



A Deeper Examination of NoSQL Database Models and Characteristics

Shaida Begum^a

^aCSE department, PDIT, Hospet, India

DOI: <https://doi.org/10.55248/genpi.2022.3.10.56>

ABSTRACT

The main problem has been managing databases because they are continuously expanding quickly and getting more complicated in terms of volume, variety, and velocity. Currently, Relational Database Management Systems (RDMS), also known as SQL, or traditional search engines are primarily used to manage and utilize data collection. Due to their highly semantic properties and applications, relational databases have emerged as powerful and superior during the previous two decades. The handling of a high volume of data and the variability in data type and structure has become a laborious task since Big Data entered the IT industry. Relational databases are inadequate to handle such big data due to their rigid data limitations, structure, relations, and other factors. As a result, data aggregation is rendered impossible. NoSQL databases offer a practical and understandable foundation for combining massive amounts of data, structures, and interactions. Data modeling and migration are now needed in order to define the issue. There isn't yet a tool for switching from relational to no-SQL databases. For this migration, relational (SQL) database queries must be converted to NoSQL database queries. Since relational databases (RDBMS) have found it difficult to keep up with modernization, NoSQL has emerged as the most practical database. This essay aims to highlight and assess the applications of NoSQL, the notion of NoSQL data modeling, and the NoSQL database migration procedure. NoSQL, NoSQL Roadmap, Data Modeling, ETL Approach for Data Migration, Data Loading, and Data Extraction.

Keywords: NoSQL database, Cassandra, Column Database, Big Data, Big Data Analytics, RDBMS, NoSQL, MongoDB, CouchDB, HBase

1. Introduction

Describe NoSQL NoSQL, short for "Not Only SQL," is a mechanism for data storage and retrieval that uses modeling techniques different from the tabular relations found in relational databases. NoSQL is best understood as a database that does not follow the conventional relational database management system (RDMS) structure. It does not alter data using SQL. Even though it may not fully guarantee atomicity, consistency, isolation, and durability, the architecture is distributed and fault tolerant. Internet giants like Facebook, Google, Amazon, and others who needed database management systems that could write and read data anywhere in the globe while scaling and delivering performance across enormous data sets and millions of users were the initial creators and consumers of NoSQL technology. NoSQL is a database system created to meet the needs of cloud applications and get beyond the limits of relational databases (RDBMSs) in terms of scale, performance, data model, and data dissemination. Today, practically every business and organization must provide cloud applications that allow customers to interact with them in a more personalized way. NoSQL is the preferred database technology for such systems. NoSQL data management for big data requirements, 1.2.1 You should use NoSQL in the first place if you have massive data projects to complete. Big data projects are typically identifiable by: Data velocity refers to the quantity and speed of data that is arriving, sometimes from many sources. Data diversity, or the storing of structured, semistructured, and unstructured data. Data volume, or data with a size of several terabytes or petabytes. Data complexity, or the complexity of data that is managed across several locations, data centers, or cloud geo-zones. A variety of diverse elements that can be combined to address the development of the right solution are essential for a successful implementation, as is the case with any sophisticated yet emerging application development framework. Big data's ecosystem is centered on a few essential architectural elements, such as parallel processing, scalable storage, and data management paradigms that are connected by an application development platform. While it's crucial to examine how these many elements function together, from the standpoint of application development, a deeper analysis of large data storage paradigms is

* Corresponding author. Tel.: +0-000-000-0000 ; fax: +0-000-000-0000.

E-mail address: shahida@pdit.ac.in

necessary due to the interaction between analytical algorithms and the underlying data management system. It is simple to understand why this is the case: when applications already rely on a standard relational database (RDBMS) model and/or a data warehousing approach to data management, it may be sufficient to port the RDBMS tools as a manner of scaling performance on a big data appliance. Many algorithms, however, are not designed to consume data in conventional RDBMS systems since they expect to benefit from a high-performance, elastic, distributed data environment. Developers must therefore take into account various data management techniques. These analytical algorithms may make use of any of a variety of alternative data management techniques that are frequently referred to as "NoSQL databases." There are two ideas that are conveyed by the phrase "NoSQL." The first proposes a non-SQL compliant data management framework. The phrase can also be understood to mean "Not Only SQL," which refers to environments that combine conventional SQL (or SQL-like query languages) with different methods of querying and access. This second, more widely accepted (that is, more frequently given) interpretation. 1.2.2 Continual Availability with NoSQL IBM categorizes database accessibility into three groups: 1) High availability: the database and application are accessible during the specified time even if the system is momentarily unavailable for maintenance. 2) Continuous Operation: System up and running constantly with no planned downtime. The combination of HA and CO ensures that data is constantly accessible and that system repair may be carried out without shutting down the entire system. Your need for continually available applications is the reason to think about a NoSQL solution[1][2][3][4][5][6][7]. Keep in mind that this is distinct from simple "high availability," where unplanned downtime is still expected even though it is not intended. Systems with continuous availability don't experience any downtime. Downtime may be fatal to a company's bottom line and reputation in today's market where the competition is only a click away. The average cost of downtime varies significantly amongst companies, from roughly \$90,000 per hour in the media sector to about \$6.48 million for large online brokerages. 1.2.3 Location Independence for NoSQL Data The requirement for real location independence with a database is a third reason to utilize NoSQL. The ability to read from and write to a database without regard to the physical location of those I/O operations, as well as having any write functionality propagated out from that location so that it is accessible to users and machines at other sites, is what is meant practically by the term "location independence." Such capability is simple to describe but challenging to construct for the majority of conventional databases. Location independent read operations can occasionally be accommodated by master/slave and typically shared systems, but putting data everywhere is a different story. The requirement for location independence is driven by a variety of factors, including the need to serve clients across numerous geographies and maintain local data at those locations for quick access. 1.2.4 A More Flexible Data Model with NoSQL: The more adaptable data model featured in most NoSQL systems is one of the main reasons IT professionals switch from a legacy RDBMS to a NoSQL database. The following are the justifications for employing a NoSQL datastore: A NoSQL data model can handle many of these use cases as well as others that don't fit well into an RDBMS, whereas the relational model works well for a few of use cases. Furthermore, compared to a relational database, a NoSQL datastore may easily accept all forms of data, including structured, semi-structured, and unstructured data. A NoSQL database is a fantastic choice for applications that use a variety of datatypes.

2. Review of Literature

Analyzing NoSQL databases:

NoSQL data systems offer what is known as schema-less modeling, a more loose approach to data modeling in which the semantics of the data are incorporated within a flexible connection topology and an associated storage model. Large data sets can be managed with more flexibility as a result, and relational database systems' insistence on a more rigid database structure is also lessened. The flexible model allows for autonomous data distribution and elasticity in the utilization of computation, storage, and network capacity without requiring precise binding of data to be persistently held in certain physical locations. Additionally, NoSQL databases offer built-in data caching, which lowers data access latency and improves performance. The relational structure has been loosened up to make it easier to apply various models to different kinds of studies. Some are designed as key-value stores, for instance, which fit in well with big data programming paradigms like MapReduce. However, it does not enforce adherence to strictly-defined structures, and the models themselves do not necessarily impose any validity rules. Despite the fact that the "relaxed" approach to modeling and management paves the way for performance increases for analytical applications. The dangers of unregulated data management activities, such as unintentional inconsistent data duplication, semantic misinterpretation, and concerns with currency and timeliness, may be introduced as a result. Four distinct NoSQL strategies are covered in this article: Document stores and key-value stores Object stores and tabular stores Key-value databases: Similar to the data structure known as a hash table, a key-value store is a schema-free NoSQL model in which data objects are linked to unique character strings called keys. Unique keys are used to both identify entities and find attribute information about those entities in many NoSQL designs, which is a variation on the key-value theme. This fundamental approach to a schema-less model gains some credence from the widespread use of unique keys. As an illustration, take the data subset shown in Table, where the key is the name of the car manufacturer and the value is a list of the names of the models that belong to that manufacturer. There are no restrictions on data structure or data typing placed by the key-value storage. The semantics of the data organization must be interpreted by the business applications that will use it[8][9][10][11][12][13]. The fundamental operations carried out on a key-value store consist of: Get(key), which gives the value connected to the key supplied. Put(key, value) links the key and value together. The function multi-get(key1, key2,..., keyN), which returns a list of values linked to a list of keys. Delete(key) eliminates the key's entry from the data store. The key needs to be unique when using a key-value store to guarantee that the values can be accessed. The developer must take into account the representations of those values and how they are to be connected to the key in order to correlate many values with a single key (such as the list of automobile models in the example. Key-value stores can be indexed by key value to speed up data queries. They are essentially lengthy, "thin" tables (in that there are not many columns associated with each row). To make finding the key during a query easier, the table's rows can be sorted by the key value. A query basically consists of two steps: calculating the unique key in the first step, and then using that key as an index into the table in the second. It is challenging to anticipate that common SQL-style queries, like "what are the most popular vehicle models based on sales," will be executed due to the requirement to calculate the key in order to access any information about the entity. Instead of employing a query engine, code is often used to address this type of inquiry. The approach does have certain potential drawbacks, despite the

fact that key-value pairs are particularly helpful for both storing and producing the results of analytical algorithms (such as the frequency of specific phrases within huge amounts of documents). The approach won't naturally offer any sort of conventional database capabilities, which is one of its weaknesses (such as atomicity of transactions, or consistency when multiple transactions are executed simultaneously). These features need to be offered by the program itself. Another potential flaw is that preserving unique values as keys may become more challenging as data volume rises; to address this problem, extra complexity must be added to the process of creating character strings that will stand out among an exceedingly large set of keys. For instance, a multinational corporation might try to handle the data related to millions of clients, many of whom have names that are the same or quite similar. The names themselves won't be enough to distinguish between various things because of duplication in the set of names. The end result is that in order to construct a unique key from the combined character string, additional data characteristics must be included. Document Storage b A document store is comparable to a key-value store in that stored items are linked to characterstring keys, which are used to access them. The distinction is that the values being saved, known as "documents," give the managed data some structure and encoding. Various popular, accepted encodings exist, such as XML (Extensible Markup Language), JSON (Java Script Object Notation), and BSON (which is a binary encoding of JSON objects). Other methods of linearizing the data values associated with a data record or object for storage reasons may be used in addition to these conventional ways to data packaging. Examples of data values that were compiled into a "document" that represents the names of particular retail stores. Though the three instances all refer to specific places, the representative models vary. The model's structure is embedded in the document representation, enabling the application to deduce the semantics of the document values[19][20][21][22][23][24]. One significant difference between a document store and a key-value store is that the latter embeds attribute metadata linked to stored content, effectively enabling data querying based on the contents. For instance a search for all documents with "MallLocation" set to "Wheaton Mall" would return a result set with all documents related to any "Retail Store" that is located in that specific shopping center. Tableau Stores Largely derived on Google's original BigTable design for managing structured data, tabular or table-based databases. A NoSQL data management system that developed from BigTable is Hadoop's HBase paradigm. The tabular model enables sparse data to be stored in a three-dimensional table that is indexed by a row key (that is used in a manner similar to the key-value and document stores), a column key that identifies the specific attribute for which a data value is stored, and a timestamp that may refer to the time at which the row's column value was stored. (For background on BigTable design, see this paper via Google's research website.) The HTML code of the page, the URLs of websites linking to it, and the author of the content are only a few examples of attributes that can be linked to a web page's URL. In a BigTable model, columns are organized into "families," and timestamps allow for the management of many versions of an entity. When the content changes, new column affiliations can be formed using the timestamp of when the item was downloaded, preserving history. Object Data Stores (d) In a sense, object data stores and object databases seem to bridge the worlds of schemaless data management and the conventional relational models. Object data stores are essentially a hybrid approach to data storage and administration. On the one hand, methods for object databases can resemble those for document stores, with the exception that object databases keep the object structures in place whereas document stores specifically serialize the object so that the data values are stored as strings. Due to their dependence on object-oriented programming languages like C++, Objective-C, Java, and Smalltalk, this is the case. Object database management solutions are more likely to provide conventional ACID compliance (that is, Atomicity, Consistency, Isolation, and Durability) than some of the other NoSQL models, which are linked to database reliability. It is crucial to keep in mind that object databases are not relational databases and are not searched using SQL, and that this is one of the few parallels between them and a conventional relational database.

Many people have used methods like normalization and de-normalization as well as logical to physical mapping. The recent rise of NoSQL databases, however, has presented additional difficulties for data modeling. In general, NoSQL practitioners prioritize the design of physical data models over the conventional conceptual or logical data model method[14][15][16][17][18]. Modern applications function in distributed, scale-out environments, in contrast to conventional, centralized scale-up systems (including the RDBMS tier). Application developers must forego the conventional conceptual, logical, and physical data model design process in order to implement scale-out since scalability and performance must be addressed first. The standard RDBMS has long been criticized for not supporting massive and unstructured data due to its rigorously set schema and limited scale-out potential. In contrast, NoSQL databases have always had the potential to store large amounts of unstructured data by utilizing flexible schemas that operate in distributed scale-out systems. 3.1 The creation of the data model: The secret to completing medium-to large-scale software projects successfully is designing relational data models. When NoSQL developers take on the responsibility for business and data model design, a new issue with data modeling tools arises. Professionals create traditional RDBMS logical and physical data models using programs like Power Designer or ER/Studio. Professional-grade data modeling tools are not available for NoSQL. As a result, stakeholders examine application source codes to find out information about data models. Non-technical users like business owners and product managers are among these stakeholders. Other methods, such as collecting actual data from databases in use, might be quite time-consuming. We will now need to invest in automation and tooling. The business and data model design methodology depicted in the following diagram should be used for NoSQL initiatives. The document-centric model of MongoDB

3. Discussion on Migration Strategy of MongoDB

MongoDB

For the majority of businesses, data gathered by information technology is one of their most valuable assets. Companies occasionally switch from one information system to another as a result of client needs and technological advancements. As a result, data must be transferred from the legacy system to the new system. Despite the importance of this process, little is known about how migration works. In their article on the data migration issue, F. Matthes and C. Schulz fixed the state of the art and reviewed the literature. According to F. Matthes and C. Schulz, data migration means: "a tool-supported one-time process that migrates formatted data from a source data structure to a target data structure when both structures differ conceptually and/or physically." Image: Data Migration Method The following general migration strategy would be used: A. Data preparation in JSON file formatB. Data extraction into flat files using the JSON file format or data extraction using customized data loaders from the processed data store [14][15][16][17][18].

Utilizing built-in or custom loaders to load data into NoSQL data structures (s) The following sections go into more detail about the numerous actions for each stage of migration. 4.1 Data Preparation and Extraction Data extraction, transformation, and loading is known as ETL, and reconciliation plays a significant role in the process's conclusion. This includes data validation using business processes. Before the data is loaded into staging tables, the ETL process also incorporates data validation and enrichment. Preparation of Data During data preparation, the following activities will be carried out: First, make database objects. 2) Based on the requirements, necessary staging tables that mimic a typical open interface or base table structure are to be established. 3) Validate and transform the data before loading it from the specified source (Dumps or Flat files).4: Data Purification Use the JSON file layout guidelines to filter out inaccurate data. 6) Using the JSON file structure guidelines, remove any unnecessary information.7) Remove any unnecessary information using the JSON file formatting guidelines.8. Put the data into the staging area. Enrichment of Data (9)10) Inadequate default data11) Determine missing data using mapping or lookups.12) data that is structured differently (1 record as-is = many records to-be). Data Extraction Activities, Section 4.3 (into JSON files) During the process of extracting data into JSON file formats, the following operations will be carried out: 1) Data selection based on JSON file structure2) SQL application development using the JSON file format3) PLSQL applications or scripts are built based on the ETL procedures and the data mapping specifications. These applications will be used for a variety of tasks, such as loading data into staging tables and open interface standard tables. 4) Data transformation prior to extraction in accordance with JSON file mapping and layout specifications.5) JSON-formatted flat files for data loading4.4 Data Loading You can access NoSQL data structures using a variety of programming languages, including (.net, C, Java, Python, Ruby, etc.). Using these programming languages, data can be loaded directly from relational databases (such as Access, SQL Server, Oracle, MySQL, IBM DB2, etc.). Depending on the enactment rules, level of customization, and kind of data processing, custom loaders may be used to load data into NoSQL data structures (s).

4. Summary

This paper's focus was on the taxonomy of NoSQL by describing its modeling and migration methodologies, as well as the reasons and benefits for using NoSQL. It also demonstrated how NoSQL is efficiently used to store massive amounts of data with great scalability. Additionally, this article explains the benefits of choosing a NoSQL database and how to use it efficiently. It is therefore anticipated that it will aid in the adoption of NoSQL databases, where databases are made to expand as they grow. The choice of which database to use depends on the type of application the system will be employing because the two databases (SQL and NoSQL) behave differently depending on the type of queries executed. The choice of a NoSQL data store over a relational model must be in line with the expectations of business users. How well will a NoSQL data store perform in comparison to their prior experiences with relational models? As should be obvious, a lot of NoSQL data management environments are designed for two main requirements: Fast accessibility, whether that means inserting data into the model or pulling it out via some query or access method, and Scalability for volume, in order to support the accumulation and management of enormous amounts of data. Distribution and parallelization can be used to fulfill each of these requirements, and the NoSQL models discussed above can be scaled, distributed, and extensible. Additionally, these distinguishing characteristics fit in well with programming paradigms like MapReduce that efficiently control the formation and operation of several parallel execution threads. Utilizing data dissemination is the key. Fortunately, several queries and data accesses can be carried out concurrently when using a distributed tabular data store or key-value store, especially when the hashing of the keys maps to various data storage nodes. Smart data allocation strategies will enable linear performance scalability in relation to data volume. NoSQL methods are designed for high performance computing for reporting and analysis. Numerous new businesses are adopting the various NoSQL concepts and releasing their own tailored iterations on the market. If you are interested in NoSQL, there is no risk in experimenting with the various methods, therefore it would make sense to create a straightforward "pilot" project model that can be implemented in a variety of ways in order to compare and contrast ease-of-use, space performance, and execution speed. NoSQL data management solutions have enticing performance characteristics; however they cannot fully take the place of relational database management systems. Choosing to use NoSQL is not always a simple choice. Before committing to the technology, one must consider the business needs as well as the capabilities required to switch from a traditional method to a NoSQL approach. The idea that you automatically made the proper decision when using SQL or NoSQL technologies is one you can have if you choose them haphazardly or because they're in demand. The best architecture will depend on the needs of the apps you design, as both SQL and NoSQL have advantages and disadvantages. That temperamental old SQL database is still incredibly powerful and capable of ably meeting your transactional demands. When you are approaching the outer limits of relational databases and the size of your data processing or scale of activities just requires a more dispersed system, look at NoSQL possibilities. Free your data and create the most incredible applications ever with careful decision-making!

REFERENCES

- [1] Abramova, V., & Bernardino, J. (2013, July). NoSQL databases: MongoDB vs Cassandra. In Proceedings of the international C* conference on computer science and software engineering (pp. 14-22).
- [2] Ali, W., Shafique, M. U., Majeed, M. A., & Raza, A. (2019). Comparison between SQL and NoSQL Databases and Their Relationship with Big Data Analytics. *Asian Journal of Research in Computer Science*, 4(2), 1-10
- [3] Becker, M. Y., & Sewell, P. (2004, June). Cassandra: Flexible trust management, applied to electronic health records. In Proceedings. 17th IEEE Computer Security Foundations Workshop, 2004. (pp. 139-154). IEEE.
- [4] Berg, K. L., Seymour, T., & Goel, R. (2013). History of databases. *International Journal of Management & Information Systems (IJMIS)*, 17(1), 29-36.

-
- [5] Bjeladinovic, S., Marjanovic, Z., & Babarogic, S. (2020). A proposal of architecture for integration and uniform use of hybrid SQL/NoSQL database components. *Journal of Systems and Software*, 168, 110633.
- [6] Chandra, D. G. (2015). BASE analysis of NoSQL database. *Future Generation Computer Systems*, 52, 13-21.
- [7] Chen, J. K., & Lee, W. Z. (2019). An introduction of NoSQL databases based on their categories and application industries. *Algorithms*, 12(5), 106.
- [8] Cuzzocrea, A., & Shahriar, H. (2017, December). Data masking techniques for NoSQL database security: A systematic review. In *2017 IEEE International Conference on Big Data (Big Data)* (pp. 4467-4473). IEEE.
- [9] de Oliveira, V. F., Pessoa, M. A. D. O., Junqueira, F., & Miyagi, P. E. (2021). SQL and NoSQL Databases in the Context of Industry 4.0. *Machines*, 10(1), 20.
- [10] Deka, G. C. (2013). A survey of cloud database systems. *IT Professional*, 16(2), 50-57. IEEE.
- [11] Di Martino, S., Fiadone, L., Peron, A., Riccabone, A., & Vitale, V. N. (2019, June). Industrial Internet of Things: Persistence for Time Series with NoSQL Databases. In *2019 IEEE 28th International Conference on Enabling Technologies: Infrastructure for Collaborative Enterprises (WETICE)* (pp. 340-345). IEEE.
- [12] dos Santos Ferreira, G., Calil, A., & dos Santos Mello, R. (2013, December). On providing DDL support for a relational layer over a document NoSQL database. In *Proceedings of International Conference on Information Integration and Web-based Applications & Services* (pp. 125-132).
- [13] Gessert, F., Wingerath, W., Friedrich, S., & Ritter, N. (2017). NoSQL database systems: a survey and decision guidance. *Computer Science-Research and Development*, 32(3), 353-365.
- [14] Guimaraes, V., Hondo, F., Almeida, R., Vera, H., Holanda, M., Araujo, A., ... & Lifschitz, S. (2015, November). A study of genomic data provenance in NoSQL document-oriented database systems. In *2015 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)* (pp. 1525-1531). IEEE.
- [15] Rodriguez, K. M., Reddy, R. S., Barreiros, A. Q., & Zehtab, M. (2012, June). Optimizing Program Operations: Creating a Web-Based Application to Assign and Monitor Patient Outcomes, Educator Productivity and Service Reimbursement. In *DIABETES* (Vol. 61, pp. A631-A631). 1701 N BEAUREGARD ST, ALEXANDRIA, VA 22311-1717 USA: AMER DIABETES ASSOC.
- [16] Kwon, D., Reddy, R., & Reis, I. M. (2021). ABCMETAapp: R shiny application for simulation-based estimation of mean and standard deviation for meta-analysis via approximate Bayesian computation. *Research synthesis methods*, 12(6), 842-848. <https://doi.org/10.1002/jrsm.1505>
- [17] Reddy, H. B. S., Reddy, R. R. S., Jonnalagadda, R., Singh, P., & Gogineni, A. (2022). Usability Evaluation of an Unpopular Restaurant Recommender Web Application Zomato. *Asian Journal of Research in Computer Science*, 13(4), 12-33.
- [18] Reddy, H. B. S., Reddy, R. R. S., Jonnalagadda, R., Singh, P., & Gogineni, A. (2022). Analysis of the Unexplored Security Issues Common to All Types of NoSQL Databases. *Asian Journal of Research in Computer Science*, 14(1), 1-12.
- [19] Singh, P., Williams, K., Jonnalagadda, R., Gogineni, A., & Reddy, R. R. (2022). International students: What's missing and what matters. *Open Journal of Social Sciences*, 10(02),
- [20] Jonnalagadda, R., Singh, P., Gogineni, A., Reddy, R. R., & Reddy, H. B. (2022). Developing, implementing and evaluating training for online graduate teaching assistants based on Addie Model. *Asian Journal of Education and Social Studies*, 1-10.
- [21] Sarmiento, J. M., Gogineni, A., Bernstein, J. N., Lee, C., Lineen, E. B., Pust, G. D., & Byers, P. M. (2020). Alcohol/illicit substance use in fatal motorcycle crashes. *Journal of surgical research*, 256, 243-250.
- [22] Brown, M. E., Rizzuto, T., & Singh, P. (2019). Strategic compatibility, collaboration and collective impact for community change. *Leadership & Organization Development Journal*.
- [23] Sprague-Jones, J., Singh, P., Rousseau, M., Counts, J., & Firman, C. (2020). The Protective Factors Survey: Establishing validity and reliability of a self-report measure of protective factors against child maltreatment. *Children and Youth Services Review*, 111, 104868