# DATUMO

The ultimate all-in-one data platform

# All-in-one data platform



| Est. 2018.11 | Clients 227+ | Investment received 11M USD | Crowd-workers 240K+ | Core staffs 100 |

Nominated as
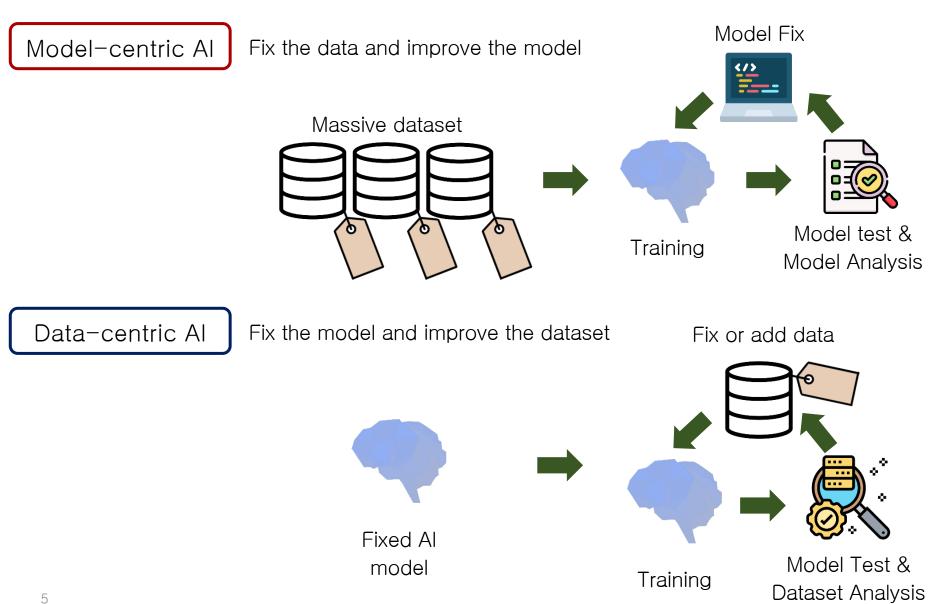2021 Forbes 30 Under 30 Asia,
Enterprise Technology

# From Model-centric to Data-centric AI

2021.03

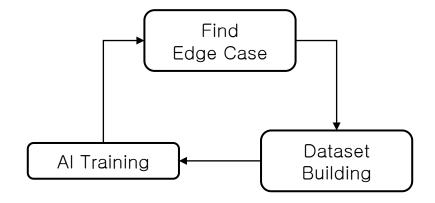Andrew Ng (World-renowned AI scholar)

It is changing from "how to make AI" to "how to make data".

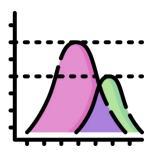Model-centric AI — Fix the data and improve the model

Model Fix

Massive dataset → Training → Model test & Model Analysis

Data-centric AI — Fix the model and improve the dataset

Fix or add data

Fixed AI model → Training → Model Test & Dataset Analysis

# Why Data-centric AI?



Find
Edge Case

Dataset
Building

AI Training

If AI performance issues
stem from data,
then the solution must
come from the data itself.

Repeating edge case
handling reduces AI
performance gap between
development and service.

Data shifts cause AI
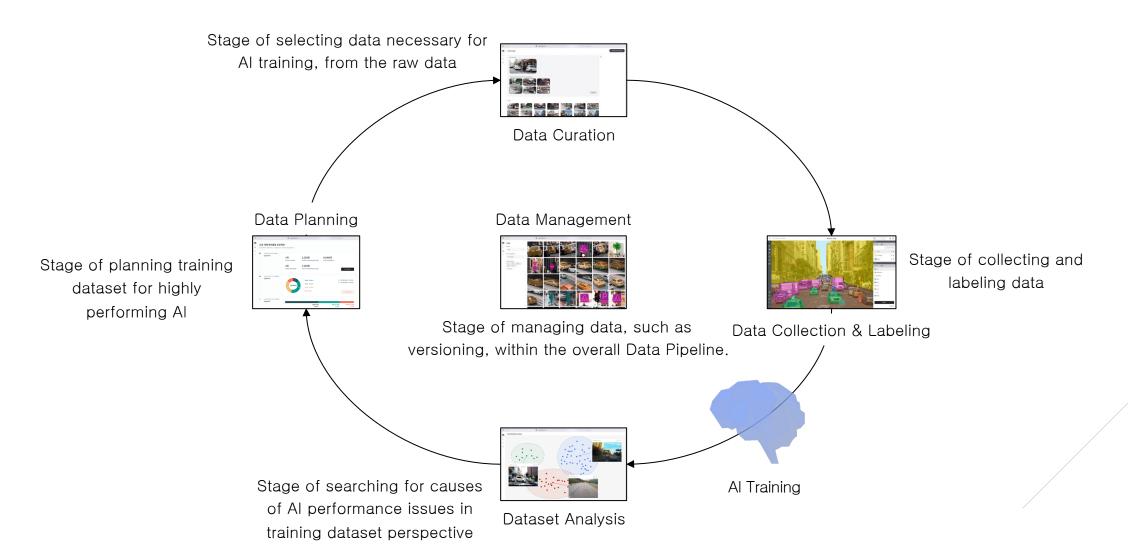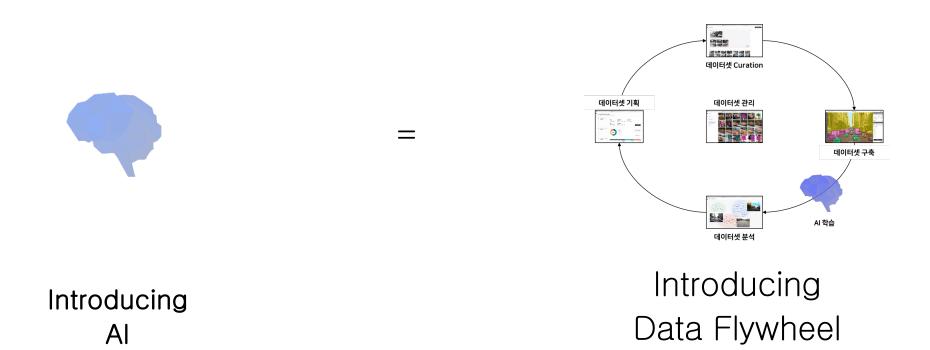performance degradation
in service operation.

Data-centric AI is essential for
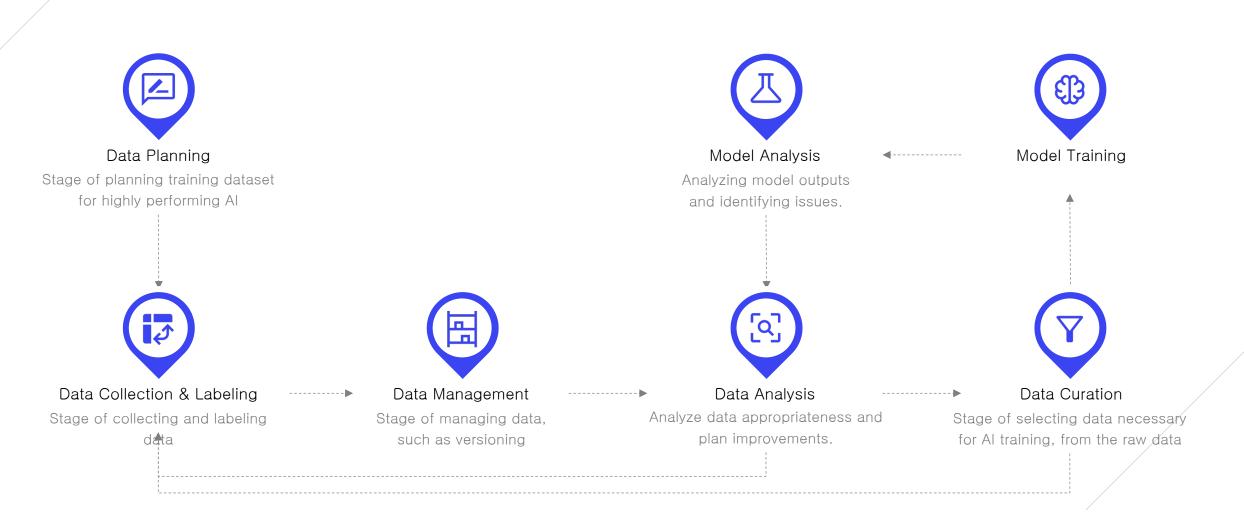both creating a custom LLM and utilizing fine-tuning.

# AI Data Flywheel

Stage of selecting data necessary for AI training, from the raw data

Data Curation

Stage of planning training dataset for highly performing AI

Data Planning

Data Management

Stage of collecting and labeling data

Stage of managing data, such as versioning, within the overall Data Pipeline.

Data Collection & Labeling

Stage of searching for causes of AI performance issues in training dataset perspective

Dataset Analysis

AI Training

As LLM and other AI API offerings continue to expand...



=



Introducing
AI

Introducing
Data Flywheel

# AI Data Flywheel

**Data Planning**
Stage of planning training dataset
for highly performing AI

**Model Analysis**
Analyzing model outputs
and identifying issues.

**Model Training**

**Data Collection & Labeling**
Stage of collecting and labeling
data

**Data Management**
Stage of managing data,
such as versioning

**Data Analysis**
Analyze data appropriateness and
plan improvements.

**Data Curation**
Stage of selecting data necessary
for AI training, from the raw data

DATUMO

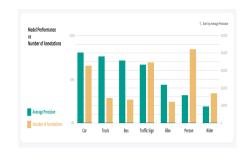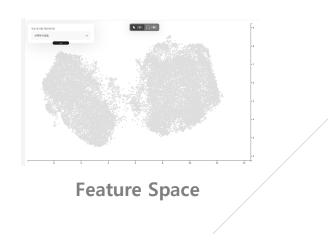# Data Analysis

Analyze data appropriateness. E.g. **coverage**

*Location of image capture*

*Time of image capture*

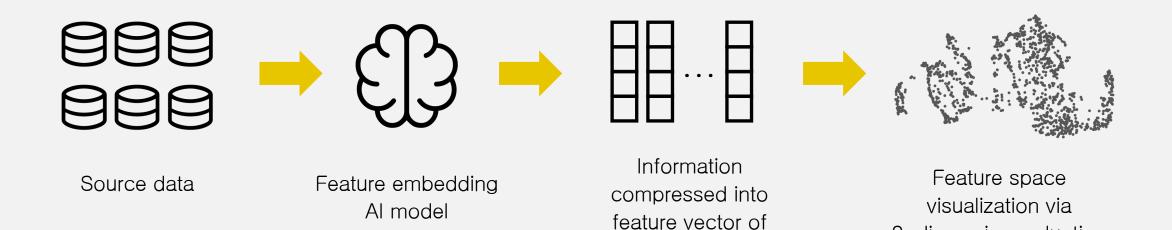*Age groups of provided vocal data*

*Type of voice recording device*

**Metadata**

**Statistics of labels per class**

**Feature Space**

# What is feature space?



Source data

Feature embedding
AI model

Information
compressed into
feature vector of
N-dimension

Feature space
visualization via
2-dimension reduction
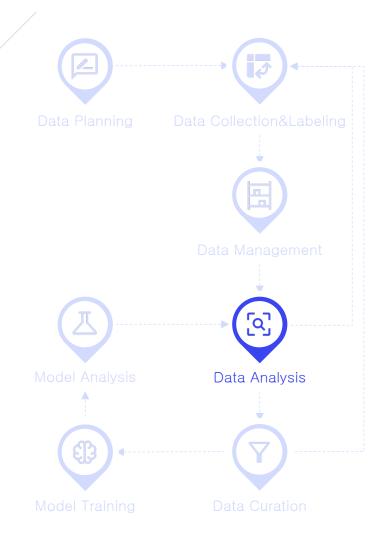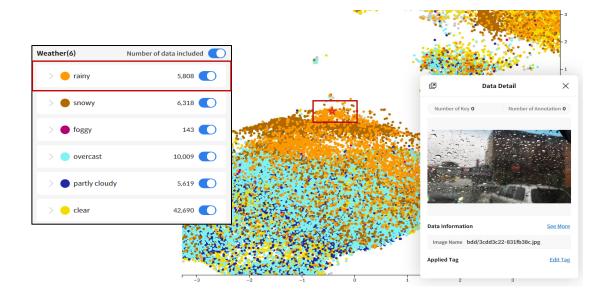
When feature is extracted from data, the data with similar characteristics
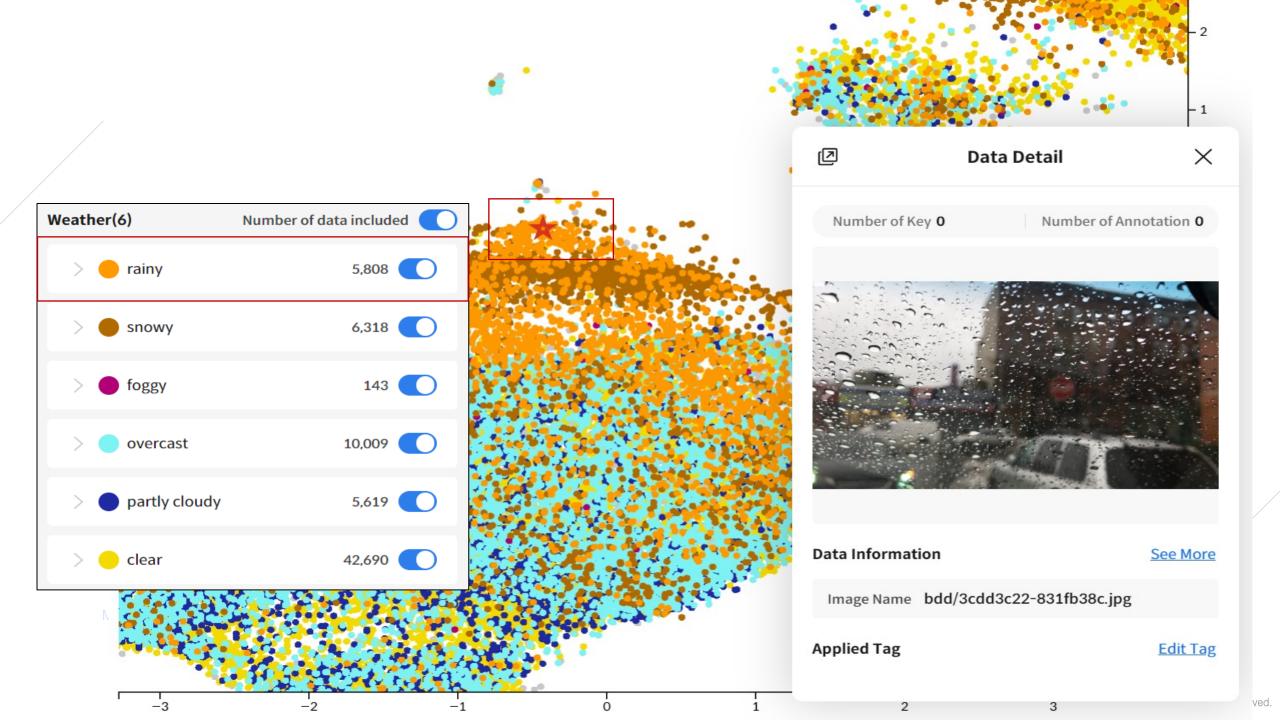in the perspective of the trained AI model tend to stay close in feature space

# Data Analysis

**AI scenario analysis according to feature space distribution based on metadata**

| Weather(6) | Number of data included | |
|---|---|---|
| rainy | 5,808 | |
| snowy | 6,318 | |
| foggy | 143 | |
| overcast | 10,009 | |
| partly cloudy | 5,619 | |
| clear | 42,690 | |

**Data Detail**

Number of Key **0** | Number of Annotation **0**

**Data Information**  See More

Image Name  bdd/3cdd3c22-831fb38c.jpg

**Applied Tag**  Edit Tag

# Data Analysis

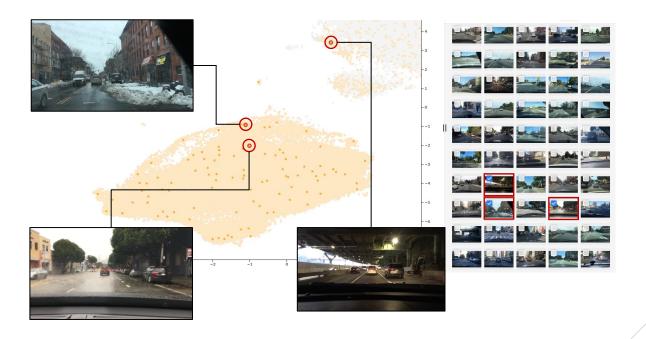**Scenario analysis by reviewing main images chosen via mathematical curation to increase coverage**

# Data Curation

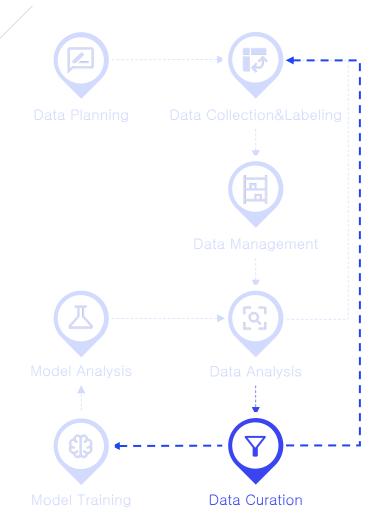## Curate dataset with coverage in mind, label, and train AI model

# Model & Data Analysis

**Combine multiple information such as model metric, data split, and others as queries and mark in color on feature space.**

**Query Detail**
Set the conditions included in query.

| data_split ⌄ | has | val ⌄ | 🗑 |
| and ⌄ | rider_AP ⌄ | ≤ ⌄ | 0.2 | 🗑 |

Add Condition                                    De

Add Condition Group

**Display-Query List**                    Edit Display
List of queries to configure display view. You can add and manage queries.

**Edge Case Analysis(3)**    Number of data included 🔵

| 🔴 val_rider AP<=0.2 | 1,091 | 🔵 |
| 🔵 train_rider | 7,842 | 🔵 |
| 🟢 unlabeled | 55,890 | 🔵 |
| Fixed Others | 15,040 | ⚪ |

*Set data with poor model metric
from rider class,
among validation set, as red*

Data Planning

Data Collection&Labeling

Data Management

Model Analysis

Data Analysis

Model Training

Data Curation

# Model & Data Analysis

Define edge cases by finding common features among data

Define edge case as:
"Dark from the shadows under bridges"

↓

Notice lack of train set around edge cases

↓

Check additional data within the area that may be additionally labeled and trained

**Display-Query List**          Edit Display

List of queries to configure display view. You can add and manage queries.

**Edge Case Analysis(3)**     Number of data included ⬤

> ⬤ val_rider AP<=0.2          1,091    ⬤

> ⬤ train_rider               7,842    ⬤

> ⬤ unlabeled                55,890    ⬤

⬤ Fixed Others               15,040    ⬤

Data Planning

Data Collection&Labeling

Data Management

Model Analysis

Data Analysis

Model Training

Data Curation

# Model & Data Analysis

**Curation of data similar to edge cases from unlabled data and labeling them**

Select edge case image

Similar data search result

# AI Data Flywheel with DATUMO SCOPE

**Data Planning**
Stage of planning training dataset for highly performing AI

**Model Analysis**
Analyzing model outputs and identifying issues.
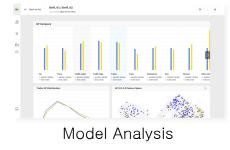
**Model Training**

**Data Collection & Labeling**
Stage of collecting and labeling data
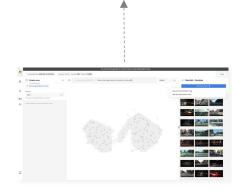
**Data Management**
Stage of managing data, such as versioning

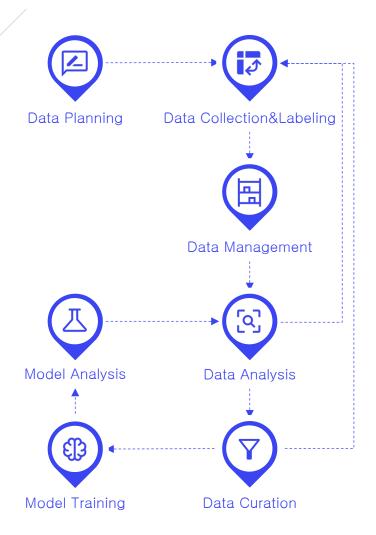**Data Analysis**
Analyze data appropriateness and plan improvements.
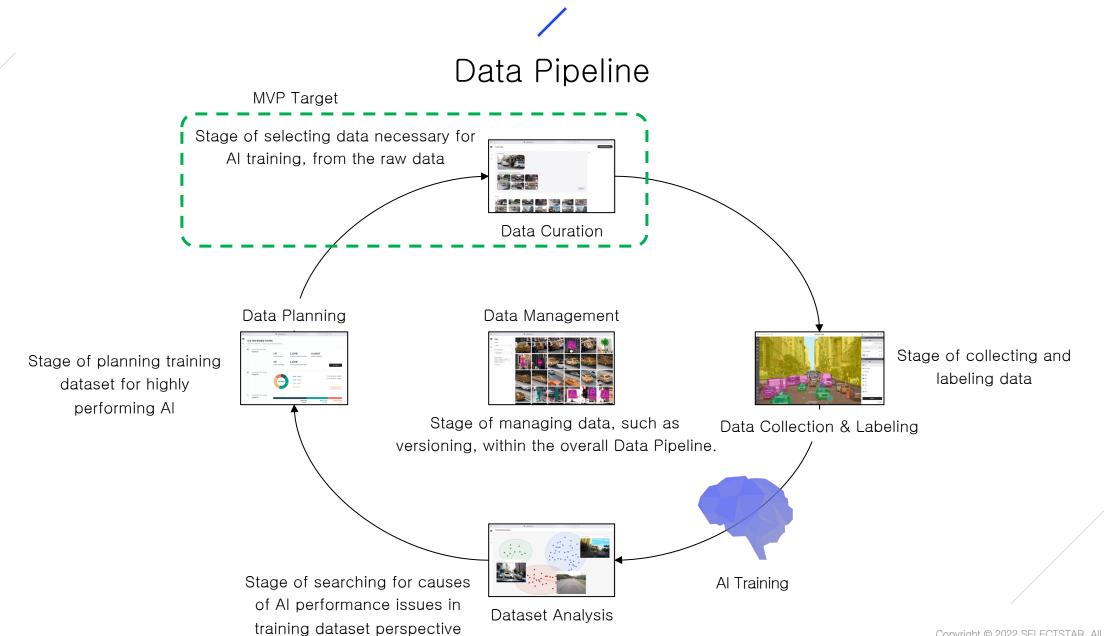
**Data Curation**
Stage of selecting data necessary for AI training, from the raw data

Build, manage, and analyze AI training data

using DATUMO SCOPE,

and deploy AI services without development.

# Data Pipeline

MVP Target

Stage of selecting data necessary for
AI training, from the raw data

Data Curation

Data Planning

Stage of planning training
dataset for highly
performing AI

Data Management

Stage of managing data, such as
versioning, within the overall Data Pipeline.

Stage of collecting and
labeling data

Data Collection & Labeling

AI Training

Stage of searching for causes
of AI performance issues in
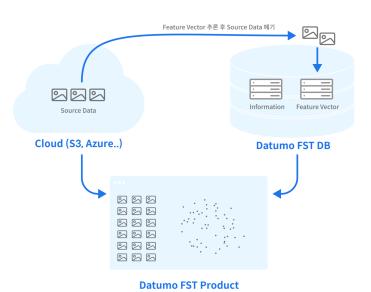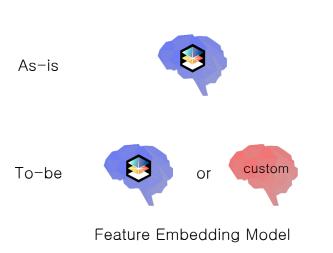training dataset perspective

Dataset Analysis

# Datumo Feature Space Tool

Integration with user's AI lifecycle

No need to upload source data
when utilizing Cloud Integration

Upload feature vector extracted
from user's AI model
('23.01 Release)

Enable the use of intranet
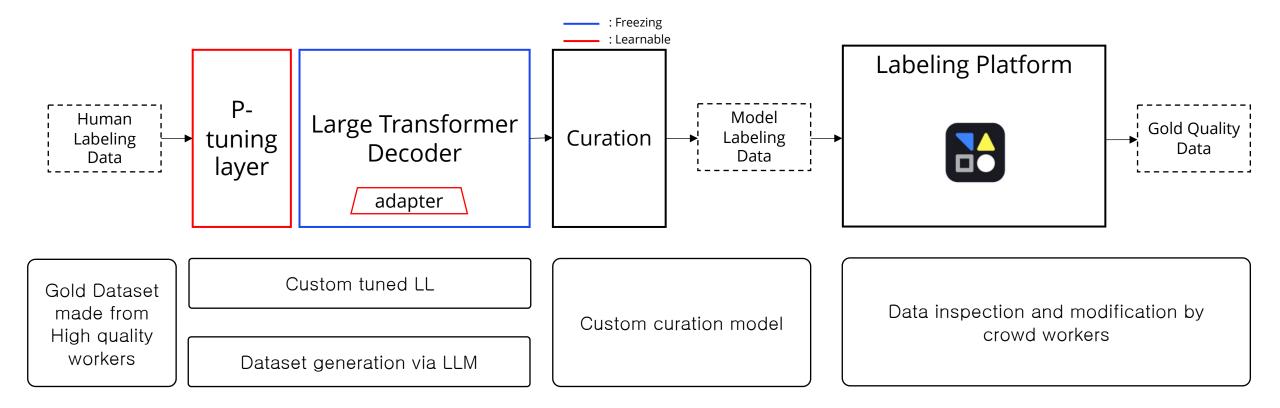environment by providing on-premise
('23.03 Release)

Feature Vector 추론 후 Source Data 폐기

Source Data

Cloud (S3, Azure..)

Information    Feature Vector

**Datumo FST DB**

**Datumo FST Product**

As-is

To-be        or    custom

Feature Embedding Model

# Filtering of similar text submissions.



Different abuse criteria can be set depending on the data by using similarity thresholds

- Similarity scores between pairs of text submitted by users are checked, with a range of 0 to 5 indicating how similar they are.

- The number of similar text pairs can be visualized for each user.

**DATUMO**

# Free Data Labeling Service for academic purposes



InstaOrder
(w. POSTECH)



KLUE : Korean Language Understanding Evaluation
(w. Upstage, NAVER, NYU, KAIST)



KOLD: Korean Offensive Language Dataset
(w. KAIST)



Analyzing Norm Violations in Real-Time Live-Streaming Chat
(w. Softly.AI, USC)



Split GCN

# Free offerings

• Try Datumo Scope Now! For FREE

• Free Data Labeling Service for Academic purposes

For more benefits, please contact contact@datumo.com

We aim to innovate the AI training data market through technology

contact@datumo.com

DATUMO.COM