# Data-centric MLOps

## : 데이터 중심 MLOps를 돕기 위한 작은 장치들

Superb AI 이정권

Superb AI

# AI / ML = Model + Data

Superb AI

# AI / ML = Model + Data

# Data centric?

# A Chat with Andrew on MLOps:
## From Model-centric to
## Data-centric AI



Task

Baseline:

70% accuracy

⬇

Target Performance:

90% accuracy

**Should the team improve**

**the code or the data?**

**: code(20%), data(80%)**

Superb AI

A Chat with Andrew on MLOps:
From Model-centric to
Data-centric AI

# Improve AI →
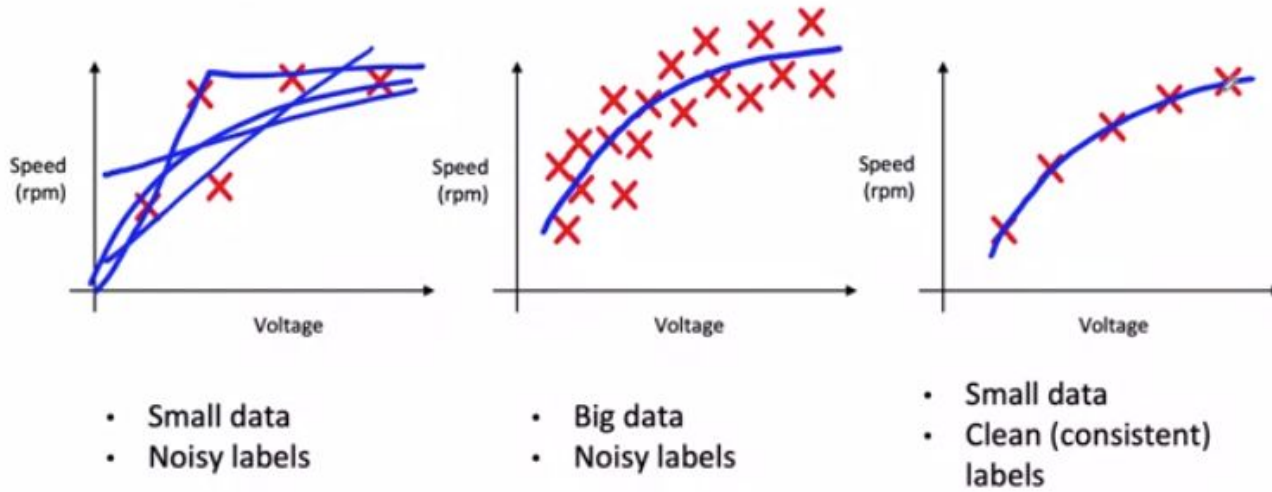# Improve the quality of the data:
consistency
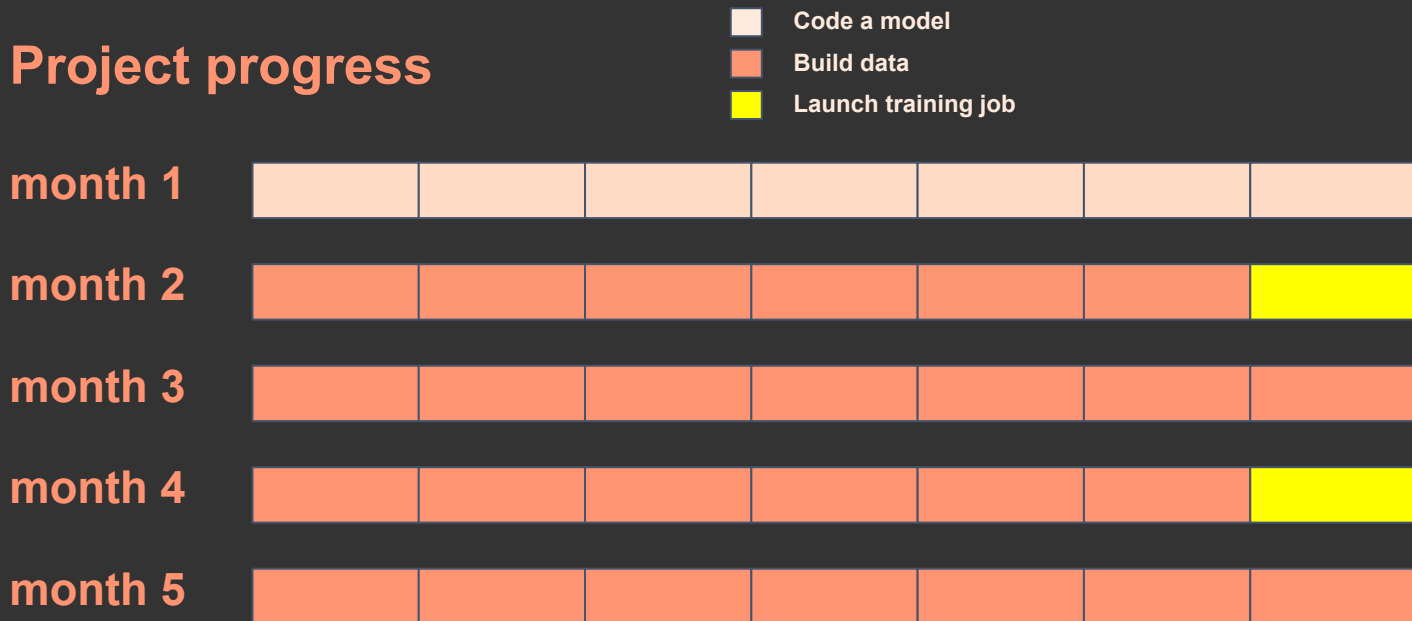error rate
diversity
coverage
feedback frequency
size

...

**☐ Superb AI**

# A Chat with Andrew on MLOps:
# From Model-centric to
# Data-centric AI

Superb AI

# 사실은, 늘 해오던 일

**Project progress**

Code a model
Build data
Launch training job

month 1

month 2

month 3

month 4

month 5

Superb AI

# 사실은, 늘 해오던 일

## Building the Software 2.0 Stack (Andrej Karpathy, 2018)

# Answer:
# 100,000 images?

# My Answer:
# I don't know. Let's start from 5,000
# WHY?

여전히, 잘 모른다
→ **Data-centric MLOps**
**Systematic & iterative way to build Data for ML**

단순히 지루한 작업을 자동화하는 과정이 아닌

**ML 문제를 해결**하기 위한 과정

저는 **Superb AI**라는 팀에서 이 문제를 풀고 있습니다.

# The Problem

## <2달

## <30명

## <20,000 Images

# Starting Point

## Labeling Tool



## Data



## Label



Superb AI

# Reusable Data Spec



```
{

   project_name:
     potato_detect_1

   data_spec:
     good_potato:
       box:
         color: red
         condition: ...
       bad_potato:
         box:
}
```



```
{

   project_name:
     potato_detect_2

   data_spec:
     good_potato:
       polygon:
         color: red
         condition: ...
       bad_potato:
         box:
}
```
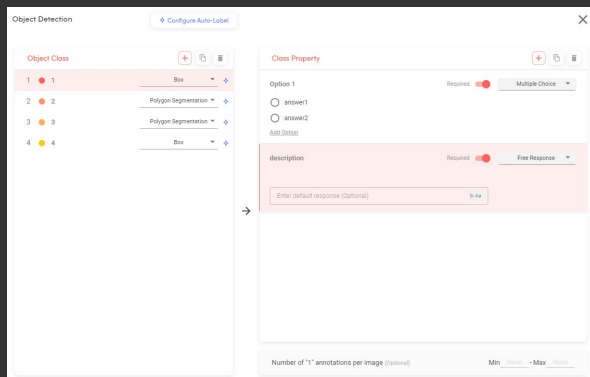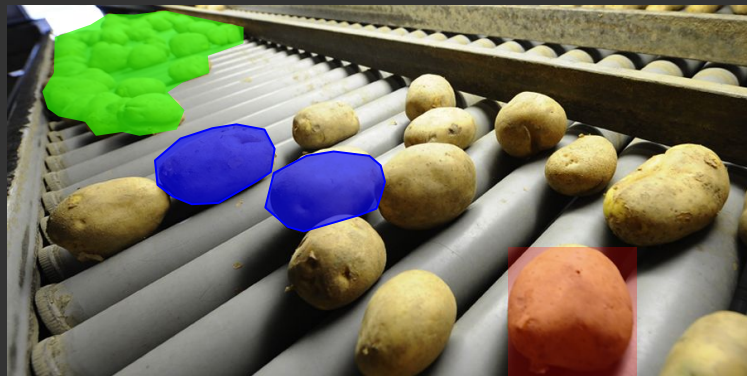
# Reusable Data Spec

```
{

    project_name:
      potato_detect_13

    data_spec:
      best_potato:
        polygon:
        direction:
        options: ...
      good_potato: {}
      normal_potato: {}
      bad_potato: {}


}
```





Goal ≠ Task

**ALWAYS configured repeatedly**

**name,
color,
type,
conditions,
options,
property,
ROI Info,
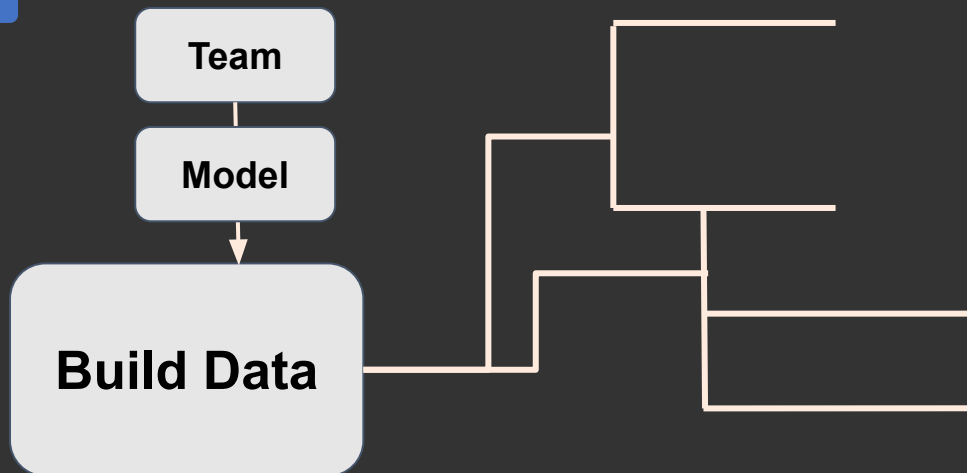...**

□ Superb AI

# Support flexible pipeline

**WORKING**  **SUBMITTED**  **REVIEWED**

## 100 different problems, 100 different datasets, 100 different ways

## To support flexible pipeline

Team

Model

**Build Data**

Superb AI

# Support flexible pipeline

# Versioning

## Set 단위, 실험 당

# Detailed Statistics & Report

ML Engineer를 위해 … ?



Superb AI

# Human in the loop ^ 2



**Human in the loop ML**

**Supervised Machine Learning**

| TRAIN MODEL | TUNE MODEL |

Data

Human Intelligence

Edge Cases — More Labels

Data Annotation

Active Learning

Low

Training Data → Model → Confident → O/P

High

□ Superb AI

# Keep labels consistent

# Keep labels consistent

# 요약

| | Data Key (Name) | Dataset |
|---|---|---|
| ☐ | | |
| ☐ | /얼굴/thai-child-116002_1920.jpg | image_face |

| | | Data Key (Name) | Dataset | Label Tag | Assignee | Issues | Last Updated | Auto-Label | QA | Status |
|---|---|---|---|---|---|---|---|---|---|---|
| ☐ ▾ | | | | | | | | | | |
| ☐ | | /vehicle/sedan2.jpg (3) | best_vehicle | | - | 0 | 11 days ago | | | ● IN PROGRESS |

2

# 마무리

**Source data analysis, User analysis,
Log, Task matching, etc
여전히 할일이 정말 많다.**

**SDK를 이용한 사용 예제!는 다음에**
https://github.com/superb-AI-Suite/
**Full-pipeline MLOps**
https://ai-infrastructure.org/