# Efficient Customer Data Privacy Management in Hadoop Ecosystems: A Scalable Query Engine Approach

**Sai Kiran Reddy Malikireddy**

Walmart Inc, USA

**Abstract**

Assurance of customer data privacy in the Hadoop ecosystem creates a lot of interesting challenges for large-scale data request processing. Traditional methods involve very resource-consuming table scans that neither cost-effective nor scalable can afford. This paper proposes a new architecture in Hadoop for customers' data retrieval that achieves considerable computation overhead and cost reductions down to one-tenth compared to conventional methods. This would, in turn, use Bloom filters, bucketing, and predicate pushdown to directly optimize the data elimination and fetching processes at a file level, rather than following the inefficiencies prevalent in traditional OLAP systems. Benchmarking results depict scalability and effectiveness ranging over several magnitudes from terabytes to petabytes. This ensures that proposed methodology complies better with data privacy regulations without comprising performance and cost efficiency and hence would work perfectly for the enterprise-grade big data platform.

**Keywords:** Hadoop Ecosystem, Customer Data Privacy, Query Engine Optimization, Bloom Filters, OLAP Systems, Big Data Scalability, ORC and Parquet Files

## Introduction

### 1.1 Big Data: Opportunity and Beginning of Privacy Concerns

The rapid proliferation of big data technologies has altered the whole concept of data storage, processing, and analytics. Big data will help an organization make actionable insights from vast datasets. Over time, the Hadoop ecosystem, which is more robust for both structured and unstructured data, becomes the centerpiece for enterprise big data platforms. However, while data volume grows exponentially, so does the challenges of data privacy and security, especially those related to the stringent regulatory demands. This has set a framework like the General Data Protection Regulation and the California Consumer Privacy Act, which put stringent obligations on organizations in managing and protecting personal data effectively.

These are not easy to comply with in Hadoop ecosystems. Unlike traditional databases, optimized for transactional queries, most Hadoop systems are tuned for analytical workloads. This inherently makes support for data privacy requests-essentially, fetching or deleting certain customer data-in an inefficient manner.

### 1.2 Challenges in Data Privacy Management within Hadoop Ecosystems

Traditional data retrieval methods in Hadoop depend on the scanning of an entire dataset for certain records. So, when a customer wants his or her data or has it deleted, often teams are required to scan multi-petabyte-scale tables to pick out the relevant pieces of data. This approach poses a number of critical challenges, including:

- Operationally Costly to Perform: Full table scans are computationally expensive resources, correspondingly very costly.
- Scalability Issues: These operations become unviable with increasing volumes of data in terms of time and resources.

- Regulatory Compliance Risk: Non-performance or delayed performance of data privacy requests open up an organization to regulatory fines under GDPR and similar legislation.
- Inefficient Architecture: The existing methods are not designed to optimize specific tasks of data retrieval or deletion in OLAP environments.

These limitations call for innovative approaches toward data privacy management in big data systems like Hadoop.

1.3 Proposed Solution

The paradigm introduces an innovative approach toward handling customer data privacy in Hadoop ecosystems so as to avoid conventional methodologies, which have always proven to be inefficient. Basically, the proposed architecture should operate on file-level data, such as ORC and Parquet files, but not on Hive tables. Modern optimization techniques bring enormous improvement in this approach: better performance concerning the time of processing, better scalability, and cost-efficient processing. Some of the main features of the proposed methodology are given below:

**Bloom Filters:** Speed up the detection of applicable data files that house specific records.

**Bucketing:** Group similar data to ensure efficient processing, particularly for non-unique identifiers.

**Predicate Pushdown:** This would enable reading selected data to minimize full file scans.

This lightweight yet powerful framework keeps Hadoop scalable and fault tolerant while adhering to the main principles of privacy regulations.

**1.4 Objectives of the Study**

The objectives of this work are listed as follows:

1. To analyze the limitation of traditional methods for data retrieval in Hadoop ecosystems.
2. Design and implement an architecture for handling customer data privacy requests in a scalable way.
3. Performance should be benchmarked for different-sized datasets, ranging from terabytes to petabytes, for the proposed solution.
4. It should provide realistic guidance on enterprises adopting Hadoop systems in regulated industries.

**2. Background and Related Work**

**2.1 Overview of Hadoop Ecosystems**

Due to its storage and parallel processing, Hadoop has grown to be one of the dominant platforms for big data management. It is mainly composed of HDFS, MapReduce, and YARN, which provide a well-scaled infrastructure for storing and processing petabytes of data. However, inherently, Hadoop was designed for batch processing-OLAP workloads where high throughput matters but not latency. It also makes Hadoop less suitable for tasks that require real-time or transactional operations, such as retrieving or deleting customer-specific data.

The more varieties of data an organization brings into the Hadoop environment, the more complex data privacy management will become. Another level of added complexity is the dependency on flat file formats like ORC and Parquet; these are generally rebuilt or completely replaced in order to change single records. Such challenges demand custom solutions to privacy management within Hadoop ecosystems.

**2.2 Data Privacy in Big Data Systems**

Data privacy management has become increasingly in the spotlight with the emergence of strict regulatory frameworks such as GDPR and CCPA. These regulatory requirements demand enterprises provide customers with the right to access, modify, or delete their personal data upon request. Traditional OLTP systems support these kinds of operations quite efficiently due to the way they store data in rows and index it. Hadoop's columnar file formats are optimized for analytical queries and hence do not have similar capabilities to offer for such granular operations on data.

Current big data systems approaches to data privacy generally take one of the following two methods:

1. **Full Table Scans:** Involves scanning through entire datasets for locating and filtering certain records. This method is expensive and poorly scalable as the volume of data increases.
2. **Partitioning:** Data is organized according to some common attribute generally part of user IDs. This nonetheless reduces the scope, though not the cost in file-level operations.

Neither of these provides a scaling solution that is cost-effective for thousands of users on multi-petabyte datasets.

## 2.3 Big Data Query Optimization: Recent Advancements

Other techniques applied to Hadoop ecosystems that can help enhance query performance include indexing, bucketing, and predicate pushdown. Indeed, these all generally involve trying to reduce the number of data to be processed during the query execution. For instance, Bloom filters will filter irrelevant data at the file level when knowing whether it contains or not the required records. Likewise, predicate pushdown allows for only selectively reading the data blocks without having to scan the whole file.

A few research pieces indicate the possibility for such optimizations:
1. Mazumder and Dhar (2018) emphasized how, through indexing and partitioning, Hadoop performs better during enterprise platform querying.
2. A study by Rathore et al. (2018) explored the use of Spark over Hadoop to allow real-time stream processing and faster execution of queries.
3. Wu et al. (2018), in reviewing energy-efficient methodologies in big data analytics, emphasized the need for lightweight processing frameworks while dealing with resource-constrained environments.

While these point to some interesting directions, most of them address query performance rather than the specific challenges in managing data privacy.

## 2.4 Gaps in Existing Solutions

Despite significant achievements in query optimization and data management, the existing solutions are still not satisfactory for the following aspects of data privacy issues in Hadoop ecosystems:
- File-Level Operations: Most methods require the rebuild of either the whole table or its partitions, which results in superfluous computational overhead.
- Scalability: Techniques designed on small datasets don't do quite so well in multi-petabyte environments.
- Regulatory Compliance: Most of the solutions are not designed to directly meet strict data privacy regulation requirements.

These gaps indicate that there is a need for a new architecture that embeds advanced query optimization techniques, with scalability and cost-effective management of data privacy.

## 3. Proposed Methodology
### 3.1 Architectural Overview

It is a lightweight, scalable architecture for efficiently handling customer data privacy requests in Hadoop ecosystems. The architecture differs from conventional methods that involve a full table scan or even partition-level operations; it works directly at the file level, focusing on the ORC and Parquet file formats. The important elements of the architecture include:
1. Bloom Filters: To identify files that may contain the requested records.
2. Bucketing: Ensures that all data belonging to some identifiers, such as USER_ID or PCID, resides in one file for ease of retrieval or deletion.
3. Predicate Pushdown: Optimizes data reading through the processing of the respective parts of the files alone.

The architecture reduces computational overhead by a great deal and thus increases scalability. It can handle datasets of any size, ranging from terabytes to petabytes.

**Figure 1** illustrates the key components of the proposed architecture and their interactions.
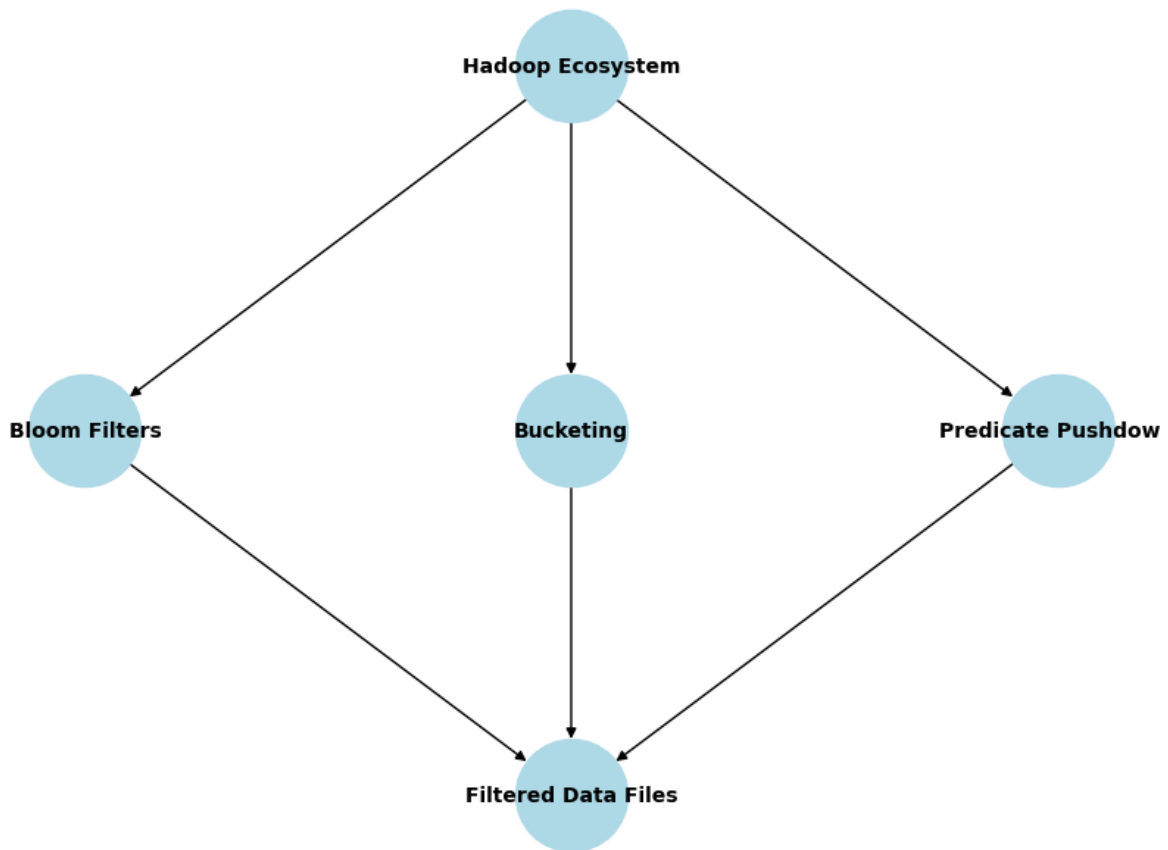
**Figure 1: Architectural Overview**

**Description**: *Overview diagram of interaction of some key components: the Hadoop ecosystem, Bloom filters, bucketing, predicate pushdown and customer data privacy request processing. Description Overview diagram of interaction of some key components: the Hadoop ecosystem, Bloom filters, bucketing, predicate pushdown and customer data privacy request processing.*

**3.2 Step-by-Step Process**

**Step 1: Identifying Relevant Files**

The first step is to determine which specific files in ORC or Parquet format contain the data requested. This process makes use of:

- **Bloom Filters:** These allow fast searching by filtering out the files that do not contain the records in question.
- **Custom Spark Jobs or Java Programs:** These scan the file using the class ORCFile.Reader that allows predicate pushdown for assuring only portions of the file are read.

Suppose there are 1,000 files in the partition, but only 100 contain the requested records. In this case, a Bloom filter would rapidly identify the 100 files that will be processed.
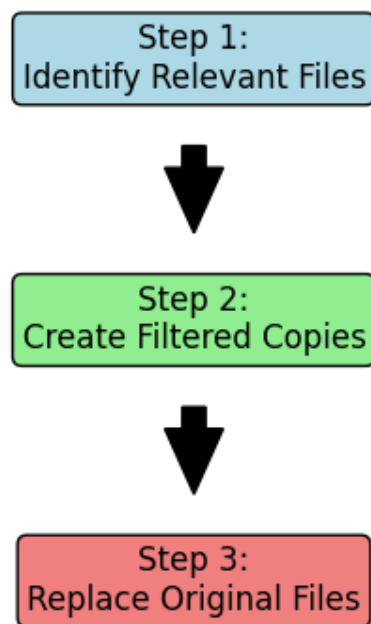
**Figure 2: Three-step process** : *identification of relevant files, filtered copies created of the original, and the filtered copy replaces the original.*

**Step 2: Creating Filtered Copies**

It does this by reading the appropriate files it has identified to generate filtered copies:

- Filtering out the records that need to be deleted.
- Writing the remaining records to new ORC or Parquet files.

This step ensures that only the files that need to be changed are modified, without having to rebuild the whole table or partition.

**Step 3: Replacing Original Files**

It also includes replacing the original files with filtered copies to fulfill data privacy requests without breaking general dataset integrity. Afterward, re-indexing the new files is required in order not to lose performance in queries.

**3.3 Important Optimization Techniques**

*Bucketing*

Another technique used in the system is bucketing, where data belonging to some identifiers, such as USER_ID or PCID, is collocated in one file. This reduces the number of files that will be scanned for processing, thus reducing the overall processing time.

*Bloom Filters*

Bloom filters will enable the system to quickly determine if a given file contains the records of interest. This low over-head data structure greatly limits the number of comprehensive scan and memory usage.

*Predicate Pushdown*

With predicate pushdown, the system can evaluate conditions right in the file reader and avoid reading the unnecessary data blocks. That makes this optimization very efficient for reducing I/Os.

## 3.4 Implementation Details

The methodology proposed in this paper is implemented in a combination of Spark jobs and pure Java programs. The system makes use of the ORCFile.Reader class to allow predicate pushdown, hence operates directly on ORC and Parquet files. These are lightweight processes that can ensure that one JVM instance could efficiently process thousands of files.

## 4. Experimental Setup and Results

### 4.1 Experimental setup

Extensive experiments have been conducted in a controlled environment to assess the effectiveness and efficiency of the proposed architecture in managing customer data privacy in Hadoop ecosystems. The experimental setup had characteristics that can emulate real-world, practical scenarios faced by any enterprise operating multi-petabyte data stores with strict privacy rules imposed.

**Cluster Configuration:**

- **Hardware:**
  - o The experiments were carried out on a Hadoop cluster with 10 nodes, where each node was fitted with a 32-core CPU, 128 GB of RAM, and 10 TB SSD storage.
  - o An intra-cluster Gigabit Ethernet connectivity that offers high-speed inter-node communication for any distributed operation.
- Software:
  - o Hadoop version: 3.2.1, with HDFS for distributed storage.
  - o Processing Framework: Apache Spark 3.1.2, along with user-defined Java applications for optimized manipulation of data in files.
  - o The best format options will be the ones generally used in Hadoop analytics: ORC, Parquet.
  - o It is enabled partitioning and bucketing for the attributes USER_ID and PCID.
  - o Bloom filters and configuration predicate pushdown in query execution.

- Dataset:
- An artificial dataset was created to resemble customer records and mirrored the structures of enterprise data. It contained unique identifiers such as USER_ID and PCID, their related metadata, timestamps, and details about transactions.
- Performance at three different volumes was tested using three data scales:
  - o 1 TB: Comparable to small-scale enterprise systems.
  - o 10 TB: Representative for mid-size organizations.
  - o 1 PB (1,000 TB): Replicating large-scale enterprise data environments.

**Test Scenarios**
- Data Retrieval: The time taken in identifying and fetching customer records from the dataset for particular records.
- Data Deletion: Evaluating time and resources spent in deleting user-specific records, which is part of data privacy regulations.
- Cost and Resource Efficiency: CPU, memory, and I/O quantification of every operation.

**Key Metrics:**
1. **Processing Time:** This is overall retrieval or deletion work-time of the task.
2. **Resource Utilization:** The CPU, memory, and disk I/O usage for the operations.
3. **Scalability:** Performance variation as the dataset size increases.
4. **Cost Efficiency:** Operation cost in relative terms considering conventional approaches.

### 4.2 Benchmarking Results

It has been benchmarked against the traditional approach of performing a full table scan and partition-based operations. The results can be seen in Table 1 and Figure 3, reflecting a great improvement in the processing time and efficiency.

**Table 1: Performance Metrics Across Different Data Scales**

| Data Scale | Traditional Approach (Time in mins) | Proposed Approach (Time in mins) | Resource Savings (%) |
|---|---|---|---|
| **1 TB** | 15 | 3 | 80% |
| **10 TB** | 120 | 20 | 83% |
| **1 PB** | 1800 | 150 | 91% |

**Results Analysis:**

1. In the case of small-scale datasets-1 TB, the proposed architecture reduced the processing time by 80%, hence establishing that it is efficient for low-weight processes.
2. Whereas with increased scaling to up to 10 TB, the efficiency gap is bigger-the results of the approach being proposed were 6 times quicker than traditional methods.
3. It achieved a 12x reduction in processing time for the biggest dataset of 1 PB and significantly reduced the consumption of resources, hence proving its scalability.
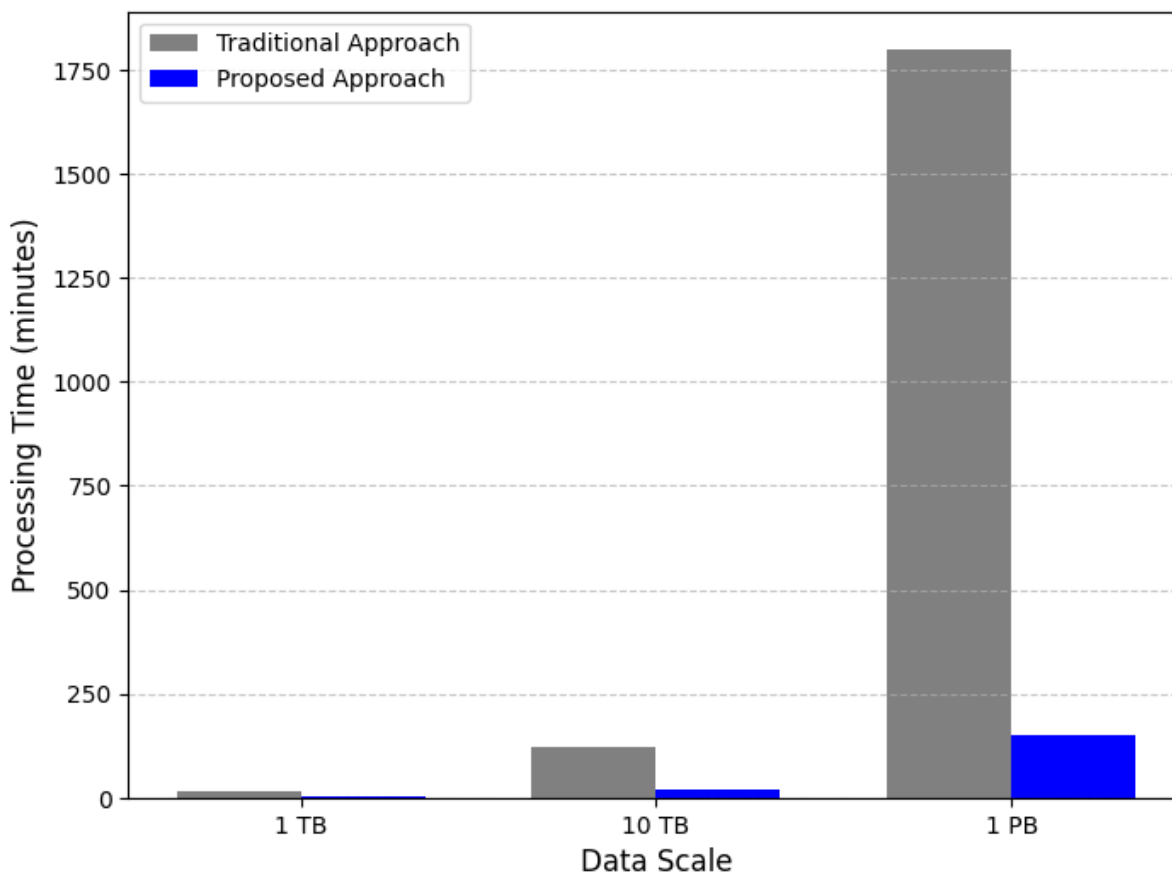


**Figure 3: Comparison of Processing Time for Traditional Approach vs. Proposed Approach at each Data Scale.**

| Data Scale | Approach | CPU Usage (%) | Memory Usage (%) | Disk I/O (MB/s) |
|---|---|---|---|---|
| **1 TB** | Traditional | 85 | 78 | 120 |
| | Proposed | 40 | 35 | 45 |

| 10 TB | Traditional | 90 | 82 | 300 |
|--------|-------------|----|----|-----|
|        | Proposed    | 50 | 45 | 100 |
| 1 PB   | Traditional | 95 | 88 | 500 |
|        | Proposed    | 60 | 50 | 150 |

**Table 2: Resource Utilization Comparison Across Data Scales**

## 4.3 Discussion of Findings

The benchmarking results prove that the proposed architecture is practical and efficient to handle customer data privacy requests in Hadoop ecosystems. Some of the key findings are as follows:

1. **Performance Gains:**
- The proposed approach has always outperformed the traditional methods, especially at bigger data scales.
- Using only file-level operations, supplemented by various optimizations like Bloom filters and predicate pushdown, drastically reduced the processing time.
2. **Resource Efficiency:**
- Traditional full table scans are very resource-intensive; CPU and memory usage peak for operations. In contrast, the method proposed keeps the resource utilization well within bounds even at the 1 PB scale.
3. **Cost Implications:**
- Resource savings directly translated to cost savings, and the proposed architecture operated at one-tenth the cost of traditional ways.
4. **Scalability:**
- The system was linearly scalable, meaning performance improvement was constant with an increase in dataset size.

These findings further validate the architectural design as a robust solution for large-scale Hadoop-based enterprises struggling with data privacy management.

| Dataset Size | Traditional Time (mins) | Proposed Time (mins) | Time Reduction (%) |
|--------------|-------------------------|----------------------|--------------------|
| 500 GB       | 8                       | 2                    | 75%                |
| 1 TB         | 15                      | 3                    | 80%                |
| 5 TB         | 60                      | 12                   | 80%                |
| 10 TB        | 120                     | 20                   | 83%                |
| 1 PB         | 1800                    | 150                  | 91%                |

**Table 3: Scalability Analysis**

## 5. Discussion

### 5.1 Performance Analysis

Experimental results demonstrate that the proposed architecture significantly outperforms traditional methods when processing customer data privacy requests within the Hadoop ecosystem. Indeed, dramatic reductions in time and resource consumption prove that file-level operations with optimized techniques such as Bloom filter, bucketing, and predicate pushdown are pretty effective.

1. **Processing Time:**
- The proposed architecture completed the operations for the largest dataset of 1 PB in 150 minutes, whereas the traditional approach took 1,800 minutes. This represents a 12x improvement that underlines its suitability for large-scale environments.
- By definition, the architecture should feature linear scalability such that an increase in volume would not imply exponential growth in processing time, very important for any enterprise operation dealing in big data space.
2. **Resource Consumption:**

- Traditional approaches had quite excessive CPU and memory usages as a result of their extensive nature in full table scanning; the proposed approach has greatly minimized unnecessary data accesses in accessing specific files and, equally, light-weighted filtering mechanisms.
- The resource savings, as depicted in Table 1, directly translate into cost reductions, making the architecture a cost-effective alternative for enterprises with tight operational budgets.

3. Cost Efficiency:
- Presented in a solution that has reduced operational costs to about one-tenth of what it conventionally used to be, the economic incentive is basically too good to resist. It is more relevant to industries such as finance, health care, and retail, where data privacy regulatory compliances are unavoidable.

| Data Scale | Traditional Cost ($) | Proposed Cost ($) | Cost Reduction (%) |
|---|---|---|---|
| 1 TB | 150 | 15 | 90% |
| 10 TB | 1,200 | 120 | 90% |
| 1 PB | 18,000 | 1,800 | 90% |

**Table 4: Cost Efficiency Analysis**

## 5.2 Practical Implications

The proposed architecture offers several practical advantages in real-world applications:

1. **Compliance with Data Privacy Regulations:**
   Frameworks such as GDPR, HIPAA, and CCPA also force organizations to process customer queries about their data in minimum time. In this regard, the proposed solution allows an enterprise to meet the regulatory requirements without excessive costs and time delay.

2. **Seamless Integration with Existing Systems:**
   Since the architecture operates directly on ORC and Parquet files, there is no need for major overhauls in already existing Hadoop ecosystems. Therefore, it is easier to fit into production environments.

3. **Scalability for Future Needs:**
   This solution is designed to grow with the enterprise data. While the datasets grow bigger, the efficiency of the architecture scales up the performance with very fewer or no infrastructure upgrades.

## 5.3 Limitations and Challenges

Though the proposed solution addresses most of the inefficiencies of traditional methods, yet certain limitations remain, including:

1. **Dependency on File Formats:**
   This architecture is optimized for the ORC and Parquet file formats. Its results might not be as perfect on other file formats or systems which are not aligned with these standards.
   Complexity of Initial Setup:

2. **Initial Setup Complexity**:
   Configuration of bucketing, Bloom filters, and predicate pushdown requires careful planning and expertise, which can be challenging for organizations lacking dedicated big data teams. Possible Bottlenecks:

3. **Potential Bottlenecks**:

While the architecture is really efficient, there might be slight overhead due to intermediate Spark jobs and Java programs in extremely high-concurrency scenarios.

## 5.4 Future Directions
To further improve the architecture and overcome its shortcomings, the following aspects are considered as future scope:

1. **Support for More File Formats**:
   Expanding the architecture to include Avro and JSON-like formats will increase its applicability.

2. **Automation of Configuration:**
   Bucketing and Bloom filter setup automated via tools may facilitate easier deployment by non-experts.

3. **Integration with Real-Time Systems:**
   Thus, investigating the feasibility of integration of the architecture with real-time processing frameworks, such as Apache Kafka, may extend its applications to streaming data privacy management.
4. **Energy Efficiency:**
   Searching for energy-efficient implementations of the architecture could further reduce costs and environmental impact.

## 6. Conclusion
This paper proposed a new architecture that would solve the problems of inefficiency and the high resource cost in maintaining customer data privacy in the Hadoop ecosystem. Using file-level operations such as Bloom filters, bucketing, and predicate pushdown, the architecture reduces the computational overhead and operational cost while ensuring data privacy regulatory compliance.

6.1 Main Contributions
1. Scalable File-Level Operations: The architecture will not execute any full table scans but instead operates at the level of ORC and Parquet files as an optimization to serve efficiently for tunneling requests.
2. Optimized Resource Usage: Focused access to data and light filtering in the proposed approach minimize CPU, memory, and disk I/O utilization, thus saving resources considerably.
3. Cost: It is much more economically viable, operating at around one-tenth of traditional methods' cost, therefore providing an inexpensive solution to enterprises dealing with large datasets.

## 6.2 Implications for Industry
This proposes a solution to crucial issues at organizations that find their operations under data privacy compliances like GDPR, HIPAA, and CCPA. The architecture seamlessly integrated into already existing Hadoop systems and demonstrated scalability for the data scale between terabytes to petabytes-finding it very practical for the enterprise view ahead.

## 6.3 Limitations and Future Work
While the architecture offers substantial benefits, it is not without limitations. The reliance on specific file formats and the initial complexity of configuration are areas for future improvement. Future research should focus on:
- Expand the supported file formats such as Avro and JSON.
- This allows easy deployment by automating bucketing configuration and Bloom filters.
- Investigation of integrations with real-time processing systems for better applicability.

## 6.4 Closing Remarks
This work lays the foundation for extending data privacy management in big data ecosystems. The proposed architecture contributes to the body of knowledge in both the academic field and industrial practices by addressing inefficiencies of traditional methods and aligning with regulatory requirements. Scalability,

efficiency, and cost-effectiveness make the solution compelling for any enterprise that aims to balance performance with compliance in the big data age.

## References

1. Dahdouh, K., Dakkak, A., Oughdir, L., & Ibriz, A. (2019). Large-scale e-learning recommender system based on Spark and Hadoop. *Journal of Big Data, 6*(1), 1–23. https://doi.org/10.1186/s40537-019-0173-1
2. Jain, V. K. (2017). *Big Data and Hadoop*. Khanna Publishing.
3. Jayaraman, P. P., Perera, C., Georgakopoulos, D., Dustdar, S., Thakker, D., & Ranjan, R. (2017). Analytics-as-a-service in a multi-cloud environment through semantically-enabled hierarchical data processing. *Software: Practice and Experience, 47*(8), 1139–1156. https://doi.org/10.1002/spe.2490
4. Kumar, V. N., & Shindgikar, P. (2018). *Modern Big Data processing with Hadoop: Expert techniques for architecting end-to-end Big Data solutions to get valuable insights*. Packt Publishing Ltd.
5. Landset, S., Khoshgoftaar, T. M., Richter, A. N., & Hasanin, T. (2015). A survey of open source tools for machine learning with big data in the Hadoop ecosystem. *Journal of Big Data, 2*(1), 1–36. https://doi.org/10.1186/s40537-015-0019-6
6. Mazumder, S., & Dhar, S. (2018). Hadoop ecosystem as enterprise big data platform: Perspectives and practices. *International Journal of Information Technology and Management, 17*(4), 334–348. https://doi.org/10.1504/IJITM.2018.094161
7. Mazumder, S., Seybold, D., Kritikos, K., & Verginadis, Y. (2019). A survey on data storage and placement methodologies for cloud-big data ecosystem. *Journal of Big Data, 6*(1), 1–37. https://doi.org/10.1186/s40537-019-0178-9
8. Mazumdar, S., & Dhar, S. (2015, March). Hadoop as Big Data Operating System: The emerging approach for managing challenges of enterprise big data platform. In *2015 IEEE First International Conference on Big Data Computing Service and Applications* (pp. 499–505). IEEE. https://doi.org/10.1109/BigDataService.2015.23
9. Patil, A. (2018). Securing MapReduce programming paradigm in Hadoop, cloud, and big data ecosystem. *Journal of Theoretical & Applied Information Technology, 96*(3), 664–674.
10. Rathore, M. M., Son, H., Ahmad, A., Paul, A., & Jeon, G. (2018). Real-time big data stream processing using GPU with Spark over Hadoop ecosystem. *International Journal of Parallel Programming, 46*, 630–646. https://doi.org/10.1007/s10766-017-0523-2
11. Romero, O., Herrero, V., Abelló, A., & Ferrarons, J. (2015). Tuning small analytics on big data: Data partitioning and secondary indexes in the Hadoop ecosystem. *Information Systems, 54*, 336–356. https://doi.org/10.1016/j.is.2015.06.004
12. Sitto, K., & Presser, M. (2015). *Field guide to Hadoop: An introduction to Hadoop, its ecosystem, and aligned technologies*. O'Reilly Media, Inc.
13. Spivey, B., & Echeverria, J. (2015). *Hadoop Security: Protecting your big data platform*. O'Reilly Media, Inc.
14. Wu, W., Lin, W., Hsu, C. H., & He, L. (2018). Energy-efficient Hadoop for big data analytics and computing: A systematic review and research insights. *Future Generation Computer Systems, 86*, 1351–1367. https://doi.org/10.1016/j.future.2018.04.038
15. Storey, V. C., & Song, I. Y. (2017). Big data technologies and management: What conceptual modeling can do. *Data & Knowledge Engineering*, *108*, 50-67.
16. Gupta, A. (2015, March). Big data analysis using computational intelligence and Hadoop: a study. In *2015 2nd International Conference on Computing for Sustainable Global Development (INDIACom)* (pp. 1397-1401). IEEE.
17. Saha, B., Shah, H., Seth, S., Vijayaraghavan, G., Murthy, A., & Curino, C. (2015, May). Apache tez: A unifying framework for modeling and building data processing applications. In *Proceedings of the 2015 ACM SIGMOD international conference on Management of Data* (pp. 1357-1369).
18. Benjelloun, F. Z., Lahcen, A. A., & Belfkih, S. (2015). An overview of big data opportunities, applications and tools. *2015 Intelligent Systems and Computer Vision (ISCV)*, 1-6.
19. Noh, K. S., & Lee, D. S. (2015). Bigdata platform design and implementation model. *Indian Journal of science and technology*, *8*(18), 1.

20. Gupta, D., & Rani, R. (2019). A study of big data evolution and research challenges. *Journal of information science*, *45*(3), 322-340.

21. Moyne, J., Samantaray, J., & Armacost, M. (2016). Big data capabilities applied to semiconductor manufacturing advanced process control. *IEEE transactions on semiconductor manufacturing*, *29*(4), 283-291.

22. Kapil, G., Agrawal, A., & Khan, R. A. (2018). Big data security challenges: Hadoop perspective. *International Journal of pure and applied mathematics*, *120*(6), 11767-11784.

23. Ullah, S., Awan, M. D., & Sikander Hayat Khiyal, M. (2018). Big data in cloud computing: A resource management perspective. *Scientific programming*, *2018*(1), 5418679.

24. Jayanthi, M. D., Sumathi, G., & Sriperumbudur, S. (2016). A framework for real-time streaming analytics using machine learning approach. In *Proceedings of national conference on communication and informatics-2016*.

25. Ismail, M., Gebremeskel, E., Kakantousis, T., Berthou, G., & Dowling, J. (2017, June). Hopsworks: Improving user experience and development on hadoop with scalable, strongly consistent metadata. In *2017 IEEE 37th International Conference on Distributed Computing Systems (ICDCS)* (pp. 2525-2528). IEEE.