

---

# [Re] AdaBelief Optimizer: Adapting Stepsizes by the Belief in Observed Gradients

---

Anonymous Author(s)

Affiliation

Address

email

## Reproducibility Summary

1

### 2 **Scope of Reproducibility**

3 The proposed optimizer: AdaBelief, claims to achieve three goals: fast convergence as in adaptive methods, good  
4 generalization as in SGD, and training stability. We perform experiments to validate the claims of the paper [28].

### 5 **Methodology**

6 To validate these claims, we reproduce experiments on **Image Classification** with CIFAR-10, CIFAR-100 and ImageNet  
7 datasets, **Language Modeling** with Penn Treebank, **Generative Modeling** with WGAN, WGAN-GP and SN-GAN  
8 architectures. We use the code provided by the author<sup>1</sup>. All experiments were performed on 8 NVIDIA V100 GPUs  
9 and took about 1096 GPU hours in total. Our entire code is provided in the supplementary material.

### 10 **Results**

11 The image classification experiments on CIFAR-10, CIFAR-100 and ImageNet are reproduced to within 0.29%, 0.18%  
12 and 0.25% of reported values respectively. The language modeling experiments produce an average deviation of 0.22%,  
13 while the generative modeling experiments on WGAN, WGAN-GP and SN-GAN are replicated to within 2.2%, 1.8%  
14 and 0.33% of reported value.

15 We perform ablation studies for change of dataset in language modeling and for effect of weight decay on ImageNet.  
16 We also perform analysis of generalization ability of optimizers and of training stability of GANs. All of the results  
17 largely support the claims made in the paper [28].

### 18 **What was easy**

19 The authors provide implementation for most of the experiments presented in the paper. Well documented code and  
20 lucid paper helped understand the experiments clearly.

### 21 **What was difficult**

22 The challenging aspects in our study were: (1) Grid search for optimal hyperparameters (HP) in cases where HP were  
23 not provided or results did not match, (2) time and resource intensive experiments like ImageNet (~ 22 hrs.) and  
24 SN-GAN (~ 15 hrs.), (3) writing code to evaluate claims of the AdaBelief paper.

### 25 **Communication with original authors**

26 We communicated the original author Juntang Zhuang on multiple occasions for doubts related to hyperparameters and  
27 code, to which he promptly replied and helped us.

---

<sup>1</sup><https://github.com/juntang-zhuang/Adabelief-Optimizer>

## 28 1 Introduction

29 Optimization is at the heart of machine learning. Training of neural networks aims to find the optimal solution (deepest  
30 valley on the loss surface) using gradient descent. The variation in method to traverse the loss landscape gives rise to  
31 different optimizers. Discovering different optimizers is an active area of research in machine learning. In this report,  
32 we reproduce and add on to the experimental analysis of a recent optimizer, AdaBelief [28], introduced in 2020 at  
33 NeurIPS conference.

34 The proposed optimizer, AdaBelief, claims to outperform its counterparts on various real world deep learning tasks. As  
35 a part of the ML Reproducibility Challenge, we replicate all the experiments mentioned in the AdaBelief paper [28],  
36 comparing it with other optimizers, and also perform additional experiments to investigate the efficacy of AdaBelief.

## 37 2 Details of Optimizers

38 Optimizers are of two types: **(1) accelerated Stochastic Gradient Descent (SGD) family** [24] that includes SGD with  
39 momentum [26] & Nestrov Accelerated Gradient (NAG) [22], and **(2) adaptive methods** like Adam [12], RAdam [13],  
40 AdamW [14], RMSProp [7], Yogi [27], AdaBound [15], AdaBelief [28], MSVAG [2], Fromage [3], Apollo [16].

41 SGD [24] family uses the same learning rate for all parameters, whereas, adaptive methods update their parameters  
42 as a function of gradients. While this has shown success in faster convergence due to a more streamlined trajectory,  
43 it has raised questions regarding the generalization ability of adaptive methods. RMSProp [7] builds over SGD by  
44 penalizing updates in directions that have high gradients. The intuition behind this is to prevent drastic updates in  
45 particular directions. It does so by damping magnitude of update by factor of exponential moving average (EMA)  
46 computed for squares of gradients. Adam [12] improves over RMSProp by introducing a momentum term that helps  
47 prevent over-damping of step size in case of RMSProp. RAdam [13] seeks to tackle the convergence problem of  
48 Adam by proposing to use a small learning rate during initial stages of training when variance is high, while AdamW  
49 [14] and MSVAG [2] address the generalization problem in Adam. AdamW does this by introducing a weight decay  
50 regularization term and MSVAG decomposes Adam as a sign update and magnitude scaling. Yogi [27] considers  
51 the effect of mini-batch size and proposes an update equation that has shown to outperform Adam with very little  
52 hyperparameter tuning. AdaBelief [28] amplifies (or dampens) its updates by a factor proportional to the 'belief' in  
53 observed gradient i.e. square of difference between the observed gradient and EMA of the gradient. AdaBound [15]  
54 bridges the gap between SGD family and Adaptive methods by making use of an update that smoothly transitions from  
55 Adam to SGD. Fromage [3] takes a different path to optimization - it accounts for the network structure by looping  
56 in weight matrices into the update equation. Apollo [16] takes a step forward from the aforementioned first order  
57 optimizers by approximating the Hessian via a diagonal matrix, keeping computations in-line with first-order schemes.

## 58 3 Scope of reproducibility

59 AdaBelief [28] claims to perform better than existing optimizers. To evaluate the validity of its claims, we investigate  
60 the following target questions:

- 61 • Does AdaBelief produce better scores in comparison to other optimizers on real world tasks of image  
62 classification, language modeling, generative modeling and reinforcement learning?
- 63 • Does AdaBelief converge fast like adaptive methods, e.g. Adam?
- 64 • Does AdaBelief generalize well like the accelerated gradient methods, e.g. SGD?
- 65 • Adaptive methods like Adam are stable in complex settings like training of Generative Adversarial Networks  
66 (GANs) [6]. How does AdaBelief compare to them?

## 67 4 Methodology

### 68 4.1 Experimental setup and model description

69 We perform experiments on many real world tasks: **(a) Image Classification:** CIFAR-10, CIFAR-100 & ImageNet  
70 datasets are used. On CIFAR-10 & CIFAR-100, we train using VGG11 [25], ResNet34 [9] and DenseNet121 [11]. In  
71 the case of ImageNet we use a ResNet18 architecture. **(b) Language Modeling:** Penn Treebank [17] and WikiText-2  
72 [18] datasets are used. Both are used to train 1, 2, 3-layer LSTM [10]. The HP of the LSTM model were taken

73 from here<sup>2</sup>. **(c) Generative Modeling:** CIFAR-10 dataset is used with Wasserstein-GAN (WGAN) [1], with the  
74 improved gradient penalty version WGAN-GP) [8] & with spectral normalization GAN (SN-GAN) [20] architectures,  
75 where generator and discriminator use same HP. WGAN is a smaller model with a vanilla CNN generator, whereas  
76 the SN-GAN is a bigger model with spectral normalization in the discriminator. For SN-GAN we make use of this  
77 repository<sup>3</sup> **(d) Reinforcement Learning:** An agent is trained by Adam and AdaBelief optimizers to play Space  
78 Invaders (Atari Game) using Deep Q-Network (DQN) [21] architecture. Implementation was taken from here<sup>4</sup>. The  
79 code for experiments on image classification, language modeling, WGAN, WGAN-GP was taken from here<sup>5</sup>.

S. No.	Task	Dataset	Setup	Rep. Status	Our Contribution	No. of Exp.	GPU HPR	Total GPU hours
1.	Image Classification	CIFAR-10	VGG, RN, DN	✓	Exp. on Apollo; bias-variance anal.	30	2.5	75
2.		CIFAR-100	VGG, RN, DN	✓	Exp. on Apollo; bias-variance anal.	30	2.5	75
3.		ImageNet	ResNet18	✓	Analysis of weight decay	3	22	66
4.	Language Modeling	PTB, WT2	LSTM (1 layer)	✓	Fromage LRS; WT2	11	1.33	14.63
5.			LSTM (2 layer)	✓	AdamW & RAdam LRS; WT2	11	2.5	27.5
6.			LSTM (3 layer)	✓	AdamW & RAdam LRS; WT2	11	3.75	41.25
7.	Generative Modeling	CIFAR-10	WGAN [1]	✓	N/A (only reproduced paper's [28] exp.)	70	0.89	53.55
8.			WGAN-GP [8]	✓	N/A (only reproduced paper's [28] exp.)	70	1	66.5
9.			SN-GAN [20]	✓	HP search; training stability anal.	45	15	675
10.	Reinforcement Learning	N/A	Space Invaders (Atari)	✓	Beyond AdaBelief paper [28]	2	1	2

Table 1: Summary of our contributions and reproducibility details of performed experiments. Exp. 1 to 9 are mentioned in the AdaBelief paper [28] and have been reproduced successfully along with some additional contribution to each experiment. We also perform exp. 10 which is not a part of AdaBelief paper. **[Legend - Rep.:** Reproducibility, **Exp.:** "Experiment(s)", **HPR:** "hours per run", **RPO:** "runs per optimizer", **anal.:** "analysis", **HP:** "hyperparameter", **LRS:** "Learning Rate Search", **WT2:** "WikiText-2", **DN:** "DenseNet121", **RN:** "ResNet34", **VGG:** "VGG11", **PTB:** "Penn Treebank"]

## 80 4.2 Datasets

81 The following datasets were used in the experiments - **(a) CIFAR-10:** It consists of 60,000 images of size  $32 \times 32$ ,  
82 grouped into 10 classes (6000 images per class). We use the default train-test split of 50,000 : 10,000. **(b) CIFAR-100:**  
83 It is same as CIFAR-10 but the images are grouped into 100 classes (600 images per class). **(c) ImageNet** [5]: We use  
84 ILSVRC 2012 dataset<sup>6</sup> which consists of  $\sim 1.35M$  images of size  $256 \times 256$  split into 1000 classes. Train-val-test  
85 split is 1,281,167 : 50,000 : 100,000. As part of pre-processing we remove mis-labelled data<sup>7</sup> **(d) Penn Treebank**<sup>8</sup>  
86 (PTB) [17]: The train-val-test split of tokens is 887,521 : 70,390 : 78,669. **(e) WikiText-2** (WT2) [18]: It is a subset  
87 of WikiText-103, features a larger vocabulary and retains the punctuation, original case and numbers which are omitted  
88 in PTB dataset. We ran experiments on WT2<sup>9</sup> using the train-val-test token split of 2,045,059 : 213,119 : 240,498.

## 89 4.3 Hyperparameters

90 In this section we mention the HP used by optimizers in our experiments. Optimal values of commonly used HP are  
91 listed in Table 2. Below we mention the source of these values and details of HP search.

92 For most experiments, we use the optimizer-specific HP as mentioned in the original repository<sup>5</sup> since searching the HP  
93 for all experiments is computationally infeasible. However, the repository does not mention the HP for SN-GAN &  
94 Fromage, and the mentioned HP for 2- and 3-layer AdamW & RAdam resulted in large deviation. So, we perform  
95 learning rate (LR) search for **Fromage** and 2- & 3-layer **AdamW** and **RAdam** over the interval  $[10^{-3}, 10^{-2}]$  (5 values).  
96 For **SN-GAN**, we search  $\beta_1$  (3 values in  $[0.4, 0.9]$ ) and  $\epsilon$  (3 values in  $[10^{-12}, 10^{-6}]$ ). For **Reinforcement Learning**,  
97 we use LR of  $10^{-4}$  and  $\epsilon = 10^{-10}$  for AdaBelief and Adam, as mentioned on the RL repository<sup>4</sup>.

<sup>2</sup><https://github.com/salesforce/awd-lstm-lm>

<sup>3</sup><https://github.com/juntang-zhuang/SNGAN-AdaBelief>

<sup>4</sup><https://github.com/juntang-zhuang/rainbow-adabelief>

<sup>5</sup><https://github.com/juntang-zhuang/Adabelief-Optimizer>

<sup>6</sup>ImageNet dataset (Kaggle)

<sup>7</sup>Blacklisted images (GitHub)

<sup>8</sup>Penn Treebank Dataset

<sup>9</sup>WikiText-2 dataset

98 Now we list the HP which are specific to each optimizer. The LR decays to  $1/10^{th}$  of its value  
99 at  $150^{th}$  epoch for image classification on CIFAR-10 and CIFAR-100, and at epoch 70 & 80 on Ima-  
100 geNet. **AdaBelief** uses `weight_decouple=False`, `fixed_decay=False`, `rectify=False` for all the experi-  
101 ments and `weight_decouple=True` on ImageNet. **SGD** uses `momentum=0.9`, and **Apollo** uses `warmup=200`,  
102 `weight_decay_type='L2'` for image classification on CIFAR-10 and CIFAR-100. **AdaBound** uses `final_lr=30`  
103 on PTB and `final_lr=0.01` with GAN experiments.

Task	Setup	Learning Rate	$\beta_1$	$\beta_2$	$\epsilon$	Weight Decay	Epochs
Image Classification	CIFAR	$10^{-3}$ ( $10_{S,M}^{-1}, 1_L$ )	0.9	0.999	$10^{-8}$ ( $10_Y^{-3}, 10_L^{-4}$ )	$5 \times 10^{-4}$ ( $10_W^{-2}, 2.5 \times 10_L^{-4}$ )	200
	ImageNet	$10^{-3}$	0.9	0.999	$10^{-8}$	$10^{-2}$	90
Language Modeling	1 layer	$10^{-3}$ ( $30_{S,M}, 10_{Y,D,F}^{-2}$ )	0.9	0.999	$10^{-8}$ ( $10_B^{-16}, 10_Y^{-3}$ )	$1.2 \times 10^{-6}$	200
	2 layer	$10^{-2}$ ( $30_{S,M}, 10_{W,R}^{-3}$ )	0.9	0.999	$10^{-8}$ ( $10_B^{-12}, 10_Y^{-3}$ )	$1.2 \times 10^{-6}$	200
	3 layer	$10^{-2}$ ( $30_{S,M}, 10_{W,R}^{-3}$ )	0.9	0.999	$10^{-8}$ ( $10_B^{-12}, 10_Y^{-3}$ )	$1.2 \times 10^{-6}$	200
Generative Modeling	WGAN	$2 \times 10^{-4}$	0.5	0.999	$10^{-8}$ ( $10_B^{-12}$ )	0 ( $5 \times 10_P^{-4}$ )	100
	WGAN-GP	$2 \times 10^{-4}$	0.5	0.999	$10^{-8}$ ( $10_B^{-12}$ )	0 ( $5 \times 10_P^{-4}$ )	100
	SN-GAN	$2 \times 10^{-4}$	0.5	0.999	$10^{-8}$ ( $10_A^{-6}, 10_B^{-12}$ )	0	100000

Table 2: Optimizer specific hyperparameter (HP) values and epochs for experiments performed. Each cell follows a format  $X(Y)$  where  $X$  is the optimal value of the HP unless stated otherwise and  $Y$  contains elements of the form  $v_o$  where  $v$  is the value of HP for optimizer  $o$ . The abbreviations used for optimizers are (S)GD, (A)dam, Adam(W), Ada(B)elief, (Y)ogi, (M)SVAG, (R)Adam, (F)romage, AdaBoun(D), Apo(L)lo, (P)adam

## 104 4.4 Computational requirements

105 We run experiments on a Portable Batch System (PBS) managed cluster. We used 8 NVIDIA V100 GPUs and 384 GB  
106 RAM. All experiments except ImageNet use a single GPU. GPU runtime of all experiments are listed in table 1.

## 107 5 Experiments and Results

### 108 5.1 Experiments reproducing original paper

109 To evaluate the performance of AdaBelief and to validate the aforesaid claims, we perform experiments on various  
110 tasks like Image Classification, Language Modeling, Generative Modeling, Reinforcement Learning and compare our  
111 results with those stated in the paper [28]. HP details can be found in Table 2

#### 112 5.1.1 Image classification

113 We run experiments on CIFAR-10 and CIFAR-100 using VGG11 [25], Resnet34 [9] and DenseNet121 [11] architectures,  
114 performing 3 independent runs on 9 optimizers<sup>10</sup>. Additionally, we perform experiments using Apollo optimizer [16],  
115 that has claimed to outperform AdaBelief on CIFAR datasets with ResNet110 architecture. Fig. 1 plots test accuracy  
116 results. Plots for train accuracies are in supplementary (Fig. 5). All the obtained results agree with those reported in the  
117 AdaBelief paper [28].

118 To assess the performance on large scale datasets, we ran experiments on ImageNet [5]. We follow a similar setting as  
119 the author and run experiments on AdaBelief [28] and MSVAG [2] and report results for remaining optimizers from  
120 literature (Table 3). The top-1 accuracy lags by 0.32% and 0.18% respectively in case of AdaBelief and MSVAG. Other  
121 optimizers from literature use weight decay of  $10^{-4}$  while the author performs experiments on AdaBelief using a value  
122 of  $10^{-2}$ . We analyse the effect of weight decay in section 6.2.

Adabelief	SGD	Adabound	Yogi	Adam	MSVAG	RAdam	AdamW
69.76	<b>70.23<sup>†</sup></b>	68.13 <sup>†</sup>	68.23 <sup>†</sup>	63.79 <sup>†</sup> (66.54 <sup>‡</sup> )	65.81	67.62 <sup>‡</sup>	67.93 <sup>†</sup>

Table 3: Top-1 accuracy of ResNet18 on ImageNet. <sup>†</sup> is reported in [4], and <sup>‡</sup> is reported in [13]

<sup>10</sup>SGD, Adam, AdamW, AdaBelief, Yogi, MSVAG, RAdam, Fromage, AdaBound

### 123 5.1.2 Language Modeling

124 We ran experiments on Penn Treebank (PTB) dataset [17] using 1,2,3-layer LSTM models. We report test perplexities  
 125 (ppl) (Fig. 2) for 3 independent runs on 9 optimizers<sup>10</sup>. Plots for train ppl are in supplementary (Fig. 1). For Fromage,  
 126 the author does not provide HP, hence we grid search and find the optimal  $LR = 10^{-2}$ . In case of 2 layer LSTM using  
 127 AdamW & RAdam, we find that an  $LR = 10^{-3}$  gives a ppl of 73.78 & 74.05, while  $LR = 10^{-2}$  gives a ppl of 93.61 &  
 128 90.49 respectively. The author reports a ppl  $\sim 73$ ,  $\sim 73.5$  at  $LR = 10^{-2}$ . Similarly, in 3-layer LSTM,  $LR = 10^{-3}$   
 129 for AdamW and RAdam works better than  $LR = 10^{-2}$ . PTB is a small dataset, so, we additionally experiment on  
 130 WikiText-2 (section 6.1) for Adam and AdaBelief (top performers in case of PTB) on the setting reported here<sup>11</sup>.

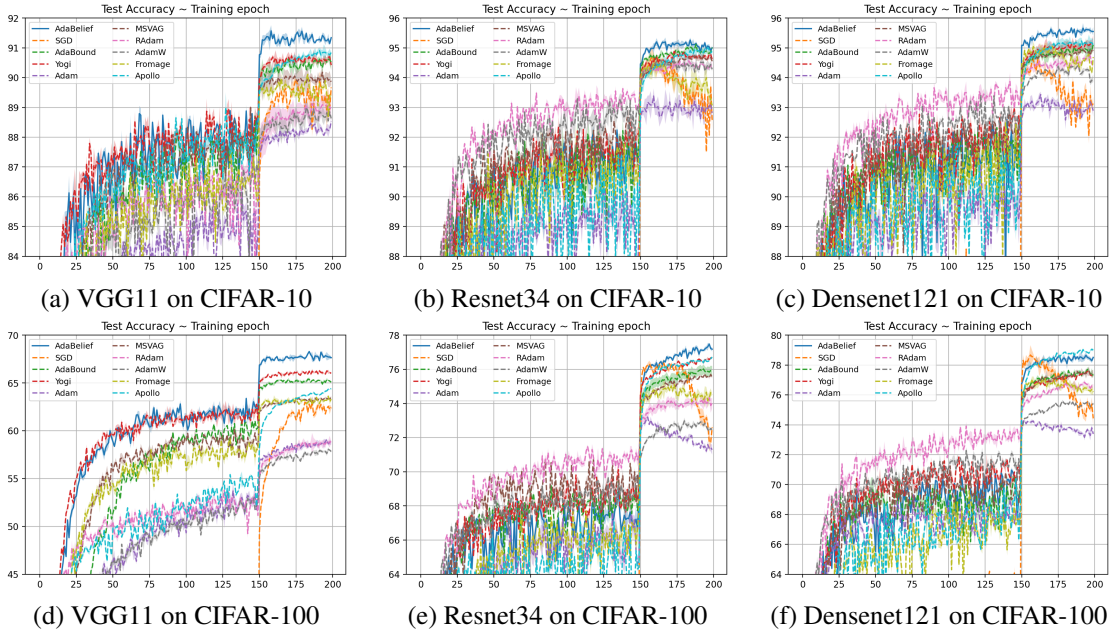


Figure 1: Test accuracy ( $[\mu \pm \sigma]$ ) on CIFAR-10 and CIFAR-100

Adabelief	RAdam	RMSProp	Adam	Fromage	Yogi	SGD	MSVAG	AdaBound
12.98 $\pm$ 0.22	13.10 $\pm$ 0.20	<b>12.86 <math>\pm</math> 0.08</b>	13.01 $\pm$ 0.15	46.31 $\pm$ 0.86	14.16 $\pm$ 0.05	48.94 $\pm$ 2.88	56.89 $\pm$ 2.61	16.84 $\pm$ 0.10

Table 4: FID ( $[\mu \pm \sigma]$ ) of a SN-GAN with ResNet generator on CIFAR-10.

### 131 5.1.3 Generative Modeling

132 We run experiments on WGAN [1], WGAN-GP [8] & SN-GAN [20]. SN-GAN makes use of a ResNet generator with  
 133 spectral normalization in the discriminator and is trained for 100,000 steps. Five independent runs on 9 optimizers<sup>12</sup> are  
 134 performed. We also perform these experiments using the Padam [19] optimizer on WGAN and WGAN-GP. FID values  
 135 for SN-GAN and Padam (Table 4, 5). Fig. 4 shows the variation in FID during training, giving an idea of stability and  
 136 convergence of different optimizers. Boxplots of FID values corresponding to multiple runs on WGAN and WGAN-GP  
 137 are shown in Fig. 3. Collages of generated images for all optimizers can be found in supplementary (Fig. 7, 8, 9).

138 **(a) SN-GAN:** In case of Fromage [3] and MSVAG [2], we obtain  $\sim 4$  and  $\sim 8$  worse FID than what is reported, while  
 139 for AdaBound [15] we obtain a  $\sim 40$  better FID. We suspect the reason for this large deviation to be a difference in  
 140 HP value being used. Since we performed a HP search for SN-GAN, our HP (Table 2) are optimal. The results of  
 141 remaining optimizers were comparable to what was reported in the paper. **(b) WGAN:** We observe that AdaBelief  
 142 outperforms other optimizers with a median FID of  $\sim 80$  which agrees with reported value. We observe a significantly  
 143 worse FID with Fromage. **(c) WGAN-GP:** AdaBelief and AdaBound achieve comparable results  $\sim 67$  FID which are

<sup>11</sup> <https://github.com/salesforce/awd-lstm-lm>

<sup>12</sup> SGD, Adam, RMSProp, AdaBelief, Yogi, MSVAG, RAdam, Fromage, AdaBound

144 better than the other optimizers. Fromage shows similar deviation like in WGAN. With Padam, we find that for both  
 145 WGAN and WGAN-GP, increasing the partial ( $p$ ) i.e. moving from SGD towards Adam, decreases the FID. The FIDs  
 146 obtained are found to agree with or are marginally better than what was stated in the paper.

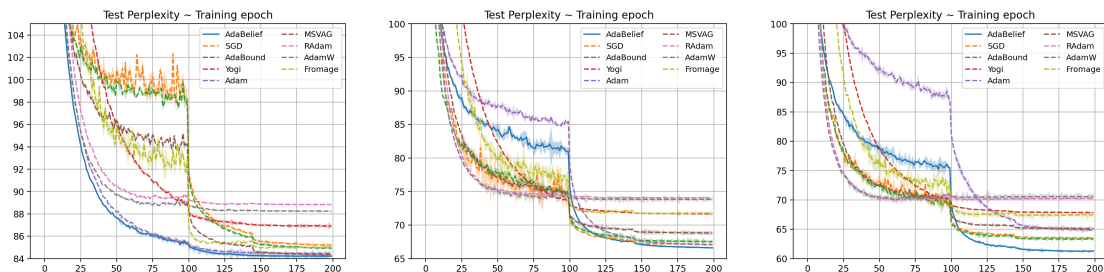
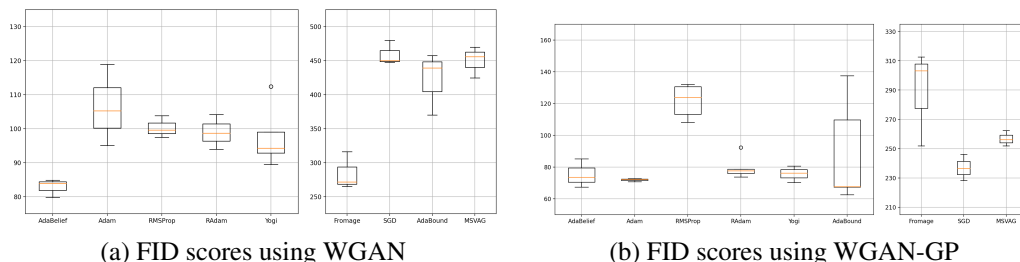


Figure 2: Left to right: Test perplexity ( $[\mu \pm \sigma]$ ) on Penn Treebank for 1,2,3-layer LSTM

	AdaBelief	Padam						
		p=1/2 (Adam)	p=2/5	p=1/4	p=1/5	p=1/8	p=1/16	p=0 (SGD)
FID (WGAN)	<b>82.85 ± 2.21</b>	106.38 ± 9.76	95.66 ± 3.76	422.62 ± 35.68	396.69 ± 24.91	330.44 ± 26.62	357.26 ± 32.39	459.01 ± 14.62
FID (WGAN-GP)	75.37 ± 7.37	<b>71.87 ± 0.83</b>	85.42 ± 5.15	152.34 ± 17.49	170.80 ± 20.43	205.57 ± 13.79	228.40 ± 18.24	236.99 ± 7.26

Table 5: FID values ( $[\mu \pm \sigma]$ ) using AdaBelief and Padam on WGAN and WGAN-GP, Lower FID is better.



(a) FID scores using WGAN

(b) FID scores using WGAN-GP

Figure 3: FID score of WGAN and WGAN-GP using a vanilla CNN generator on CIFAR-10. Lower is better. For each model, successful and failed optimizers (i.e. ones with higher FID values) are shown in the left and right respectively, with different y-axis ranges.

## 147 5.2 Experiments beyond original paper

### 148 5.2.1 RL toy

149 To investigate the efficacy of AdaBelief in use cases beyond text and images we train an agent to play Space Invaders  
 150 (Atari Game). We report  $Q$  value and reward function for Adam and AdaBelief in supplementary (Fig. 10, 11). We  
 151 compare our results with author’s results from here<sup>13</sup> and find that both results agree.

### 152 5.2.2 Image Classification on CIFAR-10 and CIFAR-100 using Apollo

153 Apollo [16] is another optimizer that claims to achieve better convergence speed and generalization than SGD and  
 154 variants of Adam. To investigate this, we experiment with Apollo on CIFAR-10 and CIFAR-100. Fig. 5, 6 (in  
 155 supplementary) show the train, test accuracies on VGG11, ResNet34 and DenseNet121 for the 3 independent runs.  
 156 AdaBelief outperforms Apollo in all settings except DenseNet121 on CIFAR-100. It can also be seen that as we move  
 157 from a simpler (VGG11) to a complex architecture (DenseNet121) the gap between Apollo and AdaBelief closes out.  
 158 We made use of official implementation of Apollo in our experiments<sup>14</sup>.

<sup>13</sup><https://github.com/juntang-zhuang/rainbow-adabelief>

<sup>14</sup><https://github.com/XuezheMax/apollo>

159 **5.2.3 Evaluating GAN training stability**

160 To assess stability of AdaBelief while training GANs, we look into difference between SN-GAN’s generator and  
 161 discriminator training losses on CIFAR-10. We do this for AdaBelief, Adam and RMSProp (since they have top-2  
 162 FID scores on SN-GAN) in the adaptive family, and with SGD for a comparison. Fig. 12 in supplementary plots the  
 163 generator and discriminator training losses. We observe that the adaptive methods are more stable than SGD and within  
 164 the adaptive family the order of stability from most stable to least stable varies as RMSProp, AdaBelief, Adam.

165 **5.2.4 Evaluating generalization ability**

166 To evaluate AdaBelief’s ability to generalize, we analyze the bias and variance of image classification models trained  
 167 using SGD, Adam, AdaBelief and Apollo optimizers on CIFAR-10 and CIFAR-100. We use the method outlined here  
 168 [23] for bias-variance analysis. For each optimizer, we note its train and test accuracy (Fig. 1) corresponding to the  
 169 epoch with best test accuracy (acc), and compute their difference. This data is stated as 3-tuples in Table 6. Lower  
 170 training acc denotes high bias and vice-versa. The difference between the train and test acc is a measure of variance.  
 171 Based on Table 6, we observe that AdaBelief models have the least bias on all configurations, while they have 2<sup>nd</sup>, 3<sup>rd</sup>  
 172 or 4<sup>th</sup> lowest variance. SGD has the least variance on most configurations (highlighted in red), but their bias is high  
 173 (mostly ranked 3<sup>rd</sup> or 4<sup>th</sup> in low bias).

Optimizer	CIFAR-10			CIFAR-100		
	VGG11	ResNet34	DenseNet121	VGG11	ResNet34	DenseNet121
SGD	95.88, 89.95, <b>5.93</b>	98.77, 94.72, <b>4.05</b>	98.72, 94.61, <b>4.11</b>	78.87, 63.09, 15.78	98.94, 76.35, 22.59	94.67, 78.67, <b>16.00</b>
Adam	94.68, 88.54, 6.14	98.36, 93.38, 4.98	99.23, 93.43, 5.80	67.63, 59.08, <b>8.55</b>	92.73, 73.20, <b>19.53</b>	96.69, 74.28, 22.41
AdaBelief	<b>99.36</b> , 91.57, 7.79	<b>99.96</b> , 95.26, 4.70	<b>99.97</b> , 95.67, 4.30	<b>98.84</b> , 68.29, 30.55	<b>99.97</b> , 77.48, 22.49	<b>99.96</b> , 78.66, 21.30
Apollo	98.79, 90.91, 7.88	99.74, 95.01, 4.73	99.82, 95.23, 4.59	74.80, 64.42, 10.38	99.54, 76.72, 22.82	99.68, 79.06, 20.62

Table 6: Analysis of generalization capability of AdaBelief on CIFAR-10 and CIFAR-100 for VGG11, ResNet34 and DenseNet121 architectures using **bias** and **variance**. Each cell denotes a 3-tuple of the form (train acc, test acc, difference b/w train and test acc) corresponding to the model which achieves best test acc (out of 3 runs) for each configuration. For each column, the value in **red** denotes the optimizer with least **variance** (i.e. the least train-test acc difference) and the value in **blue** denotes the optimizer with least **bias** (i.e. with most training acc). AdaBelief models achieve the least bias on all configurations, while they lag behind in terms of variance.

174 **5.2.5 Evaluating convergence speed**

175 **Definition 5.1** (Epoch of Convergence (EC)). Let  $m_k$  denote the metric (acc or ppl) at  $k^{th}$  epoch.  $EC$  is then defined as  
 176 the smallest epoch  $x$  such that  $|m_y - m_x| < \delta \forall y \in [x, x + w]$ , where  $w$  and  $\delta$  are chosen as 15 and 0.05 respectively.  
 177 In other words,  $EC$  is the smallest epoch for which there exists at least  $w(= 15)$  epochs to its right with accuracies (or  
 178 perplexities) within a fixed tolerance  $\delta(= 0.05)$ . If such  $x$  cannot be found, the said optimizer is said to have *failed to*  
 179 *converge (FTC)*.

180 To address the claim on convergence ability different optimizers (section 3) we make use of Def. 5.1. We perform  
 181 the analysis for Image Classification and Language Modeling (section 5.1) experiments. We smoothen the accuracy  
 182 (or perplexity) curves for all optimizers by finding the exponential moving average (EMA) with a smoothing factor  
 183  $\beta = 0.7$ . Analyzing the computed  $ECs$  yield that the convergence speed of AdaBelief is comparable to other members  
 184 of Adaptive family for experiments performed on CIFAR datasets (Fig. 1). For Language Modeling experiments, we  
 185 find that Adam and AdaBelief show similar convergence trends but considerably lag behind in comparison to RAdam,  
 186 AdamW and Fromage (Fig. 2) that are unaffected by learning rate decay which takes place at 100<sup>th</sup> epoch. For exact  
 187  $EC$  values refer Table 1 in supplementary.

188 **6 Ablation studies**

189 **6.1 WikiText-2 on LSTM**

190 To study the performance change due to a larger dataset, we ran Language Modeling experiments on WikiText-2 [18]  
 191 using AdaBelief and Adam optimizers with 1, 2, 3 layer LSTM models. Fig. 3, 4 (in supplementary) show train and test  
 192 perplexity for 3 independent runs. It can be seen that the performance of Adam and AdaBelief is comparable on 1 and 2  
 193 layer LSTM models, while in the 3 layer case AdaBelief outperforms Adam by  $\sim 5$  ppl.



194 **6.2 Effect of weight decay on ImageNet**

195 The paper [28] uses a weight decay of  $10^{-2}$  while experimenting with AdaBelief on ImageNet. However, the results for  
 196 other optimizers are from the literature that typically use a (smaller) weight decay of  $10^{-4}$ . To evaluate the effect of  
 197 weight decay, we experiment with AdaBelief using weight decay =  $10^{-4}$  and find  $\sim 2\%$  drop in top-1 accuracy. So, it  
 198 may be interesting to see the effect of weight decay on other optimizers.

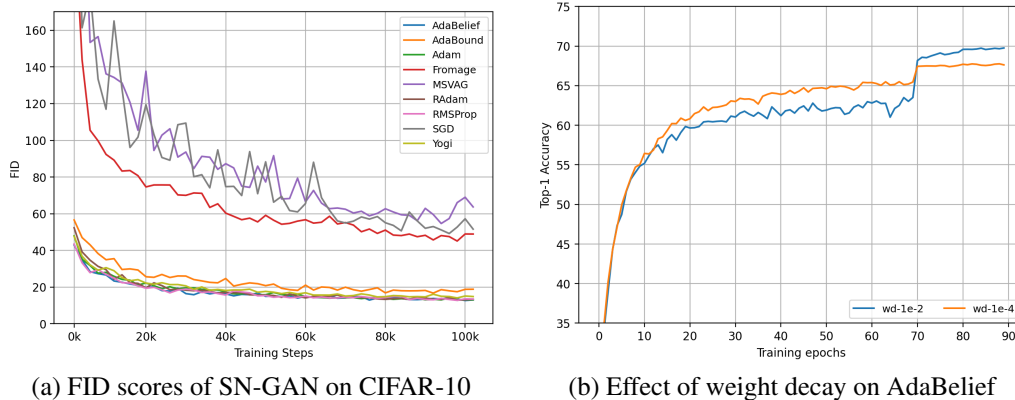


Figure 4: (a) FID values of SN-GAN over training steps for different optimizers (best run plotted out of 5). AdaBelief fares second after RMSProp. (b) AdaBelief performs better when run on larger weight decay of  $10^{-2}$ .

199 **7 Discussion**

200 We now summarize the validity of claims from section 3: (a) Results in section 5.1 show that **AdaBelief outperforms**  
 201 **other optimizers in most use cases**. (b) From section 5.2.5, we find that **the convergence speed of AdaBelief is**  
 202 **largely in line with adaptive methods**. (c) Based on the analysis in section 5.2.4, we infer that AdaBelief generalizes  
 203 well, which is evident by its models having lowest bias and relatively low variance. However, it does not uniformly  
 204 outperform SGD. Therefore, **we fail to completely validate the ability of AdaBelief generalizing as well as SGD**.  
 205 (d) Even though in section 5.2.3, the least difference between generator and discriminator loss is in case of RMSProp,  
 206 AdaBelief does outperform other members of the adaptive family. It defeats SGD by a significant margin. Thus, we  
 207 find that **AdaBelief has stability comparable to adaptive methods in complex settings like GANs**.

208 **What was easy** The authors provide implementation for most of the experiments presented in the paper. Well  
 209 documented code and lucid paper helped understand the experiments clearly.

210 **What was difficult** While hyperparameters (HP) of some experiments were absent (section 5.1.3), some had discrepan-  
 211 cies (section 5.1.2). We had to perform grid search for these cases. Training SN-GAN and ImageNet was a resource  
 212 intensive process which increased the computational burden (Table 1). Formulating the analysis to evaluate the claims  
 213 of the paper was also challenging 5.2.

214 **Communication with original authors** We are thankful to the author Juntang Zhuang. He helped us with the  
 215 implementation and HP details for various experiments. We confirmed the HP for WGAN, SN-GAN, and LSTM  
 216 experiments. We also clarified the source of Penn Treebank dataset and blacklisting of images in ImageNet.

217 **Recommendations for reproducibility** Given the time and resource constraints, we performed only a basic analysis  
 218 of bias-variance trade-off to evaluate the generalization ability of AdaBelief. A more advanced analysis might help in  
 219 revealing the exact weakness of AdaBelief models in terms of ability to generalize.

220 Based on our experiments, ablation studies and analysis, we find that AdaBelief is a promising optimizer combining the  
 221 best of both worlds - accelerated and adaptive gradient methods.

222 **References**

223 [1] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein gan, 2017.



- 224 [2] Lukas Balles and Philipp Hennig. Dissecting adam: The sign, magnitude and variance of stochastic gradients,  
225 2020.
- 226 [3] Jeremy Bernstein, Arash Vahdat, Yisong Yue, and Ming-Yu Liu. On the distance between two neural networks  
227 and the stability of learning. In *Neural Information Processing Systems*, 2020.
- 228 [4] Jinghui Chen, Dongruo Zhou, Yiqi Tang, Ziyang Yang, Yuan Cao, and Quanquan Gu. Closing the Generalization  
229 Gap of Adaptive Gradient Methods in Training Deep Neural Networks. *arXiv:1806.06763 [cs, stat]*, June 2020.  
230 arXiv: 1806.06763.
- 231 [5] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image  
232 database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- 233 [6] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron  
234 Courville, and Yoshua Bengio. Generative adversarial networks, 2014.
- 235 [7] Alex Graves. Generating sequences with recurrent neural networks, 2014.
- 236 [8] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron Courville. Improved training of  
237 wasserstein gans, 2017.
- 238 [9] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *2016*  
239 *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016.
- 240 [10] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural Computation*, 9:1735–1780, 1997.
- 241 [11] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q. Weinberger. Densely connected convolutional  
242 networks. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2261–2269,  
243 2017.
- 244 [12] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017.
- 245 [13] Liyuan Liu, Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and Jiawei Han. On the  
246 Variance of the Adaptive Learning Rate and Beyond. *arXiv:1908.03265 [cs, stat]*, April 2020. arXiv: 1908.03265.
- 247 [14] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019.
- 248 [15] Liangchen Luo, Yuanhao Xiong, Yan Liu, and Xu Sun. Adaptive gradient methods with dynamic bound of  
249 learning rate, 2019.
- 250 [16] Xuezhe Ma. Apollo: An adaptive parameter-wise diagonal quasi-newton method for nonconvex stochastic  
251 optimization, 2021.
- 252 [17] Mitchell P. Marcus, Mary Ann Marcinkiewicz, and Beatrice Santorini. Building a large annotated corpus of  
253 english: The penn treebank. *Comput. Linguist.*, 19(2):313–330, June 1993.
- 254 [18] Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models. *CoRR*,  
255 abs/1609.07843, 2016.
- 256 [19] Harshal Mittal, Kartikey Pandey, and Yash Kant. Iclr reproducibility challenge report (padam : Closing the  
257 generalization gap of adaptive gradient methods in training deep neural networks), 2019.
- 258 [20] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative  
259 adversarial networks. *CoRR*, abs/1802.05957, 2018.
- 260 [21] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and  
261 Martin Riedmiller. Playing atari with deep reinforcement learning, 2013.
- 262 [22] Y. E. NESTEROV. A method for solving the convex programming problem with convergence rate  $o(1/k^2)$ . *Dokl.*  
263 *Akad. Nauk SSSR*, 269:543–547, 1983.
- 264 [23] Andrew Ng. Diagnosing bias vs. variance. *Coursera*, 2017.
- 265 [24] H. Robbins and S. Monro. A stochastic approximation method. *Annals of Mathematical Statistics*, 22:400–407,  
266 1951.

- 267 [25] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition,  
268 2015.
- 269 [26] Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. On the importance of initialization and  
270 momentum in deep learning. In Sanjoy Dasgupta and David McAllester, editors, *Proceedings of the 30th*  
271 *International Conference on Machine Learning*, volume 28 of *Proceedings of Machine Learning Research*, pages  
272 1139–1147, Atlanta, Georgia, USA, 17–19 Jun 2013. PMLR.
- 273 [27] Manzil Zaheer, Sashank J. Reddi, Devendra Sachan, Satyen Kale, and Sanjiv Kumar. Adaptive methods for  
274 nonconvex optimization. In *Proceedings of the 32nd International Conference on Neural Information Processing*  
275 *Systems*, NIPS’ 18, page 9815–9825, Red Hook, NY, USA, 2018. Curran Associates Inc.
- 276 [28] Juntang Zhuang, Tommy Tang, Yifan Ding, Sekhar C. Tatikonda, Nicha Dvornek, Xenophon Papademetris, and  
277 James Duncan. AdaBelief Optimizer: Adapting Stepsizes by the Belief in Observed Gradients. *Advances in*  
278 *Neural Information Processing Systems*, 33:18795–18806, 2020.