L-MSA: LAYER-WISE FINE-TUNING USING THE METHOD OF SUCCESSIVE APPROXIMATIONS

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027 028 029

031

Paper under double-blind review

ABSTRACT

With the emergence of large-scale models, the machine learning community has witnessed remarkable advancements. However, the substantial memory consumption associated with these models has emerged as a significant obstacle to large-scale training. To mitigate this challenge, an increasing emphasis has been placed on parameter-efficient fine-tuning methodologies, which adapt pre-trained models by fine-tuning only a subset of parameters. We observe that in various scenarios, fine-tuning different layers could lead to varying performance outcomes, and selectively fine-tuning certain layers has the potential to yield favorable performance results. Drawing upon this insight, we propose L-MSA, a novel layer-wise finetuning approach that integrates two key components: a metric for layer selection and an algorithm for optimizing the fine-tuning of the selected layers. By leveraging the principles of the Method of Successive Approximations, our method enhances model performance by targeting specific layers based on their unique characteristics and fine-tuning them efficiently. We also provide a theoretical analysis within deep linear networks, establishing a strong foundation for our layer selection criterion. Empirical evaluations across various datasets demonstrate that L-MSA identifies layers that yield superior training outcomes and fine-tunes them efficiently, consistently outperforming existing layer-wise fine-tuning methods.

1 INTRODUCTION

With the increasing application of large-scale models across diverse task domains(Devlin et al., 2019; Dosovitskiy et al., 2021), domain-specific fine-tuning has emerged as a pivotal strategy to bolster their effectiveness in downstream tasks(Käding et al., 2017; Raffel et al., 2020). However, these fine-tuning methods are often resource-intensive, presenting significant challenges in the development of largescale models. Efforts to address these challenges have led to the development of Parameter-Efficient Fine-Tuning (PEFT) techniques, which aim to mitigate computational costs. These techniques encompass various approaches, such as prompt-based methods(Diao et al., 2022; Hambardzumyan et al., 2021; Lester et al., 2021; Liu et al., 2023), adapter methods(Diao et al., 2023; Houlsby et al., 2019; Hu et al., 2021), and selective methods(Li et al., 2024; Liu et al., 2021; Zaken et al., 2021).

Among the array of Parameter-Efficient Fine-Tuning (PEFT) techniques, layer-wise fine-tuning algorithms have emerged as a promising solution(Lee et al., 2022; Pan et al., 2024). Instead of updating all parameters simultaneously, layer-wise fine-tuning focuses on iteratively fine-tuning individual layers of the model. This approach not only reduces computational costs but also allows for more targeted adjustments, potentially leading to improved performance on downstream tasks.

However, the specific layer to fine-tune may vary based on the relationship between the source and target datasets. To explore this, we conduct experiments with a Data-efficient Image Transformer (DeiT)-Tiny (Touvron et al., 2021) in three scenarios:

049

 Pre-training on ImageNet(Deng et al., 2009a) and fine-tuning on CIFAR-100(Krizhevsky, 2009).

2. Pre-training on CIFAR-100 and transforming the input data by element-wise multiplication with a matrix, where each entry is e^x and x follows a standard normal distribution. Finetuning is then performed on the transformed data.



Figure 1: Layer-wise Fine-tuning in Different Scenarios

070 In each case, we visualize the outcomes of layer-wise fine-tuning compared to full fine-tuning after a 071 single epoch, with consistent observations even over extended training. In case 1, where the dataset shares similar low-level features but different high-level features compared to the original data, fine-072 tuning later layers outperforms earlier layers. Conversely, in case 2, with similar high-level features 073 but different low-level features, fine-tuning earlier layers yields better performance. Finally, in case 3 074 involving random tasks, all layers are equally significant, and fine-tuning individual layers alone may 075 not suffice. This variability raises the question of whether we can algorithmically determine which 076 layer(s) to fine-tune and how to perform effective layer-wise fine-tuning. 077

To address the aforementioned challenge, we propose L-MSA, a novel layer-wise fine-tuning approach that consists of two core components: a metric for layer selection and an algorithm for optimizing the fine-tuning of the selected layer. This targeted optimization seeks to enhance overall model performance by leveraging the specific strengths of different layers.

We leverage the principles of the Method of Successive Approximations (MSA) (Chernousko & Lyubushin, 1982; Li et al., 2018) within our L-MSA framework, addressing both layer selection and layer fine-tuning. The first component of our approach introduces a novel metric, derived from the state and co-state variables in MSA, which serves as the criterion for selecting layers. The second component focuses on utilizing the MSA to optimize the fine-tuning of the selected layers. This integrated approach ensures efficient optimization by systematically refining the layer-wise fine-tuning process, ultimately leading to improved performance.



Figure 2: Overview of our proposed L-MSA method: We begin with a feed-forward pass to compute the state x_i at each layer, followed by a back-propagation step to determine the co-state p_i . Utilizing both x_i and p_i , we compute a metric \hat{J}_n for every layer, as defined in equation 8, to guide layer selection. We then select the layer with the smallest metric, denoting its index as m, and maximize H_m over the parameters θ_m . By fixing m and repeating these steps iteratively, we refine the layer parameters, converging toward a desired solution that enhances model performance.

106

090

092

095

096

098

099

068

069

107 Furthermore, we provide a comprehensive theoretical analysis of our L-MSA approach within the context of deep linear networks(Arora et al., 2018b). This analysis clarifies the metric we utilize for

108 optimal layer selection, framed within a greedy one-step optimization framework. By establishing a 109 solid theoretical foundation, we pave the way for understanding how our method effectively enhances 110 model training.

111 Finally, we evaluate our approach across various datasets and tasks, utilizing multiple model architec-112 tures, and compare L-MSA with baseline layer-wise fine-tuning methods. Our results demonstrate 113 that the metric proposed in L-MSA effectively identifies the layers that will yield better training 114 outcomes. L-MSA consistently outperforms most baselines, achieving top performance in several 115 tasks and ranking as the most effective method overall, reinforcing the practical applicability of 116 our approach in real-world scenarios. We also conduct ablation tests, highlighting the metric's 117 effectiveness in layer selection and the advantages of using MSA to optimize the chosen layers.

118 119 120

121

122

123

124

125

127

128

129

130

131 132

133

134

141 142

143

144

145

147

148 149 150

156 157

158

159

We summarize our key contributions as follows:

- We experimentally show that in various scenarios, fine-tuning different layers could lead to varying performance outcomes, and selectively fine-tuning certain layers has the potential to yield favorable performance results.
- We propose the L-MSA method, which introduces a new criterion for selecting layers to fine-tune, and we also propose utilizing the method of successive approximations for layer-wise fine-tuning within our L-MSA approach, ensuring efficient optimization and improved learning outcomes.
 - We provide a theoretical analysis of our approach in the context of deep linear networks, clarifying the metric for optimal layer selection within a greedy one-step framework.
 - We empirically validate the effectiveness of our methodology in accurately identifying and efficiently fine-tuning the target layer across diverse datasets.
- L-MSA: LAYER-WISE FINE-TUNING USING THE METHOD OF SUCCESSIVE 2 APPROXIMATION

135 Adopting the control viewpoint for layer-wise fine-tuning offers a structured optimization process 136 through Pontryagin's Maximum Principle (PMP)(Pontryagin et al., 1962). This perspective treats 137 each layer as part of a controlled dynamical system, enabling precise adjustments to specific layers by 138 assessing their impact on the overall loss via the Hamiltonian. Consequently, this method facilitates 139 efficient fine-tuning by focusing on layers that offer the most significant performance improvement, 140 thereby making the optimization process more systematic and effective.

2.1 BACKGROUND: PONTRYAGIN'S MAXIMUM PRINCIPLE AND METHOD OF SUCCESSIVE **APPROXIMATION**

In supervised learning, given a collection of K sample input-label pairs $\{x^i, y^i\}_{i=1}^K$, our objective is to infer and approximate a function $F: \mathcal{X} \to \mathcal{Y}$ that accurately maps input data instances 146 x^i to their corresponding target outputs y^i . To view supervised learning within the dynamical systems framework, particularly relevant to deep and residual architectures, we consider the inputs $x = (x^1, x^2, \dots, x^K) \in \mathbb{R}^{d \times K}$ as the initial condition of a system of ordinary equations

$$\dot{x}_t^i = f(t, x_t^i, \theta_t), \quad x_0^i = x^i, \quad 0 \le t \le T, \quad i = 1, \dots, K,$$
(1)

(2)

151 where $\theta: [0,T] \to \Theta$ is the control parameters and $x_t = (x_t^1, \cdots, x_t^K) \in \mathbb{R}^{d \times K}$. In this context, 152 $f(t, x_t^i, \theta_t)$ encapsulates the transformation process within the neural network, while θ_t represents 153 the parameters at time t that govern this transformation.

154 The supervised learning problem can be formulated as 155

- $\min_{\theta \in \mathcal{U}} \sum_{i=1}^{K} \Phi_i\left(x_T^i\right) + \int_0^\top L\left(\theta_t\right) dt,$ $\dot{x}_{t}^{i} = f(t, x_{t}^{i}, \theta_{t}), \quad x_{0}^{i} = x^{i}, \quad 0 \le t \le T, \quad i = 1, \dots, K,$
- where $\Phi_i(\cdot) := \phi(\cdot, y^i)$ is the loss function, and $L: \Theta \to \mathbb{R}$ is a running cost, or the regularization 161 term.

We define the Hamiltonian $H: [0,T] \times \mathbb{R}^d \times \mathbb{R}^d \times \Theta$ given by

$$H(t, x, p, \theta) = p \cdot f(t, x, \theta) - L(\theta)$$
(3)

Pontryagin's Maximum Principle(PMP)(Pontryagin et al., 1962) shows a set of necessary conditions
 for optimal solutions to equation 2, which provides an alternative numerical algorithm to train
 equation 2 and its discrete-time formulation.

Theorem 2.1 (Pontryagin's Maximum Principle) Let $\theta^* \in U$ be an essentially bounded optimal control, i.e. a solution to equation 2 with $ess sup_{t \in [0,T]} \|\theta_t^*\|_{\infty} < \infty$ (ess sup denotes the essential supremum). Denote by x^* the corresponding optimally controlled state process. Then, there exists an absolutely continuous co-state process $P^* : [0,T] \to \mathbb{R}^d$ such that the Hamilton's equations

$$\dot{x}_{t}^{*} = \nabla_{p} H\left(t, x_{t}^{*}, P_{t}^{*}, \theta_{t}^{*}\right), \qquad x_{0}^{*} = x, \dot{P}_{t}^{*} = -\nabla_{x} H\left(t, x_{t}^{*}, P_{t}^{*}, \theta_{t}^{*}\right), \quad P_{T}^{*} = -\nabla\Phi\left(x_{T}^{*}\right),$$
(4)

164

169

are satisfied. Moreover, for each $t \in [0, T]$, we have the Hamiltonian maximization condition

179 180 181

183

185 186 187

188 189

190 191 192

197

199

200 201

202 203 204

205 206 207

208

210

211 212 Consider an *N*-layer deep neural network, which can be interpreted as a discrete-time formulation of equation 2. Within this framework, the supervised learning problem can be expressed as follows:

 $H(t, x_t^*, P_t^*, \theta_t^*) \geq H(t, x_t^*, P_t^*, \theta)$ for all $\theta \in \Theta$.

$$\min \sum_{i=1}^{K} \Phi_i \left(x_N^i \right) + \sum_{n=0}^{N-1} \delta_t L \left(\theta_n \right)$$

$$x_{n+1}^i = g_n(x_n^i, \theta_n), \ x_0^i = x^i, \ i = 0, 1, \cdots, K.$$
(6)

Here $g_n(x_n^i, \theta_n) = x_n^i + \delta_t f_n(x_n^i, \theta_n)$. Similar to equation 3, define the scaled discrete Hamiltonian

$$H_n(x, p, \theta) = p \cdot g_n(x, \theta) - \delta_t L(\theta)$$
(7)

(5)

In the following algorithms, we employ an augmented variant of Hamiltonian(Li et al., 2018), which additionally subtracts a regularization term of $\frac{1}{2}\rho ||x_{n+1} - g_n(x_n, \theta_n)||_2^2 + \frac{1}{2}\rho ||p_n - p_{n+1}\nabla_x g_n(x_n, \theta_n)||_2^2$ from the Hamiltonian discussed in equation 7. Here ρ serves as a hyperparameter, with its reciprocal $1/\rho$ exerting a similar effect as the learning rate.

A modification of the successive approximations method can be employed to address the Pontryagin Maximum Principle (PMP), thereby yielding an alternative training algorithm for deep learning(Li et al., 2018). We present the extended method of successive approximation in Figure 3.



Figure 3: Extended Method of Successive Approximation(E-MSA)

In each iteration, we commence with a feed-forward pass to compute the state x_i for $i = 0, 1, \dots, N$, followed by a back-propagation step to compute the co-state p_i for $i = N, N-1, \dots, 0$. Subsequently, we calculate the Hamiltonian for each layer using both the state and co-state, seeking to maximize H_n over θ_n . We iteratively perform these steps to converge towards the desired solution.

216 2.2 L-MSA: LAYER-WISE FINE-TUNING VIA MSA

240 241 242

253 254

255 256

257

258

259

As indicated in Section 1, it is often the case that fine-tuning the entire network is unnecessary. Rather,
 the focus lies in fine-tuning only a single layer or a small subset of layers. In response, we introduce
 L-MSA, a novel approach for layer-wise fine-tuning that consists of two key components: a metric
 for layer selection and an optimization algorithm for fine-tuning the selected layer.

We leverage the principles of the Method of Successive Approximations (MSA) (Chernousko & Lyubushin, 1982; Li et al., 2018) as the foundational principle for both layer selection and layer fine-tuning. Our method aims to improve model performance by focusing on the specific strengths of individual layers, targeting the most impactful layers for fine-tuning.

226 Denote $\Phi(x_N) = \sum_{i=1}^{K} \Phi_i(x_N^i)$ and $G_n(\cdot) = \Phi \circ g_{N-1}(\cdot, \theta_{N-1}) \circ \cdots \circ g_{n+1}(\cdot, \theta_{n+1})$, which maps 228 the state of the (n+1)-th layer to the terminal loss using the latter part of the model. Denote the 229 terminal loss $J = \sum_{i=1}^{K} \Phi_i(x_T^i)$ as a function of the (n+1)-th layer $J_n(\theta_n)$.

Layer Selection: A natural approach to layer selection is to choose the layer for fine-tuning that
 minimizes the loss and thus maximizes performance. In our proposed method, for the layer selection
 process, we approximate the optimal updated loss resulting from fine-tuning individual layers and
 use this approximated loss as the criterion for selecting layers.

To approximate the updated loss, we employ the principle of MSA. We begin with a feed-forward pass to compute the state x_i for $i = 0, 1, \dots, N$, followed by a back-propagation step to compute the co-state p_i for $i = N, N - 1, \dots, 0$. By leveraging both x_n and p_n , we approximate the optimal updated loss incurred by fine-tuning individual layers. This process effectively computes the greedy one-step loss for each layer, guiding the fine-tuning to the layer that promises the most immediate improvement in performance. The updated loss after fine-tuning θ_n can be approximated by

$$\hat{J}_{n}(\theta_{n}) = G_{(n+1)}\left(x_{n+1} + \frac{1}{\hat{\rho}_{n}}p_{n+1}x_{n}^{\top}x_{n}\right)$$
(8)

243 We'll justify this approximation in Section 3. Here $\frac{1}{\hat{\rho}_n}$ acts similarly to the learning rate, and we 244 aim to provide a well-estimated value of the optimal learning rate at an appropriate scale. Notably, 245 the optimal learning rates can vary significantly across different scenarios, even within the same 246 network, where different layers may require distinct values. The accuracy of $\hat{\rho}_n$ plays a crucial role 247 in estimating the updated loss.

Thus, we aim to provide a reasonably accurate estimate of $\hat{\rho}_n$ at the order-of-magnitude level to achieve a precise approximation of the optimal updated loss. We set $\hat{\rho}_n$ as defined in equation 9, computed using the state x_n and co-state p_n , with $r_n = p_{n+1}x_n^{\top}$ and d' being the output dimension. In practice, the terminal loss J may vary in scale. Thus, we sometimes modify $\hat{\rho}_n$ by multiplying it by a constant for all layers.

$$\hat{\rho}_n = \frac{d'}{2J} \cdot \frac{\|p_{n+1}^\top r_n x_n\|_F^2}{\|r_n\|_F^2} \tag{9}$$

We'll demonstrate in Section 3 that it approximates the optimal ρ_n^* in equation 12 to achieve the minimal updated loss within the deep linear network setting. To guide our layer selection process, we utilize $\hat{J}_n(\theta_n)$ in equation 8 as our metric for layer selection, opting to select the layer characterized by the minimal approximated loss. In other words, we select the layer of $g_m(\cdot, \theta_m)$ such that

$$m = \underset{n=0,1,\cdots,N-1}{\arg\min} J_n(\theta_n)$$

Layer Fine-tuning: Following the layer selection process, we utilize the Method of Successive Approximations (MSA) for fine-tuning the selected layer, with the primary objective of maximizing H_m with respect to θ_m . The MSA process is structured to enhance the optimization of the chosen parameters systematically.

In each iteration, we start with a feed-forward pass through the network to compute the state x_i for each layer, where *i* ranges from 0 to *N*, capturing the current output based on the input data. Once the state is computed, we proceed to a back-propagation step to derive the co-state p_i for each layer, starting from the last layer N and moving backward to layer 0. The co-state represents the sensitivity of the Hamiltonian with respect to the states, providing valuable information for optimization. Next, we compute the Hamiltonian specifically for the layer with $g_m(\cdot, \theta_m)$ using both the state and co-state variables, aiming to maximize H_m over the parameters θ_m of the selected layer. By repeating these steps iteratively, we progressively refine the layer parameters, converging toward a desired solution that enhances model performance.

Additionally, we have the flexibility to employ alternative optimization algorithms, such as Adam, during this process, which allows us to explore various strategies.

The methodology outlined is visually depicted in Figure 2, offering a comprehensive illustration of the layer-wise fine-tuning process. In Section 3, we will provide a detailed rationale and justification for our chosen metric utilized in the selection of layers.

281 282 283

293

294 295

299 300

301 302

303

304

305

306 307

313

314 315

316

323

3 Theoretical Analysis

In this section, we undertake a theoretical examination of our methodology within the idealized framework of the deep linear network. Given that deep neural networks are composed of linear and activation layers, an analysis of the deep linear network serves as a valuable avenue for gaining insight into our approach. Previous analyses(Arora et al., 2018a;b; Cohen et al., 2023) have provided significant insights into the behavior and properties of deep linear networks, underscoring the importance of this simplified model in understanding more complex architectures.

For simplicity, we employ a simplified variant of the augmented Hamiltonian and consider the maximization step of the (n + 1)-th layer as follows:

$$\max_{\theta_n^*} p_{n+1} \cdot g_n(x_n, \theta_n^*) - \frac{1}{2} \rho_n \|\theta_n^* - \theta_n\|_2^2$$
(10)

Given a collection of K sample input-label pairs $\{x^i, y^i\}_{i=1}^K$, with the inputs $x = (x^1, x^2, \dots, x^K) \in \mathbb{R}^{d \times K}$ and the labels $y = (y^1, y^2, \dots, y^K) \in \mathbb{R}^{d' \times K}$. Consider an N-layer deep linear network

 $x_{n+1} = g_n(x_n, \theta_n) = \theta_n x_n, n = 0, 1, \cdots, N - 1.$

with the input $x_0 = x$ and the loss function $J = \sum_{i=1}^{K} \Phi_i(x_N^i) = \frac{1}{2} \sum_{i=1}^{K} ||y^i - x_N^i||_2^2$.

Proposition 3.1 With given ρ_n , the updated loss after fine-tuning θ_n for one iteration is exactly $\hat{J}_n(\theta_n)$ in Equation 8, i.e.,

 $J^{update} = G_{(n+1)} \left(x_{n+1} + \frac{1}{\rho_n} p_{n+1} x_n^{\top} x_n \right)$ (11)

308 Due to space constraints, the proof details are provided in Appendix A.1.

For simplicity of expression, denote $\beta_n = \theta_{N-1} \cdots \theta_{n+1}$, and $r_n = p_{n+1} x_n^{\top} = \theta_n^{\top} \cdots \theta_{N-1}^{\top} (y - x_N) x_n^{\top}$ for $n = 0, 1, \cdots, N-1$. Below we show the relationship between the optimal ρ_n^* and our approximated $\hat{\rho}_n$.

Proposition 3.2 The optimal ρ_n^* to achieve the minimal updated loss is

$$\rho_n^* = \frac{\|\beta_n r_n x_n\|_F^2}{\|r_n\|_F^2} \tag{12}$$

and it satisfies $\rho_n^* \ge \frac{1}{d'} \hat{\rho}_n$ for the $\hat{\rho}_n$ determined in equation 9.

In addition, denote $\hat{\alpha}_n = \frac{1}{\hat{\rho}_n}$ and $\alpha_n^* = \frac{1}{\rho_n^*}$. Let θ be the 1-dimensional vectorization of all parameters. If $\theta \sim Uniform(B(0, r))$, $\forall r$, i.e., θ follows a uniform distribution in the neighborhood centered at the origin with radius r, we have $E_{\theta}\alpha_n^* = E_{\theta}\hat{\alpha}_n$, i.e., we provide an unbiased estimation for α_n^* , which functions similarly to a learning rate.

Due to space constraints, the proof details are provided in Appendix A.1.

324 4 **EXPERIMENTAL RESULTS** 325

In this section, we evaluate the performance of our proposed L-MSA method across various datasets. We compare L-MSA against established baseline methods to highlight its effectiveness in selecting 328 optimal layers and improving fine-tuning results. Further details about the datasets and the models are provided in Appendix A.2.

331 4.1 **BASELINE METHODS**

To compare with other baselines, we follow the setups from prior work(Lee et al., 2022). We employ 333 full fine-tuning as a baseline and focus on the comparison with layer-wise methods such as LIFT(Zhu 334 et al., 2023), LISA(Pan et al., 2024), and surgical fine-tuning(Lee et al., 2022). Among these methods, 335 surgical fine-tuning provides a metric, RGN, for selecting layers. We include a comparison between 336 our proposed metric and theirs to evaluate performance. 337

338 Full Fine-tuning is a widely used approach for adaptation. The model is initialized with pretrained weights and biases, and all parameters undergo gradient updates during fine-tuning. In our 339 experiments, we use the Adam optimizer to update all layers of the model. 340

341 LIFT(Zhu et al., 2023) is a layer-wise method where only one layer(or transform block) is updated in 342 each iteration. The selection policy for updating the layers can follow one of three strategies: (i) front 343 to end, (ii) end to front, or (iii) random. In our experiments, we test all three strategies and report the 344 average performance.

345 LISA(Pan et al., 2024) applies the idea of importance sampling to different layers in LLMs and 346 randomly freezes most middle layers during optimization. LISA consistently fine-tunes the first and 347 last layers, while updating each middle layer with a fixed probability. 348

Surgical Fine-tuning (Lee et al., 2022) shows that selectively fine-tuning a subset of layers matches 349 or outperforms commonly used fine-tuning approaches. The authors propose two criteria for au-350 tomatically selecting which layers to freeze, with the Relative Gradient Norm (RGN), defined as 351 $RGN = \frac{\|g\|_2}{\|\theta\|_2}$, showing better performance according to their findings. We compare our metric with 352 RGN and also evaluate the performance of our L-MSA method against Auto-RGN, which fine-tunes 353 the layer selected based on the highest RGN value. 354

355 356

326

327

330

332

4.2 EFFECTIVENESS OF OUR METRIC

357 We first conducted experiments to compute our proposed metric, the approximated optimal updated 358 loss J_n , and compared it with the true loss after training. In the case of pre-training on ImageNet and 359 fine-tuning on CIFAR-100, represented on the left side of Figure 4, the later layers exhibit smaller 360 approximated updated losses J_n . 361

Conversely, when pre-training is done on CIFAR-100 and fine-tuning is applied to a transformed 362 version of the dataset, shown on the right side of Figure 4, the earlier layers show smaller approximated 363 updated losses J_n . This transformed dataset is created by applying element-wise multiplication to 364 the input data of CIFAR-100 with a matrix, where each entry is e^x , and x follows a standard normal 365 distribution. These findings align with the actual training results shown in 1. 366



375 Figure 4: Comparison of our L-MSA metric and RGN with the true training loss. Due to differences in scale, where smaller values are preferred for both our metric and loss while larger values are preferred for RGN, all 376 values are normalized. A darker color indicates a better metric, suggesting that the corresponding layer will be 377 selected for fine-tuning.

We present the comparison of our L-MSA metric and RGN with the true training loss in 4. The results illustrate that our L-MSA metric consistently identifies layers associated with improved training loss, effectively pinpointing those that contribute to better training outcomes. However, in these two cases, RGN assigned the highest metric to the fifth layer, yet it was unable to assist in selecting the more effective layers for fine-tuning.

We also evaluated our metric on four real-data tasks: CIFAR-C, CIFAR-Flip, Living-17, and ImageNet-C. Due to space constraints, the results are provided in Appendix A.3, while the results of fine-tuning the selected layers using our L-MSA method are presented in Section 4.3.

4.3 FINE-TUNING RESULTS

We present the results of our L-MSA method in Figure 5, comparing it with Auto-RGN and full fine-tuning for DeiT models fine-tuning from ImageNet to CIFAR-100 and from CIFAR-100 to a transformed dataset. In the case of Auto-RGN, the layer selected by the RGN metric is updated using the Adam optimizer.



Figure 5: The Performance of L-MSA on DeiT-Tiny

In both scenarios, our findings show that the L-MSA method outperforms Auto-RGN and achieves performance comparable to full fine-tuning. Notably, using L-MSA for layer-wise fine-tuning results in performance improvements of up to 20% compared to full fine-tuning and up to 30% compared to Auto-RGN in the initial stages of training. Specifically, we observed a rapid decrease in training loss within the first few batches, underscoring the method's effectiveness, especially in cases where the amount of data is limited.

To further assess the performance of our L-MSA method, we evaluated the performance of our L-MSA method on four real-data tasks with a limited amount of data. For CIFAR-C(Hendrycks & Dietterich, 2019) and CIFAR-Flip(Lee et al., 2022), the models were pre-trained on CIFAR-10(Krizhevsky, 2009) using Wide ResNet-28-10(He et al., 2016). For Living-17(Santurkar et al., 2020) and ImageNet-C(Kar et al., 2022), the models were pre-trained on ImageNet(Deng et al., 2009a) using ResNet-50(He et al., 2016).

		CIEAD Ella	Lining 17	In a Net C	Assessed David
	CIFAR-C	CIFAR-Flip	Living-17	ImageNet-C	Average Rank
No Adaptation	60.3	0.0	73.2	18.1	-
Full Fine-tuning	81.1	86.2	78.2	49.0	2.5
LIFT	80.5	86.44	76.2	43.6	4.25
LISA	80.2	81.6	77.4	48.2	4.0
Auto-RGN	82.5	88.7	77.1	48.6	2.25
L-MSA	81.3	92.7	79.1	47.4	2.0

Table 1: We report the test accuracy on the target distribution across four real-data tasks. Our results show that L-MSA outperforms all other layer-wise fine-tuning methods, including Full Fine-tuning, LISA, LIFT, and Auto-RGN. The best-performing method for each distribution shift is highlighted in bold.

The results, presented in Table 1, compare L-MSA against other fine-tuning approaches, including Full Fine-tuning, LIFT, LISA, and Auto-RGN. Further details on the experimental setup can be found in the Appendix A.2.

432 The "No Adaptation" baseline provides a reference point for model performance without fine-tuning. 433 L-MSA consistently outperforms other methods, achieving the highest test accuracy on CIFAR-Flip 434 and Living-17, along with the best overall ranking across tasks. Notably, we also observe that L-MSA 435 achieves these results using fewer epochs. Auto-RGN proposed in surgical fine-tuning(Lee et al., 436 2022) also achieves a competitive average rank.

437 Overall, L-MSA's strong performance highlights its effectiveness in selecting layers for fine-tuning 438 and utilizing the MSA method to optimize the chosen layer during subsequent fine-tuning. The results 439 emphasize L-MSA's robustness and adaptability, demonstrating its ability to maintain high accuracy 440 across various types of distribution shifts.

441 442

443

4.4 EMPIRICAL ANALYSIS

To assess the effectiveness of our proposed 444 L-MSA method, we conducted an ablation 445 study comparing it against other fine-tuning ap-446 proaches, specifically (i) Full Fine-tuning, (ii) 447 Full fine-tuning using MSA, and (iii) L-MSA 448 Metric + Adam. This comprehensive compari-449 son aimed to evaluate not only the performance 450 of the L-MSA method but also to understand 451 how each approach influences model perfor-452 mance. The average test accuracies across four 453 datasets are plotted in Figure 6.



The results indicate that L-MSA significantly enhances performance compared to other fine-tuning

455 approaches. Notably, the Full Fine-tuning + MSA method underperforms because it optimizes each 456 layer's Hamiltonian individually for multiple steps, which is less effective in the context of full 457 fine-tuning. However, using only the L-MSA metric for layer-wise fine-tuning with Adam achieves 458 performance comparable to that of Full Fine-tuning, demonstrating the metric's effectiveness in layer 459 selection and the advantages of layer-wise fine-tuning. Furthermore, L-MSA outperforms the L-MSA Metric + Adam approach, emphasizing the benefits of utilizing MSA to optimize the selected layers. 460

461 462 463

454

5 LIMITATIONS AND FUTURE DIRECTIONS

464 While our layer-wise fine-tuning algorithm shows promising results, it is important to acknowledge 465 its limitations. Firstly, we select layers based on the approximated updated loss, which provides a good estimation of the training loss. However, this does not always guarantee strong generalization to 466 the test data. Additionally, while layer-wise fine-tuning reduces the computational burden compared 467 to full fine-tuning, it may still demand substantial computational resources due to performing both 468 forward and backward propagation, especially in large-scale models. 469

470 Future work could explore periodically reselecting layers and adjusting the training configuration 471 after a certain training period, allowing for continuous optimization and more efficient resource use, potentially enhancing performance over time. 472

473

475

474 6 RELATED WORK

476 6.1 LARGE-SCALE MODELS 477

The emergence of large-scale models has revolutionized various domains, ranging from natural lan-478 guage processing to computer vision. These models, characterized by their extensive parameterization 479 and sophisticated architectures, have demonstrated remarkable capabilities in capturing complex 480 patterns and representations from vast amounts of data. 481

482 In natural language processing, models like BERT(Devlin et al., 2018) and GPT(Radford et al., 483 2018) have set new benchmarks in a variety of tasks, such as language understanding and generation. By leveraging vast text corpora, these models learn rich semantic representations, excelling in 484 various downstream tasks. Similarly, in computer vision, models like ResNet(He et al., 2016) and 485 EfficientNet(Tan & Le, 2019) have demonstrated unprecedented performance in image classification,

object detection, and semantic segmentation tasks. By leveraging large datasets like ImageNet(Deng et al., 2009b), these models learn hierarchical features essential for understanding visual content.

Despite their impressive performance, these models pose significant computational challenges, particularly due to high training costs. Addressing these issues is a key research focus, with ongoing efforts aimed at developing more efficient techniques for both training and inference.

492 493

6.2 PARAMETER-EFFICIENT FINE-TUNING

494 Parameter-efficient fine-tuning (PEFT) techniques are designed to adapt pre-trained models by
495 selectively fine-tuning only a subset of parameters. In general, PEFT methods can be categorized into
496 three classes:

497
498
498
498
499
499
499
499
499
500
501
498
497
498
499
499
499
499
490
490
490
491
491
492
493
494
494
495
495
496
497
498
499
498
499
499
499
490
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400
400

Adapter methods typically introduce an auxiliary module with much fewer parameters than the original model. During training, updates are exclusively applied to the adapter module, allowing for more efficient parameter fine-tuning(Diao et al., 2023; Houlsby et al., 2019; Hu et al., 2021). These approaches require manual design and many of them also do not effectively reduce back-propagation costs.

Selective methods focus on the optimization of a subset of the model's parameters without the addition of extra modules. For instance, Exclusively fine-tuning bias terms can yield competitive performance comparable to fine-tuning the entire model(Zaken et al., 2021). Recently several noteworthy techniques have been developed in this area, particularly through the concept of layer freezing(Li et al., 2024; Liu et al., 2021). Compared with previous ones, Selective methods are more closely related to our approach.

513

514 6.3 TRANSFER LEARNING 515

Previous research in transfer learning has extensively explored the efficacy of fine-tuning to adapt 516 pre-trained features to a target distribution(Oquab et al., 2014; Sharif Razavian et al., 2014; Yosinski 517 et al., 2014). To maintain the valuable information obtained during pre-training, numerous studies 518 have proposed various methods to regularize the fine-tuning process(Li et al., 2020; Shen et al., 519 2021; Zhang et al., 2020). These methods aim to strike a balance between retaining the learned 520 features from the pre-trained model and adapting to the new target domain, thus ensuring effective 521 knowledge transfer. Notably, several works have demonstrated that freezing certain parameters within 522 the pre-trained model can significantly reduce overfitting during fine-tuning(Kirkpatrick et al., 2017; 523 Lee et al., 2019).

Contrary to most of the prevailing approaches, our work presents a counterintuitive finding: performing fine-tuning on the early layers can yield superior performance in specific scenarios. This intriguing finding resonates with recent investigations in the field(Lee et al., 2022), further undermining the prevailing notion that fine-tuning endeavors should predominantly concentrate on later layers, which are assumed to be more intricately tied to task-specific features.

529 530

531

7 CONCLUSION

532 In conclusion, we have presented L-MSA, a novel layer-wise fine-tuning approach that integrates 533 a metric for layer selection with an optimization algorithm based on the Method of Successive 534 Approximations (MSA). This framework allows for efficient and targeted fine-tuning of individual layers, significantly enhancing model performance while reducing computational costs. Our experi-535 ments across various datasets and tasks validate the effectiveness of L-MSA, demonstrating that our 536 method consistently outperforms baseline techniques. By algorithmically determining which layers 537 to fine-tune, we provide a practical solution to the challenges posed by large-scale models. Overall, 538 our work advances the field of layer-wise fine-tuning, offering new insights into optimizing model training and setting the stage for future research in this area.

540 **Reproducibility Statement:** Source codes for our experiments are provided in the supplementary 541 materials of the paper. The details of our experimental settings and computational infrastructure are 542 given in Section 4 and the Appendix A.2. All datasets that we used in the paper are published, and 543 they are easy to find in the Internet.

544 **Ethics Statement:** Given the nature of the work, we do not foresee any negative societal and ethical impacts of our work. 546

548 REFERENCES

547

549

550

551

556

559

561

562

563

564

565

566

569

- Sanjeev Arora, Nadav Cohen, Noah Golowich, and Wei Hu. A convergence analysis of gradient descent for deep linear neural networks. arXiv preprint arXiv:1810.02281, 2018a.
- 552 Sanjeev Arora, Nadav Cohen, Noah Golowich, and Wei Hu. Optimization of deep linear neural networks: Theory and algorithms. In International Conference on Learning Representations, 553 2018b. 554
 - Felix L Chernousko and AA Lyubushin. Method of successive approximations for solution of optimal control problems. Optimal Control Applications and Methods, 3(2):101–114, 1982.
- 558 Nadav Cohen, Govind Menon, and Zsolt Veraszto. Deep linear networks for matrix completion-an infinite depth limit. SIAM Journal on Applied Dynamical Systems, 22(4):3208–3232, 2023.
 - J. Deng, W. Dong, R. Socher, L. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2009a.
 - Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. 2009 IEEE conference on computer vision and pattern recognition, 2009b.
- 567 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep 568 bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805, 2018.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep 570 bidirectional transformers for language understanding, 2019. 571
- 572 Shizhe Diao, Zhichao Huang, Ruijia Xu, Xuechun Li, Yong Lin, Xiao Zhou, and Tong Zhang. 573 Black-box prompt learning for pre-trained language models. arXiv preprint arXiv:2201.08531, 574 2022.
- Shizhe Diao, Tianyang Xu, Ruijia Xu, Jiawei Wang, and Tong Zhang. Mixture-of-domain-adapters: 576 Decoupling and injecting domain knowledge to pre-trained language models memories. arXiv 577 preprint arXiv:2306.05406, 2023. 578
- 579 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas 580 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale, 581 2021. 582
- 583 Karen Hambardzumyan, Hrant Khachatrian, and Jonathan May. Warp: Word-level adversarial 584 reprogramming. arXiv preprint arXiv:2101.00121, 2021. 585
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image 586 recognition. Proceedings of the IEEE conference on computer vision and pattern recognition, 587 2016. 588
- 589 Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common 590 corruptions and perturbations. arXiv preprint arXiv:1903.12261, 2019. 591
- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, 592 Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp. In International conference on machine learning, pp. 2790–2799. PMLR, 2019.

605

612

618

639

594	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang	2,
595	and Weizhu Chen. Lora: Low-rank adaptation of large language models. arXiv preprint	t
596	arXiv:2106.09685, 2021.	
597		

- Christoph Käding, Erik Rodner, Alexander Freytag, and Joachim Denzler. Fine-tuning deep neural networks in continuous learning scenarios. In *Computer Vision–ACCV 2016 Workshops: ACCV 2016 International Workshops, Taipei, Taiwan, November 20-24, 2016, Revised Selected Papers, Part III 13*, pp. 588–605. Springer, 2017.
- Oğuzhan Fatih Kar, Teresa Yeo, Andrei Atanov, and Amir Zamir. 3d common corruptions and data augmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 18963–18974, 2022.
- 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.
- Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical Report TR-2009-08, University of Toronto, April 2009.
- Ananya Kumar, Aditi Raghunathan, Robbie Jones, Tengyu Ma, and Percy Liang. Fine-tuning can distort pretrained features and underperform out-of-distribution. *arXiv preprint arXiv:2202.10054*, 2022.
- Jaejun Lee, Raphael Tang, and Jimmy Lin. What would elsa do? freezing layers during transformer
 fine-tuning. *arXiv preprint arXiv:1911.03090*, 2019.
- Yoonho Lee, Annie S Chen, Fahim Tajwar, Ananya Kumar, Huaxiu Yao, Percy Liang, and
 Chelsea Finn. Surgical fine-tuning improves adaptation to distribution shifts. *arXiv preprint arXiv:2210.11466*, 2022.
- Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. *arXiv preprint arXiv:2104.08691*, 2021.
- Hao Li, Pratik Chaudhari, Hao Yang, Michael Lam, Avinash Ravichandran, Rahul Bhotika, and
 Stefano Soatto. Rethinking the hyperparameters for fine-tuning. *arXiv preprint arXiv:2002.11770*, 2020.
- Qianxiao Li, Long Chen, Cheng Tai, and E Weinan. Maximum principle based algorithms for deep learning. *Journal of Machine Learning Research*, 18(165):1–29, 2018.
- Sheng Li, Geng Yuan, Yue Dai, Youtao Zhang, Yanzhi Wang, and Xulong Tang. Smartfrz: An
 efficient training framework using attention-based layer freezing. *arXiv preprint arXiv:2401.16720*, 2024.
- Kiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang. Gpt understands, too. *AI Open*, 2023.
- Yuhan Liu, Saurabh Agarwal, and Shivaram Venkataraman. Autofreeze: Automatically freezing
 model blocks to accelerate fine-tuning. *arXiv preprint arXiv:2102.01386*, 2021.
- Maxime Oquab, Leon Bottou, Ivan Laptev, and Josef Sivic. Learning and transferring mid-level
 image representations using convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1717–1724, 2014.
- Rui Pan, Xiang Liu, Shizhe Diao, Renjie Pi, Jipeng Zhang, Chi Han, and Tong Zhang. Lisa: Layerwise importance sampling for memory-efficient large language model fine-tuning. *arXiv* preprint arXiv:2403.17919, 2024.
- 647 L. S. Pontryagin, V. G. Boltyanskij, R. V. Gamkrelidze, and E. F. Mishchenko. *The Mathematical Theory of Optimal Processes*. Wiley, New York, 1962.

648 649	Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language under- standing by generative pre-training. <i>OpenAI Blog</i> , 2018.			
651	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi			
652	Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text			
653	transformer. Journal of machine learning research, 21(140):1–67, 2020.			
654	Shibani Santurkar, Dimitris Tsipras, and Aleksander Madry. Breeds: Benchmarks for subpopulation			
655	shift. arXiv preprint arXiv:2008.04859, 2020.			
656	Ali Sharif Razavian Hossein Aziznour Josenhine Sullivan and Stefan Carlsson. Cnn features			
658 650	off-the-shelf: an astounding baseline for recognition. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition workshops</i> , pp. 806–813, 2014.			
660	Zhisiang Chan Zashan Lin Lin Oin Marine Semidae and Kauna Ting Chang. Destinis hatter than			
661 662	all: Revisiting fine-tuning strategy for few-shot learning. In <i>Proceedings of the AAAI conference</i> <i>on artificial intelligence</i> , volume 35, pp. 9594–9602, 2021.			
663				
664 665	Mingxing Tan and Quoc V Le. Efficientnet: Rethinking model scaling for convolutional neural networks. <i>International Conference on Machine Learning</i> , 2019.			
666	Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé			
667 668	Jégou. Training data-efficient image transformers & distillation through attention. In <i>International conference on machine learning</i> , pp. 10347–10357. PMLR, 2021.			
669	Jason Vosinski, Jeff Clune, Voshua Bengio, and Hod Linson. How transferable are features in deep			
670	neural networks? Advances in neural information processing systems, 27, 2014.			
671				
672	Elad Ben Zaken, Shauli Ravfogel, and Yoav Goldberg. Bitht: Simple parameter-efficient fine-tuning for transformer based masked language models. arXiv praprint arXiv:2106.10100.2021			
674	for transformer-based masked language-models. <i>urxiv preprint urxiv.2100.10199</i> , 2021.			
675	Jeffrey O Zhang, Alexander Sax, Amir Zamir, Leonidas Guibas, and Jitendra Malik. Side-tuning: a			
676 677	baseline for network adaptation via additive side networks. In <i>Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part III 16</i> , pp. 698–714.			
678	Springer, 2020.			
679 680	Ligeng Zhu, Lanxiang Hu, Ji Lin, and Song Han. Lift: Efficient layer-wise fine-tuning for large models. 2023.			
681				
682				
683				
684				
685				
686				
687				
680				
690				
691				
692				
693				
694				
695				
696				
697				
698				
099 700				
701				

Supplement to "L-MSA: Layer-wise Fine-tuning using the Method of Successive Approximations"

Table of Contents

A	Appendix					
	A.1	Theore	etical Analysis	14		
	A.2	Experi	mental Setup	16		
		A.2.1	Datasets	16		
		A.2.2	Models	16		
		A.2.3	Pretraining	17		
		A.2.4	Training Details	17		
		A.2.5	Environment	17		
	A.3	Additi	onal Experimental Results	17		

A APPENDIX

A.1 THEORETICAL ANALYSIS

In this section, we provide detailed proofs for propositions in Section 3. For simplicity, we consider the case where the batch size is 1. similar results can be derived for a larger batch size.

Proposition A.1 With given ρ_n , the updated loss after fine-tuning θ_n for one iteration is exactly $\hat{J}_n(\theta_n)$ in equation 8.

Proof: Given an input-label pair $\{x_0, y\}$, with the inputs $x_0 \in \mathbb{R}^d$ and the labels $y \in \mathbb{R}^{d'}$. Consider an *N*-layer deep linear network

$$x_{n+1} = g_n(x_n, \theta_n) = \theta_n x_n, n = 0, 1, \cdots, N-1.$$

with the input $x_0 = x$ and the loss function $J = \frac{1}{2} ||y - x_N||_2^2$.

Then we could compute the co-state

$$p_n = \theta_n^\top \cdots \theta_{N-1}^\top (y - x_N), \quad n = 0, 1, \cdots, N.$$

In every iteration, our objective is to maximize the augmented Hamiltonian for a single layer, as shown in equation 10. By taking the derivative of the Hamiltonian with respect to θ_n^* and setting it to zero, we can derive the updated parameters θ_n^* as follows.

$$\theta_n^* = \theta_n + \frac{1}{\rho_n} p_{n+1} x_n^\top = \theta_n + \frac{1}{\rho_n} \theta_{n+1}^\top \cdots \theta_{N-1}^\top (y - x_N) x_n^\top$$
(13)

which demonstrates that the E-MSA method using the simplified augmented Hamiltonian is equivalent to the gradient descent method with a learning rate of $\alpha_n = \frac{1}{\rho_n}$.

Then we can derive the updated loss as follows.

$$J^{update} = \frac{1}{2} \|\beta_n \theta_n^* x_n - y\|_2^2 = G_{(n+1)} \left(\theta_n^* x_n\right)$$

= $G_{(n+1)} \left(x_{n+1} + \frac{1}{\rho_n} p_{n+1} x_n^\top x_n \right)$ (14)

which is exactly $\hat{J}_n(\theta_n)$ in equation 8.

Proposition A.2 The optimal ρ_n^* to achieve the minimal updated loss is

$$\rho_n^* = \frac{\|\beta_n r_n x_n\|_F^2}{\|r_n\|_F^2} \tag{15}$$

and it satisfies $\rho_n^* \geq \frac{1}{d'}\hat{\rho}_n$ for the $\hat{\rho}_n$ in equation 9.

In addition, denote $\hat{\alpha}_n = \frac{1}{\hat{\rho}_n}$ and $\alpha_n^* = \frac{1}{\rho_n^*}$. Let θ be the 1-dimensional vectorization of all parameters. If $\theta \sim Uniform(B(0, r))$, $\forall r$, i.e., θ follows a uniform distribution in the neighborhood centered at the origin with radius r, we have $E_{\theta}\alpha_n^* = E_{\theta}\hat{\alpha}_n$.

Proof: To get the minimum loss that can be achieved after one iteration, we need to compute the optimal ρ_n . Let $\beta_n = \theta_{N-1} \cdots \theta_{n+1}$, and $r_n = p_{n+1} x_n^\top = \theta_{n+1}^\top \cdots \theta_{N-1}^\top (y - x_N) x_n^\top$ for $n = 0, 1, \cdots, N-1$. Thus, we have $p_{n+1} = \beta_n^\top (y - x_N)$ and $\theta_n^* = \theta_n + \alpha r_n$.

$$J^{update} = J_n(\theta_n^*) = \frac{1}{2} \|\beta_n \theta_n^* x_n - y\|_2^2$$

= $\frac{1}{2} \|\beta_n \theta_n x_n - y + \alpha_n \beta_n r_n x_n\|_2^2$ (16)

By taking the derivative with respect to α and setting it to zero, we get that the optimal learning rate α_n^* satisfies

$$0 = x_n^\top r_n^\top \beta_n^\top (\beta_n \theta_n x_n - y + \alpha_n^* \beta_n r_n x_n)$$
(17)

Therefore, to achieve the minimum loss, we should set

$$\alpha_n^* = \frac{tr(r_n^\top r_n)}{tr(x_n^\top r_n^\top \beta_n^\top \beta_n r_n x_n)} = \frac{\|r_n\|_F^2}{\|\beta_n r_n x_n\|_2^2}$$
(18)

where $\|\cdot\|_F$ is the Frobenius norm. Thus, we can derive that the optimal ρ_n^*

$$\rho_n^* = \frac{1}{\alpha_n^*} = \frac{\|\beta_n r_n x_n\|_2^2}{\|r_n\|_F^2} \ge \frac{\|(y - x_N)^\top \beta_n r_n x_n\|_2^2}{\|(y - x_N)\|_2^2 \cdot \|r_n\|_F^2} = \frac{1}{2J} \cdot \frac{\|p_{n+1}^\top r_n x_n\|_F^2}{\|r_n\|_F^2} = \frac{1}{d'}\hat{\rho}_n$$
(19)

Consider

$$\gamma = \frac{\|(y - x_N)^\top \beta_n r_n x_n\|_2^2}{\|y - x_N\|_2^2 \|\beta_n r_n x_n\|_2^2} = \frac{\|(y - x_N)^\top \beta_n \beta_n^\top (y - x_N) x_n^\top x_n\|_2^2}{\|y - x_N\|_2^2 \|\beta_n \beta_n^\top (y - x_N) x_n^\top x_n\|_2^2} = \frac{\|(y - x_N)^\top \beta_n \beta_n^\top (y - x_N)\|_2^2}{\|y - x_N\|_2^2 \|\beta_n \beta_n^\top (y - x_N)\|_2^2}$$
(20)

To evaluate γ , without loss of generality, we can set $||y - x_N||_2 = 1$. There exists an orthogonal matrix $Q \in \mathbb{R}^{d' \times d'}$ such that $y - x_N = Q \cdot e_1$, where $e_1 = (1, 0, \dots, 0)^\top \in \mathbb{R}^{d'}$. We denote $Q^\top \beta_n \beta_n^\top Q = (b_{ij})_{1 \le i,j \le d'}$. Then

$$\gamma = \frac{\|e_1^\top Q^\top \beta_n \beta_n^\top Q e_1\|_2^2}{\|\beta_n \beta_n^\top Q e_1\|_2^2} = \frac{\|e_1^\top Q^\top \beta_n \beta_n^\top Q e_1\|_2^2}{\|Q^\top \beta_n \beta_n^\top Q e_1\|_2^2}$$

$$= \frac{b_{11}^2}{b_{11}^2 + b_{21}^2 + \dots + b_{d'1}^2}$$
(21)

 $\forall r, \text{ if } \theta \sim \text{Uniform}(B(0, r)), \text{ i.e., } \theta \text{ follows a uniform distribution in the neighborhood centered at the origin with radius <math>r$, we have $E_{\theta}\gamma = \frac{1}{d'}$. As for the expectation of α^* , we have

$$\alpha_n^* = \frac{\|r_n\|_F^2}{\|\beta_n r_n x_n\|_2^2} = \gamma \cdot \frac{\|(y - x_N)\|_2^2 \cdot \|r_n\|_F^2}{\|(y - x_N)^\top \beta_n r_n x_n\|_2^2}$$

$$= \gamma \cdot 2J \cdot \frac{\|r_n\|_F^2}{\|p_{n+1}^\top r_n x_n\|_F^2}$$
(22)

According to our definition

$$\hat{\alpha}_n = \frac{1}{\hat{\rho}_n} = \frac{2J}{d'} \cdot \frac{\|r_n\|_F^2}{\|p_{n+1}^T r_n x_n\|_F^2}$$
(23)

Therefore, we can get the result that

 $E_{\theta}\alpha_n^* = E_{\theta}\hat{\alpha}_n$

 $= \frac{\|x_{1} - x_{N}\|_{2}^{2}}{\|y - x_{N}\|_{2}^{2}} \|\beta_{n}\|$ of generality, we can set $\|y - x_{N}\|$

810 A.2 EXPERIMENTAL SETUP

In this section, we elaborate on the experimental setup used in our study, covering various aspectssuch as datasets, models, pretraining, training details, and environment.

814

815 A.2.1 DATASETS 816

817 CIFAR-100(Krizhevsky, 2009): A widely used benchmark dataset consisting of 100 classes, with
600 images per class. The dataset is split into 50,000 training images and 10,000 test images, with
each image being a 32x32 RGB image.

CIFAR-10-C(Hendrycks & Dietterich, 2019): A corrupted version of CIFAR-10, containing 10
 classes with various common corruptions applied to the original test set, making it ideal for evaluating
 model robustness under real-world noisy conditions. The dataset has 5 levels of severity, and we
 evaluate with the most severe level. We tune on 1000 images from CIFAR-10-C.

CIFAR-Flip(Lee et al., 2022): A synthetic task where the inputs are from the original CIFAR-10 dataset, but the target labels are flipped such that each label y is transformed to 9 - y (e.g., label 0 becomes label 9, label 1 becomes label 8, etc.). This task provides a controlled setting to assess model performance when faced with adversarially shifted labels.

Living-17(Santurkar et al., 2020): A specialized dataset for classifying living organisms, including animals and plants, with 17 distinct classes. It presents a challenging test case due to its domain-specific nature, requiring the model to differentiate between fine-grained categories of organisms. For Living-17, we tune on 850 images from the target distribution, evenly split between the 17 classes, giving 50 images per class.

ImageNet-C(Kar et al., 2022): A corrupted version of the ImageNet dataset, where common corruptions such as Gaussian noise, blur, and weather distortions have been applied to the validation set. Similar to CIFAR-10-C, the dataset has 5 levels of severity, and we evaluate with the most severe level. We tune on 5000 images from ImageNet-C, evenly split between classes, giving 5 corrupted images per class.

838 839

840

A.2.2 MODELS

Here presents the models utilized in our experiments, highlighting their architectural characteristics
and parameter counts. We examine the Tiny version of DeiT(Touvron et al., 2021), Wide ResNet-2810, and ResNet-50(He et al., 2016), each selected for its unique strengths in handling various image
classification tasks. The number of parameters for each model is given in Table 2.

Number of Parameters
5.7 million
36 million
25.6 million

- 850 851
- 852

853

 Table 2: Number of Parameters for Various Models

DeiT (Data-efficient Image Transformers - Tiny)(Touvron et al., 2021): The Tiny version of DeiT
 is a transformer-based model designed for image classification. It utilizes attention mechanisms to
 capture long-range dependencies in images and achieves high performance with less data through
 efficient training strategies. In our experiments, it is used on the CIFAR-100 dataset for its ability to
 generalize well across diverse visual features.

Wide ResNet-28-10(He et al., 2016): Wide ResNet-28-10 is a convolutional neural network (CNN)
characterized by its wider architecture, which incorporates residual connections to alleviate the
vanishing gradient problem. With 28 layers and a width factor of 10, it balances capacity and
efficiency, making it suitable for various image classification tasks. We apply Wide ResNet-28-10 to
the CIFAR-10-C and CIFAR-Flip datasets to assess its performance under corrupted and adversarially
shifted labels.

ResNet-50(He et al., 2016): ResNet-50 is a deeper CNN with 50 layers, providing greater capacity
 for learning complex patterns. It is well-regarded for its robustness in both general and domain specific challenges. In our study, ResNet-50 is employed on the Living-17 and ImageNet-C datasets,
 leveraging its strong feature extraction capabilities against corrupted data.

69 A.2.3 PRETRAINING

868

For the CIFAR-100 dataset, we utilize a DeiT-Tiny model pre-trained on ImageNet. For the CIFAR-10-C and CIFAR-Flip datasets, we employ Wide ResNet-28-10 models that have been pre-trained on CIFAR-10. For the ImageNet-C dataset, we use a ResNet-50 model pre-trained on ImageNet.

For the Living-17 dataset, we load a ResNet model pre-trained on ImageNet and train it on the source data for 5 epochs. Initially, we tune only the head for 3 epochs, followed by fine-tuning all layers for an additional 2 epochs, following the Linear Probing then Fine-tuning(LP-FT) strategy(Kumar et al., 2022), while utilizing the Adam optimizer.

878 A.2.4 TRAINING DETAILS

Layer Configuration: For the DeiT model, each transformer block is treated as a layer, resulting in
 a total of 12 layers. In the case of Wide ResNet-28-10 and ResNet-50, each block is considered a
 layer, with the input convolutional layer (conv1) combined into the first layer and the output fully
 connected layer treated as a separate layer. Consequently, Wide ResNet-28-10 consists of 3 blocks
 plus the output layer, totaling 4 layers, while ResNet-50 comprises 4 blocks and the output layer,
 resulting in 5 layers.

Fine-Tuning with Adam Optimizer: For all baseline methods, we fine-tune the model on the labeled target data for a total of 10 epochs, with a batch size of 64. We explore 3 different learning rates: 1×10^{-3} , 1×10^{-4} , and 1×10^{-5} for all methods, except that for CIFAR-Flip, we adjust the learning rates to 1×10^{-1} , 1×10^{-2} , and 1×10^{-3} for last-layer fine-tuning.

Fine-Tuning with L-MSA: We fine-tune the model on the labeled target data for 10 epochs for CIFAR-100 and 5 epochs for other datasets, using a batch size of 64 and setting the hyper-parameter ρ to 1. To optimize the maximization Hamiltonian, we utilize the Adam optimizer with a learning rate of 1×10^{-4} . The optimization process runs for a total of 5 iterations.

Validation Stategy: For all datasets and experiments, we implement early stopping based on the accuracy observed on a held-out validation subset of the labeled target data. We report the test accuracy corresponding to the epoch with the highest validation accuracy.

A.2.5 ENVIRONMENT

The experiments were conducted on a machine equipped with an NVIDIA GeForce RTX 3090 GPU, which has 24GB of memory.

By providing comprehensive details of our experimental setup, we aim to facilitate transparency
 and reproducibility in our study, enabling other researchers to validate and build upon our findings
 accurately.

905

899

906 A.3 Additional Experimental Results

We present the comparison of our L-MSA metric and RGN with the true training loss on four real-world tasks: CIFAR-C, CIFAR-Flip, Living-17, and ImageNet-C. This comparison, omitted from Section 4.2 due to space limitations, is now provided here. As noted previously, due to differences in scale, where smaller values are preferred for both our metric and loss while larger values are preferred for RGN, all values are normalized. A darker color indicates a better metric, suggesting that the corresponding layer will be selected for fine-tuning.

Figure 7 illustrates that our L-MSA metric consistently identifies layers associated with improved training loss, successfully highlighting those that enhance training outcomes. In contrast, RGN often selects layers that do not match the optimal ones, making it less effective in comparison.

917 However, it is important to emphasize that our layer selection is primarily based on the approximated updated loss, which provides a solid estimate of the training loss. While this metric offers valuable



Figure 7: Comparison of L-MSA metric and RGN with the true training loss

insights for identifying effective layers, it is important to recognize that such an approach does not always guarantee robust generalization to the test data.