

# Solving a Class of Non-Convex Min-Max Games using Adaptive Momentum Methods\*

\*Babak Barazandeh  
bbarazandeh@splunk.com

†Davoud Ataee Tarzanagh  
tarzanagh@ufl.edu

‡George Michailidis  
gmichail@ufl.edu

\*Splunk, †University of Florida

## Abstract

Adaptive momentum methods have recently attracted a lot of attention for training of deep neural networks. They use an exponential moving average of past gradients of the objective function to update both search directions and learning rates. However, these methods are not suited for solving min-max optimization problems that arise in training generative adversarial networks. In this paper, we propose an adaptive momentum min-max algorithm that generalizes adaptive momentum methods to the non-convex min-max regime. Further, we establish non-asymptotic rates of convergence for the proposed algorithm when used in a reasonably broad class of non-convex min-max optimization problems. Experimental results illustrate its superior performance vis-a-vis benchmark methods for solving such problems.

**Keywords**— Non-convex min-max games, First-order Nash equilibrium, Adaptive optimization

## 1 Introduction

Stochastic first-order methods are of core practical importance for solving numerous optimization problems including training deep neural networks (DNN). Standard stochastic gradient descent (SGD) has become a widely used technique for the latter task. However, its convergence crucially depends on the tuning and update of the learning rate over iterations in order to control the variance of the gradient in the stochastic search directions, especially for non-convex functions [1].

To alleviate these issues, several improved variants of SGD that automatically update the search directions and learning rates using a metric constructed from the history of iterates have been proposed, including adaptive methods [2, 3, 4, 5] and adaptive momentum methods [6, 7]. In particular, ADAM belonging to the second category enjoys the dual advantages of variance adaption and momentum direction [8, 9] and hence represents a popular algorithm to train DNNs.

There is a large body of literature on the theoretical and empirical benefits of adaptive momentum optimization algorithms for convex [6, 7], smooth non-convex [10, 11, 12], and non-smooth non-convex settings [13]. [14] gives an analysis of an optimistic adaptive method that uses ADAGRAD [4, 5] for non-convex min-max optimization. However, ADAGRAD-type methods are suited for sparse convex settings and their performance deteriorates in (dense) non-convex optimization problems [11]. These empirical findings necessitate the use of adaptive momentum methods that incorporate knowledge of past iterations. It is important to notice that

---

\*This arXiv submission includes the details of the proofs for the paper accepted for publication in the proceeding of the 46<sup>th</sup> International Conference on Acoustics, Speech, and Signal Processing (ICASSP).

all these methods are designed for classical minimization problems. However, training DNNs such as Generative Adversarial Networks (GANs) require solving a general class of min-max optimization problems [15, 16] which due to its difficulty, keeps other generative models attractive [17, 18].

The goal of this paper is to generalize adaptive momentum methods to solve a general class of *non-convex-non-concave min-max problems*. It develops an adaptive algorithm for solving min-max saddle point games and theoretically analyzes its convergence rate. The performance of the developed algorithm is assessed on training GANs.

The remainder of the paper is organized as follows. Section 2 provides the formulation of the min-max problem, Section 3 describes the proposed algorithm and Section 4 investigates its convergence properties. Finally, Section 5 provides numerical results for training GANs.

## 2 Formulation of the Min-Max Optimization Problem

Consider the *stochastic* min-max saddle point problem

$$\min_{\boldsymbol{\theta}} \max_{\boldsymbol{\alpha}} F(\boldsymbol{\theta}, \boldsymbol{\alpha}) = \mathbb{E}_{\boldsymbol{\xi} \sim \mathcal{D}} [f(\boldsymbol{\theta}, \boldsymbol{\alpha}; \boldsymbol{\xi})], \quad (1)$$

where  $\boldsymbol{\theta} \in \mathbb{R}^{p_1}$ ,  $\boldsymbol{\alpha} \in \mathbb{R}^{p_2}$ ,  $\boldsymbol{\xi}$  is a random variable drawn from an unknown distribution  $\mathcal{D}$ , and  $F(\boldsymbol{\theta}, \boldsymbol{\alpha})$  is a non-convex-non-concave function, i.e., it is non-convex in  $\boldsymbol{\theta}$  for any given  $\boldsymbol{\alpha}$  and is non-concave in  $\boldsymbol{\alpha}$  for any given  $\boldsymbol{\theta}$ .

Next, we introduce necessary notation and definitions. Throughout,  $\mathbf{y} := (\boldsymbol{\theta}, \boldsymbol{\alpha}) \in \mathbb{R}^{p_1} \times \mathbb{R}^{p_2}$  and denote the objective function of Game (1) and its random realization by  $F(\mathbf{y})$  and  $f(\mathbf{y}; \boldsymbol{\xi})$ , respectively. Further,  $\nabla F(\mathbf{y}) = [\nabla_{\boldsymbol{\theta}} F(\boldsymbol{\theta}, \boldsymbol{\alpha}), -\nabla_{\boldsymbol{\alpha}} F(\boldsymbol{\theta}, \boldsymbol{\alpha})]$  and  $\nabla f(\mathbf{y}; \boldsymbol{\xi}) = [\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}, \boldsymbol{\alpha}; \boldsymbol{\xi}), -\nabla_{\boldsymbol{\alpha}} f(\boldsymbol{\theta}, \boldsymbol{\alpha}; \boldsymbol{\xi})]$  denotes the corresponding gradient and stochastic gradient of the objective function, respectively.

**Definition 1** (Nash Equilibrium). *A point  $(\boldsymbol{\theta}^*, \boldsymbol{\alpha}^*) \in \mathbb{R}^{p_1} \times \mathbb{R}^{p_2}$  is a Nash equilibrium of Game (1) if*

$$F(\boldsymbol{\theta}^*, \boldsymbol{\alpha}) \leq F(\boldsymbol{\theta}^*, \boldsymbol{\alpha}^*) \leq F(\boldsymbol{\theta}, \boldsymbol{\alpha}^*), \quad \forall (\boldsymbol{\theta}, \boldsymbol{\alpha}) \in \mathbb{R}^{p_1} \times \mathbb{R}^{p_2}.$$

This definition implies that  $\boldsymbol{\theta}^*$  is a global minimum of  $F(\cdot, \boldsymbol{\alpha}^*)$  and  $\boldsymbol{\alpha}^*$  is a global maximum of  $F(\boldsymbol{\theta}^*, \cdot)$ . In the convex-concave regime with  $F(\boldsymbol{\theta}, \boldsymbol{\alpha})$  being convex in  $\boldsymbol{\theta}$  for any given  $\boldsymbol{\alpha}$  and concave in  $\boldsymbol{\alpha}$  for any given  $\boldsymbol{\theta}$ , the Nash equilibrium always exists [19] and there are several algorithms for identifying it [20, 21]. However, computing a Nash equilibrium point is NP-hard in general [22, 19], and it may not even exist [23]. As a result, since we are considering the general non-convex-non-concave regime, we settle in computing a *first-order Nash equilibrium* point [24, 25] defined next.

**Definition 2** (First-Order Nash Equilibrium (FNE)). *A point  $\mathbf{y}^* \in \mathbb{R}^{p_1} \times \mathbb{R}^{p_2}$  is a first-order Nash equilibrium point of Game (1), if  $\nabla F(\mathbf{y}^*) = 0$ .*

Note that at a FNE point, each player satisfies the first-order optimality condition of its own objective function when the strategy of the other player is fixed [26, 27]. In practice, iterative algorithms are used for computing a FNE for a stochastic problem. As a result, the performance of different iterative algorithms are evaluated based on the following *approximate* stochastic FNE definition.

**Definition 3** ( $\epsilon$ -Stochastic First-Order Nash Equilibrium (SFNE)). *A random variable  $\mathbf{y}^*$  is an approximate SFNE ( $\epsilon$ -SFNE) point of Game (1) if  $\mathbb{E} [\|\nabla F(\mathbf{y}^*)\|^2] \leq \epsilon^2$ , where the expectation is taken over the distribution of the random variable  $\mathbf{y}^*$ .*

The randomness of variable  $\mathbf{y}^*$  in Definition 3 comes from the use of iterative algorithms that have access to stochastic gradients of the objective function (see, e.g., Algorithm ADAM<sup>3</sup> below). The objective of this work is to find an  $\epsilon$ -SFNE point for Game (1) using an iterative method based on adaptive momentum.

### 3 THE ADAM<sup>3</sup> Algorithm

The proposed ADaptive Momentum Min-Max (ADAM<sup>3</sup>) algorithm comes with convergence guarantees for solving a general class of *non-convex-non-concave* saddle point games defined in (1). It is obtained by integrating AMSGRAD [7], a modified version of ADAM [6], with a stochastic extra-gradient method [28]. As seen in Algorithm 1, ADAM<sup>3</sup> generates two sequences  $\mathbf{x}_k$  and  $\mathbf{z}_k$ , where  $\mathbf{x}_k$  is an ancillary sequence

---

**Algorithm 1:** ADaptive Momentum Min-Max (ADAM<sup>3</sup>)

---

**Input** :  $\{\beta_{1,k}\}_{k=1}^N, \beta_2, \beta_3 \in [0, 1), m \in \mathbb{N}$ , and  $\eta \in \mathbb{R}_+$ ;

Initialize  $\mathbf{z}_0 = \mathbf{x}_0 = \mathbf{m}_0 = \mathbf{v}_0 = \mathbf{d}_0 = \mathbf{0}_d$ .

**for**  $k = 1 : N$  **do**

$\mathbf{z}_k = \mathbf{x}_{k-1} - \eta \mathbf{d}_{k-1}$ ;  
 Draw  $\boldsymbol{\xi}_k = (\boldsymbol{\xi}_k^1, \dots, \boldsymbol{\xi}_k^m)$  from  $\mathcal{D}$ , and set  $\widehat{\mathbf{g}}_k = \frac{1}{m} \sum_{i=1}^m \nabla f(\mathbf{z}_k; \boldsymbol{\xi}_k^i)$ ;  
 $\mathbf{m}_k = \beta_{1,k} \mathbf{m}_{k-1} + (1 - \beta_{1,k}) \widehat{\mathbf{g}}_k$ ;  
 $\mathbf{v}_k = \beta_2 \mathbf{v}_{k-1} + (1 - \beta_2) \widehat{\mathbf{g}}_k \odot \widehat{\mathbf{g}}_k$ ;  
 $\tilde{\mathbf{v}}_k = \beta_3 \tilde{\mathbf{v}}_{k-1} + (1 - \beta_3) \max(\tilde{\mathbf{v}}_{k-1}, \mathbf{v}_k)$ ;  
 $\mathbf{d}_k = \tilde{\mathbf{v}}_k^{-\frac{1}{2}} \odot \mathbf{m}_k$ ;  
 $\mathbf{x}_k = \mathbf{x}_{k-1} - \eta \mathbf{d}_k$ ;

**end**

$\odot$ : Element-wise vector multiplication

---

and the stochastic gradient is only computed over the sequence of  $\mathbf{z}_k$ 's using a mini-batch of size  $m$ , i.e.,  $\widehat{\mathbf{g}}_k = 1/m \sum_{i=1}^m \nabla f(\mathbf{z}_k; \boldsymbol{\xi}_k^i)$ . Using a mini-batch for estimating the gradient is a commonly used approach and more details are available in [29] and references therein. After estimating the gradient, the algorithm calculates the momentum direction,  $\mathbf{m}_k$ , as an exponential moving average of the past gradients. Then,  $\mathbf{m}_k$  is adaptively scaled by the square root of the exponential moving average of squared past gradients  $\tilde{\mathbf{v}}_k$ .

The following remarks about ADAM<sup>3</sup> are in order:

- (1) The square and the maximum operators are applied element-wise. In some applications, to prevent division by zero, we may add a small positive constant  $\epsilon$  to  $\mathbf{v}_k$  [12]. Further, a mini-batch of size  $m$  is used in each iteration to estimate the gradient's value.
- (2) ADAM<sup>3</sup> computes adaptive learning rates from estimates of the second moments of the gradients, similar to [12]. In particular, it uses a larger learning rate compared to AMSGRAD and yet incorporates the intuition of slowly decaying the effect of previous gradients on the learning rate. The decay parameter  $\beta_3$  is an important component of ADAM<sup>3</sup>, that enables establishing its convergence properties similar to AMSGRAD ( $\beta_3 = 0$ ), while maintaining the efficiency of ADAM.

### 4 Convergence Analysis

We start by positing the following assumptions:

**Assumption A.** For all  $\mathbf{x} \in \mathbb{R}^d$ ,

1.  $\mathbb{E}_{\boldsymbol{\xi} \sim \mathcal{D}}[\nabla f(\mathbf{x}, \boldsymbol{\xi})] = \nabla F(\mathbf{x})$ .
2. The function  $f(\mathbf{x}, \boldsymbol{\xi})$  has a  $G_\infty$ -bounded gradient, i.e.,  $\forall \boldsymbol{\xi} \sim \mathcal{D}$ , it holds that  $\|\nabla f(\mathbf{x}, \boldsymbol{\xi})\|_\infty \leq G_\infty < \infty$ .
3. The function  $F$  has bounded variance, i.e.,

$$\mathbb{E}_{\boldsymbol{\xi} \sim \mathcal{D}} [\|\nabla f(\mathbf{x}, \boldsymbol{\xi}) - \nabla F(\mathbf{x})\|^2] = \sigma^2 < \infty.$$

The above assumptions are fairly standard in the non-convex optimization literature [1, 30]. Further, Assumption A(2) is slightly stronger than the assumption  $\|\nabla f(\mathbf{x}, \boldsymbol{\xi})\| \leq G_2$  that is commonly used in the analysis of stochastic gradient descent. However, Assumption A(2) is crucial for the convergence analysis of adaptive methods [11, 14, 12, 13].

**Assumption B** (Lipschitz Gradient). *The function  $F$  is  $L$ -smooth, i.e.,*

$$\|\nabla F(\mathbf{x}) - \nabla F(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|, \quad \text{for all } \mathbf{x}, \mathbf{y} \in \mathbb{R}^d.$$

The above assumption is standard and commonly used in the optimization literature [31, 32].

**Assumption C** (Minty condition). *There exists  $\mathbf{x}_* \in \mathbb{R}^d$  such that for any  $\mathbf{x} \in \mathbb{R}^d$  we have*

$$\langle \mathbf{x} - \mathbf{x}_*, \nabla F(\mathbf{x}) \rangle \geq 0.$$

As explained in [33, 14] and references therein, the Minty condition is a commonly used assumption in the literature for analyzing non-convex min-max games and is weaker than other benchmark assumptions such as pseudo-monotonicity or monotonicity [34].

**Assumption D.** *For the point  $\mathbf{x}^*$  satisfying the Minty condition and all iterates  $k$  generated by Algorithm 1, we have  $\|\mathbf{x}_*\| \leq \frac{D}{2}$  and  $\|\mathbf{x}_k\| \leq \frac{D}{2}$ .*

This assumption is required in the analysis of min-max saddle point games and has been used in [14, 35]. This assumption holds true in the training process of DNNs that have normalization layers in their structure [36, 37, 38].

**Assumption E.** *In Algorithm 1,  $G_0^2 \leq \|\tilde{\mathbf{v}}_0\|_\infty$ .*

This assumption is required in the analysis of adaptive methods [13, 12] and can be easily satisfied in the initialization step of the proposed algorithm. Next, we introduce lemmas used to establish the main result.

**Lemma 1.** *[[39], Lemma 4.2] Let Assumption A (2) hold. Then, in Algorithm 1 we have  $\|\mathbf{m}_k\|_\infty \leq G_\infty$  and  $\|\tilde{\mathbf{v}}_k\|_\infty \leq G_\infty^2$  for all  $k \in \{1 \cdots N\}$ .*

**Lemma 2.** *Assume that  $\gamma := \beta_{1,1}/\beta_2 \leq 1$  in Algorithm 1. Then, for each  $k \in \{1 \cdots N\}$  we have*

$$\|\tilde{\mathbf{v}}_k^{-\frac{1}{2}} \circ \mathbf{m}_{k-1}\| \leq \sqrt{\frac{d}{u_c}},$$

where  $u_c := (1 - \beta_3)(1 - \beta_{1,1})(1 - \beta_2)(1 - \gamma)$ .

The following Theorem 1 establishes the main result by providing an upper bound for the average norm of the gradient of the objective function.

**Theorem 1.** *Let Assumptions A–E hold, and  $L, G_\infty, G_0, \sigma$  be defined therein. In Algorithm 1, if we choose*

$$\eta \leq \sqrt{G_0^3/(56L^2G_\infty)}, \beta_{1,k} = \beta_{1,1}\kappa^{k-1}, \beta_{1,1} \leq \frac{\sqrt{C}}{\sqrt{C} + 1}$$

where  $\kappa \in (0, 1)$  and  $C = \frac{(1+\kappa)\kappa^2 G_0^3}{168(1-\kappa)G_\infty^3}$ , then

$$\frac{1}{N} \sum_{k=1}^N \mathbb{E} \|\nabla F(\mathbf{z}_k)\|^2 \leq \frac{C_1}{N} + \frac{C_2 \sigma^2}{m}, \quad (2)$$

for some positive constants  $C_1$  and  $C_2$ .

**Corollary 1.** *Under assumptions in Theorem 1, if  $N \geq 3C_1\epsilon^{-2}$  and  $m \geq 3C_2\sigma^2\epsilon^{-2}$ , then there exists an iterate  $\mathbf{z}_k$ ,  $k \in \{1, \dots, N\}$  that is an  $\epsilon$ -SFNE point of Game (1).*

**Corollary 2.** *Algorithm 1 requires  $\mathcal{O}(\epsilon^{-4})$  gradient evaluations of the objective function to find an  $\epsilon$ -SFNE point of Game 1. This is consistent with other adaptive methods such as [14].*

## 5 Numerical Studies

### (I) A Synthetic Data Experiment:

Simultaneous ADAM (S-ADAM) is one of most commonly used approaches for solving min-max problems that are formulated using DNNs such as training GANs [15]. In this method, the minimization and the maximization parameters are updated simultaneously using the ADAM algorithm [6]. However, this method fails drastically in solving simple min-max problems. To better understand this issue, consider solving the following simple stochastic min-max problem

$$f(\theta, \alpha) = \begin{cases} c(\theta - \alpha) + (\theta^2 - \alpha^2) + k\theta\alpha, & \text{w.p } \frac{1}{3}, \\ (\theta - \alpha) + (\theta^2 - \alpha^2) + k\theta\alpha, & \text{w.p } \frac{2}{3}, \end{cases} \quad (3)$$

where  $c > 1$  and  $k \geq 0$ . Some calculations lead to

$$F(\theta, \alpha) = \frac{c+2}{3}(\theta - \alpha) + (\theta^2 - \alpha^2) + k\theta\alpha.$$

This problem has the following unique FNE

$$(\theta^*, \alpha^*) = -\frac{c+2}{3k^2+12}(2-k, 2+k).$$

Since  $\nabla_{\theta}^2 F(\theta, \alpha) = 2\mathbf{I} \succ 0$  and  $\nabla_{\alpha}^2 F(\theta, \alpha) = -2\mathbf{I} \prec 0$ , this function is strongly-convex-strongly-concave and many available algorithms [40, 41] can compute its FNE due to its special structure.

This case study shows that despite the simplicity of the problem, S-ADAM is unable to recover the single FNE point of this function. We also compare the performance of S-ADAM with our proposed algorithm. To do the comparison, we define  $e_k = \frac{\|\mathbf{z}_k - \mathbf{z}_*\|}{\|\mathbf{z}_*\|}$  such that  $\mathbf{z}_k = (\theta_k, \alpha_k)$  and  $\mathbf{z}_* = (\theta_*, \alpha_*)$  and  $\mathcal{R}_k = \frac{1}{k} \sum_{i=1}^k \|\nabla F(\mathbf{z}_i)\|^2$  to measure the performance of different methods. We set the parameters at  $c = 1010, k = 0.01, N = 10^7, \eta = 10^{-2}, \beta_1 = 0, \beta_2 = 1/(1 + c^2)$  and  $\beta_3 = 0.1$ . All other parameters are initialized at zero. Figure 1 shows the result of the experiment. We have assigned 2 different scales on the vertical dimension due to space limitations. The left axis depicts the error rate,  $e_k$ , and the right one the average norm of the gradient,  $\mathcal{R}_k$ . ADAM<sup>3</sup> converges to the only FNE point, while S-ADAM is unable to locate it. This shows that S-ADAM is unreliable even for a simple strongly-convex-strongly-concave problem.

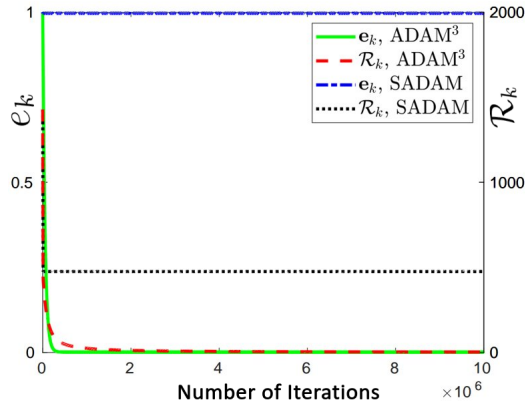


Figure 1: Left/Right y-axis: Error rate,  $e_t$  / Average norm of gradient,  $\mathcal{R}_k$ . S-ADAM misses the unique FNE point.

### (II) Training GANs with ADAM<sup>3</sup>:

Algorithm 1 is used to train GANs on the publicly available CIFAR-10 data set, containing 60000 color images of size  $32 \times 32$  in 10 different classes (see <https://www.cs.toronto.edu/~kriz/cifar.html>).

*Models and tasks:* The generator’s network consists of the input layer, 2 hidden layers and the output layer. Each of the input and hidden layers consist of a transposed convolution layer followed by batch normalization and a ReLU activation function. The output layer is a transposed convolution layer with a hyperbolic tangent activation function. The network for the discriminator also has the input layer, 2 hidden layers and the output layer. Both the input and hidden layers are convolutional layers followed by instance normalization and a Leaky ReLU activation function with slope 0.2. The output layer consists only of a convolutional layer. The scripts containing the detail design of the networks, together with the implementation of ADAM<sup>3</sup> and its competitor Optimistic AdaGrad (OAdagrad) [14] in PyTorch, will be available at <https://github.com/babakbarazandeh>.

The parameters are set to  $\eta = 0.5 \times 10^{-3}$ ,  $\beta_1 = 0.5$ ,  $\beta_2 = 0.9$  and  $\beta_3 = 0.5$ , respectively and the batch size to 64. Finally, the experiment runs for a total of 40,000 iterations. Figure 2 depicts the inception score of the generated images, a metric that evaluates their quality [42]. It can be seen that ADAM<sup>3</sup> exhibits better performance than OAdagrad at all iteration stages. Some generated samples are available at Figure 3.

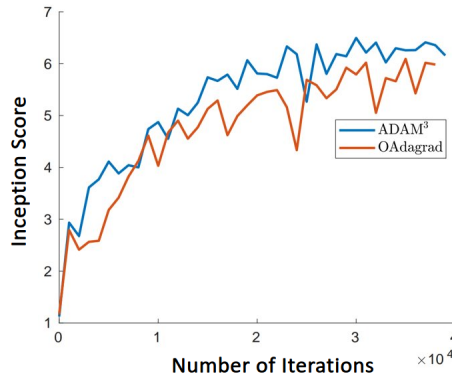


Figure 2: Inception score for generated CIFAR-10 images using ADAM<sup>3</sup> and OAdagrad.

## Acknowledgement

The work of Babak Barazandeh was supported by the UF Informatics Institute and of George Michailidis by NSF grants DMS 1854476 and DMS 1830175. The authors would also like to thank Dr. Meisam Razaviyayn for his insightful comments that helped to improve the quality of the work.

## References

- [1] Léon Bottou, Frank E Curtis, and Jorge Nocedal, “Optimization methods for large-scale machine learning,” *Siam Review*, vol. 60, no. 2, pp. 223–311, 2018.
- [2] Robert A Jacobs, “Increased rates of convergence through learning rate adaptation,” *Neural networks*, vol. 1, no. 4, pp. 295–307, 1988.
- [3] Sue Becker, Yann Le Cun, et al., “Improving the convergence of back-propagation learning with second order methods,” 1988.
- [4] John Duchi, Elad Hazan, and Yoram Singer, “Adaptive subgradient methods for online learning and stochastic optimization,” *Journal of machine learning research*, vol. 12, no. 7, 2011.
- [5] H Brendan McMahan and Matthew Streeter, “Adaptive bound optimization for online convex optimization,” *arXiv preprint arXiv:1002.4908*, 2010.

- [6] Diederik P Kingma and Jimmy Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- [7] Sashank J Reddi, Satyen Kale, and Sanjiv Kumar, “On the convergence of adam and beyond,” in *International Conference on Learning Representations*, 2018.
- [8] Yurii E Nesterov, “A method for solving the convex programming problem with convergence rate  $O(1/k^2)$ ,” in *Dokl. akad. nauk Sssr*, 1983, vol. 269, pp. 543–547.
- [9] Boris T Polyak, “Some methods of speeding up the convergence of iteration methods,” *USSR Computational Mathematics and Mathematical Physics*, vol. 4, no. 5, pp. 1–17, 1964.
- [10] Xiangyi Chen, Sijia Liu, Ruoyu Sun, and Mingyi Hong, “On the convergence of a class of adam-type algorithms for non-convex optimization,” *arXiv preprint arXiv:1808.02941*, 2018.
- [11] Manzil Zaheer, Sashank Reddi, Devendra Sachan, Satyen Kale, and Sanjiv Kumar, “Adaptive methods for nonconvex optimization,” in *Advances in Neural Information Processing Systems*, 2018, pp. 9815–9825.
- [12] Parvin Nazari, Davoud Ataee Tarzanagh, and George Michailidis, “Dadam: A consensus-based distributed adaptive gradient method for online optimization,” *arXiv preprint arXiv:1901.09109*, 2019.
- [13] Parvin Nazari, Davoud Ataee Tarzanagh, and George Michailidis, “Adaptive first-and zeroth-order methods for weakly convex stochastic optimization problems,” *arXiv preprint arXiv:2005.09261*, 2020.
- [14] Mingrui Liu, Youssef Mroueh, Jerret Ross, Wei Zhang, Xiaodong Cui, Payel Das, and Tianbao Yang, “Towards better understanding of adaptive gradient algorithms in generative adversarial nets,” *arXiv preprint arXiv:1912.11940*, 2019.
- [15] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio, “Generative adversarial nets,” in *Advances in neural information processing systems*, 2014, pp. 2672–2680.
- [16] Martin Arjovsky, Soumith Chintala, and Léon Bottou, “Wasserstein gan,” *arXiv preprint arXiv:1701.07875*, 2017.
- [17] Babak Barazandeh, Meisam Razaviyayn, and Maziar Sanjabi, “Training generative networks using random discriminators,” in *2019 IEEE Data Science Workshop (DSW)*. IEEE, 2019, pp. 327–332.
- [18] Babak Barazandeh and Meisam Razaviyayn, “On the behavior of the expectation-maximization algorithm for mixture models,” in *2018 IEEE Global Conference on Signal and Information Processing (GlobalSIP)*. IEEE, 2018, pp. 61–65.
- [19] Chi Jin, Praneeth Netrapalli, and Michael I. Jordan, “Minmax optimization: Stable limit points of gradient descent ascent are locally optimal,” *arXiv preprint arXiv:1902.00618*, 2019.
- [20] Gauthier Gidel, Tony Jebara, and Simon Lacoste-Julien, “Frank-wolfe algorithms for saddle point problems,” in *Artificial Intelligence and Statistics*, 2017, pp. 362–371.
- [21] Erfan Y. Hamedani, Afrooz Jalilzadeh, Necdet S. Aybat, and Uday V. Shanbhag, “Iteration complexity of randomized primal-dual methods for convex-concave saddle point problems,” in *arXiv preprint*, 2018, arXiv:1806.04118.
- [22] Constantinos Daskalakis, Andrew Ilyas, Vasilis Syrgkanis, and Haoyang Zeng, “Training gans with optimism,” *arXiv preprint arXiv:1711.00141*, 2017.
- [23] Farzan Farnia and Asuman Ozdaglar, “Gans may have no nash equilibria,” *arXiv preprint arXiv:2002.09124*, 2020.

- [24] Maher Nouiehed, Maziar Sanjabi, Tianjian Huang, Jason D. Lee, and Meisam Razaviyayn, “Solving a class of non-convex min-max games using iterative first order methods,” in *Advances in Neural Information Processing Systems*, 2019, pp. 14905–14916.
- [25] Babak Barazandeh and Meisam Razaviyayn, “Solving non-convex non-differentiable min-max games using proximal gradient method,” in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 3162–3166.
- [26] Jong-Shi Pang and Meisam Razaviyayn, “A unified distributed algorithm for non-cooperative games,” in *Big Data over Networks*, 2016, Cambridge University Press.
- [27] Jong-Shi Pang and Gesualdo Scutari, “Nonconvex games with side constraints,” *SIAM Journal on Optimization*, vol. 21, no. 4, pp. 1491–1522, 2011.
- [28] Alfredo N. Iusem, Alejandro Jofré, Roberto I. Oliveira, , and Philip Thompson, “Extragradient method with variance reduction for stochastic variational inequalities,” *SIAM Journal on Optimization*, vol. 27, no. 2, pp. 686–724, 2017.
- [29] Tianyi Lin, Chi Jin, and Michael I Jordan, “On gradient descent ascent for nonconvex-concave minimax problems,” *arXiv preprint arXiv:1906.00331*, 2019.
- [30] Robert Mansel Gower, Nicolas Loizou, Xun Qian, Alibek Sailanbayev, Egor Shulgin, and Peter Richtárik, “Sgd: General analysis and improved rates,” in *International Conference on Machine Learning*. PMLR, 2019, pp. 5200–5209.
- [31] Y. Nesterov, “Introductory lectures on convex programming volume i: Basic course,” in *Lecture notes*, 1998, vol. 3, p. 5.
- [32] Y. Nesterov, “Gradient methods for minimizing composite functions,” *Mathematical Programming*, vol. 140, pp. 125–161, 2013.
- [33] Meisam Razaviyayn, Tianjian Huang, Songtao Lu, Maher Nouiehed, Maziar Sanjabi, and Mingyi Hong, “Nonconvex min-max optimization: Applications, challenges, and recent theoretical advances,” *IEEE Signal Processing Magazine*, vol. 37, no. 5, pp. 55–66, 2020.
- [34] P. Mertikopoulos, H. Zenati, B. Lecouat, C.S. Foo, V. Chandrasekhar, and G. Piliouras, “Optimistic mirror descent in saddle-point problems: Going the extra (gradient) mile,” in *ICLR’19-International Conference on Learning Representations*, 2019.
- [35] Mingrui Liu, Youssef Mroueh, Wei Zhang, Xiaodong Cui, Tianbao Yang, and Payel Das, “Decentralized parallel algorithm for training generative adversarial nets,” *arXiv preprint arXiv:1910.12999*, 2019.
- [36] Hongwei Tan, Linyong Zhou, Guodong Wang, and Zili Zhang, “Improved performance of gans via integrating gradient penalty with spectral normalization,” in *International Conference on Knowledge Science, Engineering and Management*. Springer, 2020, pp. 414–426.
- [37] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen, “Progressive growing of gans for improved quality, stability, and variation,” *arXiv preprint arXiv:1710.10196*, 2017.
- [38] Karol Kurach, Mario Lučić, Xiaohua Zhai, Marcin Michalski, and Sylvain Gelly, “A large-scale study on regularization and normalization in gans,” in *International Conference on Machine Learning*. PMLR, 2019, pp. 3581–3590.
- [39] Tran Thi Phuong and Trieu Le Phong, “On the convergence proof of amsgrad and a new version,” *arXiv e-prints*, pp. arXiv–1904, 2019.
- [40] Kiran K Thekumparampil, Prateek Jain, Praneeth Netrapalli, and Sewoong Oh, “Efficient algorithms for smooth minimax optimization,” in *Advances in Neural Information Processing Systems*, 2019, pp. 12680–12691.



- [41] Dmitrii M. Ostrovskii, Andrew Lowy, and Meisam Razaviyayn, “Efficient search of first-order nash equilibria in nonconvex-concave smooth min-max problems,” *arXiv preprint arXiv:2002.07919*, 2020.
- [42] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen, “Improved techniques for training gans,” in *Advances in neural information processing systems*, 2016, pp. 2234–2242.

## Appendix

In this section, we provide proofs for Theorem 1 and required auxiliary lemmas.

**Remark 1.** For any set of vectors  $\{\mathbf{a}_i\}_{i=1}^M$ ,  $\mathbf{b}$  and  $\mathbf{c}$  in  $\mathbb{R}^d$  we have

$$1. \left\| \sum_{i=1}^M \mathbf{a}_i \right\|^2 \leq M \sum_{i=1}^M \|\mathbf{a}_i\|^2, \quad 2. \|\mathbf{a} \circ \mathbf{b}\| \leq \|\mathbf{a}\|_\infty \|\mathbf{b}\|_1, \quad 3. \|\mathbf{b} \circ \mathbf{c}\| \leq \|\mathbf{b}\|_\infty \|\mathbf{c}\|.$$

**Lemma 2.** Assume that  $\gamma := \beta_{1,1}/\beta_2 \leq 1$  in Algorithm 1. Then, for each  $k \in \{1 \cdots N\}$  we have

$$\|\tilde{v}_k^{-\frac{1}{2}} \circ m_{k-1}\| \leq \sqrt{\frac{d}{u_c}},$$

where  $u_c := (1 - \beta_3)(1 - \beta_{1,1})(1 - \beta_2)(1 - \gamma)$ .

*Proof.* For each  $k \in \{1 \cdots N\}$  and  $r \in \{1, \dots, d\}$ , let  $\tilde{v}_{r,k}^{-\frac{1}{2}}$  and  $m_{r,k}$  represent the values of the  $r^{th}$  coordinate of vectors  $\tilde{\mathbf{v}}_k^{-\frac{1}{2}}$  and  $\mathbf{m}_k$ , respectively. Then, from the update rule of Algorithm 1 we have

$$\tilde{v}_{r,k} = \beta_3 \tilde{v}_{r,k-1} + (1 - \beta_3) \max(\tilde{v}_{r,k-1}, v_{r,k}),$$

which implies that  $\tilde{v}_{r,k} \geq (1 - \beta_3)v_{r,k}$ . Besides, it can be easily seen from the update rule of  $\mathbf{m}_k$  and  $\mathbf{v}_k$  in Algorithm 1 that

$$m_{r,k} = \sum_{s=1}^k \left( \prod_{l=s+1}^k \beta_{1,l} \right) (1 - \beta_{1,s}) \hat{g}_{r,s}, \quad \text{and} \quad v_{r,k} = (1 - \beta_2) \sum_{s=1}^k \beta_2^{k-s} \hat{g}_{r,s}^2.$$

Thus,

$$\begin{aligned} |v_{r,k}^{-\frac{1}{2}} m_{r,k-1}|^2 &\leq |v_{r,k-1}^{-\frac{1}{2}} m_{r,k-1}|^2 \leq \frac{\left( \sum_{s=1}^{k-1} \left( \prod_{l=s+1}^{k-1} \beta_{1,l} \right) (1 - \beta_{1,s}) \hat{g}_{r,s} \right)^2}{(1 - \beta_2) \sum_{s=1}^{k-1} \beta_2^{k-s-1} \hat{g}_{r,s}^2} \\ &\leq \frac{\left( \sum_{s=1}^{k-1} \left( \prod_{l=s+1}^{k-1} \beta_{1,l} \right) \hat{g}_{r,s} \right)^2}{(1 - \beta_2) \sum_{s=1}^{k-1} \beta_2^{k-s-1} \hat{g}_{r,s}^2}, \end{aligned} \tag{4}$$

where the first inequality follows since  $v_{r,k}^{-\frac{1}{2}} \leq v_{r,k-1}^{-\frac{1}{2}}$  for all  $r \in [d]$  and the last inequality uses our assumption that  $\beta_{1,s} \leq 1$  for all  $s \geq 1$ .

Now, let  $\pi_s = \prod_{l=s+1}^{k-1} \beta_{1,l}$ . Since  $\beta_{1,l}$  is decreasing, we get  $\pi_s \leq \beta_{1,1}^{k-s-1}$ . This, together with

$(\sum_i a_i b_i)^2 \leq (\sum_i a_i^2)(\sum_i b_i^2)$  implies that

$$\begin{aligned}
\frac{\left(\sum_{s=1}^{k-1} \pi_s \hat{g}_{r,s}\right)^2}{(1-\beta_2) \sum_{s=1}^{k-1} \beta_2^{k-s-1} \hat{g}_{r,s}^2} &\leq \frac{\left(\sum_{s=1}^{k-1} \pi_s\right) \left(\sum_{s=1}^{k-1} \pi_s \hat{g}_{r,s}^2\right)}{(1-\beta_2) \sum_{s=1}^{k-1} \beta_2^{k-s-1} \hat{g}_{r,s}^2} \\
&\leq \frac{1}{1-\beta_2} \left(\sum_{s=1}^{k-1} \pi_s\right) \left(\sum_{s=1}^{k-1} \frac{\pi_s \hat{g}_{r,s}^2}{\beta_2^{k-s-1} \hat{g}_{r,s}^2}\right) \\
&\leq \frac{1}{1-\beta_2} \left(\sum_{s=1}^{k-1} \pi_s\right) \sum_{s=1}^{k-1} \frac{\pi_s}{\beta_2^{k-s-1}} \\
&\leq \frac{1}{1-\beta_2} \frac{1}{1-\beta_{1,1}} \frac{1}{1-\gamma},
\end{aligned}$$

where the last inequality follows from our assumption  $\gamma = \frac{\beta_{1,1}}{\beta_2} \leq 1$ . Finally, substituting the above inequality into (4) yields the desired result.  $\square$

**Lemma 3.** For each  $k \in \{1 \cdots N\}$  and  $r \in \{1, \dots, d\}$ , let  $\tilde{v}_{r,k}$  represent the value of the  $r^{th}$  coordinate of vector  $\tilde{\mathbf{v}}_k$ . Then, for the sequence of  $\tilde{\mathbf{v}}_k$ 's generated by Algorithm 1 we have

1.  $\sum_{k=1}^N \|\tilde{\mathbf{v}}_k^p - \tilde{\mathbf{v}}_{k-1}^p\|_1 \leq \sum_{r=1}^d \max(\tilde{v}_{r,0}^p, \tilde{v}_{r,N}^p)$  and
2.  $\sum_{k=1}^N \|\tilde{\mathbf{v}}_k^p - \tilde{\mathbf{v}}_{k-1}^p\|_1^2 \leq \sum_{r=1}^d \tilde{v}_{r,0}^p \max(\tilde{v}_{r,0}^p, \tilde{v}_{r,N}^p)$

where  $p \in \mathbb{R}$  and the vector powers are considered to be element-wise.

*Proof.* 1. If  $p > 0$ , from the update rule of  $\tilde{\mathbf{v}}_k$  in Algorithm 1 we have

$$\begin{aligned}
\sum_{k=1}^N \left\| \tilde{\mathbf{v}}_k^p - \tilde{\mathbf{v}}_{k-1}^p \right\|_1 &= \sum_{k=1}^N \sum_{r=1}^d (\tilde{v}_{r,k}^p - \tilde{v}_{r,k-1}^p) = \sum_{r=1}^d \sum_{k=1}^N (\tilde{v}_{r,k}^p - \tilde{v}_{r,k-1}^p) \\
&\leq \sum_{r=1}^d \tilde{v}_{r,N}^p,
\end{aligned}$$

where the first equality is due to the fact that for  $p > 0$ , each element of  $\tilde{\mathbf{v}}_k^p$  is increasing in  $k$  and the last inequality uses the telescoping sum. Now, we consider the case when  $p < 0$ . It can be easily seen that

$$\sum_{k=1}^N \left\| \tilde{\mathbf{v}}_k^p - \tilde{\mathbf{v}}_{k-1}^p \right\|_1 = \sum_{k=1}^N \sum_{r=1}^d (-\tilde{v}_{r,k}^q + \tilde{v}_{r,k-1}^q) \leq \sum_{r=1}^d \tilde{v}_{r,0}^q.$$

2. For  $p > 0$ , it follows that

$$\begin{aligned}
\sum_{k=1}^N \left\| \tilde{\mathbf{v}}_{i,k}^p - \tilde{\mathbf{v}}_{i,k-1}^p \right\|_1^2 &\leq \sum_{k=1}^N \sum_{r=1}^d \left( \tilde{v}_{r,k}^p - \tilde{v}_{r,k-1}^p \right) \tilde{v}_{r,k}^p \leq \sum_{k=1}^N \sum_{r=1}^d \left( \tilde{v}_{r,k}^p - \tilde{v}_{r,k-1}^p \right) \tilde{v}_{r,N}^p \\
&\leq \sum_{r=1}^d \left( \tilde{v}_{r,0}^p - \tilde{v}_{r,N}^p \right) \tilde{v}_{r,N}^p \\
&\leq \sum_{r=1}^d \tilde{v}_{r,0}^p \tilde{v}_{r,N}^p.
\end{aligned}$$

Now, we consider the case when  $p < 0$ . It can be easily seen that

$$\begin{aligned}
\sum_{k=1}^N \left\| \tilde{\mathbf{v}}_k^p - \tilde{\mathbf{v}}_{k-1}^p \right\|_1^2 &\leq \sum_{k=1}^N \sum_{r=1}^d (-\tilde{v}_{r,k}^p + \tilde{v}_{r,k-1}^p) (\tilde{v}_{r,k-1}^p) \leq \sum_{k=1}^N \sum_{r=1}^d (-\tilde{v}_{r,k}^p + \tilde{v}_{r,k-1}^p) \tilde{v}_{r,0}^p \\
&\leq \sum_{r=1}^d \tilde{v}_{r,0}^p \tilde{v}_{r,0}^p.
\end{aligned}$$

□

**Lemma 4.** *Under Assumptions A and D we have*

$$\sum_{k=1}^N \left( \left\| \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{x}_{k-1} - \mathbf{x}_*) \right\|^2 - \left\| \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{x}_k - \mathbf{x}_*) \right\|^2 \right) \leq 3D^2 dG_\infty.$$

*Proof.* Observe that

$$\begin{aligned}
&\sum_{k=1}^N \left( \left\| \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{x}_{k-1} - \mathbf{x}_*) \right\|^2 - \left\| \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{x}_k - \mathbf{x}_*) \right\|^2 \right) \\
&= \left\| \tilde{\mathbf{v}}_0^{\frac{1}{4}} \circ (\mathbf{x}_0 - \mathbf{x}_*) \right\|^2 + \left( - \left\| \tilde{\mathbf{v}}_0^{\frac{1}{4}} \circ (\mathbf{x}_1 - \mathbf{x}_*) \right\|^2 + \left\| \tilde{\mathbf{v}}_1^{\frac{1}{4}} \circ (\mathbf{x}_1 - \mathbf{x}_*) \right\|^2 \right) \\
&+ \left( - \left\| \tilde{\mathbf{v}}_1^{\frac{1}{4}} \circ (\mathbf{x}_2 - \mathbf{x}_*) \right\|^2 + \left\| \tilde{\mathbf{v}}_2^{\frac{1}{4}} \circ (\mathbf{x}_2 - \mathbf{x}_*) \right\|^2 \right) \\
&\quad \vdots \\
&+ \left( - \left\| \tilde{\mathbf{v}}_{N-2}^{\frac{1}{4}} \circ (\mathbf{x}_{N-1} - \mathbf{x}_*) \right\|^2 + \left\| \tilde{\mathbf{v}}_{N-1}^{\frac{1}{4}} \circ (\mathbf{x}_{N-1} - \mathbf{x}_*) \right\|^2 \right) - \left\| \tilde{\mathbf{v}}_{N-1}^{\frac{1}{4}} \circ (\mathbf{x}_N - \mathbf{x}_*) \right\|^2. \tag{5}
\end{aligned}$$

For arbitrary  $s^{th}$  pairs in (5), we have

$$\begin{aligned}
& - \left\| \tilde{\mathbf{v}}_{s-1}^{\frac{1}{4}} \circ (\mathbf{x}_s - \mathbf{x}_*) \right\|^2 + \left\| \tilde{\mathbf{v}}_s^{\frac{1}{4}} \circ (\mathbf{x}_s - \mathbf{x}_*) \right\|^2 \\
& = \left( - \left\| (\tilde{\mathbf{v}}_{s-1}^{\frac{1}{4}} - \tilde{\mathbf{v}}_s^{\frac{1}{4}} + \tilde{\mathbf{v}}_s^{\frac{1}{4}}) \circ (\mathbf{x}_s - \mathbf{x}_*) \right\|^2 + \left\| \tilde{\mathbf{v}}_s^{\frac{1}{4}} \circ (\mathbf{x}_s - \mathbf{x}_*) \right\|^2 \right) \\
& \quad \cdot \left( \left\| \tilde{\mathbf{v}}_{s-1}^{\frac{1}{4}} \circ (\mathbf{x}_s - \mathbf{x}_*) \right\|^2 + \left\| \tilde{\mathbf{v}}_s^{\frac{1}{4}} \circ (\mathbf{x}_s - \mathbf{x}_*) \right\|^2 \right) \\
& \leq \left( \left\| (\tilde{\mathbf{v}}_{s-1}^{\frac{1}{4}} - \tilde{\mathbf{v}}_s^{\frac{1}{4}}) \circ (\mathbf{x}_s - \mathbf{x}_*) \right\|^2 \right) 2\sqrt{G_\infty}D \\
& \leq 2D^2\sqrt{G_\infty}\|\tilde{\mathbf{v}}_{s-1}^{\frac{1}{4}} - \tilde{\mathbf{v}}_s^{\frac{1}{4}}\|_1,
\end{aligned} \tag{6}$$

where the first inequality follows from  $\|\mathbf{a}\| - \|\mathbf{b}\| \leq \|\mathbf{a} - \mathbf{b}\|$ , Assumption D and Lemma 1.

As a result,

$$\begin{aligned}
\sum_{k=1}^N \left( \left\| \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{x}_{k-1} - \mathbf{x}_*) \right\|^2 - \left\| \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{x}_k - \mathbf{x}_*) \right\|^2 \right) & \leq 2D^2dG_\infty + \left\| \tilde{\mathbf{v}}_0^{\frac{1}{4}} \circ (\mathbf{x}_0 - \mathbf{x}_*) \right\|^2 - \left\| \tilde{\mathbf{v}}_{N-1}^{\frac{1}{4}} \circ (\mathbf{x}_N - \mathbf{x}_*) \right\|^2 \\
& \leq 3D^2dG_\infty,
\end{aligned}$$

where the first inequality follows from Lemma 3 and last inequality uses the same lemma, Assumption D and the fact that  $d \geq 1$ .  $\square$

**Theorem 1.** *Let Assumptions A–E hold, and  $L$ ,  $G_\infty$ ,  $G_0$ ,  $\sigma$  be defined therein. In Algorithm 1, if we choose*

$$\eta \leq \sqrt{G_0^3/(56L^2G_\infty)}, \quad \text{and} \quad \beta_{1,1} \leq \frac{\sqrt{C}}{\sqrt{C} + 1},$$

where  $C = \frac{(1+\kappa)\kappa^2G_0^3}{168(1-\kappa)G_\infty^3}$ , then

$$\frac{1}{N} \sum_{k=1}^N \mathbb{E} \|\nabla F(\mathbf{z}_k)\|^2 \leq \frac{C_1}{N} + \frac{C_2\sigma^2}{m}, \tag{7}$$

for some positive constants  $C_1$  and  $C_2$ .

*Proof.* We divide the proof into four steps. In Step 1, we show that the gradient norm is bounded by the norm of search direction and auxiliary variables  $\mathbf{x}_k$  and  $\mathbf{z}_k$ . Then in Steps 2 and 3, we give upper bounds for these terms. Finally, in Step 4, we provide the convergence analysis.

**Step 1** shows that under Assumption A (2), we have

$$\frac{1}{N} \sum_{k=1}^N \|\mathbf{g}_k\|^2 \leq \frac{3}{N\eta^2(1 - \beta_{1,1})^2G_\infty^{-2}} \sum_{k=1}^N \eta^2 R_{1,k} + R_{2,k}, \tag{8}$$

where

$$\begin{aligned}
R_{1,k} & := \left\| -\mathbf{d}_k + (1 - \beta_{1,k})\tilde{\mathbf{v}}_k^{-\frac{1}{2}} \circ \mathbf{g}_k \right\|^2 \\
R_{2,k} & := \|\mathbf{z}_k - \mathbf{x}_k\|^2 + \|\mathbf{z}_k - \mathbf{x}_{k-1}\|^2.
\end{aligned} \tag{9}$$

It follows from the update rule of  $\mathbf{x}_k$  in Algorithm 1 that

$$\begin{aligned} \eta(1 - \beta_{1,k}) \left( \tilde{\mathbf{v}}_k^{-\frac{1}{2}} \circ \mathbf{g}_k \right) &= \mathbf{z}_k - \mathbf{x}_k + (\mathbf{x}_{k-1} - \eta \mathbf{d}_k) \\ &\quad - \mathbf{z}_k + \eta(1 - \beta_{1,k}) \left( \tilde{\mathbf{v}}_k^{-\frac{1}{2}} \circ \mathbf{g}_k \right). \end{aligned}$$

Now, using Remark 1, we get

$$\begin{aligned} \eta^2(1 - \beta_{1,k})^2 \left\| \tilde{\mathbf{v}}_k^{-\frac{1}{2}} \circ \mathbf{g}_k \right\|^2 &\leq 3\eta^2 \left\| -\mathbf{d}_k + (1 - \beta_{1,k}) \tilde{\mathbf{v}}_k^{-\frac{1}{2}} \circ \mathbf{g}_k \right\|^2 \\ &\quad + 3 \left( \|\mathbf{z}_k - \mathbf{x}_k\|^2 + \|\mathbf{z}_k - \mathbf{x}_{k-1}\|^2 \right). \end{aligned} \quad (10)$$

From Lemma 1, we have  $\|\tilde{\mathbf{v}}_k^{-\frac{1}{2}}\|_\infty \geq G_\infty^{-1}$  which implies that

$$\left\| \tilde{\mathbf{v}}_k^{-\frac{1}{2}} \circ \mathbf{g}_k \right\|^2 \geq G_\infty^{-2} \|\mathbf{g}_k\|^2.$$

Now, it follows from the above inequality and (9) that

$$\eta^2(1 - \beta_{1,k})^2 G_\infty^{-2} \|\mathbf{g}_k\|^2 \leq \eta^2(1 - \beta_{1,k})^2 \left\| \tilde{\mathbf{v}}_k^{-\frac{1}{2}} \circ \mathbf{g}_k \right\|^2 \leq 3\eta^2 R_{1,k} + 3R_{2,k},$$

which gives (8).

**Step 2** establishes an upper bound for  $R_{1,k}$  defined in (9). More specifically, we show that

$$\frac{1}{N} \sum_{k=1}^N R_{1,k} \leq \frac{2d\beta_{1,1}^2}{Nu_c(1 - \kappa^2)} + \frac{2}{NG_0^2} \sum_{k=1}^N \|\epsilon_k\|^2, \quad (11)$$

where  $u_c$  is defined in Lemma 2 and  $\epsilon_k = \hat{\mathbf{g}}_k - \mathbf{g}_k$ .

From the definition of  $\mathbf{d}_k$  in Algorithm 1, we have

$$\mathbf{d}_k = \beta_{1,k} \tilde{\mathbf{v}}_k^{-\frac{1}{2}} \circ \mathbf{m}_{k-1} + (1 - \beta_{1,k}) \tilde{\mathbf{v}}_k^{-\frac{1}{2}} \circ \hat{\mathbf{g}}_k. \quad (12)$$

Hence,

$$-\mathbf{d}_k + (1 - \beta_{1,k}) \tilde{\mathbf{v}}_k^{-\frac{1}{2}} \circ \mathbf{g}_k = -\beta_{1,k} \tilde{\mathbf{v}}_k^{-\frac{1}{2}} \circ \mathbf{m}_{k-1} + (1 - \beta_{1,k}) \tilde{\mathbf{v}}_k^{-\frac{1}{2}} \circ (\mathbf{g}_k - \hat{\mathbf{g}}_k),$$

which implies that

$$\begin{aligned} R_{1,k} &= \left\| -\beta_{1,k} \tilde{\mathbf{v}}_k^{-\frac{1}{2}} \circ \mathbf{m}_{k-1} + (1 - \beta_{1,k}) \tilde{\mathbf{v}}_k^{-\frac{1}{2}} \circ (\mathbf{g}_k - (\mathbf{g}_k + \epsilon_k)) \right\|^2 \\ &\leq 2\beta_{1,k}^2 \left\| \tilde{\mathbf{v}}_k^{-\frac{1}{2}} \circ \mathbf{m}_{k-1} \right\|^2 + 2(1 - \beta_{1,k})^2 \left\| \tilde{\mathbf{v}}_k^{-\frac{1}{2}} \circ \epsilon_k \right\|^2. \end{aligned} \quad (13)$$

Here, the equality is obtained since  $\epsilon_k = \hat{\mathbf{g}}_k - \mathbf{g}_k$ , and the inequality follows from Remark 1.

For the first term on the R.H.S. of (13), it follows from Lemma 2 that

$$\left\| \tilde{\mathbf{v}}_k^{-\frac{1}{2}} \circ \mathbf{m}_{k-1} \right\|^2 \leq \frac{d}{u_c}. \quad (14a)$$

Further, for the second term on the R.H.S. of (13), we have

$$\left\| \tilde{\mathbf{v}}_k^{-\frac{1}{2}} \circ \boldsymbol{\epsilon}_k \right\|^2 \leq \|\tilde{\mathbf{v}}_k^{-\frac{1}{2}}\|_\infty^2 \|\boldsymbol{\epsilon}_k\|^2 \leq \|\tilde{\mathbf{v}}_0^{-\frac{1}{2}}\|_\infty^2 \|\boldsymbol{\epsilon}_k\|^2 \leq \frac{1}{G_0^2} \|\boldsymbol{\epsilon}_k\|^2, \quad (14b)$$

where the inequality uses Remark 1, the fact that each element of  $\tilde{\mathbf{v}}_k^{-\frac{1}{2}}$  is decreasing in  $k$  and our assumption that  $\|\tilde{\mathbf{v}}_0^{-\frac{1}{2}}\|_\infty \leq 1/G_0$ .

Substituting (14a)–(14b) into (13), we obtain

$$R_{1,k} \leq \frac{2d\beta_{1,k}^2}{u_c} + \frac{2}{G_0^2} \|\boldsymbol{\epsilon}_k\|^2.$$

Summing the above inequality over  $k = 1, \dots, N$  and using the fact that  $\sum_{k=1}^N \beta_{1,k}^2 \leq \beta_{1,1}^2/(1 - \kappa^2)$ , we obtain the desired result.

**Step 3** provides an upper bound for  $R_{2,k}$  defined in (9). In particular, we show that for  $\eta \leq \sqrt{G_0^3/(56L^2G_\infty)}$ , the following holds

$$\begin{aligned} \frac{1}{N} \sum_{k=1}^N R_{2,k} &\leq \frac{6D^2dG_\infty}{NG_0} + \frac{4\eta D}{NG_0} \left( \frac{\beta_{1,1}G_\infty}{1 - \kappa} \sqrt{\frac{d}{u_c}} + \frac{G_\infty^2 d}{G_0} \right) + \frac{56\eta^2 d \beta_{1,1}^2 G_\infty}{Nu_c(1 - \kappa^2)G_0} \\ &\quad + \frac{28\eta^2 d G_\infty^3}{NG_0^3} + \frac{28\eta^2 G_\infty}{NG_0^3} \sum_{k=1}^N (\beta_{1,k} - \beta_{1,k-1})^2 \|\mathbf{g}_k\|^2 + \frac{56\eta^2 G_\infty}{G_0^3 N} \sum_{k=1}^N \|\boldsymbol{\epsilon}_k\|^2. \end{aligned} \quad (15)$$

Let

$$\tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} := \left[ \hat{v}_{1,k-1}^{\frac{1}{4}}, \hat{v}_{2,k-1}^{\frac{1}{4}}, \dots, \hat{v}_{d,k-1}^{\frac{1}{4}} \right]^\top.$$

The update rule of  $\mathbf{x}_k$  in Algorithm 1 implies that

$$\begin{aligned}
& \left\| \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{x}_k - \mathbf{x}) \right\|^2 = \left\| \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{x}_{k-1} - \eta \mathbf{d}_k - \mathbf{x}) \right\|^2 \\
&= \left\| \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{x}_{k-1} - \mathbf{x}) - \eta \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ \mathbf{d}_k \right\|^2 - \left\| \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{x}_{k-1} - \mathbf{x}_k) - \eta \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ \mathbf{d}_k \right\|^2 \\
&= \left\| \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{x}_{k-1} - \mathbf{x}) \right\|^2 - \left\| \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{x}_{k-1} - \mathbf{x}_k) \right\|^2 \\
&\quad - 2 \left\langle \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{x}_{k-1} - \mathbf{x}), \eta \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ \mathbf{d}_k \right\rangle + 2 \left\langle \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{x}_{k-1} - \mathbf{x}_k), \eta \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ \mathbf{d}_k \right\rangle \\
&\quad - 2 \left\langle \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ \mathbf{z}_k, \eta \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ \mathbf{d}_k \right\rangle + 2 \left\langle \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ \mathbf{z}_k, \eta \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ \mathbf{d}_k \right\rangle \\
&= \left\| \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{x}_{k-1} - \mathbf{x}) \right\|^2 - \left\| \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{x}_{k-1} - \mathbf{z}_k + \mathbf{z}_k - \mathbf{x}_k) \right\|^2 \\
&\quad - 2 \left\langle \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{z}_k - \mathbf{x}), \eta \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ \mathbf{d}_k \right\rangle + 2 \left\langle \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{z}_k - \mathbf{x}_k), \eta \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ \mathbf{d}_k \right\rangle \\
&= \left\| \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{x}_{k-1} - \mathbf{x}) \right\|^2 - \left\| \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{x}_{k-1} - \mathbf{z}_k) \right\|^2 - \left\| \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{z}_k - \mathbf{x}_k) \right\|^2 \\
&\quad + 2 \left\langle \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{x} - \mathbf{z}_k), \eta \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ \mathbf{d}_k \right\rangle + 2 \left\langle \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{x}_k - \mathbf{z}_k), \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{x}_{k-1} - \mathbf{z}_k) \right\rangle \\
&\quad + 2 \left\langle \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{x}_k - \mathbf{z}_k), -\eta \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ \mathbf{d}_k \right\rangle,
\end{aligned}$$

where the second equality follows since  $\mathbf{x}_{k-1} - \mathbf{x}_k - \eta \mathbf{d}_k = 0$ .

Now, substituting  $\mathbf{x} = \mathbf{x}_*$  into the above equality and rearranging the terms, we get

$$\begin{aligned}
& \left\| \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{z}_k - \mathbf{x}_k) \right\|^2 + \left\| \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{x}_{k-1} - \mathbf{z}_k) \right\|^2 \\
&= \underbrace{\left\| \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{x}_{k-1} - \mathbf{x}_*) \right\|^2 - \left\| \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{x}_k - \mathbf{x}_*) \right\|^2}_{R_{2,0,k}} \\
&\quad + 2\eta \underbrace{\left\langle \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{x}_* - \mathbf{z}_k), \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ \mathbf{d}_k \right\rangle}_{R_{2,1,k}} \\
&\quad + 2\eta \underbrace{\left\langle \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{x}_k - \mathbf{z}_k), \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{d}_{k-1} - \mathbf{d}_k) \right\rangle}_{R_{2,2,k}}. \tag{16}
\end{aligned}$$

Since by our assumption  $\|\tilde{\mathbf{v}}_0^{\frac{1}{2}}\|_\infty \geq G_0$ , we have

$$\begin{aligned}
G_0 \|\mathbf{z}_k - \mathbf{x}_k\|^2 &\leq \left\| \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{z}_k - \mathbf{x}_k) \right\|^2, \quad \text{and} \\
G_0 \|\mathbf{x}_{k-1} - \mathbf{z}_k\|^2 &\leq \left\| \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{x}_{k-1} - \mathbf{z}_k) \right\|^2.
\end{aligned}$$



We substitute the above lower bounds into (16) to get

$$\|\mathbf{z}_k - \mathbf{x}_k\|^2 + \|\mathbf{x}_{k-1} - \mathbf{z}_k\|^2 \leq \frac{R_{2,0,k}}{G_0} + \frac{2\eta}{G_0} (R_{2,1,k} + R_{2,2,k}). \quad (17)$$

Next, we provide upper bounds for the terms  $R_{2,1,k}$  and  $R_{2,2,k}$ .

**Bounding  $R_{2,1,k}$ .** It follows from the update rule of  $\mathbf{d}_k$  in (12) that

$$\begin{aligned} \mathbf{d}_k &= \mathbf{d}_k - (1 - \beta_{1,k})\tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}} \circ \widehat{\mathbf{g}}_k + (1 - \beta_{1,k})\tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}} \circ \widehat{\mathbf{g}}_k \\ &= \beta_{1,k}\tilde{\mathbf{v}}_k^{-\frac{1}{2}} \circ \mathbf{m}_{k-1} + (1 - \beta_{1,k})(\tilde{\mathbf{v}}_k^{-\frac{1}{2}} - \tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}}) \circ \widehat{\mathbf{g}}_k \\ &\quad + (1 - \beta_{1,k})\tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}} \circ \mathbf{g}_k + (1 - \beta_{1,k})\tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}} \circ (\widehat{\mathbf{g}}_k - \mathbf{g}_k). \end{aligned} \quad (18)$$

To find an upper bound for  $R_{2,1,k}$ , we first multiply each term in (18) by  $\tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}}$  and then provide an upper bound for its inner product with  $\tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{x}_* - \mathbf{z}_k)$ . From Lemmas 1, 2 and Assumption D, we get

$$\left\langle \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{x}_* - \mathbf{z}_k), \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ \tilde{\mathbf{v}}_k^{-\frac{1}{2}} \circ \mathbf{m}_{k-1} \right\rangle \leq DG_\infty \left\| \tilde{\mathbf{v}}_k^{-\frac{1}{2}} \circ \mathbf{m}_{k-1} \right\| \leq DG_\infty \sqrt{du_c^{-1}}. \quad (19a)$$

Further,

$$\begin{aligned} \left\langle \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{x}_* - \mathbf{z}_k), \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\tilde{\mathbf{v}}_k^{-\frac{1}{2}} - \tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}}) \circ \widehat{\mathbf{g}}_k \right\rangle &\leq DG_\infty \|(\tilde{\mathbf{v}}_k^{-\frac{1}{2}} - \tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}}) \circ \widehat{\mathbf{g}}_k\| \\ &\leq DG_\infty \|\widehat{\mathbf{g}}_k\|_\infty \|\tilde{\mathbf{v}}_k^{-\frac{1}{2}} - \tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}}\|_1 \\ &\leq DG_\infty^2 \|\tilde{\mathbf{v}}_k^{-\frac{1}{2}} - \tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}}\|_1, \end{aligned} \quad (19b)$$

where the second inequality is obtained from Remark 1 and the last inequality is due to Assumption A (2). From Assumption C, we have

$$\left\langle \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{x}_* - \mathbf{z}_k), \tilde{\mathbf{v}}_{k-1}^{-\frac{1}{4}} \circ \mathbf{g}_k \right\rangle = \langle \mathbf{x}_* - \mathbf{z}_k, \mathbf{g}_k \rangle \leq 0. \quad (19c)$$

Further,

$$\left\langle \tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{x}_* - \mathbf{z}_k), \tilde{\mathbf{v}}_{k-1}^{-\frac{1}{4}} \circ (\widehat{\mathbf{g}}_k - \mathbf{g}_k) \right\rangle = \langle \mathbf{x}_* - \mathbf{z}_k, \widehat{\mathbf{g}}_k - \mathbf{g}_k \rangle =: \Theta_k. \quad (19d)$$

Now, using (19d)–(19b), we obtain

$$R_{2,1,k} \leq \beta_{1,k} DG_\infty \sqrt{du_c^{-1}} + DG_\infty^2 \|\tilde{\mathbf{v}}_k^{-\frac{1}{2}} - \tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}}\|_1 + \Theta_k. \quad (20)$$

**Bounding  $R_{2,2,k}$**  From the update rule of  $\mathbf{d}_k$  in (12), we get

$$\begin{aligned}
\mathbf{d}_k - \mathbf{d}_{k-1} &= \beta_{1,k} \tilde{\mathbf{v}}_k^{-\frac{1}{2}} \circ \mathbf{m}_{k-1} + (1 - \beta_{1,k}) \tilde{\mathbf{v}}_k^{-\frac{1}{2}} \circ \hat{\mathbf{g}}_k \\
&\quad - \beta_{1,k-1} \tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}} \circ \mathbf{m}_{k-2} - (1 - \beta_{1,k-1}) \tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}} \circ \hat{\mathbf{g}}_{k-1} \\
&= \beta_{1,k} \tilde{\mathbf{v}}_k^{-\frac{1}{2}} \circ \mathbf{m}_{k-1} - \beta_{1,k-1} \tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}} \circ \mathbf{m}_{k-2} \\
&\quad + (1 - \beta_{1,k}) (\tilde{\mathbf{v}}_k^{-\frac{1}{2}} - \tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}} + \tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}}) \circ \hat{\mathbf{g}}_k - (1 - \beta_{1,k-1}) \tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}} \circ \hat{\mathbf{g}}_{k-1} \\
&= \beta_{1,k} \tilde{\mathbf{v}}_k^{-\frac{1}{2}} \circ \mathbf{m}_{k-1} - \beta_{1,k-1} \tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}} \circ \mathbf{m}_{k-2} + (1 - \beta_{1,k}) (\tilde{\mathbf{v}}_k^{-\frac{1}{2}} - \tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}}) \circ \hat{\mathbf{g}}_k \\
&\quad + (1 - \beta_{1,k}) \tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}} \circ (\mathbf{g}_k + \boldsymbol{\epsilon}_k) - (1 - \beta_{1,k-1}) \tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}} \circ (\mathbf{g}_{k-1} + \boldsymbol{\epsilon}_{k-1}) \\
&= \beta_{1,k} \tilde{\mathbf{v}}_k^{-\frac{1}{2}} \circ \mathbf{m}_{k-1} - \beta_{1,k-1} \tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}} \circ \mathbf{m}_{k-2} + (1 - \beta_{1,k}) (\tilde{\mathbf{v}}_k^{-\frac{1}{2}} - \tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}}) \circ \hat{\mathbf{g}}_k \\
&\quad + (1 - \beta_{1,k}) \tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}} \circ \boldsymbol{\epsilon}_k - (1 - \beta_{1,k-1}) \tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}} \circ \boldsymbol{\epsilon}_{k-1} \\
&\quad + (1 - \beta_{1,k-1}) \tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}} \circ (\mathbf{g}_k - \mathbf{g}_{k-1}) + (\beta_{1,k-1} - \beta_{1,k}) \tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}} \circ \mathbf{g}_k. \tag{21}
\end{aligned}$$

Next, we focus on providing upper bounds for

$$R_{2,2,k} = \eta \|\tilde{\mathbf{v}}_{k-1}^{\frac{1}{4}} \circ (\mathbf{d}_k - \mathbf{d}_{k-1})\|^2 \leq \eta G_\infty \|\mathbf{d}_k - \mathbf{d}_{k-1}\|^2.$$

Observe that

$$\begin{aligned}
&\left\| \beta_{1,k} \tilde{\mathbf{v}}_k^{-\frac{1}{2}} \circ \mathbf{m}_{k-1} \right\|^2 + \left\| -\beta_{1,k-1} \tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}} \circ \mathbf{m}_{k-2} \right\|^2 \\
&\leq 2 \max \left( \left\| \beta_{1,k} \tilde{\mathbf{v}}_k^{-\frac{1}{2}} \circ \mathbf{m}_{k-1} \right\|^2, \left\| -\beta_{1,k-1} \tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}} \circ \mathbf{m}_{k-2} \right\|^2 \right) \leq \frac{2d\beta_{1,k-1}^2}{u_c}, \tag{22a}
\end{aligned}$$

where the inequality follows from Lemma 2. Using Remark 1, we get

$$\left\| (1 - \beta_{1,k}) (\tilde{\mathbf{v}}_k^{-\frac{1}{2}} - \tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}}) \circ \hat{\mathbf{g}}_k \right\|^2 \leq \|\hat{\mathbf{g}}_k\|_\infty^2 \|\tilde{\mathbf{v}}_k^{-\frac{1}{2}} - \tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}}\|_1^2 \leq G_\infty^2 \|\tilde{\mathbf{v}}_k^{-\frac{1}{2}} - \tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}}\|_1^2, \tag{22b}$$

where the last inequality uses Assumption A (2). Similarly,

$$\begin{aligned}
\left\| (1 - \beta_{1,k-1}) \tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}} \circ (\mathbf{g}_{k-1} - \mathbf{g}_k) \right\|^2 &\leq \|\tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}}\|_\infty^2 \|\mathbf{g}_{k-1} - \mathbf{g}_k\|^2 \\
&\leq \frac{L^2}{G_0^2} \|\mathbf{z}_{k-1} - \mathbf{z}_k\|^2, \tag{22c}
\end{aligned}$$

$$\begin{aligned}
\|(\beta_{1,k} - \beta_{1,k-1}) \tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}} \circ \mathbf{g}_k\|^2 &\leq (\beta_{1,k} - \beta_{1,k-1})^2 \|\tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}}\|_\infty^2 \|\mathbf{g}_k\|^2 \\
&\leq \frac{(\beta_{1,k} - \beta_{1,k-1})^2}{G_0^2} \|\mathbf{g}_k\|^2. \tag{22d}
\end{aligned}$$

By taking the norm of (21), using Remark 1 and (22a)–(22d), we get

$$\begin{aligned}
\frac{R_{2,2,k}}{G_\infty} &\leq \eta \|\mathbf{d}_k - \mathbf{d}_{k-1}\|^2 \leq 14\eta d\beta_{1,k-1}^2 u_c^{-1} + 7\eta G_\infty^2 \|\tilde{\mathbf{v}}_k^{-\frac{1}{2}} - \tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}}\|_1^2 \\
&\quad + \frac{7\eta L^2}{G_0^2} (\|\mathbf{z}_{k-1} - \mathbf{z}_k\|^2) \\
&\quad + \frac{7\eta}{G_0^2} (\beta_{1,k} - \beta_{1,k-1})^2 \|\mathbf{g}_k\|^2 \\
&\quad + \frac{7\eta}{G_0^2} (\|\boldsymbol{\epsilon}_k\|^2 + \|\boldsymbol{\epsilon}_{k-1}\|^2). \tag{23}
\end{aligned}$$

By substituting (20) and (23) into (17), we obtain

$$\begin{aligned}
&\|\mathbf{z}_k - \mathbf{x}_k\|^2 + \|\mathbf{x}_{k-1} - \mathbf{z}_k\|^2 \leq \frac{R_{2,0,k}}{G_0} \\
&\quad + \frac{2\eta}{G_0} \left( \beta_{1,k} D G_\infty \sqrt{d u_c^{-1}} + D G_\infty^2 \|\tilde{\mathbf{v}}_k^{-\frac{1}{2}} - \tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}}\|_1 + \Theta_k \right) \\
&\quad + \frac{14\eta^2 G_\infty}{G_0} \left( 2d\beta_{1,k-1}^2 u_c^{-1} + G_\infty^2 \|\tilde{\mathbf{v}}_k^{-\frac{1}{2}} - \tilde{\mathbf{v}}_{k-1}^{-\frac{1}{2}}\|_1^2 \right) \\
&\quad + \frac{14\eta^2 G_\infty}{G_0^3} \left( L^2 \|\mathbf{z}_k - \mathbf{z}_{k-1}\|^2 + (\beta_{1,k} - \beta_{1,k-1})^2 \|\mathbf{g}_k\|^2 \right) \\
&\quad + \frac{14\eta^2 G_\infty}{G_0^3} (\|\boldsymbol{\epsilon}_k\|^2 + \|\boldsymbol{\epsilon}_{k-1}\|^2).
\end{aligned}$$

Now, summing the above inequality over  $k$ , we obtain

$$\begin{aligned}
&\left( 1 - \frac{28\eta^2 L^2 G_\infty}{G_0^3} \right) \sum_{k=1}^N \|\mathbf{x}_{k-1} - \mathbf{z}_k\|^2 + \left( 1 - \frac{28\eta^2 L^2 G_\infty}{G_0^3} \right) \sum_{k=1}^N \|\mathbf{z}_k - \mathbf{x}_k\|^2 \\
&\leq \frac{3D^2 d G_\infty}{G_0} + \frac{2\eta}{G_0} \left( \frac{D\beta_{1,1} G_\infty}{1 - \kappa} \sqrt{\frac{d}{u_c}} + \frac{D G_\infty^2 d}{G_0} + \sum_{k=1}^N \Theta_k \right) \\
&\quad + \frac{28\eta^2 d \beta_{1,1}^2 G_\infty}{u_c (1 - \kappa^2) G_0} + \frac{14\eta^2 d G_\infty^3}{G_0^3} \\
&\quad + \frac{14\eta^2 G_\infty}{G_0^3} \sum_{k=1}^N (\beta_{1,k} - \beta_{1,k-1})^2 \|\mathbf{g}_k\|^2 + \frac{28\eta^2 G_\infty}{G_0^3} \sum_{k=1}^N \|\boldsymbol{\epsilon}_k\|^2 =: \text{R.H.S.} \tag{24}
\end{aligned}$$

Here, we used Lemma 3 and the fact that

$$\begin{aligned}
\sum_{k=1}^N \|\mathbf{z}_k - \mathbf{z}_{k-1}\|^2 &\leq 2 \sum_{k=1}^N \|\mathbf{z}_k - \mathbf{x}_{k-1}\|^2 + 2 \sum_{k=1}^N \|\mathbf{x}_{k-1} - \mathbf{z}_{k-1}\|^2 \\
&= 2 \sum_{k=1}^N \|\mathbf{z}_k - \mathbf{x}_{k-1}\|^2 + 2 \sum_{k=1}^N \|\mathbf{x}_k - \mathbf{z}_k\|^2, \tag{25}
\end{aligned}$$

where the inequality follows from Remark 1 and the equality uses our assumption  $\mathbf{x}_0 = \mathbf{z}_0 = 0$ .

Now, by our choice of step size  $\eta$  in the beginning of Step 3, we have  $1 - (28\eta^2 L^2 G_\infty)/G_0^3 \geq 1/2$ . Thus, (9) together with (24) implies that

$$\frac{1}{N} \sum_{k=1}^N R_{2,k} = \sum_{k=1}^N \left( \|\mathbf{x}_{k-1} - \mathbf{z}_k\|^2 + \|\mathbf{z}_k - \mathbf{x}_k\|^2 \right) \leq \frac{2}{N} \text{R.H.S.},$$

which gives (15).

**Step 4** (Convergence Analysis) In this step, we combine the results from the previous steps to establish an error bound for  $N^{-1} \sum_{k=1}^N \mathbb{E} \|\mathbf{g}_k\|^2$ . To do so, by substituting (15) and (11) into (8) and simplifying the terms, we obtain

$$\begin{aligned} \eta^2(1 - \beta_{1,1})^2 G_\infty^{-2} \frac{1}{N} \sum_{k=1}^N \mathbb{E} \|\mathbf{g}_k\|^2 &\leq \frac{6\eta^2 d \beta_{1,1}^2}{N u_c (1 - \kappa^2)} + \frac{6\eta^2 \sigma^2}{m G_0^2} + \frac{18D^2 d G_\infty}{N G_0} \\ &\quad + \frac{12\eta D}{N G_0} \left( \frac{\beta_{1,1} G_\infty}{1 - \kappa} \sqrt{\frac{d}{u_c}} + \frac{G_\infty^2 d}{G_0} \right) \\ &\quad + \frac{168\eta^2 d \beta_{1,1}^2 G_\infty}{N u_c (1 - \kappa^2) G_0} + \frac{84\eta^2 d G_\infty^3}{N G_0^3} \\ &\quad + \frac{1}{N} \sum_{i=1}^N \frac{84(1 - \kappa) \beta_{1,1}^2 \eta^2 G_\infty}{G_0^3 \kappa^2 (1 + \kappa)} \mathbb{E} \|\mathbf{g}_k\|^2 + \frac{168\eta^2 G_\infty \sigma^2}{G_0^3 m}. \end{aligned} \tag{26}$$

Here, we used the fact that

$$\mathbb{E} \left[ \sum_{k=1}^N \Theta_k \right] = 0, \quad \text{and} \quad \mathbb{E} \left[ \sum_{k=1}^N \|\epsilon_k\|^2 \right] = \frac{\sigma^2}{m},$$

by Assumptions 1 and 3, respectively.

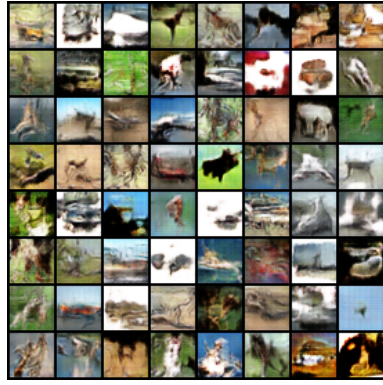
Now, it follows from (26) that

$$\frac{B_0}{N} \sum_{k=1}^N \mathbb{E} \|\mathbf{g}_k\|^2 \leq \frac{B_1}{N} + \frac{B_2 \sigma^2}{m},$$

where

$$\begin{aligned} B_0 &:= \eta^2(1 - \beta_{1,1})^2 G_\infty^{-2} - \frac{84(1 - \kappa) \beta_{1,1}^2 \eta^2 G_\infty}{G_0^3 \kappa^2 (1 + \kappa)}, \\ B_1 &:= \frac{6\eta^2 d \beta_{1,1}^2}{u_c (1 - \kappa^2)} + \frac{18D^2 d G_\infty}{G_0} + \frac{168\eta^2 d \beta_{1,1}^2 G_\infty}{u_c (1 - \kappa^2) G_0} \\ &\quad + \frac{84\eta^2 d G_\infty^3}{G_0^3} + \frac{12\eta D}{G_0} \left( \frac{\beta_{1,1} G_\infty}{1 - \kappa} \sqrt{\frac{d}{u_c}} + \frac{G_\infty^2 d}{G_0} \right), \\ B_2 &:= \frac{6\eta^2}{G_0^2} + \frac{168\eta^2 G_\infty}{G_0^3}. \end{aligned} \tag{27}$$

Next, define  $C = \frac{(1+\kappa)\kappa^2 G_0^3}{168(1-\kappa)G_\infty^3}$  and pick  $\beta_{1,1} \leq \frac{\sqrt{C}}{\sqrt{C}+1}$ . Then, dividing both sides by  $B_0$  gives us the desired result.



(a) ADAM<sup>3</sup>



(b) OAdagrad

Figure 3: Generated CIFAR-10 Samples

□