

# Separation Is for Better Reunion: Data Lake Storage at Huawei

## ABSTRACT

Abstract here.

## 1 INTRODUCTION

Intro here.

## 2 MOTIVATION

## 3 ARCHITECTURE

## 4 STREAM AND TABLE STORAGE OBJECT

In this section, we introduce the stream object and table object, purpose-built storage abstractions designed for efficient storage and access of stream and table data in the storage layer.

### 4.1 Stream Object

The stream object is a storage abstraction in the store layer that efficiently supports key-value message streaming at scale. It stores a partition of key-value pairs for continuous message streams, organized as a collection of data slices. Each slice can contain up to 256 records as depicted in Figure 2. Incoming message records are appended to a specific slice in a stream object based on its topic, key, and offset.

**Stream objects operations.** The stream object operates similarly to the block and file storage abstractions, providing read and write functionality for stream storage. Figure 1 outlines key operations supported by the stream object, including creating and destroying a stream object with functions `CreateServerStreamObject` (line 1-3) and `DestroyServerStreamObject` (line 4-5) respectively. The `*option` field (line 2) sets storage configurations, such as data redundancy methods (replicate or erasure code) and I/O quotas, so as to ensure enterprise-level reliability and performance. The assigned `objectId` (line 3) serves as a unique identifier for operating the stream object. The `AppendServerStreamObject` function appends incoming records to the stream object and returns the starting offset of the appended records. The `ReadServerStreamObject` function reads the stream object starting from a specified offset, with control conditions such as the length of the read specified in the `readCtrl` field. Since the message service is designed to support real-time streaming, it is set to respond to all subsequent messages unless specified limits by the user. `IO_CONTENT_S` (line 8 and 14) is a data structure that provides non-blocking I/O by using buffers to enhance the performance of both writing and reading operations. **Write stream messages.** We discuss how to write messages into StreamLake and endure enterprise-level load-balanced and redundant persistence for the stream objects, which is achieved on the basis of SSD and HDD storage pools. As shown in Figure 2, the messages are first assigned to stream object slices based on topics, keys, and offsets (Figure 2-a,b,c). Then, a distributed hash table is leveraged to ensure even data distribution for load-balance storage (Figure 2-d). Specifically, data slices will be distributed evenly to

```
1 int32_t CreateServerStreamObject(  
2     IN CREATE_OPTIONS_S *option,  
3     OUT object_id_t *objectId);  
4 int32_t DestroyServerStreamObject(  
5     IN object_id_t *objectId);  
6 int32_t AppendServerStreamObject(  
7     IN object_id_t *objectId,  
8     IN IO_CONTENT_S *io,  
9     OUT uint64_t *offset);  
10 int32_t ReadServerStreamObject(  
11     IN object_id_t *objectId,  
12     IN uint64_t offset,  
13     IN EAD_CTRL_S *readCtrl,  
14     INOUT IO_CONTENT_S *io);
```

Figure 1: Stream Object Operations.

4096 logical shards, each of which has the storage space managed by persistence logs (PLog, Figure 2-e)<sup>1</sup>. Each PLog unit is a collection of persistence services in OceanStor [] that controls a fixed amount of storage space on multiple disks and provides 128 MB of addresses per shard. When a message is received, the PLog unit replicates it to multiple disks for redundancy (Figure 2-f). Key-value databases<sup>2</sup> serve as indexes for PLogs for fast record lookup.

### 4.2 Table Object

We also extend the storage object layer in StreamLake to support table-like<sup>3</sup> operations for more effective data storage and management, similar to lakehouses []<sup>4</sup>. The table storage uses an open lakehouse format<sup>5</sup> with optimizations<sup>6</sup> for faster metadata access. The table abstraction is logically defined by a directory of data and metadata files, as shown an example in Figure 3.

**Data directory.** Table objects are stored in Parquet files of the data directory. In this example, the table is partitioned based on the date column, so the data objects are separated into different sub-directories by date. Each sub-directory name represents its partition range<sup>7</sup>. The data objects in each Parquet file are organized as row-groups<sup>8</sup> and stored in a columnar format for efficient data analysis. Footers in the Parquet files contain statistics to support data skipping within the file.

**Metadata directory** keeps track of the file paths of the table, table schema, and transaction commits etc., which are organized into three levels: commit, snapshot, and catalog, as shown in Figure 3-(b, c, d).

**Commits** are Arvo files that contain file-level metadata and statistics such as file paths<sup>9</sup>, record counts, and value ranges for the data

<sup>1</sup>-Plog 4096

<sup>2</sup>-can we name it specifically?

<sup>3</sup>- a term?

<sup>4</sup>\* "similar to lakehouse" is weird

<sup>5</sup>-a term? what does it mean?

<sup>6</sup>-what? a part in Fig4? catalog

<sup>7</sup>-partition key?

<sup>8</sup>-a term?

<sup>9</sup>@directory path

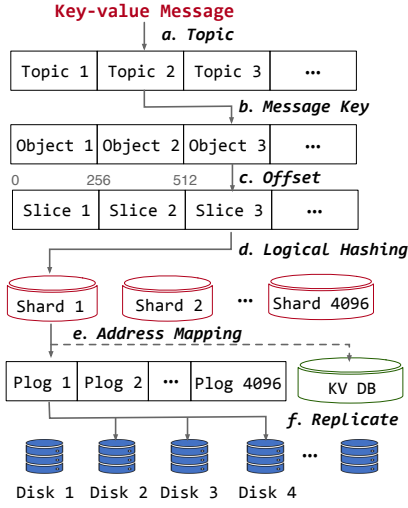


Figure 2: Write Message to StreamLake.

objects. Each data insert, update, and delete operation will generate a new commit file to record changes to the data object files.

*Snapshots* are index files that index valid commit files for a specified time period. These snapshots document commit statistics such as current files, row count and added/removed files/rows as data operation logs. Along with commits, snapshots provide snapshot-level isolation to support optimistic concurrency control. Readers can access the data by reading from the valid commit files, while changes made by a writer will not be visible to readers until they are committed and recorded in a snapshot. This allows for multiple readers and one writer to access the data simultaneously without the need for locks.

Snapshots also monitor the expiration of all commits, making them essential for supporting time travel. Time travel queries allow data to be viewed as it appeared at a specific time. By keeping old commits and snapshots, the table object enables the use of a timestamp to look up the corresponding snapshot and commits, so as to access to historical data.

*Catalog* describes the table object, including the profile data such as the table ID, directory paths, schema, snapshot descriptions, modification timestamps, etc. The data and metadata files are stored in the table directory, except for the catalog, which is stored in a distributed key-value engine<sup>10</sup> optimized for RDMA and SCM to ensure fast metadata access. The data and metadata files are converted to PLogs in the underlying storage for redundant persistence as discussed above.

## 5 STREAMLAKE DATA PROCESSING

In this Section, we present the data processing services in the data layer. Driven by practical application scenarios discussed in Section 2, these services provide a comprehensive, enterprise-level data lake storage solution to efficiently store and process log messages at scale. The StreamLake services encompass a stream storage system for message streaming (Section 5.1), lakehouse-format read/write



Figure 3: File Organization of StreamLake Table Objects.

capabilities for efficient tabular data processing (Section 5.2), and support for query operator computation pushdown(Section 5.3).

### 5.1 Message Streaming

We develop a distributed stream storage engine that facilitates message streaming at large scale. Our engine leverages the stream object storage abstraction to ensure enterprise-level reliability and scalability.

**Overall architecture of streaming service.** The high-level design of the stream service is shown in Figure 4. The stream storage system comprises of producers, consumers, stream workers, stream objects, and a stream dispatcher, which work together to provide seamless message streaming.

<sup>CC</sup>[which techniques to ensure reliability and scalability? Here, can we summarize something different (our characteristic) ?]

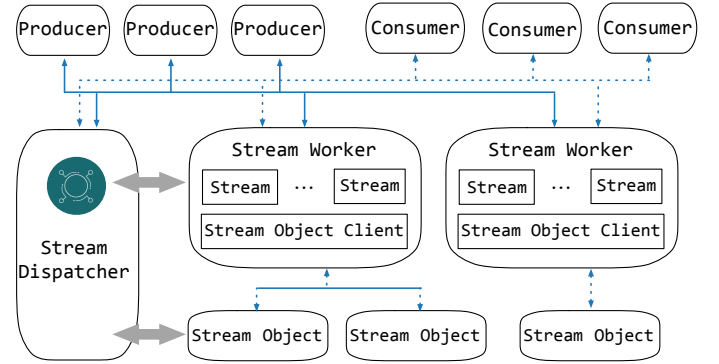


Figure 4: Write Message to StreamLake.

*Producers and Consumers.* Producers are responsible for publishing messages to topics, which are named resources for categorizing streaming messages. Consumers, located downstream, subscribe to these topics to receive and process the published messages. To ensure seamless integration with existing open-source message streaming services used by our customers in production environments, the producer and consumer message APIs are designed to be compatible with the open-source de facto standard. This maximizes connectivity with the ecosystem, allowing users to easily migrate their applications to StreamLake with minimum costs. Figure 5 demonstrates the process of writing and reading messages using the producer and consumer APIs. In this example, a producer writes a new message “Hello World” as a key-value pair to a topic named

<sup>10</sup>—can we name it? or cite

“topic\_streamlake\_test”. The consumer then subscribes to this topic and processes published messages.

```

1  /*Sample producer code*/
2  Producer producer = new Producer();
3  Message msg = new Message("Hello world");
4  producer.send("topic_streamlake_test", msg);
5  /*Sample consumer code*/
6  Consumer consumer = new Consumer();
7  consumer.subscribe("topic_streamlake_test");
8  While (true) {
9  /*Poll for new data*/ }

```

Figure 5: Sample code of Producer and Consumer.

*Stream workers* work together with stream objects discussed in Section 4.1 to tackle stream processing and message storage. The number of stream workers is determined by configurations and the physical resources allocated to the stream storage. Each stream worker is capable of handling multiple streams<sup>11</sup> and a single stream object client. When a topic is created, streams are added to the stream workers in a round-robin manner to ensure even distribution and workload balancing across the cluster.

Each stream is mapped to a unique stream object in the storage layer, which is a storage abstraction customized to key-value message streaming. The stream object offers efficient interfaces and implementations for writing and reading streams from the storage pools. The persistence process is detailed in Figure 2.

The task of message delivery is carried out by stream object clients, which monitor the stream objects. These clients unwrap messages from clients, encapsulate them in the stream object data format, and redirect them to the corresponding stream objects via RDMA. To guarantee message delivery, the clients actively monitor the health of the stream objects to which they are connected and regularly exchange critical service data with the dispatcher service. This synchronization process includes reporting the health of the stream object connections and refreshing the stream objects connected to by the client.

*Stream dispatcher.* The stream dispatcher is responsible for managing the metadata and configurations of the messaging service, and directing external and internal requests to the appropriate resources for message dispatch. The relationships among topics, streams, stream workers, and stream objects are stored as key-value pairs in a fault-tolerant key-value store within the stream dispatcher. When there is a status change (e.g., a stream worker or topic is added or removed), the metadata in the key-value store is updated immediately to refresh the topology tracking. This topology tracking aids the stream dispatcher in directing requests for message stream dispatch. When there is a producer or consumer connection request, the stream dispatcher will route the request to the appropriate stream worker based on the associated stream topic, establishing a direct message exchange channel between the producer, the stream worker, and the consumer.

The stream dispatcher also sets configurations for the messaging service in the unit of topic<sup>12</sup>. An example of configurations is shown in Figure 6.

```

1  { "stream_num" : 3,
2    "quota" : 106,
3    "scm_cache" : true,
4    "convert_2_table" : {
5      "table_schema" : { ... },
6      "table_path" : ...,
7      "split_offset" : 107,
8      "split_time" : 36000,
9      "delete_msg" : false,
10     "enabled" : true }
11   "archive" : {
12     "external_archive_url" : null,
13     "archive_size" : 262144,
14     "row_2_col" : true,
15     "enabled" : true } }

```

Figure 6: Stream Storage Configuration Example..

- The *stream\_num* configuration sets the parallelism of a topic, which should be provided during topic declaration. In the example, three streams are created for the topic and they are evenly distributed among stream workers to process messages in parallel.
- The *quota* configuration sets the maximum processing rate for each stream. In the example, each stream can process up to 10<sup>6</sup> messages per second.
- The *scm\_cache* configuration enables the use of SCM<sup>13</sup> caches.
- The *convert\_2\_table* configuration enables the automatic conversion of stream object messages to table object records. When it is set, a background process will apply the *table\_schema* to convert messages to table object records periodically and save them in *table\_path*, i.e., the table object directory. The conversion is triggered by either an accumulation of 10<sup>7</sup> messages or the passing of 36000 seconds.
- The *archive* configuration automates the archiving of historical data to meet business and regulatory requirements. Data can be stored in the cost-effective StreamLake archive storage pool or exported to an external storage system specified in the *external\_archive\_url* configuration. The *archive\_size* configuration denotes the data volume in MB that triggers archiving, and the *row\_2\_col* configuration determines whether the data is archived in a columnar format.

**CC** [The StreamLake stream storage provides guaranteed delivery, efficient transfer, and high elasticity for enterprise use.<sup>14</sup>]

**Delivery Guarantee:** Our system ensures consistent message delivery through several measures. (1) Data within a stream object is strictly ordered, ensuring that messages are consumed in the order in which they are received. (2) Message writing is idempotent, which means that for network failure, duplicate messages sent by the producer can be identified. (3) Strong data consistency is achieved by eliminating unreliable components like file systems and page caches, and storing data in stream objects that can tolerate node, network, and disk failures. (4) The system provides exactly-once<sup>15</sup> semantics through the use of a transaction manager and the two-phase commit protocol. This tracks participant actions and ensures that all results in a transaction are visible or invisible at the same time.

<sup>11</sup> streams? Above we talk about messages

<sup>12</sup>@ what is the connection with the above? I want to ask the relationship among topic, stream, messages

<sup>13</sup>no full name

<sup>14</sup>We should build connection with the above designs.

<sup>15</sup>@

**Efficient Transfer:** Our system implements several mechanisms to efficiently transfer data. First, Stream workers and stream objects are connected through a data bus<sup>16</sup> with RDMA, which reduces the switch overhead in the TCP/IP protocol stack. Second, an I/O aggregation mechanism is used to aggregate small I/O requests and increase throughput. This function can be disabled for latency-sensitive<sup>17</sup> scenarios. Finally, a local cache is implemented at the stream object client to speed up message consumption.

**High Elasticity:** Our system provides high elasticity by decoupling data storage and data serving. The number of stream workers can be adjusted without data migration, and the mapping between stream workers and stream objects can be updated to reflect the changes in a matter of seconds. This allows the message streaming service to easily scale up or down to accommodate changes in service demand.

## 5.2 Lakehouse Read and Write

StreamLake also provides support for concurrent tabular data reads and writes, similar to the architecture of lakehouse []. This section describes the storage conversion from stream messages to tabular records<sup>18</sup>, as well as the implementation of key lakehouse operations.

<sup>CC</sup>[What is the advantage of this conversion?]

**Stream-to-table conversion.** This process is performed by a background service and results in the conversion of records in stream objects to table objects, allowing for efficient downstream processing, which is triggered by the `convert_2_table` configuration in Figure 6, which include the table schema and bounds for data freshness<sup>19</sup> in the downstream processing. The table schema must be specified at the topic declaration, as it determines the expectations for field types and values across all messages. To effectively leverage the storage, users<sup>20</sup> can choose to just retain messages in crucial topics as stream objects to support real time applications while converting most <sup>CC</sup>[stream] data to table objects. The reverse conversion, from table records to stream messages, is also supported for data playback. As shown in Figure ?? and Table ??, this design provides a good trade-off between the system cost and latency, which also helps to decouple the data processing from the business logics.

<sup>CC</sup>[What is the connection between the two paragraphs?]

Our StreamLake services implement lakehouse read/write operations using a table object and high-performance caches and computation pushdowns, which eliminate unnecessary data transmission and accelerate concurrent data reads and writes.<sup>21</sup> In the rest of this subsection, we will introduce the implementation of key read/write operations in details.

**CREATE TABLE:** This operation begins by registering the table information, such as the schema, path, database, and table name, in the catalog. The `/data` and `/metadata` directories are then created under the table path. Then table configurations (schema, partition

spec<sup>22</sup>, target file size, etc) are written to the metadata directory for persistence.

**INSERT:** This operation includes the persistence of data and metadata, as well as caching of metadata to the non-volatile memory (NVM), which is introduced to combine small I/O accesses to the underlying storage pools.

(a) *Data persistence:* Records are written directly to the persistent layer as parquet files in the corresponding partition path under the table root directory.

(b) *Metadata caching:* Metadata updates are mostly small I/O operations.<sup>23</sup> To avoid generating significant number of small files, we leverage a global write cache to aggregate the metadata updates, which is achieved through the following steps: (b-1) Each added parquet file generates a commit record containing file-level metadata and descriptions. All new commit records<sup>24</sup> are written to the write cache as key-value pairs when a commit is made. (b-2) The latest snapshot will be read from the persistence layer to the cache and its commit data will be updated<sup>25</sup>. (b-3) The snapshot descriptions and version history in the catalog are also read from the persistence layer and overwritten by adding the latest snapshot description.

(c) *Metadata persistence:* Metadata in the NVM write cache is asynchronously flushed to the persistent storage pool when the buffer is full. A metadata management process (MetaFresher) transforms the commits and snapshots from key-value pairs to files and writes them to the table/metadata directory.

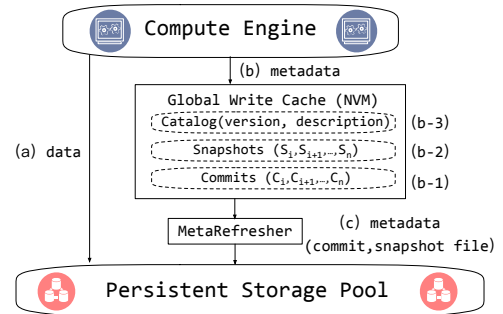


Figure 7: Write Cache Acceleration in Lakehouse Read/Write.

**SELECT:** The select operation first reads the catalog to retrieve the table profile for collecting the list of snapshot files needed for this query, such as the metadata version and snapshot descriptions. Then the corresponding snapshots and commit metadata are read from both the NVM cache and the persistent storage pool to generate the latest complete snapshots and commit metadata. When all the record file addresses are confirmed, data is read from the persistence pool by read tasks.

**DELETE:** The delete operation begins with a select operation to find files containing records that match the filtering conditions. There are two cases to consider: If the filtering conditions match all data in several partitions, only the metadata will be updated, and a new

<sup>16</sup>\*reviewer suggest to extend

<sup>17</sup>low latency?

<sup>18</sup>why discuss the conversion here

<sup>19</sup>bounds?

<sup>20</sup>\*why not automatic?

<sup>21</sup>who optimize who?

<sup>22</sup>?

<sup>23</sup>\*There seems a lack of summarization like why, benefit and connection.

<sup>24</sup>a commit includes multiples records?

<sup>25</sup>is there an arrow from storage to b2?



commit version will be generated by eliminating the information of deleted partitions. If the filtering conditions only match some files, these files will be read, and the data matching the filtering condition will be deleted. Computation pushdowns<sup>26</sup> are applied to process file reads and writes without<sup>27</sup> data transmission to/from the compute engines.

**UPDATE:** Similar to the delete operation, update operation also uses a select statement to identify records that match the specified conditions. Optimizations, such as pushdowns<sup>28</sup>, are applied to reduce<sup>29</sup> data movements during the file read and write processes.

**Drop Table:** There are two types of drop table operations: (1) Drop table soft unregisters the table from the catalog but retains the table's metadata and data in the persistent layer for potential future restoration. To restore a soft-deleted table, a new table can be created and linked to the original table path, effectively registering the deleted table back to the catalog.<sup>30</sup> (2) Drop table hard removes both the metadata (files under /metadata) and data (files under /data) of the table and clears the table from the catalog. Note that some of the metadata may have been written to the acceleration cache during the drop table hard operation and will be flushed to the persistent layer asynchronously in the background. The operation to delete the metadata will first clear it from the cache, and then delete it from the disk.

### 5.3 Query Operator Computation Push Down

In this subsection, we introduce the computation pushdown to reduce the amount of data transfer between the storage engine of StreamLake and the query engines<sup>31</sup>. It is built on top of an elastic, serverless engine<sup>32</sup> in the data service layer. Here, we choose serverless computing as our execution model due to it is lightweight and flexible, which allow us to quickly start a large number of instances for computation tasks near the data sources, and free the resources as soon as the tasks are completed. The elasticity is important since CPU resources can be scarce during critical data management jobs.

The main components of the serverless function engine are shown in Figure 9, which include function dispatcher, worker instances, a worker manager, and a function repository. The function dispatcher schedules jobs and manages workflows, and worker instances execute the tasks. The worker manager oversees server resources and the life cycle of worker instances, and the function repository registers and stores function images. These modules work together to support elastic serverless computing.

<sup>CC</sup>[pushdown what?]

To be specific, when a job request is received, the function dispatcher obtains the data location from the storage devices and selects the appropriate storage nodes based on data distribution and available resources. The worker manager is then requested to deploy worker instances to the selected node. The worker instances download the necessary functions from the repository and execute

the jobs using data from the storage infrastructure. Upon completion, a callback message is sent to the caller to notify them of the results. To ensure service quality, elastic scaling policies are used to dynamically adjust the number of nodes based on workload. For example, a load balancing method<sup>33</sup> is used to balance scheduling among instances.

To achieve maximum pushdown benefits, three categories of query operators are supported,

- Projection Pushdown: Only selected columns will be returned.
- Filter Pushdown: Only rows satisfying the filtering conditions will be returned.
- Aggregate Pushdown: The results of aggregate functions such as Count, MAX and AVG will be returned.

These operators are selected because the size of their output could be significantly smaller than the input, and thereby a large optimization opportunity can be achieved. These query operators are implemented as separate functions, which are registered and executed in the serverless engine service. The implementation allows for sharing and reuse of the same query operator function by different query engines, as long as its image<sup>34</sup> is registered in the serverless engine.

<sup>CC</sup>[need connection]

To facilitate the pushdown of query operators from the compute layer to the storage cluster in StreamLake, we have introduced the NDP Proxy<sup>35</sup>. During query planning, the optimizer of query identifies operators that can be pushed down and sends the information to the NDP Proxy. Then the NDP Proxy inserts these requests into a queue for traffic control and then sends them to the serverless engine for execution as shown in Figure 9. As soon as the query results are ready, the NDP Proxy transfers them back to the compute engine using the high-speed data exchange bus of StreamLake, so as to ensure efficient data transfer. Overall, this process leverages the flexibility and elasticity of serverless computing, which allows for efficient use and release of server resources as required, resulting in better overall performance and resource utilization.

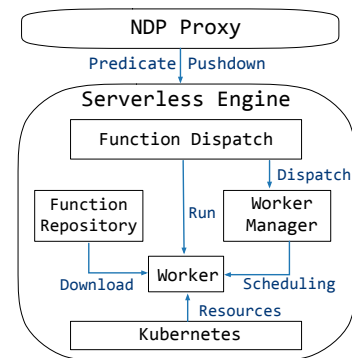


Figure 8: Serverless Function Engine.

<sup>26</sup>next section, s?

<sup>27</sup>without?

<sup>28</sup>same

<sup>29</sup>reduce or without?

<sup>30</sup>@

<sup>31</sup>we can support many query engines? Will we mention it?

<sup>32</sup>@

<sup>33</sup>not specific

<sup>34</sup>@

<sup>35</sup>why need ndp, the full name?

## 6 LAKEBRAIN OPTIMIZATION

Optimizing query processing over large-scale data is significant in data warehouse and big data systems, as discussed in []. However, for the StreamLake system with complicated storage-disaggregated architecture with multiple compute engines, it is challenging to optimize the end-to-end performance and resource usage. The reasons are two-fold. First, it is hard to capture the entire environment about the compute and storage cluster as well as the queries executed by other engines simultaneously. Second, even though all the environment data is available, it is still hard to optimize because of the large search space due to the large number of tunable and interdependent variables [].

To address the challenges, we present LakeBrain, a novel data lake storage optimizer that complements end-to-end data pipeline optimization. Unlike query engine optimizers that focus on join ordering and cardinality estimation [], LakeBrain aims to optimize data usage in storage<sup>36</sup> during query execution, which is key to improve both query performance and storage resource utilization in a storage-disaggregated design.

<sup>CC</sup>[Do you really design to achieve e2e optimization? what is the intuition of e2e?]

On the one hand, in a streaming application scenario, data ingestion and transactions often result in numerous small files, leading to low query performance on merge-on-read (MOR)<sup>37</sup> tables. LakeBrain designs the automatic compaction (Section 6.1) to combine these small files into fewer and larger ones, so as to improve inter-cluster storage and network usage as well as query performance<sup>38</sup>.

<sup>CC</sup>[Below is very hard to follow! what is talking about?where is the partition? what is the connection with the above?]

LakeBrain’s design is kept simple for ease of extension and support for different applications. It consists of three components: a statistics collector, the core optimization logic, and an executor. The statistics collector gathers system configurations, environment variables, and workload history, while the core optimization logic employs heuristic rules, probabilistic models, and machine learning algorithms to suggest the best strategy candidates. The executor then deploys the chosen strategy, with its effects being collected as feedback by the statistics collector for future optimization.

Next, we will illustrate the above aspects in detail.

### 6.1 Automatic Compaction

As discussed above, file compaction aims to find the optimal strategy for compacting files that can improve query execution time or increased block utilization<sup>39</sup> in storage. To achieve this goal, the optimization process employs two algorithms<sup>40</sup>: particle swarm optimization (PSO)<sup>41</sup> and reinforcement learning (RL).

PSO is used to search for the global optimum, a population-based method that doesn’t require assumptions about the relationship between tunable parameters and query performance. The goal is to obtain an approximately optimal solution within a limited time.

<sup>36</sup> seems coarse. what does it mean?

<sup>37</sup> @

<sup>38</sup> messy, query performance, storage resource, which one?

<sup>39</sup> not mentioned above

<sup>40</sup> two algorithms or steps?

<sup>41</sup> is there a citation?

<sup>CC</sup>[Not clear!! We should say what is the entire problem, and why and how to split it into two steps.]

<sup>CC</sup>[Also not clear!! We should introduce what is the problem and the optimization goal that fit the paradigm of RL. Then define the critical components of RL and illustrate that.] RL secondly finds a more sophisticated policy based on the states of the data lake environment. The optimization process involves considering a set of discrete compaction configurations within the action space. Since the state space is continuous, a function approximation method is preferred. The stability of the training process is critical due to the high degree of variability in query performance in a distributed environment, so proximal policy optimization (PPO) is applied.

A deep neural network (DNN) is used to approximate the policy and the value function, with a shared feature backbone network that covers both global and local characteristics of the states. The output from the feature network is processed by two fully connected networks to compute the policy output and the action value. The actor and critic are alternatively updated after collecting new trajectories using the latest policy during training.

Once a desired result is obtained, the numerical output of the compaction strategy is translated into actionable operations by the data lake connector for a specific data lake engine, facilitating the optimization process.

### 6.2 Predicate-aware Fine-grained Partitioning

Optimizing data partitioning<sup>42</sup> involves assigning records to storage blocks in the most efficient manner possible, thereby reducing the number of blocks accessed during queries. Our partitioning approach is based on the query-tree framework [43], and utilizes a sum-product network (SPN) probabilistic model [19, 26, 31] to model the distribution of the data in LakeBrain. This is done in order to ensure fast inference speed and avoid repeatedly scanning the datasets.

<sup>CC</sup>[The high level idea of query-tree framework and spn should be introduced.]

<sup>CC</sup>[how to leverage the distribution and why?]

The query-tree framework creates a tree-based partitioning strategy using pushdown predicates. Each leaf node represents a partition, and its column ranges are derived from the pushdown predicates used to split its parent nodes. By using probabilistic models to characterize the dataset, we can identify the most suitable partitioning policy. The probabilistic model-based cardinality estimation, as demonstrated in Figure 11, is used to estimate the number of records in each partition, <sup>CC</sup>[How?? query-driven or data-driven? both cases should have queries to test] instead of scanning the original data, thereby saving a significant amount of time. This greatly improves the speed of the partitioning optimization algorithm, making it suitable for large-scale systems.

Additionally, probabilistic models allow us to represent a sequence of datasets with a series of probabilistic models that have a fixed structure but varying parameters. This is achieved by representing a sequence of datasets as a series of multi-dimensional

<sup>42</sup>\* we should first say what is the data partition problem (does the previous section mention it?)

vectors, each representing the learnable variables in the probabilistic model with a fixed length, i.e. a time series. We can then use time series prediction methods to predict future probabilistic models, and use these predictions to estimate the number of records in a partition during partitioning optimization<sup>CC</sup>[like a magic].

<sup>CC</sup>[Not coherent with the Next!]

To implement the optimized data layout, we introduce a partitioning mechanism that saves data in fine-grained partitions based on the partitioning strategy. Additionally, we have implemented an evaluator that skips irrelevant partitions by checking the overlap between pushdown predicates and the column ranges in each partition. For numerical columns, the range can be represented as lower and upper bounds, which are well-handled by many data formats. For categorical columns, we either record its range or its complement using "IN" or "NOT IN" predicates. The effectiveness of this predicate-aware partitioning approach is evaluated in section 7.2, where the test results show exceptional performance.

## 7 EXPERIMENT

### 7.1 Experimental Settings

**Our Experimental Scenario.** To demonstrate the performance of the StreamLake framework, we analyze a simplified<sup>43</sup> real-world use case. The case compares StreamLake framework with an open-source storage solution to build a big data processing pipeline that can facilitate business analysis. Specifically, a mobile financial application company collaborates with a mobile carrier to collect and analyze its app usage data. The company aims to understand its app usage patterns to prevent frauds and enhance its product experience. The mobile carrier provides this analytic service through an end-to-end big data processing pipeline. This pipeline includes several jobs such as data collection, normalization, labeling, and querying, as depicted in Figure 10.

*(a) Collection:* The network carrier collects mobile app data packets in data centers across the nation via deep packet inspection (DPI) and transfers them to a centralized storage pool.

*(b) Normalization:* At the storage pool, the data packets are normalized as records in a unified schema. Data is validated to ensure accuracy and quality. Sensitive data is shielded to protect privacy.

*(c) Labeling:* Labels from knowledge bases are added, so as to classify the records and identify useful insights.

*(d) Query:* After the normalization and labeling processes are completed, the records are inserted into tables and are available for query engines. To perform analyses, the app company employs secure API calls to query the data. Figure 9 illustrates an example SQL query that counts the daily active users (DAU) in different provinces. More complicated analysis, like hidden Markov and Gaussian Mixture Models<sup>44</sup>, can also be applied to draw user profiles and identify abnormal activities.

To support both full data and real time analyses, the network carrier builds two data flows in the pipeline. One flow processes full data in batch every two hours and the other processes stream messages constantly to deliver time-sensitive logs such as new logs, payments and password modifications. This ensures that

the network carrier can effectively analyze both historical data and real-time events to make accurate and timely decisions.

```
1 Select COUNT(*) as DAU
2 From TB_DPI_LOG_HOURS '
3 Where url = 'http://streamlake_fin_app.com'
4     and start_time >= 1656806400 --July 3rd, 2022
5     and start_time < 1656892800 --July 4th, 2022
6 Group By consumer = province;
```

Figure 9: Query Example of Computing DAU.

**Settings.** This use case is evaluated in a commodity<sup>45</sup> cluster using different sizes of input data packets and the results are compared with open-source storage solution Hadoop Distributed File System (HDFS) [] and Kafka []. The reason of why we choose the two storage systems is that in reality, China Mobile has been using them for many years, which have shown stable and good performance. Hence, it is reasonable to directly compare with the systems that our customer (China Mobile) is using. Also, in practice, as we know, many other companies also use HDFS and Kafka to cope with similar application scenarios. <sup>CC</sup>[For above, do we need more justification? like why hdfs and kafka fit? or common sense?]

To be specific, the cluster hardware<sup>46</sup> consists of 3 nodes, each with 24 2.30 GHz cores and 256 GB RAM. The cluster is configured as a 3-node StreamLake when we measure it. While running the open-source solution, it is configured to host a 3-node HDFS storage and a 3-node Kafka cluster simultaneously. The number of input data packets varies: 10 million, 50 million, 100 million, 500 million, and 1 billion packets. Each packet has an average size of 1.2 KB, resulting in corresponding data volumes of 12 GB, 60 GB, 120 GB, 600 GB, and 1.2 TB, respectively.

Overall, Figure 10 shows the data processing process. Kafka and HDFS serves as independent stream storage and batch storage respectively to pass data across collection, normalization, labeling and query jobs. As a typical ETL practice, a new copy of all data is written to HDFS and Kafka after each job. In case it<sup>47</sup> fails accidentally, a job can read its input data to reproduce the results.

In our solution, StreamLake serves as a unified stream and batch processing storage. <sup>CC</sup>[It reads messages from the data collection jobs and passes messages and aggregated batches to the same stream and batch processing engines in the normalization, labeling and query jobs.]<sup>48</sup> As StreamLake supports time travel, only updated rows are written to the storage. When a job needs to re-run, it can use time travel to retrieve its input data. During the query jobs, for example, the three filters in the WHERE clause and the COUNT aggregate in Figure 9 are pushed down to compute in StreamLake, so as to accelerate the query.

### 7.2 Overall Comparison

Table 1 shows the results. The numbers of input data packets are in the top row. The storage usage and processing time for StreamLake

<sup>43</sup>why simplified?

<sup>44</sup>why mention models?

<sup>45</sup>@

<sup>46</sup>cluster hardware?

<sup>47</sup>who?

<sup>48</sup>too many and

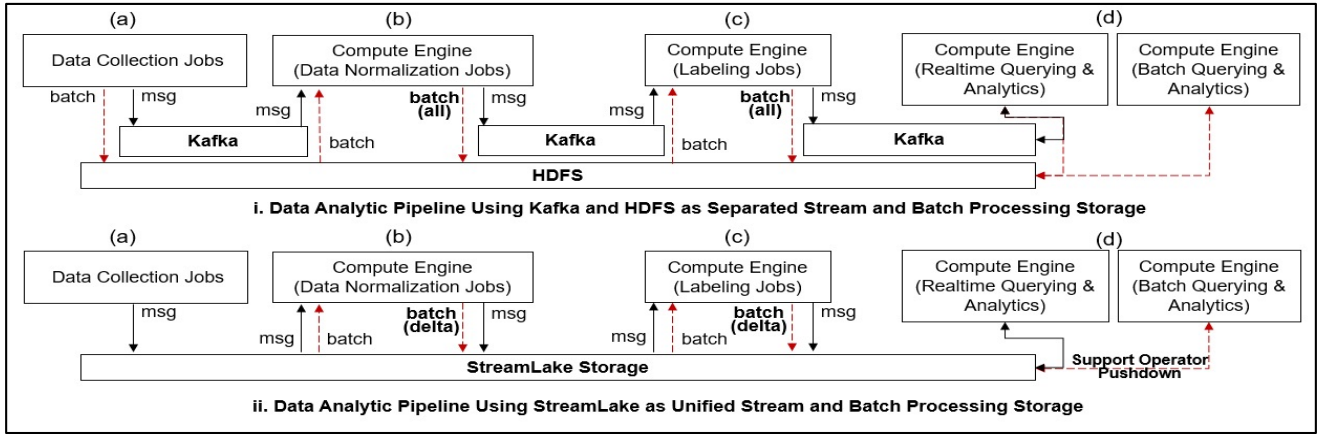


Figure 10: Data Analytic Pipelines for a Real-world Use Case.

	#-Data Packet	10,000,000	50,000,000	100,000,000	500,000,000	1,000,000,000
Storage Space Usage (GB)	StreamLake	34	166	329	1,659	3,289
	HDFS + Kafka	145	729	1451	6,901	13,816
	Ratio (HK/S)	4.33	4.38	4.40	4.16	4.20
Stream Processing Speed (Messages/Second)	StreamLake	301,522	417,303	518,065	530,077	546,987
	Kafka	302,611	413,613	527,826	531,021	539,893
	Ratio (K/S)	1.00	0.99	1.02	1.00	0.99
Batch Processing Total Time (Second)	StreamLake	259	664	1173	4868	9646
	HDFS	212	795	1548	7535	14771
	Ratio (H/S)	0.82	1.19	1.32	1.55	1.53

Table 1: StreamLake v.s. HDFS and Kafka.

(S), HDFS (H), Kafka (K) are in the following rows. The “Ratio” represents that the ratio between HDFS (Kafka) and StreamLake with respect to the storage usage or time. Note that HK denotes the sum of the storage usage in HDFS and Kafka.

The experiment demonstrates that StreamLake significantly improves the total storage usage and the batch processing time. The storage usage in the HDFS and Kafka is 4 times as much as StreamLake. The reason is that in HDFS and Kafka, full data is written into the storage when each ETL job is finished, which is a common practice to support downstream jobs restart after unexpected failures. As a result, six copies of full data are written into the storage. While for our StreamLake, since the storage natively supports time travel, we only save one copy of full data plus updates in each ETL job, saving about 75% storage usage.

The batch processing speed in StreamLake is better than HDFS when the workload is 50 million records or more. As the workload grows, the advantage of skipping irrelevant partitions becomes significant. StreamLake is 50% faster than HDFS when the workloads are 500 million and 1 billion records. On the other hand, StreamLake may not be the best choice for small workloads. When the workload is 10 million records, StreamLake is 20% slower than HDFS as it performs extra metadata management.

CC[Why not remove the first column??]

CC[For the above, can we correspond to the designs with respect to previous sections as the reasons of improvement]

The message stream processing speed in StreamLake is competitive to Kafka. StreamLake and Kafka process about 300 thousand messages per second when the workload is 10 million records. Both systems scale to process about 500 thousand messages per second when the workloads are 100 million and more.

### 7.3 Evaluation of Message Streaming

To quantitatively measure the message streaming service as an independent stream storage, we conduct an experiment to evaluate its throughput, latency, elasticity and volume. We select OpenMessaging [] as our benchmark framework as it is widely used to compare messaging platforms. A cluster with three nodes is used in this experiment for ease of reproduction. To help better understand the impact of tiered storage, two sets of hardware configurations are tested. In the first set of hardware (Set-1), each node has 10 CPU cores, 128 GB RAM and 800 GB NVMe SSD, 3 PB SAS HDD and all the nodes are connected with 10 GB <sup>CC</sup>[ethernet]. In the second set of hardware (Set-2), all the configurations are the same except that each node has additional 16 GB persistent memory to serve as an extra cache. Messages are sent from producers to consumers in a fixed size of 1 KB. The data volumes we process are 100 TB, 500 TB and 1 PB respectively.

Figure ?? shows the results. As the messages to process increase from 50000 per second to 1.5 million per second, the system throughput increases linearly, reaching a peak of 1.2 GB/s with a workload



of 1.3 million message per second. Set-1 and Set-2 achieve almost the same throughputs, indicating that it does not improve the throughput to add persistent memory as a cache. However, as shown in Figure ??(b), persist memory reduces the latency as we expect, especially when the workload is 200k messages per second or less. Figure ??(c) shows the high elasticity of the stream storage. The service gracefully scales from 1000 to 10000 partitions in less than 10 seconds. The good scalability demonstrates a significant advantage of the data centric and disaggregated storage design<sup>49</sup>.  
<sup>CC</sup>[Below is hard to follow!] Finally, Figure ??(d) compares the volumes of different storage strategies. Without scarifying the reliability, StreamLake provides the option to use erasure coding<sup>50</sup> and column-store which can offer three to five times of volume compared to standard storage with one or more replicas.

## 7.4 Evaluation of LakeBrain

In this part, we evaluate the two components in LakeBrain, *i.e.*, auto-compaction and predicate-aware partitioning.

**Auto-Compaction:** To precisely evaluate the effectiveness of our automatic compaction strategy, a TPC-H based test bed<sup>51</sup> is set up to ingest data from the message streaming platform to the data lake storage, during which a compaction strategy is tested. We run the experiment with 24 GB to 90 GB data and three compaction strategies are deployed: (1) No compaction. (2) A static strategy which simply compacts data files in a 30 second interval. (3) Auto-compaction, respectively. During the ingestion, multiple rounds of TPC-H queries are executed in parallel to obtain their end-to-end performance. As shown in Figure ??(a), the results depict how much improvement of query performance that the compaction strategies can make, compared with baselines. We can observe that the auto-compaction strategy outperforms the static one for all data volumes. As the data volume increases, the advantage becomes more significant.<sup>CC</sup>[Why?]

In addition to the query performance, we also evaluate the block utilization of the auto compaction. Specifically, we control the file ingestion speed such that we can generate different number of files to measure both the run time and the block utilization in different workloads. The run time is evaluated along with the block utilization because an ideal strategy should improve the utilization without scarifying the performance. Similar to above, we deploy three methods: (1) No compaction, (2) The static strategy compacts data files in a fixed time interval<sup>52</sup>, and (3) Auto-compaction. We can observe that the auto-compaction outperforms the static strategy in term of block utilization.

<sup>CC</sup>[why?]

<sup>CC</sup>[Ingestion speed? Below is hard to follow!]

When we deploy the auto-compaction, the system is able to identify good compaction opportunities in which there are many small files and both the file ingestion speed and the block utilization are relatively low.

File ingestion speed is important because compaction commits will fail if there are file access conflicts. As a comparison, it is hardly to avoid unnecessary or unsuccessful compactions in the static

compaction strategy hence its performance is less ideal. Figure ??(b) summarizes the results of all three test groups. Compared with no compaction and the static strategy, our method performs better in term of both block utilization and query run time.

**Predicate-Aware Partitioning:** We also tested the partitioning method on the TPC-H test bed with different scale factors. We train the probabilistic model with 3% of the data randomly sampled from the lineitem table in a dataset generated with a scale factor of 2. After that, we obtain the optimized partitioning policy with the proposed method and evaluate our system on the full dataset with scale factors of 2, 5, 10 and 100. To evaluate the performance, we compare the resulting bytes skipped for lineitem table with 1) no partition (full), 2) partitioning by the day of l\_shipdate (day), and 3) our proposed method using sum-product networks (spn). We compare the results with partitioning by the day of l\_shipdate considering it appears frequently in the pushdown predicates for lineitem table. The workload includes TPC-H query 6, 12 and 19 which involve lineitem table and include predicates other than l\_shipdate. We skipped the other TPC-H queries because their performance is driven mainly by multiple tables joining performance which is beyond our purpose.

The results presented in Figure 15(c,d) shows that the proposed method obtains non-marginal performance gains in terms of both bytes scanned and the runtime. The fine-grained partitioning is superior on the queries in terms of data skipping compared to partitioning by the day of l\_shipdate because the optimized partitioning policy split the data based on other predicates except l\_shipdate. Even though the runtime for the queries are dominated by table joining, the optimized partitioning also demonstrated some improvements for query 6 and query 19, considering we only optimize the partition of the lineitem table.

<sup>CC</sup>[Why not reverse as previous sections?]

## 7.5 Evaluation of Query Pushdown

In this experiment we evaluate the query operator pushdown method which we believe can provide stable query runtime regardless of network conditions. This is significant in the real-world deployment as it is not always an option to upgrade the data center network. In fact, many customers who we have worked with did the opposite to ask us to reuse their existing network to reduce the overall upgrade costs. Hence, it is a critical design that the query operator pushdown method ensures the query processing time and the application service level agreements even with a constraint network bandwidth.

Two groups of clusters with different network bandwidths are used in this experiment, one with 10 Gb bandwidth and the other with 1 Gb. To precisely assess the benefits, we carefully select three live queries with data intensive operations and 4.8 TB data from a China Mobile production environment. Two different query engines, Hive and Presto, are deployed to process the SQL queries for generalization. It is observed that in the test group without query operator pushdown, the query performance varies widely across engines. When the queries are executed in the 10 Gb ethernet, Presto completes the jobs in about 900 seconds while Hive takes around 1200 seconds. When network bandwidth drops to 1 Gb, all the execution time soars to over 3000 seconds. As a comparison, we applied query operator pushdown in the second test group. The

<sup>49</sup>term?

<sup>50</sup>term?

<sup>51</sup>@ a term?

<sup>52</sup>\*no need to repeat, which interval?

runtime of all the queries is close to 900 seconds with less than 10 difference, regardless of compute engines and network bandwidths. This means a 4 times performance advantage when the engine process queries in a 1 Gb bandwidth network. In summary, we can conclude that our generalized query operator pushdown method introduces stable and high performance to query processing in a storage-disaggregated architecture.

## **7.6 China Mobile Use Case**

## **8 RELATED WORK**

## **9 CONCLUSION**

## **REFERENCES**