

# What's News on Kubeflow V1.5

葉信和 / Hsin-Ho Yeh  
Software Engineer / CEO @ 信誠金融科技  
hsinho.yeh@footprint-ai.com

# Download Slides

<https://reurl.cc/Lm5lO4>

<https://github.com/FootprintAI/talks/tree/main/slides>

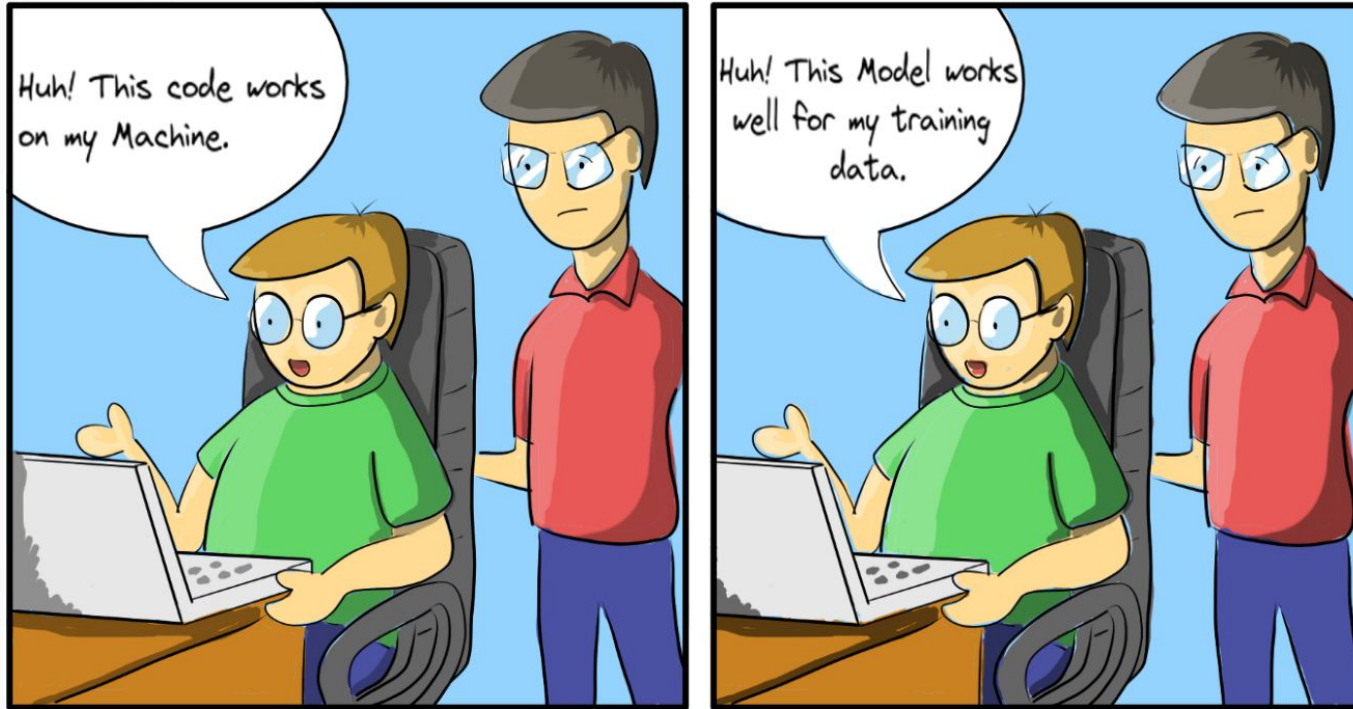


# About me


- 2020 - Present at 信誠金融科技
  - Shrimping: A data-sharing platform
    - <https://get-shrimping.footprint-ai.com>
  - Tintin: a machine learning platform for everyone
    - <https://get-tintin.footprint-ai.com>
- 2016 - 2020 at IglooInsure (16M+ in series A+ 2020)
  - Provide digital insurance for e-economic world
  - Funded in KUL, Headquartered in Singapore
  - First employee/ Engineering Lead / Regional Head/ Chief Engineer
- 2013 - 2016 at Studio Engineering @ hTC
  - Principal Engineer on Cloud Infrastructure Team
- 2009 - 2012 at IIS @ Academia Sinica
  - Computer vision, pattern recognition, and data mining
- CS@CCU, CS@NCKU alumni



A common scenario that we both experienced.



What is deployment automation in Machine learning?



DevOps + ML  
= MLOps

MLOps is the process of taking an experimental Machine Learning model into a production system by including continuous development practice of DevOps in the software field.

Ref: <https://en.wikipedia.org/wiki/MLOps>

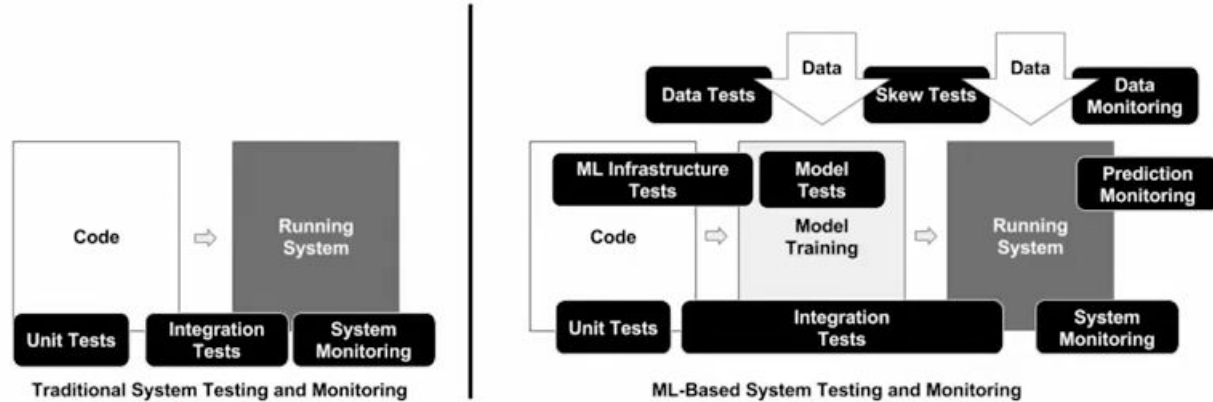
Source: <https://www.kubeflow.org/>

Building & deploying real-world ML application is *hard* and *costly* because of *lack of tooling* that covers end-to-end ML development & deployment

- CloudNext'19

# How Involving Machine Learning model could change the current software design?

## Traditional vs. ML infused systems



ML introduces two new assets into the software development lifecycle – **data** and **models**.

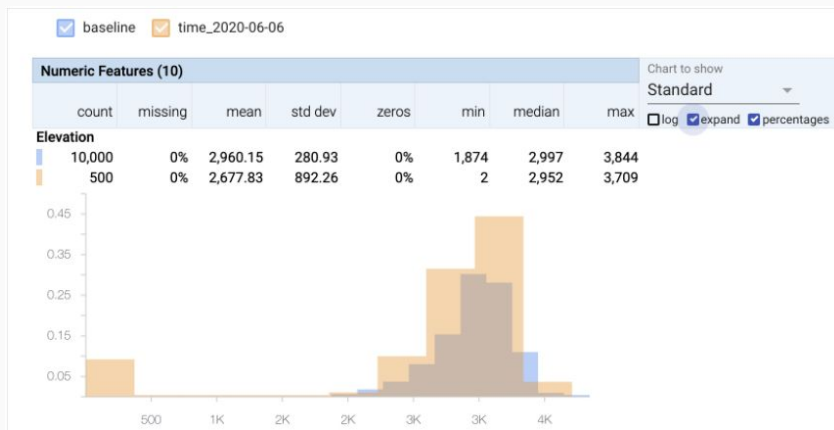
# Why we should care about drifting?

- Data drifting

- A skew grows between training data and serving data.
- The discrepancies between training data and serving data can usually be classified as schema skews or distribution skews

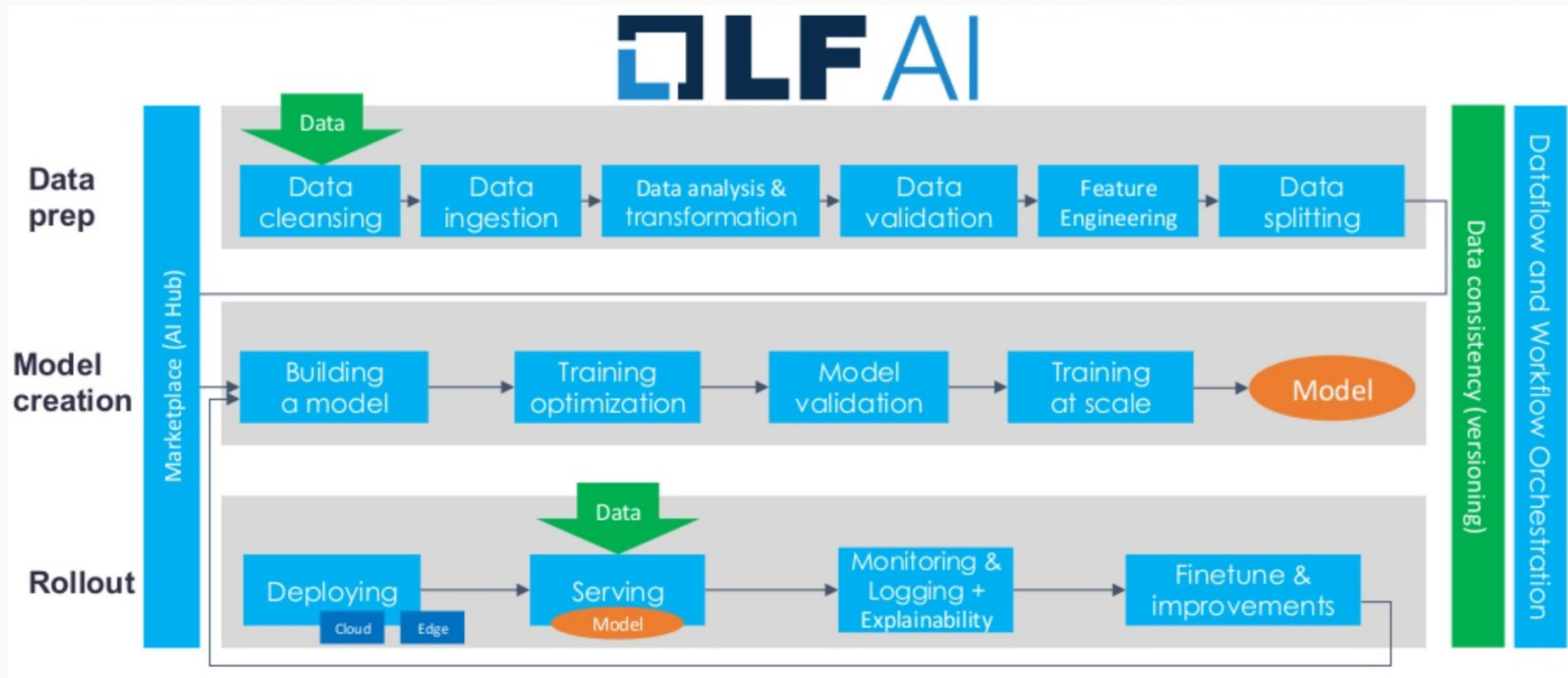
- Concept drifting

- The interpretation of the relationship between the input predictors and the target feature evolves





# Real-world Machine Learning Application - End-to-End ML LifeCycle



Source: <https://www.slideshare.net/AnimeshSingh/advanced-model-inferencing-leveraging-kubeflow-serving-knative-and-istio-196096385>

# Why machine learning on Kubernetes?

- Composability
  - Each stage are independent systems and are able to compose together
- Portability
  - Dev/Staging/Prod
  - Laptop/Edge/Cloud environment
- Scalability
  - Hyperparameter tuning, production workloads

# History Of Kubernetes

- Borg: the predecessor to Kubernetes
  - Google revealed the first time of its detail in an academic research paper, describing a “cluster manager that runs hundreds of thousands of jobs, from many thousands of different applications, across a number of clusters each with up to tens of thousands of machines.”[1]
  - A in-house cluster manager system inside Google for running every google services including Gmail, Google Maps, Google Docs...[2]
  - In a scale with ‘over 2 billion containers per week` [3]
- The very first version of Kubernetes was released in 2015
- The latest version is v1.23, released at 2022.



[1] <https://research.google/pubs/pub43438/>

[2] <https://www.wired.com/2016/04/want-build-empire-like-googles-os/>

[3] <https://cloud.redhat.com/blog/building-kubernetes-bringing-google-scale-container-orchestration-to-the-enterprise>

**Oh, you want to use ML on K8s?**

**Before that, can you become an expert in:**

- Containers
- Packaging
- Kubernetes service endpoints
- Persistent volumes
- Scaling
- Immutable deployments
- GPUs, Drivers & the GPL
- Cloud APIs
- DevOps
- ...



A close-up photograph of a person's hand holding a white marker, writing on a whiteboard. The background is blurred, showing what appears to be a workshop or office environment with some equipment and lights.

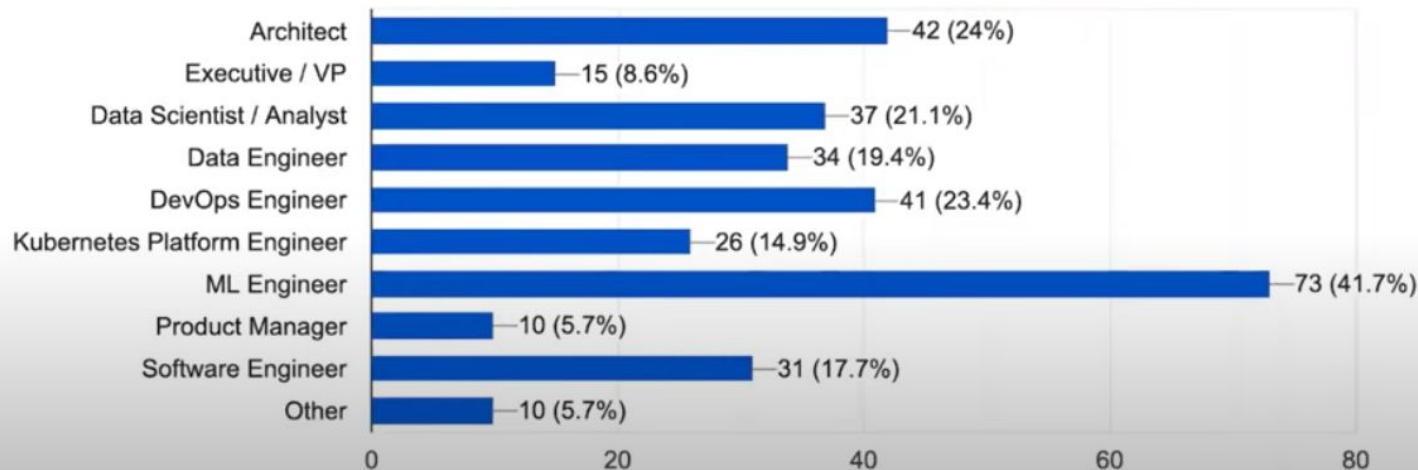
**Kubernetes + ML  
= Kubeflow**

The Kubeflow project is dedicated to making deployments of machine learning (ML) workflows on Kubernetes simple, portable and scalable.

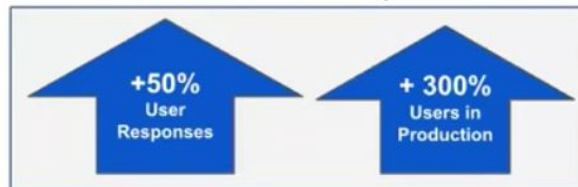
## '21 Kubeflow User Survey Results

Please identify your title?

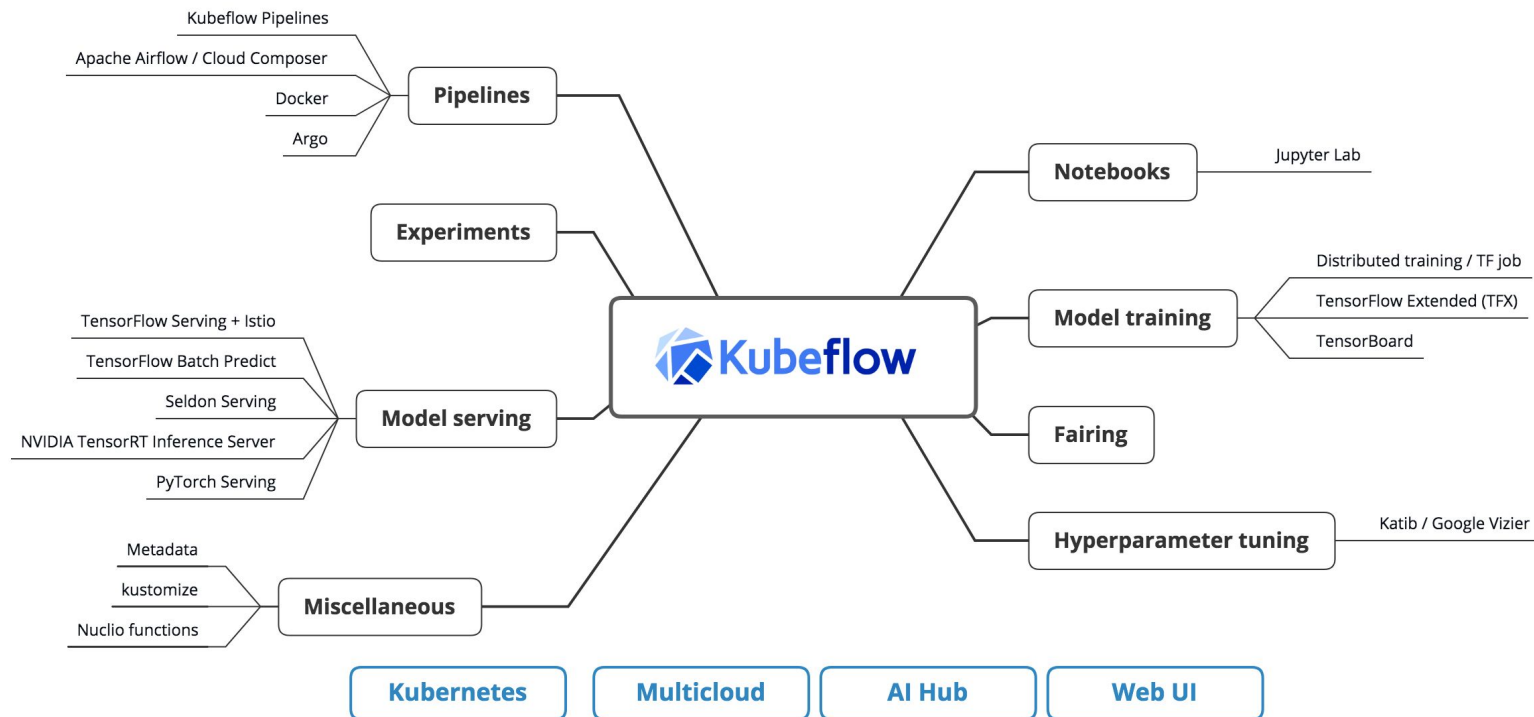
175 responses



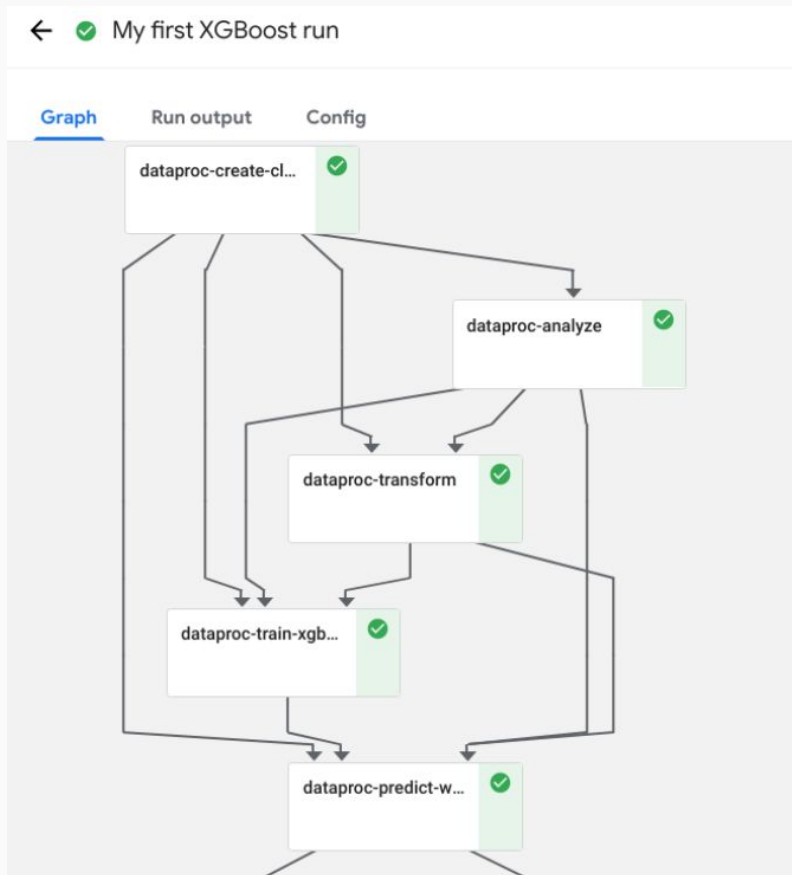
'21 vs '20 Survey



# Architectures



# Kubeflow Pipelines

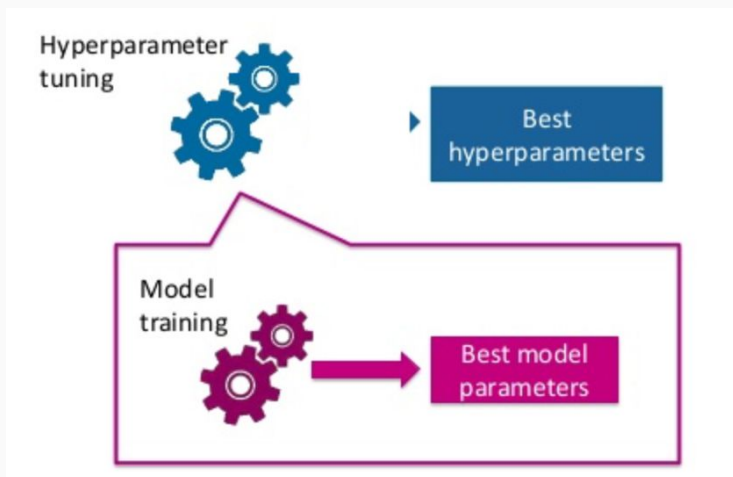


Source:  
<https://www.kubeflow.org/docs/pipelines/overview/pipelines-overview/>



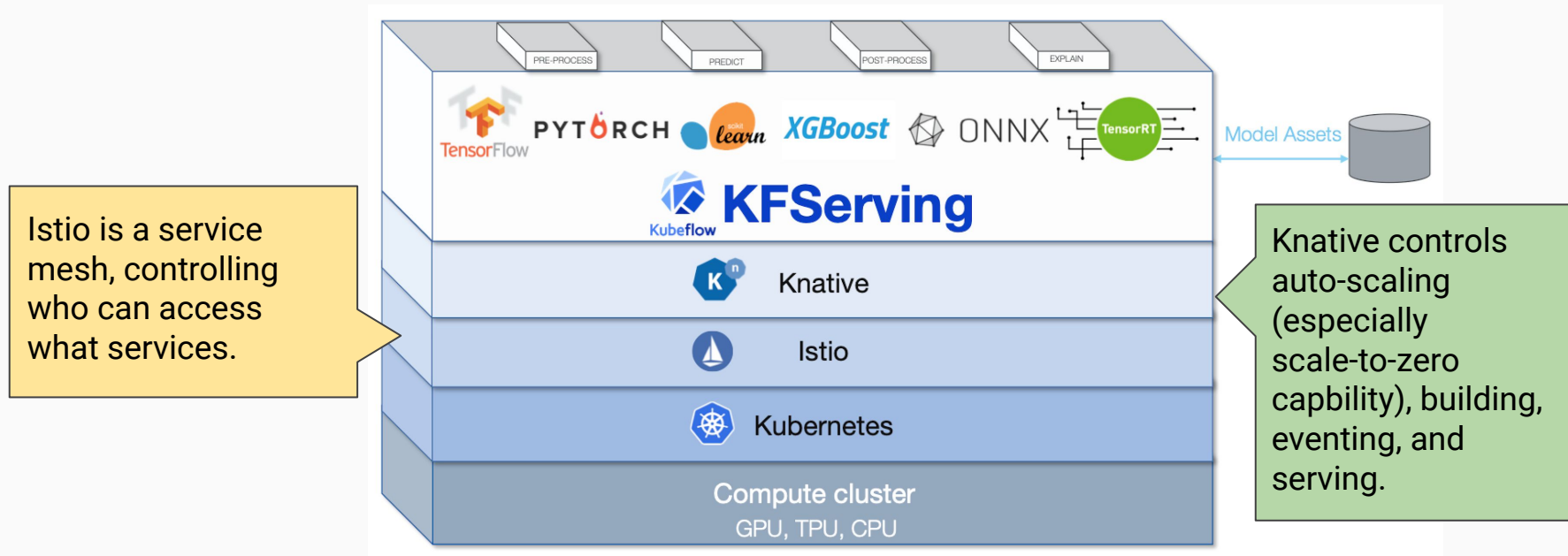
# Hyperparameter tuning

- Model Parameter vs Hyperparameter
  - Model parameters that will learn on its own during training process by the ML model, ex: weights and biases for a classifier.
  - Hyperparameter that directly control the behavior of training algorithm.



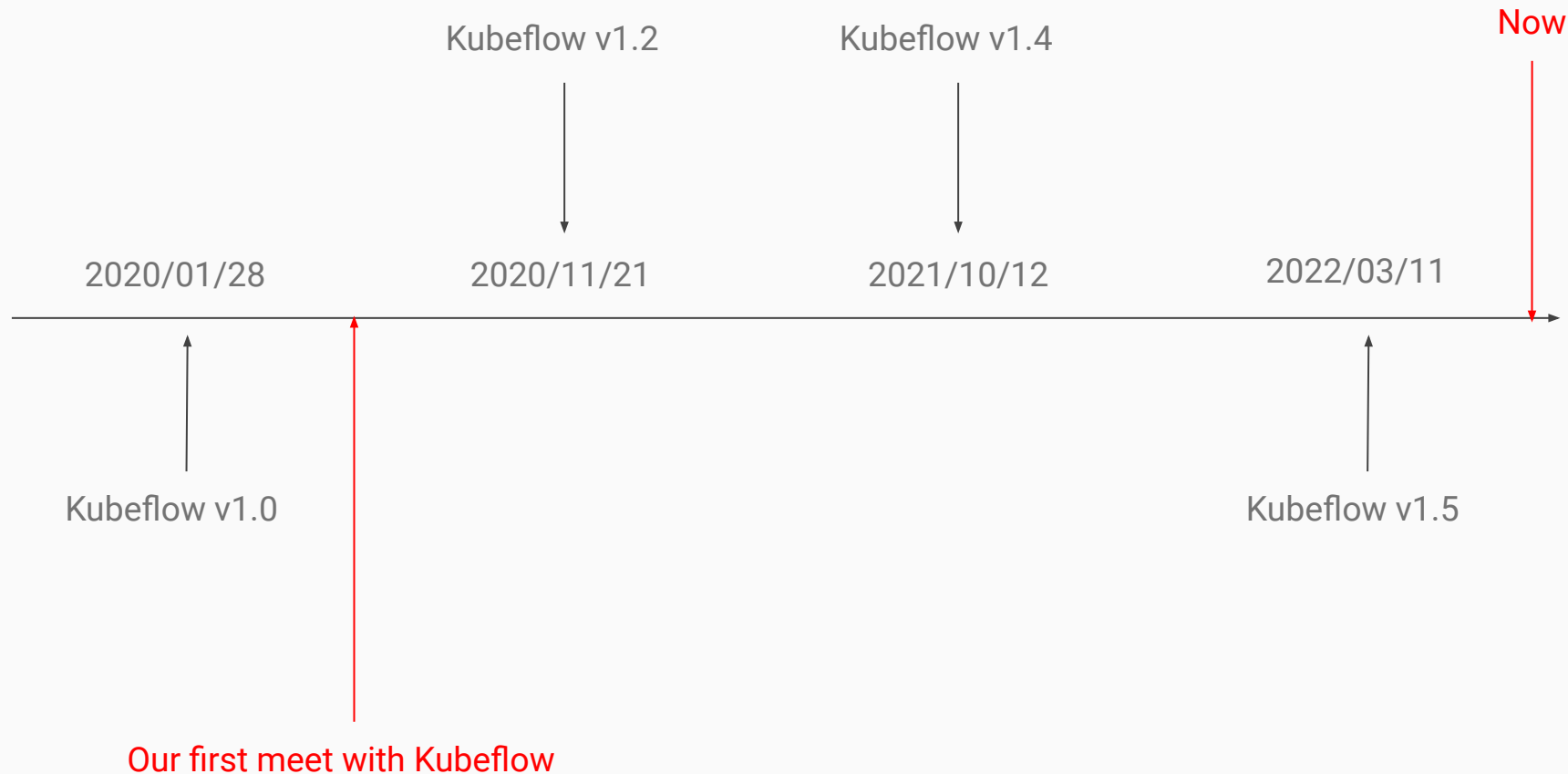
Source: <https://towardsdatascience.com/understanding-hyperparameters-and-its-optimisation-techniques-f0debba07568>

# KServe



Source: <https://towardsdatascience.com/understanding-hyperparameters-and-its-optimisation-techniques-f0debba07568>

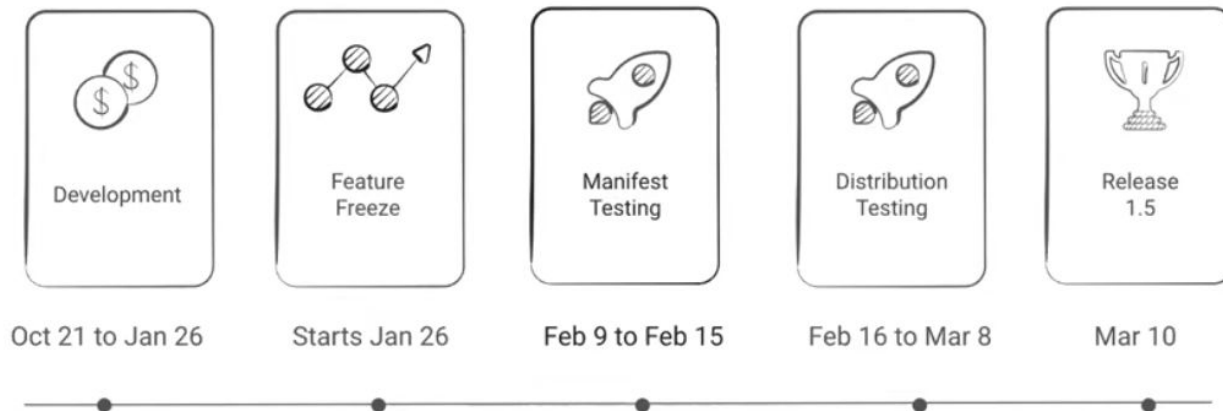
# Kubeflow Release timeline



# Kubeflow Release Cycle




## Kubeflow 1.5 Timeline - #538 & 549



Kubeflow 1.6 release: <https://github.com/kubeflow/manifests/issues/2194>

# Lessons learned in the wild - presented at 20210222 (v1.0.2)

- Kubeflow favor GKE, not as friendly on EKS/AKS/On-Prem.
  - Most features/examples are built on top of GCP (Google Cloud Platform)
- ~~Buggy~~Unstable System (Yet, we knew it is in version 1.x now)
  - Early staged components
    - Katib (v0.10.0), kfserving(v0.5.0) are not yet production-ready release.
    - Multi-Tenancy: resource sharing and access-control isolation.
  - Out-of-date dependency
    - Istio: v1.3.1 is using in kubeflow v1.2 and v1.6 will be kubeflow v1.3
- Hidden Cost
  - ~~Linux~~Open source is free if you time have no value. - [Jamie Zawinski](#)
- Off-the-shelf software
  - No way to build better UI/flow to fulfill your business requirement without digging into components
  - Require more knowledge on k8s/devops expertise to be able to master it

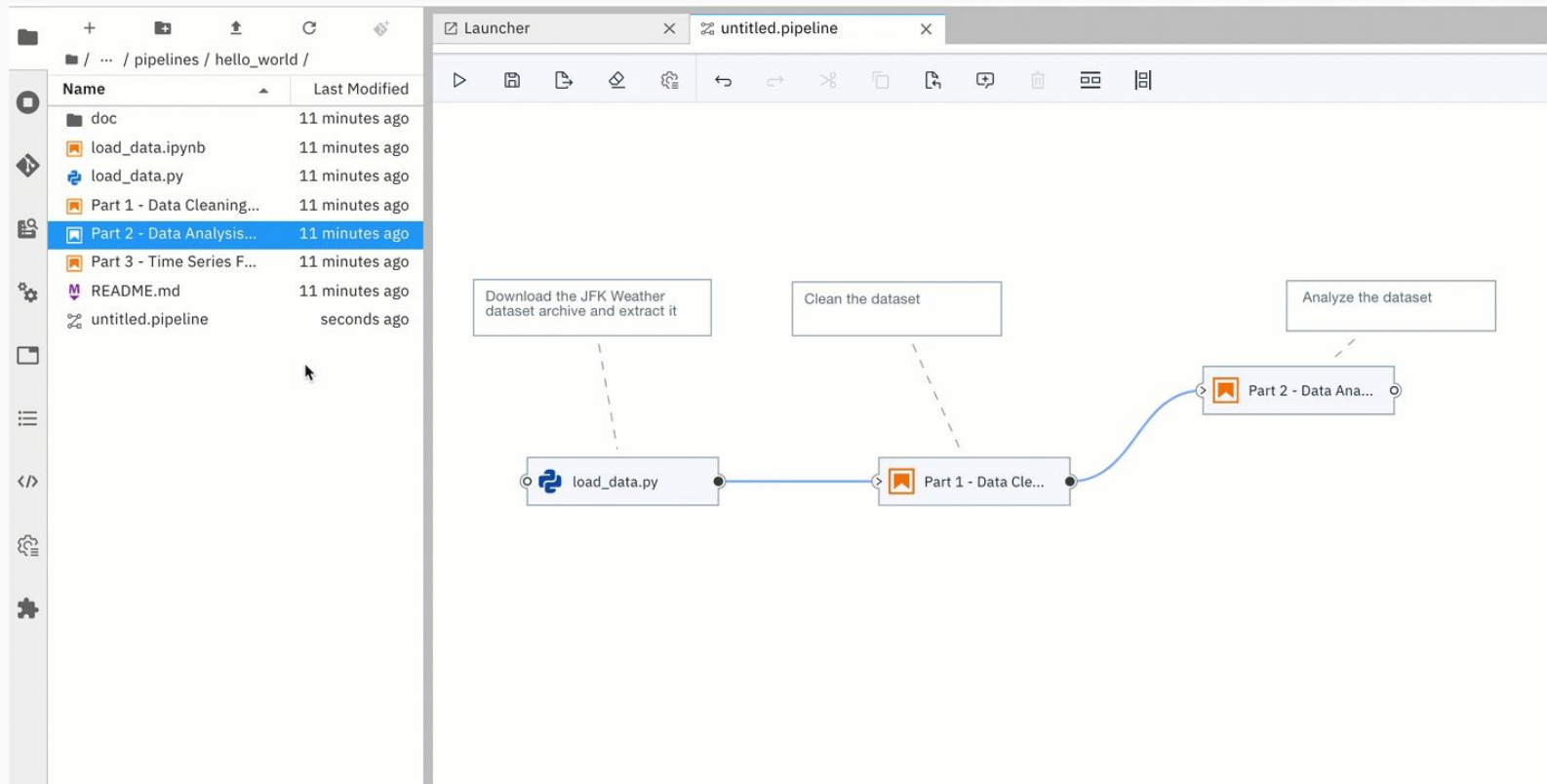
A close-up photograph of a person's hand holding a pen, poised to write on a document. The background is heavily blurred, showing indistinct shapes and colors, suggesting an office or laboratory setting. The lighting is soft, highlighting the texture of the skin and the details of the hand and pen.

What's news?

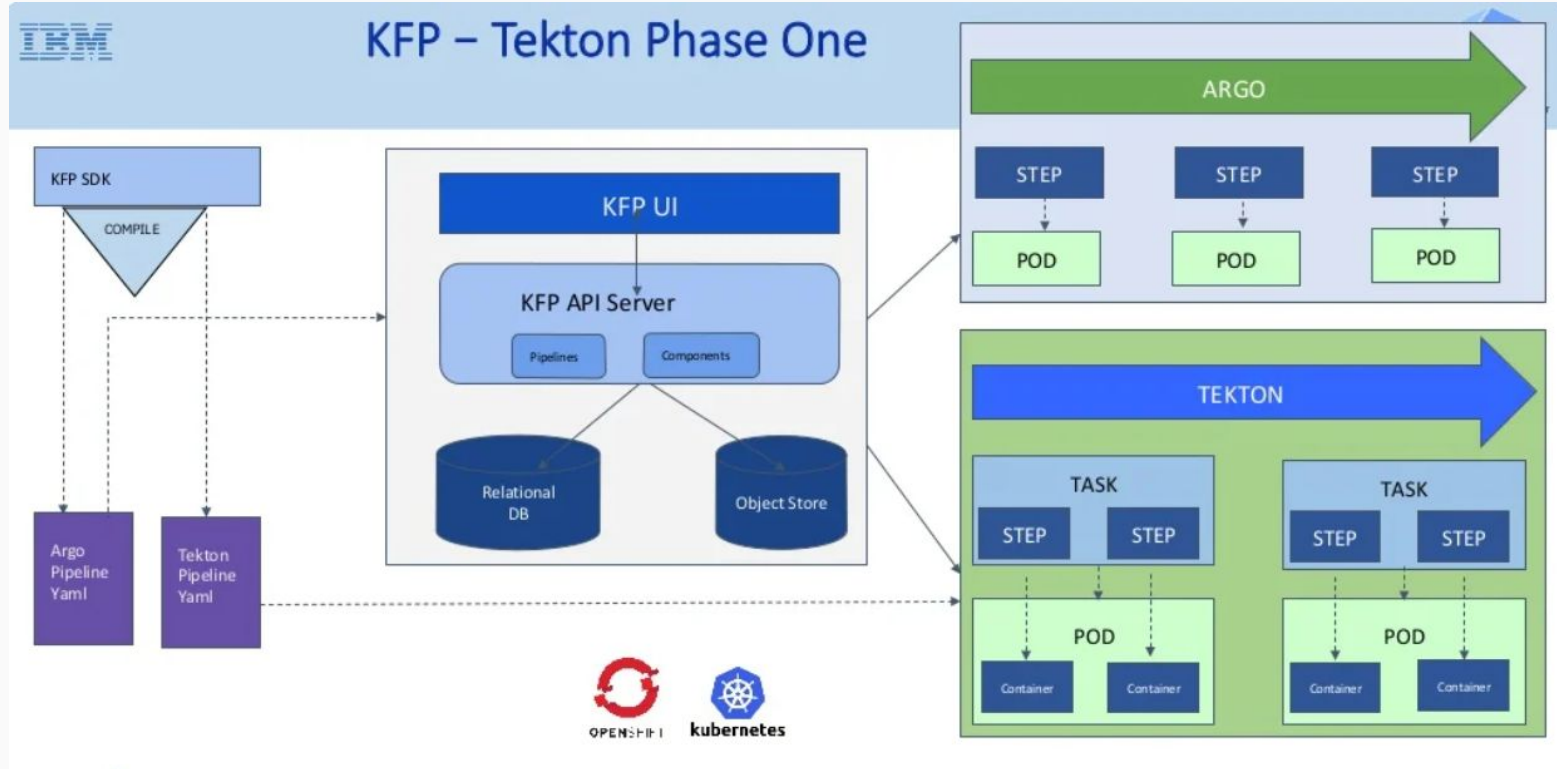
# What's news?

- Elyra: Workflow Visualization
- Tekton Pipeline
- MultiModel Serving in Kserve
- ML UI
- Katib UI

# Elyra: Workflow Visualization

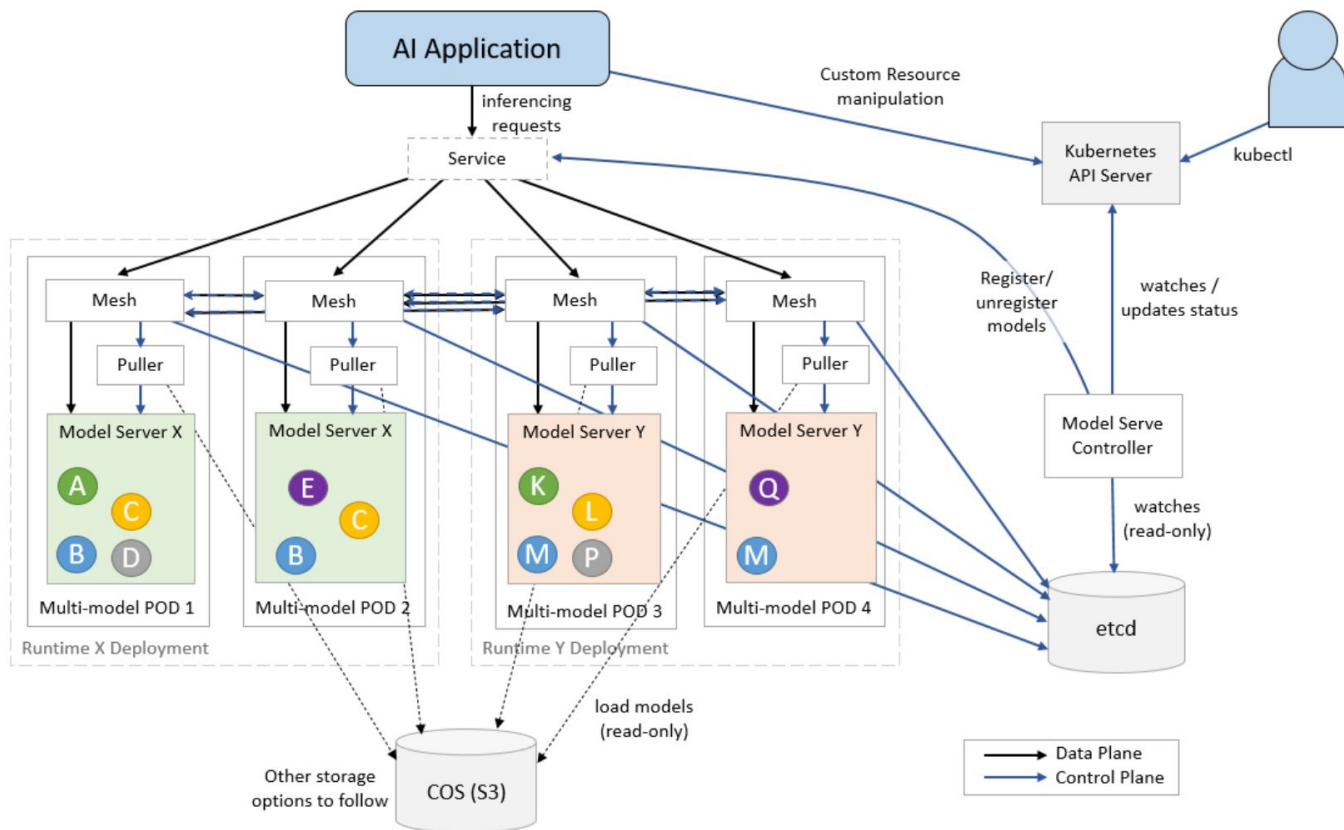


# Tekton Pipeline






# Kserve: MultitModel Serving with Model Mesh



# Kserve: Model UI (1/2)

 **Kubeflow**

Home

Notebooks

Tensorboards

**Models**

Snapshots

Volumes

Experiments (AutoML)

Experiments (KFP)

Pipelines

Runs

Recurring Runs

Artifacts

Executions

kubeflow-user (Owner)

Model server details

DELETE

flowers

OVERVIEWDETAILSMETRICSLOGSYAML

Status

Ready

URL external

http://flowers.kubeflow-user.example.com

URL internal

http://flowers.kubeflow-user.svc.cluster.local/v1/models/flowers:predict

Component

predictor

Storage URI

gs://kfserving-samples/models/tensorflow/flowers

Runtime

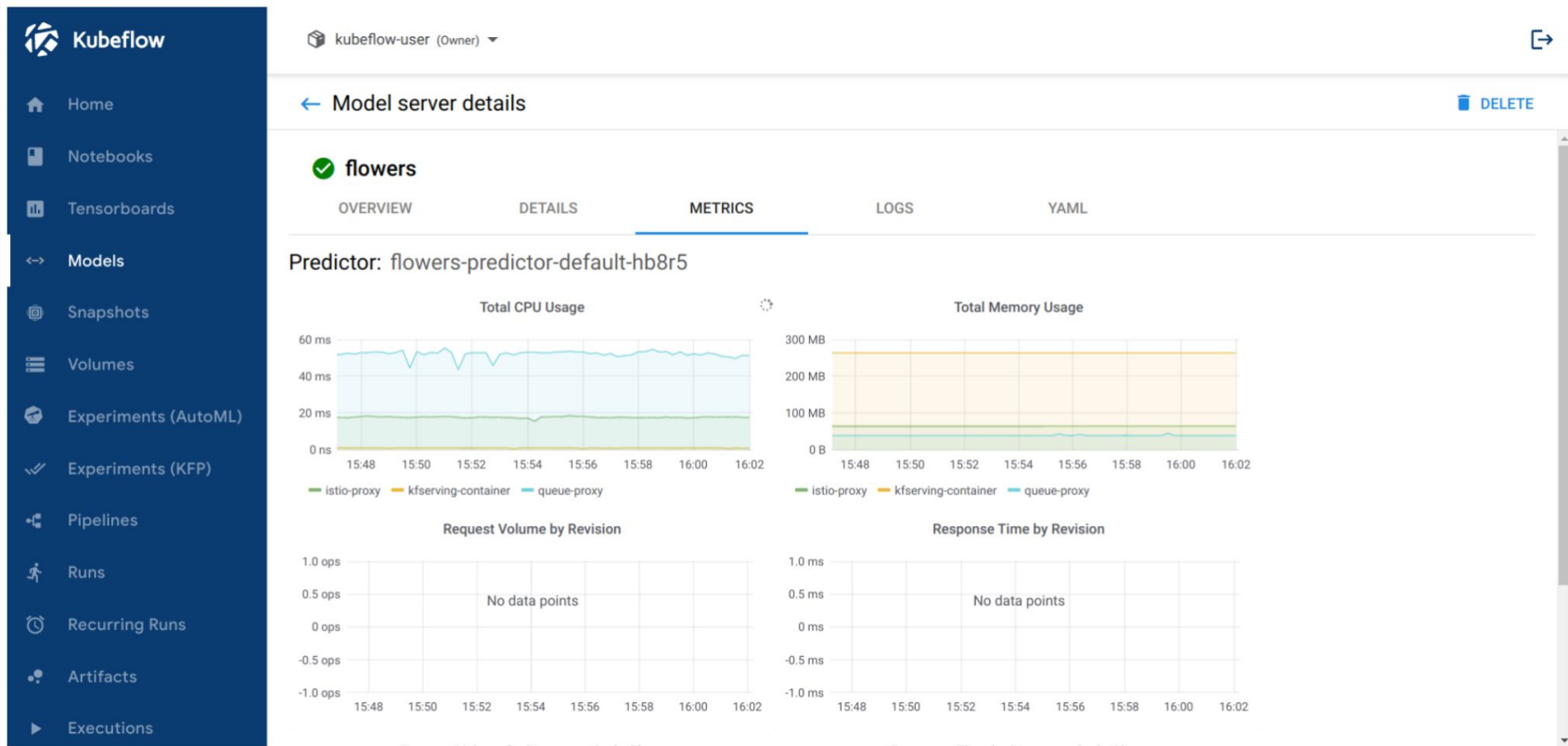
Tensorflow 1.14.0

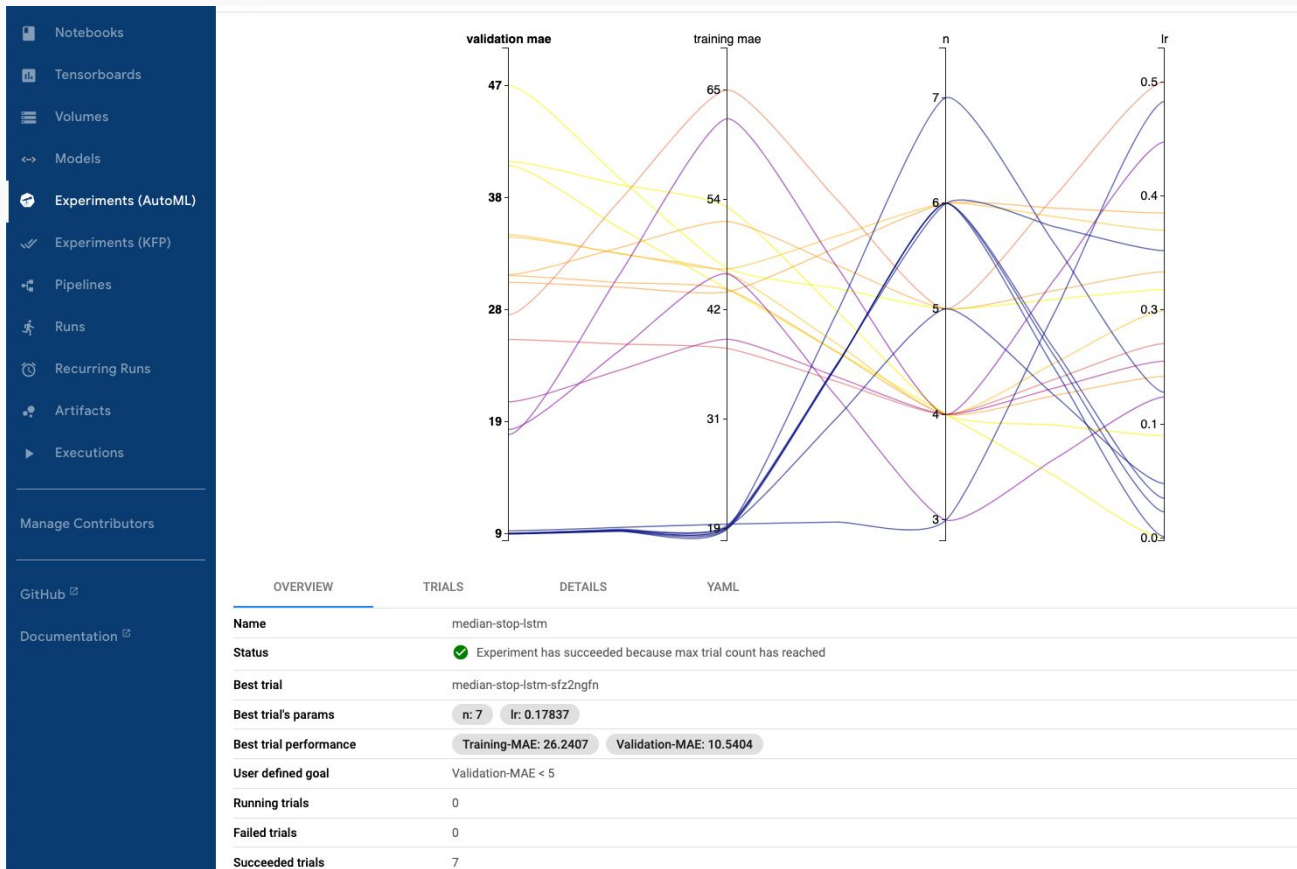
Protocol Version

InferenceService Conditions

Status	Type	Last Transition Time	Reason	Message
✓	IngressReady	23 minutes ago		
✓	PredictorConfigurationReady	23 minutes ago		

# Kserve: Model UI (2/2)



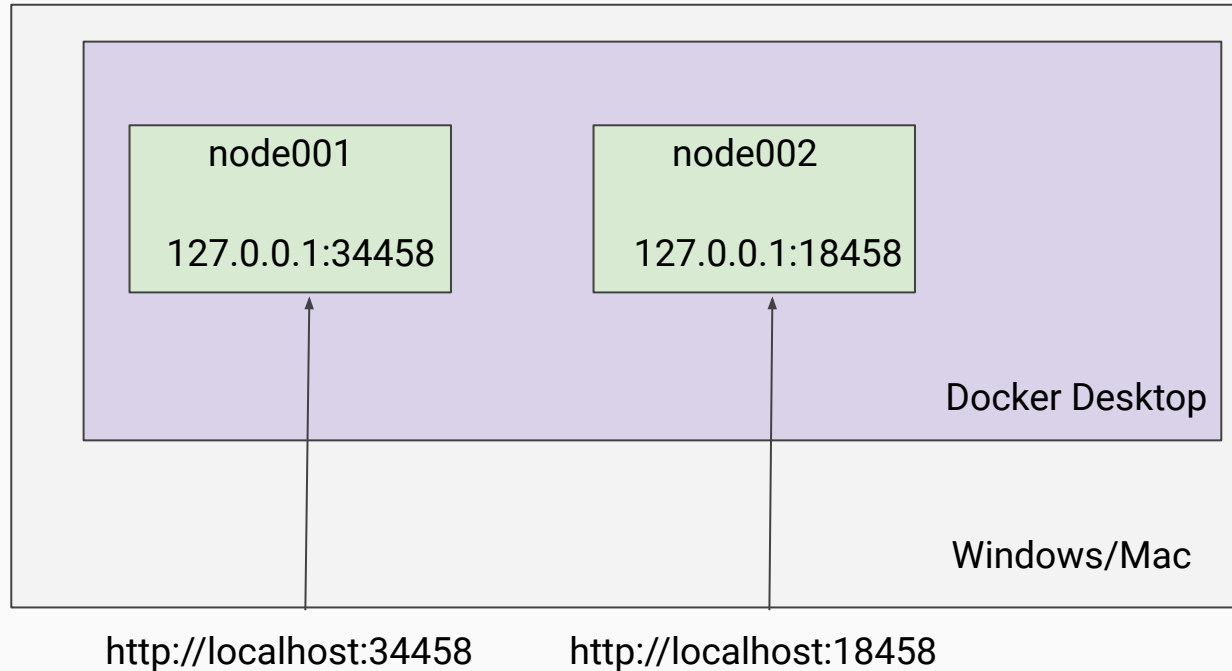


multikf(馬蒂庫夫)

Our open-sourced project  
for one-clicked running  
multiple Kubeflow instances  
within the single host  
machine.

# multikf: One-click Installation

- Multikf: <https://github.com/footprintai/multikf>



## multikf: One-click Installation

```
// install dockerd (windows)
https://github.com/FootprintAI/kubeflow-workshop/blob/main/install/windows/dockerd.bat.md

// install multikf (windows)
wget https://github.com/FootprintAI/multikf/raw/main/build/multikf.windows.exe
chmod +x multikf.windows.exe

// add an instances with port 80/443 exported
./multikf.windows.exe add node002 --export_ports 80:80,443:443

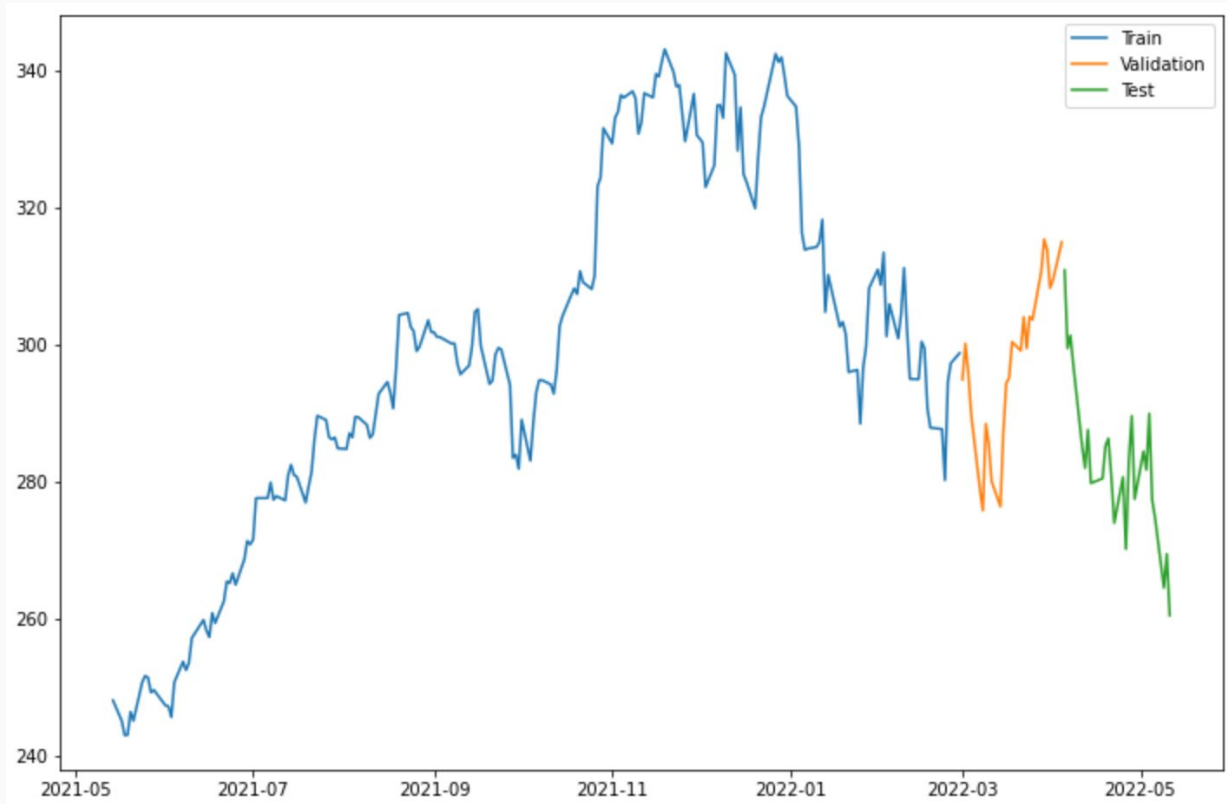
// connect kubeflow
./multikf.windows.exe connect kubeflow node002
```

# Katib Case Study

- How to predict the stock market price of tomorrow?
- How to determine the parameter is good enough?



# CastStudy: Data overview



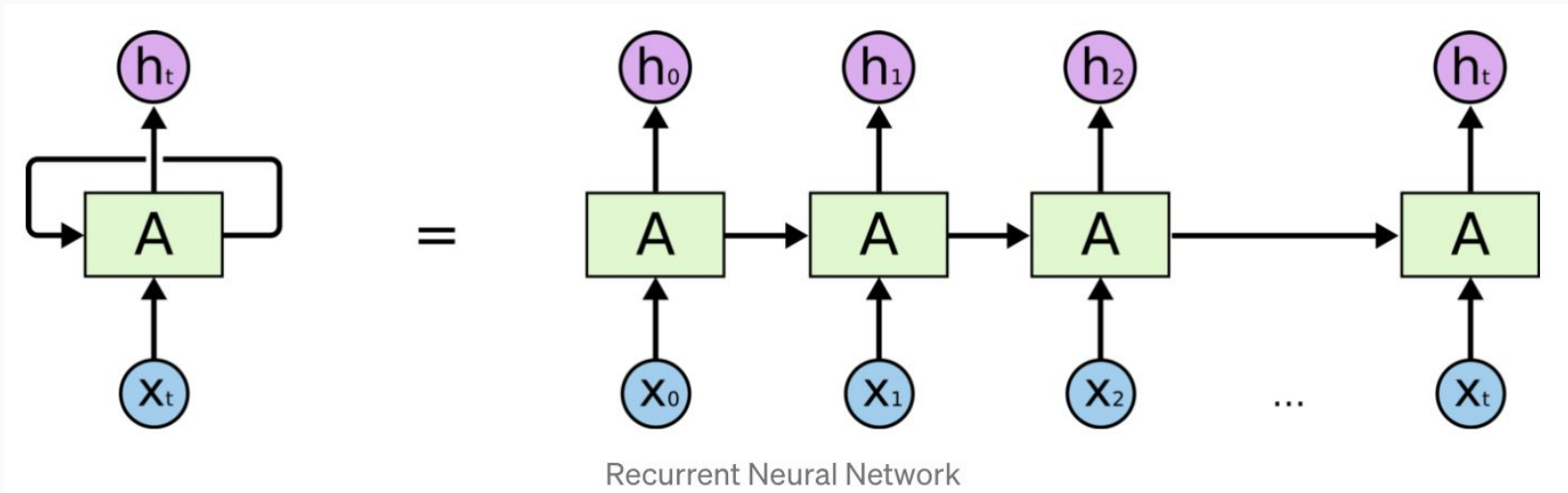
## What sequences means in a sentence?

Input: I am a good guy

t0	I
t1	am
t2	a
t3	good
t4	guy
t5	?

# What Is RNN (Recurrent Neural Network)?

Recurrent Neural Network (RNN) takes decisions on CURRENT ( $X_t$ ) and PREVIOUS ( $X_{t-1}$ ) inputs. Especially useful in topics including machine translation, speech recognition.

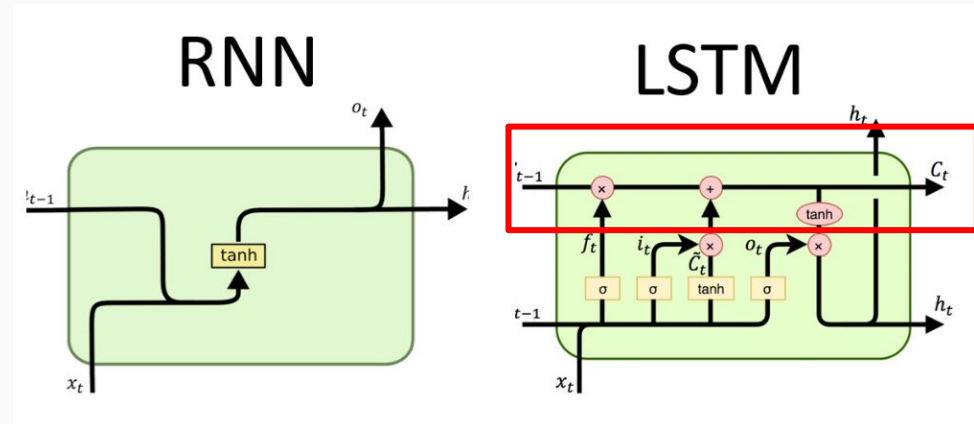


# What Is LSTM (Long-short term memory)?

RNN only remember the latest things from  $X$  and it didn't remember(no memory) anything before at the beginning.

LSTM provides an information highway to let the neuron to selectively choose

1. forget from its memory (focus on the current inputs)
2. Listens to what information it added into memory (though information highway)



# CastStudy: Build a tensorflow model with LSTM

```
# Now we build a tensorflow model with LSTM
# the network is not something fancy, it is just a common way to build the model
# And you can also find a better model online.

# Note the input tensor size should have equal length with the windows size.
# In this example, we use windows size (n=3), so the input tensor is layers.Input((3, 1)

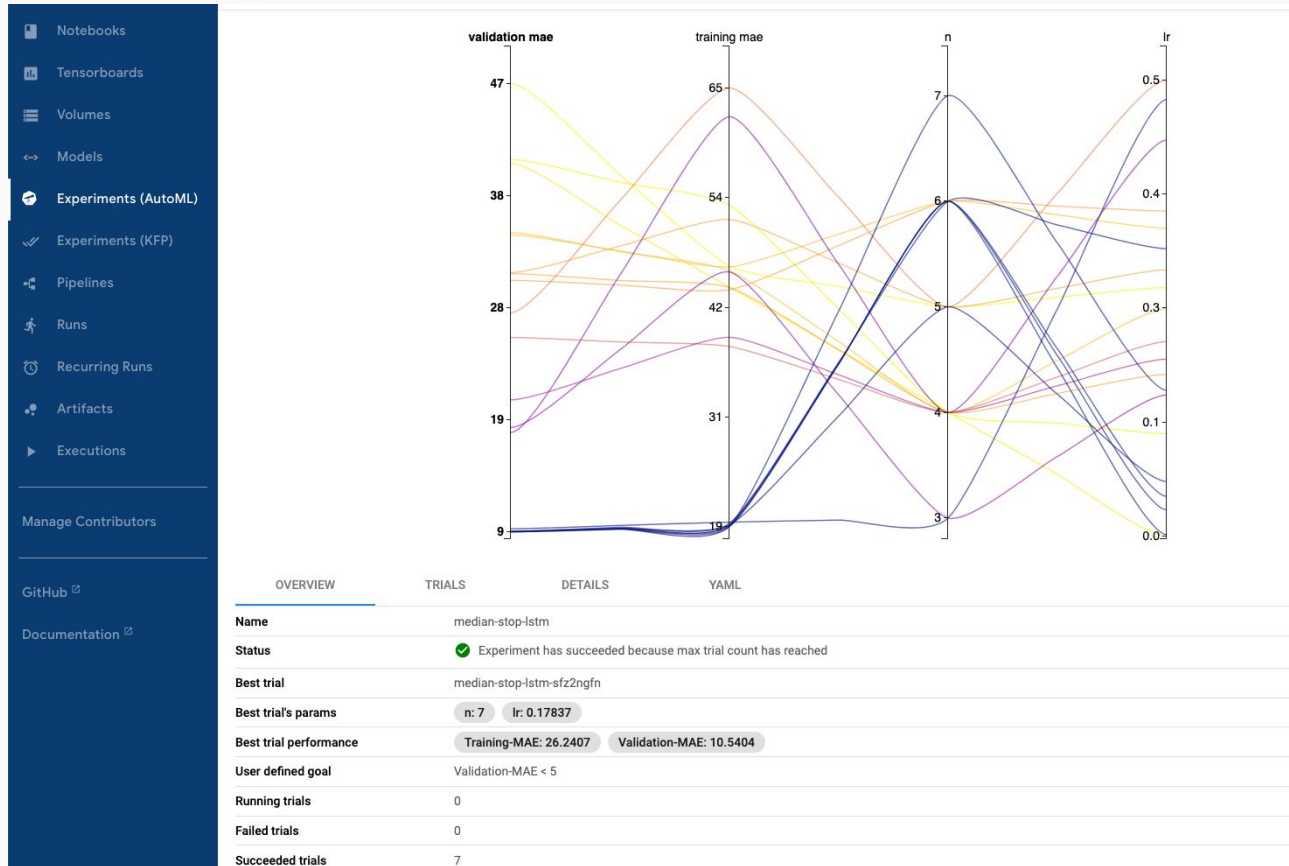
from tensorflow.keras.models import Sequential
from tensorflow.keras.optimizers import Adam
from tensorflow.keras import layers

model = Sequential([layers.Input((3, 1)),
                    layers.LSTM(64),
                    layers.Dense(32, activation='relu'),
                    layers.Dense(32, activation='relu'),
                    layers.Dense(1)])

model.compile(loss='mse',
              optimizer=Adam(learning_rate=0.001),
              metrics=['mean_absolute_error'])

model.fit(X_train, y_train, validation_data=(X_val, y_val), epochs=100)
```


# CastStudy: Optimize model parameters with Katib



Ref: <https://github.com/FootprintAI/kubeflow-workshop/blob/main/tutorials/stockprice-with-lstm/1.visualize-and-build-stockprice-model.ipynb>

## Conclusion

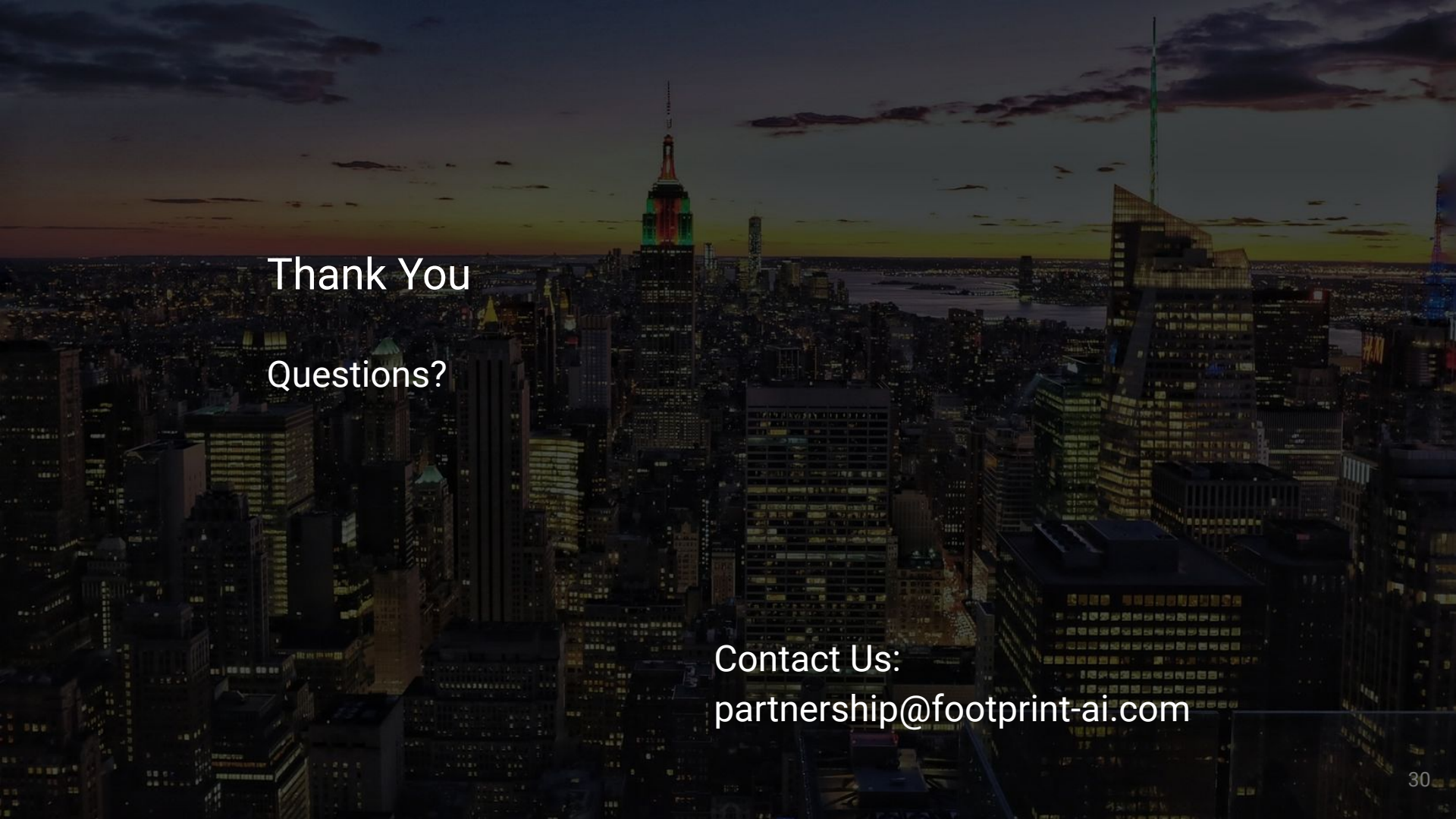
- The tooling and features on Kubeflow is more user-friendly compared to v1.0.2, all of this is contributed to the open source community and release automation.



***“The Best Engineers  
Are Lazy”***

-Ancient Engineering Proverb



An aerial photograph of the New York City skyline at dusk. The sky is a mix of dark purple, blue, and orange. The city is densely packed with skyscrapers, many of which are illuminated with their interior lights. The Empire State Building is prominent in the center, with its top lit in red and green. The Hudson River is visible on the right side of the image.

Thank You  
Questions?

Contact Us:  
[partnership@footprint-ai.com](mailto:partnership@footprint-ai.com)