002

004

010 011

012

013

014

015

016

017

018

019

021

025

026

027

028 029

031

REVISITING CRITICAL LEARNING PERIODS IN DEEP NEURAL NETWORKS

Anonymous authors

Paper under double-blind review

ABSTRACT

Deep neural networks (DNNs) exhibit critical learning periods (CLPs) during early training phases, when exposure to defective data can permanently impair model performance. The prevalent understanding of such periods, primarily based on the interpretation of Fisher Information (FI), attributes CLPs to the memorization phase. However, our theoretical and empirical study exhibits that such explanations of CLPs are inaccurate because of the misunderstanding of the relationship between FI and model memorization. As such, we revisit the CLPs in DNNs from the information theory and optimization perspectives, gaining a better and more accurate understanding of CLPs. We visualize model memorization dynamics and observe that CLPs extend beyond the memorization phase. Additionally, we introduce the concept of the effective gradient, a novel metric able to quantify the actual influence of each training epoch on the optimization trajectory. Our empirical and theoretical analyses reveal that the norm of effective gradients generally diminishes over training epochs and eventually converges to zero, highlighting the disproportionate larger impact of initial training on final model outcomes. Besides, this insight also clarifies the mechanism behind permanent performance degradation due to defective initial training: the model becomes trapped in the suboptimal region of parameter space. Our work offers novel and in-depth understandings of CLPs and sheds light on enhancing model performance and robustness through such periods.

030 1

1 INTRODUCTION

032 The critical period refers to a specific time window in the early post-natal development of humans 033 and animals during which individuals are exceptionally sensitive to certain external stimuli or 034 experiences (Kandel et al., 2000; Wiesel & Hubel, 1963; Wiesel, 1993; Konishi, 1985). Such a period is crucial for the lifelong development of specific abilities or behavioral patterns. Inappropriate 035 stimulation or experiences during this period can lead to permanent impairment of a skill. Inspired by 036 the parallels between deep neural networks (DNNs) and biological neural connections, Achille et al. 037 (2018) first observed a similar phenomenon in the training process of neural networks. Specifically, if the model is trained with defective data, such as blurred data, during this initial period, called critical learning periods (CLPs), it will suffer permanent impairment in its final performance, regardless of 040 any additional training with high-quality data successively. Thereafter, similar CLPs phenomena have 041 been observed in deep linear networks (Kleinman et al., 2023a), under multi-view setting Kleinman 042 et al. (2023b), and under the federated learning paradigm (Yan et al., 2022; 2023a;b). These periods 043 are demonstrated to be crucial windows for potentially enhancing model performance and lifting 044 robustness against attacks Yan et al. (2022; 2023a;b).

To date, the underlying mechanisms that induce the CLPs and their resulting permanent impairment in DNNs have not been fully explored. The predominant explanations are based on observations related to Fisher Information (FI), as reported in Achille et al. (2018). It was noted that the FI metric increases during the initial training periods and subsequently decreases to a certain level and a significant increase in FI is observed when the model undergoes defective training at an early stage. They interpret the dynamics in FI as indicative of a "memorization phase," when the model memorizes information from the training data, followed by a "forgetting phase," during which the model reduces its retained information. Then, the observation of FI surge during defective training is attributed to abnormal growth in "synaptic strength" of neural connections (Achille et al., 2018) as the model memorizes too much defective data, preventing the neural connections from adapting themselves effectively to the normal data afterward under high strength. Therefore, they conjecture that CLPs are centered within this memorization phase.

056 However, the prevailing narrative about FI-based model memorization and its correlation with CLPs 057 presents significant issues. First, it mistakenly conflates CLPs-an intrinsic characteristic of the model during training—with observations of model impairment caused by defective data and the corresponding FI dynamics. The performance impairment from defective training is the result of 060 interference with CLPs, with the surge in FI as a consequence of unstable learning induced by 061 defective data, not the cause of CLPs. Therefore, using FI dynamics during defective training to 062 explain CLPs confuses cause and effect. Second, FI on training data does not accurately quantify 063 how much information a model retains but instead measures the model's sensitivity to the training 064 data and how much the training data contributes to the model training process. A rise in FI likely reflects increased sensitivity to noisy information in the defective data, rather than indicating the 065 memorization phase. As such, the claim that CLPs are driven by memorization based on FI dynamics 066 warrants further scrutiny. Given these considerations, a re-examination of CLPs in light of these 067 dynamics is essential. 068

069 In this paper, we conduct an in-depth exploration of the CLPs in DNNs from the perspectives of information theory and optimization. First, we provide a theoretical analysis of FI, clarifying its correlation 071 with defective training and its inadequacy in analyzing model memorization. Following Shwartz-Ziv & Tishby (2017), we employ mutual information, which directly measures the information shared 072 between two variables, to visualize the model's memorization phase and demonstrate that CLPs are 073 not solely centered in this phase. Second, we shift our focus to the fundamental aspects of model 074 training, i.e., SGD (Stochastic Gradient Descent) optimization essence. We propose a metric, called 075 the *effective gradient*, to measure the contribution of the update in each epoch toward the optimization 076 objectives. We provide theoretical proof that the norm of the effective gradient will converge to zero, 077 exhibiting a decreasing trend, and conduct extensive experiments to validate its variations in practice. Consequently, the initial training period naturally exerts a larger contribution to the parameter updates 079 during the optimization process, dictating the final performance level and illustrating the phenomenon of initial CLPs in the model. Additionally, the irrecoverable deterioration of effective gradients 081 following initial defective training makes the model unlikely to escape from the sub-optimal area, leading to permanent performance impairment.

- 083 The contribution of this paper can be summarized as follows:084
 - We conduct a rigorous theoretical analysis and empirical reassessment of FI, unveiling significant misunderstandings for its indicative to CLPs. Our analysis is crucial as it corrects the foundational motivations that have driven past research on CLPs.
 - We explain the CLPs from the perspective of the optimization process, providing a deeper understanding of why the initial training stages are pivotal. We propose the effective gradients metrics and mathematically reveal that the effective gradients decline in the model optimization dynamics, incurring the observation of CLPs.
 - We envision several key perspectives that are vital for guiding future research on improving model robustness and security by leveraging a more informed understanding of CLPs. These new perspectives expand the scope of CLPs research, paving the way for innovative approaches to enhance model training performance.
- 095 096 097 098

099

100

085

087

090

091

092

094

2 EXISTING EXPLANATIONS TO CLPS

2.1 BACKGROUND OF FISHER INFORMATION (FI)

101 Consider a neural network $f(\cdot)$ parameterized by the parameter θ , and denote the output probability of 102 class y given input x as $p_{\theta}(y|x)$. If the weights are perturbed by $\delta\theta$, the new weights are $\theta' = \theta + \delta\theta$, 103 and the perturbed output distribution is $p_{\theta'}(y|x)$. Here, we focus on the discrepancy between $p_{\theta}(y|x)$ 104 and $p_{\theta'}(y|x)$, by calculating their KL (Kullback-Leibler) divergence and approximating the perturbed 105 output distribution $p_{\theta'}(y|x)$ via Taylor expansion of $p_{\theta}(y|x)$ around θ , yielding

$$KL(p_{\theta} \parallel p_{\theta'}) = \mathbb{E}_{x} \left[\int p_{\theta}(y|x) \log \frac{p_{\theta}(y|x)}{p_{\theta'}(y|x)} \, dy \right] \approx \frac{1}{2} \delta \theta^{T} \cdot F(\theta) \cdot \delta \theta + o(\delta \theta^{2}) ,$$

where $o(\delta\theta^2)$ includes higher-order small terms and F is the FI Matrix (FIM), calculated by

$$F(\theta) = \mathbb{E}_{x \sim \hat{Q}(x)} \left[\mathbb{E}_{y \sim p_{\theta}(y|x)} \left[\nabla_{\theta} \log p_{\theta}(y|x) \nabla_{\theta} \log p_{\theta}(y|x)^{T} \right] \right] .$$
(1)

Essentially, FIM captures how sensitive the model output probabilities are to changes in the parameter across the data distribution (Amari & Nagaoka, 2000). In practice, since FIM is too large to compute, we usually adopt the trace of FIM, denoted as $Tr(F(\theta))$, to represent its value (Achille et al., 2018), which is abbreviated as FI in the rest of this paper.

116 2.2 INTERPRETATION OF CLPS THROUGH FI

To explore the CLPs, Achille et al. (2018) proposed to observe the FI dynamics in the training. First, 118 during the training process, the Fisher Information (FI) of the training data is observed to initially 119 increase in the early stages and then decrease to a stable level. The authors state that FI can be 120 interpreted as a measure of the amount of information about the training data that the model retains. 121 Based on this, the fluctuation in FI is taken to suggest that the model first enters a "memorization 122 phase", where it memorizes information from the training data, followed by a "forgetting phase", 123 where redundant or irrelevant information is discarded. Consequently, the authors conjecture that 124 Critical Learning Periods (CLPs) are concentrated in the initial memorization phase, proposing that 125 an FI increase could be used as an indicator for detecting these periods.

Second, if the model is trained on defective data during the early epochs, an abnormally high FI on the training data value results. The authors argue that when the training data is severely corrupted in the early stages, the network is forced to memorize more information to make predictions, which substantially increases the model's neural connection strength, as reflected by the high FI values. As a result, even if normal data is provided later, the network burdened by overly strong connections may struggle to adjust its connectivity, as described in Kirkpatrick et al. (2017), leading to impaired performance in the final model.

133

135

110

117

134 2.3 FLAWS IN FI-BASED EXPLANATION

The aforementioned explanations suffer two significant flaws, illustrated below. First of all, the 136 authors mistakenly conflate Critical Learning Periods (CLPs)-an intrinsic characteristic of the 137 model during training—with observations of model impairment caused by defective data and the 138 corresponding Fisher Information (FI) dynamics. CLPs should be understood as the initial phase 139 during which the model's learning disproportionately impacts its final performance, regardless of 140 the training data quality, thereby should be explained through the dynamics of normal training. The 141 performance impairment observed after training on defective data is a result of interference with 142 CLPs, with changes in FI merely reflecting this disruption, not explaining CLPs itself. Therefore, 143 using FI dynamics to directly explain the existence or behavior of CLPs is a case of confusing cause 144 and effect: the surge in FI is a consequence of the model's unstable learning from defective data rather than a fundamental driver of CLPs. While FI on the training data and CLPs may be correlated, 145 they are not causally linked. 146

Moreover, FI on the training data cannot accurately quantify the amount of information a model
retains about the training data. Instead, it measures the model's sensitivity to the data and the amount
of information that the training data can contribute to the model training. A rise in FI likely indicates
that the model is becoming increasingly sensitive to noisy information abundant in the defective
training data, rather than entering a memorization phase. As a result, the claim that CLPs are centered
on the memorization phase based on FI dynamics observation warrants further examination.

153 154

3 REVISITING FI DYNAMICS AND MODEL MEMORIZATION

This section first theoretically explores the dynamics of FI. Then we examine the relationship between
the CLPs and the model memorization to answer whether the CLPs are centered in the model
memorization phase.

- 159 3.1 THEORETICAL ANALYSIS OF FI
- 161 Achille et al. (2018) attributes the increase in FI to the rise of data information contained in the model and proposes to observe the initial CLPs through FI. However, it's important to correct that the

definition of FI is instead the amount of information that an observable random variable X carries about an unknown parameter θ of a distribution that models X (Fisher, 1925). In the machine learning context, X is the training data, and θ , being the distribution that models X, is the model parameter, so FI measures the amount of information provided by training data about the model parameters. A mathematical interpretation is provided below.

Proposition 1. The larger the $F(\theta)$ on data X, the lower uncertainty of the model parameter θ that is estimated from data X.

Proof. The proof is based on the well-known Cramér-Rao Lower Bound (Cramér, 1999; RAO, 1945). Define the variance of any unbiased estimator $\hat{\theta}$ as $Var(\hat{\theta})$, bounded by the reciprocal of the FI, i.e.,

 $Var(\hat{\theta}) \ge 1/F(\theta)$ (2)

Hence, a larger $F(\theta)$ implies a tighter lower bound on the variance of model parameter θ . This directly means less uncertainty in estimating model parameter θ from data X.

According to proposition 1, we can comprehend why FI is aptly termed "information," as it quantifies
the uncertainty degree in estimating model parameters via data X. FI of training data essentially
quantifies the amount of information the training data provides for optimizing the model parameters,
rather than the information the model contains about the training data. Moreover, FI of the training
data decreases during normal training as follows:

Theorem 1. If the model is trained with SGD with necessary assumptions, as presented in Appendix A.1, holds and the learning rate $\eta < \frac{2}{M}$, we have:

183 184 $\lim_{t \to \infty} Tr(F(\theta_t)) = 0 , \qquad (3)$

where θ_t is the model parameter under epoch t and M is the M-Lipschitz parameter as shown in Appendix A.1.

188 The proof is deferred to Appendix A.2. Theorem 1 demonstrates that the Fisher Information (FI) of the training data in a model trained using SGD will eventually converge to zero, indicating that the model 189 parameters stabilize as training progresses. This aligns with Proposition 1, as the model continues to 190 train on the data and becomes more well-generalized, the remaining information the training data 191 provides for optimizing the model parameters decreases. If we equate model memorization with the 192 dynamics of FI during the normal training process, an unreasonable conclusion will be led: the model 193 does not memorize anything. This contradicts the reality that a model can still memorize training 194 data, even when its parameters stabilize and FI decreases. Therefore, while FI offers insight into the 195 information the training data provides for model optimization, it is insufficient to describe or quantify 196 the model's memorization behavior.

197 198

3.2 MODEL MEMORIZATION MEASUREMENT

The two-phase learning phenomenon in neural networks, where a model initially memorizes training data and then forgets label-irrelevant information, has been observed in simple neural networks (e.g., MLPs) through the *Mutual Information* dynamics (Shwartz-Ziv & Tishby, 2017). Here, we follow their work to adopt the mutual information metrics to reexamine the relationship between the memorization phase and the CLPs.

Given any two random variables X and Y, with a joint distribution p(x, y), their Mutual Information is defined as:

$$I(X;Y) = D_{KL}[p(x,y)||p(x)p(y)] = H(X) - H(X|Y),$$

208 where $D_{KL}[\cdot]$ measures the KL divergence and $H(\cdot)$ is the entropy. It measures the reduction 209 in uncertainty about X due to the knowledge of Y, thereby quantifying the information that Y 210 contains about X. In the machine learning context, we can consider the neural network parameterized 211 with θ , denoted by $f(;\theta)$, as an encoder. For the training dataset X, the model outputs of $T_t =$ 212 $f(X; \theta_t)$ represent the learned representations from training data. The mutual information of $I(X; T_t)$ 213 measures the quantity of information that learned representations contain about the training data, specifically, the information memorized by the model in epoch t. Thus, by analyzing the evolution 214 of $I(X;T_t)$ as the training epoch grows, we can understand how the model memorization phase 215 happens, further exploring its relationship to CLPs. Although advanced methods for approximating



Figure 2: Variations in performance decrease and mutual information during defective training.

precise mutual information values exist, as proposed in Kleinman et al. (2021); Belghazi et al. (2018), our primary objective is to observe the dynamics of mutual information. Therefore, we calculate mutual information using the binning method, which is widely adopted in Shwartz-Ziv & Tishby (2017); Saxe et al. (2019); Goldfeld et al. (2019).

243 3.3 EXPERIMENT ANALYSIS244

216 217 218

219

220

222

224

225

226 227

228

229 230 231

232

233

234

235 236

237

Given the above understanding, we experimentally revisit the FI explanation from Achille et al. (2018). We aim to explore three questions: 1) whether FI phenomenon is the consequence of defective training; 2) whether an FI variation indicates the model memorization phase; and 3) whether the CLPs are centered in the memorization phase.

249 Since Achille et al. (2018) has demonstrated that the CLPs and FI dynamics do not rely on any specific 250 learning rate, batch size, model structure, and datasets, we conduct experiments with general settings. Here, we exhibit the experiments for ResNet-18 training on CIFAR-10, as the approximation of the 251 FI is more stable for ResNet-18 with a smooth loss landscape (Li et al., 2018). We employ the general 252 SGD as the optimizer with a fixed learning rate of 0.01 and a total training epochs of 300. Following 253 the methodology described in Achille et al. (2018), defective data are created by applying heavy 254 Gaussian blur to the original training dataset, and FI is quantified as the trace of the FI Matrix under 255 Monte-Carlo sampling approximation. Initially, we train the models using standard training data, 256 achieving baseline testing accuracies of 91.34% for VGG-16 and 93.28% for ResNet-18, respectively. 257 For the defective training, we divide the training epochs into three equal periods: Period 1 (epochs 258 1-100), Period 2 (epochs 101-200), and Period 3 (epochs 201-300). We train three models where 259 defective data is used as the training data for one period in each model, while the remaining two 260 periods utilize the original dataset. Detailed settings, metric calculations, and the results for VGG-16 261 are provided in Appendices B.1, B.2, and B.3.

262 Finding 1: The surge in FI is the consequence of defective training: Figure 1 illustrates the 263 variations in mutual information and FI on ResNet-18. We can observe that an FI surge can occur at 264 any point during training once the model is exposed to defective data, as evidenced by the increases 265 in the blue lines during any defective training periods in Figure 1. This surge occurs because the 266 model becomes more sensitive to noise and spurious patterns present in the defective data, leading 267 to an increase in FI. In other words, despite being unreliable, these defective data still provide a significant amount of noisy information to the model. This suggests that the fluctuations in FI are 268 merely a consequence of defective training, rather than a fundamental explanation for the Critical 269 Learning Period (CLPs) that occurs during the early stages of normal training. While FI captures the 270 instability caused by noisy or low-quality data, it does not directly explain the model's ability to learn 271 meaningful patterns during the CLPs. 272

Finding 2: FI dynamics do not match the model memorization variation: First, from Figure 1, 273 we can observe that during the training periods with normal data, FI rapidly decreases to near zero. 274 For example, during epochs 0-200 in the settings shown in Figure 1(c), the model undergoes a normal 275 training process with the original dataset, achieving a test accuracy of 93.29%, while FI eventually 276 approaches 0.0026. This indicates that such a FI trending does not reflect that the model has been 277 well-trained and has memorized useful knowledge from the training data. Additionally, we observe 278 that mutual information increases rapidly when the model encounters new data, regardless of whether 279 it is defective or normal, and eventually decreases to a certain value larger than 0. This suggests that 280 the model indeed goes through a memorization phase and then forgets label-irrelevant information to retain the useful knowledge.¹ It is consistent with the findings in Shwartz-Ziv & Tishby (2017). 281 However, while mutual information begins to decrease during the defective training epochs, FI 282 continues to increase, as evidenced by the rise in the blue line while the red line drops during each 283 defective training period. This discrepancy occurs because the model starts to forget label-irrelevant 284 information while the defective data is too poor to learn from, causing the model to remain unstable, 285 which is reflected by the increasing FI. These observations demonstrate that FI does not accurately 286 capture the model's memorization phase. 287

Finding 3: CLPs do not center on the memorization phase: Figure 1 also illustrates that the model 288 memorization phase can occur at any training period when the model is exposed to new data, as 289 evident from the multiple increases in the mutual information curves. Interestingly, after training 290 with defective data in Period 2 and recovering with original data in Period 3, the final accuracy of 291 VGG-16 and ResNet-18 models is 91.46% and 93.26%, respectively, comparable to the baseline 292 performance of 91.34% and 93.28%. This suggests that memorizing defective data during the middle 293 memorization phase is reversible and does not lead to permanent damage to final model performance. 294 In contrast, training with defective data during the CLPs leads to permanent impairment, highlighting 295 the distinction between the CLPs and the memorization phases. Additionally, we also conduct the 296 experiment to reveal the variations in model performance impairment and mutual information under 297 defective training epochs, during which the model is trained on defective data, as depicted in Fig. 2. The dotted lines in the figure indicate the epochs at which maximum mutual information and model 298 performance impairment occur. We observe that for VGG-16 and ResNet-18, the impairment of final 299 model performance reaches its maximum at epochs 42 and 38, respectively, suggesting that the CLPs 300 should conclude around these epochs. However, mutual information peaks earlier, at epochs 8 and 7, 301 indicating that the memorization phase ends much earlier. Therefore, the CLPs are even unaligned 302 with the initial model memorization phase. In conclusion, the CLPs are not simply centered in the 303 memorization phases as previously thought, calling for a new explanation. 304

305 306

307

312 313

4 EXPLAINING CLPs THROUGH OPTIMIZATION

As discussed in Section 3, the current narrative that explains CLPs through FI dynamics during 308 the defective training period encounters significant limitations and issues. Hence, we pivot to the 309 essence of neural network training, i.e., the optimization process, and offer a novel perspective for 310 theoretically analyzing why the initial training epochs can dominate the overall performance of a model undergoing general training, i.e., behaving as CLPs. 311

4.1 STOCHASTIC GRADIENT DESCENT

314 The mainstream model training is essentially an optimization process in which model parameters 315 are updated according to gradient descent rules. Among different popular optimization algorithms, 316 such as momentum and Adam (Kingma & Ba, 2014; Loshchilov & Hutter, 2018), their foundation 317 lies in Stochastic Gradient Descent (SGD) (Robbins & Monro, 1951). Consider a neural network 318 parameterized by θ_t in epoch t, represented as $f(;\theta_t)$, its update process under SGD is governed by 319 the following iterative formula:

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta_t} \ell(\mathbf{x}; \theta_t; \xi) , \qquad (4)$$

³²² 323

¹However, the reason why and when the model spontaneously forgets label-irrelevant information remains unclear (Shwartz-Ziv & Tishby, 2017) and is beyond the scope of this paper.

324 where $\nabla_{\theta_{\ell}} \ell(\mathbf{x}; \theta_{\ell}; \xi)$ is the stochastic gradient calculated with the subset of the whole training dataset, 325 which incurs a noise ξ . Following this equation, the model updates its parameter to align with the 326 negative gradient direction in each iteration. Besides, to ensure the convergence of the SGD, in 327 particular, we hold some necessary assumptions (Nemirovski et al., 2009; Li & Orabona, 2019), 328 which are presented in Appendix A.1

4.2 THEORETICAL ANALYSIS OF GRADIENT DYNAMICS 330

329

345

346

347 348 349

350

351

352

353 354

355 356 357

359

361

368 369

370

371

372

373

376

377

331 The CLPs refer to the phenomenon where the initial training phase has a disproportionately significant 332 influence on the final model performance. In other words, the impact of model gains on its final 333 performance decreases as training progresses, making any defects acquired during the initial period 334 difficult to recover from and potentially leading to permanent impairment. To explore this phe-335 nomenon, our motivation is to identify a metric able to capture this diminishing impact. As suggested 336 by Eqn. (4), if the learning rate η is a constant, which is small enough to ensure convergence, the 337 magnitude of model updates is modulated by the stochastic gradient, particularly its norm. Hence, we aim to demonstrate that the influence of the stochastic gradient dynamics on the model's way to the 338 optimum decreases throughout the optimization process, confirming that the initial training epochs 339 have a greater impact on the final performance. 340

341 Considering the random noise introduced in the stochastic gradient of each training epoch will 342 inevitably result in the uncertainty of its norm, we link the expected norm of the stochastic gradient 343 to that of the true gradient, denoted as $\ell(\theta)$. The true gradient is defined as the gradient computed based on the entire data samples, without the noise caused by random sampling. 344

Lemma 1. The expectation norm of stochastic gradient is bounded as

$$\mathbb{E}[||\nabla \ell(\theta;\xi)||^2] \le \mathbb{E}[||\nabla \ell(\theta)||^2] + \sigma^2.$$
(5)

The proof is deferred to Appendix A.3. Lemma 1 indicates that the expectation norm of the stochastic gradient is bounded by the expectation norm of the true gradient and the variance of the stochastic gradient noise. As such, we can further explore its expectation norm's variation via the true gradient.

Lemma 2. If learning rate $\eta < \frac{2}{M}$, the expected norm of the stochastic gradient $\mathbb{E}[||\nabla \ell(\theta; \xi)||^2]$ will converge alongside the true gradient $\mathbb{E}[||\nabla \ell(\theta)||^2]$, that is

$$\mathbb{E}[||\nabla \ell(\theta)||^2] \to 0, \ \mathbb{E}[||\nabla \ell(\theta;\xi)||^2 \to \sigma^2 \tag{6}$$

The proof is deferred to Appendix A.4. According to Lemma 2, the expected norm of the stochastic 358 gradient converges alongside the expected norm of the true gradient, ultimately reaching the variance of the noise as the true gradient approaches zero. In other words, it implies that the expected norm of 360 the stochastic gradient will not reach zero; even at the end of the training process, stochastic gradients persist. As a result, the stochastic gradient can not precisely reflect the trajectory to the optimum 362 during the optimization process due to its noisy nature. Given this understanding, to more accurately 363 measure the actual magnitude of model updates in each epoch, we propose a new metric called the 364 effective gradient, denoted as $\ell(\theta; \xi)$ and defined as below:

365 **Definition 1.** The effective gradient is defined as the projection of the stochastic gradient on the 366 direction of the true gradient, that is 367

$$\nabla \hat{\ell}(\theta;\xi) = Proj_{\nabla \ell(\theta)} \nabla \ell(\theta;\xi) .$$
⁽⁷⁾

The norm of the effective gradient can reflect the actual contribution of the model update along the direction to the optimum brought by the stochastic gradient. Thus, it can reveal the true contribution of the stochastic gradient in achieving the optimization objectives in each epoch. Upon Lemmas 1 and 2, we have the following theorem for analyzing its variations.

Theorem 2. If learning rate $\eta < \frac{2}{M}$, the expected norm of the effective gradient $\mathbb{E}[||\hat{\ell}(\theta;\xi)||^2]$ converges alongside the true gradient $\mathbb{E}[||\nabla \ell(\theta)||^2]$, that is 374 375

$$\mathbb{E}[||\nabla \ell(\theta)||^2] \to 0, \ \mathbb{E}[||\nabla \hat{\ell}(\theta;\xi)||^2 \to 0,$$
(8)

where M is the M-Lipschitz parameter as shown in Appendix A.1



Figure 3: The variation of the norm of effective gradient, stochastic gradient, and the accumulated 389 effective gradient with normal training.

390 The proof is presented in Appendix A.5. Theorem 9 suggests that the expected norm of the effective gradient also converges alongside the expected norm of the true gradient, ultimately reaching zero. Although the non-convexity of the loss function, in general, prevents the expected norm of the 393 true gradient from decreasing monotonically, Theorem 9 still implies that the norm of the effective 394 gradient is in the decreasing tendency throughout the entire training process. We will provide the 395 empirical results to exhibit this phenomenon in Section 4.4.

396 397

391

392

4.3 EXPLAINING CLPs THROUGH EFFECTIVE GRADIENT

398 Theorem 9 provides an important observation: as the training epoch progresses, the contribution 399 of model updates, i.e., the norm of the effective gradient, to achieve the optimization objectives 400 diminishes. This gives a hint that the initial training period naturally has a larger impact on the 401 parameter update during the optimization process in dominating the final performance, thereby behaving the phenomenon that the model possesses initial CLPs. If the model is trained on defective 402 data during initial epochs, the resulting updates will be misguided but significant, causing the model 403 to converge to a suboptimal parameter space rapidly. Even if subsequent training on good data occurs, 404 the actual optimization effect of the model updates gained will have a much lower influence compared 405 to the early defective data, making it difficult for the model to escape from this suboptimal space. As 406 a result, the model may suffer permanent performance impairment.

407 408 409

4.4 EXPERIMENT ANALYSIS

410 In this section, we evaluate the effective gradient variation during the training process to explore the 411 following questions: 1) whether the expected norm of effective gradient behaves the decreasing trend; 2) what causes the model permanent impairment? 412

413 To justify the versatility of the effective gradient in probing the critical learning phase, our experiments 414 on CIFAR-10 encompass both the conventional model architectures, such as VGG-16 and ResNet-18, 415 and the transformer architecture DeiT. The optimizer adopted is SGD with a fixed learning rate with 416 configurations detailed in Appendix B.1. Note that we fixed the learning rate rather than using an 417 annealing scheme to demonstrate that the decrease of the effective gradient during the optimization process does not result from the annealed learning rate. The experimental results on models trained 418 on other popular optimizers are presented in Appendix B.3. The training epochs are 200 for ResNet 419 and VGG, and 100 for DeiT, enough for the model to converge. The defective data is generated by 420 applying heavy Gaussian blur to the original training dataset, as documented in Achille et al. (2018). 421

422 The calculation of the effective gradient depends on the actual gradient, which is derived from the entire dataset at once and is hard to measure in practice. Instead, considering the essence 423 that the effective gradient is to measure the actual contribution of the model update along the 424 direction to the optimum, we use the vector pointing from the current model parameters to the 425 optimum to approximate the true gradient. Here, to reduce the approximation error, we take the 426 parameters averaged from the model with top-10 testing accuracy as the optimum. Then, we calculate 427 the projection of the stochastic gradient in this approximated direction as the effective gradient. 428 Besides, the effective gradients are always calculated on the current training data, thus during the 429 defective/normal training period, the gradients are obtained upon defective/original data. 430

Finding 4: The dynamics of effective gradient reflect the CLPs: Figure 3 exhibits the variation in 431 the norms of the effective gradient, the stochastic gradient, and the accumulated effective gradient.

434

435

436

437 438 439

440

441

442

443 444

445 446 447



(a) 50 epochs defective training at initial stage, leading to 5.78% accuracy decrease



(b) 25 epochs defective training at initial stage, leading to 4.92% accuracy decrease



(c) 50 epochs normal training at initial stage followed by 50 epochs defective training, merely leading to 0.17% accuracy fluctuation

Figure 4: The variation of the norm of effective gradient, stochastic gradient, and the accumulated effective gradient with defective training.

We observe that the variation in these gradients aligns with our theoretical analysis: the stochastic 448 gradient tends to converge to the variance of the noise, while the effective gradient approaches zero. 449 Notably, the norm of the effective gradient exhibits a pronounced decreasing trend during the training 450 epochs, characterized by an initial sharp decline in its red bar across all model structures. This 451 pattern supports our assertion that the actual contribution of the model update in the direction of the 452 optimum decreases, resulting in the updates during the initial epochs contributing significantly more 453 to the optimization process during the whole training. Additionally, the blue line representing the 454 accumulated effective gradient showcases a rapid increase initially, followed by a plateau, further 455 indicating that in the initial epochs, the model has already made significant progress toward the 456 optimum model. Indeed, it is this disproportionate large contribution to model optimization during the initial epochs that induces us to observe the CLPs in the early stages of training. 457

458 Finding 5: Unrecoverable effective gradients deterioration after defective training leads to 459 the impairment: Next, we utilize the effective gradient to elucidate the permanent impairment 460 resulting from defective training during the CLPs. First, Figure 4(a) displays the variation of different 461 gradients of the model that is trained with defective data in the initial 50 epochs, the CIFAR-10 462 dataset is employed as the example. We can observe that the norm of effective gradient rapidly 463 decreases and converges to 0 after 25 epochs, indicating that the model parameters reach a basin, albeit sub-optimally due to training on defective data. Moreover, after the 50-th epoch, although the 464 model is exposed to normal data and experiences a significant stochastic gradient, as indicated by the 465 tallest green bar, the actual progress towards the optimum remains limited, as shown by the short 466 red bar representing the norm of the effective gradient. This suggests that gains from the stochastic 467 gradient after the initial defective training are substantially reduced, which may enable the model to 468 escape from the bad basin induced by defective data but are insufficient to propel the model toward 469 the optimal performance, leading to 5.78% accuracy impairment. Even if we shorten the defective 470 training period to 25 epochs at which the model just converges to the sub-optimal, Figure 4(b) still 471 exhibits the similar phenomenon that unrecoverable deteriorated gains on the effective gradients, as 472 shown by the short and decreasing red bar representing the norm of the effective gradient, and the 473 model undergoes 4.92% accuracy impairment.

474 Finding 6: Defective training after CLPs will not lead to permanent impairment: Figure 4(c) 475 displays the variation in different gradients of the model that is initially trained with normal data for 476 50 epochs, followed by 50 epochs of defective training. During the initial critical period, the model 477 optimizes its parameters through the normal training process, which is accompanied by smooth varia-478 tions in both stochastic and effective gradients. This process is interrupted by the onset of defective 479 training, as indicated by the sharply increased stochastic gradients. However, the deterioration of 480 the effective gradient instead acts as a safeguard, preventing the model from significantly drifting due to the defective data, as demonstrated by the short and rapidly decreasing red bar of the effective 481 gradient and the slight increase in the blue line of the accumulated effective gradient. Therefore, after 482 the defective training period, the model can readily correct the adverse movements and return to a 483 favorable optimization trajectory, yielding merely 0.17% accuracy fluctuation. 484

485 In summary, the variation in effective gradients demonstrates a robust indication for theoretically explaining and empirically interpreting the CLPs at the initial stage of training.

486 5 DISCUSSIONS

Based on our analysis, we challenged the predominant explanation of CLPs and offered an optimization-based perspective for theoretically explaining them. Our exploration provides valuable insights into understanding DNNs and leaves problems that should be explored in the future.

491 Critical learning periods are not time windows with clear boundaries. According to our analysis, 492 the importance of the initial training period is attributed to a decreasing contribution gained in each epoch throughout the training process. This implies that the CLPs are the external manifestation of 493 <u>191</u> this consistently decreasing trend. In other words, an early training epoch is always more important than its subsequent ones, rather than some certain initial epochs within time windows signify the 495 CLPs. Additionally, Figures 4(a) and 4(b) reveal why model impairment does not infinitely increase 496 but instead reaches its maximum at certain epochs: the model becomes trapped in a sub-optimal valley 497 determined by the defective data after a certain number of epochs. Therefore, existing studies Achille 498 et al. (2018); Yan et al. (2022) that focus on determining the CLPs by observing the turning point 499 at which the largest model impairment exists are misguided. The turning point actually represents 500 the epoch when the model converges on defective data rather than indicating its CLPs. Given this 501 understanding, to better track the CLPs, we should observe the effective gradient dynamics. However, 502 the effective gradient is calculated based on the true gradient of the entire training data, which is 503 practically inaccessible during the training process. Thus, our future work will explore the method 504 for estimating the effective gradient and provide theoretical bounds for the approximation methods.

505 Theoretical quantification of effective gradient deterioration under defective training. Our 506 analysis and the experiments shown in Figure 4 demonstrate that defective data can lead to a rapid 507 deterioration of the effective gradient. However, we have yet to provide a theoretical framework to 508 analyze such variation when model training with defective data. To rigorously quantify the permanent 509 impairment caused by defective training, it is essential to account for the distribution shift and information missing introduced by defective data at the beginning of the training process, which will 510 significantly drive the model optimization trajectory, leading to long-term consequences on model 511 performance. In this paper, we still follow the current studies Achille et al. (2018); Yan et al. (2022) 512 that primarily focus on the impact of blurred data which have shown the most pronounced effects, but 513 fall short of providing comprehensive modeling of various forms of defective data. Addressing this 514 gap requires developing a rigorous theoretical definition of the defective data. This is a challenging 515 task that demands a deeper analysis of the interplay between data distribution shifts and gradient 516 behavior during training. We plan to explore this direction in future research, building upon the 517 effective gradient analysis foundations laid in our current exploration. 518

Empowering attacks and defenses on DNNs with CLPs natures. Since our exploration has 519 theoretically demonstrated that training gains during CLPs can significantly impact overall model 520 performance, it is worthwhile to extend this understanding to gain deeper insights and improve 521 current attack and defense mechanisms in DNNs. However, existing explorations of CLPs-centered 522 attacks and defenses Yan et al. (2023a;b) focus solely on the federated learning paradigm, proposing 523 to allocate more attack/defense resources to the rounds during which the model is still in the CLPs. 524 While these studies have provided valuable insights, a more in-depth analysis is needed to generalize 525 this approach beyond federated learning and examine how CLPs dynamics can influence model 526 vulnerabilities and robustness in broader training settings. Additionally, exploring the interactions 527 between CLPs and different types of adversarial attacks, such as data poisoning, backdoor insertion, or model inversion attacks, could offer practical insights into improving model resilience. For 528 instance, we could examine how poisoning attacks can be more effective during CLPs or how 529 defenses could be optimized by identifying and mitigating vulnerabilities specific to these critical 530 periods. Understanding these interactions will not only enhance model robustness but also lead to the 531 development of more effective defense strategies tailored to CLPs-specific vulnerabilities. We aim to 532 continue investigating these areas in future work. 533

534 535

6 CONCLUSION

In this study, we have critically revisited the CLPs and challenged existing FI-based explanations. We
found the FI dynamics do not match the model memorization variation and CLPs do not center on
the memorization phase as well. Hence, a new metric, the effective gradient, was proposed to explain
the CLPs from the optimization perspective. Our study not only advances the exploration of CLPs
but also highlights their significant implications for training robust neural networks.

540 REFERENCES

- Alessandro Achille, Matteo Rovere, and Stefano Soatto. Critical learning periods in deep networks.
 In *International Conference on Learning Representations*, 2018.
- Shun-ichi Amari and Hiroshi Nagaoka. *Methods of information geometry*, volume 191. American
 Mathematical Soc., 2000.
- 547 Mohamed Ishmael Belghazi, Aristide Baratin, Sai Rajeshwar, Sherjil Ozair, Yoshua Bengio, Aaron
 548 Courville, and Devon Hjelm. Mutual information neural estimation. In *International conference* 549 on machine learning, pp. 531–540. PMLR, 2018.
- ⁵⁵⁰ Harald Cramér. *Mathematical methods of statistics*, volume 26. 1999.
- Ronald Aylmer Fisher. Theory of statistical estimation. In *Mathematical proceedings of the Cambridge philosophical society*, volume 22, pp. 700–725. Cambridge University Press, 1925.
- Ziv Goldfeld, Ewout Van Den Berg, Kristjan Greenewald, Igor Melnyk, Nam Nguyen, Brian Kings bury, and Yury Polyanskiy. Estimating information flow in deep neural networks. In *International Conference on Machine Learning*, pp. 2299–2308. PMLR, 2019.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Eric R Kandel, James H Schwartz, Thomas M Jessell, Steven Siegelbaum, A James Hudspeth, Sarah
 Mack, et al. *Principles of neural science*, volume 4. McGraw-hill New York, 2000.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A
 Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming
 catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114
 (13):3521–3526, 2017.
- 570 Michael Kleinman, Alessandro Achille, Daksh Idnani, and Jonathan C Kao. Usable information and
 571 evolution of optimal representations during training. In *International Conference on Learning* 572 *Representations (ICLR)*, 2021.
- 573
 574 Michael Kleinman, Alessandro Achille, and Stefano Soatto. Critical learning periods emerge even in deep linear networks. *arXiv preprint arXiv:2308.12221*, 2023a.
- 576 Michael Kleinman, Alessandro Achille, and Stefano Soatto. Critical learning periods for multisensory
 577 integration in deep networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision* 578 and Pattern Recognition, pp. 24296–24305, 2023b.
- Masakazu Konishi. Birdsong: from behavior to neuron. Annual review of neuroscience, 8(1):125–170, 1985.
- Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein. Visualizing the loss landscape
 of neural nets. *Advances in neural information processing systems*, 31, 2018.
- 584
 585
 586
 586
 587
 Xiaoyu Li and Francesco Orabona. On the convergence of stochastic gradient descent with adaptive stepsizes. In *The International conference on artificial intelligence and statistics*, pp. 983–992, 2019.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Confer- ence on Learning Representations*, 2018.
- Arkadi Nemirovski, Anatoli Juditsky, Guanghui Lan, and Alexander Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on optimization*, 19(4):1574–1609, 2009.
 - Yurii Nesterov. Introductory lectures on convex optimization: A basic course, volume 87. 2013.

594 595 596	CR RAO. Information and accuracy attanaible in the estimation of statistical parameters. <i>Bull. Calcutta Math. Soc.</i> , 37:81–91, 1945.
597 598	Herbert Robbins and Sutton Monro. A stochastic approximation method. <i>The annals of mathematical statistics</i> , pp. 400–407, 1951.
599 600 601	Andrew M Saxe, Yamini Bansal, Joel Dapello, Madhu Advani, Artemy Kolchinsky, Brendan D Tracey, and David D Cox. On the information bottleneck theory of deep learning. <i>Journal of Statistical Mechanics: Theory and Experiment</i> , 2019(12):124020, 2019.
603 604	Ravid Shwartz-Ziv and Naftali Tishby. Opening the black box of deep neural networks via information. <i>arXiv preprint arXiv:1703.00810</i> , 2017.
605 606	Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. <i>arXiv preprint arXiv:1409.1556</i> , 2014.
607 608 609 610	Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In <i>International conference on machine learning</i> , pp. 10347–10357. PMLR, 2021.
611 612	Torsten N Wiesel. The postnatal development of the visual cortex and the influence of environment nobel lecture, 8 december 1981. <i>Physiology Or Medicine: 1981-1990</i> , pp. 61, 1993.
614 615	Torsten N Wiesel and David H Hubel. Effects of visual deprivation on morphology and physiology of cells in the cat's lateral geniculate body. <i>Journal of neurophysiology</i> , 26(6):978–993, 1963.
616 617	Gang Yan, Hao Wang, and Jian Li. Seizing critical learning periods in federated learning. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 36, pp. 8788–8796, 2022.
618 619 620 621	Gang Yan, Hao Wang, Xu Yuan, and Jian Li. Criticalfl: A critical learning periods augmented client selection framework for efficient federated learning. In <i>Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining</i> , pp. 2898–2907, 2023a.
622 623 624 625 626 627 628 629 630 631 632	Gang Yan, Hao Wang, Xu Yuan, and Jian Li. Defl: defending against model poisoning attacks in federated learning via critical learning periods awareness. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 37, pp. 10711–10719, 2023b.
633 634 635	
636 637 638 639	
640 641	
642 643 644	
645 646 647	

A MATHEMATICAL PROOFS

A.1 GENERAL ASSUMPTIONS FOR SGD

Considering the SGD with the noisy gradient, in particular, we always hold the following assumptions (Nemirovski et al., 2009; Li & Orabona, 2019):

- A1: The loss function is *L*-Lipschitz, i.e., $|\ell(\theta_1) \ell(\theta_2)| \leq L ||\theta_1 \theta_2||, \forall \theta_1, \theta_2 \in \mathbb{R}^d$.
 - A2: The loss function is *M*-smooth, i.e., ℓ is differentiable and its gradient is *M*-Lipschitz: $||\nabla \ell(\theta_1) - \nabla \ell(\theta_2)|| \le M ||\theta_1 - \theta_2||, \forall \theta_1, \theta_2 \in \mathbb{R}^d.$
- A3: The stochastic gradient $\nabla \ell(\theta; \xi)$ is the unbiased estimation of true gradient $\nabla \ell(\theta)$, i.e., $\mathbb{E}[\nabla \ell(\theta; \xi)] = \nabla \ell(\theta), \forall \theta \in \mathbb{R}^d$.
- A4: The noise in the stochastic gradient is bounded by the noise variance σ , i.e., $E[||\nabla \ell(\theta;\xi) \nabla \ell(\theta)||^2] \le \sigma^2, \forall \theta \in \mathbb{R}^d$.

Note that the gradient is always calculated on the current model θ_t and data **x**, so we simplify the heavy notation by omitting θ_t and **x** in $\nabla_{\theta_t} \ell(\mathbf{x}; \theta)$

A.2 PROOF OF THEOREM 1

Proof. Since the loss function is M-smooth (Assumption A2), we have:

$$\ell(\theta_{t+1}) \le \ell(\theta_t) + \nabla \ell(\theta_t)^\top (\theta_{t+1} - \theta_t) + \frac{M}{2} \left\| \theta_{t+1} - \theta_t \right\|^2 .$$
(9)

Given the SGD update rule: $\theta_{t+1} = \theta_t - \eta_t \nabla \ell(\theta_t; \xi_t)$, we substituting back it to Eqn. (9) and have

$$\ell(\theta_{t+1}) \le \ell(\theta_t) - \eta_t \nabla \ell(\theta_t)^\top \nabla \ell(\theta_t; \xi_t) + \frac{M}{2} \eta_t^2 \left\| \nabla \ell(\theta_t; \xi_t) \right\|^2 .$$
(10)

Given the assumption A3, we take the expectation over the randomness in noise ξ_t and the conditioning on model parameters θ_t as follows

$$\mathbb{E}\left[\ell(\theta_{t+1})\right] \le \mathbb{E}\left[\ell(\theta_t)\right] - \eta_t \mathbb{E}\left[\left\|\nabla \ell(\theta_t)\right\|^2\right] + \frac{M}{2}\eta_t^2 \mathbb{E}\left[\left\|\nabla \ell(\theta_t;\xi_t)\right\|^2\right] .$$
(11)

From Assumption A4, we have:

$$\mathbb{E}\left[\left\|\nabla \ell(\theta_t;\xi_t)\right\|^2\right] = \mathbb{E}\left[\left\|\nabla \ell(\theta_t) + (\nabla \ell(\theta_t;\xi_t) - \nabla \ell(\theta_t))\right\|^2\right]$$

$$= \left\|\nabla \ell(\theta_t)\right\|^2 + \mathbb{E}\left[\left\|\nabla \ell(\theta_t;\xi_t) - \nabla \ell(\theta_t)\right\|^2\right]$$

$$< \left\|\nabla \ell(\theta_t)\right\|^2 + \sigma^2.$$
(12)

So, we substitute Eqn. (12) back into the inequality Eqn. (11) to have

$$\mathbb{E}\left[\ell(\theta_{t+1})\right] \le \mathbb{E}\left[\ell(\theta_t)\right] - \eta_t \mathbb{E}\left[\left\|\nabla \ell(\theta_t)\right\|^2\right] + \frac{M}{2}\eta_t^2\left(\left\|\nabla \ell(\theta_t)\right\|^2 + \sigma^2\right) .$$
(13)

Bring the terms involving $\mathbb{E}\left[\ell(\theta_{t+1})\right]$ to the left, we have

$$\eta_t \mathbb{E}\left[\left\|\nabla \ell(\theta_t)\right\|^2\right] \le \mathbb{E}\left[\ell(\theta_t)\right] - \mathbb{E}\left[\ell(\theta_{t+1})\right] + \frac{M}{2}\eta_t^2\left(\left\|\nabla \ell(\theta_t)\right\|^2 + \sigma^2\right) .$$
(14)

Then, we sum both sides over t = 1 to T as

$$\sum_{t=1}^{T} \eta_t \mathbb{E}\left[\left\|\nabla \ell(\theta_t)\right\|^2\right] \le \mathbb{E}\left[\ell(\theta_1)\right] - \mathbb{E}\left[\ell(\theta_{T+1})\right] + \frac{M}{2} \sum_{t=1}^{T} \eta_t^2 \left(\left\|\nabla \ell(\theta_t)\right\|^2 + \sigma^2\right) .$$
(15)

Since the loss function is bounded below (e.g., non-negative), we have:

$$\mathbb{E}\left[\ell(\theta_1)\right] - \mathbb{E}\left[\ell(\theta_{T+1})\right] \le \mathbb{E}\left[\ell(\theta_1)\right] = L .$$
(16)

Then, we can simplify the inequality Eqn. (15) as

$$\sum_{t=1}^{T} \left(\eta_t - \frac{M}{2} \eta_t^2 \right) \mathbb{E} \left[\left\| \nabla \ell(\theta_t) \right\|^2 \right] \le L + \frac{M}{2} \sigma^2 \sum_{t=1}^{T} \eta_t^2 .$$
(17)

⁷⁰⁶ ₇₀₇ Since the learning rate $\eta_t < \frac{2}{M}$, we have

$$\eta_t - \frac{M}{2}\eta_t^2 \ge \frac{\eta_t}{2} \,. \tag{18}$$

Thus, for sufficiently large *t*:

$$\sum_{t=1}^{T} \frac{\eta_t}{2} \mathbb{E}\left[\left\| \nabla \ell(\theta_t) \right\|^2 \right] \le L + \frac{M}{2} \sigma^2 \sum_{t=1}^{T} \eta_t^2 .$$
(19)

714 Divide both sides by $\sum_{t=T_0}^{T} \frac{\eta_t}{2}$, we have:

$$\frac{\sum_{t=T_0}^{T} \eta_t \mathbb{E}\left[\|\nabla \ell(\theta_t)\|^2 \right]}{\sum_{t=T_0}^{T} \eta_t} \le \frac{2L}{\sum_{t=T_0}^{T} \eta_t} + M\sigma^2 \frac{\sum_{t=1}^{T} \eta_t^2}{\sum_{t=T_0}^{T} \eta_t} .$$
 (20)

719 Taking the limit as $T \to \infty$, we have

$$\lim_{T \to \infty} \frac{\sum_{t=T_0}^{T} \eta_t \mathbb{E}\left[\left\| \nabla \ell(\theta_t) \right\|^2 \right]}{\sum_{t=T_0}^{T} \eta_t} = 0$$
(21)

Since the learning rates $\eta_t < \frac{2}{M}$ are positive, we have

$$\liminf_{t \to \infty} \mathbb{E}\left[\|\nabla \ell(\theta_t)\|^2 \right] = 0.$$
(22)

Finally, based on the definition of the trace of the Fisher Information Matrix: $\text{Tr}(I(\theta)) = \mathbb{E}\left[\|\nabla \ell(\theta; \xi)\|^2 \right]$ and the unbiased estimation assumption A3, we have

$$\lim_{t \to \infty} \operatorname{Tr}(I(\theta_t)) = 0.$$
(23)

A.3 PROOF OF LEMMA 1

Proof. We directly expand the left-hand side of the inequality, then we have

$$\mathbb{E}[||\nabla \ell(\theta;\xi)||^{2}] = \mathbb{E}[||\nabla \ell(\theta;\xi) - \nabla \ell(\theta) + \nabla \ell(\theta)||^{2}]$$

$$= \mathbb{E}[||\nabla \ell(\theta) + (\nabla \ell(\theta;\xi) - \nabla \ell(\theta))||^{2}]$$

$$= \mathbb{E}[||\nabla \ell(\theta)||^{2}] + 2\mathbb{E}[\langle \nabla \ell(\theta), \nabla \ell(\theta;\xi) - \nabla \ell(\theta) \rangle]$$

$$+ \mathbb{E}[||\nabla \ell(\theta;\xi) - \nabla \ell(\theta)||^{2}].$$
(24)

According to the unbiased estimation assumption A3, we have

$$\mathbb{E}[\nabla \ell(\theta;\xi) - \nabla \ell(\theta)] = \mathbb{E}[\nabla \ell(\theta;\xi)] - \mathbb{E}[\nabla \ell(\theta)] = 0.$$
(25)

Thus, the second term in Eqn. (24) can be

$$2\mathbb{E}[\langle \nabla \ell(\theta), \nabla \ell(\theta; \xi) - \nabla \ell(\theta) \rangle] = 2\langle \nabla \ell(\theta), \mathbb{E}[\nabla \ell(\theta; \xi) - \nabla \ell(\theta)] \rangle = 0$$
(26)

Besides, given the noise bound assumption A4, the third term in Eqn. (24) can be

 $\mathbb E$

$$[||\nabla \ell(\theta;\xi) - \nabla \ell(\theta)||^2 \le \sigma^2$$
(27)

750 As a result, Eqn. (24) can be

751
752
753
754
755
751
752

$$\mathbb{E}[||\nabla \ell(\theta;\xi)||^2] = \mathbb{E}[||\nabla \ell(\theta)||^2] + 2\mathbb{E}[\langle \nabla \ell(\theta), \nabla \ell(\theta;\xi) - \nabla \ell(\theta) \rangle]$$

 $+ \mathbb{E}[||\nabla \ell(\theta;\xi) - \nabla \ell(\theta)||^2]$
 $\leq \mathbb{E}[||\nabla \ell(\theta)||^2] + \sigma^2$
(28)

756 A.4 PROOF OF LEMMA 2

Proof. We first demonstrate that the model can converge given the learning rate. Since the stochastic
gradients is the unbiased estimation, here, we omit the notation for the noise. Given A1 and A2, we
have the following equivalent expression based on Nesterov (2013),

$$\ell(\theta_2) - \ell(\theta_1) - \langle \nabla \ell(\theta_1), \theta_2 - \theta_1 \rangle \le \frac{M}{2} ||\theta_2 - \theta_1||^2 .$$
⁽²⁹⁾

763 Considering $\theta_1 = \theta_i$ and $\theta_2 = \theta_{i+1}$, we have

$$\ell(\theta_{i+1}) - \ell(\theta_i) - \langle \nabla \ell(\theta_i), \theta_{i+1} - \theta_i \rangle \le \frac{M}{2} ||\theta_{i+1} - \theta_i||^2 .$$
(30)

We reorganize the above inequality and consider the Eqn. (4),

$$\begin{split} \ell(\theta_{i+1}) \leq & \ell(\theta_i) + \langle \nabla \ell(\theta_i), \theta_{i+1} - \theta_i \rangle + \frac{M}{2} ||\theta_{i+1} - \theta_i||^2 \\ \leq & \ell(\theta_i) - \eta \langle \nabla \ell(\theta_i), \nabla \ell(\theta_i) \rangle + \frac{M}{2} \eta^2 ||\nabla \ell(\theta_i)||^2 \\ \leq & \ell(\theta_i) - \eta ||\nabla \ell(\theta_i)||^2 + \frac{M}{2} \eta^2 ||\nabla \ell(\theta_i)||^2 \\ \leq & \ell(\theta_i) - \eta (1 - \frac{M\eta}{2}) ||\nabla \ell(\theta_i)||^2 \end{split}$$

Since
$$\eta \leq \frac{2}{M}$$
, $\eta(1 - \frac{M\eta}{2}) \geq 0$. Therefore, we have

$$\ell(\theta_{i+1}) \le \ell(\theta_i) . \tag{31}$$

780 It indicates that if learning rate $\eta \leq \frac{2}{M}$, the loss function can converge, resulting in the fact that the 781 real gradient should converge to 0, that is

$$\mathbb{E}[||\nabla \ell(\theta)||^2] \to 0.$$
(32)

784 Then, we bring it to Lemma 1, we have

$$\mathbb{E}[||\nabla \ell(\theta;\xi)||^2] \le \mathbb{E}[||\nabla \ell(\theta)||^2] + \sigma^2 \to \sigma^2.$$
(33)

A.5 PROOF OF THEOREM 2

Proof. Based on the definition of effective gradient, its norm can be expressed as

$$\begin{split} ||\hat{\ell}(\theta;\xi)||^{2} &= \left| \left| \operatorname{Proj}_{\ell(\theta)} \ell(\theta;\xi) \right| \right|^{2} \\ &= \left| \left| \frac{\nabla \ell(\theta;\xi) \cdot \nabla \ell(\theta)}{||\nabla \ell(\theta)||^{2}} \right| \right|^{2} ||\nabla \ell(\theta)||^{2} \end{split}$$

797 Take the expectation on both sides of the formula, we have

$$\mathbb{E}[||\hat{\ell}(\theta;\xi)||^2] = \mathbb{E}\Big[\Big|\Big|\frac{\nabla\ell(\theta;\xi)\cdot\nabla\ell(\theta)}{||\nabla\ell(\theta)||^2}\Big|\Big|^2||\nabla\ell(\theta)||^2\Big]$$

Based to the unbiased estimation assumption A3, we have

$$\mathbb{E}[\nabla \ell(\theta;\xi)] = \mathbb{E}[\nabla \ell(\theta)] = \nabla \ell(\theta) .$$
(34)

Thus, we have

$$\mathbb{E}\left[\left|\left|\frac{\nabla \ell(\theta;\xi) \cdot \nabla \ell(\theta)}{||\nabla \ell(\theta)||^2}\right|\right|^2\right] = 1$$

807 Then, given Lemma 2, the expected norm of effective gradient can be

$$\mathbb{E}[||\hat{\ell}(\theta;\xi)||^2] = \mathbb{E}[||\nabla \ell(\theta)||^2] \to 0$$



Figure 5: The mutual Information and the FI variations on VGG-16.

B EXPERIMENT SUPPLEMENTS

B.1 EXPERIMENT SETTINGS

827 In all the experiments, we employ the standard model structure and implementation for VGG-16 (Simonyan & Zisserman, 2014), ResNet-18 (He et al., 2016), and DeiT (Touvron et al., 2021). 828 Unless otherwise specified, the optimizer is set to SGD with a fixed learning rate. Besides, We 829 employ standard data augmentation with random horizontal flipping and Random affine to mitigate 830 overfitting. To create the defective blurred data, we use the GaussianBlur function with a strong 831 parameter of kernel size=7 and sigma=3. For VGG-16 and ResNet-18, learning rate is 0.01, weight 832 decay is $1e^{-6}$, the batch size is 256. For DeiT, the learning rate is 0.05, weight decay is $1e^{-6}$, the 833 batch size is 256. Besides, the channel size is 3, the patch size is 4, the embed size is 512, the number 834 of heads is 8, the number of heads is 4, and the hidden size is 512. The teacher model is trained by 835 VGG-16. Our experiments are conducted on a workstation equipped with NVIDIA RTX 4090 GPU 836 with 24GB of VRAM.

837 838

839

843 844

849

861

862

810 811

812

813

814

815

816 817

818

819 820

821 822 823

824 825

826

B.2 METRICS CALCULATION

840 B.2.1 APPROXIMATION OF THE TRACE OF FISHER INFORMATION MATRIX 841

842 The trace of the Fisher Information Matrix (FIM) is calculated based on its definition as

$$Tr(F(\theta)) = \mathbb{E}_{x \sim \hat{Q}(x)} \left[\mathbb{E}_{y \sim p_{\theta}(y|x)} \left[||\nabla_{\theta} \log p_{\theta}(y|x)||^2 \right] \right],$$
(35)

where x is sampled based on the Monte-Carlo methods in the training dataset X and y is the corresponding label.

848 B.2.2 APPROXIMATION OF THE MUTUAL INFORMATION

To calculate the mutual information between training data and the model output representation, we employ the binning methods following these steps:

852 First, for each class, we discretize continuous model outputs by binning the logits for each class 853 into equal intervals between -1 and 1, making it easier to calculate joint and marginal probabilities. 854 For example, according to 30 equal intervals, we have the bins $\{-1, -0.933, \dots, 0.933, 1\}$. Then the logits t = [-0.9, -0.5, 0.1, 0.4, 0.8] can be transferred to [1, 7, 16, 21, 27]. Second, we create a 855 contingency table that counts the occurrences of each combination of training data and model output 856 bins to calculate their joint distribution P(x, t). Then, we compute the marginal distributions P(x)857 and P(t) by summing the counts over the rows and columns of the contingency table, respectively. 858 Third, we use these joint and marginal distributions to calculate mutual information based on the 859 definition: 860

$$I(X;T) = D_{KL}[p(x,t)||p(x)p(t)].$$

Here, the number of intervals is defined as the hyperparameter, which will influence the approximation performance. We follow Shwartz-Ziv & Tishby (2017) to set it as 30.



Figure 6: The variation of the norm of effective gradient, stochastic gradient, and the accumulated effective gradient with normal training of SGD with momentum=0.9.



Figure 7: The variation of the norm of effective gradient, stochastic gradient, and the accumulated effective gradient with normal training of SGD with Nesterov.



Figure 8: The variation of the norm of effective gradient, stochastic gradient, and the accumulated effective gradient with normal training of Adam.

918 B.3 ADDITIONAL PLOTS 919

920 This section presents the plots for additional experiments.

First, we show the additional results referred to in Section 3. Figure 1 exhibits the variations in mutual information and FI on VGG-16. These figures display the same phenomenon observed with ResNet-18. Additionally, we present the results referred to in Section 4. Here, we train models with the optimizers: SGD with momentum, SGD with Nesterov, and Adam. Figures 6, 7, and 8 show the variations in the norms of the effective gradient, the stochastic gradient, and the accumulated effective gradient, respectively. The effective gradients exhibit a decreasing trend in all scenarios.