# Do we really have to filter out random noise in pre-training data for language models?

Anonymous ACL submission

#### Abstract

001 Web-scale pre-training datasets are the cornerstone of large language models' success. However, text data curated from the internet inevitably contains various types of noise, whose impact on language models needs to be understood. While existing research primarily focuses on low-quality or synthetic data, the ran-007 800 dom noise introduced by unregulated websites or crawler decoding errors has been largely overlooked. This paper investigates the in-011 fluence of such random noise and proposes strategies to mitigate its impact on down-012 stream tasks. Surprisingly, we observed that the performance degradation rate was significantly lower than the proportion of noise. We provide a theoretical justification for this phenomenon, which also elucidates the success of multilingual models and can be applied to other modalities. To address the adverse effects of 019 noise, we introduce a novel plug-and-play Local Gradient Matching loss, which explicitly enhances the denoising capability of the downstream task head by aligning the gradient of normal and perturbed features to improve local smoothness without requiring knowledge of the model's parameters. Extensive experiments on 027 8 language and 14 vision benchmarks validate the effectiveness of our proposed method.<sup>1</sup>

#### 1 Introduction

037

Large language models (LLMs), particularly the GPT series (Radford et al., 2019; Brown, 2020; OpenAI, 2023), have fundamentally transformed the research landscape in natural language processing. The remarkable performance and emergent capabilities of these autoregressive models (Ye and Gao, 2024; Nanda et al., 2023; Zheng et al., 2024) are largely attributed to pre-training on extensive datasets, which are gathered by crawling text from the whole internet. Given the sheer volume of these



Figure 1: Overview of the study and methodology. (a) The common scenario in which a GPT model, pretrained on filtered data  $P^c$ , demonstrates robust performance. (b) When the pre-training dataset is contaminated with random noise  $P^n$ , the resultant language model may exhibit unpredictable behavior. (c) Our approach focuses on the effective fine-tuning of black-box noisy models for downstream tasks  $P^d$ .

datasets, they inevitably encompass a wide variety of noisy data (Longpre et al., 2024; Elazar et al., 2024). Consequently, it is imperative to understand the impact of such noise, as the quality of the training data plays a decisive role in the effectiveness of LLMs (Touvron et al., 2023; Bai et al., 2023; Xie et al., 2023a). Significant research has been conducted in this area. Allen-Zhu and Li (2024a); Xie et al. (2023b) have highlighted that low-quality data can significantly diminish a model's knowledge capacity and performance. Furthermore, Shumailov et al. (2024); Seddik et al. (2024); Dohmatob et al. (2024) have provided empirical and theoretical evidence demonstrating that recursively training language models with synthetic data can lead to model

<sup>&</sup>lt;sup>1</sup>Code, data, and model checkpoint weights are available at https://anonymous.4open.science/r/lmn-acl-E9D3

collapse.

061

062

063

069

077

084

085

097

100

101

102

103

104

However, little attention has been paid to the impact of random noise within datasets. Due to anti-crawling mechanisms (Gao et al., 2023), mismatched encodings (e.g., UTF-8 and GBK), and tremendous amounts of unmaintained websites (Pletinckx et al., 2021), the raw data obtained through web crawling inevitably contains a substantial amount of unpredictable random noise (Zhou et al., 2024; Chen et al., 2022; Kang et al., 2023). Although theoretically it may not be challenging to remove such noise, practical limitations in computational resources often result in incomplete data cleaning (Albalak et al., 2024; Soldaini et al., 2024). For example, it is observed that the Chinese corpus used to train the GPT-40 tokenizer contained a considerable amount of nonsensical data<sup>2</sup>. Therefore, it is of great importance to gain a thorough understanding of the potential effects of random noise on language models, which will contribute to a deeper insight into the robustness of LLMs.

We conduct extensive experiments based on the OpenWebText dataset (Gokaslan et al., 2019) used to pre-train GPT-2. Specifically, we generate random noise with proportions of 1%, 5% and 20%, and subsequently concatenate these noise to the end of the clean data. The next-token prediction pre-training process then continues as usual. Interestingly, we observe that the presence of the random noise do not lead to a catastrophic failure in model training; instead, its effect on the autoregressive loss is disproportionately small, e.g., the increase in loss is only about 1% even with 20% of the dataset being noisy. We provide a theoretical analysis to explain these phenomena, which also sheds light on the success of multilingual models (where one language may appear as "noise" to another) and large speech models (Chen et al., 2022), indicating the broader implications of studying the effects of random noise.

Beyond the impact of noise on pre-training loss, it is also crucial to understand its effects in downstream tasks. Following Chen et al. (2024), we explore how to efficiently fine-tune language models using extracted features for downstream tasks when the pre-training data and model weights are not accessible, which reflects real-world application scenarios for large language models. To mitigate the adverse effects of noise, we propose a novel plug-and-play Local Gradient Matching (LGM)

<sup>2</sup>https://github.com/jiangyy/gpt-tokens

loss. This method involves artificially adding noise to the output features and minimizing the gradient difference between the noisy and original features. We theoretically prove that the gradient difference can be upper-bounded by local smoothness (Srebro et al., 2010), the value of the loss function and the input flatness. Comprehensive experiments on 22 vision and language understanding benchmarks further corroborate the effectiveness and robustness of our proposed method. 105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

In summary, our contributions are as follows: (1) We investigate the underexplored problem of random noise in pre-training datasets for language models. (2) We pre-train multiple GPT-2 models and the empirical results show that the influence of random noise is relatively insignificant. Then we provide a theoretical analysis, extending our findings to other domains and thus highlighting the significance of this research direction. (3) We propose a novel gradient matching loss for downstream tasks, supported by comprehensive experimental and theoretical analysis that confirm its efficacy.

The remaining part is arranged as follows. In Section 2, we summarize the related works concerning pre-training data for language models. In Section 3, we provide a detailed analysis of the impact of random noise on language models from both experimental and theoretical perspectives. To compensate for this impact, we introduce the LGM loss in Section 4 and demonstrate its effectiveness through extensive experiments and theoretical analysis. Finally, in Section 5, we conclude the paper.

#### 2 Related Works

Data Selection for Language Model Training. High-quality text corpora are essential for effective language models. Elazar et al. (2024) analyzed open-source datasets like The Pile and C4, uncovering significant amounts of low-quality content in these datasets. Allen-Zhu and Li (2024a); Seddik et al. (2024) highlighted the negative impact of such data on training. Thus, data selection is crucial for language models. Longpre et al. (2024) provides guidelines for selecting pre-training data, and Chai et al. (2024) evaluates the influence of individual samples on GPT model training dynamics. Li et al. (2024b) introduces a metric to evaluate instructiontuning data quality. Other works, including Xie et al. (2023b,a); Lee et al. (2023), focused on optimizing data composition. Despite these remarkable contributions, there remains a lack of understand-



Figure 2: Next-token prediction loss on the clean OpenWebText validation set for GPT-2 models pre-trained on synthetic OpenWebText datasets with varying levels of random noise. (a) Trend of NTP loss as training proceeds. (b) Difference in NTP loss between the noisy and clean models after the same number of training iterations. (c) Difference in loss values after undergoing the same number of training iterations on clean OpenWebText data.

ing regarding the specific effects of random noise on language model performance. This paper aims to address this gap.

155

156

157

158

159

160

162

163

164

165

167

168

169

170

171

172

173

174

175

176

177

179

181

184

186

188

189

190

192

Learning from Noisy Distributions. LLMs are trained on noisy datasets but evaluated in practical user settings under clean distribution, violating the i.i.d. assumption of traditional machine learning. Research has explored the impact of noisy labels on model performance (Song et al., 2022; Lukasik et al., 2020). For input feature noise, adding perturbations to images enhances model interpretability (Smilkov et al., 2017), while incorporating randomly generated samples into the training set helps alleviate class imbalance (Zada et al., 2022). These studies, however, mainly focus on vision modality and do not fully address the pre-training context of LLMs.

Noisy Model Learning. Our work draws significant inspiration from Noisy Model Learning (NML) proposed by Chen et al. (2024). In NML, the authors introduce noise into large datasets like ImageNet by randomly altering labels, then pretrain neural networks on these noisy datasets. The study reveals that moderate label noise enhances in-distribution (ID) sample classification, while outof-distribution (OOD) performance deteriorates with increasing noise. Focused on supervised learning in computer vision, NML modifies only the labels, leaving the image features intact. This contrasts with self-supervised learning in language models, where the text serves both as input and output, making it impossible to alter the labels without changing the corresponding inputs. This paper extends the concept of NML, presenting theoretical insights and methodologies that are applicable across multiple modalities and various problems.

Due to space limitations, the detailed related works are provided in Appendix C.

### **3** Revealing the Effect of Random Noise in Language Model Pre-training

193

194

195

196

197

198

199

200

201

202

203

204

206

208

209

210

211

212

213

214

215

216

217

218

219

220

221

223

224

225

226

In this section, we first pre-train multiple GPT-2 models on synthetic noisy OpenWebText corpus to investigate the impact of random noise in the pretraining data. We then provide a theoretical analysis of the results and validate our theory through experiments. Finally, we demonstrate that the insights gained from our investigation have broader applicability beyond the immediate scope of our study. The frequently used notation and their descriptions are shown in Appendix A.

#### 3.1 Experimental Design

**Preliminary.** Let L denote the maximum context length of the language model, and let  $\mathcal{W}$ represent the model's vocabulary with size V =We define  $\mathcal{X}$  as the set of all discrete  $|\mathcal{W}|.$ sentences that the model can represent, where  $\mathcal{X} = \bigcup_{i=1}^{L} \{0, 1, \dots, V-1\}^{i} = \bigcup_{i=1}^{L} \mathcal{W}^{i}$  and  $\{0, 1, \dots, V - 1\}^i$  represents prefixes of length *i*. For any discrete set A, let  $\Delta_A$  denote the set of all probability distributions defined on A. Given that next-token prediction (NTP) is actually a classification task given the prefix, we define joint probability distributions  $P^c, P^n, P^m \in \Delta_{\mathcal{X} \times \mathcal{W}}$  where  $P^c$  represents the distribution of clean data,  $P^n$ represents the distribution of noise data, and  $P^m$ represents the distribution of the mixed pre-training dataset which contains both clean and noise data. Since the noisy dataset can be viewed as the concatenation of clean data and random noise, it can be formalized by the Huber contamination model (Fang et al., 2022) as follows:

$$P^m = \alpha P^n + (1 - \alpha)P^c \tag{1}$$

3



Figure 3: Validation experiments. (a) Loss trends on the random noise in the *training set* of the model trained on the dataset with 5% random noise. (b) Comparison of the loss between 5% random noise and Gaussian noise. (c) The loss difference on the clean OpenWebText validation set compared to the baseline for models trained on datasets with 5% random noise and 5% Gaussian noise, respectively.

where we use  $\alpha$  to represent the noise proportion. An explanation of Equation (1) can be found in Appendix B.1. For any joint probability distribution  $P \in \Delta_{\mathcal{X} \times \mathcal{W}}$ , let  $P_X \in \Delta_{\mathcal{X}}$  and  $P_{\cdot|X} \in \Delta_{\mathcal{W}}$ represent the marginal and conditional distribution of P over  $\mathcal{X}$  and  $\mathcal{W}$ .

227

228

230

231

236

238

239

242

243

244

245

246

247

248

249

254

255

263

We use  $\mathcal{H}$  to denote the hypothesis space (e.g., all possible parameter configurations given the transformer architecture ). Define  $h : \mathcal{X} \to \mathbb{R}^V \in$  $\mathcal{H}$  as the language model and  $p^h_{\cdot|x}(w)$  as the *w*-th component of the probability distribution induced by h(x). The next-token prediction loss can be expressed as follows:

$$\mathcal{L}_{ntp}(P,h) = \mathbb{E}_{x \sim P_X} \mathbb{E}_{w \sim P_{\cdot|x}} \big[ -\log(\boldsymbol{p}^h_{\cdot|x}(w)) \big].$$
<sup>(2)</sup>

Experiment setup. We utilize the OpenWeb-Text dataset (Gokaslan et al., 2019) which comprises 8 billion tokens as an alternative to the original WebText dataset used for training GPT-2 124M models (Radford et al., 2019). Concretely, to simulate random noise in unfiltered web-crawled data, we first generate uniformly-distributed data and then increase the number of specific tokens to introduce variability. Finally, we shuffle the entire data to mimic the randomness and unpredictability of real-world web-crawled data. The generated noise is then added to the clean dataset such that  $\alpha$ is 1%, 5%, and 20% respectively. Each synthetic noisy dataset is used to pre-train a GPT-2 model. We set the context length L to be 1024 and the batch size to be 640. All models are trained for 300,000 iterations. To evaluate the performance, the resulting model checkpoints are tested on the clean OpenWebText validation set, measuring the NTP loss for comparison. Further details regarding datasets and experimental parameters can be found in Appendix D.

#### 3.2 Results

In Figure 2, we illustrate the evolution of the NTP loss throughout the training process. Although random noise has a negative effect on the model's performance as expected, experimental results yield two intriguing insights:

(1) In contrast to the low-quality or synthetic data, **the presence of random noise does not lead to training collapse**, even when the noise level reaches 20%. While increasing training time on low quality or synthetic data typically degrades model performance (Allen-Zhu and Li, 2024a; Shumailov et al., 2024), extending the training duration continues to drive down the model's loss in the case of random noise.

(2) The impact of random noise on the loss is disproportionately small. For instance, 5% of random noise only results in a 0.2% increase in the NTP loss. This discrepancy becomes even smaller if the noisy models are calibrated to match the number of training iterations with the baselines trained on clean datasets.

These positive experimental outcomes further corroborate the robustness of language models and provide insights into why pre-training on largescale datasets that inevitably contain significant amounts of noise can still yield high-performing models. These somewhat unexpected findings naturally prompt us to explore the underlying reasons.

#### 3.3 Theoretical Analysis

In the analysis below, we focus on the impact of random noise on NTP loss, as pre-training loss is crucial for the performance on downstream tasks (Saunshi et al., 2021; Wei et al., 2021; Liu et al., 2023; Zheng et al., 2023a). Specifically, we are interested in the difference of NTP Loss between

264

265

266

267



Figure 4: Loss and its difference across different types and levels of noise within the ArXiv and Wikipedia corpora.

a model  $h^*$  trained on a noise-free dataset and a model h trained with a noisy dataset. We begin by noting that sampling from the clean distribution should not yield random gibberish and vice versa. Mathematically, this implies that for any prefix rsampled from  $P_X^n$ , the probability under the clean distribution  $P_X^c(r)$  is zero. Thus, we make the following assumption:

301

302

304

310

313

315

316

317

319

321

324

325

329

331

333

336

337

340

**Assumption 1.**  $P^c$  and  $P^n$  have disjoint support sets, i.e.,  $supp(P^c) \cap supp(P^n) = \emptyset$ .

The subsequent proposition demonstrates that the error  $\epsilon$  introduced to the loss due to random noise is less than the proportion  $\alpha$  of random noise in the dataset.

**Proposition 1.** Under Assumption 1, let  $h^*$  be a model trained on  $P^c$ , with  $\mathcal{L}_{ntp}(P^c, h^*) =$  $-\log p_c$  and  $\mathcal{L}_{ntp}(P^n, h^*) = -\log p_n$ . When the model h is trained on a mixed distribution  $P^m$ which includes noise, it attempts to fit  $P^n$ , leading to an increase in the loss on the clean distribution  $P^c$ , such that  $\mathcal{L}_{ntp}(P^c, h) = -\log(p_c - \epsilon)$ and  $\mathcal{L}_{ntp}(P^n, h) = -\log(p_n + \epsilon/k)$  for some  $\epsilon > 0$  (k can be shown to be  $\Omega(e^{\mathcal{L}_{ntp}(P^n, h)})$ ). Let  $\eta = \alpha p_c - (1 - \alpha)kp_n$ . We arrive at the following conclusions:

(1) If  $\alpha \leq \frac{kp_n}{p_c+kp_n}$ , then for any  $0 < \epsilon < p_c$ , we have  $\mathcal{L}_{ntp}(P^m, h) \geq \mathcal{L}_{ntp}(P^m, h^*)$ . This means that when  $\alpha$  is sufficiently small, the global minimum on  $P^m$  will not be affected by noise.

(2) If  $\alpha > \frac{kp_n}{p_c + kp_n}$ , then for  $\epsilon < \eta$ , it holds that  $\mathcal{L}_{ntp}(P^m, h) < \mathcal{L}_{ntp}(P^m, h^*)$ . This suggests that if  $\alpha$  is large enough, the impact on the optimal hypothesis is at least as much as  $\alpha p_c - (1 - \alpha)kp_n$ . (3) When  $\alpha < \frac{1}{3}$  and  $k > \frac{\alpha(1 - 3\alpha)p_c}{(1 - \alpha)(2 - 3\alpha)p_n}$ , for  $\epsilon \ge 3\eta$  we get  $\mathcal{L}_{ntp}(P^m, h^*) < \mathcal{L}_{ntp}(P^m, h)$ . Similarly, it can be shown that  $\epsilon$  does not exceed  $2\eta$ when  $\alpha > \max\left(\frac{kp_n}{p_c + kp_n}, \frac{1}{2}\right)$  and  $k > \frac{(2\alpha - 1)p_c}{2(1 - \alpha)p_n}$ . This indicates that when k is sufficiently large, the effect of noise is at most  $\mathcal{O}(\alpha p_c - (1 - \alpha)kp_n)$ .

The proof can be found in Appendix B.2. Proposition 1 primarily investigates the performance gap between models trained on  $P^m$  and those on  $P^c$ . It is proved that when  $\alpha$  is small enough, the presence of noise has no impact on the optimal model on  $P^m$ . Even as  $\alpha$  approaches  $\frac{1}{3}$  or even  $\frac{1}{2}$ , as long as k is large enough (the analysis regarding k and other parameters is detailed in Appendix B.3), the loss induced by noise,  $\epsilon$ , does not exceed  $\mathcal{O}(\alpha p_c - (1 - \alpha)kp_n)$ . Given that k is much greater than 1, this implies  $\epsilon$  is much smaller than  $\alpha p_c$ . This explains the observed experimental results.

341

342

343

344

345

347

349

351

352

353

355

357

358

359

361

362

363

364

365

367

369

370

371

373

374

375

376

377

378

379

381

With these theoretical results in hand, we then conduct multiple experiments to substantiate their validity. First, we plot the trend of NTP loss on random noise within the training set throughout the learning process, as shown in Figure 3(a). It is evident that the loss on random noise decreases at a very slow rate, indicating that the model struggles to efficiently learn the distribution of random noise. This observation contrasts with previous findings that neural networks can easily fit random labels (Zhang et al., 2021a), which needs further investigation. Next, we add 5% Gaussian-distributed noise to the training dataset and compare the results with models trained on 5% random noise. As depicted in Figure 3(b), the loss on Gaussian noise is lower than that on the random noise. According to Proposition 1, since the Gaussian distribution corresponds to a high  $p_n$ , we can **predict** that a model trained on Gaussian noise will exhibit a lower loss on  $P^c$ . Figure 3(c) confirms our prediction, thus further validating the proportions.

#### 3.4 Experiments on Other Text Corpus

To further investigate the impact of random noise on model generalization, we evaluate the nexttoken prediction loss of the trained models on data crawled from arXiv and Wikipedia. The results are illustrated in Figure 4. Surprisingly, models trained with added noise outperformed those trained on  $P^c$ . This counterintuitive finding aligns with previous work in visual domains (Zada et al., 2022), suggest-

	SST-2		SST-fine		20newsgroup		CR		Avg	
	Linear	MLP	Linear	MLP	Linear	MLP	Linear	MLP	Linear	MLP
OpenAI's GPT-2*	87.4	/	49.2	/	63.7	/	86.8	/	71.75	/
0%	$86.71\pm0.85$	$87.36\pm0.33$	$49.19\pm0.32$	$49.18\pm0.02$	$63.12\pm0.37$	$62.70\pm0.86$	$85.65\pm0.88$	$84.86\pm0.36$	71.16	71.02
$0\% + \mathcal{L}_{gm}$	$\textbf{87.42} \pm \textbf{0.73}$	$\textbf{87.86} \pm \textbf{0.04}$	$\textbf{49.72} \pm \textbf{0.27}$	$\textbf{49.81} \pm \textbf{0.97}$	$\textbf{63.69} \pm \textbf{0.59}$	$\textbf{62.95} \pm \textbf{0.13}$	$\textbf{86.58} \pm \textbf{0.22}$	$\textbf{86.45} \pm \textbf{0.73}$	71.85	71.76
1%	$87.25\pm0.79$	$\textbf{87.53} \pm \textbf{0.27}$	$49.32\pm0.72$	$49.45\pm0.56$	$63.71\pm0.02$	$64.65\pm0.06$	$84.86 \pm 0.98$	$84.59\pm0.59$	71.28	71.55
$1\% + \mathcal{L}_{gm}$	$\textbf{87.64} \pm \textbf{0.91}$	$87.25\pm0.44$	$\textbf{49.59} \pm \textbf{0.73}$	$\textbf{50.01} \pm \textbf{0.05}$	$\textbf{63.92} \pm \textbf{0.65}$	$\textbf{64.72} \pm \textbf{0.76}$	$\textbf{85.12} \pm \textbf{0.07}$	$\textbf{85.25} \pm \textbf{0.29}$	71.56	71.80
5%	$86.92\pm0.98$	$87.23\pm0.41$	$49.04 \pm 0.11$	$\textbf{50.09} \pm \textbf{0.53}$	$63.27\pm0.79$	$62.09\pm0.28$	$85.30\pm0.63$	$\textbf{84.32} \pm \textbf{0.78}$	71.13	70.93
$5\% + \mathcal{L}_{gm}$	$\textbf{87.19} \pm \textbf{1.02}$	$\textbf{87.61} \pm \textbf{0.51}$	$\textbf{49.82} \pm \textbf{0.17}$	$48.95\pm0.89$	$\textbf{63.78} \pm \textbf{0.93}$	$\textbf{62.37} \pm \textbf{0.56}$	$\textbf{85.57} \pm \textbf{0.43}$	$84.19\pm0.69$	71.59	70.78
20%	$86.60 \pm 1.28$	$86.60\pm0.81$	$49.45\pm0.78$	$49.63\pm0.01$	$63.47\pm0.64$	$64.16\pm0.92$	$\textbf{85.32} \pm \textbf{0.60}$	$\textbf{85.45} \pm \textbf{0.86}$	71.26	71.26
$20\% + L_{gm}$	$\textbf{87.2} \pm \textbf{0.99}$	$\textbf{86.87} \pm \textbf{0.78}$	$\textbf{49.68} \pm \textbf{0.55}$	$\textbf{50.40} \pm \textbf{0.46}$	$\textbf{63.58} \pm \textbf{0.08}$	$\textbf{64.21} \pm \textbf{0.78}$	$85.25\pm0.90$	$85.52\pm0.24$	71.42	71.75
Gaussian	$85.22\pm0.24$	$86.82\pm0.72$	$46.15\pm0.51$	$49.59\pm0.76$	$63.72\pm0.35$	$\textbf{64.40} \pm \textbf{0.76}$	$84.06\pm0.74$	$\textbf{83.53} \pm \textbf{0.70}$	69.78	71.08
Gaussian + $\mathcal{L}_{gm}$	$85.94 \pm 0.55$	$\textbf{87.25} \pm \textbf{0.36}$	$  \hspace{.1cm} \textbf{48.23} \pm \textbf{0.69} \\ \\$	$\textbf{50.29} \pm \textbf{0.70}$	$\textbf{64.06} \pm \textbf{0.73}$	$64.29\pm0.94$	$\textbf{84.46} \pm \textbf{0.33}$	$83.29\pm0.47$	70.67	71.45

Table 1: Accuracy on 4 text classification benchmark. 0% represents a model trained on  $P^c$ , 1% and so on denote the proportion of random noise, and Gaussian refers to Gaussian noise. \* cited from Saunshi et al. (2021).

	BI	BBC		Balanced COPA		MRPC		WiC		/g
	Linear	MLP	Linear	MLP	Linear	MLP	Linear	MLP	Linear	MLP
Llama-3-8B	$96.90\pm0.40$	$97.50\pm0.20$	$69.00\pm0.20$	$\textbf{65.60} \pm \textbf{0.50}$	$72.00\pm0.81$	$67.53 \pm 0.93$	$64.14\pm0.56$	$59.07\pm0.34$	75.51	72.42
Llama-3-8B + $\mathcal{L}_{gm}$	$\textbf{98.00} \pm \textbf{0.50}$	$\textbf{98.20} \pm \textbf{0.40}$	$\textbf{70.80} \pm \textbf{1.70}$	$64.80\pm0.20$	$\textbf{74.89} \pm \textbf{0.40}$	$\textbf{74.14} \pm \textbf{1.49}$	$\textbf{64.71} \pm \textbf{0.94}$	$\textbf{64.21} \pm \textbf{0.83}$	77.10	75.33
Llama-3-8B-Instruct	$96.80\pm0.70$	$96.90\pm0.30$	$87.80\pm0.70$	$88.80 \pm 0.60$	$72.57\pm0.26$	$71.42\pm0.13$	$65.92\pm0.53$	$61.85\pm0.59$	80.77	79.74
Llama-3-8B-Instruct + $\mathcal{L}_{gm}$	$\textbf{97.70} \pm \textbf{0.20}$	$\textbf{97.80} \pm \textbf{0.40}$	$\textbf{88.40} \pm \textbf{0.90}$	$\textbf{89.60} \pm \textbf{0.50}$	$\textbf{77.79} \pm \textbf{0.58}$	$\textbf{76.81} \pm \textbf{0.20}$	$\textbf{68.64} \pm \textbf{0.26}$	$\textbf{67.71} \pm \textbf{0.51}$	83.13	82.98
Llama-3.2-3B-Instruct	$97.30\pm0.60$	$97.20\pm0.80$	$80.40\pm0.90$	$\textbf{79.60} \pm \textbf{0.20}$	$77.79\pm0.52$	$72.57\pm0.31$	$64.07\pm0.82$	$57.50\pm0.35$	79.89	76.71
Llama-3.2-3B-Instruct + $\mathcal{L}_{gm}$	$\textbf{97.60} \pm \textbf{0.10}$	$\textbf{97.80} \pm \textbf{0.30}$	$\textbf{81.60} \pm \textbf{1.00}$	$79.40\pm0.10$	$\textbf{78.43} \pm \textbf{0.78}$	$\textbf{76.57} \pm \textbf{1.12}$	$64.35 \pm 0.62$	$\textbf{62.64} \pm \textbf{0.07}$	80.49	79.10
Qwen2.5-1.5B-Instruct	$97.00\pm0.30$	$96.60\pm0.70$	$80.80\pm0.70$	$82.20\pm0.50$	$74.49\pm0.71$	$73.39\pm0.90$	$65.92\pm0.45$	$61.64\pm0.20$	79.55	78.45
Qwen2.5-1.5B-Instruct + $\mathcal{L}_{gm}$	$\textbf{97.40} \pm \textbf{0.10}$	$\textbf{97.20} \pm \textbf{0.80}$	$\textbf{84.00} \pm \textbf{0.90}$	$\textbf{83.40} \pm \textbf{0.30}$	$\textbf{79.65} \pm \textbf{0.62}$	$\textbf{78.37} \pm \textbf{0.84}$	$\textbf{67.71} \pm \textbf{0.49}$	$\textbf{66.92} \pm \textbf{0.55}$	82.19	81.47
Qwen2.5-7B-Instruct	$96.30\pm0.30$	$96.70\pm0.50$	$94.60\pm0.90$	$95.80\pm0.40$	$83.71\pm0.92$	$76.81\pm0.51$	$68.92 \pm 0.41$	$64.92\pm0.18$	85.88	83.55
Qwen2.5-7B-Instruct + $\mathcal{L}_{gm}$	$\textbf{97.10} \pm \textbf{0.80}$	$\textbf{97.40} \pm \textbf{0.20}$	$95.60\pm0.50$	$\textbf{96.00} \pm \textbf{0.80}$	$\textbf{84.98} \pm \textbf{0.12}$	$\textbf{83.13} \pm \textbf{0.49}$	$\textbf{72.28} \pm \textbf{0.98}$	$\textbf{70.14} \pm \textbf{0.94}$	87.49	86.66

Table 2: Accuracy of LLMs on 4 natural language understanding benchmark.

ing that incorporating random noise into training sets might enhance model robustness. Additionally, we observe that the performance of models subjected to Gaussian noise varies across different datasets. These observations warrant further investigation.

#### **3.5 Broader Impact of the Results**

382

384

388

394

398

400

401

402

403

404

405

406

407

408

409

410

411

In addition to providing explanations regarding the impact of random noise on pre-training language models, we aim to extend our proposed theory to other areas, therefore demonstrating the practical value of our research findings.

One immediate direction is the training of multilingual models (Pires et al., 2019; Chi et al., 2020; Yang et al., 2024). Clearly, tokens corresponding to different languages are distinct, and their distributions naturally satisfy Assumption 1. For example, in an English-French bilingual model, let  $P^c$  represent English and  $P^n$  represent French. Supposing the pre-training corpus consists of an equal distribution of English and French, and given that the two distributions are similar, we can assume that  $p_c \approx p_n$ , leading to  $\epsilon \approx 0$ . This provides a theoretical foundation for the success of multilingual models. See Appendix D.3 for more details.

Beyond language modality, random white noise has received increased attention in the speech domain (Chen et al., 2022, 2021; Yin et al., 2024). Since our theory applies to any cross-entropy-like loss, it can also explain why speech models pretrained on very noisy large-scale datasets, such as Gigaspeech (Chen et al., 2021), which contain significant background noise and prolonged silence at the beginning and end of a few samples, still perform remarkably well.

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

## 4 Reducing the Noise with Local Gradient Matching

In this section, we introduce a novel black-box finetuning method termed Local Gradient Matching loss. Extensive experiments across 8 natural language understanding and 14 image classification benchmark datasets demonstrate that the proposed method consistently enhances performance across different backbones and modalities. Theoretical analysis reveals that LGM is effective because minimizing it improves the local smoothness of the cross-entropy loss and reduces its value, thereby enhancing the robustness and efficiency of the model.

#### 4.1 Method

In the preceding analysis, we demonstrate that the population-level loss function is only marginally affected by random noise. However, during the SGD training process, its presence introduces certain noise into the gradients. Although stochastic gradient noise is crucial for the generalization of deep networks (Barrett and Dherin, 2021; HaoChen et al., 2021), prior studies (Chen et al., 2023; Xie et al., 2021c) have shown that artificially added noise can hurt the model's generalization. There-

Model	Efficien	tNet-B3	ResNetv	2-152x2	Swi	n-L	ConvN	lext-L	Vil	I-L
Pre-training Data	JFT-3	800M	ImageN	Vet-21K	ImageN	et-21K	Laior	1-2B	Laio	1-2B
Fine-tuning Method	Linear	MLP	Linear	MLP	Linear	MLP	Linear	MLP	Linear	MLP
w/o $\mathcal{L}_{gm}$	73.27	76.62	78.14	79.60	81.43	84.19	82.89	85.71	86.86	89.12
w/ $\mathcal{L}_{gm}$	74.02	75.90	79.49	79.94	82.70	84.42	84.07	86.27	88.03	89.31

Table 3: Average accuracy of 5 vision backbone models on 14 commonly-used vision datasets.



Figure 5: Overview of the proposed Local Gradient Mathcing scheme.

fore, inspired by Sharpness-Aware Minimization (SAM) (Foret et al., 2021; Zhang et al., 2023; Zhao et al., 2022; Wen et al., 2023) and noise-robust fine-tuning methods (Hua et al., 2023, 2021; Jiang et al., 2020), we propose explicitly enhancing the denoising capabilities of the downstream task head by aligning local gradients.

Specifically, let C denote the number of classes in the downstream task, and let  $g_{\theta} : \mathbb{R}^d \to \mathbb{R}^C$ represent the linear or MLP classification head parameterized by  $\theta$ . Let  $t^*$  be the feature extracted by  $h^*$ , t be the feature extracted by h, and y be the corresponding label. Let  $\ell(\hat{y}, y)$  be the loss function(typically cross-entropy), and  $\mathcal{L}_{ce}(\mathcal{D}, g_{\theta}) =$  $\mathbb{E}_{(t,y)\sim\mathcal{D}}\ell(g_{\theta}(t),y)$  be the population-level loss where  $\mathcal{D}$  represents the joint distribution of downstream features and labels. Due to the additional randomness introduced by h as a result of noise, t can be viewed as  $t^*$  perturbed by minor disturbances. If both  $t^*$  and t were known, their distribution could be aligned to achieve denoising. However, in practical applications, it is challenging to obtain  $t^*$ . To construct contrastive sample pairs without  $t^*$ , we add Gaussian noise to t to obtain  $\hat{t}$ :

$$\hat{t} = t + \gamma \cdot \delta \tag{3}$$

where  $\delta \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$  denotes the standard normal distribution noise. Our objective is to minimize the discrepancy between the distributions of  $g_{\theta}(t)$  and  $g_{\theta}(\hat{t})$ . Instead of the conventional regularization term  $||g_{\theta}(t) - g_{\theta}(\hat{t})||_2$ , we propose to align the gradient difference:

$$\mathcal{L}_{gm}(\theta) = ||\mathbb{E}_{(t,y)\sim\mathcal{D}}\nabla_{\theta}\ell(g_{\theta}(t), y) - \mathbb{E}_{(\hat{t},y)\sim\hat{\mathcal{D}}}\nabla_{\theta}\ell(g_{\theta}(\hat{t}), y)||_{2}$$
(4)

Intuitively, if the gradients with respect to t and  $\hat{t}$  can be perfectly aligned, then the classification

	RTE	MRPC	CoLA	STS-B
$L^2$ -SP*	70.58	87.74	60.54	89.38
$L^2$ -SP + $\mathcal{L}_{gm}$	71.25	87.62	61.79	89.62
SMART*	72.23	87.86	63.16	90.11
SMART + $\mathcal{L}_{gm}$	72.94	88.61	63.28	90.42
LNSR*	73.31	88.50	63.35	90.23
LNSR + $\mathcal{L}_{gm}$	73.95	89.42	63.82	90.47

Table 4: Evaluation of our method combined with SOTA fine-tuning techniques utilizing BERT-Large as the backbone model across 4 datasets. \* cited from Hua et al. (2021)

head is insensitive to small perturbations in the input, suggesting that it possesses some denoising capability. Consequently, it should be able to mitigate the noise in t, bringing it closer to  $t^*$ .

#### 4.2 Theoretical Analysis

To theoretically support the proposed method, we investigate the properties of Equation (4) and find that it can be upper bounded by the smoothness, input flatness, and loss function value at  $\theta$ . Concretely, since we set  $\gamma$  in Equation (3) to be small, the perturbation can be considered to distribute within an open ball  $B(0, \rho)$ . Consequently, we have the following result:

**Proposition 2.** Suppose  $\ell(g_{\theta}(t), y)$  is  $\beta$ -smooth with  $\rho$ -input flatness  $R_{\rho}(\theta)$  (c.f. Appendix B.4), for any  $\theta \in \Theta$ :

$$\mathcal{L}_{qm}(\theta) \le 2\beta + 2\mathcal{L}_{ce}(\mathcal{D}, g_{\theta}) + R_{\rho}(\theta).$$
 (5)

Proposition 2 demonstrates that  $\mathcal{L}_{gm}$  is closely associated with the smoothness of the loss function in both the parameter space and the input space. As a flat minima is widely acknowledged to benefit the generalization of neural networks (Xie et al., 2021b; Baldassi et al., 2021), it explains the effectiveness of  $\mathcal{L}_{gm}$ . The final loss function is:

$$\mathcal{L} = \mathcal{L}_{ce} + \lambda \mathcal{L}_{gm} \tag{6}$$

#### 4.3 Experiments

To evaluate the performance of the proposed loss function, we conduct extensive experiments on common vision and language tasks using various backbone models. Detailed information about the



(a) No regularization.



(b) With  $L^2$  regularization.



(c) With  $\mathcal{L}_{gm}$ .

Figure 6: Visualization of input sensitivity for models trained with (a) no (b)  $L^2$  (c)  $\mathcal{L}_{gm}$  regularization. We randomly select a sample and introduce perturbations on a two-dimensional hyperplane, where different colors represent different labels, and green indicates the correct label.

Mathad	SST-2		
Method	Linear	MLP	
0%	86.71	87.36	
$0\% +   \nabla_{\theta}\ell(g_{\theta}(t), y)  _2$	87.04	87.24	
$0\% + \cos(\nabla_{\theta} \ell(g_{\theta}(t), y), \nabla_{\theta} \ell(g_{\theta}(\hat{t}), y))$	86.89	87.52	
$0\% + \mathcal{L}_{gm}$	87.42	87.86	

Table 5: Ablation Study. To investigate the effects of reducing  $\mathcal{L}_{gm}$ , experiments are conducted to examine the impact of separately reducing the norm versus increasing the cosine similarity.

models, datasets, hyperparameters, and other experimental settings can be found in Appendix E.

506

507

509

510

511

513

514

516

517

518

519

520

521

525

526

527

529

533

We validate the performance of  $\mathcal{L}_{gm}$  on models pre-trained with noisy data using four commonly used classification datasets: SST-2, SSTfine, 20newsgroup, and CR. The training hyperparameters follow those of Saunshi et al. (2021), where  $\gamma = 0.01$  and  $\lambda = 0.15$  apply to all four experiments. In line with the approach described by Chen et al. (2024), we freeze the model parameters and only fine-tune a linear or MLP classifier head. As shown in Table 1, our model achieves competitive results without reaching the number of training iterations of GPT-2, and  $\mathcal{L}_{gm}$  consistently boosts performance.

To further test the generalizability of our method, we employ large language models, including different versions of Llama-3 (Dubey et al., 2024) and Qwen-2.5 (Hui et al., 2024), for experiments on an additional four datasets. Results in Table 2 indicate that our method provides a 3% improvement across multiple datasets.

In addition, we select five commonly used backbone models in the visual domain and conduct experiments on fourteen datasets. The results are shown in Table 3. It can be seen that our method is equally applicable to visual tasks, achieving a performance improvement of more than 1% under the linear probe setting.

γ         Λ         Linear         MLP           0.001         0.001         76.31         78.54           0.05         0.05         76.54         79.51           0.1         0.1         76.43         79.12		,	DTD			
0.001         0.001         76.31         78.54           0.05         0.05         76.54         79.51           0.1         0.1         76.43         79.12	γ	λ	Linear	MLP		
0.05         0.05         76.54         79.51           0.1         0.1         76.43         79.12	0.001	0.001	76.31	78.54		
0.1 0.1 76.43 79.12	0.05	0.05	76.54	79.51		
	0.1	0.1	76.43	79.12		

Table 6: Hyperparameter sensitivity experiments onDTD with ConvNext as the backbone.

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

555

556

557

558

559

560

561

563

Furthermore, we visualize the sensitivity of different regularization terms to input perturbations, as illustrated in Figure 6. Compared with other regularization methods, our loss function can increase the size of the region for correct decisions, thereby enhancing the model's robustness to input perturbations. We also carry out ablation studies and parameter sensitivity analyses, with results presented in Table 5 and Table 6. These experiments all demonstrate the effectiveness and robustness of our method.

#### 5 Conclusion

In this paper, we investigate the random noise present in language model pre-training datasets, which is inevitable in real-world scenarios but receives little attention. We pre-train multiple GPT-2 models under varying noise levels and find that random noise has a minor impact on the pre-training loss. We then provide a theoretical explanation for this phenomenon and discover that our theory can elucidate the success of multilingual models. Interestingly, we observe that slight noise can sometimes enhance a model's generalization ability. Then, building on the noisy model learning setup, we propose a novel local gradient matching loss. Extensive experiments across multiple datasets in both language and vision tasks, as well as with various backbone models, validate the effectiveness of our proposed method. We hope this work inspires more researchers to focus on data-centric AI.

#### 564

566

567

570

572

575

579

583

584

585

588

589

591

592

594

598

599

601

602

608

610

611

612

613

614

#### Limitations

In this section, we discuss the limitations of this paper.

Firstly, due to limitations in computational resources and costs, we pre-train only the GPT-2 124M model on the OpenWebText dataset and do not train models with other architectures on different datasets. Compared to today's large language models, both the scale of OpenWebText and that of GPT-2 are relatively small. Additionally, the types of noise considered are limited to uniform and Gaussian distributions. However, based on Proposition 1, we argue that training GPT-2 on the Synthetic OpenWebText dataset is sufficient to uncover the essence of the issue, as Proposition 1 makes no assumptions about data distribution or model architecture.

Secondly, on the theoretical front, we consider neural networks as black boxes and focus on analyzing the properties of global minima. Due to limited mathematical skills, we do not delve into the dynamical aspects to specifically examine how random noise within datasets influences model gradients, nor do we explore the differences between global and local minima obtained through stochastic gradient descent. However, experimental results indicate that neural networks trained with stochastic gradient descent do not suffer from significant disturbances.

#### References

- Alon Albalak, Yanai Elazar, Sang Michael Xie, Shayne Longpre, Nathan Lambert, Xinyi Wang, Niklas Muennighoff, Bairu Hou, Liangming Pan, Haewon Jeong, Colin Raffel, Shiyu Chang, Tatsunori Hashimoto, and William Yang Wang. 2024. A survey on data selection for language models. *Transactions on Machine Learning Research*. Survey Certification.
- Zeyuan Allen-Zhu and Yuanzhi Li. 2024a. Physics of language models: Part 3.1, knowledge storage and extraction. In *Forty-first International Conference on Machine Learning*.
- Zeyuan Allen-Zhu and Yuanzhi Li. 2024b. Physics of language models: Part 3.3, knowledge capacity scaling laws. *arXiv preprint arXiv:2404.05405*.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609*.
- Carlo Baldassi, Clarissa Lauditi, Enrico M. Malatesta, Gabriele Perugini, and Riccardo Zecchina. 2021. Un-

veiling the structure of wide flat minima in neural networks. <i>Phys. Rev. Lett.</i> , 127:278301.
David Barrett and Benoit Dherin. 2021. Implicit gradi- ent regularization. In International Conference on Learning Representations.
Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wortman Vaughan. 2010. A theory of learning from different domains. <i>Machine learning</i> , 79:151–175.
Tom B Brown. 2020. Language models are few-shot learners. <i>arXiv preprint arXiv:2005.14165</i> .
Yekun Chai, Qingyi Liu, Shuohuan Wang, Yu Sun, Qi- wei Peng, and Hua Wu. 2024. On training data in- fluence of GPT models. In <i>Proceedings of the 2024</i> <i>Conference on Empirical Methods in Natural Lan-</i> <i>guage Processing</i> , pages 3126–3150, Miami, Florida, USA. Association for Computational Linguistics.
Danqi Chen and Christopher Manning. 2014. A fast and accurate dependency parser using neural networks. In <i>Proceedings of the 2014 Conference on Empirical</i> <i>Methods in Natural Language Processing (EMNLP)</i> , pages 740–750, Doha, Qatar. Association for Com- putational Linguistics.
Feng Chen, Daniel Kunin, Atsushi Yamamura, and Surya Ganguli. 2023. Stochastic collapse: How gra- dient noise attracts SGD dynamics towards simpler subnetworks. In <i>Thirty-seventh Conference on Neu-</i> <i>ral Information Processing Systems</i> .
Guoguo Chen, Shuzhou Chai, Guanbo Wang, Jiayu Du, Wei Qiang Zhang, Chao Weng, Dan Su, Daniel Povey, Jan Trmal, Junbo Zhang, et al. 2021. Gigaspeech: An evolving, multi-domain asr corpus with 10,000 hours of transcribed audio. In 22nd Annual Confer- ence of the International Speech Communication As- sociation, INTERSPEECH 2021, pages 4376–4380. International Speech Communication.
Hao Chen, Jindong Wang, Ankit Shah, Ran Tao, Hongxin Wei, Xing Xie, Masashi Sugiyama, and Bhiksha Raj. 2024. Understanding and mitigating the label noise in pre-training on downstream tasks. In <i>The Twelfth International Conference on Learning</i> <i>Representations</i> .
Sanyuan Chen, Chengyi Wang, Zhengyang Chen, Yu Wu, Shujie Liu, Zhuo Chen, Jinyu Li, Naoyuki Kanda, Takuya Yoshioka, Xiong Xiao, et al. 2022. Wavlm: Large-scale self-supervised pre-training for

615 616

617

618

619

620

621

622

623

624

625

626

627

628

629

630 631

632

633 634

635

636

637

638

639

640

641

642

643

644

645

646 647 648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

Gong Cheng, Junwei Han, and Xiaoqiang Lu. 2017. Remote sensing image scene classification: Benchmark and state of the art. *Proceedings of the IEEE*, 105(10):1865–1883.

full stack speech processing. IEEE Journal of Se-

lected Topics in Signal Processing, 16(6):1505–1518.

Valeriia Cherepanova and James Zou. 2024. Talking nonsense: Probing large language models' understanding of adversarial gibberish inputs. *arXiv preprint arXiv:2404.17120*.

- 671 672
- 673 674
- 675 676
- 67
- 67
- 68
- 681
- 68 68
- 6

6

6

690 691

6

69

0:

69 69

700 701

702

70 70

708

710 711

712

713 714

715 716

717 718 719

720 721

72 72

723 724

- Ethan A. Chi, John Hewitt, and Christopher D. Manning. 2020. Finding universal grammatical relations<br/>in multilingual BERT. In Proceedings of the 58th<br/>Annual Meeting of the Association for Computational<br/>Linguistics, pages 5564–5577, Online. Association<br/>for Computational Linguistics.Li F<br/>La<br/>ex
- Mircea Cimpoi, Subhransu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. 2014. Describing textures in the wild. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3606–3613.
- Sayantan Dasgupta, Trevor Cohn, and Timothy Baldwin. 2023. Cost-effective distillation of large language models. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 7346–7354, Toronto, Canada. Association for Computational Linguistics.
  - Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. Imagenet: A large-scale hierarchical image database. In 2009 IEEE conference on computer vision and pattern recognition, pages 248–255. Ieee.
- Elvis Dohmatob, Yunzhen Feng, and Julia Kempe. 2024. Strong model collapse. *arXiv preprint arXiv:2410.04840*.
- Bill Dolan and Chris Brockett. 2005. Automatically constructing a corpus of sentential paraphrases. In *Third international workshop on paraphrasing* (*IWP2005*).
- Alexey Dosovitskiy. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*.
- Zhekai Du and Jingjing Li. 2023. Diffusion-based probabilistic uncertainty estimation for active domain adaptation. In *Advances in Neural Information Processing Systems*, volume 36, pages 17129–17155. Curran Associates, Inc.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Yanai Elazar, Akshita Bhagia, Ian Helgi Magnusson, Abhilasha Ravichander, Dustin Schwenk, Alane Suhr, Evan Pete Walsh, Dirk Groeneveld, Luca Soldaini, Sameer Singh, Hannaneh Hajishirzi, Noah A. Smith, and Jesse Dodge. 2024. What's in my big data? In *The Twelfth International Conference on Learning Representations.*
- Zhen Fang, Yixuan Li, Jie Lu, Jiahua Dong, Bo Han, and Feng Liu. 2022. Is out-of-distribution detection learnable? In *Advances in Neural Information Processing Systems*, volume 35, pages 37199–37213. Curran Associates, Inc.

Li Fei-Fei, Rob Fergus, and Pietro Perona. 2004a. Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories. In 2004 conference on computer vision and pattern recognition workshop, pages 178–178. IEEE. 725

726

727

728

729

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

763

764

765

767

771

772

773

774

775

- Li Fei-Fei, Rob Fergus, and Pietro Perona. 2004b. Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories. In 2004 conference on computer vision and pattern recognition workshop, pages 178–178. IEEE.
- Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. 2021. Sharpness-aware minimization for efficiently improving generalization. In International Conference on Learning Representations.
- Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, et al. 2020. The pile: An 800gb dataset of diverse text for language modeling. *arXiv preprint arXiv:2101.00027*.
- Yang Gao, Zunlei Feng, Xiaoyang Wang, Mingli Song, Xingen Wang, Xinyu Wang, and Chun Chen. 2023. Reinforcement learning based web crawler detection for diversity and dynamics. *Neurocomputing*, 520:115–128.
- Aaron Gokaslan, Vanya Cohen, Ellie Pavlick, and Stefanie Tellex. 2019. Openwebtext corpus.
- Jeff Z. HaoChen, Colin Wei, Jason Lee, and Tengyu Ma. 2021. Shape matters: Understanding the implicit bias of the noise covariance. In *Proceedings of Thirty Fourth Conference on Learning Theory*, volume 134 of *Proceedings of Machine Learning Research*, pages 2315–2357. PMLR.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770– 778.
- Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth. 2019. Eurosat: A novel dataset and deep learning benchmark for land use and land cover classification. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 12(7):2217–2226.
- Geoffrey Hinton. 2015. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*.
- Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*.

Minqing Hu and Bing Liu. 2004. Mining and summarizing customer reviews. In *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 168–177.

778

790

791

792

793

794

797

798

803

804 805

810

811 812

813

815

816

818

819

820

821 822

823

824

832 833

- Hang Hua, Xingjian Li, Dejing Dou, Cheng-Zhong Xu, and Jiebo Luo. 2023. Improving pretrained language model fine-tuning with noise stability regularization. *IEEE Transactions on Neural Networks and Learning Systems*, pages 1–15.
- Hang Hua, Xingjian Li, Dejing Dou, Chengzhong Xu, and Jiebo Luo. 2021. Noise stability regularization for improving BERT fine-tuning. In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 3229–3241, Online. Association for Computational Linguistics.
- Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, et al. 2024. Qwen2. 5-coder technical report. arXiv preprint arXiv:2409.12186.
- Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and Tuo Zhao. 2020.
  SMART: Robust and efficient fine-tuning for pretrained natural language models through principled regularized optimization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2177–2190, Online. Association for Computational Linguistics.
  - Renren Jin, Jiangcun Du, Wuwei Huang, Wei Liu, Jian Luan, Bin Wang, and Deyi Xiong. 2024. A comprehensive evaluation of quantization strategies for large language models. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 12186–12215, Bangkok, Thailand. Association for Computational Linguistics.
  - Wooyoung Kang, Jonghwan Mun, Sungjun Lee, and Byungseok Roh. 2023. Noise-aware learning from web-crawled image-text data for image captioning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2942–2952.
  - Pride Kavumba, Naoya Inoue, Benjamin Heinzerling, Keshav Singh, Paul Reisert, and Kentarou Inui. 2020.
    Balanced copa: Countering superficial cues in causal reasoning. Association for Natural Language Processing, pages 1105–1108.
- Simon Kornblith, Jonathon Shlens, and Quoc V. Le. 2019. Do better imagenet models transfer better? In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR).
- Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 2013. 3d object representations for fine-grained categorization. In *Proceedings of the IEEE international conference on computer vision workshops*, pages 554–561.
- Alex Krizhevsky, Geoffrey Hinton, et al. 2009. Learning multiple layers of features from tiny images.

Eldar Kurtic, Daniel Campos, Tuan Nguyen, Elias Frantar, Mark Kurtz, Benjamin Fineran, Michael Goin, and Dan Alistarh. 2022. The optimal BERT surgeon: Scalable and accurate second-order pruning for large language models. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 4163–4181, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics. 834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

- Jean-François Le Gall. 2022. *Measure theory, probability, and stochastic processes*. Springer.
- Alycia Lee, Brando Miranda, Sudharsan Sundar, and Sanmi Koyejo. 2023. Beyond scale: the diversity coefficient as a data quality metric demonstrates llms are pre-trained on formally diverse data. *arXiv preprint arXiv:2306.13840*.
- Jingjing Li, Zhekai Du, Lei Zhu, Zhengming Ding, Ke Lu, and Heng Tao Shen. 2021. Divergenceagnostic unsupervised domain adaptation by adversarial attacks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(11):8196–8211.
- Jingjing Li, Zhiqi Yu, Zhekai Du, Lei Zhu, and Heng Tao Shen. 2024a. A comprehensive survey on source-free domain adaptation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(8):5743–5762.
- Ming Li, Yong Zhang, Zhitao Li, Jiuhai Chen, Lichang Chen, Ning Cheng, Jianzong Wang, Tianyi Zhou, and Jing Xiao. 2024b. From quantity to quality: Boosting LLM performance with self-guided data selection for instruction tuning. In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pages 7602–7635, Mexico City, Mexico. Association for Computational Linguistics.
- Hong Liu, Sang Michael Xie, Zhiyuan Li, and Tengyu Ma. 2023. Same pre-training loss, better downstream: Implicit bias matters for language models. In Proceedings of the 40th International Conference on Machine Learning, volume 202 of Proceedings of Machine Learning Research, pages 22188–22214. PMLR.
- Ruikang Liu, Haoli Bai, Haokun Lin, Yuening Li, Han Gao, Zhengzhuo Xu, Lu Hou, Jun Yao, and Chun Yuan. 2024. IntactKV: Improving large language model quantization by keeping pivot tokens intact. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 7716–7741, Bangkok, Thailand. Association for Computational Linguistics.
- Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. 2021. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10012–10022.

Shayne Longpre, Gregory Yauney, Emily Reif, Katherine Lee, Adam Roberts, Barret Zoph, Denny Zhou, Jason Wei, Kevin Robinson, David Mimno, and Daphne Ippolito. 2024. A pretrainer's guide to training data: Measuring the effects of data age, domain coverage, quality, & toxicity. In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pages 3245–3276, Mexico City, Mexico. Association for Computational Linguistics.

900

901

902

903

906

907

908

909

910

911

912

913

914

915

916

917

918

919

921

922

923

924

925

926

927

928

930

931

932

934

935

936

937

940

941

943

- I Loshchilov. 2017. Decoupled weight decay regularization. arXiv preprint arXiv:1711.05101.
- Michal Lukasik, Srinadh Bhojanapalli, Aditya Menon, and Sanjiv Kumar. 2020. Does label smoothing mitigate label noise? In *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 6448–6458. PMLR.
- Ailong Ma, Chenyu Zheng, Junjue Wang, and Yanfei Zhong. 2023. Domain adaptive land-cover classification via local consistency and global diversity. *IEEE Transactions on Geoscience and Remote Sensing*, 61:1–17.
- Subhransu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. 2013. Fine-grained visual classification of aircraft. *arXiv preprint arXiv:1306.5151*.
- Lichao Meng, Hongzu Su, Chunwei Lou, and Jingjing Li. 2022. Cross-domain mutual information adversarial maximization. *Engineering Applications of Artificial Intelligence*, 110:104665.
- Neel Nanda, Lawrence Chan, Tom Lieberum, Jess Smith, and Jacob Steinhardt. 2023. Progress measures for grokking via mechanistic interpretability. In *The Eleventh International Conference on Learning Representations*.
- Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Baolin Wu, Andrew Y Ng, et al. 2011. Reading digits in natural images with unsupervised feature learning. In NIPS workshop on deep learning and unsupervised feature learning, volume 2011, page 4. Granada.
- Maria-Elena Nilsback and Andrew Zisserman. 2008. Automated flower classification over a large number of classes. In 2008 Sixth Indian conference on computer vision, graphics & image processing, pages 722–729. IEEE.
  - OpenAI. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and CV Jawahar. 2012. Cats and dogs. In 2012 *IEEE conference on computer vision and pattern recognition*, pages 3498–3505. IEEE.

Mohammad Taher Pilehvar and Jose Camacho-Collados. 2019. WiC: the word-in-context dataset for evaluating context-sensitive meaning representations. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers),* pages 1267–1273, Minneapolis, Minnesota. Association for Computational Linguistics. 944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

- Telmo Pires, Eva Schlinger, and Dan Garrette. 2019. How multilingual is multilingual BERT? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4996–5001, Florence, Italy. Association for Computational Linguistics.
- Stijn Pletinckx, Kevin Borgolte, and Tobias Fiebig. 2021. Out of sight, out of mind: Detecting orphaned web pages at internet-scale. In *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security*, pages 21–35.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning transferable visual models from natural language supervision. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 8748–8763. PMLR.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67.
- Maithra Raghu, Chiyuan Zhang, Jon Kleinberg, and Samy Bengio. 2019. Transfusion: Understanding transfer learning for medical imaging. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.
- Tal Ridnik, Emanuel Ben-Baruch, Asaf Noy, and Lihi Zelnik-Manor. 2021. Imagenet-21k pretraining for the masses. In Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1).
- Jinghan Ru, Jun Tian, Chengwei Xiao, Jingjing Li, and Heng Tao Shen. 2024. Imbalanced open set domain adaptation via moving-threshold estimation and gradual alignment. *IEEE Transactions on Multimedia*, 26:2504–2514.
- Antony Samuels and John Mcgonical. 2020. News sentiment analysis. *arXiv preprint arXiv:2007.02238*.

- 1000 1001
- 1002

1004

1009

1011

1012

1014

1016

1018

1019

1020

1021 1022

1023

1024

1025

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037 1038

1039

1040

1041

1042

1043

1044

1045 1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

Nikunj Saunshi, Sadhika Malladi, and Sanjeev Arora. 2021. A mathematical exploration of why language models help solve downstream tasks. In *International Conference on Learning Representations*.

 Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. 2022. Laion-5b: An open large-scale dataset for training next generation image-text models. In Advances in Neural Information Processing Systems, volume 35, pages 25278–25294. Curran Associates, Inc.

- Mohamed El Amine Seddik, Suei-Wen Chen, Soufiane Hayou, Pierre Youssef, and Merouane Abdelkader DEBBAH. 2024. How bad is training on synthetic data? a statistical analysis of language model collapse. In *First Conference on Language Modeling*.
  - Oleh Shliazhko, Alena Fenogenova, Maria Tikhonova, Anastasia Kozlova, Vladislav Mikhailov, and Tatiana Shavrina. 2024. mGPT: Few-shot learners go multilingual. *Transactions of the Association for Computational Linguistics*, 12:58–79.
  - Ilia Shumailov, Zakhar Shumaylov, Yiren Zhao, Nicolas Papernot, Ross Anderson, and Yarin Gal. 2024. Ai models collapse when trained on recursively generated data. *Nature*, 631(8022):755–759.
  - Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. 2017. Smoothgrad: removing noise by adding noise. *arXiv preprint arXiv:1706.03825*.
  - Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing, pages 1631–1642, Seattle, Washington, USA. Association for Computational Linguistics.
- Luca Soldaini, Rodney Kinney, Akshita Bhagia, Dustin Schwenk, David Atkinson, Russell Authur, Ben Bogin, Khyathi Chandu, Jennifer Dumas, Yanai Elazar, Valentin Hofmann, Ananya Jha, Sachin Kumar, Li Lucy, Xinxi Lyu, Nathan Lambert, Ian Magnusson, Jacob Morrison, Niklas Muennighoff, Aakanksha Naik, Crystal Nam, Matthew Peters, Abhilasha Ravichander, Kyle Richardson, Zejiang Shen, Emma Strubell, Nishant Subramani, Oyvind Tafjord, Evan Walsh, Luke Zettlemoyer, Noah Smith, Hannaneh Hajishirzi, Iz Beltagy, Dirk Groeneveld, Jesse Dodge, and Kyle Lo. 2024. Dolma: an open corpus of three trillion tokens for language model pretraining research. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 15725-15788, Bangkok, Thailand. Association for Computational Linguistics.

Hwanjun Song, Minseok Kim, Dongmin Park, Yooju Shin, and Jae-Gil Lee. 2022. Learning from noisy labels with deep neural networks: A survey. *IEEE transactions on neural networks and learning systems*, 34(11):8135–8153.

1057

1058

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1072

1073

1075

1076

1077

1078

1080

1081

1082

1083

1084

1085

1087

1088

1089

1090

1091

1092

1094

1095

1096

1098

1099

1100

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

- Nathan Srebro, Karthik Sridharan, and Ambuj Tewari. 2010. Smoothness, low noise and fast rates. In *Advances in Neural Information Processing Systems*, volume 23. Curran Associates, Inc.
- Chen Sun, Abhinav Shrivastava, Saurabh Singh, and Abhinav Gupta. 2017. Revisiting unreasonable effectiveness of data in deep learning era. In *Proceedings of the IEEE international conference on computer vision*, pages 843–852.
- I Sutskever. 2014. Sequence to sequence learning with neural networks. *arXiv preprint arXiv:1409.3215*.
- Mingxing Tan and Quoc Le. 2019. EfficientNet: Rethinking model scaling for convolutional neural networks. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 6105– 6114. PMLR.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Bastiaan S Veeling, Jasper Linmans, Jim Winkens, Taco Cohen, and Max Welling. 2018. Rotation equivariant cnns for digital pathology. In *Medical Image Computing and Computer Assisted Intervention–MICCAI* 2018: 21st International Conference, Granada, Spain, September 16-20, 2018, Proceedings, Part II 11, pages 210–218. Springer.
- Ziheng Wang, Jeremy Wohlwend, and Tao Lei. 2020. Structured pruning of large language models. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 6151–6162, Online. Association for Computational Linguistics.
- Colin Wei, Sang Michael Xie, and Tengyu Ma. 2021. Why do pretrained language models help in downstream tasks? an analysis of head and prompt tuning. In *Advances in Neural Information Processing Systems*, volume 34, pages 16158–16170. Curran Associates, Inc.
- Kaiyue Wen, Tengyu Ma, and Zhiyuan Li. 2023. How sharpness-aware minimization minimizes sharpness? In *The Eleventh International Conference on Learning Representations*.
- Sanghyun Woo, Shoubhik Debnath, Ronghang Hu, Xinlei Chen, Zhuang Liu, In So Kweon, and Saining Xie. 2023. Convnext v2: Co-designing and scaling convnets with masked autoencoders. In *Proceedings* of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 16133–16142.

- 1114 1115 1116 1117 1118 1119 1120 1121 1122 1123 1124 1125 1126 1127 1128 1129 1130 1131 1132 1133
- 1134 1135 1136 1137 1138 1139 1140 1141 1142 1143 1144 1145 1146
- 1147 1148 1149 1150 1151 1152 1153 1154 1155 1156 1157 1158 1159 1160 1161 1162 1163 1164 1165 1166
- 1167 1168 1169
- 1170 1171

- Kangkai Wu, Jingjing Li, Lichao Meng, Fengling Li, and Ke Lu. 2024. Online adaptive fault diagnosis with test-time domain adaptation. IEEE Transactions on Industrial Informatics, pages 1-11.
- Mingyu Xiao, Jianan Zhang, and Weibo Lin. 2022. Parameterized algorithms and complexity for the traveling purchaser problem and its variants. Journal of Combinatorial Optimization, pages 1-17.
  - Sang Michael Xie, Tengyu Ma, and Percy Liang. 2021a. Composed fine-tuning: Freezing pre-trained denoising autoencoders for improved generalization. In Proceedings of the 38th International Conference on Machine Learning, volume 139 of Proceedings of Machine Learning Research, pages 11424–11435. PMLR.
  - Sang Michael Xie, Hieu Pham, Xuanyi Dong, Nan Du, Hanxiao Liu, Yifeng Lu, Percy S Liang, Quoc V Le, Tengyu Ma, and Adams Wei Yu. 2023a. Doremi: Optimizing data mixtures speeds up language model pretraining. In Advances in Neural Information Processing Systems, volume 36, pages 69798-69818. Curran Associates, Inc.
- Sang Michael Xie, Shibani Santurkar, Tengyu Ma, and Percy S Liang. 2023b. Data selection for language models via importance resampling. In Advances in Neural Information Processing Systems, volume 36, pages 34201-34227. Curran Associates, Inc.
- Zeke Xie, Issei Sato, and Masashi Sugiyama. 2021b. A diffusion theory for deep learning dynamics: Stochastic gradient descent exponentially favors flat minima. In International Conference on Learning Representations.
- Zeke Xie, Xinrui Wang, Huishuai Zhang, Issei Sato, and Masashi Sugiyama. 2022. Adaptive inertia: Disentangling the effects of adaptive learning rate and momentum. In Proceedings of the 39th International Conference on Machine Learning, volume 162 of Proceedings of Machine Learning Research, pages 24430-24459. PMLR.
- Zeke Xie, Zhiqiang Xu, Jingzhao Zhang, Issei Sato, and Masashi Sugiyama. 2023c. On the overlooked pitfalls of weight decay and how to mitigate them: A gradient-norm perspective. In Advances in Neural Information Processing Systems, volume 36, pages 1208-1228. Curran Associates, Inc.
- Zeke Xie, Li Yuan, Zhanxing Zhu, and Masashi Sugiyama. 2021c. Positive-negative momentum: Manipulating stochastic gradient noise to improve generalization. In Proceedings of the 38th International Conference on Machine Learning, volume 139 of Proceedings of Machine Learning Research, pages 11448-11458. PMLR.
- Bang Yang, Fenglin Liu, Yuexian Zou, Xian Wu, Yaowei Wang, and David A. Clifton. 2024. Zeronlg: Aligning and autoencoding domains for zero-shot multimodal and multilingual natural language generation. IEEE Transactions on Pattern Analysis and Machine Intelligence, 46(8):5712-5724.

Shuo Yang, Zeke Xie, Hanyu Peng, Min Xu, Mingming Sun, and Ping Li. 2023. Dataset pruning: Reducing training data by examining generalization influence. In The Eleventh International Conference on Learning Representations.

1172

1173

1174

1175

1176

1177

1178

1179

1180

1181

1182

1183

1184

1185

1186

1187

1188

1189

1190

1191

1192

1193

1194

1195

1196

1197

1198

1199

1200

1201

1202

1203

1204

1205

1206

1207

1208

1209

1210

1211

1212

1213

1214

1215

1216

1217

1218

1219

1220

1221

1222

1223

1224

1225

1226

- Yuqi Ye and Wei Gao. 2024. Llm-pcgc: Large language model-based point cloud geometry compression. arXiv preprint arXiv:2408.08682.
- YongKang Yin, Xu Li, Ying Shan, and YueXian Zou. 2024. Afl-net: Integrating audio, facial, and lip modalities with a two-step cross-attention for robust speaker diarization in the wild. In Proc. Interspeech 2024, pages 42-46.
- Shiran Zada, Itay Benou, and Michal Irani. 2022. Pure noise to the rescue of insufficient data: Improving imbalanced classification by training on random noise images. In Proceedings of the 39th International Conference on Machine Learning, volume 162 of Proceedings of Machine Learning Research, pages 25817-25833. PMLR.
- Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. 2021a. Understanding deep learning (still) requires rethinking generalization. Communications of the ACM, 64(3):107–115.
- Renrui Zhang, Rongyao Fang, Wei Zhang, Peng Gao, Kunchang Li, Jifeng Dai, Yu Qiao, and Hongsheng Li. 2021b. Tip-adapter: Training-free clip-adapter for better vision-language modeling. arXiv preprint arXiv:2111.03930.
- Xingxuan Zhang, Renzhe Xu, Han Yu, Hao Zou, and Peng Cui. 2023. Gradient norm aware minimization seeks first-order flatness and improves generalization. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 20247-20257.
- Xuchao Zhang, Fanglan Chen, Chang-Tien Lu, and Naren Ramakrishnan. 2019. Mitigating uncertainty in document classification. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 3126-3136, Minneapolis, Minnesota. Association for Computational Linguistics.
- Yingtao Zhang, Haoli Bai, Haokun Lin, Jialin Zhao, Lu Hou, and Carlo Vittorio Cannistraci. 2024. Plugand-play: An efficient post-training pruning method for large language models. In The Twelfth International Conference on Learning Representations.
- Jiaqi Zhao, Miao Zhang, Chao Zeng, Ming Wang, Xuebo Liu, and Liqiang Nie. 2024. LRQuant: Learnable and robust post-training quantization for large language models. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 2240-2255, Bangkok, Thailand. Association for Computational Linguistics.

Yang Zhao, Hao Zhang, and Xiuyuan Hu. 2022. Penalizing gradient norm for efficiently improving generalization in deep learning. In *Proceedings of the* 39th International Conference on Machine Learning, volume 162 of *Proceedings of Machine Learning Research*, pages 26982–26992. PMLR.

1228

1229

1230 1231

1232

1233

1234

1235

1236

1238

1239

1240 1241

1242 1243

1244

1245 1246

1247

1248

1249

1250 1251

1252

1253

1254

1255

1256 1257

1258 1259

1260

1261

- Chenyu Zheng, Wei Huang, Rongzhen Wang, Guoqiang Wu, Jun Zhu, and Chongxuan Li. 2024. On mesa-optimization in autoregressively trained transformers: Emergence and capability. *arXiv preprint arXiv:2405.16845*.
- Chenyu Zheng, Guoqiang Wu, Fan Bao, Yue Cao, Chongxuan Li, and Jun Zhu. 2023a. Revisiting discriminative vs. generative classifiers: Theory and implications. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 42420–42477. PMLR.
  - Chenyu Zheng, Guoqiang Wu, and Chongxuan LI. 2023b. Toward understanding generative data augmentation. In *Advances in Neural Information Processing Systems*, volume 36, pages 54046–54060. Curran Associates, Inc.
- Jing Zhou, Chenglin Jiang, Wei Shen, Xiao Zhou, and Xiaonan He. 2024. Leveraging web-crawled data for high-quality fine-tuning. *arXiv preprint arXiv:2408.08003*.
- Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. 2022a. Conditional prompt learning for vision-language models. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 16816–16825.
- Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. 2022b. Learning to prompt for visionlanguage models. *International Journal of Computer Vision*, 130(9):2337–2348.

#### Α Notations

The commonly used notations and their descriptions are as follows.

Notation	Description
L	context length
d	embedding dimension
$\mathcal{W}$	vocabulary of words
$V =  \mathcal{W} $	vocabulary size
$\mathcal{X} = \cup_{i=1}^L \mathcal{W}^i$	model input space
$\mathcal{H}$	model space
$h: \mathcal{X} \to \mathbb{R}^V \in \mathcal{H}$	language model
$\Delta_A$	distribution defined on a discrete set A
$P^c \in \Delta_{\mathcal{X} \times \mathcal{W}}$	distribution of clean data
$P^n \in \Delta_{\mathcal{X} \times \mathcal{W}}$	distribution of pure noise data
$P^m \in \Delta_{\mathcal{X} \times \mathcal{W}}$	distribution of mixed noisy data
α	proportion of noise in training data
$P_X$	marginal distribution of the joint distribution P
$P_{\cdot X}$	conditional distribution of the joint distribution $P$
$p^h_{\cdot x}(w)$	the $w$ -th dimension of the probability distribution corresponding to $h(x)$
$\operatorname{supp}(P^c)$	support of distribution $P^c$
$\mathcal{L}_{ntp}(P,h)$	next-token prediction loss of model $h$ on the distribution $P$
$g_{\theta}: \mathbb{R}^d \to \mathbb{R}^C$	downstream classification head
$ heta\in\Theta$	parameters of $g$
$t\in\mathcal{T}$	feature of downstream task data extracted by backbone model
$y\in \mathcal{Y}$	label of downstream task data
$C =  \mathcal{Y} $	number of classes of the downstream task
$\ell(\hat{y},y)$	downstream task loss function, typically cross-entropy
$\mathcal{D}$	joint distribution of downstream feature and label
$\mathcal{L}_{ce}(\mathcal{D}, g_{\theta})$	population-level loss with downstream data distribution $\mathcal{D}$ and head $g_{\theta}$

Table 7: Nomenclature.

1268

1270

1271 1272

1273

1274

1275

1276

1277

1264

#### B Proofs

#### **B.1** Explanation of Equation (1)

Let  $\mathcal{M}$  be a measurable space, and let  $P_1$  and  $P_2$  be probability measures defined on this space. We assume that  $N_1$  samples are drawn from  $P_1$  and  $N_2$  samples from  $P_2$ . Define  $\mu = \frac{N_1}{N_1 + N_2}$ , so that  $1 - \mu = \frac{N_2}{N_1 + N_2}$ . We aim to show that this collection of  $N_1 + N_2$  samples can be regarded as drawn from a mixed

distribution

$$P_3 = \mu P_1 + (1 - \mu) P_2$$

First, define a new probability measure  $P_3$  as  $P_3(A) = \alpha P_1(A) + (1 - \alpha)P_2(A)$  for any measurable set  $A \subseteq \mathcal{M}$ . Here,  $P_3$  is a convex combination of  $P_1$  and  $P_2$ , and thus  $P_3$  is also a valid probability measure (Le Gall, 2022).

For any measurable set  $A \subseteq \mathcal{M}$ , we examine the probability that a single sample point falls in A by law of total probability:

- A sample from  $P_1$  is selected with probability  $\mu$ , and within this case, the probability of landing in A is  $P_1(A)$ .
- A sample from  $P_2$  is selected with probability  $1 \mu$ , and the probability of it falling in A is  $P_2(A)$ .

Thus, the probability of any given sample point falling in A is

$$\mu P_1(A) + (1 - \mu)P_2(A) = P_3(A)$$

Since  $N_1$  samples are drawn from  $P_1$  and  $N_2$  samples from  $P_2$ , these samples collectively follow the distribution  $P_3$  as each individual sample's probability of being in any measurable set A is consistent with  $P_3(A)$ . Therefore, drawing  $N_1 + N_2$  samples in this manner is equivalent to drawing  $N_1 + N_2$  samples from  $P_3$ .

#### **B.2** Proof of Proposition 1

Before procedding to the proof, we first establish a useful lemma.

**Lemma 1.** *If Assumption 1 holds, then for any*  $h \in H$ *, we have* 

$$\mathcal{L}_{ntp}(P^m, h) = \alpha \mathcal{L}_{ntp}(P^n, h) + (1 - \alpha) \mathcal{L}_{ntp}(P^c, h)$$

*Proof.* Let  $x_i$ ,  $i = 1, 2, ..., |\mathcal{X}|$  denote all prefixes, and  $w_j$ , j = 1, 2, ..., V denote all tokens. For all  $x \in \mathcal{X}$ , by Equation (1), we have: 1285

$$P_X^m(x) = \sum_{j=1}^V P^m(x, w_j) = \sum_{j=1}^V \alpha P^n(x, w_j) + (1 - \alpha) P^c(x, w_j)$$
1280

$$= \alpha \sum_{j=1}^{V} P^{n}(x, w_{j}) + (1 - \alpha) \sum_{j=1}^{V} P^{c}(x, w_{j}) = \alpha P_{X}^{n}(x) + (1 - \alpha) P_{X}^{c}(x)$$
(7) 1287

This indicates that the marginal distribution possesses additivity. Consequently,

$$\mathcal{L}_{ntp}(P^m, h) = \mathbb{E}_{x \sim P_X^m} \mathbb{E}_{w \sim P_{\cdot|x}^m} - \log(\mathbf{p}^h_{\cdot|x}(w)) = \sum_{i=1}^{|\mathcal{X}|} P_X^m(x_i) \cdot \mathbb{E}_{w \sim P_{\cdot|x_i}^m} - \log(\mathbf{p}^h_{\cdot|x_i}(w))$$

$$1289$$

$$= \sum_{i=1}^{|\mathcal{X}|} [(1-\alpha)P_X^c(x_i) + \alpha P_X^n(x_i)] \cdot \mathbb{E}_{w \sim P_{\cdot|x_i}^m} - \log(\mathbf{p}_{\cdot|x_i}^h(w))$$
 (Equation (7)) 1290

$$= (1 - \alpha) \sum_{i=1}^{|\mathcal{X}|} P_X^c(x_i) \mathbb{E}_{w \sim P_{\cdot|x_i}^m} - \log(\mathbf{p}_{\cdot|x_i}^h(w)) + \alpha \sum_{i=1}^{|\mathcal{X}|} P_X^n(x_i) \mathbb{E}_{w \sim P_{\cdot|x_i}^m} - \log(\mathbf{p}_{\cdot|x_i}^h(w))$$
(8)

The conditional distributions do not generally exhibit a linear relationship:

$$P^{m}_{\cdot|x}(w|x) = \frac{P^{m}(x,w)}{P^{m}_{X}(x)} = \frac{(1-\alpha)P^{c}(x,w) + \alpha P^{n}(x,w)}{(1-\alpha)P^{c}_{X}(x) + \alpha P^{n}_{X}(x)} \neq P^{c}_{\cdot|x}(w|x) \neq P^{n}_{\cdot|x}(w|x)$$

However, if  $\operatorname{supp}(P^c) \cap \operatorname{supp}(P^n) = \emptyset$ , it immediately follows that:

$$P^m_{\cdot|x}(w|x) = \frac{(1-\alpha)P^c(x,w) + \alpha P^n(x,w)}{(1-\alpha)P^c_X(x) + \alpha P^n_X(x)} = \begin{cases} P^c_{\cdot|x}(w|x) & \text{if } (x,w) \in \text{supp}(P^c), \\ P^n_{\cdot|x}(w|x) & \text{if } (x,w) \in \text{supp}(P^n). \end{cases}$$

1292

1294

1282

1283

1288

Consequently,

$$\sum_{i=1}^{|\mathcal{X}|} P_X^c(x_i) \mathbb{E}_{w \sim P_{\cdot|x_i}^m} - \log(\mathbf{p}_{\cdot|x_i}^h(w)) = \sum_{i=1}^{|\mathcal{X}|} P_X^c(x_i) \mathbb{E}_{w \sim P_{\cdot|x_i}^c} - \log(\mathbf{p}_{\cdot|x_i}^h(w)) = \mathcal{L}_{ntp}(P^c, h) \quad (9)$$
1293

Similarly,

$$\sum_{i=1}^{|\mathcal{X}|} P_X^n(x_i) \mathbb{E}_{w \sim P_{\cdot|x_i}^m} - \log(\boldsymbol{p}_{\cdot|x_i}^h(w)) = \mathcal{L}_{ntp}(P^n, h)$$
(10) 1295

By substituting Equation (9) and Equation (10) into Equation (8), the proof is completed.  $\Box$  1296

Now we can prove Proposition 1.

**Proposition 1.** Under Assumption 1, let  $h^*$  be a model trained on  $P^c$ , with  $\mathcal{L}_{ntp}(P^c, h^*) = -\log p_c$ and  $\mathcal{L}_{ntp}(P^n, h^*) = -\log p_n$ . When the model h is trained on a mixed distribution  $P^m$  which includes noise, it attempts to fit  $P^n$ , leading to an increase in the loss on the clean distribution  $P^c$ , such that  $\mathcal{L}_{ntp}(P^c, h) = -\log(p_c - \epsilon)$  and  $\mathcal{L}_{ntp}(P^n, h) = -\log(p_n + \epsilon/k)$  for some  $\epsilon > 0$  (k can be shown to be  $\Omega(e^{\mathcal{L}_{ntp}(P^n, h)})$ ). Let  $\eta = \alpha p_c - (1 - \alpha)kp_n$ . We arrive at the following conclusions:

(1) If  $\alpha \leq \frac{kp_n}{p_c+kp_n}$ , then for any  $0 < \epsilon < p_c$ , we have  $\mathcal{L}_{ntp}(P^m, h) \geq \mathcal{L}_{ntp}(P^m, h^*)$ . This means that when  $\alpha$  is sufficiently small, the global minimum on  $P^m$  will not be affected by noise.

(2) If  $\alpha > \frac{kp_n}{p_c+kp_n}$ , then for  $\epsilon < \eta$ , it holds that  $\mathcal{L}_{ntp}(P^m, h) < \mathcal{L}_{ntp}(P^m, h^*)$ . This suggests that if  $\alpha$  is large enough, the impact on the optimal hypothesis is at least as much as  $\alpha p_c - (1-\alpha)kp_n$ .

1307 (3) When  $\alpha < \frac{1}{3}$  and  $k > \frac{\alpha(1-3\alpha)p_c}{(1-\alpha)(2-3\alpha)p_n}$ , for  $\epsilon \ge 3\eta$  we get  $\mathcal{L}_{ntp}(P^m, h^*) < \mathcal{L}_{ntp}(P^m, h)$ . Similarly, 1308 it can be shown that  $\epsilon$  does not exceed  $2\eta$  when  $\alpha > \max(\frac{kp_n}{p_c+kp_n}, \frac{1}{2})$  and  $k > \frac{(2\alpha-1)p_c}{2(1-\alpha)p_n}$ . This indicates 1309 that when k is sufficiently large, the effect of noise is at most  $\mathcal{O}(\alpha p_c - (1-\alpha)kp_n)$ .

1310 *Proof.* We first establish that k is  $\Omega(e^{\mathcal{L}_{ntp}(P^n,h)})$ , thereby ensuring that  $\eta \ll \alpha p_c$ . Note that

$$\epsilon = \frac{1}{e^{\mathcal{L}_{ntp}(P^c,h)}} - \frac{1}{e^{\mathcal{L}_{ntp}(P^c,h)}} = \frac{e^{\mathcal{L}_{ntp}(P^c,h) - \mathcal{L}_{ntp}(P^c,h)} - 1}{e^{\mathcal{L}_{ntp}(P^c,h)}}$$
(11)

$$\frac{\epsilon}{k} = \frac{1}{e^{\mathcal{L}_{ntp}(P^n,h)}} - \frac{1}{e^{\mathcal{L}_{ntp}(P^n,h)}} = \frac{e^{\mathcal{L}_{ntp}(P^n,h) - \mathcal{L}_{ntp}(P^n,h)} - 1}{e^{\mathcal{L}_{ntp}(P^n,h)}}$$
(12)

1313 Therefore

1297

1303 1304

1305

1311

1312

1314 
$$k = \frac{\epsilon}{\frac{\epsilon}{k}} = e^{\mathcal{L}_{ntp}(P^n,h) - \mathcal{L}_{ntp}(P^c,h)} \cdot \frac{e^{\mathcal{L}_{ntp}(P^c,h) - \mathcal{L}_{ntp}(P^c,h) - 1}}{e^{\mathcal{L}_{ntp}(P^n,h) - \mathcal{L}_{ntp}(P^n,h) - 1}}$$

1315 
$$> e^{\mathcal{L}_{ntp}(P^n,h) - \mathcal{L}_{ntp}(P^c,h)} \cdot \frac{\mathcal{L}_{ntp}(P^r,h) - \mathcal{L}_{ntp}(P^r,h)}{e^{\mathcal{L}_{ntp}(P^n,h) - \mathcal{L}_{ntp}(P^n,h)}}$$

1316 
$$= e^{\mathcal{L}_{ntp}(P^{n},h)} \cdot \frac{\mathcal{L}_{ntp}(P^{c},h) - \mathcal{L}_{ntp}(P^{c},h)}{e^{\mathcal{L}_{ntp}(P^{c},h)}}$$
(13)

1317 where  $\frac{\mathcal{L}_{ntp}(P^c,h) - \mathcal{L}_{ntp}(P^c,h)}{e^{\mathcal{L}_{ntp}(P^c,h)}}$  only depends on  $P^c$ , h and h. It is worth noting that when  $P^n$  is random 1318 noise,  $e^{\mathcal{L}_{ntp}(P^n,h) - \mathcal{L}_{ntp}(P^n,h)} - 1$  is close to 0, which leads to k exceeding the lower bound established in 1319 Equation (13). Then:

(1) If 
$$\alpha \le \frac{kp_n}{p_c + kp_n}$$
, we have

1321 
$$\mathcal{L}_{ntp}(P^m, h^*) - \mathcal{L}_{ntp}(P^m, h) = (1 - \alpha)(\mathcal{L}_{ntp}(P^c, h^*) - \mathcal{L}_{ntp}(P^c, h)) + \alpha(\mathcal{L}_{ntp}(P^n, h^*) - \mathcal{L}_{ntp}(P^n, h))$$

$$= (1 - \alpha)\log\frac{p_c - \epsilon}{k} + \alpha\log\frac{p_n + \frac{\epsilon}{k}}{k}$$
(14)

1323 
$$p_c \qquad p_n \\ \leq (1-\alpha) \cdot \frac{-\epsilon}{n_c} + \alpha \cdot \frac{\epsilon}{n_c} \qquad (\log(1+t) \le t)$$

1324 
$$=\epsilon\left[\frac{(\alpha-1)}{p_c} + \frac{\alpha}{kp_n}\right] = \epsilon \frac{\alpha p_c - (1-\alpha)kp_n}{kp_c p_n}$$
(15)

1325 As 
$$\alpha \leq \frac{kp_n}{p_c + kp_n} \iff \alpha p_c - (1 - \alpha)kp_n \leq 0$$
, for  $\epsilon > 0$  we have  $\mathcal{L}_{ntp}(P^m, h^*) \leq \mathcal{L}_{ntp}(P^m, h)$ 

(2)when  $\alpha > \frac{kp_n}{p_c + kp_n}$  and  $\epsilon < \alpha p_c - (1 - \alpha)kp_n$ , we have

$$\mathcal{L}_{ntp}(P^m, h^*) - \mathcal{L}_{ntp}(P^m, h) = (1 - \alpha) \log \frac{p_c - \epsilon}{p_c} + \alpha \log \frac{p_n + \frac{\epsilon}{k}}{p_n}$$
(Equation (14)) 132

$$\geq (1-\alpha)\frac{-\epsilon}{p_c-\epsilon} + \alpha \frac{\frac{\epsilon}{k}}{p_n + \frac{\epsilon}{k}} \qquad (\log t \ge 1 - \frac{1}{t}) \qquad 132t$$

$$=\epsilon\left(\frac{\alpha-1}{p_c-\epsilon}+\frac{\alpha}{kp_n+\epsilon}\right)$$
1329

$$= \frac{\epsilon}{(p_c - \epsilon)(kp_n + \epsilon)} [\alpha(p_c - \epsilon) - (1 - \alpha)(kp_n + \epsilon)]$$
1330

$$= \frac{\epsilon}{(p_c - \epsilon)(kp_n + \epsilon)} [\alpha p_c - (1 - \alpha)kp_n - \epsilon]$$
(16) 1331

As  $\epsilon < \alpha p_c - (1 - \alpha)kp_n < \alpha p_c < p_c$ , by Equation (16) we have  $\mathcal{L}_{ntp}(P^m, h) - \mathcal{L}_{ntp}(P^m, h) > 0.$  (3) Let 1333

$$f(\epsilon) = (1-\alpha)\log\frac{p_c - \epsilon}{p_c} + \alpha\log\frac{p_n + \frac{\epsilon}{k}}{p_n} \stackrel{p'_n = kp_n}{=} (1-\alpha)\log\frac{p_c - \epsilon}{p_c} + \alpha\log\frac{p'_n + \epsilon}{p'_n}$$
(17) 1334

Take the derivative of  $f(\epsilon)$ :

$$f'(\epsilon) = (1-\alpha)\frac{-\frac{1}{p_c}}{1-\frac{\epsilon}{p_c}} + \alpha \frac{\frac{1}{p'_n}}{1+\frac{\epsilon}{p'_n}} = (1-\alpha)\frac{1}{\epsilon - p_c} + \alpha \frac{1}{p'_n + \epsilon} = \frac{[\alpha p_c - (1-\alpha)p'_n] - \epsilon}{(p_c - \epsilon)(p'_n + \epsilon)}$$
(18) 1336

Without loss of generality, assume  $\eta > 0$ , then  $f(\epsilon)$  is monotonically increasing on  $[0, \eta)$  and monotonically decreasing on  $(\eta, p_c)$ . Therefore, to prove that  $\mathcal{L}_{ntp}(P^m, h^*) < \mathcal{L}_{ntp}(P^m, h)$  for  $\epsilon \ge 3\eta$ , we only need to show  $f(3\eta) < 0$  when  $k > \frac{\alpha(1-3\alpha)p_c}{(1-\alpha)(2-3\alpha)p_n}$ . Notice that

$$f(3\eta) = (1-\alpha)\log(1 - \frac{3\alpha p_c - 3(1-\alpha)p'_n}{p_c}) + \alpha\log(1 + \frac{3\alpha p_c - 3(1-\alpha)p'_n}{p'_n})$$
1340

$$= (1 - \alpha)\log(1 - 3\alpha + \frac{3(1 - \alpha)}{\frac{p_c}{p'_n}}) + \alpha\log(3\alpha - 2 + 3\alpha\frac{p_c}{p'_n})$$
(19) 134

Let

$$g_3(t) = (1 - \alpha)\log(1 - 3\alpha + \frac{3(1 - \alpha)}{t}) + \alpha\log(3\alpha - 2 + 3\alpha t)$$
(20) 1342

Take the derivative:

$$g'_{3}(t) = (1-\alpha)\frac{1}{1-3\alpha + \frac{3(1-\alpha)}{t}}\frac{3(\alpha-1)}{t^{2}} + \alpha\frac{3\alpha}{3\alpha - 2 + 3\alpha t}$$
(21) 1349

$$= \frac{-3(1-\alpha)^2}{(1-3\alpha)t^2 + (1-\alpha)t} + \frac{3\alpha^2}{3\alpha - 2 + 3\alpha t}$$
(22) 1346

$$=\frac{-3(1-\alpha)^2(3\alpha-2+3\alpha t)+3\alpha^2[(1-3\alpha)t^2+(1-\alpha)t]}{[(1-3\alpha)t^2+(1-\alpha)t](3\alpha-2+3\alpha t)}$$
(23)

$$=\frac{[\alpha t + (\alpha - 1)][3\alpha(1 - 3\alpha)t + 3(1 - \alpha)(3\alpha - 2)]}{[(1 - 3\alpha)t^2 + (1 - \alpha)t](3\alpha - 2 + 3\alpha t)}$$
(24)

First, consider the denominator. Since  $\alpha < \frac{1}{3}$ , it is clear that  $(1 - 3\alpha)t^2 + (1 - \alpha)t > 0$ . Given that  $t = \frac{p_c}{p'_n} > \frac{1-\alpha}{\alpha}$  (because  $\eta > 0$ ), it follows that  $3\alpha - 2 + 3\alpha t > 1 > 0$ . Therefore, the denominator is always positive. Next, we consider the numerator. Since  $\eta > 0$ , it follows that  $\alpha t + (\alpha - 1) > 0$ . Therefore, when  $t = \frac{p_c}{p'_n} = \frac{p_c}{kp_n} < \frac{(1-\alpha)(2-3\alpha)}{\alpha(1-3\alpha)}$ , we have  $g'_3(t) < 0$ . This means that  $g_3(t)$  is monotonically decreasing on  $(\frac{1-\alpha}{\alpha}, \frac{(1-\alpha)(2-3\alpha)}{\alpha(1-3\alpha)}]$ . Consequently,  $f(3\eta) = g_3(t) \le g_3(\frac{1-\alpha}{\alpha}) = 0$ . 1349

1326

1335



Figure 7: Visualization of k and  $\mathcal{L}_{ntp}(P^m, h^*) - \mathcal{L}_{ntp}(P^m, h)$ . (a) The trend of k as it changes with training, plotted using the model trained on  $P^c$  as  $h^*$ . (b) Visualization of  $\mathcal{L}_{ntp}(P^m, h)$  when the parameter settings are consistent with the experiment.

Following the same line of reasoning, when  $\alpha > \frac{1}{2}$ , we have

=

1355 
$$f(2\eta) = (1-\alpha)\log(1 - \frac{2\alpha p_c - 2(1-\alpha)p'_n}{p_c}) + \alpha\log(1 + \frac{2\alpha p_c - 2(1-\alpha)p'_n}{p'_n})$$
  
1356 
$$= (1-\alpha)\log(1 - 2\alpha + \frac{2(1-\alpha)}{\frac{p_c}{p'_n}}) + \alpha\log(2\alpha - 1 + 2\alpha\frac{p_c}{p'_n})$$
(25)

Let

1354

1358

1359

1360

1362

1363

1364

1365

1366

1369

1370

1377

$$g_2(t) = (1 - \alpha)\log(1 - 2\alpha + \frac{2(1 - \alpha)}{t}) + \alpha\log(2\alpha - 1 + 2\alpha t)$$
(26)

Take the derivative:

$$g_2'(t) = (1-\alpha)\frac{1}{1-2\alpha + \frac{2(1-\alpha)}{t}}\frac{2(\alpha-1)}{t^2} + \alpha \frac{2\alpha}{2\alpha - 1 + 2\alpha t}$$
(27)

$$=\frac{-2(1-\alpha)^2}{(1-2\alpha)t^2+2(1-\alpha)t}+\frac{2\alpha^2}{2\alpha-1+2\alpha t}$$
(28)

$$=\frac{2(1-2\alpha)(\alpha t+1-\alpha)^2}{[(1-2\alpha)t^2+2(1-\alpha)t](2\alpha-1+2\alpha t)}$$
(29)

Therefore, when  $\frac{1-\alpha}{\alpha} < t < \frac{2(1-\alpha)}{2\alpha-1}$ , we have  $g'_2(t) < 0$ , which implies that  $f(2\eta) < 0$ .

#### B.3 Justification of Proposition 1

We plot the trend of k in Figure 7(a). We compare checkpoints trained for the same iterations on both  $P^c$ and  $P^m$ , where  $p_c$  is calculated based on the loss of the model trained on  $P^c$ , and  $p_n$  is determined by the loss of a model trained for 10,000 iterations on  $P^m$  when evaluated on  $P^n$ . It can be observed that the value of k corresponding to random noise is significantly greater than one, which supports the rationality of the assumption made in Proposition 1.

On the other hand, to extend the proposed theory beyond uniformly distributed random noise (for instance, in multilingual models or Gaussian noise), it is necessary to ensure that k does not become too small in these scenarios. This means that  $\mathcal{L}_{ntp}(P^n, h^*) = -\log p_n$  should not be close to  $\log V$ . One trivial way to increase  $p_n$  is to decrease V, the size of vocabulary. Apart from this, we provide two lines of reasoning to justify why  $p_n$  can be made large:

(1) Numerous studies on compressing large language models, such as pruning (Wang et al., 2020; Kurtic et al., 2022; Zhang et al., 2024), quantization (Zhao et al., 2024; Jin et al., 2024; Liu et al., 2024),

and distillation (Dasgupta et al., 2023; Hinton, 2015), have demonstrated that there exists a significant amount of redundancy within the parameters of large models. Therefore, we could first train a model 1379 on  $P^c$  and then compress it, fine-tuning the surplus parameters on  $P^n$ . This approach would allow us to 1380 improve  $p_n$  without altering  $p_c$ .

(2) A small proportion of data corresponding to  $P^n$  can be introduced into  $P^c$ , making sure that  $\alpha$  is extremely small. According to domain adaptation theory (Ben-David et al., 2010), this would only slightly increase  $\mathcal{L}_{ntp}(P^c)$ . However, existing results (Shliazhko et al., 2024; Pires et al., 2019; Chi et al., 2020) indicate that pre-trained models like BERT or GPT on English text can exhibit strong multilingual capabilities with just a very limited amount of data. Consequently, compared to a model trained solely on  $P^c$ , the resulting model has a minor difference in  $p_c$  but a relatively higher  $p_n$ .

Both thought experiments above demonstrate that there exist a lot of models within the parameter space  $\mathcal{H}$  can perform well on  $P^c$  while yielding non-trivial outcomes on  $P^n$ . Thus, we can ensure that models trained on mixed data distributions will have a sufficiently large k.

Additionally, in Figure 7(b), we illustrate how  $\mathcal{L}_{ntp}(P^m)$  varies with changes in  $\epsilon$ , under settings identical to those used during pre-training. The results depicted in the figure are consistent with our theoretical derivations.

#### **B.4** Omitted Details in Section 4.2

**Definition 1** ( $\beta$ -smooth (Zheng et al., 2023b)). A loss function  $\ell(g_{\theta}(t), y)$  is  $\beta$ -smooth, if for any  $(t, y) \in \mathcal{T} \times \mathcal{Y}$  and any  $\theta, \theta' \in \Theta$ , 1395

$$||\nabla_{\theta}\ell(g_{\theta}(t), y) - \nabla_{\theta'}\ell(g_{\theta'}(t), y)||_2 \le \beta ||\theta - \theta'||_2$$
(30)
$$(30)$$

**Definition 2** ( $\rho$ -input flatness). The  $\rho$ -input flatness  $R_{\rho}(\theta)$  of loss function  $\ell(g_{\theta}(t), y)$  is defined as:

$$R_{\rho}(\theta) = \mathbb{E}_{(t,y)\sim\mathcal{D}} \sup_{\delta'\in B(0,\rho)} \ell(g_{\theta}(t+\delta'), y) - \ell(g_{\theta}(t), y)$$
(31) (31)

*where*  $B(0, \rho) = \{\delta' : ||\delta'||_2 < \rho\}$  *is a open ball.* 

**Lemma 2.** If the loss function  $\ell(g_{\theta}(t), y)$  is  $\beta$ -smooth, then

$$||\nabla_{\theta}\ell(g_{\theta}(t), y)||_2^2 \le 4\beta\ell(g_{\theta}(t), y) \tag{32}$$

Proof. See Lemma 3.1 in Srebro et al. (2010).

**Proposition 2.** Suppose  $\ell(g_{\theta}(t), y)$  is  $\beta$ -smooth with  $\rho$ -input flatness  $R_{\rho}(\theta)$ , for any  $\theta \in \Theta$ :

$$\mathcal{L}_{gm}(\theta) \le 2\beta + 2\mathcal{L}_{ce}(\mathcal{D}, g_{\theta}) + R_{\rho}(\theta) \tag{33}$$

Proof.

$$\mathcal{L}_{gm}(\theta) = ||\mathbb{E}_{(t,y)\sim\mathcal{D}}\nabla_{\theta}\ell(g_{\theta}(t),y) - \mathbb{E}_{(\hat{t},y)\sim\hat{\mathcal{D}}}\nabla_{\theta}\ell(g_{\theta}(\hat{t}),y)||_{2}$$
(34)

$$\leq ||\mathbb{E}_{(t,y)\sim\mathcal{D}}\nabla_{\theta}\ell(g_{\theta}(t),y)||_{2} + ||\mathbb{E}_{(\hat{t},y)\sim\hat{\mathcal{D}}}\nabla_{\theta}\ell(g_{\theta}(\hat{t}),y)||_{2} \qquad \text{(Triangle Inequality)}$$

$$\leq \mathbb{E}_{(t,y)\sim\mathcal{D}}||\nabla_{\theta}\ell(g_{\theta}(t),y)||_{2} + \mathbb{E}_{(\hat{t},y)\sim\hat{\mathcal{D}}}||\nabla_{\theta}\ell(g_{\theta}(\hat{t}),y)||_{2} \qquad \text{(Jensen's Inequality)}$$
1407
1408

$$\leq \mathbb{E}_{(t,y)\sim\mathcal{D}} 2\sqrt{\beta\ell(g_{\theta}(t),y)} + \mathbb{E}_{(\hat{t},y)\sim\hat{\mathcal{D}}} 2\sqrt{\beta\ell(g_{\theta}(\hat{t}),y)}$$
(Lemma 2) 1409

$$\leq \mathbb{E}_{(t,y)\sim\mathcal{D}}(\beta + \ell(q_{\theta}(t),y)) + \mathbb{E}_{(\hat{t},y)\sim\hat{\mathcal{D}}}(\beta + \ell(q_{\theta}(\hat{t}),y))$$
(AM-GM Inequality) 1410

$$= 2\beta + 2\mathbb{F}(x) + \ell(a_0(t), y) + \ell(\mathbb{F}(x) + \ell(a_0(t), y)) + \ell(a_0(t), y))$$

$$= 2\beta + 2\mathbb{F}(x) + \ell(a_0(t), y) + \ell(\mathbb{F}(x) + \ell(a_0(t), y)) + \ell(a_0(t), y))$$
(35)
(35)

$$= 2\beta + 2\mathbb{E}(t,y) \sim \mathcal{D}^{\mathcal{E}}(g\theta(t),g) + (\mathbb{E}(t,y) \sim \mathcal{D}^{\mathcal{E}}(g\theta(t),g) = \mathbb{E}(t,y) \sim \mathcal{D}^{\mathcal{E}}(g\theta(t),g))$$
(55)

$$\leq 2\beta + 2\mathcal{L}_{ce}(\mathcal{D}, g_{\theta}) + R_{\rho}(\theta) \tag{36}$$

where the last inequality holds because  $\hat{t} - t \in B(0, \rho)$ .

1400

1401

1403

1404

1413

1398

1383

1384

1385

1386

1387

1388

1390

1391

1392

#### 1414 C Detailed Related Works

1444

1445

1446

1447

1448

1449

1450

1451

1452

1453

1454

1455

1456

1457

1458

**Data selection for language model training.** Large text corpora form the backbone of language models, 1415 with data quality being fundamental to their success. Elazar et al. (2024) conducted a systematic analysis 1416 of open-source text datasets such as The Pile (Gao et al., 2020) (used to train Pythia), C4 (Raffel et al., 1417 2020) (used to train T5) and RedPajama (used to train LLaMA), revealing that they contain a significant 1418 amount of duplicate, toxic, synthetic, and low-quality content. Therefore, it is of great importance to 1419 thoroughly understand the impact of low-quality data within these pre-training datasets on the model's 1420 performance, reliability, and safety. Allen-Zhu and Li (2024a,b) systematically investigated the effect 1421 of low-quality data and found that such data can significantly reduce the model's knowledge capacity, 1422 sometimes by up to 20 times. Another research direction primarily focuses on the synthetic data of large 1423 language models, specifically examining the impacts of using data generated by LLMs for recursive 1424 training. The study by Shumailov et al. (2024) was the first to explore this issue and introduced the concept 1425 of "model collapse", indicating that recursive training can lead to the loss of information in tail tokens, 1426 ultimately resulting in the model producing nonsensical content. Seddik et al. (2024) mainly provided 1427 a theoretical explanation for why model collapse occurs, supporting their arguments with experimental 1428 evidence. Consequently, the importance of data selection cannot be overstated. Given that data selection 1429 is an NP-hard problem in terms of combinatorial optimization (Xiao et al., 2022), numerous heuristic 1430 algorithms have been proposed to expedite the process. Longpre et al. (2024) provided a comprehensive 1431 study on pre-training data selection and optimal ratios, offering practical recommendations. Yang et al. 1432 (2023) proposed dataset pruning, an approach that assesses the impact of omitting training samples on 1433 model generalization and creates a minimal training subset with a controlled generalization gap. Chai et al. 1434 (2024) evaluated the impact of individual training samples on the dynamics of GPT model training. Li et al. 1435 (2024b) introduced the Instruction-Following Difficulty metric to assess the quality of instruction-tuning 1436 data. Xie et al. (2023b) employed importance resampling for data selection. Xie et al. (2023a); Lee 1437 et al. (2023) advocated for optimizing data composition and diversity. Despite these notable studies on 1438 data selection, they generally acknowledge that dataset noise degenerates model performance but lack 1439 a detailed understanding of how and to what extent, particularly in the case of random noise which is 1440 inevitable in large-scale datasets. Although Cherepanova and Zou (2024) investigated the influence of 1441 gibberish input, the random noise within the pre-training dataset is still underexplored. This paper aims to 1442 1443 bridge the gap.

Learning from Noisy Distributions. The majority of machine learning algorithms assume that training and test samples are independently and identically distributed (i.i.d.), a condition that is often not met in real-world scenarios. For instance, LLMs are pre-trained on datasets with all kinds of noise while their performance is evaluated by the user whose distribution is usually clean and meets real-world scenarios, which violates the i.i.d. assumption. Domain adaptation (Li et al., 2024a; Meng et al., 2022; Ma et al., 2023) addresses this issue when the distribution of the training data differs from that of the test data. Although domain adaptation methods attempt to reduce the statistical distribution discrepancy (Du and Li, 2023; Wu et al., 2024) or employ adversarial training (Li et al., 2021; Ru et al., 2024) to minimize the gap between source and target domains, they typically require access to unlabeled test data under a semi-supervised learning setup, which is impractical for LLM training. Another reason domain adaptation cannot be directly applied here is that domain adaptation theory (Ben-David et al., 2010) focuses on the performance of a model trained on one distribution when it is applied to another different but related distribution. This kind of bounds can be easily derived by Lemma 1. However, what we aim to investigate here is the extent of performance loss when comparing a model trained on one distribution (noisy dataset) to a model trained on another distribution (clean dataset).

Apart from domain adaptation, there has been extensive research directly investigating noisy training sets. Noisy label learning Song et al. (2022); Lukasik et al. (2020) have explored the impact of incorrect labels on model performance. Regarding input feature noise, Smilkov et al. (2017) added perturbations to individual image inputs to enhance model interpretability, and Zada et al. (2022) added white noise image into the training dataset to tackle the class imbalance problem. However, most of these efforts have concentrated on image classification and do not consider the pre-training paradigm.



Figure 8: The prior distribution of tokens in the data from (a) OpenWebText, (b) random noise, and (c) Gaussian noise.

**Fine-tuning Pre-trained Models.** The approach of initially pre-training model weights on large-scale 1465 datasets and subsequently fine-tuning them with downstream data has become the de facto standard in the 1466 fields of computer vision (Kornblith et al., 2019; Raghu et al., 2019) and natural language processing (Wei 1467 et al., 2021; Xie et al., 2021a). For instance, Hua et al. (2023, 2021) proposed enhancing the performance 1468 of models by increasing their resistance to minor perturbations in intermediate layers. Meanwhile, 1469 Jiang et al. (2020) improved model robustness by adding regularization terms. Besides full-parameter 1470 fine-tuning, numerous parameter-efficient fine-tuning algorithms have been extensively studied. Zhang 1471 et al. (2021b) introduced adapters into the original model architecture, optimizing only these parameters 1472 during fine-tuning. Zhou et al. (2022b,a) efficiently fine-tuned CLIP models (Radford et al., 2021) using 1473 learnable soft prompts. Hu et al. (2022) optimized models through learning low-rank residual weights. 1474 These methods achieved performance close to that of full-parameter fine-tuning while maintaining the 1475 generalization ability of the original models. However, they all require access to the model's weights and 1476 loading them into GPU memory, which can be challenging for today's large models, especially when 1477 state-of-the-art models' parameters are not publicly available. Therefore, in this paper, we follow the 1478 NML setup and explore efficient ways to fine-tune the downstream task head under a black-box scenario. 1479

**Implicit Regularization and Sharpness-aware Minimization.** Achieving good generalization in neural networks optimized using gradient descent algorithms has long been a research focus in deep learning theory. Barrett and Dherin (2021) explored the properties of stochastic gradient descent (SGD), finding that SGD implicitly constrains the gradient norm. Based on this observation, Sharpness-aware minimization (SAM) (Zhang et al., 2023; Foret et al., 2021; Wen et al., 2023; Xie et al., 2023c) improves generalization by incorporating the gradient norm as a regularization term. Our method can be seen as drawing inspiration from SAM but differs in that our optimization objective is the model's resilience to input noise rather than seeking flat minima in the parameter space.

#### **D** Experiments in Section **3**

#### **D.1** Pre-training Dataset

**OpenWebText Dataset.** The OpenWebText dataset (Gokaslan et al., 2019) is a large-scale corpus of English text data, developed to serve as an open-access alternative to proprietary dataset WebText that is utilized by OpenAI for training their GPT-2 models. This dataset originates from the analysis of outbound links clicked on Reddit, undergoing multiple stages of filtering to exclude non-English content, duplicate entries, copyrighted materials, and texts lacking in quality. These links generally direct to web pages available to the public, often shared or debated on Reddit, thereby covering a broad spectrum of subjects that mirror online popular interests and discussions. The dataset includes roughly 18 million documents, amounting to about 20GB of compressed plain text data in uint16 format. Since measures have been implemented to ensure the dataset's reliability by filtering out unsuitable content, we consider it a clean and noise-free dataset. Figure 8(a) illustrates the distribution of internal tokens.

Random Noise. To simulate the distribution of random gibberish that crawlers might retrieve from1500the Internet due to various reasons, we manually searched and collected a few websites containing such1501gibberish and also opened normally functioning websites using different decoding methods to observe1502

23

1488

1490

1491

1492

1493

1494

1495

1496

1497

1498

1499

1480

1481

1482

1483

1484

1485



Figure 9: Visualization of (a)  $\mathcal{L}_{ntp}(P^m, h^*) - \mathcal{L}_{ntp}(P^m, h)$  and (b)  $\mathcal{L}_{ntp}(P^c, h)$  with  $\alpha = 0.55, k = 1, \mathcal{L}_{ntp}(P^c, h^*) = 2.9, \mathcal{L}_{ntp}(P^n, h^*) = 2.8.$ 

the distribution of tokens. We found that, while the distribution of tokens appeared disorganized, their prior probabilities were not evenly distributed. Instead, several tokens had notably high probabilities, which is similar to that observed in the clean data as shown in Figure 8(a). Thus, on the basis of uniformly distributed random noise, we increased the frequency of certain tokens and then randomized them again. The resulting distribution is illustrated in Figure 8(b). It can be seen that while maintaining an overall uniform distribution, the frequency of tokens with IDs ranging from 0 to 1000 is higher which closely mirrors real-world scenarios.

**Gaussian Noise.** Given the diversity and unpredictability of real-world data distributions, we also artificially generated random noise with a prior probability that follows a Gaussian distribution, as shown in Figure 8(c). The rationale behind choosing the Gaussian distribution is that the noise in many real-world systems can be approximated by it. Additionally, we set the standard deviation  $\sigma = 500$  to simulate scenarios where random noise exhibits sharp peaks.

#### D.2 Training Details of GPT-2

1503

1504

1505

1506

1508

1511

1512

1513

1514

1516

1517

1518

1519

1522

1523

1524

1533

1534

1535

Our work is based on the source code of nanoGPT<sup>3</sup>. Specifically, we utilized the GPT-2 tokenizer with vocabulary size V = 50256 to process the OpenWebText dataset, and then appended randomly generated noise to the end of the training set before commencing training. The model's context length is set to L = 1024, with an embedding dimension d = 768. The GPT-2 model consists of 12 self-attention layers, totaling approximately 124 million parameters. For optimization, we employed AdamW (Loshchilov, 2017; Xie et al., 2022) with a learning rate of 6e-4, weight decay of 0.1, and  $\beta$  values of 0.9 and 0.95 for  $\beta_1$  and  $\beta_2$ , respectively. A cosine annealing scheduler was used to gradually adjust the learning rate down to 6e-5. We configured the batch size to 16, with a gradient accumulation step of 40, allowing each iteration to process 655,360 tokens (16 \* 40 \* 1024). Training proceeded for a total of 300,000 iterations.

#### 1525 D.3 Synthetic Results about Multilingual Models

To illustrate our theory's explanatory power concerning multilingual models, we have plotted the scenario where  $h^*$  is influenced by  $P^n$  under the conditions  $\alpha = 0.55$ , k = 1,  $\mathcal{L}_{ntp}(P^c, h^*) = 2.9$ , and  $\mathcal{L}_{ntp}(P^n, h^*) = 2.8$ , as shown in Figure 9. This setup simulates a model trained on a roughly 1:1 multilingual corpus, where the capacity of one language is affected by the data from another language. As can be observed from the figure, the impact on  $p_c$  does not exceed  $2\eta = 2(\alpha p_c - (1 - \alpha)p_n)$ , which translates to an increase of no more than 0.1 in  $\mathcal{L}_{ntp}(P^n, h)$ . This finding strongly supports the success of multilingual models from a theoretical perspective.

#### D.4 Hardware

We conducted the pre-training process on a server equipped with 8 NVIDIA GeForce RTX 4090 GPUs. It takes approximately 70 hours to train one model using eight 4090 GPUs, so pre-training five GPT-2

<sup>&</sup>lt;sup>3</sup>https://github.com/karpathy/nanoGPT

models in total requires 2,800 GPU hours. 1536 **Experiments in Section 4** Ε 1537 **E.1** Detailed Setup for Downstream Natural Language Understanding Experiments 1538 E.1.1 Datasets 1539 We utilize 8 commonly-used text classification benchmark: SST-2, SST-fine, 20newsgroup, CR, BBC, 1540 Balanced COPA, MRPC, WiC. The detailed information can be found in Table 8. 1541

Dataset	Classes	Train Size	Test Size
SST-2 (Socher et al., 2013)	2	6.92k	1.82k
SST-fine (Chen and Manning, 2014)	5	8.54k	2.21k
20newsgroup (Zhang et al., 2019)	20	11.3k	7.53k
CR (Hu and Liu, 2004)	2	3.39k	376
BBC (Samuels and Mcgonical, 2020)	5	1.23k	1k
Balanced COPA (Kavumba et al., 2020)	2	1k	500
MRPC (Dolan and Brockett, 2005)	2	3.67k	1.73k
WiC (Pilehvar and Camacho-Collados, 2019)	2	5.43k	1.4k

Table 8: Details of the 8 natu	ral language understanding dataset.
--------------------------------	-------------------------------------

#### E.1.2 Prompts

Since classification tasks can be processed as seq2seq tasks by adding prompts (Sutskever, 2014; Saunshi et al., 2021), we design a unique prompt for each dataset and task. This approach transforms the inputs into a format that large language models can process. The specific designs are shown in Table 9.

Dataset	Prompts
SST-2	{text} this movie is
SST-fine	{text} this movie is
20newsgroup	{text} This article is about
CR	{text} the sentiment is
BBC	Please classify the topic of the following news: {text} Answer:
Balanced COPA	Given the premise: {premise} Find the most plausible alternative
	for the {question}. Option 1: {choice1} Option 2: {choice2}
	Which option is more plausible?
MRPC	Sentence 1: {text1} Sentence 2: {text2} Is this a paraphrase?
WiC	Task: Determine if the word {phrase1} has the same meaning in
	the two sentences below. Sentence 1: {sentence1} Sentence 2:
	{sentence2} Your answer:

Table 9: Details of the prompts applied to each dataset.

#### E.1.3 Hyperparameters

For all experiments in Section 4, we utilize a two-layer MLP with hidden dimension equals to feature dimension and ReLU activation function.

For all experiments shown in Table 1, we set  $\gamma$  in Equation (3) to be 0.01 and  $\lambda$  in Equation (6) to be 1549 0.15. Following the setup as described by Saunshi et al. (2021), for each dataset, we conduct a grid search 1550 on the validation set to identify the optimal learning rate and batch size. We train for a total of ten epochs 1551 with the learning rate ranging within  $\{3e-4, 6e-4\}$  and batch size options including  $\{8, 12, 16, 32\}$ . For 1552 samples without a designated validation set, we randomly select 10% of the training set samples to form a validation set for the purpose of selecting the best parameters. 1554

1542 1543

1544

1545

1546

1547

For the experiments listed in Table 2, we set the batch size to 8 and the learning rate to 6e-4 for all linear probe tasks. For all MLP probe tasks, the learning rate is set to 1e-4. Regarding  $\gamma$  and  $\lambda$ , we conduct a grid search on the validation set to find the optimal values.

1558

# E.2 Detailed Setup for Downstream Vision Experiments

# 1559 E.2.1 Datasets

We select 14 image classification datasets, which serve as a common benchmark for evaluating model
 performance in the vision community (Zhou et al., 2022b; Chen et al., 2024). Specific information about these 14 datasets is provided in Table 10.

Dataset	Classes	Train Size	Test Size
StanfordCars (Krause et al., 2013)	196	8144	8041
Caltech101 (Fei-Fei et al., 2004a)	102	3060	6084
CIFAR-10 (Krizhevsky et al., 2009)	10	50000	10000
CIFAR-100 (Krizhevsky et al., 2009)	100	50000	10000
DTD (Cimpoi et al., 2014)	47	1880	1880
EuroSAT (Helber et al., 2019)	10	21600	5400
FGVCAircraft (Maji et al., 2013)	102	6667	3333
Flowers102 (Nilsback and Zisserman, 2008)	102	2040	6149
Food101 (Fei-Fei et al., 2004b)	101	75750	25250
OxfordPet (Parkhi et al., 2012)	37	3680	3669
PatchCamelyon (Veeling et al., 2018)	2	262144	32768
RESISC45 (Cheng et al., 2017)	45	25200	6300
Rendered SST2 (Socher et al., 2013)	2	6920	1821
SVHN (Netzer et al., 2011)	10	73257	26032

Table 10: Details of the 14 vision dataset.

# 1563 E.2.2 Models

We use five pre-trained general-purpose visual backbone models that cover a variety of architectures, datasets, and training methods. Detailed information is provided in Table 11.

Pre-training Dataset	Size	
ImageNet-1K (Deng et al., 2009)	12.3M	
and JFT-300M (Sun et al., 2017)		
ImageNet-21K (Ridnik et al., 2021)	321.7M	
ImageNet-21K	196.7M	
Laion-2B (Schuhmann et al., 2022)	200.114	
and ImageNet-1K	200.1M	
Laion-2B	428M	
	Pre-training Dataset ImageNet-1K (Deng et al., 2009) and JFT-300M (Sun et al., 2017) ImageNet-21K (Ridnik et al., 2021) ImageNet-21K Laion-2B (Schuhmann et al., 2022) and ImageNet-1K Laion-2B	

Table 11: Details of the 5 vision models.

#### 1565

1566

1562

# E.2.3 Hyperparameters

In our study, similar to the approach detailed in Chen et al. (2024), we primarily contrast our proposed method with MLP and LP tuning. For the optimization process, we employ AdamW for fine-tuning the modules over 30 epochs, utilizing a cosine learning rate scheduler. Specifically, for LP, we configure the learning rate at 0.1 without applying any weight decay. In contrast, both the MLP tuning and our method use a more conservative learning rate of 0.001 alongside a weight decay of 1e-4.

# E.2.4 Detailed Experimental Results

In Table 3, due to space limitations, we only present the average results, while detailed results are shown in Table 12.

Models	StanfordCars		Caltech101		CIFAR-10		
	Linear	MLP	Linear	MLP	Linear	MLP	
ViT-L	$93.38\pm0.76$	$94.41 \pm 1.05$	$92.07 \pm 1.19$	$95.20\pm1.12$	$97.99 \pm 0.95$	$98.35\pm0.95$	
ViT-L + $\mathcal{L}_{gm}$	$\textbf{93.71} \pm \textbf{1.37}$	$\textbf{94.56} \pm \textbf{1.50}$	$\textbf{95.01} \pm \textbf{0.97}$	$\textbf{95.29} \pm \textbf{1.27}$	$\textbf{98.07} \pm \textbf{0.74}$	$\textbf{98.48} \pm \textbf{0.60}$	
ConvNext-L	$86.01 \pm 1.48$	$88.68 \pm 0.89$	$91.02\pm0.79$	$94.47\pm0.53$	$97.49 \pm 1.36$	$98.09 \pm 0.85$	
ConvNext-L+ $\mathcal{L}_{gm}$	$\textbf{86.78} \pm \textbf{1.32}$	$\textbf{89.06} \pm \textbf{1.19}$	$\textbf{94.11} \pm \textbf{1.19}$	$\textbf{94.93} \pm \textbf{0.88}$	$\textbf{97.59} \pm \textbf{0.52}$	$\textbf{98.15} \pm \textbf{0.71}$	
EfficientNet-B3	$56.20\pm0.54$	$\textbf{58.57} \pm \textbf{1.11}$	$89.43 \pm 0.78$	$91.22\pm1.23$	$94.04 \pm 1.19$	$95.73 \pm 1.07$	
EfficientNet-B3+ $\mathcal{L}_{gm}$	$\textbf{57.02} \pm \textbf{1.28}$	$58.15 \pm 1.12$	$\textbf{90.25} \pm \textbf{1.43}$	$\textbf{91.55} \pm \textbf{0.90}$	$\textbf{94.11} \pm \textbf{0.88}$	$\textbf{95.96} \pm \textbf{0.86}$	
ResNetv2-152x2	$56.95 \pm 1.21$	$\textbf{59.18} \pm \textbf{1.34}$	$91.40 \pm 1.47$	$92.48 \pm 1.30$	$96.28 \pm 1.16$	$96.91\pm0.89$	
ResNetv2-152x2+ $\mathcal{L}_{gm}$	$\textbf{58.78} \pm \textbf{1.16}$	$58.67 \pm 1.27$	$\textbf{93.83} \pm \textbf{0.91}$	$\textbf{93.95} \pm \textbf{0.94}$	$\textbf{96.31} \pm \textbf{0.85}$	$\textbf{97.03} \pm \textbf{0.53}$	
Swin-L	$68.17 \pm 0.98$	$\textbf{74.11} \pm \textbf{0.60}$	$92.58 \pm 0.95$	$94.09 \pm 1.04$	$98.26 \pm 0.89$	$98.61 \pm 0.78$	
Swin-L+ $\mathcal{L}_{gm}$	$\textbf{69.31} \pm \textbf{1.07}$	$73.71\pm0.94$	$\textbf{93.65} \pm \textbf{1.42}$	$\textbf{94.62} \pm \textbf{0.64}$	$\textbf{98.41} \pm \textbf{0.91}$	$\textbf{98.72} \pm \textbf{1.28}$	

CIFA	CIFAR-100 EuroSAT		FGVCAircraft		OxfordPet			
Linear	MLP	Linear	MLP	Linear	MLP	Linear	MLP	
$\textbf{88.07} \pm \textbf{0.58}$	$89.49 \pm 0.52$	$97.53 \pm 1.13$	$97.75\pm0.61$	$65.76 \pm 0.73$	$68.43 \pm 0.78$	$91.65 \pm 1.18$	$93.97 \pm 1.50$	
$88.06\pm0.93$	$\textbf{89.58} \pm \textbf{0.93}$	$\textbf{97.83} \pm \textbf{0.73}$	$\textbf{98.03} \pm \textbf{0.60}$	$\textbf{66.63} \pm \textbf{1.08}$	$\textbf{68.67} \pm \textbf{1.07}$	$\textbf{93.18} \pm \textbf{0.81}$	$\textbf{94.17} \pm \textbf{1.10}$	
$\textbf{86.76} \pm \textbf{0.91}$	$87.79 \pm 1.27$	$95.57 \pm 1.22$	$96.31 \pm 1.11$	$57.18 \pm 1.12$	$62.25\pm0.66$	$94.98 \pm 0.54$	$95.80\pm0.75$	
$86.46 \pm 1.04$	$\textbf{87.88} \pm \textbf{1.24}$	$\textbf{96.05} \pm \textbf{1.46}$	$\textbf{96.74} \pm \textbf{1.27}$	$\textbf{58.35} \pm \textbf{0.68}$	$\textbf{63.61} \pm \textbf{0.83}$	$\textbf{95.39} \pm \textbf{1.21}$	$\textbf{95.99} \pm \textbf{0.92}$	
$\textbf{77.34} \pm \textbf{0.86}$	$80.28 \pm 1.02$	$94.81 \pm 1.35$	$95.90\pm0.88$	$44.73 \pm 1.31$	$46.23\pm0.92$	$93.84 \pm 1.14$	$94.79 \pm 1.14$	
$77.16 \pm 1.11$	$\textbf{80.47} \pm \textbf{1.43}$	$\textbf{95.20} \pm \textbf{0.52}$	$\textbf{96.07} \pm \textbf{1.18}$	$\textbf{45.33} \pm \textbf{0.61}$	$\textbf{47.07} \pm \textbf{0.57}$	$\textbf{94.63} \pm \textbf{1.03}$	$\textbf{94.98} \pm \textbf{1.14}$	
$\textbf{84.30} \pm \textbf{1.18}$	$\textbf{84.68} \pm \textbf{1.33}$	$97.12 \pm 1.46$	$97.46 \pm 1.28$	$42.03\pm0.72$	$48.39\pm0.85$	$91.93\pm0.68$	$92.99 \pm 1.40$	
$84.28 \pm 1.22$	$84.29 \pm 1.38$	$\textbf{97.35} \pm \textbf{1.12}$	$\textbf{97.59} \pm \textbf{0.72}$	$\textbf{45.69} \pm \textbf{0.80}$	$\textbf{48.84} \pm \textbf{0.61}$	$\textbf{92.61} \pm \textbf{0.87}$	$\textbf{93.45} \pm \textbf{1.32}$	
$89.68 \pm 1.33$	$90.74\pm0.98$	$\textbf{97.11} \pm \textbf{0.72}$	$97.59\pm0.63$	$54.96 \pm 1.24$	$\textbf{61.17} \pm \textbf{1.35}$	$92.17\pm0.64$	$94.38 \pm 1.21$	
$\textbf{89.79} \pm \textbf{0.52}$	$\textbf{91.18} \pm \textbf{1.21}$	$97.09 \pm 1.11$	$\textbf{97.71} \pm \textbf{1.08}$	$56.10 \pm 0.67$	$60.99\pm0.73$	$\textbf{93.86} \pm \textbf{0.85}$	$\textbf{94.57} \pm \textbf{1.11}$	

Food	Food101 Flowers102		D	ГD	SVHN		
Linear	MLP	Linear	MLP	Linear	MLP	Linear	MLP
$90.51 \pm 1.31$	$91.04 \pm 1.35$	$94.04 \pm 1.13$	$97.83 \pm 0.74$	$80.53 \pm 1.06$	$83.29\pm0.86$	$78.82 \pm 1.25$	$84.74 \pm 1.08$
$\textbf{90.62} \pm \textbf{1.37}$	$\textbf{91.23} \pm \textbf{0.52}$	$\textbf{96.67} \pm \textbf{1.41}$	$\textbf{98.06} \pm \textbf{1.28}$	$\textbf{82.76} \pm \textbf{0.71}$	$\textbf{83.77} \pm \textbf{0.98}$	$\textbf{79.80} \pm \textbf{0.88}$	$\textbf{84.59} \pm \textbf{1.38}$
$\textbf{89.09} \pm \textbf{1.06}$	$\textbf{90.21} \pm \textbf{0.87}$	$94.71 \pm 1.31$	$98.78 \pm 1.17$	$76.01 \pm 1.03$	$78.67 \pm 1.15$	$66.16\pm0.87$	$72.76\pm0.78$
$88.62\pm0.72$	$90.10\pm1.39$	$\textbf{97.12} \pm \textbf{0.95}$	$\textbf{98.99} \pm \textbf{0.99}$	$\textbf{77.92} \pm \textbf{0.72}$	$\textbf{80.05} \pm \textbf{1.44}$	$\textbf{68.43} \pm \textbf{1.37}$	$\textbf{73.18} \pm \textbf{0.67}$
$\textbf{76.95} \pm \textbf{0.82}$	$\textbf{81.78} \pm \textbf{0.89}$	$84.19\pm0.61$	$88.97 \pm 1.32$	$69.09 \pm 1.27$	$73.08 \pm 1.30$	$54.26 \pm 0.87$	$61.38 \pm 0.82$
$76.35\pm0.74$	$81.33 \pm 1.45$	$\textbf{86.19} \pm \textbf{1.08}$	$\textbf{89.42} \pm \textbf{0.62}$	$\textbf{71.27} \pm \textbf{1.24}$	$\textbf{73.82} \pm \textbf{1.40}$	$\textbf{56.74} \pm \textbf{0.73}$	$\textbf{63.56} \pm \textbf{0.82}$
$\textbf{84.83} \pm \textbf{1.36}$	$84.15 \pm 1.27$	$96.76\pm0.88$	$98.27\pm0.53$	$72.23 \pm 0.94$	$76.11 \pm 1.25$	$60.75\pm0.89$	$64.87 \pm 1.45$
$84.41\pm0.88$	$\textbf{84.64} \pm \textbf{1.04}$	$\textbf{98.08} \pm \textbf{1.23}$	$\textbf{98.84} \pm \textbf{0.82}$	$74.73 \pm 1.38$	$\textbf{77.12} \pm \textbf{1.39}$	$\textbf{62.04} \pm \textbf{1.01}$	$\textbf{65.06} \pm \textbf{0.62}$
$90.23\pm0.64$	$92.23\pm0.55$	$97.28 \pm 0.92$	$99.51 \pm 1.17$	$75.85\pm0.88$	$80.74 \pm 1.45$	$62.77 \pm 1.42$	$\textbf{69.53} \pm \textbf{1.22}$
$\textbf{90.26} \pm \textbf{1.14}$	$\textbf{92.32} \pm \textbf{1.00}$	$\textbf{99.12} \pm \textbf{1.05}$	$\textbf{99.60} \pm \textbf{0.51}$	$\textbf{77.44} \pm \textbf{1.00}$	$\textbf{80.91} \pm \textbf{0.83}$	$\textbf{64.83} \pm \textbf{0.98}$	$68.97 \pm 1.12$

resis	resisc45 rsst2		pca	Avg			
Linear	MLP	Linear	MLP	Linear	MLP	Linear	MLP
$95.44 \pm 1.41$	$95.79\pm0.57$	$67.65 \pm 1.12$	$73.58 \pm 1.31$	$\textbf{82.65} \pm \textbf{0.57}$	$\textbf{83.92} \pm \textbf{0.56}$	86.86	89.12
$\textbf{95.73} \pm \textbf{1.19}$	$\textbf{95.93} \pm \textbf{1.28}$	$\textbf{71.82} \pm \textbf{0.55}$	$\textbf{74.24} \pm \textbf{0.76}$	$82.55 \pm 1.49$	$83.78\pm0.97$	88.03	89.31
$92.65 \pm 1.25$	$93.09\pm0.89$	$60.73 \pm 1.30$	$66.0\pm0.55$	$72.21\pm0.72$	$77.08 \pm 0.53$	82.89	85.71
$\textbf{92.93} \pm \textbf{0.54}$	$\textbf{93.31} \pm \textbf{1.31}$	$\textbf{64.14} \pm \textbf{0.99}$	$\textbf{67.49} \pm \textbf{0.63}$	$\textbf{73.1} \pm \textbf{1.35}$	$\textbf{78.31} \pm \textbf{1.10}$	84.07	86.27
$87.19 \pm 1.05$	$89.01\pm0.70$	$\textbf{50.46} \pm \textbf{1.02}$	$50.74 \pm 1.05$	$53.32\pm0.59$	$\textbf{51.10} \pm \textbf{1.34}$	73.27	75.62
$\textbf{87.33} \pm \textbf{0.60}$	$\textbf{89.17} \pm \textbf{1.34}$	$50.19\pm0.74$	$\textbf{51.07} \pm \textbf{0.92}$	$\textbf{54.52} \pm \textbf{1.27}$	$50.10\pm0.92$	74.02	75.90
$90.96 \pm 0.81$	$91.19\pm0.57$	$50.90 \pm 0.73$	$49.91 \pm 1.08$	$77.62\pm0.85$	$77.86 \pm 1.16$	78.14	79.60
$\textbf{91.17} \pm \textbf{0.64}$	$\textbf{91.36} \pm \textbf{0.73}$	$\textbf{54.25} \pm \textbf{0.69}$	$\textbf{49.92} \pm \textbf{1.03}$	$\textbf{79.44} \pm \textbf{0.52}$	$\textbf{78.48} \pm \textbf{0.53}$	79.49	79.45
$92.79 \pm 1.25$	$94.09\pm0.99$	$50.96 \pm 1.41$	$53.59 \pm 1.20$	$77.32\pm0.68$	$78.38 \pm 1.04$	81.43	84.19
$\textbf{93.41} \pm \textbf{1.32}$	$\textbf{94.42} \pm \textbf{0.96}$	$\textbf{53.48} \pm \textbf{0.89}$	$\textbf{54.96} \pm \textbf{0.90}$	$\textbf{81.18} \pm \textbf{1.34}$	$\textbf{79.29} \pm \textbf{1.31}$	82.70	84.42

Table 12: Detailed accuracy of 5 vision backbone models on 14 commonly-used vision datasets.

1572 1573

## 1575 E.3 Runtime Analysis

Since all our models are black-box models, we first process all samples into vector-form features and then
 probe them. All models described in this paper can run on a single NVIDIA RTX 4090 GPU. Extracting
 all these features requires a total of 10 GPU hours. Subsequently, training these Linear or MLP Probes
 requires approximately 200 GPU hours in total.