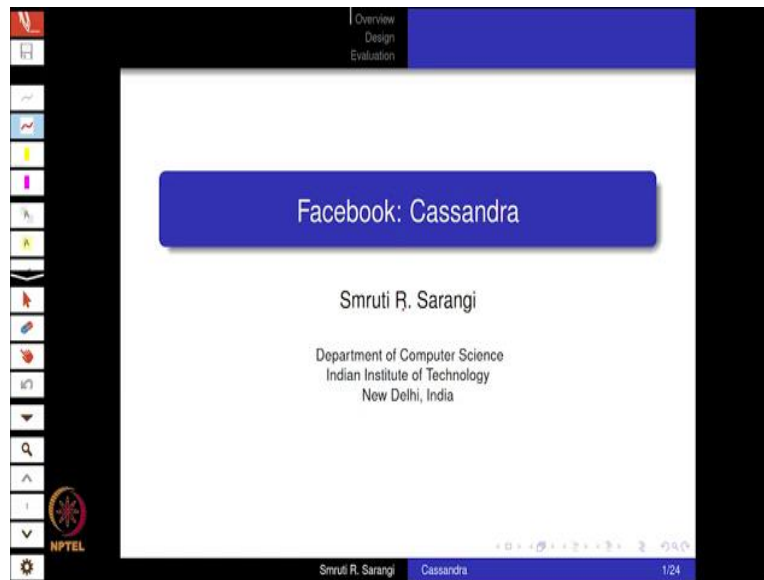


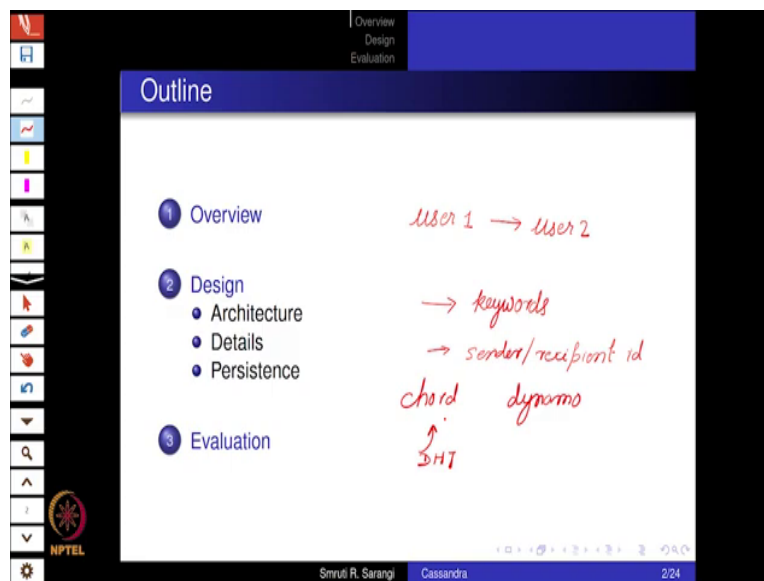
**Advanced Distributed Systems**  
**Professor Smruti R. Sarangi**  
**Department of Computer Science and Engineering**  
**Indian Institute of Technology, Delhi**  
**Lecture 21**  
**Cassandra: Distributed Storage Engine**

(Refer Slide Time: 0:17)



In this lecture, we will discuss the Facebook Cassandra system. The Cassandra system is a key component of Facebook that is used in the inbox search. So, we will discuss, what is inbox search? Over the next few slides.

(Refer Slide Time: 0:35)



Of course, this does not need any introduction, because almost all of us are very active Facebook users. So, all of us know what an inbox means. Facebook has a messaging system between users. So, it is possible for a user 1 to send a message to user 2 that appears on user 2's dashboard.

So, active users can also get hundreds of messages per day. And it is the job of Facebook to store all of this data in a database of sorts. And furthermore, it should be possible to query these messages based on keywords. And based on the recipient, and basically based on either who the message is being sent to, or who the message is being, or either the id of the sender or the id of the recipient, let us call it the sender or the recipient id.

So, actually maintaining an inbox search for a billion plus users is a rather difficult challenge. So, that is why the Cassandra system was developed. To solve this particular problem, we can think of it as a logical extension of the chord and Dynamo projects. So, that is why it is important to have at least an idea of chord, the chord DHT before proceeding through this lecture. And also, if the viewers have an idea of the Amazon Dynamo project, it would help. So, Amazon Dynamo is described in one of the previous lectures of this lecture series.

(Refer Slide Time: 02:33)

Overview  
Design  
Evaluation

## Motivation

Motivation

- Social Networking ⇒ Facebook
- Facebook is **HUGE**
  - Billion+ users
  - Millions of updates per minute
  - Many thousands of servers *[multiple data centers]*
  - Reliability is a big issue
- Cassandra: Made for the Inbox Search problem
  - Allows users to search through their Facebook Inbox
- It is also used for many other Facebook services

NPTEL

Srinu R. Sarangi Cassandra 3/24

So, nowadays, of course, social networking means implies Facebook, Facebook is the de facto social networking site. There used to be others like Google Plus and orkut, but they are not there anymore. Facebook, of course, is huge. It has billion plus users, millions of updates per minute. In this case, millions of messages are sent between users every minute, users post photographs, users write posts.

So, that is the reason internally, similar to other large systems that we have seen. Facebook will have 1000s of servers that are physically partitioned across multiple data centers. And we will need to say, in such a large distributed system, reliability is a big issue, massive issue.

So, Cassandra, was made for the inbox search problem. It allows users to search through their Facebook inbox. And it has also internally been used to operate many other Facebook services. So, as we have been kind of seeing with Google percolator, and Amazon Dynamo that many of these systems are designed to be generic. So, I said at least the company that is designing them, it can use this software for other things as well to implement other services as well.

(Refer Slide Time: 04:13)

Overview  
Design  
Evaluation

### Related Work

- Distributed File Systems: Coda, AFS: Hierarchical name space, specialized conflict resolution.
- Distributed Databases: Bayou *get → all versions*
- Dynamo: Uses vector clocks to manage versions. Conflict resolution based on the vector timestamp, and specialized business logic.
  - Gossip based membership algorithm *membership changes are rare*
- BigTable: Distributes data. However, it relies on an underlying distributed file system for durability.

*write throughput*     *GFS overheads*

Srinu R. Sarangi     Cassandra     4/24

So, a bit of related work, we have already looked at distributed file systems, we have looked at Coda and AFS. So, a distributed file system can definitely be used to implement an inbox, where pretty much every message can be written to a file, or a set of messages can return to a file. And then, so, even though this does appear kind of appealing, the main problem over here is that inbox is not really a file system. And this causes an issue.

Furthermore, file systems have a hierarchical namespace, and they have specialized conflict resolution. So, we do not need all that these things in an inbox search, all these things are not really required. And so that is the reason using a traditional file system to implement an inbox is not really a good idea.

On similar lines using a distributed database like Bayou, that is also not a good idea. Because a database has actually far more structure. So, it is okay to store airline tickets, for example.

But it is not really a good idea to implement an inbox using a database. We have seen 2 other projects we have seen Dynamo, Amazon Dynamo.

So, this is the system that Amazon uses to store its shopping carts. So, in this case, Amazon uses vector clocks to manage the versions. And whenever there is a get, we actually retrieve the entire context, which is pretty much what we get from a get message is actually that we get all the versions and then either the system or the user, they need to perform an automatic merge. The conflict resolution is based on the vector timestamp, specialized business logic. And of course, as we had discussed in Amazon, it will be the same in Facebook also.

Membership changes to a ring, are kind of infrequent. So, this is the same in Facebook, that whenever so we do not regularly add or remove servers. These are relatively infrequent operations, some membership changes are rare. So, that is a reason even in Dynamo we had a console based system. And even in Cassandra, also, we have a slow console base system.

But the main problem with Dynamo is that it is great for shopping carts is great for kind of very structured objects. But inbox is a very simple thing. And it is and for this, what we need actually, is we need a very high write throughput. Of course, latency does not matter that much for light for writes.

But we need a very, very high write throughput. And so basically, that is a prime goal. And Dynamo is not really designed for a very high write throughput, but it is designed for actually very good write latency, and also for aggregating multiple versions of an object. So, it never loses a write. We have also looked at a little bit of Google Big Table in the percolator project. So, big table is pretty pretty much a multi dimensional table where every row is indexed by the key. So, basically, it is a think of it as a key value store in that sense.

So, one problem is that percolator sorry, Bigtable relies on GFS Google File System for the underlying implementation. And the Google File System, what it does is that it divides the large table into chunks, and stores these chunks. So, we do not want structure such a strong relationship with the underlying file system, because that also adds its own overheads. So, that is the reason there is a need to create a bespoke system over here, a custom system over here.

(Refer Slide Time: 08:35)

Overview  
Design  
Evaluation

Architecture  
Details  
Persistence

## Outline

- 1 Overview
- 2 Design
  - Architecture
  - Details
  - Persistence
- 3 Evaluation

Snruti R. Sarangi | Cassandra | 5/24

Overview  
Design  
Evaluation

Architecture  
Details  
Persistence

## Data Model

700 → Locations

- Table → Distributed multi-dimensional Map
- Key → 16 to 36 byte long string to uniquely identify a row
- Value → Highly structured object
- Every row operation is **atomic** on a replica
- Columns can be grouped into **families**
  - Simple and Super (family within a family)
  - Columns can be sorted by either the time or by name

Super Column

Snruti R. Sarangi | Cassandra | 6/24

And of course, Facebook is large enough that it can afford to create a separate, fully designed system from scratch to store in Facebook inbox messages. So, we will define a couple of terms over here that will be used pretty regularly over the next 18 slides or so. So, first, the table, again to Google big table, it says distributed multi dimensional map, which means that it is a very, very large table. So, we can think of let us consider a 2d table to start with. So, every row can actually have a large number of columns. So, in a certain sense this table is partitioned can be row wise can be column wires can be block wise, and they are physically stored at different locations.

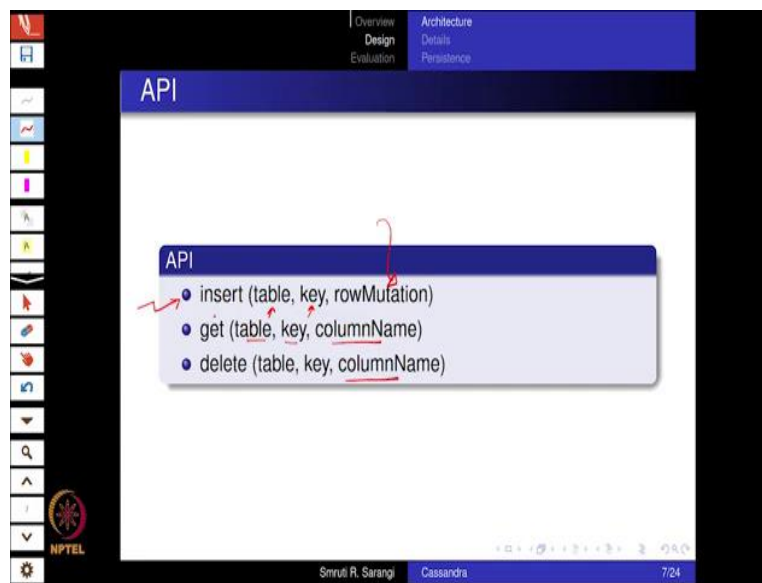
So, the physical storage is at different locations. And so, it is so this table is essentially stored in a distributed fashion. And of course, each slide a key value store for every row has a unique key, which in this case is a 16 to 36-byte long string quantity to identify a row. And a value, of

course, in this case is a highly structured object, we will see how and why. And furthermore, similar to big table, we assume that every row operation is atomic on a given replica.

So, this is important, every new operation is not atomic on like all replicas. But for a given replica, it is atomic. So, a given replica will never have a mixed row state. And for the columns, which we can have many, the columns are grouped into families. So, we can define a super column that essentially contains a lot of child columns.

So, this is possible to do, and a column themselves can either be sorted by time, time of modification time of creation, or by name. So, it is important that the columns are not just one hierarchy, but in the columns have a hierarchy within them. And so, we can group columns into families and we then we can create a family within a family. And so that would pretty much be a super column.

(Refer Slide Time: 11:23)



So, the API is a rather simple, the API is pretty much as simple as it gets. So, it is similar to broadly what we have seen in percolator, which is that inserting into the table is rather easy. We give the table id the key which uniquely identifies the row. And in this case, as we have been seeing in distributed systems up till now, anytime we create a row, we actually create a new version of it.

So, the row mutation basically means that we are creating a new version of the row. And get, well it is the same thing, but we get a certain column name. And we can also delete a certain column or we can delete the contents of a certain column. So, insert get and delete is a very

standard way that almost all of these tables operate. That includes percolator dynamo, big table, everything we have seen up till now. And the interface in that sense is rather standardized.

(Refer Slide Time: 12:30)

System Architecture

- Cassandra has a cluster of nodes
- A read/write request for a key gets routed to a node in the cluster
- The node determines the replicas that contain the key
- We write to the write quorum
- Either read from the closest replica or from the read quorum

Chord

node

10

$R+W > 10$

read quorum

NPTEL

Smit R. Sarangi | Cassandra | 8/24

And one reason for having a standardized interface is that it is useful both for inbox search, as well as Facebook can use it for other applications as well and at some point of time. Well, so some versions of Cassandra are available publicly. But let us say at some point of time, if companies decide to open source such kind of systems, developers can use them to build a wide variety of applications, which can benefit the host company as well, that is also one reason why open sourcing internal projects is typically a good idea.

So, in this case, Cassandra has a cluster of nodes, it is a distributed system, a read write request for a key will get routed to a node in a cluster. So, this cluster is of course, organized as a DHT. And this is the DHT that they use as called. So, that is why I said that chord is a understanding chord knowing chord is a prerequisite for this lecture.

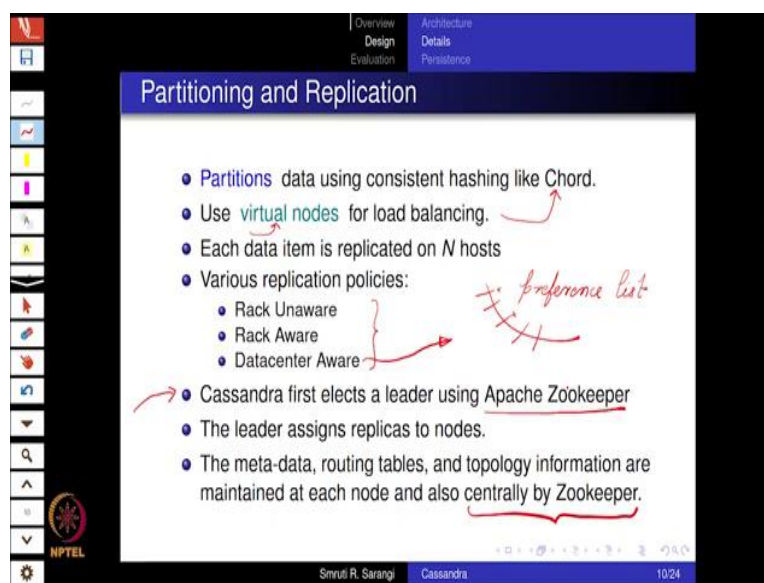
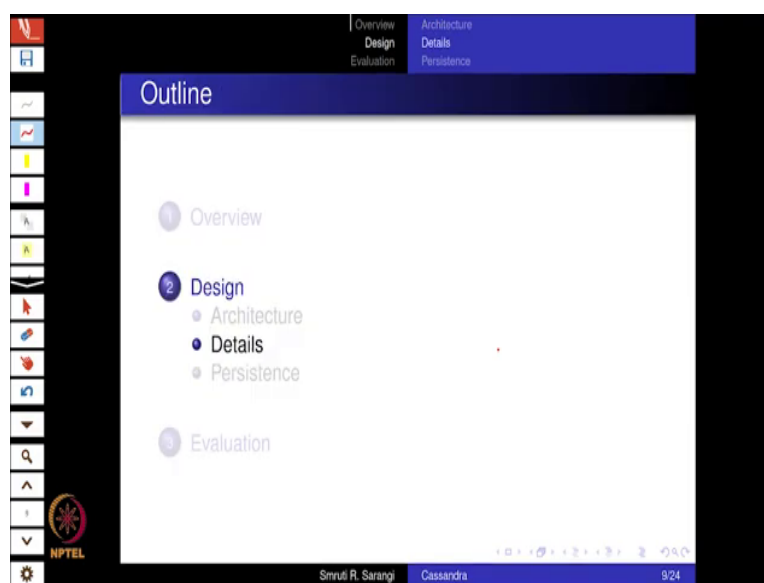
And so similar to chord, what we do is that given a key, we map it to the hash space, then we traverse the hash space clockwise till we find the nearest node. And then of course, we maintain replicas of the key value in the next n nodes subject to the fact that they are in different data centers.

So, then, of course, similar to chord we determine the replicas that contain the key. And among the replicas, we define the write quorum, and then we define the read quorum. So, write quorum, if you recall, let us say that there are 10 replicas, and the size of the quorum is W size

of the read quorum is  $R$ , then we had always said that  $r + w > 10$ . And which will basically ensure that there is at least an overlap between the read and write quorums. And so, in this case, it will mean that we will not miss an update.

So, for a read either we read from the entire read quorum and somehow merge what we were doing in Dynamo. But of course, Cassandra does not discuss the merging aspect. Or we read from the closest replica, the geographically closest replica, or the one that we find the one that is available.

(Refer Slide Time: 14:59)



The details, so, well, as we discussed, we partition data using a consistent hashing algorithm like chord. And chord, if you would recall was using virtual nodes for load balancing, we do



the same over here. So, this kind of ensures that the load per physical node is kind of balanced. And there were some formulae in code that discussed virtual nodes using some results.

So, I would refer you to this paper. So, in this case, each data item is replicated on  $N$  hosts, along the thing, of course, if these are the end hosts, but of course, we will skip a few subject to different to the replication policy. The first is it will be it will be rack unaware. So, in the rack unaware policy, we just take the end successors.

But in rack aware, we skip those successors that are in the same rack. So, at least they will be across different racks of the same data center. What is a rack, well a rack is like cabinet of servers. So, it is possible that let us see the network connection to a rack can get cut. In this case, the entire rack will go offline.

So, a rack aware replication policy will ensure that at least the replicas are in different racks. And what Facebook actually does is that it is more of data center aware in this case, you want at least in the ideal case, all  $n$  of the replicas to be in  $n$  different data centers. If that is not possible, you would at least like to maximize the span of data such that it is in different data centers.

So, just in case an entire data center gets cut off, we can Facebook will still continue to work. Cassandra will still continue to work. So, that is the broad idea over here. So the Cassandra what it does, is among the set of replicas, so similar to Dynamo they also call this the preference list. And so, Cassandra, first elects a leader using the Apache zookeeper service and Apache zookeeper service is a distributed coordination engine.

The job of the distributed coordination engine is to ensure is to provide certain key services, certain basic services. And the basic services include electing a leader and we will see it provides an implementation of a few more distributed algorithms as well. So, you do not elect a leader or do the preference list but you elected leader of all the nodes. So, leader will assign the replicas to the nodes. Along with that the meta data routing tables and topology information are maintained, of course at each node for ought to work, and also centrally by a zookeeper system to initiate repair and maintenance messages of the array. So, party Zookeeper, I would encourage viewers to take a look at this open source project, which has a wide array of distributed coordination algorithms.

(Refer Slide Time: 18:37)

The screenshot shows a presentation slide titled "Membership" with the following content:

- Uses the Scuttlebutt anti-entropy based gossip system to propagate membership information.
- Failure detection
  - $\Phi$  Accrual failure detector *probabilistic failure detector*
  - The failure detector gives a value,  $\Phi$ , for each node that is not a Boolean number. It represents the level of suspicion for the node.
  - If  $\Phi = 1$ , the chance of an error is 10%
  - If  $\Phi = 2$ , the chance of an error is 1%
  - Every node maintains a *sliding window* of inter-arrival times of gossip messages
  - $\Phi$  is calculated based on the distribution of these inter-arrival times.

The slide also features a navigation sidebar on the left, a top navigation bar with "Overview", "Design", "Evaluation", "Architecture", "Details", and "Persistence", and a footer with "Sriruti R. Sarangi", "Cassandra", and "11/24".

So, for membership well, so membership information is typically gossiped using a gossip protocol. And gossip protocols were discussed in the second or third lecture of this course. So, in this case, Cassandra uses the scuttlebutt anti-entropy-based gossip system to propagate membership information. And this is propagated relatively quickly.

And failure detection is a problem. The reason it is a problem is that we want to detect failures quickly. But to detect failures quickly, we need to actually send a lot of messages. So, what Cassandra does is that it uses a probabilistic failure detector for the accrual failure detector which kind of gives a confidence bound on the failure.

So, it gives a value  $\phi$  for each node that is not a Boolean number, which represents a level of suspicion level of suspicion that the node has actually failed. If  $\phi = 1$  the chance of an error in whatever is being predicted is 10 %. Which in this case, if you see it has failed and  $\phi = 1$  well the chance of mice being wrong is 10 %.

If  $\phi = 2$  it is 1 %. If  $\phi = 3$  it is 0.1 percent and so on. It varies exponentially, and on the log scale, every node maintains a sliding window of inter arrival times have gossip messages. And based on this sliding window of inter arrival times of gossip messages, Phi is calculated, and weavers are welcome to take a look at the details of the actual failure detected on the web.

(Refer Slide Time: 20:37)

Overview | Architecture  
Design | Details  
Evaluation | Persistence

## Bootstrapping

- When a node starts, it chooses a random **token**. The token defines its position in the ring.
- It saves the token locally and also gives it to Zookeeper.
- It then gets a list of **contact points** from Zookeeper, and contacts them such that it can join the ring.
- The **membership information** for the new ring is then gossiped in the cluster. *gossiped.*
- There might be multiple Cassandra instances running (for different types of services)
  - Tag each message with the **instance id**
  - Administrators can manually manage Cassandra instances

NPTEL | Sriniv R. Sarangi | Cassandra | 12/24

So, bootstrapping, how does the system starts well, when a node starts, it chooses a random token. So, we have used the same terminology as Dynamo. So, the token is the hash value of its IP address for example, it defines its position in the ring. The token is stored locally and also it is given to zookeeper which is acting as a global coordinator. Zookeeper in return gives it a set of contact points and known set of servers which are already on the ring.

One of the servers while other contexts, servers can be contacted, so that it will help the node join the ring. And now given that a new node has joined the DHT ring, the membership information is gossip. So, of course, there might be multiple Cassandra instances running for different types of services. Cassandra is a generic system. So, you can tag each message with the Instance ID of which Cassandra message it is. And furthermore, administrators are given a console to manually manage the instances and manually manage the process of node addition and node deletion.

(Refer Slide Time: 22:02)

The screenshot shows a presentation slide titled "Scaling" with a blue header. The slide content includes handwritten notes and a bulleted list:

- DHT ⇒ scalability*
- physical node*
- high load*
- The token assignment can be made more **intelligent**. Assign a node a position in the ring such that it can take some load off a heavily loaded node.
- Use kernel-kernel copy techniques to stream data from the old to the new nodes.
  - Parallelize this process by streaming data from multiple replicas
  - Similar to BitTorrent

At the bottom of the slide, there are handwritten sketches of a ring structure with nodes and arrows, and the text *node #11*. The slide footer includes the NPTEL logo, the name "Smruti R. Sarangi", the title "Cassandra", and the number "13/24".

Does the system scale well? Yes, of course. So, using a DHT was the best thing to do in this case, because the DHT automatically implies huge scalability. And of course, we can make so what will stop what is the main problem with DHT? Well, the main problem with a DHT is that it is possible that load balancing maybe non-uniform. So, for a given physical node, we might have a high load.

So, this is kind of taken care of the virtual nodes. This is kind of taken care of by changing the replica placement strategy. Also, the token unlike Dynamo will not assign the physical node to different points in the ring. But the initial token assignment, of course, can be made much more intelligent. In a sense, a node can be assigned a position on the ring such that it is beside a heavily loaded node, such that it can take away some of its load.

So, different combinations are possible. So, it is a job of Apache zookeeper along with a customized Cassandra engine to ensure that this indeed does happen. Once that has happened, once a node has been placed on the ring, so it will place it needs to get data. So, for all the rest of the nodes, they will be storing, they will be the owner of some key and recall that the owner is a clockwise successor of the key.

So, it will basically get replicas from a few of the other nodes which are on the ring. So, replica transfer will be a fast kernel kind of transfer, where data is streamed into the new nodes. This process, of course, can be personalized. Similar to BitTorrent in the sense, we can take the replica data, break it into small chunks. And if let us say the node has multiple network ports, then via each network port, we can sort of streaming the data in parallel, similar to what BitTorrent actually does.

So, this is again, kind of like at a lower level, it is kind of lower level in the sense, it is more at the network layer. And so, so it is very started really at the network layer, but it is pretty much at a level that is lower than the DHT. So, in this case, what essentially needs to be done is that if, let us say, there is this node over here that we just added, and another node wants to send a replica to this node.

So, clearly, there are many other nodes that also have the same replica. So, all of them can sort of divide the replica into small chunks and send the chunks to this node in parallel. So, this will speed up the process significantly if the node has multiple network interfaces and can absorb the parallel transfers.

(Refer Slide Time: 25:09)

The screenshot shows a presentation slide titled "Outline". The slide has a blue header with the title "Outline". Below the header, there is a list of sections: "1 Overview", "2 Design", and "3 Evaluation". Under "2 Design", there are three sub-items: "Architecture", "Details", and "Persistence". The "Persistence" item is highlighted with a blue dot. The slide is part of a presentation by Snruti R. Sarangi on Cassandra, slide 14/24.

The screenshot shows a presentation slide titled "Local Persistence". The slide contains a list of bullet points with handwritten annotations in red ink:

- Relies on the local file system for data persistence.
- Typical write in a fault tolerant file system:
  - Write to a **commit log** (with a handwritten arrow pointing to "journaling")
  - Update an in-memory data structure
- Dedicated disk for the commit log.
- Once the in-memory data structure **exceeds** a certain size, dump its contents to the commit log (with a handwritten arrow pointing to "file system").
- Writes are sequential to the disk. Simultaneously an **index** for disk data is created. (with handwritten note "index (database index)")
- Different files in the disk are later merged (compaction). (with a handwritten diagram showing a box splitting into two smaller boxes and a vertical bar labeled "temporal")

The slide is part of a presentation by Snruti R. Sarangi on Cassandra, slide 15/24.

So, outline a little bit of persistence will persistence basically means storing the results to disk such that if there is a memory crash, the data does not go away. So, this durable storage, so we rely on the local file system for data persistence. So, recall that this is a, this is a journaling file system, I will not describe what a journaling file system is, the expectation is that viewers will go to the web and read about it.

So, a typical and the fault tolerant file system, which uses journaling is that we write to a commit logs instead of directly writing it to the file on the disk, we write it to a commit log. This is used to update an in-memory data structure. Furthermore, it is possible that this commit log will be regularly be persisted to the disk. And we can have a dedicated disk or a dedicated non-volatile memory for the commit log.

And once this in memory data structure exceeds a certain size, we simply dump the contents to the commit log, which means that from the commit log, we actually change the underlying file system, so the commit log is kind of like a cache in the file system. And it is ordered sequentially.

So, writes are sequentially made to the disk. And, and simultaneously, what we do is that we create or update an index for this data. So, such that if we need some data on a disk, since we are we do not have a regular file system that is organized into directories, we need to create a special index similar to a database index. So, again, a database index, I am not explaining over here, because I am assuming that readers would have viewers will either already know or will kind of look it up.

So, the important point over here is that even if we have lots of data on the disk, it is important to have an index, it is important to index it. And it is important to create a database like index for doing that. And once we have a lot of files on the disk, Cassandra actually has files a variable size. So, the recent files kind of have like small sizes. And gradually, what happens is that older files are much larger sizes.

So, they are kind of merged. So, we do not need them, we create large files, and we kind of keep them. So, this is mainly because inbox messages show a high degree of temporal locality. And because of the high degree of temporal locality, well, most of the time, we will get them in the in-memory structure. If not, we will get them on some small files and the disk and sell them do we have to actually access the large files.

(Refer Slide Time: 28:11)

**Read Operation**

- First query the in-memory data structure.
- Search the files from newest to oldest.
- To speed up searches we can use a **Bloom filter**.
  - Can quickly indicate that a set of files do not contain a key
  - It is often difficult to locate all the columns in a row.
    - Solution: Use a 256 KB index

Handwritten notes and diagram:  
- A circle labeled 'S' has an arrow pointing to 'y/yes' and another to 'x false negatives'.  
- A box labeled 'Set membership' has an arrow pointing to 'y/yes'.  
- A box labeled 'BF' (Bloom Filter) has an arrow pointing to 'in-mem' (in-memory).  
- A box labeled 'Key' has an arrow pointing to 'row'.  
- A box labeled 'row' has an arrow pointing to 'BF'.  
- A box labeled 'in-mem' has an arrow pointing to 'BF'.  
- A box labeled '256 KB' has an arrow pointing to 'BF'.  
- A box labeled '256 KB' has an arrow pointing to 'row'.

So, what is what does the read operation look like? We first query the in-memory data structures. We search the files from newest to oldest. And of course, we can speed up this search process by using a Bloom filter. Again, I am not explaining but the Bloom filter, let me say, is this specific kind of a data structure that essentially indicates set membership.

So, given a set of values, and now if it is presented with and this is also listed as a set, it is a set of values. So, now it is presented with a new value, it will essentially indicate a yes or no answer. Yes, means that it is there in the set and no means that it is not there. So, bloom filter. If it says yes, it is there, it does not mean that if it is actually there, it might be there it might not be there.

So, essentially it allows for a false positive. However, if it says that the data item is not there means it is definitely not there, which means that a Bloom filter does not have any false negatives. So, false positives are there but false negatives are not there. And how exactly a Bloom filter works well, again, this is a exercise to the readers if they do not already know. Just go to Wikipedia and read about Bloom filters.

So, the advantage is that the Bloom filter can very quickly indicate that a set of files do not contain a key, which is the key for the row. Alright, so why are we doing all of this? Well, so the reason why we are doing all of this is given a key, we want to fetch the contents of the row.

So, this is required for both a get operation and an update operation. And so, for either, you know, it is like get and put in a dynamo, either for getting data or updating data. And so that is

the reason we first access the in-memory data store, and if that is not there, then we access the disk. And the disk, of course, has a range of files of different sizes.

The recent entries are stored in smaller files, and later entries are stored in very large files, hoping that they will, we will have temporal locality. And we will typically not access later entries, and just in case we are users will be willing to tolerate the additional latency. And to speed up the search, well, how do we organize it, we have to create some kind of an index.

So, we can go through a Bloom filter, which will tell you quickly whether it even makes sense to search here or not. If it is not there, we directly go to the disk. And furthermore, a row itself can have a lot of columns. So, what we do is we use a 256-kilobyte column index in the row, which allows us to exactly find out the byte offset of a given column.



(Refer Slide Time: 31:34)

The slide is titled "Implementation Details" and is part of a presentation on Cassandra. It features a navigation bar at the top with "Overview", "Design", "Evaluation", "Architecture", "Details", and "Persistence". The main content area is titled "Main Modules" and lists: "Partitioning Module, Cluster membership module, Failure detection module, Storage engine module". Below this, a list of implementation points is shown: "Written in Java using non-blocking I/O primitives", "Control messages use UDP", and "All application related messages use TCP". A handwritten note in red ink next to the third point says "do not want to lose data". The slide footer includes the NPTEL logo, the presenter's name "Srnuti R. Sarangi", the topic "Cassandra", and the slide number "17/24".

The implementation details well, we need to look, we need to consider several things we need to consider the partitioning module for partitioning the keys and the replicas a clustered membership how do we take care of membership failure detection and storage. So, the entire system has been written in Java using non-blocking I, O primitives. The control messages to manage Cassandra use UDP and all the application related messages use TCP Of course, this is required. Because this is where we need reliable communication, we do not want to lose data.

(Refer Slide Time: 32:22)

The slide is titled "Detailed: List of steps" and is part of the same presentation on Cassandra. It features the same navigation bar as the previous slide. The main content area lists five steps: "1 Identify the nodes that own the data for a key", "2 Route the requests to the nodes and wait for the responses", "3 If the replies do not arrive within a certain time, fail the request (suspect a failure)", "4 Figure out the latest version based on the timestamp", and "5 Trigger a repair (if required) and return to the client". A handwritten red checkmark is next to step 2, and a red arrow points to step 3. The slide footer includes the NPTEL logo, the presenter's name "Srnuti R. Sarangi", the topic "Cassandra", and the slide number "18/24".

So, we identify the nodes that own the data for a key, we route the request to the nodes and wait for the responses. If of course the replies do not arrive within a certain time the request

fails. And then of course, we tend to suspect a failure. So, the failure approval mechanism that we use will pretty much take such inputs, but of course, it uses far more inputs also like the inter arrival time of WhatsApp messages and so on.

However, if we do get a request, we find out a latest version based on the timestamp and just in case the repair is needed or emerge as needed, well, then we again, return the values to the client. But of course, the Cassandra paper as such, does not broach this issue in great detail. Hence, I will also not try to speculate on how exactly the versioning is done in Cassandra. Because this aspect is not really well covered in the paper.

(Refer Slide Time: 33:37)

**Journaling File System**

- Uses a **rolling** commit log
  - Create a new one after the old one has reached a certain size
  - Use a size of **128 MB** [ ... ]
- In **Cassandra** : Use an in-memory data structure, and a backing data file for each column family
- Every time the in-memory data structure is dumped to the disk, we set a bit in its commit log
- Each commit log maintains a bit vector (corresponding to disk dumps )
- We can delete commit logs for entries that have been **per-sisted**

So, it uses a rolling commit log. Well what does that mean? I mean, it is essentially a large log where data can be added and removed at any point in time. So, once a given commit log has reached a certain size, we create a new one and the old one is persisted to disk and the size threshold that is used is 128 megabytes. So, this recall that is an old paper roughly 2010 so, those days even 128 megabytes because you know we would like ideally like to keep the commit log in memory. And those days 128 megabytes was quite a bit of memory for an off the shelf processor.

So, it was 128 MB was not quite a bit of memory, but of course, recall that we are many virtual nodes and then we are storing the return for each virtual node. So, and also the additional overhead. So, keeping all of that in mind, 128 MB was the maximum that designers thought they could afford. So, I stand corrected in my previous statement.

So, in Cassandra, we use an in-memory data structure that we have seen and a backing data file for each column family. Every time the in-memory data structure is dumped to the disk, we set a bit in its commit log. And this basically says that the data has been persisted. And each commit log will maintain a bit vector corresponding to the disk dumps that have happened, which tells you what has been persisted, and what has not been persisted. And for entries that have been persisted, we can kind of delete the logs. And we do not need them. So, that will free up space for us.

(Refer Slide Time: 35:32)

The write operation

- Write operation
  - Normal mode : Unbuffered writes
  - Fast sync mode : All writes to the commit log, and the data file are buffered. *Can lose data on a disk crash*
  - Sequentialize writes to the disk. *machine crash*
    - Fast
    - Does not require costly B-tree updates

So, write operations are again of 2 types, we can have a normal, write which is unbuffered, which means that we directly persist and or we can have a fast sync mode. In fast sync modes, all the rights to the commit log and the data file are essentially buffered in memory. So, this makes the rights kind of fast, but just in case there is a crash, we will lose data. And so that is a problem.

So, we sequentialized writes to the disk, of course, which is a fast, fast operation. And unlike databases, in this case, our data structures are very simple. They are essentially like linked lists. So, we do not have costly B-tree updates. I need to make a small correction here. This is actually not a disk crash, it as a machine crash. And the machine crash is something that wipes off the state of main memory. So, that is the reason. If the machine has crashed, its main memory is gone. So, but the data on the disk still remains. And even if he ca not get the machine up, you can always move the disk from one machine to the other.

(Refer Slide Time: 36:51)

The screenshot shows a presentation slide titled "Cassandra Index". The slide content includes:

- All the data is indexed based on the primary key
- The data file on the disk is broken down into a sequence of blocks. Each block contains at most 128 keys.
- A block is demarcated by the **block index**, which captures the relative offset of a key within a block, and contains the size of its data.
- This index is also saved in memory.
- Search process:
  - First search the index in memory.
  - If it is not in memory, search index in the disk.
  - Locate the key's value in the in-memory data structure.
  - If not search in the data file.

Hand-drawn in red on the slide is a diagram of a block. It consists of two stacked rectangular boxes. A bracket on the right side of the top box is labeled "128 keys". An arrow points from the bottom of the top box to the word "value" written below it.

The slide footer includes the NPTEL logo, the name "Sruvi R. Sarangi", the title "Cassandra", and the page number "21/24".

So, let us now look at the structure of a Cassandra index. So, all the data, of course is indexed based on the primary key. So, the primary key is the index of the row. And within the row, as I said, we have columns, families of columns, super columns. So, the data file on the disk is broken down into a sequence of blocks. So, basically, what we do is if you consider the data file on the disk, we break it down into a sequence of blocks, and each block has 128 keys.

So, a block, of course, is demarcated by the block index, which captures the relative offset of a key within a block and contains the size of its data. So, basically, what we see over here is that the search process would look like this that we search, first search the index and memory. If it is not in memory, we search the index on the disk.

So, given a key, we can quickly locate it in within the index. And using this structure, of course, here we can find the block that contains the key. And from there, we locate the value of the key so, every key will have a pointer to its value. So, value is the location where the data for all of its columns is stored. So, again, the value we search for the value in memory data structure, if we do not find it, we then we search for the value in that data file. So, of course, most indices are stored in memory for quick access. Because in this case, we would like to minimize disk accesses as much and as far as possible.

(Refer Slide Time: 38:38)

The slide is titled "Arrangement of Files" and is part of a presentation on Cassandra. It features a list of four bullet points explaining file ordering and access patterns. Handwritten red text above the list reads "newest → oldest". The slide also includes a navigation sidebar on the left and a footer with the presenter's name "Snruti R. Sarangi" and the slide number "22/24".

- Files are ordered according to their time of creation.
- We start by accessing the latest file (since we want the latest update)
- Because of temporal locality of accesses, we keep recent data in small data files
- Gradually we merge them and create bigger merged files

The arrangement of files of course, files are ordered according to the time of creation. So, I start with traverse the file from the newest files from the newest to the oldest. So, since you are on the latest update, we access the latest file and because of temporal locality of accesses, we as I said we keep recent data in small files. And gradually we keep on merging the older files to create bigger and bigger files, which of course take more time to access. But we are also assuming that accesses to such files will be kind of infrequent.

(Refer Slide Time: 39:18)

The slide is titled "Lessons Learnt" and is part of a presentation on Cassandra. It features a list of four bullet points detailing lessons learned from Map-reduce jobs. Handwritten red text includes "Average failure detection time: 15 seconds" and a list of Zookeeper functions: "leader election", "replica management", and "load balancing". The slide also includes a navigation sidebar on the left and a footer with the presenter's name "Snruti R. Sarangi" and the slide number "23/24".

- Map-reduce jobs process raw data (or from a MySQL database) and send data to the Cassandra instance.
- Implemented: atomic operation per key per replica
- Average failure detection time: 15 seconds
- The usage of Zookeeper as a co-ordination service is very useful.

• leader election  
• replica management  
• load balancing

So, what are the lessons learned by the developers in this. So, of course, they use Map-reduce jobs to process all the raw data or from a MySQL database and send it to the Cassandra

instance. They did implement atomic operations per key per replica that was not hard to implement. And the average failure detection time was actually high, it was order of several minutes with a regular failure detector.

But with the accrual-based failure detection mechanism, they were able to bring it down to 15 seconds. And the usage of zookeeper as a coordination service was very useful. So, in general, in a distributed system, at least for the leader election, replica management and load balancing, so let us write it down. You will need some kind of a distributed coordination service.

(Refer Slide Time: 40:47)

The slide is titled "Facebook Inbox Search" and contains the following bullet points:

- Maintains a per user index of all messages sent and received
- Search features
  - (a) Term search
  - (b) Search by the name of the recipient
- For (a) *Term Search*
  - User id is the key, and the message words are the super column. Message ids are columns within the super column.
- For (b)
  - User id is the key and the recipient's ids are the super columns. For each super column, the message id is the column.
- Uses index prefetching

Handwritten annotations include:

- A diagram for (a) showing a tree structure with "key" at the root, "id" as a child, and "id" as a child of "id".
- A diagram for (b) showing a tree structure with "key" at the root, "id" as a child, and "id" as a child of "id".
- A diagram for (c) showing a tree structure with "key" at the root, "id" as a child, and "id" as a child of "id".

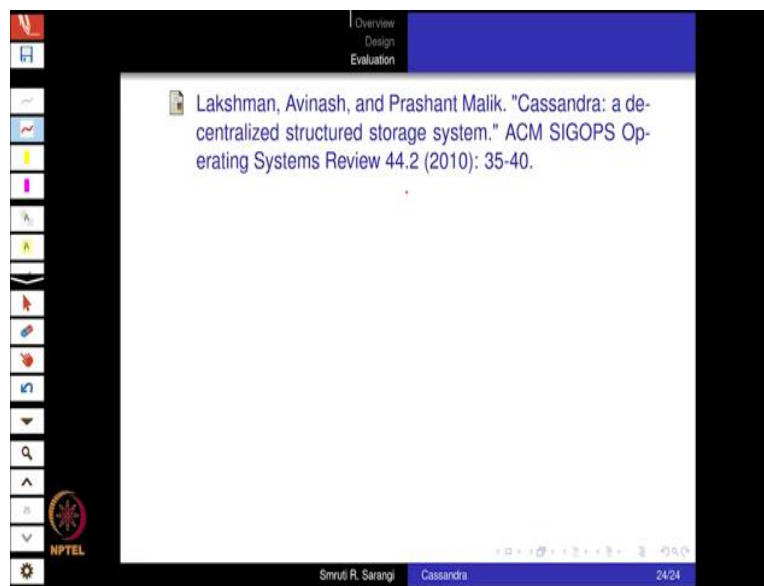
The slide footer includes "Snruti R. Sarangi", "Cassandra", and "24/24".

As the inbox search, so we maintain a per user index messages sent and received. So, we can search it as we had said in the first slide based on the terms and based on the name of the recipient. If we are to if you are doing term search, term based search. In this case, the user id is the key. And the message words are the super column.

The different message was a different message words are the super columns. And for each and then the message id is that contain that term are the columns. Similarly, if let us say I were to search for the recipient, well, the key would be the root kind of the identifier. The recipients ids are the super columns.

And for each super column, the message ids are the columns. And to speed up the execution, we use index prefetching, which means that we also prefetch the indexes and keep them in memory for very, very fast access, because we need to locate the key and given the key we need to quickly fetch its value.

(Refer Slide Time: 42:03)



So, this was a 2010 paper published by 2 Facebook engineers. And so the Cassandra system, henceforth has been used very heavily. It is very widely documented, as well, in different papers and in the technical press. So, I would ask the viewers to look at other sources of information and documentation on Cassandra.