


An Efficient NoSQL-Based Storage Schema for Large-Scale Time Series Data

Ruizhe Ma, University of Massachusetts, Lowell, USA

Weiwei Zhou, Nanjing University of Aeronautics and Astronautics, China

Zongmin Ma, Nanjing University of Aeronautics and Astronautics, China*

 <https://orcid.org/0000-0001-7780-6473>

ABSTRACT

In IoT (internet of things), most data from the connected devices change with time and have sampling intervals, which are called time-series data. It is challenging to design a time series storage model that can write massive time-series data in a short time and can query and analyze the persistent time-series data for a long time. This paper constructs the RHTSDB (Redis-HBase Time Series Database) storage model based on Redis and HBase. RHTSDB uses the memory database Redis (Remote Dictionary Server) to cache massive time-series data, providing efficient data storage and query functions. HBase is used in RHTSDB for long-term storage of time-series data to realize their persistence. The paper designs a cold and hot separation mechanism for time-series data, where the infrequently accessed cold data are stored in HBase, and the frequently accessed and latest data are stored in Redis. Experiments verify that RHTSDB has apparent advantages over Apache IoTDB and HBase in data intake and query efficiency.

KEYWORDS

HBase, Query, Redis, Storage, Time-Series Data

INTRODUCTION

With the development of the Internet of Things (IoT) (Eom & Lee, 2017), the amount of time-series data has shown explosive growth. Time-series data refers to a sequence of data points collected at fixed time intervals (Lee & Chung, 2014). Each data point is associated with a timestamp that indicates the generation time of the corresponding data. Typically, the data collected by a sensor in a particular period can be expressed as a time series $[(t_1, v_1), (t_2, v_2), \dots, (t_n, v_n)]$, where v_i refers to the value collected at t_i time (Di Martino *et al.*, 2019). Of course, complete time-series data can include the collection time and collection value as well as the source description information of the current collection value. For example, we need to include some measurement data information, such as the names of collection subject and collection index. Comprehensive use cases in the real world have generated a large amount of measurement data from millions or billions of different sources. Slack

DOI: 10.4018/JDM.339915

*Corresponding Author

This article published as an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>) which permits unrestricted use, distribution, and production in any medium, provided the author of the original work and original publication source are properly credited.

collects measurement data from 4 billion unique sources at 12 million samples per second daily, for example, generating up to 12 TB of compressed data daily. It is essential to manage and process a large amount of time-series data efficiently. Unfortunately, many off-the-shelf systems cannot scale to support these workloads, which leads to the random Patchwork and vulnerability of customized solutions (Solleza, Crotty, Karumuri, Tatbul & Zdonik, 2022). For this reason, diverse time series databases are proposed to ensure efficient ingestion performance and save storage space as much as possible. Given that time-series data in applications are generally massive and redundant data containing source description information in time-series data are enormous, efficient storage and query of massive time-series data is challenging.

We identify two major categories of time series databases: which are respectively called *native time series databases* and *common time series databases* in this paper. The native time series databases are the storage systems that are developed especially for time-series data management according to their structural and usage characteristics, such as InfluxDB¹, FluteDB (Li *et al.*, 2018), and Apache IoTDB (Wang *et al.*, 2020 & 2023). This category of time series databases can efficiently reduce the overhead of storage space and the query delay. However, for time-series data management and processing, many other functions and operations are essential in time series databases, such as flexible aggregation, data retention, multidimensional range query, among others. While the native time series databases cannot provide full support to time-series data analysis well, mature database systems are good at dealing with relationships between data and support many unnecessary operations and guarantees for time series, increasing inefficiency and unnecessary complexity (Shafer, Sambasivan, Rowe, & Ganger, 2013). The common time series databases are the storage systems that directly apply the common databases for storing and processing time-series data. Depending on what types of databases are applied, we further identify two categories of common time series databases. The first one uses relational databases as the back end of common time series databases (e.g., (Rhea *et al.*, 2017)). In recent years, NoSQL (Not only SQL) databases have attracted increasing attention from both academia and industry (Hu & Dessloch, 2015), which offer flexible data representation models and horizontal hardware scalability so that Big Data can be processed in real time (Bajaj & Bick, 2020). The second category of common time series databases uses NoSQL databases for processing time-series data (Di Martino *et al.*, 2019).

NoSQL databases contain four major types of database models: *key-value stores*, *column-family stores*, *document stores*, and *graph stores* (Grolinger *et al.*, 2013; Van Erven *et al.*, 2019). Different types of NoSQL databases, say Redis² (a key-value store), HBase³ (a column-family store), Cassandra⁴ (a column-family store), MongoDB⁵ (a document store), Couchbase⁶ (a document store), OrientDB⁷ (a graph store), have very different performances (Matallah, Belalem & Bouamrane, 2020). Among these NoSQL databases, for example, Redis (remote dictionary server) is a high-performance memory-based NoSQL database, which has excellent data writing performance and supports data persistent storage and replication of master-slave nodes (Zhou, Lu, Zhang & Qi, 2020); HBase, an open-source, distributed and versioned NoSQL database, uses Hadoop distributed file system (HDFS) to provide distributed file storage services, which has the characteristics of high availability, robust scalability and the ability to store massive data. With one of NoSQL databases as back end, some storage models for massive time-series data have been developed, for example, ModelarDB+ (Jensen, Pedersen & Thomsen, 2021) based on Cassandra, NagareDB (Calatrava, Fontal, Cucchiatti & Diví-Cuesta, 2021) built on top of MongoDB, OpenTSDB⁸ based on HBase, and KairosDB⁹ based on Cassandra.

Selecting a specific NoSQL database for massive time-series data management is more flexible and cost-saving, where the performance of the storage model is determined by the performance of the used database models (Rinaldi *et al.*, 2019). It has been demonstrated that different NoSQL databases have very different performances (Matallah, Belalem & Bouamrane, 2020). For ingestion of large-scale time-series data into the target time series database, for example, the storage consumption of HBase slowly increases along with a significant increase in data size (e.g., a 50% increase in data size requires about a 20% increase in storage consumption). However, the storage consumption of

19 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage: www.igi-global.com/article/an-efficient-nosql-based-storage-schema-for-large-scale-time-series-data/339915

Related Content

Issues and Methods for Access, Storage, and Analysis of Data From Online Social Communities

Christopher John Quinn, Matthew James Quinn, Alan Olinsky and John Thomas Quinn (2018). *Handbook of Research on Big Data Storage and Visualization Techniques* (pp. 402-432).

www.irma-international.org/chapter/issues-and-methods-for-access-storage-and-analysis-of-data-from-online-social-communities/198772

Dynamic Multidimensional Data Cubes

Mirek Riedewald, Divyakant Agrawal and Amr El Abbadi (2003). *Multidimensional Databases: Problems and Solutions* (pp. 200-221).

www.irma-international.org/chapter/dynamic-multidimensional-data-cubes/26969

Distributed Fragmented Database Set Queries

Clive C. Sanford, Robert E. Markland, Patrick R. Philipoom and Kenneth Darby-Dowman (1990). *Journal of Database Administration* (pp. 30-47).

www.irma-international.org/article/distributed-fragmented-database-set-queries/51080

The Critical Role of Information Processing in Creating an Effective Knowledge Organization

William R. King (2009). *Database Technologies: Concepts, Methodologies, Tools, and Applications* (pp. 1638-1651).

www.irma-international.org/chapter/critical-role-information-processing-creating/7996

Bringing Credibility Through Portals: Examining the Multiple Information Sources on Information Quality and Credibility

Yi-Chen Lee, Bo-Chun Shen and Choon Ling Sia (2022). *Journal of Database Management* (pp. 1-26).

www.irma-international.org/article/bringing-credibility-through-portals/313968