# *Deformable Butterfly: A Highly Structured and Sparse Linear Transform*

**Rui Lin**[1,*]     **Jie Ran**[1,*]     Kung Hung Chiu[2]     Graziano Chesi[1]     Ngai Wong[1,*]

[1] Dept. of Electrical and Electronic Engineering, The University of Hong Kong

[2] United Microelectronic Center  (Hong Kong) Limited

Video Presentation For NeurIPS 2021

NEURAL INFORMATION PROCESSING SYSTEMS

* RL, JR, and NW contributed equally to this work.          Source Code: https://github.com/RuiLin0212/DeBut

# Content

# Butterfly Matrix



right ← left

2

2

**16 * 16**

**8 * 8**

**4 * 4**

**2 * 2**

*Standard* Butterfly matrix has **Powers-of-Two limitations.**

# Butterfly Matrix



Merged Information ← Information Flow → Individual Entries
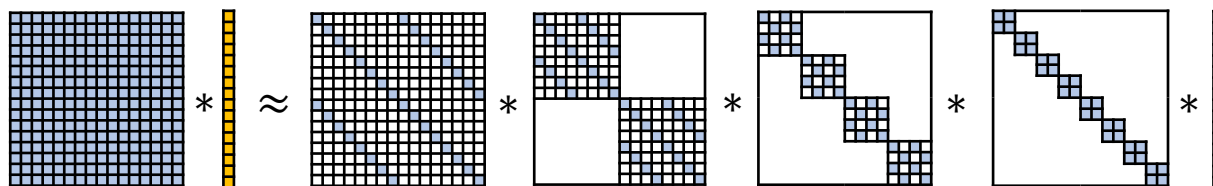
**Fully-mixed**    **Eight in a group**    **Four in a group**    **Two in a group**    **Separate numbers**

***From right to left***, the ***information carried by the input is merged***.

# Convolution as a Matrix Product



**Flattened Kernels**

$c_i$ ($c_o$): #input (output) channels
$k$ : kernel spatial size
$H_i$ ($H_o$): input (output) height
$W_i$ ($W_o$): input (output) width

Equivalent GEMM approach

DeBut substitute of $F$

Flattened filter matrix $F$

**Background**   **DeBut**   **Experiments**   **Conclusion**   **Butterfly Matrix**   **Convolution as Matrix Product**

# *Convolution as a Matrix Product*



**Flattened Input by im2col operation**

# Convolution as a Matrix Product



Equivalent GEMM approach

$c_i$ $(c_o)$: #input (output) channels
$k$ : kernel spatial size
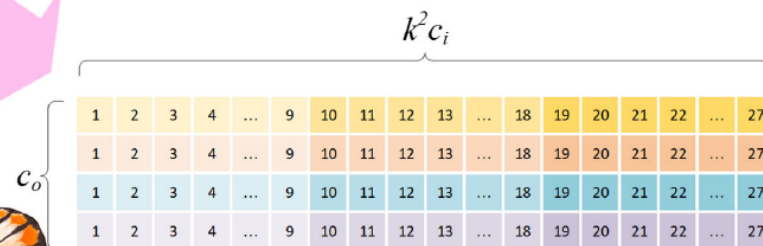$H_i$ $(H_o)$: input (output) height
$W_i$ $(W_o)$: input (output) width

DeBut substitute of $F$

Flattened filter matrix $F$

**Linear Transform**

# DeBut Chains

$R_{(2,3,1)}^{(8,12)}$

Non-zero element

$R_{(3,3,2)}^{(12,12)}$　　　$R_{(2,3,3)}^{(12,18)}$

*3*

*2*

*12*

*3*

*18*

**The spatial size of the factor**

$$R \, \frac{(p, q)}{(r,s\,,\,t)} \in \mathbb{R}^{p \times q}$$

**Number of diagonal matrices in the block**

*The size of the diagonal matrix*

**__No more__** *Powers-of-Two limitations!*

# DeBut Chains



A Bulging Chain: $128 \leftarrow_{(2, 4, 64)} 256 \leftarrow_{(2, 4, 32)} 512 \leftarrow_{(4, 5, 8)} 640 \leftarrow_{(8, 5, 1)} 400$

**A bulging DeBut chain**

Matrix Multiplication of all Matrices between Two Dashed Lines

$R^{(128,400)}_{(128,400,1)}$

$R^{(256,400)}_{(64,100,1)}$

$R^{(512,400)}_{(32,25,1)}$

$R^{(640,400)}_{(8,5,1)}$

R with dense sub-blocks

Another possible DeBut chain: $128 \leftarrow_{(2,2,64)} 128 \leftarrow_{(2,4,32)} 256 \leftarrow_{(32,50,1)} 400$

**A monotonic DeBut chain**

# DeBut Chains

**A bulging DeBut chain**

$$128 \xleftarrow[(2,4,64)]{} 256 \xleftarrow[(2,4,32)]{} \boxed{512 \xleftarrow[(4,5,8)]{} 640 \xleftarrow[(8,5,1)]{} 400}$$

**A monotonic DeBut chain**

$$128 \xleftarrow[(2,2,64)]{} \boxed{128 \xleftarrow[(2,4,32)]{} 256 \xleftarrow[(32,50,1)]{} 400}$$

# DeBut Chains



A Bulging Chain

$128 \leftarrow (2, 4, 64) \leftarrow 256 \leftarrow (2, 4, 32) \leftarrow 512 \leftarrow (4, 5, 8) \leftarrow 640 \leftarrow (8, 5, 1) \leftarrow 400$

Matrix Multiplication of all Matrices between Two Dashed Lines

$R^{(640,400)}_{(8,5,1)}$

$R^{(512,400)}_{(32,25,1)}$

$R^{(256,400)}_{(64,100,1)}$

R with dense sub-blocks

$R^{(128,400)}_{(128,400,1)}$

**Dense matrix at the end**

***No entry*** in the linear transform is ***forced zero***.

# *Alternating Least Squares (ALS)*



To Approximate $\approx$ F

$$error = \frac{\left|F - \hat{F}\right|_2}{\left|F\right|_2}$$

*Each time only one factor is updated*

# *LeNet Trained on MNIST (Baseline: 99.29%)*

| Layer | Monotonic/Bulging Chains | LC | MC | Params | Acc% (no ALS) | Acc% (w/ ALS) |
|-------|--------------------------|-----|-----|--------|---------------|---------------|
| FC1 | $\overset{256}{\underset{(1,2,32)}{\longleftarrow}}$ $\overset{256}{\underset{(2,2,16)}{\longleftarrow}}$ $\overset{400}{\underset{(16,25,1)}{\longleftarrow}}$ $\overset{128}{\underset{(2,2,64)}{\longleftarrow}}$ $\overset{128}{\underset{(2,2,32)}{\longleftarrow}}$ 128 | 85.00% | 70.78% | 17.96$K$ | 98.89($\pm$0.08) | 98.72($\pm$0.02) |
| CONV2 | $16\underset{(2,2,8)}{\longleftarrow}16\underset{(2,6,4)}{\longleftarrow}48\underset{(1,2,4)}{\longleftarrow}96\underset{(4,3,1)}{\longleftarrow}72$ | 41.67% | 1.04% | 60.84$K$ | 99.02($\pm$0.06) | 98.86($\pm$0.02) |
| CONV2 | $16\underset{(2,6,8)}{\longleftarrow}48\underset{(1,2,8)}{\longleftarrow}96\underset{(2,2,4)}{\longleftarrow}96\underset{(4,3,1)}{\longleftarrow}72$ | 41.67% | | | | |
| FC1 | $\overset{256}{\underset{(1,2,32)}{\longleftarrow}}$ $\overset{256}{\underset{(2,2,16)}{\longleftarrow}}$ $\overset{400}{\underset{(16,25,1)}{\longleftarrow}}$ $\overset{128}{\underset{(2,2,64)}{\longleftarrow}}$ $\overset{128}{\underset{(2,2,32)}{\longleftarrow}}$ 128 | 85.00% | 83.43% | 10.19$K$ | 98.64($\pm$0.15) | 97.27($\pm$0.02) |
| FC2 | $\overset{64}{\underset{(2,2,8)}{\longleftarrow}}$ $\overset{64}{\underset{(2,2,4)}{\longleftarrow}}$ $\overset{64}{\underset{(2,2,2)}{\longleftarrow}}$ $\overset{128}{\underset{(2,4,1)}{\longleftarrow}}$ $64\underset{(2,2,32)}{\longleftarrow}64\underset{(2,2,16)}{\longleftarrow}64$ | 89.06% | | | | |

*Reduce the #params significantly*

DeBut is able to ***reduce the number of parameters*** further while delivering ***promising output accuracy***.

# *LeNet Trained on MNIST (Baseline: 99.29%)*

| Layer | Monotonic/Bulging Chains | LC | MC | Params | Acc% (no ALS) | Acc% (w/ ALS) |
|---|---|---|---|---|---|---|
| FC1 | ← 256 ← 256 ← 400 (1,2,32) (2,2,16) (16,25,1)  128 ← 128 ← 128 (2,2,64) (2,2,32) | 85.00% | 70.78% | 17.96$K$ | 98.89($\pm$0.08) | 98.72($\pm$0.02) |
| CONV2 | 16 ← 16 ← 48 ← 96 ← 72 (2,2,8) (2,6,4) (1,2,4) (4,3,1) | 41.67% | 1.04% | 60.84$K$ | 99.02($\pm$0.06) | 98.86($\pm$0.02) |
| CONV2 | 16 ← 48 ← 96 ← 96 ← 72 (2,6,8) (1,2,8) (2,2,4) (4,3,1) | 41.67% | | | | |
| FC1 | ← 256 ← 256 ← 400 (1,2,32) (2,2,16) (16,25,1)  128 ← 128 ← 128 (2,2,64) (2,2,32) | 85.00% | 83.43% | 10.19$K$ | 98.64($\pm$0.15) | 97.27($\pm$0.02) |
| FC2 | ← 64 ← 64 ← 64 ← 128 (2,2,8) (2,2,4) (2,2,2) (2,4,1)  64 ← 64 ← 64 (2,2,32) (2,2,16) | 89.06% | | | | |

*w/o ALS has better performance in small examples*

In the **small example**, ALS initialization is **not always** beneficial for learning the latent information. However, its advantage will become obvious in **larger examples**.

# VGG Trained on CIFAR-10 (Baseline: 93.96%)

| Layer | Chain(s) | | | | | LC | MC | Params | Acc% (w/ ALS) |
|---|---|---|---|---|---|---|---|---|---|
| CONV13 | $4096 \xleftarrow{(2,2,16)} 4096 \xleftarrow{(2,2,8)} 4096 \xleftarrow{(8,9,1)} 4608$ $512 \xleftarrow{(2,4,256)} 1024 \xleftarrow{(2,4,128)} 2048 \xleftarrow{(2,4,64)} 4096 \xleftarrow{(2,2,32)}$ | | | | | 96.79% | 15.23% | 12.71M | 93.91(±0.08) |
| CONV8 | $2048 \xleftarrow{(2,2,32)} 2048 \xleftarrow{(2,2,16)} 2048 \xleftarrow{(2,2,8)} 2048 \xleftarrow{(8,9,1)} 2304$ $512 \xleftarrow{(2,2,256)} 512 \xleftarrow{(2,4,128)} 1024 \xleftarrow{(2,4,64)}$ | | | | | 96.79% | 83.77% | 2.43M | 93.72(±0.07) |
| CONV9~13 | $4096 \xleftarrow{(2,2,16)} 4096 \xleftarrow{(2,2,8)} 4096 \xleftarrow{(8,9,1)} 4608$ $512 \xleftarrow{(2,4,256)} 1024 \xleftarrow{(2,4,128)} 2048 \xleftarrow{(2,4,64)} 4096 \xleftarrow{(2,2,32)}$ | | | | | 96.79% | | | |

*Only a slight 0.24% accuracy drop*

Amazingly, this VGG-16-BN with DeBut layers achieves **_a remarkable MC of 83.77%_** with only **_a slight 0.24% accuracy drop_**.

# *ResNet-50 Trained on ImageNet (Baseline: 76.01%)*

| Model | MC | Params | Top-1(%) with ALS | Top-5(%) with ALS |
|---|---|---|---|---|
| ResNet-50 | – – | $25.55M$ | 76.01 | 92.93 |
| DeBut-bulging | 47.56% | $13.40M$ | 74.52 | 92.18 |
| DeBut-mono | 48.74% | $13.10M$ | 74.34 | 92.31 |

*47.56% Parameters fewer,*
*1.49% accuracy lower*

For ***DeBut-bulging***, the number of parameters reduces by ***47.56%***, the top-1 accuracy is 74.52%, ***1.47% lower*** compared with the baseline 76.01%.

Background    DeBut    Experiments    Conclusion    LeNet@MNIST  VGG@CIFAR-10  ResNet@ImageNet  Comparison

# ResNet-50 Trained on ImageNet (Baseline: 76.01%)

| Model | MC | Params | Top-1(%) with ALS | Top-5(%) with ALS |
|-------|------|--------|-------------------|-------------------|
| ResNet-50 | – – | $25.55M$ | 76.01 | 92.93 |
| DeBut-bulging | 47.56% | $13.40M$ | 74.52 | 92.18 |
| DeBut-mono | 48.74% | $13.10M$ | 74.34 | 92.31 |

*0.3M Parameters fewer than Debut-bulging,
the accuracy 74.34% is comparable*

For ***DeBut-mono***, the compressed model has ***0.3M fewer parameters than the DeBut-bulging*** model, yet still achieving ***a comparable 74.34% top-1 accuracy***.
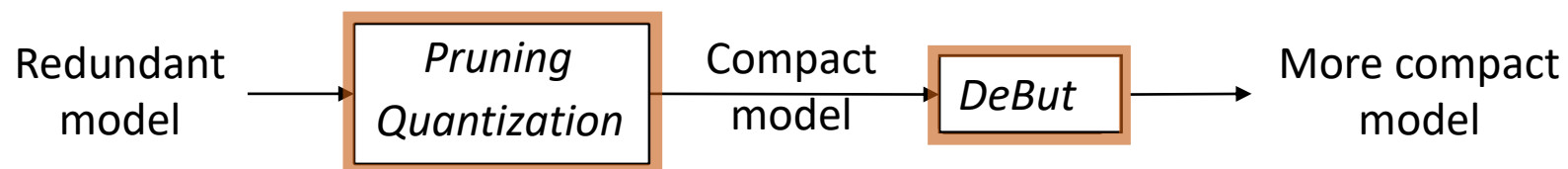
# DeBut vs. Other Linear Transform Schemes

| Layer | Method | MC | Params | Acc% | Training Time(s/epoch) | Inference Time(s) |
|-------|--------|-----|--------|------|------------------------|-------------------|
| CONV8~13 | Adaptive Fastfood | 85.65% | 2.15M | 93.60($\pm$0.02) | 2100 | 148.27 |
| | Butterfly | 85.82% | 2.13M | 93.34($\pm$0.12) | 105 | 4.58 |
| | DeBut | 83.77% | 2.43M | 93.72($\pm$0.07) | 50 | 4.01 |

*\* VGG trained on CIFAR-10 using different linear transform schemes*

**DeBut is the fastest algorithm**

**_Fast Hadamard Transform_** in Adaptive Fastfood is **exceedingly time-consuming** and makes it **_difficult_** to train Adaptive FastFood on **_multiple layers in a large CNN_**.
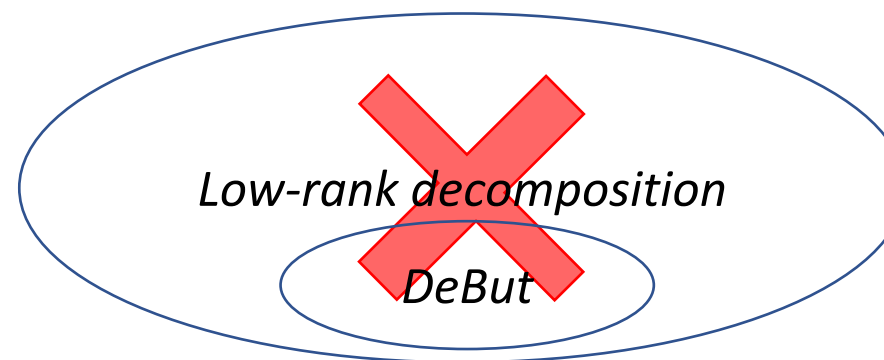
# *DeBut vs. Conventional Compression Schemes*

Redundant model → **Pruning Quantization** → Compact model → **DeBut** → More compact model

**Debut are orthogonal and complementary to Pruning & Quantization**

| Layer | Method | Params | Acc(%) |
|-------|--------|--------|--------|
| CONV8~13 | Baseline | 14.99$M$ | 93.96 |
| | Tucker-2 | 3.21$M$ | 93.36 |
| | DeBut | 2.43$M$ | 93.71 |

*\* VGG trained on CIFAR-10 using different compression methods*

**DeBut achieves better performance with fewer parameters**

*Low-rank decomposition*

*DeBut*

# *Conclusion*

1. DeBut Layers are ***truly practical*** .

2. DeBut is ***fundamentally different*** from standard butterfly.

3. DeBut layers ***further unify* and homogenize** FC and CONV layers.

4. Bulging and monotonic chains have different ***compression ability, performance and stability***.

5. The DeBut factor product chain provides an important implication for a ***pipelined DNN inference speedup***.

# *Thanks for Your Attention!*

*If you have any question, please contact us.*