- Nanda Kishore

- Sr. Principal Engineer, Platform Engineering

- https://twitter.com/none_daa

- https://www.linkedin.com/in/none-da/

# Scale

- 28 Cities.

- 12 M customers.

- 250K Orders / day.

- 40K+ Products.
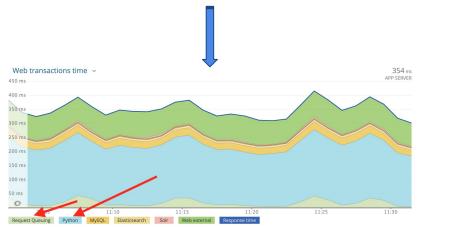
- 3x by March 2020.

- 100 Million Edge hits per day.

- 30 Million origin hits per day.

- 120K Cache Reads per sec.

- 350+ servers on AWS.

- Multiple Services (Monolith, 20+ Microservices)

- Multiple Environments for every service.

**India's largest online supermarket**

# Scalability Challenges

- Iteration 1: Scaling on CPU & Memory is not enough, nor possible always.

- Iteration 2: Scaling based on Avg response time (Newrelic)

BB | Scale: 63 + 4 = 67 Resp: ['507'] ms at - 2019-09-06 09:51:49



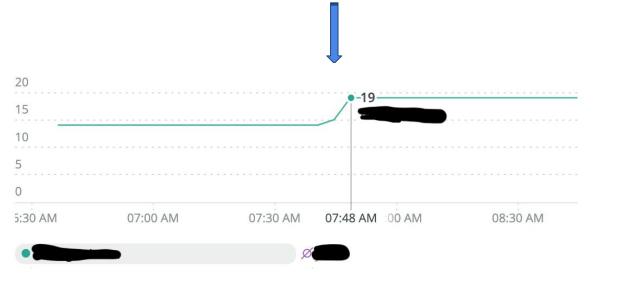**India's largest online supermarket**

# Scalability Challenges

- Iteration 3: Scaling based on Requests count



`SCALE DECISION is UPSCALE, pod count is 19`

# Scalability Challenges

- Iteration 4: Scaling based on Queue Lag (For Kafka Consumers)

```
SCALE DECISION is UPSCALE and the new workers count is: 12.
```

# Reliability Challenges with K8s ☸️

- Monolith + 🐌 Traffic = 👍 🌈 😀

- Monolith + ⚡ Traffic  = 👎 💥 😫

- Enter…..μ▲ (s)…, Enter…. K8s

- But…. K8s is NOT a SILVER BULLET!!!

- Why ?

# Reliability Challenges with K8s 🛞

- Challenge 1: Routine/scheduled deployments causing downtimes. 😱 🤦

- Note: All K8s deployments in BB are configuration driven.

# Reliability Challenges with K8s

- Usual deployment

```
 1  apiVersion: apps/v1
 2  kind: Deployment
 3  metadata:
 4    name: my-fancy-app
 5    labels:
 6      app: fancy
 7  spec:
 8    replicas: 3
 9    selector:
10      matchLabels:
11        app: fancy
12    template:
13      metadata:
14        labels:
15          app: fancy
16          version: v1
17      spec:
18        containers:
19        - name: fancy
20          image: fancy:v0.1.0
21          ports:
22          - containerPort: 80
23
```

```
 1  apiVersion: autoscaling/v2beta2
 2  kind: HorizontalPodAutoscaler
 3  metadata:
 4    name: my-fancy-hpa
 5    namespace: my-fancy-app
 6  spec:
 7    scaleTargetRef:
 8      apiVersion: apps/v1
 9      kind: Deployment
10      name: my-fancy-app
11    minReplicas: 1
12    maxReplicas: 10
13    metrics:
14    - type: Resource
15      resource:
16        name: cpu
17        target:
18          type: Utilization
19          averageUtilization: 75
20
```

**India's largest online supermarket**

# Reliability Challenges with K8s 🚢

1. Remove <mark>replicas: &lt;n&gt;</mark> setting from Deployment using :

```
1 kubectl -n my-fancy-svc apply edit-last-applied deployment/my-fancy-app
```

2. Update the config yaml (source code) accordingly.

3. Let HPA alone manage replicas count.

# Reliability Challenges with K8s ☸️

- Challenge 2: Not having a thorough understanding of maxSurge & maxUnavailable in a Deployment. (and Pdbs)

- 💡 maxSurge is rounded UP.

- 💡 maxUnavailable is rounded DOWN.

- And scaling up, scaling down happens in parallel.

# Reliability Challenges with K8s

- If a deployment has 5 pods(version=v1), with a rolling deployment of :

  - maxSurge: 30%

  - maxUnavailable: 25%

- Immediately 4 v1 pods would receive **20% extra traffic**.

# Reliability Challenges with K8s

- Challenge 3: Fair Distribution of Pods is not done by default

- Solution: Pod Anti-Affinity rules.



```
1 apiVersion: apps/v1
2 kind: Deployment
3 metadata:
4    name: my-fancy-app
5    namespace: my-fancy-svc
6 spec:
7    template:
8      spec:
9        affinity:
10         podAntiAffinity:
11           preferredDuringSchedulingIgnoredDuringExecution:
12           - podAffinityTerm:
13               labelSelector:
14                 matchExpressions:
15                 - key: app
16                   operator: In
17                   values:
18                   - my-fancy-app
19             topologyKey: kubernetes.io/hostname
```

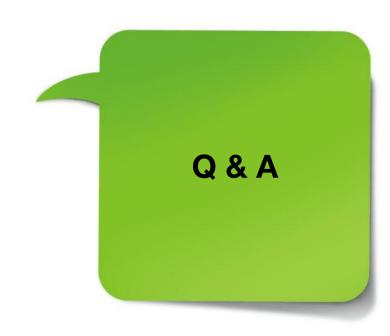# Reliability Challenges with K8s

- Challenge 4: Fair Distribution of Pods based on instance life-cycle (On-demand vs Spot).

- Solution: WORK-in-Progress

- Many more reliability challenges…

**India's largest online supermarket**

# Whats the future like ?

- Figuring out seamless upgrade strategies for fast paced projects like K8s.

- Figuring out unusual change in data patterns because of a small bug introduced by a microservice.

- Protecting services from cascade failures.

- Making critical systems like Kafka, more robust and self-healing as too many microservices evolve with light coupling architecture

Q & A

**We are hiring...**

[nandakishore.b@bigbasket.com](mailto:nandakishore.b@bigbasket.com)