I'm not robot

reCAPTCHA

**Continue**

I'm not robot

reCAPTCHA

**Continue**

# Designing data-intensive applications goodreads

Member Since March 2015 Show 1-33 Start your review of the design data of intensive web applications Lalitha rated it was normal October 22, 2013 Ekta appreciated he liked it On March 28, 2020 Japuanna appreciated it really liked January 30, 2017 Ashish appreciated it really liked November 05, 2016 Vidya appreciated it very much liked Sep 02, 2013 Mmm he appreciated disliked it on August 10 , 2017 Clodpated marked it as a reading Sep 09, 2011 Ezzo noted it as a k-read July 17, 2013 Shudhodhan marked it as a reading september 26, 2013 Chris tagged it as to read October 18, 2014 Yang Chang tagged it as a reading may 31, 2017 Darren marked it as a k-read February 15, 2018 Aristotl noted it as a reading March 19, 2018 Azza noted it as a k-read January 16 , 2019 I consider this book a mini encyclopedia of modern data engineering. Like the specialized encyclopedia, it covers a wide field in great detail. But it's not a practice or cookbook for a particular big data product, NoS'L or newS'L. The author should set out the principles of the current distributed big data systems, and he is doing a very good job. If you're after obscure details of a particular product, or some tutorials and how-tos, go elsewhere. But if you want to undie I find this book a mini encyclopedia of modern data engineering. Like the specialized encyclopedia, it covers a wide field in great detail. But it's not a practice or cookbook for a particular big data product, NoS'L or newS'L. The author should set out the principles of the current distributed big data systems, and he is doing a very good job. If you're after obscure details of a particular product, or some tutorials and how-tos, go elsewhere. But if you want to understand the basic principles, problems, and problems of an intense and distributed data system, you've come to the right place. Martin Kleppmann begins by giving the reader a conceptual basis in the first chapter: what does reliability mean? How is this defined? What is the difference between guilt and failure? How do you describe the load on the intensive data processing system? How do you talk about performance and scalability in a meaningful way? What does it mean to have a supported system? The second chapter gives a brief overview of the different data models and shows their suitability for different usages, using the current challenges faced by companies such as Twitter. This chapter is a solid basis for understanding the difference between the relational data model, the document data model, the graphic data model, and the languages used to process the data stored by these models. The third chapter details the building blocks of different types of database systems: data structures and algorithms used for different types of database systems Shown in the previous chapter. you'll learn hash indices, SSTables (sorted string tables), strings), Merging trees (LSM trees), B-trees and other data structures. After this chapter, you are injected into the database columns, as well as the basic principles and structures underlying them. Following this, the book describes data coding techniques ranging from the venerable XML and JSON, and going into detail formats such as Avro, Thrift and buffer protocol, showing the trade-offs between these elections. After the building blocks and foundations comes Part II, and this is where things start to get really interesting, because now the reader is beginning to learn about the complex topic of distributed systems: how to use basic building blocks in an environment where things can go wrong in the most unexpected ways. Part II is the most challenging part of the book: you'll learn how to reproduce your data, what happens when replication lags behind, how you provide a consistent picture for the end user or end programmer, what algorithms are used to elect leaders in consensus systems, and how replication works without leaders. One of the main goals of using a distributed system is to have an advantage over one central system, and this advantage is to provide a higher service, which means a more sustainable service with an acceptable level of response. This means that you need to distribute the load and your data, and there are many schemes for separating the data. Chapter 6 of Part II provides a lot of detailed information about separation, keys, indices, secondary indices, and how to handle data requests when sharing data using different methods. No data systems book can be complete without touching on the topic of transactions, and this book is no exception to the rule. You'll learn about the fusambience surrounding the definition of ACID, isolation levels, and serialization. The other two chapters of Part II, Chapters 8 and 9 are probably the most interesting part of the book. Now you are ready to learn details on how to deal with all kinds of network and other types of malfunctions to keep your data system in the right and consistent state, problems with cap theorem, version vectors and that they are not vector clocks, Byzantine malfunctions, how to have a sense of causality and order in a distributed system, why there are algorithms such as Paxos, Raft and AA (used in zooKeeper) distributed transactions and many other topics. The rest of the book, i.e. Part III, is devoted to the processing of batches and streams. The author describes in detail the famous Map Reduce package processing model and briefly touches on the modern framework for processing distributed data, such as Apache Spark. The final chapter looks at the flow of events and messaging systems and the problems that occur when trying to process that data in motion. You may not be in the business of creating the next generation of streaming but you definitely have to handle on these topics because you will run into described problems in the practical thread processing systems that you deal with daily as a data engineer. As I said in opening this review, consider this mini encyclopedia for a modern data engineer, and don't be surprised if you see over 100 links at the end of some chapters; if the author tried to include most of them in the text itself, the book may well go beyond 2,000 pages! At the time of writing, the book is 90% complete, according to its official website, there is only one chapter to be added (Chapter 12: Materialized Views and Cashing), so it's safe to say that I recommend this book to anyone who works with distributed big data systems when dealing with NoS'L and newS'L databases, stores, document stores, data-focused speakers. As for me, this will certainly be my go-to link to the coming years on these topics. ... More Building on the success of the bestselling Pure Coder and Clean Code, legendary software master Robert C. Uncle Bob Martin shows how to bring greater professionalism and discipline to applica... Since this Jolt-winning classic was last updated in 2008 (shortly after the release of Java 6), Java has changed dramatically. The main improvement in Java 8 was the addition of functional pr.... Go is an open source programming language that makes it easier to create clean, reliable, and efficient software. It has been winning converts from dynamic language enthusiasts as well as to tradit... Java 8 in action is a well-written guide to the new Java 8 features. The book covers lambdas, streams, and functional programming style. With the functional features of Java 8 now you can write more... Programming at Scala is the final book on Scala, a new java Platform language that combines object-oriented and functional programming concepts into a unique and powerful tool for design... I was lucky enough to really work with a fantastic team to design and implement the concurrency features added to Java's Java 5.0 and Java 6. Now this same team provides ... This title covers a wide range of algorithms in depth, but makes their design and analysis accessible to all levels of readers. Each chapter is relatively autonomous and can be used as a unit... Monolithic architecture works well for small, simple applications. However, successful apps have a habit of growing. Eventually the development team gets into what is known as monolithic... You're not alone. At any moment, somewhere in the world someone is struggling with the same Designing the software you have. You know you don't want to reinvent the wheel (or, worse, with a flat wheel) .... Spring in Action 2E is an extended, fully updated second edition of the bestselling Spring in Action. Action. Walls, one of Manning's best writers, this book covers an exciting new fairy... 2% For example, Amazon describes response time requirements for internal services in terms of 99.9 percentile, although this affects only 1 in 1,000 requests. This is because customers with the slowest queries are often the ones who have the most data on their accounts because they have made a lot of purchases - that is, they are the most valuable customers of Umesh Sai Gurram and 3 other people liked it 3% Amazon also noted that a 100 ms increase in response time reduces sales by 1% , and others report that a 1-second slump reduces customer satisfaction by 16% YMM. At Goodreads, we haven't seen anywhere near the same benefits. It should depend on how impulsive the activity is. Reading against buying a dog Halloween costume sale. Eric Franklin and 3 other people liked this 3% elastic system can be useful if the load is very unpredictable, but manually scaled the system easier and perhaps fewer operational surprises yet heard of any groups really using automatic zoom up and down. 3% For this reason, common sense until recently was to keep your database on one node (scale up) until scaling cost or high availability requirements forced you to make it distributed. In a startup at an early stage or an unproven product, it is usually more important to be able to quickly iterate on product characteristics than to scale to some hypothetical future load. Rails 3% Is well aware that most of the software's value is not in its initial development, but in its constant maintenance - fixing bugs, maintaining workable systems, investigating failures, adapting it to new platforms, changing it for new usage, paying off technical debt and adding new features. Well known on paper, but well ignored in practice. David and 4 other people liked this 5% in relational databases, it's okay to refer to lines in other ID tables because connections are simple. 5% However, if the application uses many-to-many relationships, the document model becomes less attractive. You can reduce the need for connections by denormalization, but then the application code has to do more work to keep the denormalized data consistent. 6% more accurate term schema-on-read (data structure is implicit, and interpreted only when reading data), as opposed to a diagram-on-the-record (traditional relational database approach where the diagram is explicit and the database ensures that all written data corresponds to it) 6% Schema-on-read is similar to the dynamic (time of execution) type of verification in programming languages, while the sche-chart is similar to the dynamic (time of execution) type of verification in programming languages. (time compilation) check type. Just as proponents of static and dynamic type checks have a great debate about their relative merits, the execution of circuits in the database is a controversial topic, and in there is no right or wrong answer. The 6% Hybrid relational and documentary models is a good route for databases to adopt in the future. 7% In a web browser, using the declarative style of CSS is much better than manipulating styles imperatively in JavaScript. Similarly, in databases, declarative query languages such as S'L have proven to be much better than imperative query APIs. This analogy has been very helpful to me to understand the benefits of making your code more declarative. 7% group name. called it! 8% If the same query can be written in 4 lines in one query language but requires 29 lines in another, it just shows that different data models are designed to meet different usage cases. It's important to choose a data model that fits the app. Also a good ruby comparison against Java.... Unfortunately, the semantic web was overhyped in the early 2000s, but has yet to show any signs of being implemented in practice, which has made many people cynical about it. LOL found out about it in graduate school. While tiny parts of it such as google's e-mail information collection shows that there are many promises in the idea of a 9% document database targeting the use of cases where data arrives in standalone documents and the relationship between one document and another is rare. 10% of researchers working with genome data often have to search for a similarity sequence, which means taking one very long line (representing a DNA molecule) and matching it with a large database of strings that are similar but not identical. None of the databases described here can handle this type of use, so the researchers have written specialized genome database software, such as GenBank. 11% Compaction means throwing away duplicate keys in a log and storing only the latest update for each key. 11% Control Concurrency As it writes, nailed to the magazine in a strictly consistent manner, the overall choice of implementation is to have only one stream writer. Data file segments are only applications and otherwise immutable, so they can be read in multiple threads at the same time. The 11% Application and merging segment are consistent recording operations that are usually much faster than casual writings, especially on magnetic spinning drive hard drives. 12% Most databases can fit in a B-tree that is three or four levels deep, so you don't need to follow many page links to find the page you're looking for. (A four-level 4-page KB tree with a 500 branch factor can store up to 256TB.) 12% in order to make the database resistant to failures, usually b-tree implementation additional data structure on the disk: Journal of Record Forward (WAL, also known as the journal remade). This is an application-only file to which every B-tree modification must be written before it can be applied to the pages of the tree itself. When the database returns after a failure, The magazine is used to restore the B-tree back to a consistent state of 12% (although it depends on the configuration of the storage engine and workload), we are officially in the butts ... 13% This machine can be converted into a Levenshtein machine gun, which supports effective word search within a given distance of editing wow we are in a groov for other wills now. Zzz... 13% Of Redis and Couchbase provide weak strength by writing asynchronously on the disc. 13% This process of getting data into a warehouse known as Extract-Transform-Load (ETL) is 13% In a large company, a lot of hard work is required to do something simple in a small company. 14% of Amazon RedShift is a featured version of ParAccel. 14% of the many open source projects have appeared in S'L-on-Hadoop; they are young but seek to compete with commercial storage systems. These include Apache Hive, Spark S'L, Cloudera, not to mention Snout'L 15% Although this chapter may not make you an expert in setting up any one particular storage engine, it is hopefully equipped with enough vocabulary and ideas that you can understand the documentation for the database of your choice. or ask your dba;-) 16% This is often a source of security problems: if an attacker can force your application to decipher an arbitrary byte sequence, they can instantly arbitrary classes looking at you, yml 18% of five-year data will still be there in the original coding if you've explicitly rewrote it since then. This observation sometimes adds up as the data survives the code. 18% This way of creating applications is traditionally called service-oriented architecture (SOA), recently improved and rebranded as a microservice architecture (31, 32). A 19% definition format such as OpenAPI, also known as Swagger, can be used to describe the RESTful API and obtain documentation. 19% Every time you call a local function, you usually run roughly the same time. The network's query is much slower than the function call, and its delay is also wildly variable: in good times it can complete in less than a millisecond, but when the network is overloaded or the remote service is overloaded, it can take many seconds to do exactly the same. 19% All of these factors mean that there is no point in trying to make a remote service too similar to a local object in the programming language, because it is fundamentally different. Part of the appeal of REST is that it doesn't try to hide the fact that it's a network protocol (although that doesn't seem to stop people from creating RPC libraries on top of REST). Coral FTL 19% RESTful API has other significant advantages: it's good for experimentation and debugging (you can just make requests for it using or the command line tool curl, without any code generation or installation of software), it's this all major languages and programming platforms, and there is an extensive ecosystem of available tools (servers, caches, load balancers, proxies, firewalls, monitoring, debugging tools, testing tools, etc.). 22% If the replication protocol does not allow this version to be inconsistent, as is often the case with WAL delivery, such updates require downtime. The 22% alternative is to use different magazine formats for replication and engine storage, allowing the replication log to separate from the interior of the storage engine. This type of replication log is called a logical journal to distinguish it from the (physical) representation of storage engine data. 22% This discrepancy is only a temporary state - if you stop writing in the database and wait a while, followers will eventually catch up and become according to the leader. For this reason, this effect is known as the ultimate level of consistency of 23% other criteria can be used to decide whether to read from a leader. For example, you can track the time of the last update and do all the readings from the leader within one minute of the last update. You can also control the backlog of replication from followers and prevent requests for any follower who is more than one minute behind the leader. 23% When you read the data, you can see the old value; monotony reading means only that if one user makes a few reads in sequence, they won't see the time going backwards, i.e. they won't read the old data after they've read the new data before. 23% For example, auto-installation keys, triggers and integrity limitations can be problematic. For this reason, multi-head replication is often considered dangerous territory, which should be avoided as far as possible 24% The simplest strategy to fight conflicts is to avoid them: if the application can guarantee that all entries for a particular record pass through the same leader, then conflicts cannot occur. Since many multi-liberal replication implementations handle conflicts quite poorly, avoiding conflicts is often recommended by the 24% timestamp approach used, this method is known as the last write wins (LWW). While this approach is popular, it is dangerously prone to data loss 24% Amazon is an oft-quoted example of surprising effects due to the conflict resolution handler: for some time, the logic of resolving conflicts on the basket will retain items added to the cart, but not items removed from the basket. Thus, customers sometimes see items appearing in the baskets, even if they had previously been removed 25% If they were able to write on some replicas, but not on others (for example, because it are full), and generally managed on smaller than W replicas, it's not a rollback on the replica where it succeeded. 25% The latter is known as a sloppy quorum : writes and reads still require w and g successful answers, but may include nodes that are not included in the designated n home node for value. By analogy, if you lock yourself in your house, you can knock on a neighbor's door and ask if you can stay on the couch temporarily. 26% of the siblings' union. To prevent this problem, the item cannot simply be removed from the database when it is removed; instead, the system should leave a marker with the appropriate version number to indicate that the item was removed when the siblings were merging. This marker is known as a tombstone. For example, Data support riak uses a family of data structures called CRDTs (38, 39, 55) that can automatically connect siblings in a reasonable way, including saving removals. Remove.