

Druid: A Real-time Analytical Data Store

Fangjin Yang, Eric Tschetter, Gian Merlino, Nelson Ray, Xavier Léauté, Deep Ganguli,
Michael Driscoll
{fangjin, cheddar, gian, nelson, xavier, deep, mike}@metamarkets.com

ABSTRACT

Druid is an open source¹ data store designed for exploratory real-time analytics. The system combines a column-oriented data layout, a distributed, shared-nothing architecture, and an advanced indexing structure to allow for the arbitrary exploration of billion-row tables with sub-second latencies. In this paper, we describe Druid’s architecture, and detail how it supports fast aggregations, flexible filters, and both streaming and batch data ingestion.

1. INTRODUCTION

Enterprises routinely collect diverse data sets that can contain up to terabytes of information per day. Companies are increasingly realizing the importance of efficiently storing and analyzing this data in order to increase both productivity and profitability. Numerous commercial database systems (e.g., IBM’s Netezza [30], HP’s Vertica [5], EMC’s Greenplum [25]) and several research papers [4, 7, 11] offer solutions for how to store and extract information from large data sets. Traditional data warehouse solutions often utilize Relational Database Management Systems (RDBMS) and key-value NoSQL stores to house information about business operations. The traditional data ingestion and analysis procedure is often slow and traditional storage systems may not always scale with data volume and complexity. Modern data analysis requires data insights as soon as new data is made available. Moreover, businesses increasingly rely on high performance applications on top of their data warehouses to visualize and interpret data for critical decisions.

Druid was built and open sourced to directly address the need for a real-time analytical data store to power modern data intensive applications. The system works best for exploring high volume, high dimension, event data. These data sets are typical in OLAP [28] workflows. As an example, an online advertiser may use Druid to determine ad click through rates based on user dimensions such as age, gender, and location. An online video provider may use Druid to determine the top watched videos in the last 5 minutes based on video length, content, and language. Druid is designed to answer questions about data, and because the system is freely available, users are constantly discovering new use cases for various types of data.

Druid shares some similarities with main-memory databases [13] and interactive query systems such as Dremel [24] and PowerDrill [16]. Druid’s focus is fast aggregations, arbitrarily filters, and low-latency data ingestion. Queries on in-memory data typically complete in milliseconds, and real-time data ingestion means that new events are immediately available for analysis. Furthermore, Druid

is highly configurable and allows users to adjust fault tolerance and performance properties. The system is used to power applications at several enterprises and is the core engine of the Metamarkets data analytics platform.

In this paper, we make the following contributions:

- We describe Druid’s low latency data ingestion implementation.
- We detail how the architecture enables fast multi-dimensional data exploration.
- We present Druid performance benchmarks.

The outline is as follows: Section 2 describes the Druid data model. Section 3 presents an overview of the components of a Druid cluster. Section 4 outlines the query API. Section 5 describes the data storage format in greater detail. Section 6 discusses Druid robustness and failure responsiveness. Section 7 presents performance benchmarks. Section 8 discusses related work and Section 9 presents our conclusions.

2. DATA MODEL

The fundamental storage unit in Druid is the segment and replication and distribution are always done at a segment level. Each data table in Druid (called a *data source*) is partitioned into a collection of segments, each typically comprising 5–10 million rows. A sample table containing advertising data is shown in Table 1. Many core Druid concepts can be described using this simple table.

Druid always requires a timestamp column as a method of simplifying data distribution policies, data retention policies, and first-level query pruning. Druid partitions its data sources into well-defined time intervals, typically an hour or a day, and may further partition on values from other columns to achieve the desired segment size. For example, partitioning the data in Table 1 by hour results in two segments for 2011-01-01, and partitioning the data by day results in a single segment. The time granularity to partition segments is a function of data volume and data time range. A data set with timestamps spread over a year is better partitioned by day, and a data set with timestamps spread over a day is better partitioned by hour.

Segments are uniquely identified by a data source identifier, the time interval of the data, a version string that increases whenever a new segment is created, and a partition number. This segment metadata is used by the system for concurrency control; read operations always access data in a particular time range from the segments with the latest version identifier for that time range.

¹<https://github.com/metamx/druid>

Table 1: Sample Druid data

| Timestamp | Publisher | Advertiser | Gender | Country | Impressions | Clicks | Revenue |
|----------------------|-------------------|------------|--------|---------|-------------|--------|---------|
| 2011-01-01T01:00:00Z | bieberfever.com | google.com | Male | USA | 1800 | 25 | 15.70 |
| 2011-01-01T01:00:00Z | bieberfever.com | google.com | Male | USA | 2912 | 42 | 29.18 |
| 2011-01-01T02:00:00Z | ultratrimfast.com | google.com | Male | USA | 1953 | 17 | 17.31 |
| 2011-01-01T02:00:00Z | ultratrimfast.com | google.com | Male | USA | 3194 | 170 | 34.01 |

Most segments in a Druid cluster are immutable *historical* segments. Such segments are persisted on local disk or in a distributed filesystem (“deep” storage) such as S3 [10] or HDFS [29]. These historical segments have associated metadata describing properties of the segment such as size in bytes, compression format, and location in deep storage. Data for intervals covered by historical segments can be updated by creating new historical segments that obsolete the old ones. This process is described in greater detail in Section ??.

Segments covering very recent intervals are mutable *real-time* segments. Real-time segments are incrementally updated as new events are ingested, and are available for queries throughout the incremental indexing process. Periodically, real-time segments are converted into historical segments through a finalization and handoff process described in Section 3.3.

2.1 Incremental Indexing

Druid is best used for aggregating event streams, and both historical and real-time segments are built through an incremental indexing process that takes advantage of this assumption. Incremental indexing works by computing running aggregates of interesting metrics (e.g. number of impressions, sum of revenue from the data in Table 1) across all rows that have identical attributes (e.g. publisher, advertiser). This often produces an order of magnitude compression in the data without sacrificing analytical value. Of course, this comes at the cost of not being able to support queries over the non-aggregated metrics.

3. CLUSTER

A Druid cluster consists of different types of nodes, each performing a specific function. The nodes come together to form a full working system. The name Druid comes from the Druid class in many role-playing games: it is a shape-shifter, capable of taking many different forms to fulfill various different roles in a group. The composition of a Druid cluster is shown in Figure 1.

Recall that the Druid data model has the notion of historical and real-time segments. The Druid cluster is architected to reflect this conceptual separation of data. Real-time nodes are responsible for ingesting, storing, and responding to queries for the most recent events. Similarly, historical compute nodes are responsible for loading and responding to queries for historical events. Real-time nodes and historical compute nodes are jointly classified as *storage* nodes.

Queries to access storage node data will typically first hit a layer of broker nodes. Broker nodes are responsible for finding and routing queries down to the storage nodes that host the pertinent data. The storage nodes compute their portion of the query result in parallel and return the results to the brokers. Broker nodes, compute nodes, and realtime nodes are classified as *queryable* nodes. All queryable nodes expose the same query interface.

Druid also has a set of coordination nodes to manage load assign-

ment, distribution, and replication. Coordination nodes are not queryable and instead focus on maintaining cluster stability. Coordination nodes have an external dependency on a MySQL database.

Druid relies on Apache Zookeeper [17] for coordination. Most intra-cluster communication is over Zookeeper, although queries are typically sent over HTTP.

3.1 Apache Zookeeper

Zookeeper is a service for coordinating processes of distributed applications. Zookeeper provides connecting applications an abstraction of a hierarchy of data nodes known as *znodes*. Each *znode* is part of a hierarchical namespace, similar to file systems. Zookeeper that has the concept of ephemeral and permanent *znodes*. Permanent nodes must be created and destroyed explicitly by a connecting application. Ephemeral *znodes* can be created by connecting applications and deleted either explicitly or if the session that created the *znode* is terminated (such as in the event of service failure).

3.2 Historical Compute Nodes

Historical compute nodes (or just compute nodes) are the main workers of a Druid cluster and are self-contained and self-sufficient. Historical compute nodes load historical segments from deep storage and exposes them for queries. Deep storage is used for permanent segment storage and compute nodes are used for intermediate segment storage. Segments do not ever need to be deleted from deep storage, however, compute nodes may load and drop segments at any time.

There is no single point of contention between historical compute nodes and the nodes have no knowledge of one another. Historical compute nodes are operationally simple; they only know how to perform the tasks they are assigned. To help other services discover historical compute nodes and the data they hold, every compute node maintains a constant Zookeeper connection. Compute nodes “announce” their online state and the segments they serve by creating ephemeral nodes under a specifically configured “announcement” path in Zookeeper. To instruct a compute node to load new segments or drop existing segments, ephemeral *znodes* are created under a special a “load queue” path associated with the compute node. Figure 2 illustrates a simple compute node and Zookeeper interaction. Each compute node announces themselves under an “announcements” path when they come online and each compute has a load queue path associated with it.

To enable queries on a segment, a compute node must first possess a local copy of the segment. Before a compute node downloads a segment from deep storage, it first checks a local disk directory (cache) to see if the segment already exists in local storage. If no cache information about the segment is present, the compute node will download metadata about the segment from Zookeeper. This metadata includes information about where the segment is located in deep storage and about how to decompress and process the segment. Once a compute node completes processing a segment, the

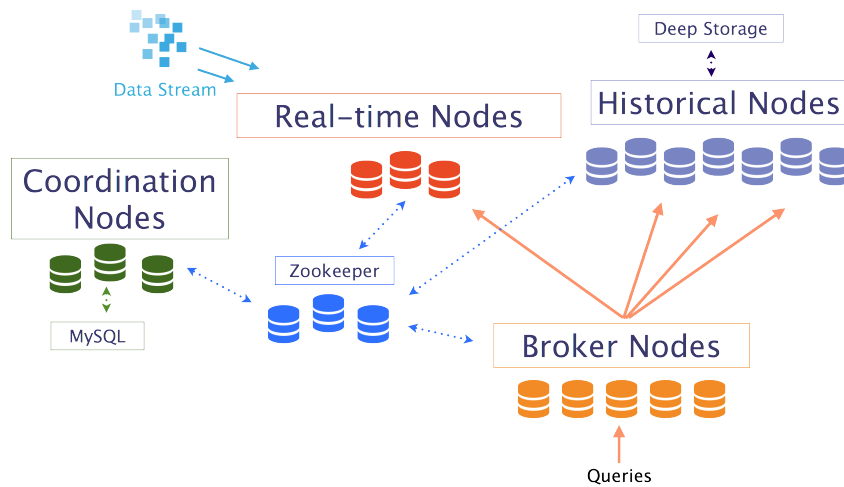


Figure 1: An overview of a Druid cluster.

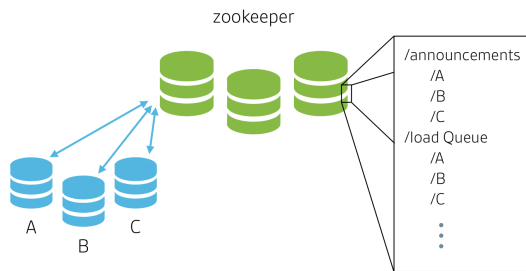


Figure 2: Compute nodes create ephemeral znodes under specifically configured Zookeeper paths.

node announces (in Zookeeper) that it is serving the segment. At this point, the segment is queryable.

3.2.1 Tiers

Compute nodes can be grouped in different tiers, where all nodes in a given tier are identically configured. Different performance and fault-tolerance parameters can be set for each tier. The purpose of tiered nodes is to enable higher or lower priority segments to be distributed according to their importance. For example, it is possible to spin up a “hot” tier of compute nodes that have a high number of cores and a large RAM capacity. The “hot” cluster can be configured to download more frequently accessed segments. A parallel “cold” cluster can also be created with much less powerful backing hardware. The “cold” cluster would only contain less frequently accessed segments.

3.3 Real-time Nodes

Real-time nodes encapsulate the functionality to ingest and query real-time data streams. Data indexed via these nodes is immediately available for querying. Real-time nodes are a consumer of data and require a corresponding producer to provide the data stream. Typically, for data durability purposes, a message bus such as Kafka [19] sits between the producer and the real-time node as shown in Figure 3.

The purpose of the message bus in Figure 3 is to act as a buffer for

incoming events. The message bus can maintain offsets indicating the position in an event stream that a real-time node has read up to and real-time nodes can update these offsets periodically. The message bus also acts as backup storage for recent events. Real-time nodes ingest data by reading events from the message bus. The time from event creation to event consumption is on the order of hundreds of milliseconds.

Real-time nodes maintain an in-memory index buffer for all incoming events. These indexes are incrementally populated as new events appear on the message bus. The indexes are also directly queryable. Real-time nodes persist their indexes to disk either periodically or after some maximum row limit is reached. After each persist, a real-time node updates the message bus with the offset of the last event of the most recently persisted index. Each persisted index is immutable. If a real-time node fails and recovers, it can simply reload any indexes that were persisted to disk and continue reading the message bus from the point the last offset was committed. Periodically committing offsets reduces the number of messages a real-time node has to rescan after a failure scenario.

Real-time nodes maintain a consolidated view of the currently updating index and of all indexes persisted to disk. This unified view allows all indexes on a node to be queried. On a periodic basis, the nodes will schedule a background task that searches for all persisted indexes of a data source. The task merges these indexes together and builds a historical segment. The nodes will upload the segment to deep storage and provide a signal for the historical compute nodes to begin serving the segment. The ingest, persist, merge, and handoff steps are fluid; there is no data loss as a real-time node converts a real-time segment to a historical one. Figure 4 illustrates the process.

Similar to compute nodes, real-time nodes announce segments in Zookeeper. Unlike historical segments, real-time segments may represent a period of time that extends into the future. For example, a real-time node may announce it is serving a segment that contains data for the current hour. Before the end of the hour, the real-time node continues to collect data for the hour. Every 10 minutes (the persist period is configurable), the node will flush and persist its in-memory index to disk. At the end of the current hour, the real-time node prepares to serve data for the next hour by creating a new in-

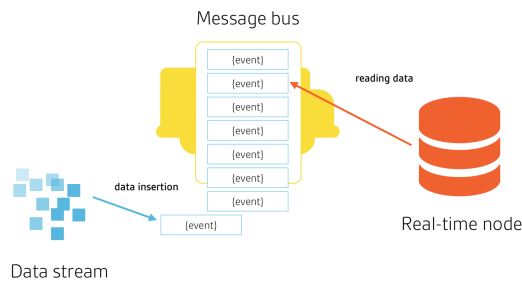


Figure 3: Real-time data ingestion.

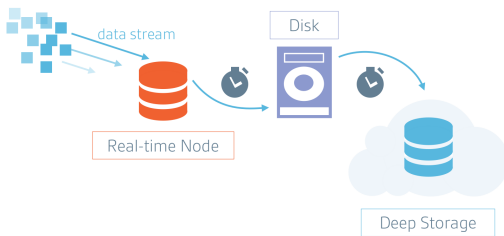


Figure 4: Real-time data durability

dex and announcing a new segment for the next hour. The node does not immediately merge and build a historical segment for the previous hour until after some window period has passed. Having a window period allows for straggling data points to come in and minimizes the risk of data loss. At the end of the window period, the real-time node will merge all persisted indexes, build a historical segment for the previous hour, and hand the segment off to historical nodes, the real-time node flushes all information about the segment and unannounces it is serving the segment.

Real-time nodes are highly scalable. If the data volume and ingestion rates for a given data source exceed the maximum capabilities of a single node, additional nodes can be added. Multiple nodes can consume events from the same stream, and every individual node only holds a portion of the total number of events. This creates natural partitions across nodes. Each node announces the real-time segment it is serving and each real-time segment has a partition number. Data from individual nodes will be merged at the broker level. To our knowledge, the largest production level real-time Druid cluster is consuming approximately 500MB/s (150,000 events/s or 2 TB/hour of raw data).

3.4 Broker Nodes

Broker nodes act as query routers to other queryable nodes such as compute and real-time nodes. Broker nodes understand the meta-data published in Zookeeper about what segments exist and on what nodes the segments are stored. Broker nodes route incoming queries such that the queries hit the right storage nodes. Broker nodes also merge partial results from storage nodes before returning a final consolidated result to the caller. Additionally, brokers provide an extra level of data durability as they maintain a cache of recent results. In the event that multiple storage nodes fail and all copies of a segment are somehow lost, it is still possible that segment results can still be returned if that information exists in the cache.

3.4.1 Timeline

To determine the correct nodes to forward queries to, brokers first build a view of the world from information in Zookeeper. Recall that Druid uses Zookeeper to maintain information about all compute and real-time nodes in a cluster and the segments those nodes are serving. For every data source in Zookeeper, the broker node builds a timeline of segments for the data source and the nodes that serve them. A timeline consists of segments and represents which segments contain data for what ranges of time. Druid may have multiple segments where the data source and interval are the same but versions differ. The timeline view will always surface segments with the most recent version identifier for a time range. If two segments intervals overlap, the segment with the more recent version always has precedence. When queries are received for a specific data source and interval, the broker node performs a lookup on the timeline associated with the query data source for the query interval and retrieves the segments that contain data for the query. The broker node maps these segments to the storage nodes that serve them and forwards the query down to the respective nodes.

3.4.2 Caching

Broker nodes employ a distributed cache with a LRU [26, 18] cache invalidation strategy. The broker cache stores per segment results. The cache can be local to each broker node or shared across multiple nodes using an external distributed cache such as memcached [15]. Recall that each time a broker node receives a query, it first maps the query to a set of segments. A subset of these segment results may already exist in the cache and the results can be directly pulled from the cache. For any segment results that do not exist in the cache, the broker node will forward the query to the compute nodes. Once the compute nodes return their results, the broker will store those results in the cache. Real-time segments are never cached and hence requests for real-time data will always be forwarded to real-time nodes. Real-time data is perpetually changing and caching the results would be unreliable.

3.5 Coordination (Master) Nodes

The Druid coordination or master nodes are primarily in charge of segment management and distribution. The Druid master is responsible for loading new segments, dropping outdated segments, managing segment replication, and balancing segment load. Druid uses a multi-version concurrency control swapping protocol for managing segments in order to maintain stable views. If any segment contains data that is wholly obsoleted by segments with newer versions, the segment is dropped from the cluster.

The Druid master runs periodically to determine the current state of the cluster. It makes decisions by comparing the expected state of the cluster with the actual state of the cluster at the time of the run. As with all Druid nodes, the Druid master maintains a connection to Zookeeper for current cluster information. The master also maintains a connection to a MySQL database that contains additional operational parameters and configurations. One of the key pieces of information located in the MySQL database is a segment table that contains a list of historical segments that should be served. This table can be updated by any service that creates historical segments. The MySQL database also contains a rule table that governs how segments are created, destroyed, and replicated in the cluster.

The master does not directly communicate with a compute node when assigning it work; instead the master creates an ephemeral znode in Zookeeper containing information about what the compute node should do. The compute node maintains a similar connection to Zookeeper to monitor for new work.

3.5.1 Rules

Rules govern how historical segments are loaded and dropped from the cluster. Rules indicate how segments should be assigned to different compute node tiers and how many replicates of a segment should exist in each tier. Rules may also indicate when segments should be dropped entirely from the cluster. Rules are usually set for a period of time. For example, a user may use rules to load the most recent one month's worth of segments into a "hot" cluster, the most recent one year's worth of segments into a "cold" cluster, and drop any segments that are older.

The master loads a set of rules from a rule table in the MySQL database. Rules may be specific to a certain data source and/or a default set of rules may be configured. The master will cycle through all available segments and match each segment with the first rule that applies to it.

3.5.2 Load Balancing

In a typical production environment, queries often hit dozens or even hundreds of data segments. Since each compute node has limited resources, historical segments must be distributed among the cluster to ensure that the cluster load is not too imbalanced. Determining optimal load distribution requires some knowledge about query patterns and speeds. Typically, queries cover recent data spanning contiguous time intervals for a single data source. On average, queries that access smaller segments are faster.

These query patterns suggest replicating recent historical segments at a higher rate, spreading out large segments that are close in time to different compute nodes, and co-locating segments from different data sources. To optimally distribute and balance segments among the cluster, we developed a cost-based optimization procedure that takes into account the segment data source, recency, and size. The exact details of the algorithm are beyond the scope of this paper and may be discussed in future literature.

4. QUERY API

Druid has its own query language and accepts queries as POST requests. All queryable Druid nodes share the same query API.

The body of the POST request is a JSON object containing key-value pairs specifying various query parameters. A typical query will contain the data source name, the granularity of the result data, time range of interest, the type of request, and the metrics to aggregate over. The result will also be a JSON object containing the aggregated metrics over the time period.

Most query types will also support a filter set. A filter set is a Boolean expression of dimension name and value pairs. Any number and combination of dimensions and values may be specified. When a filter set is provided, only the subset of the data that pertains to the filter set will be scanned. The ability to handle complex nested filter sets is what enables Druid to drill into data at any depth.

The exact query syntax depends on the query type and the information requested. A sample time series query is shown below:

```
{
  "queryType"      : "timeseries",
  "dataSource"     : "sample-data",
  "intervals"      : "2013-01-01/2013-01-08",
  "filter"         : {
```

```
    "type"         : "selector",
    "dimension"    : "poets",
    "value"        : "Ke$ha"
  },
  "granularity"    : "day",
  "aggregations"   : [
    {
      "type"        : "count",
      "name"        : "rows"
    },
    {
      "type"        : "doubleSum",
      "name"        : "blingbling",
      "fieldName"   : "revenue"
    }
  ]
}
```

The query shown above will return a count of the number of rows and the sum of the "revenue" metric (renamed as "blingbling") in the *sample-data* datasource from 2013-01-01 to 2013-01-08, filtered for only those rows where the value of the "poets" dimension is equal to "Ke\$ha". The results will be bucketed by day and will be a JSON array of the following form:

```
[
  {
    "timestamp": "2012-01-01T00:00:00.000Z",
    "result": {
      "rows": 393298,
      "blingbling": 9000.00,
    }
  },
  {
    "timestamp": "2012-01-02T00:00:00.000Z",
    "result": {
      "rows": 382932,
      "blingbling": 8594.34,
    }
  },
  ...
  {
    "timestamp": "2012-01-07T00:00:00.000Z",
    "result": {
      "rows": 1337,
      "blingbling": 1337.00
    }
  }
]
```

Other aggregations types include long sums, minimums, maximums and several more. Druid also has the ability to complex aggregations such as cardinality estimation and approximate quantile estimation. The results of aggregations can be combined in mathematical expressions to form further aggregations. Other query types include search, time boundary, and group by. The query API is highly customizable and can be extended to filter and group results based on almost any arbitrary condition. It is beyond the scope of this paper to fully describe the query API but more information can be found online². We are also in the process of extending the Druid

²<https://github.com/metamx/druid/wiki/querying>

API to understand SQL.

5. STORAGE

Druid is a column-oriented data store. When considering aggregates over a large number of events, the advantages storing data as columns rather than rows are well documented [1]. Column storage allows for more efficient CPU usage as only what is needed is actually loaded and scanned. In a row oriented data store, all columns associated with a row must be scanned as part of an aggregation. The additional scan time can introduce performance degradations as high as 250% [5].

5.1 Column Types

Druid has multiple column types to represent the various column value formats. Depending on the column type, different compression methods are used to reduce the cost of storing a column in memory and on disk. In the example given in Table 1, the publisher, advertiser, gender, and country columns only contain strings. String columns can be dictionary encoded. Dictionary encoding is a common method to compress data and has been used in other data stores such as PowerDrill [16]. In the example in Table 1, we can map each publisher to a unique integer identifier.

```
bieberfever.com    -> 0
ultratrimfast.com -> 1
```

This mapping allows us to represent the publisher column as an integer array where the array indices correspond to the rows of the original data set. For the publisher column, we can represent the unique publishers as follows:

```
[0, 0, 1, 1]
```

The resulting integer array lends itself very well to compression methods. Generic compression algorithms on top of encodings are extremely common in column-stores. Druid uses the LZ4 [21] compression algorithm.

Similar compression methods can be applied to numeric columns. For example, the clicks and revenue columns in Table 1 can also be expressed as individual arrays.

```
Clicks  -> [25, 42, 17, 170]
Revenue -> [15.70, 29.18, 17.31, 34.01]
```

In this case we compress the raw values as opposed to their dictionary representations.

5.2 Filters

To support arbitrary filter sets, Druid creates additional lookup indices for string columns. These lookup indices are compressed and Druid operates over the indices in their compressed form. Filters can be expressed as Boolean equations of multiple lookup indices. Boolean operations of indices in their compressed form is both performance and space efficient.

Let us consider the publisher column in Table 1. For each unique publisher in Table 1, we can form some representation indicating

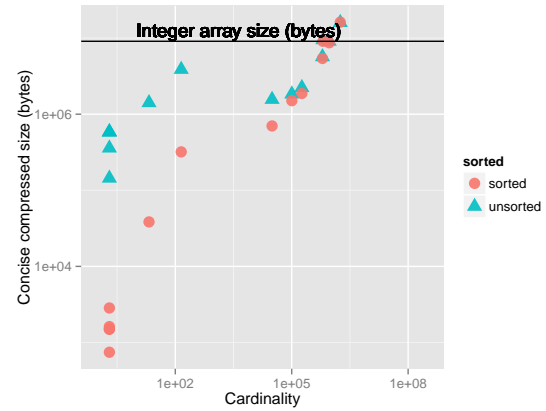


Figure 5: Integer array size versus Concise set size.

which table rows a particular publisher is seen. We can store this information in a binary array where the array indices represent our rows. If a particular publisher is seen in a certain row, that array index is marked as 1. For example:

```
bieberfever.com    -> rows [0, 1] -> [1] [1] [0] [0]
ultratrimfast.com -> rows [2, 3] -> [0] [0] [1] [1]
```

bieberfever.com is seen in rows 0 and 1. This mapping of column values to row indices forms an inverted index [32]. To know which rows contain bieberfever.com or ultratrimfast.com, we can OR together the two arrays.

```
[0] [1] [0] [1] OR [1] [0] [1] [0] = [1] [1] [1] [1]
```

This approach of performing Boolean operations on large bitmap sets is commonly used in search engines. Bitmap compression algorithms are a well-defined area of research and often utilize run-length encoding. Popular algorithms include Byte-aligned Bitmap Code [3], Word-Aligned Hybrid (WAH) code [36], and Partitioned Word-Aligned Hybrid (PWAH) compression [34]. Druid opted to use the Concise algorithm [9] as it can outperform WAH by reducing the size of the compressed bitmaps by up to 50%. Figure 5 illustrates the number of bytes using Concise compression versus using an integer array. The results were generated on a cc2.8xlarge system with a single thread, 2G heap, 512m young gen, and a forced GC between each run. The data set is a single day's worth of data collected from the Twitter garden hose [33] data stream. The data set contains 2,272,295 rows and 12 dimensions of varying cardinality. As an additional comparison, we also resorted the data set rows to maximize compression.

In the unsorted case, the total Concise size was 53,451,144 bytes and the total integer array size was 127,248,520 bytes. Overall, Concise compressed sets are about 42% smaller than integer arrays. In the sorted case, the total Concise compressed size was 43,832,884 bytes and the total integer array size was 127,248,520 bytes. What is interesting to note is that after sorting, global compression only increased minimally. The total Concise set size to total integer array size is 34%. It is also interesting to note that as the cardinality of a dimension approaches the total number of rows

in a data set, integer arrays require less space than Concise sets and become a better alternative.

5.3 Storage Engine

Druid's persistence components allows for different storage engines to be plugged in, similar to Dynamo [10]. These storage engines may store data in in-memory structures such as the JVM heap or in memory-mapped structures. The ability to swap storage engines allows for Druid to be configured depending on a particular application's specifications. An in-memory storage engine may be operationally more expensive than a memory-mapped storage engine but could be a better alternative if performance is critical. At Metamarkets, we commonly use a memory-mapped storage engine.

6. ROBUSTNESS

To achieve high system availability and data durability, Druid employs several fault recovery techniques. Druid has no single point of failure.

6.1 Replication

Druid replicates historical segments on multiple hosts. The number of replicates in each tier of the historical compute cluster is fully configurable. Setups that require high levels of fault tolerance can be configured to have a high number of replicates. Replicates are assigned to compute nodes by coordination nodes using the same load distribution algorithm discussed in Section 3.4.2. Broker nodes forward queries to the first node they find that contain a segment required for the query.

Real-time segments follow a different replication model as real-time segments are mutable. Multiple real-time nodes can read from the same message bus and event stream if each node maintains a unique offset and consumer id, hence creating multiple copies of a real-time segment. This is conceptually different than multiple nodes reading from the same event stream and sharing the same offset and consumer id, doing so would create multiple segment partitions. If a real-time node fails and recovers, it can reload any indexes that were persisted to disk and read from the message bus from the point it last committed an offset.

6.2 Failure Detection

If a compute node completely fails and becomes unavailable, the ephemeral Zookeeper znodes it created are deleted. The master node will notice that certain segments are insufficiently replicated or missing altogether. Additional replicates will be created and re-distributed throughout the cluster.

We are moving towards building out infrastructure to support programmatic creation of real-time nodes. In the near future, the master node or a similar service will notice if real-time segments are insufficiently replicated and automatically create additional real-time nodes as redundant backups.

Coordination and broker nodes always have redundant backup nodes in the event the primary fails. The backup nodes are usually idle until Zookeeper changes alert them to assume the responsibilities of their primary counterparts.

6.3 Zookeeper Failure

Zookeeper itself is a highly available system, but in the event of network disconnections with Zookeeper, Druid queries will continue

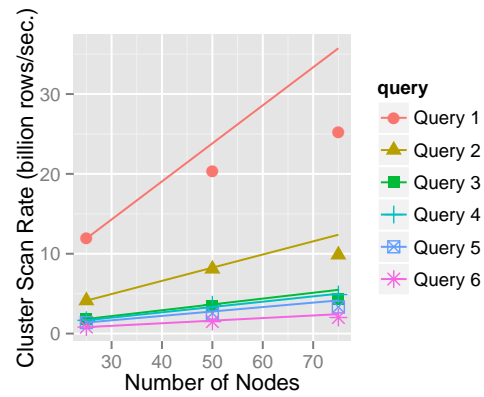


Figure 6: Druid cluster scan rate with lines indicating linear scaling from 25 nodes.

to work. The Druid cluster utilizes Zookeeper for cluster wide coordination, and without Zookeeper, the cluster maintains a steady state. Each node type interacts with Zookeeper in different ways. Compute nodes rely on Zookeeper to publish the segments they are serving and whether or not they need to load or drop segments. If Zookeeper goes down, compute nodes will no longer be able to load or drop data, however, they will still respond to queries as queries are over HTTP. Broker nodes require Zookeeper to know which compute nodes serve which segments. It caches this information locally, and in the event of Zookeeper failure, Broker nodes will use the latest view of compute nodes and segments it had prior to the outage. Broker nodes will continue to forward queries down to compute and realtime nodes. Master nodes use Zookeeper to determine how to assign, balance, and drop segments. In the event of a Zookeeper outage, these features will stop working.

6.4 MySQL Failure

The Druid master node has a dependency on MySQL. The master node loads information about how to replicate, assign, and drop segments from a segment table. If MySQL goes down, the master node effectively stops working. However, this does not stop the rest of the cluster from functioning. Queries will continue to be served, however, during the time that MySQL is not available, no new segments will be loaded in the cluster and no old segments will be dropped from the cluster. The design principle for MySQL failures is similar to Zookeeper failures; if a piece of the cluster that is responsible for coordination fails, the status maintains a status quo state. No new actions are taken, however, data is still available.

6.5 Adding and Removing Nodes

Starting and removing Druid nodes is done by starting and stopping Java processes. There is minimal operational overhead with adding nodes in batches. Scaling down the cluster is usually done one node at a time with some time lapse between shutdowns. This allows the master to have ample time to redistribute load and create additional replicates. Shutting down nodes in batches is not recommended as it may destroy all copies of a segment, which would lead to data loss.

7. PERFORMANCE BENCHMARKS

To benchmark Druid performance, we created a large test cluster with 6TB of uncompressed data, representing tens of billions of fact rows. The data set contained more than a dozen dimensions, with

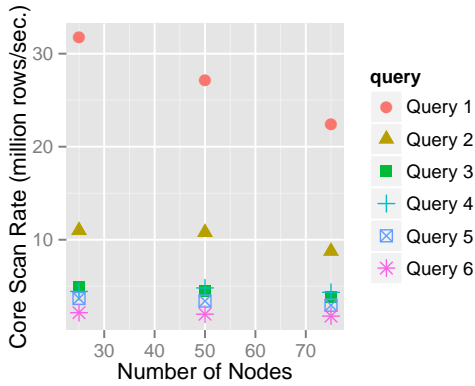


Figure 7: Druid core scan rate.

cardinalities ranging from the double digits to tens of millions. We computed four metrics for each row (counts, sums, and averages). The data was sharded first on timestamp then on dimension values, creating thousands of shards roughly 8 million fact rows apiece.

The cluster used in the benchmark consisted of 100 historical compute nodes, each with 16 cores, 60GB of RAM, 10 GigE Ethernet, and 1TB of disk space. Collectively, the cluster comprised of 1600 cores, 6TB or RAM, sufficiently fast Ethernet and more than enough disk space.

SQL statements are included in Table 2 to describe the purpose of each of the queries. Please note:

- The timestamp range of the queries encompassed all data.
- Each machine was a 16-core machine with 60GB RAM and 1TB of local disk. The machine was configured to only use 15 threads for processing queries.
- A memory-mapped storage engine was used (the machine was configured to memory map the data instead of loading it into the Java heap.)

Figure 6 shows the cluster scan rate and Figure 7 shows the core scan rate. In Figure 6 we also include projected linear scaling based on the results of the 25 core cluster. In particular, we observe diminishing marginal returns to performance in the size of the cluster. Under linear scaling, SQL query 1 would have achieved a speed of 37 billion rows per second on our 75 node cluster. In fact, the speed was 26 billion rows per second. However, queries 2-6 maintain a near-linear speedup up to 50 nodes: the core scan rates in Figure 7 remain nearly constant. The increase in speed of a parallel computing system is often limited by the time needed for the sequential operations of the system, in accordance with Amdahl’s law [2].

The first query listed in Table 2 is a simple count, achieving scan rates of 33M rows/second/core. We believe the 75 node cluster was actually overprovisioned for the test dataset, explaining the modest improvement over the 50 node cluster. Druid’s concurrency model is based on shards: one thread will scan one shard. If the number of segments on a compute node modulo the number of cores is small (e.g. 17 segments and 15 cores), then many of the cores will be idle during the last round of the computation.

When we include more aggregations we see performance degrade. This is because of the column-oriented storage format Druid employs. For the count (*) queries, Druid only has to check the timestamp column to satisfy the “where” clause. As we add metrics, it has to also load those metric values and scan over them, increasing the amount of memory scanned.

8. RELATED WORK

Cattell [6] maintains a great summary about existing Scalable SQL and NoSQL data stores. Druid feature-wise sits somewhere between Google’s Dremel [24] and PowerDrill [16]. Druid has most of the features implemented in Dremel (Dremel handles arbitrary nested data structures while Druid only allows for a single level of array-based nesting) and many of the interesting compression algorithms mentioned in PowerDrill.

Although Druid builds on many of the same principles as other distributed columnar data stores [14], most existing data stores are designed to be key-value stores [20], or document/extensible record stores [31]. Such data stores are great solutions for traditional data warehouse needs and general back-office/reporting usage. Typically, analysts will query these data stores and build reports from the results. In-memory databases such as SAP’s HANA [13] and VoltDB [35] are examples of other data stores that are highly suited for traditional data warehousing needs. Druid is a front-office system designed such that user-facing dashboards can be built on top of it. Similar to [27], Druid has analytical features built in. The main features Druid offers over traditional data warehousing solutions are real-time data ingestion, interactive queries and interactive query latencies. In terms of real-time ingestion and processing of data, Trident/Storm [23] and Streaming Spark [37] are other popular real-time computation systems, although they lack the data storage capabilities of Druid. Spark/Shark [12] are also doing similar work in the area of fast data analysis on large scale data sets. Cloudera Impala [8] is another system focused on optimizing querying performance, but more so in Hadoop environments.

Druid leverages a unique combination of algorithms in its architecture. Although we believe no other data store has the same set of functionality as Druid, some of Druid’s optimization techniques such as using inverted indices to perform fast filters are also used in other data stores [22].

9. CONCLUSIONS

In this paper, we presented Druid, a distributed, column-oriented, real-time analytical data store. Druid is a highly customizable solution that is optimized for fast query latencies. Druid ingests data in real-time and is fault-tolerant. We discussed the performance of Druid on billion row data sets. We summarized key Druid architecture aspects such as the storage format, query language and general execution. In the future, we plan to cover more in depth the different algorithms we’ve developed for Druid and how other systems may plug into Druid to achieve powerful use cases.

10. ACKNOWLEDGEMENTS

We want to thank Steve Harris for his feedback on improving this paper, Adam Smith for giving us the chance to write this paper, and Katherine Chu for helping to create all the images in this paper. Druid could not have been built without the help of many great engineers at Metamarkets and in the community. We want to thank Danny Yuan, Jae Hyeon Bae, Paul Baclace, Dave Nielsen, and Dhruv Parthasarathy for their contributions to Druid.

Table 2: Druid Queries

| Query # | Query |
|---------|---|
| 1 | SELECT count(*) FROM _table_ WHERE timestamp \geq ? AND timestamp < ? |
| 2 | SELECT count(*), sum(metric1) FROM _table_ WHERE timestamp \geq ? AND timestamp < ? |
| 3 | SELECT count(*), sum(metric1), sum(metric2), sum(metric3), sum(metric4) FROM _table_ WHERE timestamp \geq ? AND timestamp < ? |
| 4 | SELECT high_card_dimension, count(*) AS cnt FROM _table_ WHERE timestamp \geq ? AND timestamp < ? GROUP BY high_card_dimension ORDER BY cnt limit 100 |
| 5 | SELECT high_card_dimension, count(*) AS cnt, sum(metric1) FROM _table_ WHERE timestamp \geq ? AND timestamp < ? GROUP BY high_card_dimension ORDER BY cnt limit 100 |
| 6 | SELECT high_card_dimension, count(*) AS cnt, sum(metric1), sum(metric2), sum(metric3), sum(metric4) FROM _table_ WHERE timestamp \geq ? AND timestamp < ? GROUP BY high_card_dimension ORDER BY cnt limit 100 |

11. ADDITIONAL AUTHORS

12. REFERENCES

- [1] D. J. Abadi, S. R. Madden, and N. Hachem. Column-stores vs. row-stores: How different are they really? In *Proceedings of the 2008 ACM SIGMOD international conference on Management of data*, pages 967–980. ACM, 2008.
- [2] G. M. Amdahl. Validity of the single processor approach to achieving large scale computing capabilities. In *Proceedings of the April 18-20, 1967, spring joint computer conference*, pages 483–485. ACM, 1967.
- [3] G. Antoshenkov. Byte-aligned bitmap compression. In *Data Compression Conference, 1995. DCC’95. Proceedings*, page 476. IEEE, 1995.
- [4] L. A. Barroso and U. Hözl. The datacenter as a computer: An introduction to the design of warehouse-scale machines. *Synthesis Lectures on Computer Architecture*, 4(1):1–108, 2009.
- [5] C. Bear, A. Lamb, and N. Tran. The vertica database: Sql rdbms for managing big data. In *Proceedings of the 2012 workshop on Management of big data systems*, pages 37–38. ACM, 2012.
- [6] R. Cattell. Scalable sql and nosql data stores. *ACM SIGMOD Record*, 39(4):12–27, 2011.
- [7] S. Chaudhuri and U. Dayal. An overview of data warehousing and olap technology. *ACM Sigmod record*, 26(1):65–74, 1997.
- [8] Cloudera impala. <http://blog.cloudera.com/blog>, March 2013.
- [9] A. Colantonio and R. Di Pietro. Concise: Compressed ‘n’composable integer set. *Information Processing Letters*, 110(16):644–650, 2010.
- [10] G. DeCandia, D. Hastorun, M. Jampani, G. Kakulapati, A. Lakshman, A. Pilchin, S. Sivasubramanian, P. Voshall, and W. Vogels. Dynamo: amazon’s highly available key-value store. In *ACM SIGOPS Operating Systems Review*, volume 41, pages 205–220. ACM, 2007.
- [11] D. DeWitt and J. Gray. Parallel database systems: the future of high performance database systems. *Communications of the ACM*, 35(6):85–98, 1992.
- [12] C. Engle, A. Lupher, R. Xin, M. Zaharia, M. J. Franklin, S. Shenker, and I. Stoica. Shark: fast data analysis using coarse-grained distributed memory. In *Proceedings of the 2012 international conference on Management of Data*, pages 689–692. ACM, 2012.
- [13] F. Färber, S. K. Cha, J. Primsch, C. Bornhövd, S. Sigg, and W. Lehner. Sap hana database: data management for modern business applications. *ACM Sigmod Record*, 40(4):45–51, 2012.
- [14] B. Fink. Distributed computation on dynamo-style distributed storage: riak pipe. In *Proceedings of the eleventh ACM SIGPLAN workshop on Erlang workshop*, pages 43–50. ACM, 2012.
- [15] B. Fitzpatrick. Distributed caching with memcached. *Linux journal*, (124):72–74, 2004.
- [16] A. Hall, O. Bachmann, R. Büsow, S. Găncăanu, and M. Nunkesser. Processing a trillion cells per mouse click. *Proceedings of the VLDB Endowment*, 5(11):1436–1446, 2012.
- [17] P. Hunt, M. Konar, F. P. Junqueira, and B. Reed. Zookeeper: Wait-free coordination for internet-scale systems. In *USENIX ATC*, volume 10, 2010.
- [18] C. S. Kim. Lrfu: A spectrum of policies that subsumes the least recently used and least frequently used policies. *IEEE Transactions on Computers*, 50(12), 2001.
- [19] J. Kreps, N. Narkhede, and J. Rao. Kafka: A distributed messaging system for log processing. In *Proceedings of 6th International Workshop on Networking Meets Databases (NetDB)*, Athens, Greece, 2011.
- [20] R. Lerner. At the forge: Redis. *Linux Journal*, 2010(197):3, 2010.
- [21] Liblzf. <http://freecode.com/projects/liblzf>, March 2013.
- [22] R. MacNicol and B. French. Sybase iq multiplex-designed for analytics. In *Proceedings of the Thirtieth international conference on Very large data bases-Volume 30*, pages 1227–1230. VLDB Endowment, 2004.
- [23] N. Marz. Storm: Distributed and fault-tolerant realtime computation. <http://storm-project.net/>, February 2013.
- [24] S. Melnik, A. Gubarev, J. J. Long, G. Romer, S. Shivakumar, M. Tolton, and T. Vassilakis. Dremel: interactive analysis of web-scale datasets. *Proceedings of the VLDB Endowment*, 3(1-2):330–339, 2010.
- [25] D. Miner. Unified analytics platform for big data. In *Proceedings of the WICSA/ECSA 2012 Companion Volume*, pages 176–176. ACM, 2012.
- [26] E. J. O’neil, P. E. O’neil, and G. Weikum. The lru-k page replacement algorithm for database disk buffering. In *ACM SIGMOD Record*, volume 22, pages 297–306. ACM, 1993.
- [27] Paracel analytic database. <http://www.paracel.com/resources/Datasheets/ParAccel-Core-Analytic-Database.pdf>, March 2013.

- [28] S. Sarawagi, R. Agrawal, and N. Megiddo. *Discovery-driven exploration of OLAP data cubes*. Springer, 1998.
- [29] K. Shvachko, H. Kuang, S. Radia, and R. Chansler. The hadoop distributed file system. In *Mass Storage Systems and Technologies (MSST), 2010 IEEE 26th Symposium on*, pages 1–10. IEEE, 2010.
- [30] M. Singh and B. Leonhardi. Introduction to the ibm netezza warehouse appliance. In *Proceedings of the 2011 Conference of the Center for Advanced Studies on Collaborative Research*, pages 385–386. IBM Corp., 2011.
- [31] M. Stonebraker, D. J. Abadi, A. Batkin, X. Chen, M. Cherniack, M. Ferreira, E. Lau, A. Lin, S. Madden, E. O’Neil, et al. C-store: a column-oriented dbms. In *Proceedings of the 31st international conference on Very large data bases*, pages 553–564. VLDB Endowment, 2005.
- [32] A. Tomasic and H. Garcia-Molina. Performance of inverted indices in shared-nothing distributed text document information retrieval systems. In *Parallel and Distributed Information Systems, 1993., Proceedings of the Second International Conference on*, pages 8–17. IEEE, 1993.
- [33] Twitter public streams. <https://dev.twitter.com/docs/streaming-apis/streams/public>, March 2013.
- [34] S. J. van Schaik and O. de Moor. A memory efficient reachability data structure through bit vector compression. In *Proceedings of the 2011 international conference on Management of data*, pages 913–924. ACM, 2011.
- [35] L. VoltDB. Voltdb technical overview. <https://voltdb.com/>, 2010.
- [36] K. Wu, E. J. Otoo, and A. Shoshani. Optimizing bitmap indices with efficient compression. *ACM Transactions on Database Systems (TODS)*, 31(1):1–38, 2006.
- [37] M. Zaharia, T. Das, H. Li, S. Shenker, and I. Stoica. Discretized streams: an efficient and fault-tolerant model for stream processing on large clusters. In *Proceedings of the 4th USENIX conference on Hot Topics in Cloud Computing*, pages 10–10. USENIX Association, 2012.