

# Data Curation for ML: Toward a Principled Approach

**Laure Berti-Equille**

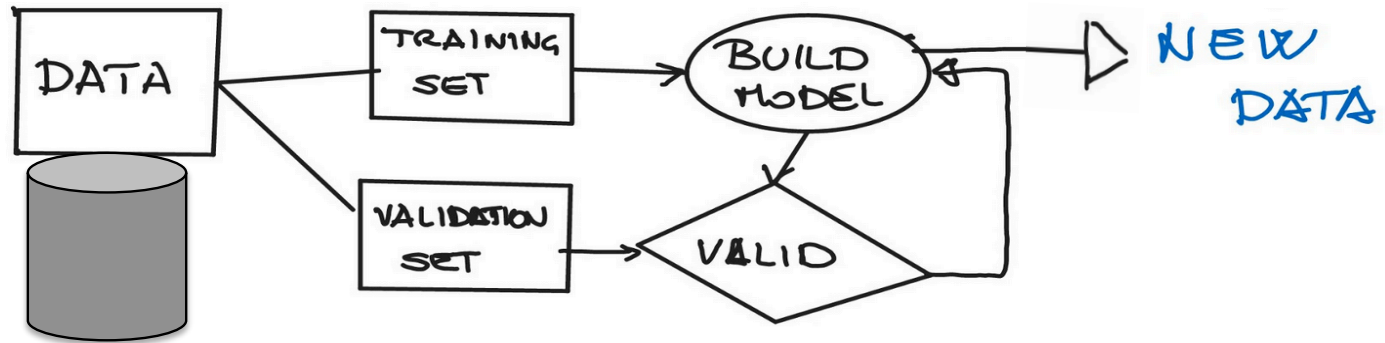
[laure.berti@ird.fr](mailto:laure.berti@ird.fr)

Espace-Dev, IRD, Univ Montpellier, Univ Guyane, Univ La Réunion, Univ Antilles,  
Univ Nouvelle Calédonie, Montpellier France

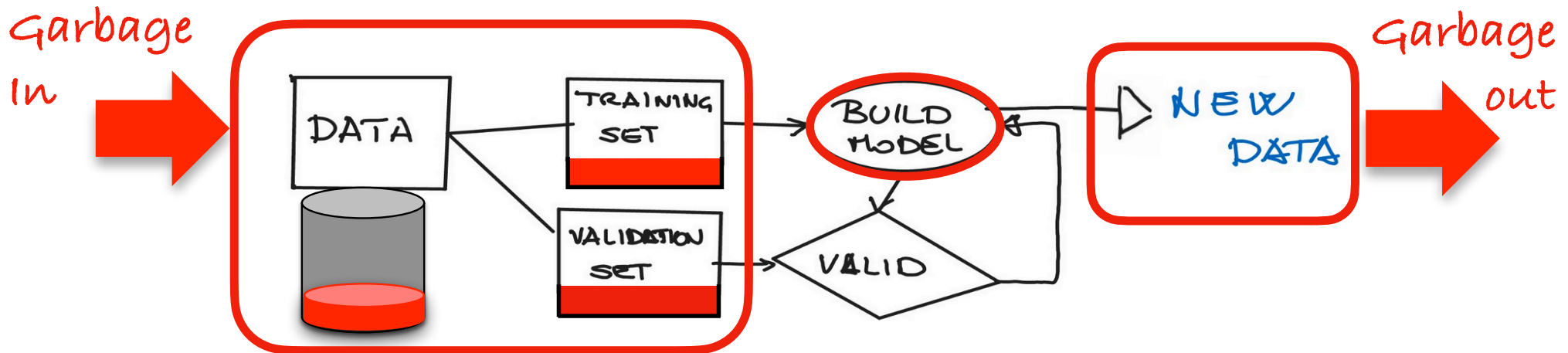
TALOS AI4SHS Mobility Grant Presentation at UCRC-KEME  
May 15, 2024



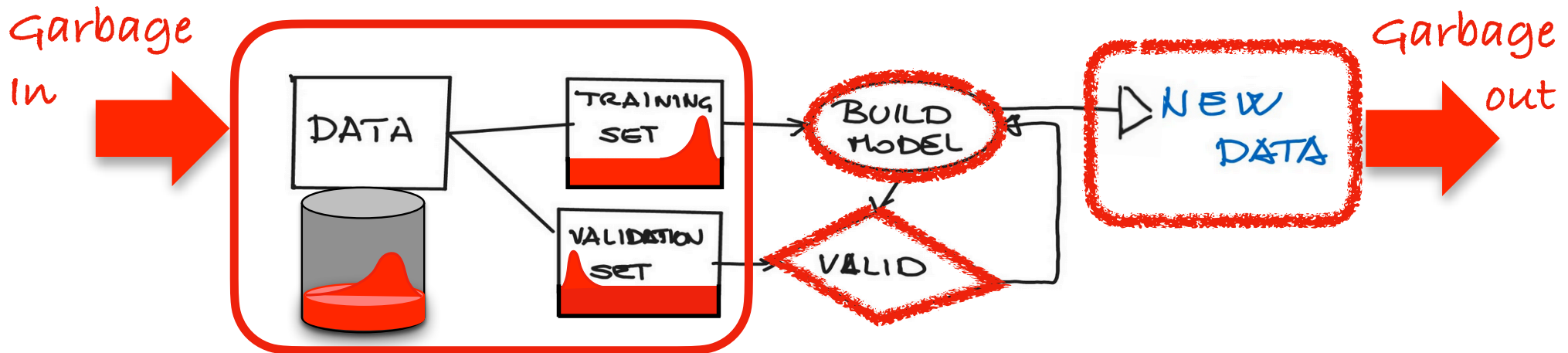
# Learning from dirty data is risky



# Learning from dirty data is risky



# Learning from dirty data is risky



Glitch types and distributions can be very different in the datasets used for training, testing, and validation and they affect accuracy of ML models in different ways.

# Two complementary approaches

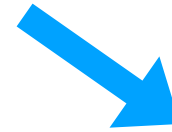
---



## **INTERVENE**

### **How to efficiently fix the data:**

- ◆ Detect the anomalies
- ◆ Correct them with minimal cost (domain expert intervention, time, external master data, etc.)
- ◆ Select the repair/preparation strategies that will maximize the ML result quality



## **MITIGATE**

### **How to reduce the impact of dirty data:**

- ◆ Robustify the ML algorithms and apply ML ensembling strategies
- ◆ Use AutoML to find optimal parameter setting
- ◆ Select portions of the data and/or augment the data

# Two complementary approaches

---



## **INTERVENE**

### **How to efficiently fix the data:**

- ◆ Detect the anomalies
- ◆ Correct them with minimal cost (domain expert intervention, time, external master data, etc.)
- ◆ Select the repair/preparation strategies that will maximize the ML result quality



## **MITIGATE**

### **How to reduce the impact of dirty data:**

- ◆ Robustify the ML algorithms and apply ML ensembling strategies
- ◆ Use AutoML to find optimal parameter setting
- ◆ Select portions of the data and/or augment the data

# Outline

---

1. **Detection of data quality problems**  
Profiling data quality
2. **Data cleaning**  
Leveraging the patterns of glitches
3. **Data preparation strategies for ML**  
Learning to clean and prepare the data

# Outline

---

## 1. **Detection of data quality problems**

Profiling data quality

## 2. **Data cleaning**

Leveraging the patterns of glitches

## 3. **Data preparation strategies**

Learning to clean and prepare the data



# Data Quality Problems

## DATA TYPES

0101010101

ACACGTGT

John Doe

High  
Medium  
Low

## RELATIONSHIPS

- Structural (record)
- Sequential
- Graph-based
- Temporal
- Spatial
- Spatio-Temporal

## DATA QUALITY PROBLEMS

### TYPE

Missing Data  
Anomalous Data  
Duplicate Data  
Inconsistent Data  
Obsolete Data  
Incorrect data

### CARDINALITY

Single-Point  
Collection

### DETECTION MODE

- Model-based
- Data distribution-based
- Constraint-based
- Pattern-based

# Data Quality Problems: Example I

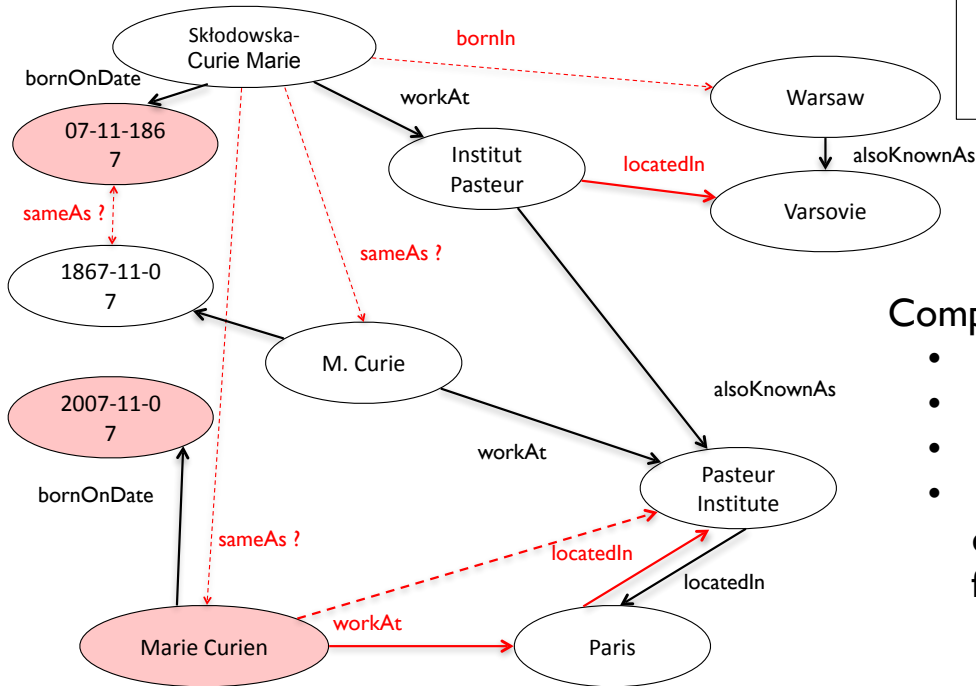
Relational data quality problems

*Nobel Laureates in Chemistry*

Name	Institution	Institution_City	DoB
Sklodowska-Curie Marie	Institut Pasteur	Varsovie	07-11-1867
M. Curie	Pasteur Institute	Paris	1867-11-07
Melvin Calvin	UC Berkeley	Berkeley	1911-04-08
Marie Curien	Paris	Pasteur Institute	2007-11-07
Avram Hershko	NULL	Haifa	NULL
Ronald Hoffman		US	00000000

# Data Quality Problems: Example 2

Knowledge Graph data problems  
Nobel Laureates in Chemistry: Excerpt



Name	Institution	Institution_City	DoB
Skłodowska-Curie Marie	Institut Pasteur	Varsovie	07-11-1867
M. Curie	Pasteur Institute	Paris	1867-11-07
Melvin Calvin	UC Berkeley	Berkeley	1911-04-08
Marie Curien	Paris	Pasteur Institute	2007-11-07
Avram Hershko	NULL	Haifa	NULL
Ronald Hoffman	US	US	00000000

Annotations on the table:

- Representation:** Points to the table header.
- Misfielded Value:** Points to the 'DoB' column.
- Duplicates:** Points to the first two rows (Skłodowska-Curie Marie and M. Curie).
- Typos:** Points to 'Haifa' and 'US' in the 'Institution\_City' column.
- Inconsistencies:** Points to 'Paris' and 'Pasteur Institute' in the 'Institution\_City' column for Marie Curien.
- Incorrect Values:** Points to 'US' in the 'Institution' column for Ronald Hoffman.
- Incorrect Value:** Points to 'US' in the 'Institution\_City' column for Ronald Hoffman.
- Missing Values:** Points to 'NULL' in the 'DoB' column for Avram Hershko.

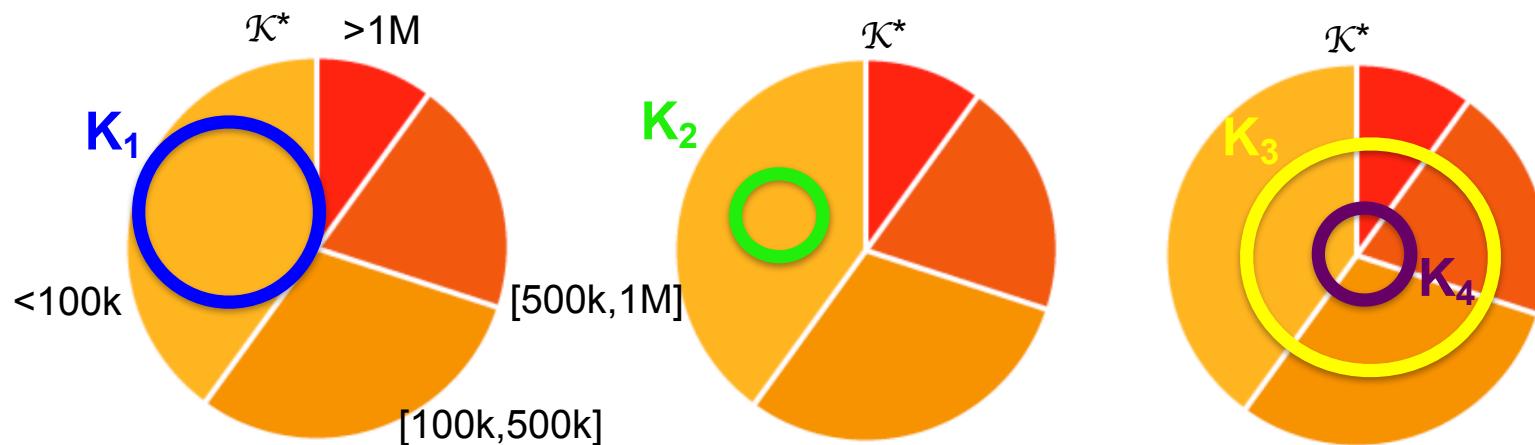
Complex combination of:

- Missing links and entities
- Spurious links : existence, type, direction
- Erroneous entity name
- Errors in literal values with various degrees of severity: formatting, up-to-dateness, veracity issues

# Data Quality Problems: Example 3

## Completeness

Suppose you have the accurate and complete knowledge of the world-wide populations per city grouped into 4 categories: e.g. (<100k, [100k,500k], [500k,1M], >1M) and 4 KBs.

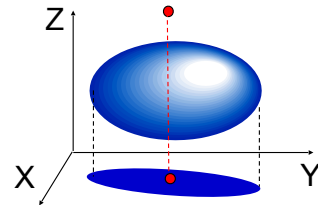


$K_1$  is more complete than  $K_2$  but both are somehow biased toward one category

$K_1$  and  $K_2$  are not as representative as  $K_3$  or  $K_4$

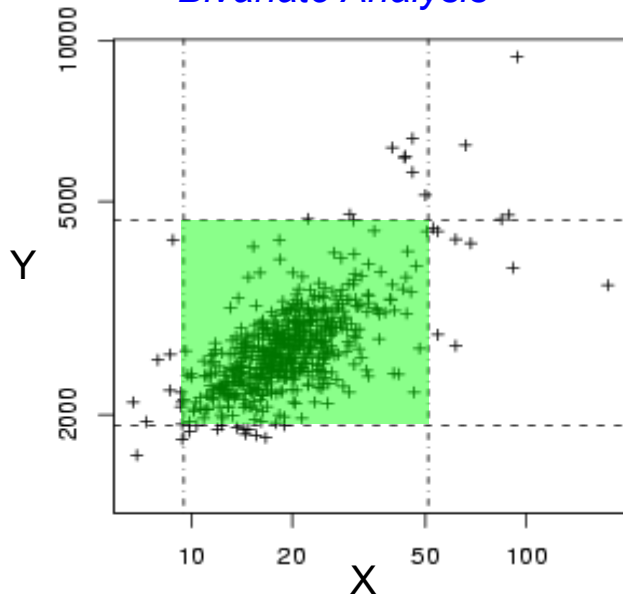
- Soulet, Giacometti, Markhoff, Suchanek: Representativeness of Knowledge Bases with the Generalized Benford's Law. *International Semantic Web Conference (I)* 2018: 374-390
- Wagner, Garcia, Jadidi, Strohmaier: It's a man's Wikipedia? Assessing gender inequality in an online encyclopedia. *ICWSM*. pp. 454-463 (2015)
- Callahan, Herring: Cultural bias in Wikipedia content on famous persons. *J. of the Association for Information Science and Technology*, 62(10), 1899-1915 (2011)
- Pitoura, Tsabaras, Flouris, Fundulaki, Pabadakos, Abiteboul, Weikum. On Measuring Bias in Online Information. *SIGMOD Record*. Vol.46 No.4. December 2017

# Example 4. Numerical Outliers



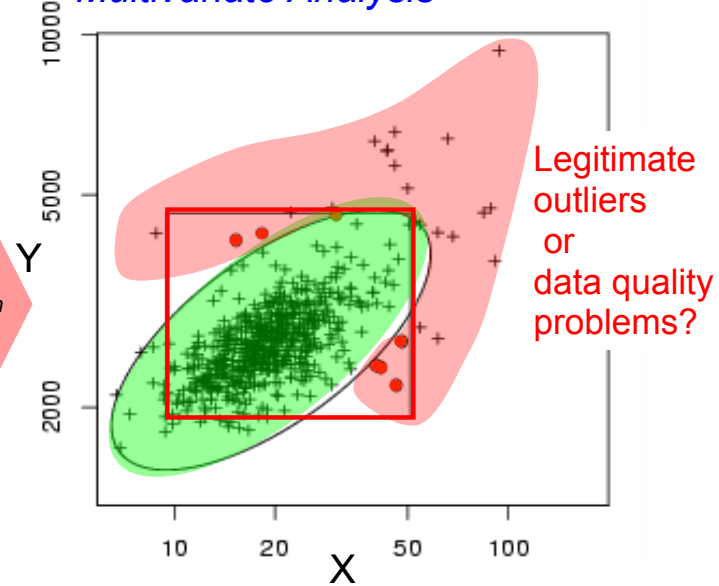
(Classical Setting)

*Bivariate Analysis*



Rejection area: Data space excluding the area defined between 2% and 98% quantiles for X and Y

*Multivariate Analysis*




Rejection area based on:  
 $\text{Mahalanobis\_dist}(\text{cov}(X,Y)) > \chi^2(.98,2)$

# Example 5: Up-to-dateness

## Asynchronous Real World and KG evolution

<https://www.dbpedia.org/resources/ontology/>

Table 1. DBpedia - Classes and Properties



Version	OWL Class				RDF Property				Object Prop.			Datatype Prop.		
	#	$\Delta$	(-)	(+)	#	$\Delta$	(-)	(+)	#	(-)	(+)	#	(-)	(+)
3.2/3	174				720				384			336		
3.4	204	30	-2	32	2168	1448	-271	1719	1144	-139	899	1024	-132	820
3.5	255	51	-6	57	1274	-894	-1198	304	601	-673	130	673	-525	174
3.6	272	17	0	17	1335	61	-37	98	629	-26	54	706	-11	44
3.7	319	47	-1	48	1643	308	-17	325	750	-6	127	893	-11	198
3.8	359	40	-1	41	1775	132	-3	135	800	-1	51	975	-2	84
3.9	529	170	-1	171	2333	558	-8	566	927	-6	133	1406	-2	433
2014	683	154	-5	159	2795	462	-46	508	1079	-9	161	1716	-37	347
2015-04	735	52	-5	57	2819	24	-103	127	1098	-23	42	1721	-80	85
2015-10	739	4	-5	9	2833	14	-9	23	1099	-3	4	1734	-6	19
2016-04	754	15	0	15	2849	16	-2	18	1103	-1	5	1746	-1	13

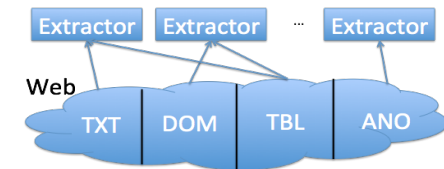
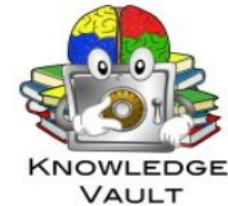
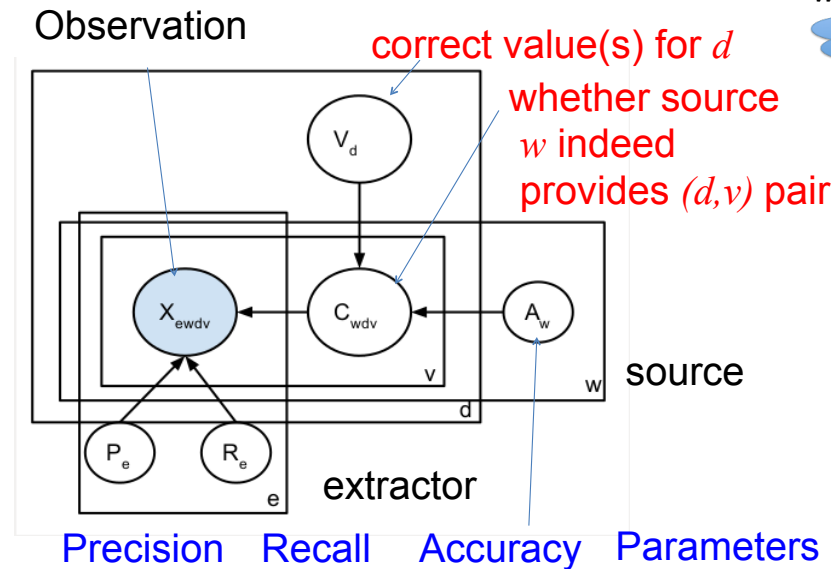
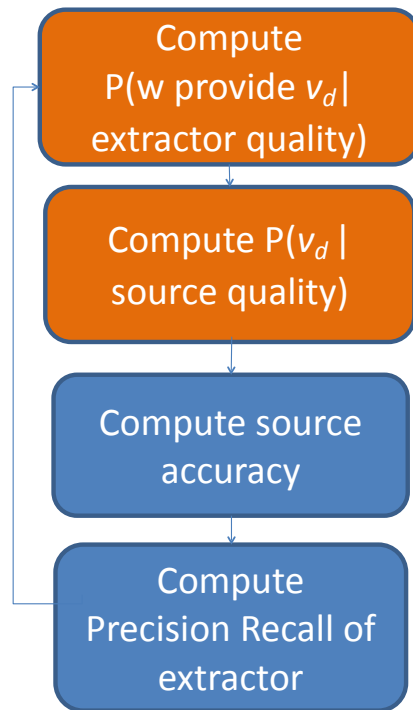
Today's DBpedia Ontology: 768 classes described by 3000 properties 4,233,000 instances.

Mihindukulasooriya, Poveda-Villalon, Garcia-Castro, Gomez-Perez. Collaborative Ontology Evolution and Data Quality -An Empirical Analysis, in OWL: Experiences and Directions – Reasoner Evaluation, Springer International Publishing, Cham, 2017, pp. 95–114. [https://www.w3.org/community/owled/files/2016/11/OWLED-ORE-2016\\_paper\\_9.pdf](https://www.w3.org/community/owled/files/2016/11/OWLED-ORE-2016_paper_9.pdf)

# Example 6. Veracity and Trustworthiness

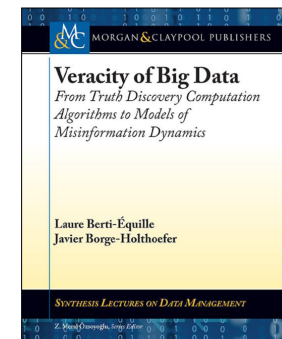
ML-based approach for knowledge-based trust:

- Multi-Layer Model based on EM and Bayesian inference
- Distinguish extractor errors from source errors



#Triples	3.0B (0.3B w. pr>=0.7)
#URLs	2.5B (28M Websites)
#Extractors	16

As of 2014



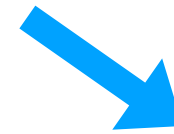
# Existing approaches for detecting/fixing DQ problems

---



## **Declarative**

- Data debugging
- Checking data assertions
- Transform



## **ML-based**

Learn from clean data and replace



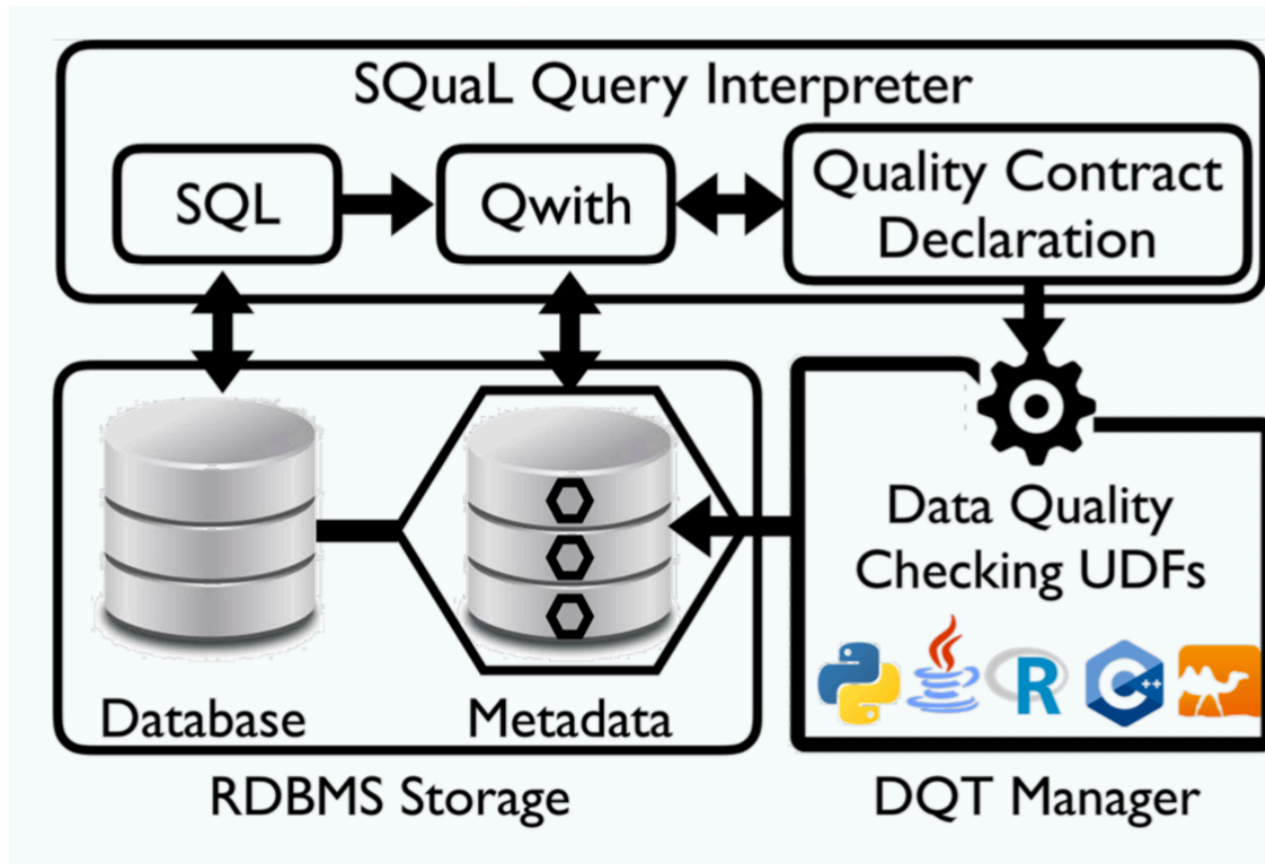
# Declarative Approaches

---

## Checking data assertions and transform

- ◆ **Deequ** [Schelter et al., VLDB 2018] requires cloud infrastructure and manual integration into training and serving systems; dependent on Apache Spark
- ◆ **TensorFlow Data Validation** (TFDV) [Caveness et al., SIGMOD 2020] integrated with Google TFX difficult to use outside of these platforms
- ◆ Lightweight Python-based approaches like **great\_expectations** (<https://greatexpectations.io>) or **hooqu** (<https://github.com/mfcabrera/hooqu>) not integrated with the ML development process

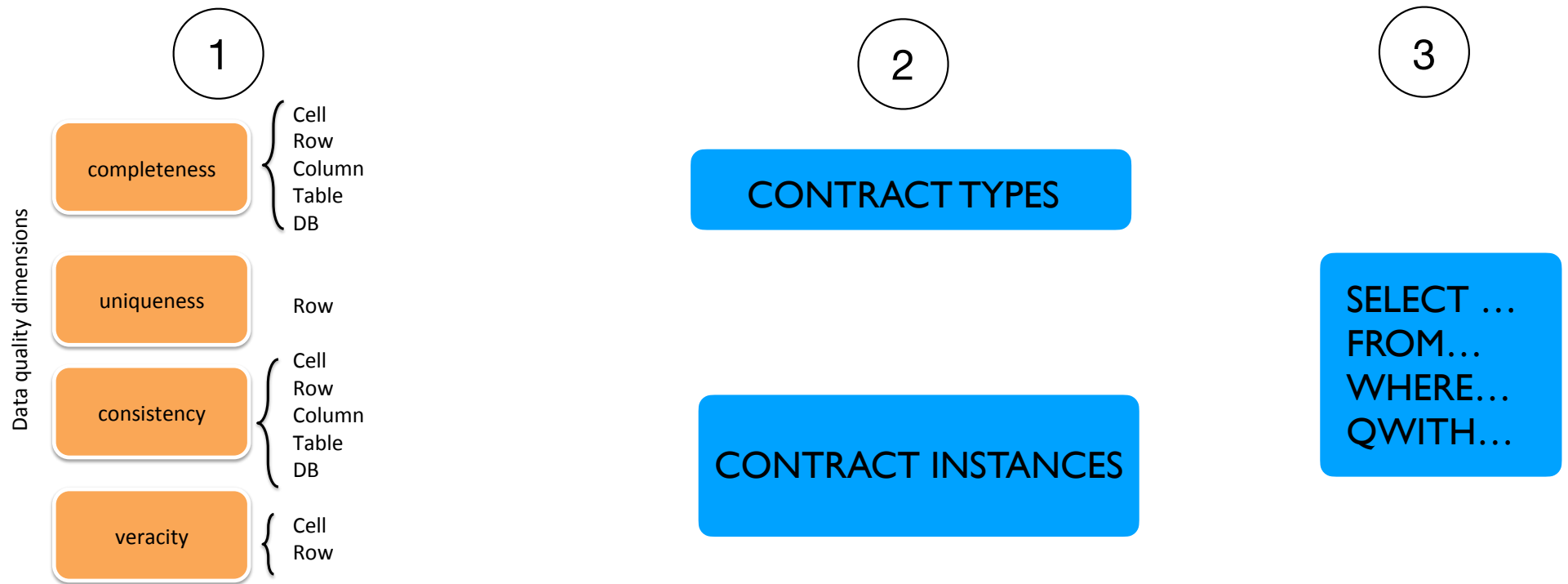
# Declarative data profiling with MeSQual



<https://github.com/ucomignani/MeSQual>

# MeSQual Key Concepts

*Flexible declarative data quality profiling with UDFs*



**Procedural approach with UDFs**

**Declarative approach**

**Extended query**

# MeSQual Examples

## DECLARATION

```
CREATE CONTRACTTYPE StatTests (
  autocorrelation BY FUNCTION 'durbinWatsonTest.py' LANGUAGE PYTHON,
  multicollinearity BY FUNCTION 'varInflationFactor.py' LANGUAGE PYTHON,
  heteroscedasticity BY FUNCTION 'BreuschPaganTest.py' LANGUAGE PYTHON,
  KErrorNormality BY FUNCTION 'KolmogorovSmirnov.py' LANGUAGE PYTHON,
  SWerrorNormality BY FUNCTION 'ShapiroWilkTest.py' LANGUAGE PYTHON);
```

```
CREATE CONTRACT RegressionAssumptions (
  StatTests.autocorrelation > 0
  AND StatTests.autocorrelation < 4
  AND StatTests.multicollinearity <= 4
  AND StatTests.heteroscedasticity < 0.05
  AND StatTests.SWerrorNormality < 0.05);
```

```
CREATE CONTRACTTYPE CheckQDB (
  completeness BY FUNCTION 'completeness.py' LANGUAGE PYTHON,
  uniqueness BY FUNCTION 'uniqueness.py' LANGUAGE PYTHON,
  consistency BY FUNCTION 'consistency.py' LANGUAGE PYTHON,
  outlyingness BY FUNCTION 'outlyingness.py' LANGUAGE PYTHON);
```

```
CREATE CONTRACT CheckBeforeAnalysis (
  RegressionAssumptions
  AND CheckQDB.consistency > 0.9
  AND CheckQDB.outlyingness < 0.2);
```

## MANIPULATION

AoT	{ SELECT * FROM ChicagoDataset } QWITH CheckQDB.completeness> 0.95;
	{ SELECT * FROM ChicagoDataset } QWITH CheckBeforeAnalysis AND RegressionAssumptions;
MIMIC-III	{ SELECT timestamp, node_id,value_raw,valuehrf FROM ChicagoDataset WHERE ChicagoDataset.sensor = 'o3' } QWITH CheckBeforeAnalysis AND CheckQDB.completeness> 0.95;
	{ SELECT * FROM Admissions } QWITH CheckQDB.completeness> 0.95;
	{ SELECT * FROM Admissions WHERE Admissions.insurance = 'Private' } QWITH CheckBeforeAnalysis AND CheckQDB.completeness> 0.95;
	{ SELECT gender, dob, admittime FROM Admissions INNER JOIN (SELECT * FROM Patients WHERE dob < '2090-12-12 00:00:00' QWITH CheckQDB.completeness> 0.95) as Pat ON Admissions.subject_id=Pat.subject_id; } QWITH CheckQDB.completeness> 0.95;

# MeSQual GUI

**A** SQual Query

```
{
  SELECT timestamp, node_id, value_raw, value_hrf
  FROM ChicagoDataset
  WHERE ChicagoDataset.sensor = 'o3'
}
QWITH CheckBeforeAnalysis AND CheckQDB.completeness > 0.95
```

Run

**C** Tables

database	table
Test	CONTRACTTYPE
Test	CONTRACT
Test	ChicagoDataset

Contracts

contractName	constraintOperator	dimensionName	comparedValue
CheckBeforeAnalysis	CONTRACT	RegressionAssumptions	-1.00
CheckBeforeAnalysis	LESSER	outlyingness	0.20
CheckBeforeAnalysis	GREATER	consistency	0.90

Contract Types

contractTypeName	dimensionName	language	functionPath
CheckQDB	outlyingness	PYTHON	outlyingness.py
CheckQDB	consistency	PYTHON	consistency.py
CheckQDB	uniqueness	PYTHON	uniqueness.py

Queries

queryId	query
c21418d8-5e3c-4814-9556-5e0a7196b502	{ SELECT timestamp, node_id, value_raw, value_hrf FROM ChicagoDataset WHERE ChicagoDataset.sensor = 'o3' } QWITH CheckBeforeAnalysis AND CheckQDB.completeness > 0.95

**B** Query Results

timestamp	node_id	value_raw	value_hrf
2019/11/18 12:55:07	001e061146cb	-629.00	0.00
2019/11/18 12:55:06	001e06117b41	970.00	0.00
2019/11/18 12:55:02	001e0510ee43	1.55 K	0.00
2019/11/18 12:54:59	001e061183f3	1.83 K	0.00
2019/11/18 12:54:54	001e061144be	1.67 K	0.00
2019/11/18 12:54:52	001e0610f6db	436.00	0.00

**D** Data Quality Checks

Check	Value
completeness db	0.95
consistency db	0.90
completeness value_raw	0.80
completeness value_hrf	0.90
consistency value_raw	0.90
consistency value_hrf	0.90
heteroscedasticity value_hrf	0.10

**E** Monitoring

Monitoring: Last Checkpoint

Check	Value
multicollinearity.att.value_raw	2.300

**F** Query & Check Logs

queryId	queryTimestamp	contractType	dimensionName	comparedValue	operator	elementGranularity	elementId	udfResult	violation
c21418d8-5e3c-4814-9556-5e0a7196b502	2019-11-26 15:00:00.000	CheckQDB	consistency	0.90	>	att	value_hrf	0.90	1.00
c21418d8-5e3c-4814-9556-5e0a7196b502	2019-11-26 15:00:00.000	StatTests	heteroscedasticity	0.05	<	att	value_hrf	0.10	1.00
c21418d8-5e3c-4814-9556-5e0a7196b502	2019-11-26 15:00:00.000	CheckQDB	completeness	0.95	>	att	value_raw	0.80	1.00

<https://github.com/ucomignani/MeSQual>

# ML-based Approaches

## Learn from clean data and replace/repair

- Pattern enforcement
  - Syntactic patterns (date formatting)
  - Semantic patterns (name/address)
- Value update to satisfy a set of rules, constraints, FDs, CFDs, Denial Constraints (DCs), Matching Dependencies (MDs) with minimal number of changes.
- Value replacement
- Entity resolution

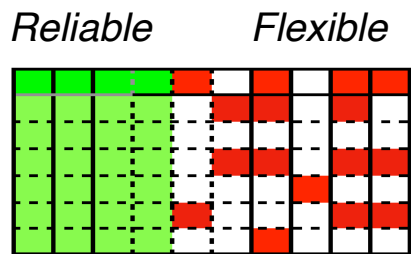
### EXAMPLES

- ◆ SCARE: Scalable Automatic Repair
- ◆ On-demand ETL [Yang et al., VLDB'15]
- ◆ ActiveClean [Krishnan et al., VLDB'16]
- ◆ HoloClean [Rekatsinas et al., VLDB 2017]
- ◆ Deep learning for Entity Resolution
- ◆ Transformers for data prep

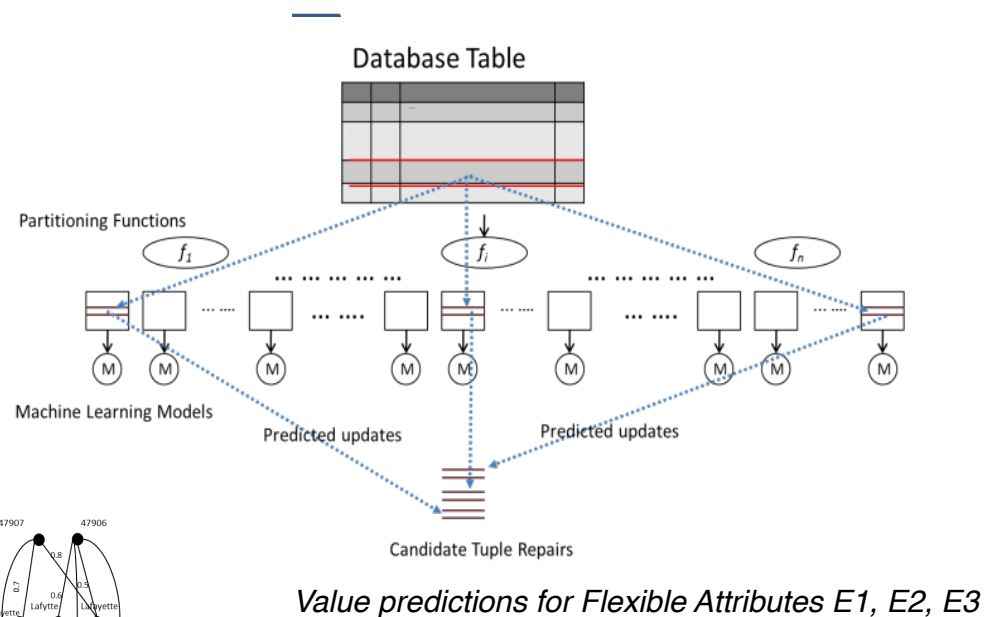
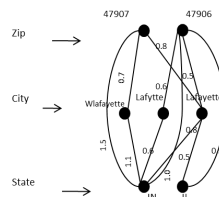
# SCARE: SCalable Automatic Repair

[Yakout, Berti-Equille, Elmagarmid, SIGMOD 2013]

Goal: Find the repair that would maximize the sum of the probabilities of the values co-occurrence (i.e., association strength between predicted and reliable values) under a certain update cost budget.



1. Modeling Dependency and Predicting Updates
2. Data Partitioning
3. Tuple Repair Selection

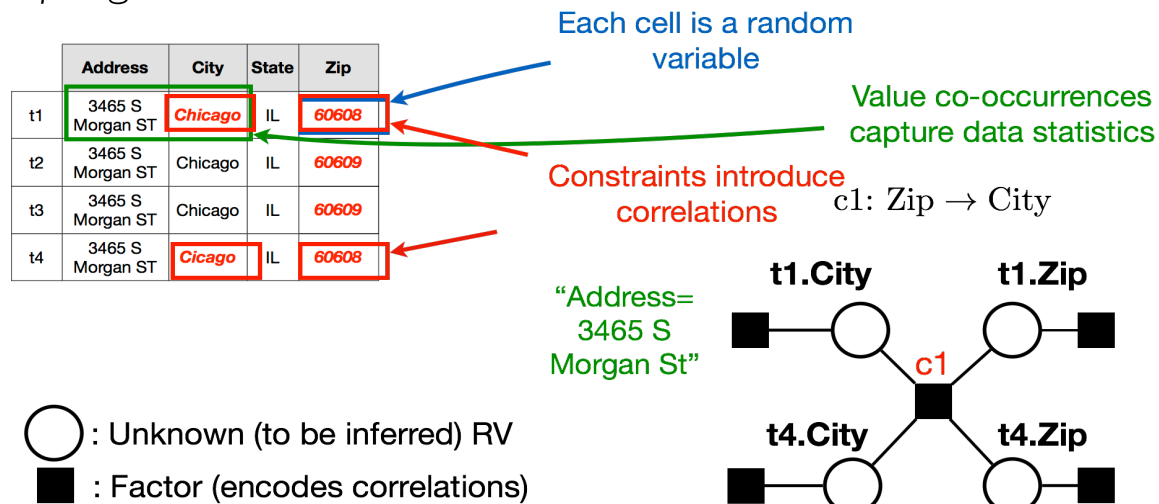


# HoloClean

[Rekatsinas et al., VLDB 2017]

<https://github.com/HoloClean/HoloClean>

HoloClean generates a factor graph capturing co-occurrences, correlations based on a set of constraints and external evidences. It uses SGD to learn parameters and infer the marginal distribution of unknown variables with Gibbs sampling.



Denial constraints:

$$\forall t_1, t_2 \in D : \neg(t_1[\text{Zip}] = t_2[\text{Zip}] \wedge t_1[\text{City}] \neq t_2[\text{City}])$$

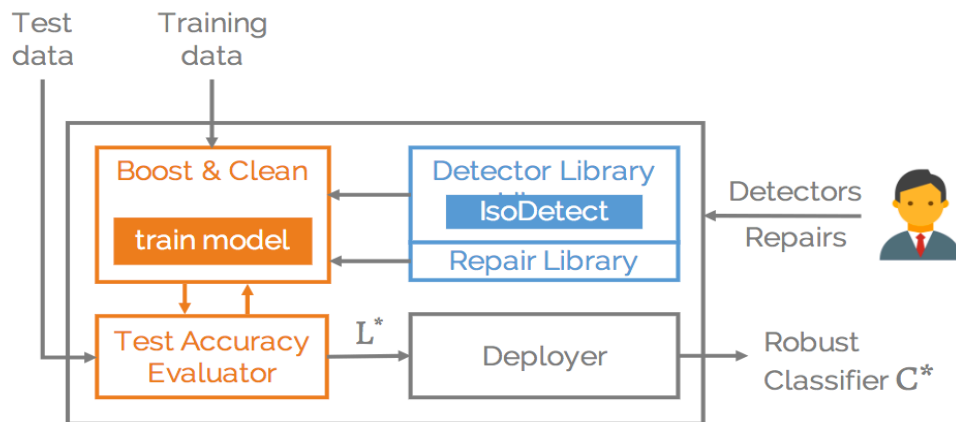
$$\forall t_1, t_2 \in D : \neg(t_1[\text{Zip}] = t_2[\text{Zip}] \wedge t_1[\text{State}] \neq t_2[\text{State}])$$



# BoostClean

[Krishnan et al., 2017]

*BoostClean selects an ensemble of methods (statistical and logic rules) for error detection and for repair combinations using statistical boosting.*



## Algorithm : Boost-and-Clean Algorithm

**Data:**  $(X, Y)$

- 1 Initialize  $W_i^{(1)} = \frac{1}{N}$
- 2  $\mathcal{L}$  generates a set of classifiers  $\mathcal{C} \{C^{(0)}, C^{(1)}, \dots, C^{(k)}\}$  where  $C^{(0)}$  is the base classifier and  $C^{(1)}, \dots, C^{(k)}$  are derived from the cleaning operations.
- 3 **for**  $t \in [1, T]$  **do**
- 4      $C_t = \text{Find } C_t \in \mathcal{C}$  that maximizes the weighted accuracy on the test set.  $\epsilon_t = \text{Calculate weighted classification error on the test set } \alpha_t = \ln\left(\frac{1-\epsilon_t}{\epsilon_t}\right)$   
    $W_i^{(t+1)} \propto W_i^{(t)} e^{-\alpha_t y_i C_t(x_i)}$ : down-weight correct predictions, up-weight incorrectly predictions.
- 5 **return**  $C(x) = \text{sign}\left(\sum_t^T \alpha_t C_t(x)\right)$

# Record Linkage (RL): Generic Workflow

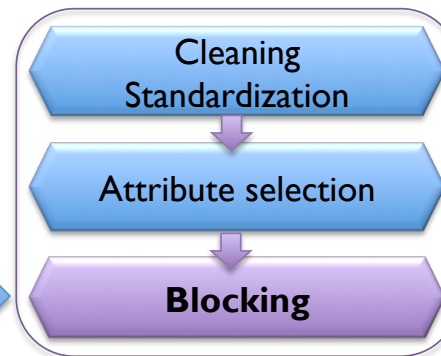
Database R

Name	SSN	Addr
Will Forth	354-564-339	Ada Bd
Jacky Khan	435-232-129	Marple Street
Dom Hack	235-575-689	Main Street
...	...	...

Database S

Name	SSN	Addr
Jack Khan	435-223-129	Marple St
Hans Ford	354-564-339	Clover Bd
Tom Hack	235-557-689	Main St
...	...	...

R X S



[Fellegi, Sunter, 1969]  
[Christen, 2012]

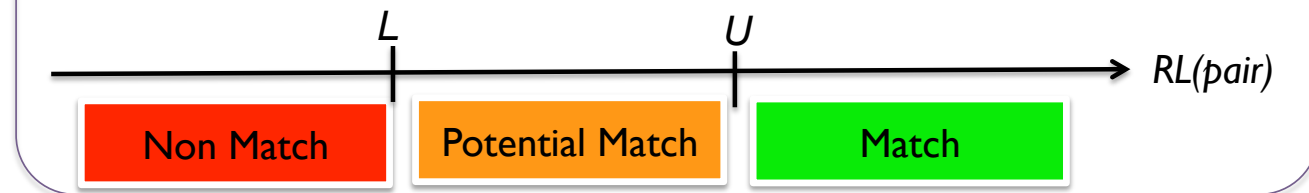
- Hashing
- Sorted keys
- Sorted NN
- (Multiple) Windowing
- Clustering



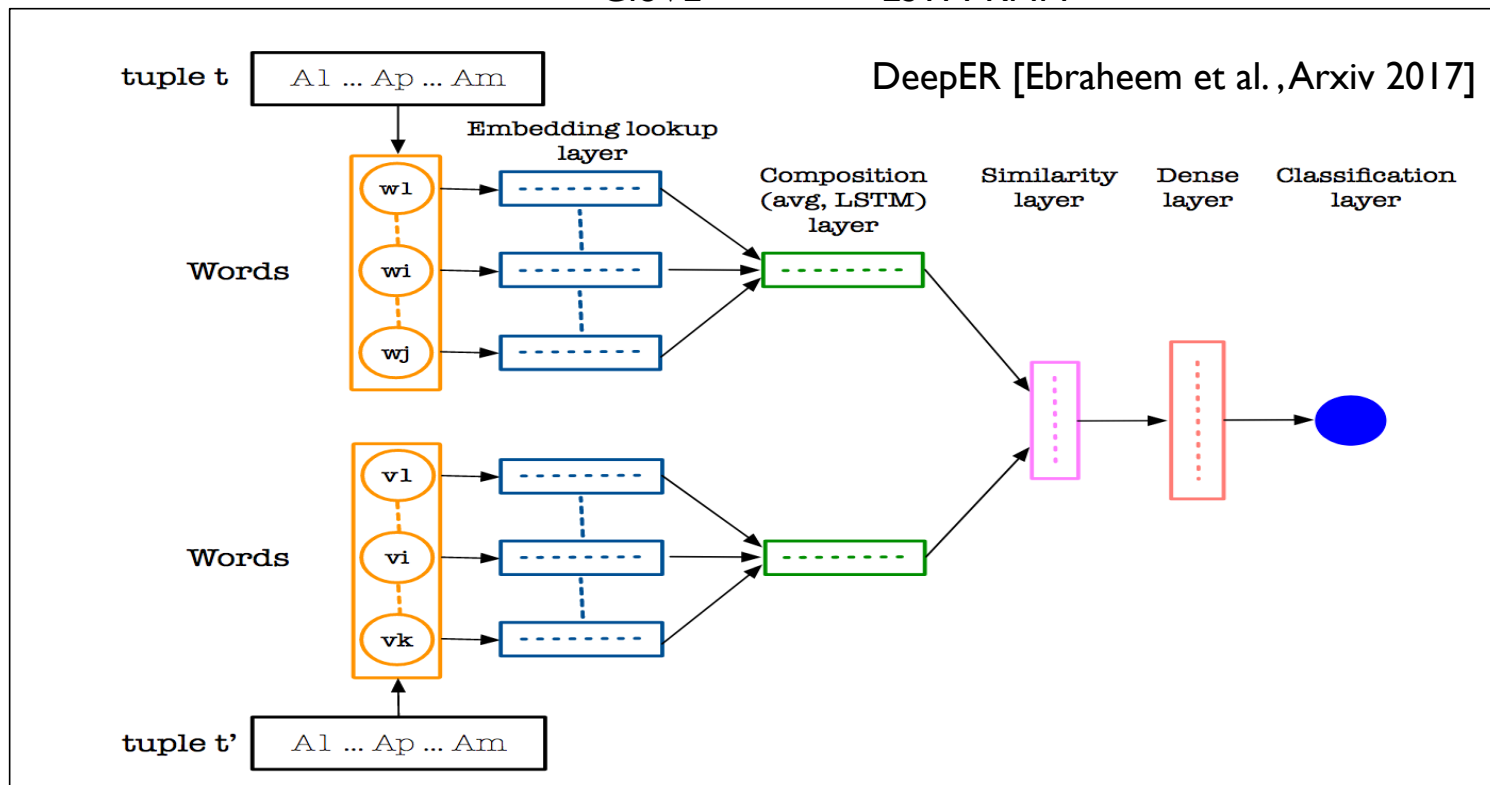
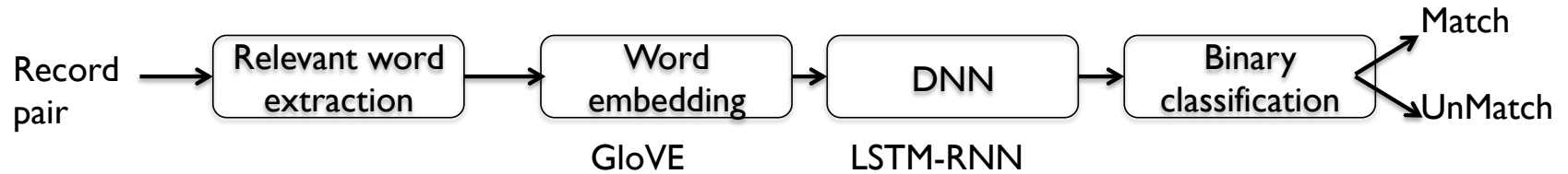
- Token-based : N-grams...
- Distance-based: Jaro, Edit, Levenshtein, Soundex
- Domain-dependent



$$\text{Linkage decision: } RL(\text{pair}) = \frac{P(\text{vector} \mid \text{pair} \in \text{Match})}{P(\text{vector} \mid \text{pair} \in \text{Non Match})}$$



# Deep learning for Entity Resolution



# Outline

---

1. **Detection of data quality problems:**

Profiling data quality

2. **Data cleaning**

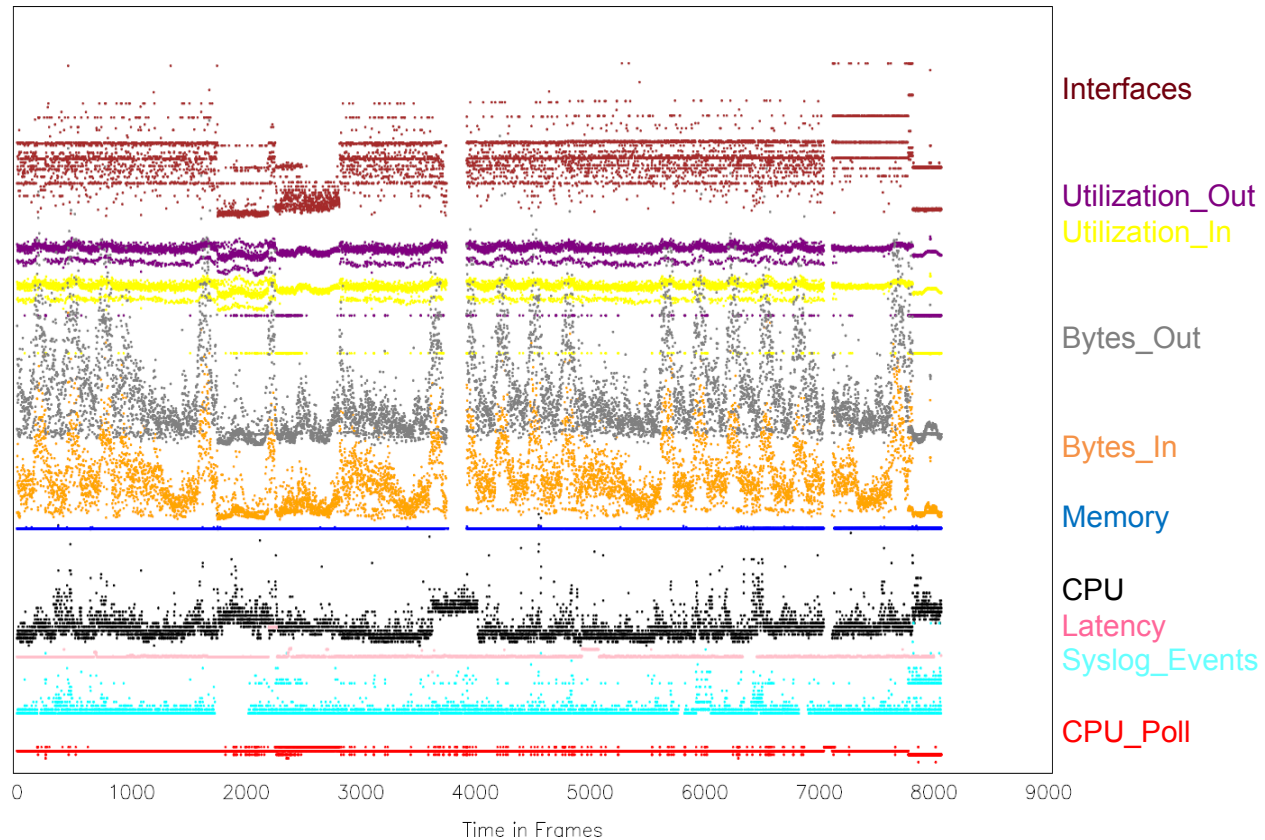
Leveraging the patterns of glitches

3. **Data preparation strategies:**

Learning to clean and prepare the data

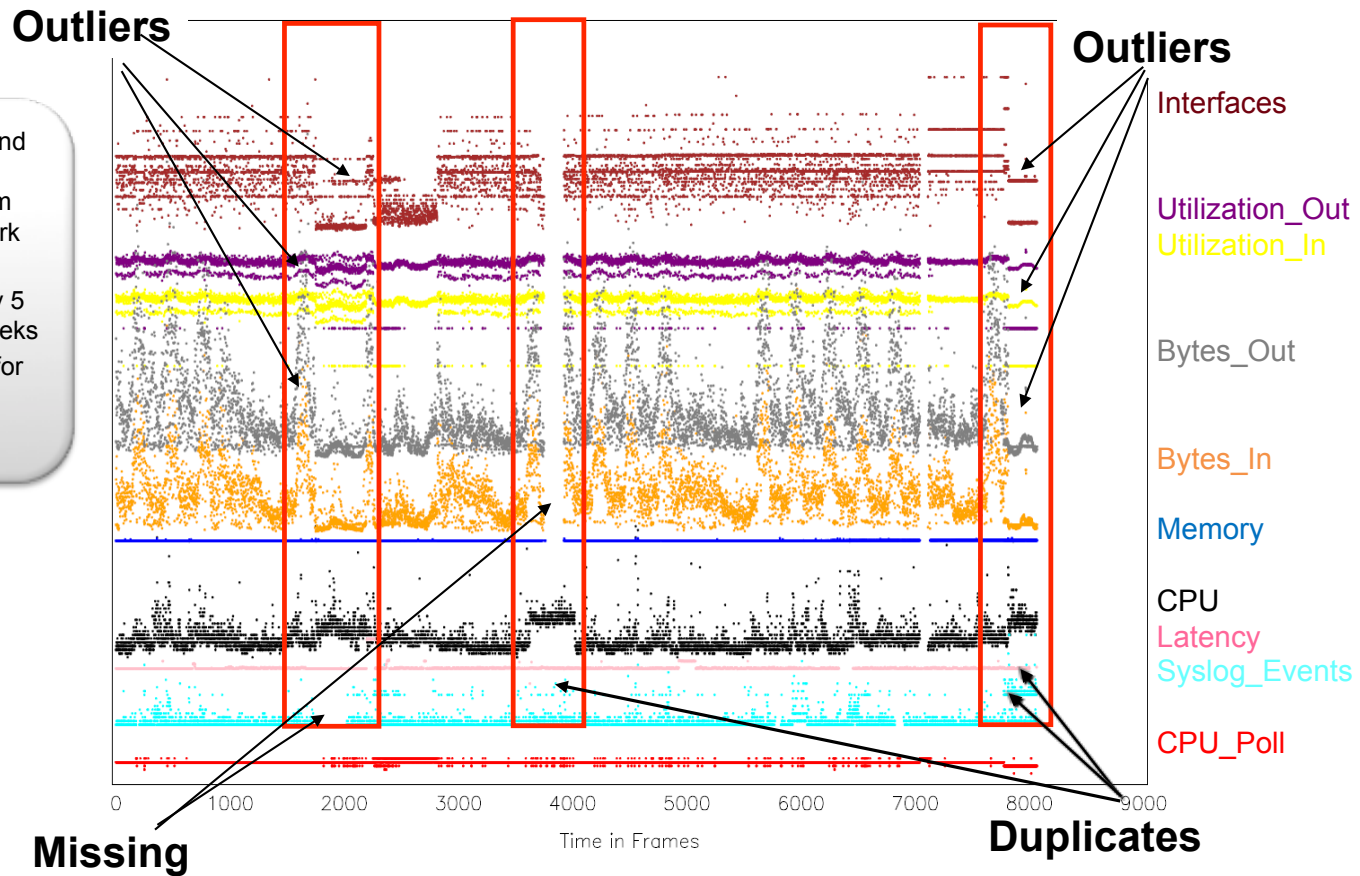
# SNMP Data Analysis

- Periodic inbound and outbound traffic measurements from interfaces of network devices
- 10 attributes, every 5 minutes, over 4 weeks
- Axes transformed for plotting

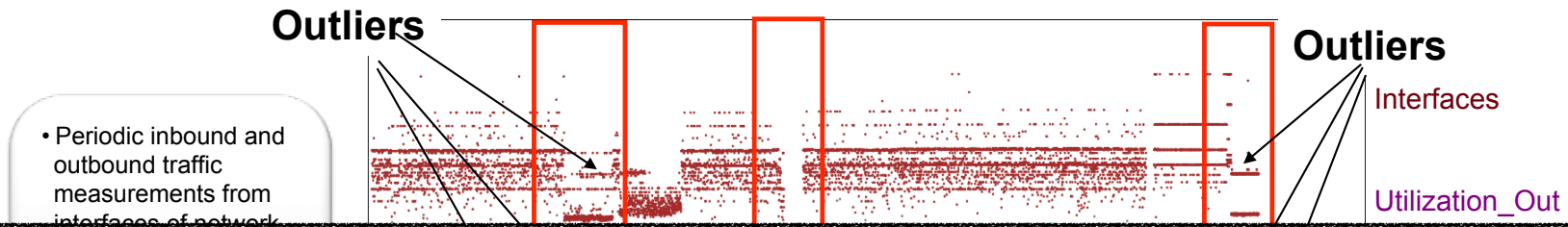


# SNMP Data Analysis

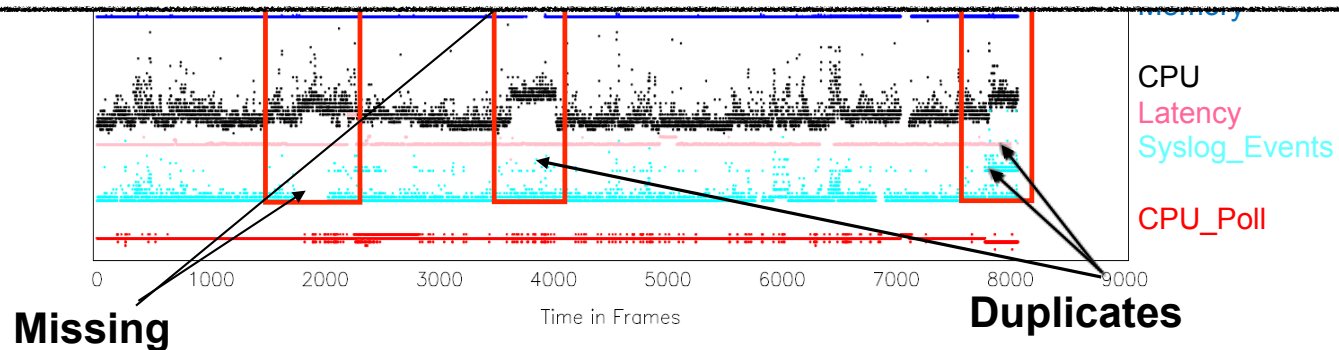
- Periodic inbound and outbound traffic measurements from interfaces of network devices
- 10 attributes, every 5 minutes, over 4 weeks
- Axes transformed for plotting



# SNMP Data Analysis



1. Detect patterns of multivariate, concomitant data anomalies
2. Use the anomaly patterns for consistent cleaning



# Understanding Complex Glitch Patterns

---

## **Benefits**

- A common root cause can generate correlated data errors
- In-depth anomaly analysis could help for:
  - Characterizing anomaly sources, processes, and propagation mechanisms
  - Systematizing data cleaning

## **Current methods**

- Make unrealistic assumptions (e.g., MAR)
- Treat glitches in isolation
- Are one-shot approaches (no reiteration between detection and cleaning)

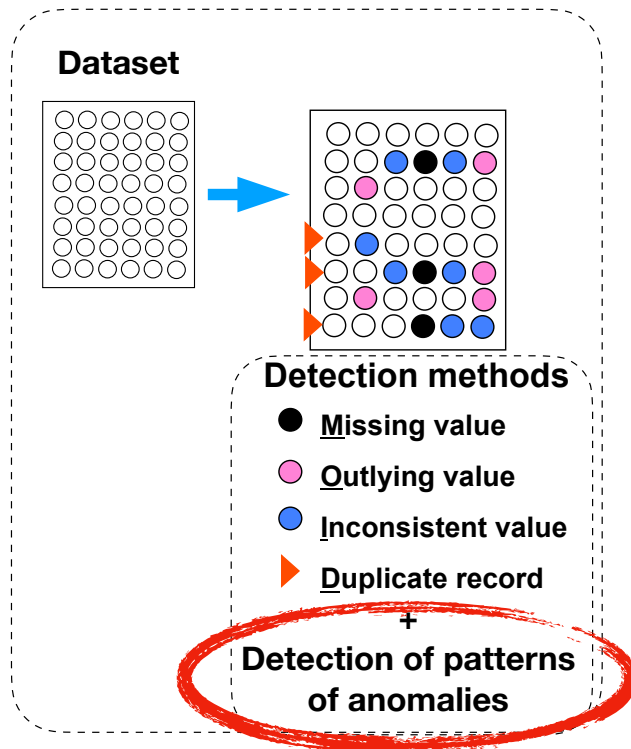
Data cleaning and preprocessing may introduce new errors and distortions.



# Detection-Exploration-Cleaning Framework

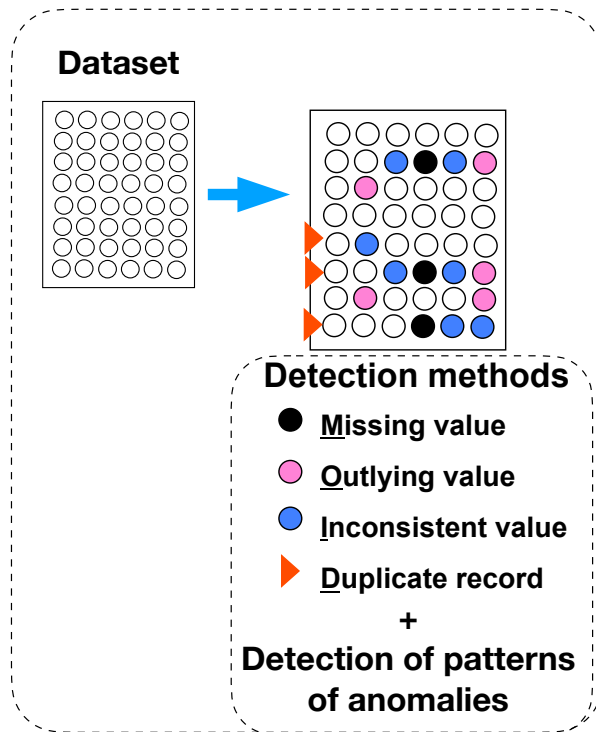
---

*Input:*

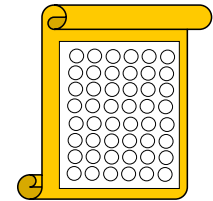


# Detection-Exploration-Cleaning Framework

*Input:*



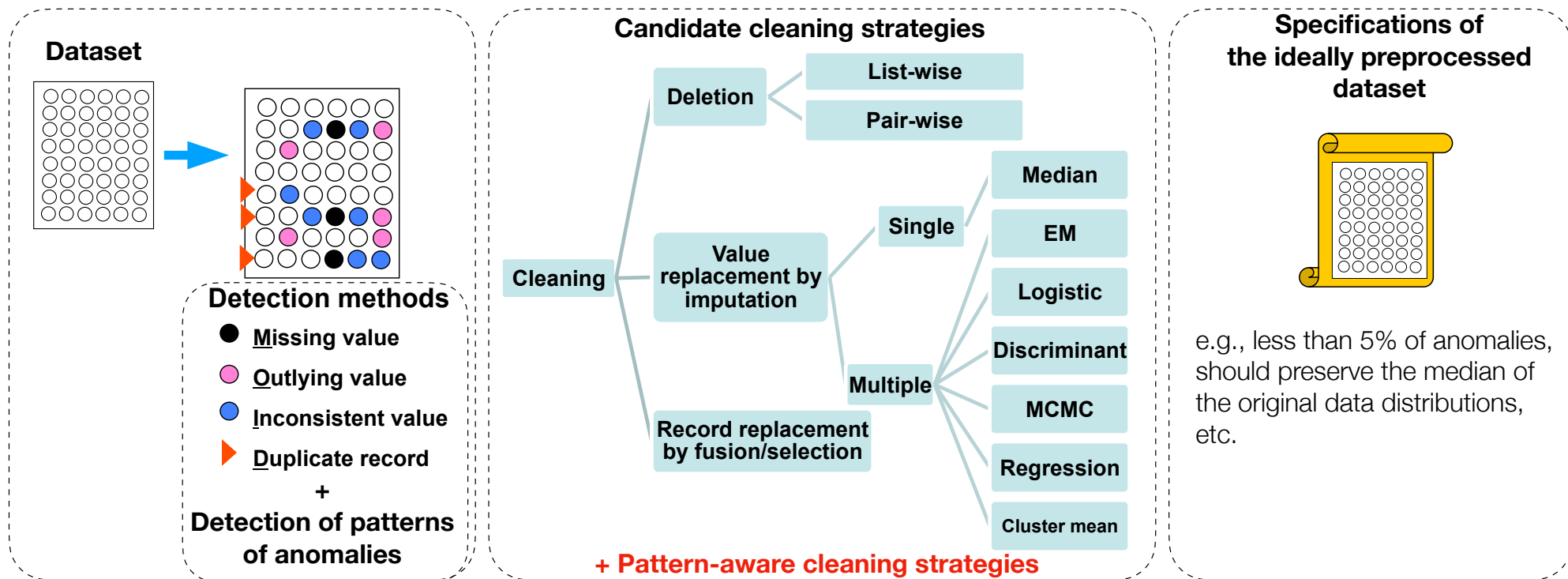
**Specifications of the ideally preprocessed dataset**



e.g., less than 5% of anomalies, should preserve the median of the original data distributions, etc.

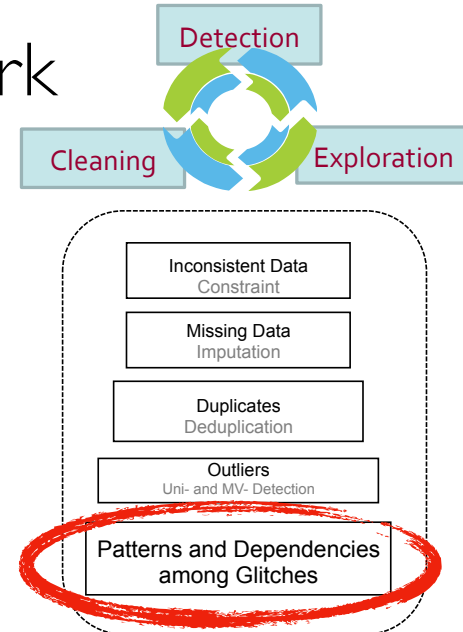
# Detection-Exploration-Cleaning Framework

*Input:*



# Detection-Exploration-Cleaning Framework

[Berti-Equille, Dasu, Srivastava, ICDE 2011]



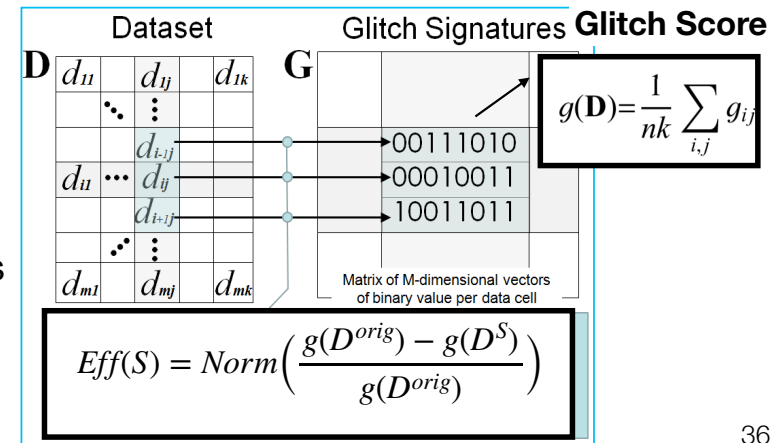
## Problem Statement:

Find the quantitative cleaning strategy  $B$  composed of  $M$  methods among the candidate strategies  $S$  such that its resulting dataset  $D^B$  is the closest to an ideal dataset  $D^*$  specified from  $D$  as

$$D^B = \arg \min_{\{s \in S\}} (\text{dist}(D^s, D^*))$$

subject to  $Cost(s) \leq U$  and  $Eff(s) \geq \Gamma > 0$

- $\mathbf{dist}$  is the Kullback-Leibler distance between two data distributions
- $\mathbf{U}$  is a pre-defined upper bound for the cost of strategy  $s$
- $\mathbf{\Gamma}$  is the lower bound of  $Eff(s)$ , the effectiveness of strategy  $s$



# Experiments

---

## Real-world and semi-synthetic data

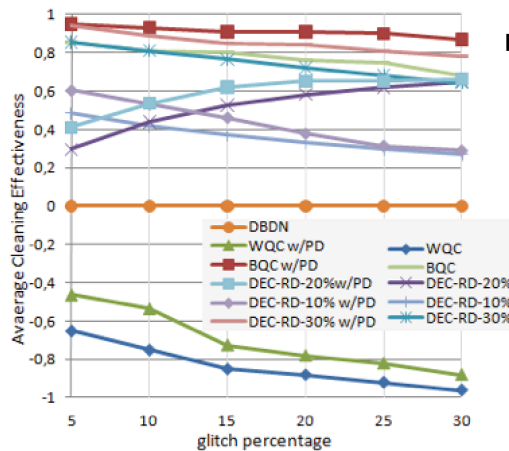
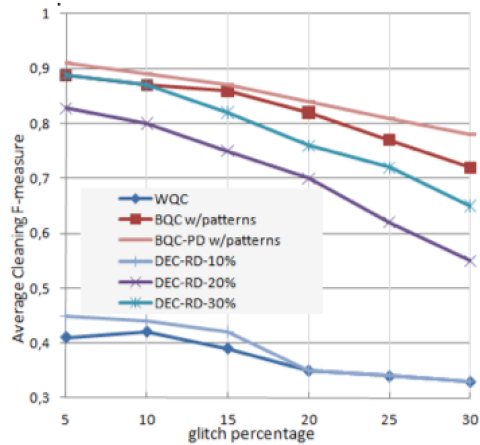
- **EPO Dataset:** 754,075 records, 4 non-key attributes (string, categorical and numerical data)
- **Intel Berkeley Research lab Dataset:** 2,313,682 million readings, 8 attributes (timestamp, sensorID, temperature, light, voltage) collected every 31 seconds from 54 sensors deployed in the between February 28th and April 5th, 2004
- **SNMP Dataset:** (8,632 tuples, 11 variables) collected every 5 minutes during one month (timestamps, categorical and numerical values)

## Comparison of various cleaning strategies

- Cost-based
- Effectiveness-based
- Resource-driven to treat just p% of glitches (DEC-RD)
- Specification-driven to treat a particular glitch type (DEC-SD)
- Pattern-based (DEC-PD)

# Experimental results

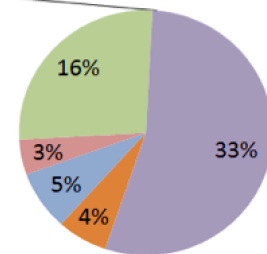
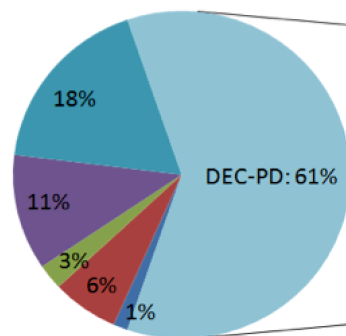
## SNMP



➔ **Pattern discovery always improves the accuracy when glitch percentage increases.**

➔ **Effectiveness is improved by + 8% in average.**

➔ **61% of the best strategies are pattern-based.**



- DEC-RD-10%
- DEC-SD (M, O, I, D)
- DEC-MD- (.3, .5, .7)
- DEC-RD-20%
- DEC-RD-30%
- DEC-RD-10%-PD
- DEC-SD (M, O, I, D)-PD
- DEC-MD- (.3, .5, .7)-PD
- DEC-RD-20%-PD
- DEC-RD-30%-PD

# Outline

---

1. **Detection of data quality problems:**

Profiling data quality with MeSQuaL

2. **Data cleaning**

Leveraging the patterns of glitches

3. **Data preparation strategies**

Learning to clean and prepare the data

# Data preprocessing is challenging



Data preparation pipeline

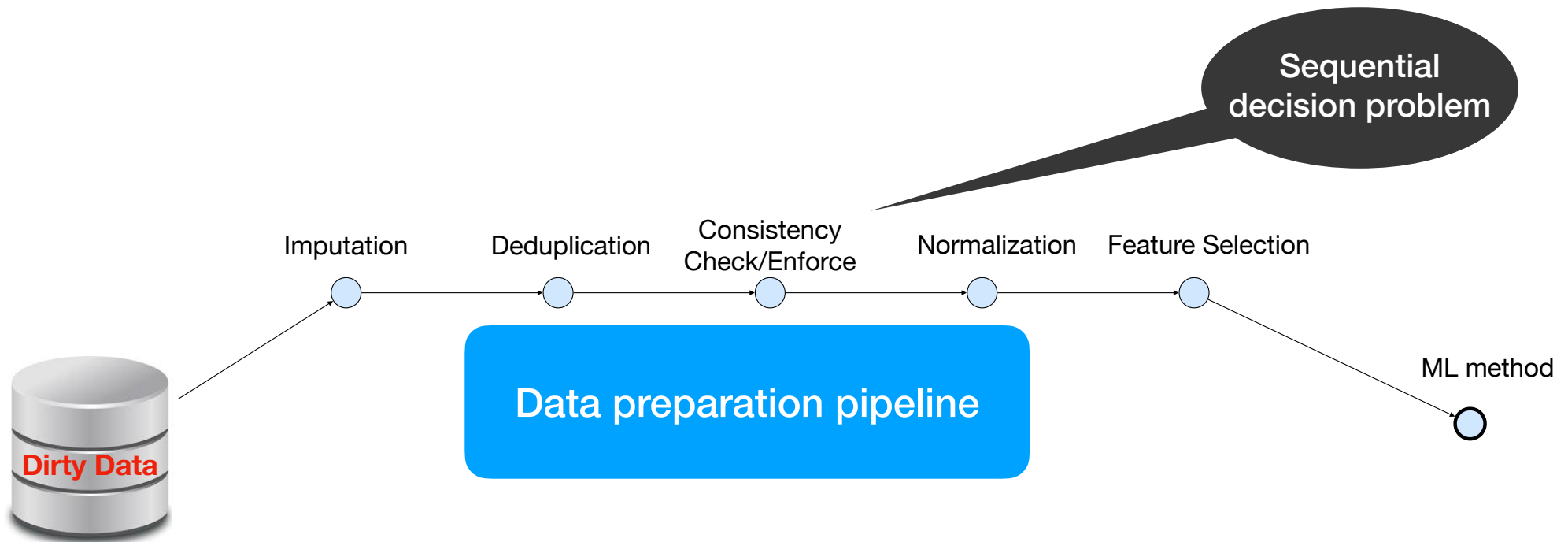


ML method

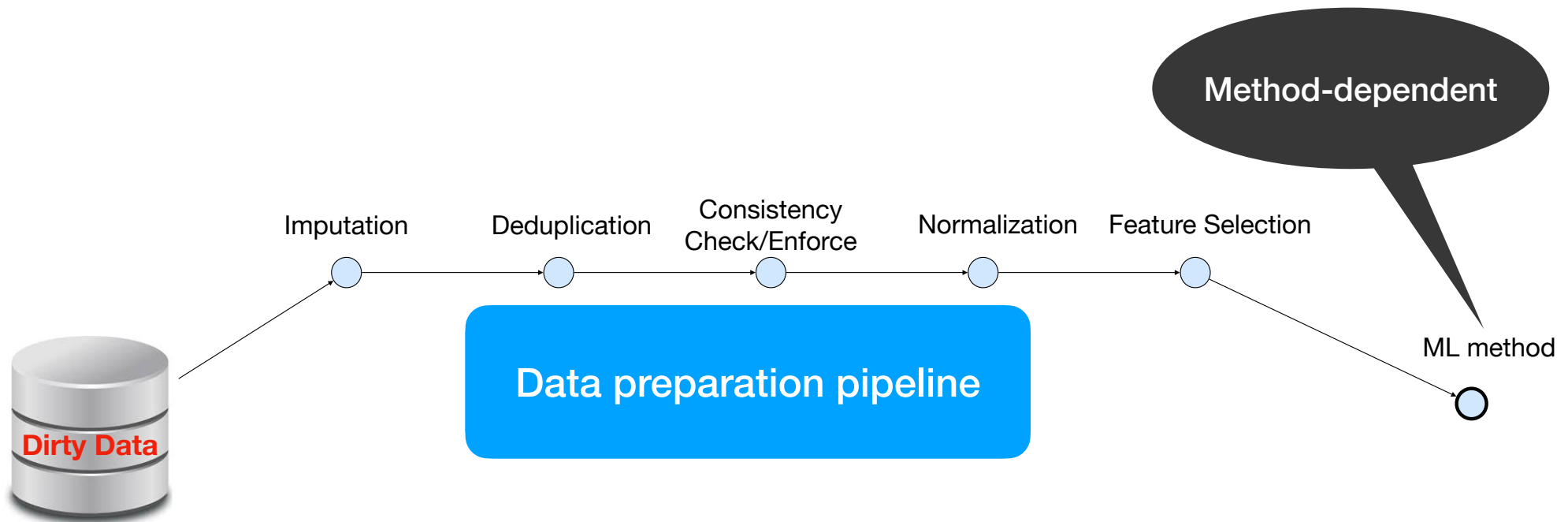




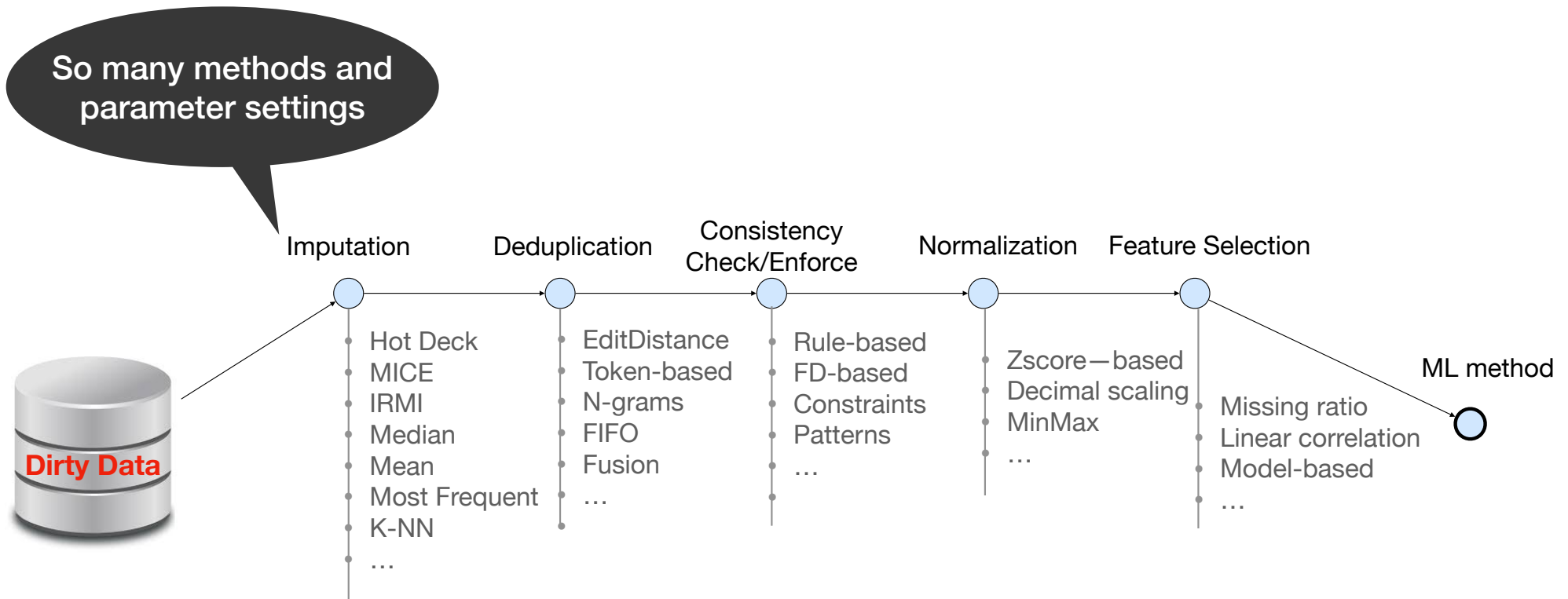
# Data preprocessing is challenging



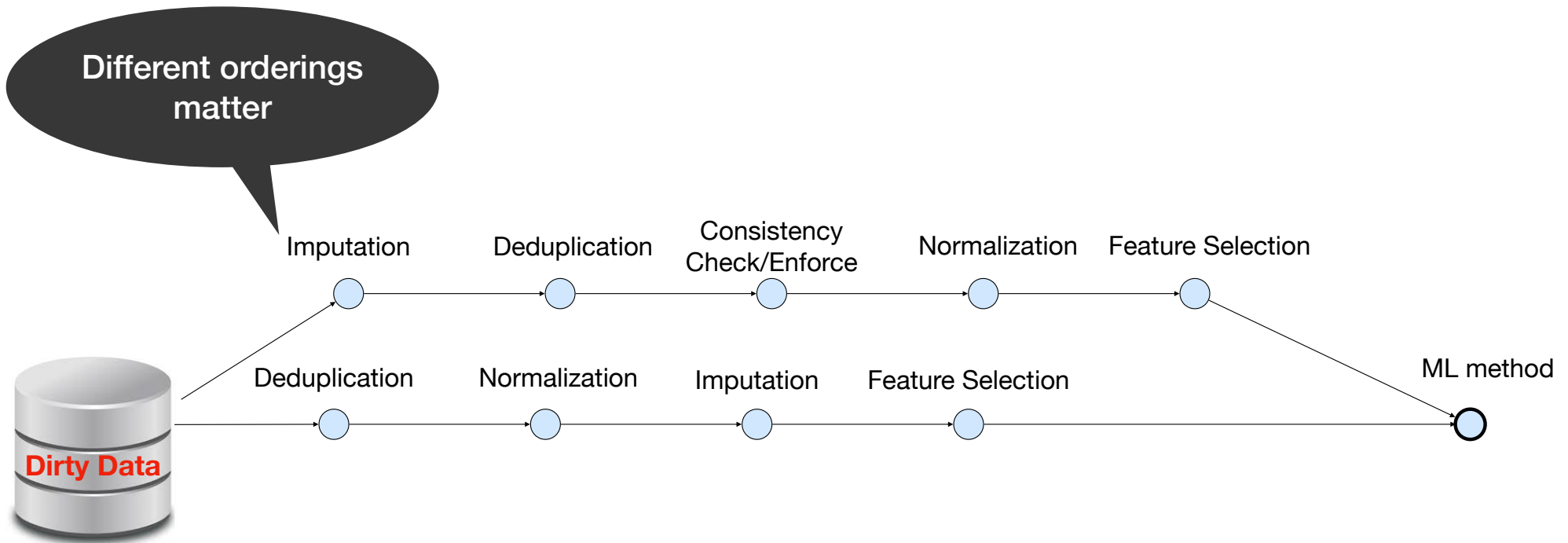
# Data preprocessing is challenging



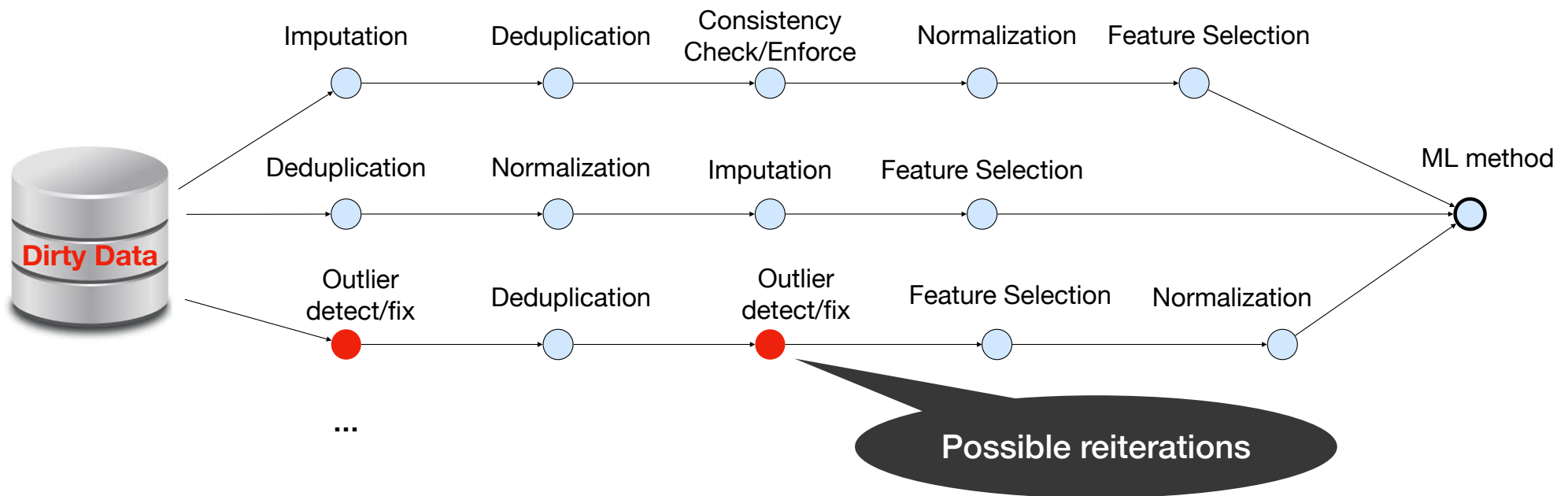
# Data preprocessing is challenging



# Data preprocessing is challenging

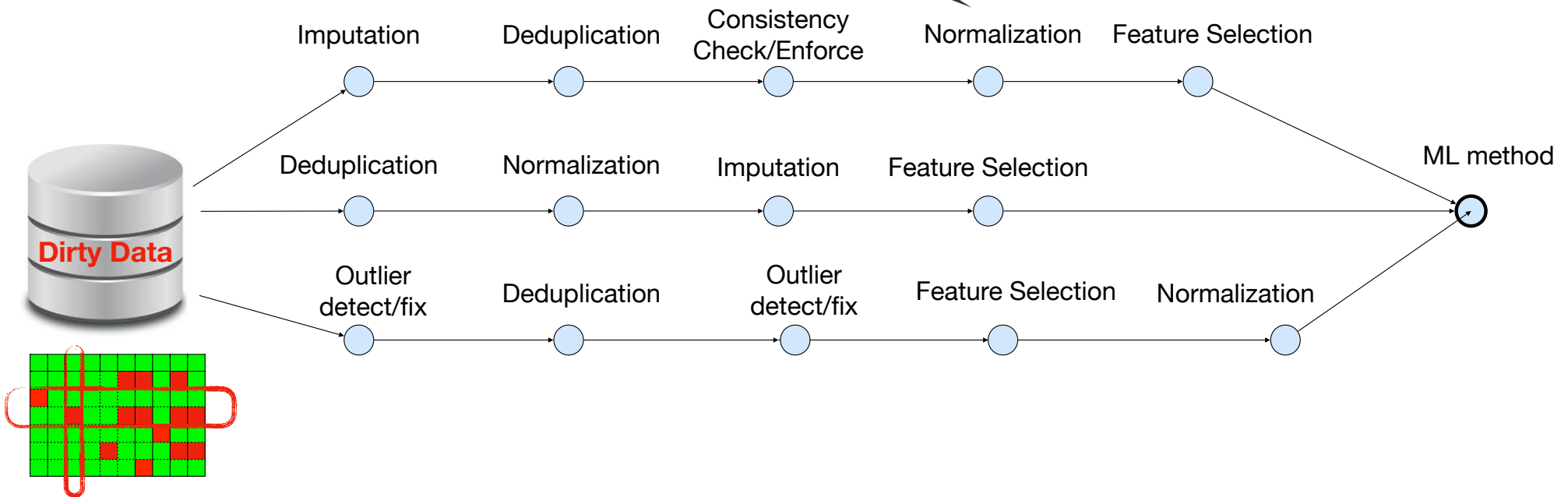


# Data preprocessing is challenging

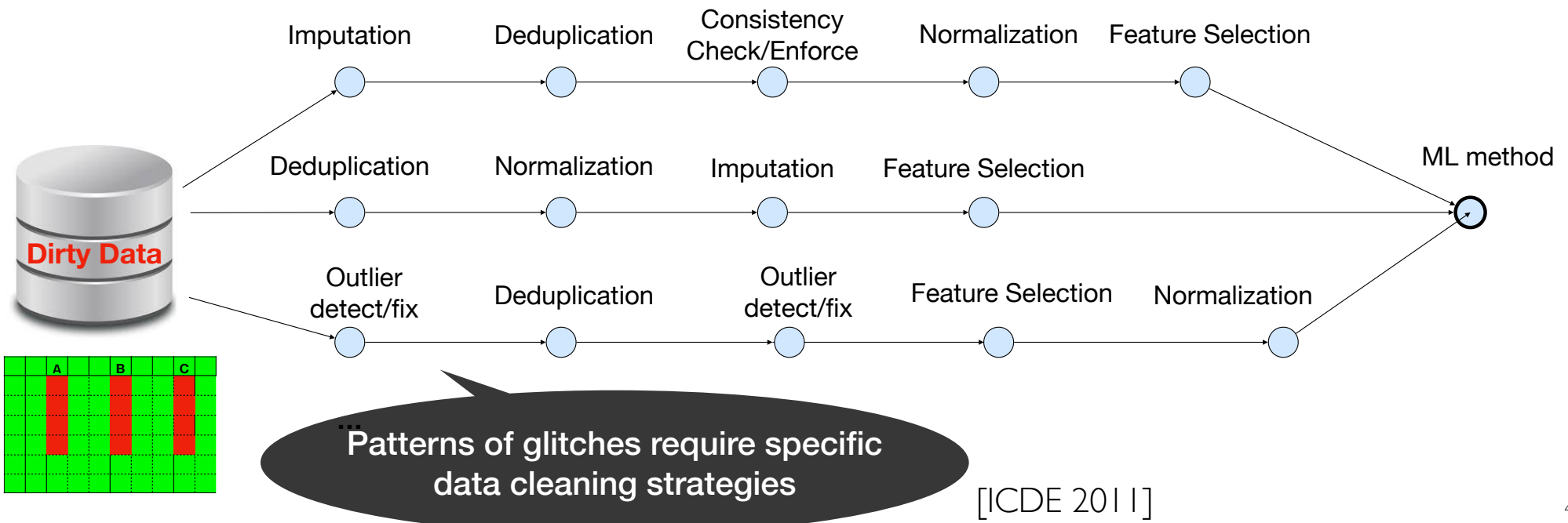


# Data preprocessing is challenging

Selective processing of some parts of the dataset

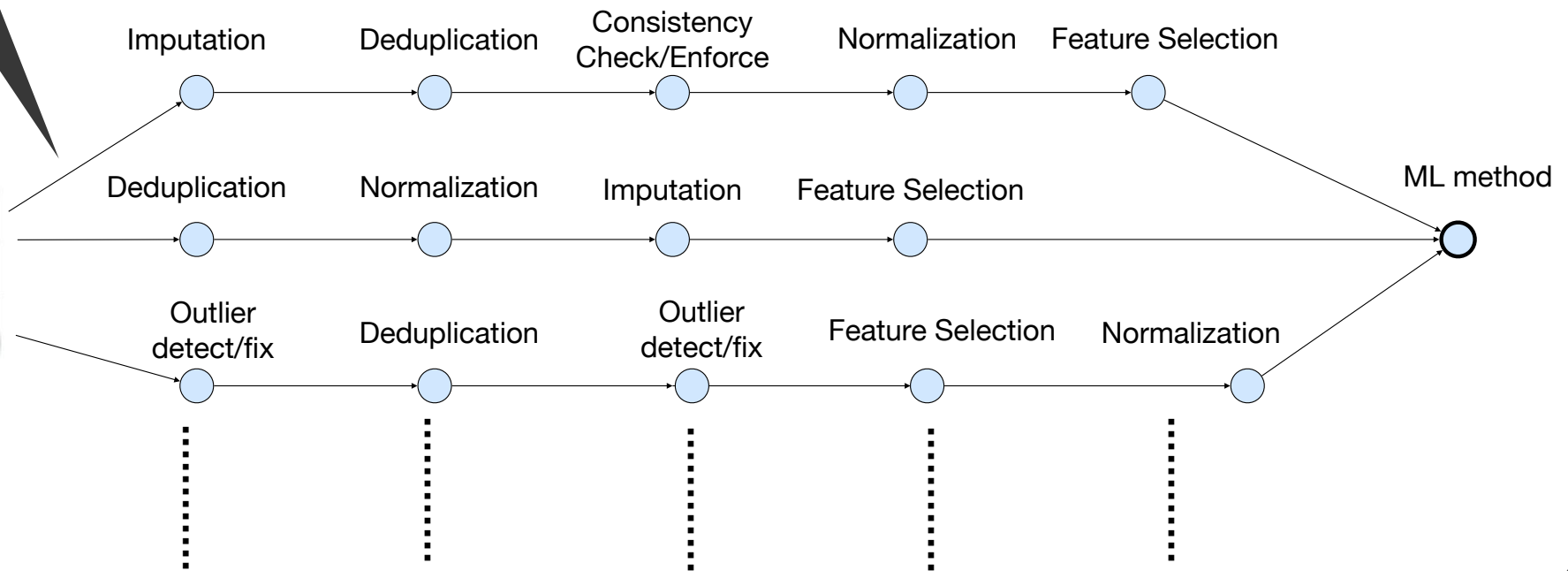


# Data preprocessing is challenging



# Data preprocessing is challenging

Infinite space of possible strategies

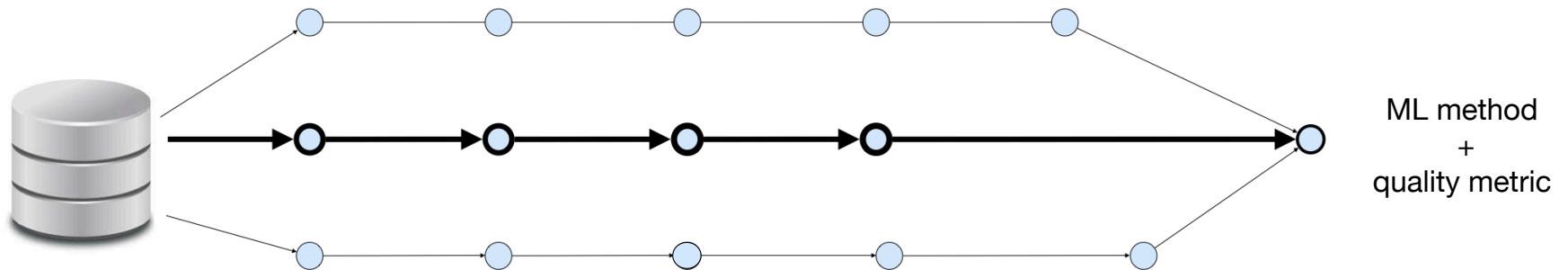




# Optimization Problem



Can we help the user in composing the data preparation pipeline that maximizes the quality performance of the ML method ?



# Optimization Problem



Can we help the user in composing the data preparation pipeline that maximizes the quality performance of the ML method ?

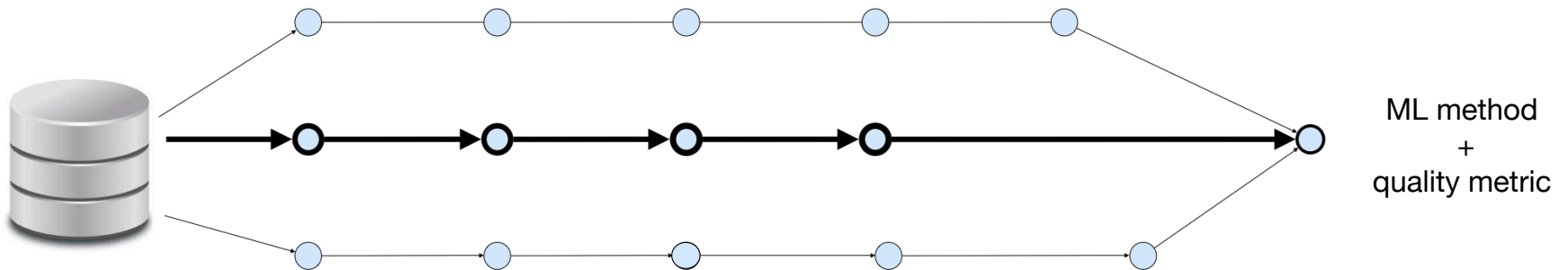
No training  
example for "good"  
data cleaning

Metric-dependent

No model a priori

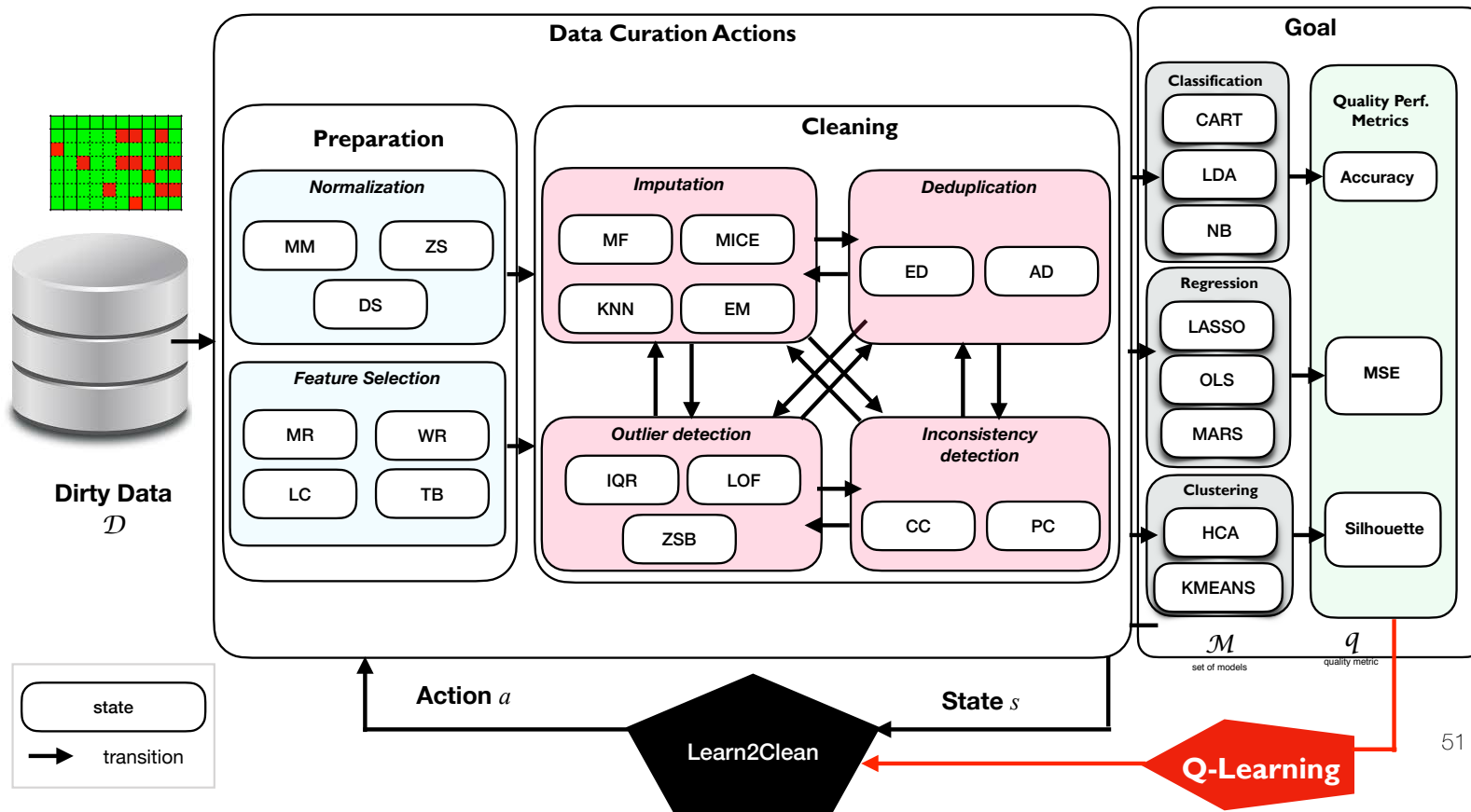
AutoML-like approach

Human-In-The Loop



# First Solution: Learn2Clean

<https://github.com/LaureBerti/Learn2Clean>



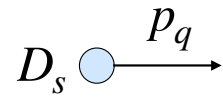
*AutoML-like approach for data curation*

# Reinforcement Learning Framework

Markov Decision Process



**Learn2Clean**



# Reinforcement Learning Framework

Markov Decision Process

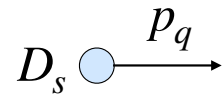
State

Action

Transition

Reward

Learn2Clean

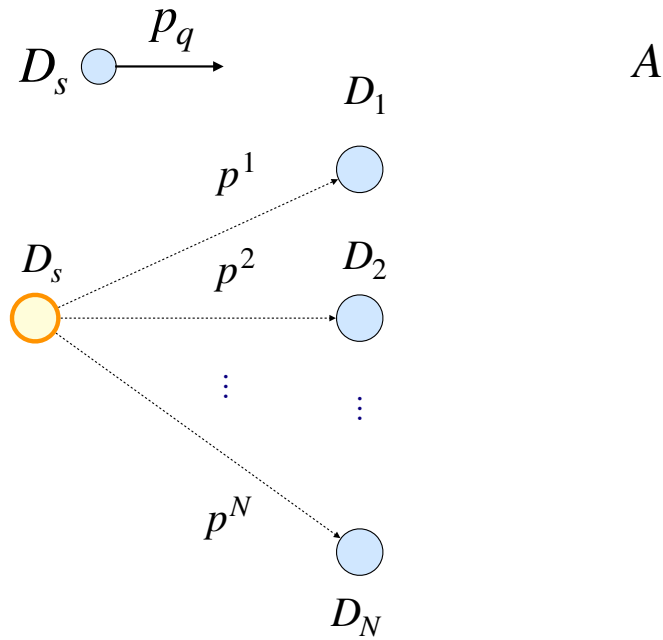


# Reinforcement Learning Framework

Markov Decision Process



**Learn2Clean**

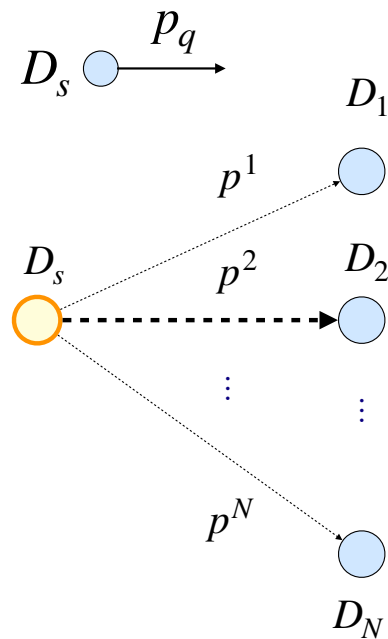


# Reinforcement Learning Framework

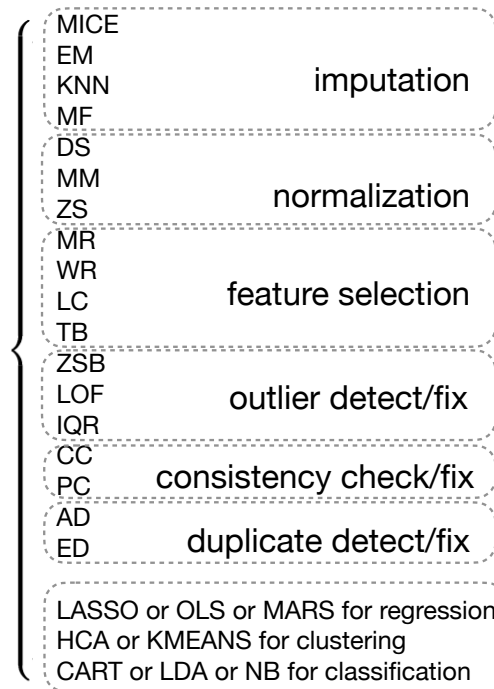
Markov Decision Process



Learn2Clean



A

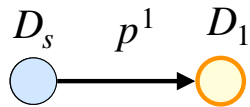
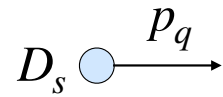


# Reinforcement Learning Framework

Markov Decision Process



Learn2Clean



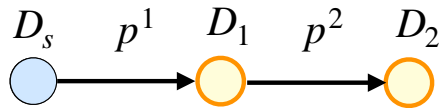
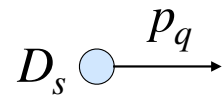


# Reinforcement Learning Framework

Markov Decision Process



Learn2Clean

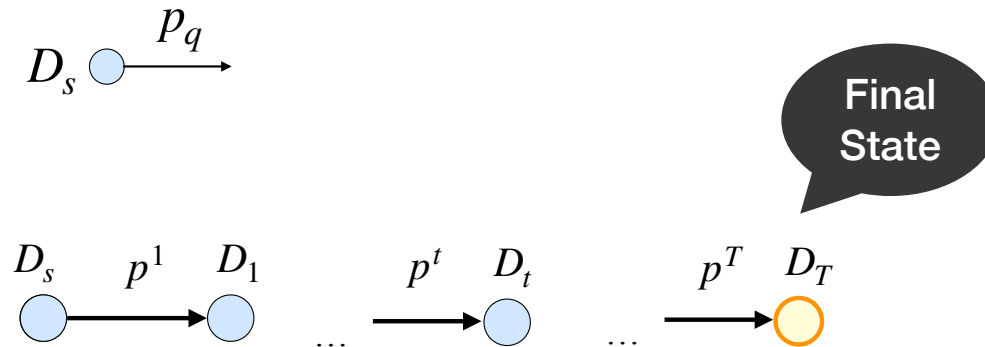


# Reinforcement Learning Framework

Markov Decision Process



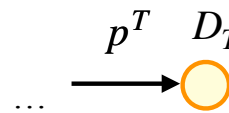
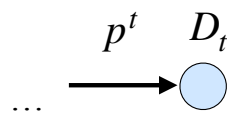
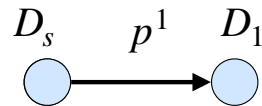
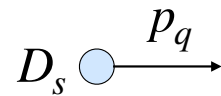
Learn2Clean



LASSO or OLS or MARS for regression  
HCA or KMEANS for clustering  
CART or LDA or NB for classification

# Reinforcement Learning Framework

Markov Decision Process



*deterministic*

$R =$

MICE EM KNN MF DS MM ZS MR WR LC TB ZSB LOF IQR CC PC AD ED LASSO

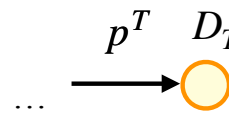
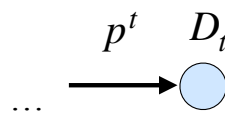
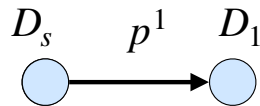
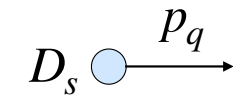
-1	-1	-1	-1	0	0	0	-1	0	0	0	0	0	0	0	0	0	0	0	-1	-1	0	0	100	
-1	-1	-1	-1	0	0	0	-1	0	0	0	0	0	0	0	0	0	0	0	-1	-1	0	0	100	
-1	-1	-1	-1	0	0	0	-1	0	0	0	0	0	0	0	0	0	0	0	-1	-1	0	0	100	
-1	-1	-1	-1	-1	-1	-1	0	0	0	0	0	0	0	0	0	0	0	0	-1	-1	0	0	-1	
-1	-1	-1	-1	-1	-1	-1	0	0	0	0	0	0	0	0	0	0	0	0	-1	-1	0	0	-1	
0	0	0	0	-1	-1	-1	0	0	0	0	0	0	0	0	0	0	0	0	-1	-1	0	0	-1	
0	0	0	0	0	0	0	-1	-1	-1	0	0	0	0	0	0	0	0	0	-1	-1	0	0	-1	
0	0	0	0	0	0	0	0	-1	-1	-1	0	0	0	0	0	0	0	0	-1	-1	0	0	-1	
0	0	0	0	-1	-1	-1	0	0	0	0	-1	-1	-1	0	0	0	0	0	-1	-1	0	0	-1	
-1	-1	-1	-1	-1	-1	-1	0	0	0	0	-1	-1	-1	-1	-1	0	0	0	-1	-1	0	0	100	
-1	-1	-1	-1	-1	-1	-1	0	0	0	0	-1	-1	-1	-1	-1	0	0	0	-1	-1	0	0	100	
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	-1	-1	100
0	0	0	0	-1	-1	-1	0	0	0	0	-1	-1	-1	-1	-1	0	0	0	-1	-1	0	0	100	
-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	0	0	0	-1	-1	0	0	100	
-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	0	0	0	-1	-1	0	0	100	
-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	0	0	0	-1	-1	0	0	100	

# Reinforcement Learning Framework

Markov Decision Process



Learn2Clean



Final State

deterministic

R =

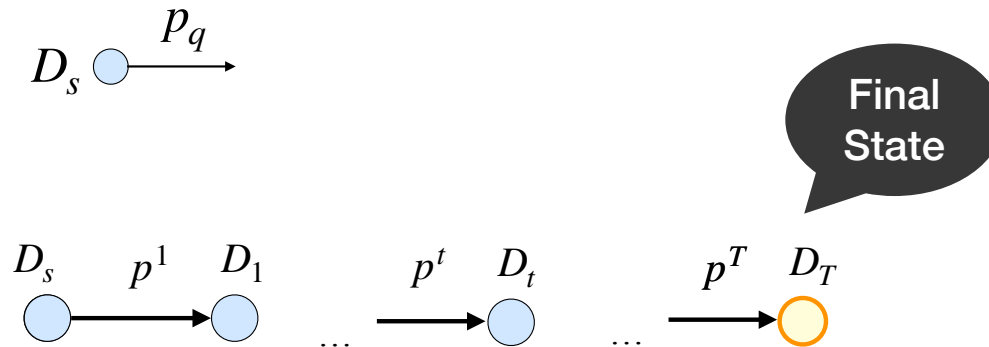
	MICE	EM	KNN	MF	DS	MM	ZS	MR	WR	LC	TB	ZSB	LOF	IQR	CC	PC	AD	ED	LASSO	
-1	-1	-1	-1	0	0	0	-1	0	0	0	0	0	0	0	-1	-1	0	0	100	
-1	-1	-1	-1	0	0	0	-1	0	0	0	0	0	0	0	-1	-1	0	0	100	
-1	-1	-1	-1	0	0	0	-1	0	0	0	0	0	0	0	-1	-1	0	0	100	
-1	-1	-1	-1	0	0	0	-1	0	0	0	0	0	0	0	-1	-1	0	0	100	
-1	-1	-1	-1	-1	-1	-1	0	0	0	0	0	0	0	0	-1	-1	0	0	-1	
-1	-1	-1	-1	-1	-1	-1	0	0	0	0	0	0	0	0	-1	-1	0	0	-1	
0	0	0	0	-1	-1	-1	0	0	0	0	0	0	0	0	-1	-1	0	0	-1	
0	0	0	0	0	0	0	-1	-1	-1	0	0	0	0	0	-1	-1	0	0	-1	
0	0	0	0	0	0	0	-1	-1	-1	-1	-1	0	0	0	-1	-1	0	0	-1	
0	0	0	0	-1	-1	-1	0	0	0	0	-1	-1	-1	-1	0	0	-1	-1	0	-1
-1	-1	-1	-1	-1	-1	-1	0	0	0	0	-1	-1	-1	-1	-1	-1	0	0	100	
-1	-1	-1	-1	-1	-1	-1	0	0	0	0	-1	-1	-1	-1	-1	-1	0	0	100	
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	-1	-1	100
0	0	0	0	-1	-1	-1	0	0	0	0	-1	-1	-1	-1	0	0	0	-1	-1	100
-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	100

- LASSO or OLS or MARS for regression
  - HCA or KMEANS for clustering
  - CART or LDA or NB for classification
- MSE
  - Silhouette
  - Accuracy

Quality metric

# Reinforcement Learning Framework

Markov Decision Process



*deterministic*

R =

	MICE	EM	KNN	MF	DS	MM	ZS	MR	WR	LC	TB	ZSB	LOF	IQR	CC	PC	AD	ED	LASSO
-1	-1	-1	-1	0	0	0	-1	0	0	0	0	0	0	0	-1	-1	0	0	100
-1	-1	-1	-1	0	0	0	-1	0	0	0	0	0	0	0	-1	-1	0	0	100
-1	-1	-1	-1	0	0	0	-1	0	0	0	0	0	0	0	-1	-1	0	0	100
-1	-1	-1	-1	-1	-1	-1	0	0	0	0	0	0	0	0	-1	-1	0	0	-1
-1	-1	-1	-1	-1	-1	-1	0	0	0	0	0	0	0	0	-1	-1	0	0	-1
0	0	0	0	-1	-1	-1	0	0	0	0	0	0	0	0	-1	-1	0	0	-1
0	0	0	0	0	0	0	-1	-1	-1	0	0	0	0	0	-1	-1	0	0	-1
0	0	0	0	0	0	0	-1	-1	-1	-1	0	0	0	0	-1	-1	0	0	-1
0	0	0	0	0	0	0	-1	-1	-1	-1	-1	0	0	0	-1	-1	0	0	-1
0	0	0	0	-1	-1	-1	0	0	0	0	0	-1	-1	-1	-1	0	0	0	-1
-1	-1	-1	-1	-1	-1	-1	0	0	0	0	-1	-1	-1	-1	-1	-1	0	0	100
0	0	0	0	-1	-1	-1	0	0	0	0	-1	-1	-1	-1	-1	-1	0	0	100
-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	0	0	0	0	-1	-1	0	0	100
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	-1	-1	100
0	0	0	0	-1	-1	-1	0	0	0	0	-1	-1	-1	-1	-1	-1	0	0	100
-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1

- > MSE
- > Silhouette
- > Accuracy

Quality metric

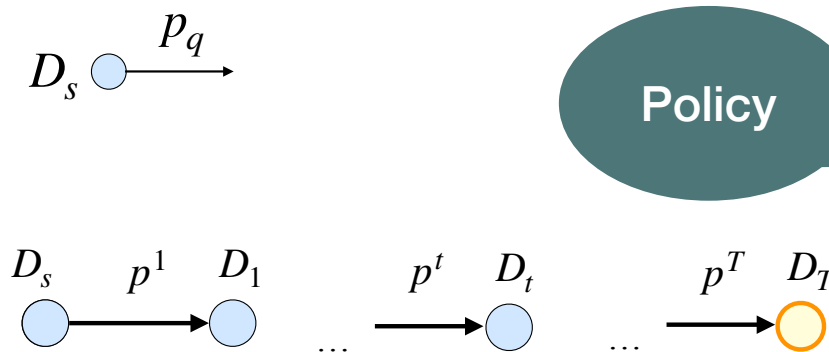
Learn2Clean selects the sequence of preprocessing actions that maximizes the quality metric (or minimizes the error)

# Reinforcement Learning Framework

Markov Decision Process



Learn2Clean



**Softmax action selection**

$$\pi = P(a | s) = \frac{e^{Q(s,a)/k}}{\sum_j e^{Q(s,a_j)/k}}$$

**Q-table**

*value iteration update*

$$Q^\pi(s, a) \leftarrow (1 - \alpha) \cdot Q(s, a) + \alpha \cdot \left( R(s, a) + \gamma \cdot \max_{a'} Q(s', a') \right)$$

↑ new value      ↗ learning rate      ↑ old value      ↑ reward      ↑ discount factor      ↖ optimal future value

learned value

# Experiments

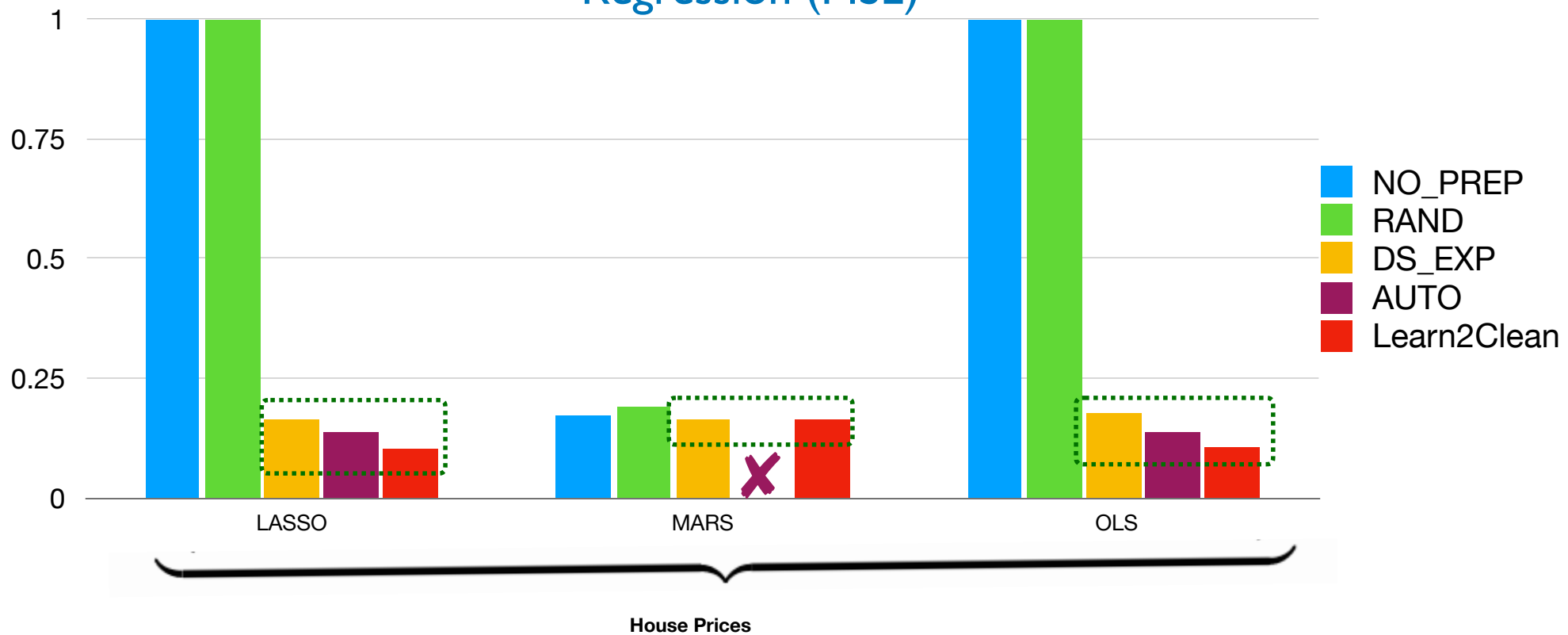
## Datasets

Name	# Att.	# Rows	Clustering	Regression	Classification
House Prices	81	1.46k	✓	✓	✓
Google Playstore Users	5	64.3k	✓		
Google Playstore Apps	13	10.8k	✓		✓

**Evaluation** : Silhouette for Clustering  
MSE for Regression  
Accuracy for Classification

# Experimental Results

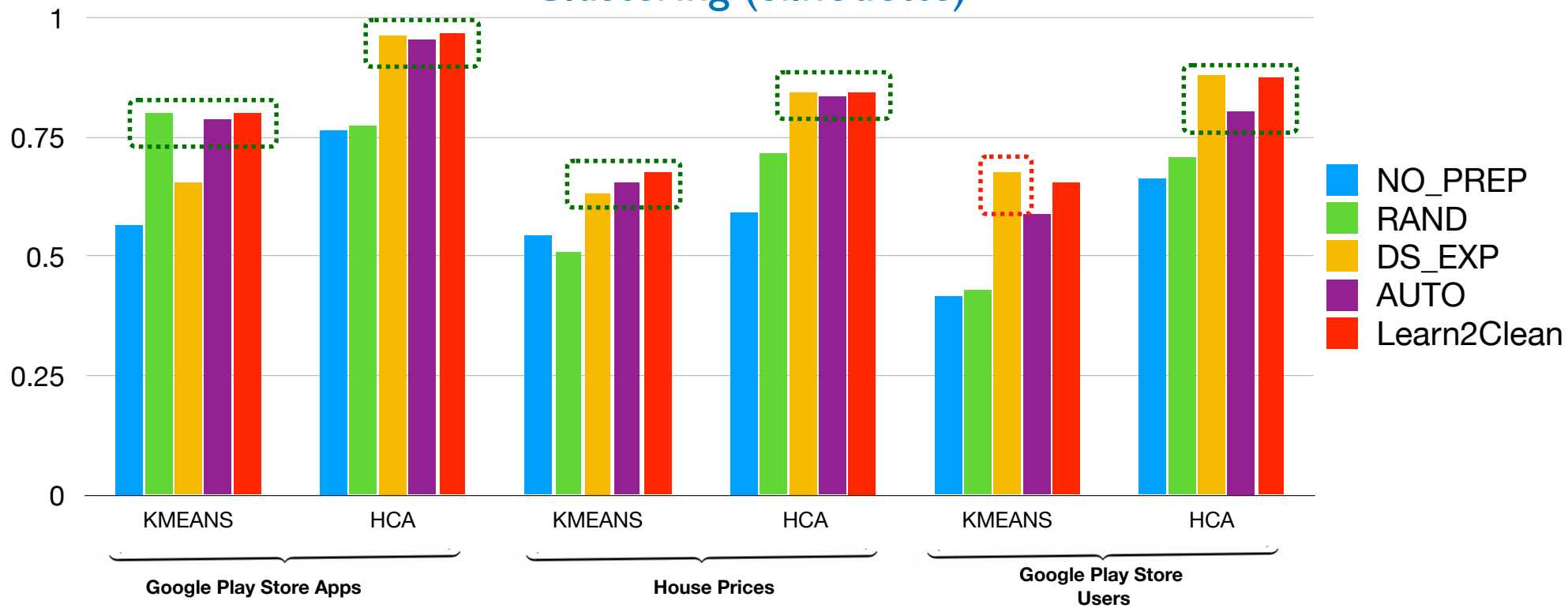
Regression (MSE)





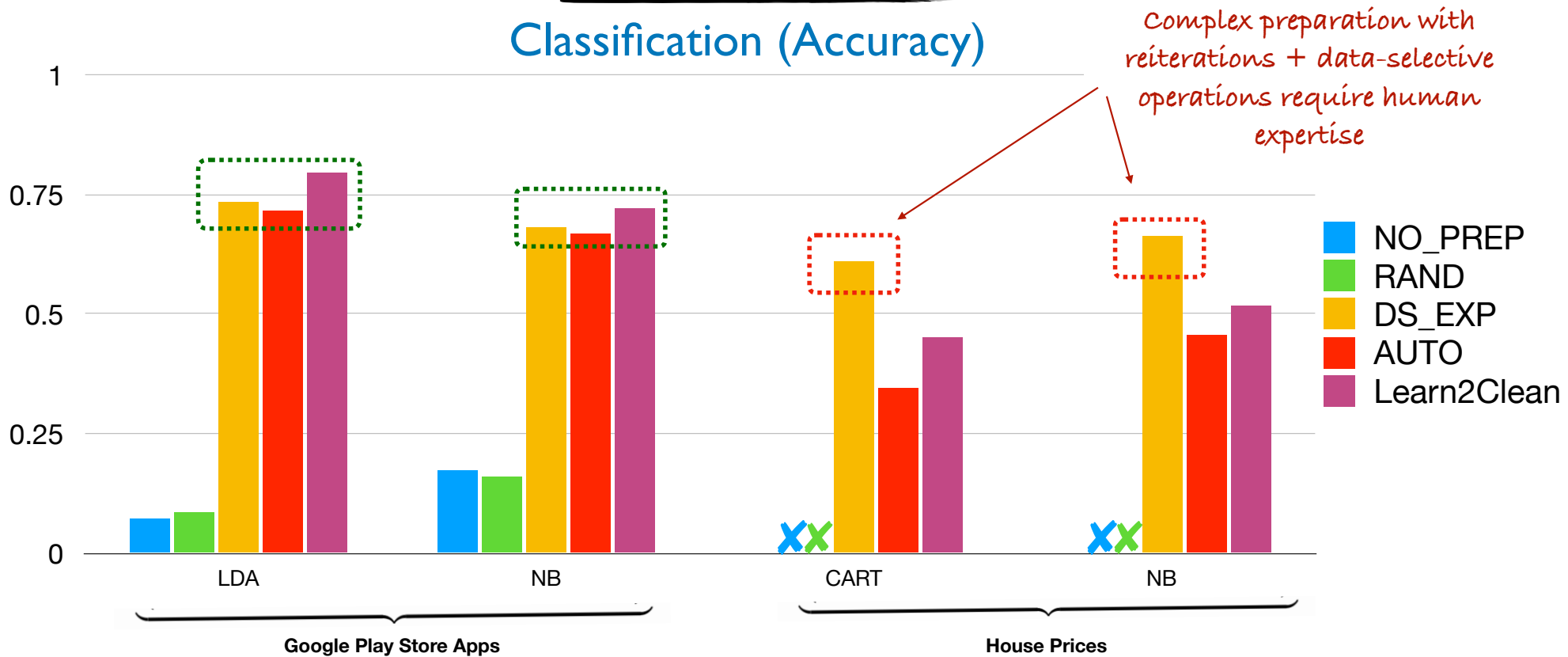
# Experimental Results

## Clustering (Silhouette)

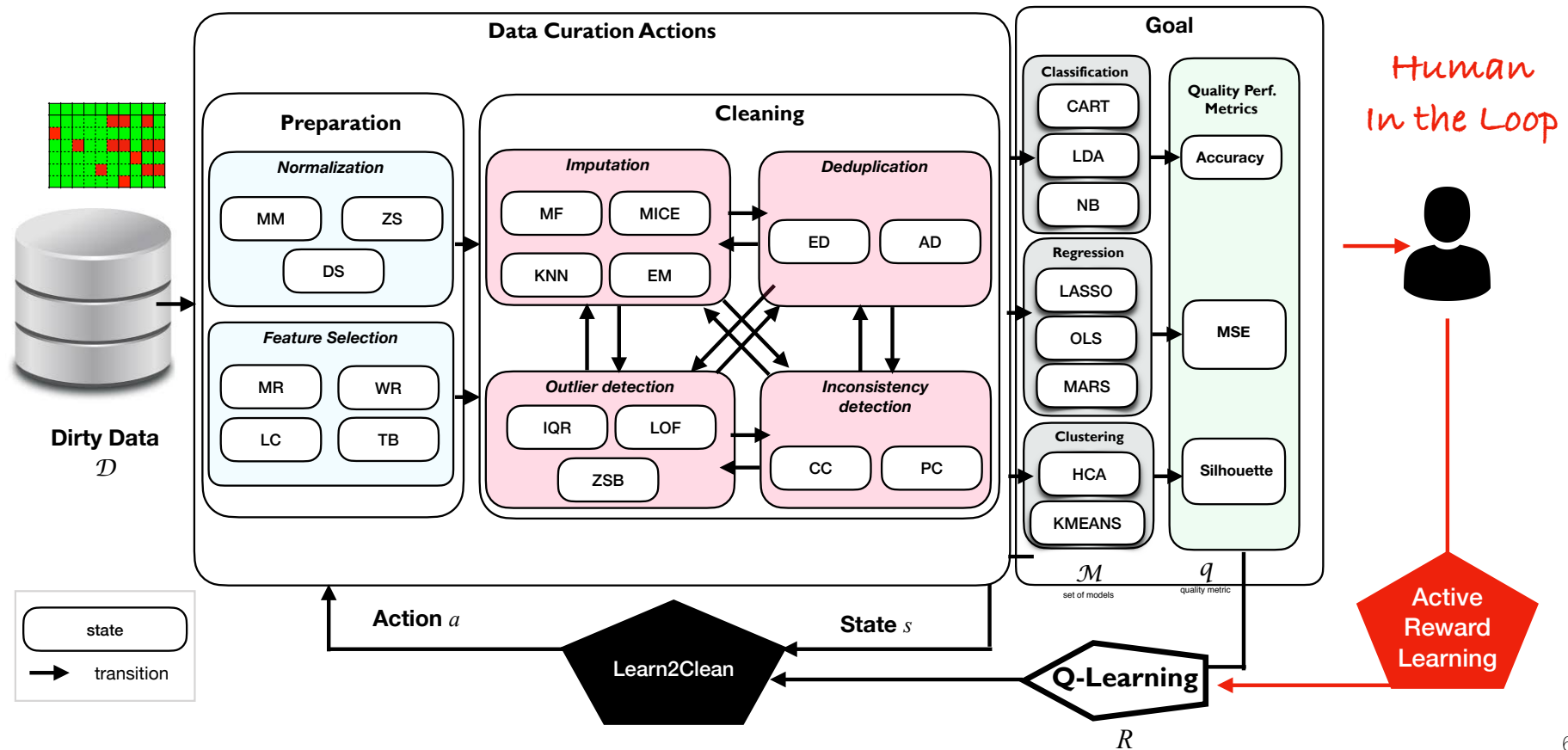


# Experimental Results

## Classification (Accuracy)

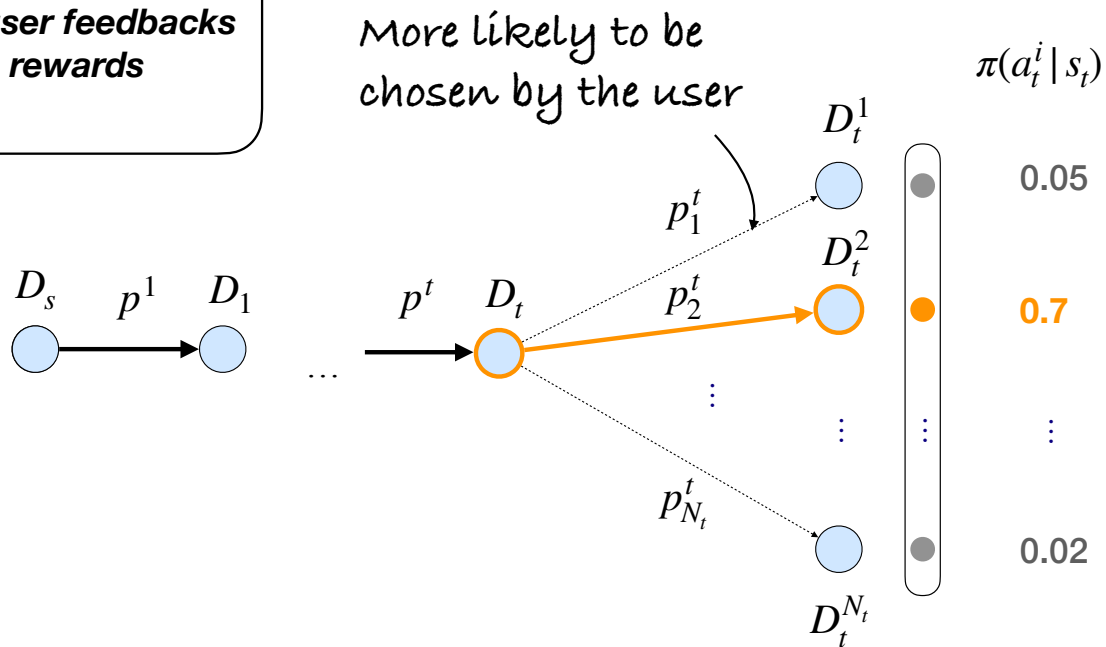


# HIL with Active Reward Learning



# Active Reward Learning

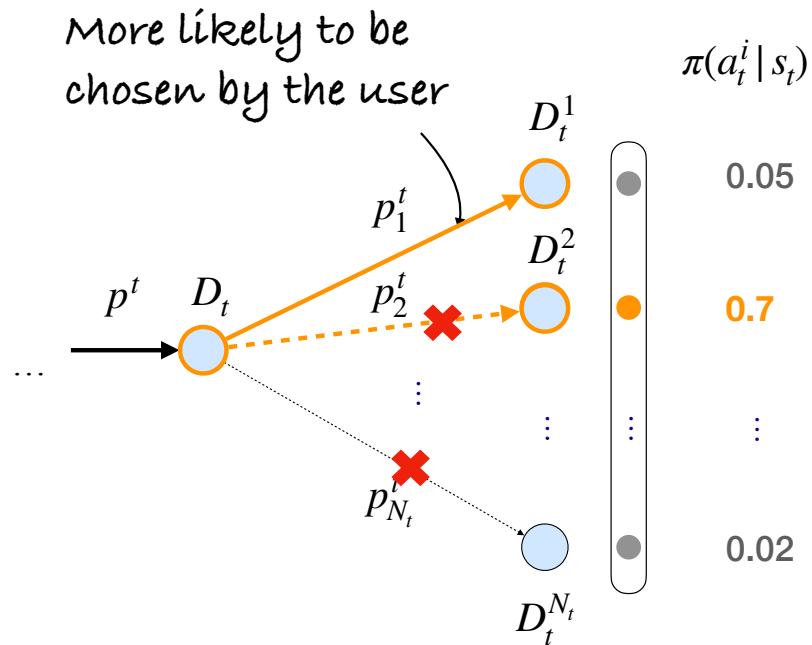
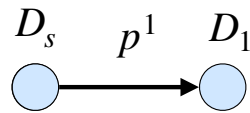
Goal: **learn from user feedbacks to adapt the rewards**



# Active Reward Learning

Learn2Clean  
+  
HIL

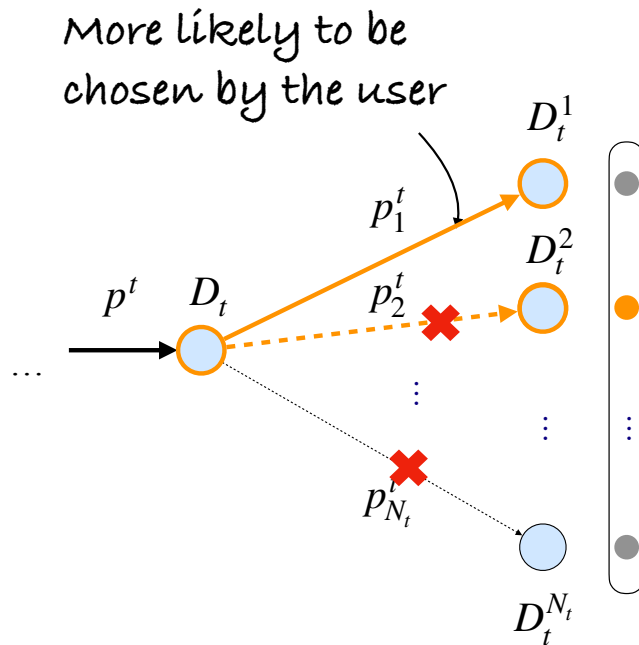
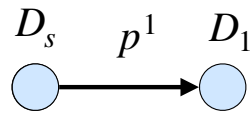
Goal: **learn from user feedbacks**  
**to adapt the rewards**



# Active Reward Learning

Learn2Clean  
+  
HIL

Goal: **learn from user feedbacks to adapt the rewards**

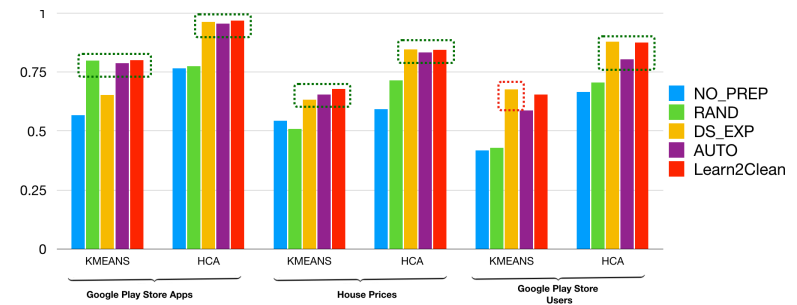
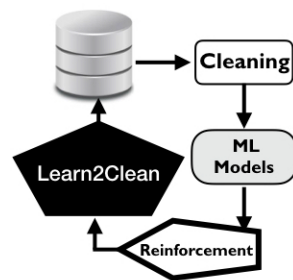


$\tilde{\pi}(a_t^i   s_t)$	$\pi(a_t^i   s_t)$
0.9	0.05
0	0.7
⋮	⋮
0	0.02

Force exploration

# Ongoing work

- New version of Learn2Clean with deep RL agents
- Combine AutoML, AutoCuration, and HIL
- Learn better reward functions
- Extend the library of ML and data preparation methods
- Extend experiments with more intricate data glitches and various glitch distributions



Code: <https://github.com/LaureBerti/Learn2Clean>

# Concluding Remarks

- ML crucially needs principled data curation and preparation, adequate tooling, and user assistance
- The impact of data preprocessing variability is largely underestimated in ML
- Many data preprocessing tasks require seamless integration of Human-in-the-Loop and automated ML-based solutions
- Perfect timing for many R&D opportunities:
  - Manage and orchestrate human/machine resources
  - Challenge and transfer research ideas to operational and very large-scale contexts



Thank you!