

Randomized Algorithms for Rounding the Tensor Train Format

Grey Ballard

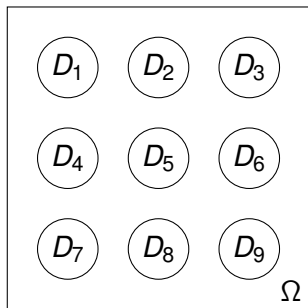
joint with Hussam Al Daas, Paul Cazeaux, Eric Hallman,
Agnieszka Miedlar, Mirjeta Pasha, Tim Reid, and Arvind Saibaba

SIAM Applied Linear Algebra
May 21, 2021



WAKE FOREST
UNIVERSITY

Tensor Train (TT) can make very high dimensional problems tractable



Consider the parameter-dependent PDE:

$$\begin{aligned} -\operatorname{div}(\sigma(x, y; \rho) \nabla(u(x, y; \rho))) &= f(x, y) && \text{in } \Omega, \\ u(x, y; \rho) &= 0 && \text{on } \partial\Omega, \end{aligned}$$

where σ is defined as:

$$\sigma(x, y; \rho) = \begin{cases} 1 + \rho_i & \text{if } (x, y) \in D_i \\ 1 & \text{elsewhere} \end{cases}$$

known as *cookies* problem [Tob12]

- Solving for all parameter values simultaneously, u is 11-D
- With mild assumptions, solution u has low TT ranks
- TT-based iterative linear solver exploits low-rank structure
 - can solve problem for high spatial and parameter resolution

Given a tensor in TT format, often need to compress the ranks

- algebraic operations on TT formats over-extend ranks
- recompression (rank truncation) subject to error threshold
 - or subject to target ranks
- analogous to floating point rounding

Given a tensor in TT format, often need to compress the ranks

- algebraic operations on TT formats over-extend ranks
- recompression (rank truncation) subject to error threshold
 - or subject to target ranks
- analogous to floating point rounding

Low-rank matrix addition example

- consider $\mathbf{A}_1 \mathbf{B}_1^T + \mathbf{A}_2 \mathbf{B}_2^T$, where each factor has r columns
- can represent this in low-rank format $[\mathbf{A}_1 \quad \mathbf{A}_2] [\mathbf{B}_1 \quad \mathbf{B}_2]^T$ which now has rank $2r$
- goal is to compute low-rank approximation with rank $k < 2r$

Randomized algorithm for low-rank approximation of matrix \mathbf{X} :

function $[\mathbf{U}, \mathbf{V}] = \text{RAND-LOW-RANK}(\mathbf{X})$

$$\mathbf{Y} = \mathbf{X}\Omega$$

▷ Ω is random matrix with k columns

$$[\mathbf{U}, \sim] = \text{QR}(\mathbf{Y})$$

▷ (tall-skinny) QR decomposition

$$\mathbf{V}^T = \mathbf{U}^T \mathbf{X}$$

▷ $\mathbf{X} \approx \mathbf{U}\mathbf{V}^T$

Randomized algorithm for low-rank approximation of matrix \mathbf{X} :

function $[\mathbf{U}, \mathbf{V}] = \text{RAND-LOW-RANK}(\mathbf{X})$

$\mathbf{Y} = \mathbf{X}\Omega$ ▷ Ω is random matrix with k columns

$[\mathbf{U}, \sim] = \text{QR}(\mathbf{Y})$ ▷ (tall-skinny) QR decomposition

$\mathbf{V}^T = \mathbf{U}^T \mathbf{X}$ ▷ $\mathbf{X} \approx \mathbf{UV}^T$

Same algorithm tailored for rank- r matrix $\mathbf{X} = \mathbf{AB}^T$ ($r > k$):

function $[\mathbf{U}, \mathbf{V}] = \text{RAND-ROUNDING}(\mathbf{A}, \mathbf{B})$

$\mathbf{Y} = \mathbf{A}(\mathbf{B}^T \Omega)$ ▷ Ω is random matrix with k columns

$[\mathbf{U}, \sim] = \text{QR}(\mathbf{Y})$ ▷ (tall-skinny) QR decomposition

$\mathbf{V}^T = (\mathbf{U}^T \mathbf{A})\mathbf{B}^T$ ▷ $\mathbf{AB}^T \approx \mathbf{UV}^T$

Matrix Rounding

Randomized algorithm for rank- r matrix $\mathbf{X} = \mathbf{A}\mathbf{B}^T$:

function $[\mathbf{U}, \mathbf{V}] = \text{RAND-ROUNDING}(\mathbf{A}, \mathbf{B})$

$\mathbf{Y} = \mathbf{A}(\mathbf{B}^T \Omega)$ $\triangleright \Omega$ is random matrix with k columns

$[\mathbf{U}, \sim] = \text{QR}(\mathbf{Y})$ \triangleright (tall-skinny) QR decomposition

$\mathbf{V}^T = (\mathbf{U}^T \mathbf{A})\mathbf{B}^T$ $\triangleright \mathbf{A}\mathbf{B}^T \approx \mathbf{U}\mathbf{V}^T$

Deterministic algorithm (from Hussam's talk):

function $[\mathbf{U}, \mathbf{V}] = \text{DET-ROUNDING}(\mathbf{A}, \mathbf{B})$

$[\mathbf{Q}_A, \mathbf{R}_A] = \text{QR}(\mathbf{A})$ \triangleright (tall-skinny) QR decomposition

$[\mathbf{Q}_B, \mathbf{R}_B] = \text{QR}(\mathbf{B})$ \triangleright (tall-skinny) QR decomposition

$[\hat{\mathbf{U}}_R, \hat{\Sigma}_R, \hat{\mathbf{V}}_R] = \text{TSVD}(\mathbf{R}_A \mathbf{R}_B^T, k)$ \triangleright k th truncated SVD

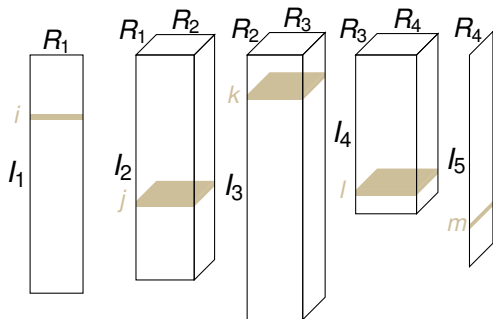
$\mathbf{U} = \mathbf{Q}_A \hat{\mathbf{U}}_R$

$\mathbf{V} = \mathbf{Q}_B (\hat{\mathbf{V}}_R \hat{\Sigma}_R)$ $\triangleright \mathbf{A}\mathbf{B}^T \approx \mathbf{U}\mathbf{V}^T$

Benefits of randomization for matrix rounding

- 1 if \mathbf{A} is $m \times r$, \mathbf{B} is $n \times r$, and they are rounded to rank k , reduces computation from $O((m+n)r^2)$ to $O((m+n)rk)$
- 2 shifts some computational burden from QR to matrix multiplication, which often has higher performance
- 3 opens possibility of choosing Ω for faster multiplication
- 4 enables cheaper rounding of sums of s low-rank matrices:
$$\mathbf{A}_1 \mathbf{B}_1^T + \mathbf{A}_2 \mathbf{B}_2^T + \cdots + \mathbf{A}_s \mathbf{B}_s^T$$
 - sketch of the sum is the sum of the sketches
 - cost of randomized algorithm is linear rather than quadratic in s

Tensor Train (TT) Notation



$$\mathcal{X} \approx \{\mathcal{T}_{X,n}\}, \mathcal{X} \in \mathbb{R}^{l_1 \times l_2 \times l_3 \times l_4 \times l_5}$$

$$\mathcal{T}_{X,n} \in \mathbb{R}^{R_{n-1} \times l_n \times R_n}$$

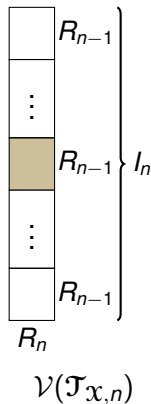
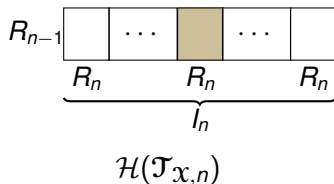
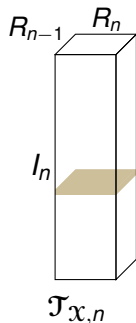
are *TT cores*

$$x_{ijklm} \approx \sum_{\alpha=1}^{R_1} \sum_{\beta=1}^{R_2} \sum_{\gamma=1}^{R_3} \sum_{\delta=1}^{R_4} \mathcal{T}_{X,1}(i, \alpha) \mathcal{T}_{X,2}(\alpha, j, \beta) \mathcal{T}_{X,3}(\beta, k, \gamma) \mathcal{T}_{X,4}(\gamma, l, \delta) \mathcal{T}_{X,5}(\delta, m)$$

Important core unfoldings

$$\mathcal{H}(\mathcal{J}_{x,n}) \in \mathbb{R}^{R_{n-1} \times I_n R_n} \text{ and } \mathcal{V}(\mathcal{J}_{x,n}) \in \mathbb{R}^{R_{n-1} I_n \times R_n}$$

are horizontal and vertical unfoldings of n th core



Deterministic TT Rounding Algorithm [Ose11]

function $\{\mathcal{J}_{\mathbf{z},n}\} = \text{TT-ROUNDING}(\{\mathcal{J}_{\mathbf{x},n}\})$

▷ Orthogonalization Phase

for $n = N$ down to 2 **do**

$$[\mathbf{Y}_n, \mathbf{R}_n] = \text{QR}(\mathcal{H}(\mathcal{J}_{\mathbf{x},n})^T)$$

$$\mathcal{V}(\mathcal{J}_{\mathbf{x},n-1}) = \mathcal{V}(\mathcal{J}_{\mathbf{x},n-1}) \cdot \mathbf{R}_n^T$$

▷ (tall-skinny) QR factorization

▷ Apply \mathbf{R} to previous core

▷ Truncation Phase

$$\mathbf{z} = \mathbf{x}$$

for $n = 1$ to $N - 1$ **do**

$$[\mathbf{Y}_n, \mathbf{R}_n] = \text{QR}(\mathcal{V}(\mathcal{J}_{\mathbf{z},n}))$$

$$[\hat{\mathbf{U}}_R, \hat{\Sigma}, \hat{\mathbf{V}}] \approx \text{TSVD}(\mathbf{R}_n)$$

$$\mathcal{V}(\mathcal{J}_{\mathbf{z},n}) = \text{APPLY-Q}(\mathbf{Y}_n, \hat{\mathbf{U}}_R)$$

$$\mathcal{H}(\mathcal{J}_{\mathbf{z},n+1})^T = \text{APPLY-Q}(\mathbf{Y}_{n+1}, \hat{\mathbf{V}}\hat{\Sigma})$$

▷ (tall-skinny) QR factorization

▷ Truncated SVD of \mathbf{R}

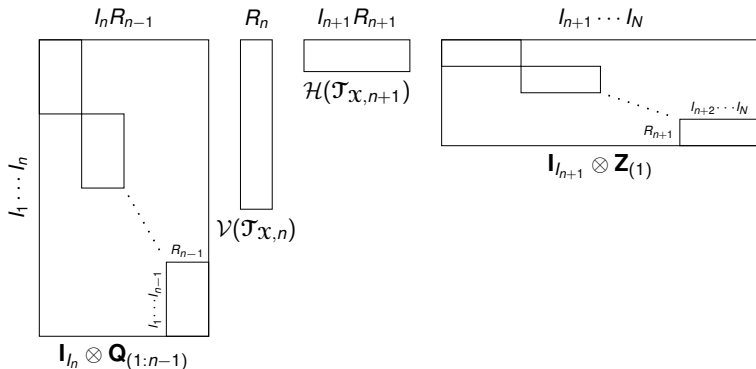
▷ Form explicit $\hat{\mathbf{U}}$

▷ Apply $\hat{\Sigma}\hat{\mathbf{V}}^T$ to next core

More details on TT rounding...

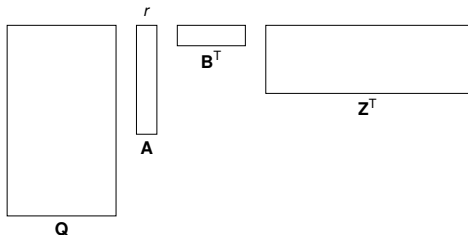
TT rounding does truncated SVDs on $\mathbf{X}_{(1)}$, $\mathbf{X}_{(1:2)}$, $\mathbf{X}_{(1:3)}$, etc., and we have matrix expressions of those unfoldings [ADBB20]

$$\mathbf{X}_{(1:n)} = (\mathbf{I}_{l_n} \otimes \mathbf{Q}_{(1:n-1)}) \cdot \mathcal{V}(\mathcal{J}_{\mathbf{x},n}) \cdot \mathcal{H}(\mathcal{J}_{\mathbf{x},n+1}) \cdot (\mathbf{I}_{l_{n+1}} \otimes \mathbf{Z}_{(1)})$$



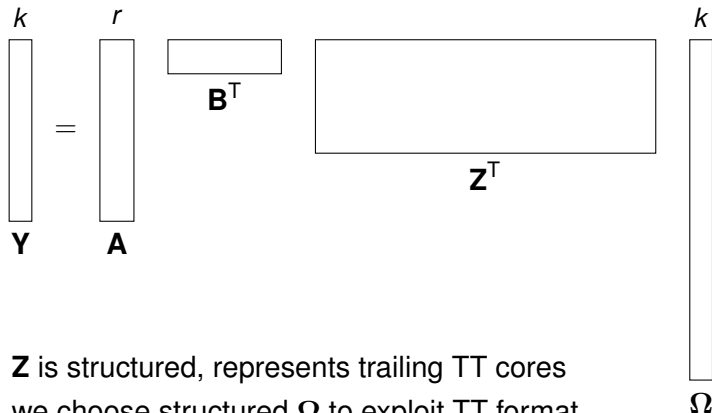
Opportunities for randomized low-rank approximation

Suppose you want to compute the (randomized) truncated SVD of a rank- r matrix $\mathbf{QAB}^T\mathbf{Z}^T$, where \mathbf{A} and \mathbf{B} are tall and skinny with r columns, and



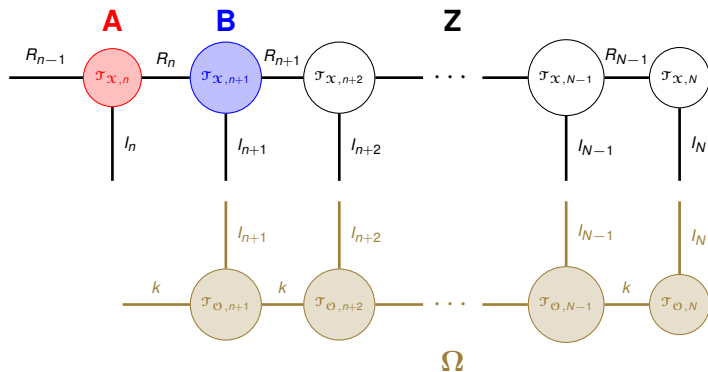
- 1 \mathbf{Q} , \mathbf{B} , and \mathbf{Z} are all column orthonormal
 - use $\text{RAND-LOW-RANK}(\mathbf{A})$
- 2 \mathbf{Q} is column orthonormal but \mathbf{B} and \mathbf{Z} are not
 - use $\text{RAND-ROUNDING}(\mathbf{A}, \mathbf{ZB})$
- 3 none of \mathbf{Q} , \mathbf{B} , and \mathbf{Z} is column orthonormal
 - use $\text{GEN-NYSTRÖM}(\mathbf{QAB}^T\mathbf{Z}^T)$ [Nak20]

Most important computation: $\mathbf{Y} = \mathbf{A}\mathbf{B}^T\mathbf{Z}^T\Omega$



- \mathbf{Z} is structured, represents trailing TT cores
- we choose structured Ω to exploit TT format

TT-like structure of Ω



Tensor network diagram: vertices represent tensors, edges represent modes, connected edges represent contractions

Summary of Algorithms

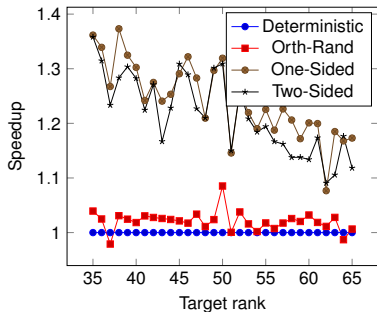
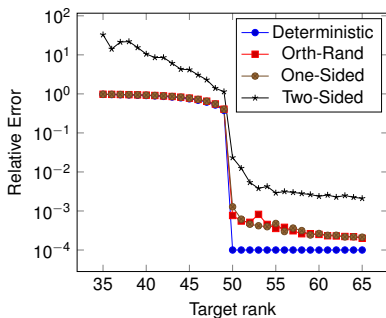
- Deterministic
 - uses orthogonalization phase and truncation phase
- Orthogonalize-then-Randomize (Orth-Rand)
 - uses same orthogonalization phase and randomizes truncation phase with Gaussian projection
 - can use adaptive range-finder algorithm
- Randomize-then-Orthogonalize (One-Sided)
 - avoids TT orthogonalization, uses TT-structured random projection
 - can exploit linearity in sums of TT tensors
- Generalized Nyström (Two-Sided)
 - avoids orthogonalization phase and uses TT-structured random projection on left and right

Experimental results for single synthetic tensor

We round a synthetic TT tensor:

$$\mathbf{y} = \mathbf{x} + 10^{-4} \cdot \mathcal{N},$$

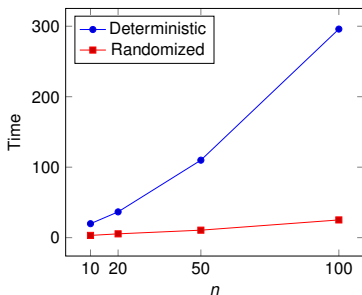
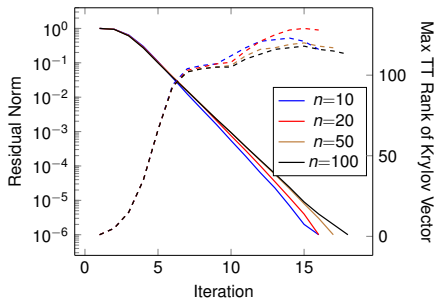
where $\|\mathbf{x}\| = \|\mathcal{N}\| = 1$ and each has 10 modes of dimension 1000 with all TT ranks equal to 50, using target ranks between 35 and 65



Experimental results for cookies problem

We solve the cookies problem with 4 cookies using a tensor with dimension $1781 \times n \times n \times n \times n$ using TT-GMRES

- bottleneck is rounding sum of TT tensors
- deterministic algorithm forms the sum and rounds via orthogonalization
- randomized One-Sided algorithm exploits the sum of sketches, gets same answers



- TT rounding is key operation for TT arithmetic
 - efficient deterministic algorithms exploit low-rank structure
- Randomized algorithms can reduce arithmetic cost and maintain sufficient accuracy
 - benefits depend on ratio of two low ranks
- Randomized approaches yield more benefits for higher-level problems, like rounding sums of TT tensors
 - bottleneck within Krylov solver that exploits TT structure

Thanks for your attention!

`ballard@wfu.edu`



Hussam Al Daas, Grey Ballard, and Peter Benner.
Parallel algorithms for tensor train arithmetic.
Technical Report 2011.06532, arXiv, 2020.



N. Halko, P. Martinsson, and J. Tropp.
Finding structure with randomness: Probabilistic algorithms for
constructing approximate matrix decompositions.
SIAM Review, 53(2):217–288, 2011.



Yuji Nakatsukasa.
Fast and stable randomized low-rank matrix approximation.
Technical Report 2009.11392, arXiv, 2020.



Ivan Oseledets.
Tensor-train decomposition.
SIAM Journal on Scientific Computing, 33(5):2295–2317, 2011.



Christine Tobler.

Low-rank Tensor Methods for Linear Systems and Eigenvalue Problems.

PhD thesis, ETH Zurich, 2012.

Motivation: what if you have to solve many PDEs?

A single PDE simulation can already create a ton of data...
what if we have design/uncertain parameters?

Suppose you have 10 parameters, each with 10 possible values

- now you have to run your simulation 10^{10} times...
- and store all this data...

If the resulting data could be compressed, why not compute the compressed representation from the start?

More details on TT rounding...

TT rounding does SVDs on $\mathbf{X}_{(1)}$, $\mathbf{X}_{(1:2)}$, $\mathbf{X}_{(1:3)}$, etc., so we seek similar matrix expressions of those unfoldings

The unfolding of \mathcal{X} that maps the first n tensor dimensions to rows can be expressed as a product of four matrices:

$$\mathbf{X}_{(1:n)} = (\mathbf{I}_n \otimes \mathbf{Q}_{(1:n-1)}) \cdot \mathcal{V}(\mathcal{X}_{n,n}) \cdot \mathcal{H}(\mathcal{X}_{n,n+1}) \cdot (\mathbf{I}_{n+1} \otimes \mathbf{Z}_{(1)})$$

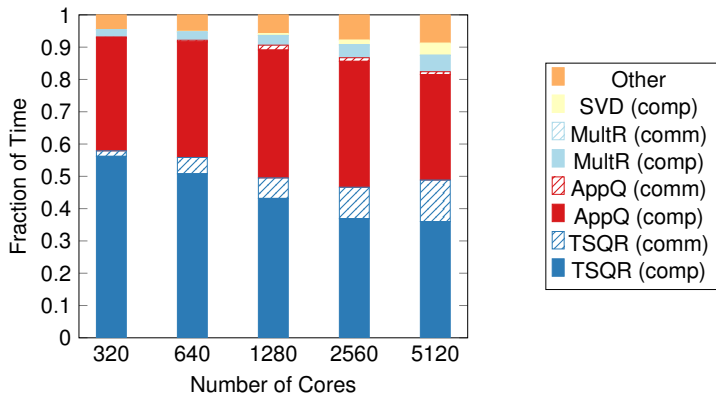
where \mathbf{Q} is $I_1 \times \cdots \times I_{n-1} \times R_{n-1}$ with

$$\mathcal{Q}(i_1, \dots, i_{n-1}, r_{n-1}) = \mathcal{X}_{1,1}(i_1, :) \cdot \mathcal{X}_{2,2}(:, i_2, :) \cdots \mathcal{X}_{n-1,n-1}(:, i_{n-1}, r_{n-1}),$$

and \mathbf{Z} is $R_{n+1} \times I_{n+2} \times \cdots \times I_N$ with

$$\mathcal{Z}(r_{n+1}, i_{n+2}, \dots, i_N) = \mathcal{X}_{n+2,n+2}(r_{n+1}, i_{n+2}, :) \cdot \mathcal{X}_{n+3,n+3}(:, i_{n+3}, :) \cdots \mathcal{X}_{N,N}(:, i_N).$$

Time breakdown of (parallel) TT rounding



- TT tensor: $I_n = 512K$, $R_n = 60 \rightarrow 30$, $N = 50$
- 70-80% of time spent in QR computations