

Consistent Hashing

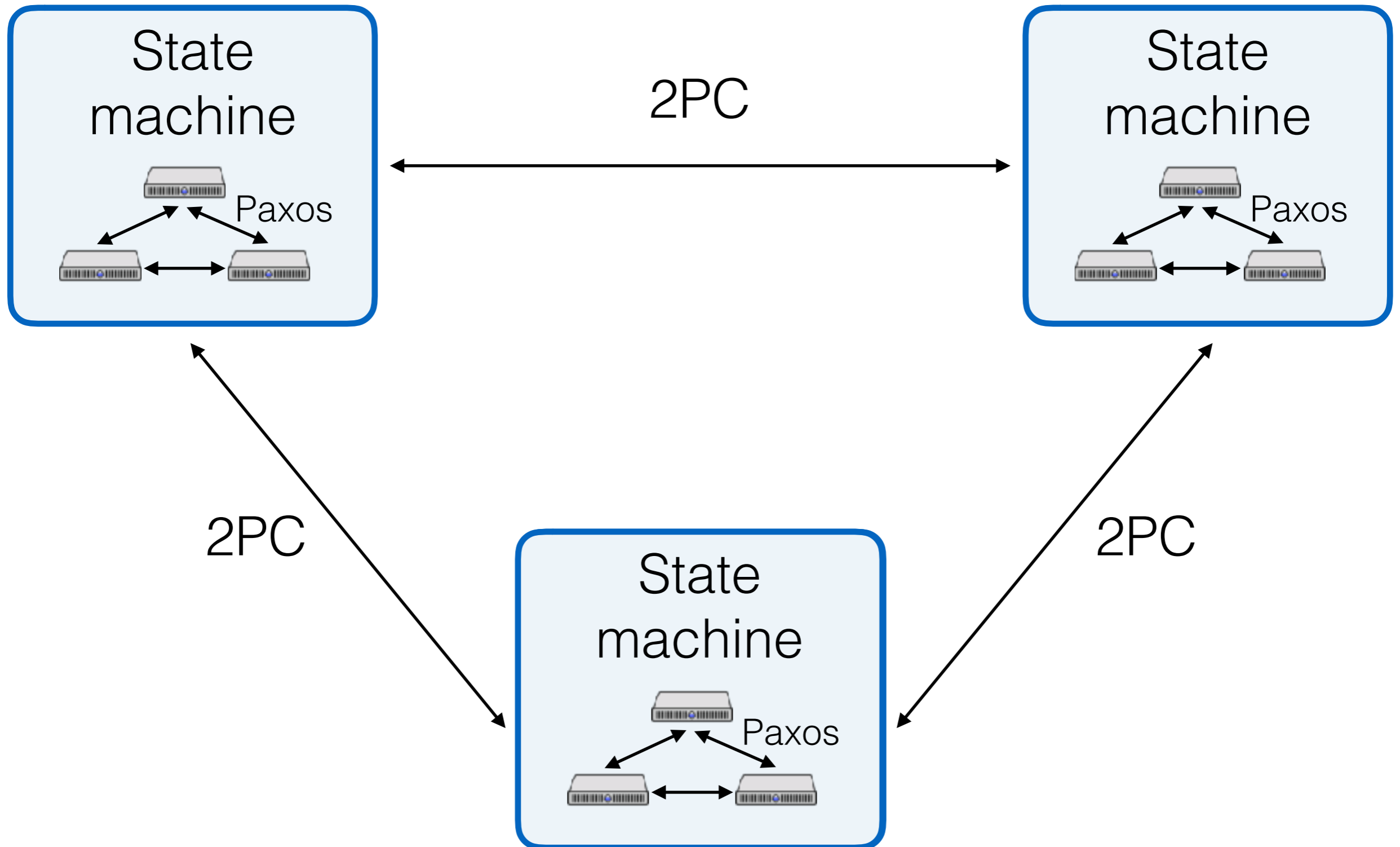
Doug Woos

Logistics notes

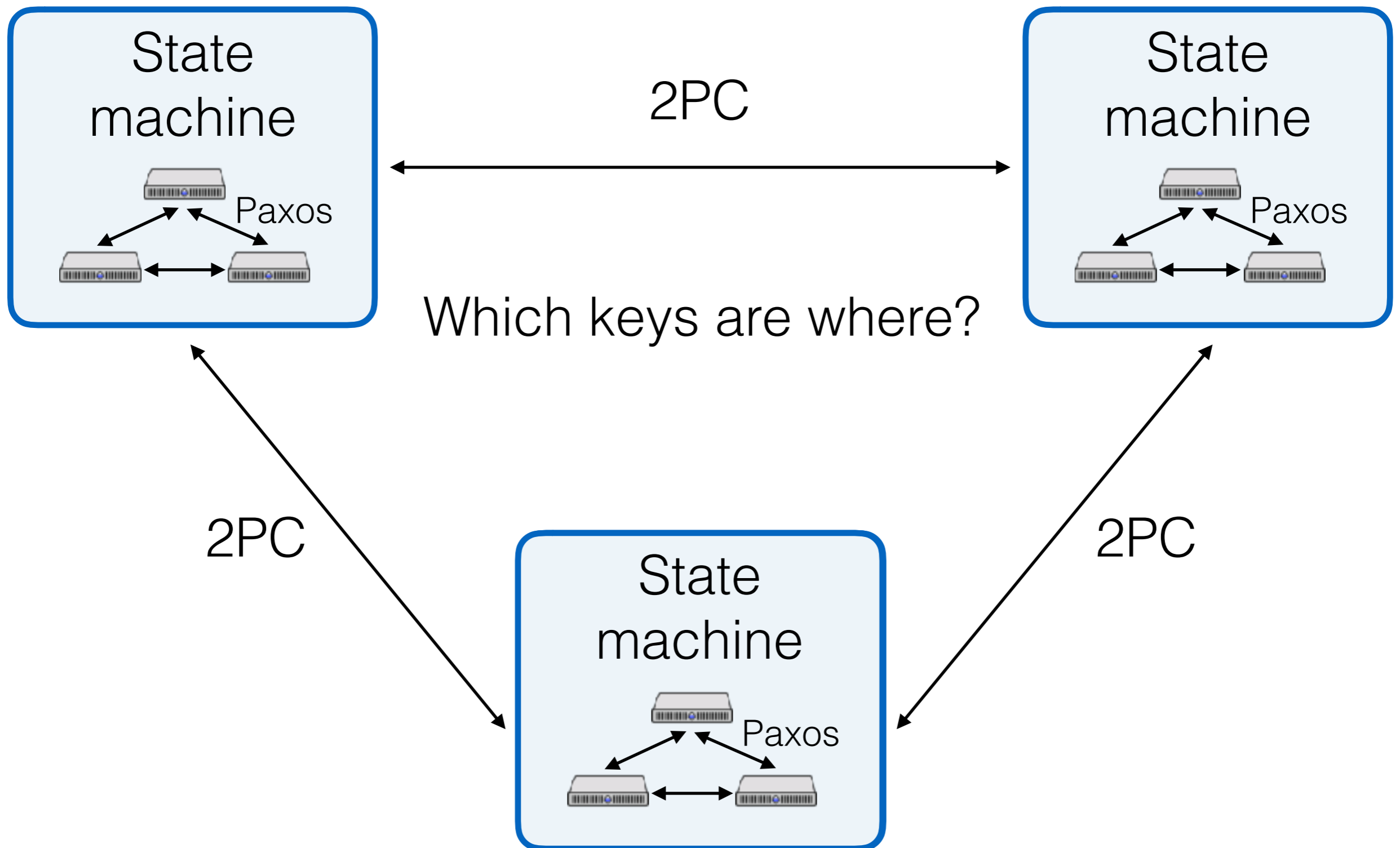
PS3 out later today

Lab 3a due Monday

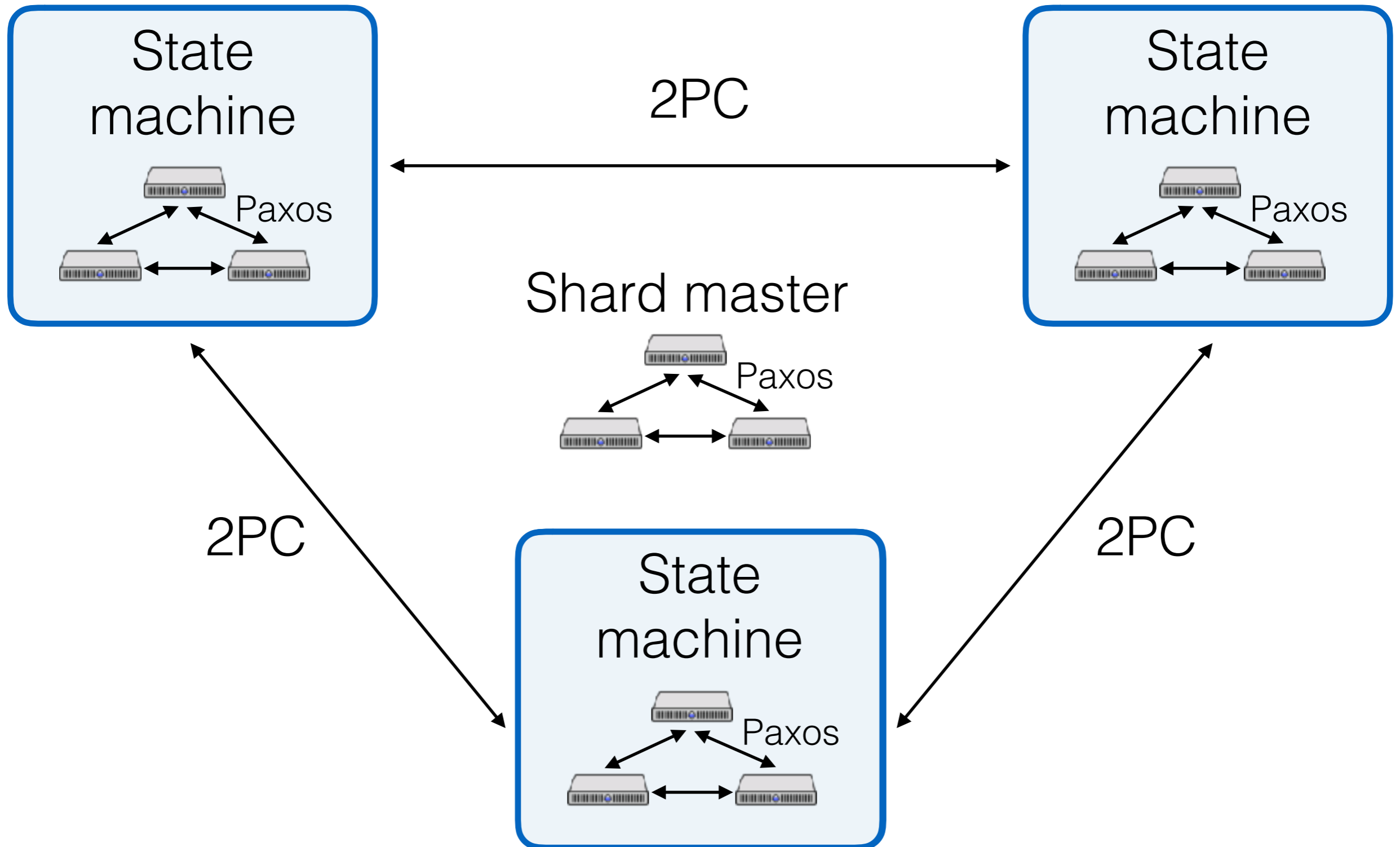
Replicated, Sharded Database



Replicated, Sharded Database



Lab 4 (and others!)



Highly available stores

Want to avoid communication with shard master

Shards operate independently

So: how to ensure clients know who to talk to?

Another scenario



View editing options | mail | write | like

Doug Woos



Ph.D. (2008) and M.S. (2006) in Computer Science at the University of Washington, where I was advised by [Tom Anderson](#), [John Storm](#), and [Dan Zeng](#). My research interests are in distributed algorithms, data structures, and the theory of networks. My other interests include philosophy, science fiction, and recreational mathematics. My research has been published in *SIAM Journal on Computing*, *SIAM Review*, *SIAM Journal on Applied Mathematics*, *SIAM Journal on Algebraic Combinatorics*, *SIAM Journal on Matrix Analysis and Applications*, and *SIAM Journal on Numerical Analysis*.

My work at INRIA has been primarily supported by the [ANR-11-01-0001-01](#) grant.

The views expressed in articles and technical publications are those of the author and do not necessarily reflect those of INRIA. More generally, the views in this blog are those of the author and do not necessarily reflect those of INRIA.

Advising

I've advised the following students and interns at INRIA:

- [Alexandre Gallet](#)
- [Sébastien Goulet](#)

Publications

Client



Another scenario



LinkedIn profile for Doug Woos, showing a profile picture, name, and a bio. The bio mentions his education at the University of Washington and his research interests in distributed systems and algorithms. It also lists his current position at Amazon and his previous roles at Microsoft and Google.

GET index.html

Client



Another scenario



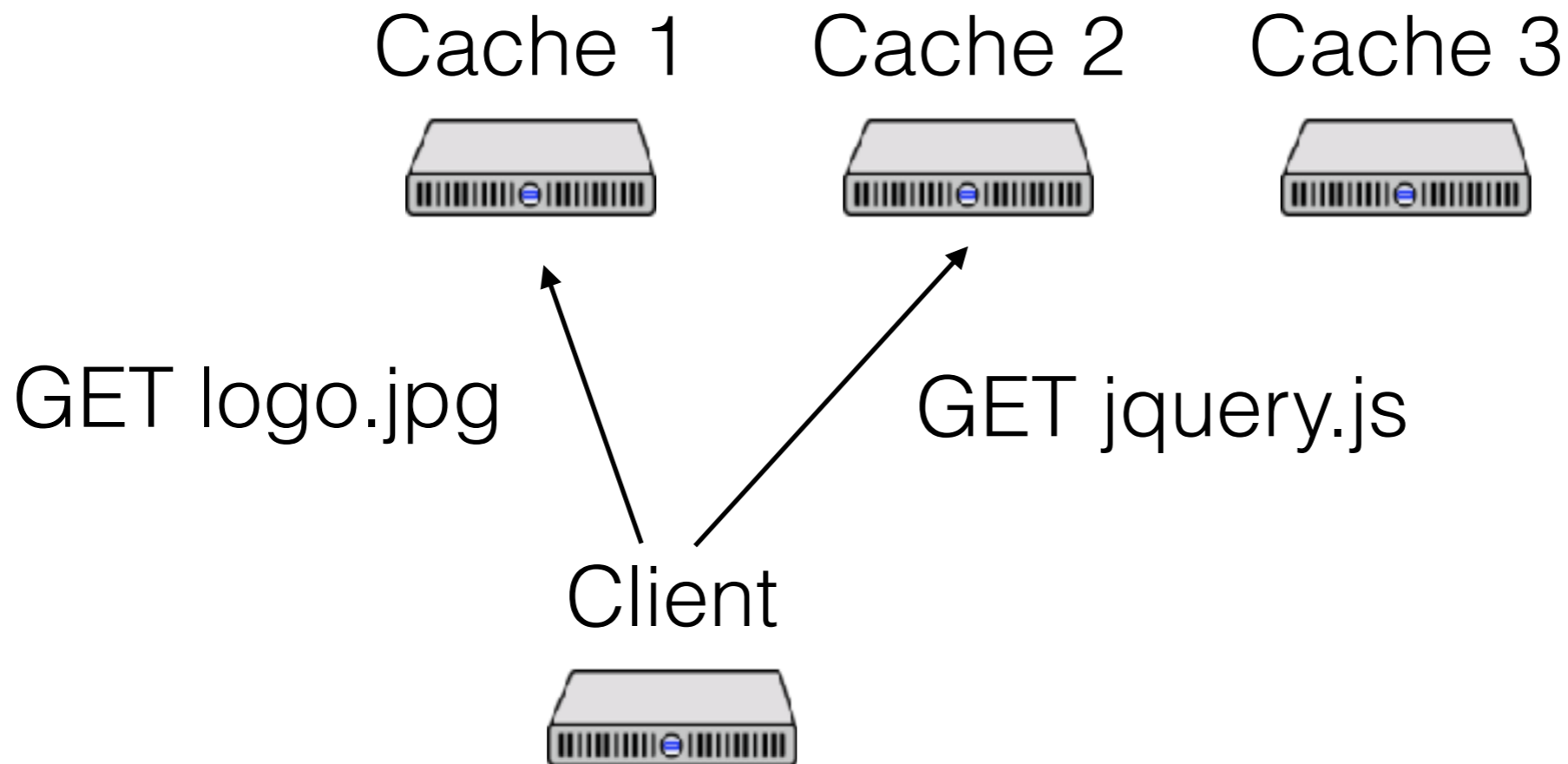
index.html

Links to: logo.jpg, jquery.js, ...

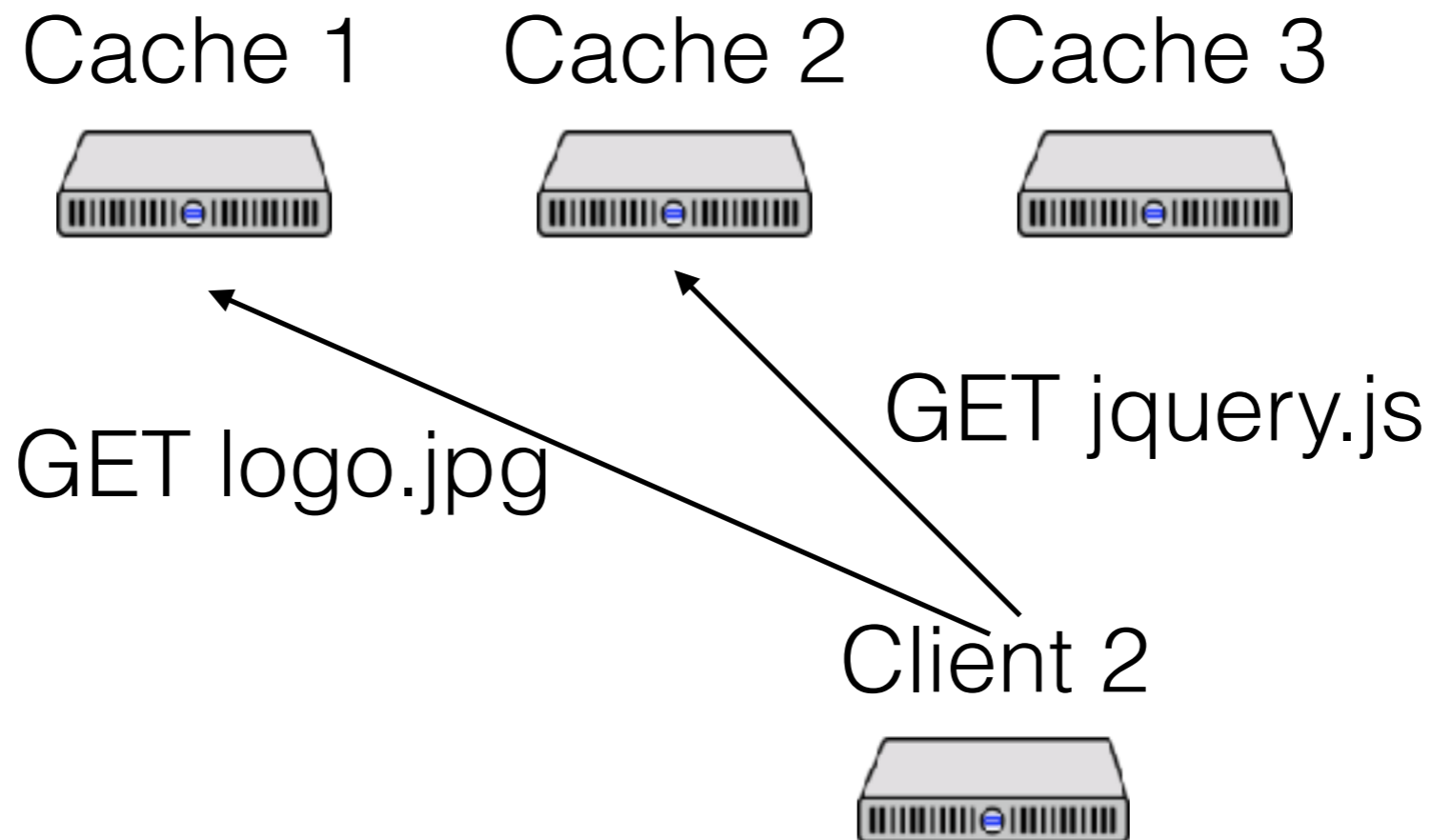
Client



Another scenario



Another scenario



What's in common?

Want to assign keys to servers w/o communication

Requirement 1: clients all have same assignment

Proposal 1

For n nodes, a key k goes to $k \bmod n$

Cache 1



“a”, “d”, “ab”

Cache 2



“b”

Cache 3



“c”

Proposal 1

For n nodes, a key k goes to $k \bmod n$

Cache 1



“a”, “d”, “ab”

Cache 2



“b”

Cache 3



“c”

Problems with this approach?

Proposal 1

For n nodes, a key k goes to $k \bmod n$

Cache 1



“a”, “d”, “ab”

Cache 2



“b”

Cache 3



“c”

Problems with this approach?

- Likely to have distribution issues

Requirements, revisited

Requirement 1: clients all have same assignment

Requirement 2: keys uniformly distributed

Proposal 2

For n nodes, a key k goes to $hash(k) \bmod n$

Cache 1



Cache 2



Cache 3



$h("a")=1$

$h("abc")=2$

$h("b")=3$

Hash distributes keys uniformly

Proposal 2

For n nodes, a key k goes to $hash(k) \bmod n$

Cache 1



Cache 2



Cache 3



$$h("a")=1$$

$$h("abc")=2$$

$$h("b")=3$$

Hash distributes keys uniformly

But, new problem: what if we add a node?

Proposal 2

For n nodes, a key k goes to $hash(k) \bmod n$

Cache 1



Cache 2



Cache 3



Cache 4



$$h("a")=1$$

$$h("abc")=2$$

$$h("b")=3$$

Hash distributes keys uniformly

But, new problem: what if we add a node?

Proposal 2

For n nodes, a key k goes to $hash(k) \bmod n$

Cache 1



Cache 2



Cache 3



Cache 4



$$h(\text{"abc"})=2 \quad h(\text{"a"})=3 \quad h(\text{"b"})=3$$

Hash distributes keys uniformly

But, new problem: what if we add a node?

Proposal 2

For n nodes, a key k goes to $hash(k) \bmod n$

Cache 1



Cache 2



Cache 3



Cache 4



$h(\text{"abc"})=2$ $h(\text{"a"})=3$ $h(\text{"b"})=4$

Hash distributes keys uniformly

But, new problem: what if we add a node?

- Redistribute a lot of keys! (on average, all but K/n)

Requirements, revisited

Requirement 1: clients all have same assignment

Requirement 2: keys uniformly distributed

Requirement 3: can add/remove nodes w/o redistributing too many keys

Proposal 3

First, hash the node ids

Proposal 3

First, hash the node ids

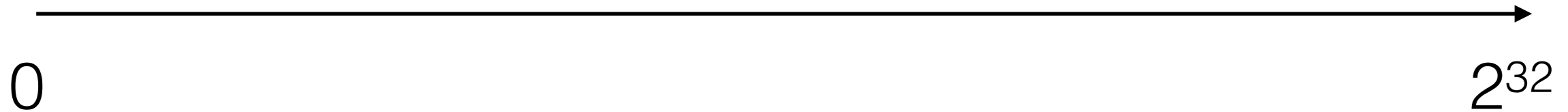
Cache 1



Cache 2

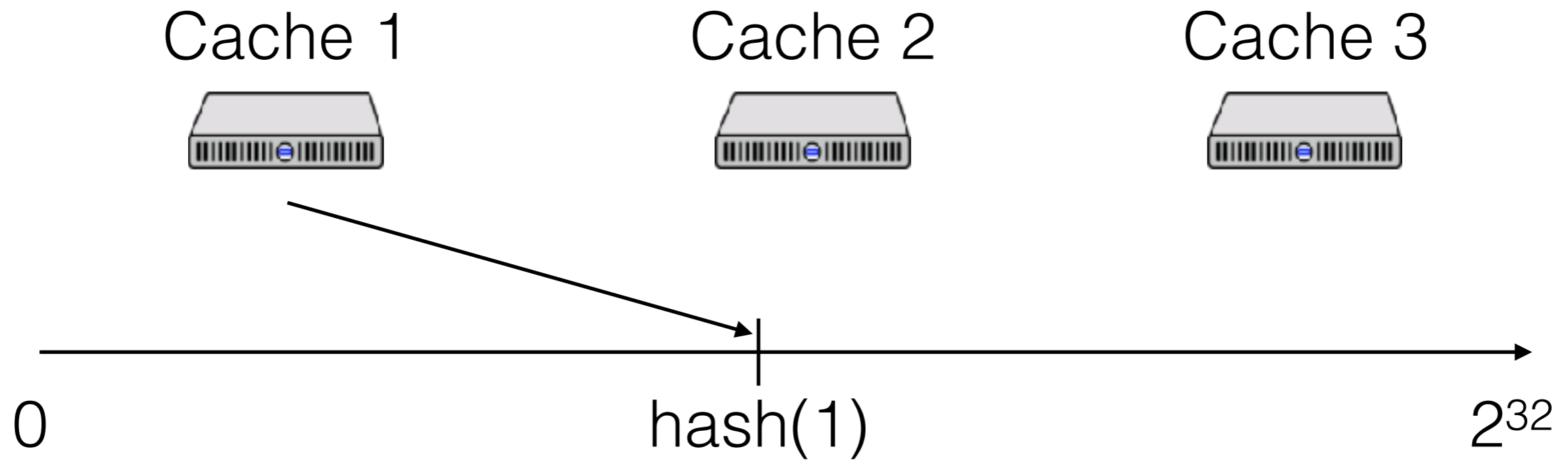


Cache 3



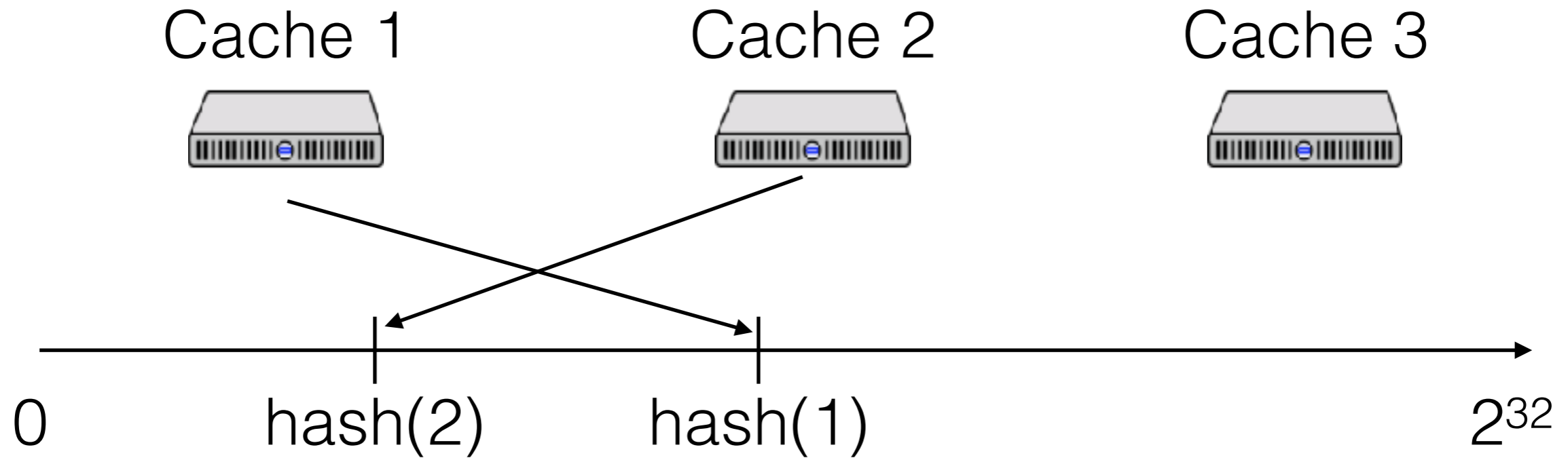
Proposal 3

First, hash the node ids



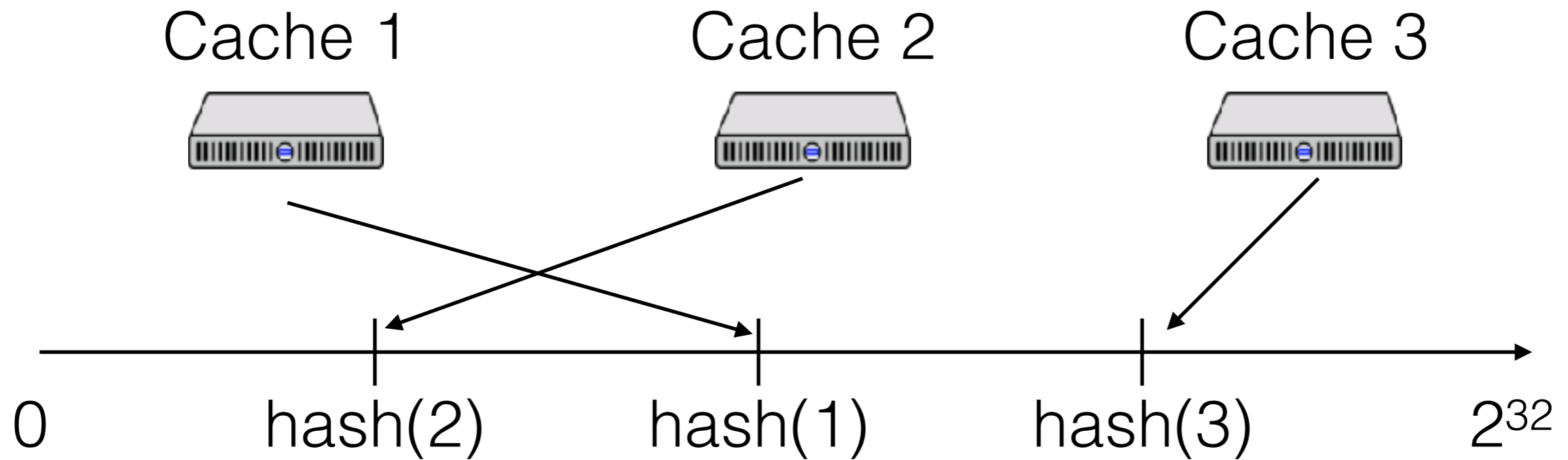
Proposal 3

First, hash the node ids



Proposal 3

First, hash the node ids



Proposal 3

First, hash the node ids

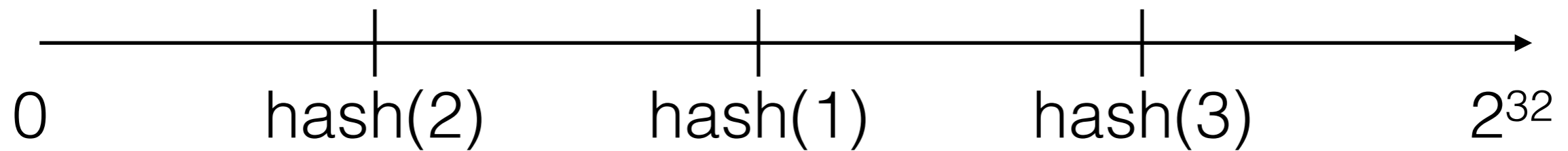
Cache 1



Cache 2

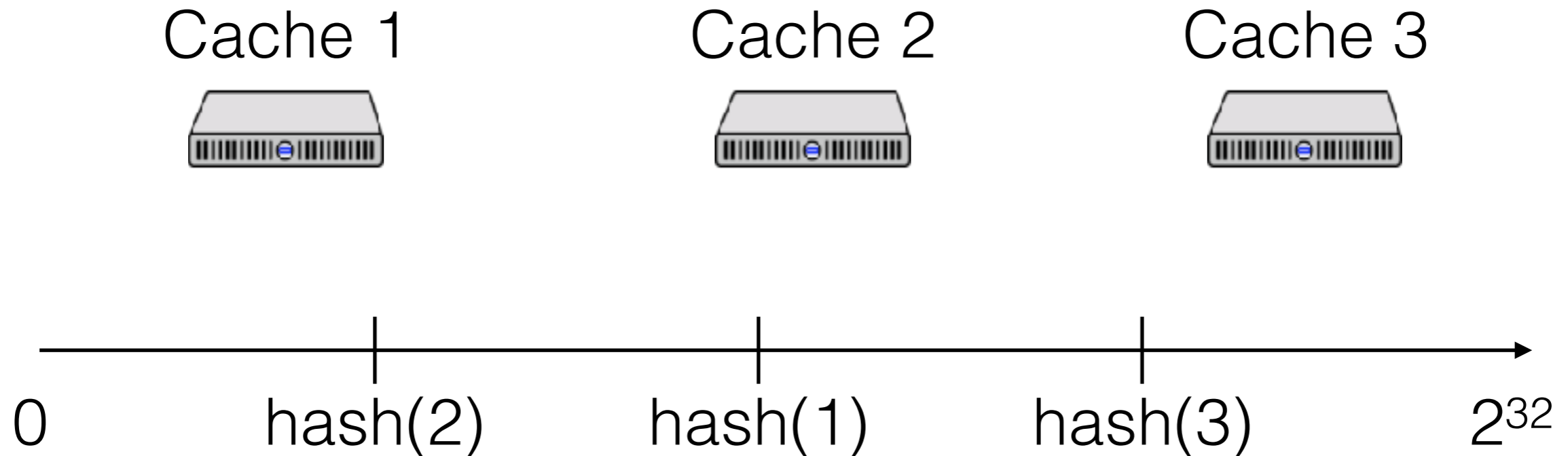


Cache 3



Proposal 3

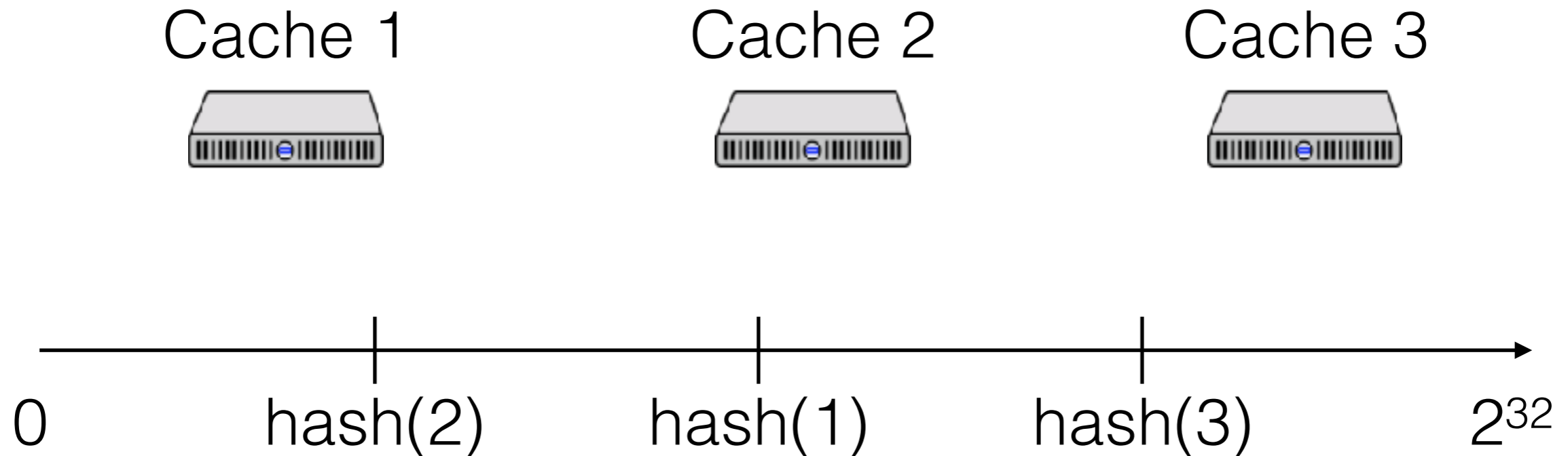
First, hash the node ids



Keys are hashed, then go to the “next” node

Proposal 3

First, hash the node ids

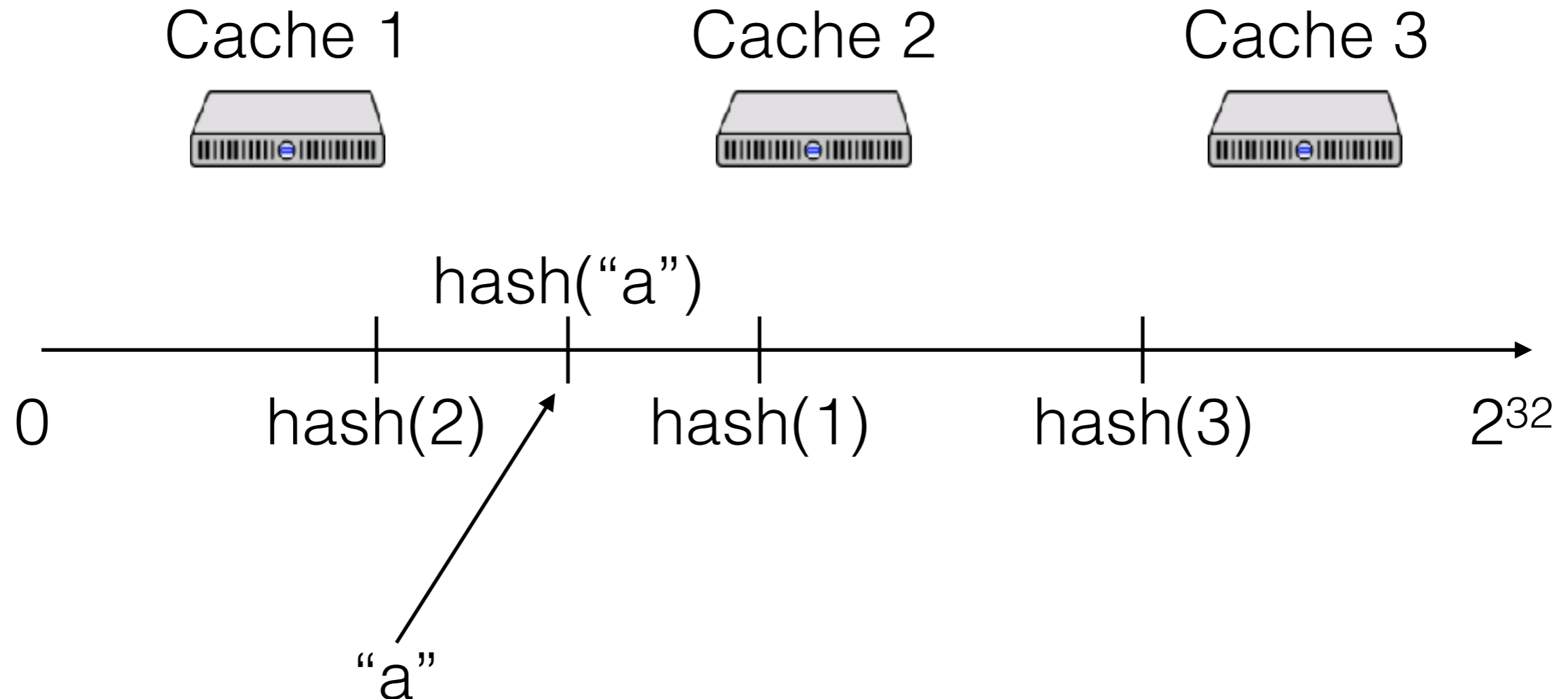


“a”

Keys are hashed, then go to the “next” node

Proposal 3

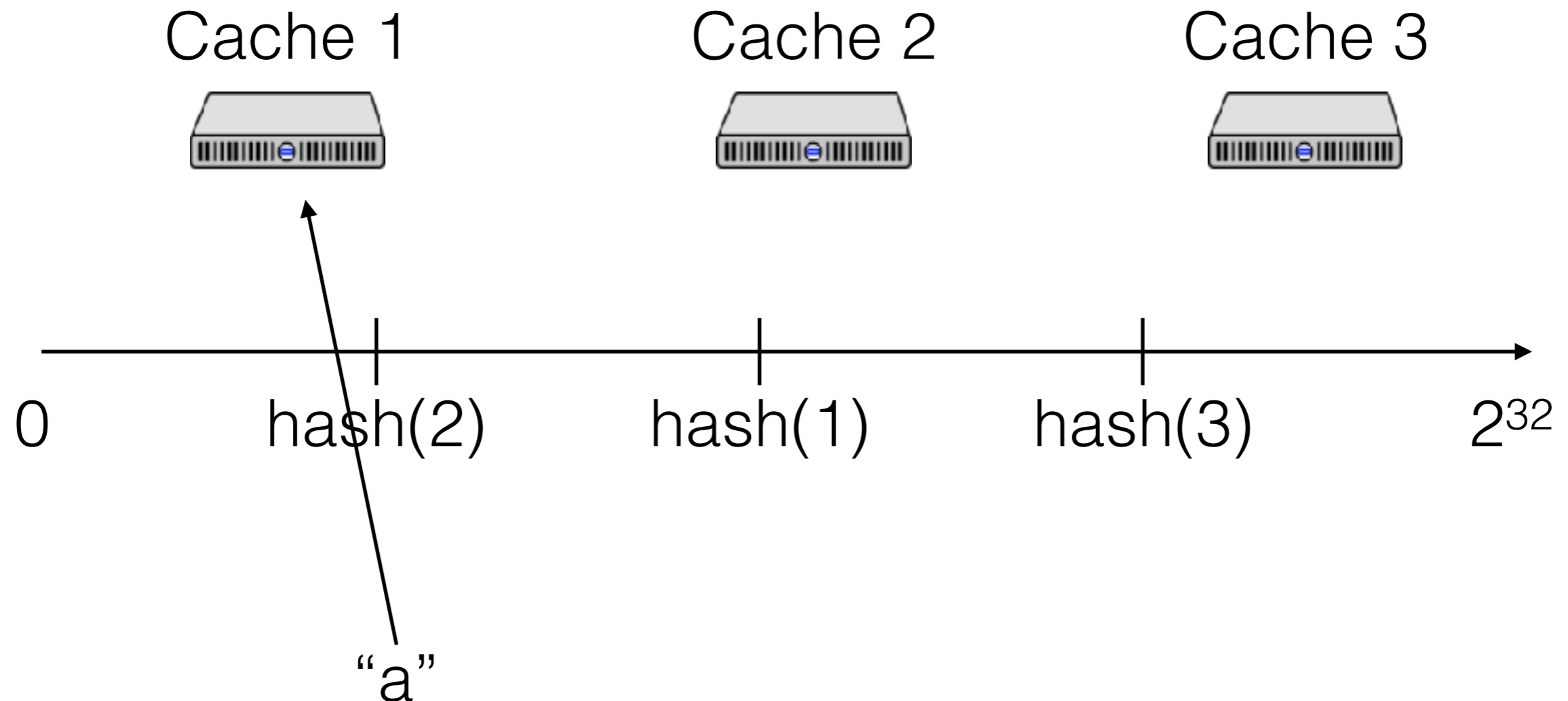
First, hash the node ids



Keys are hashed, then go to the "next" node

Proposal 3

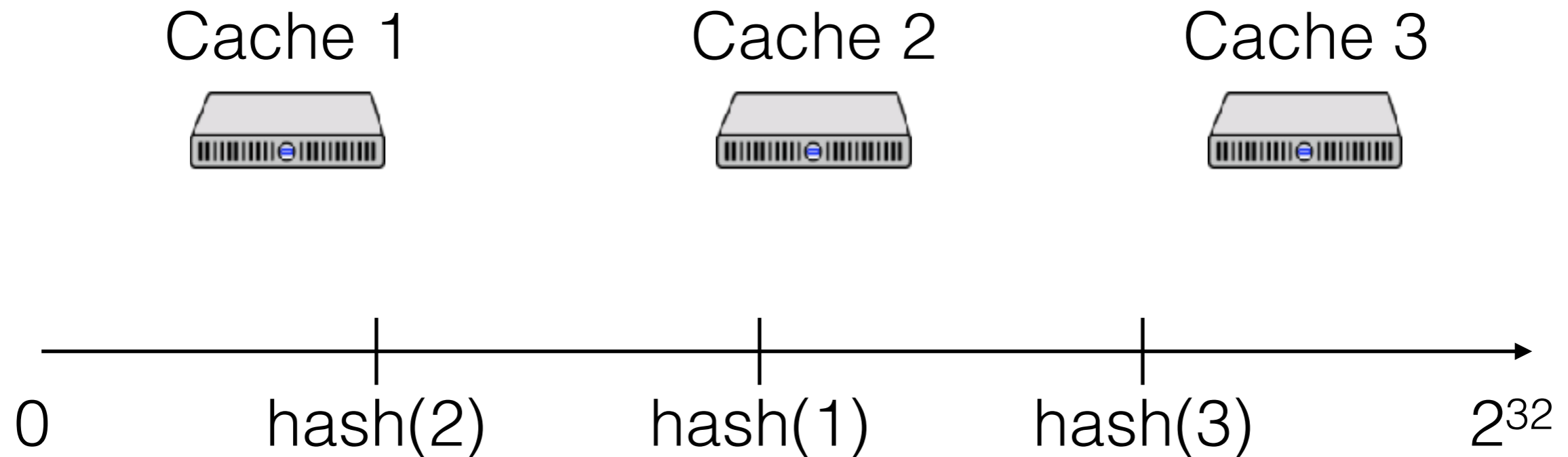
First, hash the node ids



Keys are hashed, then go to the "next" node

Proposal 3

First, hash the node ids

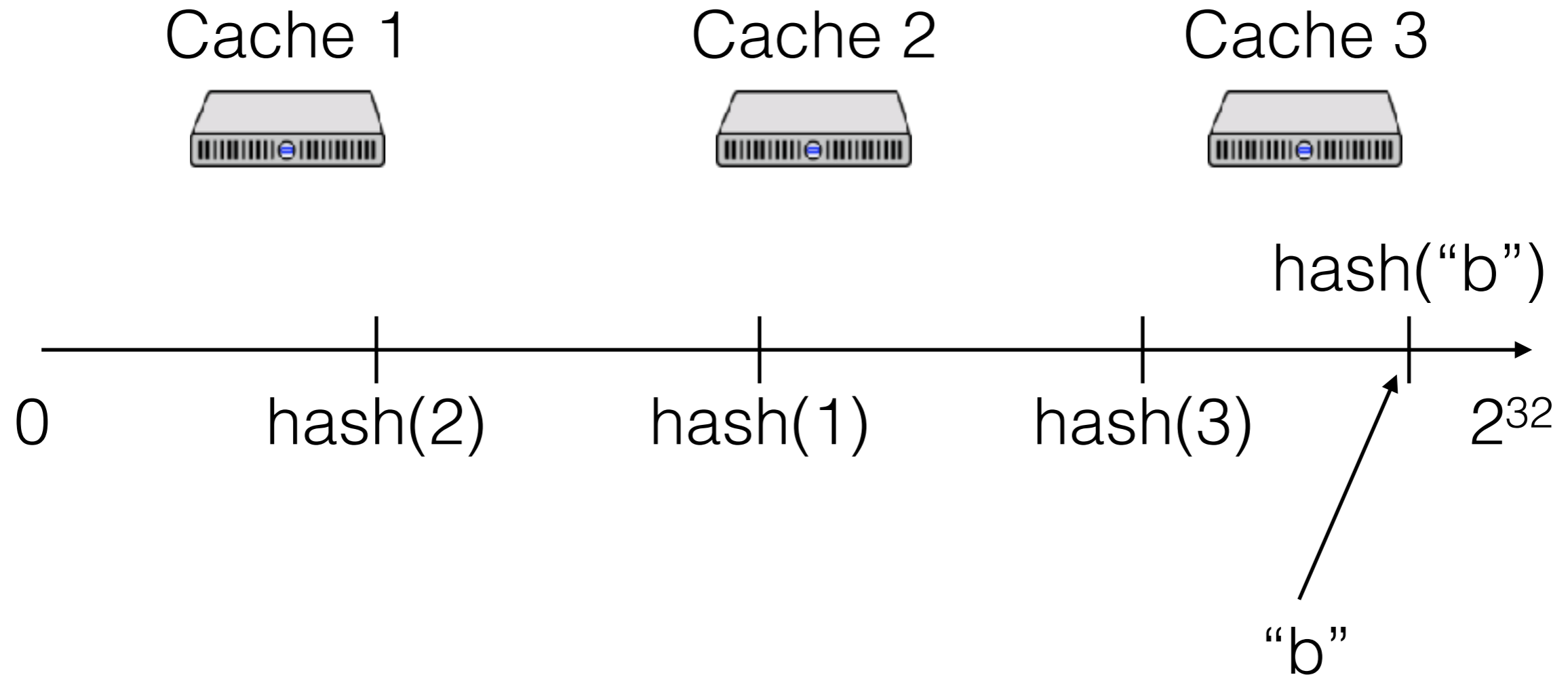


“b”

Keys are hashed, then go to the “next” node

Proposal 3

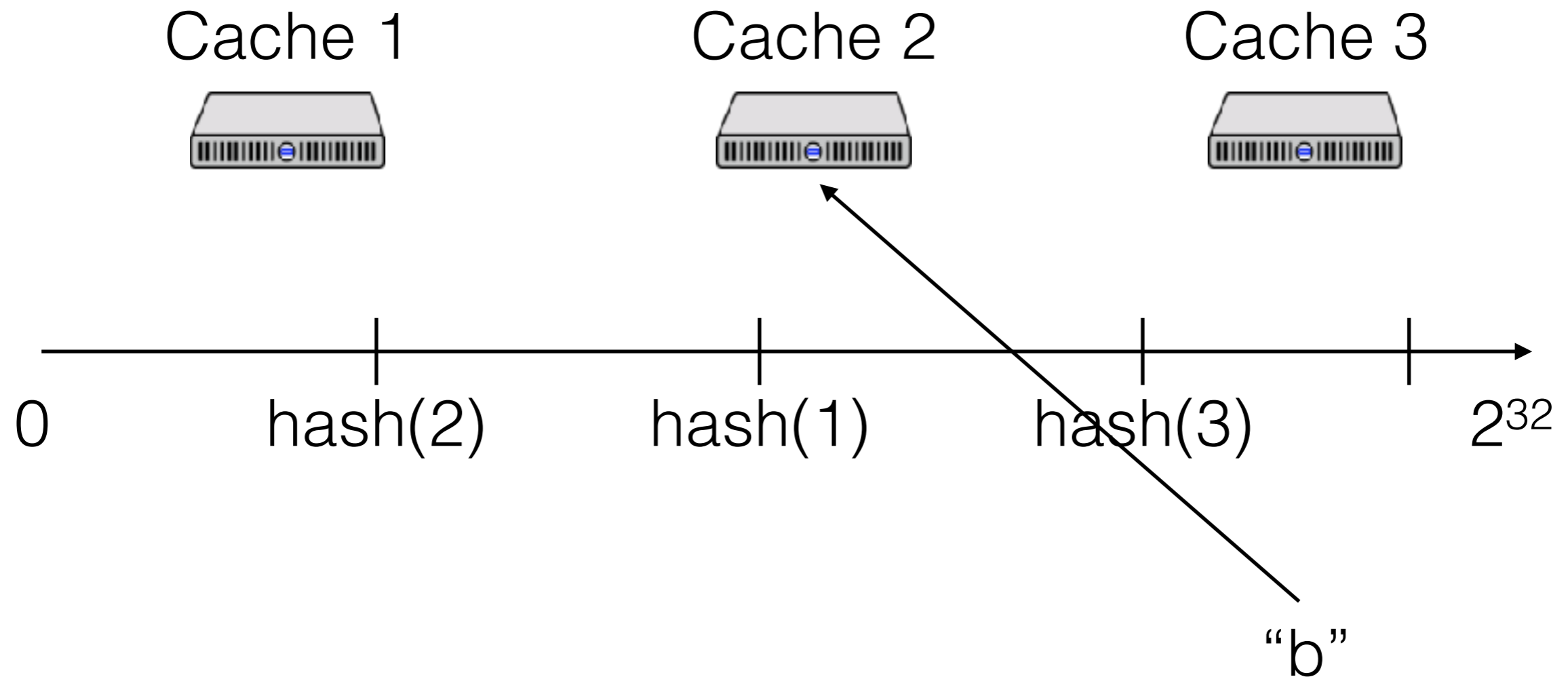
First, hash the node ids



Keys are hashed, then go to the "next" node

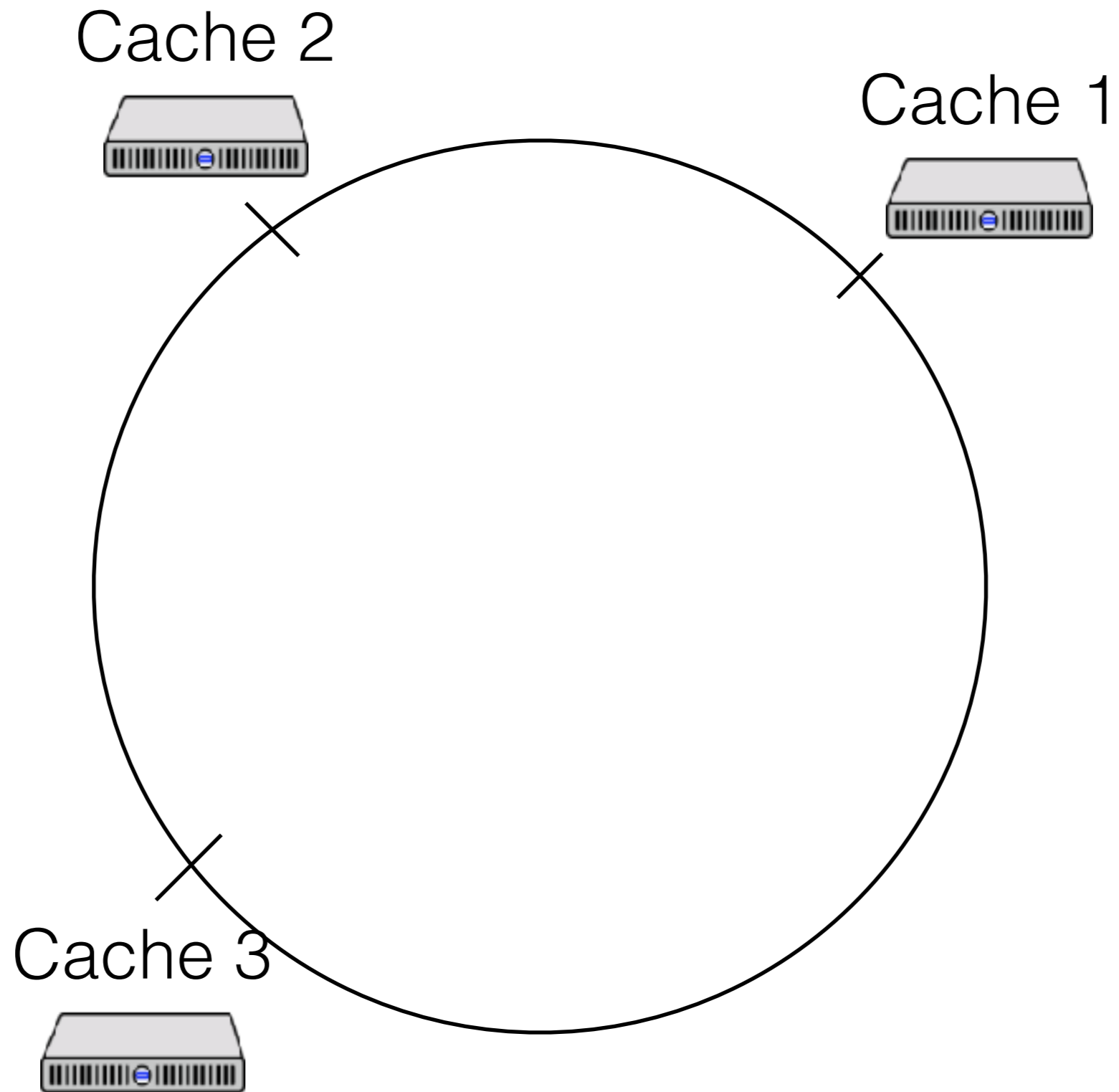
Proposal 3

First, hash the node ids

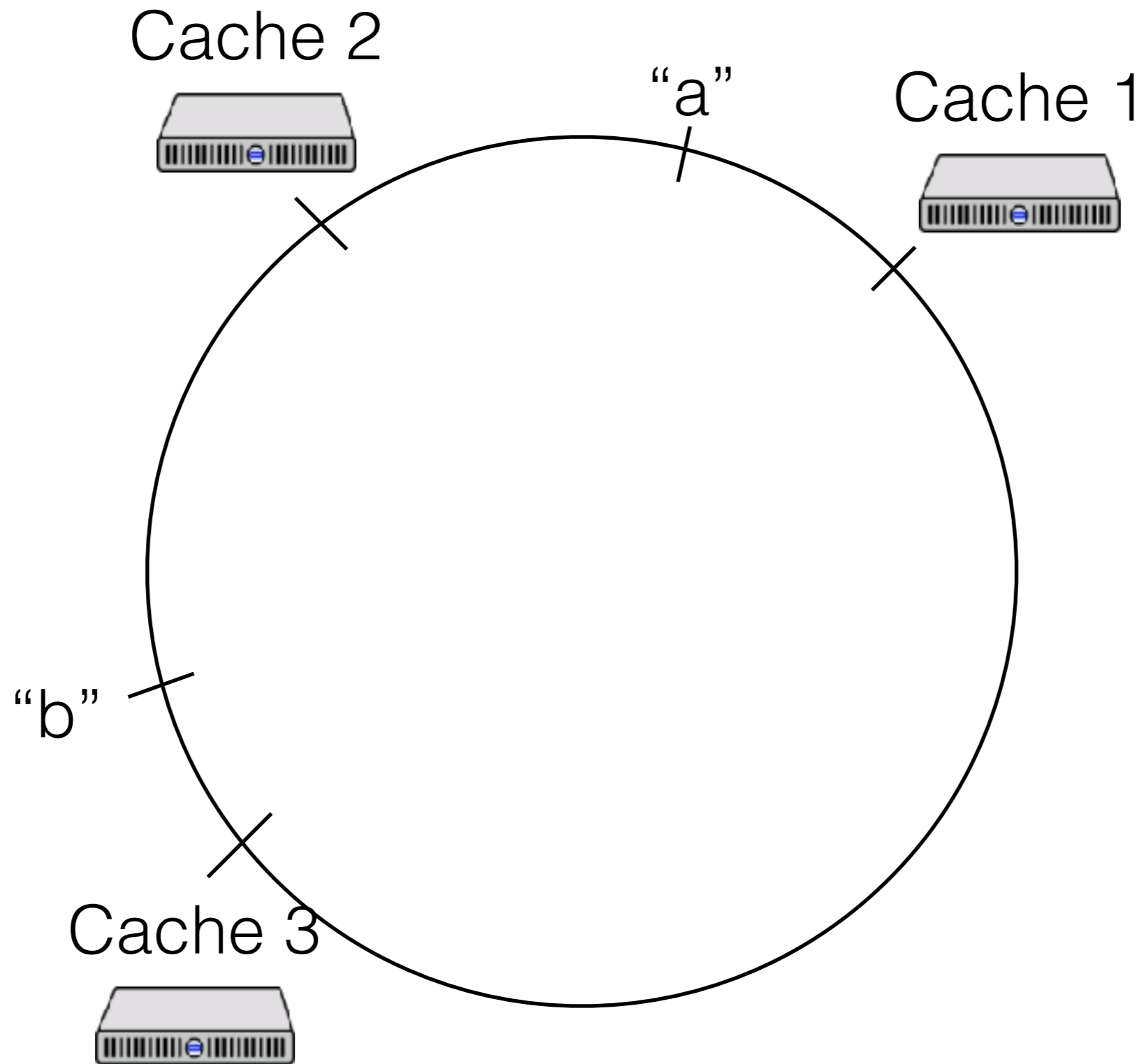


Keys are hashed, then go to the "next" node

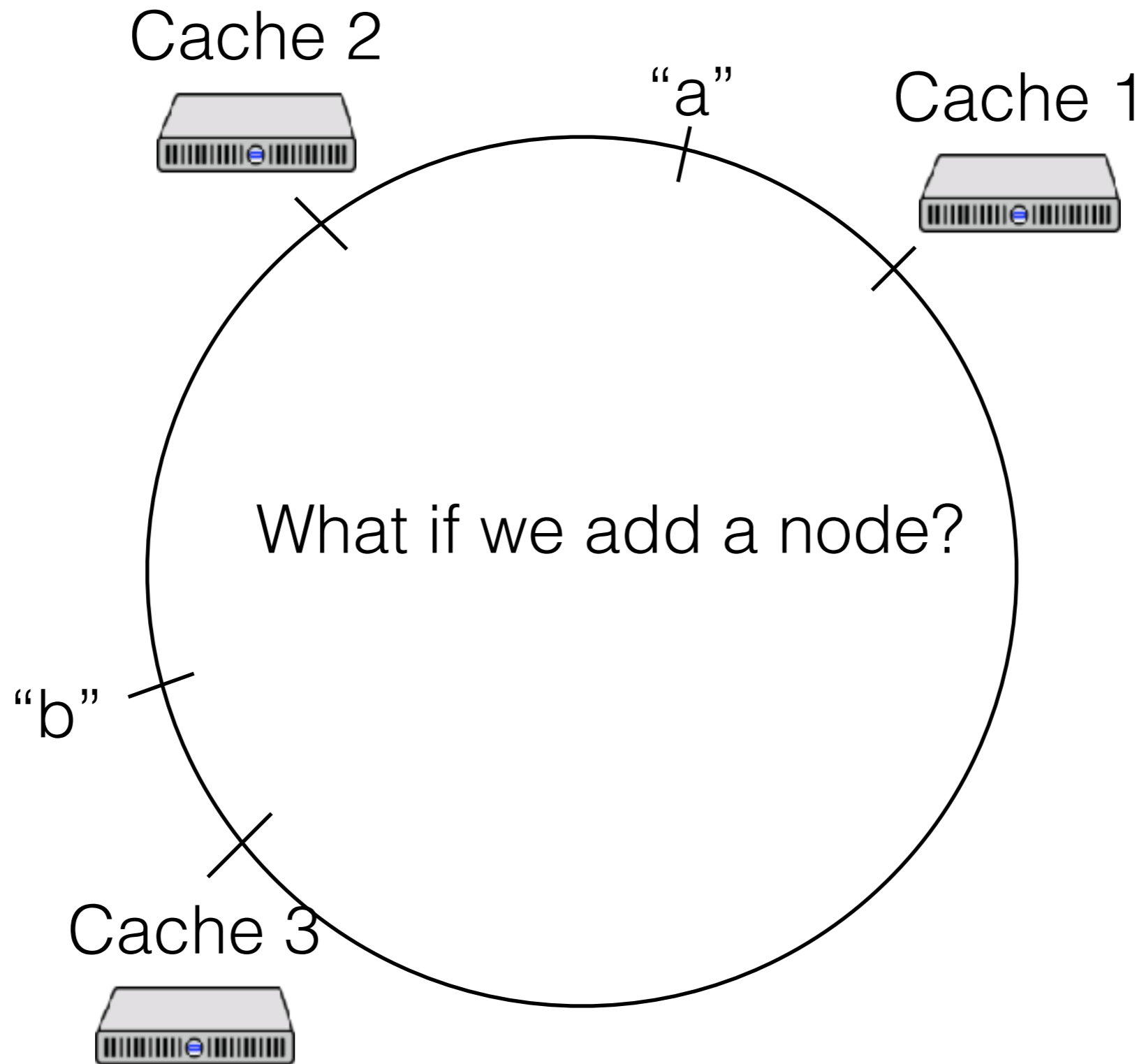
Proposal 3



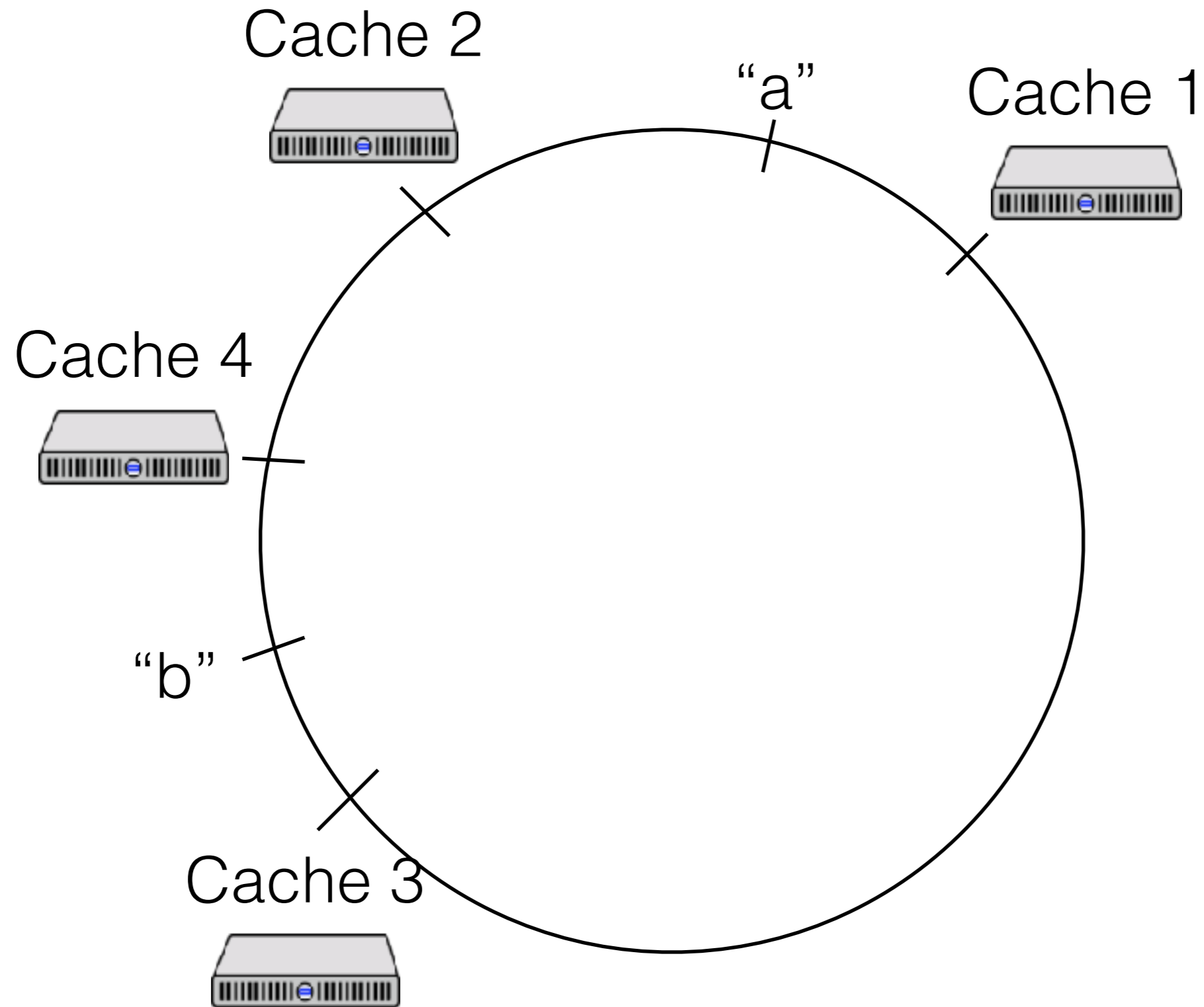
Proposal 3



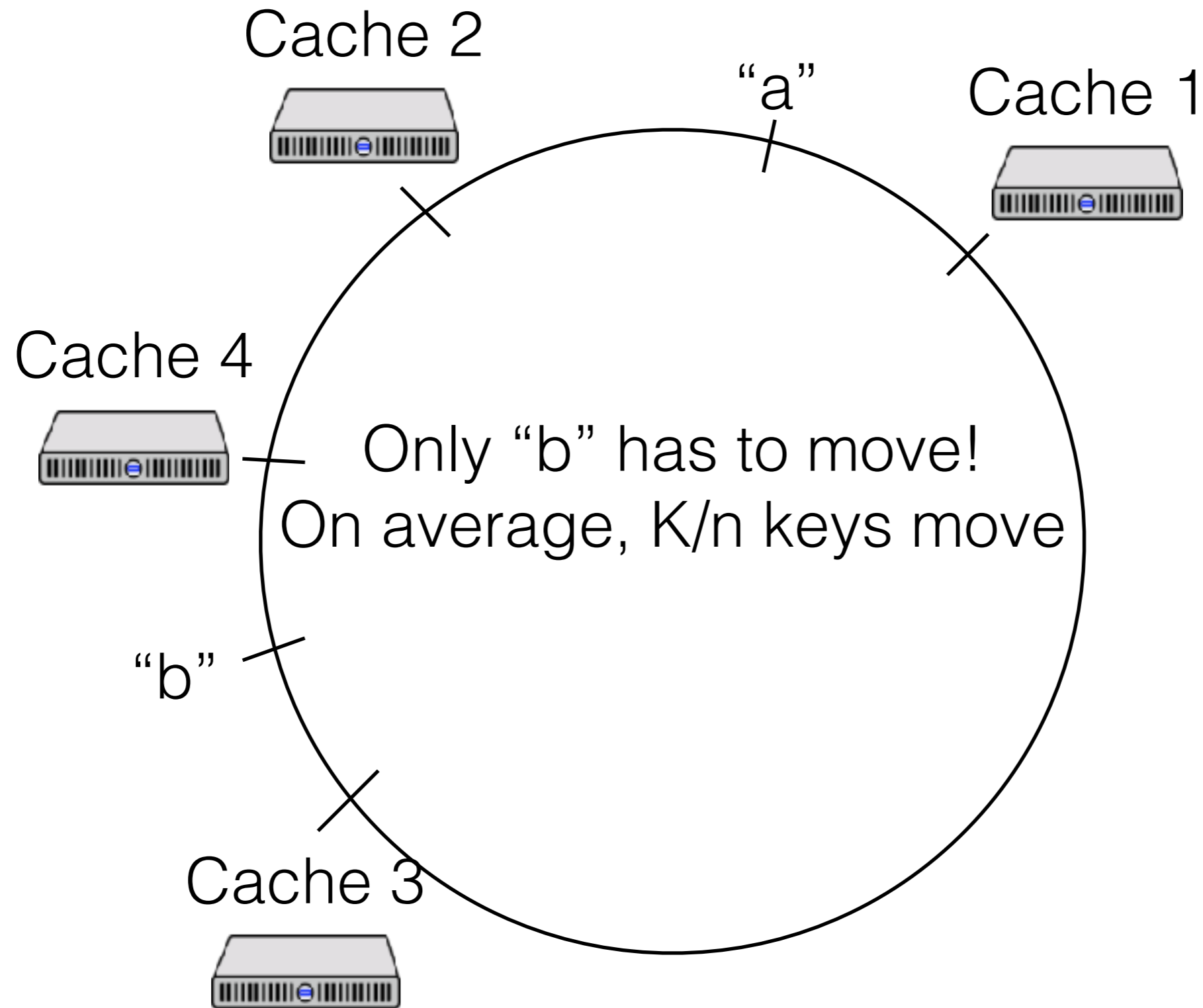
Proposal 3



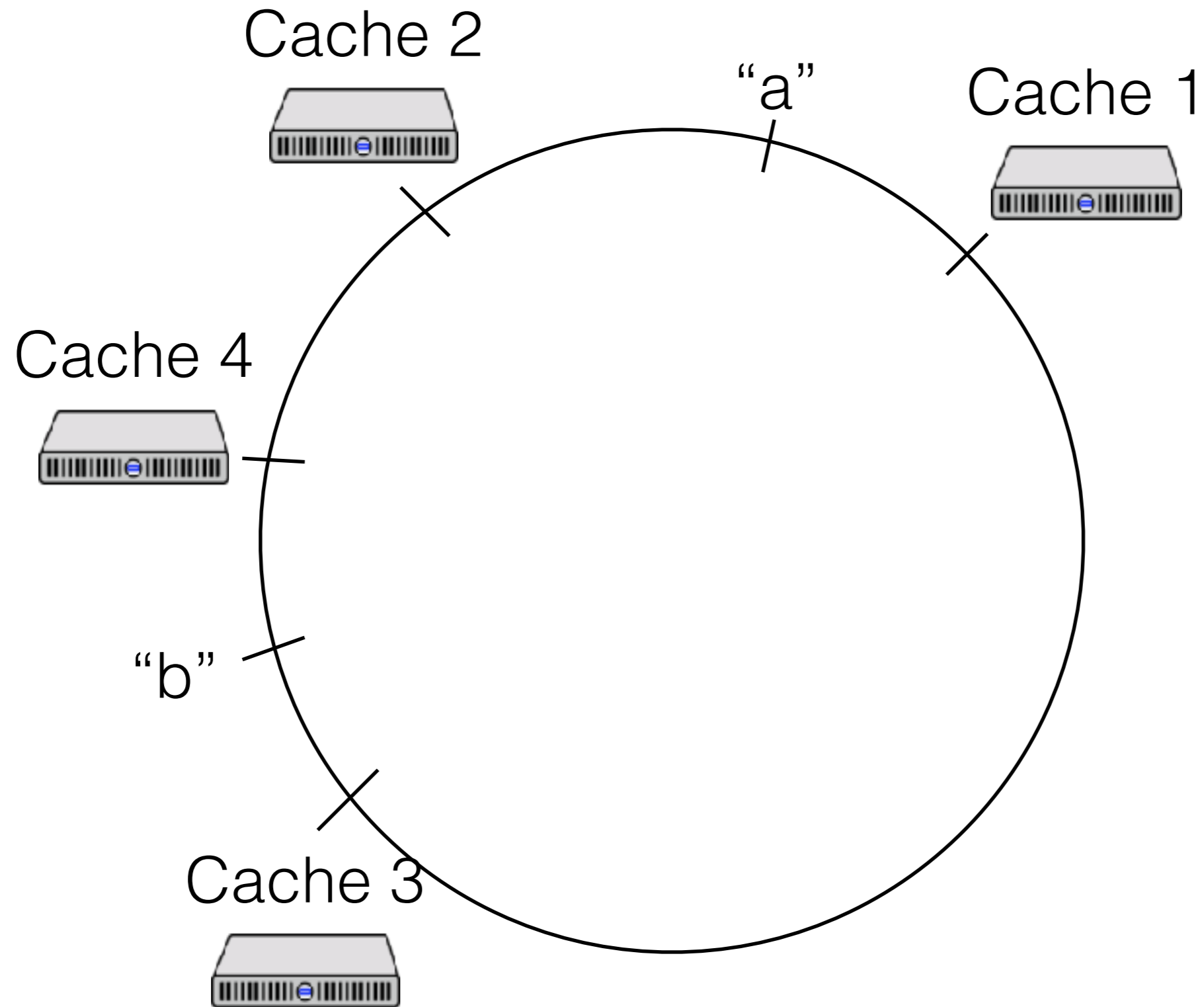
Proposal 3



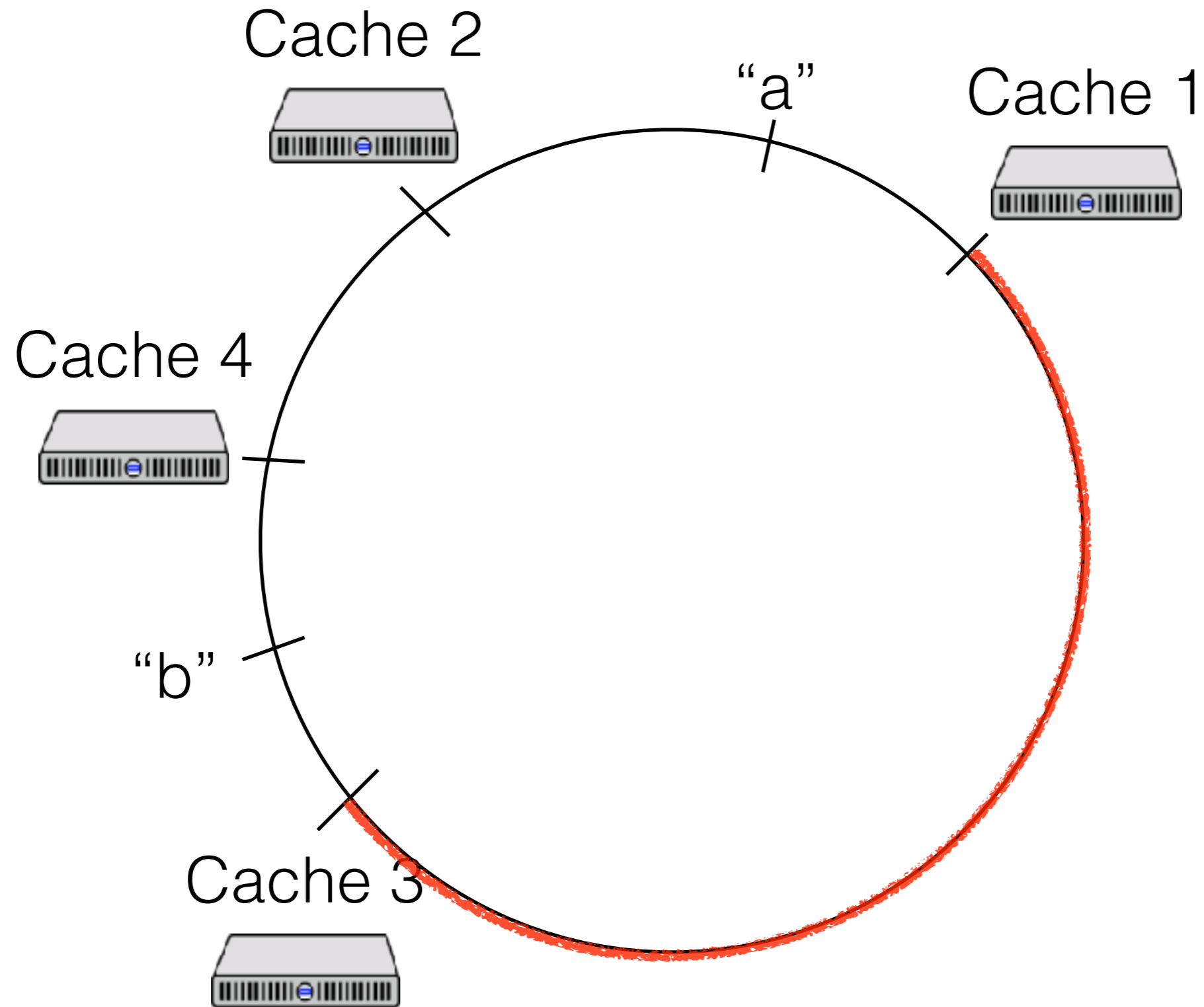
Proposal 3



Proposal 3



Proposal 3



Requirements, revisited

Requirement 1: clients all have same assignment

Requirement 2: keys uniformly distributed

Requirement 3: can add/remove nodes w/o redistributing too many keys

Proposal 4

First, hash the node ids to *multiple locations*

Cache 1



Cache 2

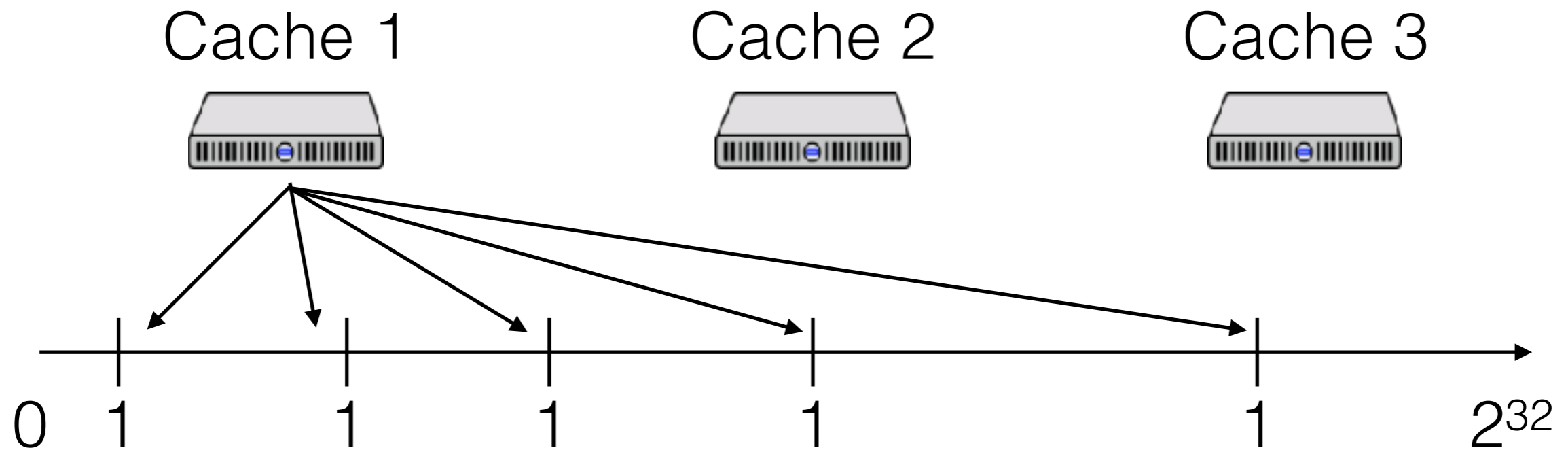


Cache 3



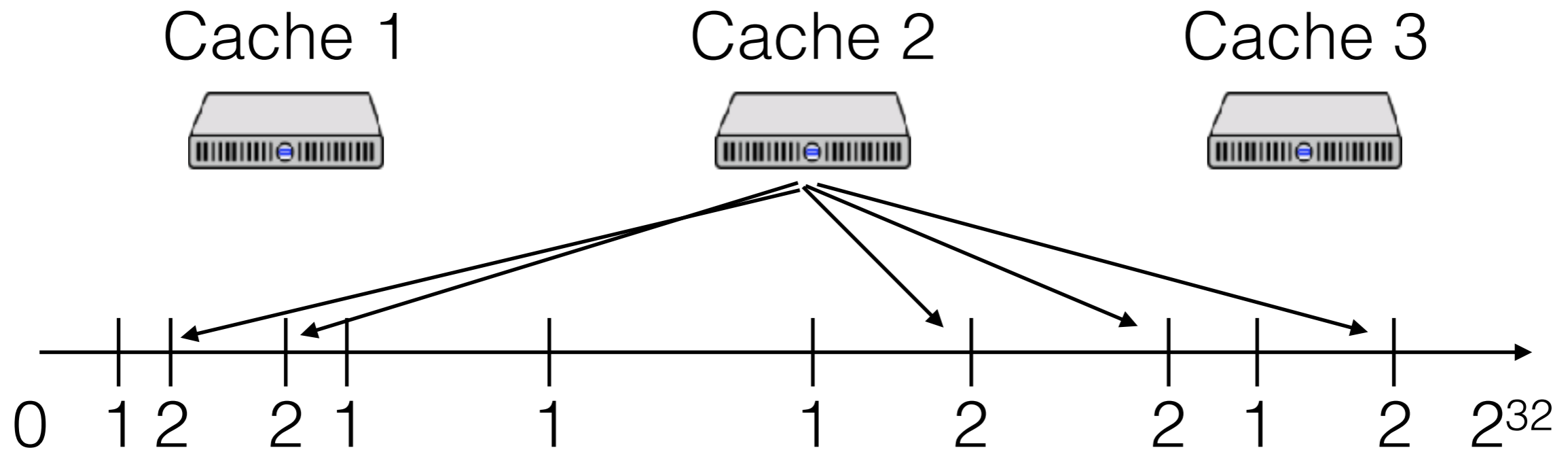
Proposal 4

First, hash the node ids to *multiple locations*



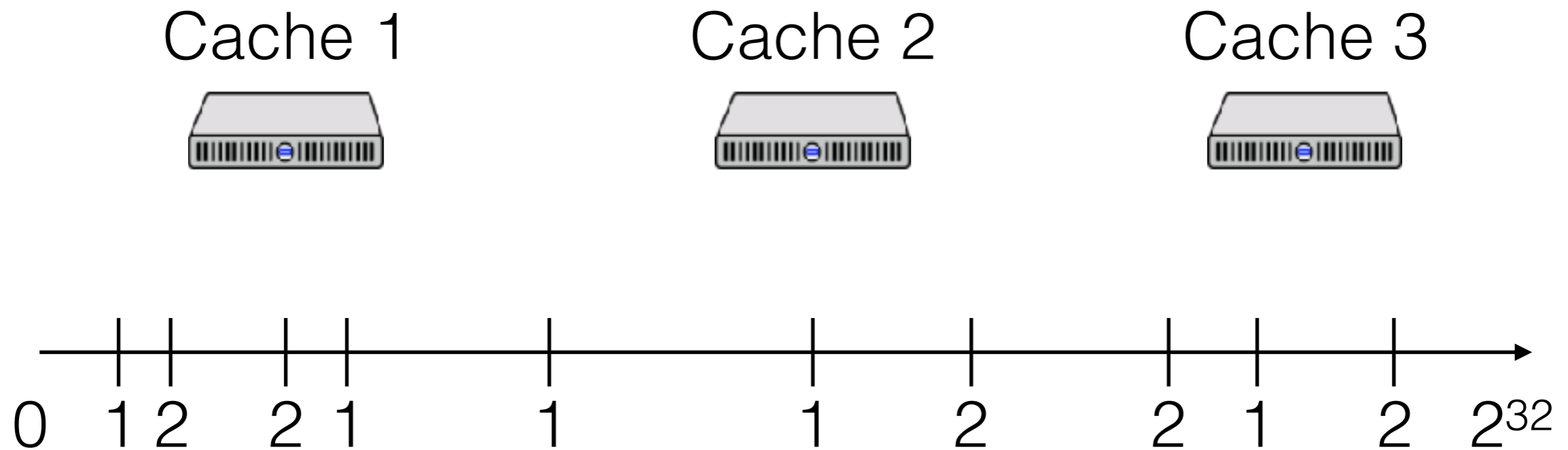
Proposal 4

First, hash the node ids to *multiple locations*



Proposal 4

First, hash the node ids to *multiple locations*



As it turns out, hash functions come in families s.t. their members are independent. So this is easy!

Proposal 4

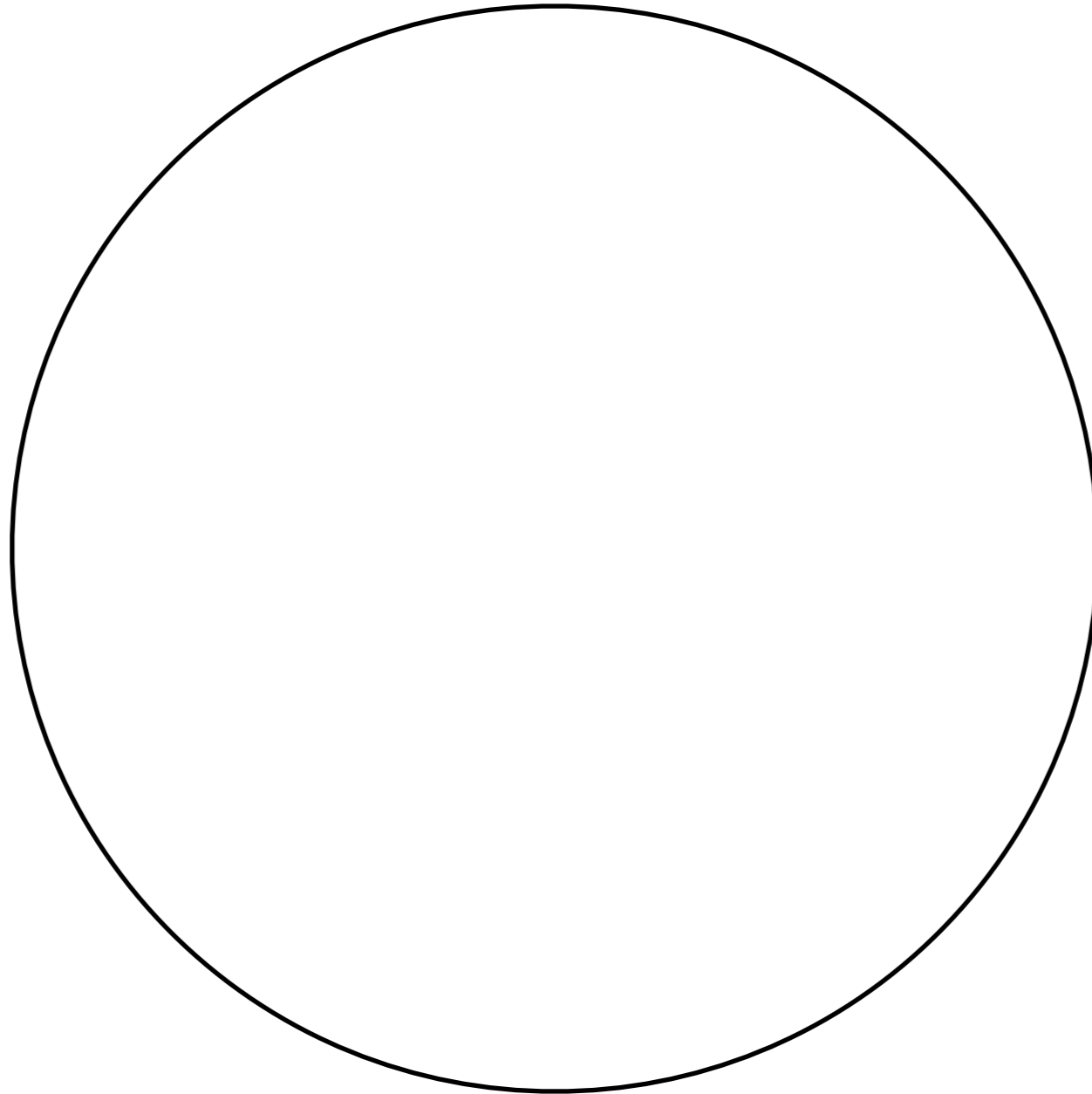
Cache 1 ●



Cache 2 ●



Cache 3 ●



Proposal 4

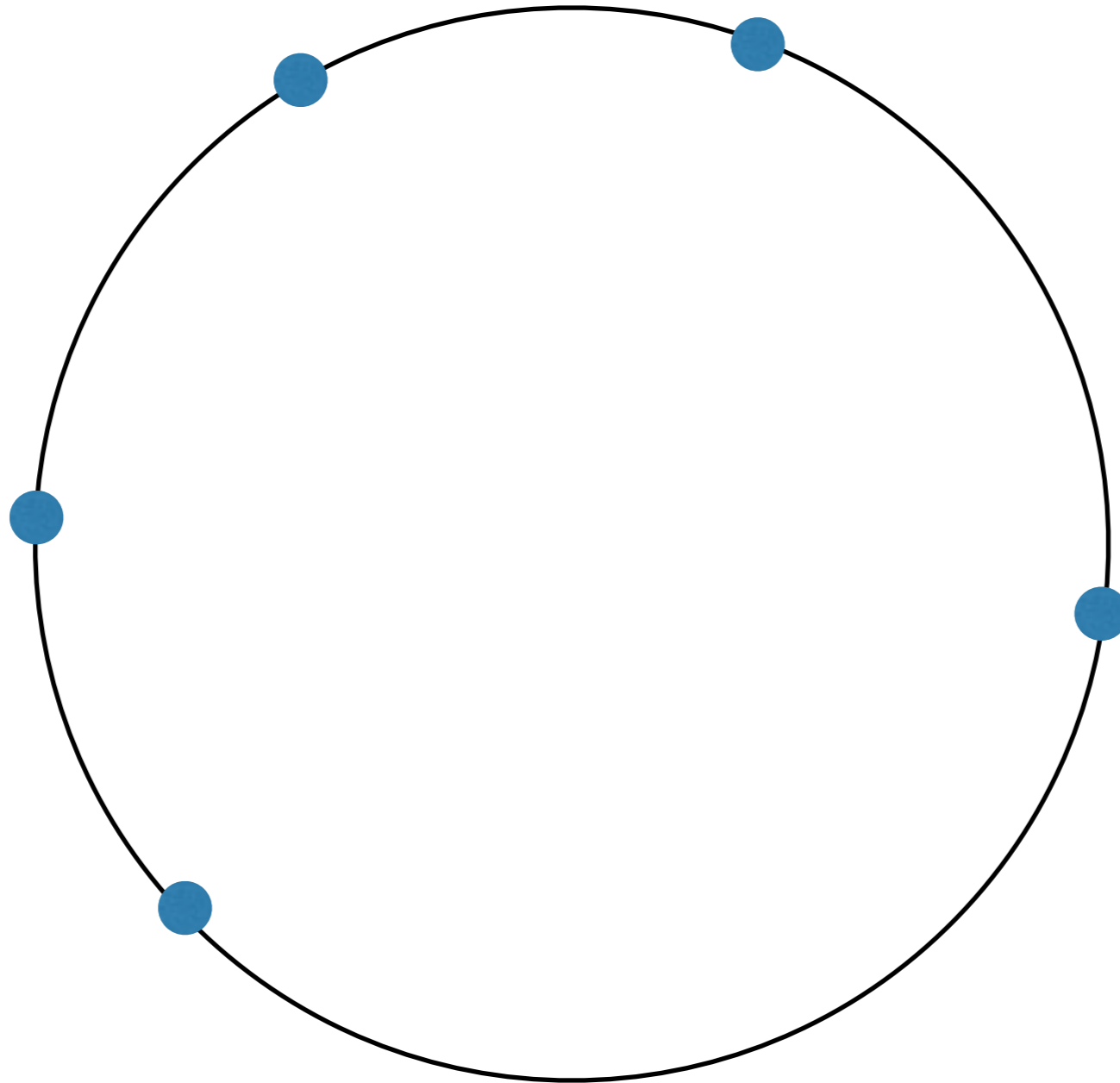
Cache 1 ●



Cache 2 ●



Cache 3 ●



Proposal 4

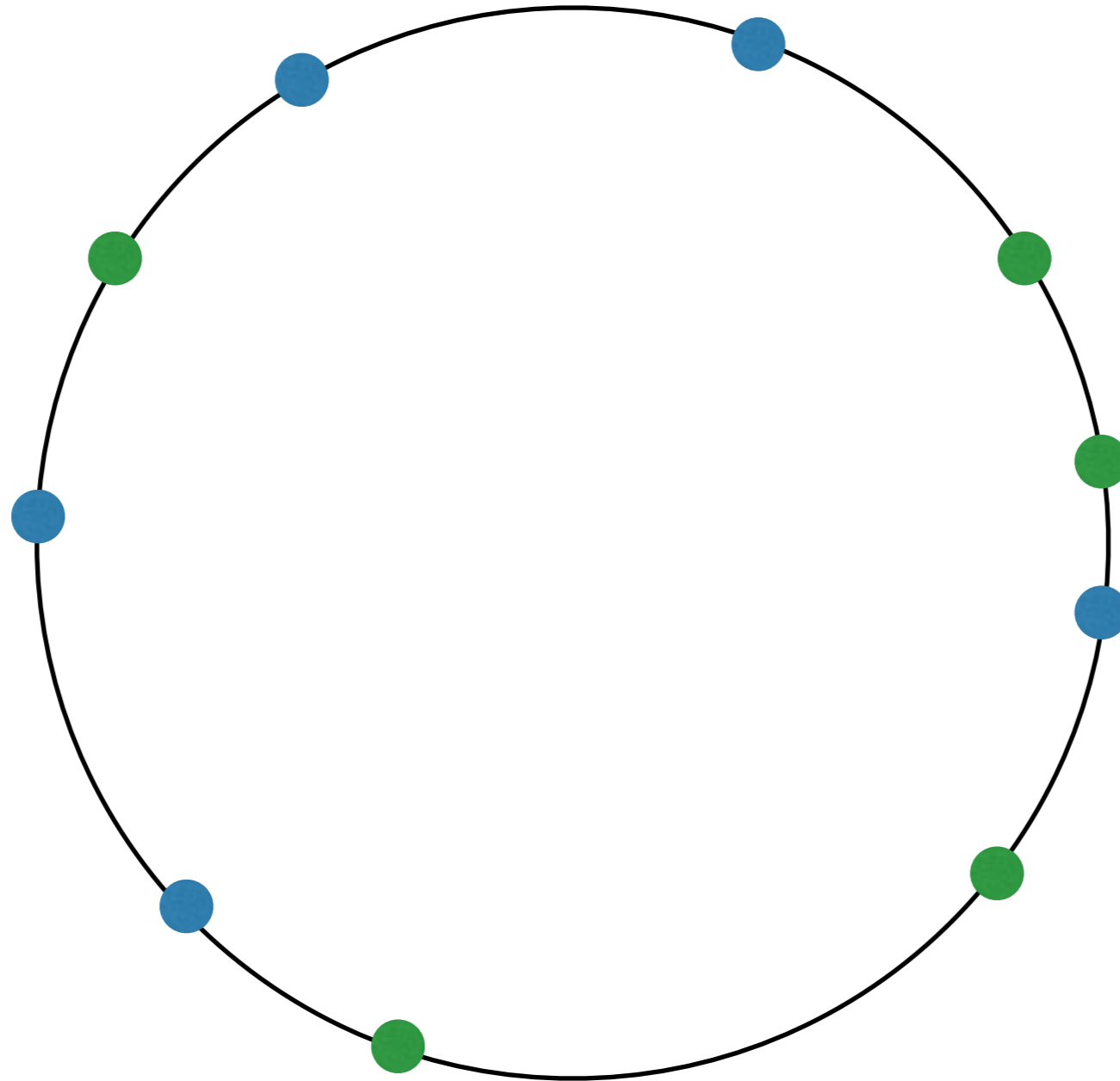
Cache 1 ●



Cache 2 ●



Cache 3 ●



Proposal 4

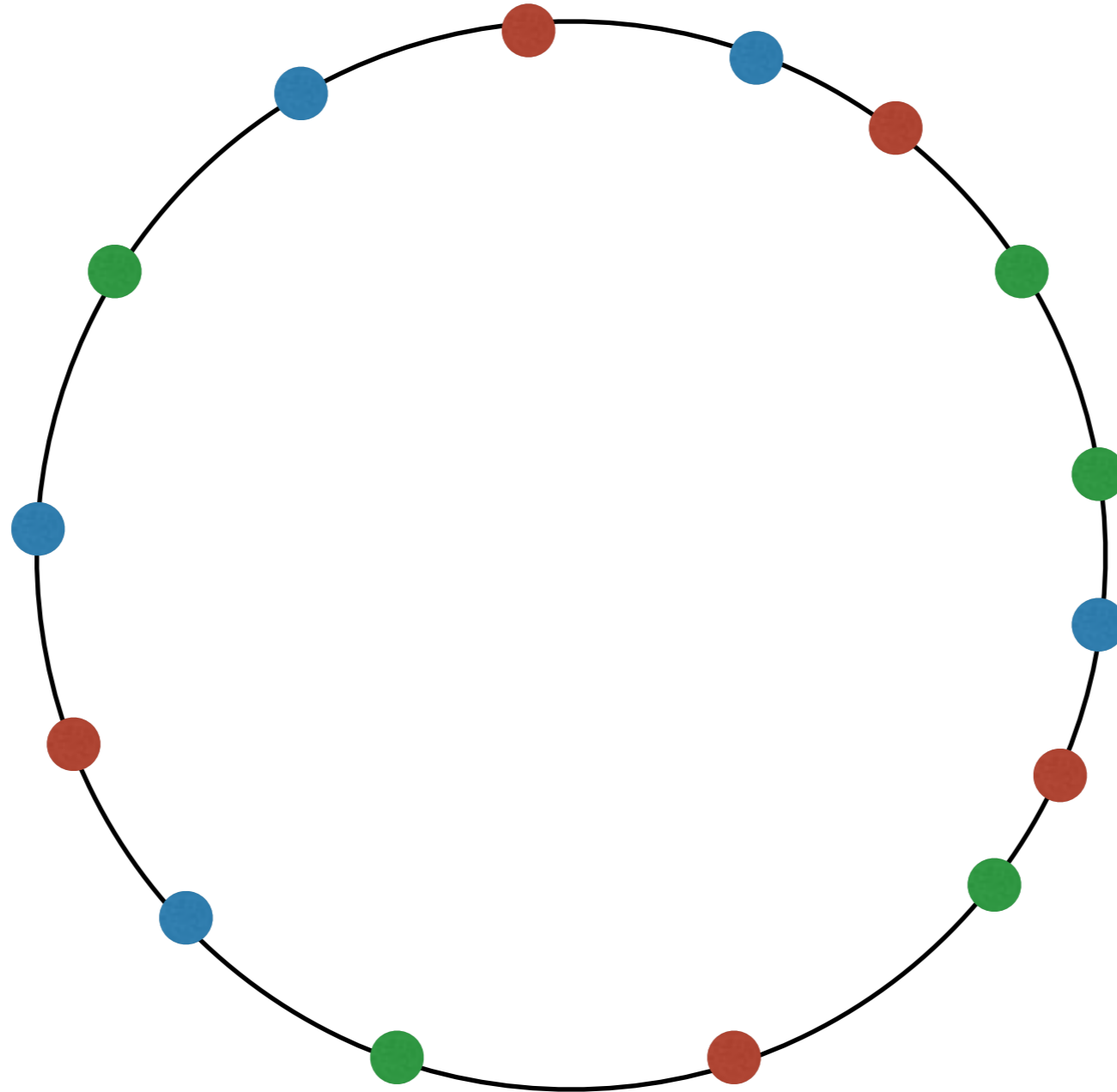
Cache 1 ●



Cache 2 ●



Cache 3 ●



Requirements, revisited

Requirement 1: clients all have same assignment

Requirement 2: keys uniformly distributed

Requirement 3: can add/remove nodes w/o redistributing too many keys

Requirements, revisited

Requirement 1: clients all have same assignment

Requirement 2: keys uniformly distributed

Requirement 3: can add/remove nodes w/o redistributing too many keys

Recap: consistent hashing

Node ids hashed to many pseudorandom points on a circle

Keys hashed onto circle, assigned to “next” node

Idea used widely:

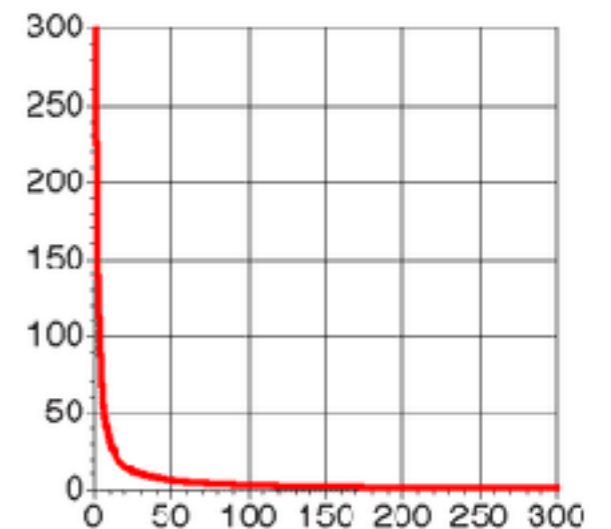
- Developed for Akamai CDN
- Used in Chord distributed hash table
- Used in Dynamo distributed DB

Further issues in load balancing

Accesses to keys are often Zipf-distributed

- Most popular accessed twice as often as 2nd-most
- Accessed 3x as often as 3rd-most, etc.

How to load balance “hot” keys?



Next week

Start of 3 weeks on “distributed systems in practice”

Lots of papers and discussion

Next week

Monday: Yegge on Service-Oriented Architectures

- Steve Yegge, prolific programmer and blogger
- Moved from Amazon to Google
- Monday's reading is an accidentally-leaked memo about differences between Amazon's and Google's system architectures (at that time)
- Advocates for SOA: separating applications (e.g. Google Search, Amazon) into many primitive services, run internally as products

