

# Matrix Compression via Randomized Low Rank and Low Precision Factorization

**Rajarshi Saha**, Varun Srivastava, Mert Pilanci

**Stanford** | ENGINEERING  
Electrical Engineering



Dec 10 - Dec 16, 2023  
New Orleans

# Matrix Compression

- Matrices are ubiquitous – can involve billions of elements making their storage and processing quite demanding in terms of computational resources and memory usage.

Eg.: Vector databases, Kernel matrices, LLM weight matrices, etc.

# Matrix Compression

- Matrices are ubiquitous – can involve billions of elements making their storage and processing quite demanding in terms of computational resources and memory usage.  
Eg.: Vector databases, Kernel matrices, LLM weight matrices, etc.
- Several real-world matrices exhibit **approximately low-rank structure** due to inherent redundancy or patterns. [\[Udell & Townsend, 2019\]](#)

# Matrix Compression

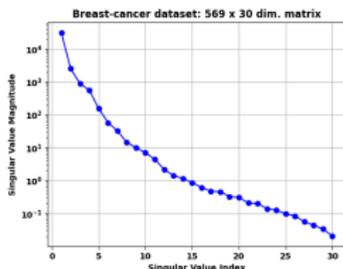
- Matrices are ubiquitous – can involve billions of elements making their storage and processing quite demanding in terms of computational resources and memory usage.

Eg.: Vector databases, Kernel matrices, LLM weight matrices, etc.

- Several real-world matrices exhibit **approximately low-rank structure** due to inherent redundancy or patterns. [Udell & Townsend, 2019]
- **Singular value decomposition**: Any matrix  $\mathbf{A} \in \mathbb{R}^{n \times d}$  can be written as:

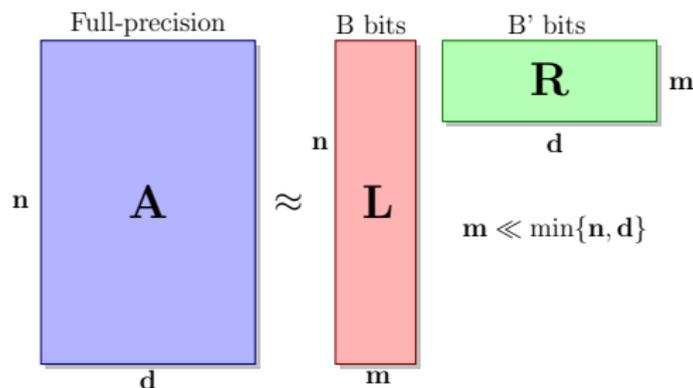
$$\mathbf{A} = \sum_{i=1}^{\text{rank}(\mathbf{A})} \sigma_i \mathbf{u}_i \mathbf{v}_i^{\top},$$

where  $\{\sigma_i\}$  are the singular values, and  $\mathbf{u}_i \in \mathbb{R}^n$ ,  $\mathbf{v}_i \in \mathbb{R}^d$  are singular vectors.



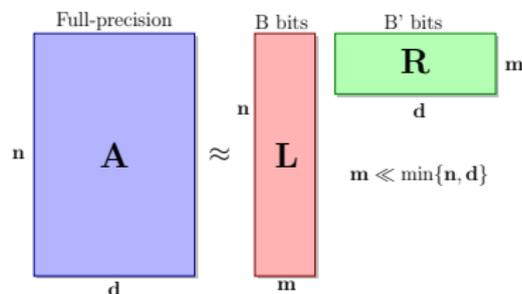
- Question: How to compress matrices via **dimensionality reduction** and **quantization**?
- Our solution uses **Randomized Embeddings**.

# Matrix Compression: Low-precision and Low-Rank



- We obtain a **randomized factorization**,  $\mathbf{A} \approx \mathbf{L}\mathbf{R}$ , where the entries of left factor ( $\mathbf{L}$ ) and right factor ( $\mathbf{R}$ ) are quantized with  $B$  and  $B'$  bits per entry respectively.
- Total bit requirement is  $mnB + mdB'$ .
- By tuning sketch-size  $m$  we can ensure compression while letting  $B$  and  $B'$  to take values allowed by current hardware-primitives, e.g., 4-bits, 8-bits, etc.
- Low-precision computations also have **low latency**: Computing  $\mathbf{A}\mathbf{x}$  versus  $\mathbf{L}(\mathbf{R}\mathbf{x})$ .

# Matrix Compression: Low-precision and Low-Rank



Our **LPLR** algorithm:

- Computes randomized rangefinder  $\mathbf{AS}$ , and quantizing it.
- Computes approximate projection of the columns of  $\mathbf{A}$  onto this quantized basis.

Algorithm 1: **LPLR**: Randomized Low-Precision Low-Rank factorization

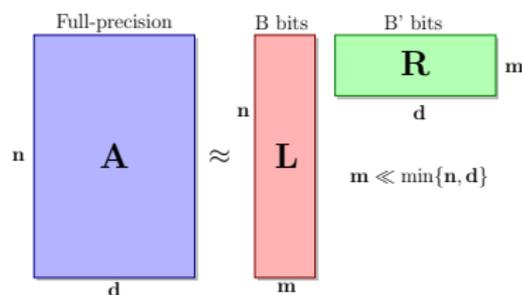
---

**Input** : Matrix  $\mathbf{A} \in \mathbb{R}^{n \times d}$ , sketch size  $m$ , Quantizers  $Q, Q'$  with dynamic ranges  $R_Q, R_{Q'}$  and bit-budgets  $B, B'$  respectively.

**Output** : Factorization:  $\mathbf{LR}$  where  $\mathbf{L} \in \mathbb{R}^{n \times m}, \mathbf{R} \in \mathbb{R}^{m \times d}$

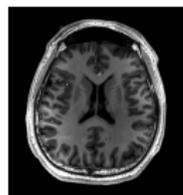
- 1 Sample a Gaussian sketching matrix  $\mathbf{S} \in \mathbb{R}^{d \times m}$  with entries  $S_{ij} \sim \mathcal{N}(0, \frac{1}{m})$ .
  - 2 Compute an approximate basis of column space of  $\mathbf{A}$  by forming the sketch:  $\mathbf{AS}$ .
  - 3 Quantize the approximate basis with  $Q$  to get  $Q(\mathbf{AS})$ .
  - 4 Find  $\mathbf{W}^* = \arg \min_{\mathbf{W}} \|Q(\mathbf{AS})\mathbf{W} - \mathbf{A}\|_{\mathbb{F}}^2$ .
  - 5 Quantize  $\mathbf{W}^*$  using quantizer  $Q'$  to get  $Q'(\mathbf{W}^*)$ .
  - 6 **return** Low-rank and low-precision approximation  $\mathbf{LR}$  where  $\mathbf{L} = Q(\mathbf{AS}), \mathbf{R} = Q'(\mathbf{W}^*)$ .
-

# Matrix Compression: Low-precision and Low-Rank

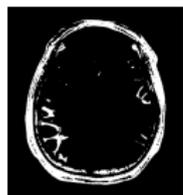


- We obtain a factorization,  $\mathbf{A} \approx \mathbf{L}\mathbf{R}$ , where the entries of  $\mathbf{L}$  and  $\mathbf{R}$  are quantized with  $B$  and  $B'$  bits per entry respectively.
- (A popular benchmark) **Naive quantization**: Quantize each entry of  $\mathbf{A} \in \mathbb{R}^{n \times d}$  uniformly with a  $B_{\text{nq}}$  – bit quantizer.
- Compression ratio with respect to naive quantization is  $\frac{m n B + m d B'}{n d B_{\text{nq}}}$ .
- By tuning sketch-size  $m$  we can ensure compression ratio  $\leq 1$  for  $B_{\text{nq}} = 1$ , while letting  $B$  and  $B'$  to take values allowed by current hardware-primitives, e.g., 4-bits, 8-bits, etc.
- **Direct-SVD** quant. benchmark: Compute the best rank- $k$  approximation  $(\mathbf{U}\Sigma)_k \mathbf{V}_k^T$  by retaining the top- $k$  singular vectors, and subsequently quantize:  $\mathbf{A} \approx \mathbf{Q}((\mathbf{U}\Sigma)_k) \mathbf{Q}'(\mathbf{V}_k^T)$ .

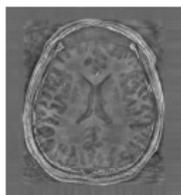
# Image compression



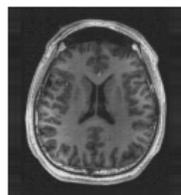
Original



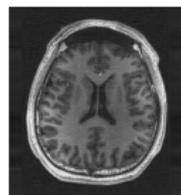
Naive



DSVD



LPLR (ours)



LSVD (ours)

Compressing a brain MRI.  $B = 4$ ,  $B' = 8$ ,  $B_{nq} = 1$ ,  $m = 124$ ,  $n = 1534$ ,  $d = 1433$



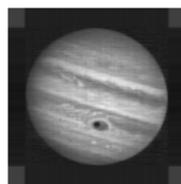
Original



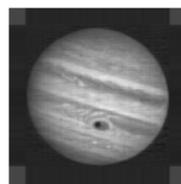
Naive



DSVD



LPLR (ours)



LSVD (ours)

Compression of a Jupiter image showing its Great Red Spot and Ganymede's shadow (NASA/ESA Hubble Space Telescope).  $B = 2$ ,  $B' = 8$ ,  $B_{nq} = 1$ ,  $m = 110$ . Orig. image dim.:  $1102 \times 1102$

# Approximate nearest neighbor search

- For a given data matrix  $\mathbf{A} \in \mathbb{R}^{n \times d}$  and a query  $\mathbf{x} \in \mathbb{R}^d$ , retrieve

$$i^* = \operatorname{argmax}_{i \in [n]} (\mathbf{A}\mathbf{x})_i \approx \operatorname{argmax}_{i \in [n]} (\mathbf{L}\mathbf{R}\mathbf{x})_i$$

- Applications: Semantic search over vector databases (music recommendation), In-context learning for LLMs, etc.

Table 5: CIFAR100 embeddings generated by MobileNetV3 with an unquantized accuracy and F1 score 76%: Results on LPLR and LPLR-SVD with  $B = B' = 8$  bits

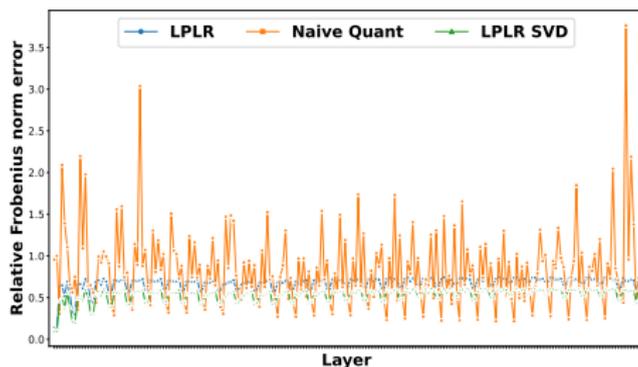
$B_{\text{nq}}$	Frobenius Norm Error				Accuracy (%)				Weighted F1 Score (%)			
	LPLR	LSVD	DSVD	NQ	LPLR	LSVD	DSVD	NQ	LPLR	LSVD	DSVD	NQ
1	<b>1.04</b>	1.08	1.09	6.75	79	<b>82</b>	<b>82</b>	1	79	<b>82</b>	<b>82</b>	0
2	<b>1.08</b>	1.1	1.12	2.18	<b>80</b>	<b>80</b>	<b>80</b>	1.7	<b>80</b>	<b>80</b>	<b>80</b>	1.3
4	<b>1.11</b>	1.12	1.14	1.17	<b>79</b>	78	77	75	<b>79</b>	78	78	75

Table 6: IMDB embeddings generated by BERT with an unquantized accuracy and F1 score 75% and 74% respectively: Results on LPLR and LPLR-SVD with  $B = B' = 8$  bits

$B_{\text{nq}}$	Frobenius Norm Error				Accuracy (%)				Weighted F1 Score (%)			
	LPLR	LSVD	DSVD	NQ	LPLR	LSVD	DSVD	NQ	LPLR	LSVD	DSVD	NQ
1	0.313	<b>0.241</b>	0.229	6.63	73	74	<b>75</b>	50	74	74	<b>75</b>	33
2	0.235	0.178	<b>0.161</b>	1.016	<b>74</b>	<b>74</b>	<b>74</b>	50	<b>74</b>	<b>74</b>	<b>74</b>	50
4	0.148	0.122	<b>0.098</b>	0.417	<b>75</b>	74	<b>75</b>	73	74	74	<b>75</b>	73

# Compressing weight matrices of LLMs

- **LLaMa 7b** [Touvron et. al, 2023]: An LLM with several layers (difficult to deploy on GPUs)



Comparison of LPLR and LPLR-SVD on LLaMa weight matrices with  $B = B' = 8$  bits,  $B_{nq} = 4$  bits, ordered by the original sequence of layers on the “Layer” – axis. We observe consistently better Frobenius norm error using LPLR and LPLR-SVD, with the exception of specific layers which lend themselves to naive quantization.

B = B' = 8 bits, B <sub>nq</sub> = 4 bits			
Metric	LPLR	LPLR-SVD	Naive Quant.
Mean	0.672	<b>0.537</b>	0.836
Std Dev	0.080	<b>0.079</b>	0.470

Average relative Frobenius norm error on LLaMa weight matrices

# Theoretical analysis

- Approximation error upper bounds on the Frobenius norm  $\|\mathbf{LR} - \mathbf{A}\|_F^2$ .
- Bit requirement: How many bits are required per matrix coordinate to achieve the corresponding approximation error?
- Computation requirement: No. of floating point multiplications of the rate determining step.  $O(ndm)$  for LPLR vs.  $O(nd^2)$  for direct-SVD quant.
- Properties of randomized embeddings useful for LPLR factorization:
  - **Subspace approximation:** For approximately low-rank matrices  $\mathbf{A} \in \mathbb{R}^{n \times d}$ , randomly sketching the columns, i.e.,  $\mathbf{AS} \in \mathbb{R}^{n \times m}$  constitutes a basis for  $\text{range}(\mathbf{A})$  with high probability.  
[Halko et. al, 2011; Witten & Candes, 2015; Tropp et. al, 2017, ...]
  - **Democratic equalization:**  $\|\mathbf{AS}\|_{\max} \triangleq \max_{i,j} |A_{ij}|$  is "small" with high probability.  
[Charikar, 2002; Boufounos & Baraniuk, 2008; Plan & Vershynin, 2014, Lyubarskii & Vershynin, 2006; Studer et. al. 2015, ...]

Details in paper

# Conclusions

1. Randomized-embedding based matrix compression: Low precision and low rank representations.
2. Computationally efficient:  $O(ndm)$ ,  $m \ll \min\{n, d\}$  for randomized-embedding based LPLR vs.  $O(nd^2)$  for direct-SVD based methods.
3. Sketch size  $m$  is a tunable knob. Allows flexible compression ratios that achieve parity with (aggressive) quantization as low as a *single bit* (using current hardware primitives).
4. Provably better approximation error guarantees (Details in paper).
5. Applications in compressing datasets, neural network weights, approximate nearest neighbors, etc.

# Thank you!

Reach out for questions or discussions:

[rajsaha@stanford.edu](mailto:rajsaha@stanford.edu)

## Poster Session:

Tue 12 Dec 10:45 a.m. CST — 12:45 p.m. CST, Great Hall & Hall B1+B2 #1824

<https://neurips.cc/virtual/2023/poster/70291>

Paper: <https://openreview.net/forum?id=rxsCTtkqA9>

GitHub: <https://github.com/pilancilab/matrix-compressor>