

삼성 오픈소스 컨퍼런스

Use Case로 살펴보는 Production Deep Learning 서비스 인프라 Eco System

Production AI Pain Point 및 그 해결 방법 사례 공유

SK텔레콤 | AI센터 | 김훈동, 박찬엽



We are ...

- Hoon Dong Kim.

- SK Telecom , AI Center , AI Product DevOps Team Leader.
- Microsoft BigData MVP(Most Valuable Professional) : 2016년 ~ 2017년
- Microsoft AI MVP(Most Valuable Professional) : 2018년 ~ 2019년
- Korea Spark User Group ([스파크 사용자 모임](#)) 운영진

<http://hoondongkim.blogspot.kr> [블로그]

<https://www.slideshare.net/ssusere94328/> [슬라이드 쉐어]

<https://www.facebook.com/kim.hoondong> [SNS]

- Chan Yup Park

- SK Telecom , AI Center , AI Product DevOps Team, AI DevOps Engineer
- 툴링에 관심이 많아 docker, kubernetes의 운영/관리/배포등의 과정을 개선하는 일을 맡고 있습니다.
- R로 개발 경험을 시작해서 관련 생태계에 기여하기 위해 노력하고 있습니다.

<https://mrchypark.github.io/>[블로그]

<https://github.com/mrchypark> [github]

Agenda

1. Production AI Pain Point & Solution Use Case
2. Production AI Open Source Eco System
3. Production AI DevOps
4. Our Open Source Contribution

1. Production AI Pain Point & Solution Use Case

Our Pain Point 1 – ML Serving Dilemma

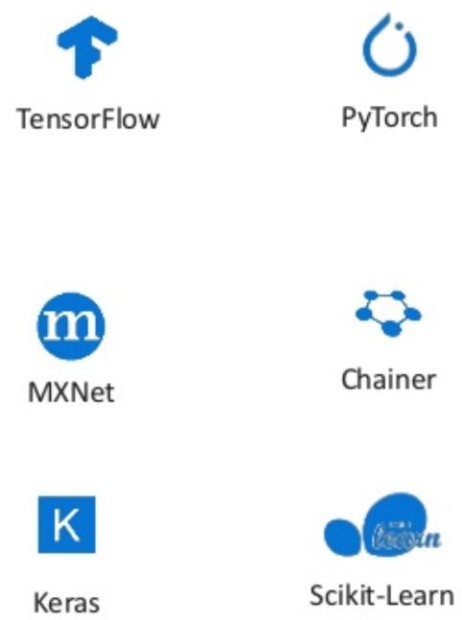
- [방법1] Tensorflow Serving
 - bazel 빌드 , C++ code, gRPC.
 - python serving performance 좋지 않음.
 - 요즘은 pytorch, mxnet 하시는 분들도 급격히 늘고 있음.
 - Scikit learn 전처리 모델은?
- [방법2] Cloud PaaS
 - Azure Machine learning Service , AWS SageMaker , GCP CloudML
 - 3사가 다 10%가 부족한 부분이 있음.
 - 매우 비쌘.
- [방법3] Flask 등을 이용한 범용 서빙 아키텍처 구성
 - 빠르고 쉽게 prototype 할 수 있으나, production을 하기 위해선 험난한 Engineering Art 가 필요함.
 - Python 은 너무 너무 * 10 느린 언어임.
- [방법4] Full 사양 GPU 서버 혹은 VM 으로 Serving?
 - 월 4만원 짜리 CPU Docker vs 월 1000만원 짜리 최신 GPU VM

Our Pain Point 1 – ML Serving Dilemma

- ML Serving for Lots of Frameworks

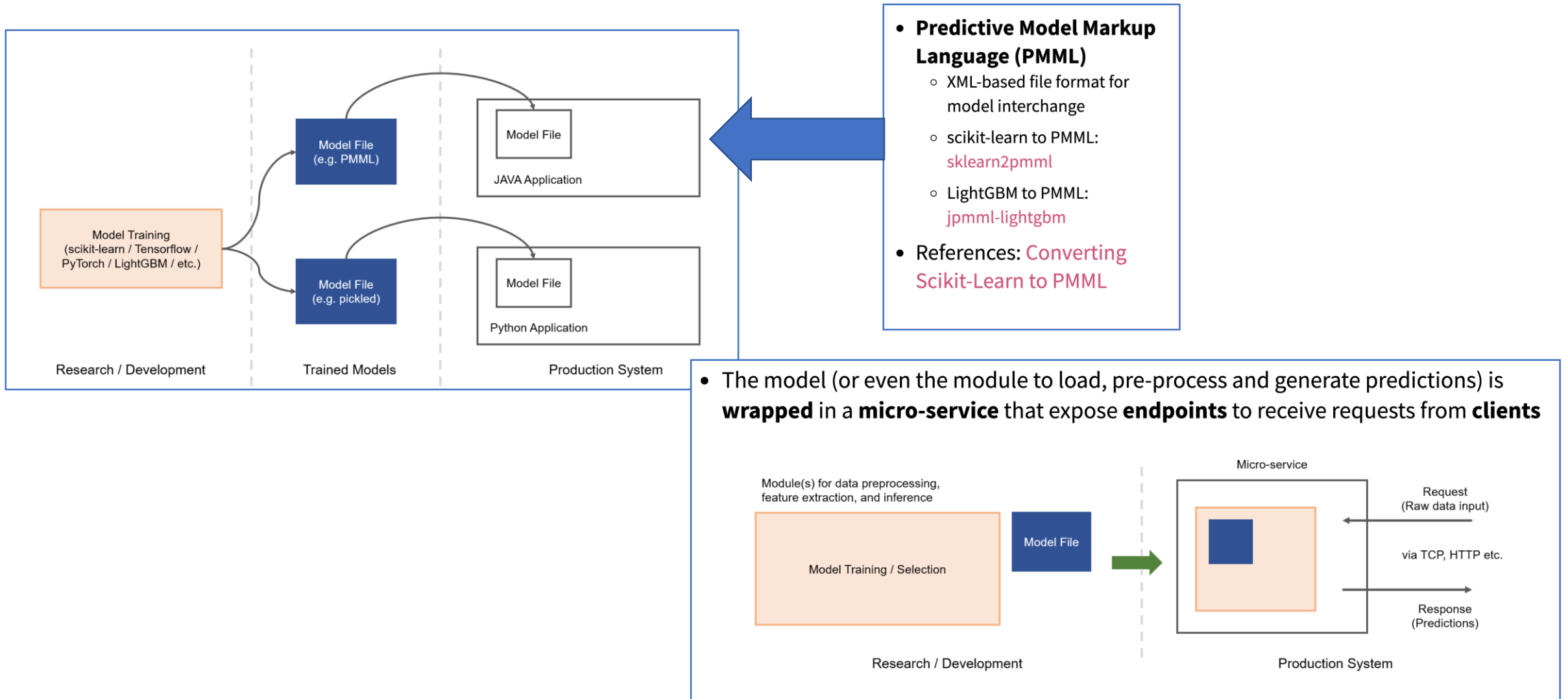
Tensorflow Serving 만 가지고는...

Lots of ML Frameworks



The image displays a collection of logos for various machine learning frameworks arranged in a 3x2 grid. The logos are: TensorFlow (top-left), PyTorch (top-right), MXNet (middle-left), Chainer (middle-right), Keras (bottom-left), and Scikit-Learn (bottom-right). The text 'Lots of ML Frameworks' is written in green on the left side of the grid.

Our Pain Point 1 – 해결 팁



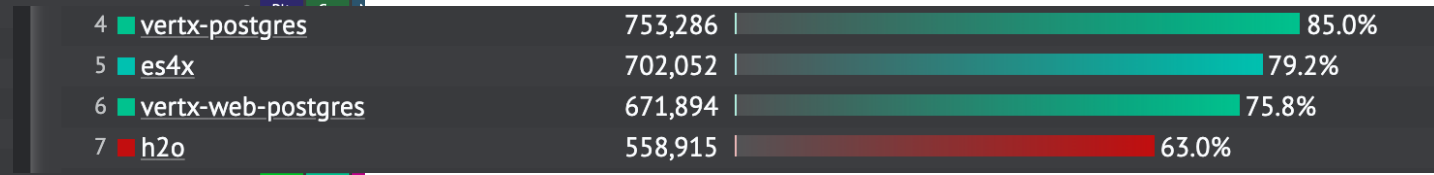
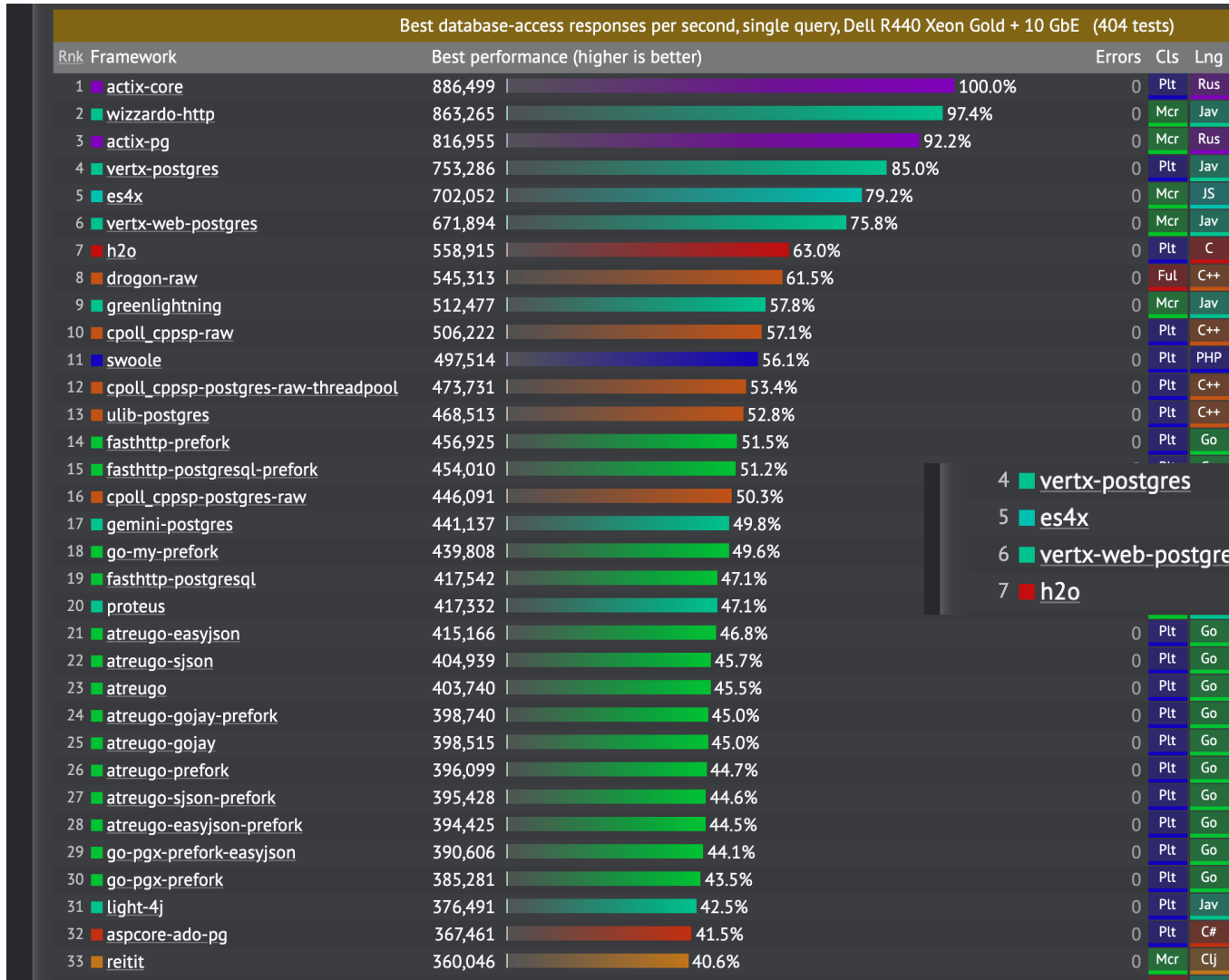
Our Pain Point 1 – ML Serving Dilemma 해결팁

- [방법1] Tensorflow Serving
 - bazel 빌드 , C++ code, gRPC.
 - python serving performance 좋지 않음.
 - 요즘은 pytorch, mxnet 하시는 분들도 급격히 늘고 있음.
 - Scikit learn 전처리 모델은?
- [방법2] Cloud PaaS
 - Azure Machine learning Service , AWS SageMaker , GCP CloudML
 - 3사가 다 10%가 부족한 부분이 있음.
 - 매우 비쌘.
- **[방법3] Flask 등을 이용한 범용 서빙 아키텍처 구성**
 - 빠르고 쉽게 prototype 할 수 있으나, production을 하기 위해선 **힘난한 Engineering Art**가 필요함.
 - Python 은 너무 너무 * 10 느린 언어 임.
- [방법4] Full 사양 GPU 서버 혹은 VM 으로 Serving?
 - 월 4만원 짜리 CPU Docker vs 월 1000만원 짜리 최신 GPU VM



본 슬라이드
전체에서 설명...

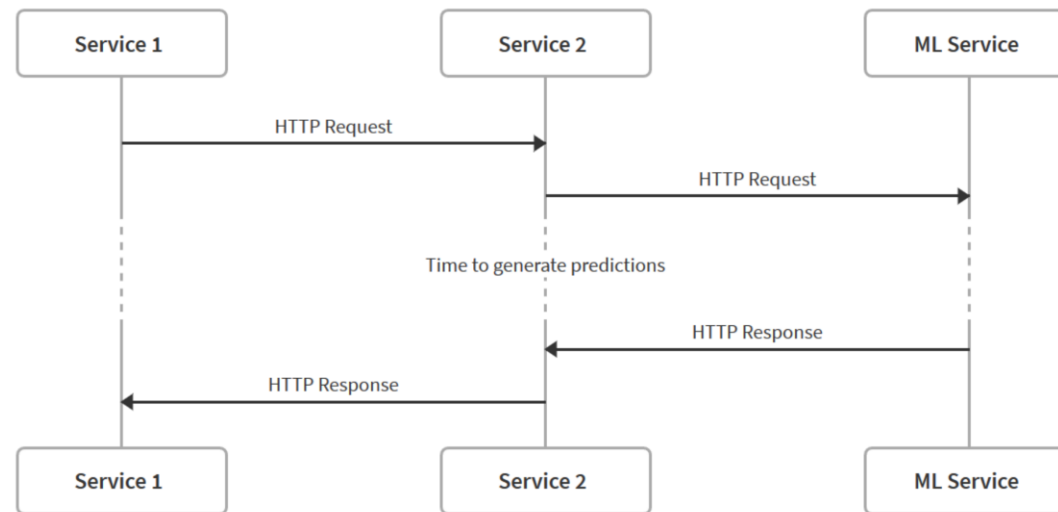
Our Pain Point 2- Poor Python Performance



Our Pain Point 2 – Poor Python Performance

- Disadvantage of Using HTTP APIs

- HTTP REST APIs are simple to implement, however the process is **synchronous**
- The client must **wait** until the ML service has finished the process of generating predictions
- In a complex system involving a lot of components, this may not be efficient



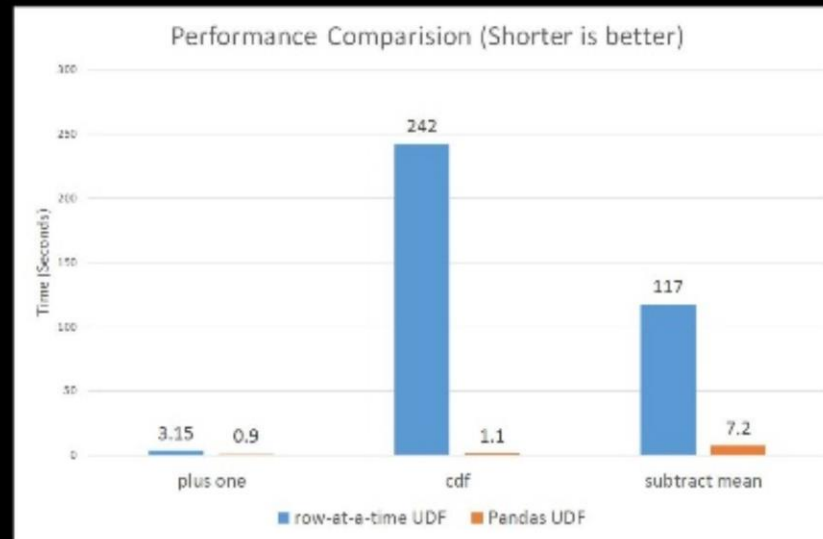
Our Pain Point 2 – 해결 팁

- Poor Python Performance
 - Pandas Data 함수 전달
 - > Data Copy 과정에서의 병목.
 - > Thread, Multi Process 에서 Thread Safe 하지 않음.
 - 함수형 언어 처럼 가상함수를 사용하고 싶지만...
 - > Python 에서 Map, Lambda를 : `functiontools.partial`
 - Python Thread 는 GIL 때문에 엄청나게 느림.
 - > 차라리 muti-process 가 낫다 : Go 의 Goroutine 쓰듯이 쓰고자 한다면... `cotyledon`
- MomoryView 활용
 - Data 핸들링 시, pointer 접근 하듯...
 - 대용량 Data Memory 복사 방지
- Microservice 로 잘게 쪼갬다.
- 모든 것은 비동기로...
- PySpark, ScalaSpark...
- PMML + Java 컨버전...

Our Pain Point 2 – 해결 팁

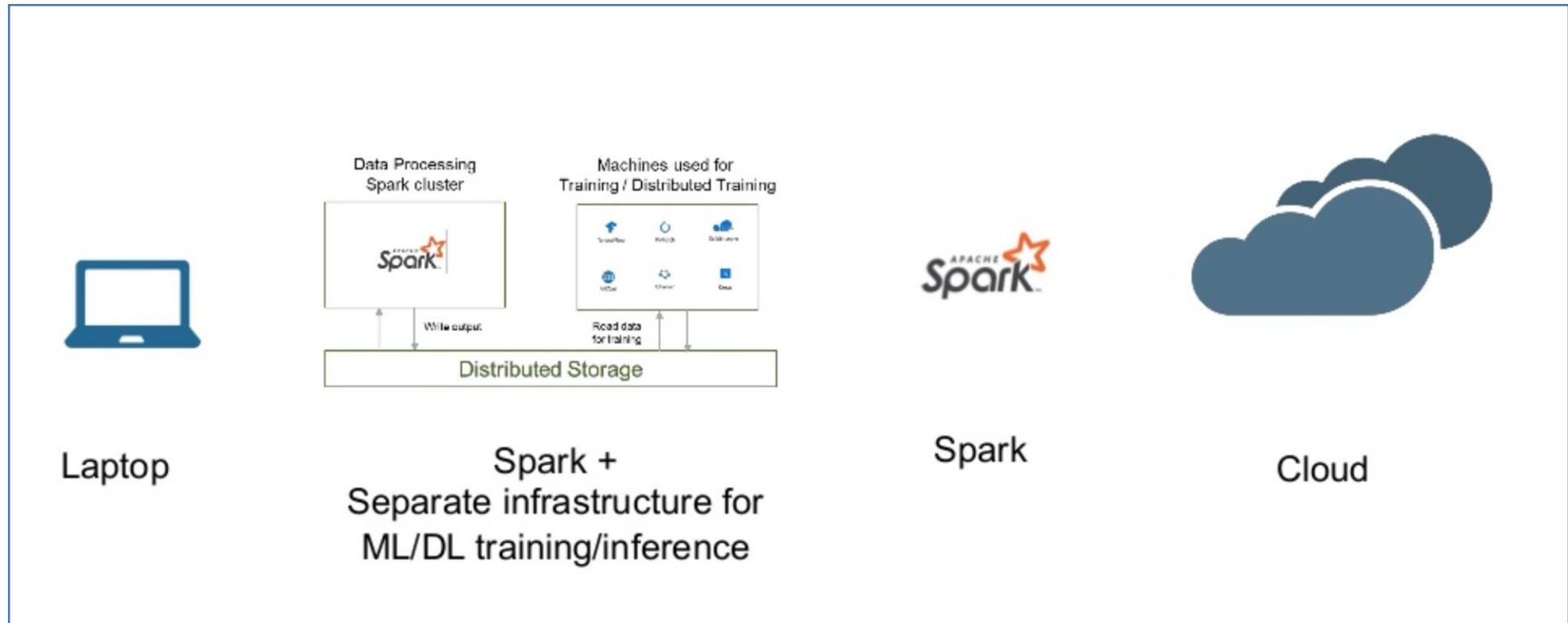
- Pandas : Poor Python Performance -> Pandas UDF

Pandas UDF was introduced in Spark 2.3, which uses Arrow for data exchange and utilizes Pandas for vectorized computation.



Our Pain Point 3

- Multiple System, Multiple Dev-env. , Multiple Production-env.



Our Pain Point 3 – 해결 팁

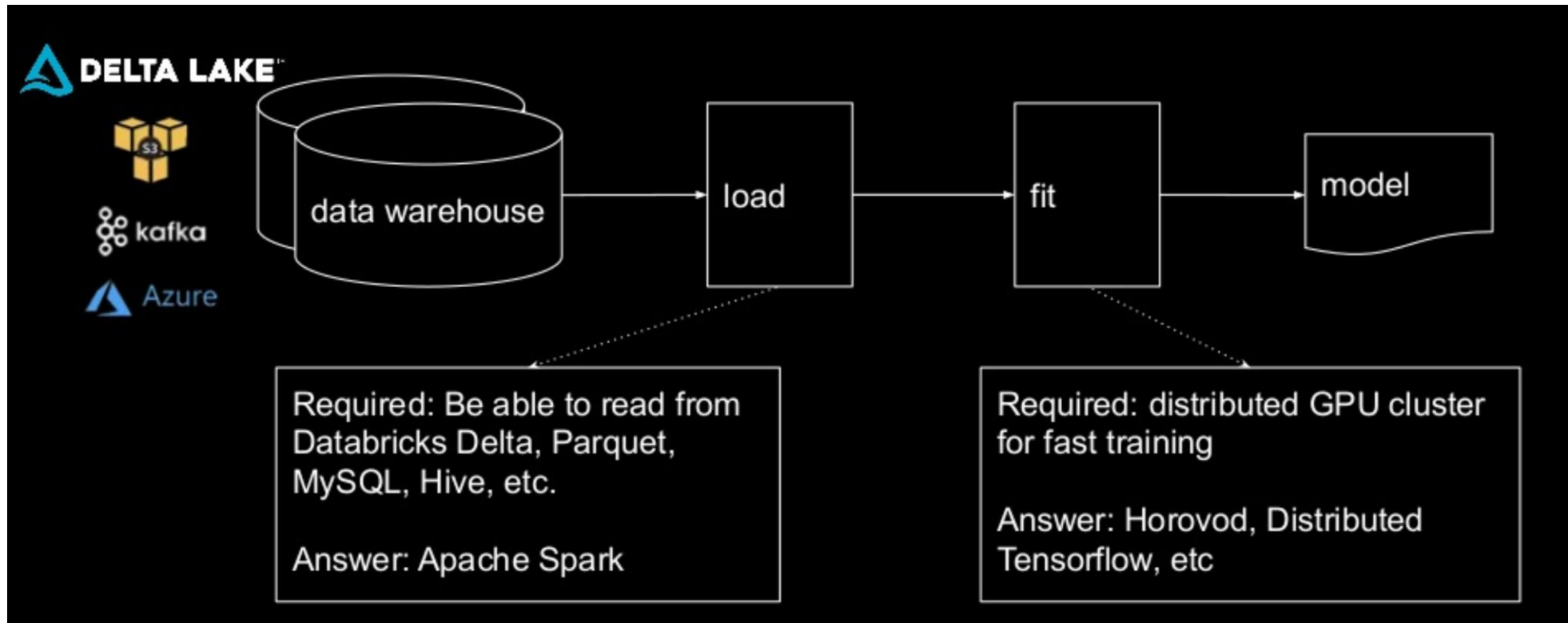
- Auto Scale Out, Auto Scale Down & Env. Abstraction.

Kubernetes is All You Needs.

But, 현실이 녹록치 않은 않습니다.
뒤에서 AI DevOps 부분으로 따로 언급~

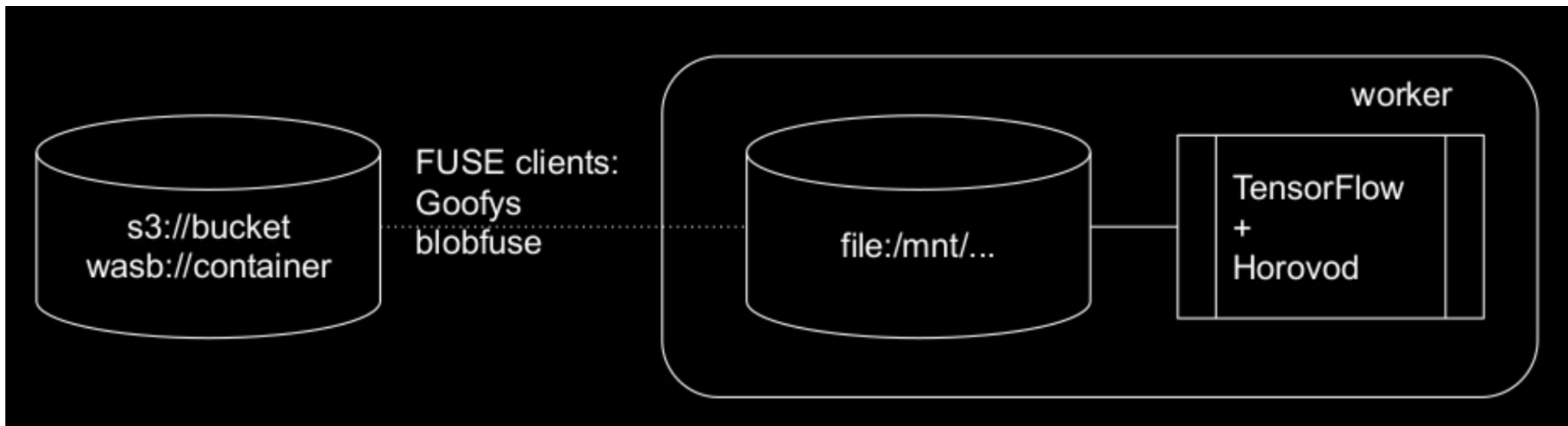
Our Pain Point 4

- BigData Scale Data Pre-processing & Distributed Training & Training Performance



Our Pint Point 4 – 해결 팁

- BigData Scale Data Pre-processing & Distributed Training & Training Performance
- **Recommendation** : High-performance FUSE clients to mount remote storage as local files so DL frameworks can load/save data directly.



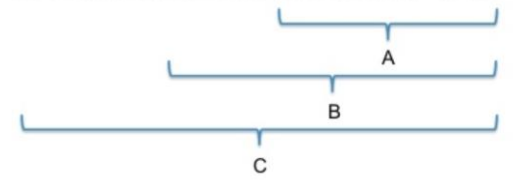
Our Pint Point 4 – 해결 팁

XGBoost in BigData Scale Production

- **XGBoost4J: Distributed XGBoost for Scala/Java**
- XGBoost4J-Spark : XGBoost on Apache Spark
 - **USE Case : Airbnb, Alibaba, eBay, Microsoft, Snapshots, Tencent, Uber, etc...**
- <https://github.com/dmlc/xgboost/tree/master/jvm-packages>

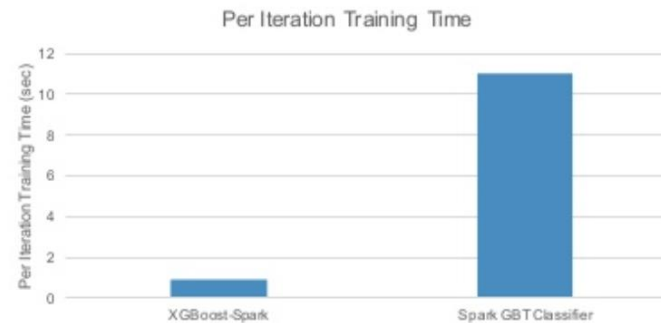
dmlc
XGBoost

XGBoost: Gradient Boost Decision Tree



Performance Evaluation

Airline Dataset (22M examples), 48 Workers in XGBoost/Tasks in Spark
Hardware: 6 D4V2 VMs on Azure serving Spark Executors



XGBoost + Spark + GPU

- XGBoost + Distributed mode via Spark
- XGBoost + GPU support via CUDA
- XGBoost + Multi-GPU via NCCL2
- XGBoost + Spark + GPU : Multi-GPU + Multi Node

기타,
이런 류의
좀더 많은
Open Source
는 뒤에서
소개...

Spark Cluster

- Standalone cluster on GCP
- 5 virtual machines, each has:
 - 64 vCPUs (32 physical cores)
 - 416 GB memory
 - 4 x NVIDIA Tesla T4
 - 400 GB SSD persistent disk
 - Default networking
- 4 Spark workers per VM

Sample Code

```
import ml.dmlc.xgboost4j.scala.spark.XGBoostClassifier
val xgbParam = Map("eta" -> 0.1f,
  "max_depth" -> 20,
  "max_leaves" -> 256,
  "grow_policy" -> "lossguide",
  "num_round" -> 100,
  "num_workers" -> 20,
  "nthread" -> 16,
  "tree_method" -> "gpu_hist")
val xgbClassifier = new XGBoostClassifier(xgbParam).
  setFeaturesCol("features").
  setLabelCol("labels")
```

Preliminary Results

		Accuracy (AUC)	Training Loop (Seconds)
Max Tree Depth = 8	CPU	0.832	1071.002
	GPU	0.832	139.641
	Speedup		766.97%
Max Tree Depth = 20	CPU	0.833	1088.662
	GPU	0.833	165.868
	Speedup		656.34%

Our Pain Point 5 - Serving Cost & 해결 팁

- Low Cost & Agile 접근 방법에 대하여

딥러닝 모델러는 Graph DB, Microservice 를 모르고...,
NoSQL 개발자는 딥러닝을 모르고...

Deep Learning Framework 자체는 여기 까지만 구동 가능

이 단계로 포팅하려면, Deep Learning Output 그래프 노드 값을 Graph DB 화 하거나, Hash DB 등으로 inverted-Index 화 해야 하며, 고 난이도의 Engineering 작업이 수반 됨.

이 단계로 되면, Allibaba Scale의 Front AI 서비스가 매우 빠르고 매우 저렴하게 적용 가능.

GPU VM > CPU VM > GPU Docker > CPU Docker > K8S PaaS > Docker PaaS > Docker BaaS > ML BaaS > Serverless Microservice

GPU on-premise
1 GPU VM,
1 달에 얼마?

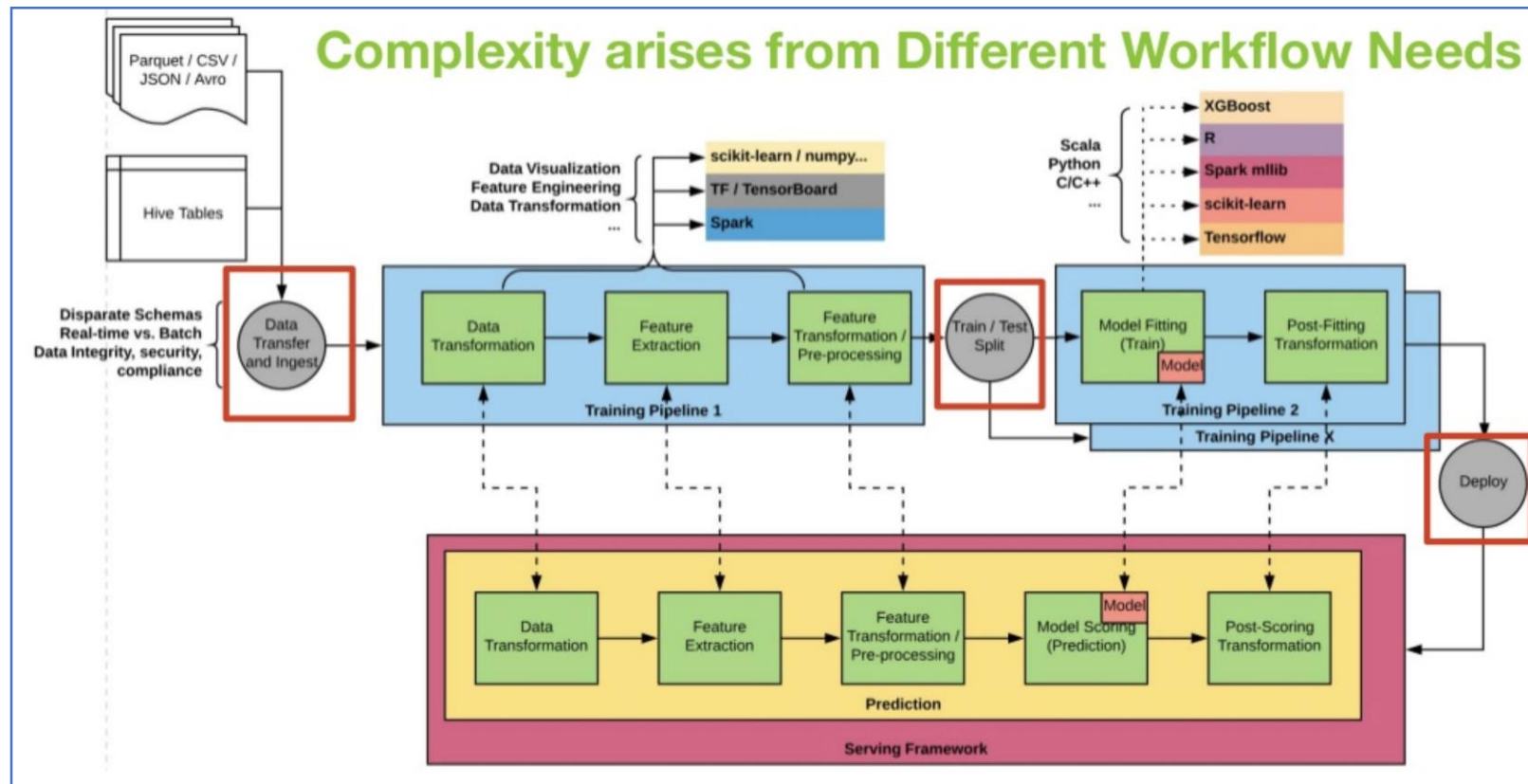
1 Docker,
1 시간에 얼마?
1 달에 얼마?

1,000,000 Transaction 에
얼마???

**Production AI 는 모델링 보다 Engineering 이 더 어렵다.
Production AI 는 Engineering Art 에 가깝다.
Cloud 를 잘 활용한다면, 많은 도움이 된다!**

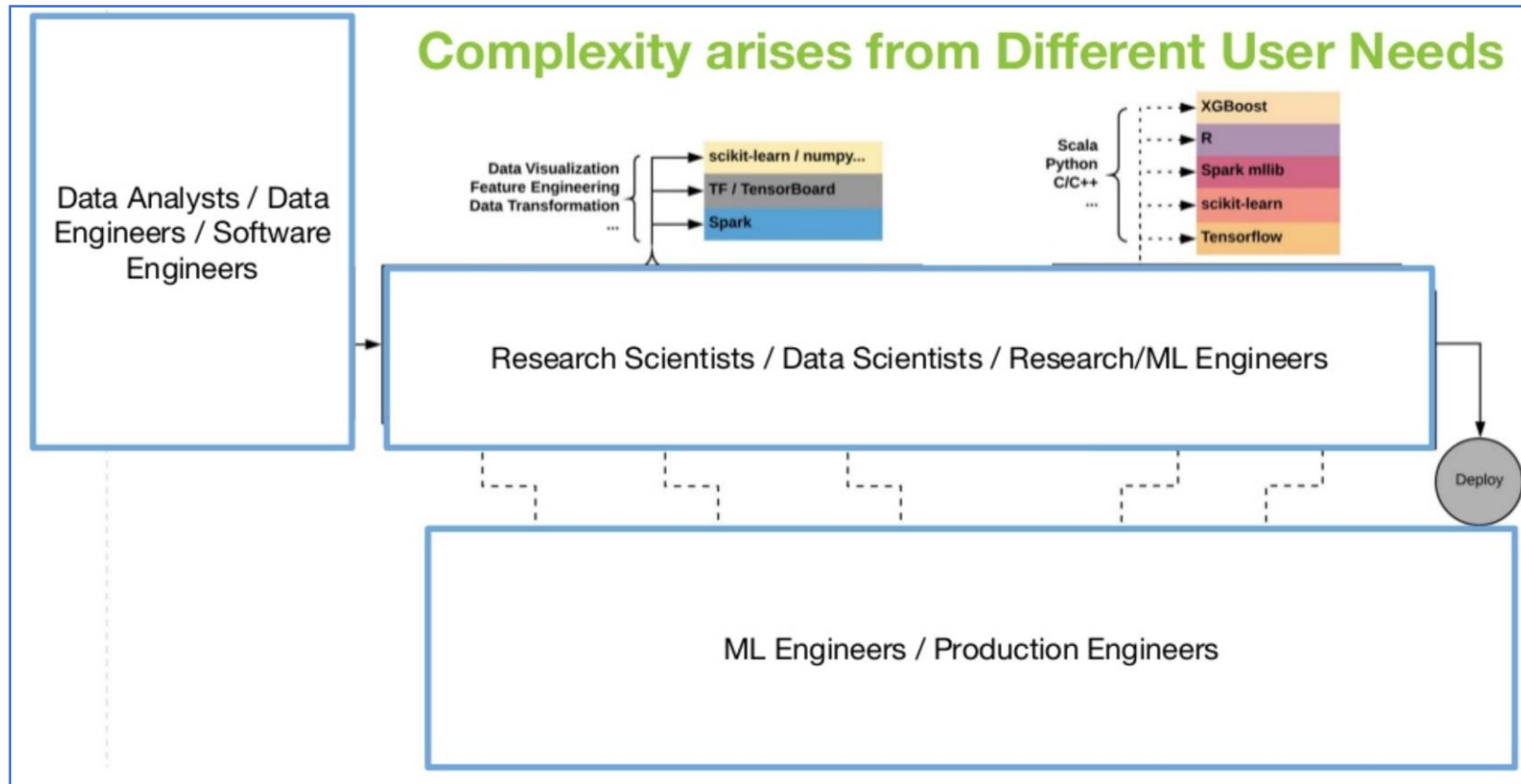
Our Pain Point 6

- None-Stop CI/CD deploy & DevOps & MLOps



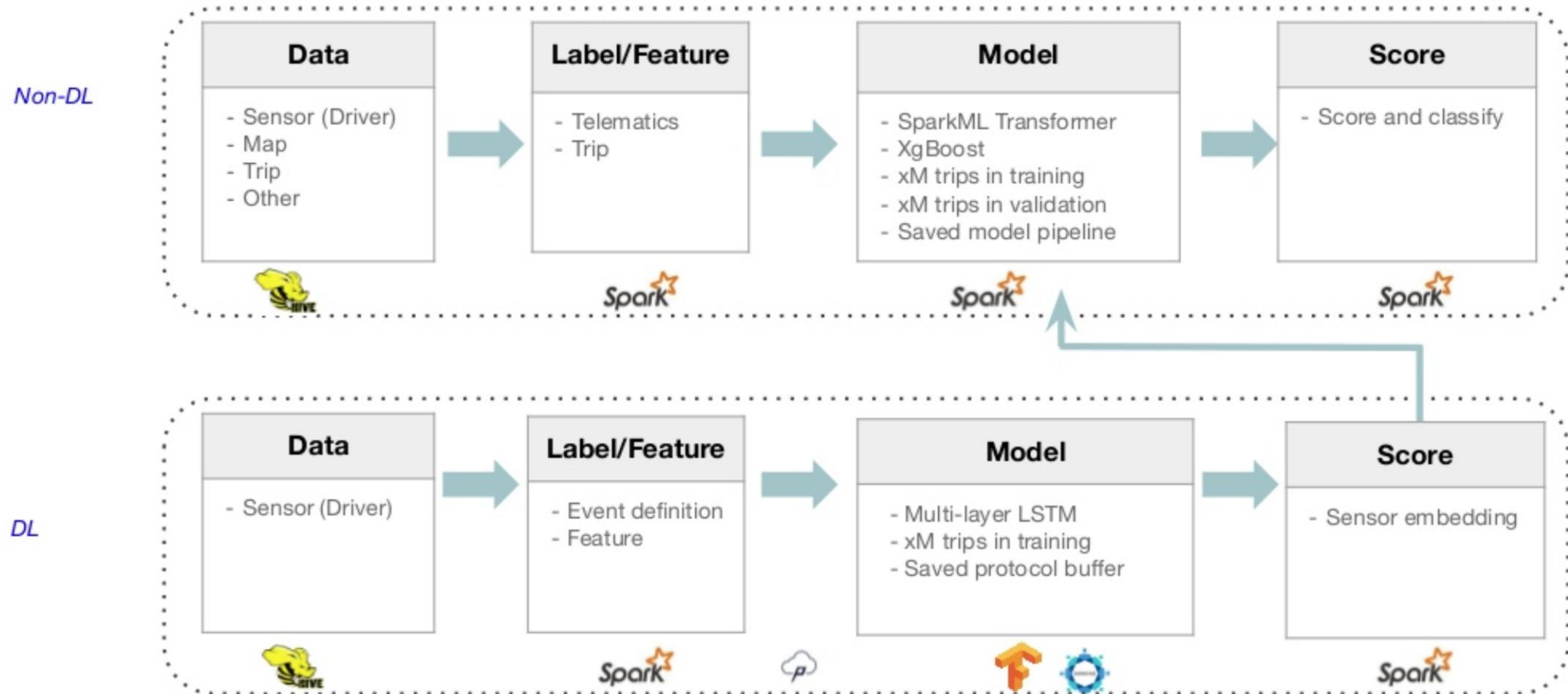
Our Pain Point 6 – 해결 팁 (AI DevOps)

- None-Stop CI/CD deploy & DevOps & MLOps



2. Production AI Open Source Eco System

BigData Scale {Deep Learning + None Deep Learning} Production Model Pipe Line



Horovod (Distributed Deep Learning at Scale)



- Open source library developed at Uber
- **Distributed** training for TensorFlow, Keras & PyTorch
- Uses bandwidth-optimal communication protocols & makes use of advanced networking
- Seamlessly installs via `pip install horovod`

Horovod (Distributed Deep Learning at Scale)

horovod.spark

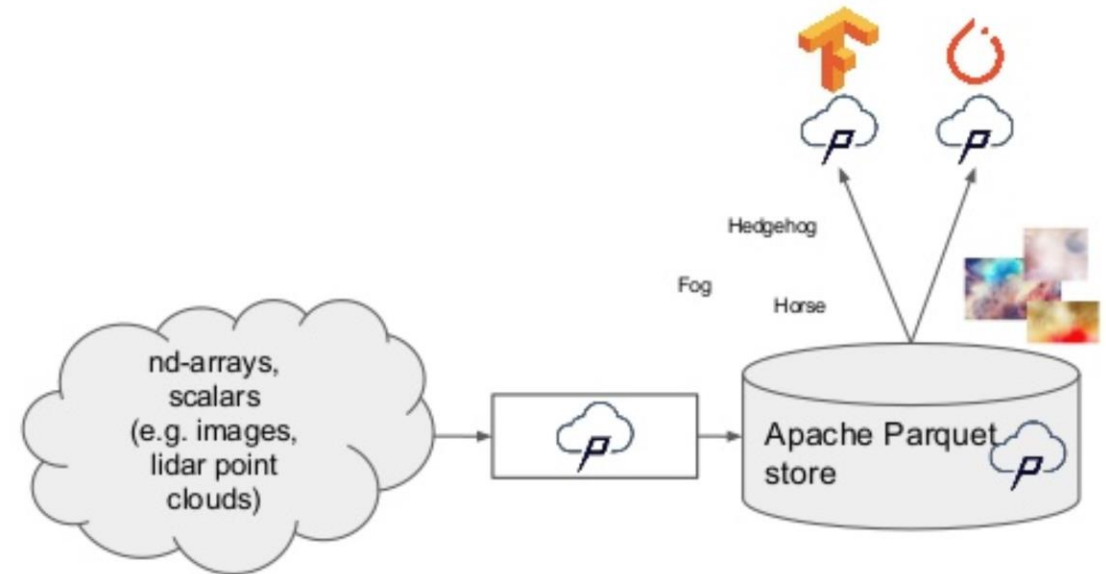
horovod.spark is a new feature in Horovod 0.16 release. Similar to HorovodRunner, it runs Horovod as a Spark job and takes python train functions. Its assumption is more general:

- no dependency on SSH,
- system-independent process termination,
- multiple Spark versions,
- and more ... also check out horovodrun:)

Petastorm (BigData Scale Data Deep Learning)

- Open source library developed at Uber ATG
- Enables deep learning directly from Parquet
- Supports Tensorflow, PyTorch, and PySpark

Apache Parquet as a dataframe with tensors



Horizon (Deep Reinforcement Learning at Scale)

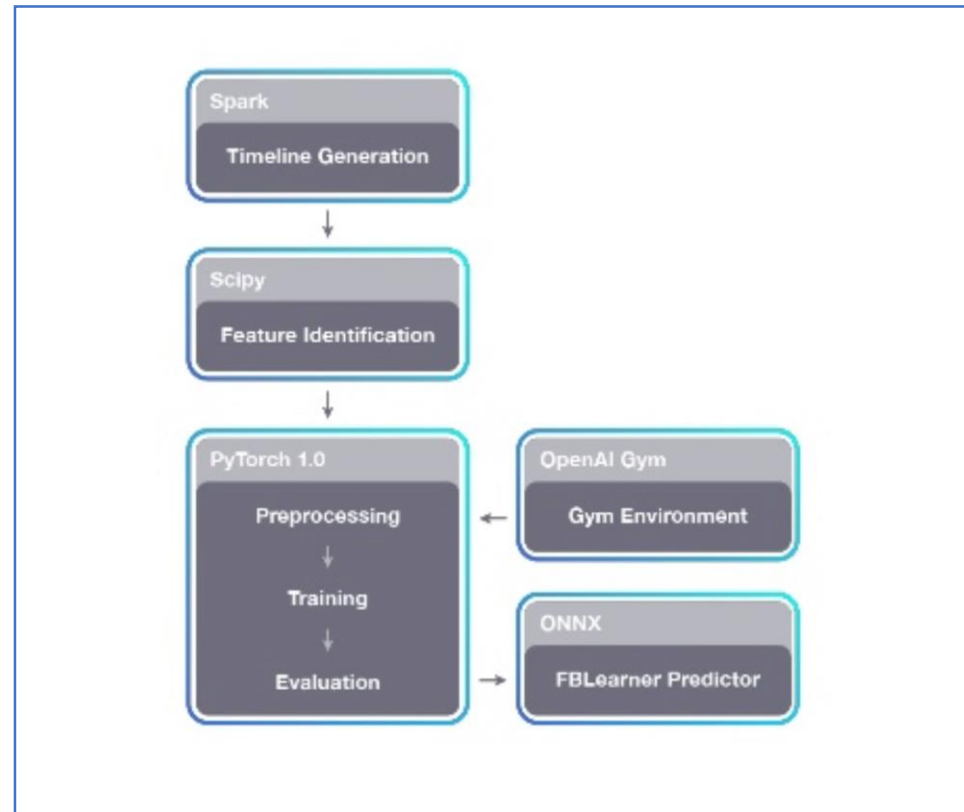
- Open Source End-To-End Large-scale RL framework By Facebook
 - distributed popular deep RL algorithms training.
 - workflow management.
 - data preprocessing.
 - feature transformation.
 - counterfactual policy evaluation.
 - optimized serving.
- Reinforcement Learning & Contextual Bandits.
- PyTorch for Modeling and Training.
- Caffe2 for Model Serving.
- <https://reagent.ai/>
- <https://github.com/facebookresearch/ReAgent>

Horizon (Deep Reinforcement Learning at Scale)



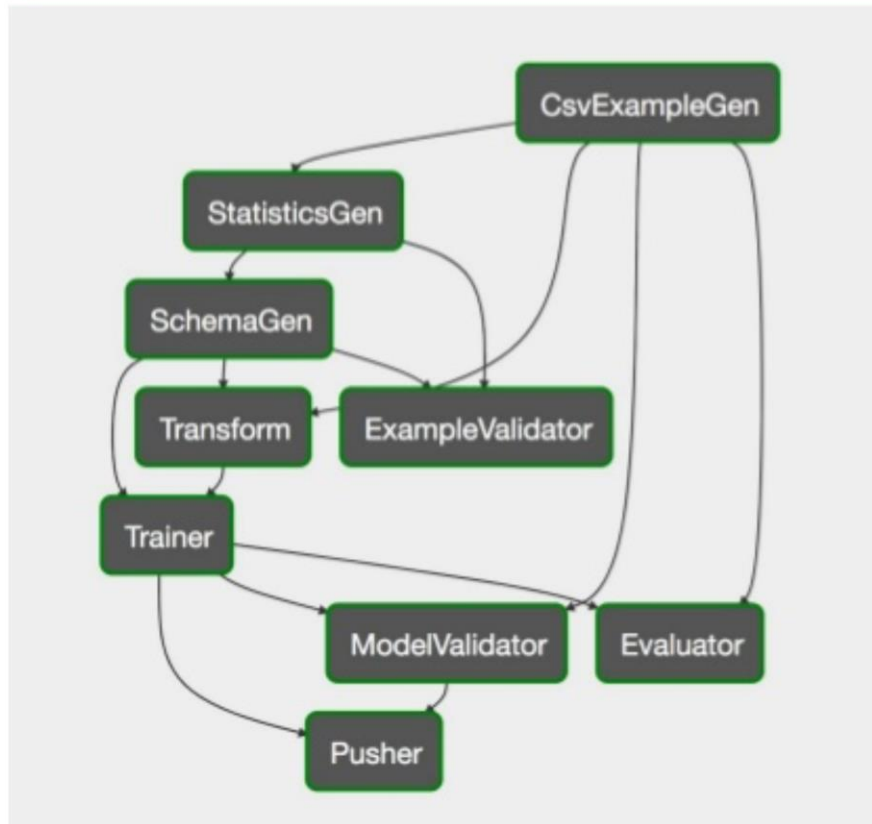
Horizon: Applied Reinforcement Learning Platform

- Robust
- Massively Parallel
- Open Source
- Built on high-performance platforms
 - Spark
 - PyTorch
 - ONNX
- OpenAI Gym & Gridworld Integration tests

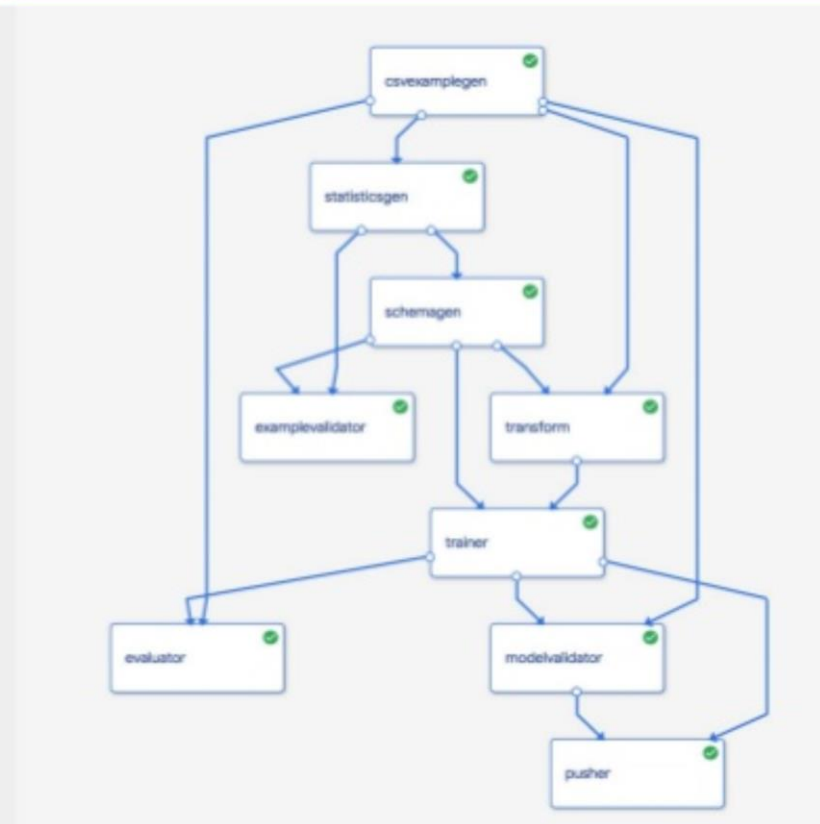


ML WorkFlow

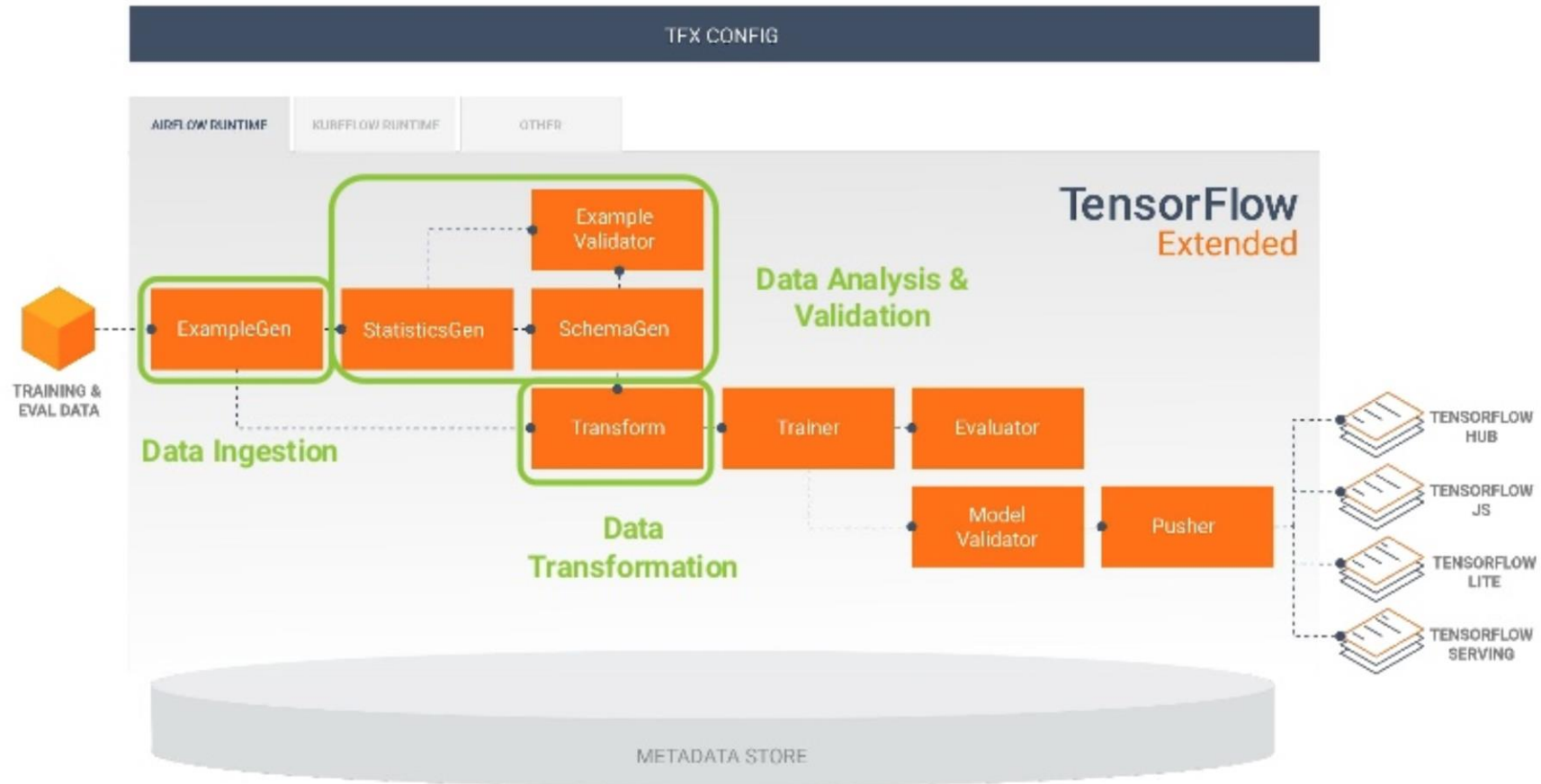
Airflow



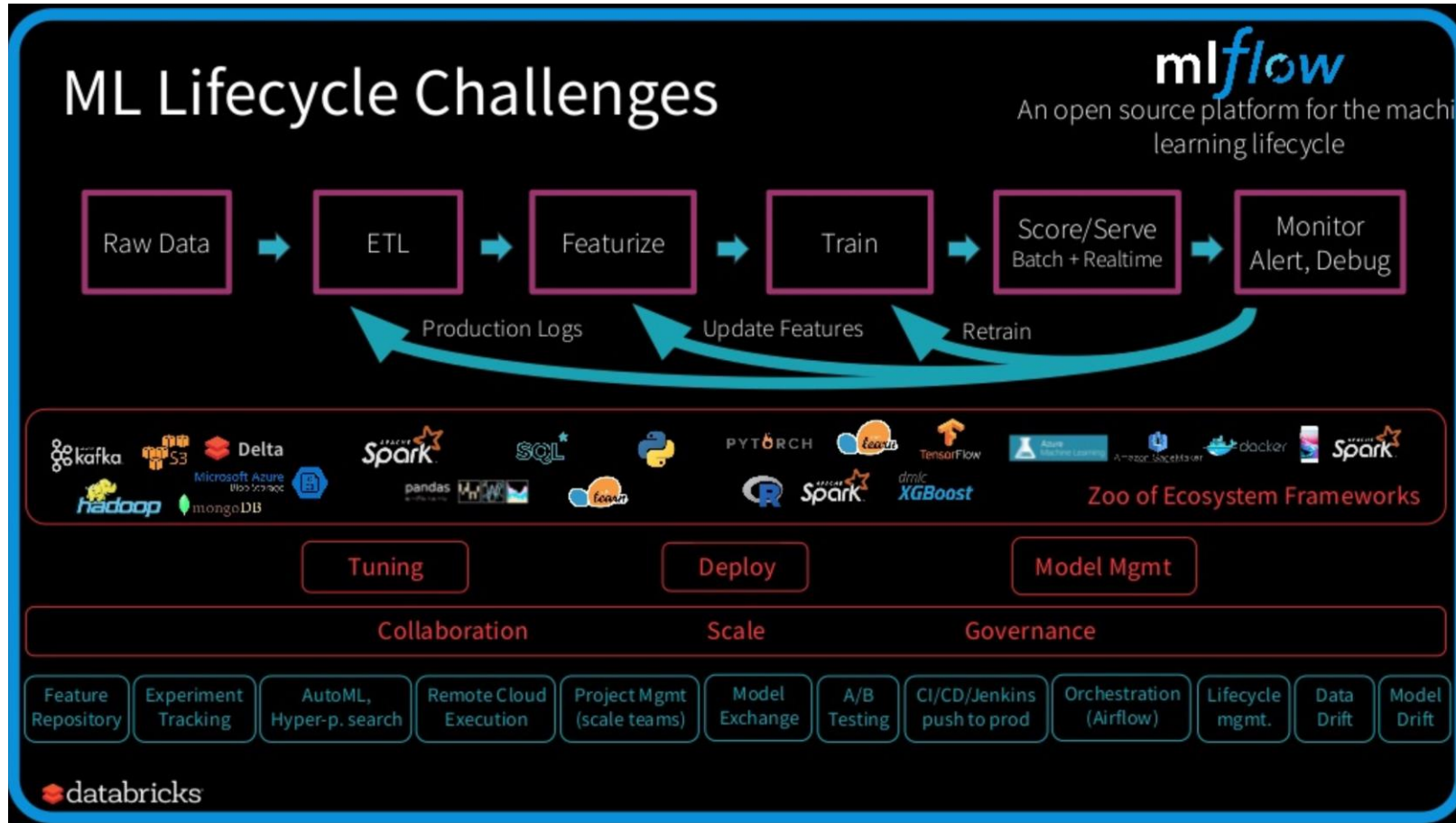
Kubeflow Pipelines



TFX



mflow



mflow

Supported Integrations: June '19

Programming Languages



ML Libraries



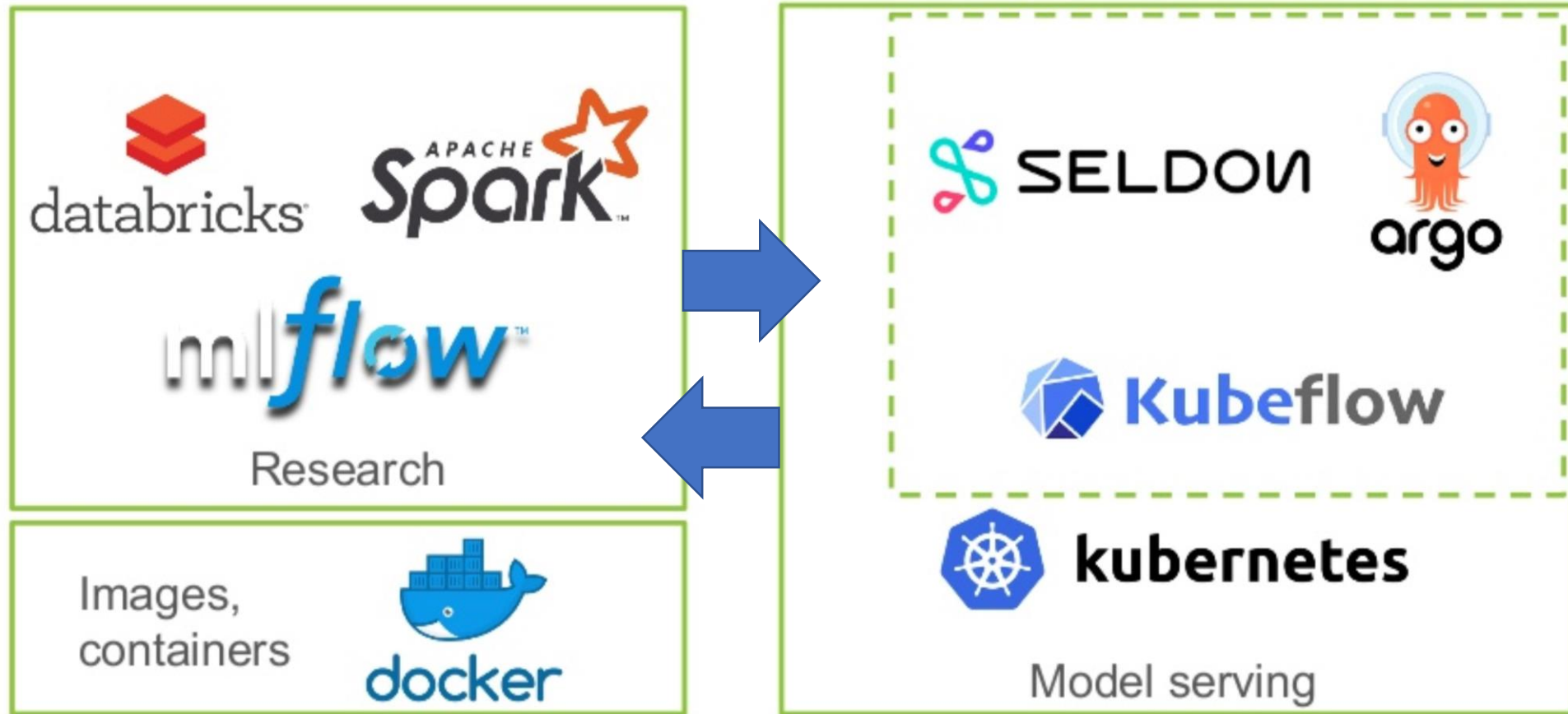
Storage Backends



Deployment Systems

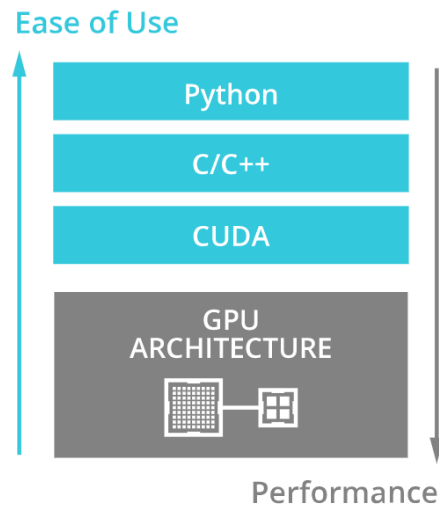


Serving at Scale



Rapids.ai

- End-to-End Data Analytics Pipeline with GPUs
- Apache Arrow , cuDF(on Spark) , cuML(like Pandas)



GPU DATA SCIENCE

i ACCELERATED DATA SCIENCE

The RAPIDS suite of open source software libraries gives you the freedom to execute end-to-end data science and analytics pipelines entirely on GPUs.

[Learn more >>](#)

x SCALE OUT ON GPUS

Seamlessly scale from GPU workstations to multi-GPU servers and multi-node clusters with Dask.

[Learn more about Dask >>](#)

+ PYTHON INTEGRATION

Accelerate your Python data science toolchain with minimal code changes and no new tools to learn.

🎯 TOP MODEL ACCURACY

Increase machine learning model accuracy by iterating on models faster and deploying them more frequently.

🕒 REDUCED TRAINING TIME

Drastically improve your productivity with more interactive data science.

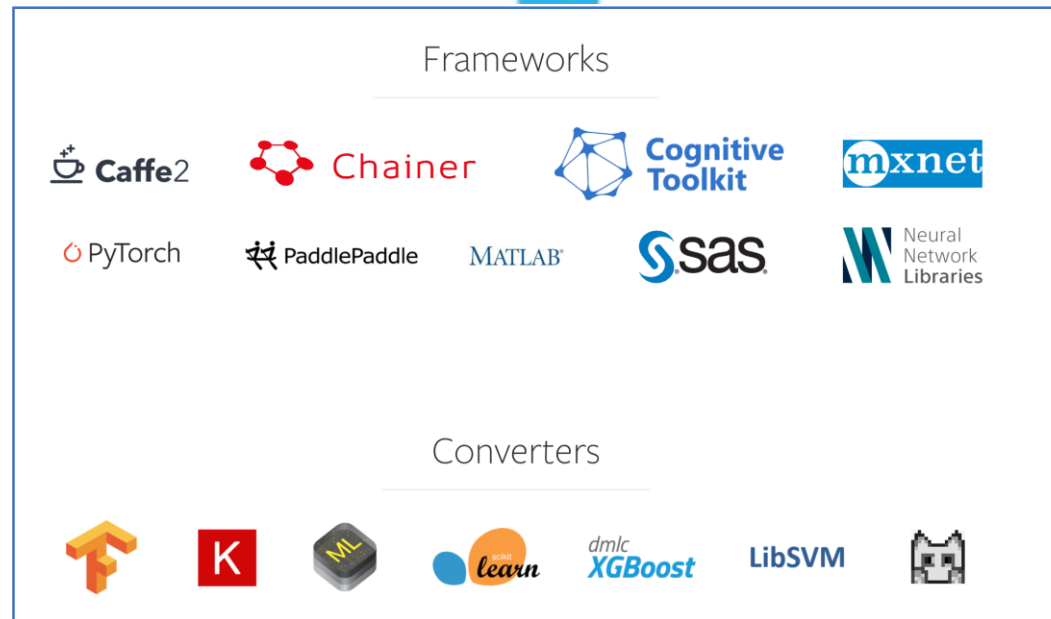
[Learn more about XGBoost >>](#)

🔓 OPEN SOURCE

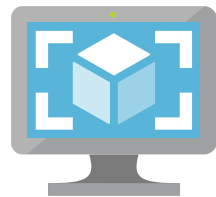
RAPIDS is an open source project. Supported by NVIDIA, it also relies on numba, apache arrow, and many more open source projects.

[Learn more >>](#)

ONNX - Multiple ML/DL Framework Collaboration

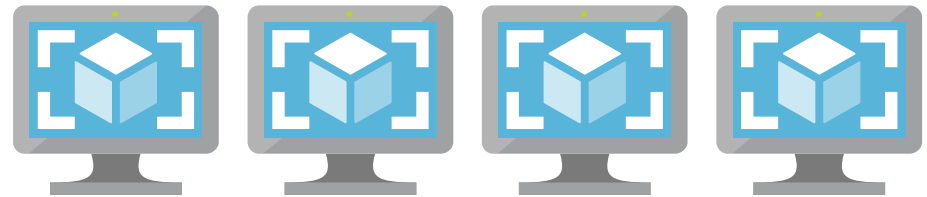
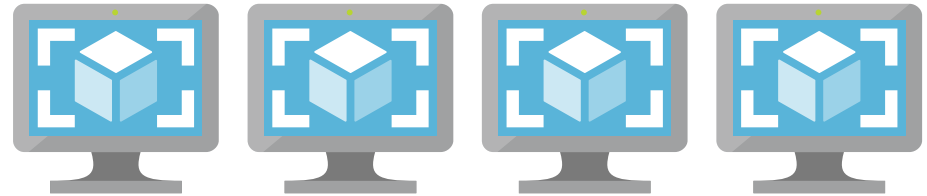
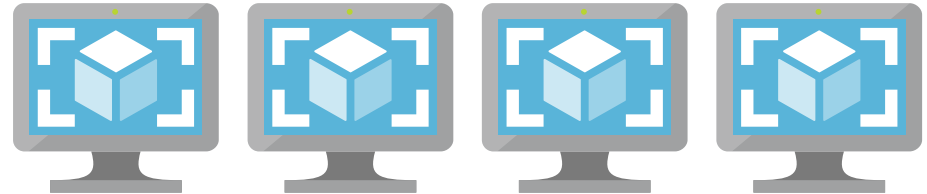


3. Production AI DevOps



현재 직면한 문제

- 여러 개발자의 컴퓨터, 여러 배포 환경에서 실행시켜야 함.



재현성

"한 연구로부터 도출된 결과와 결론을 서로 독립적인 다른 연구를 통해 입증하는 것"

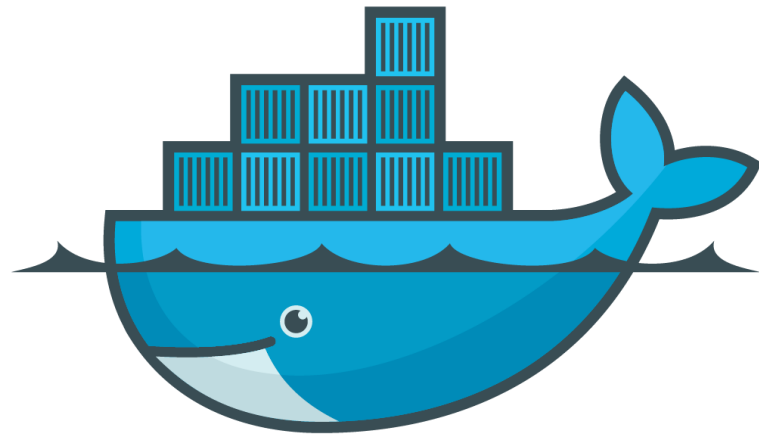
(The confirmation of results and conclusions from one study obtained independently in another)

-출처 : <https://www.sciencemag.org/>

재현성

"한 코드로부터 도출된 **결과와 동작**을 서로 독립적인 다른 머신을 통해 **같은 동작**을 하게 하는 것"

재현성을 위한 강력한 도구 container

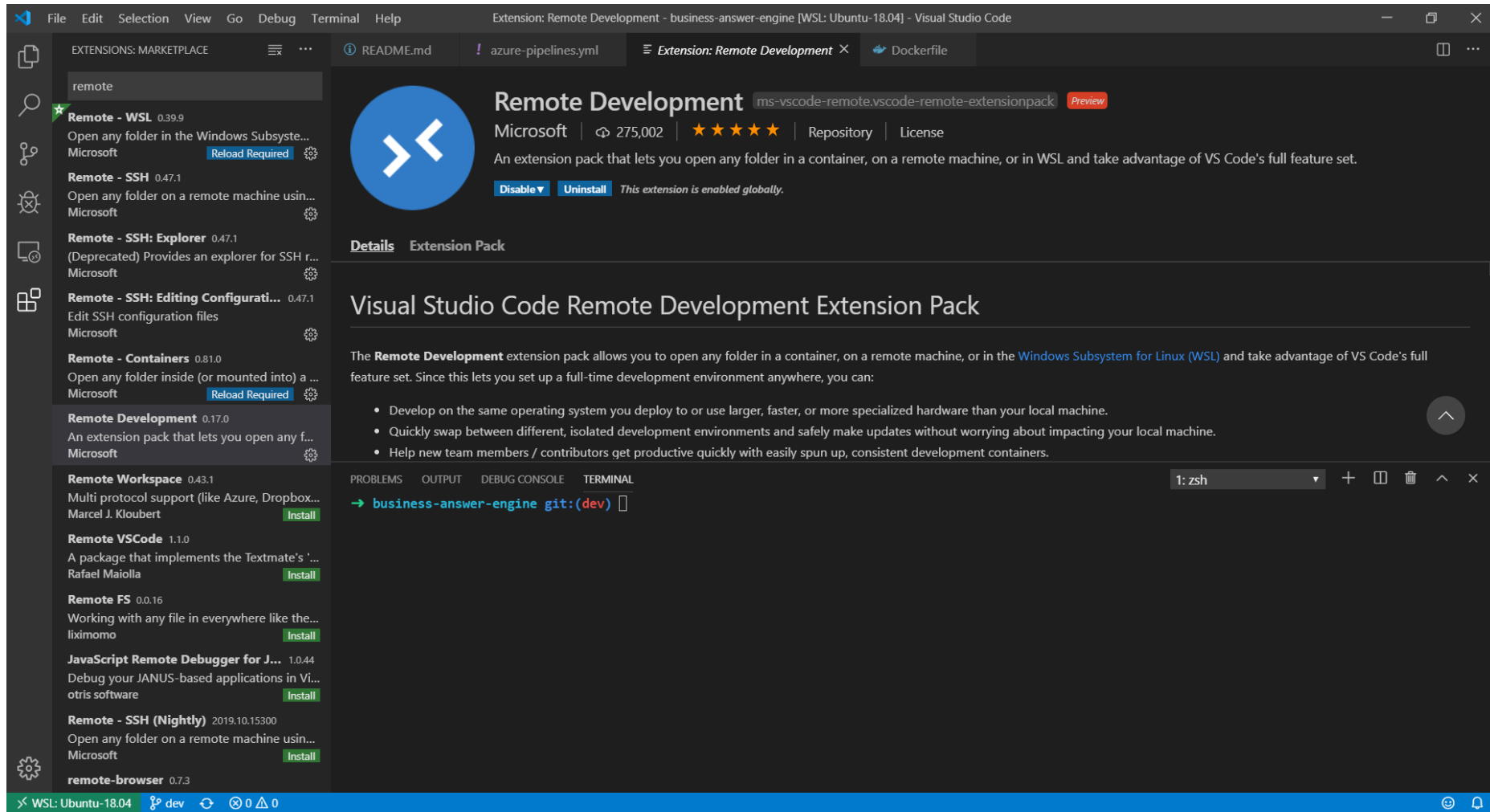


docker

재현성으로 해결 가능한 문제

- 개발 환경 공유
- 운영 환경 동작 보장

개발 환경 재현성을 위한 도구

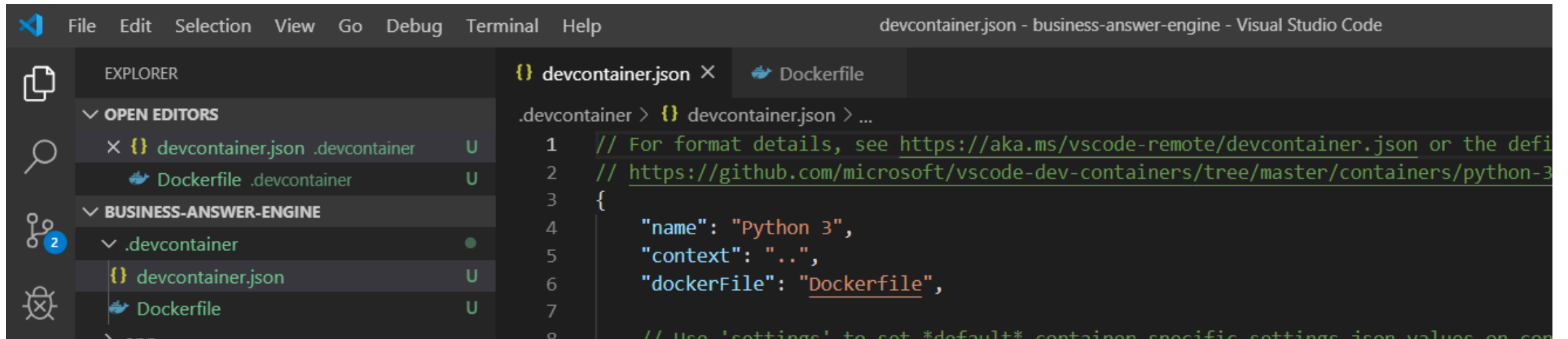


Remote target

- SSH
- Container
- WSL

Remote-container

Devcontainer 설정과 Dockerfile을 통해 개발 환경을 코드로 관리



The screenshot shows the Visual Studio Code interface with a devcontainer.json file open. The Explorer sidebar on the left shows the project structure, including the .devcontainer folder with devcontainer.json and Dockerfile files. The main editor displays the content of devcontainer.json, which is a JSON configuration for a remote container. The configuration includes a name, context, and dockerFile path.

```
devcontainer.json - business-answer-engine - Visual Studio Code
File Edit Selection View Go Debug Terminal Help
EXPLORER
OPEN EDITORS
x {} devcontainer.json .devcontainer U
Dockerfile .devcontainer U
BUSINESS-ANSWER-ENGINE
.devcontainer
{} devcontainer.json U
Dockerfile U
.devcontainer > {} devcontainer.json > ...
1 // For format details, see https://aka.ms/vscode-remote/devcontainer.json or the defi
2 // https://github.com/microsoft/vscode-dev-containers/tree/master/containers/python-3
3 {
4     "name": "Python 3",
5     "context": "..",
6     "dockerFile": "Dockerfile",
7
8     // Use 'settings' to set *default* container specific settings. Values on con
```

DevOps 란

- 지속적 통합과 배포
- 마이크로서비스
- 코드형 인프라스트럭처
- 마이크로 서비스를 위한 모니터링
- 커뮤니케이션

-> 코드형 인프라스트럭처까지 방법론이 성숙되면서 GitOps라는 용어도 발생

출처 : <https://aws.amazon.com/ko/devops/what-is-devops/>

도커가 해주는 것

- 지속적 통합과 배포
- 마이크로서비스
- 코드형 인프라스트럭처
- 마이크로 서비스를 위한 모니터링
- 커뮤니케이션

-> 코드형 인프라스트럭처까지 방법론이 성숙되면서 GitOps라는 용어도 발생

출처 : <https://aws.amazon.com/ko/devops/what-is-devops/>

개발자의 영역

- 지속적 통합과 배포
- 마이크로서비스
- 코드형 인프라스트럭처
- 마이크로 서비스를 위한 모니터링
- 커뮤니케이션

-> 코드형 인프라스트럭처까지 방법론이 성숙되면서 GitOps라는 용어도 발생

출처 : <https://aws.amazon.com/ko/devops/what-is-devops/>

슬랙?

- 지속적 통합과 배포
- 마이크로서비스
- 코드형 인프라스트럭처
- 마이크로 서비스를 위한 모니터링
- 커뮤니케이션

-> 코드형 인프라스트럭처까지 방법론이 성숙되면서 GitOps라는 용어도 발생

출처 : <https://aws.amazon.com/ko/devops/what-is-devops/>

이것을 위한 도구가 필요

- 지속적 통합과 배포
- 마이크로서비스
- 코드형 인프라스트럭처
- 마이크로 서비스를 위한 모니터링
- 커뮤니케이션

-> 코드형 인프라스트럭처까지 방법론이 성숙되면서 GitOps라는 용어도 발생

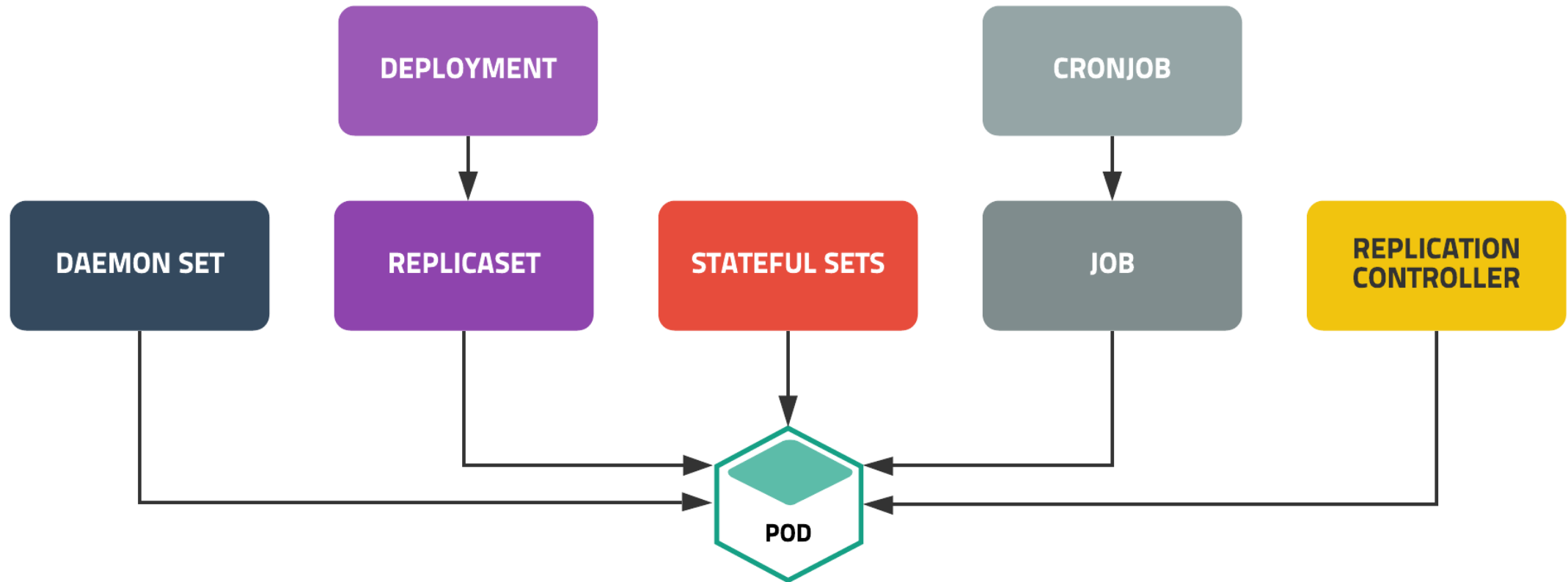
출처 : <https://aws.amazon.com/ko/devops/what-is-devops/>

자동 운영환경을 위한 kubernetes



kubernetes

Pod을 관리하는 상위 개념들



선언적 설정으로 인프라 관리

```
1  apiVersion: apps/v1beta1
2  kind: Deployment
3  metadata:
4    name: cookiemap
5  spec:
6    replicas: 1
7    template:
8      metadata:
9        labels:
10         app: cookiemap
11      spec:
12        hostAliases:
13         - ip: "52.231.34.29"
14           hostnames:
15             - "internal.commbot.nginx"
16        containers:
17         - name: cookiemap
18           image: acrcommbotprdinternal.azurecr.io/cookiemap:latest
19           imagePullPolicy: Always
20           imagePullSecrets:
21             - name: commbotprdsecret
22           ports:
23             - containerPort: 8080
24           envFrom:
25             - configMapRef:
26               name: configmap-cookiemap
28  apiVersion: v1
29  kind: Service
30  metadata:
31    name: cookiemap
32  spec:
33    sessionAffinity: ClientIP
34    ports:
35     - port: 80
36       targetPort: 8080
37    selector:
38     app: cookiemap
--
40  apiVersion: autoscaling/v1
41  kind: HorizontalPodAutoscaler
42  metadata:
43    name: cookiemap
44  spec:
45    scaleTargetRef:
46     apiVersion: apps/v1
47     kind: Deployment
48     name: cookiemap
49    minReplicas: 1
50    maxReplicas: 10
51    targetCPUUtilizationPercentage: 70
```

배포

- Azure Devops의 Pipelines의 기능을 활용

The screenshot shows the Azure DevOps interface for a pipeline named 'cookiemap-prd-k8s'. The left sidebar contains navigation options: Overview, Boards, Repos, Pipelines (selected), Builds, Releases, Library, Task groups, Deployment groups, Test Plans, and Artifacts. The main content area displays the pipeline configuration for 'cookiemap-prd-k8s' under the 'Pipeline' tab. It shows two artifact boxes: '_cookiemap' and '_k8s-setting', both with a lightning bolt icon indicating they are artifacts. A 'Stage 1' box is connected to the '_cookiemap' artifact, showing '1 job, 2 tasks'. A 'Schedule not set' icon is also visible at the bottom of the artifact section.

배포

- 새로운 이미지 배포를 위해 latest 태그 사용 지양

All pipelines > 1_prd > cookiemap-prd-k8s Save Create release View releases ...

Pipeline **Tasks** ▾ Variables Retention Options History

Stage 1
Deployment process ...

Agent job
 Run on agent +

Command Line Script
Command line ✓ ⋮

kubectl apply
Deploy to Kubernetes

Command line ⓘ View YAML Remove

Task version ▾

Display name *

Script *

```
sed -i "s:/latest/:$(Version)/g" $(System.DefaultWorkingDirectory)/_k8s-setting/cookiemap/cookiemap-prd.yml
```

배포

- Kubectl 명령을 devops 내에서만 실행

Pipeline **Tasks** ▾ Variables Retention Options History

Stage 1
Deployment process

Agent job
Run on agent

Command Line Script
Command line

kubectl apply
Deploy to Kubernetes

Deploy to Kubernetes ⓘ [View YAML](#) [Remove](#)

Task version 1.* ▾

Display name *
kubectl apply

Kubernetes Cluster ^

Service connection type * ⓘ
Kubernetes Service Connection ▾

Kubernetes service connection * ⓘ | [Manage](#)

CON-AKS-COMMBOT-PRD-INTERNAL ▾ ↻ + New

Namespace ⓘ

Commands ^

배포

- 배포 기록 관리
- 기록 기반 롤백

Releases Deployments Analytics

Releases	Created	Stages
Release-29 🏰 20191016.2 🌐 dev	2019. 10. 16. 오후 10:55:15	✔ Stage 1
Release-28 🏰 20191016.1 🌐 dev	2019. 10. 16. 오후 10:51:06	✔ Stage 1
Release-27 🏰 20191014.10 🌐 dev	2019. 10. 14. 오후 10:57:26	✔ Stage 1
Release-26 🏰 20191014.8 🌐 dev	2019. 10. 14. 오후 10:31:31	✔ Stage 1
Release-25 🏰 20191014.7 🌐 master	2019. 10. 14. 오후 9:46:04	✔ Stage 1
Release-24 🏰 20191014.6 🌐 master	2019. 10. 14. 오후 9:40:20	✔ Stage 1
Release-23 🏰 20191014.5 🌐 dev	2019. 10. 14. 오후 9:40:19	✔ Stage 1
Release-22 🏰 20191014.4 🌐 master	2019. 10. 14. 오후 9:29:56	✔ Stage 1
Release-21 🏰 20191014.3 🌐 master	2019. 10. 14. 오후 9:22:39	✔ Stage 1

환경별 변수 관리

- Library에 글로벌 변수 작성, 파이프라인 별 적용

The screenshot displays the 'Variables' configuration page in Azure DevOps. The left sidebar contains a navigation menu with the following items: Overview, Boards, Repos, Pipelines (selected), Builds, Releases, Library, Task groups, Deployment groups, Test Plans, and Artifacts. The main content area has tabs for Pipeline, Tasks, Variables (selected), Retention, Options, and History. Under the 'Variables' tab, there are sections for Pipeline variables, Variable groups (highlighted), and Predefined variables. The 'Variable groups' section lists two groups:

Name	Value
▼ k8s_service_alias (11)	Scopes: Release
▼ database_config_dev (21)	Scopes: Release

At the bottom of the variable groups list, there are two buttons: 'Link variable group' and 'Manage variable groups'.

환경별 변수 관리

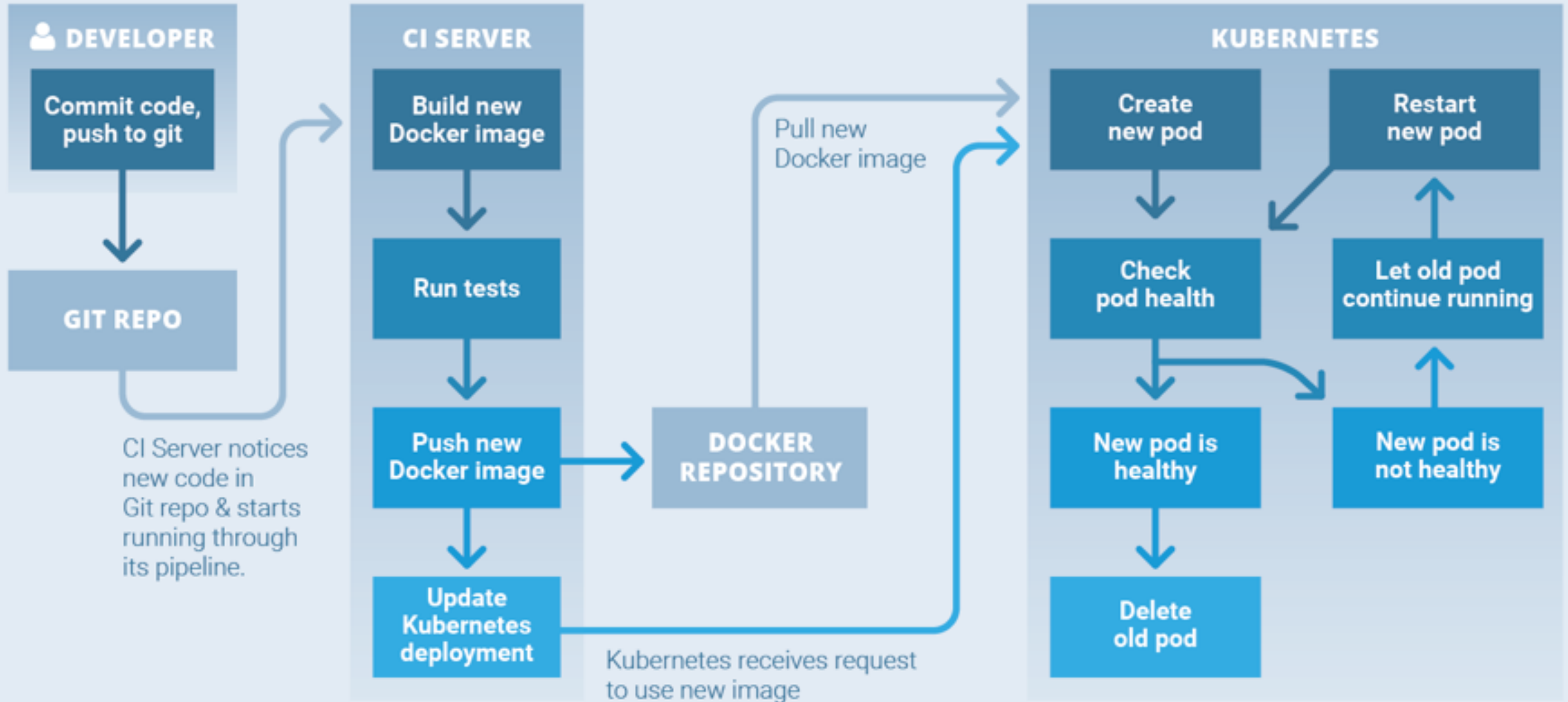
- 코드 내에서는 env 호출로 사용

config / **config.py**

[Contents](#) [History](#) [Compare](#) [Blame](#) | [Edit](#) [Rename](#) [Delete](#) [Download](#)

```
1 import os
2
3 class log_connection_info:
4     port = os.environ['ES_PORT']
5     host = os.environ['ES_URL']
6     elastic_user = os.environ['ES_USER']
7     elastic_pw = os.environ['ES_PASSWORD']
8     multi_turn_log_index = os.environ['ES_INDEX_MULTITURN']
9     user_info_log_index = os.environ['ES_USERINFO']
10    doc_type='_doc'
11
12
13 class redis_connection_info:
14     host = os.environ['REDIS_CACHE_URL']
15     port = os.environ['REDIS_CACHE_PORT']
16     db = os.environ['REDIS_CACHE_DB']
17     password = os.environ['REDIS_CACHE_KEY_PRIMARY']
18     if os.environ['REDIS_CACHE_SSL'] == "TRUE":
19         ssl = True
20     else:
21         ssl = False
22
23
```

CI/CD Pipeline Workflow with Kubernetes



Function as a Service Platform

- Kubernetes 위에서 동작하는 opensource serverless 도구들
- 개발자가 더욱 로직에 집중할 수 있게 자동화 가능



출처 : <https://landscape.cncf.io/format=serverless&fullscreen=yes&zoom=150>

Our Open Source

Korean BERT pre-trained cased (KoBERT)

SKTBrain / KoBERT

Unwatch 10 Unstar 230 Fork 32

Code Issues 0 Pull requests 0 Projects 0 Wiki Security Insights

Korean BERT pre-trained cased (KoBERT)

korean-nlp language-model

10 commits 1 branch 0 releases 1 contributor Apache-2.0

Branch: master New pull request Create new file Upload files Find file Clone or download

haven-jeon Update README.md Latest commit 21c6de6 16 days ago

imgs	initial commit	22 days ago
kobert	add pkg install script	22 days ago
logs	initial commit	22 days ago
scripts/NSMC	add MXNet example	16 days ago
LICENSE	Update LICENSE	22 days ago
README.md	Update README.md	16 days ago
requirements.txt	add pkg install script	22 days ago
setup.py	add pkg install script	22 days ago

README.md

- Korean BERT pre-trained cased (KoBERT)
 - Why?'
 - Training Environment
 - Requirements
- How to install
- Using with PyTorch
- Using with ONNX
- Using with MXNet-Gluon

<https://github.com/SKTBrain/KoBERT>

USE Case 소개(상세 내용은 ...)

- 4 CPU * 10 Docker on Cloud
- 1.5GB LSTM Tensorflow NLP 모델 실시간 Serving (batchsize = 1)
- 랜덤 쿼리 (캐쉬 의미 없음)
- Deep Learning Inference 1만 TPS 구현해 보기.
 - NginX + ASGI + FastAPI + asyncio + aiohttp + tornado
 - Pandas 건너내기, MemoryView 활용하기.
 - Distributed Shared Session으로 Redis 사용.

Q&A

김훈동 : <https://www.facebook.com/kim.hoondong>

박찬엽 : <https://www.facebook.com/mrchypark>