# High Performance  Hierarchical Tucker Tensor Learning Using GPU Tensor Cores

Hao Huang[1], Xiao-Yang Liu[2], Weiqin Tong[1], Tao Zhang[1], Anwar Walid[2]

[1]Shanghai University, Shanghai, China;  [2] Columbia University, New York, USA.
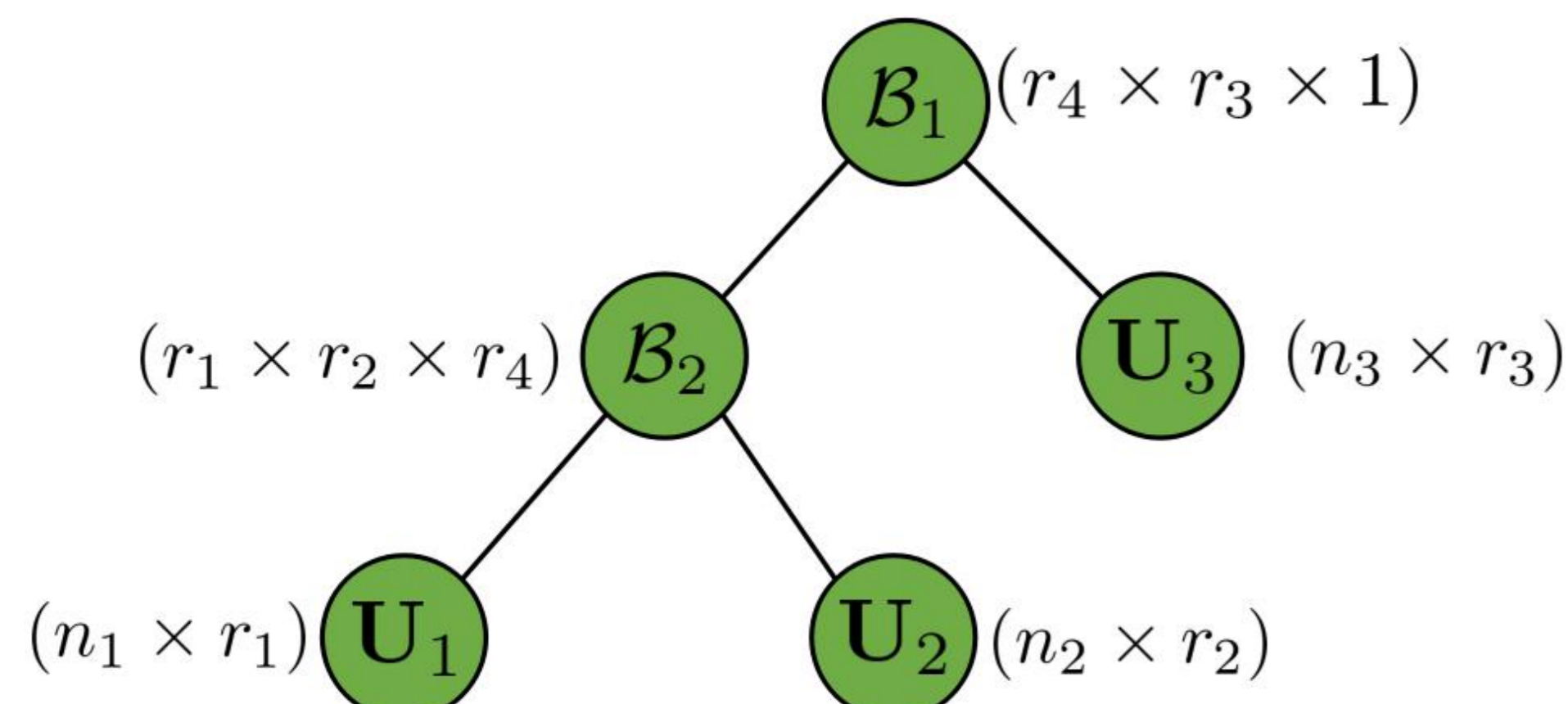
## Introduction

In many real-world applications, data are usually represented as Hierarchical Tucker (HT) tensor model, such as big data analysis, neural network compression and quantum machine learning. However, it is compute intensive due to the time complexity grows exponentially with the order of the data tensor. In this work, we present efficient HT tensor learning primitives using GPU tensor cores.

## Hierarchical Tucker Decomposition

The HT format is stored in the form of a binary tree T, where each branch is a hierarchical division of the tensor mode set.
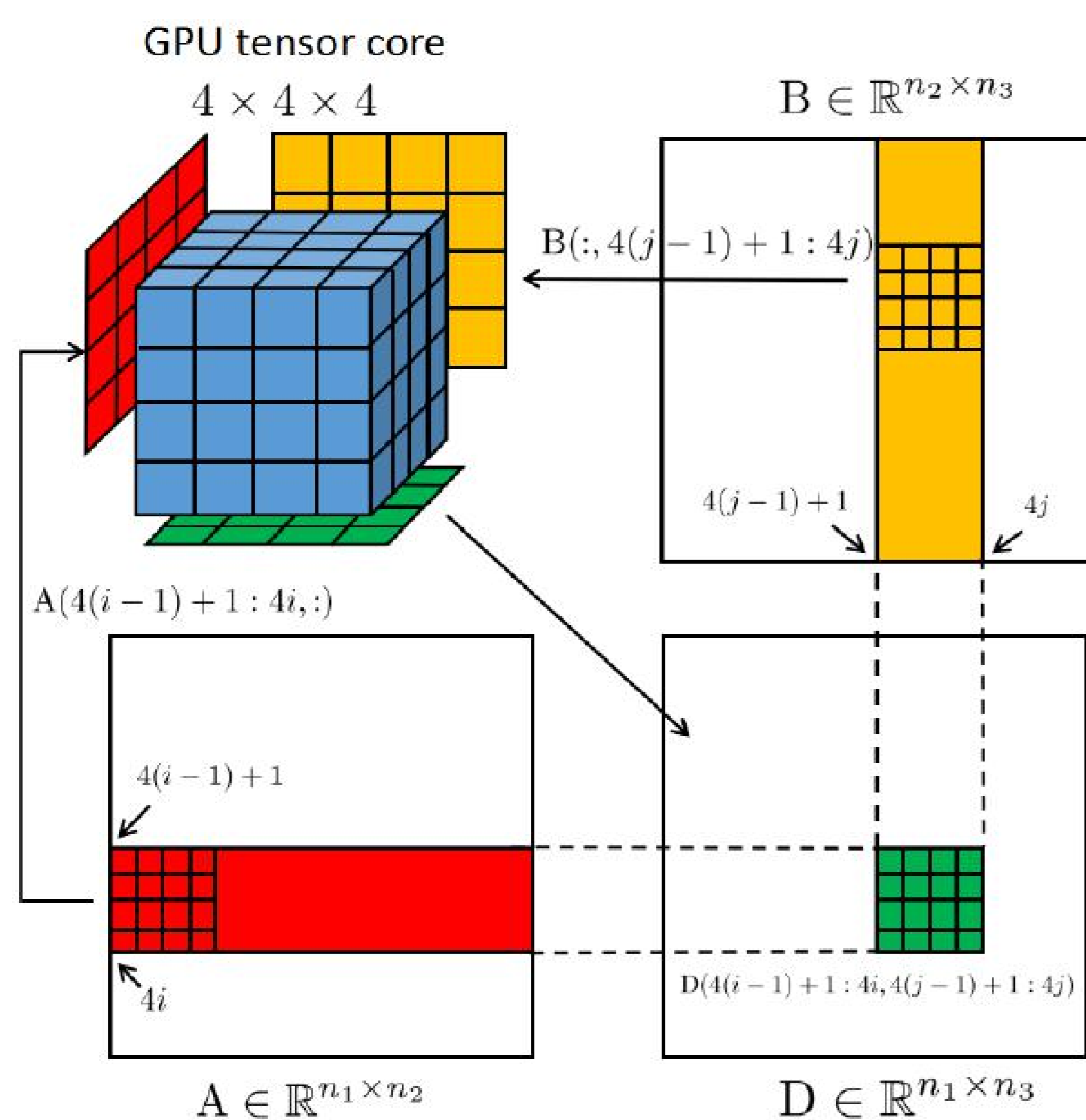


## Parallel HT Decomposition on the GPU

- **Overview of the algorithm**

**Algorithm 2** Optimized hierarchical Tucker tensor decomposition using GPU tensor cores

1: **Input:** Tensor $\mathcal{X} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$, ranks $r_1, r_2, r_3, r_4$.
2: **for** $i = 1$ to 3 **do** in parallel
3:     $H_i = X_{(i)} X_{(i)}^T$,
4: **end for**
5: **for** $i = 1$ to 3 **do** in parallel
6:     $U_i \leftarrow (r_i$ leading eigen vectors) of $H_i$,
7: **end for**
8: $U_4 \leftarrow (r_4$ leading left singular vectors) of $X_{(3)}^T$ using the TSQR algorithm,
9: $\mathcal{B}_2 = \mathcal{U}_4 \times_1 U_1^T \times_2 U_2^T$,
10: $\mathcal{B}_1 = \mathcal{X} \times_1 U_4^T \times_2 U_3^T$,
11: **Output:** $\mathcal{B}_1 \in \mathbb{R}^{r_4 \times r_3 \times 1}, \mathcal{B}_2 \in \mathbb{R}^{r_1 \times r_2 \times r_4}, U_3 \in \mathbb{R}^{n_3 \times r_3}, U_2 \in \mathbb{R}^{n_2 \times r_2}, U_1 \in \mathbb{R}^{n_1 \times r_1}$.
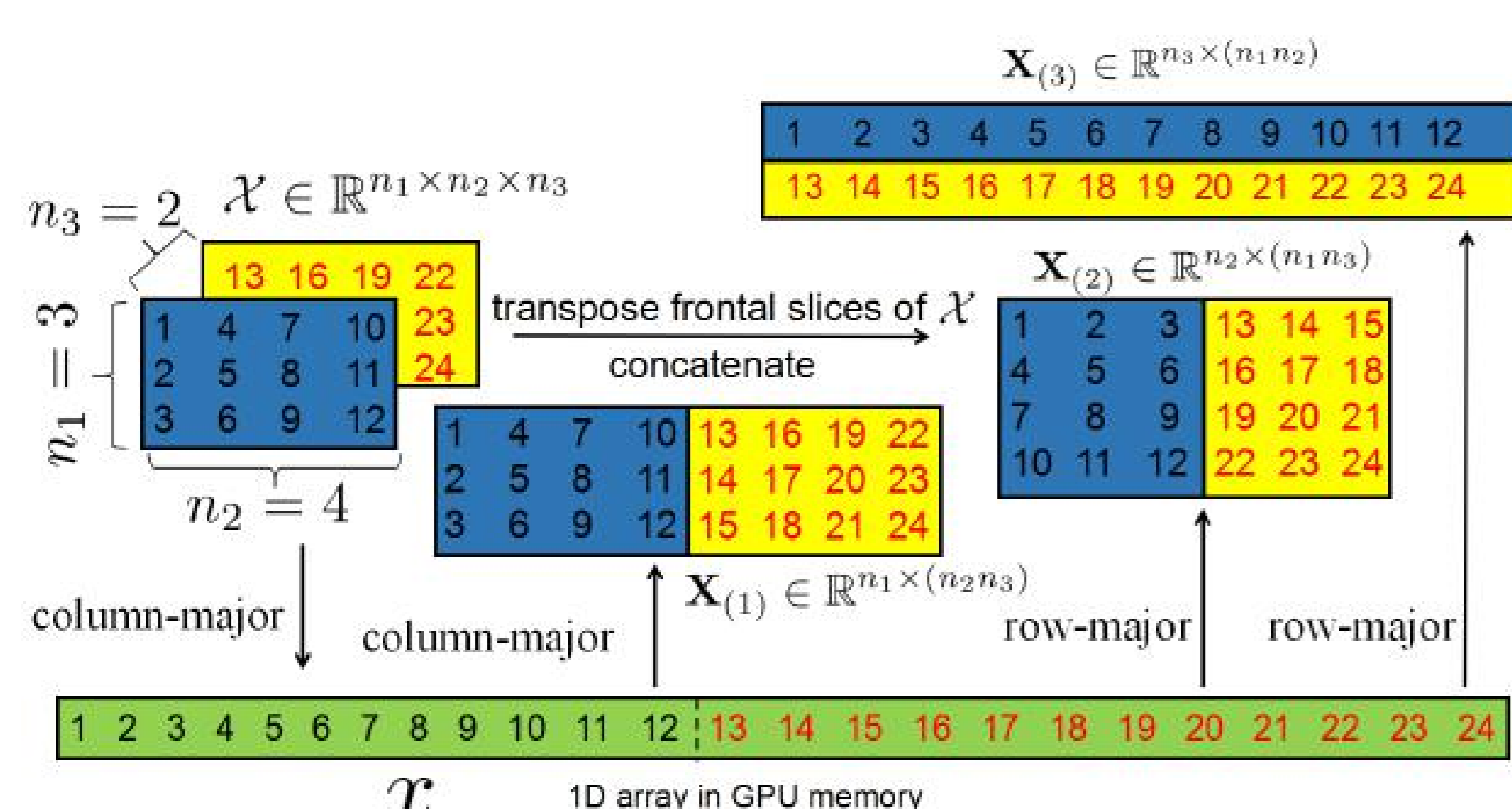
- **Matrix Multiplication Using GPU Tensor Cores**

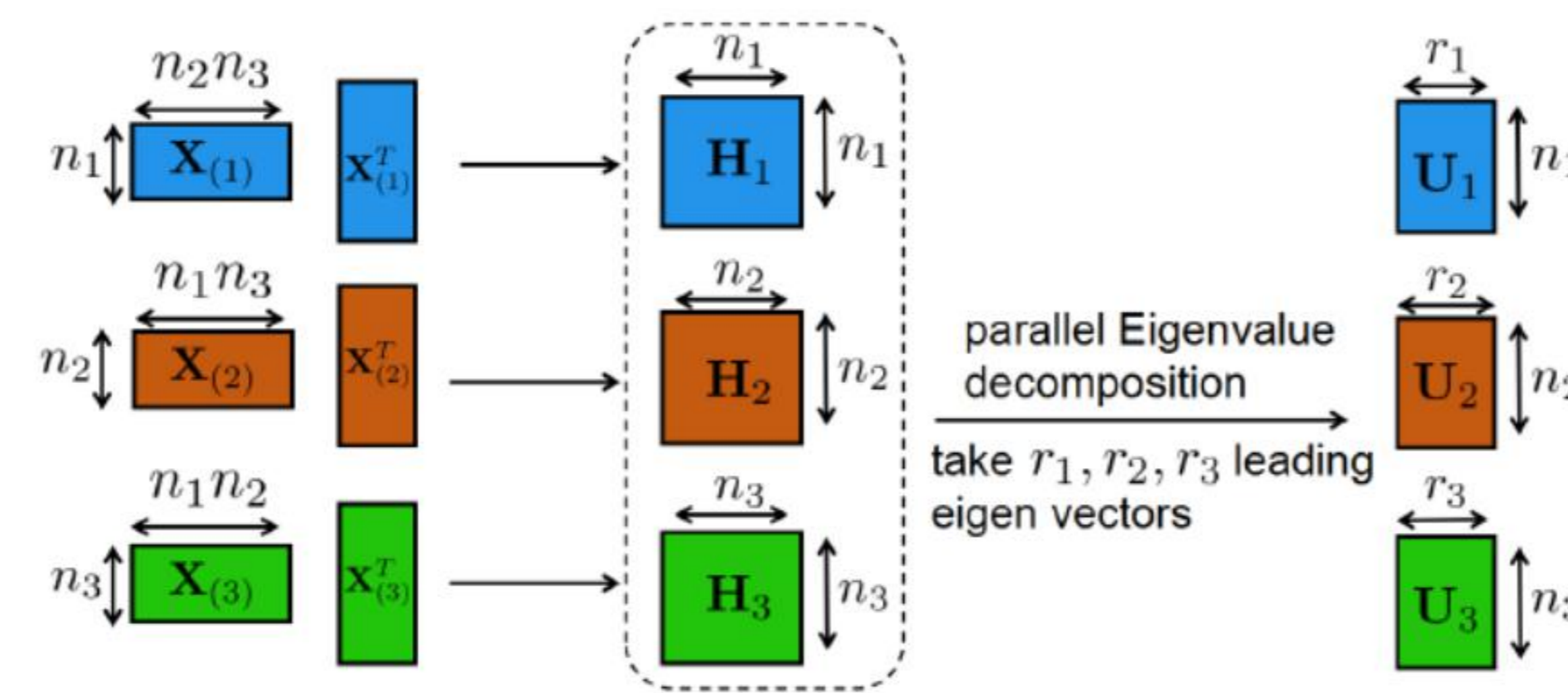We schedule GPU tensor cores to optimize matrix multiplication.



- **Matricization-free Memory Access**

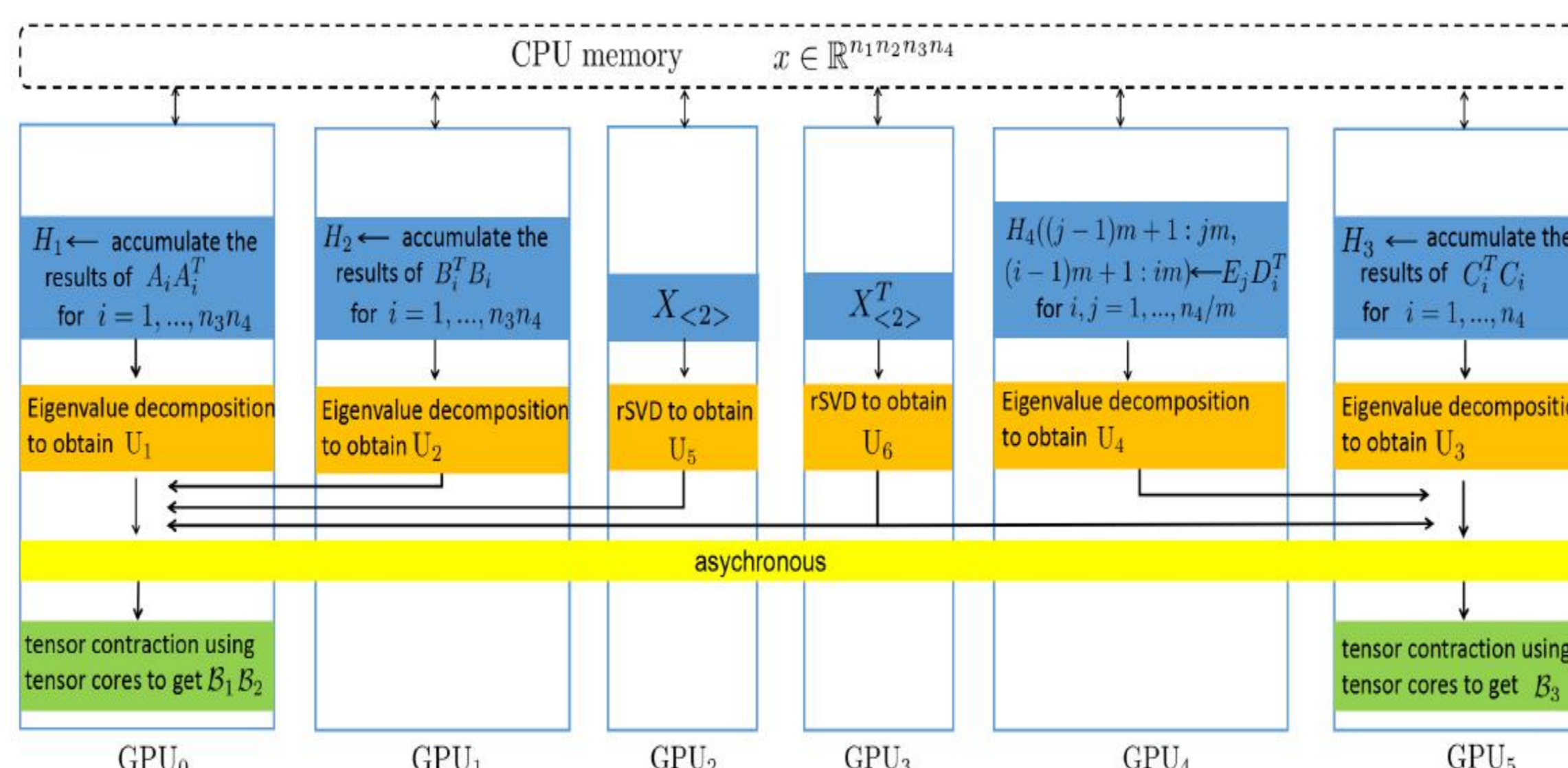We use a matricization-free memory access and avoid consumptions of conversions between tensors and matrices.



- **Batch operation**

Use Eigenvalue decomposition instead of SVD and take the batch solution to calculate intermediate variables in parallel.
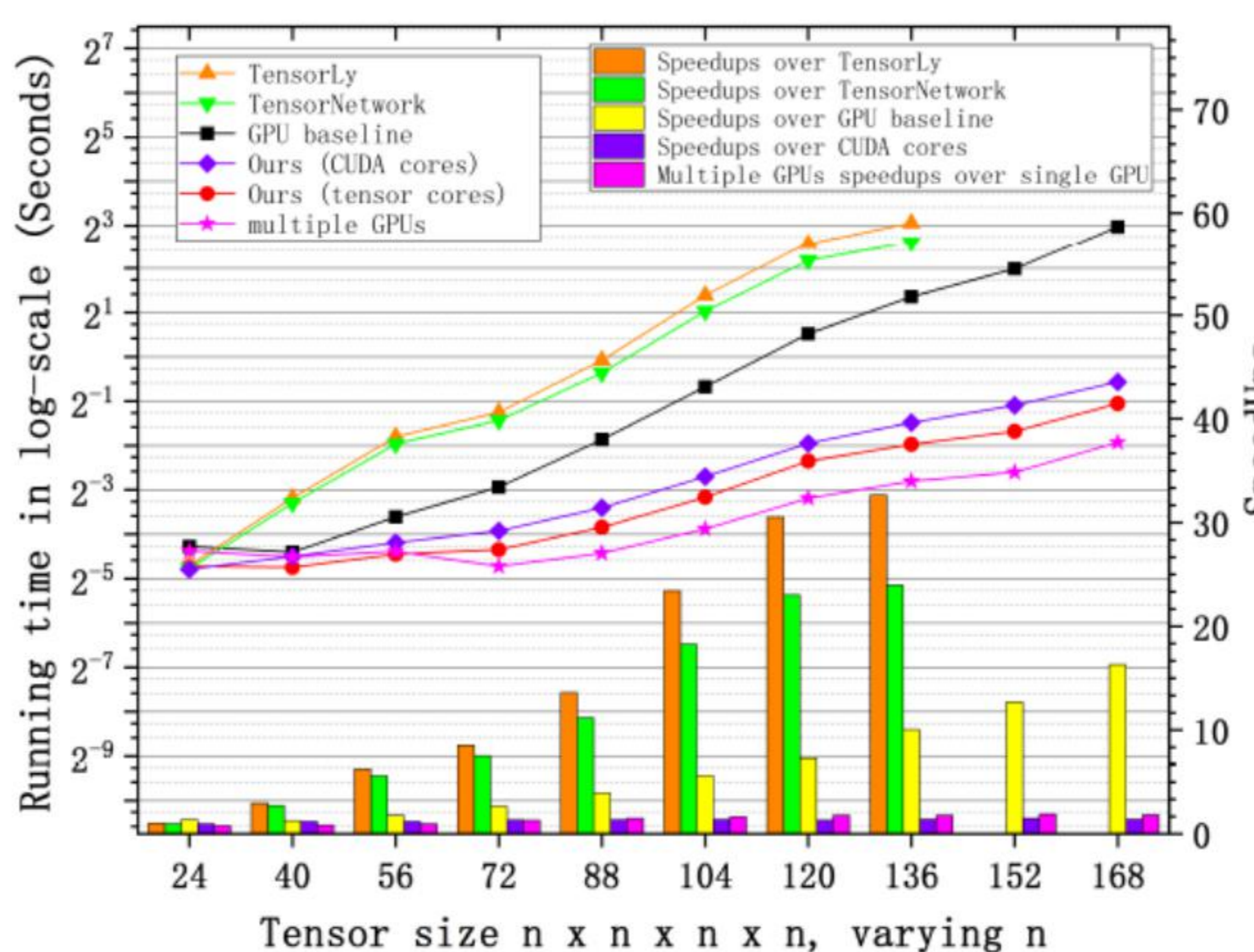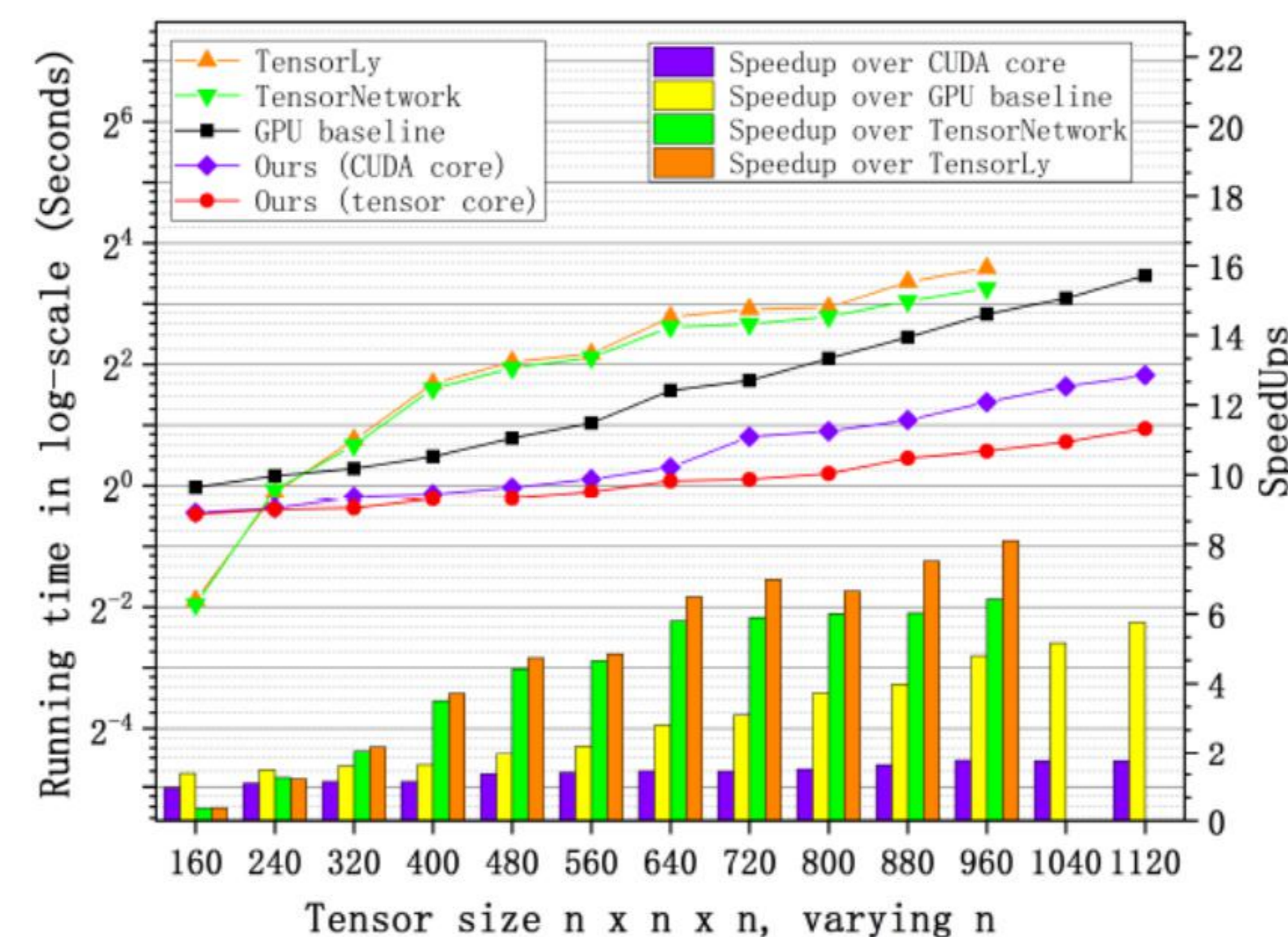


- **Shard mode on multiple GPUs**

High-order HT tensor decomposition on multiple GPUs using shard mode. We divide the original input tensor to shards and transfer each shard to different GPU memory.
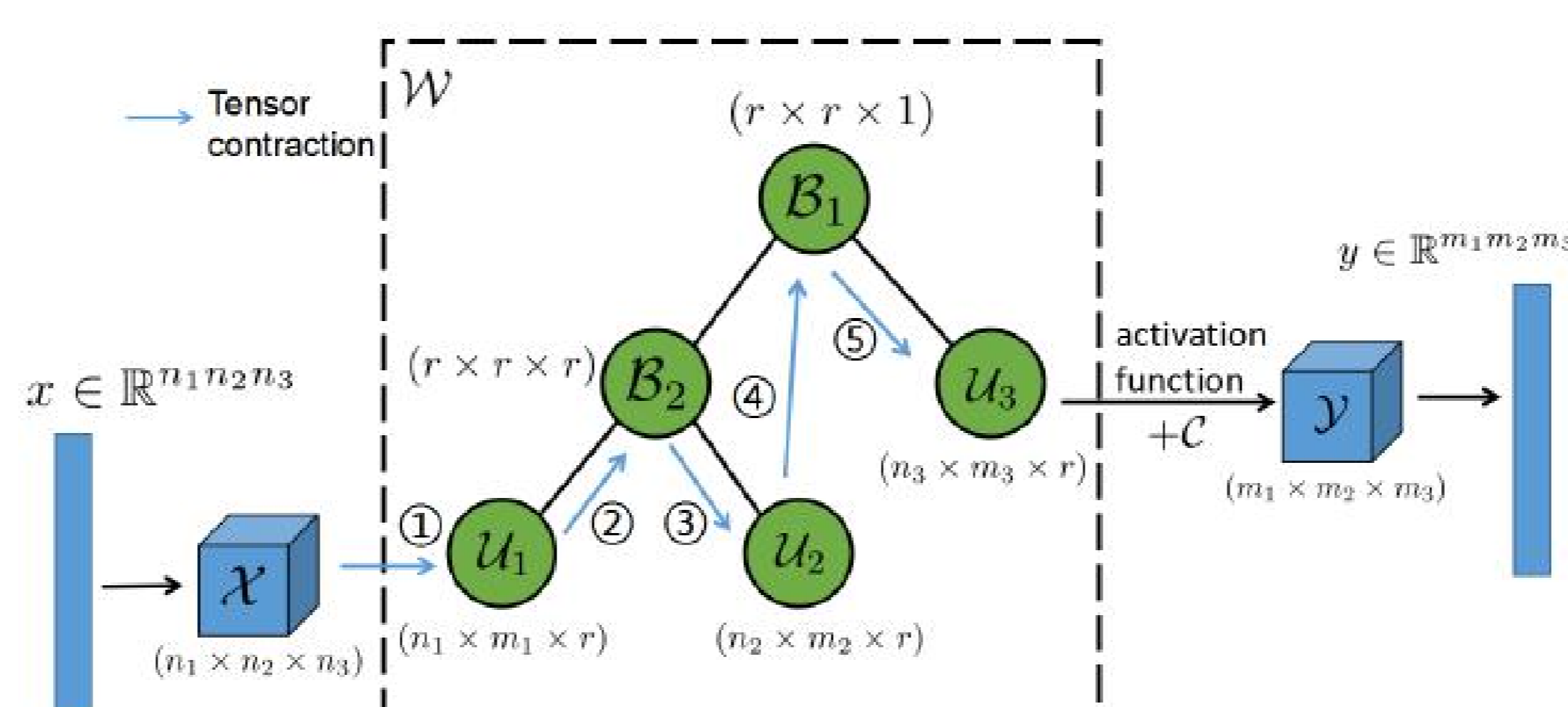


- **Experimental results:**





## HT Tensor Layer for Deep Neural Networks Compression
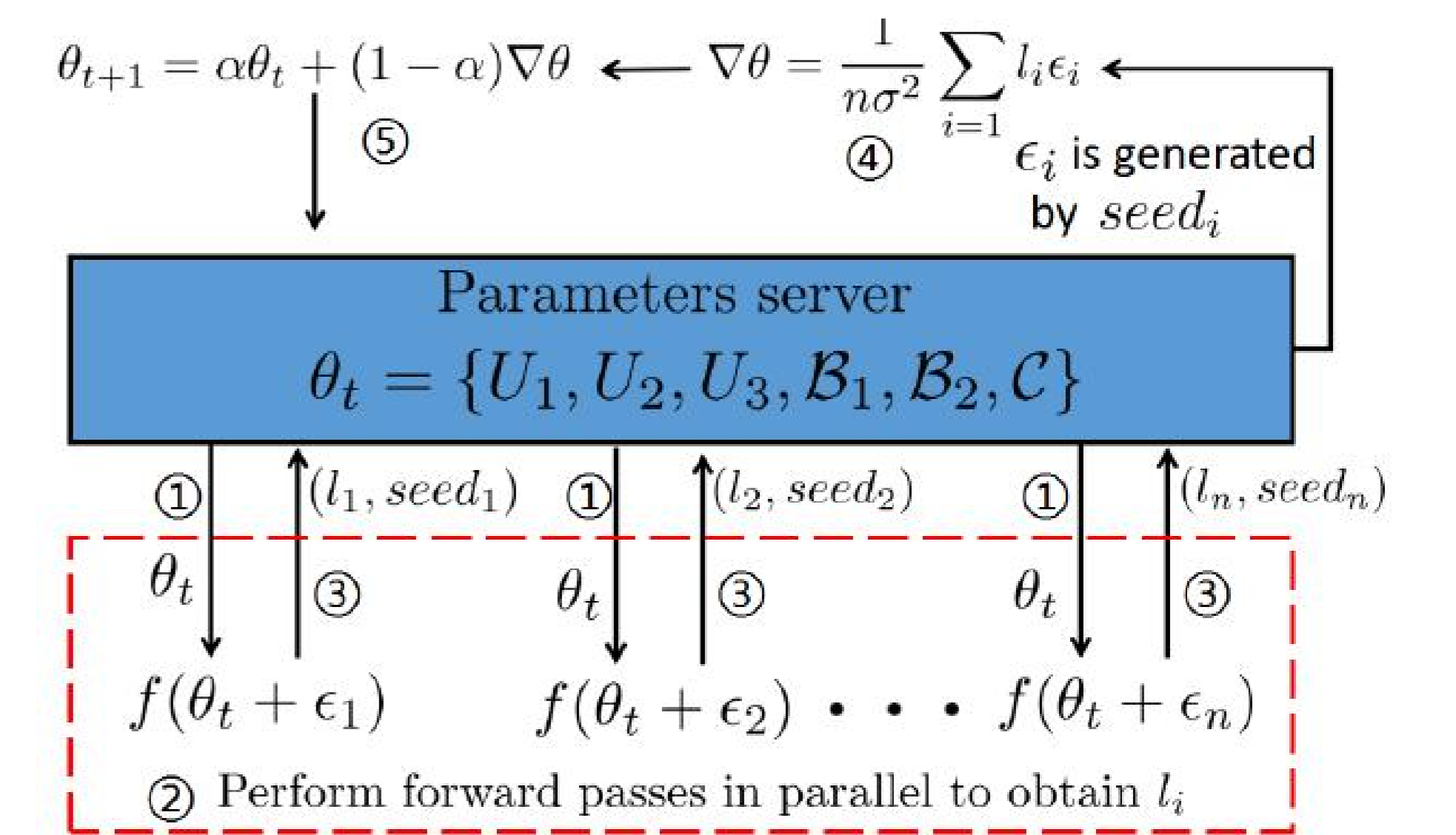
- **forward pass**

For forward pass, we use a HT tensor layer to represent W in a HT format and perform optimized tensor contractions.
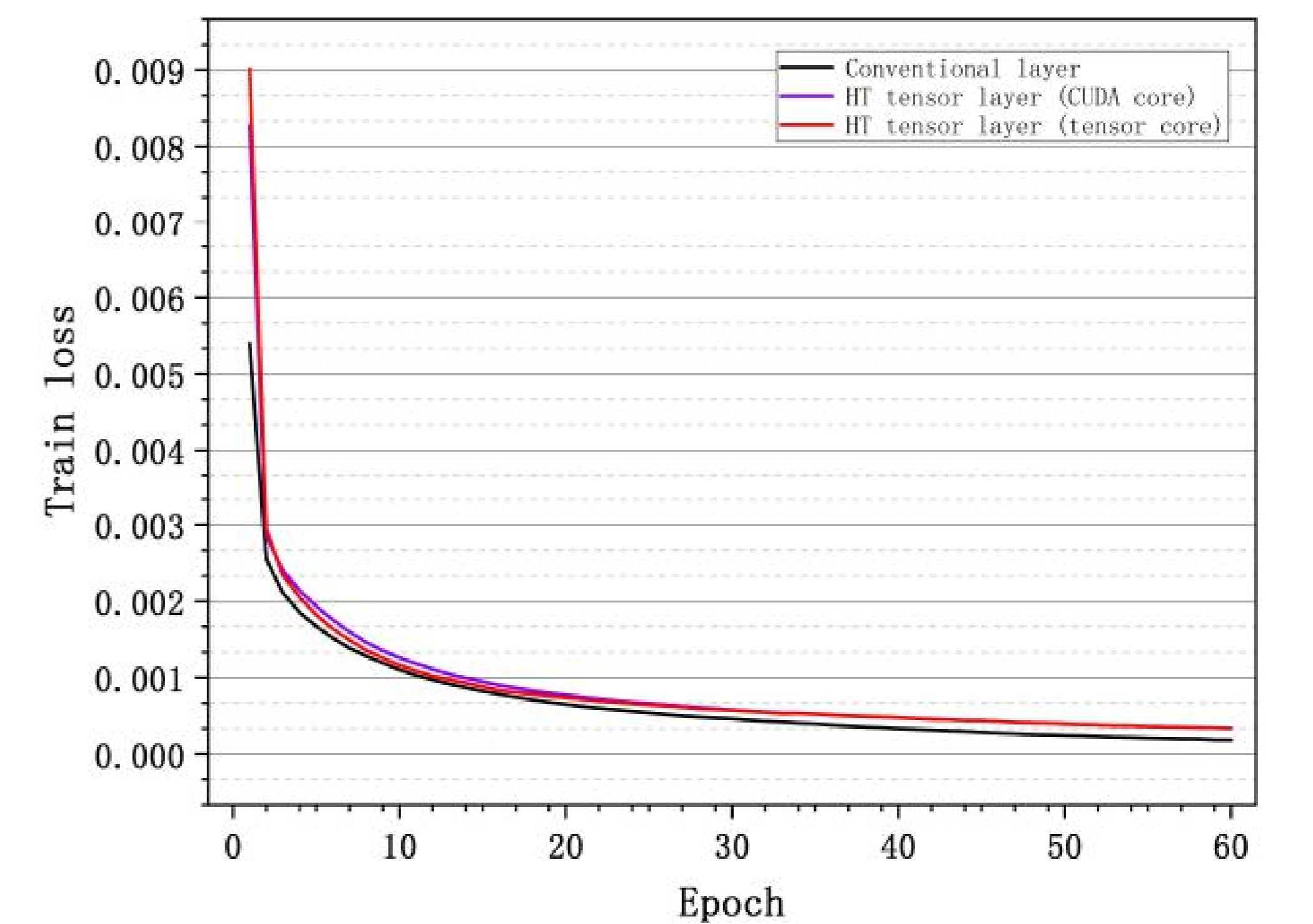


- **Parameters update process**

Use natural gradients to compute in parallel and update the parameters.
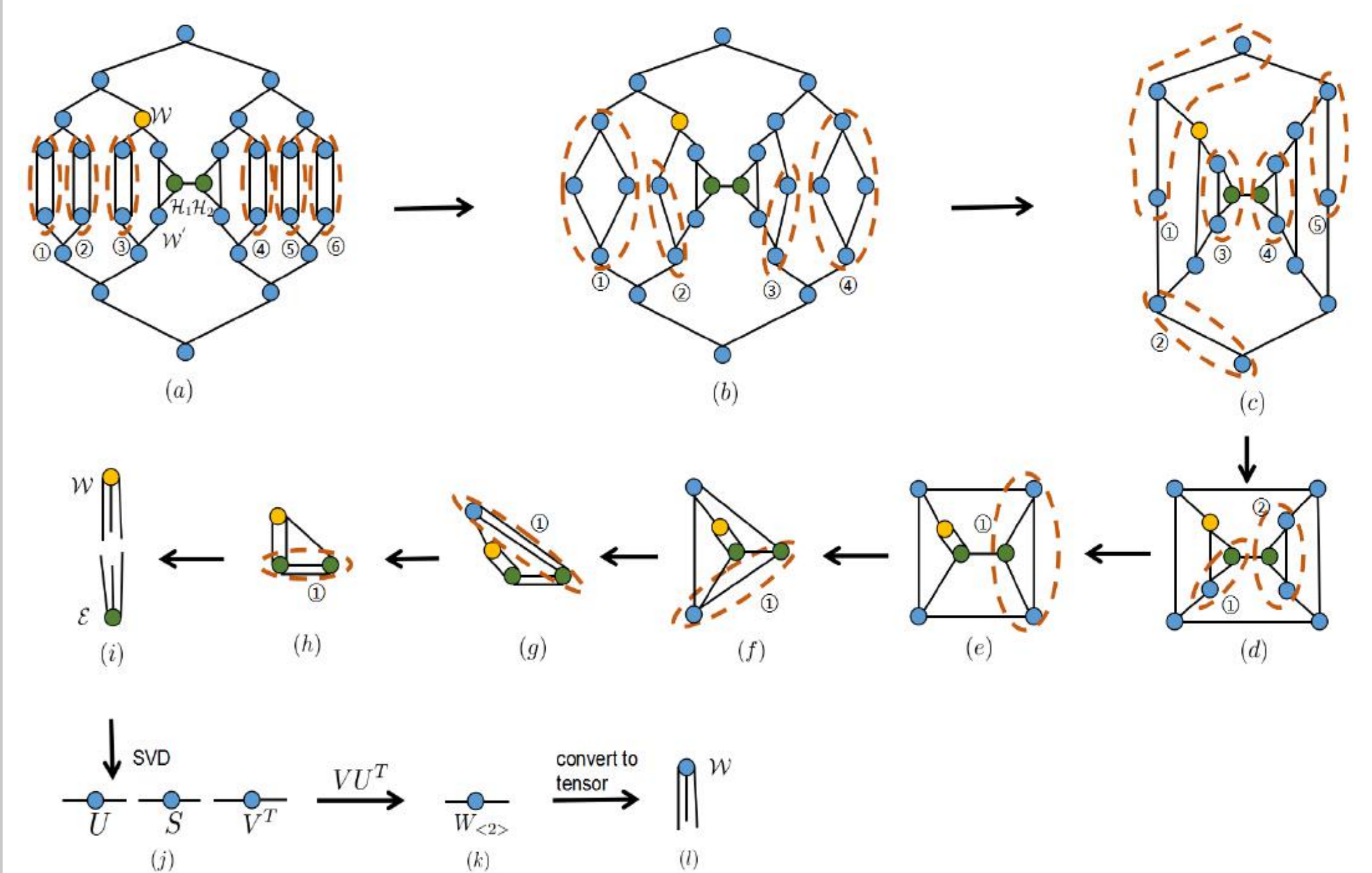


- **Experimental results:**



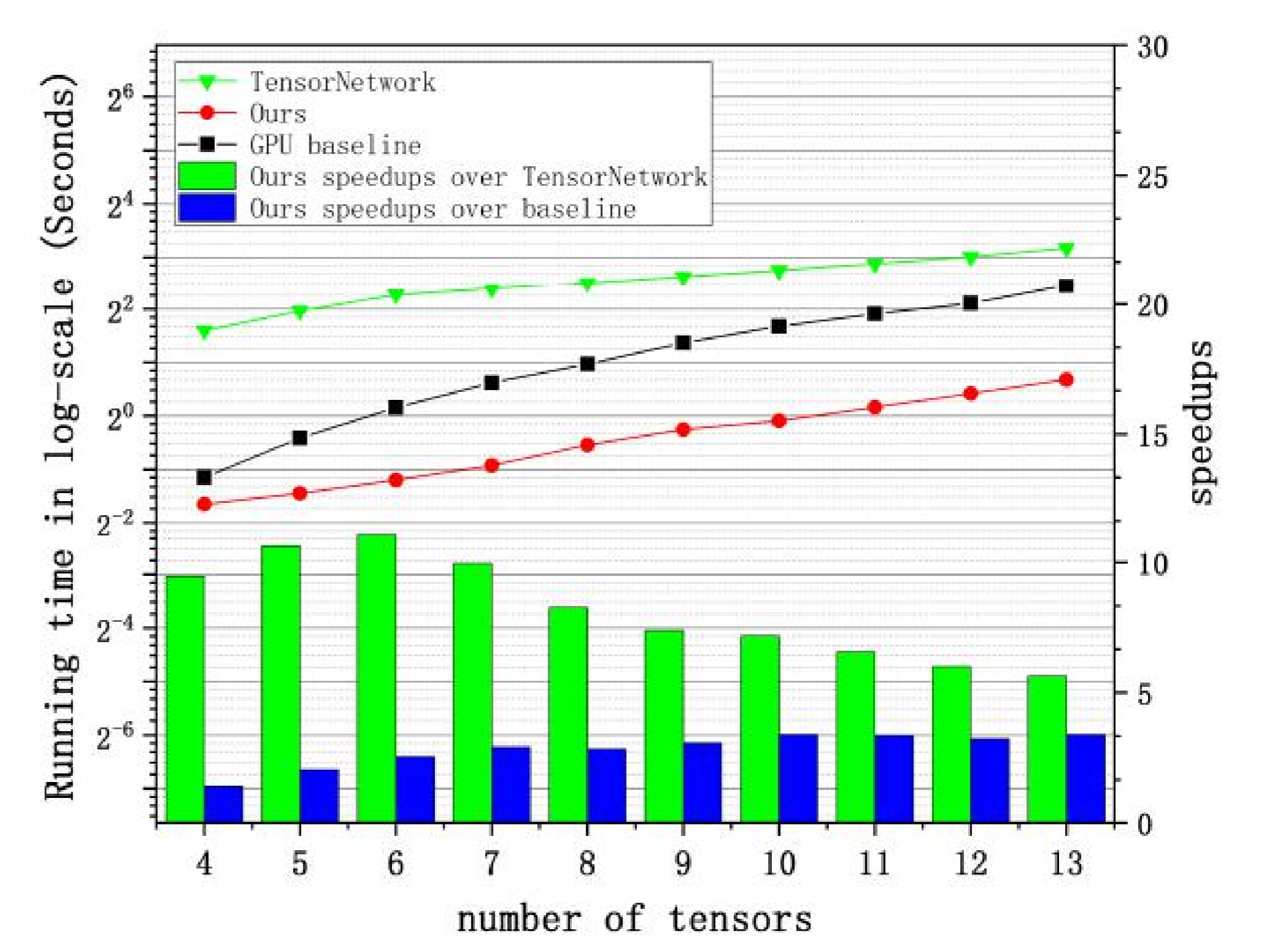| Methods | Parameters | Accuracy | Running time (s) |
|---|---|---|---|
| Conventional layer | 1, 071, 370 | 98.8% | 38, 306 |
| HT tensor layer (CUDA cores) | 21, 770 (49.2×) | 98.3% | 102, 705 |
| HT tensor layer (tensor cores) (ours) | 21, 770 (49.2×) | 98.3% | 72, 327 |

## Efficient TTN Algorithm Using GPU Tensor Cores

- **Algorithm optimization:**

TTN algorithm consists of multiple tensor contractions, and we optimize the algorithm according to the parallelism.



- **Experimental results:**



## Conclusions

In this work, we implement the optimized tensor learning primitives into HT decompositions, a HT tensor layer for deep neural network and TTN algorithm. Experimental results demonstrated that our optimized algorithms have higher efficiency.

[1] H. Huang, X. Liu, W. Tong, T. Zhang and A. Walid. High Performance  Hierarchical Tucker Tensor Learning Using GPU Tensor Cores. IEEE Transactions on Computers. Under review.