





# Model-Driven Sparse CP Decomposition for Higher-Order Tensors

Jiajia Li<sup>1</sup>, Jee Choi<sup>2</sup>, Ioakeim Perros<sup>1</sup>, Jimeng Sun<sup>1</sup>, Richard Vuduc<sup>1</sup>

 $^1$  Computational Science & Engineering, Georgia Institute of Technology, GA, USA  $^2$  IBM Thomas J. Watson Research Center, Yorktown Heights, NY, USA

SIAM AN'17, July 12nd 2017

### The problem



A 4th-Matriced Tensor Times Khatri-Rao Product (MTTKRP) sequence from a tensor decomposition.

## The problem



A 4th-MTTKRP sequence from a tensor decomposition.

#### The problem



A 4th-MTTKRP sequence from a tensor decomposition.

The Problem

#### The problem



A 4th-MTTKRP sequence from a tensor decomposition.

The Problem

## The problem



#### Outline

- Background
- Motivation
- Properties and Formats of Sparse Tensors
- Adaptive Tensor Memoization (ADATM)
- Experiments
- Conclusion

#### Background

#### Tensors

- Tensors, multi-way arrays, provide a natural way to represent multidimensional data.
  - Special cases: matrices (U) 2D tensors, vectors (x) - 1D tensors.
  - Tensor mode (*N*): also called dimension or order.
- A sparse tensor, a tensor consisting mostly of zero entries, widely exist in real applications.
- Tensor analysis is usually factorizing a tensor into interpretable components.
  - $\bullet~$  E.g. CP decomposition, where  $\rm MTTKRP$  is a critical computational kernel.



A 3D CP decomposition on a sparse tensor from healthcare data.

## Basic Tensor Operations



J. Li et.al. (CSE, GaTech)

# Matriced Tensor Times Khatri-Rao Product (MTTKRP)

• Matriced Tensor Times Khatri-Rao Product (MTTKRP)



#### **CP** Decomposition

**Input:** An  $N^{th}$ -order sparse tensor  $\mathbf{X} \in \mathbb{R}^{I \times \cdots \times I}$  and an integer rank  $\mathbb{R}$ : **Output:** Dense factors  $\mathbf{A}^{(1)}, \ldots, \mathbf{A}^{(N)}, \mathbf{A}^{(i)} \in \mathbb{R}^{I \times R}$  and weights  $\lambda$ : 1: Initialize  $\mathbf{A}^{(1)}, \ldots, \mathbf{A}^{(N)}$ ; 2: **do** 3: for  $n = 1, \ldots, N$  do  $\mathbf{V} \leftarrow \mathbf{A}^{(1)\dagger} \mathbf{A}^{(1)} \mathbf{*} \dots \mathbf{A}^{(n-1)\dagger} \mathbf{A}^{(n-1)} \mathbf{*}$ 4  $\mathbf{\Delta}^{(n+1)\dagger}\mathbf{\Delta}^{(n+1)} \ast \ldots \ast \mathbf{\Delta}^{(N)\dagger}\mathbf{\Delta}^{(N)}.$  $\tilde{\mathbf{A}}^{(n)} \leftarrow \mathbf{X}_{(n)}(\mathbf{A}^{(N)} \odot \cdots \odot \mathbf{A}^{(n+1)} \odot \mathbf{A}^{(n-1)} \odot \cdots \odot \mathbf{A}^{(1)}):$ 5:  $\mathbf{\Delta}^{(n)} \leftarrow \widetilde{\mathbf{\Delta}}^{(n)} \mathbf{V}^{\dagger}$ . 6. 7: Normalize columns of  $A^{(n)}$  and store the norms as  $\lambda$ ; 8. end for 9: while Fit ceases to improve or maximum iterations exhausted. 10: Return:  $[\lambda, A^{(1)}, \dots, A^{(N)}]$ :

MTTKRP is the performance bottleneck.

 $T_{CP} \approx N(N^{\epsilon}mR + NIR^2) \approx NT_M,$  $IR \ll m,$ where  $T_M = \mathcal{O}(N^{\epsilon}mR), \epsilon \in (0, 1]$  is the time for a single MTTKRP. Motivation

#### Motivation

$$\tilde{\mathbf{A}} \leftarrow \mathbf{X}_{(1)} \left( \mathbf{D} \odot \mathbf{C} \odot \mathbf{B} \right) \Leftrightarrow \begin{cases} \underline{\mathbf{Y}}^{(1)} = \underline{\mathbf{X}} \times_{4} \mathbf{D}; \\ \underline{\mathbf{Y}}^{(2)} = \underline{\mathbf{Y}}^{(1)} \diamond_{3} \mathbf{C}; \\ \tilde{\mathbf{A}} = \underline{\mathbf{Y}}^{(2)} \diamond_{2} \mathbf{B}; \end{cases}$$

$$\tilde{\mathbf{B}} \leftarrow \mathbf{X}_{(2)} \left( \mathbf{D} \odot \mathbf{C} \odot \tilde{\mathbf{A}} \right) \Leftrightarrow \begin{cases} \underline{\mathbf{Y}}^{(1)} = \underline{\mathbf{X}} \times_{4} \mathbf{D}; \\ \underline{\mathbf{Y}}^{(2)} = \underline{\mathbf{Y}}^{(1)} \diamond_{3} \mathbf{C}; \\ \tilde{\mathbf{B}} = \underline{\mathbf{Y}}^{(2)} \diamond_{1} \tilde{\mathbf{A}}; \end{cases}$$

An  $\operatorname{Mttkrp}$  sequence has arithmetic redundancy.

Motivation

#### Motivation

$$\begin{split} \tilde{\mathbf{A}} &\leftarrow \mathbf{X}_{(1)} \left( \mathbf{D} \odot \mathbf{C} \odot \mathbf{B} \right) \Leftrightarrow \begin{cases} & \underline{\mathbf{Y}}^{(2)} = \underline{\mathbf{X}} \times_4 \mathbf{D}; \\ & \underline{\mathbf{Y}}^{(2)} = \underline{\mathbf{Y}}^{(1)} \diamond_3 \mathbf{C}; \\ & \tilde{\mathbf{A}} = \underline{\mathbf{Y}}^{(2)} \diamond_2 \mathbf{B}; \end{cases} \\ \tilde{\mathbf{B}} &\leftarrow \mathbf{X}_{(2)} \left( \mathbf{D} \odot \mathbf{C} \odot \tilde{\mathbf{A}} \right) \Leftrightarrow \begin{cases} & \underline{\mathbf{Y}}^{(1)} = \underline{\mathbf{X}} \times_4 \mathbf{D}; \\ & \underline{\mathbf{Y}}^{(2)} = \underline{\mathbf{Y}}^{(1)} \diamond_3 \mathbf{C}; \\ & \tilde{\mathbf{B}} = \underline{\mathbf{Y}}^{(2)} \diamond_1 \tilde{\mathbf{A}}; \end{cases} \end{split}$$

( - (1))

An  $\operatorname{Mttkrp}$  sequence has arithmetic redundancy.



The time of an  $\ensuremath{\mathrm{MTTKRP}}$  sequence grows with tensor order.

# Special properties of Sparse $\mathrm{T}\mathrm{T}\mathrm{M}$ and q- $\mathrm{T}\mathrm{T}\mathrm{M}$

#### $\text{Sparse } \mathrm{T}\mathrm{T}\mathrm{M}$

Sparse  $T{\rm\scriptscriptstyle TM}$  outputs a semi-sparse tensor:

- Its product mode becomes dense;
- Its index modes are unchanged.



# Special properties of Sparse $\mathrm{T}\mathrm{T}\mathrm{M}$ and q- $\mathrm{T}\mathrm{T}\mathrm{M}$

#### $\text{Sparse } \mathrm{T}\mathrm{T}\mathrm{M}$

Sparse  $\ensuremath{\mathrm{TTM}}$  outputs a semi-sparse tensor:

- Its product mode becomes dense;
- Its index modes are unchanged.



#### Sparse q- $\mathrm{T}\mathrm{T}\mathrm{M}$

The q- ${\rm TTM}$  of a semi-sparse tensor and a dense matrix yields another semi-sparse tensor:

- Its index modes are unchanged;
- Its product mode disappears.



### **Tensor Formats**



#### vCSF

- The dashed indices are not actually stored, but reuse the indices in CSF tree [Smith et al].
- vCSF is associated with CSF format.

Memoized Intermediate Tensors:



















#### Scheme B



Another scheme to save flops but minimize the amount of storage.

### Schemes Comparison



# Adaptive Tensor Memoization (ADATM)

- ADATM has an analytical performance model, which estimates both the time and storage, given some choice of intermediate tensors.
- $\bullet$  User tells AdaTM how much extra storage is allowed, and  $\rm ADATM$  selects the memoization strategy that approximately minimizes the time predicted by the model.
- For details, check the full paper:
  - "Model-Driven Sparse CP Decomposition for Higher-Order Tensors, 31st IEEE International Parallel & Distributed Processing Symposium (IPDPS17)"

### The model-driven framework



## Platforms and Datasets

| Experimental Plat | forms Configu       | uration      | Sparse ten | sors  |               |      |         |
|-------------------|---------------------|--------------|------------|-------|---------------|------|---------|
|                   | Intel               | Intel        | Dataset    | Order | Ma× Mode size | NNZ  | Density |
| Parameters        | Core i7-4770K       | Xeon E7-4820 | nell2      | 3     | 30K           | 77M  | 1.3e-05 |
| Microarchitecture | Haswell             | Westmere     | nell1      | 3     | 25M           | 144M | 3.1e-13 |
| Frequency         | 3.5 GHz             | 2.0 GHz      | deli       | 3     | 17M           | 140M | 6.1e-12 |
| #Physical cores   | 4                   | 16           | ehr36      | 36    | 19            | 11K  | 4.7e-26 |
| Memory size       | 32 GiB              | 512 GB       | ehr71      | 71    | 21            | 221K | 1.4e-55 |
| Memory bandwidth  | $25.6\mathrm{GB/s}$ | 34.2 GB/s    | ehr85      | 85    | 21            | 920K | 7.9e-68 |
| Compiler          | gcc 4.7.3           | gcc 4.4.7    |            |       |               |      |         |

#### Tensor source:

- Never Ending Language Learning (NELL) project, "nell1, nell2 with noun-verb-noun".
- Data crawled from tagging systems, "deli with user-item-tag".
- Electronic Health Records (EHR) by considering a specific group of similar diseases as one mode and the co-occurrence counts of different diagnoses as values to build the higher-order tensors.

#### Performance



Splatt [Smith et al.] Multi-threaded, using CSF format. Tensor Toolbox [Bader and Kolda] Sequential, using COO format.

nell2 nell1

1000

deli ehr36 ehr71 ehr85

**Sparse Tensors** 

### Storage

|         | Stora | Storage Space (MBytes) |          |      | Ratios |  |  |
|---------|-------|------------------------|----------|------|--------|--|--|
| Dataset | COO   | CSF                    | CSF+vCSF | /CSF | /COO   |  |  |
| nell2   | 2290  | 2540                   | 2581     | 102% | 113%   |  |  |
| nell1   | 4280  | 6430                   | 8510     | 132% | 199%   |  |  |
| deli    | 4180  | 5570                   | 11090    | 199% | 265%   |  |  |
| ehr36   | 3.04  | 1.94                   | 7.97     | 411% | 262%   |  |  |
| ehr71   | 121   | 62                     | 205      | 333% | 169%   |  |  |
| ehr85   | 604   | 200                    | 470      | 236% | 78%    |  |  |

Storage range:

- /CSF: 1-4×;
- /COO: 0.8-2.7×

## Scalability



Comparable multi-threading Scalability



Better scalability in dimensionality.

#### Model Analysis



### **CPD** Application



The speedup of  $\operatorname{AdaTM}$  over  $\operatorname{Splatt}$  on CP-ALS.

### Conclusion

#### Summary

- $\bullet$  We consider the  ${\rm Mttkrp}$  sequence as it arises in the context of CPD.
- We identify a memoization technique that permits a gradual tradeoff of storage for time.
- We parameterize our algorithm and build a model-driven and user-guided framework for it.

#### Future

- Apply our adaptive tensor memoization algorithm to other tensor decompositions;
- We also believe a closer inspection of not just the arithmetic but also communication properties of our method coupled with more architecture-specific tuning are ripe opportunities.

#### Source code: https://github.com/hpcgarage/AdaTM.

#### References

#### References

- B. W. Bader and T. G. Kolda. Efficient MATLAB computations with sparse and factored tensors, SIAM Journal on Scientific Computing 30(1):205-231, December 2007.
- S. Smith, N. Ravindran, N. Sidiropoulos, and G. Karypis, "Splatt: Efficient and parallel sparse tensor-matrix multiplication," IPDPS, 2015.
- O. Kaya and B. Ucar, "Scalable sparse tensor decompositions in distributed memory systems," SC'15. New York, NY, USA: ACM, 2015, pp. 77:1–77:11.
- O. Kaya and B. Uar, "High performance parallel algorithms for the tucker decomposition of sparse tensors," ICPP, Aug 2016, pp. 103–112.
- J. H. Choi and S. Vishwanathan, "Dfacto: Distributed factorization of tensors," NIPS, 2014, pp. 1296–1304.
- M. Baskaran, B. Meister, N. Vasilache, and R. Lethin, "Efficient and scalable computations with sparse tensors," HPEC, Sept 2012, pp. 1–6.
- ... and so on

Backup Slides

#### Two example 4-D tensor memoization algorithms



Comparison:

Storage:  $\underline{\mathbf{Y}}^{(1)} + \underline{\mathbf{Y}}^{(2)}$  vs  $\underline{\mathbf{Y}}^{(2)} + \underline{\mathbf{Z}}^{(2)}$ + permuted  $\underline{\mathbf{X}}$ . – Depend on the input sparse tensor. #Products: 9 vs 8.

#### Performance Analysis of an MTTKRP sequence

**Problem**: Find the number of memoized MTTKRPS  $n_p^*$  that minimizes the total number of products (TTM and q-TTM)  $n_O$  in an N<sup>th</sup>-order MTTKRP sequence, given infinite storage space.

Suppose the input tensor  $\underline{X} \in \mathbb{R}^{I \times \cdots \times I}$  is hypercubical and dense,

#### Lemma

 $n_p^* = \sqrt{N/2}$  minimizes the number of products  $n_O$  for an  $N^{th}$ -order MTTKRP sequence.

$$\begin{cases} n_p = 1, & n_O = N(N+1)/2 = \mathcal{O}(N^2) \\ n_p = N/2, & n_O = N^2/2 = \mathcal{O}(N^2) \\ n_p = n_p^*, & n_O = n_O^* = \mathcal{O}(N^{1.5}) \end{cases}$$

which is asymptotically better for higher-order tensors.

## The model-driven framework



- $\mathbf{n}_{\mathbf{p}}$ : The number of producer modes, with one per memoized MTTKRP. Its range is  $n_{p} \in \{1, \dots, \sqrt{N/2}\}.$
- $\bullet~m_o\colon$  The order of modes of each sparse tensor.
- $n_i$ : The number of intermediate semi-sparse tensors saved from each memoized MTTKRP. Its range is  $\{1, \ldots, N/n_p 1\}$ .
- $s = s(n_p, m_o, n_i)$ : predicated total storage.
- $t = t(n_p, m_o, n_i)$ : predicated total execution time. J. Li et.al. (CSE, GaTech)

#### References

#### Predictive model

$$t = 2\sum_{i=1}^{n_p} \left( \sum_{l=2}^{N} m_l R + \sum_{l=1}^{\frac{N}{n_p}-1} \sum_{j=2}^{l+1} m_j \right) R \triangleq 2\tilde{N}mR; \quad s = \sum_{i=1}^{n_p} \left( m_{CSF}^i + 8\sum_{l=\frac{N}{n_p}-n_i+1}^{\frac{N}{n_p}} m_l R \right)$$

| Algorithms       |                 | #Flops  | Tensor Storage Space (Bytes)   |  |  |
|------------------|-----------------|---|--|--|--|
| Product          | Ттм<br>q-Ттм    | 2mR<br>2mR  | m <sub>CSF</sub><br>8m   |  |  |
| One MTTKRP group | Memoized MTTKRP | $2\sum_{l=2}^{N}m_{l}R$   | $m_{CSF} + 8 \sum_{l=\frac{N}{n_p} - n_i + 1}^{\frac{N}{n_p}} m_l R$   |  |  |
|                  | Partial MTTKRPS | $2\sum_{l=1}^{\frac{n}{p}-1}\sum_{j=2}^{l+1}m_jR$   | -  |  |  |
| MTTKRP sequence  | AdaTM<br>Splatt | $2\sum_{i=1}^{n_{p}} \left( \sum_{l=2}^{N} m_{l}R + \sum_{l=1}^{\frac{N}{n_{p}}-1} \sum_{j=2}^{l+1} m_{j} \right) R$ $2NmR$ | $\left  \sum_{i=1}^{n_p} \left( m_{CSF}^i + 8 \sum_{l=\frac{N}{n_p}-n_i+1}^{\frac{N}{n_p}} m_l R \right) \right. \\ \left. m_{CSF}^i \right  $ |  |  |
|                  |                 | **  | — N  |  |  |

Indices and values use "uint64-t" and "double" respectively.  $m_l$  is the number of fibers at the  $l^{th}$ -level of a CSF tree,  $m_{CSF} = 16 \sum_{l=1}^{N} m_l$ .

J. Li et.al. (CSE, GaTech)

.