adform

# HOW WE MANAGE OUR KAFKA CLUSTERS

**Robert Fabisiak**
**Tomasz Gintowt**

**21.11.2019 Warszawa**

FORWARD. TOGETHER.
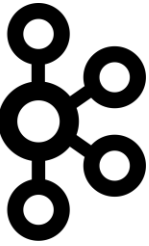
Robert Fabisiak
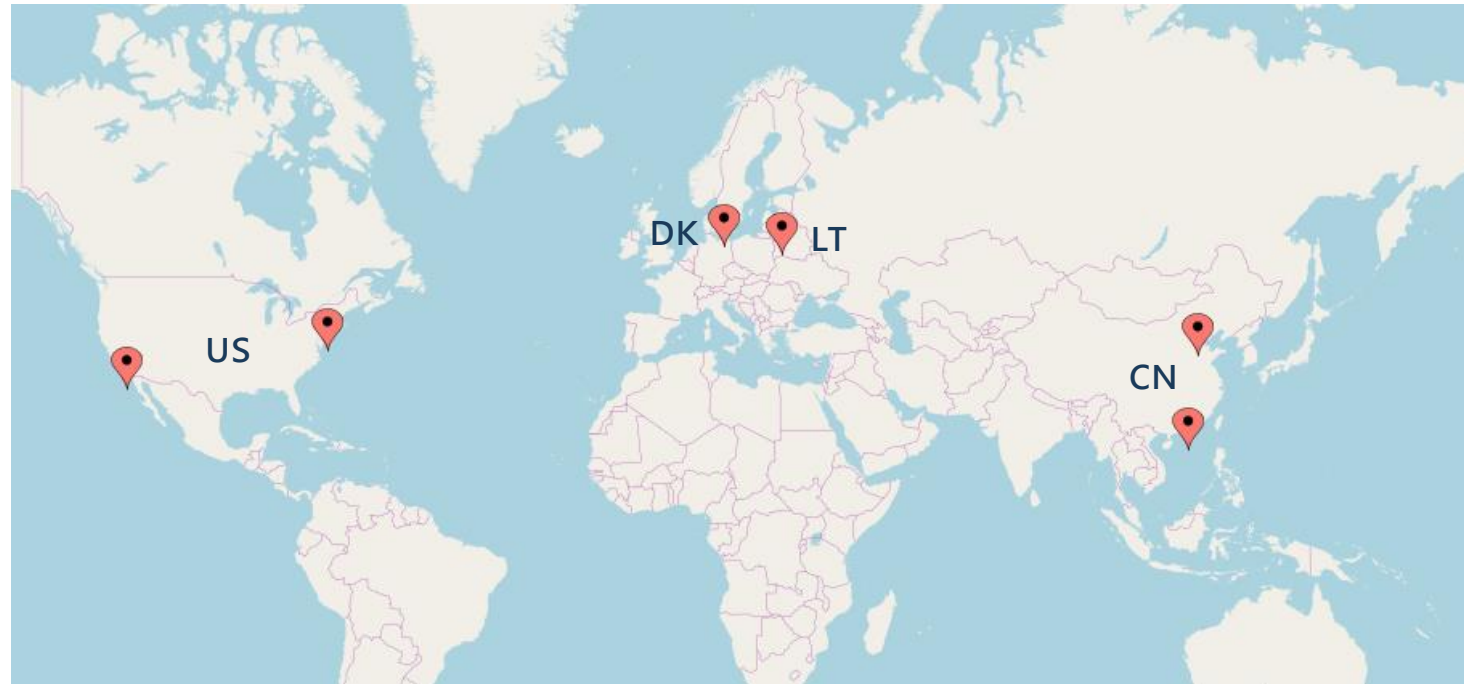
Tomasz Gintowt

# Agenda

- Our clusters
- Manage clusters with Ansible
- Tools and Metrics
- Problems and Solutions

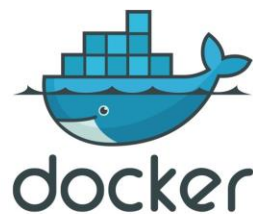adform **FORWARD. TOGETHER.**

# Kafka Clusters

- More than 100 servers (50% Bare Metal, 50% VM)
- Biggest cluster
  - 51 HW nodes
  - 460TiB storage
  - 16 TiB RAM
  - 2400 VCores

- **181 000 000 000 messages per day**
- Mirror data from remote clusters

# Technologies

# Automate upgrades

- **Follow official guidelines**
- Rolling upgrade with serial=1
    - stop and upgrade package
    - change configuration
    - start service
    - verify cluster health
- **Upgrade time reduced from weeks to days**
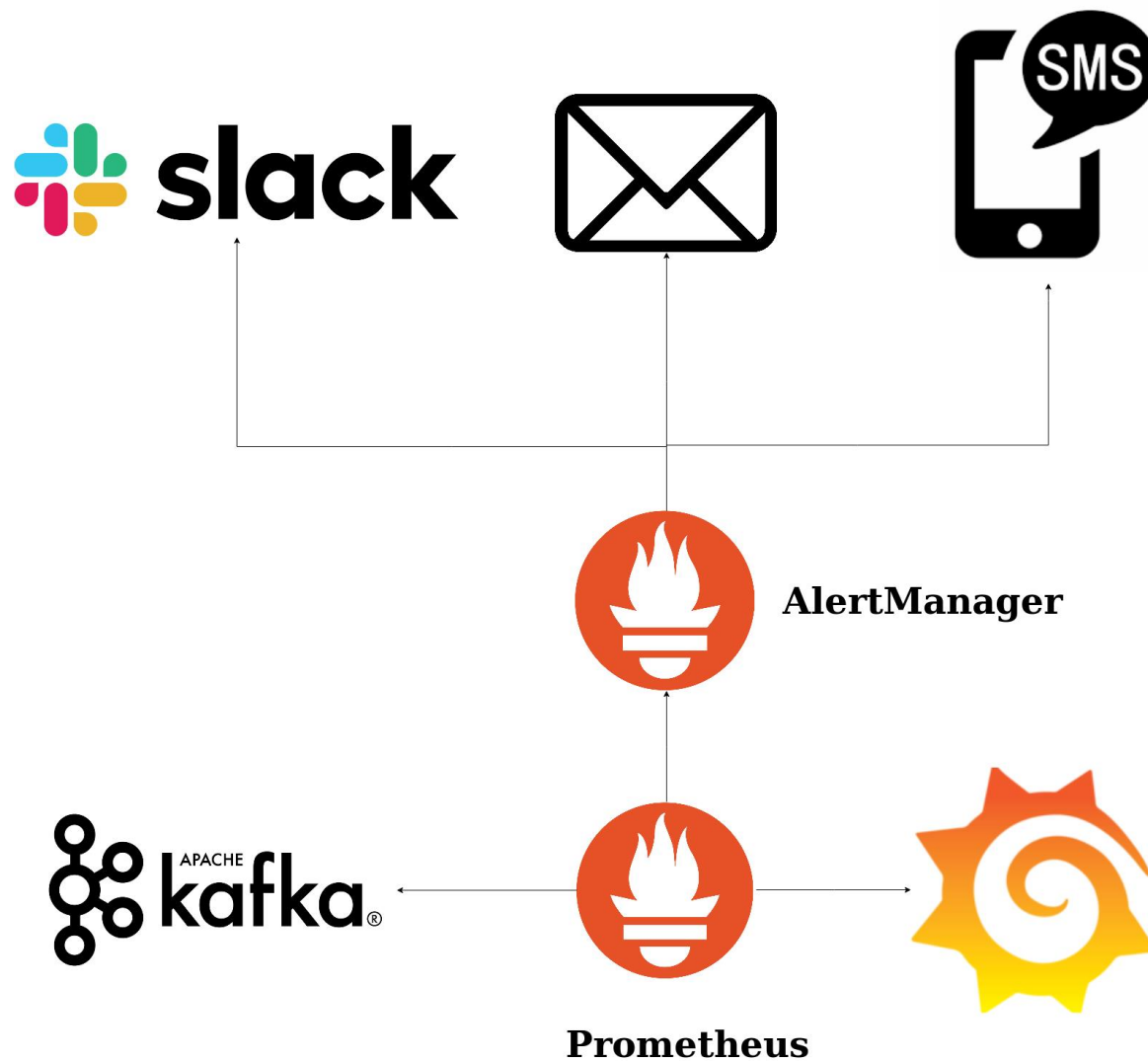
# Automate topic management

- deploy request
- CI build job
  - verify and validate YAML
  - check cluster utilization
  - check user quota
  - send notify

```yaml
---

kafka_topics:
  manager_host: control.node
  zookeeper: zookeeper.node:2181
  envs:
    KAFKA_OPTS: -Djava.security.auth.login.config=/etc/kafka/jaas.conf
  topics:
    high:
      partitions: 10
      replication_factor: 2
    load:
      partitions: 10
      replication_factor: 1
    strategy:
      partitions: 20
      replication_factor: 1
    adform:
      partitions: 19
      replication_factor: 9
      config:
        - max.message.bytes=107374182400 # 100GiB
        - retention.ms=604800000 # 7d
```

# Monitoring - Grafana



FORWARD. TOGETHER.

8

# Monitoring - alerting

# Monitoring - metrics

**Kafka**
- Messages in/out
- Offline partitions
- Broker up/down
- Partition Leader per broker
- Prefered replica imbalance
- URP
- LogCleaner
- Consumer Lag

**JVM**
- Heap Memory used
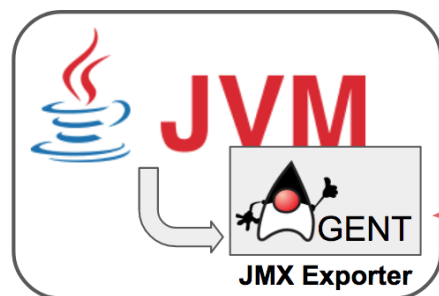- GC Process time

**Zookeeper**
- Nodes
- Followers
- In-Sync Followers
- Pending Syncs

# JMX exporter

- Installed on all nodes
- ~4k metrics per node
- Whitelisted MBeans objects:
  - Java
  - Cluster
  - Controller
  - Log
  - Network
  - Server

# Kafka exporter

- Single exporter per cluster
- Most important metrics:
  - kafka_consumergroup_lag
  - kafka_topic_partition_leader
  - under_replicated_partition

**JVM**

GENT

**JMX Exporter**

HTTP scrape

**Prometheus**

**FORWARD. TOGETHER.**

# Kafka-tools

Partition reassignment and preferred replica election.
Modules:

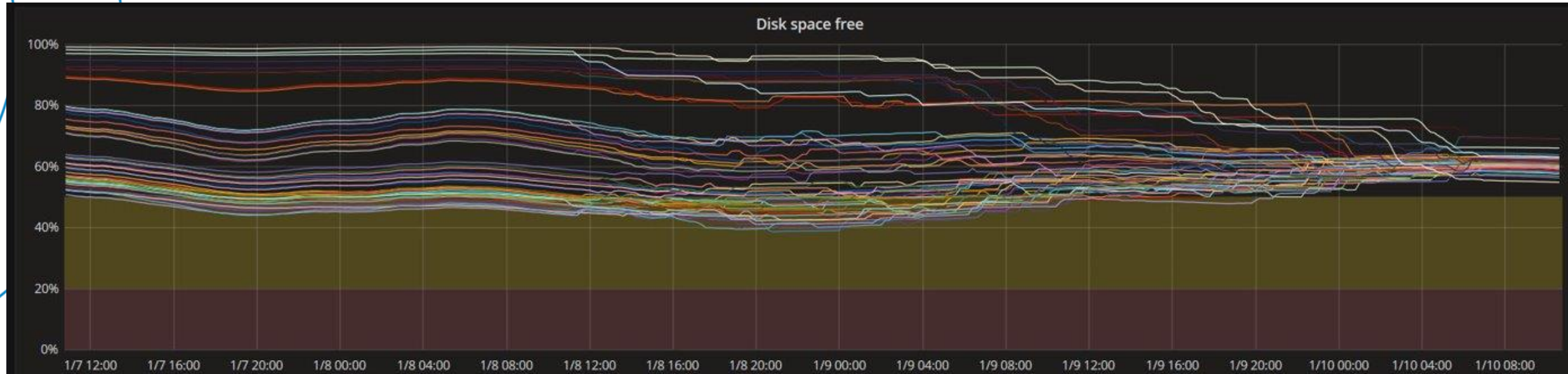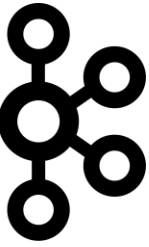- **Balance**
- Clone
- Elect
- Remove
- Trim

Type of balance:
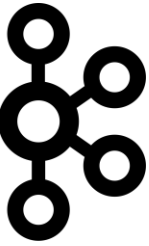
- Count
- **Size**
- Even
- Leader

# Kafka Manager

Kafka Manager  `cpkfla`  Cluster ▾  Brokers  Topic ▾  Consumers

Clusters / cpkfla / Brokers

← Brokers

| Id | Host | Port | JMX Port | Bytes In | Bytes Out |
|----|------|------|----------|----------|-----------|
| 1 | cpkfla001prvla2.lin.pr.adform.zone | PLAINTEXT:9092 | 9998 | 734k | 10m |
| 2 | cpkfla002prvla2.lin.pr.adform.zone | PLAINTEXT:9092 | 9998 | 305k | 3.8m |
| 3 | cpkfla003prvla2.lin.pr.adform.zone | PLAINTEXT:9092 | 9998 | 497k | 996k |
| 4 | cpkfla004prvla2.lin.pr.adform.zone | PLAINTEXT:9092 | 9998 | 170k | 9.7m |
| 5 | cpkfla005prvla2.lin.pr.adform.zone | PLAINTEXT:9092 | 9998 | 437k | 3.4m |
| 6 | cpkfla006prvla2.lin.pr.adform.zone | PLAINTEXT:9092 | 9998 | 199k | 6.6m |

adform

**FORWARD. TOGETHER.**

# Problems and Solutions

Data rebalance results

# Problems and Solutions

## Data rebalance with offline partitions

| Partition Information | | | | |
| --- | --- | --- | --- | --- |
| **Partition** | **Latest Offset** | **Leader** | **Replicas** | **In Sync Replicas** |
| 0 | 100 | -1 | (1) | (1) |

```
$ kafka-reassign-partitions --zookeeper zookeeper.service:2181/kafka --reassignment-json-file /tmp/a.json --verify
Status of partition reassignment:
Reassignment of partition high_load_strategy_adform-0 is still in progress
```

**Solution**
- reelect kafka controller with zookeper CLI
- always configure replicas for topics

# Problems and Solutions

## Wrong record timestamp



**Solution**
- Configure
  `log.message.timestamp.difference.max.ms`

# Problems and Solutions

## Consumer group constant rebalance



**Solution**
- `Monitor consumer performance`
- `session.timeout.ms`
- `max.poll.interval.ms`
- `heartbeat.interval.ms`

**FORWARD. TOGETHER.**

# Problems and Solutions

## NULL key is not null

### Partition Information

| Partition | Latest Offset | Leader | Replicas | In Sync Replicas | Preferred Leader? | Under Replicated? |
|-----------|---------------|--------|----------|------------------|-------------------|-------------------|
| 0 | 0 | 1 | (1) | (1) | true | false |
| 1 | 11,111,111 | 2 | (2) | (2) | true | false |
| 2 | 0 | 3 | (3) | (3) | true | false |
| 3 | 0 | 1 | (1) | (1) | true | false |

**Solution:**
- `Not always NULL mean null, know your`
  `libs.`

**THANK YOU!**

adform

FORWARD. TOGETHER.