

# Model-driven Sparse CP Decomposition for High-Order Tensors

Jiajia Li, Ioakeim Perros, Jimeng Sun, Richard Vuduc  
Georgia Institute of Technology

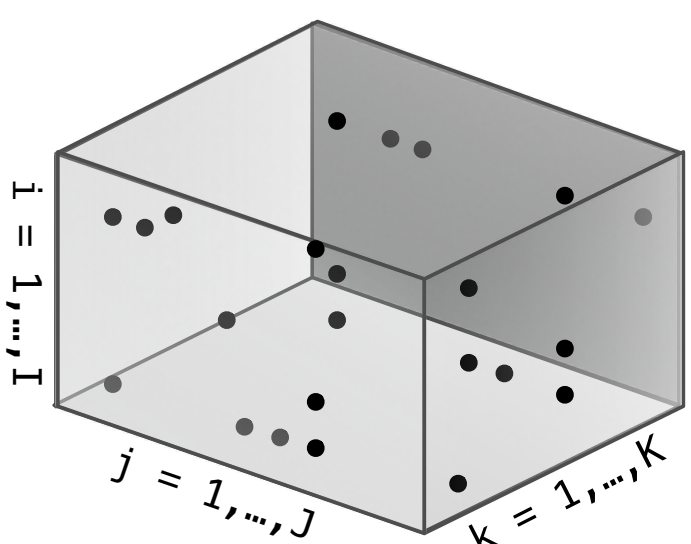


## Background

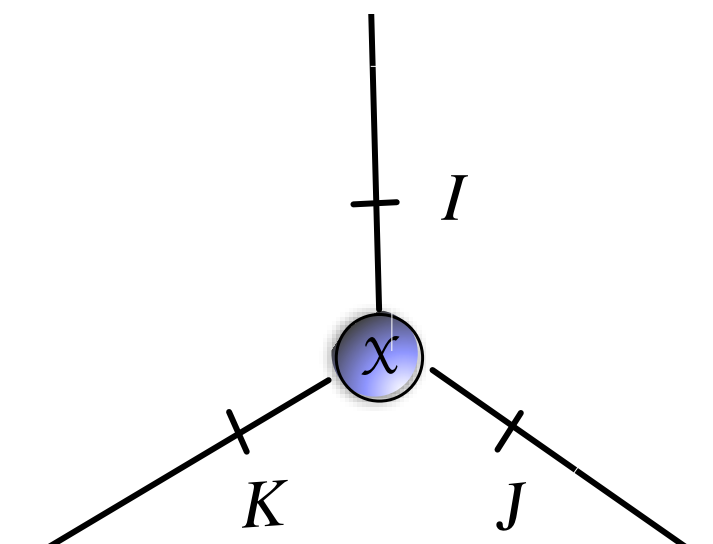
Tensor decomposition is a set of unsupervised methods to analyze and extract knowledge from tensors, which is widely used in healthcare analytics, image processing, machine learning, and social network analytics.

### Sparse tensors

Many real-world tensors are hyper-sparse and have specific features. To discover useful knowledge, efficient sparse algorithms are critical to performance and scalability.



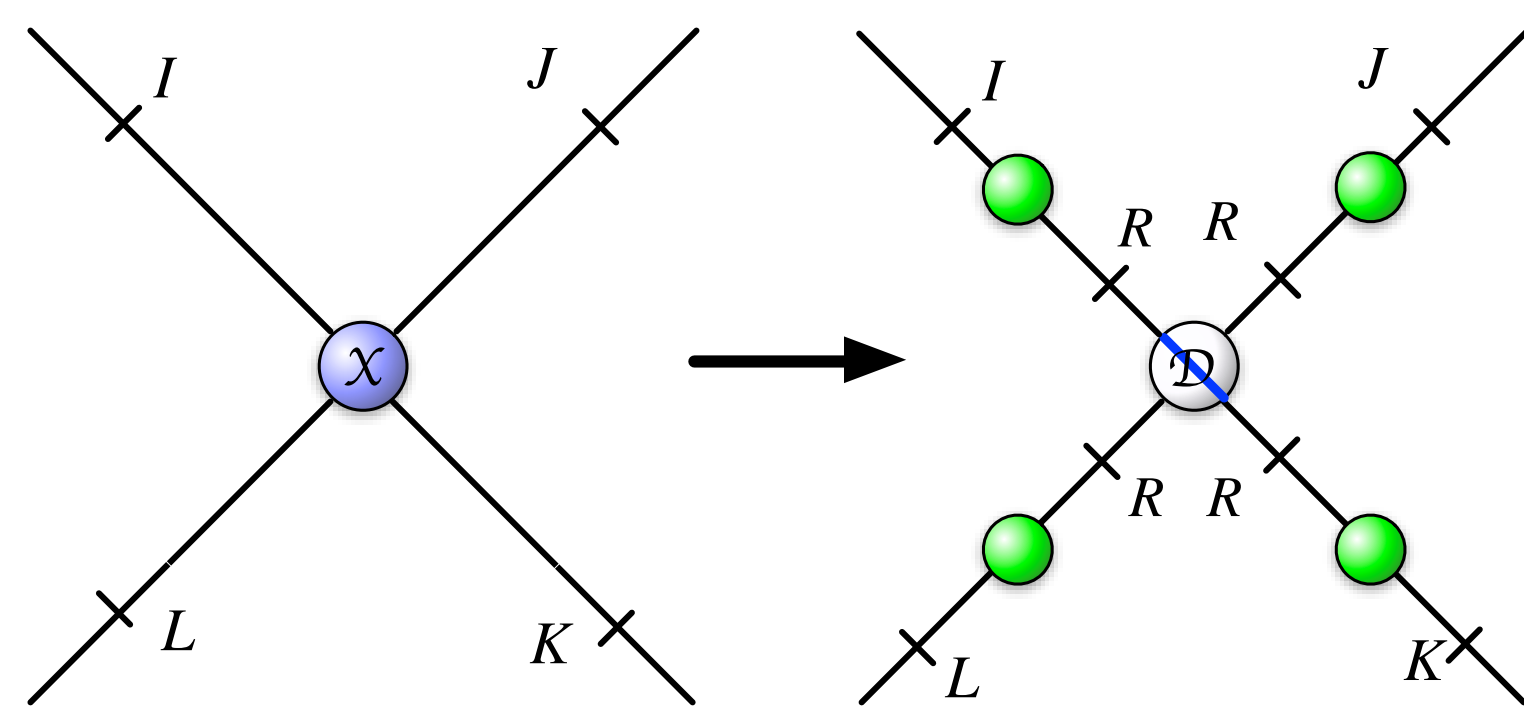
A 3<sup>rd</sup>-order sparse tensor



Tensor network diagram

### CPD

CANDECOMP/PARAFAC decomposition (CPD) approximates an input tensor as a sum of component rank-one tensors, given the number of desired components. CPD is scalable in time and space compared to other low-rank methods.



CPD on a 4<sup>th</sup>-order tensor

## Contributions

- This work proposes an adaptive and efficient algorithm for the computational kernel, Matricized Tensor Times Khatri-Rao Product operation (MTTKRP) of the classical CANDECOMP/PARAFAC decomposition (CPD), by minimizing redundant computations.
- We also build a model-driven procedure to determine the adaptable algorithmic parameters for different input sparse tensors and enable the trade-off between time and space based on the user need and memory resource.
- Our adaptive algorithm shows up to 10x speedup compared to state of the art on real-world high order tensors. Our method also shows near constant scalability with respect to the tensor order, while using acceptable storage space.

## Algorithms

- Decouple an MTTKRP operation into a TTM and a TMHR.

### Algorithm 3 Memorized Full MTTKRP algorithm (memFM).

**Input:** A fourth-order sparse tensor  $\mathcal{X} \in \mathbb{R}^{I \times J \times K \times L}$ , dense factor matrices  $\mathbf{B} \in \mathbb{R}^{J \times R}$ ,  $\mathbf{C} \in \mathbb{R}^{K \times R}$ ,  $\mathbf{D} \in \mathbb{R}^{L \times R}$ , the product order  $p_o$ ;  
**Output:** Updated dense factor matrix  $\tilde{\mathbf{A}} \in \mathbb{R}^{I \times R}$ .

$\tilde{\mathbf{A}} \leftarrow \mathcal{X}_{(1)}(\mathbf{D} \odot \mathbf{C} \odot \mathbf{B})$   
// Compute as  $p_o = [4, 3, 2]$ .  
1: Save  $\mathcal{Y}^{(1)}$ ;  $\mathcal{Y}^{(1)} \leftarrow \mathcal{X} \times_4 \mathbf{D}$  TTM  
2: Save  $\mathcal{Y}^{(2)}$ ;  $\mathcal{Y}^{(2)} = \text{TMHR}(\mathcal{Y}^{(1)}, \mathbf{C})$  TMHR  
3:  $\tilde{\mathbf{A}} = \text{TMHR}(\mathcal{Y}^{(2)}, \mathbf{B})$  return  $\tilde{\mathbf{A}}$ ;

1<sup>st</sup> MTTKRP  $\tilde{\mathbf{A}} \leftarrow \mathcal{X}_{(1)}(\mathbf{D} \odot \mathbf{C} \odot \mathbf{B}) \Leftrightarrow \tilde{\mathbf{A}} \leftarrow \mathcal{X} \begin{matrix} \bullet \\ \mathbf{D} \end{matrix} \begin{matrix} \bullet \\ \mathbf{C} \end{matrix} \begin{matrix} \bullet \\ \mathbf{B} \end{matrix}$

$\mathcal{Y}^{(1)}$   
 $\mathcal{Y}^{(2)}$

## Algorithms

- Design a simple memoization algorithm for MTTKRP sequence.

$$\begin{aligned} \tilde{\mathbf{A}} \leftarrow \mathcal{X}_{(1)}(\mathbf{D} \odot \mathbf{C} \odot \mathbf{B}) &\Leftrightarrow \tilde{\mathbf{A}} \leftarrow \mathcal{X} \begin{matrix} \bullet \\ \mathbf{D} \end{matrix} \begin{matrix} \bullet \\ \mathbf{C} \end{matrix} \begin{matrix} \bullet \\ \mathbf{B} \end{matrix} \\ &\quad \text{Reuse degree: 2} \quad \mathcal{Y}^{(1)} \quad \mathcal{Y}^{(2)} \\ &\Leftrightarrow \tilde{\mathbf{B}} \leftarrow \mathcal{Y}^{(2)} \begin{matrix} \bullet \\ \mathbf{A} \end{matrix} \\ \tilde{\mathbf{C}} \leftarrow \mathcal{X}_{(3)}(\mathbf{D} \odot \tilde{\mathbf{B}} \odot \tilde{\mathbf{A}}) &\Leftrightarrow \tilde{\mathbf{C}} \leftarrow \mathcal{Y}^{(1)} \begin{matrix} \bullet \\ \mathbf{B} \end{matrix} \begin{matrix} \bullet \\ \mathbf{A} \end{matrix} \\ \tilde{\mathbf{D}} \leftarrow \mathcal{X}_{(4)}(\tilde{\mathbf{C}} \odot \tilde{\mathbf{B}} \odot \tilde{\mathbf{A}}) &\Leftrightarrow \tilde{\mathbf{D}} \leftarrow \mathcal{X} \begin{matrix} \bullet \\ \tilde{\mathbf{C}} \end{matrix} \begin{matrix} \bullet \\ \tilde{\mathbf{B}} \end{matrix} \begin{matrix} \bullet \\ \tilde{\mathbf{A}} \end{matrix} \end{aligned}$$

One Group

- Design an **adaptive** memoization algorithm for MTTKRP sequence, to flexibly choose number of memoized MTTKRPs and reuse degrees.

$$\begin{aligned} \tilde{\mathbf{A}} \leftarrow \mathcal{X}_{(1)}(\mathbf{D} \odot \mathbf{C} \odot \mathbf{B}) &\Leftrightarrow \tilde{\mathbf{A}} \leftarrow \mathcal{X} \begin{matrix} \bullet \\ \mathbf{D} \end{matrix} \begin{matrix} \bullet \\ \mathbf{C} \end{matrix} \begin{matrix} \bullet \\ \mathbf{B} \end{matrix} \\ &\quad \text{Reuse degree: 1} \quad \mathcal{Y}^{(2)} \\ &\Leftrightarrow \tilde{\mathbf{B}} \leftarrow \mathcal{X}_{(2)}(\mathbf{D} \odot \mathbf{C} \odot \tilde{\mathbf{A}}) \Leftrightarrow \tilde{\mathbf{B}} \leftarrow \mathcal{Y}^{(2)} \begin{matrix} \bullet \\ \mathbf{A} \end{matrix} \\ \tilde{\mathbf{C}} \leftarrow \mathcal{X}_{(3)}(\tilde{\mathbf{B}} \odot \tilde{\mathbf{A}} \odot \mathbf{D}) &\Leftrightarrow \tilde{\mathbf{C}} \leftarrow \mathcal{X} \begin{matrix} \bullet \\ \tilde{\mathbf{B}} \end{matrix} \begin{matrix} \bullet \\ \mathbf{A} \end{matrix} \begin{matrix} \bullet \\ \mathbf{D} \end{matrix} \\ &\quad \text{Reuse degree: 1} \quad \mathcal{Z}^{(2)} \\ \tilde{\mathbf{D}} \leftarrow \mathcal{X}_{(4)}(\tilde{\mathbf{B}} \odot \tilde{\mathbf{A}} \odot \tilde{\mathbf{C}}) &\Leftrightarrow \tilde{\mathbf{D}} \leftarrow \mathcal{Z}^{(2)} \begin{matrix} \bullet \\ \tilde{\mathbf{C}} \end{matrix} \end{aligned}$$

Group 1

Group 2

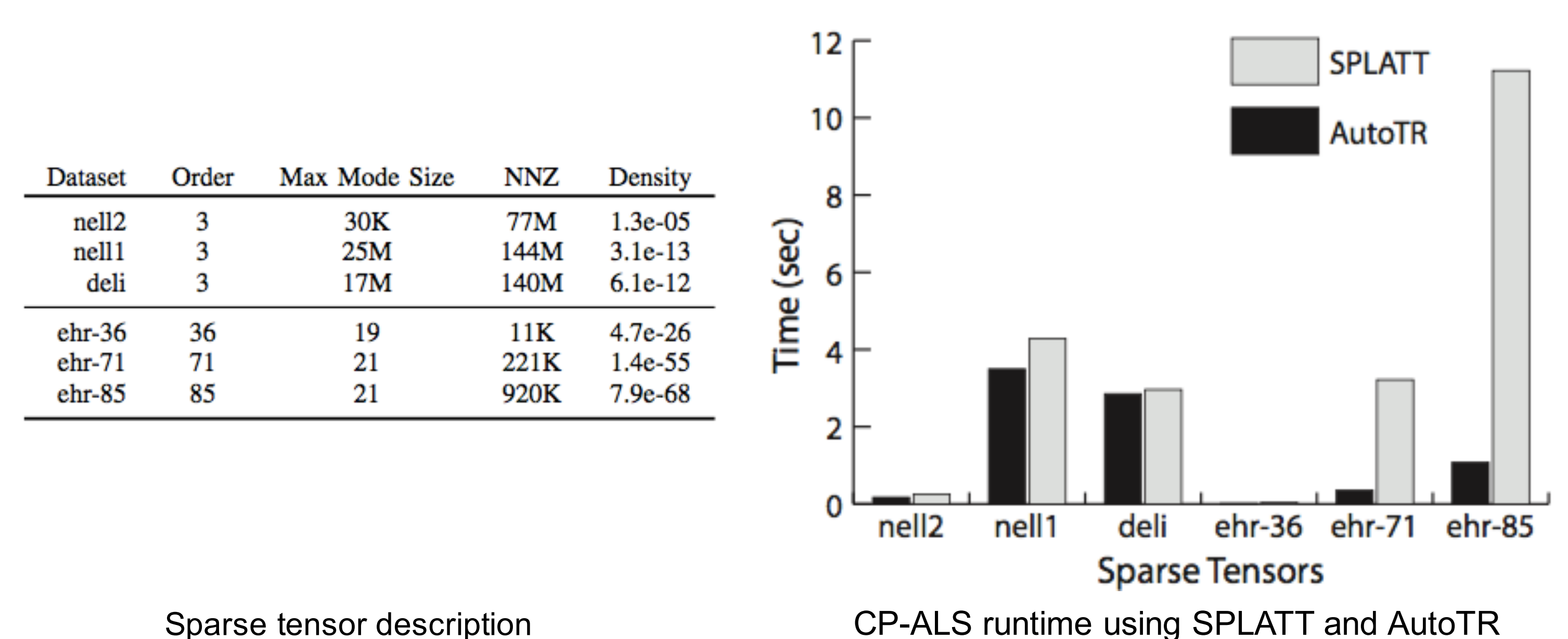
- Decrease the operation number.
- Probably increase the storage space.

Tradeoff

- Build a mode-driven auto-tuner to accurately predict the algorithm parameters for sparse tensors.

## Results

We test our algorithm on Intel Xeon E7-4820 platform using 16 threads.



## Conclusion

Our optimization to increase the reuse of MTTKRP sequence is beneficial for high-order tensors. We use acceptable space to trade for higher performance.

## References

- B. W. Bader, T. G. Kolda et al., "Matlab tensor toolbox (Version 2.6)," Available online, February 2015.
- S. Smith, N. Ravindran, N. Sidiropoulos, and G. Karypis, "Splatt: Efficient and parallel sparse tensor-matrix multiplication," in Proceedings of the 29th IEEE International Parallel & Distributed Processing Symposium, ser. IPDPS, 2015.