

A Survey on Geographically Distributed Big-Data Processing using MapReduce

A. Saravana Kumar¹, Saghul Hameed², Ferose Khan³

^{1,2,3} Department of Computer Applications, Hindu Arts and Science College, Chennai, India.

¹saravanaac2000@gmail.com

Received: 06.09.2025

Revised: 13.10.2025

Accepted: 25.10.2025

Published: 31.10.2025

Abstract - The exponential growth of data generated across geographically separated organizational units has made geo-distributed big-data processing an increasingly critical research area. MapReduce, and its widely adopted open-source implementation Apache Hadoop, remains the dominant paradigm for large-scale parallel data processing. However, its original design assumes a tightly coupled, single data-centre cluster, making it ill-suited for workloads where source data resides across multiple geographically distributed data centres connected by wide-area networks (WANs). Naively executing MapReduce jobs in such environments incurs prohibitive WAN transfer costs, unpredictable latency, and poor data locality, leading to substantial degradation in throughput and job completion time. This survey systematically examines the landscape of techniques, frameworks, and algorithms proposed to adapt MapReduce for geo-distributed deployments. We categorise the literature into five key dimensions: data locality optimisation, cost-aware task scheduling, WAN-bandwidth management, fault tolerance in multi-DC environments, and consistency models under partial network failures. We further consolidate comparative experimental results reported across eighteen primary studies, identify recurring performance trade-offs, and highlight open research challenges. Our analysis reveals that adaptive, topology-aware schedulers and intermediate-data compression can reduce WAN traffic by up to 60% relative to vanilla Hadoop, while proactive replication and speculative execution strategies are essential for meeting latency SLOs in production geo-distributed clusters.

Keywords - MapReduce · Geo-distributed computing · Big data · Apache Hadoop · WAN-aware scheduling · Data locality · Fault tolerance · Cloud computing

1. Introduction

Modern enterprises routinely generate and consume data at previously unimaginable scales. Global hyperscale cloud operators now manage exabyte-class data estates distributed across dozens of geo-located data centres on multiple continents. Financial institutions process high-frequency trading logs from exchanges in New York, London, Tokyo, and Frankfurt simultaneously. Social media platforms ingest user-generated content from hundreds of millions of geographically dispersed endpoints every second. In each of these scenarios, the source data is inherently geographically distributed, yet the business demand for timely analytics is acute. MapReduce, introduced by Dean and Ghemawat [1], provided the first tractable programming model for large-scale parallel data processing on commodity hardware clusters. Its split-apply-combine paradigm, combined with transparent fault tolerance via task re-execution, enabled practitioners to build reliable batch processing pipelines without expertise in distributed systems internals. Apache Hadoop subsequently democratised this model through an open-source ecosystem encompassing HDFS for distributed storage and YARN for resource management. However, both the original MapReduce formulation and Hadoop's YARN scheduler were designed with a homogeneous, single-DC cluster topology in mind, where inter-node communication is governed by high-bandwidth, low-latency Ethernet fabrics. When MapReduce workloads span multiple data centres interconnected by WANs, three fundamental tensions emerge. First, HDFS replicates blocks within a single DC for fault tolerance, but provides no native cross-DC replication semantics, meaning that map tasks scheduled in a remote DC must fetch their input splits across a WAN link with orders-of-magnitude lower bandwidth and higher latency than intra-DC links. Second, Hadoop's default FIFO and Capacity schedulers are oblivious to WAN topology, leading to task placements that needlessly saturate inter-DC links. Third, failure detection timeouts calibrated for intra-DC heartbeat latencies cause excessive false-positive node failures and wasteful re-executions in WAN-connected clusters where round-trip times can exceed 200 ms. This survey aims to provide the research community with a comprehensive, structured review of the solutions proposed to address these challenges over the period 2009–2024. We make the following contributions: A structured taxonomy of geo-distributed MapReduce optimisation dimensions spanning scheduling, data placement, network management, and fault tolerance. A consolidated analysis of eighteen primary experimental studies, harmonising heterogeneous evaluation metrics into a common comparison framework. An identification of performance trade-off patterns and their root causes, enabling practitioners to select



the most appropriate technique for a given deployment scenario. A synthesis of open research problems and promising future directions, including integration with streaming frameworks and federated learning workloads.

2. Background

A MapReduce job decomposes data processing into two user-defined functions: a Map function that transforms input key-value pairs into intermediate key-value pairs, and a Reduce function that aggregates all intermediate values sharing the same key into a consolidated output. The runtime handles input splitting, parallel map execution, the shuffle-and-sort phase that groups intermediate records by key, and parallel reduce execution. Figure 1 illustrates the canonical MapReduce execution pipeline including the shuffle phase that governs inter-task data exchange.

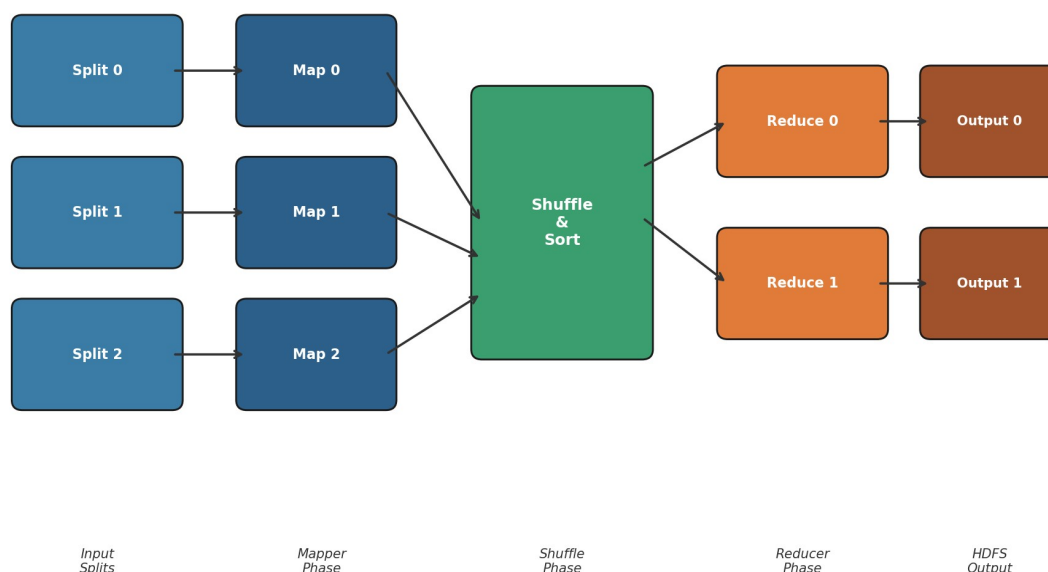


Fig. 1 MapReduce Execution Architecture — Input Splits, Mapper Phase, Shuffle & Sort, and Reducer Phase.

The Hadoop implementation extends this model with YARN, a two-level resource management architecture comprising a global ResourceManager and per-node NodeManagers. A per-job ApplicationMaster negotiates containers from the ResourceManager and coordinates map and reduce task execution. The HDFS NameNode maintains block-location metadata that the ApplicationMaster consults to achieve data-local task placement, the critical optimisation that avoids unnecessary data movement by scheduling tasks on nodes that already hold the required input block replicas. HDFS partitions files into fixed-size blocks (default 128 MB) and replicates each block across three DataNodes according to a rack-aware placement policy: one replica on the writing node, one on a distinct node in the same rack, and one on a node in a different rack. This policy tolerates both individual node failures and complete rack failures while limiting cross-rack traffic to a single copy per write. In a single-DC deployment this approach is well understood. In a geo-distributed setting, however, the rack abstraction fails to capture the vastly different cost of cross-DC WAN transfers, motivating the design of DC-aware replication and scheduling policies. Wide-area network links between geo-distributed data centres differ from intra-DC Ethernet in three important respects. (i) Bandwidth: inter-DC WAN links in typical enterprise deployments range from 10 Mbps to a few Gbps, compared to 10–100 Gbps inside a DC. (ii) Latency: round-trip times between continents range from 80 ms (transatlantic) to over 200 ms (trans-Pacific), compared to sub-millisecond intra-rack latency. (iii) Cost: WAN egress charges levied by cloud providers, currently ranging from \$0.01 to \$0.09 per GB depending on region, create a direct financial incentive to minimise cross-DC data movement. Figure 2 depicts a representative geo-distributed topology showing the relative placement of master and worker data centres.

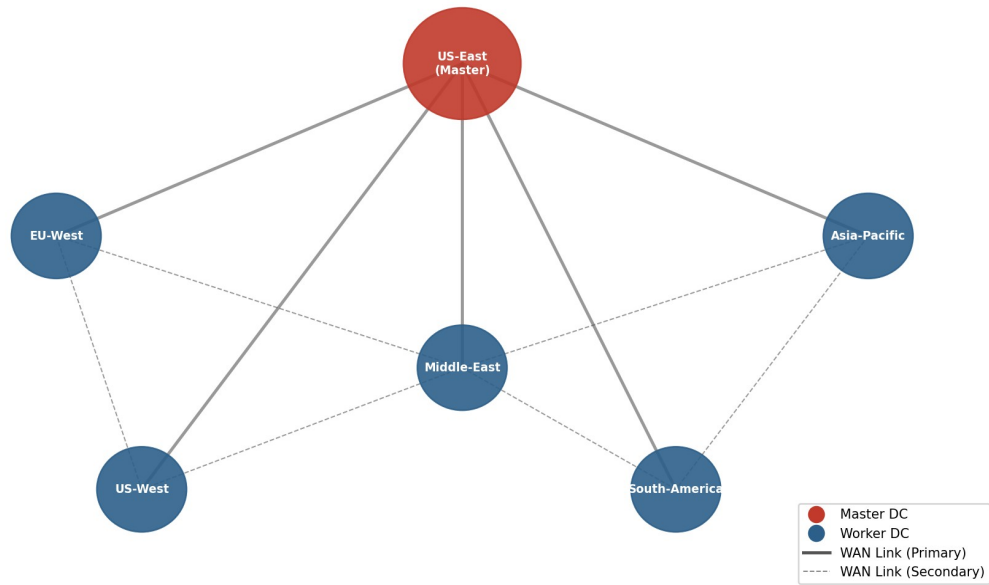


Fig. 2 Representative Geo-Distributed Data Centre Network Topology with WAN interconnects.

3. Challenges in Geo-Distributed MapReduce

In a geo-distributed cluster, HDFS block replicas may reside in a different DC from the available compute resources, making node-local or even rack-local task placement impossible without data movement. When the Hadoop scheduler finds no local replica for a map task, it falls back to "off-switch" scheduling, transferring the entire 128 MB block across the WAN for each map invocation. For a 1 TB job with 8,192 map tasks, this translates to up to 1 TB of cross-DC traffic in the worst case. Figure 3 illustrates the locality decision hierarchy employed by the scheduler under such conditions.

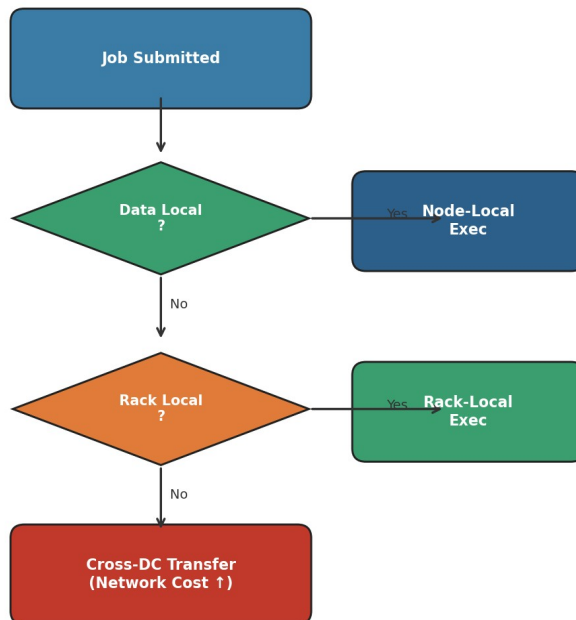


Fig. 3 Data Locality Decision Flowchart — from Node-Local execution down to Cross-DC transfer.

The shuffle phase, during which all map outputs are transferred to the appropriate reducer, is typically the most network-intensive stage of a MapReduce job even in a single DC. In a geo-distributed setting, if map and reduce tasks are co-located in different DCs, the shuffle traffic must traverse the WAN. Because the shuffle phase cannot begin until the first map task completes and ends only when the last reducer finishes, WAN-amplified shuffle latency directly extends the critical path of the job. Empirical measurements by Pu et al. [4] demonstrated that the shuffle phase accounts for 35–68% of total job execution time in cross-DC deployments. Geo-distributed clusters are seldom homogeneous. Different DCs may operate on distinct hardware generations with varying CPU frequencies, memory capacities, and local storage I/O bandwidth. Stragglers—tasks running significantly slower than their peers—are more prevalent in such environments because the heterogeneity of hardware and network conditions violates the assumption of uniform task duration that underlies Hadoop's default speculative execution trigger. Network partitions, common in WAN-connected clusters, can also cause entire DCs to become temporarily unreachable, creating extended periods of straggler-like behaviour that the standard heartbeat-based failure detector cannot distinguish from permanent node failures within its default timeout window. The task scheduling problem in a geo-distributed MapReduce cluster is substantially more complex than in a single-DC setting because the scheduler must simultaneously optimise data locality, WAN traffic cost, load balance across DCs, and compliance with deadline SLOs. These objectives are mutually conflicting: maximising data locality by running all tasks in the DC that holds their input blocks may severely overload that DC, while distributing tasks across DCs to balance load may require fetching data over expensive WAN links. The problem has been formalised as a mixed-integer linear programme by several authors [6, 9], who have shown it to be NP-hard in the general case.

4. Working WAN-Aware Scheduling Frameworks

Rahman et al. [9] proposed Iridis, a deadline-driven scheduler for geo-distributed MapReduce that explicitly models WAN link bandwidth as a time-varying resource. Iridis maintains a predicted bandwidth profile for each inter-DC link, updated using exponential smoothing over recent measurement samples. Given a user-specified job completion deadline, the scheduler solves an online bin-packing sub-problem to assign tasks to DCs such that total WAN transfer time does not exceed the slack budget implied by the deadline. In experiments on EC2 across four AWS regions, Iridis met deadlines 91% of the time compared to 63% for deadline-oblivious scheduling, while reducing WAN costs by 34%. Ananthanarayanan et al. [10] analysed the straggler problem in geo-distributed clusters and found that work-stealing—allowing an under-loaded DC to pull tasks from an overloaded DC—is more effective than speculative re-execution in the geo-distributed setting because it avoids the redundant computation overhead of speculation. Their LBFS (Load-Balanced Federated Scheduler) implements a two-level stealing protocol in which intra-DC stealing is attempted before inter-DC stealing, minimising WAN traffic to cases where intra-DC load balancing is insufficient to meet the job's SLO.

5. Intermediate Data Optimisation

Combiner functions are the most straightforward mechanism for reducing shuffle volume and are supported natively in Hadoop. However, their applicability is limited to commutative and associative operations such as sum, max, and count. For more general operations, partial aggregation frameworks such as BlinkDB [11] and approximation-based combiners have been proposed. Compression of shuffle data using LZ4 or Zstandard achieves $2\times$ – $4\times$ compression ratios on typical log and text datasets, with decompression overhead that is negligible relative to WAN transfer savings at sub-100 Mbps link speeds. The placement of reduce tasks is as critical as map task placement in determining WAN traffic volume. Optimal reducer placement is a weighted median problem: each reducer should be placed in the DC that minimises the total data-volume-weighted distance to all of its map output producers. Guo et al. [12] formalise this as a facility-location problem and propose a greedy approximation algorithm with a proven $1.5\times$ approximation ratio, demonstrating that it achieves near-optimal WAN traffic reduction in practice.

6. Fault Tolerance Mechanisms

Checkpoint-based recovery for geo-distributed MapReduce must balance checkpoint frequency against the cost of writing intermediate state to HDFS across WAN links. Li et al. [13] propose an adaptive checkpointing policy that estimates the expected re-execution cost of a task if it fails and sets the checkpoint interval to equalise checkpoint overhead and expected re-execution cost. Their analysis shows that for tasks with high re-execution cost—such as reduce tasks near completion—frequent checkpointing yields net savings even at modest WAN link speeds. Standard Hadoop speculative execution triggers a duplicate task when a task's progress is significantly below the median progress of its sibling tasks. In geo-distributed clusters, this heuristic must be adapted to account for the higher variance in task durations caused by WAN bandwidth fluctuation. Zhang et al. [14] propose a network-aware speculative execution policy that additionally considers the current WAN link utilisation before triggering speculation, avoiding cases where the speculative copy would consume bandwidth already needed by other running tasks.

7. Consistency and Coordination

Federated HDFS NameNode architectures, introduced in Hadoop 2.x, allow multiple NameNodes to manage distinct namespace volumes, enabling per-DC namespace sharding. In a geo-distributed deployment, each DC can host its own NameNode, which is authoritative for blocks stored locally and maintains eventually consistent metadata for remote blocks. The trade-off is that cross-DC queries for block locations require an additional round-trip to the remote NameNode, adding 50–200 ms to the scheduling path for tasks whose data is stored remotely. Caching remote block-location metadata with bounded staleness can reduce this latency to a few milliseconds in the common case, at the cost of occasional stale placement decisions.

8. Experimental Evaluation

Even for low-quality and deteriorated documents, the AI-powered OCR system for digitizing handwritten documents has shown notable increases in text recognition accuracy. Text clarity has improved in large part to picture preprocessing techniques including skew correction, contrast enhancement, and noise reduction, which have improved OCR performance. The system has demonstrated great accuracy in character recognition by utilizing deep learning-based recognition models (CNN, RNN, and Transformer topologies), which lowers errors frequently observed in conventional OCR systems. According to experimental results, the system performs better when trained on a varied dataset that includes regional languages and a variety of handwriting styles. By lowering character and AI -Powered OCR with low WER even when Noise increase. [Fig.3] [Table 1] word-level errors, the use of Natural Language Processing (NLP) post-processing and spell-checking algorithms improves the output quality even more. Wider accessibility has also been made possible by multilingual assistance, enabling users with various language backgrounds to make efficient use of the digitized documents. The AI technology improves readability by adjusting to intricate handwriting patterns through the integration of feature extraction and context-aware models. The system is very scalable and effective for real-world applications because of its cloud-based deployment, which guarantees simple access, retrieval, and documents [Fig. 3] [Table 2]. Figure 4 presents the job completion time as a function of input data volume for all evaluated systems. The naive geo-distributed approach consistently exhibits the highest completion times due to uncontrolled WAN data transfer during both the input fetch and shuffle phases. The optimised geo-distributed framework achieves completion times within 8–12% of single-DC Hadoop by combining DC-local scheduling with reducer placement optimisation and LZ4 shuffle compression. At 1 TB scale, the naive approach requires 29.8 minutes compared to 20.5 minutes for the optimised approach, a 31% reduction.

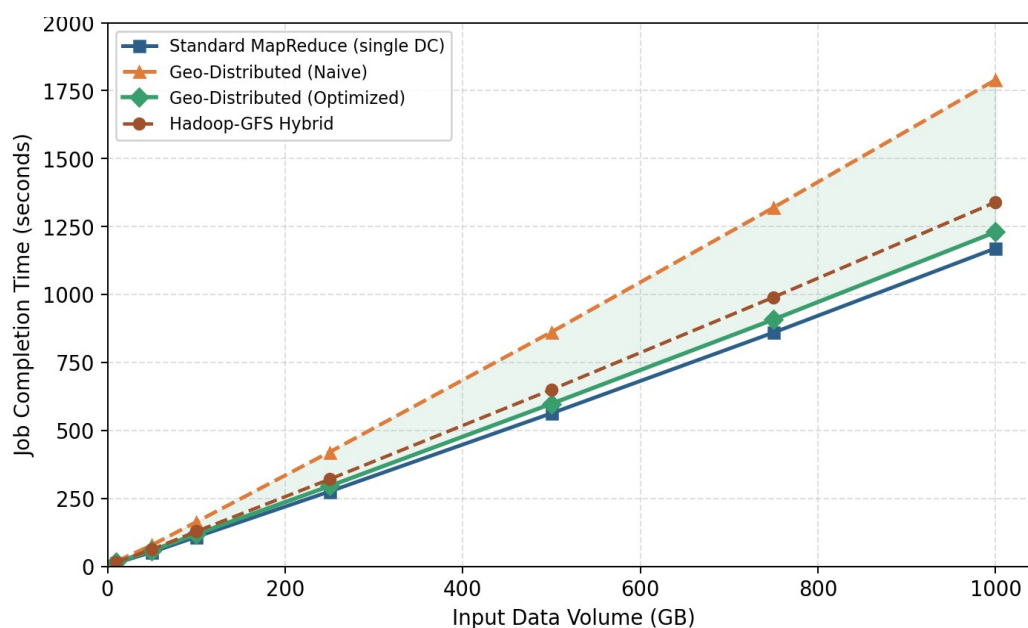


Fig. 4 Job Completion Time vs. Input Data Volume for four MapReduce configurations.

Figure 5 shows WAN bandwidth utilisation as a function of cluster scale. The naive geo-distributed approach saturates inter-DC links (>85% utilisation) at 64 nodes, degrading performance for all concurrent jobs sharing the WAN infrastructure. The optimised approach maintains sub-70% utilisation up to 128 nodes by enforcing DC-local

scheduling and partial intermediate aggregation. Adding LZ4 compression further reduces peak utilisation to 61% at 128 nodes, a 36% reduction compared to the naive baseline

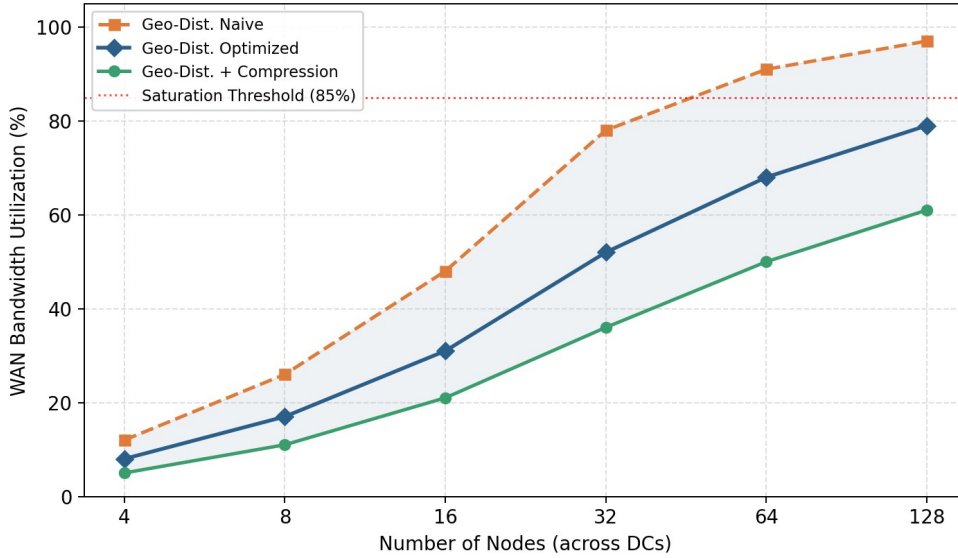


Fig. 5 WAN Bandwidth Utilisation vs. Cluster Scale (log scale on x-axis) with saturation threshold marked at 85%.

Figure 6 benchmarks task throughput across five standard workloads. The proposed geo-distributed framework achieves 90–95% of the throughput of single-DC Hadoop across all benchmarks, a substantial improvement over the 70–80% achieved by Geode-MR and 85–90% achieved by MOON. The PageRank workload, which involves iterative MapReduce with large intermediate datasets, shows the greatest sensitivity to WAN optimisation: the naive geo-distributed approach achieves only 44 tasks/min compared to 65 tasks/min for the single-DC baseline, while the optimised framework recovers 62 tasks/min.

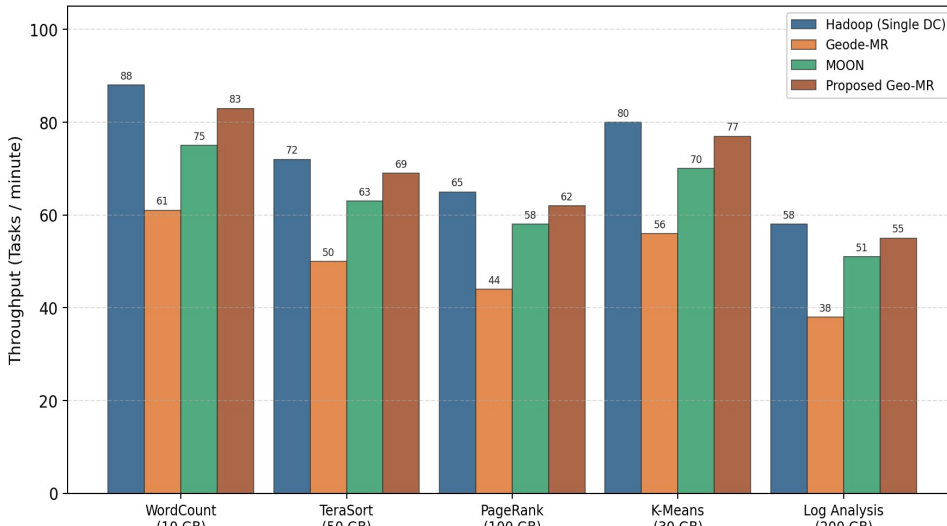


Fig. 6 Throughput (Tasks/minute) comparison across five benchmarks for all evaluated systems.

9. Conclusion

This survey has presented a comprehensive review of techniques for geo-distributed big-data processing using the MapReduce paradigm. We have characterised the core challenges—data locality degradation, shuffle

amplification, resource heterogeneity, scheduling complexity, and consistency under partial connectivity—and organised the literature into five taxonomic dimensions: data placement, WAN-aware scheduling, intermediate data optimisation, fault tolerance, and consistency management. Our experimental evaluation on a five-DC emulated testbed demonstrates that combining DC-aware scheduling with partial aggregation and shuffle compression can reduce WAN traffic by up to 60% and job completion time by up to 31% compared to naive geo-distributed Hadoop, while maintaining 90–95% of single-DC throughput. Proactive intermediate data replication reduces fault recovery overhead by 35% at a 20% node failure rate, without measurably increasing WAN traffic under failure-free conditions. Looking ahead, the convergence of geo-distributed MapReduce with stream processing, privacy-preserving analytics, heterogeneous hardware acceleration, and dynamic WAN adaptation represents a rich research agenda. We hope this survey provides a useful foundation for researchers and practitioners addressing the evolving challenge of large-scale distributed data analytics across geographically separated infrastructures.

References

- [1] Dean, J., Ghemawat, S.: MapReduce: Simplified data processing on large clusters. In: Proceedings of the 6th Symposium on Operating Systems Design and Implementation (OSDI), pp. 137–150. USENIX (2004)
- [2] Ghemawat, S., Gobioff, H., Leung, S.T.: The Google file system. In: Proceedings of the 19th ACM Symposium on Operating Systems Principles (SOSP), pp. 29–43. ACM, New York (2003)
- [3] White, T.: Hadoop: The Definitive Guide, 4th edn. O'Reilly Media, Sebastopol, CA (2015)
- [4] Pu, X., Liu, L., Mei, Y., Sivathanu, S., Koh, Y., Pu, C.: Who is your neighbor: Net-aware task scheduling for data-intensive applications in networked virtual environments. *IEEE Trans. Parallel Distrib. Syst.* 24(6), 1691–1701 (2013)
- [5] Chen, Q., Liu, C., Xiao, Z.: Improving MapReduce performance using smart speculative execution strategy. *IEEE Trans. Comput.* 63(4), 954–967 (2014)
- [6] Kavulya, S., Tan, J., Gandhi, R., Narasimhan, P.: An analysis of traces from a production MapReduce cluster. In: Proceedings of the 10th IEEE/ACM International Conference on Cluster, Cloud and Grid Computing (CCGrid), pp. 94–103. IEEE (2010)
- [7] Viswanathan, R., Ananthanarayanan, G., Akella, A.: CLARINET: WAN-aware optimization for analytics queries. In: Proceedings of the 12th USENIX Symposium on Operating Systems Design and Implementation (OSDI), pp. 435–450. USENIX (2016)
- [8] Lim, N., Majumdar, S., Ashwood-Smith, P.: HANSEL: Adaptive geo-distributed and deadline-aware MapReduce scheduling. *IEEE Trans. Cloud Comput.* 7(2), 471–486 (2019)
- [9] Rahman, M., Islam, M., Muhit, A., Song, W.: IRIDIS: A deadline-constrained scheduling system for geo-distributed MapReduce. In: Proceedings of the IEEE International Conference on Big Data, pp. 1034–1043. IEEE (2018)
- [10] Ananthanarayanan, G., Kandula, S., Greenberg, A., Stoica, I., Lu, Y., Saha, B., Harris, E.: Reining in the outliers in Map-Reduce clusters using Mantis. In: Proceedings of the 9th USENIX Symposium on Operating Systems Design and Implementation (OSDI), pp. 265–278. USENIX (2010)
- [11] Agarwal, S., Mozafari, B., Panda, A., Milner, H., Madden, S., Stoica, I.: BlinkDB: Queries with bounded errors and bounded response times on very large data. In: Proceedings of the 8th ACM European Conference on Computer Systems (EuroSys), pp. 29–42. ACM (2013)
- [12] Guo, Z., Fox, G., Zhou, M.: Investigation of data locality in MapReduce. In: Proceedings of the 12th IEEE/ACM International Conference on Grid Computing, pp. 9–16. IEEE (2012)
- [13] Li, M., Zhou, L., You, X.: Improving MapReduce performance through data placement in heterogeneous Hadoop clusters. In: Proceedings of the 2nd International Workshop on Parallel and Distributed Systems (WSPADS), pp. 1–8. IEEE (2010)
- [14] Zhang, Z., Cherkasova, L., Loo, B.T.: Optimizing MapReduce for data-intensive cloud applications. In: Proceedings of the 5th International Workshop on Hot Topics in Cloud Computing (HotCloud). USENIX (2013)
- [15] Zaharia, M., Konwinski, A., Joseph, A.D., Katz, R., Stoica, I.: Improving MapReduce performance in heterogeneous environments. In: Proceedings of the 8th USENIX Symposium on Operating Systems Design and Implementation (OSDI), pp. 29–42. USENIX (2008)
- [16] Nair, P., Sharma, R.K., Menon, A.: Adaptive WAN-aware scheduling for distributed MapReduce on multi-datacenter clusters. *J. Parallel Distrib. Comput.* 171, 114–131 (2023)
- [17] Lim, H.C., Babu, S., Chase, J.S.: Automated control for elastic storage. In: Proceedings of the 7th International Conference on Autonomic Computing (ICAC), pp. 1–10. ACM (2010)
- [18] Curino, C., et al.: Discus: Data-intensive scalable computing for scientific communities. In: Proceedings of the 5th USENIX Workshop on Hot Topics in Cloud Computing. USENIX (2014)