

# STOCHASTIC ADAPTIVE SEQUENTIAL BLACK-BOX OPTIMIZATION FOR DIFFUSION TARGETED GENERATION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Diffusion models have demonstrated great potential in generating high-quality content for images, natural language, protein domains, etc. However, how to perform user-preferred targeted generation via diffusion models with only black-box target scores of users remains challenging. To address this issue, we first formulate the fine-tuning of the inference phase of a pre-trained diffusion model as a sequential black-box optimization problem. Furthermore, we propose a novel stochastic adaptive sequential optimization algorithm to optimize cumulative black-box scores under unknown transition dynamics. Theoretically, we prove a  $O(\frac{d^2}{\sqrt{T}})$  convergence rate for cumulative convex functions without smooth and strongly convex assumptions. Empirically, we can naturally apply our algorithm for diffusion black-box targeted generation. Experimental results demonstrate the ability of our method to generate target-guided images with high target scores.

## 1 INTRODUCTION

Diffusion models have shown great success in generating high-quality content in various domains, such as image generation (Rombach et al., 2022; Ramesh et al., 2022), video generation (Ho et al., 2022), speech generation (Kim et al., 2022b; Kong et al., 2020), and natural language generation (Hu et al., 2023; He et al., 2023). Thanks to the super-promising power of the diffusion models, guided sampling via diffusion models to achieve desired properties recently emerged and shown fantastic potential in many applications, e.g., text-to-image generation (Kim et al., 2022a), image-to-image translation (Tumanyan et al., 2023), protein design (Lee et al., 2023; Gruver et al., 2023).

Despite the popularity and success of diffusion models, how to employ diffusion models to generate user-preferred content with black-box target scores **while avoiding re-training from scratch** is still challenging **and unexplored**. One direct idea is to treat this problem as a black-box optimization problem and employ black-box optimization techniques (Audet & Hare, 2017; Alarie et al., 2021; Doerr & Neumann, 2019) to perform the **fine-tuning of a pre-trained diffusion model** with only black-box target scores. However, naively applying black-box optimization methods to optimize diffusion model parameters faces high-dimensional optimization challenges, which are prohibitive to achieving a meaningful solution in a feasible time.

More importantly, current black-box optimization techniques, e.g., Bayesian optimization techniques (Srinivas et al., 2010; Gardner et al., 2017; Nayebi et al., 2019), Evolution strategies (ES) or stochastic zeroth-order optimization (Back et al., 1991; Hansen, 2006; Wierstra et al., 2014; Lyu & Tsang, 2021; Liu et al., 2018; Wang et al., 2018) and genetic algorithms (Srinivas & Patnaik, 1994; Mirjalili & Mirjalili, 2019), are designed for single objective without considering the dynamic nature **of sequential functions**. As a result, we can not directly apply them to diffusion models due to ignoring the sequential nature of the generation process of diffusion models.

In this paper, we dig into the transition dynamic of the inference of diffusion models. By leveraging the relationship between the inference of diffusion models and the Stochastic Differential Equation (SDE) solver (Song et al., 2020; Lu et al., 2022a), we naturally formulate the fine-tuning of the inference of diffusion models as a black-box sequential optimization problem.

To solve the black-box sequential optimization problem, we propose a novel stochastic adaptive black-box sequential optimization algorithm by explicitly handling the history trajectory dependency in the cumulative black-box target functions. Our method performs full covariance matrix adaptive

updates that can take advantage of second-order information to deal with ill-conditioned problems. Theoretically, we prove a  $O(\frac{d^2}{\sqrt{T}})$  convergence rate for convex functions without smooth and strongly convex assumptions. Thus, our method can handle non-smooth problems. Our contributions are listed as follows:

- We formulate the fine-tuning of the inference of the diffusion model as a black-box sequential optimization problem for black-box diffusion targeted generation.
- We proposed a novel stochastic adaptive black-box sequential optimization (SABSO) algorithm. Our SABSO can perform a full covariance matrix update to exploit the second-order information. Theoretically, we prove a  $O(\frac{d^2}{\sqrt{T}})$  convergence rate for convex functions without smooth and strongly convex assumptions. Thus, our theoretical analysis can handle non-smooth problems. The convergence analysis of full matrix adaptive black-box optimization for convex functions without the smooth and strongly convex assumptions is technically challenging. Technically, we add a  $\gamma_t$  enlargement term in the gradient update. This technique enables us to construct feasible solution sets of the adaptive update matrix during the whole algorithm running process to ensure convergence. To the best of our knowledge, our SABSO algorithm is the first full covariance matrix adaptive black-box optimization method that achieves a provable  $O(\frac{d^2}{\sqrt{T}})$  convergence rate for convex functions without smooth and strongly convex assumptions.
- Empirically, we can naturally apply our algorithm for diffusion black-box targeted generation. Experimental results demonstrate the ability of our method to generate target-guided images with high black-box target scores.

## 2 PRELIMINARY BACKGROUND

### 2.1 DIFFUSION MODEL SAMPLING

The sampling phase of the diffusion model from noise to image can be implemented by solving the stochastic differential equation (SDE) (Song et al., 2020; Kingma et al., 2021) as in Eq.(1):

$$d\mathbf{x}_s = [\hat{f}(s)\mathbf{x}_s + \frac{g(s)^2}{\sigma_s}\epsilon_\phi(\mathbf{x}_s, s)]ds + g(s)d\bar{\mathbf{w}}_s \quad (1)$$

where  $\bar{\mathbf{w}}_s$  is the reverse-time Wiener process,  $s$  denotes time changing from  $K$  to 0,  $\epsilon_\phi(\mathbf{x}_s, s)$  denotes the diffusion model noise prediction with input  $\mathbf{x}_s$  and time  $s$ . And  $\sigma_s$  denotes the standard deviation of the diffusion noise scheme at time  $s$ . And  $\hat{f}(s) := \frac{d \log \alpha_s}{ds}$ , where  $\alpha_s$  denotes the scaling parameter scheme in the diffusion model. And  $g(s)^2 := \frac{d\sigma_s^2}{ds} - 2\hat{f}(s)\sigma_s^2$  (Kingma et al., 2021).

Recently, Lu et al. (2022a;b) proposed a DPM solver for solving the diffusion SDE with a small number of samples. The first-order SDE DPM solver is given in Eq.(2).

$$\mathbf{x}_s = \frac{\alpha_s}{\alpha_{s'}}\mathbf{x}_{s'} - 2\sigma_s(e^h - 1)\epsilon_\phi(\mathbf{x}_{s'}, s') + \sigma_s\sqrt{e^{2h} - 1}\mathbf{z} \quad (2)$$

where  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ ,  $h = \lambda_s - \lambda_{s'}$ , and  $\lambda_s = \log(\alpha_s/\sigma_s)$  and  $\lambda_{s'} = \log(\alpha_{s'}/\sigma_{s'})$ . And  $\alpha_s, \alpha_{s'}$  denote the scaling parameter at step  $s$  and  $s'$  in diffusion model, respectively.

### 2.2 BLACK-BOX OPTIMIZATION

Given a proper function  $f(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}$  such that  $f(\mathbf{x}) > -\infty$ , black-box optimization is to minimize  $f(\mathbf{x})$  by using function queries only. Instead of optimizing the original problem directly, ES or stochastic zeroth-order optimization methods optimize a **relaxation** of the problem  $J(\theta) := \mathbb{E}_{p(\mathbf{x};\theta)}[f(\mathbf{x})]$  w.r.t. the parameter  $\theta$  of the sampling distribution of the relaxed problem.

**Evolution strategies** (Rechenberg & Eigen, 1973; Nesterov & Spokoiny, 2017; Liu et al., 2018) employ a Gaussian distribution  $\mathcal{N}(\boldsymbol{\mu}, \sigma^2 \mathbf{I})$  with a fixed variance for candidate sampling. The approximate gradient descent update is given as

$$\boldsymbol{\mu}_{t+1} = \boldsymbol{\mu}_t - \frac{\beta}{N\sigma} \sum_{i=1}^N \epsilon_i f(\boldsymbol{\mu}_t + \sigma \epsilon_i), \quad (3)$$

where  $\epsilon_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  and  $\beta$  denotes the step-size, and  $\boldsymbol{\mu}_t$  denotes the mean parameter of the Gaussian distribution for candidate sampling at  $t^{\text{th}}$  black-box optimization iteration.

The ES methods only perform the first-order approximate gradient update, the convergence speed is limited. Wierstra et al. (2014) proposed the natural evolution strategies (NES), which perform the approximate natural gradient update, in which a Gaussian distribution  $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  is employed for sampling. Besides the updating of parameter  $\boldsymbol{\mu}$ , the covariance matrix  $\boldsymbol{\Sigma}$  is also updated. Lyu & Tsang (2021) proposed an implicit natural gradient optimizer (INGO) for black-box optimization, which provides an alternative way to compute the natural gradient update. In INGO update rule, the inverse covariance matrix  $\boldsymbol{\Sigma}^{-1}$  is updated instead of the covariance matrix  $\boldsymbol{\Sigma}$ . Moreover, CMAES (Hansen, 2006) provides a more sophisticated update rule and performs well on a wide range of black-box optimization problems. Despite the success of these methods, all these methods ignore the dynamic nature of the target function.

The recent work (Krishnamoorthy et al., 2023) introduces Denoising Diffusion Optimization Models (DDOM) for solving offline black-box optimization tasks using diffusion models. This method can also be naturally extended to black-box targeted generation tasks. The DDOM relies on an offline conditional model trained with reweighted data sampling. The generation is performed conditioned on a high target score. The pre-collected data set has a crucial influence on DDOM generation.

### 3 NOTATION AND SYMBOLS

Denote  $\|\cdot\|_2$  and  $\|\cdot\|_F$  as the spectral norm and Frobenius norm for matrices, respectively. Define  $\text{tr}(\cdot)$  as the trace operation for matrix. Notation  $\|\cdot\|_2$  will also denote  $l_2$ -norm for vectors. Symbol  $\langle \cdot, \cdot \rangle$  denotes inner product under  $l_2$ -norm for vectors and inner product under Frobenius norm for matrices. For a positive semi-definite matrix  $C$ , define  $\|\mathbf{x}\|_C := \sqrt{\langle \mathbf{x}, C\mathbf{x} \rangle}$ . Denote  $\mathcal{S}^+$  as the set of positive semi-definite matrices. Denote  $\Sigma^{\frac{1}{2}}$  as the symmetric positive semi-definite matrix such that  $\Sigma = \Sigma^{\frac{1}{2}}\Sigma^{\frac{1}{2}}$  for  $\Sigma \in \mathcal{S}^+$ . Denote  $\bar{\mathbf{x}}_k = [\mathbf{x}_1^\top, \dots, \mathbf{x}_k^\top]^\top \in \mathcal{R}^{kd}$ , where  $\mathbf{x}_i \in \mathcal{R}^d$ ,  $d$  denotes the dimension of the data. Denote  $\bar{\boldsymbol{\mu}}_k = [\boldsymbol{\mu}_1^\top, \dots, \boldsymbol{\mu}_k^\top]^\top$  and  $\bar{\boldsymbol{\Sigma}}_k = \text{diag}(\boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_k) \in \mathcal{R}^{kd \times kd}$  as the mean and diagonal block-wise covariance matrix for Gaussian distribution, respectively. Denote  $\theta_k := \{\bar{\boldsymbol{\mu}}_k, \bar{\boldsymbol{\Sigma}}_k\}$  as the parameter of the distribution for candidate sampling in black-box optimization and  $\theta_k := \{\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}$  as its  $k^{\text{th}}$  component. Denote  $\bar{\theta}_k^t := \{\bar{\boldsymbol{\mu}}_k^t, \bar{\boldsymbol{\Sigma}}_k^t\}$  as the parameter at  $t^{\text{th}}$  iteration.

## 4 METHOD

### 4.1 BLACK-BOX OPTIMIZATION FOR DIFFUSION TARGETED GENERATION

Diffusion model sampling can be implemented by solving Diffusion SDEs (Song et al., 2020). From SDE solvers in Eq.(2), the inference phase of the diffusion SDE model can be rewritten as

$$\mathbf{x}_k = \hat{\boldsymbol{\mu}}_\phi(\mathbf{x}_{k-1}, k-1) + \tilde{\sigma}_k \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (4)$$

for  $k \in \{1, \dots, K\}$ ,  $k = K-s$ , and  $\mathbf{x}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ ,  $\tilde{\sigma}_k$  denotes the coefficient of the SDE solver. And  $\hat{\boldsymbol{\mu}}_\phi(\mathbf{x}_{k-1}, k-1)$  is the prediction of a pre-trained diffusion model with fixed parameter  $\phi$ . The concrete  $\hat{\boldsymbol{\mu}}_\phi(\mathbf{x}_{k-1}, k-1)$  depends on the solver type.

To perform guided sampling from the diffusion model with the black-box target score function  $F(\mathbf{x})$  (e.g., CLIP model evaluates on the input image  $\mathbf{x}$ ), we generalize the inference of SDE solver as

$$\mathbf{x}_k = \hat{\boldsymbol{\mu}}_\phi(\mathbf{x}_{k-1}, k-1) + \tilde{\sigma}_k \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (5)$$

for  $k \in \{1, \dots, K\}$ , and fine-tuning the parameter  $\theta_k := \{\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}$  for  $k \in \{1, \dots, K\}$ . This naturally leads to a black-box optimization of the cumulative target score as

$$\tilde{J}(\bar{\theta}_K) := \mathbb{E}_{\mathbf{x}_{0:K} \sim p_{\bar{\theta}_K}} \left[ \sum_{k=1}^K F(\mathbf{x}_k) \right] \quad (6)$$

where  $\mathbf{x}_k$  transition as  $\mathbf{x}_k = \hat{\boldsymbol{\mu}}_\phi(\mathbf{x}_{k-1}, k-1) + \tilde{\sigma}_k \boldsymbol{\epsilon}_k$  with  $\boldsymbol{\epsilon}_k \sim \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ , and  $\bar{\theta}_K := \{\bar{\boldsymbol{\mu}}_K, \bar{\boldsymbol{\Sigma}}_K\}$ ,  $F(\mathbf{x}_k)$  denotes the score function evaluated on input data  $\mathbf{x}_k$ . It is worth to remark that  $\mathbf{x}_k$  depends on the output state  $\mathbf{x}_{k-1}$ . Problem (6) is essentially a sequential optimization problem. The strategy of choosing  $\theta_k := \{\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}$  depends on the strategy of choosing  $\{\theta_{k-1}, \theta_{k-2}, \dots, \theta_1\}$  in a nested manner.

Instead of directly optimizing the objective Eq.(6), we optimize an augmented objective Eq.(7) in the noise space.

$$J(\bar{\theta}_K) := \mathbb{E}_{\mathbf{x}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \mathbb{E}_{\bar{\epsilon}_K \sim \mathcal{N}(\bar{\boldsymbol{\mu}}_K, \bar{\boldsymbol{\Sigma}}_K)} \left[ \sum_{k=1}^K f_k(\bar{\epsilon}_k) \right] \quad (7)$$

where  $\bar{\epsilon}_k = [\epsilon_1^\top, \dots, \epsilon_k^\top]^\top$  for  $k \in \{1, \dots, K\}$ , and  $f_k(\bar{\epsilon}_k) = F_k(\mathbf{x}_k)$  for  $\mathbf{x}_k$  transitioned as  $\mathbf{x}_k = \hat{\boldsymbol{\mu}}_\phi(\mathbf{x}_{k-1}, k-1) + \tilde{\sigma}_k \epsilon_k$  with the trajectory  $[\epsilon_1, \dots, \epsilon_k]$ .  $F_k(\mathbf{x}_k)$  performs a deterministic sampling process  $\mathbf{x}_{k+1} = \hat{\boldsymbol{\mu}}_\phi(\mathbf{x}_k, k) + \tilde{\sigma}_{k+1} \boldsymbol{\mu}_{k+1}$  for the future steps to achieve a predicted  $\mathbf{x}_K$  and evaluate on the predicted  $\mathbf{x}_K$ .

The objective Eq.(7) can be rewritten as

$$J(\bar{\theta}_K) = \sum_{k=1}^K J_k(\bar{\theta}_k) = \sum_{k=1}^K \mathbb{E}_{\mathbf{x}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \mathbb{E}_{\bar{\epsilon}_k \sim \mathcal{N}(\bar{\boldsymbol{\mu}}_k, \bar{\boldsymbol{\Sigma}}_k)} [f_k(\bar{\epsilon}_k)] \quad (8)$$

where  $J_k(\bar{\theta}_k) = \mathbb{E}_{\mathbf{x}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \mathbb{E}_{\bar{\epsilon}_k \sim \mathcal{N}(\bar{\boldsymbol{\mu}}_k, \bar{\boldsymbol{\Sigma}}_k)} [f_k(\bar{\epsilon}_k)]$ .

## 4.2 CLOSED-FORM UPDATE RULE

In this section, we derive the update rule of the parameter to optimize Eq.(7). Without loss of generality, we assume minimization in this paper.

Given a parameter  $\bar{\theta}_K^t$  at  $t^{th}$  iteration, we aim to find a new parameter  $\bar{\theta}_K^{t+1}$  by minimizing the objective difference as

$$\min_{\bar{\theta}_K} J(\bar{\theta}_K) - J(\bar{\theta}_K^t) \quad (9)$$

However, it is challenging to solve the optimization (9) directly. We thus optimize an approximation by first order Taylor expansion. We add a KL-divergence regularization to ensure  $q_{\bar{\theta}_K}$  and  $q_{\bar{\theta}_K^t}$  close, thus to keep the approximation accurate. The new optimization problem is given as

$$\min_{\bar{\theta}_K} J(\bar{\theta}_K) - J(\bar{\theta}_K^t) + \text{KL}(q_{\bar{\theta}_K} | q_{\bar{\theta}_K^t}) = \sum_{k=1}^K J_k(\bar{\theta}_k) - J_k(\bar{\theta}_k^t) + \text{KL}(q_{\bar{\theta}_K} | q_{\bar{\theta}_K^t}) \quad (10)$$

$$\approx \sum_{k=1}^K \left\langle \bar{\theta}_k - \bar{\theta}_k^t, \beta_k \nabla_{\bar{\theta}_k^t} J_k(\bar{\theta}_k^t) \right\rangle + \text{KL}(q_{\bar{\theta}_K} | q_{\bar{\theta}_K^t}) \quad (11)$$

where  $q_{\bar{\theta}_K} := \mathcal{N}(\bar{\boldsymbol{\mu}}_K, \bar{\boldsymbol{\Sigma}}_K)$  and  $q_{\bar{\theta}_K^t} := \mathcal{N}(\bar{\boldsymbol{\mu}}_K^t, \bar{\boldsymbol{\Sigma}}_K^t)$ .

Note that the problem (11) is convex w.r.t.  $\bar{\theta}_K := \{\bar{\boldsymbol{\mu}}_K, \bar{\boldsymbol{\Sigma}}_K\}$ , by setting the derivative to zero, we can achieve a closed-form update as

$$\boldsymbol{\mu}_k^{t+1} = \boldsymbol{\mu}_k^t - \sum_{i=k}^K \beta_i \mathbb{E}_{\mathbf{x}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \mathbb{E}_{\mathcal{N}(\bar{\boldsymbol{\mu}}_i, \bar{\boldsymbol{\Sigma}}_i)} [(\epsilon_i - \boldsymbol{\mu}_i^t) f_i(\bar{\epsilon}_i)] \quad (12)$$

$$\boldsymbol{\Sigma}_k^{t+1} = \boldsymbol{\Sigma}_k^{t-1} + \sum_{i=k}^K \beta_i \mathbb{E}_{\mathbf{x}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \mathbb{E}_{\mathcal{N}(\bar{\boldsymbol{\mu}}_i, \bar{\boldsymbol{\Sigma}}_i)} [(\boldsymbol{\Sigma}_i^{t-1} (\epsilon_i - \boldsymbol{\mu}_i^t) (\epsilon_i - \boldsymbol{\mu}_i^t)^\top \boldsymbol{\Sigma}_i^{t-1} - \boldsymbol{\Sigma}_i^{t-1}) f_i(\bar{\epsilon}_i)] \quad (13)$$

Detailed derivation can be found in Appendix B.

To compute the update in Eq.(12) and Eq.(13), we perform Monte Carlo sampling by taking  $N$  i.i.d. sequence  $\{\mathbf{x}_0^j, \epsilon_1^j, \dots, \epsilon_K^j\}$  for  $j \in \{1, \dots, N\}$ , where  $\mathbf{x}_0^j \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  and  $\epsilon_k^j \sim \mathcal{N}(\boldsymbol{\mu}_k^t, \boldsymbol{\Sigma}_k^t)$  for  $k \in \{1, \dots, K\}$ . This leads to unbiased estimators of the RHS of Eq.(12) and Eq.(13) as follows:

$$\boldsymbol{\mu}_k^{t+1} = \boldsymbol{\mu}_k^t - \frac{1}{N} \sum_{i=k}^K \beta_i \sum_{j=1}^N [(\epsilon_i^j - \boldsymbol{\mu}_i^t) (f_i(\bar{\epsilon}_i^j) - f_i(\bar{\boldsymbol{\mu}}_i^t))] \quad (14)$$

$$\boldsymbol{\Sigma}_k^{t+1} = \boldsymbol{\Sigma}_k^{t-1} + \frac{1}{N} \sum_{i=k}^K \beta_i \sum_{j=1}^N \left( \boldsymbol{\Sigma}_i^{t-1} (\epsilon_i^j - \boldsymbol{\mu}_i^t) (\epsilon_i^j - \boldsymbol{\mu}_i^t)^\top \boldsymbol{\Sigma}_i^{t-1} - \boldsymbol{\Sigma}_i^{t-1} \right) (f_i(\bar{\epsilon}_i^j) - f_i(\bar{\boldsymbol{\mu}}_i^t)) \quad (15)$$

**Algorithm 1** BDTG

---

**Input:** Number of Batch Size  $N$ , step-size  $\beta$ , a pre-trained diffusion model  $\widehat{\mu}_\phi(\mathbf{x}_k, k)$ , number of sampling step  $K$ , SDE solver coefficient  $\tilde{\sigma}_k$  for  $k \in \{1, \dots, K\}$ . **Number of total iteration**  $T$ .  
**Initialization:** Initialize  $\boldsymbol{\mu}_k^1 = \mathbf{0}, \boldsymbol{\Sigma}_k^1 = \mathbf{I}$ , and set  $\beta_k = \beta \tilde{\sigma}_k$  for  $k \in \{1, \dots, K\}$ .  
**for**  $t = 1, \dots, T$  **do**  
  Take i.i.d. samples  $\mathbf{x}_0^1, \dots, \mathbf{x}_0^N \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ .  
  **for**  $k = 1, \dots, K$  **do**  
    Take i.i.d. samples  $\boldsymbol{\epsilon}_k^1, \dots, \boldsymbol{\epsilon}_k^N \sim \mathcal{N}(\boldsymbol{\mu}_k^t, \boldsymbol{\Sigma}_k^t)$   
    Set  $\mathbf{x}_k^j = \widehat{\mu}_\phi(\mathbf{x}_{k-1}^j, k-1) + \tilde{\sigma}_k \boldsymbol{\epsilon}_k^j$  for all  $j \in \{1, \dots, N\}$ .  
    Query black-box target function score  $f_k(\bar{\boldsymbol{\epsilon}}_k^1) = F_k(\mathbf{x}_k^1), \dots, f_k(\bar{\boldsymbol{\epsilon}}_k^N) = F_k(\mathbf{x}_k^N)$ .  
  **end for**  
  Update  $\boldsymbol{\mu}_k^{t+1}$  for all  $k \in \{1, \dots, K\}$  using Eq. (16)  
  Update  $\boldsymbol{\Sigma}_k^{t+1}$  for all  $k \in \{1, \dots, K\}$  using Eq. (17)  
**end for**

---

where the offset term  $f_i(\boldsymbol{\mu}_i^t)$  is employed to reduce variance.

In practice, to avoid the numeric scale problem, we normalize the score by  $h(f_i(\bar{\boldsymbol{\epsilon}}_i^j)) = \frac{f_i(\bar{\boldsymbol{\epsilon}}_i^j) - \widehat{\mu}_i}{\widehat{\sigma}_i}$ , where  $\widehat{\mu}_i$  and  $\widehat{\sigma}_i$  denote mean and standard deviation of function values  $[f_i(\bar{\boldsymbol{\epsilon}}_i^1), \dots, f_i(\bar{\boldsymbol{\epsilon}}_i^N)]$ . We thus obtain the following update rule for practical updates.

$$\boldsymbol{\mu}_k^{t+1} = \boldsymbol{\mu}_k^t - \frac{1}{N} \sum_{i=k}^K \beta_i \sum_{j=1}^N [(\boldsymbol{\epsilon}_i^j - \boldsymbol{\mu}_i^t) \frac{f_i(\bar{\boldsymbol{\epsilon}}_i^j) - \widehat{\mu}_i}{\widehat{\sigma}_i}] \quad (16)$$

$$\boldsymbol{\Sigma}_k^{t+1} = \boldsymbol{\Sigma}_k^{t-1} + \frac{1}{N} \sum_{i=k}^K \beta_i \sum_{j=1}^N \left( \boldsymbol{\Sigma}_i^{t-1} (\boldsymbol{\epsilon}_i^j - \boldsymbol{\mu}_i^t) (\boldsymbol{\epsilon}_i^j - \boldsymbol{\mu}_i^t)^\top \boldsymbol{\Sigma}_i^{t-1} \right) \frac{f_i(\bar{\boldsymbol{\epsilon}}_i^j) - \widehat{\mu}_i}{\widehat{\sigma}_i} \quad (17)$$

Note that  $\boldsymbol{\epsilon}_i^j = \boldsymbol{\mu}_i^t + \boldsymbol{\Sigma}_i^{t \frac{1}{2}} \mathbf{z}_i^j$  for  $\mathbf{z}_i^j \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , Eq. (17) can be rewritten as

$$\boldsymbol{\Sigma}_k^{t+1} = \boldsymbol{\Sigma}_k^{t-1} (\mathbf{I} + \beta \mathbf{H}_k^t) \boldsymbol{\Sigma}_k^{t-1} \quad (18)$$

where  $\mathbf{H}_k^t$  is constructed as Eq. (19).

$$\mathbf{H}_k^t = \frac{1}{N} \sum_{i=k}^K \frac{\beta_i}{\beta} \sum_{j=1}^N \mathbf{z}_i^j \mathbf{z}_i^{j \top} \frac{f_i(\bar{\boldsymbol{\epsilon}}_i^j) - \widehat{\mu}_i}{\widehat{\sigma}_i} \quad (19)$$

The property of  $\mathbf{H}_k^t$  has a crucial impact on the convergence speed.

We summarize our algorithm for black-box diffusion target generation algorithm (BDTG) into Algorithm 1. In Algorithm 1, the user preference can be incorporated via the black-box target score. In addition,  $\tilde{\sigma}_k$  is the SDE solver coefficient. For example, when DPM solver (Lu et al., 2022a) is employed,  $\tilde{\sigma}_k = \sigma_k \sqrt{e^{2h} - 1}$ . More details about different solvers can be found in (Lu et al., 2022b).

## 5 CONVERGENCE ANALYSIS

In this section, we provide the convergence analysis of our general algorithm framework. Without loss of generality, we focus on minimizing the following sequential optimization problem:

$$\widehat{F}(\bar{\mathbf{x}}_K) = \sum_{k=1}^K f_k(\bar{\mathbf{x}}_k) \quad (20)$$

**Remark:** The black-box function  $f_k(\bar{\mathbf{x}}_k)$  with  $\bar{\mathbf{x}}_k = [\mathbf{x}_1^\top, \dots, \mathbf{x}_k^\top]^\top$  explicitly shows the dependency of the whole history trajectory  $[\mathbf{x}_1, \dots, \mathbf{x}_k]$ . The sequential black-box optimization in Eq.(20) is general enough to include many interesting scenarios as special cases. One particular interesting problem is  $f_k(\bar{\mathbf{x}}_k) = F(\mathbf{y}_k)$  where  $\mathbf{y}_k$  is obtained by an unknown transition dynamic  $\mathbf{y}_k = Q(\mathbf{y}_{k-1}, \mathbf{x}_k)$ .

**Algorithm 2** SABSO Framework

---

**Input:** Batch-size  $N$ . **Parameter**  $\beta_t, \alpha_t, \gamma_t,$  and  $\omega_t$ . **Number of total iteration**  $T$  and the step  $K$ .  
**Initialization:** Initialize  $\boldsymbol{\mu}_k^1 = \mathbf{0}, \boldsymbol{\Sigma}_k^1 = \tau \mathbf{I}$  for  $k \in \{1, \dots, K\}$ .  
**for**  $t = 1, \dots, T$  **do**  
  **for**  $k = 1, \dots, K$  **do**  
    Take i.i.d. samples  $\mathbf{z}_k^1, \dots, \mathbf{z}_k^N \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$   
    Set  $\mathbf{x}_k^j = \boldsymbol{\mu}_k^t + \boldsymbol{\Sigma}_k^{t \frac{1}{2}} \mathbf{z}_k^j$  for  $j \in \{1, \dots, N\}$   
    Query black-box objective function value  $f_k(\bar{\mathbf{x}}_k^1), \dots, f_k(\bar{\mathbf{x}}_k^N)$ .  
    Construct unbiased estimator  $\hat{g}_{k1}, \dots, \hat{g}_{kN}$  as Eq. (22).  
    Compute  $\hat{\boldsymbol{\mu}}_k, \hat{\boldsymbol{\sigma}}_k$  as the mean and std of  $\{f_k(\bar{\mathbf{x}}_k^1), \dots, f_k(\bar{\mathbf{x}}_k^N)\}$ .  
  **end for**  
  **for**  $k = 1, \dots, K$  **do**  
    Construct  $\mathbf{H}_k^t = c_1 \frac{1}{N} \sum_{i=k}^K \sum_{j=1}^N \mathbf{z}_i^j \mathbf{z}_i^{j \top} \frac{f_i(\bar{\mathbf{x}}_i^j) - \hat{\boldsymbol{\mu}}_i}{\hat{\boldsymbol{\sigma}}_i} + c_2 \mathbf{I}$  with constants  $c_1 > 0, c_2 > 0$   
    such that  $\mathbf{H}_k^t \preceq \frac{1}{\alpha_t} (\frac{\beta_{t+1}}{\beta_t} - \omega_t) \mathbf{I} + \frac{\beta_{t+1} \gamma_t}{\alpha_t} \boldsymbol{\Sigma}_k^t$  and  $\nu \mathbf{I} \preceq \hat{\mathbf{G}}_k^t = \boldsymbol{\Sigma}_k^{t - \frac{1}{2}} \mathbf{H}_k^t \boldsymbol{\Sigma}_k^{t - \frac{1}{2}}$   
    Set  $\hat{\mathbf{G}}_k^t = \boldsymbol{\Sigma}_k^{t - \frac{1}{2}} \mathbf{H}_k^t \boldsymbol{\Sigma}_k^{t - \frac{1}{2}}$ .  
    Set  $\boldsymbol{\mu}_k^{t+1} = \boldsymbol{\mu}_k^t - \beta_t \boldsymbol{\Sigma}_k^t (\gamma_t \boldsymbol{\mu}_k^t + (\sum_{i=k}^K \hat{g}_{ik}))$   
    Set  $\boldsymbol{\Sigma}_k^{t+1} = \omega_t \boldsymbol{\Sigma}_k^{t-1} + \alpha_t \hat{\mathbf{G}}_k^t$ .  
  **end for**  
**end for**

---

Instead of directly optimizing problem (20), we optimize an auxiliary problem (21) as

$$J(\bar{\boldsymbol{\mu}}_K, \bar{\boldsymbol{\Sigma}}_K) = \sum_{i=1}^K J_i(\bar{\boldsymbol{\mu}}_i, \bar{\boldsymbol{\Sigma}}_i) = \sum_{i=1}^K \mathbb{E}_{\bar{\mathbf{x}}_i \sim \mathcal{N}(\bar{\boldsymbol{\mu}}_i, \bar{\boldsymbol{\Sigma}}_i)} [f_i(\bar{\mathbf{x}}_i)] \quad (21)$$

where  $J_i(\bar{\boldsymbol{\mu}}_i, \bar{\boldsymbol{\Sigma}}_i) = \mathbb{E}_{\bar{\mathbf{x}}_i \sim \mathcal{N}(\bar{\boldsymbol{\mu}}_i, \bar{\boldsymbol{\Sigma}}_i)} [f_i(\bar{\mathbf{x}}_i)]$ .

Denote gradient estimator  $\hat{g}_{ik}^t$  for the  $i^{\text{th}}$  objective w.r.t. the  $k^{\text{th}}$  component  $\boldsymbol{\mu}_k$  at  $t^{\text{th}}$  iteration as

$$\hat{g}_{ik}^t = \frac{1}{N} \sum_{j=1}^N \hat{g}_{ik}^{tj} = \frac{1}{N} \sum_{j=1}^N \boldsymbol{\Sigma}_k^{t - \frac{1}{2}} \mathbf{z}_k^j (f_i(\bar{\boldsymbol{\mu}}_i^t + \bar{\boldsymbol{\Sigma}}_i^{t \frac{1}{2}} \mathbf{z}_i^j) - f_i(\bar{\boldsymbol{\mu}}_i^t)), \quad (22)$$

where  $\hat{g}_{ik}^{tj}$  is the gradient estimator using  $j^{\text{th}}$  i.i.d. sample:

$$\hat{g}_{ik}^{tj} = \boldsymbol{\Sigma}_k^{t - \frac{1}{2}} \mathbf{z}_k^j (f_i(\bar{\boldsymbol{\mu}}_i^t + \bar{\boldsymbol{\Sigma}}_i^{t \frac{1}{2}} \mathbf{z}_i^j) - f_i(\bar{\boldsymbol{\mu}}_i^t)), \quad (23)$$

where  $\mathbf{z}_1^j, \dots, \mathbf{z}_i^j \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$  and  $\bar{\mathbf{z}}_i^j = [\mathbf{z}_1^\top, \dots, \mathbf{z}_i^\top]^\top$  for  $i \geq k$ .

We show our Stochastic Adaptive Black-box Sequential Optimization algorithm (SABSO) in Algorithm 2. Our SABSO can perform full matrix updates to take advantage of second-order information. We add a  $\gamma_t \boldsymbol{\mu}_k^t$  term in the update step of  $\boldsymbol{\mu}_k^{t+1}$  in Algorithm 2. The sequence  $\gamma_t = O(\frac{1}{\sqrt{t+1}})$  decreasing to zero.

We now list the assumptions employed in our convergence analysis. All the assumptions are common in the literature. The assumptions are weak because neither smooth assumptions nor strongly convex assumptions are involved. Thus, our algorithm can handle non-smooth cases. More importantly, we do not add any additional assumptions of the auxiliary problem (21). This is important for practical use because we can not check whether the auxiliary problem satisfies the assumptions given a black-box original problem. To the best of our knowledge, our algorithm is the first full matrix adaptive black-box optimization algorithm that achieves a provable  $O(\frac{d^2 K^4}{\sqrt{T}})$  convergence for convex functions without smooth and strongly convex assumptions, and any assumptions of the auxiliary problems.

**Assumption 5.1.**  $f_1(\bar{\mathbf{x}}_1), \dots, f_K(\bar{\mathbf{x}}_K)$  are all convex functions.

**Assumption 5.2.**  $f_i(\bar{\mathbf{x}}_i)$  is a  $L_i$ -Lipschitz continuous function for  $\forall i \in \{1, \dots, K\}$ , i.e.,  $|f_i(\bar{\mathbf{x}}_i) - f_i(\bar{\mathbf{y}}_i)| \leq L_i \|\bar{\mathbf{x}}_i - \bar{\mathbf{y}}_i\|_2$ .

**Assumption 5.3.** The initialization  $\bar{\theta}_K^1 := \{\bar{\mu}_K^1, \bar{\Sigma}_K^1\}$  is bounded, i.e.,  $\sum_{k=1}^K \|\mu_k^1 - \mu_k^*\|_{\Sigma_k^1}^2 \leq R$  and  $\bar{\Sigma}_K^1 \in \mathcal{S}^+$ , and  $\underline{\tau}\mathbf{I} \preceq \bar{\Sigma}_K^1 \preceq \bar{\tau}\mathbf{I}$  for  $\bar{\tau} \geq \underline{\tau} > 0$ , and  $\sum_{k=1}^K \|\mu_k^*\|_2^2 \leq B$ .

**Theorem 5.4.** Suppose the assumptions 5.1 5.2 5.3 holds. Set  $\beta_t = t\beta$  with  $\beta > 0$ ,  $\alpha_t = \sqrt{t+1}\alpha$  with  $\alpha > 0$ , and  $\gamma_t = \frac{\alpha\nu}{\beta\sqrt{t+1}}$ , and  $\nu > 0$ , and  $\omega_t = 1$ . Initialize  $\Sigma_k^1$  such that  $\|\Sigma_k^1\|_2^{-1} \geq \frac{5}{3}\alpha\nu$  for  $\forall k \in \{1, \dots, K\}$ . Then, running Algorithm 2 with  $T$ -steps, we have

$$\frac{1}{T} \sum_{t=1}^T \sum_{k=1}^K f_k(\bar{\mu}_k^t) - \sum_{k=1}^K f_k(\bar{\mu}_k^*) \leq \frac{\sum_{k=1}^K \|\mu_k^1 - \mu_k^*\|_{\Sigma_k^1}^2}{2\beta T} + \frac{2\sqrt{T+1}C_1}{T} + \frac{4(T+1)^{\frac{1}{4}}C_2}{T} + \frac{\sqrt{T+2}C_3}{T} \quad (24)$$

$$\leq O\left(\frac{d^2 K^4}{\sqrt{T}}\right) \quad (25)$$

where  $\bar{\mu}_k^t = [\mu_1^{t\top}, \dots, \mu_k^{t\top}]^\top$  and  $\bar{\mu}_k^* = [\mu_1^{*\top}, \dots, \mu_k^{*\top}]^\top$ . And  $C_1 = \frac{3\beta \sum_{i=1}^K K L_i^2 (id+1)^2}{2\alpha\nu}$  and  $C_2 = \frac{\sum_{i=1}^K \sqrt{3id}L_i}{\sqrt{\alpha\nu}}$ ,  $C_3 = \frac{\alpha\nu B}{\beta}$ .

Detailed proof can be found in Appendix D. In Theorem 5.4, the error term  $\frac{\sqrt{T+2}C_3}{T}$  in Eq.(24) is due to the  $\gamma_t$ -enlargement. Error term  $\frac{4(T+1)^{\frac{1}{4}}C_2}{T}$  results from the covariance matrix term  $\bar{\Sigma}_K$  in Gaussian-smooth relaxation of the original problem. The first two error terms result from the stochastic gradient update.

Note that for a convex function  $f(x)$ , we have  $f(\frac{1}{T} \sum_{t=1}^T x_t) \leq \frac{1}{T} \sum_{t=1}^T f(x_t)$ . Then, we can directly obtain the solution of the original problem (20) by averaging the auxiliary variable  $\mu_k^t$  for the auxiliary problem (21).

**Corollary 5.5.** Suppose the assumptions 5.1 5.2 5.3 holds. Set  $\beta_t = t\beta$  with  $\beta > 0$ ,  $\alpha_t = \sqrt{t+1}\alpha$  with  $\alpha > 0$ , and  $\gamma_t = \frac{\alpha\nu}{\beta\sqrt{t+1}}$ , and  $\nu > 0$ , and  $\omega_t = 1$ . Initialize  $\Sigma_k^1$  such that  $\|\Sigma_k^1\|_2^{-1} \geq \frac{5}{3}\alpha\nu$  for  $\forall k \in \{1, \dots, K\}$ . Then, running Algorithm 2 with  $T$ -steps, set  $\bar{x}_K^T = \frac{1}{T} \sum_{t=1}^T \bar{\mu}_K^t$ , we have the cumulative regret as:

$$\sum_{k=1}^K f_k(\bar{x}_k^T) - \sum_{k=1}^K f_k(\bar{\mu}_k^*) \leq \frac{\sum_{k=1}^K \|\mu_k^1 - \mu_k^*\|_{\Sigma_k^1}^2}{2\beta T} + \frac{2\sqrt{T+1}C_1}{T} + \frac{4(T+1)^{\frac{1}{4}}C_2}{T} + \frac{\sqrt{T+2}C_3}{T} \quad (26)$$

$$\leq O\left(\frac{d^2 K^4}{\sqrt{T}}\right) \quad (27)$$

where  $\bar{x}_k^T = [\mathbf{x}_1^{T\top}, \dots, \mathbf{x}_k^{T\top}]^\top$  and  $\bar{\mu}_k^* = [\mu_1^{*\top}, \dots, \mu_k^{*\top}]^\top$ . And  $C_1 = \frac{3\beta \sum_{i=1}^K K L_i^2 (id+1)^2}{2\alpha\nu}$  and  $C_2 = \frac{\sum_{i=1}^K \sqrt{3id}L_i}{\sqrt{\alpha\nu}}$ ,  $C_3 = \frac{\alpha\nu B}{\beta}$ .

**Remark:** Note that for convex problems, the optimum  $\{\bar{\mu}_K^*, 0\}$  of the auxiliary problem (21) is also the optimum of the original problem (20), i.e.,  $J(\bar{\mu}_K^*, 0) = \hat{F}(\bar{\mu}_K^*)$ . Thus, the solution  $\bar{x}_k^T$  achieve a  $O(\frac{d^2 K^4}{\sqrt{T}})$  cumulative regret of the original problem (20). In addition, our algorithm can handle non-smooth problems without expert designing of proximal operators for different types of non-smooth functions. This can be remarkably interesting when the unknown function involves compositions of lots of different types of non-smooth functions, in which case human experts can not derive the operators explicitly.

## 6 EXPERIMENTS

### 6.1 EMPIRICAL STUDY ON NUMERICAL TEST PROBLEM

We first evaluate our algorithm on minimizing the numerical cumulative summation problem  $\hat{F}(\bar{x}_K) = \sum_{i=1}^K f(\mathbf{x}_i)$  with black-box transition dynamic  $\mathbf{x}_k = \mathbf{x}_{k-1} + 0.1$ . We test two cases:  $L_2$ -norm Ellipsoid  $f(\mathbf{x}) := \sum_{m=1}^d 10^{\frac{6(m-1)}{d-1}} x_m^2$  and  $L_1$ -norm Ellipsoid  $f(\mathbf{x}) := \sum_{m=1}^d 10^{\frac{6(m-1)}{d-1}} |x_m|$

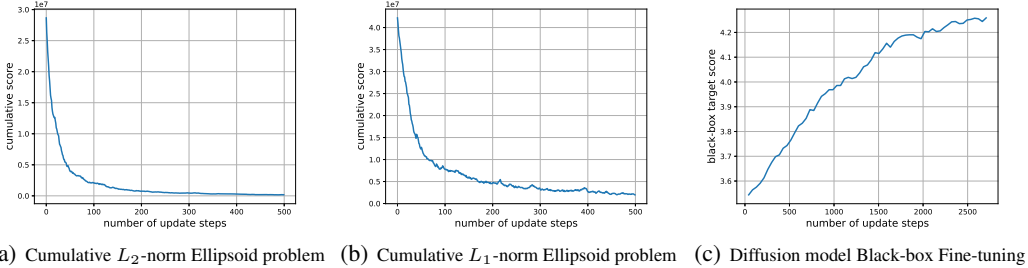


Figure 1: Target score v.s. the number of update steps on different problems

Table 1: Black-box function scores for each method evaluated on 2700 generated images

	Mean	Std	50%ile	80%ile	95%ile
Dataset	0.0000	1.0000	0.0511	0.7876	1.7414
Conditional model	3.5048	0.6672	3.5065	4.0601	4.5863
DDOM	3.6571	0.6252	3.6764	4.1807	4.6631
Fine-tune (steps=1596) (ours)	4.0998	0.5976	4.1259	4.6083	5.0468
Fine-tune (steps=2646) (ours)	4.2287	0.5642	4.2739	4.7069	5.0797
Extend dataset (ours)	4.7948	0.4897	4.8166	5.2222	5.5511

The plot of the target score is shown in Figure 1 (a) and (b), respectively. We can see our converge fast in the sequential optimization test problem.

## 6.2 EMPIRICAL STUDY ON BLACK-BOX DIFFUSION TARGET GENERATION

**Black-box Target Score:** We employ the CLIP model<sup>1</sup> (Radford et al., 2021) to compute the black-box target score. Specifically, we compute the cosine distance between the latent embedding of the generated image and the latent embedding of the target text and employ the normalized cosine distance  $\frac{y - \mu_y}{\sigma_y}$  as the black-box target score, where  $\mu_y$  and  $\sigma_y$  denotes the mean and standard deviation of the cosine distances between the target and the images in the dataset. We chose "a close-up of a man with long hair" as our target text to compute the black-box target score. The larger the score is the better.

**Baselines:** We compare our methods with the DDOM generation (Krishnamoorthy et al., 2023), the conditional model generation conditioned on the pre-computed max score from the dataset, and the trivial max score from the dataset as baselines. In addition to our black-box methods optimization for diffusion fine-tuning, we employ the generated images from our fine-tuned diffusion model to extend the dataset and employ the extended dataset to retrain a conditional model.

**Dataset and Implementation:** In our experiments, we use the dataset CelebA-HQ<sup>2</sup> that contains 30,000 face images. We employ DPM-Solver++ (Lu et al., 2022a) as the sampler for all experiments in both the training and evaluation phases. In all experiments, we set the number of sampling steps as  $K = 14$ . More implementation details can be found in the Appendix E.

For each method, we generate 2700 images and evaluate the normalized scores. The mean, std, 50%-Percentile score, 80%-Percentile score and 95%-Percentile score of each method are shown in Table 1. We can observe that both our fine-tuning method and extended dataset retraining method achieve significantly higher scores compared with baselines consistently across all metrics. Moreover, we can see that re-training with our extended dataset achieves the highest scores among all the compared methods, which demonstrates the ability and potential of our method for target generation. In addition, all the methods obtain significantly higher scores compared with the training dataset.

We further show the plot of the relationship between the black-box function score value and the number of update steps of our method in Figure 1 (c). We observe that the score increases almost linearly with the number of update steps, which demonstrates the potential of our optimization method in diffusion target generation.

<sup>1</sup><https://huggingface.co/sentence-transformers/clip-ViT-L-14>

<sup>2</sup><https://huggingface.co/datasets/huggan/CelebA-HQ>



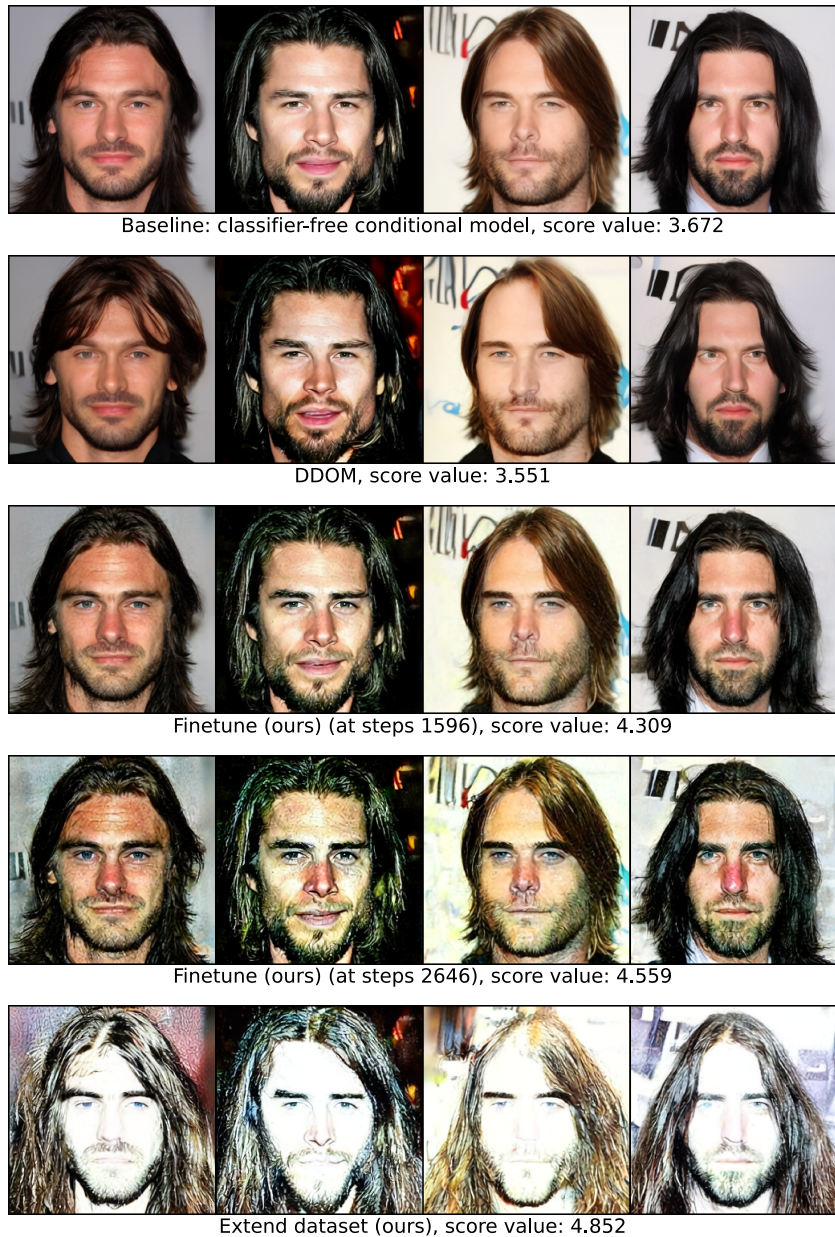


Figure 2: Generated Images of all methods

We further visualize the images generated by different methods with the same initialized noise  $x_0$  in Figure 2. We observe that our method trades off optimizing the score at the cost of image quality. This issue may be mitigated by incorporating quality measurement into the black-box function. We leave it to future work.

## 7 CONCLUSION

In this paper, we formulated the fine-tuning of the diffusion model for black-box target generation as a sequential black-box optimization problem. We proposed a novel stochastic adaptive sequential black-box optimization (SASBO) algorithm to address this problem. Theoretically, we prove a  $O(\frac{d^2 K^4}{\sqrt{T}})$  convergence rate of SASBO without smooth and strongly convex assumptions. Thus, our theoretical results hold true for all non-smooth/smooth convex function families that are of great challenge for full matrix adaptive black-box algorithms to converge. Empirically, our method enables the fine-tuning of the diffusion model to generate targeted images with a high black-box target score.

## REFERENCES

- Stéphane Alarie, Charles Audet, Aïmen E Gheribi, Michael Kokkolaras, and Sébastien Le Digabel. Two decades of blackbox optimization applications. *EURO Journal on Computational Optimization*, 9:100011, 2021.
- Charles Audet and Warren Hare. Derivative-free and blackbox optimization. 2017.
- Thomas Back, Frank Hoffmeister, and Hans-Paul Schwefel. A survey of evolution strategies. In *Proceedings of the fourth international conference on genetic algorithms*, volume 2. Morgan Kaufmann Publishers San Mateo, CA, 1991.
- Benjamin Doerr and Frank Neumann. Theory of evolutionary computation: Recent developments in discrete optimization. 2019.
- Jacob Gardner, Chuan Guo, Kilian Weinberger, Roman Garnett, and Roger Grosse. Discovering and exploiting additive structure for bayesian optimization. In *Artificial Intelligence and Statistics*, pp. 1311–1319. PMLR, 2017.
- Nate Gruver, Samuel Stanton, Nathan C Frey, Tim GJ Rudner, Isidro Hotzel, Julien Lafrance-Vanasse, Arvind Rajpal, Kyunghyun Cho, and Andrew Gordon Wilson. Protein design with guided discrete diffusion. *arXiv preprint arXiv:2305.20009*, 2023.
- Nikolaus Hansen. The cma evolution strategy: a comparing review. In *Towards a new evolutionary computation*, pp. 75–102. Springer, 2006.
- Zhengfu He, Tianxiang Sun, Qiong Tang, Kuanning Wang, Xuanjing Huang, and Xipeng Qiu. DiffusionBERT: Improving generative masked language models with diffusion models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics*, pp. 4521–4534, 2023.
- Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022.
- Minghui Hu, Chuanxia Zheng, Zuopeng Yang, Tat-Jen Cham, Heliang Zheng, Chaoyue Wang, Dacheng Tao, and Ponnuthurai N. Suganthan. Unified discrete diffusion for simultaneous vision-language generation. In *The Eleventh International Conference on Learning Representations*, 2023.
- Gwanghyun Kim, Taesung Kwon, and Jong Chul Ye. Diffusionclip: Text-guided diffusion models for robust image manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2426–2435, 2022a.
- Heeseung Kim, Sungwon Kim, and Sungroh Yoon. Guided-tts: A diffusion model for text-to-speech via classifier guidance. In *International Conference on Machine Learning*, pp. 11119–11133. PMLR, 2022b.
- Diederik Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. *Advances in neural information processing systems*, 34:21696–21707, 2021.
- Zhifeng Kong, Wei Ping, Jiaji Huang, Kexin Zhao, and Bryan Catanzaro. Diffwave: A versatile diffusion model for audio synthesis. *arXiv preprint arXiv:2009.09761*, 2020.
- Siddarth Krishnamoorthy, Satvik Mehul Mashkaria, and Aditya Grover. Diffusion models for black-box optimization. *arXiv preprint arXiv:2306.07180*, 2023.
- Jin Sub Lee, Jisun Kim, and Philip M Kim. Score-based generative modeling for de novo protein design. *Nature Computational Science*, pp. 1–11, 2023.
- Sijia Liu, Bhavya Kailkhura, Pin-Yu Chen, Paishun Ting, Shiyu Chang, and Lisa Amini. Zeroth-order stochastic variance reduction for nonconvex optimization. *Advances in Neural Information Processing Systems*, 31, 2018.

- Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps. *Advances in Neural Information Processing Systems*, 35:5775–5787, 2022a.
- Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver++: Fast solver for guided sampling of diffusion probabilistic models. *arXiv preprint arXiv:2211.01095*, 2022b.
- Yueming Lyu and Ivor W Tsang. Black-box optimizer with stochastic implicit natural gradient. In *Machine Learning and Knowledge Discovery in Databases. Research Track: European Conference, ECML PKDD 2021, Bilbao, Spain, September 13–17, 2021, Proceedings, Part III 21*, pp. 217–232. Springer, 2021.
- Seyedali Mirjalili and Seyedali Mirjalili. Genetic algorithm. *Evolutionary Algorithms and Neural Networks: Theory and Applications*, pp. 43–55, 2019.
- Amin Nayebi, Alexander Munteanu, and Matthias Poloczek. A framework for bayesian optimization in embedded subspaces. In *International Conference on Machine Learning*, pp. 4752–4761. PMLR, 2019.
- Yurii Nesterov and Vladimir Spokoiny. Random gradient-free minimization of convex functions. *Foundations of Computational Mathematics*, 17(2):527–566, 2017.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.
- Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022.
- Ingo Rechenberg and M. Eigen. *Optimierung technischer Systeme nach Prinzipien der biologischen Evolution*. PhD thesis, 1973.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020.
- Mandavilli Srinivas and Lalit M Patnaik. Genetic algorithms: A survey. *computer*, 27(6):17–26, 1994.
- Niranjan Srinivas, Andreas Krause, Sham M Kakade, and Matthias Seeger. Gaussian process optimization in the bandit setting: No regret and experimental design. In *ICML*, 2010.
- Narek Tumanyan, Michal Geyer, Shai Bagon, and Tali Dekel. Plug-and-play diffusion features for text-driven image-to-image translation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 1921–1930, 2023.
- Yining Wang, Simon Du, Sivaraman Balakrishnan, and Aarti Singh. Stochastic zeroth-order optimization in high dimensions. In *International conference on artificial intelligence and statistics*, pp. 1356–1365. PMLR, 2018.
- Daan Wierstra, Tom Schaul, Tobias Glasmachers, Yi Sun, Jan Peters, and Jürgen Schmidhuber. Natural evolution strategies. *The Journal of Machine Learning Research (JMLR)*, 15(1):949–980, 2014.

## APPENDIX

## A ADDITIONAL EXPERIMENTS

## A.1 TARGETED IMAGE GENERATION

We evaluate our revised method on two targeted image generation cases: (1) "long hair man" (2)"Asian face". The "Asian face" is rare in the dataset. As a result, the targeted generation focuses more on out-of-distribution generation, which is more challenging than the "long hair man" case. For the "long hair man" case, we keep the target text as "a close-up of a man with long hair", which is the same as in our previous submission version. For the "Asian face" case, we set the target text as "a high quality close up of an asian".

For both cases, we set the number of iterations  $T$  of our method as  $T = 120$ . The experimental results reported are at  $T = 120$ . We perform independent draws to generate 3,000 images for evaluation. The same set consists of 3,000 i.i.d. sampled initial noise  $x_0 \sim \mathcal{N}(0, I)$  is employed for all the methods to generate images for comparison.

## A.1.1 TARGETED IMAGE GENERATION: LONG HAIR MAN

In this experiment, the target text is set as "a close-up of a man with long hair". Figure 3 shows the comparison of generated images. Figure 4 shows 64 randomly generated images by our method. Table 2 shows the score values.

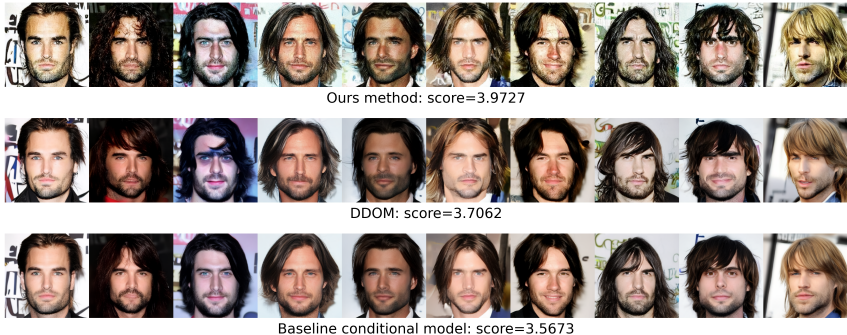


Figure 3: Comparison of different methods on "long hair man"

Table 2: Score values evaluated on 3000 generated images on "long hair man"

	Mean	Std	50%ile	80%ile	95%ile
Dataset	0.0000	1.0000	0.0511	0.7876	1.7414
Conditional model	3.5673	0.6762	3.5672	4.1299	4.6761
DDOM	3.7062	0.6372	3.7110	4.2618	4.7319
Our method	<b>3.9727</b>	<b>0.5902</b>	<b>3.9937</b>	<b>4.4598</b>	<b>4.9197</b>

## A.1.2 TARGETED IMAGE GENERATION: ASIAN FACE

In this experiment, the target text is set as "a high quality close up of an asian". Figure 5 shows the comparison of generated images. Figure 6 shows 64 randomly generated images by our method. Table 3 shows the score values.

Table 3: Score values evaluated on 3000 generated images on "Asian Face"

	Mean	Std	50%ile	80%ile	95%ile
Dataset	0.0000	1.0000	0.0511	0.7876	1.7414
Conditional model	2.8784	0.5654	2.8870	3.3682	3.7874
DDOM	3.2397	0.4824	3.2693	3.6346	4.0078
Our method	<b>3.2435</b>	<b>0.5453</b>	<b>3.2733</b>	<b>3.7074</b>	<b>4.0789</b>

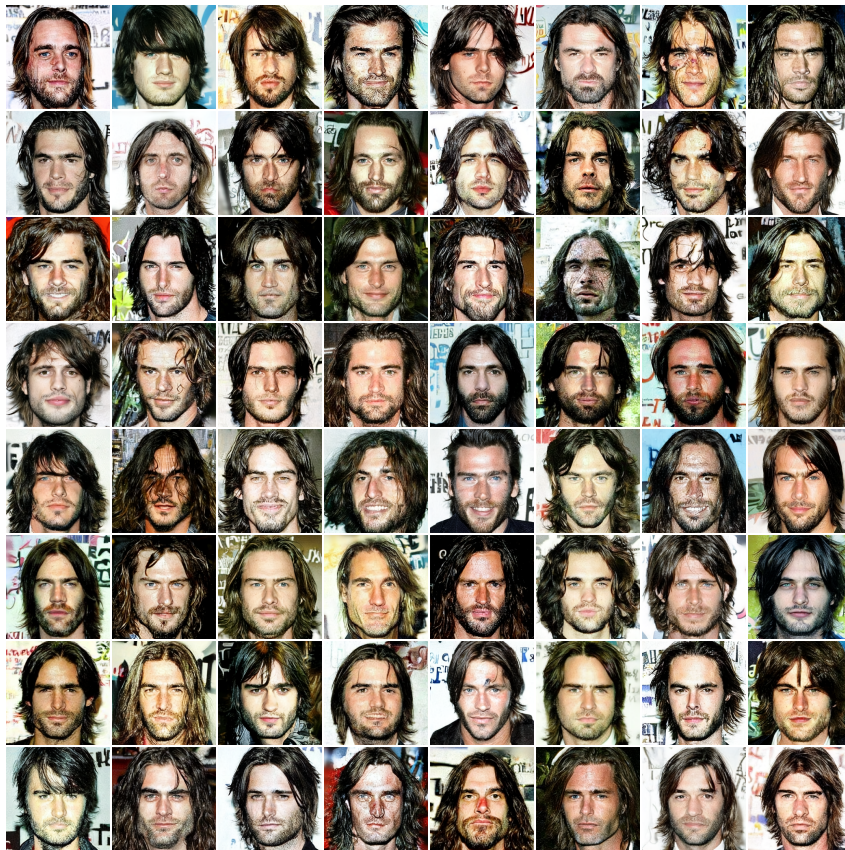


Figure 4: 64 images generated by our method for "long hair man"

Table 4: Normalized maximum score on Design Bench evaluated on 5 independent runs

	SUPERCON.	CHEMBL
Conditional model	0.4824 $\pm$ 0.0466	0.6345 $\pm$ 0.0026
DDOM	0.4777 $\pm$ 0.0350	0.6344 $\pm$ 0.0027
Our method	<b>0.5123 <math>\pm</math> 0.0145</b>	<b>0.6392 <math>\pm</math> 0.0055</b>

We can see that the visual quality of our method’s generated images is better than DDOM’s on the "Asian face" cases. DDOM employs a reweighed sampling of training samples to train a conditional diffusion model from scratch. This training scheme focuses on the tail of the distribution, which is more vulnerable to overfitting, especially for out-of-distribution target generation cases where the target image is rare in the training set.

## A.2 DESIGN-BENCH: SUPERCONDUCTOR

We compare our method with the DDOM’s on the design-bench tasks. We train our model for  $T = 200$  iterations. We follow Krishnamoorthy et al. (2023) to perform 5 independent runs with different seed, and report the normalized maximum score along with standard deviation. The normalization method is the same as Krishnamoorthy et al. (2023). We experiment on two tasks, Superconductor and ChEMBL. Experimental results are shown in Table 4. Our method outperforms the DDOM on both tasks, which shows the potential of our method on different domains beyond the targeted image generation.

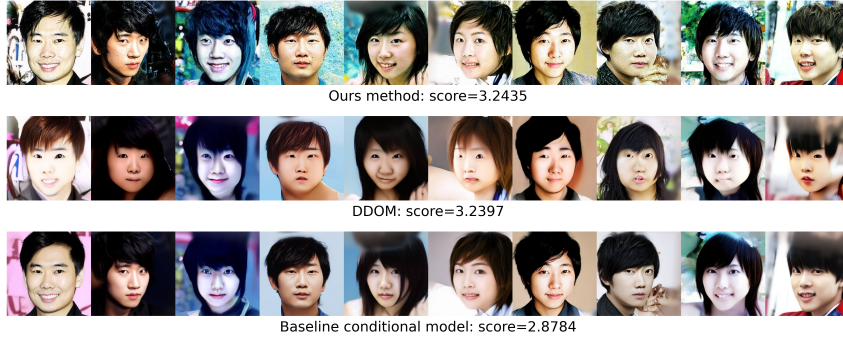


Figure 5: Comparison of different methods on "Asian Face"



Figure 6: 64 images generated by our method for "Asian Face"

## B DERIVATION OF UPDATE RULE

The minimization can be rewritten as

$$\begin{aligned}
 & \sum_{k=1}^K \left\langle \bar{\theta}_k - \bar{\theta}_k^t, \beta_k \nabla_{\bar{\theta}_k^t} J_k(\bar{\theta}_k^t) \right\rangle + \text{KL}(q_\theta \| q_{\theta^t}) = \mathbb{E}_{\mathbf{x}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \sum_{k=1}^K \beta_k \bar{\boldsymbol{\mu}}_k^\top \nabla_{\bar{\boldsymbol{\mu}}_k^t} \mathbb{E}_{q_{\bar{\theta}_k^t}} [f_k(\bar{\mathbf{x}}_k)] + \\
 & \sum_{k=1}^K \beta_k \text{tr}(\bar{\Sigma}_k \nabla_{\bar{\Sigma}_k^t} \mathbb{E}_{q_{\bar{\theta}_k^t}} [f_k(\bar{\mathbf{x}}_k)]) + \frac{1}{2} [\text{tr}(\bar{\Sigma}_K^{-1} \bar{\Sigma}_K) + (\bar{\boldsymbol{\mu}}_K - \bar{\boldsymbol{\mu}}_K^t)^\top \bar{\Sigma}_K^{-1} (\bar{\boldsymbol{\mu}}_K - \bar{\boldsymbol{\mu}}_K^t) + \log \frac{|\bar{\Sigma}_K^t|}{|\bar{\Sigma}_K|} - d],
 \end{aligned} \tag{28}$$

where  $\nabla_{\bar{\boldsymbol{\mu}}_k} \mathbb{E}_{q_{\bar{\theta}_k^t}} [f_k(\bar{\boldsymbol{x}}_k)]$  and  $\nabla_{\bar{\boldsymbol{\Sigma}}_k} \mathbb{E}_{q_{\bar{\theta}_k^t}} [f_k(\bar{\boldsymbol{x}}_k)]$  denotes the derivative w.r.t  $\bar{\boldsymbol{\mu}}_k$  and  $\bar{\boldsymbol{\Sigma}}_k$  taking at  $\bar{\boldsymbol{\mu}}_k = \bar{\boldsymbol{\mu}}_k^t$  and  $\bar{\boldsymbol{\Sigma}}_k = \bar{\boldsymbol{\Sigma}}_k^t$ . The above problem is convex with respect to  $\bar{\boldsymbol{\mu}}_K = [\boldsymbol{\mu}_1^\top, \dots, \boldsymbol{\mu}_K^\top]^\top$  and  $\bar{\boldsymbol{\Sigma}}_K = \text{diag}(\boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_K)$ . Taking the derivative w.r.t  $\bar{\boldsymbol{\mu}}_K$  and  $\bar{\boldsymbol{\Sigma}}_K$  and setting them to zero, for  $k^{\text{th}}$  component, we can obtain that

$$\mathbb{E}_{\boldsymbol{x}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \sum_{i=k}^K \beta_i \nabla_{\boldsymbol{\mu}_k^t} \mathbb{E}_{q_{\bar{\theta}_i^t}} [f_i(\bar{\boldsymbol{x}}_i)] + \boldsymbol{\Sigma}_k^{t-1} (\boldsymbol{\mu}_k - \boldsymbol{\mu}_k^t) = 0 \quad (29)$$

$$\mathbb{E}_{\boldsymbol{x}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \sum_{i=k}^K \beta_i \nabla_{\boldsymbol{\Sigma}_k^t} \mathbb{E}_{q_{\bar{\theta}_i^t}} [f_i(\bar{\boldsymbol{x}}_i)] + \frac{1}{2} [\boldsymbol{\Sigma}_k^{-1} - \boldsymbol{\Sigma}_k^{t-1}] = 0. \quad (30)$$

for  $k \in \{1, \dots, K\}$ .

Set  $\boldsymbol{\mu}_k^{t+1} = \boldsymbol{\mu}_k^t$  and  $\boldsymbol{\Sigma}_k^{t+1-1} = \boldsymbol{\Sigma}_k^{t-1}$  in the above equation. We then have the update rule as

$$\boldsymbol{\mu}_k^{t+1} = \boldsymbol{\mu}_k^t - \mathbb{E}_{\boldsymbol{x}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \sum_{i=k}^K \beta_i \boldsymbol{\Sigma}_k^t \nabla_{\boldsymbol{\mu}_k^t} \mathbb{E}_{q_{\bar{\theta}_i^t}} [f_i(\bar{\boldsymbol{x}}_i)] \quad (31)$$

$$\boldsymbol{\Sigma}_k^{t+1-1} = \boldsymbol{\Sigma}_k^{t-1} + \mathbb{E}_{\boldsymbol{x}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \sum_{i=k}^K 2\beta_i \nabla_{\boldsymbol{\Sigma}_k^t} \mathbb{E}_{q_{\bar{\theta}_i^t}} [f_i(\bar{\boldsymbol{x}}_i)]. \quad (32)$$

for  $k \in \{1, \dots, K\}$ .

In addition, note that the gradient has the following closed-form (Wierstra et al., 2014)

$$\nabla_{\boldsymbol{\mu}_k^t} \mathbb{E}_{q_{\bar{\theta}_i^t}} [f_i(\bar{\boldsymbol{x}}_i)] = \boldsymbol{\Sigma}_k^{t-1} \mathbb{E}_{q_{\bar{\theta}_i^t}} [(\boldsymbol{x}_k - \boldsymbol{\mu}_k) f_i(\bar{\boldsymbol{x}}_i)] \quad (33)$$

$$\nabla_{\boldsymbol{\Sigma}_k^t} \mathbb{E}_{q_{\bar{\theta}_i^t}} [f_i(\bar{\boldsymbol{x}}_i)] = \frac{1}{2} \mathbb{E}_{q_{\bar{\theta}_i^t}} [(\boldsymbol{\Sigma}_k^{t-1} (\boldsymbol{x}_k - \boldsymbol{\mu}_k) (\boldsymbol{x}_k - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{t-1} - \boldsymbol{\Sigma}_k^{t-1}) (f_i(\bar{\boldsymbol{x}}_i))] \quad (34)$$

Then, we have that

$$\boldsymbol{\mu}_k^{t+1} = \boldsymbol{\mu}_k^t - \sum_{i=k}^K \beta_i \mathbb{E}_{\boldsymbol{x}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \mathbb{E}_{q_{\bar{\theta}_i^t}} [(\boldsymbol{x}_k - \boldsymbol{\mu}_k) f_i(\bar{\boldsymbol{x}}_i)] \quad (35)$$

$$\boldsymbol{\Sigma}_k^{t+1-1} = \boldsymbol{\Sigma}_k^{t-1} + \sum_{i=k}^K \beta_i \mathbb{E}_{\boldsymbol{x}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \mathbb{E}_{q_{\bar{\theta}_i^t}} [(\boldsymbol{\Sigma}_k^{t-1} (\boldsymbol{x}_k - \boldsymbol{\mu}_k) (\boldsymbol{x}_k - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{t-1} - \boldsymbol{\Sigma}_k^{t-1}) (f_i(\bar{\boldsymbol{x}}_i))]. \quad (36)$$

## C TECHNICAL LEMMAS

In this section, we introduce the following technical lemmas for convergence analysis.

**Lemma C.1.** *Given a positive definite matrix  $\boldsymbol{\Sigma}$ , we have  $\|\boldsymbol{\Sigma}(\boldsymbol{x} + \boldsymbol{y})\|_{\boldsymbol{\Sigma}^{-1}}^2 \leq 2(\|\boldsymbol{\Sigma}\boldsymbol{x}\|_{\boldsymbol{\Sigma}^{-1}}^2 + \|\boldsymbol{\Sigma}^{\frac{1}{2}}\boldsymbol{y}\|_2^2)$*

*Proof.*

$$\|\boldsymbol{\Sigma}(\boldsymbol{x} + \boldsymbol{y})\|_{\boldsymbol{\Sigma}^{-1}}^2 = \langle \boldsymbol{\Sigma}^{-1} \boldsymbol{\Sigma}(\boldsymbol{x} + \boldsymbol{y}), \boldsymbol{\Sigma}(\boldsymbol{x} + \boldsymbol{y}) \rangle \quad (37)$$

$$= \langle (\boldsymbol{x} + \boldsymbol{y}), \boldsymbol{\Sigma}(\boldsymbol{x} + \boldsymbol{y}) \rangle \quad (38)$$

$$= \|\boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{x} + \boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{y}\|_2^2 \quad (39)$$

$$\leq 2(\|\boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{x}\|_2^2 + \|\boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{y}\|_2^2) \quad (40)$$

Note that  $\|\boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{x}\|_2^2 = \langle \boldsymbol{x}, \boldsymbol{\Sigma}\boldsymbol{x} \rangle = \langle \boldsymbol{\Sigma}^{-1} \boldsymbol{\Sigma}\boldsymbol{x}, \boldsymbol{\Sigma}\boldsymbol{x} \rangle = \|\boldsymbol{\Sigma}\boldsymbol{x}\|_{\boldsymbol{\Sigma}^{-1}}^2$ , we achieve that

$$\|\boldsymbol{\Sigma}(\boldsymbol{x} + \boldsymbol{y})\|_{\boldsymbol{\Sigma}^{-1}}^2 \leq 2(\|\boldsymbol{\Sigma}\boldsymbol{x}\|_{\boldsymbol{\Sigma}^{-1}}^2 + \|\boldsymbol{\Sigma}^{\frac{1}{2}}\boldsymbol{y}\|_2^2) \quad (41)$$

□

**Lemma C.2.** Suppose the gradient estimator  $\hat{g}_{ik}^t$  for the  $i^{\text{th}}$  objective w.r.t. the  $k^{\text{th}}$  component  $\mu_k$  at  $t^{\text{th}}$  iteration as

$$\hat{g}_{ik}^t = \Sigma_k^{t-\frac{1}{2}} \mathbf{z}_k (f_i(\bar{\mu}_i^t + \bar{\Sigma}_i^{\frac{1}{2}} \bar{\mathbf{z}}_i) - f_i(\bar{\mu}_i^t)), \quad (42)$$

where  $\mathbf{z}_1, \dots, \mathbf{z}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$  and  $\bar{\mathbf{z}}_i = [\mathbf{z}_1^\top, \dots, \mathbf{z}_i^\top]^\top$ ,  $i \geq k$ . Suppose assumptions 5.2 hold, using the parameter setting in Theorem 5.4, and the covariance matrix update  $\hat{G}_k^t = \Sigma_k^{t-\frac{1}{2}} \mathbf{H}_k^t \Sigma_k^{t-\frac{1}{2}}$  are positive semi-definite matrix that satisfies  $\nu \mathbf{I} \preceq \hat{G}_k^t$ . Apply the update rule  $\Sigma_k^{t+1-1} = \omega_t \Sigma_k^{t-1} + \alpha_t \hat{G}_k^t$ , we have

- (a)  $\hat{g}_{ik}^t$  is an unbiased estimator of  $g_{ik}^t = \nabla_{\mu_k} \mathbb{E}_{\bar{\mathbf{x}} \sim \mathcal{N}(\bar{\mu}_i^t, \bar{\Sigma}_i^t)} [f_i(\bar{\mathbf{x}})]$ .
- (b)  $\|\Sigma_k^{t+1}\|_2 \leq \frac{1}{\|\Sigma_k^t\|_2^{-1} + \sqrt{t+1}\alpha\nu} \leq \dots \leq \frac{3}{2\alpha\nu} \frac{1}{(t+1)^{\frac{3}{2} + \frac{3}{2}}} \leq \frac{3}{2\alpha\nu} \frac{1}{(t+1)^{\frac{3}{2}}}$ .
- (c)  $\mathbb{E} \sum_{k=1}^K \|\Sigma_k^t (\sum_{i=k}^K \hat{g}_{ik}^t)\|_{\Sigma_k^{t-1}}^2 \leq \frac{3 \sum_{i=1}^K K L_i^2 (id+1)^2}{2t^{\frac{3}{2}} \alpha\nu}$

For the average of i.i.d. sampled unbiased gradient estimators (each one has the same form as Eq.(42)), the results (a),(b),(c) still hold true.

*Proof.* (a). We first show that  $\hat{g}_{ik}^t$  is an unbiased estimator of  $\nabla_{\mu_k} \mathbb{E}_{\bar{\mathbf{x}} \sim \mathcal{N}(\bar{\mu}_i^t, \bar{\Sigma}_i^t)} [f_i(\bar{\mathbf{x}})]$ .

$$\mathbb{E}_{\bar{\mathbf{z}}_i} [\hat{g}_{ik}^t] = \mathbb{E}_{\bar{\mathbf{z}}_i} [\Sigma_k^{t-\frac{1}{2}} \mathbf{z}_k f_i(\bar{\mu}_i^t + \bar{\Sigma}_i^{\frac{1}{2}} \bar{\mathbf{z}}_i)] - \mathbb{E}_{\bar{\mathbf{z}}_i} [\Sigma_k^{t-\frac{1}{2}} \mathbf{z}_k f_i(\bar{\mu}_i^t)] \quad (43)$$

$$= \mathbb{E}_{\bar{\mathbf{z}}_i} [\Sigma_k^{t-\frac{1}{2}} \mathbf{z}_k f_i(\bar{\mu}_i^t + \bar{\Sigma}_i^{\frac{1}{2}} \bar{\mathbf{z}}_i)] \quad (44)$$

$$= \mathbb{E}_{\bar{\mathbf{x}} \sim \mathcal{N}(\bar{\mu}_i^t, \bar{\Sigma}_i^t)} [\Sigma_k^{t-1} (\mathbf{x}_k - \mu_k^t) f_i(\bar{\mathbf{x}})] \quad (45)$$

Note that  $\bar{\Sigma}_i^t = \text{diag}(\Sigma_1^t, \dots, \Sigma_i^t)$  is a block-wise diagonal matrix, and  $\bar{\mu}_i = [\mu_1^\top, \dots, \mu_i^\top]^\top$ ,  $i \geq k$ , we then have that

$$\mathbb{E}_{\bar{\mathbf{z}}_i} [\hat{g}_{ik}^t] = \mathbb{E}_{\bar{\mathbf{x}} \sim \mathcal{N}(\bar{\mu}_i^t, \bar{\Sigma}_i^t)} [\Sigma_k^{t-1} (\mathbf{x}_k - \mu_k^t) f_i(\bar{\mathbf{x}})] = \nabla_{\mu_k} \mathbb{E}_{\bar{\mathbf{x}} \sim \mathcal{N}(\bar{\mu}_i^t, \bar{\Sigma}_i^t)} [f_i(\bar{\mathbf{x}})]. \quad (46)$$

For  $N$  i.i.d. sampled unbiased gradient estimator, the average is still an unbiased gradient estimator.  $\square$

*Proof.* (b) We now prove the decay of the spectral norm of covariance matrix.

From the update rule of  $\Sigma_k^t$ , we know that

$$\Sigma_k^{t+1-1} = \omega_t \Sigma_k^{t-1} + \alpha_t \hat{G}_k^t \quad (47)$$

Note that  $\nu \mathbf{I} \preceq \hat{G}_k^t$ , we have that

$$\lambda_{\min}(\Sigma_k^{t+1-1}) = \lambda_{\min}(\omega_t \Sigma_k^{t-1} + \alpha_t \hat{G}_k^t) \quad (48)$$

$$\geq \omega_t \lambda_{\min}(\Sigma_k^{t-1}) + \alpha_t \nu \quad (49)$$

Note that  $\|\Sigma_k^{t+1}\|_2 = \frac{1}{\lambda_{\min}(\Sigma_k^{t+1-1})}$  and  $\alpha_t = \sqrt{t+1}\alpha$ ,  $\omega_t = 1$ , we have that

$$\|\Sigma_k^{t+1}\|_2 \leq \frac{1}{\omega_t \lambda_{\min}(\Sigma_k^{t-1}) + \alpha_t \nu} = \frac{1}{\|\Sigma_k^t\|_2^{-1} + \sqrt{t+1}\alpha\nu} \quad (50)$$

It follows that

$$\lambda_{\min}(\Sigma_k^{t+1-1}) \geq \lambda_{\min}(\Sigma_k^{t-1}) + \sqrt{t+1}\alpha\nu \quad (51)$$

$$\geq \lambda_{\min}(\Sigma_k^{t-1-1}) + \sqrt{t}\alpha\nu + \sqrt{t+1}\alpha\nu \quad (52)$$

$$\geq \lambda_{\min}(\Sigma_k^{1-1}) + \left(\sum_{i=1}^t \sqrt{i+1}\right)\alpha\nu \quad (53)$$

$$\geq \lambda_{\min}(\Sigma_k^{1-1}) + \frac{2\alpha\nu}{3} ((t+1)^{\frac{3}{2}} - 1) \quad (54)$$



Note that the initialization such that  $\lambda_{\min}(\Sigma_k^1)^{-1} = \|\Sigma_k^1\|_2^{-1} \geq \frac{5}{3}\alpha\nu$ , we have that

$$\lambda_{\min}(\Sigma_k^{t+1}) \geq \frac{2\alpha\nu}{3}(t+1)^{\frac{3}{2}} + \alpha\nu = \frac{2\alpha\nu}{3}\left((t+1)^{\frac{3}{2}} + \frac{3}{2}\right) \quad (55)$$

Note that  $\|\Sigma_k^{t+1}\|_2 = \frac{1}{\lambda_{\min}(\Sigma_k^{t+1})}$ , we then have that

$$\|\Sigma_k^{t+1}\|_2 \leq \frac{3}{2\alpha\nu} \frac{1}{(t+1)^{\frac{3}{2}} + \frac{3}{2}} \quad (56)$$

□

*Proof.* (c). We now prove the upper bound of  $\mathbb{E} \sum_{k=1}^K \|\Sigma_k^t(\sum_{i=k}^K \hat{g}_{ik})\|_{\Sigma_k^{t-1}}^2$ .

Note that  $\hat{g}_{ik}^t = \Sigma_k^{t-\frac{1}{2}} \mathbf{z}_k (f_i(\bar{\mu}_i^t + \bar{\Sigma}_i^{t-\frac{1}{2}} \bar{\mathbf{z}}_i) - f_i(\bar{\mu}_i^t))$ , we have that

$$\|\Sigma_k^t(\sum_{i=k}^K \hat{g}_{ik})\|_{\Sigma_k^{t-1}}^2 = \|\Sigma_k^t \Sigma_k^{t-\frac{1}{2}} \mathbf{z}_k (\sum_{i=k}^K (f_i(\bar{\mu}_i^t + \bar{\Sigma}_i^{t-\frac{1}{2}} \bar{\mathbf{z}}_i) - f_i(\bar{\mu}_i^t)))\|_{\Sigma_k^{t-1}}^2 \quad (57)$$

$$= \|\Sigma_k^{t-\frac{1}{2}} \mathbf{z}_k (\sum_{i=k}^K (f_i(\bar{\mu}_i^t + \bar{\Sigma}_i^{t-\frac{1}{2}} \bar{\mathbf{z}}_i) - f_i(\bar{\mu}_i^t)))\|_{\Sigma_k^{t-1}}^2 \quad (58)$$

$$= \|\mathbf{z}_k (\sum_{i=k}^K (f_i(\bar{\mu}_i^t + \bar{\Sigma}_i^{t-\frac{1}{2}} \bar{\mathbf{z}}_i) - f_i(\bar{\mu}_i^t)))\|_2^2 \quad (59)$$

$$= (\sum_{i=k}^K (f_i(\bar{\mu}_i^t + \bar{\Sigma}_i^{t-\frac{1}{2}} \bar{\mathbf{z}}_i) - f_i(\bar{\mu}_i^t)))^2 \|\mathbf{z}_k\|_2^2 \quad (60)$$

Note that  $f_i(\bar{x})$  is  $L_i$ -Lipschitz continuous function, we then have that

$$\left(\sum_{i=k}^K (f_i(\bar{\mu}_i^t + \bar{\Sigma}_i^{t-\frac{1}{2}} \bar{\mathbf{z}}_i) - f_i(\bar{\mu}_i^t))\right)^2 \leq (K-k+1) \left(\sum_{i=k}^K (f_i(\bar{\mu}_i^t + \bar{\Sigma}_i^{t-\frac{1}{2}} \bar{\mathbf{z}}_i) - f_i(\bar{\mu}_i^t))^2\right) \quad (61)$$

$$\leq (K-k+1) \left(\sum_{i=k}^K L_i^2 \|\bar{\mu}_i^t + \bar{\Sigma}_i^{t-\frac{1}{2}} \bar{\mathbf{z}}_i - \bar{\mu}_i^t\|_2^2\right) \quad (62)$$

$$\leq (K-k+1) \left(\sum_{i=k}^K L_i^2 \|\bar{\Sigma}_i^{t-\frac{1}{2}}\|_2^2 \|\bar{\mathbf{z}}_i\|_2^2\right) \quad (63)$$

$$= (K-k+1) \left(\sum_{i=k}^K L_i^2 \|\bar{\Sigma}_i^t\|_2 \|\bar{\mathbf{z}}_i\|_2^2\right) \quad (64)$$

Plug Eq.(64) into Eq.(60), we have that

$$\|\Sigma_k^t(\sum_{i=k}^K \hat{g}_{ik})\|_{\Sigma_k^{t-1}}^2 \leq (K-k+1) \|\mathbf{z}_k\|_2^2 \left(\sum_{i=k}^K L_i^2 \|\bar{\Sigma}_i^t\|_2 \|\bar{\mathbf{z}}_i\|_2^2\right) \quad (65)$$

It follows that

$$\sum_{k=1}^K \|\Sigma_k^t (\sum_{i=k}^K \hat{g}_{ik})\|_{\Sigma_k^{t-1}}^2 \leq \sum_{k=1}^K (K-k+1) \|\mathbf{z}_k\|_2^2 (\sum_{i=k}^K L_i^2 \|\bar{\Sigma}_i^t\|_2 \|\bar{\mathbf{z}}_i\|_2^2) \quad (66)$$

$$= \sum_{i=1}^K \sum_{k=1}^i ((K-k+1) \|\mathbf{z}_k\|_2^2 L_i^2 \|\bar{\Sigma}_i^t\|_2 \|\bar{\mathbf{z}}_i\|_2^2) \quad (67)$$

$$= \sum_{i=1}^K L_i^2 \|\bar{\Sigma}_i^t\|_2 \|\bar{\mathbf{z}}_i\|_2^2 \left( \sum_{k=1}^i (K-k+1) \|\mathbf{z}_k\|_2^2 \right) \quad (68)$$

$$\leq \sum_{i=1}^K L_i^2 \|\bar{\Sigma}_i^t\|_2 \|\bar{\mathbf{z}}_i\|_2^2 K \|\bar{\mathbf{z}}_i\|_2^2 \quad (69)$$

$$= \sum_{i=1}^K K L_i^2 \|\bar{\Sigma}_i^t\|_2 \|\bar{\mathbf{z}}_i\|_2^4 \quad (70)$$

In addition, note that for  $z \sim \mathcal{N}(0, \sigma^2)$ , we have  $\mathbb{E}[z^4] = 3\sigma^4$ . It follows that

$$\mathbb{E} \|\bar{\mathbf{z}}_i\|_2^4 = \sum_{j=1}^{id} \mathbb{E}[z_j^4] + \sum_{j_1=1}^{id} \sum_{j_2 \neq j_1}^{id} \mathbb{E}[z_{j_1}^2 z_{j_2}^2] = 3id + id(id-1) = i^2 d^2 + 2id \quad (71)$$

We then have that

$$\mathbb{E} \sum_{k=1}^K \|\Sigma_k^t (\sum_{i=k}^K \hat{g}_{ik})\|_{\Sigma_k^{t-1}}^2 \leq \sum_{i=1}^K K L_i^2 \|\bar{\Sigma}_i^t\|_2 \mathbb{E} \|\bar{\mathbf{z}}_i\|_2^4 \quad (72)$$

$$\leq \sum_{i=1}^K K L_i^2 \|\bar{\Sigma}_i^t\|_2 (id+1)^2 \quad (73)$$

Note that  $\bar{\Sigma}_i^t = \text{diag}(\Sigma_1^t, \dots, \Sigma_i^t)$  is a block-wise diagonal matrix, we have  $\|\bar{\Sigma}_i^t\|_2 \leq \max_{k \in \{1, \dots, i\}} \|\Sigma_k^t\|_2$ . Then we know that

$$\mathbb{E} \sum_{k=1}^K \|\Sigma_k^t (\sum_{i=k}^K \hat{g}_{ik})\|_{\Sigma_k^{t-1}}^2 \leq \sum_{i=1}^K K L_i^2 \|\bar{\Sigma}_i^t\|_2 (id+1)^2 \quad (74)$$

$$\leq \max_{k \in \{1, \dots, K\}} \|\Sigma_k^t\|_2 \sum_{i=1}^K K L_i^2 (id+1)^2 \quad (75)$$

From Lemma C.2 (b), we know that  $\|\Sigma_k^t\|_2 \leq \frac{3}{2\alpha\nu} \frac{1}{t^{\frac{3}{2}}}$  for  $\forall k \in \{1, \dots, K\}$ . Then, we have

$$\mathbb{E} \sum_{k=1}^K \|\Sigma_k^t (\sum_{i=k}^K \hat{g}_{ik})\|_{\Sigma_k^{t-1}}^2 \leq \max_{k \in \{1, \dots, K\}} \|\Sigma_k^t\|_2 \sum_{i=1}^K K L_i^2 (id+1)^2 \quad (76)$$

$$\leq \frac{3 \sum_{i=1}^K K L_i^2 (id+1)^2}{2t^{\frac{3}{2}} \alpha\nu} \quad (77)$$

Note that the square norm  $\|\cdot\|_{\Sigma_k^{t-1}}^2$  is a convex function, then for the average of  $N$  i.i.d. sampled gradient estimator  $\hat{g}_{ik}^j, j \in \{1, \dots, N\}$ , we have

$$\mathbb{E} \sum_{k=1}^K \|\Sigma_k^t \left( \sum_{i=k}^K \frac{1}{N} \sum_{j=1}^N \hat{g}_{ik}^j \right)\|_{\Sigma_k^{t-1}}^2 \leq \frac{1}{N} \sum_{j=1}^N \mathbb{E} \sum_{k=1}^K \|\Sigma_k^t \left( \sum_{i=k}^K \hat{g}_{ik}^j \right)\|_{\Sigma_k^{t-1}}^2 \quad (78)$$

$$= \frac{N}{N} \mathbb{E} \sum_{k=1}^K \|\Sigma_k^t \left( \sum_{i=k}^K \hat{g}_{ik} \right)\|_{\Sigma_k^{t-1}}^2 \quad (79)$$

$$\leq \frac{3 \sum_{i=1}^K K L_i^2 (id + 1)^2}{2t^{\frac{3}{2}} \alpha \nu} \quad (80)$$

□

**Lemma C.3.** Denote  $G_i^t = \nabla_{\bar{\Sigma}_i = \bar{\Sigma}_i^t} J_i(\bar{\mu}_i^t, \bar{\Sigma}_i^t)$ . Suppose assumption 5.2 holds, using the parameter setting in Theorem 5.4, and the covariance matrix update  $\hat{G}_k^t = \Sigma_k^{t-\frac{1}{2}} \mathbf{H}_k^t \Sigma_k^{t-\frac{1}{2}}$  are positive semi-definite matrix that satisfies  $\nu \mathbf{I} \preceq \hat{G}_k^t$ . Apply the update rule  $\Sigma_k^{t+1} = \omega_t \Sigma_k^t + \alpha_t \hat{G}_k^t$ , for  $k \in \{1, \dots, K\}$ . Then we have

$$\text{tr}(G_i^t \bar{\Sigma}_i^t) \leq |\text{tr}(G_i^t \bar{\Sigma}_i^t)| \leq \frac{L_i id}{2t^{\frac{3}{4}}} \sqrt{\frac{3}{\alpha \nu}} \quad (81)$$

*Proof.*

$$\text{tr}(G_i^t \bar{\Sigma}_i^t) = \text{tr}(\bar{\Sigma}_i^{t\frac{1}{2}} G_i^t \bar{\Sigma}_i^{t\frac{1}{2}}) \quad (82)$$

$$= \frac{1}{2} \text{tr} \left( \bar{\Sigma}_i^{t\frac{1}{2}} \mathbb{E}_{\mathcal{N}(\bar{\mu}_i^t, \bar{\Sigma}_i^t)} [(\bar{\Sigma}_i^{t-1}(\bar{\mathbf{x}}_i - \bar{\mu}_i^t)(\bar{\mathbf{x}}_i - \bar{\mu}_i^t)^\top \bar{\Sigma}_i^{t-1} - \bar{\Sigma}_i^{t-1}) f_i(\bar{\mathbf{x}}_i)] \bar{\Sigma}_i^{t\frac{1}{2}} \right) \quad (83)$$

$$= \frac{1}{2} \text{tr} \left( \mathbb{E}_{\mathcal{N}(\bar{\mu}_i^t, \bar{\Sigma}_i^t)} [(\bar{\Sigma}_i^{t-\frac{1}{2}}(\bar{\mathbf{x}}_i - \bar{\mu}_i^t)(\bar{\mathbf{x}}_i - \bar{\mu}_i^t)^\top \bar{\Sigma}_i^{t-\frac{1}{2}} - \mathbf{I}) f_i(\bar{\mathbf{x}}_i)] \right) \quad (84)$$

$$= \frac{1}{2} \text{tr} \left( \mathbb{E}_{\bar{\mathbf{z}} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [(\bar{\mathbf{z}} \bar{\mathbf{z}}^\top - \mathbf{I}) f_i(\bar{\mu}_i^t + \bar{\Sigma}_i^{t\frac{1}{2}} \bar{\mathbf{z}})] \right) \quad (85)$$

$$= \frac{1}{2} \text{tr} \left( \mathbb{E}_{\bar{\mathbf{z}} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [(\bar{\mathbf{z}} \bar{\mathbf{z}}^\top - \mathbf{I}) (f_i(\bar{\mu}_i^t + \bar{\Sigma}_i^{t\frac{1}{2}} \bar{\mathbf{z}}) - f_i(\bar{\mu}_i^t))] \right) \quad (86)$$

$$= \frac{1}{2} \mathbb{E}_{\bar{\mathbf{z}} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[ \left( \sum_{j=1}^{id} (z_j^2 - 1) \right) (f_i(\bar{\mu}_i^t + \bar{\Sigma}_i^{t\frac{1}{2}} \bar{\mathbf{z}}) - f_i(\bar{\mu}_i^t)) \right] \quad (87)$$

where  $z_j$  denotes the  $j^{\text{th}}$  element in  $\bar{\mathbf{z}}$ .

From Cauchy–Schwarz inequality  $|\mathbb{E}[XY]| \leq \sqrt{\mathbb{E}[X^2] \mathbb{E}[Y^2]}$ , we know that

$$|\text{tr}(G_i^t \bar{\Sigma}_i^t)| = \frac{1}{2} \mathbb{E}_{\bar{\mathbf{z}} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[ \left( \sum_{j=1}^{id} (z_j^2 - 1) \right) (f_i(\bar{\mu}_i^t + \bar{\Sigma}_i^{t\frac{1}{2}} \bar{\mathbf{z}}) - f_i(\bar{\mu}_i^t)) \right] \quad (88)$$

$$\leq \frac{1}{2} \sqrt{\mathbb{E}_{\bar{\mathbf{z}} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[ \left( \sum_{j=1}^{id} (z_j^2 - 1) \right)^2 \right] \mathbb{E}_{\bar{\mathbf{z}} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[ (f_i(\bar{\mu}_i^t + \bar{\Sigma}_i^{t\frac{1}{2}} \bar{\mathbf{z}}) - f_i(\bar{\mu}_i^t))^2 \right]} \quad (89)$$

We first check the term  $\mathbb{E}_{\bar{\mathbf{z}} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [(\sum_{j=1}^{id} (z_j^2 - 1))^2]$ .

$$\mathbb{E}_{\bar{\mathbf{z}} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [(\sum_{j=1}^{id} (z_j^2 - 1))^2] = \sum_{j=1}^{id} \mathbb{E}(z_j^2 - 1)^2 + \sum_{j_1=1}^{id} \sum_{j_2 \neq j_1}^{id} \mathbb{E}(z_{j_1}^2 - 1)(z_{j_2}^2 - 1) \quad (90)$$

$$= \sum_{j=1}^{id} \mathbb{E}(z_j^4 - 2z_j^2 + 1) + \sum_{j_1=1}^{id} \sum_{j_2 \neq j_1}^{id} \mathbb{E}(z_{j_1}^2 - 1)\mathbb{E}(z_{j_2}^2 - 1) \quad (91)$$

$$= \sum_{j=1}^{id} \mathbb{E}(z_j^4 - 2z_j^2 + 1) = \sum_{j=1}^{id} [3 - 2 + 1] = 2id \quad (92)$$

We now check the term  $\mathbb{E}_{\bar{\mathbf{z}} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [(f_i(\bar{\boldsymbol{\mu}}_i^t + \bar{\boldsymbol{\Sigma}}_i^{t \frac{1}{2}} \bar{\mathbf{z}}) - f_i(\bar{\boldsymbol{\mu}}_i^t))^2]$ . Note that  $f_i(\mathbf{x})$  is  $L_i$ -Lipschitz continuous function, we then have that

$$\mathbb{E}_{\bar{\mathbf{z}} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [(f_i(\bar{\boldsymbol{\mu}}_i^t + \bar{\boldsymbol{\Sigma}}_i^{t \frac{1}{2}} \bar{\mathbf{z}}) - f_i(\bar{\boldsymbol{\mu}}_i^t))^2] \leq L_i^2 \mathbb{E}_{\bar{\mathbf{z}} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [\|\bar{\boldsymbol{\Sigma}}_i^{t \frac{1}{2}} \bar{\mathbf{z}}\|_2^2] \quad (93)$$

$$\leq L_i^2 \|\bar{\boldsymbol{\Sigma}}_i^{t \frac{1}{2}}\|_2^2 \mathbb{E}_{\bar{\mathbf{z}} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \|\bar{\mathbf{z}}\|_2^2 \quad (94)$$

$$= L_i^2 \|\bar{\boldsymbol{\Sigma}}_i^t\|_2 id \quad (95)$$

From Lemma C.2 (b) we know that

$$\|\bar{\boldsymbol{\Sigma}}_i^t\|_2 = \max_{k \in \{1, \dots, i\}} \|\boldsymbol{\Sigma}_k^t\|_2 \leq \frac{3}{2\alpha\nu} \frac{1}{t^{\frac{3}{2}}} \quad (96)$$

Together with Eq.(95) and Eq.(96), we know that

$$\mathbb{E}_{\bar{\mathbf{z}} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [(f_i(\bar{\boldsymbol{\mu}}_i^t + \bar{\boldsymbol{\Sigma}}_i^{t \frac{1}{2}} \bar{\mathbf{z}}) - f_i(\bar{\boldsymbol{\mu}}_i^t))^2] \leq \frac{3L_i^2 id}{2t^{\frac{3}{2}} \alpha\nu} \quad (97)$$

Plug Eq.(97) and Eq.(92) into Eq.(89), we have that

$$|\text{tr}(G_i^t \bar{\boldsymbol{\Sigma}}_i^t)| \leq \frac{L_i id}{2t^{\frac{3}{4}}} \sqrt{\frac{3}{\alpha\nu}} \quad (98)$$

□

**Lemma C.4.** Given a convex function  $f(x)$ , for Gaussian distribution with parameters  $\boldsymbol{\theta} := \{\boldsymbol{\mu}, \boldsymbol{\Sigma}^{\frac{1}{2}}\}$ , let  $\bar{J}(\boldsymbol{\theta}) := \mathbb{E}_{p(\mathbf{x}; \boldsymbol{\theta})} [f(\mathbf{x})]$ . Then  $\bar{J}(\boldsymbol{\theta})$  is a convex function with respect to  $\boldsymbol{\theta}$ .

*Proof.* For  $\lambda \in [0, 1]$ , we have

$$\lambda \bar{J}(\boldsymbol{\theta}_1) + (1 - \lambda) \bar{J}(\boldsymbol{\theta}_2) = \lambda \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [f(\boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_1^{\frac{1}{2}} \mathbf{z})] + (1 - \lambda) \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [f(\boldsymbol{\mu}_2 + \boldsymbol{\Sigma}_2^{\frac{1}{2}} \mathbf{z})] \quad (99)$$

$$= \mathbb{E}[\lambda f(\boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_1^{\frac{1}{2}} \mathbf{z}) + (1 - \lambda) f(\boldsymbol{\mu}_2 + \boldsymbol{\Sigma}_2^{\frac{1}{2}} \mathbf{z})] \quad (100)$$

$$\geq \mathbb{E}[f(\lambda \boldsymbol{\mu}_1 + (1 - \lambda) \boldsymbol{\mu}_2 + (\lambda \boldsymbol{\Sigma}_1^{\frac{1}{2}} + (1 - \lambda) \boldsymbol{\Sigma}_2^{\frac{1}{2}}) \mathbf{z})] \quad (101)$$

$$= \bar{J}(\lambda \boldsymbol{\theta}_1 + (1 - \lambda) \boldsymbol{\theta}_2) \quad (102)$$

□

**Lemma C.5.** Given a convex function  $f(x)$ , let  $J(\boldsymbol{\mu}, \boldsymbol{\Sigma}) := \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})} [f(\mathbf{x})]$ . Then, we have

$$f(\boldsymbol{\mu}) - f(\boldsymbol{\mu}^*) \leq J(\boldsymbol{\mu}, \boldsymbol{\Sigma}) - J(\boldsymbol{\mu}^*, \mathbf{0}) \quad (103)$$

*Proof.* From the definition of  $J(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , we know that  $f(\boldsymbol{\mu}^*) = J(\boldsymbol{\mu}^*, \mathbf{0})$ .

Note that  $f(\mathbf{x})$  is a convex function, we have that

$$f(\boldsymbol{\mu}) = f(\mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})} [\mathbf{x}]) \leq \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})} [f(\mathbf{x})] = J(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad (104)$$

It follows that

$$f(\boldsymbol{\mu}) - f(\boldsymbol{\mu}^*) \leq J(\boldsymbol{\mu}, \boldsymbol{\Sigma}) - J(\boldsymbol{\mu}^*, \mathbf{0}) \quad (105)$$

□

## D PROOF OF THEOREM 5.4

In this section, we prove our main Theorem 5.4. We decompose the proof into two parts. The proof of Theorem D.1 and the proof of Theorem D.2. Together with Theorem D.1 and the Theorem D.2, we achieve our main Theorem 5.4.

**Theorem D.1.** *Suppose the assumptions 5.1 5.2 5.3 holds. Set  $\beta_t = t\beta$  with  $\beta > 0$ ,  $\alpha_t = \sqrt{t+1}\alpha$  with  $\alpha > 0$ , and  $\gamma_t = \frac{\alpha\nu}{\beta\sqrt{t+1}}$ , and  $\nu > 0$ , and  $\omega_t = 1$ . Initialize  $\Sigma_k^1$  such that  $\|\Sigma_k^1\|_2^{-1} \geq \frac{5}{3}\alpha\nu$  for  $\forall k \in \{1, \dots, K\}$ . Suppose the constraints  $\mathbf{H}_k^t \preceq \frac{1}{\alpha_t}(\frac{\beta_{t+1}}{\beta_t} - \omega_t)\mathbf{I} + \frac{\beta_{t+1}\gamma_t}{\alpha_t}\Sigma_k^t$  and  $\nu\mathbf{I} \preceq \hat{G}_k^t = \Sigma_k^t{}^{-\frac{1}{2}}\mathbf{H}_k^t\Sigma_k^t{}^{-\frac{1}{2}}$  always have feasible solutions. Then, running Algorithm 2 with  $T$ -steps, we have*

$$\frac{1}{T} \sum_{t=1}^T \sum_{k=1}^K f_k(\bar{\mu}_k^t) - \sum_{k=1}^K f_k(\bar{\mu}_k^*) \leq \frac{\sum_{k=1}^K \|\mu_k^1 - \mu_k^*\|_{\Sigma_k^{1-1}}^2}{2\beta T} + \frac{2\sqrt{T+1}C_1}{T} + \frac{4(T+1)^{\frac{1}{4}}C_2}{T} + \frac{\sqrt{T+2}C_3}{T} \quad (106)$$

$$\leq O\left(\frac{d^2 K^4}{\sqrt{T}}\right) \quad (107)$$

where  $\bar{\mu}_k^t = [\mu_1^{t\top}, \dots, \mu_k^{t\top}]^\top$  and  $\bar{\mu}_k^* = [\mu_1^{*\top}, \dots, \mu_k^{*\top}]^\top$ . And  $C_1 = \frac{3\beta \sum_{i=1}^K K L_i^2 (id+1)^2}{2\alpha\nu}$  and  $C_2 = \frac{\sum_{i=1}^K \sqrt{3idL_i}}{\sqrt{\alpha\nu}}$ ,  $C_3 = \frac{\alpha\nu B}{\beta}$

*Proof.* For  $\forall k \in \{1, \dots, K\}$ , we have

$$\begin{aligned} & \|\mu_k^{t+1} - \mu_k^*\|_{\Sigma_k^{t-1}}^2 \\ &= \|\mu_k^t - \beta_t \Sigma_k^t ((\sum_{i=k}^K \hat{g}_{ik}^t) + \gamma_t \mu_k^t) - \mu_k^*\|_{\Sigma_k^{t-1}}^2 \end{aligned} \quad (108)$$

$$= \|\mu_k^t - \mu_k^*\|_{\Sigma_k^{t-1}}^2 - 2\beta_t \left\langle \Sigma_k^t ((\sum_{i=k}^K \hat{g}_{ik}^t) + \gamma_t \mu_k^t), \mu_k^t - \mu_k^* \right\rangle_{\Sigma_k^{t-1}} + \beta_t^2 \|\Sigma_k^t ((\sum_{i=k}^K \hat{g}_{ik}^t) + \gamma_t \mu_k^t)\|_{\Sigma_k^{t-1}}^2 \quad (109)$$

$$= \|\mu_k^t - \mu_k^*\|_{\Sigma_k^{t-1}}^2 - 2\beta_t \left\langle \gamma_t \mu_k^t + \sum_{i=k}^K \hat{g}_{ik}^t, \mu_k^t - \mu_k^* \right\rangle + \beta_t^2 \|\Sigma_k^t ((\sum_{i=k}^K \hat{g}_{ik}^t) + \gamma_t \mu_k^t)\|_{\Sigma_k^{t-1}}^2 \quad (110)$$

Note that

$$\gamma_t \langle \mu_k^t, \mu_k^t - \mu_k^* \rangle = \frac{\gamma_t}{2} \|\mu_k^t - \mu_k^*\|_2^2 - \frac{\gamma_t}{2} \|\mu_k^*\|_2^2 + \frac{\gamma_t}{2} \|\mu_k^t\|_2^2 \quad (111)$$

Plug Eq.(111) into Eq.(110), we have that

$$\begin{aligned} & \|\mu_k^{t+1} - \mu_k^*\|_{\Sigma_k^{t-1}}^2 \\ &= \|\mu_k^t - \mu_k^*\|_{\Sigma_k^{t-1}}^2 - 2\beta_t \left\langle \sum_{i=k}^K \hat{g}_{ik}^t, \mu_k^t - \mu_k^* \right\rangle - \beta_t \gamma_t (\|\mu_k^t - \mu_k^*\|_2^2 - \|\mu_k^*\|_2^2 + \|\mu_k^t\|_2^2) + \beta_t^2 \|\Sigma_k^t ((\sum_{i=k}^K \hat{g}_{ik}^t) + \gamma_t \mu_k^t)\|_{\Sigma_k^{t-1}}^2 \end{aligned} \quad (112)$$

From Lemma C.1, we then have that

$$\begin{aligned} \|\mu_k^{t+1} - \mu_k^*\|_{\Sigma_k^{t-1}}^2 &\leq \|\mu_k^t - \mu_k^*\|_{\Sigma_k^{t-1}}^2 - 2\beta_t \left\langle \sum_{i=k}^K \hat{g}_{ik}^t, \mu_k^t - \mu_k^* \right\rangle - \beta_t \gamma_t (\|\mu_k^t - \mu_k^*\|_2^2 - \|\mu_k^*\|_2^2 + \|\mu_k^t\|_2^2) \\ &\quad + 2\beta_t^2 \|\Sigma_k^t ((\sum_{i=k}^K \hat{g}_{ik}^t))\|_{\Sigma_k^{t-1}}^2 + 2\beta_t^2 \|\gamma_t \Sigma_k^t{}^{\frac{1}{2}} \mu_k^t\|_2^2 \end{aligned} \quad (113)$$

From Lemma C.2 (b), we know that  $\|\Sigma_k^t\|_2 \leq \frac{3}{2\alpha\nu} \frac{1}{t\sqrt{t+3/2}}$ , together with the setting  $\beta_t = t\beta$  and  $\gamma_t = \frac{\alpha\nu}{\beta\sqrt{t+1}}$ , we know that

$$-\gamma_t \|\mu_k^t\|_2^2 + 2\beta_t \|\gamma_t \Sigma_k^t \mu_k^t\|_2^2 = -\gamma_t \|\mu_k^t\|_2^2 + 2\beta_t \gamma_t^2 \|\Sigma_k^t \mu_k^t\|_2^2 \quad (114)$$

$$\leq -\gamma_t \|\mu_k^t\|_2^2 + 2\beta_t \gamma_t^2 \|\Sigma_k^t\|_2^2 \|\mu_k^t\|_2^2 \quad (115)$$

$$= -\gamma_t \|\mu_k^t\|_2^2 + 2\beta_t \gamma_t^2 \|\Sigma_k^t\|_2 \|\mu_k^t\|_2^2 \quad (116)$$

$$= \gamma_t \|\mu_k^t\|_2^2 (-1 + 2\beta_t \gamma_t \|\Sigma_k^t\|_2) \quad (117)$$

$$\leq \gamma_t \|\mu_k^t\|_2^2 \left(-1 + 2t\beta \frac{\alpha\nu}{\beta\sqrt{t+1}} \frac{3}{2\alpha\nu} \frac{1}{3/2 + t\sqrt{t}}\right) \quad (118)$$

$$= \gamma_t \|\mu_k^t\|_2^2 \left(-1 + \frac{3t}{\frac{3}{2}\sqrt{(t+1)} + t\sqrt{t(t+1)}}\right) \quad (119)$$

We now check the term  $\left(-1 + \frac{3t}{\frac{3}{2}\sqrt{(t+1)} + t\sqrt{t(t+1)}}\right)$ . For  $t = 1$  and  $t = 2$ , it is easy to see the term  $\left(-1 + \frac{3t}{\frac{3}{2}\sqrt{(t+1)} + t\sqrt{t(t+1)}}\right) \leq 0$ . For  $t \geq 3$ , we have that

$$\frac{3}{2}\sqrt{(t+1)} + t\sqrt{t(t+1)} - 3t \geq \frac{3}{2}\sqrt{(t+1)} + t^2 - 3t \geq 0 \quad (120)$$

It follows that  $\left(-1 + \frac{3t}{\frac{3}{2}\sqrt{(t+1)} + t\sqrt{t(t+1)}}\right) \leq 0$ . Thus, we have that

$$-\gamma_t \|\mu_k^t\|_2^2 + 2\beta_t \|\gamma_t \Sigma_k^t \mu_k^t\|_2^2 \leq 0 \quad (121)$$

Plug the inequality (121) into inequality (113), we know that

$$\begin{aligned} \|\mu_k^{t+1} - \mu_k^*\|_{\Sigma_k^{t-1}}^2 &\leq \|\mu_k^t - \mu_k^*\|_{\Sigma_k^{t-1}}^2 - 2\beta_t \left\langle \sum_{i=k}^K \hat{g}_{ik}^t, \mu_k^t - \mu_k^* \right\rangle - \beta_t \gamma_t (\|\mu_k^t - \mu_k^*\|_2^2 - \|\mu_k^*\|_2^2) \\ &\quad + 2\beta_t^2 \|\Sigma_k^t\|_2 \left( \sum_{i=k}^K \hat{g}_{ik}^t \right)_{\Sigma_k^{t-1}}^2 \end{aligned} \quad (122)$$

It follows that

$$\begin{aligned} \sum_{k=1}^K \|\mu_k^{t+1} - \mu_k^*\|_{\Sigma_k^{t-1}}^2 &\leq \sum_{k=1}^K \|\mu_k^t - \mu_k^*\|_{\Sigma_k^{t-1}}^2 - 2\beta_t \sum_{k=1}^K \left\langle \sum_{i=k}^K \hat{g}_{ik}^t, \mu_k^t - \mu_k^* \right\rangle + 2\beta_t^2 \sum_{k=1}^K \|\Sigma_k^t\|_2 \left( \sum_{i=k}^K \hat{g}_{ik}^t \right)_{\Sigma_k^{t-1}}^2 \\ &\quad - \sum_{k=1}^K \beta_t \gamma_t (\|\mu_k^t - \mu_k^*\|_2^2 - \|\mu_k^*\|_2^2) \end{aligned} \quad (123)$$

Then, we have that

$$\begin{aligned} \mathbb{E} \sum_{k=1}^K \|\mu_k^{t+1} - \mu_k^*\|_{\Sigma_k^{t-1}}^2 &\leq \sum_{k=1}^K \mathbb{E} \|\mu_k^t - \mu_k^*\|_{\Sigma_k^{t-1}}^2 - 2\beta_t \mathbb{E} \sum_{k=1}^K \left\langle \sum_{i=k}^K \hat{g}_{ik}^t, \mu_k^t - \mu_k^* \right\rangle + 2\beta_t^2 \mathbb{E} \sum_{k=1}^K \|\Sigma_k^t\|_2 \left( \sum_{i=k}^K \hat{g}_{ik}^t \right)_{\Sigma_k^{t-1}}^2 \\ &\quad - \sum_{k=1}^K \beta_t \gamma_t (\mathbb{E} \|\mu_k^t - \mu_k^*\|_2^2 - \|\mu_k^*\|_2^2) \end{aligned} \quad (124)$$

Note that  $\bar{\mu}_k^t = [\mu_1^{t\top}, \dots, \mu_k^{t\top}]^\top$  and  $\bar{\mu}_k^* = [\mu_1^{*\top}, \dots, \mu_k^{*\top}]^\top$ , together with Lemma C.2 (a), we have that

$$\mathbb{E} \sum_{k=1}^K \left\langle \sum_{i=k}^K \hat{g}_{ik}^t, \mu_k^t - \mu_k^* \right\rangle = \mathbb{E} \sum_{i=1}^K \left\langle \sum_{k=1}^i \hat{g}_{ik}^t, \mu_k^t - \mu_k^* \right\rangle = \sum_{i=1}^K \langle g_i^t, \bar{\mu}_i^t - \bar{\mu}_i^* \rangle \quad (125)$$

where  $g_i^t = \nabla_{\bar{\mu}_i^t} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\bar{\mu}_i^t, \bar{\Sigma}_i^t)} [f_i(\mathbf{x})] = \nabla_{\bar{\mu}_i^t} J_i(\bar{\mu}_i^t, \bar{\Sigma}_i^t)$

From Eq. (125) and Eq. (124), we have that

$$\begin{aligned} \sum_{i=1}^K \langle g_i^t, \bar{\boldsymbol{\mu}}_i^t - \bar{\boldsymbol{\mu}}_i^* \rangle &\leq \frac{\sum_{k=1}^K \mathbb{E} \|\boldsymbol{\mu}_k^t - \boldsymbol{\mu}_k^*\|_{\Sigma_k^{t-1}}^2 - \sum_{k=1}^K \mathbb{E} \|\boldsymbol{\mu}_k^{t+1} - \boldsymbol{\mu}_k^*\|_{\Sigma_k^{t-1}}^2}{2\beta_t} + \beta_t \sum_{k=1}^K \mathbb{E} \|\Sigma_k^t (\sum_{i=k}^K \hat{g}_{ik})\|_{\Sigma_k^{t-1}}^2 \\ &\quad - \sum_{k=1}^K \frac{\gamma_t}{2} (\mathbb{E} \|\boldsymbol{\mu}_k^t - \boldsymbol{\mu}_k^*\|_2^2 - \|\boldsymbol{\mu}_k^*\|_2^2) \end{aligned} \quad (126)$$

From Lemma C.4, we know that for  $\forall i \in \{1, \dots, K\}$ ,  $J_i(\bar{\boldsymbol{\mu}}_i, \bar{\Sigma}_i)$  is convex function w.r.t.  $\bar{\boldsymbol{\mu}}_i$  and  $\bar{\Sigma}_i^{\frac{1}{2}}$ . Then, we have that

$$J_i(\bar{\boldsymbol{\mu}}_i^t, \bar{\Sigma}_i^t) - J_i(\bar{\boldsymbol{\mu}}_i^*, 0) \leq \langle g_i^t, \bar{\boldsymbol{\mu}}_i^t - \bar{\boldsymbol{\mu}}_i^* \rangle + \langle \nabla_{\bar{\Sigma}_i^{\frac{1}{2}} = \bar{\Sigma}_i^{\frac{1}{2}}} J_i(\bar{\boldsymbol{\mu}}_i^t, \bar{\Sigma}_i^t), \bar{\Sigma}_i^{\frac{1}{2}} - 0 \rangle \quad (127)$$

Denote  $G_i^t = \nabla_{\bar{\Sigma}_i = \bar{\Sigma}_i^t} J_i(\bar{\boldsymbol{\mu}}_i^t, \bar{\Sigma}_i^t)$ . Note that  $\nabla_{\bar{\Sigma}_i^{\frac{1}{2}}} J_i = \bar{\Sigma}_i^{\frac{t}{2}} \nabla_{\bar{\Sigma}_i^t} J_i + \nabla_{\bar{\Sigma}_i^t} J_i \bar{\Sigma}_i^{\frac{t}{2}}$ , and  $G_i^t, \nabla_{\bar{\Sigma}_i^t} J_i$  and  $\bar{\Sigma}_i^{\frac{t}{2}}$  are symmetric matrix, it follows that

$$\langle \nabla_{\bar{\Sigma}_i^{\frac{1}{2}} = \bar{\Sigma}_i^{\frac{1}{2}}} J_i(\bar{\boldsymbol{\mu}}_i^t, \bar{\Sigma}_i^t), \bar{\Sigma}_i^{\frac{1}{2}} - 0 \rangle = \langle \bar{\Sigma}_i^{\frac{t}{2}} G_i^t + G_i^t \bar{\Sigma}_i^{\frac{t}{2}}, \bar{\Sigma}_i^{\frac{1}{2}} \rangle \quad (128)$$

$$= 2 \langle G_i^t, \bar{\Sigma}_i^t \rangle = 2 \text{tr}(G_i^t \bar{\Sigma}_i^t) \quad (129)$$

Plug Eq.(129) into Eq. (127), we have that

$$J_i(\bar{\boldsymbol{\mu}}_i^t, \bar{\Sigma}_i^t) - J_i(\bar{\boldsymbol{\mu}}_i^*, 0) \leq \langle g_i^t, \bar{\boldsymbol{\mu}}_i^t - \bar{\boldsymbol{\mu}}_i^* \rangle + 2 \text{tr}(G_i^t \bar{\Sigma}_i^t) \quad (130)$$

It follows that

$$\sum_{i=1}^K J_i(\bar{\boldsymbol{\mu}}_i^t, \bar{\Sigma}_i^t) - \sum_{i=1}^K J_i(\bar{\boldsymbol{\mu}}_i^*, 0) \leq \sum_{i=1}^K \langle g_i^t, \bar{\boldsymbol{\mu}}_i^t - \bar{\boldsymbol{\mu}}_i^* \rangle + 2 \sum_{i=1}^K \text{tr}(G_i^t \bar{\Sigma}_i^t) \quad (131)$$

Plug Eq.(126) into Eq.(131), we have that

$$\begin{aligned} &\sum_{i=1}^K J_i(\bar{\boldsymbol{\mu}}_i^t, \bar{\Sigma}_i^t) - \sum_{i=1}^K J_i(\bar{\boldsymbol{\mu}}_i^*, 0) \\ &\leq \frac{\sum_{k=1}^K \mathbb{E} \|\boldsymbol{\mu}_k^t - \boldsymbol{\mu}_k^*\|_{\Sigma_k^{t-1}}^2 - \sum_{k=1}^K \mathbb{E} \|\boldsymbol{\mu}_k^{t+1} - \boldsymbol{\mu}_k^*\|_{\Sigma_k^{t-1}}^2 - \beta_t \gamma_t \sum_{i=1}^K \mathbb{E} \|\boldsymbol{\mu}_i^t - \boldsymbol{\mu}_i^*\|_2^2}{2\beta_t} \\ &\quad + \beta_t \sum_{k=1}^K \mathbb{E} \|\Sigma_k^t (\sum_{i=k}^K \hat{g}_{ik})\|_{\Sigma_k^{t-1}}^2 + 2 \sum_{i=1}^K \text{tr}(G_i^t \bar{\Sigma}_i^t) + \frac{\gamma_t}{2} \sum_{k=1}^K \|\boldsymbol{\mu}_k^*\|_2^2 \end{aligned} \quad (132)$$

In addition, we have that

$$\begin{aligned} &\frac{1}{\beta_{t+1}} \mathbb{E} \|\boldsymbol{\mu}_k^{t+1} - \boldsymbol{\mu}_k^*\|_{\Sigma_k^{t+1-1}}^2 - \frac{1}{\beta_t} \mathbb{E} \|\boldsymbol{\mu}_k^{t+1} - \boldsymbol{\mu}_k^*\|_{\Sigma_k^{t-1}}^2 - \gamma_t \mathbb{E} \|\boldsymbol{\mu}_i^t - \boldsymbol{\mu}_i^*\|_2^2 \\ &= \frac{1}{\beta_{t+1}} \mathbb{E} \langle \Sigma_k^{t+1-1} (\boldsymbol{\mu}_k^{t+1} - \boldsymbol{\mu}_k^*), \boldsymbol{\mu}_k^{t+1} - \boldsymbol{\mu}_k^* \rangle - \frac{1}{\beta_t} \mathbb{E} \langle \Sigma_k^{t-1} (\boldsymbol{\mu}_k^{t+1} - \boldsymbol{\mu}_k^*), \boldsymbol{\mu}_k^{t+1} - \boldsymbol{\mu}_k^* \rangle - \gamma_t \mathbb{E} \|\boldsymbol{\mu}_i^t - \boldsymbol{\mu}_i^*\|_2^2 \end{aligned} \quad (133)$$

$$= \mathbb{E} \langle (\frac{1}{\beta_{t+1}} \Sigma_k^{t+1-1} - \frac{1}{\beta_t} \Sigma_k^{t-1} - \gamma_t \mathbf{I}) (\boldsymbol{\mu}_k^{t+1} - \boldsymbol{\mu}_k^*), \boldsymbol{\mu}_k^{t+1} - \boldsymbol{\mu}_k^* \rangle \quad (134)$$

Note that  $\Sigma_k^{t+1-1} = \omega_t \Sigma_k^{t-1} + \alpha_t \hat{G}_k^t = \Sigma_k^{t-1} (\omega_t \mathbf{I} + \alpha_t \mathbf{H}_k^t) \Sigma_k^{t-1}$ , we have that

$$\frac{1}{\beta_{t+1}} \Sigma_k^{t+1-1} - \frac{1}{\beta_t} \Sigma_k^{t-1} - \gamma_t \mathbf{I} = \frac{1}{\beta_{t+1}} \Sigma_k^{t-1} (\omega_t \mathbf{I} + \alpha_t \mathbf{H}_k^t) \Sigma_k^{t-1} - \frac{1}{\beta_t} \Sigma_k^{t-1} - \gamma_t \mathbf{I} \quad (135)$$

$$= \Sigma_k^{t-1} (\frac{\omega_t}{\beta_{t+1}} \mathbf{I} + \frac{\alpha_t}{\beta_{t+1}} \mathbf{H}_k^t - \frac{1}{\beta_t} \mathbf{I} - \gamma_t \Sigma_k^t) \Sigma_k^{t-1} \quad (136)$$

Because of  $\mathbf{H}_k^t \succeq \frac{1}{\alpha_t} \left( \frac{\beta_{t+1}}{\beta_t} - \omega_t \right) \mathbf{I} + \frac{\beta_{t+1} \gamma_t}{\alpha_t} \Sigma_k^t$  in algorithm 2, we have that

$$\frac{1}{\beta_{t+1}} \Sigma_k^{t+1} - \frac{1}{\beta_t} \Sigma_k^{t-1} - \gamma_t \mathbf{I} = \Sigma_k^t \text{pre}^{-\frac{1}{2}} \left( \frac{\omega_t}{\beta_{t+1}} \mathbf{I} + \frac{\alpha_t}{\beta_{t+1}} \mathbf{H}_k^t - \frac{1}{\beta_t} \mathbf{I} - \gamma_t \Sigma_k^t \right) \Sigma_k^t \text{pre}^{-\frac{1}{2}} \quad (137)$$

$$\succeq \Sigma_k^t \text{pre}^{-\frac{1}{2}} \left( \frac{\omega_t}{\beta_{t+1}} \mathbf{I} + \frac{1}{\beta_{t+1}} \left( \frac{\beta_{t+1}}{\beta_t} - \omega_t \right) \mathbf{I} + \gamma_t \Sigma_k^t - \frac{1}{\beta_t} \mathbf{I} - \gamma_t \Sigma_k^t \right) \Sigma_k^t \text{pre}^{-\frac{1}{2}} \quad (138)$$

$$\succeq \mathbf{0} \quad (139)$$

Plug Eq.(139) into Eq.(134), we know that

$$\frac{1}{\beta_{t+1}} \mathbb{E} \|\boldsymbol{\mu}_k^{t+1} - \boldsymbol{\mu}_k^*\|_{\Sigma_k^{t+1}}^2 - \frac{1}{\beta_t} \mathbb{E} \|\boldsymbol{\mu}_k^{t+1} - \boldsymbol{\mu}_k^*\|_{\Sigma_k^{t-1}}^2 - \gamma_t \mathbb{E} \|\boldsymbol{\mu}_i^t - \boldsymbol{\mu}_i^*\|_2^2 \leq 0 \quad (140)$$

Telescope with Eq.(132), we have that

$$\begin{aligned} & \sum_{t=1}^T \sum_{i=1}^K J_i(\bar{\boldsymbol{\mu}}_i^t, \bar{\Sigma}_i^t) - \sum_{t=1}^T \sum_{i=1}^K J_i(\bar{\boldsymbol{\mu}}_i^*, 0) \\ & \leq \frac{\sum_{k=1}^K \mathbb{E} \|\boldsymbol{\mu}_k^1 - \boldsymbol{\mu}_k^*\|_{\Sigma_k^{1-1}}^2}{2\beta_1} - \frac{\sum_{k=1}^K \mathbb{E} \|\boldsymbol{\mu}_k^{T+1} - \boldsymbol{\mu}_k^*\|_{\Sigma_k^{T-1}}^2}{2\beta_T} + \sum_{t=1}^T \beta_t \sum_{k=1}^K \mathbb{E} \|\Sigma_k^t (\sum_{i=k}^K \hat{g}_{ik})\|_{\Sigma_k^t}^2 \\ & + 2 \sum_{t=1}^T \sum_{i=1}^K \text{tr}(G_i^t \bar{\Sigma}_i^t) + \frac{1}{2} \sum_{t=1}^T \gamma_t \sum_{k=1}^K \|\boldsymbol{\mu}_k^*\|_2^2 \end{aligned} \quad (141)$$

We now show the upper bound of term  $\sum_{t=1}^T \beta_t \sum_{k=1}^K \mathbb{E} \|\Sigma_k^t (\sum_{i=k}^K \hat{g}_{ik})\|_{\Sigma_k^t}^2$ .

Note that  $\beta_t = t\beta$ , together with Lemma C.2 (c), we know that

$$\sum_{t=1}^T \beta_t \sum_{k=1}^K \mathbb{E} \|\Sigma_k^t (\sum_{i=k}^K \hat{g}_{ik})\|_{\Sigma_k^t}^2 \leq \sum_{t=1}^T t \frac{C_1}{t^{\frac{3}{2}}} \leq 2\sqrt{T+1}C_1 \quad (142)$$

where  $C_1 = \frac{3\beta \sum_{i=1}^K KL_i^2 (id+1)^2}{2\alpha\nu}$ .

We now show the upper bound of term  $2 \sum_{t=1}^T \sum_{i=1}^K \text{tr}(G_i^t \bar{\Sigma}_i^t)$ .

From Lemma C.3, we know that

$$2 \sum_{t=1}^T \sum_{i=1}^K \text{tr}(G_i^t \bar{\Sigma}_i^t) \leq 2 \sum_{t=1}^T \sum_{i=1}^K \frac{L_i id}{2t^{\frac{3}{4}}} \sqrt{\frac{3}{\alpha\nu}} \quad (143)$$

$$= \sum_{t=1}^T C_2 \frac{1}{t^{\frac{3}{4}}} \quad (144)$$

$$\leq 4(T+1)^{\frac{1}{4}} C_2 \quad (145)$$

where  $C_2 = \frac{\sum_{i=1}^K \sqrt{3idL_i}}{\sqrt{\alpha\nu}}$ .

We now show the upper bound of the term  $\frac{1}{2} \sum_{t=1}^T \gamma_t \sum_{k=1}^K \|\boldsymbol{\mu}_k^*\|_2^2$ .

Note that the optimal solution is bounded from Assumption 5.3, i.e.,  $\sum_{k=1}^K \|\boldsymbol{\mu}_k^*\|_2^2 \leq B$ . Together with the setting  $\gamma_t = \frac{\alpha\nu}{\beta\sqrt{t+1}}$ , we can achieve that

$$\frac{1}{2} \sum_{t=1}^T \gamma_t \sum_{k=1}^K \|\boldsymbol{\mu}_k^*\|_2^2 \leq \frac{B}{2} \sum_{t=1}^T \gamma_t \leq \frac{\alpha\nu B}{2\beta} 2\sqrt{T+2} = \sqrt{T+2}C_3 \quad (146)$$

where  $C_3 = \frac{\alpha\nu B}{\beta}$



Plug Eq.(142) and Eq.(145) into Eq.(141), we have that

$$\begin{aligned} & \sum_{t=1}^T \sum_{i=1}^K J_i(\bar{\boldsymbol{\mu}}_i^t, \bar{\boldsymbol{\Sigma}}_i^t) - \sum_{t=1}^T \sum_{i=1}^K J_i(\bar{\boldsymbol{\mu}}_i^*, 0) \\ & \leq \frac{\sum_{k=1}^K \mathbb{E} \|\boldsymbol{\mu}_k^1 - \boldsymbol{\mu}_k^*\|_{\boldsymbol{\Sigma}_k^{1-1}}^2}{2\beta_1} - \frac{\sum_{k=1}^K \mathbb{E} \|\boldsymbol{\mu}_k^{T+1} - \boldsymbol{\mu}_k^*\|_{\boldsymbol{\Sigma}_k^{T-1}}^2}{2\beta_T} + 2\sqrt{T+1}C_1 + 4(T+1)^{\frac{1}{4}}C_2 + \sqrt{T+2}C_3 \end{aligned} \quad (147)$$

From Lemma C.5, we then have that

$$\begin{aligned} \sum_{t=1}^T \left( \sum_{k=1}^K (f_k(\bar{\boldsymbol{\mu}}_k^t) - f_k(\bar{\boldsymbol{\mu}}_k^*)) \right) & \leq \sum_{t=1}^T \sum_{i=1}^K J_i(\bar{\boldsymbol{\mu}}_i^t, \bar{\boldsymbol{\Sigma}}_i^t) - \sum_{t=1}^T \sum_{i=1}^K J_i(\bar{\boldsymbol{\mu}}_i^*, 0) \\ & \leq \frac{\sum_{k=1}^K \mathbb{E} \|\boldsymbol{\mu}_k^1 - \boldsymbol{\mu}_k^*\|_{\boldsymbol{\Sigma}_k^{1-1}}^2}{2\beta_1} - \frac{\sum_{k=1}^K \mathbb{E} \|\boldsymbol{\mu}_k^{T+1} - \boldsymbol{\mu}_k^*\|_{\boldsymbol{\Sigma}_k^{T-1}}^2}{2\beta_T} \\ & \quad + 2\sqrt{T+1}C_1 + 4(T+1)^{\frac{1}{4}}C_2 + \sqrt{T+2}C_3 \end{aligned} \quad (148)$$

Finally, divide  $T$  on both sides of Eq.(149), we have that

$$\frac{1}{T} \sum_{t=1}^T \sum_{k=1}^K f_k(\bar{\boldsymbol{\mu}}_k^t) - \sum_{k=1}^K f_k(\bar{\boldsymbol{\mu}}_k^*) \leq \frac{\sum_{k=1}^K \|\boldsymbol{\mu}_k^1 - \boldsymbol{\mu}_k^*\|_{\boldsymbol{\Sigma}_k^{1-1}}^2}{2\beta T} + \frac{2\sqrt{T+1}C_1}{T} + \frac{4(T+1)^{\frac{1}{4}}C_2}{T} + \frac{\sqrt{T+2}C_3}{T} \quad (150)$$

□

**Theorem D.2.** Set  $\beta_t = t\beta$  with  $\beta > 0$ ,  $\alpha_t = \sqrt{t+1}\alpha$  with  $\alpha > 0$ , and  $\gamma_t = \frac{\alpha\nu}{\beta\sqrt{t+1}}$ , and  $\nu > 0$ , and  $\omega_t = 1$ . Then, during the running process of Algorithm 2, the constraints  $\mathbf{H}_k^t \preceq \frac{1}{\alpha_t} \left( \frac{\beta_{t+1}}{\beta_t} - \omega_t \right) \mathbf{I} + \frac{\beta_{t+1}\gamma_t}{\alpha_t} \boldsymbol{\Sigma}_k^t$  and  $\nu \mathbf{I} \preceq \hat{\mathbf{G}}_k^t = \boldsymbol{\Sigma}_k^t^{-\frac{1}{2}} \mathbf{H}_k^t \boldsymbol{\Sigma}_k^t^{-\frac{1}{2}}$  always have feasible solutions.

*Proof.* To ensure the constraints  $\mathbf{H}_k^t \preceq \frac{1}{\alpha_t} \left( \frac{\beta_{t+1}}{\beta_t} - \omega_t \right) \mathbf{I} + \frac{\beta_{t+1}\gamma_t}{\alpha_t} \boldsymbol{\Sigma}_k^t$  and  $\nu \mathbf{I} \preceq \hat{\mathbf{G}}_k^t = \boldsymbol{\Sigma}_k^t^{-\frac{1}{2}} \mathbf{H}_k^t \boldsymbol{\Sigma}_k^t^{-\frac{1}{2}}$  always feasible during the algorithm, it is equivalent to show the constraints Eq.(151) always has feasible solutions  $\mathbf{H}_k^t$ .

$$\nu \boldsymbol{\Sigma}_k^t \preceq \mathbf{H}_k^t \preceq \frac{1}{\alpha_t} \left( \frac{\beta_{t+1}}{\beta_t} - \omega_t \right) \mathbf{I} + \frac{\beta_{t+1}\gamma_t}{\alpha_t} \boldsymbol{\Sigma}_k^t \quad (151)$$

It is equivalent to show that the inequality (152) always holds true.

$$\nu \boldsymbol{\Sigma}_k^t \preceq \frac{1}{\alpha_t} \left( \frac{\beta_{t+1}}{\beta_t} - \omega_t \right) \mathbf{I} + \frac{\beta_{t+1}\gamma_t}{\alpha_t} \boldsymbol{\Sigma}_k^t \quad (152)$$

Then, it is equivalent to show that the inequality (153) always holds true.

$$\left( \nu - \frac{\beta_{t+1}\gamma_t}{\alpha_t} \right) \boldsymbol{\Sigma}_k^t \preceq \frac{1}{\alpha_t} \left( \frac{\beta_{t+1}}{\beta_t} - \omega_t \right) \mathbf{I} \quad (153)$$

We first check the left hand side of the inequality (153).

Note that the setting  $\beta_t = t\beta$  with  $\beta > 0$ ,  $\alpha_t = \sqrt{t+1}\alpha$  with  $\alpha > 0$ , and  $\gamma_t = \frac{\alpha\nu}{\beta\sqrt{t+1}}$ , and  $\nu > 0$ , and  $\omega_t = 1$ . We can achieve that

$$\left( \nu - \frac{\beta_{t+1}\gamma_t}{\alpha_t} \right) \boldsymbol{\Sigma}_k^t = \left( \nu - \frac{(t+1)\beta}{\sqrt{t+1}\alpha} \frac{\alpha\nu}{\beta\sqrt{t+1}} \right) \boldsymbol{\Sigma}_k^t \quad (154)$$

$$= (\nu - \nu) \boldsymbol{\Sigma}_k^t = 0 \quad (155)$$

We now check the right hand side of the inequality (153).

$$\frac{1}{\alpha_t} \left( \frac{\beta_{t+1}}{\beta_t} - \omega_t \right) = \frac{1}{\alpha_t} \left( \frac{(t+1)\beta}{t\beta} - 1 \right) = \frac{1}{\alpha_t t} > 0 \quad (156)$$

Thus, the inequality (153) always hold true. As a result, the constraints set always have feasible solutions. □

## E DETAILS OF DIFFUSION TARGET GENERATION EXPERIMENTS

### E.1 BLACK-BOX TARGET FUNCTION

CLIP (Contrastive Language-Image Pre-Training) (Radford et al., 2021) is trained on various (image, text) pairs to capture their similarity. Let the generated image be  $x$ , the target text be  $y$ , and the image and text encoders of CLIP be  $\phi_{\text{image}}$  and  $\phi_{\text{text}}$ , respectively. These encoders share an aligned embedding space, allowing us to measure similarity using cosine distance. The black-box function is defined as:

$$f(x; y) = \frac{\phi_{\text{image}}(x)^T \phi_{\text{text}}(y)}{\|\phi_{\text{image}}(x)\| \|\phi_{\text{text}}(y)\|}$$

We then employ the normalized cosine distance  $\frac{y - \mu_y}{\sigma_y}$  as the black-box target score, where  $\mu_y$  and  $\sigma_y$  denotes the mean and standard deviation of the cosine distances between the target and the images in the dataset.

For our implementation, we choose to use the pre-trained publicly available CLIP model ViT-L/14<sup>3</sup>.

### E.2 SAMPLER: DPM-SOLVER++

We choose to use DPM-Solver++ (Lu et al., 2022a) as our sampler for all experiments in both the training and evaluation phases. We use the 2nd-order SDE solver and the data evaluation formulation. For all experiments, we use  $K = 14$  sampling steps.

### E.3 DATASET: CELEBA-HQ

We use dataset CelebA-HQ in our experiments, it contains 30,000 face images. It is public available at <https://huggingface.co/datasets/hugan/CelebA-HQ>

### E.4 BASELINE: CLASSIFIER-FREE CONDITIONAL DIFFUSION MODEL

Denotes the dataset as  $\mathcal{X} = \{x_1, x_2, \dots, x_n\}$ . We evaluate the black-box function target score value, that is  $\mathcal{Y} = \{f(x_i) \mid \forall x_i \in \mathcal{X}\}$ , and denotes the mean and standard deviation as  $\mu_y$  and  $\sigma_y$ . We then obtain the normalized target score value as  $\tilde{\mathcal{Y}} = \{\frac{y_i - \mu_y}{\sigma_y} \mid \forall y_i \in \mathcal{Y}\}$ .

We use the dataset and normalized score value  $\mathcal{X} \times \tilde{\mathcal{Y}}$  to train a classifier-free conditional diffusion model  $\hat{\mu}_\phi(x, t, y)$ .

An unconditional diffusion model trained on the CelebA-HQ dataset is provided by the authors of latent diffusion<sup>4</sup>. In practice, instead of training the model from scratch, we load the weights of the unconditional layers from the pre-trained model and continue training from there. This approach can speed up the training time.

<sup>3</sup><https://huggingface.co/sentence-transformers/clip-ViT-L-14>

<sup>4</sup><https://ommer-lab.com/files/latent-diffusion/celeba.zip>

### E.5 BASELINE: DDOM

The recent work (Krishnamoorthy et al., 2023) introduces Denoising Diffusion Optimization Models (DDOM) for solving offline black-box optimization tasks using diffusion models. This method can also be naturally extended to black-box targeted generation tasks.

During the pre-processing phase, the black-box function score  $y_i = f(x_i)$  is evaluated for each sample  $x_i \in \mathcal{D}$ . The offline dataset  $\mathcal{D}$  is then partitioned into  $N_B$  bins of equal width based on  $y$ . Each bin is assigned a weight proportional to both the number of points in the bin and the average score value of the bin. Specifically, the weight of the  $i$ -th is given by:

$$w_i = \frac{|B_i|}{|B_i| + C} \exp\left(\frac{-|\hat{y} - y_{b_i}|}{\tau}\right) \quad (157)$$

where  $\hat{y}$  is the best function value in the offline dataset  $\mathcal{D}$ ,  $|B_i|$  denotes the number of points in the  $i$ -th bin, and  $y_{b_i}$  is the midpoint of the interval corresponding to the bin  $B_i$ . The parameters  $K$  and  $\tau$  are hyper-parameters.

During training, this weight is used to compute the weighted loss, which is given by:

$$\mathbb{E}_t \left[ \lambda(t) \mathbb{E}_{\mathbf{x}_0, y} \left[ w(y) \mathbb{E}_{\mathbf{x}_t | \mathbf{x}_0} \left[ \|\epsilon_\theta(\mathbf{x}_t, t, y) - \nabla_{\mathbf{x}} \log p_t(\mathbf{x}_t | \mathbf{x}_0)\|_2^2 \right] \right] \right] \quad (158)$$

where  $w(y) = w_i$  if  $y \in B_i$ . In practice, the score values are normalized to fit a standard normal distribution to ensure that the  $y_i$  values are well behaved.

The authors demonstrate DDOM for solving black-box optimization (BBO) tasks. Since DDOM uses a diffusion model, it can also be naturally implemented for image generation tasks. In our experiment, we implement the DDOM for black-box targeted image generation tasks. In the weight function defined in equation (157), we set the hyper-parameters  $C = 0.01$ ,  $\tau = 0.1$ ,  $N_B = 64$ .

### E.6 OUR METHOD: FINE-TUNE THE CLASSIFIER-FREE CONDITIONAL DIFFUSION MODEL

We use our method to improve the baseline conditional diffusion model. Instead of fine-tune the parameters for all time steps, we only apply it for the second half of the time steps. The rationale is that we should place more emphasis on the time steps closer to the final output image  $x_K$ . That is, we have the fine-tuning parameter sets  $\theta_k = \{\mu_k, \Sigma_k\}$  for  $k \in \{1, \dots, K\}$ .

We revise the sequential black-box objective  $\sum_{k=1}^K f_k(\bar{\epsilon}_k)$ . In our previous experiments, we set  $f_k(\bar{\epsilon}_k) = F(\mathbf{x}_k)$  for  $k \in \{1, \dots, K\}$ , which call CLIP to evaluate the input (noised) image  $\mathbf{x}_k$  at each diffusion sampling step  $k$ . This scheme is not effective enough because we care more about the generated image at the last step, i.e.,  $\mathbf{x}_K$ , than the generated image at the inner steps. Thus, we set the black-box function at the inner step as  $f_k(\bar{\epsilon}_k) = F_k(\mathbf{x}_k)$  for  $k \in \{1, \dots, K-1\}$ . The black-box function  $F_k(\mathbf{x}_k)$  takes the noised image  $\mathbf{x}_k$  at the inner step as input and performs a deterministic sampling process  $\mathbf{x}_{k+1} = \hat{\mu}_\phi(\mathbf{x}_k, k) + \tilde{\sigma}_{k+1} \boldsymbol{\mu}_{k+1}$  for the future steps to achieve a predicted  $\mathbf{x}_K$ , and call CLIP to evaluate the predicted  $\mathbf{x}_K$ .