

Services are a core component in Kubernetes that are used to manage networking and traffic flow within a cluster. They provide a **stable IP address** and **DNS name** for a set of pods and allow for communication between different components within and outside of the application. Services also enable **load balancing**, **service discovery** and **traffic management**, making them a critical component for building scalable and resilient applications in Kubernetes.

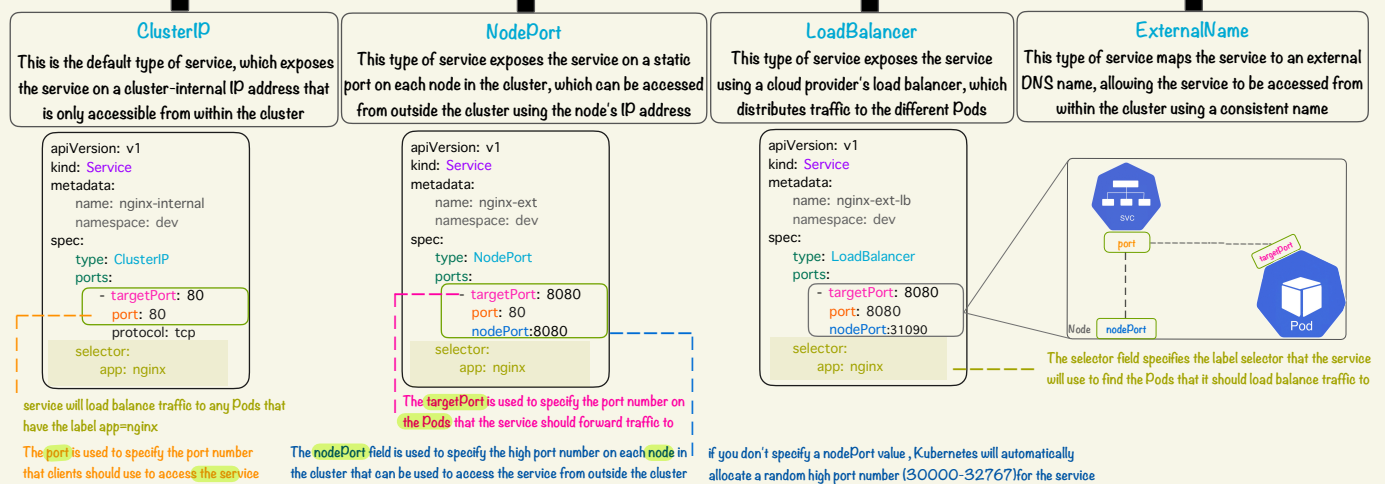
When a service is created, it is assigned a virtual IP address (known as a ClusterIP), which is used to route traffic to the pods that are part of the service. The service also has a DNS name, which can be used to access the service from within the cluster

- 1 If a pod fails or is removed from the service, controller will automatically remove it from the list of endpoints for the service. This ensures that traffic is not sent to a non-existent pod.
- 2 When a pod managed by a deployment fails, the controller creates a new pod to replace the failed pod
- 3 Once the new pod is running and ready, service's endpoint controller will add it back to the list of endpoints for the service, allowing traffic to be routed to it (used labels and selectors to discover)

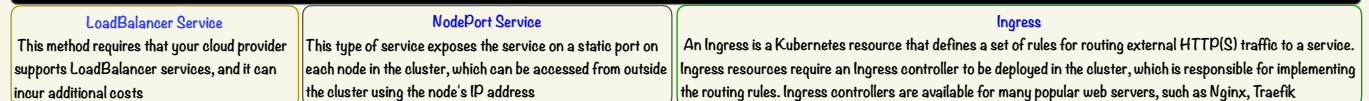
services use **labels** and **selectors** to **discover and route** traffic to the pods that are part of the service

Each service has a unique IP address and DNS name that can be used to access the pods that provide the service.

Service types in k8s



To expose a service to the outside world in k8s, you can use one of the following methods



While a load balancer service can provide a stable IP address and port for accessing the service, it still requires manual intervention to update the endpoints, which can be time-consuming and error-prone. Therefore, a better solution to this problem would be to use Kubernetes Ingress, which provides a more flexible and automated way of managing external access to the services in a k8s cluster

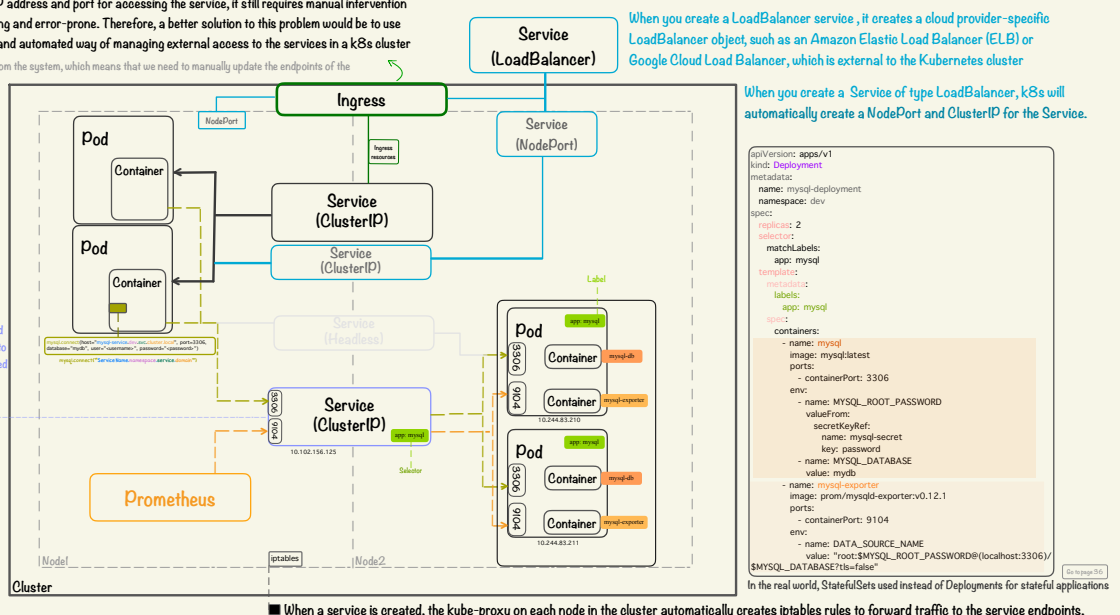
It is possible that a node might be added or removed from the system, which means that we need to manually update the endpoints of the load balancer or recreate the service.

When you create a LoadBalancer service, it creates a cloud provider-specific LoadBalancer object, such as an Amazon Elastic Load Balancer (ELB) or Google Cloud Load Balancer, which is external to the Kubernetes cluster

Headless services are used for direct access to pods

For services that we do not want to expose to the outside world (such as database clusters like mysql), we set the service type to ClusterIP. If we do not specify the service type, it will be selected as ClusterIP by default

```
apiVersion: v1
kind: Service
metadata:
  name: mysql-service
spec:
  selector:
    matchLabels:
      app: mysql
  ports:
    - name: mysql
      targetPort: 3306
    - name: mysql-exporter
      port: 9104
      targetPort: 9104
```



When a service is created, the kube-proxy on each node in the cluster automatically creates iptables rules to forward traffic to the service endpoints.

NAME	TYPE	CLUSTER-IP	EXTERNAL-IP	PORT(S)	AGE	SELECTOR	NODE-PORT	ENDPOINTS
mysql-service	ClusterIP	10.102.156.125	<none>	3306/TCP,9104/TCP	1d	app=mysql	<none>	10.244.83.210:3306, 10.244.83.211:3306

```
iptables -A KUBE-SERVICES -d 10.102.156.125/32 -p tcp -m comment --comment "dev/mysql: cluster IP" -m tcp --dport 3306 -j KUBE-SVC-ABC123
iptables -A KUBE-SVC-ABC123 -m comment --comment "dev/mysql-service" -j KUBE-SEP-abc123
iptables -A KUBE-SVC-ABC123 -m comment --comment "dev/mysql-service" -j KUBE-SEP-abc123
```