



# 优步流处理应用

袁泳 @ UBER

U B E R

# 优步简介





# 便捷交通，弹指之间





# 流数据 - 城市的脉搏

# 市场动态

当前全球有多少车可用



Query

Table

Heatmap

Stats





过去**10分钟**里旧金山里每个区域有多少车载客？

Query

Download as JSON

Download as CSV

Table

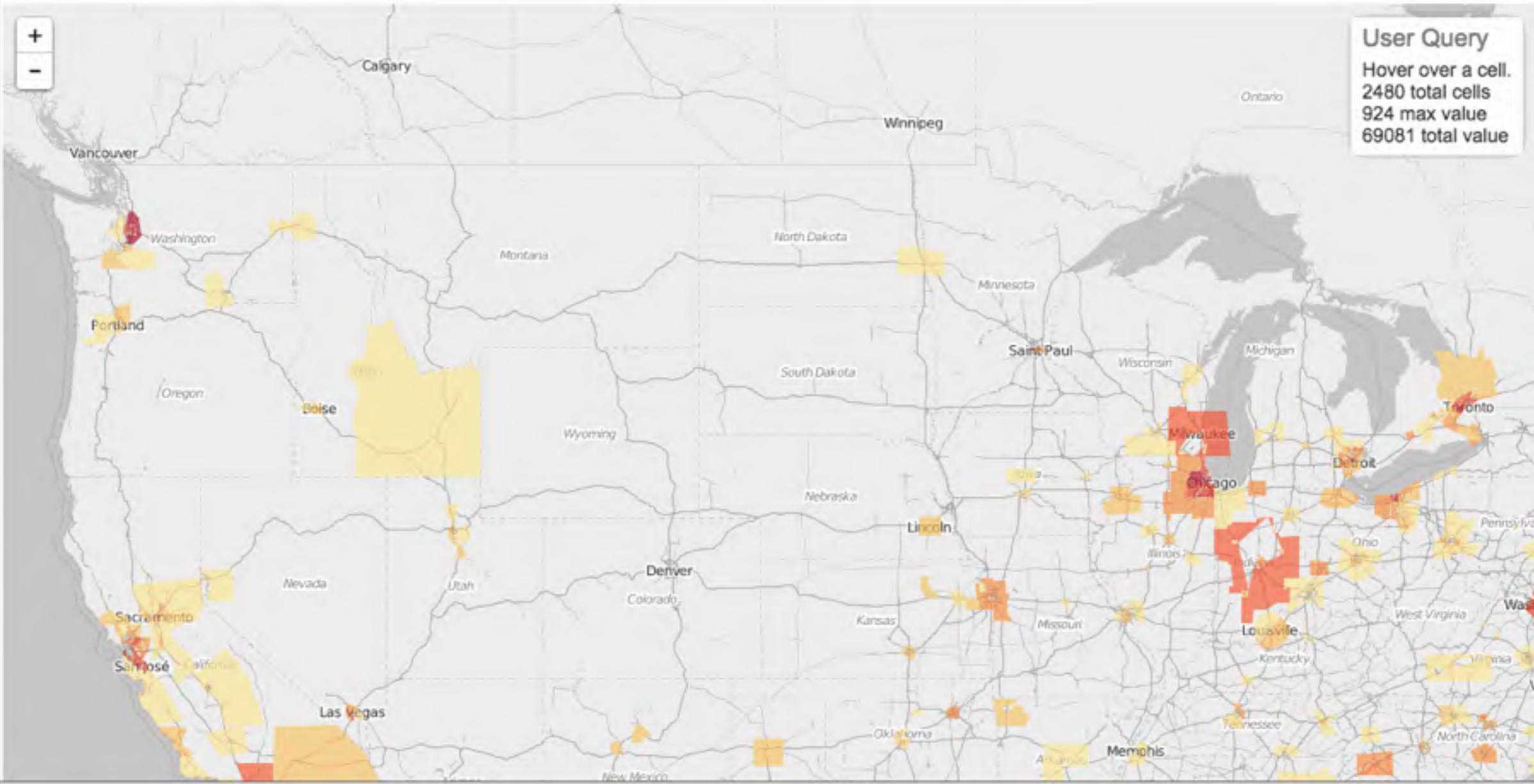
Heatmap

Stats



User Query

Hover over a cell.  
2480 total cells  
924 max value  
69081 total value



过去**10分钟**里**旧金山**里**每个六边形**有多少**车载客**？





# 历史变化

SAN FRANCISCO ALLEN KEY

## MARKETPLACE HEALTH (ALPHA)

LOCATION: San Francisco Bay Area

VEHICLE TYPE: UberX 8

FROM DATE: 11/15/15 8:00pm

TO DATE: 11/16/15 8:00pm

INTERVAL: 1h

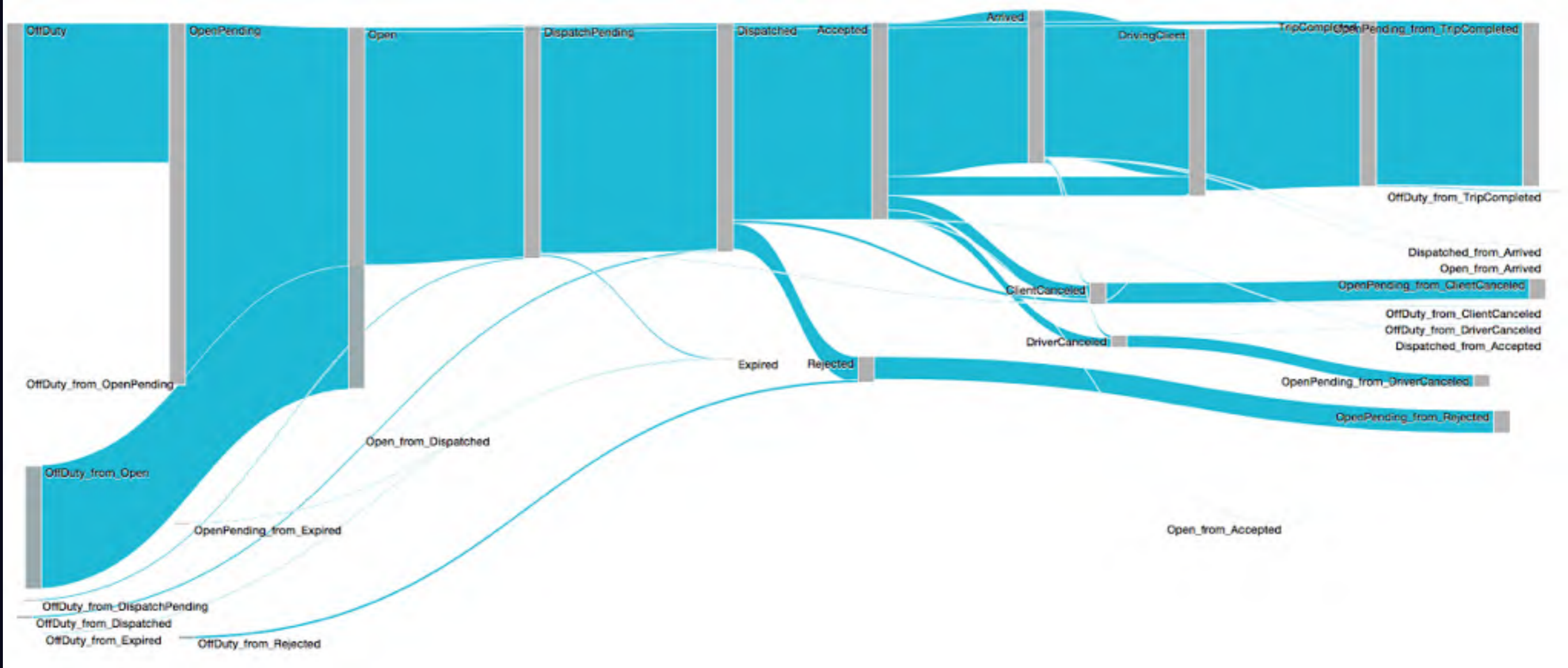
QUERY

SELECTION RADIUS: 0 1 2 3 4 5 6 7

No metrics loaded. Select hexagons to load metrics

Note: Hexagon k-ring selection performs a second query on your selection of marketplace health related metrics aggregating values over the selected hexagons only.

# 状态跟踪

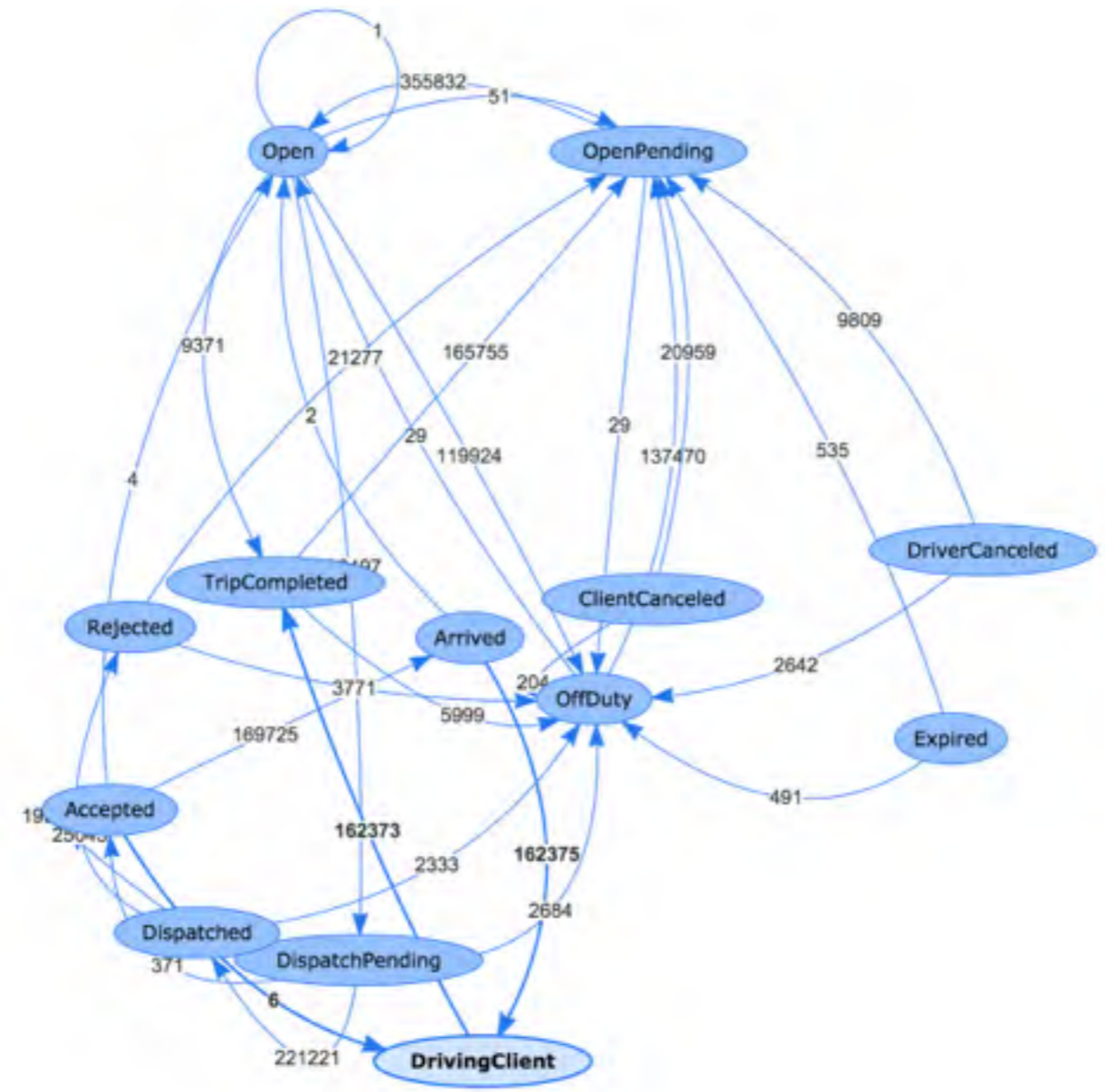


# Driver States

Choose a city

Driver ID

Click to Query



Dispatch Query Service

# Driver States

Choose a city:

Driver ID:

[Click to Query](#)

```
graph TD; OpenPending((OpenPending)) -- 7 --> Open((Open)); Open -- 4 --> OpenPending; Open -- 3 --> OffDuty((OffDuty)); OffDuty -- 1 --> Open; OffDuty -- 1 --> OpenPending; OffDuty -- 1 --> DispatchPending((DispatchPending)); DispatchPending -- 1 --> OffDuty; DispatchPending -- 1 --> Rejected((Rejected)); Rejected -- 1 --> DispatchPending; Rejected -- 1 --> Dispatched((Dispatched)); Dispatched -- 1 --> Rejected; Dispatched -- 2 --> Accepted((Accepted)); Accepted -- 2 --> Dispatched; Accepted -- 2 --> Arrived((Arrived)); Arrived -- 2 --> Accepted; Arrived -- 1 --> ClientCanceled((ClientCanceled)); ClientCanceled -- 1 --> Arrived; ClientCanceled -- 1 --> TripCompleted((TripCompleted)); TripCompleted -- 1 --> ClientCanceled; TripCompleted -- 1 --> OffDuty; TripCompleted -- 1 --> OpenPending; OpenPending -- 1 --> TripCompleted; TripCompleted -- 5 --> Open; Open -- 1 --> TripCompleted; DrivingClient((DrivingClient)) -- 2 --> Arrived; Arrived -- 2 --> DrivingClient;
```

The diagram illustrates the state transitions for a driver. States are represented by blue ovals, and transitions are shown as blue arrows with numerical labels. The states and their transitions are:

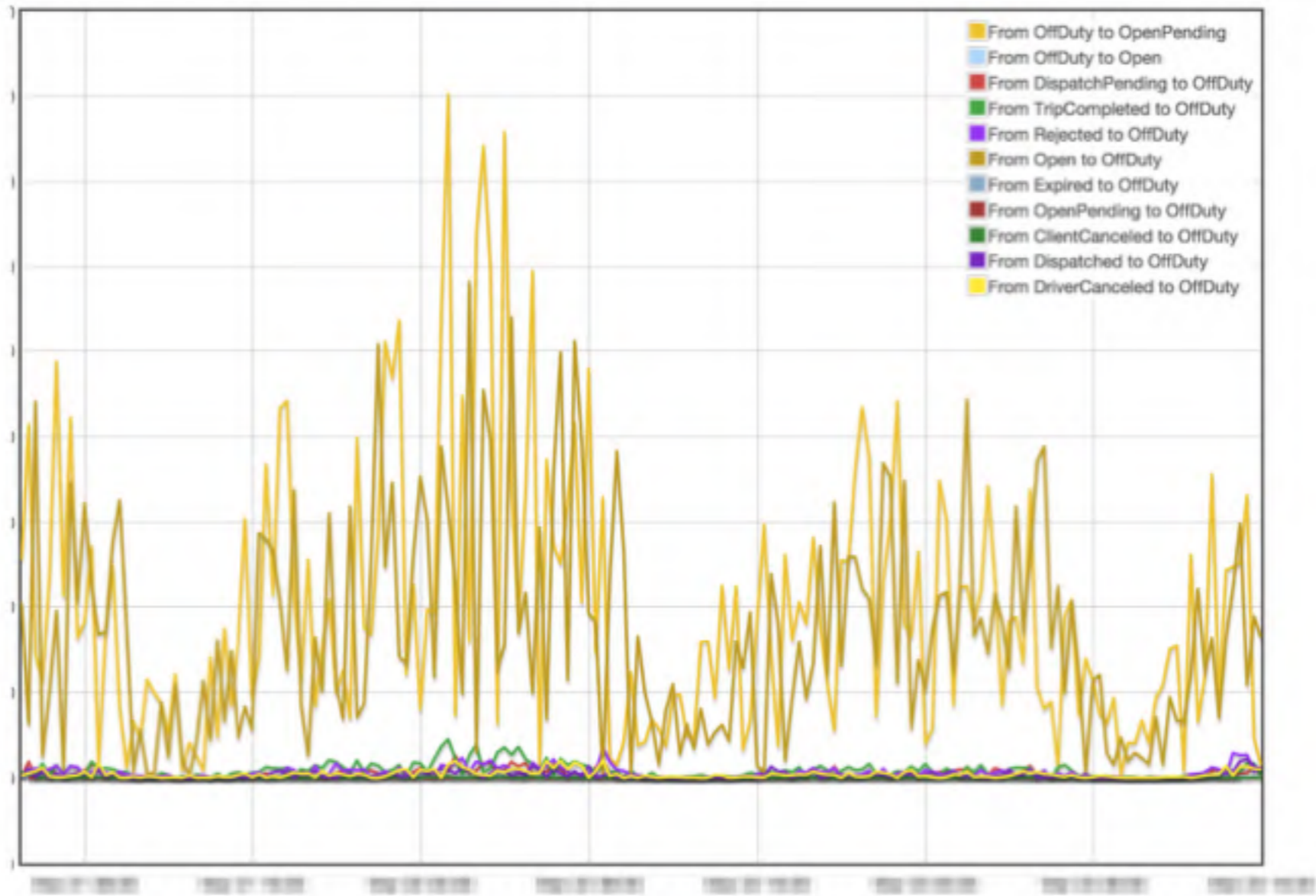
- OpenPending** transitions to **Open** (7) and **TripCompleted** (1).
- Open** transitions to **OpenPending** (4), **OffDuty** (3), and **TripCompleted** (1).
- OffDuty** transitions to **Open** (1), **OpenPending** (1), and **DispatchPending** (1).
- DispatchPending** transitions to **OffDuty** (1) and **Rejected** (1).
- Rejected** transitions to **DispatchPending** (1) and **Dispatched** (1).
- Dispatched** transitions to **Rejected** (1) and **Accepted** (2).
- Accepted** transitions to **Dispatched** (2) and **Arrived** (2).
- Arrived** transitions to **Accepted** (2) and **ClientCanceled** (1).
- ClientCanceled** transitions to **Arrived** (1) and **TripCompleted** (1).
- TripCompleted** transitions to **ClientCanceled** (1), **OffDuty** (1), and **OpenPending** (1).
- DrivingClient** transitions to **Arrived** (2).
- Arrived** transitions to **DrivingClient** (2).



59

Driver S

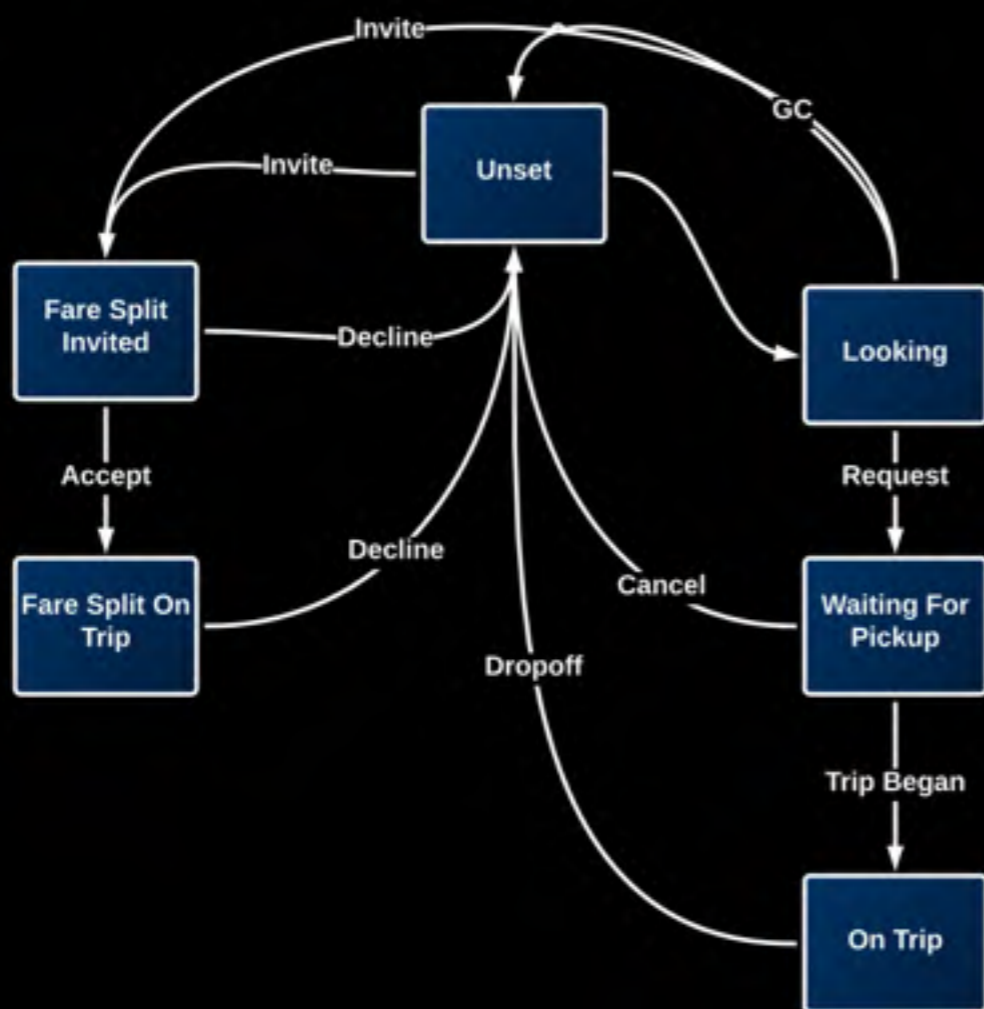
### Count of State Change Over Time



出发点

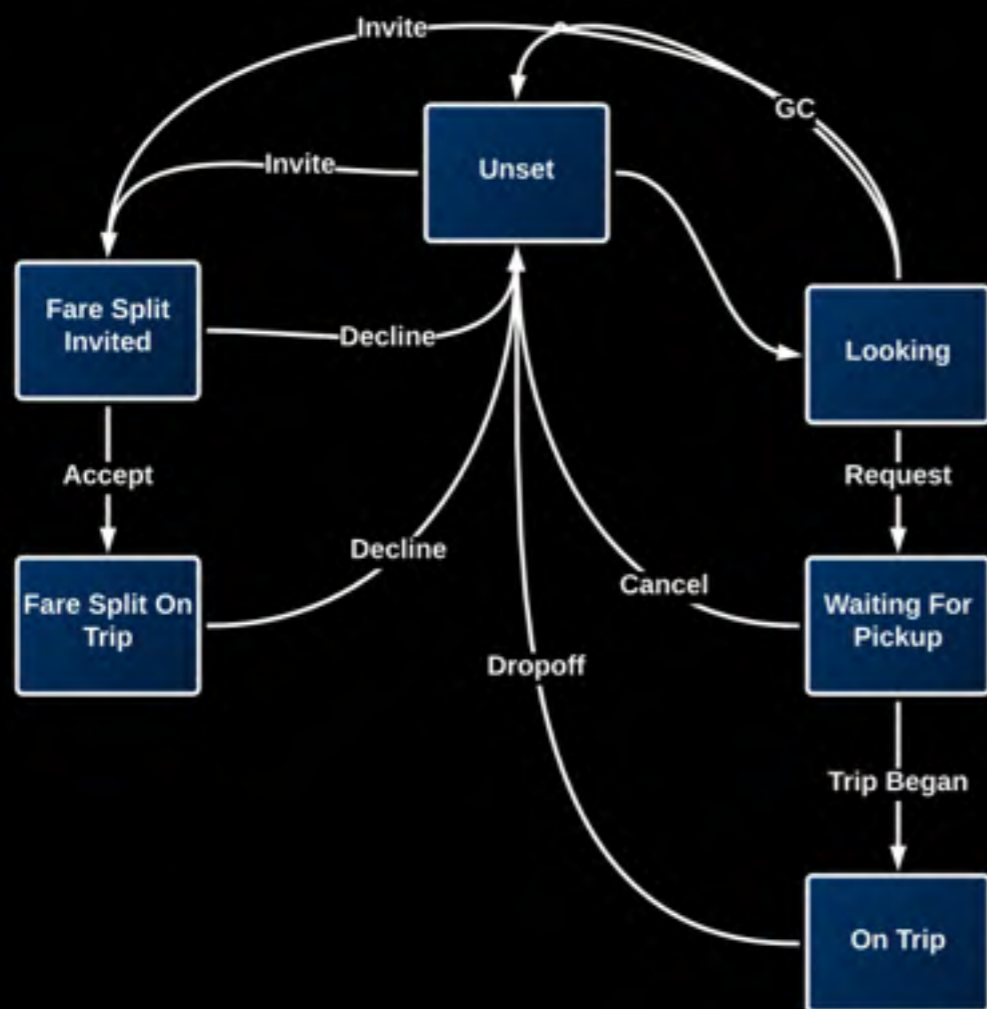
# 优步平台本质：分布式状态机

## 乘客状态

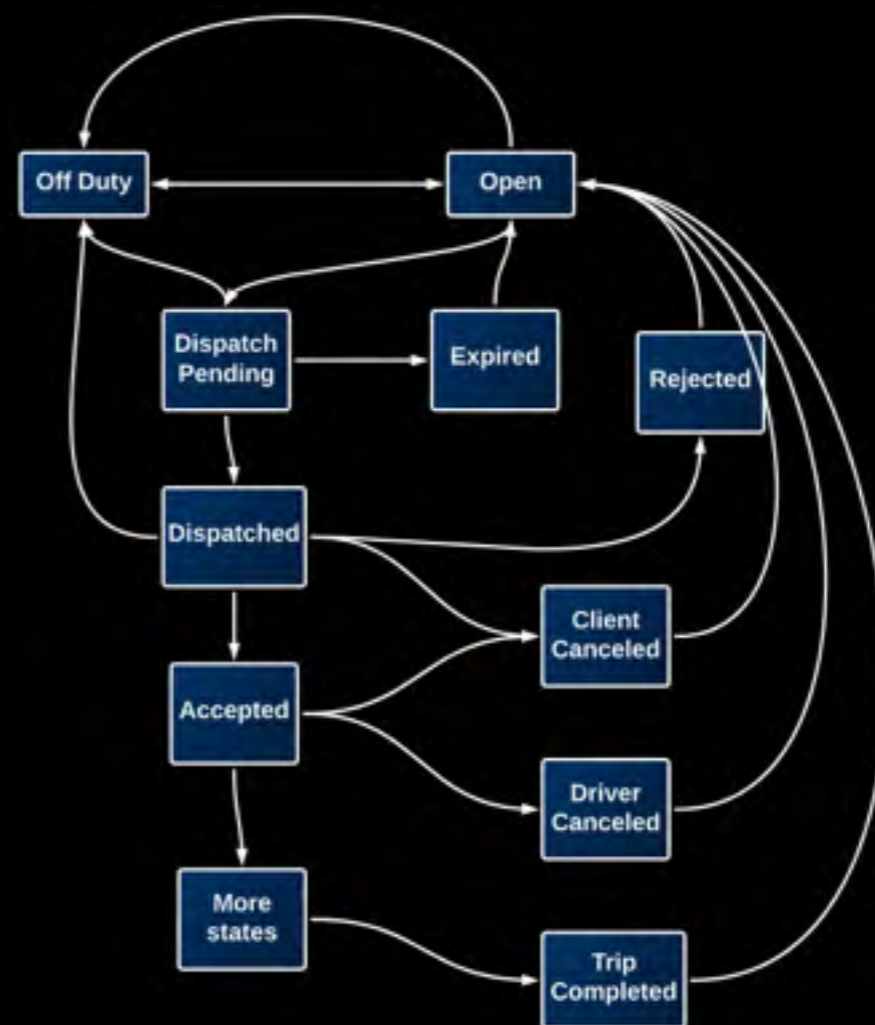


# 优步平台本质：分布式状态机

## 乘客状态



## 司机状态





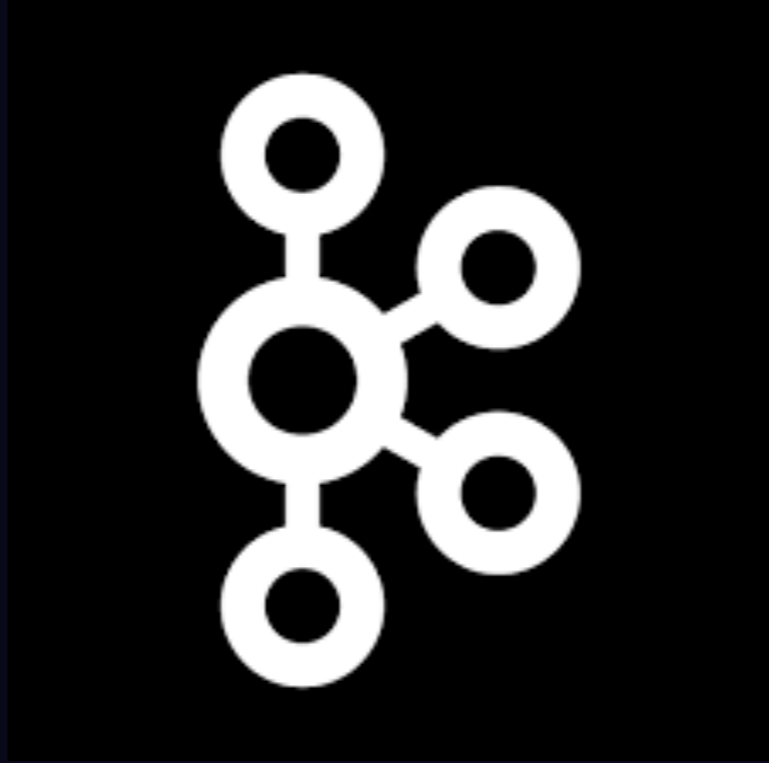
各类应用生成事件

# 秒级事件延迟

事件绝少丢失

收集事件消息 - 低开销、易扩展



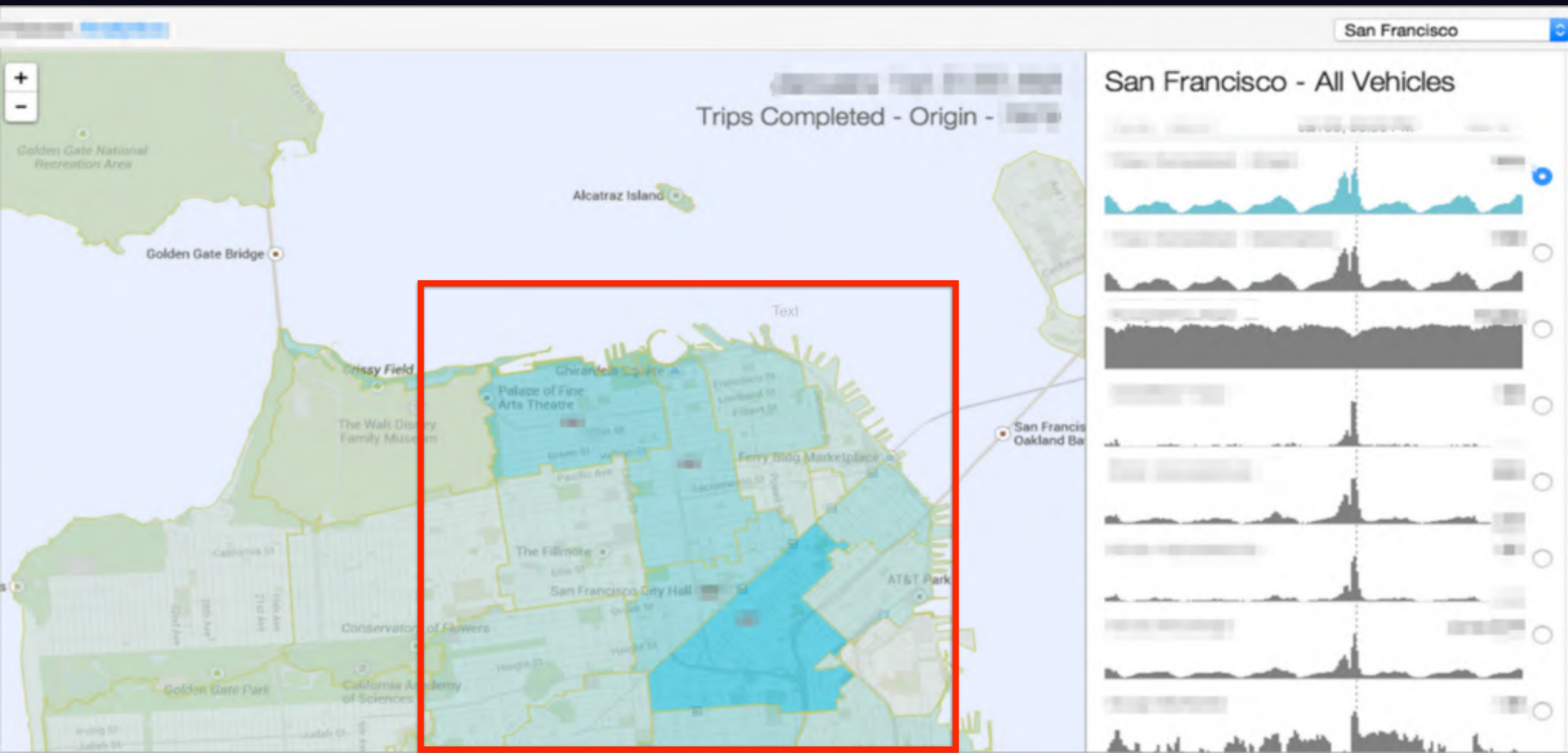


挑战在哪里？

# 多维度

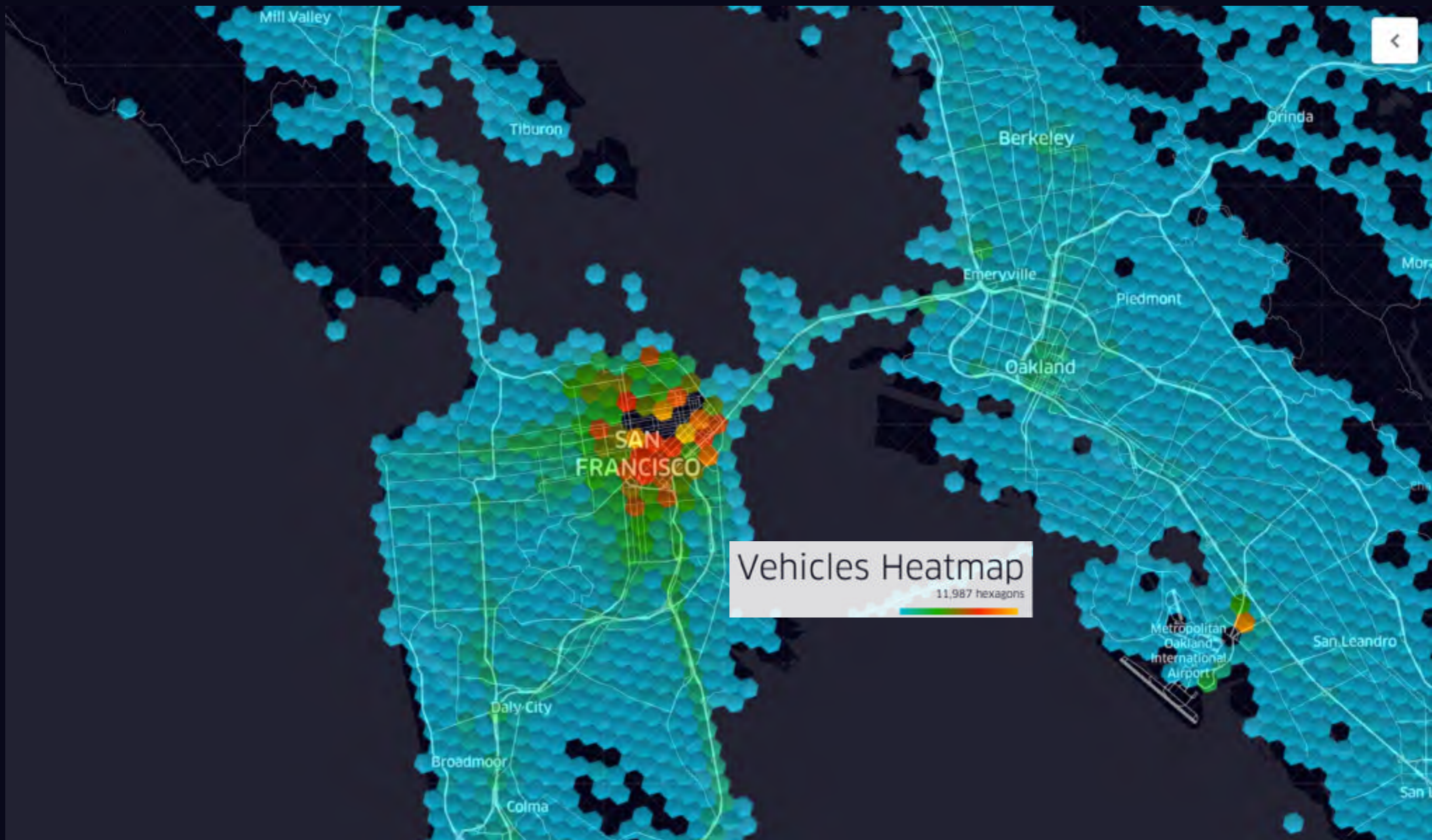
每个事件消息包含数十字段

# 细粒度数据



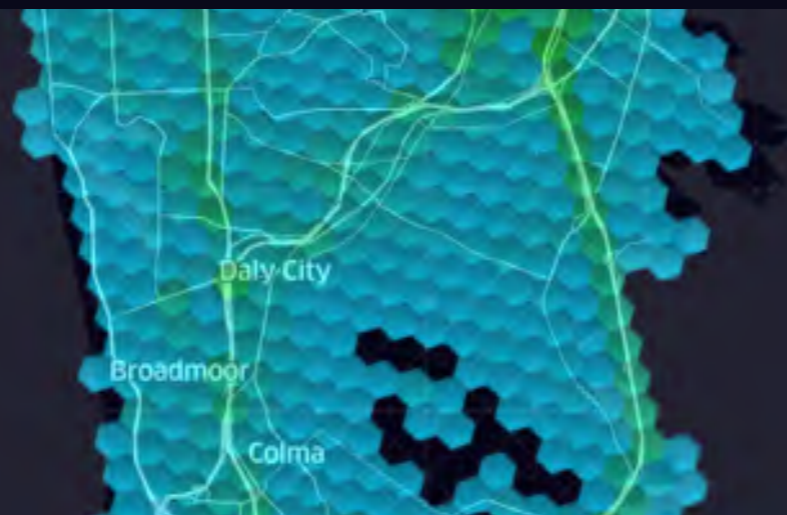


# 细粒度数据



# 细粒度数据

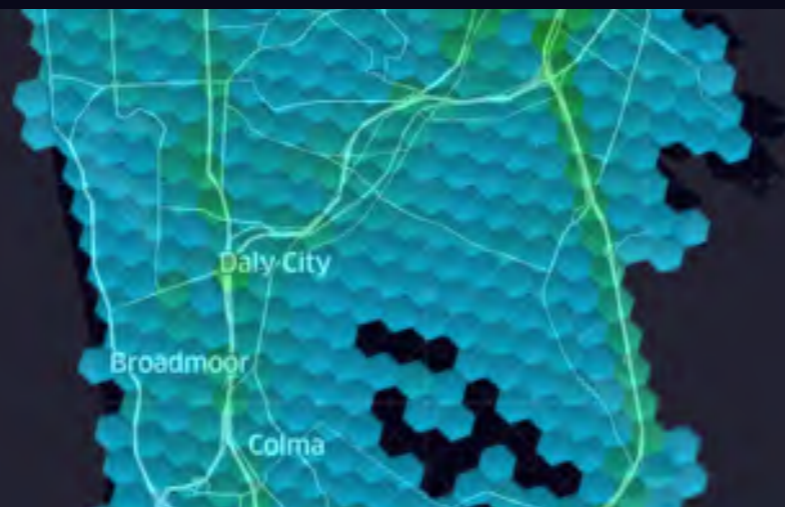
每个城市多于 10,000 六边形





# 细粒度数据

## 7种车型



# 细粒度数据

一天1440分钟





# 细粒度数据

## 13 种司机状态



# 细粒度数据

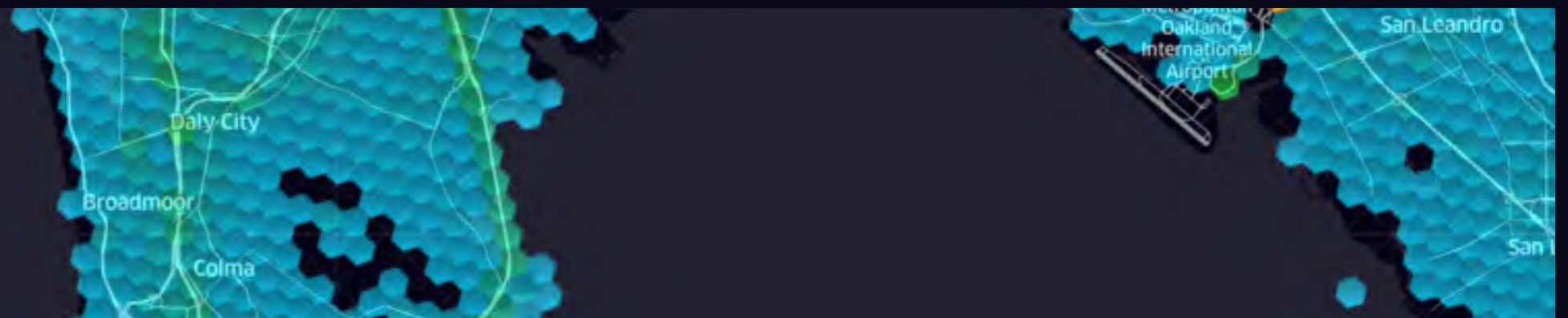
300 座城市



# 细粒度数据

一天数据量:

$300 \times 10,000 \times 7 \times 1440 \times 13 = 3930$ 亿可能组合



查询模式不定

任意维度的组合



# 多种聚合查询

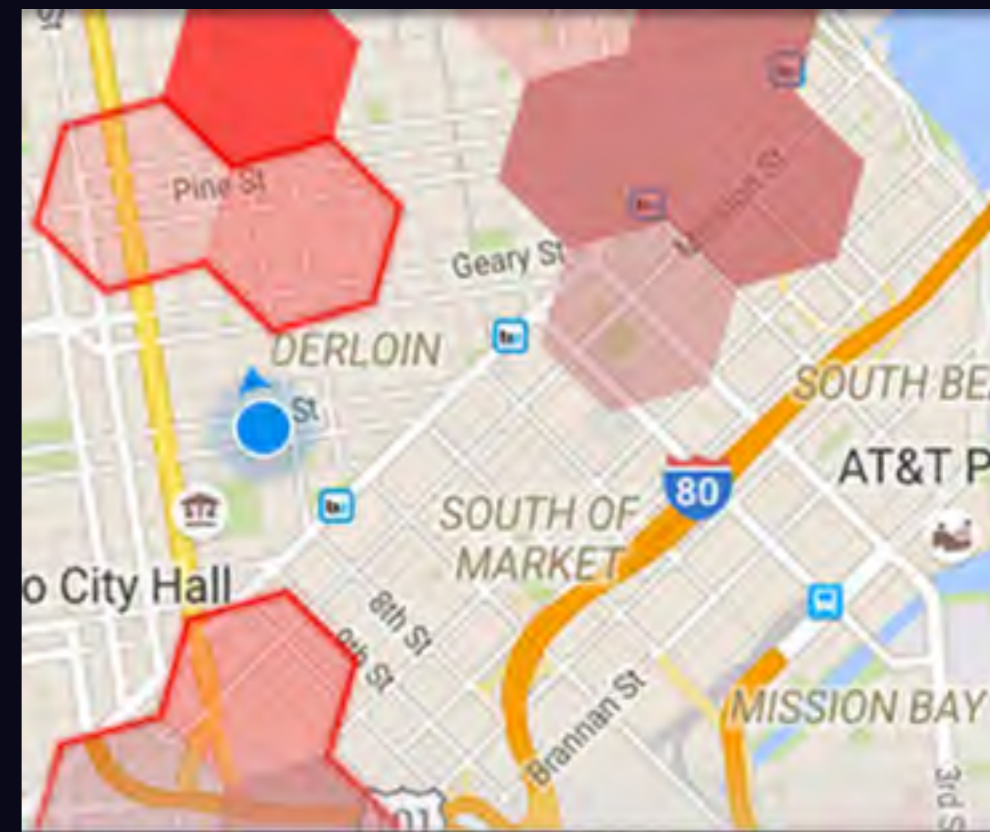
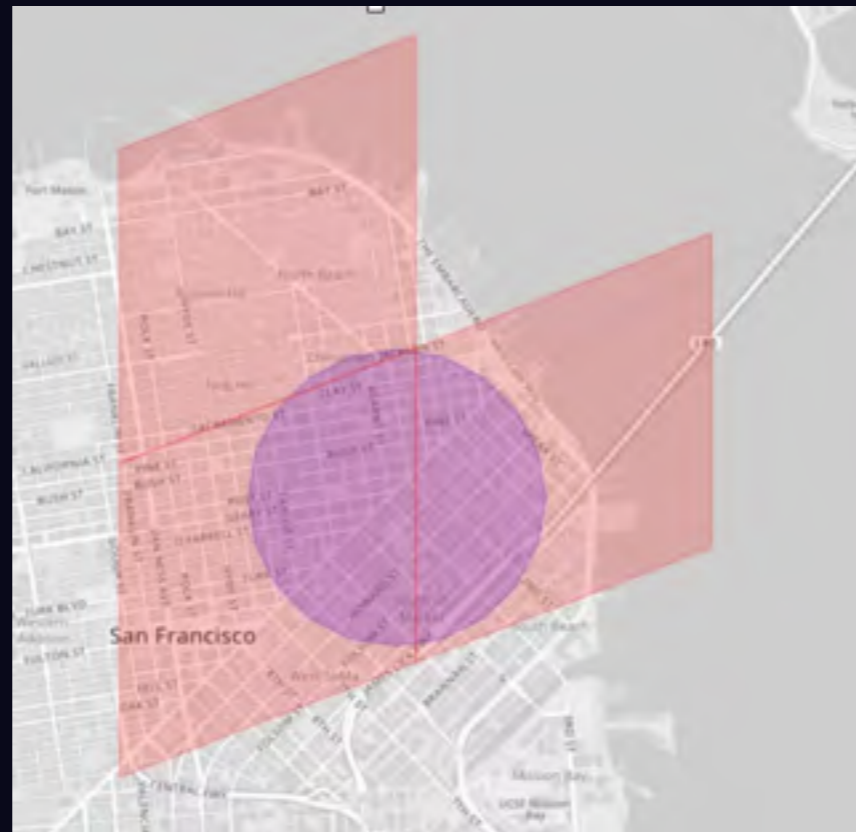
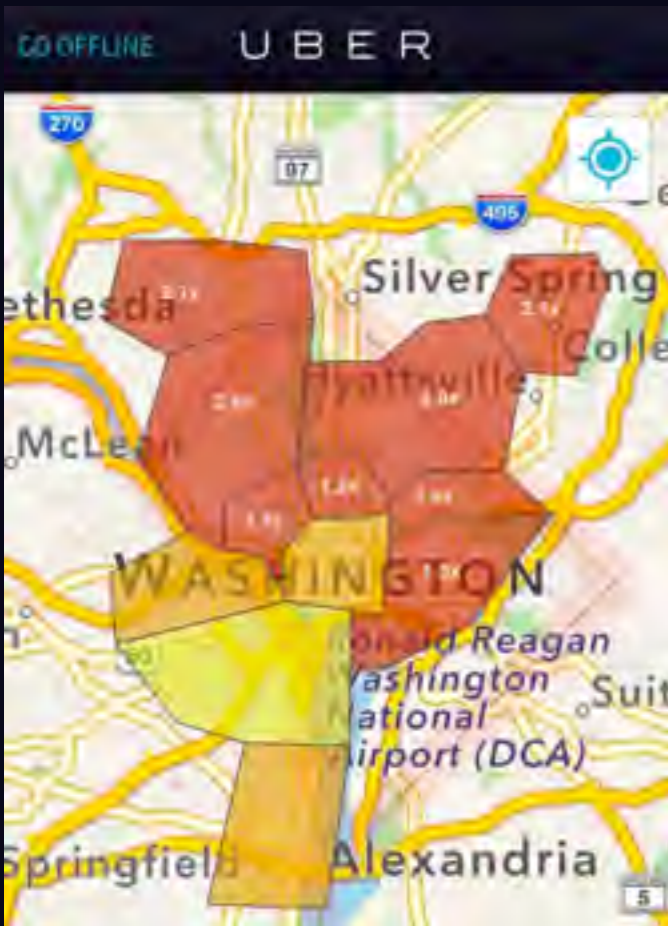
Heatmap

Top N

Histogram

`count()`, `avg()`, `sum()`, `percent()`, `geo`

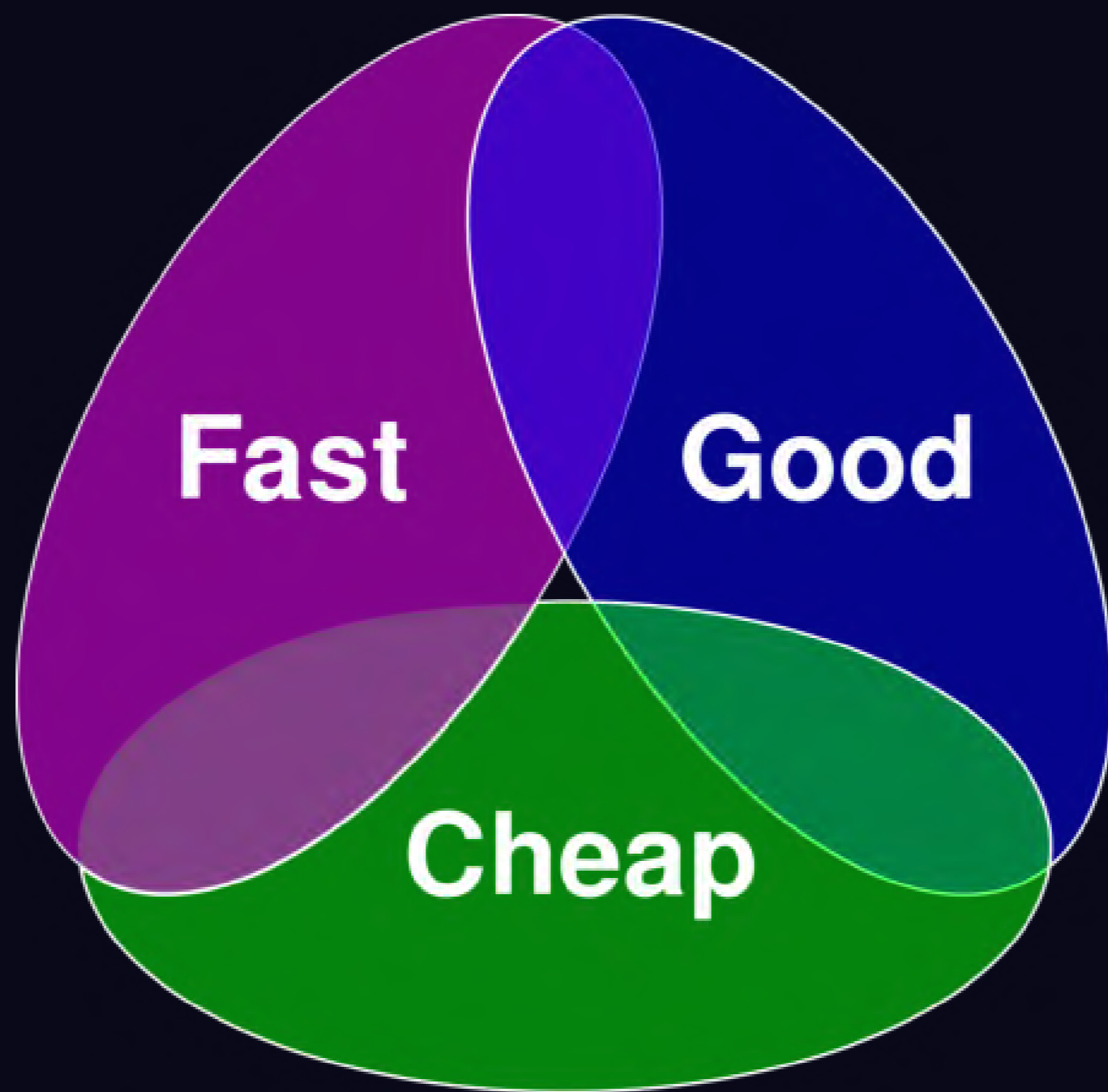
# 多变的地理位置聚合查询



# 高流量

- 数十万条消息每秒，数百亿一天
- 每条消息包含数十字段

短时交货



**Fast**

**Good**

**Cheap**

关键: 把问题一般化



# 数据类型

- 多维时序数据

维度	值
state	driver_arrived
vehicle type	uber X
timestamp	13244323342
latitude	12.23
longitude	30.00

# 数据查询

- 基于单表时空数据的OLAP

```
SELECT <agg functions>, <dimensions>  
FROM <data_source>  
WHERE <boolean filter>  
GROUP BY <dimensions>  
HAVING <boolean filter>  
ORDER BY <sorting criterial>  
LIMIT <n>  
DO <post aggregation>
```

# 选择存储系统

# 最低要求

- 支持时序和地理空间的OLAP
- 支持大流量数据
- 支持秒级查询
- 支持原始数据查询

# 键值数据库



# 键值数据库

# 一键一值：预算所有组合

维度	值
A	a
B	b

# 一键一值：预算所有组合

- 布尔操作符: AND, OR, NOT

维度	值
A	a
B	b

# 一键一值：预算所有组合

维度	值
A	a
B	b

- 布尔操作符: AND, OR, NOT
  - A and (not B)
  - B and (not A)
  - A or B
  - not (A or B)

# 一键一值：预算所有组合

维度
A
B

- {A}
- {B}
- {A, B}
- {}

# 一键一值：预算所有组合

维度
A
B

- {A}
- {B}
- {A, B}
- {}

—> 计算幂集



# 键值数据库的局限

预算所有键值组合的时空复杂度： $O(2^n)$

# 关系数据库

# 关系数据库

---

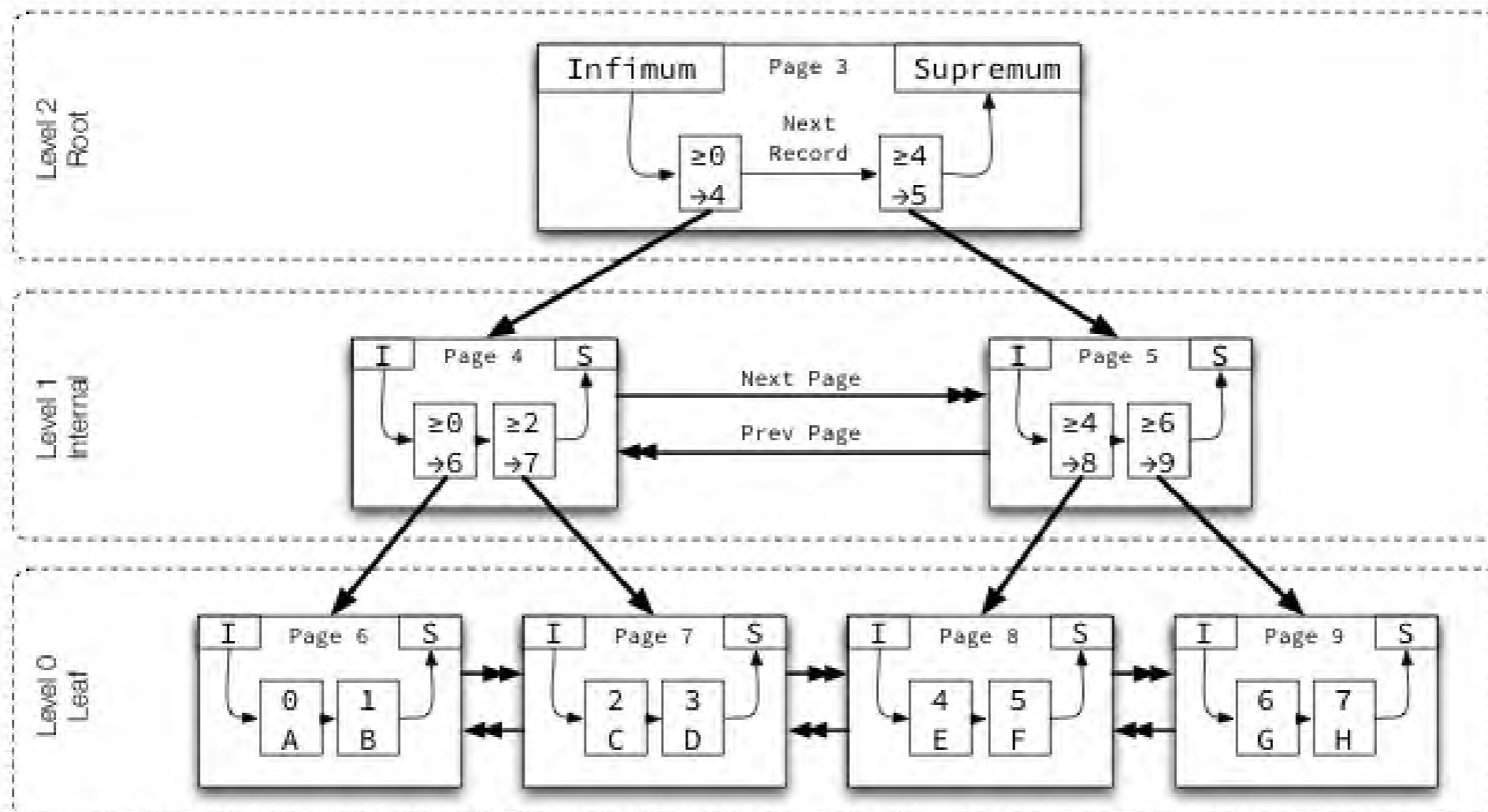
# 关系数据库的局限

- 不易管理多项索引

# 关系数据库的局限

- 扫描速度不够

## B+Tree Structure



Levels are numbered starting from 0 at the leaf pages, incrementing up the tree.

Pages on each level are doubly-linked with previous and next pointers in ascending order by key.

Records within a page are singly-linked with a next pointer in ascending order by key.

Infimum represents a value lower than any key on the page, and is always the first record in the singly-linked list of records.

Supremum represents a value higher than any key on the page, and is always the last record in the singly-linked list of records.

Non-leaf pages contain the minimum key of the child page and the child page number, called a "node pointer".

但我的KV系统速度奇快



# 光快不行

一个城市里每个六边形里车的数目 => 18,000 次查询

平均延迟: 1ms

99.99百分位延迟: 2s

失败率: 0.001%

# 光快不行

一次查询延迟超过99.99百分位的概率:  $(1 - 0.9999^{18000}) \times 100\% = 83\%$

一次查询成功的概率:  $(1 - 0.000001)^{18000} = 84\%$

# 系统必备功能

- 快速扫描
- 布尔查询
- 原始数据
- 各类聚合

Elasticsearch

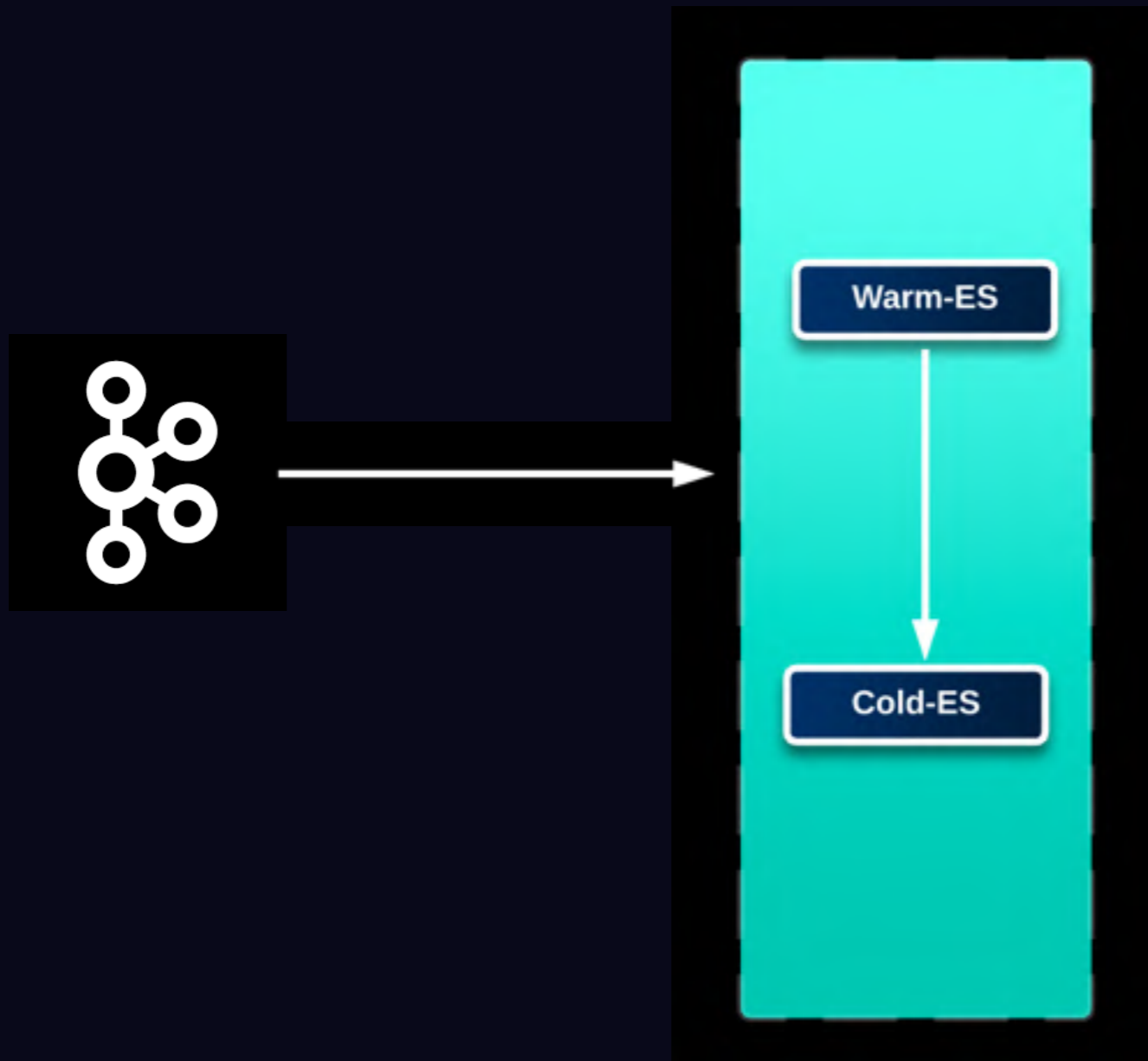
# 基于高效倒排索引的布尔查询

# 内建分布式查询



快速扫描，灵活聚合

# 存储

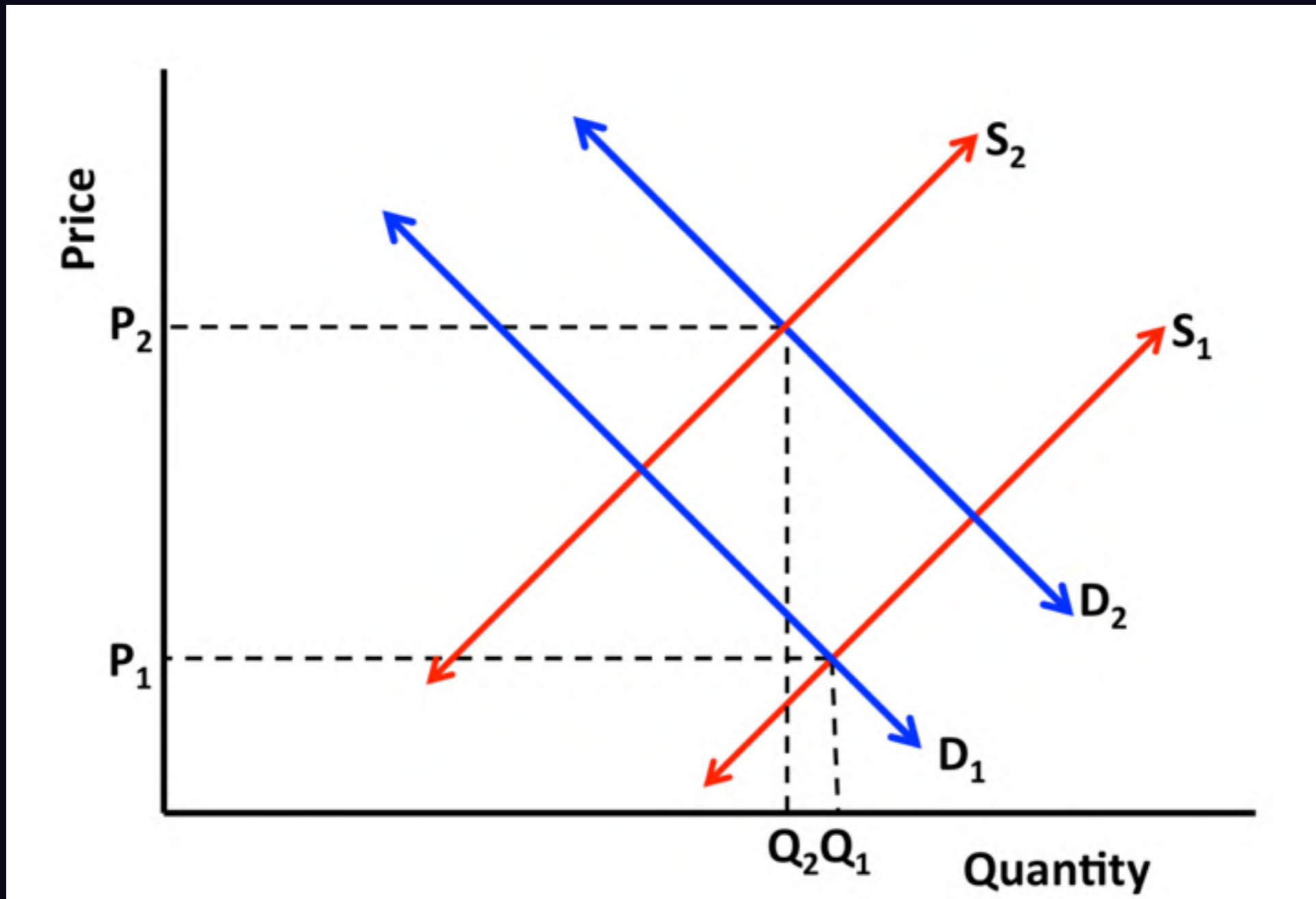


搞定没？

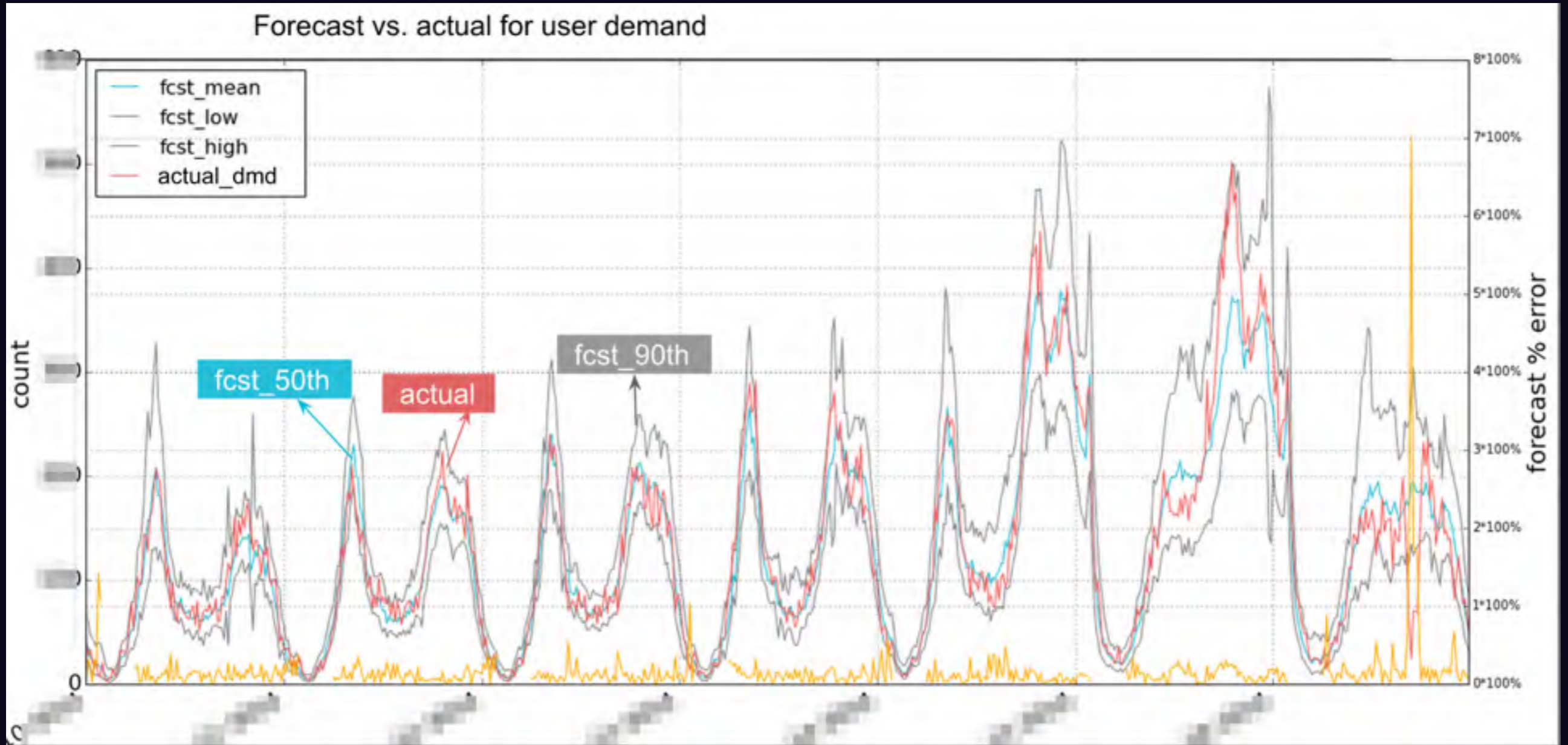
# 数据转换

e.g. (Lat, Long) -> (zipcode, hexagon)

# 动态定价

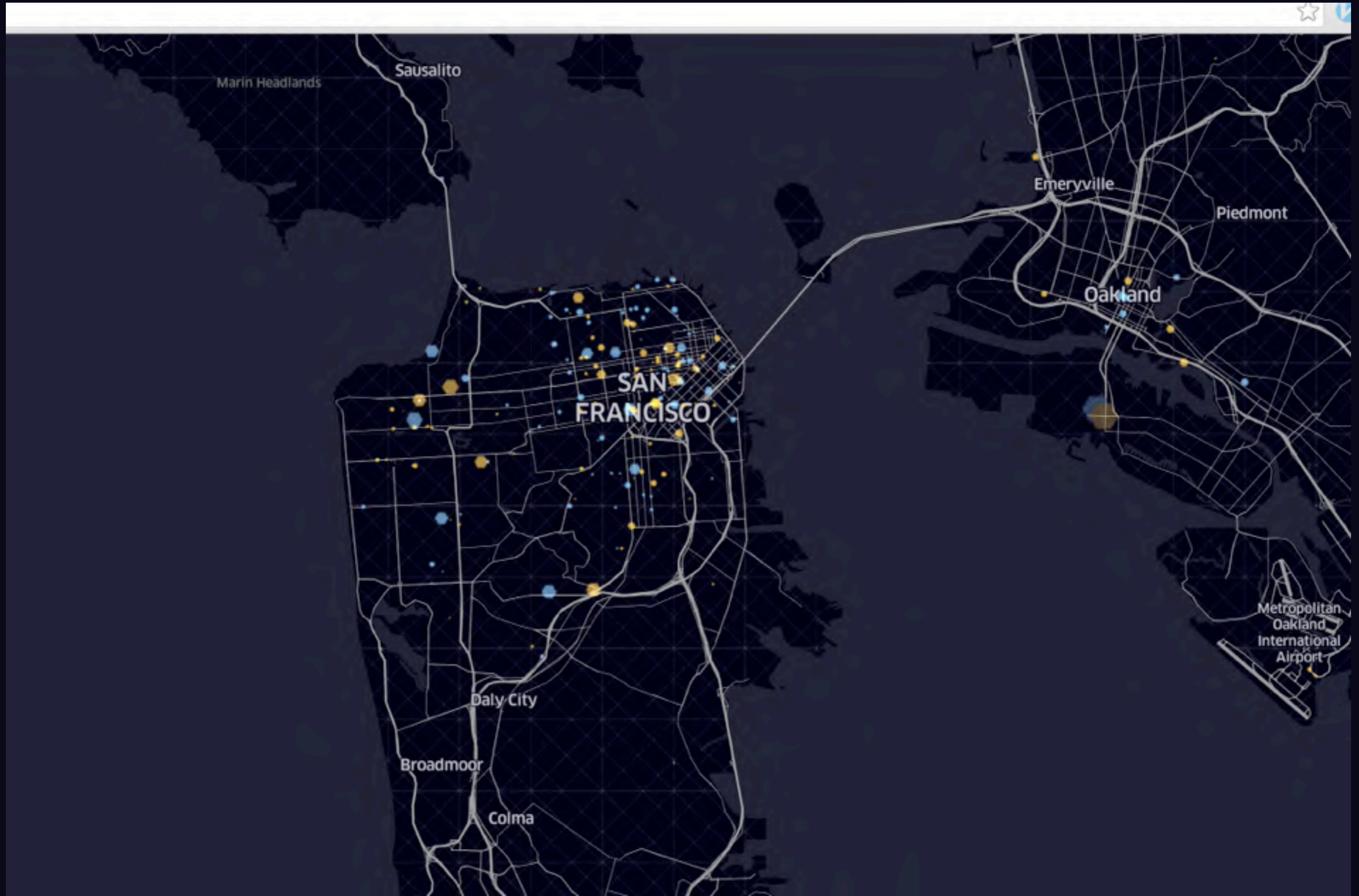


# 趋势预测





# 供求分布







新场景 → 新需求

# 预处理

# Joining Multiple Streams

# Sessionization

# 多级处理

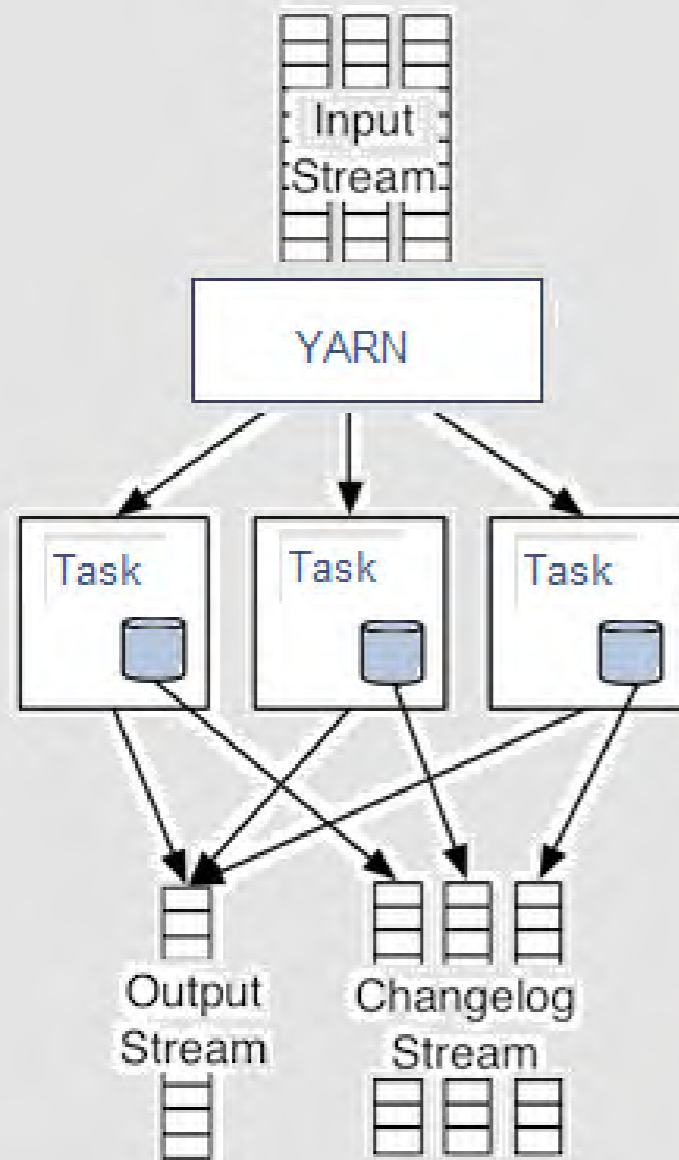
# 状态管理



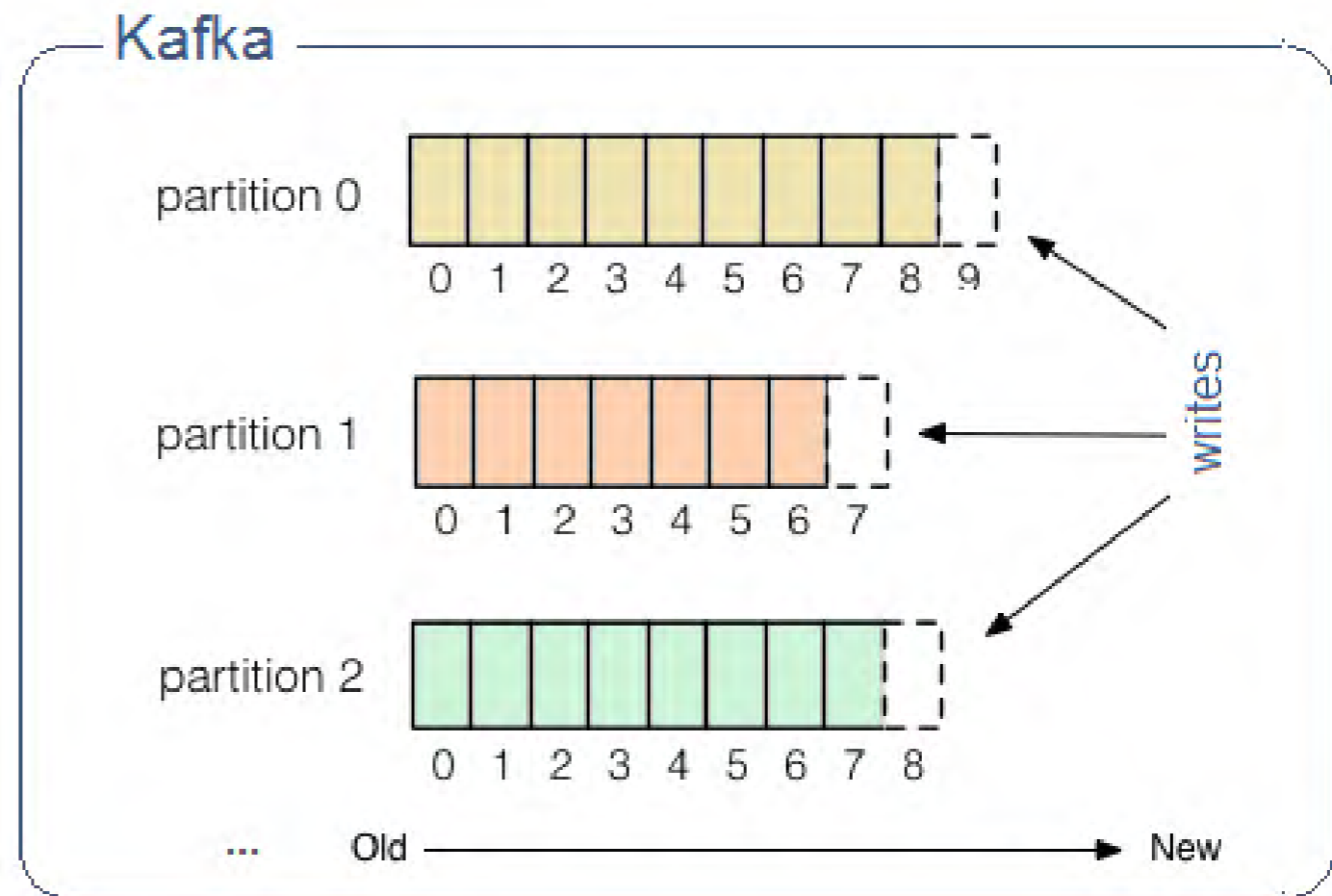
# Apache Samza

# Why Apache Samza?

Samza Job



Partitioned Stream



samza

# DAG on Kafka

# 与Kafka的一流整合

# 内置检查点

# 内置状态管理



# 处理



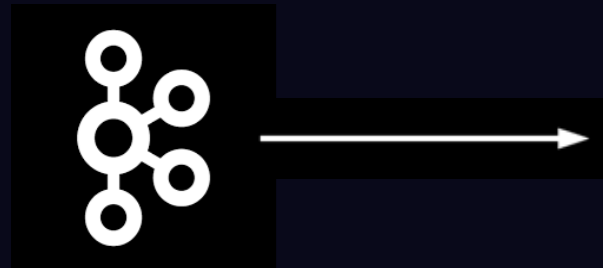
# 存储



存储层当掉怎么办？

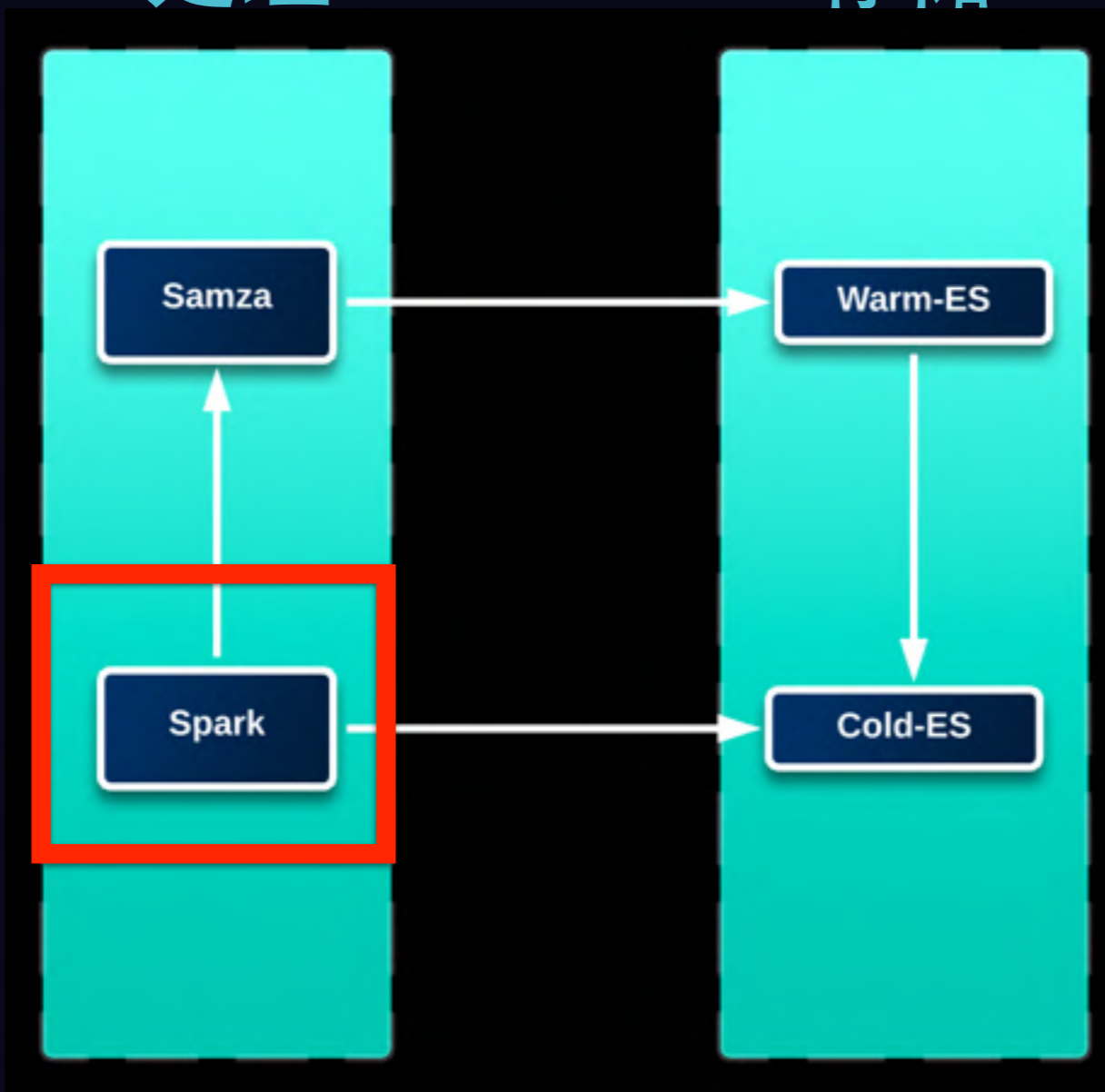


预处理耗时太久怎么办?



处理

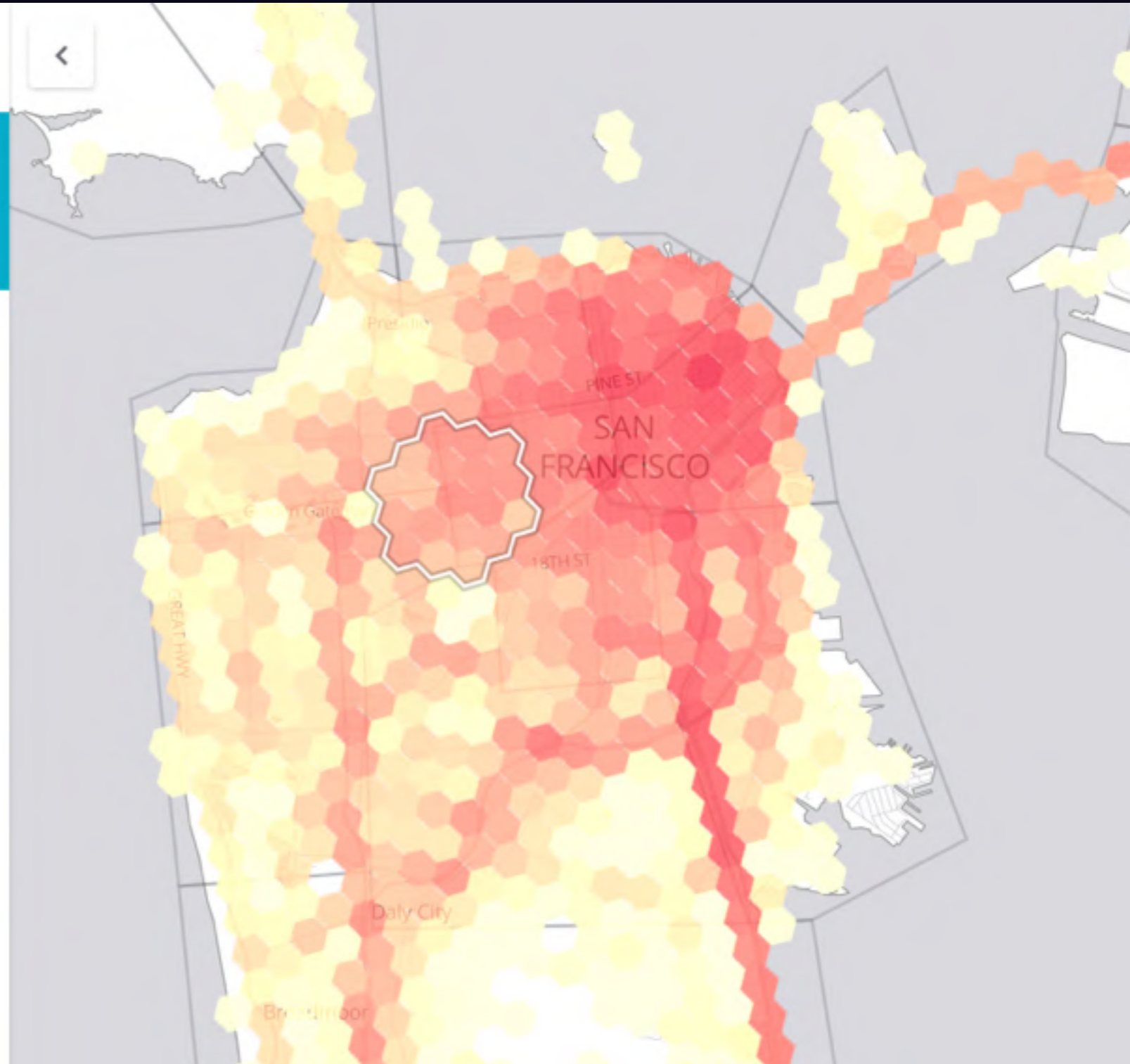
存储



终于搞定了？

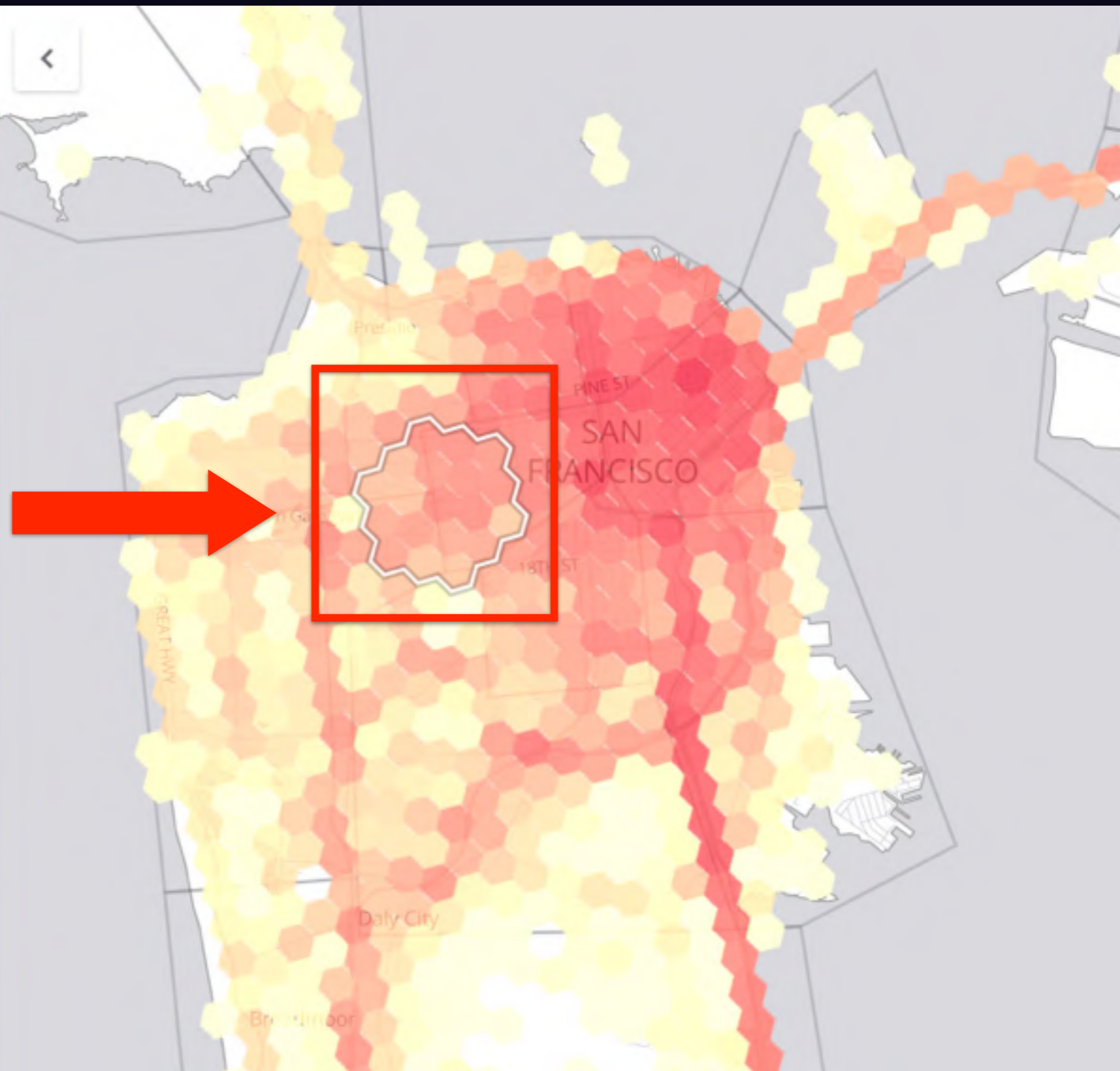
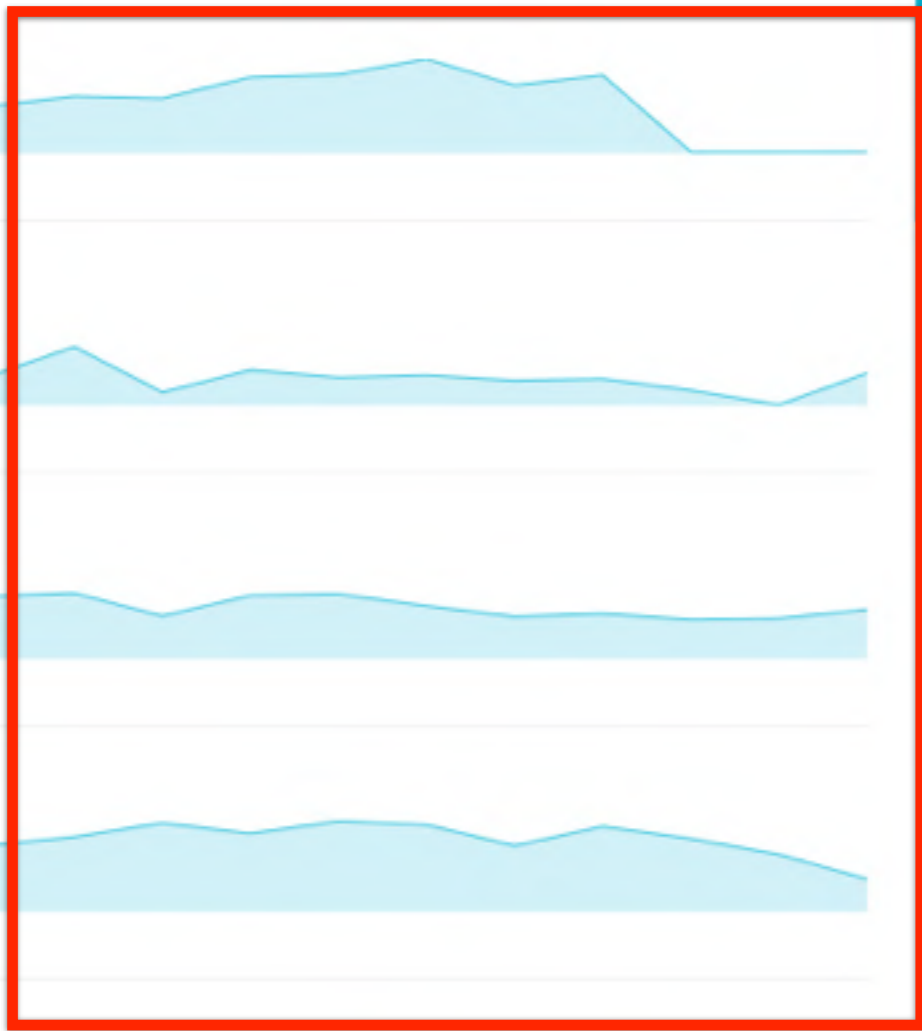
SELECTION RADIUS

0 1 2 3 4 5 6 7



SELECTION RADIUS

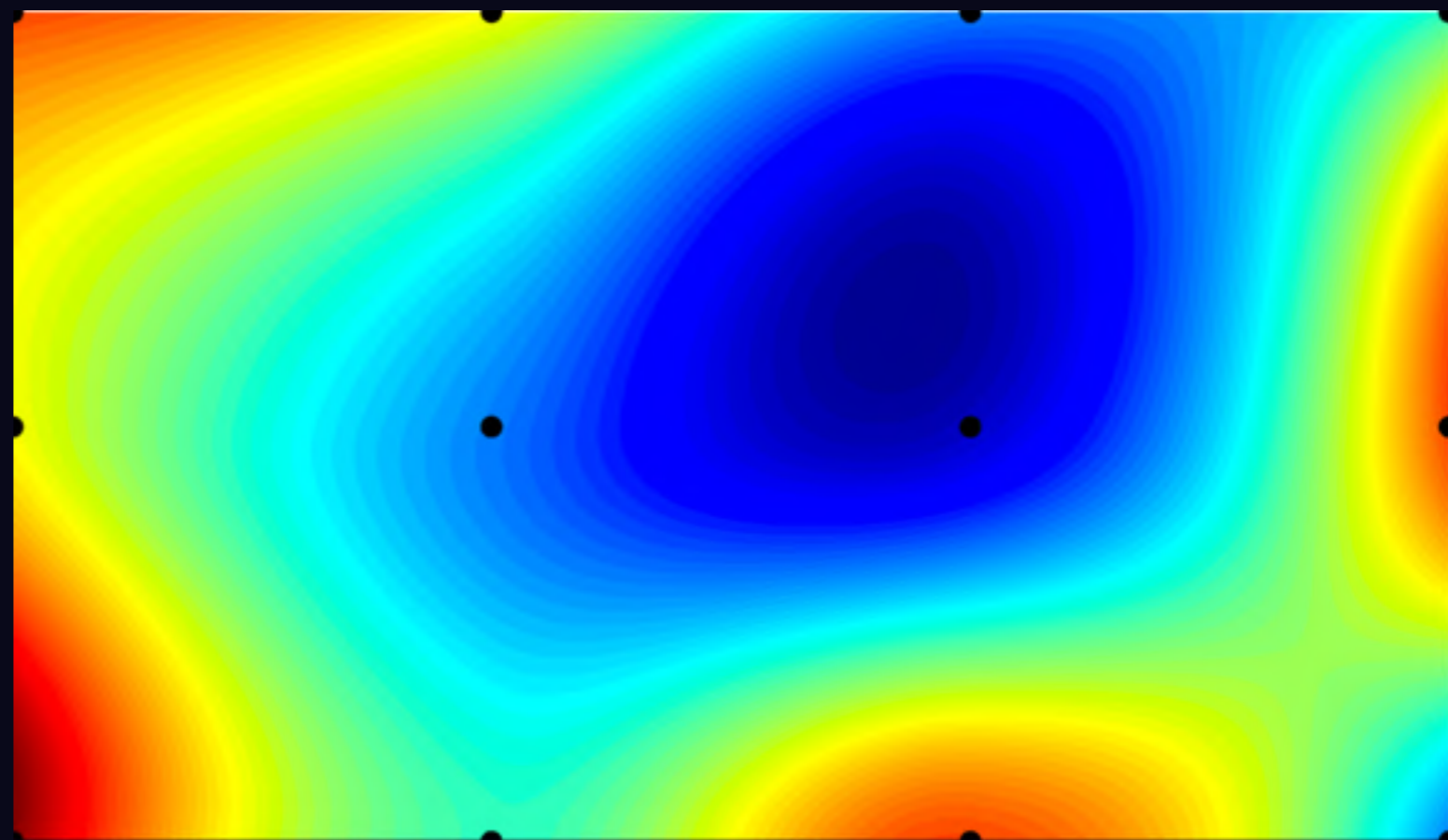
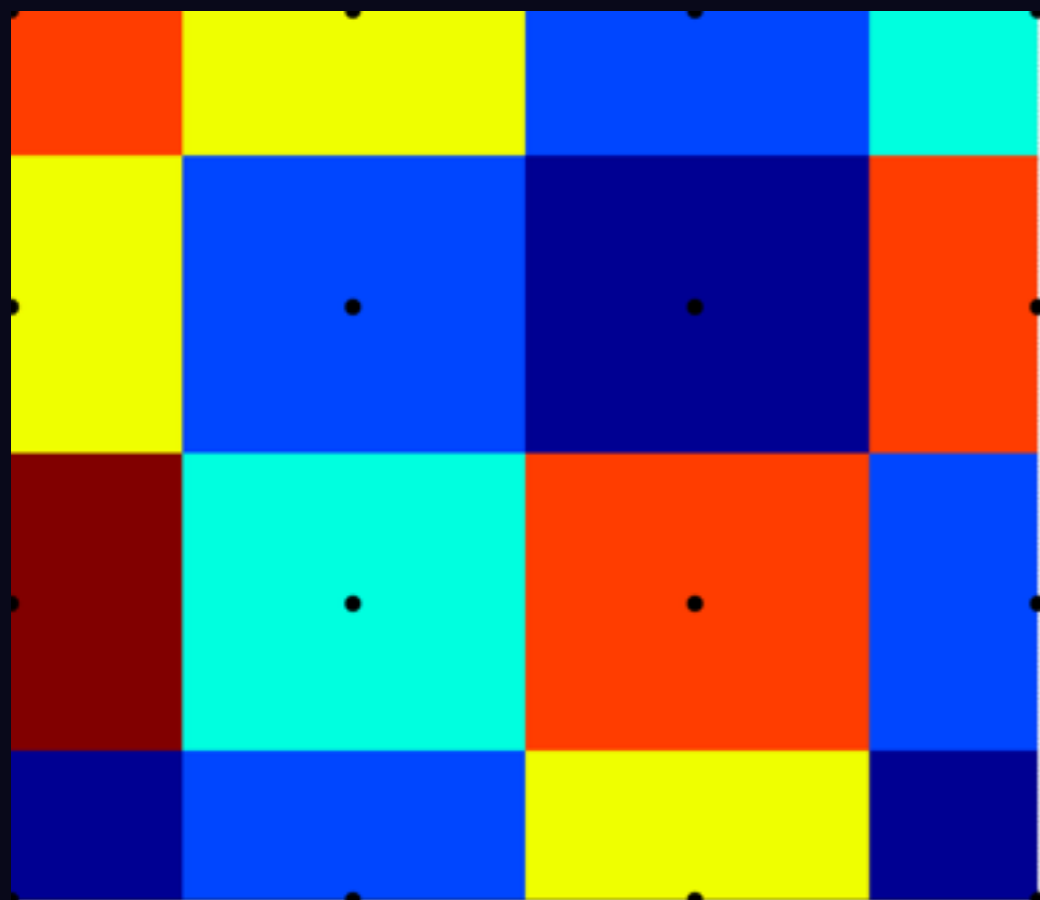
0 1 2 3 4 5 6 7



# 后期处理

# 查询结果转换和平滑处理

# 查询结果转换和平滑处理





# 计算规模

一个城市至少10,000六边形

# 计算规模

每个六边形**331**个邻居需要处理

# 计算规模

一次查询:  $331 \times 10,000 = 310$ 万六边形

# 计算规模

99%-ile 处理时间: 70ms

# 简单架构

# Scalability! But at what COST?

Frank McSherry  
Unaffiliated

Michael Isard  
Unaffiliated\*

Derek G. Murray  
Unaffiliated†

## Abstract

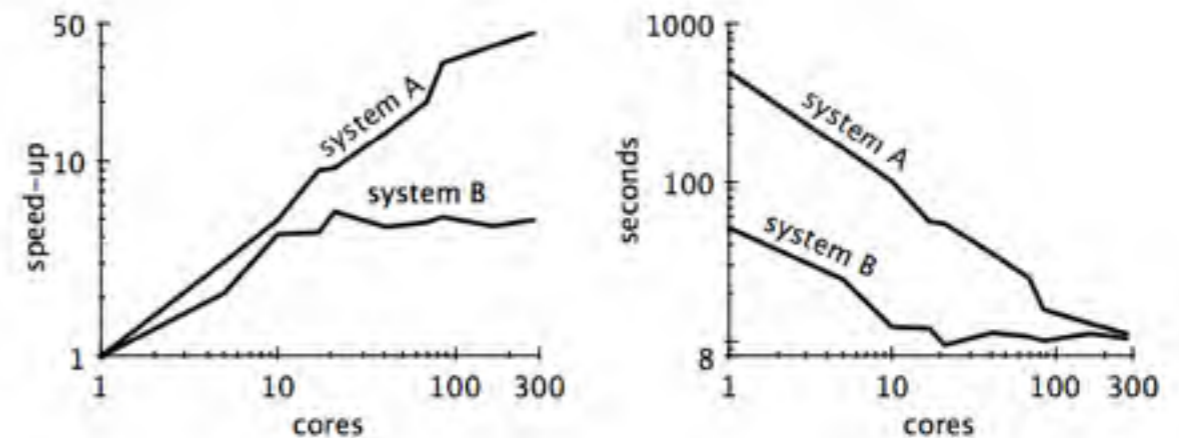
We offer a new metric for big data platforms, COST, or the Configuration that Outperforms a Single Thread. The COST of a given platform for a given problem is the hardware configuration required before the platform outperforms a competent single-threaded implementation. COST weighs a system’s scalability against the overheads introduced by the system, and indicates the actual performance gains of the system, without rewarding systems that bring substantial but parallelizable overheads.

We survey measurements of data-parallel systems recently reported in SOSOP and OSDI, and find that many systems have either a surprisingly large COST, often hundreds of cores, or simply underperform one thread for all of their reported configurations.

## 1 Introduction

“You can have a second computer once you’ve shown you know how to use the first one.”

—Paul Barham



**Figure 1: Scaling and performance measurements for a data-parallel algorithm, before (system A) and after (system B) a simple performance optimization. The unoptimized implementation “scales” far better, despite (or rather, because of) its poor performance.**

While this may appear to be a contrived example, we will argue that many published big data systems more closely resemble system A than they resemble system B.

### 1.1 Methodology

In this paper we take several recent graph processing papers from the systems literature and compare their re-



<b>Connectivity</b>	<b>cores</b>	<b>twitter_rv</b>	<b>uk_2007_05</b>
<b>Spark</b>	128	1784s	8000s+
<b>Giraph</b>	128	200s	8000s+
<b>GraphLab</b>	128	242s	714s
<b>GraphX</b>	128	251s	800s

Connectivity	cores	twitter_rv	uk_2007_05
Spark	128	1784s	8000s+
Giraph	128	200s	8000s+
GraphLab	128	242s	714s
GraphX	128	251s	800s
Laptop	1	153s	417s



Connectivity	cores	twitter_rv	uk_2007_05
Spark	128	1784s	8000s+
Giraph	128	200s	8000s+
GraphLab	128	242s	714s
GraphX	128	251s	800s
Laptop	1	<del>153s</del> 15s	<del>417s</del> 30s

**“You can have a second computer once you’ve shown you know how to use the first one.”**

**- Paul Graham**

# 后期处理

- 每个处理单元都是纯函数
- 通过组合算子组合处理单元

# 后期处理

- 高度并行化的执行
- 流水线

# 务实考量

# 数据发现

Elasticsearch 查询语句过度复杂

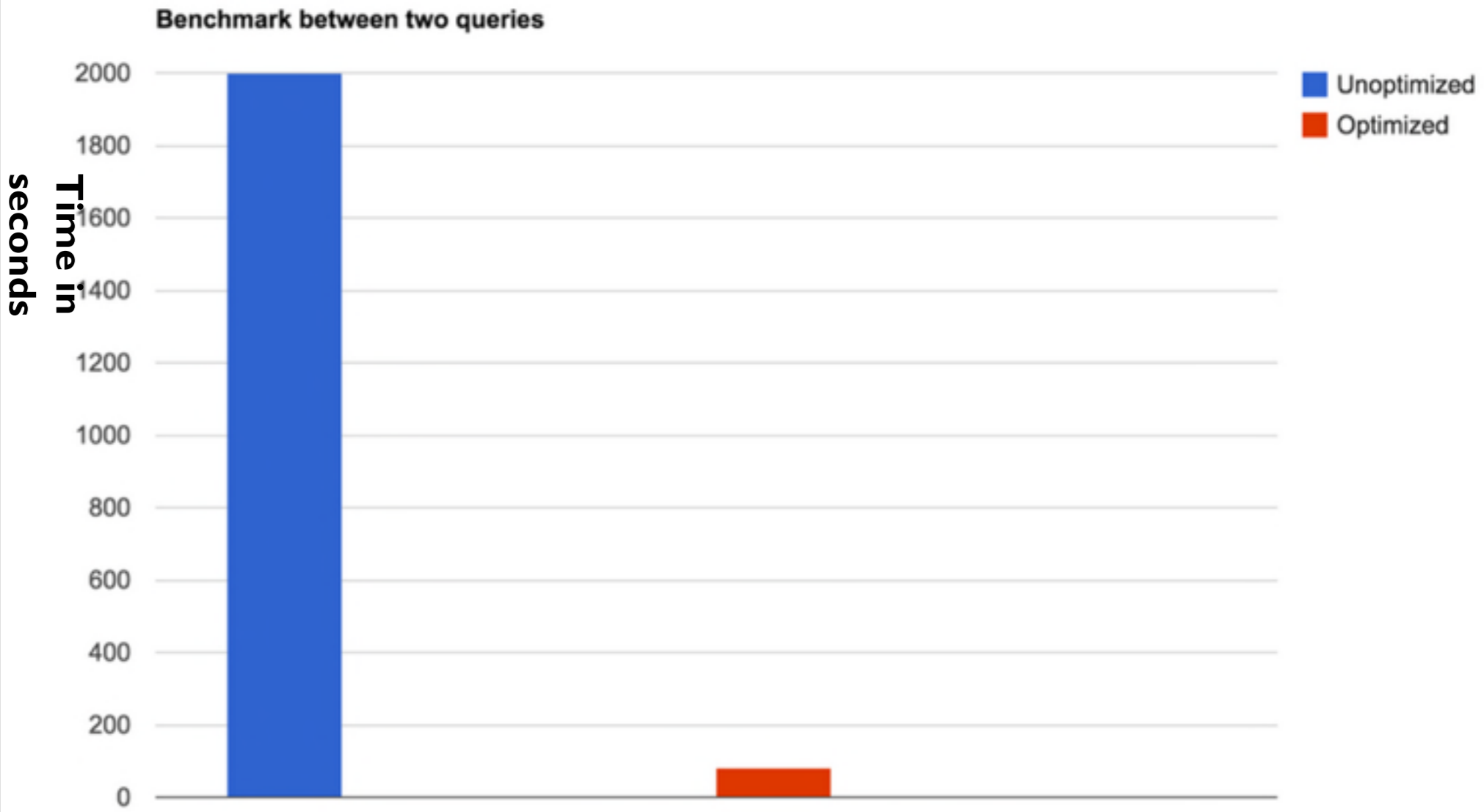
```
SELECT timeseries(7d)
FROM driverAcceptanceRate
WHERE geo_dist(10, [37, 22])
AND time IN (2015-02-04, 2015-03-06)
AND msg.driverId = 1
```

```
}
  }
    }
```



# Elasticsearch 查询可以优化

- 流水线
- 查询验证
- 查询限速

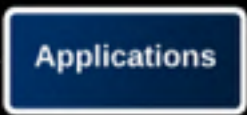
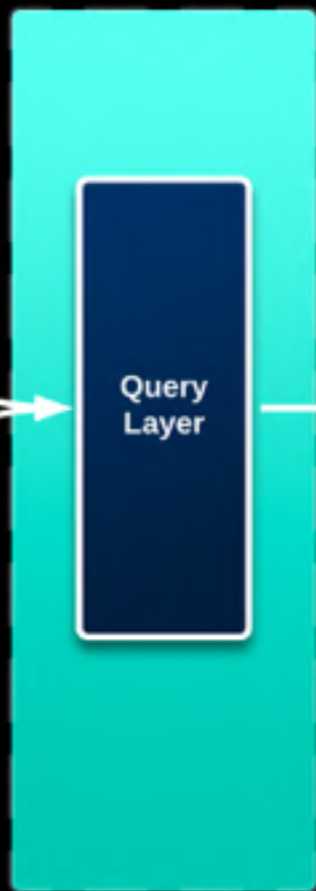
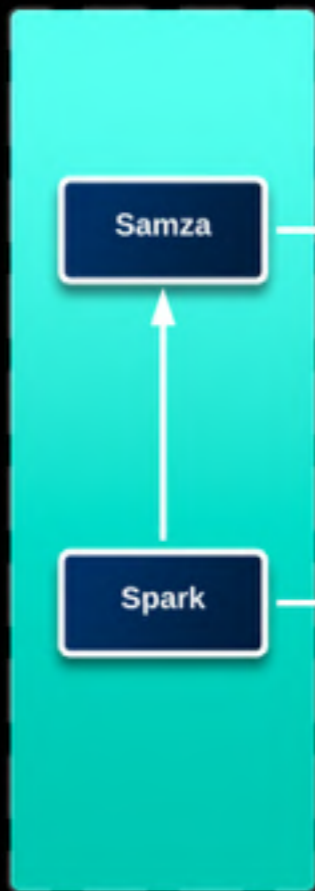


Elasticsearch 也许会被替换

# Processing

# Storage

# Query



还有一件事

数据流里总有不同模式

总有快速探索发现模式的需要

多少司机在5分钟内连续取消请求5次以上



哪些乘客半小时内相距超过100公里的地方叫车？

# Complex Event Processing

```
FROM driver_canceled#window.time(10 min)
SELECT clientUUID, count(clientUUID) as cancelCount
GROUP BY clientUUID HAVING cancelCount > 10
INSERT INTO hipchat(room);
```



简单实现



谢谢!

U B E R

