

---

# Multi-block Min-max Bilevel Optimization with Applications in Multi-task Deep AUC Maximization

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 In this paper, we study multi-block min-max bilevel optimization problems, where  
2 the upper level is non-convex strongly-concave minimax objective and the lower  
3 level is a strongly convex objective, and there are multiple blocks of dual variables  
4 and lower level problems. Due to the intertwined multi-block min-max bilevel  
5 structure, the computational cost at each iteration could be prohibitively high,  
6 especially with a large number of blocks. To tackle this challenge, we present a  
7 single-loop randomized stochastic algorithm, which requires updates for only a  
8 constant number of blocks at each iteration. Under some mild assumptions on the  
9 problem, we establish its sample complexity of  $\mathcal{O}(1/\epsilon^4)$  for finding an  $\epsilon$ -stationary  
10 point. This matches the optimal complexity for solving stochastic nonconvex  
11 optimization under a general unbiased stochastic oracle model. Moreover, we  
12 provide two applications of the proposed method in multi-task deep AUC (area  
13 under ROC curve) maximization and multi-task deep partial AUC maximization.  
14 Experimental results validate our theory and demonstrate the effectiveness of our  
15 method on problems with hundreds of tasks.

## 16 1 Introduction

17 We consider multi-block min-max bilevel optimization problem of the following formulation

$$\min_{\mathbf{x} \in \mathbb{R}^{d_x}} \max_{\alpha \in \mathcal{A}^m} F(\mathbf{x}, \alpha) := \frac{1}{m} \sum_{i=1}^m \{f_i(\mathbf{x}, \alpha_i, \mathbf{y}_i(\mathbf{x})) := \mathbb{E}_{\xi \in \mathcal{P}_i} [f_i(\mathbf{x}, \alpha_i, \mathbf{y}_i(\mathbf{x}); \xi)]\} \quad (\text{upper}) \quad (1)$$

$$\text{s.t. } \mathbf{y}_i(\mathbf{x}) = \arg \min_{\mathbf{y}_i \in \mathbb{R}^{d_y}} g_i(\mathbf{x}, \mathbf{y}_i) := \mathbb{E}_{\zeta \in \mathcal{Q}_i} [g_i(\mathbf{x}, \mathbf{y}_i; \zeta)], \quad \text{for } i = 1, 2, \dots, m. \quad (\text{lower})$$

18 where  $f_i$  and  $g_i$  are smooth functions and  $\mathcal{A} \subset \mathbb{R}^{d_\alpha}$  is a convex set. In particular, in this paper  
19 we assume that for each  $i \in \{1, \dots, m\}$ ,  $f_i(\mathbf{x}, \alpha_i, \mathbf{y}_i)$  is strongly concave in the dual variable  $\alpha_i$   
20 but can be nonconvex in the primal variable  $\mathbf{x}$ , and  $g_i(\mathbf{x}, \mathbf{y}_i)$  is strongly convex in  $\mathbf{y}_i$ . The *upper*  
21 *problem*  $\min_{\mathbf{x} \in \mathbb{R}^{d_x}} \max_{\alpha \in \mathcal{A}^m} F(\mathbf{x}, \alpha)$  is a min-max optimization problem where the *lower problems*  
22  $\{\mathbf{y}_i(\mathbf{x}) = \arg \min_{\mathbf{y}_i \in \mathbb{R}^{d_y}} g_i(\mathbf{x}, \mathbf{y}_i)\}_{i=1}^m$  are involved as variables. For each *block*  $i$ , the *upper-level*  
23 *objective*  $f_i(\mathbf{x}, \alpha_i, \mathbf{y}_i)$  and the *lower-level objective*  $g_i(\mathbf{x}, \mathbf{y}_i)$  depend only on its corresponding block  
24 of variables  $\alpha$  and  $\mathbf{y}$ , i.e.  $\alpha_i$  and  $\mathbf{y}_i$ . This problem has important applications in machine learning,  
25 e.g., multi-task deep AUC maximization as presented in section 3.

26 Tackling problem (1) is challenging as it involves solving a min-max problem with coupled multiple  
27 minimization problems simultaneously. The naive way for solving it is to do multiple gradient ascents  
28 and descents for  $\alpha$  and  $\mathbf{y}$ , respectively, to ensure a good estimation of the gradient for updating  $\mathbf{x}$ .  
29 However, this approach has two major drawbacks. As the algorithm involves a double loop structure,  
30 it can be computationally expensive and give suboptimal theoretical complexity. On the other hand,  
31 the multi-block structure requires data sampling from distributions  $\mathcal{P}_i, \mathcal{Q}_i$  for all blocks, which may  
32 lead to an impractical demand for memory.

Submitted to 36th Conference on Neural Information Processing Systems (NeurIPS 2022). Do not distribute.

## 33 1.1 Related work

34 **Min-max Bilevel optimization.** To the best of our knowledge, the only existing work that provides a  
35 stochastic algorithm with provable convergence guarantee on min-max bilevel problems is [8]. They  
36 propose a single loop bi-time scale stochastic algorithm based on gradient descent ascent, and prove  
37 that it converges to an  $\epsilon$ -stationary point with an oracle complexity of  $\mathcal{O}(\epsilon^{-5})$ . Nevertheless, this  
38 convergence result is established for a special case where  $f(\mathbf{x}, \cdot, \mathbf{y})$  is a linear function.

39 **Stochastic Nonconvex Strongly Concave Min-max Problems.** The considered problem is also  
40 closely related to non-convex strongly concave min-max problems, which have been studied ex-  
41 tensively recently. To the best of our knowledge, [28] establishes the first result on non-smooth  
42 nonconvex concave min-max problems. They prove a convergence to nearly stationary point of the  
43 primal objective function with an oracle complexity in the order of  $\mathcal{O}(\epsilon^{-4})$  for non-convex strongly  
44 concave min-max problems with a certain special structure. The same order of oracle complexity is  
45 achieved in [34] without relying on any special structure. These two works use two-loop algorithms.  
46 There are some studies focusing on single-loop algorithms. [23] analyzes a single-loop stochastic  
47 gradient descent ascent (SGDA) method for smooth nonconvex strongly concave problem, which  
48 achieves  $\mathcal{O}(\epsilon^{-4})$  complexity but with a large mini-batch size. In [11], the same order of complexity  
49 is achieved without large mini-batch by employing the stochastic moving average estimator. Some  
50 recent works improve the sample complexity to  $\mathcal{O}(\epsilon^{-3})$  under the Lipschitz continuous oracle model  
51 for the stochastic gradient using less practical variance reduction techniques [15, 26, 29]. [39]  
52 establishes lower complexity bounds for non-convex strongly concave min-max problem under  
53 both general and finite-sum setting and proposes accelerated algorithms that nearly match the lower  
54 bounds.

55 **Stochastic Nonconvex Bilevel optimization.** The considered problem belongs to a general family  
56 of non-convex bilevel optimization problems. Non-asymptotic convergence results for nonconvex  
57 stochastic bilevel optimization (SBO) with a strongly convex lower problem has been established  
58 in several recent studies [4, 7, 11, 13, 18]. As the one who gives the first results for this problem,  
59 [7] proposes a double-loop algorithm with  $\mathcal{O}(\epsilon^{-6})$  oracle complexity for finding an  $\epsilon$ -stationary  
60 point of the objective function. [18] improves the complexity order to  $\mathcal{O}(\epsilon^{-4})$ , but suffers from  
61 a large mini-batch size. [13] proposes a single-loop algorithm with two time-scale updates that  
62 achieves an oracle complexity of  $\tilde{\mathcal{O}}(\epsilon^{-5})$ . Recently, [11] improves the oracle complexity to the  
63 state-of-the-art oracle complexity  $\tilde{\mathcal{O}}(\epsilon^{-4})$  by proposing a single-loop algorithm based on a moving-  
64 average estimator. [5] presents a new analysis for (double-loop) SGD-type updates showing that  
65 an improved sample complexity  $\mathcal{O}(\epsilon^{-4})$  can be achieved. There are studies that further improve  
66 the complexity to  $\mathcal{O}(\epsilon^{-3})$  by leveraging the Lipschitz continuous conditions of stochastic oracles  
67 [4, 9, 20]. [22] considers bilevel optimization under distributed setting and proposed algorithms  
68 achieving state-of-the-art complexities. However, none of these works tackle multi-block min-max  
69 bilevel optimization problems directly.

70 **Multi-block Bilevel Optimization.** There are some recent studies considering bilevel optimization  
71 with multi-block structure. [9] extends their single-block bilevel optimization algorithm to multi-block  
72 structure. Their algorithm requires two independently sampled block batches and for all sampled  
73 blocks, each variable needs update using variance reduction technique STORM [6]. For unsampled  
74 blocks, an update involving constant factor multiplication is also required. Under Lipschitz continuous  
75 conditions on stochastic oracles, the complexity is no worse than  $\mathcal{O}(m/\epsilon^3)$  with  $m$  blocks. A more  
76 recent work [27] considers top-K NDCG optimization, which is formulated as a compositional bilevel  
77 optimization with multi-block structure. Their method simplifies the updates by sampling only one  
78 block batch in each iteration and requires updates only for the sampled blocks. Their method achieves  
79 complexity of  $\mathcal{O}(m/\epsilon^4)$ . We use a similar approach as the latter work for estimating the hessian  
80 inverse in a block-wise manner. However, this paper differs from [27] in that we tackle a more general  
81 multi-block min-max bilevel problems without assuming a particular form of the objective.

## 82 1.2 Our Contributions

83 We present a simple single loop single timescale stochastic method with randomized block-sampling  
84 for solving a general form of multi-block min-max bilevel optimization problem under the nonconvex  
85 strongly concave (upper) strongly convex (lower) setting. It employs SGD for updating selected  $\mathbf{y}_i$  for  
86 their corresponding lower-level problems, employs SGA for updating the selected  $\alpha_i$ , and employs a  
87 momentum update for the primal variable  $\mathbf{x}$  based on the sampled  $\alpha_i, \mathbf{y}_i$ . Then we show, theoretically,  
88 that it converges to  $\epsilon$ -stationary point with complexity  $\mathcal{O}(\epsilon^{-4})$  under a general unbiased stochastic  
89 oracle model. This oracle complexity matches the optimal complexity order for solving stochastic

90 nonconvex optimization under a general unbiased stochastic oracle model [1]. Finally, we present  
 91 two applications of multi-block min-max bilevel optimization in deep AUC maximization: multi-task  
 92 deep AUC maximization and multi-task deep partial AUC maximization. Empirical results show the  
 93 effectiveness of the proposed method.

94

## 95 2 Preliminaries

96 **Notations.** Let  $\|\cdot\|$  denote the Euclidean norm of a vector or the spectral norm of a matrix.  
 97 For a twice differentiable function  $f : X \times Y \rightarrow \mathbb{R}$ ,  $\nabla_x f(x, y)$  (resp.  $\nabla_y f(x, y)$ ) denotes its  
 98 partial gradient taken w.r.t  $x$  (resp.  $y$ ), and  $\nabla_{xy}^2 f(x, y)$  (resp.  $\nabla_{yy}^2 f(x, y)$ ) denotes the Jacobian  
 99 of  $\nabla_x f(x, y)$  w.r.t  $x$  (resp.  $\nabla_y f(x, y)$  w.r.t  $y$ ). We let  $f(\cdot; \mathcal{B})$  represent the unbiased stochastic  
 100 oracle of  $f(\cdot)$  with a sample batch  $\mathcal{B}$  as the input. The unbiased stochastic oracle is said to have  
 101 bounded variance  $\sigma^2$  if  $\mathbb{E}[\|f(\cdot; \mathcal{B}) - f(\cdot)\|^2] \leq \sigma^2$ . A mapping  $f : X \rightarrow \mathbb{R}$  is  $C$ -Lipschitz  
 102 continuous if  $\|f(x) - f(x')\| \leq C\|x - x'\| \forall x, x' \in X$ . Function  $f$  is  $L$ -smooth if its gradient  
 103  $\nabla f(\cdot)$  is  $L$ -Lipschitz continuous. A function  $g : X \rightarrow \mathbb{R}$  is  $\lambda$ -strongly convex if  $\forall x, x' \in X$ ,  
 104  $g(x) \geq g(x') + \nabla g(x')^T(x - x') + \frac{\lambda}{2}\|x - x'\|^2$ . A function  $g : X \rightarrow \mathbb{R}$  is  $\lambda$ -strongly concave if  
 105  $-g(x)$  is  $\lambda$ -strongly convex. Let  $\Pi_{\mathcal{A}}$  denote a projection function onto a convex set  $\mathcal{A}$ . For notation  
 106 simplicity, we use  $\mathcal{S}$  to denote the set of all block indices, *i.e.*  $\mathcal{S} = \{1, \dots, m\}$ .

107

108 We state the definition of  $\epsilon$ -stationary point as following.

109 **Definition 2.1.** Consider a differentiable function  $F(\mathbf{x})$ , a point  $\mathbf{x}$  is called  $\epsilon$ -stationary if  $\|\nabla F(\mathbf{x})\| \leq$   
 110  $\epsilon$ . A stochastic algorithm is said to achieve an  $\epsilon$ -stationary point if  $\mathbb{E}[\|\nabla F(\bar{\mathbf{x}}_t)\|] \leq \epsilon$ , where  $\bar{\mathbf{x}}_t$  is  
 111 the algorithm output at the  $t$ -th iteration and the expectation is taken over the randomness of the  
 112 algorithm until the iteration  $t$ .

113 **Assumptions.** Before presenting our algorithm, we make the following well-behaving assumptions.

114 **Assumption 2.2.** For functions  $f_i$  and  $g_i$ , we assume that the following conditions hold for all  $i \in \mathcal{S}$

- 115 •  $f_i(\mathbf{x}, \alpha_i, \mathbf{y}_i)$  is  $\mu_f$ -strongly concave in terms of  $\alpha_i$ .  $g_i(x, \mathbf{y}_i)$  is  $\mu_g$ -strongly convex in terms  
 116 of  $\mathbf{y}_i$ .
- 117 •  $f_i, g_i$  are  $C_f, C_g$ -Lipschitz continuous and  $L_f, L_g$ -smooth respectively.
- 118 •  $\|\nabla_{xy}^2 g_i(\mathbf{x}, \mathbf{y}_i)\|^2 \leq C_{gxy}^2, \nabla_{yy}^2 g_i(\mathbf{x}, \mathbf{y}_i; \zeta) \succeq \mu_g I$ .
- 119 •  $\nabla_{xy}^2 g_i(\mathbf{x}, \mathbf{y}_i), \nabla_{yy}^2 g_i(\mathbf{x}, \mathbf{y}_i)$  are  $L_{gxy}, L_{gyy}$ -Lipschitz continuous respectively.

120 Moreover, the gradients of functions  $f_i$  and  $g_i$  can only be accessed through unbiased oracles with  
 121 bounded variance.

122 **Assumption 2.3.** The unbiased stochastic oracles  $\nabla_x f_i(\mathbf{x}, \alpha_i, \mathbf{y}_i; \mathcal{B}), \nabla_{\alpha} f_i(\mathbf{x}, \alpha_i, \mathbf{y}_i; \mathcal{B}),$   
 123  $\nabla_y f_i(\mathbf{x}, \alpha_i, \mathbf{y}_i; \mathcal{B}), \nabla_y g_i(\mathbf{x}, \mathbf{y}_i; \mathcal{B}), \nabla_{xy}^2 g_i(\mathbf{x}, \mathbf{y}_i; \mathcal{B}), \nabla_{yy}^2 g_i(\mathbf{x}, \mathbf{y}_i; \mathcal{B})$  have variances bounded by  $\frac{\sigma^2}{|\mathcal{B}|}$   
 124 for all  $i \in \mathcal{S}$ , where  $|\mathcal{B}|$  denotes the size of the sampled batch  $\mathcal{B}$ .

125 These assumptions are similar to those made in many existing works for SBO [4, 7, 13, 18].

126 **Moving average gradient estimator.** Algorithms based on moving average estimators have achieved  
 127 the state-of-the-art oracle complexity in both min-max and bilevel optimizations [10]. Here we  
 128 give a brief introduction to the moving average estimator. For solving a nonconvex minimization  
 129 problem  $\min_{\mathbf{x} \in \mathbb{R}^d} F(\mathbf{x})$  through an unbiased oracle  $\mathcal{O}_F(\mathbf{x})$ , *i.e.*  $\mathbb{E}[\mathcal{O}_F(\mathbf{x})] = \nabla F(\mathbf{x})$ , the stochastic  
 130 momentum method (stochastic heavy-ball method) that employs moving average updates is given by

$$\mathbf{v}_{t+1} = (1 - \beta)\mathbf{v}_t + \beta\mathcal{O}_F(\mathbf{x}_t), \quad \mathbf{x}_{t+1} = \mathbf{x}_t - \eta\mathbf{v}_{t+1},$$

131 where  $\beta$  and  $\eta$  are momentum parameter and learning rate, respectively. As a moving average of the  
 132 historical gradient estimator, the sequence of  $\mathbf{v}_{t+1}$  could achieve an effect of variance diminishing  
 133 across a long run ([30]).

### 134 2.1 The Proposed Algorithm

135 In this section, we propose a simple single loop stochastic algorithm 1 to solve the multi-block  
 136 min-max bilevel optimization problem. At the beginning of each iteration, we first sample a set of

---

**Algorithm 1** A Stochastic Algorithm for Multi-block Min-max Bilevel Optimization
 

---

**Require:**  $\alpha^0, \mathbf{y}^0, H^0, \mathbf{z}_0, \mathbf{x}_0$

- 1: **for**  $t = 0, 1, \dots, T$  **do**
  - 2:   Sample tasks  $I_t \in \mathcal{S}$ .
  - 3:   Sample data batch  $\mathcal{B}_i^t$  for each  $i \in I_t$ .
  - 4:    $\alpha_i^{t+1} = \begin{cases} \Pi_{\mathcal{A}}[\alpha_i^t + \eta_1 \nabla_{\alpha} f_i(\mathbf{x}_t, \alpha_i^t, \mathbf{y}_i^t; \mathcal{B}_i^t)] & \text{if } i \in I_t \\ \alpha_i^t & \text{o.w.} \end{cases}$
  - 5:    $\mathbf{y}_i^{t+1} = \begin{cases} \mathbf{y}_i^t - \eta_2 \nabla_{\mathbf{y}} g_i(\mathbf{x}_t, \mathbf{y}_i^t; \mathcal{B}_i^t) & \text{if } i \in I_t \\ \mathbf{y}_i^t & \text{o.w.} \end{cases}$
  - 6:   Update estimator  $H^{t+1}$  of  $[\nabla_{yy}^2 g_i(\mathbf{x}_t, \mathbf{y}_i^t)]^{-1}$  by (2)
  - 7:   Update gradient estimator  $\Delta^{t+1}$  of  $\nabla_x F(\mathbf{x}_t)$  by (3)
  - 8:    $\mathbf{z}_{t+1} = (1 - \beta_0)\mathbf{z}_t + \beta_0 \Delta^{t+1}$
  - 9:    $\mathbf{x}_{t+1} = \mathbf{x}_t - \eta_0 \mathbf{z}_{t+1}$
  - 10: **end for**
- 

137 blocks  $I_t$  and a data batch  $\mathcal{B}_i^t$  for each selected block  $i \in I_t$ . Then we update estimators of  $\alpha_i(\mathbf{x}_t)$   
 138 and  $\mathbf{y}_i(\mathbf{x}_t)$  for all selected blocks  $i \in I_t$  using one step of SGA and SGD. Then, we compute an  
 139 estimator of the hessian inverse  $H^{t+1}$  of the lower-level objective and compute a gradient estimator  
 140  $\Delta^{t+1}$  of the upper-level objective. Finally, we compute the moving average estimator  $\mathbf{z}_{t+1}$  of  $\nabla F(\mathbf{x}_t)$   
 141 and update  $\mathbf{x}_{t+1}$ . Note that the design of our algorithm on the min-max bilevel optimization part is  
 142 inspired by [10], and it is similar to their momentum-based algorithms PDSM (for min-max problem)  
 143 and SMB (for bilevel problem) in their paper. In fact, if we set the number of blocks to be one and  
 144 remove  $\mathbf{y}$  and the lower level problems, then Algorithm 1 is the same as PDSM. Similarly, if the  
 145 number of blocks is one and the dual variables in the upper level problem are removed, then the  
 146 proposed algorithm becomes similar as SMB except for the hessian inverse update with reasons  
 147 explain shortly. In other words, our proposed method is a generalized form of momentum-based  
 148 algorithm for min-max and bilevel optimization problems. Additionally, Algorithm 1 only updates  
 149  $O(1)$  blocks of dual variables  $\alpha_i$  and the variables  $\mathbf{y}_i$  of the lower-level problems. These make the  
 150 analysis of Algorithm 1 much more involved.

151 To further understand Algorithm 1, we first define the objective function

$$F(\mathbf{x}) := \frac{1}{m} \sum_{i \in \mathcal{S}} f_i(\mathbf{x}, \alpha_i(\mathbf{x}), \mathbf{y}_i(\mathbf{x})),$$

152 where  $\mathbf{y}_i(\mathbf{x}) = \arg \min_{\mathbf{y}_i} g_i(\mathbf{x}, \mathbf{y}_i)$  and  $\alpha_i(\mathbf{x}) := \arg \max_{\alpha_i} f_i(\mathbf{x}, \alpha_i, \mathbf{y}_i(\mathbf{x}))$ , so that the Problem (1)  
 153 can be rewritten as  $\min_{\mathbf{x}} F(\mathbf{x})$ . The updates for  $\alpha_i^{t+1}$ 's and  $\mathbf{y}_i^{t+1}$ 's are intuitive since the gradient  
 154 estimations of  $\nabla_{\alpha} f_i(\mathbf{x}_t, \alpha_i^t, \mathbf{y}_i^t)$  and  $\nabla_{\mathbf{y}} g_i(\mathbf{x}_t, \mathbf{y}_i^t)$  are directly available from the unbiased stochastic  
 155 oracles. However, since functions  $\mathbf{y}_i(\mathbf{x})$  and  $\alpha_i(\mathbf{x})$  are implicit, estimating the gradient  $\nabla F(\mathbf{x})$  is  
 156 difficult. In fact, one may apply the corollary of Theorem 1 in [3] to get:

$$\nabla F(\mathbf{x}) = \frac{1}{m} \sum_{i \in \mathcal{S}} (\nabla_x f_i(\mathbf{x}, \alpha_i(\mathbf{x}), \mathbf{y}_i(\mathbf{x})) + \nabla_{\mathbf{y}_i}(\mathbf{x}) \nabla_{\mathbf{y}} f_i(\mathbf{x}, \alpha_i(\mathbf{x}), \mathbf{y}_i(\mathbf{x}))).$$

157 A standard approach in bilevel optimization literature [7] for computing  $\nabla_{\mathbf{y}_i}(\mathbf{x})$  is to derive  $\nabla_{\mathbf{y}_i}(\mathbf{x}) =$   
 158  $-\nabla_{xy}^2 g_i(\mathbf{x}, \mathbf{y}_i(\mathbf{x})) [\nabla_{yy}^2 g_i(\mathbf{x}, \mathbf{y}_i(\mathbf{x}))]^{-1}$  from the optimality condition of  $\mathbf{y}_i(\mathbf{x})$ . Therefore, the gradient  
 159 we are looking for is given by

$$\nabla F(\mathbf{x}) = \frac{1}{m} \sum_{i \in \mathcal{S}} \nabla_x f_i(\mathbf{x}, \alpha_i(\mathbf{x}), \mathbf{y}_i(\mathbf{x})) - \nabla_{xy}^2 g_i(\mathbf{x}, \mathbf{y}_i(\mathbf{x})) [\nabla_{yy}^2 g_i(\mathbf{x}, \mathbf{y}_i(\mathbf{x}))]^{-1} \nabla_{\mathbf{y}} f_i(\mathbf{x}, \alpha_i(\mathbf{x}), \mathbf{y}_i(\mathbf{x})).$$

160 All components in this gradient can be easily obtained from unbiased stochastic oracles except for  
 161 the inverse of hessian  $[\nabla_{yy}^2 g_i(\mathbf{x}, \mathbf{y}_i(\mathbf{x}))]^{-1}$  for all blocks. This could be problematic in the sense of  
 162 theory and practical implementation. For practical implementation, we do not want to update the  
 163 hessian inverse estimators for all blocks, which is prohibitive when the number of blocks is large. A  
 164 common approach used in the literature of SBO is to use Neumann series [7]:

Choose  $q \in \{1, \dots, k_t\}$  randomly,

$$H_i^{t+1} = \frac{k_t}{C_{gyy}} \prod_{j=1}^q \left( I - \frac{1}{C_{gyy}} \nabla_{yy}^2 g_i(\mathbf{x}_t, \mathbf{y}_i^t; \xi_i^t) \right),$$

165 where  $k_t$  is the number of samples  $\{\xi_i^t\}_{i=1}^{k_t}$  for estimating the hessian inverse. This estimator is  
 166 a biased one and its error w.r.t to  $\nabla_{yy}^2 g_i(\mathbf{x}_t, \mathbf{y}_i^t)^{-1}$  is controlled by the number of samples  $k_t$  [7].  
 167 However, it is problematic to employ the above estimator for only the sampled blocks  $i \in I_t$  because  
 168 the error for those not sampled cannot be controlled.

169 To address this issue, we use a different approach for estimating the hessian inverse by only updating  
 170 the estimators for those sampled blocks [27]. The idea is to maintain a momentum term  $s_i^{t+1}$  for each  
 171 block that stores historical information on the hessian estimator  $\nabla_{yy}^2 g_i(\mathbf{x}_t, \mathbf{y}_i^t; \mathcal{B}_i^t)$ . And the hessian  
 172 inverse is approximated by directly computing the inverse of  $s_i^{t+1}$ , i.e.,

$$s_i^{t+1} = \begin{cases} (1 - \beta_1)s_i^t + \beta_1 \nabla_{yy}^2 g_i(\mathbf{x}_t, \mathbf{y}_i^t; \mathcal{B}_i^t) & \text{if } i \in I_t \\ s_i^t & \text{o.w.} \end{cases} \quad (2)$$

$$H_i^{t+1} = [s_i^{t+1}]^{-1}, \quad \forall i \in I_t.$$

173 In terms of theoretical analysis, we are not bounding the individual error  $H_i^{t+1} - \nabla_{yy}^2 g_i(\mathbf{x}_t, \mathbf{y}_i^t)^{-1}$  for  
 174 all blocks, but the cumulative error for all blocks across all iterations. This is exhibited in Lemma 2.6.  
 175 As the conclusion of the above discussion, the gradient estimator of  $\nabla F(\mathbf{w}_t)$  is given by

$$\Delta^{t+1} = \frac{1}{|I_t|} \sum_{i \in I_t} \{ \nabla_x f_i(\mathbf{x}_t, \alpha_i^t, \mathbf{y}_i^t; \mathcal{B}_i^t) - \nabla_{xy} g_i(\mathbf{x}_t, \mathbf{y}_i^t; \mathcal{B}_i^t) H_i^t \nabla_y f_i(\mathbf{x}_t, \alpha_i^t, \mathbf{y}_i^t; \mathcal{B}_i^t) \}. \quad (3)$$

176 We maintain a moving average estimator  $\mathbf{z}_{t+1}$  for  $\Delta^{t+1}$  and finally update  $\mathbf{x}_{t+1}$  using  $\mathbf{z}_{t+1}$ .

## 177 2.2 Convergence Analysis

178 In this section, we present a brief convergence analysis of Algorithm 1. The detailed theorems and  
 179 proofs are deferred to the appendix. The key point of this analysis is the gap between the true gradient  
 180  $\nabla F(\mathbf{x}_t)$  and its estimator  $\mathbf{z}^{t+1}$ . To this end, we define

$$\nabla F(\mathbf{x}_t, \alpha^t, \mathbf{y}^t) := \frac{1}{m} \sum_{i \in \mathcal{S}} \{ \nabla_x f_i(\mathbf{x}_t, \alpha_i^t, \mathbf{y}_i^t) - \nabla_{xy}^2 g_i(\mathbf{x}_t, \mathbf{y}_i^t) \mathbb{E}_t[H_i^t] \nabla_y f_i(\mathbf{x}_t, \alpha_i^t, \mathbf{y}_i^t) \}.$$

181 One may notice that the estimator  $\Delta^{t+1}$  is in fact approximating  $\nabla_x F(\mathbf{x}_t, \alpha^t, \mathbf{y}^t)$  instead of  $\nabla_x F(\mathbf{x}_t)$ .  
 182 We exploit the moving average formulation of  $\mathbf{z}^{t+1}$  and decompose the gap into two parts,  $\|\Delta^{t+1} -$   
 183  $\nabla_x F(\mathbf{x}_t, \alpha^t, \mathbf{y}^t)\|^2$  and  $\|\nabla_x F(\mathbf{x}_t, \alpha^t, \mathbf{y}^t) - \nabla_x F(\mathbf{x}_t)\|^2$ . These two gaps are determined by how  
 184 well  $\mathbf{y}^t$ ,  $\alpha^t$  and  $H_i^t$  approximate  $\mathbf{y}(\mathbf{x}_t)$ ,  $\alpha(\mathbf{x}_t)$  and  $[\nabla_{yy}^2 g_i(\mathbf{x}_t, \mathbf{y}_i^t)]^{-1}$ , respectively. In other words,  
 185 we aim to bound the following three errors,  $\|\mathbf{y}(\mathbf{x}_t) - \mathbf{y}^t\|^2 =: \delta_{y,t}$ ,  $\|\alpha(\mathbf{x}_t) - \alpha^t\|^2 =: \delta_{\alpha,t}$  and  
 186  $\|[\nabla_{yy}^2 g_i(\mathbf{x}_t, \mathbf{y}_i^t)]^{-1} - H_i^t\|^2$ .

187 We first bound the variance  $\mathbb{E}[\delta_{y,t}]$  by proving the following lemma.

188 **Lemma 2.4.** Consider the updates for  $\mathbf{y}^t$  in Algorithm 1, under Assumption 2.2 and 2.3, with  
 189  $\eta_2 \leq \min\{\frac{\mu_g}{L_g^2}, \frac{2m}{|I_t|\mu_g}\}$  we have

$$\sum_{t=0}^T \mathbb{E}[\delta_{y,t}] \leq \frac{2m}{|I_t|\eta_2\mu_g} \delta_{y,0} + \frac{4m\eta_2 T \sigma^2}{\mu_g |\mathcal{B}_i^t|} + \frac{8m^3 C_y^2 \eta_0^2}{|I_t|^2 \eta_2^2 \mu_g^2} \sum_{t=0}^{T-1} \mathbb{E}[\|\mathbf{z}_{t+1}\|^2]$$

190 One may also bound the second variance  $\mathbb{E}[\delta_{\alpha,t}]$  based on the previous variance  $\mathbb{E}[\delta_{y,t}]$  following a  
 191 similar strategy.

192 **Lemma 2.5.** Consider the updates for  $\alpha^t$  in Algorithm 1, under Assumption 2.2, 2.3, with  $\eta_1 \leq$   
 193  $\min\{\frac{\mu_f}{L_f^2}, \frac{1}{\mu_f}, \frac{4m}{\mu_f |I_t|}\}$ , we have

$$\sum_{t=0}^T \mathbb{E}[\delta_{\alpha,t}] \leq \frac{4m}{\eta_1 \mu_f |I_t|} \delta_{\alpha,0} + \frac{24L_f^2}{\mu_f^2} \sum_{t=0}^{T-1} \mathbb{E}[\delta_{y,t}] + \frac{8m\mu_f \eta_1 \sigma^2 T}{|\mathcal{B}_i^t|} + \frac{32m^3 C_\alpha^2 \eta_0^2}{\eta_1^2 \mu_f^2 |I_t|^2} \sum_{t=0}^{T-1} \mathbb{E}[\|\mathbf{z}_{t+1}\|^2].$$

194 Due to the lower bound assumption of  $\nabla_{yy}^2 g_i(\mathbf{x}, \mathbf{y}_i; \zeta)$  in Assumption 2.2, the error in Hessian  
 195 approximation can be bounded by bounding  $\delta_{g_{yy},t} := \sum_{i \in \mathcal{S}} \|s_i^t - \nabla_{yy}^2 g_i(\mathbf{x}_t, \mathbf{y}_i(\mathbf{x}_t))\|^2$ . We prove  
 196 the following lemma.

197 **Lemma 2.6.** Under Assumption 2.2,2.3, considering the momentum method (2) for the update of  
 198 hessian, with  $\beta_1 \leq 1$  we have

$$\sum_{t=0}^T \mathbb{E}[\delta_{g_{yy},t}] \leq \frac{4m\delta_{g_{yy},0}}{|I_t|\beta_1} + 32L_{g_{yy}}^2 \sum_{t=0}^{T-1} \mathbb{E}[\delta_{y,t}] + \frac{8m\beta_1 T \sigma^2}{|\mathcal{B}_i^t|} + \frac{32m^3 L_{g_{yy}}^2 (1 + C_y^2) \eta_0^2}{|I_t|^2 \beta_1^2} \sum_{t=0}^{T-1} \mathbb{E}[\|\mathbf{z}_{t+1}\|^2].$$

199 It then follows the convergence theorem for Algorithm 1.

200 **Theorem 2.7.** *Under Assumption 2.2, 2.3 and with a proper settings of parameters  $\eta_1, \eta_2, \beta_1 =$*   
 201  *$\mathcal{O}(|\mathcal{B}_i^t| \epsilon^2)$ ,  $\beta_0 = \mathcal{O}(\min\{|I_t|, |\mathcal{B}_i^t|\} \epsilon^2)$  and  $\eta_0 = \mathcal{O}\left(\min\left\{\min\{|I_t|, |\mathcal{B}_i^t|\} \epsilon^2, \frac{|\mathcal{B}_i^t| |I_t| \epsilon^2}{m}\right\}\right)$ , Algo-*  
 202 *rithm 1 ensures that after  $T = \mathcal{O}\left(\max\left\{\frac{m}{|I_t| |\mathcal{B}_i^t| \epsilon^4}, \frac{1}{\min\{|I_t|, |\mathcal{B}_i^t|\} \epsilon^4}\right\}\right)$  iterations we can find an*  
 203  *$\epsilon$ -stationary solution of  $F(\mathbf{x})$ , i.e.,  $\mathbb{E}[\|\nabla F(\mathbf{x}_\tau)\|^2] \leq \epsilon^2$  for a randomly selected  $\tau \in \{0, \dots, T\}$ .*

204 **Remark.** In Theorem 2.7, there is no condition on the sizes of data batch  $\mathcal{B}_i^t$  nor block batch  $I_t$   
 205 for the algorithm to converge. Hence, their sizes can be as small as one. The order of complexity  
 206 is  $\mathcal{O}(1/\epsilon^4)$ , which matches the optimal complexity for nonconvex optimization under a general  
 207 unbiased stochastic oracle model. In addition, there is parallel speed up by increasing batch sizes for  
 208 data samples and task samples due to the scaling in terms of  $|I_t|$  and  $|\mathcal{B}_i^t|$  in the iteration complexity.

### 209 3 Applications in Multi-task Deep (Partial) AUC Maximization

210 In this section, we present two applications of multi-block min-max bilevel optimization: multi-task  
 211 deep AUC maximization and deep partial AUC maximization. (partial) AUC is a performance measure  
 212 of classifiers for imbalanced data. Recent studies have shown great success of deep AUC maximization  
 213 in various domains (e.g., medical image classification and molecular property prediction) [24, 31, 38].  
 214 However, efficient algorithms for multi-task deep (partial) AUC maximization have not been well  
 215 developed. For multi-task deep AUC maximization, we solve an existing formulation by our algorithm.  
 216 For multi-task deep partial AUC maximization, we propose a new bilevel formulation and solve it by  
 217 our algorithm.

#### 218 3.1 Muti-task Deep AUC maximization

219 Following the previous work [24, 38], deep AUC maximization problem can be formulated as a  
 220 non-convex strongly concave min-max optimization problem  $\min_{\mathbf{w}, a, b} \max_{\alpha \in \mathcal{A}} L_{\text{AUC}}(\mathbf{w}, a, b, \alpha)$ .  
 221 However, training a deep neural network from scratch by optimizing AUC loss does not necessarily  
 222 lead to a good performance [37]. To address this issue, [37] proposed a compositional training strategy  
 223 for deep AUC maximization:

$$\min_{\mathbf{w}, a, b} \max_{\alpha \in \mathbb{R}_+} L_{\text{AUC}}(\mathbf{w} - \tilde{\eta} \nabla L_{\text{CE}}(\mathbf{w}), a, b, \alpha),$$

224 where  $L_{\text{CE}}$  denotes the cross-entropy loss. The outer objective remains to be the AUC loss, while  
 225 the inner objective is a gradient descent step of minimizing the traditional cross-entropy loss. This  
 226 method has shown superior performance on various datasets [37]. We extend this formulation to  
 227 multi-task problems and reformulate it into a multi-block min-max bilevel optimization:

$$\begin{aligned} & \min_{(\mathbf{w}_l, \mathbf{w}_h), \mathbf{a}, \mathbf{b}} \max_{\alpha \in \mathbb{R}_+^m} \sum_{i=1}^m L_{\text{AUC}}(\mathbf{u}^i(\mathbf{w}_l, \mathbf{w}_h), a^i, b^i, \alpha^i) \\ & \text{s.t. } \mathbf{u}^i(\mathbf{w}_l, \mathbf{w}_h) = \arg \min_{\mathbf{u}^i} \frac{1}{2} \left\| \mathbf{u}^i - ((\mathbf{w}_l, \mathbf{w}_h^i) - \tilde{\eta} \nabla L_{\text{CE}}(\mathbf{w}_l, \mathbf{w}_h^i)) \right\|^2, \end{aligned}$$

228 where  $\mathbf{w}_l$  denotes the weight for the encoder network that is shared for all tasks, and  $\mathbf{w}_h =$   
 229  $(\mathbf{w}_h^1, \dots, \mathbf{w}_h^m)$  denote the task-owned classification heads. The upper objective is strongly concave  
 230 in terms of dual variables  $\alpha_i$  and the lower level objective is strongly convex in terms of  $\mathbf{u}^i$ .  
 231 The hessian of the lower-level objective is the identity matrix. Hence, there is no need to track and  
 232 estimate the hessian matrix.

#### 233 3.2 Multi-task Deep Partial AUC Maximization

234 Some real-world applications (e.g., medical diagnosis [2]) cannot tolerate a model with a high False  
 235 Positive Rate (FPR) even though it has significant performance in AUC. Hence, a measure of interest  
 236 is one-way partial AUC (pAUC), which puts a restriction on the range of FPR (i.e.,  $\text{FPR} \in [\rho_l, \rho]$ ,  
 237 where  $0 \leq \rho_l \leq \rho \leq 1$ ). Below, we focus on the case  $\rho_l = 0$ . However, our method can be easily  
 238 extended for handling  $\rho_l > 0$ . Let  $\mathcal{D}_+, \mathcal{D}_-$  denote the set of positive and negative data for a binary  
 239 classification task, respectively. Let  $\mathcal{D}_-[K]$  denote the top-K negative examples according to their  
 240 prediction scores. Let  $n_+, n_-$  denote the number of positive and negative samples respectively. Then

---

**Algorithm 2** Min-Max Bilevel Optimization for pAUC Maximization (MMB-pAUC)
 

---

**Require:**  $\alpha^0, \lambda^0, H^0, \mathbf{z}^0, \mathbf{w}^0, \mathbf{a}^0, \mathbf{b}^0$ 

```

1: for  $t = 0, 1, \dots, T$  do
2:   Draw task batch  $I_t \subset \mathcal{S}$ .
3:   Draw data sample batch  $\mathcal{B}_k^t$  for each  $k \in I_t$ 
4:   For sampled tasks  $k \in I_t$ , update
5:      $\alpha_k^{t+1} \leftarrow \alpha_k^t + \eta_1 G_\alpha(\mathbf{w}^t, \alpha_k^t, \lambda_k^t; \mathcal{B}_k^t)$   $\diamond G_\alpha(\cdot)$  denotes a stochastic gradient w.r.t  $\alpha_k$ 
6:      $\lambda_k^{t+1} \leftarrow \lambda_k^t - \eta_2 \nabla_{\lambda} L_k(\lambda_k^t, \mathbf{w}^t; \mathcal{B}_k^t)$ 
7:      $H_k^{t+1} \leftarrow (1 - \beta_1)H_k^t + \beta_1 \nabla_{\lambda\lambda}^2 L_k(\lambda_k^t, \mathbf{w}^t; \mathcal{B}_k^t)$ 
8:   Compute loss  $G^{t+1}$  according to (25)  $\diamond G^{t+1}$  denotes an appropriate loss
9:   Update gradient estimator  $\Delta^{t+1} \leftarrow \text{autograd}(G^{t+1})$ 
10:   $\mathbf{z}^{t+1} \leftarrow (1 - \beta_0)\mathbf{z}^t + \beta_0 \Delta^{t+1}$ 
11:   $(\mathbf{w}^{t+1}, \mathbf{a}^{t+1}, \mathbf{b}^{t+1}) \leftarrow (\mathbf{w}^t, \mathbf{a}^t, \mathbf{b}^t) - \eta_0 \mathbf{z}^{t+1}$ 
12: end for

```

---

241 we have partial AUC optimization with a pairwise square loss formulated as following [35]:

$$\min_{\mathbf{w}} \frac{1}{n_+} \sum_{\mathbf{x}_i \in \mathcal{D}_+} \frac{1}{n_- \rho} \sum_{\mathbf{x}_j \in \mathcal{D}_-[K]} (h_{\mathbf{w}}(\mathbf{x}_j) - h_{\mathbf{w}}(\mathbf{x}_i) + c)^2,$$

242 where  $K = n_- \rho$ ,  $c$  is a constant and  $h_{\mathbf{w}}(\cdot)$  denotes the prediction score on a data. A key challenge for  
 243 solving the above problem is to deal with the non-differentiable top-K selector  $\mathbf{x}_j \in \mathcal{D}_-[K]$ , which  
 244 depends on the model parameters  $\mathbf{w}$ . This challenge has been recently tackled in [36, 40]. We focus  
 245 on the comparison with the first work as it is optimization oriented similar to ours and also has the  
 246 state-of-the-art performance. They formulate the problem into either a weakly convex minimization  
 247 or approximate it by a smooth objective in a compositional form. A caveat of their algorithms (named  
 248 SOPA, SOPA-s) is that they need to maintain and update  $n_+$  auxiliary variables with one for each  
 249 positive data. If we apply their algorithms for multi-task problems, one needs to maintain and update  
 250  $\sum_{i=1}^m n_+^i$  auxiliary variables, which could dramatically slow down the convergence.

251 To address this problem, we first transform it into a min-max optimization problem. Let  $a(\mathbf{w}) =$   
 252  $\frac{1}{n_+} \sum_{\mathbf{x}_i \in \mathcal{D}_+} h_{\mathbf{w}}(\mathbf{x}_i)$  and  $b(\mathbf{w}) = \frac{1}{n_- \rho} \sum_{\mathbf{x}_j \in \mathcal{D}_-[K]} h_{\mathbf{w}}(\mathbf{x}_j)$ . Then we can write the problem as (cf.  
 253 Appendix B for a derivation)

$$\min_{\mathbf{w}, a, b} \frac{1}{n_+} \sum_{\mathbf{x}_i \in \mathcal{D}_+} (h_{\mathbf{w}}(\mathbf{x}_i) - a)^2 + \frac{1}{n_- \rho} \sum_{\mathbf{x}_j \in \mathcal{D}_-} \mathbb{I}(\mathbf{x}_j \in \mathcal{D}_-[K]) (h_{\mathbf{w}}(\mathbf{x}_j) - b)^2 + (b(\mathbf{w}) - a(\mathbf{w}) + c)^2.$$

254 To tackle non-continuous non-differentiable indicator function  $\mathbb{I}(\cdot)$  we can replace it by a sigmoid  
 255 function. To tackle non-differentiability of the top-K selector  $\mathbf{x}_j \in \mathcal{D}_-[K]$ , we follow [27] and  
 256 formulate it as lower-level optimization problem, i.e.,  $\mathbf{x}_j \in \mathcal{D}_-[K]$  is equivalent to  $h_{\mathbf{w}}(\mathbf{x}_j) > \lambda(\mathbf{w})$ ,  
 257 where  $\lambda(\mathbf{w})$  represents the  $K + 1$ -th largest scores among all negative examples, which can be  
 258 approximated by a solution from a smooth strongly convex minimization problem as following:

$$\lambda(\mathbf{w}) = \arg \min_{\lambda \in \mathbb{R}} L(\lambda, \mathbf{w}) := \frac{K + \varepsilon}{n_-} \lambda + \frac{\tau_2}{2} \lambda^2 + \frac{1}{n_-} \sum_{\mathbf{x} \in \mathcal{D}_-} \tau_1 \ln(1 + \exp((h_{\mathbf{w}}(\mathbf{x}) - \lambda)/\tau_1)),$$

259 where  $\varepsilon, \tau_1, \tau_2$  are small constants. Based on above, the multi-task deep partial AUC minimization  
 260 problem can be formulated as a multi-block min-max bilevel optimization problem give by:

$$\min_{\mathbf{w}, \mathbf{a}, \mathbf{b}} \sum_{k=1}^m \left\{ \frac{1}{n_+^k} \sum_{\mathbf{x}_i \in \mathcal{D}_+^k} (h_{\mathbf{w}}(\mathbf{x}_i; k) - a_k)^2 + \frac{1}{n_- \rho} \sum_{\mathbf{x}_j \in \mathcal{D}_-^k} \phi(h_{\mathbf{w}}(\mathbf{x}_j; k) - \lambda_k(\mathbf{w})) (h_{\mathbf{w}}(\mathbf{x}_j; k) - b_k)^2 \right\}$$

$$+ \max_{\alpha \in \mathbb{R}^m} \left\{ 2\alpha_k \left( \frac{1}{n_- \rho} \sum_{\mathbf{x}_j \in \mathcal{D}_-^k} \phi(h_{\mathbf{w}}(\mathbf{x}_j; k) - \lambda_k(\mathbf{w})) h_{\mathbf{w}}(\mathbf{x}_j; k) - \frac{1}{n_+^k} \sum_{\mathbf{x}_i \in \mathcal{D}_+^k} h_{\mathbf{w}}(\mathbf{x}_i; k) + c \right) - \alpha_k^2 \right\}$$

$$\lambda_k(\mathbf{w}) = \arg \min_{\lambda \in \mathbb{R}} L_k(\lambda, \mathbf{w}) := \frac{K + \varepsilon}{n_-} \lambda + \frac{\tau_2}{2} \lambda^2 + \frac{1}{n_-} \sum_{\mathbf{x}_j \in \mathcal{D}_-^k} \tau_1 \ln(1 + \exp((h_{\mathbf{w}}(\mathbf{x}_j; k) - \lambda)/\tau_1)),$$

Table 1: The testing AUC scores on four datasets.

Method\DataSet	CIFAR100	CheXpert	CelebA	ogbg-molpcba
mAUC (baseline)	0.9044 (0.0015)	0.8084( 0.1455)	0.9062 (0.0042)	0.7793(0.0028)
mAUC-CT (ours)	<b>0.9272 (0.0014)</b>	<b>0.8198(0.1495)</b>	<b>0.9192 (0.0004)</b>	<b>0.8406(0.0044)</b>

261 where  $\mathcal{D}_{+/-}^k$  denote the positive/negative data set of the  $k$ -th task,  $h_{\mathbf{w}}(\mathbf{x}; k)$  denote the prediction  
 262 score for the  $k$ -th classifier, and  $\phi(s) = \frac{1}{1+\exp(-s)}$  is the sigmoid function. The upper-level objective  
 263 is strongly concave in terms of  $\alpha_k$  and the lower-level objective is strongly convex in terms of  $\lambda_k$ .

264 We develop a tailored algorithm based on Algorithm 1 for solving the above formulation of multi-task  
 265 pAUC maximization as shown in Algorithm 2. For the hessian update, the momentum update (2)  
 266 is efficient due to that the each lower-level problem is only one-dimensional. For simplicity of  
 267 implementation, we define a loss  $G^{t+1}$  (25) in the Appendix B, on which auto-differentiation can be  
 268 directly applied for computing a gradient estimator. We refer readers to Appendix B for a detailed  
 269 explanation and derivation of  $G^{t+1}$ .

## 270 4 Experiments

### 271 4.1 Multi-task Deep AUC Maximization with Compositional Training

272 **Data.** We use four datasets, namely CIFAR100, CheXpert, CelebA and ogbg-molpcba. CIFAR-100  
 273 [21] is an image dataset consisting of 60,000  $32 \times 32$  color images in 100 classes. Hence, there are 100  
 274 tasks for CIFAR100. We follow 45,000/5,000/10,000 split to construct training/validation/testing  
 275 datasets. CelebA [25] is a large-scale face attributes dataset with more than 200K celebrity images,  
 276 each with 40 attribute annotations (i.e., 40 tasks). We use the recommended training/validation/testing  
 277 split as 162,770/19,866/19,961. CheXpert [17] is a dataset that contains 224,316 chest radiographs  
 278 with 14 observations. Since the official testing dataset is not open to public, we take the official  
 279 validation set as the testing data, and take the last 1000 images in the training dataset for validation.  
 280 Due to the absence of positive samples for the observation *Fracture* in the testing dataset, we ignore  
 281 this label and only consider the rest 13 observations (i.e., 13 tasks). The last dataset ogbg-molpcba is  
 282 a molecular property prediction graph dataset [14]. It consists of 437,929 graphs with 128 labels (i.e.,  
 283 128 tasks). We follow scaffold splitting procedure as recommended in [32].

284 **Models.** We use ResNet18 [12] for CIFAR-100 and CelebA, and ImageNet pretrained DenseNet121  
 285 [16] for CheXpert. For ogbg-molpcba, we use Graph Isomorphism Network (GIN) [33].

286 **Setup.** We compare our method for optimizing the multi-task AUC maximization with compositional  
 287 training denoted by mAUC-CT (ours) with a baseline that directly optimizes multi-task min-max  
 288 AUC loss denoted by mAUC (baseline). We do not compare with other straightforward baselines  
 289 (e.g., optimizing the CE loss and the focal loss) since they have been shown to be inferior than AUC  
 290 maximization methods for imbalanced data in many previous works [38, 40]. For both methods, the  
 291 learning rates  $\eta_2, \eta_1, \eta_0$  are set to be the same and tuned in  $\{0.01, 0.03, 0.05, 0.07, 0.1\}$ . The learning  
 292 rates decay by a factor of 10 at the 4th and 30th epoch for CheXpert and CelebA, respectively. No  
 293 learning rate decay is applied for CIFAR-100 and ogbg-molpcba. The moving average parameter  
 294  $\beta_0$  and  $\tilde{\eta}$  in the lower level problem of mAUC-CT (ours) are tuned in  $\{0.1, 0.5, 0.9\}$ . Regarding the  
 295 task sampling, for datasets CIFAR-100 and ogbg-molpcba, 10 tasks are sampled to be updated in  
 296 each iteration, and for each sampled task, we independently sample a data batch of size 128. For the  
 297 other two datasets with fewer tasks, CheXpert and CelebA, we sample one task at each iteration. The  
 298 batch size for data samples is 32 for CheXpert, and 128 for CelebA. We run both methods the same  
 299 number of epochs which varies on different data, 2000 epochs for CIFAR100, 6 epochs for CheXpert,  
 300 40 epochs for CelebA and 100 epochs for ogbg-molpcba.

301 **Results.** In Table 1, we report the testing AUC score with the model selected according to the  
 302 best performance on validation datasets. Comparing with optimizing AUC loss directly, mAUC-CT  
 303 (ours) achieves better performance on all tested datasets. We show the results of an ablation study in  
 304 Figure 1, which verifies convergence has a parallel speed-up effect on both the batch sizes of data  
 305 samples and task samples. The algorithm converges faster as either the data or task sample batch size  
 306 increases. In Figure 2 (left two), we compare our method with the baseline in terms of convergence  
 307 speed on the training data of two datasets, which demonstrate that our method converges faster. More  
 308 results on training convergence are included in the appendix.

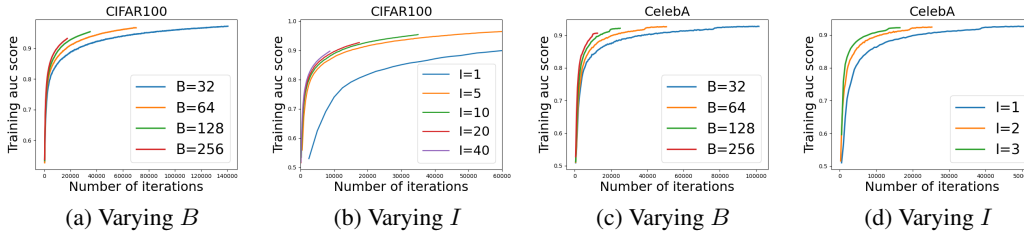


Figure 1: Convergence of our method vs data sample batch sizes  $B := |\mathcal{B}_t^t|$  and vs task sample batch size  $I := |I_t^t|$  for multi-task deep AUC maximization.

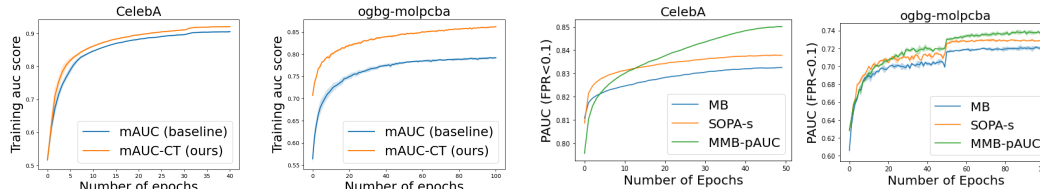


Figure 2: Comparison of Convergence on training data for multi-task deep AUC maximization (left two) and multi-task deep pAUC maximization (right two) on the CelebA and ogbg-molpcba datasets.

Table 2: The testing partial AUC scores on the four datasets.

Method\DataSet	CIFAR100	CelebA	CheXpert	ogbg-molpcba
CE	0.8895 (0.0009)	0.8024 (0.0026)	0.6606 (0.0159)	0.6576 (0.0010)
MB	0.9188 (0.0006)	0.8304 (0.0005)	0.6759 (0.0160)	0.7213 (0.0018)
SOPA-s	0.9251 (0.0003)	0.8336 (0.0001)	0.6682 (0.0156)	0.7290 (0.0019)
Ours	<b>0.9262 (0.0005)</b>	<b>0.8360 (0.0003)</b>	<b>0.6827 (0.0183)</b>	<b>0.7374 (0.0015)</b>

## 309 4.2 Multi-task Deep Partial AUC Maximization

310 **Setup.** For this task, we use the same datasets and the same networks for the image datasets and  
 311 graph datasets as in the previous subsection. For baselines, we compare with a naive mini-batch based  
 312 method (MB) for pAUC maximization [19], and a state-of-the-art pAUC maximization method SOPA-  
 313 s [40]. Following previous works [38, 40], for pAUC maximization methods we use a pretrained  
 314 encoder network by optimizing the CE loss as the initial encoder and learn the whole network by  
 315 maximizing pAUC. We also report the performance of optimizing the CE loss for a reference. For  
 316 all methods, the learning rate is tuned in  $\{0.0001, 0.0005, 0.001, 0.005, 0.01\}$ . The hyperparameters  
 317 selection of MMB-pAUC are:  $\eta_1$  and  $\eta_2 \in \{0.5, 0.1, 0.01\}$ ,  $\beta_1 \in \{0.99, 0.9, 0.5, 0.1, 0.01\}$  and  
 318  $\beta_0 \in \{0.9, 0.99\}$ . The momentum parameters in SOPA-s are tuned in the same range and their  $\lambda$   
 319 parameter in  $\{0.1, 1, 10\}$  as in [40]. The margin parameter in the surrogate loss (e.g.,  $c$ ) is set to be 1.  
 320 Regarding the task sampling, we sample one task at each iteration for ogbg-molpcba and CheXpert,  
 321 sample 10 tasks for CIFAR100, and sample 4 tasks for CelebA. The data sample batch size is 32 for  
 322 CheXpert, and 64 for others. For smaller datasets (CIFAR100 and ogbg-molpcba), we run 100 epochs  
 323 for each, and we decay the learning rate by a factor of 10 at the 50-th epoch. For larger datasets  
 324 (CelebA, CheXpert), we run 50 and 5 epochs respectively.

325 **Results.** The partial AUC scores with  $FPR \leq 0.1$  on the testing data of different methods are shown  
 326 in Table 2. From the results, we can see that our methods perform better than baseline methods with  
 327 a significant margin. In Figure 2 (right two), we compare our method with the baselines in terms of  
 328 convergence speed on the training data of two datasets, which demonstrate that our method converges  
 329 faster. More results on training convergence are included in the appendix.

## 330 5 Conclusion and Future Work

331 We have developed a simple single loop randomized stochastic algorithm for solving multi-block  
 332 min-max bilevel optimization problems. This algorithm requires updates for only constant number of  
 333 blocks in each iteration. We showed that it achieves an oracle complexity of  $\mathcal{O}(\epsilon^{-4})$ , which matches  
 334 the optimal complexity order for solving stochastic nonconvex optimization under a general unbiased  
 335 stochastic oracle model. One limitation of the proposed algorithm is that the computation of hessian  
 336 inverse estimator could be high for high dimensional problems. We plan to address this issue in the  
 337 future work. At the same time, we hope our work inspires others to find more novel applications of  
 338 our idea.

## References

- 339
- 340 [1] Yossi Arjevani, Yair Carmon, John C Duchi, Dylan J Foster, Nathan Srebro, and Blake Wood-  
341 worth. Lower bounds for non-convex stochastic optimization. *arXiv preprint arXiv:1912.02365*,  
342 2019.
- 343 [2] Stuart G. Baker and Paul F. Pinsky. A proposed design and analysis for comparing digital and  
344 analog mammography: Special receiver operating characteristic methods for cancer screening.  
345 *Journal of the American Statistical Association*, 96(454):421–428, 2001.
- 346 [3] Pierre Bernhard and Alain Rapaport. On a theorem of danskin with an application to a theorem  
347 of von neumann-sion. *Nonlinear Analysis: Theory, Methods Applications*, 24(8):1163–1181,  
348 1995.
- 349 [4] Tianyi Chen, Yuejiao Sun, and Wotao Yin. A single-timescale stochastic bilevel optimization  
350 method. *arXiv preprint arXiv:2102.04671*, 2021.
- 351 [5] Tianyi Chen, Yuejiao Sun, and Wotao Yin. Tighter analysis of alternating stochastic gradient  
352 method for stochastic nested problems, 2021.
- 353 [6] Ashok Cutkosky and Francesco Orabona. Momentum-based variance reduction in non-convex  
354 SGD. In *Advances in Neural Information Processing Systems 32 (NeurIPS)*, pages 15236–15245,  
355 2019.
- 356 [7] Saeed Ghadimi and Mengdi Wang. Approximation methods for bilevel programming. *arXiv*  
357 *preprint arXiv:1802.02246*, 2018.
- 358 [8] Alex Gu, Songtao Lu, Parikshit Ram, and Lily Weng. Min-max bilevel multi-objective opti-  
359 mization with applications in machine learning, 2022.
- 360 [9] Zhishuai Guo, Quanqi Hu, Lijun Zhang, and Tianbao Yang. Randomized stochastic  
361 variance-reduced methods for multi-task stochastic bilevel optimization. *arXiv preprint*  
362 *arXiv:2105.02266*, 2021.
- 363 [10] Zhishuai Guo, Yi Xu, Wotao Yin, Rong Jin, and Tianbao Yang. A novel convergence analysis  
364 for algorithms of the adam family and beyond, 2021.
- 365 [11] Zhishuai Guo, Yi Xu, Wotao Yin, Rong Jin, and Tianbao Yang. On stochastic moving-average  
366 estimators for non-convex optimization. *CoRR*, abs/2104.14840, 2021.
- 367 [12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
368 recognition. In *CVPR*, pages 770–778. IEEE Computer Society, 2016.
- 369 [13] Mingyi Hong, Hoi-To Wai, Zhaoran Wang, and Zhuoran Yang. A two-timescale framework  
370 for bilevel optimization: Complexity analysis and application to actor-critic. *arXiv preprint*  
371 *arXiv:2007.05170*, 2020.
- 372 [14] Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele  
373 Catasta, and Jure Leskovec. Open graph benchmark: Datasets for machine learning on graphs,  
374 2020.
- 375 [15] Feihu Huang, Shangqian Gao, Jian Pei, and Heng Huang. Accelerated zeroth-order momentum  
376 methods from mini to minimax optimization. *arXiv preprint arXiv:2008.08170*, 2020.
- 377 [16] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected  
378 convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern*  
379 *recognition*, pages 4700–4708, 2017.
- 380 [17] Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silvana Ciurea-Ilcus, Chris Chute,  
381 Henrik Marklund, Behzad Haghgoo, Robyn Ball, Katie Shpanskaya, et al. Chexpert: A large  
382 chest radiograph dataset with uncertainty labels and expert comparison. In *Proceedings of the*  
383 *AAAI Conference on Artificial Intelligence*, volume 33, pages 590–597, 2019.
- 384 [18] Kaiyi Ji, Junjie Yang, and Yingbin Liang. Provably faster algorithms for bilevel optimization  
385 and applications to meta-learning. *arXiv preprint arXiv:2010.07962*, 2020.

- 386 [19] Purushottam Kar, Harikrishna Narasimhan, and Prateek Jain. Online and stochastic gradi-  
387 ent methods for non-decomposable loss functions. In *Proceedings of the 27th International*  
388 *Conference on Neural Information Processing Systems - Volume 1, NIPS' 14*, page 694–702,  
389 Cambridge, MA, USA, 2014. MIT Press.
- 390 [20] Prashant Khanduri, Siliang Zeng, Mingyi Hong, Hoi-To Wai, Zhaoran Wang, and Zhuoran Yang.  
391 A near-optimal algorithm for stochastic bilevel optimization via double-momentum, 2021.
- 392 [21] Alex Krizhevsky. Learning multiple layers of features from tiny images. 2009.
- 393 [22] Junyi Li, Feihu Huang, and Heng Huang. Local stochastic bilevel optimization with momentum-  
394 based variance reduction, 2022.
- 395 [23] Tianyi Lin, Chi Jin, and Michael I Jordan. On gradient descent ascent for nonconvex-concave  
396 minimax problems. *arXiv preprint arXiv:1906.00331*, 2019.
- 397 [24] Mingrui Liu, Zhuoning Yuan, Yiming Ying, and Tianbao Yang. Stochastic AUC maximization  
398 with deep neural networks. In *8th International Conference on Learning Representations*  
399 *(ICLR)*, 2020.
- 400 [25] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the  
401 wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.
- 402 [26] Luo Luo, Haishan Ye, and Tong Zhang. Stochastic recursive gradient descent ascent for  
403 stochastic nonconvex-strongly-concave minimax problems. *CoRR*, abs/2001.03724, 2020.
- 404 [27] Zi-Hao Qiu, Quanqi Hu, Yongjian Zhong, Lijun Zhang, and Tianbao Yang. Large-scale  
405 stochastic optimization of ndcg surrogates for deep learning with provable convergence, 2022.
- 406 [28] Hassan Rafique, Mingrui Liu, Qihang Lin, and Tianbao Yang. Non-convex min-max opti-  
407 mization: Provable algorithms and applications in machine learning. *CoRR*, abs/1810.02060,  
408 2018.
- 409 [29] Quoc Tran-Dinh, Deyi Liu, and Lam M. Nguyen. Hybrid variance-reduced sgd algorithms  
410 for minimax problems with nonconvex-linear function. In *Advances in Neural Information*  
411 *Processing Systems 33 (NeurIPS)*, 2020.
- 412 [30] Mengdi Wang, Ethan X Fang, and Han Liu. Stochastic compositional gradient descent: algo-  
413 rithms for minimizing compositions of expected-value functions. *Mathematical Programming*,  
414 161(1-2):419–449, 2017.
- 415 [31] Zhengyang Wang, Meng Liu, Youzhi Luo, Zhao Xu, Yaochen Xie, Limei Wang, Lei Cai, Qi Qi,  
416 Zhuoning Yuan, Tianbao Yang, and Shuiwang Ji. Advanced graph and sequence neural networks  
417 for molecular property prediction and drug discovery, 2020.
- 418 [32] Zhenqin Wu, Bharath Ramsundar, Evan N. Feinberg, Joseph Gomes, Caleb Geniesse, Aneesh S.  
419 Pappu, Karl Leswing, and Vijay Pande. Moleculenet: A benchmark for molecular machine  
420 learning, 2017.
- 421 [33] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural  
422 networks?, 2018.
- 423 [34] Yan Yan, Yi Xu, Qihang Lin, Wei Liu, and Tianbao Yang. Optimal epoch stochastic gradient de-  
424 scent ascent methods for min-max optimization. In *Advances in Neural Information Processing*  
425 *Systems 33 (NeurIPS)*, 2020.
- 426 [35] Tianbao Yang and Yiming Ying. AUC maximization in the era of big data and AI: A survey.  
427 *CoRR*, abs/2203.15046, 2022.
- 428 [36] Zhiyong Yang, Qianqian Xu, Shilong Bao, Yuan He, Xiaochun Cao, and Qingming Huang.  
429 When all we need is a piece of the pie: A generic framework for optimizing two-way partial auc.  
430 In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference*  
431 *on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages  
432 11820–11829. PMLR, 18–24 Jul 2021.

- 433 [37] Zhuoning Yuan, Zhishuai Guo, Nitesh Chawla, and Tianbao Yang. Compositional training for  
434 end-to-end deep AUC maximization. In *International Conference on Learning Representations*,  
435 2022.
- 436 [38] Zhuoning Yuan, Yan Yan, Milan Sonka, and Tianbao Yang. Robust deep auc maximization:  
437 A new surrogate loss and empirical studies on medical image classification. *arXiv preprint*  
438 *arXiv:2012.03173*, 2020.
- 439 [39] Siqu Zhang, Junchi Yang, Cristóbal Guzmán, Negar Kiyavash, and Niao He. The complexity  
440 of nonconvex-strongly-concave minimax optimization. In Cassio de Campos and Marloes H.  
441 Maathuis, editors, *Proceedings of the Thirty-Seventh Conference on Uncertainty in Artificial*  
442 *Intelligence*, volume 161 of *Proceedings of Machine Learning Research*, pages 482–492. PMLR,  
443 27–30 Jul 2021.
- 444 [40] Dixian Zhu, Gang Li, Bokun Wang, Xiaodong Wu, and Tianbao Yang. When AUC meets  
445 DRO: optimizing partial AUC for deep learning with non-convex convergence guarantee. *CoRR*,  
446 abs/2203.00176, 2022.

447 **Checklist**

- 448 1. For all authors...
- 449 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's  
450 contributions and scope? [Yes]
- 451 (b) Did you describe the limitations of your work? [Yes]
- 452 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 453 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
454 them? [Yes]
- 455 2. If you are including theoretical results...
- 456 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 457 (b) Did you include complete proofs of all theoretical results? [Yes]
- 458 3. If you ran experiments...
- 459 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
460 mental results (either in the supplemental material or as a URL)? [Yes]
- 461 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
462 were chosen)? [Yes]
- 463 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
464 ments multiple times)? [Yes]
- 465 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
466 of GPUs, internal cluster, or cloud provider)? [N/A]
- 467 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 468 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 469 (b) Did you mention the license of the assets? [Yes]
- 470 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 471 (d) Did you discuss whether and how consent was obtained from people whose data you're  
472 using/curating? [N/A]
- 473 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
474 information or offensive content? [N/A]
- 475 5. If you used crowdsourcing or conducted research with human subjects...
- 476 (a) Did you include the full text of instructions given to participants and screenshots, if  
477 applicable? [N/A]
- 478 (b) Did you describe any potential participant risks, with links to Institutional Review  
479 Board (IRB) approvals, if applicable? [N/A]
- 480 (c) Did you include the estimated hourly wage paid to participants and the total amount  
481 spent on participant compensation? [N/A]