

Motivation: Reduce memory footprint by lowering numerical precision in numerical algebra (e.g. FFT) by exploiting **rescaling invariances**.

Rank-One Quantization problem

$$\min_{\hat{x} \in \mathbb{C}\mathbb{F}_t^m, \hat{y} \in \mathbb{C}\mathbb{F}_t^n} \left\| \underset{\text{unquantized}}{xy^H} - \underset{\text{quantized}}{\hat{x}\hat{y}^H} \right\|_F^2$$

Q: how to efficiently solve the problem with **quantized** unknowns?

Optimally solved in the real case [1] by exploiting the **rescaling invariance** property $xy^H = (\lambda x) \left(\frac{1}{\lambda} y \right)^H$

Round-To-Nearest (RTN)

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

$$\hat{y} = \text{round}(y)$$

Separately maps each coefficient to their nearest neighbor in $\mathbb{C}\mathbb{F}_t$

Fast but potentially not **optimal**

Key characterization

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

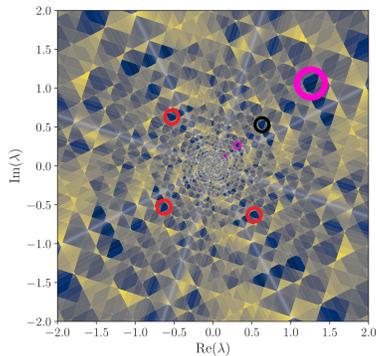
$$\hat{y}^* = \text{round}(\mu(\lambda^*) y)$$

where $\lambda^* \in \arg \inf_{\lambda \in \mathbb{C}} f(\lambda)$

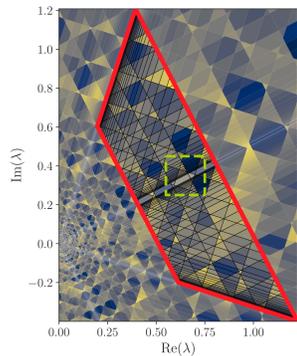
Reduction of a problem with $2(m+n)$ quantized variables to a **one scalar** problem

Geometrical properties of f

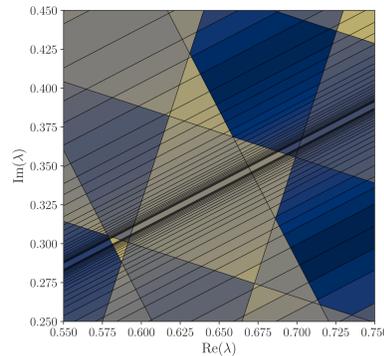
1. Invariance by multiplication by 2 and by i



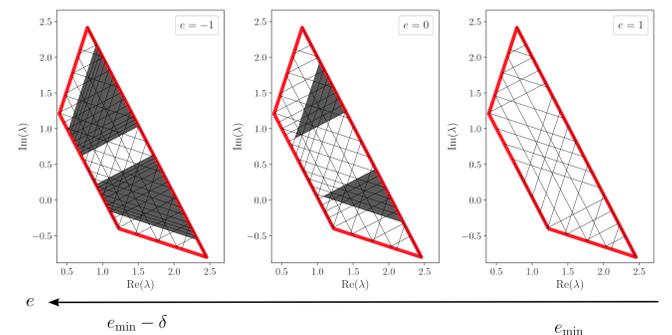
2. Piecewise constant with polygonal pieces



3. But with an **accumulation phenomenon** that generates an **infinite number of pieces**



Manage the accumulation phenomenon



Introduction of a parameter δ that controls the number of polygonal pieces and their distance to the accumulation lines

Construction of CROQuant

Algorithm 1 CROQuant: Complex Rank-One Quantization Algorithm

Input: $x \in \mathbb{R}^m, y \in \mathbb{R}^n, t \geq 1, \delta \in \mathbb{N}$
 1: Initialize $\hat{\lambda} \leftarrow 1$
 2: Build the tiling domain Ω
 3: Build \mathcal{R}_δ , the set of polygonal pieces inside Ω associated with δ
 4: Build Λ_δ , the set of centroids from the polygonal pieces \mathcal{R}_δ
 5: **for** $\lambda \in \Lambda_\delta$ **do**
 6: **if** $f(\lambda) < f(\hat{\lambda})$ **then** $\hat{\lambda} \leftarrow \lambda$
 7: **return** $\text{round}(\hat{\lambda}x), \text{round}(\mu(\hat{\lambda})y)$

Application to butterfly Factorizations

$$Z = \begin{bmatrix} \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots \end{bmatrix} \approx \begin{bmatrix} B_1 & & & & \\ & B_2 & & & \\ & & \dots & & \\ & & & \dots & \\ & & & & B_L \end{bmatrix}$$

Butterfly Quantization problem

$$\min_{\hat{B}_1, \dots, \hat{B}_L \in \mathbb{C}\mathbb{F}_t^{n \times n}} \left\| B_1 \cdots B_L - \hat{B}_1 \cdots \hat{B}_L \right\|_F^2$$

Key property

$$B_1 \cdots B_j \underbrace{B_{j+1} \cdots B_k}_X \underbrace{B_{k+1} \cdots B_l}_{Y^H} B_{l+1} \cdots B_L \rightarrow \text{use } n \text{ times CROQuant}$$

$$XY^H = \sum_{i=1}^n x_i y_i^H \text{ where } x_i y_i^H \text{ have disjoint supports [2]}$$

For $L > 2$: need **heuristics** to decide how to **order/group** the factors

Left-To-Right (LTR): writing $B_1(B_2(\cdots(B_{L-1}B_L)\cdots))$
 Pairwise : writing $(B_1B_2)(B_3B_4)\cdots(B_{L-1}B_L)$

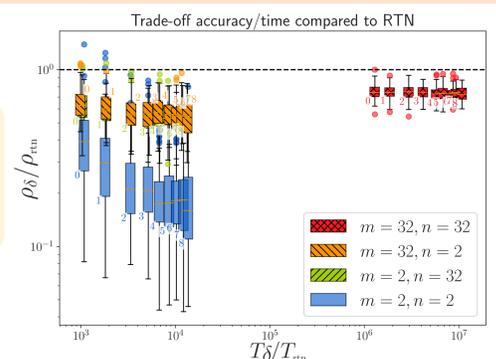
Experiments

On random rank-one matrices

$$\text{Re}(x), \text{Im}(x) \sim \mathcal{U}_{[0,1]}$$

$$\text{Re}(y), \text{Im}(y) \sim \mathcal{U}_{[0,1]}$$

$$\rho := \frac{\|xy - \hat{x}\hat{y}\|_F^2}{\|xy\|_F^2}$$



1. CROQuant **improves accuracy** compared to RTN when $\delta \geq 2$
2. As m and n increase, increasing δ is less worthwhile

On random butterfly matrices

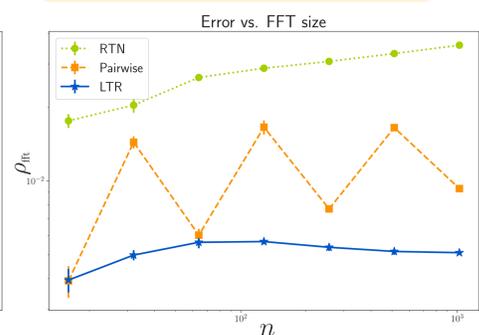
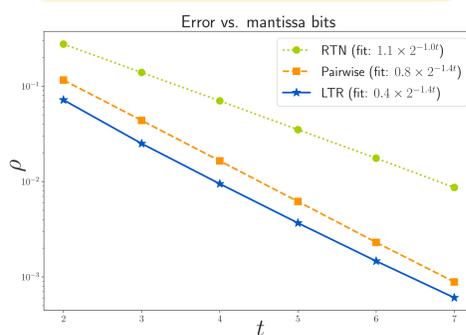
$$B_1 \cdots B_L \sim \mathcal{CN}(0, I_n) \quad \rho := \frac{\|B_1 \cdots B_L - \hat{B}_1 \cdots \hat{B}_L\|_F^2}{\|B_1 \cdots B_L\|_F^2}$$

On the Fast Fourier transform

$$x \sim \mathcal{N}(0, I_n)$$

$$y = \mathcal{F}x$$

$$\rho_{\text{fit}} := \frac{\|y - \hat{y}\|_2^2}{\|y\|_2^2}$$



→ Reduction of **30% bits** compared to RTN

→ LTR and Pairwise **outperform** RTN on the FFT

