

# KUBERNETES IN A GROWN ENVIRONMENT AND INTEGRATION INTO CONTINUOUS DELIVERY

Stephan Fudeus, Expert Continuous Delivery  
Dr. Sascha Mühlbach, Expert Infrastructure Architect



# United Internet / 1&1 Mail & Media



## United Internet

- Is a leading European internet specialist
- > 9000 employees
- 90k servers in 10 data centers
- Access
  - DSL and Mobile
- Applications
  - Business (Server, Hosting)
  - Consumer (WEB.DE etc.)

## 1&1 Mail & Media

**GMX**



**mail.com**

- Main brands GMX, WEB.DE and MAIL.COM
- Various services around a free or paid mail account (calendar, news portal, cloud storage)
- 33 million active users / month

## Speakers



- **Stephan Fudeus**
  - Joined 1&1 in 2005
  - Long-term experience in building highly scalable multi-tenant applications
  - Product Owner and TechLead for our Kubernetes Clusters
  - Twitter: @der\_sfu



- **Dr. Sascha Mühlbach**
  - Expert Infrastructure Architect
  - 15 years professional experience
  - Responsible for the global operations strategy of the applications and systems infrastructure

# Agenda

- Motivation / Environment
- Cluster-Design
- Network-Setup / Ingress
- Git-driven cluster operations
- Multi-Tenancy
- Build processes
- Continuous delivery environment
- Onboarding / Training

## Motivation

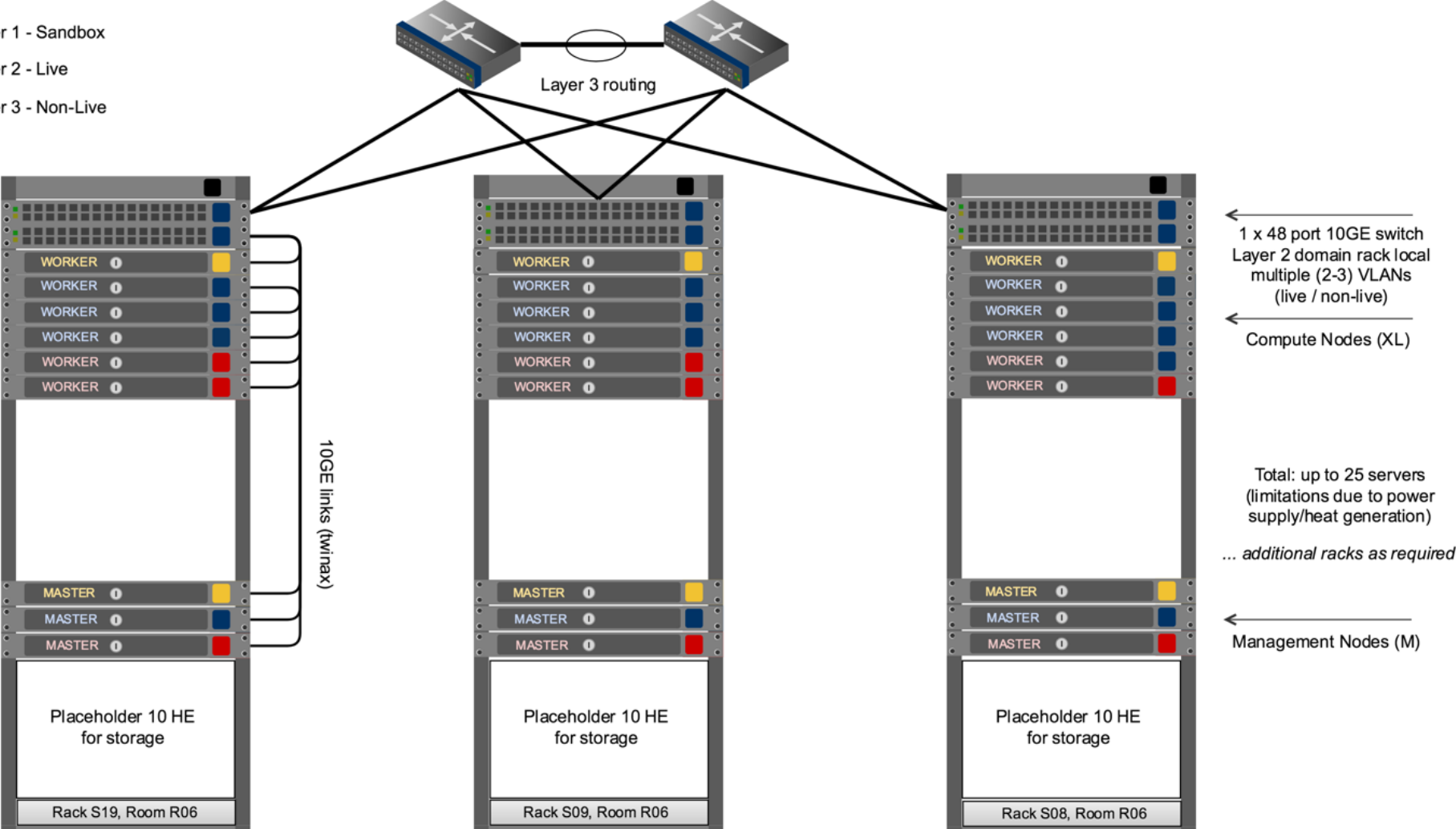
- Why Container?
  - Strong coupling between code and application runtime environment
  - One build responsibility
  - Hide core infrastructure from application
  - Reproducibility in development
  - Follow new standards in software development
- Requirements
  - Reliable platform that provides the same level of availability that our existing environment is delivering
  - Efficient deployment for geo-redundant services in multiple data centers
  - Self-service for the development and application teams
  - Multi-tenancy with strong separation for security reasons
  - Must fit into the existing network environment
  - Mostly automated operation of the base Kubernetes platform

## Environment

- Organizational Environment
  - Approx. 25 Dev Teams with 10 Ops Teams
  - Strong organizational separation between PM / Dev / Ops
  - 24/7
  - Central Ops Team to build and run the Kubernetes platform
- Technical Environment
  - 3 datacenters (2 in DE, 1 in US) that are owned by us
  - bare metal and virtual machines (KVM, ESX)
  - All servers are Puppet managed
    - Our infrastructure has ~15.000 Puppet clients
  - Majority of services are written in Java
  - Ongoing transition to CD and microservices

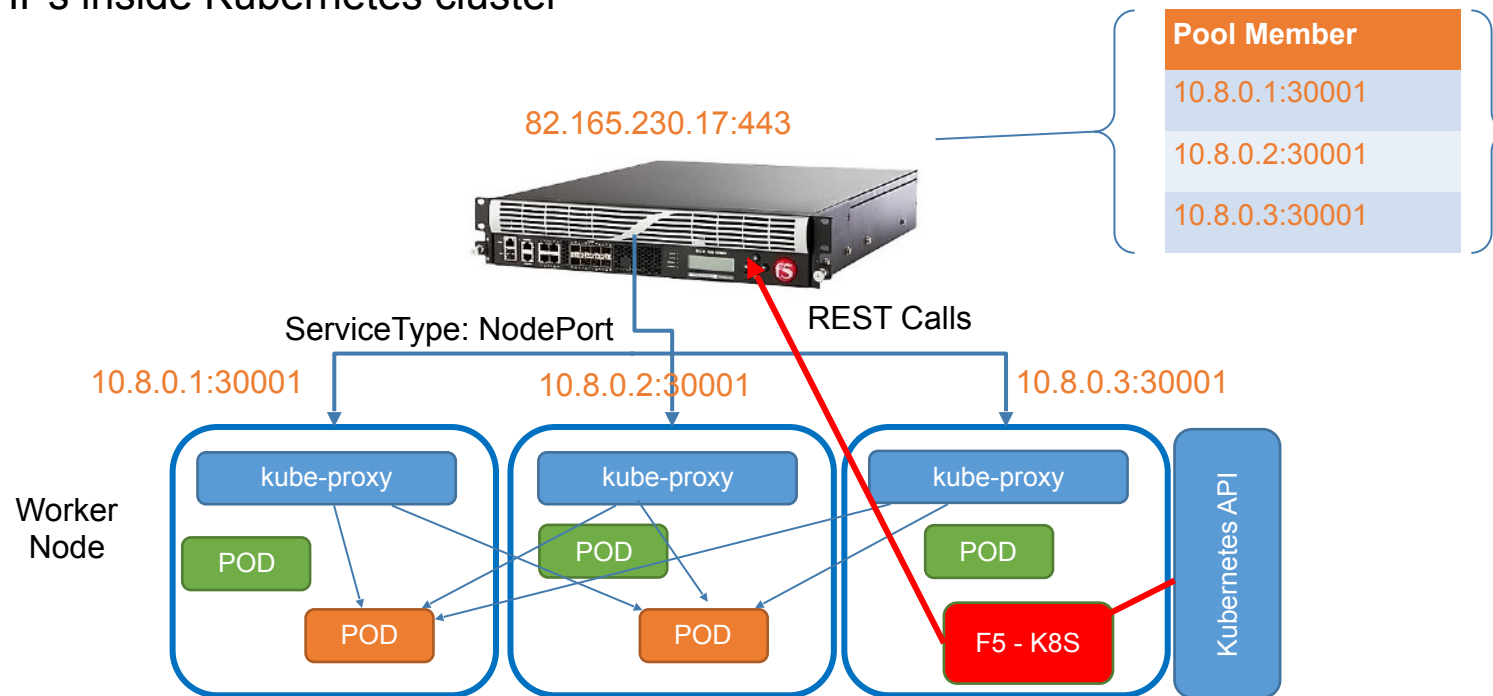
# Cluster Design

- Cluster 1 - Sandbox
- Cluster 2 - Live
- Cluster 3 - Non-Live



## Network Setup for Frontend Zone

- Integration of existing F5 BigIP load balancing platform with their features
- Service IPs are BGP-routed to Balancer and then forwarded with SNAT to NodePorts
- BGP enables global redundancy
- No public IPs inside Kubernetes cluster



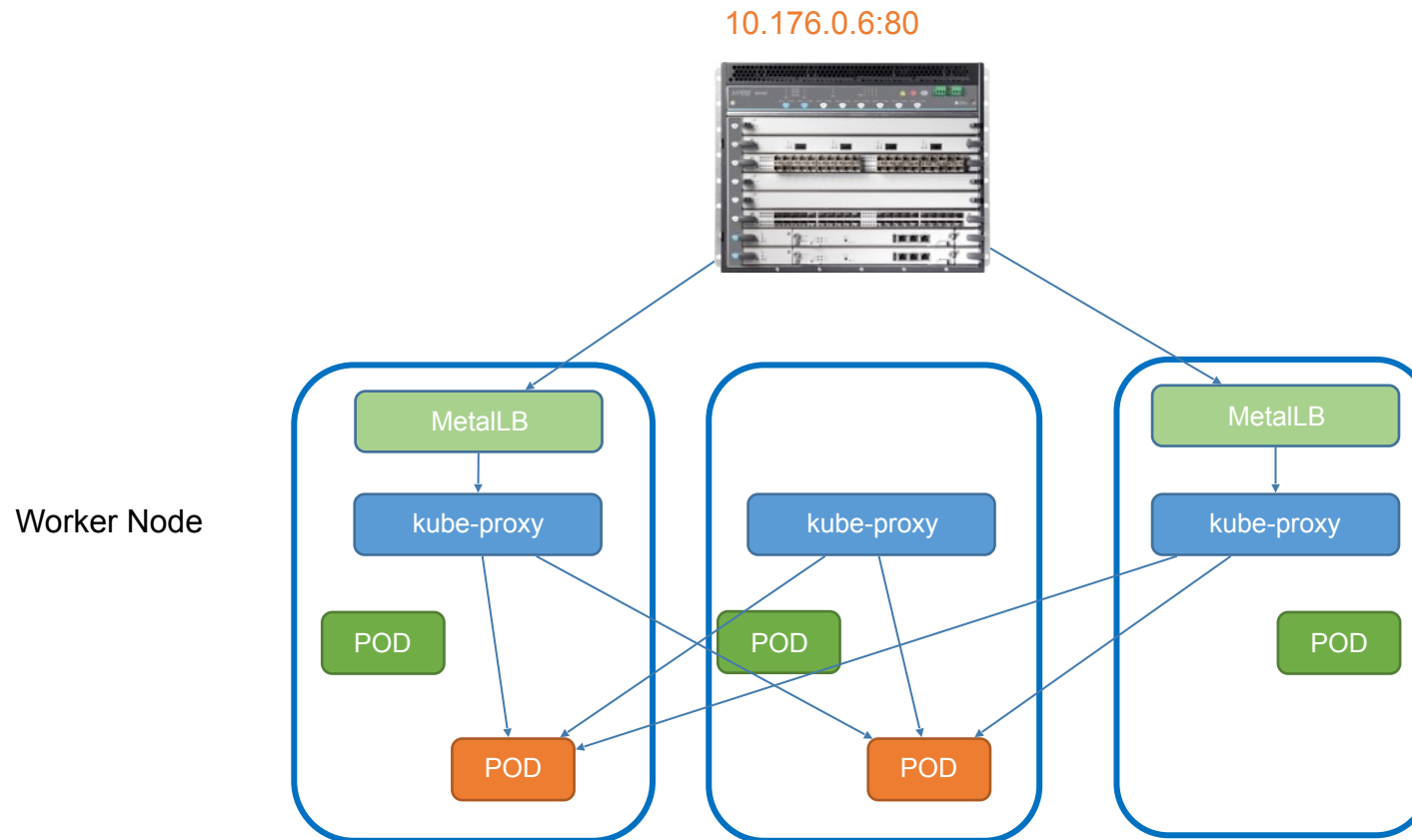


# Network configuration via ConfigMap for F5

```
1  apiVersion: v1
2  kind: ConfigMap
3  metadata:
4    labels:
5      app: prometheus
6      f5type: virtual-server
7      pipeline-managed: "true"
8  name: prometheus-ingress
9  namespace: monitoring
10 data:
11  schema: "f5schemadb://bigip-virtual-server_v0.1.7.json"
12  data: |
13    {
14      "virtualServer": {
15        "backend": {
16          "serviceName": "prometheus",
17          "servicePort": 9090
18        },
19        "frontend": {
20          "partition": "k8s-be-qa-iz1-bs",
21          "iapp": "/Common/K8S_iApp_010",
22          "iappPoolMemberTable": {
23            "name": "Pool_Definition__Pool_Members",
24            "columns": [
25              {"name": "Pool_Member_IP", "kind": "IPAddress"},
26              {"name": "Pool_Member_Port", "kind": "Port"}
27            ]
28          },
29          "iappOptions": {
30            "description": "Prometheus monitoring for Kubernetes"
31          },
32          "iappVariables": {
33            "Virtual_Definition__VS_Type": "{{ .Values.global.ingress.type }}",
34            "Virtual_Definition__VS_VLANS_Enabled": "{{ .Values.global.ingress.vlan }}",
35            "Virtual_Definition__VS_SMATPool": "{{ .Values.global.ingress.snat }}",
36            "Virtual_Definition__VS_TCP_Client_Profile": "{{ .Values.global.ingress.tcp_client_profile }}",
37            "Virtual_Definition__VS_HTTP_Profile": "{{ .Values.global.ingress.http_profile }}",
38            "Virtual_Definition__VS_SSL_Client_Profiles": "{{ .Values.global.ingress.ssl_client_profiles }}",
39            "Virtual_Definition__VS_FastL4_Client_Profile": "{{ .Values.global.ingress.fastl4_client_profile }}",
40            "Virtual_Definition__VS_Persistence": "{{ .Values.global.ingress.stickyness }}",
41            "Virtual_Definition__VS_SSL_Server_Profiles": "{{ .Values.global.ingress.ssl_server_profiles }}",
42            "Virtual_Definition__VS_TCP_Server_Profile": "{{ .Values.global.ingress.tcp_server_profile }}",
43            "Virtual_Definition__VS_iRules": "{{ .Values.global.ingress.irules }}",
44            "Pool_Definition__Pool_Health_Monitors": "{{ .Values.global.ingress.healthmonitor }}",
45            "Pool_Definition__Pool_Load_Balancing_Method": "{{ .Values.global.ingress.loadbalancing_method }}",
46            "Virtual_Definition__VS_IP": "{{ .Values.global.ingress.vip }}",
47            "Virtual_Definition__VS_Port": "{{ .Values.global.ingress.port }}"
48          }
49        }
50      }
51    }
```

## Network Setup for Backend Zone

- In backend networks, we use MetalLB (no specific Layer 7 requirements)
- Service IPs are BGP-announced with ECMP distribution (easy scaling)
- LoadBalancing only with K8S base algorithms or ingress controller features



# Network configuration via Service for MetalLB

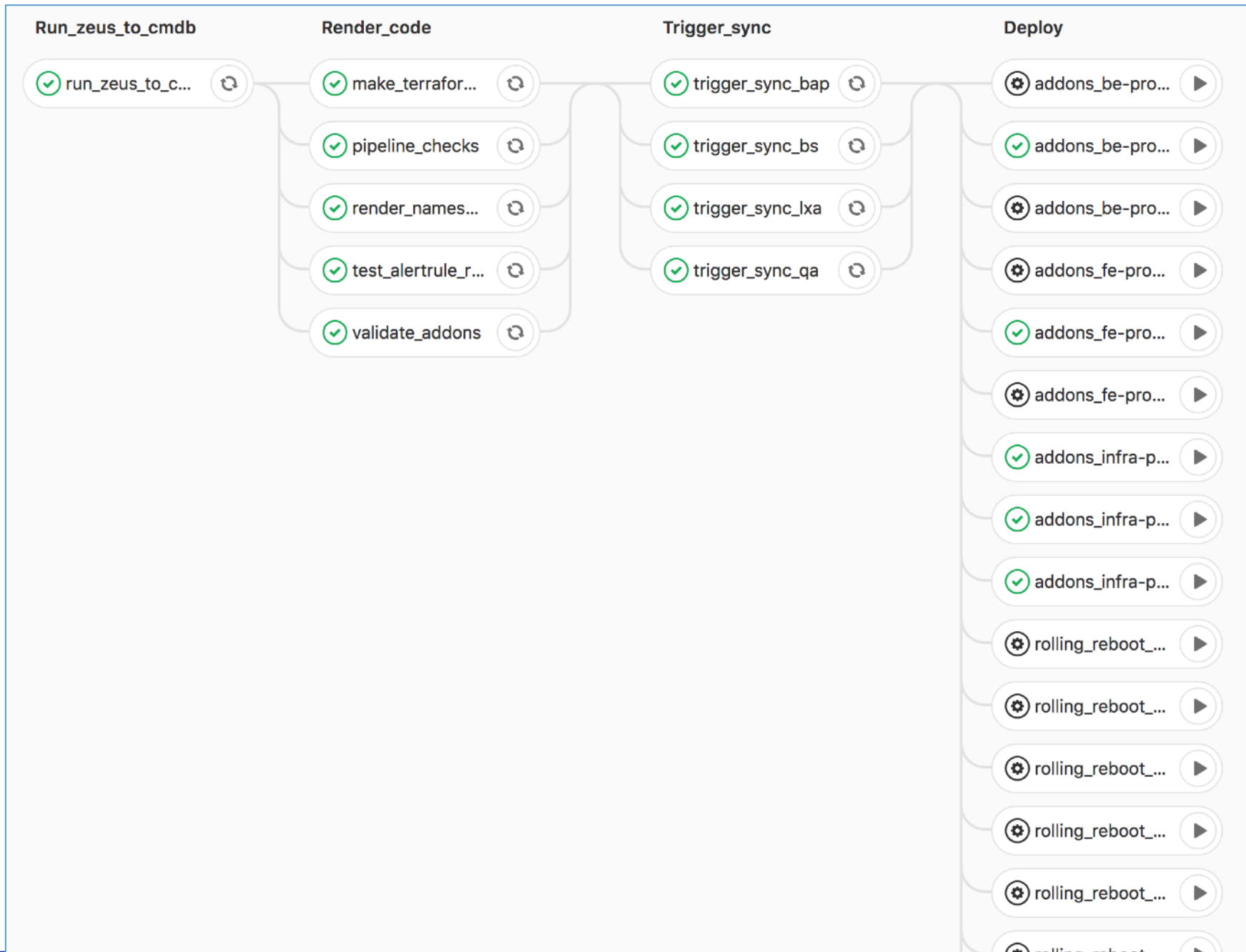
```
1  apiVersion: v1
2  kind: Service
3  metadata:
4    labels:
5      name: prometheus
6  name: prometheus
7  namespace: monitoring
8  spec:
9    clusterIP: 100.72.112.153
10   externalTrafficPolicy: ClusterIP
11   type: LoadBalancer
12   loadBalancerIP: 10.176.0.6
13
14   - name: prometheus
15     port: 80
16     protocol: TCP
17     targetPort: 9090
18  selector:
19    app: prometheus
20  sessionAffinity: ClientIP
21  sessionAffinityConfig:
22    clientIP:
23      timeoutSeconds: 10800
```

## Git-driven cluster operations

- Maturity level via 3 branches (master, integration, production)
- All cluster operations are triggered based on Gitlab-CI pipelines
  - automatically on git-pushes to relevant branch
  - manually triggered jobs for cluster changes
  - scheduled jobs for periodic changes (namespace updates / purges)

## Git-driven operations use cases

- Full redeployment of clusters
  - Only if cluster is broken, will wipe everything
  - Will redeploy all nodes in parallel
- Rolling upgrade of clusters
  - Usually done on a weekly basis
  - Will wipe and reset nodes one by one
- Namespaces update
  - Nightly updates for production
  - on-push for integration
- Addon update
  - Addons as helm charts, rendered via helm template and injected via kubectl apply
  - Done ad-hoc for addon changes without redeployment



## Multi Tenancy

- Common platform for several teams
  - PodSecurityPolicies (no-root, no host-net, r/o layers)
  - Dedicated resources for teams
    - Dedicated in-cluster prometheus for scraping
    - Configurable log-sink (Elasticsearch, Kafka)
  - Authentication via OIDC <-> Dex <-> LDAP
- Maximum separation between teams targeted
  - Namespaces are a „managed“ resource
  - Resource constraints defined centrally per namespace
  - Users are restricted to their namespaces via RBAC
  - Network policies
  - Team-centric „helper“ namespace
    - e.g. \$team-helper
    - Used for managed resources, e.g. team-prometheus
  - Individual namespaces per (group of) application and stage
    - \$team-\$app-live, \$team-\$app-prelive

## Multi Tenancy

- Dedicated namespaces for individuals
  - Purpose: Training, PoC, Experiments
  - Daily process to read users from LDAP and generate and flush namespaces
  - Service exposure via central ingress controller (traefik)



# Namespace-Config via yaml

```
1 .qa-ns: &qa-default-resources
2   cpu_total: 20
3   memory_total: 130Gi
4   max_pod_size: L
5 .dev-ns: &dev-default-resources
6   cpu_total: 10
7   memory_total: 65Gi
8   max_pod_size: S
9 team:
10  ams:
11   log_sink:
12     address: "kafka1,kafka2,kafka3"
13     type: "kafka"
14   prometheus_vip: "10.176.0.7"
15   app_namespaces:
16     ams-tooling-qa: *qa-default-resources
17     ams-testing-qa: *qa-default-resources
18     ams-testing-dev: *dev-default-resources
19   admins:
20     - "ams-admins"
21     - "ams-users"
22 personal_namespaces:
23   - cn=ams-users
24   - cn=ams-admin
25   - cn=other-users
26 cluster_name: "sandbox-intg-iz2-bap"
27 mam_dc: "bap"
28 activate_network_policies: true
```

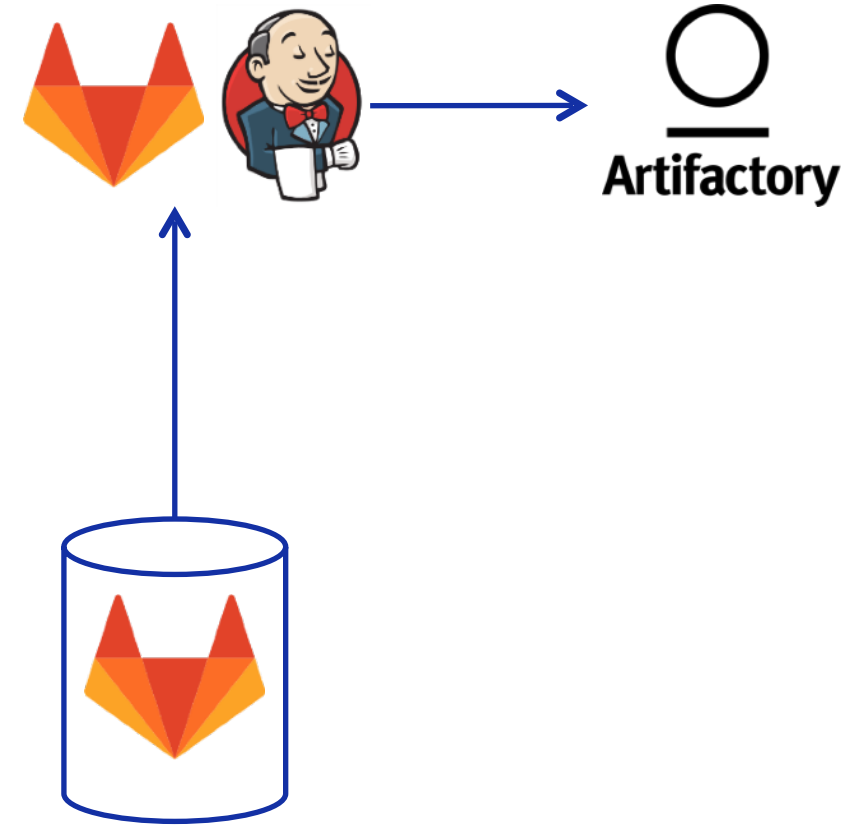
Rendered via helm

36 resulting manifests:

- 1 kind: Deployment
- 1 kind: Ingress
- 1 kind: Service
- 2 kind: ConfigMap
- 2 kind: ServiceAccount
- 4 kind: LimitRange
- 4 kind: Namespace
- 4 kind: ResourceQuota
- 8 kind: NetworkPolicy
- 9 kind: RoleBinding

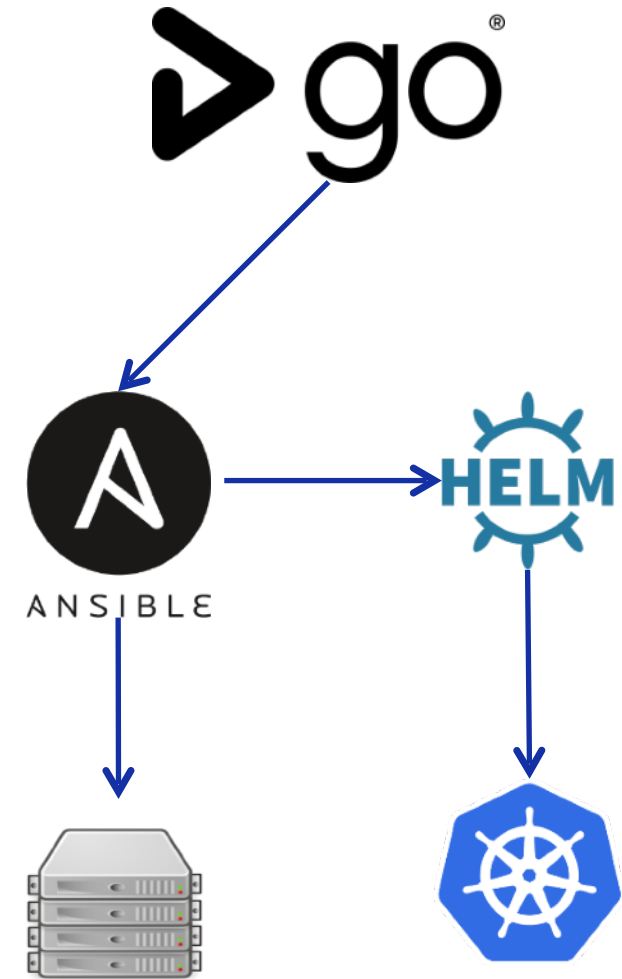
## Build Processes

- Fully automated builds
- High degree of standardization
  - e.g. central maven POM
- Parallel builds for classical and container deployments
  - Containers use a centrally provided base image
- Build processes are triggered upon base image changes
- Policy: updates / rebuilds are enforced every 4 weeks

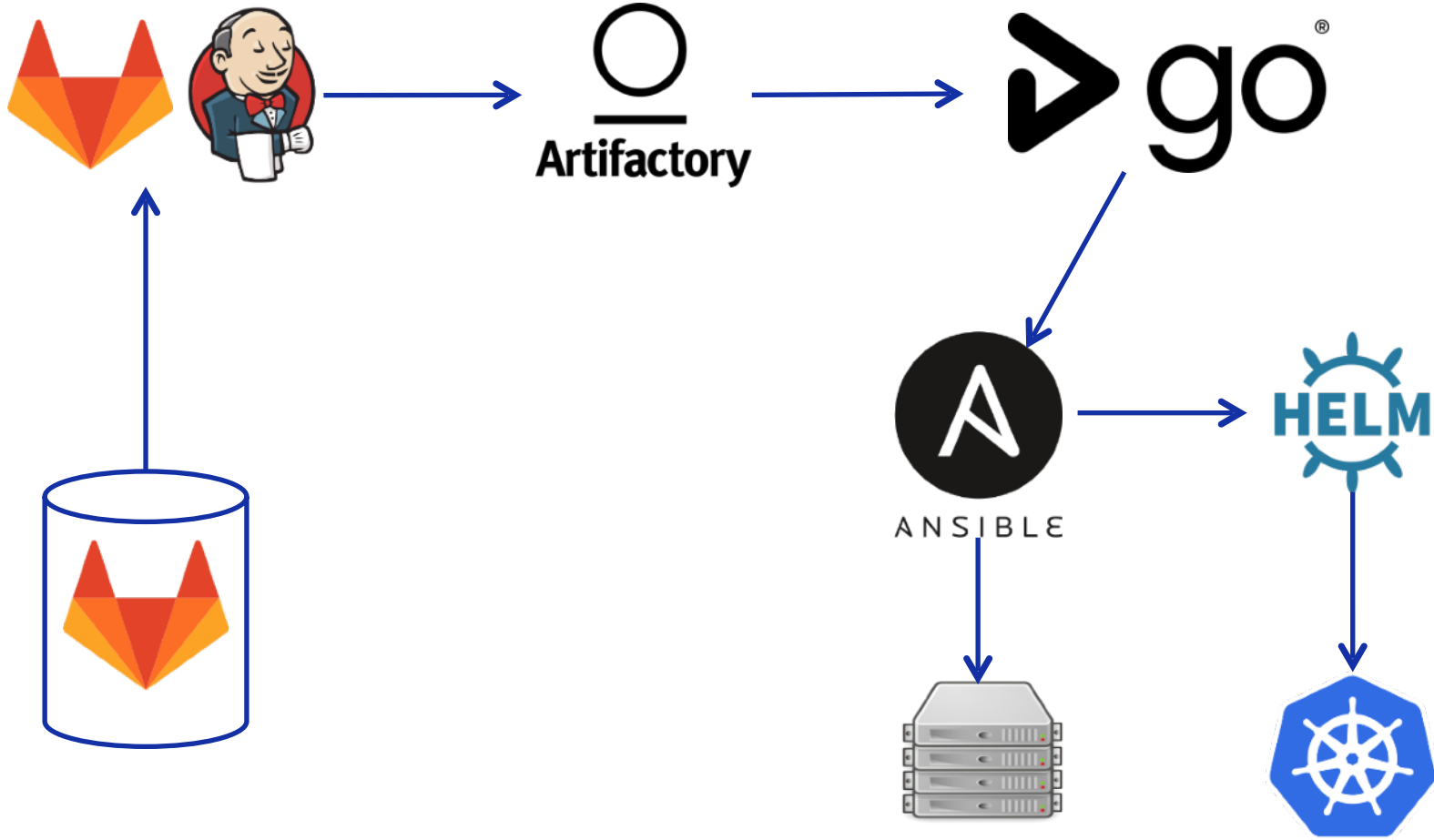


## Continuous Delivery Environment

- GoCD maps business processes
  - Dedicated instance per team
  - Standardized pipeline templates
- Technical processes are mapped separately
  - Ansible for host based deployments
  - Helm/Kubectl for k8s deployments
- Supports hybrid deployments
  - Container and Hosts in parallel
  - Hybrid usage via loadbalancer
  - Assists during transition phase



# Fully automated deployment chain



## Onboarding & Training

- 4 training blocks for system administrators (1-2 days each)
  - Docker & Kubernetes
  - GoCD & Helm
    - Pipeline Design
    - Helm Templating
  - Development Techniques for Ops
    - Repositories and versioning
    - Secure Software Development Lifecycle
  - Operating Container Applications
    - Monitoring, Logging and Failure Handling
    - Operations Lifecycle

## Links

- F5-Ctrl (<https://github.com/F5Networks/k8s-bigip-ctrl>)
  - MetalLB (<https://metallb.universe.tf/>)
  - Dex (<https://github.com/coreos/dex>)
  - GoCD (<https://www.gocd.org>)
- 
- <https://jobs.1und1.de/>
  - <https://web.de>
  - <https://www.gmx.net>
  - <https://www.mail.com>
  - <https://www.united-internet.de/>