

(Kubernetes - k8s)

## Index

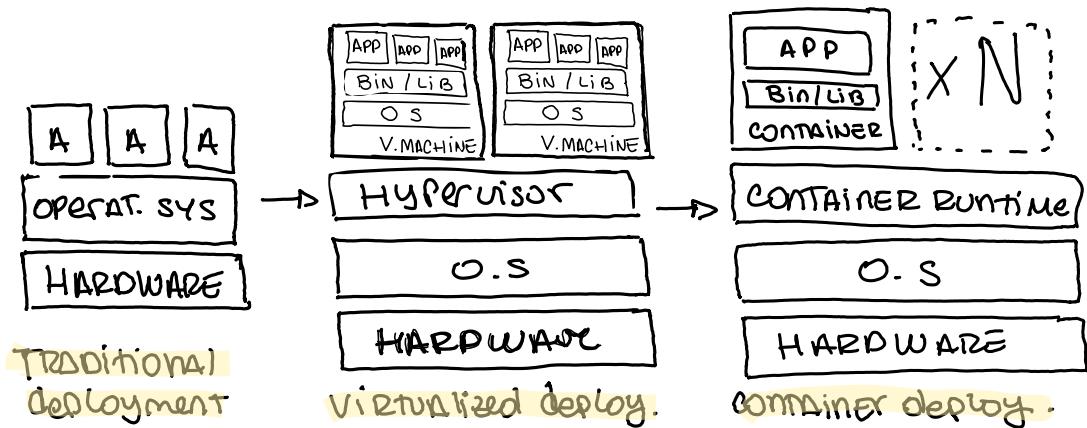
- K-0.0 → What's Kubernetes?
- K-1.0 → Components
- K-2.0 → K8S Objects
  - | K-2.1 → Objects Management
  - | K-2.2 → Names
  - | K-2.3 → Namespaces
  - | K-2.4 → Label and Selectors
- K-3.0 → Architecture.
- K-4.0 → PODs
- K-5.0 → CONTROLLERS.

by Lucas Coutureos

# Kubernetes (K8s)

What's K8s: Open-source container orchestration system for containers. It automates deployment, scaling & management.

## Containerization 101



What does K8s Provide: You need to know what it provides if you wish to know where to aim...

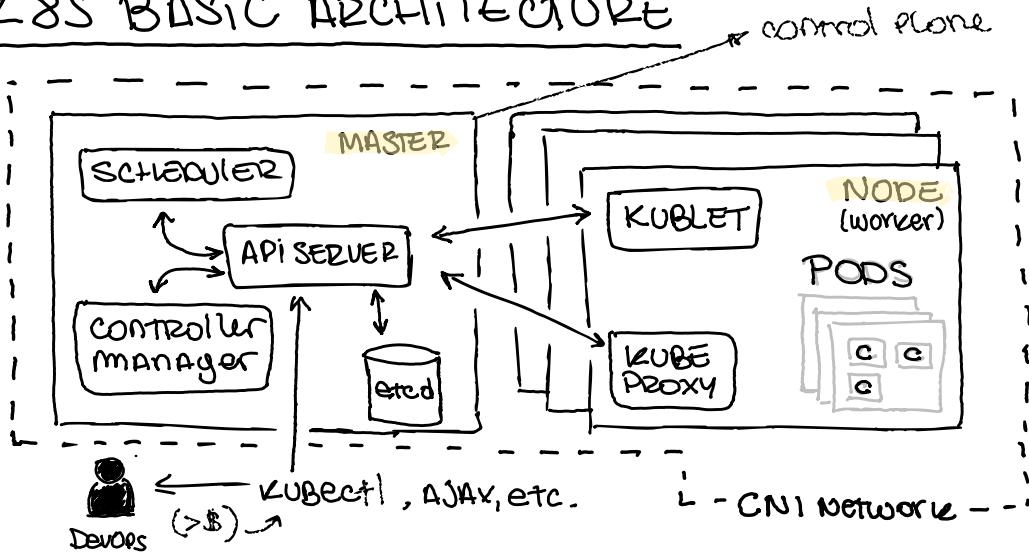
- Service discovery & load balancing.
- Storage orchestration
- Automated rollouts & rollbacks → # SEMANTIC
- Automated bin packing
- Self healing
- Secret and config. management.

It's important to know what does K8s provides and what it does not, this lets you foresee any possible problem you think K8s will solve when it actually doesn't. Besides you might find that already addressed issues will reappear since you are changing the whole ARCH.

# Components

Before describing each component we'll take a look to the k8s architecture.

## K8S BASIC ARCHITECTURE



## MASTER COMPONENTS

**Kube-api server**: is the front-end for the k8s control plane. It's design to scale horizontally (by deploying more instances).

**Etc**: consistent and highly-available key value store used as k8s' backing store for all the cluster data.

**Kube-scheduler**: watches newly created PODS that have no node assigned and selects a node for them to run.

**Kube-controller-manager**: runs controllers (a control loop that watches the state of the cluster)

Controllers: (they all run under the same process)

- NODE
- ENDPOINTS
- REPLICATION
- SERVICE ACCOUNT & TOKENS

## NODE COMPONENTS

**Kubelet:** is the k8s agent that runs on each node. It makes sure that containers are running in the pods but it doesn't interfere with other containers.

**Kube-proxy:** network proxy that runs on each node in the cluster implementing part of the k8s Service concept. allows network communication to your pods from inside or outside the cluster.

**Container-runtime:** software responsible for running the containers (i.e Docker). K8S supports several container runtimes.

## The Kubernetes API (Introducing kubectl)

This is where we are going to manage our cluster. We are going to do this mostly from a command-line tool called **kubectl** with which we'll be able to create, update, delete and get our API objects. By making calls to the API

## Kubernetes Objects

K-2.0

Just for you not to panic and to get familiar with k8s obj, these are the guys that you'll be able to express in a **.yaml** file

Kubernetes Objects are persistent entities in the k8s system. These entities represent the state of your cluster. Specifically:

- + What containerized apps you are running (and on which nodes)
- + The resources available to those applications.
- + The policies around how those apps behaves (i.e. restart...)

A K8S object is a "record of intent" - once the object is created, K8S will always ensure that objects fit. By creating an object you are telling K8S what you want your cluster to look like.

To CRUD objects in K8S we use the K8S API, as we mentioned early the kubectl cli makes those API calls for us.

(POD properties?)

## Object Specs n' Status

Every object includes two nested "object fields" that governs the object configuration: Specs and Status

The SPEC, which you must provide describe the desire state for the object.

The STATUS describes the actual state of the object. This last is supplied and updated by Kubernetes.

#> kubectl api-resources → List all kubernetes objects

When creating an object you send the specs in a JSON format.  
(Most often you provide this information to kubectl in a .yaml file.)

ApiVersion: apps/v1

. /deploy.yaml

Kind: Deployment

#> is kubectl api-resources

metadata:

name: nginx-deployment

→ .yaml example.

spec:

selector:

matchLabels:

APP:nginx

replicas:

# Kubernetes Object Management.

There're several management techniques and the kubectl supports different way to create and manage k8s objects.

**#WARNING:** A k8s object should be managed using only one technique. Mixing and matching technique for the same object may result in undefined behavior.

## A WONDERFUL TABLE I CAME THROUGH

MANAGEMENT TECHN	OPERATES ON	RECC. ENV	SW	LEARNING CURVE
Imperative commands	Live Objects	Develop.	1+	Lowest.
Imperative object-conf.	Individual files	Production	1	Moderate
Decorative obj. conf.	Directories	Production	1+	Highest

### Imperative commands.

When using imperatives commands the user operates directly on live obj. of a cluster passing option to the kubectl command as arguments or flags.

It provides no history of previous configuration.

Examples: \$> kubectl run nginx --image nginx  
\$> kubectl create deployment nginx --image nginx.

### TRADE-OFFS: (compared 2 Imp/Dec obj configuration)

#### PROs

- Simple, easy to learn
- Single step to make changes to the cluster

#### CONS

- no template for new obj.
- no integration with change review process
- no source of record
- no audit trail

## Imperative object configuration

In imperatives object conf. commands you specify the operation (create, replace, etc), optional flags and at least one filename. The file must contain a full def. of the object. (YAML or JSON).

**#WARNING:** the "replace" command replaces the file specs with write the new one, dropping all changes to the object missing from the conf. file. This approach should not be used with resources whose specs are update independently from the configuration file.

Examples:

kubectl create -f nginx.yaml

kubectl delete -f nginx.yaml -f redis.yaml

kubectl replace -f nginx.yaml

## TRADE-OFFS:

### PROs (comp. with Imp. obj)

- + can be stored in Git
- + integrated with reviewing processes.
- + provides a template.

### CONS (comp. with Imp. obj)

- requires understand. obj. sem.
- you have to write a YAML

### PROs (comp. with Dec. obj. conf.)

- + is simpler and easier to understand
- + as of k8s 1.8, Imp. obj is more mature.

### CONS (comp. with Dec. obj. conf.)

- doesn't work well on files
- updates to live obj. must be reflected on conf. files or they will be lost on the next replacement.

## Declarative Object Configuration.

When using DOC, the user operates on conf. files stored locally, however the user does not define the operations to be taken on the files. Create, Update and Delete are automatically detected per-object by kubectl. This enables working on directories where # operations might be needed for # objects.

### Example:

```
$> kubectl diff -f configs/ (1)
$> kubectl apply -f configs/ (2)
```

- (2) Process all conf. files in the configs directory.
- (1) Check changes by applying.

## TRADE-OFFS (compared to Imp. Object Configuration)

### PROs

- + Changes to live obj. are required even if not merged to the conf. file.
- + Better support for directories and auto-detecting oper. types.

### COWs

- Harder to debug and understand when results are unexpected.
- Partial updates using diff create complex merge and patch operations

## Names and UIDs

K-2.20

All objects are unambiguously identified by a Name and an UID

Each object has an unique Name and it's a string specified under the tag "metadata:" in the yaml file.

# Check Conventions. (max 253 char, lowercase + "-" and ".")

UID are unique strings assigned and managed by k8s

## NAMESPACES (Virtual Clusters)

NAMESPACES is how K8S manage to distinguish between multiple virtual clusters. These virtual clusters ARE called NAMESPACES

NAMESPACES are intended for USE in env. with many users spread across multiple teams or projects. for clusters with a few to tens of users, you should not need to create OR think about NAMESPACES.

- NAMESPACES provide a scope for Names (K-2.2)
- NAMESPACES can not be nested inside one another
- EACH K8S RESOURCE can be in only ONE NAMESPACE

## Default NAMESPACES

default → where resources with no NS. go.

kube-system → NS for k8s resources.

kube-public → readable by all users. mostly for cluster usage

## Labels and Selectors

K-2.4

Labels are key/value pairs attached to obj. below the metadata tag. They are used to organize and select a subset of objects.

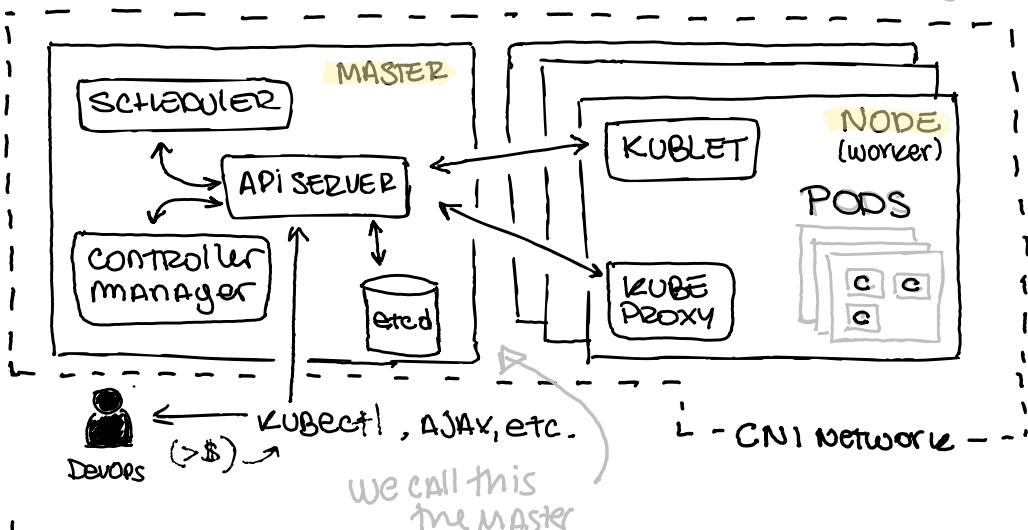
# Unlike name and UID label are not unique

- There are two types of Selector: Equality-based and Set-based
- both labels and selectors are specified on the YAML file or via command line. (with the flag -l)

(for + info on k8s please read about annotations, field selectors)

# Kubernetes Architecture.

## K8S BASIC ARCHITECTURE



Nodes: A node is a worker machine.

A node may be a VM or a bare metal machine.

Each node contains the services necessary to run Pods and is managed by the master components.

Services in a node (worker) includes

- container runtime - }
- kube-proxy }
- kubelet } Previously defined components.

Node Status A node status contains the following info.

- Addresses (Hostname, Internal IP, External IP)
- Conditions (Out of Disk, Ready, MemoryPressure, PIDPressure, etc.)
- Capacity and Allocatable (Describes CPU, Mem & max Pods available)
- Info. (General info about the node)

# Management

Unlike pods and services nodes are not created by k8s. They are created by external cloud providers or they are in your pool of physical or virtual machines.

Nevertheless, what k8s does is, it creates a node (internally) and then checks if that node exists. like.

```
{ "kind": "Node",
  "apiVersion": "v1",
  "metadata": {
    "name": "10.240.79.167",
    ...
  }
}
```

k8s will check if the node is based on the [metadata.name]

# k8s will keep the node created as "invalid" and will keep on performing healthchecks on the node until it becomes valid.

Node controller. The node controller is a k8s master component which manages various aspects of nodes. It has 3 main roles on a node life.

1. Assign a CIDR block to the node when it is registered.
2. Keep the node controller's list of nodes up to date with the cloud provider's list of available machines. Whenever a node is unhealthy is the node-controller the one who asks the cloud provider if the VM is still available.
3. Monitoring node's health - it updates the NodeReady condition on the NodeStatus

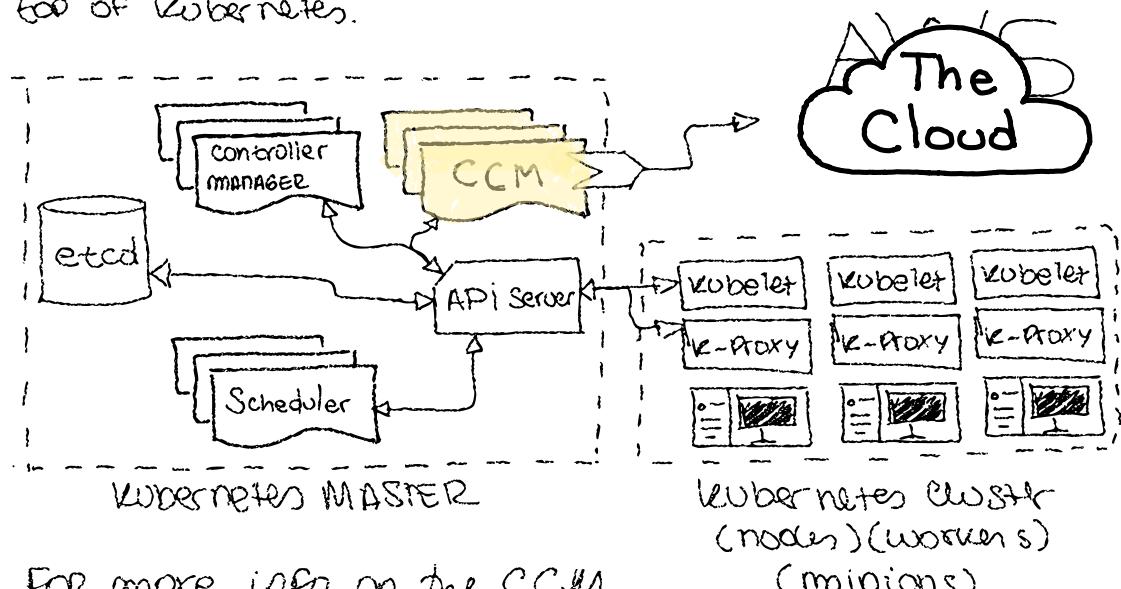
Self-registration of Nodes: (Nodes can be manually admin.)

The kubelet comp. will register itself with the API unless the flag --register-node is set to FALSE explicitly

## The Cloud controller Manager

The **CCM** was originally created to allow specific vendor code and the k8s core to evolve independent of one another.

The **CCM** runs alongside the Master components such as the controller manager, API server and scheduler. It can also be started as a add-on in which case it runs on top of Kubernetes.



For more info on the CCM  
Please refer to the k8s site.

## Vendor implementations of the CCM

- AWS
- Azure
- Oracle :/
- Digital Ocean
- GCP
- Linode
- Baidu Cloud
- Google Cloud ?

# POD

"THE SMALLEST DEPLOYABLE OBJ!"

A POD: IS A GROUP OF ONE OR MORE CONTAINERS.

- A POD IS THE BASIC EXECUTION UNIT OF K8S APPS.
- A POD REPRESENTS PROCESSES RUNNING ON YOUR CLUSTER.
- A POD ENCAPSULATES AN APPLICATION'S CONTAINER(S) OR MORE STORAGE RESOURCES, A UNIQUE NETWORK IP AND OPTIONS ON HOW THOSE CONTAINERS SHOULD RUN.

TWO MAIN WAY OF USAGE.

**RUN A SINGLE CONTAINER:** THE "ONE-CONTAINER-PER-POD" MODEL IS THE MOST COMMON K8S USE CASE. YOU CAN THINK A POD AS A WRAPPER FOR A SINGLE CONTAINER.

**RUN MULTIPLE CONTAINERS THAT NEEDS TO WORK TOGETHER:** A POD MIGHT CONTAIN AN APP. THAT CONSIST OF MULTIPLE CO-LOCATED CONTAINERS THAT ARE TIGHTLY COUPLED AND NEED TO SHARE RESOURCES.

Each POD IS MEANT TO RUN A SINGLE INSTANCE OF AN APP. IF YOU WISH TO SCALE HORIZONTALLY YOU WOULD LIKE TO HAVE SEVERAL PODS. IN K8S, THESE ARE KNOWN AS **REPLICAS**.

PODS Share resources → Networking.  
→ Storage.

# PODS ARE DESIGNED AS RELATIVELY EPHEMERAL, DISPOSABLE ENTITIES.

# PODS DO NOT RUN BY THEMSELF, IS AN ENVIRONMENT WHERE THE CONTAINERS RUN IN AND PERSIST.

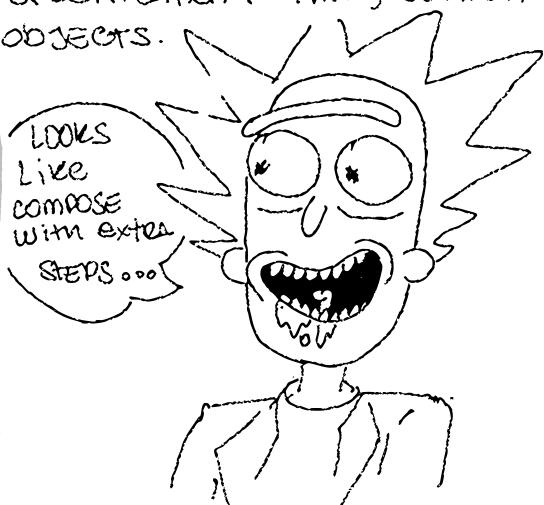
## Pods And Controllers.

A controller can create and manage multiple pods for you. For example, if a node fails, the controller can replace the pod by scheduling an identical one on a different node.

**Deployments, StatefulSet and DaemonSet** are some example of controllers that contain one or more pods.

**Pod templates.** Are pod specification (.yaml) which are included in another objects.

```
apiVersion: v1
kind: Pod
metadata:
  name: myAPP-Pod
spec:
  containers:
    - name: myAPP-container
      image: busybox
      command: ["sh -c"]
      /Pod.yaml
```



FOR FURTHER READING ON PODS I DO SUGGEST YOU TO READ ABOUT:

- **POD CONDITIONS** → this could help.
- **RESTART POLICY** → THESE ARE A MUST
- **CONTAINER STATES**
- **INIT CONTAINERS**
- **POD PRESSETS** → this is interesting

# Controllers

A controller is a control loop that ensures the desired state of the cluster (PODs) is the initially specified.

As we mentioned before controllers are responsible for PODs. In this section we are going to list some main controllers and how they work.

## Replica Sets (before learning controllers we need to understand Replica Sets)

A ReplicaSet purpose is to maintain a stable set of replica PODs at a given time. Is often used to guarantee the availability of a specific number of identical PODs.

A ReplicaSet ensures the availability of PODs. However, a Deployment is a higher-level concept that manages ReplicaSets and do more. K8S recommends to use Deployments instead of using ReplicaSets directly unless you require custom update orchestration or don't require updates at all.

## Deployments

#NOTE: A "Deployment" that configures a "ReplicaSet" is now the way to set up replication.

K-2.1

A Deployment controller provides declarative updates for PODs and ReplicaSets.

You can use Deployments to describe a desire state, create new Replica Set or remove existing Deployments and adopt all its resources.

I SSSeee no Python...sss.

## Use cases for deployments:

- + **rollout a replicaset** → There's a gtx example on the k8s site to get started with.
- + declare new state of the PODs
- + scale up a deployment
- + pause the deployment
- + clean Replica Sets.

## Statefull Set

Statefull Set is the workload API Obj used to manage stateful applications.

Like Deployments, StatefullSets manage Pods that are based on identical container specs. Unlike Deployment, it maintains a sticky identity for each of their pods.

These pods are created from the same specs, but are not interchangeable: each has a unique identifier that maintains across any rescheduling.

For further reading on Statefull Sets please consider reading following concepts from the k8s site:

- StatefullSet usage (using StatefullSet)
- Statefull Set Limitations (Limitations)
- Pod Identity.



## Daemon Sets

A DaemonSet ensures that all (or some) nodes run a copy of a Pod. As nodes are added to the cluster, pods are added. As nodes are removed from the cluster the pods are garbage collected.

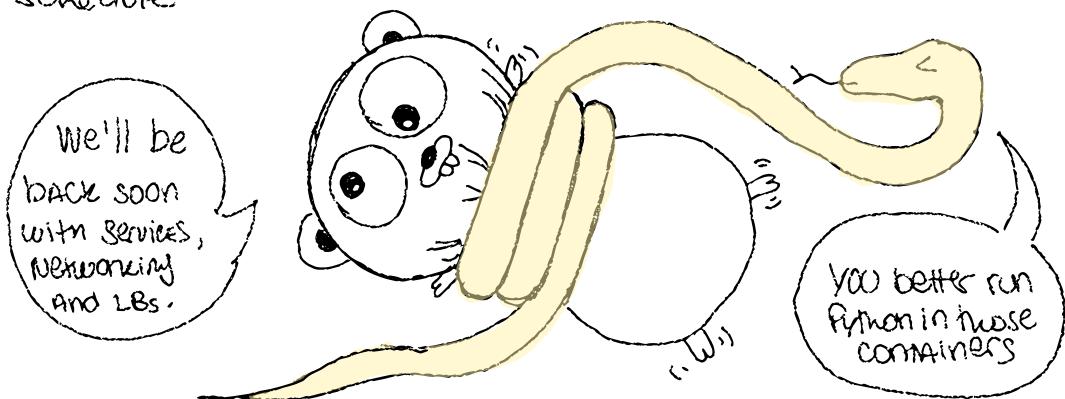
Use cases are.

- + Running a cluster storage daemon, such as **glusterd** or **ceph** on each node.
- + Running logs collection daemon on every node, such as **fluentd** or **logstash**.
- + Running a node monitoring daemon on every node such as: **Prometheus Node Exporter**, **Sysdig Agent**, **collectd**, **gmond**, **NewRelic Agent**

## Jobs

A job controller creates one or more pods and ensures that a specified number of them are terminated successfully. When a specified number of successful completions is reached the job is complete.

Cronjob. A cron job creates a job on a time-based schedule.



# Services

A Service is what you wanna read about when you are trying to expose an APP that is running on a set of PODs as a network service.

K8S Give PODs their own IP Address and a single DNS for a set of PODs, you can of course load balance across them.

To set a set of PODs to be targeted by a Service you use Selectors (k-2.4) but Services is independently from Selectors.

## Define a service. (guess what? It's an Object too)

Suppose you have a set of PODs that each carry the label app=My-APP and they all listen to port 9376 TCP.

```
apiVersion: v1           ./svc.yaml
kind: Service
metadata:
  name: my-service
spec:
  selector:
    app: My-APP
  ports:
    - protocol: TCP
      port: 80
      targetPort: 9376
```

### Services without Selectors

Provide abstract access to K8S, here's why you want to do that:

- X You want to have an external db in Prod, but your test env uses your own db.
- X You want to point your Svc to a ≠ Svc on another Name Sp. or to another cluster.
- X Migration purposes.

# Publishing Services (ServiceTypes)

We are not going to describe each type in particular nor explain how they work. For more info on that please refer to the doc.

# If not explicitly specified, the default type is ClusterIP

**ClusterIP**: Exposes the service on a cluster internal IP. Choosing this value makes the service only reachable from within the cluster.

**NodePort**: Exposes the service on each Node's IP at a static port (the NodePort). A ClusterIP service to which the nodeport routes, is created automatically. You will be able to contact the service from outside the cluster by requesting <Node's IP>;<NodePort>

**LoadBalancer**: Exposes the service externally using a cloud provider's load balancer. NodePort & ClusterIP are created automatically.

**ExternalName**: Maps the service to the content of the external name field. (e.g. foo.bar.example.com) by returning a CNAME record with its value. No proxying of any kind is setup.

**DNS** Every service defined in the cluster gets a DNS name

Assume a service named **foo**, in the k8s namespace **bar**. A POD running under under **bar** can look up this service by **"foo"**. A POD running in a namespace **qux** can find this service by **"foo.bar"**.

