

Pixelated Butterfly: Fast ML

with sparsity

Beidi Chen

NN Training Bottlenecks

- Exponentially ↑ Data volume
- ↑ model size
- ↑ resources

Popular direction: Sparsity

Not new

- lots of history

Existing approaches

- Pruning: Deep Comp
- Approx matmul: Reformer, SLIDE, Mongoose, Scatterbrain

Hard to ↑ speed without ↓ accuracy

Challenge & Goals

Challenges

- Dynamic sparsity can maintain accuracy but slow down training time
 - SOTA requires up to 5x more epochs
- Unstructured sparsity is NOT hardware efficient (Hooker et. al)
- Sparse attention target one module & hence does not speed up all layers
 - In many apps, MLP layer is the bottleneck

Ideal Sparsity pattern

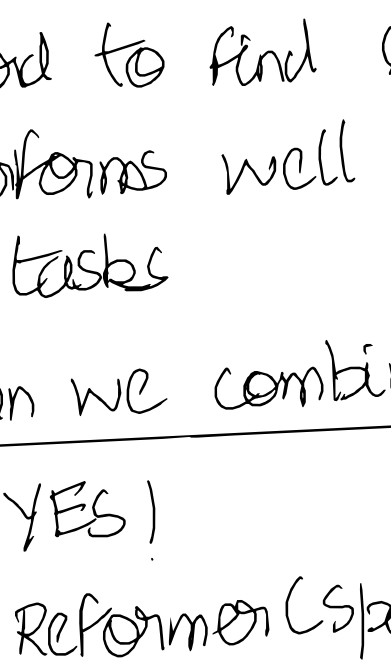
- static, simple yet accurate
- Aligned with available hardware
- Applied to most NN layers

ILDR: Butterfly + Low rank is a simple & effective sparsity pattern

① Attention Approximation: Sparse or low-rank

Trade accuracy for efficiency

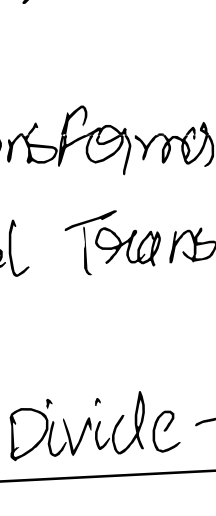
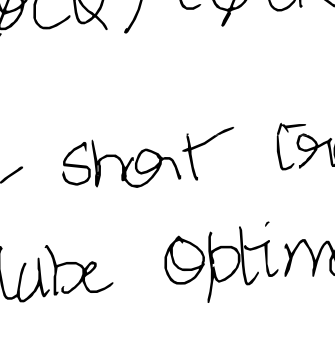
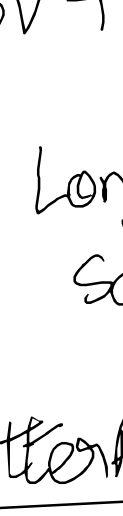
Sparse



Not all cells have values

Reformer, Sparse transform

Low Rank



use smaller matrix & only learn that

Linformer, Performer

↓
founder started a startup.

Hard to find robust approx. that performs well on wide variety of tasks

Can we combine sparse & low-rank?

YES!

Reformer (Sparse) + Performer (Low Rank)

$$SV + (U(Q) (U(K^T)V) \approx softmax(K^T)V$$

Long-short Transformers

Scalable Optimal Transport

Butterfly matrices: Divide-and-conquer

similar of FFT

Transform of size N

size N/2

size N/2

size N/4

size N/4

size N/4

size N/4

Class of fast linear maps that encodes the above pattern called butterfly

Each butterfly matrix is a product of log N sparse factors

Red dots indicates location of non zeros. The pattern is fixed & represents the recursion tree.

These matrices are trainable with grad descent

We can compose butterfly matrices to capture any sparse matrix with near optimal space & time complexity.

Butterfly + Low-Rank can

- avoid dynamic overhead (fixed pattern)
- apply to most matmul-based layers

HOWEVER, IT IS NOT HARDWARE EFFICIENT

Pixelated butterfly

Issues	Pixelated Butterfly
→ slow speed: sparsity patterns are not block aligned → not friendly to hardware	Block butterfly: Block-aligned sparsity pattern
→ Difficulty of parallelisation: Products of many ops → sequential ops	Flat butterfly: First order approx Turn product to sum
→ Reduced expressiveness: Flat butterfly are necessarily high rank	Low rank term: ↑ expressiveness of Flat block Butterfly matrices
→ cannot represent low rank matrices	

Combine Flat block butterfly + Low Rank

Workflow:

- Give model schema & budget allocation
- Pick a simple sparsity pattern based on flat block butterfly. Use this to produce mask on attention & also MLP layers

Theoretical properties

Theorem 1: Block butterfly retains the expressiveness of Butterfly & Flat butterfly can approx residual form of butterfly

Theorem 2: FBB ↑ expressive than sparse or low rank alone

Theorem 3: Training wide & sparse NN converge locally

Applications

- On CIFAR & Imagenet, train upto 2.3x faster & no accuracy loss
- 2.5x faster than GPT-2 for LM tasks
- ↑ model size → ↑ savings using FBB

NTK suggests that FBB generalizes close to best model

Conclusion

- Butterfly + low-rank works well
- Pixelated butterfly → static & block sparsity direction

Future directions

Go beyond dense models i.e FDE solving MRI reconstruction

Pipeline-hardware co-design

Efficient sparsity/butterfly on next gen ML accelerators

Data sparsity

Find structures in data & speed up training from a different angle