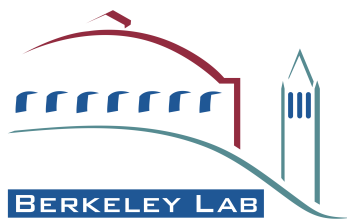# Efficient Leverage Score sampling for Tensor Train Decomposition

Vivek Bharadwaj[*,1]
Beheshteh T. Rakhshan[*,2]
Osman Asif Malik[3]
Guillaume Rabusseau[2,4]

[1] UC Berkeley
[2] Mila, Udem
[3] Lawrence National Lab
[4] CIFAR

# How to compute Tensor Train Decomposition?

- TT_SVD
- Randomized TT_SVD
- TT_ALS
- Randomized TT_ALS (proposal)

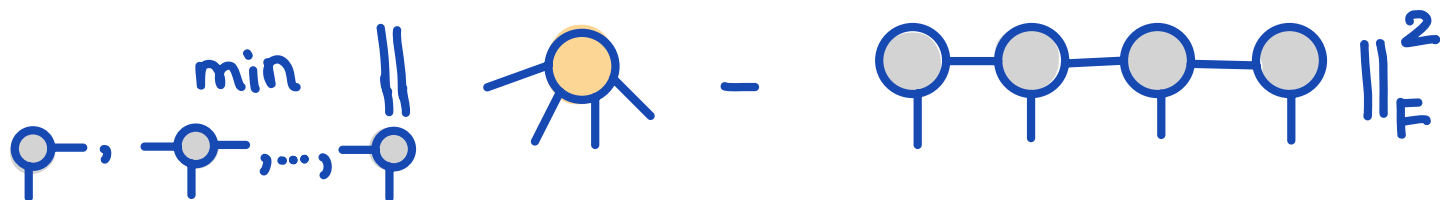# why Tensor Train Decomposition?

- #of parameters $O(NIR^2)$ instead of $O(I^N)$

  order of a $\swarrow$ $\downarrow$ $\hookrightarrow$ rank of decomposition
  Tensor

  $\downarrow$

  dimension size $(I_1 = I_2 = \ldots = I_N = I)$

- Finding a good approximation is feasible.

- Numerically Stable.
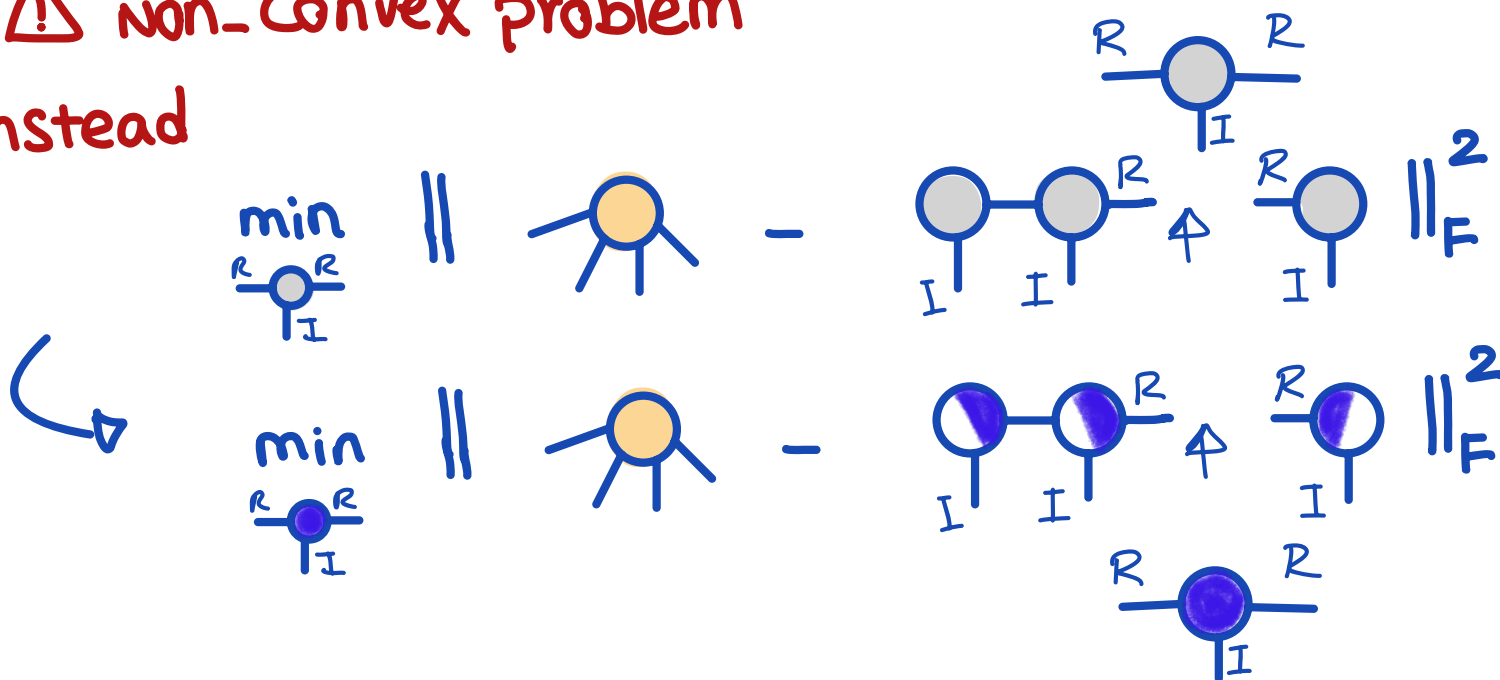
# why Randomized Alternating Least Squares (ALS)?

- Lack of randomized ALS methods to find

  a Tensor Train Decomposition.

# Goal:

$$\min_{\bullet,\ \bullet,\ \ldots,\ \bullet} \left\| \ \bullet \ - \ \bullet\!-\!\bullet\!-\!\bullet\!-\!\bullet \ \right\|_F^2$$

⚠️ **Non-convex problem**

## Instead

$$\min_{R-\bullet-R \atop I} \left\| \ \bullet \ - \ \bullet\!-\!\bullet \overset{R}{\underset{I\ \ I}{\ }} \!\!\Uparrow\ \overset{R}{\underset{I}{\bullet}} \ \right\|_F^2$$

$$\min_{R-\bullet-R \atop I} \left\| \ \bullet \ - \ \bullet\!-\!\bullet \overset{R}{\underset{I'\ \ I}{\ }} \!\!\Uparrow\ \overset{R}{\underset{I}{\bullet}} \ \right\|_F^2$$

## How?

✓ Initialize $\bullet\!-,\ -\!\bullet\!-,\ -\!\bullet\!-,\ -\!\bullet$ randomly.

✓ Update one core at a time until convergence.

$$\min_{A_j} \left\| \left( A^{<j} \otimes A^{>j^T} \right) \left( A_j \right)^T_{(2)} - X^T_{(j)} \right\|^2_F$$

⚠️ Cost $O(I^N)$ to solve Least-Squares

## Randomized Tensor Train ALS (proposal)

- General Randomized Leas-Squares

$$\min_X \| AX - b \|^2 \xrightarrow{\text{instead}} \min \| SAX - Sb \|^2$$

- Randomized Tensor Train Least-Squares

$$\min_{A_j} \left\| S \left( A^{<j} \otimes A^{>j^T} \right) \left( A_j^T \right)_{(2)} - S X^T_{(j)} \right\|^2_F$$

S is a sampling matrix

↳ How to construct S?

📝 with leverage scores ; $\mathbb{P}_i \propto A[i,:](A^T A)^+ A[i,:]^T$

$\underbrace{\phantom{(A^T A)^+}}$
orthogonal
↝ no cost to
compute

In Tensor Train Case:

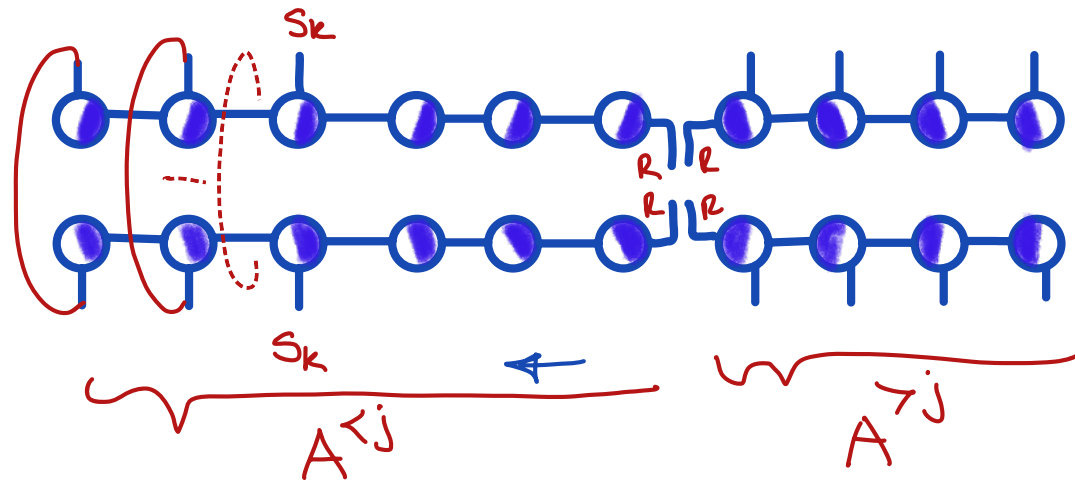Compute $\mathbb{P}_i \propto A^{\neq j}[i^{\neq j},:] A^{\neq j}[i^{\neq j},:]^T$

↳ Still is a challenge.

🌟 Solution

Suppose we have drawn $\hat{S}_{j-1} = S_{j-1}, \ldots, \hat{S}_{k+1} = S_{k+1}$
and now want to draw $\hat{S}_k = S_k$

$$\mathbb{P}(\hat{S}_k = S_k \mid \hat{S}_{j-1} = S_{j-1}, \ldots, \hat{S}_{k+1} = S_{k+1}) = \frac{\mathbb{P}(\hat{S}_k = S_k, \ldots, \hat{S}_{j-1} = S_{j-1})}{\mathbb{P}(\hat{S}_{k+1} = S_{k+1}, \ldots, \hat{S}_{j-1} = S_{j-1})}$$

4

$$\rightsquigarrow \quad \mathbb{P}(\hat{S}_K = S_k \mid \hat{S}_{>k} = S_{>k}) \propto$$

$$\text{Tr}\left( H_{>k}^{\mathsf{T}} A_k[:, S_k, :]^{\mathsf{T}} A[:, S_k, :] H_{>k} \right)$$

where
$$H_{>k} = A_{k+1}[:, S_{k+1}, :] \ldots A_{j-1}[:, S_{j-1}, :]$$

⚠ Updating $H_{>k}$ cost $O(R^3)$

Define
$$q = \frac{1}{R}\left( L[:, 1]^2 + \ldots + L[:, R]^2 \right) \quad \text{where } L = A^{<j}$$

5

✓ Sample a column uniformly, $\hat{t} = t$

✓ Sample a row from $L[:,t]^2$

↳ Reduce the cost to $O(R^2)$

For any $\varepsilon, \delta \in (0,1)$ the sampling procedure above guarantees that with $J = \tilde{O}(R^2/\varepsilon\delta)$ sample per Least-Square problem
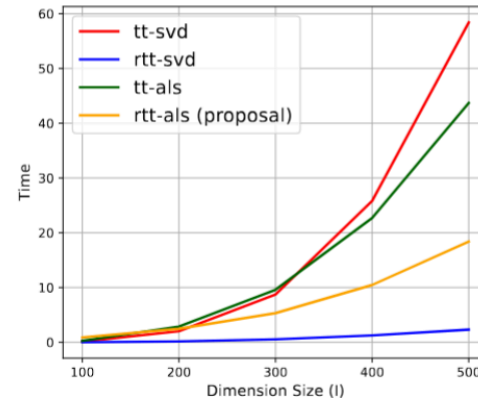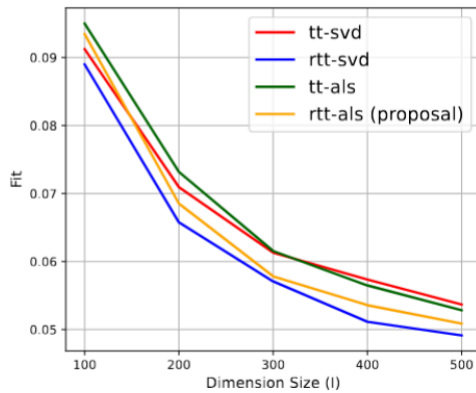
$$\| A^{\neq j}(\tilde{A}_j)^T_{(2)} - X^T_{(j)} \| \leq (1+\varepsilon) \min \| A^{\neq j}(A_j)^T_{(2)} - X^T_{(j)} \|$$

The overall complexity

$$O\left( \frac{\#iter}{\varepsilon\delta} R^4 \sum_{j=1}^{N} N\log I_j + I_j \right)$$

6

## Experiments:

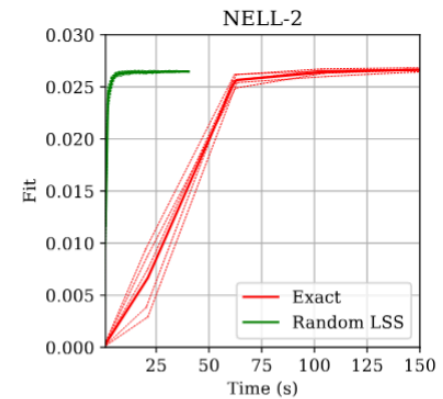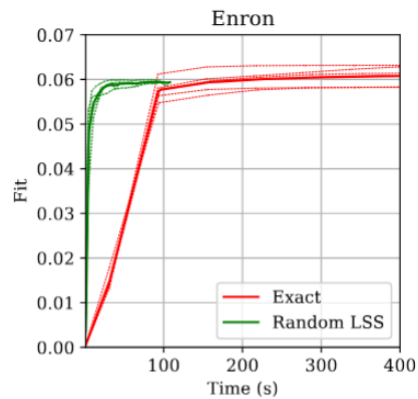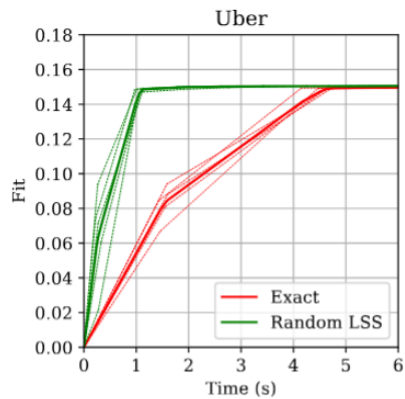### ✓ Synthetic Dense Tensors (J=5000, 5 trials)



### ✓ Real Data (J=2000, target rank $\tilde{R}$=5)

| Method | Pavia Uni. | | Tabby Cat | | MNIST | | DC Mall | |
|---|---|---|---|---|---|---|---|---|
| | Fit | Time | Fit | Time | Fit | Time | Fit | Time |
| TT-ALS | 0.61 | 4.16 | 0.65 | 44.570 | 0.46 | 8.29 | 0.59 | 21.86 |
| **rTT-ALS (proposal)** | 0.60 | 0.82 | 0.65 | 7.360 | 0.45 | 2.20 | 0.59 | 2.81 |
| TT-SVD | 0.61 | 6.65 | 0.65 | 136.189 | 0.46 | 17.19 | 0.59 | 41.45 |
| rTT-SVD | 0.61 | 0.33 | 0.65 | 4.285 | 0.46 | 0.65 | 0.59 | 0.46 |

7

✓ **Sparse Tensors**

Fit vs time  R=6, J=16



| | Uber | | | Enron | | | NELL-2 | | |
|---|---|---|---|---|---|---|---|---|---|
| $R$ | rTT-ALS | TT-ALS | Speedup | rTT-ALS | TT-ALS | Speedup | rTT-ALS | TT-ALS | Speedup |
| 4 | 0.1332 | 0.1334 | 4.0x | 0.0498 | 0.0507 | 17.8x | 0.0213 | 0.0214 | 26.0x |
| 6 | 0.1505 | 0.1510 | 3.5x | 0.0594 | 0.0611 | 12.4x | 0.0265 | 0.0269 | 22.8x |
| 8 | 0.1646 | 0.1654 | 3.0x | 0.0669 | 0.0711 | 10.5x | 0.0311 | 0.0317 | 22.2x |
| 10 | 0.1747 | 0.1760 | 2.4x | 0.0728 | 0.0771 | 8.5x | 0.0350 | 0.0359 | 20.5x |
| 12 | 0.1828 | 0.1846 | 1.5x | 0.0810 | 0.0856 | 7.4x | 0.0382 | 0.0394 | 15.8x |

J=16
40 iters



Fit vs # of samples

8

Thank you !