

# MapReduce 和 Non-MapReduce 分布式计算框架综述

Jianhua Guo

## 摘要

随着信息技术的不断发展，越来越多的数据分析问题使用传统的算法难以解决，迫切需要新的方法来填补这方面的空缺。MapReduce 是一种分布式并行计算范式，它凭借独特的任务调度和 shuffle 机制，极大地提升了大数据分析计算的性能。尽管如此，基于 MapReduce 的计算模型在支持大数据分析计算的过程中仍然存在一些限制条件，如执行分析算法的效率较低，缺少数据的可扩展性等，因此，Non-MapReduce 分布式计算框架被提出并用以解决已知的这些问题。文章针对 MapReduce 和 Non-MapReduce 分布式计算框架进行概述，并归纳总结了两种计算范式的步骤、特点及优缺点。

**关键词：**MapReduce；Non-MapReduce；分布式计算框架；大数据分析

## 1. 引言

在大数据时代，数据规模和复杂性呈现爆炸式增长，推动着数据处理技术的迅猛发展。大数据不仅指数据量的庞大，还涵盖了数据类型的多样化和数据生成速度的提升。大数据通常被定义为具有“4V”特征的海量数据集合，即数据量庞大（Volume）、数据类型多样（Variety）、数据价值巨大（Value）和数据生成与处理速度快（Velocity）。随着物联网、社交网络、移动互联网等技术的发展，数据的生成和采集渠道变得更加多样，涵盖了结构化数据、半结构化数据和非结构化数据等多种类型，涉及科学、医疗、金融、商业等各个领域。这些海量而复杂的数据蕴含着巨大的潜在价值，但要充分挖掘和利用这些数据，就需要先进的大数据处理技术和方法。

大数据处理指的是在特定的硬件和软件资源范围内，利用高效算法对海量数据进行搜索、清洗、聚集、连接、排序、挖掘等一系列操作，从中提取出对决策有用的信息和知识。由于数据量巨大且无法全部加载至内存，大数据处理需要分布式系统和外存算法的支持。在大数据处理的实践中，常用的算法通常遵循分治策略，将复杂的计算任务分解为多个小任务并行处理，最终将计算结果汇总输出。这种处理方式显著提升了数据处理的速度和效率，使得科学研究、商业分析和社会管理等领域能够从大数据中快速获取洞察和决策支持。

在大数据处理领域，分布式计算框架起到了关键作用。分布式计算框架是一种软件架构，用于将复杂计算任务分解成小任务，并分发到多个计算节点上并行处理，以提高计算效率和处理能力。它通过任务调度、数据分片与存储、节点间通信和容错机制来协调各个节点的工作，确保系统在节点故障的情况下仍能继续运行，同时具备良好的扩展性，常见的分布式计算框架包括 Hadoop、Apache Spark 和 Flink，广泛应用于大数据处理 and 数据分析等领域。MapReduce 作为一种经典的分布式计算框架，基于分治策略将任务分解为 Map 和 Reduce 两个阶段。Map 阶段将任务分散至多个节点，进行并行数据处理；Reduce 阶段则对 Map 阶段的中间结果进行汇总与处理，得到最终结果。MapReduce 框架因其简单、高效、高容错和易扩展的特点，被广泛应用于科研、商业、社交网络等大数据场景，尤其适合在廉价服务器集群上进行大规模数据的批处理。然而，MapReduce 在一些特定场景下也存在局限性，特别是对于需要快速响应和迭代计算的任务，其计算效率和资源利用率相对较低。

随着应用需求的多样化,新的 Non-MapReduce 分布式计算框架不断涌现。非 MapReduce 分布式计算框架是一种新型的数据处理架构,旨在有效支持大数据分析,特别是近似计算。与传统的 MapReduce 框架不同,该框架通过 RSP (随机样本分块) 数据模型来处理数据,允许对数据集进行更灵活的分析。在这个框架中,计算操作分为本地操作和全局操作,先在本地节点上执行所有迭代算法并生成本地结果,再进行全局操作,从而显著降低数据通信成本。由于本地算法可以独立运行,无需重新改写为 Map 和 Reduce 操作序列,任何串行算法都能直接在本地节点的 RSP 数据块上应用。这种方法避免了传统方法中对大数据集进行切割和分配的复杂性。

非 MapReduce 框架通过块级抽样方法随机选择 RSP 数据块并进行分析,从而提高了大数据分析的准确性和效率。由于只需读取少量随机样本,抽样时间显著缩短,允许在 TB 级和更大的超大规模数据上进行处理。此外,框架还支持多个随机样本的近似计算,能将多个本地近似结果集成,从而提供更为可靠的统计估计和置信区间。这使得非 MapReduce 框架在进行大数据统计时能够达到更好的效果,而不仅仅是精确计算,且在节点故障的情况下,依然能够确保结果的可靠性与容错性。

非 MapReduce 分布式计算框架通过引入 RSP 数据模型和灵活的计算机制,为大数据分析提供了更高效、更准确的解决方案,能够应对现代数据处理中的挑战。

本文将详细综述 MapReduce 和 Non-MapReduce 分布式计算框架的工作机制、应用步骤及其特点和优劣势,分析不同框架在大数据处理中的适用性和性能表现,为在实际应用中选择合适的分布式计算框架提供依据。

本文的其余部分组织如下。第 2 节概述了 MapReduce 和 Non-MapReduce 分布式计算框架的相关工作。第 3 节在讨论了大数据计算处理的基本流程。第 4 节和第 5 节分别细致阐述了两种分布式计算框架的步骤、特点和优缺点。最后,在第 6 节对全文进行了高度总结。

## 2. 相关工作

本文主要关注 MapReduce 和 Non-MapReduce 分布式计算框架的相关文献。MapReduce 是一种自上而下的编程模型,广泛应用于大规模数据处理,其基本思想是将任务分为 Map 和 Reduce 两个阶段。相关工作主要集中在对 MapReduce 框架的优化、任务调度和资源管理等方面。

在早期的研究中,Dean 和 Ghemawat (2004) 提出了 MapReduce 的基本框架,奠定了其在分布式计算中的基础。随后,许多学者对 MapReduce 进行了一系列改进。例如,Apache Hadoop 作为实现 MapReduce 模型的开源框架,提出了针对大数据集的高效存储和计算解决方案。针对数据倾斜问题,Wang 等人 (2015) 提出了改进的负载均衡算法,旨在优化任务调度并提高资源利用率。

此外,一些研究者探讨了 MapReduce 在特定领域的应用,如 Zaharia 等人 (2010) 在机器学习中的应用,以及 Dyer 和 Shirky (2015) 在社交网络分析中的使用。这些研究表明,MapReduce 模型在处理海量数据方面的有效性,特别是在批处理和离线分析任务中。

尽管 MapReduce 模型在大数据处理领域取得了巨大成功,但它在执行迭代算法时却仍然面临着效率低下的问题。这是因为 MapReduce 在算法的每次迭代中都需要进行大量的数据传输和节点间通信,导致了显著的性能开销。为了克服这些限制,研究者们提出了非 MapReduce 计算框架,该框架通过预先将大数据文件分割成一系列随机样本文件,即 RSP 数据块,存储在集群节点上,从而减少了数据通信的需求。这种方法不仅提高了计算效率,还增强了对大数据集的可扩展性,特别是在处理需要迭代处理的复杂算法时。

非 MapReduce 框架的另一个关键优势是其对串行算法的直接支持。这意味着在分布式环境中,无需将串行算法改写为 MapReduce 风格,从而简化了开发过程,并允许更多的分析算法在分布式系统中使用。此外,非 MapReduce 框架通过 RSP 数据模型支持近似计算,使得在用户控制的精度范围内,

仅使用数据集的一个子集来计算近似结果成为可能。这种统计感知的方法不仅提高了计算效率，还减少了对存储资源的需求，因为不需要将整个数据集加载到内存中。

在智能大数据分析领域，非 MapReduce 计算框架显示出了其独特的优势。通过使用 Spark 实现的机器学习算法，研究者们比较了 MapReduce 和非 MapReduce 在处理大规模数据集时的性能。实验结果表明，非 MapReduce 在计算效率和数据可扩展性方面明显优于 MapReduce，能够高效地处理高达 10TB 的数据集。这些发现为大数据分析提供了新的视角，并指出了非 MapReduce 框架在未来可能成为主流的分布式计算技术。

在总结这些框架的相关研究时，可以将其大致分类为两大类：基于传统的 MapReduce 模型和基于内存计算的 Non-MapReduce 模型。基于 MapReduce 的方法多侧重于任务的分解和并行处理，适用于大规模批量数据处理。而基于 Non-MapReduce 的框架则通过内存计算和流处理的结合，适应了日益增长的实时数据处理需求。

这些研究不仅提高了分布式计算框架的性能，也推动了大数据处理技术的发展，促进了在各行业中更高效的数据分析和挖掘。随着技术的不断进步，未来的工作可能会集中在如何融合这两种框架的优势，以满足更为复杂的应用需求和处理挑战。

### 3. 大数据处理的基本流程

数据处理是指对收集到的原始数据进行一系列操作，以转化为结构化、可理解的信息的过程。其目的是提升数据的质量和可用性，通过数据清洗、整理、分析、转换和呈现等步骤，使数据更适合应用需求，并为决策提供支持。数据处理广泛应用于各类领域，是数据分析和数据挖掘等工作的基础。

最初，数据来源于多个渠道，包括数据库、网络日志和 Web 文档等。这些数据通常是原始的、未经处理的，可能包含噪声和不一致性，因此需要经过数据预处理阶段。在数据预处理阶段，主要进行三个关键操作：选择、清洗和转换。选择是指从大量数据中筛选出对分析有用的部分；清洗是指去除错误、

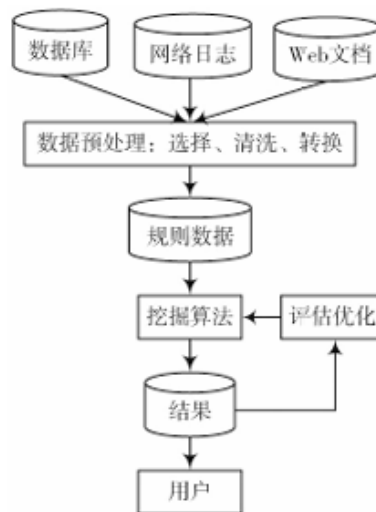


图 1. 数据处理的基本流程

重复或不完整的数据记录；转换则是将数据转换成适合分析的格式或结构。

经过预处理的数据变为规则数据，这些数据更加干净、有序，为后续的分析 and 挖掘做好了准备。接下来，进入评估优化阶段，这里涉及到对数据挖掘算法的选择和调整。根据数据的特点和分析目标，选择合适的挖掘算法，并对其进行参数调优，以确保算法能够在数据上达到最佳的挖掘效果。

最后，经过评估和优化的算法被应用于数据上，挖掘出有价值的信息和知识。这些结果可以用于支持决策、发现模式、预测趋势等，最终以报告、图表或其他形式呈现给用户，供用户进一步分析和决策使用。整个流程是一个连续的、迭代的过程，可能需要根据结果反馈回到前面的步骤进行调整和优化，以提高数据处理的质量和效果。

### 4. MapReduce 分布式计算框架

MapReduce 作为比较适用于进行大数据量处理、计算环节简单的并行计算框架，把 MapReduce 应用到数据挖掘方面成为有效解决大数据挖掘难题的一种需求 [6]。在这一节内容中，我们主要介绍了 MapReduce 分布式计算框架支持复杂算法大数据计算的步骤、特点和优缺点，包括它是如何运行的，以及在实际应用中具有哪些限制，并据此提出

Non-MapReduce 分布式计算框架存在的必要性。

#### 4.1. MapReduce 的步骤

首先需要了解 Hadoop，它是一个开源的 MapReduce 实现，同时提供了一个分布式文件系统 HDFS (Hadoop Distributed File System) [5]。HDFS 可以将超大文件自动分片并进行负载均衡，将这些分片存储在集群的各个节点上，并在存储数据的节点上执行 MapReduce 任务。这样一来，数据和计算任务在本地节点上进行，减少了网络传输带来的 I/O 开销，从而提升了运行效率。Hadoop 集群通常由一个主节点和多个从节点组成，可以运行在普通的商用计算机上，具有高性价比 [7]。

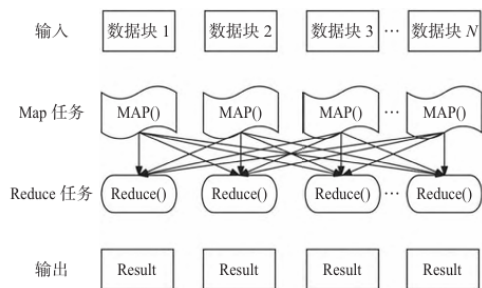


图 2. MapReduce 的步骤流程和原理

而基于 Hadoop 云平台的 MapReduce 是一种并行计算模型，最早由 Google 于 2004 年提出，专门用于在计算机集群上处理大规模数据集的分布式计算任务。它主要将数据处理任务分为 Map 和 Reduce 两个阶段，此外还包括数据输入和分割以及 Shuffle 阶段。

首先，数据输入是指将数据从不同来源导入到数据处理系统的过程，以便进行分析和处理。数据输入是数据处理的起点，确保数据能够被读取、解析并加载到内存或存储系统中，为后续的数据清洗、转换和分析提供基础，而数据分割则是指将完整数据集划分为较小的部分或子集的过程，以便能够更高效地进行处理和分析。数据分割通常用于并行计算或模型训练中，以便分散数据负载、降低存储需求或优化计算资源，例如在分布式计算中，数据集会被分割为若干数据块，由不同节点并行处理；在机器学习中，数据集也常被分割为训练集、验证集

和测试集，以便进行模型训练和评估。

第一阶段结束后是 Map 阶段，由 Map 函数在多个节点上并行运行，将输入数据划分为若干分片并进行处理，生成中间结果，Map 阶段算法可以开启一个或多个 Map 任务，常见的 Map 阶段算法包括数据搜索、数据清洗和数据变换；接着，在 Shuffle 阶段的过程中，Map 阶段生成的键值对会按照键进行分组和排序，以确保具有相同键的所有值能够传送到同一个 Reducer；最后，Reduce 函数会接收并处理这些中间结果，最终合并成全局输出 [10]。

MapReduce 的输入和输出数据通常以键值对的格式呈现，这种结构化的数据便于分布式存储和处理。MapReduce 的优势在于它的自动并行化、容错机制、以及负载均衡等特性，使得它在大数据处理和数据挖掘领域得到广泛应用。许多复杂的数据分析和数据挖掘算法已经被转化为基于 MapReduce 的模型，以便更高效地处理大数据。

#### 4.2. MapReduce 的特点

MapReduce 是专为分布式环境设计的大数据处理框架，其核心特点是通过将大规模数据任务分解为 Map 和 Reduce 两个阶段实现并行处理，从而显著提升数据处理速度和效率。该框架可在数百甚至数千台计算机上协同运行，各计算节点利用自身计算能力共同完成海量数据的处理工作。MapReduce 具备自动错误检测和恢复机制，保障任务的容错性，即便某些节点出现故障，任务也能在其他节点上重新执行。此外，开发者只需专注于编写 Map 和 Reduce 函数，任务的分配、执行及结果收集均由框架自动管理，这大大简化了并行编程的复杂度 [8]。MapReduce 尤其适用于 TB 或 PB 级数据的处理，采用分而治之的方式将任务分解成小块并行处理，从而提升整体效率。

该框架还具备较高的灵活性，适用于数据挖掘、日志分析和数据聚合等多种场景。此外，MapReduce 强调数据本地化处理，尽量让计算在数据所在位置执行，以减少网络传输成本并提升性能。MapReduce 提供可扩展的接口，允许开发者根据特定需求自定义 Map 和 Reduce 函数，进一步提升框架的适用性。



尽管 MapReduce 非常适合批处理任务，但在需要实时响应和执行复杂查询的场景下存在一定限制 [9]。Apache Hadoop 提供了 MapReduce 的开源实现，广泛应用于企业和科研领域，同时与 Hadoop 生态系统中的其他组件（如 HDFS、HBase、Hive）紧密集成，形成了功能全面的大数据处理解决方案。

### 4.3. MapReduce 的优缺点

MapReduce 是处理大规模数据集的分布式编程模型，具有一定的可扩展性和容错性，允许开发者通过简单的 Map 和 Reduce 函数来实现数据并行处理，从而显著降低了编程复杂度。在批量数据处理时，MapReduce 提供了较高的吞吐量，尤其是在其开源实现 Apache Hadoop 的支持下，易于获取、扩展和定制。此外，MapReduce 通过优化数据本地性，尽量让计算发生在数据所在位置，减少了网络传输成本并提高了整体处理效率。

然而，MapReduce 在一些方面也存在不足。首先，MapReduce 的数据 Shuffle 阶段可能导致资源利用率下降，尤其在网络和 I/O 成为瓶颈时表现尤为明显。此外，对于需要快速响应的实时任务，MapReduce 的批处理模式会导致一定的延迟，难以满足低延时要求。对于需要多次迭代的复杂算法，例如机器学习中的迭代计算，MapReduce 并不理想，因为每次迭代都需重新执行完整的 Map 和 Reduce 阶段，导致效率降低。MapReduce 也可能出现数据倾斜问题，导致部分节点负载过重而影响整体平衡性，增加了系统负担。

在数据处理流程的灵活性上，MapReduce 主要适合简单的键值对数据处理，对于复杂的数据处理需求和多步骤操作的支持不够灵活。同时，对于不熟悉 MapReduce 模型的开发者来说，编写和维护高效的 MapReduce 程序也具有一定难度。分布式系统的调试和维护相较于单机系统更为复杂。因此，尽管 MapReduce 在批量数据处理方面表现出色，但在实时处理、迭代计算和复杂数据处理场景中存在局限性，这也促使了像 Apache Spark 等新一代框架的出现，旨在弥补 MapReduce 的不足，满足更广泛的大数据处理需求。

## 5. Non-MapReduce 分布式计算框架

### 5.1. Non-MapReduce 的步骤

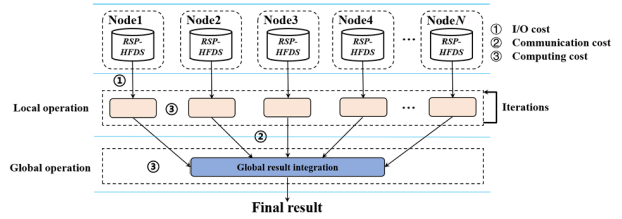


图 3. 非 MapReduce 分布式计算框架

首先需要了解 RSP 数据模型，RSP 数据模型是一种新的大数据管理和分析模型，它将大数据文件表示为一系列 RSP 数据块文件的集合，然后这些数据块分布存储在集群节点上。RSP 模型的生成操作确保每个 RSP 数据块的分布与整个大数据的分布在统计意义上保持一致，因此每个 RSP 数据块可以视为大数据的随机样本。这些数据块可以用来估计大数据的统计特征，或用于建立分类和回归模型。基于 RSP 模型的分析任务可以通过对这些 RSP 数据块的分析来完成，而无需对整个大数据进行计算，这大大减少了计算量，降低了对计算资源的需求，同时提高了集群系统的计算能力和扩展能力。

RSP 模型的核心思想在于将大数据划分为许多小的随机样本数据块文件，这样的划分方式使得随机样本数据可以直接通过选择数据块文件获得，避免了对大数据的单个记录进行抽样的复杂操作。此外，通过对少量数据块文件的分析和建模，可以得到大数据的统计估计结果和模型，从而极大地提高了大数据分析的效率和能力。

Non-MapReduce 则是一种由 RSP 数据模型支持的分布式计算框架，如图 3 所示 [3]，非 MapReduce 分布式计算框架通过利用 RSP 数据模型，在多个节点上分布式地处理大数据集，该计算框架整体分为两个步骤，一个是本地操作，一个是全局操作，在这个框架中，每个节点如 Node1、Node2、Node3 等，都存储着来自 HDFS 的 RSP 数据块，这些数据块是大数据集的随机样本。当执行迭代算法时，每个节点独立地从本地存储中读取其 RSP 数据块，减

少了网络 I/O 开销。由于数据处理主要在本地进行，节点间不需要传输中间结果，这大幅降低了通信成本。每个节点计算其 RSP 数据块的局部结果，这些计算并行发生，提高了整体的计算效率。最终，所有节点的局部结果被传输到一个主节点，在那里进行全局结果集成，得到最终的分析结果。这个过程减少了传统 MapReduce 框架中常见的数据混洗和跨节点通信的需求，使得非 MapReduce 框架在处理大数据时更为高效和可扩展，特别适用于需要迭代计算的场景。

## 5.2. Non-MapReduce 的特点

非 MapReduce 分布式计算框架以其处理大数据的高效性和可扩展性引起了广泛关注。与传统的 MapReduce 框架不同，这种框架通过减少数据通信和内存成本，为分布式大数据处理提供了一种创新的解决方案。通过避免对所有分布式数据块的直接处理，而是通过抽样的方式选择性地分析数据，这一框架在保持计算精度的同时有效降低了资源消耗。

该框架的一个核心特点在于依赖于 RSP 数据模型 [4]，这是一种新颖的数据表示方法。通过对分布式数据块进行随机抽样，RSP 数据模型实现了统计分析中的随机抽样效果，使得在仅分析一部分数据的情况下仍能近似描述整个数据集的特征。相比传统的 MapReduce 框架，该方法使得非 MapReduce 框架能更灵活地管理数据，避免了对所有数据块进行全面扫描的高资源需求，从而提高了计算效率。

非 MapReduce 框架还引入了“近似计算”概念，这是其高效能计算的关键之一。通过计算近似结果而非精确结果，该框架能够在低计算成本的条件下满足许多任务需求。这种近似计算方法在数据挖掘和机器学习等大数据分析领域尤为适用，能够在满足精度要求的前提下大幅提升计算速度。对于许多统计分析和模式识别任务而言，近似结果已经足够准确，因此这种方法减少了不必要的计算，使得大数据处理更加高效。

此外，非 MapReduce 框架支持在分布式环境中直接执行串行算法，这为用户提供了显著的便捷性。在 MapReduce 模式中，通常需要将串行算法

并行化并重新设计为 Map 和 Reduce 阶段，而非 MapReduce 框架则允许直接在分布式数据块上运行串行算法，降低了算法并行化的复杂性。通过这一特性，用户在分布式计算中应用传统的分析方法变得更为轻松，同时也大大降低了技术门槛 [1]。

实验结果表明，非 MapReduce 框架在计算效率和数据扩展性方面优于 MapReduce，能够高效处理高达 10TB 的数据集。这一框架不仅在处理大规模分布式数据集方面表现优异，还能用于机器学习、社交网络分析、不平衡分类等多种应用场景。因此，非 MapReduce 分布式计算框架凭借其 RSP 数据模型和近似计算方法，成为大数据分析的一种新型高效、可扩展的选择，为构建低能耗、广泛应用的大数据处理提供了坚实的技术支持。

## 5.3. Non-MapReduce 的优缺点

非 MapReduce 分布式计算框架具有显著的优缺点，其独特的设计使其在某些场景中表现优异，但也面临一些局限性。

首先，非 MapReduce 框架通过减少数据通信和计算资源的占用，显著提升了处理效率。它的设计使得大规模数据处理任务能够更高效地分解到各个节点上，特别适合需要对数据进行局部操作的场景。这种框架通常支持更复杂的算法和多种大数据分析任务，不像 MapReduce 只能处理较为简单的键值对数据结构。此外，非 MapReduce 框架允许在分布式环境中直接执行串行算法，减少了将串行算法改写为并行化模式的工作量，从而降低了算法实现的复杂性和维护成本。对于迭代计算等场景，非 MapReduce 框架能够避免 MapReduce 中频繁的数据写入和读取操作，显著提升了计算效率和整体性能。

非 MapReduce 框架的另一个优势是其灵活性，特别是在近似计算和数据采样方面。该框架通过随机样本划分数据模型来支持大数据集的近似计算。这种模型允许在处理超大规模数据集时仅利用数据的随机样本来估计最终结果，减少了对整个数据集的依赖 [2]。通过近似计算，这种框架能够高效地完成许多任务，不需要对全部数据进行完整处理，大

幅降低了计算资源需求和执行时间，这对于海量数据分析具有重要意义。

然而，非 MapReduce 框架也存在一定的局限性。尽管它能够提供近似计算结果，但在一些应用中精确性非常关键，近似计算的结果可能不满足实际需求。对于需要精确处理和复杂数据流的应用，非 MapReduce 框架的近似计算方式可能带来一定误差，进而影响分析的可信度和准确性。此外，非 MapReduce 框架依赖于内存来存储和处理数据，这在处理极其庞大的数据集时可能导致内存消耗过高，特别是在有限硬件环境下，系统可能无法承载过大的数据负载。

总体而言，非 MapReduce 框架在批处理任务和复杂计算场景中表现优异，其高效的数据处理和灵活的算法支持使其成为大数据分析的有效工具。尽管如此，针对需要精确计算的大数据场景，该框架的近似计算方法可能不适用；此外，对于需要实时响应的任务，非 MapReduce 框架的性能仍不及一些新兴的数据处理框架。

## 6. 结论

在本文中，我们对 MapReduce 和 Non-MapReduce 分布式计算框架进行了全面的综述。MapReduce 作为一种经典的分布式计算框架，以其简单、高效、高容错和易扩展的特点，在大规模数据批处理领域发挥了重要作用。然而，随着大数据应用需求的多样化，MapReduce 在实时处理、迭代计算和复杂数据处理场景中的局限性逐渐显现，这促使了 Non-MapReduce 框架的发展。

Non-MapReduce 框架通过引入 RSP 数据模型和近似计算方法，为大数据分析提供了更高效、更准确的解决方案。这种框架通过减少数据通信和内存成本，显著提升了处理效率，特别是在需要迭代计算的场景中表现出色。此外，Non-MapReduce 框架支持在分布式环境中直接执行串行算法，简化了算法并行化的复杂性，降低了技术门槛。

尽管 Non-MapReduce 框架在处理大规模数据集和复杂计算任务方面具有明显优势，但也存在一定的局限性。例如，在需要精确计算的场景中，近似

计算的结果可能不满足实际需求。此外，对于需要实时响应的任务，Non-MapReduce 框架的性能可能不及一些新兴的数据处理框架。

总的来说，MapReduce 和 Non-MapReduce 框架各有优势和适用场景。随着技术的不断发展，未来的研究可能会集中在如何融合这两种框架的优势，以满足更为复杂的应用需求和处理挑战。同时，新的分布式计算框架和算法的不断涌现，将进一步推动大数据处理技术的发展，促进各行业更高效的数据分析和挖掘。[?]

## 参考文献

- [1] X. S. et al. Alternate framework for distributed computing tames big data's ever growing costs. *journal Big Data Mining and Analytics*, 2023. 6
- [2] X. Sun, Y. C. Dingming Wu, L. Z. Xiao, and J. Z. Huang. Mapreduce vs. non-mapreduce efficiency and scalability in big data computing. 2023. 6
- [3] X. Sun, D. W. Yulin He, and J. Z. Huang. Survey of distributed computing frameworks for supporting big data analysis, big data mining and analytics. *IEEE*, 6(2):154–169, 2023. 5
- [4] X. Sun, C. J. Zhao Lingxiang, C. Y. Dingming, and J. Z. Huang. Non-mapreduce computing for intelligent big data analysis. *ENGINEERING APPLICATIONS OF ARTIFICIAL INTELLIGENCE*, 2024. 6
- [5] . . 周晓峰. 基于 mapreduce 计算模型的并行关联规则挖掘算法研究综述. *计算机应用研究*, 35(1):13–23, 2018. 4
- [6] . . 周翼. Mapreduce 模型在大规模数据并行挖掘中的应用. *智能物联技术*, 56(2):38–42, 2024. 3
- [7] . . . 孟凡兴. 大数据挖掘中的 mapreduce 并行聚类优化算法研究. *现代电子技术*, 42(11):161–164, 2019. 4
- [8] . 嵇圣国. Mapreduce 与 spark 用于大数据分析之比较. *软件学报*, 29(6):1770–1791, 2018. 4
- [9] . . . 王慧娇. 基于 mapreduce 改进 k-nn 的大数据分类算法研究. *微电子学与计算机*, 35(10):36–40+45, 2018. 5
- [10] 门威. 基于 mapreduce 的大数据处理算法综述. *濮阳职业技术学院学报*, 30(5):85–88, 2017. 4