# DISTRIBUTED SYSTEMS AT UBER

**MATT RANNEY**

UBER

# OPEN SOURCE

partners

riders

dispatch

money

maps / ETA

services

post trip processing

storage

data

partners

riders

dispatch

money

maps / ETA

services

post trip processing

storage

data

partners

riders

dispatch

money

maps / ETA

services

post trip processing

storage

data

partners

riders

dispatch

money

maps / ETA

services

post trip processing

storage

data

partners

riders

dispatch

money

maps / ETA

services

post trip processing

storage

data

partners

riders

dispatch

money

maps / ETA

services

post trip processing

storage

data

partners    riders

dispatch    money

maps / ETA    services    post trip processing

storage    data

partners

riders

dispatch

money

maps / ETA

services

post trip
processing

storage

data

partners

riders

dispatch

money

maps / ETA

services

post trip processing

storage

data

| | dispatch | money |
|---|---|---|
| maps / ETA | services | post trip processing |
| | storage | data |

supply
humans

demand
humans

supply

demand

**Dispatch**

| supply humans | demand humans |
| --- | --- |

| supply | demand |
| --- | --- |

**Dispatch**

supply
humans

demand
humans

supply

demand

**Dispatch**

supply humans | demand humans

supply | demand

DISCO

Dispatch

supply humans

demand humans

supply

demand

DISCO

geo by supply

routing / ETA

geo by demand

Dispatch

supply humans

demand humans

supply

demand

DISCO

geo by supply

routing / ETA

geo by demand

**Dispatch**

| supply | demand |
| DISCO |
| geo by supply | routing / ETA | geo by demand |

Dispatch

# SCALING NODE.JS

Scalable, fault-tolerant application-layer sharding — Edit

🕐 **447** commits  |  ⅄ **19** branches  |  🏷 **97** releases  |  👥 **18** contributors

⟳ | Branch: **master** ▾  **ringpop-node** / +  ☰

Merge pull request **#175** from uber/move-validate-props  ⋯

jwolski authored 3 hours ago  |  latest commit **387f81d1f6** 📋

| 📁 benchmarks | not to expose membership functions | 2 months ago |
| 📁 docs | Fix typos in docs sample code | 7 days ago |
| 📁 examples | add an example for RingpopHandler | a month ago |
| 📁 lib | Add Member::getId | a day ago |
| 📁 scripts | Flap damp scoring | a month ago |
| 📁 server | Get rid of a function off Ringpop prototype | a day ago |
| 📁 test | Add Member::getId | a day ago |
| 📄 .gitignore | Initial commit | 9 months ago |
| 📄 .jshintrc | Initial commit | 9 months ago |
| 📄 .travis.yml | Travis CI integration | 8 months ago |

<> **Code**

⊙ **Issues**  19

⅄ **Pull requests**  8

📖 **Wiki**

⚡ **Pulse**

📊 **Graphs**

⚙ **Settings**

**SSH** clone URL

git@github.com:uber,  📋

You can clone with **HTTPS**, **SSH**, or **Subversion**. ⊘

⬇ **Clone in Desktop**

⬇ **Download ZIP**
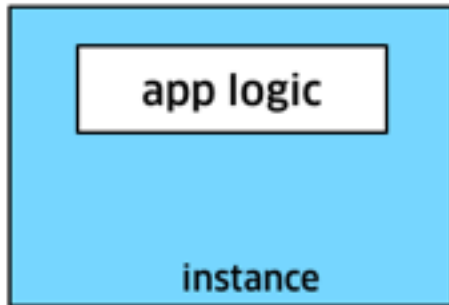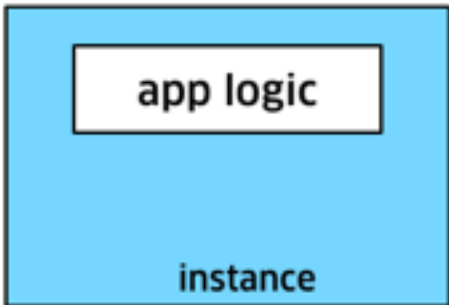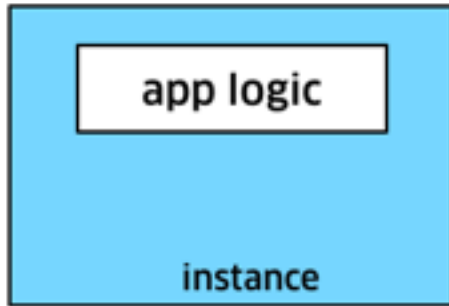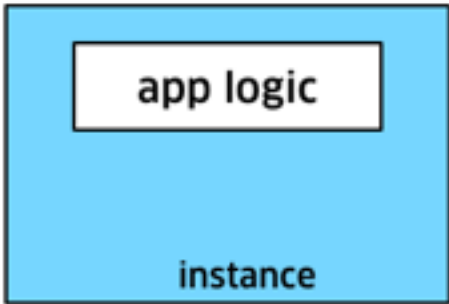
# How Ringpop Works

## Joining a Cluster

1. The first node, A, checks a bootstrap list and finds that no other nodes are running.
2. Next, B starts up and has A to join. B reads the file from disk, then selects a random number of members. It will find A and start to form a consistent hash ring in the background, running within memory in Ringpop.
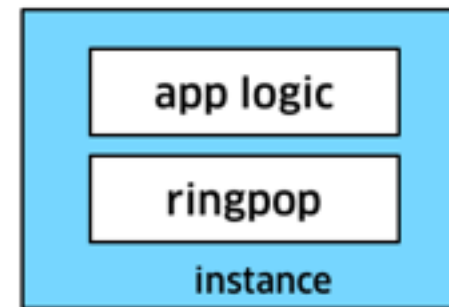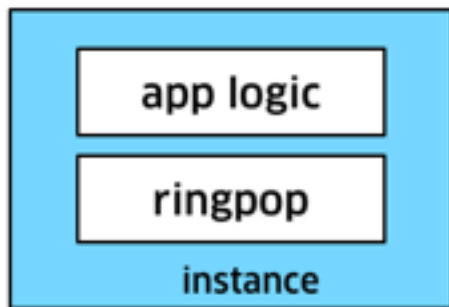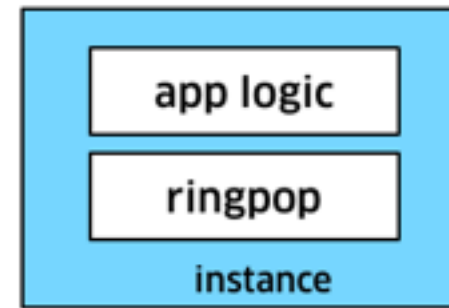3. The nodes are positioned along the ring and exchange information with one another, forming a two-node cluster and pinging each other back and forth.
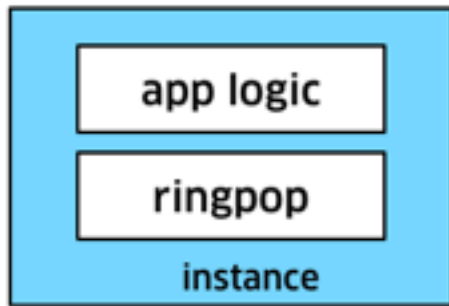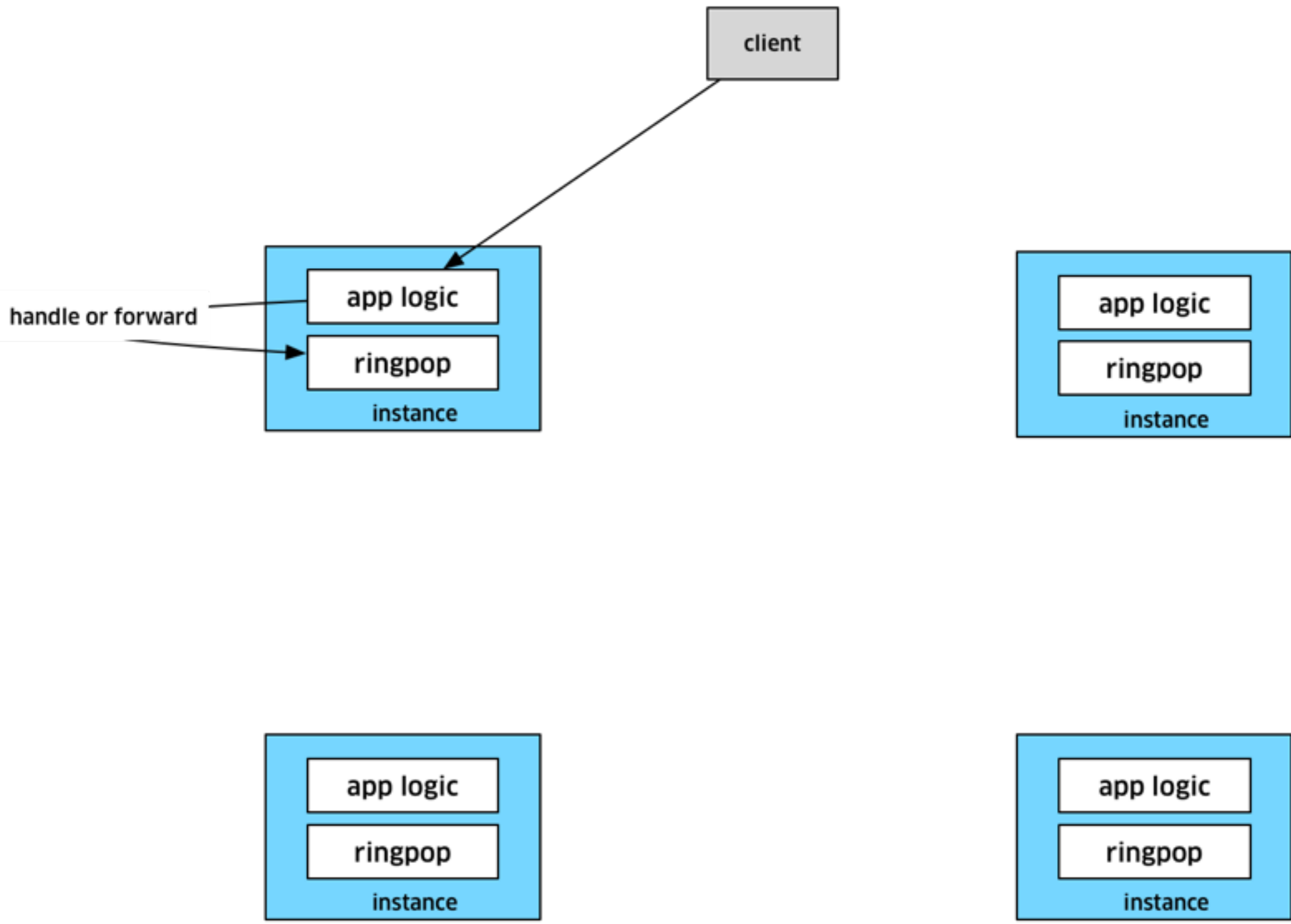
## Handle or Forward

Upon arrival of a proxied request at its destination, membership checksums of the sender and receiver will be compared. The request will be refused if checksums differ. Mismatches are expected when nodes are entering or exiting the cluster due to deploys, added/removed capacity, or failures. The cluster will eventually converge on one membership checksum, therefore refused requests are best handled by retrying them.
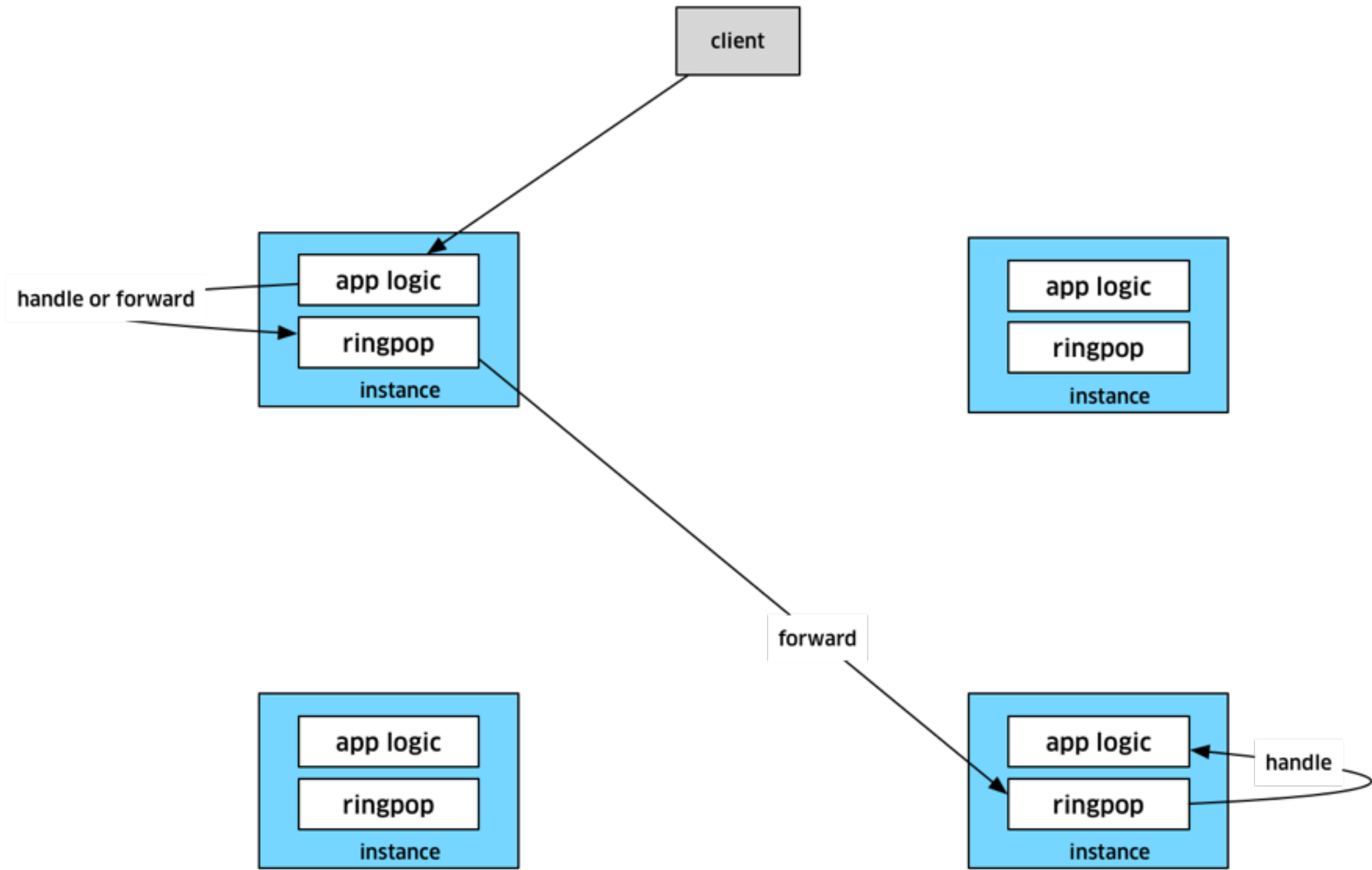
Ringpop's request proxy has retries built in and can be tuned using two parameters provided at the time Ringpop is instantiated: `requestProxyMaxRetries` and `requestProxyRetrySchedule` or per-request with: `maxRetries` and `retrySchedule`. The first parameter is an integer representing the number of times a particular request is retried. The second parameter is an array of integer or floating point values representing the delay in-between consecutive retries.

svc B

svc B

svc B

svc B

forward

svc A

req sent to
random healthy node

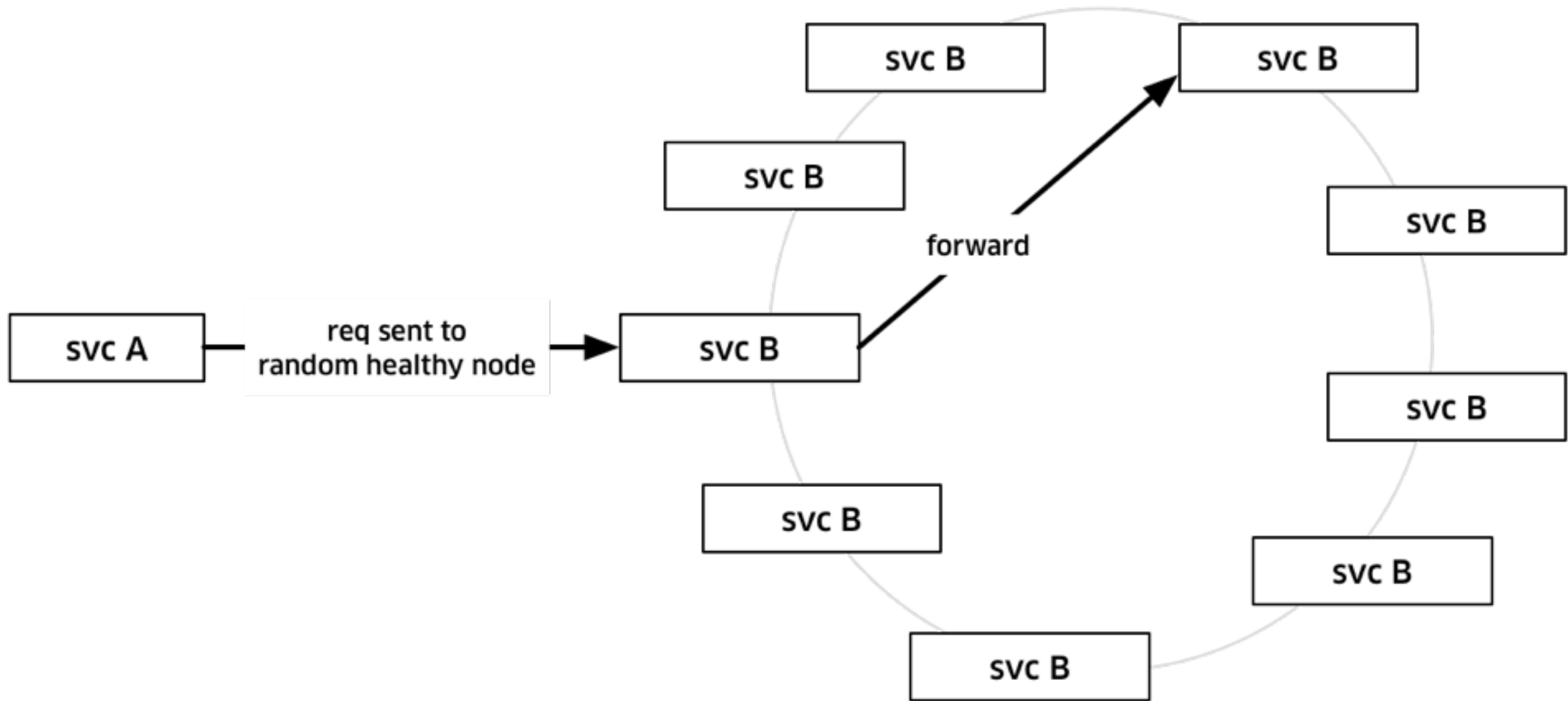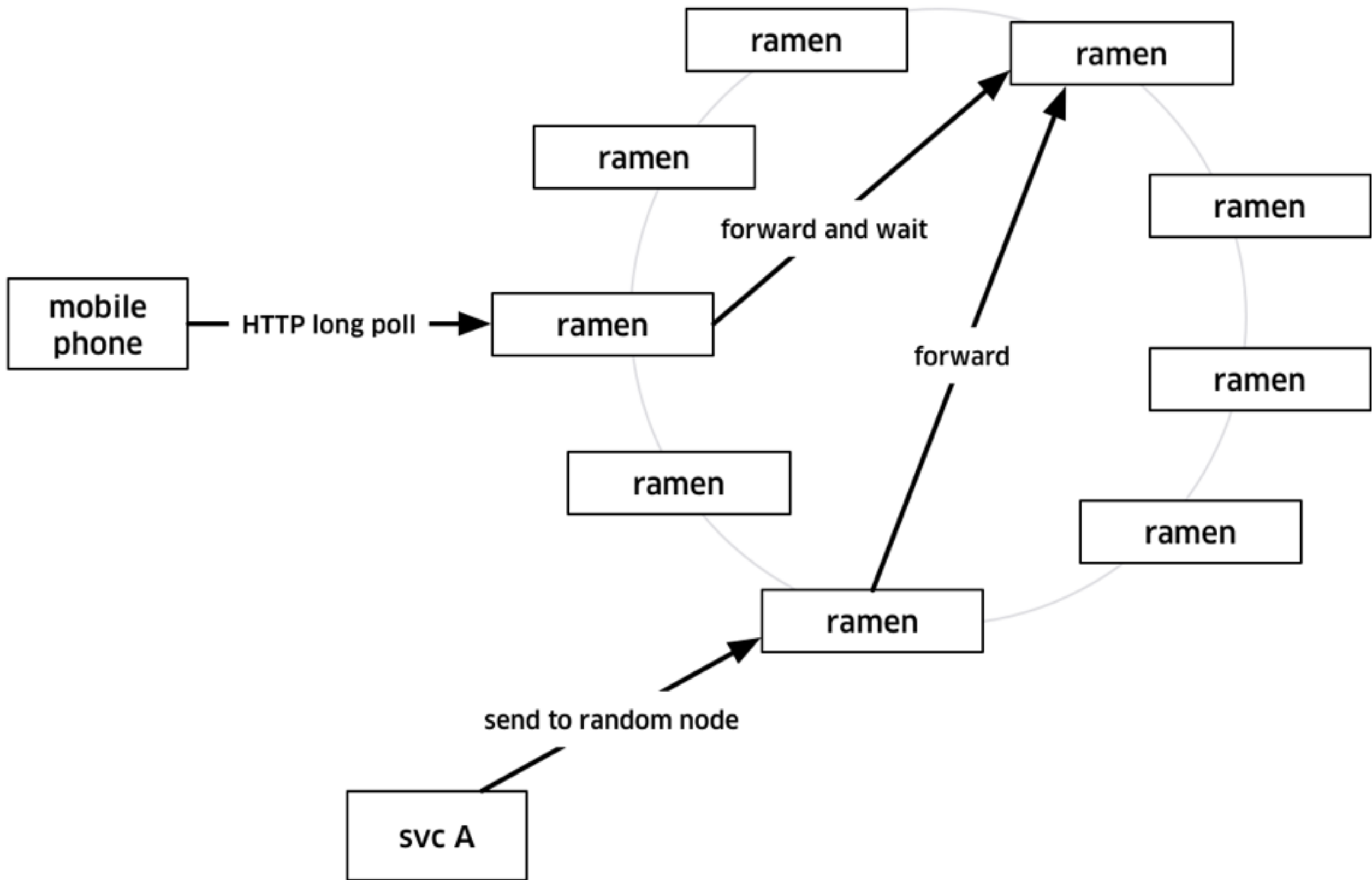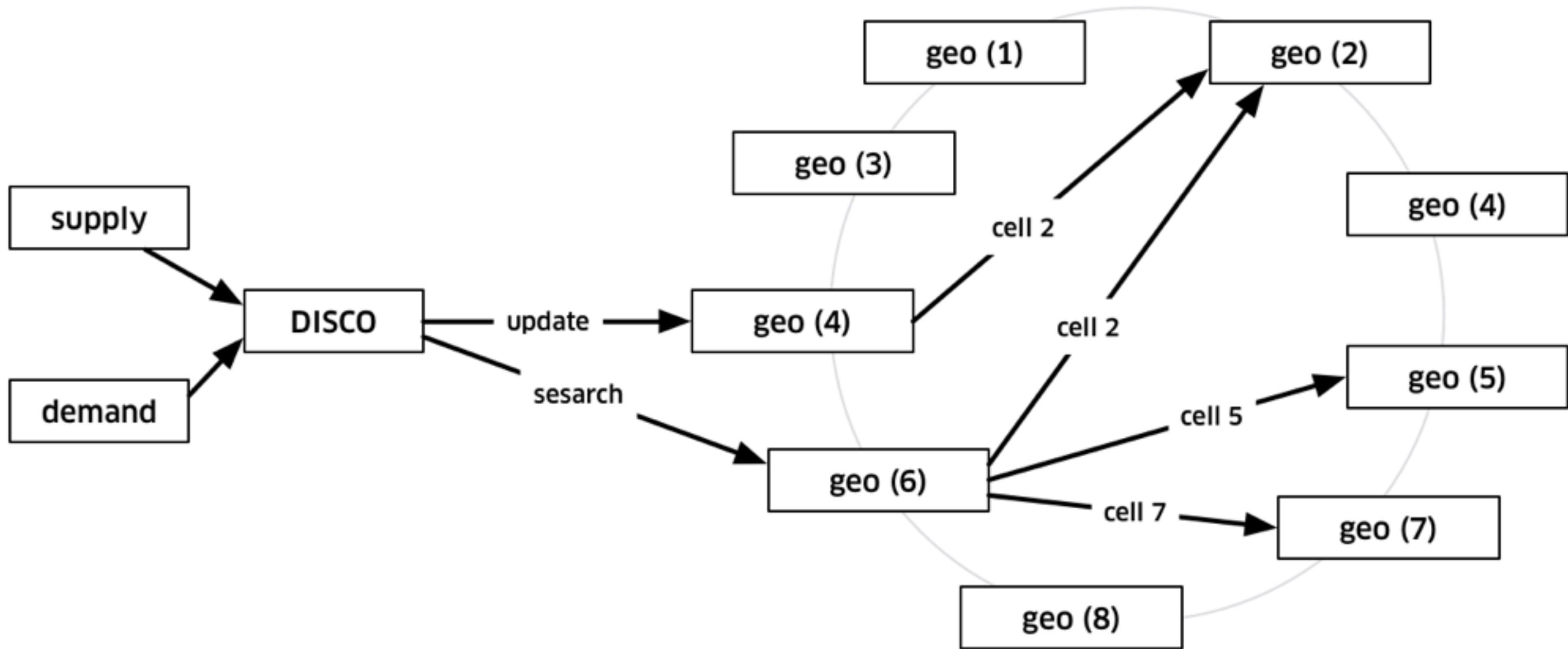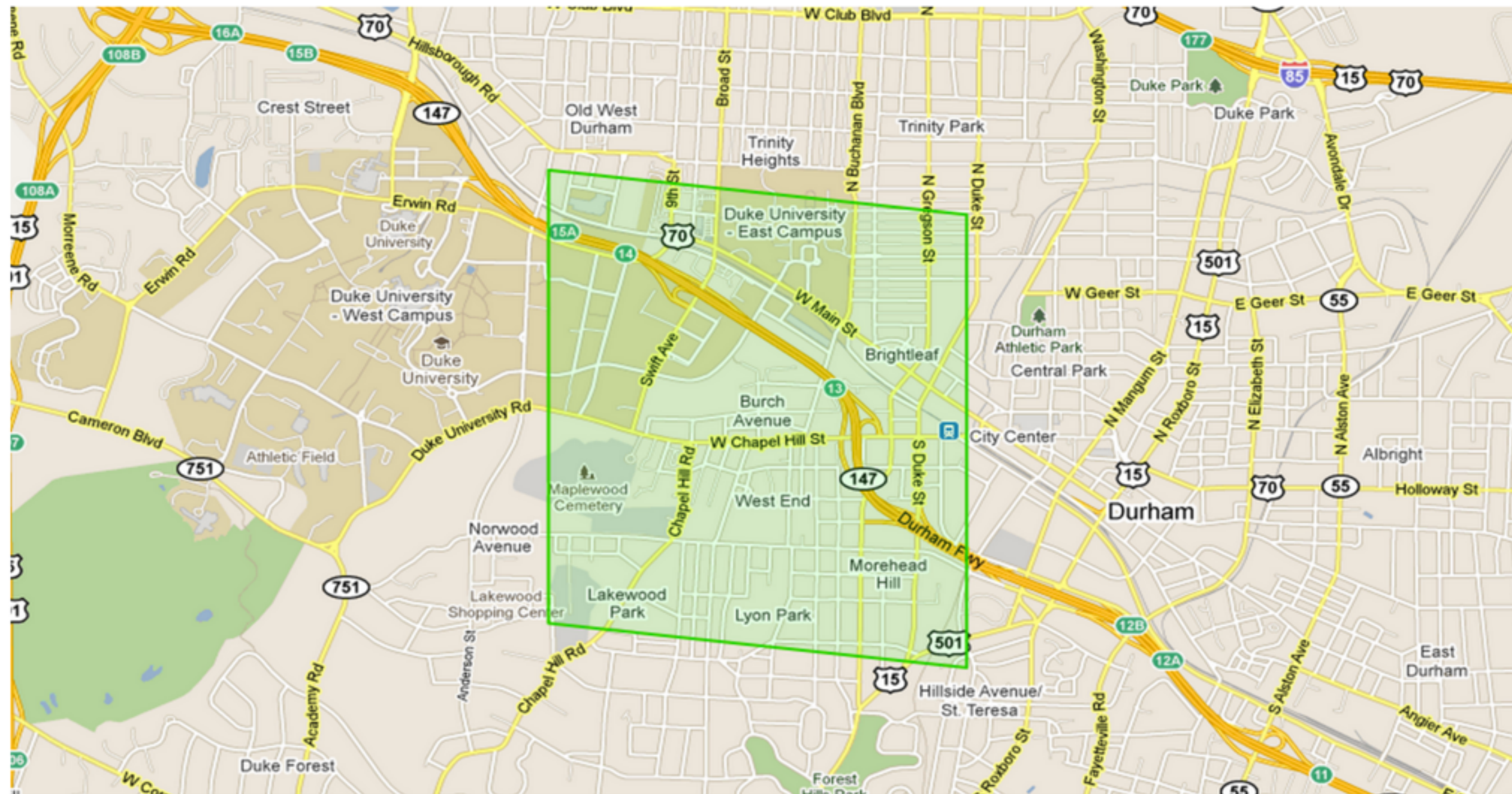svc B

svc B

svc B

svc B

svc B

# One S2 Cell

Id: 0x89ace41000000000 (0b1000100110101100111001000001000...), Level: 12

# S2 Cells - Stats

| Level | Min Area | Max Area |
|---|---|---|
| 0 | 85,011,012 km$^2$ | 85,011,012 km$^2$ |
| 1 | 21,252,753 km$^2$ | 21,252,753 km$^2$ |
| 12 | 3.31 km$^2$ | 6.38 km$^2$ |
| 30 | 0.48 cm$^2$ | 0.93 cm$^2$ |

⇧
smallest cell

Every cm$^2$ on Earth can be represented using a 64-bit integer.

Source: Geometry on the Sphere: Google's S2 Library
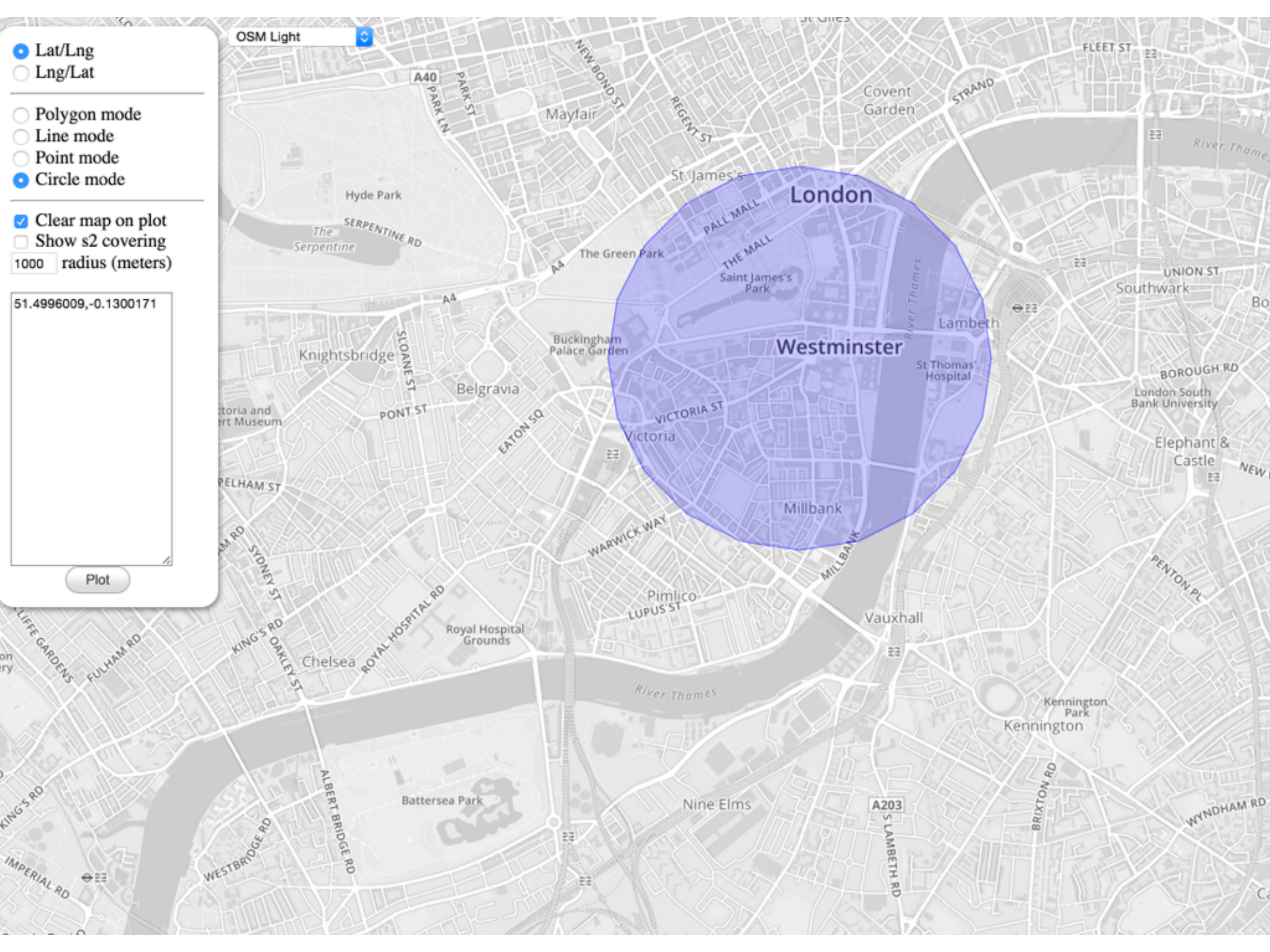
OSM Light

- Lat/Lng
- Lng/Lat

- Polygon mode
- Line mode
- Point mode
- Circle mode

☑ Clear map on plot
☐ Show s2 covering
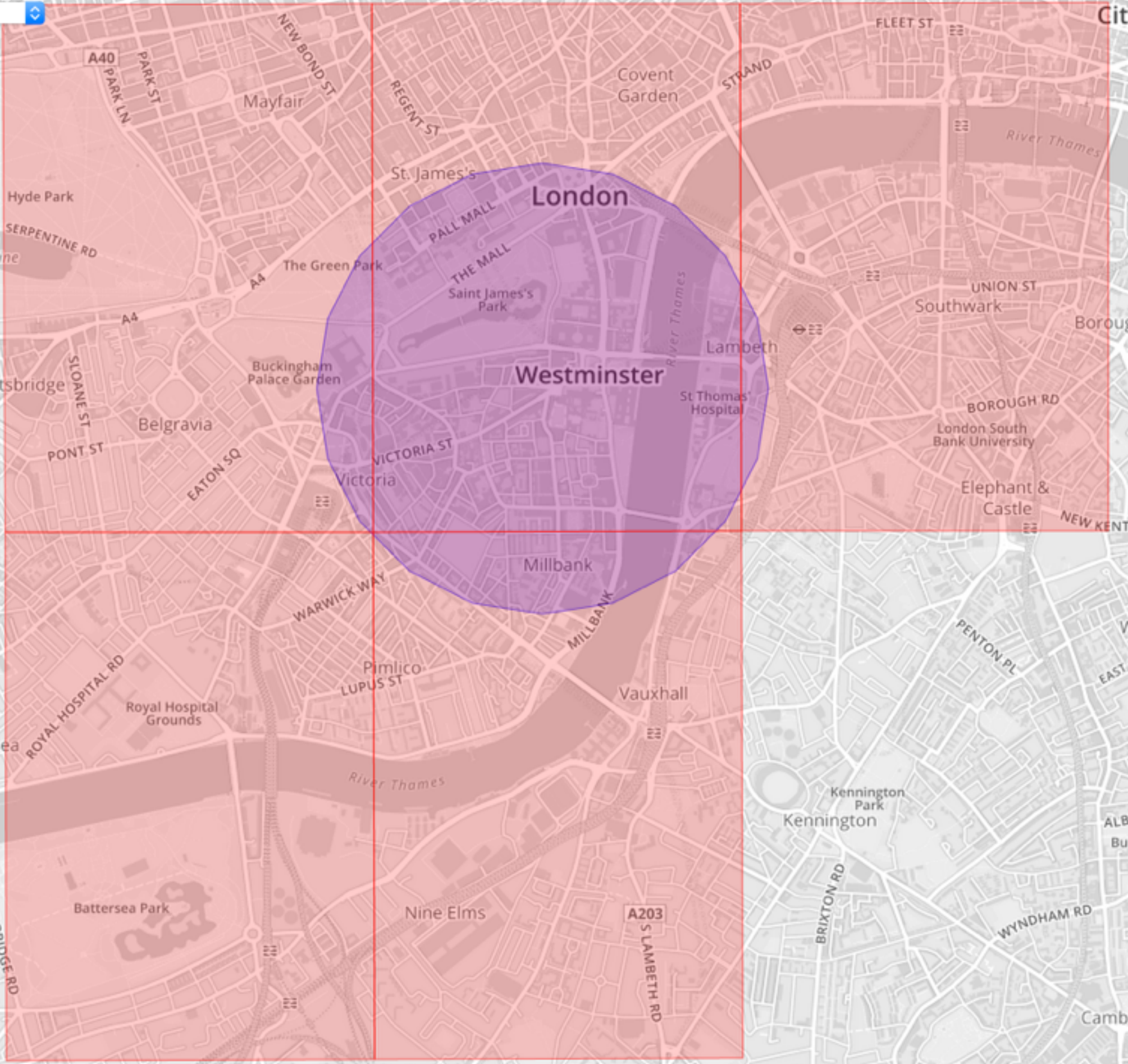1000  radius (meters)

51.4996009,-0.1300171

Plot
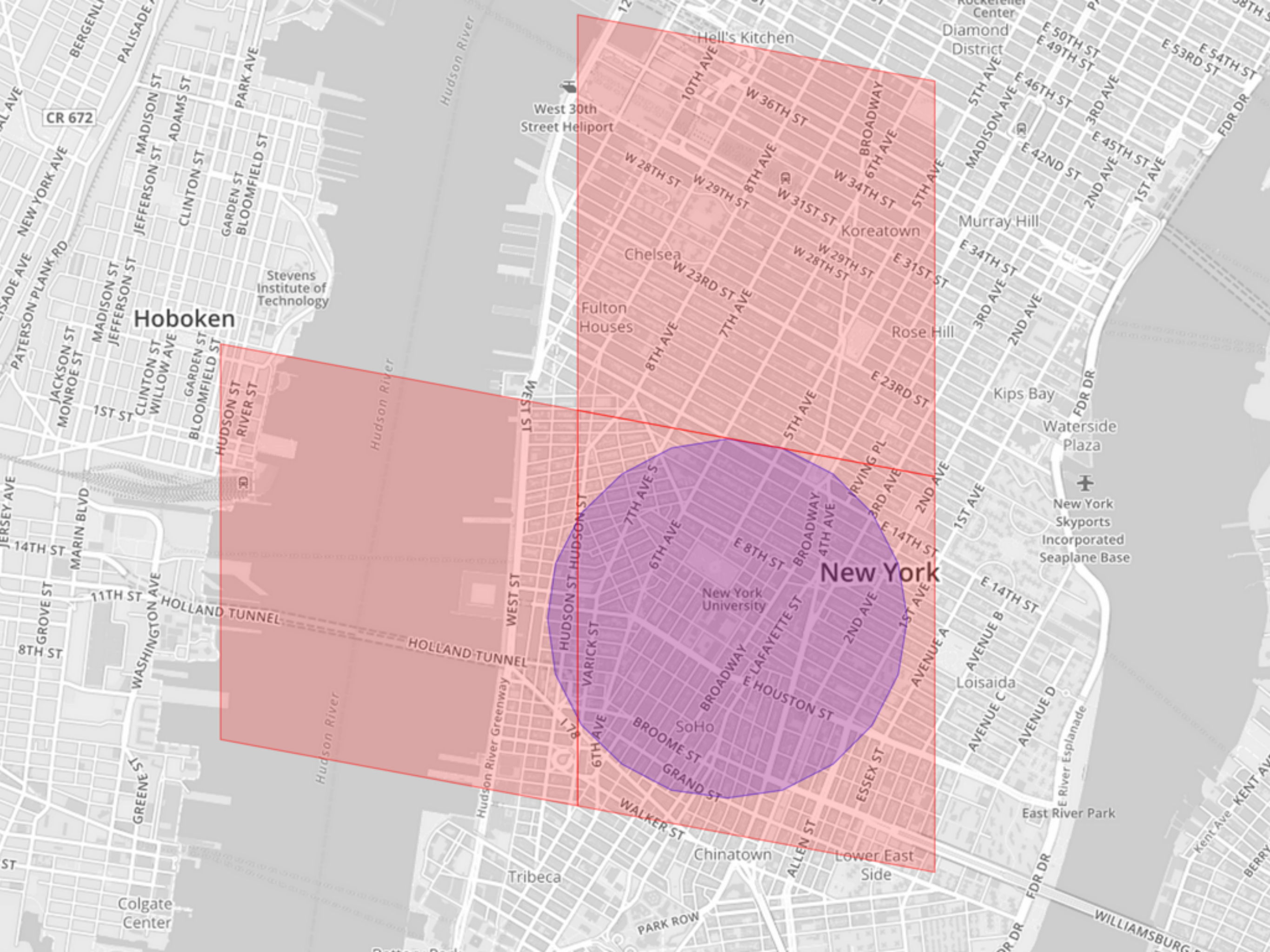
Lat/Lng
Lng/Lat

Polygon mode
Line mode
Point mode
Circle mode

☑ Clear map on plot
☑ Show s2 covering
12 min_level
12 max_level
200 max_cells
1 level_mod
1000 radius (meters)

51.4996009,-0.1300171

Plot

OSM Light

# SWIM: *Scalable Weakly-consistent Infection-style Process Group Membership Protocol*

Abhinandan Das, Indranil Gupta, Ashish Motivala*
Dept. of Computer Science, Cornell University
Ithaca NY 14853 USA
{asdas,gupta,ashish}@cs.cornell.edu

## Abstract

*Several distributed peer-to-peer applications require weakly-consistent knowledge of process group membership information at all participating processes. SWIM is a generic software module that offers this service for large-scale process groups. The SWIM effort is motivated by the unscalability of traditional heart-beating protocols, which either impose network loads that grow quadratically with group size, or compromise response times or false positive frequency w.r.t. detecting process crashes. This paper reports on the design, implementation and performance of the SWIM sub-system on a large cluster of commodity PCs.*

## 1. Introduction

*As you swim lazily through the milieu,*
*The secrets of the world will infect you.*

Several large-scale peer-to-peer distributed process groups running over the Internet rely on a distributed membership maintenance sub-system. Examples of existing middleware systems that utilize a membership protocol include reliable multicast [3, 11], and epidemic-style information dissemination [4, 8, 13]. These protocols in turn find use in applications such as distributed databases that need to reconcile recent disconnected updates [14], publish-subscribe systems, and large-scale peer-to-peer systems[15]. The performance

# Ringpop UI

## Datacenter

| Supply | > | SJC9 |

*Last data fetch: 2015/8/11 11:3:32*

*Connected Node: 10.32.162.16:3000*

*TChannel version: ^1.2.5, Ringpop version: 9.8.18*

15 nodes

■ 12 alive

**Cluster Health**

### Key lookup

fake-uuid

Search

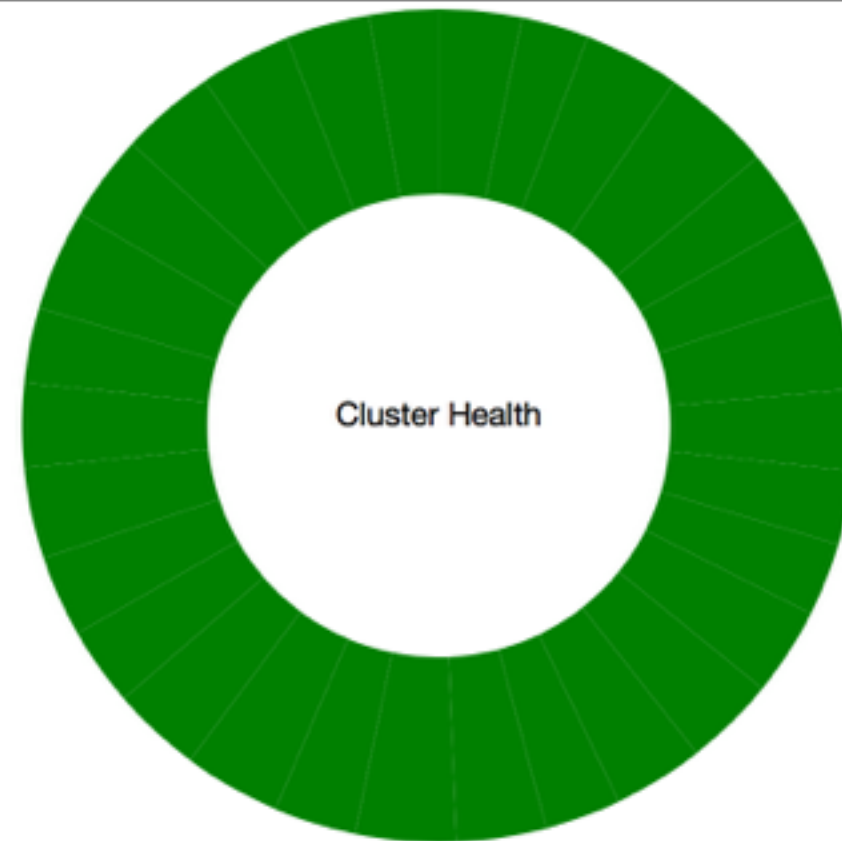|  | 11AM | 10AM | 9AM | 8AM | 7AM | 6AM | 5AM | 4AM | 3AM | 2AM | 1AM | 12PM | 11PM | 10PM | 9PM | 8PM | 7PM | 6PM | 5PM | 4PM | 3PM | 2PM | 1PM | 12AM |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 10.32.162.16:3000 | | | | | | | | | | | | | | | | | | | | | | | | |
| 10.32.162.16:3004 | | | | | | | | | | | | | | | | | | | | | | | | |
| 10.32.162.16:3001 | | | | | | | | | | | | | | | | | | | | | | | | |
| 10.32.162.16:3002 | | | | | | | | | | | | | | | | | | | | | | | | |
| 10.32.162.16:3003 | | | | | | | | | | | | | | | | | | | | | | | | |
| 10.32.162.16:3005 | | | | | | | | | | | | | | | | | | | | | | | | |
| 10.32.162.16:3006 | | | | | | | | | | | | | | | | | | | | | | | | |

# Ringpop UI

## Datacenter

Supply > SJC9

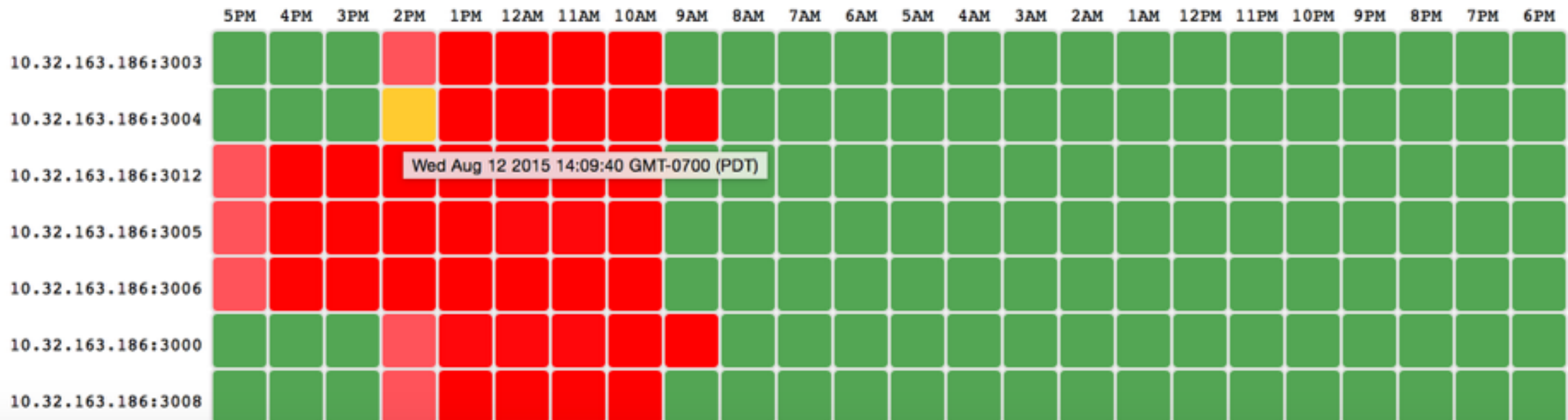*Last data fetch: 2015/8/12 17:4:1*

*Connected Node: 10.32.163.186:3000*

*TChannel version: ^1.2.5, Ringpop version: 9.8.18*

30 nodes

| | 30 alive |

Cluster Health

### Key lookup

fake-uuid

Search

| | 5PM | 4PM | 3PM | 2PM | 1PM | 12AM | 11AM | 10AM | 9AM | 8AM | 7AM | 6AM | 5AM | 4AM | 3AM | 2AM | 1AM | 12PM | 11PM | 10PM | 9PM | 8PM | 7PM | 6PM |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 10.32.163.186:3003 | | | | | | | | | | | | | | | | | | | | | | | | |
| 10.32.163.186:3004 | | | | | | | | | | | | | | | | | | | | | | | | |
| 10.32.163.186:3012 | | | | | Wed Aug 12 2015 14:09:40 GMT-0700 (PDT) | | | | | | | | | | | | | | | | | | | |
| 10.32.163.186:3005 | | | | | | | | | | | | | | | | | | | | | | | | |
| 10.32.163.186:3006 | | | | | | | | | | | | | | | | | | | | | | | | |
| 10.32.163.186:3000 | | | | | | | | | | | | | | | | | | | | | | | | |
| 10.32.163.186:3008 | | | | | | | | | | | | | | | | | | | | | | | | |

| | | |
|---|---|---|
| 📄 .travis.yml | .travis.yml: simplify make chdir'ing | a day ago |
| 📄 LICENSE | Add LICENSE | 3 days ago |
| 📄 README.md | Add Travis build badge to readme | 2 days ago |

📖 **README.md**

# TChannel build passing

Network multiplexing and framing protocol for RPC

## Design goals

- Easy to implement in multiple languages, especially JS and Python.
- High performance forwarding path. Intermediaries can make a forwarding decision quickly.
- Request / response model with out of order responses. Slow requests will not block subsequent faster requests at head of line.
- Large requests/responses may/must be broken into fragments to be sent progressively.
- Optional checksums.
- Can be used to transport multiple protocols between endpoints, eg. HTTP+JSON and Thrift.

## MIT Licenced

# GOALS

- performance

- forwarding

- language support

- proper pipelining

- checksums / tracing

- encapsulation

# RPC

Getting out of the HTTP and JSON business

HTTP is slow, complex, and inconsistent

JSON is hard to validate and awkward in non-node

Thrift is OK, but generated code is bad

📖 **README.md**

# tcurl

A command line utility to talk to a tchannel server

```
tcurl -p host:port <service> <endpoint> [options]

  Options:
    -2 [data] send an arg2 blob
    -3 [data] send an arg3 blob
    --shardKey send ringpop shardKey transport header
    --depth=n configure inspect printing depth
    -j print JSON
    -J [indent] print JSON with indentation
    -t [dir] directory containing Thrift files
```

# Installation

📖 **README.md**

# tcap

Uses pcap to inspect tchannel traffic over a network interface.

```
Usage: tcap [options]

Options:

  -h, --help                  output usage information
  -V, --version               output the version number
  -i --interface <interface>  network interface interfaces (defaults to first with an addres:
  -p --port <port>            a port to track or use "port1-port2" for a range of ports to t
  -f --filter <filter>        packet filter in pcap-filter(7) syntax (default: all TCP packe
  -s --service <service-name> service name or names to show (default: all services shown), o
                              use "~service-name" to exclude the service
  -t --thrift <thrift>        path of the directory for thrift spec files
  -1 --arg1 <arg1-method>     arg1 method or methods to show (default: all arg1 methods show
                              use "~arg1-method" to exclude the arg1
  --m1                        show arg1 name in call responses
```
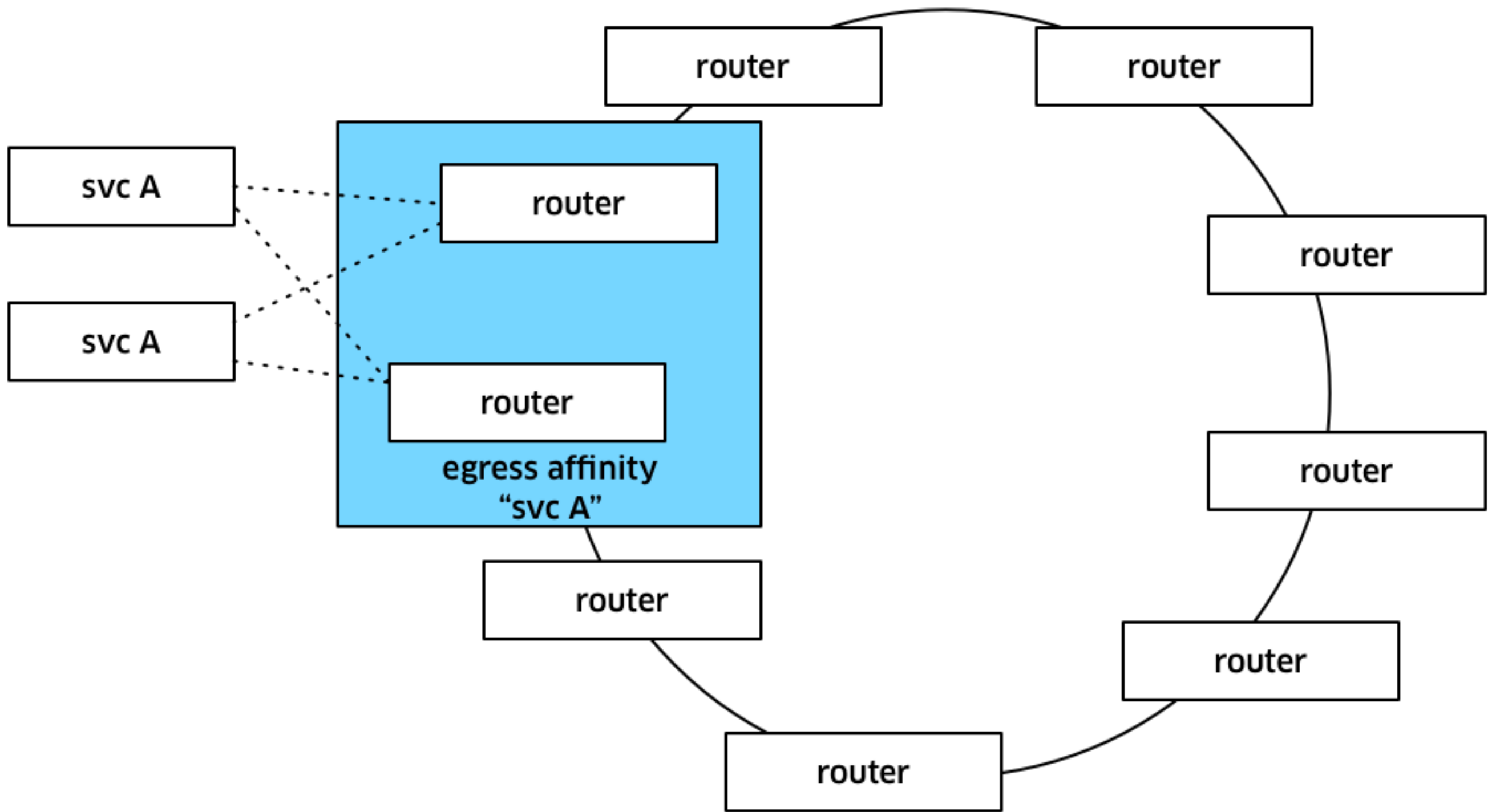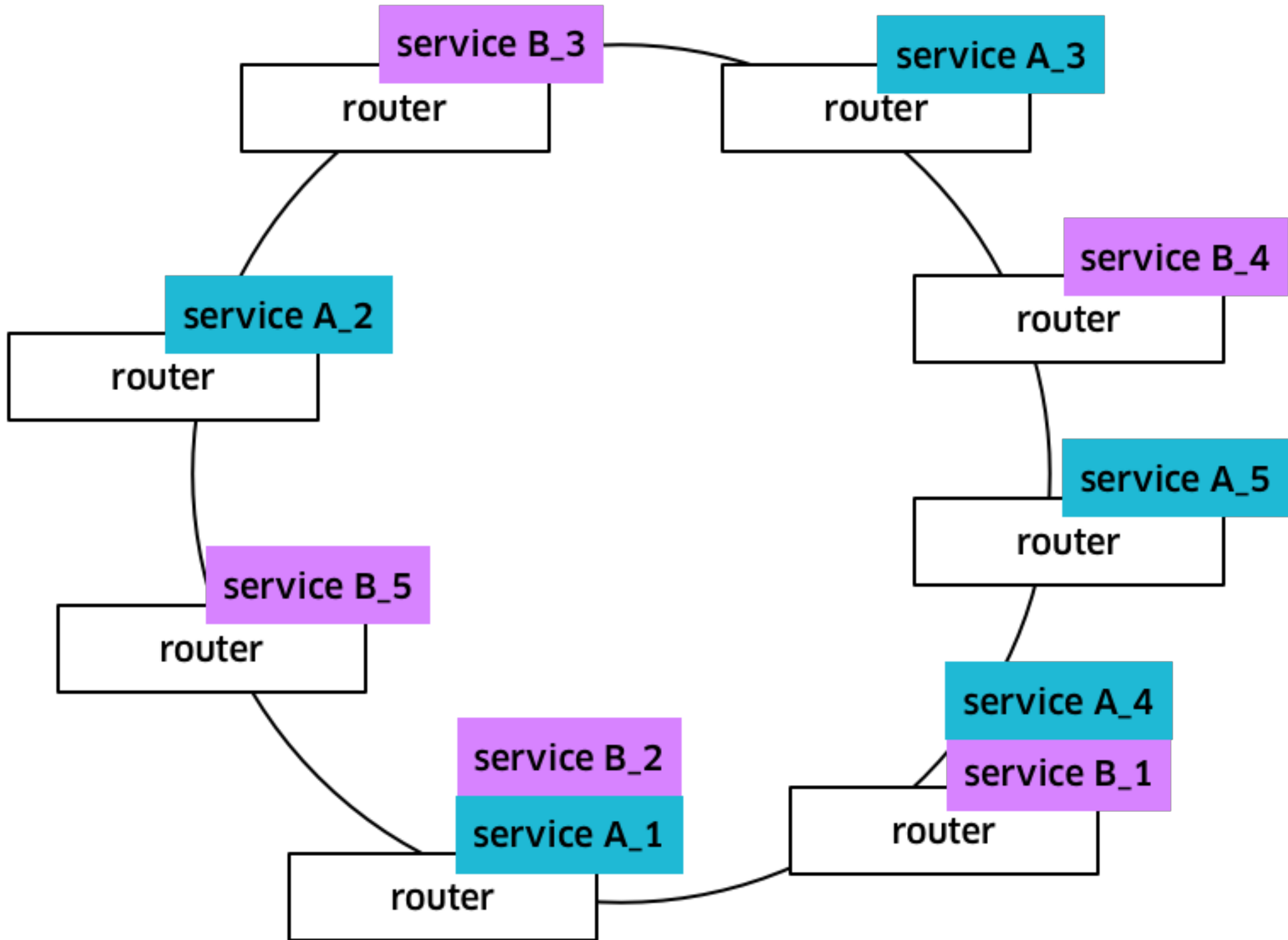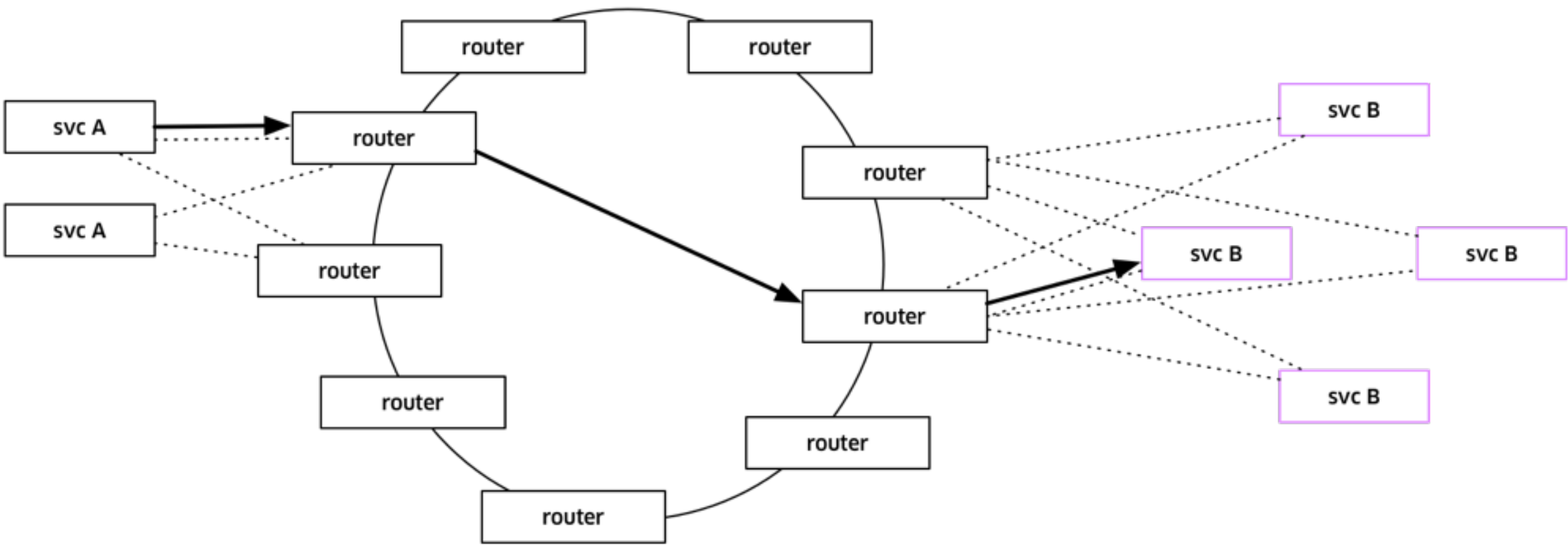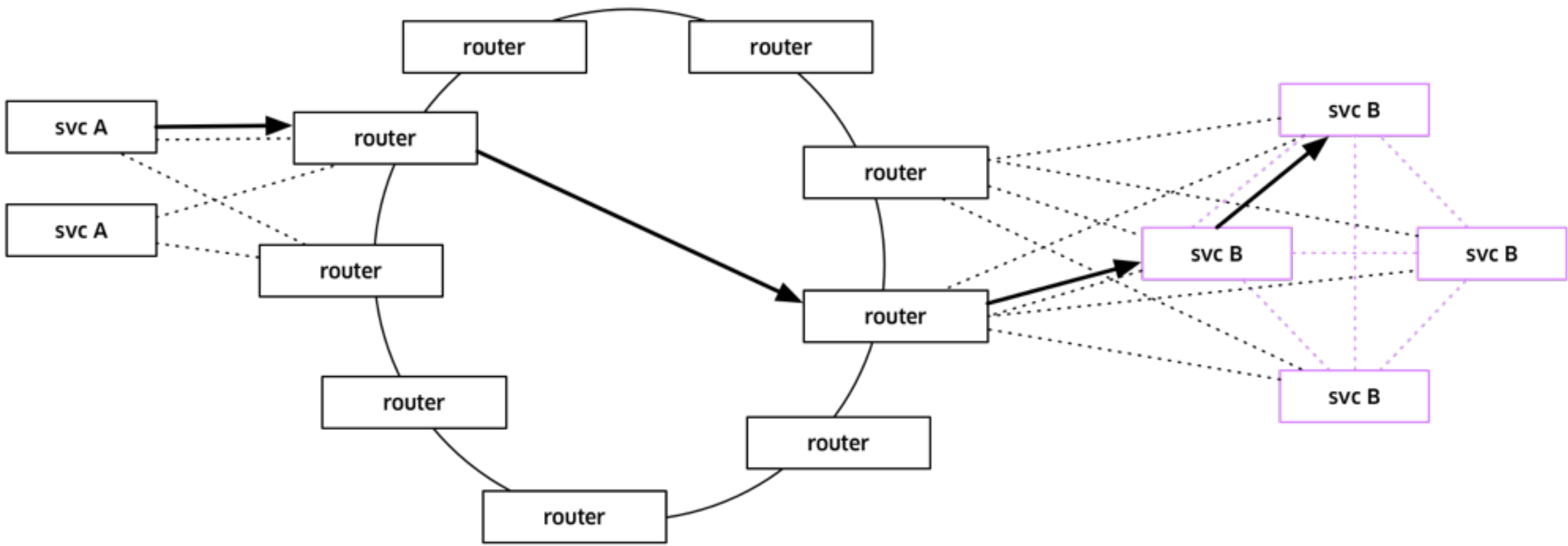
```
service A  →  load balancer  →  service B
```
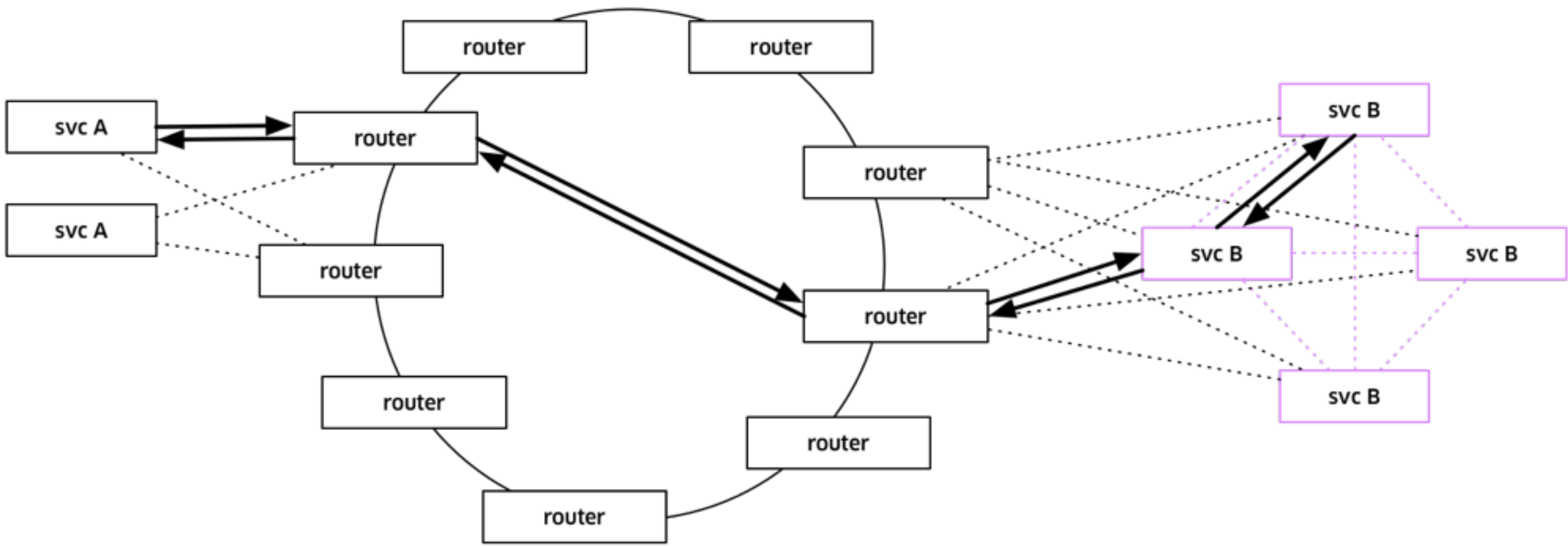
# HYPERBAHN

# HYPERBAHN

scalable registry and health checks

zipkin tracing

circuit breaking

rate limiting

failure testable

# AVAILABILITY

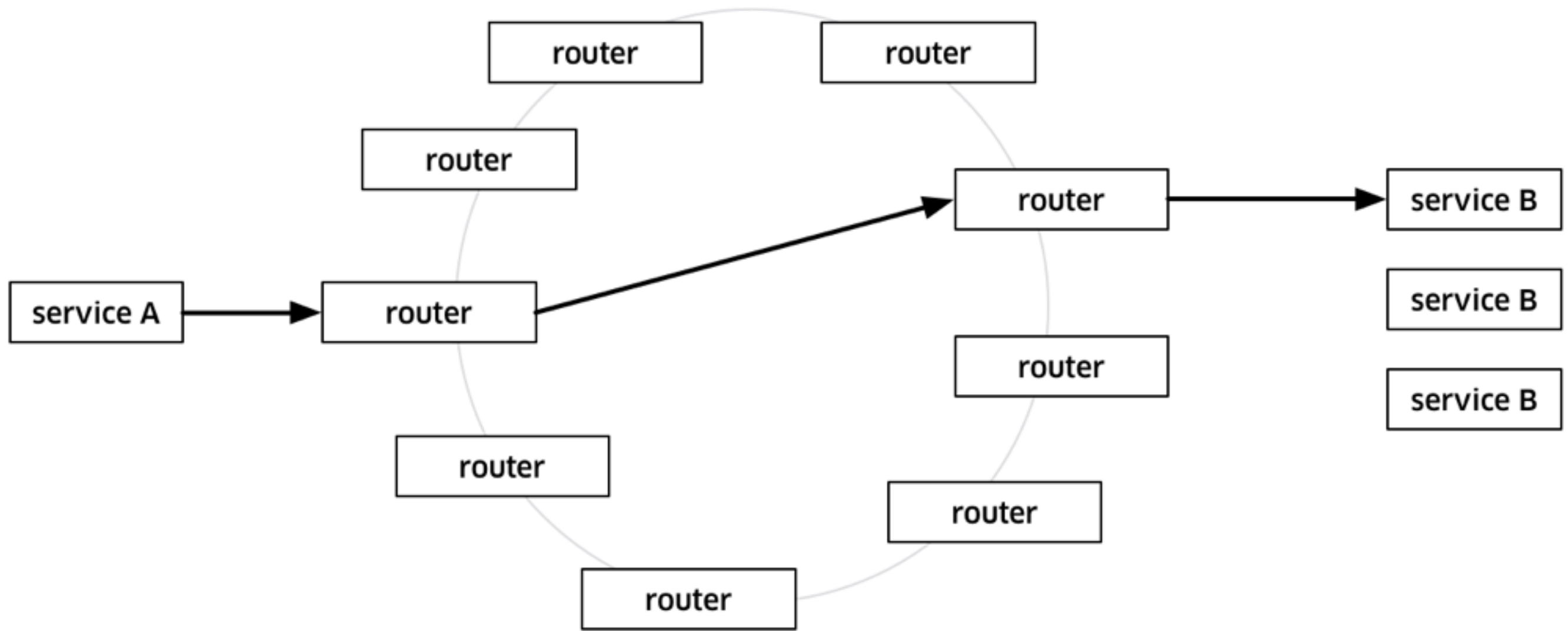everything retryable

everything killable
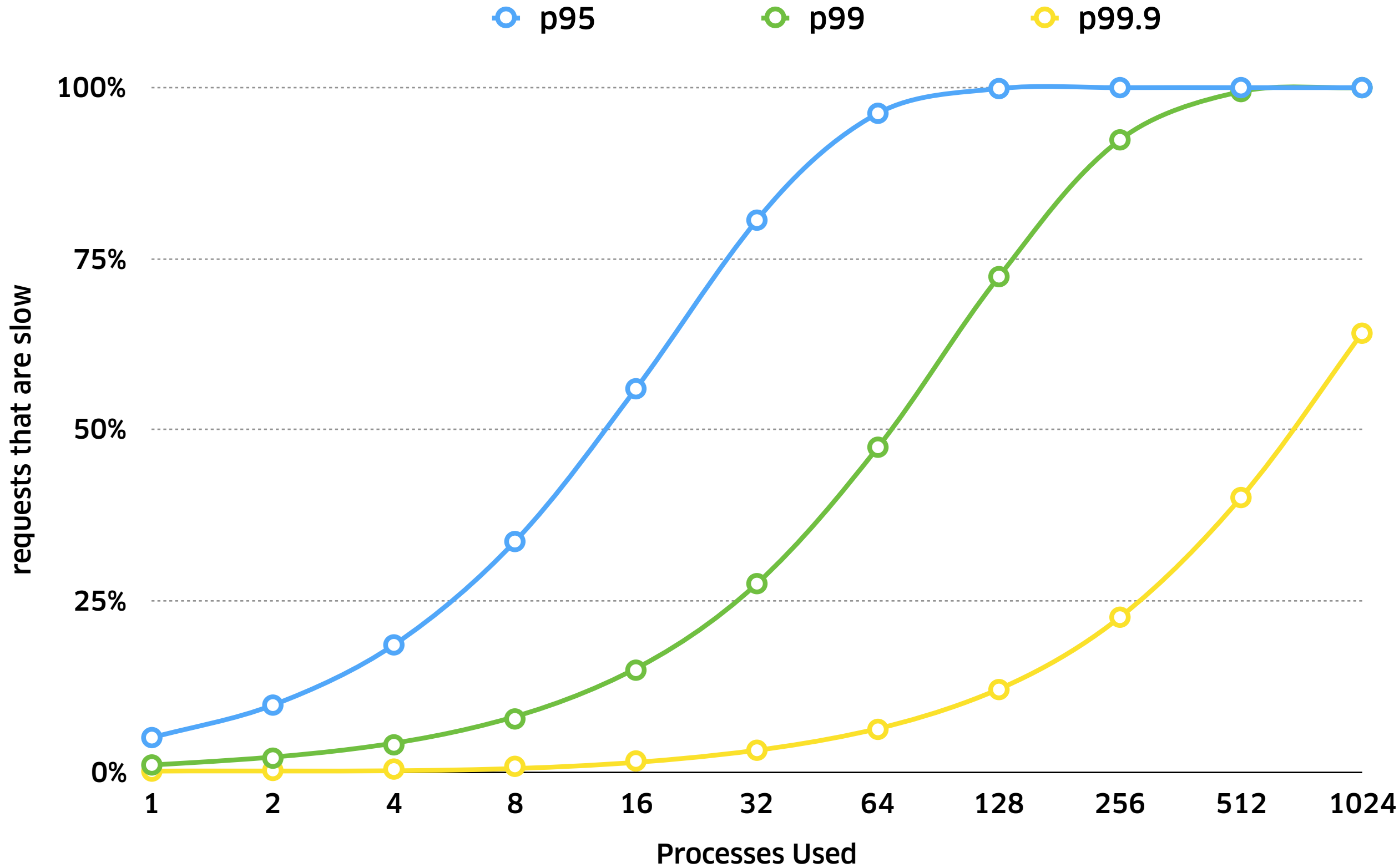
crash only

# CULTURAL CHANGES

no pairs

kill everything

even databases

# LATENCY
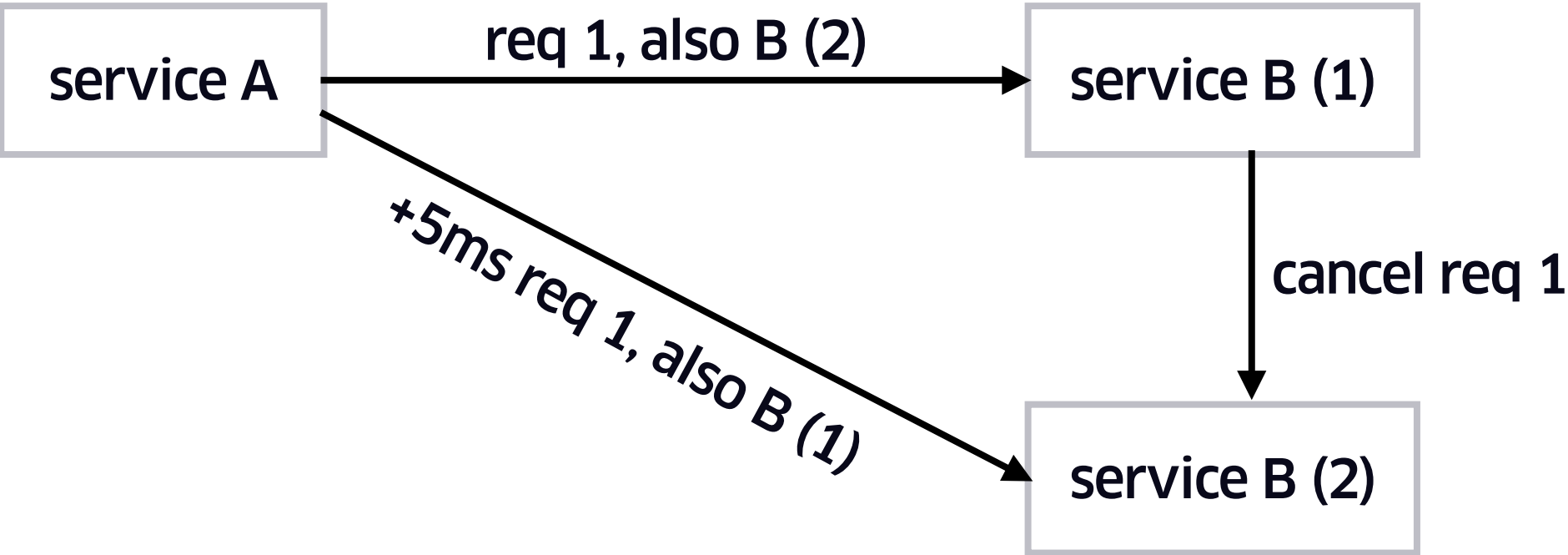
- overall latency ≥ latency of slowest component

- 1ms avg, 1000ms p99

  - use 1: 1% at least 1000ms

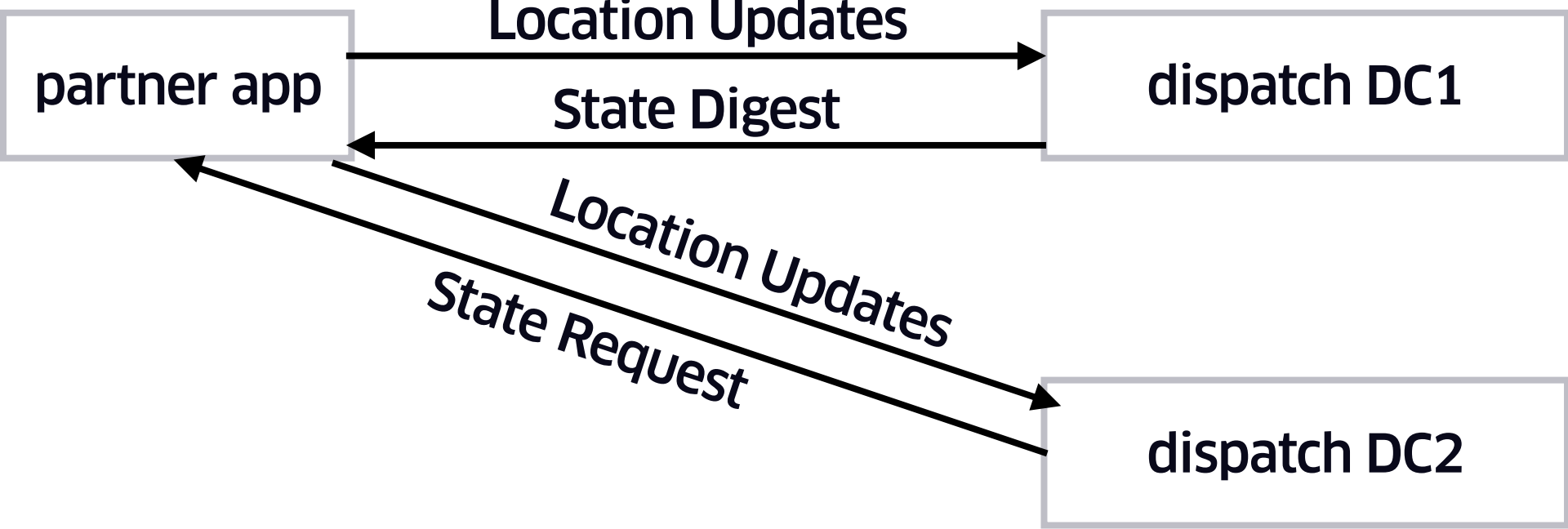  - use 100: 63% at least 1000ms

  - $1.0 - 0.99^{100} = 0.634 = 63.4\%$

# Google™

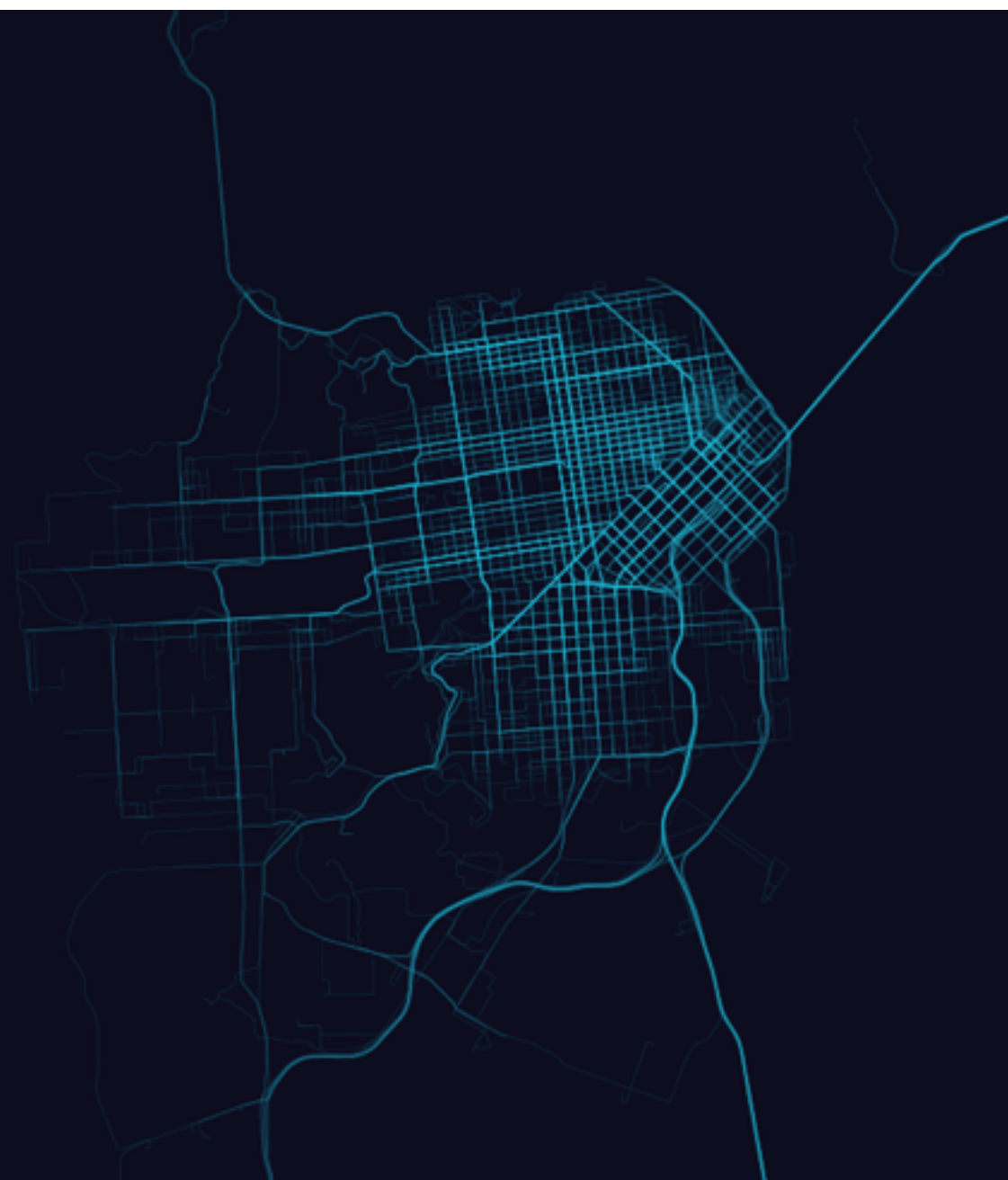# Achieving Rapid Response Times in Large Online Services

## Jeff Dean
## Google Fellow
jeff@google.com

# DATACENTER FAILURE

THANKS

UBER