



CS514: Intermediate Course in Computer Systems

Lecture 30: April 4, 2003
“All about DHTs”



What is a Distributed Hash Table (DHT)?

CS514

- Exactly that ☺
- A service, distributed over multiple machines, with hash table semantics
 - *Put(key, value)*, *Value = Get(key)*
- Designed to work in a peer-to-peer (P2P) environment
 - No central control
 - Nodes under different administrative control
- But of course can operate in an “infrastructure” sense



More specifically

CS514

- Hash table semantics:
 - *Put*(key, value),
 - Value = *Get*(key)
 - Key is a single flat string
 - Limited semantics compared to keyword search
- *Put*() causes value to be stored at one (or more) peer(s)
- *Get*() retrieves value from a peer
- *Put*() and *Get*() accomplished with unicast routed messages
 - In other words, it scales
- Other API calls to support application, like notification when neighbors come and go



P2P “environment”

CS514

- Nodes come and go at will (possibly quite frequently---a few minutes)
- Nodes have heterogeneous capacities
 - Bandwidth, processing, and storage
- Nodes may behave badly
 - Promise to do something (store a file) and not do it (free-loaders)
 - Attack the system



Several flavors, each with variants

CS514

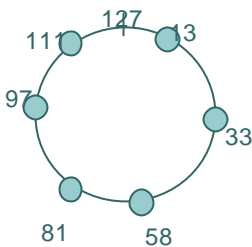
- Tapestry (Berkeley)
 - Based on Plaxton trees---similar to hypercube routing
 - The first* DHT
 - Complex and hard to maintain (hard to understand too!)
- CAN (ACIRI), Chord (MIT), and Pastry (Rice/MSR Cambridge)
 - Second wave of DHTs (contemporary with and independent of each other)

* Landmark Routing, 1988, used a form of DHT called Assured Destination Binding (ADB)



Basics of all DHTs

CS514

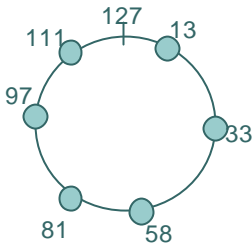


- Goal is to build some “structured” overlay network with the following characteristics:
 - Node IDs can be mapped to the hash key space
 - Given a hash key as a “destination address”, you can route through the network to a given node
 - Always route to the same node no matter where you start from



Simple example (doesn't scale)

CS514

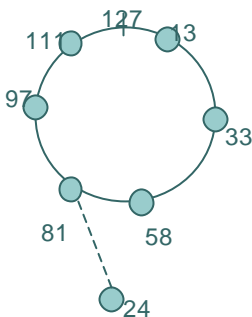


- Circular number space 0 to 127
- Routing rule is to move counter-clockwise until current node ID \geq key, and last hop node ID $<$ key
- Example: key = 42
- Obviously you will route to node 58 from no matter where you start



Building any DHT

CS514

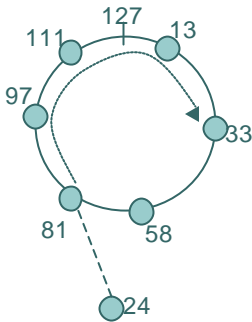


- Newcomer always starts with at least one known member



Building any DHT

CS514

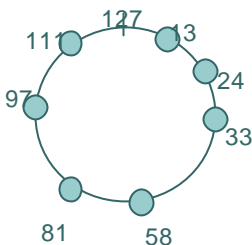


- Newcomer always starts with at least one known member
- Newcomer searches for “self” in the network
 - hash key = newcomer’s node ID
 - Search results in a node in the vicinity where newcomer needs to be



Building any DHT

CS514

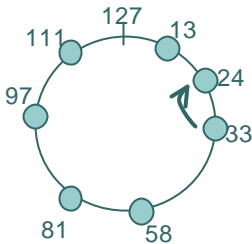


- Newcomer always starts with at least one known member
- Newcomer searches for “self” in the network
 - hash key = newcomer’s node ID
 - Search results in a node in the vicinity where newcomer needs to be
- Links are added/removed to satisfy properties of network



Building any DHT

CS514

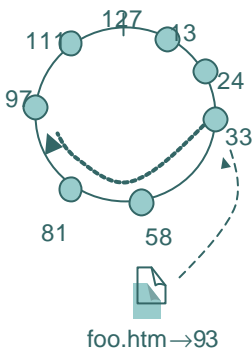


- Newcomer always starts with at least one known member
- Newcomer searches for “self” in the network
 - hash key = newcomer’s node ID
- Search results in a node in the vicinity where newcomer needs to be
- Links are added/removed to satisfy properties of network
- Objects that now hash to new node are transferred to new node



Insertion/lookup for any DHT

CS514



- Hash name of object to produce key
 - Well-known way to do this
- Use key as destination address to route through network
 - Routes to the target node
- Insert object, or retrieve object, at the target node



Properties of all DHTs

CS514

- Memory requirements grow (something like) logarithmically with N
- Routing path length grows (something like) logarithmically with N
- Cost of adding or removing a node grows (something like) logarithmically with N
- Has caching, replication, etc...



DHT Issues

CS514

- Resilience to failures
- Load Balance
 - Heterogeneity
 - Number of objects at each node
 - Routing hot spots
 - Lookup hot spots
- Locality (performance issue)
- Churn (performance and correctness issue)
- Security



We're going to look at four DHTs

CS514

- At varying levels of detail...
 - CAN (Content Addressable Network)
 - ACIRI (now ICIR)
 - Chord
 - MIT
 - Kelips
 - Cornell
 - Pastry
 - Rice/Microsoft Cambridge



Things we're going to look at

CS514

- What is the structure?
- How does routing work in the structure?
- How does it deal with node departures?
- How does it scale?
- How does it deal with locality?
- What are the security issues?



CAN structure is a cartesian
coordinate space in a D
dimensional torus

CS514

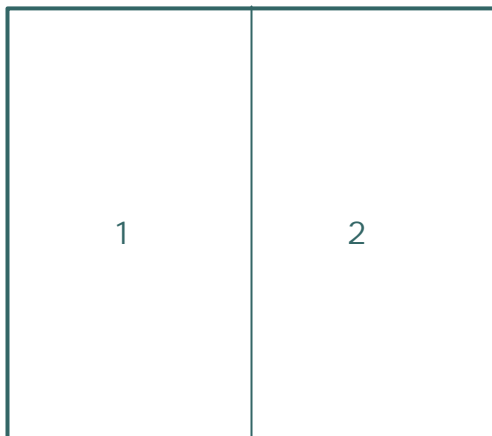


CAN graphics care of Santashil PalChaudhuri, Rice Univ



Simple example in two
dimensions

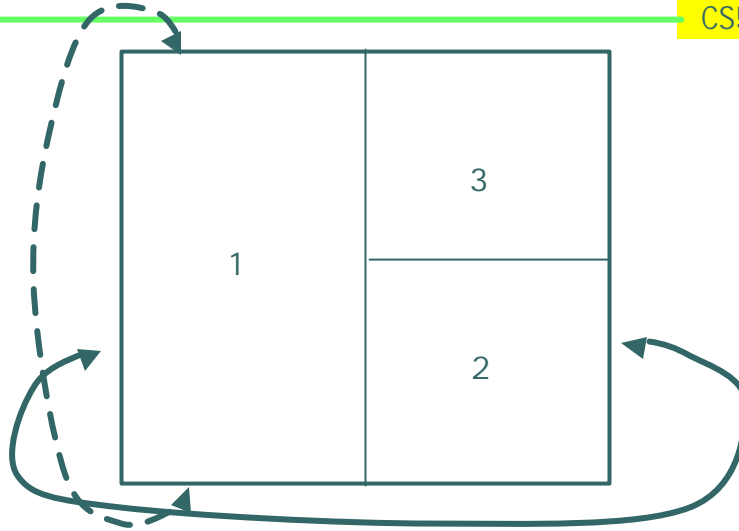
CS514





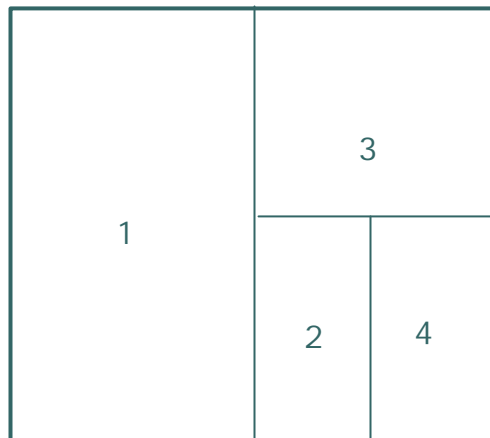
Note: torus wraps on “top”
and “sides”

CS514



Each node in CAN network
occupies a “square” in the
space

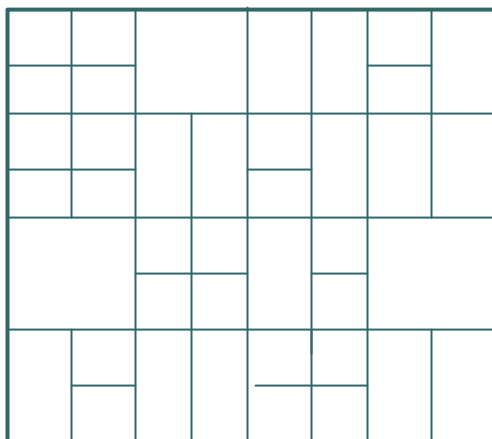
CS514





With relatively uniform square sizes

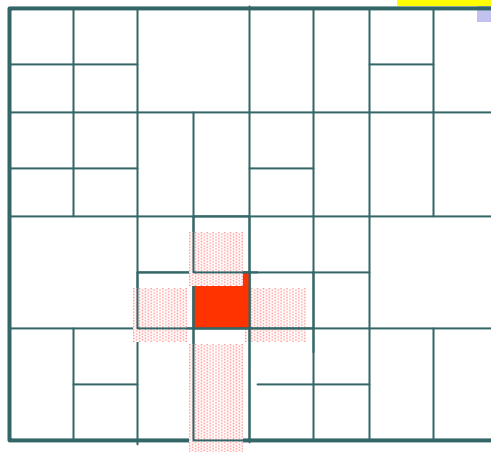
CS514



Neighbors in CAN network

CS514

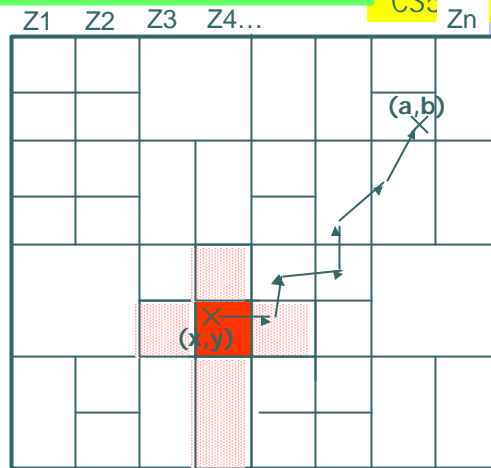
- Neighbor is a node that:
- Overlaps d-1 dimensions
- Abuts along one dimension



Route to neighbors that are closer to target

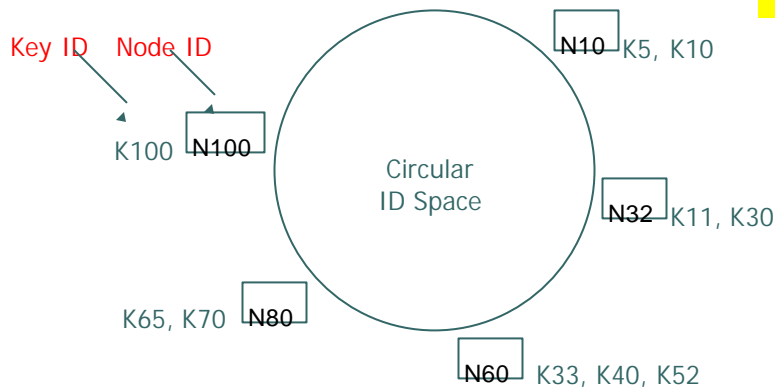
CS514

- d-dimensional space
- n zones
 - Zone is space occupied by a "square" in one dimension
- Avg. route path length
 - $(d/4)(n^{1/d})$
- Number neighbors = $O(d)$
- Tunable (vary d or n)
- Can factor proximity into route decision



Chord uses a circular ID space

CS514



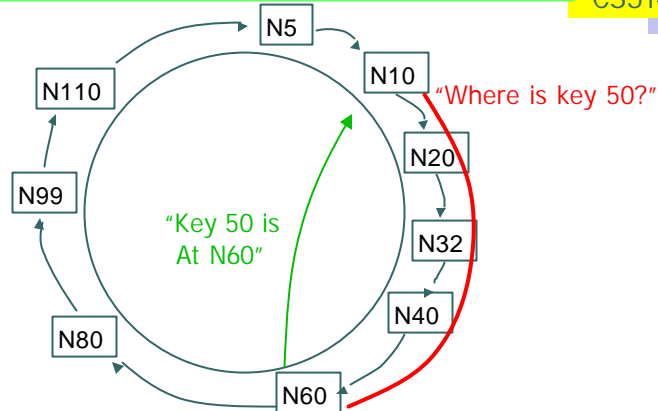
- Successor: node with next highest ID

Chord slides care of Robert Morris, MIT



Basic Lookup

CS514

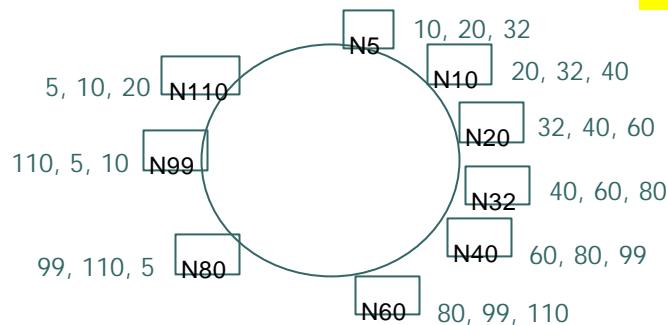


- Lookups find the ID's predecessor
- Correct if successors are correct



Successor Lists Ensure Robust Lookup

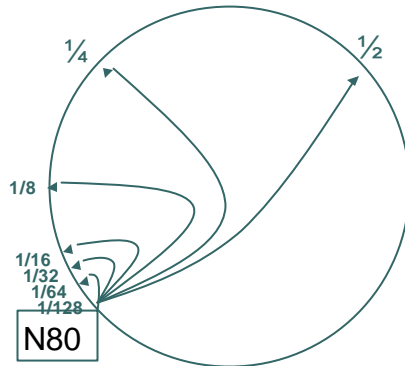
CS514



- Each node remembers r successors
- Lookup can skip over dead nodes to find blocks
- Periodic check of successor and predecessor links

Chord "Finger Table" Accelerates Lookups

CS514

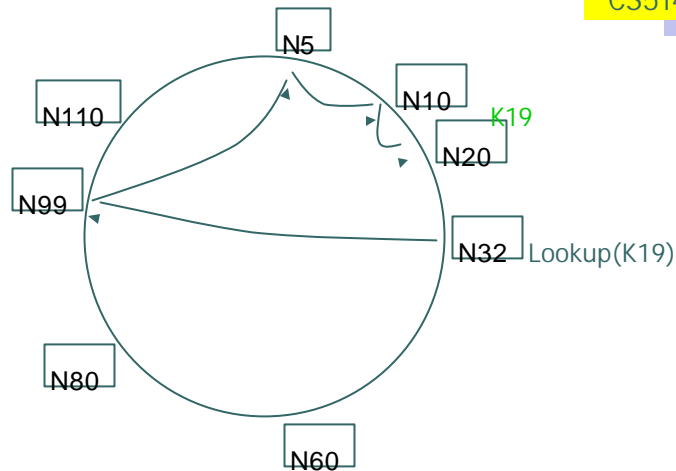


To build finger tables, new node searches for the key values for each finger

To do it efficiently, new nodes obtain successor's finger table, and use as a hint to optimize the search

Chord lookups take $O(\log N)$ hops

CS514





Drill down on Chord reliability

CS514

- Interested in maintaining a correct routing table (successors, predecessors, and fingers)
- Primary invariant: correctness of successor pointers
 - Fingers, while important for performance, do not have to be exactly correct for routing to work
 - Algorithm is to “get closer” to the target
 - Successor nodes always do this



Maintaining successor pointers

CS514

- Periodically run “stabilize” algorithm
 - Finds successor’s predecessor
 - Repair if this isn’t self
- This algorithm is also run at join
- Eventually routing will repair itself
- Fix_finger also periodically run
 - For randomly selected finger



Initial: 25 wants to join
correct ring (between 20 and
30)

CS514



25



20

30

25



20

25

30

25 finds successor,
and tells successor
(30) of itself

20 runs "stabilize":
20 asks 30 for 30's predecessor
30 returns 25
20 tells 25 of itself



This time, 28 joins before 20
runs "stabilize"

CS514

20

30

28



20

30

25

28



20

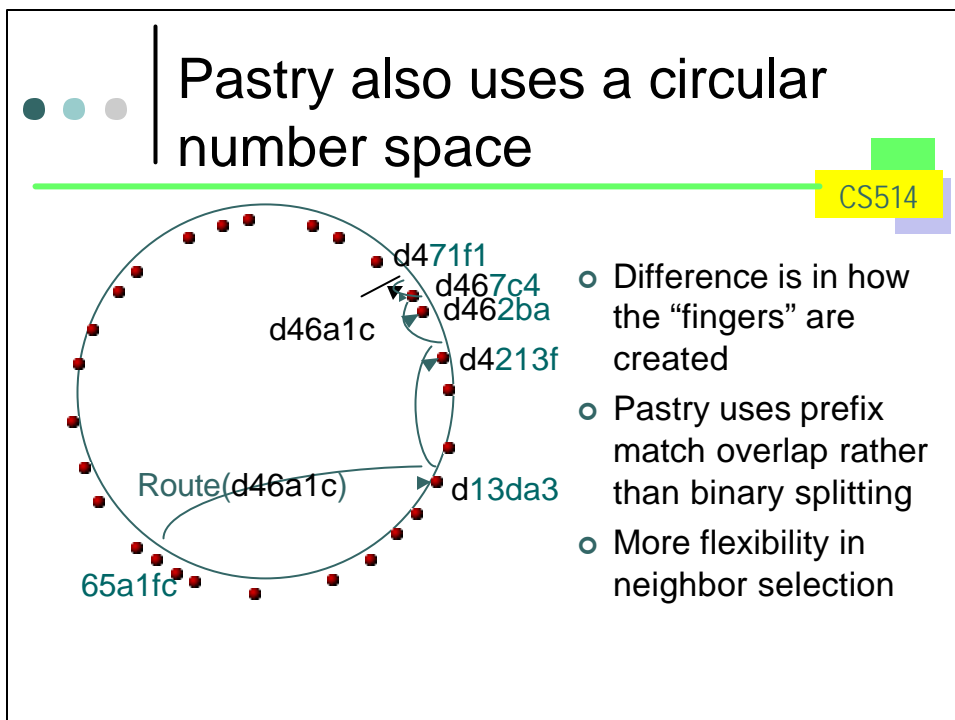
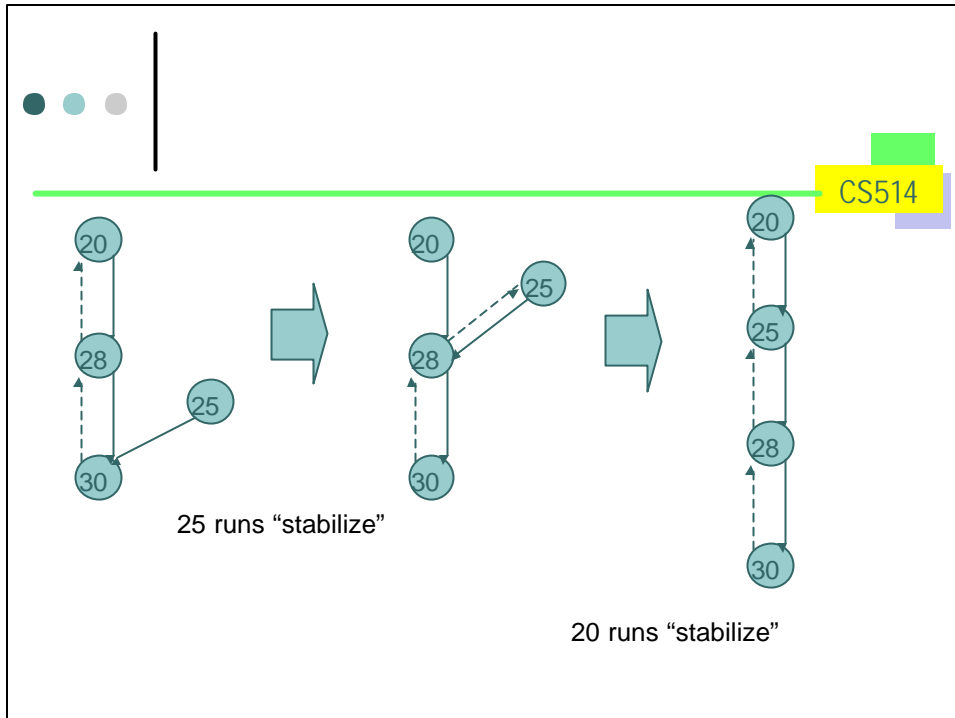
28

30

25

28 finds successor,
and tells successor
(30) of itself

20 runs "stabilize":
20 asks 30 for 30's predecessor
30 returns 28
20 tells 28 of itself



Pastry routing table (for node 65a1fc)

0	1	2	3	4	5	6	7	8	9	a	b	c	d	e	f
x	x	x	x	x	x		x	x	x	x	x	x	x	x	x
6	6	6	6	6		6	6	6	6	6	6	6	6	6	6
0	1	2	3	4		6	7	8	9	a	b	c	d	e	f
x	x	x	x	x		x	x	x	x	x	x	x	x	x	x
6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6
5	5	5	5	5	5	5	5	5	5	5	5	5	5	5	5
0	1	2	3	4	5	6	7	8	9	b	c	d	e	f	
x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x
6		6	6	6	6	6	6	6	6	6	6	6	6	6	6
5		5	5	5	5	5	5	5	5	5	5	5	5	5	5
a		a	a	a	a	a	a	a	a	a	a	a	a	a	a
0		2	3	4	5	6	7	8	9	a	b	c	d	e	f
x		x	x	x	x	x	x	x	x	x	x	x	x	x	x

CS514

Pastry nodes also have a “leaf set” of immediate neighbors up and down the ring

Similar to Chord's list of successors

Pastry join

CS514

- X = new node, A = bootstrap, Z = nearest node
- A finds Z for X
- In process, A, Z, and all nodes in path send state tables to X
- X settles on own table
 - Possibly after contacting other nodes
- X tells everyone who needs to know about itself
- Pastry paper doesn't give enough information to understand how concurrent joins work
 - 18th IFIP/ACM, Nov 2001

Ask other neighbors

0	1	2	3	4	5	6	7	8	9	a	b	c	d	e	f
x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x
6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6
0	1	2	3	4	5	6	7	8	9	a	b	c	d	e	f
x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x
6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6
5	5	5	5	5	5	5	5	5	5	5	5	5	5	5	5
0	1	2	3	4	5	6	7	8	9	a	b	c	d	e	f
x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x
6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6
5	5	5	5	5	5	5	5	5	5	5	5	5	5	5	5
a	a	a	a	a	a	a	a	a	a	a	a	a	a	a	a
0	1	2	3	4	5	6	7	8	9	a	b	c	d	e	f
x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x

CS514

Try asking some neighbor in the same row for its 655x entry

If it doesn't have one, try asking some neighbor in the row below, etc.

CAN, Chord, Pastry differences

CS514

- CAN, Chord, and Pastry have deep similarities
- Some (important???) differences exist
 - CAN nodes tend to know of multiple nodes that allow equal progress
 - Can therefore use additional criteria (RTT) to pick next hop
 - Pastry allows greater choice of neighbor
 - Can thus use additional criteria (RTT) to pick neighbor
 - In contrast, Chord has more determinism
 - Harder for an attacker to manipulate system?



Security issues

CS514

- In many P2P systems, members may be malicious
- If peers untrusted, all content must be signed to detect forged content
 - Requires certificate authority
 - Like we discussed in secure web services talk
 - This is not hard, so can assume at least this level of security



Security issues: Sybil attack

CS514

- Attacker pretends to be multiple system
 - If surrounds a node on the circle, can potentially arrange to capture all traffic
 - Or if not this, at least cause a lot of trouble by being many nodes
- Chord requires node ID to be an SHA-1 hash of its IP address
 - But to deal with load balance issues, Chord variant allows nodes to replicate themselves
- *A central authority must hand out node IDs and certificates to go with them*
 - Not P2P in the Gnutella sense



General security rules

CS514

- Check things that can be checked
 - Invariants, such as successor list in Chord
- Minimize invariants, maximize randomness
 - Hard for an attacker to exploit randomness
- Avoid any single dependencies
 - Allow multiple paths through the network
 - Allow content to be placed at multiple nodes
- But all this is expensive...



Load balancing

CS514

- Query hotspots: given object is popular
 - Cache at neighbors of hotspot, neighbors of neighbors, etc.
 - Classic caching issues
- Routing hotspot: node is on many paths
 - Of the three, Pastry seems most likely to have this problem, because neighbor selection more flexible (and based on proximity)
 - This doesn't seem adequately studied



Load balancing

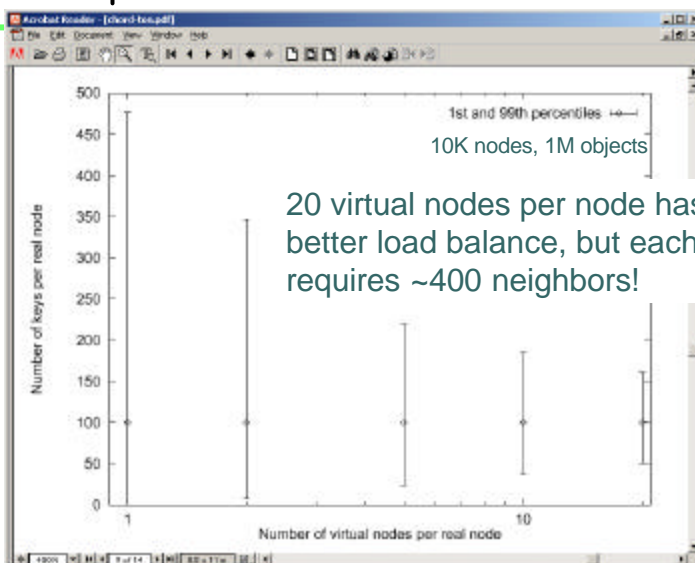
CS514

- Heterogeneity (variance in bandwidth or node capacity)
- Poor distribution in entries due to hash function inaccuracies
- One class of solution is to allow each node to be multiple virtual nodes
 - Higher capacity nodes virtualize more often
 - But security makes this harder to do



Chord node virtualization

CS514



20 virtual nodes per node has much better load balance, but each node requires ~400 neighbors!



Primary concern: churn

CS514

- Churn: nodes joining and leaving frequently
- Join or leave requires a change in some number of links
- Those changes depend on correct routing tables in other nodes
 - Cost of a change is higher if routing tables not correct
 - In chord, ~6% of lookups fail if three failures per stabilization
- But as more changes occur, probability of incorrect routing tables increases



Control traffic load generated by churn

CS514

- Chord and Pastry appear to deal with churn differently
- Chord join involves some immediate work, but repair is done periodically
 - Extra load only due to join messages
- Pastry join and leave involves immediate repair of all effected nodes' tables
 - Routing tables repaired more quickly, but cost of each join/leave goes up with frequency of joins/leaves
 - Scales quadratically with number of changes???
 - Can result in network meltdown???



Kelips takes a different approach

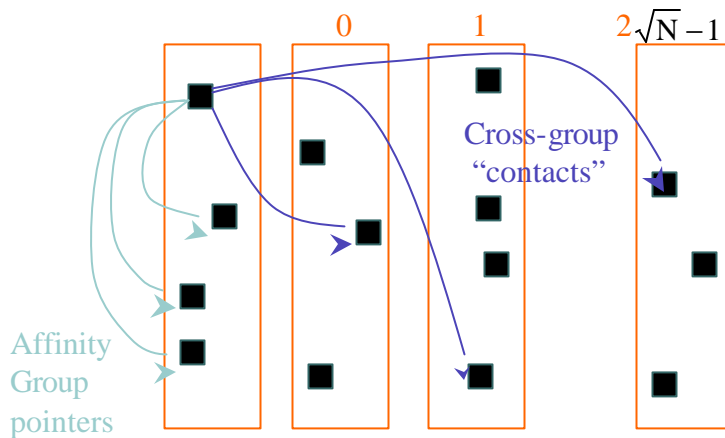
CS514

- Network partitioned into \sqrt{N} “affinity groups”
- Hash of node ID determines which affinity group a node is in
- Each node knows:
 - One or more nodes in each group
 - All objects and nodes in own group
- *But this knowledge is soft-state, spread through peer-to-peer “gossip” (epidemic multicast)!*



Kelips

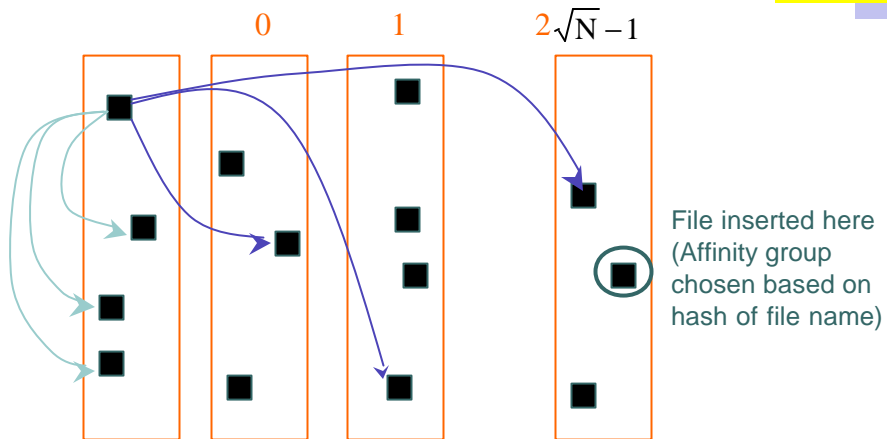
CS514



Graphics courtesy Indranil Gupta, Cornell

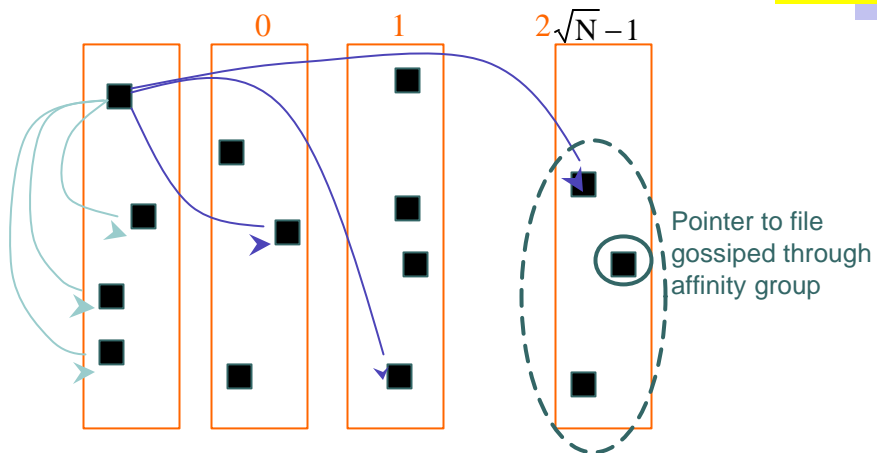
Kelips

CS514



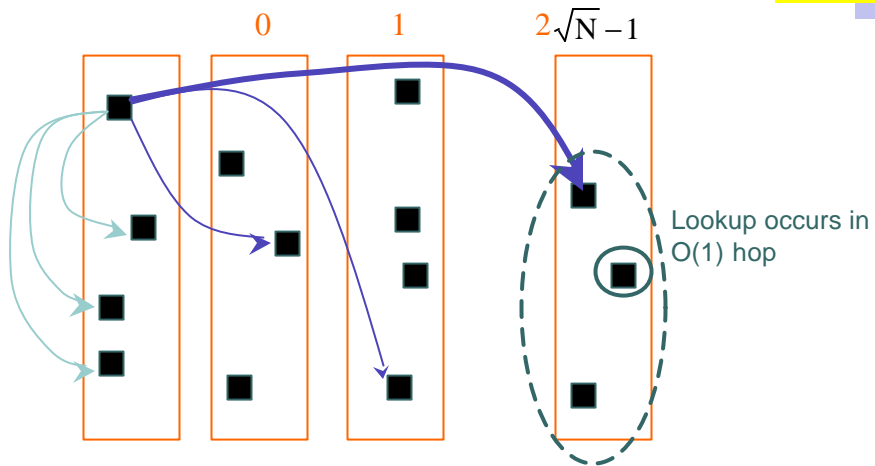
Kelips

CS514



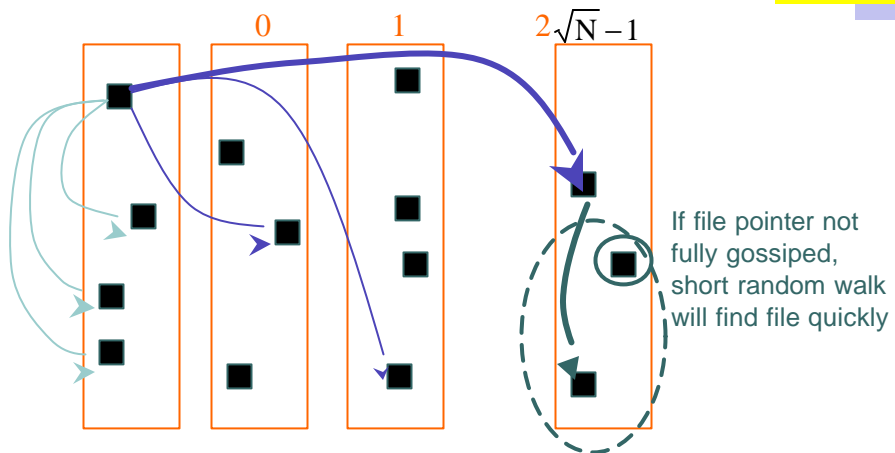
Kelips

CS514



Kelips

CS514





Kelips gossip

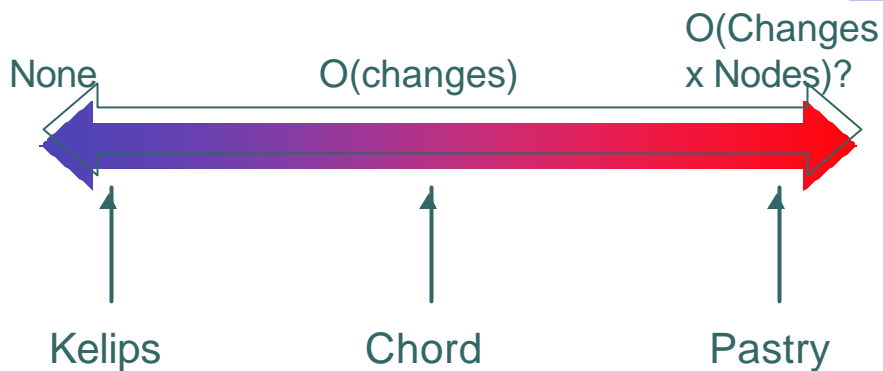
CS514

- Operates at constant “background” rate
 - Independent of frequency of changes in the system
 - Average overhead may be higher than other DHTs, but not bursty
- If churn too high, system performs poorly (failed lookups), *but does not collapse...*



Control traffic load generated by churn

CS514





DHT applications

CS514

- Ken will talk about file system and archival storage applications for DHT
- DNS as a DHT application
 - No structure required in name!
 - Fewer administration errors
 - No DoS target
- DNS over Chord
 - Median lookup time increased from 43ms (DNS) to 350ms



To finish up

CS514

- Various applications have been designed over DHTs
 - File system, DNS-like service, pub/sub system
- DHTs are elegant and promising tools
- Concerns about churn and security