
Supplementary Materials:

Asynchronous Decentralized Online Learning

Supplementary materials are organized as follows. Appendix A gives a more detailed review of related work. Appendix B analyzes the feasibility of the iterations of AD-OGP. Appendix C provides more detailed experimental setup and supplements additional empirical results. Appendix D and Appendix E provide detailed proofs of Theorem 1 and Corollary 1, which have been omitted in our main paper due to the space limit.

A More Related Work

In this section, we provide a more detailed review of related work in two fields, namely decentralized online learning and decentralized stochastic optimization.

A.1 Decentralized Online Learning

In many distributed online learning scenarios such as learning on mobile networks or sensor networks, centralized architectures are usually unsuitable due to the very high computation and communication overhead on the central node [7]. As a general solution to reduce such overhead, decentralized online learning has received much research attention in recent years.

As mentioned in the main paper, most existing studies in this field are conducted in the synchronous setting, namely, they require some global coordination among learners during the learning process. Such coordination is used in either feedback processing, local update, or node communication (also see [3] for a more detailed explanation of synchronicity). Many previous works have investigated exemplar online algorithms in the synchronous decentralized setting [36, 13, 15, 24, 38, 16, 18]. Other recent researches have focused on developing synchronous algorithms in specific scenarios, such as learning on dynamic networks [1, 28, 37, 20], more efficient information exchange [32, 39, 33], or optimization under constraints [4].

Different from these previous works, we propose to conduct decentralized online learning in a **fully asynchronous** manner [3]. In particular, our work does not need any global synchronization mechanism to coordinate the learners throughout the learning process. In other words, in our framework, each learner makes predictions, processes feedback, performs local updates, sends out messages, receives messages, and performs model averaging independently; in particular, it does not need to know the states of other learners. As far as we are concerned, the recently proposed [14] is the only work that investigates delay and asynchronization in decentralized online learning. However, their work only allows learners to communicate specific gradients, hence it is restricted to dual averaging method where the decision only depends on historic gradients. By taking a completely different approach, our work allows each learner to communicate its model parameters with other learners, which is more efficient and suitable to more general online update methods.

A.2 Decentralized Stochastic Optimization

Decentralized stochastic optimization aims to minimize the objective function $\sum_{i \in V} f_i/|V|$, where each learner $i \in V$ only has the access to a local loss $f_i(\mathbf{w}) = \mathbb{E}_{\xi \sim \mathcal{D}} F_i(\mathbf{w}; \xi)$. In this setting, each instance point ξ are i.i.d. sampled from some fixed but unknown distribution \mathcal{D} ; this is a core assumption in the stochastic optimization setting. Previous studies in this field can be classified into

two categories: synchronous algorithms and asynchronous algorithms. In synchronous algorithms, at each iteration, all learners compute stochastic gradients locally and are synchronized to update their local models at the same pace [7, 26, 29, 30, 23, 17]. However, as stated before, synchronous algorithms suffer from the straggler problem. To resolve this problem, many recent works have developed asynchronous algorithms for stochastic optimization [21, 8, 35, 12, 2, 31].

Despite the existence of many previous works for asynchronous decentralized stochastic optimization, we note that, as a first systematic study for asynchronous decentralized online learning, our work is largely different from them in the following three aspects. *First*, we contribute the first framework of asynchronous decentralized online learning. In particular, we propose an event indexing system by taking both prediction and local update into account, which is distinct from previous stochastic frameworks using the update indexing system [21] in both formulation and analysis. *Second*, our algorithm has two specific and novel designs in the online setting, namely weighted projection and instantaneous model averaging. Both designs make our algorithm more effective than simply tailoring standard stochastic optimization algorithms [21, 2, 31] to the online setting. *Third*, the regret analysis for decentralized online learning is intrinsically different from the convergence analysis for decentralized stochastic optimization.

B The Feasibility of AD-OGP

In this section, we strictly prove that the predictions generated by our proposed AD-OGP algorithm are guaranteed to lie in the feasible region \mathcal{K} . The proof is not straightforward because the predictions \mathbf{w}_i of learner i are determined by its model \mathbf{x}_i and push-sum weight y_i together, namely $\mathbf{w}_i = \mathbf{x}_i/y_i$; during the learning process, both local update and model averaging can change the parameters \mathbf{x}_i and y_i , which possibly induces predictions that lie outside of the feasible region.

We first introduce some notations to specify each learner i 's parameters \mathbf{x}_i and y_i at some certain event time points $t \in \{1, \dots, 2T\}$. Specifically, we use $\mathbf{x}_i(t)$ and $y_i(t)$ to denote its most recent local model and push-sum weight *before* event t ; and use $\mathbf{x}'_i(t)$ and $y'_i(t)$ to denote its local model and push-sum weight *immediately after* event t . Then its predictions immediately before and after event t are $\mathbf{w}_i(t) = \mathbf{x}_i(t)/y_i(t)$ and $\mathbf{u}_i(t) = \mathbf{x}'_i(t)/y'_i(t)$, respectively. To help better understand our algorithm, we describe the full procedure of AD-OGP in Algorithm 1, which is aligned with our proposed AD-OCO framework (Protocol 1 in our main paper).

Proposition 1. *In AD-OGP, for any $i \in V, t \in \mathcal{P}_T$, it holds that*

$$y_i(t) > 0, \quad \mathbf{w}_i(t) = \frac{\mathbf{x}_i(t)}{y_i(t)} \in \mathcal{K}.$$

Proof. As a prerequisite of proving feasibility, we first show that, for each learner $i \in V$, its push-sum weight $y_i(t)$ is always positive; otherwise, the division $\mathbf{x}_i(t)/y_i(t)$ may be invalid.

Lemma 1. *In AD-OGP, for each learner $i \in V$, its push-sum weight is positive throughout the learning process, i.e., $y_i(t) > 0, y'_i(t) > 0, \forall t \in \{1, \dots, 2T\}$.*

Proof. There are only two types of actions that may change the value of push-sum weights, namely local update \mathcal{U} and model averaging \mathcal{A} . Consider an arbitrary learner $i \in V$. At each of its local update event, its push-sum weight y_i is multiplied by a factor of γ_i ($1/m \leq \gamma_i < 1$). Moreover, each time it executes a model averaging operation, its push-sum weight y_i is added with some weight y^j of another learner. In fact, the positiveness of $y_i(t)$ can be proved by inducting on $t \in \{1, \dots, 2T\}$.

At the beginning of the learning process ($t = 1$), for each learner $i \in V$, its push-sum weight is initialized as $y_i(1) = 1$. Since the first event must be a prediction event, we must have $y'_i(1) = y_i(1) = 1$ for any $i \in V$. Now we suppose that, for some $t \leq 2T$, we have $y_i(l) > 0, y'_i(l) > 0, \forall i \in V, \forall l \in \{1, \dots, t\}$. We now check the positiveness of $y_i(t+1)$ and $y'_i(t+1)$.

(i) During the interval between events t and $t+1$, y_i can only change during model averaging on learner i . Recall the notations in our framework, namely, learner i performs model averaging using the copies in $\mathcal{M}_i(t)$ during this interval. If $\mathcal{M}_i(t) = \emptyset$, we directly have $y_i(t+1) = y'_i(t) > 0$. If $\mathcal{M}_i(t) \neq \emptyset$, consider any single copy $(\mathbf{x}', y') = (\gamma_j \mathbf{x}'_j, \gamma_j y_j) \in \mathcal{M}_i(t)$. We suppose that its corresponding message delay is d' . Then the weight y'_j is exactly learner j 's push-sum weight

Algorithm 1 Asynchronous Decentralized Online Gradient-Push (**AD-OGP**)

1: **Input:** Time horizon T , convex set \mathcal{K} , and learning rate η .
2: **Initialize:** Local model $\mathbf{x}_j(1) \leftarrow \mathbf{w}_0 \in \mathcal{K}$ and push-sum weight $y_j(1) \leftarrow 1, \forall j \in V$.
3: **for** $t = 1, \dots, 2T$ **do**
4: **if** $\delta_t = 0$ **then** *// for a prediction event*
5: Learner i_t predicts with $\mathbf{w}_{i_t}(t) = \mathbf{x}_{i_t}(t)/y_{i_t}(t)$.
6: Learner i_t starts to compute the gradient $\nabla f_t(\mathbf{w}_{i_t}(t))$.
7: (Parameters of *all* learners $j \in V$ stay unchanged: $\mathbf{x}'_j(t) = \mathbf{x}_j(t), y'_j(t) = y_j(t)$.)
8: **else** *// for a local update event*
9: Learner i_t performs weighted projected gradient descent:

$$\mathbf{x}'_{i_t}(t) = \gamma_{i_t} y_{i_t}(t) \Pi_{\mathcal{K}}\left(\frac{\mathbf{x}_{i_t}(t) - \eta \nabla f_t(\mathbf{w}_{i_t}(t))}{y_{i_t}(t)}\right), \quad y'_{i_t}(t) = \gamma_{i_t} y_{i_t}(t).$$

10: Learner i_t sends $\mathbf{x}'_{i_t}(t)$ and $y'_{i_t}(t)$ to each neighbor $j \in \mathcal{N}(i_t)$.
11: (Parameters of *other* learners $j \neq i_t$ stay unchanged: $\mathbf{x}'_j(t) = \mathbf{x}_j(t), y'_j(t) = y_j(t)$.)
12: **end if**
13: **for** $i \in V$ **do**
14: Learner i performs model averaging as long as it is *not* occupied:

$$\mathbf{x}_i(t+1) = \mathbf{x}'_i(t) + \sum_{\mathbf{x}' \in \mathcal{M}_i(t)} \mathbf{x}', \quad y_i(t+1) = y'_i(t) + \sum_{y' \in \mathcal{M}_i(t)} y'.$$

15: **end for**
16: **end for**

immediately after the event $t - d'$, namely $y'_j(t - d')$. Since we have assumed that $y'_j(t - d') > 0$, we thus have $y_i(t+1) = y'_i(t) + \sum_{(\mathbf{x}', y') \in \mathcal{M}_i(t)} y' > 0$.

(ii) At the event $t + 1$, from Algorithm 1 we know that either $y'_i(t+1) = y_i(t+1)$ or $y'_i(t+1) = \gamma_i y_i(t+1)$. Since $\gamma_i > 0$, in both cases we have $y'_i(t+1) > 0$.

In summary, we have $y_i(t+1) > 0$ and $y'_i(t+1) > 0$ for any $i \in V$. Hence we prove this lemma. \square

Now we can formally analyze the feasibility of AD-OGP. In intuition, we expect that both operations of local update and model averaging will not violate the feasibility of the generated predictions. Recall that, in our algorithm, for each local update, the feasibility is always preserved by our proposed weighted projection mechanism. Hence we only need to prove that, during model averaging, the feasibility of AD-OGP is also preserved. Indeed, such a property is implied by the following lemma.

Lemma 2. *Suppose the parameters (\mathbf{x}_i, y_i) of all learners satisfy $y_i > 0, \mathbf{x}_i/y_i \in \mathcal{K}, \forall i \in V$. Then given any averaging weights $\{\alpha_i\}_{i \in V}$ such that $\alpha_i \geq 0, \forall i \in V$ and $\sum_{i \in V} \alpha_i > 0$, it holds that*

$$\frac{\sum_{i \in V} \alpha_i \mathbf{x}_i}{\sum_{i \in V} \alpha_i y_i} \in \mathcal{K}.$$

Proof. We first consider the simple case of two learners with the same weights $\alpha_1 = \alpha_2 > 0$. In this case, we directly have

$$\frac{x_1 + x_2}{y_1 + y_2} = \frac{y_1}{y_1 + y_2} \cdot \frac{x_1}{y_1} + \frac{y_2}{y_1 + y_2} \cdot \frac{x_2}{y_2},$$

thus $(x_1 + x_2)/(y_1 + y_2)$ is a convex combination of x_1/y_1 and x_2/y_2 . If $x_1/y_1, x_2/y_2 \in \mathcal{K}$, the parameters after model averaging will still generate the prediction $(x_1 + x_2)/(y_1 + y_2)$ that lies in \mathcal{K} .

Now we prove this lemma by inducting on the cardinality m of V . It trivially holds for the single-learner setting where $m = 1$. Now we suppose that it is satisfied under $m \leq k$ for some $k \geq 1$. We now intend to prove that it still holds when $m = k + 1$. Denote $V = \{1, \dots, k + 1\}$. Since $\sum_{i \in V} \alpha_i > 0$, α_i cannot be all zero for $i \in V$. We now consider the following two cases.

(i) There is only a single α_i being positive. Then, without loss of generality, we can assume $\alpha_1 > 0$ and $\alpha_i = 0$ for any $i > 1$. In this case, we trivially have $(\alpha_1 x_1)/(\alpha_1 y_1) \in \mathcal{K}$.

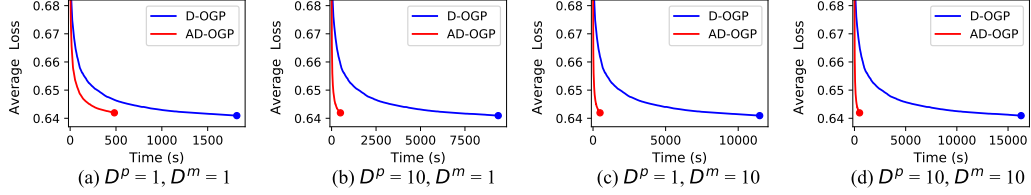


Figure 5: More results to verify the benefit of asynchronization. The plots compare AD-OGP and D-OGP on a 64-node *ring* graph under varying levels of processing delays D^P and message delays D^m on *higgs*.

(ii) There are at least two α_i that are positive. Without loss of generality, we can assume $\alpha_1, \alpha_2 > 0$. Since $\sum_{i=2}^{k+1} \alpha_i > 0$, from our assumption we have $(\sum_{i=2}^{k+1} \alpha_i x_i) / (\sum_{i=2}^{k+1} \alpha_i y_i) \in \mathcal{K}$. Set $x'_1 = \alpha_1 x_1, y'_1 = \alpha_1 y_1, x'_2 = \sum_{i=2}^{n+1} \alpha_i x_i, y'_2 = \sum_{i=2}^{k+1} \alpha_i y_i$, then from the initial case of two learners, we have $(x'_1 + x'_2) / (y'_1 + y'_2) \in \mathcal{K}$.

Combining the above two cases, we check the setting of $|V| = k + 1$. Hence, by inducting on k , we prove the lemma. \square

Finally, we can prove the feasibility of AD-OGP by inducting on $t \in \{1, \dots, 2T\}$. Recall that, at the beginning of the learning process $t = 1$, each learner $i \in V$ is initialized as $\mathbf{x}_i(1) = \mathbf{w}_0 \in \mathcal{K}$ and $y_i(1) = 1$; and $\mathbf{x}'_i(1) = \mathbf{x}_i(1), y'_i(1) = y_i(1)$.

We now suppose that for some $t \leq 2T$, it holds that $\mathbf{x}_i(l)/y_i(l) \in \mathcal{K}, \mathbf{x}'_i(l)/y'_i(l) \in \mathcal{K}, \forall i \in V, l \in \{1, \dots, t\}$. Then during the interval between events t and $t + 1$, each learner $i \in V$ may perform model averaging using the copies in $\mathcal{M}_i(t)$. From Lemma 2 we know that $\mathbf{x}_i(t+1)/y_i(t+1) \in \mathcal{K}$. In addition, recall our designed mechanism of weighted projection, we also have $\mathbf{x}'_i(t+1)/y'_i(t+1) \in \mathcal{K}$. Consequently, we prove the proposition by induction. \square

C More Details of Experiments

In this section, we first provide more details of our experimental setup and algorithm configurations, then supplement additional empirical results.

C.1 A Summary of Baselines

In our experiments, we have compared AD-OGP with three baseline algorithms. We now give a summary of the compared baselines, as the readers might expect.

1. *Synchronous online gradient push* [2]. This is the synchronous counterpart of AD-OGP. It is compared in Figure 1 of Section 4.1, which verifies the efficiency of AD-OGP compared to its synchronous counterpart.
2. *Asynchronous stochastic gradient descent* [4]. This is a classic asynchronous decentralized algorithm from stochastic optimization. We tailor it to the pure online setting, and present the comparison results in Figure 2 of Section 4.3.1. It verifies the effectiveness of asymmetric gossiping in AD-OGP.
3. *Asynchronous stochastic gradient-push* [5]. This is the state-of-the-art asynchronous algorithm from stochastic optimization. We also tailor it to the pure online setting, and present the comparison results in Figure 3 of Section 4.3.2 (Figure 3 is on the right of Figure 2). It verifies the effectiveness of instantaneous model averaging in AD-OGP.

C.2 Logistic Loss for Binary Classification

In binary classification, the class set is $\mathcal{C} = \{\pm 1\}$. Each learner is parametrized by the (vector) model $\mathbf{x} \in \mathbb{R}^n$ and the push-sum weight $y > 0$. Given the feature vector of an instance $\boldsymbol{\xi} \in \mathbb{R}^n$, the learner

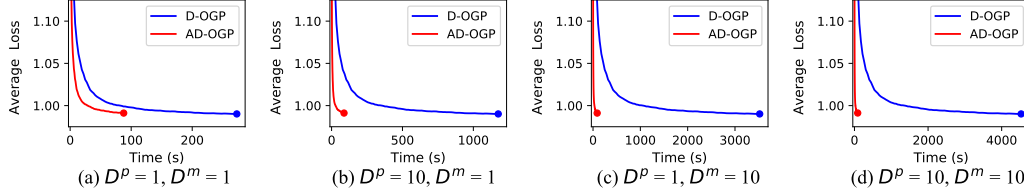


Figure 6: More results to verify the benefit of asynchronization. The four plots compare AD-OGP and D-OGP on a 64-node *complete* graph under varying levels of processing delays D^p and message delays D^m on *pokerhand*.

predicts its class label as $\text{sgn}(\mathbf{w}^\top \boldsymbol{\xi})$, where $\mathbf{w} = \mathbf{x}/y$ is restricted in \mathcal{K} . Suppose the true class label of such instance is $l \in \{\pm 1\}$. We adopt the logistic loss, which is defined as

$$f(\mathbf{w}) = \log(1 + \exp(-l\mathbf{w}^\top \boldsymbol{\xi})).$$

C.3 Multivariate Logistic Loss for Multi-Class Classification

In H -class classification ($H \geq 3$), the class set is $\mathcal{C} = \{1, \dots, H\}$. Each learner is parametrized by the model (in a matrix form) $\mathbf{X} = [\mathbf{x}_1; \dots; \mathbf{x}_H] \in \mathbb{R}^{n \times H}$ and the push-sum weight $y > 0$. Given the feature vector of an instance $\boldsymbol{\xi} \in \mathbb{R}^n$, the learner predicts its class label as $\arg \max_{h \in \mathcal{C}} \mathbf{w}_h^\top \boldsymbol{\xi}$, where $\mathbf{w}_h = \mathbf{x}_h/y, \forall h \in \mathcal{C}$, is restricted within \mathcal{K} . Suppose the true class label of the instance is $l \in \mathcal{C}$. We adopt the multivariate logistic loss [11], which is defined as

$$f(\mathbf{W}) = \log\left(\sum_{h \in \mathcal{C}} \exp(\mathbf{w}_h^\top \boldsymbol{\xi} - \mathbf{w}_l^\top \boldsymbol{\xi})\right),$$

where $\mathbf{W} = \mathbf{X}/y$ is the decision (in a matrix form). Note that, the matrix form is compatible with our framework and theoretical analysis. Specifically, for each learner, we can concatenate all column vectors $\mathbf{x}_1, \dots, \mathbf{x}_h$ of the matrix model \mathbf{X} into a $nh \times 1$ vector $\mathbf{x} \in \mathbb{R}^{nh}$, and all column vectors $\mathbf{w}_1, \dots, \mathbf{w}_h$ of the matrix prediction \mathbf{W} into a $nh \times 1$ vector $\mathbf{w} \in \mathbb{R}^{nh}$. Then we have $\mathbf{w} = \mathbf{x}/y$, and the loss function f is convex w.r.t. \mathbf{w} .

C.4 Network Topology

We here provide detailed descriptions of the various types of network topology examined in our experiments, namely the complete graph, the Watts-Strogatz graph, and the ring graph. All these types of graphs are commonly used in previous studies for decentralized algorithms [7, 38], which represent different levels of connectivity.

1. *Complete graph*. In this kind of graph, each node is connected to any other node. Any complete graph has a diameter of $\mathcal{D} = 1$, which represents a high level of connectivity.
2. *Ring graph*. In this kind of graph, all nodes are arranged in the shape of a ring, and each node is only connected to its two immediate neighbors in the ring. The ring graph with m nodes has a diameter of $\mathcal{D} = \lfloor m/2 \rfloor$, which represents a low level of connectivity.
3. *Watts-Strogatz graph*. This is a kind of random graph, which has two parameters that control its topology, namely the average degree k and the rewiring probability p . In general, a higher average degree or a higher rewiring probability will result in better connectivity of the generated random graph [6]. In our experiments, we follow [34] and set $k = 4, p = 0.3$, to generate random graphs that represent a medium level of connectivity.

C.5 More Details of Algorithm Configurations

We adopt L2-norm balls as the decision set in our experiments. Specifically, for binary classification, it is defined as $\mathcal{K} = \{\mathbf{w} \in \mathbb{R}^n \mid \|\mathbf{w}\| \leq F/2\}$, where $\|\cdot\|$ denotes the L2-norm of any vector. For multi-class classification, it is defined as $\mathcal{K} = \{\mathbf{W} \in \mathbb{R}^{n \times H} \mid \|\mathbf{W}\|_F \leq F/2\}$, where $\|\mathbf{W}\|_F$ is the

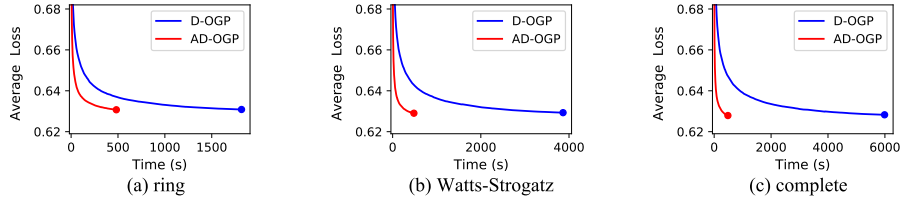


Figure 7: Illustration of the efficiency of AD-OGP in the non-convex setting. The plots compare FMs using AD-OGP and D-OGP on 64-node graphs of varying topologies on *higgs*.

Frobenius norm of the matrix $\mathbf{W} = [\mathbf{w}_1; \dots; \mathbf{w}_H] \in \mathbb{R}^{n \times H}$, i.e., $\|\mathbf{W}\|_F = (\sum_{h=1}^H \|\mathbf{w}_h\|^2)^{1/2}$. We set their diameters $F = 10$.

To configure the learning rate, we follow the common practice in asynchronous online learning [22] by setting $\eta = \alpha\eta^*$, where η^* is the appropriate learning rate theoretically suggested by the given algorithm, and α is chosen from a grid search over $\{\alpha_0(1.25)^i \mid i \in \mathbb{N}\}$ (α_0 is an initial guess for each dataset). In particular, as suggested by Corollary 1 in our main paper, the theoretically optimal η^* is determined by several quantities, namely T , D^{proc} , D^{msg} , Γ_u and G . In our experiments, we use $G = 1$, which has been suggested as a generally good choice in practice [22]. To specify the quantities related to delays (D^{proc} , D^{msg} and Γ_u), we first simulate a small number of rounds (say $T_0 = 10000$) and measure these quantities within the first T_0 rounds [19], and then use them to estimate the quantities over the entire time horizon T .

Note that, although our proposed framework characterizes the number of predictions N_i for each learner $i \in V$, we do NOT need to know these values. Instead, we only need to know the time horizon $T = \sum_{i \in V} N_i$. In some specific scenarios where T is unknown, we can adopt the commonly used doubling trick technique [5] or tune the learning rates by hand [35].

In our experimental setup, the time complexity of each local update or model averaging operation is in $O(n)$, where n is the dimension of model parameters. Hence the overall time complexity is $O(nT)$. All runs are deployed on Xeon(R) E5-2699 @ 2.2GHz.

C.6 More Results to Verify the Benefit of Asynchronization

In our main paper, we report the results of the experiment on Watts-Strogatz graphs on *pokerhand*. Here we also compare AD-OGP and D-OGP on *higgs* as well as other types of networks (i.e., ring graphs or complete graphs). Specifically, we measure the performance of both algorithms on a 64-node ring graph on *higgs*, as well as on a 64-node complete graph on *pokerhand*, under varying levels of processing delays D^p and message delays D^m (recall that, we use D^p (D^m) = 1 to represent a low delay level and D^p (D^m) = 10 to represent a high delay level). The results are plotted in Figure 5 and Figure 6 respectively. They are consistent with the empirical results presented in our main paper. In particular, AD-OGP runs significantly faster than its synchronous counterpart while incurring hardly any loss in the performance.

C.7 Experiments in the Non-Convex Setting

Recall that, our framework and theoretical analysis originate from online convex optimization, which relies on the convexity assumption [10]. Nevertheless, in principle, our proposed AD-OGP algorithm can also be applied to the non-convex setting. Here we conduct additional experiments to verify this point, namely, the efficiency of AD-OGP against its synchronous counterpart [28] is consistent in the non-convex setting.

In this experiment, we adopt the Factorization Machine (FM) [27] as the non-convex model. Specifically, a FM model of degree $d \in \mathbb{N}_+$ is parametrized by a weight vector $\mathbf{w} \in \mathbb{R}^n$, a bias $w_0 \in \mathbb{R}$, and an interaction matrix $\mathbf{V} \in \mathbb{R}^{n \times d}$. Given the feature vector $\boldsymbol{\xi} = (\xi_1, \dots, \xi_n) \in \mathbb{R}^n$ of any instance,

FM predicts its value as

$$h(\boldsymbol{\xi}; \mathbf{w}, w_0, \mathbf{V}) = w_0 + \mathbf{w}^\top \boldsymbol{\xi} + \sum_{i=1}^n \sum_{j=i+1}^n \langle \mathbf{v}_i, \mathbf{v}_j \rangle \xi_i \xi_j,$$

where \mathbf{v}_i represents the i -th row of \mathbf{V} .

We conduct online binary classification experiments on *higgs* and adopt the logistic loss. Specifically, the class label of $\boldsymbol{\xi}$ is predicted as $\text{sgn}(h(\boldsymbol{\xi}; \mathbf{w}, w_0, \mathbf{V}))$. Suppose its true class label is $l \in \{\pm 1\}$, then the loss function is given as $f(\mathbf{w}, w_0, \mathbf{V}) = -\log \sigma(l \cdot h(\boldsymbol{\xi}; \mathbf{w}, w_0, \mathbf{V}))$ where $\sigma(u) = 1/(1 + e^{-u})$.

We employ AD-OGP and its synchronous counterpart D-OGP in the aforementioned non-convex setting. In the experiment, we set the FM's degree d to be 8, which is chosen via a grid search. Other configurations are set to be the same as those for the convex experiments in our main paper. Figure 7 presents the results on 64-node graphs with different topologies. The results show that, in the online non-convex setting, AD-OGP runs significantly faster than its synchronous counterpart, with negligible loss in performance, which accords well with our intuition that AD-OGP can also be applied to the non-convex setting.

D Omitted Proof of Theorem 1

In this section, we present the detailed proof of Theorem 1, which is omitted in our main paper.

Our proof consists of three steps. First, we extend the graph augmentation technique [9] to our online setting to handle message delays, then explicate the transition matrices $\mathbf{P}(t)$ at each round $t \in \{1, \dots, 2T\}$. Second, we decouple the effect of prediction, local update and model averaging, then derive a general regret bound in terms of the transition matrices $\mathbf{P}(t)$ and the associated gradients \mathbf{g}_t and $\hat{\mathbf{g}}_t$. Third, we look deep into the mechanism of push-sum and analyze its convergence rate; plugging it into the above general bound, we successfully derive the bound required in Theorem 1.

D.1 Formulate Transition Matrices on the Augmented Graph

As the first step of regret analysis, we characterize the transition matrices $\mathbf{P}(t)$ at each iteration $t \in \{1, \dots, 2T\}$. Recall the message sending after any local update event $t \in \mathcal{Q}_T$, namely, the updated learner i_t will send out a copy of its updated parameters $(\mathbf{x}'_{i_t}(t), y'_{i_t}(t))$ to each of its neighbor in $\mathcal{N}(i_t)$. Consider its arbitrary neighbor $j \in \mathcal{N}(i_t)$, and suppose that the message $(\mathbf{x}'_{i_t}(t), y'_{i_t}(t))$ will experience a message delay of $d_{i_t j}^m(t)$. For simplicity, abbreviate $(\mathbf{x}'_{i_t}(t), y'_{i_t}(t))$ into (\mathbf{x}', y') and $d_{i_t j}^m(t)$ into d . Then such message will be used for learner j 's model averaging between events $t + d$ and $t + d + 1$. Notably, it is captured by NONE of the learners during the interval $(t, t + d]$.

Such a property will make the recursive relation of local models at different time points extremely complex. In the above example, due to the message delay, a copy of learner i_t 's model $(\mathbf{x}'_{i_t}(t), y'_{i_t}(t))$ after its local update event t will be added to learner j 's model $(\mathbf{x}_j(t + d + 1), y_j(t + d + 1))$ at event $t + d + 1$. Since the message delay d is at most D^{msg} (see Assumption 1), the order of the recursion can be as large as $D^{msg} + 1$. Moreover, each message has its specific delay. Hence the recursion of parameters will be very complex, which makes the analysis almost intractable.

To alleviate this complexity, we introduce the graph augmentation technique [9]. Its main idea is to append some virtual nodes to the decentralized network G , so that any message that undergoes delays will be captured by exactly one of the virtual nodes during its transmission. With this approach, the algorithm can be viewed as running on the augmented graph with zero message delay. Then we can derive a first-order recursion of parameters, which is relatively easier in the subsequent analysis.

Formally, the augmented version $\tilde{G} = (\tilde{V}, \tilde{E})$ of $G = (V, E)$ is constructed in the following two steps. First, regarding each actual node $i \in V$ in the original graph, we add D^{msg} virtual nodes $b_i^1, \dots, b_i^{D^{msg}}$. Second, we add two types of virtual edges to connect these virtual nodes and some actual nodes: (i) Corresponding to each actual node $i \in V$, we add D^{msg} virtual edges, namely (b_i^l, i) and $(b_i^{l+1}, b_i^l), \forall l \in \{1, \dots, D^{msg} - 1\}$; then any actual node $i \in V$ and its corresponding D^{msg} virtual nodes $b_i^1, \dots, b_i^{D^{msg}}$ form a chain $(i, b_i^1, \dots, b_i^{D^{msg}})$ in the augmented graph \tilde{G} . (ii) Corresponding to each actual edge $(i, j) \in E$, we add $2D^{msg}$ virtual edges $(b_i^l, j), (b_j^l, i), \forall l \in$

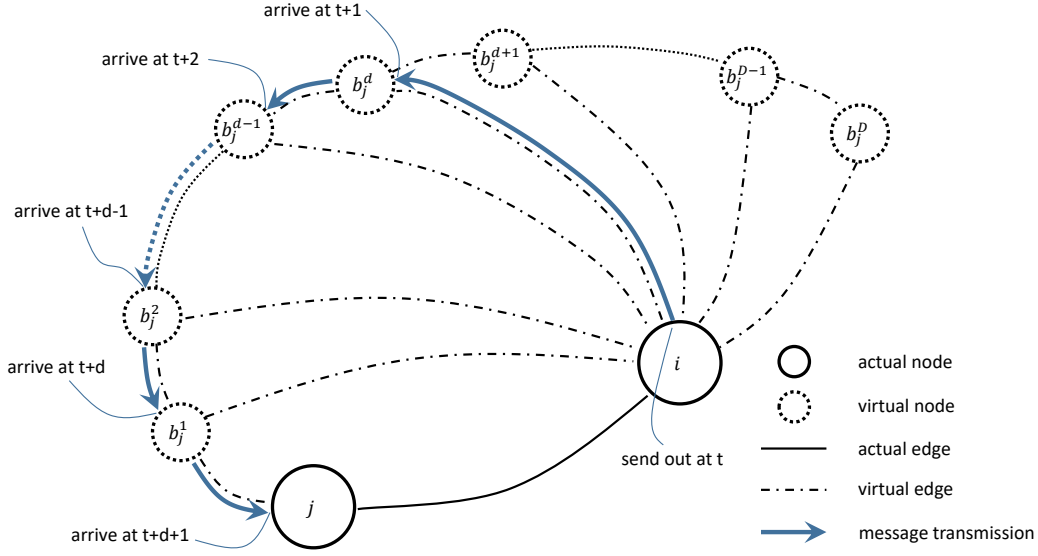


Figure 8: Illustration of message passing on the augmented graph. The two solid circles i and j represent two actual nodes. The solid curve represent the actual edge (i, j) . The D dash circles denote all virtual nodes corresponding to the actual node j , namely b_j^1, \dots, b_j^D . The dash curves represent two types of virtual edges, i.e., those in the chain (b_j^D, \dots, b_j^1, j) as well as those corresponding to the actual edge (i, j) . The blue arrows illustrate the transmission of the message sent from node i to node j with a message delay of d . In particular, it transmits along $i \rightarrow b_j^d \rightarrow \dots \rightarrow b_j^1 \rightarrow j$.

$\{1, \dots, D^{msg}\}$. Mathematically, the augmented graph \tilde{G} is constituted by the augmented vertex set

$$\tilde{V} = V \cup \{b_i^l \mid i \in V, 1 \leq l \leq D^{msg}\},$$

and the augmented edge set

$$\begin{aligned} \tilde{E} = E \cup & \{(b_i^{D^{msg}}, b_i^{D^{msg}-1}), \dots, (b_i^2, b_i^1), (b_i^1, i) \mid i \in V\} \\ & \cup \{(b_i^l, j), (b_j^l, i) \mid (i, j) \in E, 1 \leq l \leq D^{msg}\}. \end{aligned}$$

To ensure that the message is captured by exactly one virtual node during its transmission, we can imagine an information flow along the chain $(b_j^{D^{msg}}, \dots, b_j^1, j)$. Specifically, suppose that some message is sent from learner i to learner j after some event $t \in \{1, \dots, 2T\}$, which has a message delay of $d \geq 0$. Then we can imagine that it first arrives at the virtual node b_j^d at event $t + 1$, then is passed to b_j^{d-1} at event $t + 2$, and so forth until it arrives at b_j^1 at event $t + d$, and finally reaches its destination j at event $t + d + 1$. In this imaginary process, this message is captured by exactly one node (either actual or virtual) at each step during its transmission. Therefore, the message can be viewed as transmitting via the path $i \rightarrow b_j^d \rightarrow \dots \rightarrow b_j^1 \rightarrow j$, with zero message delay at each hop. Note that, when $d = 0$, the transmission trivially reduces to $i \rightarrow j$, namely, the message reaches node j immediately and is processed by node j before the incoming event. To better explicate this point, in Figure 8 we give a demonstration of how a message is diffusing over the augmented graph. Note that, for conciseness, we abbreviate the maximal message delay D^{msg} into D throughout the proof of Theorem 1.

Now we can characterize the evolving dynamics of the whole process via formulating the recursive relation of parameters. For each virtual node $b_i^l \in \tilde{V} - V$, we assume that it also maintains model parameters \mathbf{x}_i^l and push-sum weight y_i^l during the process, and initialize these parameters as $\mathbf{x}_i^l(1) = \mathbf{0} \in \mathbb{R}^n$ and $y_i^l(1) = 0$. In addition, to give a more concise expression of the nodes in \tilde{G} , for each actual node $i \in V$, we also represent it by b_i^0 and rewrite its parameters as $(\mathbf{x}_i^0(t), y_i^0(t)) = (\mathbf{x}_i(t), y_i(t))$. In the following, we will can formulate the evolution of model parameters $(\mathbf{x}_i^l(t), y_i^l(t))$ of either actual ($l = 0$) or virtual ($l \geq 1$) node throughout the learning process.

(A) First we analyze the evolution of model parameters *at any specific event* $t \in \{1, \dots, 2T\}$. If event t is a local update event ($t \in \mathcal{P}_T$), then the local model of any learner will not change at this event. If event t is a local update event ($t \in \mathcal{Q}_T$), then i_t is the only learner whose model parameters changes at this event, with the update form

$$\mathbf{x}'_{i_t}(t) = \gamma_{i_t} y_{i_t}(t) \Pi_{\mathcal{K}}\left(\frac{\mathbf{x}_{i_t}(t) - \eta \mathbf{g}_t}{y_{i_t}(t)}\right), \quad y'_{i_t}(t) = \gamma_{i_t} y_{i_t}(t).$$

In the above update form, we first cast aside the multiplication operation with γ_{i_t} (which is regarded as part of model averaging and will be tackled later). Denote the local model immediately after the projected gradient descent step as

$$\mathbf{z}_{i_t}(t) = y_{i_t}(t) \Pi_{\mathcal{K}}\left(\frac{\mathbf{x}_{i_t}(t) - \eta \mathbf{g}_t}{y_{i_t}(t)}\right).$$

Since the projection operation $\Pi_{\mathcal{K}}$ will introduce nonlinearity to the local update, we denote the residual of weighted projection as $\mathbf{r}_t = \mathbf{z}_{i_t}(t) - (\mathbf{x}_{i_t}(t) - \eta \mathbf{g}_t)$ and linearize the local update as

$$\mathbf{z}_{i_t}(t) = \mathbf{x}_{i_t}(t) - \eta \mathbf{g}_t + \mathbf{r}_t.$$

More compactly, we denote the incremental term associated with this local update as $\boldsymbol{\epsilon}_t = -\eta \mathbf{g}_t + \mathbf{r}_t \in \mathbb{R}^n$. In summary, for each learner $i \in V$, the change of its local model (without multiplication with γ_i) at any event $t \in \{1, \dots, 2T\}$ takes

$$\mathbf{z}_i^l(t) = \begin{cases} \mathbf{x}_i^l(t) + \boldsymbol{\epsilon}_t, & t \in \mathcal{Q}_T, i = i_t, l = 0; \\ \mathbf{x}_i^l(t), & \text{else.} \end{cases}$$

In addition, its push-sum weight $y_i^l(t)$ stays unchanged.

(B) Then we give a matrix form of model averaging *between any two consecutive events* t and $t + 1$. Specifically, for the transition matrix $\mathbf{P}(t)$ between events t and $t + 1$ (which is now defined on the augmented graph \tilde{G}), we identify each of its entries $p_{ij}(t)$. Recall that, for each entry $p_{ij}(t)$ of the transition matrix $\mathbf{P}(t)$, it specifies how much portion of node j 's push-sum weight $y_j^l(t)$ immediately before event t will be added to node i 's push-sum weight $y_i^l(t + 1)$ until the next event $t + 1$. Note that, as the averaging operation on model parameters $\mathbf{x}_i^l(t)$ and push-sum weight $y_i^l(t)$ are totally the same, the same transition matrix $\mathbf{P}(t)$ is also applied to model parameters.

As a prerequisite, we represent the parameters of all nodes in the augmented graph via a compact matrix form. Specifically, for each node $b_i^l \in \tilde{V}$, we view its local model \mathbf{x}_i^l as a *row vector* in $\mathbb{R}^{1 \times n}$. Then we can formulate the matrix model $\mathbf{X}(t)$ by stacking the local models of all nodes (both actual and virtual) together into a $m(D + 1) \times n$ matrix, namely

$$\mathbf{X}(t) = [\mathbf{x}_1(t)^\top, \dots, \mathbf{x}_m(t)^\top, \mathbf{x}_1^1(t)^\top, \dots, \mathbf{x}_m^1(t)^\top, \dots, \mathbf{x}_1^D(t)^\top, \dots, \mathbf{x}_m^D(t)^\top]^\top \in \mathbb{R}^{m(D+1) \times n}.$$

Similarly, we can denote the compact form of push-sum weights $\mathbf{y}(t)$ by stacking the push-sum weights $y_i^l(t)$ of all nodes (both actual and virtual) into a $m(D + 1) \times 1$ *column vector*, i.e.,

$$\mathbf{y}(t) = [y_1(t), \dots, y_m(t), y_1^1(t), \dots, y_m^1(t), \dots, y_1^D(t), \dots, y_m^D(t)]^\top \in \mathbb{R}^{m(D+1)}.$$

Now the parameters of each node can be specified via $\mathbf{X}(t)$ and $\mathbf{y}(t)$. Specifically, (i) For each actual node $i \in V$, its local model constitutes the i -th row of $\mathbf{X}(t)$, i.e., $\mathbf{x}_i(t) = \mathbf{e}_i^\top \mathbf{X}(t) = \mathbf{X}_{i \cdot}(t)$ ¹, and its push-sum weight $y_i(t)$ is exactly the i -th entry of $\mathbf{y}(t)$. (ii) For any virtual node $b_i^l, i \in V, l \in \{1, \dots, D\}$, its local model constitutes the $(lm + i)$ -th row of $\mathbf{X}(t)$, i.e., $\mathbf{x}_i^l(t) = \mathbf{e}_{lm+i}^\top \mathbf{X}(t) = \mathbf{X}_{(lm+i) \cdot}(t)$, and its push-sum weight $y_i^l(t)$ is the $(lm + i)$ -th entry of $\mathbf{y}(t)$. Since the index b_i^l is a bit lengthy, we simplify the notations by using an integer $h = lm + i$ to identify this virtual node b_i^l . Then its parameters can be rewritten as $\mathbf{x}_h(t) = \mathbf{e}_h^\top \mathbf{X}(t) = \mathbf{X}_{h \cdot}(t)$ and $y_h(t) = \mathbf{e}_h^\top \mathbf{y}(t)$. Note that, now the vertex set of \tilde{G} is equivalent to $\tilde{V} = \{1, \dots, m(D + 1)\}$, and the virtual nodes are specified in $\tilde{V} - V = \{m + 1, \dots, m(D + 1)\}$.

¹We use \mathbf{e}_i to represent the i -th unit vector in $\mathbb{R}^{m(D+1)}$, whose i -th entry takes 1 and other entries take 0. We also use \mathbf{A}_i to represent the i -th row of any matrix \mathbf{A} .

The transition matrix $\mathbf{P}(t)$ now acts on the model parameters as

$$\mathbf{y}(t+1) = \mathbf{P}(t)\mathbf{y}(t).$$

The entries $p_{ij}(t)$ of $\mathbf{P}(t)$ formulate the multiplication and the averaging operation on push-sum weights, in the form of $y_i(t+1) = \sum_{j \in \tilde{V}} p_{ij}(t)y_j(t)$. To specify the matrix $\mathbf{P}(t)$, first recall that, after a local update event $t \in \mathcal{Q}_T$, learner i will first multiple its model and push-sum weight by a factor of γ_{i_t} , then sent out copies of its updated parameters (e.g., the push-sum weight $y'_{i_t}(t) = \gamma_{i_t}y_{i_t}(t)$). In the following, we investigate the two types of events separately.

(i) If event t is a prediction event, i.e., $t \in \mathcal{P}_T$, then $y'_i(t) = y_i(t), \forall i \in \tilde{V}$.

(i-a) For each actual learner $i \in V$, it can only perform model averaging during the interval between events t and $t+1$. In Algorithm 1, the model averaging takes

$$y_i(t+1) = y'_i(t) + \sum_{(\mathbf{x}', y') \in \mathcal{M}_i(t)} y' = y_i(t) + \sum_{(\mathbf{x}', y') \in \mathcal{M}_i(t)} y',$$

where the set $\mathcal{M}_i(t)$ contains all messages used for model averaging on learner i during the interval between events t and $t+1$. Recall that, in the graph augmentation technique, any message sent to learner i with a delay of d is viewed as transmitting along the chain $b_i^d \rightarrow b_i^{d-1} \rightarrow \dots \rightarrow b_i^1 \rightarrow i$. Consequently, $\mathcal{M}_i(t)$ inherits all messages that are captured by node b_i^1 immediately before event t , and we simply have

$$y_i(t+1) = y_i(t) + y_i^1(t).$$

(i-b) For each virtual node $h = b_i^l, i \in V, l \in \{1, \dots, D\}$, it will pass its parameters to the successor b_i^{l-1} in the chain, and may inherit the parameters of the predecessor b_i^{l+1} in the chain (if b_i^{l+1} does exist, namely $l < D$). Therefore, we have

$$y_i^l(t+1) = \begin{cases} y_i^{l+1}(t), & l \in \{1, \dots, D-1\}, \\ 0, & l = D. \end{cases}$$

Combining (i-a) and (i-b), the transition matrix $\mathbf{P}(t)$ after each prediction event $t \in \mathcal{P}_T$ and before the next event can be formulated as²

$$\mathbf{P}(t) = \begin{bmatrix} I_m & I_m & O_m & \cdots & O_m \\ O_m & O_m & I_m & \cdots & O_m \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ O_m & O_m & O_m & \cdots & I_m \\ O_m & O_m & O_m & \cdots & O_m \end{bmatrix}. \quad (1)$$

(ii) If event t is a local update event, i.e., $t \in \mathcal{Q}_T$, the averaging effect will be more complex, because the updated node i_t will also send out messages.

(ii-a) We first consider actual nodes $i \in V$. For the updated node i_t , it will multiply its push-sum weight by γ_{i_t} immediately after the projected gradient descent step, namely $y'_{i_t}(t) = \gamma_{i_t}y_{i_t}(t)$. For any other actual node $j \in V - \{i_t\}$, its push-sum weight does not need to multiply γ_j , it may receive the message sent from learner i_t if it is a neighbor of i_t and the message has zero message delay, i.e., $j \in \mathcal{N}(i_t)$ and $d_{i_t j}^m(t) = 0$. Then, just like the above case (i-a), each actual node $i \in V$ will inherit the parameters from the virtual node b_i^1 . To give a more compact expression, similar to [31], we introduce an indicator variable $\delta_i^l(t)$ for $i \in V, 0 \leq l \leq D$, i.e.,

$$\delta_i^l(t) = \begin{cases} 1, & i \in \mathcal{N}(i_t) \cup \{i_t\}, l = d_{i_t i}^m(t), \\ 0, & \text{otherwise.} \end{cases}$$

$\delta_i^l(t)$ equals to 1 if and only if node i_t sends a message to node i and such message exactly has a delay of l . Note that, we assume $d_{i_t i_t}^m(t) = 0$, which makes the indicator variable $\delta_i^l(t)$ well-defined. In summary, we have

$$y_i(t+1) = \begin{cases} \gamma_{i_t}y_{i_t}(t) + y_i^1(t), & i = i_t, \\ y_i(t) + y_i^1(t) + \delta_i^0(t)\gamma_{i_t}y_{i_t}(t), & i \neq i_t. \end{cases}$$

²We use I_m to represent the $m \times m$ identity matrix and O_m to represent the $m \times m$ all-zero matrix.

(ii-b) We then consider virtual nodes $h = b_i^l, i \in V, l \in \{1, \dots, D\}$. If its corresponding actual node i is a neighbor of the updated node i_t and the message sent from node i_t has a delay of l , namely $i \in \mathcal{N}(i_t)$ and $d_{i_t i}^m(t) = l$, then this virtual node will receive the model copy of learner i_t . Moreover, just like the above case (i-b), it will pass its parameters to the successor b_i^{l-1} in the chain and may inherit the parameters of the predecessor b_i^{l+1} in the chain (if node b_i^{l+1} exists). In summary, we have

$$y_i^l(t+1) = \begin{cases} \delta_i^l(t) \gamma_{i_t} y_{i_t}(t) + y_i^{l+1}(t), & l \in \{1, \dots, D-1\}, \\ \delta_i^l(t) \gamma_{i_t} y_{i_t}(t), & l = D. \end{cases}$$

Combining (ii-a) and (ii-b), the transition matrix $\mathbf{P}(t)$ after each local update event $t \in \mathcal{Q}_T$ takes

$$\mathbf{P}(t) = \begin{bmatrix} P^0(t) & I_m & O_m & \cdots & O_m \\ P^1(t) & O_m & I_m & \cdots & O_m \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ P^{D-1}(t) & O_m & O_m & \cdots & I_m \\ P^D(t) & O_m & O_m & \cdots & O_m \end{bmatrix}, \quad (2)$$

where each block $P^l(t) \in \mathbb{R}^{m \times m}$ for $l \in \{0, \dots, D\}$ takes

$$P_{ij}^l(t) = \begin{cases} 1, & i = j \neq i_t, l = 0, \\ \delta_i^l(t) \gamma_{i_t}, & j = i_t, \\ 0, & \text{else,} \end{cases}$$

for any $i, j \in \{1, \dots, m\}$.

The transition matrix $\mathbf{P}(t)$ derived above has the following properties:

Lemma 3. *In AD-OGP, for any $t \in \{1, \dots, 2T\}$, the transition matrix $\mathbf{P}(t)$ defined on the augmented graph satisfies:*

- (i) $\mathbf{P}(t)$ is non-negative and column stochastic, i.e., $\mathbf{1}^\top \mathbf{P}(t) = \mathbf{1}^\top$.
- (ii) Its first m diagonal entries $\mathbf{P}_{ii}(t), 1 \leq i \leq m$ are positive. In addition, its positive entries are at least γ_{i_t} if $t \in \mathcal{Q}_T$, and 1 if $t \in \mathcal{P}_T$.

Proof. Since this lemma obviously holds for $t \in \mathcal{P}_T$, we only need to check the case of $t \in \mathcal{Q}_T$. From the above formulation in (B), $\mathbf{P}(t)$ is non-negative, and its minimal positive entry is exactly γ_{i_t} . Moreover, From the formulation of $P^0(t)$ and the fact that $\delta_{i_t}^0(t) = 1$, we have $\mathbf{P}_{jj}(t) = 1$ for $j \in \{1, \dots, m\} - \{i_t\}$ and $\mathbf{P}_{i_t i_t}(t) = \gamma_{i_t}$, which are all positive.

Now we prove the column stochasticity of $\mathbf{P}(t)$. As it trivially holds for the last mD columns, we only need to check the first m columns. For any $j \in \{1, \dots, m\} - \{i_t\}$, the only non-zero entry in the j -th column is $\mathbf{P}_{jj}(t) = 1$. For the i_t -th column, the entries in this column sum up to $\sum_{l=0}^D \sum_{i \in V} \delta_i^l(t) \gamma_{i_t}$. From Assumption 1, the message delay $d_{i_t i}^m(t)$ must take a value within $\{0, \dots, D_m\}$ for any $i \in \mathcal{N}(i_t) \cup \{i_t\}$, and particularly $d_{i_t i_t}^m(t) = 0$. Hence

$$\sum_{l=0}^D \delta_i^l(t) = \begin{cases} 1, & i \in \mathcal{N}(i_t) \cup \{i_t\}, \\ 0, & \text{otherwise.} \end{cases}$$

Recall that $\gamma_{i_t} = 1/(|\mathcal{N}(i_t)| + 1)$. Thus the sum of the entries in its i_t -th column is exactly

$$\sum_{l=0}^D \sum_{i \in V} \delta_i^l(t) \gamma_{i_t} = \sum_{i \in \mathcal{N}(i_t) \cup \{i_t\}} \gamma_{i_t} = (|\mathcal{N}(i_t)| + 1) \gamma_{i_t} = 1.$$

Hence, we prove that $\mathbf{P}(t)$ is column stochastic. \square

D.2 Derive a General Bound in Terms of Transition Matrices and Gradients

Now that we have formulated the transition matrix at each round, we can derive the compact matrix form of the learning procedure.

We define a matrix $\Delta(t) \in \mathbb{R}^{m(D+1) \times n}$ at each round $t \in \{1, \dots, 2T\}$: if $t \in \mathcal{P}_T$, we let $\Delta(t)$ be the $m(D+1) \times n$ all-zero matrix; if $t \in \mathcal{Q}_T$, we let $\Delta(t) = e_{i_t} \epsilon_t$, whose i_t -th row takes ϵ_t and the other rows are the all-zero vector. Then for $t \in \mathcal{Q}_T$, the weighted projected gradient step (without multiplication with γ_{i_t}) takes $\mathbf{X}(t) + \Delta(t)$; and this also holds for $t \in \mathcal{P}_T$ where $\Delta(t)$ is an all-zero matrix. Hence we can call $\Delta(t)$ the *incremental matrix* at round t . Combining it with the transition matrix $\mathbf{P}(t)$ derived before, we can write out the full recursive relations of model parameters or push-sum weights of all learners, namely,

$$\begin{aligned}\mathbf{X}(t+1) &= \mathbf{P}(t)(\mathbf{X}(t) + \Delta(t)), \\ \mathbf{y}(t+1) &= \mathbf{P}(t)\mathbf{y}(t).\end{aligned}\tag{3}$$

From the above recursions, we can now express $\mathbf{X}(t)$ and $\mathbf{y}(t)$ in terms of the initial parameters $\mathbf{X}(0)$, $\mathbf{Y}(0)$ and the incremental terms $\Delta(s)$ at previous rounds $s < t$. For conciseness, for any $s < t$, we define the product of the $t - s$ consecutive transition matrices (termed the *multi-hop transition matrix*) as

$$\mathbf{Q}(s, t) = \mathbf{P}(t-1) \cdots \mathbf{P}(s).$$

We also set $\mathbf{Q}(t, t) = I_{m(D+1)}$ to make the notation $\mathbf{Q}(s, t)$ well-defined. Then, it is easy to verify that $\mathbf{Q}(s, t)$ is column stochastic for any $s \leq t$. Specifically, since all $\mathbf{P}(t)$ are column stochastic, we have

$$\mathbf{1}^\top \mathbf{Q}(s, t) = \mathbf{1}^\top \mathbf{P}(t-1) \cdots \mathbf{P}(s) = \mathbf{1}^\top.$$

Recall that $\Delta(t)$ is an all-zero matrix if $t \in \mathcal{P}_T$ and also notice that $\{1, \dots, t-1\} - \mathcal{P}_T = \{1, \dots, t-1\} \cap \mathcal{Q}_T = \mathcal{Q}_{1,t}$, hence we have

$$\begin{aligned}\mathbf{X}(t) &= \mathbf{Q}(1, t)\mathbf{X}(1) + \sum_{s=1}^{t-1} \mathbf{Q}(s, t)\Delta(s) \\ &= \mathbf{Q}(1, t)\mathbf{X}(1) + \sum_{s \in \mathcal{Q}_{1,t}} \mathbf{Q}(s, t)\Delta(s), \\ \mathbf{y}(t) &= \mathbf{Q}(1, t)\mathbf{y}(1).\end{aligned}\tag{4}$$

Based on the above recursions in the matrix form, we now give a general analysis of the regret. Since the regret is evaluated by the actual predictions made at prediction events, we consider each prediction event $t \in \mathcal{P}_T$. For each learner $i \in V$, since $\mathbf{x}_i(t)$ and $y_i(t)$ denote its most recent parameters immediately before event t , the actual prediction made by learner i is exactly $\mathbf{w}_i(t) = \mathbf{x}_i(t)/y_i(t)$. To analyze the regret, we need to compare the loss $f_i(\mathbf{w}_j(t))$ evaluated at the reference learner j with the loss $f_i(\mathbf{w}^*)$ evaluated by the best fixed prediction \mathbf{w}^* in hindsight.

Inlighted by the analysis of synchronous online algorithms [36], we introduce the notion of "average model $\bar{\mathbf{x}}(t)$ " as a communication-invariant quantity. However, such a quantity in the asynchronous setting is much more complex than that in the synchronous setting due to the existence of message delay. Specifically, in the synchronous setting, $\bar{\mathbf{x}}(t)$ can be computed by simply averaging the local models of all actual learners, i.e., $\bar{\mathbf{x}}(t) = \sum_{i \in V} \mathbf{x}_i(t)/m$. In the asynchronous setting, however, due to the message delay, each message will not be captured by any actual learner during its transmission. Suppose some message is in transmission at some round t , then the average model $\sum_{i \in V} \mathbf{x}_i(t)/m$ of all actual learners will not cover this message, and consequently, the recursion of the average model will crash. To resolve this issue, recall that, each message is captured by exactly one virtual node in the augmented graph \tilde{G} during its transmission. Therefore, when calculating the average model, the parameters captured by any virtual node should also be taken into account. Consequently, the average model can be defined as

$$\bar{\mathbf{x}}(t) = \frac{1}{m} \sum_{i \in V} \sum_{l=0}^D \mathbf{x}_i^l(t) = \frac{1}{m} \sum_{h \in \tilde{V}} \mathbf{x}_h(t) = \frac{1}{m} \mathbf{1}^\top \mathbf{X}(t),$$

where \tilde{V} is the aforementioned augmented vertex set. Plug the formulation (4) of $\mathbf{X}(t)$ into the above equation, and recall that $\mathbf{Q}(s, t)$ is column stochastic for all $s \leq t$, then we have

$$\bar{\mathbf{x}}(t) = \frac{1}{m} \mathbf{1}^\top (\mathbf{X}(1) + \sum_{s \in \mathcal{Q}_{1,t}} \mathbf{\Delta}(s)).$$

Also recall that, the local model of each node is initialized as $\mathbf{x}_i(1) = \mathbf{w}_0$ and $\mathbf{x}_i^l(1) = \mathbf{0}$ for all $i \in V, l \in \{1, \dots, D\}$. Hence the first m rows of $\mathbf{X}(1)$ are \mathbf{w}_0 , and the last mD rows of $\mathbf{X}(1)$ are all-zero vectors. Moreover, for any $t \in \mathcal{Q}_T$, the i_t -th row of $\mathbf{\Delta}(t)$ takes $\boldsymbol{\epsilon}_t = -\eta \mathbf{g}_t + \mathbf{r}_t$, and the other rows of $\mathbf{\Delta}(t)$ are all-zero vectors. Therefore, the average model $\bar{\mathbf{x}}(t)$ has a simple form, namely,

$$\bar{\mathbf{x}}(t) = \mathbf{w}_0 + \frac{1}{m} \sum_{s \in \mathcal{Q}_{1,t}} (\mathbf{r}_s - \eta \mathbf{g}_s).$$

Similarly, we can define the average push-sum weight $\bar{y}(t)$ at each round $t \in \{1, \dots, 2T\}$. Specifically, we sum up the push-sum weights $y_i^l(t)$ of all nodes b_i^l for $i \in V, 0 \leq l \leq D$ and divide it by m , which turns out to be

$$\bar{y}(t) = \frac{1}{m} \mathbf{1}^\top \mathbf{y}(t) = \frac{1}{m} \mathbf{1}^\top \mathbf{Q}(1, t) \mathbf{y}(1) = \frac{1}{m} \mathbf{1}^\top \mathbf{y}(1),$$

where the last step utilizes the column stochasticity of $\mathbf{Q}(1, t)$. Moreover, since the push-sum weight is initialized as $y_i(1) = 1$ and $y_i^l(1) = 0$ for $i \in V, l \in \{1, \dots, D\}$, we have $\bar{y}(t) = 1$ for any $t \in \{1, \dots, 2T\}$, which means that the average push-sum weight stays unchanged during the learning process. For now, at each event $t \in \{1, \dots, 2T\}$ (either a prediction event or a local update event), the average model $\bar{\mathbf{x}}(t)$ and the average push-sum weight $\bar{y}(t)$ can induce a prediction termed the *system prediction*³

$$\bar{\mathbf{w}}(t) = \frac{\bar{\mathbf{x}}(t)}{\bar{y}(t)},$$

which is equivalent to the average model $\bar{\mathbf{x}}(t)$ since $\bar{y}(t) = 1$. Using the formulation (4) of $\mathbf{X}(t)$, we can also write out the recursive relation of the system prediction $\bar{\mathbf{w}}(t)$ at different events

$$\begin{aligned} \bar{\mathbf{w}}(t+1) &= \bar{\mathbf{x}}(t+1) = \mathbf{w}_0 + \frac{1}{m} \sum_{s \in \mathcal{Q}_{1,t+1}} (\mathbf{r}_s - \eta \mathbf{g}_s) \\ &= \bar{\mathbf{x}}(t) + \delta_t \left(-\frac{\eta}{m} \mathbf{g}_t + \frac{1}{m} \mathbf{r}_t \right) = \bar{\mathbf{w}}(t) + \delta_t \left(-\frac{\eta}{m} \mathbf{g}_t + \frac{1}{m} \mathbf{r}_t \right). \end{aligned}$$

Recall that, the indicator variable δ_t equals to 0 if $t \in \mathcal{P}_T$ and to 1 if $t \in \mathcal{Q}_T$.

We further investigate the distance between the system prediction $\bar{\mathbf{w}}(t+1)$ and the best fixed prediction \mathbf{w}^* in hindsight. Specifically, for any $t \in \mathcal{P}_T$, we simply have

$$\|\bar{\mathbf{w}}(t+1) - \mathbf{w}^*\|^2 = \|\bar{\mathbf{w}}(t) - \mathbf{w}^*\|^2. \quad (5)$$

For any $t \in \mathcal{Q}_T$, the form is slightly more complex

$$\begin{aligned} \|\bar{\mathbf{w}}(t+1) - \mathbf{w}^*\|^2 &= \|\bar{\mathbf{w}}(t) - \mathbf{w}^*\|^2 + \frac{1}{m^2} \|\mathbf{r}_t - \eta \mathbf{g}_t\|^2 \\ &\quad + \frac{2}{m} \mathbf{r}_t^\top (\bar{\mathbf{w}}(t) - \mathbf{w}^*) - \frac{2\eta}{m} \mathbf{g}_t^\top (\bar{\mathbf{w}}(t) - \mathbf{w}^*). \end{aligned} \quad (6)$$

We further investigate the latter case $t \in \mathcal{Q}_T$. Notice that, the recursion (6) is related with the residual \mathbf{r}_t of the weighted projection operation. We show that, the norm of the residual is upper bounded by the norm of the associated gradient \mathbf{g}_t .

Lemma 4. *In AD-OGP, the residual \mathbf{r}_t of the weighted projection operation at each local update event $t \in \mathcal{Q}_T$ is defined as above. Then its norm is bounded as*

$$\|\mathbf{r}_t\| \leq \eta g_t.$$

³Remind that, in our definition, the system prediction $\bar{\mathbf{w}}(t)$ is induced by the average parameters $\bar{\mathbf{x}}(t)$ and $\bar{y}(t)$ of all learners at event t , which does NOT necessarily equal to the average of the predictions generated by all actual learners at event t , i.e., $\bar{\mathbf{w}}(t) \neq \sum_{i \in V} \mathbf{w}_i(t)/m$.

Proof. Recall our definition of weighted projection operation (note that, Lemma 1 ensures that $y > 0$)

$$\Pi_{y\mathcal{K}}(\mathbf{x}) = \operatorname{argmin}_{\mathbf{z} \in y\mathcal{K}} \|\mathbf{z} - \mathbf{x}\| = y \operatorname{argmin}_{\mathbf{w} \in \mathcal{K}} \|\mathbf{w} - \frac{\mathbf{x}}{y}\|,$$

or equivalently

$$\frac{\Pi_{y\mathcal{K}}(\mathbf{x})}{y} = \operatorname{argmin}_{\mathbf{w} \in \mathcal{K}} \|\mathbf{w} - \frac{\mathbf{x}}{y}\|.$$

Denote learner i_t 's local model after gradient descent and *before weighted projection* as $\mathbf{z}_{i_t}(t) = \mathbf{x}_{i_t}(t) - \eta \mathbf{g}_t$, and its induced prediction as $\mathbf{v}_{i_t}(t) = \mathbf{z}_{i_t}(t)/y_{i_t}(t)$. From Proposition 1 we have $\mathbf{x}_{i_t}(t) \in y_{i_t}(t)\mathcal{K}$, and from Lemma 1 we know $y_{i_t}(t) > 0$. Hence

$$\begin{aligned} \|\mathbf{r}_t\| &= \|\Pi_{y_{i_t}(t)\mathcal{K}}(\mathbf{z}_{i_t}(t)) - \mathbf{z}_{i_t}(t)\| = y_{i_t}(t) \left\| \frac{\Pi_{y_{i_t}(t)\mathcal{K}}(\mathbf{z}_{i_t}(t))}{y_{i_t}(t)} - \frac{\mathbf{z}_{i_t}(t)}{y_{i_t}(t)} \right\| \\ &\leq y_{i_t}(t) \left\| \frac{\mathbf{x}_{i_t}(t)}{y_{i_t}(t)} - \frac{\mathbf{z}_{i_t}(t)}{y_{i_t}(t)} \right\| = \|\mathbf{x}_{i_t}(t) - \mathbf{z}_{i_t}(t)\| = \|\eta \mathbf{g}_t\| = \eta g_t, \end{aligned}$$

which proves the lemma. \square

Now we turn back to analyze the recursion (6). Specifically, in the right-hand side of the recursion:

- (i) The first term $\|\bar{\mathbf{w}}(t) - \mathbf{w}^*\|$ is preserved for successive elimination.
- (ii) The second term is bounded as

$$\|\mathbf{r}_t - \eta \mathbf{g}_t\| \leq \|\mathbf{r}_t\| + \eta \|\mathbf{g}_t\| \leq 2\eta g_t.$$

- (iii) For the third term $\mathbf{r}_t^\top (\bar{\mathbf{w}}(t) - \mathbf{w}^*)$, recall the definition of \mathbf{r}_t and the property of standard projection operation onto a convex set

$$(\Pi_{\mathcal{K}}(\mathbf{a}) - \mathbf{a})^\top (\mathbf{a} - \mathbf{b}) \leq 0, \quad \forall \mathbf{a} \in \mathbb{R}^n, \mathbf{b} \in \mathcal{K}.$$

Set $\mathbf{a} = \mathbf{v}_{i_t}(t)$, $\mathbf{b} = \mathbf{w}^*$ in the above inequality, and utilize $y_{i_t}(t) > 0$, then we have

$$\mathbf{r}_t^\top (\mathbf{v}_{i_t}(t) - \mathbf{w}^*) = y_{i_t}(t) (\Pi_{\mathcal{K}}(\mathbf{v}_{i_t}(t)) - \mathbf{v}_{i_t}(t))^\top (\mathbf{v}_{i_t}(t) - \mathbf{w}^*) \leq y_{i_t}(t) \cdot 0 = 0.$$

Therefore, we can bound the third term as

$$\begin{aligned} \mathbf{r}_t^\top (\bar{\mathbf{w}}(t) - \mathbf{w}^*) &= \mathbf{r}_t^\top (\bar{\mathbf{w}}(t) - \mathbf{v}_{i_t}(t)) + \mathbf{r}_t^\top (\mathbf{v}_{i_t}(t) - \mathbf{w}^*) \\ &\leq \|\mathbf{r}_t\| \|\bar{\mathbf{w}}(t) - \mathbf{v}_{i_t}(t)\| \leq \eta g_t \|\bar{\mathbf{w}}(t) - \mathbf{v}_{i_t}(t)\|, \end{aligned}$$

where the last inequality utilizes Lemma 4.

- (iv) For the last term $\mathbf{g}_t^\top (\mathbf{w}^* - \bar{\mathbf{w}}(t))$, recall that $\mathbf{g}_t = \nabla f_{l_t}(\mathbf{w}_{i_t}(l_t))$ and f_{l_t} is convex, then we have

$$\begin{aligned} \mathbf{g}_t^\top (\mathbf{w}^* - \bar{\mathbf{w}}(t)) &= \mathbf{g}_t^\top (\mathbf{w}_{i_t}(l_t) - \bar{\mathbf{w}}(t)) + \mathbf{g}_t^\top (\mathbf{w}^* - \mathbf{w}_{i_t}(l_t)) \\ &\leq g_t \|\mathbf{w}_{i_t}(l_t) - \bar{\mathbf{w}}(t)\| + f_{l_t}(\mathbf{w}^*) - f_{l_t}(\mathbf{w}_{i_t}(l_t)) \end{aligned}$$

Plug the above terms (i)-(iv) into the recursion (6) and rearrange the inequality, then we obtain

$$\begin{aligned} f_{l_t}(\mathbf{w}_{i_t}(l_t)) - f_{l_t}(\mathbf{w}^*) &\leq \frac{m}{2\eta} (\|\bar{\mathbf{w}}(t) - \mathbf{w}^*\|^2 - \|\bar{\mathbf{w}}(t+1) - \mathbf{w}^*\|^2) \\ &\quad + g_t (\|\mathbf{w}'_{i_t}(t) - \bar{\mathbf{w}}(t)\| + \|\mathbf{w}_{i_t}(l_t) - \bar{\mathbf{w}}(t)\|) + \frac{2\eta}{m} g_t. \end{aligned}$$

Moreover, to relate the above inequality to loss evaluated at the reference learner j , namely $f_{l_t}(\mathbf{w}_j(l_t))$, we again utilize the convexity of f_{l_t} and recall the definition $\hat{\mathbf{g}}_t = \nabla f_{l_t}(\mathbf{w}_j(l_t))$. Then we have

$$f_{l_t}(\mathbf{w}_j(l_t)) - f_{l_t}(\mathbf{w}_{i_t}(l_t)) \leq (\hat{\mathbf{g}}_t)^\top (\mathbf{w}_j(l_t) - \mathbf{w}_{i_t}(l_t)) \leq \hat{g}_t \|\mathbf{w}_j(l_t) - \mathbf{w}_{i_t}(l_t)\|.$$

Summing the above two inequality together, we derive an upper bound for $f_{l_t}(\mathbf{w}_j(l_t)) - f_{l_t}(\mathbf{w}^*)$. Notably, the upper bound is in terms of $\|\mathbf{w}'_{i_t}(t) - \bar{\mathbf{w}}(t)\|$, $\|\mathbf{w}_{i_t}(l_t) - \bar{\mathbf{w}}(t)\|$ and $\|\mathbf{w}_j(l_t) - \mathbf{w}_{i_t}(l_t)\|$, which capture the distance between some specific pair of predictions and need to be further analyzed. In fact, we can decompose these distance terms as

$$\begin{aligned} \|\mathbf{w}'_{i_t}(t) - \bar{\mathbf{w}}(t)\| &\leq \|\mathbf{w}_{i_t}(t) - \bar{\mathbf{w}}(t)\| + \frac{\eta g_t}{y_{i_t}(t)}, \\ \|\mathbf{w}_{i_t}(l_t) - \bar{\mathbf{w}}(t)\| &\leq \|\mathbf{w}_{i_t}(l_t) - \bar{\mathbf{w}}(l_t)\| + \|\bar{\mathbf{w}}(t) - \bar{\mathbf{w}}(l_t)\|, \\ \|\mathbf{w}_j(l_t) - \mathbf{w}_{i_t}(l_t)\| &\leq \|\mathbf{w}_{i_t}(l_t) - \bar{\mathbf{w}}(l_t)\| + \|\mathbf{w}_j(l_t) - \bar{\mathbf{w}}(l_t)\|. \end{aligned}$$

In summary, for each $t \in \mathcal{Q}_t$, the quantity $f_{l_t}(\mathbf{w}_j(l_t)) - f_{l_t}(\mathbf{w}^*)$ can be bounded as

$$\begin{aligned} f_{l_t}(\mathbf{w}_j(l_t)) - f_{l_t}(\mathbf{w}^*) &\leq \frac{m}{2\eta} (\|\bar{\mathbf{w}}(t) - \mathbf{w}^*\|^2 - \|\bar{\mathbf{w}}(t+1) - \mathbf{w}^*\|^2) \\ &\quad + g_t(\|\mathbf{w}_{i_t}(t) - \bar{\mathbf{w}}(t)\| + \|\mathbf{w}_{i_t}(l_t) - \bar{\mathbf{w}}(l_t)\|) \\ &\quad + g_t(\|\bar{\mathbf{w}}(t) - \bar{\mathbf{w}}(l_t)\| + (\frac{2}{m} + \frac{1}{y_{i_t}(t)})\eta g_t) \\ &\quad + \hat{g}_t(\|\mathbf{w}_{i_t}(l_t) - \bar{\mathbf{w}}(l_t)\| + \|\mathbf{w}_j(l_t) - \bar{\mathbf{w}}(l_t)\|). \end{aligned} \quad (7)$$

To further analyze the above quantity, we are particularly interested in $\|\bar{\mathbf{w}}(s) - \bar{\mathbf{w}}(t)\|$ and $\|\mathbf{w}_i(t) - \bar{\mathbf{w}}(t)\|$ for any $i \in V, 0 \leq s \leq t$. In the following, we analyze these terms, respectively.

We first derive a bound for $\|\bar{\mathbf{w}}(s) - \bar{\mathbf{w}}(t)\|$ as follows. Recall that we use $\mathcal{Q}_{s,t} = \mathcal{Q}_T \cap \{s+1, \dots, t-1\}$ to denote all update events that happen within (s, t) .

Lemma 5. *In AD-OGP, the average prediction $\bar{\mathbf{w}}(t)$ is defined above. Then for any $1 \leq s \leq t \leq 2T$,*

$$\|\bar{\mathbf{w}}(s) - \bar{\mathbf{w}}(t)\| \leq \frac{2\eta}{m} \sum_{k \in \mathcal{Q}_{s-1,t}} g_k.$$

Proof. Recall the recursion (6) of the system predictions, i.e.,

$$\bar{\mathbf{w}}(t+1) = \bar{\mathbf{w}}(t) + \delta_t(-\frac{\eta}{m}\mathbf{g}_t + \frac{1}{m}\mathbf{r}_t).$$

We then have

$$\|\bar{\mathbf{w}}(s) - \bar{\mathbf{w}}(t)\| \leq \sum_{k=s}^{t-1} \delta_k(\frac{\eta}{m}\|\mathbf{g}_k\| + \frac{1}{m}\|\mathbf{r}_k\|) \leq \sum_{k \in \mathcal{Q}_{s-1,t}} (\frac{\eta}{m}\|\mathbf{g}_k\| + \frac{1}{m}\|\mathbf{r}_k\|).$$

From Lemma 4 we also know that $\|\mathbf{r}_t\| \leq \eta g_t$. Hence we prove the lemma. \square

We then analyze the term $\|\mathbf{w}_i(t) - \bar{\mathbf{w}}(t)\|$, which is a bit more complex. In fact, this term measures the distance between the prediction of the specific learner i and the system prediction, which is thus affected by both local update and communication. The following lemma shows that such term has a bound in terms of the incremental term ϵ_s at local weight updates as well as the multiple-hop transition matrices $\mathbf{Q}(s, t)$.

Lemma 6. *In AD-OGP, the average prediction $\bar{\mathbf{w}}(t)$, the residual of projection \mathbf{r}_t and the multi-hop transition matrix $\mathbf{Q}(s, t)$ are defined above. Then for any $t \in \{1, \dots, 2T\}$, it holds that*

$$\|\mathbf{w}_i(t) - \bar{\mathbf{w}}(t)\| \leq \sum_{s \in \mathcal{Q}_{1,t}} R_i(s, t),$$

where the quantity $R_i(s, t)$ is defined as⁴

$$R_i(s, t) = \left\| \frac{\mathbf{e}_i^\top \mathbf{Q}(s, t) \Delta(s)}{\mathbf{e}_i^\top \mathbf{Q}(s, t) \mathbf{y}(s)} - \frac{\epsilon_s}{m} \right\|.$$

Proof. Recall that, for any actual learner $i \in V$, its prediction $\mathbf{w}_i(t)$ at round t is determined by $\mathbf{x}_i(t)/y_i(t) = \mathbf{e}_i^\top \mathbf{X}(t)/\mathbf{e}_i^\top \mathbf{y}(t)$. From the above formulations (4) of $\mathbf{X}(t)$ and $\mathbf{y}(t)$, the prediction can be expressed in terms of the initial local model $\mathbf{X}(1)$ and the update matrix $\Delta(s)$, namely,

$$\mathbf{w}_i(t) = \frac{\mathbf{e}_i^\top \mathbf{X}(t)}{\mathbf{e}_i^\top \mathbf{y}(t)} = \frac{\mathbf{e}_i^\top \mathbf{Q}(1, t) \mathbf{X}(1)}{\mathbf{e}_i^\top \mathbf{Q}(1, t) \mathbf{y}(1)} + \sum_{s \in \mathcal{Q}_{1,t}} \frac{\mathbf{e}_i^\top \mathbf{Q}(s, t) \Delta(s)}{\mathbf{e}_i^\top \mathbf{Q}(s, t) \mathbf{y}(s)}.$$

Since all actual learners are initialized as $\mathbf{x}_i(1) = \mathbf{w}_0, y_i(1) = 1, i \in V$, and all virtual learners are initialized as $\mathbf{x}_i^l(1) = \mathbf{0}, y_i^l(1) = 0, i \in V, l \in \{1, \dots, D\}$, we can verify that (here we view $\mathbf{w}_0 \in \mathbb{R}^{1 \times n}$ as a row vector)

$$\mathbf{X}(1) = \begin{bmatrix} \mathbf{1}_m \mathbf{w}_0 \\ \mathbf{0}_{mD \times n} \end{bmatrix} = \begin{bmatrix} \mathbf{1}_m \\ \mathbf{0}_{mD} \end{bmatrix} \mathbf{w}_0 = \mathbf{y}(1) \mathbf{w}_0.$$

⁴Recall that, \mathbf{e}_i denotes the i -th unit vector in $\mathbb{R}^{m(D+1)}$; for any $t \in \mathcal{Q}_T$, we have also denoted the update matrix at round t as $\Delta(s) = \mathbf{e}_{i_s} \epsilon_t$ and the update term at round t as $\epsilon_t = \mathbf{r}_t - \eta \mathbf{g}_t$.

Since $e_i^\top \mathbf{Q}(1, t) \mathbf{y}(1)$ is a scalar, we directly have

$$\frac{e_i^\top \mathbf{Q}(1, t) \mathbf{X}(1)}{e_i^\top \mathbf{Q}(1, t) \mathbf{y}(1)} = \frac{(e_i^\top \mathbf{Q}(1, t) \mathbf{y}(1)) \mathbf{w}_0}{e_i^\top \mathbf{Q}(1, t) \mathbf{y}(1)} = \mathbf{w}_0,$$

and thus

$$\mathbf{w}_i(t) = \mathbf{w}_0 + \sum_{s \in \mathcal{Q}_{1,t}} \frac{e_i^\top \mathbf{Q}(s, t) \boldsymbol{\Delta}(s)}{e_i^\top \mathbf{Q}(s, t) \mathbf{y}(s)}.$$

A similar deduction on $\bar{\mathbf{w}}(t)$ gives

$$\bar{\mathbf{w}}(t) = \frac{1}{m} \mathbf{1}^\top \mathbf{X}(t) = \mathbf{w}_0 + \sum_{s \in \mathcal{Q}_{1,t}} \frac{\boldsymbol{\epsilon}_s}{m}.$$

Plugging the above expressions of $\mathbf{w}_i(t)$ and $\bar{\mathbf{w}}(t)$ into the term $\|\mathbf{w}_i(t) - \bar{\mathbf{w}}(t)\|$, we derive

$$\|\mathbf{w}_i(t) - \bar{\mathbf{w}}(t)\| \leq \sum_{s \in \mathcal{Q}_{1,t}} R_i(s, t),$$

which proves the lemma. \square

Utilizing the above Lemma 5 and Lemma 6, we can derive a general bound in terms of $R_i(s, t)$ (which is a function of transition matrix $\mathbf{Q}(s, t)$ and update term $\boldsymbol{\epsilon}_t$) and the associated gradients $\mathbf{g}_t, \hat{\mathbf{g}}_t$. Specifically, for any $t \in \mathcal{Q}_T$, we can apply both lemmas to inequality (7), which derives

$$\begin{aligned} f_{l_t}(\mathbf{w}_j(l_t)) - f_{l_t}(\mathbf{w}^*) &\leq \frac{m}{2\eta} (\|\bar{\mathbf{w}}(t) - \mathbf{w}^*\|^2 - \|\bar{\mathbf{w}}(t+1) - \mathbf{w}^*\|^2) \\ &\quad + \sum_{s \in \mathcal{Q}_{1,t}} g_t R_{i_t}(s, t) + \sum_{s \in \mathcal{Q}_{1,l_t}} g_t R_{i_t}(s, l_t) \\ &\quad + \frac{2\eta}{m} g_t (g_t + \sum_{s \in \mathcal{Q}_{l_t-1,t}} g_s) + \frac{1}{y_{i_t}(t)} \eta g_t^2 \\ &\quad + \sum_{s \in \mathcal{Q}_{1,l_t}} \hat{g}_t (R_{i_t}(s, l_t) + R_j(s, l_t)). \end{aligned} \tag{8}$$

Recall the definition $\mathcal{Q}_{s,t} = \{k \mid s < k < t, \delta_k = 1\}$. Since event l_t is a prediction event ($\delta_{l_t} = 0$) and event t is a local update event ($\delta_t = 1$), we have $\mathcal{Q}_{l_t-1,t} \cup \{t\} = \mathcal{Q}_{l_t,t+1}$. Hence $g_t + \sum_{s \in \mathcal{Q}_{l_t-1,t}} g_s = \sum_{s \in \mathcal{Q}_{l_t,t+1}} g_s$.

Moreover, for any $t \in \mathcal{P}_T$, from the equation (5), we trivially have

$$0 = \frac{m}{2\eta} (\|\bar{\mathbf{w}}(t) - \mathbf{w}^*\|^2 - \|\bar{\mathbf{w}}(t+1) - \mathbf{w}^*\|^2). \tag{9}$$

We sum the above inequality (8) and equation (9) over $t \in \{1, \dots, 2T\} = \mathcal{P}_T \cup \mathcal{Q}_T$. The left-hand side of the inequalities sum up to $\sum_{t \in \mathcal{Q}_T} f_{l_t}(\mathbf{w}_j(l_t)) - f_{l_t}(\mathbf{w}^*)$. Consider the index set $\{l_t \mid t \in \mathcal{Q}_T\}$. Recall that, in our event indexing system, for any local update event $t \in \mathcal{Q}_T$, l_t is defined as the index of its corresponding prediction event. Therefore, the index set $\{l_t \mid t \in \mathcal{Q}_T\}$ is exactly a *permutation* of the indices of the prediction events \mathcal{P}_T . Therefore, the left-hand side has an equivalent expression, namely,

$$\sum_{t \in \mathcal{Q}_T} f_{l_t}(\mathbf{w}_j(l_t)) - f_{l_t}(\mathbf{w}^*) = \sum_{t \in \mathcal{P}_T} f_t(\mathbf{w}_j(t)) - f_t(\mathbf{w}^*),$$

which is exactly the definition of our regret.

Therefore, to analyze the regret, it suffices to bound the right-hand side of the summation of the inequalities. From $\bar{\mathbf{w}}(1) = \mathbf{w}_0$, the first terms in the right-hand side of the inequalities ($t \in \mathcal{Q}_T$) and equations ($t \in \mathcal{P}_T$) sum up to

$$\frac{m}{2\eta} (\|\bar{\mathbf{w}}(1) - \mathbf{w}^*\|^2 - \|\bar{\mathbf{w}}(2T+1) - \mathbf{w}^*\|^2) \leq \frac{m}{2\eta} \|\mathbf{w}_0 - \mathbf{w}^*\|^2.$$

Now we can derive a general bound, i.e.,

$$\begin{aligned} \text{Regret}_j &\leq \frac{m}{2\eta} \|\mathbf{w}_0 - \mathbf{w}^*\|^2 + \frac{2\eta}{m} \sum_{t \in \mathcal{Q}_T} \sum_{s \in \mathcal{Q}_{i_t, t+1}} g_s g_t + \frac{1}{y_{i_t}(t)} \eta g_t^2 \\ &\quad + \sum_{t \in \mathcal{Q}_T} \sum_{s \in \mathcal{Q}_{1, t}} g_t R_{i_t}(s, t) + \sum_{t \in \mathcal{Q}_T} \sum_{s \in \mathcal{Q}_{1, t_t}} (g_t R_{i_t}(s, l_t) + \hat{g}_t R_{i_t}(s, l_t) + \hat{g}_t R_j(s, l_t)), \end{aligned} \quad (10)$$

where the quantity $R_i(s, t)$ for any $i \in V, s \in \mathcal{Q}_T, t \in \{s, s+1, \dots, 2T\}$ is defined as

$$R_i(s, t) = \left\| \frac{\mathbf{e}_i^\top \mathbf{Q}(s, t) \boldsymbol{\Delta}(s)}{\mathbf{e}_i^\top \mathbf{Q}(s, t) \mathbf{y}(s)} - \frac{\epsilon_s}{m} \right\|.$$

D.3 Analyze the Effect of Push-Sum in Asynchronous Decentralized Online Learning

The final step of our analysis is to derive an upper bound of the quantity $R_i(s, t)$, which specialize the regret to our proposed algorithm. Recall that, the multi-hop transition matrix $\mathbf{Q}(s, t)$ characterizes the overall effect of push-sum among all nodes (actual or virtual) during the interval between events s and t . Also notice that $\epsilon_s/m = \mathbf{1}^\top \boldsymbol{\Delta}(s)/\mathbf{1}^\top \mathbf{y}(s)$ for all $s \leq 2T$, since $\mathbf{1}^\top \mathbf{y}(s) = \mathbf{1}^\top \mathbf{y}(1) = m$.

Now the physical meaning of the quantity $R_i(s, t)$ can be interpreted as below. We imagine a $(t-s)$ -step *distributed averaging* process operated on the augmented graph \tilde{G} (notably, the nodes communicate with each other at each step $k \in \{1, \dots, t-s\}$, but they never perform local update). At the beginning of the averaging process, the local model \mathbf{x}_h and push-sum weight y_h of each node $h \in \tilde{V} = \{1, \dots, m(D+1)\}$ are initialized to be the h -th row of $\boldsymbol{\Delta}(s)$ and the h -th entry of $\mathbf{y}(s)$, respectively. At each step $k = 1, \dots, t-s$, the nodes mix their model parameters according to the transition matrix $\mathbf{P}(s+k-1)$. Consequently, $\mathbf{e}_i^\top \mathbf{Q}(s, t) \boldsymbol{\Delta}(s)$ and $\mathbf{e}_i^\top \mathbf{Q}(s, t) \mathbf{y}(s)$ are exactly node i 's local model \mathbf{x}_i and push-sum weight y_i after the whole imaginary averaging process, respectively. Moreover, $\mathbf{1}^\top \boldsymbol{\Delta}(s)$ and $\mathbf{1}^\top \mathbf{y}(s)$ are exactly the average parameters of all nodes in \tilde{V} , which actually stay unchanged during the averaging process as $\mathbf{Q}(s, t)$ is column stochastic. Therefore, the quantity $R_i(s, t)$ measures the distance between the prediction of node i and the system prediction induced by the average parameters after the $(t-s)$ -step averaging process.

From the above analysis, we know that the quantity $R_i(s, t)$ is determined by the initial states of nodes $h \in \tilde{V}$ and the multi-hop communication matrix $\mathbf{Q}(s, t)$. Therefore, to give a bound for $R_i(s, t)$, we need to investigate the property of multi-hop transition matrix $\mathbf{Q}(s, t)$. Our analysis is enlightened by prior work for asynchronous stochastic optimization [31].

By its definition, each entry $\mathbf{Q}_{ij}(s, t)$ of the transition matrix measures the portion of j 's initial model \mathbf{u}_j attributing to i 's model after $t-s$ push-sum steps. If this entry is positive, namely $\mathbf{Q}_{ij}(s, t) > 0$, it then implies that node i is able to receive information from node j during the push-sum steps. In intuition, after sufficient number of communication steps, each actual node $i \in V$ can access the information on any other node (either actual or virtual) $h \in \tilde{V}$ via communication; alternatively speaking, when $t-s$ is sufficiently large, $\mathbf{Q}_{ih}(s, t)$ is expected to be positive for all $i \in V, h \in \tilde{V}$.

The following lemma gives a lower bound of $t-s$, which guarantees that the entries in first m row of $\mathbf{Q}(s, t)$ are all positive. Note that, [31] also derives a similar result in their Lemma 7, namely $t-s \geq m(D^{msg} + \Gamma_d)$. In our work, via a finer-grained analysis, we only requires $t-s \geq (\mathcal{D}+1)(D^{msg} + \Gamma_d)$ where \mathcal{D} denote the diameter of the original graph G . In other words, our bound is much smaller than the bound in [31], since the diameter \mathcal{D} of a decentralized architecture is usually much smaller than the number of learners m . For example, the diameter of any complete graph is only 1.

Lemma 7. *Denote $\Gamma = D^{msg} + \Gamma_d$. In AD-OGP, for any $s \geq 1$, the entries in the first m rows of $\mathbf{Q}(s, s + (\mathcal{D}+1)\Gamma)$ are positive. In addition, any positive entry of $\mathbf{Q}(s, s + (\mathcal{D}+1)\Gamma)$ is at least*

$$\alpha = \left(\frac{1}{m}\right)^{(\mathcal{D}+1)\Gamma}.$$

Proof. Recall that, in our derived forms (8) and (9) of the transition matrix $\mathbf{P}(t)$ at each push-sum step, the first $m \times m$ block of $\mathbf{P}(t)$ for any event $t \in \{1, \dots, 2T\}$ has all-positive diagonals.

Therefore, for any multi-hop transition matrix $\mathbf{Q}(p, q)$ ($p \leq q$), if its (i, h) -th entry located at the first m rows (i.e., $i \leq m$) is positive, then the corresponding (i, h) -th entry of $\mathbf{Q}(p, q + 1)$ is exactly

$$\mathbf{Q}_{ih}(p, q + 1) = \sum_{j=1}^{m(D+1)} \mathbf{P}_{ij}(q) \mathbf{Q}_{jh}(p, q) \geq \mathbf{P}_{ii}^0(q + 1) \mathbf{Q}_{ih}(p, q),$$

which must be positive as well. This suggests that the positiveness of the entries in the first m row of the multi-hop transition matrix will be preserved as the push-sum process proceeds on. Therefore, to prove the positiveness of the first m rows of $\mathbf{Q}(s, s + (\mathcal{D} + 1)\Gamma)$, it suffices to prove that, for any $i \leq m, h \leq m(D + 1)$, there exists some intermediate step $t \in [s, s + (\mathcal{D} + 1)\Gamma]$ such that the (i, h) -th entry of $\mathbf{Q}(s, t)$ is positive. Equivalently, we only need to prove that, the information on any node $h \in \tilde{V}$ is able to reach all actual nodes $i \in V$ within at most $(\mathcal{D} + 1)\Gamma$ steps.

To this end, we consider an arbitrary node $h = b_j^l \in \tilde{V}$ for any $j \in V, 0 \leq l \leq \mathcal{D}$. The information on node h will reach its corresponding actual node j within l steps, i.e., along the path $b_j^l \rightarrow b_j^{l-1} \rightarrow \dots \rightarrow j$. In addition, since the underlying decentralized network G is connected, there must exist some path linking j and i in G , and the path length r is no more than the diameter \mathcal{D} of the graph. We assume the path to be $j \rightarrow u_1 \rightarrow \dots \rightarrow u_{r-1} \rightarrow i$ where $u_k \in V, k \in \{1, \dots, r - 1\}$ are the intermediate nodes; also denote the endpoints of the path $u_0 = j, u_r = i$.

Consider each segment (u_k, u_{k+1}) of the path for any $k \in \{0, \dots, r - 1\}$. From Assumption 1 and the graph augmentation technique, we know that: (i) the node u_k will perform local update at least once every Γ_d steps, after which it will send out a message to its neighbor u_{k+1} ; (ii) such a message will take at most D steps to reach its destination u_{k+1} . Therefore, starting from any step, it takes at most $\Gamma = D + \Gamma_d$ event steps for the information on u_k to diffuse to u_{k+1} . Since such threshold Γ holds for all segments (u_k, u_{k+1}) , the information on j will reach i within at most $\mathcal{D}\Gamma$ steps. In summary, for an arbitrary node $h \in \tilde{V}$ in the augmented graph, it takes at most $D + \mathcal{D}\Gamma \leq (\mathcal{D} + 1)\Gamma$ communication steps in AD-OGP for its information to diffuse to any actual node $i \in V$. Therefore, all entries in the first m row of $\mathbf{Q}(s, s + (\mathcal{D} + 1)\Gamma)$ are positive.

In addition, from the formulation (8) or (9) of $\mathbf{P}(t)$ for any $t \in \{1, \dots, 2T\}$, its positive entries are at least $\lambda_{i_i} \geq 1/m$. Therefore, since $\mathbf{Q}(s, s + (\mathcal{D} + 1)\Gamma) = \mathbf{P}(s + (\mathcal{D} + 1)\Gamma - 1) \cdots \mathbf{P}(s)$, the positive entries of $\mathbf{Q}(s, s + (\mathcal{D} + 1)\Gamma)$ are at least $\prod_{k=0}^{(\mathcal{D}+1)\Gamma-1} \lambda_{i_{s+k}} \geq (1/m)^{(\mathcal{D}+1)\Gamma}$. \square

The above property of multi-hop transition matrices also implies the range of the push-sum weights $\mathbf{y}(t)$. In fact, the following lemma gives an upper and a lower bound on the push-sum weights $y_h(t)$ of all nodes $h \in \tilde{V}$ throughout the learning process. Note that, [31] also investigates the range of the push-sum weight $y_h(t)$ in their Lemma 11, where the lower bound of the weight on the virtual node is given as $m\alpha^2$. Compared to their bound, we give a lower bound α by looking into the physical communication process; it is much tighter than their bound $m\alpha^2$, as $\alpha = (1/m)^{(\mathcal{D}+1)\Gamma}$ is usually much smaller than $1/m$.

Lemma 8. *In AD-OGP, for any $t \in \{1, \dots, 2T\}$ and $i \in V$, it holds that*

$$m\alpha \leq y_i(t) \leq m.$$

Moreover, for any $l \in \{1, \dots, \mathcal{D}\}$, we have either $y_i^l(t) = 0$ or

$$\alpha \leq y_i^l(t) \leq m.$$

Proof. Recall our initialization of the push-sum weight on all nodes, namely $y_i(1) = 1$ and $y_i^l(1) = 0$ for all $i \in V, l \in \{1, \dots, \mathcal{D}\}$. Therefore, the stacked vector $\mathbf{y}(t)$ (which is constructed by stacking the parameter y of all nodes in the augmented graph at step t together, as we have defined before) can be expressed as

$$\mathbf{y}(t) = \mathbf{Q}(1, t)\mathbf{y}(1) = \mathbf{Q}(1, t) \begin{bmatrix} \mathbf{1}_m \\ \mathbf{0}_{m\mathcal{D}} \end{bmatrix}.$$

From Lemma 3, for any $t > 0, i \in \{1, \dots, m\}$, $\mathbf{P}_{ii}(t) > 0$, and $\mathbf{P}(t)$ is non-negative. Since $\mathbf{y}(1)$ is non-negative, $\mathbf{y}(t) = \mathbf{Q}(1, t)\mathbf{y}(1)$ is also non-negative. Moreover, since $\mathbf{P}(t)$ is column stochastic, i.e., $\mathbf{1}^\top \mathbf{P}(t) = \mathbf{1}^\top$, the sum of the parameters $y_h(t)$ of all nodes $h \in \tilde{V}$ at any step t is exactly

$$\sum_{h \in \tilde{V}} y_h(t) = \mathbf{1}^\top \mathbf{y}(t) = \mathbf{1}^\top \mathbf{Q}(1, t)\mathbf{y}(1) = \mathbf{1}^\top \mathbf{y}(0) = m.$$

Therefore, the parameter $y_i^l(t)$ for any node $b_i^l \in \tilde{V}$ (recall that, we use b_i^0 to represent the actual node $i \in V$) is at most m , which gives the common upper bound m for the push-sum weight $y_h(t)$ throughout the learning process.

To analyze the lower bound, notice that, for any actual node $i \in V$, we particularly have

$$y_i(t) = e_i^\top \prod_{s=0}^{t-1} \mathbf{P}(t-s) \mathbf{y}(1) \geq \prod_{s=0}^{t-1} \mathbf{P}_{ii}(s) \cdot y_i(1) > 0.$$

To further investigate the range of the push-sum weight $y_h(t)$, we consider the following two cases.

(i) When $t < (\mathcal{D} + 1)\Gamma$, the positive entries of $\mathbf{Q}(1, t)$ are at least $(1/m)^t \geq m\alpha$. Therefore, the positive entries of $\mathbf{y}(t)$ are at least $m\alpha$.

(ii) When $t \geq (\mathcal{D} + 1)\Gamma$, we can split the multi-hop transition matrix into $\mathbf{Q}(1, t) = \mathbf{Q}(t-s, t) \mathbf{Q}(1, t-s)$ where $s = (\mathcal{D} + 1)\Gamma$. From Lemma 7, the first m rows of $\mathbf{Q}(t-s, t)$ are positive, and its positive entries are at least α . Recall that $\mathbf{Q}(1, t-s)$ is column stochastic, i.e., $\sum_{j=1}^{m(D+1)} \mathbf{Q}_{jh}(1, t-s) = 1$ for any h . Therefore, for any $i \in V, h \in \tilde{V}$, $\mathbf{Q}_{ih}(1, t) = \sum_{j=1}^{m(D+1)} \mathbf{Q}_{jh}(1, t-s) \mathbf{Q}_{ij}(t-s, t)$ is a convex combination of entries in the i -th row of $\mathbf{Q}(t-s, t)$. Since all these entries are no less than α from Lemma 7, $\mathbf{Q}_{ih}(1, t)$ is also at least α . From the arbitrariness of $i \in V$ and $h \in \tilde{V}$, the entries in the first m rows of $\mathbf{Q}(1, t)$ are at least α . Therefore, for any actual node $i \in V$, we further have

$$y_i(t) \geq \sum_{j=1}^m \mathbf{Q}_{ij}(1, t) y_j(1) \geq m\alpha.$$

As for any virtual node $b_i^l, i \in V, l \in \{1, \dots, D\}$, we assume its parameter $y_i^l(t)$ is positive at some step t . Recall the transmission of message on the augment graph \tilde{G} (you can also see the demonstration in Figure 8). Specifically, for any virtual node $h = b_i^l \in \tilde{V} - V$, we have $y_i^l(t) > 0$ only when this node is carrying at least one delayed message that is sent to node i . We suppose some specific message currently carried by node b_i^l was sent by node $j \in V$ at some previous step s . From our communication strategy, the weight parameter captured in such message is exactly a γ_j -fraction of node j 's parameter $y_j(s)$. From the former case (i), the weight parameters carried by other messages are all positive. As we have also proved that $y_j(s) \geq m\alpha$, we directly have

$$y_i^l(t) \geq \gamma_j y_j(s) \geq \gamma_j m\alpha \geq \alpha.$$

Combining the above derived bounds, we successfully prove the lemma. \square

From Lemma 8, we know that the parameter y_i^l of any virtual node b_i^l is either zero, or restricted to the interval $[\alpha, m]$. The following lemma explores the first case: when $y_i^l(t) = 0$, its corresponding local model $\mathbf{x}_i^l(t)$ must also be zero.

Lemma 9. *For any virtual node $b_i^l, i \in V, l \in \{1, \dots, D\}$, if its push-sum weight $y_i^l(t) = 0$ at some step $t \in \{1, \dots, 2T\}$, then its local model parameters $\mathbf{x}_i^l(t) = \mathbf{0}$.*

Proof. In intuition, $y_i^l(t) = 0$ implies that the virtual node b_i^l is carrying no delayed message at step t . Since parameters \mathbf{x} and \mathbf{y} are coupled in each message, node b_i^l must not be carrying any local model as well, which means that its local model parameters $\mathbf{x}_i^l(t)$ is also zero.

In the following, we provide a strict proof of this lemma by inducting on t . The lemma obviously holds for $t = 1$, since the parameters of virtual nodes are initialized as $\mathbf{x}_i^l(1) = \mathbf{0}$ and $y_i^l(1) = 0$.

We suppose the lemma holds for some $t \geq 0$. From Lemma 8, the first m entries of $\mathbf{y}(t)$ are all positive. Therefore, for any $h \in \tilde{V}$, $y_h(t) = 0$ implies that the h -th row in $\mathbf{X}(t)$ equals to $\mathbf{0}$.⁵ Denote $\mathbf{X}'(t) = \mathbf{X}(t) + \mathbf{\Delta}(t)$. From the definition of $\mathbf{\Delta}(t)$, the entries in its last mD rows all equal to zero. Therefore, $y_h(t) = 0$ also implies that the h -th row of $\mathbf{X}'(t)$ equals to $\mathbf{0}$. Recall the recursion

$$\mathbf{X}(t+1) = \mathbf{P}(t) \mathbf{X}'(t), \quad \mathbf{y}(t+1) = \mathbf{P}(t) \mathbf{y}(t).$$

⁵Recall that, we slightly abuse the notations of any virtual node $h = b_i^l \in \tilde{V} - V$. Notably, its push-sum weight at step t can be represented by either $y_i^l(t)$ or $y_h(t)$ (here we mean $h = lm + i$), which refers to the h -th entry of $\mathbf{y}(t)$.

Suppose $y_h(t+1) = 0$ for an arbitrary $h \in \tilde{V} = \{1, \dots, m(D+1)\}$, namely,

$$y_h(t+1) = \sum_{k=1}^{m(D+1)} \mathbf{P}_{hk}(t) y_k(t) = 0.$$

Since the entries of $\mathbf{P}(t)$ and $\mathbf{y}(t)$ are all non-negative, for any $k \in \tilde{V}$ such that $y_k(t) > 0$, we must have $\mathbf{P}_{hk}(t) = 0$. For simplicity, we use \mathbf{A}_h to denote the h -th row of any given matrix \mathbf{A} . Therefore, the h -th row of $\mathbf{X}(t+1)$ can be expressed as

$$\mathbf{X}_{h\cdot}(t+1) = \sum_{k=1}^{m(D+1)} \mathbf{P}_{hk}(t) \mathbf{X}_k(t) = \sum_{k: y_k(t)=0} \mathbf{P}_{hk}(t) \mathbf{X}_k(t).$$

Since $y_k(t) = 0$ leads to $\mathbf{X}_k(t) = \mathbf{0}$, we directly have $\mathbf{X}_{h\cdot}(t+1) = \mathbf{0}$. Hence the lemma also holds at the next step $t+1$. \square

With the above preparations, we can now begin to analyze the quantity $R_i(s, t)$ in Lemma 6, which is a crucial term for regret analysis. Recall that $R_i(s, t)$ is defined as

$$R_i(s, t) = \left\| \frac{\mathbf{e}_i^\top \mathbf{Q}(s, t) \mathbf{\Delta}(s)}{\mathbf{e}_i^\top \mathbf{Q}(s, t) \mathbf{y}(s)} - \frac{\epsilon_s}{m} \right\|.$$

As discussed before, we can imagine a $(t-s)$ -step distributed averaging process operated on the augmented graph \tilde{G} . In the averaging process, the parameters of all nodes are initialized as $\tilde{\mathbf{X}}(1) = \mathbf{\Delta}(s)$ and $\tilde{\mathbf{y}}(1) = \mathbf{y}(s)$. The initial predictions are $\tilde{\mathbf{w}}(1)$ such that $\tilde{w}_i(1) = \tilde{x}_i(1)/\tilde{y}_i(1)$. At each step $k = 1, \dots, \tau$ (here $\tau = t-s$), each node mixes its parameters with other nodes according to the transition matrix $\mathbf{P}(s+k-1)$, i.e.,

$$\begin{aligned} \tilde{\mathbf{X}}(k+1) &= \mathbf{P}(s+k-1) \tilde{\mathbf{X}}(k), \\ \tilde{\mathbf{y}}(k+1) &= \mathbf{P}(s+k-1) \tilde{\mathbf{y}}(k). \end{aligned}$$

At the end of the process, the parameters of node i are exactly $\tilde{x}_i(\tau+1) = \mathbf{e}_i^\top \mathbf{Q}(s, t) \mathbf{\Delta}(s)$ and $\tilde{y}_i(\tau+1) = \mathbf{e}_i^\top \mathbf{Q}(s, t) \mathbf{y}(s)$, and its prediction $\tilde{w}_i(\tau+1) = \tilde{x}_i(\tau+1)/\tilde{y}_i(\tau+1)$.

With the above notations, we now analyze the effect of the imaginary τ -step averaging process. As a prerequisite, we can verify that, the previously derived lemmas on model parameters, namely Lemma 8 and Lemma 9, also hold for \tilde{x} and \tilde{y} . Specifically, for Lemma 8, we directly utilize $\tilde{\mathbf{y}}(k) = \mathbf{y}(k-s)$; for Lemma 9, we know it holds for $k=1$ and then induct on $k=2, \dots, \tau$.

We first try to give a matrix form of the evolution of predictions, just like our previously derived matrix form (4) of the recursion of model parameters $\tilde{\mathbf{X}}(k)$ and $\tilde{\mathbf{y}}(k)$. To this end, we first give a definition of the predictions $\tilde{w}_i^l(k)$ on virtual nodes b_i^l . However, a problem will arise in the formulation, namely, some $\tilde{y}_i^l(k)$ may be zero during the averaging process, so directly computing $\tilde{x}_i^l(k)/\tilde{y}_i^l(k)$ does not necessarily make sense. To tackle this issue, from Lemma 9 we know that $\tilde{y}_i^l(k) = 0$ gives $\tilde{x}_i^l(k) = \mathbf{0}$, we can thus set its prediction $\tilde{w}_i^l(k) = \mathbf{0}$ if $\tilde{y}_i^l(k) = 0$. Now the predictions of all nodes in \tilde{V} are well-defined. For a more compact expression, we define the vector $\tilde{\mathbf{y}}^{-1}(k) = (\tilde{y}_1^{-1}(k), \dots, \tilde{y}_{m(D+1)}^{-1}(k)) \in \mathbb{R}^{m(D+1)}$ as

$$\tilde{y}_h^{-1}(k) = \begin{cases} 1/\tilde{y}_h(k), & \tilde{y}_h(k) > 0, \\ 0, & \tilde{y}_h(k) = 0. \end{cases}$$

For any $1 \leq l \leq k \leq \tau$, the matrix form of predictions $\tilde{\mathbf{W}}(k+1)$ can be expressed via $\tilde{\mathbf{X}}(l)$ as

$$\begin{aligned} \tilde{\mathbf{W}}(k+1) &= \text{diag}(\tilde{\mathbf{y}}^{-1}(k+1)) \tilde{\mathbf{X}}(k+1) \\ &= \text{diag}(\tilde{\mathbf{y}}^{-1}(k+1)) \mathbf{P}(s+k-1) \tilde{\mathbf{X}}(k) \\ &= \dots \\ &= \text{diag}(\tilde{\mathbf{y}}^{-1}(k+1)) \mathbf{Q}(s+l-1, s+k) \tilde{\mathbf{X}}(l), \end{aligned}$$

where $\text{diag}(\mathbf{u})$ denotes the diagonal matrix, whose (i, i) -th entry takes u_i for all i , and other entries are zero.

From the above expression, we can directly derive the recursive function of the stacked predictions. The following lemma gives an arbitrary-order recursive relation on $\tilde{\mathbf{W}}(k)$.

Lemma 10. *In the above push-sum averaging process, for any $1 \leq l \leq k \leq \tau + 1$, the stacked predictions satisfy the recursive relation*

$$\tilde{\mathbf{W}}(k) = \mathbf{A}(l, k)\tilde{\mathbf{W}}(l),$$

where the (multi-hop) transition matrix $\mathbf{A}(l, k)$ for decisions is defined as

$$\mathbf{A}(l, k) = \text{diag}(\tilde{\mathbf{y}}^{-1}(k))\mathbf{Q}(s+l-1, s+k-1)\text{diag}(\tilde{\mathbf{y}}(l))$$

when $l < k$, and $\mathbf{A}(l, k) = I_{m(D+1)}$ when $l = k$.

Proof. Since the lemma trivially holds for $l = k$, we only need to investigate the case of $l < k$. We first prove that $\tilde{\mathbf{X}}(k) = \text{diag}(\tilde{\mathbf{y}}(k))\tilde{\mathbf{W}}(k)$ for any $k \in \{1, \dots, \tau + 1\}$. Specifically, we check each row of the matrices on both sides. Recall our definition of predictions $\tilde{\mathbf{w}}_i^l(k)$ on virtual nodes b_i^l . For any $h \leq m(D+1)$, the h -th row of $\tilde{\mathbf{W}}(k)$ is exactly (remind that, we use $\mathbf{A}_{i\cdot} = \mathbf{e}_i^\top \mathbf{A}$ to represent the i -th row vector of any matrix \mathbf{A})

$$\tilde{\mathbf{W}}_{h\cdot}(k) = \begin{cases} \tilde{\mathbf{X}}_{h\cdot}(k)/\tilde{y}_h(k), & \tilde{y}_h(k) > 0, \\ \mathbf{0}, & \tilde{y}_h(k) = 0. \end{cases}$$

From Lemma 9, when $\tilde{y}_h(k) = 0$ for some virtual node, its corresponding local model parameters $\tilde{\mathbf{x}}_h(k)$ always satisfies $\tilde{\mathbf{X}}_{h\cdot}(k) = \mathbf{0}$. Besides, from Lemma 8, the push-sum weight parameter $y_i(k)$ of any actual node $i \in V$ is positive. Therefore, it is easy to check that, for any $h \in \tilde{V} = \{1, \dots, m(D+1)\}$, the h -th row of $\tilde{\mathbf{X}}(k)$ satisfies

$$\tilde{\mathbf{X}}_{h\cdot}(k) = \tilde{\mathbf{W}}_{h\cdot}(k) \cdot \tilde{y}_h(k).$$

From the arbitrariness of $h \in \tilde{V}$, we derive $\tilde{\mathbf{X}}(k) = \text{diag}(\tilde{\mathbf{y}}(k))\tilde{\mathbf{W}}(k)$. Plugging $\tilde{\mathbf{X}}(l) = \text{diag}(\tilde{\mathbf{y}}(l))\tilde{\mathbf{W}}(l)$ into the above relation $\tilde{\mathbf{X}}(k) = \text{diag}(\tilde{\mathbf{y}}(k))\tilde{\mathbf{W}}(k)$, we prove the lemma. \square

We now focus on the transition matrix $\mathbf{A}(l, k)$ of predictions, which is determined by the multi-hop communication matrix $\mathbf{Q}(s+l-1, s+k-1)$ induced by push-sum and the stacked weights $\tilde{\mathbf{y}}(l), \tilde{\mathbf{y}}(k)$, as suggested in the above lemma. We already know that $\mathbf{Q}(s, t)$ is non-negative and column stochastic, and expect $\mathbf{A}(l, k)$ to inherit some properties from $\mathbf{Q}(s, t)$. In addition, the transition matrix $\mathbf{A}(l, k)$ of predictions is analogous to the mixing matrices \mathbf{P} in other decentralized online algorithms [36, 38]. However, their mixing matrices \mathbf{P} are conventionally doubly stochastic (i.e., $\mathbf{P}\mathbf{1} = \mathbf{P}^\top\mathbf{1} = \mathbf{1}$), whereas our transition matrices are not. Therefore, their analysis which relies on the doubly stochastic property and spectral gaps can not be directly applied to our setting.

To tackle this issue, we introduce the technique in [25], but make a refinement on the analysis which *saves a factor of \sqrt{n}* in the convergence rate, where n is the dimension of the decision set. Specifically, we utilize the internal structure of incremental term $\Delta(s)$, which is a rank-one matrix, namely $\Delta(s) = \mathbf{e}_{i_s}\epsilon_s$, since there is only a single learner performing local update at each step. This enables us to tackle the n dimensions of ϵ_s simultaneously by directly analyzing its $L2$ -norm $\|\epsilon_s\|$ (in our analysis, $\|\cdot\|$ represents the $L2$ -norm). In contrast, [25] tackles each of the n dimensions separately. In their approach, the k -th dimension contributes a coefficient $|\epsilon_{s,k}|$ ($\epsilon_{s,k}$ is the k -th entry of ϵ_s) to the convergence rate; summing them together gives a bound $\|\epsilon_s\|_1$ in terms of $L1$ -norm (see Corollary 1 in [25]). To transform it into $L2$ -norm, one need to utilize $\|\epsilon_s\|_1 \leq \sqrt{n}\|\epsilon_s\|$, which will incur an extra coefficient \sqrt{n} as the expense of expressing the final bound in terms of $L2$ -norm $\|\epsilon_s\|$.

More specifically, since $\Delta(s) = \mathbf{e}_{i_s}\epsilon_s$, the quantity $R_i(s, t)$ can be written as

$$R_i(s, t) = \left\| \frac{\mathbf{e}_i^\top \mathbf{Q}(s, t) \mathbf{e}_{i_s} \epsilon_s}{\mathbf{e}_i^\top \mathbf{Q}(s, t) \mathbf{y}(s)} - \frac{\epsilon_s}{m} \right\|.$$

Since both $\mathbf{e}_i^\top \mathbf{Q}(s, t) \mathbf{e}_{i_s}$ and $\mathbf{e}_i^\top \mathbf{Q}(s, t) \mathbf{y}(s)$ are scalars, we further have

$$R_i(s, t) \leq \left| \frac{\mathbf{e}_i^\top \mathbf{Q}(s, t) \mathbf{e}_{i_s}}{\mathbf{e}_i^\top \mathbf{Q}(s, t) \mathbf{y}(s)} - \frac{1}{m} \right| \|\epsilon_s\|. \quad (11)$$

Concretely, we investigate the effect of transition matrix $\mathbf{A}(l, k)$ on an arbitrary vector $\mathbf{u} \in \mathbb{R}^{m(D+1)}$, which satisfies $\mathbf{u}_h = 0$ for any $h \in \tilde{V}$ such that $\tilde{y}_h(k) = 0$ (then we can set $\mathbf{u} = \text{diag}(\tilde{\mathbf{y}}^{-1}(1))\mathbf{e}_{i_s} =$

$e_{i_s}/y_{i_s}(s)$ to restore the averaging process). Enlightened by the analysis in [31], we use the index set $I^k = \{h \leq m(D+1) \mid \tilde{y}_h(k) > 0\}$ to denote the positive entries of $\tilde{\mathbf{y}}(k)$. Lemma 8 implies that $V \subset I^k$. From Lemma 9, only those nodes inside I^k will give valid predictions (in fact, for any other node $h \notin I^k$, its parameters $\tilde{\mathbf{x}}_h(k), \tilde{y}_h(k)$ are all zero, which cannot give any meaningful predictions). Therefore, we are mainly interested in the nodes inside I^k at each step k . In this spirit, we call $\{u_j \mid j \in I^l\}$ the *valid entries* of $\mathbf{u} = (u_1, \dots, u_{m(D+1)})$ at step l . Denote $\mathbf{v} = \mathbf{A}(l, k)\mathbf{u}$, we similarly call $\{v_j \mid j \in I^k\}$ the valid entries of \mathbf{v} at step k . We can now give some general properties of $\mathbf{A}(l, k)$ as below.

Lemma 11. *For all $1 \leq l \leq k \leq \tau + 1$, the transition matrix $\mathbf{A}(l, k)$ of decisions satisfies:*

(a) *If $i \notin I^k$ or $j \notin I^l$, then $\mathbf{A}_{ij}(l, k) = 0$.*

(b) *The entries in the h -th row of $\mathbf{A}(l, k)$ sum up to*

$$\sum_{j=1}^{m(D+1)} \mathbf{A}_{hj}(l, k) = \sum_{j \in I^l} \mathbf{A}_{hj}(l, k) = \begin{cases} 1, & h \in I^k, \\ 0, & h \notin I^k. \end{cases}$$

Proof. (a) From the definition of $\mathbf{A}(l, k)$ in Lemma 10, its (i, j) -th entry

$$\mathbf{A}_{ij}(l, k) = \tilde{y}_i^{-1}(k) \mathbf{Q}_{ij}(s+l-1, s+k-1) \tilde{y}_j(l),$$

which equals to zero if $i \notin I^k$ or $j \notin I^l$.

(b) Since $\tilde{\mathbf{y}}(k)$ denotes the parameter after $k-1$ push-sum steps starting from step s , we actually have $\tilde{\mathbf{y}}(k) = \mathbf{y}(s+k-1)$. Therefore, the entries in the h -th row of $\mathbf{A}(l, k)$ sum up to

$$\begin{aligned} \mathbf{e}_h^\top \mathbf{A}(l, k) \mathbf{1} &= \mathbf{e}_h^\top \text{diag}(\tilde{\mathbf{y}}^{-1}(k)) \mathbf{Q}(s+l-1, s+k-1) \tilde{\mathbf{y}}(l) \\ &= \tilde{y}_h^{-1}(k) \cdot \mathbf{e}_h^\top \mathbf{Q}(s+l-1, s+k-1) \tilde{\mathbf{y}}(l) \\ &= \tilde{y}_h^{-1}(k) \cdot \mathbf{e}_h^\top \mathbf{Q}(s+l-1, s+k-1) \mathbf{y}(s+l-1) \\ &= \tilde{y}_h^{-1}(k) \cdot \mathbf{e}_h^\top \mathbf{y}(s+k-1) = \tilde{y}_h^{-1}(k) \cdot \tilde{y}_h(k) \\ &= \begin{cases} 1, & h \in I^k, \\ 0, & h \notin I^k. \end{cases} \end{aligned}$$

Moreover, we have $\mathbf{A}_{hj}(l, k) = 0$ if $j \notin I^l$. Therefore, for any $h \in I^k$, it holds that

$$1 = \mathbf{A}_h(l, k) \mathbf{1} = \sum_{j=1}^{m(D+1)} \mathbf{A}_{hj}(l, k) = \sum_{j \in I^l} \mathbf{A}_{hj}(l, k),$$

which proves the lemma. \square

Let's return to $\mathbf{v} = \mathbf{A}(l, k)\mathbf{u}$. The following lemma shows that, the transition matrix of predictions induced by push-sum steps will shrink the range of the valid entries of a vector, which is actually the core step to the analysis of push-sum.

Lemma 12. *Let $\mathbf{u} = (u_1, \dots, u_{m(D+1)})$ be an arbitrary vector in $\mathbb{R}^{m(D+1)}$. Set $\mathbf{v} = \mathbf{A}(l, k)\mathbf{u}$ and denote $\mathbf{v} = (v_1, \dots, v_{m(D+1)})$. The index set of valid entries $I^k = \{h : \tilde{y}_h(k) > 0\}$ is defined above. Then it holds that*

$$\max_{h \in I^k} v_h - \min_{h \in I^k} v_h \leq \max_{h \in I^l} u_h - \min_{h \in I^l} u_h.$$

Proof. For any $h \in I^k$, the h -th entry of \mathbf{v} takes

$$v_h = \sum_{j=1}^{m(D+1)} \mathbf{A}_{hj}(l, k) u_j = \sum_{j \in I^l} \mathbf{A}_{hj}(l, k) u_j.$$

From Lemma 10 and Lemma 11 we know that $\mathbf{A}_{hj}(l, k) \geq 0, \forall j \in I^l$ and $\sum_{j \in I^l} \mathbf{A}_{hj}(l, k) = 1$, respectively. Hence v_h is a convex combination of the elements in $\{u_j \mid j \in I^l\}$. Consequently, we have

$$\min_{j \in I^l} u_j \leq v_h \leq \max_{j \in I^l} u_j.$$

Since h is arbitrarily chosen from I^k , we have

$$\min_{j \in I^l} u_j \leq \min_{h \in I^k} u_h \leq \max_{h \in I^k} u_h \leq \max_{j \in I^l} u_j,$$

which proves the lemma. \square

The above lemma implies that, any push-sum step will not expand the range of the entries in any given vector. To show the convergence to consensus, we expect that, after operating sufficient number of push-sum steps, the *range* of a vector will be scaled down. Such convergence effect is proved in the following lemma, in which the minimal number of steps $(\mathcal{D} + 1)\Gamma$ stems from Lemma 7, namely the minimal number of hops $t - s$ for $\mathbf{Q}(s, t)$ to guarantee the positiveness of its first m rows. We note that, a similar approach to derive the convergence effect can be seen in the analysis of symmetric communication [36, 21], which relies on the double stochasticity of the transition matrices. However, their methods are unsuitable to our setting, where $\mathbf{A}(l, k)$ is actually not doubly stochastic.

As a prerequisite, we derive some properties of any $r = (\mathcal{D} + 1)\Gamma$ -hop transition matrix $\mathbf{A}(l, l + r)$ in the lemma below. Remind that, the vertex set and the augmented vertex set are specified as $V = \{1, \dots, m\}$ and $\tilde{V} = \{1, \dots, m(\mathcal{D} + 1)\}$, respectively.

Lemma 13. *Set $r = (\mathcal{D} + 1)\Gamma$. For any $l \geq 1$, the transition matrix $\mathbf{A}(l, l + r)$ satisfies:*

(a) *The first $m \times m$ block of $\mathbf{A}(l, l + r)$ is positive, namely $\mathbf{A}_{ij}(l, l + r) > 0, \forall i, j \in V$.*

(b) *The positive entries in the first m columns of $\mathbf{A}(l, l + r)$ are at least α^2 . Moreover, the positive entries in the last $m\mathcal{D}$ columns of $\mathbf{A}(l, l + r)$ are at least α^2/m .*

(c) *For any $m + 1 \leq h \leq m(\mathcal{D} + 1)$, if $h \in I^{l+r}$, then there must exist some $i \in V$ such that $\mathbf{A}_{hi}(l, l + r) \geq \alpha^2$.*

Proof. From Lemma 10, for any $i, j \in \tilde{V}$, $\mathbf{A}_{ij}(l, l + r) = \tilde{y}_i^{-1}(l + r)\mathbf{Q}_{ij}(s + l - 1, s + l + r - 1)\tilde{y}_j(l)$.

(a) When $i, j \in V$, from Lemma 8, $\tilde{y}_i(l + r), \tilde{y}_j(l) > 0$; also from Lemma 7, $\mathbf{Q}_{ij}(s + l - 1, s + l + r - 1) > 0$. Therefore, $\mathbf{A}_{ij}(l, l + r) > 0$.

(b) Suppose $\mathbf{A}_{hi}(l, l + r) > 0$ for some $h \in \tilde{V}, i \in V$. Then we have $\tilde{y}_h(t + s) > 0$ and $\mathbf{Q}_{hi}(s + l - 1, s + l + r - 1) > 0$. From Lemma 8, we have $1/\tilde{y}_h(l + r) \geq 1/m, \tilde{y}_i(t) \geq m\alpha$; in addition, from Lemma 7, we have $\mathbf{Q}_{hi}(s + l - 1, s + l + r - 1) \geq \alpha$, therefore $\mathbf{A}_{hi}(l, l + r) \geq (1/m)\alpha(m\alpha) = \alpha^2$. Similarly, we can prove that, if $\mathbf{A}_{hj}(l, l + r) > 0$ for some $h \in \tilde{V}, j > m$, then $\mathbf{A}_{hj}(l, l + r) \geq (1/m)\alpha\alpha = \alpha^2/m$.

(c) From (b), we only need to prove that $\mathbf{A}_{hi}(l, l + r) > 0$ for some $i \in V$. Since $\mathbf{A}(l, l + r)$ is non-negative, and from Lemma 8 $\tilde{y}_i(l) > 0$ for all $i \in V$, it suffices to prove that $\mathbf{Q}_{hi}(s + l - 1, s + l + r - 1) > 0$ for some $i \in V$. For simplicity we denote $p = s + l - 1$. We now assume $\mathbf{Q}_{hi}(p, p + r) = 0$ for all $i \in V$. From the definition of multi-hop communication matrix, we have $\mathbf{Q}_{ii}(p + k, p + r) > 0$ for all $i \in V, k \leq r$. Since $\mathbf{Q}_{hi}(p, p + r) \geq \mathbf{Q}_{hi}(p, p + k)\mathbf{Q}_{ii}(p + k, p + r)$ for all $k \leq r$, we must have $\mathbf{Q}_{hi}(p, p + k) = 0$ for all $i \in V, k \leq r$, which means that the information on any actual node $i \in V$ will not reach the virtual node h during step l and $l + r$. Since the message delay is at most $D < r - 1$, then the node h will not capture any message at step $l + r$, which means $\tilde{y}_h(l + r) = 0$ and will lead to contradiction. \square

Given the above lemma, we now return to the relation $\mathbf{v} = \mathbf{A}(l, l + r)\mathbf{u}$. The following lemma shows that $\mathbf{A}(l, l + r)$ will shrink the range of the valid entries of \mathbf{u} . Then, after sufficient number of push-sum steps, all valid entries of the vector tend to reach a consensus (i.e., attain a fixed value).

Lemma 14. *Given $1 \leq l \leq k \leq \tau + 1$. Assume $\mathbf{u} = (u_1, \dots, u_{m(\mathcal{D} + 1)}) \in \mathbb{R}^{m(\mathcal{D} + 1)}$ satisfies $\tilde{y}_h(l) = 0 \Rightarrow u_h = 0, \forall h \in \tilde{V} = \{1, \dots, m(\mathcal{D} + 1)\}$. Define*

$$\beta(\mathbf{w}; p) = \max_{h \in I^p} w_h(p) - \min_{h \in I^p} w_h(p)$$

as the range of the valid entries of $\mathbf{w} = (w_1, \dots, w_{m(\mathcal{D} + 1)})$ at step p . Set $r = (\mathcal{D} + 1)\Gamma$, then we have

$$\beta(\mathbf{A}(l, k)\mathbf{u}; k) \leq (1 - m\alpha^4)^{\lfloor \frac{k-l}{2r} \rfloor} \beta(\mathbf{u}; l).$$

Proof. Denote $\mathbf{v} = \mathbf{A}(l, k)\mathbf{u}$ and $q = \lfloor (k-l)/r \rfloor$. We construct a sequence of vectors $\{\mathbf{u}(p) \mid 0 \leq p \leq q+1\}$, namely, $\mathbf{u}(0) = \mathbf{u}$, $\mathbf{u}(p) = \mathbf{A}(l + (p-1)r, l + pr)\mathbf{u}(p-1)$ for $p \in \{1, \dots, q\}$, and $\mathbf{u}(q+1) = \mathbf{A}(l + qr, k)\mathbf{u}(q)$. We now specify $\mathbf{u}(p) = (u_1(p), \dots, u_{m(D+1)}(p))$ for any $p \in \{1, \dots, q+1\}$.

As assumed in this lemma, $\tilde{y}_h(l) = 0$ leads to $u_h = 0$. From the definition of $\mathbf{A}(l, r)$ in Lemma 10, we can inductively prove that $\tilde{y}_h(l+pr) = 0$ leads to $u_h(p) = 0$, and further $\mathbf{u}(p) = \mathbf{A}(l, l+pr)\mathbf{u}(0)$ for all $p \in \{1, \dots, q\}$. Similarly, $\tilde{y}_h(k) = 0$ leads to $u_h(p+1) = 0$; and we also have $\mathbf{u}(p+1) = \mathbf{A}(l, k)\mathbf{u}(0) = \mathbf{v}$.

For any $0 \leq p \leq q$, denote $\beta_p = \beta(\mathbf{u}(p); l + pr)$. We now derive the recursive relation of β_p . Specifically, we introduce an auxiliary sequence $\{\gamma_p\}_{p \in \{0, \dots, q\}}$, i.e.,

$$\gamma_p = \max_{i \in V} u_i(p) - \min_{i \in V} u_i(p),$$

which measures the range of the first m -entries of $\mathbf{u}(p)$. From Lemma 8, the first m entries of $\tilde{\mathbf{y}}(k)$ must be positive for any k . Therefore, for any $p \in \{1, \dots, q\}$, we have $V \subset I^{l+pr}$ and further

$$\begin{aligned} \max_{i \in V} u_i(p) &\leq \max_{h \in I^{l+pr}} u_i(p), \\ \min_{i \in V} u_i(p) &\geq \min_{h \in I^{l+pr}} u_i(p), \end{aligned}$$

which implies $\gamma_p \leq \beta_p$. Now we investigate the relation between γ_p, β_p and $\gamma_{p-1}, \beta_{p-1}$.

(i) First we consider an arbitrary actual node $i \in V$. From Lemma 8, $\tilde{y}_i(l + pr) > 0$. We abbreviate $\mathbf{A}_{ih}(l + (p-1)r, l + pr)$ into c_h for $h \in \tilde{V}$. Since $\mathbf{A}(l + (p-1)r, l + pr)$ is non-negative, we have $c_h \geq 0$. From Lemma 11 we further have $\sum_{h \in \tilde{V}} c_h = 1$. Thus $u_i(p) = \sum_{h \in \tilde{V}} c_h u_h(p-1)$, which means that $u_i(p)$ is a convex combination of the entries in $\mathbf{u}(p-1)$. From Lemma 13 we also know that $c_j \geq \alpha^2$ for any $j \in V$. Denote the maximal valid entry of $\mathbf{u}(p-1)$ as $u' = \max_{h \in I^{l+(p-1)r}} u_h(p-1)$, then we have

$$\begin{aligned} u_i(p) &= \sum_{h=1}^{m(D+1)} c_h u_h(p-1) = \sum_{j=1}^m c_j u_j(p-1) + \sum_{h=m+1}^{m(D+1)} c_h u_h(p-1) \\ &\leq \sum_{j=1}^m c_j u_j(p-1) + u'(1 - \sum_{j=1}^m c_h) \leq \sum_{j=1}^m \alpha^2 u_j(p-1) + u'(1 - m\alpha^2), \end{aligned}$$

where the last step utilizes $c_j \geq \alpha^2$ and $u_j(p-1) \leq u'$ for $j \in V$. Similarly, denote the minimal valid entry of $\mathbf{u}(p-1)$ as $u'' = \min_{h \in I^{l+(p-1)r}} u_h(p-1)$, then we have

$$\begin{aligned} u_i(p) &= \sum_{h=1}^{m(D+1)} c_h u_h(p-1) = \sum_{j=1}^m c_j u_j(p-1) + \sum_{h=m+1}^{m(D+1)} c_h u_h(p-1) \\ &\geq \sum_{j=1}^m c_j u_j(p-1) + u''(1 - \sum_{j=1}^m c_h) \geq \sum_{j=1}^m \alpha^2 u_j(p-1) + u''(1 - m\alpha^2). \end{aligned}$$

Since the above two inequalities hold for all $i \in V$, we set $i = \arg \max_{j \in V} u_j(p)$ in the first inequality, and set $i = \arg \min_{j \in V} u_j(p)$ in the second inequality, to derive

$$\gamma_p = \max_{j \in V} u_j(p) - \min_{j \in V} u_j(p) \leq (1 - m\alpha^2)(u' - u'') = (1 - m\alpha^2)\beta_{p-1}.$$

(ii) Then we consider an arbitrary "valid" node $i' \in I^{l+pr}$, and denote $c_h = \mathbf{A}_{i'h}(l + (p-1)r, l + pr)$ for $h \in \tilde{V}$. Via a very similar deduction in (i), we can prove that $u_{i'}(p) = \sum_{h \in \tilde{V}} c_h u_h(p-1)$ is a convex combination of the entries in $\mathbf{u}(p-1)$. However, since c_h may locate at the last mD rows of $\mathbf{A}(l + (p-1)r, l + pr)$, from Lemma 13 we only have $c_j \geq \alpha^2$ for some $j \in V$, which leads to $\sum_{j \in V} c_j \geq \alpha^2$. Therefore, using the same notations u' and u'' as the first case (i) to denote the

maximal and minimal valid entry of $\mathbf{u}(p-1)$, respectively, we have

$$\begin{aligned} u_{i'}(p) &= \sum_{h=1}^{m(D+1)} c_h u_h(p-1) = \sum_{j=1}^m c_j u_j(p-1) + \sum_{h=m+1}^{m(D+1)} c_h u_h(p-1) \\ &\leq \sum_{j=1}^m c_j u_j(p-1) + u'(1 - \sum_{j=1}^m c_h) \leq \alpha^2 \max_{j \in V} u_j(p-1) + u'(1 - \alpha^2), \end{aligned}$$

and similarly

$$\begin{aligned} u_{i'}(p) &= \sum_{h=1}^{m(D+1)} c_h u_h(p-1) = \sum_{j=1}^m c_j u_j(p-1) + \sum_{h=m+1}^{m(D+1)} c_h u_h(p-1) \\ &\geq \sum_{j=1}^m c_j u_j(p-1) + u''(1 - \sum_{j=1}^m c_h) \geq \alpha^2 \min_{j \in V} u_j(p-1) + u''(1 - \alpha^2). \end{aligned}$$

Set $i' = \arg \max_{h \in I^{l+pr}} u_h(p)$ in the first inequality, and set $i' = \arg \min_{h \in I^{l+pr}} u_h(p)$ in the second inequality, then we have

$$\beta_p \leq \alpha^2 \gamma_{p-1} + (1 - \alpha^2) \beta_{p-1}.$$

From Lemma 12, we have $\beta_{p-1} \leq \beta_{p-2}$. Hence we can combine both cases (i) and (ii) to derive

$$\beta_p \leq \alpha^2 (1 - m\alpha^2) \beta_{p-2} + (1 - \alpha^2) \beta_{p-1} \leq (1 - m\alpha^4) \beta_{p-2},$$

which gives the recursive relation of β_p .

Now we can prove the desired lemma. Specifically, from the above recursion and the fact that $\beta_p \leq \beta_{p-1}, \forall p \in \{1, \dots, q\}$, we have

$$\begin{aligned} \beta(\mathbf{u}(q); l + qr) &= \beta_q \leq \beta_{2 \lfloor \frac{q}{2} \rfloor} \leq (1 - m\alpha^4) \beta_{2(\lfloor \frac{q}{2} \rfloor - 1)} \leq \dots \\ &\leq (1 - m\alpha^4)^{\lfloor \frac{q}{2} \rfloor} \beta_0 = (1 - m\alpha^4)^{\lfloor \frac{q}{2} \rfloor} \beta(\mathbf{u}, l), \end{aligned}$$

and consequently,

$$\beta(\mathbf{v}; k) \leq \beta(\mathbf{u}(q); l + qr) \leq (1 - m\alpha^4)^{\lfloor \frac{q}{2} \rfloor} \beta(\mathbf{u}, l).$$

Since $q = \lfloor (k - l)/r \rfloor$, we have $\lfloor q/2 \rfloor \leq \lfloor (k - l)/2r \rfloor$, which proves the lemma. \square

To relate the above lemma with the imaginary push-sum averaging process on $\tilde{\mathbf{X}}(1)$ and $\tilde{\mathbf{y}}(1)$, we set $l = 1, k = \tau + 1$, and $\mathbf{u} = \mathbf{e}_{i_s} / \tilde{y}_{i_s}(1)$ in the above lemma. Also denote $\mathbf{z} = \mathbf{A}(1, \tau + 1) \mathbf{e}_{i_s} / \tilde{y}_{i_s}(1) \in \mathbb{R}^{m(D+1)}$ and specify its entries as $\mathbf{z} = (z_1, \dots, z_{m(D+1)})$. Recall the inequality (11) where the bound of $R_i(s, t)$ is related with the term $|z_i - 1/m|$. In other words, we need to compare the quantity z_i with $1/m$. Notice that $\min_{h \in I^{\tau+1}} z_h \leq z_i \leq \max_{h \in I^{\tau+1}} z_h$. We now verify that $\min_{h \in I^{\tau+1}} z_h \leq 1/m \leq \max_{h \in I^{\tau+1}} z_h$.

Lemma 15. *The transition matrix of predictions $\mathbf{A}(1, \tau + 1)$ is defined as above. Set $\mathbf{z} = (z_1, \dots, z_{m(D+1)}) = \mathbf{A}(1, \tau + 1) \mathbf{e}_{i_s} / \tilde{y}_{i_s}(1)$, then we have*

$$\min_{h \in I^{\tau+1}} z_h \leq \frac{1}{m} \leq \max_{h \in I^{\tau+1}} z_h.$$

Proof. It suffices to prove that, $1/m$ is a convex combination of the elements in $\{z_h \mid h \in I^{\tau+1}\}$. To this end, we observe that, z_h can be viewed as the prediction of learner $h \in I^{\tau+1}$ after τ push-sum steps, where the scalar models and push-sum weights of each learner $k \in \tilde{V}$ are initialized as the k -th entry of \mathbf{e}_{i_s} and $\tilde{\mathbf{y}}(1)$, respectively. Concretely, from the definition of $\mathbf{A}(1, \tau + 1)$ in Lemma 10, we have

$$z_h = \frac{\mathbf{e}_h^\top \mathbf{Q}(s, t) \mathbf{e}_{i_s}}{\tilde{y}_h(\tau + 1)} = \frac{\mathbf{Q}_{h \cdot}(s, t) \mathbf{e}_{i_s}}{\tilde{y}_h(\tau + 1)}.$$

Recall that $\mathbf{1}^\top \mathbf{e}_{i_s} = 1$ and $\mathbf{1}^\top \tilde{\mathbf{y}}(\tau + 1) = \mathbf{1}^\top \tilde{\mathbf{y}}(t) = m$; in addition, $\mathbf{Q}(s, t)$ is column stochastic. Hence we have

$$\frac{1}{m} = \frac{\mathbf{1}^\top \mathbf{Q}(s, t) \mathbf{e}_{i_s}}{m} = \sum_{h \in \tilde{V}} \frac{\mathbf{Q}_{h \cdot}(s, t) \mathbf{e}_{i_s}}{m}.$$

From Lemma 9, for any node $h \in \tilde{V} - I^{\tau+1}$, its push-sum weight $\tilde{y}_h(\tau+1) = 0$ and local scalar model $\mathbf{Q}_{h\cdot}(s, t)\mathbf{e}_{i_s} = 0$. Therefore, we further have

$$\frac{1}{m} = \sum_{h \in I^{\tau+1}} \frac{\mathbf{Q}_{h\cdot}(s, t)\mathbf{e}_{i_s}}{m} = \sum_{h \in I^{\tau+1}} \frac{\tilde{y}_h(\tau+1)}{m} \cdot \frac{\mathbf{Q}_{h\cdot}(s, t)\mathbf{e}_{i_s}}{\tilde{y}_h(\tau+1)} = \sum_{h \in I^{\tau+1}} \frac{\tilde{y}_h(\tau+1)}{m} z_h,$$

which is a convex combination of $\{z_h \mid h \in I^{\tau+1}\}$, because $\tilde{y}_h(\tau+1) \geq 0$ and $\sum_{h \in I^{\tau+1}} \tilde{y}_h(\tau+1)/m = \sum_{h \in \tilde{V}} \tilde{y}_h(\tau+1)/m = 1$. Therefore, $1/m$ must lie between the maximal and the minimal values of the elements in $\{z_h \mid h \in I^{\tau+1}\}$, which proves the lemma. \square

Set $\lambda = 1 - m\alpha^4$, $B = (\mathcal{D} + 1)\Gamma = (\mathcal{D} + 1)(D^{msg} + \Gamma_d)$, and $\mathbf{u} = \mathbf{e}_{i_s}/\tilde{y}_{i_s}(1)$. From the above lemma, the term $|z_i - 1/m|$ can be bounded as

$$\begin{aligned} |z_i - \frac{1}{m}| &\leq \max_{h \in I^{\tau+1}} z_h - \min_{h \in I^{\tau+1}} z_h \\ &\leq \lambda^{\lfloor \frac{\tau}{2B} \rfloor} (\max_{h \in I^1} u_h - \min_{h \in I^1} u_h) = \frac{1}{\tilde{y}_{i_s}(1)} \lambda^{\lfloor \frac{\tau-s}{2B} \rfloor} = \frac{1}{y_{i_s}(s)} \lambda^{\lfloor \frac{\tau-s}{2B} \rfloor}. \end{aligned}$$

Now, we successfully derive a bound for the term $R_i(s, t)$ for any $i \in V, s \in \mathcal{Q}_T, t \in \{s, s+1, \dots, 2T\}$, i.e.,

$$R_i(s, t) = |z_i - \frac{1}{m}| \|\boldsymbol{\epsilon}_s\| \leq \lambda^{\lfloor \frac{\tau-s}{2B} \rfloor} \frac{\|\boldsymbol{\epsilon}_s\|}{y_{i_s}(s)}.$$

Moreover, recall our definition of the incremental term $\boldsymbol{\epsilon}_s = -\eta \mathbf{g}_s + \mathbf{r}_s$. From Lemma 4, the residual term $\|\mathbf{r}_s\| \leq \eta g_s$. Therefore, for any $i \in V, s \in \mathcal{Q}_T, t \in \{s, s+1, \dots, 2T\}$, we have

$$R_i(s, t) \leq 2\eta \lambda^{\lfloor \frac{\tau-s}{2B} \rfloor} \frac{g_s}{y_{i_s}(s)}.$$

Finally, we plug the above inequality of $R_i(s, t)$ into the general regret bound (10), then derive the final bound in the main theorem.

E Omitted Proof of Corollary 1

In this section, we provide the detailed proof of Corollary 1 in our main paper.

Starting from the regret bound in Theorem 1. We analyze each regret term under the bounded gradient norm assumption separately. The first term is preserved. For the second term, we have

$$\sum_{t \in \mathcal{Q}_T} \sum_{s \in \mathcal{Q}_{l_t, t+1}} g_t g_s \leq G^2 \sum_{t \in \mathcal{Q}_T} |\mathcal{Q}_{l_t, t+1}|.$$

We first focus on the elements in $\mathcal{Q}_{l_t, t+1}$, which is defined as $\{l_t < s < t+1 \mid s \in \mathcal{Q}_T\}$ in our main paper. Specifically, for any $t \in \mathcal{Q}_T$, recall our definition of processing delay $d^p(t) = |\mathcal{Q}_{l_t, t}|$, we directly have

$$|\mathcal{Q}_{l_t, t+1}| = |\mathcal{Q}_{l_t, t} \cup \{t\}| = d^p(t) + 1.$$

Summing the above equation over $t \in \mathcal{Q}_T$, we thus have

$$\sum_{t \in \mathcal{Q}_T} |\mathcal{Q}_{l_t, t+1}| = \left(\sum_{t \in \mathcal{Q}_T} d^p(t) \right) + |\mathcal{Q}_T| = D^{proc} + T.$$

To analyze the last three terms in the regret bound, we first utilize Lemma 8 which gives $1/y_{i_s}(t) \leq 1/m\alpha$ for any $t \in \mathcal{Q}_T$. Then we investigate the geometric series of λ . Specifically, for any $t \in \mathcal{Q}_T$, we have

$$\begin{aligned} \sum_{s \in \mathcal{Q}_{1, t+1}} \lambda^{\lfloor \frac{t-s}{2B} \rfloor} &\leq \sum_{s=1}^t \lambda^{\lfloor \frac{t-s}{2B} \rfloor} = \sum_{s=0}^{t-1} \lambda^{\lfloor \frac{s}{2B} \rfloor} \leq \sum_{s=0}^{T-1} \lambda^{\lfloor \frac{s}{2B} \rfloor} \leq \sum_{s=0}^{2B \lceil \frac{T}{2B} \rceil - 1} \lambda^{\lfloor \frac{s}{2B} \rfloor} \\ &= \sum_{t=0}^{\lceil \frac{T}{2B} \rceil - 1} \sum_{s=0}^{2B-1} \lambda^{\lfloor \frac{2tB+s}{2B} \rfloor} = \sum_{t=0}^{\lceil \frac{T+1}{2B} \rceil - 1} 2B \lambda^t \leq \frac{2B}{1-\lambda}. \end{aligned}$$

Therefore, the last three terms in the regret bound in Theorem 1 can be bounded by

$$8\eta \sum_{t \in \mathcal{Q}_T} \frac{2B}{1-\lambda} \frac{G^2}{m\alpha} + \eta \sum_{t \in \mathcal{Q}_T} \frac{G^2}{m\alpha} = \frac{\eta}{m} \left(\frac{16B}{m\alpha^5} + \frac{1}{\alpha} \right) G^2 |\mathcal{Q}_T| \leq \frac{2\eta}{m} \left(\frac{8}{\alpha^5} + \frac{1}{\alpha} \right) BG^2 T.$$

Combining the above analysis for each regret term, we prove this corollary.

References

- [1] Mohammad Akbari, Bahman Ghahsifard, and Tamás Linder. Distributed online convex optimization on time-varying directed graphs. *IEEE Transactions on Control of Network Systems*, 4(3):417–428, 2015.
- [2] Mahmoud Assran and Michael Rabbat. Asynchronous gradient-push. *IEEE Transactions on Automatic Control*, 2020.
- [3] Dimitri P Bertsekas and John N Tsitsiklis. *Parallel and distributed computation: numerical methods*, volume 23. Prentice hall Englewood Cliffs, NJ, 1989.
- [4] Xuanyu Cao, Junshan Zhang, and H Vincent Poor. Constrained online convex optimization with feedback delays. *IEEE Transactions on Automatic Control*, 2020.
- [5] Nicolo Cesa-Bianchi, Yoav Freund, David Haussler, David P Helmbold, Robert E Schapire, and Manfred K Warmuth. How to use expert advice. *Journal of the ACM (JACM)*, 44(3):427–485, 1997.
- [6] Igor Colin, Aurelien Bellet, Joseph Salmon, and Stéphan Cléménçon. Gossip dual averaging for decentralized optimization of pairwise functions. In *International Conference on Machine Learning*, pages 1388–1396, 2016.
- [7] John C Duchi, Alekh Agarwal, and Martin J Wainwright. Dual averaging for distributed optimization: Convergence analysis and network scaling. *IEEE Transactions on Automatic control*, 57(3):592–606, 2011.
- [8] Hamid Reza Feyzmahdavian, Arda Aytekin, and Mikael Johansson. An asynchronous mini-batch algorithm for regularized stochastic optimization. *IEEE Transactions on Automatic Control*, 61(12):3740–3754, 2016.
- [9] Christoforos N Hadjicostis and Themistoklis Charalambous. Average consensus in the presence of delays in directed graph topologies. *IEEE Transactions on Automatic Control*, 59(3):763–768, 2013.
- [10] Elad Hazan et al. Introduction to online convex optimization. *Foundations and Trends® in Optimization*, 2(3-4):157–325, 2016.
- [11] Elad Hazan and Haipeng Luo. Variance-reduced and projection-free stochastic optimization. In *International Conference on Machine Learning*, pages 1263–1271, 2016.
- [12] Hadrien Hendrikx, Francis Bach, and Laurent Massoulié. Accelerated decentralized optimization with local updates for smooth and strongly convex objectives. In *AISTATS 2019-22nd International Conference on Artificial Intelligence and Statistics*, 2019.
- [13] Saghar Hosseini, Airlie Chapman, and Mehran Mesbahi. Online distributed optimization via dual averaging. In *52nd IEEE Conference on Decision and Control*, pages 1484–1489. IEEE, 2013.
- [14] Yu-Guan Hsieh, Franck Iutzeler, Jérôme Malick, and Panayotis Mertikopoulos. Multi-agent online optimization with delays: Asynchronicity, adaptivity, and optimism. *arXiv preprint arXiv:2012.11579*, 2020.
- [15] Alec Koppel, Felicia Y Jakubiec, and Alejandro Ribeiro. A saddle point algorithm for networked online convex optimization. *IEEE Transactions on Signal Processing*, 63(19):5149–5164, 2015.
- [16] Alec Koppel, Santiago Paternain, Cédric Richard, and Alejandro Ribeiro. Decentralized online learning with kernels. *IEEE Transactions on Signal Processing*, 66(12):3240–3255, 2018.
- [17] Guanghui Lan, Soomin Lee, and Yi Zhou. Communication-efficient algorithms for decentralized and stochastic optimization. *Mathematical Programming*, pages 1–48, 2017.
- [18] Jinlong Lei, Peng Yi, Yiguang Hong, Jie Chen, and Guodong Shi. Online convex optimization over erdos-renyi random networks. *Advances in Neural Information Processing Systems*, 33, 2020.
- [19] Bingcong Li, Tianyi Chen, and Georgios B Giannakis. Bandit online learning with unknown delays. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 993–1002. PMLR, 2019.

- [20] Jueyou Li, Chuanye Gu, and Zhiyou Wu. Online distributed stochastic learning algorithm for convex optimization in time-varying directed networks. *Neurocomputing*, 416:85–94, 2020.
- [21] Xiangru Lian, Wei Zhang, Ce Zhang, and Ji Liu. Asynchronous decentralized parallel stochastic gradient descent. In *International Conference on Machine Learning*, pages 3049–3058, 2018.
- [22] Brendan McMahan and Matthew Streeter. Delay-tolerant algorithms for asynchronous distributed online learning. In *Advances in Neural Information Processing Systems*, pages 2915–2923, 2014.
- [23] Aryan Mokhtari and Alejandro Ribeiro. Dsa: Decentralized double stochastic averaging gradient algorithm. *The Journal of Machine Learning Research*, 17(1):2165–2199, 2016.
- [24] Angelia Nedić, Soomin Lee, and Maxim Raginsky. Decentralized online optimization with global objectives and local communication. In *2015 American Control Conference (ACC)*, pages 4497–4503. IEEE, 2015.
- [25] Angelia Nedić and Alex Olshevsky. Stochastic gradient-push for strongly convex functions on time-varying directed graphs. *IEEE Transactions on Automatic Control*, 61(12):3936–3947, 2016.
- [26] Angelia Nedic and Asuman Ozdaglar. Distributed subgradient methods for multi-agent optimization. *IEEE Transactions on Automatic Control*, 54(1):48, 2009.
- [27] Steffen Rendle. Factorization machines. In *2010 IEEE International Conference on Data Mining*, pages 995–1000. IEEE, 2010.
- [28] S Shahrapour and A Jadbabaie. Distributed online optimization in dynamic environments using mirror descent. *IEEE Transactions on Automatic Control*, 63(3):714–725, 2018.
- [29] Wei Shi, Qing Ling, Gang Wu, and Wotao Yin. Extra: An exact first-order algorithm for decentralized consensus optimization. *SIAM Journal on Optimization*, 25(2):944–966, 2015.
- [30] Wei Shi, Qing Ling, Gang Wu, and Wotao Yin. A proximal gradient algorithm for decentralized composite optimization. *IEEE Transactions on Signal Processing*, 63(22):6013–6023, 2015.
- [31] Artin Spiridonoff, Alex Olshevsky, and Ioannis Ch Paschalidis. Robust asynchronous stochastic gradient-push: Asymptotically optimal and network-independent performance for strongly convex functions. *Journal of Machine Learning Research*, 21(58):1–47, 2020.
- [32] Konstantinos I Tsianos and Michael G Rabbat. Efficient distributed online prediction and stochastic optimization with approximate distributed averaging. *IEEE Transactions on Signal and Information Processing over Networks*, 2(4):489–506, 2016.
- [33] Yuanyu Wan, Wei-Wei Tu, and Lijun Zhang. Projection-free distributed online convex optimization with $o(\sqrt{t})$ communication complexity. In *International Conference on Machine Learning*, pages 9818–9828. PMLR, 2020.
- [34] Duncan J Watts and Steven H Strogatz. Collective dynamics of ‘small-world’ networks. *nature*, 393(6684):440, 1998.
- [35] Tianyu Wu, Kun Yuan, Qing Ling, Wotao Yin, and Ali H Sayed. Decentralized consensus optimization with asynchrony and delays. *IEEE Transactions on Signal and Information Processing over Networks*, 4(2):293–307, 2017.
- [36] Feng Yan, Shreyas Sundaram, SVN Vishwanathan, and Yuan Qi. Distributed autonomous online learning: Regrets and intrinsic privacy-preserving properties. *IEEE Transactions on Knowledge and Data Engineering*, 25(11):2483–2493, 2013.
- [37] Xinlei Yi, Xiuxian Li, Lihua Xie, and Karl H Johansson. Distributed online convex optimization with time-varying coupled inequality constraints. *IEEE Transactions on Signal Processing*, 68:731–746, 2020.
- [38] Wenpeng Zhang, Peilin Zhao, Wenwu Zhu, Steven CH Hoi, and Tong Zhang. Projection-free distributed online learning in networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 4054–4062. JMLR. org, 2017.
- [39] Yawei Zhao, Chen Yu, Peilin Zhao, and Ji Liu. Decentralized online learning: Take benefits from others’ data without sharing your own to track global trend. *arXiv preprint arXiv:1901.10593*, 2019.

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
 - (b) Did you describe the limitations of your work? [Yes]
 - (c) Did you discuss any potential negative societal impacts of your work? [No] We study a mathematical problem. The societal impact largely depends on specific real-world application scenarios.
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [Yes] Stated in Assumption 1.
 - (b) Did you include complete proofs of all theoretical results? [Yes] See Appendix.
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] Code and data source are provided in supplementary materials.
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] (i) The online learning setup does not require data splits. (ii) How the hyperparameters were chosen is clarified in Appendix.
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [No] In distributed online learning, it is very common not to report the error bar, as the average loss plot is usually very stable in different runs (see reference [22, 38, 33] for more examples). We do find such stability in our experiments.
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] See Appendix.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? [Yes] We cited the creators of datasets.
 - (b) Did you mention the license of the assets? [Yes] We mention it in Appendix.
 - (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
 - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A] We only use open source datasets.
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A] We only use open source datasets.
5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]