

Hashing in Networked Systems

COS 461: Computer Networks
Spring 2011

Mike Freedman

<http://www.cs.princeton.edu/courses/archive/spring11/cos461/>

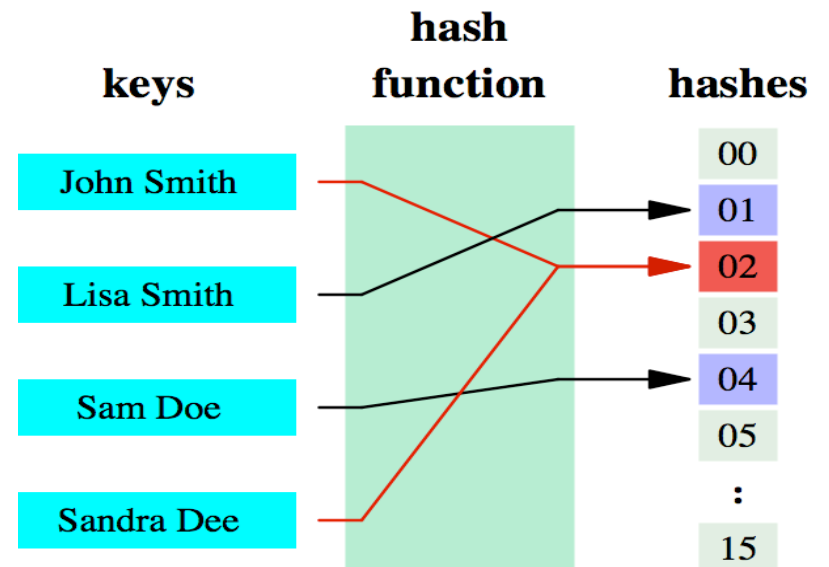
Hashing

- **Hash function**

- Function that maps a large, possibly variable-sized datum into a small datum, often a single integer that serves to index an associative array
- In short: maps n -bit datum into k buckets ($k \ll 2^n$)
- Provides time- & space-saving data structure for lookup

- **Main goals:**

- Low cost
- Deterministic
- Uniformity (load balanced)

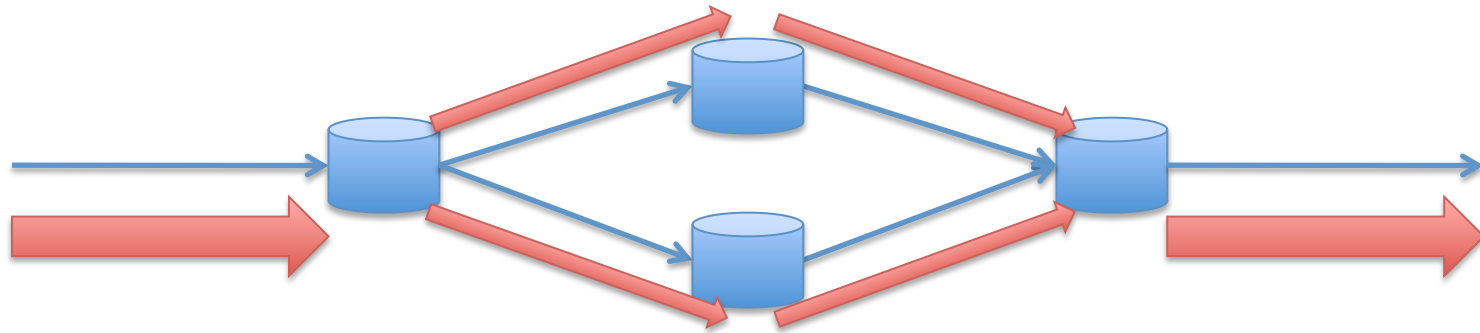


Today's outline

- **Uses of hashing**
 - Equal-cost multipath routing in switches
 - Network load balancing in server clusters
 - Per-flow statistics in switches (QoS, IDS)
 - Caching in cooperative CDNs and P2P file sharing
 - Data partitioning in distributed storage services
- **Various hashing strategies**
 - Modulo hashing
 - Consistent hashing
 - Bloom Filters

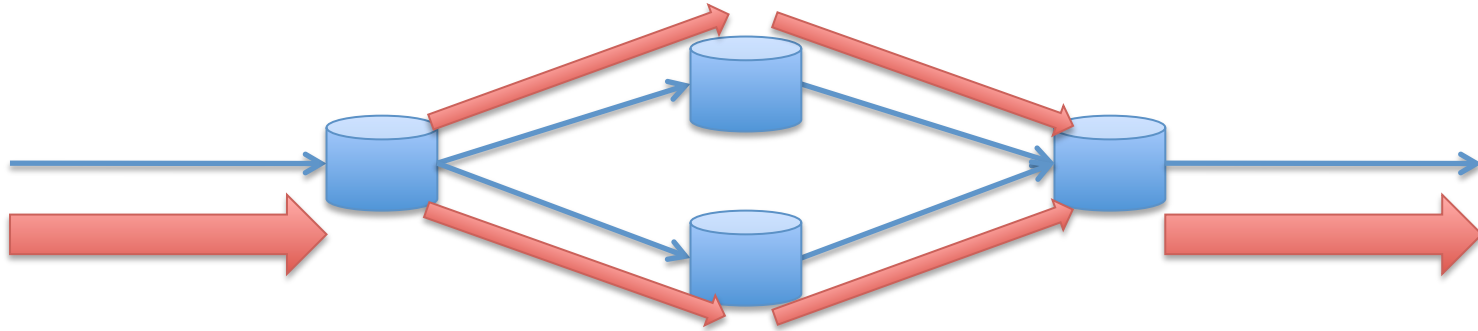
Uses of Hashing

Equal-cost multipath routing (ECMP)



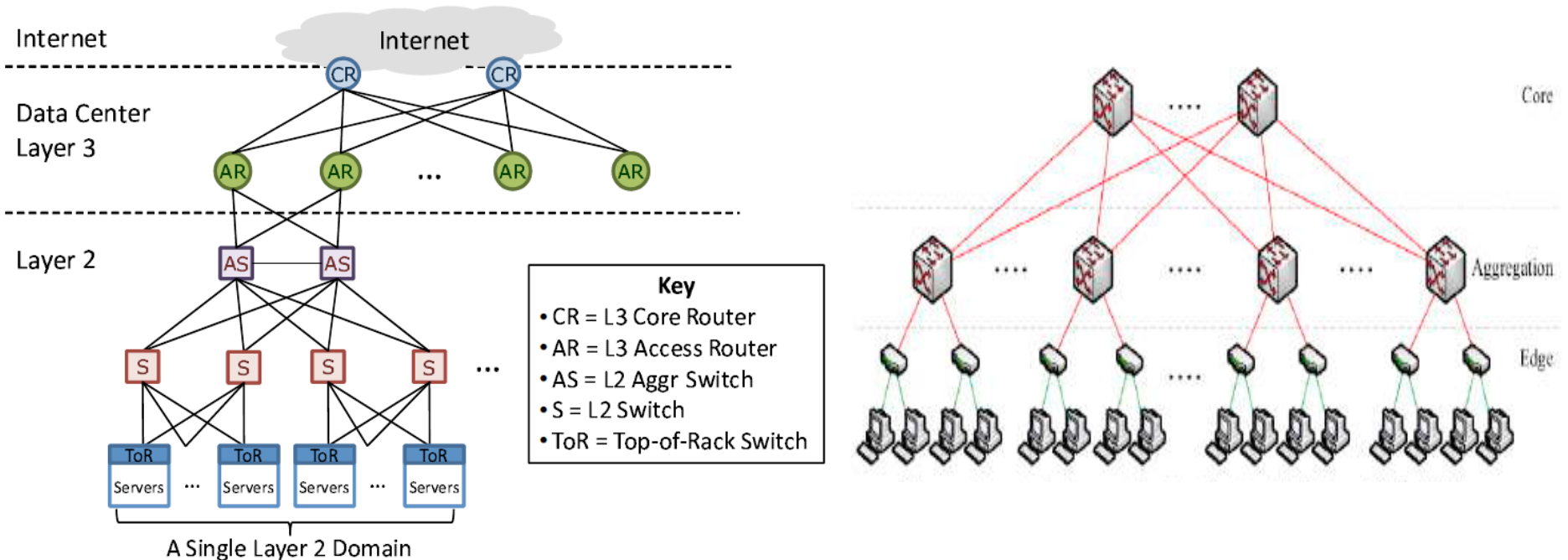
- **ECMP**
 - Multipath routing strategy that splits traffic over multiple paths for load balancing
- **Why not just round-robin packets?**
 - Reordering (lead to triple duplicate ACK in TCP?)
 - Different RTT per path (for TCP RTO)...
 - Different MTUs per path

Equal-cost multipath routing (ECMP)



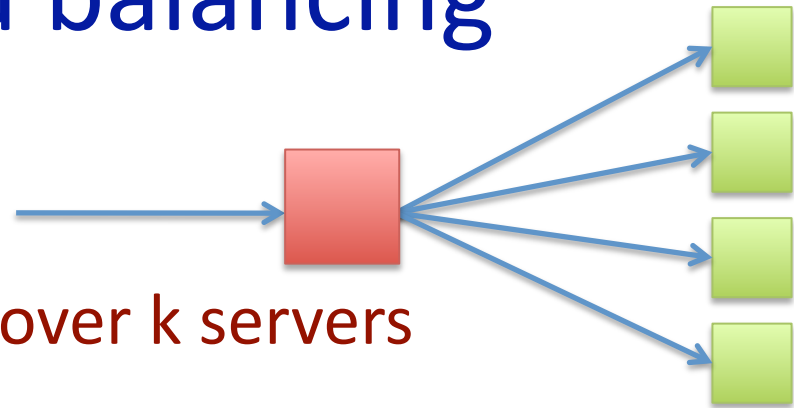
- Path-selection via hashing
 - # buckets = # outgoing links
 - Hash network information (source/dest IP addrs) to select outgoing link: preserves flow affinity

Now: ECMP in datacenters



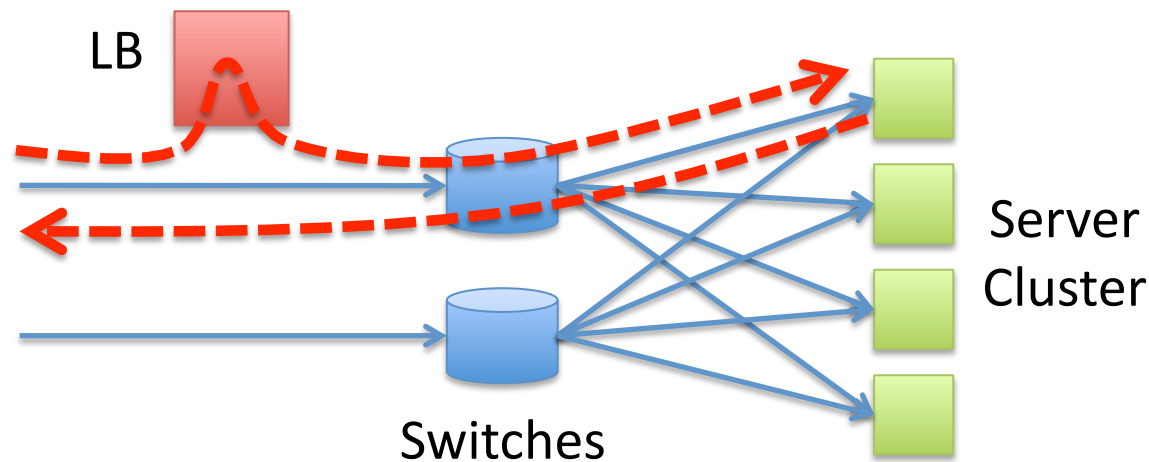
- **Datacenter networks are multi-rooted tree**
 - Goal: Support for 100,000s of servers
 - Recall Ethernet spanning tree problems: No loops
 - L3 routing and ECMP: Take advantage of multiple paths

Network load balancing



- **Goal: Split requests evenly over k servers**
 - Map new flows to any server
 - Packets of existing flows continue to use same server
- **3 approaches**
 - Load balancer terminates TCP, opens own connection to server
 - Virtual IP / Dedicated IP (VIP/DIP) approaches
 - One global-facing virtual IP represents all servers in cluster
 - Hash client's network information (source IP:port)
 - **NAT approach:** Replace virtual IP with server's actual IP
 - **Direct Server Return (DSR)**

Load balancing with DSR



- Servers bind to both virtual and dedicated IP
- Load balancer just replaces dest MAC addr
- Server sees client IP, responds directly
 - Packet in reverse direction do not pass through load balancer
 - Greater scalability, particularly for traffic with asymmetric bandwidth (e.g., HTTP GETs)

Per-flow state in switches

- Switches often need to maintain connection records or per-flow state
 - Quality-of-service for flows
 - Flow-based measurement and monitoring
 - Payload analysis in Intrusion Detection Systems (IDSs)
- On packet receipt:
 - Hash flow information (packet 5-tuple)
 - Perform lookup if packet belongs to known flow
 - Otherwise, possibly create new flow entry
 - Probabilistic match (false positives) may be okay

Cooperative Web CDNs

- **Tree-like topology of cooperative web caches**

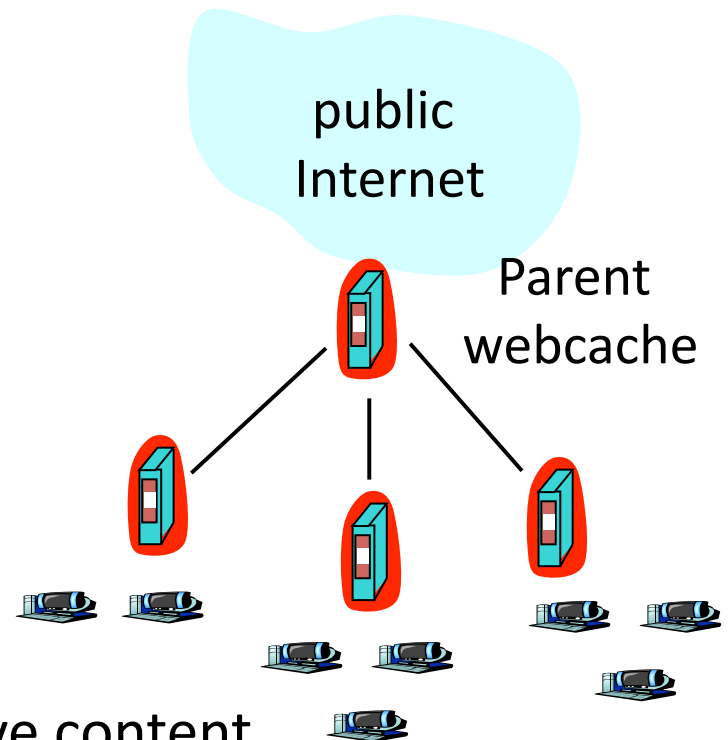
- Check local
- If miss, check siblings / parent

- **One approach**

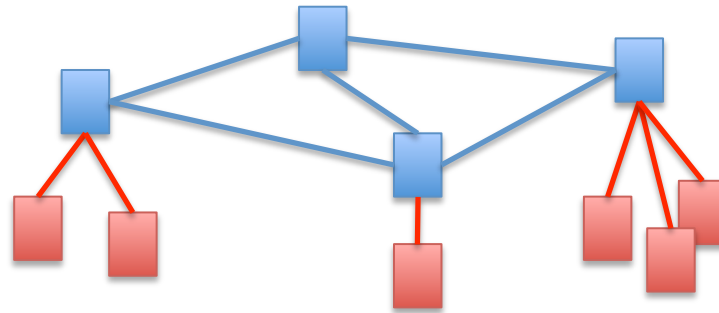
- Internet Cache Protocol (ICP)
- UDP-based lookup, short timeout

- **Alternative approach**

- A priori guess is siblings/children have content
- Nodes share hash table of cached content with parent / siblings
- Probabilistic check (false positives) okay, as actual ICP lookup to neighbor could just return false

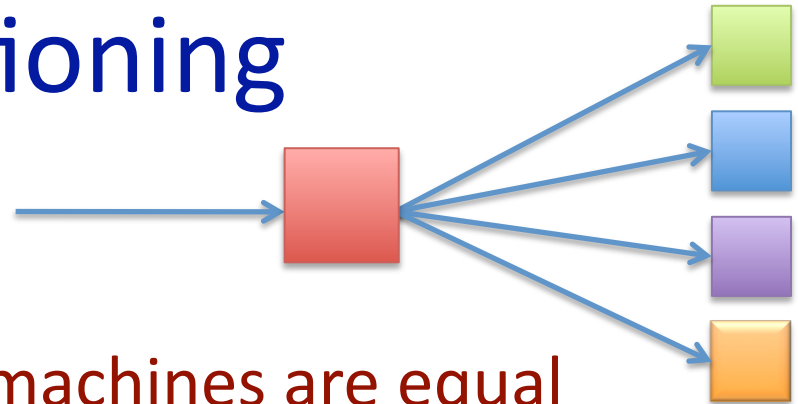


Hash tables in P2P file-sharing



- **Two-layer network (e.g., Gnutella, Kazaa)**
 - Ultrapeers are more stable, not NATted, higher bandwidth
 - Leaf nodes connect with 1 or more ultrapeers
- **Ultrapeers handle content searchers**
 - Leaf nodes send hash table of content to ultrapeers
 - Search requests flooded through ultrapeer network
 - When ultrapeer gets request, checks hash tables of its children for match

Data partitioning

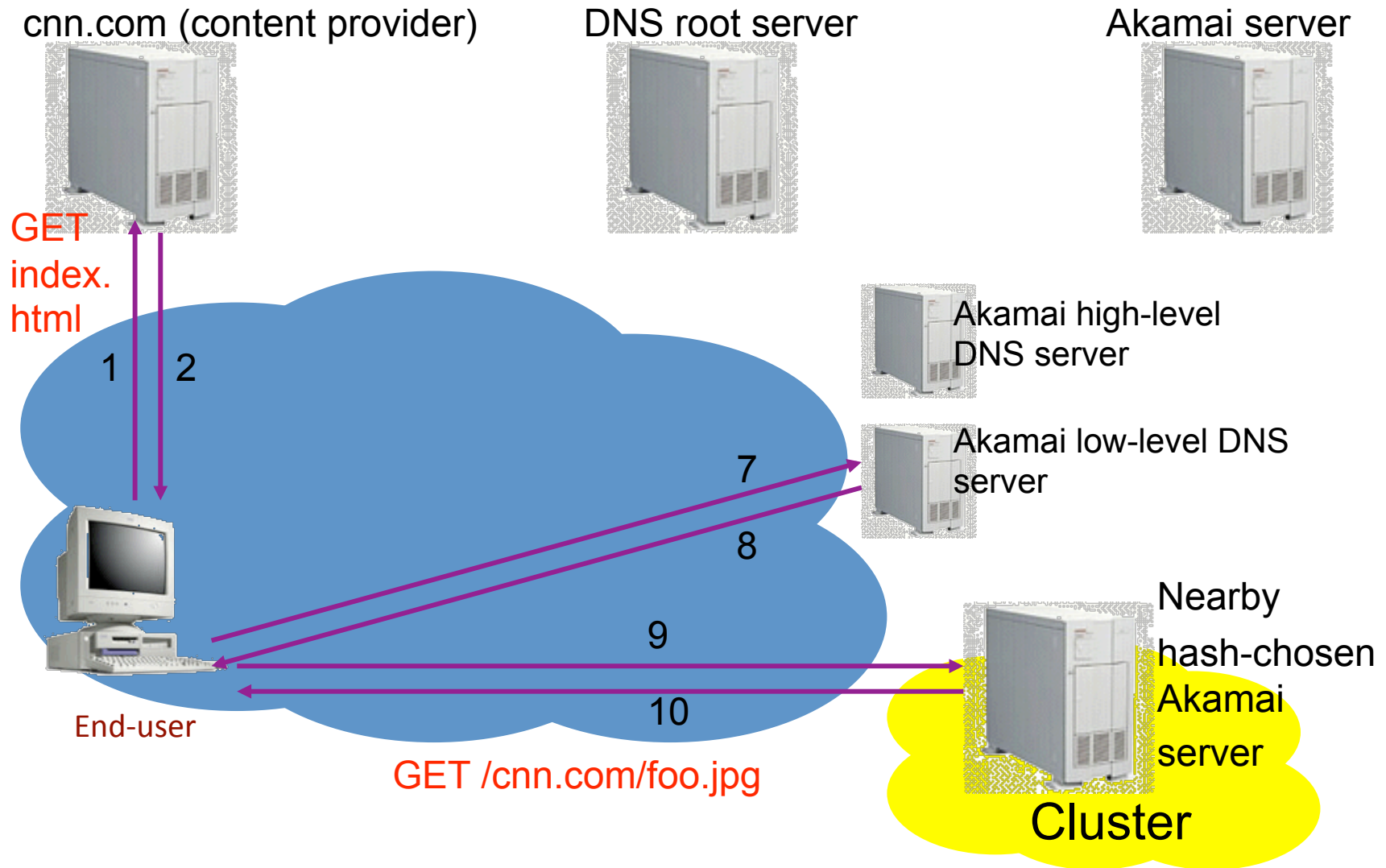


- Network load balancing: All machines are equal
- Data partitioning: Machines store different content
- Non-hash-based solution
 - “Directory” server maintains mapping from $O(\text{entries})$ to machines (e.g., Network file system, Google File System)
 - Named data can be placed on any machine
- Hash-based solution
 - Nodes maintain mappings from $O(\text{buckets})$ to machines
 - Data placed on the machine that owns the name’s bucket

Examples of data partitioning

- **Akamai**
 - 1000 clusters around Internet, each ≥ 1 servers
 - Hash (URL's domain) to map to one server
 - Akamai DNS aware of hash function, returns machine that
 1. is in geographically-nearby cluster
 2. manages particular customer domain
- **Memcached (Facebook, Twitter, ...)**
 - Employ k machines for in-memory key-value caching
 - On read:
 - Check memcache
 - If miss, read data from DB, write to memcache
 - On write: invalidate cache, write data to DB

How Akamai Works – Already Cached



Hashing Techniques

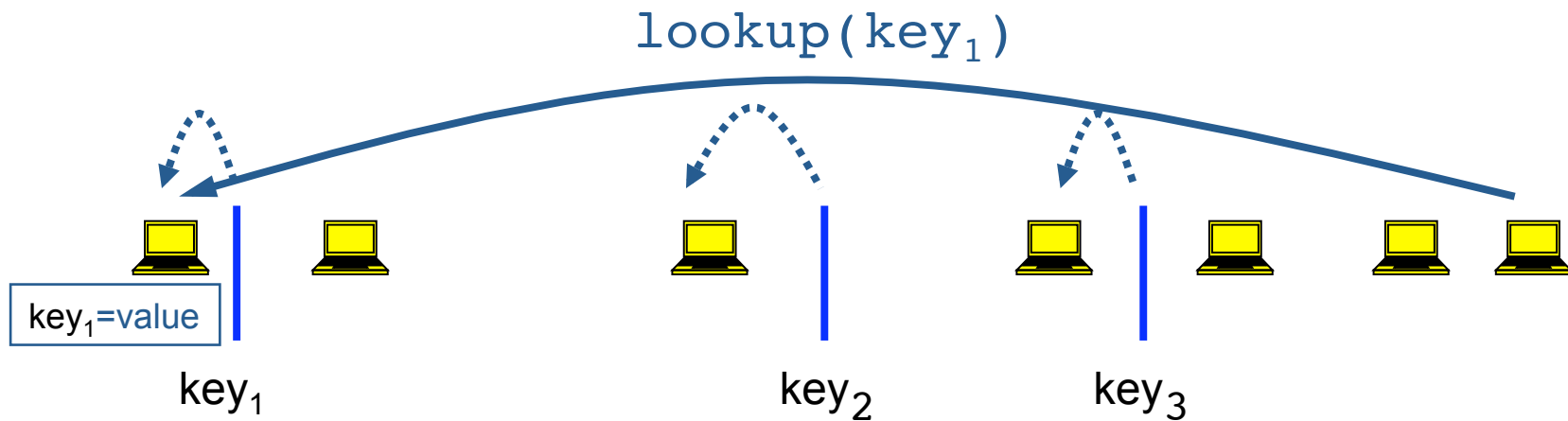
Basic Hash Techniques

- **Simple approach for uniform data**
 - If data distributed uniformly over N , for $N \gg n$
 - Hash fn = $\langle \text{data} \rangle \bmod n$
 - Fails goal of uniformity if data not uniform
- **Non-uniform data, variable-length strings**
 - Typically split strings into blocks
 - Perform rolling computation over blocks
 - CRC32 checksum
 - Cryptographic hash functions (SHA-1 has 64 byte blocks)

Applying Basic Hashing

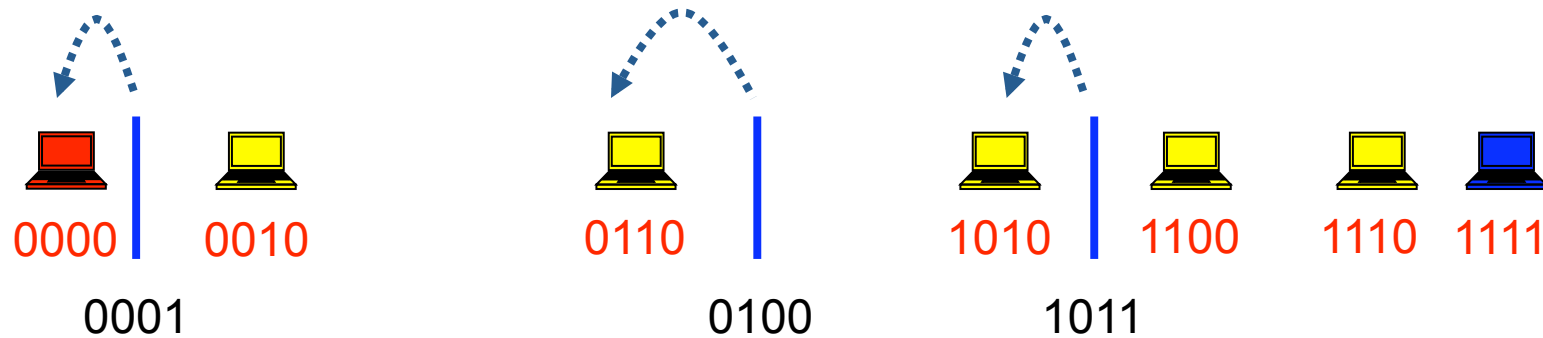
- Consider problem of data partition:
 - Given document X , choose one of k servers to use
- Suppose we use modulo hashing
 - Number servers $1..k$
 - Place X on server $i = (X \bmod k)$
 - Problem? Data may not be uniformly distributed
 - Place X on server $i = \text{hash}(X) \bmod k$
 - Problem?
 - What happens if a server fails or joins ($k \rightarrow k \pm 1$)?
 - What is different clients has different estimate of k ?
 - Answer: All entries get remapped to new nodes!

Consistent Hashing



- Consistent hashing partitions key-space among nodes
- Contact appropriate node to lookup/store key
 - Blue node determines red node is responsible for key_1
 - Blue node sends lookup or insert to red node

Consistent Hashing



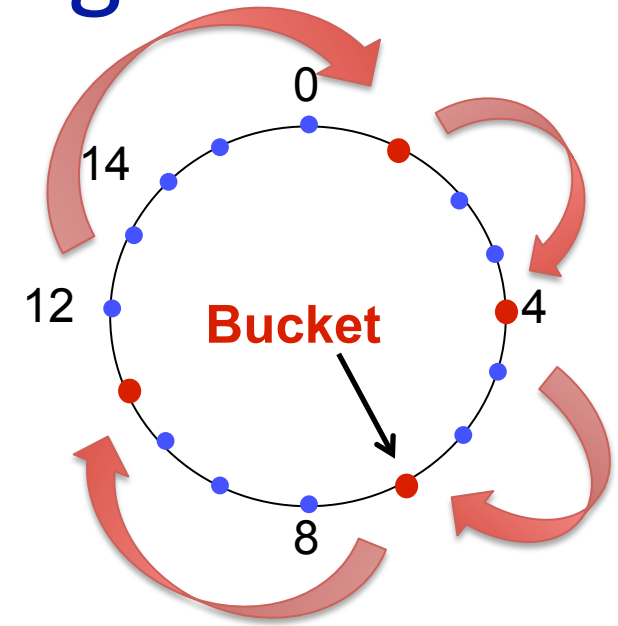
- Partitioning key-space among nodes

- Nodes choose random identifiers: e.g., `hash(IP)`
- Keys randomly distributed in ID-space: e.g., `hash(URL)`
- Keys assigned to node “nearest” in ID-space
- Spreads ownership of keys evenly across nodes

Consistent Hashing

- **Construction**

- Assign C hash buckets to random points on mod 2^n circle; hash key size = n
- Map object to random position on circle
- Hash of object = closest clockwise bucket

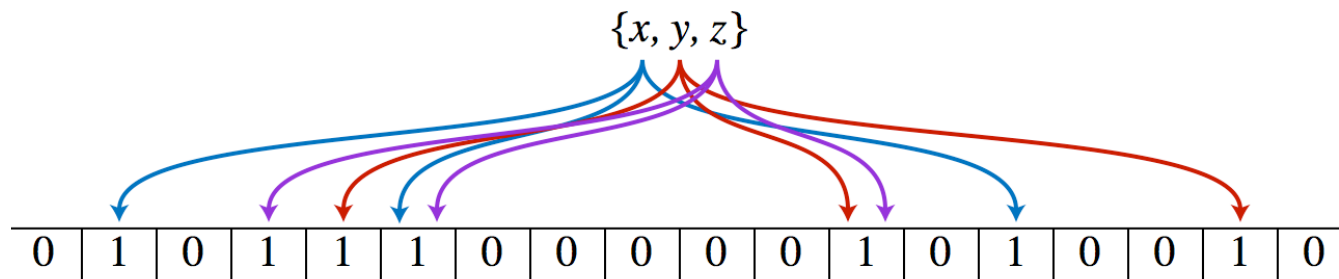


- **Desired features**

- **Balanced:** No bucket has disproportionate number of objects
- **Smoothness:** Addition/removal of bucket does not cause movement among existing buckets (only immediate buckets)
- **Spread and load:** Small set of buckets that lie near object

Bloom Filters

- **Data structure for probabilistic membership testing**
 - Small amount of space, constant time operations
 - False positives possible, no false negatives
 - Useful in per-flow network statistics, sharing information between cooperative caches, etc.
- **Basic idea using hash fn's and bit array**
 - Use k independent hash functions to map item to array
 - If all array elements are 1, it's present. Otherwise, not



Bloom Filters

Start with an m bit array, filled with 0s.

0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

To insert, hash each item k times. If $H_i(x) = a$, set $Array[a] = 1$.

0	1	0	0	1	0	1	0	0	1	1	1	0	1	1	0
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

To check if y is in set, check array at $H_i(y)$. All k values must be 1.

0	1	0	0	1	0	1	0	0	1	1	1	0	1	1	0
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

Possible to have a false positive: all k values are 1, but y is not in set.

0	1	0	0	1	0	1	0	0	1	1	1	0	1	1	0
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

Today's outline

- **Uses of hashing**
 - Equal-cost multipath routing in switches
 - Network load balancing in server clusters
 - Per-flow statistics in switches (QoS, IDS)
 - Caching in cooperative CDNs and P2P file sharing
 - Data partitioning in distributed storage services
- **Various hashing strategies**
 - Modulo hashing
 - Consistent hashing
 - Bloom Filters