

Distributed Systems

CS425/ECE428

Instructor: Radhika Mittal

Acknowledgements for some of the materials: Indy Gupta

Logistics

- HW4 is due today.
- HW5 has been released.
 - you should be able to solve first two questions right-away
 - you should be able to solve last two questions in a week.

Our agenda for the next 3-4 classes

- Brief overview of key-value stores
- Distributed Hash Tables
 - Peer-to-peer protocol for efficient insertion and retrieval of key-value pairs.
- Key-value stores in the cloud
 - How to run large-scale distributed computations over key-value stores?
 - Map-Reduce Programming Abstraction
 - How to schedule jobs in the cloud?
 - How to design a large-scale distributed key-value store?
 - Case-study: Facebook's Cassandra

Recap

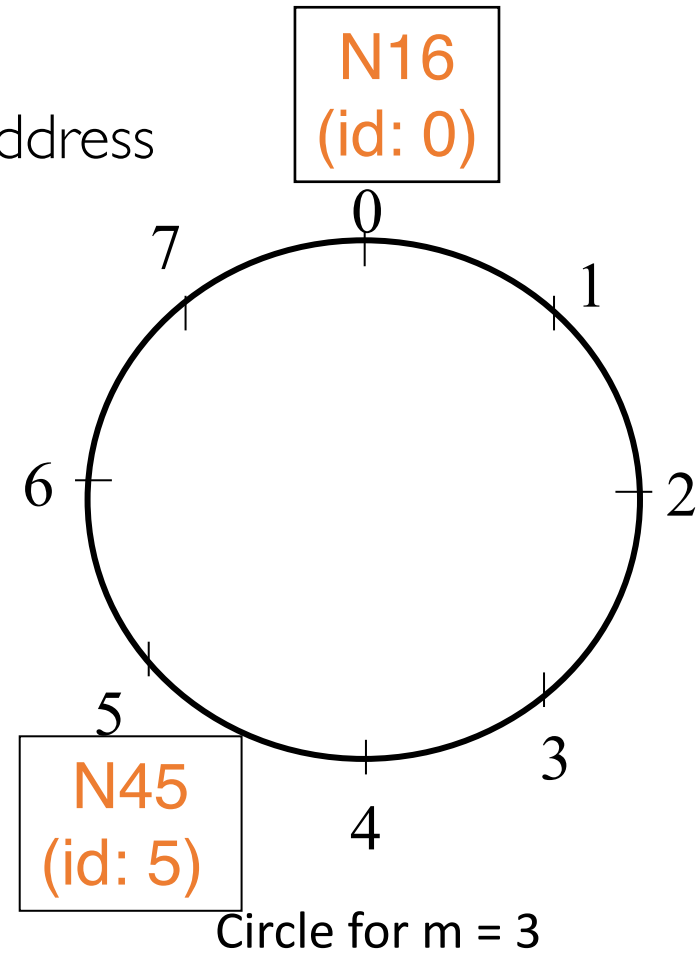
- Distributed Hash Tables
 - Peer-to-peer protocol for efficient insertion and retrieval of key-value pairs.
 - Other required properties: load balancing, fault tolerance.
- Case-study: Chord

Recap: Chord

- Uses **consistent hashing** to map nodes on a ring with m -bits identifiers.
- Uses consistent hashing to map a key to a node.
 - stored at **successor(key)**
- Each node maintains a **finger table** with m fingers.
 - With high probability, results in $O(\log N)$ hops for a look-up.

Recap: Chord Consistent Hashing

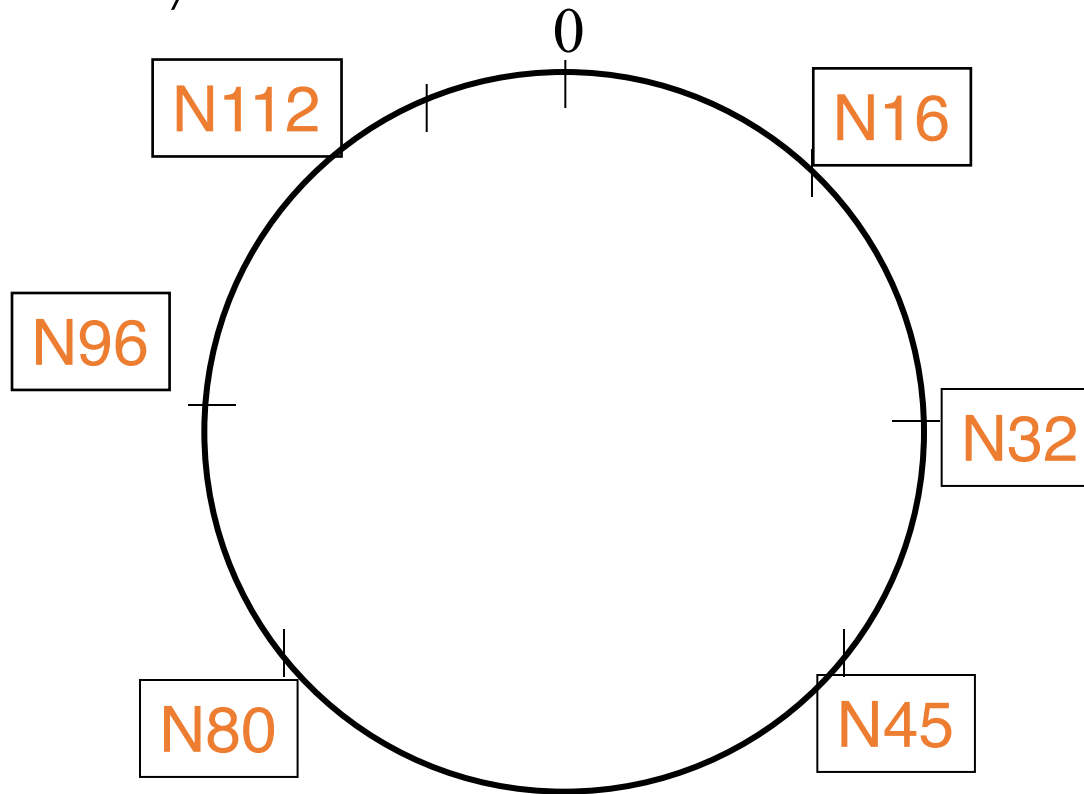
- Uses *Consistent Hashing* on node's (peer's) address
 - $\text{SHA-1}(\text{ip_address, port}) \rightarrow 160$ bit string
 - Truncated to m bits (modulo 2^m)
 - Called peer id (number between 0 and $2^m - 1$)
 - m chosen such that negligible chance of id conflicts
 - Can then map peers to one of 2^m logical points on a circle



Where will N45 be placed on this circle?

Ring of Peers: Running Example

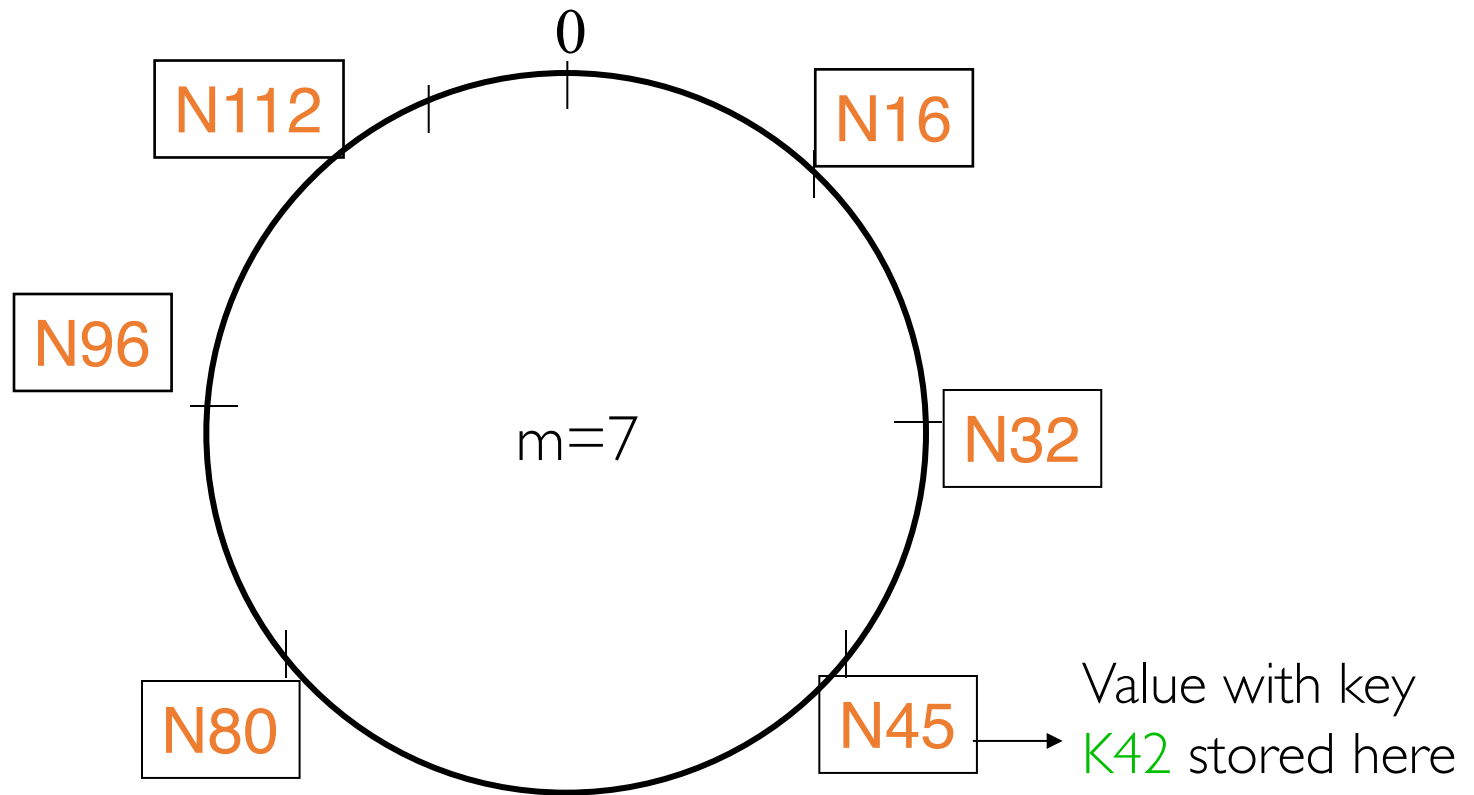
- Say $m=7$ (128 possible points on the circle – not shown)
- 6 nodes in the system.



Recap: Mapping Keys to Nodes

- Use the same consistent hash function
 - $\text{SHA-1}(\text{key}) \rightarrow 160$ bit string (key identifier)
 - Henceforth, we refer to $\text{SHA-1}(\text{key})$ as *key*.
 - The key-value pair stored at the key's *successor* node.
 - $\text{successor}(\text{key}) = \text{first peer with id greater than or equal to } (\text{key mod } 2^m)$
 - *Cross-over the ring when you reach the end.*
 - $0 < 1 < 2 < 3 \dots \dots < 127 < 0$ (for $m=7$)
- Consistent Hashing \Rightarrow with K keys and N peers, each peer stores $O(K/N)$ keys. (i.e., $< c.K/N$, for some constant c)

Recap: Ring of Peers: Running Example



Where will the value with key 42 be stored?

Recap: Performing Lookups

- Option 1: Each node is aware of (can route to) any other node in the system.
 - Need a very large routing table.
 - Poor scalability with 1000s of nodes.
 - Any node failure and join will require a *necessary* update at all nodes.
- Option 2: Each node is aware of only its ring successor.
 - $O(N)$ lookup. Not very efficient.
- Chord chooses a sweet middle-ground.

Recap: Performing Lookups

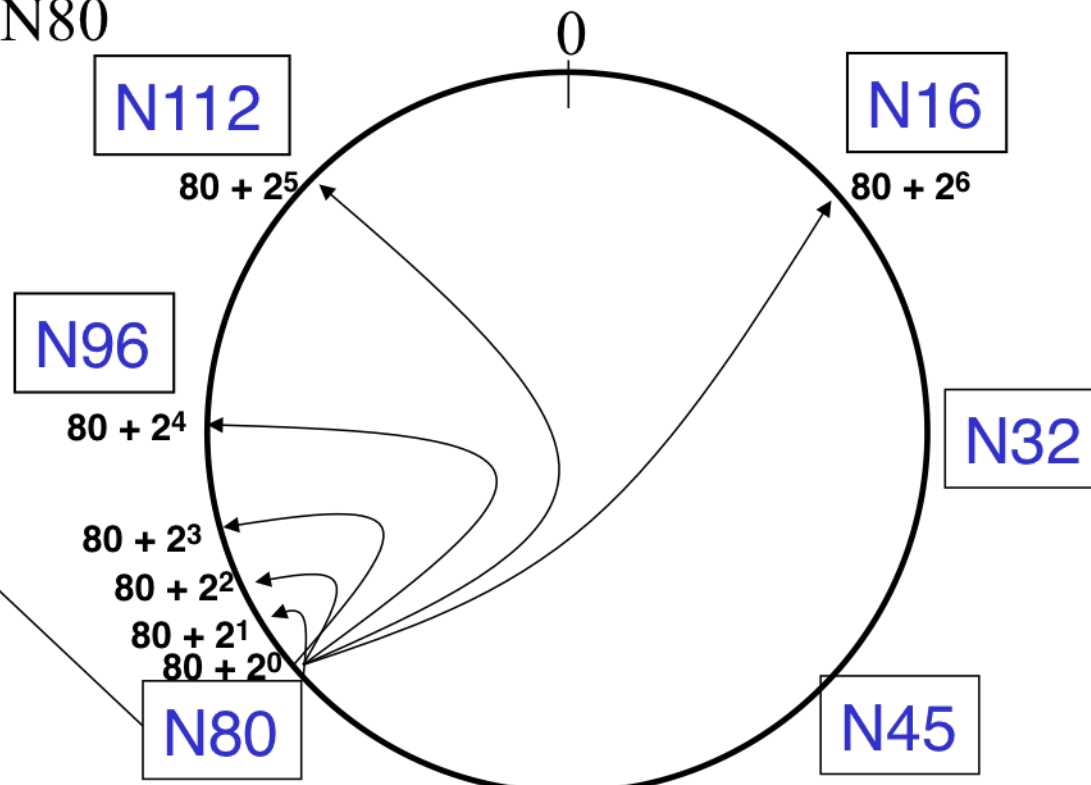
- Chord chooses a sweet middle-ground.
 - Each node is aware of $\sim m$ other nodes.
 - Maintains a *finger table* with m entries.
 - The i th entry of node n 's finger table = $\text{successor}(n + 2^i)$
 - i ranges from 0 to $m-1$

Recap: Finger Tables

Say $m=7$

Finger Table at N80

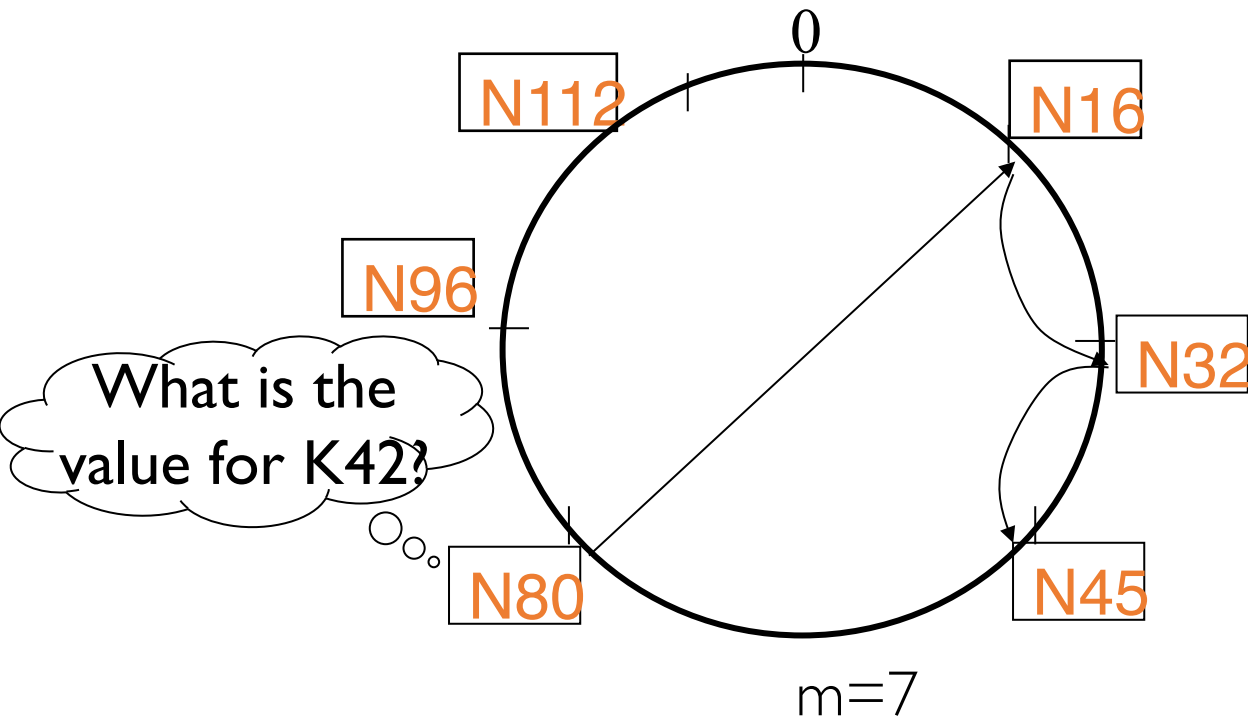
i	$ft[i]$
0	96
1	96
2	96
3	96
4	96
5	112
6	16



i th entry of node n 's finger table = $\text{successor}(n + 2^i)$

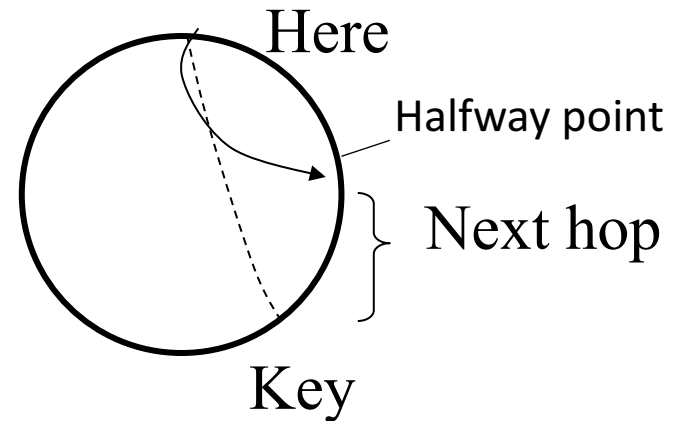
Recap: Search for key k at node n

At node n , if k lies in range $(n, \text{next}(n)]$, where $\text{next}(n)$ is n 's *ring successor* then $\text{next}(n) = \text{successor}(key)$. Send query to $\text{next}(n)$
Else, send query for k to largest finger entry $\leq k$



Recap: Analysis

Search takes $O(\log(N))$ time



Proof Intuition:

- (intuition): at each step, distance between query and peer-with-file reduces by a factor of at least 2 (why?)
- (intuition): after $\log(N)$ forwardings, distance to key is at most $2^m / 2^{\log(N)} = 2^m / N$
- Expected number of node identifiers in a range of $2^m / N$:
 - ideally one
 - $O(\log(N))$ with high probability (by properties of consistent hashing)

So using ring successors in that range will use another $O(\log(N))$ hops. Overall lookup time stays $O(\log(N))$.

Recap: Chord

- Uses **consistent hashing** to map nodes on a ring with m -bits identifiers.
- Uses consistent hashing to map a key to a node.
 - stored at **successor(key)**
- Each node maintains a **finger table** with m fingers.
 - With high probability, results in $O(\log N)$ hops for a look-up.
 - $O(\log(N))$ hops true only if finger and successor entries correct.
 - What happens when nodes fails or when new nodes join in?
 - Our focus today.

Recap: Chord

- Handling node failures:
 - For lookups:
 - maintain r multiple ring successor entries
 - in case of failure, use another successor entries
 - if r is sufficiently large ($\log N$), can handle a high rate of failures (50%) with a very high probability.
 - When storing keys:
 - replicate key-value at r successors and predecessors

Need to deal with dynamic changes

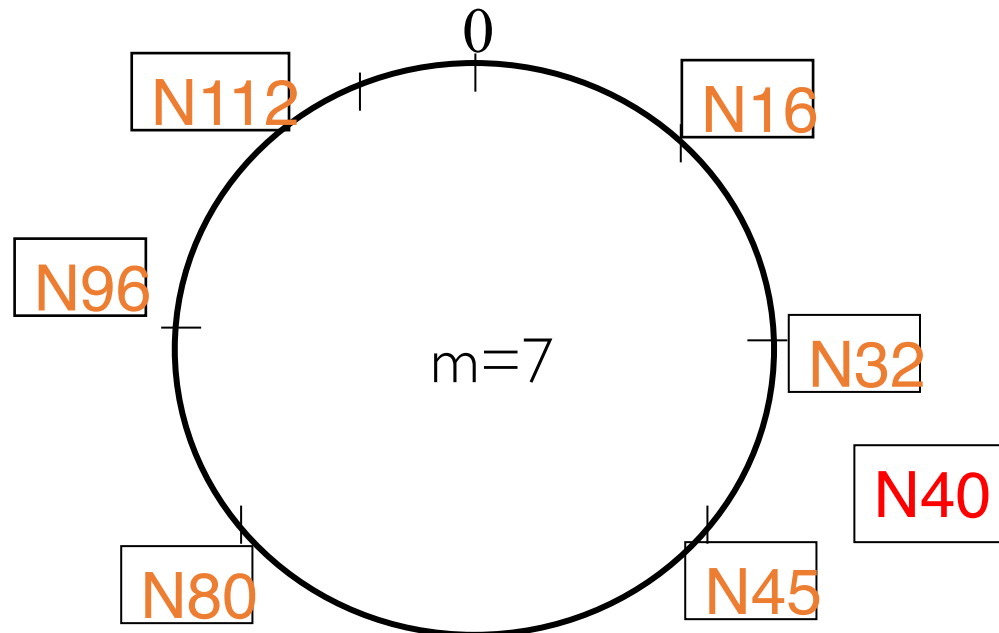
- Nodes fail
- New nodes join
- Nodes leave

So, all the time, need to:

→ Need to update successors and fingers, and copy keys

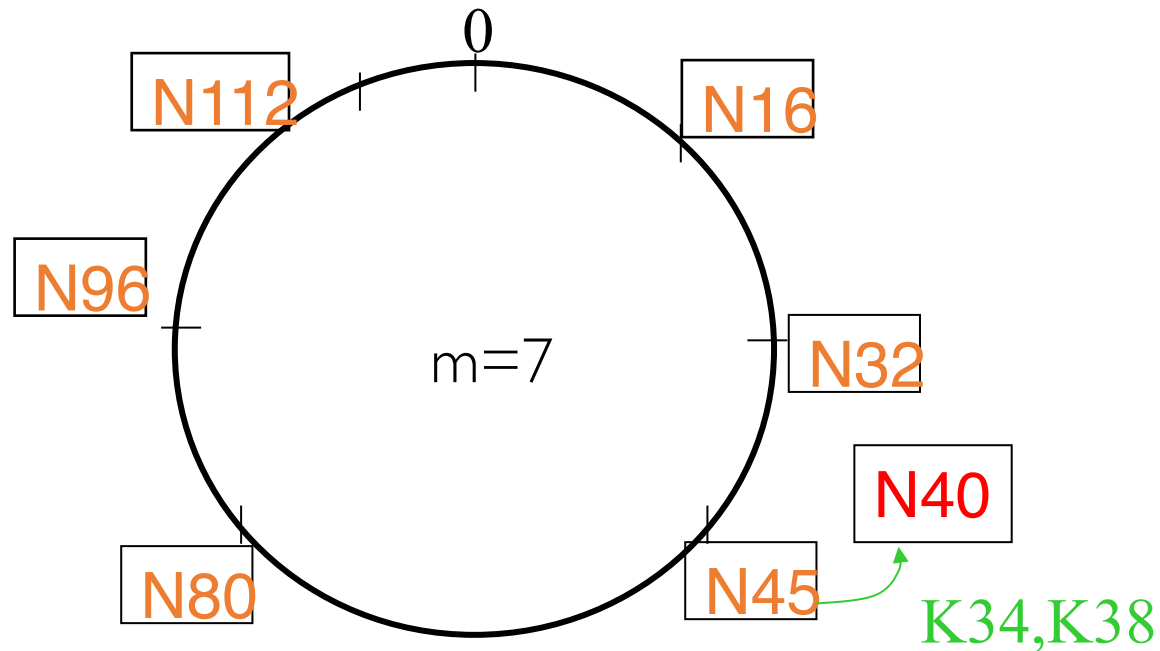
New node joins

New node contacts an existing Chord node (introducer).
Introducer directs N40 to N45 (and N32).
N32 updates its ring successor to N40.
N40 initializes its ring successor to N45, and initializes its finger table.
Other nodes also update their finger table.



New node joins

N40 may need to copy some files/keys from N45
(files with fileid between 32 and 40)



Concurrent Joins

- Aggressively maintaining and updating finger tables each time a node joins can be difficult under high *churn*.
 - E.g. when new nodes are concurrently added.
- Correctness of lookup does not require all nodes to have fully “correct” finger table entries.
- Need two invariants:
 - Each node, n , correctly maintains its ring successor ($next(n)$)
 - First entry in the finger table.
 - The node, $successor(k)$, is responsible for key k .

Stabilization Protocol

- When a node n joins (via an introducer)
 - initialize $\text{next}(n)$, i.e. the ring successor
 - notify $\text{next}(n)$.
- When node n gets notified by node n' :
 - // update $\text{prev}(n)$, i.e. the ring predecessor of n
 - if ($\text{prev}(n) == \text{nil}$ or n' is in $(\text{prev}(n), n)$)
 - $\text{prev}(n) = n'$.

Stabilization Protocol (contd)

- Each node n will periodically run stabilization:
 - $x = \text{prev}(\text{next}(n))$
 - if x in $(n, \text{next}(n))$, then $\text{next}(n) = x$.
 - notify $\text{next}(n)$.
- Each node n periodically updates a random finger entry.
 - Pick a random i in $[0, m-1]$
 - Lookup $\text{successor}(n + 2^i)$

New node joins

New node contacts an existing Chord node (introducer).

Introducer informs N40 of N45.

N40 initializes its ring successor to N45.

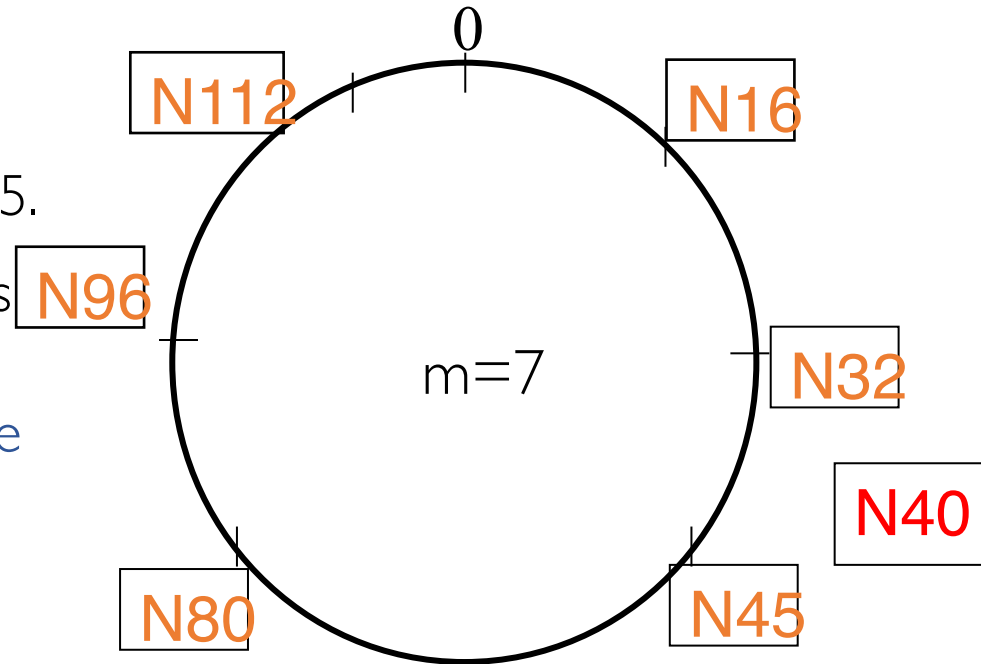
N40 notifies N45, and N45 updates its ring predecessor to N40.

N45 can begin transferring appropriate key-value pairs to N40.

N32 realizes its new successor is N40 when it runs stabilization.

N32 notifies N40, and N40 initializes its ring predecessor to N32.

Periodically and eventually, each node update their finger table entries.



Stabilization Protocol (contd)

- Failures can be handled in a similar way.
 - *Also need failure detectors (you've seen them!)*
 - Maintain knowledge of r ring successors.
 - The predecessor of a failed node update's its ring successor, and notifies it.

Stabilization Protocol (contd)

- Look-ups may fail while the Chord system is getting stabilized.
 - Such failures are transient.
 - Eventually ring successors and finger-table entries will get updated.
 - Application can then try again after a timeout.
 - Such failures are also unlikely in practice
 - Multiple key-value replicas and ring successors.

Chord Summary

- Consistent hashing for load balancing.
- $O(\log n)$ lookups via correct finger tables.
- *Correctness* of lookups requires correctly maintaining ring successors.
- As nodes join and leave a Chord network, runs a stabilization protocol to periodically update ring successors and finger table entries.
- Fault tolerance: Maintain r ring successors and r key replicas.

Our agenda for the next 3-4 classes

- Brief overview of key-value stores
- Distributed Hash Tables
 - Peer-to-peer protocol for efficient insertion and retrieval of key-value pairs.
- Key-value stores in the cloud
 - How to run large-scale distributed computations over key-value stores?
 - Map-Reduce Programming Abstraction
 - How to schedule jobs in the cloud?
 - How to design a large-scale distributed key-value store?
 - Case-study: Facebook's Cassandra

Our agenda for the next 3-4 classes

- Brief overview of key-value stores
- Distributed Hash Tables
 - Peer-to-peer protocol for efficient insertion and retrieval of key-value pairs.
- Key-value stores in the cloud
 - How to run large-scale distributed computations over key-value stores?
 - Map-Reduce Programming Abstraction
 - How to schedule jobs in the cloud?
 - How to design a large-scale distributed key-value store?
 - Case-study: Facebook's Cassandra

Cloud Computing

Many Cloud Providers

- AWS: Amazon Web Services
 - EC2: Elastic Compute Cloud
 - S3: Simple Storage Service
- Microsoft Azure
- Google Cloud/Compute Engine/AppEngine
- Rightscale, Salesforce, EMC, Gigaspaces, I0gen, Datastax, Oracle, VMWare, Yahoo, Cloudera
- And many many more!

What is a cloud?

- Cloud = Lots of storage + compute cycles nearby



- Cloud services provide:
 - managed *clusters* for distributed computing.
 - managed *distributed datastores*.

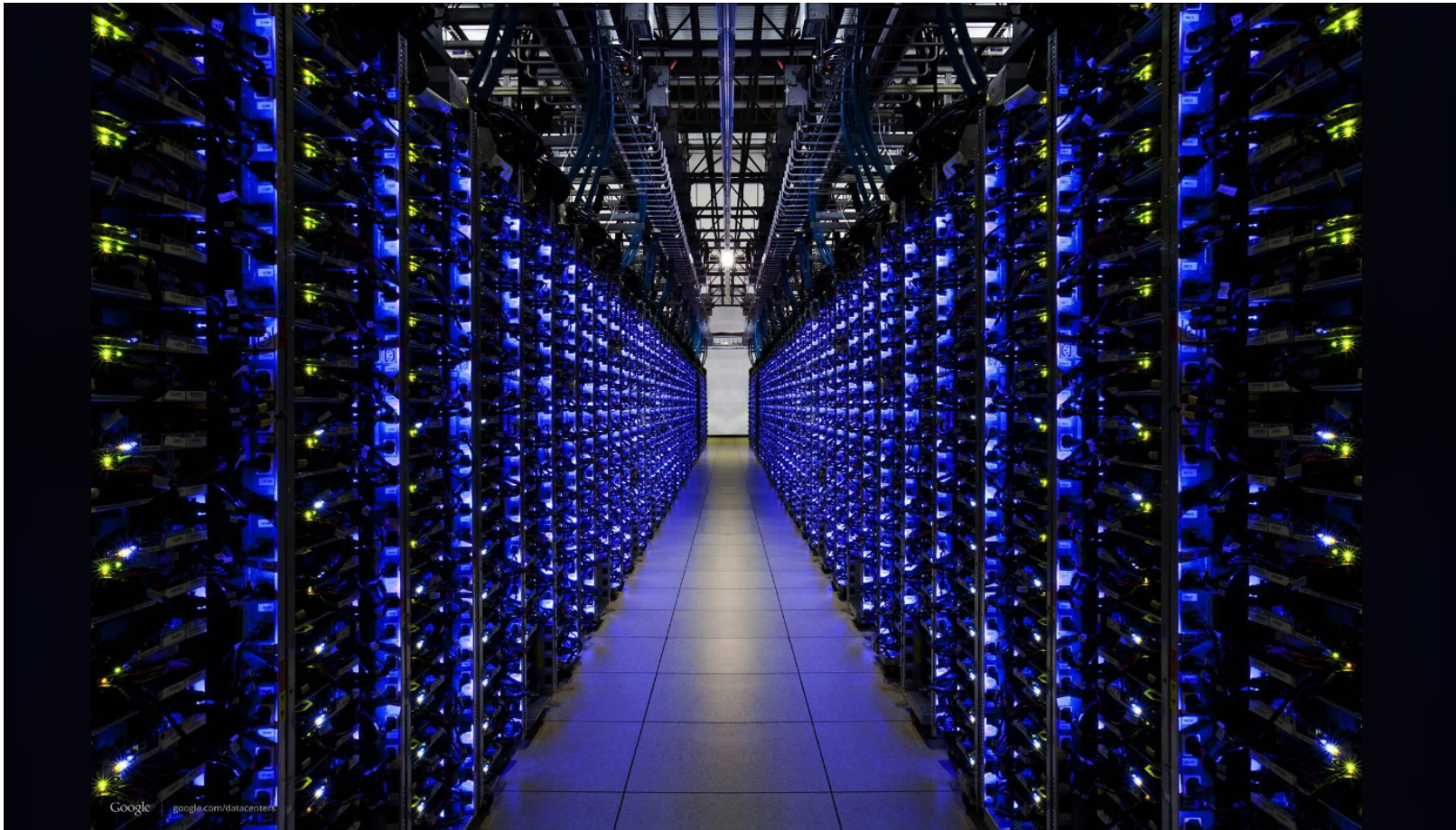
What is a cloud?

- A single cloud-site (aka “Datacenter”) consists of
 - Compute nodes (grouped into racks) (2)
 - Switches, connecting the racks in a hierarchical network topology.
 - Storage (backend) nodes connected to the network (3)
 - Front-end for submitting jobs and receiving client requests (1)
 - (1-3: Often called “three-tier architecture”)
- A geographically distributed cloud consists of
 - Multiple such sites
 - Each site perhaps with a different structure and services

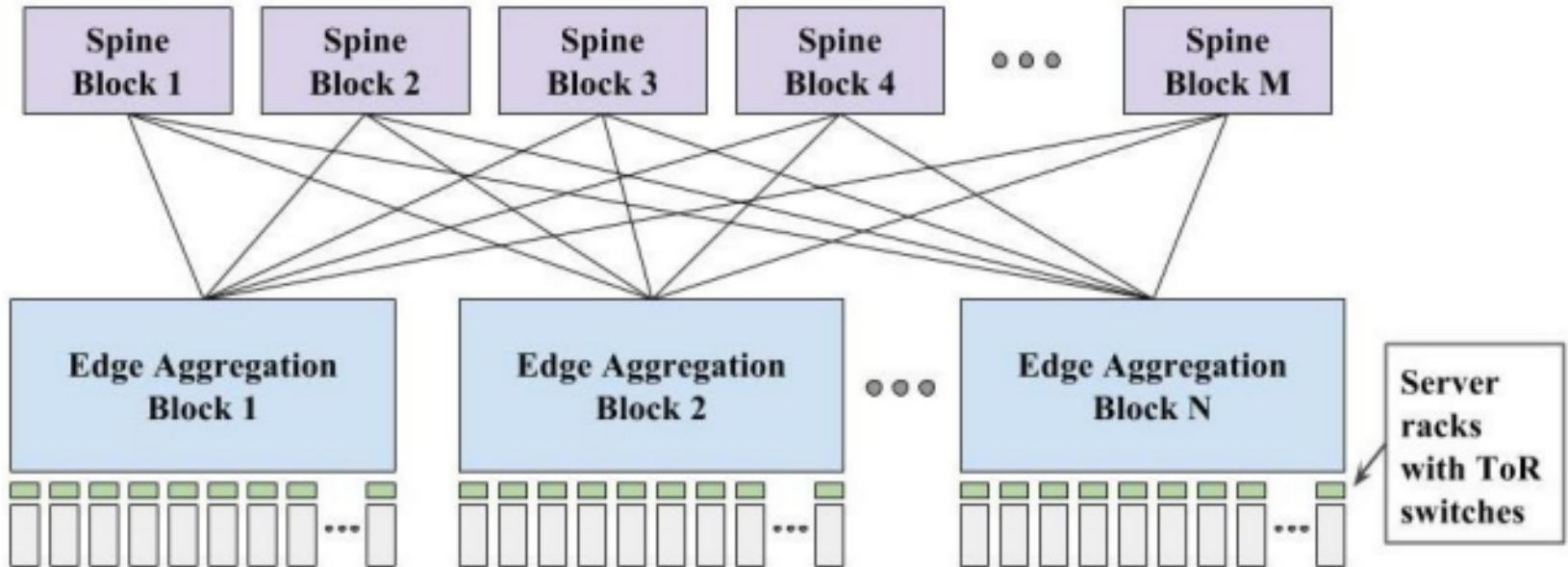
Picture of Google's Datacenter



Picture of Google's Datacenter



Typical Datacenter Topology



Features of cloud

I. Massive scale.

- Tens of thousands of servers and cloud tenants, and hundreds of thousands of VMs.

II. On-demand access:

- Pay-as-you-go, no upfront commitment, access to anyone.

III. Data-intensive nature:

- What was MBs has now become TBs, PBs and XBs.
 - Daily logs, forensics, Web data, etc.

Must deal with immense complexity!

- Fault-tolerance and failure-handling
- Replication and consensus
- Cluster scheduling

- How would a cloud user deal with such complexity?
 - **Powerful abstractions and frameworks**
 - Provide **easy-to-use** API to users.
 - Deal with the complexity of distributed computing under the hood.

MapReduce
is one such powerful
abstraction.

MapReduce Abstraction

- Map/Reduce
 - Programming model inspired from LISP (and other functional languages).
- Expressive: many problems can be phrased as map/reduce.
- Easy to distribute across nodes.
 - High-level job divided into multiple independent “map” tasks, followed by multiple independent “reduce” tasks.
- Nice retry/failure semantics.

MapReduce Architecture

- *MapReduce programming abstraction:*
 - Easy to program distributed computing tasks.
- MapReduce programming abstraction offered by multiple open-source *application frameworks*:
 - Handle creation of “map” and “reduce” tasks.
 - e.g. *Hadoop: one of the earliest map-reduce frameworks.*
 - e.g. *Spark: easier API and performance optimizations.*
- Application frameworks use *resource managers*.
 - Deal with the hassle of distributed cluster management.
 - e.g. *Kubernetes, YARN, Mesos, etc.*

MapReduce Architecture

- *Map/Reduce abstraction:*
 - Easy to program distributed computing tasks.
- MapReduce features:
 - Automatic parallelization & distribution
 - Fault tolerance
 - Scheduling
 - Monitoring & status updates
- Application frameworks use *resource managers*.
 - Deal with the hassle of distributed cluster management.
 - e.g. *Kubernetes, YARN, Mesos, etc.*

MapReduce Architecture

- *Map/Reduce abstraction:*
 - Easy to program distributed computing tasks.
- MapReduce programming abstraction offered by multiple open-s
 - Cre
 - e.g. Hadoop: one of the earliest map reduce frameworks.
 - e.g. Spark: easier API and performance optimizations.
- Application frameworks use *resource managers*.
 - Deal with the hassle of distributed cluster management.
 - e.g. Kubernetes, YARN, Mesos, etc.

To be continued in next class....