NLoRA: Nyström-Initiated Low-Rank Adaptation for Large Language Models

Anonymous ACL submission

Abstract

Parameter-efficient fine-tuning (PEFT) is essential for adapting large language models (LLMs), 003 with low rank adaptation (LoRA) being the most popular approach. However, LoRA suffers from slow convergence, and some recent LoRA variants, such as PiSSA, primarily rely on Singular Value Decomposition (SVD) for initialization, leading to expensive computation. To mitigate these problems, we resort to Nyström method, which follows a threematrix manipulation. Therefore, we first introduce StructuredLoRA (SLoRA), investigat-013 ing to introduce a small intermediate matrix between the low-rank matrices A and B. Secondly, we propose NyströmLoRA (NLoRA), which leverages Nyström-based initialization for SLoRA to improve its effectiveness and effi-017 ciency. Finally, we propose IntermediateTune (IntTune) to explore fine-tuning exclusively the intermediate matrix of NLoRA to furthermore boost LLMs' efficiency. We evaluate our meth-022 ods on 5 natural language generation (NLG) tasks and 8 natural language understanding (NLU) tasks. On GSM8K, SLoRA and NLoRA achieve accuracies of 56.48% and 57.70%, surpassing LoRA by 33.52% and 36.41% with 027 only 3.67M additional trainable parameters. IntTune boosts average NLG performance over LoRA by 7.45% while using only 1.25% of its parameters. These results demonstrate the efficiency and effectiveness of our approach in enhancing model performance with minimal parameter overhead.

1 Introduction

034

037

041

Fine-tuning large language models (LLMs) has emerged as a fundamental approach to enhancing model capabilities (Yu et al., 2023; Li et al., 2023; Xia et al., 2024) and aligning models with specific application requirements (Zheng et al., 2023; Ouyang et al., 2022). However, the growing scale of LLMs introduces significant challenges to LLM



Figure 1: The comparison among LoRA and our models

042

043

044

047

048

054

056

060

061

062

063

064

065

066

067

068

070

development, with fine-tuning requiring substantial computational and memory resources (Hu et al., 2021; Chang et al., 2024). For example, fine-tuning a LLaMA-65B model requires more than 780 GB of GPU memory (Dettmers et al., 2023), while training GPT-3 175B requires 1.2 TB of VRAM (Hu et al., 2021). Such resource-intensive processes are infeasible for many researchers and institutions, driving the development of parameterefficient fine-tuning (PEFT) methods. Among these methods, Low-Rank Adaptation (LoRA) (Hu et al., 2021) has received widespread attention due to its ability to achieve competitive performance compared to full parameter fine-tuning, while significantly reducing memory consumption and avoiding additional inference latency.

LoRA enables the indirect training of dense layers in a neural network by optimizing low-rank decomposition matrices that represent changes in the dense layers during adaptation, all while keeping the pre-trained weights fixed. For a pre-trained weight matrix $W \in \mathbb{R}^{m \times n}$, LoRA introduces a low-rank decomposition $\Delta W = AB$, where $A \in \mathbb{R}^{m \times r}$, $B \in \mathbb{R}^{r \times n}$, and the rank $r \ll \min(m, n)$. This modifies the forward pass of a layer as follows:

$$Y = X(W + \Delta W) = X(W + AB), \quad (1)$$

where $X \in \mathbb{R}^{b \times m}$, $Y \in \mathbb{R}^{b \times n}$, and b represents the batch size. For initialization, A is randomly ini-



Figure 2: The comparison among Full Fine-tuning, LoRA, and SLoRA

tialized with Gaussian values and *B* is set to zero, ensuring that injection of the low-rank adaptation does not alter the model predictions at the start of training. Unlike traditional fine-tuning methods that require updating and storing gradients for the full weight matrix *W*, LoRA optimizes only the smaller matrices *A* and *B*, significantly reducing the number of trainable parameters and memory usage. Furthermore, LoRA often achieves performance comparable or superior to full fine-tuning, demonstrating that adapting only a small subset of parameters suffices for many downstream tasks.

071

077

084

091

100

101

102

104

106

Despite the above benefits, LoRA suffers from slow convergence (Ding et al., 2023). To address this issue, some recent LoRA variants, such as PiSSA (Meng et al., 2024), choose to conduct initialization of the low rank matrices by using Singular Value Decomposition (SVD). However, SVDbased initialization is computationally expensive and requires a long time. To mitigate this issue, we investigate using Nyström method, which approximates a matrix as a product of three matrices, to approximate SVD. To fit the three-matrix structure, we first propose StructuredLoRA (SLoRA), where an additional $r \times r$ matrix is inserted between the low-rank matrices A and B, as shown in Figure 2. Furthermore, we explore whether an extra matrix can influence the language model's performance, experimental results indicate that SLoRA effectively enhances performance with only a minor increase in the number of parameters, demonstrating the potential of the three-matrix structure for PEFT.

Secondly, inspired by NyströmFormer (Xiong et al., 2021), we proposed NyströmLoRA (NLoRA) to leverage Nyström method, which conducts SVD approximation by sampling a subset of rows and columns of the pre-trained parameter matrix to reduce the computational cost, for weight initialization. NLoRA is supposed to bypass the computational cost of SVD's eigenvalue decomposition, reducing time complexity to $O(mr+r^2+rn)$ compared to the $O(mn^2)$ complexity of SVD-based methods. 107

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

Finally, to explore whether we can further compress the trainable parameters of NLoRA, we propose **Int**ermediate**Tune** (IntTune), which exclusively adjusts the intermediate matrix of NLoRA. This method significantly reduces the number of trainable parameters. Specifically, on the evaluation of LLaMA 2-7B across five NLG benchmarks, LoRA uses 320M parameters, while our IntTune method only requires tuning 4M parameters. In the meantime, IntTune outperforms LoRA by 7.45% on average across NLG benchmarks. The comparison of our proposed methods with LoRA in terms of performance and trainable parameters is illustrated in Figure 1.

In summary, our contributions are as follows:

- 1. We propose SLoRA, an extension to the LoRA framework, incorporating an additional intermediate matrix to enhance model expressiveness, achieving improved performance with minimal parameter overhead.
- 2. We introduce NLoRA, leveraging Nyström approximation for efficient and effective initialization, particularly excelling in natural language generation (NLG) and natural language understanding (NLU) tasks.
- 3. We propose IntTune to fulfill supervised finetuning (SFT) LLaMA 2-7B by tuning 4M 141

142 143

- 144

- 145

146

147

148

149

150

152

153

154

155

156

158

159

160

162

164

165

166

167

168

171

172

173

174

175

176

177

178

179

181

184

185

186

187

189

190

parameters, achieving superior performance compared to LoRA on average, offering a lightweight and efficient alternative for SFT LLMs in resource-constrained scenarios.

2 **Related Works**

2.1 LoRA's variants

With the introduction of LoRA (Hu et al., 2021), many derivative methods have emerged. AdaLORA (Zhang et al., 2023) highlights that LoRA ignores the importance of different layer parameters based on a uniform setting of the rank, and proposes an adaptive allocation strategy based on parameter importance to improve fine-tuning efficiency. DoRA (Liu et al., 2024) introduces a decomposation of weight matrices into magnitude and direction components, leveraging LoRA to update only the directional component, thereby reducing the number of trainable parameters. ReLoRA (Lialin et al., 2023) achieves high-rank training through iterative low-rank updates, periodically merging parameters into the main model. LoRA+ (Hayou et al., 2024) further improves efficiency by applying different learning rates to the two matrices in LoRA, assigning a higher learning rate to matrix B to accelerate convergence and enhance performance. Other works have focused on improving the initialization of the AB matrix, such as PiSSA (Meng et al., 2024), which suggests initializing Aand B by performing SVD on the pre-trained matrix W to accelerate the convergence speed. LoRA-GA (Wang et al., 2024) initializes A and B using the eigenvectors of the full-gradient matrix, aligning the gradient direction of the low-rank product BA with the gradient direction of the pretrained weight matrix W.

2.2 Nyström-like methods

Nyström-like methods approximate matrices by sampling a subset of columns, a technique widely used in kernel matrix approximation (Baker and Taylor, 1979; Williams and Seeger, 2000). Numerous variants have been proposed to enhance the basic Nyström method, including Nyström with k-means clustering (Wang et al., 2019), Nyström with spectral problems (Vladymyrov and Carreira-Perpinan, 2016), randomized Nyström (Li et al., 2010; Persson et al., 2024), ensemble Nyström method (Kumar et al., 2009), fast-Nys (Si et al., 2016).

The Nyström method has also been extended to

general matrix approximation beyond symmetric matrices (Nemtsov et al., 2016). Some methods (Wang and Zhang, 2013; Xiong et al., 2021) explicitly address general matrix approximation by sampling both rows and columns to reconstruct the full matrix. Inspired by such strategies, we propose NLoRA method by to optimize the approximation for efficient matrix reconstruction.

191

192

193

194

195

196

197

200

201

202

203

204

206

207

208

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

3 Method

The Nyström method (Baker and Taylor, 1979), originating from the field of integral equations, is a approach for discretizing integral equations using a quadrature technique. It is commonly employed for out-of-sample extension problems. Specifically, given an eigenfunction problem of the form:

$$\lambda f(x) = \int_{a}^{b} M(x, y) f(y) \, dy, \qquad (2)$$

the Nyström method utilizes a set of s sample points y_1, y_2, \ldots, y_s to approximate f(x) as follows:

$$\lambda \tilde{f}(x) \triangleq \frac{b-a}{s} \sum_{j=1}^{s} M(x, y_j) f(y_j).$$
 (3)

This approach effectively converts the continuous integral equation into a discrete summation, facilitating numerical computation and enabling out-ofsample extensions.

For the pre-trained matrix $W \in \mathbb{R}^{m \times n}$, we assume that it can be decomposed as follows:

$$W = \begin{bmatrix} A_W & B_W \\ F_W & C_W \end{bmatrix}, \tag{4}$$

where, $A_W \in \mathbb{R}^{r \times r}$ is designated to be our sample matrix, $B_W \in \mathbb{R}^{r \times (n-r)}$ and $F_W \in \mathbb{R}^{(m-r) \times r}$ represent the remaining sampled column and row components, respectively, and $C_W \in \mathbb{R}^{(m-r) \times (n-r)}$ corresponds to the remainder of the matrix W. The matrix W can be efficiently approximated using the Nyström method's basic quadrature technique. Starting with the singular value decomposition (SVD) of the sample matrix A_W , represented as $A_W = U\Lambda V^T$, where $U, V \in \mathbb{R}^{r \times r}$ are unitary matrices and $\Lambda \in \mathbb{R}^{r \times r}$ is diagonal. The Nyström approximation reconstructs W based on the outof-sample approximation strategy (Nemtsov et al., 2016). This strategy utilizes the entries of F_W and B_W as interpolation weights for extending the singular vector, resulting in the full approximations of



Figure 3: The diagram of the Nyström-based initialization

the left and right singular vectors of W:

234

235

238

239

240

241

242

244

245

247

248

249

257

258

262

266

$$\hat{U} = \begin{bmatrix} U \\ F_W V \Lambda^{-1} \end{bmatrix}, \quad \hat{V} = \begin{bmatrix} V \\ B_W^T U \Lambda^{-1} \end{bmatrix}, \quad (5)$$

Using the Nyström method, the pretrained matrix W can be approximated as:

$$\widehat{W} = \widehat{U}\Lambda\widehat{V}^{T} = \begin{bmatrix} A_{W} & B_{W} \\ F_{W} & F_{W}A_{W}^{+}B_{W} \end{bmatrix}$$
$$= \begin{bmatrix} A_{W} \\ F_{W} \end{bmatrix} A_{W}^{+} \begin{bmatrix} A_{W} & B_{W} \end{bmatrix}, \qquad (6)$$

where A_W^+ is the Moore-Penrose pseudoinverse of the sampled core matrix A_W . The remaining block C_W is approximated as $F_W A_W^+ B_W$. This approximation demonstrates that W can be effectively reconstructed using only A_W , B_W , and F_W , significantly reducing computational complexity. For the detailed derivation, please refer to Appendix A.

In this way, the matrix W can be approximated as the product of three matrices. Based on this finding, we propose an improvement to LoRA by introducing an intermediate matrix, named as StructuredLoRA (SLoRA). Specifically, we introduce an intermediate matrix $N \in \mathbb{R}^{r \times r}$ between the low-rank matrices A and B, as illustrated in Figure 2. This modification transforms the weight update into:

$$\Delta W = ANB,\tag{7}$$

where $A \in \mathbb{R}^{m \times r}$, $B \in \mathbb{R}^{r \times n}$, $N \in \mathbb{R}^{r \times r}$, and $r \ll \min(m, n)$.

Building on the three-matrix structure, we further enhance SLoRA's effectiveness by employing a Nyström-based initialization. Specifically, by sampling r rows and r columns—corresponding to the rank of LoRA—we efficiently approximate W through matrix decomposition. The resulting submatrices are then directly utilized to initialize the three components of SLoRA, specifically: • The component $\begin{bmatrix} A_W \\ F_W \end{bmatrix}$ is used to initialize the 267 matrix A in SLORA. 268

270

271

272

273

275

276

277

278

279

281

282

284

285

287

288

289

290

291

292

293

294

295

297

298

299

300

302

- The component A_W^+ , representing the Moore-Penrose pseudoinverse of A_W , is used to initialize the matrix N in SLoRA.
- The component $\begin{bmatrix} A_W & B_W \end{bmatrix}$ is used to initialize the matrix B in SLoRA.

While the pseudoinverse can be computed using singular value decomposition (SVD), the process is computationally inefficient on GPUs. To overcome this challenge, we simplify the initialization by directly employing A_W instead of its pseudoinverse, thereby reducing computational overhead while preserving the effectiveness of the initialization. The diagram of the Nyström-based initialization is shown in Figure 3.

By employing this decomposition based on the Nyström approximation method, we propose an initialization strategy for SLoRA, which we term as NyströmLoRA (NLoRA). Additionally, we explore fine-tuning only the intermediate matrix while keeping the other two matrices fixed, which we term IntermediateTune (IntTune).

4 **Experiments**

The experiments were performed on NVIDIA L20 GPUs. For these experiments, we follow the experimental setting given by (Meng et al., 2024), we employ the AdamW optimizer with a batch size of 4, a learning rate of 2E-4, and a cosine annealing schedule with a warmup ratio of 0.03, all while avoiding weight decay. The parameter lora_alpha is consistently set equal to lora_r, with lora_dropout fixed at 0. Adapters are integrated into all linear layers of the base model, and both the base model and adapters utilized Float32 precision for computation. We take the convenience to directly cite

Model	Strategy	Parameters	GSM8K	MATH	HumanEval	MBPP	MT-Bench
	Full FT	6738M	49.05	7.22	21.34	35.59	4.91
	LoRA	320M	42.30	5.50	18.29	35.34	4.58
LLaMA 2-7B	PiSSA	320M	53.07	7.44	21.95	37.09	4.87
	SLoRA	323M	56.48	10.68	23.78	42.32	4.85
	NLoRA	323M	57.70	9.94	25.00	43.12	4.82
	Full FT	7242M	67.02	18.6	45.12	51.38	4.95
	LoRA	168M	67.70	19.68	43.90	58.39	4.90
Mistral-7B	PiSSA	168M	72.86	21.54	46.95	62.66	5.34
	SLoRA	169M	73.01	21.88	47.6	60.3	5.12
	NLoRA	169M	73.92	22.00	44.5	60.3	5.21

Table 1: Experimental results on NLG tasks

Strategy	MNLI	SST-2	MRPC	CoLA	QNLI	QQP	RTE	STS-B		
	DeBERTa-v3-base									
Full FT	89.90	95.63	89.46	69.19	94.03	92.40	83.75	91.60		
LoRA	90.65	94.95	89.95	69.82	93.87	91.99	85.20	91.60		
PiSSA	90.43	95.87	91.67	72.64	94.29	92.26	87.00	91.88		
SLoRA	90.43	96.10	91.91	70.82	93.94	92.11	88.09	91.86		
NLoRA	90.74	96.22	91.91	73.41	94.45	92.03	88.09	92.14		
			RoBE	RTa-larg	ge					
Full FT	90.2	96.4	90.9	68.0	94.7	92.2	86.6	91.5		
LoRA	90.6	96.2	90.9	68.2	94.9	91.6	87.4	92.6		
PiSSA	90.7	96.7	91.9	69.0	95.1	91.6	91.0	92.9		
SLoRA	90.8	96.8	91.7	68.5	94.9	91.6	90.3	92.7		
NLoRA	90.7	96.6	91.9	69.7	95.2	91.6	90.3	92.7		

Table 2: Experimental results on NLU tasks

the baseline performance values from (Meng et al.,2024).

305

307

309

311

312

313

In this section, we evaluate the performance of SLoRA and NLoRA across various benchmark datasets. We compare them with the following baselines:

- Full Fine-tune: which updates all model parameters;
- LoRA (Hu et al., 2021): which approximates weight updates with low-rank matrices while freezing the base model;
- PiSSA (Meng et al., 2024): which initializes adapters using principal singular components and freezes residuals while retaining LoRA's architecture.

We evaluate the capabilities of natural language generation (NLG) using the LLaMA 2-7B (Touvron et al., 2023) and Mistral-7B (Jiang et al., 2023) models through mathematical reasoning, coding proficiency, and dialogue tasks. Additionally, natural language understanding (NLU) tasks were evaluated using the GLUE dataset (Wang, 2018) with DeBERTa-v3-base (He et al., 2021) and RoBERTalarge (Liu, 2019). Finally, we analyze the empirical effects of exclusively fine-tuning the intermediate matrix on both NLU and NLG tasks. 318

319

320

321

322

323

324

325

326

327

328

329

330

331

332

333

334

4.1 Experiments on Natural Language Generation

We conduct experiments using LLaMA 2-7B and Mistral-7B-v0.1. To evaluate mathematical reasoning abilities, we perform fine-tuning using the MetaMathQA dataset and evaluated their perfor-

Strategy	Parameters	GSM8K	MATH	HumanEval	MBPP	MT-Bench
LoRA	320M	42.30	5.50	18.29	35.34	4.58
IntTune	4 M	44.28	6.86	20.70	34.40	4.46

Table 3: IntTune	performance of	n NLG tasks
------------------	----------------	-------------

Strategy	Parameters	MNLI	SST-2	MRPC	CoLA	QNLI	QQP	RTE	STS-B
LoRA	1.33M	90.65	94.95	89.95	69.82	93.87	91.99	85.20	91.60
IntTune	3.07K	81.93	92.20	85.29	65.38	89.13	85.18	76.90	88.37

Table 4: IntTune performance on NLU tasks

mance on GSM8K (Cobbe et al., 2021) and MATH (Yu et al., 2023). In terms of coding capability, we perform fine-tuning on the CodeFeedback dataset (Zheng et al., 2024) and evaluated them using the HumanEval (Chen et al., 2021) and MBPP (Austin et al., 2021) benchmarks. To measure session capabilities, the model is fine-tuned on the WizardLM-Evol-Instruct dataset (Xu et al., 2024) and tested using the MT-Bench dataset (Zheng et al., 2023). All experiments use a subset of 100K data points.

335

337

340

341

342

343

347

353

354

356

357

As shown in Table 1, SLoRA consistently outperforms LoRA, which is labeled with a blue background in Table 1, and even outperforms PiSSA in most tasks. In most cases, NLoRA further enhances the performance of SLoRA. Both methods maintain high parameter efficiency, with only slight increases in trainable parameters (1.15% for LLaMA 2-7B and 0.55% for Mistral-7B compared to LoRA), yet deliver significant performance gains. On these two models, SLoRA achieves average improvements of 38.68%, 15.37%, and 5.19% in mathematical reasoning, coding, and conversational tasks, respectively, relative to LoRA's performance, while NLoRA achieves improvements of 34.53%, 15.83%, and 5.78% over LoRA.

Although the addition of intermediate matrices results in additional matrix multiplication opera-361 tions, the time overhead increases only slightly 362 compared to LoRA. In the MetaMathQA dataset, the training time for SLoRA increases to 27,690.03 364 seconds, which is an increase of 10. 13% compared to LoRA (25142.26 seconds). The training time for NLoRA increases to 25,323.34 seconds, which is almost identical to LoRA's training time. As for initialization time, as shown in Table 5, SLoRA incurs only an 11.95% increase in initialization time compared to LoRA, while NLoRA adds just 12.66 seconds. Both are significantly lower than the time 372

Strategy	Time (seconds)
SLoRA	14.21
NLoRA	25.35
LoRA	12.69
PiSSA	106903.20

Table 5: In	nitialization	time of	different	strategies
-------------	---------------	---------	-----------	------------

cost of PiSSA. Subsequently, we further discuss the effects under different ranks (Section 4.4), learning rates (Appendix C), and optimizers (Appendix D).

373

374

375

378

379

381

382

383

384

386

388

389

390

391

392

393

395

396

397

399

4.2 Experiments on Natural Language Understanding

We also assess the NLU capabilities of RoBERTalarge and DeBERTa-v3-base on the GLUE benchmark. Table 2 summarizes the results of eight tasks performed using these two base models.

SLoRA demonstrates consistent improvements over the baseline LoRA across all tasks, as highlighted in blue. In addition, SLoRA surpasses PiSSA in several cases, showcasing the potential of incorporating an intermediate matrix in LoRA. NLoRA further enhances the performance of SLoRA in most tasks, achieving superior results in tasks such as QNLI, MRPC, and CoLA. For instances where NLoRA does not outperform PiSSA, NLoRA consistently achieves a lower training loss in these scenarios, suggesting its potential for further optimization and efficient fine-tuning. Details can be found in Appendix E.

4.3 NLoRA's Intermediate Matrix Fine-Tuning: A Minimalist Approach

To further improve the computational efficiency of NLoRA, we try to investigate reducing its trainable parameters without sacrificing much perfor-



Figure 4: Compare the performance of different ranks for NLoRA on NLG tasks



Figure 5: Comparison of GPU memory allocation and trainable parameters between IntTune and LoRA

mance. Therefore, we propose **Int**ermediate **Tune** (IntTune), which exclusively fine-tune the intermediate matrix in SFT. To validate the effectiveness of IntTune, we conduct experiments using LLaMA-2-7B and DeBERTa-v3-base for NLG and NLU tasks, respectively. For NLG tasks, we set the learning rate to 2E-3 while keeping other settings unchanged. For NLU tasks, the specific parameter settings are detailed in Appendix E. The results are shown in Table 3 and Table 4.

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

494

425

For NLG tasks, IntTune achieves competitive performance, surpassing LoRA on the GSM8K, MATH, and HumanEval tasks, and attaining comparable results on MBPP and MT-Bench. Overall, the average performance of IntTune across all tasks exceeds that of LoRA, surpassing LoRA's average performance by 7.45%. The comparison of training parameters and memory allocation between Int-Tune and LoRA is shown in Figure 5, with all measurements recorded on the MetaMathQA dataset. In terms of computational efficiency, IntTune significantly reduces the number of trainable parameters to 4M, accounting for only 0.05% of the total model parameters and just 1.13% of LoRA's trainable parameters. Despite this substantial reduction, the training time is shortened to 85.2% of LoRA's. Specifically, LoRA's training time is 25,142.27s, IntTune's training time is reduced to 21,439.26s. Additionally, IntTune enables GPU memory allocation to decrease as well. The percentage of GPU memory allocated drops from 80.9% to 72.5%, with the average memory usage reduced from 36.42 GB to 32.78 GB, a reduction of 9.98%. These results highlight the method's potential for improving performance while optimizing computational resources, making it particularly suitable for SFT LLMs in resource-constrained scenarios. 426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

For NLU tasks, the number of trainable parameters was reduced to 3.07K, representing 0.002% of the total model parameters. Despite this significant reduction, the approach achieved 92.61% of LoRA's average performance across all tasks. Specifically, it attained 96.2% of LoRA's performance on SST-2, 94.5% on QNLI, and 96.2% on STS-B, demonstrating comparable performance across various GLUE tasks, underscoring its robustness and effectiveness in diverse scenarios.

These results demonstrate the potential of the Nyström initialization, as fine-tuning only the intermediate matrix can still yield competitive performance.

4.4 Experiments on Various Ranks

In this section, we examine the impact of progressively increasing the rank of NLoRA and SLoRA from 1 to 128 to assess their ability to consistently outperform the baseline across different ranks. Training is performed on the MetaMathQA dataset for a single epoch, with validation conducted on the GSM8K and MATH datasets.

The experimental results are presented in Figure 4. On the GSM8K dataset, NLoRA performs relatively better at higher ranks, surpassing LoRA

Strategy	Parameters	GSM8K	MATH	HumanEval	MBPP	MT-Bench
LoRA	320M	42.30	5.50	18.29	35.34	4.58
IntTune(Rank=256)	15M	49.51	6.62	21.30	33.90	3.59
IntTune(Rank=128)	4M	44.28	6.86	20.70	34.40	4.46
IntTune(Rank=64)	0.9M	37.98	5.56	14.60	34.70	4.55



Table 6: Compare the performance of different ranks for IntTune on NLG tasks

Figure 6: Compare the performance of different ranks for IntTune on NLU tasks

by 43.08% and 36.41% at ranks 64 and 128, respectively. SLoRA, on the other hand, exhibits relatively stronger performance at lower ranks, outperforming LoRA by 107.45%, 77.31%, 53.54%, and 76.13% at ranks 1, 2, 4, and 8, respectively. On the MATH dataset, SLoRA shows a slight overall advantage, while NLoRA continues to deliver strong performance, particularly at higher ranks.

462

463

464

465

466

467

468

469

For IntTune, we compared ranks of 64, 128, and 470 256 in the NLG tasks, following the same experi-471 mental setup as shown in Section 4.1. In the NLU 472 experiments, we evaluated ranks of 4, 8, and 16. 473 The results of these experiments are presented in 474 Table 6 and Figure 6. On NLG tasks, IntTune does 475 not exhibit a strictly increasing performance trend 476 with higher ranks. Instead, different ranks excel in 477 different tasks. Specifically, rank 128 and rank 256 478 achieve 7.45% and 5.62% higher performance than 479 LoRA on average, both outperforming LoRA over-480 all. Meanwhile, rank 64, though slightly lower, still 481 reaches 93.66% of LoRA's performance, demon-482 strating the feasibility of fine-tuning with even 483 fewer parameters while maintaining competitive re-484 sults. On NLU tasks, the model performance grad-485 486 ually improves with increasing rank. For ranks 4, 8, and 16, the average performance reaches 86.20%, 487 92.61%, and 95.80% of LoRA's performance, re-488 spectively, while the number of parameters is only 489 1.35K, 3.07K, and 9.99K, respectively. 490

5 Conclusion

This work advances parameter-efficient fine-tuning strategies for large language models by introducing SLoRA and NLoRA, along with an exploration of an intermediate matrix fine-tuning method, IntTune. SLoRA incorporates a small intermediate matrix, enhancing expressiveness with minimal parameter overhead, while NLoRA leverages Nyström-based initialization to bypass the computational complexity of SVD, achieving competitive downstream performance. IntTune, by fine-tuning only the intermediate matrix in NLoRA, even boosts average NLG performance over LoRA while maintaining high parameter efficiency. Extensive experiments on NLG and NLU tasks demonstrate the robustness and adaptability of our methods, providing practical solutions for optimizing large models under resource constraints.

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

508

509

510

511

512

513

514

515

516

517

518

519

6 Limitaion

While our method demonstrates strong performance in both NLG and NLU tasks, its applicability to ultra-low parameter fine-tuning approaches, such as IntTune, warrants further exploration. Additionally, extending our approach to visual tasks could provide valuable insights into its generalization and versatility across modalities. Furthermore, integrating SLoRA with advanced LoRA variants presents a compelling direction for future research to further enhance fine-tuning efficacy.

References

520

522

523

524

525

526

528

530

532

533

534

535

537

539

540 541

542

543

544

548

552

554

558

569

570

571

572

573

- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. 2021. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*.
- Christopher TH Baker and RL Taylor. 1979. The numerical treatment of integral equations. *Journal of Applied Mechanics*, 46(4):969.
- Yupeng Chang, Yi Chang, and Yuan Wu. 2024. Balora: Bias-alleviating low-rank adaptation to mitigate catastrophic inheritance in large language models. *arXiv preprint arXiv:2408.04556*.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: efficient finetuning of quantized llms (2023). *arXiv preprint arXiv:2305.14314*, 52:3982–3992.
- Ning Ding, Yujia Qin, Guang Yang, Fuchao Wei, Zonghan Yang, Yusheng Su, Shengding Hu, Yulin Chen, Chi-Min Chan, Weize Chen, et al. 2023.
 Parameter-efficient fine-tuning of large-scale pretrained language models. *Nature Machine Intelligence*, 5(3):220–235.
- Soufiane Hayou, Nikhil Ghosh, and Bin Yu. 2024. Lora+: Efficient low rank adaptation of large models. *arXiv preprint arXiv:2402.12354.*
- Pengcheng He, Jianfeng Gao, and Weizhu Chen. 2021. Debertav3: Improving deberta using electra-style pretraining with gradient-disentangled embedding sharing. arXiv preprint arXiv:2111.09543.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. *arXiv preprint arXiv:2310.06825*.
- Sanjiv Kumar, Mehryar Mohri, and Ameet Talwalkar. 2009. Ensemble nystrom method. *Advances in Neural Information Processing Systems*, 22.

Mu Li, James Tin-Yau Kwok, and Baoliang Lü. 2010. Making large-scale nyström approximation possible. In *Proceedings of the 27th International Conference on Machine Learning, ICML 2010*, page 631. 574

575

576

578

579

580

581

582

583

584

585

587

589

591

592

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

- Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, et al. 2023. Starcoder: may the source be with you! *arXiv preprint arXiv:2305.06161*.
- Vladislav Lialin, Sherin Muckatira, Namrata Shivagunde, and Anna Rumshisky. 2023. Relora: Highrank training through low-rank updates. In *The Twelfth International Conference on Learning Representations*.
- Shih-Yang Liu, Chien-Yi Wang, Hongxu Yin, Pavlo Molchanov, Yu-Chiang Frank Wang, Kwang-Ting Cheng, and Min-Hung Chen. 2024. Dora: Weightdecomposed low-rank adaptation. *arXiv preprint arXiv:2402.09353*.
- Yinhan Liu. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 364.
- Fanxu Meng, Zhaohui Wang, and Muhan Zhang. 2024. Pissa: Principal singular values and singular vectors adaptation of large language models. *arXiv preprint arXiv:2404.02948*.
- Arik Nemtsov, Amir Averbuch, and Alon Schclar. 2016. Matrix compression using the nyström method. *Intelligent Data Analysis*, 20(5):997–1019.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744.
- David Persson, Nicolas Boullé, and Daniel Kressner. 2024. Randomized nystr\" om approximation of non-negative self-adjoint operators. *arXiv preprint arXiv:2404.00960*.
- Si Si, Cho-Jui Hsieh, and Inderjit Dhillon. 2016. Computationally efficient nyström approximation using fast transforms. In *International conference on machine learning*, pages 2655–2663. PMLR.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Max Vladymyrov and Miguel Carreira-Perpinan. 2016. The variational nystrom method for large-scale spectral problems. In *International Conference on Machine Learning*, pages 211–220. PMLR.

684

691 692 693

690

- 694 695 696
- 698

697

- 699 700
- 701

702

706

707 708 709

- 710 711
- 713 714
- 715

- 716

717

718

719

720

721

722

723

analysis platform for natural language understanding. arXiv preprint arXiv:1804.07461. Shaowen Wang, Linxi Yu, and Jian Li. 2024. Lora-ga:

627

628

630

631 632

634

635

636

637

642

647

649

653

658

664

667

668

673

674

675

679

Low-rank adaptation with gradient approximation. arXiv preprint arXiv:2407.05000.

Alex Wang. 2018. Glue: A multi-task benchmark and

- Shusen Wang, Alex Gittens, and Michael W Mahoney. 2019. Scalable kernel k-means clustering with nystrom approximation: Relative-error bounds. Journal of Machine Learning Research, 20(12):1–49.
- Shusen Wang and Zhihua Zhang. 2013. Improving cur matrix decomposition and the nyström approximation via adaptive sampling. The Journal of Machine Learning Research, 14(1):2729–2769.
- Christopher Williams and Matthias Seeger. 2000. Using the nyström method to speed up kernel machines. Advances in neural information processing systems, 13.
- Tingyu Xia, Bowen Yu, Kai Dang, An Yang, Yuan Wu, Yuan Tian, Yi Chang, and Junyang Lin. 2024. Rethinking data selection at scale: Random selection is almost all you need. arXiv preprint arXiv:2410.09335.
- Yunyang Xiong, Zhanpeng Zeng, Rudrasis Chakraborty, Mingxing Tan, Glenn Fung, Yin Li, and Vikas Singh. 2021. Nyströmformer: A nyström-based algorithm for approximating self-attention. In *Proceedings of* the AAAI Conference on Artificial Intelligence, volume 35, pages 14138-14148.
- Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, Qingwei Lin, and Daxin Jiang. 2024. Wizardlm: Empowering large pre-trained language models to follow complex instructions. In The Twelfth International Conference on Learning Representations.
- Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. 2023. Metamath: Bootstrap your own mathematical questions for large language models. arXiv preprint arXiv:2309.12284.
- Qingru Zhang, Minshuo Chen, Alexander Bukharin, Nikos Karampatziakis, Pengcheng He, Yu Cheng, Weizhu Chen, and Tuo Zhao. 2023. Adalora: Adaptive budget allocation for parameter-efficient finetuning. arXiv preprint arXiv:2303.10512.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. Advances in Neural Information Processing Systems, 36:46595-46623.
- Tianyu Zheng, Ge Zhang, Tianhao Shen, Xueling Liu, Bill Yuchen Lin, Jie Fu, Wenhu Chen, and Xiang Yue. 2024. Opencodeinterpreter: Integrating code generation with execution and refinement. arXiv preprint arXiv:2402.14658.

Α **Detailed Derivation for Nyström** Approximation

This section provides a detailed derivation of the Nyström approximation presented in Section 3, following the approach proposed in (Nemtsov et al., 2016). Specifically, the quadrature technique is applied to the sample matrix of W, followed by an out-of-sample extension to approximate W.

The basic quadrature technique of the Nyström method is used to approximate the Singular Value Decomposition (SVD) of a matrix. In this context, no eigen-decomposition is required. Specifically, denote the matrix $W \in \mathbb{R}^{m \times n}$ can be decomposed as:

$$W = \begin{bmatrix} A_W & B_W \\ F_W & C_W \end{bmatrix}.$$
 (8)

where, $A_W \in \mathbb{R}^{r \times r}$ is designated to be the sample matrix, $B_W \in \mathbb{R}^{r \times (n-r)}$ and $F_W \in \mathbb{R}^{(m-r) \times r}$ represent the remaining sampled column and row components, respectively, and $C_W \in \mathbb{R}^{(m-r) \times (n-r)}$ corresponds to the remainder of the matrix W.

The derivation begins with the SVD of A_W , expressed as:

$$A_W = U\Lambda V^T, \tag{9}$$

where $U, V \in \mathbb{R}^{r \times r}$ are unitary matrices, and $\Lambda \in$ $\mathbb{R}^{r \times r}$ is a diagonal matrix. Assuming that zero is not a singular value of A_W , the decomposition can be further approximated. Accordingly, the matrix U is formulated as:

$$U = A_W V \Lambda^{-1}.$$
 (10)

Let $u^i, h^i \in \mathbb{R}^r$ represent the *i*-th columns of U and V, respectively. Denote $u^i = \{u_l^i\}_{l=1}^r$ as the individual elements of the i-th column of U. Using Eq. (10), each element u_l^i is expressed as the sum:

$$u_l^i = \frac{1}{\lambda_i} \sum_{j=1}^n W_{lj} \cdot h_j^i. \tag{11}$$

The elements of F_W can be used as interpolation weights to extend the singular vector u^i to the k^{th} row of W, where $s + 1 \leq k \leq n$. Let $\tilde{u}^i = {\tilde{u}^i_{k-s}}_{k=s+1}^n \in \mathbb{R}^{n-s \times 1}$ denote a column vector comprising all the approximated entries. Each element \tilde{u}_k^i is computed as:

$$\tilde{u}_{k}^{i} = \frac{1}{\lambda_{i}} \sum_{j=1}^{n} W_{kj} \cdot h_{j}^{i}.$$
 (12) 724

Thus, the matrix form of \tilde{u}^i is given by $\tilde{u}^i = \frac{1}{\lambda_i} F_W \cdot h^i$. By arranging all the \tilde{u}^i 's into a matrix $\tilde{U} = \begin{bmatrix} \tilde{u}^1 & \tilde{u}^2 & \dots & \tilde{u}^r \end{bmatrix} \in \mathbb{R}^{n-s \times r}$, the following expression is obtained:

725

726 727

730

731

733

735

737

739

740

741

742

743

744

745

746

747

748

$$\tilde{U} = F_W H \Lambda^{-1}.$$
 (13)

The Eq. (9) can also be written as $V = A_W^T U \Lambda^{-1}$. To approximate the right singular vectors of the out-of-sample columns, a symmetric argument is applied, yielding:

$$\tilde{H} = B_W^T U \Lambda^{-1}.$$
 (14)

In that case, the full approximations of the left and right singular vectors of \widehat{W} , represented by \widetilde{U} and \widetilde{H} , respectively, are then obtained as follows:

$$\widehat{U} = \begin{bmatrix} U \\ F_W V \Lambda^{-1} \end{bmatrix}, \quad \widehat{V} = \begin{bmatrix} V \\ B_W^T U \Lambda^{-1} \end{bmatrix}.$$
(15)

The explicit Nyström form of M is given by:

where A_W^+ denotes the pseudo-inverse of W. In this approximation, \widehat{W} does not modify A_W, B_W and F_W but approximates C_W by $F_W A_W^+ B_W$. This approach achieves a matrix approximation using only the selected rows and columns, effectively capturing the essential structure with reduced computational complexity.

B Experiments on Various Initializations

For SLoRA, we kept the initialisation of the A and B matrices the same as for LoRA, and in turn 753 explored the effect of different methods of initiali-754 sation of the intermediate matrices on the results. 755 Specifically, we experimented with Kaiming initial-756 ization and Gaussian initialization on all the NLG tasks of LLaMA 2-7B, with the same experimental 758 setup as in Section 4. The performance of the models under these settings is shown in Table 7. The results indicate that Kaiming initialization consis-762 tently achieves better performance across all tasks. Gaussian initialization also achieves competitive results, which demonstrates the robustness of our method. In our experiments, we use kaiming to initialize SLoRA. 766

Tasks	Kaiming	Gaussian
GSM8K	56.48	56.10
MATH	10.68	9.56
HumanEval	23.78	23.2
MBPP	42.32	40.5
MT-Bench	4.85	3.93

Table 7: Different Initialization on SLoRA

Strategy	LR	GSM8K	MATH
	2E-4	56.48	10.68
	5E-4	59.51	11.04
SLOKA	2E-5	51.02	6.94
	5E-5	52.84	8.36
	2E-4	57.70	9.94
	5E-4	54.81	10.60
NLOKA	2E-5	45.11	6.42
	5E-5	52.39	7.58

Table 8: Comparasion of different learning rate onSLoRA and NLoRA

767

768

769

770

772

773

774

775

776

777

778

779

780

781

782

783

784

785

786

C Experiments on Various Learning Rates

We evaluated the impact of four learning rates: 2E-4, 2E-5, 5E-4 and 5E-5 on the model's performance. The experimental setup remains the same as described earlier. The results of these experiments are presented in Table 8. Among the evaluated learning rates, 5E-4 achieved the best overall performance. However, we opted for 2E-4 in our experiments, as its performance, while slightly lower than that of 5E-4, remained comparable and still exceeded the original baseline. Moreover, at the learning rate of 2E-4, NLoRA exhibited lower loss and better convergence behavior, making it a more appropriate choice for our experimental setup.

For the case of fine-tuning only the intermediate matrix, we tested the performance under different learning rates. The results indicate that a learning rate of 2E-3 achieved the best performance. The result is shown in Figure 9.

LR	GSM8K	MATH
2E-4	43.29	5.74
5E-4	44.20	5.70
2E-3	44.28	6.86
5E-3	40.86	6.08

 Table 9: Comparasion of Different Learning Rates on

 IntTune

Parameters	GSM8K	MATH	HumanEval	MBPP	MT-Bench
320M	42.30	5.50	18.29	35.34	4.58
323M	57.70	9.94	25.00	43.12	4.82
323M	58.10	10.82	25.60	43.40	4.99
	Parameters 320M 323M 323M	ParametersGSM8K320M42.30323M57.70323M58.10	ParametersGSM8KMATH320M42.305.50323M57.709.94323M58.1010.82	ParametersGSM8KMATHHumanEval320M42.305.5018.29323M57.709.9425.00323M58.1010.8225.60	ParametersGSM8KMATHHumanEvalMBPP320M42.305.5018.2935.34323M57.709.9425.0043.12323M58.1010.8225.6043.40

Table 10: Comparision of Adamw and RMSProp on NLG

Strategy	MNLI	SST-2	MRPC	CoLA	QNLI	QQP	RTE	STS-B
LoRA	90.65	94.95	89.95	69.82	93.87	91.99	85.20	91.60
NLoRA	90.74	96.22	91.91	73.41	94.45	92.03	88.09	92.14
NLoRA+RMSProp	90.41	96.22	91.91	68.61	94.18	92.03	88.09	91.86

Table 11: Comparision of Adamw and RMSProp on NLU

D Experiments on Various Optimizers

787

788

790

793

796

809

810

811

We experimented with different optimizers on both NLG and NLU tasks. In addition to the default AdamW optimizer, we also evaluated the RMSProp optimizer. Other experimental setups are the same as Section 4. The experimental results are shown in Table 10 and Table 11.

On NLG tasks, we observed that the RMSProp optimizer further improved the model's performance. However, its performance on NLU tasks was relatively mediocre. This discrepancy might stem from the underlying differences in the nature of NLG and NLU tasks. NLG tasks typically involve generating coherent sequences of text, which require more stable gradient updates over longer contexts. RMSProp's adaptive learning rate mechanism, which emphasizes recent gradients, may help maintain stability and enhance performance in such scenarios. In contrast, NLU tasks often involve classification or regression over shorter input sequences, where AdamW's weight decay and bias correction might be more effective in avoiding overfitting and ensuring generalization, thus outperforming RMSProp in these tasks.

E Experimental Settings on NLU

We evaluate the performance on the GLUE bench-812 mark, which includes two single-sentence tasks 813 (CoLA and SST-2), three natural language infer-814 ence tasks (MNLI, QNLI, and RTE), and three 815 816 similarity and paraphrase tasks (MRPC, QQP, and STS-B). For evaluation metrics, we report over-817 all accuracy (matched and mismatched) for MNLI, 818 Matthew's correlation for CoLA, Pearson's correlation for STS-B, and accuracy for the remaining 820

datasets.

In DeBERTa-v3-base, SLoRA and NLoRA were applied to the W_Q , W_K , and W_V matrices, while in RoBERTa-large, they were applied to the W_Q and W_V matrices. The experiments for natural language understanding (NLU) were conducted using the publicly available LoRA codebase. For MRPC, RTE, and STS-B tasks, we initialized RoBERTalarge with a pretrained MNLI checkpoint. The rank of SLoRA and NLoRA in these experiments was set to 8. Optimization was performed using AdamW with a cosine learning rate schedule. Table 12 and Table 13 outline the hyperparameters used for the GLUE benchmark experiments.

For IntTune, we set both the LoRA rank and LoRA alpha to 8. The remaining parameter configurations are provided in Table 14.

837

821

Dataset		ERTa-v3	-base	RoBERTa-large				
	LR	BS	Epoch	LoRA alpha	LR	BS	Epoch	LoRA alpha
CoLA	3E-04	16	40	16	4E-04	8	20	8
SST-2	5E-04	16	10	8	5E-04	16	10	8
MRPC	5E-04	32	100	16	2E-04	32	50	16
MNLI	3E-04	32	10	16	3E-04	32	10	16
QNLI	2E-04	32	20	16	6E-04	16	10	8
QQP	6E-04	32	20	8	6E-04	16	10	16
RTE	3E-04	32	40	16	5E-04	32	30	16
STS-B	5E-04	16	10	16	3E-04	16	30	16

Table 12: Hyperparameters of NLoRA on GLUE

Dataset		ERTa-v3	-base	RoBERTa-large				
	LR	BS	Epoch	LoRA alpha	LR	BS	Epoch	LoRA alpha
CoLA	3E-04	16	40	16	4E-04	8	20	8
SST-2	5E-04	16	10	8	5E-04	16	10	8
MRPC	5E-04	32	100	16	2E-04	32	50	16
MNLI	3E-04	32	10	16	3E-04	32	20	16
QNLI	2E-04	32	20	16	6E-04	16	10	8
QQP	6E-04	32	20	8	6E-04	16	10	16
RTE	3E-04	32	40	16	5E-04	32	30	16
STS-B	5E-04	16	10	16	3E-04	16	30	16

Table 13: Hyperparameters of SLoRA on GLUE

Dataset	LR	BS	Epoch
CoLA	7E-03	16	40
SST-2	6E-03	32	30
MRPC	4E-03	16	50
MNLI	6E-03	64	20
QNLI	8E-03	64	20
QQP	6E-03	32	20
RTE	6E-03	16	25
STS-B	6E-03	16	60

Table 14: Hyperparameters for IntTune