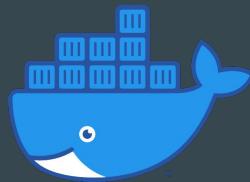




# Gentle introduction to scaling up ML service with Kubernetes + Mlflow



**mlflow**<sup>TM</sup>

PyData NYC, November 2022

Kei Nemoto

# About me

- Kei Nemoto
- Data Scientist @ Montefiore Einstein Center for Health Data Innovations
  - MLOps
  - CI/CD
  - ETL
  - Backend Development
  - NLP



[linkedin.com/in/kei-nemoto](https://linkedin.com/in/kei-nemoto)



[github.com/box-key](https://github.com/box-key)

# All code is available on Github!

**[github.com/box-key/pydata-kubernetes-mlflow](https://github.com/box-key/pydata-kubernetes-mlflow)**

# ML Service Tech Stack



- Machine learning framework



- Web framework



**gunicorn**

- WSGI Http server

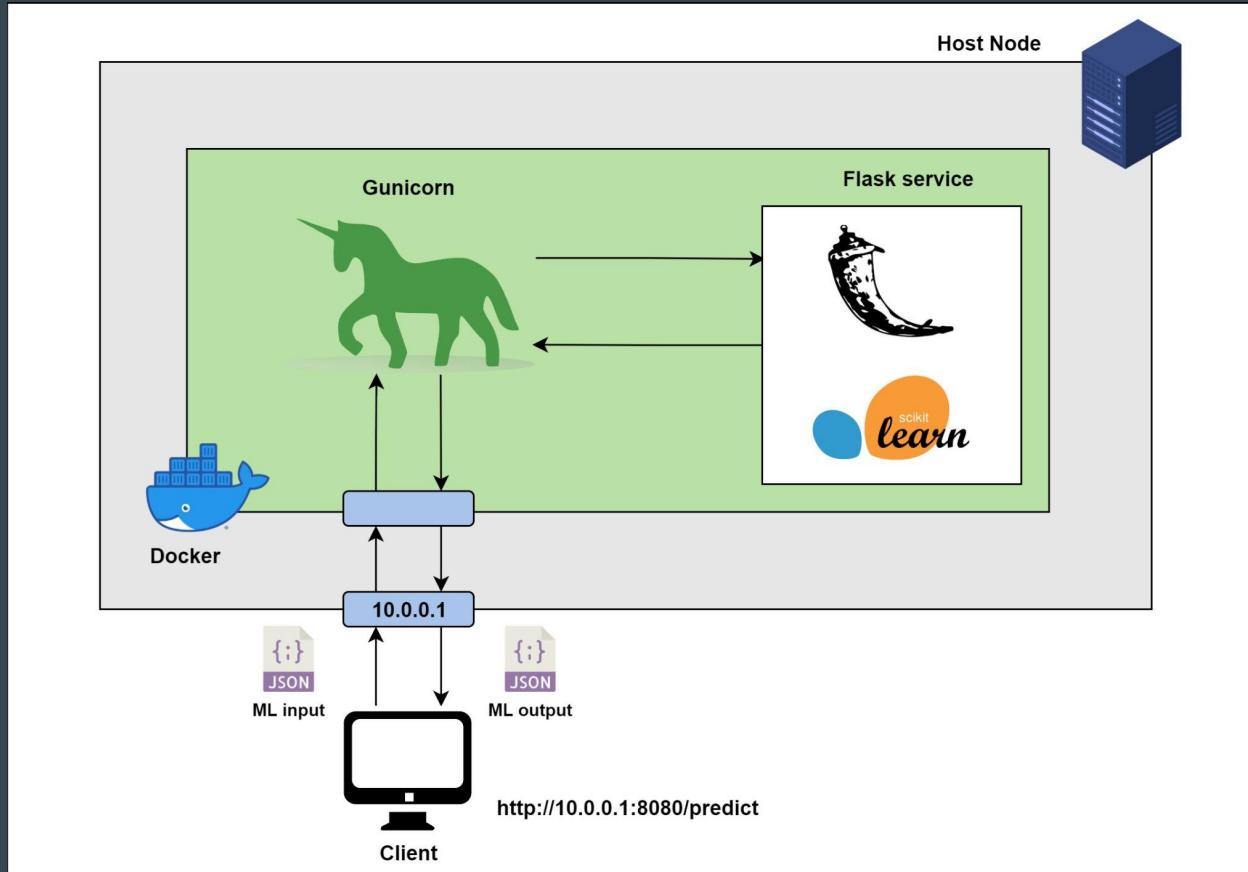


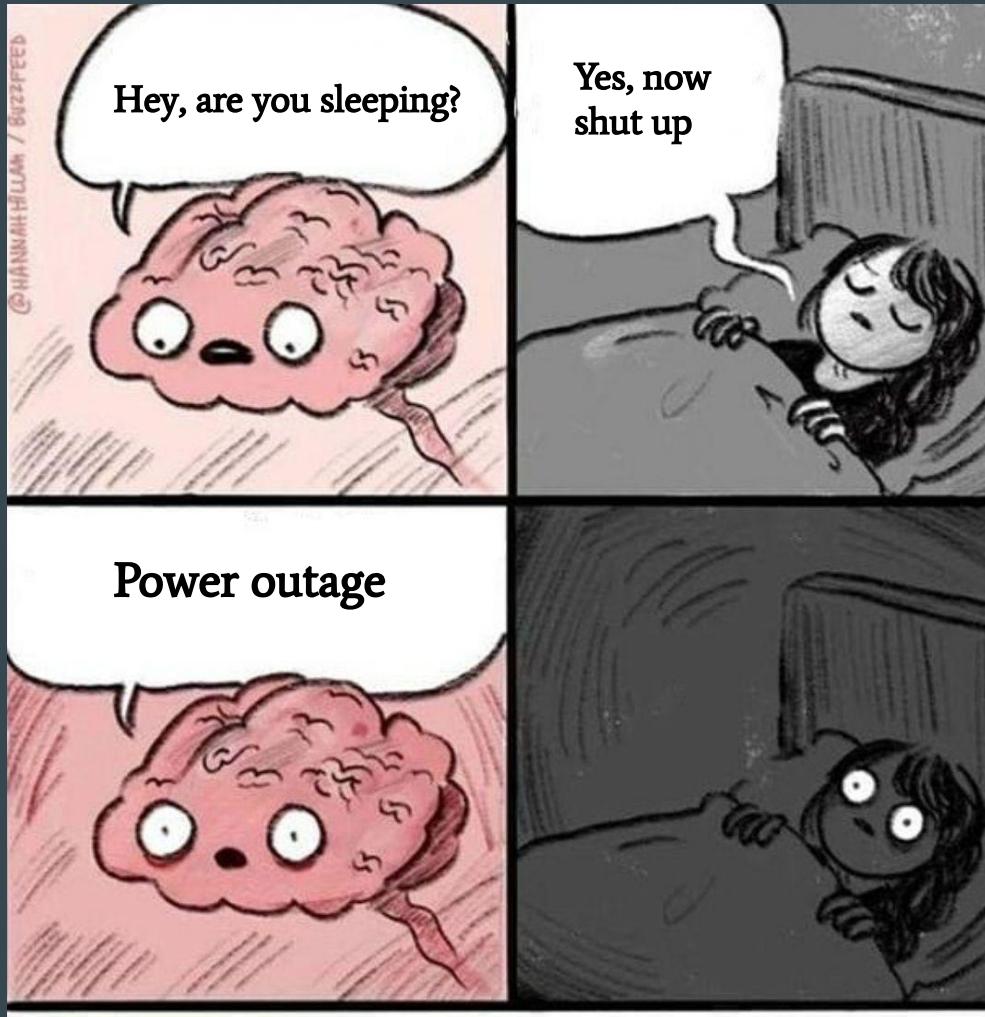
**docker**

- Container development tool

```
1  from flask import Flask, request
2  from flask_restful import Resource, Api
3  from sklearn.ensemble import RandomForestClassifier
4
5
6  class IrisRandomForestPredictor(Resource):
7      def __init__(self, **kwargs):
8          self.model: RandomForestClassifier = kwargs["model"]
9
10     def post(self):
11         payload = request.get_json()
12         input_arr = [
13             float(payload["sepal_length"]),
14             float(payload["sepal_width"]),
15             float(payload["petal_length"]),
16             float(payload["petal_width"])
17         ]
18         prediction = int(self.model.predict([input_arr])[0])
19         return {
20             "status": 200,
21             "message": "Successfully predicted!",
22             "data": {
23                 "prediction": prediction
24             }
25         }
```

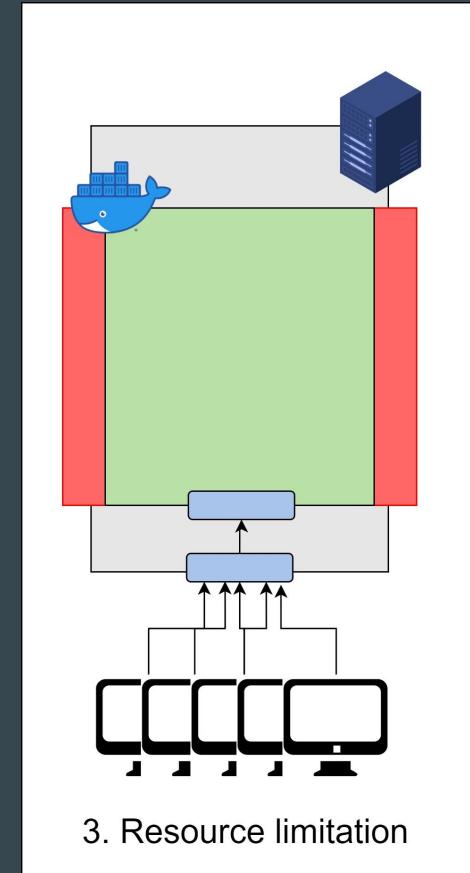
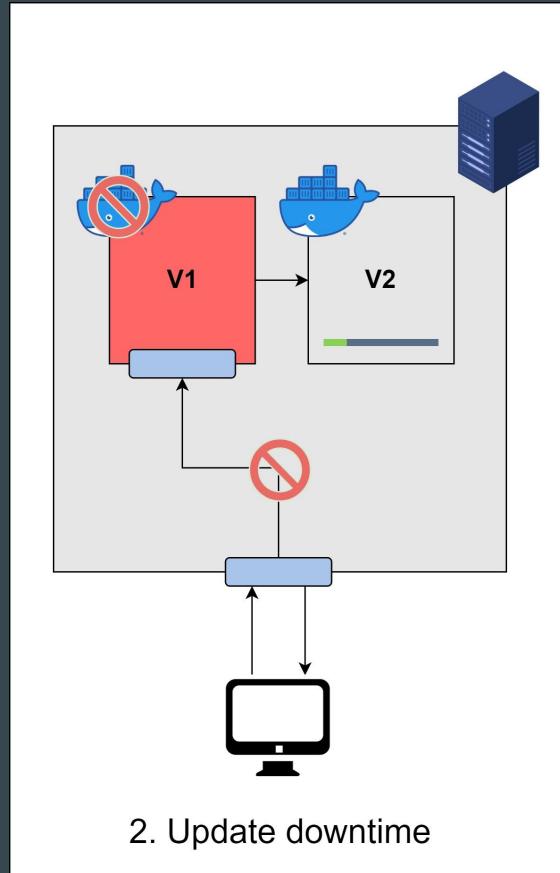
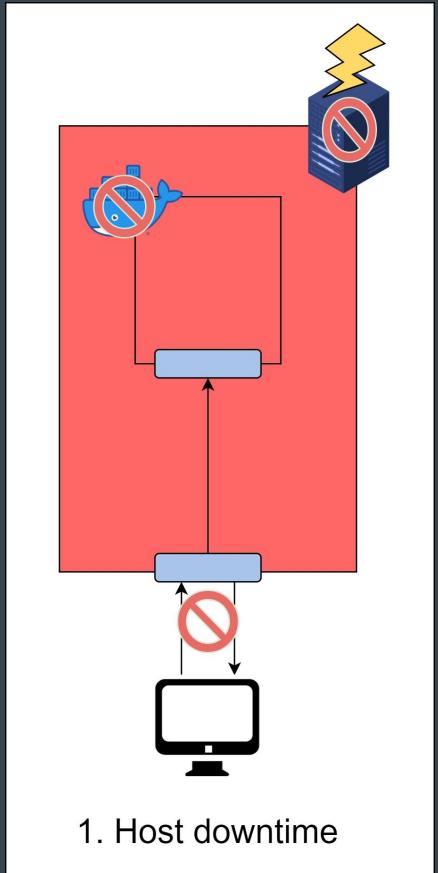
# Single Container Deployment



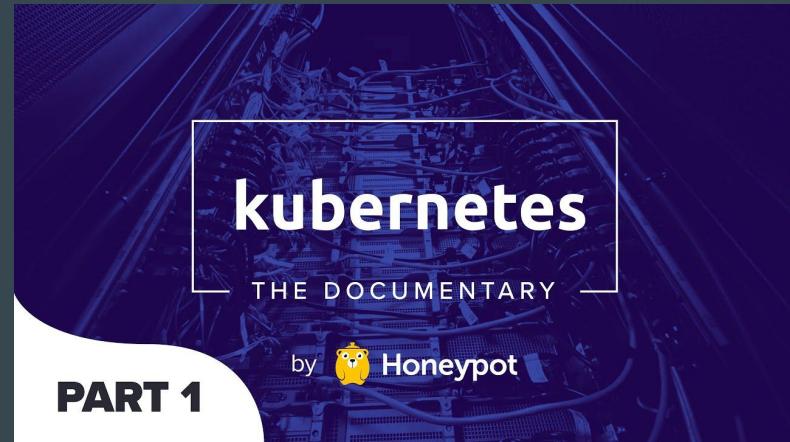
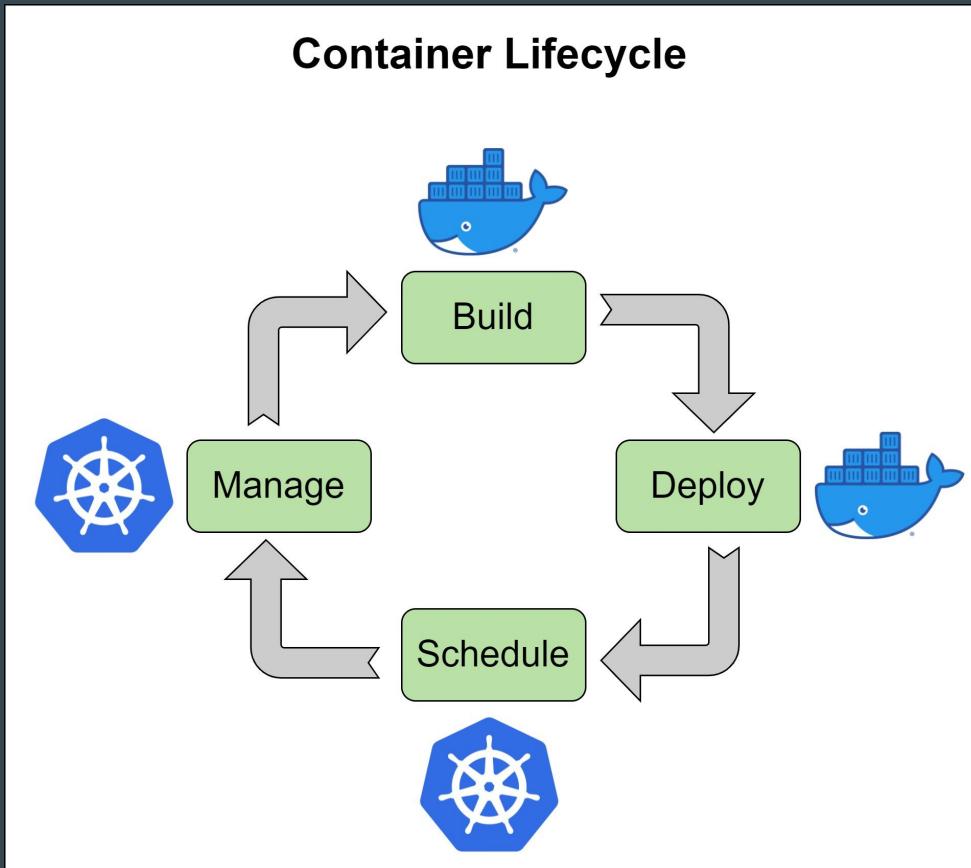


432.J2259 / ພັດທະນາຄະດີ

# Problems of Single Container Deployment



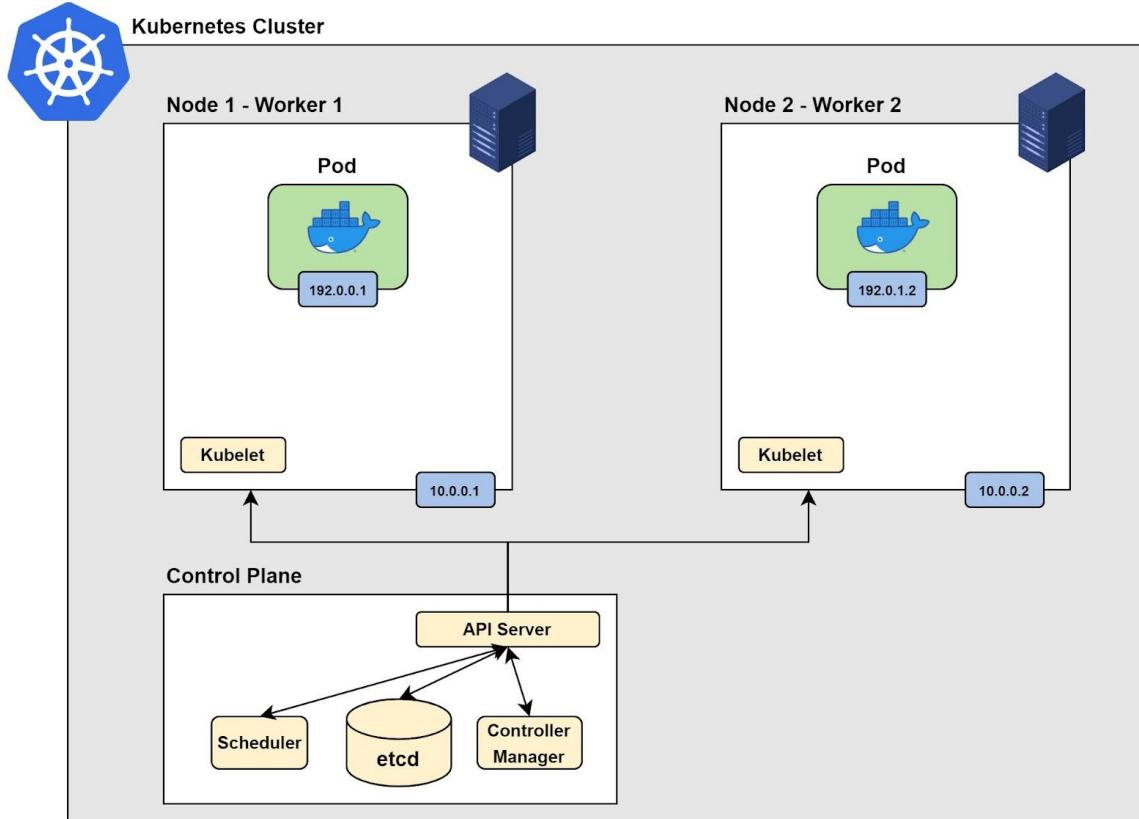
# Kubernetes = Container Management System



Let's invent the post office. Docker creates an envelope (container).  
Kubernetes is the postal system that delivers envelopes to the destination.

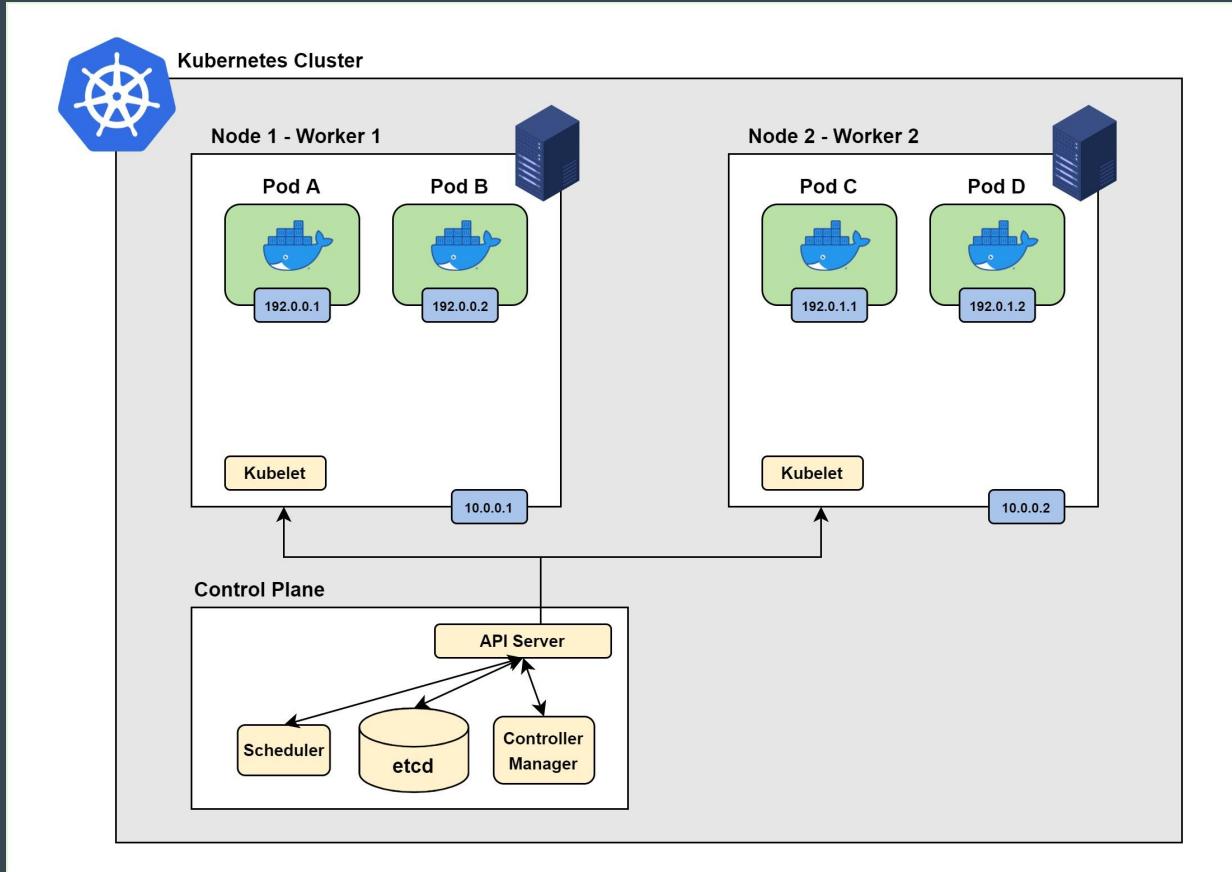
Reference: "Kubernetes: The Documentary [PART 1]"  
by Honeypot on Jan 21, 2022.  
<https://www.youtube.com/watch?v=BE77h7dmoQU>

# Kubernetes 101

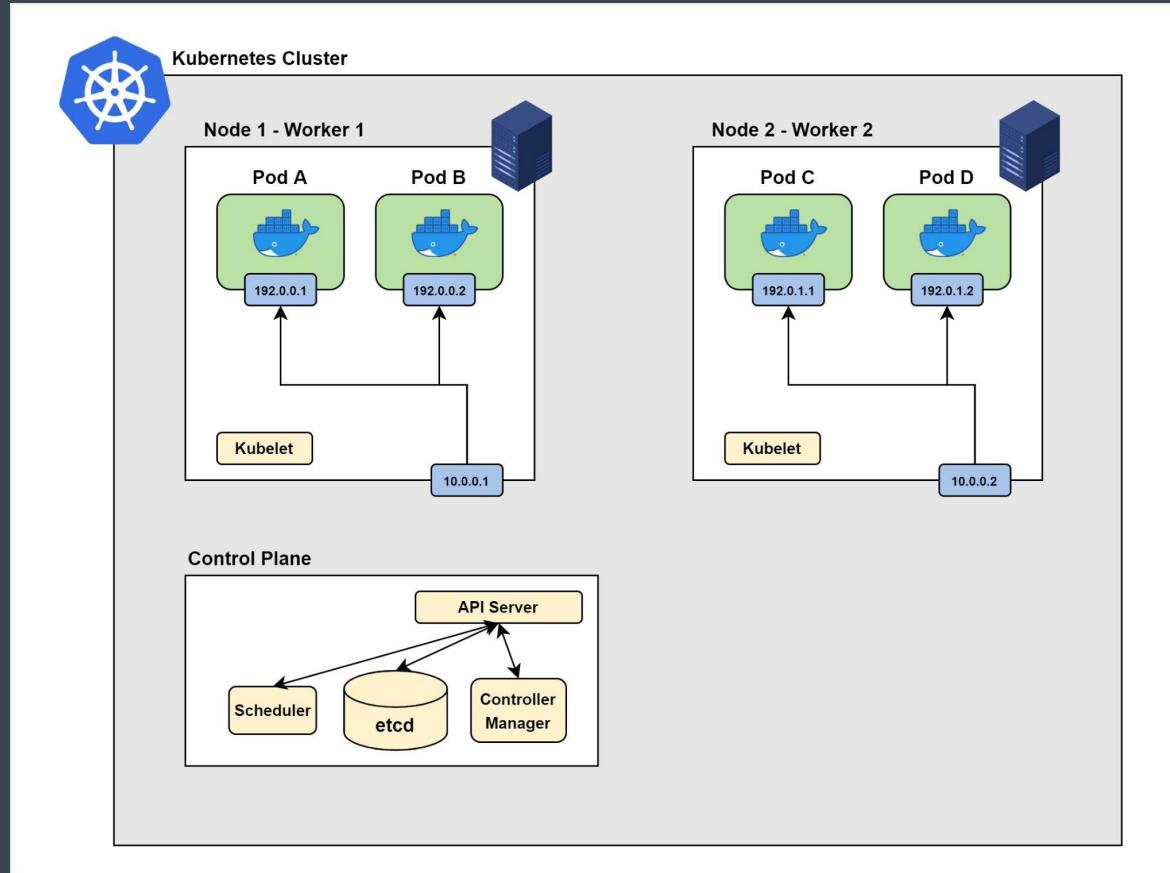


- Kubernetes cluster consists of *Control Plane* and *Workers*
- *Pod* executes container
- Each *pod* gets unique IP address

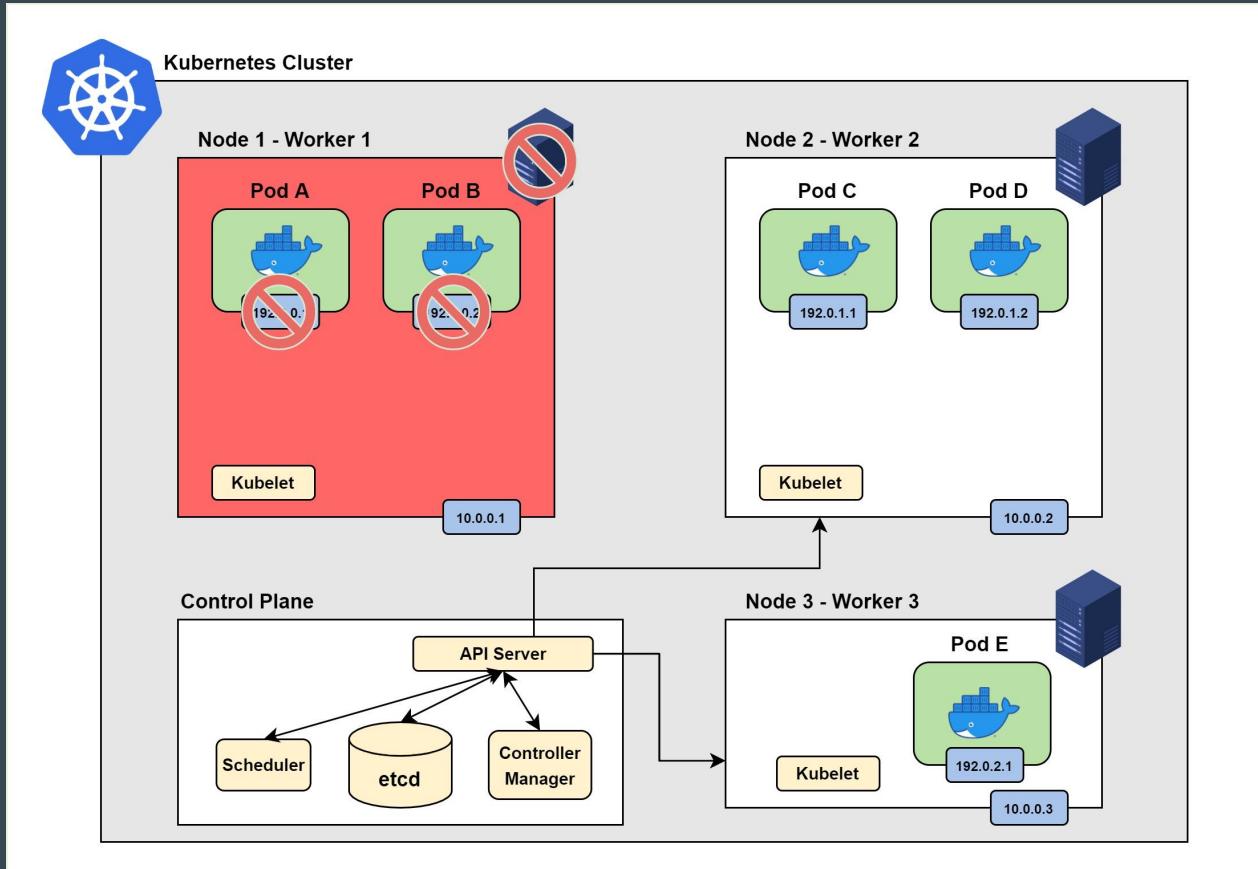
# Deployment



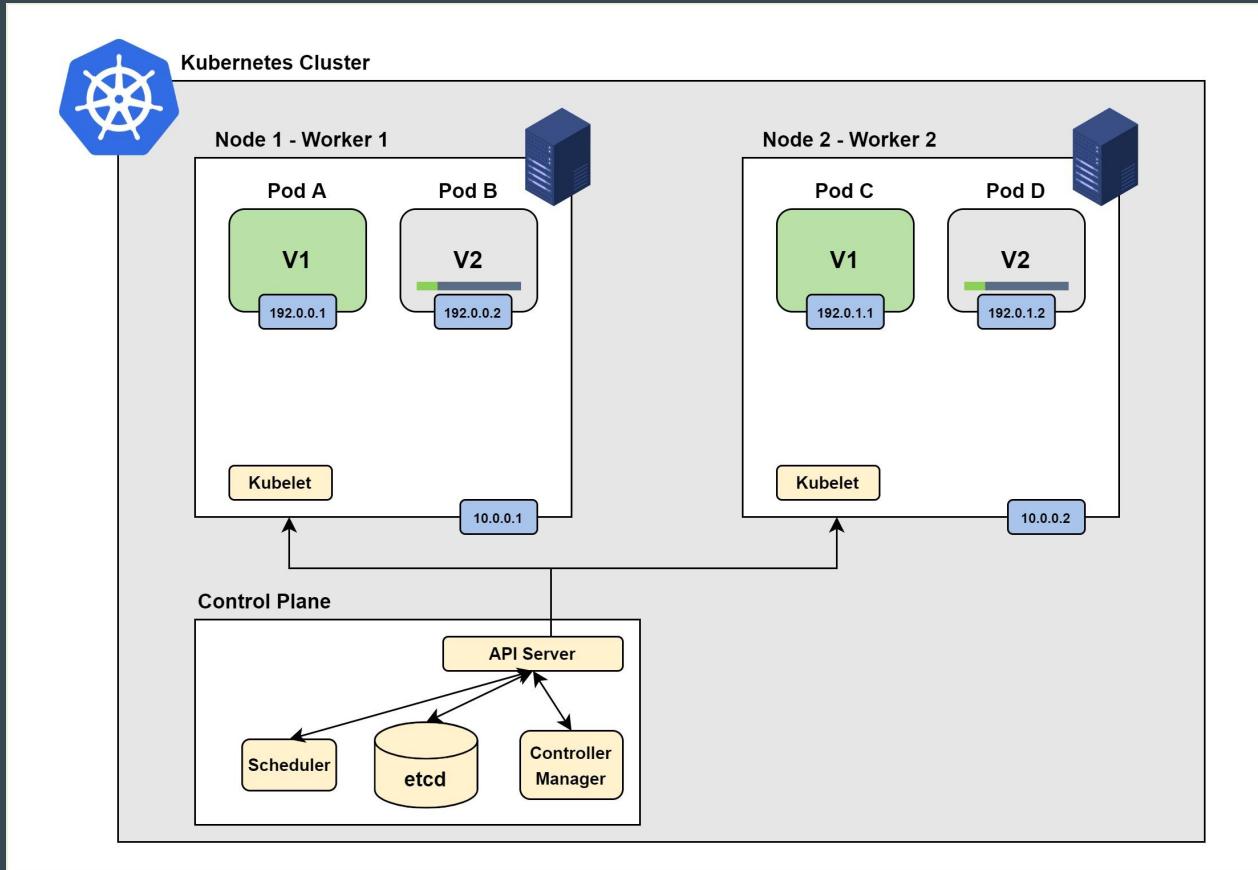
# Deployment - Scalability



# Deployment - High Availability



# Deployment - Rolling update



```
1   apiVersion: apps/v1
2   kind: Deployment
3   metadata:
4     name: iris-rf
5     namespace: deployment-demo
6   labels:
7     app: iris-rf
8   spec:
9     replicas: 3
10    selector:
11      matchLabels:
12        app: iris-rf
13    template:
14      metadata:
15        labels:
16          app: iris-rf
17      spec:
18        containers:
19          - name: iris-rf
20            image: boxkey/iris-rf:beta
21        ports:
22          - containerPort: 5000
```

# Create deployment

```
> kubectl apply -f kubernetes/deployment-demo.yaml
```

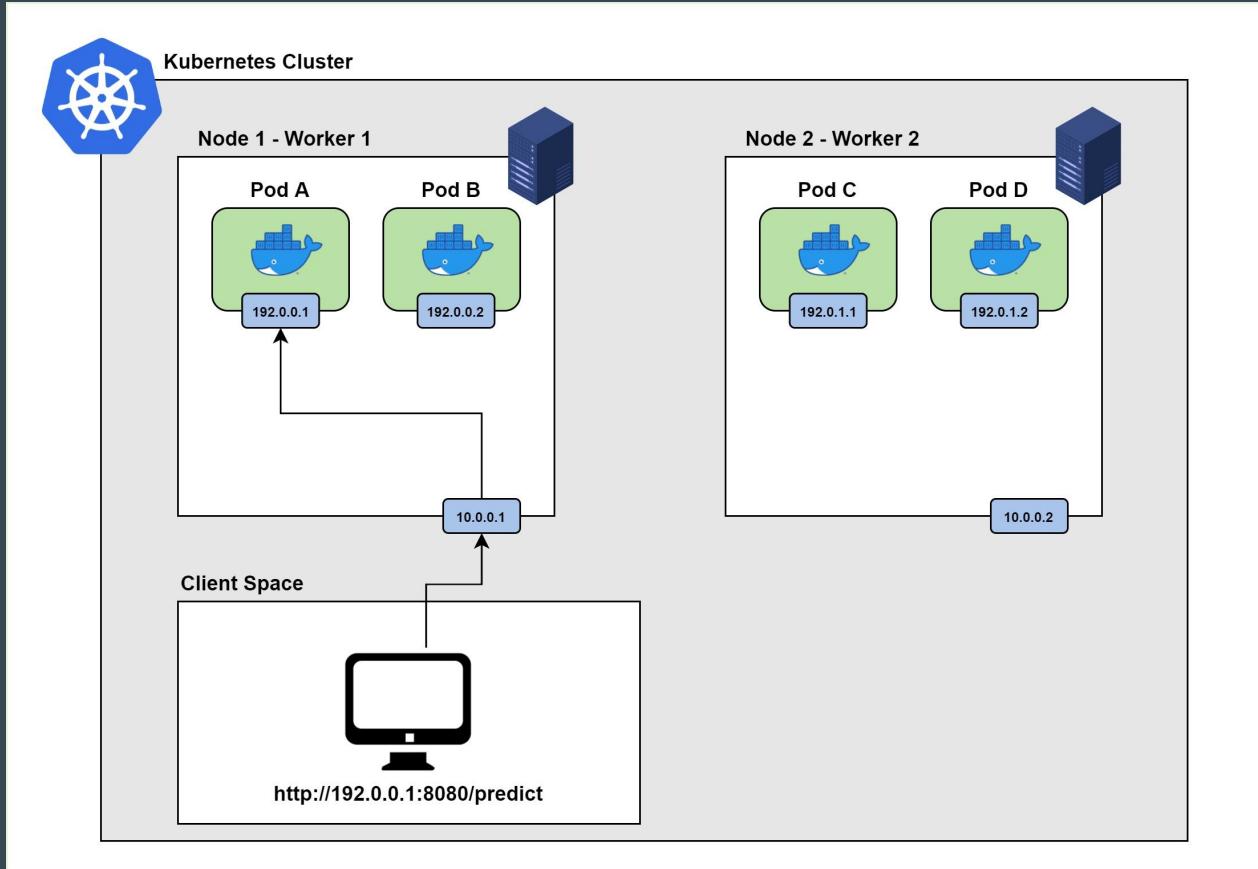
```
deployment.apps/iris-rf created
```

# Print pods in deployment

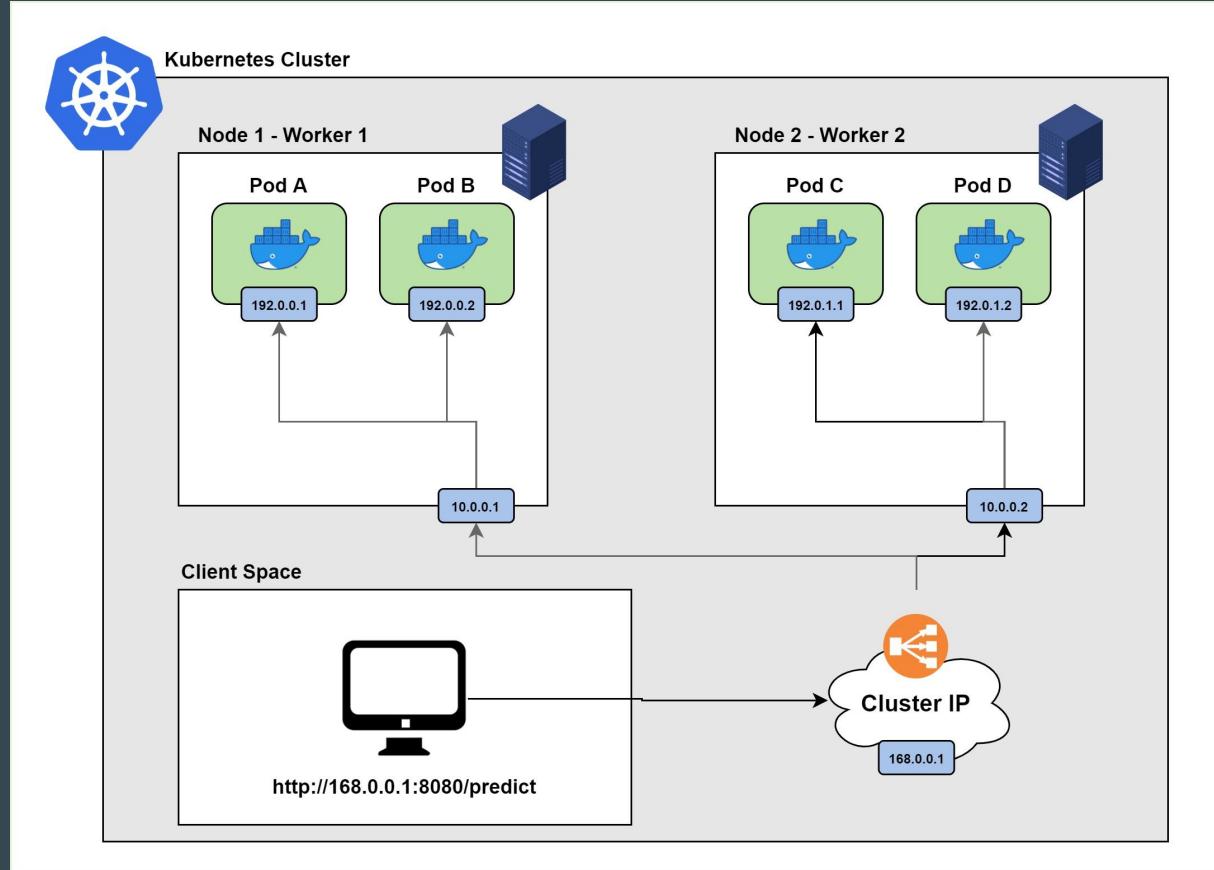
```
> kubectl get pods -n deployment-demo \
> -o=custom-columns=NAME:.metadata.name, \
> IP:.status.podIP,STATUS:.status.phase, \
> IMAGE:.spec.containers[0].image
```

NAME	IP	STATUS	IMAGE
iris-rf-8c7857fd6-v52qc	10.42.1.80	Running	boxkey/iris-rf:beta
iris-rf-8c7857fd6-s9dmp	10.42.0.85	Running	boxkey/iris-rf:beta
iris-rf-8c7857fd6-kgx72	10.42.2.70	Running	boxkey/iris-rf:beta

# Which pod to send request to?



# Cluster IP Service



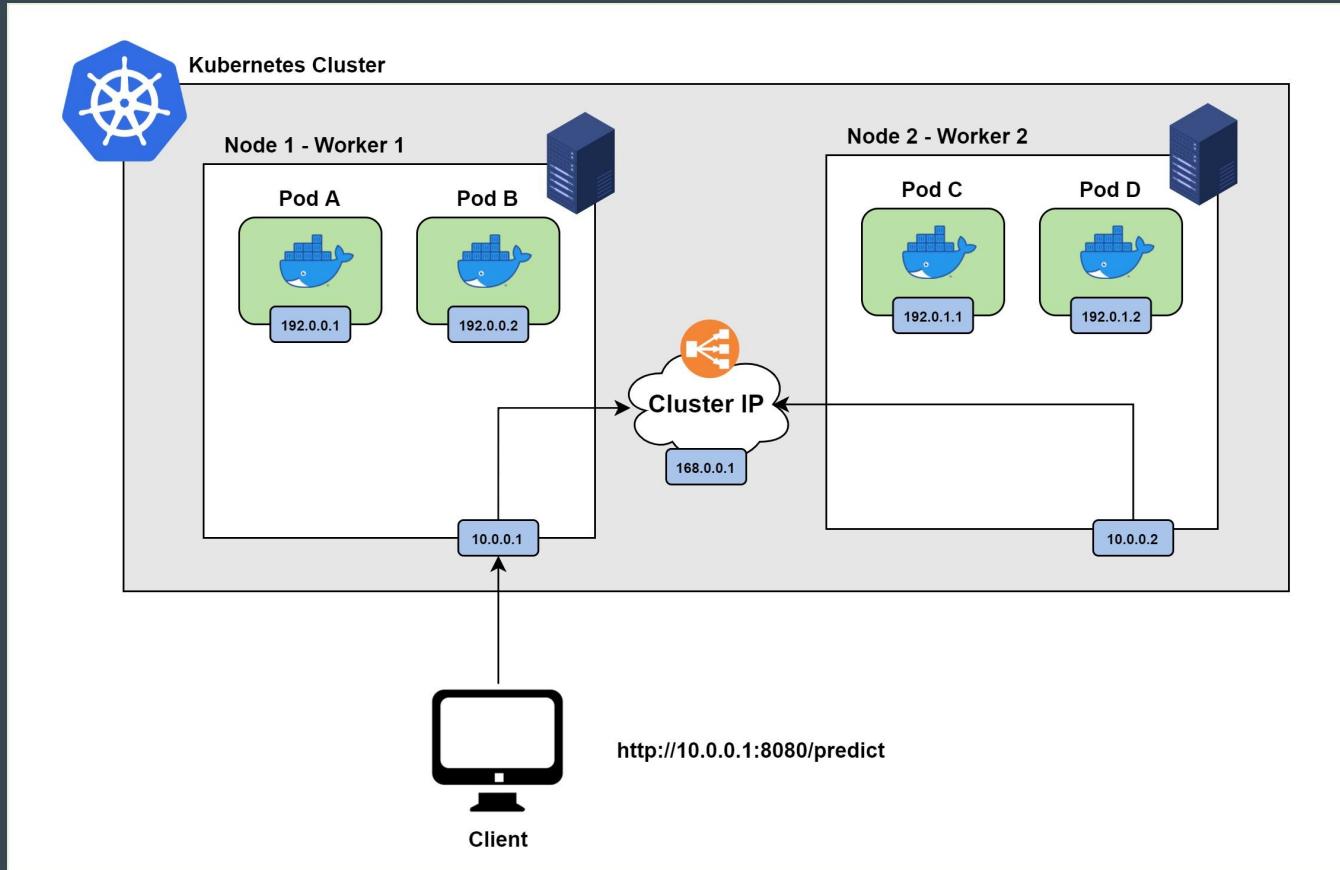
```
1  apiVersion: apps/v1
2  kind: Deployment
3  metadata:
4    name: iris-rf
5    namespace: cluster-ip-demo
6    labels:
7      app: iris-rf
8  spec:
9    replicas: 3
10   selector:
11     matchLabels:
12       app: iris-rf
13   template:
14     metadata:
15       labels:
16         app: iris-rf
17     spec:
18       containers:
19         - name: iris-rf
20           image: boxkey/iris-rf:beta
21           ports:
22             - containerPort: 5000
```

```
37  apiVersion: v1
38  kind: Service
39  metadata:
40    name: iris-svc
41    namespace: cluster-ip-demo
42  spec:
43  ports:
44    - port: 5000
45      protocol: TCP
46      targetPort: 5000
47  selector:
48    app: iris-rf
49  type: ClusterIP
```

```
> kubectl get services -n cluster-ip-demo
```

NAME	TYPE	CLUSTER-IP	EXTERNAL-IP	PORT(S)	AGE
iris-svc	ClusterIP	10.43.251.45	<none>	5000/TCP	26s

# NodePort Service

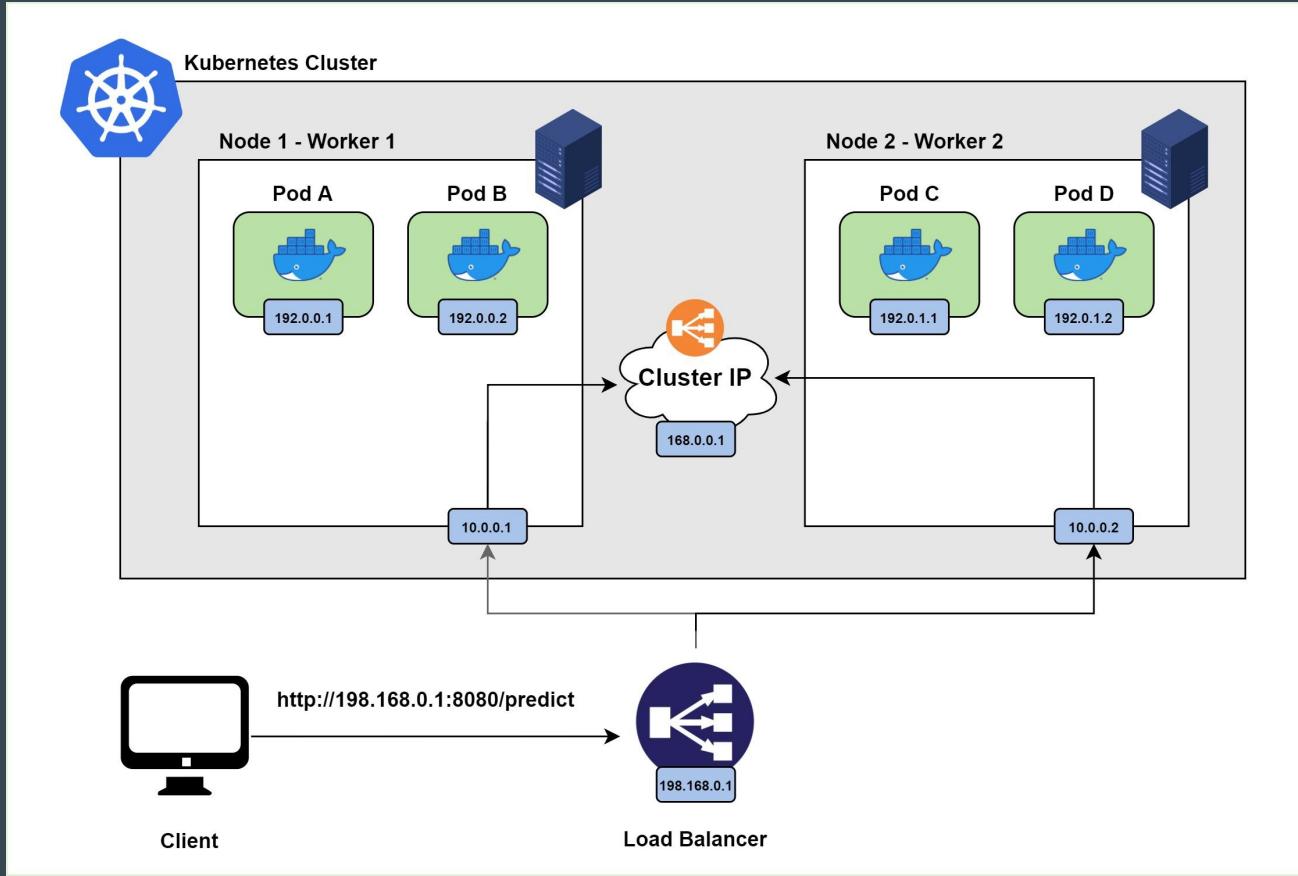


```
37    apiVersion: v1
38    kind: Service
39  < metadata:
40      name: iris-svc
41      namespace: node-port-demo
42  < spec:
43  <   ports:
44    <     - port: 5000
45    <       protocol: TCP
46    <       targetPort: 5000
47  <   selector:
48    <       app: iris-rf
49    type: NodePort
```

```
> kubectl get services -n node-port-demo
```

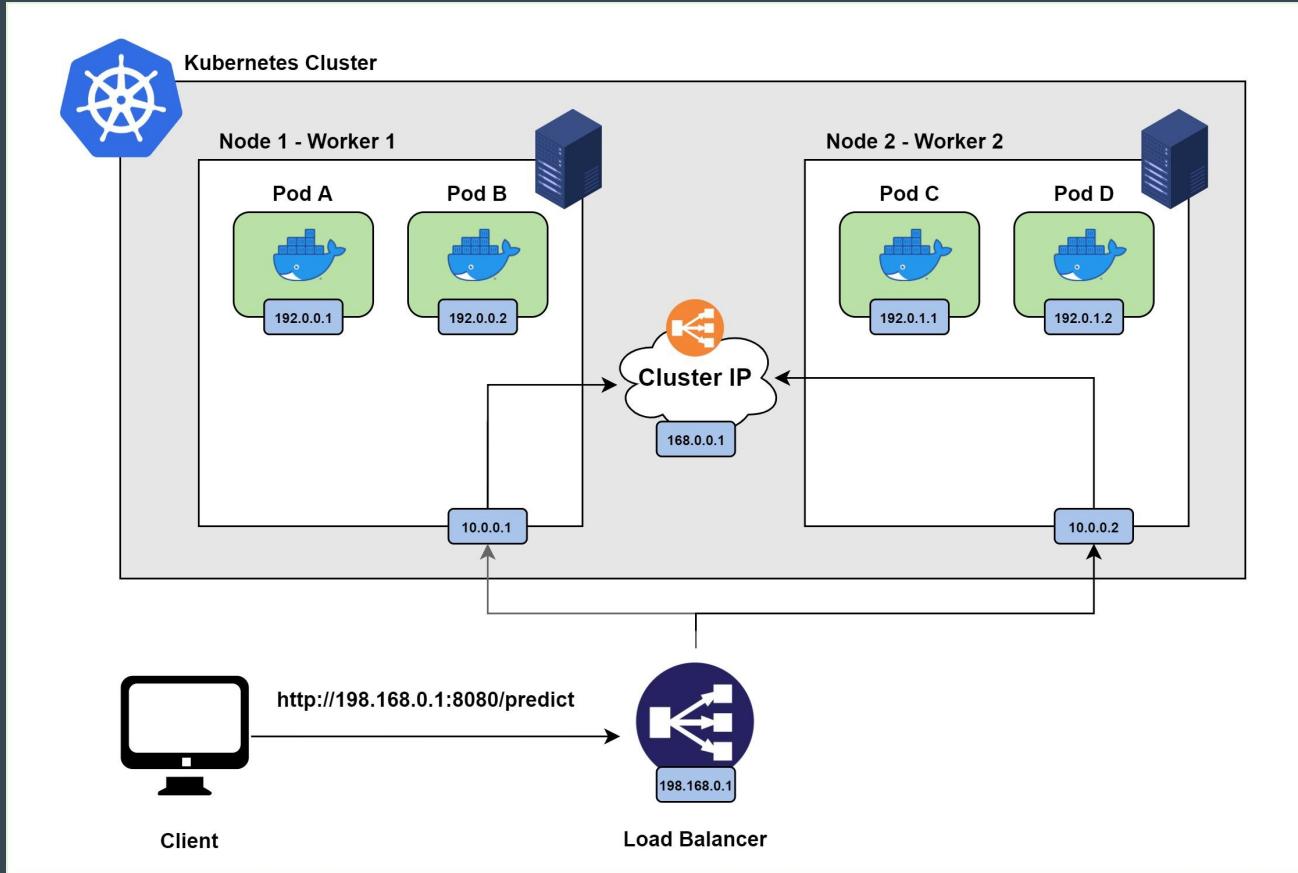
NAME	TYPE	CLUSTER-IP	EXTERNAL-IP	PORT(S)	AGE
iris-svc	NodePort	10.43.11.78	<none>	5000:31304/TCP	4h48m

# LoadBalancer Service

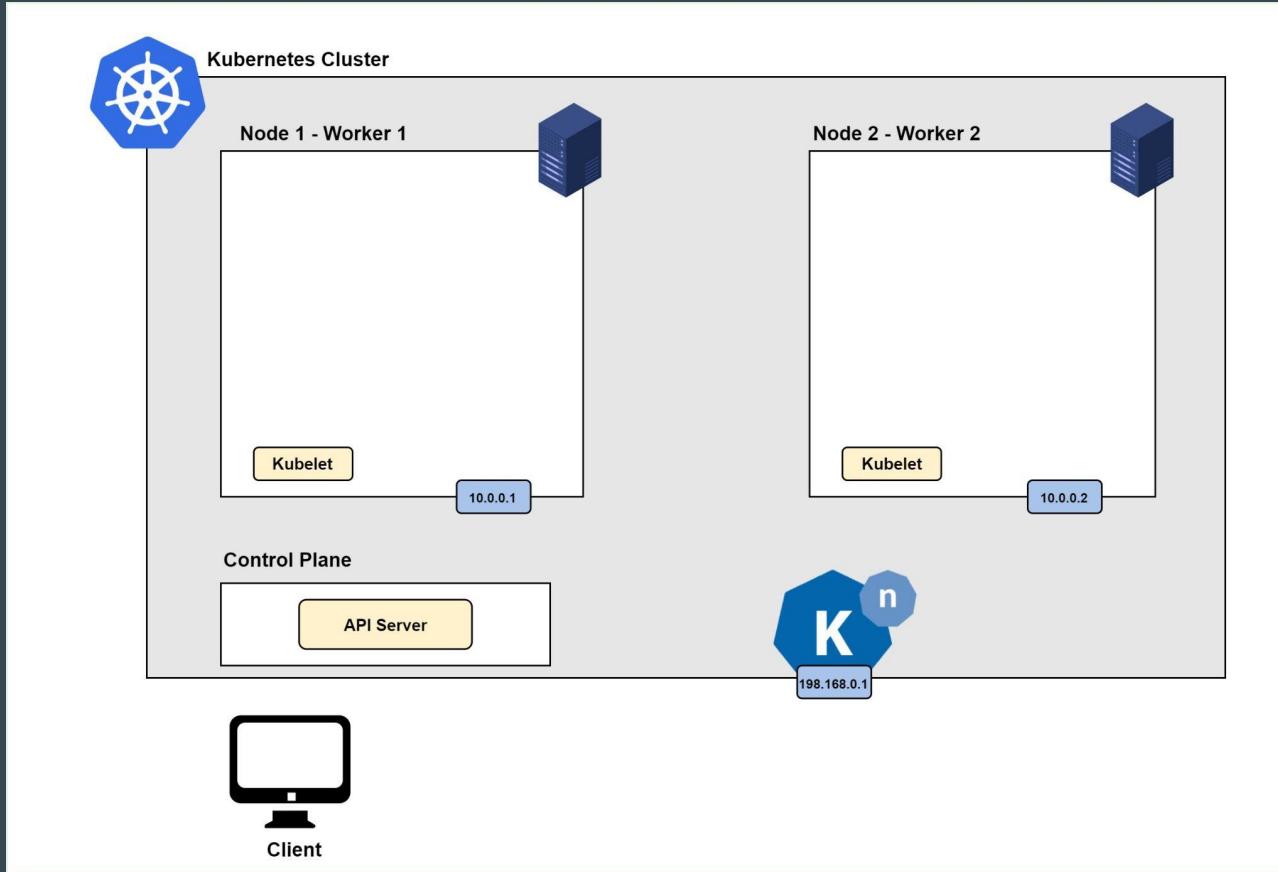


```
37    apiVersion: v1
38    kind: Service
39    metadata:
40      name: iris-svc
41      namespace: load-balancer-demo
42    spec:
43      ports:
44        - port: 5000
45          protocol: TCP
46          targetPort: 5000
47      selector:
48        app: iris-rf
49      type: LoadBalancer
```

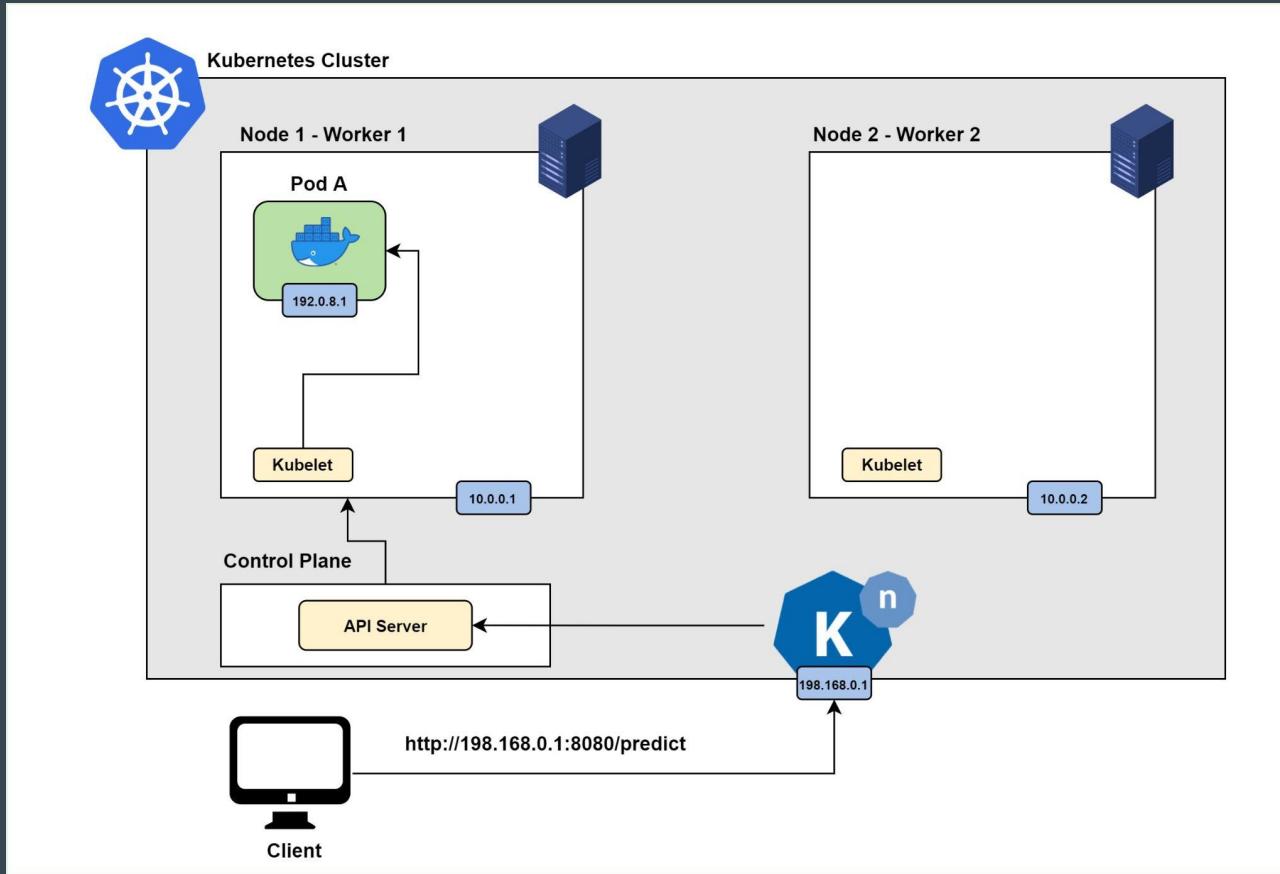
# Should we run 4 pods?



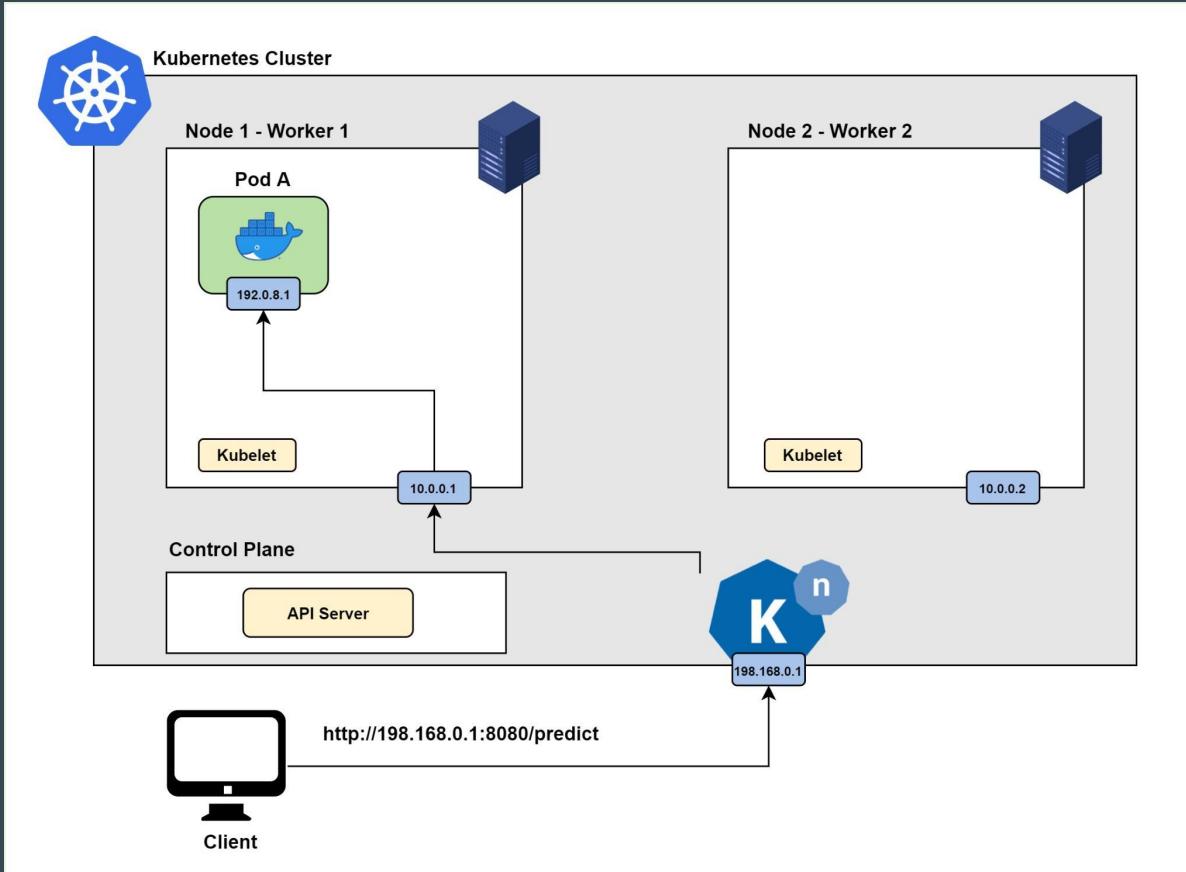
# Serverless Service by Knative



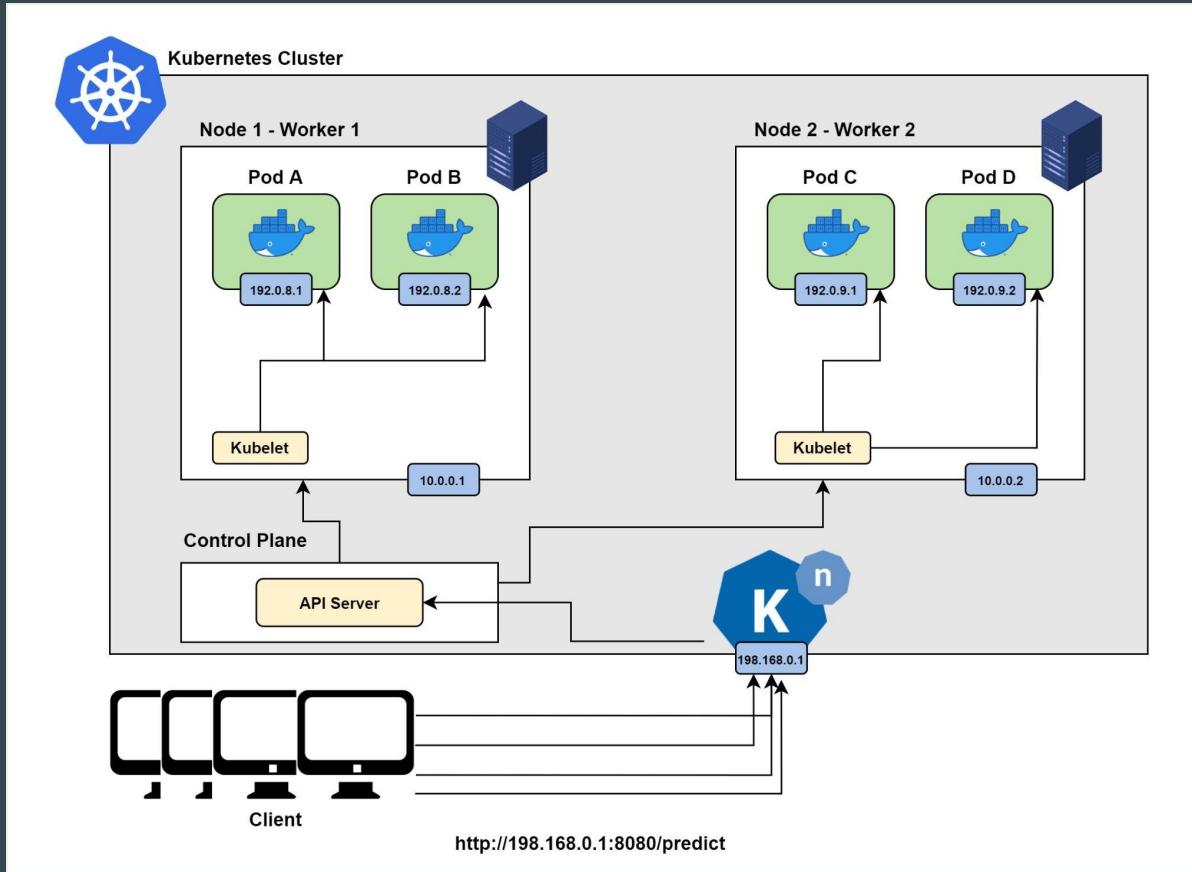
# Serverless Service by Knative - Pod Creation



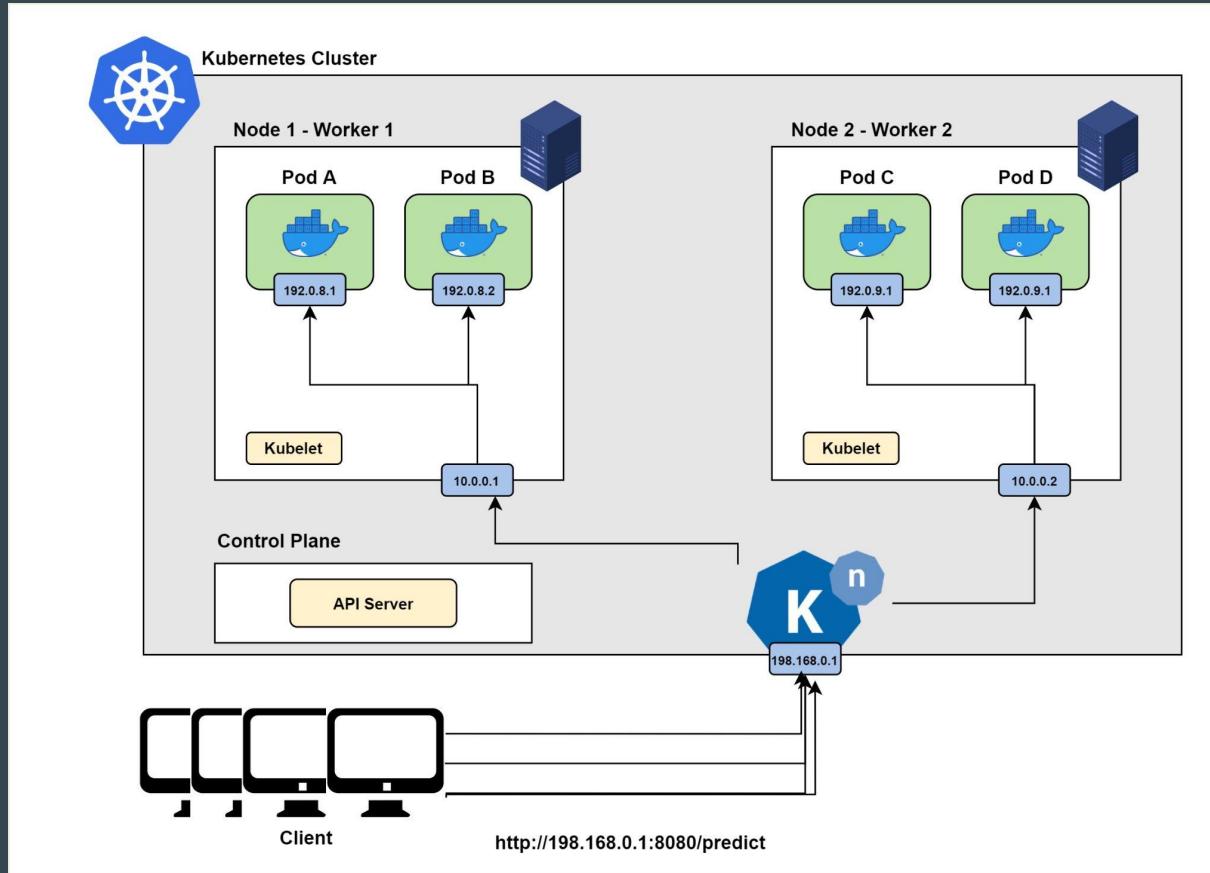
# Serverless Service by Knative - Forward request



# Automatic Scaling - Pod Creation

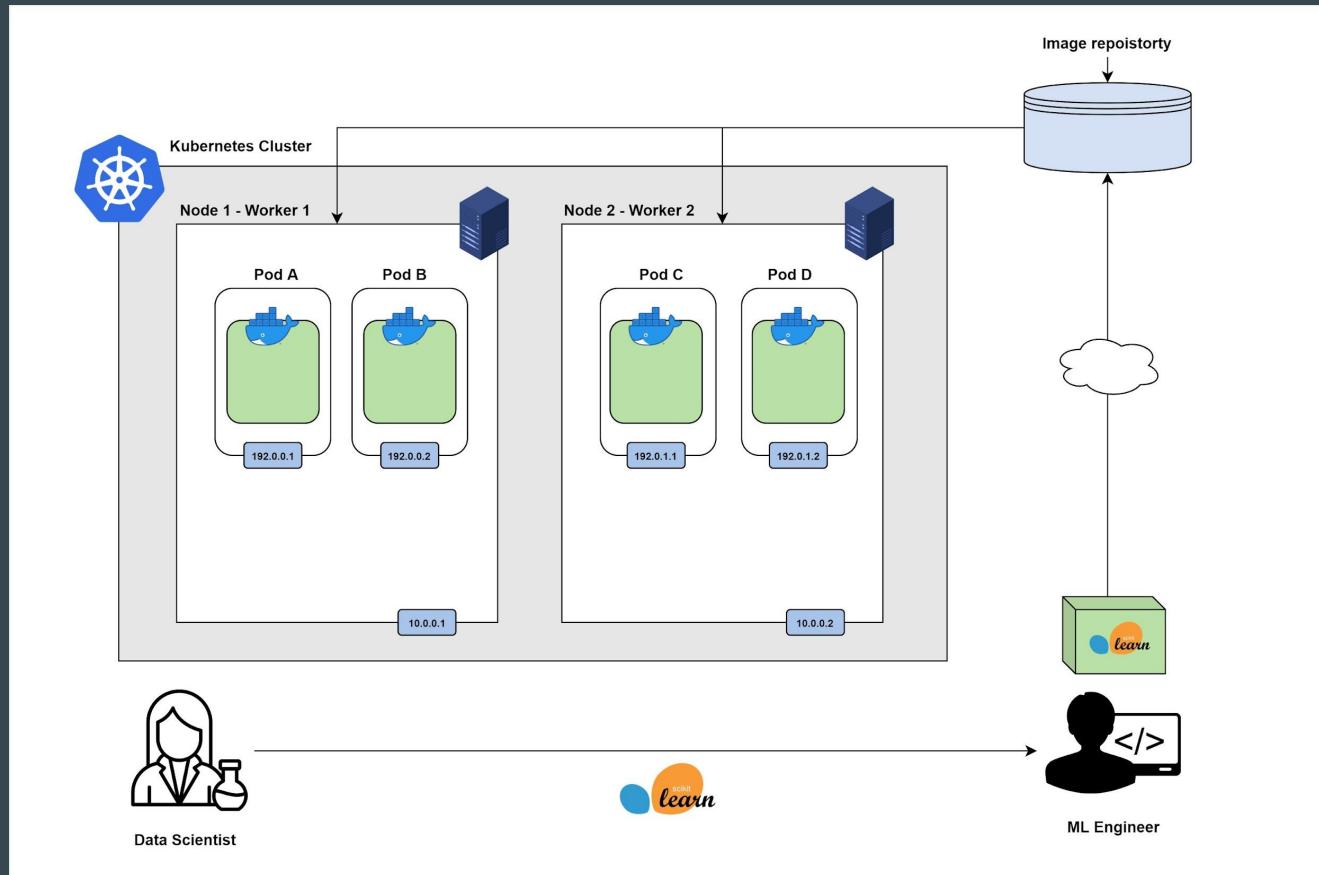


# Automatic Scaling - Distribute Requests

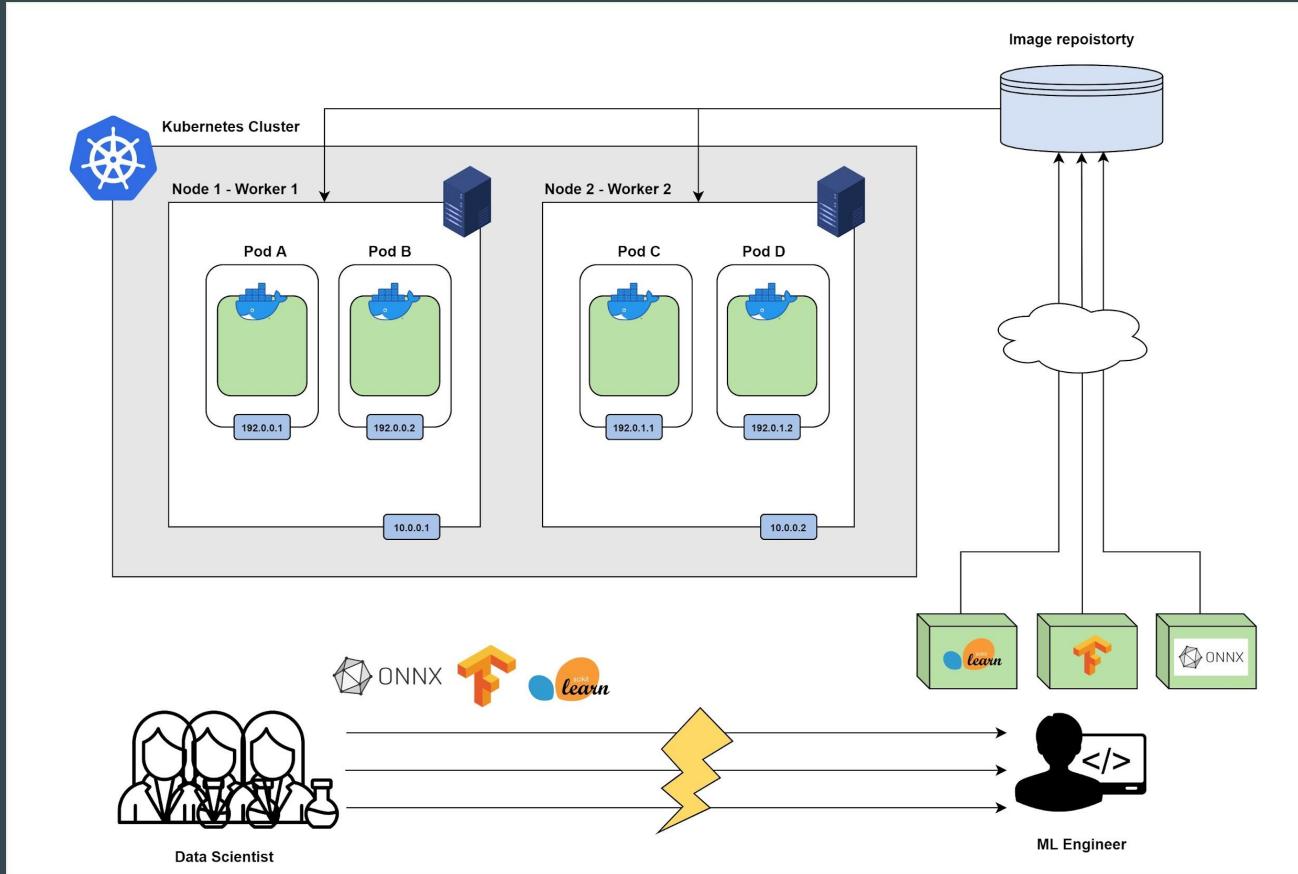


```
1  apiVersion: serving.knative.dev/v1
2  kind: Service
3  metadata:
4    name: iris-rf
5    namespace: knative-demo
6  spec:
7    template:
8      spec:
9        containers:
10       - name: iris-rf
11         image: boxkey/iris-rf:beta
12       ports:
13         - containerPort: 5000
```

# Model in Kubernetes

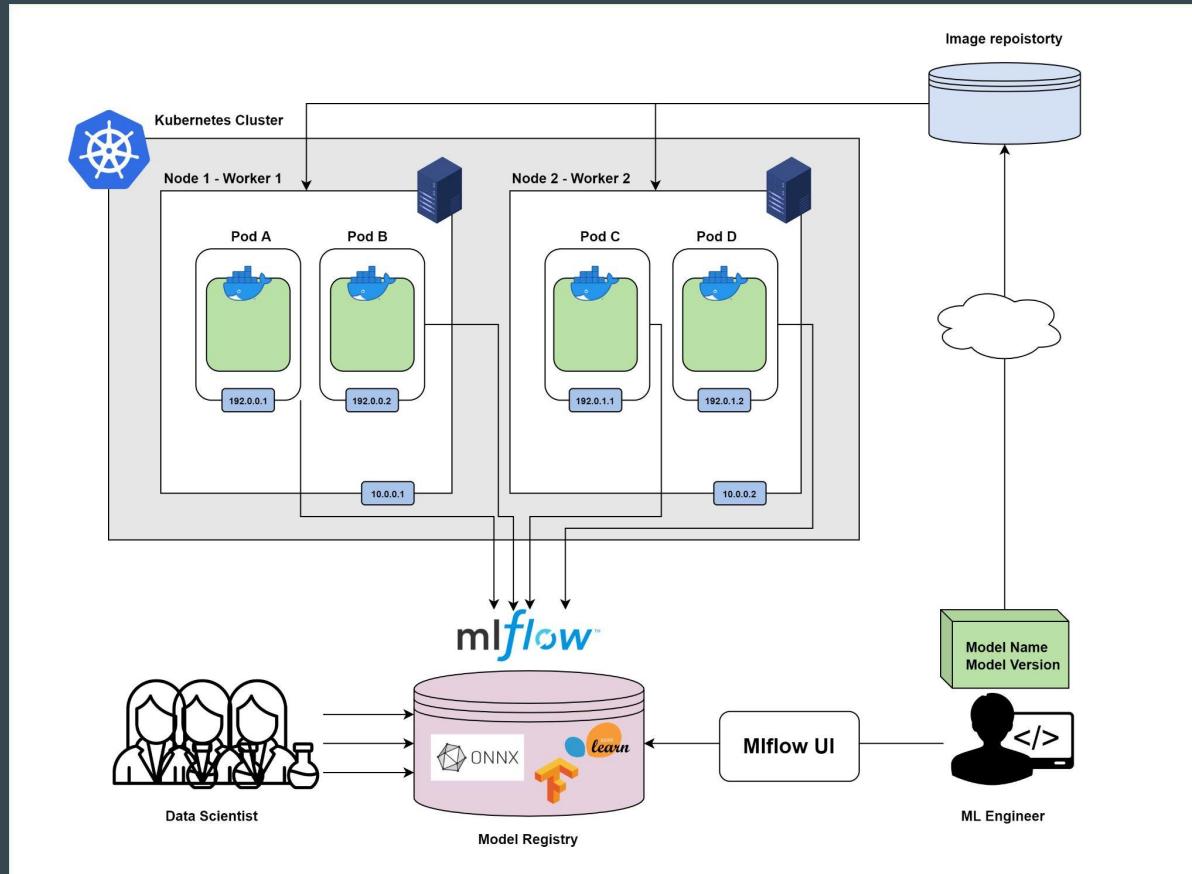


# Disconnection between Data Scientist and Developer





# Mlflow Model Registry + Kubernetes



Registered Models &gt; Iris\_Random\_Forest

## Iris\_Random\_Forest

Created Time: 2022-10-30 13:57:45

Last Modified: 2022-11-03 18:03:43

▼ Description [Edit](#)

Random forest model for iris prediction

▶ Tags

▼ Versions

[All](#)[Active 2](#)[Compare](#)

<input type="checkbox"/>	Version	Registered at	Created by	Stage
<input type="checkbox"/>	Version 2	2022-11-03 18:02:52		<span>Staging</span>
<input checked="" type="checkbox"/>	Version 1	2022-10-30 13:57:45		<span>Production</span>

# Register Model in Mlflow

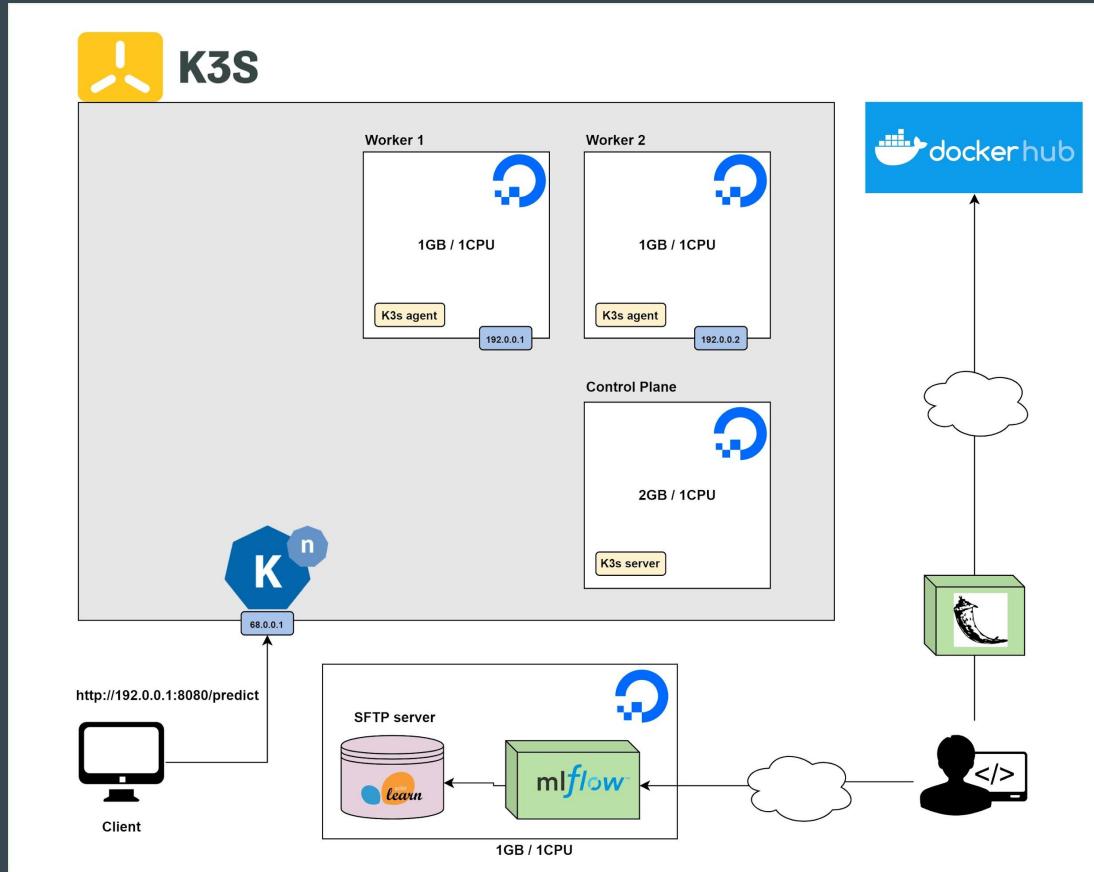
```
17 # train rf classifier
18 logger.info("Training random forest model..")
19 iris = datasets.load_iris()
20 clf = rf()
21 clf.fit(iris.data, iris.target)
22 # save model in mlflow
23 logger.info("Saving trained model in mlflow..")
24 name = "iris-rf"
25 mlflow_uri = os.getenv("MLFLOW_URI")
26 artifact_location = os.getenv("ARTIFACT_LOCATION")
27 mlflow.set_tracking_uri(mlflow_uri)
28 # create experiment if not exist
29 if mlflow.get_experiment_by_name(name) is None:
30     mlflow.create_experiment(
31         name, artifact_location=artifact_location
32     )
33 mlflow.set_experiment(name)
34 mlflow.sklearn.log_model(clf, artifact_path="sk_models")
35 logger.info("All done!")
36
```

# Load Model for inference

```
44 def create_app():
45     import mlflow
46     import os
47
48     # get model name and version
49     model_name = os.getenv("MLFLOW_MODEL_NAME")
50     model_version = os.getenv("MLFLOW_MODEL_VERSION")
51     mlflow_uri = os.getenv("MLFLOW_URI")
52     # load model
53     logger.info(f"Loading {model_name} | {model_version}")
54     mlflow.set_tracking_uri(mlflow_uri)
55     model = mlflow.sklearn.load_model(
56         f"models:/{{model_name}}/{{model_version}}"
57     )
58     # init flask
59     app = Flask(__name__)
60     api = Api(app)
61     api.add_resource(
62         IrisRandomForestPredictor,
63         "/predict",
64         resource_class_kwargs={"model": model},
65     )
66     logger.info("Service ready!")
67     return app
```

```
1  apiVersion: serving.knative.dev/v1
2  kind: Service
3  metadata:
4    name: iris-rf
5    namespace: knative-demo
6  spec:
7    template:
8      spec:
9        containers:
10       - name: iris-rf
11         image: boxkey/iris-rf:beta
12         ports:
13           - containerPort: 5000
14         env:
15           - name: MLFLOW_MODEL_NAME
16             value: "Iris_Random_Forest"
17           - name: MLFLOW_MODEL_VERSION
18             value: "1"
19           - name: MLFLOW_URI
20             valueFrom:
21               secretKeyRef:
22                 name: mlflow
23                 key: mlflowuri
24                 optional: false
```

# Kubernetes Setup - \$30/month



Thank you!