

Trillion-Tensor: Trillion-Scale CP Tensor Decomposition

Zeliang Zhang, Xiao-Yang Liu, Pan Zhou

Huazhong University of Science and Technology

Department of Electrical Engineering, Columbia University, USA



Motivations

- Tensor decomposition is the basis of many machine learning applications, including graphic analysis, image classification, data mining, etc. There are two major tensor decomposition models, CP tensor decomposition and Tucker tensor decomposition.
- In recent years, the size of tensors becomes increasingly large, approaching millions to trillions of nonzero elements. Due to the limitation of the main memory, both CP and Tucker tensor decomposition algorithms are impractical for large-scale tensors.
- Many proposed methods may handle tensors with each mode up to millions. However, the input tensors are required to be extremely sparse, where the total number of nonzero elements is up to millions.
- Exploiting the high-performance large scale tensor decomposition scheme.

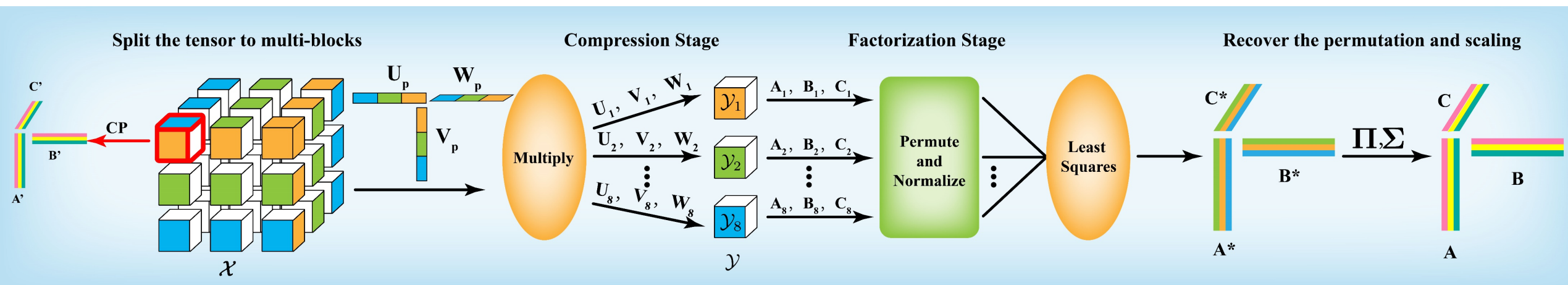


Challenges

- Designing a compression method to process an “out-of-memory” tensor;
- Exploring an appropriate trick to recover the permutations and scaling of the full mode factorization matrices;
- implementing efficient parallel schemes both on CPUs and GPUs, respectively.



Trillion-Tensor Decomposition



Overview of the proposed trillion-tensor decomposition

- The trillion-tensor algorithm mainly consists of three stages, namely compression stage, factorization stage, and a scaling and permutation stage, as given in the above Figure.
- First, a tensor $X \in \mathbb{R}^{300 \times 300 \times 300}$, consisting of 27 tensor blocks $B_n \in \mathbb{R}^{100 \times 100 \times 100}$, $n=1,2, \dots,27$, is first compressed into 8 replicas $Y_p \in \mathbb{R}^{50 \times 50 \times 50}$, $p= 1,2, \dots,8$.
- Second, each replica is factored to mode matrices independently.
- Thirdly, with a proper permutation and normalization, solving a least squares problem for each mode will result in the full mode matrices (A^*,B^*,C^*) . Using the factorization of the first block B_n , we get the permutation matrix Π and the scaling matrix Σ . Then, we use Π and Σ to obtain the full mode matrices (A,B,C)



Optimization

Bottleneck:

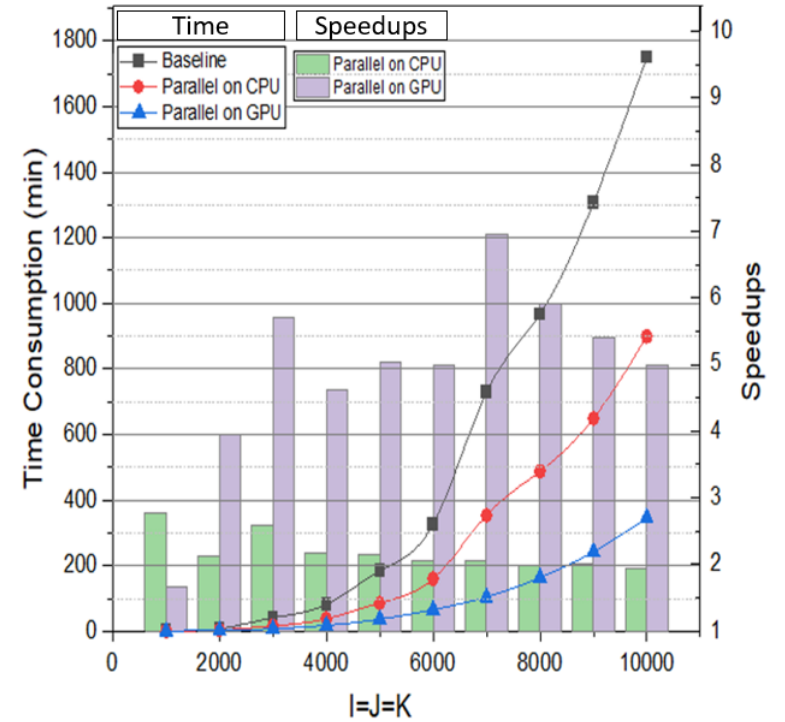
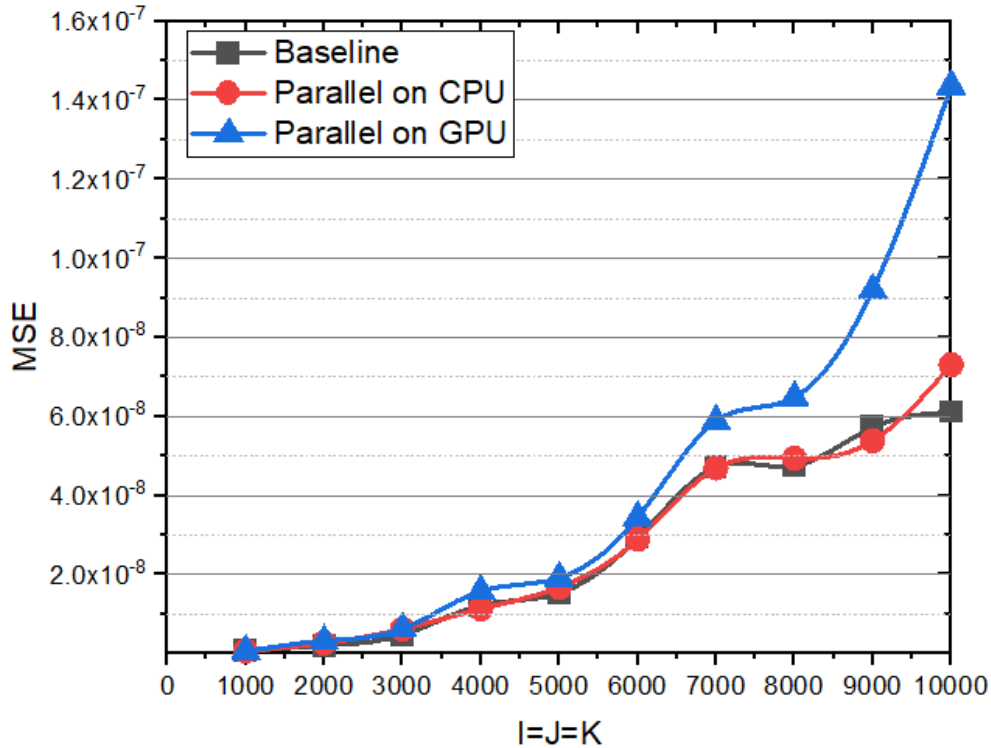
The computational complexity of the compression stage is $O(P \times IJK \times \min(L, M, N))$ versus $O(IJF)$ of the factorization stage. Huge time consumption of the compression stage is the bottleneck of this algorithm with the size of tensor increasing.

Methods:

- Parallel implementation on CPUs: This scheme advocates exploiting the property of multi-core processors on CPUs.
- Parallel implementation on GPUs: This scheme advocates exploiting GPUs' high performance in matrix multiplication.



Performance Evaluations



In our experiments, we test tensors ranging from million-scale to trillion-scale and obtain a relatively low mean squared error. Evaluation results show that trillion-tensor supports up to large scale tensors up to $1w * 1w * 1w$, and a speedup up to $6.95\times$, compared with the baseline on CPU.



THANK YOU!

