

# Optimal Quantization of Rank-One Matrices in Floating-Point Arithmetic

Elisa Riccietti

LIP – École Normale Supérieure de Lyon

*RAIM meeting 2025, 3-7 November 2025*



Rémi Gribonval (INRIA - LIP)



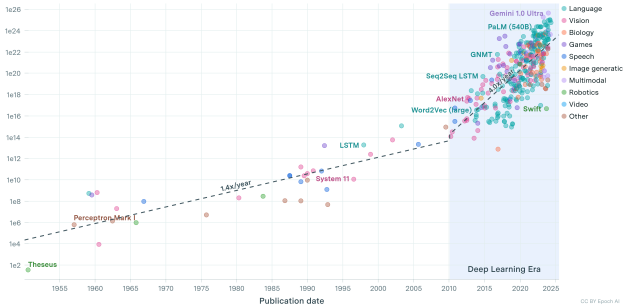
Théo Mary (CNRS - LIP6)

# Motivation

## Growing size of models and datasets

Notable AI models

Training compute (FLOP)

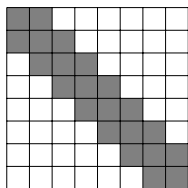


[J. Sevilla et al. 2022]

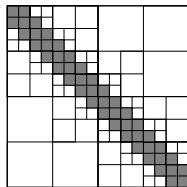
→ approximate computing

# Approximate computing (I)

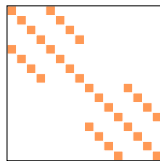
- Low-rank, structured, data sparse matrices



BLR matrix



$\mathcal{H}$ -matrix



Butterfly matrix

# Sparse matrix factorization

Given a dense matrix  $A$ , find *multiple* factors  $S_1, S_2, \dots S_L$  such that:

$$A \approx S_1 S_2 \dots S_L$$

where  $S_i$  are *sparse* matrices.

# Sparse matrix factorization

Given a dense matrix  $A$ , find *multiple* factors  $S_1, S_2, \dots, S_L$  such that:

$$A \approx S_1 S_2 \dots S_L$$

where  $S_i$  are *sparse* matrices.

$$\underbrace{A}_{\text{dense}} \approx \underbrace{S_1 S_2 \dots S_L}_{\text{sparse}} \Rightarrow Ax \approx S_1(S_2(\dots(S_L x)))$$

# Sparse matrix factorization

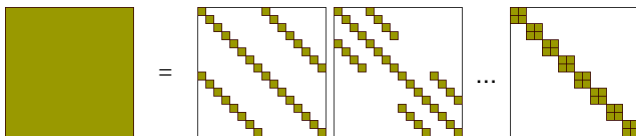
Given a dense matrix  $A$ , find *multiple* factors  $S_1, S_2, \dots, S_L$  such that:

$$A \approx S_1 S_2 \dots S_L$$

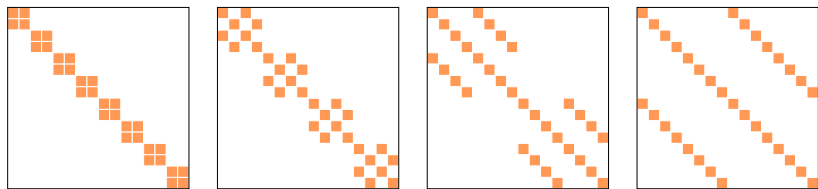
where  $S_i$  are *sparse* matrices.

$$\underbrace{A}_{\text{dense}} \approx \underbrace{S_1 S_2 \dots S_L}_{\text{sparse}} \Rightarrow Ax \approx S_1(S_2(\dots(S_L x)))$$

Application: fast transforms, networks compression...



# The butterfly factorization



- Butterfly matrices are extremely sparse yet highly expressive, they appear in many fast linear transforms
- Butterfly factorization: given dense  $A$   $n \times n$  matrix :

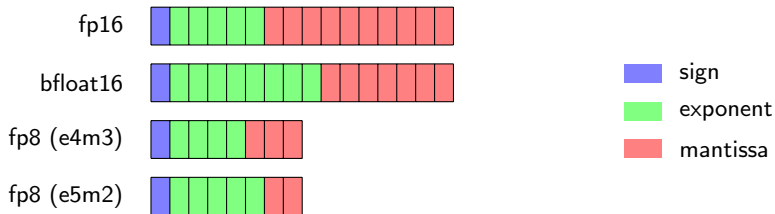
$$A \approx B_1 \dots B_L$$

with  $L = \log_2 n \Rightarrow O(n \log n)$  complexity



# Approximate computing (II)

- Quantization to low precision floating-point arithmetic



# Quantized butterfly factorization

Butterfly factorization + quantization

$$A \approx \hat{B}_1 \dots \hat{B}_L$$

**Problem:** Given a butterfly factorization  $B_1 \dots B_L$ , compute  $\hat{B}_1 \dots \hat{B}_L$  with optimal error

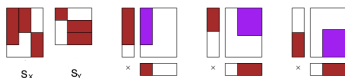
$$\frac{\|B_1 \dots B_L - \hat{B}_1 \dots \hat{B}_L\|}{\|B_1 \dots B_L\|}$$

# Key property: optimal two factors quantization

- If  $L = 2$  the problem can be optimally solved.
- Given two butterfly factors  $B_1, B_2$

$$B_1 B_2^T = \sum_{i=1}^n b_1(:, i) b_2(:, i)^T$$

where the rank-one matrices  $b_1(:, i) b_2(:, i)^T$  have disjoint support.

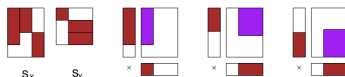


# Key property: optimal two factors quantization

- If  $L = 2$  the problem can be optimally solved.
- Given two butterfly factors  $B_1, B_2$

$$B_1 B_2^T = \sum_{i=1}^n b_1(:, i) b_2(:, i)^T$$

where the rank-one matrices  $b_1(:, i) b_2(:, i)^T$  have disjoint support.



- We can optimally quantize the product  $B_1 B_2^T$  by quantizing each  $b_1(:, i) b_2(:, i)^T$  optimally and independently

Part I  
Quantization of rank-one matrices

# Quantization of rank-one matrices

Goal: quantize the rank-one matrix

$$xy^T \rightarrow \widehat{x}\widehat{y}^T \quad (x \in \mathbb{R}^m, y \in \mathbb{R}^n)$$

where the coefficients of  $\widehat{x}$ ,  $\widehat{y}$  have  $t$  bits of mantissa

- The standard approach uses round-to-nearest (RTN) and leads to an error of order  $u = 2^{-t}$ : if  $\widehat{x} = \text{round}(x)$ ,  $\widehat{y} = \text{round}(y)$  then

$$\|\widehat{x} - x\| \leq u\|x\|$$

$$\|\widehat{y} - y\| \leq u\|y\|$$

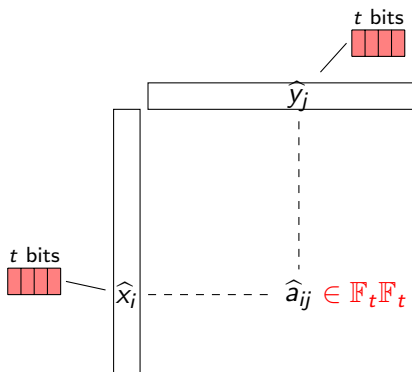
$$\Rightarrow \|\widehat{x}\widehat{y}^T - xy^T\| \leq (2u + u^2)\|x\|\|y\|$$

- We will show this is far from optimal!

# An analogy

Let  $\mathbb{F}_t$  be the set of  $t$ -bit floating-point numbers and

$$\mathbb{F}_t \mathbb{F}_t = \{a = xy, x \in \mathbb{F}_t, y \in \mathbb{F}_t\}$$

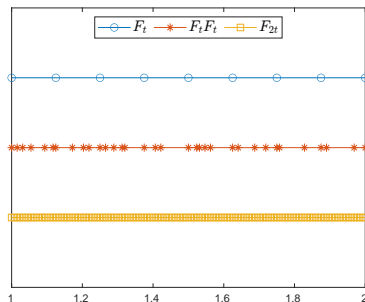


- What we really care about is the accuracy of  $\hat{a}_{ij} = \hat{x}_i \hat{y}_j$
- Think of multiword arithmetic:  $a \approx \hat{x} + \hat{y}$  with  $\hat{x} = \text{round}(a)$  and  $\hat{y} = \text{round}(a - \hat{x}) \rightarrow 2t$ -bit accuracy
- What about  $a = xy$ ? (Which  $\hat{x}, \hat{y}$  yields the best approximation  $\hat{x}\hat{y}$ ?)<sub>12/27</sub>

# The set $\mathbb{F}_t\mathbb{F}_t$

- Let  $\mathbb{F}_t$  be the set of  $t$ -bit floating-point numbers. We are interested in the set

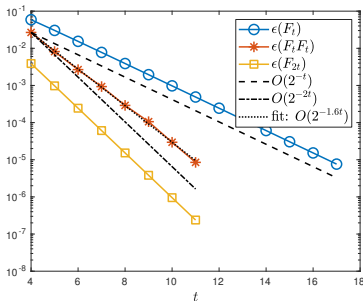
$$\mathbb{F}_t\mathbb{F}_t = \{a = xy, x \in \mathbb{F}_t, y \in \mathbb{F}_t\}$$



- No closed form expression of its elements, but we can simply enumerate all of them for small  $t$



# The set $\mathbb{F}_t \mathbb{F}_t$

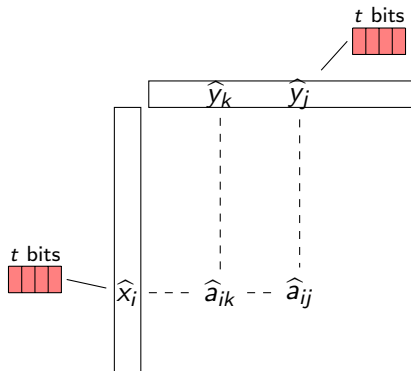


- $\epsilon(S) = \sup_{z \neq 0} \frac{d(z, S)}{|z|}$ : the **worst-case** relative error of quantizing an element  $z \in \mathbb{R}$  on  $S$

- $\epsilon(\mathbb{F}_t) = \frac{2^{-t}}{1+2^{-t}}$

$\Rightarrow$   $\epsilon(\mathbb{F}_t \mathbb{F}_t)$  error of order  $2^{-1.6t}$

# A constrained combinatorial problem



We don't just have one scalar, but a rank-one matrix  $\Rightarrow$

- We have constraints:  $\hat{x}_i$  must be the same in  $\hat{a}_{ij} = \hat{x}_i \hat{y}_j$  and  $\hat{a}_{ik} = \hat{x}_i \hat{y}_k$
- How can we find the optimal quantization? Combinatorial problem!

$$\min_{\hat{x} \in \mathbb{F}_t^m, \hat{y} \in \mathbb{F}_t^n} \|xy^T - \hat{x}\hat{y}^T\|$$

# Intuition to simplify the problem

In exact arithmetic:

$$xy^T = (\lambda x) \left( \frac{1}{\lambda} y \right)^T$$

In floating point arithmetic

$$\text{round}(xy^T) \neq \text{round}(\lambda x) \text{round} \left( \frac{1}{\lambda} y \right)^T$$

- Can we find the optimal scaling  $\lambda^*$ ?
- Can we reduce the problem to a scalar problem?

# Characterization of the optimum

## Theorem

$$\min_{\hat{x} \in \mathbb{F}_t^m, \hat{y} \in \mathbb{F}_t^n} \|xy^T - \hat{x}\hat{y}^T\| = \min_{\lambda \in \mathbb{R}} \|xy^T - \text{round}(\lambda x) \text{round}(\mu(\lambda)y)^T\|$$

The optimal quantization  $\hat{x}\hat{y}^T$  is given by

$$\hat{x} = \text{round}(\lambda x)$$

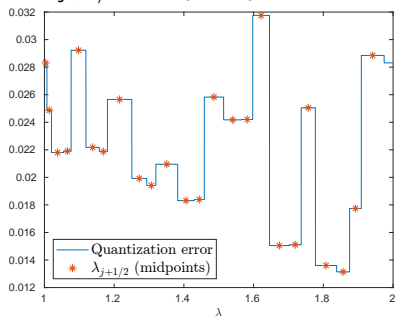
$$\hat{y} = \text{round}(\mu(\lambda)y^T)$$

where  $\lambda \in \mathbb{R}$  and  $\mu(\lambda) = \frac{x^T \hat{x}}{\|\hat{x}\|^2}$ .

- It suffices to find the optimal  $\lambda$  to find the optimal  $\hat{x}\hat{y}^T$  !

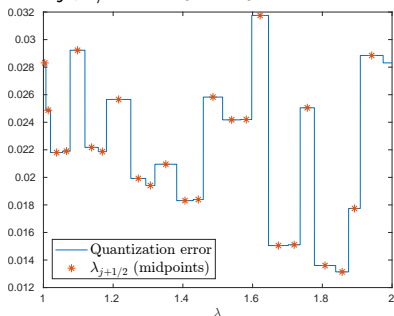
# Finding $\lambda$

- How do we find the optimal  $\lambda \in \mathbb{R}$  ?
- The optimum is stable under sign flip and multiplication by powers of two  $\rightarrow$  restrict the search to  $\lambda \in [1, 2]$
- Only a finite number of values of  $\lambda$  change the value of  $\text{round}(\lambda x)$ . Denoting these “breakpoints” as  $\lambda_j$ , we can enumerate the midpoints  $\lambda_{j+1/2} = (\lambda_j + \lambda_{j+1})/2$



# Finding $\lambda$

- How do we find the optimal  $\lambda \in \mathbb{R}$  ?
- The optimum is stable under sign flip and multiplication by powers of two  $\rightarrow$  restrict the search to  $\lambda \in [1, 2]$
- Only a finite number of values of  $\lambda$  change the value of  $\text{round}(\lambda x)$ . Denoting these “breakpoints” as  $\lambda_j$ , we can enumerate the midpoints  $\lambda_{j+1/2} = (\lambda_j + \lambda_{j+1})/2$

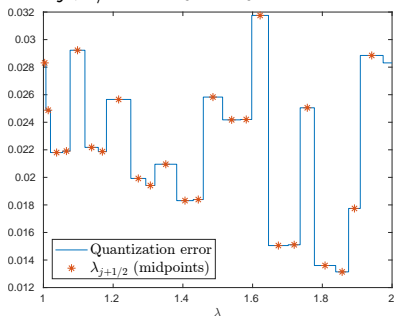


## Algorithm:

- Build the set of midpoints
- For each midpoint  $\lambda_{j+1/2}$ :
  - Build  $\hat{x} = \text{round}(\lambda_{j+1/2} x)$
  - Compute  $\mu(\hat{x}) = x^T \hat{x} / \|\hat{x}\|^2$
  - Build  $\hat{y} = \text{round}(\mu y)$
  - Test the accuracy of  $\hat{x} \hat{y}^T$

# Finding $\lambda$

- How do we find the optimal  $\lambda \in \mathbb{R}$  ?
- The optimum is stable under sign flip and multiplication by powers of two  $\rightarrow$  restrict the search to  $\lambda \in [1, 2]$
- Only a finite number of values of  $\lambda$  change the value of  $\text{round}(\lambda x)$ . Denoting these “breakpoints” as  $\lambda_j$ , we can enumerate the midpoints  $\lambda_{j+1/2} = (\lambda_j + \lambda_{j+1})/2$

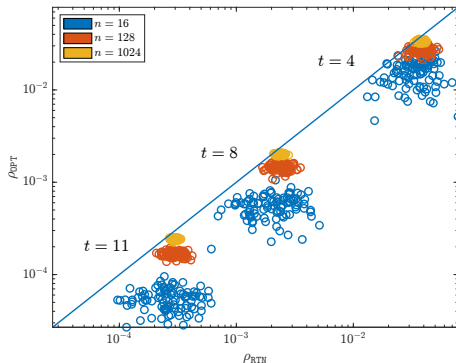


## Algorithm:

- Build the set of midpoints
- For each midpoint  $\lambda_{j+1/2}$ :
  - Build  $\hat{x} = \text{round}(\lambda_{j+1/2} x)$
  - Compute  $\mu(\hat{x}) = x^T \hat{x} / \|\hat{x}\|^2$
  - Build  $\hat{y} = \text{round}(\mu y)$
  - Test the accuracy of  $\hat{x} \hat{y}^T$

$O(mn2^t)$  complexity  $\Rightarrow$  tractable for large matrices and low precisions

# Experiments



$$\rho_{OPT} = \|xy^T - \widehat{x}\widehat{y}^T\| / \|xy^T\|,$$

$$\rho_{RTN} = \|xy^T - \text{round}(x) \text{round}(y)^T\| / \|xy^T\|.$$

100 randomly chosen couples  $(x, y) \in \mathbb{R}^n \times \mathbb{R}^n$



## Part II

### Application to butterfly factorization

# Quantization of butterfly factorization

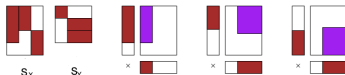
**Aim:** Given  $B_1 \dots B_L$ , compute  $\hat{B}_1 \dots \hat{B}_L$

- Key property: optimal two factor quantization
- For any partial product  $XY^T$  of consecutive factors

$$B_1 \dots B_j \underbrace{B_{j+1} \dots B_k}_X \underbrace{B_{k+1} \dots B_\ell}_{Y^T} B_{\ell+1} \dots B_L$$

$$XY^T = \sum_{i=1}^n x_i y_i^T$$

where the rank-one matrices  $x_i y_i^T$  have disjoint support.



# Quantization of butterfly factorization

- We can use our optimal algorithm to quantize each  $x_i y_i^T$  optimally and independently:  $\hat{x}_i = \text{round}(\lambda_i x_i)$ ,  $\hat{y}_i = \text{round}(\mu_i y_i)$
- We then obtain

$$\begin{aligned}\hat{X} &= \text{round}(X\Lambda), & \Lambda &= \text{diag}(\lambda_i) \\ \hat{Y} &= \text{round}(YM), & M &= \text{diag}(\mu_i)\end{aligned}$$

# Heuristics for the $L$ -factor butterfly quantization

When  $L > 2$ , need heuristics to decide how to order/group the factors

- Pairwise heuristic:

$$B_1 B_2 B_3 B_4 \dots B_L$$

# Heuristics for the $L$ -factor butterfly quantization

When  $L > 2$ , need heuristics to decide how to order/group the factors

- Pairwise heuristic:

$$\underbrace{B_1 B_2}_{XY^T} B_3 B_4 \dots B_L$$

# Heuristics for the $L$ -factor butterfly quantization

When  $L > 2$ , need heuristics to decide how to order/group the factors

- Pairwise heuristic:

$$\underbrace{\hat{B}_1 \hat{B}_2}_{XY^T} B_3 B_4 \dots B_L$$

# Heuristics for the $L$ -factor butterfly quantization

When  $L > 2$ , need heuristics to decide how to order/group the factors

- Pairwise heuristic:

$$\underbrace{\hat{B}_1 \hat{B}_2}_{XY^T} \underbrace{B_3 B_4}_{XY^T} \dots B_L$$

# Heuristics for the $L$ -factor butterfly quantization

When  $L > 2$ , need heuristics to decide how to order/group the factors

- Pairwise heuristic:

$$\underbrace{\hat{B}_1 \hat{B}_2}_{XY^T} \underbrace{\hat{B}_3 \hat{B}_4}_{XY^T} \dots B_L$$



# Heuristics for the $L$ -factor butterfly quantization

When  $L > 2$ , need heuristics to decide how to order/group the factors

- Pairwise heuristic:

$$\underbrace{\hat{B}_1 \hat{B}_2}_{XY^T} \underbrace{\hat{B}_3 \hat{B}_4}_{XY^T} \dots \underbrace{\hat{B}_L}_{RTN}$$

# Heuristics for the $L$ -factor butterfly quantization

When  $L > 2$ , need heuristics to decide how to order/group the factors

- Pairwise heuristic:

$$\underbrace{\hat{B}_1 \hat{B}_2}_{XY^T} \underbrace{\hat{B}_3 \hat{B}_4}_{XY^T} \dots \underbrace{\hat{B}_L}_{RTN}$$

- Left-to-right heuristic:

$$B_1 B_2 B_3 \dots B_{L-1} B_L$$

# Heuristics for the $L$ -factor butterfly quantization

When  $L > 2$ , need heuristics to decide how to order/group the factors

- Pairwise heuristic:

$$\underbrace{\hat{B}_1 \hat{B}_2}_{XY^T} \underbrace{\hat{B}_3 \hat{B}_4}_{XY^T} \dots \underbrace{\hat{B}_L}_{RTN}$$

- Left-to-right heuristic:

$$\underbrace{B_1}_X \underbrace{B_2 B_3 \dots B_{L-1} B_L}_{Y^T}$$

# Heuristics for the $L$ -factor butterfly quantization

When  $L > 2$ , need heuristics to decide how to order/group the factors

- Pairwise heuristic:

$$\underbrace{\hat{B}_1 \hat{B}_2}_{XY^T} \underbrace{\hat{B}_3 \hat{B}_4}_{XY^T} \dots \underbrace{\hat{B}_L}_{RTN}$$

- Left-to-right heuristic:

$$\underbrace{\hat{B}_1}_X \underbrace{M_2 B_2 B_3 \dots B_{L-1} B_L}_{Y^T}$$

# Heuristics for the $L$ -factor butterfly quantization

When  $L > 2$ , need heuristics to decide how to order/group the factors

- Pairwise heuristic:

$$\underbrace{\hat{B}_1 \hat{B}_2}_{XY^T} \underbrace{\hat{B}_3 \hat{B}_4}_{XY^T} \dots \underbrace{\hat{B}_L}_{RTN}$$

- Left-to-right heuristic:

$$\underbrace{\hat{B}_1}_X \underbrace{\hat{M}_2 \hat{B}_2}_X \underbrace{B_3 \dots B_{L-1} B_L}_{Y^T}$$

# Heuristics for the $L$ -factor butterfly quantization

When  $L > 2$ , need heuristics to decide how to order/group the factors

- Pairwise heuristic:

$$\underbrace{\hat{B}_1 \hat{B}_2}_{XY^T} \underbrace{\hat{B}_3 \hat{B}_4}_{XY^T} \dots \underbrace{\hat{B}_L}_{RTN}$$

- Left-to-right heuristic:

$$\underbrace{\hat{B}_1}_{X} \underbrace{\hat{B}_2}_X \underbrace{M_3 B_3 \dots B_{L-1} B_L}_{Y^T}$$

# Heuristics for the $L$ -factor butterfly quantization

When  $L > 2$ , need heuristics to decide how to order/group the factors

- Pairwise heuristic:

$$\underbrace{\hat{B}_1 \hat{B}_2}_{XY^T} \underbrace{\hat{B}_3 \hat{B}_4}_{XY^T} \dots \underbrace{\hat{B}_L}_{RTN}$$

- Left-to-right heuristic:

$$\begin{array}{ccccccc} \hat{B}_1 & \hat{B}_2 & \hat{B}_3 & \dots & M_{L-1} & B_{L-1} & B_L \\ & & & & \underbrace{\hspace{1.5cm}} & & \underbrace{\hspace{1cm}} \\ & & & & & X & Y^T \\ & \underbrace{\hspace{1cm}} & & \underbrace{\hspace{3cm}} & & & \\ & X & & Y^T & & & \\ \underbrace{\hspace{1cm}} & \underbrace{\hspace{5cm}} & & & & & \\ X & Y^T & & & & & \end{array}$$

# Heuristics for the $L$ -factor butterfly quantization

When  $L > 2$ , need heuristics to decide how to order/group the factors

- Pairwise heuristic:

$$\underbrace{\hat{B}_1 \hat{B}_2}_{XY^T} \underbrace{\hat{B}_3 \hat{B}_4}_{XY^T} \dots \underbrace{\hat{B}_L}_{RTN}$$

- Left-to-right heuristic:

$$\begin{array}{ccccccc} \hat{B}_1 & \hat{B}_2 & \hat{B}_3 & \dots & \hat{B}_{L-1} & \hat{B}_L \\ & & & & \underbrace{\hspace{1cm}}_X & \underbrace{\hspace{1cm}}_{Y^T} \\ & & \underbrace{\hspace{1cm}}_X & & \underbrace{\hspace{1cm}}_{Y^T} & & \\ \underbrace{\hspace{1cm}}_X & & \underbrace{\hspace{1cm}}_{Y^T} & & & & \end{array}$$



# Heuristics for the $L$ -factor butterfly quantization

When  $L > 2$ , need heuristics to decide how to order/group the factors

- Pairwise heuristic:

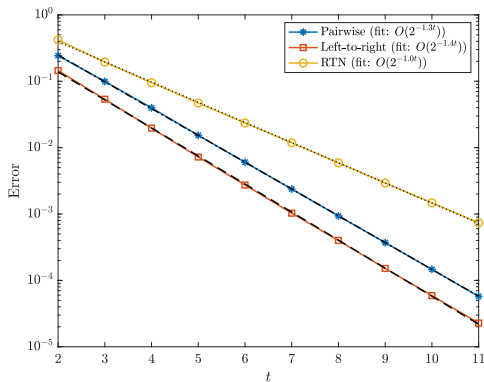
$$\underbrace{\hat{B}_1 \hat{B}_2}_{XY^T} \underbrace{\hat{B}_3 \hat{B}_4}_{XY^T} \dots \underbrace{\hat{B}_L}_{RTN}$$

- Left-to-right heuristic:

$$\begin{array}{ccccccc} \hat{B}_1 & \hat{B}_2 & \hat{B}_3 & \dots & \hat{B}_{L-1} & \hat{B}_L \\ & & & & \underbrace{\hspace{1.5cm}}_X & \underbrace{\hspace{1cm}}_{Y^T} \\ & & \underbrace{\hspace{1cm}}_X & \underbrace{\hspace{2.5cm}}_{Y^T} & & \\ \underbrace{\hspace{1cm}}_X & \underbrace{\hspace{4.5cm}}_{Y^T} & & & & \end{array}$$

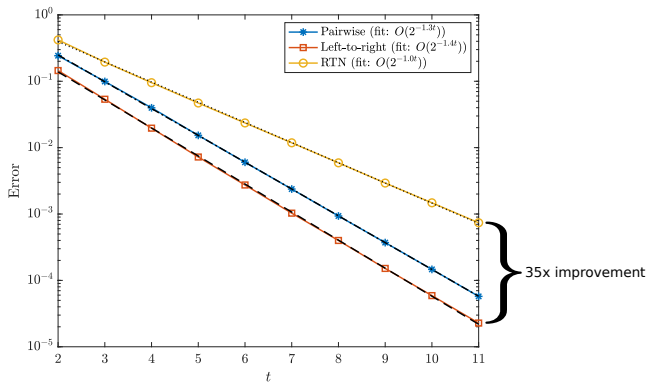
- L2R more expensive because it densifies the factors
- If  $L > 2$  the optimality is lost

# Experimental results



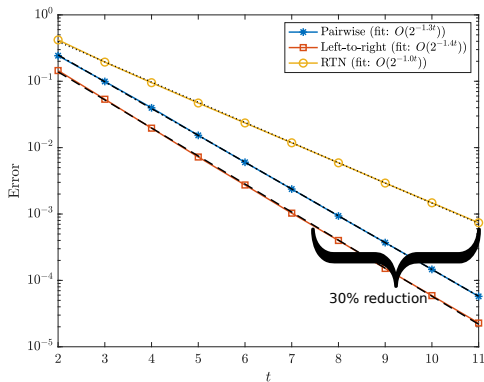
- Randomly generated butterfly factors

# Experimental results



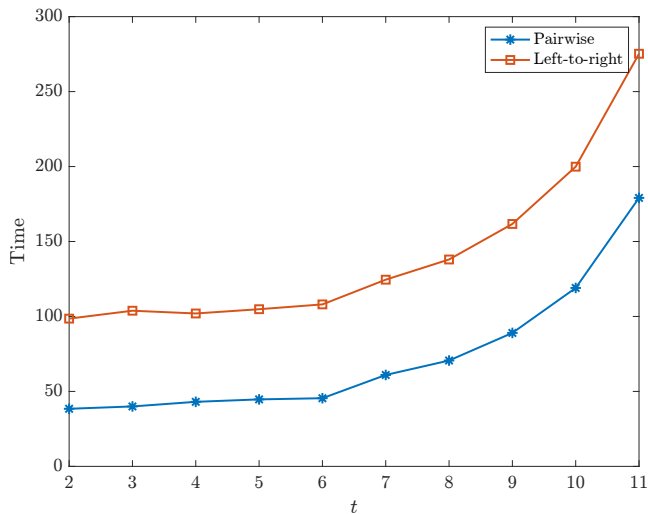
- Randomly generated butterfly factors
- Significant accuracy improvement. . .

# Experimental results



- Randomly generated butterfly factors
- Significant accuracy improvement. . .
- . . . or, equivalently, **can reduce storage by about 30% with no loss of accuracy**

# Experimental results



# Conclusion

Key results:

- Characterized optimal quantization of  $xy^T$  as  $\text{round}(\lambda x) \text{round}(\mu y)^T$
- Proposed algorithm to find the optimal  $\lambda$  in  $O(mn2^t)$  complexity
- Proposed two heuristics to apply method to butterfly factorization and obtained **storage reductions of 30% with no loss of accuracy**

**Limitation:** most of the fast transforms involving butterflies are complex-valued  $\rightarrow$  Maël Chaumette's talk tomorrow !

**Perspectives:** Butterfly matrices are only one possible application, many other perspectives: rank- $r$  matrices, tensors, DNNs, ...

# In practice: the FA $\mu$ ST library

**FA $\mu$ ST library:** an implementation of the hierarchical algorithm, fast GPU matrix vector multiplication of butterfly matrices, quantization algorithm in C++ via Python and Matlab wrappers **FA $\mu$ ST 3.25 toolbox** at <https://faust.inria.fr/>.

To know more:



R. Gribonval, T. Mary, E. Riccietti (2024), Optimal quantization of rank-one matrices in floating point arithmetic - with applications to butterfly factorizations, in revision for SISC.



Q.-T. Le, E. Riccietti, and R. Gribonval (2023), Spurious Valleys, Spurious Minima and NP-hardness of Sparse Matrix Factorization With Fixed Support, SIMAX.

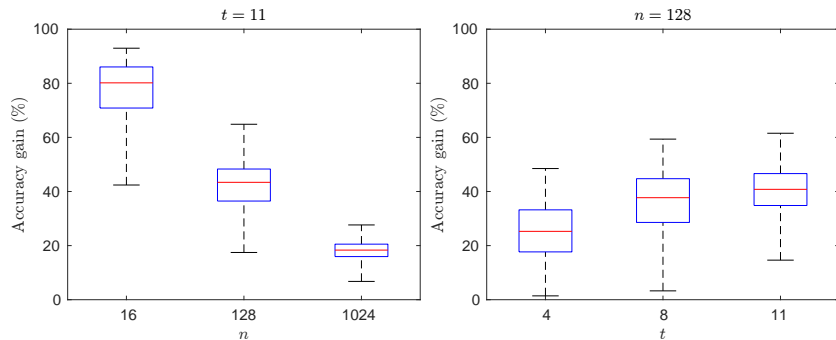


L. Zheng, E. Riccietti, and R. Gribonval (2023), Efficient Identification of Butterfly Sparse Matrix Factorizations, SIMODS.



Q.-T. Le, L. Zheng, E. Riccietti, and R. Gribonval (2022), Fast learning of fast transforms, with guarantees, ICASSP 2022

# Experiments



$$100 \times \left( 1 - \frac{\rho_{OPT}}{\rho_{RTN}} \right)$$



# Experiments

