

# GRADIENT-BASED CONSTRAINED MULTI-OBJECTIVE OPTIMIZATION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

There is more and more attention on constrained multi-objective optimization (CMOO) problems, however, most of them are based on gradient-free methods. This paper proposes a constraint gradient-based algorithm for multi-objective optimization (MOO) problems based on multi-gradient descent algorithms. We first establish a framework for the CMOO problem. Then, we provide a Moreau envelope-based Lagrange Multiplier (MLM-CMOO) algorithm to solve the formulated CMOO problem, and the convergence analysis shows that the proposed algorithm convergence to Pareto stationary solutions with a rate of  $\mathcal{O}(\frac{1}{\sqrt{T}})$ . Finally, the MLM-CMOO algorithm is tested on several CMOO problems and has shown superior results compared to some chosen state-of-the-art designs.

## 1 INTRODUCTION

Multi-objective optimization (MOO) is widely used in many real-world application scenarios, such as, in online advertising, the models need to maximize both the Click-Through Rate and the Post-Click Conversion Rate. In MOO, one attempts to simultaneously optimize several, potentially conflicting functions. MOO has wide applications in all industry sectors where decision-making is involved due to the natural appearance of conflicting objectives or criteria. Applications span across applied engineering, operations management, finance, economics, and social sciences, agriculture, green logistics, and health systems. When the individual objectives are conflicting, no single solution exists that optimizes all of them simultaneously. In such cases, the goal of MOO is then to find Pareto optimal solutions (also known as efficient points), roughly speaking points for which no other combination of variables leads to a simultaneous improvement in all objectives. The determination of the set of Pareto optimal solutions helps decision makers to determine the best trade-offs among the several competing criteria.

MOO research can be divided into 2 categories, which are gradient-free and gradient-based methods. For the gradient-free method, people focus on evolutionary and Bayesian MOO algorithms, which are suitable for small-scale problems but less practical for high-dimensional MOO models and can not provide a convergence guarantee. On the contrary, the gradient-based method can provide a convergence guarantee in strongly convex, convex, and non-convex functions for MOO problems with different assumptions. The CMOO problem in the gradient-free method is well-developed. However, there is no gradient-based method for the CMOO problem.

Compared to conventional single-objective optimization, one key difference in MOO is the coupling and potential conflicts between different objective functions. As a result, there may not exist a common solution that minimizes all objective functions. Rather, the goal in MOO is to find a *Pareto stationary solution* that is not improvable for all objectives without sacrificing some objectives. The gradient-based method MOO has 2 lines, single-objective transformation, and the conflict gradients alleviating method, where the latter has garnered more attention in recent years due to their better performances. The single-objective transformation is the first step for the gradient-based method MOO method. It first transfers a MOO problem into a single-objective optimization (SOO) problem with a given fixed coefficient. With the sufficient algorithm in SOO, it is easy to solve. However, this transformation can not give a stable guarantee for the convergence rate as it may give the farthest Pareto front for the given coefficient. Then, the conflict gradients alleviating method is proposed to resolve the conflicting gradients in MOO. However, none of them pay attention to the gradient-based CMOO problem.

Although many gradient-free methods can provide solutions to CMOO problems, very few of them can provide the convergence guarantee. In addition, gradient-free methods are more suitable for small-scale and low-dimension MOO problems, which limits the application of gradient-free algorithms. Thus, we provide a Moreau envelope-based Lagrange Multiplier (MLM-CMOO) algorithm for the CMOO problem via the gradient method. Our contributions are summarized as follows.

- We propose MLM-CMOO, which solves the CMOO problem using the gradient-based method. MLM-CMOO first divides the CMOO problem into 2 parts, maximizes the minimum decrease across the losses, and makes the decrease obey the constraints. To maximize the minimum decrease across the losses, we use a similar method in MGDA (Sener & Koltun, 2018). To limit the decrease, we use the Moreau envelope-based proximal gradient method.
- We provide convergence analyses for MLM-CMOO with convex multi-objectives and convex multi-constraints. The convergence rate of MLM-CMOO is  $\mathcal{O}(\frac{1}{\sqrt{T}})$ .
- We conduct numerical experiments to verify the effectiveness of MLM-CMOO. The experimental results demonstrate the efficiency of the MLM-CMOO.

The remainder of this paper is organized as follows. Section 2 reviews related work. In Section 3, we present the system model and algorithm design of MLM-CMOO. In Section 4, we provide the convergence analysis of the MLM-CMOO algorithm. Numerical results and conclusions are provided in Section 5 and Section 6, respectively.

## 2 RELATED WORK

**MOO.** MOO algorithms can be grouped into two main categories. The first line of work is gradient-free methods (e.g., evolutionary MOO algorithms and Bayesian MOO algorithms (Lin et al., 2022; Zhang & Li, 2007; Laumanns & Ocenasek, 2002; Deb et al., 2002; Belakaria et al., 2020)). These methods are more suitable for small-scale problems but less practical for high-dimensional MOO models (e.g., deep neural networks). (Do et al., 2023; Zheng et al., 2022) provided the convergence analysis for gradient-free methods to solve MOO problems. The second line of work focuses on gradient-based approaches (Liu & Vicente, 2024; J. Fliege & Vicente, 2019; Momma et al., 2022; Peitz & Dellnitz, 2018; Désidéri, 2012), which are more practical for high-dimensional MOO problems. However, while having received increasing attention from the community in recent years, the Pareto-stationary convergence analysis of these gradient-based MOO methods attracts much more attention.

Various works explored the convergence rates under different assumptions in strongly convex, convex, and non-convex functions for MOO problems. Using full gradient, MGD (J. Fliege & Vicente, 2019) could achieve tight convergence rates in strongly-convex and non-convex cases, i.e., linear rate  $\mathcal{O}(r^T)$ ,  $r \in (0, 1)$  and sub-linear rate  $\mathcal{O}(1/T)$ . However, it requires a linear search of the learning rate in the algorithm and sequence convergence ( $\{x_t\}$  converges to  $x^*$ ). The linear search of learning rate is a classic technique but does not fit in gradient-based algorithms in deep learning. Moreover, the sequence convergence assumption is too strong. If using a stochastic gradient, SMGD (Liu & Vicente, 2024) methods make a further complicated case. The stochastic gradient noise would complicate the analysis. an  $\mathcal{O}(1/T)$  rate analysis for SMGD was provided in (Liu & Vicente, 2024) based on rather strong assumptions on a first-moment bound and Lipschitz continuity of common descent direction. On the other hand, (Liu & Vicente, 2024) and (Zhou et al., 2022) showed that the common descent direction provided by the SMGD method is likely to be a biased estimation, which may cause divergence issues. Recently, by utilizing momentum, oCo (Fernando et al., 2024) and CR-MOGM (Zhou et al., 2022) were proposed with corresponding convergence guarantees. (Xiao et al., 2023) utilized a direction-oriented approach by a preference direction. (Yang et al., 2023) proposed a federated MOO algorithm with GD and SGD matching previous centralized MOO algorithms.

**CMOO.** Most existing CMOO research focuses on gradient-free methods. SaE-CMO (Song et al., 2024) proposed a cooperative evolutionary algorithm with a dual-population approach to enhance search progress. PAC-MOO was proposed in (Ahmadianshalchi et al., 2024) based on Bayesian optimization. (Zhang et al., 2024) introduced a dynamic assistant population to search direction for CMOO. (Yang et al., 2024) proposed a feasibility tracking strategy to explore all feasible regions

for CMOO. (Belaiche et al., 2023) proposed PCMOEA/D-DMA based on a multi-population mechanism and implemented under a synchronous master-slave parallel model to select the best Pareto front based on an elitism mechanism. (Li et al., 2023) proposed a surrogate-ensemble-assisted co-evolutionary algorithm to improve the search efficiency.

**Constraint Handling Techniques.** Penalty function methods, decoders, special operators, and separation techniques are a simple taxonomy of the constraint handling methods in nature-inspired optimization algorithms. There are several types of penalty functions used with evolutionary algorithms (EAs), the most important ones include (Kramer, 2010) Death penalty, Dynamic penalty, Static penalty, Adaptive penalty, and Stochastic ranking. As an example of decoders, (Koziel & Michalewicz, 1998) proposed a homomorphous mapping (HM) method between an  $n$ -dimensional cube and feasible space. The feasible region can be mapped onto a sample space where a population-based algorithm could run a comparative performance (Koziel & Michalewicz, 1998; Kim & Husbands, 1998a;b; Koziel & Michalewicz, 1999). However, this method requires high computational costs. A special operator is used to preserve the feasibility of a solution or move within a special region (Michalewicz, 1996; Schoenauer & Michalewicz, 1996; 1997). Nevertheless, this method is hindered by the initialization of feasible solutions in the initial population, which is challenging with highly constrained optimization problems. Unlike the penalty function technique, another approach separates the values of objective functions and constraints in the nature-inspired algorithms (NIAs) (Powell & Skolnick, 1993), which is known as the separation of objective function and constraints. The authors of (Hinterding & Michalewicz, 1998) initially proposed dividing the search space into two phases. In the first phase, feasible solutions are found, and optimizing the objective function is considered in the second phase. Representative methods of this type of CHT are the Constraint dominance principle (CDP), Epsilon CHT, and Feasibility rules.

### 3 SYSTEM AND PROBLEM FORMULATION

This section introduces the basic background knowledge of MOO, typical algorithms, and their convergence analysis. Then, we proposed the objective of CMOO and the algorithm to solve the formulated problem.

#### 3.1 MULTI-OBJECTIVE OPTIMIZATION

Multi-objective optimization (MOO) is concerned with solving the problems of optimizing multi-objective functions simultaneously, which can be formulated as

$$\min_x F(x) = (f_1(x), f_2(x), \dots, f_n(x))^T, \quad (1)$$

where  $f_i$  are real-valued functions, and  $\mathcal{N}$  represents the set of the total  $n$  objectives ( $n > 2$ ). The MOO problem is smooth if all objective functions  $f_i$  are continuously differentiable. Different from single objective optimization where 2 solutions  $x, y$  can be ordered by  $f(x) < f(y)$  or  $f(x) \geq f(y)$ . MOO could have two parameter vectors where one performs better for task  $i$  and the other performs better for task  $j$ , where  $i \neq j$ . Therefore, Pareto optimality is defined to deal with such an incomparable case.

**Definition 3.1 (Pareto optimality).** For any two solutions  $x_1, x_2 \in \mathcal{X}$ , we say that  $x_1$  dominates  $x_2$ , denoted as  $x_1 \prec x_2$ , if  $f_i(x_1) \leq f_i(x_2)$  for all  $i$ , and there exists one  $i$  such that  $f_i(x_1) < f_i(x_2)$ ; otherwise, we say that  $x_1$  does not dominate  $x_2$ , denoted as  $x_1 \not\prec x_2$ . A solution  $x^* \in \mathcal{X}$  is called Pareto optimal if it is not dominated by any other solution in  $\mathcal{X}$ .

Note that a set of Pareto optimal solutions is called a Pareto set. The goal of MOO is to find a Pareto optimal solution, which must be Pareto critical (Custódio et al., 2011).

**Definition 3.2 (Pareto criticality).** A solution  $x^* \in \mathcal{X}$  is called Pareto critical if there is no common descent direction  $\mathbf{d}$  such that  $\nabla f_i(x^*)^\top \mathbf{d} < 0$ ,  $i \in \mathcal{M}$  for all objectives.

This definition indicates that if  $x$  is not Pareto critical, such direction  $\mathbf{d}$  will be a local descent direction for  $\mathbf{F}$  at point  $x$ . Optimizing through  $\mathbf{d}$  in the local neighborhood of  $x$  can get a better solution that dominates  $x$  (J. Fliege & Vicente, 2019). Since Pareto criticality reflects the local

property compared with Pareto optimality, it is often used as the local minimal condition for MOO with non-convex objectives (Fliege & Svaiter, 2000). We then present sufficient conditions for determining Pareto criticality/optimality, which appear as metrics to study the convergence for the MOO algorithm (Fliege & Svaiter, 2000; H. Tanabe & Yamashita, 2023).

Similar to single-objective optimization, MOO can be solved by running iteratively with gradient-based algorithms. For example, in MGDA (Sener & Koltun, 2018), it directly optimizes towards the Pareto criticality in Definition 3.2. Specifically, in each iteration, MGDA aims to find a direction  $\mathbf{d}$  to maximize the minimum decrease across the losses by solving the following subproblem,

$$\max_{\mathbf{d}} \min_i (f_i(x) - f_i(x + \eta \mathbf{d})) \approx \eta \max_{\mathbf{d}} \min_i \nabla f_i(x)^\top \mathbf{d}.$$

By regularizing the norm of  $\mathbf{d}$  on the right side, it computes the direction by

$$\mathbf{d} = \arg \min_{\mathbf{d}} \left\{ \max_i \lambda_i \nabla f_i(x) + \frac{1}{2} \|\mathbf{d}\|^2 \right\}.$$

This sub-problem can be rewritten equivalently as the following differentiable quadratic optimization

$$\mathbf{d}, \mu = \arg \min_{\mathbf{d}, \mu} \left( \frac{1}{2} \|\mathbf{d}\|^2 + \mu \right), \text{ s.t. } \lambda_i \nabla f_i(x) \leq \mu.$$

If  $\mu < 0$ , then  $\nabla f_i(x)^\top \mathbf{d} < 0$ , which means  $x$  is not Pareto critical from Definition 3.2, and  $\mathbf{d}$  is the direction to descent all the objectives simultaneously (Fliege & Svaiter, 2000; J. Fliege & Vicente, 2019). To simplify the optimization, such a primal problem has a dual objective as a min-norm oracle

$$\lambda = \arg \min_{\lambda} \left\| \sum_{i=1}^m \lambda_i \nabla f_i(x) \right\|.$$

The direction is then calculated by  $\mathbf{d} = -\sum_{i=1}^m \lambda_i \nabla f_i(x)$ .

**Convergence analysis.** MGDA has been shown to converge to an arbitrary Pareto critical/optimal point with the same rate as single-objective optimization (J. Fliege & Vicente, 2019). A similar result has been proved with PCGrad (Yu et al., 2020). CAGrad has been shown to converge to the minimizer or stationary point of the averaging loss  $\frac{1}{n} \sum_{i=1}^n f_i(x)$  when  $c \in [0, 1)$ , or an arbitrary Pareto critical/optimal point when  $c \geq 1$  (Liu et al., 2021).

### 3.2 CONSTRAINED MULTI-OBJECTIVE OPTIMIZATION

Followed by previous research, this paper considers a MOO problem, where each objective has its constraints. It is formulated as

$$\begin{aligned} \min_x F(x) &= (f_1(x), f_2(x), \dots, f_n(x))^\top, \\ \text{s.t. } g_i(x) &\leq 0, \forall i \in \mathcal{N}. \end{aligned}$$

Similar to MGDA, we aim to find a direction  $\mathbf{d}$  to maximize the minimum decrease across the losses. Thus, our problem can be reformulated as

$$\begin{aligned} \min_x F(x, \lambda) &:= \sum_{i=1}^m \lambda_i f_i(x), \\ \text{s.t. } \lambda &= \arg \min_{\lambda'} \left\| \sum_{i=1}^m \lambda'_i \nabla f_i(x) \right\|, \\ g_i(x) &\leq 0, \forall i \in \mathcal{N}. \end{aligned} \tag{2}$$

To simplify the expression, we note  $H(x, \lambda) := \left\| \sum_{i=1}^m \lambda_i \nabla f_i(x) \right\|$ , due to the non-smooth of  $H(x, \lambda)$ , we use  $\nabla H(x, \lambda)$  to express the proximal gradient of  $H(x, \lambda)$ , where  $\nabla_x H(x, \lambda) := \arg \min_u \{H(u, \lambda) + \frac{1}{2} \|u - x\|^2\}$  and  $\nabla_\lambda H(x, \lambda) := \arg \min_v \{H(x, v) + \frac{1}{2} \|v - \lambda\|^2\}$ .

To solve the above problem 2, we first find the suitable coefficient ( $\lambda$ ) of the problem 2, where we get a subproblem.

$$\begin{aligned} & \min_{\lambda} H(x, \lambda), \\ & \text{s.t. } g_i(x) \leq 0, \forall i \in \mathcal{N}. \end{aligned} \quad (3)$$

Then, we transfer problem 3 into an unconstrained optimization problem via the Lagrange Multiplier method, which is expressed as

$$L = H(x, \lambda) + \sum_{i=1}^m \mu_i g_i(x). \quad (4)$$

Due to the absolute value of  $\sum_{i=1}^m \lambda_i \nabla f_i(x)$ , the Lagrange function ( $L$ ) may not be smooth, thus we use a Moreau envelope-based Lagrange Multiplier function to solve above problems, which can be expressed as

$$L_s(x, \lambda, z) := \min_{\theta} \max_{\mu} \left\{ H(x, \theta) + \sum_{i=1}^N \mu_i g_i(x) + \frac{1}{2\gamma_1} \sum_{i=1}^N \|\theta_i - \lambda_i\|^2 - \frac{1}{2\gamma_2} \sum_{i=1}^N \|z_i - \mu_i\|^2 \right\},$$

where  $\gamma_1, \gamma_2$  are the proximal parameter and  $\gamma_1 \geq 0, \gamma_2 \geq 0$ .

Employing the function of  $L_s$ , we reformulated the problem 2 as

$$\begin{aligned} & \min_x F(x, \lambda), \\ & \text{s.t. } H(x, \lambda) - L_s \leq 0. \end{aligned} \quad (5)$$

To guarantee the theoretical convergence of the proposed method, instead of directly solving reformulation 5, we consider its variant using a truncated Lagrangian function,

$$L_{s,r}(x, \lambda, z) := \min_{\theta} \max_{\mu \in Z} \left\{ H(x, \theta) + \sum_{i=1}^N \mu_i g_i(x) + \frac{1}{2\gamma_1} \sum_{i=1}^N \|\theta_i - \lambda_i\|^2 - \frac{1}{2\gamma_2} \sum_{i=1}^N \|z_i - \mu_i\|^2 \right\}.$$

where  $Z := [0, r]^p$  and  $r > 0$ . We define  $\theta^* := \theta^*(x, \lambda, z)$  and  $\mu^* := \mu^*(x, \lambda, z)$  is the unique saddle point of the above min-max problem. Compared with  $L_s(x, \mu)$ , the truncated version  $L_{s,r}(x, \mu)$  is defined by maximizing  $z$  over a bounded set  $Z$ . The truncated Lagrangian value function gives us the following variant to a reformulation of problem 5

$$\begin{aligned} & \min_x F(x, \lambda), \\ & \text{s.t. } H(x, \lambda) - L_{s,r} \leq 0. \end{aligned} \quad (6)$$

Note that  $\|\sum_{i=1}^m \lambda_i \nabla f_i(x)\| - L_s \leq 0$  for any  $x, \lambda$  in their domain. If  $r$  is sufficiently large, the solution of reformulation of problem 5 can be obtained by solving variant problem 6. A comprehensive proof is presented in Theorem A.2 within Appendix A.3.

Then, problem 6 is solved by introducing a penalty parameter  $\{c^{(t)}\}_{t=0}^{T-1}$ , where  $t$  is the round index,

$$\min_{x, \lambda} \frac{1}{c^{(t)}} F(x, \lambda) + H(x, \lambda) - L_{s,r}.$$

The detailed steps are provided in Algorithm 1.

## 4 CONVERGENCE ANALYSIS

### 4.1 CONVERGENCE RESULTS OF MLM-CMOO

**Assumption 4.1** *In the CMOO, suppose the objectives  $f_i$  and the constraints  $g_i$  are convex for any  $i \in \mathcal{N}$ .*

**Algorithm 1** Moreau envelope-based Lagrange Multiplier Constrained MOO (MLM-CMOO)

---

```

1: Input: Initial point  $x^{(0)}$ .  $\lambda^{(0)}$ ,  $\theta^{(0)}$ .  $\mu^{(0)}$ , penalty parameter  $\{c^{(t)}\}_{t=0}^{T-1}$ ;
2: for  $t = 0$  to  $T - 1$  do
3:    $\theta^{(t+1)} = \theta^{(t)} - \eta \nabla_{\theta} L_{s,r}(x^{(t)}, \lambda^{(t)}, \mathbf{z}^{(t)})$ ;
4:    $\mu^{(t+1)} = \mu^{(t)} - \eta \nabla_{\mu} L_{s,r}(x^{(t)}, \lambda^{(t)}, \mathbf{z}^{(t)})$ ;
5:    $U = \arg \min_u \{H(u, \lambda^{(t)}) + \frac{1}{2} \|u - x^{(t)}\|^2\} - \arg \min_u \{H(u, \theta^{(t+1)}) + \frac{1}{2} \|u - x^{(t)}\|^2\}$ ;
6:    $x^{(t+1)} = x^{(t)} - \alpha \left( \frac{1}{c^{(t)}} \nabla_x F(x^{(t)}, \lambda^{(t)}) + U + \sum_{i=1}^n \mu_i^{(t+1)} \nabla_x g_i(x^{(t)}) \right)$ ;
7:    $\lambda' = \arg \min_{\lambda} \left\| \sum_{i=1}^m \lambda_i \nabla f_i(x^{(t)}) \right\|$ ;
8:    $V = \arg \min_v \{H(x^{(t)}, v) + \frac{1}{2} \|v - \lambda^{(t)}\|^2\} - \arg \min_v \{H(x^{(t)}, \theta^{(t+1)}) + \frac{1}{2} \|v - \lambda^{(t)}\|^2\}$ ;
9:    $\lambda^{(t+1)} = \lambda^{(t)} - \alpha \left( \frac{1}{c^{(t)}} (\lambda^{(t)} - \lambda') + V - \frac{1}{\gamma_1} (\lambda^{(t)} - \theta)^{(t+1)} \right)$ ;
10:   $\mathbf{z}^{(t+1)} = \mathbf{z}^{(t)} + \frac{\beta}{\gamma_2} (\mu^{(t+1)} - \mathbf{z}^{(t)})$ 
11: end for

```

---

**Assumption 4.2** For the general MOO, there exists a finite constant  $B \in \mathbb{R}$ , such that  $0 \leq \lambda_i^{(t)} \leq B$ ,  $\sum_{i=1}^N \lambda_i^{(t)} = 1$ , for all  $t = 0, \dots, T - 1$ .

**Assumption 4.3**  $f_1(x), \dots, f_N(x)$  are all differentiable,  $S_f$ -Lipschitz and  $L_f$ -smoothness, suggesting that for all  $x, y$  and  $i \in \mathcal{N}$ , it holds  $\|\nabla f_i(x)\| \leq S_f$  and  $\|\nabla f_i(x) - \nabla f_i(y)\| \leq L_f \|x - y\| \leq 2S_f$ .

**Assumption 4.4**  $g_1(x), \dots, g_N(x)$  are all differentiable,  $S_g$ -Lipschitz and  $L_g$ -smoothness, suggesting that for all  $x, y$  and  $i \in \mathcal{N}$ , it holds  $\|\nabla f_i(x)\| \leq S_g$  and  $\|\nabla f_i(x) - \nabla f_i(y)\| \leq L_g \|x - y\| \leq 2S_g$ .

We first define a new function and then demonstrate the decreasing properties of this new function to show the convergence rate of the proposed algorithm.

$$V_t := \phi_{c_t}(x^{(t)}, \lambda^{(t)}, \mathbf{z}^{(t)}) + C_{\theta, \mu} \left\| (\theta^{(t)}, \mu^{(t)}) - (\theta^*(x^{(t)}, \lambda^{(t)}, \mathbf{z}^{(t)}), \mu^*(x^{(t)}, \lambda^{(t)}, \mathbf{z}^{(t)})) \right\|^2,$$

where  $C_{\theta, \mu} := \max\{(L_g + C_Z L_g)^2 + 1/(2\gamma_1^2) + L_g^2, 1/\gamma_2^2\}$ , and  $\phi_{c_t}(x^{(t)}, \lambda^{(t)}, \mathbf{z}^{(t)}) := \frac{1}{c^{(t)}}(F(x, \lambda) - F^*) + H(x, \lambda) - L_{s,r}$ , where  $F^*$  is the optimal value of function  $F(x, \lambda)$ .

**Lemma 4.5** Under Assumptions 4.2, 4.3 and 4.4 hold, let  $\gamma_1 \in (0, 1/\rho_T)$ ,  $\gamma_2 > 0$ ,  $c_t \leq c_{t+1}$  and  $\eta_t \in (\eta, \rho_\gamma/L_B^2)$  with  $\eta > 0$ , then there exist constants  $c_\alpha, c_\beta > 0$  such that when  $0 < \alpha \leq c_\alpha$  and  $0 < \beta \leq c_\beta$ , the sequence of  $(x^{(t)}, \lambda^{(t)}, \mathbf{z}^{(t)})$  generated by Algorithm 1: MLM-CMOO satisfies

$$\begin{aligned} V_{t+1} - V_t &\leq -\frac{1}{4\alpha} \left\| x^{(t+1)} - x^{(t)} \right\|^2 - \frac{1}{4\alpha} \left\| \lambda^{(t+1)} - \lambda^{(t)} \right\|^2 - \frac{1}{4\beta} \left\| \mathbf{z}^{(t+1)} - \mathbf{z}^{(t)} \right\|^2 \\ &\quad - \eta \rho_T C_{\theta, \mu} \left\| (\theta^{(t)}, \mu^{(t)}) - (\theta^*(x^{(t)}, \lambda^{(t)}, \mathbf{z}^{(t)}), \mu^*(x^{(t)}, \lambda^{(t)}, \mathbf{z}^{(t)})) \right\|^2. \end{aligned}$$

The step sizes are carefully chosen to guarantee the sufficient descent property of  $V_t$ . This is essential for the non-asymptotic convergence analysis.

Given the decreasing property of  $V_t$ , we establish the non-asymptotic convergence analysis. The standard KKT conditions are inappropriate as necessary optimality conditions for problem equation 6. Motivated by the approximate KKT condition presented by (Andreani et al., 2010), which is characterized as an optimality condition for nonlinear program, regardless of constraint qualifications' fulfillment, we consider the following residual function  $R_t := R_t(x^{(t)}, \lambda^{(t)}, \mathbf{z}^{(t)})$  as a stationarity measure, we define the residual function as  $R_t := \text{dist}(0, (\nabla F(x^{(t)}, \lambda^{(t)}), 0)) + c_t((\nabla H(x^{(t)}, \lambda^{(t)}), 0) - \nabla L_{s,r}(x^{(t)}, \lambda^{(t)}, \mathbf{z}^{(t)})) + \mathcal{M}_{C \times Z}(x^{(t)}, \lambda^{(t)}, \mathbf{z}^{(t)})$ , where  $\mathcal{M}_\Omega(s)$  denotes the normal cone to  $\Omega$  at  $s$ . This residual function  $R_t$  also serves as a stationarity measure for the penalized problem of equation 6, with  $c_t$  serving as the penalty parameter,

$$\min_{x, \lambda, \mathbf{z} \in Z} \phi_{c_t}(x^{(t)}, \lambda^{(t)}, \mathbf{z}^{(t)}) := F(x^{(t)}, \lambda^{(t)}) + c_t \left( H(x^{(t)}, \lambda^{(t)}) - L_{s,r}(x^{(t)}, \lambda^{(t)}, \mathbf{z}^{(t)}) \right). \quad (7)$$

Evidently,  $R_t = 0$  if and only if  $(x^{(t)}, \boldsymbol{\lambda}^{(t)}, \mathbf{z}^{(t)})$  is a stationary point for the problem equation 7, meaning  $0 \in \nabla \phi_{c_t}(x^{(t)}, \boldsymbol{\lambda}^{(t)}, \mathbf{z}^{(t)}) + \mathcal{M}_{C \times Z}(x^{(t)}, \boldsymbol{\lambda}^{(t)}, \mathbf{z}^{(t)})$ .

**Theorem 4.6** *If Assumptions of Assumptions 4.2, 4.3 and 4.4 hold, let  $\gamma_1 \in (0, 1/\rho_\gamma)$ ,  $\gamma_2 > 0$ ,  $c_t = \underline{c}(t+1)^p$  with  $p \in (0, 1/2)$  and  $\underline{c} > 0$ . Pick  $\eta_t \in (0, \rho_\gamma/L_B^2)$ , then there exists  $c_\alpha, c_\beta > 0$  such that when  $\alpha \in (\underline{\alpha}, c_\alpha)$  and  $\beta \in (\underline{\beta}, c_\beta)$ , with  $\underline{\alpha}, \underline{\beta} > 0$ , the sequence of  $(x^{(t)}, \boldsymbol{\lambda}^{(t)}, \mathbf{z}^{(t)}, \boldsymbol{\theta}^{(t)}, \boldsymbol{\mu}^{(t)})$  generated by Algorithm 1: MLM-CMOO satisfies*

$$\min_t \left\| (\boldsymbol{\theta}^{(t)}, \boldsymbol{\mu}^{(t)}) - (\boldsymbol{\theta}^*(x^{(t)}, \boldsymbol{\lambda}^{(t)}, \mathbf{z}^{(t)}), \boldsymbol{\mu}^*(x^{(t)}, \boldsymbol{\lambda}^{(t)}, \mathbf{z}^{(t)})) \right\| = \mathcal{O}\left(\frac{1}{\sqrt{T}}\right),$$

and

$$\min_t R_t(x^{(t)}, \boldsymbol{\lambda}^{(t)}, \mathbf{z}^{(t)}) = \mathcal{O}\left(\frac{1}{\sqrt{T^{1-2p}}}\right).$$

**Remark 1:** Theorem 4.6 first shows that the reformulated Lagrange Multiplier function  $L_{s,r}(x, \boldsymbol{\lambda}, \mathbf{z})$  reaches its KKT stationary point with a convergence rate of  $\mathcal{O}(\frac{1}{\sqrt{T}})$ , which matches general constrained optimization method. In addition, MLM-CMOO can converge to the stationary point of the problem equation 7, where problem equation 7 is the penalized form of the original problem equation 2. Furthermore, the last statement shows that the convergence rate of each client's truncated proximal Lagrangian value function is related to the selection of the penalized parameter. The maximum convergence rate is  $\mathcal{O}(\frac{1}{\sqrt{T}})$ .

## 4.2 PROOF SKETCH

**Lemma 4.7** *Suppose the assumption of 4.3 and 4.4 hold, and let  $\gamma_1 \in (0, 1/\rho_g)$ ,  $\gamma_2 > 0$ . Pick  $\eta_t \in (0, \rho_T/L_B^2)$  with  $L_B := \max\{(2 + C_z)L_g + 1/\gamma_1, L_g + 1/\gamma_2\}$  then the sequence of  $(x^{(t)}, \boldsymbol{\lambda}^{(t)}, \mathbf{z}^{(t)}, \boldsymbol{\theta}^{(t)}, \boldsymbol{\mu}^{(t)})$  generated by Algorithm 1: MLM-CMOO satisfies*

$$\begin{aligned} & \left\| (\boldsymbol{\theta}^{(t+1)}, \boldsymbol{\mu}^{(t+1)}) - (\boldsymbol{\theta}^*(x^{(t)}, \boldsymbol{\lambda}^{(t)}, \mathbf{z}^{(t)}), \boldsymbol{\mu}^*(x^{(t)}, \boldsymbol{\lambda}^{(t)}, \mathbf{z}^{(t)})) \right\| \\ & \leq (1 - \eta_t \rho_T) \left\| (\boldsymbol{\theta}^{(t)}, \boldsymbol{\mu}^{(t)}) - (\boldsymbol{\theta}^*(x^{(t)}, \boldsymbol{\lambda}^{(t)}, \mathbf{z}^{(t)}), \boldsymbol{\mu}^*(x^{(t)}, \boldsymbol{\lambda}^{(t)}, \mathbf{z}^{(t)})) \right\| \end{aligned}$$

Lemma 4.7 shows that the decay rate of the reformulated Lagrange Multiplier function  $L_{s,r}(x, \boldsymbol{\lambda}, \mathbf{z})$  is related to the step size and problem parameters  $((1 - \eta_t \rho_T))$ . If we select the proper step size (i.e.,  $\eta_t \in (0, \rho_\gamma/L_B^2)$ ), the reformulated Lagrange Multiplier function converges to the KKT stationary point.

**Lemma 4.8** *Suppose the assumption of 4.2, 4.3 and 4.4 hold, and let  $\gamma_1 \in (0, 1/\rho_g)$ ,  $\gamma_2 > 0$ . Pick  $\eta \in (0, \rho_\gamma/L_B^2)$  with  $L_B := \max\{2L_g + C_z L_g + 1/\gamma_1, L_g + 1/\gamma_2\}$  then the sequence of  $(x^{(t)}, \boldsymbol{\lambda}^{(t)}, \mathbf{z}^{(t)})$  generated by Algorithm 1: MLM-CMOO satisfies*

$$\begin{aligned} \phi_{c_t}(x^{(t+1)}, \boldsymbol{\lambda}^{(t+1)}, \mathbf{z}^{(t+1)}) & \leq \phi_{c_t}(x^{(t)}, \boldsymbol{\lambda}^{(t)}, \mathbf{z}^{(t)}) - \left(\frac{1}{2\beta} - \frac{L_{vz}}{2}\right) \left\| \mathbf{z}^{(t+1)} - \mathbf{z}^{(t)} \right\|^2 \\ & \quad - \left(\frac{1}{2\alpha} - \frac{L_{\phi_k}}{2} - \frac{\beta L_{\theta, \mu}^2}{\gamma_2^2}\right) \left( \left\| x^{(t+1)} - x^{(t)} \right\|^2 + \left\| \boldsymbol{\lambda}^{(t+1)} - \boldsymbol{\lambda}^{(t)} \right\|^2 \right) \\ & \quad + \frac{\alpha}{2} \left( 2(L_g + C_z L_g)^2 + \frac{1}{\gamma_1^2} \right) \left\| \boldsymbol{\theta}^{(t+1)} - \boldsymbol{\theta}^*(x^{(t)}, \boldsymbol{\lambda}^{(t)}, \mathbf{z}^{(t)}) \right\|^2 \\ & \quad + \left( \alpha L_g^2 + \frac{\beta}{\gamma_2^2} \right) \left\| \boldsymbol{\mu}^{(t+1)} - \boldsymbol{\mu}^*(x^{(t)}, \boldsymbol{\lambda}^{(t)}, \mathbf{z}^{(t)}) \right\|^2. \end{aligned}$$

where  $L_{\phi_t} := L_f/c_t + L_g + \rho_v$ .

Lemma 4.8 shows that a client's Lagrangian reformulation of objects decreases with the distance between its current parameters  $(x^{(t)}, \boldsymbol{\lambda}^{(t)}, \mathbf{z}^{(t)})$ . In addition, the bias of the reformulated Lagrange Multiplier function affects the convergence rate. The more accurate the reformulated Lagrange

Multiplier function is, the faster the reformulated Lagrangian reformulation of CMOO objectives decreases. Then, the proof of Lemma 4.8 lies in 2 parts: the descent of the MOO objectives ( $F(x, \lambda)$ ) and the descent of the reformulated Lagrange Multiplier function  $L_{s,r}(x, \lambda, z)$ . The challenge of the proof lies in 2 aspects: the decreasing of the MOO objectives based on dynamic coefficient ( $\lambda$ ) and binding the reformulated Lagrange Multiplier function  $L_{s,r}(x, \lambda, z)$ .

## 5 EXPERIMENTS

This section introduces the experimental setups, including the datasets and models, multi-objective optimization algorithm setup, and experimental settings. We use multi-task learning experiments to verify the effectiveness of the proposed method. A typical MTL system is given a collection of input points and sets of targets for various tasks per point. A common way to set up the inductive bias across tasks is to design a parametrized hypothesis class that shares some parameters across tasks. One effective solution for MTL is finding solutions that are not dominated by any others, which is the same objective as MOO problems (Sener & Koltun, 2018).

### 5.1 EXPERIMENTS SETUP

**1. Datasets and Models.** 1). MultiMNIST Datasets and Learning Tasks: We test the convergence performance of MLM-CMOO using the “MultiMNIST” dataset (Sabour et al., 2017), which is a multi-task learning version of the MNIST dataset (LeCun et al., 1998) from LIBSVM repository. Specifically, to convert the hand-written classification problem into a multi-task problem, we randomly chose 60000 images. Images are divided into 2 tasks, and each task has  $m = 30000$  samples. In our experiments, a group of images is positioned in the top left corner, while another group of images is positioned in the bottom right. The two tasks are task “L” (to categorize the top-left digit) and task “R” (to classify the bottom-right digit). The overall problem is to classify the images of different tasks. All algorithms use the same randomly generated initial point. The learning rates are chosen as  $\eta = \beta = \alpha = 0.01$ . we directly apply existing single-task MNIST models.

2). CelebA Dataset and Learning Tasks: We utilize the CelebA dataset (Liu et al., 2015), consisting of 200K facial images annotated with 40 attributes. We approach each attribute as a binary classification task, resulting in a 40-way multi-task learning (MTL) problem. To create a shared representation function, we implement ResNet-18 (He et al., 2016) without the final layer, attaching a linear layer to each attribute for classification. In this experiment, we set  $\eta = 0.0005, \alpha = \beta = 0.1$ .

3). River Flow Dataset and Learning Tasks: We further test our algorithms on MOO problems of larger sizes. In this experiment, we use the River Flow dataset (Nie et al., 2017), which is for flow prediction flow at eight locations within the Mississippi River network. Thus, there are eight tasks in this problem. In this experiment, we set  $\eta = 0.001, \alpha = \beta = 0.1$ . To better visualize 8 different tasks, we illustrate the normalized loss in radar charts.

**2. Baseline.** The NSGA-II (Deb et al., 2002) and PSL (Lin et al., 2022) are considered as our baselines. NSGA-II is a well-known MOO evolutionary algorithm, while PSL is a novel MOO Bayesian optimization algorithm. NSGA-II and PSL do not handle constraints. For a fair comparison, we report the Pareto optimal solutions satisfying constraints generated from NSGA-II and PSL, thus NSGA-II and PSL can work towards the Pareto optimal solutions without considering constraints.

**NSGA-II Setup.** For binary chromosomes, we apply a single-point crossover with a probability of 0.9 and a bit-flip mutation with a probability of 0.1. For real-valued chromosomes, we apply a simulated binary crossover (SBX) [14] with a probability of 0.9 and  $n_c = 2$  and a polynomial mutation with a probability of 0.1 and  $n_m = 20$ , where  $n_c$  and  $n_m$  denote spread factor distribution indices for crossover and mutation, respectively.

**PSL Setup.** We follow literature (Lin et al., 2022) to set PSL parameters. At each iteration, we train the Pareto set model  $h_\theta$  with 1000 update steps using Adam optimizer with a learning rate of  $1e^{-5}$  and no weight decay. At each iteration, we generate 1000 candidate solutions using  $h_\theta$  and select the population size of solutions from the 1000 candidates.



5.2 EXPERIMENTS RESULTS

For the MultiMNIST Datasets and Learning Tasks. Table 5.2 shows that the MLM-CMOO outperforms than the rest 2 baselines. This is because that NSGA-II and PSL are designed for unconstrained MOO problems, even we have generated the Pareto optimal solutions satisfying constraints from those algorithm, their evolutionary or Bayesian Optimisation are not the most fit algorithm. On the contrary, MLM-CMOO is specifically designed for constrained MOO problems, where it uses an increasing Lagrange coefficient ( $\mu$ ) to make it follow the limits. In this way, MLM-CMOO is more time-efficient. The Fig 1(a) and Fig 1(b) shows the results of 3 algorithm on CelebA Dataset and Learning Tasks. It shows that MLM-CMOO matches the selected baselines. The Fig 1(c) shows that MLM-CMOO’s loss on River Flow Dataset and Learning Tasks is better than selected algorithms.

Table 1: Completion time are taken to reach Pareto stationary point with a specified loss for different algorithms using MultiMNIST Datasets.

Loss	$10^{-2}$					
Task	Task L			Task R		
Algorithm	NSGA-II	PSL	MLM-CMOO	NSGA-II	PSL	MLM-CMOO
Time (s)	$\times 2.2$	$\times 1.8$	1302	$\times 2.3$	$\times 1.7$	1322
Loss	$10^{-3}$					
Task	Task L			Task R		
Algorithm	NSGA-II	PSL	MLM-CMOO	NSGA-II	PSL	MLM-CMOO
Time (s)	$\times 3.2$	$\times 2.03$	2057	$\times 3.1$	$\times 2.1$	2104

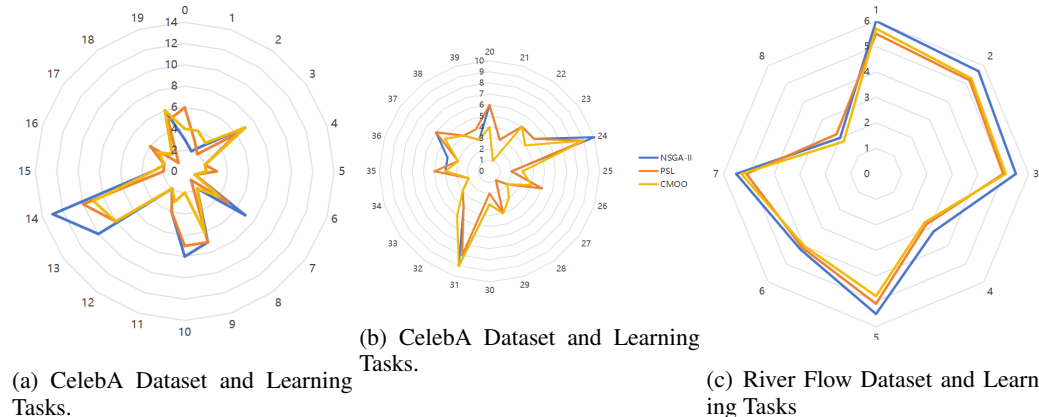


Figure 1: Experiments on CelebA dataset and River Flow Dataset.

6 CONCLUSION

This paper studies the constrained multi-objective optimization problem. We first establish a framework for the CMOO problem, which is suitable for gradient descent algorithms. For the formulated CMOO problem, we use the Lagrange Multiplier method to make a decrease in the overall objective obey the constraints. Then, due to the non-smoothness of the coefficient function, we use a Moreau envelope to make it smooth. Next, the convergence analysis shows that the proposed algorithm (MLM-CMOO) convergence to Pareto stationary solutions with a rate of  $\mathcal{O}(\frac{1}{\sqrt{T}})$ . Finally, we conduct experiments to verify the effectiveness of the MLM-CMOO algorithm.

REFERENCES

Alaleh Ahmadianshalchi, Syrine Belakaria, and Janardhan Rao Doppa. Preference-aware constrained multi-objective bayesian optimization (student abstract). *Proceedings of the AAAI Con-*

- 486 *ference on Artificial Intelligence (AAAI)*, 2024.
- 487
- 488 Roberto Andreani, J. M. Martínez, and B. F. Svaiter. A new sequential optimality condition for con-  
489 strained optimization and algorithmic consequences. *SIAM Journal on Optimization*, 20, 2010.
- 490 Amir Beck. First-order methods in optimization. In *MOS-SIAM Series on Optimization*, 2017.
- 491
- 492 Leyla Belaiche, Laid Kahloul, Maroua Grid, Nedjma Abidallah, and Saber Benharzallah. Parallel  
493 multi-objective evolutionary algorithm for constrained multi-objective optimization. In *2023 24th*  
494 *International Arab Conference on Information Technology (ACIT)*, 2023.
- 495
- 496 Syrine Belakaria, Aryan Deshwal, Nitthilan Kannappan Jayakodi, and Janardhan Rao Doppa.  
497 Uncertainty-aware search framework for multi-objective bayesian optimization. *Proceedings of*  
498 *the AAAI Conference on Artificial Intelligence (AAAI)*, 2020.
- 499
- 500 A. L. Custódio, J. F. A. Madeira, A. I. F. Vaz, and L. N. Vicente. Direct multisearch for multiobjec-  
501 tive optimization. *SIAM Journal on Optimization*, 21(3), 2011.
- 502
- 503 K. Deb, A. Pratap, S. Agarwal, and T. Meyarivan. A fast and elitist multiobjective genetic algorithm:  
504 Nsga-ii. *IEEE Transactions on Evolutionary Computation*, 6(2), 2002.
- 505
- 506 Anh Viet Do, Aneta Neumann, Frank Neumann, and Andrew M. Sutton. Rigorous runtime anal-  
507 ysis of MOEA/d for solving multi-objective minimum weight base problems. In *Thirty-seventh*  
508 *Conference on Neural Information Processing Systems (NeurIPS)*, 2023.
- 509
- 510 Jean-Antoine Désidéri. Multiple-gradient descent algorithm (mgda) for multiobjective optimization.  
511 *Comptes Rendus Mathématique*, 350(5), 2012.
- 512
- 513 Heshan Fernando, Han Shen, Miao Liu, Subhajit Chaudhury, Keerthiram Murugesan, and Tianyi  
514 Chen. Mitigating gradient bias in multi-objective learning: A provably convergent stochastic  
515 approach, 2024.
- 516
- 517 J. Fliege and B. F. Svaiter. Steepest descent methods for multicriteria optimization. *Mathematical*  
518 *Methods of Operations Research*, 51(3), 2000.
- 519
- 520 E. H. Fukuda H. Tanabe and N. Yamashita. Convergence rates analysis of a multiobjective proximal  
521 gradient method. *Optimization Letters*, 17(2), 2023.
- 522
- 523 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-  
524 nition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- 525
- 526 R. Hinterding and Z. Michalewicz. Your brains and my beauty: parent matching for constrained  
527 optimisation. In *1998 IEEE International Conference on Evolutionary Computation Proceedings.*  
528 *IEEE World Congress on Computational Intelligence (Cat. No.98TH8360)*, 1998.
- 529
- 530 A. I. F. Vaz J. Fliege and L. N. Vicente. Complexity of gradient descent for multiobjective optimiza-  
531 tion. *Optimization Methods and Software*, 34(5), 2019.
- 532
- 533 Dae Gyu Kim and Phil Husbands. Mapping based constraint handling for evolutionary search;  
534 thurston’s circle packing and grid generation. In *Adaptive Computing in Design and Manufacture*,  
535 pp. 161–173, 1998a.
- 536
- 537 Dae Gyu Kim and Phil Husbands. Landscape changes and the performance of mapping based  
538 constraint handling methods. In *Parallel Problem Solving from Nature — PPSN V*, pp. 221–229,  
539 1998b.
- 534
- 535 Slawomir Koziel and Zbigniew Michalewicz. A decoder-based evolutionary algorithm for con-  
536 strained parameter optimization problems. In *Parallel Problem Solving from Nature*, 1998.
- 537
- 538 Slawomir Koziel and Zbigniew Michalewicz. Evolutionary algorithms, homomorphous mappings,  
539 and constrained parameter optimization. *Evolutionary Computation*, 7(1):19–44, 1999.
- Oliver Kramer. A review of constraint-handling techniques for evolution strategies. 2010, 2010.

- 540 Marco Laumanns and Jiri Ocenasek. Bayesian optimization algorithms for multi-objective opti-  
541 mization. In *Parallel Problem Solving from Nature — PPSN VII*, 2002.
- 542
- 543 Y. LeCun, C. Cortes, and C. Burges. Mnist handwritten digit database, 1998. URL [http://](http://yann.lecun.com/exdb/mnist)  
544 [yann.lecun.com/exdb/mnist](http://yann.lecun.com/exdb/mnist).
- 545
- 546 Wenji Li, Ruitao Mai, Pengxiang Ren, Zhaojun Wang, Qinchang Zhang, Ning Xu, Biao Xu, Zhun  
547 Fan, and Zhifeng Hao. A surrogate-ensemble assisted coevolutionary algorithm for expensive  
548 constrained multi-objective optimization problems. In *2023 IEEE Congress on Evolutionary*  
549 *Computation (CEC)*, 2023.
- 550 Xi Lin, Zhiyuan Yang, Xiaoyuan Zhang, and Qingfu Zhang. Pareto set learning for expensive multi-  
551 objective optimization. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- 552
- 553 Bo Liu, Xingchao Liu, Xiaojie Jin, Peter Stone, and Qiang Liu. Conflict-averse gradient descent for  
554 multi-task learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- 555
- 556 S. Liu and L. N. Vicente. The stochastic multi-gradient algorithm for multi-objective optimization  
557 and its application to supervised machine learning. *Annals of Operations Research*, 339(3), 2024.
- 558
- 559 Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild.  
560 In *2015 IEEE International Conference on Computer Vision (ICCV)*, 2015.
- 561
- 562 Zbigniew Michalewicz. *Genetic Algorithms + Data Structures = Evolution Programs*. Springer  
563 Berlin, Heidelberg, 1996. ISBN 978-3-540-60676-5.
- 564
- 565 Michinari Momma, Chaosheng Dong, and Jia Liu. A multi-objective / multi-task learning frame-  
566 work induced by pareto stationarity. In *Proceedings of the 39th International Conference on*  
567 *Machine Learning (ICML)*, 2022.
- 568
- 569 Lin Nie, Keze Wang, Wenxiong Kang, and Yuefang Gao. Image retrieval with attribute-associated  
570 auxiliary references. In *2017 International Conference on Digital Image Computing: Techniques*  
571 *and Applications (DICTA)*, 2017.
- 572
- 573 Sebastian Peitz and Michael Dellnitz. Gradient-based multiobjective optimization with uncertain-  
574 ties. In *Results of the Numerical and Evolutionary Optimization Workshop NEO 2016 and the*  
575 *NEO Cities 2016 Workshop (NEO)*, 2018.
- 576
- 577 David Powell and Michael M. Skolnick. Using genetic algorithms in engineering design optimiza-  
578 tion with non-linear constraints. In *Proceedings of the 5th International Conference on Genetic*  
579 *Algorithms*, 1993. ISBN 1558602992.
- 580
- 581 Sara Sabour, Nicholas Frosst, and Geoffrey E. Hinton. Dynamic routing between capsules. In  
582 *Proceedings of the 31st International Conference on Neural Information Processing Systems*  
583 *(NeurIPS)*, 2017.
- 584
- 585 Marc Schoenauer and Zbigniew Michalewicz. Evolutionary computation at the edge of feasibility.  
586 In *Parallel Problem Solving from Nature — PPSN IV*, pp. 245–254, 1996.
- 587
- 588 Marc Schoenauer and Zbigniew Michalewicz. Boundary operators for constrained parameter opti-  
589 mization problems. In *International Conference on Genetic Algorithms*, 1997.
- 590
- 591 Ozan Sener and Vladlen Koltun. Multi-task learning as multi-objective optimization. In *Advances*  
592 *in Neural Information Processing Systems (NeurIPS)*, 2018.
- 593
- 594 Shiquan Song, Kai Zhang, Ling Zhang, and Ni Wu. A dual-population algorithm based on self-  
595 adaptive epsilon method for constrained multi-objective optimization. *Information Sciences*, 655:  
596 119906, 2024.
- 597
- 598 Peiyao Xiao, Hao Ban, and Kaiyi Ji. Direction-oriented multi-objective learning: Simple and prov-  
599 able stochastic algorithms, 2023.
- 600
- 601 Haibo Yang, Zhuqing Liu, Jia Liu, Chaosheng Dong, and Michinari Momma. Federated multi-  
602 objective learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.

594 Lei Yang, Jinglin Tian, Jiale Cao, Kangshun Li, and Chaoda Peng. An evolutionary algorithm with  
595 feasibility tracking strategy for constrained multi-objective optimization problems. In *2024 IEEE*  
596 *Congress on Evolutionary Computation (CEC)*, 2024.

597 Wei Yao, Chengming Yu, Shangzhi Zeng, and Jin Zhang. Constrained bi-level optimization: Proxi-  
598 mal lagrangian value function approach and hessian-free algorithm. In *The Twelfth International*  
599 *Conference on Learning Representations (ICLR)*, 2024.

600

601 Tianhe Yu, Saurabh Kumar, Abhishek Gupta, Sergey Levine, Karol Hausman, and Chelsea Finn.  
602 Gradient surgery for multi-task learning. In *Advances in Neural Information Processing Systems*  
603 *(NeurIPS)*, 2020.

604

605 Qingfu Zhang and Hui Li. Moea/d: A multiobjective evolutionary algorithm based on decomposi-  
606 tion. *IEEE Transactions on Evolutionary Computation*, 11(6), 2007.

607 Xuerui Zhang, Zhongyang Han, Zhiyuan Wang, Jun Zhao, and Wei Wang. Ingredient planning for  
608 copper industry: A deep reinforcement learning-based  $\epsilon$ -constrained multi-objective optimization  
609 framework. In *2024 IEEE Congress on Evolutionary Computation (CEC)*, 2024.

610

611 Weijie Zheng, Yufei Liu, and Benjamin Doerr. A first mathematical runtime analysis of the non-  
612 dominated sorting genetic algorithm ii (nsga-ii). *Proceedings of the AAAI Conference on Artificial*  
613 *Intelligence(AAAI)*, 2022.

614 Shiji Zhou, Wenpeng Zhang, Jiyan Jiang, Wenliang Zhong, Jinjie GU, and Wenwu Zhu. On the  
615 convergence of stochastic multi-objective gradient manipulation and beyond. In *Advances in*  
616 *Neural Information Processing Systems (NeurIPS)*, 2022.

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

## Supplementary Material for "Constrained Multi-Objective Optimization"

### A SOME USEFUL LEMMA

To simplify the notions, we use  $\omega^{(t)} = (x^{(t)}, \lambda^{(t)}, z^{(t)})$  in the following proof.

**Lemma A.1** Yao et al. (2024) **Lemma A.1.** Under Assumption 4.3 and 4.4, for the Moreau envelope-based Lagrange Multiplier function  $L_s(x, \lambda, z)$  with  $\gamma_1 \in (0, 1/\rho_f)$  and  $\gamma_2 > 0$ . That is,

(1) The function  $L_s(x, \lambda, z)$  is continuously differentiable;

(2) The gradient of  $L_s(x, \lambda, z)$  has closed-form given by

$$\begin{aligned}\nabla_x L_s(x, \lambda, z) &= \arg \min_u \left\{ H(u, \theta^*) + \frac{1}{2} \|u - x\|^2 \right\} + \sum_{i=1}^n \mu_i^* \nabla g_i(x), \\ \nabla_\lambda L_s(x, \lambda, z) &= \frac{\lambda - \theta^*}{\gamma_1}, \\ \nabla_z L_s(x, \lambda, z) &= \frac{\mu^* - z}{\gamma_2},\end{aligned}$$

where  $\theta^* := \theta^*(x, \lambda, z)$  and  $\mu^* := \mu^*(x, \lambda, z)$  is the unique saddle point of the following min-max problem:

$$\min_{\theta} \max_{\mu} \left\{ H(x, \theta) + \sum_{i=1}^N \mu_i g_i(x) + \frac{1}{2\gamma_1} \sum_{i=1}^N \|\theta_i - \lambda_i\|^2 - \frac{1}{2\gamma_2} \sum_{i=1}^N \|z_i - \mu_i\|^2 \right\}.$$

(3) Furthermore, for any  $\rho_v \geq \rho_f/(1 - \gamma_1\rho_f)$ ,  $L_s(x, \lambda, z)$  is  $\rho_v$ -weakly convex with respect to variables  $(x, \lambda)$  on for any fixed  $z$ .

*Proof:* The proof is similar to the proof of **Lemma A.1** in Yao et al. (2024).

**Lemma A.2** Yao et al. (2024) **Lemma A.2** and **Lemma A.4.**

Under Assumption 4.3 and 4.4, let  $\gamma_1 \in (0, 1/\rho_f)$  and  $\gamma_2 > 0$ . Then, for any  $\rho_v \geq \rho_f/(1 - \gamma_1\rho_f)$ , the following inequality holds:

$$\begin{aligned}-L_s(x_1, \lambda, z) &\leq -L_s(x_2, \lambda, z) - \langle \nabla_x L_s(x_2, \lambda, z), x_1 - x_2 \rangle + \frac{\rho_v}{2} \|x_1 - x_2\|^2, \\ -L_s(x, \lambda_1, z) &\leq -L_s(x, \lambda_2, z) - \langle \nabla_\lambda L_s(x, \lambda_2, z), \lambda_1 - \lambda_2 \rangle + \frac{\rho_v}{2} \|\lambda_1 - \lambda_2\|^2, \\ -L_s(x, \lambda, z_1) &\leq -L_s(x, \lambda, z_2) - \langle \nabla_z L_s(x, \lambda, z_2), z_1 - z_2 \rangle + \frac{L_z}{2} \|z_1 - z_2\|^2,\end{aligned}$$

where  $L_z := (\gamma_2\rho_T + 1)/(\gamma_2^2\rho_T)$ .

*Proof:* The first 2 conclusions follow directly from **Lemma A.2** that  $L_s(x, \lambda, z)$  is  $\rho_v$ -weakly convex with respect to variables  $(x, \lambda)$  on for any fixed  $z$ , and the third conclusion is similar to the proof of **Lemma A.4** in Yao et al. (2024).

**Lemma A.3** Yao et al. (2024) **Lemma A.3.** Under Assumption 4.3 and 4.4, let  $\gamma_1 \in (0, 1/\rho_f)$  and  $\gamma_2 > 0$ . Then, for any  $(x_1, \lambda_1, z_1)$  and  $(x_2, \lambda_2, z_2)$ , the following Lipschitz property holds:

$$\begin{aligned}& \|(\theta^*(x_1, \lambda_1, z_1), \mu^*(x_1, \lambda_1, z_1)) - (\theta^*(x_2, \lambda_2, z_2), \mu^*(x_2, \lambda_2, z_2))\| \\ & \leq \frac{L_f + L_g + C_Z L_g}{\rho_T} \|x_1 - x_2\| + \frac{1}{\gamma_1\rho_T} \|\lambda_1 - \lambda_2\| + \frac{1}{\gamma_2\rho_T} \|z_1 - z_2\| \\ & \leq L_{\theta, \mu} \|(x_1, \lambda_1, z_1) - (x_2, \lambda_2, z_2)\|,\end{aligned}$$

where  $\rho_T := \min\{1/\gamma_1 - \rho_f, 1/\gamma_2\}$ ,  $C_Z = \max_{z \in Z} \|z\|$ , and  $L_{\theta, \mu} := \sqrt{3} \max\{L_f + L_g + C_Z L_g, 1/\gamma_1, 1/\gamma_2\}/\rho_T$ .

*Proof:* The proof is similar to the proof of **Lemma A.3** in Yao et al. (2024).

## B PROOF OF MAIN THEOREM AND LEMMAS

### B.1 PROOF OF THEOREM 4.6

**Theorem 4.6** If Assumptions of Assumptions 4.2, 4.3 and 4.4 hold, let  $\gamma_1 \in (0, 1/\rho_\gamma)$ ,  $\gamma_2 > 0$ ,  $c_t = \underline{c}(t+1)^p$  with  $p \in (0, 1/2)$  and  $\underline{c} > 0$ . Pick  $\eta_t \in (0, \rho_\gamma/L_B^2)$ , then there exists  $c_\alpha, c_\beta > 0$  such that when  $\alpha \in (\underline{\alpha}, c_\alpha)$  and  $\beta \in (\underline{\beta}, c_\beta)$ , with  $\underline{\alpha}, \underline{\beta} > 0$ , the sequence of  $(x^{(t)}, \lambda^{(t)}, z^{(t)}, \theta^{(t)}, \mu^{(t)})$  generated by Algorithm 1: MLM-CMOO satisfies

$$\min_t \left\| (\theta^t, \mu^t) - (\theta_r^*(x^{(t)}, \lambda^{(t)}, z^{(t)}), \mu_r^*(x^{(t)}, \lambda^{(t)}, z^{(t)})) \right\| = \mathcal{O}\left(\frac{1}{\sqrt{T}}\right),$$

and

$$\min_t R_t(x^{(t)}, \lambda^{(t)}, z^{(t)}) = \mathcal{O}\left(\frac{1}{\sqrt{T^{1-2p}}}\right).$$

*Proof:* First, using the descent lemma in Lemma 4.5 and its condition, telescoping the inequality for  $t = 0, 1, \dots, T-1$ , we get

$$\begin{aligned} V_T - V_0 &\leq -\frac{1}{4\alpha} \sum_{t=0}^{T-1} \left( \|x^{(t+1)} - x^{(t)}\|^2 + \|\lambda^{(t+1)} - \lambda^{(t)}\|^2 \right) - \frac{1}{4\beta} \sum_{t=0}^{T-1} \|z^{(t+1)} - z^{(t)}\|^2 \\ &\quad - \eta \rho_T C_{\theta, \mu} \sum_{t=0}^{T-1} \left\| (\theta^{(t)}, \mu^{(t)}) - (\theta^*(x^{(t)}, \lambda^{(t)}, z^{(t)}), \mu^*(x^{(t)}, \lambda^{(t)}, z^{(t)})) \right\|^2. \end{aligned}$$

From assumptions, we have  $\sum_{t=0}^{T-1} \left\| (\theta^{(t)}, \mu^{(t)}) - (\theta^*(x^{(t)}, \lambda^{(t)}, z^{(t)}), \mu^*(x^{(t)}, \lambda^{(t)}, z^{(t)})) \right\|^2$  is upper bounded, which is

$$\sum_{t=0}^{T-1} \left\| (\theta^{(t)}, \mu^{(t)}) - (\theta^*(x^{(t)}, \lambda^{(t)}, z^{(t)}), \mu^*(x^{(t)}, \lambda^{(t)}, z^{(t)})) \right\|^2 \leq +\infty.$$

Thus, we have

$$\min_t \left\| (\theta^{(t)}, \mu^{(t)}) - (\theta^*(x^{(t)}, \lambda^{(t)}, z^{(t)}), \mu^*(x^{(t)}, \lambda^{(t)}, z^{(t)})) \right\| = \mathcal{O}\left(\frac{1}{\sqrt{T}}\right).$$

Secondly, according to the update rule of variables  $(x, y, z)$ , we have

$$\begin{aligned} 0 &\in c_t(d_x^{(t)}, d_\lambda^{(t)}) + \mathcal{M}_C(x^{(t)}, \lambda^{(t)}) + \frac{c_t}{\alpha}((x^{(t+1)}, \lambda^{(t+1)}) - (x^{(t)}, \lambda^{(t)})), \\ 0 &\in c_t d_z^{(t)} + \mathcal{M}_Z(z^{(t+1)}) + \frac{c_t}{\beta}(z^{(t+1)} - z^{(t)}). \end{aligned}$$

where  $d_x^{(t)} = \frac{1}{c^{(t)}} \nabla_x F(x^{(t)}, \lambda^{(t)}) + U + \sum_{i=1}^n \mu_i^{(t+1)} \nabla_x g_i(x^{(t)})$ ,  $d_\lambda^{(t)} = \frac{1}{c^{(t)}}(\lambda^{(t)} - \lambda') + \mathbf{V} - \frac{1}{\gamma_1}(\lambda^{(t)} - \theta^{(t+1)})$ , and  $d_z^{(t)} = \mu^{(t+1)} - z^{(t)}$ . Note,  $U = \arg \min_u \{H(u, \lambda^{(t)}) + \frac{1}{2} \|u - x^{(t)}\|^2\} - \arg \min_u \{H(u, \theta^{(t+1)}) + \frac{1}{2} \|u - x^{(t)}\|^2\}$ ,  $\lambda' = \arg \min_\lambda \|\sum_{i=1}^m \lambda_i \nabla f_i(x^{(t)})\|$ , and  $\mathbf{V} = \arg \min_v \{H(x^{(t)}, v) + \frac{1}{2} \|v - \lambda^{(t)}\|^2\} - \arg \min_v \{H(u, \theta^{(t+1)}) + \frac{1}{2} \|u - x^{(t)}\|^2\}$ .

By the meanings of  $d_x^{(t)}$ ,  $d_\lambda^{(t)}$ , and  $d_z^{(t)}$ , we obtain

$$\begin{aligned} (e_{x, \lambda}^{(t)}, e_z^{(t)}) &\in (\nabla F(x^{(t+1)}, \lambda^{(t+1)}), 0) + c_t \left( \sum_{i=1}^n \mu_i \nabla g_i(x^{(t+1)}), 0 \right) \\ &\quad - c_t (\nabla L_{i, s, r}(x^{(t+1)}, \lambda^{(t+1)}, z^{(t+1)}) + \mathcal{M}_{C \times Z}(x^{(t+1)}, \lambda^{(t+1)}, z^{(t+1)})), \end{aligned}$$

where

$$e_{x, \lambda}^{(t)} := \nabla_{x, \lambda} \phi_{c_t}(x^{(t)}, \lambda^{(t)}, z^{(t)}) - c_t(d_x^{(t)}, d_\lambda^{(t)}) - \frac{c_t}{\alpha}((x^{(t+1)}, \lambda^{(t+1)}) - ((x^{(t)}, \lambda^{(t)})),$$

$$e_{\mathbf{z}}^{(t)} := \nabla_{\mathbf{z}} \phi_{c_t}(x^{(t)}, \boldsymbol{\lambda}^{(t)}, \mathbf{z}^{(t)}) - c_t(d_x^{(t)}, d_{\boldsymbol{\lambda}}^{(t)}) - c_t d_{\mathbf{z}}^{(t)} - \frac{c_t}{\beta} (\mathbf{z}^{(t+1)} - \mathbf{z}^{(t)}).$$

Next, we estimate  $\|e_{x, \boldsymbol{\lambda}}^{(t)}\|$ . Using the estimates in Yao et al. (2024), we have

$$\begin{aligned} \|e_{x, \boldsymbol{\lambda}}^{(t)}\| &\leq c_t L_{\phi_1} \left\| (x^{(t+1)}, \boldsymbol{\lambda}^{(t+1)}, \mathbf{z}^{(t+1)}) - (x^{(t)}, \boldsymbol{\lambda}^{(t)}, \mathbf{z}^{(t)}) \right\| \\ &\quad + \frac{c_t}{\alpha} \left\| (x^{(t+1)}, \boldsymbol{\lambda}^{(t+1)}) - (x^{(t)}, \boldsymbol{\lambda}^{(t)}) \right\| + c_t C_{\phi_1} \\ &\quad + \left\| (\boldsymbol{\theta}^{(t)}, \boldsymbol{\mu}^{(t)}) - (\boldsymbol{\theta}^*(x^{(t)}, \boldsymbol{\lambda}^{(t)}, \mathbf{z}^{(t)}), \boldsymbol{\mu}^*(x^{(t)}, \boldsymbol{\lambda}^{(t)}, \mathbf{z}^{(t)})) \right\|, \end{aligned}$$

where  $C_{\phi_1} := \sqrt{\max\{2(L_g + C_z L_g)^2, 2L_g^2\}}$ .

For  $\|e_{\mathbf{z}}^{(t)}\|$ , we have

$$\|e_{\mathbf{z}}^{(t)}\| \leq \left(\frac{c_t}{\beta} + \frac{c_t}{\gamma_2}\right) \|\mathbf{z}^{(t+1)} - \mathbf{z}^{(t)}\| + \frac{c_t}{\gamma_2} \|\boldsymbol{\mu}^{(t)} - \boldsymbol{\mu}^*(x^{(t)}, \boldsymbol{\lambda}^{(t)}, \mathbf{z}^{(t)})\|.$$

Thus,

$$\begin{aligned} R_t(x^{(t)}, \boldsymbol{\lambda}^{(t)}, \mathbf{z}^{(t)}) &\leq \left(\frac{c_t}{\beta} + \frac{c_t}{\gamma_2}\right) \|\mathbf{z}^{(t+1)} - \mathbf{z}^{(t)}\| + \frac{c_t}{\alpha} \left\| (x^{(t+1)}, \boldsymbol{\lambda}^{(t+1)}) - (x^{(t)}, \boldsymbol{\lambda}^{(t)}) \right\| \\ &\quad + c_t \left(C_{\phi_1} + \frac{1}{\gamma_2}\right) \left\| (\boldsymbol{\theta}^{(t)}, \boldsymbol{\mu}^{(t)}) - (\boldsymbol{\theta}^*(x^{(t)}, \boldsymbol{\lambda}^{(t)}, \mathbf{z}^{(t)}), \boldsymbol{\mu}^*(x^{(t)}, \boldsymbol{\lambda}^{(t)}, \mathbf{z}^{(t)})) \right\| \\ &\quad + c_t L_{\phi_1} \left\| (x^{(t+1)}, \boldsymbol{\lambda}^{(t+1)}, \mathbf{z}^{(t+1)}) - (x^{(t)}, \boldsymbol{\lambda}^{(t)}, \mathbf{z}^{(t)}) \right\|. \end{aligned}$$

Let  $\alpha_t \geq \underline{\alpha}$  and  $\beta_t \geq \underline{\beta}$  for some positive constants  $\underline{\alpha}$  and  $\underline{\beta}$ , we can show that there exists  $C_R > 0$  such that

$$\begin{aligned} \frac{1}{c_t^2} R_t^2(x^{(t)}, \boldsymbol{\lambda}^{(t)}, \mathbf{z}^{(t)}) &\leq C_R \left( \frac{1}{4\underline{\alpha}} \left\| (x^{(t+1)}, \boldsymbol{\lambda}^{(t+1)}) - (x^{(t)}, \boldsymbol{\lambda}^{(t)}) \right\|^2 + \frac{1}{4\underline{\beta}} \left\| \mathbf{z}^{(t+1)} - \mathbf{z}^{(t)} \right\|^2 \right. \\ &\quad \left. + \eta \rho_T C_{\boldsymbol{\theta}, \boldsymbol{\mu}} \left\| (\boldsymbol{\theta}^{(t)}, \boldsymbol{\mu}^{(t)}) - (\boldsymbol{\theta}^*(x^{(t)}, \boldsymbol{\lambda}^{(t)}, \mathbf{z}^{(t)}), \boldsymbol{\mu}^*(x^{(t)}, \boldsymbol{\lambda}^{(t)}, \mathbf{z}^{(t)})) \right\|^2 \right). \end{aligned}$$

This completes the proof.

## B.2 PROOF OF LEMMA 4.7

**Lemma 4.7.** Under Assumption 4.3 and 4.4, let  $\gamma_1 \in (0, 1/\rho_f)$ ,  $\gamma_2 > 0$  and pick  $\eta \in (0, \rho_T/L_b^2)$ , where  $L_b := \max\{L_f + L_g + C_z L_g + 1/\gamma_1, L_g + 1/\gamma_2\}$ . Then, the sequence generated by Algorithm 1 satisfies

$$\begin{aligned} &\left\| (\boldsymbol{\theta}^{(t+1)}, \boldsymbol{\mu}^{(t+1)}) - (\boldsymbol{\theta}^*(\boldsymbol{\omega}^{(t)}), \boldsymbol{\mu}^*(\boldsymbol{\omega}^{(t)})) \right\| \\ &\leq (1 - \eta \rho_T) \left\| (\boldsymbol{\theta}^{(t)}, \boldsymbol{\mu}^{(t)}) - (\boldsymbol{\theta}^*(\boldsymbol{\omega}^{(t)}), \boldsymbol{\mu}^*(\boldsymbol{\omega}^{(t)})) \right\|. \end{aligned}$$

*Proof:* The proof is similar to the proof of **Lemma A.5** in Yao et al. (2024).

## B.3 PROOF OF LEMMA 4.8

**Lemma 4.8.** Suppose the assumption of 4.2, 4.3 and 4.4 hold, and let  $\gamma_1 \in (0, 1/\rho_g)$ ,  $\gamma_2 > 0$ . Pick  $\eta \in (0, \rho_\gamma/L_B^2)$  with  $L_B := \max\{2L_g + C_z L_g + 1/\gamma_1, L_g + 1/\gamma_2\}$  then the sequence of  $(\boldsymbol{\omega}^{(t)})$  generated by Algorithm 1: MLM-CMOO satisfies

$$\phi_{c_t}(\boldsymbol{\omega}^{(t+1)}) \leq \phi_{c_t}(\boldsymbol{\omega}^{(t)}) - \left(\frac{1}{2\beta} - \frac{L_{v_z}}{2}\right) \|\mathbf{z}^{(t+1)} - \mathbf{z}^{(t)}\|^2$$

$$\begin{aligned}
& - \left( \frac{1}{2\alpha} - \frac{L_{\phi_k}}{2} - \frac{\beta L_{\theta, \mu}^2}{\gamma_2^2} \right) \left( \|x^{(t+1)} - x^{(t)}\|^2 + \|\boldsymbol{\lambda}^{(t+1)} - \boldsymbol{\lambda}^{(t)}\|^2 \right) \\
& + \frac{\alpha}{2} \left( 2(L_g + C_z L_g)^2 + \frac{1}{\gamma_1^2} \right) \|\boldsymbol{\theta}^{(t+1)} - \boldsymbol{\theta}^*(\boldsymbol{\omega}^{(t)})\|^2 \\
& + \left( \alpha L_g^2 + \frac{\beta}{\gamma_2^2} \right) \|\boldsymbol{\mu}^{(t+1)} - \boldsymbol{\mu}^*(\boldsymbol{\omega}^{(t)})\|^2,
\end{aligned}$$

where  $L_{\phi_t} := L_f/c_t + L_g + \rho_v$ .

*Proof:* Given Assumptions 4.2, 4.3, and 4.4 that  $\nabla F$  and  $\nabla g$  are  $L_F$ - and  $L_g$ -Lipschitz continuous on their domain, respectively, and applying **Lemma 5.7** in Beck (2017) and previous Lemmas, we obtain

$$\begin{aligned}
\phi_{c_t}(\boldsymbol{\omega}^{(t+1)}) & \leq \phi_{c_t}(\boldsymbol{\omega}^{(t)}) + \left\langle \nabla_{x, \boldsymbol{\lambda}} \phi_{c_t}(\boldsymbol{\omega}^{(t)}), (x^{(t+1)}, \boldsymbol{\lambda}^{(t+1)}) - (x^{(t)}, \boldsymbol{\lambda}^{(t)}) \right\rangle \\
& + \frac{L_{\phi_t}}{2} \left\| (x^{(t+1)}, \boldsymbol{\lambda}^{(t+1)}) - (x^{(t)}, \boldsymbol{\lambda}^{(t)}) \right\|^2,
\end{aligned}$$

with  $L_{\phi_t} := L_F/c_t + L_g + \rho_v$ . Based on the update rule of variable  $x^{(t)}, \boldsymbol{\lambda}^{(t)}$ , the convexity and the property of the proximal operator, we have

$$\left\langle (x^{(t)}, \boldsymbol{\lambda}^{(t)}) - \alpha(d_x^{(t)}, d_\lambda^{(t)}) - (x^{(t+1)}, \boldsymbol{\lambda}^{(t+1)}), (x^{(t)}, \boldsymbol{\lambda}^{(t)}) - (x^{(t+1)}, \boldsymbol{\lambda}^{(t+1)}) \right\rangle \leq 0,$$

thus, we have

$$\left\langle (d_x^{(t)}, d_\lambda^{(t)}), (x^{(t+1)}, \boldsymbol{\lambda}^{(t+1)}) - (x^{(t)}, \boldsymbol{\lambda}^{(t)}) \right\rangle \leq -\frac{1}{\alpha} \left\| (x^{(t+1)}, \boldsymbol{\lambda}^{(t+1)}) - (x^{(t)}, \boldsymbol{\lambda}^{(t)}) \right\|^2.$$

Considering the formula of  $\nabla_{x, \boldsymbol{\lambda}} L_{s, r}$  derived in **Lemma A.2** and the meanings of  $d_x^{(t)}, d_\lambda^{(t)}$  provided in the previous proof, we can obtain that

$$\begin{aligned}
& \left\| \nabla_{x, \boldsymbol{\lambda}} L_{s, r}(\boldsymbol{\omega}^{(t)}) - (d_x^{(t)}, d_\lambda^{(t)}) \right\|^2 \\
& = \left\| \nabla_x H(x^{(t)}, \boldsymbol{\theta}^*(\boldsymbol{\omega}^{(t)})) + \sum_{i=1}^n \mu_i^*(\boldsymbol{\omega}^{(t)}) \nabla_x g(x^{(t)}) - \nabla_x H(x^{(t)}, \boldsymbol{\theta}^{(t+1)}) - \sum_{i=1}^n \mu_i^{(t+1)} \nabla_x g(x^{(t)}) \right\|^2 \\
& + \frac{1}{\gamma_1^2} \left\| \boldsymbol{\theta}^{(t+1)} - \boldsymbol{\theta}^*(\boldsymbol{\omega}^{(t)}) \right\|^2 \\
& \leq 2 \left\| \nabla_x H(x^{(t)}, \boldsymbol{\theta}^*(\boldsymbol{\omega}^{(t)})) + \sum_{i=1}^n \mu_i^*(\boldsymbol{\omega}^{(t)}) \nabla_x g(x^{(t)}) - \nabla_x H(x^{(t)}, \boldsymbol{\theta}^*(\boldsymbol{\omega}^{(t)})) - \sum_{i=1}^n \mu_i^{(t+1)} \nabla_x g(x^{(t)}) \right\|^2 \\
& + 2 \left\| \nabla_x H(x^{(t)}, \boldsymbol{\theta}^*(\boldsymbol{\omega}^{(t)})) + \sum_{i=1}^n \mu_i^{(t+1)} \nabla_x g(x^{(t)}) - \nabla_x H(x^{(t)}, \boldsymbol{\theta}^{(t+1)}) - \sum_{i=1}^n \mu_i^{(t+1)} \nabla_x g(x^{(t)}) \right\|^2 \\
& + \frac{1}{\gamma_1^2} \left\| \boldsymbol{\theta}^{(t+1)} - \boldsymbol{\theta}^*(\boldsymbol{\omega}^{(t)}) \right\|^2 \\
& \leq \left( 2(L_f + C_z L_g + \frac{1}{\gamma_1^2}) \right) \left\| \boldsymbol{\theta}^{(t+1)} - \boldsymbol{\theta}^*(\boldsymbol{\omega}^{(t)}) \right\|^2 + 2L_g^2 \left\| \boldsymbol{\mu}^{(t+1)} - \boldsymbol{\mu}^*(\boldsymbol{\omega}^{(t)}) \right\|^2,
\end{aligned}$$

which yields

$$\begin{aligned}
& \left\langle \nabla_{x, \boldsymbol{\lambda}} L_{s, r}(\boldsymbol{\omega}^{(t)}) - (d_x^{(t)}, d_\lambda^{(t)}), \boldsymbol{\lambda}^{(t+1)}, (x^{(t+1)}, \boldsymbol{\lambda}^{(t+1)}) - (x^{(t)}, \boldsymbol{\lambda}^{(t)}) \right\rangle \\
& \leq \frac{\alpha}{2} \left( 2(L_f + C_z L_g + \frac{1}{\gamma_1^2}) \right) \left\| \boldsymbol{\theta}^{(t+1)} - \boldsymbol{\theta}^*(\boldsymbol{\omega}^{(t)}) \right\|^2 + \alpha L_g^2 \left\| \boldsymbol{\mu}^{(t+1)} - \boldsymbol{\mu}^*(\boldsymbol{\omega}^{(t)}) \right\|^2 \\
& + \frac{1}{2\alpha} \left\| (x^{(t+1)}, \boldsymbol{\lambda}^{(t+1)}) - (x^{(t)}, \boldsymbol{\lambda}^{(t)}) \right\|^2,
\end{aligned}$$



864 Combing with the above inequalities, we have

$$865 \phi_{c_t}(\boldsymbol{\omega}^{(t+1)}) \leq \phi_{c_t}(\boldsymbol{\omega}^{(t)}) + \left(\frac{1}{2\alpha} - \frac{L_{\phi_t}}{2}\right) \left\| (x^{(t+1)}, \boldsymbol{\lambda}^{(t+1)}) - (x^{(t)}, \boldsymbol{\lambda}^{(t)}) \right\|^2$$

$$866 + \frac{\alpha}{2} \left( 2(L_f + C_Z L_g + \frac{1}{\gamma_1^2}) \left\| \boldsymbol{\theta}^{(t+1)} - \boldsymbol{\theta}^*(\boldsymbol{\omega}^{(t)}) \right\|^2 + \alpha L_g^2 \left\| \boldsymbol{\mu}^{(t+1)} - \boldsymbol{\mu}^*(\boldsymbol{\omega}^{(t)}) \right\|^2 \right)$$

871 For variable  $\mathbf{z}$ , we have

$$872 \phi_{c_t}(\boldsymbol{\omega}^{(t+1)}) \leq \phi_{c_t}(\boldsymbol{\omega}^{(t)}) + \left\langle \nabla_{\mathbf{z}} \phi_{c_t}(\boldsymbol{\omega}^{(t)}), \mathbf{z}^{(t+1)} - \mathbf{z}^{(t)} \right\rangle + \frac{L_z}{2} \left\| \mathbf{z}^{(t+1)} - \mathbf{z}^{(t)} \right\|^2.$$

876 According to the property of the proximal gradient, we have

$$877 \left\langle d_{\mathbf{z}}^{(t)}, \mathbf{z}^{(t+1)} - \mathbf{z}^{(t)} \right\rangle \leq -\frac{1}{\beta} \left\| \mathbf{z}^{(t+1)} - \mathbf{z}^{(t)} \right\|^2$$

881 Thus, we have

$$882 \phi_{c_t}(\boldsymbol{\omega}^{(t+1)}) \leq \phi_{c_t}(\boldsymbol{\omega}^{(t)}) + \left\langle \nabla_{\mathbf{z}} \phi_{c_t}(\boldsymbol{\omega}^{(t)}) - d_{\mathbf{z}}^{(t)}, \mathbf{z}^{(t+1)} - \mathbf{z}^{(t)} \right\rangle + \left(\frac{L_z}{2} - \frac{1}{\beta}\right) \left\| \mathbf{z}^{(t+1)} - \mathbf{z}^{(t)} \right\|^2.$$

885 Based on the definition of  $d_{\mathbf{z}}^{(t)}$  provided in the previous section, we have

$$886 \left\| \boldsymbol{\omega}^{(t)} - d_{\mathbf{z}}^{(t)} \right\|^2 \leq \frac{1}{\gamma_2^2} \left\| \boldsymbol{\mu}^{(t+1)} - \boldsymbol{\mu}^*(x^{(t+1)}, \boldsymbol{\lambda}^{(t+1)}, \mathbf{z}^{(t)}) \right\|^2,$$

889 and

$$890 \left\langle \nabla_{\mathbf{z}} \phi_{c_t}(\boldsymbol{\omega}^{(t)}) - d_{\mathbf{z}}^{(t)}, \mathbf{z}^{(t+1)} - \mathbf{z}^{(t)} \right\rangle \leq \frac{\beta}{2\gamma_2^2} \left\| \boldsymbol{\mu}^{(t+1)} - \boldsymbol{\mu}^*(x^{(t+1)}, \boldsymbol{\lambda}^{(t+1)}, \mathbf{z}^{(t)}) \right\|^2 + \frac{1}{2\beta} \left\| \mathbf{z}^{(t+1)} - \mathbf{z}^{(t)} \right\|^2$$

894 The, for variable  $\mathbf{z}$ , we can get

$$895 \phi_{c_t}(\boldsymbol{\omega}^{(t+1)}) \leq \phi_{c_t}(\boldsymbol{\omega}^{(t)}) + \frac{\beta}{2\gamma_2^2} \left\| \boldsymbol{\mu}^{(t+1)} - \boldsymbol{\mu}^*(x^{(t+1)}, \boldsymbol{\lambda}^{(t+1)}, \mathbf{z}^{(t)}) \right\|^2 + \left(\frac{L_z}{2} - \frac{1}{2\beta}\right) \left\| \mathbf{z}^{(t+1)} - \mathbf{z}^{(t)} \right\|^2$$

$$896 \leq \phi_{c_t}(\boldsymbol{\omega}^{(t)}) + \frac{\beta}{2\gamma_2^2} \left\| \boldsymbol{\mu}^{(t+1)} - \boldsymbol{\mu}^*(x^{(t)}, \boldsymbol{\lambda}^{(t)}, \mathbf{z}^{(t)}) \right\|^2 + \left(\frac{L_z}{2} - \frac{1}{2\beta}\right) \left\| \mathbf{z}^{(t+1)} - \mathbf{z}^{(t)} \right\|^2$$

$$897 + \frac{\beta L_{\boldsymbol{\theta}, \boldsymbol{\mu}}^2}{2\gamma_2^2} \left\| (x^{(t+1)}, \boldsymbol{\lambda}^{(t+1)}) - (x^{(t)}, \boldsymbol{\lambda}^{(t)}) \right\|^2.$$

903 Combining the inequities for variable  $(x, \boldsymbol{\lambda})$  and  $\mathbf{z}$ , we can get Lemma 4.8.

#### 905 B.4 PROOF OF LEMMA 4.5

906 **Lemma 4.5** Under Assumptions 4.2, 4.3 and 4.4 hold, let  $\gamma_1 \in (0, 1/\rho_T)$ ,  $\gamma_2 > 0$ ,  $c_t \leq c_{t+1}$  and  $\eta_t \in (\eta, \rho_\gamma/L_B^2)$  with  $\eta > 0$ , then there exist constants  $c_\alpha, c_\beta > 0$  such that when  $0 < \alpha \leq c_\alpha$  and  $0 < \beta \leq c_\beta$ , the sequence of  $(x^{(t)}, \boldsymbol{\lambda}^{(t)}, \mathbf{z}^{(t)})$  generated by Algorithm 1: MLM-CMOO satisfies

$$910 V_{t+1} - V_t \leq -\frac{1}{4\alpha} \left\| x^{(t+1)} - x^{(t)} \right\|^2 - \frac{1}{4\alpha} \left\| \boldsymbol{\lambda}^{(t+1)} - \boldsymbol{\lambda}^{(t)} \right\|^2 - \frac{1}{4\beta} \left\| \mathbf{z}^{(t+1)} - \mathbf{z}^{(t)} \right\|^2$$

$$911 - \eta \rho_T C_{\boldsymbol{\theta}, \boldsymbol{\mu}} \left\| (\boldsymbol{\theta}^{(t)}, \boldsymbol{\mu}^{(t)}) - (\boldsymbol{\theta}^*(x^{(t)}, \boldsymbol{\lambda}^{(t)}, \mathbf{z}^{(t)}), \boldsymbol{\mu}^*(x^{(t)}, \boldsymbol{\lambda}^{(t)}, \mathbf{z}^{(t)})) \right\|^2.$$

915 *Proof:* From Lemma 4.8 and server aggregation rule, we have

$$916 \phi_{c_t}(\boldsymbol{\omega}^{(t)}) \leq \phi_{c_t}(\boldsymbol{\omega}^{(t)}) - \left(\frac{1}{2\beta} - \frac{L_{v_z}}{2}\right) \left\| \mathbf{z}^{(t+1)} - \mathbf{z}^{(t)} \right\|^2 \quad (8)$$

$$\begin{aligned}
& - \left( \frac{1}{2\alpha} - \frac{L_{\phi_k}}{2} - \frac{\beta L_{\theta, \mu}^2}{\gamma_2^2} \right) \left( \|x^{(t+1)} - x^{(t)}\|^2 + \|\boldsymbol{\lambda}^{(t+1)} - \boldsymbol{\lambda}^{(t)}\|^2 \right) \\
& + \frac{\alpha}{2} \left( 2(L_g + C_z L_g)^2 + \frac{1}{\gamma_1^2} \right) \|\boldsymbol{\theta}^{(t+1)} - \boldsymbol{\theta}^*(\boldsymbol{\omega}^{(t)})\|^2 \\
& + (\alpha L_g^2 + \frac{\beta}{\gamma_2^2}) \|\boldsymbol{\mu}^{(t+1)} - \boldsymbol{\mu}^*(\boldsymbol{\omega}^{(t)})\|^2.
\end{aligned}$$

Since  $c_{t+1} \geq c_t$ , we can infer that  $(F(x^{(t)}, \boldsymbol{\lambda}^{(t)}) - \underline{F})/c_{t+1} \leq (F(x^{(t)}, \boldsymbol{\lambda}^{(t)}) - \underline{F})/c_t$ . Combining with inequality equation 8 leads to

$$\begin{aligned}
V_{t+1} - V_t &= \phi_{c_{t+1}}(\boldsymbol{\omega}^{(t+1)}) - \phi_{c_t}(\boldsymbol{\omega}^{(t)}) + C_{\theta, \mu} \left\| (\boldsymbol{\theta}^{(t+1)}, \boldsymbol{\mu}^{(t+1)}) - (\boldsymbol{\theta}^*(\boldsymbol{\omega}^{(t+1)}), \boldsymbol{\mu}^*(\boldsymbol{\omega}^{(t+1)})) \right\|^2 \\
& - C_{\theta, \mu} \left\| (\boldsymbol{\theta}^{(t)}, \boldsymbol{\mu}^{(t)}) - (\boldsymbol{\theta}^*(\boldsymbol{\omega}^{(t)}), \boldsymbol{\mu}^*(\boldsymbol{\omega}^{(t)})) \right\|^2 \\
& \leq - \left( \frac{1}{2\alpha} - \frac{L_{\phi_t}}{2} - \frac{\beta L_{\theta, \mu}^2}{\gamma_2^2} \right) \left\| (x^{(t+1)}, \boldsymbol{\lambda}^{(t+1)}) - (x^{(t)}, \boldsymbol{\lambda}^{(t)}) \right\|^2 - \left( \frac{1}{2\beta} - \frac{L_{v_z}}{2} \right) \left\| \mathbf{z}^{(t+1)} - \mathbf{z}^{(t)} \right\|^2 \\
& + (\alpha L_g^2 + \frac{\beta}{\gamma_2^2}) \left\| \boldsymbol{\mu}^{(t+1)} - \boldsymbol{\mu}^*(\boldsymbol{\omega}^{(t)}) \right\|^2 + \frac{\alpha}{2} \left( 2(L_g + C_z L_g)^2 + \frac{1}{\gamma_1^2} \right) \left\| \boldsymbol{\theta}^{(t+1)} - \boldsymbol{\theta}^*(\boldsymbol{\omega}^{(t)}) \right\|^2 \\
& + C_{\theta, \mu} \left\| (\boldsymbol{\theta}^{(t+1)}, \boldsymbol{\mu}^{(t+1)}) - (\boldsymbol{\theta}^*(\boldsymbol{\omega}^{(t+1)}), \boldsymbol{\mu}^*(\boldsymbol{\omega}^{(t+1)})) \right\|^2 \\
& - C_{\theta, \mu} \left\| (\boldsymbol{\theta}^{(t)}, \boldsymbol{\mu}^{(t)}) - (\boldsymbol{\theta}^*(\boldsymbol{\omega}^{(t)}), \boldsymbol{\mu}^*(\boldsymbol{\omega}^{(t)})) \right\|^2 \\
& \leq - \left( \frac{1}{2\alpha} - \frac{L_{\phi_t}}{2} - \frac{\beta L_{\theta, \mu}^2}{\gamma_2^2} \right) \left\| (x^{(t+1)}, \boldsymbol{\lambda}^{(t+1)}) - (x^{(t)}, \boldsymbol{\lambda}^{(t)}) \right\|^2 - \left( \frac{1}{2\beta} - \frac{L_{v_z}}{2} \right) \left\| \mathbf{z}^{(t+1)} - \mathbf{z}^{(t)} \right\|^2 \\
& + C_{\theta, \lambda} \left\{ - \left\| (\boldsymbol{\theta}^{(t)}, \boldsymbol{\mu}^{(t)}) - (\boldsymbol{\theta}^*(\boldsymbol{\omega}^{(t)}), \boldsymbol{\mu}^*(\boldsymbol{\omega}^{(t)})) \right\|^2 + \left\| (\boldsymbol{\theta}^{(t)}, \boldsymbol{\mu}^{(t)}) - (\boldsymbol{\theta}^*(\boldsymbol{\omega}^{(t)}), \boldsymbol{\mu}^*(\boldsymbol{\omega}^{(t)})) \right\|^2 \right. \\
& \left. + 2 \max\{\alpha, \beta\} \left\| (\boldsymbol{\theta}^{(t+1)}, \boldsymbol{\mu}^{(t+1)}) - (\boldsymbol{\theta}^*(\boldsymbol{\omega}^{(t)}), \boldsymbol{\mu}^*(\boldsymbol{\omega}^{(t)})) \right\|^2 \right\},
\end{aligned}$$

where the last inequality follows from the fact that  $C_{\theta, \lambda} := \max\{(L_g + C_z L_g)^2 + 1/(2\gamma_1^2) + L_g^2, 1/\gamma_2^2\}$ .

Then, for the last 3 terms in the previous equation, we have

$$\begin{aligned}
& - \left\| (\boldsymbol{\theta}^{(t)}, \boldsymbol{\mu}^{(t)}) - (\boldsymbol{\theta}^*(\boldsymbol{\omega}^{(t)}), \boldsymbol{\mu}^*(\boldsymbol{\omega}^{(t)})) \right\|^2 + \left\| (\boldsymbol{\theta}^{(t)}, \boldsymbol{\mu}^{(t)}) - (\boldsymbol{\theta}^*(\boldsymbol{\omega}^{(t)}), \boldsymbol{\mu}^*(\boldsymbol{\omega}^{(t)})) \right\|^2 \\
& + 2\alpha \left\| (\boldsymbol{\theta}^{(t+1)}, \boldsymbol{\mu}^{(t+1)}) - (\boldsymbol{\theta}^*(\boldsymbol{\omega}^{(t)}), \boldsymbol{\mu}^*(\boldsymbol{\omega}^{(t)})) \right\|^2 \\
& \stackrel{a}{\leq} \left(1 + \frac{1}{\epsilon_t}\right) \left\| (\boldsymbol{\theta}^*(\boldsymbol{\omega}^{(t+1)}), \boldsymbol{\mu}^*(\boldsymbol{\omega}^{(t+1)})) - (\boldsymbol{\theta}^*(\boldsymbol{\omega}^{(t)}), \boldsymbol{\mu}^*(\boldsymbol{\omega}^{(t)})) \right\|^2 \\
& - \left\| (\boldsymbol{\theta}^{(t)}, \boldsymbol{\mu}^{(t)}) - (\boldsymbol{\theta}^*(\boldsymbol{\omega}^{(t)}), \boldsymbol{\mu}^*(\boldsymbol{\omega}^{(t)})) \right\|^2 \\
& + (1 + \epsilon_t + 2\alpha) \left\| (\boldsymbol{\theta}^{(t+1)}, \boldsymbol{\mu}^{(t+1)}) - (\boldsymbol{\theta}^*(\boldsymbol{\omega}^{(t)}), \boldsymbol{\mu}^*(\boldsymbol{\omega}^{(t)})) \right\|^2 \\
& \stackrel{b}{\leq} \left(1 + \frac{1}{\epsilon_t}\right) L_{\theta, \mu} \left\| \boldsymbol{\omega}^{(t+1)} - \boldsymbol{\omega}^{(t)} \right\|^2 - \left\| (\boldsymbol{\theta}^{(t)}, \boldsymbol{\mu}^{(t)}) - (\boldsymbol{\theta}^*(\boldsymbol{\omega}^{(t)}), \boldsymbol{\mu}^*(\boldsymbol{\omega}^{(t)})) \right\|^2 \\
& + (1 + \epsilon_t + 2\alpha)(1 - \eta\rho_T)^2 \left\| (\boldsymbol{\theta}^{(t+1)}, \boldsymbol{\mu}^{(t+1)}) - (\boldsymbol{\theta}^*(\boldsymbol{\omega}^{(t)}), \boldsymbol{\mu}^*(\boldsymbol{\omega}^{(t)})) \right\|^2 \\
& \leq \left(1 + \frac{2}{\eta\rho_T}\right) L_{\theta, \mu}^2 \left\| \boldsymbol{\omega}^{(t+1)} - \boldsymbol{\omega}^{(t)} \right\|^2 - \eta\rho_T \left\| (\boldsymbol{\theta}^{(t)}, \boldsymbol{\mu}^{(t)}) - (\boldsymbol{\theta}^*(\boldsymbol{\omega}^{(t)}), \boldsymbol{\mu}^*(\boldsymbol{\omega}^{(t)})) \right\|^2,
\end{aligned}$$

where  $a$  from Lemma A.5 and A.7 for  $\epsilon > 0$ , and  $b$  from setting  $\epsilon = \eta\rho_T/2$  and picking  $\alpha \leq \eta\rho_T/4$  where holds that  $(1 + \epsilon + 2\alpha)(1 - \eta\rho_T) \leq 1$ .

Similarly, we can show that when  $\beta \leq \eta\rho_T/4$ , it holds that

$$\begin{aligned} & - \left\| (\boldsymbol{\theta}^{(t)}, \boldsymbol{\mu}^{(t)}) - (\boldsymbol{\theta}^*(\boldsymbol{\omega}^{(t)}), \boldsymbol{\mu}^*(\boldsymbol{\omega}^{(t)})) \right\|^2 + \left\| (\boldsymbol{\theta}^{(t)}, \boldsymbol{\mu}^{(t)}) - (\boldsymbol{\theta}^*(\boldsymbol{\omega}^{(t)}), \boldsymbol{\mu}^*(\boldsymbol{\omega}^{(t)})) \right\|^2 \\ & \leq (1 + \frac{2}{\eta\rho_T}) L_{\boldsymbol{\theta}, \boldsymbol{\mu}}^2 \left\| \boldsymbol{\omega}^{(t+1)} - \boldsymbol{\omega}^{(t)} \right\|^2 - \eta\rho_T \left\| (\boldsymbol{\theta}^{(t)}, \boldsymbol{\mu}^{(t)}) - (\boldsymbol{\theta}^*(\boldsymbol{\omega}^{(t)}), \boldsymbol{\mu}^*(\boldsymbol{\omega}^{(t)})) \right\|^2. \end{aligned}$$

Combining the above inequities, we have

$$\begin{aligned} V_{t+1} - V_t & \leq - \left( \frac{1}{2\alpha} - \frac{L_{\phi_t}}{2} - \frac{\beta L_{\boldsymbol{\theta}, \boldsymbol{\mu}}^2}{\gamma_2^2} - (1 + \frac{2}{\eta\rho_T}) L_{\boldsymbol{\theta}, \boldsymbol{\mu}}^2 C_{\theta, \lambda} \right) \left\| (x^{(t+1)}, \boldsymbol{\lambda}^{(t+1)}) - (x^{(t)}, \boldsymbol{\lambda}^{(t)}) \right\|^2 \\ & \quad - \left( \frac{1}{2\beta} - \frac{L_{v_z}}{2} - (1 + \frac{2}{\eta\rho_T}) L_{\boldsymbol{\theta}, \boldsymbol{\mu}}^2 C_{\theta, \lambda} \right) \left\| \mathbf{z}^{(t+1)} - \mathbf{z}^{(t)} \right\|^2 \\ & \quad + \eta\rho_T C_{\theta, \lambda} \left\| (\boldsymbol{\theta}^{(t)}, \boldsymbol{\mu}^{(t)}) - (\boldsymbol{\theta}^*(\boldsymbol{\omega}^{(t)}), \boldsymbol{\mu}^*(\boldsymbol{\omega}^{(t)})) \right\|^2. \end{aligned}$$

When  $c_{t+1} \geq c_t$ ,  $\eta \geq \underline{\eta} > 0$ ,  $\alpha \leq \eta\rho_T/4$  and  $\beta \leq \eta\rho_T/4$ , then  $\frac{L_{\phi_t}}{2} + \frac{\beta L_{\boldsymbol{\theta}, \boldsymbol{\mu}}^2}{\gamma_2^2} + (1 + \frac{2}{\eta\rho_T}) L_{\boldsymbol{\theta}, \boldsymbol{\mu}}^2 C_{\theta, \mu} \leq \frac{L_{\phi_0}}{2} - \frac{\eta\rho_T L_{\boldsymbol{\theta}, \boldsymbol{\mu}}^2}{\gamma_2^2} - (1 + \frac{2}{\eta\rho_T}) L_{\boldsymbol{\theta}, \boldsymbol{\mu}}^2 C_{\theta, \mu} =: C_\alpha$  and  $\frac{L_{v_z}}{2} + (1 + \frac{2}{\eta\rho_T}) L_{\boldsymbol{\theta}, \boldsymbol{\mu}}^2 C_{\theta, \mu} \leq \frac{L_{v_z}}{2} + (1 + \frac{2}{\eta\rho_T}) L_{\boldsymbol{\theta}, \boldsymbol{\mu}}^2 C_{\theta, \mu} =: C_\beta$

Consequently, if  $C_\alpha, C_\beta > 0$  satisfies  $C_\alpha \leq \min \left\{ \frac{\eta\rho_T}{4}, \frac{1}{4C_\alpha} \right\}$  and  $C_\beta \leq \min \left\{ \frac{\eta\rho_T}{4}, \frac{1}{4C_\beta} \right\}$ , it holds that  $\frac{L_{\phi_t}}{2} + \frac{\beta L_{\boldsymbol{\theta}, \boldsymbol{\mu}}^2}{\gamma_2^2} + (1 + \frac{2}{\eta\rho_T}) L_{\boldsymbol{\theta}, \boldsymbol{\mu}}^2 C_{\theta, \mu} \geq \frac{1}{4\alpha}$  and  $\frac{L_{v_z}}{2} + (1 + \frac{2}{\eta\rho_T}) L_{\boldsymbol{\theta}, \boldsymbol{\mu}}^2 C_{\theta, \mu} \geq \frac{1}{4\beta}$

This completes the proof.