

Virtual

### MLOps at Snapchat: Continuous Machine Learning with Kubeflow & Spinnaker

Kevin Dela Rosa



# Kevin Dela Rosa

@perhaps\_ai

KubeCon CloudNativeCon North America 2020



# Agenda



- Snapchat
- Anatomy of a production ML system
- MLOps best practices
- Journey to production grade ML



# **Unlocking lenses**

### **Traditional Mechanisms**

- Snapcode
- Lens Link
- Lens Explorer



### **Discovery through Scan**

- Scan Triggers
- Marker Images





Virtual

CloudNativeCon

North America 2020

KubeCon

# Scan triggers





North America 2020 —



# **Marker lenses**





North America 2020 —



# KubeCon CloudNativeCon Virtua

### Machine Learning Models

- Image Classification
- Object Detection
- Semantic Segmentation
- Content Based Information Retrieval
- Nearest Neighbor Search
- Ranking

and more...

## **MLOps Setup: CI/CD Automation**



From - https://cloud.google.com/solutions/machine-learning/mlops-continuous-delivery-and-automation-pipelines-in-machine-learning

**CloudNativeCon** 

North America 2020

KubeCon

Virtual

# Everybody has to start somewhere with the cloud Native Con Virtual

	local	nost	Ċ	
Home			Untitled	+
Jupyter Untitled (unsaved changes)	)			Logout
File Edit View Insert Cell	Kernel Help			Python 3 O
₽ + % 4 6 1 + 4	Code	•	CellToolbar	
In [ ]:				

# This is the way



### MLOps = DevOps + ML

### **Special considerations**

- ML is Model + Code + Data
- Data changes

### **Unique to ML systems**

- Data pipelines / ETL
- Feature store
- ML Pipeline

## **MLOps Principles**<sup>1</sup>

- Versioning
- Testing
- Automation
- Reproducibility
- Monitoring

Incremental process, level of automation will increase as your system matures

### **Steps**

- 1. Notebook decomposition, source control
- 2. ML pipeline automation
- 3. Continuous Integration
- 4. Continuous delivery of ML pipeline
- 5. Continuous training, triggering of ML pipeline
- 6. Continuous delivery of model / serving
- 7. Monitoring

# **1: Notebooks to Containers**



Break up monolithic Jupyter notebooks into modular code / containerized programs. Add to source control.



### **Characteristics**

<u>Versioning</u> – track changes to algorithms & transform code <u>Testing</u> – modular code = easier unit testing <u>Reproducibility</u> – environment and package dependencies

# 2: Kubeflow for ML pipeline



Connect Docker images via Kubeflow pipelines to automate the ML pipeline.



# **Kubeflow**

#### Kubeflow

**Pipelines** 

Terminate Retry Clone run Archive



Experiments > My XGBoost experiment



**Pipeline**: Graph describing ML workflow, it's components and how they relate with each other

**Component**: A "step" in pipeline that launches one or more Kubernetes pod, like a function it has: name, parameters, return values, and "body" (docker image)

Serialized data (strings, files) is passed between components; e.g. return value of parent component is used as parameter to a child component.

**Run**: Single execution of a pipeline

Build commit: ee207f2

(i) Runtime execution graph. Only steps that are currently running or have already completed are shown.

### **2: Kubeflow for ML pipeline**



North America 2020



data-validation features train-model model-evaluation model-validation

### **Characteristics**

<u>Automation</u> – automate the steps in your ML process with a pipeline(s) <u>Reproducibility</u> – store metadata and model information <u>Versioning</u> – track changes to your ML process, pipeline parameters, hyper parameters & training configuration <u>Testing</u> – data validation, model evaluation/validation, model spec unit tested, integration test ML pipeline <u>Monitoring</u> – log model performance on holdout validation set

# **3: Continuous integration**



Automate the building of code, unit tests, and integration tests.

**Publish artifacts**: docker images, compiled Kubeflow pipeline YAML, Kubernetes configuration



# **3: Continuous integration**



### **Characteristics**

<u>Versioning</u> – publish docker images, Kubeflow component and pipeline specification YAML

<u>Automation</u> – automated unit and integration testing of docker images and ML pipelines



Orchestrate the deployment of ML pipeline(s) with a continuous delivery system.



### **Characteristics**

<u>Automation</u> – trigger from CI to automatically react to changes in training process / ML pipeline <u>Reproducibility</u> – deployment of ML pipeline in repeatable fashion, rollback to previous versions

<ul> <li>Promote to Pr</li> </ul>	od		Permalink 🖇	Create	🌣 Configure 👻	Pipeline Actions 🗸
Configuration Fir Sta	nd Image from Deploy Canary A	pproval (Red/Black)	Tear Down (	Canary Wait 2	hrs Des	troy Old Prod
	C Add stage			ľ	Copy an existing stage	
Deploy Prod (Red/Black) Stage Na		Deploy Prod (Red/Black)				Remove stage
Deploys the previously baked or found image Depends On C	Cutover Manual Approval		<b>a</b>			
Pipeline: Deployment mana	agement construct consisting	g of a sequence of				
actions known as <b>Stages</b> .						
Stage: An action for pipelin	e to perform, like Deploy, Ru	in Job, Manual	Region	Strategy Ca	pacity	Actions

### 4: Spinnaker for CD of ML Pipeline



Virtual CloudNativeCon

North America 2020

KubeCon

# 5: Continuous training



Trigger ML Pipeline Deployment on a set schedule (e.g. daily, weekly, monthly, etc.) or upon availability of new data.

### **Characteristics**

<u>Automation</u> – training/new model generation process reacts to new data



Trigger ML model server deployment upon availability of new model or changes to serving configuration.



### **Characteristics**

<u>Automation</u> – continuous delivery/deployment of new models and/or serving infrastructure <u>Reproducibility</u> – deployment of model in repeatable fashion, rollback to previous versions

# 6: Spinnaker for model delivery



 $\sim$ 

CloudNativeCon

North America 2020

KubeCon

Virtual

# 7: Monitoring



Log telemetry about server and model performance. In addition to typical server metrics (latency, traffic, errors, saturation), monitor and alert on unexpected data/prediction changes.

### **Characteristics**

<u>Monitoring</u> – unexpected prediction changes can signal change in user behavior or bad model behavior

# "Continuous Machine Learning"











### Snapchat

- <u>https://lensstudio.snapchat.com/guides/sharing/scan/</u>
- <u>https://lensstudio.snapchat.com/templates/marker/</u>

### **MLOps Resources**

- <u>https://ml-ops.org/</u>
- <u>https://cloud.google.com/solutions/machine-learning/mlops-continuo</u> <u>us-delivery-and-automation-pipelines-in-machine-learning</u>
- <u>https://github.com/cdfoundation/sig-mlops</u>

### Tools

- <u>https://www.kubeflow.org/</u>
- <u>https://spinnaker.io/</u>



Kevin Dela Rosa



# Thanks for listening!

