
Federated Multi-Objective Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In recent years, multi-objective optimization (MOO) emerges as a foundational
2 problem underpinning many multi-agent multi-task learning applications. How-
3 ever, existing algorithms in MOO literature remain limited to centralized learning
4 settings, which do not satisfy the distributed nature and data privacy needs of
5 such multi-agent multi-task learning applications. This motivates us to propose a
6 new federated multi-objective learning (FMOL) framework with multiple clients
7 distributively and collaboratively solving an MOO problem while keeping their
8 training data private. Notably, our FMOL framework allows a different set of objec-
9 tive functions across different clients to support a wide range of applications, which
10 advances and generalizes the MOO formulation to the federated learning paradigm
11 for the first time. For this FMOL framework, we propose two new federated multi-
12 objective optimization (FMOO) algorithms called federated multi-gradient descent
13 averaging (FMGDA) and federated stochastic multi-gradient descent averaging
14 (FSMGDA). Both algorithms allow local updates to significantly reduce commu-
15 nication costs, while achieving the *same* convergence rates as those of the their
16 algorithmic counterparts in the single-objective federated learning. Our extensive
17 experiments also corroborate the efficacy of our proposed FMOO algorithms.

18 1 Introduction

19 In recent years, multi-objective optimization (MOO) has emerged as a foundational problem un-
20 derpinning many multi-agent multi-task learning applications, such as training neural networks for
21 multiple tasks [1], hydrocarbon production optimization [2], and tissue engineering [3]. MOO aims
22 at optimizing multiple objectives simultaneously, which can be mathematically cast as:

$$\min_{\mathbf{x} \in \mathcal{D}} \mathbf{F}(\mathbf{x}) := [f_1(\mathbf{x}), \dots, f_S(\mathbf{x})], \quad (1)$$

23 where $\mathbf{x} \in \mathcal{D} \subseteq \mathbb{R}^d$ is the model parameter, and $f_s : \mathbb{R}^d \rightarrow \mathbb{R}$, $s \in [S]$ is one of the objective
24 functions. Compared to conventional single-objective optimization, one key difference in MOO is the
25 coupling and potential conflicts between different objective functions. As a result, there may not exist
26 a common \mathbf{x} -solution that minimizes all objective functions. Rather, the goal in MOO is to find a
27 *Pareto stationary solution* that is not improvable for all objectives without sacrificing some objectives.
28 For example, in recommender system designs for e-commerce, the platform needs to consider different
29 customers with substantially conflicting shopping objectives (price, brand preferences, delivery speed,
30 etc.). Therefore, the platform’s best interest is often to find a Pareto-stationary solution, where one
31 cannot deviate to favor one consumer group further without hurting any other group. MOO with
32 conflicting objectives also has natural incarnations in many competitive game-theoretic problems,
33 where the goal is to determine an equilibrium among the conflicting agents in the Pareto sense.

34 Since its inception dating back to the 1950s, MOO algorithm design has evolved into two major
35 categories: gradient-free and gradient-based methods, with the latter garnering increasing attention

36 in the learning community in recent years due to their better performances (see Section 2 for more
 37 detailed discussions). However, despite these advances, all existing algorithms in the current MOO
 38 literature remain limited to centralized settings (i.e., training data are aggregated and accessible to
 39 a centralized learning algorithm). Somewhat ironically, such centralized settings do *not* satisfy the
 40 distributed nature and data privacy needs of many multi-agent multi-task learning applications, which
 41 motivates application of MOO in the first place. This gap between the existing MOO approaches and
 42 the rapidly growing importance of distributed MOO motivates us to make the first attempt to pursue a
 43 new **federated multi-objective learning** (FMOL) framework, with the aim to enable multiple clients
 44 to distributively solve MOO problems while keeping their computation and training data private.

45 So far, however, developing distributed optimization algorithms for FMOL with provable Pareto-
 46 stationary convergence remains uncharted territory. There are several key technical challenges that
 47 render FMOL far from being a straightforward extension of centralized MOO problems. First of
 48 all, due to the distributed nature of FMOL problems, one has to consider and model the *objective*
 49 *heterogeneity* (i.e., different clients could have different sets of objective functions) that is unseen in
 50 centralized MOO. Moreover, with local and private datasets being a defining feature in FMOL, the
 51 impacts of *data heterogeneity* (i.e., datasets are non-i.i.d. distributed across clients) also need to be
 52 mitigated in FMOL algorithm design. Last but not least, under the combined influence of objective
 53 and data heterogeneity, FMOL algorithms could be extremely sensitive to small perturbations in the
 54 determination of common descent direction among all objectives. This makes the FMOL algorithm
 55 design and the associated convergence analysis far more complicated than those of the centralized
 56 MOO. Toward this end, a fundamental question naturally arises:

57 *Under both objective and data heterogeneity in FMOL, is it possible to design effective and efficient*
 58 *algorithms with Pareto-stationary convergence guarantees?*

59 In this paper, we give an affirmative answer to the above question. Our key contribution is that
 60 we propose a new FMOL framework that captures both objective and data heterogeneity, based on
 61 which we develop two gradient-based algorithms with provable Pareto-stationary convergence rate
 62 guarantees. To our knowledge, our work is the first systematic attempt to bridge the gap between
 63 federated learning and MOO. Our main results and contributions are summarized as follows:

- 64 • We formalize the first federated multi-objective learning (FMOL) framework that supports both
 65 *objective and data heterogeneity* across clients, which significantly advances and generalizes the
 66 MOO formulation to the federated learning paradigm. As a result, our FMOL framework becomes
 67 a generic model that covers existing MOO models and various applications as special cases (see
 68 Section 3.2 for further details). This new FMOL framework lays the foundation to enable us to
 69 systematically develop FMOO algorithms with provable Pareto-stationary convergence guarantees.
- 70 • For the proposed FMOL framework, we first propose a federated multi-gradient descent averaging
 71 (FMGDA) algorithm based on the use of local full gradient evaluation at each client. Our analysis
 72 reveals that FMGDA achieves a linear $\mathcal{O}(\exp(-\mu T))$ and a sublinear $\mathcal{O}(1/T)$ Pareto-stationary
 73 convergence rates for μ -strongly convex and non-convex settings, respectively. Also, FMGDA
 74 employs a two-sided learning rates strategy to significantly lower communication costs (a key
 75 concern in the federated learning paradigm). It is worth pointing out that, in the single-machine
 76 special case where FMOL degenerates to a centralized MOO problem and FMGDA reduces to the
 77 traditional MGD method [4], our results improve the state-of-the-art analysis of MGD by eliminating
 78 the restrictive assumptions on the linear search of learning rate and extra sequence convergence.
 79 Thus, our results also advance the state of the art in general MOO theory.
- 80 • To alleviate the cost of full gradient evaluation in the large dataset regime, we further propose
 81 a federated stochastic multi-gradient descent averaging (FSMGDA) algorithm based on the use
 82 of stochastic gradient evaluations at each client. We show that FSMGDA achieves $\tilde{\mathcal{O}}(1/T)$ and
 83 $\mathcal{O}(1/\sqrt{T})$ Pareto-stationary convergence rate for μ -strongly convex and non-convex settings, re-
 84 spectively. We establish our convergence proof by proposing a new (α, β) -Lipschitz continuous
 85 stochastic gradient assumption (cf. Assumption 4), which relaxes the strong assumptions on first
 86 moment bound and Lipschitz continuity on common descent directions in [5]. We note that this new
 87 (α, β) -Lipschitz continuous stochastic gradient assumption can be viewed as a natural extension of
 88 the classical Lipschitz-continuous gradient assumption and could be of independent interest.

89 The rest of the paper is organized as follows. In Section 2, we review related works. In Section 3,
 90 we introduce our FMOL framework and two gradient-based algorithms (FMGDA and FSMGDA),
 91 which are followed by their convergence analyses in Section 4. We present the numerical results in

Table 1: Convergence rate results (shaded parts are our results) comparisons.

Methods	Strongly Convex		Non-convex	
	Rate	Assumption*	Rate	Assumption*
MGD [4]	$\mathcal{O}(r^T)$ #	Linear search & sequence convergence	$\mathcal{O}(1/T)$	Linear search & sequence convergence
SMGD [5]	$\mathcal{O}(1/T)$	First moment bound & Lipschitz continuity of λ	Not provided	Not provided
FMGDA	$\mathcal{O}(\exp(-\mu T))$ #	Not needed	$\mathcal{O}(1/T)$	Not needed
FSMGDA	$\tilde{\mathcal{O}}(1/T)$	(α, β) -Lipschitz continuous stochastic gradient	$\mathcal{O}(1/\sqrt{T})$	(α, β) -Lipschitz continuous stochastic gradient

#Notes on constants: μ is the strong convexity modulus; r is a constant depends on μ , s.t., $r \in (0, 1)$.

Assumption short-hands: “Linear search”: learning rate linear search [4]; “Sequence convergence”: $\{\mathbf{x}_t\}$ converges to \mathbf{x}^ [4]; “First moment bound” (Asm. 5.2(b) [5]): $\mathbb{E}[\|\nabla f(\mathbf{x}, \xi) - \nabla f(\mathbf{x})\|] \leq \eta(a + b\|\nabla f(\mathbf{x})\|)$; “Lipschitz continuity of λ ” (Asm. 5.4 [5]): $\|\lambda_k - \lambda_t\| \leq \beta \left\| \left[(\nabla f_1(\mathbf{x}_k) - \nabla f_1(\mathbf{x}_t))^T, \dots, (\nabla f_S(\mathbf{x}_k) - \nabla f_S(\mathbf{x}_t))^T \right] \right\|$; “ (α, β) -Lipschitz continuous stochastic gradient”: see Asm. 4.

92 Section 5 and conclude the work in Section 6. Due to space limitations, we relegate all proofs and
 93 some experiments to supplementary material.

94 **2 Related work**

95 In this section, we will provide an overview on algorithm designs for MOO and federated learning
 96 (FL), thereby placing our work in a comparative perspective to highlight our contributions and novelty.

97 **1) Multi-objective Optimization (MOO):** As mentioned in Section 1, since federated/distributed
 98 MOO has not been studied in the literature, all existing works we review below are centralized MOO
 99 algorithms. Roughly speaking, MOO algorithms can be grouped into two main categories. The first
 100 line of works are gradient-free methods (e.g., evolutionary MOO algorithms and Bayesian MOO
 101 algorithms [6, 7, 8, 9]). These methods are more suitable for small-scale problems but less practical
 102 for high-dimensional MOO models (e.g., deep neural networks). The second line of works focus on
 103 gradient-based approaches [10, 11, 4, 12, 5], which are more practical for high-dimensional MOO
 104 problems. However, while having received increasing attention from the community in recent years,
 105 Pareto-stationary convergence analysis of these gradient-based MOO methods remains in its infancy.

106 Existing gradient-based MOO methods can be further categorized as i) multi-gradient descent (MGD)
 107 algorithms with full gradients and ii) stochastic multi-gradient descent (SMGD) algorithms. It has
 108 been shown in [4] that MGD methods achieve $\mathcal{O}(r^T)$ for some $r \in (0, 1)$ and $\mathcal{O}(1/T)$ Pareto-
 109 stationary convergence rates for μ -strongly convex and non-convex functions, respectively. However,
 110 these results are established under the unconventional linear search of learning rate and sequence
 111 convergence assumptions, which are difficult to verify in practice. In comparison, FMGDA achieves a
 112 linear rate without needing such assumptions. For SMGD methods, the Pareto-stationary convergence
 113 analysis is further complicated by the stochastic gradient noise. Toward this end, an $\mathcal{O}(1/T)$ rate
 114 analysis for SMGD was provided in [5] based on rather strong assumptions on a first-moment bound
 115 and Lipschitz continuity of common descent direction. As a negative result, it was shown in [5]
 116 and [13] that the common descent direction needed in the SMGD method is likely to be a biased
 117 estimation, which may cause divergence issues.

118 In contrast, our FSMGDA achieves state-of-the-art $\tilde{\mathcal{O}}(1/T)$ and $\mathcal{O}(1/\sqrt{T})$ convergence rates for
 119 strongly-convex and non-convex settings, respectively, under a much milder assumption on Lipschitz
 120 continuous stochastic gradients. For easy comparisons, we summarize our results and the existing
 121 works in Table 1. It is worth noting recent works [13, 14] established faster convergence rates in
 122 the centralized MOO setting by using acceleration techniques, such as momentum, regularization
 123 and bi-level formulation. However, due to different settings and focuses, these results are orthogonal
 124 to ours and thus not directly comparable. Also, since acceleration itself is a non-trivial topic and
 125 could be quite brittle if not done right, in this paper, we focus on the basic and more robust stochastic
 126 gradient approach in FMOL. But for a comprehensive comparison on assumptions and main results
 127 of accelerated centralized MOO, we refer readers to Appendix A for further details.

128 **Federated Learning (FL) :** Since the seminal work by [15], FL has emerged as a popular distributed
 129 learning paradigm. Traditional FL aims at solving single-objective minimization problems with a large
 130 number of clients with decentralized data. Recent FL algorithms enjoy both high communication

131 efficiency and good generalization performance [15, 16, 17, 18, 19, 20]. Theoretically, many
 132 FL methods have the same convergence rates as their centralized counterparts under different FL
 133 settings [21, 22, 23, 24]. Recent works have also considered FL problems with more sophisticated
 134 problem structures, such as min-max learning [25, 26], reinforcement learning [27], multi-armed
 135 bandits [28], and bilevel and compositional optimization [29]. Although not directly related, classic
 136 FL has been reformulated in the form of MOO[30], which allows the use of a MGD-type algorithm
 137 instead of vanilla local SGD to solve the standard FL problem. We will show later that this MOO
 138 reformulation is a special case of our FMOL framework. So far, despite a wide range of applications
 139 (see Section 3.2 for examples), there remains a lack of a general FL framework for MOO. This
 140 motivates us to bridge the gap by proposing a general FMOL framework and designing gradient-based
 141 methods with provable Pareto-stationary convergence rates.

142 3 Federated multi-objective learning

143 3.1 Multi-objective optimization: A primer

144 As mentioned in Section 1, due to potential conflicts among the objective functions in MOO problem
 145 in (1), MOO problems adopt the notion of Pareto optimality:

146 **Definition 1** ((Weak) Pareto Optimality). *For any two solutions \mathbf{x} and \mathbf{y} , we say \mathbf{x} dominates \mathbf{y} if
 147 and only if $f_s(\mathbf{x}) \leq f_s(\mathbf{y}), \forall s \in [S]$ and $f_s(\mathbf{x}) < f_s(\mathbf{y}), \exists s \in [S]$. A solution \mathbf{x} is Pareto optimal if
 148 it is not dominated by any other solution. One solution \mathbf{x} is weakly Pareto optimal if there does not
 149 exist a solution \mathbf{y} such that $f_s(\mathbf{x}) > f_s(\mathbf{y}), \forall s \in [S]$.*

150 Similar to solving single-objective non-convex optimization problems, finding a Pareto-optimal
 151 solution in MOO is NP-Hard in general. As a result, it is often of practical interest to find a solution
 152 satisfying Pareto-stationarity (a necessary condition for Pareto optimality) stated as follows [10, 31]:

153 **Definition 2** (Pareto Stationarity). *A solution \mathbf{x} is said to be Pareto stationary if there is no common
 154 descent direction $\mathbf{d} \in \mathbb{R}^d$ such that $\nabla f_s(\mathbf{x})^\top \mathbf{d} < 0, \forall s \in [S]$.*

155 Note that for strongly convex functions, Pareto stationary solutions are also Pareto optimal. Following
 156 Definition 2, gradient-based MOO algorithms typically search for a common descent direction $\mathbf{d} \in \mathbb{R}^d$
 157 such that $\nabla f_s(\mathbf{x})^\top \mathbf{d} \leq 0, \forall s \in [S]$. If no such a common descent direction exists at \mathbf{x} , then
 158 \mathbf{x} is a Pareto stationary solution. For example, MGD [11] searches for an optimal weight $\boldsymbol{\lambda}^*$ of
 159 gradients $\nabla \mathbf{F}(\mathbf{x}) \triangleq \{\nabla f_s(\mathbf{x}), \forall s \in [S]\}$ by solving $\boldsymbol{\lambda}^*(\mathbf{x}) = \operatorname{argmin}_{\boldsymbol{\lambda} \in C} \|\boldsymbol{\lambda}^\top \nabla \mathbf{F}(\mathbf{x})\|^2$. Then,
 160 a common descent direction can be chosen as: $\mathbf{d} = \boldsymbol{\lambda}^\top \nabla \mathbf{F}(\mathbf{x})$. MGD performs the iterative
 161 update rule: $\mathbf{x} \leftarrow \mathbf{x} - \eta \mathbf{d}$ until a Pareto stationary point is reached, where η is a learning rate.
 162 SMGD [5] also follows the same process except for replacing full gradients by stochastic gradients.
 163 For MGD and SMGD methods, it is shown in [4] and [13] show that if $\|\boldsymbol{\lambda}^\top \nabla \mathbf{F}(\mathbf{x})\| = 0$ for some
 164 $\boldsymbol{\lambda} \in C$, where $C \triangleq \{\mathbf{y} \in [0, 1]^S, \sum_{s \in [S]} y_s = 1\}$, then \mathbf{x} is a Pareto stationary solution. Thus,
 165 $\|\mathbf{d}\|^2 = \|\boldsymbol{\lambda}^\top \nabla \mathbf{F}(\mathbf{x})\|^2$ can be used as a metric to measure the convergence of non-convex MOO
 166 algorithms [4, 13, 14]. On the other hand, for more tractable strongly convex MOO problems, the
 167 optimality gap $\sum_{s \in [S]} \lambda_s [f_s(\mathbf{x}) - f_s(\mathbf{x}^*)]$ is typically used as the metric to measure the convergence
 168 of an algorithm [5], where \mathbf{x}^* denotes the Pareto optimal point. We summarize and compare different
 169 convergence metrics as well as assumptions in MOO, detailed in Appendix A.

170 3.2 A general federated multi-objective learning framework

171 With the MOO preliminaries in Section 3.1, we now formalize our general federated multi-objective
 172 learning (FMOL) framework. For a system with M clients and S tasks (objectives), our FMOL
 173 framework can be written as:

$$\begin{aligned}
 & \min_{\mathbf{x}} \operatorname{Diag}(\mathbf{F}\mathbf{A}^\top), \tag{2} \\
 \mathbf{F} \triangleq & \begin{bmatrix} f_{1,1} & \cdots & f_{1,M} \\ \vdots & \ddots & \vdots \\ f_{S,1} & \cdots & f_{S,M} \end{bmatrix}_{S \times M}, \mathbf{A} \triangleq \begin{bmatrix} a_{1,1} & \cdots & a_{1,M} \\ \vdots & \ddots & \vdots \\ a_{S,1} & \cdots & a_{S,M} \end{bmatrix}_{S \times M},
 \end{aligned}$$

174 where matrix \mathbf{F} groups all potential objectives $f_{s,i}(\mathbf{x})$ for each task s at each client i , and $\mathbf{A} \in$
175 $\{0, 1\}^{S \times M}$ is a *binary* objective indicator matrix, with each element $a_{s,i} = 1$ if task s is of client
176 i 's interest and $a_{s,i} = 0$ otherwise. For each task $s \in [S]$, the global objective function $f_s(\mathbf{x})$
177 is the average of local objectives over all related clients, i.e., $f_s(\mathbf{x}) \triangleq \frac{1}{|R_s|} \sum_{i \in R_s} f_{s,i}(\mathbf{x})$, where
178 $R_s = \{i : a_{s,i} = 1, i \in [M]\}$. Note that, for notation simplicity, here we use simple average in $f_s(\mathbf{x})$,
179 which corresponds to the balanced dataset setting. Our FMLO framework can be directly extended to
180 imbalanced dataset settings by using weighted average proportional to dataset sizes of related clients.
181 For a client $i \in [M]$, its objectives of interest are $\{f_{s,i}(\mathbf{x}) : a_{s,i} = 1, s \in [S]\}$, which is a subset of $[S]$.

182 We note that FMOL generalizes MOO to the FL paradigm, which includes many existing MOO
183 problems as special cases and corresponds to a wide range of applications.

- 184 • If each client has only one distinct objective, i.e., $\mathbf{A} = \mathbb{I}_M$, $S = M$, then $\text{Diag}(\mathbf{FA}^\top) =$
185 $[f_1(\mathbf{x}), \dots, f_S(\mathbf{x})]^\top$, where each objective $f_s(\mathbf{x}), s \in [S]$ is optimized only by client s . This
186 special FMOL setting corresponds to the conventional multi-task learning and federated learning.
187 Indeed, [1] and [32] formulated a multi-task learning problem as MOO and considered Pareto
188 optimal solutions with various trade-offs. [30] also formulated FL as as distributed MOO problems.
189 Other examples of this setting include bi-objective formulation of offline reinforcement learning [33]
190 and decentralized MOO [34].
- 191 • If all clients share the same S objectives, i.e., \mathbf{A} is an all-one matrix, then $\text{Diag}(\mathbf{FA}^\top) =$
192 $[\frac{1}{M} \sum_{i \in [M]} f_{1,i}(\mathbf{x}), \dots, \frac{1}{M} \sum_{i \in [M]} f_{S,i}(\mathbf{x})]^\top$. In this case, FMOL reduces to federated MOO
193 problems with decentralized data that jointly optimizing fairness, privacy, and accuracy [35, 36, 37],
194 as well as MOO with decentralized data under privacy constraints (e.g., machine reassignment
195 among data centres [38] and engineering problems [39, 40]).
- 196 • If each client has a different subset of objectives (i.e., objective heterogeneity), FMLO allows
197 distinct preferences at each client. For example, each customer group on a recommender system in
198 e-commerce platforms might have different combinations of shopping preferences, such as product
199 price, brand, delivery speed, etc.

200 3.3 Federated Multi-Objective Learning Algorithms

201 Upon formalizing our FMOL frame-
202 work, our next goal is to develop
203 gradient-based algorithms for solving
204 large-scale high-dimensional FMOL
205 problems with *provable* Pareto station-
206 ary convergence guarantees and low
207 communication costs. To this end,
208 we propose two FMOL algorithms,
209 namely federated multiple gradient de-
210 scent averaging (FMGDA) and feder-
211 ated stochastic multiple gradient de-
212 scent averaging (FSMGDA) as shown
213 in Algorithm 1. We summarize our
214 key notation in Table 3 in Appendix
215 to allow easy references for readers.

216 As shown in Algorithm 1, in each
217 communication round $t \in [T]$, each
218 client synchronizes its local model
219 with the current global model \mathbf{x}_t
220 from the server (cf. Step 1). Then
221 each client runs K local steps based
222 on local data for all effective objec-
223 tives (cf. Step 2) with two options:
224 i) for FMGDA, each local step per-
225 forms local full gradient descent, i.e.,
226 $\mathbf{x}_{s,i}^{t,k+1} = \mathbf{x}_{s,i}^{t,k} - \eta_L \nabla f_{s,i}(\mathbf{x}_{s,i}^{t,k}), \forall s \in S_i$; ii) For FSMGDA, the local step performs stochastic
227 gradient descent, i.e., $\mathbf{x}_{s,i}^{t,k+1} = \mathbf{x}_{s,i}^{t,k} - \eta_L \nabla f_{s,i}(\mathbf{x}_{s,i}^{t,k}, \xi_i^{t,k}), \forall s \in S_i$, where $\xi_i^{t,k}$ denotes a random

Algorithm 1 Federated (Stochastic) Multiple Gradient De-
scendent Averaging (FMGDA/FSMGDA).

At Each Client i :

1. Synchronize local models $\mathbf{x}_{s,i}^{t,0} = \mathbf{x}_t, \forall s \in S_i$.
2. Local updates: for all $s \in S_i$, for $k = 1, \dots, K$,
(FMGDA): $\mathbf{x}_{s,i}^{t,k} = \mathbf{x}_{s,i}^{t,k-1} - \eta_L \nabla f_{s,i}(\mathbf{x}_{s,i}^{t,k-1})$.
(FSMGDA): $\mathbf{x}_{s,i}^{t,k} = \mathbf{x}_{s,i}^{t,k-1} - \eta_L \nabla f_{s,i}(\mathbf{x}_{s,i}^{t,k-1}, \xi_i^{t,k})$.
3. Return accumulated updates to server $\{\Delta_{s,i}^t, s \in S_i\}$:
(FMGDA): $\Delta_{s,i}^t = \sum_{k \in [K]} \nabla f_{s,i}(\mathbf{x}_{s,i}^{t,k})$.
(FSMGDA): $\Delta_{s,i}^t = \sum_{k \in [K]} \nabla f_{s,i}(\mathbf{x}_{s,i}^{t,k}, \xi_i^{t,k})$.

At the Server:

4. Receive accumulated updates $\{\Delta_{s,i}^t, \forall s \in S_i, \forall i \in [M]\}$.
5. Compute $\Delta_s^t = \frac{1}{|R_s|} \sum_{i \in R_s} \Delta_{s,i}^t, \forall s \in [S]$, where
 $R_s = \{i : a_{s,i} = 1, i \in [M]\}$.
6. Compute $\lambda_t^* \in [0, 1]^S$ by solving

$$\min_{\lambda_t \geq \mathbf{0}} \left\| \sum_{s \in [S]} \lambda_s^t \Delta_s^t \right\|^2, \quad \text{s.t.} \sum_{s \in [S]} \lambda_s^t = 1.$$

7. Let $\mathbf{d}_t = \sum_{s \in [S]} \lambda_s^{t,*} \Delta_s^t$ and update the global model
as: $\mathbf{x}_{t+1} = \mathbf{x}_t - \eta_t \mathbf{d}_t$, with a global learning rate η_t .
-

228 sample in local step k and round t at client i . Upon finishing K local updates, each client returns
 229 the accumulated update $\Delta_{s,i}^t$ for each effective objective to the server (cf. Step 3). Then, the server
 230 aggregates all returned Δ -updates from the clients to obtain the overall updates Δ_s^t for each objective
 231 $s \in [S]$ (cf. Steps 4 and 5), which will be used in solving a convex quadratic optimization problem
 232 with linear constraints to obtain an approximate common descent direction \mathbf{d}_t (cf. Step 6). Lastly, the
 233 global model is updated following the direction \mathbf{d}_t with global learning rate η_t (cf. Step 7).

234 Two remarks on Algorithm 1 are in order. First, we note that a two-sided learning rates strategy is
 235 used in Algorithm 1, which decouples the update schedules of local and global model parameters at
 236 clients and server, respectively. As shown in Section 4 later, this two-sided learning rates strategy
 237 enables better convergence rates by choosing appropriate learning rates. Second, to achieve low
 238 communication costs, Algorithm 1 leverages K local updates at each client and infrequent periodic
 239 communications between each client and the server. By adjusting the two-sided learning rates
 240 appropriately, the K -value can be made large to further reduce communication costs.

241 4 Pareto stationary convergence analysis

242 In this section, we analyze the Pareto stationary convergence performance for our FMGDA and
 243 FSMGDA algorithms in Sections 4.1 and 4.2, respectively, each of which include non-convex and
 244 strongly convex settings.

245 4.1 Pareto stationary convergence of FMGDA

246 **1) FMGDA: The Non-convex Setting.** Before presenting our Pareto stationary convergence results
 247 for FMGDA, we first state several assumptions as follows:

248 **Assumption 1.** (*L-Lipschitz continuous*) There exists a constant $L > 0$ such that $\|\nabla f_s(\mathbf{x}) -$
 249 $\nabla f_s(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|, \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d, s \in [S]$.

250 **Assumption 2.** (*Bounded Gradient*) The gradient of each objective at any client is bounded, i.e.,
 251 there exists a constant $G > 0$ such that $\|\nabla f_{s,i}(\mathbf{x})\|^2 \leq G^2, \forall s \in [S], i \in [M]$.

252 With the assumptions above, we state the Pareto stationary convergence of FMGDA as follows:

253 **Theorem 1** (FMGDA for Non-convex FMOL). Let $\eta_t = \eta \leq \frac{3}{2(1+L)}$. Under Assumptions 1 and 2,
 254 if at least one function $f_s, s \in [S]$ is bounded from below by f_s^{\min} , then the sequence $\{\mathbf{x}_t\}$ output by
 255 FMGDA satisfies: $\min_{t \in [T]} \|\mathbf{d}_t\|^2 \leq \frac{4(f_s^0 - f_s^{\min})}{T\eta} + \delta$, where $\delta \triangleq (8\eta_L^2 K^2 L^2 G^2)/\eta$.

256 The convergence bound in Theorem 1 contains two parts. The first part is an optimization error, which
 257 depends on the initial point and vanishes as T increases. The second part is due to local update steps
 258 K and data heterogeneity G , which can be mitigated by carefully choosing the local learning rate η_L .
 259 Specifically, the following Pareto stationary convergence rate of FMGDA follows immediately from
 260 Theorem 1 with an appropriate choice of local learning rate η_L :

261 **Corollary 2.** With a constant global learning rate $\eta_t = \eta, \forall t$, and a local learning rate $\eta_L =$
 262 $\mathcal{O}(1/\sqrt{T})$, the Pareto stationary convergence rate of FMGDA is $(1/T) \sum_{t \in [T]} \|\mathbf{d}_t\|^2 = \mathcal{O}(1/T)$.

263 Several interesting insights of Theorem 1 and Corollary 2 are worth pointing out: **1)** We note that
 264 FMGDA achieves a Pareto stationary convergence rate $\mathcal{O}(1/T)$ for non-convex FMOL, which is the
 265 *same* as the Pareto stationary rate of MGD for centralized MOO and the *same* convergence rate of
 266 gradient descent (GD) for single objective problems. This is somewhat surprising because FMGDA
 267 needs to handle more complex objective and data heterogeneity under FMOL; **2)** The two-sided
 268 learning rates strategy decouples the operation of clients and server by utilizing different learning
 269 rate schedules, thus better controlling the errors from local updates due to data heterogeneity; **3)**
 270 Note that in the single-client special case, FMGDA degenerates to the basic MGD algorithm. Hence,
 271 Theorem 1 directly implies a Pareto stationary convergence bound for MGD by setting $\delta = 0$ due
 272 to no local updates in centralized MOO. This convergence rate bound is consistent with that in [4].
 273 However, we note that our result is achieved *without* using the linear search step for learning rate [4],
 274 which is much easier to implement in practice (especially for deep learning models); **4)** Our proof is
 275 based on standard assumptions in first-order optimization, while previous works require strong and
 276 unconventional assumptions. For example, a convergence of \mathbf{x} -sequence is assumed in [4].

277 **2) FMGDA: The Strongly Convex Setting.** Now, we consider the strongly convex setting for FMOL,
 278 which is more tractable but still of interest in many learning problems in practice. In the strongly
 279 convex setting, we have the following additional assumption:

280 **Assumption 3.** (μ -Strongly Convex Function) Each objective $f_s(\mathbf{x})$, $s \in [S]$ is a μ -strongly convex
 281 function, i.e., $f_s(\mathbf{y}) \geq f_s(\mathbf{x}) + \nabla f_s(\mathbf{x})(\mathbf{y} - \mathbf{x}) + \frac{\mu}{2} \|\mathbf{y} - \mathbf{x}\|^2$ for some $\mu > 0$.

282 For more tractable strongly-convex FMOL problems, we show that FMGDA achieves a stronger
 283 Pareto stationary convergence performance as follows:

284 **Theorem 3** (FMGDA for μ -Strongly Convex FMOL). Let $\eta_t = \eta$ such that $\eta \leq \frac{3}{2(1+L)}$, $\eta \leq \frac{1}{2L+\mu}$
 285 and $\eta \geq \frac{1}{\mu T}$. Under Assumptions 1- 3, pick \mathbf{x}_t as the final output of the FMGDA algorithm with
 286 weights $w_t = (1 - \frac{\mu\eta}{2})^{1-t}$. Then, it holds that $\mathbb{E}[\Delta_Q^t] \leq \|\mathbf{x}_0 - \mathbf{x}_*\|^2 \mu \exp(-\frac{\eta\mu T}{2}) + \delta$, where
 287 $\Delta_Q^t \triangleq \sum_{s \in [S]} \lambda_s^{t,*} [f_s(\mathbf{x}_t) - f_s(\mathbf{x}_*)]$ and $\delta = \frac{8\eta_L^2 K^2 L^2 G^2 S^2}{\mu} + 2\eta_L^2 K^2 L^2 G^2$.

288 Theorem 3 immediately implies following Pareto stationary convergence rate for FMGDA with a
 289 proper choice of local learning rate:

290 **Corollary 4.** If η_L is chosen sufficiently small such that $\delta = \mathcal{O}(\mu \exp(-\mu T))$, then the Pareto
 291 stationary convergence rate of FMGDA is $\mathbb{E}[\Delta_Q^t] = \mathcal{O}(\mu \exp(-\mu T))$.

292 Again, several interesting insights can be drawn from Theorem 3 and Corollary 4. First, for strongly
 293 convex FMOL, FMGDA achieves a linear convergence rate $\mathcal{O}(\mu \exp(-\mu T))$, which again matches
 294 those of MGD for centralized MOO and GD for single-objective problems. Second, compared with
 295 the non-convex case, the convergence bounds suggest FMGDA could use a larger local learning rate
 296 for non-convex functions thanks to our two-sided learning rates design. A novel feature of FMGDA
 297 for strongly convex FMOL is the randomly chosen output x_t with weight w_t from the \mathbf{x}_t -trajectory,
 298 which is inspired by the classical work in stochastic gradient descent (SGD) [41]. Note that, for
 299 implementation in practice, one does not need to store all \mathbf{x}_t -values. Instead, the algorithm can be
 300 implemented by using a random clock for stopping [41].

301 4.2 Pareto stationary convergence of FSMGDA

302 While enjoying strong performances, FMGDA uses local full gradients at each client, which could be
 303 costly in the large dataset regime. Thus, it is of theoretical and practical importance to consider the
 304 stochastic version of FMGDA, i.e., federated stochastic multi-gradient descent averaging (FSMGDA).

305 **1) FSMGDA: The Non-convex Setting.** A fundamental challenge in analyzing the Pareto stationarity
 306 convergence of FSMGDA and other stochastic multi-gradient descent (SMGD) methods stems
 307 from bounding the error of the common descent direction estimation, which is affected by both λ_t^*
 308 (obtained by solving a quadratic programming problem) and the stochastic gradient variance. In fact,
 309 it is shown in [5] and [13] that the stochastic common descent direction in SMGD-type methods
 310 could be biased, leading to divergence issues. To address these challenges, in this paper, we propose
 311 to use a *new* assumption on the stochastic gradients, which is stated as follows:

312 **Assumption 4** ((α, β) -Lipschitz Continuous Stochastic Gradient). A function f has (α, β) -Lipschitz
 313 continuous stochastic gradients if there exist two constants $\alpha, \beta > 0$ such that, for any two indepen-
 314 dent training samples ξ and ξ' , $\mathbb{E}[\|\nabla f(\mathbf{x}, \xi) - \nabla f(\mathbf{y}, \xi')\|^2] \leq \alpha \|\mathbf{x} - \mathbf{y}\|^2 + \beta \sigma^2$.

315 In plain language, Assumption 4 says that the stochastic gradient estimation of an objective does not
 316 change too rapidly. We note that the (α, β) -Lipschitz continuous stochastic gradient assumption is a
 317 natural extension of the classic L -Lipschitz continuous gradient assumption (cf. Assumption 1) and
 318 generalizes several assumptions of SMGD convergence analysis in previous works. We note that
 319 Assumption 1 is not necessarily too hard to satisfy in practice. For example, when the underlying
 320 distribution of training samples ξ has a bounded support (typically a safe assumption for most
 321 applications in practice due to the finite representation limit of computing systems), suppose that
 322 Assumption 1 holds (also a common assumption in the optimization literature), then for any given
 323 \mathbf{x} and \mathbf{y} , the left-hand-side of the inequality in Assumption 4 is bounded due to the L -smoothness
 324 in Assumption 1. In this case, there always exist a sufficiently large α and a β such that the right-
 325 hand-side of the inequality in Assumption 1 holds. Please see Appendix A for further details. In
 326 addition, we need the following assumptions for the stochastic gradients, which are commonly used
 327 in standard SGD-based analyses [41, 42, 43, 44].

328 **Assumption 5.** (Unbiased Stochastic Estimation) The stochastic gradient estimation is unbiased for
 329 each objective among clients, i.e., $\mathbb{E}[\nabla f_{s,i}(\mathbf{x}, \xi)] = \nabla f_{s,i}(\mathbf{x}), \forall s \in [S], i \in [M]$.

330 **Assumption 6.** (Bounded Stochastic Gradient) The stochastic gradients satisfy $\mathbb{E}[\|\nabla f_{s,i}(\mathbf{x}, \xi)\|^2] \leq$
 331 $D^2, \forall s \in [S], i \in [M]$ for some constant $D > 0$.

332 With the assumptions above, we now state the Pareto stationarity convergence of FSMGDA as follows:
 333

334 **Theorem 5** (FSMGDA for Non-convex FMOL). Let $\eta_t = \eta \leq \frac{3}{2(1+L)}$. Under Assumptions 4–6, if
 335 an objective f_s is bounded from below by f_s^{\min} , then the sequence $\{\mathbf{x}_t\}$ output by FSMGDA satisfies:
 336 $\min_{t \in [T]} \mathbb{E} \|\mathbf{d}_t\|^2 \leq \frac{2S(f_s^0 - f_s^{\min})}{\eta T} + \delta$, where $\delta = L\eta S^2 D^2 + S(\alpha\eta_L^2 K^2 D^2 + \beta\sigma^2)$.

337 Theorem 5 immediately implies an $\mathcal{O}(1/\sqrt{T})$ convergence rate of FSMGDA for non-convex FMOL:

338 **Corollary 6.** With a constant global learning rate $\eta_t = \eta = \mathcal{O}(1/\sqrt{T}), \forall t$ and a local learning
 339 rate $\eta_L = \mathcal{O}(1/T^{1/4})$, and if $\beta = \mathcal{O}(\eta)$, the Pareto stationarity convergence rate of FSMGDA is
 340 $\min_{t \in [T]} \mathbb{E} \|\mathbf{d}_t\|^2 = \mathcal{O}(1/\sqrt{T})$.

341 **2) The Strongly Convex Setting:** For more tractable strongly convex FMOL problems, we can show
 342 that FSMGDA achieve stronger convergence results as follows:

343 **Theorem 7** (FSMGDA for μ -Strongly Convex FMOL). Let $\eta_t = \eta = \Omega(\frac{1}{\mu T})$. Under Assumptions 3,
 344 5 and 6, pick \mathbf{x}_t as the final output of the FSMGDA algorithm with weight $w_t = (1 - \frac{\mu\eta}{2})^{1-t}$. Then,
 345 it holds that: $\mathbb{E}[\Delta_Q^t] \leq \|\mathbf{x}_0 - \mathbf{x}_*\|^2 \mu \exp(-\frac{\eta}{2}\mu T) + \delta$, where $\Delta_Q^t = \sum_{s \in [S]} \lambda_s^{t,*} [f_s(\mathbf{x}_t) - f_s(\mathbf{x}_*)]$
 346 and $\delta = \frac{1}{\mu} S^2 (\alpha\eta_L^2 K^2 D^2 + \beta\sigma^2) + \frac{\eta S^2 D^2}{2}$.

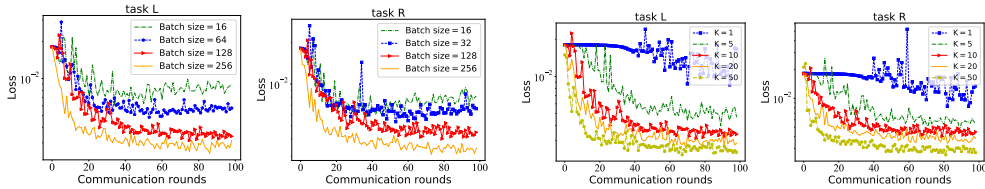
347 The following Pareto station convergence rate of FSMGDA follows immediately from Theorem 7:

348 **Corollary 8.** Choose $\eta_L = \mathcal{O}(\frac{1}{\sqrt{T}})$ and $\eta = \Theta(\frac{\log(\max(1, \mu^2 T))}{\mu T})$. If $\beta = \mathcal{O}(\eta)$, then the Pareto
 349 stationary convergence rate of FSMGDA is $\mathbb{E}[\Delta_Q^t] \leq \tilde{\mathcal{O}}(1/T)$.

350 Corollary 8 says that, With proper learning rates, FSMGDA achieves $\tilde{\mathcal{O}}(1/T)$ Pareto stationary
 351 convergence rate (i.e., ignoring logarithmic factors) for strongly convex FMOL. Also, in the single-
 352 client special case with no local updates, FSMGDA reduces to the SMGD algorithm and $\delta =$
 353 $\frac{4}{\mu} \beta S^2 \sigma^2 + \frac{\eta S^2 D^2}{2}$ in this case. Then, Theorem 7 implies an $\tilde{\mathcal{O}}(\frac{1}{T})$ Pareto stationarity convergence
 354 rate for SMGD for strongly convex MOO problems, which is consistent with previous works [5].
 355 However, our convergence rate proof uses a more conventional (α, β) -Lipschitz stochastic gradient
 356 assumption, rather than the unconventional assumptions on the first moment bound and Lipschitz
 357 continuity of common descent directions in [5].

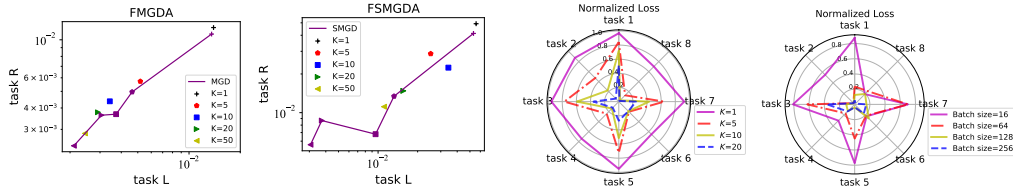
358 5 Numerical results

359 In this section, we show the main numerical experiments of our FMGDA and FSMGDA algorithms
 360 in different datasets, while relegating the experimental settings and details to the appendix.



(a) Training loss convergence in terms of communication rounds with different batch-sizes under non-i.i.d. data partition in MultiMNIST. (b) The impacts of local update number K on training loss convergence in terms of communication rounds.

Figure 1: Training loss convergence comparison.



(a) 100 communication rounds with various local steps K , corresponding federated and centralized settings share the same marker shape. (b) Normalized loss with the River Flow datasets.

Figure 2: Training losses comparison

361 **1) Ablation Experiments on Two-**
 362 **Tasks FMOL: 1-a) Impacts of Batch**
 363 **Size on Convergence:** First, we compare
 364 the convergence results in terms of the
 365 number of communication rounds using
 366 the “MultiMNIST” dataset [45] with two
 367 tasks (L and R) as objectives. We test our
 368 algorithms with four different cases with
 369 batch sizes being [16, 64, 128, 256]. To
 370 reduce computational costs in this experiment,
 371 the dataset size for each client is limited to
 372 256. Hence, the batch size 256 corresponds to
 373 FMGDA and all other batch sizes correspond to
 374 FSMGDA. As shown in Fig. 1(a), under non-i.i.d.
 data partition, both FMGDA and FSMGDA
 algorithms converge. Also, the convergence
 speed of the FSMGDA algorithm increases as
 the batch size gets larger. These results are
 consistent with our theoretical analyses as
 outlined in Theorems 1 and 5.

375 **1-b) Impacts of Local Update Steps on**
 376 **Convergence:** Next, we evaluate our
 377 algorithms with different numbers of local
 378 update steps K . As shown in Fig. 1(b) and
 Table 2, both algorithms converge faster as
 the number of the local steps K increases.
 This is because both algorithms effectively run
 more iterative updates as K gets large.

379 **1-c) Comparisons between FMOL and**
 380 **Centralized MOO:** Since this work is the
 381 first that investigates FMOL, it is also
 382 interesting to empirically compare the
 383 differences between FMOL and centralized
 384 MOO methods. In Fig. 2(a), we compare the
 385 training loss of FMGDA and FSMGDA with
 386 those of the centralized MGD and SMGD
 methods after 100 communication rounds. For
 fair comparisons, the centralized MGD and
 SMGD methods use $\sum_i^M |S_i|$ batch-sizes
 and run $K \times T$ iterations. Our results
 indicate that FMGDA and MGD produce
 similar results, while the performance of
 FSMGDA is slightly worse than that of
 SMGD due to FSMGDA’s sensitivity to
 objective and data heterogeneity in
 stochastic settings. These numerical results
 confirm our theoretical convergence analysis.

387 **2) Experiments on Larger FMOL:** We
 388 further test our algorithms on FMOL
 389 problems of larger sizes. In this experiment,
 390 we use the River Flow dataset[46], which
 391 contains *eight* tasks in this problem. To
 392 better visualize 8 different tasks, we
 illustrate the normalized loss in radar charts
 in Fig. 2(b). In this 8-task setting, we can
 again verify that more local steps K and a
 larger training batch size lead to faster
 convergence. In the appendix, we also
 vary the effectiveness of our FMGDA and
 FSMGDA algorithms in CelebA [47] (40
 tasks), alongside with other hyperparameter
 tuning results.

393 6 Conclusion and discussions

394 In this paper, we proposed the first general
 395 framework to extend multi-objective
 396 optimization to the federated learning
 397 paradigm, which considers both objective
 398 and data heterogeneity. We showed that,
 399 even under objective and data heterogeneity,
 400 both of our proposed algorithms enjoy the
 401 same Pareto stationary convergence rate as
 their centralized counterparts. In our future
 work, we will go beyond the limitation in
 the analysis of MOO that an extra
 assumption on the stochastic gradients (and
 λ). In this paper, we have proposed a
 weaker assumption (Assumption 4). We
 conjecture that using acceleration
 techniques, e.g., momentum, variance
 reduction, and regularization, could relax
 such assumption and achieve better
 convergence rate, which is a promising
 direction for future works.

Table 2: Communication rounds needed for 10^{-2} loss.

	i.i.d.		non-i.i.d.	
	Task L	Task R	Task L	Task R
$K = 1$	82	84	96	82
$K = 5$	18(4.6 \times)	20(4.2 \times)	24(4.0 \times)	20(4.1 \times)
$K = 10$	10(8.2 \times)	9(9.3 \times)	13(7.4 \times)	10(8.2 \times)
$K = 20$	5(16.4 \times)	5(16.8 \times)	6(16.0 \times)	5(16.4 \times)

402 **References**

- 403 [1] O. Sener and V. Koltun, “Multi-task learning as multi-objective optimization,” *Advances in*
404 *neural information processing systems*, vol. 31, 2018.
- 405 [2] J. You, W. Ampomah, and Q. Sun, “Development and application of a machine learning based
406 multi-objective optimization workflow for co2-eor projects,” *Fuel*, vol. 264, p. 116758, 2020.
- 407 [3] J. Shi, J. Song, B. Song, and W. F. Lu, “Multi-objective optimization design through machine
408 learning for drop-on-demand bioprinting,” *Engineering*, vol. 5, no. 3, pp. 586–593, 2019.
- 409 [4] J. Fliege, A. I. F. Vaz, and L. N. Vicente, “Complexity of gradient descent for multiobjective
410 optimization,” *Optimization Methods and Software*, vol. 34, no. 5, pp. 949–959, 2019.
- 411 [5] S. Liu and L. N. Vicente, “The stochastic multi-gradient algorithm for multi-objective optimiza-
412 tion and its application to supervised machine learning,” *Annals of Operations Research*, pp.
413 1–30, 2021.
- 414 [6] Q. Zhang and H. Li, “Moea/d: A multiobjective evolutionary algorithm based on decomposition,”
415 *IEEE Transactions on evolutionary computation*, vol. 11, no. 6, pp. 712–731, 2007.
- 416 [7] K. Deb, A. Pratap, S. Agarwal, and T. Meyarivan, “A fast and elitist multiobjective genetic
417 algorithm: Nsga-ii,” *IEEE transactions on evolutionary computation*, vol. 6, no. 2, pp. 182–197,
418 2002.
- 419 [8] S. Belakaria, A. Deshwal, N. K. Jayakodi, and J. R. Doppa, “Uncertainty-aware search frame-
420 work for multi-objective bayesian optimization,” in *Proceedings of the AAAI Conference on*
421 *Artificial Intelligence*, vol. 34, 2020, pp. 10 044–10 052.
- 422 [9] M. Laumanns and J. Ocenasek, “Bayesian optimization algorithms for multi-objective optimiza-
423 tion,” in *International Conference on Parallel Problem Solving from Nature*. Springer, 2002,
424 pp. 298–307.
- 425 [10] J. Fliege and B. F. Svaiter, “Steepest descent methods for multicriteria optimization,” *Mathemat-*
426 *ical methods of operations research*, vol. 51, no. 3, pp. 479–494, 2000.
- 427 [11] J.-A. Désidéri, “Multiple-gradient descent algorithm (mgda) for multiobjective optimization,”
428 *Comptes Rendus Mathématique*, vol. 350, no. 5-6, pp. 313–318, 2012.
- 429 [12] S. Peitz and M. Dellnitz, “Gradient-based multiobjective optimization with uncertainties,” in
430 *NEO 2016*. Springer, 2018, pp. 159–182.
- 431 [13] S. Zhou, W. Zhang, J. Jiang, W. Zhong, J. GU, and W. Zhu, “On the convergence of stochastic
432 multi-objective gradient manipulation and beyond,” in *Advances in Neural Information*
433 *Processing Systems*, A. H. Oh, A. Agarwal, D. Belgrave, and K. Cho, Eds., 2022. [Online].
434 Available: <https://openreview.net/forum?id=ScwfQ7hdwyP>
- 435 [14] H. Fernando, H. Shen, M. Liu, S. Chaudhury, K. Murugesan, and T. Chen, “Mitigating gradient
436 bias in multi-objective learning: A provably convergent stochastic approach,” *arXiv preprint*
437 *arXiv:2210.12624*, 2022.
- 438 [15] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, “Communication-efficient
439 learning of deep networks from decentralized data,” in *Artificial intelligence and statistics*.
440 PMLR, 2017, pp. 1273–1282.
- 441 [16] T. Li, A. K. Sahu, M. Zaheer, M. Sanjabi, A. Talwalkar, and V. Smith, “Federated optimization
442 in heterogeneous networks,” in *Proceedings of Machine Learning and Systems*, I. Dhillon,
443 D. Papailiopoulos, and V. Sze, Eds., vol. 2, 2020, pp. 429–450.
- 444 [17] D. A. E. Acar, Y. Zhao, R. M. Navarro, M. Mattina, P. N. Whatmough, and V. Saligrama,
445 “Federated learning based on dynamic regularization,” in *International Conference on Learning*
446 *Representations*, 2021.
- 447 [18] J. Wang, Q. Liu, H. Liang, G. Joshi, and H. V. Poor, “Tackling the objective inconsistency
448 problem in heterogeneous federated optimization,” *Advances in Neural Information Processing*
449 *Systems*, vol. 33, 2020.
- 450 [19] T. Lin, S. U. Stich, K. K. Patel, and M. Jaggi, “Don’t use large mini-batches, use local
451 sgd,” in *International Conference on Learning Representations*, 2020. [Online]. Available:
452 <https://openreview.net/forum?id=B1eyO1BFPr>

- 453 [20] H. Yang, P. Qiu, and J. Liu, “Taming fat-tailed (“heavier-tailed” with potentially infinite
454 variance) noise in federated learning,” in *Advances in Neural Information Processing
455 Systems*, A. H. Oh, A. Agarwal, D. Belgrave, and K. Cho, Eds., 2022. [Online]. Available:
456 <https://openreview.net/forum?id=8SiIFGuXgmk>
- 457 [21] S. P. Karimireddy, S. Kale, M. Mohri, S. Reddi, S. Stich, and A. T. Suresh, “SCAFFOLD:
458 Stochastic controlled averaging for federated learning,” in *Proceedings of the 37th International
459 Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, H. D. III
460 and A. Singh, Eds., vol. 119. PMLR, 13–18 Jul 2020, pp. 5132–5143.
- 461 [22] H. Yang, M. Fang, and J. Liu, “Achieving linear speedup with partial worker participation in
462 non-IID federated learning,” in *International Conference on Learning Representations*, 2021.
- 463 [23] H. Yang, X. Zhang, P. Khanduri, and J. Liu, “Anarchic federated learning,” in *International
464 Conference on Machine Learning*. PMLR, 2022, pp. 25 331–25 363.
- 465 [24] X. Zhang, M. Fang, Z. Liu, H. Yang, J. Liu, and Z. Zhu, “Net-fleet: achieving linear convergence
466 speedup for fully decentralized federated learning with heterogeneous data,” *Proceedings of
467 the Twenty-Third International Symposium on Theory, Algorithmic Foundations, and Protocol
468 Design for Mobile Networks and Mobile Computing*, 2022.
- 469 [25] H. Yang, Z. Liu, X. Zhang, and J. Liu, “SAGDA: Achieving $\mathcal{O}(\epsilon^{-2})$ communication
470 complexity in federated min-max learning,” in *Advances in Neural Information Processing
471 Systems*, A. H. Oh, A. Agarwal, D. Belgrave, and K. Cho, Eds., 2022. [Online]. Available:
472 <https://openreview.net/forum?id=wTp4KgVIJ5>
- 473 [26] P. Sharma, R. Panda, G. Joshi, and P. Varshney, “Federated minimax optimization: Improved
474 convergence analyses and algorithms,” in *International Conference on Machine Learning*.
475 PMLR, 2022, pp. 19 683–19 730.
- 476 [27] S. Khodadadian, P. Sharma, G. Joshi, and S. T. Maguluri, “Federated reinforcement learning:
477 Linear speedup under markovian sampling,” in *International Conference on Machine Learning*.
478 PMLR, 2022, pp. 10 997–11 057.
- 479 [28] C. Shi, C. Shen, and J. Yang, “Federated multi-armed bandits with personalization,” in *International
480 Conference on Artificial Intelligence and Statistics*. PMLR, 2021, pp. 2917–2925.
- 481 [29] D. A. Tarzanagh, M. Li, C. Thrampoulidis, and S. Oymak, “FedNest: Federated bilevel,
482 minimax, and compositional optimization,” in *Proceedings of the 39th International Conference
483 on Machine Learning*, ser. Proceedings of Machine Learning Research, K. Chaudhuri, S. Jegelka,
484 L. Song, C. Szepesvari, G. Niu, and S. Sabato, Eds., vol. 162. PMLR, 17–23 Jul 2022, pp.
485 21 146–21 179.
- 486 [30] Z. Hu, K. Shaloudegi, G. Zhang, and Y. Yu, “Federated learning meets multi-objective optimiza-
487 tion,” *IEEE Transactions on Network Science and Engineering*, 2022.
- 488 [31] K. Miettinen, *Nonlinear multiobjective optimization*. Springer Science & Business Media,
489 2012, vol. 12.
- 490 [32] X. Lin, H.-L. Zhen, Z. Li, Q.-F. Zhang, and S. Kwong, “Pareto multi-task learning,” *Advances
491 in neural information processing systems*, vol. 32, 2019.
- 492 [33] Y. Yang, J. Jiang, T. Zhou, J. Ma, and Y. Shi, “Pareto policy pool for model-based offline
493 reinforcement learning,” in *International Conference on Learning Representations*, 2022.
494 [Online]. Available: <https://openreview.net/forum?id=OqcZu8JIzS>
- 495 [34] M. J. Blondin and M. Hale, “A decentralized multi-objective optimization algorithm,” *Journal
496 of Optimization Theory and Applications*, vol. 189, no. 2, pp. 458–485, 2021.
- 497 [35] L. T. Bui, H. A. Abbass, and D. Essam, “Local models—an approach to distributed multi-
498 objective optimization,” *Computational Optimization and Applications*, vol. 42, no. 1, pp.
499 105–139, 2009.
- 500 [36] S. Cui, W. Pan, J. Liang, C. Zhang, and F. Wang, “Addressing algorithmic disparity and
501 performance inconsistency in federated learning,” *Advances in Neural Information Processing
502 Systems*, vol. 34, pp. 26 091–26 102, 2021.
- 503 [37] N. Mehrabi, C. de Lichy, J. McKay, C. He, and W. Campbell, “Towards multi-objective
504 statistically fair federated learning,” *arXiv preprint arXiv:2201.09917*, 2022.

- 505 [38] T. Saber, X. Gandibleux, M. O'Neill, L. Murphy, and A. Ventresque, "A comparative study
506 of multi-objective machine reassignment algorithms for data centres," *Journal of Heuristics*,
507 vol. 26, no. 1, pp. 119–150, 2020.
- 508 [39] L. Yin, T. Wang, and B. Zheng, "Analytical adaptive distributed multi-objective optimization
509 algorithm for optimal power flow problems," *Energy*, vol. 216, p. 119245, 2021.
- 510 [40] Y. Jin, *Multi-objective machine learning*. Springer Science & Business Media, 2006, vol. 16.
- 511 [41] S. Ghadimi and G. Lan, "Stochastic first-and zeroth-order methods for nonconvex stochastic
512 programming," *SIAM Journal on Optimization*, vol. 23, no. 4, pp. 2341–2368, 2013.
- 513 [42] L. Bottou, F. E. Curtis, and J. Nocedal, "Optimization methods for large-scale machine learning,"
514 *Siam Review*, vol. 60, no. 2, pp. 223–311, 2018.
- 515 [43] H. B. McMahan *et al.*, "Advances and open problems in federated learning," *Foundations and
516 Trends® in Machine Learning*, vol. 14, no. 1, 2021.
- 517 [44] J. Wang, Z. Charles, Z. Xu, G. Joshi, H. B. McMahan, M. Al-Shedivat, G. Andrew, S. Aves-
518 timehr, K. Daly, D. Data *et al.*, "A field guide to federated optimization," *arXiv preprint
519 arXiv:2107.06917*, 2021.
- 520 [45] S. Sabour, N. Frosst, and G. E. Hinton, "Dynamic routing between capsules," *Advances in
521 neural information processing systems*, vol. 30, 2017.
- 522 [46] L. Nie, K. Wang, W. Kang, and Y. Gao, "Image retrieval with attribute-associated auxiliary
523 references," in *2017 International Conference on Digital Image Computing: Techniques and
524 Applications (DICTA)*. IEEE, 2017, pp. 1–6.
- 525 [47] Z. Liu, P. Luo, X. Wang, and X. Tang, "Deep learning face attributes in the wild," in *Proceedings
526 of the IEEE international conference on computer vision*, 2015, pp. 3730–3738.
- 527 [48] Q. Mercier, F. Poirion, and J.-A. Désidéri, "A stochastic multiple gradient descent algorithm,"
528 *European Journal of Operational Research*, vol. 271, no. 3, pp. 808–817, 2018.
- 529 [49] Y. LeCun, C. Cortes, and C. Burges, "Mnist handwritten digit database," Available: [http://yann.
530 lecun.com/exdb/mnist](http://yann.lecun.com/exdb/mnist), 1998.
- 531 [50] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition. cvpr. 2016,"
532 *arXiv preprint arXiv:1512.03385*, 2016.

533 **A Gradient-based methods in MOO**

534 (Stochastic) Gradient-based methods in MOO have attracted much attention owing to simple update
 535 rules and less intensive computation recently, thus rendering them perfect candidates to underpin
 536 MOO applications in deep learning under first-oracle. However, their theoretical understandings
 537 remain less explored relative to their counterparts of single objective optimization. Hence, we
 538 highlight the existing works and corresponding assumptions alongside with convergence metrics.

539 **Existing Works.** Various works managed to explore the convergence rates under different assump-
 540 tions in strongly-convex, convex, and non-convex functions as listed in Table 4. Using full gradient,
 541 MGD [4] could achieve tight convergence rates in strongly-convex and non-convex cases, i.e., linear
 542 rate $\mathcal{O}(r^T)$, $r \in (0, 1)$ and sub-linear rate $\mathcal{O}(\frac{1}{T})$. However, it requires linear search of learning rate
 543 in the algorithm and sequence convergence $\{\mathbf{x}_t\}$ converges to \mathbf{x}_* . The linear search of learning
 544 rate is a classic technique, but does not fits in gradient-based algorithms in deep learning. Moreover,
 545 sequence convergence assumption is a too strong assumption. With no local step, our FMGDA
 546 degenerates to MGD. As a result, our analysis also provide the same order convergence rates in both
 547 strongly-convex and non-convex functions while avoiding strong and unpractical assumptions. If
 548 using stochastic gradient, SMGD methods makes a further complicated case. The stochastic gradient
 549 noise would complicate the analysis and thus it is still unclear whether SMGD is guaranteed to
 550 converge. [5] provided convergence rate for SMGD but extra assumptions and/or unreasonably large
 551 batch requirements were needed. On the other hand, [5] and [13] showed that the common descent
 552 direction provided by SMGD method is likely to be a biased estimation, rendering non-convergence
 553 issues. Recently, by utilizing momentum, MoCo [14] and CR-MOGM [13] were proposed with
 554 corresponding convergence guarantees. However, these analyses do not shed light on pure SMGD
 555 despite its widespread application.

556 **Assumptions.** When applying stochastic gradient to MOO, common descent direction estimation
 557 $\lambda^T \nabla \mathbf{F}(\mathbf{x}, \xi)$ is a biased estimation and thus rendering potential non-convergence issues [5, 13].
 558 This is a limitation for SMGD. However, SMGD does work well with a wide range of applications
 559 in practice. Understanding under what conditions can SMGD have convergence guarantee is thus
 560 an important problem. [48] assumes convexity property(H5): $f(\mathbf{x}, \xi) - f(\mathbf{x}^*, \xi) \geq \frac{\alpha}{2} \|\mathbf{x} - \mathbf{x}^*\|^2$
 561 almost sure. [5] utilizes weaker assumptions but still needs first moment bound (Assumption 5.2(b)):
 562 $\mathbb{E}[\|\nabla f(\mathbf{x}, \xi) - \nabla f(\mathbf{x})\|] \leq \eta(a + b\|\nabla f(\mathbf{x})\|)$ and Lipschitz continuity of λ (Assumption 5.4):
 563 $\|\lambda_k - \lambda_t\| \leq \beta \left\| \left[(\nabla f_1(\mathbf{x}_k) - \nabla f_1(\mathbf{x}_t))^T, \dots, (\nabla f_S(\mathbf{x}_k) - \nabla f_S(\mathbf{x}_t))^T \right] \right\|$.

564 In this paper, we use (α, β) -Lipschitz continuous stochastic gradient (Assumption 4). In essence, we
 565 need the stochastic gradient estimation satisfying $\mathbb{E}[\|\nabla f(\mathbf{x}, \xi) - \nabla f(\mathbf{y}, \xi')\|^2] \leq \alpha \|\mathbf{x} - \mathbf{y}\|^2 + \beta \sigma^2$
 566 for any two independent samples ξ and ξ' . For the inequality $\mathbb{E}[\|\nabla f(\mathbf{x}, \xi) - \nabla f(\mathbf{y}, \xi')\|^2] \leq$
 567 $\alpha \|\mathbf{x} - \mathbf{y}\|^2 + \beta \sigma^2$ in Assumption 4, the notation σ^2 just represents a general positive constant. This
 568 σ^2 does not denote the variance of the stochastic gradient variance. Thus, this inequality does not

Table 3: List of key notation.

Notation	Definition
i	Client index
M	Total number of clients
s	Objective/task index
S	Total number of Objectives/tasks
S_i	Number of objectives/tasks of client i 's interest
k	Local step index
K	Total number of local steps
t	Communication round index
T	Total number of communication rounds
$\mathbf{x} \in \mathbb{R}^d$	Global model parameters of FMOL in Problem (2)
$\mathbf{x}_0 \in \mathbb{R}^d$	Initial solution of FMOL in Problem (2)
$\mathbf{x}_* \in \mathbb{R}^d$	A Pareto optimal solution of FMOL in Problem (2)
η_L	The learning rate on the client side
η_t	The learning rate on the server side in round t

569 depend on the batch size of the stochastic gradient. More specifically, unlike the assumption in [5]
 570 that characterizes the difference between a stochastic gradient and its full gradient (hence depending
 571 on the batch size), our Assumption 4 only measures the average norm square of two stochastic
 572 gradient difference $\nabla f(\mathbf{x}, \xi) - \nabla f(\mathbf{y}, \xi')$ given any two points \mathbf{x} and \mathbf{y} and any two samples ξ and
 573 ξ' . In other words, Assumption 4 does not involve any full gradient, and hence no dependence on
 574 batch size. In the revised version of this paper, we will replace σ^2 by a C to signify that it is a general
 575 constant.

576 It is a natural extension of the classic Lipschitz continuous gradient assumption and could generalize
 577 existing assumptions.

578 1. If ξ and ξ' are the whole dataset, by setting $\alpha = L^2$ and $\beta = 0$, (α, β) -Lipschitz continuous
 579 stochastic gradient condition generalizes the traditional Lipschitz continuous gradient assumption
 580 $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|$.

581 2. If ξ is one data sample, ξ' are the whole dataset and $\mathbf{x} = \mathbf{y}$, by setting $\alpha = 0$ and $\beta = 1$,
 582 (α, β) -Lipschitz continuous stochastic gradient condition generalizes the traditional bounded variance
 583 assumption $\|\nabla f(\mathbf{x}, \xi) - \nabla f(\mathbf{x})\|^2 \leq \sigma^2$.

584 3. If ξ is one data sample, ξ' are the whole dataset and $\mathbf{x} = \mathbf{y}$, by setting $\beta = \alpha_k$, (α, β) -Lipschitz
 585 continuous stochastic gradient condition generalizes the bound on the first moment assumption
 586 (assumption 5.2(b)) and bounded sets assumption (assumption 5.3) [5] ($\mathbb{E}[\|\nabla f(\mathbf{x}, \xi) - \nabla f(\mathbf{x})\|] \leq$
 587 $\alpha_k(C_i + \hat{C}_i\|\nabla f_i(\mathbf{x}_k)\|)$ and $\|\nabla f_i(\mathbf{x})\| \leq M_\nabla + L\Theta$).

588 **Metrics.** For strongly-convex functions, we use $\Delta_Q^t = \sum_{s \in [S]} \lambda_s^{t,*} [f_s(\mathbf{x}_t) - f_s(\mathbf{x}_*)]$ as
 589 the metrics. We note similar metrics are used in other works. For example, [5] uses
 590 $\min_{t=1, \dots, T} \sum_{s \in [S]} [\lambda_s^t f_s(\mathbf{x}_t) - \bar{\lambda}_T f_s(\mathbf{x}_*)]$ where $\bar{\lambda}_T = \sum_{t=1}^T \frac{t}{\sum_{t=1}^T t} \lambda_t$, and [14] utilizes
 591 $\sum_{s \in [S]} \lambda_s^{t,*} [f_s(\mathbf{x}_t) - f_s(\mathbf{x}_*)]$ as the metrics. In non-convex functions, $\|\mathbf{d}_t\|^2$ are used as the metrics
 592 for FMOO, where $\mathbf{d}_t = \lambda_t^T \nabla \mathbf{F}(\mathbf{x}_t)$ and λ_t is calculated based on accumulated (stochastic) gradients
 593 Δ_t . We note, directly extended from MOO [13, 14], $\mathbf{d}_t^* = \lambda_t^{*T} \nabla \mathbf{F}(\mathbf{x}_t)$ could also be used as the
 594 metrics in FMOO, where λ_t^* is calculated based on full gradients $\nabla \mathbf{F}(\mathbf{x}_t)$. However, we prefer \mathbf{d}_t
 595 for the following reasons: i). For applying gradient descent with no local steps, \mathbf{d}_t degenerates to
 596 \mathbf{d}_t^* . ii). Clearly, $\|\mathbf{d}_t\|^2 \leq \|\mathbf{d}_t^*\|^2$ as λ_t^* is calculated based on gradients $\nabla \mathbf{F}(\mathbf{x}_t)$. Hence, $\|\mathbf{d}_t\|^2$ is
 597 stronger convergence measure for FMOO. iii). λ_t is calculated in the algorithm and thus being more
 598 practical to use in practice, while λ_t^* is unknown. Also, the convergence of \mathbf{d}_t implicitly indicates λ_t
 599 converges to λ_t^* .

600 B Proof of gradient descent type methods

601 For gradient descent type methods, each step utilizes a full gradient to update and the corresponding
 602 parameter λ is deterministic. For clarity of notation, we drop $*$ for λ , that is, we use λ_t^s to represent
 603 the solution of quadratic problem (Step 6 in the algorithm) for task s in the t -th round.

604 **Lemma 1.** *Under bounded gradient assumption, the local model updates for any client s could be*
 605 *bounded*

$$G_{s,i}^{t,k} = \|\mathbf{x}_{s,i}^{t,k} - \mathbf{x}_t\|^2 \leq 4\eta_L^2 K^2 G^2, \quad (3)$$

$$H_{t,s} = \|\nabla f_s(\mathbf{x}_t) - \Delta_t^s\|^2 \leq 4\eta_L^2 K^2 L^2 G^2. \quad (4)$$

606 *Proof.* For one task $s \in [S]$ and one client $i \in R_s$, the local update $\|\mathbf{x}_t - \mathbf{x}_{s,i}^{t,k}\|^2$ could be further
 607 bounded.

$$\left\| \mathbf{x}_t - \mathbf{x}_{s,i}^{t,k} \right\|^2 = \left\| \mathbf{x}_t - \mathbf{x}_{s,i}^{t,k-1} + \eta_L \nabla f_{s,i}(\mathbf{x}_{s,i}^{t,k-1}) \right\|^2 \quad (5)$$

$$\leq \left(1 + \frac{1}{K-1}\right) \left\| \mathbf{x}_t - \mathbf{x}_{s,i}^{t,k-1} \right\|^2 + \eta_L^2 K \left\| \nabla f_{s,i}(\mathbf{x}_{s,i}^{t,k-1}) \right\|^2 \quad (6)$$

$$\leq \left(1 + \frac{1}{K-1}\right) \left\| \mathbf{x}_t - \mathbf{x}_{s,i}^{t,k-1} \right\|^2 + \eta_L^2 K G^2 \quad (7)$$

Table 4: Convergence rate (shaded parts are our results) for strongly-convex and non-convex functions, respectively:

Methods		Rate		Assumption
Setting	Algorithm	SC	NC	
Vanilla Gradient	MGD [4]	$\mathcal{O}(r^T), r \in (0, 1)$	$\mathcal{O}(\frac{1}{T})$	Sequence convergence
	MGD	$\mathcal{O}(\exp(-\mu T))$	$\mathcal{O}(\frac{1}{T})$	-
	SMGD [5]	$\mathcal{O}(\frac{1}{T})$	-	Lipschitz continuity of λ
	SMGD [48]	$\mathcal{O}(\frac{1}{T})$	-	Convexity property
	SMGD [33]	-	$\mathcal{O}(\frac{1}{\sqrt{T}})$	Given exact solution λ^*
	SMGD	$\tilde{\mathcal{O}}(\frac{1}{T})$	$\mathcal{O}(\frac{1}{\sqrt{T}})$	Asm. 4
Momentum	MoCo [14]	-	$\mathcal{O}(\frac{1}{\sqrt{T}})$	-
	CR-MOGM [13]	-	$\mathcal{O}(\frac{1}{\sqrt{T}})$	-
Federated Settings	FMGDA	$\mathcal{O}(\exp(-\mu T))$	$\mathcal{O}(\frac{1}{T})$	-
	FMSGDA	$\tilde{\mathcal{O}}(\frac{1}{T})$	$\mathcal{O}(\frac{1}{\sqrt{T}})$	Asm. 4

Assumptions. Linear search [4]: stepsize linear search; sequence convergence [4]: $\{\mathbf{x}_t\}$ converges to \mathbf{x}_* ; first moment bound (Asm. 5.2(b) [5]): $\mathbb{E}[\|\nabla f(\mathbf{x}, \xi) - \nabla f(\mathbf{x})\|] \leq \eta(a + b\|\nabla f(\mathbf{x})\|)$; Lipschitz continuity of λ (Asm. 5.4 [5]): $\|\lambda_k - \lambda_s\| \leq \beta \|[(\nabla f_1(\mathbf{x}_k) - \nabla f_1(\mathbf{x}_t))^T, \dots, (\nabla f_m(\mathbf{x}_k) - \nabla f_m(\mathbf{x}_t))^T]\|$; convexity property(H5) [48]: $f(\mathbf{x}, \xi) - f(\mathbf{x}^*, \xi) \geq \frac{\epsilon}{2} \|\mathbf{x} - \mathbf{x}^*\|^2$ almost sure; (α, β) -Lipschitz continuous stochastic gradient (Asm. 4).

$$\leq \sum_{\tau \in [k-1]} (2\eta_L^2 K G^2) \left(1 + \frac{1}{K-1}\right)^\tau \quad (8)$$

$$\leq (K-1) \left[\left(1 + \frac{1}{K-1}\right)^K - 1 \right] (\eta_L^2 K G^2) \quad (9)$$

$$\leq 4\eta_L^2 K^2 G^2, \quad (10)$$

608 where the first inequality comes from Young's inequality, the second inequality follows from bounded
609 gradient assumption, and the last inequality follows if $\left(1 + \frac{1}{K-1}\right)^K - 1 \leq 4$ for $K > 1$.

610 We have the bound for local update for each task s , $H_{t,s}$, as follows:

$$H_{t,s} = \|\nabla f_s(\mathbf{x}_t) - \Delta_t^s\|^2 \quad (11)$$

$$= \left\| \frac{1}{K} \sum_{k \in [K]} \frac{1}{|R_s|} \sum_{i \in R_s} [\nabla f_{s,i}(\mathbf{x}_t) - \nabla f_{s,i}(\mathbf{x}_{s,i}^{t,k})] \right\|^2 \quad (12)$$

$$\leq \frac{1}{K} \sum_{k \in [K]} \frac{1}{|R_s|} \sum_{i \in R_s} \|\nabla f_{s,i}(\mathbf{x}_t) - \nabla f_{s,i}(\mathbf{x}_{s,i}^{t,k})\|^2 \quad (13)$$

$$\leq \frac{1}{K} L^2 \sum_{k \in [K]} \frac{1}{|R_s|} \sum_{i \in R_s} \|\mathbf{x}_t - \mathbf{x}_{s,i}^{t,k}\|^2 \quad (14)$$

$$\leq 4\eta_L^2 K^2 L^2 G^2. \quad (15)$$

611

□

612 **Lemma 2.** For general L -smooth functions $\{f_s, s \in [S]\}$, choose the learning rate η_t s.t. $\eta_t \leq$
613 $\frac{3}{2(1+L)}$, the update d_t of the algorithm satisfies:

$$\frac{\eta_t}{4} \|\mathbf{d}_t\|^2 \leq -f_s(\mathbf{x}_{t+1}) + f_s(\mathbf{x}_t) + 6\eta_L^2 K^2 L^2 G^2 \quad (16)$$

Proof.

$$f_s(\mathbf{x}_{t+1}) \leq f_s(\mathbf{x}_t) + \langle \nabla f_s(\mathbf{x}_t), -\eta_t \mathbf{d}_t \rangle + \frac{1}{2} L \|\eta_t \mathbf{d}_t\|^2 \quad (17)$$

$$= f_s(\mathbf{x}_t) + \langle \nabla f_s(\mathbf{x}_t) - \Delta_t^s, -\eta_t \mathbf{d}_t \rangle - \eta_t \langle \Delta_t^s, \mathbf{d}_t \rangle + \frac{1}{2} L \|\eta_t \mathbf{d}_t\|^2 \quad (18)$$

$$\leq f_s(\mathbf{x}_t) + \langle \nabla f_s(\mathbf{x}_t) - \Delta_t^s, -\eta_t \mathbf{d}_t \rangle - \eta_t \|\mathbf{d}_t\|^2 + \frac{1}{2} L \|\eta_t \mathbf{d}_t\|^2 \quad (19)$$

$$\leq f_s(\mathbf{x}_t) + \frac{1}{2} \|\nabla f_s(\mathbf{x}_t) - \Delta_t^s\|^2 + \frac{1}{2} \eta_t^2 \|\mathbf{d}_t\|^2 - \eta_t \|\mathbf{d}_t\|^2 + \frac{1}{2} L \eta_t^2 \|\mathbf{d}_t\|^2 \quad (20)$$

$$= f_s(\mathbf{x}_t) + \frac{1}{2} \|\nabla f_s(\mathbf{x}_t) - \Delta_t^s\|^2 - \eta_t \left(1 - \frac{1}{2} \eta_t - \frac{1}{2} L \eta_t\right) \|\mathbf{d}_t\|^2 \quad (21)$$

$$\leq f_s(\mathbf{x}_t) + 2\eta_L^2 K^2 L^2 G^2 - \eta_t \left(1 - \frac{1}{2} \eta_t - \frac{1}{2} L \eta_t\right) \|\mathbf{d}_t\|^2. \quad (22)$$

614 The third inequality follows from $\langle \Delta_t^s, \mathbf{d}_t \rangle \geq \|\mathbf{d}_t\|^2$ since \mathbf{d}_t is a general solution in the convex hull
615 of the family of vectors $\{\Delta_t^s, s \in [S]\}$ (see Lemma 2.1 [11])

616 By setting $(1 - \frac{1}{2} \eta_t - \frac{1}{2} L \eta_t) \geq \frac{1}{4}$, that is, $\eta_t \leq \frac{3}{2(1+L)}$, we have

$$\frac{\eta_t}{4} \|\mathbf{d}_t\|^2 \leq -f_s(\mathbf{x}_{t+1}) + f_s(\mathbf{x}_t) + 2\eta_L^2 K^2 L^2 G^2. \quad (23)$$

617

□

618 B.1 Strongly Convex Functions

619 **Theorem 3** (FMGDA for μ -Strongly Convex FMOL). *Let $\eta_t = \eta$ such that $\eta \leq \frac{3}{2(1+L)}$, $\eta \leq \frac{1}{2L+\mu}$
620 and $\eta \geq \frac{1}{\mu T}$. Under Assumptions 1- 3, pick \mathbf{x}_t as the final output of the FMGDA algorithm with
621 weights $w_t = (1 - \frac{\mu\eta}{2})^{1-t}$. Then, it holds that $\mathbb{E}[\Delta_Q^t] \leq \|\mathbf{x}_0 - \mathbf{x}_*\|^2 \mu \exp(-\frac{\eta\mu T}{2}) + \delta$, where
622 $\Delta_Q^t \triangleq \sum_{s \in [S]} \lambda_s^{t,*} [f_s(\mathbf{x}_t) - f_s(\mathbf{x}_*)]$ and $\delta = \frac{8\eta_L^2 K^2 L^2 G^2 S^2}{\mu} + 2\eta_L^2 K^2 L^2 G^2$.*

Proof.

$$f_s(\mathbf{x}_{t+1}) \leq f_s(\mathbf{x}_t) + \langle \nabla f_s(\mathbf{x}_t), -\eta_t \mathbf{d}_t \rangle + \frac{1}{2} L \|\eta_t \mathbf{d}_t\|^2 \quad (24)$$

$$\leq f_s(\mathbf{x}_*) + \langle \nabla f_s(\mathbf{x}_t), \mathbf{x}_t - \mathbf{x}_* \rangle - \frac{\mu}{2} \|\mathbf{x}_t - \mathbf{x}_*\|^2 \quad (25)$$

$$+ \langle \nabla f_s(\mathbf{x}_t), -\eta_t \mathbf{d}_t \rangle + \frac{1}{2} L \|\eta_t \mathbf{d}_t\|^2, \quad (26)$$

623 where the first inequality is due to L -smoothness, the second inequality follows from μ -strongly
624 convex.

$$\sum_{s \in [S]} \lambda_s^s [f_s(\mathbf{x}_{t+1}) - f_s(\mathbf{x}_*)] \quad (27)$$

$$\leq \left\langle \sum_{s \in [S]} \lambda_s^s \nabla f_s(\mathbf{x}_t), \mathbf{x}_t - \mathbf{x}_* \right\rangle - \frac{\mu}{2} \|\mathbf{x}_t - \mathbf{x}_*\|^2 + \left\langle \sum_{s \in [S]} \lambda_s^s \nabla f_s(\mathbf{x}_t), -\eta_t \mathbf{d}_t \right\rangle + \frac{1}{2} L \|\eta_t \mathbf{d}_t\|^2 \quad (28)$$

$$= \left\langle \sum_{s \in [S]} \lambda_s^s \nabla f_s(\mathbf{x}_t), \mathbf{x}_t - \mathbf{x}_* - \eta_t \mathbf{d}_t \right\rangle - \frac{\mu}{2} \|\mathbf{x}_t - \mathbf{x}_*\|^2 + \frac{1}{2} L \|\eta_t \mathbf{d}_t\|^2 \quad (29)$$

$$= \langle \mathbf{d}_t, \mathbf{x}_t - \mathbf{x}_* - \eta_t \mathbf{d}_t \rangle - \frac{\mu}{2} \|\mathbf{x}_t - \mathbf{x}_*\|^2 + \frac{1}{2} L \|\eta_t \mathbf{d}_t\|^2 + \left\langle \sum_{s \in [S]} \lambda_s^s \nabla f_s(\mathbf{x}_t) - \mathbf{d}_t, \mathbf{x}_t - \mathbf{x}_* - \eta_t \mathbf{d}_t \right\rangle \quad (30)$$

$$= \langle \mathbf{d}_t, \mathbf{x}_t - \mathbf{x}_* \rangle - \eta_t \|\mathbf{d}_t\|^2 - \frac{\mu}{2} \|\mathbf{x}_t - \mathbf{x}_*\|^2 + \frac{1}{2} L \eta_t^2 \|\mathbf{d}_t\|^2 + \left\langle \sum_{s \in [S]} \lambda_t^s \nabla f_s(\mathbf{x}_t) - \mathbf{d}_t, \mathbf{x}_{t+1} - \mathbf{x}_* \right\rangle \quad (31)$$

$$\leq \frac{1}{2\eta_t} (\|\mathbf{x}_t - \mathbf{x}_*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}_*\|^2) - \frac{1}{2} \eta_t \|\mathbf{d}_t\|^2 - \frac{\mu}{2} \|\mathbf{x}_t - \mathbf{x}_*\|^2 + \frac{1}{2} L \eta_t^2 \|\mathbf{d}_t\|^2 \quad (32)$$

$$+ \frac{1}{4\epsilon} \left\| \underbrace{\sum_{s \in [S]} \lambda_t^s \nabla f_s(\mathbf{x}_t) - \mathbf{d}_t}_{H_t} \right\|^2 + \epsilon \|\mathbf{x}_{t+1} - \mathbf{x}_*\|^2 \quad (33)$$

$$\leq \frac{1}{2\eta_t} (\|\mathbf{x}_t - \mathbf{x}_*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}_*\|^2) - \frac{1}{2} \eta_t \|\mathbf{d}_t\|^2 - \frac{\mu}{2} \|\mathbf{x}_t - \mathbf{x}_*\|^2 + \frac{1}{2} L \eta_t^2 \|\mathbf{d}_t\|^2 \quad (34)$$

$$+ \frac{1}{4\epsilon} H_t + \epsilon (2 \|\mathbf{x}_t - \mathbf{x}_*\|^2 + 2\eta_t^2 \|\mathbf{d}_t\|^2) \quad (35)$$

$$\leq \frac{1}{2\eta_t} \left((1 - \frac{\mu}{2} \eta_t) \|\mathbf{x}_t - \mathbf{x}_*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}_*\|^2 \right) - \left(\frac{1}{2} \eta_t - \frac{1}{2} L \eta_t^2 - \frac{\mu}{4} \eta_t^2 \right) \|\mathbf{d}_t\|^2 + \frac{2}{\mu} H_t, \quad (36)$$

625 where $\|\mathbf{x}_t - \mathbf{x}_*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}_*\|^2 = -\eta_t^2 \|\mathbf{d}_t\|^2 + 2 \langle \eta_t \mathbf{d}_t, \mathbf{x}_t - \mathbf{x}_* \rangle$, and we choose $\epsilon = \frac{\mu}{8}$ in the
626 last inequality.

627 From Lemma 2, it is clear that

$$|f_s(\mathbf{x}_{t+1}) - f_s(\mathbf{x}_t)| \leq |2\eta_L^2 K^2 L^2 G^2 - \frac{\eta_t}{4} \|\mathbf{d}_t\|^2| \quad (37)$$

$$\leq 2\eta_L^2 K^2 L^2 G^2 + \frac{\eta_t}{4} \|\mathbf{d}_t\|^2. \quad (38)$$

$$\Delta_Q^t = \sum_{s \in [S]} \lambda_t^s [f_s(\mathbf{x}_t) - f_s(\mathbf{x}_*)] \leq \sum_{s \in [S]} \lambda_t^s [f_s(\mathbf{x}_{t+1}) - f_s(\mathbf{x}_*)] + |f_s(\mathbf{x}_{t+1}) - f_s(\mathbf{x}_t)| \quad (39)$$

$$\leq \frac{1}{2\eta_t} \left((1 - \frac{\mu}{2} \eta_t) \|\mathbf{x}_t - \mathbf{x}_*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}_*\|^2 \right) - \left(\frac{1}{4} \eta_t - \frac{1}{2} L \eta_t^2 - \frac{\mu}{4} \eta_t^2 \right) \|\mathbf{d}_t\|^2 + \frac{2}{\mu} H_t + 2\eta_L^2 K^2 L^2 G^2. \quad (40)$$

$$H_t = \left\| \sum_{s \in [S]} \lambda_t^s \nabla f_s(\mathbf{x}_t) - \mathbf{d}_t \right\|^2 \quad (41)$$

$$\leq S \sum_{s \in [S]} (\lambda_t^s)^2 H_{t,s} \quad (42)$$

$$\leq 4\eta_L^2 K^2 L^2 G^2 S^2. \quad (43)$$

628 By setting $\eta_t \leq \frac{1}{2L+\mu}$, we have

$$\Delta_Q^t = \sum_{s \in [S]} \lambda_t^s [f_s(\mathbf{x}_{t+1}) - f_s(\mathbf{x}_*)] \quad (44)$$

$$\leq \frac{1}{2\eta_t} \left((1 - \frac{\mu}{2} \eta_t) \|\mathbf{x}_t - \mathbf{x}_*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}_*\|^2 \right) + \underbrace{\frac{8\eta_L^2 K^2 L^2 G^2 S^2}{\mu} + 2\eta_L^2 K^2 L^2 G^2}_{\delta}. \quad (45)$$

629 Averaging using weight $w_t = (1 - \frac{\mu\eta}{2})^{1-t}$ and using such weight to pick output \mathbf{x} . By using Lemma
630 1 in [21] with $\eta \geq \frac{1}{uR}$, we have

$$\mathbb{E}[\Delta_Q] \leq \|\mathbf{x}_0 - \mathbf{x}_*\|^2 \mu \exp\left(-\frac{\eta\mu T}{2}\right) + \delta \quad (46)$$

$$= \mathcal{O}(\mu \exp(-\mu T)) + \mathcal{O}(\delta). \quad (47)$$

631 If we set η_L sufficiently small such that $\delta = \mathcal{O}(\mu \exp(-\mu T))$, then we have the convergence rate
 632 $\mathbb{E}[\Delta_Q] = \mathcal{O}(\mu \exp(-\mu T))$. \square

633 B.2 Non-Convex Functions

634 **Theorem 1** (FMGDA for Non-convex FMOL). *Let $\eta_t = \eta \leq \frac{3}{2(1+L)}$. Under Assumptions 1 and 2,*
 635 *if at least one function $f_s, s \in [S]$ is bounded from below by f_s^{\min} , then the sequence $\{\mathbf{x}_t\}$ output by*
 636 *FMGDA satisfies: $\min_{t \in [T]} \|\mathbf{d}_t\|^2 \leq \frac{4(f_s^0 - f_s^{\min})}{T\eta} + \delta$, where $\delta \triangleq (8\eta_L^2 K^2 L^2 G^2)/\eta$.*

637 *Proof.* From Lemma 2, we have

$$\frac{\eta_t}{4} \|\mathbf{d}_t\|^2 \leq -f_s(\mathbf{x}_{t+1}) + f_s(\mathbf{x}_t) + 2\eta_L^2 K^2 L^2 G^2. \quad (48)$$

638 With constant learning rate $\eta_t = \eta$,

$$\frac{1}{T} \sum_{t \in [T]} \|\mathbf{d}_t\|^2 \leq \frac{4(f_s^0 - f_s^{\min})}{T\eta} + \frac{8\eta_L^2 K^2 L^2 G^2}{\eta}. \quad (49)$$

639 With constant learning rate η and local learning rate $\eta_L = \mathcal{O}(\frac{1}{\sqrt{TKLG}})$, we have

$$\frac{1}{T} \sum_{t \in [T]} \|\mathbf{d}_t\|^2 \leq \mathcal{O}\left(\frac{1}{T}\right) \quad (50)$$

640 \square

641 C Proof of stochastic gradient descent type methods

642 For stochastic gradient descent type methods, each step utilizes a stochastic gradient to update and
 643 the corresponding parameter λ is stochastic, depending on the random samples in each client. For
 644 clarity of notation, we drop $*$ for λ , that is, we use λ_i^s to represent the solution of quadratic problem
 645 (Step 6 in the algorithm) for task s in the t -th round.

646 **Lemma 3.** *Under bounded stochastic gradient assumption, the local model updates could be bounded*

$$G_{s,i}^{t,k} = \mathbb{E} \|\mathbf{x}_{s,i}^{t,k} - \mathbf{x}_t\|^2 \leq 6\eta_L^2 k^2 \|\nabla f_{s,i}(\mathbf{x}_t)\|^2, \quad (51)$$

$$\mathbb{E} \left\| \sum_{s \in [S]} \lambda_s^t \Delta_s^t \right\|^2 \leq S^2 D^2. \quad (52)$$

647 *Further with assumption 4, we have*

$$H_{t,s} = \mathbb{E} \|\nabla f_s(\mathbf{x}_t, \xi_t) - \Delta_s^t\|^2 \leq \alpha \eta_L^2 K^2 D^2 + \beta \sigma^2. \quad (53)$$

648 *Proof.* For one task $s \in [S]$ and one client $i \in R_s$, the local update $\|\mathbf{x}_t - \mathbf{x}_{s,i}^{t,k}\|^2$ could be further
 649 bounded.

$$G_{s,i}^{t,k} = \mathbb{E} \left\| \mathbf{x}_t - \mathbf{x}_{s,i}^{t,k} \right\|^2 \quad (54)$$

$$= \mathbb{E} \left\| \sum_{\tau \in [k]} \eta_L \nabla f_{s,i}(\mathbf{x}_{s,i}^{t,\tau}, \xi_{s,i}^{t,\tau}) \right\|^2 \quad (55)$$

$$\leq \eta_L^2 k^2 D^2. \quad (56)$$

$$\mathbb{E} \left\| \sum_{s \in [S]} \lambda_s^t \Delta_s^t \right\|^2 \leq S \sum_{s \in [S]} \mathbb{E} \left[(\lambda_s^t)^2 \|\Delta_s^t\|^2 \right] \quad (57)$$

$$\leq S \sum_{s \in [S]} \mathbb{E} \left[\|\Delta_s^t\|^2 \right] \quad (58)$$

$$\leq S \sum_{s \in [S]} \mathbb{E} \left\| \frac{1}{R_s} \sum_{i \in R_s} \frac{1}{K} \sum_{\tau \in [K]} \nabla f_{s,i}(\mathbf{x}_{s,i}^{t,\tau}, \xi_{s,i}^{t,\tau}) \right\|^2 \quad (59)$$

$$\leq S \sum_{s \in [S]} \frac{1}{R_s} \sum_{i \in R_s} \frac{1}{K} \sum_{\tau \in [K]} \mathbb{E} \|\nabla f_{s,i}(\mathbf{x}_{s,i}^{t,\tau}, \xi_{s,i}^{t,\tau})\|^2 \quad (60)$$

$$\leq S^2 D^2. \quad (61)$$

$$H_t = \mathbb{E} \|\nabla f_s(\mathbf{x}_t, \xi_t) - \Delta_s^t\|^2 \quad (62)$$

$$\leq \mathbb{E} \left\| \frac{1}{K} \sum_{k \in [K]} \frac{1}{|R_s|} \sum_{i \in R_s} \left(\nabla f_{s,i}(\mathbf{x}_t, \xi_t) - \nabla f_{s,i}(\mathbf{x}_{s,i}^{t,k}, \xi_{s,i}^{t,k}) \right) \right\|^2 \quad (63)$$

$$\leq \frac{1}{K} \sum_{k \in [K]} \frac{1}{|R_s|} \sum_{i \in R_s} \mathbb{E} \|\nabla f_{s,i}(\mathbf{x}_t, \xi_t) - \nabla f_{s,i}(\mathbf{x}_{s,i}^{t,k}, \xi_{s,i}^{t,k})\|^2 \quad (64)$$

$$\leq \frac{1}{K} \sum_{k \in [K]} \frac{1}{|R_s|} \sum_{i \in R_s} \left(\alpha \mathbb{E} \|\mathbf{x}_t - \mathbf{x}_{s,i}^{t,k}\|^2 + \beta \sigma^2 \right) \quad (65)$$

$$\leq \alpha \eta_L^2 K^2 D^2 + \beta \sigma^2. \quad (66)$$

650

□

651 C.1 Strongly Convex Functions

652 **Theorem 7** (FSMGDA for μ -Strongly Convex FMOL). *Let $\eta_t = \eta = \Omega(\frac{1}{\mu T})$. Under Assumptions 3,*
 653 *5 and 6, pick \mathbf{x}_t as the final output of the FSMGDA algorithm with weight $w_t = (1 - \frac{\mu\eta}{2})^{1-t}$. Then,*
 654 *it holds that: $\mathbb{E}[\Delta_Q^t] \leq \|\mathbf{x}_0 - \mathbf{x}_*\|^2 \mu \exp(-\frac{\eta}{2}\mu T) + \delta$, where $\Delta_Q^t = \sum_{s \in [S]} \lambda_s^{t,*} [f_s(\mathbf{x}_t) - f_s(\mathbf{x}_*)]$*
 655 *and $\delta = \frac{1}{\mu} S^2 (\alpha \eta_L^2 K^2 D^2 + \beta \sigma^2) + \frac{\eta S^2 D^2}{2}$.*

656 *Proof.* Taking expectation over random samples conditioning on \mathbf{x}_t , we have

$$\mathbb{E} \|\mathbf{x}_{t+1} - \mathbf{x}_*\|^2 = \mathbb{E} \left\| \mathbf{x}_t - \eta_t \sum_{s \in [S]} \lambda_s^t \Delta_s^t - \mathbf{x}_* \right\|^2 \quad (67)$$

$$= \|\mathbf{x}_t - \mathbf{x}_*\|^2 - \mathbb{E} \left\langle \mathbf{x}_t - \mathbf{x}_*, 2\eta_t \sum_{s \in [S]} \lambda_s^t \Delta_s^t \right\rangle + \mathbb{E} \left\| \eta_t \sum_{s \in [S]} \lambda_s^t \Delta_s^t \right\|^2 \quad (68)$$

$$= \|\mathbf{x}_t - \mathbf{x}_*\|^2 - \mathbb{E} \left\langle \mathbf{x}_t - \mathbf{x}_*, 2\eta_t \sum_{s \in [S]} \lambda_s^t \nabla f_s(\mathbf{x}_t, \xi_t) \right\rangle \quad (69)$$

$$+ \mathbb{E} \left\langle \mathbf{x}_t - \mathbf{x}_*, 2\eta_t \sum_{s \in [S]} \lambda_s^t (\nabla f_s(\mathbf{x}_t, \xi_t) - \Delta_s^t) \right\rangle + \mathbb{E} \left\| \eta_t \sum_{s \in [S]} \lambda_s^t \Delta_s^t \right\|^2 \quad (70)$$

$$= \|\mathbf{x}_t - \mathbf{x}_*\|^2 - \left\langle \mathbf{x}_t - \mathbf{x}_*, 2\eta_t \sum_{s \in [S]} \mathbb{E}[\lambda_s^t] \nabla f_s(\mathbf{x}_t) \right\rangle \quad (71)$$

$$+ \mathbb{E} \left\langle \mathbf{x}_t - \mathbf{x}_*, 2\eta_t \sum_{s \in [S]} \lambda_s^t (\nabla f_s(\mathbf{x}_t, \xi_t) - \Delta_s^t) \right\rangle + \mathbb{E} \left\| \eta_t \sum_{s \in [S]} \lambda_s^t \Delta_s^t \right\|^2 \quad (72)$$

$$\leq \|\mathbf{x}_t - \mathbf{x}_*\|^2 - 2\eta_t \left(\frac{\mu}{2} \|\mathbf{x}_t - \mathbf{x}_*\|^2 + \sum_{s \in [S]} \mathbb{E}[\lambda_s^t] (f_s(\mathbf{x}_t) - f_s(\mathbf{x}_*)) \right) + \epsilon \|\mathbf{x}_t - \mathbf{x}_*\|^2 \quad (73)$$

$$+ \frac{1}{4\epsilon} 4\eta_t^2 \mathbb{E} \left\| \sum_{s \in [S]} \lambda_s^t (\nabla f_s(\mathbf{x}_t, \xi_t) - \Delta_s^t) \right\|^2 + \eta_t^2 \mathbb{E} \left\| \sum_{s \in [S]} \lambda_s^t \Delta_s^t \right\|^2 \quad (74)$$

$$\leq \|\mathbf{x}_t - \mathbf{x}_*\|^2 - 2\eta_t \left(\frac{\mu}{2} \|\mathbf{x}_t - \mathbf{x}_*\|^2 + \sum_{s \in [S]} \mathbb{E}[\lambda_s^t] (f_s(\mathbf{x}_t) - f_s(\mathbf{x}_*)) \right) + \epsilon \|\mathbf{x}_t - \mathbf{x}_*\|^2 \quad (75)$$

$$+ \frac{1}{4\epsilon} 4\eta_t^2 S \sum_{s \in [S]} \mathbb{E} \left[(\lambda_s^t)^2 \|\nabla f_s(\mathbf{x}_t, \xi_t) - \Delta_s^t\|^2 \right] + \eta_t^2 \mathbb{E} \left\| \sum_{s \in [S]} \lambda_s^t \Delta_s^t \right\|^2 \quad (76)$$

$$\leq \|\mathbf{x}_t - \mathbf{x}_*\|^2 - 2\eta_t \left(\frac{\mu}{2} \|\mathbf{x}_t - \mathbf{x}_*\|^2 + \sum_{s \in [S]} \mathbb{E}[\lambda_s^t] (f_s(\mathbf{x}_t) - f_s(\mathbf{x}_*)) \right) + \epsilon \|\mathbf{x}_t - \mathbf{x}_*\|^2 \quad (77)$$

$$+ \frac{1}{4\epsilon} 4\eta_t^2 S \sum_{s \in [S]} \mathbb{E} \|\nabla f_s(\mathbf{x}_t, \xi_t) - \Delta_s^t\|^2 + \eta_t^2 \mathbb{E} \left\| \sum_{s \in [S]} \lambda_s^t \Delta_s^t \right\|^2 \quad (78)$$

$$\leq \|\mathbf{x}_t - \mathbf{x}_*\|^2 - 2\eta_t \left(\frac{\mu}{2} \|\mathbf{x}_t - \mathbf{x}_*\|^2 + \sum_{s \in [S]} \mathbb{E}[\lambda_s^t] (f_s(\mathbf{x}_t) - f_s(\mathbf{x}_*)) \right) + \epsilon \|\mathbf{x}_t - \mathbf{x}_*\|^2 \quad (79)$$

$$+ \frac{1}{4\epsilon} 4\eta_t^2 S^2 (\alpha \eta_L^2 K^2 D^2 + \beta \sigma^2) + \eta_t^2 S^2 D^2 \quad (80)$$

$$\leq (1 - \frac{\eta_t \mu}{2}) \|\mathbf{x}_t - \mathbf{x}_*\|^2 - 2\eta_t \left(\sum_{s \in [S]} \mathbb{E}[\lambda_s^t] (f_s(\mathbf{x}_t) - f_s(\mathbf{x}_*)) \right) \quad (81)$$

$$+ \frac{2}{\mu} \eta_t S^2 (\alpha \eta_L^2 K^2 D^2 + \beta \sigma^2) + \eta_t^2 S^2 D^2, \quad (82)$$

657 where the first equality is due to strongly-convex objective functions, and we set $\epsilon = \frac{\eta_t \mu}{2}$.

$$\sum_{s \in [S]} \mathbb{E}[\lambda_s^t] (f_s(\mathbf{x}) - f_s(\mathbf{x}_*)) \leq \frac{1}{2\eta_t} (1 - \frac{\eta_t \mu}{2}) \|\mathbf{x}_t - \mathbf{x}_*\|^2 - \frac{1}{2\eta_t} \|\mathbf{x}_{t+1} - \mathbf{x}_*\|^2 \quad (83)$$

$$+ \underbrace{\frac{1}{\mu} S^2 (\alpha \eta_L^2 K^2 D^2 + \beta \sigma^2) + \frac{\eta_t S^2 D^2}{2}}_{\delta} \quad (84)$$

658 Averaging using weight $w_t = (1 - \frac{\mu \eta_t}{2})^{1-t}$ and using such weight to pick output \mathbf{x} . By using Lemma
659 1 in [21] with constant learning rate $\eta_t = \eta = \Omega(\frac{1}{\mu T})$, we have

$$\mathbb{E}[\Delta_Q] \leq \|\mathbf{x}_0 - \mathbf{x}_*\|^2 \mu \exp(-\frac{\eta}{2} \mu T) + \mathcal{O}(\delta) \quad (85)$$

660 where $\delta = \frac{1}{\mu} S^2 (\alpha \eta_L^2 K^2 D^2 + \beta \sigma^2) + \frac{\eta S^2 D^2}{2}$.

661 By letting $\beta = \eta$, $\eta_L = \mathcal{O}(\frac{1}{\sqrt{T}})$ and $\eta = \Theta(\frac{\log(\max(1, \mu^2 T))}{\mu T})$,

$$\mathbb{E}[\Delta_Q] \leq \tilde{\mathcal{O}}(\frac{1}{T}). \quad (86)$$

662

□

663 C.2 Non-convex Functions

664 **Theorem 5** (FSMGDA for Non-convex FMOL). *Let $\eta_t = \eta \leq \frac{3}{2(1+L)}$. Under Assumptions 4–6, if*

665 *an objective f_s is bounded from below by f_s^{\min} , then the sequence $\{\mathbf{x}_t\}$ output by FSMGDA satisfies:*

666 $\min_{t \in [T]} \mathbb{E} \|\mathbf{d}_t\|^2 \leq \frac{2S(f_s^0 - f_s^{\min})}{\eta T} + \delta$, where $\delta = L\eta S^2 D^2 + S(\alpha\eta_L^2 K^2 D^2 + \beta\sigma^2)$.

667 *Proof.* Taking expectation on the random data samples conditioning on \mathbf{x}_t , we have

$$\mathbb{E} f_s(\mathbf{x}_{t+1}) \leq f_s(\mathbf{x}_t) + \mathbb{E} \left\langle \nabla f_s(\mathbf{x}_t), -\eta_t \sum_{j \in [S]} \lambda_j^t \Delta_j^t \right\rangle + \frac{1}{2} L \mathbb{E} \left\| \eta_t \sum_{j \in [S]} \lambda_j^t \Delta_j^t \right\|^2 \quad (87)$$

$$= f_s(\mathbf{x}_t) + \mathbb{E} \left\langle \nabla f_s(\mathbf{x}_t), -\eta_t \sum_{j \in [S]} \lambda_j^t \nabla f_j(\mathbf{x}_t, \xi_t) \right\rangle \quad (88)$$

$$+ \eta_t \mathbb{E} \left\langle \nabla f_s(\mathbf{x}_t), \sum_{j \in [S]} \lambda_j^t [-\Delta_j^t + \nabla f_j(\mathbf{x}_t, \xi_t)] \right\rangle + \frac{1}{2} L \eta_t^2 \mathbb{E} \left\| \sum_{j \in [S]} \lambda_j^t \Delta_j^t \right\|^2 \quad (89)$$

$$\leq f_s(\mathbf{x}_t) - \eta_t \sum_{j \in [S]} \mathbb{E}[\lambda_j^t] \|\nabla f_j(\mathbf{x}_t)\|^2 \quad (90)$$

$$+ \eta_t \mathbb{E} \left\langle \nabla f_s(\mathbf{x}_t), \sum_{j \in [S]} \lambda_j^t [-\Delta_j^t + \nabla f_j(\mathbf{x}_t, \xi_t)] \right\rangle + \frac{1}{2} L \eta_t^2 \mathbb{E} \left\| \sum_{j \in [S]} \lambda_j^t \Delta_j^t \right\|^2 \quad (91)$$

$$\leq f_s(\mathbf{x}_t) - \eta_t \sum_{j \in [S]} \mathbb{E}[\lambda_j^t] \|\nabla f_j(\mathbf{x}_t)\|^2 + \frac{\eta_t}{2} S \mathbb{E} \|\lambda_s^t \nabla f_s(\mathbf{x}_t)\|^2 \quad (92)$$

$$+ \frac{\eta_t}{2} \sum_{j \in [S]} \mathbb{E} \|\nabla f_j(\mathbf{x}_t, \xi_t) - \Delta_j^t\|^2 + \frac{1}{2} L \eta_t^2 \mathbb{E} \left\| \sum_{j \in [S]} \lambda_j^t \Delta_j^t \right\|^2 \quad (93)$$

$$\leq f_s(\mathbf{x}_t) - \frac{\eta_t}{2} \sum_{j \in [S]} \mathbb{E} \|\lambda_j^t \nabla f_j(\mathbf{x}_t)\|^2 + \frac{\eta_t}{2} \sum_{j \in [S]} \mathbb{E} \|\nabla f_j(\mathbf{x}_t, \xi_t) - \Delta_j^t\|^2 + \frac{1}{2} L \eta_t^2 \mathbb{E} \left\| \sum_{j \in [S]} \lambda_j^t \Delta_j^t \right\|^2. \quad (94)$$

668 Here we construct a virtual stochastic gradient $\nabla f_s(\mathbf{x}_t, \xi_t)$ with an independent sample. As

669 λ_s^t only depends on Δ_s^t , so λ_s^t and $\nabla f_s(\mathbf{x}_t, \xi_t)$ are independent, from which the first in-

670 equality follows. The second inequality is due to $ab \leq \frac{1}{2}a^2 + \frac{1}{2}b^2$. Specifically,

671 $\mathbb{E} \left\langle \nabla f_s(\mathbf{x}_t), \sum_{j \in [S]} \lambda_j^t (-\Delta_j^t + \nabla f_j(\mathbf{x}_t, \xi_t)) \right\rangle = \sum_{j \in [S]} \mathbb{E} \left\langle \lambda_s^t \nabla f_s(\mathbf{x}_t), -\Delta_j^t + \nabla f_j(\mathbf{x}_t, \xi_t) \right\rangle \leq$

672 $\frac{S}{2} \mathbb{E} \|\lambda_s^t \nabla f_s(\mathbf{x}_t)\|^2 + \frac{1}{2} \sum_{s \in [S]} \mathbb{E} \|\nabla f_s(\mathbf{x}_t, \xi_t) - \Delta_s^t\|^2$. Also, following the fact that

673 $\lambda_s^t \in [0, 1]$, we have $\eta_t \sum_{s \in [S]} \mathbb{E}[\lambda_s^t] \|\nabla f_s(\mathbf{x}_t)\|^2 \geq \eta_t \sum_{s \in [S]} \mathbb{E}[(\lambda_s^t)^2] \|\nabla f_s(\mathbf{x}_t)\|^2 =$

674 $\eta_t \sum_{s \in [S]} \mathbb{E} \|\lambda_s^t \nabla f_s(\mathbf{x}_t)\|^2$. We also note that there exist a task s , such that $\frac{S}{2} \mathbb{E} \|\lambda_s^t \nabla f_s(\mathbf{x}_t)\|^2 \leq$

675 $\frac{1}{2} \sum_{j \in [S]} \mathbb{E} \|\lambda_j^t \nabla f_j(\mathbf{x}_t)\|^2$, which leads to the last inequality.

676 Rearranging the terms, we have

$$\sum_{s \in [S]} \mathbb{E} \|\lambda_s^t \nabla f_s(\mathbf{x}_t)\|^2 \leq \frac{2(f_s(\mathbf{x}_t) - \mathbb{E}f_s(\mathbf{x}_{t+1}))}{\eta_t} + \sum_{s \in [S]} \mathbb{E} \|\nabla f_s(\mathbf{x}_t, \xi_t) - \Delta_s^t\|^2 + L\eta_t \mathbb{E} \left\| \sum_{j \in [S]} \lambda_j^t \Delta_j^t \right\|^2. \quad (95)$$

677 With constant learning rate $\eta_t = \eta$ and averaging from T communication rounds, we have

$$\frac{1}{T} \sum_{t \in [T]} \mathbb{E} \|\mathbf{d}_t\|^2 \leq \frac{1}{T} \sum_{t \in [T]} S \sum_{s \in [S]} \mathbb{E} \|\lambda_s^t \nabla f_s(\mathbf{x}_t)\|^2 \quad (96)$$

$$\leq \frac{2S(f_s(\mathbf{x}_1) - \mathbb{E}f_s(\mathbf{x}_{T+1}))}{\eta T} + S(\alpha\eta_L^2 K^2 D^2 + \beta\sigma^2) + L\eta S^2 D^2. \quad (97)$$

678 With constant learning rate $\eta = \frac{1}{\sqrt{T}}$, local learning rate $\eta_L = \mathcal{O}(\frac{1}{T^{1/4}})$ and $\beta = \eta$,

$$\frac{1}{T} \sum_{t \in [T]} \mathbb{E} \|\mathbf{d}_t\|^2 = \mathcal{O}\left(\frac{1}{\sqrt{T}}\right). \quad (98)$$

679

□

680 D Further Experiments and Additional Results

681 In the following, we provide the detailed machine learning models for our experiments:

682 **1) MultiMNIST Datasets and Learning Tasks:** We test the convergence performance of our
 683 algorithms using the ‘‘MultiMNIST’’ dataset [45], which is a multi-task learning version of the
 684 MNIST dataset [49] from LIBSVM repository. Specifically, to convert the hand-written classification
 685 problem into a multi-task problem, we randomly chose 60000 images and divided them into M
 686 agents. Each agent has two tasks, where each task has $n = 60000/(2 * M)$ samples. Due to space
 687 limitations, we only present the convergence results for the case of non-i.i.d. data partition (i.e., data
 688 heterogeneity) and relegate the results of the i.i.d. data case to the appendix. For the non-i.i.d. data
 689 partition, we use the same data partition strategy as in [22], where each client can access data with
 690 at most two labels. In our experiments, a group of images is positioned in the top left corner, while
 691 another group of images is positioned in the bottom right. The two tasks are task ‘‘L’’ (to categorize
 692 the top-left digit) and task ‘‘R’’ (to classify the bottom-right digit). The overall problem is to classify
 693 the images of different tasks at different agents. All algorithms use the same randomly generated
 694 initial point. Here, we present experiments with $M = 10$ agents, where each agent has two tasks (i.e.,
 695 $\mathbf{A} \in \mathbb{R}^{M \times 2}$ is an all-one matrix). We set the local update rounds $K = 10$. Experiments with a larger
 696 number of agents ($M = 5, 10, 30$) are provided here. The learning rates are chosen as $\eta_L = 0.1$ and
 697 $\eta_t = 0.1, \forall t$.

698 **2): River Flow Dataset and Learning Tasks:** We further test our algorithms on FMOL problems of
 699 larger sizes. In this experiment, we use the River Flow dataset[46], which is for flow prediction flow
 700 at eight locations within the Mississippi River network. Thus, there are *eight* tasks in this problem.
 701 In this experiment, we set $\eta_L = 0.001$, $\eta_t = 0.1$, $M = 10$, and keep the batch size = 256 while
 702 comparing K , and keep $K = 30$ while comparing the batch size. To better visualize 8 different tasks,
 703 we illustrate the normalized loss in radar charts in Fig. 2(b). We again verify that utilizing a larger
 704 training batch size and conducting additional local steps K results in accelerated convergence.

705 **3): CelebA Dataset and Learning Tasks:** We utilize the CelebA dataset [47], consisting of 200K
 706 facial images annotated with 40 attributes. We approach each attribute as a binary classification task,
 707 resulting in a 40-way multi-task learning (MTL) problem. To create a shared representation function,
 708 we implement ResNet-18 [50] without the final layer, attaching a linear layer to each attribute for
 709 classification. In this experiment, we set $\eta_L = 0.0005$, $\eta_t = 0.1$, $M = 10$, and $K = 10$. Figure
 710 3 displays a radar chart depicting the loss value of each binary classification task. In Figure 3, we
 711 demonstrate the efficacy of our FMGDA and FSMGDA algorithms in both i.i.d. case and non-i.i.d.
 712 case.

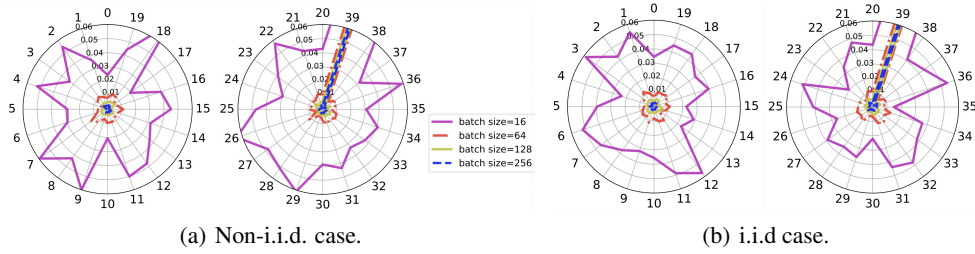


Figure 3: Experiments on CelebA dataset.

713 **Experiments on i.i.d. data:** First, we compare the convergence results with the same experimental
 714 settings in our Section. 5 but tested on the i.i.d data. As shown in Fig. 4, both FMGDA and FSMGDA
 715 successfully converged in i.i.d. data, and the algorithm with a larger training batch size and more
 716 local updates K may converge faster.

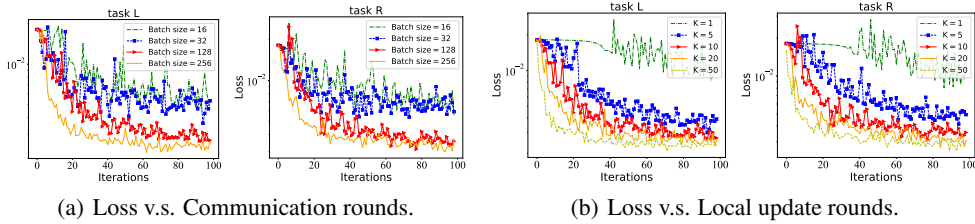


Figure 4: Experiments on i.i.d. data.

717 **Impact of the number of clients:** In this experiment, we choose the different number of clients from
 718 the discrete set $\{5, 10, 30\}$ and fix learning rates at 0.1 and local update rounds at 10. As shown in
 719 Fig. 5, a larger number of workers leads to faster convergence rates of our proposed algorithms both
 720 in i.i.d. case and non-i.i.d. case; this is mainly because more samples have been used while training
 721 while having more workers.

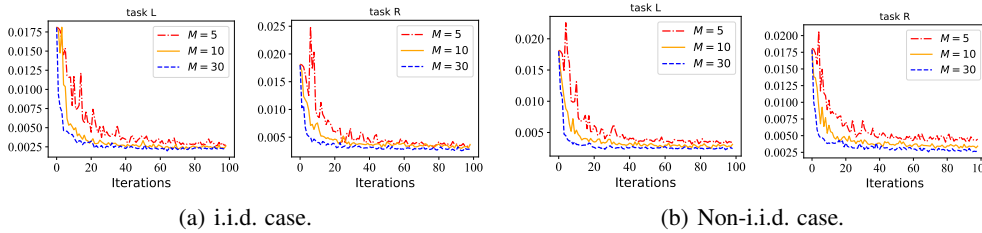


Figure 5: Loss value comparisons of algorithms on a different numbers of clients M .

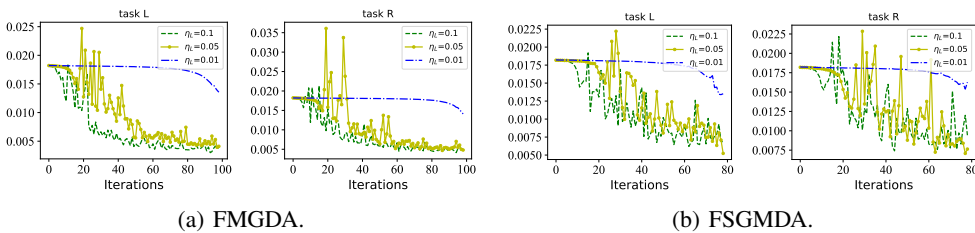


Figure 6: Comparisons of different step-sizes.

722 **Impact of the Step-size:** In this experiment, we choose the value of the learning rate η_L from
723 the discrete set $\{0.05, 0.01, 0.1\}$ and fix worker number at 5, local update rounds at 10. As shown
724 in Fig. 6, larger local step-sizes lead to faster convergence rates on both FMGDA algorithm and
725 FSMGDA algorithm.