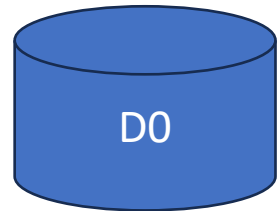
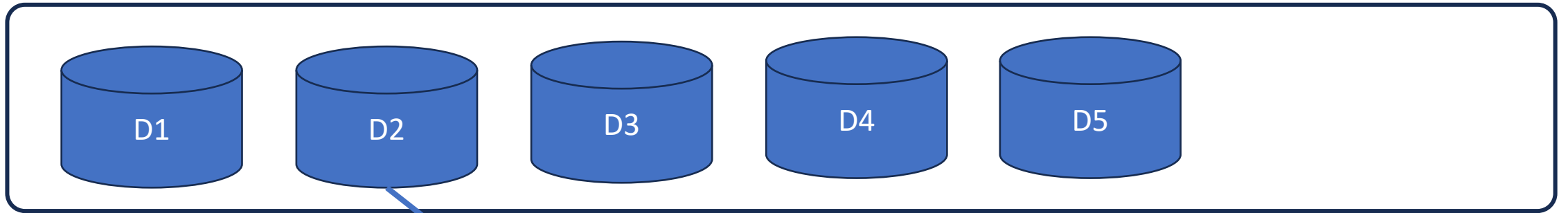


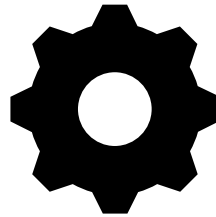
Finding Testdataset for Transfer Learning

Repository of labeled datasets



Unlabeled dataset

1. train/optimize

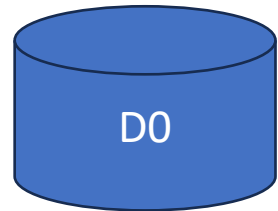
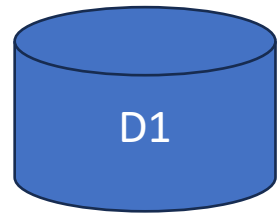


2. apply

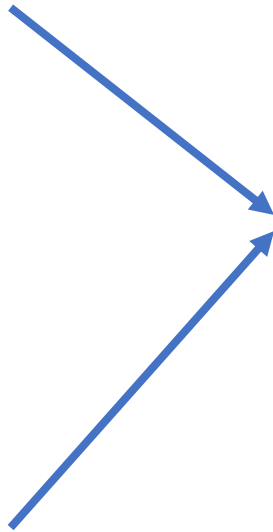
Can we use any of the labeled dataset to configure a data matching approach for D0?

First Approach

Labeled dataset



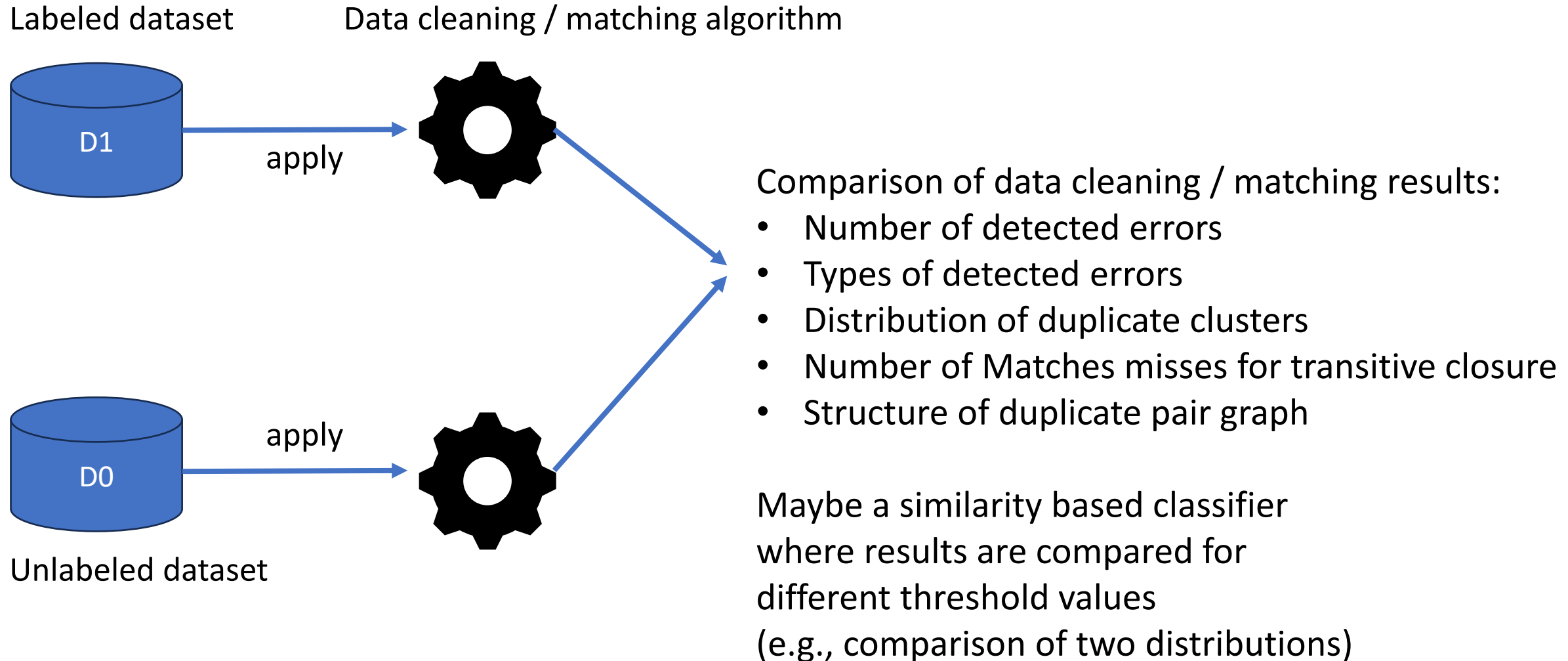
Unlabeled dataset



Comparison of dataset characteristics, such as:

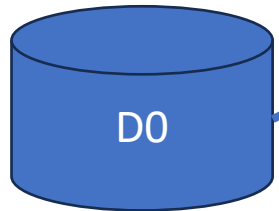
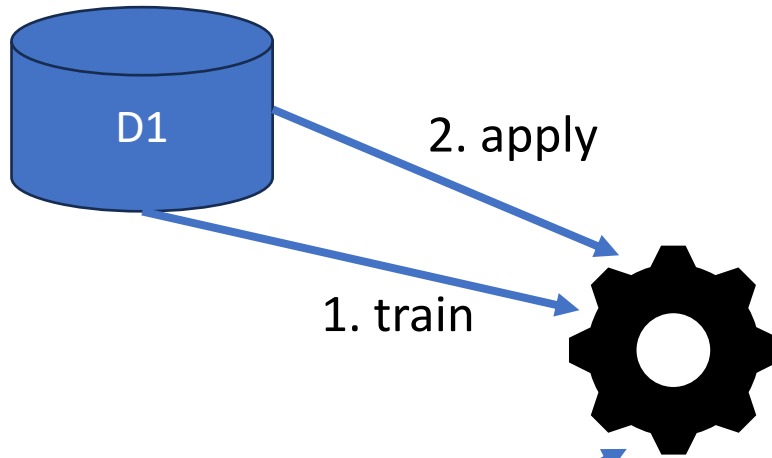
- Number attributes
- Number records
- Data types (numerical, categorical, text)
- Sparsity
- Integrity constraints

Second Approach



Third Approach

Labeled dataset



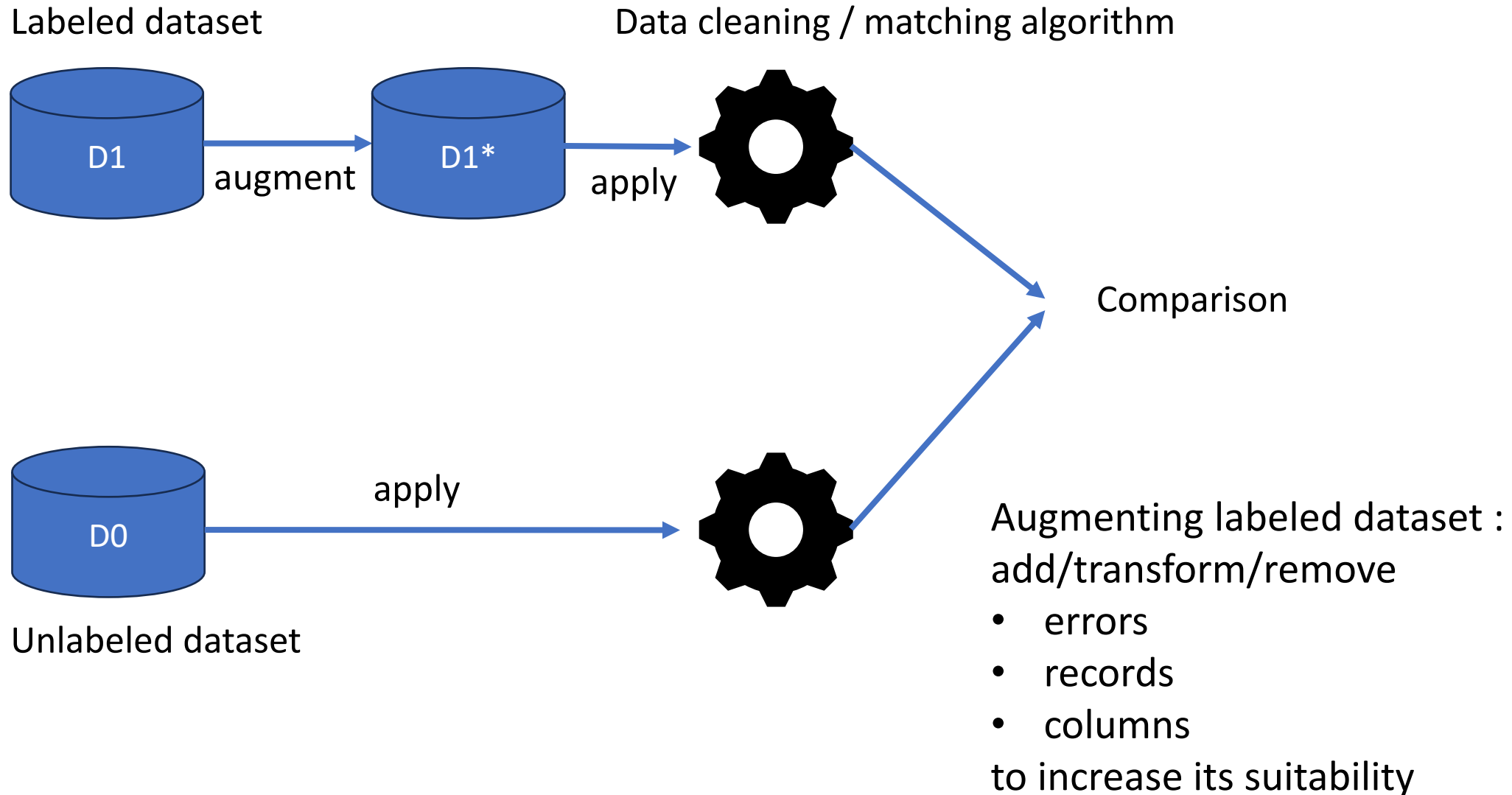
Unlabeled dataset

Comparison of matching results:

- Distribution of duplicate clusters
- Number of Matches misses for transitive closure
- Structure of duplicate pair graph

Maybe with different train/test splits of D1

Augmentation



First Steps

1. Building a suitable repository of 10-15 labeled datasets
2. Building a suitable repository of 3-5 data cleaning/matching algorithms
3. Evaluating suitability among these datasets
(would be interesting to see how much the quality differs
and depends on dataset used)