

# AdaInject: Injection Based Adaptive Gradient Descent Optimizers for Convolutional Neural Networks

Shiv Ram Dubey, *Senior Member, IEEE*, S.H. Shabbeer Basha, Satish Kumar Singh, *Senior Member, IEEE*, and Bidyut Baran Chaudhuri, *Life Fellow, IEEE*

**Abstract**—The convolutional neural networks (CNNs) are generally trained using stochastic gradient descent (SGD) based optimization techniques. The existing SGD optimizers generally suffer with the overshooting of the minimum and oscillation near minimum. In this paper, we propose a new approach, hereafter referred as AdaInject, for the gradient descent optimizers by injecting the second order moment into the first order moment. Specifically, the short-term change in parameter is used as a weight to inject the second order moment in the update rule. The AdaInject optimizer controls the parameter update, avoids the overshooting of the minimum and reduces the oscillation near minimum. The proposed approach is generic in nature and can be integrated with any existing SGD optimizer. The effectiveness of the AdaInject optimizer is explained intuitively as well as through some toy examples. We also show the convergence property of the proposed injection based optimizer. Further, we depict the efficacy of the AdaInject approach through extensive experiments in conjunction with the state-of-the-art optimizers, namely AdamInject, diffGradInject, RadamInject, and AdaBeliefInject on four benchmark datasets. Different CNN models are used in the experiments. A highest improvement in the top-1 classification error rate of 16.54% is observed using diffGradInject optimizer with ResNeXt29 model over the CIFAR10 dataset. Overall, we observe very promising performance improvement of existing optimizers with the proposed AdaInject approach. The code is available at: <https://github.com/shivram1987/AdaInject>.

**Impact Statement**—Adaptive moment based optimizers are among the popular gradient descent optimization techniques for the training of deep learning models. They try to control the step size based on the gradient behavior. However, the existing gradient descent optimization techniques either overshoot the “steep and narrow” valley (i.e., minimum) or oscillate near it, due to large step size caused by the exponential moving average of gradients used for parameter updates. The AdaInject optimization technique we introduce in this paper tackled this problem by incorporating the immediate parameter change weighted second order moment injection for the parameter updates. Using the proposed optimization technique, a significant improvement is observed in the performance of image classification using different

CNN models. Moreover, the proposed AdaInject approach can be used with any existing adaptive moment based optimization technique. Hence, it can provide the alternative optimizers with better step size control to train different deep learning models for diverse applications.

**Index Terms**—Adaptive Optimizers, Convolutional Neural Networks, Deep Learning, Image Recognition, Parameter Update History, Second Order Moment Injection, Stochastic Gradient Descent.

## I. INTRODUCTION

DEEP learning has shown a great impact over the performance of the neural networks for a wide range of problems [1]. In recent past, convolutional neural networks (CNNs) have shown very promising results for different computer vision applications, such as object recognition [2], [3], [4], [5]; object detection [6], [7]; face recognition [8], [9]; image quality assessment [10]; gesture recognition [11]; Covid-19 grading [12]; and many more. CNNs have also been used as basic building blocks in other networks like Autoencoder [13], [14], [15], Siamese Network [16], [17], Generative Adversarial Networks [18], [19], etc.

The training of different types of deep neural networks (DNNs) is mainly performed with the help of stochastic gradient descent (SGD) based optimization [20]. SGD optimizer updates the parameters of the network based on the gradient of objective function w.r.t. the corresponding parameters [21]. The vanilla SGD optimization suffers from three problems, including 1) zero gradient in local minimum and saddle regions leading to no update in the parameters, 2) a jittering effect along steep dimensions due to the inconsistent changes in the loss caused by the different parameters, and 3) noisy updates due to the gradient computed from the batch of data. SGD with moment (SGDM) [22] considers the first order moment (i.e., velocity) as an exponential moving average (EMA) of gradient for each parameter while training progresses [23]. The parameter is updated in SGDM based on the EMA of gradient which resolves the problem of zero gradient.

Several SGD based optimization techniques have been proposed in the recent past [24], [25], [26], [27], [28], [29], [30], and etc. AdaGrad [24] controls the learning rate by dividing it with the root of the sum of the squares of the past gradients. However, it makes the learning rate very small after certain iterations and kills the parameter update. AdaDelta [25] resolves the diminishing learning rate issue of AdaGrad by

S.R. Dubey and S.K. Singh are with the Computer Vision and Biometrics Laboratory (CVBL), Indian Institute of Information Technology, Allahabad, Prayagraj, Uttar Pradesh-211015, India (e-mail: srubey@iitaa.ac.in, sk.singh@iitaa.ac.in).

S.H.S. Basha is with the PathPartner Technology Pvt. Ltd., Bangalore, India (e-mail: shabbeer.sh@pathpartnertech.com).

B.B. Chaudhuri was with the Indian Statistical Institute, Kolkata, India and now associated with Techno India University, Kolkata, India (e-mail: bidyutbaranchaudhuri@gmail.com).

This paper is accepted for publication by IEEE Transactions on Artificial Intelligence. Copyright © 2022 IEEE. Personal use of this material is permitted. However, permission to use this material for any other purposes must be obtained from the IEEE by sending an email to [pubs-permissions@ieee.org](mailto:pubs-permissions@ieee.org).

**Algorithm 1: Adam Optimizer**


---

**Initialize:**  $\theta_0, m_0 \leftarrow 0, v_0 \leftarrow 0, t \leftarrow 0$   
**Hyperparameters:**  $\alpha, \beta_1, \beta_2$   
**While**  $\theta_t$  not converged  
 $t \leftarrow t + 1$   
 $g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$   
 $m_t \leftarrow \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t$   
 $v_t \leftarrow \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2$   
**Bias Correction**  
 $\widehat{m}_t \leftarrow m_t / (1 - \beta_1^t), \widehat{v}_t \leftarrow v_t / (1 - \beta_2^t)$   
**Update**  
 $\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \widehat{m}_t / (\sqrt{\widehat{v}_t} + \epsilon)$

---

considering only a few immediate past gradients. However, it is not able to exploit the global information. In another attempt to resolve the problem of AdaGrad, RMSProp [26] divides the learning rate by the root of the exponentially decaying average of squared gradients. In 2015, Kingma and Ba [27] proposed the adaptive moment based Adam optimizer. Adam combines the ideas of SGDM and RMSprop and uses first order and second order moments. Adam computes the first order moment as the EMA of gradients and uses it to update the parameter. Adam also computes the second order moment as the EMA of the square of gradients and uses it to control the learning rate. Adam performs well in practice to train the convolutional neural networks (CNNs) [27]. However, it suffers from overshooting and oscillations near minimum and varying gradient variance due to batch wise computation. diffGrad [28] resolved the issues as posed by Adam by introducing a friction term in parameter update using the rate of change in gradients. Radam [29] resolved the variance issue as posed by Adam by rectifying the variance of gradients during parameter update. AdaBelief [30] uses the belief in gradients to compute the second order moment. The belief in gradients is computed as the difference between the gradient and the first order moment of the gradient. Other recently proposed and notable gradient descent optimizers are Proportional Integral Derivative (PID) [31], Nesterov’s Moment Adam (NADAM) [32], Nostalgic Adam (NosAdam) [33], YOGI [34], Adaptive Bound (AdaBound) [35], Adaptive and Momental Bound (AdaMod) [36], Aggregated Moment (AggMo) [37], Lamb [38], Adam Projection (AdamP) [39], Gradient Centralization (GC) [40], AdaHessian [41], and AngularGrad [42].

The adaptive SGD optimization techniques have led to a promising performance on deep CNN models. The majority of the above mentioned adaptive gradient descent optimizers suffer due to the overshooting of the minimum and oscillation near minimum. However, it is evident that a robust online stepsize adaptation in optimization plays an important role in gradient descent optimization [43]. We resolve the above issues by injecting the second order moment in first order for the parameter update, which is weighted by the short-term parameter update history to incorporate the robust adaptation of step size. The major contribution of this work is summarized as follows:

**Algorithm 2: AdamInject (i.e., Adam + AdaInject) Optimizer**


---

**Initialize:**  $\theta_0, s_0 \leftarrow 0, v_0 \leftarrow 0, t \leftarrow 0$   
**Hyperparameters:**  $\alpha, \beta_1, \beta_2, k$   
**While**  $\theta_t$  not converged  
 $t \leftarrow t + 1$   
 $g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$   
**If**  $t = 1$   
 $s_t \leftarrow \beta_1 \cdot s_{t-1} + (1 - \beta_1) \cdot g_t$   
**Else**  
 $\Delta\theta \leftarrow \theta_{t-2} - \theta_{t-1}$   
 $s_t \leftarrow \beta_1 \cdot s_{t-1} + (1 - \beta_1) \cdot (g_t + \Delta\theta \cdot g_t^2) / k$   
 $v_t \leftarrow \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2$   
**Bias Correction**  
 $\widehat{s}_t \leftarrow s_t / (1 - \beta_1^t), \widehat{v}_t \leftarrow v_t / (1 - \beta_2^t)$   
**Update**  
 $\theta_t \leftarrow \theta_{t-1} - \alpha \widehat{s}_t / (\sqrt{\widehat{v}_t} + \epsilon)$

---

- We propose AdaInject for the adaptive optimizers by injecting the short-term parameter change weighted second order moment in EMA of gradient used for parameter update.
- We provide an intuitive explanation in support of the effectiveness of the proposed AdaInject in different optimization scenarios.
- We show the effect of the proposed approach using toy examples. The convergence analysis is also conducted using regret bound which shows the convergence property of the proposed approach.
- We validate the superiority of the proposed injection concept with the recent state-of-the-art optimizers, including Adam [27], diffGrad [28], Radam [29] and AdaBelief [30] using a wide range of CNN models for image classification over four benchmark datasets.
- The proposed concept is generic and can be easily integrated with any existing adaptive moment based SGD optimizer.

The remaining paper is structured by presenting the proposed Injection based optimizers in Section 2; Intuitive explanation and empirical analysis in Section 3; Convergence analysis in Section 4; Experimental analysis in Section 5; and Concluding remarks in Section 6.

## II. PROPOSED INJECTION BASED OPTIMIZERS

As per the conventions used in Adam [27], the aim of gradient descent optimization is to minimize the loss function  $f(\theta) \in \mathbb{R}$  where  $\theta \in \mathbb{R}^d$  is the parameter. The gradient ( $g_t$ ) at  $t^{th}$  step is computed as  $g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$ . Adam computes the first order moment ( $m_t$ ) and second order moment ( $v_t$ ) as the exponential moving average (EMA) of  $g_t$  and  $g_t^2$ , respectively, which can be written as,

$$m_t = \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t \quad (1)$$

$$v_t = \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2 \quad (2)$$

where  $\beta_1$  and  $\beta_2$  are the smoothing hyperparameters, typically set as  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$ . The  $g_t^2$  is computed as  $g_t \cdot g_t$  as in the Adam. A bias correction is performed as  $\widehat{m}_t \leftarrow m_t / (1 - \beta_1^t)$ ,  $\widehat{v}_t \leftarrow v_t / (1 - \beta_2^t)$  to avoid very large step size in the initial iterations. The parameter update rule in Adam [27] is given as,

$$\theta_t \leftarrow \theta_{t-1} - \alpha \widehat{m}_t / (\sqrt{\widehat{v}_t} + \epsilon) \quad (3)$$

where  $\alpha$  is the learning rate and  $\epsilon = 1e^{-8}$  is a small number for numerical stability to avoid division by zero. A detailed algorithm of Adam optimizer is summarized in Algorithm 1. The first order moment  $m_t$  is used to update the parameters in Adam wherein, the second order moment  $v_t$  is used to control the learning rate. It can be noticed that Adam mainly relies on the gradients.

However, the SGDM considers only the momentum to update the parameters as follows:

$$\theta_t \leftarrow \theta_{t-1} - \alpha m_t. \quad (4)$$

In order to utilize the parameter update history information during optimization, we propose a novel concept named AdaInject. Basically, we inject the short-term parameter change weighted second order moment into first order moment to compute the injected moment using the EMA of  $(g_t + \Delta\theta \cdot g_t^2)/k$  as,

$$s_t = \beta_1 \cdot s_{t-1} + (1 - \beta_1) \cdot (g_t + \Delta\theta \cdot g_t^2)/k \quad (5)$$

where  $\Delta\theta = \theta_{t-2} - \theta_{t-1}$  is the short-term change in parameter  $\theta$  to utilize the parameter history information and  $k$  is an injection controlling hyperparameter, typically set to 2 in the experiment. The injection of parameter history guided second order moment helps the optimizers to perform the smaller updates near minimum (i.e., “steep and narrow” valley) to avoid the overshooting and oscillation, while reasonably large updates are used in the small curvature regions. This phenomenon is depicted in Fig. 1 with a detailed analysis in the next section. We perform the bias correction of injected moment and second order moment as  $\widehat{s}_t \leftarrow s_t / (1 - \beta_1^t)$  and  $\widehat{v}_t \leftarrow v_t / (1 - \beta_2^t)$ , respectively.

The parameter ( $\theta$ ) update of AdamInject optimizer is given as,

$$\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \widehat{s}_t / (\sqrt{\widehat{v}_t} + \epsilon) \quad (6)$$

where  $\alpha$  is the learning rate and  $\epsilon = 1e^{-8}$  is a small number for numerical stability to avoid the division by zero. We refer to Adam optimizer with the proposed second order moment injection as AdamInject optimizer. A detailed algorithm of AdamInject optimizer is presented in Algorithm 2 with highlighted changes in blue color as compared to vanilla Adam optimizer which is shown in Algorithm 1.

Basically, we use the proposed AdaInject concept with four existing state-of-the-art optimizers, including Adam [27], diffGrad [28], Radam [29] and AdaBelief [30], and propose the corresponding AdamInject (i.e., Adam + AdaInject), diffGradInject (i.e., diffGrad + AdaInject), RadamInject (i.e., Radam and AdaInject) and AdaBeliefInject (i.e., AdaBelief + AdaInject) optimizers, respectively. The algorithms for different optimizers (i.e., without and with AdaInject), such as

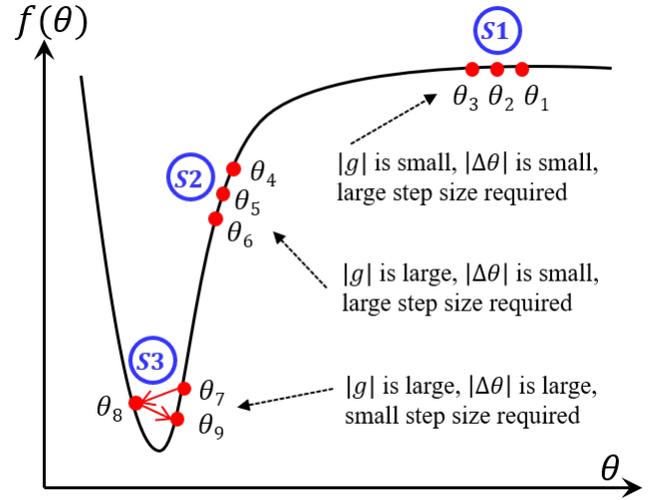


Fig. 1. A typical scenario in the optimization depicting the importance of adaptive parameter update in optimization [43], [30].

diffGrad, diffGradInject, Radam, RadamInject, AdaBelief, and AdaBeliefInject, are provided in Supplementary. Though we test the proposed injection concept with four optimizers, it can be extended to any EMA based gradient descent optimization technique. In the next section, we analyze the property of the proposed approach.

### III. INTUITIVE EXPLANATION AND EMPIRICAL ANALYSIS

In this section, we present an intuitive explanation using a one dimensional optimization landscape having three scenarios and an empirical analysis using three toy examples.

#### A. Intuitive Explanation

The existing gradient descent optimizers such as Adam, diffGrad, Radam, etc. only consider the EMA of gradient for parameter update. However, the consideration of parameter history is important as the gradient behavior and required stepsize are different for different regions of loss optimization landscape [43], [30]. We explain the advantage of the proposed optimizer by considering three typical scenarios using a one dimensional optimization curvature (i.e., S1, S2 and S3) as depicted in Fig. 1. The bias correction step is ignored in the explanation for simplicity.

**S1:** This scenario depicts a flat region on the optimization landscape. An ideal optimizer is expected to perform large parameter updates in this scenario. The  $|g_t|$  and  $|\Delta\theta|$  in flat region are small. Thus, the EMA of gradient (i.e.,  $m_t$ ) as well as the EMA of proposed injected gradient (i.e.,  $s_t$ ) are small. It leads to a small stepsize in SGD. However, the step size is sufficiently large in both Adam and AdaInject due to the small value of  $\sqrt{\widehat{v}_t}$  in the denominator.

**S2:** The “large gradient, small curvature” is another scenario in optimization landscape. The gradient  $|g_t|$  is higher in such regions. An ideal optimizer is expected to take the large parameter updates in such regions. The EMA of gradient (i.e.,  $m_t$ ) as well as squared gradient (i.e.,  $v_t$ ) are large. Moreover,

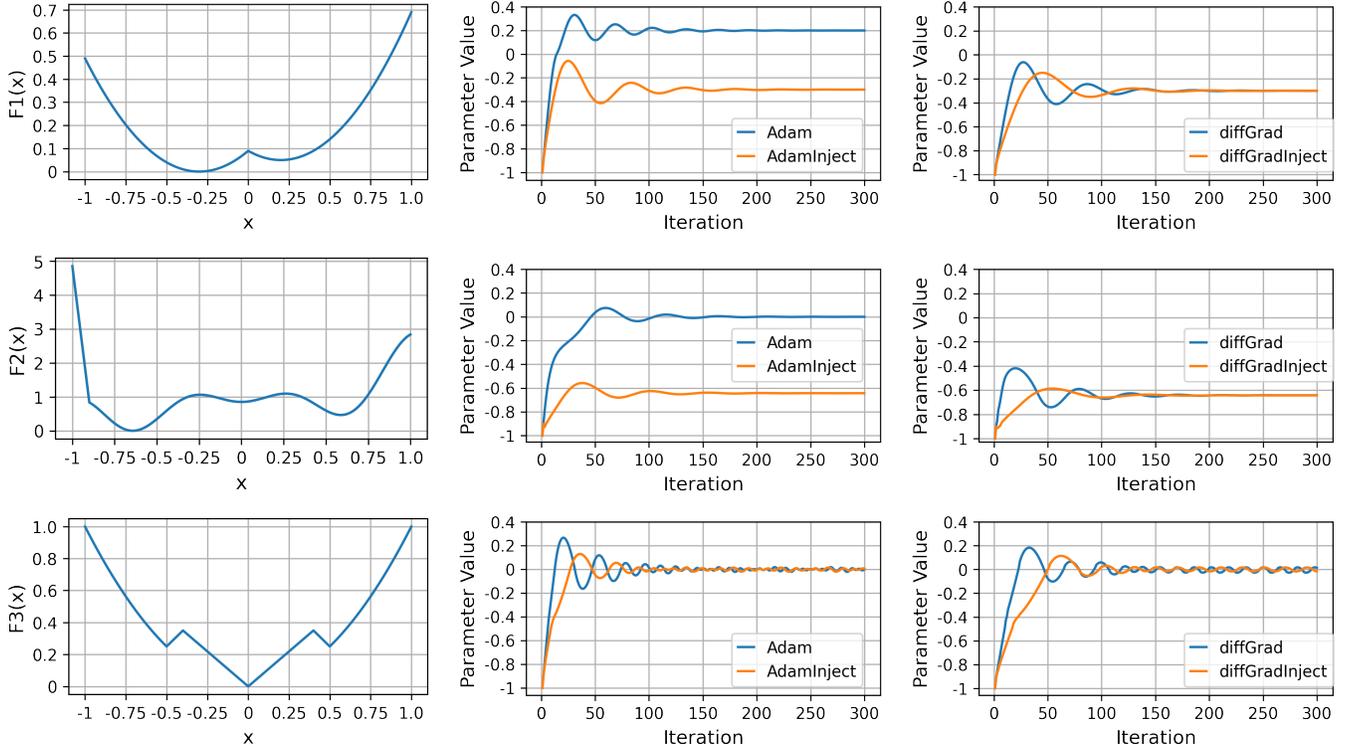


Fig. 2. The empirical results computed over three synthetic, non-convex functions as toy examples. Each row corresponds to a function. The 1<sup>st</sup> column shows the used functions. The 2<sup>nd</sup> and 3<sup>rd</sup> columns show the parameter value updates for 300 iterations using Adam and diffGrad optimizers, respectively, with and without the proposed second order injection. The regression loss based objective function is used to update the parameters. The initialization is done at  $x = -1$ .

the EMA of the proposed injected gradient (i.e.,  $s_t$ ) is also sufficiently large as  $|\Delta\theta|$  is small. Hence, the SGD takes large step in this scenario. However, both the Adam and AdaInject take relatively smaller step due to large value of  $\sqrt{v_t}$  in the denominator. But, we show experimentally that this problem can be reduced by considering AdaBelief concept [30] with the proposed injection idea (i.e., AdaBeliefInject).

**S3:** The third scenario is parameter updates near “steep and narrow” valley (i.e., minimum). It is expected for an ideal optimizer to decrease the step size for parameter updates in this scenario to avoid the overshooting as well as to reduce the oscillation near the valley. The proposed AdaInject optimizer is very beneficial in this scenario too. The gradient  $|g_t|$  is large in this scenario, hence  $m_t$  and  $\sqrt{v_t}$  are also large. The SGD suffers due to large value of  $m_t$ . This problem is reduced to a certain extent in Adam due to large value of  $\sqrt{v_t}$  in the denominator. In this scenario,  $|\Delta\theta|$  is large,  $\Delta\theta < 0$  when  $g_t > 0$  and  $\Delta\theta > 0$  when  $g_t < 0$ , leading to  $|s_t| < |m_t|$  (Note that  $t$  is expected not to be the initial iterations near minimum, rather sufficiently large). Hence, the proposed AdaInject method reduces  $s_t$  while enjoying the benefits of Adam (i.e., large  $\sqrt{v_t}$  in denominator) leading to a reduced step size, which avoids the overshooting and oscillation near minimum to a greater extent. In order to show this effect using toy examples, we conduct an empirical study with the help of synthetic, non-convex functions in the next subsection.

### B. Empirical Analysis using Toy Examples

We perform the empirical analysis using three synthetic, one-dimensional, non-convex functions by following the protocol of diffGrad [28]. These functions are given as:

$$F1(x) = \begin{cases} (x + 0.3)^2, & \text{for } x \leq 0 \\ (x - 0.2)^2 + 0.05, & \text{for } x > 0 \end{cases} \quad (7)$$

$$F2(x) = \begin{cases} -40x - 35.15, & \text{for } x \leq -0.9 \\ x^3 + x \sin(8x) + 0.85, & \text{for } x > -0.9 \end{cases} \quad (8)$$

$$F3(x) = \begin{cases} x^2, & \text{for } x \leq -0.5 \\ 0.75 + x, & \text{for } -0.5 < x \leq -0.4 \\ -7x/8, & \text{for } -0.4 < x \leq 0 \\ 7x/8, & \text{for } 0 < x \leq 0.4 \\ 0.75 - x, & \text{for } 0.4 < x \leq 0.5 \\ x^2, & \text{for } 0.5 < x \end{cases} \quad (9)$$

where  $-\infty < x < +\infty$  is the input. Functions  $F1$ ,  $F2$ , and  $F3$  are illustrated in Fig. 2 in the 1<sup>st</sup> column and in the 1<sup>st</sup>, 2<sup>nd</sup>, and 3<sup>rd</sup> rows, respectively, for  $-1 < x < +1$ . The parameter  $x$  is initialized at  $-1$ . The experiment is performed by computing the regression loss as the objective function. The 2<sup>nd</sup> column shows the parameter values at different iterations using Adam and AdamInject optimizers. Similarly, the 3<sup>rd</sup> column illustrates the parameter values at different iterations using diffGrad and diffGradInject optimizers. It can

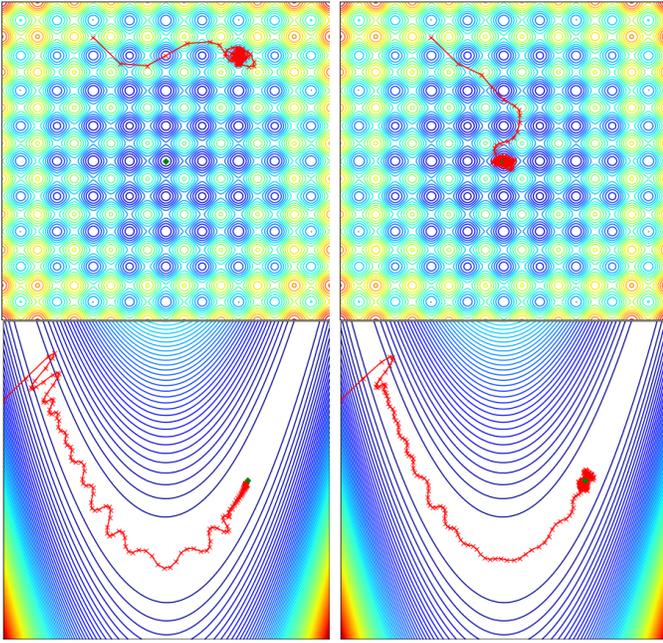


Fig. 3. The optimization illustration for Rastrigin (upper row) and Rosenbrock (lower row) functions using Adam (left column) and AdamInject (right column).

be noticed that Adam overshoots the minimum for both  $F1$  and  $F2$  functions, whereas AdamInject is able to avoid the overshooting due to the small step size caused by the proposed parameter change weighted second order moment injection in parameter update. In other cases, including Adam and AdamInject for  $F3$  function and diffGrad and diffGradInject for all three functions, the effect of the proposed optimizer can be easily observed in terms of the smooth parameter updates and less oscillations near minimum by accumulating the injected momentum in an accurate direction. It is noticed that AdaInject is more effective with Adam than diffGrad as diffGrad utilizes the short-term gradient change as friction coefficient. These results confirm that the proposed parameter change guided second moment injection leads to accurate and precise parameter updates, especially near “steep and narrow” valley.

In order to demonstrate the effect of the proposed optimizer on 2-dimensional optimization, we consider non-convex Rastrigin and Rosenbrock functions<sup>1</sup>, which are widely used to show the optimization characteristics. The Rastrigin function has one global minimum at  $(0.0, 0.0)$ . However, the Rosenbrock has one global minimum at  $(1.0, 1.0)$ . The optimization trajectories using Adam and AdamInject optimizers under the same experimental setup are depicted in Fig. 3. It can be noticed that Adam is not able to converge over the Rastrigin function due to the presence of several local minima. Whereas, the AdamInject is able to converge over the Rastrigin function due to the improved parameter updates caused by the second order moment injection. It is also observed that the Adam optimizer takes more steps to reach the minimum over the Rosenbrock function due to irregular parameter updates caused

by long, narrow, parabolic shaped flat valley. However, the AdamInject optimizer is able to tackle this issue and reaches the minimum in less number of steps over the Rosenbrock function.

#### IV. CONVERGENCE ANALYSIS

We use the online learning framework to show the convergence property of the proposed injection based AdamInject optimizer, similar to Adam [27]. Let’s represent the unknown sequence of convex cost functions as  $f_1(\theta)$ ,  $f_2(\theta), \dots, f_T(\theta)$ . We want to estimate parameter  $\theta_t$  at each iteration  $t$  and assess over  $f_t(\theta)$ . The regret bound is commonly used in such scenarios to assess the algorithm where the information of the sequence is not known in advance. The sum of the difference between all the previous online guesses  $f_t(\theta_t)$  and the best fixed point parameter  $f_t(\theta^*)$  from a feasible set  $\chi$  of all the previous iterations are used to compute the regret bound. The regret bound is given as,

$$R(T) = \sum_{t=1}^T [f_t(\theta_t) - f_t(\theta^*)] \quad (10)$$

where  $\theta^* = \arg \min_{\theta \in \chi} \sum_{t=1}^T f_t(\theta)$ . The regret bound for the proposed injection based AdamInject is noticed as  $O(\sqrt{T})$  which is the same as Adam and is comparable to general convex online learning approaches. We provide the proof in the Supplementary. We consider  $g_{t,i}$  as the gradient for the  $i^{th}$  element in the  $t^{th}$  iteration,  $g_{1:t,i} = [g_{1,i}, g_{2,i}, \dots, g_{t,i}] \in \mathbb{R}^t$  is the gradient vector in the  $i^{th}$  dimension up to  $t^{th}$  iterations, and  $\gamma \triangleq \frac{\beta_1^2}{\sqrt{\beta_2}}$ .

**Theorem 1.** Assume that the gradients for function  $f_t$  (i.e.,  $\|g_{t,\theta}\|_2 \leq G$  and  $\|g_{t,\theta}\|_\infty \leq G_\infty$ ) are bounded for all  $\theta \in \mathbb{R}^d$ . Let also consider that the bounded distance is generated by the proposed optimizer between any  $\theta_t$  (i.e.,  $\|\theta_n - \theta_m\|_2 \leq D$  and  $\|\theta_n - \theta_m\|_\infty \leq D_\infty$  for any  $m, n \in \{1, \dots, T\}$ ). Let  $\gamma \triangleq \frac{\beta_1^2}{\sqrt{\beta_2}}$ ,  $\beta_1, \beta_2 \in [0, 1)$  satisfy  $\frac{\beta_1^2}{\sqrt{\beta_2}} < 1$ ,  $\alpha_t = \frac{\alpha}{\sqrt{t}}$ , and  $\beta_{1,t} = \beta_1 \lambda^{t-1}$ ,  $\lambda \in (0, 1)$  with  $\lambda$  is around 1, e.g  $1 - 10^{-8}$ . For all  $T \geq 1$ , the proposed injection based AdamInject shows the following guarantee as derived in the Supplementary:

$$\begin{aligned} R(T) \leq & \frac{D^2}{\alpha(1-\beta_1)} \sum_{i=1}^d \sqrt{T \hat{v}_{T,i}} \\ & + \frac{2\alpha(1+\beta_1)G_\infty}{(1-\beta_1)\sqrt{1-\beta_2}(1-\gamma)^2} \sum_{i=1}^d \|g_{1:T,i}\|_2 \\ & + \sum_{i=1}^d \frac{D_\infty^2 G_\infty \sqrt{1-\beta_2}}{\alpha(1-\beta_1)(1-\lambda)^2} + 4D_\infty G_\infty^2 \sum_{i=1}^d \|g_{1:T,i}\|_2^2 \end{aligned} \quad (11)$$

The aggregation terms over the dimension ( $d$ ) can be very small as compared to the corresponding upper bounds, such as  $\sum_{i=1}^d \|g_{1:T,i}\|_2 \ll dG_\infty \sqrt{T}$ ,  $\sum_{i=1}^d \|g_{1:T,i}\|_2^2 \ll dG_\infty \sqrt{T}$  and  $\sum_{i=1}^d \sqrt{T \hat{v}_{T,i}} \ll dG_\infty \sqrt{T}$ . The adaptive methods such as the proposed optimizers and Adam show the upper bound as  $O(\log d \sqrt{T})$ , which is better than  $O(\sqrt{dT})$  of non-adaptive optimizers. By following the convergence analysis of Adam [27], we also use the decay of  $\beta_{1,t}$ .

<sup>1</sup><https://github.com/jettify/pytorch-optimizer>

TABLE I

THE EXPERIMENTAL RESULTS OF DIFFERENT CNNs IN TERMS OF TOP-1 CLASSIFICATION ERROR (%) OVER THE CIFAR10 DATASET USING DIFFERENT OPTIMIZERS, WITHOUT AND WITH THE PROPOSED ADAINJECT. THE RESULTS WITH THE PROPOSED APPROACH ARE HIGHLIGHTED IN BOLD. THE IMPROVEMENT IN THE ERROR DUE TO THE PROPOSED INJECTION CONCEPT IS ALSO MENTIONED. THE HIGHEST INCREASE FOR AN OPTIMIZER IS ALSO HIGHLIGHTED IN BOLD. THE SYMBOLS  $\uparrow$  AND  $\downarrow$  REPRESENT THE IMPROVEMENT AND DEGRADATION IN %, RESPECTIVELY, IN THE TOP-1 ERROR. WE FOLLOW THE SAME CONVENTION IN THE RESULTS REPORTED IN TABLE II AND III ALSO. THESE RESULTS ARE COMPUTED AS THE AVERAGE OVER THREE INDEPENDENT TRIALS.

CNN Models	Classification error (%) using different optimizers without and with AdaInject							
	Adam		diffGrad		Radam		AdaBelief	
	Adam	AdamInject	diffGrad	diffGradInject	Radam	RadamInject	AdaBelief	AdaBeliefInject
VGG16	7.45	<b>7.20</b> ( $\uparrow$ 3.36)	7.24	<b>7.04</b> ( $\uparrow$ 2.76)	7.06	<b>6.88</b> ( $\uparrow$ 2.55)	7.29	<b>7.07</b> ( $\uparrow$ 3.02)
ResNet18	6.46	<b>6.20</b> ( $\uparrow$ 4.02)	6.51	<b>6.10</b> ( $\uparrow$ 6.30)	6.18	<b>5.87</b> ( $\uparrow$ 5.02)	6.37	<b>6.30</b> ( $\uparrow$ 1.10)
SENet18	6.61	<b>6.29</b> ( $\uparrow$ 4.84)	6.44	<b>6.21</b> ( $\uparrow$ 3.57)	6.05	<b>5.83</b> ( $\uparrow$ 3.64)	6.59	<b>6.23</b> ( $\uparrow$ 5.46)
ResNet50	6.17	<b>5.89</b> ( $\uparrow$ 4.54)	6.19	<b>5.73</b> ( $\uparrow$ 7.43)	5.86	<b>5.29</b> ( $\uparrow$ 9.73)	5.90	<b>5.78</b> ( $\uparrow$ 2.03)
ResNet101	6.90	<b>6.01</b> ( $\uparrow$ <b>12.90</b> )	6.45	<b>5.69</b> ( $\uparrow$ 11.78)	6.29	<b>5.76</b> ( $\uparrow$ 8.43)	6.37	<b>6.03</b> ( $\uparrow$ 5.34)
ResNeXt29	6.79	<b>6.16</b> ( $\uparrow$ 9.28)	6.83	<b>5.70</b> ( $\uparrow$ <b>16.54</b> )	6.00	<b>5.67</b> ( $\uparrow$ 5.50)	6.43	<b>5.99</b> ( $\uparrow$ <b>6.84</b> )
DenseNet121	6.30	<b>5.63</b> ( $\uparrow$ 10.63)	5.90	<b>5.43</b> ( $\uparrow$ 7.97)	5.25	<b>5.10</b> ( $\uparrow$ 2.86)	6.05	<b>5.64</b> ( $\uparrow$ 6.78)

We show the convergence of average regret of AdamInject in below corollary with the help of the above theorem and  $\sum_{i=1}^d \|g_{1:T,i}\|_2 \ll dG_\infty\sqrt{T}$ ,  $\sum_{i=1}^d \|g_{1:T,i}\|_2^2 \ll dG_\infty\sqrt{T}$  and  $\sum_{i=1}^d \sqrt{T}\hat{v}_{T,i} \ll dG_\infty\sqrt{T}$ .

**Corollary 1.** Assume that the gradients for function  $f_t$  (i.e.,  $\|g_{t,\theta}\|_2 \leq G$  and  $\|g_{t,\theta}\|_\infty \leq G_\infty$ ) are bounded for all  $\theta \in R^d$ . Let also consider that the bounded distance is generated by the proposed optimizer between any  $\theta_t$  (i.e.,  $\|\theta_n - \theta_m\|_2 \leq D$  and  $\|\theta_n - \theta_m\|_\infty \leq D_\infty$  for any  $m, n \in \{1, \dots, T\}$ ). For all  $T \geq 1$ , the proposed injection based AdamInject optimizer shows the following guarantee:

$$\frac{R(T)}{T} = O\left(\frac{1}{\sqrt{T}}\right). \quad (12)$$

Thus,  $\lim_{T \rightarrow \infty} \frac{R(T)}{T} = 0$ .

## V. EXPERIMENTAL ANALYSIS

In this section, first we describe the experimental setup. Then, we present the detailed results using different optimizers. Finally, we analyze effects of the hyperparameters.

### A. Experimental Setup

We use a wide range of CNN models (i.e., VGG16 [44], ResNet18, ResNet50, ResNet101 [2], SENet18 [45], ResNeXt29 [3] and DenseNet121 [4]) to test the suitability of the proposed AdaInject concept for optimizers. We follow the publicly available Pytorch implementation<sup>2</sup> of these CNN models. For ResNeXt29 model, we set the cardinality as 4 and bottleneck width as 64. We train all the CNN models using all the optimizers under the same experimental setup. The training is performed for 100 epochs with a batch size of 64 for CIFAR10/100 and FashionMNIST and 256 for TinyImageNet dataset. The learning rate is set to 0.001 for the first 80 epochs and 0.0001 for the last 20 epochs. Different computers are used for the experiments, including Google colab<sup>3</sup>. We performed a random crop and random horizontal flip over

training data. The normalization is performed for both training and test data.

In order to demonstrate the efficacy of the proposed AdaInject based optimizers experimentally, we use four benchmark object recognition dataset, including CIFAR10 [46], CIFAR100 [46], FashionMNIST [47], and TinyImageNet<sup>4</sup> [48]. We use CIFAR and FashionMNIST datasets directly from the PyTorch library. CIFAR10 dataset consists of a total 60,000 images of dimension  $32 \times 32 \times 3$  from 10 object classes with 6,000 images per class. The training set contains 50,000 images with 5,000 images per class and the test set contains 10,000 images with 1,000 images per class in CIFAR10. CIFAR100 dataset contains all the images of CIFAR10, but is categorized into 100 classes. Thus, CIFAR100 dataset contains 50,000 training images and 10,000 test images with 500 and 100 images per class, respectively. FashionMNIST dataset contains 70,000 labeled fashion images of dimension  $28 \times 28$  from 10 categories. The training and test sets consist of 60,000 and 10,000 images, respectively. TinyImageNet dataset [48] is a subset of the large-scale visual recognition ImageNet challenge [49]. This dataset consists of the images from 200 object classes with 1,00,000 images in the training set (i.e., 500 images in each class) and 10,000 images in the validation set (i.e., 50 images in each class).

### B. Experimental Results

We compare the performance using four recent state-of-the-art adaptive gradient descent optimizers (i.e., Adam [27], diffGrad [28], Radam [29] and AdaBelief [30]), without and with the proposed injection approach. We consider VGG16 [44], ResNet18, ResNet50, ResNet101 [2], SENet18 [45], ResNeXt29 [3] and DenseNet121 [4] CNN models. The experimental results over the CIFAR10 dataset are depicted in Table I in terms of the error rate. It is observed that the performance of all CNN models is improved with AdaInject based optimizers as compared to its performance with corresponding vanilla optimizers. The RadamInject optimizer leads to best performance using the DenseNet121 model with a 5.10% error rate in classification. The highest improvement is reported by the ResNeXt29 model using diffGradInject.

<sup>2</sup><https://github.com/kuangliu/pytorch-cifar>

<sup>3</sup><https://colab.research.google.com/>

<sup>4</sup><http://cs231n.stanford.edu/tiny-imagenet-200.zip>

TABLE II

THE EXPERIMENTAL RESULTS OF DIFFERENT CNNs IN TERMS OF TOP-1 CLASSIFICATION ERROR (%) OVER THE CIFAR100 DATASET USING DIFFERENT OPTIMIZERS, WITHOUT AND WITH THE PROPOSED ADAINJECT. THESE RESULTS ARE COMPUTED AS THE AVERAGE OVER THREE INDEPENDENT TRIALS.

CNN Models	Classification error (%) using different optimizers without and with AdaInject							
	Adam	AdamInject	diffGrad	diffGradInject	Radam	RadamInject	AdaBelief	AdaBeliefInject
VGG16	32.71	<b>31.81</b> (↑ 2.75)	31.81	<b>30.80</b> (↑ 3.18)	<b>29.31</b>	30.07 (↓ 2.59)	31.08	<b>30.04</b> (↑ 3.35)
ResNet18	28.91	<b>27.28</b> (↑ 5.64)	26.50	<b>26.23</b> (↑ 1.02)	26.78	<b>25.84</b> (↑ 3.51)	27.28	<b>26.31</b> (↑ 3.56)
SENet18	29.15	<b>28.74</b> (↑ 1.41)	28.60	<b>27.64</b> (↑ 3.36)	27.66	<b>26.63</b> (↑ 3.72)	26.90	<b>26.52</b> (↑ 1.41)
ResNet50	28.12	<b>25.44</b> (↑ 9.53)	24.94	<b>24.18</b> (↑ 3.05)	25.05	<b>24.13</b> (↑ 3.67)	24.47	<b>24.25</b> (↑ 0.90)
ResNet101	25.78	<b>23.98</b> (↑ 6.98)	26.58	<b>24.17</b> (↑ <b>9.07</b> )	25.74	<b>23.83</b> (↑ 7.42)	<b>24.12</b>	24.24 (↓ 0.50)
ResNeXt29	28.78	<b>24.96</b> (↑ <b>13.27</b> )	25.47	<b>24.53</b> (↑ 3.69)	24.66	<b>22.74</b> (↑ 7.79)	24.61	<b>23.63</b> (↑ <b>3.98</b> )
DenseNet121	26.40	<b>24.33</b> (↑ 7.84)	24.14	<b>23.66</b> (↑ 1.99)	25.17	<b>23.06</b> (↑ <b>8.38</b> )	24.68	<b>24.06</b> (↑ 2.51)

TABLE III

THE EXPERIMENTAL RESULTS OF DIFFERENT CNNs IN TERMS OF TOP-1 CLASSIFICATION ERROR (%) OVER THE FASHIONMNIST DATASET USING DIFFERENT OPTIMIZERS, WITHOUT AND WITH THE PROPOSED ADAINJECT. THESE RESULTS ARE COMPUTED AS THE AVERAGE OVER THREE INDEPENDENT TRIALS.

CNN Models	Classification error (%) using different optimizers without and with AdaInject							
	Adam	AdamInject	diffGrad	diffGradInject	Radam	RadamInject	AdaBelief	AdaBeliefInject
VGG16	5.15	<b>5.01</b> (↑ 2.72)	5.13	<b>5.03</b> (↑ 1.95)	5.11	<b>5.07</b> (↑ 0.78)	5.12	<b>4.97</b> (↑ 2.93)
ResNet18	4.76	<b>4.74</b> (↑ 0.42)	4.82	<b>4.65</b> (↑ 3.53)	4.78	<b>4.67</b> (↑ 2.30)	4.95	<b>4.75</b> (↑ 4.04)
SENet18	5.14	<b>4.95</b> (↑ 3.70)	5.11	<b>4.79</b> (↑ 6.26)	5.08	<b>4.79</b> (↑ 5.71)	5.06	<b>4.91</b> (↑ 2.96)
ResNet50	5.10	<b>4.76</b> (↑ 6.67)	4.93	<b>4.77</b> (↑ 3.25)	4.98	<b>4.84</b> (↑ 2.81)	5.10	<b>4.78</b> (↑ 6.27)
ResNet101	4.94	<b>4.65</b> (↑ 5.87)	5.05	<b>4.73</b> (↑ 6.34)	4.91	<b>4.64</b> (↑ 5.50)	5.21	<b>4.69</b> (↑ <b>9.98</b> )
ResNeXt29	6.16	<b>5.59</b> (↑ <b>9.25</b> )	5.92	<b>5.16</b> (↑ <b>12.84</b> )	5.78	<b>5.37</b> (↑ <b>7.09</b> )	5.25	<b>4.90</b> (↑ 6.67)
DenseNet121	4.88	<b>4.69</b> (↑ 3.89)	4.77	<b>4.70</b> (↑ 1.47)	4.89	<b>4.68</b> (↑ 4.29)	4.68	<b>4.56</b> (↑ 2.56)

TABLE IV

THE EXPERIMENTAL RESULTS OF VGG16, RESNET18, AND SENET MODELS IN TERMS OF TOP-1 CLASSIFICATION ACCURACY (%) OVER THE TINYIMAGENET DATASET USING DIFFERENT OPTIMIZERS. THESE RESULTS ARE COMPUTED AS THE AVERAGE OVER THREE INDEPENDENT TRIALS.

CNN Models	Accuracy (%) using different optimizers without and with AdaInject							
	Adam	AdamInject	diffGrad	diffGradInject	Radam	RadamInject	AdaBelief	AdaBeliefInject
VGG16	44.05	<b>44.58</b> (↑ 1.20)	46.00	<b>47.18</b> (↑ 2.57)	45.92	<b>46.38</b> (↑ 1.00)	47.88	<b>48.25</b> (↑ 0.77)
ResNet18	50.58	<b>51.90</b> (↑ 2.61)	52.04	<b>52.37</b> (↑ 0.63)	52.12	<b>52.50</b> (↑ 0.73)	52.05	<b>52.74</b> (↑ 1.33)
SENet18	48.04	<b>49.52</b> (↑ 3.08)	49.51	<b>50.28</b> (↑ 1.56)	50.73	<b>51.01</b> (↑ 0.55)	51.76	<b>51.94</b> (↑ 0.35)

Moreover, the performance of the ResNeXt29 model is also significantly improved using AdaBeliefInject. In general, we observe better performance gain by heavy CNN models.

The results over the CIFAR100 dataset are illustrated in Table II. The best performance of 77.26% accuracy is achieved by the RadamInject optimizer using the ResNeXt29 model. The performance of ResNeXt29 is improved significantly using the proposed injection for optimizers with highest improvement by AdamInject. The results due to the proposed injection based optimizers are improved using all the CNN models except RadamInject using VGG16 and AdaBeliefInject using ResNet101. Note that Radam does not use second order moment if rectification criteria is not met and AdaBelief reduces the second order moment. These could be the possible reasons that the performance of RadamInject and AdaBeliefInject is marginally down in some cases. A very similar trend is also noticed over FashionMNIST (FMNIST) dataset in Table III, where the performance using the proposed approach is improved in all the cases. The best accuracy of 95.44% is observed for the AdaBeliefInject optimizer using the DenseNet121 model. An outstanding improvement in top-1 error is perceived for the ResNeXt29 model over the FashionMNIST dataset using the optimizers with the proposed AdaInject concept. The performance of other models is also significantly improved due to the proposed injection approach.

TABLE V

ACCURACY (%) USING ADAMINJECT OPTIMIZER WITH DIFFERENT VALUES OF  $k$ . RESULTS ARE COMPUTED AS THE AVERAGE OVER THREE INDEPENDENT TRIALS. NOTE THAT BEST AND SECOND BEST RESULTS ARE HIGHLIGHTED IN BOLD AND UNDERLINE, RESPECTIVELY.

Model	Dataset	$k=1$	$k=2$	$k=3$	$k=4$	$k=5$	$k=10$	$k=20$	$k=50$
VGG16	CIFAR10	92.68	<b>92.80</b>	92.78	92.50	92.61	91.85	90.75	87.76
	CIFAR100	67.04	68.19	<u>67.53</u>	<u>68.29</u>	<b>68.48</b>	68.25	66.97	61.08
	MNIST	<u>94.94</u>	<b>94.99</b>	94.91	94.91	94.93	94.8	94.47	94.15
	TinyImageNet	42.46	<b>44.58</b>	43.14	42.71	<u>44.21</u>	44.04	44.05	40.55
ResNet18	CIFAR10	93.71	<u>93.80</u>	<b>93.88</b>	93.75	93.76	93.41	92.24	89.36
	CIFAR100	71.97	72.72	<u>73.30</u>	73.26	<b>73.59</b>	72.76	69.99	64.43
	MNIST	95.15	<u>95.26</u>	95.19	<b>95.33</b>	95.18	95.16	94.84	94.52
	TinyImageNet	49.09	<b>51.90</b>	48.47	50.43	<u>51.17</u>	50.11	49.37	43.64

We also perform the experiment over the TinyImageNet dataset using VGG16, ResNet18 and SENet18 models and show the results in terms of the classification accuracy % in Table IV for different optimizers with and without the proposed injection concept. It is observed from this experiment that the proposed approach is able to improve the performance of the existing optimizers over large scale dataset as well. These results confirm that the proposed injection updates the parameter in an optimal way by utilizing the short-term parameter update information with second order moment.

TABLE VI

ACCURACY (%) USING ADAMINJECT OPTIMIZER WITH DIFFERENT BATCH SIZE (BS) AND LEARNING RATE (LR). RESULTS ARE COMPUTED AS THE AVERAGE OVER THREE INDEPENDENT TRIALS. NOTE THAT BEST RESULTS ARE HIGHLIGHTED IN BOLD.

Model	Dataset	Batch Size (BS)			Learning Rate (LR)		
		32	64	128	0.0001	0.001	0.01
VGG16	CIFAR10	92.46	<b>92.80</b>	92.45	91.16	<b>92.80</b>	92.45
	CIFAR100	67.00	<b>68.19</b>	68.03	66.92	<b>68.19</b>	66.97
	MNIST	94.90	<b>94.99</b>	94.96	94.75	<b>94.99</b>	94.85
	TinyImageNet	41.67	<b>44.58</b>	42.69	<b>45.18</b>	44.58	39.33
ResNet18	CIFAR10	93.94	<b>94.13</b>	93.71	92.36	<b>94.13</b>	93.65
	CIFAR100	73.05	<b>74.16</b>	72.76	70.79	<b>74.16</b>	67.31
	MNIST	95.27	<b>95.33</b>	95.25	95.13	<b>95.33</b>	95.26
	TinyImageNet	49.63	<b>51.90</b>	50.18	49.30	<b>51.90</b>	45.96

### C. Effect of Injection Hyperparameter ( $k$ )

In the previous results, we use the value of the injection hyperparameter ( $k$ ) as 2. We show a performance comparison by considering the value of  $k$  as 1, 2, 3, 4, 5, 10, 20, and 50 in Table V. The results are presented using the AdamInject optimizer for VGG16 and ResNet18 models over the CIFAR10, CIFAR100, FMNIST, and TinyImageNet datasets. It is noticed that  $k = 2$  is better suitable for the VGG16 model on CIFAR10, MNIST and TinyImageNet datasets. Moreover, the accuracy using  $k = 2$  is also either best or second best for the ResNet18 model on CIFAR10, MNIST and TinyImageNet datasets. It is also evident that the results on fine-grained CIFAR100 dataset are best using  $k = 5$  for both VGG16 and ResNet18 models. It is suggested to consider the value of  $k \in \{2, 3, 4, 5\}$ . The original selection of the value of  $k$  as 2 is also justified from this analysis.

### D. Effect of Batch Size and Learning Rate

In the previous experiments, the batch size (BS) and learning rate (LR) was set to 64 and 0.001, respectively. In this experiment, we analyze the impact of batch size and learning rate as detailed in Table VI. The results are reported for VGG16 and ResNet18 models on CIFAR10, CIFAR100, MNIST and TinyImageNet datasets. The batch size is considered as 32, 64 and 128, respectively. It is evident from the results that the batch size as 64 is better suitable with the proposed AdamInject optimizer in all the cases. The learning rate is considered as 0.0001, 0.001 and 0.01, respectively. Note that the learning rate is divided by 10 once in all the cases after 80 epochs of training for a fair comparison. It is noticed that the proposed optimizer performs best for 0.001 learning rate in almost all the cases. This analysis confirms the suitability of original batch size (i.e., 64) and learning rate (i.e., 0.001) choices used for the experiments.

## VI. CONCLUSION

In this paper, we present a novel and generic injection based EMA of gradients for parameter update by utilizing the parameter change information along with the second order moment. The proposed injection approach leads to an accurate and precise update by performing smaller updates near minimum to avoid the overshooting as well as oscillation and reasonably higher updates in the small curvature regions. The effect of

the proposed injection based optimizers is observed using toy examples. The convergence property of the proposed optimizer is also analyzed. The object recognition results for different CNN models over benchmark datasets using four optimizers show the superiority of the proposed injection concept. It is noticed that the injection hyperparameter as 2 yields to better results in majority of the cases using the AdamInject optimizer. It is also noted that the batch size as 64 and learning rate as 0.001 are better suitable with the proposed AdamInject optimizer. The intuitive explanation, empirical, convergence, and experimental analyses are evident that the proposed injection based optimizers lead to better optimization of CNNs by avoiding the overshooting of the minimum and reducing the oscillation near minimum to a greater extent.

## REFERENCES

- [1] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [2] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [3] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He, "Aggregated residual transformations for deep neural networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1492–1500.
- [4] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 4700–4708.
- [5] J. Kantipudi, S. R. Dubey, and S. Chakraborty, "Color channel perturbation attacks for fooling convolutional neural networks and a defense against such attacks," *IEEE Transactions on Artificial Intelligence*, vol. 1, no. 2, pp. 181–191, 2020.
- [6] R. Girshick, "Fast r-cnn," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 1440–1448.
- [7] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *Proceedings of the Advances in neural information processing systems*, 2015, pp. 91–99.
- [8] F. Schroff, D. Kalenichenko, and J. Philbin, "Facenet: A unified embedding for face recognition and clustering," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 815–823.
- [9] Y. Taigman, M. Yang, M. Ranzato, and L. Wolf, "Deepface: Closing the gap to human-level performance in face verification," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2014, pp. 1701–1708.
- [10] Z. Pan, F. Yuan, X. Wang, L. Xu, S. Xiao, and S. Kwong, "No-reference image quality assessment via multi-branch convolutional neural networks," *IEEE Transactions on Artificial Intelligence*, 2022.
- [11] Y. Zou and L. Cheng, "A transfer learning model for gesture recognition based on the deep features extracted by cnn," *IEEE Transactions on Artificial Intelligence*, vol. 2, no. 5, pp. 447–458, 2021.
- [12] C. De Vente, L. Boulogne, K. V. Venkatesh, C. Sital, N. Lessmann, C. Jacobs, C. I. S. Gutierrez, and B. Van Ginneken, "Automated covid-19 grading with convolutional neural networks in computed tomography scans: A systematic comparison," *IEEE Transactions on Artificial Intelligence*, vol. 1, no. 01, pp. 1–1, 2021.
- [13] K. Zeng, J. Yu, R. Wang, C. Li, and D. Tao, "Coupled deep autoencoder for single image super-resolution," *IEEE transactions on cybernetics*, vol. 47, no. 1, pp. 27–37, 2015.
- [14] Z. Chen, K. Yin, M. Fisher, S. Chaudhuri, and H. Zhang, "Bae-net: Branched autoencoder for shape co-segmentation," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 8490–8499.
- [15] G. Dewangan and S. Maurya, "Fault diagnosis of machines using deep convolutional beta-variational autoencoder," *IEEE Transactions on Artificial Intelligence*, 2021.
- [16] X. Dong and J. Shen, "Triplet loss in siamese network for object tracking," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 459–474.
- [17] H. Fan and H. Ling, "Siamese cascaded region proposal networks for real-time visual tracking," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 7952–7961.

- [18] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2223–2232.
- [19] K. K. Babu and S. R. Dubey, "Pcsgan: Perceptual cyclic-synthesized generative adversarial networks for thermal and nir to visible image transformation," *Neurocomputing*, vol. 413, pp. 41–50, 2020.
- [20] L. Bottou, "Large-scale machine learning with stochastic gradient descent," in *Proceedings of the COMPSTAT*, 2010, pp. 177–186.
- [21] H. Robbins and S. Monro, "A stochastic approximation method," *The annals of mathematical statistics*, pp. 400–407, 1951.
- [22] B. 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.
- [23] I. Sutskever, J. Martens, G. Dahl, and G. Hinton, "On the importance of initialization and momentum in deep learning," in *Proceedings of the International Conference on Machine Learning*, 2013, pp. 1139–1147.
- [24] J. Duchi, E. Hazan, and Y. Singer, "Adaptive subgradient methods for online learning and stochastic optimization," *Journal of Machine Learning Research*, vol. 12, no. Jul, pp. 2121–2159, 2011.
- [25] M. D. Zeiler, "Adadelata: an adaptive learning rate method," *arXiv preprint arXiv:1212.5701*, 2012.
- [26] G. Hinton, N. Srivastava, and K. Swersky, "Neural networks for machine learning," Lecture 6a overview of mini-batch gradient descent course, 2012.
- [27] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *Proceedings of the 3rd International Conference on Learning Representations*, 2015.
- [28] S. R. Dubey, S. Chakraborty, S. K. Roy, S. Mukherjee, S. K. Singh, and B. B. Chaudhuri, "Diffgrad: an optimization method for convolutional neural networks," *IEEE Transactions on Neural Networks and Learning Systems*, 2019.
- [29] L. Liu, H. Jiang, P. He, W. Chen, X. Liu, J. Gao, and J. Han, "On the variance of the adaptive learning rate and beyond," in *Proceedings of the International Conference on Learning Representations*, 2019.
- [30] J. Zhuang, T. Tang, S. Tatikonda, N. Dvornik, Y. Ding, X. Papademetris, and J. S. Duncan, "Adabelief optimizer: Adapting stepsizes by the belief in observed gradients," in *Proceedings of the Conference on Neural Information Processing Systems*, 2020.
- [31] W. An, H. Wang, Q. Sun, J. Xu, Q. Dai, and L. Zhang, "A pid controller approach for stochastic optimization of deep networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 8522–8531.
- [32] T. Dozat, "Incorporating nesterov momentum into adam," in *Proceedings of the International Conference on Learning Representations Workshops*, 2016.
- [33] H. Huang, C. Wang, and B. Dong, "Nostalgic adam: weighting more of the past gradients when designing the adaptive learning rate," in *Proceedings of the 28th International Joint Conference on Artificial Intelligence*. AAAI Press, 2019, pp. 2556–2562.
- [34] M. Zaheer, S. Reddi, D. Sachan, S. Kale, and S. Kumar, "Adaptive methods for nonconvex optimization," in *Advances in neural information processing systems*, 2018, pp. 9793–9803.
- [35] L. Luo, Y. Xiong, Y. Liu, and X. Sun, "Adaptive gradient methods with dynamic bound of learning rate," in *International Conference on Learning Representations*, 2018.
- [36] J. Ding, X. Ren, R. Luo, and X. Sun, "An adaptive and momental bound method for stochastic learning," *arXiv preprint arXiv:1910.12249*, 2019.
- [37] J. Lucas, S. Sun, R. Zemel, and R. Grosse, "Aggregated momentum: Stability through passive damping," in *Proceedings of the International Conference on Learning Representations*, 2019.
- [38] Y. You, J. Li, S. Reddi, J. Hseu, S. Kumar, S. Bhojanapalli, X. Song, J. Demmel, K. Keutzer, and C.-J. Hsieh, "Large batch optimization for deep learning: Training bert in 76 minutes," in *Proceedings of the International Conference on Learning Representations*, 2019.
- [39] B. Heo, S. Chun, S. J. Oh, D. Han, S. Yun, Y. Uh, and J.-W. Ha, "Slowing down the weight norm increase in momentum-based optimizers," *arXiv preprint arXiv:2006.08217*, 2020.
- [40] H. Yong, J. Huang, X. Hua, and L. Zhang, "Gradient centralization: A new optimization technique for deep neural networks," in *Proceedings of the European Conference on Computer Vision*, 2020.
- [41] Z. Yao, A. Gholami, S. Shen, M. Mustafa, K. Keutzer, and M. Mahoney, "Adahessian: An adaptive second order optimizer for machine learning," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 12, 2021, pp. 10665–10673.
- [42] S. Roy, M. Paoletti, J. Haut, S. Dubey, P. Kar, A. Plaza, and B. Chaudhuri, "Angulargrad: A new optimization technique for angular convergence of convolutional neural networks," *arXiv preprint arXiv:2105.10190*, 2021.
- [43] M. Toussaint, "Lecture notes: Some notes on gradient descent," 2012. [Online]. Available: <https://ipvs.informatik.uni-stuttgart.de/mlr/marc/notes/gradientDescent.pdf>
- [44] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [45] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7132–7141.
- [46] A. Krizhevsky, "Learning multiple layers of features from tiny images," *Master's thesis, University of Tront*, 2009.
- [47] H. Xiao, K. Rasul, and R. Vollgraf, (2017) Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms.
- [48] Y. Le and X. Yang, "Tiny imagenet visual recognition challenge," *CS 231N*, vol. 7, p. 7, 2015.
- [49] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, "Imagenet large scale visual recognition challenge," *International Journal of Computer Vision*, vol. 115, no. 3, pp. 211–252, 2015.



**Shiv Ram Dubey** is with the Indian Institute of Information Technology (IIT), Allahabad since July 2021, where he is currently the Assistant Professor of Information Technology. He was with IIT Sri City as Assistant Professor from Dec 2016 to July 2021 and Research Scientist from June 2016 to Dec 2016. He received the PhD degree from IIT Allahabad in 2016. Before that, from 2012 to 2013, he was a Project Officer at Indian Institute of Technology (IIT), Madras. He was a recipient of several awards including Best PhD Award in PhD Symposium at IEEE-CICT2017, Early Career Research Award from SERB, Govt. of India and NVIDIA GPU Grant Award Twice from NVIDIA. His research interest includes Computer Vision and Deep Learning.



**S.H. Shabbeer Basha** is associated with PathPartner Technology Pvt. Ltd., Bangalore, as a lead computer vision engineer. At PathPartner, he is involved in R&D activities on neural network compression and deep learning. He received the PhD degree from IIT Sri City. His research interests include Computer Vision, Deep Learning, Deep Model Compression, Unsupervised Domain Adaptation, Transfer Learning, and Multi-Task Learning.



**Satish Kumar Singh** is with the Indian Institute of Information Technology Allahabad, as an Associate Professor from 2013 and heading the Computer Vision and Biometrics Lab (CVBL). Earlier, he served at Jaypee University of Engineering and Technology Guna, India from 2005 to 2012. His areas of interest include Image Processing, Computer Vision, Biometrics, Deep Learning, and Pattern Recognition. Dr. Singh is proactively offering his volunteer services to IEEE for the last many years in various capacities. He is the senior member of IEEE. Presently Dr. Singh is Section Chair IEEE Uttar Pradesh Section (2021-2022) and a member of IEEE India Council (2021). He also served as the Vice-Chair, Operations, Outreach and Strategic Planning of IEEE India Council (2020) & Vice-Chair IEEE Uttar Pradesh Section (2019 & 2020). Prior to that Dr. Singh was Secretary, IEEE UP Section (2017 & 2018), Treasurer, IEEE UP Section (2016 & 2017), Joint Secretary, IEEE UP Section (2015), Convener Web and Newsletters Committee (2014 & 2015). Dr. Singh is also the technical committee affiliate of IEEE SPS IVMSP and MMSP and presently the Chair of IEEE Signal Processing Society Chapter of Uttar Pradesh Section.



**Bidyut Baran Chaudhuri** received the PhD degree from IIT Kanpur, in 1980. He was a Leverhulme Postdoctoral Fellow with Queen's University, U.K., from 1981 to 1982. He joined the Indian Statistical Institute, in 1978, where he worked as an INAE Distinguished Professor and a J C Bose Fellow at Computer Vision and Pattern Recognition Unit. He is now affiliated to Techno India University, Kolkata as Pro-Vice Chancellor (Academic). His research interests include Pattern Recognition, Image Processing, Computer Vision, and Deep learning, etc.

He pioneered the first workable OCR system for printed Indian scripts Bangla, Assamese and Devnagari. He also developed computerized *Bharati Braille system* with speech synthesizer and has done statistical analysis of Indian language. Prof. Chaudhuri received Leverhulme fellowship award, Sir J. C. Bose Memorial Award, M. N. Saha Memorial Award, Homi Bhabha Fellowship, Dr. Vikram Sarabhai Research Award, C. Achuta Menon Award, Homi Bhabha Award: Applied Sciences, Ram Lal Wadhwa Gold Medal, Jawaharlal Nehru Fellowship, J C Bose fellowship, Om Prakash Bhasin Award, etc. Prof. Chaudhuri is the fellow of INSA, NASI, INAE, IAPR, The World Academy of Sciences (TWAS) and life fellow of IEEE (2015).

## SUPPLEMENTARY

### A. Convergence Proof

**Theorem 2.** Assume that the gradients for function  $f_t$  (i.e.,  $\|g_{t,\theta}\|_2 \leq G$  and  $\|g_{t,\theta}\|_\infty \leq G_\infty$ ) are bounded for all  $\theta \in \mathbb{R}^d$ . Let also consider that the bounded distance is generated by the proposed optimizer between any  $\theta_t$  (i.e.,  $\|\theta_n - \theta_m\|_2 \leq D$  and  $\|\theta_n - \theta_m\|_\infty \leq D_\infty$  for any  $m, n \in \{1, \dots, T\}$ ). Let  $\gamma \triangleq \frac{\beta_1^2}{\sqrt{\beta_2}}$ ,  $\beta_1, \beta_2 \in [0, 1)$  satisfy  $\frac{\beta_1^2}{\sqrt{\beta_2}} < 1$ ,  $\alpha_t = \frac{\alpha}{\sqrt{t}}$ , and  $\beta_{1,t} = \beta_1 \lambda^{t-1}$ ,  $\lambda \in (0, 1)$  with  $\lambda$  is around 1, e.g  $1 - 10^{-8}$ . For all  $T \geq 1$ , the proposed injection based AdamInject optimizer shows the following guarantee:

$$\begin{aligned} R(T) &\leq \frac{D^2}{\alpha(1-\beta_1)} \sum_{i=1}^d \sqrt{T \hat{v}_{T,i}} \\ &+ \frac{2\alpha(1+\beta_1)G_\infty}{(1-\beta_1)\sqrt{1-\beta_2}(1-\gamma)^2} \sum_{i=1}^d \|g_{1:T,i}\|_2 \\ &+ \sum_{i=1}^d \frac{D_\infty^2 G_\infty \sqrt{1-\beta_2}}{\alpha(1-\beta_1)(1-\lambda)^2} + 4D_\infty G_\infty^2 \sum_{i=1}^d \|g_{1:T,i}\|_2^2 \end{aligned} \quad (13)$$

*Proof:* Following can be written from Lemma 10.2 of Adam [27],

$$f_t(\theta_t) - f_t(\theta^*) \leq g_t^T(\theta_t - \theta^*) = \sum_{i=1}^d g_{t,i}(\theta_{t,i} - \theta_{i,i}^*)$$

Following can be also written by utilizing the update formula of the proposed injection with Adam (i.e., AdamInject) with  $k=2$  and after discarding  $\epsilon$ ,

$$\begin{aligned} \theta_{t+1} &= \theta_t - \frac{\alpha_t \hat{s}_t}{\sqrt{\hat{v}_t}} \\ &= \theta_t - \frac{\alpha_t}{(1-\beta_1^t)} \left( \frac{\beta_{1,t}}{\sqrt{\hat{v}_t}} s_{t-1} + \frac{(1-\beta_{1,t})(g_t + \Delta\theta g_t^2)}{2} \right) \end{aligned} \quad (14)$$

where  $\beta_{1,t}$  is the  $1^{\text{st}}$  order moment coefficient at  $t^{\text{th}}$  iteration. We can write the following w.r.t. the  $i^{\text{th}}$  dimension of parameter vector  $\theta_t \in \mathbb{R}^d$ ,

$$\begin{aligned} (\theta_{t+1,i} - \theta_{i,i}^*)^2 &= (\theta_{t,i} - \theta_{i,i}^*)^2 - \frac{2\alpha_t}{1-\beta_1^t} \left( \frac{\beta_{1,t}}{\sqrt{\hat{v}_{t,i}}} s_{t-1,i} \right. \\ &+ \left. \frac{(1-\beta_{1,t})(g_{t,i} + \Delta\theta g_{t,i}^2)}{2} \right) (\theta_{t,i} - \theta_{i,i}^*) \\ &+ \alpha_t^2 \left( \frac{\hat{s}_{t,i}}{\sqrt{\hat{v}_{t,i}}} \right)^2 \end{aligned} \quad (15)$$

We can reorganize the above equation as,

$$\begin{aligned} g_{t,i}(\theta_{t,i} - \theta_{i,i}^*) &= \frac{(1-\beta_1^t)\sqrt{\hat{v}_{t,i}}}{\alpha_t(1-\beta_{1,t})} \left( (\theta_{t,i} - \theta_{i,i}^*)^2 \right. \\ &- (\theta_{t+1,i} - \theta_{i,i}^*)^2 \left. \right) + \frac{2\beta_{1,t}}{1-\beta_{1,t}} (\theta_{i,i}^* - \theta_{t,i}) s_{t-1,i} \\ &+ \frac{\alpha_t(1-\beta_1^t)}{(1-\beta_{1,t})} \frac{(\hat{s}_{t,i})^2}{\sqrt{\hat{v}_{t,i}}} - \Delta\theta g_{t,i}^2 (\theta_{t,i} - \theta_{i,i}^*). \end{aligned} \quad (16)$$

We can rewrite it as follows:

$$\begin{aligned} g_{t,i}(\theta_{t,i} - \theta_{i,i}^*) &= \frac{(1-\beta_1^t)\sqrt{\hat{v}_{t,i}}}{\alpha_t(1-\beta_{1,t})} \left( (\theta_{t,i} - \theta_{i,i}^*)^2 - (\theta_{t+1,i} - \theta_{i,i}^*)^2 \right) \\ &+ \sqrt{\frac{2\beta_{1,t}}{\alpha_{t-1}(1-\beta_{1,t})}} (\theta_{i,i}^* - \theta_{t,i})^2 \sqrt{\hat{v}_{t-1,i}} \sqrt{\frac{2\beta_{1,t}\alpha_{t-1}(s_{t-1,i})^2}{(1-\beta_{1,t})\sqrt{\hat{v}_{t-1,i}}}} \\ &+ \frac{\alpha_t(1-\beta_1^t)}{(1-\beta_{1,t})} \frac{(\hat{s}_{t,i})^2}{\sqrt{\hat{v}_{t,i}}} - \Delta\theta g_{t,i}^2 (\theta_{t,i} - \theta_{i,i}^*) \end{aligned} \quad (17)$$

Next, we use the Young's inequality,  $ab \leq a^2/2 + b^2/2$  as well as the information that  $\beta_{1,t} \leq \beta_1$ . We also replace  $\Delta\theta$  with  $\theta_{t-1} - \theta_t$ . Thus, we can write the above equation as,

$$\begin{aligned} g_{t,i}(\theta_{t,i} - \theta_{i,i}^*) &\leq \frac{1}{\alpha_t(1-\beta_1)} \left( (\theta_{t,i} - \theta_{i,i}^*)^2 \right. \\ &- (\theta_{t+1,i} - \theta_{i,i}^*)^2 \left. \right) \sqrt{\hat{v}_{t,i}} \\ &+ \frac{\beta_{1,t}}{\alpha_{t-1}(1-\beta_{1,t})} (\theta_{i,i}^* - \theta_{t,i})^2 \sqrt{\hat{v}_{t-1,i}} \\ &+ \frac{\beta_1 \alpha_{t-1} (s_{t-1,i})^2}{(1-\beta_1)\sqrt{\hat{v}_{t-1,i}}} + \frac{\alpha_t}{(1-\beta_1)} \frac{(\hat{s}_{t,i})^2}{\sqrt{\hat{v}_{t,i}}} \\ &- (\theta_{t-1,i} - \theta_{t,i})(\theta_{t,i} - \theta_{i,i}^*) g_{t,i}^2 \end{aligned} \quad (18)$$

In order to compute the regret bound, we aggregate it as per the Lemma 10.4 of Adam [27] across the dimensions for  $i \in \{1, \dots, d\}$  and the convex function sequence for

$t \in \{1, \dots, T\}$  in the upper bound of  $f_t(\theta_t) - f_t(\theta^*)$  as,

$$\begin{aligned}
R(T) &\leq \sum_{i=1}^d \frac{1}{\alpha_1(1-\beta_1)} (\theta_{1,i} - \theta_{,i}^*)^2 \sqrt{\hat{v}_{1,i}} \\
&+ \sum_{i=1}^d \sum_{t=2}^T \frac{1}{(1-\beta_1)} (\theta_{t,i} - \theta_{,i}^*)^2 \left( \frac{\sqrt{\hat{v}_{t,i}}}{\alpha_t} - \frac{\sqrt{\hat{v}_{t-1,i}}}{\alpha_{t-1}} \right) \\
&+ \sum_{i=1}^d \sum_{t=1}^T \frac{\beta_{1,t}}{\alpha_t(1-\beta_{1,t})} (\theta_{,i}^* - \theta_{t,i})^2 \sqrt{\hat{v}_{t,i}} \\
&+ \frac{2\beta_1\alpha G_\infty}{(1-\beta_1)\sqrt{1-\beta_2}(1-\gamma)^2} \sum_{i=1}^d \|g_{1:T,i}\|_2 \\
&+ \frac{2\alpha G_\infty}{(1-\beta_1)\sqrt{1-\beta_2}(1-\gamma)^2} \sum_{i=1}^d \|g_{1:T,i}\|_2 \\
&+ \sum_{i=1}^d \sum_{t=1}^T (\theta_{t,i} - \theta_{t-1,i})(\theta_{t,i} - \theta_{,i}^*) g_{t,i}^2
\end{aligned} \tag{19}$$

It can be further refined with the assumptions that  $\alpha = \alpha_t \sqrt{t}$ ,  $\|\theta_t - \theta^*\|_2 \leq D$ ,  $\|\theta_m - \theta_n\|_\infty \leq D_\infty$  and  $\Delta\theta = (\theta_t - \theta_{t-1})$  is very small. Moreover,  $\Delta\theta \approx 0$  as step size is very small when  $t$  is large. Then, we can approximate  $\Delta\theta$  with an upper bound of  $1/t^2$ . Thus, the above equation can be written as,

$$\begin{aligned}
R(T) &\leq \frac{D^2}{\alpha(1-\beta_1)} \sum_{i=1}^d \sqrt{T\hat{v}_{T,i}} \\
&+ \frac{2\alpha(1+\beta_1)G_\infty}{(1-\beta_1)\sqrt{1-\beta_2}(1-\gamma)^2} \sum_{i=1}^d \|g_{1:T,i}\|_2 \\
&+ \frac{D_\infty^2}{\alpha} \sum_{i=1}^d \sum_{t=1}^t \frac{\beta_{1,t}}{(1-\beta_{1,t})} \sqrt{t\hat{v}_{t,i}} \\
&+ D_\infty \sum_{i=1}^d \sum_{t=1}^T \frac{\sqrt{t}}{t^2} g_{t,i}^2 \\
&\leq \frac{D^2}{\alpha(1-\beta_1)} \sum_{i=1}^d \sqrt{T\hat{v}_{T,i}} \\
&+ \frac{2\alpha(1+\beta_1)G_\infty}{(1-\beta_1)\sqrt{1-\beta_2}(1-\gamma)^2} \sum_{i=1}^d \|g_{1:T,i}\|_2 \\
&+ \frac{D_\infty^2 G_\infty \sqrt{1-\beta_2}}{\alpha} \sum_{i=1}^d \sum_{t=1}^t \frac{\beta_{1,t}}{(1-\beta_{1,t})} \sqrt{t} \\
&+ D_\infty \sum_{i=1}^d \sum_{t=1}^T \frac{g_{t,i}^2}{t}
\end{aligned} \tag{20}$$

As per the finding of Adam [27], i.e.,  $\sum_{t=1}^t \frac{\beta_{1,t}}{(1-\beta_{1,t})} \sqrt{t} \leq$

$\frac{1}{(1-\beta_1)(1-\gamma)^2}$ , the regret bound can be further rewritten as,

$$\begin{aligned}
R(T) &\leq \frac{D^2}{\alpha(1-\beta_1)} \sum_{i=1}^d \sqrt{T\hat{v}_{T,i}} \\
&+ \frac{2\alpha(1+\beta_1)G_\infty}{(1-\beta_1)\sqrt{1-\beta_2}(1-\gamma)^2} \sum_{i=1}^d \|g_{1:T,i}\|_2 \\
&+ \sum_{i=1}^d \frac{D_\infty^2 G_\infty \sqrt{1-\beta_2}}{\alpha(1-\beta_1)(1-\lambda)^2} + D_\infty \sum_{i=1}^d \sum_{t=1}^T \frac{g_{t,i}^2}{t}
\end{aligned} \tag{21}$$

Finally, we utilize Lemma 10.3 of Adam [27] to approximate the upper bound as  $\sum_{t=1}^T \frac{g_{t,i}^2}{t} \leq 4G_\infty^2 \|g_{1:T,i}\|_2^2$ . Thus, the final regret bound is given as,

$$\begin{aligned}
R(T) &\leq \frac{D^2}{\alpha(1-\beta_1)} \sum_{i=1}^d \sqrt{T\hat{v}_{T,i}} \\
&+ \frac{2\alpha(1+\beta_1)G_\infty}{(1-\beta_1)\sqrt{1-\beta_2}(1-\gamma)^2} \sum_{i=1}^d \|g_{1:T,i}\|_2 \\
&+ \sum_{i=1}^d \frac{D_\infty^2 G_\infty \sqrt{1-\beta_2}}{\alpha(1-\beta_1)(1-\lambda)^2} + 4D_\infty G_\infty^2 \sum_{i=1}^d \|g_{1:T,i}\|_2^2
\end{aligned} \tag{22}$$

■

## B. Algorithms

This section provides the Algorithms for different optimization techniques, including diffGrad (Algorithm 3), diffGradInject (Algorithm 4), Radam (Algorithm 5), RadamInject (Algorithm 6), AdaBelief (Algorithm 7) and AdaBeliefInject (Algorithm 8).

---

### Algorithm 3: diffGrad Optimizer

---

**Initialize:**  $\theta_0, m_0 \leftarrow 0, v_0 \leftarrow 0, t \leftarrow 0$

**Hyperparameters:**  $\alpha, \beta_1, \beta_2$

**While**  $\theta_t$  not converged

$t \leftarrow t + 1$

$g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$

$\xi_t \leftarrow 1/(1 + e^{-|g_t - g_{t-1}|})$

$m_t \leftarrow \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t$

$v_t \leftarrow \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2$

**Bias Correction**

$\widehat{m}_t \leftarrow m_t / (1 - \beta_1^t), \widehat{v}_t \leftarrow v_t / (1 - \beta_2^t)$

**Update**

$\theta_t \leftarrow \theta_{t-1} - \alpha \xi_t \widehat{m}_t / (\sqrt{\widehat{v}_t} + \epsilon)$

---

---

**Algorithm 4:** diffGradInject (diffGrad + AdaInject) Optimizer

---

**Initialize:**  $\theta_0, s_0 \leftarrow 0, v_0 \leftarrow 0, t \leftarrow 0$

**Hyperparameters:**  $\alpha, \beta_1, \beta_2, k$

**While**  $\theta_t$  not converged

$t \leftarrow t + 1$

$g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$

$\xi_t \leftarrow 1/(1 + e^{-|g_t - g_{t-1}|})$

**If**  $t = 1$

$s_t \leftarrow \beta_1 \cdot s_{t-1} + (1 - \beta_1) \cdot g_t$

**Else**

$\Delta\theta \leftarrow \theta_{t-2} - \theta_{t-1}$

$s_t \leftarrow \beta_1 \cdot s_{t-1} + (1 - \beta_1) \cdot (g_t + \Delta\theta \cdot g_t^2)/k$

$v_t \leftarrow \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2$

**Bias Correction**

$\widehat{s}_t \leftarrow s_t/(1 - \beta_1^t), \widehat{v}_t \leftarrow v_t/(1 - \beta_2^t)$

**Update**

$\theta_t \leftarrow \theta_{t-1} - \alpha \xi_t \widehat{s}_t / (\sqrt{\widehat{v}_t} + \epsilon)$

---



---

**Algorithm 5:** Radam Optimizer

---

**Initialize:**  $\theta_0, m_0 \leftarrow 0, v_0 \leftarrow 0, t \leftarrow 0$

**Hyperparameters:**  $\alpha, \beta_1, \beta_2$

**While**  $\theta_t$  not converged

$t \leftarrow t + 1$

$\rho_{\infty} \leftarrow 2/(1 - \beta_2) - 1$

$g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$

$m_t \leftarrow \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t$

$v_t \leftarrow \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2$

$\rho_t = \rho_{\infty} - 2t\beta_2^t/(1 - \beta_2^t)$

**If**  $\rho_t \geq 5$

$\rho_u = (\rho_t - 4) \times (\rho_t - 2) \times \rho_{\infty}$

$\rho_d = (\rho_{\infty} - 4) \times (\rho_{\infty} - 2) \times \rho_t$

$\rho = \sqrt{(1 - \beta_2) \times \rho_u / \rho_d}$

$\alpha_1 = \rho \times \alpha / (1 - \beta_1^t)$

**Update**

$\theta_t \leftarrow \theta_{t-1} - \alpha_1 \times m_t / (\sqrt{v_t} + \epsilon)$

**Else**

$\alpha_2 = \alpha / (1 - \beta_1^t)$

**Update**

$\theta_t \leftarrow \theta_{t-1} - \alpha_2 \times m_t$

---



---

**Algorithm 6:** RadamInject (i.e., Radam + Inject) Optimizer

---

**Initialize:**  $\theta_0, s_0 \leftarrow 0, v_0 \leftarrow 0, t \leftarrow 0$

**Hyperparameters:**  $\alpha, \beta_1, \beta_2, k$

**While**  $\theta_t$  not converged

$t \leftarrow t + 1$

$\rho_{\infty} \leftarrow 2/(1 - \beta_2) - 1$

$g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$

**If**  $t = 1$

$s_t \leftarrow \beta_1 \cdot s_{t-1} + (1 - \beta_1) \cdot g_t$

**Else**

$\Delta\theta \leftarrow \theta_{t-1} - \theta_{t-2}$

$s_t \leftarrow \beta_1 \cdot s_{t-1} + (1 - \beta_1) \cdot (g_t - \Delta\theta \cdot g_t^2)/k$

$v_t \leftarrow \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2$

$\rho_t = \rho_{\infty} - 2t\beta_2^t/(1 - \beta_2^t)$

**If**  $\rho_t \geq 5$

$\rho_u = (\rho_t - 4) \times (\rho_t - 2) \times \rho_{\infty}$

$\rho_d = (\rho_{\infty} - 4) \times (\rho_{\infty} - 2) \times \rho_t$

$\rho = \sqrt{(1 - \beta_2) \times \rho_u / \rho_d}$

$\alpha_1 = \rho \times \alpha / (1 - \beta_1^t)$

**Update**

$\theta_t \leftarrow \theta_{t-1} - \alpha_1 \times s_t / (\sqrt{v_t} + \epsilon)$

**Else**

$\alpha_2 = \alpha / (1 - \beta_1^t)$

**Update**

$\theta_t \leftarrow \theta_{t-1} - \alpha_2 \times s_t$

---



---

**Algorithm 7:** AdaBelief Optimizer

---

**Initialize:**  $\theta_0, m_0 \leftarrow 0, v_0 \leftarrow 0, t \leftarrow 0$

**Hyperparameters:**  $\alpha, \beta_1, \beta_2$

**While**  $\theta_t$  not converged

$t \leftarrow t + 1$

$g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$

$m_t \leftarrow \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t$

$v_t \leftarrow \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot (g_t - m_t)^2$

**Bias Correction**

$\widehat{m}_t \leftarrow m_t / (1 - \beta_1^t), \widehat{v}_t \leftarrow v_t / (1 - \beta_2^t)$

**Update**

$\theta_t \leftarrow \theta_{t-1} - \alpha \widehat{m}_t / (\sqrt{\widehat{v}_t} + \epsilon)$

---

---

**Algorithm 8:** AdaBeliefInject (AdaBelief + AdaInject)  
Optimizer

---

**Initialize:**  $\theta_0, s_0 \leftarrow 0, v_0 \leftarrow 0, t \leftarrow 0$

**Hyperparameters:**  $\alpha, \beta_1, \beta_2, k$

**While**  $\theta_t$  not converged

$t \leftarrow t + 1$

$g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$

**If**  $t = 1$

$s_t \leftarrow \beta_1 \cdot s_{t-1} + (1 - \beta_1) \cdot g_t$

**Else**

$\Delta\theta \leftarrow \theta_{t-2} - \theta_{t-1}$

$s_t \leftarrow \beta_1 \cdot s_{t-1} + (1 - \beta_1) \cdot (g_t + \Delta\theta \cdot g_t^2)/k$

$v_t \leftarrow \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot (g_t - s_t)^2$

**Bias Correction**

$\hat{s}_t \leftarrow s_t / (1 - \beta_1^t), \hat{v}_t \leftarrow v_t / (1 - \beta_2^t)$

**Update**

$\theta_t \leftarrow \theta_{t-1} - \alpha \hat{s}_t / (\sqrt{\hat{v}_t} + \epsilon)$

---