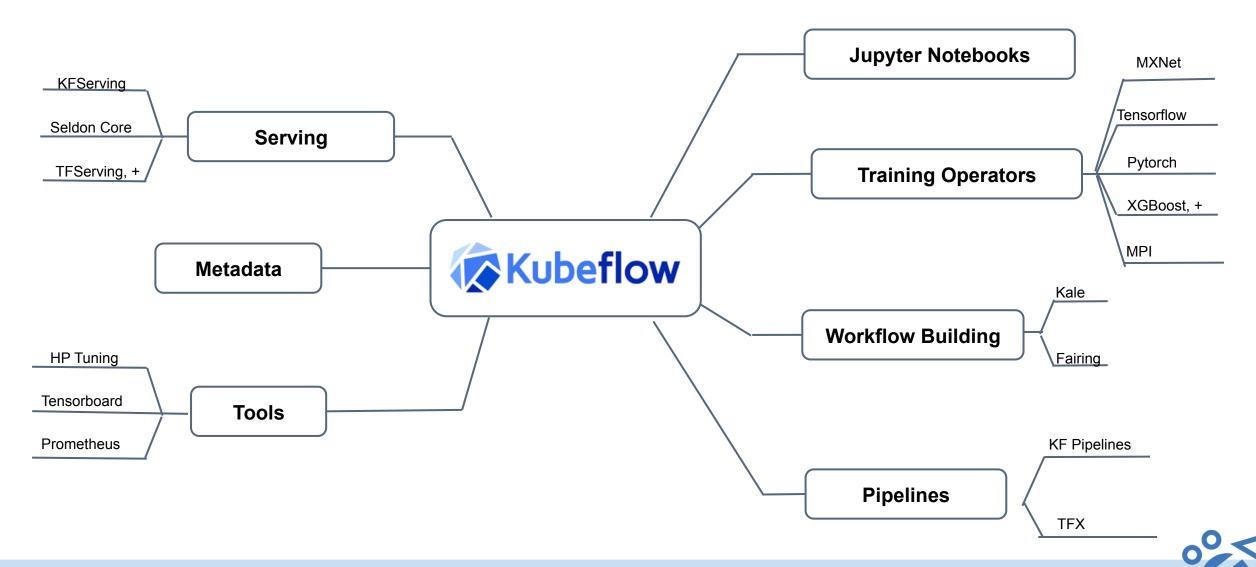


IBM KFServing

Animesh Singh, Tommy Li







Production Model Serving? How hard could it be?



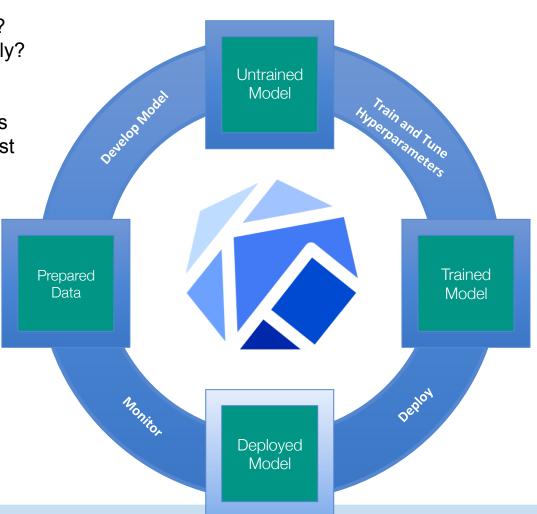


Is the model over or under scaled?
Are resources being used efficiently?

☐ Monitoring:

Are the endpoints healthy? What is the performance profile and request trace?

- ☐ Rollouts:
 - Is this rollout safe? How do I roll back? Can I test a change without swapping traffic?
- □ Protocol Standards: How do I make a prediction? GRPC? HTTP? Kafka?



- ☐ How do I handle batch predictions?
- ☐ How do I leverage standardized Data Plane protocol so that I can move my model across MLServing platforms?
- ☐ Frameworks:

 How do I serve on Tensorflow?

 XGBoost? Scikit Learn? Pytorch?

 Custom Code?
- ☐ Features: How do I explain the predictions? What about detecting outliers and skew? Bias detection? Adversarial Detection?
- ☐ How do I wire up custom pre and post processing





Experts fragmented across industry



- Seldon Core was pioneering Graph Inferencing.
- IBM and Bloomberg were exploring serverless ML lambdas. IBM gave a talk on the ML Serving with Knative at last KubeCon in Seattle
- Google had built a common Tensorflow HTTP API for models.
- Microsoft Kubernetizing their Azure ML Stack















Putting the pieces together



- Kubeflow created the conditions for collaboration.
- A promise of open code and open community.
- Shared responsibilities and expertise across multiple companies.
- Diverse requirements from different customer segments







Bloomberg









Introducing KFServing





KFServing



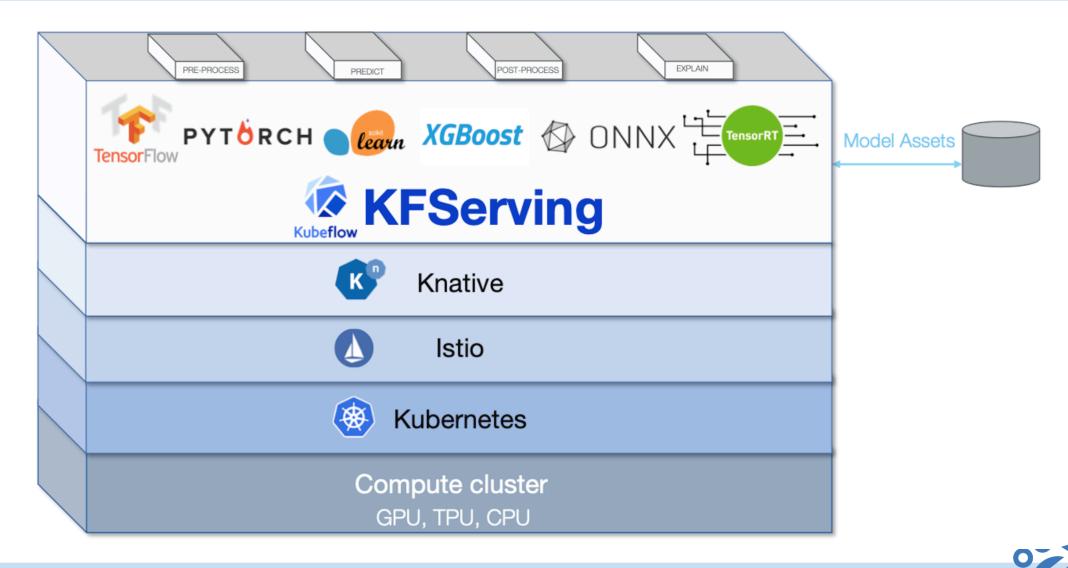
- Founded by Google, Seldon, IBM, Bloomberg and Microsoft
- Part of the Kubeflow project
- Focus on 80% use cases single model rollout and update
- Kfserving 1.0 goals:
 - Serverless ML Inference
 - Canary rollouts
 - Model Explanations
 - Optional Pre/Post processing





KFServing Stack





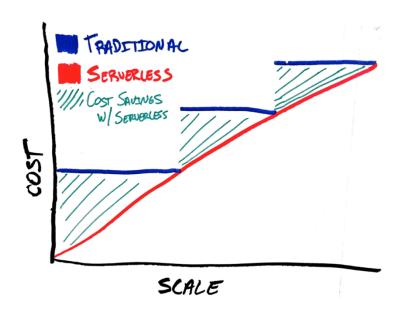
IBM

KNative





IBM is 2nd largest contributor



Knative provides a set of building blocks that enable declarative, container-based, serverless workloads on Kubernetes. Knative Serving provides primitives for serving platforms such as:

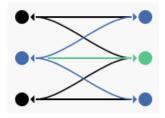
- Event triggered functions on Kubernetes
- Scale to and from zero
- Queue based autoscaling for GPUs and TPUs. KNative autoscaling by default provides inflight requests per pod
- Traditional CPU autoscaling if desired. Traditional scaling hard for disparate devices (GPU, CPU, TPU)



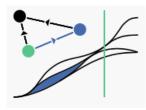
IBM Istio



An open service mesh platform to **connect**, **observe**, **secure**, and **control** microservices. Founded by Google, IBM and Lyft. IBM is the 2nd largest contributor



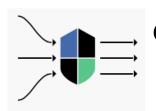
Connect: Traffic Control, Discovery, Load Balancing, Resiliency



Observe: Metrics, Logging, Tracing



Secure: Encryption (TLS), Authentication, and Authorization of service-to-service communication



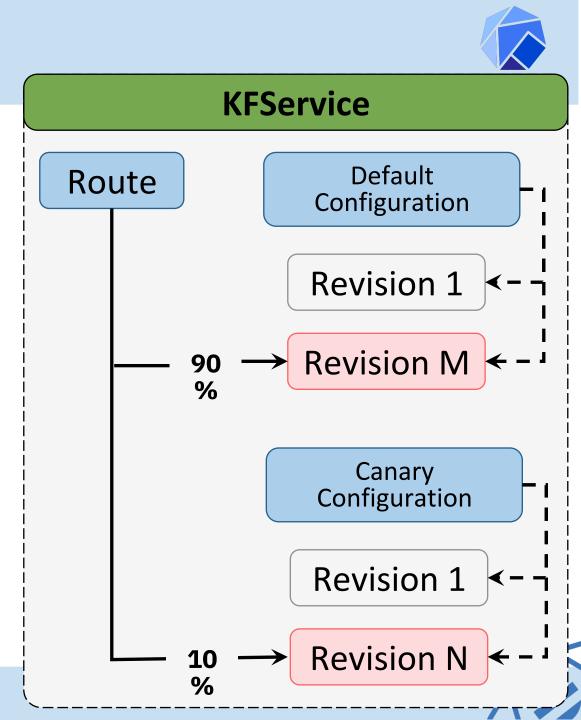
Control: Policy Enforcement



KFServing: Default and Canary Configurations

Manages the hosting aspects of your models

- InferenceService manages the lifecycle of models
- Configuration manages history of model deployments. Two configurations for default and canary.
- Revision A snapshot of your model version
- Route Endpoint and network traffic management





IBM Supported Frameworks, Components and **Storage Subsystems**



Model Servers

- TensorFlow
- Nvidia TRTIS
- PyTorch
- XGBoost
- SKLearn
- ONNX

Components:

- Predictor, Explainer, Transformer (pre-processor, post-processor)

Storage

- AWS/S3
- GCS
- Azure Blob
- PVC

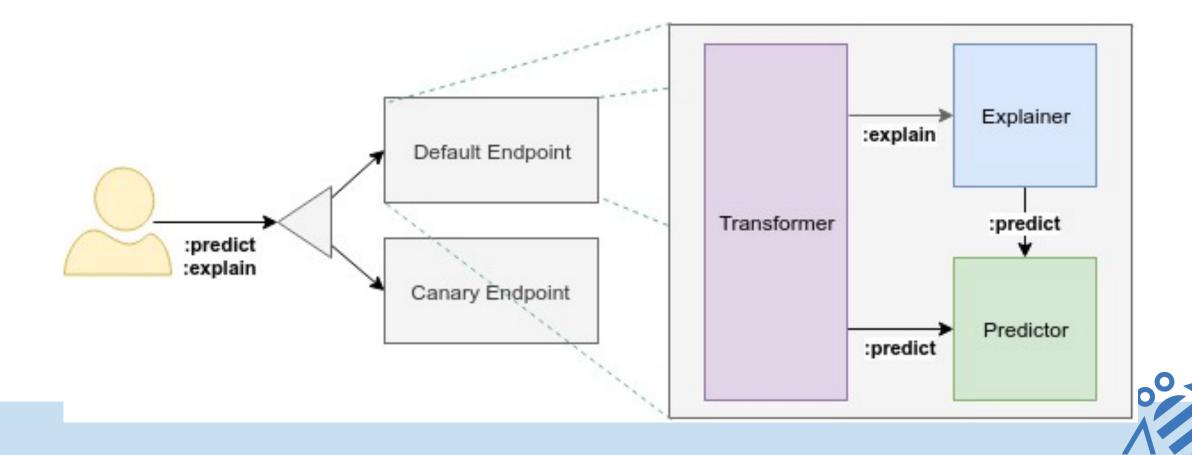




Inference Service Control Plane



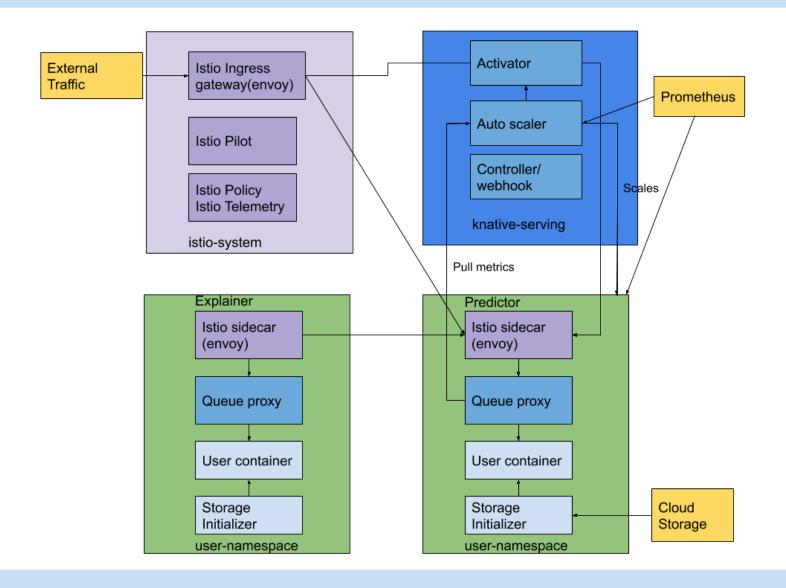
The InferenceService architecture consists of a static graph of components which coordinate requests for a single model. Advanced features such as Ensembling, A/B testing, and Multi-Arm-Bandits should compose InferenceServices together.





KFServing Deployment View







KFServing Data Plane Unification



- Today's popular model servers, such as TFServing, ONNX, Seldon, TRTIS, all
 communicate using similar but non-interoperable HTTP/gRPC protocol
- KFServing v1 data plane protocol uses TFServing compatible HTTP API and introduces explain verb to standardize between model servers, punt on v2 for gRPC and performance optimization.





KFServing Data Plane v1 protocol



API	Verb	Path	Payload
List Models	GET	/v1/models	[model_names]
Readiness	GET	/v1/models/ <model_name></model_name>	
Predict	POST	/v1/models/ <model_name>:predict</model_name>	Request: {instances:[]} Response: {predictions:[]}
Explain	POST	/v1/models <model_name>:explain</model_name>	Request: {instances:[]} Response: {predictions:[], explanations:[]}



KFServing Examples



apiVersion: "serving.kubeflow.org/v1alpha1"

kind: "InferenceService"

metadata:

name: "sklearn-iris"

spec: default:

sklearn:

modelUri: "gs://kfserving-samples/models/sklearn/iris"



kind: "InferenceService"

metadata:

name: "flowers-sample"

spec: default:

tensorflow:

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

apiVersion: "serving.kubeflow.org/v1alpha1"

kind: "InferenceService"

metadata:

name: "pytorch-iris"

spec: default:

pytorch:

modelUri: "gs://kfserving-samples/models/pytorch/iris"











Canary/Pinned Examples



```
apiVersion: "serving.kubeflow.org/v1alpha1"
kind: "KFService"
metadata:
name: "my-model"
spec:
default:

# 90% of traffic is sent to this model
tensorflow:
modelUri: "gs://mybucket/mymodel-2"
canaryTrafficPercent: 10
canary:

# 10% of traffic is sent to this model
tensorflow:
```

modelUri: "gs://mybucket/mymodel-3"





Pinned







TensorFlow Serving(TFServing)

- Flexible, high performance serving system for TensorFlow
- https://www.tensorflow.org/tfx/guide/serving
- Stable and Google has been using it since 2016
- Saved model format and graphdef
- Written in C++, support both REST and gRPC
- KFServing allows you to easily spin off an Inference Service with TFServing to serve your tensorflow model on CPU or GPU with serverless features like canary rollout, autoscaling.





IBM Inference Service with Transformer and TFServing



```
apiVersion: serving.kubeflow.org/vlalpha2
kind: InferenceService
metadata:
  name: bert-serving
spec:
  default
                      Pre/Post Processing
    transformer:
      custom:
        container:
          image: bert-transformer:v1
    predictor:
      tensorflow:
        storageUri: s3://examples/bert
        runtimeVersion: 1.14.0-gpu
        resources:
          limits:
            nvidia.com/qpu: 1
                                  Tensorflow Model
                                       Server
```

```
class BertTransformer(kfserving.KFModel):
  def init (self, name):
       super(). init (name)
       self.bert_tokenizer = BertTokenizer(vocab file)
  def preprocess(self, inputs: Dict) -> Dict:
      encoded features = bert tokenizer.encode plus(
       text=text a, text pair=text b)
      return {"input ids": encoded features["input ids"],
         "input mask": encoded features["attention mask"],
         "segment ids": encoded features["segment ids"],
         "label ids": 1}
  def postprocess(self, inputs: Dict) -> Dict:
      return inputs
```



IBM

NVIDIA Triton Inference Server



- NVIDIA's highly-optimized model runtime on GPUs
- https://docs.nvidia.com/deeplearning/sdk/tensorrt-inference-server-guide/docs
- Supports model repository, versioning
- Dynamic batching
- Concurrent model execution
- Supports TensorFlow, PyTorch, ONNX models
- Written in C++, support both REST and gRPC
- TensorRT Optimizer can further bring down the BERT inference latency





Inference Service with Triton Inference Service



```
apiVersion: serving.kubeflow.org/v1alpha2
kind: InferenceService
metadata:
  name: bert-serving
                        Pre/Post Processing
spec:
  default
    transformer:
      custom:
        container:
          image: bert-transformer:v1
          env:
            name: STORAGE URI
            value: s3://examples/bert transformer
    predictor:
      tensorrt:
        storageUri: s3://examples/bert
        runtimeVersion: r20.02
        resources:
          limits:
            nvidia.com/gpu: 1
                                  Triton Inference Server
```

```
infer ctx = InferContext(url, protocol, model name,
model version)
unique ids = np.int32([1])
segment ids = features["segment ids"]
input ids = features["input ids"]
input mask = features["input mask"]
result = infer ctx.run({ 'unique ids' : (unique ids,),
                         'segment ids' : (segment ids,),
                          'input ids' : (input ids,),
                          'input mask' : (input mask,) },
                       { 'end logits' :
InferContext.ResultFormat.RAW,
                         'start logits' :
InferContext.ResultFormat.RAW }, batch size)
```



PyTorch Model Server



- PyTorch model server maintained by KFServing
- https://github.com/kubeflow/kfserving/tree/master/python/pytorchserver
- Implemented in Python with Tornado server
- Loads model state dict and model class python file
- GPU Inference is supported in KFServing 0.3 release
- Alternatively you can export PyTorch model in ONNX format and serve on TensorRT Inference Server or ONNX Runtime Server.





ONNX Runtime Server



 ONNX Runtime is a performance-focused inference engine for ONNX models



- https://github.com/microsoft/onnxruntime
- Supports Tensorflow, Pytorch models which can be converted to ONNX
- Written in C++, support both REST and gRPC
- ONNX Runtime optimized BERT transformer network to further bring down the latency

https://github.com/onnx/tutorials/blob/master/tutorials/Inference-TensorFlow-Bert-Model-for-High-Performance-in-ONNX-Runtime.ipynb



Inference Service with PyTorch/ONNX Runtime



```
apiVersion: serving.kubeflow.org/vlalpha2
kind: InferenceService
metadata:
  name: bert-serving
                        Pre/Post Processing
spec:
  default
    transformer:
      custom:
        container:
          image: bert-transformer:v1
          env:
            name: STORAGE URI
            value: s3://examples/bert transformer
    predictor:
      pytorch:
        storageUri: s3://examples/bert
        runtimeVersion: v0.3.0-qpu
        resources:
          limits:
            nvidia.com/gpu: 1
                                   Pytorch Model Server
```

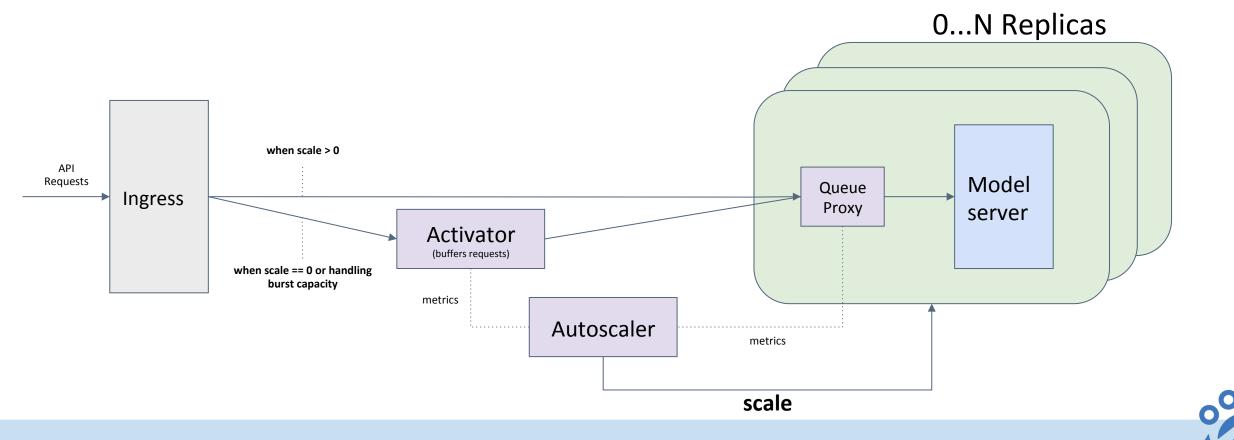
```
apiVersion: serving.kubeflow.org/v1alpha2
kind: InferenceService
metadata:
  name: bert-serving-onnx
spec:
                             Pre/Post Processing
  default
    transformer:
      custom:
        container:
          image: bert-transformer:v1
          env:
            name: STORAGE URI
            value: s3://examples/bert transformer
    predictor:
      onnx:
        storageUri: s3://examples/bert
        runtimeVersion: 0.5.1
                               ONNX Runtime Server
```



GPU Autoscaling - KNative solution



- Scale based on # in-flight requests against expected concurrency
- Simple solution for heterogeneous ML inference autoscaling





But the Data Scientist Sees...



apiVersion: "serving.kubeflow.org/v1alpha2"

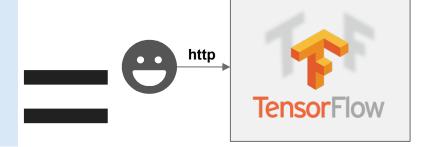
kind: "InferenceService"

metadata:

name: "flowers-sample"

spec: default: predictor: tensorflow:

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



- A pointer to a Serialized Model File
- 9 lines of YAML
- A live model at an HTTP endpoint

- Scale to Zero
- GPU Autoscaling
- Safe Rollouts
- Optimized Serving Containers
- Network Policy and Auth
- HTTP APIs (gRPC soon)
- Tracing
- Metrics





Production Experience (Early 2020)



Production users include:

Bloomberg





- CPU Throttling
 - Queue proxy can be throttled unfairly, linux kernel bug
 - Result: high tail latency
 - Solutions: (monitor carefully)
 - Don't set resource limit on queue proxy if no resource quotas on your cluster
 - Upgrade linux kernel on cluster nodes
- Cold start latency
 - Scale 0 means initial requests need to wait for model infrastructure to be created
 - ML models can be very large.
 - Solution: MinReplicas set to 1



Summary



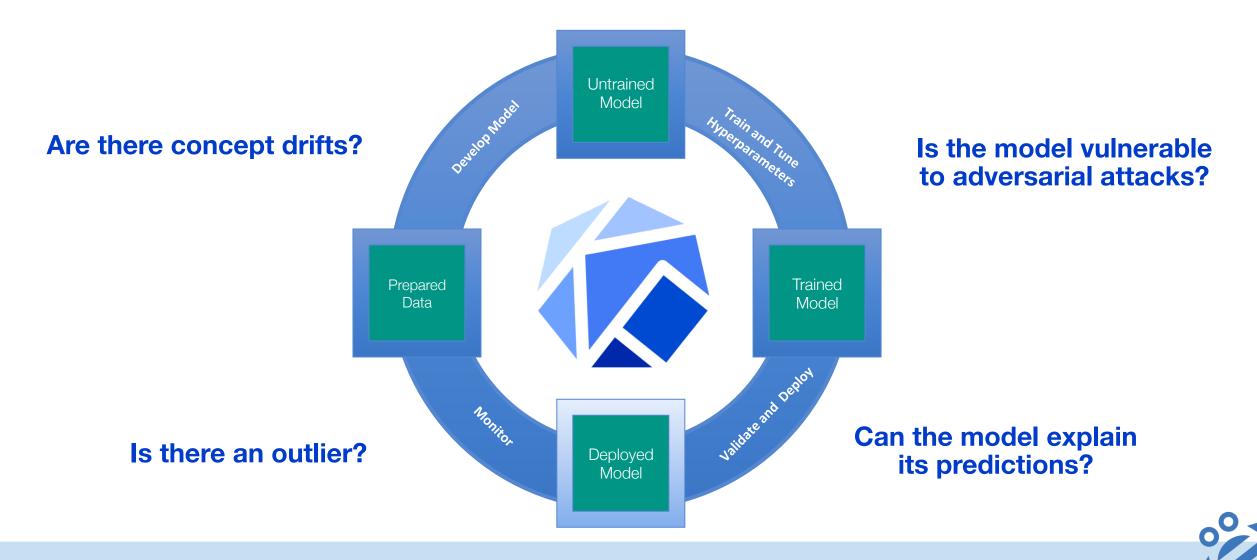
- With KFServing user can easily deploy the service for inference on GPU with performant industry leading model servers as well as benefiting from all the serverless features.
- Autoscale the inference workload based on your QPS, much better resource utilization.
- gRPC can provide better performance over REST which allows multiplexing and protobuf is a efficient and packed format than JSON.
- Transformer can work seamlessly with different model servers thanks to KFServing's data plane standardization.
- GPUs benefit a lot from batching the requests.





Model Serving is accomplished. Can the predictions be trusted?

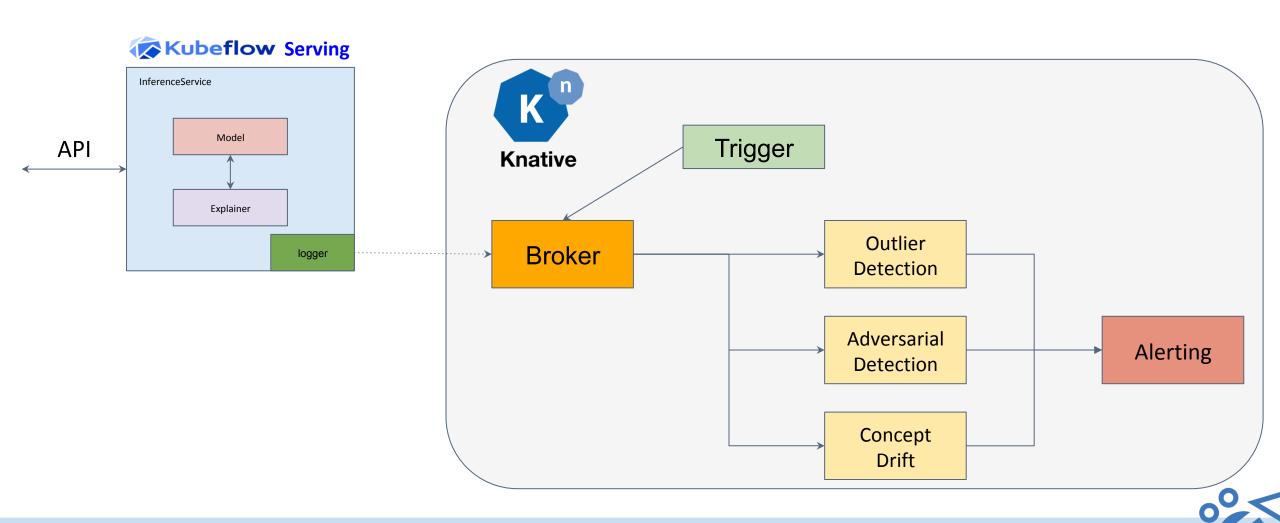






Production ML Architecture









Payload Logging





Payload Logging



Why:

- Capture payloads for analysis and future retraining of the model
- Perform offline processing of the requests and responses

KfServing Implementation (alpha):

- Add to any InferenceService Endpoint: Predictor, Explainer, Transformer
- Log Requests, Responses or Both from the Endpoint
- Simple specify a URL to send the payloads
- URL will receive CloudEvents



POST /event HTTP/1.0 Host: example.com

Content-Type: application/json

ce-specversion: 1.0 ce-type: repo.newItem

ce-source: http://bigco.com/repo

ce-id: 610b6dd4-c85d-417b-b58f-3771e532

<payload>





Payload Logging



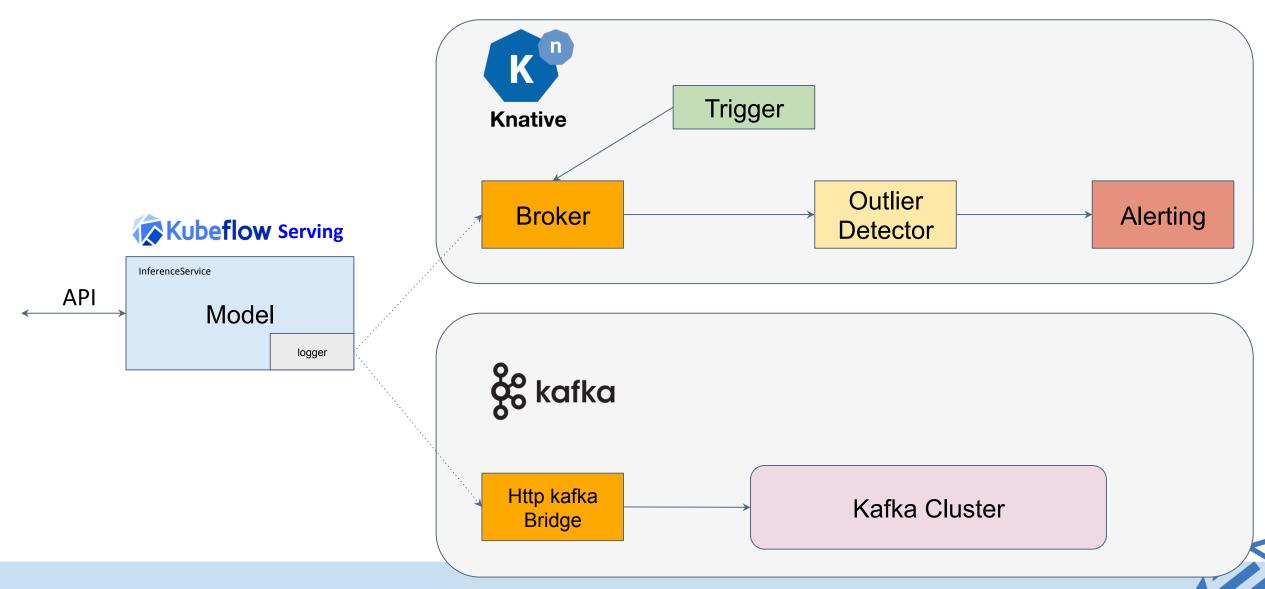
```
apiVersion: "serving.kubeflow.org/v1alpha2"
kind: "InferenceService"
metadata:
name: "sklearn-iris"
spec:
default:
 predictor:
  minReplicas: 1
  logger:
   url: http://message-dumper.default/
   mode: all
  sklearn:
   storageUri: "gs://kfserving-samples/models/sklearn/iris"
   resources:
    requests:
     cpu: 0.1
```





Payload Logging Architecture Examples









Machine Learning Explanations





KfServing Explanations



```
apiVersion: "serving.kubeflow.org/v1alpha2"
kind: "InferenceService"
metadata:
name: "income"
spec:
default:
predictor:
sklearn:
storageUri: "gs://seldon-models/sklearn/income/model"
explainer:
alibi:
type: AnchorTabular
storageUri: "gs://seldon-models/sklearn/income/
explainer"
```

```
apiVersion: "serving.kubeflow.org/v1alpha2"
kind: "InferenceService"
metadata:
name: "moviesentiment"
spec:
default:
predictor:
sklearn:
storageUri: "gs://seldon-models/sklearn/moviesentiment"
explainer:
alibi:
type: AnchorText
```





Seldon Alibi: Explain



https://github.com/SeldonIO/alibi



Giovanni Vacanti



Janis Klaise



Arnaud Van Looveren



Alexandru Coca

State of the art implementations:

- Anchors
- Counterfactuals
- Contrastive explanations
- Trust scores







Explanations: Resources



AI Explainability 360

→ (AIX360)

https://github.com/IBM/AIX360

AIX360 toolkit is an open-source library to help explain AI and machine learning models and their predictions. This includes three classes of algorithms: local post-hoc, global post-hoc, and directly interpretable explainers for models that use image, text, and structured/tabular data.

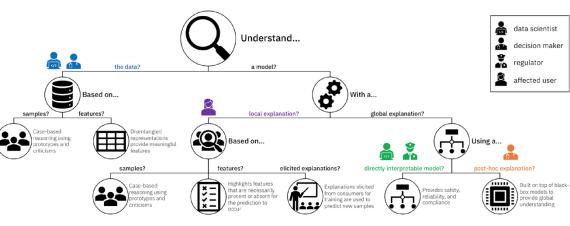
The AI Explainability360 Python package includes a comprehensive set of explainers, both at global and local level.

Toolbox

Local post-hoc Global post-hoc Directly interpretable

http://aix360.mybluemix.net

AIX360





AIX360 Explainability in KFServing



```
apiVersion: serving.kubeflow.org/v1alpha2
kind: InferenceService
metadata:
 labels:
   controller-tools.k8s.io: "1.0"
  name: aixserver
spec:
  default:
   predictor:
     minReplicas: 1
      custom:
       container:
         name: predictor
         image: aipipeline/rf-predictor:0.2.2
         command: ["python", "-m", "rfserver", "--model_name", "aixserver"]
          imagePullPolicy: Always
          resources:
           requests:
             memory: "2Gi"
             cpu: "1"
            limits:
             memory: "2Gi"
             cpu: "1"
    explainer:
     minReplicas: 1
      custom:
        container:
         name: explainer
          image: aipipeline/aix-explainer:0.2.2
         command: ["python", "-m", "aixserver", "--predictor_host", "aixserver-predictor-default.default.svc.cluster.local", "explainer_type", "LimeImages"]
          imagePullPolicy: Always
          resources:
           requests:
             memory: "4Gi"
             cpu: "2"
            limits:
             memory: "4Gi"
             cpu: "2"
```







ML Inference Analysis





ML Inference Analysis



Don't trust predictions on instances outside of training distribution!

- Outlier Detection
- Adversarial Detection
- Concept Drift





Outlier Detection

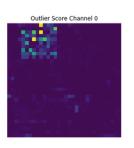


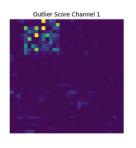
Don't trust predictions on instances outside of training distribution!

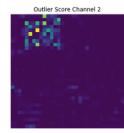
→ Outlier Detection

Detector types:









- stateful online vs. pretrained offline
- feature vs. instance level detectors

Data types:

- tabular, images & time series

Outlier types:

- global, contextual & collective outliers



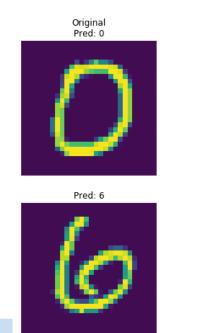


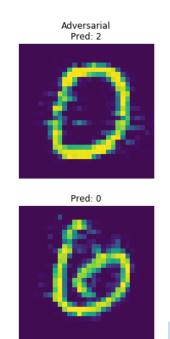
Adversarial Detection

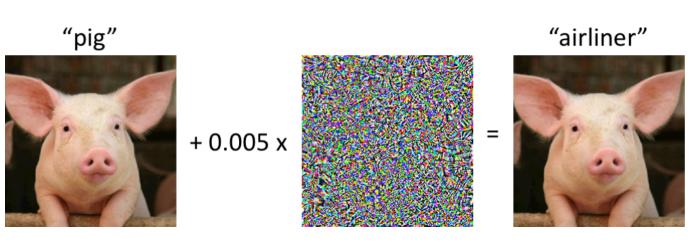


Don't trust predictions on instances outside of training distribution!

- → Adversarial Detection
 - Outliers w.r.t. the model prediction
 - Detect small input changes with a big impact on predictions!











Concept Drift



Production data distribution != training distribution?

→ Concept Drift! Retrain!

Need to track the right distributions:

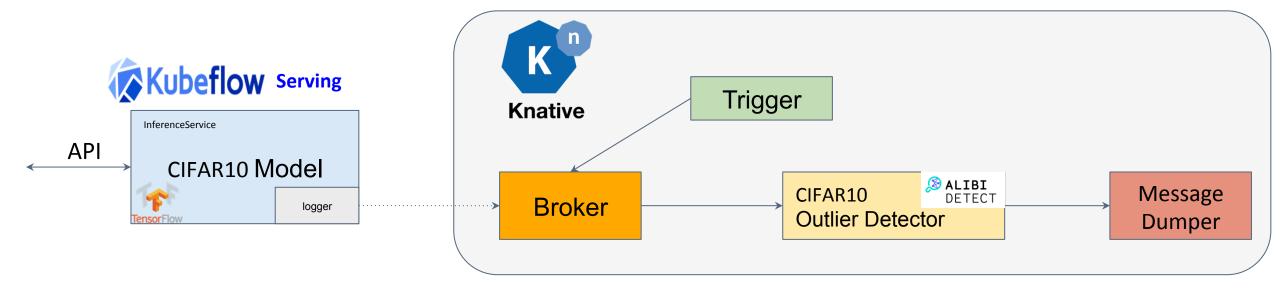
- feature vs. instance level
- continuous vs. discrete
- online vs. offline training data
- track streaming number of outliers





Outlier Detection on CIFAR10









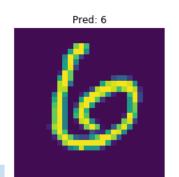
Adversarial Detection Demos

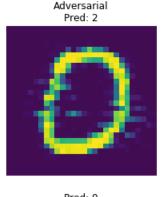


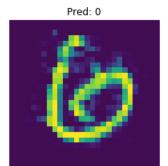
KFServing MNIST Model with Alibi:Detect VAE Adversarial Detector

https://github.com/SeldonIO/alibi-detect/tree/master/integrations/ samples/kfserving/ad-mnist

Original Pred: 0







KFServing Traffic Signs Model with Alibi:Detect VAE Adversarial Detector









Adversarial









Pred adversarial: 18





Adversarial Robustness 360 \(\text{ART} \)



https://github.com/IBM/adversarial-robustness-toolbox

ART is a library dedicated to adversarial machine learning. Its purpose is to allow rapid crafting and analysis of **attack**, **defense and detection methods** for machine learning models. Applicable domains include finance, self driving vehicles etc.

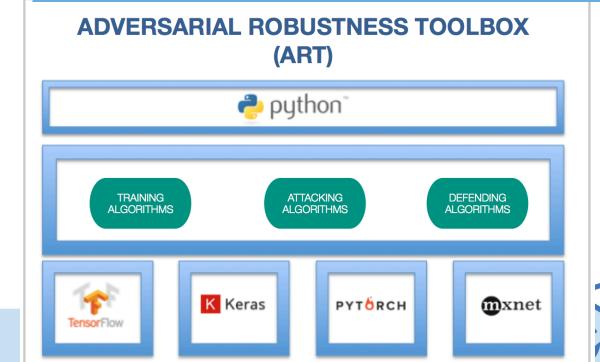
The Adversarial Robustness Toolbox provides an implementation for many state-of-the-art methods for attacking and defending classifiers.

Toolbox: Attacks, defenses, and metrics

Evasion attacks Defense methods Detection methods Robustness metrics

https://art-demo.mybluemix.net/

ART



Open Source Projects



ML Inference KFServing	https://github.com/kubeflow/kfserving
Seldon Core	https://github.com/SeldonIO/seldon-core
 Model Explanations Seldon Alibi 	https://github.com/seldonio/alibi
∘ IBM AI Explainability 360	https://github.com/IBM/AIX360
 Outlier and Adversarial Detection and Concept Drift Seldon Alibi-detect 	https://github.com/seldonio/alibi-detect
 Adversarial Attack, Detection and Defense IBM Adversarial Robustness 360 	https://github.com/IBM/adversarial-robustness-toolbox

