# Tales from Lastminute.com machine room: our journey towards a full on-premise kubernetes architecture in production

michele.orsi@lastminute.com
manuel.ranieri@lastminute.com

lastminute.comgroup

# An inspiring travel company ..

# A tech company to the core

Tech department: **300+ people**
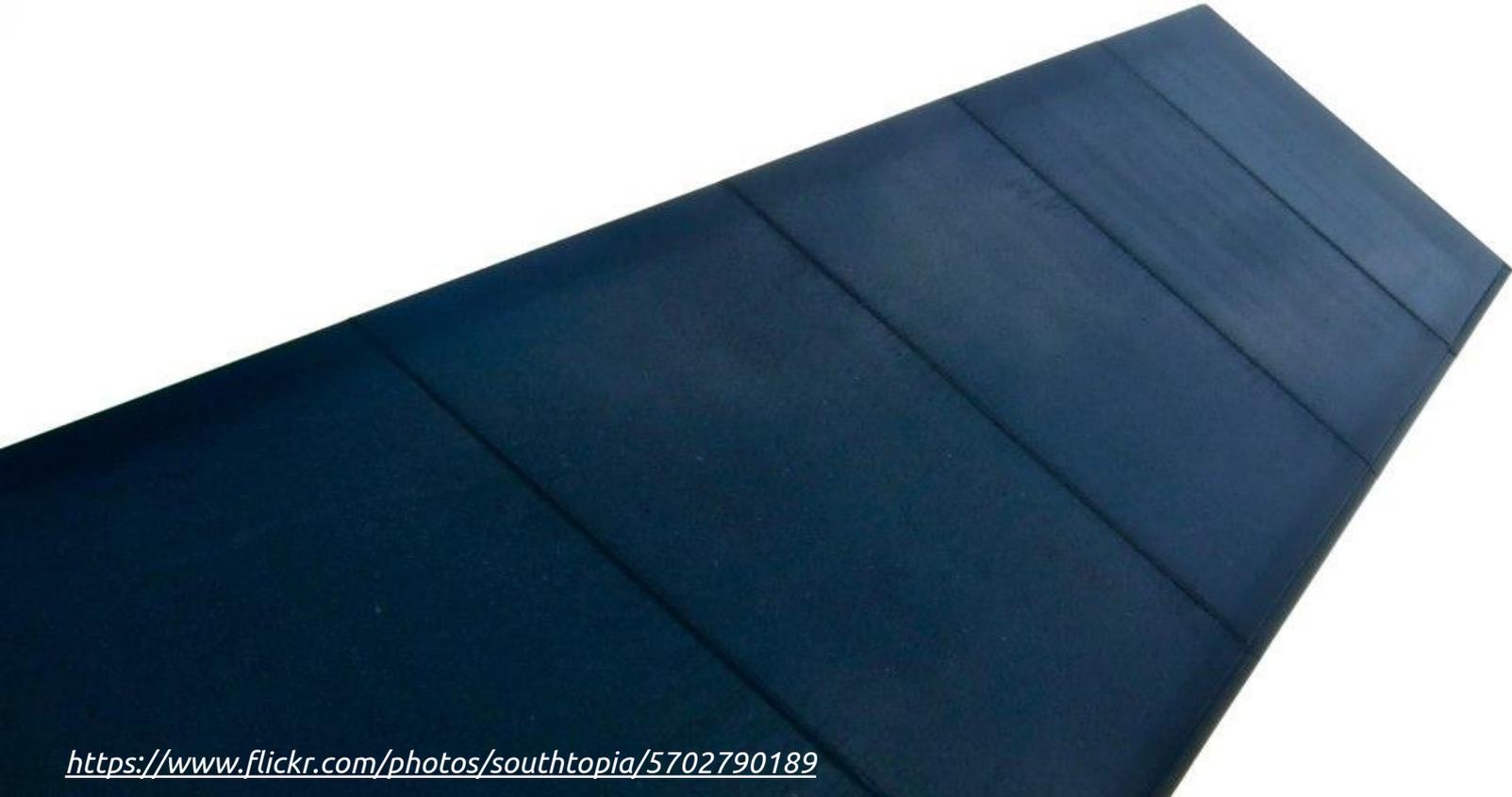
Applications: **~100**

Database: **4 TB of data**

Servers: **1400 VMs, 300 physical machines**

Locations: **Chiasso, Milan, Madrid, London, Bengaluru**

# Business: "technology is slow"

# Technology: "the monolith is the problem"

"... let's break into microservices!"

# A lot of issues

- **LONG** provisioning time

- **LACK OF** alignment across environments

- **LACK OF** alignment across applications

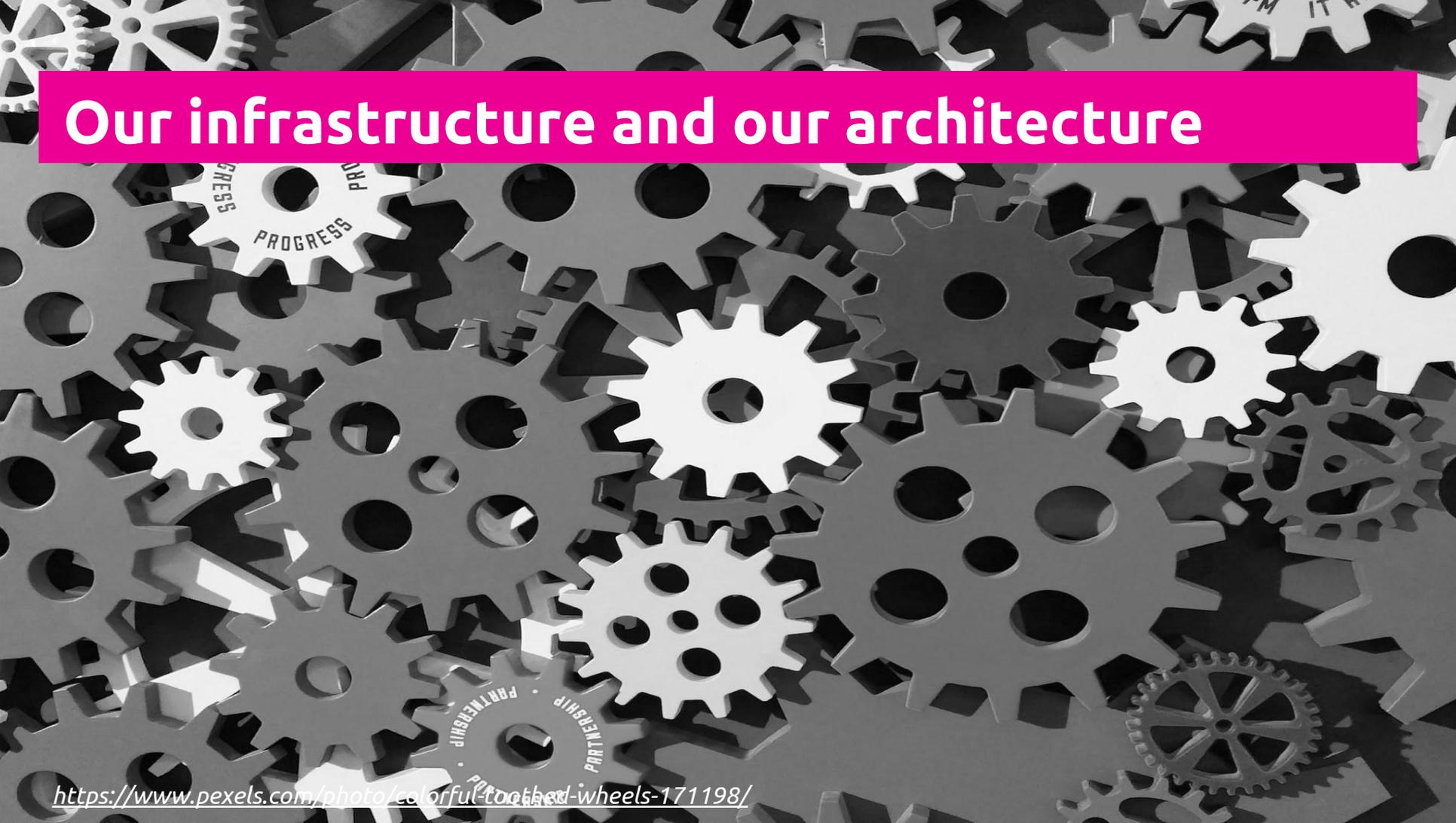- **LACK OF** awareness about *ops*

# A year-long endeavour

- build a **new, modern infrastructure**

- migrate the **search (flight/hotel) product** there

*... without:*

- impacting the business
- throwing away our whole datacenter

# Our infrastructure and our architecture

# TONS

# OF

# VIRTUAL MACHINES

# Virtualization platform

# Engage

- CoreOS, the all-in-one choice

  - Cloudconfig configuration

  - Automatable in a shot

  - Really simple patch management

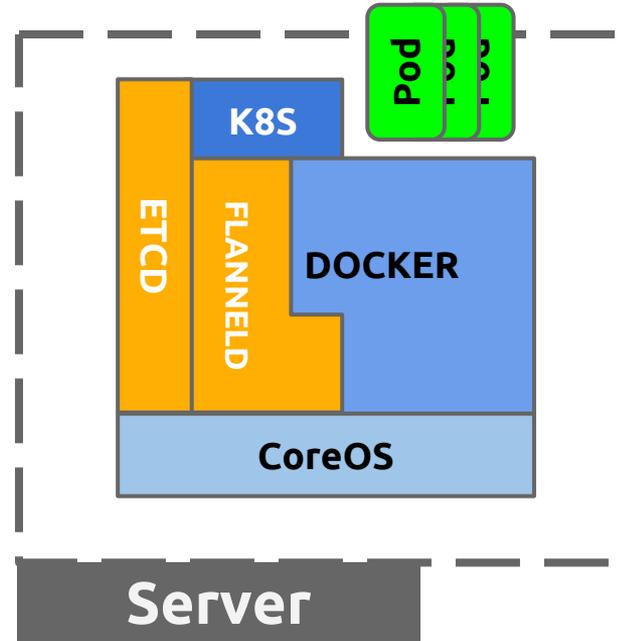# Our Kubernetes on CoreOS architecture is born

- The stack
  - ETCD
  - FLANNELD
  - DOCKER
- KUBERNETES (Google!)



ETCD FLANNELD DOCKER K8S Pod CoreOS Server

lastminute.comgroup

# How to talk with pods

# In the name of service

```
      - host: awesomeservice.prd.mykubecluster.intra
        http:
          paths:
          - path: /
            backend:
              serviceName: awesomeservice
              servicePort: 8081
```

lastminute.comgroup

# In the name of service

```
*.[prd|qa|dev].mykubecluster.intra.  IN CNAME  kubef5ingress
```

# The return of NodePort

# The registry brought another question...

lastminute.comgroup

# Seriously?

# Rear window on kubernetes



Nagios first
Grafana 4 now

Kube API

OS

collectd
image

graphite

Server

lastminute.comgroup

*icons from http://www.flaticon.com*

# We were happy!

# Not happy anymore

# Seriously?

# The change… It's a kind of magic

KEEP

CALM

and

TRUST KUBERNETES

# Lots of things!

# The final architecture (so far...)



OUTSIDE KUBERNETES

F5

ingress

ETCD | FLANNELD | K8S | DOCKER

Pod

Ubuntu

**INSIDE  KUBERNETES:**

3 different environments
7 MASTERS
2 REGISTRYs
+ 70 PHYSICAL NODES
+ 47 ETCDs
+ 7 DNS
+ 140 Namespaces
+ 1300 PODs

lastminute.comgroup

# Our infrastructure and our architecture

# Our core axioms

- **same architecture** across environments
- a **common framework** to align software
- **centralized monitoring/logging**, with alerts
- **zero downtime** deployment
- **automation** everywhere

# Kubernetes: our architecture

**production**

**APP1-PRODUCTION**

**APP1-PREVIEW**

**nonproduction**

**APP1-DEVELOPMENT**

**APP1-QA**

lastminute.comgroup

# Kubernetes: our architecture and choices

production

**APP1-PRODUCTION**

app1-production.prd.mykubecluster.intra

replica-set

deployment

POD-1  POD-2  POD-3

secret  configmap

lastminute.comgroup

# "To ingress or not to ingress? .."

**NODE 1**
NGINX

**NODE 2**
NGINX

**NODE 3**

F5

- easier DNS management
- customizable **proxy server**

- 3rd party tool
- requires **external sync**
- **all requests** go through it
- reload risks

# Kubernetes: our architecture and choices

production

APP1-PRODUCTION

POD

fluentd

*application*

collectd

carbon

lastminute.comgroup

# Monitoring and alerting: grafana/graphite

**cluster**

**APP1-PRODUCTION**

**POD**

**application**

**collectd**

**carbon**

Grafana 4

graphite

*icons from http://www.flaticon.com*

# Zero downtime (1): graceful shutdown

## deployment.yaml

```
lifecycle:
  preStop:
    exec:
      command: ["/stop_helper.sh"]
```

## stop_helper.sh

```
#!/bin/bash

wget http://localhost:8002/stop
```

# Zero downtime (2): graceful startup

**JobsExecutor.java**

```java
private CompletableFuture run(Stream<CompletableFuture> startupJobs)
 {

   return allOf(startupJobs.toArray(CompletableFuture[]::new))
     .thenAccept(this::raiseReadinessUp)
     .exceptionally(this::shutdown);

 }
```

lastminute.comgroup

# Automate everything: pipeline DSL

**pipeline**

```
microservice = factory.newDeployRequest()
    .withArtifact("com.lastminute.application1",2)
    .fromGitRepo("git.lastminute.com/team/application")

lmn_deployCanaryStrategy(microservice,"qa")
lmn_deployCanaryStrategy(microservice,"preview")
lmn_deployCanaryStrategy(microservice,"production")
```
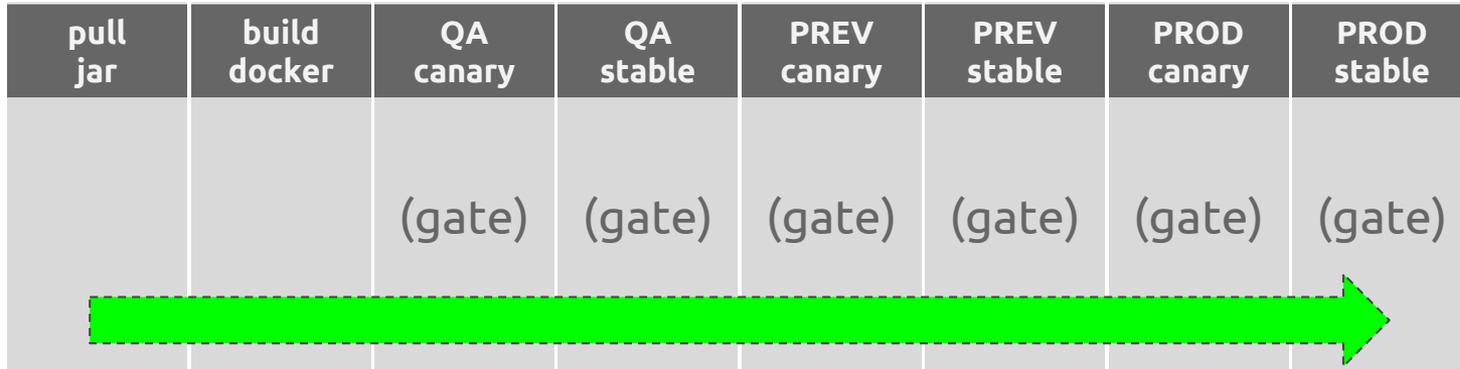
lastminute.comgroup
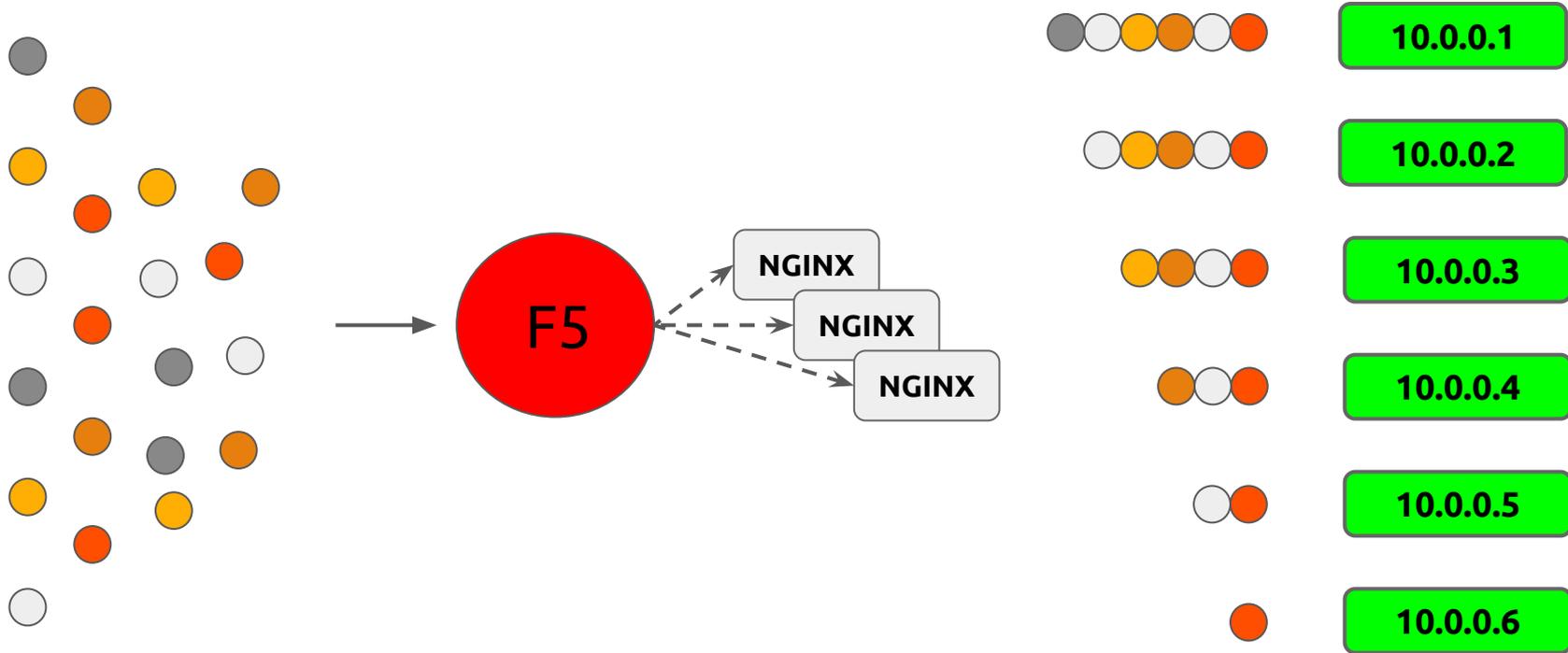
# Automate everything: pipeline

- git push
  - continuous integration
  - **continuous delivery**

| pull jar | build docker | QA canary | QA stable | PREV canary | PREV stable | PROD canary | PROD stable |
|----------|--------------|-----------|-----------|-------------|-------------|-------------|-------------|
|          |              | (gate)    | (gate)    | (gate)      | (gate)      | (gate)      | (gate)      |

# .. failure ..

# nginx ingress controller problem

# There's light .. at the end

# Give me the numbers .. again!

- **20K req/sec** in the new cluster
- **10 minutes** to create a new environment
- whole pipeline runs in **16 minutes**
  - **4 minutes** to release 100 instances of a new version
- **2M metrics/minute** flows

# Yes, we're hiring!



## THANKS

www.lastminutegroup.com