

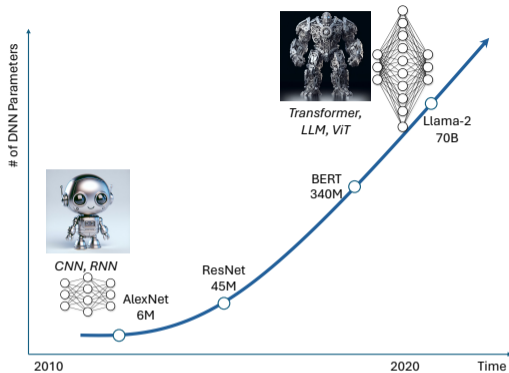
Differentiable Learning of Generalized Structured Matrices for Efficient Deep Neural Networks

Changwoo Lee, Hun-Seok Kim

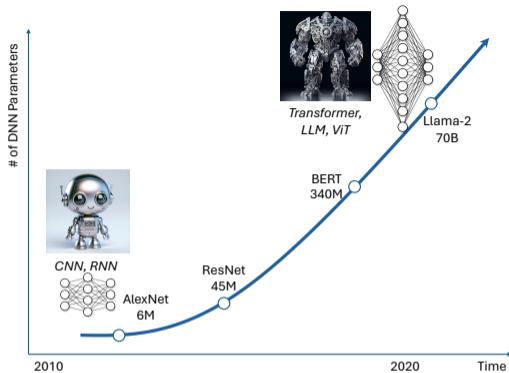
University of Michigan

ICLR 2024





Larger Deep Neural Networks (DNNs) { Better Accuracy 🦹



Larger Deep Neural Networks (DNNs)

{ Better Accuracy 🦾

{ Slower Inference, More Resources 💰

Problem

Goal: find a Deep Neural Network (DNN) with good performance and low cost of inference in floating-point operations (FLOPs)

$$\min_f \sum_{\mathbf{x}, \mathbf{y} \sim \mathcal{D}} \text{error}(f(\mathbf{x}), \mathbf{y}) \quad \text{s.t.} \quad \text{cost}(f) \leq B,$$

where f is a Deep Neural Network with L layers.

Problem

Goal: find a Deep Neural Network (DNN) with good performance and low cost of inference in floating-point operations (FLOPs)

$$\min_f \sum_{\mathbf{x}, \mathbf{y} \sim \mathcal{D}} \text{error}(f(\mathbf{x}), \mathbf{y}) \quad \text{s.t.} \quad \text{cost}(f) \leq B,$$

where f is a Deep Neural Network with L layers.

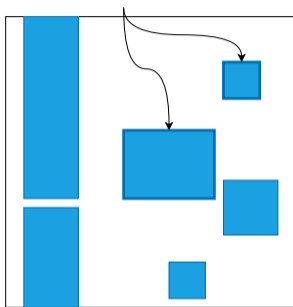
Use *structured matrices* for linear layers $\mathbf{W}\mathbf{x}$ to satisfy $\text{cost}(f) \leq B$.

- $\text{cost}(\mathbf{W}_{\text{structured}}\mathbf{x}) \ll \text{cost}(\mathbf{W}_{\text{dense}}\mathbf{x})$

$$\begin{matrix}
 \boxed{W_{LR}} & = & \boxed{U_r} \times & \boxed{V_r^T} \\
 n \times n & & n \times r & r \times n
 \end{matrix}$$

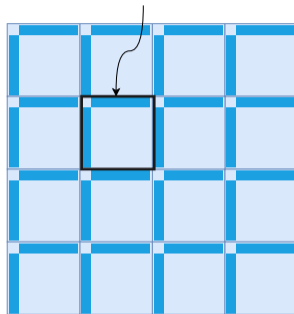
Low-Rank (LR)

High-rank Blocks



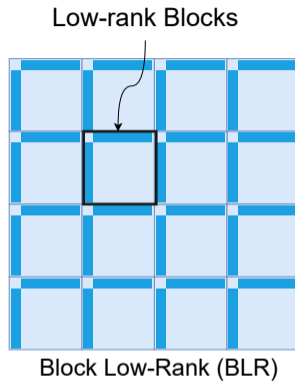
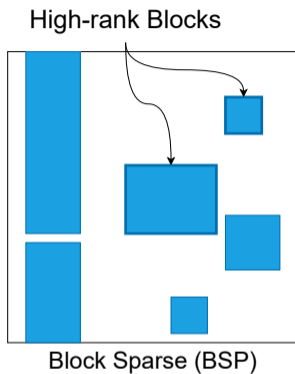
Block Sparse (BSP)

Low-rank Blocks



Block Low-Rank (BLR)

$$\begin{array}{ccc}
 \begin{array}{c} \boxed{W_{LR}} \\ n \times n \end{array} & = & \begin{array}{c} \boxed{U_r} \\ n \times r \end{array} \times \begin{array}{c} \boxed{V_r^T} \\ r \times n \end{array} \\
 \text{Low-Rank (LR)} & &
 \end{array}$$



Optimal *Layer-wise* Structure Format for Better Accuracy-Efficiency Trade-off?

Challenges

Challenge 1. Discrete Exponential Search Space

- Search space size exponential to the number of weights.

Challenges

Challenge 1. Discrete Exponential Search Space

- Search space size exponential to the number of weights.

Challenge 2. Lack of Structured Matrix Format

- Structured matrices are hand-designed from human insight.

Challenges

Challenge 1. Discrete Exponential Search Space

- Search space size exponential to the number of weights.

Challenge 2. Lack of Structured Matrix Format

- Structured matrices are hand-designed from human insight.

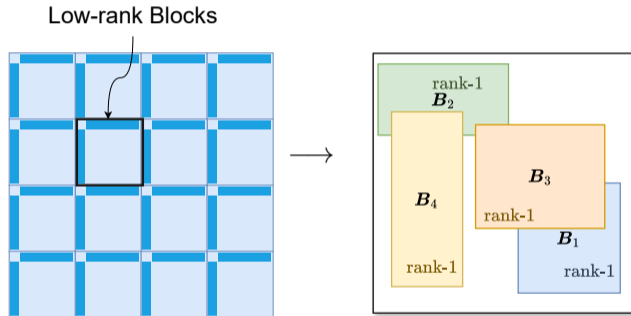
Solution: Differentiable Learning Approach

Generalize, Parameterize, then Optimize Structure by Gradient Descent!



Gaudi-GBLR Matrix

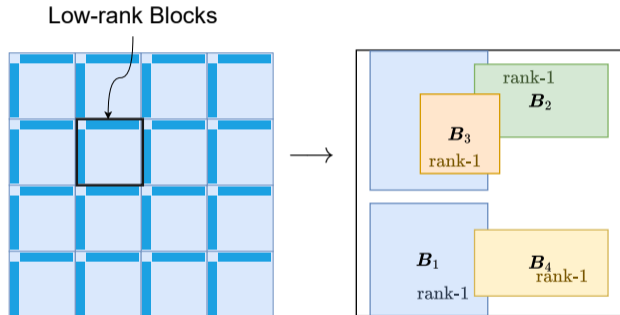
Generalized Block Low-Rank (GBLR) Matrix



Block Low-Rank

Generalized Block Low-Rank

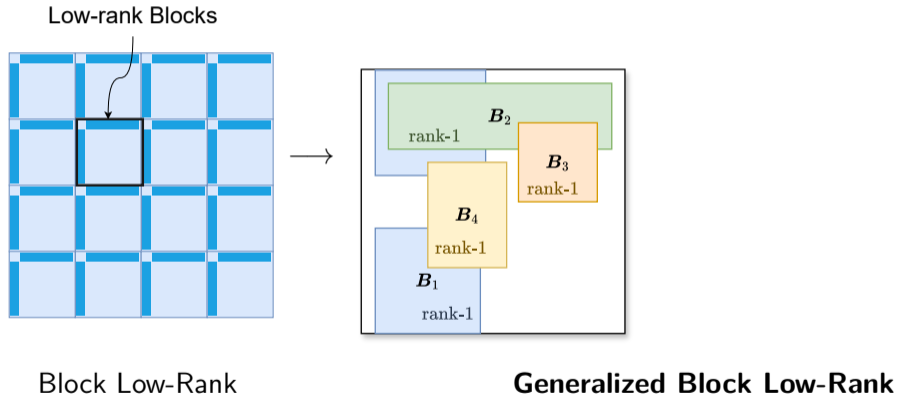
Generalized Block Low-Rank (GBLR) Matrix



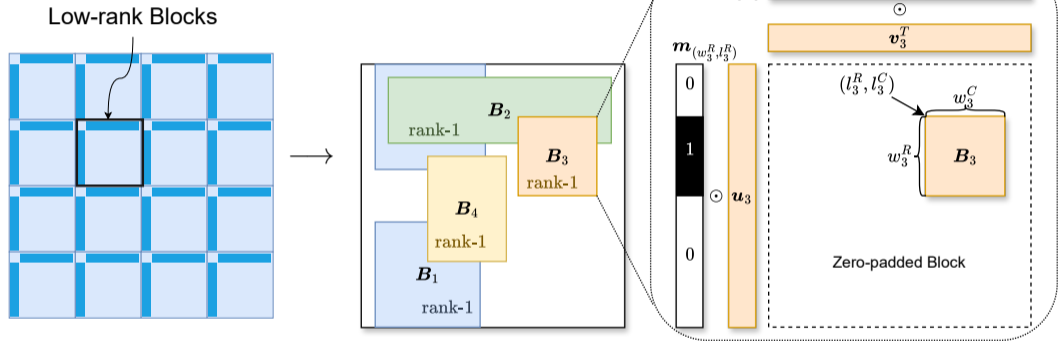
Block Low-Rank

Generalized Block Low-Rank

Generalized Block Low-Rank (GBLR) Matrix



Block by Mask



Block Low-Rank

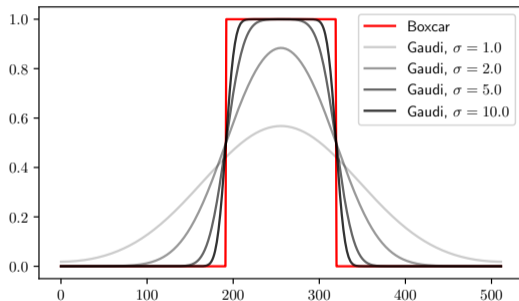
Faster Mult when
Low Rank Blocks

Generalized Block Low-Rank

Faster Mult when
Small Mask Widths

Differentiable Mask

Position of a Block = Two Gaussian-Dirichlet (Gaudi) Masks



- A Gaudi mask is fully differentiable with respect to the width and the starting point of the non-zero elements.

Gaudi-GBLR Matrix

A GBLR matrix with Gaudi masks $\tilde{\mathbf{m}}_{(w_k^R, l_k^R)}^\sigma, \tilde{\mathbf{m}}_{(w_k^C, l_k^C)}^\sigma$:

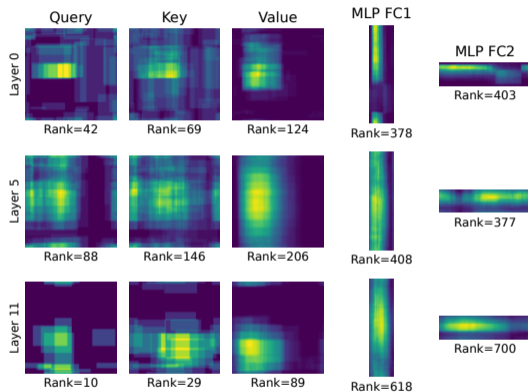
$$\mathbf{W}^\theta = \sum_{k=1}^K \left(\tilde{\mathbf{m}}_{(w_k^R, l_k^R)}^\sigma \odot \mathbf{u}_k \right) \left(\tilde{\mathbf{m}}_{(w_k^C, l_k^C)}^\sigma \odot \mathbf{v}_k \right)^T .$$

$\text{cost}(\mathbf{W}^\theta \mathbf{x}) \propto$ sum of the widths of the masks.

Training Efficient DNN by Gaudi-GBLR Weights

1. Replace weights of a DNN f to Gaudi-GBLR matrices.
2. Penalize the sum of the width of the Gaudi masks during the training at each gradient descent step.

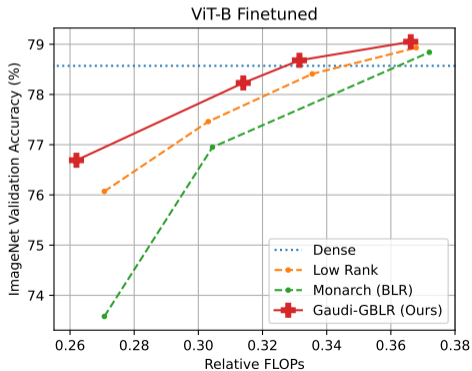
Learned Block Layout



Learned Block Layout of ViT-Base on ImageNet
The brighter, the more overlaps among blocks.

Vision Tasks

- ImageNet Classification (1,000 classes)
- Replace weights of ViT-Base to Gaudi-GBLR matrices and fine-tune.
- Better accuracy-efficiency trade-off than fixed structured matrices.

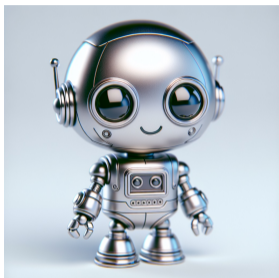


Language Task

- Language Generation-WikiText103
- Better performance while using fewer computations.

Perplexity by weight type of GPT-2 after fine-tuning on WikiText103.

Weight Type	Perplexity (\downarrow)	Relative FLOPs
Dense	19.36	100%
Low Rank	19.48	43.75%
Monarch	20.56	43.75%
Gaudi-GBLR	19.24	43.7%



Thank you!

Visit us at our poster on
Wed 8 May 10:45 a.m.-12:45 p.m. CEST



Paper



Code