

SCALING UBER'S REALTIME MARKET PLATFORM

QCON LONDON 2015

U B E R





WORLD SPORTS









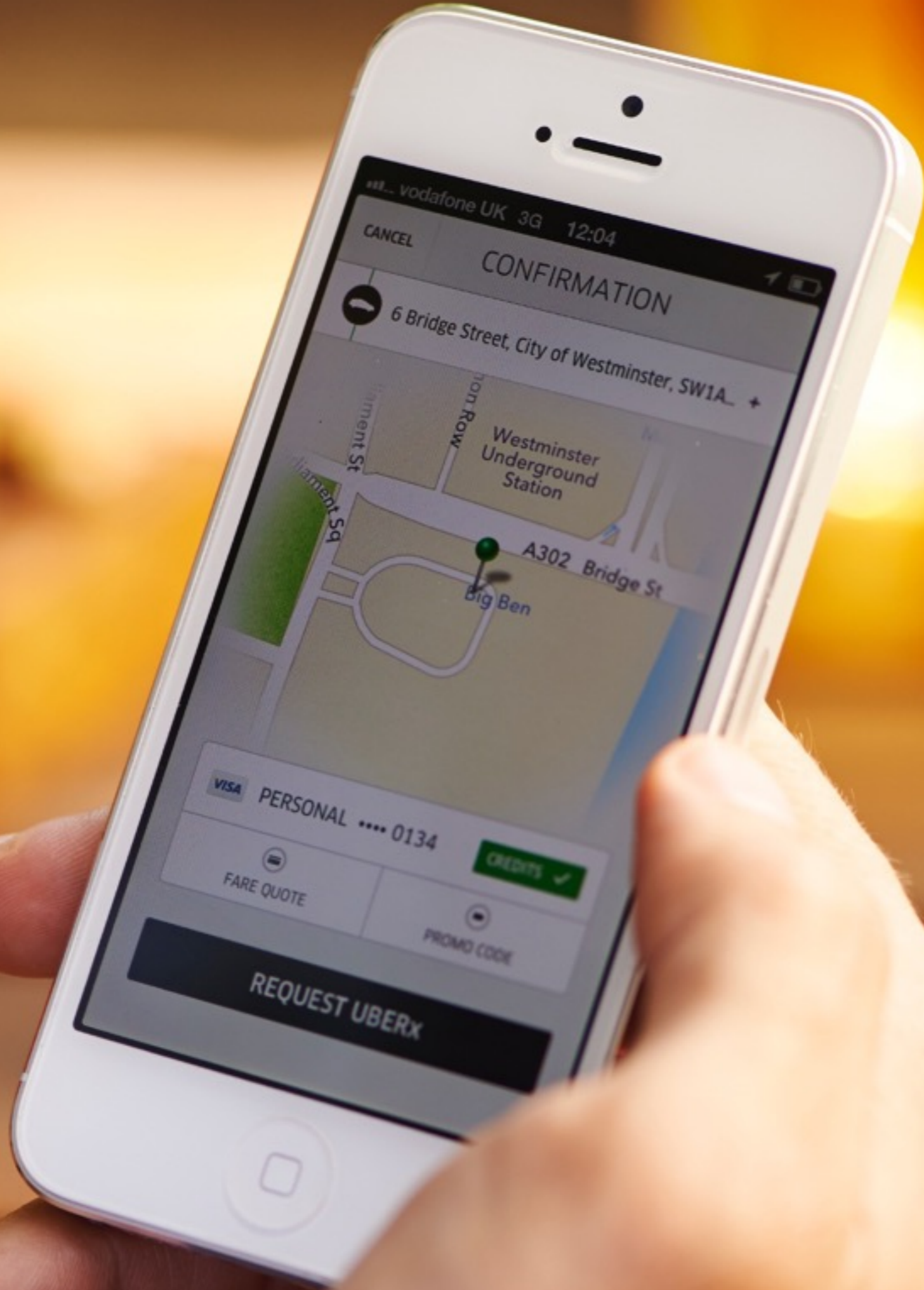








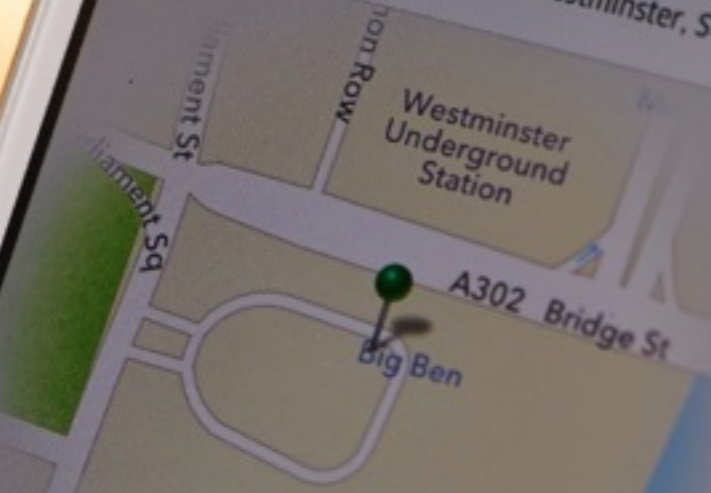




CANCEL

CONFIRMATION

6 Bridge Street, City of Westminster, SW1A



VISA PERSONAL **** 0134

CREDITS ✓

FARE QUOTE

PROMO CODE

REQUEST UBERX

partners

riders

dispatch

maps / ETA

services

**post trip
processing**

databases

money

partners

riders

dispatch

maps / ETA

services

**post trip
processing**

databases

money

partners

riders

dispatch

maps / ETA

services

**post trip
processing**

databases

money

partners

riders

dispatch

maps / ETA

services

post trip
processing

databases

money

partners

riders

dispatch

maps / ETA

services

**post trip
processing**

databases

money

MICROSERVICES

partners

riders

dispatch

maps / ETA

services

**post trip
processing**

databases

money

partners

riders

dispatch

maps / ETA

services

post trip
processing

databases

money

partners

riders

dispatch

maps / ETA

services

post trip
processing

databases

money

partners

riders

dispatch

maps / ETA

services

post trip
processing

databases

money

dispatch

maps / ETA

services

**post trip
processing**

databases

money

Things You Should Never Do, Part I

by Joel Spolsky

Thursday, April 06, 2000

Netscape 6.0 is finally going into its first public beta. There never was a version 5.0. The last major release, version 4.0, was released almost three years ago. Three years is an *awfully* long time in the Internet world. During this time, Netscape sat by, helplessly, as their market share plummeted.

It's a bit smarmy of me to criticize them for waiting so long between releases. They didn't do it *on purpose*, now, did they?

Well, yes. They did. They did it by making the **single worst strategic mistake** that any software company can make:

They decided to rewrite the code from scratch.

Netscape wasn't the first company to make this mistake. Borland made the same mistake when they bought Arago and tried to make it into dBase for Windows, a doomed project



that took so long that Microsoft Access ate their lunch, then they made

PROBLEMS

- **1 rider, 1 vehicle**
- **Moving people**
- **Sharding by city**
- **MPOF**

dispatch

maps / ETA

services

**post trip
processing**

databases

money

**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

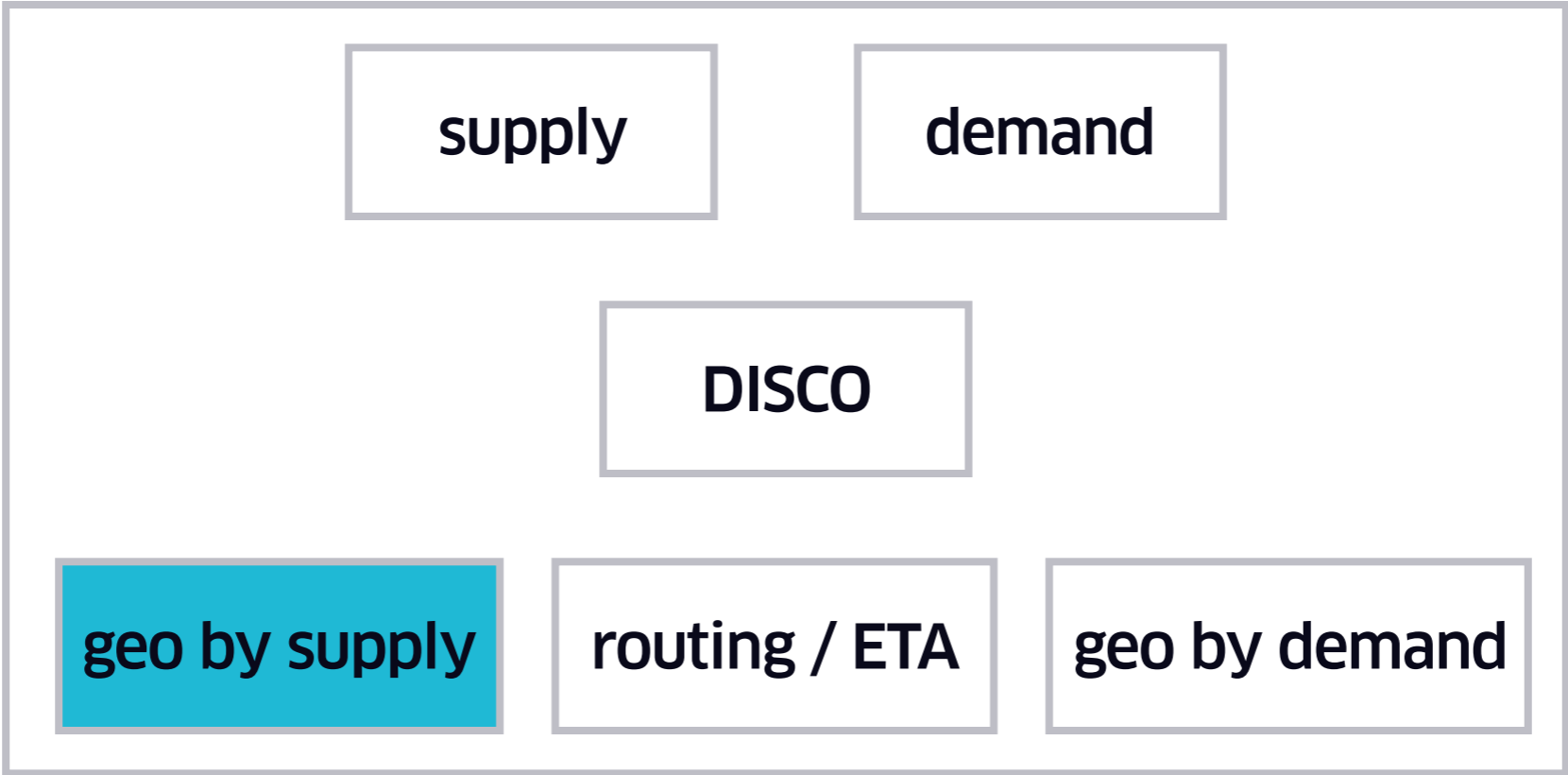
DISCO

geo by supply

routing / ETA

geo by demand

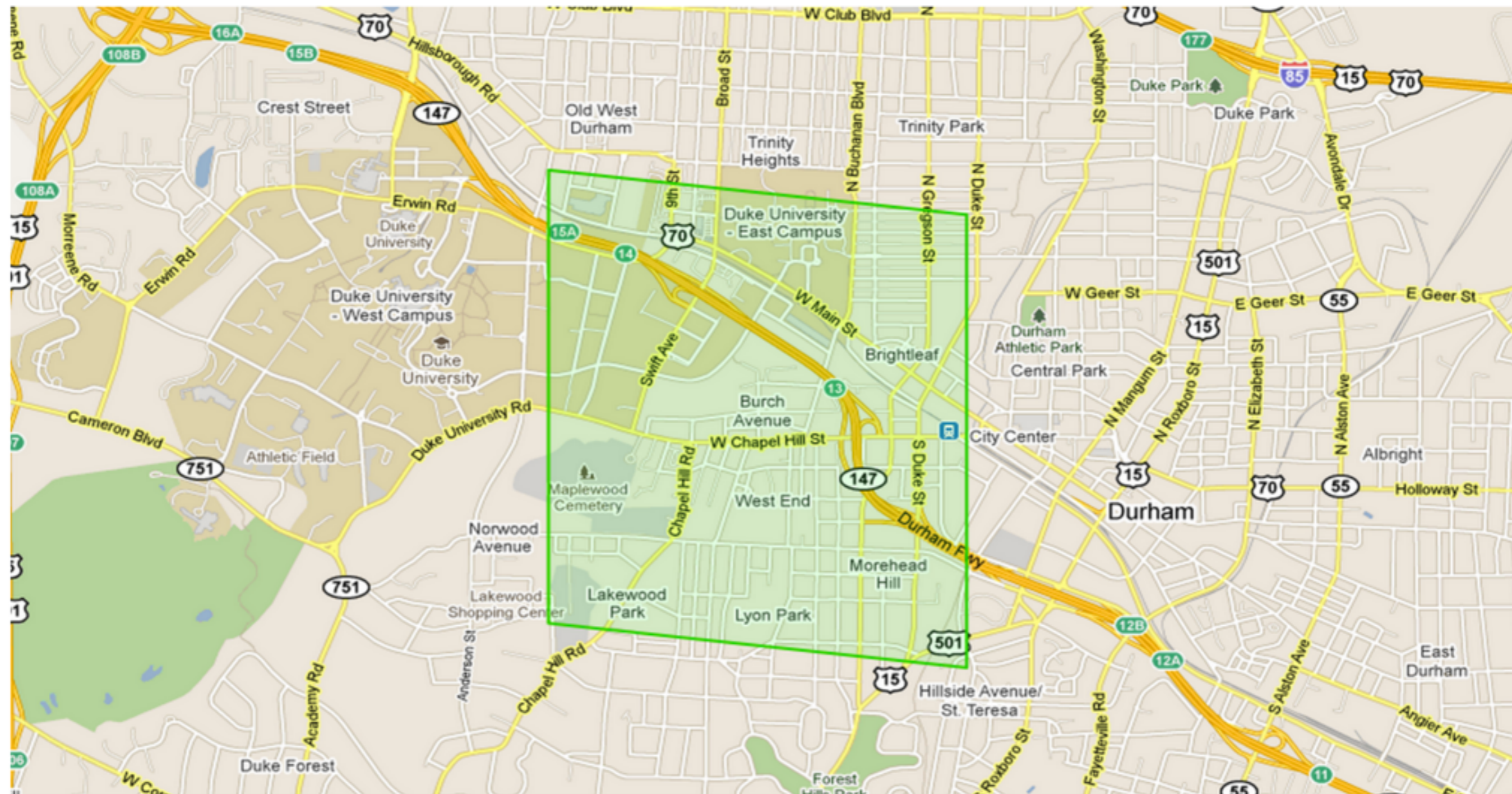
Dispatch



Dispatch

One S2 Cell


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



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

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 cm² on Earth can be represented using a 64-bit integer.

OSM Light

- Lat/Lng
- Lng/Lat

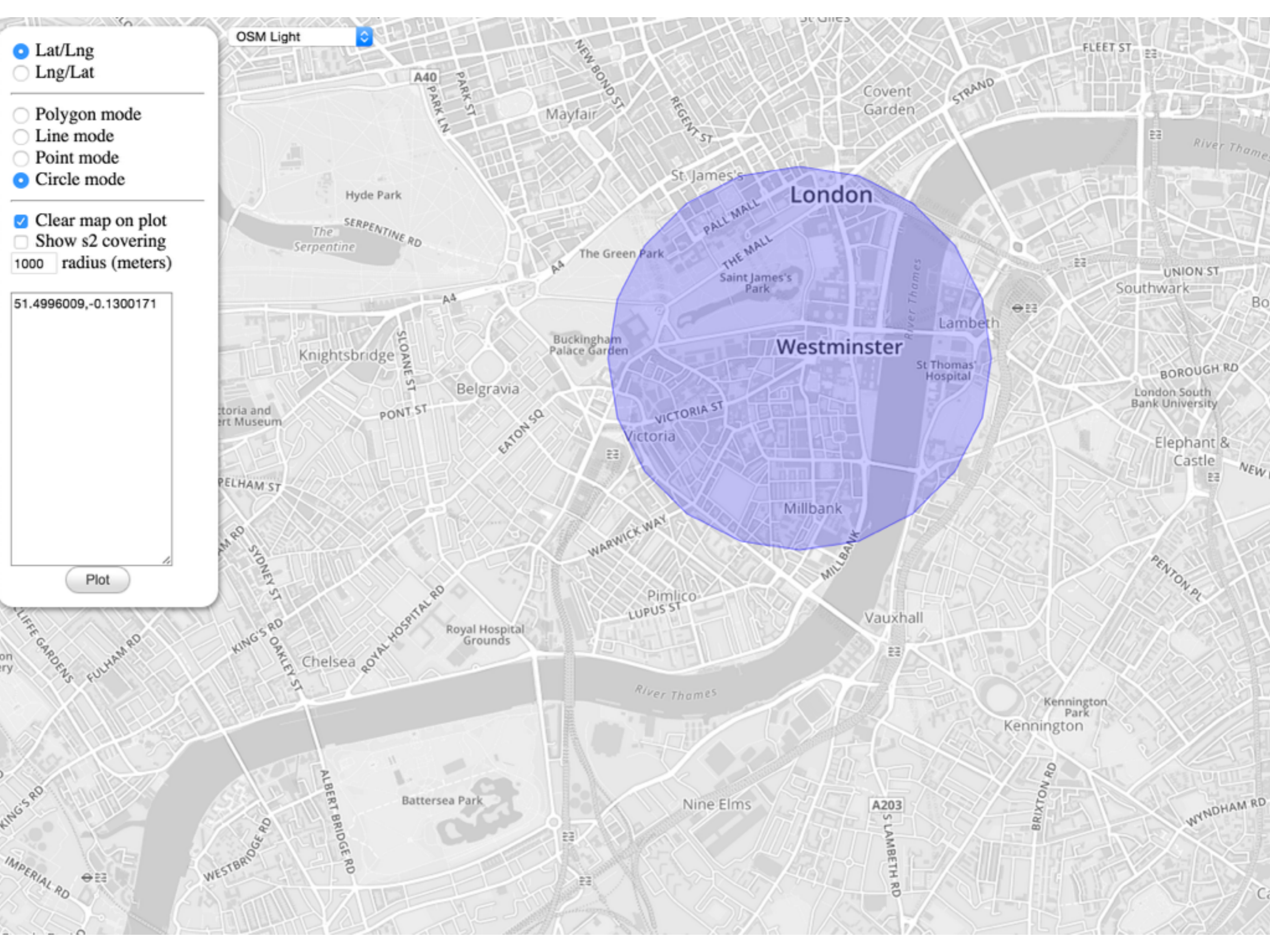
- Polygon mode
- Line mode
- Point mode
- Circle mode

- Clear map on plot
- Show s2 covering

1000 radius (meters)

51.4996009,-0.1300171

Plot



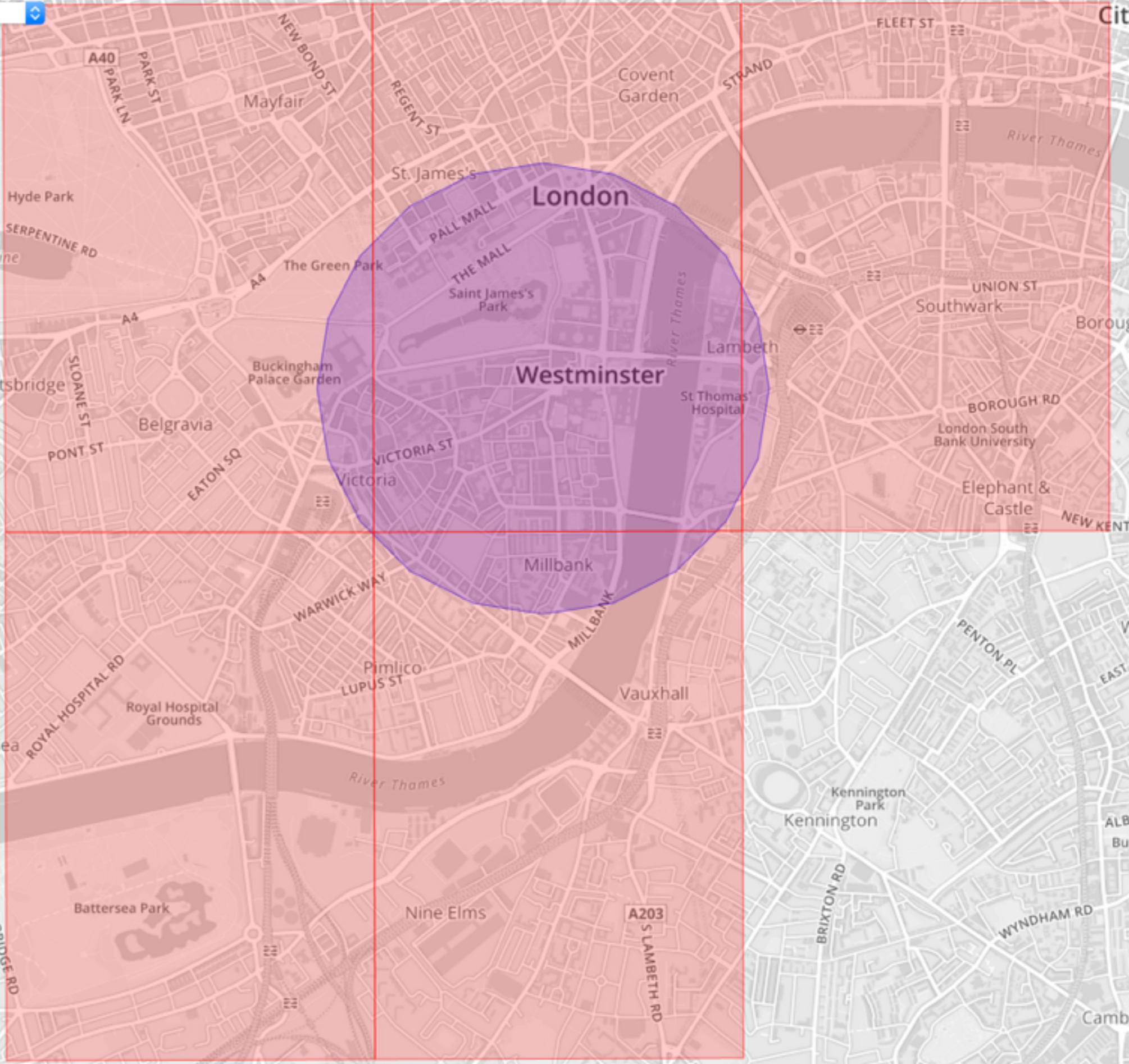
- Lat/Lng
- Lng/Lat
- Polygon mode
- Line mode
- Point mode
- Circle mode

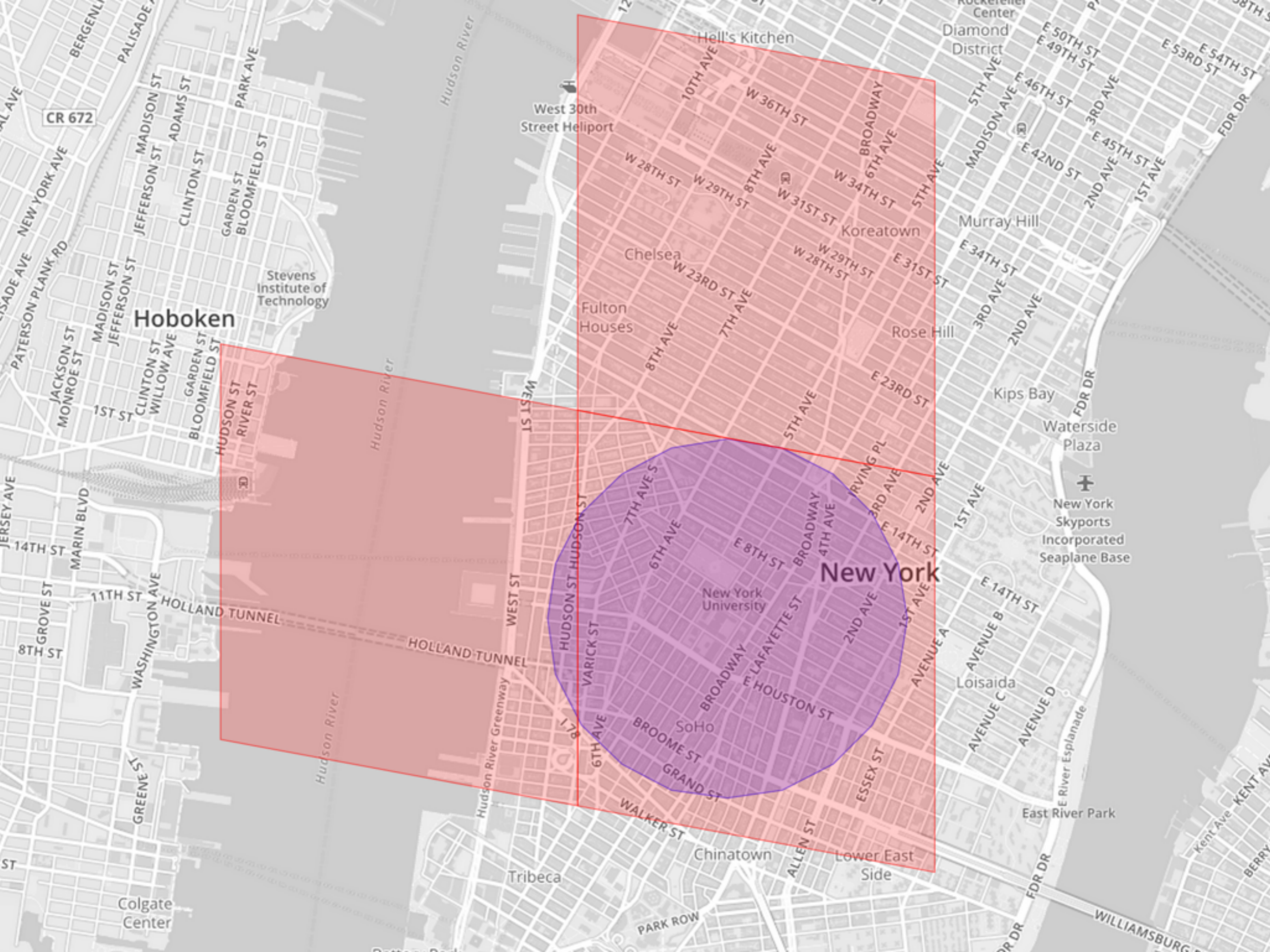
- Clear map on plot
- Show s2 covering
- 12 min_level
- 12 max_level
- 200 max_cells
- 1 level_mod
- 1000 radius (meters)

51.4996009,-0.1300171

Plot

OSM Light





Hoboken

New York

CR 672

West 30th Street Heliport

Stevens Institute of Technology

Hell's Kitchen

Rocket Center

Diamond District

Koreatown

Murray Hill

Chelsea

Fulton Houses

Rose Hill

Kips Bay

Waterside Plaza

New York Skyports Incorporated

Seaplane Base

New York University

SoHo

Loisaida

East River Park

Tribeca

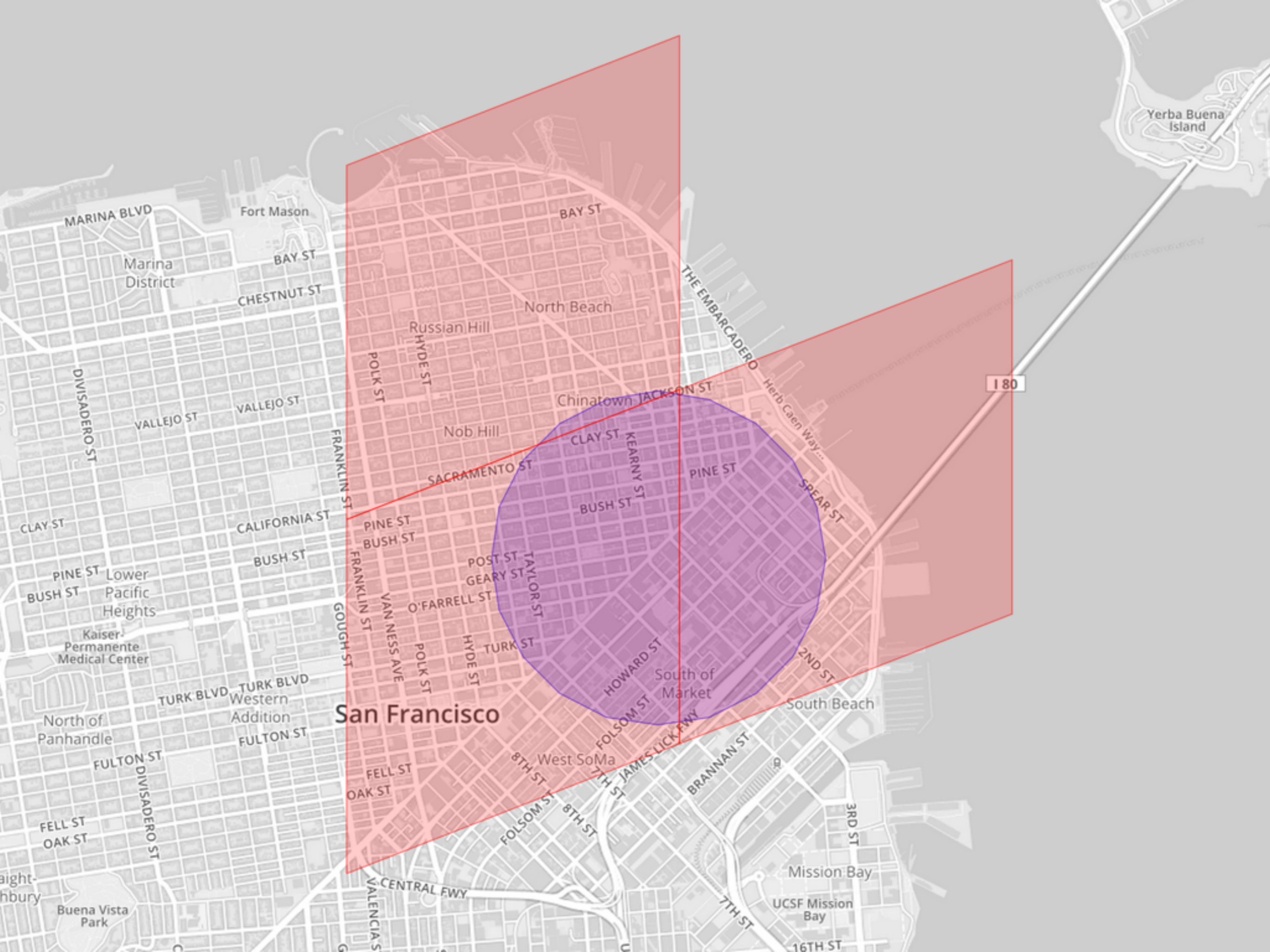
Chinatown

Lower East Side

Colgate Center

Park Row

Williamsburg



San Francisco

I 80

Yerba Buena Island

Fort Mason

Marina District

North Beach

Russian Hill

Chinatown

Nob Hill

South of Market

South Beach

Western Addition

West SoMa

Mission Bay

UCSF Mission Bay

MARINA BLVD

BAY ST

CHESTNUT ST

BAY ST

VALLEJO ST

VALLEJO ST

FRANKLIN ST

POLK ST

HYDE ST

THE EMBARCADERO

JACKSON ST

CLAY ST

KEARNY ST

PINE ST

SACRAMENTO ST

BUSH ST

SPEAR ST

CLAY ST

CALIFORNIA ST

PINE ST

BUSH ST

POST ST

TAYLOR ST

O'FARRELL ST

TURK ST

Lower Pacific Heights

Kaiser-Permanente Medical Center

FRANKLIN ST

VAN NESS AVE

POLK ST

HYDE ST

HOWARD ST

2ND ST

North of Panhandle

TURK BLVD

TURK BLVD

FULTON ST

San Francisco

FELL ST

OAK ST

8TH ST

7TH ST

8TH ST

FOLSOM ST

BRANNAN ST

JAMES LICK FWY

3RD ST

CENTRAL FWY

7TH ST

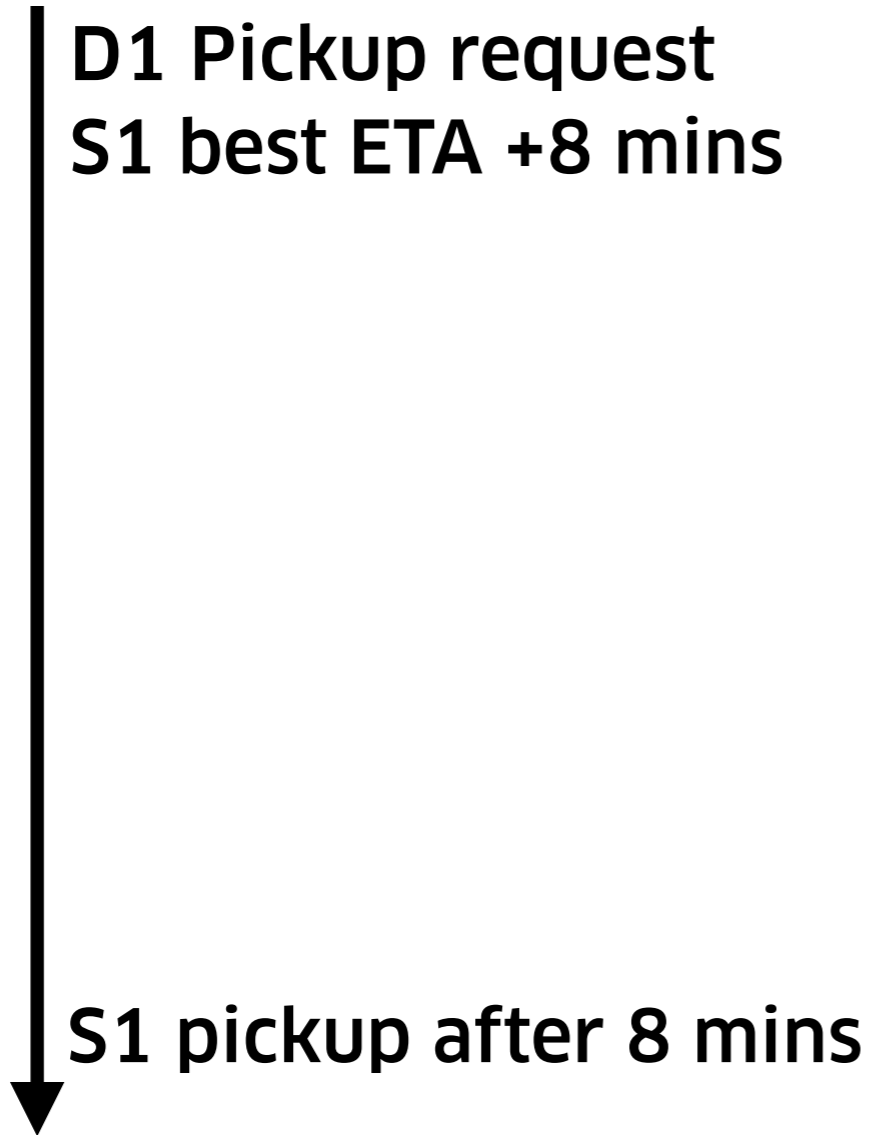
16TH ST

Buena Vista Park

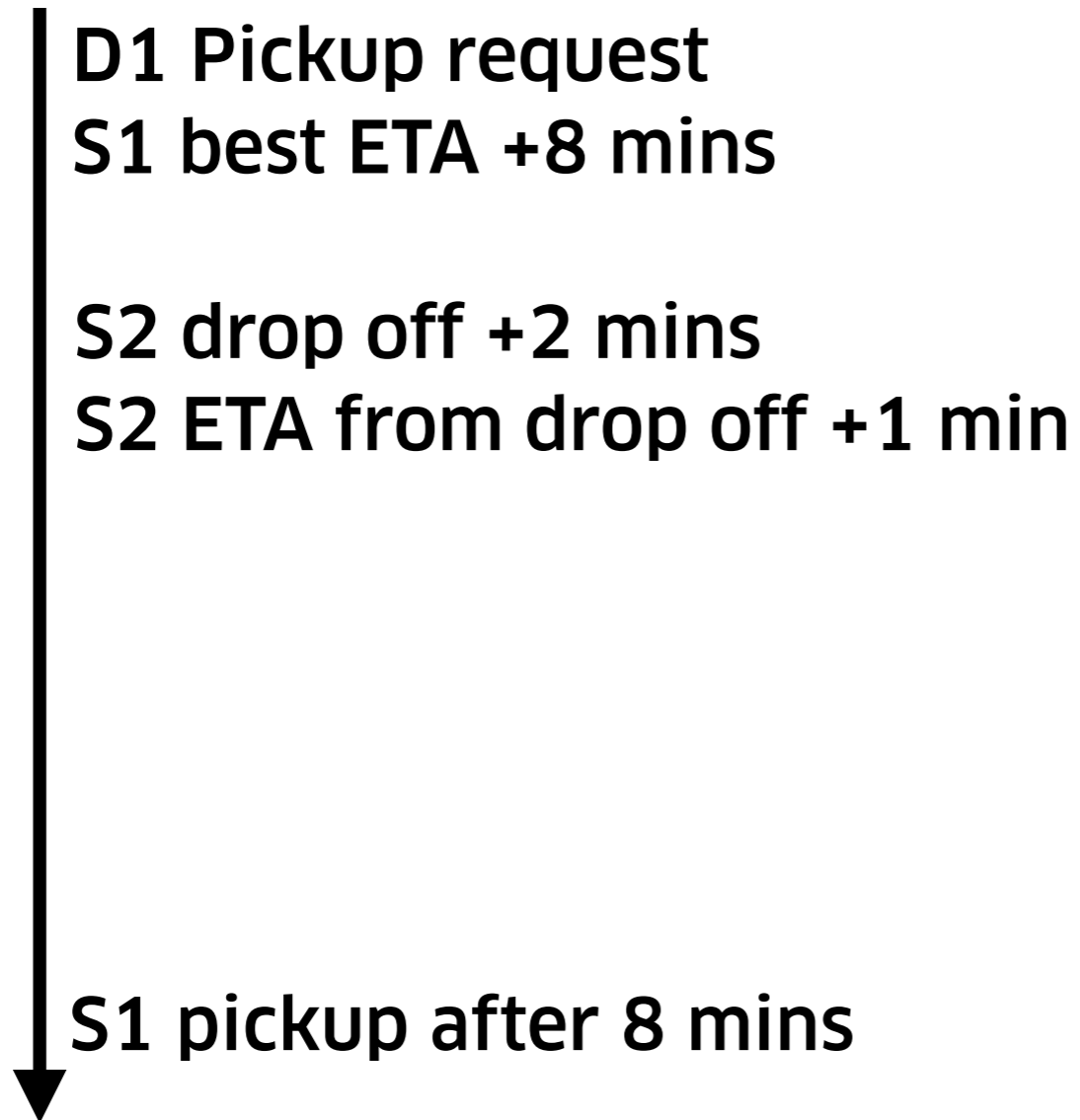
GOALS

- **reduce waiting**
- **reduce extra driving**
- **lowest overall ETAs**

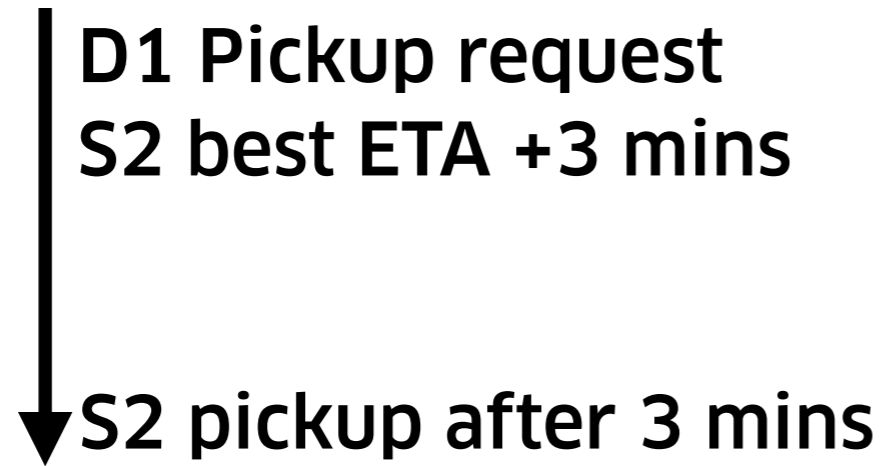
time



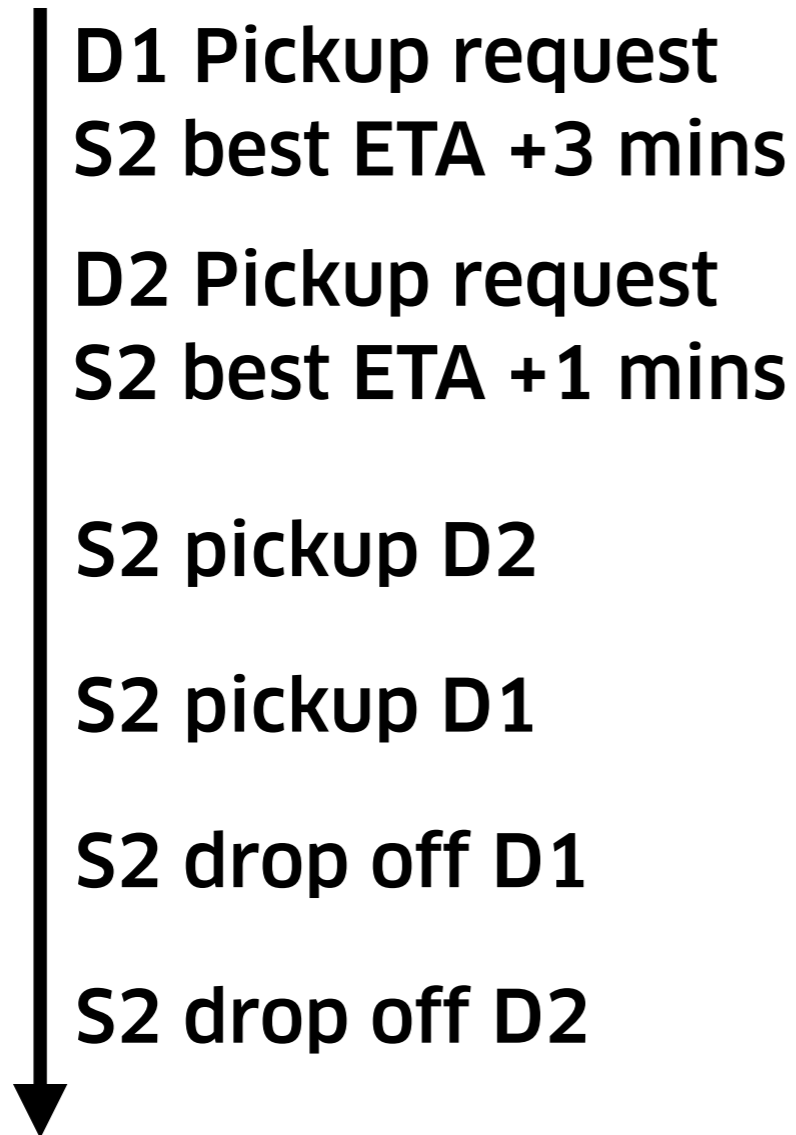
time

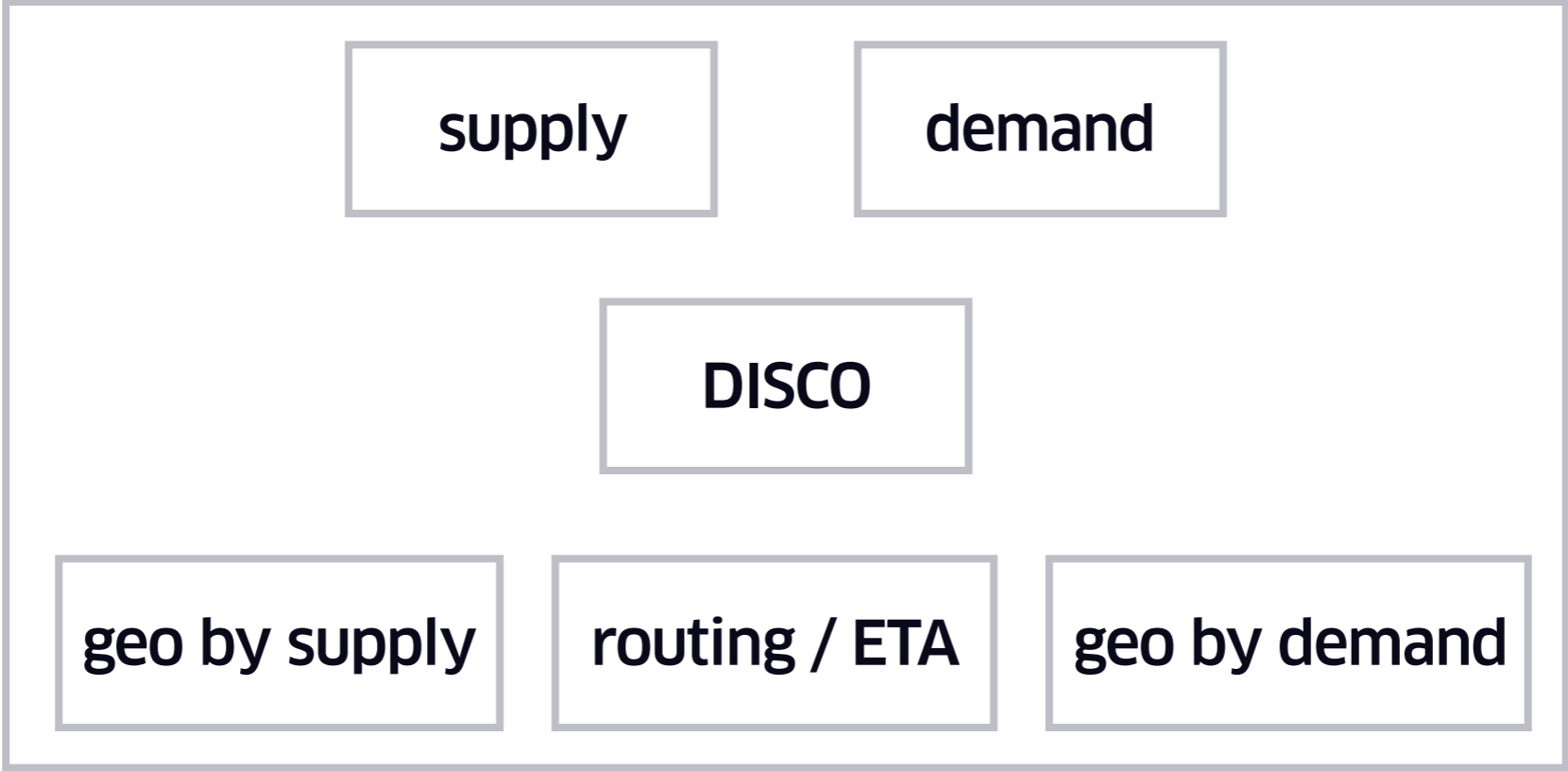


time



time





Dispatch

ringpop

ringpop brings application-layer sharding to your services in a fault tolerant and scalable manner. It is an embeddable server that reliably partitions your data, detects node failures and easily integrates new nodes into your application cluster when they become available. For more information about the techniques applied within ringpop, see the Concepts section below.

Table of Contents

- [Motivation](#)
- [Concepts](#)
- [Developer's Guide](#)
- [Operator's Guide](#)
- [Community](#)
- [References](#)
- [Installation](#)

Motivation

As an organization's architecture grows in complexity engineers must find a way to make their services more resilient while keeping operational overhead low. ringpop is a step in that direction and an effort to generalize the sharding needs of various services by providing a simple hash ring abstraction. We've found that the use cases to which ringpop can be applied are numerous and that new ones are discovered often.

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

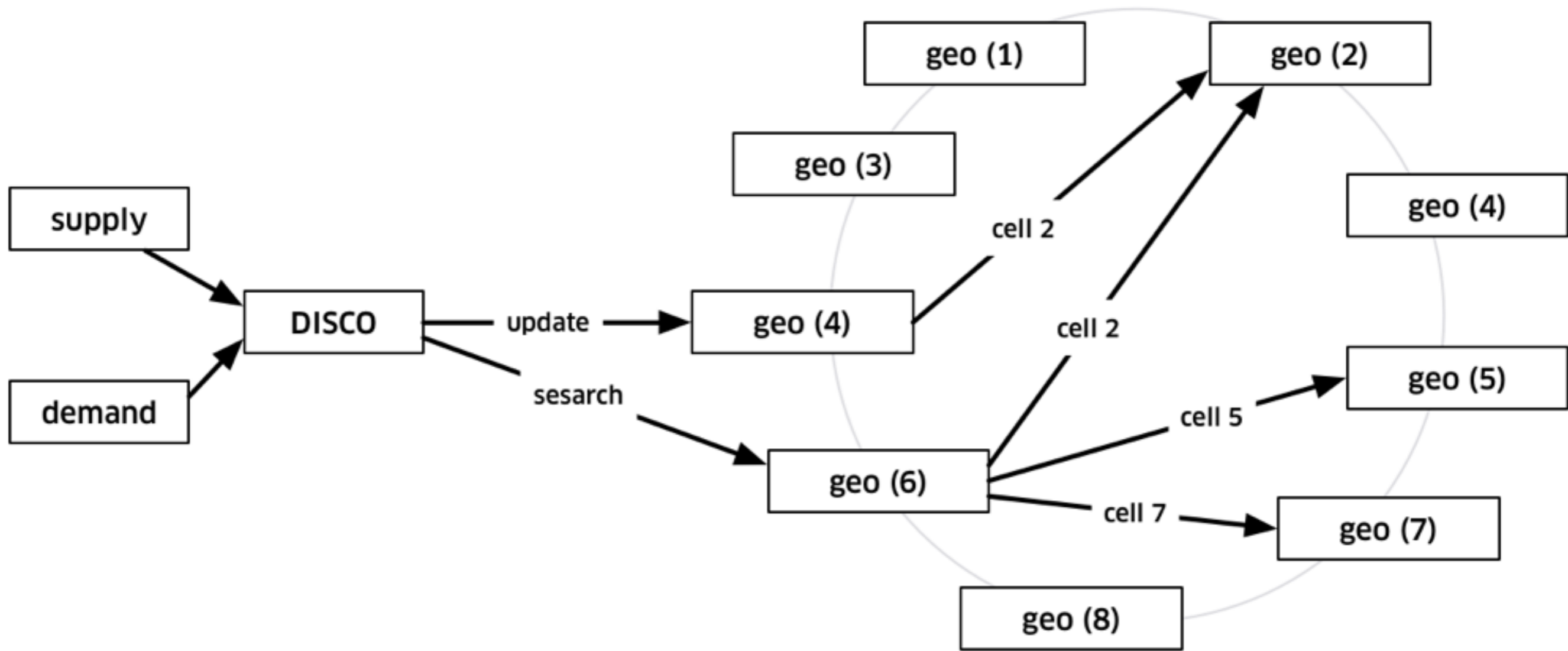
Abstract




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



 .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

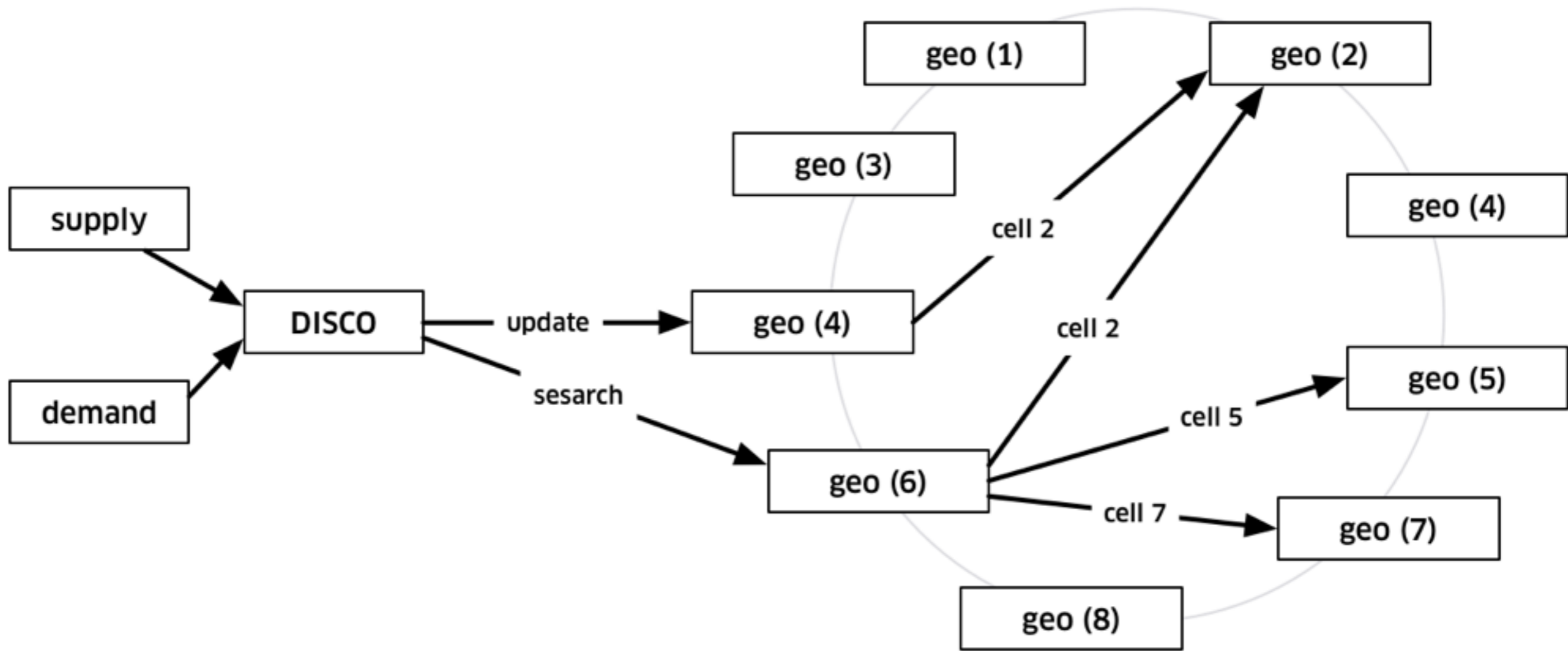
Design goals

- 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**



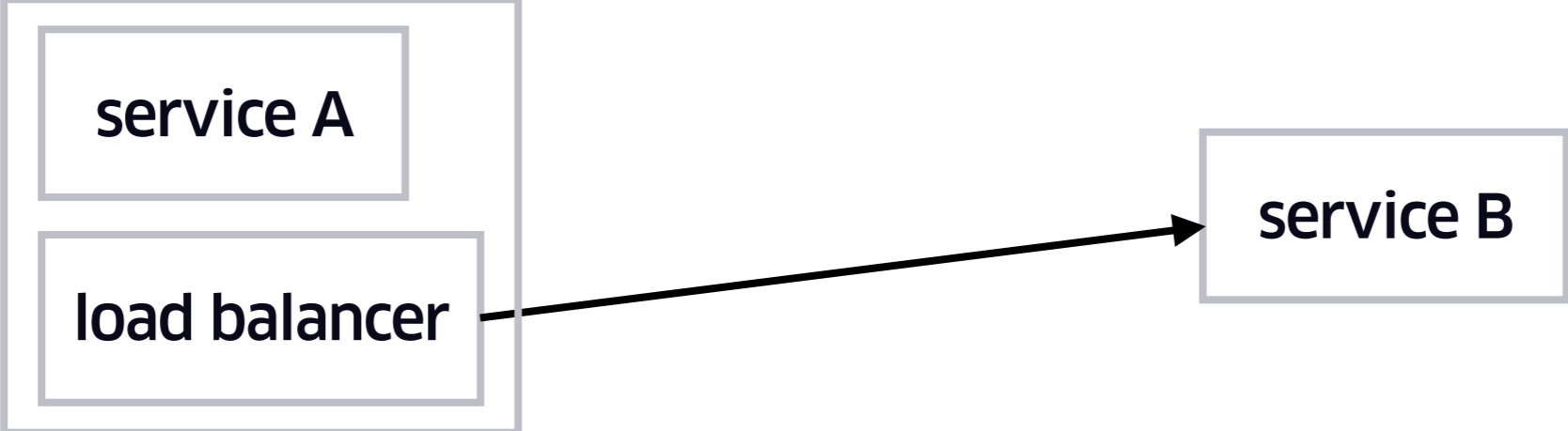
AVAILABILITY

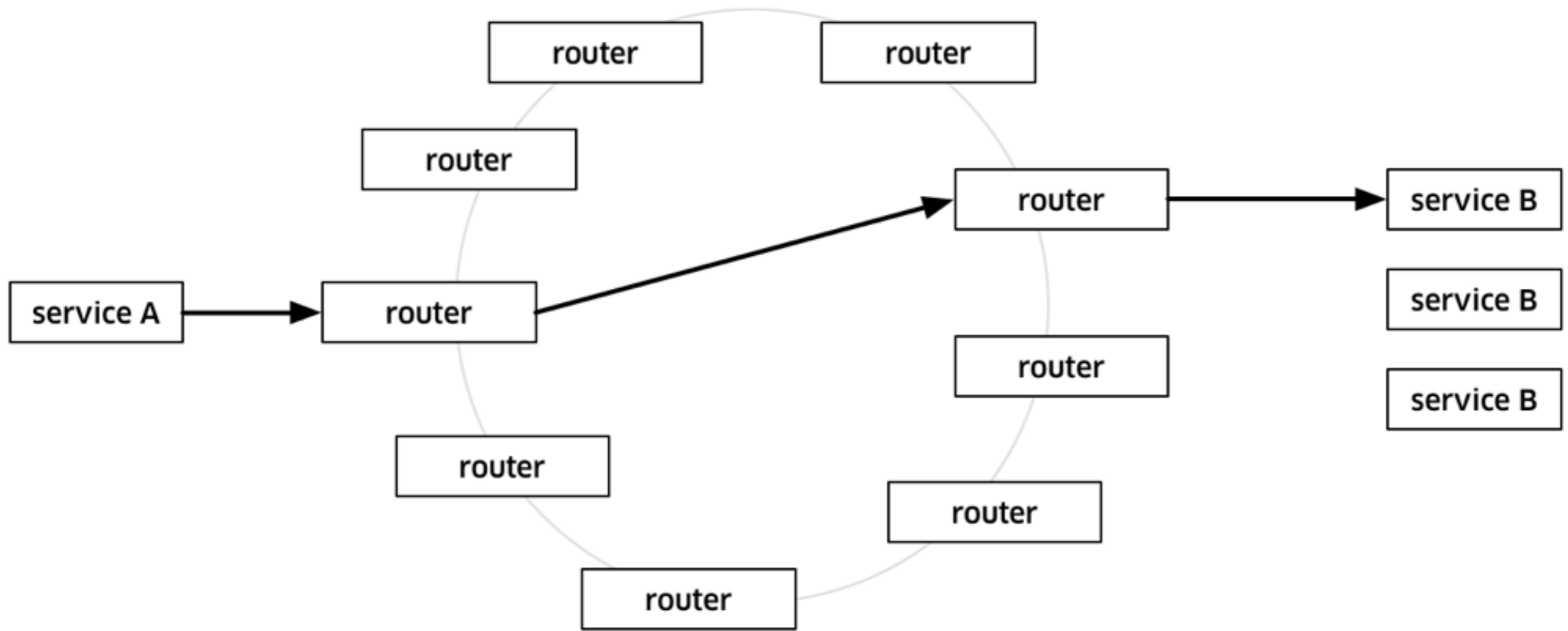
- **everything retryable**
- **everything killable**
- **crash only**
- **small pieces**

CULTURAL CHANGES

- **no pairs**
- **kill everything**
- **even databases**







backup requests cross server cancellation



Web

Shopping

Videos

News

Images

More ▾

Search tools

About 33,500,000 results (0.61 seconds)

[\[PDF\] Achieving Rapid Response Times in Large Online Services](#)

research.google.com/people/jeff/Berkeley-Latency-Mar2012.pdf ▾

by J Dean - [Cited by 18](#) - [Related articles](#)

26 Mar 2012 - **Backup Requests w/ Cross-Server Cancellation**. Server 1. Client.

Server 2 req 3 req 6 req 5. Monday, March 26, 2012 ...

[Google on Latency Tolerant Systems: Making a Predictable ...](#)

highscalability.com/.../google-on-latency-tolerant-systems-making-a-pre... ▾

18 Jun 2012 - 3 posts - 3 authors

If a **request** has to access 100 servers, now 63% of all **requests** will take over a second.

... **Backup Requests with Cross-Server Cancellation**.

[\[PDF\] Jeffrey Dean \(Google, Inc.\) - Computing Research Associ...](#)

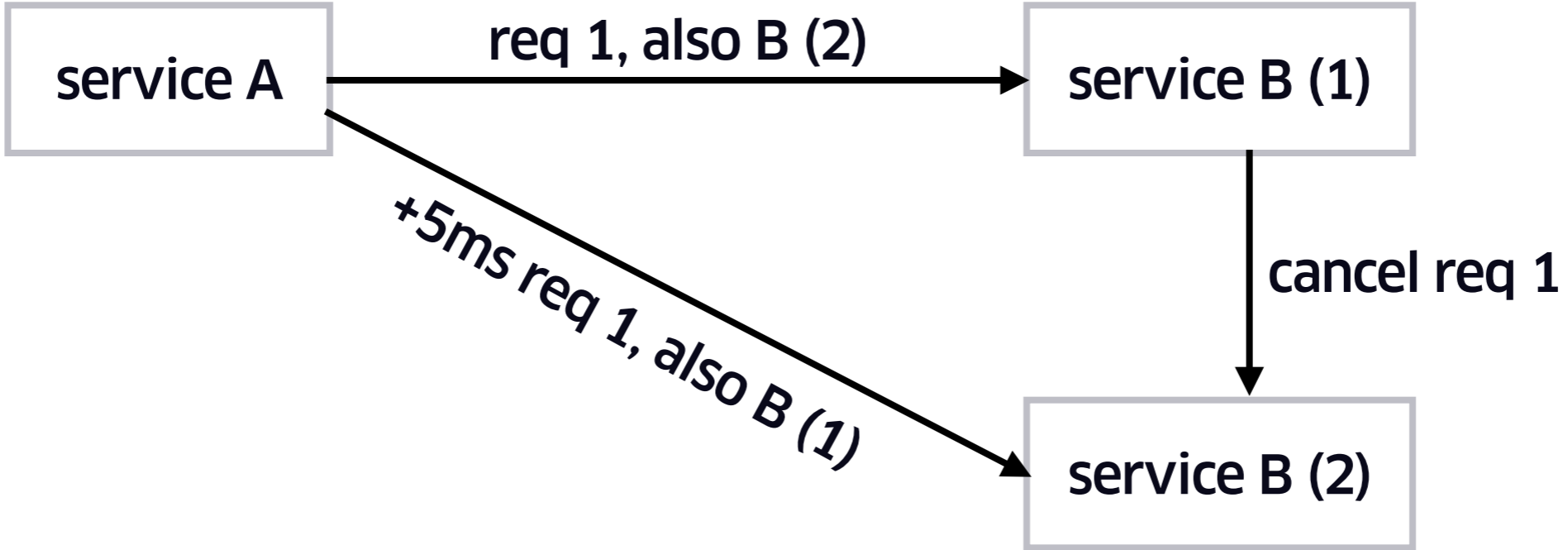
cra.org/uploads/documents/resources/snowbird2012_slides/dean.pdf ▾

Backup Requests w/ Cross-Server Cancellation. Server 1. Client. Server 2 req 3 req 6

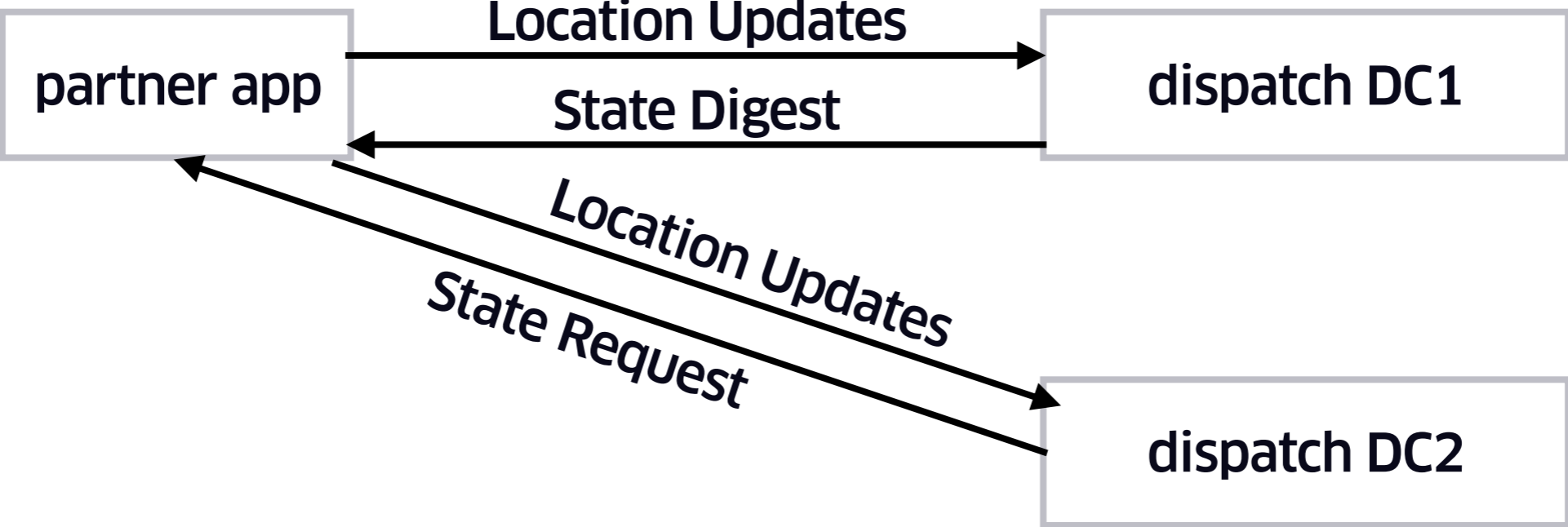
req 5. Similar to Michael Mitzenmacher's work on "The Power of Two.

LATENCY

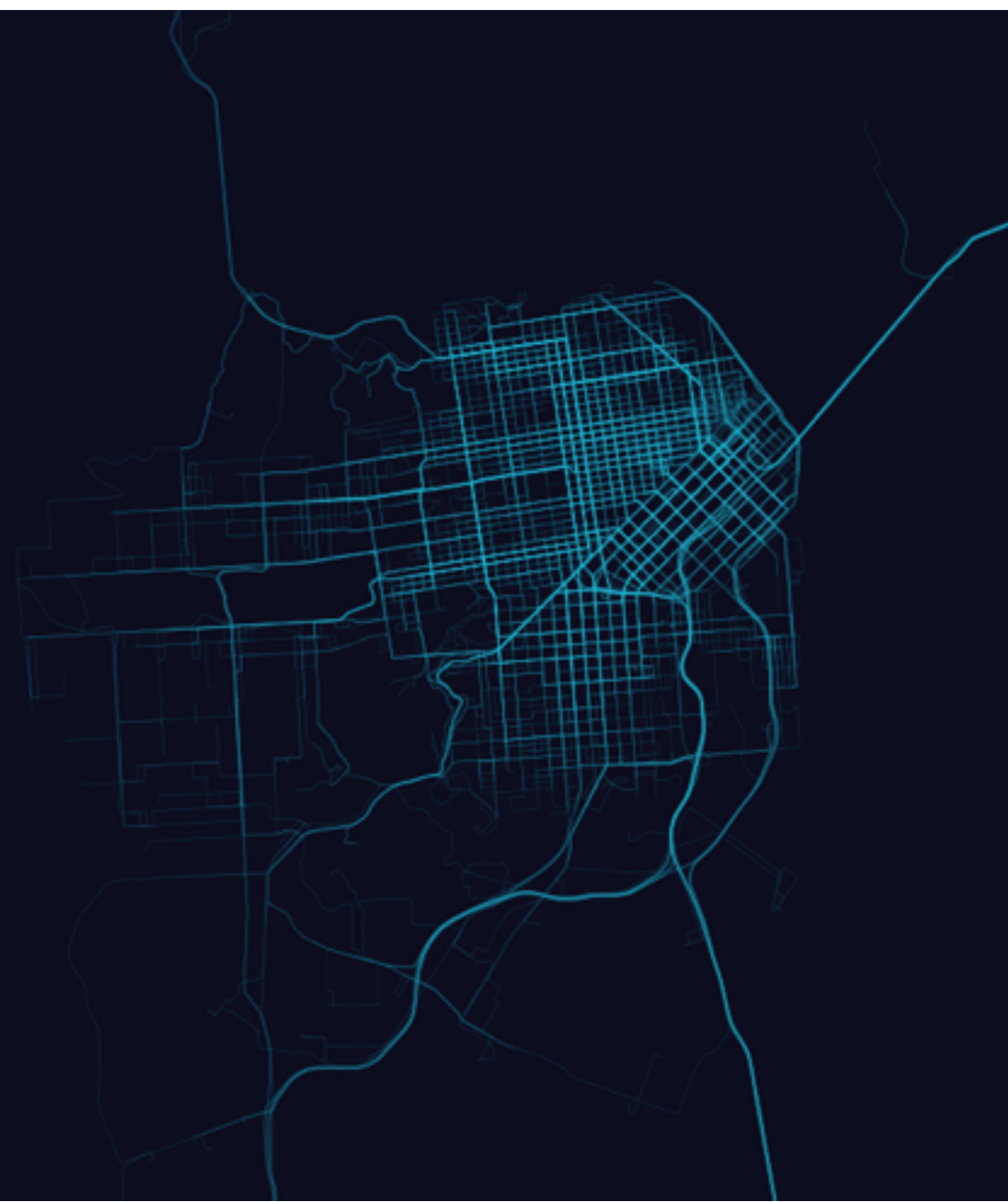
- **overall latency \geq latency of slowest component**
- **1ms avg, 1000ms p99**
- **use 1: 1% at least 1000ms**
- **use 100: 63% at least 1000ms**
- **$1.0 - 0.99^{100} = 0.634 = 63.4\%$**



DATACENTER FAILURE



THANKS



U B E R