

# 基于异构系统的多级并行稀疏张量向量乘算法

陈玥丹<sup>1,4)</sup> 肖国庆<sup>2,3,4)</sup> 阳王东<sup>2,3)</sup> 金纪勇<sup>5)</sup> 龙军<sup>1)</sup> 李肯立<sup>2,3)</sup>

<sup>1)</sup>(中南大学 大数据研究院 长沙 410083)

<sup>2)</sup>(湖南大学 信息科学与工程学院 长沙 410082)

<sup>3)</sup>(国家超级计算长沙中心 长沙 410082)

<sup>4)</sup>(湖南大学深圳研究院 广东 深圳 518000)

<sup>5)</sup>(之江实验室 基础理论研究院一应用数学与机器智能研究中心 杭州 311100)

**摘要** 张量在许多实际应用中被用来表示大规模、多源、高维、多模态的数据. 稀疏张量分解作为挖掘数据中隐藏信息的有效方法之一, 已被广泛应用于机器学习、文本分析、生物医疗等研究领域. 稀疏张量向量乘 (Sparse Tensor-Vector Multiplication, SpTV) 是张量分解中最基础、耗时最多的运算之一. 为加速大数据和人工智能相关应用的运行效率, 本文提出了基于 CPU-GPU 异构系统的多级并行 SpTV 加速算法. 首先, 为了将 SpTV 运算映射到混合、多级并行的分布式 CPU-GPU 异构多/众核构架, 本文设计了一种多维并行 SpTV 划分方法, 采用面向节点级并行的  $N-1$  维张量划分和面向 GPU 线程级并行的矩阵划分, 充分利用计算节点间和节点内的多级并行计算能力. 其次, 设计了一种基于稀疏张量纤维的压缩存储格式, 压缩稀疏张量的内存占用, 优化 SpTV 运算的计算和访存模式. 最后, 提出了基于多流并行的异构高效 SpTV 算法, 进一步设计了稀疏张量的细粒度划分方法、多流并行运行机制和基于张量块排序的多流并行优化技术, 实现了 SpTV 运算中通信开销和计算开销的相互重叠与隐藏. 实验结果表明, 与相关工作 aeSpTV 相比, 所提出的 SpTV 算法在所有测试数据集上最高能够获得 3.28 倍的加速比.

**关键词** CPU-GPU; 异构并行计算; 多级并行; 稀疏张量; 张量运算

中图法分类号 TP301

## Exploiting Hierarchical Parallelism for Sparse Tensor-Vector Multiplication on Heterogeneous Parallel Systems

CHEN Yue-Dan<sup>1,4)</sup> XIAO Guo-Qing<sup>2,3,4)</sup> YANG Wang-Dong<sup>2,3)</sup> JIN Ji-Yong<sup>5)</sup> LONG Jun<sup>1)</sup> LI Ken-Li<sup>2,3)</sup>

<sup>1)</sup>(Big Data Institute, Central South University, Changsha 410083)

<sup>2)</sup>(College of Computer Science and Electronic Engineering, Hunan University, Changsha 410082)

<sup>3)</sup>(National Supercomputing Center in Changsha, Changsha 410082)

<sup>4)</sup>(Shenzhen Research Institute, Hunan University, Shenzhen, Guangdong 518000)

<sup>5)</sup>(Research Center for Applied Mathematics and Machine Intelligence, Zhejiang Lab, Hangzhou 311100)

本课题得到广东省重点领域研究计划(2021B0101190004)、国家自然科学基金(62172157、62202149)、湖南省科技项目(2023GK2002、2021RC3062、2023JJ60002)、广东省自然科学基金(2023A1515012915)、深圳市基础研究面上项目(JCYJ20210324135409026)、之江实验室开放课题(2022RC0AB03)资助. 陈玥丹(通信作者), 博士, 副教授, 中国计算机学会(CCF)会员, 主要研究领域为高性能计算、并行与分布式处理、人工智能与大数据计算. E-mail: chenyeddan@hnu.edu.cn. 肖国庆(通信作者), 博士, 副教授, 中国计算机学会(CCF)高级会员, 主要研究领域为高性能计算与智能计算. E-mail: xiaoguoqing@hnu.edu.cn. 阳王东, 博士, 教授, 主要研究领域为并行计算. 金纪勇, 硕士, 助理研究员, 主要研究领域为分布式计算. 龙军, 博士, 教授, 中国计算机学会(CCF)杰出会员, 主要研究领域为大数据计算与智能软件系统. 李肯立, 博士, 教授, 中国计算机学会(CCF)会士, 主要研究领域为高性能计算系统软件与应用.

**Abstract** Many application domains give rise to multidimensional data that can be naturally represented via tensors. The tensors used in most real-world applications that are extremely large and very sparse. The sparse tensor decomposition is an effective approach to predict the unobserved data and is commonly used in machine learning, text analysis, healthcare analytics, and numerous other applications. Sparse tensor-vector multiplication (SpTV) is one of the most fundamental and time-intensive operations in computing tensor decomposition. In order to improve the efficiency of related applications, this paper exploits the hierarchical parallelism for SpTV on CPU-GPU heterogeneous parallel computing systems. First of all, we propose a multidimensional partitioning method to map parallel SpTV to the underlying CPU-GPU heterogeneous and parallel computing architectures. It utilizes the  $N-1$ -dimensional tensor partitioning to exploit the inter-node parallelism and the matricized tensor partitioning to exploit the intra-node parallelism. Second, based on the multidimensional data partitioning, we design a fiber-wise compressed storage format for sparse tensors to reduce the memory footprint and optimize the computing and memory accessing patterns in parallel SpTV. Third, we design the parallel streaming SpTV algorithm, by adopting the fine-grained data partitioning method, the parallel streaming execution scheme, and the tensor block sorting technique, to overlap the data swapping cost and the computation overhead and further leverage the computing power of GPUs. The experimental results show that the parallel and efficient SpTV algorithm achieves the speedup of up to 3.28 compared to state-of-the-art (aeSpTV) on a CPU-GPU system.

**Key words** CPU-GPU; heterogeneous and parallel computing; hierarchical parallelism; sparse tensors; tensor operations

## 1 引言

张量是多维度数据表现形式, 稀疏张量中的大部分元素为零, 因此零元素不需要存储并参与相关计算. 稀疏张量通常在人工智能<sup>[1-3]</sup>、数据挖掘与分析<sup>[4-6]</sup>、医疗保健<sup>[7-9]</sup>等很多实际应用中表示大规模、多源、多模态的数据, 因此获得了许多研究者的关注. 例如, 电子邮件的 4 个属性{主题, 发件人, 收件人, 时间}可以用一个 4 维张量来表示, 其中每封邮件的某个属性值对应于该 4 维张量中每个非零元素在相应维度上的取值. 张量的分解可以挖掘出张量数据中的隐藏信息, 因此广泛地应用于推荐系统、会话检测等实际应用.

张量 CANDECOMP/PARAFAC (CP)分解<sup>[10-12]</sup>是应用最广泛的张量分解方法之一. 计算 CP 分解最常用的方法是基于交替最小二乘法(Alternating Least Squares, ALS)进行的, 该方法简称为 CP-ALS. CP-ALS 算法采用一系列张量运算, 通过每次迭代为张量的每个维度计算一个新的因子矩阵. 具体来说, 对于一个 3 阶张量  $\mathbf{X} \in \mathbf{R}^{I \times J \times K}$  (秩为  $R$ ) 的 CP 分解计算, 需要获得 3 个稠密的因子矩阵  $\mathbf{A} \in \mathbf{R}^{I \times R}$ 、 $\mathbf{B} \in \mathbf{R}^{J \times R}$  和  $\mathbf{C} \in \mathbf{R}^{K \times R}$ . 张量  $\mathbf{X}$  的展开矩阵可以由  $\mathbf{A}$ 、 $\mathbf{B}$  和  $\mathbf{C}$  表示:

$$\mathbf{X}_{(1)} \approx \mathbf{A}(\mathbf{C} \odot \mathbf{B}), \mathbf{X}_{(2)} \approx \mathbf{B}(\mathbf{C} \odot \mathbf{A}), \mathbf{X}_{(3)} \approx \mathbf{C}(\mathbf{B} \odot \mathbf{A}) \quad (1)$$

其中,  $\mathbf{X}_{(1)}$ 、 $\mathbf{X}_{(2)}$  和  $\mathbf{X}_{(3)}$  分别为张量  $\mathbf{X}$  在三个维度上的展开矩阵(如第 2.1 节介绍), 运算符  $\odot$  被称为 Khatri-Rao 乘积(Khatri-Rao Product), 下式给出了  $\mathbf{A} \odot \mathbf{B}$  运算示例:

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}, \mathbf{B} = \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \\ b_{31} & b_{32} \end{bmatrix},$$

$$\mathbf{A} \odot \mathbf{B} = \begin{bmatrix} a_{11}b_{11} & a_{12}b_{12} \\ a_{11}b_{21} & a_{12}b_{22} \\ a_{11}b_{31} & a_{11}b_{32} \\ a_{21}b_{11} & a_{22}b_{12} \\ a_{21}b_{21} & a_{22}b_{22} \\ a_{21}b_{31} & a_{22}b_{32} \end{bmatrix} \quad (2)$$

基于 ALS 的 CP 分解通过迭代运算求得  $\mathbf{A}$ 、 $\mathbf{B}$  和  $\mathbf{C}$ . 例如, 在计算  $\mathbf{A}$  时, 每次迭代相当于进行一个最小二乘问题的求解, 该问题的最优解形式为:

$$\hat{\mathbf{A}} = \mathbf{X}_{(1)}(\mathbf{C} \odot \mathbf{B})(\mathbf{C} \mathbf{C}^* \mathbf{B} \mathbf{B})^\dagger \quad (3)$$

其中,  $(\mathbf{C} \mathbf{C}^* \mathbf{B} \mathbf{B})$  是一个  $I \times R$  的矩阵, 运算

$(\mathbf{C} \mathbf{C}^* \mathbf{B} \mathbf{B})^\dagger$  是求该矩阵的逆. 因此, 与运算  $\mathbf{X}_{(1)}(\mathbf{C} \odot \mathbf{B})$  (被称为 MTTKRP (matricized tensor times Khatri-Rao product)) 相比, 运算  $(\mathbf{C} \mathbf{C}^* \mathbf{B} \mathbf{B})^\dagger$  涉及了很少的计算量. 可以看出, MTTKRP 运算  $\mathbf{X}_{(1)}(\mathbf{C} \odot \mathbf{B})$  是 CP 分解中的主要性能瓶颈. 同时, MTTKRP 的本质就是进行  $I$  次 SpTV 运算, 因此实际应用中张量 CP 分解的运算效率取决于 SpTV 的性能. 由于 SpTV 运算中较低的计算/访存比和不规则的数据索引, 已成为计算张量 CP 分解所面临的最基础的性能挑战.

在现实应用中, 稀疏张量的规模非常大. 为了支持大规模实际应用中稀疏张量数据的存储与计算, 越来越多的工作利用分布式高性能计算系统, 将大规模张量划分为较小的数据块, 并分布到多个计算节点上处理. 目前主流的高性能计算系统设计趋势是采用“主处理器+加速器”的异构多/众核计算构架(例如 CPU-GPU), 这种高性能计算系统的大部分计算能力都来自加速器, 因此, 在设计高效、可扩展算法时, 充分挖掘加速器中强大的计算性能, 以实现目标计算任务的效率是至关重要的.

大规模 SpTV 运算在分布式异构多/众核系统上的并行加速主要面临着 3 个主要的挑战:

(1) 如何有效地将 SpTV 运算映射到目标分布式异构多/众核构架.

分布式异构多/众核系统提供了混合、多级的并行计算能力, 包括计算节点间的并行性、以及节点内主处理器与加速器提供的混合同行性. 数据划分是并行化 SpTV 运算、开发平台构架混合、多级并行性的重要方法之一.

许多研究提出了面向分布式异构多/众核构架的稀疏张量划分模型, 降低稀疏张量的通信成本的同时, 保证其在所有计算节点上的计算成本. Karsavuran 等人<sup>[13]</sup>提出一种针对稀疏张量数据的中粒度划分的超图模型, 该模型将张量划分为非零元素不相交的子张量, 然后对每个子张量构造一个模相关的粗粒超图, 从而将通信量最小化. aeSpTV<sup>[14]</sup> 基于目标平台的计算构架和内存结构采用一种自适应的稀疏张量划分方法, 缓解 SpTV 运算中不规则访存的高开销、并行写冲突、额外的计算量和中间结果爆炸等问题. ALTO<sup>[15]</sup> 将张量的非零元素划分到不同分块中, 这些分块的空间坐标可能存在重叠, 为了解决这些重叠可能产生的更新冲突, ALTO

在多维空间中定位它们的位置, 并根据张量数据的重用性自动选择适当的同步机制. CASpMV<sup>[16]</sup> 采用了一种基于描述稀疏矩阵结构特征统计模型的自动调节的稀疏矩阵划分方法, 通过分析非零元素分布特征和目标计算平台的构架特征, 将稀疏矩阵均匀地分布到不同的计算节点上.

然而, 现有针对 SpTV 加速的稀疏张量划分模型没有考虑通过细粒度张量数据划分来发挥软件流水并行技术, 隐藏对细粒度数据块的数据传输时间和并行计算时间, 进一步利用目标平台的计算能力.

(2) 如何更高效地压缩并存储大规模稀疏张量.

实际应用中的张量数据具有非常稀疏、非零元素分布不规则的特点. 一方面, 在基于异构并行计算系统的 SpTV 运算中, 稀疏张量的存储方式决定了内存占用量, 因此影响了计算节点间数据传输和节点内主处理器与加速器之间数据交换的通信效率. 另一方面, 适合的稀疏张量存储结构能够优化 SpTV 运算中的不规则计算和访存模式, 从而提高并行 SpTV 算法的总体运算效率.

很多相关学者提出了稀疏张量存储格式, 以及基于这种存储格式实现的稀疏张量运算算法. Qin 等人<sup>[17]</sup> 设计了能够支持多种稀疏张量存储格式(包括 COO (Coordinate), CSR (Compressed Sparse Row), CSC (Compressed Sparse Column)<sup>[18]</sup>, DIA(Diagonal)<sup>[19]</sup>, ZVC(Zero-Value Compression)<sup>[20]</sup>, HICOO (Hierarchical Coordinate)<sup>[21]</sup>, CSF (Compressed Sparse Fiber)<sup>[22,23]</sup> 的加速器硬件扩展.

其中, ZVC<sup>[20]</sup> 为每  $N$  个非零元素生成一个  $N$  位掩码, 其中掩码为 ‘0’ 表示对应的张量元素值为零, 而 ‘1’ 表示该元素为非零值. 生成  $N$  位掩码后, 再将追加存储张量非零元素. 因此,  $N$  个连续的张量零元素可以压缩为一个  $N$  位全零掩码, 进而获得  $N$  倍的压缩比. 然而, ZVC 只考虑了张量零元素的压缩存储, 没有考虑非零元素和张量分块相关索引数据的压缩存储. HICOO<sup>[21]</sup> 对稀疏张量分块中索引数据进行压缩, 并使用更短的整数类型来表示张量分块的偏移量. 但 HICOO 格式没有进一步对张量纤维维度的索引进行压缩. CSF<sup>[22,23]</sup> 根据稀疏张量运算基于纤维的计算特性, 利用压缩、分层和基于纤维的树形结构表示稀疏张量. 但是基于 CSF 数据结构的张量运算总计算量会随着纤维片段的

增加而增加<sup>[24]</sup>. Dun 等人<sup>[25]</sup>基于分块和位映射设计了一种新的稀疏矩阵压缩存储格式 TB-COO, TB-COO 压缩了张量切面维度的索引,但在还原压缩索引时,需要进行一系列位操作. Sparta<sup>[26]</sup>使用多维、高效的哈希表表示 SpTC 运算的输入稀疏张量.然而,基于哈希表的稀疏张量表示增加了运算中的索引开销,在不擅长访存密集型任务和逻辑运算的加速器(如 GPU 等)上,很难发挥良好的运算加速效果. Qiu 等人<sup>[27]</sup>假设给定的大张量具有低多线性秩,提出使用 Tucker 模型对原始大张量进行压缩表示,从而有效地执行张量分解.

(3) 基于给定的任务划分方法和稀疏张量压缩存储格式,如何充分利用高性能异构并行计算系统的计算能力,优化 SpTV 运算性能.

因此,需要基于并行任务划分和数据存储结构设计,考虑计算平台上的并行计算效率与节点间、节点内的通信效率,为 SpTV 算法进行并行化设计与优化,充分发挥目标平台的计算能力.

为应对和缓解上述挑战,本文基于分布式异构 CPU-GPU 计算系统,设计了多级并行的 SpTV 算法.本文主要包括四个贡献:

(1) 针对上述第 1 个挑战,本文设计了一种多维并行 SpTV 划分方法.面向分布式异构 CPU-GPU 计算系统,采用面向节点级并行的  $N-1$  维张量划分和面向 GPU 线程级并行的矩阵划分,降低张量的稀疏度的同时,将计算任务更均匀地映射到多级并行计算架构上,避免了节点间通信和额外的累加计,减少了节点内的通信开销.

(2) 针对上述第 2 个挑战,本文为稀疏张量设计了一种基于张量纤维的压缩存储格式,根据 SpTV 运算属性、以及设计的多维并行 SpTV 划分方法,按稀疏张量纤维对非零元素进行压缩存储,进而减少稀疏张量的内存占用,提高了多维并行 SpTV 运算的数据索引效率.

(3) 针对上述第 3 个挑战,本文面向分布式异构 CPU-GPU 计算系统,设计了基于张量块排序的多流并行优化策略.首先,为了充分利用多 CUDA 流的并行性,本文基于多维并行 SpTV 划分方法,进一步设计了对稀疏张量的细粒度划分;并设计了多流并行运行机制,将对基于细粒度数据块的并行 SpTV 运算所产生的通信开销与计算开销进行相互重叠与隐藏;最后,提出了基于张量块排序策略,进一步提高张量块处理开销之间的重叠程度.

(4) 本文选取了 5 个实际应用中的真实稀疏

张量数据集,在 CPU-GPU 构架上测试了多级并行 SpTV 算法的性能.与现有相关 aeSpTV 工作相比,多级并行 SpTV 算法的吞吐量(每秒 10 亿次的浮点运算数(Giga Floating-point Operations Per Second, GFLOPS))最高获得了 3.28 倍的加速比.

## 2 研究背景

### 2.1 张量及其相关基本操作

张量的阶(Order)或模(Mode)表示其维度的数量.向量属于阶为 1 的张量,矩阵属于阶为 2 的张量.阶数不小于 3 的张量被称为高阶张量.

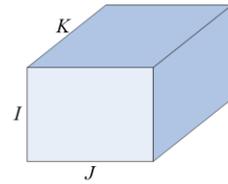


图 1 3阶张量  $X \in \mathbb{R}^{I \times J \times K}$

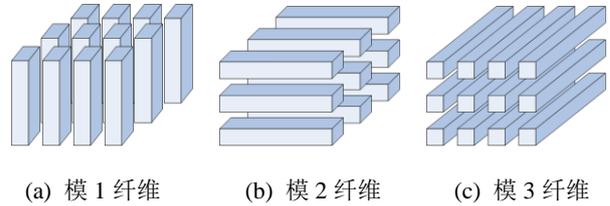


图 2 张量纤维

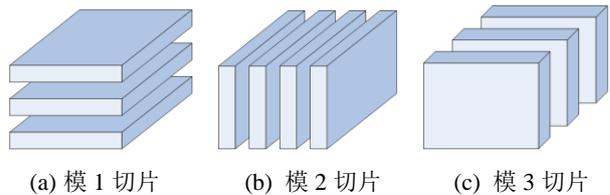


图 3 张量切片

一个  $N$  阶张量的纤维(Fiber)通过固定  $N-1$  个维度上的索引值得到.图 2 分别展示了图 1 中的 3 阶张量  $X \in \mathbb{R}^{I \times J \times K}$  在 3 个维度上的纤维.其中,模 1 纤维表示为  $X_{jk}$ ,模 2 纤维表示为  $X_{ik}$ ,模 3 纤维表示为  $X_{ij}$ ,其中,  $i \in \{1, 2, \dots, I\}$ ,  $j \in \{1, 2, \dots, J\}$ ,  $k \in \{1, 2, \dots, K\}$ .

一个  $N$  阶张量的切片(Slice)通过固定  $N-2$  个维度上的索引获得.图 3 分别展示了图 1 中的 3 阶张量  $X \in \mathbb{R}^{I \times J \times K}$  在 3 个维度上的切片.其中,第  $i$  个模 1 切片表示为  $X_{i:}$ ,第  $j$  个模 2 切片表示为  $X_{:j}$ ,第  $k$  个模 3 切片表示为  $X_{::k}$ ,其中,  $i \in \{1, 2, \dots, I\}$ ,  $j \in \{1, 2, \dots, J\}$ ,  $k \in \{1, 2, \dots, K\}$ .

张量的矩阵化(Matricization)是将张量展开为一个 2 阶矩阵的操作.一个  $N$  阶张量的模  $n$  矩阵化

操作以模  $n$  纤维作为矩阵的列将张量展开, 展开的张量记为  $\mathbf{X}_{(n)}$ , 张量中的元素  $(i_1, i_2, \dots, i_N)$  对应于模  $n$  展开矩阵中的元素  $(i_n, j)$ , 其中

$$j = \sum_{k=1, k \neq n}^N i_k \prod_{m=1, m \neq n}^{k-1} I_m \quad (4)$$

以一个  $2 \times 3 \times 2$  的 3 阶张量  $\mathbf{X}$  为例.  $\mathbf{X}$  的模 3 切片分别为

$$\mathbf{X}_{:0} = \begin{bmatrix} 1 & 3 & 5 \\ 2 & 4 & 6 \end{bmatrix}, \mathbf{X}_{:1} = \begin{bmatrix} 7 & 9 & 11 \\ 8 & 10 & 12 \end{bmatrix} \quad (5)$$

那么,  $\mathbf{X}$  的模 1、模 2 和模 3 展开矩阵分别为

$$\mathbf{X}_{(1)} = \begin{bmatrix} 1 & 3 & 5 & 7 & 9 & 11 \\ 2 & 4 & 6 & 8 & 10 & 12 \end{bmatrix} \quad (6)$$

$$\mathbf{X}_{(2)} = \begin{bmatrix} 1 & 2 & 7 & 8 \\ 3 & 4 & 9 & 10 \\ 5 & 6 & 11 & 12 \end{bmatrix} \quad (7)$$

$$\mathbf{X}_{(3)} = \begin{bmatrix} 1 & 2 & 3 & 4 & 5 & 6 \\ 7 & 8 & 9 & 10 & 11 & 12 \end{bmatrix} \quad (8)$$

一个  $N$  阶张量  $\mathbf{X} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$  与一个长度为  $I_n$  的向量  $\mathbf{v}$  的模  $n$  乘 (Mode- $n$  Multiplication) 记为  $\mathbf{X} \times_n \mathbf{v}$ , 计算结果是一个  $N-1$  阶张量  $\mathbf{Y} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_{n-1} \times I_{n+1} \times \dots \times I_N}$ , 即对  $\mathbf{X}$  的每个模  $n$  纤维与  $\mathbf{v}$  执行内积运算而获得. 例如, 式(5)中表示的 3 阶张量  $\mathbf{X}$  与向量  $\mathbf{v} = \{1, 2\}$  的模 1 乘积为

$$\mathbf{X} \times_1 \mathbf{v} = \begin{bmatrix} 5 & 23 \\ 11 & 27 \\ 17 & 30 \end{bmatrix} \quad (9)$$

## 2.2 异构编程模型

随着机器学习和数据挖掘等实际应用中计算负载的剧增, 越来越多的计算系统配备了 GPU 加速器. GPU 通常包含了超过 2000 个计算核(线程), 远远超过了现有的多核 CPU 处理器. 在典型的 CPU-GPU 异构构架中, CPU 和 GPU 都配备了各自的私有内存, 分别称为主机内存(Host Memory)和设备内存(Device Memory), 并通过 PCIe (Peripheral Component Interconnect-Express)互联互通.

统一计算设备架构 CUDA (Compute Unified

Device Architecture) 为开发人员提供了利用 NVIDIA GPU 计算资源的并行编程框架. CUDA 程序会将运行的 kernel 函数中设置的线程块调度到 GPU 的流多处理器(Stream Multiprocessor, SM)上并发进行. CUDA 处理流程主要包括 4 个步骤:

- (1) 将数据从主机内存复制到设备内存;
- (2) CPU 启动 GPU 上的 kernel 函数;
- (3) GPU 的 CUDA 线程并行处理 kernel 函数;
- (4) 将设备内存上的计算结果传回主机内存.

本文针对一个  $N$  阶稀疏张量  $\mathbf{X} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$ , 介绍其 SpTV 模  $n$  运算在大规模异构众核 CPU-GPU 集群上的高效算法, 其中  $n \in \{1, 2, \dots, N\}$ .

## 3 多维并行 SpTV 划分方法

大规模异构众核 CPU-GPU 计算系统提供了多级并行计算能力: 大量的 CPU-GPU 节点提供了节点级并行计算能力, 每个节点中的 GPU 能够提供线程级并行计算能力. 为了充分利用计算系统的多级并行性, 本文对稀疏张量  $\mathbf{X}$  的 SpTV 模  $n$  运算提出了一种多维并行划分方法.

多维并行 SpTV 划分方法主要包括面向节点级并行的  $N-1$  维张量划分和面向节点内 GPU 线程级并行的模  $n$  展开矩阵划分两个主要的部分.

### 3.1 面向节点级并行的 $N-1$ 维张量划分

如第 2.1 节所述, SpTV 模  $n$  运算是依次对输入张量  $\mathbf{X}$  的每个模  $n$  纤维与输入向量  $\mathbf{v}$  执行内积运算. 在对运算进行面向  $\alpha$  个计算节点的并行数据划分时, 如果在模  $n$  切片的维度将  $\mathbf{X}$  划分为  $\alpha$  个子张量 (其中第  $i$  个子张量块记为  $\text{subX}_i \in \mathbb{R}^{I_1, I_2, \dots, I_{n-1}, I_n^{(i)}, I_{n+1}, \dots, I_N}$  ( $i = \{1, 2, \dots, \alpha\}$ ), 并且

$$I_n = \sum_{i=1}^{\alpha} I_n^{(i)},$$

那么每个节点上的计算结果需进一步传输到主节点进行最后的累加操作, 这样的划分方法造成了节点间通信和额外的计算操作. 为了解决这个问题, 面向节点级并行的  $N-1$  维张量划分方法分别在张量模  $m$  切片的维度对  $\mathbf{X}$  进行  $N-1$  个维度的划分, 其中  $m = \{1, 2, \dots, n-1, n+1, \dots, N\}$ , 那么每个节点上获得的是结果张量的一个相应分块, 避免了节点间通信和额外的累加计算.

具体而言,面向节点级并行的 $N-1$ 维张量划分方法根据 $X$ 各个维度的大小和计算节点的数量(记为 $\alpha$ ),将 $N$ 阶稀疏张量 $X$ 划分为 $\alpha$ 个 $N$ 阶子张量块,其中第 $i$ 个子张量块记为 $subX_i \in \mathbb{R}^{I_1^{(i)} \times I_2^{(i)} \times \dots \times I_{n-1}^{(i)} \times I_n \times I_{n+1}^{(i)} \times \dots \times I_N^{(i)}} (i=\{1, 2, \dots, \alpha\})$ ,并且 $I_1^{(1)}=I_1^{(2)}=\dots=I_1^{(\alpha)}$ ,  $I_2^{(1)}=I_2^{(2)}=\dots=I_2^{(\alpha)}$ , ...,  $I_N^{(1)}=I_N^{(2)}=\dots=I_N^{(\alpha)}$ . 每个计算节点中的CPU读取一个子张量块 $subX_i$ 至主存,该节点负责对 $subX_i$ 与输入向量 $v$ 进行SpTV模 $n$ 乘运算,结果为一个 $N-1$ 维张量块 $subY_i \in \mathbb{R}^{I_1^{(i)} \times I_2^{(i)} \times \dots \times I_{n-1}^{(i)} \times I_{n+1}^{(i)} \times \dots \times I_N^{(i)}}$ .

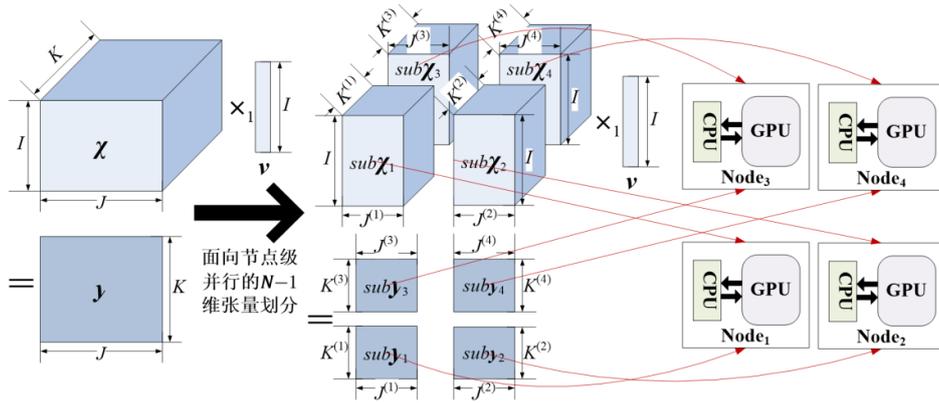


图4 面向节点级并行的 $N-1$ 维张量划分,其中假设计算节点数量设为 $\alpha=4$

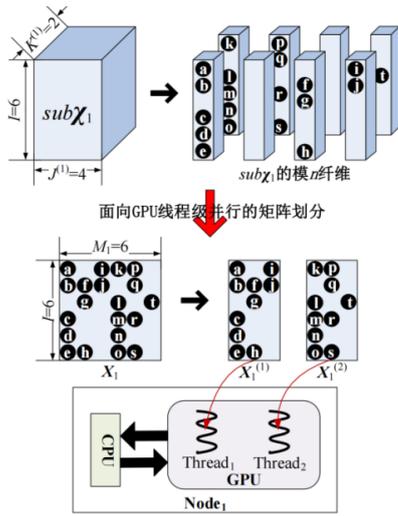


图5 图4所示第1个计算节点内面向GPU线程级并行的矩阵划分示例,其中假设该节点内的GPU线程数设为 $\beta=2$

### 3.2 面向GPU线程级张量划分

如第2.2节所述,为了利用GPU资源,首先需要将输入张量和向量数据从CPU主机内存传输到GPU设备内存,接着将计算任务分配给GPU线程,待GPU上的计算完成后,计算结果需从设备内存

图4以一个3阶张量 $X \in \mathbb{R}^{I \times J \times K}$ 与向量 $v \in \mathbb{R}^I$ 的SpTV模1乘为例,其中计算节点数量设为 $\alpha=4$ .如图所示,分别在模2切片和模3切片的维度将 $X$ 划分为4个子张量 $subX_1 \in \mathbb{R}^{I \times J^{(1)} \times K^{(1)}}$ ,  $subX_2 \in \mathbb{R}^{I \times J^{(2)} \times K^{(2)}}$ ,  $subX_3 \in \mathbb{R}^{I \times J^{(3)} \times K^{(3)}}$ 和 $subX_4 \in \mathbb{R}^{I \times J^{(4)} \times K^{(4)}}$ ,其中 $J^{(1)}=J^{(2)}=J^{(3)}=J^{(4)}$ ,  $K^{(1)}=K^{(2)}=K^{(3)}=K^{(4)}$ .每个计算节点对一个子张量块和向量 $v$ 执行SpTV运算.

传回主机内存.此外,根据SpTV的模 $n$ 乘运算属性,张量中每个模 $n$ 纤维与向量 $v$ 的内积结果对应的是结果张量 $X$ 中的一个元素.但稀疏张量中经常存在空模 $n$ 纤维,即所有元素都为零的模 $n$ 纤维,在运算中这些空纤维对应于结果张量 $Y$ 中的元素也为零.为了减少主机内存与设备内存之间的数据传输量、更均匀地为GPU划分计算任务,本文提出了面向GPU线程级并行的矩阵划分方法,通过删除空纤维、降低张量的稀疏度,减少主机内存与设备内存之间的数据传输量,同时更均匀地划分输入张量.

本文设计面向GPU线程级并行的矩阵划分方法,为每个计算节点内的GPU线程划分并行任务.该划分方法主要包括3个步骤:

(1) 将分配给每个计算节点的张量数据块 $subX_i$ 展开为模 $n$ 展开矩阵;

(2) 将该模 $n$ 展开矩阵中所有元素都为零的模 $n$ 纤维(称为空纤维)删除,并记该矩阵为 $X_i \in I_n \times M_i$ ;

(3) 根据  $X_i$  中矩阵列的数量( $M_i$ )和 GPU 线程的数量(记为  $\beta$ ), 继续将  $X_i$  在矩阵列的维度进行均匀地划分, 得到  $\beta$  个矩阵块, 每个矩阵块中包含了  $M_i/(\alpha \times \beta)$  列. 本文记  $X_i$  中的第  $j$  个矩阵块为  $X_i^{(j)}=I_n \times M_i/(\alpha \times \beta)$ , 其中  $j \in \{1, 2, \dots, \beta\}$ . 第  $i$  个计算节点中每个 GPU 线程对  $X_i^{(j)}$  与输入向量  $v$  进行 SpTV 模  $n$  乘运算.

图 5 展示了图 4 所示第 1 个计算节点内面向 GPU 线程级并行的矩阵划分示例, 其中该节点内的 GPU 线程数设为  $\beta=2$ . 首先  $subX_1$  展开为模 1 展开矩阵  $X_1 \in \mathbb{R}^{I \times M_1}$ , 其中  $I=6$  并且  $M_1=6$ ; 接着对  $X_1$  按列划分为  $\beta=2$  个矩阵块  $X_1^{(j)}$ , 每个  $X_1^{(j)}$  中包含了 3 列. 该计算节点内的每个 GPU 线程对一个矩阵块和向量  $v$  进行 SpTV 运算.

## 4 稀疏张量纤维压缩存储格式

稀疏张量的存储取决于数据的划分策略和计算模式, 适配于划分策略的张量存储结构能够提高计算节点和线程对数据块的索引效率. 根据本文提出的多维并行 SpTV 划分方法, 各计算节点和 GPU 线程都是按模  $n$  纤维对张量数据进行访问. 因此, 本文采用一种稀疏张量纤维压缩存储格式, 按模  $n$  纤维对张量中的非零元素进行存储, 从而适配设计的多维并行 SpTV 划分方法划分策略、提高数据索引效率以及减少内存占用.

稀疏张量纤维压缩存储格式采用 4 个数组对每个计算节点上的子稀疏张量块  $subX_i$  进行存储:

- (1) 存储  $subX_i$  中每个非零模  $n$  纤维中非零元素数量的数组  $P_f$ ;
- (2) 存储每个非零元素在相应模  $n$  纤维中的索引的数组  $Ind$ ;
- (3) 存储每个非零元素数值的数组  $Val$ ;
- (4) 存储  $subX_i$  的每个模  $m$  切片中非零模  $n$  纤维数量的数组  $N_f(m \in \{1, 2, \dots, n-1, n+1, \dots, N\})$ .

以图 5 中的子张量块  $subX_1$  为例, 稀疏张量纤维压缩存储格式使用 4 个数组存储  $subX_1$ :

(1)  $P_f[10]=\{0, 5, 5, 8, 10, 15, 19, 19, 20\}$  存储非零元素的数量, 其中  $P_f[r]-P_f[r-1]$  表示第  $r$  个非零模 1 纤维中非零元素的数量;

(2)  $Ind[20]=\{1, 2, 4, 5, 6, 2, 3, 6, 1, 2, 1, 3, 4, 5, 6, 1, 2, 4, 6, 3\}$  存储 20 个非零元素在模 1 纤维中的索引;

(3)  $Val[20]=\{a, b, c, d, e, f, g, h, i, j, k, l, m, n, o, p, q, r, s, t\}$  存储 20 个非零元素的数值;

(4)  $N_f[3]=\{0, 3, 6\}$  存储  $subX_1$  的每个模 3 切片中非零模 1 纤维的数量, 其中  $N_f[s]-N_f[s-1]$  表示第  $s$  个模 3 切片中非零模 1 纤维的数量.

基于张量纤维的多维并行 SpTV 算法如算法 1 所述, CPU 首先为输入稀疏张量  $X$ 、输入向量  $v$  和输出张量  $Y$  在主机内存中开辟存储空间(第 1-6 行), 再读取输入数据  $X$  和  $v$  (第 7 行), 并采用稀疏张量纤维压缩存储格式对稀疏张量  $X$  进行存储(第 8 行). 接着为  $X$ 、 $v$  和  $Y$  在 GPU 设备内存中开辟存储空间(第 9-13 行), 再将存储稀疏张量  $X$  的 3 个数组  $P_f$ 、 $Ind$  和  $Val$  以及存储输入向量的数组  $v$  从主机内存传到设备内存(第 14-17 行). GPU 获得数据后, 开始执行并行 SpTV 运算, 计算结果存储到设备内存中的数组  $Y_{Device}$  中(第 18 行). 并行计算结束后, 计算结果再从设备内存传输至主机内存中(第 19 行).

**算法 1.** 基于张量纤维的多维并行 SpTV 算法.

输入: 输入稀疏张量  $X$ , 输入向量  $v$ , GPU 并行线程数  $\beta$   
输出: 输出张量  $Y$

1. `cudaHostAlloc(P_f, ...)`; //在 CPU 主机内存中分配  $X$  的存储空间
2. `cudaHostAlloc(Ind, ...)`;
3. `cudaHostAlloc(Val, ...)`;
4. `cudaHostAlloc(N_f, ...)`;
5. `cudaHostAlloc(v, ...)`; //在 CPU 主机内存中分配  $v$  的存储空间
6. `cudaHostAlloc(Y, ...)`; //在 CPU 主机内存中分配  $Y$  的存储空间
7. 读取输入数据  $X$  和  $v$ ;
8. 以稀疏张量纤维压缩存储格式将稀疏张量  $X$  存储到主机内存中;
9. `cudaMalloc(P_f_Dev, ...)`; //在 GPU 设备内存中分配  $X$  的存储空间
10. `cudaHostAlloc(Ind_Dev, ...)`;

- 10. `cudaHostAlloc(Val_Dev, ...)`;
- 11. `cudaHostAlloc(v_Dev, ...)`; //在 GPU 设备内存中分配  $v$  的存储空间
- `cudaHostAlloc(Y_Dev, ...)`; //在 GPU 设备内存中分配  $Y$  的存储空间
- 12. `cudaMemcpy(Pf, Pf_Dev, ..., HostToDevice)`; //将并行计算所需数据从主机内存传输到设备内存中
- 13. `cudaMemcpy(Ind, Ind_Dev, ..., HostToDevice)`;
- 14. `cudaMemcpy(Val, Val_Dev, ..., HostToDevice)`;
- 15. `cudaMemcpy(v, v_Dev, ..., HostToDevice)`;
- 16. `SpTV( $\beta, Pf_Dev, Ind_Dev, Val_Dev, v_Dev, Y_Dev, ...$ )`;  
//GPU 上的并行 SpTV
- 17. `cudaMemcpy(Y_Dev, Y, ..., DeviceToHost)`; //将计算结果从设备内存传输到主机内存

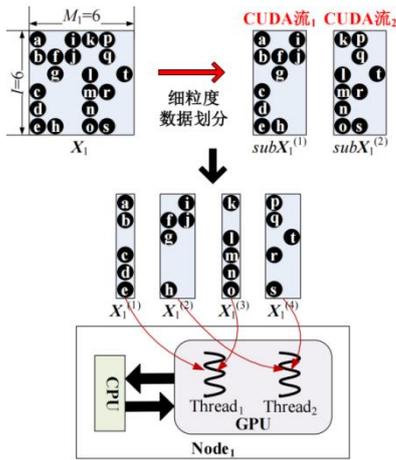


图6 图4所示第1个计算节点内多流并行运行机制中的细粒度数据划分示例,其中假设GPU线程数为 $\beta=2$ 以及CUDA流的数量为2

## 5 基于张量块排序的多流并行优化

根据算法1所述,CUDA处理流程主要包括数据传输和并行计算两个主要部分,其中数据传输部分包括从主存加载数据到设备内存以及从设备内存返回结果数据到主存,并行计算部分为GPU上进行的计算.对于CPU-GPU构架,CPU主机端发出CUDA操作的命令后,不会等待该命令执行完毕,而是立刻执行后续命令,这个特征支持了kernel函数中数据传输和并行计算两个主要部分之前的并行化.为了利用这种软件可并行性,本文基于细粒度数据划分,采用多流并行运行机制,在每个CPU-GPU计算节点上用多个CUDA流处理细粒度数据块,在并行SpTV中的数据运输部分和并行计算部分之间产生并行性,从而使通信时间与计算时间相互隐藏,减少并行SpTV的运行时间.

同时,由于SpTV中张量数据的稀疏性,各细粒度数据块的内存占用不均匀,因此,CUDA流每次对各数据块的处理时间存在差异,从而也导致CUDA流之间通信时间与计算时间相互隐藏的效果不佳.为了缓解这个问题,本文进一步提出了基于张量块排序的多流并行优化技术,通过调整各数据块的处理顺序,提高CUDA流之间通信与计算之间的重叠程度,获得更好的多流并行优化效果.

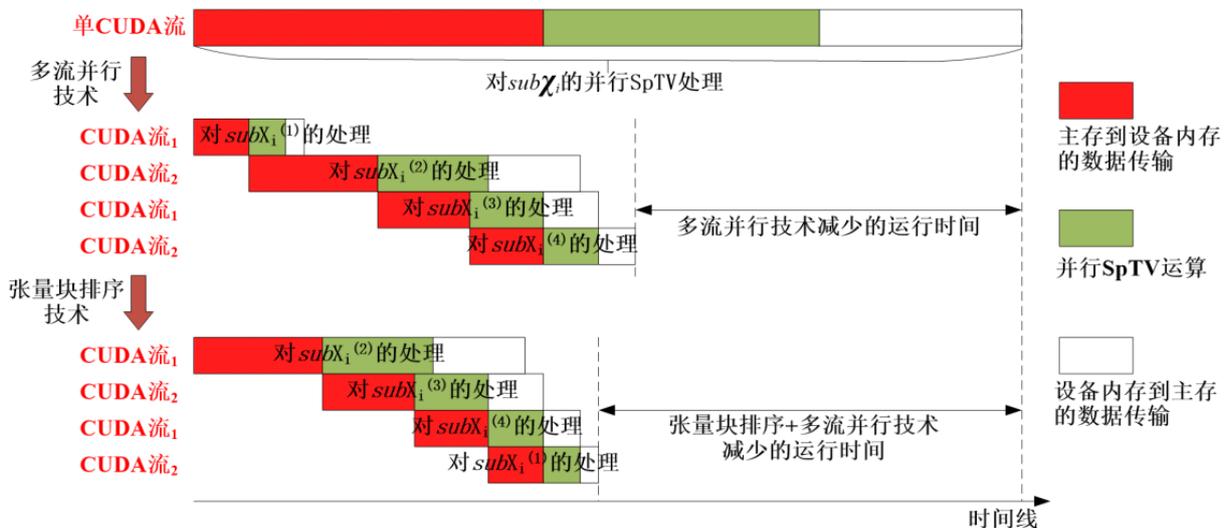


图7 基于多流并行的异构高效优化策略示例,其中假设CUDA流的数量为2

### 5.1 细粒度数据划分

为了保证每个 CUDA 流都被分配了并行计算任务,提高多流处理效率,多流并行运行机制需要对张量数据进行更细粒度的划分.细粒度数据划分进行在面向 GPU 线程级并行划分的步骤 1 之后,将模  $n$  展开矩阵  $X_i$  在矩阵列的维度按非零元素的个数划分为  $P$  个子矩阵块,记为  $subX_i^{(p)}$ ,每个  $subX_i^{(p)}$  包含  $M_i/P$  列,其中  $p=\{1, 2, \dots, P\}$ .细粒度数据划分完成后,再分别对每个  $subX_i^{(p)}$  进行面向 GPU 线程级并行划分中的步骤 2 操作,将  $subX_i^{(p)}$  进一步划分为  $\beta$  个细粒度矩阵块  $X_i^{(j)}$ ,每个  $X_i^{(j)}$  包含  $M_i/(P \times \beta)$  列.因此,模  $n$  展开矩阵  $X_i$  最终被划分为  $P \times \beta$  个矩阵块  $X_i^{(j)}$ ,其中  $j \in \{1, 2, \dots, P \times \beta\}$ ,在第  $j \% P$  个 CUDA 流的任务中,第  $j/P$  个 GPU 线程对  $X_i^{(j)}$  进行 SpTV 运算.

以图 5 中所展示的面向 GPU 线程级并行的张量划分为例,图 6 展示了第 1 个计算节点内面向多 CUDA 流并行的细粒度数据划分,其中假设 GPU 线程数为  $\beta=2$  以及 CUDA 流的数量为 2.如图 6 所示,细粒度数据划分将张量的模  $n$  展开矩阵  $X_1$  划分为  $P=2$  个子矩阵块( $subX_1^{(1)}$  和  $subX_1^{(2)}$ ),第一个 CUDA 流负责对  $subX_1^{(1)}$  进行并行处理,第二个 CUDA 流负责对  $subX_1^{(2)}$  进行并行处理.接着,进一步将  $subX_1^{(1)}$  划分为  $\beta=2$  个细粒度矩阵块( $X_1^{(1)}$  和  $X_1^{(2)}$ ),将  $subX_1^{(2)}$  划分为  $\beta=2$  个细粒度矩阵块( $X_1^{(3)}$  和  $X_1^{(4)}$ ),每个细粒度矩阵块中包含了差不多数量的矩阵列.因此,在第一个 CUDA 流的任务中,第一个 GPU 线程对  $X_1^{(1)}$  进行计算,第二个 GPU 线程对  $X_1^{(2)}$  进行计算;在第二个 CUDA 流的任务中,第一个 GPU 线程对  $X_1^{(3)}$  进行计算,第二个 GPU 线程对  $X_1^{(4)}$  进行计算.

### 5.2 多流并行运行机制

在多流运行机制中,一个 CUDA 流对一个子矩阵块  $subX_i^{(p)}$  进行并行 SpTV 运算处理的同时,另一个 CUDA 流从主存加载下一个子矩阵块  $subX_i^{(p+1)}$  到 GPU 设备内存中.因此,多流运行机制在对  $subX_i^{(p)}$  进行的并行计算任务与传输  $subX_i^{(p+1)}$  的通信任务之间生成了并行性,并行 SpTV 算法的总体性能获得提高.

每个计算节点使用单个 CUDA 流时,并行 SpTV 处理主要包括 3 个步骤:1) 主存至设备内存:

主存中的张量数据  $subX_i$  首先被传输至 GPU 设备内存,2) 并行计算:接着 GPU 启用并行线程对  $subX_i$  和向量  $v$  进行并行 SpTV 模  $n$  乘运算,3) 设备内存至主存:最后将设备内存中的结果张量块  $subY_i$  传输回主存中.图 7 显示了第  $i$  个计算节点内的 GPU 多流并行处理机制的示例,其中 CUDA 流的数量设为 2,且在每个 CUDA 流的每次执行过程中,“主存到设备内存的数据传输”、“并行 SpTV 运算”和“设备内存到主存的数据传输”三个步骤的耗时并不相等.在使用多流并行处理机制时,张量数据  $subX_i$  被划分为  $P$  个子矩阵块  $subX_i^{(p)}$ ,启用的 2 个 CUDA 流轮流对  $subX_i^{(p)}$  进行主存至设备内存、并行计算和设备内存至主存的并行 SpTV 处理.如图 7 所示,除了将  $subX_i^{(1)}$  从主存加载到设备内存的通信无法与并行计算部分相互重叠,对其他子矩阵块的通信都可以与并行计算之间产生并行性,因此减少了并行 SpTV 的运行时间.

### 5.3 基于张量块排序的多流并行优化技术

在多流并行运行机制中,各数据块  $subX_i^{(p)}$  的处理时间存在差异,因此,对各数据块的处理顺序影响了 CUDA 流之间通信开销与计算开销的重叠效果.为了获得更好的多流并行优化效果,本文进一步提出了一种基于张量块排序的多流并行优化技术,通过调整各数据块  $subX_i^{(p)}$  的处理顺序,提高 CUDA 流之间通信开销与计算开销的并行度.

如图 7 中示例所示,4 个矩阵块的处理实现存在差异.只采用多流并行技术时,4 个矩阵块的处理顺序为:  $subX_i^{(1)}$ ,  $subX_i^{(2)}$ ,  $subX_i^{(3)}$ ,  $subX_i^{(4)}$ .张量块排序技术将 4 个矩阵块的处理顺序按照各数据块的处理时间从大到小进行排序,调整后的处理顺序为:  $subX_i^{(2)}$ ,  $subX_i^{(3)}$ ,  $subX_i^{(4)}$ ,  $subX_i^{(1)}$ .如图所示,张量块排序技术能够进一步优化多流并行运行机制的性能.

由于各数据块的处理时间取决于非零元素的数量,因此张量块排序技术按照各数据块中非零元素数量从大到小的顺序对数据块的处理顺序进行排序,提高并行运行机制的优化效果.

表 1 数据集

来源	数据集	$I_1$	$I_2$	$I_3$	$I_4$	NNZ
----	-----	-------	-------	-------	-------	-----

MovieLens	ratings-m1	6K	4K	5	/	1M
MovieLens	ratings-m20	72K	131K	10	/	1M
MovieLens+IMDb/Rotten Tomatoes	user_ratedmovies	72K	65K	10	/	856K
Last.fm	user_taggedartist	2K	19K	13K	/	186K
Kaggle	Submissions	3M	347K	2	/	2M
NYC Taxi & Limousine Commission (TLC)	Uber_Pickups	2K	1K	24	183	3M

## 6 性能测试与分析

### 6.1 实验环境

本文基于异构并行计算构架对多级并行的 SpTV 算法性能进行测试与分析, 实验采用的异构并行计算平台中包含了一个 14 核的 E5-2680 CPU@2.40GHz 和一个内存为 16GB 的 NVIDIA Tesla P100 GPU. 所有的实验中, 多级并行 SpTV 算法测试的是张量与向量的模 1 乘积.

多级并行的 SpTV 算法在 6 个多阶的真实张量数据集上进行测试, 包括 5 个 3 阶张量和 1 个 4 阶张量. 表 1 中展示了这 6 个数据集的信息, 其中  $I_1$ 、 $I_2$ 、 $I_3$  和  $I_4$  分别表示数据集在 4 个维度上的维度大小, NNZ 表示数据集中非零元素的个数. 数据集 ratings-m1 和 ratings-m20 来自电影推荐系统 MovieLens 中的用户对电影的评分数据, 进而形成表示用户-电影-评分的 3 阶张量<sup>[28]</sup>. 数据集 user\_ratedmovies 结合了 MovieLens 数据集及其相应的网页 Internet Movie Database (IMDb)<sup>1</sup> 和电影评价系统 Rotten Tomatoes<sup>2</sup> 中的真实数据, 构建了用户-电影-评分的 3 阶张量. 数据集 user\_taggedartists 包含了在线音乐系统 Last.fm<sup>3</sup> 中 2 千个用户对音乐艺术家的定义的标签数据, 即一个表示用户-音乐艺术家-标签的 3 阶张量. 数据集 Submissions 来自 Kaggle<sup>4</sup> 社区和活动中收集的丰富数据(包括竞赛、用户、提交分数和内核的公共数据). 数据集 Uber\_Pickups 来自 NYC Taxi&Limousine Commission (TLC)<sup>5</sup>, 是一个 4 阶张量, 包含了纽约市 uber 订单发生的日期、时间、经度和纬度信息.

### 6.2 实验结果

为了分析多级并行 SpTV 算法在异构并行计算平台上的整体优化效果, 图 8 展示了与串行 SpTV 运算相比, 多级并行 SpTV 算法的并行加速比. 多级并行 SpTV 算法在所有测试数据集上平均获得了 7.34 倍并行加速比(最高在 Uber\_Pickups 上获得 10.60 倍并行加速比, 最低在 ratings-m1 上获得 5.33 倍并行加速比).

#### 6.2.1 多维并行 SpTV 划分的评估

为了测试多级并行 SpTV 算法对并行计算节点数量的可扩展性, 图 10 展示了算法在不同计算节点数量上的加速比. 实现结果显示, 与采用 1 个计算节点相比, 多级并行 SpTV 算法在 8 个计算节点上平均获得了 3.57 倍加速比, 其中, 在数据集 user\_taggedartists 和 Uber\_Pickups 上获得更高的加速比(分别为 3.95 倍和 4.58 倍), 在 Submissions 上表现出较差的可扩展性(1.64 倍加速比). 这是因为对于 SpTV 模  $n$  运算, 面向节点级并行的  $N-1$  维张量划分不会在模  $n$  切面的维度对张量进行划分, 即  $I_n$  的大小始终保持不变; 同时, 在每个计算节点中, CPU 主存到 GPU 设备内存的传输数据都包括了长度为  $I_n$  的输入向量  $\mathbf{v}$ , 在进行主存到设备内存的数据传输时, 除了向量  $\mathbf{v}$  以外, 其他数据的传输量都随着计算节点的增加而减少; 因此, 即使面向节点级并行的  $N-1$  维张量划分避免了节点间通信, 但在每个计算节点内, 对向量  $\mathbf{v}$  的数据传输开销是限制节点级并行可扩展性的主要原因. 根据表 1 中的数据, user\_taggedartists 和 Uber\_Pickups 的  $I_n$  值很小(实验测试的是 SpTV 模 1 乘积, 即  $n=1$ ), 即传输向量  $\mathbf{v}$  的开销较小, 从而对节点级并行可扩展性的限制较小; 然而 Submissions 的  $I_n$  值很大, 传输向量  $\mathbf{v}$  的开销在总运算开销中的占比很大, 因此限制了多级并行 SpTV 算法在该数据集上的节点级并行可扩展性.

1 <http://www.imdb.com/>

2 <http://www.rottentomatoes.com/>

3 <http://www.last.fm/>

4 <https://www.kaggle.com/>

5 <https://www.nyc.gov/site/tlc/index.page/>

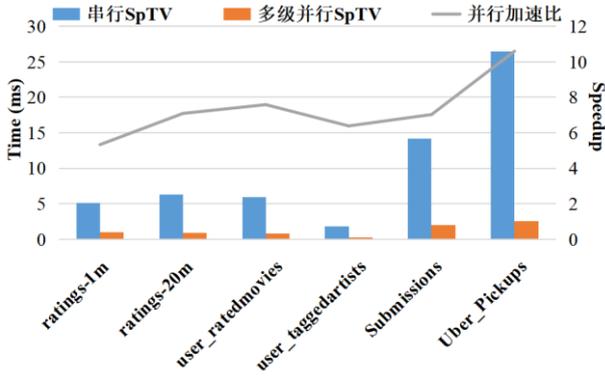


图 8 多级并行 SpTV 算法的并行加速比

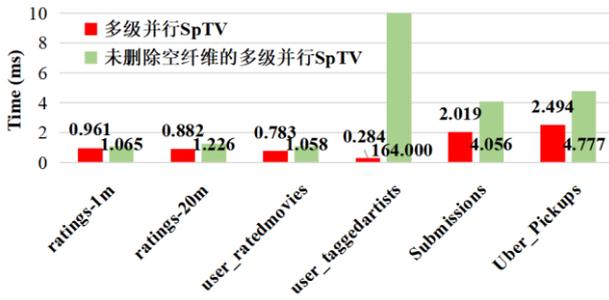


图 9 面向 GPU 线程级并行的矩阵划分中删除空纤维操作对多级并行 SpTV 性能的影响

为了研究面向 GPU 线程级并行的矩阵划分中删除空纤维操作对多级并行 SpTV 性能的影响，表

表 2 面向 GPU 线程级并行的矩阵划分中删除空纤维操作对张量展开矩阵  $X_i$  大小和密度的影响

数据集	$I_n \times M_i$ (不删除空纤维)	$NNZ/(I_n \times M_i)$ (不删除空纤维)	$I_n \times M_i$ (删除空纤维)	$NNZ/(I_n \times M_i)$ (不删除空纤维)
ratings-m1	6040×19760	8.38E-03	6040×16912	9.79E-03
ratings-m20	7120×1306420	1.13E-04	7120×70147	2.10E-03
userRatedmovies	71534×651330	1.84E-05	71534×63015	1.90E-04
userTaggedartist	2100×23705536	3.75E-07	2100×109750	8.09E-04
Submissions	3253292×693430	7.08E-07	1596660×121318	8.24E-06
Uber_Pickups	1717×5006880	3.85E-04	1717×692597	2.78E-03

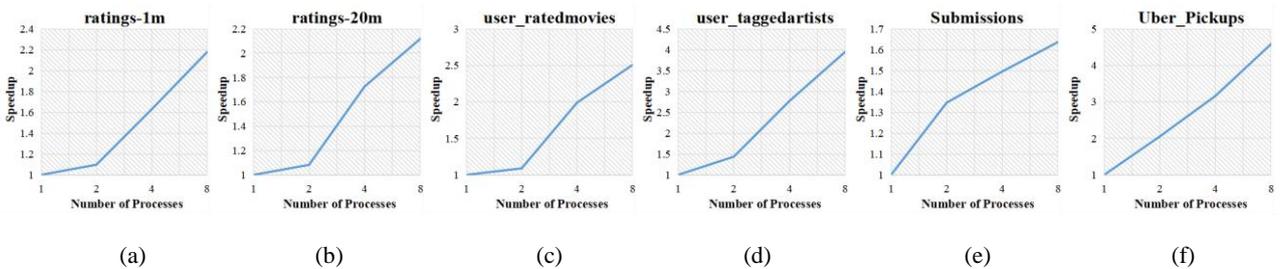


图 10 多级并行的 SpTV 算法对计算节点数量的可扩展性

2 展示了删除空纤维操作对张量展开矩阵  $X_i$  大小的影响,图 9 展示了删除空纤维操作对多级并行 SpTV 运行时间的影响. 删除空纤维操作能够使多级并行 SpTV 算法的运行时间平均减少 43.61%, 其中, 在数据集 user\_taggedartists 上的优化效果最佳(运行时间减少 99.83%), 在数据集 ratings-m1 上的优化效果最差(运行时间减少 9.77%). 根据表 2 中的数据, 删除空纤维操作缩减了 user\_taggedartists 展开矩阵  $X_i$  在行向和列向两个维度的大小,  $X_i$  的密度增加, 非零元素分布的稀疏度降低. 一方面, 行向维度大小(即  $I_n$ )的缩减能够相应地缩短向量  $v$  的长度, 同时, 根据本文所采用的稀疏张量纤维压缩存储格式, 列向维度大小(即  $M_i$ )的缩减能相应地减小数组  $P_f$  的大小和结果张量的大小, 从而减少了主机内存与设备内存之间的数据通信量. 另一方面, 删除空纤维使得 user\_taggedartists 展开矩阵的密度大大增加, 非零元素能够更集中地分布在矩阵中, 进而在进行多维并行划分时, 获得更好的并行负载均衡. 因此, 删除空纤维操作在数据集 user\_taggedartists 上的优化效果最佳. 同理, 相比之下, 删除空纤维操作在 ratings-m1 上的优化效果就没有那么明显.

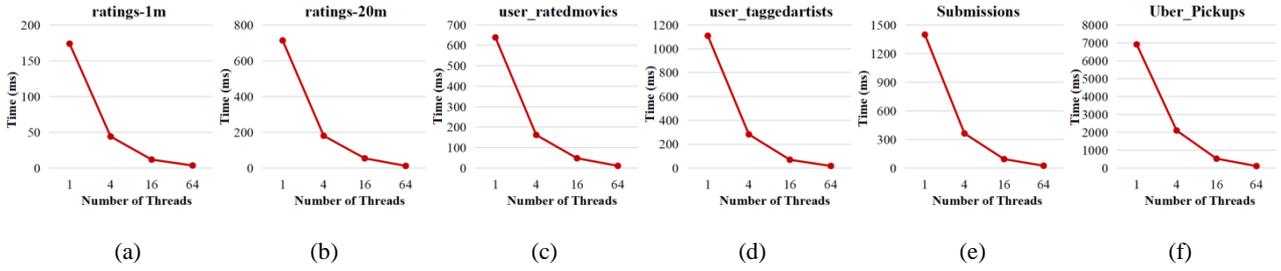


图 11 多级并行的 SpTV 算法对 GPU 并行线程数量的可扩展性

为了测试多级并行 SpTV 对 GPU 并行线程数量的可扩展性, 图 11 展示了 GPU 并行线程数量分别为 1、4、16 和 64 时算法的运行时间. 由于每个 GPU 线程的计算任务相互独立, 且没有数据通信, 因此图 11 显示, 多级并行 SpTV 在 GPU 线程上表现出良好的并行可扩展性. 与采用 1 个 GPU 线程相比, 采用 64 个 GPU 并行线程的多级并行 SpTV 算法的运行时间在 6 个测试数据集上平均获得了 56.91 倍的加速比.

### 6.2.2 多流并行异构优化策略的评估

图 12 比较了未采用基于张量块排序的多流并行技术、仅采用多流并行技术与采用了基于张量块排序的多流并行技术的多级并行 SpTV 算法的运行时间. 实验结果证明, 基于张量块排序的多流并行策略能够对多级并行 SpTV 产生良好的优化效果. 与未采用基于张量块排序的多流并行技术相比, 仅采用多流并行技术的算法运行时间平均减少了 40.09%. 其中, 在数据集 `user_taggedartists` 和 `Submissions` 上的优化效果较差, 运行时间减少 16.80% 和 27.14%. 从图 7 中分析其原因, 多流并行技术重叠了通信时间和并行计算时间, 当相互重叠的通信时间与并行计算时间相差较小时, 多流并行策略的重叠效果更好; 相反, 则优化效果更差. 同时, 根据算法设计, 通信时间取决于  $I_n$ ,  $M_i$  和  $NNZ$  的大小, 计算时间取决于每个线程处理的非零元数量. 如表 2 所示, 对于 `user_taggedartists` 和 `Submissions`, 一方面, 这两个数据集的  $I_n$  和  $M_i$  非常大, 这导致通信时间在总计算时间里的占比很大; 另一方面, 这两个数据集的密度相对较小(删除空纤维后), 非零元素的分布更加稀疏, 这导致每个线程处理的非零元数量较少, 计算时间在总计算时间里的占比较小. 这两方面的原因造成通信时间和计算时间相差较大, 因此多流并行策略的优化效果较差.

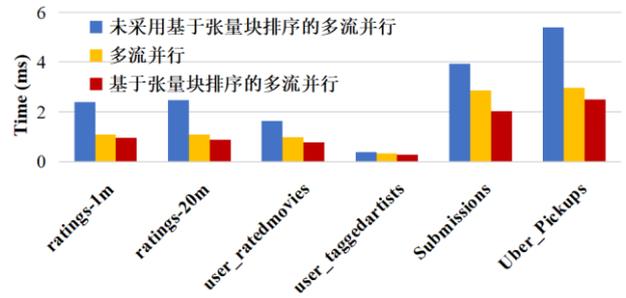


图 12 基于张量块排序的多流并行策略的优化效果

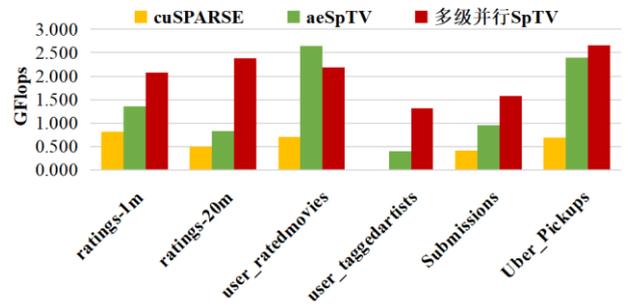


图 13 多级并行 SpTV 算法与现有相关工作的吞吐量对比. 与仅采用多流并行技术的算法相比, 采用了基于张量块排序的多流并行技术的算法在运行时间上平均减少了 17.48%. 其中, 在数据集 `Submissions` 上的优化效果最佳(运行时间减少 29.65%). 根据表 2 中的数据, `Submissions` 是所有数据集中非零元素分布最稀疏的数据集(删除空纤维后), 稀疏度较大时, 更可能发生各张量分块  $subX_i^{(p)}$  中非零元的数量相差较大、对各张量分块处理时间相差较大的情况. 如图 7 所示, 当各张量块  $subX_i^{(p)}$  的大小相差较大时, 张量排序技术能够使较大张量块的计算、通信开销与多个较小张量块的计算、通信开销进行重叠, 从而获得更好的优化效果. 因此张量排序技术在 `Submissions` 上能表现出更好的优化效果.

### 6.2.3 与现有工作的对比

图 13 展示了多级并行的 SpTV 算法与 cuSPARSE 库<sup>6</sup>和 aeSpTV<sup>[14]</sup>在 6 个测试数据集上的

<sup>6</sup> <https://docs.nvidia.com/cuda/cusparse/>

吞吐量对比. 由于张量  $X$  与向量  $v$  的 SpTV 模  $n$  运算相当于对张量模  $n$  展开矩阵的转置矩阵  $X_i^{(n)}$  与  $v$  的稀疏矩阵乘向量 (Sparse Matrix-Vector Multiplication, SpMV) 运算, 且 cuSPARSE 提供了一组用于处理稀疏矩阵基本线性代数的算法, 由 NVIDIA 公司基于 CUDA 编程模型在 GPU 上设计与实现. 因此, 本实验将多级并行的 SpTV 运算与 cuSPARSE 库中提供的 SpMV 算法进行性能对比.

如图 13 所示, 由于多维并行 SpTV 划分方法和基于张量块排序的多流并行策略的优化作用, 多级并行的 SpTV 算法在所有测试数据集上的运行效率都优于 cuSPARSE 库, 特别是在数据集 user\_taggedartists 上获得了超高的加速比, 我们根据表 2 中的数据, 分析其主要原因是多级并行 SpTV 算法中的空纤维删除操作极大地减小了 user\_taggedartists 展开矩阵的列向维度 ( $M_i$ ), 同时也极大地减小了展开矩阵的列索引量、计算结果的大小、以及主机内存与设备内存之间的通信量, 因此与 cuSPARSE 库相比, 多级并行 SpTV 算法在该数据集上表现出极大的性能优势.

与现有工作 aeSpTV 相比, 多级并行 SpTV 算法的吞吐量平均获得 1.88 倍的加速比, 但在数据集 userRatedmovies 上, aeSpTV 的吞吐量略高于多级并行 SpTV 算法, 我们分析了可能的原因. 第一, aeSpTV 与多级并行 SpTV 的数据划分和存储方法不同, 多级并行 SpTV 按张量展开矩阵的列向维度对非零元素进行并行任务划分和压缩存储, 当想要利用更多 GPU 线程时, 不需要额外增加索引量和通信量; 然而 aeSpTV 按张量展开矩阵的行向维度对非零元素进行压缩存储, 因此张量细粒度分块中非空行的数量影响了 aeSpTV 中的通信效率, 当想要利用更多 GPU 线程时, 需要对张量进行更细粒度的划分, 随之而来的代价是细粒度分块中非空行的数量增多, 增加了 aeSpTV 中的通信量; 但对于数据集 userRatedmovies, 随着 GPU 线程数的增加, aeSpTV 并行计算效率的提升抵消了通信量增加所带来的负面影响. 第二, aeSpTV 采用了自适应张量划分选择器, 选择最优的数据划分方法, 从而获得更好的并行负载均衡.

## 7 总结

本文基于 CPU-GPU 异构并行计算构架设计了多级并行的高效 SpTV 算法, 为相关应用提供快速、

高效的运算基础. 首先, 本文设计了一种多维并行 SpTV 划分方法, 采用面向节点级并行的  $N-1$  维张量划分和面向 GPU 线程级并行的矩阵划分, 将 SpTV 运算映射到混合、多级并行的分布式 CPU-GPU 异构多/众核构架, 充分利用计算节点间和节点内的多级并行计算能力. 其次, 基于这种多维划分方法, 本文设计了一种基于稀疏张量纤维的压缩存储格式, 优化 SpTV 运算的计算和访存模式. 再次, 本文设计了基于多流并行的异构高效 SpTV 算法, 进一步采用稀疏张量的细粒度划分方法、多流并行运行机制和基于张量块排序的多流并行优化技术, 对并行 SpTV 运算中的通信开销与计算开销进行相互隐藏, 充分开发 GPU 的计算能力. 在 CPU-GPU 构架上的实验证明, 与现有相关工作 aeSpTV 相比, 并行高效 SpTV 算法在所有测试数据集上最高获得 3.28 倍的加速比.

致 谢 感谢所有评审人员的建议和帮助!

## 参 考 文 献

- [1] Poulernard A, Guibas L J. A functional approach to rotation equivariant non-linearities for tensor field networks//Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, Virtual Event, 2021: 13174-13183
- [2] Huang C. Ringcnn: Exploiting algebraically-sparse ring tensors for energy-efficient cnn-based computational imaging//Proceedings of the ACM/IEEE Annual International Symposium on Computer Architecture, ISCA 2021, Valencia, Spain, 2021: 1096-1109
- [3] Zheng Y, Huang T, Zhao X, et al. Fully-connected tensor network decomposition and its application to higher-order tensor completion//Proceedings of the Conference on Artificial Intelligence, AAAI 2021, Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, 2021: 11071-11078
- [4] Ioannidis V N, Zamzam A S, Giannakis G B, et al. Coupled graphs and tensor factorization for recommender systems and community detection. IEEE Transactions on Knowledge and Data Engineering, 2021, 33(3):909-920
- [5] Kwon Y, Lee Y, Rhu M. Tensor casting: Co-designing algorithm-architecture for personalized recommendation training//Proceedings of the IEEE International Symposium on High-Performance Computer Architecture, HPCA 2021, Seoul, South Korea, 2021: 235-248
- [6] Chou S, Jang J R, Yang Y. Fast tensor factorization for large-scale context-aware recommendation from implicit feedback. IEEE

- Transactions on Big Data, 2020, 6(1):201-208
- [7] Ho J C, Ghosh J, Sun J. Marble: high-throughput phenotyping from electronic health records via sparse nonnegative tensor factorization//Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD'14, New York, USA, 2014: 115-124
- [8] Ren Y, Lou J, Xiong L, et al. Robust irregular tensor factorization and completion for temporal health data analysis//Proceedings of the ACM International Conference on Information and Knowledge Management, Virtual Event, Ireland, 2020: 1295-1304
- [9] Wang Y, Chen R, Ghosh J, et al. Rubik: Knowledge guided tensor factorization and completion for health data analytics//Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Sydney, Australia, 2015: 1265-1274
- [10] Kolda T G, Bader B W. Tensor decompositions and applications. SIAM Review, 2009, 51(3):455-500
- [11] Ranavive T M, Baskaran M M. An all-at-once CP decomposition method for count tensors//Proceedings of the IEEE High Performance Extreme Computing Conference, HPEC 2021, Waltham, USA, 2021: 1-8
- [12] Won T, Park I, Lee D, et al. Slicenstitch: Continuous CP decomposition of sparse tensor streams//Proceedings of the IEEE International Conference on Data Engineering, ICDE 2021, Chania, Greece, 2021: 816-827
- [13] Karsavuran M O, Acer S, Aykanat C. Partitioning models for general medium-grain parallel sparse tensor decomposition. IEEE Transactions on Parallel Distributed Systems, 2021, 32(1):147-159
- [14] Chen Y, Xiao G, Özsü M T, et al. aesptv: An adaptive and efficient framework for sparse tensor-vector product kernel on a high-performance computing platform. IEEE Transactions on Parallel Distributed Systems, 2020, 31(10):2329-2345
- [15] Helal A E, Laukemann J, Checconi F, et al. ALTO: adaptive linearized storage of sparse tensors//Proceedings of the ACM International Conference on Supercomputing, ICS'21, Virtual Event, USA, 2021: 404-416
- [16] Xiao G, Li K, Chen Y, et al. Caspmv: A customized and accelerative spmv framework for the sunway taihulight. IEEE Transactions on Parallel Distributed Systems, 2021, 32(1):131-146
- [17] Qin E, Jeong G, Won W, et al. Extending sparse tensor accelerators to support multiple compression formats//Proceedings of the IEEE International Parallel and Distributed Processing Symposium, IPDPS 2021, Portland, USA, 2021: 1014-1024
- [18] Chen Y, Li K, Yang W, et al. Performance-aware model for sparse matrix-matrix multiplication on the sunway taihulight supercomputer. IEEE Transactions on Parallel Distributed Systems, 2019, 30(4):923-938
- [19] Bell N, Garland M. Implementing sparse matrix-vector multiplication on throughput-oriented processors//Proceedings of the ACM/IEEE Conference on High Performance Computing, SC'09, Portland, USA, 2009: 1-11
- [20] Rhu M, O'Connor M, Chatterjee N, et al. Compressing DMA engine: Leveraging activation sparsity for training deep neural networks//Proceedings of the IEEE International Symposium on High Performance Computer Architecture, HPCA 2018, Vienna, Austria, 2018: 78-91
- [21] Li J, Sun J, Vuduc R W. Hicoo: hierarchical storage of sparse tensors//Proceedings of the International Conference for High Performance Computing, Networking, Storage, and Analysis, SC 2018, Dallas, USA, 2018: 19:1-19:15
- [22] Smith S, Ravindran N, Sidiropoulos N D, et al. SPLATT: efficient and parallel sparse tensor-matrix multiplication//Proceedings of the IEEE International Parallel and Distributed Processing Symposium, IPDPS 2015, Hyderabad, India, 2015: 61-70
- [23] Smith S, Karypis G. Tensor-matrix products with a compressed sparse tensor//Proceedings of the Workshop on Irregular Applications - Architectures and Algorithms, IA3 2015, Austin, USA, 2015: 5:1-5:7
- [24] Abubaker N, Acer S, Aykanat C. True load balancing for matricized tensor times khatri-rao product. IEEE Transactions on Parallel Distributed Systems, 2021, 32(8):1974-1986
- [25] Dun M, Li Y, Yang H, et al. An optimized tensor completion library for multiple gpus//Proceedings of the International Conference on Supercomputing, ICS'21, Virtual Event, USA, 2021: 417-430
- [26] Liu J, Ren J, Gioiosa R, et al. Sparta: high-performance, element-wise sparse tensor contraction on heterogeneous memory//Proceedings of the ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming, PPOPP'21, Virtual Event, Republic of Korea, 2021: 318-333
- [27] Qiu Y, Zhou G, Zhang Y, et al. Canonical polyadic decomposition (CPD) of big tensors with low multilinear rank. Multimedia Tools and Applications, 2021, 80(15): 22987-23007
- [28] Harper F M, Konstan J A. The movielens datasets: History and context. ACM Transactions on Interactive Intelligent Systems, 2016, 5(4):19:1-19:19



**CHEN Yue-Dan**, Ph.D., associate professor. Her research interests include high-performance computing,

parallel and distributed processing, AI and big data computing.

**XIAO Guo-Qing**, Ph.D., associate professor. His research

interests mainly include high-performance computing and AI computing.

**YANG Wang-Dong**, Ph.D., professor. His research interest mainly is parallel computing.

**JIN Ji-Yong**, MS D., assistant professor. His research interest mainly is distributed computing.

**LONG Jun**, Ph.D., professor. His research interests include big data computing and intelligent software systems.

**LI Ken-Li**, Ph.D., professor. His research interests include system software and applications for high-performance computing.

## Background

Many applications give rise to multidimensional data that can be naturally represented via tensors. The tensors used in most real-world applications that are extremely large and very sparse. The sparse tensor decomposition is an effective approach to predict the unobserved data and is commonly used in machine learning, text analysis, healthcare analytics, and numerous other applications. SpTV is one of the most fundamental and time-intensive operations in computing tensor decomposition. Therefore, there have been a considerable number of researches on accelerating SpTV on heterogeneous parallel computing systems.

Parallelizing and accelerating large-scale SpTV on distributed CPU-GPU architectures faces three main challenges. First, how to effectively map SpTV operations to target distributed heterogeneous multi-core architectures and leverage the multi-level and hybrid parallelism. Second, how to compress the storage of sparse tensors effectively. The data structure determines the memory footprints and data access patterns in SpTV. Third, how to fully leverage the computing power of high-performance heterogeneous parallel computing

systems for SpTV, based on the given task division method and sparse tensor compression storage format.

To alleviate the above-mentioned challenges, this paper exploits the hierarchical parallelism for SpTV on CPU-GPU heterogeneous parallel computing systems. First, we propose a multidimensional partitioning method to exploit the inter- and intra-node parallelism for SpTV. Second, based on the multidimensional data partitioning, we design a fiber-wise compressed storage format for sparse tensors to optimize the data access patterns for parallel SpTV. Third, we design the parallel streaming SpTV algorithm to overlap the data swapping cost and the computation overhead on GPUs.

This work has been supported in part by the Key-Area R&D Program of Guangdong Province (2021B0101190004), the Programs of the National Natural Science Foundation of China (No.62172157, 62202149), the Programs of the Hunan Province (No.2023GK2002, 2021RC3062, 2023JJ60002), the Programs of the Guangdong Province and Shenzhen (No.2023A1515012915, JCYJ20210324135409026), and the Program of Zhejiang Lab (No.2022RC0AB03).