# LaRD

# Automating 01 Weak Supervision to Find Missing Labels for Big Data

### Presented by: Mona Nashaat

# ① Motivation: Why Labels?

## Derive Value from Business Data

- Essential to build supervised machine learning models.
- The quality and the size of training data limits the performance of predictive systems.
- Labeled training datasets do not exist.

> About 70% of complex analytical tasks today are related to data preparation. There have to be people who are preparing and labeling data for machines to understand. Here's a situation in which human labor automation driven by ML creates new job opportunities.

**Guru Banavar, IBM data scientist**

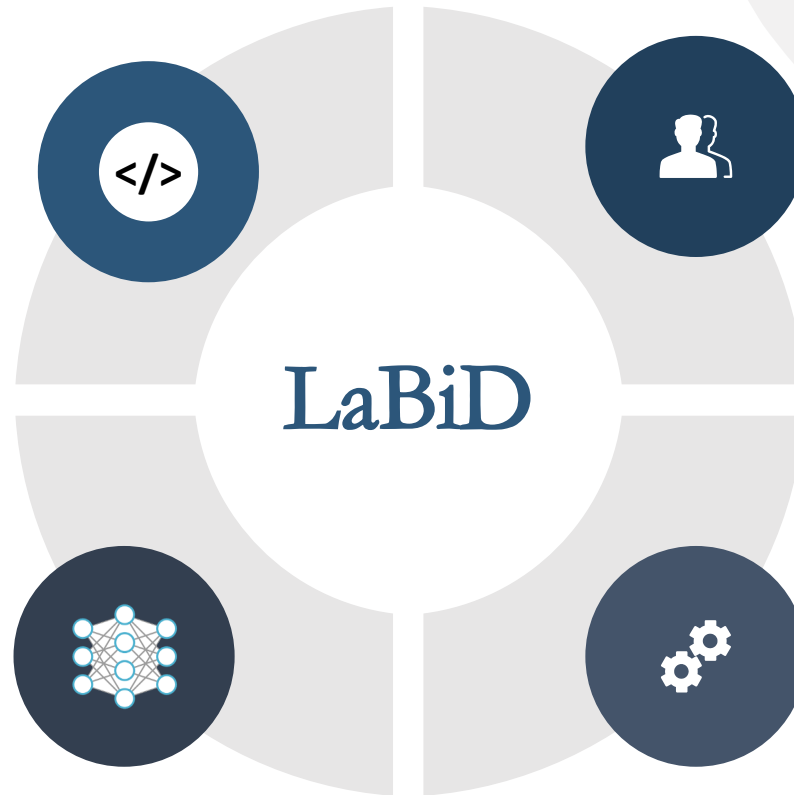## Obtaining Accurate Labels is Expensive

**Data Programming**

Learns a model of the training set that includes labeling functions.

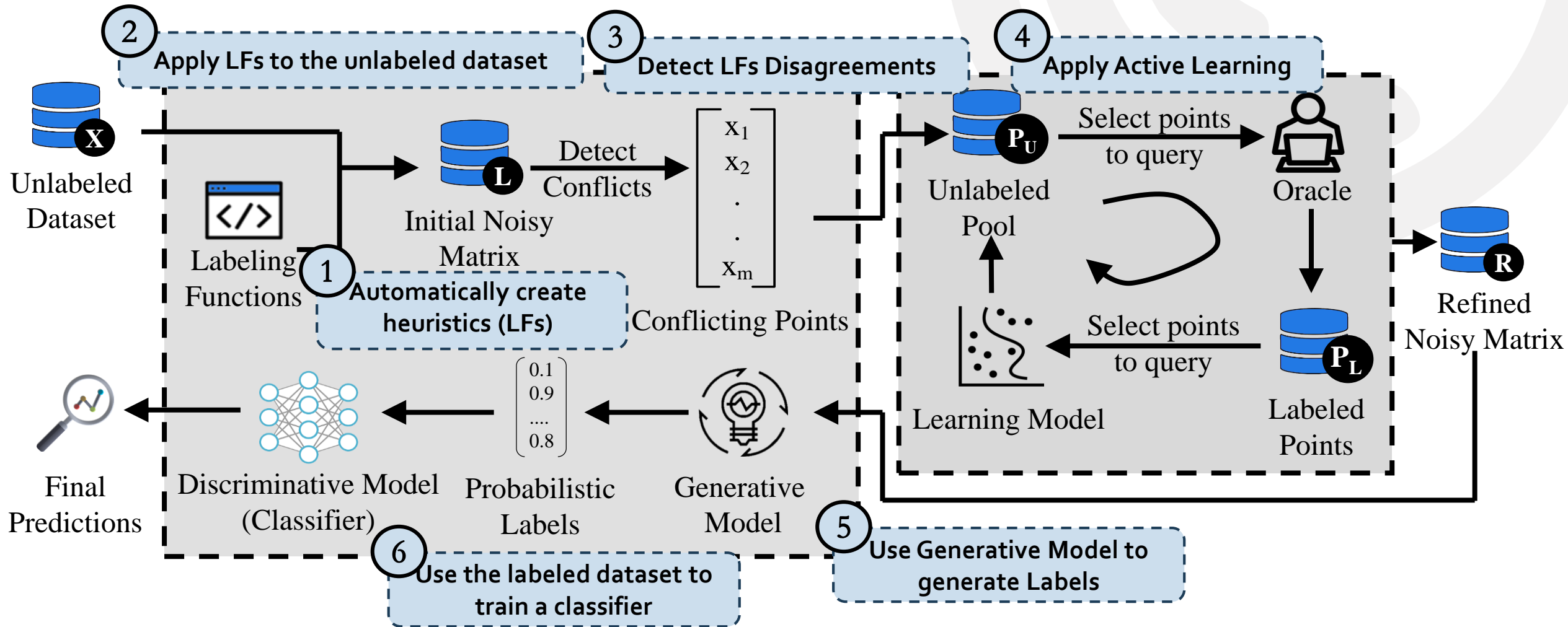**Meta Active Learning**

Treats active learning algorithm design as a meta-learning problem and learn the best criterion from data

LaBiD

Gets a lower-quality labels more efficiently and/or at a higher abstraction level

**Weak Supervision**

Automating the process of generating heuristics that assign training labels to unlabeled data

**Automating Weak Supervision**

② Apply LFs to the unlabeled dataset
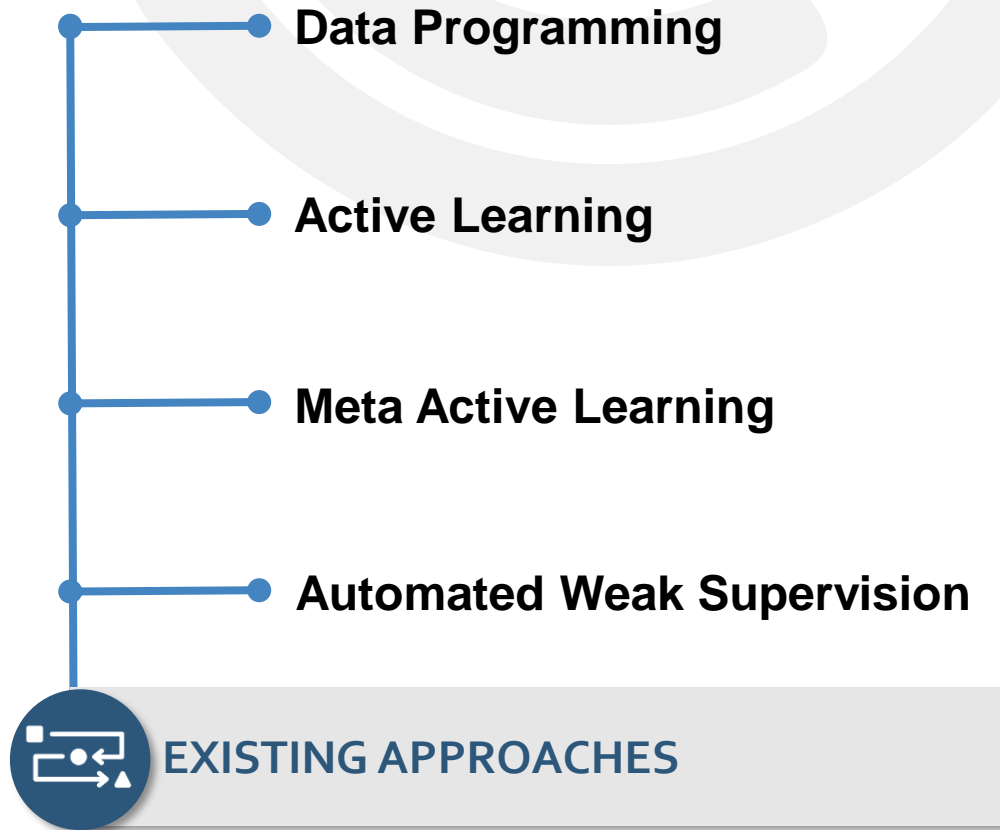
③ Detect LFs Disagreements

④ Apply Active Learning

X — Unlabeled Dataset

Labeling Functions

① Automatically create heuristics (LFs)

L — Initial Noisy Matrix

Detect Conflicts

$$\begin{bmatrix} x_1 \\ x_2 \\ . \\ . \\ . \\ x_m \end{bmatrix}$$

Conflicting Points

$P_U$ — Unlabeled Pool

Select points to query

Oracle

R — Refined Noisy Matrix

Select points to query

Learning Model

$P_L$ — Labeled Points

$$\begin{pmatrix} 0.1 \\ 0.9 \\ .... \\ 0.8 \end{pmatrix}$$

Final Predictions

Discriminative Model (Classifier)

Probabilistic Labels

Generative Model

⑥ Use the labeled dataset to train a classifier

⑤ Use Generative Model to generate Labels

4

# ④ Experimental Results

| Dataset | # of records | # of attributes |
| --- | --- | --- |
| Higgs | 11,000,000 | 28 |
| Renewal Sales | 1,354,704 | 11 |
| Rain Prediction | 142,000 | 24 |
| Travel Insurance | 63,300 | 11 |
| Bank | 45,211 | 17 |
| News | 39,797 | 61 |
| Credit Card | 30,000 | 24 |
| Tenancy Detection | 20,560 | 7 |
| Magic | 19,020 | 12 |

- **Data Programming**
- **Active Learning**
- **Meta Active Learning**
- **Automated Weak Supervision**

**EXISTING APPROACHES**

# ⑤ Application to Database

## Analytics challenges

- Never assume the data is clean.
- Automatically create heuristics.
- Apply the LaBiD flow and compare the results with ground truth.
- Double check with the user to detect outliers and missing values.

Bad data is bad for business. Poor quality data is costing businesses at least 30% of revenues.

**Reported by Ovum Research**

**Next steps…**

7

# THANK YOU

✉ nashaata@ualberta.ca