST ( (a), (b) and (c) ) and CIFAR-10 ( (d), (e) and (f) )

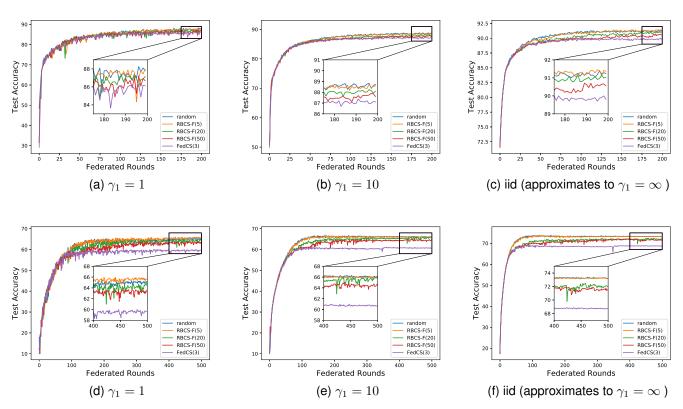


Fig. 6. Accuracy vs. federated rounds for fashion-MNIST ( (a), (b), (c) ) and CIFAR-10 ( (d), (e), (f) )

### 6.2.3 Accuracy vs. federated round

Fig. 6 depicts how our proposed RBCS-F with different settings of V performs in the real training, being compared

6. Both Cifar and fashion-Mnist have 10 targets (or classes)

with the baselines, random and FedCS(3). The result is consistent with our former conjecture that RBCS-F with a smaller V, which is known to be fairer, would achieve a higher final accuracy after rounds of training. Random, a categorically fair scheme, yields the best performance in terms of final model accuracy, while the FedCS(3), another extreme in terms of fairness, does not promise us a commensurate result.

Another point we are interested in is that RBCS-F with a higher penalty seems to spend more rounds to reach a certain accuracy, which we refer to a lower *round efficiency*. This phenomenon can be justified by the result from Fig. 2, which indicates that RBCS-F with a higher penalty tends only to consider fairness when the queue length is large, or in other words, only during a big number of training rounds. Correspondingly, the delay of fairness consideration would make the global model having the chance of aggregating some seldom access data only when the number of rounds is large, and thereby, causing postpone on convergence. Besides, we also found that RBCS-F generally outperforms FedCS(3) in terms of round efficiency, which is conspicuously depicted by subfigures (e) and (f).

## 6.2.4 Accuracy vs. training time

It is also of interest to see how the strategies perform in terms of time efficiency (i.e. the elapsed time to reach a certain accuracy). As such, we take advantage of Fig. 7 to show our investigation on how the model test accuracy evolves according to the elapse of training time. There are a few observations we can derive here. First, a fairer scheme generally achieves a higher final accuracy (or a higher convergence value), which is not a new result as we have already corroborated it in Section 6.2.3. Second, in terms of time efficiency, FedCS(3) achieves an outstanding performance when accuracy is low, which critically outperforms all other schemes, while random, shown to have the worst performance. It is not surprising to see that FedCS(3) could reach the highest time efficiency in the first few rounds as it has complete access to the client-information and has no regard for the fairness factor, but we also see that our proposed algorithm RBCS-F has an acceptable performance gap compared with FedCS(3) and achieves a significant improvement compared with random.

Finally, it is interesting to compare the time efficiency of RBCS-F with different penalties. During our investigation, we see that a higher penalty does not necessarily bring us a higher time efficiency. Specifically, from Figs. 7(e) and 7(f), we see that RBCS-F(20) outperforms RBCS-F(50) during the whole training process. The phenomenon sounds weird as we can derive from Fig. 3 that the average time for each round is strictly decreasing with penalty V. But actually, it can be explained if taking into account the result from Section 6.2.3 that the convergence round would also increase with penalty V. The tradeoff between convergence round and training time of each round is what we believe to cause the unexpected phenomenon.

## 7 CONCLUSION AND FUTURE PROSPECT

In this paper, we have investigated the client selection problem for federated learning. Our concern mainly focuses on the tradeoff between fairness factor and training efficiency. In light of the experiment on our proposed method, we found that fairness is indeed playing a critical role in the training process. In particular, we show that a fairer strategy could promise us a higher final accuracy while inevitably sacrificing a few training efficiency. In terms of how the fairness factor would affect the final achieved accuracy, as well as the convergence speed, however, we could not figure out a rigorous way to quantify their relation. And neither could we track down from the existing literature any theoretical analysis of the fairness factor for FL, making this particular issue quite worthy of investigation. Our future effort would be mainly on this emerging issue.

#### **ACKNOWLEDGMENT**

This work is supported by National Natural Science Foundation of China (Grant Nos. 61872084, 61772205), Guangzhou Science and Technology Program key projects (Grant Nos. 202007040002 and 201902010040, Guangdong Major Project of Basic and Applied Basic Research (2019B030302002).

#### REFERENCES

- [1] T. Nishio and R. Yonetani, "Client selection for federated learning with heterogeneous resources in mobile edge," in ICC 2019-2019 IEEE International Conference on Communications (ICC). IEEE, 2019, pp. 1–7.
- [2] Q. Zeng, Y. Du, K. K. Leung, and K. Huang, "Energy-Efficient Radio Resource Allocation for Federated Edge Learning," arXiv:1907.06040 [cs, math], Jul. 2019, arXiv: 1907.06040. [Online]. Available: http://arxiv.org/abs/1907.06040
- Available: http://arxiv.org/abs/1907.06040
  [3] F. Li, J. Liu, and B. Ji, "Combinatorial sleeping bandits with fairness constraints," *IEEE Transactions on Network Science and Engineering*, 2019.
- [4] S. Deng, H. Zhao, J. Yin, S. Dustdar, and A. Y. Zomaya, "Edge Intelligence: the Confluence of Edge Computing and Artificial Intelligence," arXiv:1909.00560 [cs], Sep. 2019, arXiv: 1909.00560. [Online]. Available: http://arxiv.org/abs/1909.00560
- [5] X. Wang, Y. Han, C. Wang, Q. Zhao, X. Chen, and M. Chen, "In-Edge AI: Intelligentizing Mobile Edge Computing, Caching and Communication by Federated Learning," arXiv:1809.07857 [cs], Sep. 2018, arXiv: 1809.07857. [Online]. Available: http://arxiv.org/abs/1809.07857
- [6] Z. Zhou, X. Chen, E. Li, L. Zeng, K. Luo, and J. Zhang, "Edge Intelligence: Paving the Last Mile of Artificial Intelligence with Edge Computing," arXiv:1905.10083 [cs], May 2019, arXiv: 1905.10083. [Online]. Available: http://arxiv.org/abs/1905.10083
- [7] Y. Kang, J. Hauswald, C. Gao, A. Rovinski, T. Mudge, J. Mars, and L. Tang, "Neurosurgeon: Collaborative intelligence between the cloud and mobile edge," ACM SIGARCH Computer Architecture News, vol. 45, no. 1, pp. 615–629, 2017.
- [8] H. B. McMahan, É. Moore, D. Ramage, S. Hampson et al., "Communication-efficient learning of deep networks from decentralized data," arXiv preprint arXiv:1602.05629, 2016.
- [9] A. Hard, K. Rao, R. Mathews, S. Ramaswamy, F. Beaufays, S. Augenstein, H. Eichner, C. Kiddon, and D. Ramage, "Federated learning for mobile keyboard prediction," arXiv preprint arXiv:1811.03604, 2018.
- [10] S. Ramaswamy, R. Mathews, K. Rao, and F. Beaufays, "Federated learning for emoji prediction in a mobile keyboard," arXiv preprint arXiv:1906.04329, 2019.
- [11] Y. Liu, A. Huang, Y. Luo, H. Huang, Y. Liu, Y. Chen, L. Feng, T. Chen, H. Yu, and Q. Yang, "Fedvision: An online visual object detection platform powered by federated learning," arXiv preprint arXiv:2001.06202, 2020.
- [12] E. Jeong, S. Oh, H. Kim, J. Park, M. Bennis, and S.-L. Kim, "Communication-efficient on-device machine learning: Federated distillation and augmentation under non-iid private data," arXiv preprint arXiv:1811.11479, 2018.
- [13] L. Liu, J. Zhang, S. Song, and K. B. Letaief, "Edge-assisted hierarchical federated learning with non-iid data," arXiv preprint arXiv:1905.06641, 2019.

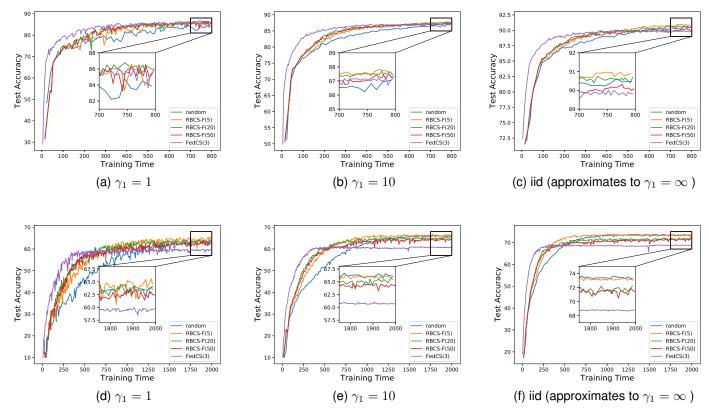


Fig. 7. Accuracy vs. training time for fashion-MNIST ((a), (b), (c)) and CIFAR-10 ((d), (e), (f))

- [14] N. Yoshida, T. Nishio, M. Morikura, K. Yamamoto, and R. Yonetani, "Hybrid-fl: Cooperative learning mechanism using non-iid data in wireless networks," arXiv preprint arXiv:1905.07210, 2019.
- [15] D. Ye, R. Yu, M. Pan, and Z. Han, "Federated learning in vehicular edge computing: A selective model aggregation approach," *IEEE Access*, vol. 8, pp. 23 920–23 935, 2020.
- [16] H. H. Yang, Z. Liu, T. Q. Quek, and H. V. Poor, "Scheduling policies for federated learning in wireless networks," *IEEE Transactions* on Communications, 2019.
- [17] C. Xie, S. Koyejo, and I. Gupta, "Asynchronous federated optimization," arXiv preprint arXiv:1903.03934, 2019.
- [18] W. Wu, L. He, W. Lin, S. Jarvis *et al.*, "Safa: a semi-asynchronous protocol for fast federated learning with low overhead," *arXiv* preprint arXiv:1910.01355, 2019.
- [19] J. Kang, Z. Xiong, D. Niyato, H. Yu, Y.-C. Liang, and D. I. Kim, "Incentive design for efficient federated learning in mobile networks: A contract theory approach," in 2019 IEEE VTS Asia Pacific Wireless Communications Symposium (APWCS). IEEE, 2019, pp. 1–5.
- [20] L. U. Khan, N. H. Tran, S. R. Pandey, W. Saad, Z. Han, M. N. Nguyen, and C. S. Hong, "Federated learning for edge networks: Resource optimization and incentive mechanism," arXiv preprint arXiv:1911.05642, 2019.
- [21] E. Bagdasaryan, A. Veit, Y. Hua, D. Estrin, and V. Shmatikov, "How to backdoor federated learning," arXiv preprint arXiv:1807.00459, 2018
- [22] L. Li, W. Chu, J. Langford, and R. E. Schapire, "A contextual-bandit approach to personalized news article recommendation," in *Proceedings of the 19th international conference on World wide web*, 2010, pp. 661–670.
- [23] Y. Abbasi-Yadkori, D. Pál, and C. Szepesvári, "Improved algorithms for linear stochastic bandits," in Advances in Neural Information Processing Systems, 2011, pp. 2312–2320.
- mation Processing Systems, 2011, pp. 2312–2320.
  [24] L. Qin, S. Chen, and X. Zhu, "Contextual combinatorial bandit and its application on diversified online recommendation," in Proceedings of the 2014 SIAM International Conference on Data Mining. SIAM, 2014, pp. 461–469.
- [25] S. Li, B. Wang, S. Zhang, and W. Chen, "Contextual combinatorial cascading bandits." in *ICML*, vol. 16, 2016, pp. 1245–1253.

- [26] M. J. Neely, "Stochastic network optimization with application to communication and queueing systems," Synthesis Lectures on Communication Networks, vol. 3, no. 1, pp. 1–211, 2010.
- [27] T.-M. H. Hsu, H. Qi, and M. Brown, "Measuring the Effects of Non-Identical Data Distribution for Federated Visual Classification," arXiv:1909.06335 [cs, stat], Sep. 2019, arXiv: 1909.06335. [Online]. Available: http://arxiv.org/abs/1909.06335

## APPENDIX A PROOF OF THEOREM 2

*Proof.* Taking square of (6), we have

$$Z_{t+1,n}^{2} = Z_{t,n}^{2} + 2Z_{t,n} (\beta - x_{t,n}) + (\beta - x_{t,n})^{2}$$
 (25)

Then the difference between  $\frac{1}{2}Z_{t+1,n}^2$  and  $\frac{1}{2}Z_{t,n}^2$  becomes:

$$\frac{1}{2} \left( Z_{t+1,n}^2 - Z_{t,n}^2 \right) 
= \frac{1}{2} \left( \beta - x_{t,n} \right)^2 + Z_{t,n} (\beta - x_{t,n}) 
\leq \frac{1}{2} \left( x_{t,n}^2 + \beta^2 \right) + Z_{t,n} (\beta - x_{t,n}) 
\stackrel{(a)}{\leq} \frac{1}{2} (1 + \beta^2) + Z_{t,n} (\beta - x_{t,n})$$
(26)

Among which, (a) is valid since  $x_{t,n} \in \{0,1\}$ , trivially we have  $x_{t,n}^2 < 1$ .

Now combining (8), (9) and (26), it yields:

$$\Delta(\Theta(t)) \le \Gamma + \sum_{n \in \mathcal{N}} Z_{t,n} \mathbb{E}[\beta - x_{t,n} | \mathbf{\Theta}(t)]$$
 (27)

where  $\Gamma = N \left(1 + \beta^2\right)/2$ .

Plugging  $V\mathbb{E}[\hat{f}(\mathcal{S}_t, \boldsymbol{\tau}_t)|\Theta(t))]$  into (27), it can smoothly transform to the form in (11). This completes the proof.

# APPENDIX B PROOF OF LEMMA 1

We first give Lemma 2 as a fundamental preliminary, and then we will show our justification of Lemma 1.

**Lemma 2** ([23], Theorem 2). Assume that  $\epsilon_t$  is conditionally R-sub-Gaussian for  $R \geq 0$ ,  $\|\boldsymbol{\theta}_n^*\|_2 \leq S$  and  $\|\mathbf{c}_{t,n}\|_2 \leq L$  for all  $t \geq 1$  and  $n \in \mathcal{N}$ . Here S and L are both positive finite constants. Define  $\mathbf{H}_{T,n} = \mathbf{H} + \sum_{t=1}^{T} x_{t,n} \mathbf{c}_{t,n} \mathbf{c}_{t,n}^{\mathsf{T}}$  and set  $\mathbf{H} = \lambda \mathbf{I}$ . Then, with probability at least  $1 - \delta$ , for all rounds  $t \geq 1$ ,  $\hat{\boldsymbol{\theta}}_{t,n}$  satisfies:

$$\left\|\hat{\boldsymbol{\theta}}_{t,n} - \boldsymbol{\theta}_n^*\right\|_{\mathbf{H}_{t-1,n}} \le R\sqrt{3\log\left(\frac{1 + tL^2/\lambda}{\delta}\right)} + \lambda^{1/2}S \tag{28}$$

where we denote  $\|\mathbf{b}\|_{\mathbf{M}} \triangleq \sqrt{\mathbf{b}^T \mathbf{M} \mathbf{b}}$ . **b** is a vector and **M** is a positive definite matrix.

The proof is ommited here for brevity. We refer to [23] for a more dedicated justification.

*Remark.* Recall that  $\epsilon_t$  in our formulation is assumed to be an stochastic variable following a conditionally R-sub-Gaussian, the first assumption trivially holds. In addition, as we have assumed that there exist concrete maximum bounds for  $\tau_n^b$ ,  $\tau_n^s$  and  $1/\eta$ , we can always find a positive  $S < \infty$  making  $\|\boldsymbol{\theta}_n^*\|_2 \leq S$ . Also, we can ensure  $\|\mathbf{c}_{t,n}\|_2 \leq L$  due to the finite volume of contexts. The above analysis justifies the applicability of Lemma 2 to our system model.

Now we are going to show the proof of Lemma 1.

*Proof.* From our estimation rule shown in (19) we can derive:

$$\tau_{t,n}^* - \bar{\tau}_{t,n}$$

$$= \mathbf{c}_{t,n}^{\top} \boldsymbol{\theta}_n^* - \max \left\{ 0, \mathbf{c}_{t,n}^{\top} \hat{\boldsymbol{\theta}}_{t,n} - \alpha_t \| \mathbf{c}_{t,n} \|_{\mathbf{H}_{t-1,n}^{-1}} \right\}$$

$$= \min \left\{ \mathbf{c}_{t,n}^{\top} \boldsymbol{\theta}_n^*, \quad \mathbf{c}_{t,n}^{\top} \boldsymbol{\theta}_n^* - (\mathbf{c}_{t,n}^{\top} \hat{\boldsymbol{\theta}}_{t,n} - \alpha_t \| \mathbf{c}_{t,n} \|_{\mathbf{H}_{t-1,n}^{-1}}) \right\}$$
(29)

Now we focus on the second term in the min function.

$$\mathbf{c}_{t,n}^{\top} \boldsymbol{\theta}_{n}^{*} - (\mathbf{c}_{t,n}^{\top} \hat{\boldsymbol{\theta}}_{t,n} - \alpha_{t} \| \mathbf{c}_{t,n} \|_{\mathbf{H}_{t-1,n}^{-1}})$$

$$= \mathbf{c}_{t,n}^{\top} \left( \boldsymbol{\theta}_{n}^{*} - \hat{\boldsymbol{\theta}}_{t,n} \right) + \alpha_{t} \| \mathbf{c}_{t,n} \|_{\mathbf{H}_{t-1,n}^{-1}}$$

$$\geq - \left\| \boldsymbol{\theta}_{n}^{*} - \hat{\boldsymbol{\theta}}_{t,n} \right\|_{\mathbf{H}_{t-1,n}} \| \mathbf{c}_{t,n} \|_{\mathbf{H}_{t-1,n}^{-1}} + \alpha_{t} \| \mathbf{c}_{t,n} \|_{\mathbf{H}_{t-1,n}^{-1}}$$

$$\geq - \alpha_{t} \| \mathbf{c}_{t,n} \|_{\mathbf{H}_{t-1,n}^{-1}} + \alpha_{t} \| \mathbf{c}_{t,n} \|_{\mathbf{H}_{t-1,n}^{-1}}$$

$$= 0$$
(30)

where (a) can be derived from Lemma 2. Combining the result with the fact  $\mathbf{c}_{t,n}^{\mathsf{T}}\boldsymbol{\theta}_{n}^{*}>0$  and plugging it into (29) yields  $\tau_{t,n}^{*}-\bar{\tau}_{t,n}\geq0$ .

From a different perspective, we have:

$$\tau_{t,n}^{*} - \bar{\tau}_{t,n}$$

$$= \mathbf{c}_{t,n}^{\mathsf{T}} \boldsymbol{\theta}_{n}^{*} - \max \left\{ 0, \mathbf{c}_{t,n}^{\mathsf{T}} \hat{\boldsymbol{\theta}}_{t,n} - \alpha_{t} \| \mathbf{c}_{t,n} \|_{\mathbf{H}_{t-1,n}^{-1}} \right\}$$

$$\leq \mathbf{c}_{t,n}^{\mathsf{T}} \boldsymbol{\theta}_{n}^{*} - \mathbf{c}_{t,n}^{\mathsf{T}} \hat{\boldsymbol{\theta}}_{t,n} + \alpha_{t} \| \mathbf{c}_{t,n} \|_{\mathbf{H}_{t-1,n}^{-1}}$$

$$\leq \left| \mathbf{c}_{t,n}^{\mathsf{T}} \left( \boldsymbol{\theta}_{n}^{*} - \hat{\boldsymbol{\theta}}_{t,n} \right) \right| + \alpha_{t} \| \mathbf{c}_{t,n} \|_{\mathbf{H}_{t-1,n}^{-1}}$$

$$\stackrel{(a)}{\leq} \left\| \boldsymbol{\theta}_{n}^{*} - \hat{\boldsymbol{\theta}}_{t,n} \right\|_{\mathbf{H}_{t-1,n}} \| \mathbf{c}_{t,n} \|_{\mathbf{H}_{t-1,n}^{-1}} + \alpha_{t} \| \mathbf{c}_{t,n} \|_{\mathbf{H}_{t-1,n}^{-1}}$$

$$\leq 2\alpha_{t} \| \mathbf{c}_{t,n} \|_{\mathbf{H}_{t-1,n}^{-1}}$$

$$(31)$$

where (a) is valid by Cauchy-Schwarz inequality. The above results complete the proof.  $\Box$ 

## APPENDIX C PROOF OF THEOREM 3

*Proof.* The definition of the time average regret gives:

$$R(T) = \frac{1}{T} \sum_{t=1}^{T} \mathbb{E}\left[f(\mathcal{S}_t, \boldsymbol{\tau}_t) - f(\mathcal{S}_t^*, \boldsymbol{\tau}_t)\right]$$

$$= \frac{1}{T} \sum_{t=1}^{T} \mathbb{E}\left[\max_{n \in \mathcal{S}_t} \{\tau_{t,n}\} - \max_{n \in \mathcal{S}_t^*} \{\tau_{t,n}\}\right]$$
(32)

Let  $\Delta R(t) = \mathbb{E}\left[\max_{n \in \mathcal{S}_t} \{\tau_{t,n}\} - \max_{n \in \mathcal{S}_t^*} \{\tau_{t,n}\}\right]$  capture the estimated rewards gap between optimal policy and the policy obtained by RBCS-F.

Taking expectation of (27) with respect to  $\Theta(t)$ , it yields:

$$\mathbb{E}\left[\mathcal{L}(\boldsymbol{\Theta}(t+1)) - \mathcal{L}(\boldsymbol{\Theta}(t))\right] \leq \Gamma + \sum_{n \in \mathcal{N}} \mathbb{E}[Z_{t,n}(\beta - x_{t,n})]$$

(33)

Combining the definition of  $\Delta R(t)$  and (33), we have

$$\mathbb{E}\left[\mathcal{L}(\mathbf{\Theta}(t+1)) - \mathcal{L}(\mathbf{\Theta}(t))\right] + V\Delta R(t)$$

$$\leq \Gamma + \sum_{n \in \mathcal{N}} \mathbb{E}[Z_{t,n}(\beta - x_{t,n})] + V\mathbb{E}\left[\max_{n \in \mathcal{S}_t} \{\tau_{t,n}\} - \max_{n \in \mathcal{S}_t^*} \{\tau_{t,n}\}\right]$$

$$= \Gamma + \mathbb{E}\left[V\max_{n \in \mathcal{S}_t} \{\tau_{t,n}\} - \sum_{n \in \mathcal{S}_t^*} Z_{t,n}\right]$$

$$- \mathbb{E}\left[V\max_{n \in \mathcal{S}_t^*} \{\tau_{t,n}\} - \sum_{n \in \mathcal{S}_t^*} Z_{t,n}\right] + \sum_{n \in \mathcal{N}} \mathbb{E}\left[Z_{t,n}(\beta - x_{t,n}^*)\right]$$

$$\stackrel{(a)}{\leq} \Gamma + \mathbb{E}\left[V\max_{n \in \mathcal{S}_t} \{\tau_{t,n}\} - \sum_{n \in \mathcal{S}_t^*} Z_{t,n}\right]$$

$$- \mathbb{E}\left[V\max_{n \in \mathcal{S}_t^*} \{\tau_{t,n}\} - \sum_{n \in \mathcal{S}_t^*} Z_{t,n}\right]$$

$$(34)$$

where (a) is valid since  $\mathbb{E}\left[Z_{t,n}(\beta-x_{t,n}^*)\right] \leq 0$ . This is trivially true because the optimal policy has to ensure the mean rate stability of the fairness queue, which implies that its expected input rate must be smaller than its service rate.

Summing (34) over  $t \in \{1, 2, ..., T\}$  for some T > 0 using telescope yields:

$$\mathbb{E}[\mathcal{L}(\mathbf{\Theta}(T))] - \mathbb{E}[\mathcal{L}(\mathbf{\Theta}(0))] + V \sum_{t=1}^{T} \Delta R(t)$$

$$\leq T\Gamma + \sum_{t=1}^{T} \mathbb{E}[G_1(t)]$$
(35)

where

$$G_1(t) = \left(V \max_{n \in \mathcal{S}_t} \{\tau_{t,n}\} - \sum_{n \in \mathcal{S}_t} Z_{t,n}\right) - \left(V \max_{n \in \mathcal{S}_t^*} \{\tau_{t,n}\} - \sum_{n \in \mathcal{S}_t^*} Z_{t,n}\right)$$
(36)

Recall that  $\mathbb{E}[L(\Theta(T))] \ge 0$  according to its definition, then we can reformulate (35) into the following form:

$$\frac{1}{T} \sum_{t=1}^{T} \Delta R(t) \le \frac{\Gamma}{V} + \sum_{t=1}^{T} \frac{\mathbb{E}\left[G_1(t)\right]}{TV} + \frac{\mathbb{E}\left[L(\boldsymbol{\Theta}(0))\right]}{TV}$$
(37)

Since  $\mathbb{E}[L(\mathbf{\Theta}(0))] = 0$ , we have

$$R(T) = \frac{1}{T} \sum_{t=1}^{T} \Delta R(t) \le \frac{\Gamma}{V} + \sum_{t=1}^{T} \frac{\mathbb{E}[G_1(t)]}{TV}$$
 (38)

Then we proceed to bound  $G_1(t)$ . Assume a policy  $\pi'$  makes her decision  $S'_t$  following such a manner:

$$S'_{t} = \underset{S'_{t} \in \mathcal{C}_{t}}{\operatorname{argmin}} \quad V \max_{n \in S'_{t}} \{ \tau_{t,n} \} - \sum_{n \in \mathcal{S}'_{t}} Z_{t,n}$$
(39)

Here  $C_t$  captures all the possible solutions confined by constraints (explicitly, constraints in P4). The only difference between policy  $\pi'$  and our proposed one is that  $\pi'$  can be fully aware of the real model exchange time. But recall that our algorithm has minimized the objective function in P4,

trivially we know that:

$$V \max_{n \in \mathcal{S}_t} \{\bar{\tau}_{t,n}\} - \sum_{n \in \mathcal{S}_t} Z_{t,n} \le V \max_{n \in \mathcal{S}_t'} \{\bar{\tau}_{t,n}\} - \sum_{n \in \mathcal{S}_t'} Z_{t,n}$$
(40)

Likewise, not a single policy, including the optimal policy  $\pi^*$ , can outperform  $\pi'$  in her objective as shown in (39), trivially we can derive:

$$V \max_{n \in \mathcal{S}'_t} \{ \tau_{t,n} \} - \sum_{n \in \mathcal{S}'_t} Z_{t,n} \le V \max_{n \in \mathcal{S}^*_t} \{ \tau_{t,n} \} - \sum_{n \in \mathcal{S}^*_t} Z_{t,n}$$
 (41)

Therefore, according to (40) and (41), it follows that:

$$G_{1}(t)$$

$$= \left(V \max_{n \in \mathcal{S}_{t}} \{\tau_{t,n}\} - \sum_{n \in \mathcal{S}_{t}} Z_{t,n}\right) - \left(V \max_{n \in \mathcal{S}_{t}^{*}} \{\tau_{t,n}\} - \sum_{n \in \mathcal{S}_{t}^{*}} Z_{t,n}\right)$$

$$\leq \left(V \max_{n \in \mathcal{S}_{t}} \{\tau_{t,n}\} - \sum_{n \in \mathcal{S}_{t}} Z_{t,n}\right) - \left(V \max_{n \in \mathcal{S}_{t}^{'}} \{\tau_{t,n}\} - \sum_{n \in \mathcal{S}_{t}^{'}} Z_{t,n}\right)$$

$$\leq \left(V \max_{n \in \mathcal{S}_{t}} \{\tau_{t,n}\} - \sum_{n \in \mathcal{S}_{t}^{'}} Z_{t,n}\right) - \left(V \max_{n \in \mathcal{S}_{t}^{'}} \{\tau_{t,n}\} - \sum_{n \in \mathcal{S}_{t}^{'}} Z_{t,n}\right)$$

$$+ \left(V \max_{n \in \mathcal{S}_{t}^{'}} \{\bar{\tau}_{t,n}\} - \sum_{n \in \mathcal{S}_{t}^{'}} Z_{t,n}\right) - \left(V \max_{n \in \mathcal{S}_{t}^{'}} \{\bar{\tau}_{t,n}\} - \sum_{n \in \mathcal{S}_{t}^{'}} Z_{t,n}\right)$$

$$\leq V \left[\left(\max_{n \in \mathcal{S}_{t}} \{\tau_{t,n}\} - \max_{n \in \mathcal{S}_{t}^{'}} \{\bar{\tau}_{t,n}\}\right) + \left(\max_{n \in \mathcal{S}_{t}^{'}} \{\bar{\tau}_{t,n}\} - \max_{n \in \mathcal{S}_{t}^{'}} \{\tau_{t,n}\}\right)\right]$$

$$\leq V \left[\left(\tau_{t,s} - \bar{\tau}_{t,s}\right) + \left(\bar{\tau}_{t,s'} - \tau_{t,s'}\right)\right]$$

$$(42)$$

where arm s is the arm that get the maximum  $\tau_{t,n}$  among the set  $n \in \mathcal{S}_t$ . Likewise, s' is another arm that gets the maximum  $\bar{\tau}_{t,n}$  among the set  $n \in \mathcal{S}'_t$ . Notice that both  $\bar{\tau}_{t,s}$  and  $\bar{\tau}_{t,s'}$  are not coupled with the stochastic value  $\tau_{t,s}$  and we also have  $\mathbb{E}[\tau_{t,s}] = \tau^*_{t,s'}$ ,  $\mathbb{E}[\tau_{t,s'}] = \tau^*_{t,s'}$ . Subsequently, we have:

$$\mathbb{E}[G_1(t)] \le V \left[ \underbrace{(\tau_{t,s}^* - \bar{\tau}_{t,s})}_{J_1(t)} + \underbrace{(\bar{\tau}_{t,s'} - \tau_{t,s'}^*)}_{J_2(t)} \right]$$
(43)

Combining Lemma 1, we can conclude that with probability at least  $(1-\delta)^2$ ,  $J_1(t) \leq 2\alpha_t \|\mathbf{c}_{t,s}\|_{\mathbf{H}_{t-1,s}^{-1}}$  and  $J_2(t) \leq 0$ . Now we know that  $\mathbb{E}[G_1(t)] \leq VJ_1(t)$  with large probability.

Note that we have a concret bound for  $\tau_{t,s}^*$  since elements of  $c_{t,n}$  and  $\theta_n^*$  are all bounded. Now assume that we have  $\tau_{t,s}^* \leq K$ , then combining  $\bar{\tau}_{t,s} \geq 0$ , it follows that  $J_1(t) \leq K$ . Consequently, we have:

$$\begin{split} J_{1}(t) &\leq \min \left\{ 2\alpha_{t} \, \| \mathbf{c}_{t,s} \|_{\mathbf{H}_{t-1,s}^{-1}}, K \right\} \\ &\leq \max\{K,1\} \min \left\{ 2\alpha_{t} \, \| \mathbf{c}_{t,s} \|_{\mathbf{H}_{t-1,s}^{-1}}, 1 \right\} \\ &\leq \max\{K,1\} \cdot \max\{2\alpha_{t},1\} \cdot \min \left\{ \| \mathbf{c}_{t,s} \|_{\mathbf{H}_{t-1,s}^{-1}}, 1 \right\} \end{split}$$

where the last two inequalities hold by  $\min\{a,b\} \leq \max\{b,1\} \cdot \min\{a,1\}$  and  $\min\{ab,1\} \leq \max\{a,1\} \cdot \min\{b,1\}$ , respectively. Let  $\zeta_t = \max\{K,1\} \cdot \max\{2\alpha_t,1\}$ . Now  $\mathbb{E}[G_1(t)]$  has been successfully bounded into a closed form:  $\mathbb{E}[G_1(t)] \leq V\zeta_t \min\left\{\|\mathbf{c}_{t,s}\|_{\mathbf{H}_{t-1,s}^{-1}},1\right\}$ .

Then we continue our proof by given  $\Upsilon(t)$  $\frac{1}{T}\sum_{t=1}^{T}\zeta_{t}\min\left\{\left\|\mathbf{c}_{t,s}\right\|_{\mathbf{H}_{t-1,s}^{-1}},1\right\}$ , which is the upper bound on the second term of (38) that we have derived so far. Before our further bounding on  $\Upsilon(t)$ , we first familiarize the readers with Lemma 3.

**Lemma 3.** Assume  $\mathbf{H}_{T,n} = \mathbf{H} + \sum_{t=1}^{T} x_{t,n} \mathbf{c}_{t,n} \mathbf{c}_{t,n}^{\top}$  and  $\mathbf{H} = \lambda \mathbf{I}$ . Then, if  $\lambda \geq 1$  and  $\|\mathbf{c}_{t,n}\|_2 \leq L$  for all t and n, we have

$$\sum_{t=1}^{T} \min \left\{ \left\| \mathbf{c}_{t,n} \right\|_{H_{t-1,n}^{-1}}^{2}, 1 \right\} \\
\leq 2 \left( \log \det \left( \mathbf{H}_{T,n} \right) - \log \det (\mathbf{H}) \right) \\
\leq 6 \log \left( \left( \operatorname{trace}(\mathbf{H}) + TL^{2} \right) / 3 \right) - 2 \log \det (\mathbf{H}) \tag{45}$$

for any n and T.

The proof is omitted here for briefness. For more detail, we refer the interested readers to Lemma 11 in [23].

With this lemma introduced, now we shall introduce our further deduction, as presented in the following:

$$\Upsilon(t) = \frac{1}{T} \sum_{t=1}^{T} \zeta_{t} \min \left\{ \| \mathbf{c}_{t,s} \|_{\mathbf{H}_{t-1,s}^{-1}}, 1 \right\}$$

$$\stackrel{(a)}{\leq} \sqrt{\frac{\sum_{t=1}^{T} \zeta_{t}^{2} \min \left\{ \| \mathbf{c}_{t,s} \|_{\mathbf{H}_{t-1,s}^{-1}}^{2}, 1 \right\}}{T}}$$

$$\leq \zeta_{T} \sqrt{\frac{\sum_{t=1}^{T} \min \left\{ \| \mathbf{c}_{t,s} \|_{\mathbf{H}_{t-1,s}^{-1}}^{2}, 1 \right\}}{T}}$$

$$\stackrel{(b)}{\leq} \zeta_{T} \sqrt{\frac{6 \log((\operatorname{trace}(\mathbf{H}) + TL^{2})/3) - 2 \log \det(\mathbf{H})}{T}}$$

$$\stackrel{(c)}{=} \zeta_{T} \sqrt{\frac{6 \log(\lambda + TL^{2}/3) - 6 \log \lambda}{T}}$$

$$= \zeta_{T} \sqrt{\frac{6 \log(1 + TL^{2}/3\lambda)}{T}}$$
(46)

where (a) can be justified by arithmetic means inequality and Theorem 3 gives (b). Using the facts trace( $\mathbf{H}$ ) =  $3\lambda$  and  $\det(\mathbf{H}) = \lambda^3$  we have (c).

Plugging the above result into (37) and combining the fact that  $\Gamma = N(1 + \beta^2)/2$ , finally we reach a high probability (i.e. with probability  $(1 - \delta)^2$ ) upper bound on the regret, as presented in the following:

$$R(T) \le \frac{N\left(1 + \beta^2\right)}{2V} + \zeta_T \sqrt{\frac{6\log(1 + TL^2/3\lambda)}{T}} \tag{47}$$

where 
$$\zeta_T = \max\{K,1\} \cdot \max\left\{2R\sqrt{3\log\left(\frac{1+TL^2/\lambda}{\delta}\right)} + \lambda^{1/2}S,1\right\}$$
 This completes the proof.  $\Box$ 

## APPENDIX D **PROOF OF THEOREM 4**

Here we first prepare the readers with Lemma 4, which we will use in our proof of Theorem 4.

**Lemma 4.** For any  $\delta' > 0$ , there exists an  $\omega$ -only policy  $\pi$ , which makes independent, stationary, and randomized decisions in every round t based only on the observed stochastic events  $\omega(t)$ , gives the following results:

$$\mathbb{E}\left[V \max_{n \in \mathcal{N}} \{x_{t,n}^{\pi} \bar{\tau}_{t,n}\} | \mathbf{\Theta}(t)\right] = \mathbb{E}\left[V \max_{n \in \mathcal{N}} \{x_{t,n}^{\pi} \bar{\tau}_{t,n}\}\right]$$

$$\leq V \max_{n \in \mathcal{N}} \{x_{t,n}^{opt} \bar{\tau}_{t,n}\} + \delta'$$

$$\mathbb{E}\left\{\beta - x_{t,n}^{\pi} | \mathbf{\Theta}(t)\right\} = \mathbb{E}\left\{\beta - x_{t,n}^{\pi}\right\} \leq \delta' \quad \forall n \in \mathcal{N}$$
(48)

where the expectations here are with respect to random actions achieved by the policy and the stochastic event  $\omega(t)$ .  $x_{t,n}^{opt}$  is the optimal policy that minimizes over  $\max_{n \in \mathcal{N}} \{x_{t,n} \bar{\tau}_{t,n}\}$  while meeting the long-term fairness constraint.

The proof of Lemma 4 is omitted here, we refer the readers to Theorem 4.5 in [26].

Formally, now we begin our justification of Theorem 4.

*Proof.* Since RBCS-F minimizes the R.H.S of inequation (11), straightforwardly not other policies could match up with RBCS-F in its objective, namely, it gives:

$$\sum_{n \in \mathcal{N}} Z_{t,n}(\beta - x_{t,n}) + V \max_{n \in \mathcal{N}} \{x_{t,n} \bar{\tau}_{t,n}\}$$

$$\leq \sum_{n \in \mathcal{N}} Z_{t,n}(\beta - x_{t,n}^{\pi}) + V \max_{n \in \mathcal{N}} \{x_{t,n}^{\pi} \bar{\tau}_{t,n}\}$$
(49)

The above inequality is valid for any  $\omega$ -only policy  $\pi$  and under any  $\omega(t)$  and  $\Theta(t)$ . Therefore we have:

$$\mathbb{E}\left[\sum_{n\in\mathcal{N}} Z_{t,n}(\beta - x_{t,n}) + V \max_{n\in\mathcal{N}} \{x_{t,n}\bar{\tau}_{t,n}\} | \mathbf{\Theta}(t)\right]$$

$$\leq \mathbb{E}\left[\sum_{n\in\mathcal{N}} Z_{t,n}(\beta - x_{t,n}^{\pi}) + V \max_{n\in\mathcal{N}} \{x_{t,n}^{\pi}\bar{\tau}_{t,n}\} | \mathbf{\Theta}(t)\right]$$
(50)

Plugging this into (11), it yields:

$$\Delta(\Theta(t)) + V \mathbb{E}[\max_{n \in \mathcal{N}} \{x_{t,n} \bar{\tau}_{t,n}\} | \Theta(t)]$$

$$\leq \Gamma + \sum_{n \in \mathcal{N}} Z_{t,n} \mathbb{E}[\beta - x_{t,n}^{\pi} | \Theta(t)]$$

$$+ V \mathbb{E}[\max_{n \in \mathcal{N}} \{x_{t,n}^{\pi} \bar{\tau}_{t,n}\} | \Theta(t)]$$

$$\stackrel{(a)}{\leq} \Gamma + \sum_{n \in \mathcal{N}} Z_{t,n} \delta' + V \max_{n \in \mathcal{N}} \{x_{t,n}^{opt} \bar{\tau}_{t,n}\} + \delta'$$
(51)

where inequality (a) could be derived from Lemma 4. Then taking  $\delta' \to 0$  yields:

$$\Delta(\Theta(t)) + V \mathbb{E}[\max_{n \in \mathcal{N}} \{x_{t,n} \bar{\tau}_{t,n}\} | \Theta(t)]$$

$$\leq \Gamma + V \max_{n \in \mathcal{N}} \{x_{t,n}^{opt} \bar{\tau}_{t,n}\}$$
(52)

Taking expectation on the above inequality with repect to  $\Theta(t)$  and then summing them over  $t \in \{1, 2, ..., T\}$  for  $T \to \infty$  using telescope yields:

$$\mathbb{E}[L(\mathbf{\Theta}(T))] - \mathbb{E}[L(\mathbf{\Theta}(0))] + V \sum_{t=1}^{T} \mathbb{E}[\max_{n \in \mathcal{N}} \{x_{t,n} \bar{\tau}_{t,n}\}]$$

$$\leq T\Gamma + V \sum_{t=1}^{T} \max_{n \in \mathcal{N}} \{x_{t,n}^{opt} \bar{\tau}_{t,n}\}$$
(53)

In light of the law of large numbers, we have

$$\sum_{t=1}^{T} \max_{n \in \mathcal{N}} \{ x_{t,n}^{opt} \bar{\tau}_{t,n} \} = \sum_{t=1}^{T} \mathbb{E} [\max_{n \in \mathcal{N}} \{ x_{t,n}^{opt} \bar{\tau}_{t,n} \}]$$
 (54)

In addtion, recall that not a policy can obtain smaller objective than the optimal scheme, therefore, it yields:

$$\mathbb{E}[\max_{n \in \mathcal{N}} \{x_{t,n}^{opt} \bar{\tau}_{t,n}\}] \le \mathbb{E}[\max_{n \in \mathcal{N}} \{x_{t,n} \bar{\tau}_{t,n}\}] \quad \forall t \in \mathcal{T}$$
 (55)

Reranging (53), it gives:

$$\mathbb{E}[L(\mathbf{\Theta}(T))] - \mathbb{E}[L(\mathbf{\Theta}(0))] \le T\Gamma \tag{56}$$

Plugging the definition of  $\mathcal{L}(\mathbf{\Theta}(T))$  in (56) yields:

$$\sum_{n \in \mathcal{N}} \mathbb{E}[Z_{T,n}^2)] \le 2T\Gamma + 2\mathbb{E}[L(\mathbf{\Theta}(0))] \tag{57}$$

Plugging the fact  $\mathbb{E}[Z_{T,n}]^2 \leq \mathbb{E}[Z_{T,n}^2]$  into the above inequality, it yields:

$$\sum_{n \in \mathcal{N}} \mathbb{E}\left[Z_{T,n}\right] \le \sqrt{2T\Gamma + 2\mathbb{E}[L(\mathbf{\Theta}(0))]}$$
 (58)

Dividing by T yields:

$$\lim_{T \to \infty} \sum_{n \in \mathcal{N}} \frac{\mathbb{E}\left[Z_{T,n}\right]}{T} \le \lim_{T \to \infty} \sqrt{\frac{2\Gamma}{T} + \frac{2\mathbb{E}\left[L(\boldsymbol{\Theta}(0))\right]}{T^2}} = 0 \quad (59)$$

Since  $\mathbb{E}[Z_{T,n}] \geq 0$  for any T and n, finally we conclude that

$$\lim_{T \to \infty} \frac{\mathbb{E}\left[Z_{T,n}\right]}{T} = 0 \quad \forall n \in \mathcal{N}$$
 (60)

which explicitly marks the mean rate stability of all queues regardless of the real setting of V. Combining the result from Theorem 1, we can ensure no violation on the long-term average fairness constraint.  $\Box$