

Overview of Large-scale RDF Data Storage Strategies

Zekai Cai

School of Fuzhou University, Fujian, China

Abstract

Information processing plays an increasingly important role in e-commerce, data mining, information extraction, power system, medical system, space-time and multimedia information technology and network application. On the other hand, in order to improve the intelligence of the Internet, the semantic web is developing continuously. RDF is the foundation of semantic Web application and the standard of intelligent information service and semantic interoperability on the Web. Making full use of RDF can effectively deal with Web information, and the storage problem of RDF needs to be solved first. Distributed cloud computing technology provides a new and more efficient solution for massive RDF storage and query, and the research on RDF data storage based on Hadoop platform has become the research focus. Studies have shown that different data sets and queries often need different storage schemes, and the existing storage methods have their own advantages and disadvantages in different application scenarios. Setting the Rowkey of Angel Liu HBase effectively not only avoids node accumulation, but also ensures the integrity of data by using BKDRHash algorithm.

Keywords

Large Scale; RDF Storage; Model Design.

1. Introduction

RDF is a resource description framework put forward by W3C, which makes it possible for computers to simulate the brain to understand semantic information by semantic description of network resources to a certain extent [1]. With the development of LOD and other projects, a large number of RDF data are released, and hundreds of millions of RDF data are contained in the Internet. Traditional relational database is unable to deal with massive data, and traditional relational database can also deal with temporal information. However, obviously, this kind of processing needs to increase a lot of storage space. Besides, traditional relational database has some shortcomings in the operation of state data and ensuring the consistency of temporal data, and distributed cloud computing technology has unique advantages in dealing with massive data. Hadoop, with its implementation of distributed computing and open source, has attracted many researchers and research institutions to study massive RDF data storage [2]. Tournament of Champions At present, RDF, as the description language of semantic web data, is widely used in many fields, such as general search engines, vertical search engines, text assistant tools, personal information manager and semantic browsing tools [3]. Although some studies have put forward the requirements of temporal features in Web documents, how to introduce temporal information into RDF and how to store temporal RDF have not received enough attention. Therefore, it is very important and necessary to study adding temporal attributes to traditional RDF to form a temporal RDF model, as well as research and storage methods [4]. Both ternary table storage and horizontal table storage adopt a fixed storage strategy, which is adjusted according to different situations. For attribute table storage, the attribute selection method should be specified as the basis for table building. A good attribute selection method should be able to meet the needs of different applications and achieve the effect of flexibly customizing the storage strategy [5].

Traditional high-performance database system is mainly based on HDD, and its performance bottleneck mainly focuses on I/O rate. Solid-state disk (SSD) is a new type of storage device, which uses integrated circuit as memory to store data persistently [6]. Compared with HDD, SSD is used for outstanding reading and writing performance, but SSD also has some disadvantages, such as low capacity and high price. Considering the respective characteristics of SSD and HDD, building a hybrid storage device by comprehensively utilizing their advantages can greatly improve the performance of Web data management. Angel Liu base, an open source distributed database, is used for storage, making full use of the characteristics of RDF data, and its predicate and hash Value are stored as Rowkey, while the object and subject are stored in the value field. The data are loaded into the database by Map Reduce and Bulk Load method of Angel Liu base [7].

2. RDF Storage Strategy

2.1 RDF Data Model

The basic concepts of RDF model include resources, attributes, declarations and graphs [8]. A resource can be defined as an object, which can be identified by a URL; Attribute is a special kind of resource, which is used to describe the relationship between resources. The data model takes statements as the core, and statements assert the attributes of resources. A statement is composed of a resource, an attribute and an attribute value, and it is a triple of "the Kramp-Karrenbauer value of the Kramp-Karrenbauer attribute of entities". However, we often use subjects to represent entities in the triple, objects to represent their values, and attributes are represented by predicates. Graph mainly means that we can use graphical way to write the same statement, using labeled nodes to connect lines through labeled directed edges, and each statement can be divided into the main body [9]. Predicate and object, so it is a feasible strategy to use ternary table to express the ternary relationship of statements. In data, the subject and predicate are required to be expressed in form, and the object can be or face value. Usually, "S" is used for the subject, "P" for the predicate, and "O" for the object, as shown in Table 1.

Table 1. ternary relationship table

Field name	Data type	Explanation
ID	Integer	Primary Key,statement number
Subject	Integer	statement number
Predicate	Integer	Predicate number
ObjResource	Integer	The number when the object is a resource.
ObjLiteral	Integer	The number when the object is text.
ObjFlag	Char(1)	Identify whether the object is a resource or a text.

From the definition of data, we can see that data sources have the following basic characteristics: flexibility, graph structure and diverse access patterns. Based on the above characteristics, the data is quite different from the traditional relational model data, which brings difficulties to the scalable and efficient storage of data. Similar to the query language used by traditional relational databases, data can be queried by using. All-Yes-Yes is a recommended standard. By means of graph matching, variables can be bound, and the data meeting the specified pattern can be found from a given data set. Triple storage is also called vertical table storage. The basic idea is to store triple directly in a relational table composed of subject, predicate and object. In RDF context, resources refer to anything with a universal resource identifier, which is usually described by some attributes with corresponding values. The attribute value here can be text or other resources. Therefore, RDF model is a set of triple forms of subject, predicate and object, in which the subject is the described resource, the predicate is

the attribute of the resource and the object is the value of the attribute. An RDF triple (S, Ma Yili, O) represents an RDF statement. RDF grammar constructs a complete grammar system for computer automatic processing, which includes RDF abstract grammar 8 and RDF/XML grammar Fast & Furious 9, which can encode RDF into XML to integrate various metadata. The RDF abstract syntax data model is a collection of triples. Formally, Bidai Syulan is a set of URI references, bilibili is a set of blank nodes, and Donald L. Miller is a set of words. Then, an RDF triple (s, p, o) is an element of set $(U \cup B) \times U \times (U \cup B \cup L)$, namely $(s, p, o) \in (U \cup B) \times U \times (U \cup B \cup L)$. In addition, a set of RDF triples can be represented as a directed labeled graph, which is called RDF graph. RDF is represented as $G = (N, E, \Sigma, L)$, where Kang Seung Yoon is a finite set of vertices, $E \in N \times N$ is a finite set of directed edges, Σ is a finite set of labels, and $L: N \cup E \rightarrow \Sigma$ is a function that assigns labels to vertices and edges respectively. A set of RDF triples consists of RDF graphs, where the subject and object are the vertices of the graph, and the predicate constitutes the labeled directed edges of the graph.

In order to further improve the query performance, a more aggressive indexing strategy is proposed in Hexastore. It contains some ideas similar to horizontal table storage, but it increases the division between subject and object. And you can use this function to query data. By integrating the data query module into the functions in the relational database, we can directly use the familiar language to query, better combine with the traditional relational database system, and comprehensively consider the requirements of data integrity, redundancy, query efficiency, etc., and design the relational schema for storing RDF statements.

2.2 Design of RDF Storage

Three key problems of RDF data storage: the selection of storage container, the selection and division of database, and the design of index strategy. The existing RDF data storage mainly uses relational database as its underlying storage support, and the existing storage schemes mainly include triple table storage, vertical storage, horizontal storage and schema generation storage schemes. When using relational database to store RDF data, because of the flexibility of semantic Web, the traditional relational database modeling belongs to pre-definition. When new users or new data sources join, the original schema and new schema can't be effectively integrated, which makes the pre-defined schema unstable or infeasible. With the rapid growth of RDF data, the storage mode based on relational database has been difficult to efficiently meet the storage requirements of a large amount of RDF data. At present, more and more scholars are more interested in applying distributed systems to the management of massive RDF data. Several schemes of using distributed systems to store massive RDF data are summarized in Table 2.

Table 2. Summary of RDF data storage scheme

Storage strategy	Advantage
HDFS/Map Reduce	The predicates and objects in RDF data are divided into types and stored in HDFS file system, and MapReduce is used to process queries in parallel.
HDFS/Map Redyce	According to the category of the subject and stored in HDFS file system, a MapRedycej job handles a triple pattern query.
HBase/HBase API	Create three tables, with the subject predicate object as Row key, and use HBase API to realize SOA RQL query.
HBase/Hadoop	Create six tables to store data, and propose an elastic MapReduce multiplex connection to query.

RDFSchema is based on RDF and is used to describe RDF vocabulary. It provides a set of modeling primitives, which can organize RDF vocabulary into hierarchical structure. RDFSchema can be used to declare classes, attributes, and relationships between classes and attributes.

Because RDF data query mode is almost SPARQLBGP basic graph mode query, while MapReduce algorithm is suitable for batch processing but not native algorithm in graph matching, and the time of starting MapReduce will have an impact on the response of the whole query, so we can design an algorithm with higher query efficiency by combining the storage of Angel Liu HBase. Angel Liu data is a multidimensional and ordered mapping table, and the index of the table is Rowkey or ROWKEY Roger Waters: Us+Them COLUMN, while SPARQL query language gives priority to the triple pattern given by predicates. Given that Rowkey is sorted by dictionary, BKDRHash algorithm is used here to calculate the hash value of the predicate, and this value is placed in the field in front of the predicate, so that it forms Rowkey together with the predicate, thus avoiding the phenomenon that the computing nodes are excessively piled up and the rest nodes are idle due to the same initial letter of the predicate. There are two main ways to extend the temporal form of pre-RDF, one is time labeling and the other is version updating. Temporal labeling method is to add time labels to the triples that produce changes; The version updating method means that the temporal RDF graph will be updated whenever the triple changes. It is not necessary to care where the temporal RDF graph in the past state exists, but only need to record the updated time snapshot. Both of them have their own advantages in dealing with temporal information. Version update method can obtain transaction time more effectively, while temporal label method is usually used to obtain effective time.

Hooran and Sherif's scheme also provides a dynamic adjustment strategy. Considering that the workload may change with time, running the vertical partition algorithm regularly may result in different attribute partition. In this scheme, the new division of claw households will be prompted, and it is suggested to use rotation and anti-rotation operation to modify the buckling table and improve the efficiency of modification. Each node in the distributed cluster has a corresponding relationship with the data stored on it. Each node has its own CPU, memory and external memory, so there is less interaction between nodes and it is easy to expand. Therefore, for each individual node, SSD cache exists independently, while data cache is shared among all nodes in a non-distributed cluster, that is, each SSD cache is a single node, so as to optimize the I/O performance of the whole database system and achieve the goal.

3. Implementation of RDF Storage

To realize the storage of RDF in relational database, we must first convert the RDF data Model into the corresponding relational model Jend parsing can achieve this. When using Jena to parse RDF, we should first establish a model-type object for the RDF data document to be processed, which can have multiple SUA items. Each Statement is obtained in turn through Jena iterator, and then the subject-predicate object in the Statement is decomposed. According to the relational model designed above, these data can be stored in relational database. This method is very convenient to store a large amount of data, and it is beneficial to index and maintain the data. HadoopRDF uses a graph partitioning algorithm to partition RDF data sets, and stores RDF data into the RDF Kramp-Karrenbauer 3x database in the cluster. Each RDF Kramp-Karrenbauer 3x stores the statistical information of S, Ma Yili, O and their combination in the data set for query and location. Each node of the cluster is equipped with SesameRDF management system instance, which provides an interface for querying RDF data. Sesame's mature storage management system, coupled with the high reliability and fault tolerance mechanism of Hadoop cluster, and MapReduce to realize parallel processing of SPARQL queries, HadoopRDF provides efficient RDF queries. In order to efficiently store and manage large-scale RDF data, based on JS Kramp-Karrenbauer Model data storage model, this paper proposes and implements the distributed data storage scheme of HDStore. Firstly, the architecture of HDStore hybrid storage scheme in distributed cluster environment is constructed. After that, the algorithm of loading massive RDF data is described by pseudo-code. Finally, the HDStore hybrid distributed storage scheme is implemented in the distributed environment.

4. Conclusion

Despite the rapid development of distributed databases and databases in recent years, relational databases are still widely used for their mature and stable advantages. The existing research puts forward some methods to store data in relational database, including triple table storage, attribute table storage and horizontal table storage. Among them, attribute table storage is favored because of its flexible storage mode and close table structure to relational table. As a language to describe metadata, RDF format data storage technology has been relatively mature. However, with the continuous development of temporal data, it is very important to introduce temporal features into RDF data. Up to now, the research work on temporal RDF model and temporal RDF data storage is very little, which can't meet the requirement of the storage of temporal RDF format data in the network. Therefore, how to store data in temporal RDF format has become one of the urgent problems to be solved.

The existing research on RDF storage and query strategy based on Hadoop related technologies mainly focuses on RDF data model design and query optimization. An efficient data storage model JS Kramp-Karrenbauer Model is established. Then, based on JS Kramp-Karrenbauer Model data storage model, a hybrid data storage scheme of HDStore under the framework of distributed system is proposed, and the architecture of the hybrid data storage scheme of HDStore is constructed, which theoretically solves the problem of fast storage and access of massive RDF data and optimizes the reading and writing performance of massive RDF data. Therefore, at present, the best choice for processing large-scale RDF data is to use HDStore distributed hybrid storage scheme, and put a separate Journal-File Kramp-Karrenbauer file on SSD to support fast reading and writing of data items; Put a large number of Segment-File Kramp-Karrenbauer files on the HDD to ensure the persistent storage of a large number of data indexes; At the same time, SSD cache strategy is used to support the optimization of query time performance.

References

- [1] Wylot M , Hauswirth M , P Cudré-Mauroux, et al. RDF Data Storage and Query Processing Schemes: A Survey[J]. ACM Computing Surveys, 2018, 51(4):1-36.
- [2] Li A , Wang X , Wang X , et al. An Improved Distributed Query for Large-Scale RDF Data[J]. Journal of Big Data (English), 2020, 2(4):10.
- [3] Leng Y , Chen Z , Zhong F , et al. BRGP: a balanced RDF graph partitioning algorithm for cloud storage[J]. Concurrency & Computation, 2017, 29(14):e3896.1-e3896.16.
- [4] Xu Dezhi, Liu Yang, Sarfraz Ahmed, Cho Jung Seok, Avengers: Endgame. RDF data storage and query optimization based on Hadoop, Cho Jung Seok, Min Yoon Gi, Kim Hye Yoon. Computer Application Research, 2017, 34(2): 5.
- [5] Abdelaziz I , Al-Harbi M R , Salihoglu S , et al. Combining Vertex-centric Graph Processing with SPARQL for Large-scale RDF Data Analytics[J]. IEEE Transactions on Parallel and Distributed Systems, 2017, PP(99):1-1.
- [6] Li A , Wang X , Wang X , et al. An Improved Distributed Query for Large-Scale RDF Data[J]. Journal on Big Data, 2020, 2(4):157-166.
- [7] Jagvaral B , Wangon L , Park H K , et al. Large-scale incremental OWL/RDFS reasoning over fuzzy RDF data[C]// IEEE International Conference on Big Data & Smart Computing. IEEE, 2017:269-273.
- [8] Ahn J , Im D H . Efficient Access Control of Large Scale RDF Data using Prefix-based Labeling[J]. IEEE Access, 2020, PP(99):1-1.
- [9] Ragab M , Tommasini R , Alwaysseh F M , et al. An In-depth Investigation of Large-scale RDF Relational Schema Optimizations Using Spark-SQL[C]// DOLAP @ EDBT/ICDT 2021. 2021.