# Scaling for the Known Unknown

Suhail Patel

# March 2016

**1,861**
Investors
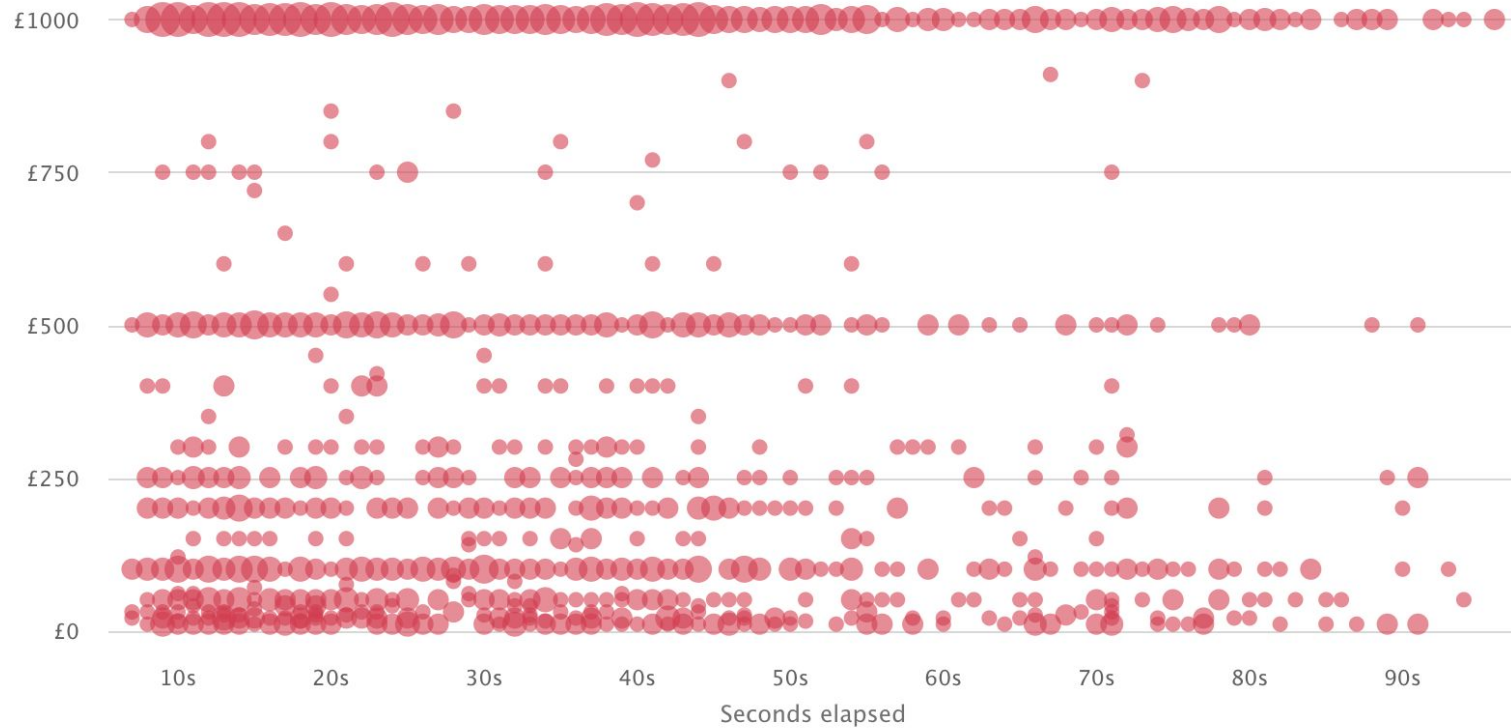
**£1,000,000**
Raised

**96**
Seconds

# March 2016

# February 2017
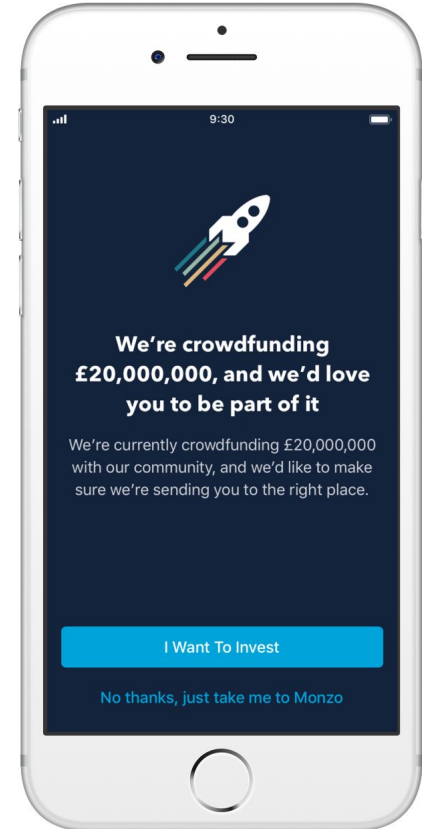
**41,267**

**Pledges to invest**

**£2,500,000**

**Raised**

# Late 2018

Monzo is raising £20,000,000 and all our customers will be eligible to participate

👋

Hi, i'm Suhail

I'm a Platform Engineer at **Monzo**. I work on the
Infrastructure and Reliability squad. We help build
the base so other engineers can ship their
services and applications.

- Email:   hi@suhailpatel.com
- Twitter:  @suhailpatel

monzo

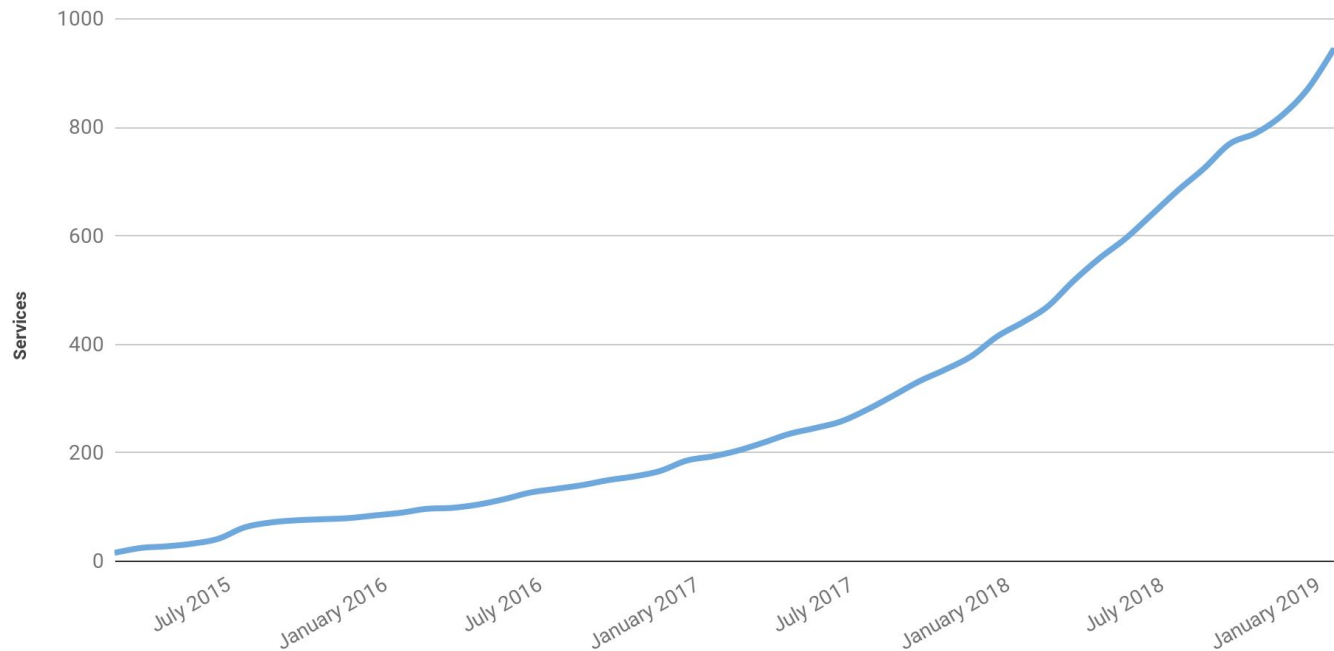5354 2000 0000 0000

EXP END 08/1

469-624-882

Introduction

# A brief overview of our Platform

Building a Crowdfunding Backend

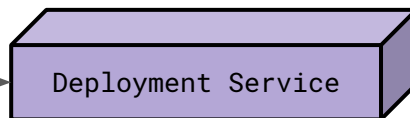Load testing + Finding bottlenecks
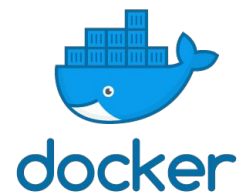
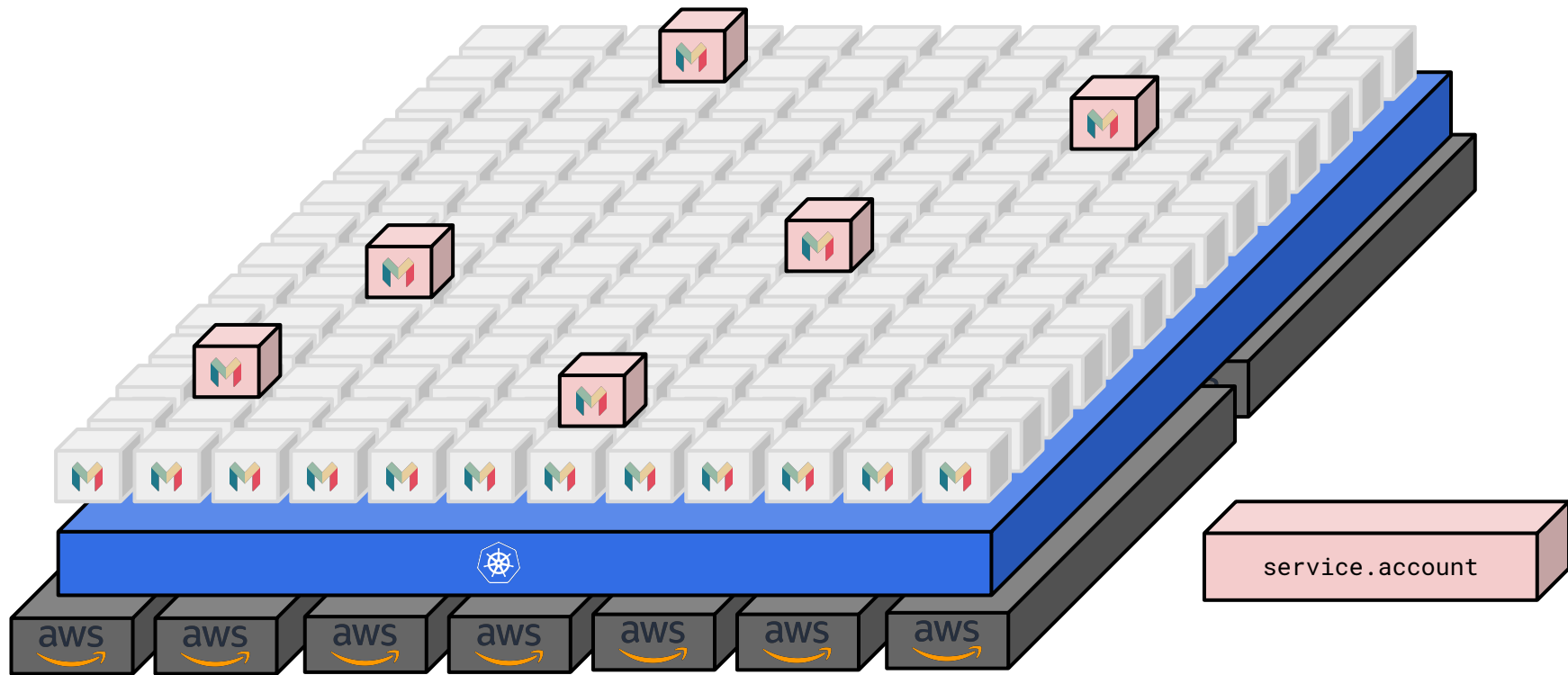# Number of services

Deployment Service

Please deploy
service.account at
revision b32a9e64
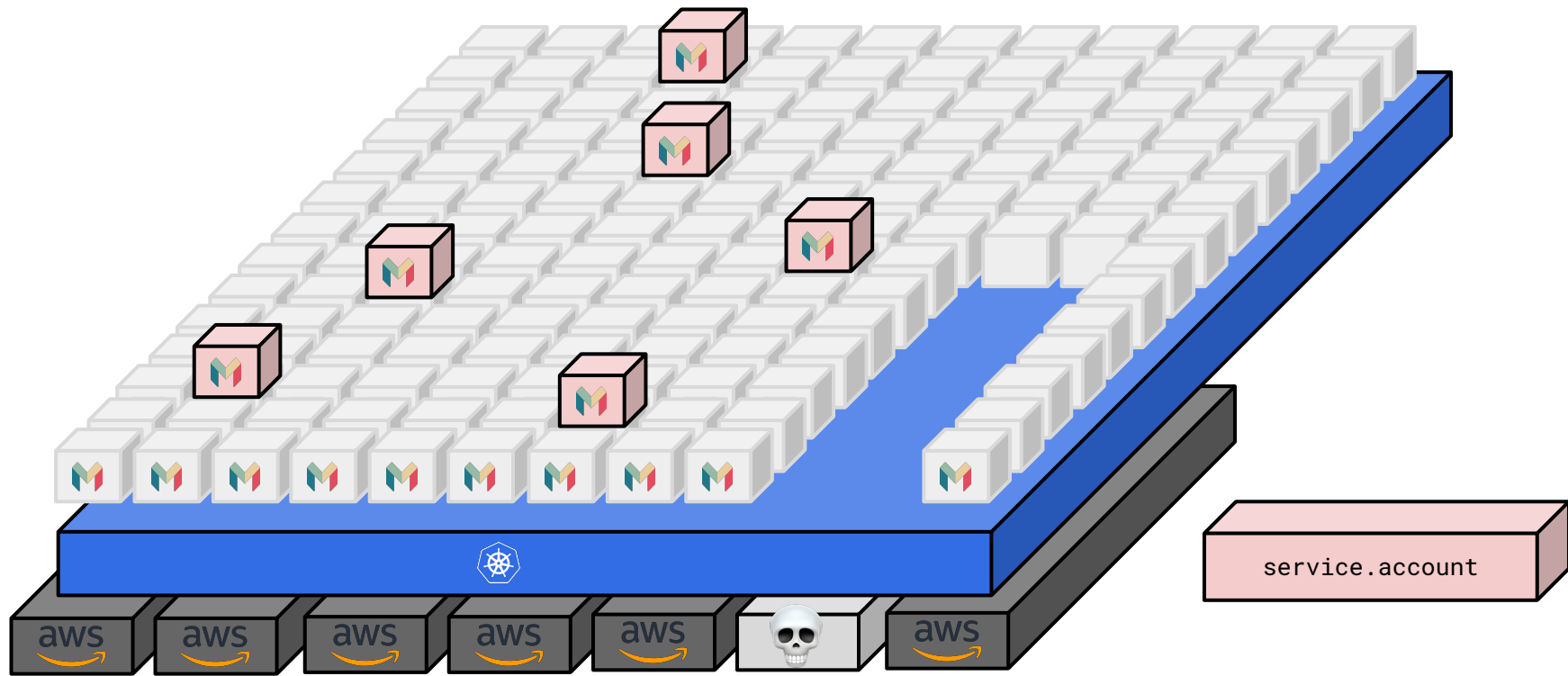
Review checks
Static analysis
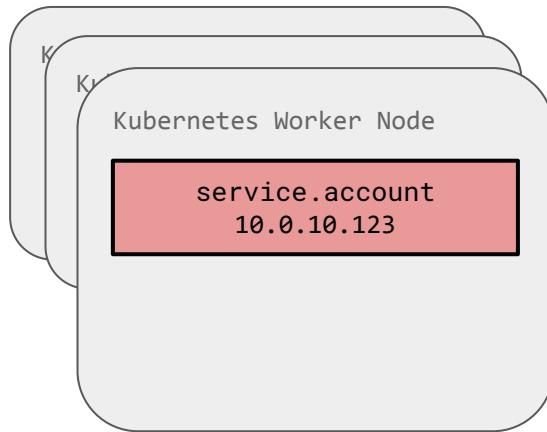Build checks

# Running services

# Running services

What we want from services:
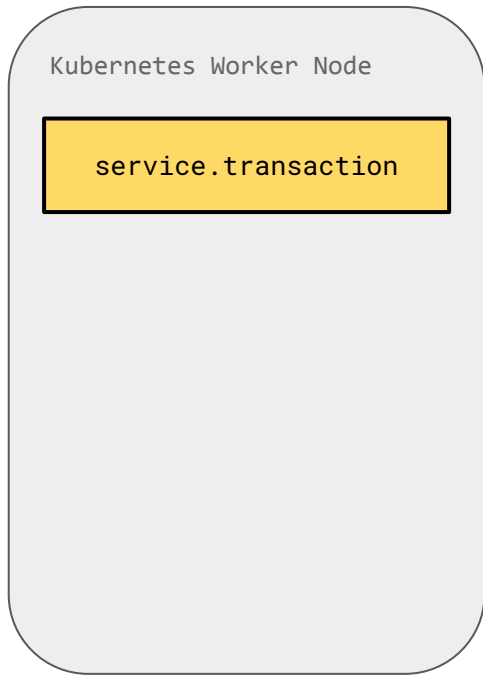- Self-contained
- Scalable
- Stateless
- Fault tolerance

# Running services

# Running services

# Running services



Kubernetes Worker Node

service.transaction

Host: service.account
Proxy: 127.0.0.1:4140
HTTP GET /account

Service Mesh

Kubernetes Worker Node

service.account
10.0.10.123

Service Mesh

Route request to a
service.account replica, let's
try the one at 10.0.10.123

# Service Mesh

The Service Mesh ties the microservices together.
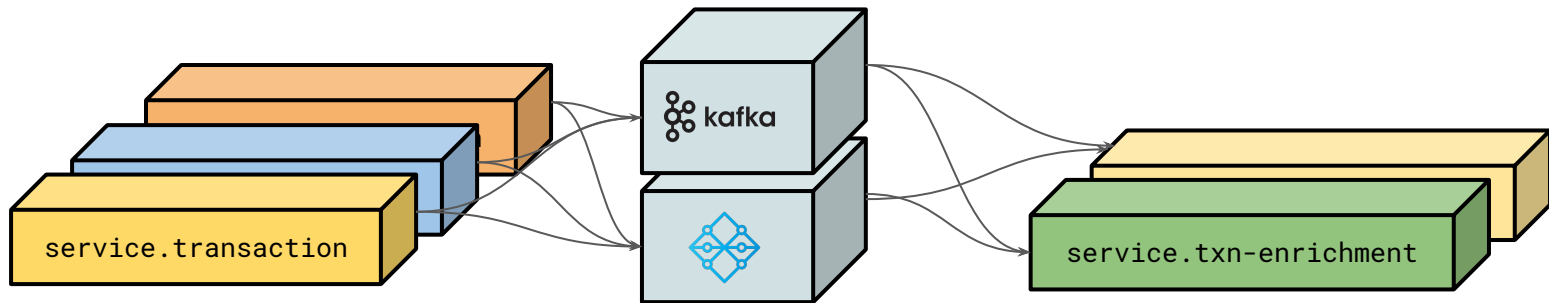It acts as the RPC proxy.

- Handles service discovery and routing
- Retries / Timeouts / Circuit Breaking
- Observability

# Asynchronous messaging

Many things can occur asynchronously rather than a direct blocking RPC.
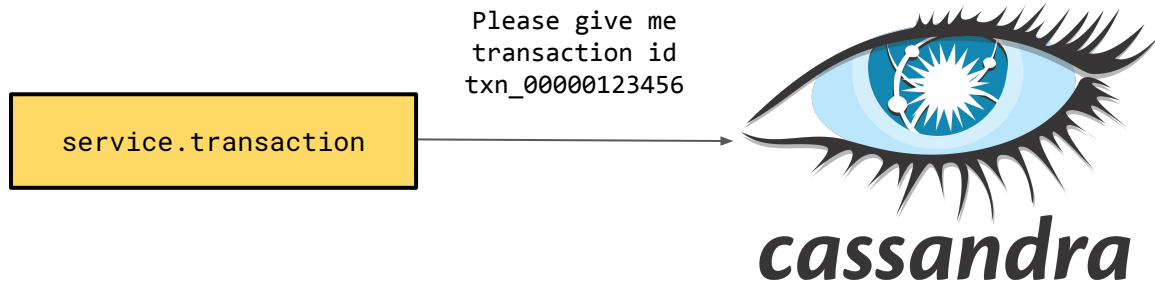
Message queues like NSQ and Kafka provide asynchronous flows with at least once message delivery semantics.

# Asynchronous messaging

# Storing data with Cassandra



service.transaction

Please give me
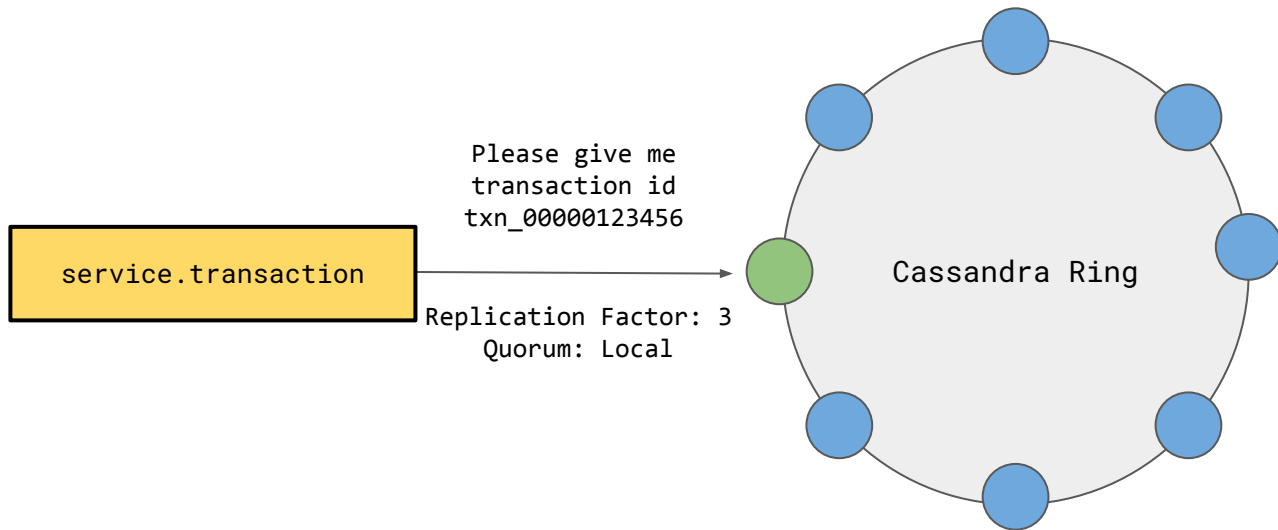transaction id
txn_00000123456

cassandra
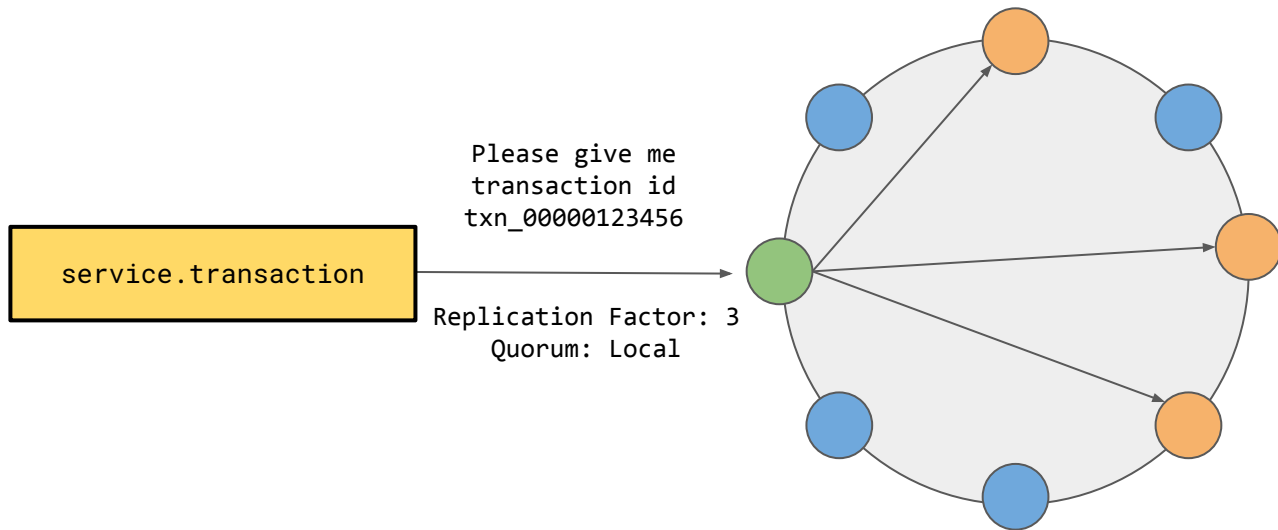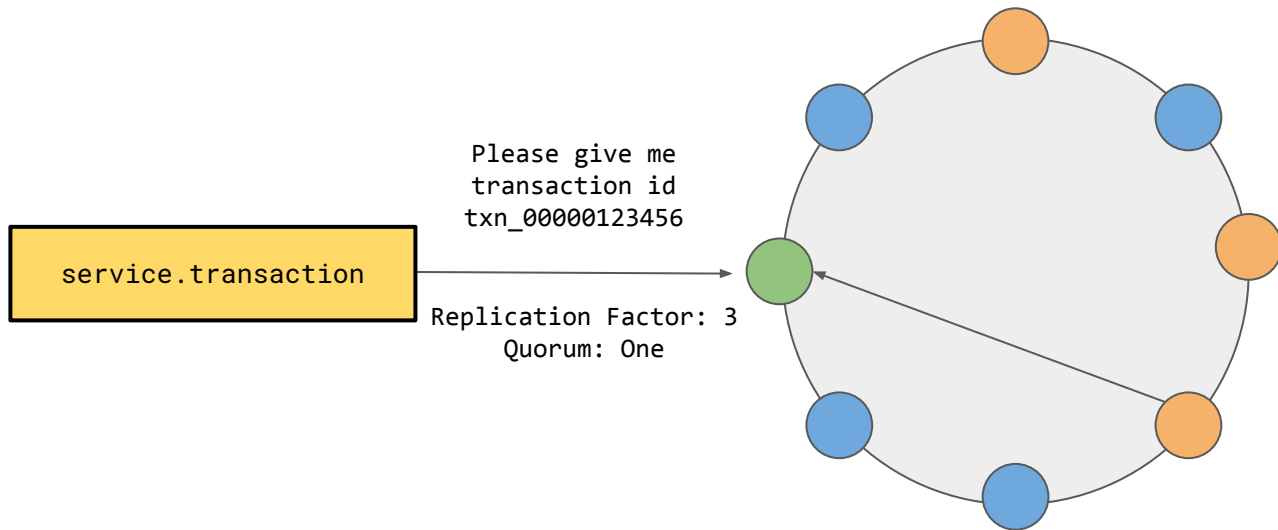
# Storing data with Cassandra

# Storing data with Cassandra

# Storing data with Cassandra

# Storing data with Cassandra

service.transaction

Please give me
transaction id
txn_00000123456

Replication Factor: 3
Quorum: Local

# Distributed Locking with etcd

service.transaction

Please can I get a
lock on transaction
txn_00000123456
so I have sole access

# Distributed Locking with etcd

# Monitoring with Prometheus

Prometheus is a flexible time-series data store and query engine

Each of our services expose metrics in Prometheus format at `/metrics`

Monitor all the things
- RPC Request/Response cycles
- CPU / Memory / Network use
- Asynchronous processing
- C* and Distributed Locking

# Requirements
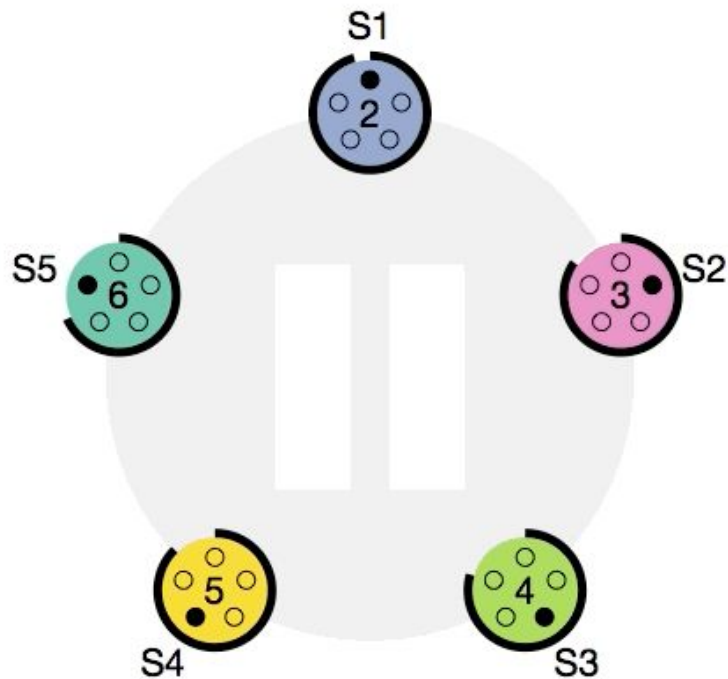
**1. Raise at most £20,000,000**
We'd agreed with institutional investors leading the funding round that £20M was the cap

**2. Ensure users have enough money**
Users should have the money they are pledging. We need to verify this before accepting the investment.

**3. Handle lots of traffic**
It was first-come-first-serve so we expected a lot of interest at the start of the crowdfunding round

**4. Don't bring down the bank**
All banking functions should continue to work whilst we're running the crowdfunding

# Requirements

**1. Raise at most £20,000,000**
We'd agreed with institutional investors leading the funding round that £20M was the cap

**2. Ensure users have enough money**
Users should have the money they are pledging. We need to verify this before accepting the investment.

**3. Handle lots of traffic**
It was first-come-first-serve so we expected a lot of interest at the start of the crowdfunding round

**4. Don't bring down the bank**
All banking functions should continue to work whilst we're running the crowdfunding

# Counters / Transactions

What if we used as Cassandra counter?

"In Cassandra, at any given moment, the counter value may be stored in the Memtable, commit log, and/or one or more SSTables. Replication between nodes can cause consistency issues in certain edge cases"

Source: https://docs.datastax.com/en/cql/3.3/cql/cql_using/useCountersConcept.html

Edge Proxy

service.crowdfunding-pre-investment

rate limited consumption

service.crowdfunding-investment

Ledger checks, confirm transaction

# Requirements

**1. Raise at most £20,000,000**
We'd agreed with institutional investors leading the funding round that £20M was the cap

**2. Ensure users have enough money**
Users should have the money they are pledging. We need to verify this before accepting the investment.

**3. Handle lots of traffic**
It was first-come-first-serve so we expected a lot of interest at the start of the crowdfunding round

**4. Don't bring down the bank**
All banking functions should continue to work whilst we're running the crowdfunding

Introduction

A brief overview of our Platform

Building a Crowdfunding Backend

**Load testing + Finding bottlenecks**

# Building our own load tester

There's some great off-the-shelf solutions for load testing
- Bees with Machine Guns
- Locust
- ApacheBench (ab)
- Gatling

# Building our own load tester

😅

At one point, we saw really high error rates in the load testing metrics. We didn't see load test requests make it to our our AWS Load Balancer.

The load test nodes were using internal DNS provided by Amazon Route 53. We were constantly resolving *.monzo.com subdomains.

## Request Timings (Load Test)

- 99th
- 95th
- 75th
- mean

## API Timings (Edge Proxy)

# Load testing in production

For our testing to create realistic load and give us useful results, we needed to test against our production systems – the real bank.

# Load testing in production

We set up our load testing system as a third "app" alongside our iOS and Android apps, and we gave it read-only access to the data we needed to test.

**Target:** Reach 1,000 app launches per second

# Scaling services

**Target:** Reach 1,000 app launches per second

Request Timings

# Scaling services

**Target:** Reach 1,000 app launches per second

```
replicas: 9
template:
  spec:
    containers:
      resources:
        limits:
          cpu: 30m
          memory: 40Mi
        requests:
          cpu: 10m
          memory: 20Mi
```

# Scaling services

**Target:** Reach 1,000 app launches per second

```
replicas: 9
template:
  spec:
    containers:
      resources:
        limits:
          cpu: 100m
          memory: 40Mi
        requests:
          cpu: 50m
          memory: 20Mi
```

"But wait, you are re-inventing autoscaling, manually?"

# Cassandra Bottlenecks

We got to around 500-600 app launches before we found a major Platform bottleneck

**Daniel Cannon** 4:34 PM

For transparency it looks like the issue might be Cassandra, we are seeing very high load averages around the time of load tests. We are now spending some time digging into this to try and find out exactly what the cause of this is.

# The numbers

21 x i3.4xlarge EC2 machines
- 16 cores
- 122GiB memory
- 2 * 1.9TiB of NVMe disks

Each node holds about 500GB of data

# Cassandra Bottlenecks

Our profiling identified three key areas
- Generating Prometheus metrics
- LZ4 Decompression
- CQL Statement Processing

CPU

# LZ4 Decompression

Read Latency (per table and instance) - p95

# CQL Statement Parsing

We saw a significant amount of time being spent in parsing CQL statements.

The majority of our applications had a fixed model during the service pod lifetime so we would've been processing the same statement over and over again.

# Prepared Statements

Cassandra supports prepared statements! Our gocql library which runs Cassandra queries was actively using them too for the majority of queries.

# Prepared Statements

```sql
SELECT id, accountid, userid, amount, currency
FROM transaction.transaction_map_Id
WHERE id = ?
```

```sql
SELECT currency, accountid, userid, id, amount
FROM transaction.transaction_map_Id
WHERE id = ?
```

# Service Mesh Bottlenecks

**Target:** Reach 1,000 app launches per second

At around 800 app launches per second, we saw our RPCs take a really long time across our Platform.

# linkerd-m4-2xlarge



## CPU Throttling (Services Only)



## CPU Usage (Cores)

# What we ended up with

- A comprehensive spreadsheet of all the services involved and how much we'd need to scale them (replicas/resource requests/limits)
- An idea of how many EC2 Kubernetes Worker Nodes we need, so we could provision them before it started
- Much more knowledge of where things can fail at this scale
- Confidence!
  - Knowing what levers you can pull when things go wrong

# Levers

No matter how much preparation we did beforehand, we wanted to ensure we could recover the Platform if anything went wrong

- Feature Toggles
  - Gracefully degrading the less critical app features
- Shedding traffic
  - Stopping the traffic before it even enters our edge

Branch: master ▾   **runbooks** / runbooks / crowdfunding / **k8s-out-of-capacity.md**    Find file   Copy path

**evnsio** Add section on changing via the AWS Console in a pinch    4223216 on 29 Nov 2018

2 contributors 

17 lines (11 sloc)   922 Bytes    Raw   Blame   History   ⌗   ✎   🗑

# Kubernetes Cluster is out of capacity

## Symptoms

New pods are not able to schedule because there is not enough capacity (CPU or RAM) on the existing Kubernetes cluster.

## Pre-checks

- Have a peek at the Kubernetes Cluster Overview dashboard to see the state of the cluster

## Resolution

Add more workers to the Kubernetes worker pool. You can do this change and apply it via Terraform or by manually changing the min, max and current to a higher value in the AWS console.

Ensure the number of nodes is divisible by 3 so we have equal capacity across our availability zones.

| App Launches | API Throughput | App Reported Errors |
|:---:|:---:|:---:|
| **20**/s | **3822**/rps | **9**/s |

### ⌄ Error Rate

| service.api.oauth2 | service.api.accounts | service.api.balance | service.api.card | service.api.config | service.api.feed |
|:---:|:---:|:---:|:---:|:---:|:---:|
| **0.13**/rps | **0**/rps | **0**/rps | **0.07**/rps | **0.3**/rps | **0.13**/rps |

| service.api.help | service.api.news | service.api.overdraft | service.api.pots | service.api.profile | service.api.spending-breakdo... |
|:---:|:---:|:---:|:---:|:---:|:---:|
| **0**/rps | **0.07**/rps | **0**/rps | **0**/rps | **0.53**/rps | **0**/rps |

### ⌄ Latency

| service.api.oauth2 | service.api.accounts | service.api.balance | service.api.card | service.api.config | service.api.feed |
|:---:|:---:|:---:|:---:|:---:|:---:|
| **100**ms | **92**ms | **207**ms | **77**ms | **104**ms | **402**ms |

| service.api.help | service.api.news | service.api.overdraft | service.api.pots | service.api.profile | service.api.spending-breakdo... |
|:---:|:---:|:---:|:---:|:---:|:---:|
| **473**ms | **51**ms | **97**ms | **163**ms | **63**ms | **443**ms |

### ⌄ Request Throughput

| service.api.oauth2 | service.api.accounts | service.api.balance | service.api.card | service.api.config | service.api.feed |
|:---:|:---:|:---:|:---:|:---:|:---:|

| 95th API Response Timing | API Error Rate | Successful Attempts | Failed Attempts | Confirmed Investments | NSQ Depth |
|---|---|---|---|---|---|
| **19** ms | **0.07** /rps | **929** /min | **16** /min | **1055** /min | **0** |

| Total Invested (pre-investment) | Total investments (confirmed) | Invested Amount (estimate) | Average investment amount |
|---|---|---|---|
| 1231013/2592520 | 1226509/2592520 | £9.503M/£20m | £761 |

> Crowdfunding Total  (7 panels)

> Crowdfunding Investments  (8 panels)

∨ Crowdfunding API

**RPC Throughput**

Success / Error

| | |
|---|---|
| 2.0K rps | 0.30 rps |
| 1.5K rps | 0.20 rps |
| 1.0K rps | |
| 500 rps | 0.10 rps |
| 0 rps | 0 rps |

09:30 — 10:00

**Server RPC Timing**

50 ms
40 ms
30 ms
20 ms
10 ms
0 ns

09:30 — 10:00

**CPU Usage (Cores)**

2.5
2.0
1.5
1.0
0.5
0

09:30 — 10:00

**Heap in use**

14 MiB
10 MiB
5 MiB
0 B

09:30 — 10:00

∨ Crowdfunding Pre Investment

**RPC Throughput**

Success / Error

| | |
|---|---|
| 80 rps | 5 rps |
| 60 rps | 4 rps |
| 40 rps | 3 rps |
| 20 rps | 2 rps |
| | 1 rps |

**Server RPC Timing**

250 ms
200 ms
150 ms
100 ms
50 ms

**CPU Usage (Cores)**

2.5
2.0
1.5
1.0
0.5

**Heap in use**

14 MiB
10 MiB
5 MiB

# Things went well

**36,006**
Investors

**£20M**
Raised

**£6.8M**
first 5 minutes

# What we learned

Here are the key takeaways and what we learnt as a result of this exercise
- Horizontal scaling has limits
- Treat software as just that, software
- Continuously load test

# Thanks!

Email:    hi@suhailpatel.com
Twitter:  @suhailpatel / @monzo

monzo