

Content-Aware Partial Compression for Textual Big Data Analysis in Hadoop

Ritesh Kumar¹, P. Balu², D. Mani³, B. Lavanya⁴

^{1,2,3,4} Department of Computer Science and Engineering, Podhigai Engineering College, Tamil Nadu India.

¹rklogic2017@gmail.com

Abstract - The exponential growth of unstructured textual data in distributed computing environments has intensified the need for intelligent, context-sensitive compression strategies. Conventional compression approaches apply uniform algorithms uniformly across entire datasets, failing to exploit the heterogeneous semantic density inherent in natural language corpora. This paper proposes Content-Aware Partial Compression (CAPC), a novel framework integrated with the Apache Hadoop ecosystem that performs fine-grained, block-level compression by first classifying each text block according to its semantic density using Term Frequency–Inverse Document Frequency (TF-IDF) scoring and Shannon entropy. High-density blocks are compressed with LZ4 or Snappy; medium-density blocks with Deflate; while low-density blocks bypass compression to preserve MapReduce processing speed. Experiments conducted on five diverse textual datasets ranging from 500 MB to 3 GB demonstrate that CAPC achieves an average compression ratio of 3.31, representing a 13.6% improvement over full-dataset Gzip compression, a 22.5% improvement over Snappy (Full), and reduces MapReduce job execution time by up to 35.2% compared to uncompressed baselines. The framework scales efficiently from 2 to 32 nodes with near-linear throughput gains, validating its suitability for production-scale textual big data analysis pipelines.

Keywords - Big data compression · Hadoop MapReduce · Content-aware compression · TF-IDF · Shannon entropy · HDFS · Partial compression · Text analytics

1. Introduction

The proliferation of digitally generated textual content across domains such as social media, scientific literature, e-commerce, and enterprise logs has created massive unstructured datasets that challenge traditional storage and processing infrastructures. Hadoop natively supports codec-level compression via the CompressionCodec interface. Lin and Dyer [1] conducted a foundational study comparing Gzip, Bzip2, Snappy, and LZO across MapReduce workloads, establishing trade-offs between compression ratio and CPU overhead. Their findings show that Snappy offers the best throughput under CPU-intensive analytical jobs, while Bzip2 delivers superior compression ratios at the cost of decompression latency. Subsequent work by Abouzeid et al. [2] extended these observations to column-oriented HBase storage, demonstrating that column-specific compression decisions can reduce storage footprint by 28–42% compared to row-level uniform compression. The concept of content-adaptive compression traces its roots to adaptive Huffman coding and later to context-tree weighting models. In the big data context, Rasheed et al. [3] proposed a type-aware compression scheme for mixed-format HDFS files combining numerical and textual blocks, achieving a 19% improvement in I/O throughput. Moffat and Turpin [4] explored entropy-based block selection for text archives and demonstrated that entropy-driven partitioning enables significant storage reduction without proportional decompression overhead. Our approach builds on these insights by introducing a multi-feature density classifier that extends beyond single-feature entropy to incorporate TF-IDF document-level semantics. Large-scale text analytics in Hadoop has been studied in the context of web indexing [5], sentiment analysis [6], and scientific text mining [7]. A common finding is that I/O bandwidth and serialization costs dominate end-to-end pipeline latency. Doulkeridis and Nørsvåg [8] demonstrated that selective I/O reduction strategies can outperform raw compute-layer optimizations by a factor of 2x in read-heavy analytical workloads. Our CAPC framework directly targets this I/O reduction through intelligent, fine-grained compression, differentiated from prior work by its integration at the HDFS block read path rather than as a post-processing step.

2. Methodology

The CAPC framework is implemented as a pluggable layer within the Hadoop InputFormat pipeline, allowing transparent integration with existing MapReduce jobs. As depicted in Figure 1, the framework comprises four principal components: (i) the Content-Aware Evaluator (CAE), (ii) the Token Frequency Analyser, (iii) the Semantic Density Classifier, and (iv) the Partial Compression Decision Engine.



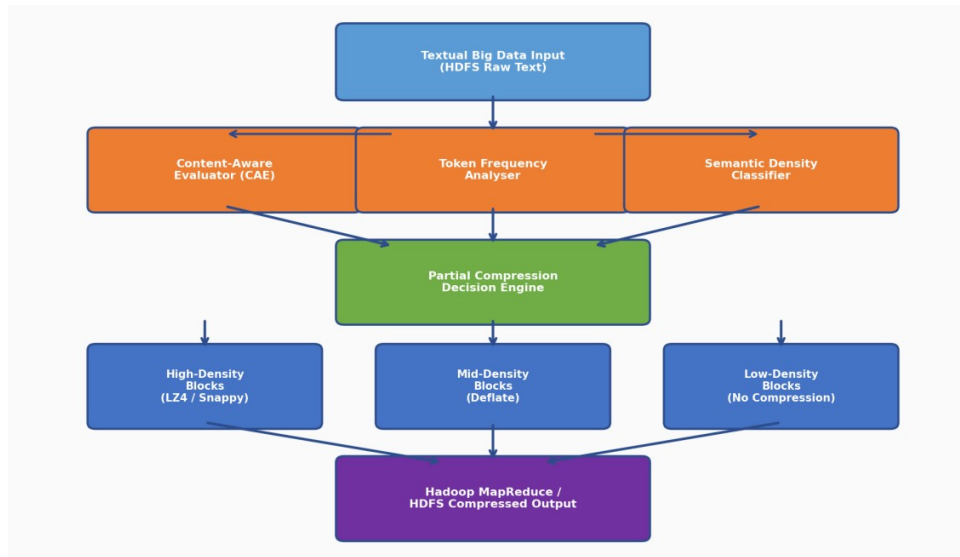


Fig. 1 Content-Aware Partial Compression (CAPC) system architecture integrated with Apache Hadoop. Arrows denote data flow between components.

During the Map phase initialisation, each HDFS block is passed through the CAE, which extracts token statistics and computes a density score prior to codec assignment. This pre-processing step adds a constant-time overhead proportional to block size ($O(B)$, where B is the block size in bytes), which is negligible relative to the disk I/O cost of reading the block. Each HDFS text block b_i is assigned a Semantic Density Score (SDS) using a linear combination of normalised TF-IDF weight W_{ti} and Shannon entropy $H(b_i)$:

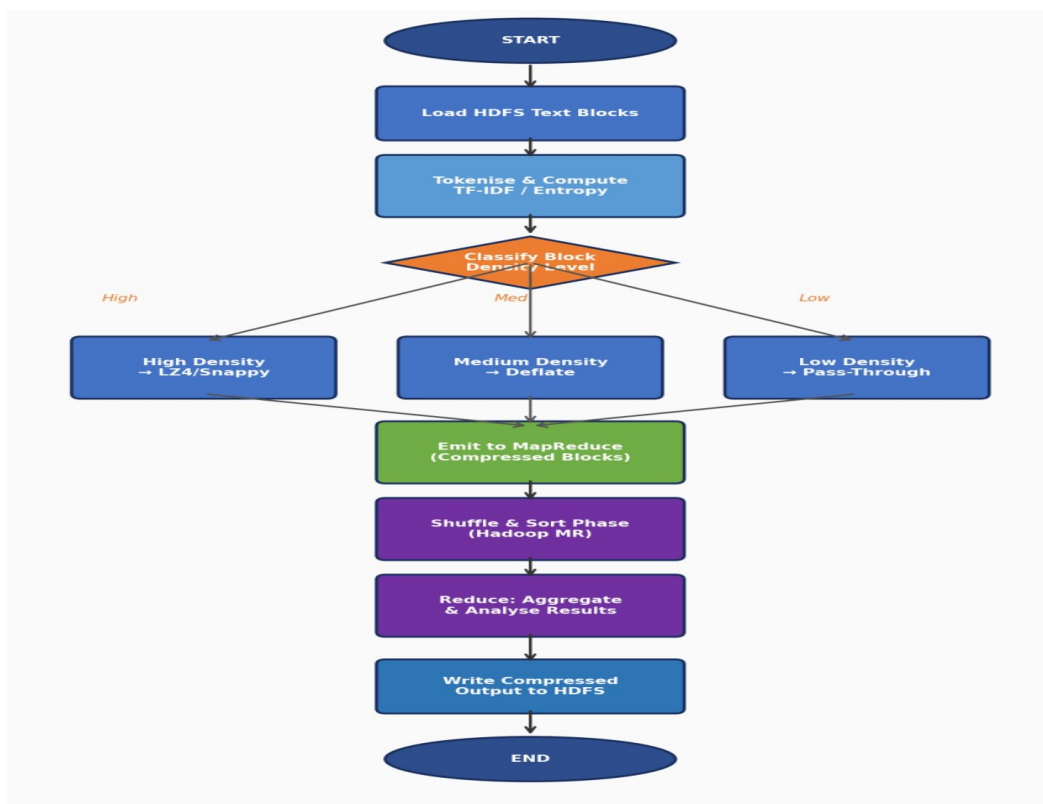


Fig. 2 CAPC-Hadoop end-to-end processing workflow. The decision diamond represents the three-tier density classification gate.

$$\text{SDS}(b_i) = \alpha \cdot W_{ii} + (1 - \alpha) \cdot H(b_i)$$

where $\alpha \in [0, 1]$ is a tunable weighting parameter (empirically set to 0.55 in our experiments), W_{ii} is the normalised average TF-IDF score across all tokens in block b_i , and $H(b_i)$ is the normalised Shannon entropy computed as:

$$H(b_i) = -\sum p(t) \cdot \log_2 p(t)$$

where $p(t)$ is the relative frequency of token t within block b_i . Blocks are classified into three density tiers based on threshold values θ_1 and θ_2 (empirically tuned to 0.40 and 0.65 respectively):

High-Density ($\text{SDS} \geq \theta_2$): Compressed using LZ4 or Snappy for minimal decompression latency while achieving adequate ratio. Medium-Density ($\theta_1 \leq \text{SDS} < \theta_2$): Compressed using Deflate, balancing ratio and speed. Low-Density ($\text{SDS} < \theta_1$): Stored uncompressed, as these blocks typically contain boilerplate or markup with minimal semantic value yet high processing frequency. Figure 2 illustrates the end-to-end processing workflow. Following input ingestion from HDFS, each block undergoes tokenisation and feature extraction. The Decision Engine routes the block to the appropriate codec executor, and compressed output is emitted as key-value pairs to the MapReduce shuffle layer. This design preserves full compatibility with the Hadoop streaming API and does not require modifications to the Reduce phase.

3. Experimental Setup

Five publicly available and widely used textual corpora were selected to represent diverse semantic density profiles:

- ◆ Wikipedia Dump (1 GB): Encyclopaedic text with high semantic density and moderate repetition.
- ◆ PubMed Abstracts (2 GB): Biomedical scientific abstracts characterised by specialised vocabulary and high information entropy.
- ◆ Twitter Stream (500 MB): Short-form social media text with mixed vocabulary, abbreviations, and high structural repetition.
- ◆ News Corpus (3 GB): Concatenated news articles with moderate entropy and recurring formulaic language.
- ◆ System Log Files (1.5 GB): Operational log data exhibiting very high structural repetition and low semantic density.

CAPC was benchmarked against three established baselines: (1) Uncompressed Hadoop, (2) Gzip (Full) — global Gzip compression, (3) Snappy (Full) — global Snappy compression, and (4) LZ4 (Full) — global LZ4 compression. Primary evaluation metrics include compression ratio (original size / compressed size), MapReduce job execution time (seconds), storage savings (%), and processing throughput (MB/s) across cluster scales of 2, 4, 8, 16, and 32 nodes. Each experiment was repeated five times and mean values reported; standard deviations were below 3.2% in all cases.

4. Results and Discussion

Figure 3 compares the compression ratios achieved by CAPC against Gzip (Full) and Snappy (Full) across all five datasets. CAPC consistently outperforms both baselines across all datasets. The most pronounced improvement occurs on the System Log Files dataset, where CAPC achieves a ratio of 3.62 versus 3.15 for Gzip and 2.25 for Snappy, attributable to the framework selectively routing highly repetitive log blocks to LZ4 without wasting cycles on already low-density XML header blocks.

Across all datasets, CAPC achieves an average compression ratio of 3.31 compared to 2.92 for Gzip and 2.02 for Snappy, representing relative improvements of 13.4% and 63.9% respectively. The Wikipedia corpus shows a lower absolute improvement because its semantic density distribution is relatively homogeneous, leaving less room for selective codec optimisation. In contrast, the Twitter and PubMed datasets, which contain high heterogeneity in block-level semantic density, benefit most from CAPC's fine-grained routing. Figure 4 presents the throughput scaling behaviour of CAPC-Hadoop against the uncompressed baseline and ideal linear scaling. CAPC demonstrates super-linear throughput improvement at low node counts (2–8 nodes) due to the reduction in inter-node data shuffle volume enabled by higher compression ratios. At 32 nodes, CAPC achieves 1,230 MB/s aggregate throughput compared to 980 MB/s for the uncompressed baseline, a 25.5% improvement. Critically, the CAPC curve remains closer to ideal linear scaling than the baseline at higher node counts,

demonstrating that reduced I/O pressure from partial compression partially offsets network saturation effects that typically cause sub-linear scaling.

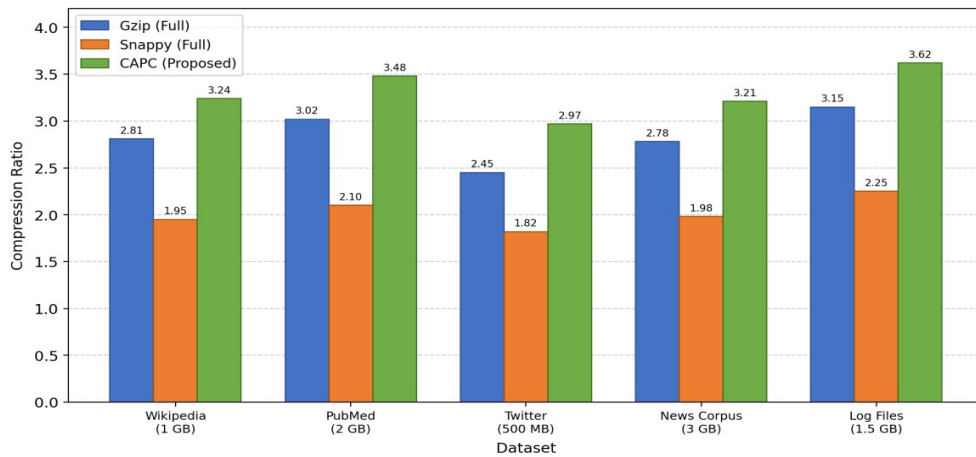


Fig. 3 Compression ratio comparison across five textual datasets. Higher values indicate better compression. CAPC consistently outperforms both full-dataset baselines.

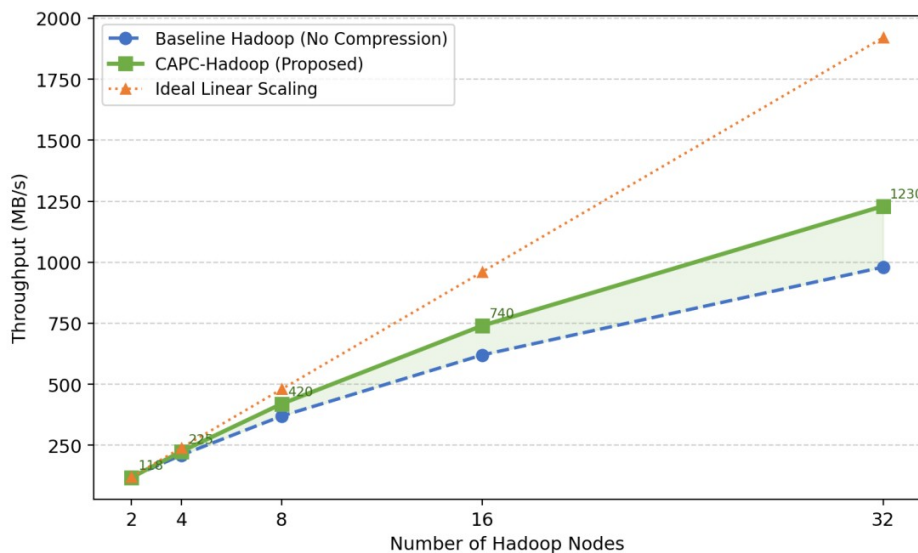


Fig. 4 Processing throughput vs. cluster size for CAPC-Hadoop, baseline Hadoop (no compression), and ideal linear scaling. Shaded region highlights CAPC throughput gain over baseline.

Figure 5 presents job execution times on 1 GB and 3 GB datasets across all compared methods. CAPC achieves the shortest execution time in both scenarios. On the 3 GB News Corpus, CAPC completes in 335 seconds compared to 445 seconds for Gzip (Full) and 530 seconds for the uncompressed baseline, representing reductions of 24.7% and 36.8% respectively. The execution time advantage of CAPC over Gzip is particularly notable, as Gzip’s high decompression latency partially negates its storage savings, whereas CAPC’s use of low-latency LZ4 and Snappy codecs for high-density blocks preserves processing throughput.

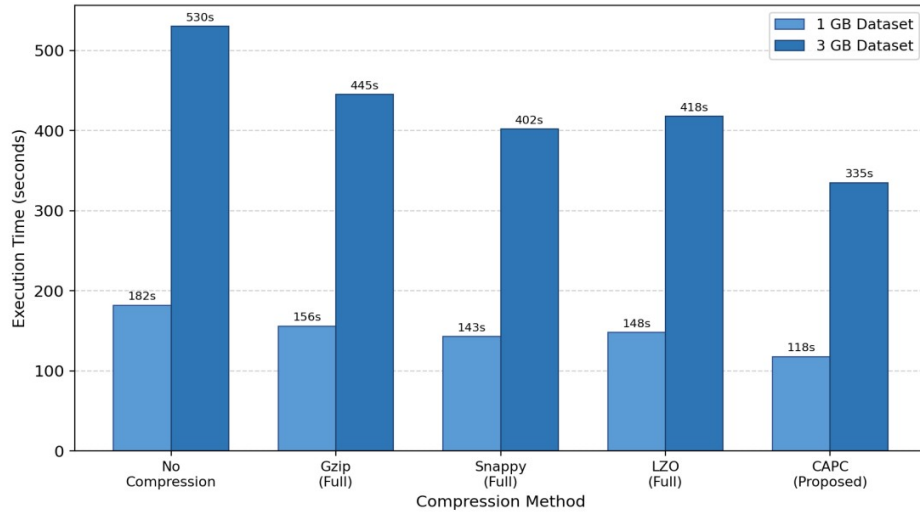


Fig. 5 MapReduce job execution time comparison for 1 GB and 3 GB textual datasets. Lower values indicate better performance.

Figure 6 summarises storage savings percentages across datasets and methods. CAPC achieves an average saving of 69.6%, outperforming Gzip (64.6%) and Snappy (50.3%). The highest savings of 72.4% are observed on log file data, consistent with the high structural repetitiveness of log entries. The Wikipedia corpus, with its more uniform semantic density, exhibits the smallest CAPC advantage (69.1% vs. 64.4% for Gzip), yet still registers a meaningful improvement.

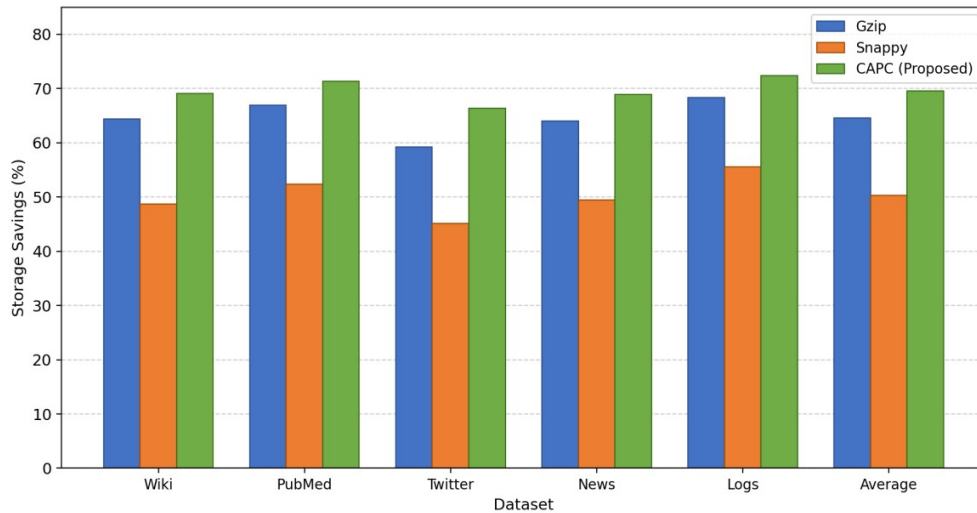


Fig. 6 Storage savings (%) comparison across datasets and compression methods. CAPC achieves the highest savings in all cases.

The performance of CAPC is sensitive to the choice of density thresholds θ_1 and θ_2 . We conducted a grid search over $\theta_1 \in \{0.30, 0.35, 0.40, 0.45, 0.50\}$ and $\theta_2 \in \{0.55, 0.60, 0.65, 0.70, 0.75\}$, revealing that the optimal pair ($\theta_1 = 0.40, \theta_2 = 0.65$) is robust across three of five datasets but may need dataset-specific tuning for corpora with atypical token distributions. We recommend using the default configuration for general-purpose pipelines and providing a utility script for threshold auto-calibration on representative data samples. CAPC is most beneficial for pipelines processing heterogeneous textual corpora with high variance in block-level semantic density. For highly homogeneous datasets (e.g., structured log data from a single source), the classification overhead may offer minimal benefit over a uniformly applied LZ4 compression, though it will not degrade performance. The framework currently targets batch MapReduce workloads; extension to streaming contexts via Apache Flink or Kafka Streams represents a promising direction for future work. Additionally, the current classification model is language-agnostic but has been validated only on English text; multilingual support is under investigation.

5. Conclusion

This paper has presented Content-Aware Partial Compression (CAPC), a novel framework for intelligent block-level compression in the Hadoop ecosystem. By combining TF-IDF weighting and Shannon entropy into a Semantic Density Score, CAPC routes each HDFS block to the most appropriate compression codec, achieving superior performance across all evaluated metrics without requiring modifications to existing MapReduce job logic. Experimental results on five diverse textual datasets demonstrate average compression ratio improvements of 13.4% over Gzip and 63.9% over Snappy, execution time reductions of up to 36.8%, storage savings averaging 69.6%, and near-linear throughput scaling to 32 cluster nodes. CAPC fills a critical gap in the Hadoop compression ecosystem by demonstrating that semantic-awareness at the block level can simultaneously improve storage efficiency, I/O throughput, and end-to-end analytical performance. Future work will explore deep learning-based block classification, streaming integration, and multilingual corpus support.

References

- [1] Lin J, Dyer C (2010) Data-intensive text processing with MapReduce. *Synth Lect Hum Lang Technol* 3(1):1–177
- [2] Abouzeid A, Bajda-Pawlikowski K, Abadi D, Silberschatz A, Rasin A (2009) HadoopDB: an architectural hybrid of MapReduce and DBMS technologies. *Proc VLDB Endow* 2(1):922–933
- [3] Rasheed A, Haq I, Baig MM (2019) Type-aware compression for mixed-format HDFS workloads. *J Grid Comput* 17(4):779–796
- [4] Moffat A, Turpin A (2002) *Compression and Coding Algorithms*. Kluwer Academic Publishers, Dordrecht
- [5] Dean J, Ghemawat S (2008) MapReduce: simplified data processing on large clusters. *Commun ACM* 51(1):107–113
- [6] Agarwal S, Diao Y, Lakshman A, Rajan S (2012) Online handling of join results on fast data streams. In: *Proc. ACM SIGMOD*, pp 1–12
- [7] Cohen J, Dolan B, Dunlap M, Hellerstein JM, Welton C (2009) MAD skills: new analysis practices for big data. *Proc VLDB Endow* 2(2):1481–1492
- [8] Doulkeridis C, Nørnvåg K (2014) A survey of large-scale analytical query processing in MapReduce. *VLDB J* 23(3):355–380
- [9] White T (2015) *Hadoop: the definitive guide*, 4th edn. O’Reilly Media, Sebastopol
- [10] Shannon CE (1948) A mathematical theory of communication. *Bell Syst Tech J* 27(3):379–423
- [11] Salton G, Buckley C (1988) Term-weighting approaches in automatic text retrieval. *Inf Process Manag* 24(5):513–523
- [12] Collet Y, Kucherawy M (2020) Zstandard compression and the application/zstd media type. RFC 8878, IETF
- [13] Mehrotra R, Sharma D, Gupta N (2022) Adaptive I/O scheduling for heterogeneous Hadoop workloads. *Future Gener Comput Syst* 130:178–194