## LSR-Adapt: Ultra-Efficient Parameter Tuning with Matrix Low Separation Rank Kernel Adaptation

**Anonymous ACL submission** 

#### Abstract

001 Imposing an effective structural assumption on neural network weight matrices has been 003 the major paradigm for designing Parameter-Efficient Fine-Tuning (PEFT) systems for adapting modern large pre-trained models to various downstream tasks. However, low rank based adaptation has become increasingly chal-007 800 lenging due to the sheer scale of modern large language models. In this paper, we propose an effective kernelization to further reduce the number of parameters required for adaptation 012 tasks. Specifically, from the classical idea in numerical analysis regarding matrix Low-Separation-Rank (LSR) representations, we de-014 velop a kernel using this representation for the low rank adapter matrices of the linear layers from large networks, named the Low Sepa-017 ration Rank Adaptation (LSR-Adapt) kernel. With the ultra-efficient kernel representation of the low rank adapter matrices, we manage to achieve state-of-the-art performance with even higher accuracy with almost half the number of parameters as compared to conventional low rank based methods. This structural assumption also opens the door to further GPU-side optimizations due to the highly parallelizable 027 nature of Kronecker computations.

### 1 Introduction

037

041

Effectively designing structural assumptions is the key to the parameter-efficient approximation of network weight matrices. Low-Rank Adaptation (LoRA) (Hu et al., 2021) has been the pioneering method in PEFT that assumes a low-rank structure of the network weight matrices. However, despite its earlier successes, with an increasing number of parameters of modern day large language models, such simple structural assumption simply cannot effectively reduce the number of parameters into a manageable size with decent accuracy. To address this issue, various other structural assumptions have been proposed over the years (Dettmers et al., 2023; Liu et al., 2024; Edalati et al., 2022; He et al., 2023; Xu et al., 2023). Most of these methods, however, despite the effectiveness, lack solid theoretical reasoning of their structure design choices, and thus do not offer fine-grained control on the model performance. In this work, we provide a PEFT kernel based on the separable representation of matrices derived from the ideas in high dimensional numerical analysis to further decompose the factor matrices in various PEFT methods, coined as the Low Separation Rank Adaptation (LSR-Adapt) kernel, which not only yields higher fine-tuning accuracy with even less trainable parameters, but also provides a solid theoretical foundation of this structural assumption for more control over the fine-tuning process.

042

043

044

047

048

053

054

056

058

059

060

061

062

063

064

065

066

067

068

069

071

072

073

074

076

077

In summary, the major contributions of this paper are as follows.

- 1. Developing a structural assumption for PEFT based on a separable representation of matrices, which can be used as a kernel to further decompose the factor matrices of various PEFT methods, such as LoRA family methods (Hu et al., 2021; Dettmers et al., 2023).
- 2. Providing a theoretical analysis of the structure choice to give more insight for fine-tuning performance control.
- 3. Experimental evaluations of our method as compared to other state-of-the-art PEFT methods against GLUE and SuperGLUE benchmarks (Wang, 2018; Wang et al., 2019).
- 4. Discussions on how this kernel structured computation can be parallelized using GPU, which can be interesting for further research in high-performance computing (Jangda and Yadav, 2024).

#### 2 Related Works

078

079

084

091

099

100

101

102

103

104

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

Numerous attempts have been done regarding Parameter-Efficient Fine-Tuning to adapt modern large language models to various applications. LoRA (Hu et al., 2021) has been one of the first major attempts in imposing efficient structural assumption on the neural network weight matrices of large models, subsequent research based on LoRA involves utilizing lower-precision quantization to harness the advantages of efficient calculations on lower-precision numbers offered by contemporary tensor core-based GPUs (Dettmers et al., 2023), and other form of weight decompositions with better semantic understanding of the weight matrices (Liu et al., 2024). Further more, Kronecker product based factorizations of the weight matrices have also been studied to further reduce the parameter counts (Edalati et al., 2022; He et al., 2023), and He et al. provides a mixture of low-rank and Kronecker factorization to achieve parameter-efficient tuning for vision models.

#### **3** Preliminaries

To develop an efficient kernel to supercharge parameter-efficient tuning using separated representation of factor matrices, we first recall the generic definition of the separated representation in high-dimensional numerical analysis (Beylkin and Mohlenkamp, 2005),

**Definition 3.1** (The Separated Representation). Given an equation in r dimensions (r independent variables), we can try to approximate its solution f by the following separation of variables,

$$f(x_1, \cdots, x_r)$$
  
=  $\sum_{k=1}^{s} \lambda_k \cdot g_k^{(1)}(x_1) \cdots g_k^{(r)}(x_r) + \mathcal{O}(\epsilon)$  (1)

which is called a *separated representation*, where  $\mathcal{O}(\epsilon)$  is the desired asymptotic error proportional to  $\epsilon$ ,  $\{g_k^{(i)}(x_i)\}$  is the factor function for the *r*-dimensional variable  $\boldsymbol{x} = \{x_1, \dots, x_r\}$  at each dimension  $i \in \{1, \dots, r\}$  and  $\lambda_k$  is a scaling factor for the *k*-th summation term where  $k \in \{1, 2, \dots, s\}$ , and *s* is called the *separation rank*.

This formulation effectively allows one to approximate a high-dimensional function f with a linear complexity of  $\mathcal{O}(r)$ . Using this idea, we define the separated representation of matrices, by thinking an matrix  $M \in \mathbb{R}^{m_1 \times m_2}$  of dimension / rank

 $r \leq m_2$  (may not be full rank) as a discrete representation of an *r*-dimensional linear operator  $\mathcal{M}$  on a rectangular domain of indices  $(i, j) \in \mathbb{R}^{m_1 \times m_2}$ , *i.e.*, the matrix entries  $M_{i,j} = \mathcal{M}(i, j)$ , we can effectively extend the separated representation for *r*-dimensional functions to matrices. 124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

47

148

149

150

151

152

153

154

155

156

157

158

159

160

161

**Definition 3.2** (The Matrix Separated Representation). For a given approximation error  $\epsilon$ , we can represent the matrix  $M \in \mathbb{R}^{m_1 \times m_2}$  as,

$$M = \sum_{k=1}^{s} \lambda_k M_k^{(1)} \otimes \cdots \otimes M_k^{(r)} + \mathcal{O}(\epsilon)$$
 (2)

with scalars  $\lambda_1 \geq \cdots \geq \lambda_r > 0$ , the integer *s* the matrix *separation rank*, and the factor matrix  $M_k^{(i)}$  is of dimension  $m_{k,1}^{(i)} \times m_{k,2}^{(i)}$  and  $\prod_{i=1}^r m_{k,1}^{(i)} = m_1$ ,  $\prod_{i=1}^r m_{k,2}^{(i)} = m_2$  for all  $k = 1, 2, \cdots, s$ . In practice, we would like this separation rank term to be low for a parameter-efficient approximation, which leads to the matrix *Low Separation Rank* (*LSR*) structure.

The operator " $\otimes$ " in the definition above denotes the Kronecker product. Specifically, for two matrices  $U \in \mathbb{R}^{u_1 \times u_2}$ ,  $V \in \mathbb{R}^{v_1 \times v_2}$ , their Kronecker product, denoted as  $U \otimes V \in \mathbb{R}^{(u_1v_1) \times (u_2v_2)}$ , takes the format,

$$\boldsymbol{U} \otimes \boldsymbol{V} = \begin{bmatrix} U_{1,1}\boldsymbol{U} & U_{1,2}\boldsymbol{V} & \cdots & U_{1,u_2}\boldsymbol{V} \\ U_{2,1}\boldsymbol{U} & U_{2,2}\boldsymbol{V} & \cdots & U_{2,u_2}\boldsymbol{V} \\ \vdots & \vdots & \ddots & \vdots \\ U_{u_1,1}\boldsymbol{U} & U_{u_1,2}\boldsymbol{V} & \cdots & U_{u_1,u_2}\boldsymbol{V} \end{bmatrix}$$

To gain a fine-grained control for the accuracy of the approximation, Beylkin and Mohlenkamp proposed to use a condition number for the separated representation.

**Definition 3.3** (Condition Number of A Separated Representation). The condition number of (3) is the ratio

$$\gamma = \frac{\left(\sum_{k=1}^{s} \lambda_k^2\right)^{1/2}}{\|\boldsymbol{M}\|_F}$$
(3)

where  $\|\cdot\|_F$  denotes the Frobenius norm.

In a numerical computing system, we do not want  $\gamma$  to be too large, a good rule of thumb to set the condition number would be to make it satisfy (Beylkin and Mohlenkamp, 2005),

$$\gamma \mu \|\boldsymbol{M}\|_F \le \epsilon \tag{4}$$

where  $\mu$  is the machine round-off, for instance, in a 16-bit precision machine, the round-off number 163 164 is  $\mu = 2^{-11} \approx 4.88 \times 10^{-4}$ . With this condition 165 number, we can gain a fine-grained control over 166 desired accuracy and dimensions of the factor ma-167 trices constrained by the numerical precision used 168 during the finetuning process.

#### 4 Our Approach

169

170

171

172

173

175

176

177

178

179

180

181

185

186

187

189

190

191

192

193

196

197

198

201

To develop an even more parameter efficient tuning mechanism, we are looking at a more parameterefficient representation of the weight update matrix  $\Delta W \in \mathbb{R}^{w_1 \times w_2}$  for the weight matrix  $W \in \mathbb{R}^{w_1 \times w_2}$  of the target network layer  $\ell$ ,

$$\boldsymbol{W} = \boldsymbol{W} + \Delta \boldsymbol{W},\tag{5}$$

while a naive approach to adopt the matrix separable representation is to simply do,

$$\Delta \boldsymbol{W} \approx \sum_{k=1}^{s} \lambda_k \boldsymbol{W}_k^{(1)} \otimes \dots \otimes \boldsymbol{W}_k^{(r)} \qquad (6)$$

where  $r = \operatorname{rank}(\Delta W)$ , while we can follow the common hypothesis in LoRA (Hu et al., 2021) where the weight update matrix is approximately low rank and use a rather small r, chaining even r > 3 Kronecker products can still be computationally expensive. Hence, we choose to take the low rank adapter matrices from LoRA,

$$\Delta \boldsymbol{W} \approx \boldsymbol{A}\boldsymbol{B} \tag{7}$$

for

$$\boldsymbol{W}' = \boldsymbol{W} + \alpha \Delta \boldsymbol{W}, \tag{8}$$

where  $A \in \mathbb{R}^{w_1 \times r}$  and  $B \in \mathbb{R}^{r \times w_2}$  are the factor matrices,  $\alpha$  is a scalar controlling the impact of the weight update matrix during adaptation, and structure the factor matrices A, B with their matrix low-separation rank representations, *i.e.*, the LSR-Adapt Kernel. Note that since r is already a small value from the LoRA assumption, we can simply use two Kronecker factor matrices at each summation term of the LSR-Adapt kernel to achieve a decent reduction of parameter counts as compared to original LoRA,

$$\boldsymbol{A} \approx \sum_{k=1}^{s_A} \lambda_k^A \boldsymbol{A}_k^{(1)} \otimes \boldsymbol{A}_k^{(2)}$$
(9)

$$\boldsymbol{B} \approx \sum_{k=1}^{s_B} \lambda_k^B \boldsymbol{B}_k^{(1)} \otimes \boldsymbol{B}_k^{(2)}, \qquad (10)$$

where  $s_A$  and  $s_B$  are respective separation ranks for factor matrices A and B, which for simplified evaluation we set  $s_A = s_B = s$ ,  $\lambda_k^A$  and  $\lambda_k^B$  are the corresponding scalar factors at summation 205 term  $k = 1, 2, \dots, s$ , which we will drop in the 206 actual implementation and merge them into the  $\alpha$  207 factor of the final low-rank factorization, the small 208 Kronecker factor matrices take the shape  $A_k^{(i)} \in$  209  $\mathbb{R}^{a_{k,1}^{(i)} \times a_{k,2}^{(i)}}$  for  $i = \{1, 2\}$  and  $B_k^{(j)} \in \mathbb{R}^{b_{k,1}^{(j)} \times b_{k,2}^{(j)}}$  210 for  $j = \{1, 2\}$ , where, 211

$$a_{k,1}^{(1)} \times a_{k,1}^{(2)} = w_1, \qquad a_{k,2}^{(1)} \times a_{k,2}^{(2)} = r$$

$$b_{k,1}^{(1)} \times b_{k,1}^{(2)} = r, \qquad b_{k,2}^{(1)} \times b_{k,2}^{(2)} = w_2, \quad (11)$$
213

$$b_{k,1}^{(1)} \times b_{k,1}^{(2)} = r, \qquad b_{k,2}^{(1)} \times b_{k,2}^{(2)} = w_2, \quad (11)$$
 213

214

215

216

223

224

225

226

227

228

229

230

231

232

234

235

236

237

238

239

240

241

242

243

in practice we simply set  $a_{k,2}^{(1)} = b_{k,1}^{(1)} = r^{(1)}$  and  $a_{k,2}^{(2)} = b_{k,1}^{(2)} = r^{(2)}$  such that  $r^{(1)} \times r^{(2)} = r$ . Thus the weight update matrix takes the format,

for

$$W' = W + \alpha \Delta W$$
 220

$$\approx \boldsymbol{W} + \alpha \left(\sum_{k=1}^{s} \boldsymbol{A}_{k}^{(1)} \otimes \boldsymbol{A}_{k}^{(2)}\right) \times$$
 22

$$\left(\sum_{k=1}^{s} \boldsymbol{B}_{k}^{(1)} \otimes \boldsymbol{B}_{k}^{(2)}\right). \quad (13)$$

A simple diagram of this adaptation mechanism is shown in Figure 1. Note that this is much more parameter-efficient as compared to the original lowrank factorization. Take a 768 × 768 network weight matrix for instance, if we set r = 8, we are looking at  $2 \times 768 \times 8 = 12,288$  parameters, with an even higher rank r = 16, and assume balanced dimensions for the small kernel weight matrices, say,  $A_k^{(1)} \in \mathbb{R}^{32 \times 4}, A_k^{(2)} \in \mathbb{R}^{24 \times 4}$ ,  $B_k^{(1)} \in \mathbb{R}^{32 \times 4}, B_k^{(2)} \in \mathbb{R}^{24 \times 4}$ , and separation rank s = 16, we can achieve a much lower parameter count  $2 \times (32 \times 4 + 24 \times 4) \times 16 = 5,632$  while still maintaining a higher accuracy as we have found in fine-tuning experiments.

One can also show that this Kronecker-product based structure is amenable to parallel computation on modern power GPUs (Golub and Van Loan, 2013; Jangda and Yadav, 2024). This enables potential development of custom CUDA kernels to further improve the training runtime, which is left to the future work.



Figure 1: Overview of the working mechanism of LSR-Adapt kernel.

Method	GLUE			SuperGLUE					Avorago
	MRPC	SST-2	CoLA	RTE	CB	COPA	WSC	BoolQ	Average
LoRA	74.51	94.43	83.32	68.23	76.79	57.39	63.46	75.41	74.19
KronA	76.57	94.12	82.59	67.92	79.23	56.48	63.46	74.97	74.42
KAdaptation	77.68	93.81	83.16	68.73	78.63	57.94	63.46	75.22	74.69
LSR-Adapt	80.88	94.27	83.41	68.95	82.14	60.32	63.46	75.72	76.14

Table 1: Performance comparison of different adaptation methods on GLUE and SuperGLUE benchmark tasks. The best results for each task are bolded.

#### **5** Experiments

245

247

248

249

251

252

260

261

262

263

265

For our experiments <sup>1</sup>, we test our kernel for PEFT against both GLUE (Wang, 2018) and SuperGLUE benchmarks (Wang et al., 2019) with RoBERTa model (Liu, 2019), the results are summarized in Table 1. We train our model along with other baseline models using Hugging Face's Trainer framework with the default learning rate scheduler provided by the Transformers library, which is a linear scheduler with warmup (Wolf et al., 2020). The model is optimized with a batch size of 256 for training and 64 for evaluation. For GLUE benchmark experiments, we train all the models for 20 epochs and for the more challenging SuperGLUE benchmark experiments, we train all the models for 50 epochs to get a more faithful comparison. Regarding the model hyperparameter set up, we set the LoRA rank as 8 for all of our fine-tuning experiments, which leads to a parameter count of  $2 \times 768 \times 8 = 12,288$  for the attention layer of dimension  $768 \times 768$  with  $\alpha = 32$ . For our LSR-Kernel experiments, we

set r = 4 and  $A_k^{(1)} \in \mathbb{R}^{32 \times 2}, A_k^{(2)} \in \mathbb{R}^{24 \times 2}$ ,  $B_k^{(1)} \in \mathbb{R}^{32 \times 2}, B_k^{(2)} \in \mathbb{R}^{24 \times 2}$ , and separation rank s = 16, we can achieve a much lower parameter count  $2 \times (32 \times 2 + 24 \times 2) \times 16 = 3,584$ . All the other baseline methods follow the optimal settings given in the original papers (Edalati et al., 2022; He et al., 2023). From the results in Table 1, we can see that our method still maintains a high performance for the PEFT benchmark tasks with almost 25% of the LoRA parameters.

#### 6 Conclusion and Future Works

In this paper, we have demonstrated the effectiveness of adopting the separable representations in PEFT tasks. Specifically, we have shown that by restructuring the LoRA factor matrices using matrix low separation rank representations, we can not only drastically reduce the number of trainable parameters, but also provide more robust fine-tuning accuracy. However, in this study, we did not fully utilize the favorable computational attributes of Kronecker products. This could enhance the efficiency of computation during training, and we plan to explore this in future research.

287

266 267

268

<sup>&</sup>lt;sup>1</sup>For detailed experimental setups and implementations, please feel free to check out our GitHub Repository: https://anonymous.4open.science/r/lsr-adapt-7707.

289

7

Limitations

References

2159.

# 291

- 298

303

- 307
- 309

310

- 312
- 313 314

315

316 317 318

319

321

332 333

335 336 As discussed in the main paper, this work does

not exploit the amenable computation properties of

Kronecker products on modern tensor-core based

GPUs, which might lead to further memory effi-

ciency and faster training runtime, and potential

Gregory Beylkin and Martin J Mohlenkamp. 2005. Al-

Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and

of quantized llms. Preprint, arXiv:2305.14314.

Ali Edalati, Marzieh Tahaei, Ivan Kobyzev, Vahid Partovi Nia, James J Clark, and Mehdi Rezagholizadeh.

necker adapter. arXiv preprint arXiv:2212.10650.

G.H. Golub and C.F. Van Loan. 2013. Matrix Compu-

Xuehai He, Chunyuan Li, Pengchuan Zhang, Jianwei

Yang, and Xin Eric Wang. 2023. Parameter-efficient

model adaptation for vision transformers. In Proceed-

ings of the AAAI Conference on Artificial Intelligence,

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan

Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,

and Weizhu Chen. 2021. Lora: Low-rank adap-

tation of large language models. arXiv preprint

Abhinav Jangda and Mohit Yadav. 2024. Fast kronecker

Shih-Yang Liu, Chien-Yi Wang, Hongxu Yin, Pavlo Molchanov, Yu-Chiang Frank Wang, Kwang-Ting

Weight-decomposed low-rank adaptation. Preprint,

Roberta:

mized bert pretraining approach. arXiv preprint

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

analysis platform for natural language understanding.

Cheng, and Min-Hung Chen. 2024.

matrix-matrix multiplication on gpus. In Proceedings of the 29th ACM SIGPLAN Annual Symposium on Principles and Practice of Parallel Programming,

Sciences. Johns Hopkins University Press.

volume 37, pages 817-825.

arXiv:2106.09685.

pages 390-403.

arXiv:2402.09353.

arXiv:1907.11692, 364.

arXiv preprint arXiv:1804.07461.

Yinhan Liu. 2019.

tations. Johns Hopkins Studies in the Mathematical

2022. Krona: Parameter efficient tuning with kro-

Luke Zettlemoyer. 2023. Qlora: Efficient finetuning

gorithms for numerical analysis in high dimensions. SIAM Journal on Scientific Computing, 26(6):2133-

robustness to low-precision training.

Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2019. Superglue: A stickier benchmark for general-purpose language understanding systems. Advances in neural information processing systems, 32.

337

338

340

341

343

345

346

347

348

349

350

351

352

353

354

356

359

360

362

363

364

365

367

369

370

375

376

377

378

379

380

381

384

- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, and Jamie Brew. 2020. Transformers: State-of-theart natural language processing. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, pages 38–45, Online. Association for Computational Linguistics.
- Lingling Xu, Haoran Xie, Si-Zhao Joe Oin, Xiaohui Tao, and Fu Lee Wang. 2023. Parameter-efficient fine-tuning methods for pretrained language models: A critical review and assessment. arXiv preprint arXiv:2312.12148.

#### Appendix А

#### **Basic Properties of Kronecker Products** A.1

In this section we review some basic properties of the Kronecker products, which can be helpful in the case that we treat some matrix A as a block matrix whose the entries are all scalar multiplies of the same matrix (Golub and Van Loan, 2013). Denote  $A \in \mathbb{R}^{m \times n}$  where  $m = m_1 m_2$ ,  $n = n_1 n_2$ , then the matrix A is a Kronecker product means that there exist two Kronecker factor matrices  $B \in$  $\mathbb{R}^{m_1 \times n_1}$  and  $C \in \mathbb{R}^{m_2 \times n_2}$  such that,

$$\boldsymbol{A} = \boldsymbol{B} \otimes \boldsymbol{C}. \tag{14}$$

Some of the important Kronecker product properties include,

Transpose: $(\boldsymbol{B} \otimes \boldsymbol{C})^{\top} = \boldsymbol{B}^{\top} \otimes \boldsymbol{C}^{\top}$	371
Product: $(\boldsymbol{B}\otimes \boldsymbol{C})(\boldsymbol{D}\otimes \boldsymbol{E}) = \boldsymbol{B}\boldsymbol{D}\otimes \boldsymbol{C}\boldsymbol{E}$	372
Associativity: $B \otimes (C \otimes D) = (B \otimes C) \otimes D$ .	373

As for multiple Kronecker products, say A = $B \otimes C \otimes D$ , one can regard it as a block matrix whose entries are block matrices. Specifically, for (i, j)-th block of A, the value  $B_{i,j}C_{k,l}D$  is its (k, l)-th block.

### A.2 Understanding the Matrix Low **Separation Rank Representation**

Here we provide a mathematical analysis on why a matrix low-separation rank representation is an effective approximation mechanism for matrices. Suppose we have a matrix  $A \in \mathbb{R}^{m \times n}$  with

Dora:

A robustly opti-

386

387

 $rank(\mathbf{A}) = r$ , we would like to show that it admits the following approximation,

$$\boldsymbol{A} = \sum_{k=1}^{s} \lambda_k \boldsymbol{A}_k^{(1)} \otimes \cdots \otimes \boldsymbol{A}_k^{(r)} + \mathcal{O}(\epsilon), \quad (15)$$

with a separation rank s and  $A_k^{(i)} \in \mathbb{R}^{m_i \times n_i}$ , where  $\prod_{i=1}^r m_i = m, \prod_{i=1}^r n_i = n$ . We start from the fact that rank(A) = r, and thus there exist vectors  $u_1, \dots, u_r \in \mathbb{R}^m$  and  $v_1, \dots, v_r \in \mathbb{R}^n$  such that,

$$\boldsymbol{A} = \sum_{k=1}^{r} \boldsymbol{u}_k \boldsymbol{v}_k^{\top}.$$
 (16)

Then for each rank-1 matrix  $\boldsymbol{u}_k \boldsymbol{v}_k^\top \in \mathbb{R}^{m \times n}$ , we can do the reshaping,

$$\mathbf{u}_k = \mathbf{u}_k^{(1)} \otimes \cdots \otimes \mathbf{u}_k^{(r)}, \qquad (17)$$

$$\boldsymbol{v}_k = \boldsymbol{v}_k^{(1)} \otimes \cdots \otimes \boldsymbol{v}_k^{(r)}, \quad (18)$$

where each vector  $\boldsymbol{u}_{k}^{(i)} \in \mathbb{R}^{m_{i}}, \boldsymbol{v}_{k}^{(i)} \in \mathbb{R}^{n_{i}}$ . Thus from the basic Kronecker properties mentioned in A.1 we have,

401  

$$\boldsymbol{u}_{k}\boldsymbol{v}_{k}^{\top} = \left(\otimes_{i=1}^{r}\boldsymbol{u}_{k}^{(i)}\right)\left(\otimes_{i=1}^{r}\boldsymbol{v}_{k}^{(i)}\right)^{\top}$$
402  

$$= \left(\otimes_{i=1}^{r}\boldsymbol{u}_{k}^{(i)}\right)\left(\otimes_{i=1}^{r}\left(\boldsymbol{v}_{k}^{(i)}\right)^{\top}\right)$$

403

404

405

406

407

408

409

410

412

$$=\bigotimes_{i=1}^{r} \left( \boldsymbol{u}_{k}^{(i)} \left( \boldsymbol{v}_{k}^{(i)} \right)^{\mathsf{T}} \right).$$
 (20)

To see why the last equality in the above derivations works, consider the simpler example where we wish to compute

$$\left(oldsymbol{u}^{(1)}\otimesoldsymbol{u}^{(2)}\otimesoldsymbol{u}^{(3)}
ight)\left(oldsymbol{v}^{(1)}\otimesoldsymbol{v}^{(2)}\otimesoldsymbol{v}^{(3)}
ight),$$

if we define the substitutions  $U = u^{(1)} \otimes u^{(2)}$  and  $V = v^{(1)} \otimes v^{(2)}$ , with the Kronecker properties in A.1 we have the above equation becomes,

411 
$$\left( oldsymbol{U} \otimes oldsymbol{u}^{(3)} 
ight) \left( oldsymbol{V} \otimes oldsymbol{v}^{(3)} 
ight) = \left( oldsymbol{U} oldsymbol{V} 
ight) \otimes \left( oldsymbol{u}^{(3)} oldsymbol{v}^{(3)} 
ight)$$

where

413  

$$UV = \left(u^{(1)} \otimes u^{(2)}\right) \left(v^{(1)} \otimes v^{(2)}\right)$$

$$= \left(u^{(1)}v^{(1)}\right) \otimes \left(u^{(2)}v^{(2)}\right). \quad (21)$$

Then we substitute this back, yielding

$$\left(oldsymbol{u}^{(1)}\otimesoldsymbol{u}^{(2)}\otimesoldsymbol{u}^{(3)}
ight)\left(oldsymbol{v}^{(1)}\otimesoldsymbol{v}^{(2)}\otimesoldsymbol{v}^{(3)}
ight)$$
 416

$$= \left(\boldsymbol{u}^{(1)}\boldsymbol{v}^{(1)}\right) \otimes \left(\boldsymbol{u}^{(2)}\boldsymbol{v}^{(2)}\right) \otimes \left(\boldsymbol{u}^{(3)}\boldsymbol{v}^{(3)}\right). \tag{22}$$

Or in simplified notations,

$$\left(\otimes_{i=1}^{3}\boldsymbol{u}^{(i)}\right)\left(\otimes_{i=1}^{3}\boldsymbol{v}^{(i)}\right) = \bigotimes_{i=1}^{3}\left(\boldsymbol{u}^{(i)}\boldsymbol{v}^{(i)}\right).$$
(23)

Then if we define,

$$\boldsymbol{A}_{k}^{(i)} \triangleq \boldsymbol{u}_{k}^{(i)} \left(\boldsymbol{v}_{k}^{(i)}\right)^{\top} \in \mathbb{R}^{m_{i} \times n_{i}}, \qquad (24)$$

we essentially have

$$\boldsymbol{u}_k \boldsymbol{v}_k^\top = \boldsymbol{A}_k^{(1)} \otimes \cdots \otimes \boldsymbol{A}_k^{(r)}.$$
 (25) 42

To make the computations more controllable, we can set all the factor matrices  $A_k^{(i)}$  to be of unit norm (*e.g.*, Frobenius norm or operator norm) and factor out a scalar factor  $\lambda_k$ ,

$$\boldsymbol{u}_k \boldsymbol{v}_k^{\top} = \lambda_k \boldsymbol{A}_k^{(1)} \otimes \cdots \otimes \boldsymbol{A}_k^{(r)}.$$
 (26)

Then we have the form

$$oldsymbol{A} = oldsymbol{u}_k oldsymbol{v}_k^ op$$
 430

$$=\sum_{k=1}^{r}\lambda_k \boldsymbol{A}_k^{(1)}\otimes\cdots\otimes \boldsymbol{A}_k^{(r)},\qquad(27)$$

however, in practice, A might not be low rank and attaining the actual r can be expensive, hence if we instead set the approximate rank  $r \leq \operatorname{rank}(A)$ , we have the following approximation

$$\boldsymbol{A} = \sum_{k=1}^{s} \lambda_k \boldsymbol{A}_k^{(1)} \otimes \dots \otimes \boldsymbol{A}_k^{(r)} + \mathcal{O}(\epsilon), \quad (28)$$
436

where integer  $s \ge r$  is the separation rank and  $\epsilon$  is the approximation error.

418

419

420

421

422

424

425 426

427

428

429

431

432

433

434

435

437

438

415