



Group and Shuffle: Efficient Structured Orthogonal Parametrization

Mikhail Gorbunov¹ Nikolay Yudin¹ Vera Soboleva^{2, 1} Aibek Alanov^{2, 1}
Alexey Naumov^{1, 3} Maxim Rakhuba¹

¹HSE University ²AIRI ³Steklov Mathematical Institute RAS



1. Introduction

We propose a new low-dimensional subset of orthogonal matrices and a way to efficiently parametrize it. In particular, our contributions are:

- A new class of \mathcal{GS} -matrices that generalizes multiple previous structures.
- An optimal way to form a dense matrix in the \mathcal{GS} -class.
- Efficient structured orthogonal parametrization of \mathcal{GS} -matrices.
- Comparison of our class and existing methods in orthogonal fine-tuning (LMs and diffusion models) and orthogonal convolutional architectures.

2. Orthogonal Fine-Tuning

Orthogonal Fine-Tuning framework modifies forward pass of pre-trained linear layers:

$$\mathbf{y} = (\mathbf{W}^0)^\top \mathbf{x} \rightarrow \mathbf{y} = (\mathbf{Q}\mathbf{W}^0)^\top \mathbf{x}; \quad \mathbf{Q} - \text{orthogonal matrix.}$$

- **OFT** method [Qiu et al., 2023] uses block-diagonal structure for \mathbf{Q} , parametrizing it as

$$\mathbf{Q} = \text{diag}(\mathbf{Q}_1, \mathbf{Q}_2, \dots, \mathbf{Q}_r),$$

where $\mathbf{Q}_i \in \mathbb{R}^{b \times b}$ are small orthogonal matrices.

Problem: restrictive structure.

- **BOFT** method [Liu et al., 2024] parametrizes \mathbf{Q} as a product of several orthogonal sparse matrices, aiming to form a dense orthogonal matrix:

$$\mathbf{Q} = \mathbf{B}_m \mathbf{B}_{m-1} \dots \mathbf{B}_1.$$

\mathbf{B}_i is (up to permutations) a block-diagonal matrix with r orthogonal $\mathbf{b} \times \mathbf{b}$ blocks.

Problem: requires computing a product of multiple matrices. For example, for \mathbf{Q} of the size 1024×1024 and for $r = 32$, we have $m = 6$.

4. Orthogonal \mathcal{GS} -matrices

We can achieve orthogonality through enforcing it for every block of block-diagonal matrix \mathbf{B}_i . This is theoretically justified, as shown by **Theorem 2**.

Theorem 2

If each block of every \mathbf{B}_i in $\mathcal{GS}(\mathbf{P}_{m+1}, \dots, \mathbf{P}_1)$ is orthogonal then $\mathcal{GS}(\mathbf{P}_{m+1}, \dots, \mathbf{P}_1)$ is orthogonal matrix. In case $m = 2$, any orthogonal matrix from $\mathcal{GS}(\mathbf{P}_3, \mathbf{P}_2, \mathbf{P}_1)$ admits $\mathbf{P}_3(\mathbf{B}_2\mathbf{P}_2\mathbf{B}_1)\mathbf{P}_1$ representation with the matrices $\mathbf{B}_1, \mathbf{B}_2$ consisting of orthogonal blocks.

Remark: This theorem allows us to maintain orthogonality of smaller blocks instead of preserving orthogonality of whole matrix.

Parametrizing blocks: We mostly use *Cayley transform* to maintain orthogonality of each smaller block in \mathbf{B}_i : for any skew-symmetric matrix $\mathbf{K}^\top = -\mathbf{K}$, we know that \mathbf{Q} , given by

$$\mathbf{Q} = (\mathbf{I} + \mathbf{K})(\mathbf{I} - \mathbf{K})^{-1}$$

is orthogonal matrix. Alternatively, one can enforce orthogonality with the help of matrix exponent: for any skew-symmetric $\mathbf{K}^\top = -\mathbf{K}$, we have

$$\mathbf{Q} = \exp(\mathbf{K}), \quad \exp(\mathbf{K}) = \sum_{n=0}^{\infty} \frac{\mathbf{K}^n}{n!},$$

where \mathbf{Q} is orthogonal.

3. \mathcal{GS} -matrices

Definition 1

\mathbf{A} is said to be in $\mathcal{GS}(\mathbf{P}_{m+1}, \dots, \mathbf{P}_1)$ if

$$\mathbf{A} = \mathbf{P}_{m+1}\mathbf{B}_m\mathbf{P}_m \dots \mathbf{B}_1\mathbf{P}_1,$$

where each matrix \mathbf{B}_i is a block-diagonal matrix with k_i blocks of size $b_i^1 \times b_i^2$, matrices \mathbf{P}_i are permutation matrices and $b_i^1 \cdot k_i = b_{i+1}^2 \cdot k_{i+1}$.

Remark: \mathcal{GS} -class contains OFT, BOFT, order- p Monarch matrices and more.

Crucial question

How to choose \mathbf{P}_i to minimize m and ensure that \mathbf{A} is dense?

Answer

Make \mathbf{P}_i a *perfect shuffle* matrix $\mathbf{P}_{(r,br)}$! To visualize $\mathbf{y} = \mathbf{P}_{(r,br)}\mathbf{x}$ imagine splitting a deck (vector \mathbf{x}) into b piles and taking sequentially upper card from each pile to form \mathbf{y} .

The answer in a more formal way:

Theorem 1

Let $k_i = r, b_i^1 = b_i^2 = b$. Then using $m = 1 + \lceil \log_b(r) \rceil$ is sufficient for the class $\mathcal{GS}(\mathbf{P}_L, \mathbf{P}_{(r,br)}, \dots, \mathbf{P}_{(r,br)}, \mathbf{P}_R)$ to form a dense matrix for any $\mathbf{P}_L, \mathbf{P}_R$. Moreover, the choice of $\mathbf{P}_2 = \dots = \mathbf{P}_m = \mathbf{P}_{(r,br)}$ is optimal in the sense that all matrices from $\mathcal{GS}(\mathbf{P}_{m+1}, \dots, \mathbf{P}_1)$ contain zero blocks for any integer $m < 1 + \lceil \log_b(r) \rceil$ and any permutations $\mathbf{P}_1, \dots, \mathbf{P}_{m+1}$.

We also find that class $\mathcal{GS}(\mathbf{I}, \mathbf{P}, \mathbf{I})$ consists of block matrices with low-rank blocks, whose ranks are defined by permutation matrix \mathbf{P} .

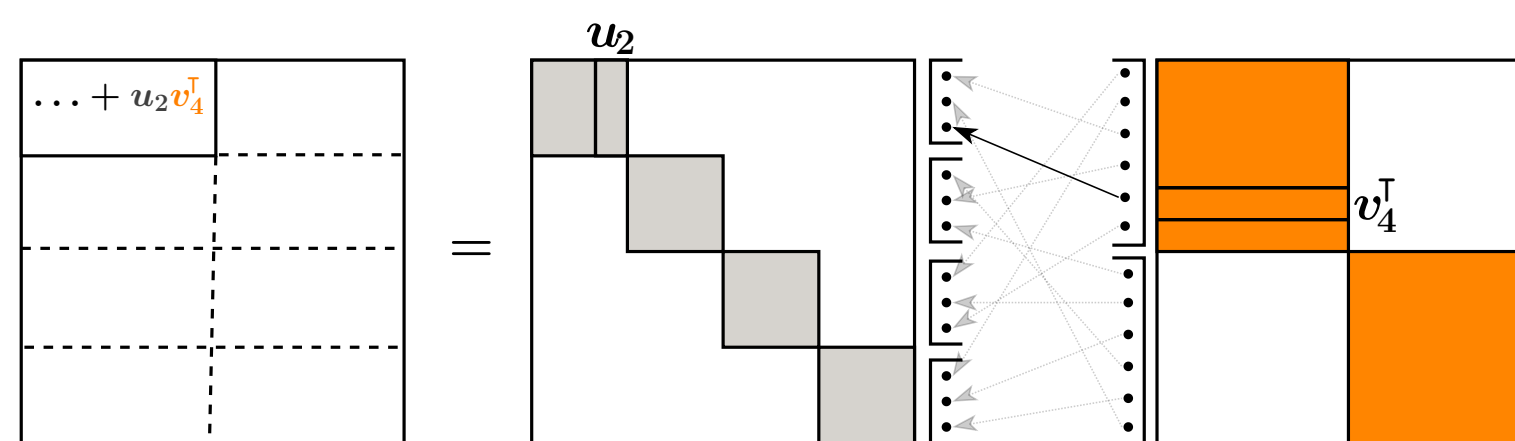


Figure 1. Illustration of block low-rank interpretation of $\mathcal{GS}(\mathbf{I}, \mathbf{P}, \mathbf{I})$ matrices.

Applications

- **NLP** Learnable orthogonal matrices for PEFT methods prevents training instabilities and overfitting that alternative methods like LoRA suffer from.
- **Diffusion models** Orthogonal fine-tuning of diffusion models helps to impose specific properties to a model without losing quality of generation.
- **Convolutional architectures** Equipping convolution operator with particular properties allows to obtain convolution with Lipschitz constant equal to 1.

Bibliography

Weiyang Liu, Zeju Qiu, Yao Feng, Yuliang Xiu, Yuxuan Xue, Longhui Yu, Haiwen Feng, Zhen Liu, Jueyong Heo, Songyou Peng, Yandong Wen, Michael J. Black, Adrian Weller, and Bernhard Schölkopf. Parameter-efficient orthogonal finetuning via butterfly factorization. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=7WzgkEdGyr>.

Zeju Qiu, Weiyang Liu, Haiwen Feng, Yuxuan Xue, Yao Feng, Zhen Liu, Dan Zhang, Adrian Weller, and Bernhard Schölkopf. Controlling text-to-image diffusion by orthogonal finetuning. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=K30wTdIIYc>.

Sahil Singla and Soheil Feizi. Skew orthogonal convolutions. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 9756–9766. PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/singla21a.html>.

Application 1: NLP

We utilize the pipeline of orthogonal fine-tuning parametrizing multiplicative matrix \mathbf{Q} with $\mathcal{GS}(\mathbf{P}_{(r,br)}^\top, \mathbf{P}_{(r,br)}, \mathbf{I})$ -matrices:

$$\mathbf{Q} = \mathbf{P}_{(r,br)}^\top \mathbf{B}_2 \mathbf{P}_{(r,br)} \mathbf{B}_1.$$

Method	# Params	MNLI	SST-2	CoLA	QQP	QNLI	RTE	MRPC	STS-B	ALL
FT	125M	<u>87.62</u>	94.38	61.97	91.5	93.06	80.14	88.97	90.91	86.07
LoRA _{r=8}	1.33M	87.82	95.07	64.02	<u>90.97</u>	92.81	81.95	88.73	<u>90.84</u>	<u>86.53</u>
OFT _{b=16}	1.41M	87.21	95.07	<u>64.37</u>	90.6	92.48	79.78	<u>89.95</u>	90.71	86.27
BOFT _{b=8} ^{m=2}	1.42M	87.14	94.38	<u>62.57</u>	90.48	92.39	80.14	88.97	90.67	85.84
GSOFT _{b=8}	1.42M	87.16	95.06	65.3	90.46	92.46	81.95	90.2	90.76	86.67

Table 1. Results on GLUE benchmark with RoBERTa-base model. We report Pearson correlation for STS-B, Matthew's correlation for CoLA and accuracy for other tasks.

Application 2: Subject-driven generation

For pre-trained diffusion models we investigate an approach that multiplies weight matrices from both sides with $\mathbf{Q}_U, \mathbf{Q}_V$ parametrized as orthogonal \mathcal{GS} -matrices:

$$\mathbf{y} = (\mathbf{W}^0)^\top \mathbf{x} \rightarrow \mathbf{y} = (\mathbf{Q}_U \mathbf{W}^0 \mathbf{Q}_V)^\top \mathbf{x}$$

Compared to one-sided adaptation, this modification allows to adapt both left and right singular vectors of \mathbf{W}^0 , improving method's flexibility.

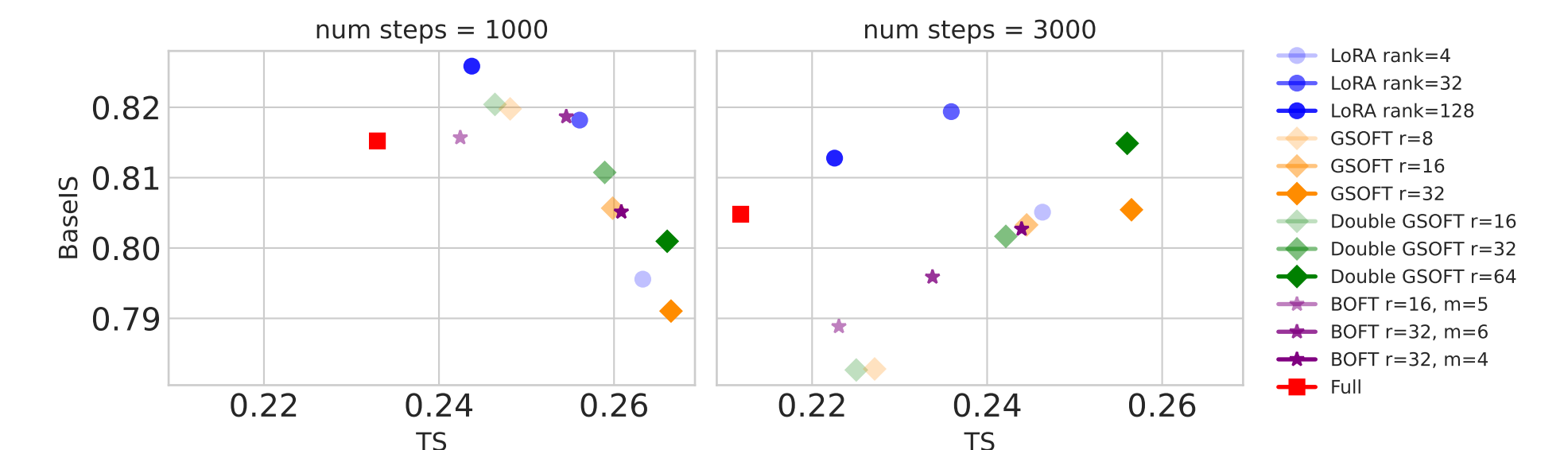


Figure 2. Image and text similarity visualization for different methods on subject-driven generation.

Application 3: Orthogonal Convolutions

We can also impose orthogonality to the convolution layer using skew-symmetric filters. Such constraint helps maintaining Lipschitz constant equal to 1 during training.

$$\mathbf{Y} = \text{ChShuffle}_2(\mathbf{L}_{grouped}^{(2)} \star_e (\text{ChShuffle}_1(\mathbf{L}_{grouped}^{(1)} \star_e \mathbf{X})))$$

$$\mathbf{L} \star_e \mathbf{X} = \mathbf{X} + \frac{\mathbf{L} \star \mathbf{X}}{1!} + \dots + \frac{\mathbf{L} \star^n \mathbf{X}}{n!} + \dots, \quad \star^i - \text{convolution applied } i \text{ times}$$

Table 2. Results of LipConvnet-15 on CIFAR-100. (a, b) in “Groups” denotes that we have two grouped exponential convolutions (first with *kernel_size* = 3, second with *kernel_size* = 1). If $b = “-”$, we only use one \mathcal{GS} orthogonal convolution. We use **ChShuffle** before each grouped convolution.

Conv. Layer	# Params	Groups	Speedup	Accuracy	Robust Accuracy
SOC [Singla and Feizi, 2021]	24.1M	-	1	43.15%	29.18%
GS-SOC (Ours)	6.81M	(4, -)	1.64	43.48%	29.26%
GS-SOC (Ours)	8.91M	(4, 1)	1.21	43.42%	29.56%
GS-SOC (Ours)	7.86M	(4, 2)	1.22	42.86%	28.98%
GS-SOC (Ours)	7.3M	(4, 4)	1.23	42.75%	28.7%