Large Learning Rates without the Agonizing PAIN: DISPELLING THE CURSE OF SINGULARITIES IN DEEP NEURAL NETWORKS

Anonymous authors

Paper under double-blind review

ABSTRACT

Employing large learning rates (LRs) in deep learning can accelerate convergence and improve generalization, but it can also cause training instability and loss explosion: determining an appropriate LR is an often laborious and painful art. Our study into the fine-grained behaviors of parametric singularities, specifically the stable ranks of weight matrices of network components, reveals a strong connection between these singularities and training instability. As training progresses, parametric singularities trend upward, a phenomenon that is directly aggravated by large LRs. Crucially, several training steps before prominent instabilities such as gradient explosions, we observe unusually high parametric singularities across the network components, leading to rank-deficient representations. These representations, in turn, amplify parametric singularities during backpropagation, creating a vicious cycle that eventually results in loss explosions. We refer to this phenomenon as the curse of singularities. Building on this understanding, we propose a lightweight and robust stabilization method called Parametric Singularity Smoothing (PSS), which allows for early intervention and mitigates impending instability by smoothing the singular spectra of weight matrices, thereby preventing the curse of singularities. This approach is easy to implement, works at any stage of training by restoring stable training even after instability, has neglectable computational overhead, and, most importantly, frees us from the painful LR finetunings to avoid instabilities. Experimental results across various datasets, networks, and optimizers demonstrate that our approach allows a 5-10 \times increase in LR without producing instability, attaining better training efficiency and generalization. We release our code for everyone to use our methods and reproduce the experiments, available at https://anonymous.4open.science/r/ ICLR_stability-C69C.

036

038

006

008 009 010

011

013

014

015

016

017

018

019

021

023

024

025

026

027

028

029

031

032

034

1 INTRODUCTION

Applying large learning rates (LRs) in deep learning is widely known to have benefits, e.g., escaping 040 local minima (Li et al., 2019; Jastrzebski et al., 2020), accelerated convergence (Smith, 2017; Smith 041 & Topin, 2019), and better generalization (Lu et al., 2023). However, large LRs also increase the 042 probability of training instability (Lewkowycz et al., 2020; Cohen et al., 2021; Wang & Roberts, 043 2023), which often leads to loss explosions that waste computational resources, a serious issue 044 in large models. Although recent work has investigated the properties of large LRs from various aspects, there still lacks an effective solution to address the instabilities that are more prominent with large LRs. Nowadays, fine-tuning LRs for stabilities remains commonly indispensable and 046 can be painful, especially for real-world foundation models. For instance, it takes multiple runs 047 to determine a proper LR free of instabilities when training OPT-175B (Zhang et al., 2022), each 048 consuming dozens of days on thousands of GPUs. Our goal is to enable a reliable, stable, automated, restart-free approach to using large LRs. 050

We start by investigating the behaviors of the stable ranks (Rudelson & Vershynin, 2007), i.e., parametric singularities, of weight matrices across model components, such as the key, query, and value matrices of self-attention heads in transformers (Dong et al., 2021). Recent studies (Mousavi-Hosseini et al., 2022; Zangrando et al., 2024; Feng et al., 2022; Huh et al., 2021; Arora et al., 2019)

have revealed the low-rank properties of parametric matrices in different neural networks. These
matrices tend to be increasingly "singular" as training progresses. Such singularities are usually
considered to be linked with feature spaces of decreased rankings (Dong et al., 2021; Noci et al.,
2022); pronounced parametric singularities generally imply that a few singular vectors dominate the
feature space, leading to rank-deficient representations with poor expressiveness (Dong et al., 2021;
Noci et al., 2022). Previous work did not discuss or investigate the relationship between singularities
and training instabilities; we are the first to study how the network's fine-grained, dynamic singularity behaviors cause loss explosions that prevent further learning, which happens quite commonly
for large LRs.

063 Specifically, we investigate the parametric singularity distributions across network components of 064 training sessions with exploding losses. We find that in practice parametric singularities increase monotonically as training proceeds, with deep layers having higher singularities than shallow lay-065 ers. Notably, growing singularity trends can be observed across varying LR values and are more 066 pronounced with large LRs. Inspecting the few training steps preceding loss explosions reveals a 067 particular step in which growing singularities in the parametric space yield rank-deficient represen-068 tations with unusually high similarity in the feature spaces. Specifically, a few dominant singular 069 values in the parametric space overshadow the contributions of other singular directions, particularly 070 in deep layers, leading to a sharp reduction in representational diversity. Updating such networks 071 will further increase the parametric singularities, leading to yet more severe rank deficiencies in the 072 next training step. Within a few steps, this vicious cycle produces loss explosion and puts an end 073 to productive learning. We call this phenomenon the curse of singularities, and argue that it is a 074 primary cause of training instability.

075 Based on these findings, we propose a computationally efficient and robust method called Parametric 076 Singularity Smoothing (PSS) for using large LRs without training instability. PSS enables early 077 detection of instability and restores stability by smoothing the singular spectrums of each weight 078 matrix across the network components, raising stable ranks and avoiding the curse of singularities. 079 This method is simple, easy to implement, and does not require restarts or manual LR tuning when 080 countering training instabilities. It enables the use of large LRs without producing loss explosions, 081 and it can even restore stability after instability has occurred. We evaluated this method using various 082 datasets and networks and found that it increases usable LR by $5-10\times$, providing potentially better 083 training efficiency and generalized performance.

- 084 This work's contributions follow:
 - It is the first to describe the dynamic patterns of network singularities and reveal their tight associations with loss explosion, a phenomenon we call *the curse of singualrites*.
 - It describes a lightweight, easy to implement and robust method enabling large LRs without loss explosion that does not depend on LR tuning or training restarts.
 - The PSS method significantly increases the maximum stable LR by 5-10×, providing potentially better training efficiency and generalized performance.
- 092 093 094

096

097

090

086

2 ANALYSIS

2.1 STABLE RANK AND STABLE JACOBIAN ENERGY

Define a feed-forward neural network as $f(\mathbf{x}; \boldsymbol{\theta}) = f_L(f_{L-1}(\dots f_1(\mathbf{x}; \boldsymbol{\mathcal{W}}^1); \dots; \boldsymbol{\mathcal{W}}^{L-1}); \boldsymbol{\mathcal{W}}^L)$ where \mathbf{x} and $\boldsymbol{\theta}$ denote the input and learnable parameters, respectively, L is the number of layers, f_l and $\boldsymbol{\mathcal{W}}^l$ denote the function and the collection of weight matrices of the *l*-th layer. Let $\boldsymbol{\mathcal{W}}^l = \{\mathbf{W}_i^l\}_{i=1}^{k^l}$ where \mathbf{W}_i^l denotes the *i*-th weight matrix of the *l*-th layer and k^l denotes the total number of weight matrices in the *l*-th layer. We use the stable rank (Rudelson & Vershynin, 2007) to represent the degree of singularity of a weight matrix \mathbf{W} .

Formally, for a matrix $\mathbf{W} \in \mathbb{R}^{m \times n}$ with rank $r \leq \min(m, n)$, we perform the Singular Value Decomposition (SVD) to obtain its singular values $\{\sigma_i\}_{i=1}^{\min(m,n)}$, left singular vectors $\{\mathbf{u}_i\}_{i=1}^m$ where $\mathbf{u}_i \in \mathbb{R}^m$, and right singular vectors $\{\mathbf{v}_i\}_{i=1}^n$ where $\mathbf{v}_i \in \mathbb{R}^n$, such that $\mathbf{W} = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^\top$. Throughout the paper, we assume all singular values are sorted in descending order, namely $\sigma_1 \geq$ 108 $\sigma_2 \ge \dots \sigma_r > 0$, and we denote the *i*-th singular value of **W** as $\sigma_i(\mathbf{W})$. The stable rank (SR) is defined as follows.

Definition 1. (Stable Rank). The stable rank of a matrix \mathbf{W} is defined as the squared ratios of its Frobenius norm $\|\mathbf{W}\|_F^2$ to its matrix 2-norm $\|\mathbf{W}\|_2^2$, which can be written as

$$\operatorname{SR}(\mathbf{W}) = \frac{\|\mathbf{W}\|_F^2}{\|\mathbf{W}\|_2^2} = \frac{\sum_{i=1}^r \sigma_i^2(\mathbf{W})}{\sigma_1^2(\mathbf{W})}.$$
(1)

SR is a "softer" version of the matrix rank and possesses several desirable properties. It better numerically reflects the singular degree of a weight matrix and while its insensitivity to matrix scales and small perturbations enables use across varying parametric spaces. We therefore use it as a measure of the singularities of weight matrices.

A small SR value intuitively indicates the emergence of pronounced singularities, where a few singular values dominate the others in the weight matrix. These overly dominant singular values can lead to rank-deficient representations, which in turn produce Jacobian matrices that, upon updates, further exacerbate the weight singularities and reduce SR values even more.

We define a new and intuitive descriptor to quantify how the gradient updates will change the parametric singularities.

Formally, for a matrix $\mathbf{W} \in \mathbb{R}^{m \times n}$, let $\mathbf{J}_{\mathbf{W}}(\mathbf{x}) = \frac{\partial \mathcal{L}(\mathbf{x})}{\partial \mathbf{W}} \in \mathbb{R}^{m \times n}$ be the Jacobian matrix of the loss function \mathcal{L} w.r.t. \mathbf{W} computed on \mathbf{x} , and we will write $\mathbf{J}_{\mathbf{W}}$ for simplicity. For a singular value $\sigma_i(\mathbf{W})$ with corresponding singular vectors \mathbf{u}_i and \mathbf{v}_i , we could project the Jacobian matrix along \mathbf{u}_i and \mathbf{v}_i to find its degree of "energy" projected on this singular component, which is defined as

4

$$\varphi_i(\mathbf{W}) = (\mathbf{u}_i^\top \mathbf{J}_{\mathbf{W}} \mathbf{v}_i)^2.$$
⁽²⁾

(3)

 $\varphi_i(\mathbf{W}) \text{ can be interpreted as the length of the projection of } \mathbf{J}_{\mathbf{W}} \text{ along the singular vector direction of } \sigma_i(\mathbf{W}). In general, the higher <math>\varphi_i(\mathbf{W})$ is, the more energy of $\mathbf{J}_{\mathbf{W}}$ is distributed along the singular vector direction of $\sigma_i(\mathbf{W})$, indicating more substantial parametric updates along that singular component. To capture this, we introduce the concept of Stable Jacobian Energy (SJE), which quantifies the proportion of Jacobian energy concentrated on the top singular vectors within the SR.

Definition 2. (Stable Jacobian Energy). The Stable Jacobian Energy (SJE) of a matrix W is the ratio of its Jacobian matrix J_W 's energy projected on the singular vectors within the SR to the energy projected on all singular vectors, which is defined as

 $SJE(\mathbf{W}) = \frac{\sum_{i=1}^{\lfloor SR(\mathbf{W}) \rfloor} \varphi_i(\mathbf{W})}{\sum_{i=1}^{\min(m,n)} \varphi_j(\mathbf{W})}$

141

131

132

113 114 115

142

143 144

145

Where $|SR(\mathbf{W})|$ denotes the floor of $SR(\mathbf{W})$.

A large SJE value typically indicates that those top singular vectors within the SR will receive high
 Jacobian energy, and thus, the parametric updates along those singular components can be larger.

SR and SJE are the two key measures used to describe the behaviors of fine-grained network singularities in this work, and we present the empirical observations in the next section.

150 151 152

2.2 THE CURSE OF SINGULARITIES: EMPIRICAL EVIDENCE

We start by illustrating the overall trend of parametric singularities across the training session and discuss the role of LRs during the process. We then illustrate the cyclic interaction between parametric singularities and SJE, investigating network behaviors in the steps leading up to prominent training instabilities, ultimately revealing the existence of *the curse of singularities*. Without loss of generality, we plot figures in this section using data collected from a BERT-base model trained on Wikitext dataset, while the observations are prevalent across networks and datasets.

Parametric singularities and LR impacts. In Fig. 1, we present the layer-wise average parametric singularities, measured by SR, across the weight matrices of the query modules W_Q and feed-forward networks (FFN) W_{fc} as training advances. The SR values remain stable during the warm-up phase of around 1000 steps, with weights matrices consisting of significant singular values 172

173

174

188 189

190

191

192

193

194

195

196

197



Figure 1: SR curves of 12 layers across training steps. The left and right columns demonstrate the average SRs computed on weight matrices of the query modules and FFNs, respectively. All components of the model share the same trend (See in Appendix A).

175 on more than 35 singular directions in \mathbf{W}_Q or 180 in \mathbf{W}_{fc} . We can reasonably assume that the 176 expressivities of the network at this stage are intact. After the warm-up phase, we could observe 177 that the SR curves demonstrate a continuous downward trend despite some differences in layers and 178 network components. This trend can be explained as when the network learns to fit the dataset, 179 most irrelevant singular components are gradually discarded, and weight matrices focus on several 180 dominating singular directions. Notably, the monotonic increasing trend of singularities is more pronounced in deep layers. We believe that the increased singularities could somehow impair the 181 expressivity capabilities and lead to rank-deficiency representations, which is in accordance with the 182 insights of previous works (Bao et al., 2024). 183

We further discuss the impacts of LRs on parametric singularities. As shown in Fig. 2(a), a larger
LR will generally lead to SR curves of lower values (only an FFN of one layer is shown but this is a
common observation). That is, large LRs will typically amplify the effects of network singularities,
which are potentially associated with higher probabilities of training instabilities.



Figure 2: a) *Left*: SR curves of FFN in layer 6 under varying LRs; More layers/modules can be found in Appendix A with similar trends. b) *Right*: SR-SJE relationships across training steps.

201 SR-SJE Dynamics Driving Singularities. We delve further into the cyclic interaction between 202 Stable Jacobian Energy (SJE) and Stable Rank (SR), as visualized in Fig. 2(b). Intuitively, the SJE 203 values stay less than 0.4 during the warm-up stage, indicating < 40% Jacobian energy is projected 204 on the singular vectors within the stable ranks. As learning proceeds, we witness a constant upward 205 trend of SJEs that reaches approximately 0.8 in late training stages, implying that around 80% of 206 Jacobian energy is sprayed on the dominating singular directions, which is rather significant considering the stable ranks are even smaller than the warm-up stage. The SJE's increasing trend is 207 complementary to the declining tendency of SR curves. As network singularity intensifies and dom-208 inant singular vectors gain prominence in the weight matrices, Jacobian matrices of high SJE will 209 enhance the singular trend by further pushing parametric updates along those directions of signif-210 icant singular values. This cycle of escalating singularities, captured by the SR-SJE relationship, 211 underscores the curse of singularities, which ultimately leads to unstable losses and the breakdown 212 of training stability. 213

Training Breakdown. Then, we investigate the network's behavior leading up to training instabilities, focusing on how escalating singularity drives loss explosions and training collapse. Fig. 3 (top row, blue) shows a case of loss explosion with an LR of 5×10^{-4} , while the bottom row (red)



Figure 3: The top four panels depict the evolution of NTK- λ_{max} (for measuring gradient alignment), average token cosine similarity (between tokens from different sequences), Stable Rank(SR), training loss, and gradient norm during a case of loss explosion with a learning rate of 5×10^{-4} . The bottom four panels illustrate the same metrics during normal training with a learning rate of 1×10^{-4} .

represents normal training with an LR of 1×10^{-4} . The first column highlights an unusual plunge in the SR for the blue case, in contrast to the steady decrease in the red case. This sharp rise in parametric singularities leads to representations with overly high similarities in the feature space as shown in the second column of Fig. 3. The average cosine similarity between tokens reaches an unusually high 0.6 in shallow layers, climbing to an extreme 0.9 in deeper layers for the blue case, compared to typical values around 0.1 in the red one. Such similarity patterns suggest that network representations become rank-deficient and fail to retain sample-wise distinguishability.

This escalation of singularity, driven by the vicious cycle between Stable Jacobian Energy (SJE) 244 and SR, culminates in severe rank deficiency and heightened gradient alignment. We quantify this 245 alignment using the principal eigenvalue of a modified Neural Tangent Kernel (NTK) matrix (Jacot 246 et al., 2018), where each element captures the dot product of normalized gradients between pairs of 247 data points. The sharp rise of NTK λ_{max} , as shown in the third column of Fig. 3 for the blue case, 248 indicates that most gradients are increasingly aligning in the same direction, signaling the model's 249 diminishing ability to learn diverse patterns from different data. Ultimately, these metrics reach 250 critical values: token cosine similarities converge to 1.0 across all layers, and NTK λ_{max} hits its 251 upper bound, meaning the gradients for all data points are fully aligned. This marks the complete 252 collapse of the model's ability to differentiate between data points and capture distinct, nuanced 253 patterns. The training loss escalates sharply and remains irreversibly high, signaling the breakdown of the learning process and the model's failure to recover. 254

3 Method

236

255 256

257

To dispel the curse of singularities, we propose a training strategy named Parametric Singularity Smoothing (PSS) to *detect* approaching instability and *prevent* its occurrence by smoothing the singular spectra of weight matrices when necessary. PSS is an easy-to-implement, efficient, and robust method capable of both *rescuing* network from suffering instabilities and *restoring* the trainability even after the loss explosions.

263 264 **Detection**. The curse of singularities is observed to constantly befall with gradients of significantly 265 increasing magnitudes. Thus, we adopt gradient norms that can be fetched off-the-shelf during 266 backpropagation as an effective indicator to launch the smoothing operations. Specifically, we track 267 the ratio of the current step's gradient norms to the average gradient norms of all training steps, 268 denoted as $\mu_{t+1} = \frac{\|\mathbf{g}^{t+1}\|_F}{\|\mathbf{g}^t_{avg}\|_F}$, where \mathbf{g}^{t+1} is the gradient at step t + 1 of the model, and \mathbf{g}^t_{avg} is 269 the historical average gradient over the previous t steps, calculated using an exponentially weighted 269 moving average: $\mathbf{g}^t_{avg} = (1 - \alpha)\mathbf{g}^{t-1}_{avg} + \alpha \mathbf{g}^t$. Here α is the smoothing coefficient, typically 0 < 270 $\alpha \leq 1$. If it is spotted that $\mu_{t+1} \geq \tau$ where τ is a threshold (we empirically set τ to 2.5 throughout 271 the experiments), the model is considered to be stepping into the zone of potential instabilities, and 272 hence, the smoothing operation is required to meditate the parametric singularities. We emphasize 273 that the choice of τ is robust: a lower τ may cause false positives but without affecting normal 274 training, while a higher τ may miss instabilities, but PSS can still rescue and restore stable training even after instability has occurred. 275

276 Protection. Upon detecting instability, we mitigate the singularity of the parameter matrix W 277 in each linear module by smoothing the singular spectrum. To reduce computational cost, we de-278 compose only the dominant singular values and their corresponding singular vectors, limited in 279 number by the SR, reduce the dominant singular values, and reparameterize the matrix using the 280 corresponding singular vectors. Specifically, we first compute the largest singular value σ_1 using the Power Iteration Method (Golub & Van Loan, 2013), then calculate the SR(W). With this, we apply 281 Dominant Direction Decomposition (DDD) (Halko et al., 2009) to decompose the dominant part of 282 the parameter matrix: 283

$$\mathbf{W}_{dominant} = \sum_{i=1}^{\lfloor \mathrm{SR}(\mathbf{W}) \rfloor} \sigma_i \mathbf{u}_i \mathbf{v}_i^{\top},$$

where $\boldsymbol{\sigma} = \{\sigma_i\}_{i=1}^{\lfloor SR(\mathbf{W}) \rfloor}$ denotes dominant singular values of the parameter matrix \mathbf{W} .

We apply a smooth function f_{smooth} to σ to get the smoothed singular values

$$\boldsymbol{\sigma}^* = \{\sigma_i^*\}_{i=1}^{\lfloor \operatorname{SR}(\mathbf{W}) \rfloor} = f_{\operatorname{smooth}}(\boldsymbol{\sigma}),$$

and we then reparameterize the dominant part $\mathbf{W}_{dominant}$ without altering the singular vectors:

$$\mathbf{W}_{dominant}^* = \sum_{i=1}^{\lfloor \mathrm{SR}(\mathbf{W}) \rfloor} \sigma_i^* \mathbf{u}_i \mathbf{v}_i^\top,$$

resulting in the smoothed parameter matrix:

 $\mathbf{W}^* = \mathbf{W}^*_{dominant} + \mathbf{W} - \mathbf{W}_{dominant}.$

299 Given that the dominant singular vectors of the parameter matrix capture essential features from 300 the data, the goal during smoothing is to avoid overly diminishing their magnitude. We select $\sigma_{|SR(\mathbf{W})|+1}$ as a threshold, ensuring that all dominant singular values remain above this level.

302 The choice of f_{smooth} remains flexible and robust. Sub-linear functions such as Logarithmic with 303 Scaling, Softplus, Softmax, or smoothing operations like convolution that maintain the relative order 304 of dominant singular values are effective. Additionally, we show that an aggressive clipping-based 305 method, where all dominant singular values σ_i are reduced to $\sigma_{|SR(\mathbf{W})|+1}$, also demonstrates effec-306 tiveness. Ultimately, the objective is to increase the SR of weight matrices, with the specific choice 307 of smoothing function adaptable to the task.

By integrating detection and protection, PSS enables early intervention to prevent impending instability, stabilizes the parameter matrix under large LRs, and can rescue and restore stability even after instability has occurred. This leads to robust and reliable training dynamics.

311 312 313

308

309

310

284 285 286

287 288

289 290 291

292 293

295 296

297

298

301

4 **EXPERIMENTS**

314 4.1 EXPERIMENTAL SETUP 315

316 Model configuration. We conduct experiments on both the BERT (Devlin et al., 2019) and GPT-317 2 (Radford et al., 2019) models. We evaluate the effectiveness and robustness of PSS method by 318 exploring its performance across different model scales and LRs. We present the main experimen-319 tal results on BERT-base (110M parameters) and GPT-2-Medium (345M). We validate our method 320 on larger models, including BERT-large (340M), GPT-2-Large (774M) and GPT-2-XL (1.5B), and 321 details are in the appendix C.2 For the Masked Language Modeling (MLM) task, we train BERT from scratch using the Wikitext (Merity et al., 2016) dataset, and for the Causal Language Mod-322 eling (CLM) task we train GPT-2 using the mixed dataset of Amazon-review and the OpenWeb-323 Text (Gokaslan & Cohen, 2019). Both are trained using a LR schedule with warm-up followed by 324 decay. We present experimental results across a wide range of LRs for each model scale, spanning 325 from rates that enable stable convergence to those that lead to instability. More detailed configura-326 tions are provided in the appendix C.1. Our approach is compared against two stabilization methods, 327 including Gradient Clipping (GC) and Orthogonal Regularization (OR) (Brock et al., 2017; Brock, 328 2018).

Evaluation metrics. To assess the effectiveness of the stabilization method, we track the fre-330 quency of loss explosions across multiple trials under the same configurations following Zhai et al. (2023). While maintaining stable training, we analyze the expansion of the LR range for which 332 stable training occurs by presenting convergence performance under various LR settings. We eval-333 uate the impact of each method on the overall model training efficiency duration by evaluating the 334 required floating-point computations(FLOPs) (Rasley et al., 2020) and time overhead.

4.2 MAIN RESULTS

331

335 336

350

	Model	BERT-base		GPT-2-Medium			
Method		lr=1e-4	lr=2e-4	lr=4e-4	lr=5e-4	lr=1e-3	lr=2e-3
N	Jaive	$0/3_{3.0\pm0.2}$	$3/3_{7.5\pm0.1}$	$3/3_{7.5\pm0.1}$	$0/3_{3.8\pm0.1}$	$3/3_{7.4\pm0.1}$	$3/3_{7.5\pm0.1}$
	GC	$0/3_{2.9\pm0.1}$	$0/3_{2.7\pm0.1}$	$3/3_{7.4\pm0.1}$	$0/3_{3.7\pm0.1}$	$0/3_{3.9\pm0.1}$	$0/3_{5.0\pm0.1}$
	OR	$0/3_{3.1\pm0.1}$	$1/3_{2.9\pm0.1}$	$3/3_{7.4\pm0.1}$	$0/3_{3.8\pm0.1}$	$3/3_{7.5\pm0.1}$	$3/3_{7.4\pm0.1}$
]	PSS	$0/3_{3.0\pm0.1}$	$0/3_{2.8\pm0.1}$	$0/3_{2.6\pm0.1}$	$0/3_{3.8\pm0.1}$	$0/3_{4.0\pm0.1}$	$0/3_{4.4\pm0.1}$

346 Table 1: Results of stability test on BERT-base and GPT-2-Medium models. In the context of the x/y 347 (z) entry, y denotes the number of evaluations performed for the considered method, x represents instances where the model exhibited instability, and z signifies the perplexity of the model at the 348 conclusion of the training process. 349

To evaluate the effectiveness of the stabilization method, we explore different LRs 351 Stability. from stable convergence to loss explosion. For each fixed experimental configuration, we repeat the 352 experiments three times using different fixed random seeds and measure the frequency of loss explo-353 sions to assess the reliability of the method, ensuring that the stabilization approach can definitively 354 resolve instability issues. 355

356 Table 1 presents the comparison result, indicating that the GC and OR methods can only handle scenarios with small LRs, but at the cost of reducing the convergence speed and degrading the final 357 performance. At moderately higher LRs, they can only slightly alleviate instability issues but cannot 358 reliably and definitively stable the model training, and at even larger LRs, they completely fail. In 359 contrast, our method ensures stable training across all configurations without compromising the 360 model's convergence performance. 361

362 Fig. 4(a) illustrates the test loss curves for different methods on the BERT-base model at LRs of 363 1e-4 and 4e-4. Our method effectively detect two instances of loss explosion and quickly stabilize the training, with only minor, short-lived spikes in loss. Furthermore, it ensures stability without 364 significant additional training steps to recover to the pre-explosion loss levels. 365

366 We investigate the usable LR range once stable training is ensured. Fig. 4(b) presents the final test 367 loss across different LRs for the BERT-base and GPT-2-Medium models. On the BERT-base model, 368 the GC and OR methods expand the maximum LR from 1e-4 to 2e-4, a twofold increase, while our method pushed it to 2e-3, achieving a 10-fold improvement. Similarly, on the GPT-2-Medium 369 model, our method extends the LR range by by 10 times, outperforming other approaches. We 370 further validate our method on larger models, including BERT-large, GPT-2-Large and GPT-2-XL, 371 finding that PSS expand the LR range by up to 5 times without compromising convergence. Details 372 are in the appendix C.2. This broader range allows model engineers to set LR hyper-parameters more 373 flexibly without the risk of instability, enabling the use of larger LRs and reducing the complexity 374 and time required for hyper-parameter tuning. 375

Another potential benefit of PSS expanding the LR range is improved performance at convergence. 376 Typically, LR and final performance exhibit a U-shaped relationship. However, the training process 377 may suffer from instability when the LR is near the bottom of the curve. As shown in Fig. 4(b),



Figure 4: a) *Left*: Test loss curves of each stabilization method on two LRs: 1e-4 and 4e-4 in BERT-base model. b) *right*: The varying LRs and corresponding convergence performance of each stabilization method in BERT-base and GPT2-Medium models.

Method	CPU+GPU co	ost @ 1 step	CPU+GPU cost @ 100k step		
	Fp computations	Time	Fp computations	Time	
Naive	9.33 TFLOPs	$781\pm10~\mathrm{ms}$	933 PFLOPs	$9.15\pm0.08~\mathrm{h}$	
GC	9.33 TFLOPs	$784 \pm 12 \text{ ms}$	936 PFLOPs	$9.14\pm0.10~\mathrm{h}$	
OR	11.46 TFLOPs	$864\pm15~\mathrm{ms}$	1146 PFLOPs	$10.04\pm0.10~\mathrm{h}$	
PSS	9.36 TFLOPs	$1979\pm25~\mathrm{ms}$	936 PFLOPs	$9.17\pm0.09~\mathrm{h}$	

Table 2: Computational cost comparison of different methods in terms of floating-point (Fp) computations and execution time for both a single step and 100k steps.

407 408 409

406

394

395

396 397

under the BERT-base configuration, our method's expanded LR range enables it to achieve a lower
 test loss, which other methods fail to stabilize.

412 **Computation overhead.** For a parameter matrix $W \in \mathbb{R}^{m \times n}$, in the detection policy, we use 413 gradient norms as early indicators of instability, with a negligible computational cost of O(mn). In 414 the protection policy, the computational complexity of DDD on the matrix with $|SR(\mathbf{W})| = k$ is $\mathcal{O}(mn \log k)$. For a feature matrix $X \in \mathbb{R}^{n \times b}$, where n is the feature dimension and b the batch 415 size, the forward and backward pass complexity is $\mathcal{O}(mnb)$. In practice, since b and $\log k$ are of 416 the same order of magnitude, the cost of DDD is comparable to that of a forward-backward step. 417 PSS's protection is triggered fewer than 10 times in most cases, and even in extreme LR settings, 418 the activation frequency is less than 0.1% of the total steps, adding virtually little overhead to the 419 training process. 420

We experimentally demonstrate the computational cost of our method in Table 2. In the "1 step" section, we show the cost of using the protection strategy when PSS detects instability. It can be seen that a single invocation of PSS incurs up to 2.40 times the CPU + GPU time of a single step in the naive baseline. However, since the protection policy is rarely triggered, over the full training process, PSS introduces nearly no additional overhead. In the BERT-base model, the additional computational overhead introduced by our method accounts for only 0.21% of the baseline training time. Moreover, compared to the regularization methods that similarly adjust singularity, PSS achieves a %8.5 improvement in overall training time.

428 429

430 431 4.3 ROBUSTNESS OF PSS

Our method demonstrates exceptional robustness, specifically in the following three aspects:



Figure 5: The experiment is conducted on the BERT-base model, trained with a learning rate of 4e-4. a) Left: The test loss curve in the naive baseline, which experiences instability, compared to the test loss curve when PSS is applied at different stages. b) *Right*: The test loss curve using different smoothing policies upon detecting instability.

449 Fig. 5(a) highlights our method's effectiveness across different instability **Rescue Robustness.** stages. When applied preemptively, it causes only brief, minor loss spikes while preventing future 450 instability. If applied after divergence, it quickly halts gradient vanishing and restores normal loss 451 descent. In contrast, prior methods fail in such cases, requiring trial-and-error tuning from saved 452 checkpoints. This demonstrates the robustness of our detection metric: false positives do not disrupt 453 training, and even after full divergence, the method reliably restores stability. 454

455 **Smoothing Policy Robustness.** Fig. 5 demonstrates the effectiveness of various smoothing poli-456 cies, indicating that all are effective. Details of each policy are provided in the Appendix B.1. This suggests that as long as parametric singularity is reduced, the choice of smoothing policy becomes 457 less critical, highlighting the robustness of the policy selection. 458

459 Compatibility Robustness. Our method imposes no specific requirements on model architec-460 ture, optimizer, or dataset, and is highly flexible and easy to integrate, requiring only minimal code changes. Furthermore, it is orthogonal to other stabilization methods, allowing for simultaneous use 462 to further enhance training stability.

4.4 TRAINING DYNAMICS WITH OUR METHOD



Figure 6: The effect of our method on the analytical metrics before, including SR, NTK- λ_{max} and token-level cosine similarity. Each metric varies at the time our method is adapted

479 480

477

478

443

444

445

446 447 448

461

463 464 465

466

481 To better understand how our method stabilizes model training, Fig. 6 shows the changes in key 482 metrics from the analysis section, including SR, NTK- λ_{max} and token-level cosine similarity, after 483 applying our method. When instability occurs, both NTK- λ_{max} and token-level cosine similarity spike sharply, as noted in the analysis section, while the SR remains low. After applying our method, 484 these metrics return to normal levels, and the SR increases, allowing the model to resume stable 485 training.

486 5 **RELATED WORK**

487

488 **Large LRs.** Large LRs can offer significant benefits in deep learning. Jastrzebski et al. (2020) 489 shows that large LRs improve the conditioning of the covariance of gradients, leading to flatter min-490 ima that generalize better (Lewkowycz et al., 2020; Lu et al., 2023; Huang et al., 2020; Leclerc 491 & Madry, 2020). Additionally, Li et al. (2019) demonstrates from a feature learning perspective 492 that small LRs cause models to first memorize low-noise, hard-to-fit patterns, resulting in poorer generalization. Smith (2017) and Smith & Topin (2019) show that cyclical large LRs can sig-493 nificantly accelerate model convergence. However, the instability caused by large LRs has been 494 widely observed in empirical studies. Theoretical analyses often focus on sharpness bounds (Cohen 495 et al., 2021), and Wang & Roberts (2023) explains this instability through landscape flattening and 496 landscape shift. We establish a link between large LRs and parametric singularity, introducing the 497 concept of the "curse of singularities" to explain the resulting instability. 498

Training Stability. Training instability has been widely studied from a theoretical perspective. 499 Techniques like meticulous weight initialization (Glorot & Bengio, 2010; He et al., 2015; Mishkin 500 & Matas, 2015; Hu et al., 2020) and normalization (Ba et al., 2016; Ioffe & Szegedy, 2015; Ulyanov 501 et al., 2016) are critical for stabilizing neural networks. Additionally, research has examined the 502 effects of activation non-linearity (Hu, 2016), skip connections (Szegedy et al., 2017), and model depth and width (Hanin, 2018; Yang & Schoenholz, 2018) on stability. In practice, gradient clip-504 ping (Allen-Zhu et al., 2019) is a common method for managing instability by capping gradient 505 magnitudes. Adaptive methods, such as those proposed by Tang et al. (2023), improve on this by 506 incorporating optimizers like Adagrad. Other approaches include selectively freezing parameters 507 in ViT models (Chen et al., 2021) and using spectral normalization to prevent attention entropy 508 collapse (Zhai et al., 2023), further stabilizing training. Despite the aforementioned stability meth-509 ods, models still encounter instability under large LRs, and these methods prove ineffective once instability arises. Our PSS method significantly broadens the usable LR range and ensures reliable 510 training. 511

512 513

514

CONCLUSION 6

515 This work tackles the challenge of training instability from large LRs by identifying parametric 516 singularity as a key factor. Excessive singularity in weight matrices leads to rank deficiencies and 517 loss explosions. We uncover a vicious cycle, the "curse of singularities", that worsens instability as 518 training progresses, ultimately leading to rank-deficient representation and resulting in loss explosion. Our method smooths singular spectra, enabling stable use of large LRs. This foundation offers 519 insights into managing singularities and stabilizing training. We show our approach's robustness 520 across networks, extending LR ranges and reducing training times. 521

- 522 523 524 525 526 527 528 529 530 531 532 533 534 535
- 536
- 538

540 REFERENCES

558

559

- Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over parameterization. In *International conference on machine learning*, pp. 242–252. PMLR, 2019.
- Sanjeev Arora, Nadav Cohen, Wei Hu, and Yuping Luo. Implicit regularization in deep matrix factorization. *Advances in Neural Information Processing Systems*, 32, 2019.
- Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- Han Bao, Ryuichiro Hataya, and Ryo Karakida. Self-attention networks localize when qkeigenspectrum concentrates. *arXiv preprint arXiv:2402.02098*, 2024.
- Andrew Brock. Large scale gan training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*, 2018.
- Andrew Brock, Theodore Lim, James Millar Ritchie, and Nicholas J Weston. Neural photo editing with introspective adversarial networks. In *5th International Conference on Learning Representations 2017*, 2017.
 - X Chen, S Xie, and K He. An empirical study of training self-supervised vision transformers. In 2021 IEEE/CVF International Conference on Computer Vision (ICCV), pp. 9620–9629, 2021.
- Jeremy Cohen, Simran Kaur, Yuanzhi Li, J Zico Kolter, and Ameet Talwalkar. Gradient descent on neural networks typically occurs at the edge of stability. In *International Conference on Learning Representations*, 2021.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding, 2019.
- 566
 567
 568
 568
 569
 569
 569
 569
 569
 569
 560
 560
 560
 561
 562
 563
 563
 564
 564
 565
 565
 565
 565
 566
 566
 567
 568
 569
 568
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
 569
- Ruili Feng, Kecheng Zheng, Yukun Huang, Deli Zhao, Michael Jordan, and Zheng-Jun Zha. Rank diminishing in deep neural networks. *Advances in Neural Information Processing Systems*, 35: 33054–33065, 2022.
- 573 Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural
 574 networks. In *Proceedings of the thirteenth international conference on artificial intelligence and* 575 *statistics*, pp. 249–256. JMLR Workshop and Conference Proceedings, 2010.
- 577 Aaron Gokaslan and Vanya Cohen. Openwebtext corpus. http://Skylion007.github.io/
 578 OpenWebTextCorpus, 2019.
- Gene H Golub and Charles F Van Loan. *Matrix computations*. JHU press, 2013.
- Nathan Halko, Per-Gunnar Martinsson, and Joel A Tropp. Finding structure with randomness:
 Stochastic algorithms for constructing approximate matrix decompositions. *arXiv preprint arXiv:0909.4061*, 909, 2009.
- Boris Hanin. Which neural net architectures give rise to exploding and vanishing gradients? Advances in neural information processing systems, 31, 2018.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pp. 1026–1034, 2015.
- Kai Hu. Revisiting exploding gradient: A ghost that never leaves. *Revisiting_Exploding_Gradient*.
 pdf, 2016.
- 593 Wei Hu, Lechao Xiao, and Jeffrey Pennington. Provable benefit of orthogonal initialization in optimizing deep linear networks. arXiv preprint arXiv:2001.05992, 2020.

594 Wei Huang, Weitao Du, Richard Yi Da Xu, and Chunrui Liu. Implicit bias of deep linear networks 595 in the large learning rate phase. arXiv preprint arXiv:2011.12547, 2020. 596 Minyoung Huh, Hossein Mobahi, Richard Zhang, Brian Cheung, Pulkit Agrawal, and Phillip Isola. 597 The low-rank simplicity bias in deep networks. arXiv preprint arXiv:2103.10427, 2021. 598 Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by 600 reducing internal covariate shift. In International conference on machine learning, pp. 448–456. 601 pmlr, 2015. 602 603 Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. Advances in neural information processing systems, 31, 2018. 604 605 Stanislaw Jastrzebski, Maciej Szymczak, Stanislav Fort, Devansh Arpit, Jacek Tabor, Kyunghyun 606 Cho, and Krzysztof Geras. The break-even point on optimization trajectories of deep neural 607 networks. arXiv preprint arXiv:2002.09572, 2020. 608 609 Guillaume Leclerc and Aleksander Madry. The two regimes of deep network training. arXiv preprint 610 arXiv:2002.10376, 2020. 611 Aitor Lewkowycz, Yasaman Bahri, Ethan Dyer, Jascha Sohl-Dickstein, and Guy Gur-Ari. The large 612 learning rate phase of deep learning: the catapult mechanism. arXiv preprint arXiv:2003.02218, 613 2020. 614 615 Yuanzhi Li, Colin Wei, and Tengyu Ma. Towards explaining the regularization effect of initial large 616 learning rate in training neural networks. Advances in neural information processing systems, 32, 617 2019. 618 Miao Lu, Beining Wu, Xiaodong Yang, and Difan Zou. Benign oscillation of stochastic gradient 619 descent with large learning rates. arXiv preprint arXiv:2310.17074, 2023. 620 621 Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture 622 models, 2016. 623 624 Dmytro Mishkin and Jiri Matas. All you need is a good init. arXiv preprint arXiv:1511.06422, 2015. 625 626 Alireza Mousavi-Hosseini, Sejun Park, Manuela Girotti, Ioannis Mitliagkas, and Murat A Erdogdu. 627 Neural networks efficiently learn low-dimensional representations with sgd. arXiv preprint 628 arXiv:2209.14863, 2022. 629 630 Lorenzo Noci, Sotiris Anagnostidis, Luca Biggio, Antonio Orvieto, Sidak Pal Singh, and Aurelien 631 Lucchi. Signal propagation in transformers: Theoretical perspectives and the role of rank collapse. Advances in Neural Information Processing Systems, 35:27198–27211, 2022. 632 633 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language 634 models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019. 635 636 Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. Deepspeed: System opti-637 mizations enable training deep learning models with over 100 billion parameters. In *Proceedings* 638 of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, 639 pp. 3505–3506, 2020. 640 Mark Rudelson and Roman Vershynin. Sampling from large matrices: An approach through geo-641 metric functional analysis. Journal of the ACM (JACM), 54(4):21-es, 2007. 642 643 Leslie N Smith. Cyclical learning rates for training neural networks. In 2017 IEEE winter conference 644 on applications of computer vision (WACV), pp. 464–472. IEEE, 2017. 645 Leslie N Smith and Nicholay Topin. Super-convergence: Very fast training of neural networks using 646 large learning rates. In Artificial intelligence and machine learning for multi-domain operations 647

applications, volume 11006, pp. 369-386. SPIE, 2019.

Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alexander Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 31, 2017.
Jiaxi Tang, Yoel Drori, Daryl Chang, Maheswaran Sathiamoorthy, Justin Gilmer, Li Wei, Xinyang Yi, Lichan Hong, and Ed H. Chi. Improving training stability for multitask ranking models in recommender systems. In <i>Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining</i> , KDD '23, pp. 4882–4893, New York, NY, USA, 2023. Associa- tion for Computing Machinery. ISBN 9798400701030. doi: 10.1145/3580305.3599846. URL https://doi.org/10.1145/3580305.3599846.
Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Instance normalization: The missing in- gredient for fast stylization. <i>arXiv preprint arXiv:1607.08022</i> , 2016.
Lawrence Wang and Stephen Roberts. The instabilities of large learning rate training: a loss land-scape view. <i>arXiv preprint arXiv:2307.11948</i> , 2023.
Greg Yang and Sam S Schoenholz. Deep mean field theory: Layerwise variance and width variation as methods to control gradient explosion. 2018.
Emanuele Zangrando, Piero Deidda, Simone Brugiapaglia, Nicola Guglielmi, and Francesco Tud- isco. Neural rank collapse: Weight decay and small within-class variability yield low-rank bias. <i>arXiv preprint arXiv:2402.03991</i> , 2024.
Shuangfei Zhai, Tatiana Likhomanenko, Etai Littwin, Dan Busbridge, Jason Ramapuram, Yizhe Zhang, Jiatao Gu, and Joshua M Susskind. Stabilizing transformer training by preventing attention entropy collapse. In <i>International Conference on Machine Learning</i> , pp. 40770–40803. PMLR, 2023.
Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christo- pher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language models. <i>arXiv preprint arXiv:2205.01068</i> , 2022.

702 A DETAILED ANALYSIS

We conducted a detailed parametric singularity analysis on the BERT-base model. In Figs8, 9 and 10, we present the dynamic changes in the average stable rank of the Attention-head, Attention-Dense and FFN modules across different layers and learning rates. It can be observed that, as training progresses, the stable rank of various modules decreases across layers, with deeper layers showing consistently lower stable ranks compared to shallower ones. Additionally, regarding the learning rate, we observe that larger learning rates accelerate the growing singularity trends, and the parameter matrices tend to converge to a lower stable rank.





Figure 7: The dynamic changes in the average stable rank of the Attention-head modules across different layers and learning rates in BERT-base model.



Figure 8: The dynamic changes in the average stable rank of the Attention-head modules across different layers and learning rates in BERT-base model.

B METHOD DETAILS

806 807 B.1 Smoothing Policy

802

803 804

805

Given a weight matrix **W** and its dominant singular values $\boldsymbol{\sigma} = \{\sigma_i\}_{i=1}^{\lceil \text{SR}(\mathbf{W}) \rceil}$, we introduce several smoothing functions below.



Figure 9: The dynamic changes in the average stable rank of the Attention-Dense modules across different layers and learning rates in BERT-base model.

 Convolution. We apply a convolution kernel $\mathbf{k} = \{k_j\}_{j=-m}^m$ to the dominant singular values $\boldsymbol{\sigma}$ to smooth the transition between them. The smoothed singular values $\boldsymbol{\sigma}^{\text{smooth}} = \{\sigma_i^{\text{smooth}}\}$ are computed by convolving the kernel with the original singular values:

$$\sigma_i^{\text{smooth}} = \sum_{j=-m}^m k_j \cdot \sigma_{i+j}, \quad \forall i \in \{m+1, \dots, \lceil \text{SR}(\mathbf{W}) \rceil - m\}$$



Figure 10: The dynamic changes in the average stable rank of the FFN modules across different layers and learning rates in BERT-base model. 910

where σ_{i+j} refers to neighboring singular values, and the kernel k is typically normalized to ensure 912 that the sum of the kernel elements equals 1: 913

914
915
916
$$\sum_{j=-m}^{m} k_j = 1.$$

915

909

911

916 This convolution operation smooths the singular value spectrum by averaging over a local neigh-917 borhood of values. It reduces abrupt changes between adjacent singular values while preserving their relative order, ensuring that the most dominant singular values retain their influence in the parametric space.

Softmax. We apply the Softmax function to the dominant singular values $\boldsymbol{\sigma} = \{\sigma_i\}_{i=1}^{\lceil SR(\mathbf{W}) \rceil}$ to ensure that the singular values are normalized while retaining their relative importance. The Softmax-smoothed singular values $\sigma^{\text{softmax}} = \{\sigma_i^{\text{softmax}}\}$ are computed as follows:

$$\sigma_i^{\text{softmax}} = \frac{\exp(\beta \cdot \sigma_i)}{\sum_{i=1}^{\lceil \text{SR}(\mathbf{W}) \rceil} \exp(\beta \cdot \sigma_i)}$$

where β is a scaling parameter that controls the sharpness of the distribution. A higher β places more emphasis on the larger singular values, while a lower β spreads the emphasis more evenly across all singular values.

The Softmax function smooths the singular values by normalizing them in a way that maintains the hierarchy of dominance, but it prevents any one singular value from dominating excessively. This ensures that the singular value distribution remains balanced and well-behaved during the smoothing process.

С **EXPERIMENT DEATILS**

MODEL CONFIGURATION DETAILS C.1

The detailed training configuration of BERT and GPT-2 in shown in Table 3 and Table 4.

942	Configurations	BERT-base	BERT-large
943	Hidden activation	GELU	
944	Max position ambaddings	512	
945			
946	lokenizer	ROBER	la-base
047	Vocabulary size	502	65
947	Sequence length	128	
948	Layernorm ϵ	1e-12	
949	Dropout probability	0.1	
950	Optimizer	AdamW	
951	AdamW weight decay	0.01	
952	AdamW (β_1, β_2)	(0.9, 0.999)	
953	AdamW ϵ	0.01	
954	LR warm-up steps	1000	
955	LR schedular	Cosine	
956	Hidden size	768	1024
957	Intermediate size	$4 \times$ Hidden size	
958	Hidden Layers	12	24
959	Num. attention heads	12	16
960	Batch size	64	32

Table 3:	Configuration	1 of BERT	training.
----------	---------------	-----------	-----------

C.2 EXPERIMENT RESULTS ON LARGER MODELS

We validated the effectiveness of our method on larger models, including BERT-large (340M), GPT-2-Large (774M), and GPT-2-XL (1.5B). For these larger models, we primarily compared with the naive baseline to demonstrate our method's ability to stabilize training. In all three models, we successfully used large learning rates, which would cause instability in the naive baseline, to achieve stable training.

978						
979	Configurations	GPT-2-medium	GPT-2-large	GPT-2-XL		
980	Hidden activation		GELU			
981	Max position embeddings		512			
982	Tokenizer		r50-k			
983	Vocabulary size		50257			
984	Sequence length	256				
985	Layernorm ϵ	1e-5				
986	Dropout probability	0.5				
987	Optimizer	AdamW				
988	AdamW weight decay	0.1				
989	Adam $W(\beta_1,\beta_2)$	(0.9, 0.995)				
990	AdamW ϵ	0.01				
001	LR warm-up steps	1000				
991	LR schedular	Linear				
992	Hidden size	1024	1280	1600		
993	Intermediate size	4 ×	Hidden size			
994	Hidden Lavers	24	36	48		
995	Num. attention heads	16	20	25		
996	Batch size	64	32	16		
997						

Table 4: Configuration of GPT-2 training.

Model	BERT-Large		GPT-2-Large		GPT-2-XL	
Method	lr=5e-5	lr=1e-4	5e-5	1e-4	1e-5	5e-5
Naive	$0/2_{2.4\pm0.1}$	$2/2_{-}$	$0/1_{3.9}$	$1/1_{-}$	$0/1_{3.7}$	1/1_
PSS	$0/2_{2.4\pm0.1}$	$0/2_{2.7\pm0.1}$	$0/1_{4.0}$	$0/1_{4.2}$	$0/1_{3.8}$	$0/1_{3.9}$

Table 5: Results on BERT-Large, GPT-2-Large and GPT-2-XL