

ON THE EFFECTIVENESS OF ADAPTING PRE-TRAINED TRANSFORMER MODELS VIA ADVERSARIAL NOISE

Anonymous authors

Paper under double-blind review

ABSTRACT

Pretraining Transformer-based language models followed by adapting the pre-trained models to a downstream task is an effective transfer mechanism in NLP. While it is well-known that the pretraining stage is computationally expensive, the downstream adaptation also becomes costly as Transformers grow in size rapidly and the wide usage scenarios of fine-tuning pre-trained Transformers. In this work, we find that techniques that have demonstrated success in accelerating the pre-training tasks, such as large-batch optimizations, lead to severe accuracy degradation. We find strong regularization techniques such as adversarial training help to close the accuracy gap. However, the computational complexity associated with this approach, due to the high cost of generating adversaries, prevents it from reducing adaptation costs even with a large number of GPUs. As such, we systematically study both the computation efficiency and generalization of adversarial training for adapting pre-trained transformers, under a large-batch optimization regime. Our investigation yields simple yet effective algorithms for adapting transformer models. We show in experiments that our proposed method attains up to $9.8\times$ adaptation speedups over the baseline on $BERT_{base}$ and $RoBERTa_{large}$, while achieving comparable and sometimes higher accuracy than fine-tuning using existing baselines.

1 INTRODUCTION

In the past few years, we have witnessed the success of transformer models Bommasani et al. (2021), such as BERT Devlin et al. (2019), RoBERTa Liu et al. (2019), T5 Raffel et al. (2019), and GPT-3 Brown et al. (2020). These models are trained on massive open-domain data and subsequently adapted to various downstream tasks, which have led to accuracy breakthroughs in many NLP applications (Wang et al., 2019a). Despite their remarkable performance in accuracy, training these models is extremely time-consuming given their huge model sizes, ranging from a few hundred million parameters to over billions of parameters. As a result, optimizations for faster training speed with high accuracy are the focus of a highly active research area and have a clear, practical impact.

To accelerate the training speed of large models, one of the most popular approaches is to leverage distributed training, where a mini-batch is partitioned across multiple processors (e.g., GPUs) to compute gradients locally in parallel and then aggregate the local updates (Li et al., 2020; Liu et al., 2019; Huang et al., 2019; Shazeer et al., 2018; Shoeybi et al., 2019; Rajbhandari et al., 2019). Under such a paradigm, increasing the batch size has the benefit of improved training throughput per iteration. However, increasing the batch size has a non-trivial impact on model convergence and generation in practice. To close the accuracy gap, prior works proposed to either increase the number of training iterations, which limits the performance benefits of large-batch optimizations Hoffer et al. (2017) or variants of adaptive optimizers such as LAMB You et al. (2019a). It has been empirically observed that LAMB (You et al., 2019a) is able to speed up BERT pre-training by using considerably larger batch sizes on massive GPUs. Despite showing promising results, prior work primarily focuses on large-batch optimizations for accelerating pre-training Transformers. However, as the size of Transformers increases rapidly, reducing the training overhead at the adaption stage starts to become more prominent, e.g., with the active research that has been pushing the training time of BERT models to only a few hours or less than one hour You et al. (2019a); Zheng et al. (2020); ber, it takes tens of hours or even days to fine-tune these models on downstream tasks (Liu et al., 2019). Furthermore, since the adaptation of these large transformer models has been used

by major players in the industry, many model scientists have to perform adaptation more frequently than pre-training the Transformers. As a result, the excessively long adaptation time hinders the turnaround time, and the aggregated training cost for adaptation is also quite high.

We aim to accelerate the adaption of pre-trained Transformer models. For this purpose, we introduce ScaLA, a method that achieves similar model adaptation quality but with significantly shorter optimization time. Especially, the contributions of our paper consist of:

- We analyze projected gradient descent based adversarial training under the large-batch adaptation regime, which shows strong generalization results but leads to a severe computation-vs-generalization dilemma: adversarial large-batch optimization adds significant overhead, making it difficult to actually reduce training time even with a large number of processors.
- We perform systematic studies of how different training strategies affect the computational efficiency and generalization for adapting Transformers. We find that many computations in adversarial training can be avoided with only a minimal impact on the final model accuracy.
- We present ScaLA, a simple yet effective algorithm that injects lightweight adversaries into large batch optimization to speed up the adaptation of pre-trained transformer networks.
- We theoretically quantify the convergence rate of adversarial large-batch optimization using techniques for analyzing non-convex saddle-point problems.
- We conduct extensive evaluation, and our results show that ScaLA accelerates the adaptation of pre-trained Transformer-networks by up to 9.8 times over the baseline on BERT (Devlin et al., 2019), RoBERTa (Liu et al., 2019), and T5 (Raffel et al., 2019) over a wide range of natural language understanding (NLU) tasks. We conduct ablation studies to assess the impact of our approach on both generalization and computational efficiency under various conditions.

2 BACKGROUND AND RELATED WORK

Despite the great success of pre-trained transformer networks such as BERT (Devlin et al., 2019), a big challenge, in general, comes from the training efficiency – even with self-attention and parallelizable recurrence (Vaswani et al., 2017), and high-performance hardware (Jouppi et al., 2017), training transformer networks can still take a significant amount of time. One effective approach to reducing training time is through data parallelism (Devlin et al., 2019; Liu et al., 2019; Shoeybi et al., 2019), which motivates studies on large-batch stochastic non-convex optimizations for transformer networks (You et al., 2019a). These studies have raised concerns with respect to its convergence, generalizability, and training stability by observing that training with a large batch could be difficult (Keskar et al., 2017; Hoffer et al., 2017; Nado et al., 2021). Different from prior works, which mostly focus on reducing the pre-training time (You et al., 2019a; Zhang & He, 2020; Gong et al., 2019; Clark et al., 2020), this work shows an effective approach to accelerate the adaptation of pre-trained models while preserving the accuracy of downstream tasks.

There has also been an increasing interest in developing efficient adaptation methods for pre-trained Transformer models (Houlsby et al., 2019; Wang et al., 2021; Hu et al., 2021; He et al., 2021). For example, (Houlsby et al., 2019) inserts small modules called adapters to each layer of the pre-trained model, and only the adapters are trained during adaptation. (Hu et al., 2021) adds low-rank matrices to approximate parameter updates. (Pfeiffer et al., 2021) shows that it is possible to quickly adapt to new tasks by collectively learning knowledge from multiple tasks. These methods have achieved comparable performance to standard fine-tuning on different sets of tasks. However, their focus is on reducing memory consumption of adaptation by reducing the trainable parameters needed per task, whereas our study focuses on speeding up the adaptation from the perspective of accelerating turnaround time.

Adversarial training has been proposed and studied extensively in the computer vision literature mainly for improving the robustness against adversarial attacks (Goodfellow et al., 2015; Madry et al., 2018). There are some studies show that adversarial training helps improve the generalizability of language modeling Cheng et al. (2019); Wang et al. (2019b); Jiang et al. (2020); Liu et al. (2020). However, very few works examine how adversarial learning works helps improve the adaptation speed of pre-trained Transformer models. The most related work is FreeLb (Zhu et al., 2020), which uses adversarial training as a data augmentation technique to enrich the training set. As such, FreeLb

solely focused on improving the generalization without including an evaluation of its computation efficiency. We provide a comparison with FreeLb in Section 5.

3 CHALLENGES AND OPPORTUNITIES OF ADAPTING PRE-TRAINED TRANSFORMERS WITH LARGE-BATCH OPTIMIZATION

We are interested in extending the large-batch optimization to the Transformation adaptation phase using pre-trained BERT_{base} model on GLUE as an example. This part presents several studies that motivate the design of the lightweight adversarial large-batch optimization approach in Section 4. The detailed hardware/software setup is described in Section 5.

Scalability analysis. First, we carry out a scalability test by varying the number of GPUs from 1 to 32, with and without communication. Different from pre-training, the adaptation stage often employs a much smaller batch size (e.g., 16, 32) than pre-training (e.g., 4096) (Devlin et al., 2019; Liu et al., 2019). We choose a batch size 32, as suggested by most literature for BERT fine-tuning Devlin et al. (2019); Liu et al. (2019), and we divide the samples in the mini-batch among $P=\{1,2,4,8,16,32\}$ GPUs. If the per-worker batch size (e.g., 16) is larger than the maximum admissible per-worker batch size (e.g., 8), we use local gradient accumulation Goyal et al. (2017) to avoid running out of memory. Figure 1(a) shows the scalability results. For batch size 32, the training time decreases when P increases from 1 to 4. However, it quickly plateaus and even decreases with more GPUs. We find that this is because when the batch size is small, the communication overhead dominates the total execution time (e.g., B=32 vs. B=32 (no comm)). The communication overhead is huge, especially when there is cross-machine communication (e.g., from 16 to 32), hindering the scalability of multi-GPU training. In contrast, by increasing the batch size (e.g., to 1K), the training time keeps decreasing as the number of GPUs increases because an increased batch size reduces the number of all-reduce communications to process the same amount of data and also increases the compute resource utilization per GPU (i.e., increased computation-vs-communication ratio).

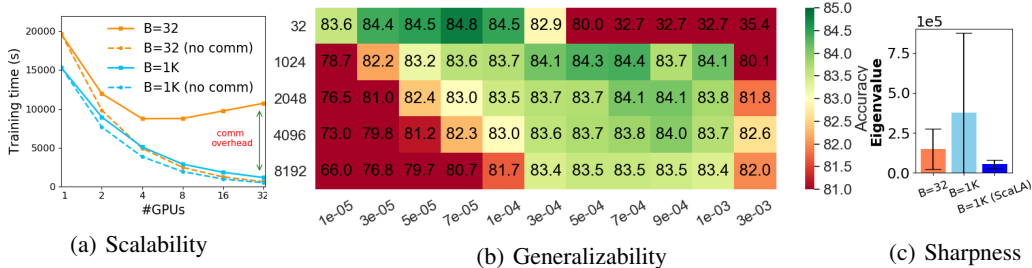


Figure 1: Scalability, generalizability, and curvature analysis results by adapting BERT_{base} to the MNLI task.

Generalizability analysis. Increasing the batch size leads to accelerated per-epoch execution time due to the efficient utilization of hardware. However, how would increasing the batch size affect the generalizability in adapting transformer networks? Since prior works on batch size scaling often focus on computer vision tasks and pre-training Smith et al. (2018); Goyal et al. (2017); Smith (2018); You et al. (2019a), we conduct an analysis of large-batch adaptation on pre-trained Transformers by performing a hyperparameter sweep on batch sizes {1K, 2K, 4K, 8K} and learning rates {1e-4, 3e-4, 5e-4, 7e-4, 9e-4, 1e-3, 3e-3}, where the learning rate range covers both linear scaling (Goyal et al., 2017) and sqrt scaling (You et al., 2019a). We report the validation accuracy in Figure 1(b). We make two observations: (1) the learning rate scales roughly with the square root of the increase of the mini-batch size, although the best learning rates do not always follow the sqrt rule; (2) there is a generalization gap between the small batch and large batch accuracies, and the gap becomes larger when the batch size increases. Furthermore, methods, such as LAMB You et al. (2019a), works well on pre-training with extremely large batch sizes ($\log_2 B = \{15, 16\}$) but do not close the generalization gap in adaptation (as shown in Section 5). These results pose the question: can we increase the batch size during adaptation in the interest of making adaptation more efficient while preserving generalization?

Hessian analysis. To further examine the generalization gap, we resort to the Hessian-based curvature analysis. Prior work Keskar et al. (2017); You et al. (2020) correlates the low generalization with sharp minima (which are characterized by a positive curvature of large magnitude in the pa-

parameter space). The indication is that a sharp local minimum also reflects a higher sensitivity of the loss even within the neighborhood of training data points and can attribute to the difficulty in generalization. Their hypothesis was that a larger noise due to the higher variance in gradient estimates computed using small mini-batches, in contrast to gradient estimates computed using large mini-batches, encourages the parameter weights to exit out of the basin of sharp minima and towards flatter minima which have better generalization.

To verify this hypothesis, we quantitatively measure the steepness of loss landscape by loading the checkpoint of an adapted model and computing the curvature, i.e., properties of the second derivative of the model, with respect to its parameters, for a fixed batch of samples. Following Yao et al. (2018), for a model $\Phi(x)$, we compute the largest eigenvalue of the model’s Hessian, $L_{\max}[\nabla_x^2 \Phi(x)]$, using the Hessian-vector product primitive and the power method. We use the largest eigenvalue as a measure of sharpness since the corresponding (top) eigenvector characterizes the direction of the largest change in gradient at a given point in the parameter space. From Figure 1(c), the largest eigenvalue of the model trained with a large batch (e.g., 1K) is much larger (e.g., 2.6x) than the small-batch baseline and with higher deviations (e.g., 3.9x). This result confirms that large-batch adaptation makes the loss landscape of the model more prone to ill-conditioning and less robust to perturbation, which helps explain the loss in generalization.

4 ADVERSARIAL NOISE FOR CLOSING LARGE-BATCH GENERALIZATION GAP

Given the generalization challenge of large-batch adaptation, we investigate strategies that may help close the generalization gap. We first describe our setup and then our analysis results.

Basic setup. Since language expressions are quite sensitive to individual words or clauses, where noises against those would likely generate incorrect or biased training data with wrong labels (Zhang & Yang, 2018). We follow prior success in applying adversarial training to NLP models (Miyato et al., 2017; Zhu et al., 2020) to have an adversarial training setup by applying noises to the continuous word embeddings instead of directly to discrete words or tokens.

$$\min_{x \in \mathbb{X}} \mathbb{E}_{\xi \sim Q} [g(x, \xi)] = \min_{x \in \mathbb{X}} \max_{y \in \mathbb{Y}} \mathbb{E}_{\xi \sim Q} [f(x, \xi) + \lambda r(x, y)] \quad (1)$$

where $g : \mathbb{X} \times \mathbb{Y} \rightarrow \mathbb{R}$ denotes the overall training objective, $r : \mathbb{X} \rightarrow \mathbb{R}$ denotes the augmented regularization and ξ denotes samples drawn from Q (for simplicity, we slightly abuse the notation in using ξ to denote the random variable, e.g. $\mathbb{E}_{\xi} [g(x, \xi)]$, or its empirical realizations, e.g. $\frac{1}{K} \sum_{k=1}^K g(x, \xi_k)$ for any K). The overall (outer) training objective involves a minimization problem in the parameter space while being stochastic with respect to the data space. The adversarial regularization (inner) term is a deterministic maximization problem operating in the data space conditioned on a fixed parameter configuration. We emphasize that this formulation is a two-player sequential (Jin et al., 2020), not simultaneous, game wherein the goal is to optimize a transformer network that is insensitive to adversarial noise.

Why Is Adversarial Noise Expensive? The generation of adversarial noise requires an extra PGA inner loop that standard training does not have. Figure 3 provides the time breakdown of optimization using PGA with $\mathcal{T} = 1$ (denoted as PGA-1). PGA-1 performs the perturbation and takes approximately the same time as making three forward passes (Fwd) through the network. This is because one step of PGA requires making one forward and backward pass (Bwd) over the entire network. The backward pass of the optimization takes roughly twice the amount of time as the standard backward step because the back-propagation is triggered twice to calculate the noise and the gradients. The time spent on the optimizer step function remains the same. In total, the optimization would slow down training by at least 2 times, even with $\mathcal{T}=1$. This motivates us to look at the effectiveness of different perturbation steps as well as the usefulness of perturbation from the initial epochs.

The Usefulness of Maximization Steps. Prior works often do multiple gradient computation steps ($\mathcal{T} > 1$) and take several times longer training time to produce adversaries (Madry et al., 2018; Zhu et al., 2020), likely because their focus is on generalization instead of computational efficiency. Subsequently, researchers presented Curriculum Adversarial Training (CAT) (Cai et al., 2018) and Annealing-based Adversarial Training (Ye et al., 2020), which progressively increase the perturbation with various strengths, cutting the adversarial training cost while maintaining good accuracy.

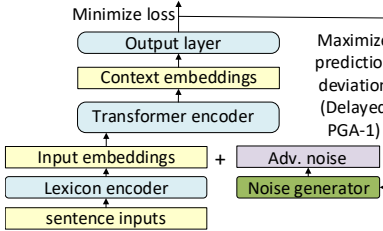


Figure 2: The architecture of the proposed method.

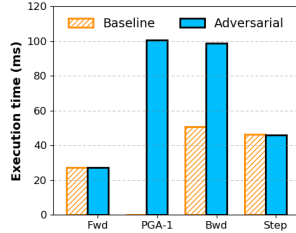


Figure 3: Time breakdown without and with PGA-1.

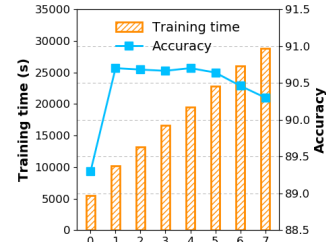


Figure 4: Impact of perturbation steps.

To investigate how CAT and similar methods affect large-scale NLP problems involving transformers, we evaluate the final accuracy and training cost of QNLI, varying the number of perturbation steps \mathcal{T} and report the results in Figure 4. Interestingly, although using a large \mathcal{T} helps to produce stronger noises, we find that this does not lead to improved accuracy, despite the fact that the training overhead still increases almost linearly. In fact, the best accuracy is achieved with $\mathcal{T} = 1$.

Why Is One-shot Perturbation Sufficient? We note that the model has two components, namely, the parameter space and data space. First, unlike the minimization in the parameter space, which is stochastic, the maximization in the data space is deterministic. Second, with respect to the testing phase, the numerical convergence in the model’s parameter space is of primary importance rather than the numerical convergence in the data space, i.e., the maximization is an auxiliary procedure that augments the training phase to make the parameter space “aware” of effects of the batch size across epochs. Due to these two points, at a certain epoch, for a given batch, the marginal utility of an additional PGA step is low, and we are able to get away with inexact deterministic maximization. Therefore, we apply PGA-1 in our large-batch optimization scheme, given that it produces sufficiently good solutions while being much more computationally efficient.

The Usefulness of Adversarial Noise in the Adaptation Process. Given that PGA-1 still adds an overhead factor of 2, we are motivated to further reduce the overhead of adversarial noise. In particular, we investigate how useful adversarial noises are in the whole large-batch optimization process. We conduct additional experiments to measure the final accuracy corresponding to starting from a regular fine-tuning and then enabling PGA-1 for $t \geq t_s$ where $t_s \in [T]$. Our observation is that enabling PGA-1 from the beginning does not offer much improvement in accuracy, whereas adversarial noise becomes more potent as the model begins to stabilize towards the end of training (more detailed results in Appendix A.1). In general, at initialization, the model’s parameters are relatively far from their final values and are less likely to get stuck at local minima. Therefore the adversarial noises generated in the initial training iterations are quite different from the noises towards the end of training because they would not maximize the adversarial loss in Equation 1. This hypothesis suggests that we might be able to inject adversarial noise in the later training process while still leveraging it to improve generalizability. We remark that this phenomenon has been observed by prior work on computer vision tasks (Cai et al., 2018; Gupta et al., 2020).

Would Alternative Noise Injection Strategies Work? One may wonder what would happen if we use a different type of noise. To answer this question, we have performed additional experiments by adding Gaussian noise to the embeddings as well as comparing one-hot labels and label probability for generating adversarial noises. Overall, we find that explicitly enforcing the smoothness of the loss landscape via adversarial noise can result in better improvement and label probability (soft) consistently leads to higher accuracy than using the ground truth (hard). Section 5 provides more detailed results.

Proposed Method. Combining the formulation with the above investigations, we propose ScaLA, which adopts both one-shot perturbation and delayed noise injection as they help reduce the computation cost of adversarial noise while still allowing it to preserve the generalization benefit. The full procedure of ScaLA is provided in Algorithm 1, whose convergence rate is characterized in Theorem 4.1.

Theorem 4.1 (Complexity of Algorithm 1; Informal – Details in Appendix D). *Consider the problem in Equation 1. Let $t_s = 0$. Setting the outer learning rate as $\eta = O(1/\sqrt{T})$ and scaling batch size as $b = O(T)$, for Algorithm 1, we have $\mathbb{E} [\|\nabla_{g_{1/2\alpha}}(\bar{x})\|^2] \leq O(\epsilon + \kappa_\alpha/\sqrt{T})$ where \bar{x} is the*

Algorithm 1**ScaLA**

```

1: Input: Epochs  $T$ , delay  $t_s$ , perturbation (inner) step size  $\rho$ , clipping radius  $\omega$ , regularization
   strength  $\lambda$ , (outer) learning rate  $\eta$ 
2: Output:  $h$ -layer transformer model  $\Phi$  with converged robust parameters  $\bar{x} := x_T$ 
3: for  $t \in [T]$  do
4:   for worker  $p \in [P]$  do
5:     for mini-batch  $\xi_p \sim Q$  do
6:        $r(x_t) \leftarrow 0, \gamma \leftarrow \Phi(x, \xi_p)$ , select  $y_0$ 
7:       if  $t \geq t_s$  then
8:         ▷ Check delay condition
9:          $y_1 \leftarrow \Pi_\omega(y_0 + \rho \nabla_y r(x_t, y))$  ▷ Generate adversarial noise with PGA-1
10:         $\underline{r}(x_t) \leftarrow \text{KL}_{\text{sym}}(\gamma, \Phi(x_{t-1}, y_1))$ 
11:        end if
12:         $g(x_t, \xi_p) \leftarrow f(x_{t-1}, \xi_p) + \lambda \underline{r}(x_t)$ 
13:         $\nabla_x g(x_t, \xi_p) \leftarrow \text{Backward pass on } \Phi$ 
14:      end for
15:    end for
16:     $\widehat{\nabla}_x g(x_t) \leftarrow \frac{1}{B} \sum_{p=1}^P \nabla_x g(x_t, \xi_p)$ 
17:     $x_t^i \leftarrow x_{t-1}^i - \eta_t \widehat{\nabla}_x g(x_t)$ 
18:  end for

```

estimator obtained from running T steps of Algorithm 1 and picking x_t uniformly at random for $t \in [T]$. Here, ϵ is the error due to the approximate inner maximization oracle, α characterizes the smoothness of $f(x, \cdot)$, $g_{1/2\alpha}$ is the Moreau-envelope of g and $\kappa_\alpha = \max_i \alpha_i / \min_i \alpha_i$.

5 EVALUATION

We evaluate the effectiveness of ScaLA in adapting pre-trained transformer networks over a set of NLP tasks.

Hardware. We conduct the evaluation using 2 NVIDIA DGX-2 nodes. Each node consists of 16 NVIDIA V100 GPUs. The nodes are connected with InfiniBand using a 648-port Mellanox MLNX-OS CS7500 switch. **Model/Dataset.** We study adaptation on pre-trained BERT_{base} model RoBERTa_{large}, and T5_{base} hosted by HuggingFace (Wolf et al., 2020). We use the GLUE benchmark (Wang et al., 2019a), which is a collection of sentence or sentence-pair natural language understanding tasks including question answering, sentiment analysis, and textual entailment. We also include SQuAD-v2 in our evaluation. We exclude tasks in GLUE that have very small datasets (e.g., CoLA, RTE). We report the details about the hyperparameters in Appendix B.

5.1 MAIN RESULTS – ADAPTATION TIME ACCELERATION

We first compare the following schemes: (1) **Single GPU + SB:** This is the existing PyTorch implementation of Transformer fine-tuning from HuggingFace (HF), using small batch (SB) sizes (e.g., 32). (2) **Multi-GPU + SB:** This is multi-GPU PyTorch implementation using DistributedDataParallel (Li et al., 2020), (3) **Multi-GPU + LB + FreeLB:**, this is the work described in Zhu et al. (2020) using large minibatches (LB), e.g., 1K, and perturbation step $K = 5$ for adaptation, and (4) **Multi-GPU + LB + ScaLA:** This is our approach as described in Algorithm 1. Table 1 shows results on MNLI, QNLI, QQP, and SST2, which are larger datasets and less sensitive to random seeds. $n \times g$ refers to P_n nodes each with P_g GPUs for a total of $P = P_n P_g$ homogeneous workers (e.g., 32 GPUs on 2 NVIDIA DGX-2 nodes). For a fair comparison, we reproduce BERT and RoBERTa baseline. Our reproduced baseline achieves the same or slightly higher accuracy than the originally reported results in (Devlin et al., 2019) and (Liu et al., 2019). We now discuss our results and observations.

Adaptation time analysis. Compared with single-GPU training, the multi-GPU baseline leads to only modest training speedup improvements, e.g., with $1.5 - 2.4 \times$ faster training speed for both BERT and RoBERTa, even with $32 \times$ more compute resources. The speedup is limited because of the small mini-batches (e.g., 32) used for adaptation, which do not provide a sufficient workload to fully utilize the underlying hardware. Thus, communication overhead becomes the dominant part, and the adaptation often struggles to obtain speedups even with more workers. In contrast, ScaLA achieves

Table 1: The adaptation time and accuracy results on GLUE benchmark. ScaLA achieves the same average accuracy as the baseline while providing up to $18\times$ speedups than a single GPU, and up to $9.8\times$ speedups with the same amount of hardware.

| BERT _{base} | n×g | bsz | MNLi-m | | | QNLI | | | QQP | | | SST-2 | | | Avg. |
|--------------------------|------|-----|--------|-------------|------|-------|-------------|------|-------|-------------|-----------|-------|------------|------|-------------|
| | | | Steps | Time | Acc. | Steps | Time | Acc. | Steps | Time | Acc/F1 | Steps | Time | Acc. | |
| Devlin et al. 2019 | | | | | 84.4 | | | 88.4 | | | - | | | 92.7 | - |
| Baseline (B=32) | 1x1 | 32 | 73632 | 19635 | 84.8 | 19644 | 5535 | 90.6 | 68226 | 16494 | 91/88.0 | 12630 | 2736 | 93.1 | 89.4 |
| Baseline (B=32) | 2x16 | 32 | 73632 | 8848 | 84.8 | 19644 | 2408 | 90.6 | 68226 | 11311 | 91/88.0 | 12630 | 1494 | 93.1 | 89.4 |
| FreeLb (B=1K) | 2x16 | 1K | 2301 | 5953 | 85.2 | 615 | 1944 | 90.3 | 2133 | 19030 | 91.2/88.2 | 396 | 680 | 92.8 | 89.5 |
| ScaLA (B=1K) | 2x16 | 1K | 2301 | 1323 | 85.1 | 615 | 432 | 90.0 | 2133 | 4229 | 90.9/87.7 | 396 | 151 | 93.5 | 89.4 |
| RoBERTa _{large} | n×g | bsz | MNLi-m | | | QNLI | | | QQP | | | SST-2 | | | Avg. |
| | | | Steps | Time | Acc. | Steps | Time | Acc. | Steps | Time | Acc/F1 | Steps | Time | Acc. | |
| Liu et al. 2020 | | | | | 90.2 | | | 94.7 | | | 92.2/- | | | 96.4 | - |
| Baseline (B=32) | 1x1 | 32 | 73632 | 43090 | 90.5 | 19644 | 14188 | 94.7 | 68226 | 40945 | 92.0/89.4 | 12630 | 4940 | 96.4 | 92.5 |
| Baseline (B=32) | 2x16 | 32 | 73632 | 18114 | 90.5 | 19644 | 4842 | 94.7 | 68226 | 16614 | 92.0/89.4 | 12630 | 3072 | 96.4 | 92.5 |
| FreeLb (B=1K) | 2x16 | 1K | 2301 | 15133 | 91.2 | 615 | 5256 | 95.2 | 2133 | 10818 | 92.5/90.0 | 396 | 1804 | 96.9 | 93.3 |
| ScaLA (B=1K) | 2x16 | 1K | 2301 | 3363 | 90.9 | 615 | 1168 | 95.1 | 2133 | 2404 | 92.3/89.8 | 396 | 401 | 96.7 | 92.9 |

up to $18\times$ speedups over the single-GPU baseline with 32 GPUs. When using the same number of GPUs (e.g., 32), ScaLA is 2.7–9.8× faster. The speedups come from three aspects: (1) the improved hardware efficiency for each worker from increased per-worker micro-batch size; (2) the reduced all-reduce communication overhead since it takes fewer communication rounds to process the same number of samples in one epoch; (3) the lightweight adversarial noise incurs only a small portion of the total training overhead. Finally, ScaLA obtains the speedups while achieving the same accuracy (88.4 vs. 88.4) average accuracy for BERT and higher accuracy (92.9 vs. 92.5) for RoBERTa as the baselines. ScaLA is 4.5 times faster than FreeLb while achieving similar accuracy on BERT (89.4 vs. 89.5) and RoBERTa (92.9 vs. 93.5). ScaLA is faster than FreeLb because FreeLb does not consider much about the training cost and performs multiple ascent steps to calculate adversaries across the full training process. As a matter of fact, FreeLb is even slower to run than vanilla baseline (e.g., QNLI on RoBERTa). In contrast, ScaLA analyzes the computational efficiency of adversarial large-batch optimization and introduces several simple yet effective approaches to reduce the adversarial noise cost, which leads to overall improved computational efficiency.

Generalizability analysis. Since there are very few works on large-batch adaptation, we create several baselines to compare with ScaLA: (1) Multi-GPU + LB + Tuning LR: This configuration uses large mini-batches (e.g., 1K), and applies a heuristic-based scheduling rule (e.g., square root) combined with an extensive grid search for learning rates; (2) Multi-GPU + LB + LAMB: Uses LAMB (You et al., 2019a) optimizer for large-batch adaptation. We make several observations from the results in Table 2. First, compared with the baseline accuracy reported in the paper, the accuracy of Multi-GPU + LB drops by close to 1 point (88.4 vs. 89.4, and 92.1 vs. 92.9) on average and close to 2 points for some tasks (e.g., QQP on BERT), indicating that it is challenging to obtain on-par accuracy with large-batch optimizations for adaptation despite with heavy hyperparameter tuning. Second, since LAMB is designed primarily for improving the convergence of pre-training instead of the adaptation, its ability to accelerate adaptation has yet to be proven. In our experiments, LAMB leads to only marginal improvements (88.6 vs. 88.4, and 92.1 vs. 92.1) than the baseline and is 0.8 points lower than the small-batch baseline. This is because LAMB does not directly minimize the sharpness of the loss landscape, so it can still lead to poor generalizability during adaptation. With ScaLA, we are able to close the generalization gap from large-batch optimization (89.4 vs. 89.4, and 92.5 vs. 92.9) and achieve 0.8 points higher accuracy (89.4 vs. 88.6, 92.9 vs. 92.1) than LAMB on both BERT and RoBERTa. ScaLA improves generalizability because it introduces adversarial noise in the large-batch optimization process, which serves as a regularizer. By training the network to be robust to such perturbations, the model loss landscape is smoothed out, leading to improved generalization.

5.2 EXPERIMENT – ANALYSIS RESULTS

Ablation analysis: We study the importance of components in ScaLA. We set t_s to 0, which denotes as *w/o Delaying PGA-1*. We replace the outer minimization to use ADAM (Kingma & Ba, 2015), which is noted as *w/o Groupwise LR*. We set λ to 0, which denotes as *w/o PGA-1*. The results are reported in Table 3.

The results in Table 3 show that the removal of either design element would result in a performance drop. For example, removing PGA-1 leads to 0.8 points accuracy drop (88.6 vs. 89.4), indicating that adversarial noise is crucial for improving the generalizability of large-batch adaptation. More-

Table 2: The comparison results between ScaLA and alternative methods for large-batch adaptation on the GLUE benchmark, which show that ScaLA achieves higher accuracy than baselines after training the same number of samples and steps.

| BERT _{base} | n×g | Batch size | MNLi-m | | | QNLI | | | QQP | | | SST-2 | | | Avg. |
|----------------------|------|------------|--------|------|-------------|-------|------|-------------|-------|------|------------------|-------|------|-------------|-------------|
| | | | Steps | Time | Acc. | Steps | Time | Acc. | Steps | Time | Acc/F1 | Steps | Time | Acc. | |
| Vanilla (B=1K) | 2x16 | 1K | 2301 | 1148 | 84.3 | 615 | 349 | 89.3 | 2133 | 2892 | 89.6/86.1 | 396 | 134 | 93 | 88.4 |
| LAMB (B=1K) | 2x16 | 1K | 2301 | 1180 | 84.1 | 615 | 359 | 89.6 | 2133 | 2978 | 90.5/87.0 | 396 | 139 | 92.4 | 88.6 |
| ScaLA (B=1K) | 2x16 | 1K | 2301 | 1323 | 85.1 | 615 | 432 | 90.0 | 2133 | 4229 | 90.9/87.7 | 396 | 151 | 93.5 | 89.4 |

| RoBERTa _{large} | n×g | Batch size | MNLi-m | | | QNLI | | | QQP | | | SST-2 | | | Avg. |
|--------------------------|------|------------|--------|------|-------------|-------|------|-------------|-------|------|------------------|-------|------|-------------|-------------|
| | | | Steps | Time | Acc. | Steps | Time | Acc. | Steps | Time | Acc/F1 | Steps | Time | Acc. | |
| Vanilla (B=1K) | 2x16 | 1K | 2301 | 2514 | 90.1 | 615 | 936 | 94.3 | 2133 | 1874 | 91.7/89.1 | 396 | 317 | 95.9 | 92.1 |
| LAMB (B=1K) | 2x16 | 1K | 2301 | 2646 | 90.5 | 615 | 973 | 94.5 | 2133 | 1998 | 91.3/88.5 | 396 | 324 | 96.2 | 92.1 |
| ScaLA (B=1K) | 2x16 | 1K | 2301 | 3363 | 90.9 | 615 | 1168 | 95.1 | 2133 | 2404 | 92.3/89.8 | 396 | 401 | 96.7 | 92.9 |

Table 3: Ablation study of ScaLA using BERT_{base} on GLUE tasks.

| | MNLi-m | | QNLI | | QQP | | SST-2 | | Avg. | Speedup |
|--------------------|--------|------|------|------|-------|-----------|-------|------|------|---------|
| | Time | Acc. | Time | Acc. | Time | Acc/F1 | Time | Acc. | | |
| Baseline | 19635 | 84.8 | 5535 | 90.6 | 16494 | 91/88.0 | 2736 | 93.1 | 89.4 | 1 |
| ScaLA | 1323 | 85.1 | 432 | 90 | 4229 | 90.9/87.7 | 151 | 93.5 | 89.4 | 12.4 |
| w/o Delaying PGA-1 | 2503 | 85.2 | 726 | 90.2 | 6407 | 91.3/88.3 | 272 | 93.1 | 89.5 | 7.0 |
| w/o Groupwise LR | 1290 | 85.0 | 422 | 89.9 | 4212 | 90.7/87.6 | 146 | 93.0 | 89.2 | 12.7 |
| w/o PGA-1 | 1180 | 84.1 | 359 | 89.6 | 2978 | 90.5/87.0 | 139 | 92.4 | 88.6 | 14.3 |

over, if we perform PGA-1 without delayed injection, the average accuracy increases by 0.1 points (89.5 vs. 89.4), but the execution time is increased by 1.5–1.9x, indicating the importance of having lightweight adversarial noise for speeding up the adaptation. Finally, removing group-wise learning rates leads to a small 0.2 points accuracy drop (89.2 vs. 89.4), indicating that ScaLA still achieves benefits without group-wise learning rates (89.2 vs. 88.6), but they are complementary to each other.

Table 4: Evaluation results of T5_{base} on SQuAD-v2. Table 5: Alternatives to generate perturbations using random noise, ground-truth, and label probability.

| T5/SQuAD-v2 | n | bsz | EM | F1 | Time |
|-------------|----|------|-------|-------|-------|
| Baseline | 1 | 12 | 78.53 | 81.71 | 2h10m |
| Baseline | 16 | 12 | N/A | N/A | N/A |
| Baseline | 16 | 1024 | 78.27 | 81.47 | 8m |
| FreeLb | 16 | 1024 | 79.04 | 82.65 | 66m |
| ScaLA | 16 | 1024 | 78.98 | 82.55 | 18m |

| Model | MNLi-m | QNLI | QQP | SST-2 | Avg |
|----------------|-------------|-----------|------------------|-------------|-------------|
| Baseline | 84.3 | 89.3 | 89.6/86.1 | 93 | 88.4 |
| Gaussian noise | 84.5 | 89.4 | 90.3/87.0 | 92.6 | 88.7 |
| ScaLA (GT) | 84.1 | 89.6 | 90.7/87.6 | 93.2 | 89.0 |
| ScaLA (LP) | 85.1 | 90 | 90.9/87.7 | 93.5 | 89.4 |

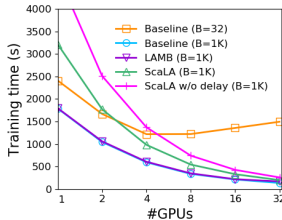


Figure 5: Comparison of scalability using different large-batch optimization methods on SST-2.

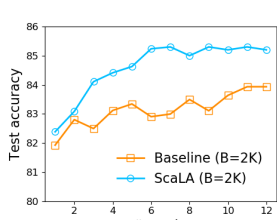
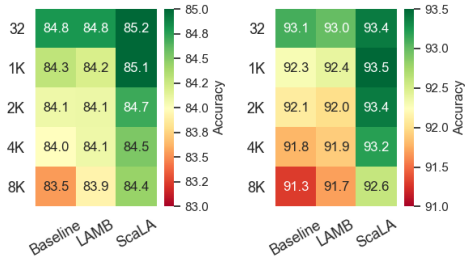


Figure 6: Comparison of test accuracy by training the baseline longer.



(a) MNLi-m (b) SST-2

Figure 7: Comparison of accuracy under even larger batch sizes.

Application to Encoder-decoder Architectures. We used T5-base (220M) model from the HuggingFace model zoo’s checkpoint. We applied ScaLA to adapt pre-trained T5 on SQuAD-v2 with a batch size of 1k. Overall, the observations on this task are in line with the observations we obtained from the pre-trained BERT/RoBERTa and the GLUE benchmarks in the manuscript: First, ScaLA achieves 0.37/0.84 points higher EM/F1 scores (78.98/82.55 vs. 78.53/81.71) with a 7.2x faster adaptation speed (2 hours 10 minutes vs. 18 minutes) than the baseline. This is achieved by (1) using large batch sizes to improve the aggregated training throughput, and (2) using lightweight injected adversarial noises to improve generalization. Second, compared with FreeLb, ScaLA achieves a similar EM/F1 score but with a 3.6x faster speed. ScaLA is faster because it removes redundant perturbation steps and leverages delayed perturbations, greatly reducing the overhead from adversarial noise. Third, the scalability of the baseline is limited by the small batch sizes it uses, causing

severe under-utilization when scaling out the adaptation to multiple GPUs. On the other hand, although increasing the batch sizes and learning rates allows the adaptation to finish in a much shorter time (from 2 hours 10 minutes to 8 minutes), the final accuracy is lower than the baseline. In contrast, ScaLA allows the adaptation to match and even sometimes exceed the baseline accuracy while offering a much faster adaptation speed.

Curvature analysis. We measure the steepness of the loss landscape again after applying ScaLA. As shown in Fig. 1(c), the largest eigenvalue of the model becomes much smaller ($6.9\times$) with lower deviations with ScaLA and is slightly better than the small batch baseline, which is a strong indication that our approach enforces the smoothness of the model that leads to the accuracy improvement.

Comparison with random noise. We have performed additional experiments by adding Gaussian noise to the embeddings. Table 5 that random noise indeed can improve the accuracy for MNLI-m (84.3 vs. 84.5), QNLI (89.3 vs. 89.4), and QQP (90.3/87.0 vs. 89.6/86.1) over the baseline, but it also leads to worse results on SST-2 (93. vs. 92.6). Compared with ScaLA, random noise consistently falls behind ScaLA in its ability to reduce the generalization error on all tested tasks and is on average 0.7 points lower than ScaLA (88.7 vs. 89.4). These results indicate that ScaLA’s approach of explicitly enforcing the smoothness of the loss landscape can result in better improvement.

Perturbations via ground-truth vs. label probability. We also create one-hot labels and use those to generate perturbations instead of using label probability generated by the network. Table 5 shows that using label probability (LP) consistently leads to higher accuracy than using the ground-truth (GT), e.g., 89.4 vs. 89.0 on average. Label probability leads to better generalization, probably because it provides a better measurement of the adversarial direction, which is the direction in the input space in which the label probability of the model is most sensitive to small perturbations.

Scalability analysis varying GPUs. Figure 5 shows the scalability comparison on SST-2 after optimizations. While the speedup still plateaus at 4 GPUs with a small batch size (e.g., $B = 32$), the four large-batch configurations are able to scale well up to 32 GPUs and take a similar amount of time with 32 GPUs. ScaLA scales better than ScaLA without delaying PGA-1, and achieves a much faster training speed, especially in the 1-16 GPU range.

Train longer, generalize better? Despite improved adaptation speed, one may still wonder whether simply performing large-batch adaptation longer would also close the generalization gap. Figure 6 shows the comparison between ScaLA and the baseline on a batch size of 2K. ScaLA obtains an accuracy of 85.2 after 6 epochs of training, whereas the baseline has difficulty to reach 84 after training twice longer (e.g., 12 epochs). ScaLA achieves better accuracy because it explicitly penalizes model weights from getting stuck at sharp minima, leading to better generalizability.

Generalizability under different batch sizes. We also evaluate how different batch sizes affect the generalizability of adapting transformers. Figure 7 shows the results on MNLI-m and SST-2. We make two major observations: (1) The accuracy tends to drop as the batch size increases. (2) While both the baseline and LAMB suffer from significant accuracy drop by drastically increasing the batch size (e.g., from 32 to 8K), ScaLA is able to mitigate the generalization gap and consistently achieves higher accuracy than the baseline (e.g., 84.4 vs. 83.5 for MNLI, and 92.6 vs. 91.3 for SST-2 at batch size 8K) and LAMB (e.g., 84.4 vs. 83.9 for MNLI, and 92.6 vs. 91.7 for SST-2 at batch size 8K). These results indicate the benefit of ScaLA is maintained by further increasing the batch size, which could bring even greater speedups when increasing the data parallelism degree.

6 CONCLUSIONS AND FUTURE DIRECTIONS

In this paper, we study how to accelerate the adaptation speed of pre-trained Transformer models for NLU tasks. We introduce ScaLA, an efficient large-batch adaptation method using carefully injected lightweight adversarial noises. The experiment results show that ScaLA obtains up to $9.8\times$ speedups on adapting transformer networks and outperforms state-of-the-art large-batch optimization methods in generalizability. Given the promising results of ScaLA on accelerating the adaptation speed, it opens new research opportunities for applying ScaLA to accelerate the more expensive pre-training tasks as well as emerging pre-trained transformer networks for computer vision domains tasks.

REFERENCES

- Microsoft DeepSpeed achieves the fastest BERT training time. <https://www.deepspeed.ai/news/2020/05/27/fastest-bert-training.html>. Accessed: 05-20-2021.
- Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. *CoRR*, abs/2005.14165, 2020. URL <https://arxiv.org/abs/2005.14165>.
- Qi-Zhi Cai, Chang Liu, and Dawn Song. Curriculum adversarial training. In Jérôme Lang (ed.), *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJ-CAI 2018, July 13-19, 2018, Stockholm, Sweden*, pp. 3740–3747. ijcai.org, 2018.
- Yong Cheng, Lu Jiang, and Wolfgang Macherey. Robust neural machine translation with doubly adversarial inputs. In Anna Korhonen, David R. Traum, and Lluís Màrquez (eds.), *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pp. 4324–4333. Association for Computational Linguistics, 2019.
- Kevin Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. ELECTRA: pre-training text encoders as discriminators rather than generators. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020.
- Damek Davis and Dmitriy Drusvyatskiy. Stochastic subgradient method converges at the rate $o(k^{-1/4})$ on weakly convex functions. *arXiv preprint arXiv:1802.02988*, 2018.
- Damek Davis and Dmitriy Drusvyatskiy. Stochastic model-based minimization of weakly convex functions. *SIAM Journal on Optimization*, 29(1):207–239, 2019.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019*, pp. 4171–4186, 2019.
- Saeed Ghadimi and Guanghui Lan. Stochastic first-and zeroth-order methods for nonconvex stochastic programming. *SIAM Journal on Optimization*, 23(4):2341–2368, 2013.
- Linyuan Gong, Di He, Zhuohan Li, Tao Qin, Liwei Wang, and Tie-Yan Liu. Efficient training of BERT by progressively stacking. In *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, pp. 2337–2346, 2019.
- Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In Yoshua Bengio and Yann LeCun (eds.), *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.
- Priya Goyal, Piotr Dollár, Ross B. Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large minibatch SGD: training imagenet in 1 hour. *CoRR*, abs/1706.02677, 2017.
- Sidharth Gupta, Parijat Dube, and Ashish Verma. Improving the affordability of robustness training for dnns. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR Workshops 2020, Seattle, WA, USA, June 14-19, 2020*, pp. 3383–3392. IEEE, 2020.

- Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. Towards a unified view of parameter-efficient transfer learning. *CoRR*, abs/2110.04366, 2021.
- Elad Hoffer, Itay Hubara, and Daniel Soudry. Train longer, generalize better: closing the generalization gap in large batch training of neural networks. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pp. 1731–1741, 2017.
- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin de Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for NLP. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pp. 2790–2799. PMLR, 2019.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *CoRR*, abs/2106.09685, 2021.
- Yanping Huang, Youlong Cheng, Ankur Bapna, Orhan Firat, Dehao Chen, Mia Xu Chen, HyoukJoong Lee, Jiquan Ngiam, Quoc V. Le, Yonghui Wu, and Zhifeng Chen. Gpipe: Efficient training of giant neural networks using pipeline parallelism. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019*, pp. 103–112, 2019.
- Prateek Jain and Purushottam Kar. Non-convex optimization for machine learning. *arXiv preprint arXiv:1712.07897*, 2017.
- Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and Tuo Zhao. SMART: robust and efficient fine-tuning for pre-trained natural language models through principled regularized optimization. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel R. Tetreault (eds.), *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pp. 2177–2190. Association for Computational Linguistics, 2020.
- Chi Jin, Praneeth Netrapalli, and Michael Jordan. What is local optimality in nonconvex-nonconcave minimax optimization? In *International Conference on Machine Learning*, pp. 4880–4889. PMLR, 2020.
- Norman P. Jouppi, Cliff Young, Nishant Patil, David Patterson, Gaurav Agrawal, Raminder Bajwa, Sarah Bates, Suresh Bhatia, Nan Boden, Al Borchers, Rick Boyle, Pierre-luc Cantin, Clifford Chao, Chris Clark, Jeremy Coriell, Mike Daley, Matt Dau, Jeffrey Dean, Ben Gelb, Tara Vazir Ghaemmaghami, Rajendra Gottipati, William Gulland, Robert Hagmann, C. Richard Ho, Doug Hogberg, John Hu, Robert Hundt, Dan Hurt, Julian Ibarz, Aaron Jaffey, Alek Jaworski, Alexander Kaplan, Harshit Khaitan, Daniel Killebrew, Andy Koch, Naveen Kumar, Steve Lacy, James Laudon, James Law, Diemthu Le, Chris Leary, Zhuyuan Liu, Kyle Lucke, Alan Lundin, Gordon MacKean, Adriana Maggiore, Maire Mahony, Kieran Miller, Rahul Nagarajan, Ravi Narayanaswami, Ray Ni, Kathy Nix, Thomas Norrie, Mark Omernick, Narayana Penukonda, Andy Phelps, Jonathan Ross, Matt Ross, Amir Salek, Emad Samadiani, Chris Severn, Gregory Sizikov, Matthew Snelham, Jed Souter, Dan Steinberg, Andy Swing, Mercedes Tan, Gregory Thorson, Bo Tian, Horia Toma, Erick Tuttle, Vijay Vasudevan, Richard Walter, Walter Wang, Eric Wilcox, and Doe Hyun Yoon. In-Datacenter Performance Analysis of a Tensor Processing Unit. In *Proceedings of the 44th Annual International Symposium on Computer Architecture, ISCA '17*, pp. 1–12, 2017. ISBN 978-1-4503-4892-8. doi: 10.1145/3079856.3080246. URL <http://doi.acm.org/10.1145/3079856.3080246>.
- Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Peter Tang. On large-batch training for deep learning: Generalization gap and sharp minima. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017.
- Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. In *3rd International Conference on Learning Representations, ICLR 2015*, 2015.

- Shen Li, Yanli Zhao, Rohan Varma, Omkar Salpekar, Pieter Noordhuis, Teng Li, Adam Paszke, Jeff Smith, Brian Vaughan, Pritam Damania, and Soumith Chintala. Pytorch distributed: Experiences on accelerating data parallel training. *Proc. VLDB Endow.*, 13(12):3005–3018, 2020.
- Tianyi Lin, Chi Jin, and Michael Jordan. On gradient descent ascent for nonconvex-concave min-max problems. In *International Conference on Machine Learning*, pp. 6083–6093. PMLR, 2020.
- Xiaodong Liu, Hao Cheng, Pengcheng He, Weizhu Chen, Yu Wang, Hoifung Poon, and Jianfeng Gao. Adversarial training for large neural language models. *CoRR*, abs/2004.08994, 2020.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized BERT pretraining approach. *CoRR*, abs/1907.11692, 2019.
- Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018.
- Takeru Miyato, Andrew M. Dai, and Ian J. Goodfellow. Adversarial training methods for semi-supervised text classification. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017.
- Zachary Nado, Justin Gilmer, Christopher J. Shallue, Rohan Anil, and George E. Dahl. A large batch optimizer reality check: Traditional, generic optimizers suffice across batch sizes. *CoRR*, abs/2102.06356, 2021.
- Jonas Pfeiffer, Aishwarya Kamath, Andreas Rücklé, Kyunghyun Cho, and Iryna Gurevych. Adapterfusion: Non-destructive task composition for transfer learning. In Paola Merlo, Jörg Tiedemann, and Reut Tsarfaty (eds.), *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, EACL 2021, Online, April 19 - 23, 2021*, pp. 487–503. Association for Computational Linguistics, 2021.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *CoRR*, abs/1910.10683, 2019. URL <http://arxiv.org/abs/1910.10683>.
- Hassan Rafique, Mingrui Liu, Qihang Lin, and Tianbao Yang. Weakly-convex-concave min-max optimization: provable algorithms and applications in machine learning. *Optimization Methods and Software*, pp. 1–35, 2021.
- Samyam Rajbhandari, Jeff Rasley, Olatunji Ruwase, and Yuxiong He. Zero: Memory optimization towards training A trillion parameter models. *CoRR*, abs/1910.02054, 2019.
- Ralph Tyrell Rockafellar. *Convex analysis*. Princeton university press, 2015.
- Noam Shazeer, Youlong Cheng, Niki Parmar, Dustin Tran, Ashish Vaswani, Penporn Koanantakool, Peter Hawkins, HyoukJoong Lee, Mingsheng Hong, Cliff Young, Ryan Sepassi, and Blake A. Hechtman. Mesh-tensorflow: Deep learning for supercomputers. In *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018*, pp. 10435–10444, 2018.
- Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan Catanzaro. Megatron-lm: Training multi-billion parameter language models using model parallelism. *CoRR*, abs/1909.08053, 2019.
- Leslie N. Smith. A disciplined approach to neural network hyper-parameters: Part 1 - learning rate, batch size, momentum, and weight decay. *CoRR*, abs/1803.09820, 2018. URL <http://arxiv.org/abs/1803.09820>.
- Samuel L. Smith, Pieter-Jan Kindermans, Chris Ying, and Quoc V. Le. Don’t decay the learning rate, increase the batch size. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018.

- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017*, pp. 5998–6008, 2017.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In *7th International Conference on Learning Representations*, 2019a.
- Dilin Wang, ChengYue Gong, and Qiang Liu. Improving neural language modeling via adversarial training. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pp. 6555–6565. PMLR, 2019b.
- Ruize Wang, Duyu Tang, Nan Duan, Zhongyu Wei, Xuanjing Huang, Jianshu Ji, Guihong Cao, Daxin Jiang, and Ming Zhou. K-adapter: Infusing knowledge into pre-trained models with adapters. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (eds.), *Findings of the Association for Computational Linguistics: ACL/IJCNLP 2021, Online Event, August 1-6, 2021*, volume ACL/IJCNLP 2021 of *Findings of ACL*, pp. 1405–1418. Association for Computational Linguistics, 2021.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Transformers: State-of-the-art natural language processing. In Qun Liu and David Schlangen (eds.), *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, EMNLP 2020 - Demos, Online, November 16-20, 2020*, pp. 38–45. Association for Computational Linguistics, 2020.
- Zhewei Yao, Amir Gholami, Qi Lei, Kurt Keutzer, and Michael W. Mahoney. Hessian-based analysis of large batch training and robustness to adversaries. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pp. 4954–4964, 2018.
- Nanyang Ye, Qianxiao Li, Xiao-Yun Zhou, and Zhanxing Zhu. Amata: An annealing mechanism for adversarial training acceleration. *CoRR*, abs/2012.08112, 2020.
- Yang You, Jing Li, Jonathan Hseu, Xiaodan Song, James Demmel, and Cho-Jui Hsieh. Reducing BERT pre-training time from 3 days to 76 minutes. *CoRR*, abs/1904.00962, 2019a.
- Yang You, Jing Li, Sashank Reddi, Jonathan Hseu, Sanjiv Kumar, Srinadh Bhojanapalli, Xiaodan Song, James Demmel, Kurt Keutzer, and Cho-Jui Hsieh. Large batch optimization for deep learning: Training bert in 76 minutes. *arXiv preprint arXiv:1904.00962*, 2019b.
- Yang You, Yuhui Wang, Huan Zhang, Zhao Zhang, James Demmel, and Cho-Jui Hsieh. The limit of the batch size. *CoRR*, abs/2006.08517, 2020.
- Dongxu Zhang and Zhichao Yang. Word embedding perturbation for sentence classification. *CoRR*, abs/1804.08166, 2018.
- Minjia Zhang and Yuxiong He. Accelerating training of transformer-based language models with progressive layer dropping. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020.
- Shuai Zheng, Haibin Lin, Sheng Zha, and Mu Li. Accelerated large batch optimization of BERT pretraining in 54 minutes. *CoRR*, abs/2006.13484, 2020.
- Chen Zhu, Yu Cheng, Zhe Gan, Siqi Sun, Tom Goldstein, and Jingjing Liu. Freelib: Enhanced adversarial training for natural language understanding. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020.

A ADDITIONAL ANALYSIS RESULTS

In the part, we present results that are not included in the main text due to the space limit.

A.1 THE USEFULNESS OF ADVERSARIAL NOISES AT DIFFERENT EPOCHS

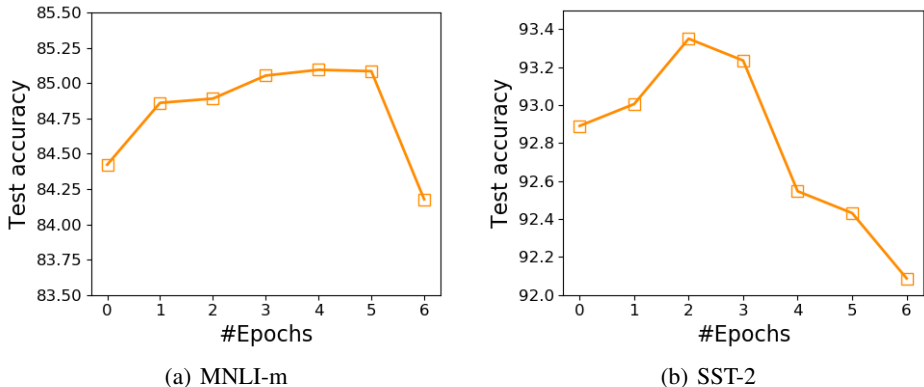


Figure 8: Accuracy results from delaying the injection of adversarial noises at different epochs.

In Section 4, we mention that no adversaries are needed at the initial epochs of adaptation. To verify, we conduct experiments to measure the final accuracy corresponding to starting from regular training and switching to PGA-1 after t_s epochs, where $t_s \in [T]$. Figure 8 shows that enabling PGA-1 from the very beginning does not offer much improvement on accuracy. However, as we delay the injection of adversarial noises, the model accuracy starts to increase. By delaying the injection of adversarial noises, we observe improved test accuracy on downstream tasks. However, it also seems that the adversarial noise should not be injected too late, which may inadvertently affect the accuracy. It is possible that a more advanced method to adaptively choose the value of t_s is desired. However, given that (1) the primary focus of this work is to demonstrate that it is possible and effective to accelerate the adaptation of transformer networks via large-batch adaptation and adversarial noises and (2) the search space of is quite small for most downstream tasks, we leave this as an interesting research question for future exploration.

B HYPERPARAMETERS

For all configurations, we fine-tune against the GLUE datasets and set the maximum number of epochs to 6. We use a linear learning rate decay schedule with a warm-up ratio of 0.1. For ScaLA, we set $\lambda = 1$, perturbation clipping radius $\omega = 10^{-5}$, step size $\rho = 10^{-4}$, and $t_s = \{3, 5\}$. These values worked well enough that we did not feel the need to explore more. For fairness, we perform a grid search of learning rates in the range of $\{1e-5, 3e-5, 5e-5, 7e-5, 9e-5, 1e-4, 3e-4\}$ for small batch sizes and $\{5.6e-5, 8e-5, 1e-4, 1.7e-4, 2.4e-4, 2.8e-4, 4e-4, 5.6e-4, 1e-3\}$ for large batch sizes. We keep the remaining hyperparameters unchanged.

C HYPERPARAMETER TUNING COST FOR LARGE-BATCH ADAPTATION WITH ScaLA

In this part, we investigate how large-batch adaptation affects the generalizability of transformer networks on downstream tasks. As there are various heuristics for setting the learning rates (Smith et al., 2018; Goyal et al., 2017; Smith, 2018; You et al., 2019a), and because few work studies the learning rate scaling effects on adapting pre-trained Transformer networks, we perform a grid search on learning rates $\{1e-4, 3e-4, 5e-4, 7e-4, 9e-4, 1e-3, 3e-3\}$ and batch sizes $\{1K, 2K, 4K, 8K\}$ while keeping the other hyperparameters the same to investigate how ScaLA affects the hyperparameter tuning effort.

Table 6 shows the results of using the square root scaling rule to decide the learning rates for large batch sizes vs. accuracy results with tuned learning rate results, without and with ScaLA. The first

row represents the best accuracy found through fine-tuning with a small batch size 32. The next two rows correspond to fine-tuning with batch size 1024 using tuned learning rates vs. using the scaling rule. The last two rows represent fine-tuning using ScaLA with batch size 1024, also using tuned learning rates vs. the scaling rule. Even with square-root scaling, the large-batch baseline still cannot reach the small-batch accuracy (88.7 vs. 89.4). Moreover, although tuning the learning rates lead to better results on some datasets such as MNLI-m (84.9 vs. 85.1) and SST-2 (92.9 vs. 93.5), the square-root scaling rule leads to better results on other tasks such as QNLI (90.8 vs. 90) and QQP (91.4/88.4 vs. 90.9/87.7). So the best learning rates on fine-tuning tasks are not exactly sqrt. However, given that ScaLA with square-root learning rate scaling achieves on average better results than the grid search of learning rates (89.4 vs. 89.7), we suggest to use sqrt scaling for learning rates to simplify the hyperparameter tuning effort for ScaLA.

Table 6: Evaluation results on hyperparameter tuning vs. using square-root learning rate scaling.

| | MNLI-m | QNLI | QQP | SST-2 | Avg |
|-----------------------------------|--------|------|-----------|-------|------|
| Bsz=32 (tuned, baseline) | 84.8 | 90.6 | 91/88 | 93.1 | 89.4 |
| Bsz=1024 (tuned, baseline) | 84.3 | 89.3 | 89.6/86.1 | 93 | 88.5 |
| Bsz=1024 (scaling rule, baseline) | 83.9 | 89.2 | 90.6/87.4 | 92.5 | 88.7 |
| Bsz=1024 (tuned, ScaLA) | 85.1 | 90 | 90.9/87.7 | 93.5 | 89.4 |
| Bsz=1024 (scaling rule, ScaLA) | 84.9 | 90.8 | 91.4/88.4 | 92.9 | 89.7 |

D CONVERGENCE ANALYSIS

In this section, we provide the formal statements and detailed proofs for the convergence rate. The convergence analysis builds on techniques and results in (Davis & Drusvyatskiy, 2018; You et al., 2019b). We consider the general problem of a two-player sequential game represented as nonconvex-nonconcave minimax optimization that is stochastic with respect to the outer (first) player playing $x \in \mathbb{X}$ while sampling ξ from Q and deterministic with respect to the inner (second) player playing $y \in \mathbb{Y}$, i.e.,

$$\min_x \max_y \mathbb{E}_{\xi \sim Q} [f(x, y, \xi)] := \min_x \mathbb{E}_{\xi \sim Q} [g(x, \xi)] \quad (2)$$

Since finding the Stackelberg equilibrium, i.e., the global solution to the saddle point problem, is NP-hard, we consider the optimality notion of a *local minimax* point (Jin et al., 2020). Since maximizing over y may result in a non-smooth function even when f is smooth, the norm of the gradient is not particularly a suitable metric to track the convergence progress of an iterative minimax optimization procedure. Hence, we use the gradient of the *Moreau envelope* (Davis & Drusvyatskiy, 2019) as the appropriate potential function. Let $\mu \in \mathbb{R}_+^h$. The μ -Moreau envelope for a function $g : \mathbb{X} \rightarrow \mathbb{R}$ is defined as $g_\mu(x) := \min_z g(z) + \sum_{i=1}^h \frac{1}{2\mu^i} \|x^i - z^i\|^2$. Another reason for the choice of this potential function is due to the special property (Rockafellar, 2015) of the Moreau envelope that if its gradient $\nabla_x [g_\mu(x)]$ almost vanishes at x , such x is close to a stationary point of the original function g .

Assumptions: We assume that $\mathbb{X} = \bigsqcup_{i=1}^h \mathbb{X}^i$ is partitioned into h disjoint groups, i.e., in terms of training a neural network, we can think of the network having the parameters partitioned into h (hidden) layers. The measure Q characterizes the training data. Let $\widehat{\nabla}_x f(x, y)$ denote the noisy estimate of the true gradient $\nabla_x f(x, y)$. We assume that the noisy gradients are unbiased, i.e., $\mathbb{E}[\widehat{\nabla}_x f(x, y)] = \nabla_x f(x, y)$. For each group $i \in [h]$, we make the standard (groupwise) boundedness assumption (Ghadimi & Lan, 2013) on the variance of the stochastic gradients, i.e., $\mathbb{E}\|\widehat{\nabla}_x^i f(x, y) - \nabla_x^i f(x, y)\|^2 \leq \sigma_i^2, \forall i \in [h]$. We assume that $f(x, y)$ has Lipschitz continuous gradients. Specifically, let $f(x, y)$ be α -smooth in x where $\alpha := (\alpha_1, \dots, \alpha_h)$ denotes the h -dimensional vector of (groupwise) Lipschitz parameters, i.e., $\|\nabla_x^i f(x_a, y) - \nabla_x^i f(x_b, y)\| \leq \alpha_i \|x_a^i - x_b^i\|, \forall i \in [h]$ and $x_a, x_b \in \mathbb{X}, y \in \mathbb{Y}$. Let $\kappa_\alpha := \frac{\max_i \alpha_i}{\min_i \alpha_i}$.

Super-scripts are used to index into a vector (i denotes the group index and j denotes an element in group i). For any $c \in \mathbb{R}$, the function $\nu : \mathbb{R} \rightarrow [\mathcal{L}, \mathcal{U}]$ clips its values, i.e., $\nu(c) := \max(\mathcal{L}, \min(c, \mathcal{U}))$ where $\mathcal{L} < \mathcal{U}$. Let $\|\cdot\|, \|\cdot\|_1$ and $\|\cdot\|_\infty$ denote the ℓ_2, ℓ_1 , and ℓ_∞ norms. We assume that the true gradients are bounded, i.e., $\|\nabla_x f(x, y)\|_\infty \leq \mathcal{G}$.

First, we begin with relevant supporting lemmas. The following lemma characterizes the convexity of an additive modification of g .

Lemma D.1 ((Lin et al., 2020; Jin et al., 2020; Rafique et al., 2021)). *Let $g(x) := \max_y f(x, y)$ with f being α -smooth in x where $\alpha \in \mathbb{R}_+^h$ is the vector of groupwise Lipschitz parameters. Then, $g(x) + \sum_{i=1}^h \frac{\alpha_i}{2} \|x^i\|^2$ is convex in x .*

The following property of the Moreau envelope relates it to the original function.

Lemma D.2 (Rockafellar, 2015). *Let g be defined as in Lemma D.1. Let $\widehat{x} = \arg \min_{\widehat{x}} g(\widehat{x}) + \sum_{i=1}^h \frac{1}{2\mu^i} \|\widehat{x}^i - x^i\|^2$. Then, $\|g_\mu(x)\| \leq \epsilon$ implies $\|\widehat{x} - x\| \leq \|\mu\|_\infty \epsilon$ and $\min_h \|h\| \leq \epsilon$ with $h \in \partial g$ where ∂g denotes the subdifferential of g .*

We now present the formal version of Theorem 4.1 in Theorem D.3. Note that Lemma D.2 facilitates giving the convergence guarantees in terms of the gradient of the Moreau envelope. Recall that

$t \in [T]$ denotes the epochs corresponding to the outer maximization. Without loss of generality, we set the delay parameter for injection of the adversarial perturbation in Algorithm 1 as $t_s = 0$. Here, we assume that the PGA provides an ϵ -approximate maximizer.

Theorem D.3 (Groupwise outer minimization with an ϵ -approximate inner maximization oracle). *Let us define relevant constants as $\mathcal{D} := (g_{1/2\alpha}(x_0) - \mathbb{E}(\min_x g(x)))$ being the optimality gap due to initialization, $\kappa_\alpha := \frac{\max_i \alpha_i}{\min_i \alpha_i}$ being the condition number, $\|\nabla_x f(x, y)\|_\infty \leq \mathcal{G}$ being gradient bound, $\mathcal{Z} := \max_{i,j,t} \frac{(\hat{x}_t^{i,j} - x_t^{i,j})}{(\nabla_t^{i,j})} \sigma_i$ being the variance term, \mathcal{L}, \mathcal{U} being clipping constants such that $\mathcal{L} \leq \mathcal{U}$. For the outer optimization, setting the learning rate as $\eta = \frac{1}{\mathcal{U}\sqrt{T}}$ and scaling batch size as $b = \frac{16T\mathcal{L}^2\mathcal{Z}^2}{\mathcal{U}^2}$, we have*

$$\mathbb{E} [\|\nabla g_{1/2\alpha}(\bar{x})\|^2] \leq 4\epsilon\|\alpha\|_\infty + \frac{2\kappa_\alpha \mathcal{D}\mathcal{G}}{\sqrt{T}} \quad (3)$$

where \bar{x} is the estimator obtained from running T steps of Algorithm 1 and picking x_t uniformly at random for $t \in [T]$.

Proof. In this proof, for brevity, we define the vector $\nabla_t := \nabla_x f(x, y)$, i.e., the gradient of the objective with respect to x , evaluated at the outer step t . Since evaluating gradients using mini-batches produces noisy gradients, we use $\widehat{\nabla}$ to denote the noisy version of a true gradient ∇ , i.e., $\widehat{\nabla} = \nabla + \Delta$ for a noise vector Δ . For any outer step t , we have $f(x_t, \widehat{y}) \geq g(x_t) - \epsilon$ where \widehat{y} is an ϵ -approximate maximizer. For any $\tilde{x} \in \mathbb{X}$, using the smoothness property (Lipschitz gradient) of f , we have

$$\begin{aligned} g(\tilde{x}) &\geq f(\tilde{x}, y_t) \\ &\geq f(x_t, y_t) + \sum_{i=1}^h \langle \nabla_t^i, \tilde{x}^i - x_t^i \rangle - \sum_{i=1}^h \frac{\alpha_i}{2} \|\tilde{x}^i - x_t^i\|^2 \\ &\geq g(x_t) - \epsilon + \sum_{i=1}^h \langle \nabla_t^i, \tilde{x}^i - x_t^i \rangle - \sum_{i=1}^h \frac{\alpha_i}{2} \|\tilde{x}^i - x_t^i\|^2 \end{aligned} \quad (4)$$

Let $\phi_\mu(x, z) := g(z) + \sum_{i=1}^h \frac{1}{2\mu^i} \|x^i - z^i\|^2$. Recall that the μ -Moreau envelope for g is defined as $g_\mu(x) := \min_z \phi_\mu(x, z)$ and its gradient is the groupwise proximal operator given by $\nabla_x [g_\mu(x)] = \left[\frac{1}{\mu^1} (x^1 - \arg \min_{z^1} \phi_\mu(x, z)), \dots, \frac{1}{\mu^h} (x^h - \arg \min_{z^h} \phi_\mu(x, z)) \right]$.

Now, let $\widehat{x}_t = \arg \min_x \phi_{1/2\alpha}(x_t, x) = \arg \min_x \left(g(x) + \sum_{i=1}^h \alpha_i \|x_t^i - x^i\|^2 \right)$. Then, plugging in the update rule for x at step $t+1$ in terms of quantities at step t , using the shorthand $\nu_t^i := \nu(\|x_t^i\|)$

and conditioning on the filtration up to time t , we have

$$\begin{aligned}
g_{1/2\alpha}(x_{t+1}) &\leq g(\hat{x}_t) + \sum_{i=1}^h \alpha_i \|x_{t+1}^i - \hat{x}_t^i\|^2 \\
&\leq g(\hat{x}_t) + \sum_{i=1}^h \alpha_i \left\| x_t^i - \eta_t \nu_t^i \frac{\widehat{\nabla}_t^i}{\|\widehat{\nabla}_t^i\|} - \hat{x}_t^i \right\|^2 \\
&\leq g(\hat{x}_t) + \sum_{i=1}^h \alpha_i \|x_t^i - \hat{x}_t^i\|^2 + \sum_{i=1}^h 2\alpha_i \eta_t \left\langle \nu_t^i \frac{\widehat{\nabla}_t^i}{\|\widehat{\nabla}_t^i\|}, \hat{x}_t^i - x_t^i \right\rangle + \sum_{i=1}^h \alpha_i \eta_t^2 (\nu_t^i)^2 \\
&\leq g_{1/2\alpha}(x_t) + \sum_{i=1}^h 2\alpha_i \eta_t \left\langle \nu_t^i \frac{\widehat{\nabla}_t^i}{\|\widehat{\nabla}_t^i\|}, \hat{x}_t^i - x_t^i \right\rangle + \sum_{i=1}^h \alpha_i \eta_t^2 (\nu_t^i)^2 \\
&\leq g_{1/2\alpha}(x_t) + 2\eta_t \sum_{i=1}^h \alpha_i \nu_t^i \sum_{j=1}^{d_i} \left(\frac{\widehat{\nabla}_t^{i,j}}{\|\widehat{\nabla}_t^i\|} - \frac{\nabla_t^{i,j}}{\|\nabla_t^i\|} + \frac{\nabla_t^{i,j}}{\|\nabla_t^i\|} \right) \times (\hat{x}_t^{i,j} - x_t^{i,j}) + \sum_{i=1}^h \alpha_i \eta_t^2 (\nu_t^i)^2 \\
&\leq g_{1/2\alpha}(x_t) + 2\eta_t \sum_{i=1}^h \alpha_i \nu_t^i \sum_{j=1}^{d_i} \left(\frac{\nabla_t^{i,j}}{\|\nabla_t^i\|} \right) \times (\hat{x}_t^{i,j} - x_t^{i,j}) \\
&\quad + 2\eta_t \sum_{i=1}^h \alpha_i \nu_t^i \sum_{j=1}^{d_i} \left(\frac{\widehat{\nabla}_t^{i,j}}{\|\widehat{\nabla}_t^i\|} - \frac{\nabla_t^{i,j}}{\|\nabla_t^i\|} \right) \times (\hat{x}_t^{i,j} - x_t^{i,j}) + \sum_{i=1}^h \alpha_i \eta_t^2 (\nu_t^i)^2 \\
&\leq g_{1/2\alpha}(x_t) + 2\eta_t \sum_{i=1}^h \frac{\alpha_i \nu_t^i}{\|\nabla_t^i\|} \langle \nabla_t^i, \hat{x}_t^i - x_t^i \rangle \\
&\quad + 2\eta_t \sum_{i=1}^h \alpha_i \nu_t^i \sum_{j=1}^{d_i} \left(\frac{\nabla_t^{i,j} + \Delta_t^{i,j}}{\|\nabla_t^i + \Delta_t^i\|} - \frac{\nabla_t^{i,j}}{\|\nabla_t^i\|} \right) \times (\hat{x}_t^{i,j} - x_t^{i,j}) + \sum_{i=1}^h \alpha_i \eta_t^2 (\nu_t^i)^2
\end{aligned}$$

$$\begin{aligned}
&\leq g_{1/2\alpha}(x_t) + 2\eta_t \mathcal{U} \sum_{i=1}^h \frac{\alpha_i}{\|\nabla_t^i\|} \langle \nabla_t^i, \widehat{x}_t^i - x_t^i \rangle \\
&+ 2\eta_t \sum_{i=1}^h \alpha_i \nu_t^i \sum_{j=1}^{d_i} \left(\frac{\|\nabla_t^i\| (\nabla_t^{i,j}) (\nabla_t^{i,j} + \Delta_t^{i,j}) - \|\nabla_t^i + \Delta_t^i\| (\nabla_t^{i,j})^2}{\|\nabla_t^i + \Delta_t^i\| \|\nabla_t^i\|} \right) \times \frac{(\widehat{x}_t^{i,j} - x_t^{i,j})}{(\nabla_t^{i,j})} \\
&+ \sum_{i=1}^h \alpha_i \eta_t^2 (\nu_t^i)^2 \\
&\stackrel{E_1}{\leq} g_{1/2\alpha}(x_t) + 2\eta_t \mathcal{U} \max_i \frac{\alpha_i}{\|\nabla_t^i\|} \left(g(\widehat{x}_t) - g(x_t) + \epsilon + \sum_{i=1}^h \frac{\alpha_i}{2} \|\widehat{x}^i - x_t^i\|^2 \right) \\
&+ 2\eta_t \sum_{i=1}^h \alpha_i \nu_t^i \max_j \frac{(\widehat{x}_t^{i,j} - x_t^{i,j})}{(\nabla_t^{i,j})} \left(\frac{\langle \nabla_t^i, \nabla_t^i + \Delta_t^i \rangle - \|\nabla_t^i + \Delta_t^i\| \|\nabla_t^i\|}{\|\nabla_t^i + \Delta_t^i\|} \right) + \sum_{i=1}^h \alpha_i \eta_t^2 (\nu_t^i)^2 \\
&\leq g_{1/2\alpha}(x_t) + 2\eta_t \mathcal{U} \max_i \frac{\alpha_i}{\|\nabla_t^i\|} \left(g(\widehat{x}_t) - g(x_t) + \epsilon + \sum_{i=1}^h \frac{\alpha_i}{2} \|\widehat{x}^i - x_t^i\|^2 \right) \\
&- 2\eta_t \sum_{i=1}^h \alpha_i \nu_t^i \max_j \frac{(\widehat{x}_t^{i,j} - x_t^{i,j})}{(\nabla_t^{i,j})} \left(\frac{\|\nabla_t^i + \Delta_t^i\| \|\nabla_t^i\| - \|\nabla_t^i + \Delta_t^i\|^2 + \langle \Delta_t^i, \nabla_t^i + \Delta_t^i \rangle}{\|\nabla_t^i + \Delta_t^i\|} \right) \\
&+ \sum_{i=1}^h \alpha_i \eta_t^2 (\nu_t^i)^2 \tag{5} \\
&\leq g_{1/2\alpha}(x_t) + 2\eta_t \mathcal{U} \max_i \frac{\alpha_i}{\|\nabla_t^i\|} \left(g(\widehat{x}_t) - g(x_t) + \epsilon + \sum_{i=1}^h \frac{\alpha_i}{2} \|\widehat{x}^i - x_t^i\|^2 \right) \\
&- 2\eta_t \sum_{i=1}^h \alpha_i \nu_t^i \max_j \frac{(\widehat{x}_t^{i,j} - x_t^{i,j})}{(\nabla_t^{i,j})} \left(\|\nabla_t^i\| - \|\nabla_t^i + \Delta_t^i\| - \frac{|\langle \Delta_t^i, \nabla_t^i + \Delta_t^i \rangle|}{\|\nabla_t^i + \Delta_t^i\|} \right) + \sum_{i=1}^h \alpha_i \eta_t^2 (\nu_t^i)^2 \\
&\stackrel{E_2}{\leq} g_{1/2\alpha}(x_t) + 2\eta_t \mathcal{U} \max_i \frac{\alpha_i}{\|\nabla_t^i\|} \left(g(\widehat{x}_t) - g(x_t) + \epsilon + \sum_{i=1}^h \frac{\alpha_i}{2} \|\widehat{x}^i - x_t^i\|^2 \right) \\
&- 2\eta_t \sum_{i=1}^h \alpha_i \nu_t^i \max_j \frac{(\widehat{x}_t^{i,j} - x_t^{i,j})}{(\nabla_t^{i,j})} (\|\nabla_t^i\| - \|\nabla_t^i + \Delta_t^i\| - \|\Delta_t^i\|) + \sum_{i=1}^h \alpha_i \eta_t^2 (\nu_t^i)^2 \\
&\stackrel{E_3}{\leq} g_{1/2\alpha}(x_t) + 2\eta_t \mathcal{U} \max_i \frac{\alpha_i}{\|\nabla_t^i\|} \left(g(\widehat{x}_t) - g(x_t) + \epsilon + \sum_{i=1}^h \frac{\alpha_i}{2} \|\widehat{x}^i - x_t^i\|^2 \right) \\
&- 4\eta_t \sum_{i=1}^h \alpha_i \nu_t^i \max_j \frac{(\widehat{x}_t^{i,j} - x_t^{i,j})}{(\nabla_t^{i,j})} \|\Delta_t^i\| + \sum_{i=1}^h \alpha_i \eta_t^2 (\nu_t^i)^2 \\
g_{1/2\alpha}(x_T) &\stackrel{E_4}{\leq} g_{1/2\alpha}(x_0) + 2\mathcal{U} \sum_{t=0}^{T-1} \eta_t \max_i \frac{\alpha_i}{\|\nabla_t^i\|} \left(g(\widehat{x}_t) - g(x_t) + \epsilon + \sum_{i=1}^h \frac{\alpha_i}{2} \|\widehat{x}^i - x_t^i\|^2 \right) \\
&- 4 \sum_{t=0}^{T-1} \eta_t \sum_{i=1}^h \alpha_i \nu_t^i \max_j \frac{(\widehat{x}_t^{i,j} - x_t^{i,j})}{(\nabla_t^{i,j})} \|\Delta_t^i\| + \sum_{t=0}^{T-1} \sum_{i=1}^h \alpha_i \eta_t^2 (\nu_t^i)^2
\end{aligned}$$

where we have used Hölder's inequality along with bound equation 4 in E_1 , Cauchy-Schwarz inequality in E_2 , triangle inequality in E_3 , telescoping sum in E_4 . Rearranging and using $\eta_t = \eta$ in

E_5 along with Hölder's inequality,

$$\begin{aligned} \frac{1}{2\eta\mathcal{U}} (g_{1/2\alpha}(x_T) - g_{1/2\alpha}(x_0)) &\leq \sum_{t=0}^{T-1} \max_i \frac{\alpha_i}{\|\nabla_t^i\|} \left(g(\hat{x}_t) - g(x_t) + \epsilon + \sum_{i=1}^h \frac{\alpha_i}{2} \|\hat{x}^i - x_t^i\|^2 \right) \\ &\quad - \frac{2}{\mathcal{U}} \sum_{t=0}^{T-1} \sum_{i=1}^h \alpha_i \nu_t^i \max_j \frac{(\hat{x}_t^{i,j} - x_t^{i,j})}{(\nabla_t^{i,j})} \|\Delta_t^i\| + \frac{\eta}{2\mathcal{U}} \sum_{t=0}^{T-1} \sum_{i=1}^h \alpha_i (\nu_t^i)^2 \\ \frac{1}{2\eta\mathcal{U}} (g_{1/2\alpha}(x_T) - g_{1/2\alpha}(x_0)) &\stackrel{E_5}{\leq} \max_{i,t} \frac{\alpha_i}{\|\nabla_t^i\|} \sum_{t=0}^{T-1} \left(g(\hat{x}_t) - g(x_t) + \epsilon + \sum_{i=1}^h \frac{\alpha_i}{2} \|\hat{x}^i - x_t^i\|^2 \right) \\ &\quad - \frac{2}{\mathcal{U}} \sum_{t=0}^{T-1} \sum_{i=1}^h \alpha_i \nu_t^i \max_j \frac{(\hat{x}_t^{i,j} - x_t^{i,j})}{(\nabla_t^{i,j})} \|\Delta_t^i\| + \frac{\eta}{2\mathcal{U}} \sum_{t=0}^{T-1} \sum_{i=1}^h \alpha_i (\nu_t^i)^2 \end{aligned}$$

Dividing by T and rearranging,

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} \left(g(x_t) - g(\hat{x}_t) - \sum_{i=1}^h \frac{\alpha_i}{2} \|\hat{x}^i - x_t^i\|^2 \right) &\leq \epsilon - \frac{1}{2\eta\mathcal{U}\zeta T} (g_{1/2\alpha}(x_T) - g_{1/2\alpha}(x_0)) \\ &\quad - \frac{2}{\mathcal{U}\zeta T} \sum_{t=0}^{T-1} \sum_{i=1}^h \alpha_i \nu_t^i \max_j \frac{(\hat{x}_t^{i,j} - x_t^{i,j})}{(\nabla_t^{i,j})} \|\Delta_t^i\| \\ &\quad + \frac{\eta}{2\mathcal{U}\zeta T} \sum_{i=1}^h \alpha_i \sum_{t=0}^{T-1} (\nu_t^i)^2 \end{aligned}$$

where we define $\zeta := \max_{i,t} \frac{\alpha_i}{\|\nabla_t^i\|}$. Defining $\mathcal{D} := (g_{1/2\alpha}(x_0) - \mathbb{E}(\min_x g(x)))$ and taking expectation with respect to ξ on both sides, we have

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \left(g(x_t) - g(\hat{x}_t) - \sum_{i=1}^h \frac{\alpha_i}{2} \|\hat{x}^i - x_t^i\|^2 \right) &\leq \epsilon + \frac{\mathcal{D}}{2\eta\mathcal{U}\zeta T} \\ &\quad - \frac{2\mathcal{L}}{\mathcal{U}\zeta T} \sum_{t=0}^{T-1} \sum_{i=1}^h \alpha_i \max_j \frac{(\hat{x}_t^{i,j} - x_t^{i,j})}{(\nabla_t^{i,j})} \mathbb{E} \|\Delta_t^i\| + \frac{\eta\mathcal{U}\|\alpha\|_1}{2\zeta} \\ &\stackrel{E_6}{\leq} \epsilon + \frac{\mathcal{D}}{2\eta\mathcal{U}\zeta T} \\ &\quad - \frac{2\mathcal{L}}{\mathcal{U}\zeta T} \sum_{t=0}^{T-1} \sum_{i=1}^h \alpha_i \max_j \frac{(\hat{x}_t^{i,j} - x_t^{i,j})}{(\nabla_t^{i,j})} \frac{\sigma_i}{\sqrt{b}} + \frac{\eta\mathcal{U}\|\alpha\|_1}{2\zeta} \\ &\stackrel{E_7}{\leq} \epsilon + \frac{\mathcal{D}}{2\eta\mathcal{U}\zeta T} \\ &\quad - \frac{2\mathcal{L}\|\alpha\|_1}{\mathcal{U}\zeta\sqrt{b}} \max_{i,j,t} \frac{(\hat{x}_t^{i,j} - x_t^{i,j})}{(\nabla_t^{i,j})} \sigma_i + \frac{\eta\mathcal{U}\|\alpha\|_1}{2\zeta} \\ &\stackrel{E_8}{=} \epsilon + \frac{\mathcal{D}}{2\eta\mathcal{U}\zeta T} - \frac{2\mathcal{L}\|\alpha\|_1\mathcal{Z}}{\mathcal{U}\zeta\sqrt{b}} + \frac{\eta\mathcal{U}\|\alpha\|_1}{2\zeta} \quad (6) \end{aligned}$$

where we have used the assumption on the variance of stochastic gradients in E_6 , Hölder's inequality in E_7 and we define $\mathcal{Z} := \max_{i,j,t} \frac{(\hat{x}_t^{i,j} - x_t^{i,j})}{(\nabla_t^{i,j})} \sigma_i$ in E_8 ; b denotes batch size. Now, we lower bound

the left hand side using the convexity of the additive modification of g .

$$\begin{aligned}
& g(x_t) - g(\hat{x}_t) - \sum_{i=1}^h \frac{\alpha_i}{2} \|\hat{x}^i - x_t^i\|^2 \\
& \geq g(x_t) + 0 - g(\hat{x}_t) - \sum_{i=1}^h \alpha_i \|\hat{x}^i - x_t^i\|^2 + \sum_{i=1}^h \frac{\alpha_i}{2} \|\hat{x}^i - x_t^i\|^2 \\
& \geq \left(\left(g(x_t) + \sum_{i=1}^h \alpha_i \|x_t^i - x_t^i\|^2 \right) - \min_x \left(g(x_t) + \sum_{i=1}^h \alpha_i \|x^i - x_t^i\|^2 \right) \right) + \sum_{i=1}^h \frac{\alpha_i}{2} \|\hat{x}^i - x_t^i\|^2 \\
& \geq \sum_{i=1}^h \frac{\alpha_i}{2} \|\hat{x}^i - x_t^i\|^2 + \sum_{i=1}^h \frac{\alpha_i}{2} \|\hat{x}^i - x_t^i\|^2 = \sum_{i=1}^h \frac{4\alpha_i^2}{4\alpha_i} \|\hat{x}^i - x_t^i\|^2 \\
& \stackrel{E_9}{\geq} \frac{1}{4 \max_i \alpha_i} \|\nabla g_{1/2\alpha}(x_t)\|^2 \tag{7}
\end{aligned}$$

where we have used the expression for the gradient of the Moreau envelope in E_9 . Combining the inequalities from Equation equation 7 and Equation equation 6, we have

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \left(\frac{1}{4 \max_i \alpha_i} \|\nabla g_{1/2\alpha}(x_t)\|^2 \right) \leq \epsilon + \frac{\mathcal{D}}{2\eta\mathcal{M}\zeta T} + \left(\frac{\eta\mathcal{U}}{2\zeta} - \frac{2\mathcal{L}\mathcal{Z}}{\mathcal{U}\zeta\sqrt{b}} \right) \|\alpha\|_1$$

Setting the learning rate as $\eta = \frac{1}{\mathcal{U}\sqrt{T}}$ and batch size as $b = \frac{16T\mathcal{L}^2\mathcal{Z}^2}{\mathcal{U}^2}$,

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} [\|\nabla g_{1/2\alpha}(x_t)\|^2] \leq 4\epsilon \max_i \alpha_i + \frac{2\mathcal{D} \max_i \alpha_i}{\zeta\sqrt{T}}$$

Now, to simplify ζ , using the inequality that $\max_k (a_k \cdot b_k) \geq \min_{k_a} a_{k_a} \cdot \min_{k_b} b_{k_b}$ for two finite sequences $\{a, b\}$ with positive values, along with the bounded gradients assumption, we have

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} [\|\nabla g_{1/2\alpha}(x_t)\|^2] \leq 4\epsilon \max_i \alpha_i + \frac{2\mathcal{D}\mathcal{G} \max_i \alpha_i}{\sqrt{T} \min_i \alpha_i} = 4\epsilon \|\alpha\|_\infty + \frac{2\kappa_\alpha \mathcal{D}\mathcal{G}}{\sqrt{T}}$$

where $\kappa_\alpha := \frac{\max_i \alpha_i}{\min_i \alpha_i}$. \square

In analyzing inexact version, as in Theorem D.3, we assumed the availability of an adversarial oracle. Next, we open up the adversarial oracle to characterize the oracle-free complexity. In order to do this, we will assume additional properties about the function f as well as our deterministic perturbation space, $\mathbb{Y}_t \subseteq \mathbb{Y}$, $\forall t \in [T]$. Note that, for any given t , $y_\tau \in \mathbb{Y}_t$, $\forall \tau \in \mathcal{T}$. We recall the following guarantee for generalized non-convex projected gradient ascent.

Lemma D.4 ((Jain & Kar, 2017)). *For every t , Let $f(x_t, \cdot)$ satisfy restricted strong convexity with parameter \mathcal{C} and restricted strong smoothness with parameter \mathcal{S} over a non-convex constraint set with $\mathcal{S}/\mathcal{C} < 2$, ie, $\frac{\mathcal{C}}{2} \|z - y\|^2 \leq f(x_t, y) - f(x_t, z) - \langle \nabla_z f(x_t, z), y - z \rangle \leq \frac{\mathcal{S}}{2} \|z - y\|^2$ for $y, z \in \mathbb{Y}_t$. For any given t , let the PGA- \mathcal{T} algorithm $y_\tau \leftarrow \Pi_\epsilon [y_{\tau-1} + \rho \nabla_y f(x_t, y)]$ be executed with step size $\rho = 1/\mathcal{S}$. Then after at most $\mathcal{T} = O\left(\frac{\mathcal{C}}{2\mathcal{C}-\mathcal{S}} \log \frac{1}{\epsilon}\right)$ steps, $f(x_t, y_\mathcal{T}) \geq \max_y f(x_t, y) - \epsilon$.*

Using Theorem D.3 and Lemma D.4 (together with the additional restricted strong convexity/smoothness assumptions), we have the following theorem on the full oracle-free rates for Algorithm 1.

Theorem D.5 (Groupwise outer minimization with inner maximization using projected gradient ascent). *Setting the inner iteration count as $\mathcal{T} = O\left(\frac{\mathcal{C}}{2\mathcal{C}-\mathcal{S}} \log \frac{8\|\alpha\|_\infty}{\epsilon}\right)$ and the outer iteration count as $T = \frac{16\kappa_\alpha \mathcal{D}^2 \mathcal{G}^2}{\epsilon^2}$, for a combined total of $O\left(\frac{1}{\epsilon^2} \log \frac{1}{\epsilon}\right)$ adaptive adversarial iterations, Algorithm 1 achieves $\mathbb{E} [\|\nabla g_{1/2\alpha}(\bar{x})\|^2] \leq \epsilon$.*