# Infrastructure testing using Kubernetes

@rantav

# OR:
# Testing **Kafka replication** over the Atlantic using **Golang, Kubernetes** and friends

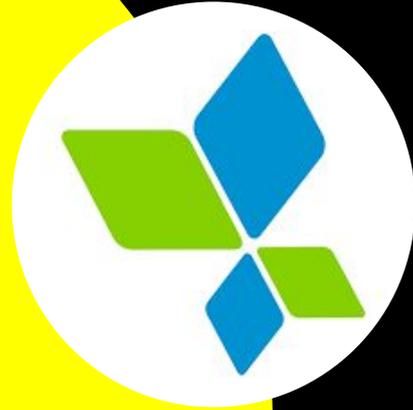# Hello!

**I am Ran Tavory**

I work on **multi-region at appsflyer**

You can find me **@rantav** (twitter/gmail/facebook)

**The Goal: Test Kafka Replication**

- ⊙ Replicate 100MB per second
- ⊙ Validate **correctness**
- ⊙ Measure **latency**
- ⊙ Run **test cases** (failure scenarios)
  - ○ Broker crash, cluster resize, packet loss...

**Does this remind anyone anything?**

- Run **test cases** (failure scenarios)
  - **Broker crash**
  - **cluster resize**
  - **packet loss**
  - **etc...**

" Jepsen: Call me maybe

// Kyle Kingsbury

# HOW ?

High level design

Kafka Replication

From https://github.com/AppsFlyer/kafka-mirror-tester

```
 ------------------------                    --------------------------------------
|                        |                  |                                      |
| Source DC              |                  | Destination DC                       |
|                        |                  |                                      |
|  ----------------      |                  |  --------------     ---------------- |
| | Source Kafka |       | - - - - - - - - ->|  | Replicator | -> | Target Kafka | |
|  ----------------      |                  |  --------------     ---------------- |
|      ↑                 |                  |                            |         |
|      |                 |                  |                            |         |
|      |                 |                  |                            ↓         |
|  ------------          |                  |                      ------------    |
| | producer |           |                  |                     | consumer |    |
|  ------------          |                  |                      ------------    |
 ------------------------                    --------------------------------------
```

With a home-grown **producer** and **consumer**
**written in Golang**

**WHAT DO WE WANT?**

DOCUMENTATION

**HOW DO WE DO IT?**

With Code

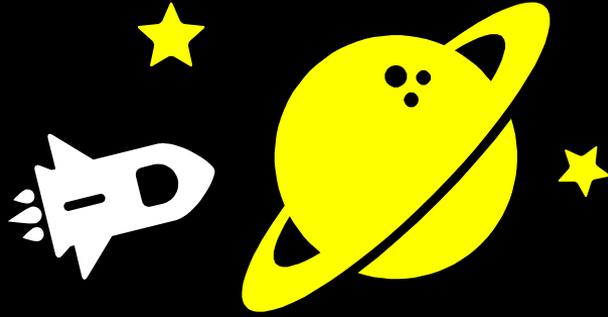# WHAT DO WE WANT?

REPRODUCIBILITY

## HOW DO WE DO IT?

With Code

# WHAT DO WE WANT?

## MONITORING

# HOW DO WE DO IT?

With Code

# KUBERNETES

# A brief introduction to **Kubernetes**

Kubernetes is a

# CLUSTER MANAGER

It manages

# NODES

The **host** running a kubelet and a set of Pods

Nodes contain

# PODS

Pods consist of 1 or more **containers**

Pods run

# DEPLOYMENTS

A set of **stateless** pods

Pods also run

# STATEFUL SET

A set of **stateful** pods

You control k8s using

# KUBECTL

The k8s CLI

```
$ make k8s-all
```

# What does that do?

- Provision VMs in AWS
- In two regions. Ireland and Virginia.
- Setup VPCs, Subnets, Routing Tables
- Create Security Groups
- Setup Load Balancers
- Install Kubernetes (etcd, masters, nodes)
- Setup Monitoring (Prometheus & Grafana and install dashboard)
- Install Weave Scope
- Install Kafkas (and test them)
- Install **uReplicator / Brooklin** (and test them)
- Install test programs (written by my in Go)
- And more... (ASGs, DHCP etc)

💡

**What are we building?**

us-east-1

eu-west-1

**What are we building?**

us-east-1

**40** nodes
k8s cluster

eu-west-1

**48** nodes
k8s cluster

## What are we building?

**us-east-1**

**40** nodes
k8s cluster

**30** brokers
kafka cluster

**eu-west-1**

**48** nodes
k8s cluster

**30** brokers
kafka cluster

**What are we building?**

**us-east-1**

**40** nodes
k8s cluster

**30** brokers
kafka cluster

**eu-west-1**

**48** nodes
k8s cluster

**30** brokers
kafka cluster

**8** workers
uReplicator
/
32 Brooklin

**What are we building?**

## us-east-1

**40** nodes
k8s cluster

**30** brokers
kafka cluster

**10** producers

## eu-west-1

**48** nodes
k8s cluster

**30** brokers
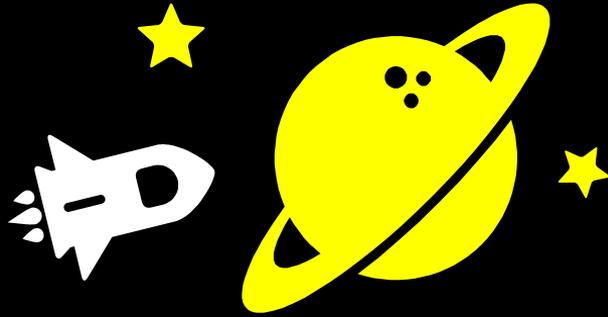kafka cluster

uReplicator /
Brooklin

**4** consumers

# End Result

5FxkdGcuNdiNwMt7bALgUenZx7IAAj04-YO3qgRNtwCSU5F23EQBmN4-3Dbeb9w9TWkY4Ey6GRVU6-VC7Ifz8t0CM4KtNxf4gyPIQ8z7DmioC6YuBNkvwMW06BEexKY9F7S_iqaC6Yrb
NRoY_Gg

```
        Now run: kubectl --context us-east-1.k8s.local proxy
        And then open http://localhost:8001/api/v1/namespaces/kube-system/services/https:kubernetes-dashboard:/proxy/
        🔍 Weave Scope us-east-1: http://a2f7e98f3380511e9bc0c1234b713cd3-37681181.us-east-1.elb.amazonaws.com
        ✅ Prometheus us-east-1: http://a051a4719380411e9bc0c1234b713cd3-1186349614.us-east-1.elb.amazonaws.com:9090

Monitor low level cluster events:
        kubectl --context us-east-1.k8s.local get events --watch --all-namespaces -o wide


>>> Admin for eu-west-1:
Name:           eks-admin-token-57pf9
Namespace:      kube-system
Labels:         <none>
Annotations:    kubernetes.io/service-account.name: eks-admin
                kubernetes.io/service-account.uid: 3b3a3cd2-3804-11e9-a5eb-0695644d7ac0

Type:  kubernetes.io/service-account-token

Data
====
token:          eyJhbGciOiJSUzI1NiIsImtpZCI6IiJ9.eyJpc3MiOiJrdWJlcm5ldGVzL3NlcnZpY2VhY2NvdW50Iiwia3ViZXJuZXRlcy5pby9zZXJ2aWNlYWNjb3VudC9uYW1lc3B
bmV0ZXMuaW8vc2VydmljZWFjY291bnQvc2VjcmV0Lm5hbWUiOiJla3MtYWRtaW4tdG9rZW4tNTdwZjkiLCJrdWJlcm5ldGVzLmlvL3NlcnZpY2VhY2NvdW50L3NlcnZpY2UtYWNjb3Vu
XJuZXRlcy5pby9zZXJ2aWNlYWNjb3VudC9zZXJ2aWNlLWFjY291bnQudWlkIjoiM2IzYTNjZDItMzgwNC0xMWU5LWE1ZWItMDY5NTY0NGQ3YWMwIiwic3ViIjoic3lzdGVtOnNlcnZpY
1hZG1pbiJ9.b-ERTE0fw-XBKVL9he7amprJVOKLkedeH1wcR3Q5Mp6gmNi7fn0n7ZIegINMOOJEN7GUSEh4YHLAtGA9yfuA0yNJF884hbWwdr_Zl355daKzf4HJ_-5SL5YQyLqrVwS24
7kuHRiwN8l_XwhiDNypa_X7NbEEpkmlgFZ5cHWS4RwvLmsJXI4KEOZ7RhiUAA8rL8OClDEZr7P1NQb7CuaZj7n0TTdBcwbQ7ChZovM16ILskeqLWk2ucWY8gtYnlWHCPHygSUZaHE5db
HEFPw-Q
ca.crt:         1042 bytes
namespace:      11 bytes
        Now run: kubectl --context eu-west-1.k8s.local proxy --port 8002
        And then open http://127.0.0.1:8002/api/v1/namespaces/kube-system/services/https:kubernetes-dashboard:/proxy/

Monitor low level cluster events:
        kubectl --context eu-west-1.k8s.local get events --watch --all-namespaces -o wide

        🔍 Weave Scope eu-west-1: http://a34fbd169380511e9a5eb0695644d7ac-1299023419.eu-west-1.elb.amazonaws.com
        ✅ Prometheus eu-west-1: http://a32782681380411e9a5eb0695644d7ac-1477431009.eu-west-1.elb.amazonaws.com:9090
        ✅ Grafana (user/pass: admin/admin): http://a2d0544df380411e9a5eb0695644d7ac-1371947342.eu-west-1.elb.amazonaws.com:3000
→  kafka-mirror-tester git:(master)
→  kafka-mirror-tester git:(master)
```

# Grafana Dashboard

30

# Weave Scope

# Weave Scope



kafka-mirror-te...
10 pods

The Internet
Inbound connecti...

kubernetes
0 pods

pzoo    0 pods

broker    29 pods

zookeeper
1 pod

## kubectl example

```
kubectl --context eu-west-1.k8s.local -n kafka-destination get po -o wide
NAME                    READY   STATUS     RESTARTS   AGE    IP                NODE
kafka-destination-0     2/2     Running    2          16m    100.97.161.67     ip-172-20-77-119.eu-west-1.compute.internal
kafka-destination-1     2/2     Running    0          15m    100.99.203.131    ip-172-20-77-91.eu-west-1.compute.internal
kafka-destination-10    2/2     Running    0          10m    100.124.129.131   ip-172-20-44-138.eu-west-1.compute.internal
kafka-destination-11    2/2     Running    0          10m    100.115.218.2     ip-172-20-123-96.eu-west-1.compute.internal
kafka-destination-12    2/2     Running    0          9m     100.117.24.2      ip-172-20-108-204.eu-west-1.compute.internal
kafka-destination-13    2/2     Running    0          8m     100.97.229.66     ip-172-20-97-229.eu-west-1.compute.internal
kafka-destination-14    2/2     Running    0          8m     100.109.13.194    ip-172-20-91-204.eu-west-1.compute.internal
kafka-destination-15    2/2     Running    0          7m     100.99.164.66     ip-172-20-114-89.eu-west-1.compute.internal
kafka-destination-16    2/2     Running    0          7m     100.114.149.195   ip-172-20-97-100.eu-west-1.compute.internal
kafka-destination-17    2/2     Running    0          6m     100.101.233.3     ip-172-20-73-176.eu-west-1.compute.internal
kafka-destination-18    2/2     Running    0          6m     100.107.249.67    ip-172-20-47-107.eu-west-1.compute.internal
kafka-destination-19    2/2     Running    0          5m     100.111.184.67    ip-172-20-109-255.eu-west-1.compute.internal
kafka-destination-2     2/2     Running    0          14m    100.125.0.130     ip-172-20-41-241.eu-west-1.compute.internal
kafka-destination-20    2/2     Running    0          5m     100.116.89.66     ip-172-20-75-179.eu-west-1.compute.internal
kafka-destination-21    2/2     Running    0          4m     100.99.191.130    ip-172-20-39-5.eu-west-1.compute.internal
kafka-destination-22    2/2     Running    0          3m     100.97.200.2      ip-172-20-77-210.eu-west-1.compute.internal
kafka-destination-23    2/2     Running    0          3m     100.106.224.67    ip-172-20-111-29.eu-west-1.compute.internal
kafka-destination-24    2/2     Running    0          2m     100.117.212.66    ip-172-20-45-107.eu-west-1.compute.internal
kafka-destination-25    2/2     Running    0          1m     100.126.211.195   ip-172-20-104-30.eu-west-1.compute.internal
kafka-destination-26    2/2     Running    0          1m     100.113.255.66    ip-172-20-79-147.eu-west-1.compute.internal
kafka-destination-27    2/2     Running    0          46s    100.100.86.2      ip-172-20-108-2.eu-west-1.compute.internal
kafka-destination-28    0/2     Init:0/1   0          1s     <none>            ip-172-20-62-131.eu-west-1.compute.internal
kafka-destination-3     2/2     Running    0          14m    100.117.8.130     ip-172-20-40-131.eu-west-1.compute.internal
kafka-destination-4     2/2     Running    0          13m    100.110.127.2     ip-172-20-38-150.eu-west-1.compute.internal
kafka-destination-5     2/2     Running    0          13m    100.126.198.67    ip-172-20-124-136.eu-west-1.compute.internal
kafka-destination-6     2/2     Running    0          12m    100.119.5.2       ip-172-20-118-231.eu-west-1.compute.internal
kafka-destination-7     2/2     Running    0          12m    100.125.161.194   ip-172-20-111-190.eu-west-1.compute.internal
kafka-destination-8     2/2     Running    0          11m    100.114.9.131     ip-172-20-47-40.eu-west-1.compute.internal
kafka-destination-9     2/2     Running    0          11m    100.101.32.67     ip-172-20-46-70.eu-west-1.compute.internal
pzoo-destination-0      2/2     Running    0          16m    100.111.184.66    ip-172-20-109-255.eu-west-1.compute.internal
```
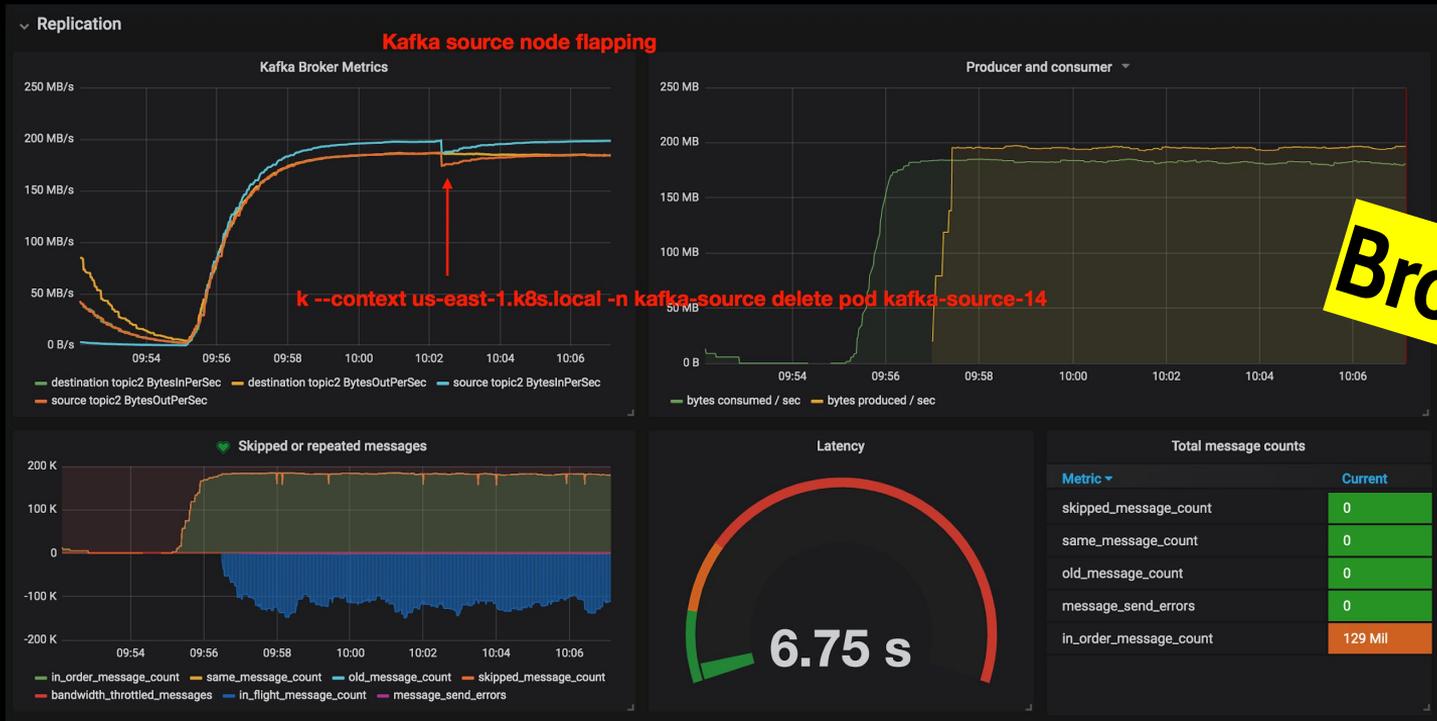
# TEST IT

Let's put all this to test now

# KILL A BROKER

```
$ kubectl \
--context us-east-1.k8s.local \
-n kafka-source \
delete pod kafka-source-2
```
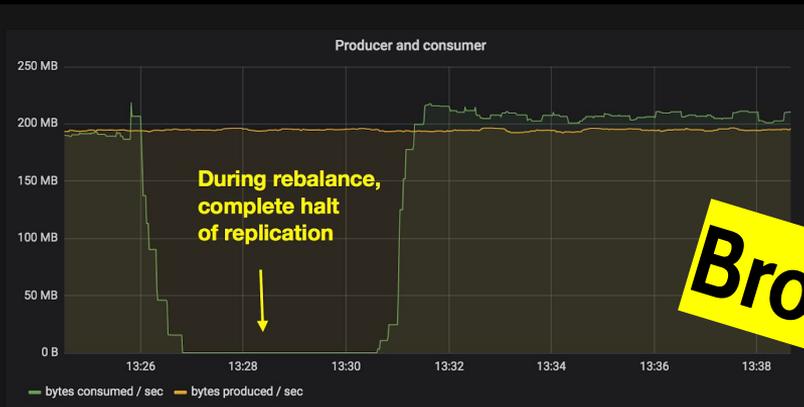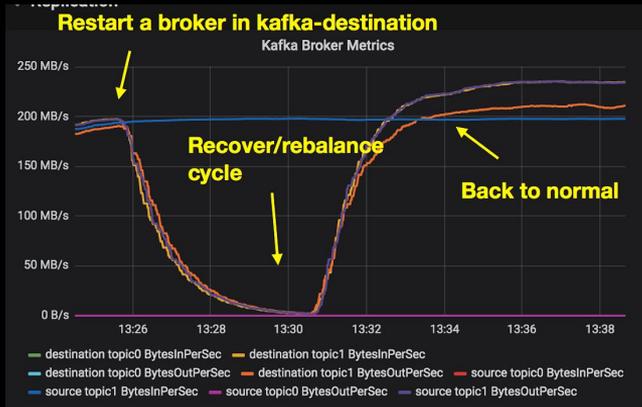
# KILL A BROKER - SRC

uReplicator

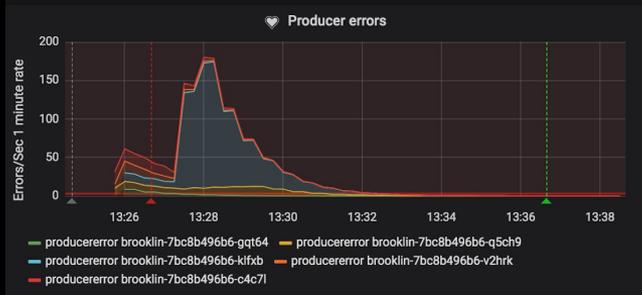# KILL A BROKER - DST



**Restart a broker in kafka-destination**

Kafka Broker Metrics

**Recover/rebalance cycle**

**Back to normal**

- destination topic0 BytesInPerSec — destination topic1 BytesInPerSec
- destination topic0 BytesOutPerSec — destination topic1 BytesOutPerSec — source topic0 BytesInPerSec
- source topic1 BytesInPerSec — source topic0 BytesOutPerSec — source topic1 BytesOutPerSec

Producer and consumer

**During rebalance, complete halt of replication**

- bytes consumed / sec — bytes produced / sec

Brooklin

♡ Producer errors

- producererror brooklin-7bc8b496b6-gqt64 — producererror brooklin-7bc8b496b6-q5ch9
- producererror brooklin-7bc8b496b6-klfxb — producererror brooklin-7bc8b496b6-v2hrk
- producererror brooklin-7bc8b496b6-c4c7l

Latency

**4.27 min**

Total message counts

| Metric ▾ | Current |
| --- | --- |
| skipped_message_count | 29 K |
| same_message_count | 0 |
| old_message_count | 2 K |
| message_send_errors | 0 |
| in_order_message_count | 196 Mil |

38

# RESIZE KAFKA CLUSTER

```
$ kubectl \
--context us-east-1.k8s.local \
-n kafka-source \
scale \
statefulset kafka-source \
--replicas 29
```

# RESIZE KAFKA CLUSTER

# ADD uREPLICATOR WORKER

```
$ kubectl \
    --context eu-west-1.k8s.local \
    -n ureplicator \
    scale deployment \
    ureplicator-worker \
    --replicas 9
```

# ADD uReplicator WORKER

# ADD Brooklin WORKER



**brooklin transmitted bytes**

New worker joined

Brooklin

# ADD NEW TOPIC

`$ make k8s-redeploy-tests`

```
args:
  - produce
  - --bootstrap-servers
  - broker.kafka-source.svc.cluster.local:9092
  - --id
  - $(ID)
  - --message-size
  - "1000"
  - --throughput
  - "20000"
  - --topics
  - topic0
  - --retention
  - "300000"
  - --num-partitions
```

# ADD NEW TOPIC

# ADD PARTITIONS

```
$ make k8s-kafka-shell-source
```

```
$ bin/kafka-topics.sh --zookeeper
zookeeper:2181 --alter --topic
topic5 --partitions 300
```

# ADD PARTITIONS

Problem: **uReplicator** does not see the new partitions

uReplicator

# ADD PARTITIONS (fix)

uReplicator

```
$ kubectl --context eu-west-1.k8s.local \
    -n ureplicator port-forward \
    ureplicator-controller-76ff85b889-l9mzl 9000

$ curl -X DELETE http://localhost:9000/topics/topic5

$ curl -X POST -d \
    '{"topic":"topic5", "numPartitions":"300"}' \
    http://localhost:9000/topics
```

# ADD PARTITIONS



Add partitions dest first then source

brooklin takes a while to start replicating the partitions

all partitions are replicated

Brooklin

Kafka Broker Metrics

Producer and consumer

There is some message loss, possibly due to low retention and the time it takes brooklin to start

Producer errors

Latency

3.74 min

| Metric | Current |
| --- | --- |
| skipped_message_count | 4 Mil |
| same_message_count | 0 |
| old_message_count | 0 |
| message_send_errors | 0 |
| in_order_message_count | 409 Mil |

# PACKET LOSS

# PACKET LOSS

uReplicator

# PACKET LOSS



**Kafka Broker Metrics**

Packet loss deactivated

Packet loss activated on 3 source brokers

**Producer and consumer**

No message loss

Brooklin

**Skipped or repeated messages**

Increased 99%ile latency

**Message arrival latency**

**Total message counts**

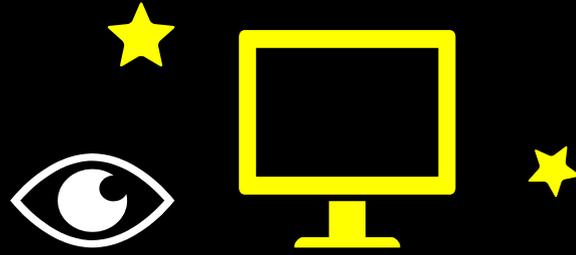| Metric | Current |
| --- | --- |
| skipped_message_count | 0 |
| same_message_count | 0 |
| old_message_count | 0 |
| message_send_errors | 0 |
| in_order_message_count | 218 Mil |

# This is so much FUN!!!

But how do we actually do that?

# CODE DEEP DIVE

Golang Producer/Consumer

(just a few bites)

**Producer**

**Main loop**

```go
    // the rate limiter regulates the producer by limiting its
    limiter := rate.NewLimiter(rate.Limit(throughput), int(math

    // Sequence number per message
    seq := initialSequence

    go eventsProcessor(p, errorCounter)

    topicString := string(topic)
    tp := kafka.TopicPartition{Topic: &topicString, Partition:
    for ; ; seq++ {
      err := limiter.Wait(ctx)
      if err != nil {
        log.Errorf("Error waiting %+v", err)
        continue
      }
      messageKey := types.MessageKey(uint(seq) % numPartitions)
      scopedSeq := seq / types.SequenceNumber(numPartitions)
      produceMessage(ctx, p, tp, producerID, messageKey, scopedS
    }
}
```

# Consumer

# Main loop

```go
sigchan := make(chan os.Signal, 1)
signal.Notify(sigchan, syscall.SIGINT, syscall.SIGTERM)

for {
  select {
  case sig := <-sigchan:
    log.Infof("Caught signal %v: terminating", sig)
    return
  case <-ctx.Done():
    log.Infof("Done. %s", ctx.Err())
    return
  case ev := <-c.Events():
    // Most events are typically juse messages, still w
    // Partition changes, EOF and Errors
    switch e := ev.(type) {
    case kafka.AssignedPartitions:
      log.Infof("AssignedPartitions %v", e)
      c.Assign(e.Partitions)
    case kafka.RevokedPartitions:
      log.Infof("RevokedPartitions %v", e)
      c.Unassign()
    case *kafka.Message:
      processMessage(e, useMessageHeaders)
    case kafka.PartitionEOF:
      log.Debugf("PartitionEOF Reached %v", e)
    case kafka.Error:
      // Errors should generally be considered as infor
      log.Errorf("Error: %+v", e)
    }
  }
}
```

**Consumer**

**Process message**

```go
// Process a single message, ke
func processMessage(
  msg *kafka.Message,
  useMessageHeaders bool,
) {
  data := message.Extract(msg,
  log.Tracef("Data: %s", data)
  validateSequence(data)
  collectThroughput(data)
  collectLatencyStats(data)
}
```

" Question:  How can multiple consumers validate message arrival order?

Redis?

DynamoDB?

No! There's a trick!

# Message format

```
+------------------------------------------------------------------+
| producer-id;sequence-number;timestamp;payload...|
+------------------------------------------------------------------+
```

**Producer**

**Sequence numbers**

```
messageKey = seq % partitions
perKeySeq  = seq / partitions
```

But - there is still a bug here...
Hint: 🎁 🎂 🎉

the BIRTHDAY PARADOX

**Producer**

**Sequence numbers**

Fix:

```
partitions =  partitions * 17

messageKey = seq % partitions
perKeySeq  = seq / partitions
```

# CODE DEEP DIVE

Kubernetes, head first

(just a few bites)

**Kafka**

**There is also today a Kafka operator**

```yaml
apiVersion: apps/v1beta2
kind: StatefulSet
metadata:
  name: kafka-source
  namespace: kafka-source
spec:
  selector:
    matchLabels:
      app: kafka-source
  serviceName: "broker"
  replicas: 30
  updateStrategy:
    type: OnDelete
  template:
    metadata:
      labels:
        app: kafka-source
    spec:
      terminationGracePeriodSeconds: 30
      initContainers:
```

kind: StatefulSet

replicas: 30

# Metrics sidecar container

```
          mountPath: /var/lib/kafka/data
        - name: metrics
          image: solsson/kafka-prometheus-jmx-exporter
          command:
          - java
          - -XX:+UnlockExperimentalVMOptions
          - -XX:+UseCGroupMemoryLimitForHeap
          - -XX:MaxRAMFraction=1
          - -XshowSettings:vm
          - -jar
          - jmx_prometheus_httpserver.jar
          - "5556"
          - /etc/jmx-kafka/jmx-kafka-prometheus.yml
          ports:
          - name: prometheus
            containerPort: 5556
          resources:
            requests:
              cpu: 100m
              memory: 500Mi
          volumeMounts:
          - name: jmx-config
            mountPath: /etc/jmx-kafka
```

# Kafka anti-affinity

```
affinity:
  podAntiAffinity:
    requiredDuringSchedulingIgnoredDuringExecution:
    - labelSelector:
        matchExpressions:
        - key: app
          operator: In
          values:
          - kafka-source
```

**uReplicator deployment**

```yaml
apiVersion: extensions/v1beta1
kind: Deployment
metadata:
  namespace: ureplicator
  name: ureplicator-worker
  labels:
    app: ureplicator
    component: worker
spec:
  replicas: 8
  selector:
    matchLabels:
      app: ureplicator
      component: worker
```

## Service monitoring

```yaml
apiVersion: monitoring.coreos.com/v1
kind: ServiceMonitor
metadata:
  labels:
    k8s-app: kafka-mirror-tester-producer
  name: kafka-mirror-tester-producer
  namespace: monitoring
spec:
  endpoints:
  - port: metrics
  jobLabel: k8s-app
  namespaceSelector:
    matchNames:
    - default
  selector:
    matchLabels:
      app: kafka-mirror-tester-producer
```

`$ make k8s-delete-all`

```
→ kafka-mirror-tester git:(master) make k8s-delete-all
make k8s-delete-cluster-eu-west-1& make k8s-delete-cluster-us-east-1
kops delete cluster --state s3://us-east-1.k8s.local  us-east-1.k8s.local --yes
kops delete cluster --state s3://eu-west-1.k8s.local  eu-west-1.k8s.local --yes
TYPE                    NAME                                                              ID
autoscaling-config      master-eu-west-1c.masters.eu-west-1.k8s.local-20190224065453      master-eu-west-1c.masters.eu-west-1.k8s.local-20190224065453
autoscaling-config      nodes.eu-west-1.k8s.local-20190224065453                          nodes.eu-west-1.k8s.local-20190224065453
autoscaling-group       master-eu-west-1c.masters.eu-west-1.k8s.local                     master-eu-west-1c.masters.eu-west-1.k8s.local
autoscaling-group       nodes.eu-west-1.k8s.local                                         nodes.eu-west-1.k8s.local
dhcp-options            eu-west-1.k8s.local                                               dopt-06269241b3a1ceaa9
iam-instance-profile    masters.eu-west-1.k8s.local                                       masters.eu-west-1.k8s.local
iam-instance-profile    nodes.eu-west-1.k8s.local                                         nodes.eu-west-1.k8s.local
iam-role                masters.eu-west-1.k8s.local                                       masters.eu-west-1.k8s.local
iam-role                nodes.eu-west-1.k8s.local                                         nodes.eu-west-1.k8s.local
instance                master-eu-west-1c.masters.eu-west-1.k8s.local                     i-06badabd6f115f8cf
instance                nodes.eu-west-1.k8s.local                                         i-004a39091a7d9f8fe
instance                nodes.eu-west-1.k8s.local                                         i-01c0b3b098667d1ed
instance                nodes.eu-west-1.k8s.local                                         i-01d8f59fb81b21850
instance                nodes.eu-west-1.k8s.local                                         i-01de61741cec1804a
instance                nodes.eu-west-1.k8s.local                                         i-0222b59a4f37e3903
instance                nodes.eu-west-1.k8s.local                                         i-026747ceca5b33040
instance                nodes.eu-west-1.k8s.local                                         i-02ddff43751175be2
instance                nodes.eu-west-1.k8s.local                                         i-02e65755a0a49da67
instance                nodes.eu-west-1.k8s.local                                         i-0354fc79b5d4d8549
instance                nodes.eu-west-1.k8s.local                                         i-03d8b2261106f9a07
instance                nodes.eu-west-1.k8s.local                                         i-04f3a2cf1413abb00
instance                nodes.eu-west-1.k8s.local                                         i-0587c8d514c9e1b91
instance                nodes.eu-west-1.k8s.local                                         i-0623ac0f93fbcb75e
instance                nodes.eu-west-1.k8s.local                                         i-06757c8a72405efd6
instance                nodes.eu-west-1.k8s.local                                         i-06d804fbaa5edc3ed
instance                nodes.eu-west-1.k8s.local                                         i-06e643b30f0898357
instance                nodes.eu-west-1.k8s.local                                         i-0723f19aedc715ebb
instance                nodes.eu-west-1.k8s.local                                         i-072c417566d4321d7
instance                nodes.eu-west-1.k8s.local                                         i-07479f8b27cf25ce8
instance                nodes.eu-west-1.k8s.local                                         i-078526c2abe80e110
instance                nodes.eu-west-1.k8s.local                                         i-084121996022da5c2
instance                nodes.eu-west-1.k8s.local                                         i-084a82289d67496fc
instance                nodes.eu-west-1.k8s.local                                         i-087f7ddaa6bf2cd09
instance                nodes.eu-west-1.k8s.local                                         i-08f14b87fe19ba471
```

# Thanks!

## Any questions?

You can find me at @rantav &
rantav@appsflyer.com

This presentation:
https://speakerdeck.com/rantav/infrastructure-testing-using-kubernetes-and-go

The project: https://github.com/AppsFlyer/kafka-mirror-tester