

# Event Log Data Quality in Process Mining

## Modelling, Assessing, Improving

**Marco Comuzzi**

Department of Industrial Engineering  
Ulsan National Institute of Science and Technology (UNIST)  
[mcomuzzi@unist.ac.kr](mailto:mcomuzzi@unist.ac.kr)

# About me

CS Undergraduate, master degree, Politecnico di Milano (밀라노 공대)

PhD Information Technology, Politecnico di Milano (밀라노 공대), 2007

2009 – 2013 Assistant Professor, Eindhoven University of Technology (네델란드)

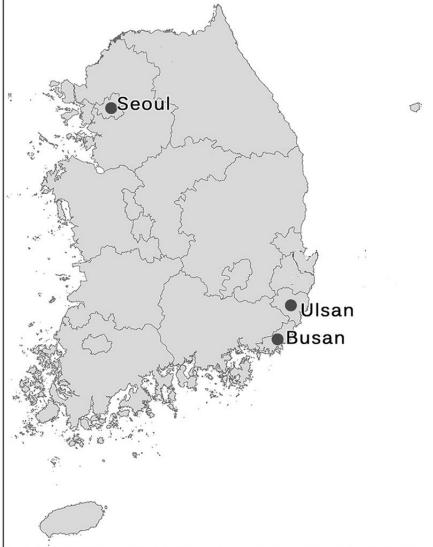
2013 – 2016 Assistant Professor, City, University of London (영국)

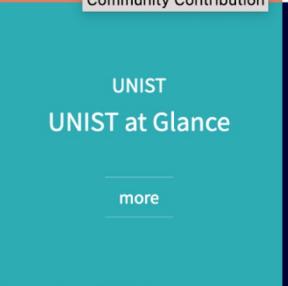
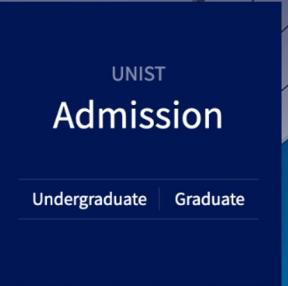
2016 – Associate Professor, Department of Industrial Engineering, UNIST

## Research interests

Process mining

Blockchain



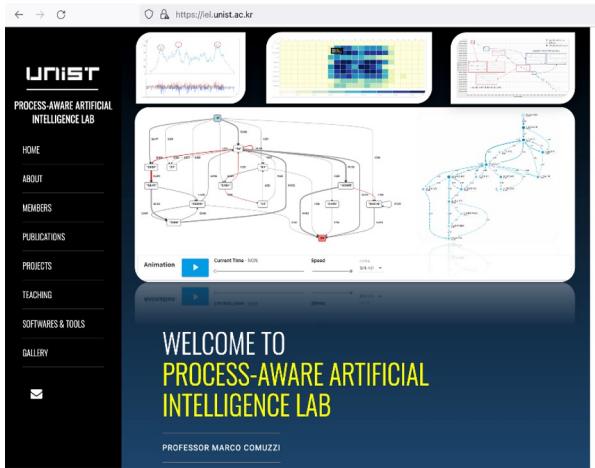




연구  
연구활동

연구활동

연구실소개



UNIST  
PROCESS-AWARE ARTIFICIAL  
INTELLIGENCE LAB

HOME  
ABOUT  
MEMBERS  
PUBLICATIONS  
PROJECTS  
TEACHING  
SOFTWARES & TOOLS  
GALLERY

WELCOME TO  
PROCESS-AWARE ARTIFICIAL  
INTELLIGENCE LAB

PROFESSOR MARCO COMUZZI

Animations Current time: 10:00 Speed: 100%

INTRODUCTION OF LAB (SINCE 2016 - )

This lab was established in 2016. It focuses on the application of machine learning techniques (mainly classification) to the analysis of business process event logs. As a lab, we are always trying to contribute to the main international academic conferences in the field: Int. Conf. on Business Process Management, Int. Conf. on Process Mining and Int. Conf. on Advanced Information Systems Engineering. Our research has always practical relevance. For instance we have applied the results of our research on event logs obtained from the Ulsan Port Authority and from a large university hospital in Korea.

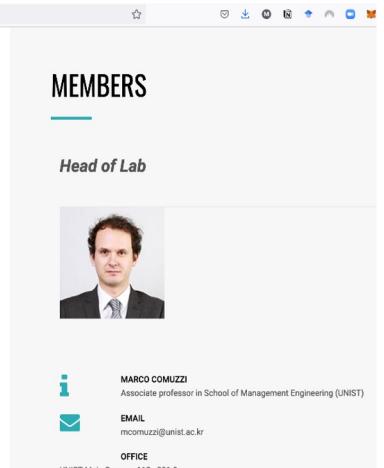


교수  
연구그룹  
전문분야

Comuzzi Marco

스마트서비스, 스마트 비즈니스 & 파이낸스

데이터마이닝, 기계학습 알고리즘(연합학습, 표현학습), 산업인공지능



ABOUT

PAAI Lab.

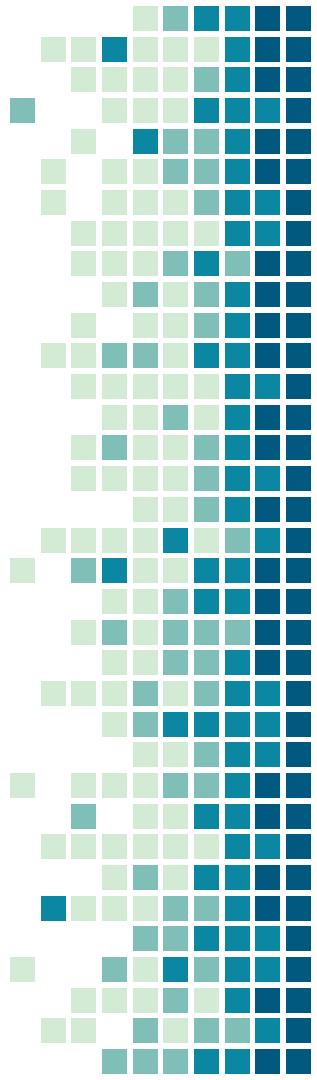
Office: UNIST Main Campus - 112 - Building 302-2 room

MEMBERS

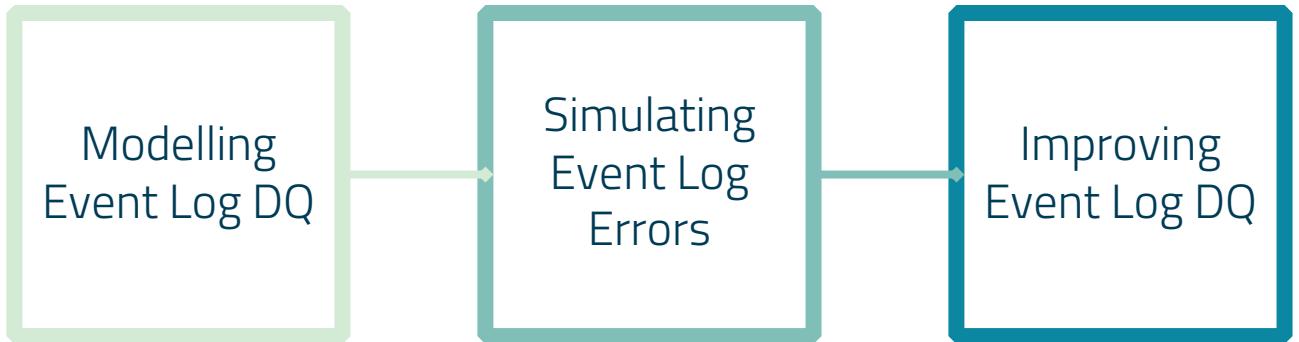
Head of Lab



MARCO COMUZZI  
Associate professor in School of Management Engineering (UNIST)  
EMAIL  
mcomuzzi@unist.ac.kr  
OFFICE

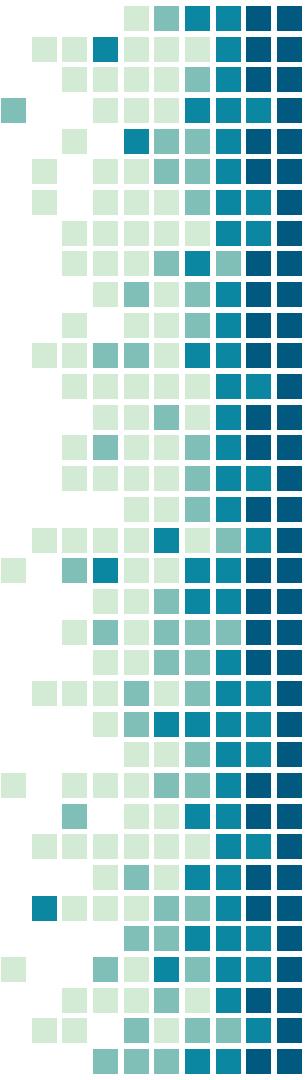


# Plan for today



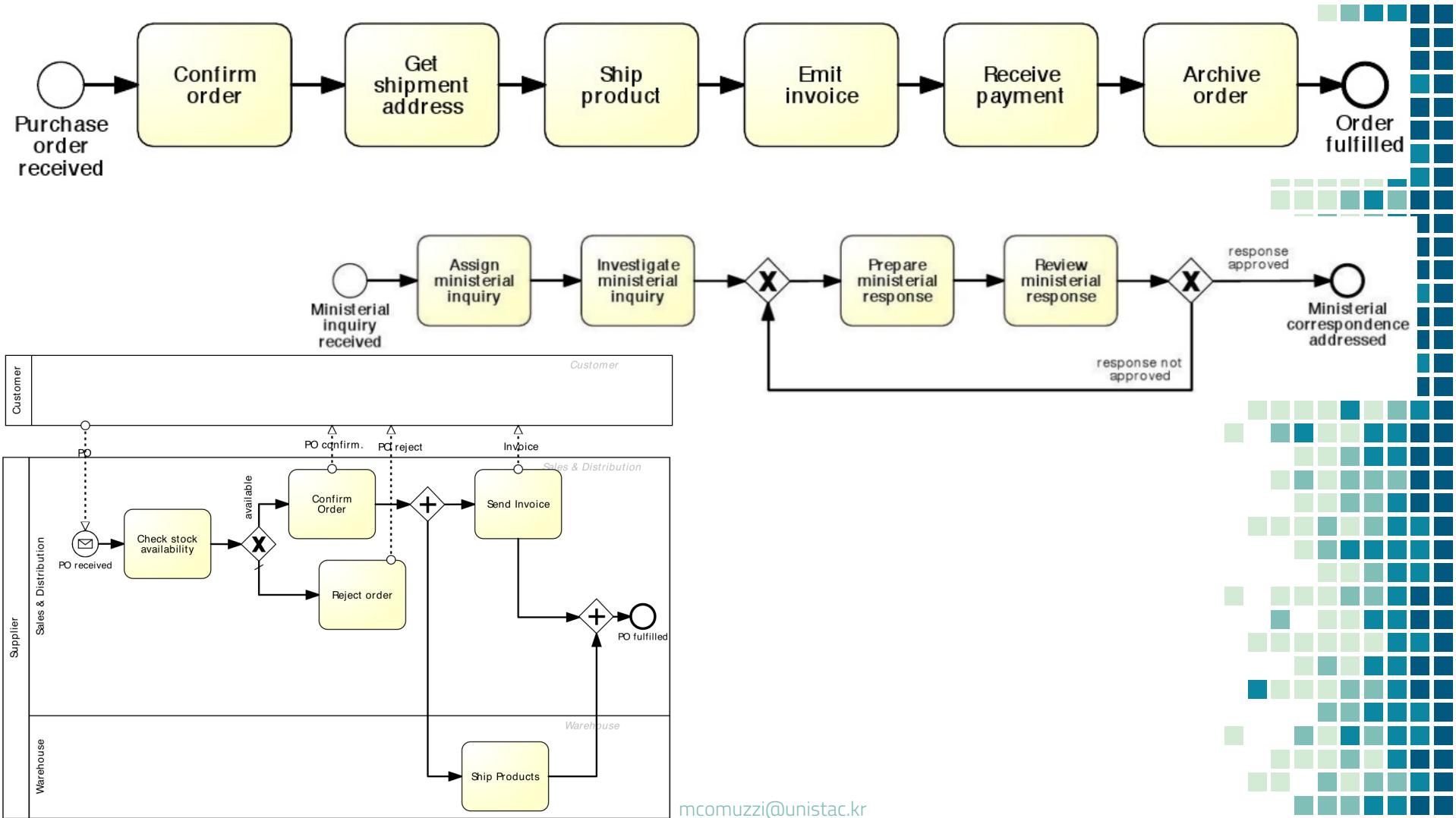
# Business Processes and Process Mining 101

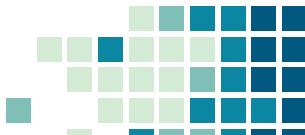




# What is a business process?

"A collection of related events, activities and decisions that involve a number of actors and resources, and that collectively lead to an outcome that is of value to an organisation and/or its customers"





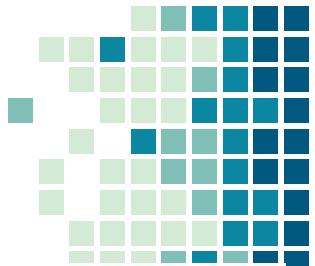
# Processes are everywhere



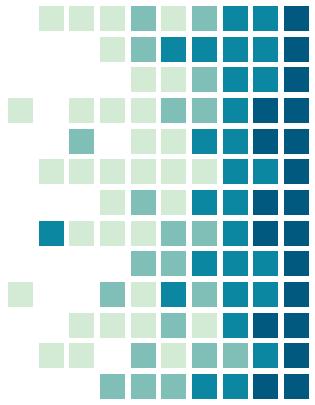
Service Administration Public

Role Madding  
Design Work Comparing  
Languages  
Comparing  
Value  
Growth  
Challenges  
Motivation  
Dilemma  
Separate  
Assume  
Reform  
Budget  
Processes  
Rationality  
Power  
Leadership  
Research  
Toward  
Participation  
Management  
Policy  
New  
Approach  
Analysis  
Trends  
Manufacturing  
Innovation  
Management  
Market  
Trend  
Advising  
Adhesive  
Choice  
Transfessional  
Organizations  
Delivery  
Systems  
Control  
Implementation  
Citizen  
Transformation  
Study  
Sector  
Theory  
Google  
Change  
Different





# A bridge between data and process science



# Nokia transforms business processes with Celonis



By **Derek du Preez** July 25, 2022



Audio mode



**SUMMARY:** Multinational communications company Nokia has also created Excellence for Celonis' Execution Management System, which allows the business to scale the platform.

## Process mining become the terrain

Enterprise technology companies like SAP and Microsoft have adopted process mining technology, but the current rate of acquisition is rapid

Home > Cloud > Why Process Mining Is the top skill to learn in 2022

Cloud News

## Why Process Mining Is the top skill to learn in 2022

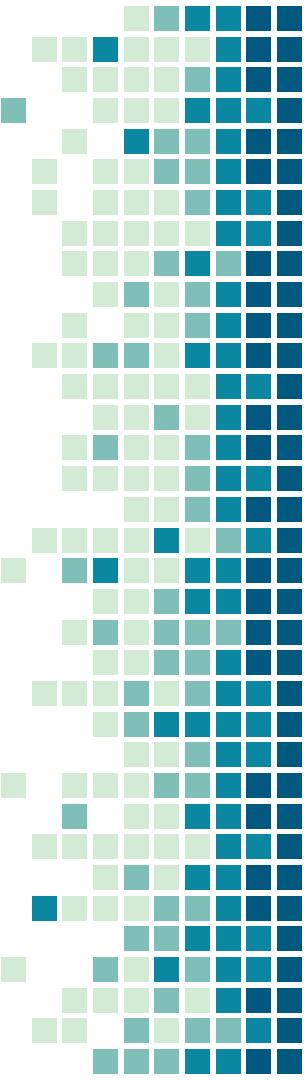
By **CRN Team** - July 22, 2022

2



Like 0





# What are “process data”?

Data logged by information systems supporting the process execution



# Event logs for process mining

A process consists of cases  
(or "instances")

A case consists of a "trace" of events

Each event relates precisely to one case

Events in a case are ordered in time

Events can have attributes:

Activity label

Timestamp

Cost

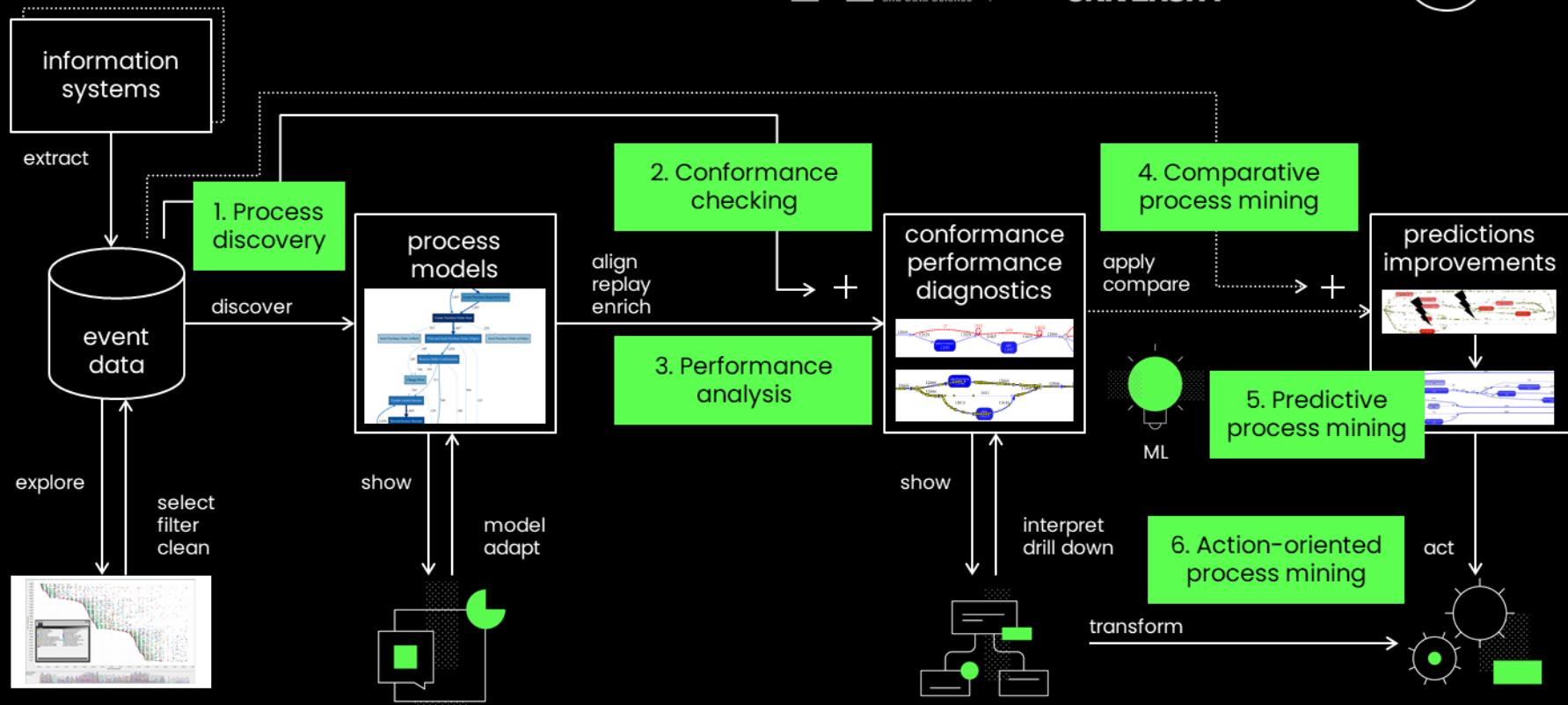
Resource

...

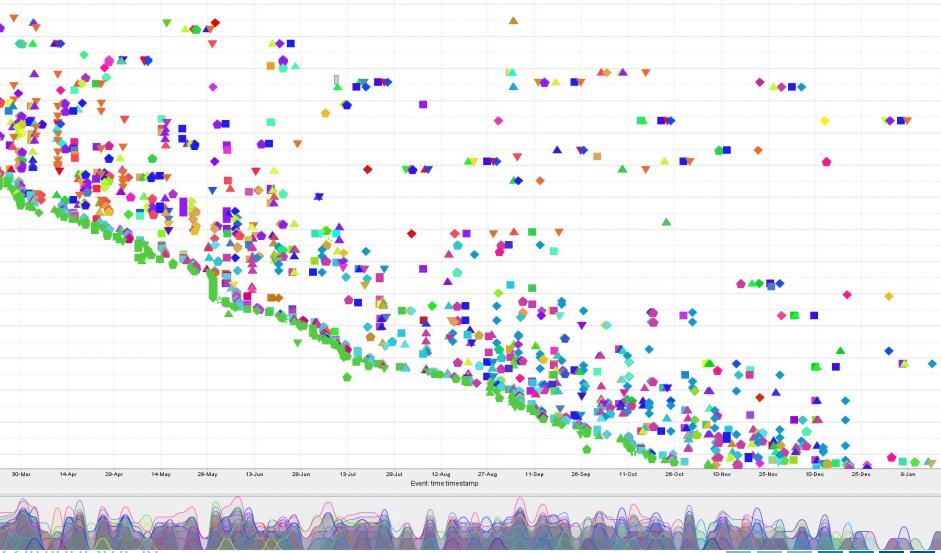
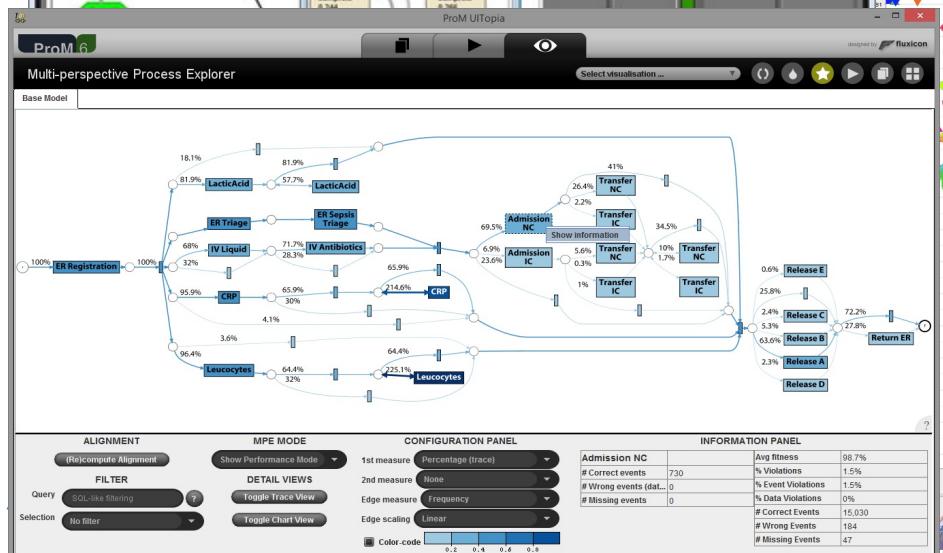
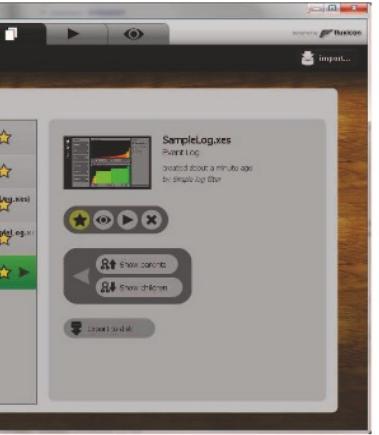
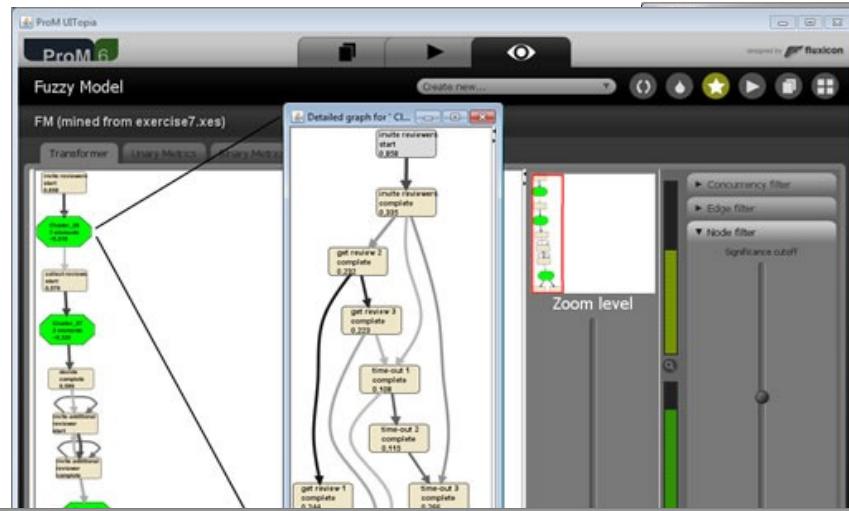
	case id	event id	properties				
			timestamp	activity	resource	cost	...
1	35654423	30-12-2010:11.02	register request	Pete	50	...	
	35654424	31-12-2010:10.06	examine thoroughly	Sue	400	...	
	35654425	05-01-2011:15.12	check ticket	Mike	100	...	
	35654426	06-01-2011:11.18	decide	Sara	200	...	
	35654427	07-01-2011:14.24	reject request	Pete	200	...	
2	35654483	30-12-2010:11.32	register request	Mike	50	...	
	35654485	30-12-2010:12.12	check ticket	Mike	100	...	
	35654487	30-12-2010:14.16	examine casually	Pete	400	...	
	35654488	05-01-2011:11.22	decide	Sara	200	...	
	35654489	08-01-2011:12.05	pay compensation	Ellen	200	...	
3	35654521	30-12-2010:14.32	register request	Pete	50	...	
	35654522	30-12-2010:15.06	examine casually	Mike	400	...	
	35654524	30-12-2010:16.34	check ticket	Ellen	100	...	
	35654525	06-01-2011:09.18	decide	Sara	200	...	
	35654526	06-01-2011:12.18	reinitiate request	Sara	200	...	
	35654527	06-01-2011:13.06	examine thoroughly	Sean	400	...	
	35654530	08-01-2011:11.43	check ticket	Pete	100	...	
	35654531	09-01-2011:09.55	decide	Sara	200	...	
	35654533	15-01-2011:10.45	pay compensation	Ellen	200	...	
4	35654641	06-01-2011:15.02	register request	Pete	50	...	
	35654643	07-01-2011:12.06	check ticket	Mike	100	...	
	35654644	08-01-2011:14.43	examine thoroughly	Sean	400	...	
	35654645	09-01-2011:12.02	decide	Sara	200	...	
	35654647	12-01-2011:15.44	reject request	Ellen	200	...	
...							



Case_ID	Event_ID	Activity_Name	CompleteTimeStamp	Resource_ID	Department	Reserved_Timestamp	Doctor_ID	System
P001	e_1042	Reserve a consultation	2020/10/27 10:25:21	R0002	Administration	2020/10/29 12:30:00		Sys_admin
	e_1049	Receipt treatment	2020/10/29 12:31:48	R0104	Neurology		D1023	Sys_neuro
	e_1050	Start treatment	2020/10/29 13:01:13	D1023	Neurology		D1023	Sys_neuro
	e_1051	Finish treatment	2020/10/29 13:14:09	D1023	Neurology		D1023	Sys_neuro
	e_1053	Start an examination	2020/10/29 13:29:44	R3211	Radiology			Sys_radi
	e_1054	Finish an examination	2020/10/29 14:02:12	R3211	Radiology			Sys_radi
	e_1057	Payment	2020/10/29 14:11:02	R2062	Administration			Sys_admin
	e_1062	Print a prescription	2020/10/29 14:31:52	R0912	Pharmacy			Sys_phar
P002	e_1045	Reserve a consultation	2020/10/27 14:11:32	R0004	Administration	2020/10/29 13:15:00		Sys_admin
	e_1052	Receipt treatment	2020/10/29 13:18:21	R0422	Cardiology		D1030	Sys_cardi
	e_1055	Start treatment	2020/10/29 14:02:42	D1030	Cardiology		D1030	Sys_cardi
	e_1056	Finish treatment	2020/10/29 14:11:01	D1030	Cardiology		D1030	Sys_cardi
	e_1061	Payment	2020/10/29 14:27:23	R0426	Administration			Sys_admin
	e_1071	Print a prescription	2020/10/29 14:47:07	R0912	Pharmacy			Sys_phar

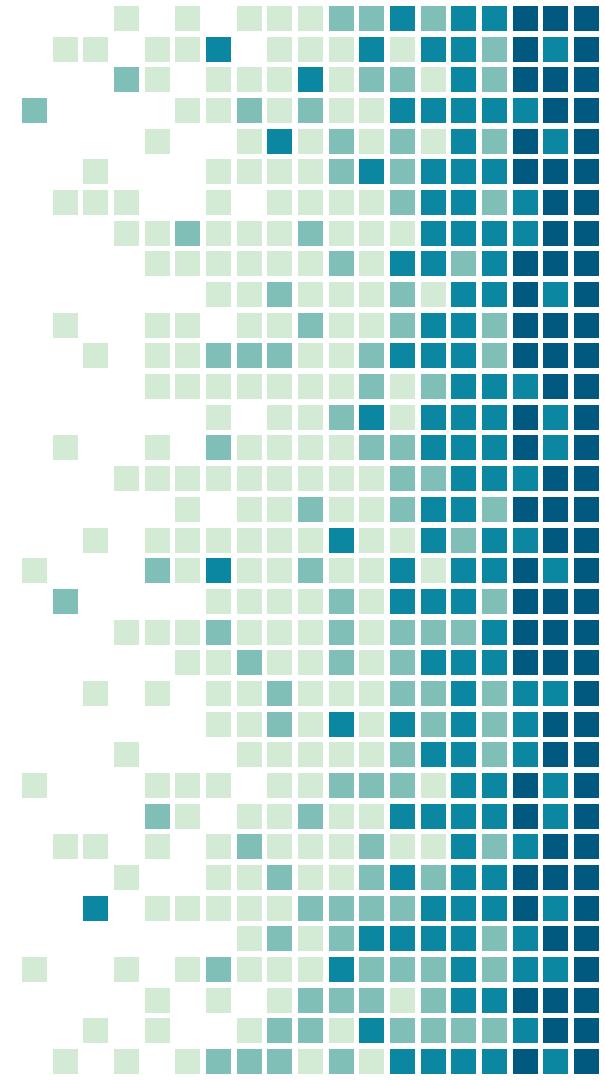


## Process Mining: From Theory to Execution



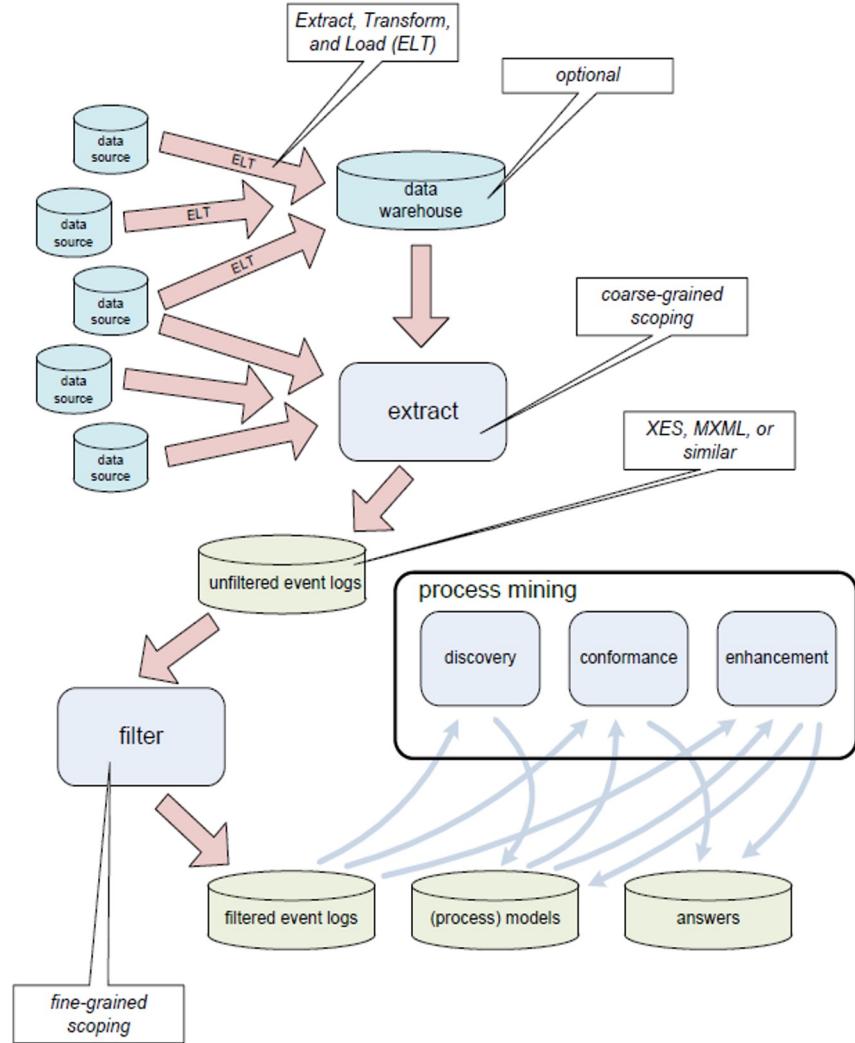
Ter Hofstede, A.H., Koschmider, A., Marrella, A., Andrews, R., Fischer, D.A., Sadeghianasl, S., Wynn, M.T., Comuzzi, M., De Weerdt, J., Goel, K. and Martin, N., 2023. Process-data quality: the true frontier of process mining. *ACM Journal of Data and Information Quality*, 15(3), pp.1-21.

# Event Log Data Quality



Any system used in a process  
can be a source of process data:

ERP, CRM, Sharepoint,  
spreadsheets,  
accounting/management  
software, IoT data streams, ...



# Event log preparation is critical, indeed

Ensuring event log has sufficient quality is crucial

Event log extraction and preparation can take up to 80% of the time in a process mining project

Lack of structured approaches to deal with this phase

# Data Quality

## Accuracy

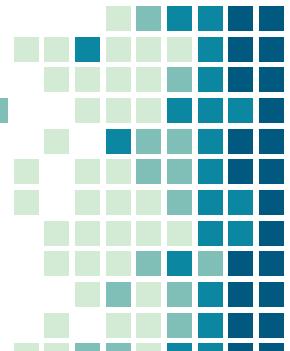
Degree of adherence of the data to reality

## Completeness

Capability of the data of representing all relevant aspects of reality

## Reliability

Capability of the data to comply with the properties of reality (integrity constraints and business rules)



# DQ of event logs

Bose et al. 2012

Top-down approach

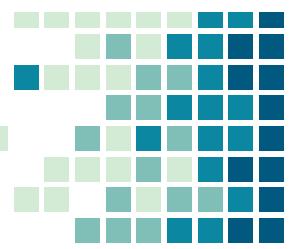
Missing Data

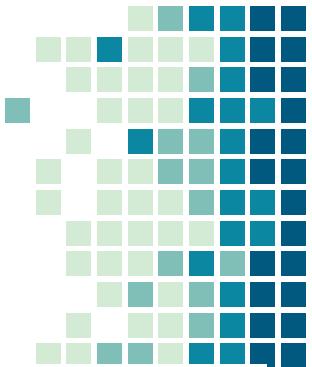
Incorrect Data

Imprecise Data

Irrelevant Data

	case	event	belongs to	c_attribute	position	activity name	timestamp	resource	e_attribute
Missing Data	I1	I2	I3	I4	I5	I6	I7	I8	I9
Incorrect Data	I10	I11	I12	I13	I14	I15	I16	I17	I18
Imprecise Data			I19	I20	I21	I22	I23	I24	I25
Irrelevant Data	I26	I27							





# Event Log Imperfection Patterns

A different, bottom-up approach

Derive patterns from existing process mining case studies

11 patterns identified

1. [Form-based Event Capture](#)
2. [Inadvertent Time Travel](#)
3. [Unanchored Event](#)
4. [Scattered Event](#)
5. [Elusive Case](#)
6. [Scattered Case](#)
7. [Collateral Events](#)
8. [Polluted Label](#)
9. [Distorted Label](#)
10. [Synonymous Labels](#)
11. [Homonymous Label](#)

# Form-based event capture

Episode ID	Event	Timestamp	Description	...
ID1	Primary Survey	2012-11-23 15:42:38	.....	....
ID1	Airway Clear	2012-11-23 15:42:38	.....	....
	...	2012-11-23 15:42:38	.....	..
	Primary Survey	2012-11-24 09:58:33	.....	..
ID2	Airway Clear	2012-11-24 09:58:33	.....	..
	...	2012-11-24 09:58:33	.....	....
	Procedure 1	2012-11-24 09:58:33	Completed on 2012-11-24 06:58:34	....

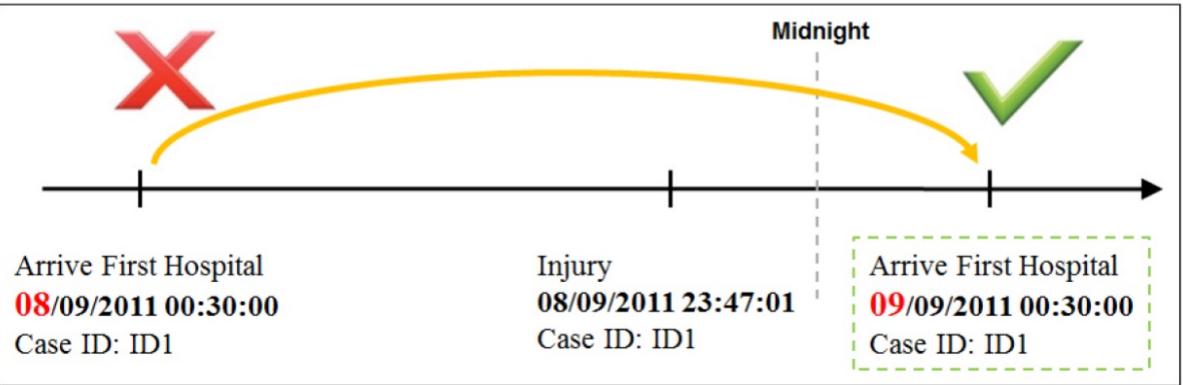
These events are recorded on a form ...

... and all have the same timestamp.

# Inadvertent time travel

Episode ID	Activity	Timestamp	...
ID1	Arrival first hospital	2011-09-08 00:30:00	
ID1	Injury	2011-09-08 23:47:01	.....
...	.....	.....	.....
ID1	Operation	2011-09-09 16:30:00	.....

'Midnight' problem.  
Time portion correct  
but date part in  
error.



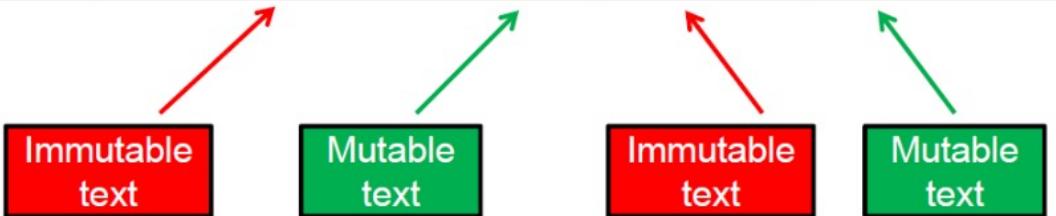
# Collateral events

caseID	Activity	Timestamp
1234567	Adjust recovery cost	19/06/2014 12:15:18
1234567	Adjust recovery cost	19/06/2014 12:16:53
1234567	Email	19/06/2014 12:19:25
....	....	.....
1234567	Pay assessor fee	19/06/2014 12:20:00
1234567	Adjust admin cost	19/06/2014 12:22:48

All events refer to  
single process step  
'Pay Insurance  
Claim Assessor'.

# Polluted labels

caseID	activity	timestamp	.....
xxxx	Notification of Loss - AAAA Incident No. aaaa	xxxx-xx-xx xx:xx:xx	.....
xxxx	Notification of Loss - BBBB Incident No. bbbb	yyyy-yy-yy yy:yy:yy	.....
xxxx	Notification of Loss - CCCC Incident No. cccc	zzzz-zz-zz zz:zz:zz	.....
.....	<u>Notification of Loss - DDDD Incident No. dddd</u>	.....	.....



# Distorted labels

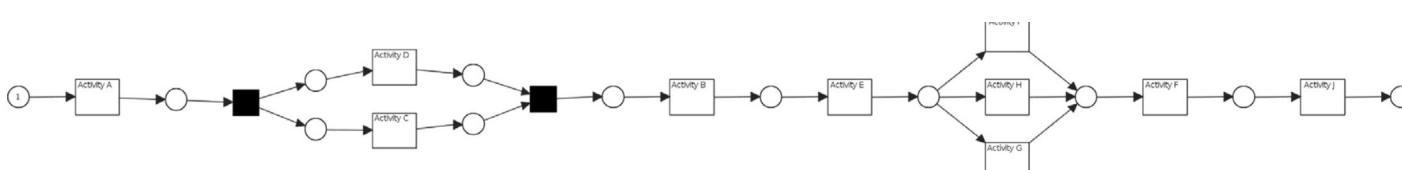
caseID	activity	timestamp	Description
1234567	a/w inv to cls.	06/09/2013 12:33:17	.....
8912345	a/w inv to cls	06/09/2013 13:10:23	.....
1234567	XX – Further Information Required	06/09/2013 13:15:00	.....
8912345	XX – Further Infomation Required	13/09/2013 07:24:36	.....

# Homonymous labels

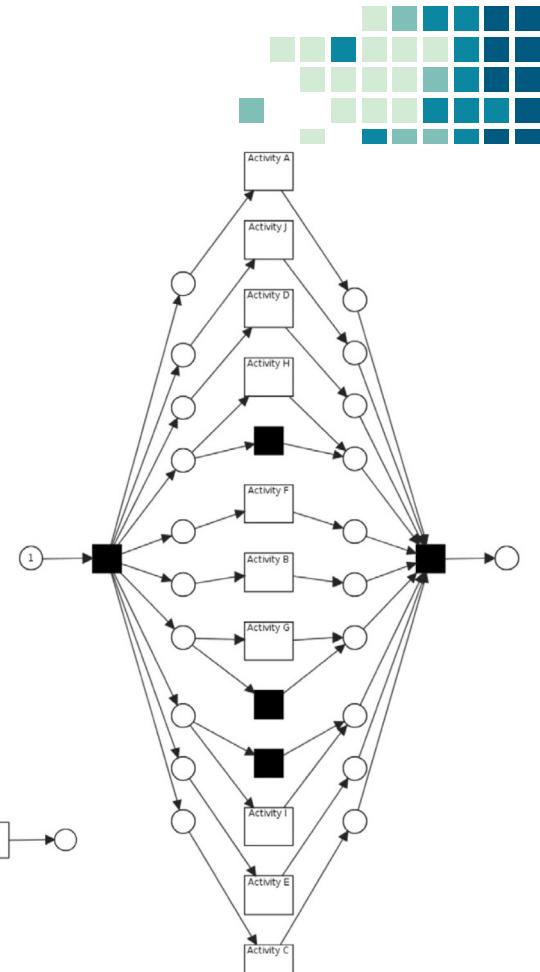
caseID	activity	timestamp	Description
1234567	Triage Assessment	06/09/2013 12:33:17	.....
1234567	Progress Note	06/09/2013 13:10:23	.....
1234567	Discharged	06/09/2013 13:15:00	.....
1234567	Triage Assessment	13/09/2013 07:24:36	.....
1234567	Triage Assessment	13/09/2013 07:28:51	.....

# Effects of poor DQ of event logs

Process model discovered from a synthetic log perturbed with 10% “wrong” activity labels and timestamps



(a) Process model mined from original (clean) event log



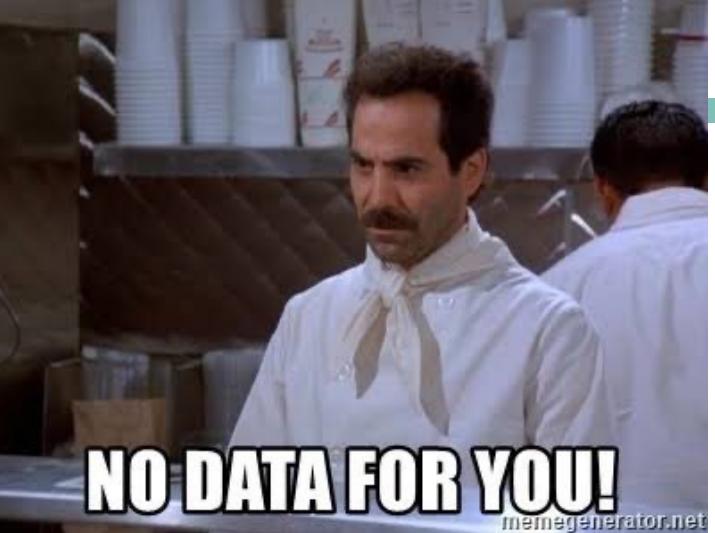
(b) Process model mined from log with injected anomalous values

# Effects of poor DQ of event logs

Evident to everybody in practice...

...but rarely researched in academia, why?

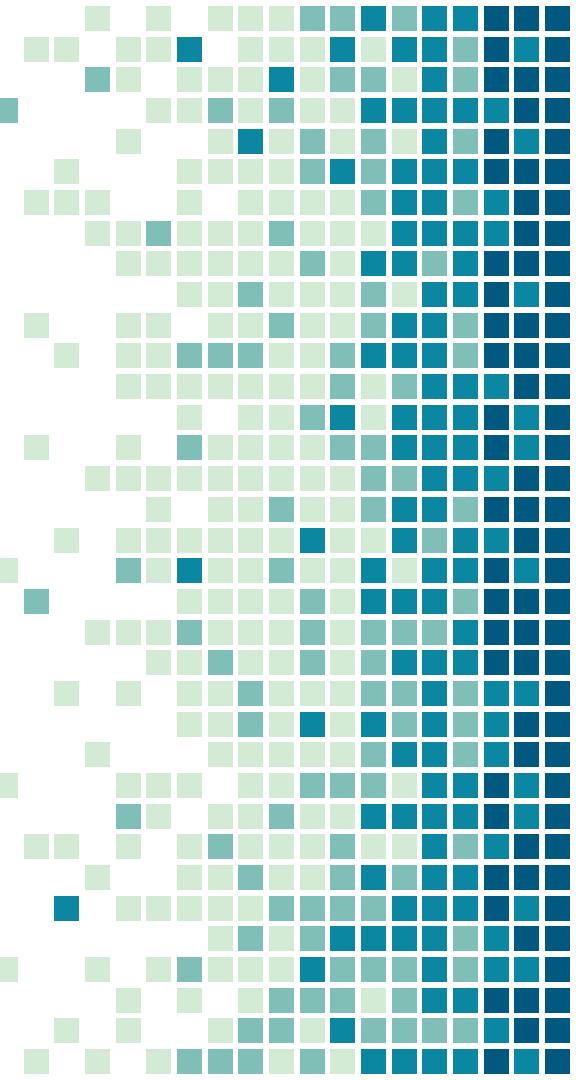
No (labeled) data!



Some imperfection patterns can be easily simulated, e.g. synonym labels, distorted labels

A general, extensible approach to generate data for testing event log cleaning tools is missing

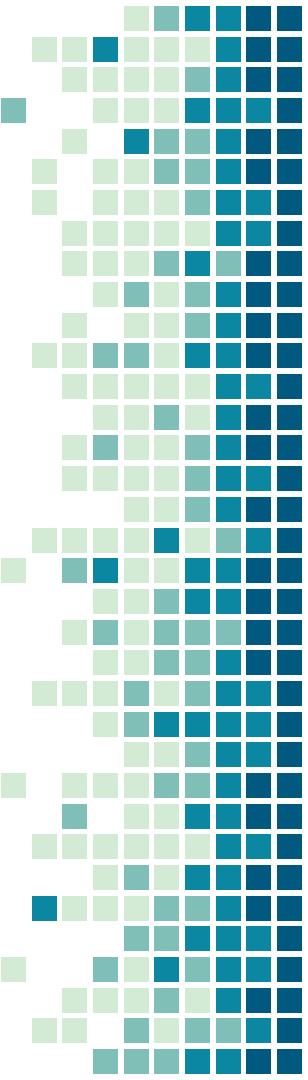
# FLAWD: a language for simulating event log errors

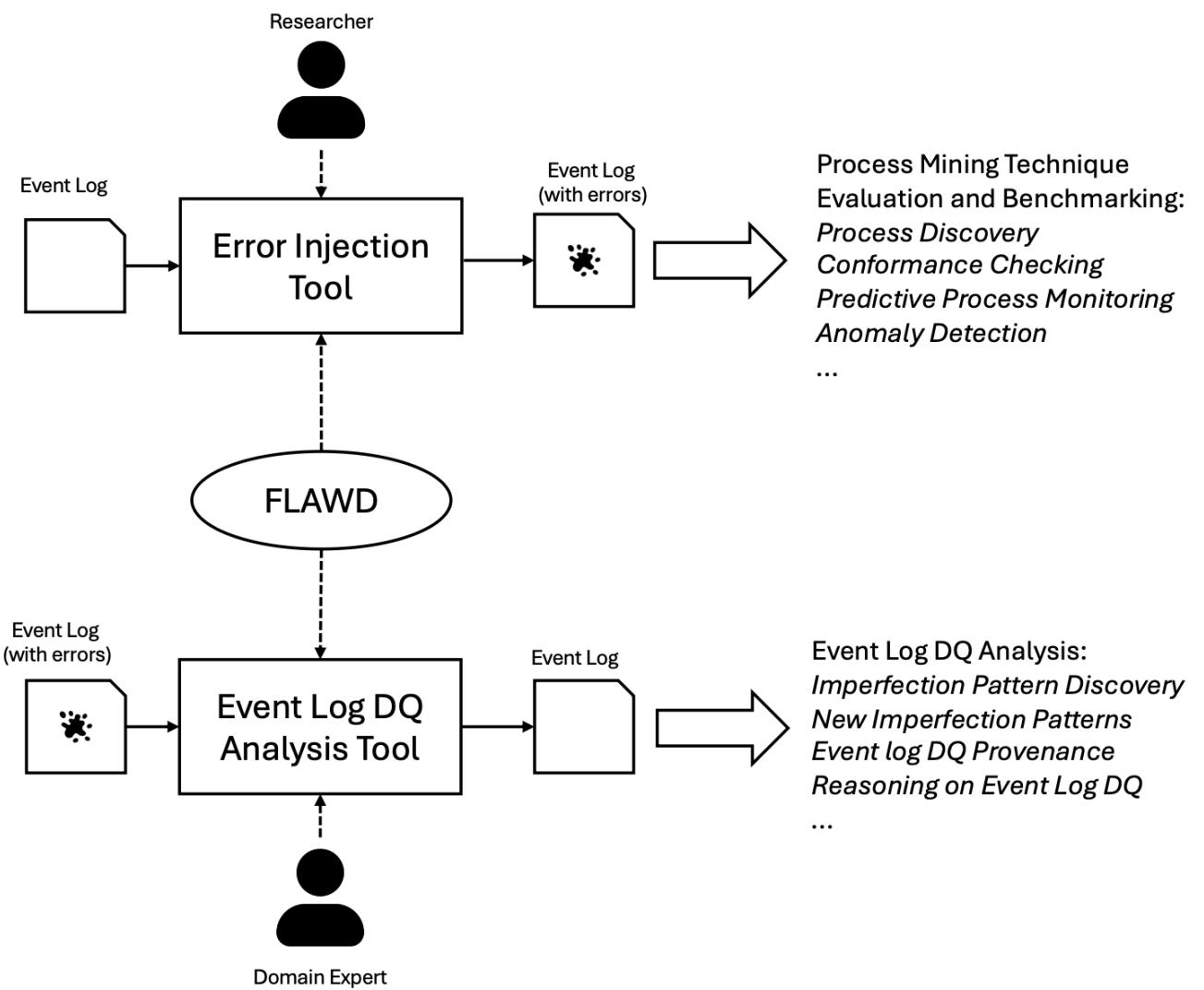


# Introducing FLAWD

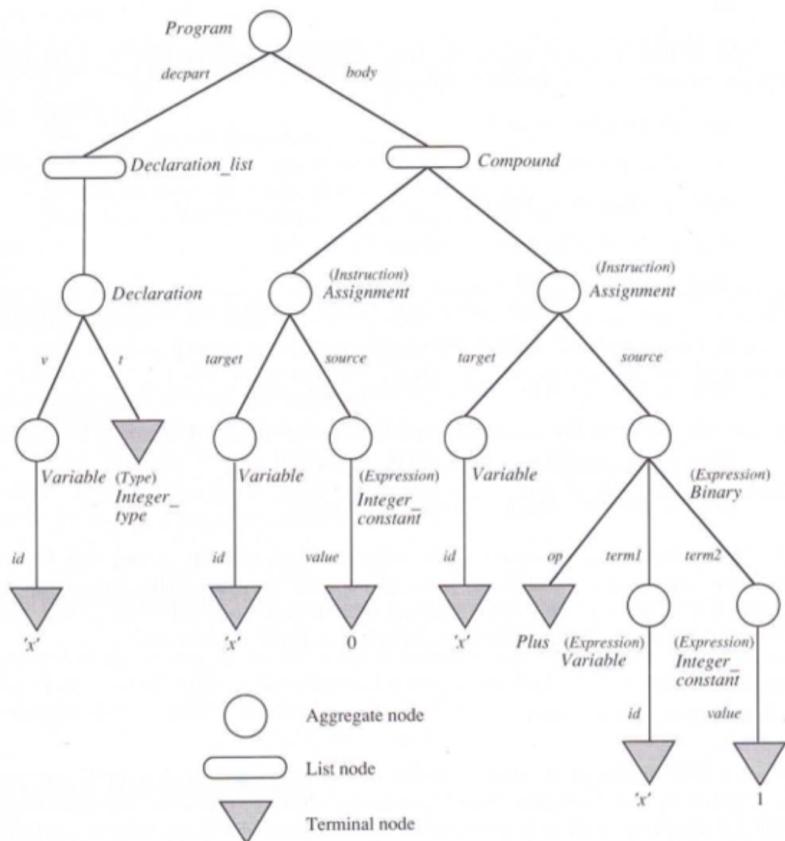
A language to model event log imperfection patterns

An implementation of the language as an event log error injection tool





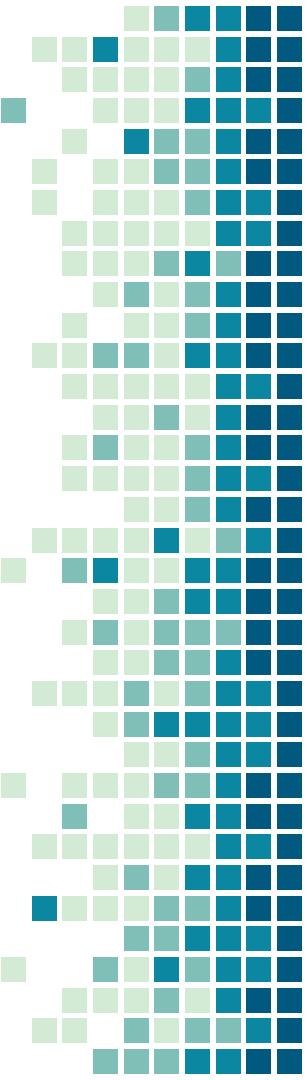
# How to describe FLAWD?



Bertrand Meyer  
Introduction  
to the Theory  
of Programming  
Languages

PRENTICE HALL  
INTERNATIONAL  
SERIES IN  
COMPUTER  
SCIENCE

C.A.R. HOARE SERIES EDITOR



# Why an abstract syntax grammar?

Allows to focus on the elements of a language, instead of the external appearance of programs written in that language

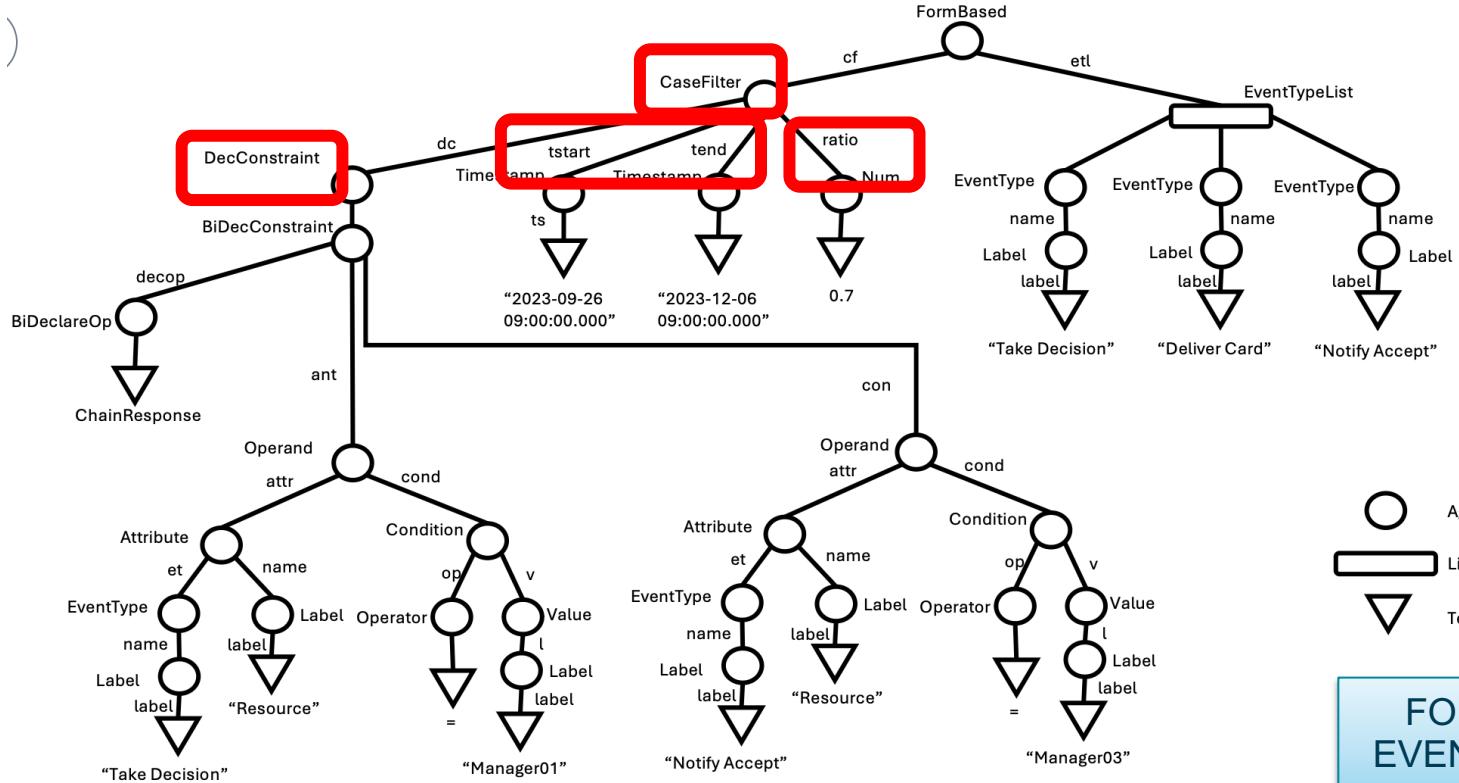
Facilitates understanding of properties, expressiveness, and concrete syntax generation

Case ID	Activity	Resource	Timestamp
1...	...	...	...
1 Take decision	Manager01	2023-09-27 11:56:18	
1 Deliver card	Manager01	2023-09-28 15:56:20	
1 Notify accept	Manager03	2023-09-29 09:56:20	
1...	...	...	...

Event log with injected errors

Case ID	Activity	Resource	Timestamp
1...	...	...	...
1 Take decision	Manager01	2023-09-27 11:56:18	
1 Deliver card	Manager01	2023-09-27 11:56:18	
1 Notify accept	Manager03	2023-09-27 11:56:18	
1...	...	...	...

)

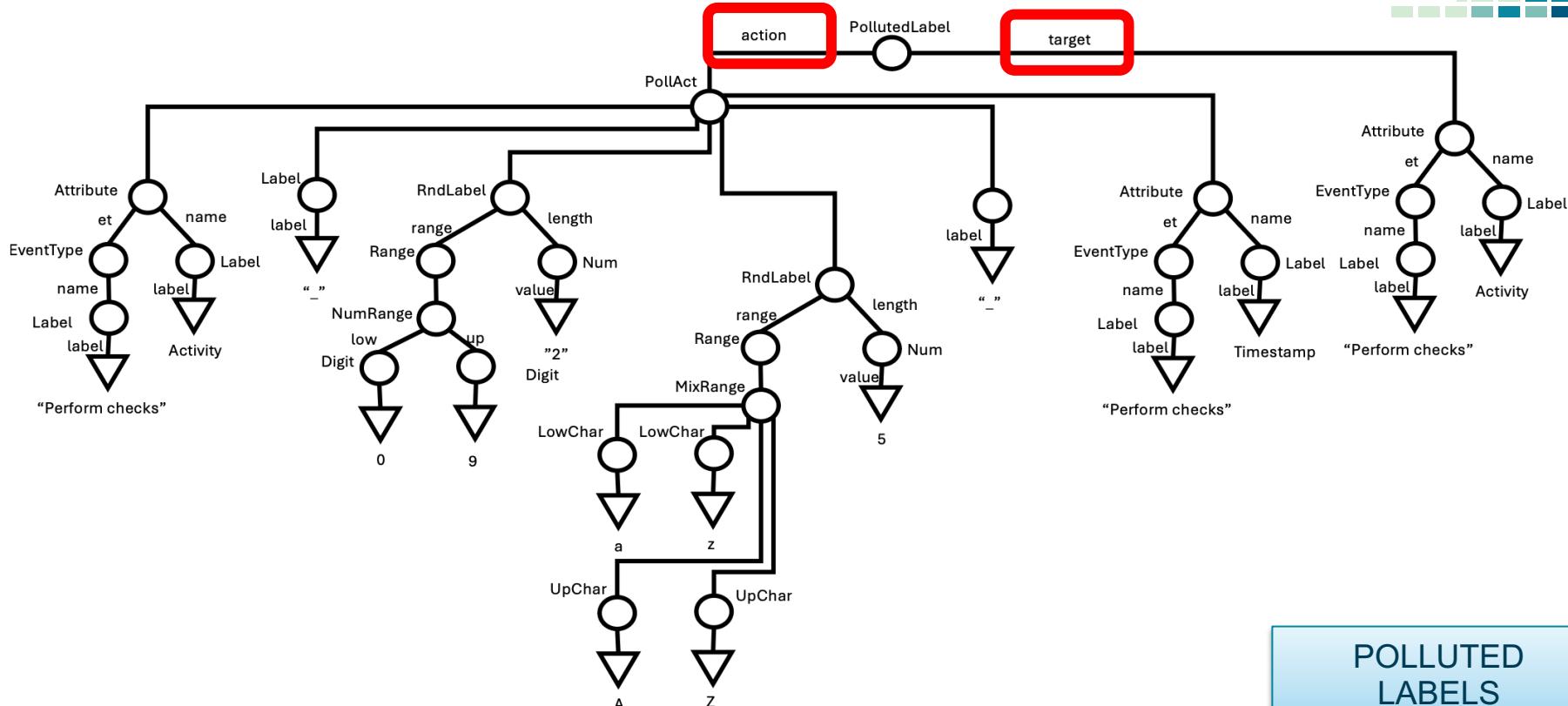


Event log

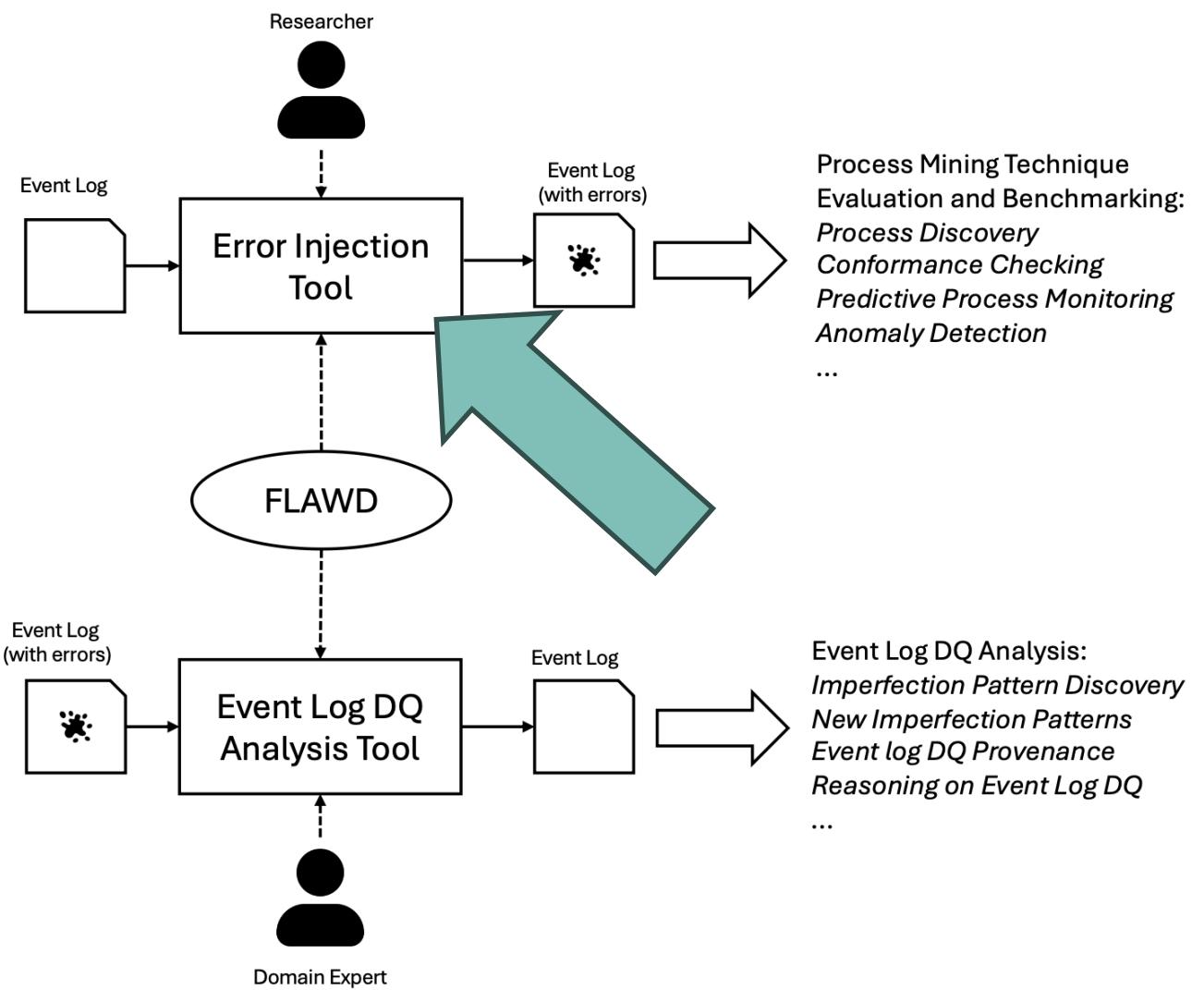
Case ID	Activity	Resource	Timestamp
1	New application received	Clerk01	2011-10-03 14:56:18
1	Perform checks	Clerk01	2011-10-04 11:56:18
1	Take decision	Manager01	2011-10-05 15:56:20
1	Notify accept	Manager02	2011-10-06 09:56:20
1	Deliver Card	Clerk03	2011-10-06 09:58:20

Event log with injected errors

Case ID	Activity	Resource	Timestamp
1	New application received	Clerk01	2011-10-03 14:56:18
1	Perform checks 91HTfsG_20111004 115618	Clerk01	2011-10-04 11:56:18
1	Take decision	Manager01	2011-10-05 15:56:20
1	Notify accept	Manager02	2011-10-06 09:56:20
1	Deliver Card	Clerk03	2011-10-06 09:58:20



POLLUTED  
LABELS



# Implementation

Python scripts

PM4Py for event log import/export

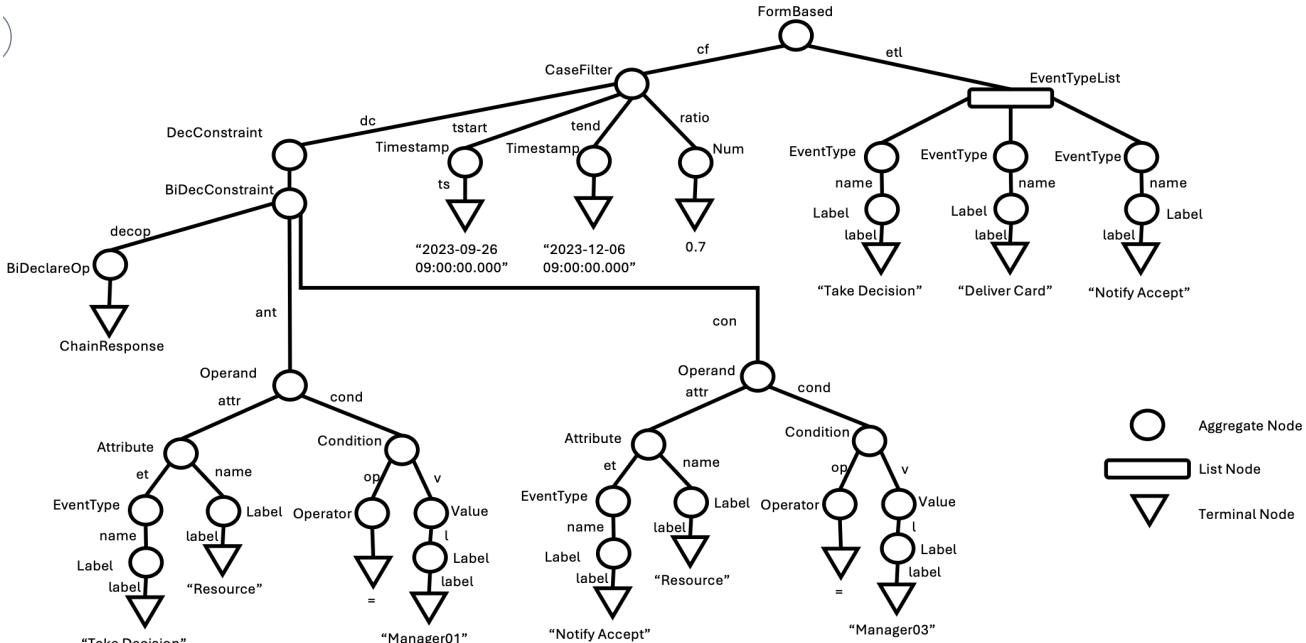
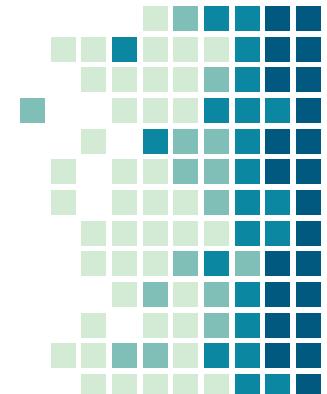
Declare4Py for case filtering based on Declare constraints

Generation of meaningful labels

```

el_polluted = form_based(el,
    etl = "['Take\u201cDecision', 'Deliver\u201cCard', 'Notify\u201cAccept']",
    ratio = 0.7,
    time_start = "2023-09-26\u09:00:00.000",
    time_end = "2023-12-06\u09:00:00.000",
    decConstraint = "ChainResponse[Take\u201cDecision,\u201cNotify\u201cAccept]\u201cA.
        \u2192 Resource\u201cManager01\u201cT.\u201cResource\u201cManager03']"
)

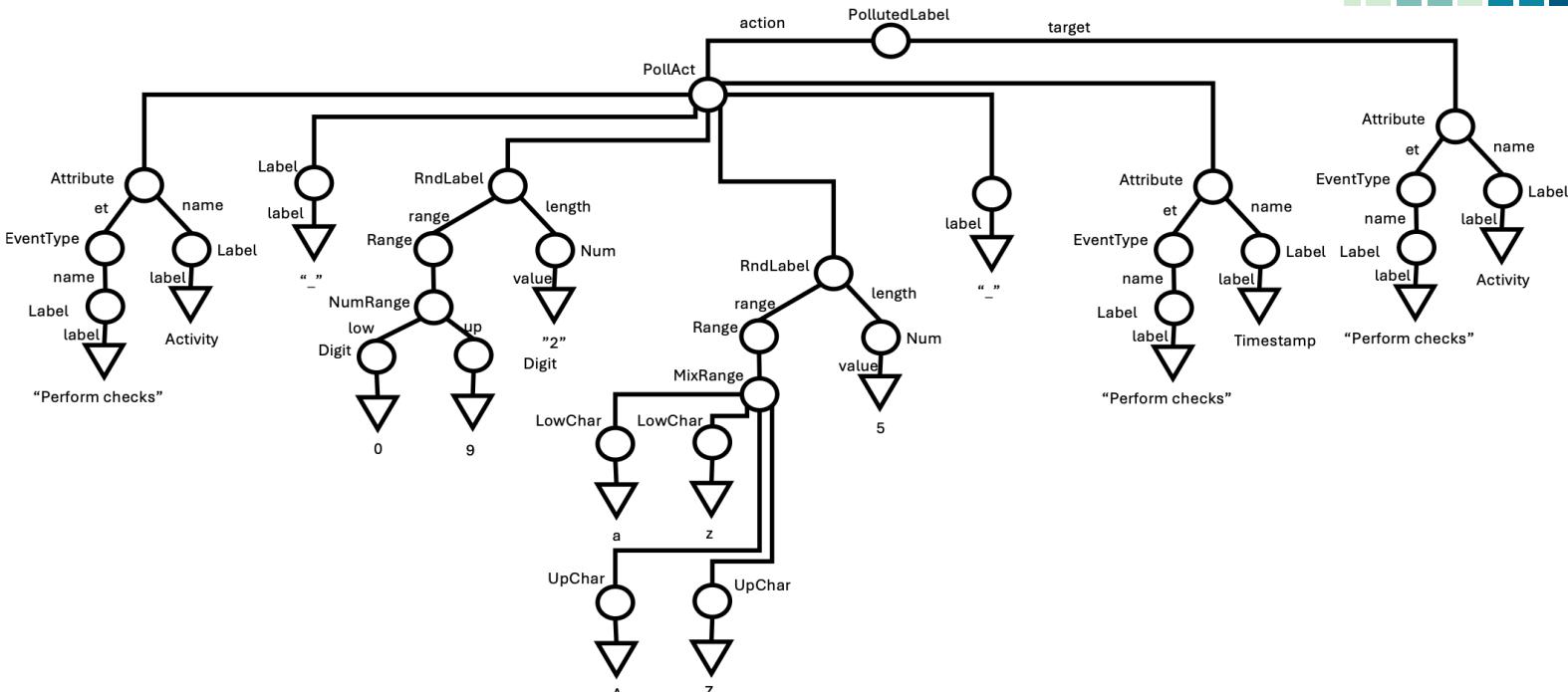
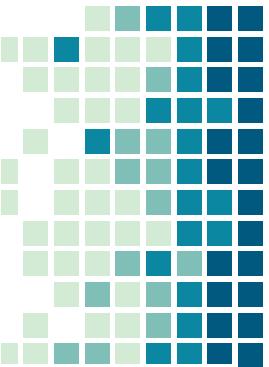
```



```

el_polluted = polluted_label(el,
    target = "[Activity:(Perform\u2022checks', 'Take\u2022decision')]"
    action = "[Activity]_[0-9:{2}][a-zA-Z:{5}]_[Timestamp*(%Y%m%d\u2022%H%M%
        ↪ S%f)]",
    # ratio = ...
    # filters = ...

```



# Generated labels

Trace the “provenance” of error injection

Case	Activity	Timestamp	Resource	label
1113	Check for completeness	2023-12-04 07:02:40.068	Clerk-000005	
1113	Perform checks	2023-12-04 07:14:40.451	Clerk-000001	
1113	Make decision	2023-12-04 07:15:07.899	Manager-000001	form-based events(2023-12-04 07:15:07.899000)
1113	Notify accept	2023-12-04 07:15:07.899	Manager-000003	form-based events(2023-12-04 08:02:31.947000)
1113	Deliver card	2023-12-04 07:15:07.899	Manager-000002	form-based events(2023-12-04 08:08:22.361000)

# Integration into the Praeclarus Toolkit

The screenshot shows the Praeclarus Process Data Quality Framework (PDQ) interface running in a browser window at localhost:8080. The main area displays a process flow diagram and an event log table.

**Process Flow Diagram:**

```
graph LR; CSVReader[CSV Reader<br/>63980 rows] --> ImperfectionInjector[Imperfection Injector<br/>620 rows]
```

**Event Log Table:**

Case	Activity	Timestamp	Resource	label
C0	into received	2023-12-04 07:02:30.068		
23	Check for completeness	2023-12-04 07:02:40.068	Clerk-000005	
24	Perform checks	2023-12-04 07:14:40.451	Clerk-000001	
25	Make decision	2023-12-04 07:15:07.899	Manager-000001	form-based events(2023-12-04...)
26	Notify accept	2023-12-04 07:15:07.899	Manager-000003	form-based events(2023-12-04...)
27	Deliver card	2023-12-04 07:15:07.899	Manager-000002	form-based events(2023-12-04...)
28	EVENT 13 END	2023-12-04 08:12:16.164		
29	Check for completeness	2023-12-06 03:16:56.084	Clerk-000003	
30	New online application received	2023-12-06 03:16:56.084		
31				

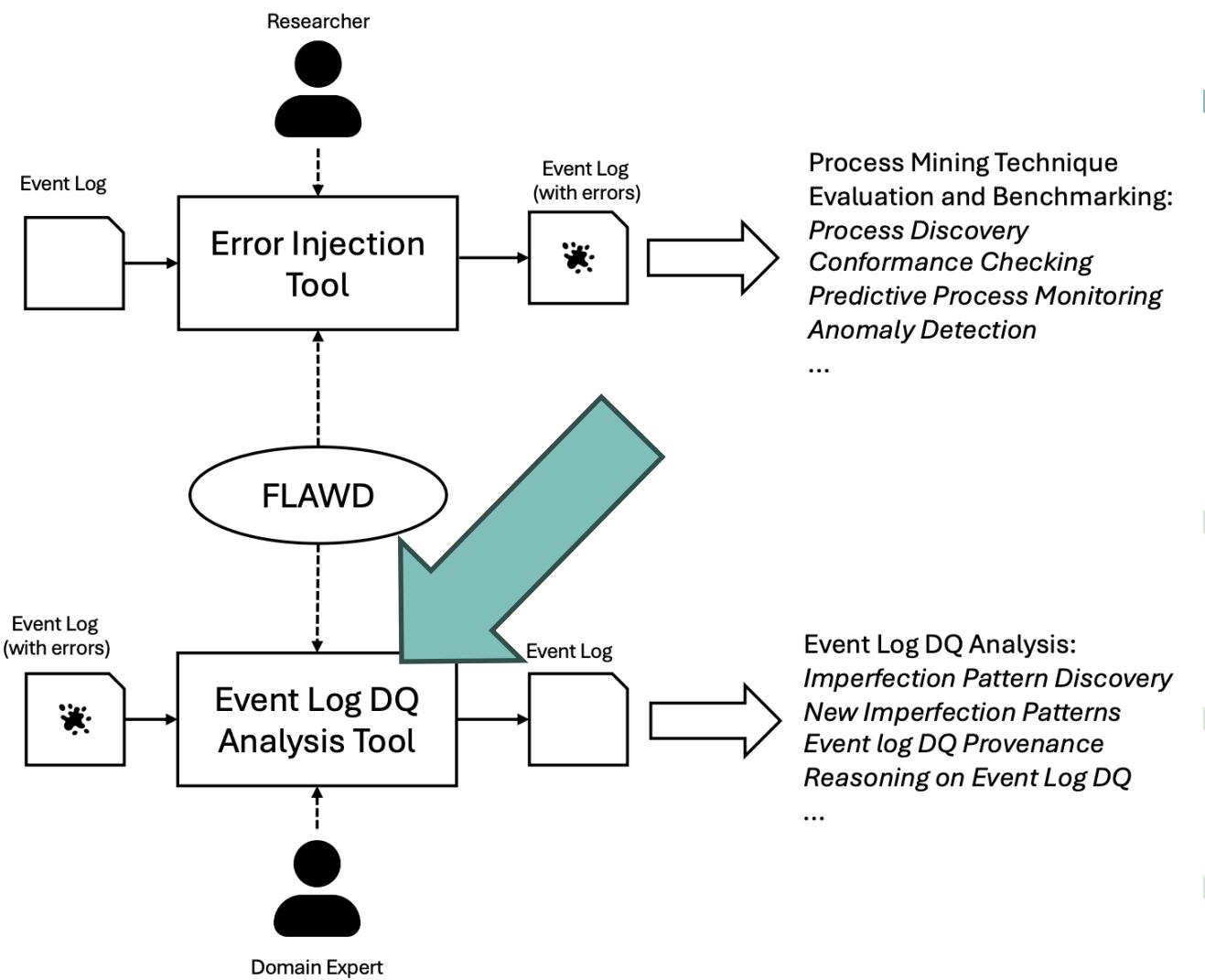
**Left Sidebar:**

- Plugins:
  - Imperfection Injector
  - Leverage M1
  - Leverage M1 (User Thresh...)
  - > Distorted Label
  - > Polluted Label
  - > Synonymous Labels
- Parameters:
  - Declare: none
  - Output File Name: none
  - Ratio: 0.85
  - Target: none
  - Time end: none
  - Time start: none

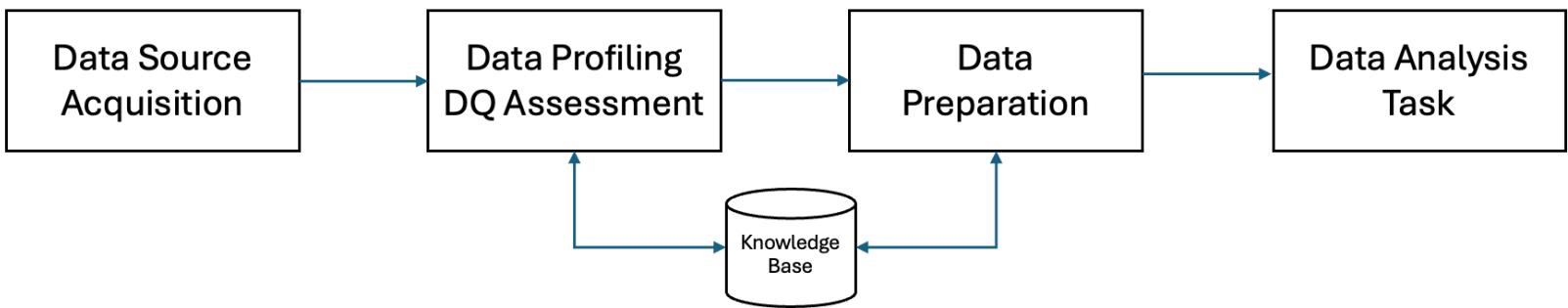
**Bottom Navigation:**

4! Output Log Diff Detected Events

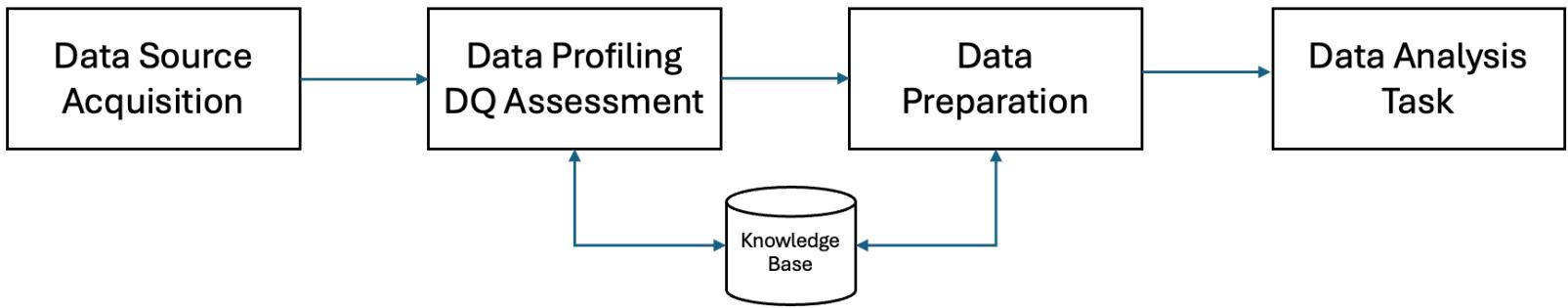
# Improving Event Log Data Quality



# Data preparation pipelines in data science



# Data preparation pipelines in process mining

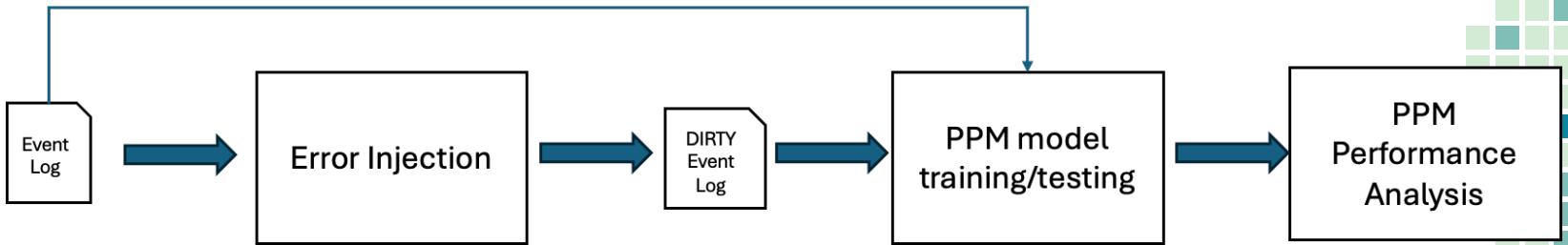


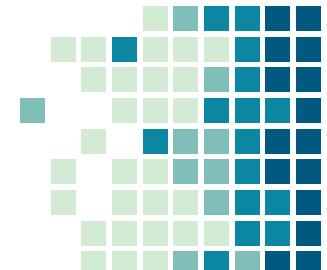
Data Source: Event Log

Data Preparation: Methods to Improve Event Log Quality

Data Analysis Task: Process Discovery, Conformance Checking, Predictive Process Monitoring (PPM)

# Step 1: Empirical analysis for building knowledge base





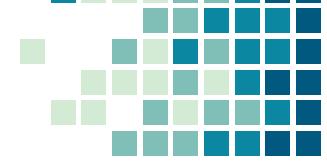
What is business process predictive monitoring?

Event log is a dataset

Can be used to train machine learning predictive models...

...which can then be used to predict “the future” of (running) process cases

case id	event id	properties				
		timestamp	activity	resource	cost	...
1	35654423	30-12-2010:11.02	register request	Pete	50	...
	35654424	31-12-2010:10.06	examine thoroughly	Sue	400	...
	35654425	05-01-2011:15.12	check ticket	Mike	100	...
	35654426	06-01-2011:11.18	decide	Sara	200	...
2	35654427	07-01-2011:14.24	reject request	Pete	200	...
	35654483	30-12-2010:11.32	register request	Mike	50	...
	35654485	30-12-2010:12.12	check ticket	Mike	100	...
	35654487	30-12-2010:14.16	examine casually	Pete	400	...
	35654488	05-01-2011:11.22	decide	Sara	200	...
	35654489	08-01-2011:12.05	pay compensation	Ellen	200	...
3	35654521	30-12-2010:14.32	register request	Pete	50	...
	35654522	30-12-2010:15.06	examine casually	Mike	400	...
	35654524	30-12-2010:16.34	check ticket	Ellen	100	...
	35654525	06-01-2011:09.18	decide	Sara	200	...
	35654526	06-01-2011:12.18	reinitiate request	Sara	200	...
	35654527	06-01-2011:13.06	examine thoroughly	Sean	400	...
	35654530	08-01-2011:11.43	check ticket	Pete	100	...
	35654531	09-01-2011:09.55	decide	Sara	200	...
	35654533	15-01-2011:10.45	pay compensation	Ellen	200	...
4	35654641	06-01-2011:15.02	register request	Pete	50	...
	35654643	07-01-2011:12.06	check ticket	Mike	100	...
	35654644	08-01-2011:14.43	examine thoroughly	Sean	400	...
	35654645	09-01-2011:12.02	decide	Sara	200	...
	35654647	12-01-2011:15.44	reject request	Ellen	200	...
...	...	...	...	...	...	...



What can we predict about a process case?

“Which will be the next activity that will be executed?”

Predicting activities that will be executed

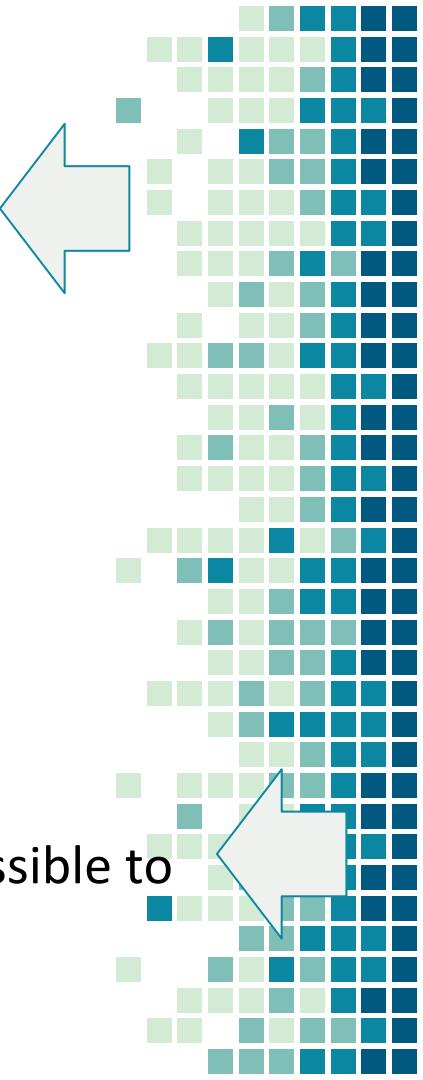
“When will the next activity be executed?

What will be the timestamp of the last activity of this case?”

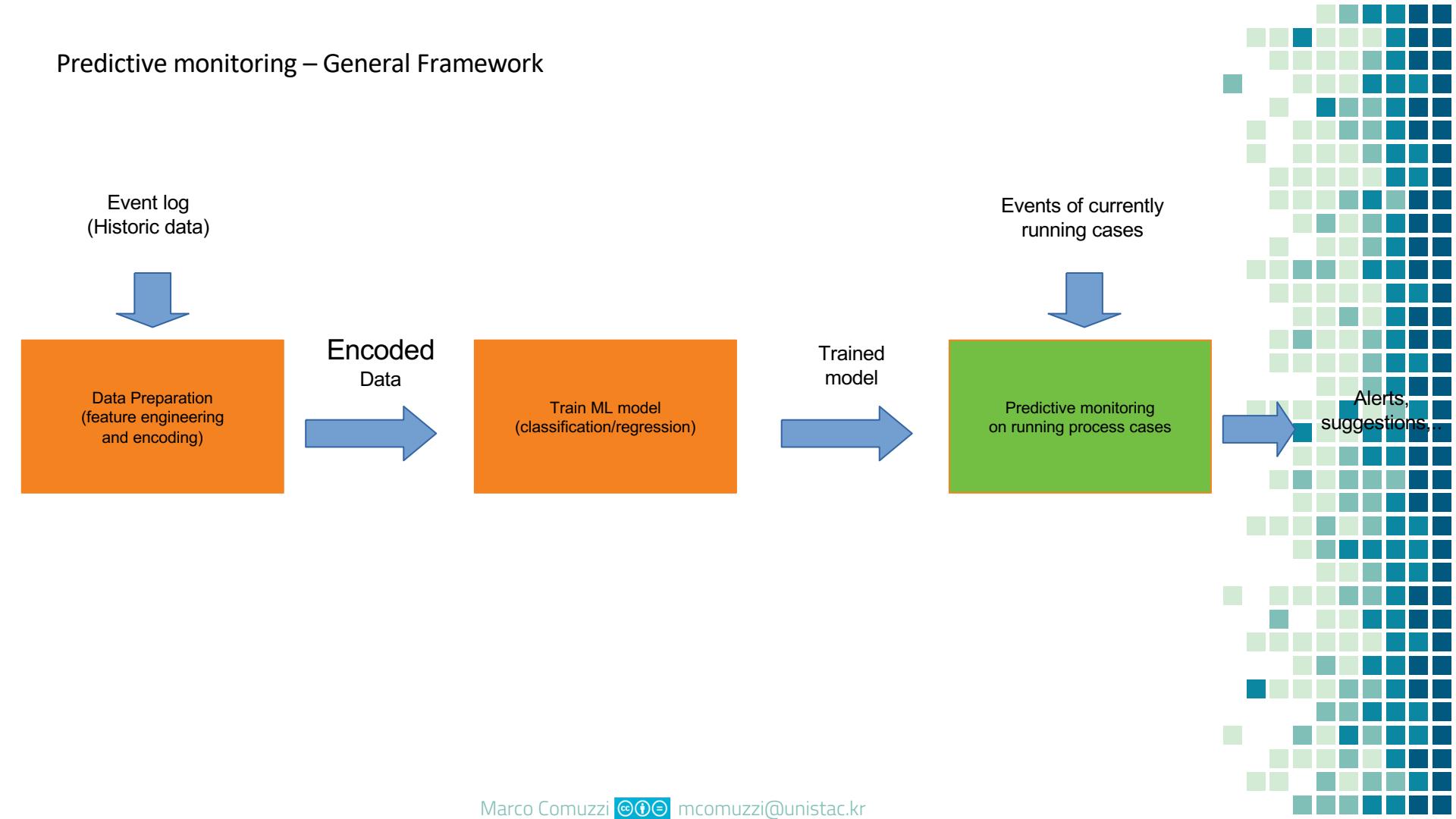
Predicting timestamps

“What would be the outcome of this case, e.g., would it be possible to discharge this patient by tomorrow?”

Predict process outcomes



## Predictive monitoring – General Framework



# Index encoding: example

	Activity, timestamp, resource, cost				Outcome
Trace 1	A, t=1, Alice, 100	B, t=4, Alice, 50	C, t=6, Bob, 60	B, t=7, Chris, 70	1
Trace 2	A, t=3, Alice, 90	C, t=8, Chris, 100	-	-	0

	Act_1*	Res_1*	t_1	c_1	Act_2*	Res_2*	t_2	c_2	Act_3*	Res_3*	t_3	c_3	Act_4*	Res_4*	t_4	c_4	Outcome
Trace 1	A	Alice	1	100	B	Alice	4	50	C	Bob	6	60	B	Chris	7	70	1
Trace 2	A	Alice	3	90	C	Chris	8	100	0	0	0	0	0	0	0	0	0

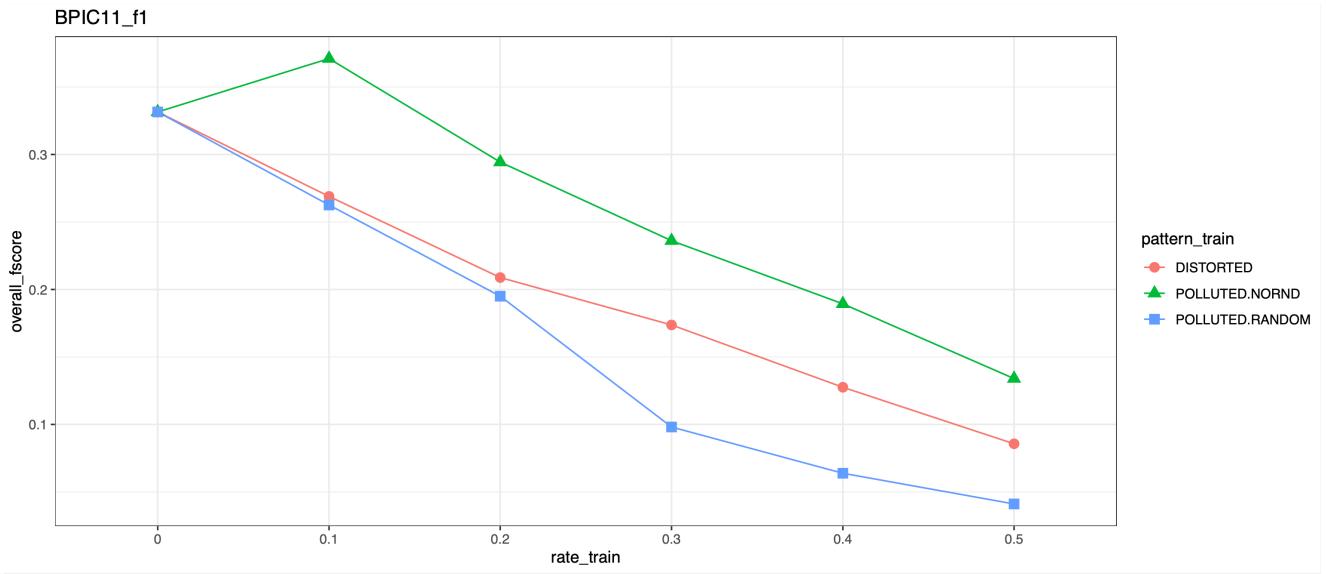
\* Features from one-hot encoding of activity and resource omitted for readability

Error Category	Error Type	Imperfection Pattern?	How to inject error details
Errors on categorical labels (e.g. activity, resources)	Missing labels	NO	<p>Randomly delete activity label of X% events in a log</p> <p>Randomly delete resource label of X% events in a log (if available)</p> <p>X=5, 10, 20, 30, 40, 50%</p> <p>File Names: &lt;log&gt;-MISSING-&lt;label-type&gt;-&lt;X&gt;</p> <p>Examples: "BPIC12-MISSING-activity-20", "Pub-MISSING-resource-50"</p>
Syntactic accuracy: values outside of a pre-defined domain (e.g. typos)	#5 Polluted Label #10 Distorted Label		<p>A) POLLUTED LABELS</p> <p>Randomly pollute activity labels of X% of events in a log X=5, 10, 20, 30, 40, 50%</p> <p>A1) With random string: How="([Activity]_[a-zA-Z{8}]_[Resource]"</p> <p>A2) Without random string: How="([Activity]_[Resource]"</p> <p>B) DISTORTED LABELS</p> <p>Randomly distort activity labels of X% events in a log, X=5, 10, 20, 30, 40, 50%</p> <p>Type of distortion is chosen randomly: How="([Activity:random(Skip,Insert,proximity,interchange,uplow)]")</p> <p>File Names: &lt;log&gt;-&lt;type&gt;-&lt;label-type&gt;-&lt;X&gt;</p> <p>Example: BPIC12-DISTORT-activity-30, Pub-POLLUTED-RANDOM-activity-40, BPIC11-POLLUTED-NORND-activity-5</p>
Semantic accuracy: homonym labels	"#7 Homonymous labels"		<p>Define homonym labels (using ontology), ~X% of labels are substituted by a homonym</p> <p>For now: fixed homonym set, only for credit-card and pub log (see file)</p> <p>Do not use the same homonym: if X = 30%, then 5% of errors concern one homonym, 5% another one etc.</p> <p>File names: Credit-HOMONYM-30</p>
Semantic accuracy: synonym labels	"#6 Synonym labels"		<p>Define synonym labels (using ontology), ~X% of labels are substituted by a synonym</p> <p>For now: fixed synonym set, only for credit-card and pub log (see file)</p> <p>Do not use the same activity synonym: if X = 30%, then 5% of errors concern one activity, 5% another one etc.</p> <p>File names: Credit-SYNONYM-30</p>

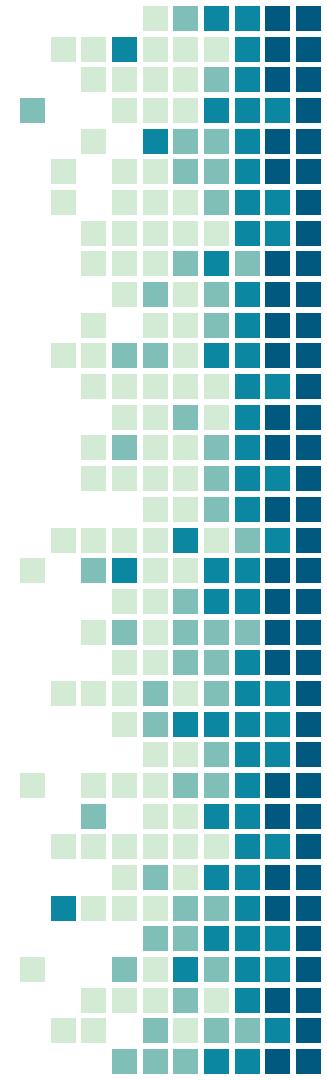
Next activity prediction (SoTA LSTM architecture)

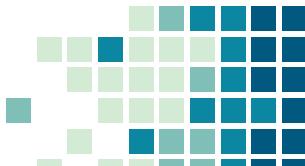
Up to 50% of events in a log affected by pattern

Performance on a “clean” test set



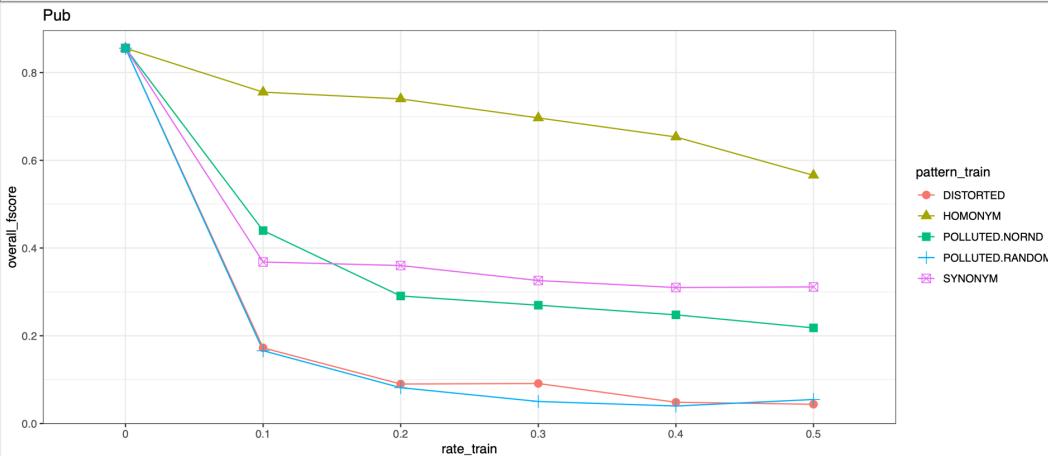
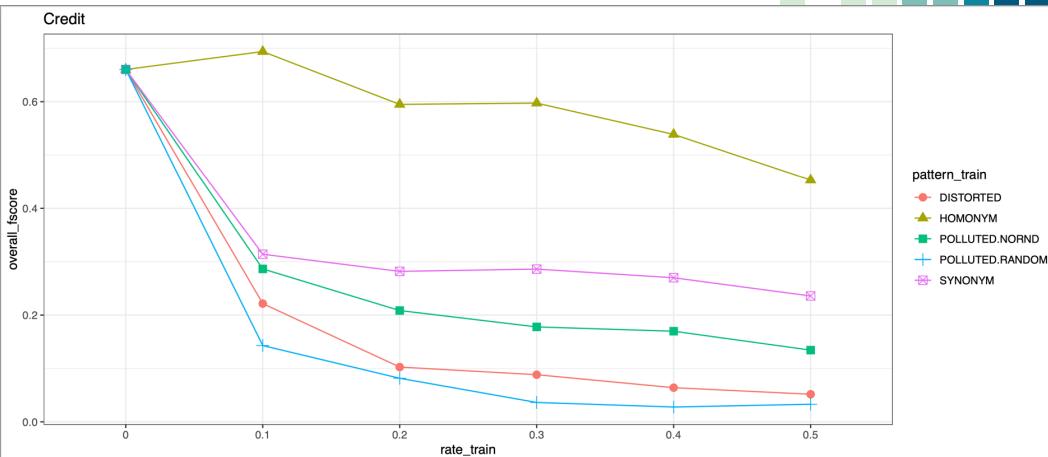
The higher the number of new activity labels introduced by errors, the lower the performance





## Next activity prediction, synthetic logs

Impact of homonym labels  
on PPM performance is  
limited (reduced activity  
labels variability)

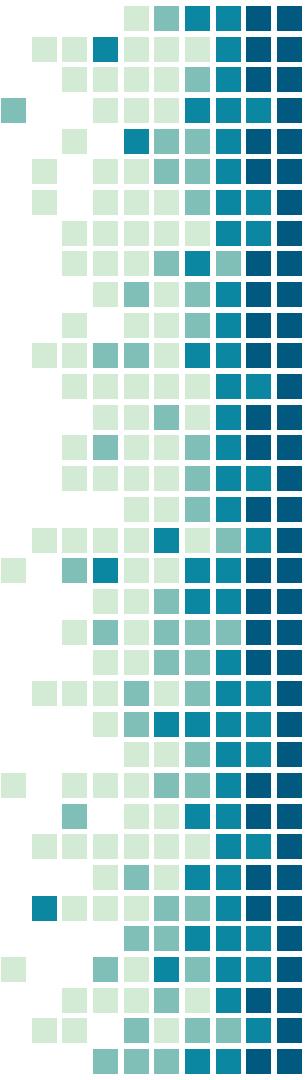


# Some (very preliminary!) insights for a knowledge base

Focus first on errors that increase label variability  
(typos, random errors)

Consider synonym labels only later

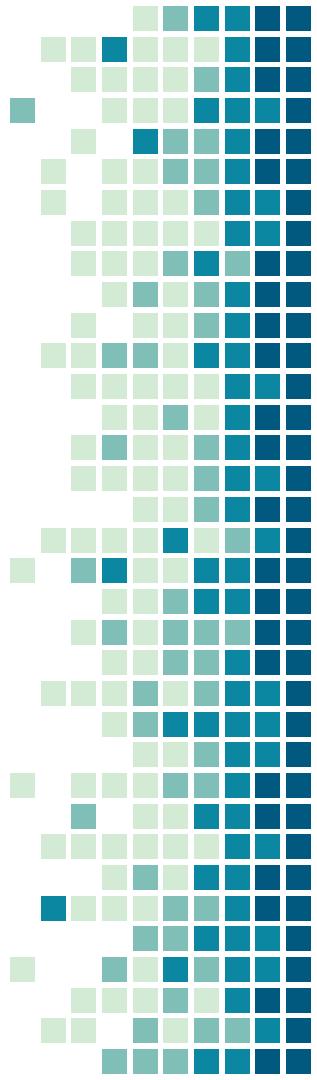
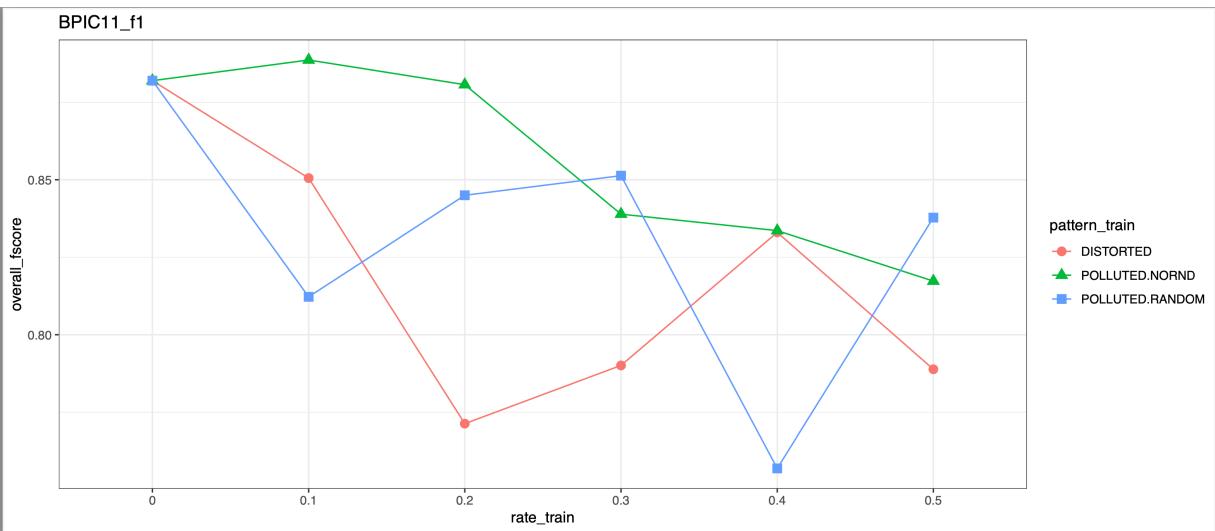
Homonym labels less critical



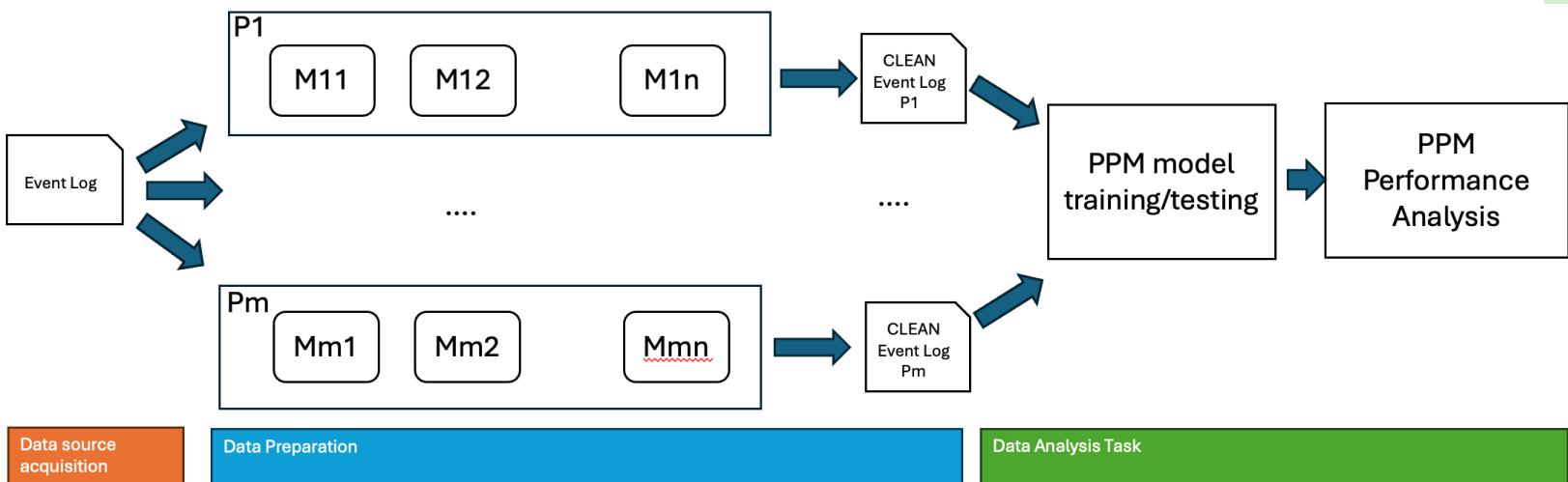
# Outcome prediction?

Hard to see patterns

Outcome label does not depend on which activities are executed?



# Step 2: Building an (optimized) pipeline



# To conclude...

Process mining is now a multi-million dollar industry worldwide

Data quality can be a major hurdle in process mining projects

We need tools to create labelled logs with real-world data quality problems to approach systematically the event log preparation phase



# THANKS!

<https://marco-comuzzi.github.io>

<http://iel.unist.ac.kr/>

@dr\_bsad

mcomuzzi@unist.ac.kr