

# DATA INTENSIVE COMPUTING

## RESEARCH PAPER REVIEW

Name: Aayushi Pandey  
UBID: 50610515

Date: 03/29/2025

### Introduction:

The paper that I have chosen for this paper review is “Parallel data intensive computing in scientific and commercial applications”. The authors of this research paper are Mario Cannataro, Domenico Talia, and Pradip K. Srimani. This paper explores data-intensive applications that process very large datasets. It discusses the main issues in parallel data-intensive computing and also explains how parallel and distributed computational environments will produce major improvements.

### 1. Why is this paper related to Data-Intensive Computing (DIC)?

This paper is related to data intensive computing because it considers challenges in parallel and distributed computing environments. The traditional computing mechanism is outdated and inadequate because of the complexity and volume of data. This paper talks about the following aspects of data intensive computing -

- **Scalability and performance:** The authors talk about how parallel and distributed systems handle a large amount of data generated commercially and scientifically.
- **Technical and Architecture:** It covers aspects of the data world such as data warehousing, data mining and database systems. It also highlights the techniques that are used to optimize these for high-performance computing environments.
- **Emerging Technologies:** The paper also explains data grids and peer-to-peer systems, which are critical for distributed data systems and preprocessing in DIC.

### 2. Your understanding of this paper.

This paper is a well-written overview of how parallel data-intensive computing focuses on storing and accessing data properly. The authors argue that data-intensive applications help to analyze, query and visualize large-scale datasets. The exponential growth of data in fields like astronomy, genomics, and e-commerce calls for the use of parallel and distributed computing to achieve practical processing times and meaningful insights.

### 3. This paper has at least two main contributions and at least two limitations.

Of the many insightful researches in the paper, the most impactful in opinion are as listed below-

#### a. **Comprehensive Survey of Parallel Data- Intensive Techniques:**

The paper explains the fundamentals of parallel processing in depth and covers a wide range of topics like data mining, data processing and its connectivity to web applications. Hence, it gives a holistic view of the field. For example, in section 2.1, the paper talks about how relational databases can be parallelized for OLTP and OLAP.

**b. Architectural insights for Scalable Systems:**

Authors have also drawn a comparison between shared-nothing and shared-memory parallel architecture. Also, they have shed light on how suitable they are for various data-intensive tasks. They have also talked about data grids as a way to manage scientific databases. For example, the Large Hadron Collider (LHC) at CERN generated petabytes of data, which is stored in Datagrids.

Limitations of the paper:

**a. Lack of Empirical Evaluation:**

This paper provides a decent theoretical overview, but it fails to show case studies or standardised practical results. For example, real-world benchmarks of performance of parallel data mining or data grids could have strengthened the arguments regarding these techniques.

**b. Dated Perspectives on Emerging Technologies:**

Since the paper was published back in 2002, it does not account for recent advancements in its examples. Technologies such as cloud computing, big data frameworks like Hadoop, Spark are also not included in it. Machine Learning techniques for optimization and analytics are also missing.

**c. Limited Discussion on Fault Tolerance and Reliability:**

It is a very critical issue in data-intensive computing where node failure, network issues, and data corruption are very common faulty phenomena that occur. This paper doesn't talk much in deep about these issues and just compares parallel and distributed techniques for managing and analysing large datasets.