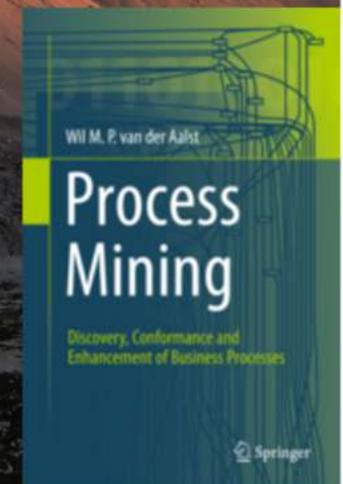


# Process Mining: A historical perspective

prof.dr.ir. Wil van der Aalst

**2013**  
PROCESS MINING CAMP





When did process mining start?



How did PM tooling develop over time?



Three key observations



What are the main research challenges?



How about data mining and business process management?



What are the main PM developments in this century?



Why is process discovery so difficult?

Conclusion

How did PM tooling develop over time?



When did process mining start?



Three key observations



What are the main research challenges?



How about data mining and business process management?

What are the main PM developments in this century?



Why is process discovery so difficult?



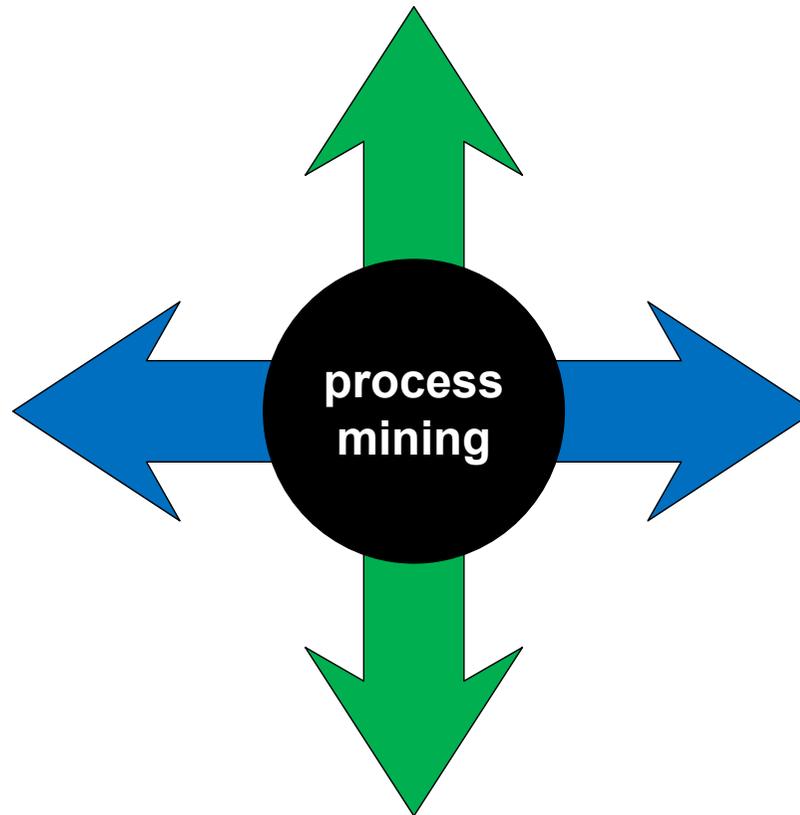
Conclusion

# Positioning Process Mining

## Business Process Management (BPM)

(process analysis/modeling, enactment, verification, etc.)

**performance-oriented questions,  
problems and solutions**

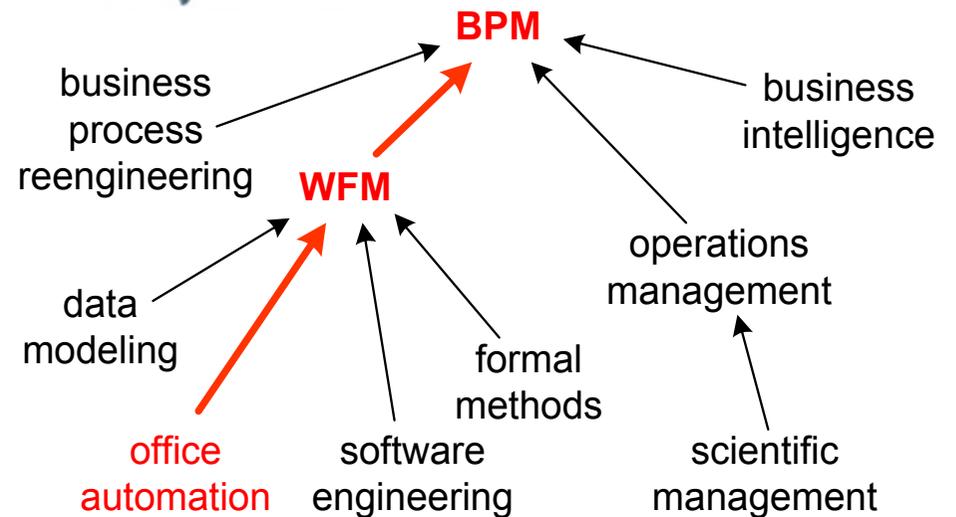
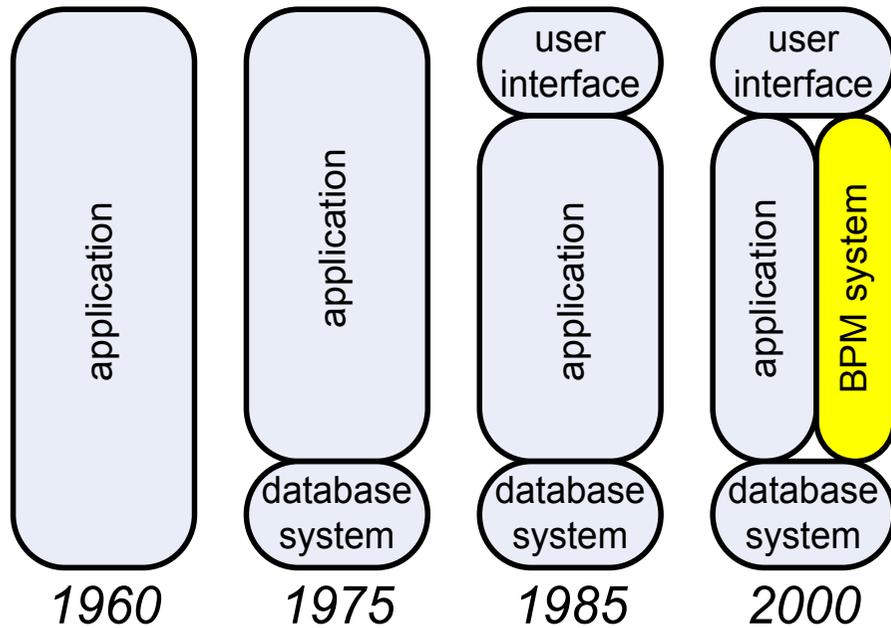


**compliance-oriented questions,  
problems and solutions**

## Data Mining (DM)

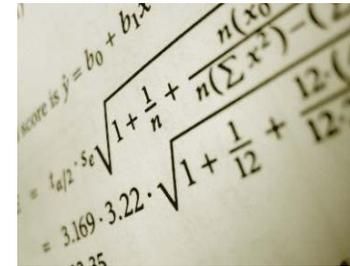
(clustering, classification, rule discovery, etc.)

# History and Origins of BPM

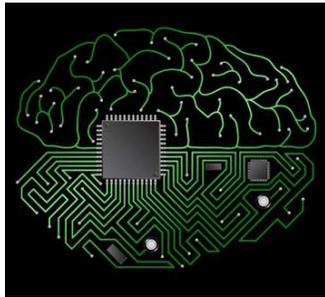


# History and Origins of Data Mining

**Classical statistics (since 500 BC):** descriptive statistics (e.g., sample mean) statistical inference (e.g., confidence interval, regression, hypothesis testing).

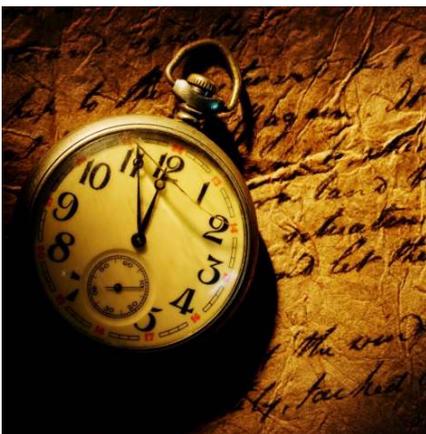


data dredging, data fishing, data snooping



**Artificial intelligence (since 1950):** making intelligent machines by applying human-thought-like processing to statistical problems.

**Machine learning (since 1950):** construction and study of systems that can learn from data.



Many other terms: knowledge discovery, (predictive) analytics, ...



# Data Mining: Supervised Learning

- Labeled data, i.e., there is a **response variable** that labels each instance.
- Goal: explain **response variable** (dependent variable) in terms of **predictor variables** (independent variables).
- **Classification techniques** (e.g., decision tree learning) assume a categorical response variable and the goal is to classify instances based on the predictor variables.
- **Regression techniques** assume a numerical response variable. The goal is to find a function that fits the data with the least error.



# Unsupervised Learning

- Unsupervised learning assumes **unlabeled** data, i.e., the variables are not split into response and predictor variables.
- Examples: **clustering** (e.g., k-means clustering and agglomerative hierarchical clustering) and **pattern discovery** (association rules)

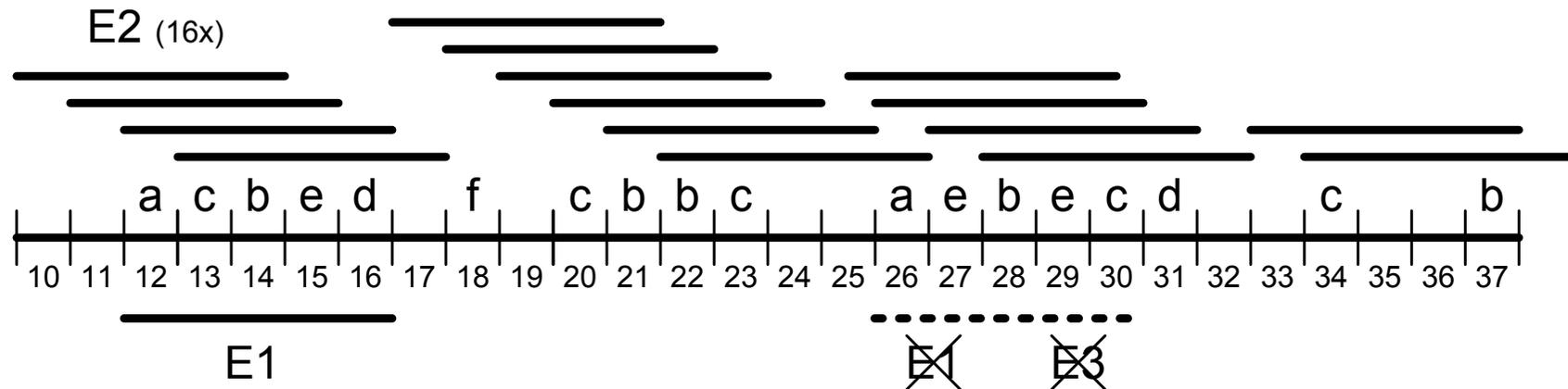
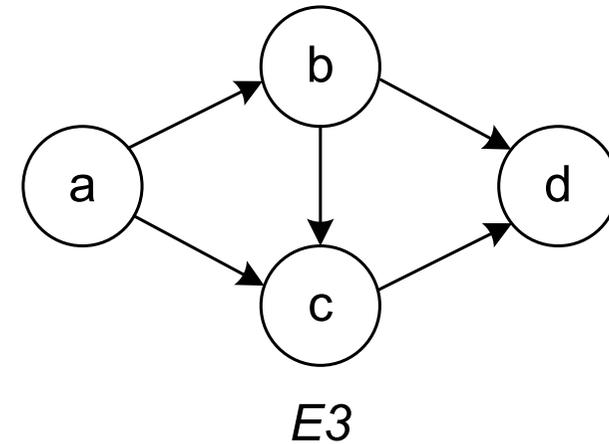
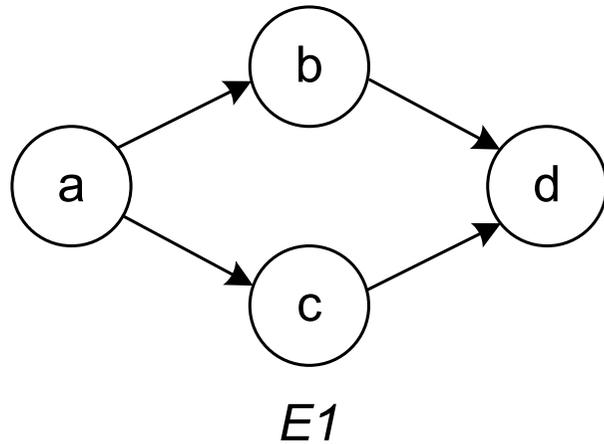
# Example: Association rules

cappuccino	latte	espresso	americano	ristretto	tea	muffin	bagel
1	0	0	0	0	0	1	0
0	2	0	0	0	0	1	1
0	0	1	0	0	0	0	0
1	0	0	0	0	0	0	0
0	0	0	0	0	1	2	0
0	0	0	1	1	0	0	0
...	...	...	...	...	...	...	...

$tea \wedge latte \Rightarrow muffin$

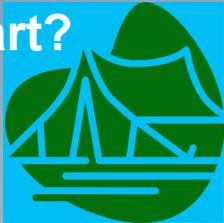
$tea \Rightarrow muffin \wedge bagel$

# Example: Episode Mining



$E2 \Rightarrow E1$  has a confidence of  $1/16$

When did process mining start?



How did PM tooling develop over time?



Three key observations



What are the main research challenges?



How about data mining and business process management?



What are the main PM developments in this century?



Why is process discovery so difficult?



Conclusion



# Language identification in the limit (Mark Gold 1967)

- Mother uses sentences from some language  $\{aab, ab, ab, abc, \dots\}$ .
- "Perfect child" listens to mother and hypothesizes what the full language is like (given all sentences so far).
- Eventually the perfect child's hypothesis is correct and never changes again (without knowing), i.e., only finitely many wrong hypotheses are generated.
- A language is **learnable in the limit** if such a perfect child exists.



# Language identification in the limit (E. Mark Gold 1967)

- Gold showed that most languages cannot be learned in the limit (including the most simple ones like regular languages  $(ab^*(c|d))$ ).
- He noted that it matters whether the child gets **positive** and negative examples (corrections), whether the mother is evil, etc.
- Frequencies matter!
- Representational bias matters!



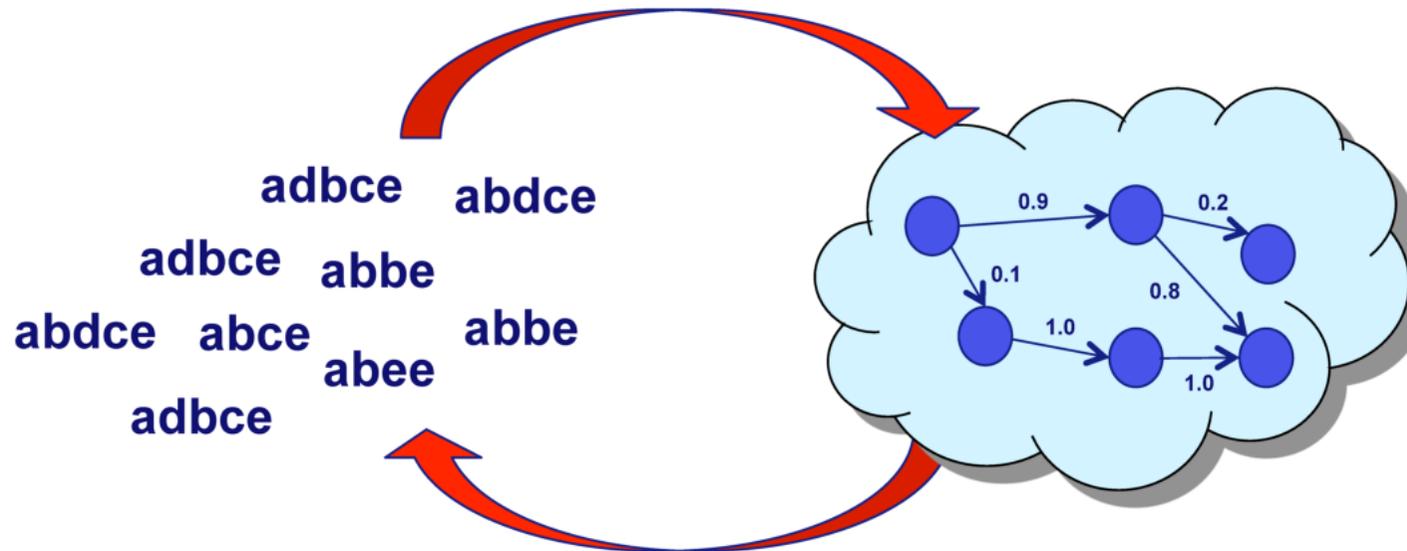
**sentence  $\cong$  trace in event log**

**language  $\cong$  process model**

# Myhill-Nerode Theorem (1958) and the Biermann/Feldman Algorithm (1972)

- There is a **unique minimal** deterministic finite automaton recognizing a **regular language**  $L$  ( shown by John Myhill and Anil Nerode in 1958).
- The equivalence classes defined by  $\cong$  determine the states of the automaton:  $x \cong y$  if there is no  $z$  such that  $xz \notin L$  and  $yz \in L$ .
- Cannot be applied to example traces: overfitting and no generalization.
- Alan W. Biermann and Jerome A. Feldman propose in 1972 techniques to **learn finite state machines from examples** (e.g., considering  $k$ -tails).

# Baum–Welch (1970) and Viterbi (1967) Algorithms to learn Hidden Markov Models



- The **Viterbi algorithm** finds the most likely sequence of hidden states – called the Viterbi path – that results in a sequence of observed events (Andrew Viterbi, 1967).
- The **Baum–Welch algorithm** is an expectation-maximization algorithm that constructs a HMM (Leonard E. Baum and Lloyd R. Welch, 1970).

# Where/when did process mining start?

- Myhill/Nerode (1958)?

- Gold (1965)?

- Baum/Welch (1970)?

- Biermann/Feldman (1972)?

- Rakesh Agrawal (1998)?

– Apriori algorithm for frequent patterns, extended to sequences, episodes, ...

- Jonathan Cook and Alexander Wolf (1998)?

– "Discovering Models of Software Processes from Event-Based Data"

– Using techniques similar to Biermann/Feldman (k-tails) and Baum/Welch (Markov models)

- Rakesh Agrawal, Dimitrios Gunopoulos, Frank Leymann?

– "Mining Process Models from Workflow Logs" (1998)

– Flowmark process models without covering type of splits and joins, no loops, etc.

- Anindya Datta (1998)?

– Automating the Discovery of AS-IS Business Process Models

– Biermann/Feldman style work, embedded in BPM



no concurrency

unable to handle noise

unable to handle incompleteness

no end-to-end models

informal, no precise semantics

# How did process mining start at TU/e?

- Paper and research proposal: **"Process Design by Discovery: Harvesting Workflow Knowledge from Ad-hoc Executions"** (1999)
  - Upcoming move to Technology Management department to lead the IS group (working at CU-Boulder at the time).
  - Collaboration with Ton Weijters stimulated by BETA (linking Petri nets and workflow to Ton's expertise in machine learning).
- First PhDs on process mining (many followed):
  - Laura Maruster
  - Ana Karla Alves de Medeiros
  - Boudewijn van Dongen
- Initial work on **alpha algorithm** (formal limits) and **heuristic** and **genetic** mining (dealing with noise).



# Initial team



When did process mining start?



How did PM tooling develop over time?



Three key observations



What are the main research challenges?



How about data mining and business process management?



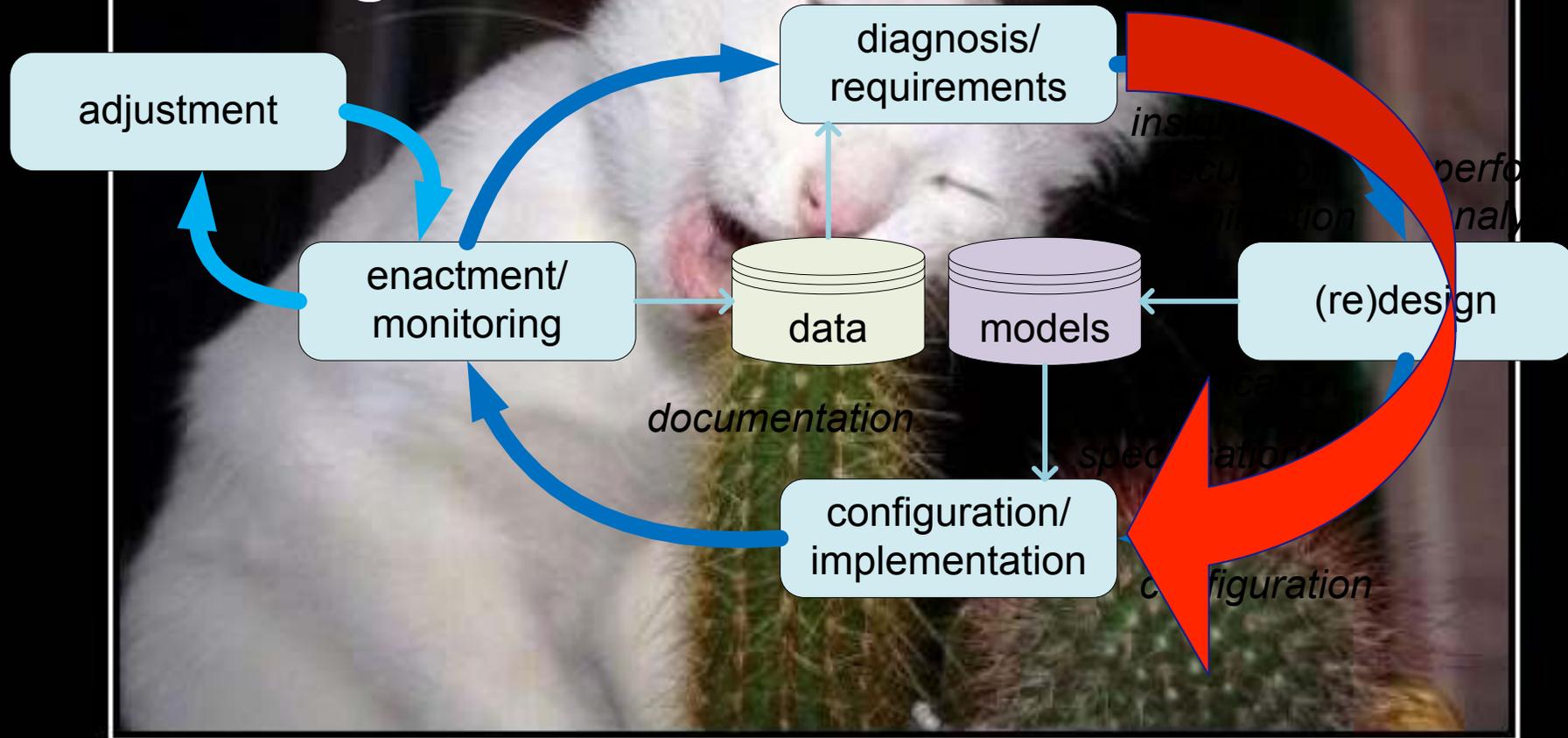
What are the main PM developments in this century?



Why is process discovery so difficult?

Conclusion

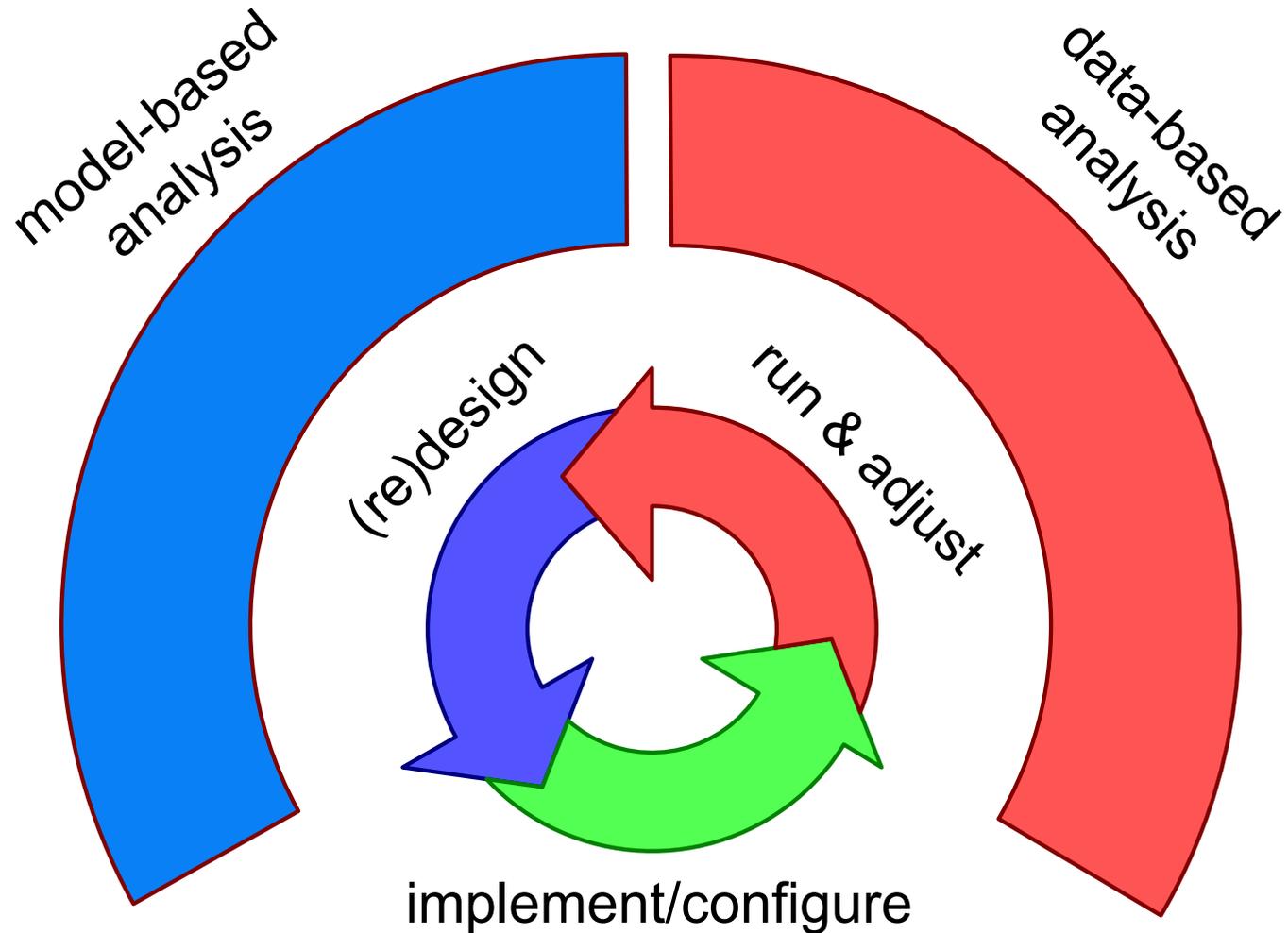
# Workflow Mining



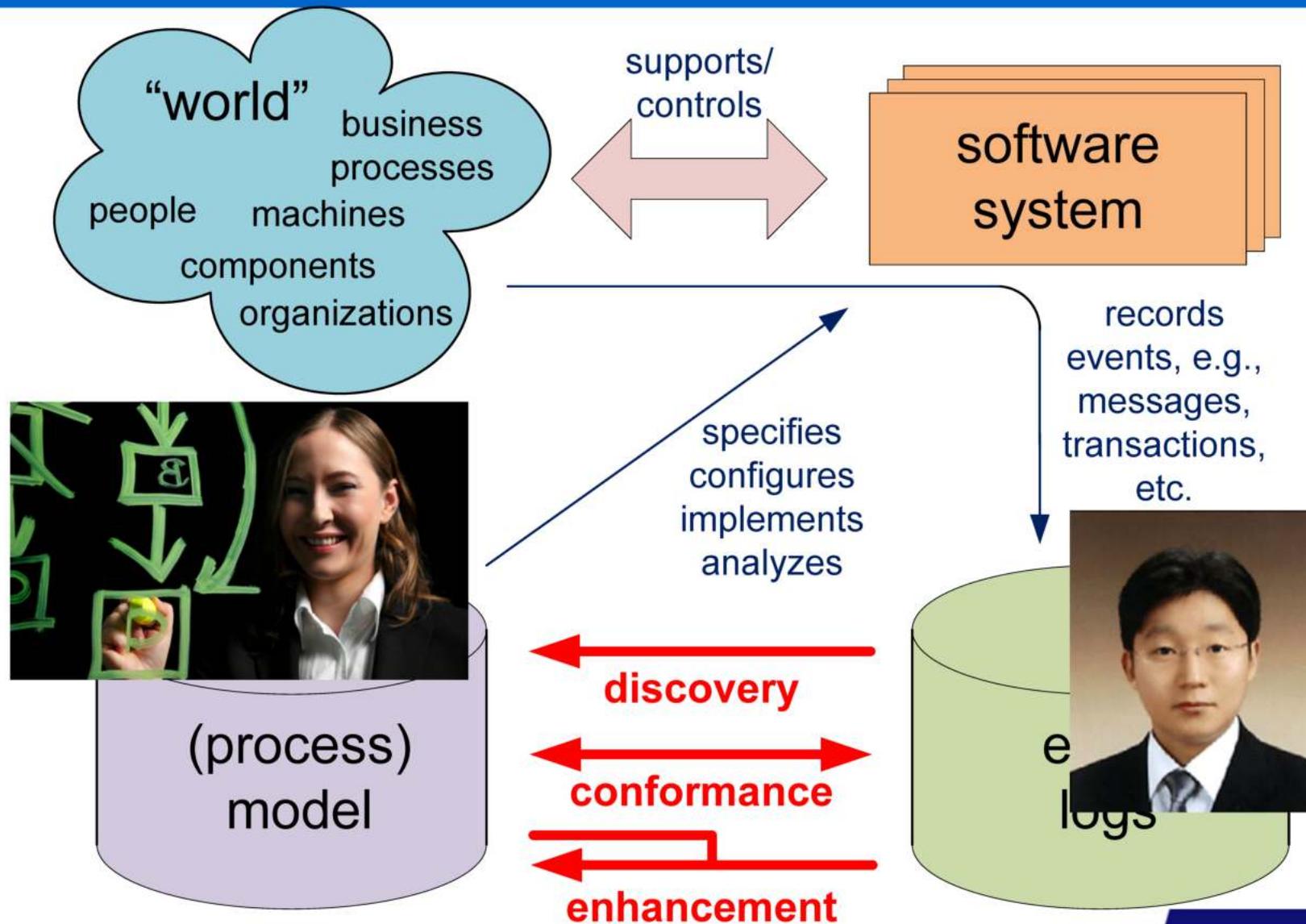
# BAD IDEAS

We All Have Them

# Models, data, and systems coexist



# Process mining spectrum in 2007: Beyond control-flow discovery



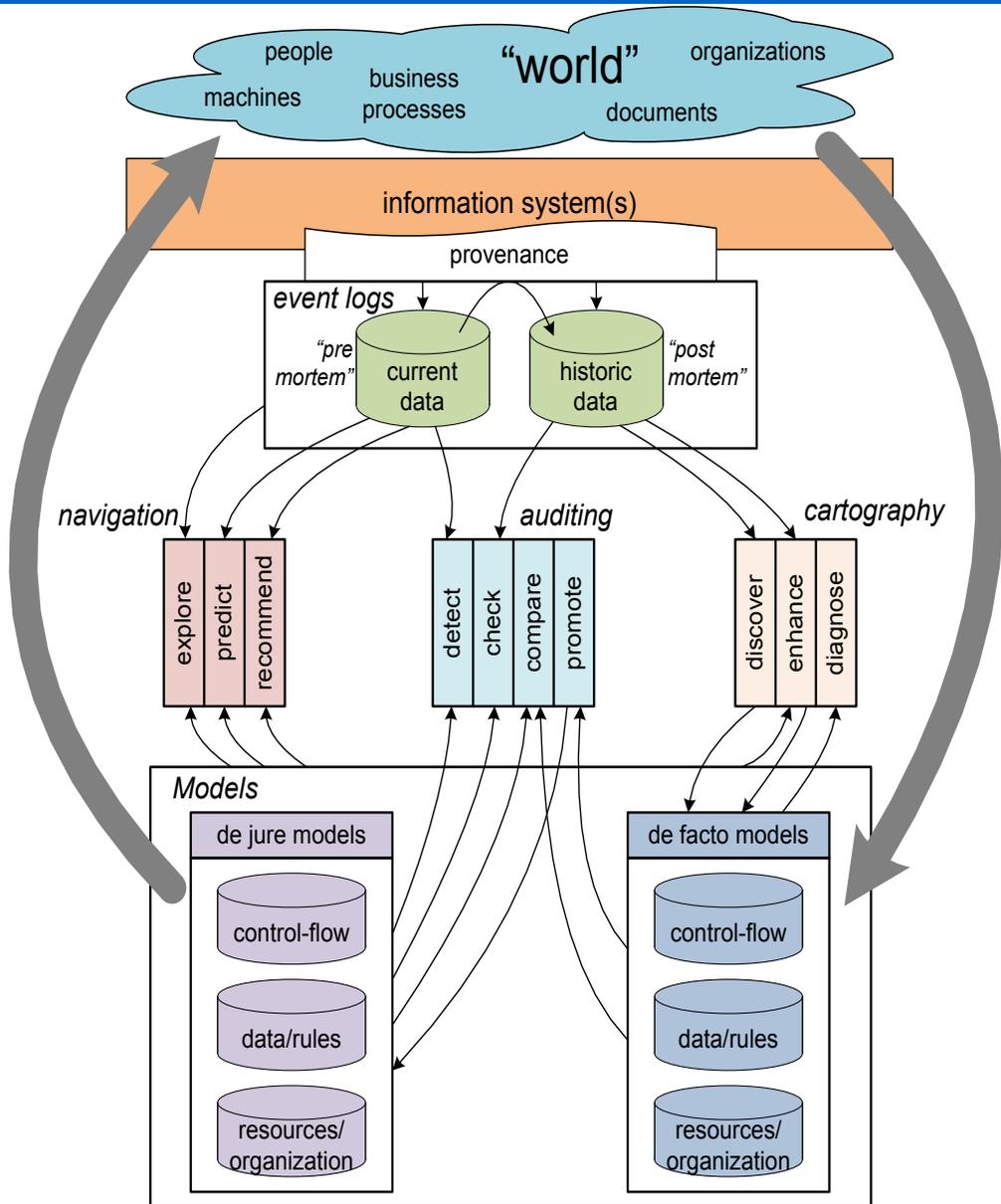
# Team in November 2007



Some people are missing, e.g., Peter van den Brand.

# Current process mining spectrum

(including alignments, operational support, and multiple perspectives)



How did PM tooling develop over time?



When did process mining start?



Three key observations



What are the main research challenges?



Conclusion

How about data mining and business process management?



What are the main PM developments in this century?

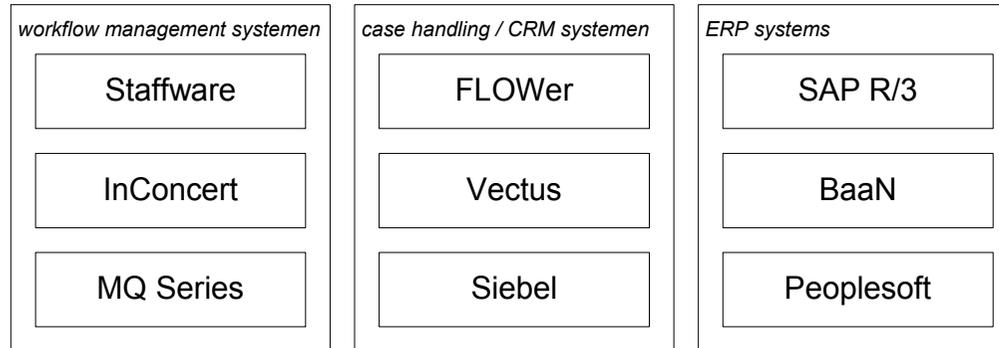


Why is process discovery so difficult?



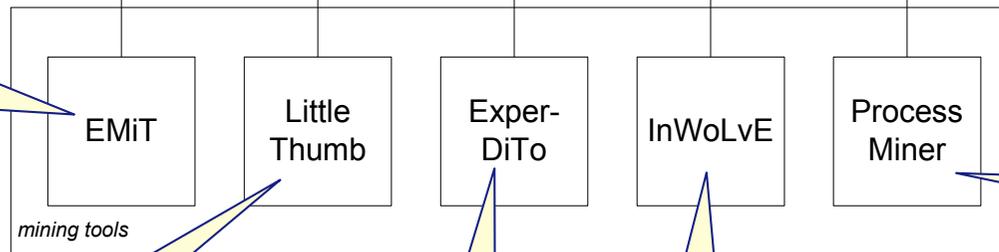
# Pre-ProM

(figure from March 2002!)



The first tool to support the alpha algorithm for process mining was the MiMo (Mining Module) tool based on ExSpect. Later it was implemented in EMiT and ProM.

gemeenschappelijk XML formaat voor het opslaan van workflow logs



mining block structured models (Guido Schimm)



alpha algorithm including time analysis (BvD)



predecessor of ProM's heuristic miner (TW)



evaluation tool (Laura Maruster)



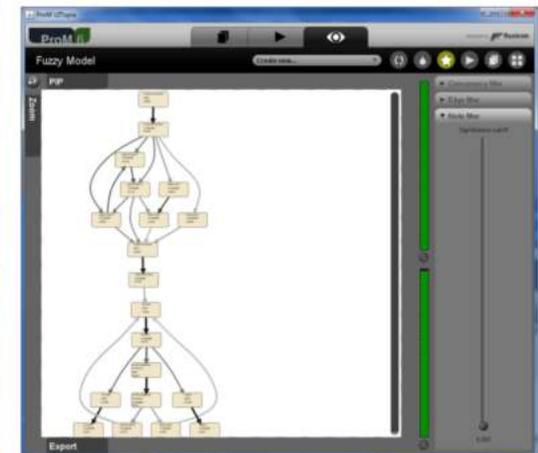
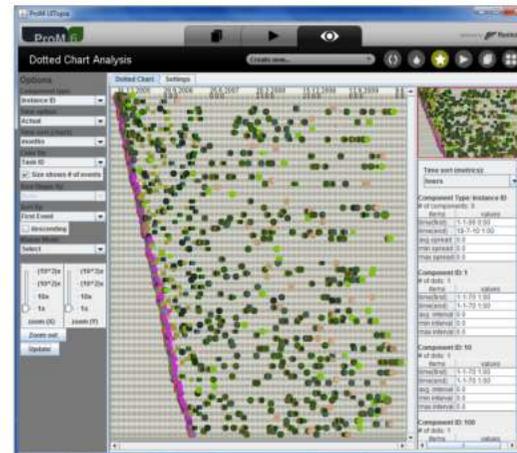
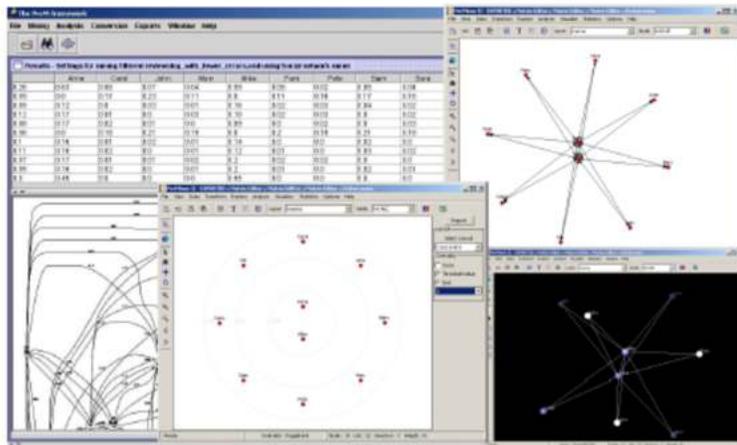
mining with duplicate tasks (Joachim Herbst)



Tobias Blickle (ARIS PPM)

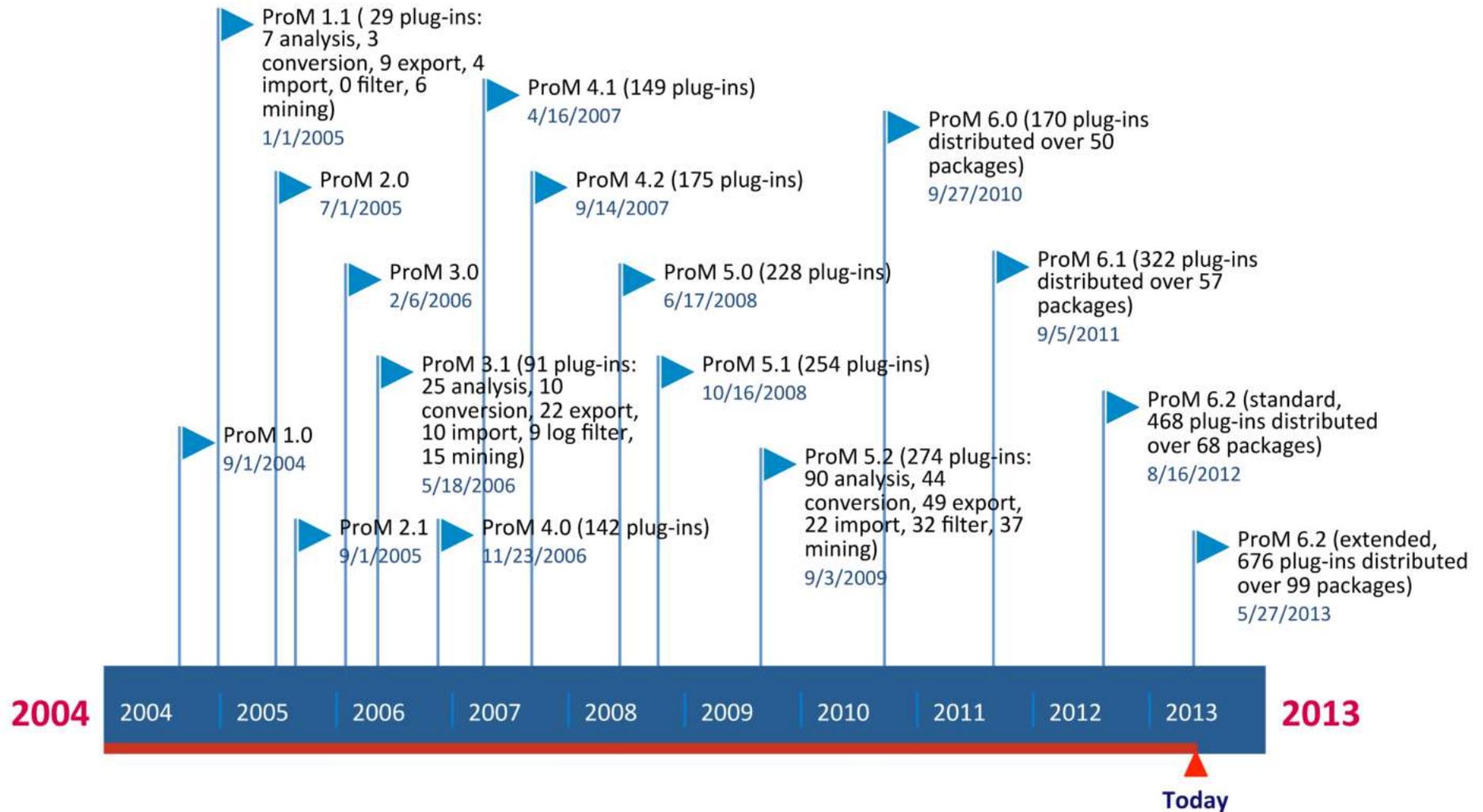


# ProM (2004 – now)



See [www.processmining.org](http://www.processmining.org)

# Overview of ProM releases



# ProM 1.1

**ProM 1.1 ( 29 plug-ins:  
7 analysis, 3  
conversion, 9 export, 4  
import, 0 filter, 6  
mining)**  
1/1/2005

ProM 2.0  
7/1/2005

ProM 3.0  
2/6/2006

ProM 3.1 (91 plug-ins:  
25 analysis, 10  
conversion, 22 export,  
10 import, 9 log filter,  
15 mining)  
5/18/2006

ProM 1.0  
9/1/2004

ProM 2.1  
9/1/2005

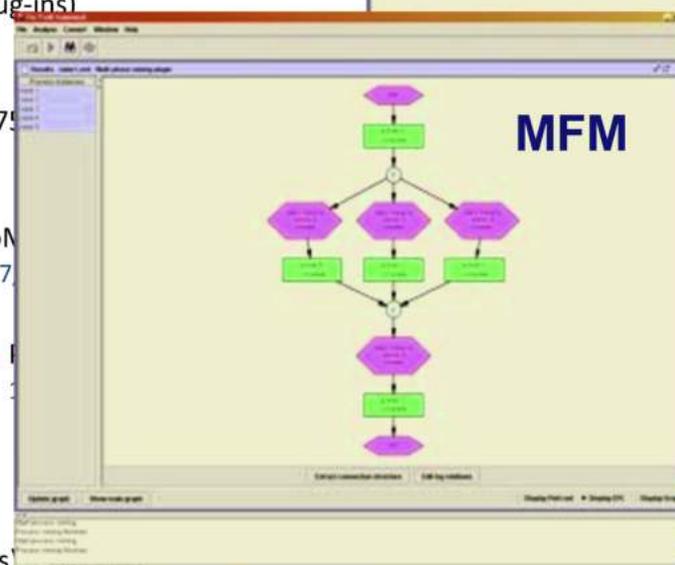
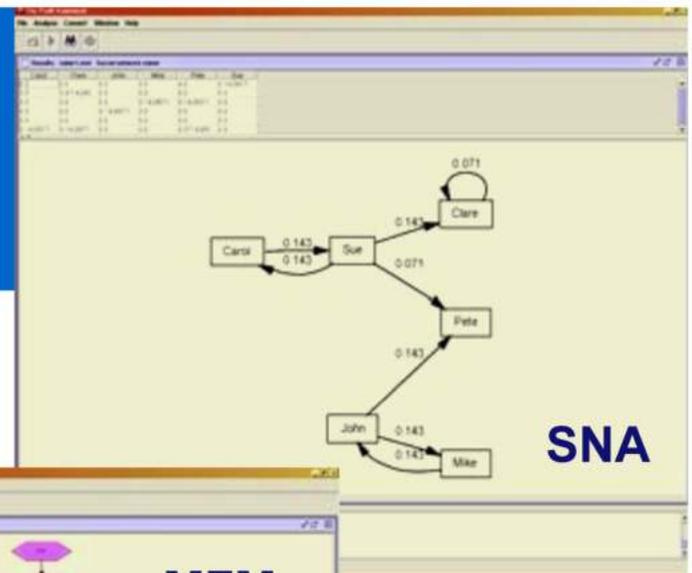
ProM 4.0 (142 plug-ins)  
11/23/2006

ProM 4.1 (149 plug-ins)  
4/16/2007

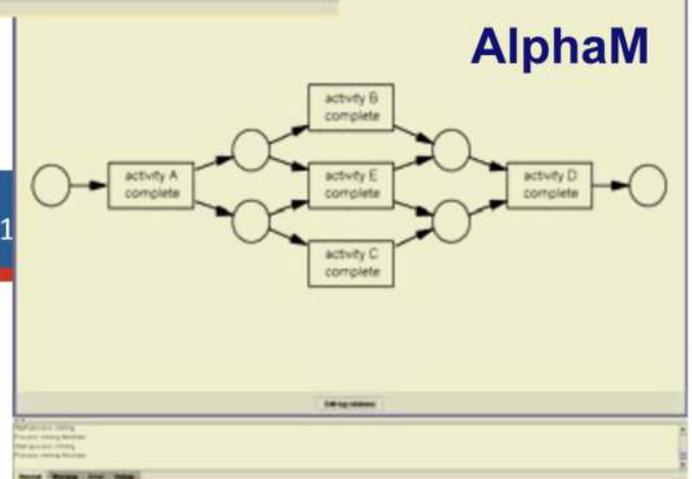
ProM 4.2 (175 plug-ins)  
9/14/2007

ProM 4.3 (182 plug-ins)  
6/17/2008

ProM 4.4 (190 plug-ins:  
mining)  
9/3/2009



(standard, plug-ins distributed packages)



2004

2004

2005

2006

2007

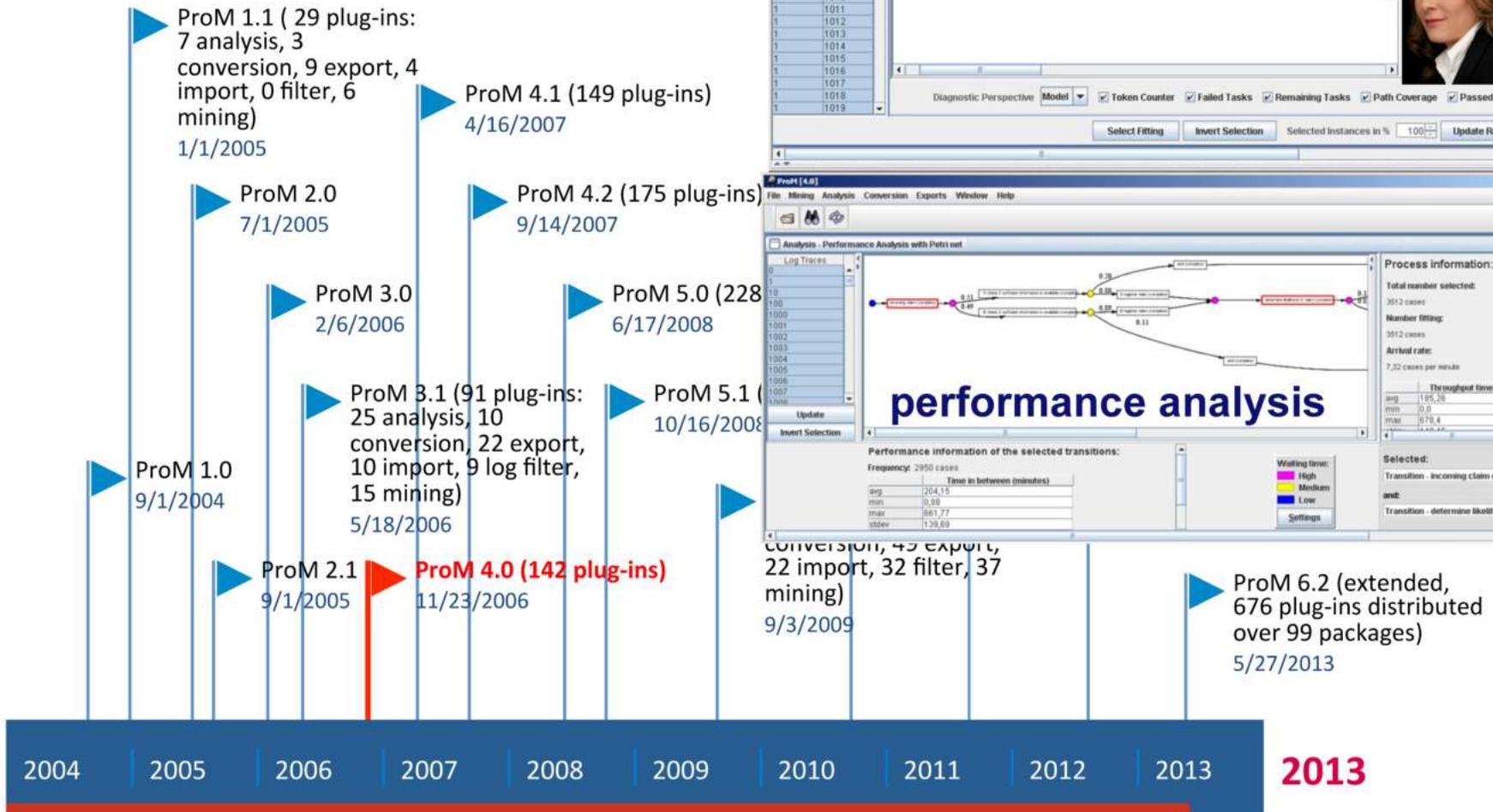
2008

2009

2010

2011

# ProM 4.0



2004

2004

2005

2006

2007

2008

2009

2010

2011

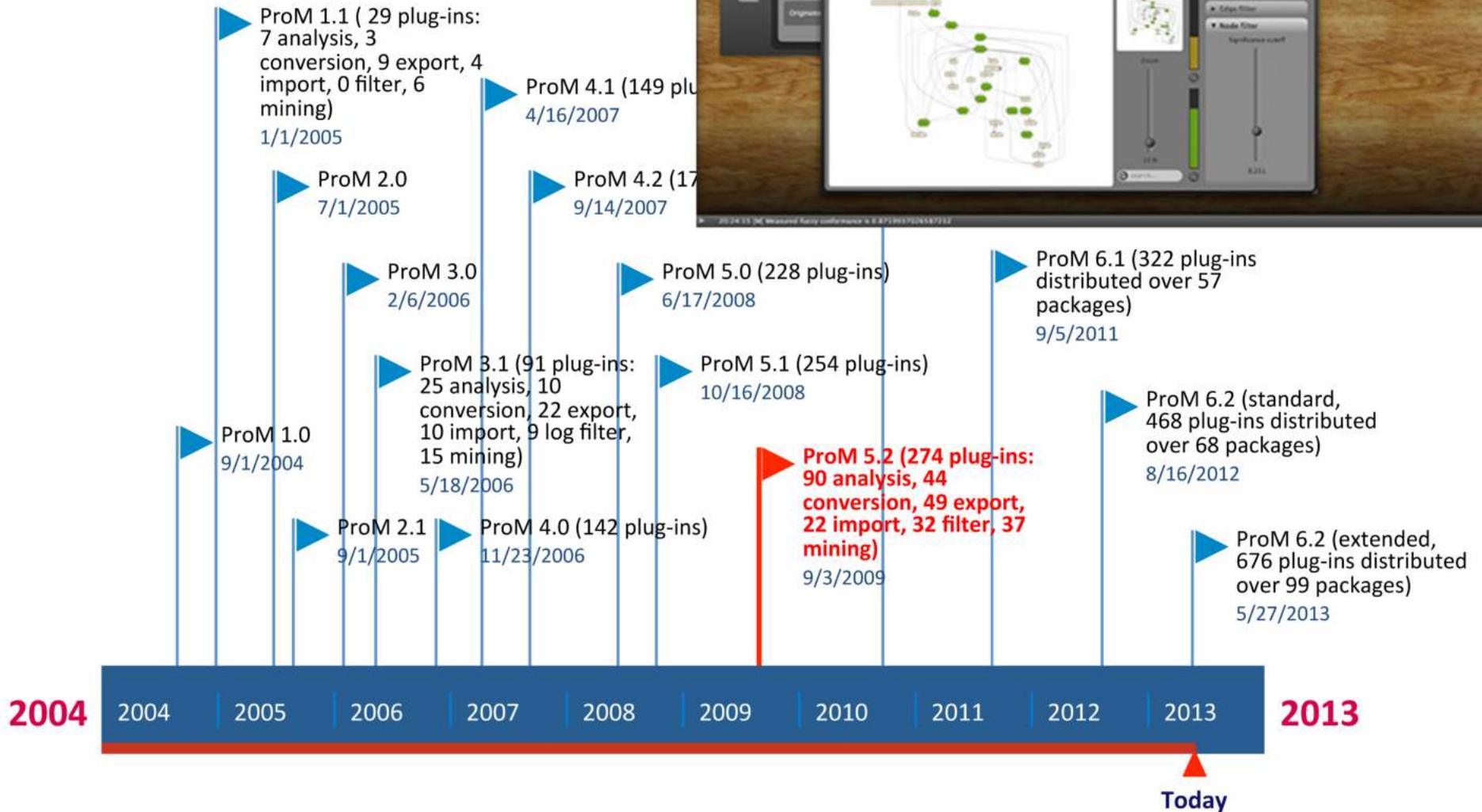
2012

2013

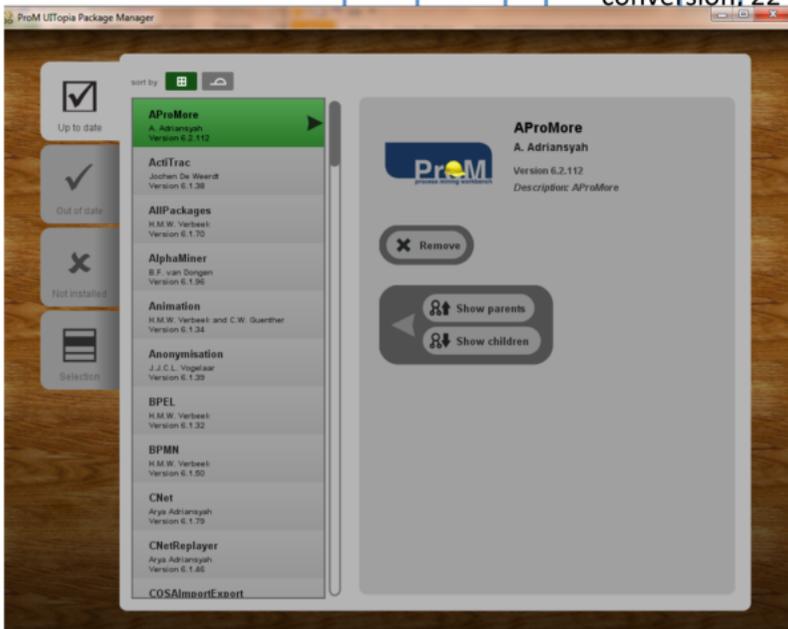
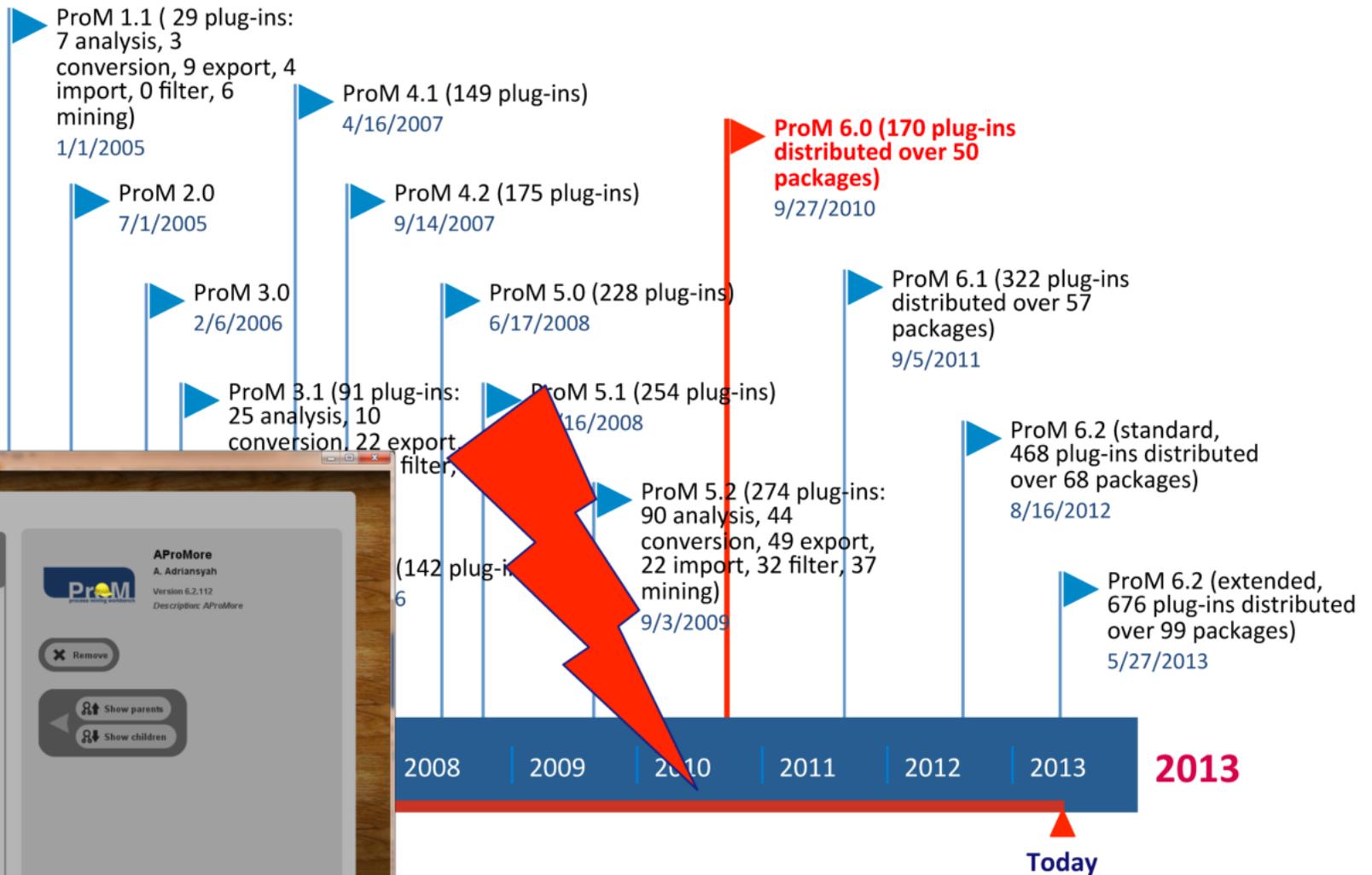
2013

Today

# ProM 5.2



# ProM 6.0: A new start ...



# ProM Today

ProM 1.1 ( 29 plug-ins:  
7 analysis, 3  
conversion, 9 export, 4  
import, 0 filter, 6  
mining)  
1/1/2005

ProM 2.0  
7/1/2005

ProM 4.1 (149 plug-ins)  
4/16/2007

ProM 4.2 (175 plug-ins)  
9/14/2007

ProM 6.0 (170 plug-ins  
distributed over 50  
packages)  
9/27/2010

ProM 6.1 (322 plug-ins  
distributed over 57  
packages)  
9/5/2011

ProM 5.0 (228 plug-ins)

ProM 5.1 (254 plug-ins)  
9/3/2008

ProM 5.2 (274 plug-ins:  
90 analysis, 44  
conversion, 49 export,  
22 import, 32 filter, 37  
mining)  
9/3/2009

ProM 6.2 (standard,  
468 plug-ins distributed  
over 68 packages)  
8/16/2012

**ProM 6.2 (extended, 676  
plug-ins distributed over  
99 packages)  
5/27/2013**

2010

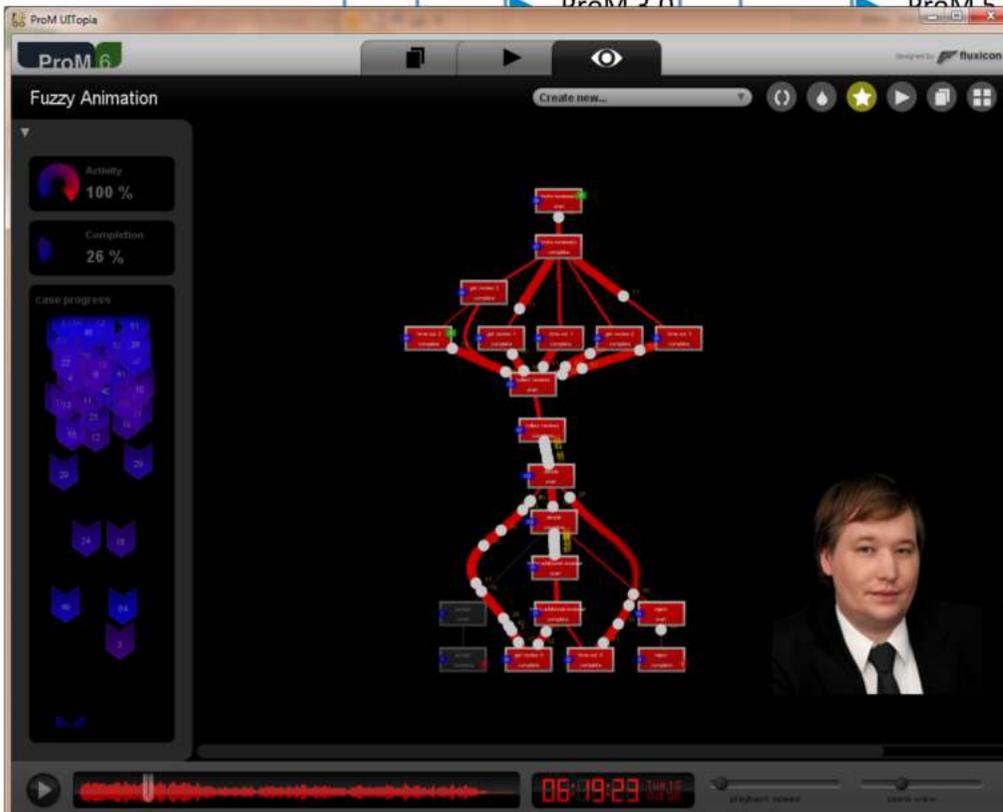
2011

2012

2013

**2013**

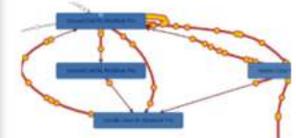
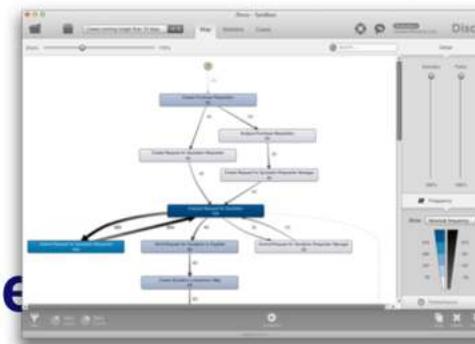
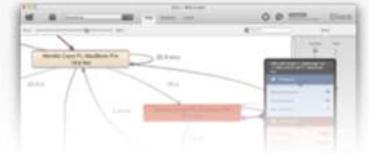
Today



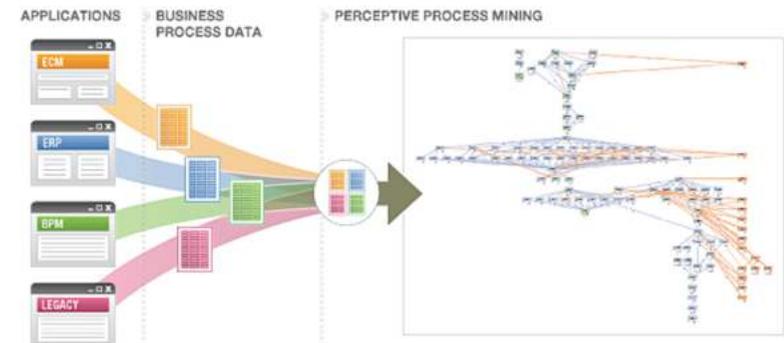
# Commercial PM tools



Process Mining has arrived.  
Finally.



- **Disco (Fluxicon)**
- **Perceptive Process Mining**  
(before Futura Reflect and BPM|one)
- **ARIS Process Performance Manager**
- **QPR ProcessAnalyzer**
- **Interstage Process Discovery**  
(Fujitsu)
- **Discovery Analyst (StereoLOGIC)**
- **XMAalyzer (XMPPro)**
- ...



When did process mining start?



How did PM tooling develop over time?



Three key observations



What are the main research challenges?



How about data mining and business process management?



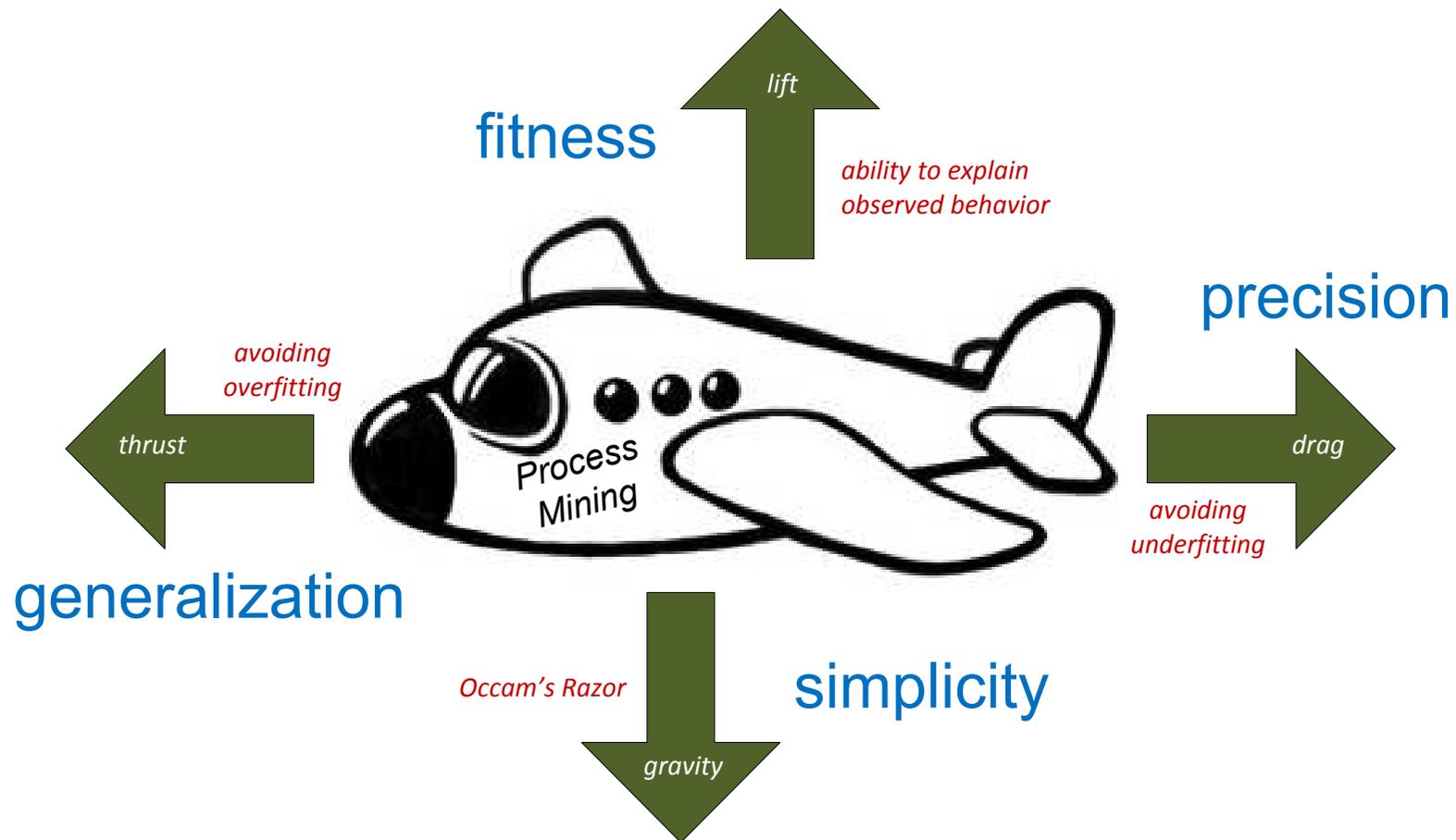
What are the main PM developments in this century?



Why is process discovery so difficult?

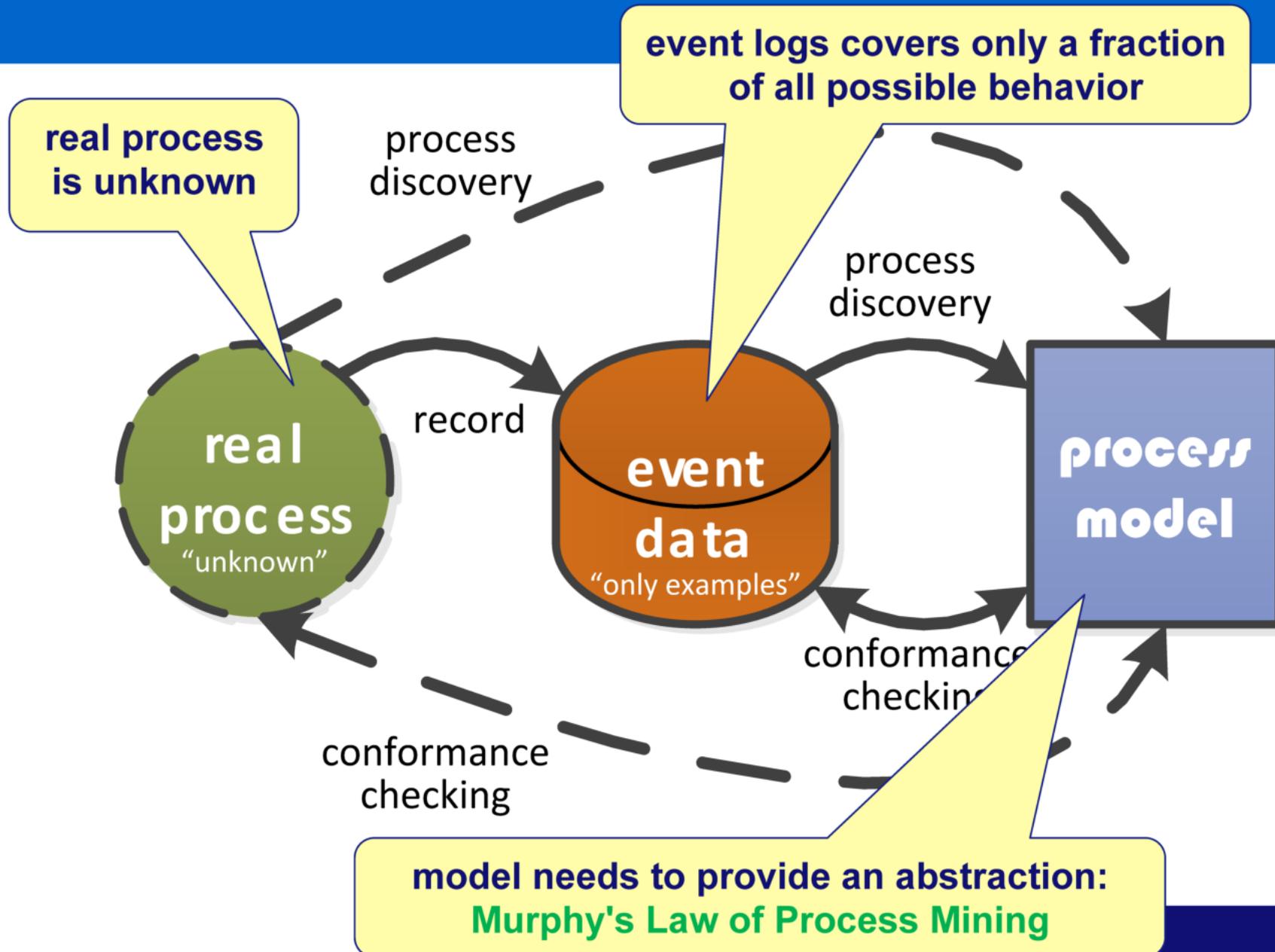
Conclusion

# How good is my model: Four forces



Leaving out one of these dimensions during discovery will lead to degenerate cases!

# Problem





1

formal (not just a picture)

2

fast (should not take years)

ability to balance all conformance dimensions (fitness, precision, generalization, and simplicity) incl. noise

3

4

sound (result should at least be free of deadlocks, etc.)

5

provide guarantees (not just a best effort)

When did process mining start?



How did PM tooling develop over time?



Three key observations



What are the main research challenges?



How about data mining and business process management?



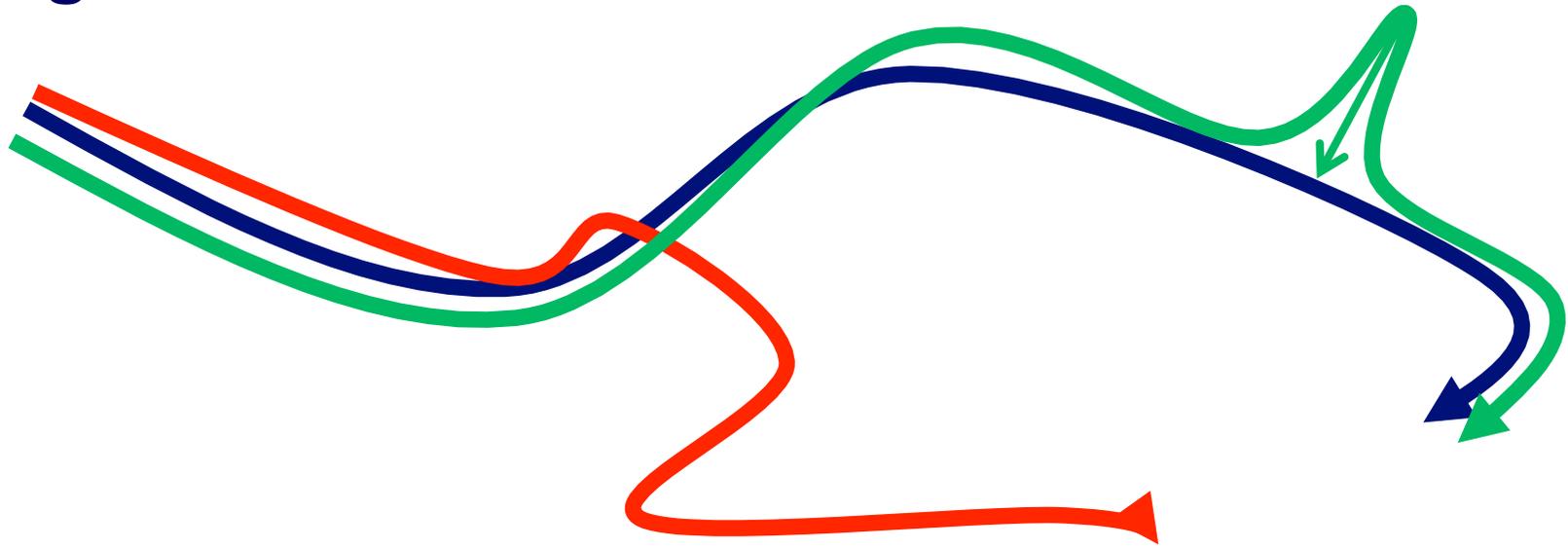
What are the main PM developments in this century?



Why is process discovery so difficult?

Conclusion

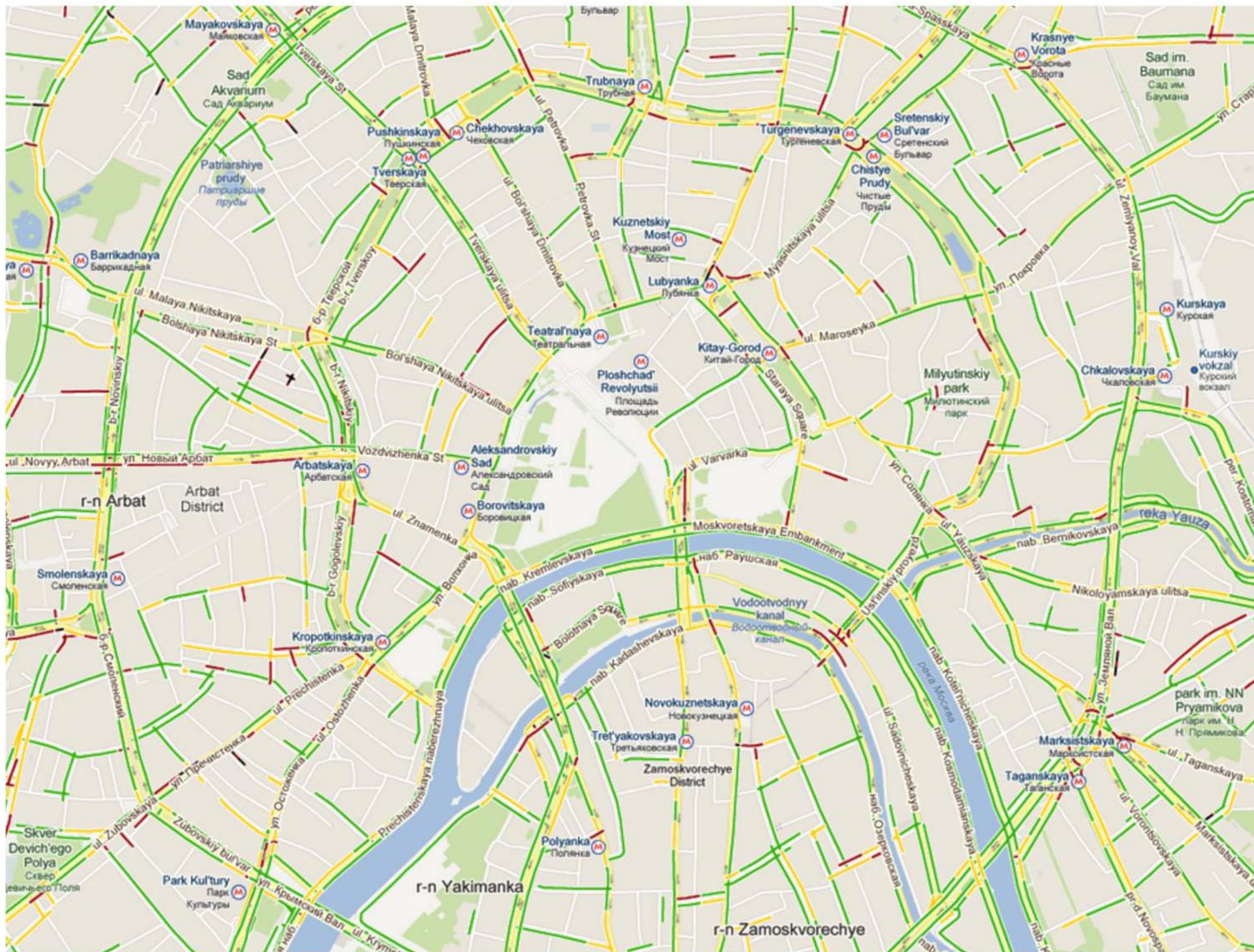
# #1 Alignments are essential!



- conformance checking to diagnose deviations
- squeezing reality into the model to do model-based analysis

$a$	$c$	$\gg$	$d$	$\gg$	$f$	$\gg$
$a$	$c$	$b$	$d$	$\tau$	$\gg$	$h$
$t1$	$t4$	$t3$	$t5$	$t7$		$t10$





When did process mining start?



How did PM tooling develop over time?



Three key observations



What are the main research challenges?



How about data mining and business process management?



What are the main PM developments in this century?



Why is process discovery so difficult?

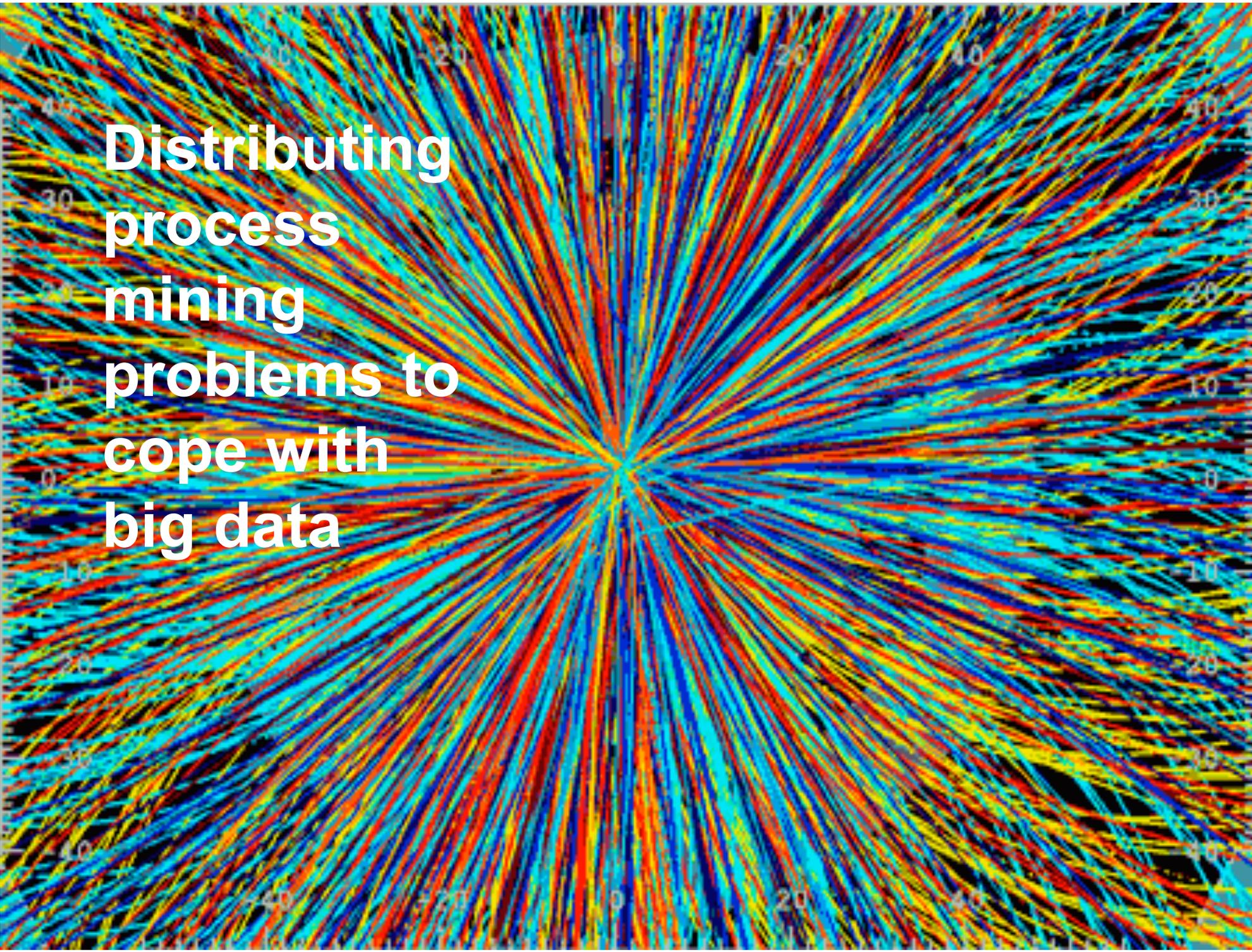
Conclusion



**Finding  
sheep with  
five legs**

we are getting close...

**Distributing  
process  
mining  
problems to  
cope with  
big data**



# On-the-fly process mining



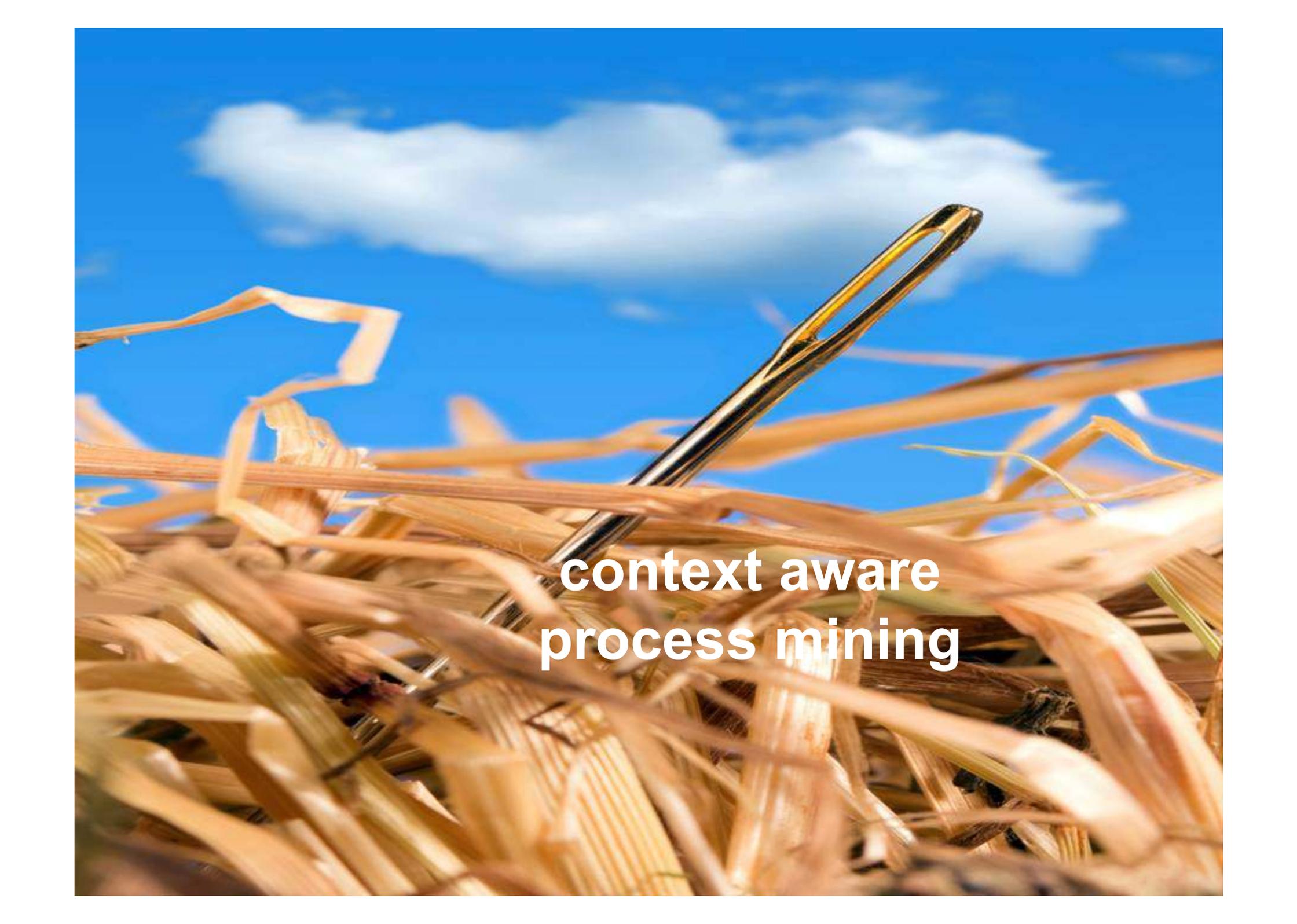
Operational  
support

Concept drift



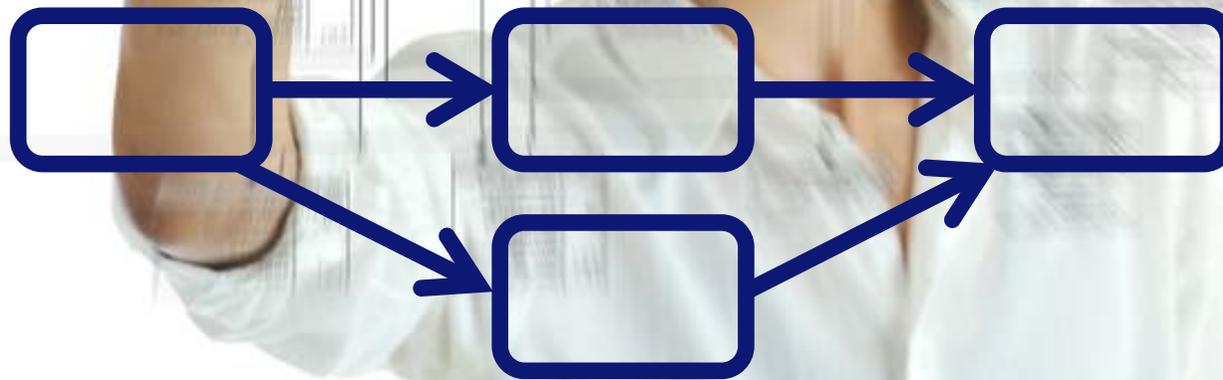


**cross-organizational /  
comparative process mining**



**context aware  
process mining**

# Supporting the process of process mining



When did process mining start?



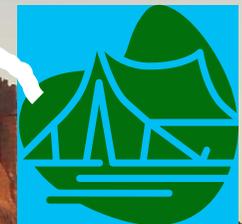
How did PM tooling develop over time?



Three key observations



What are the main research challenges?



Conclusion

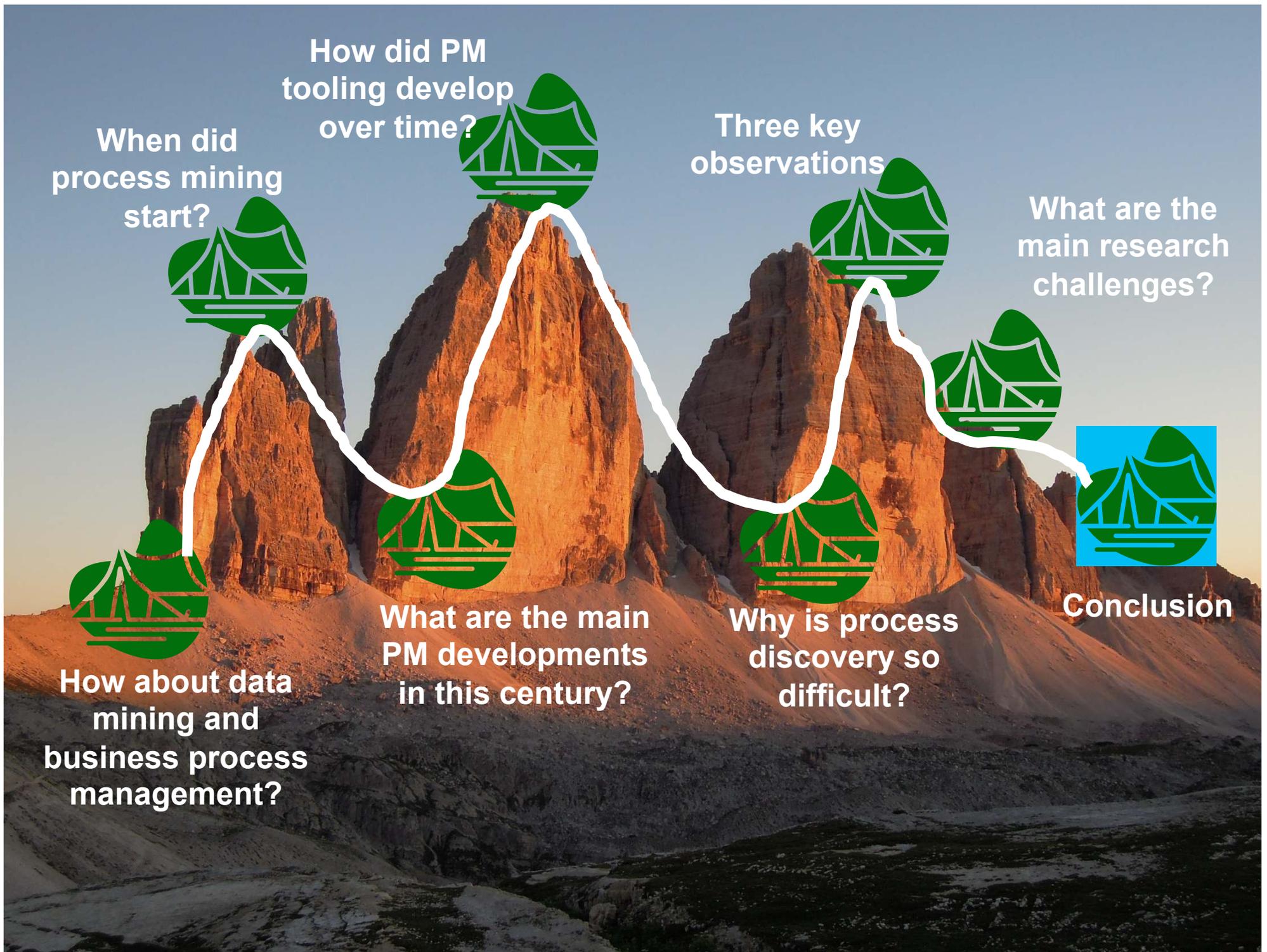
How about data mining and business process management?



What are the main PM developments in this century?



Why is process discovery so difficult?



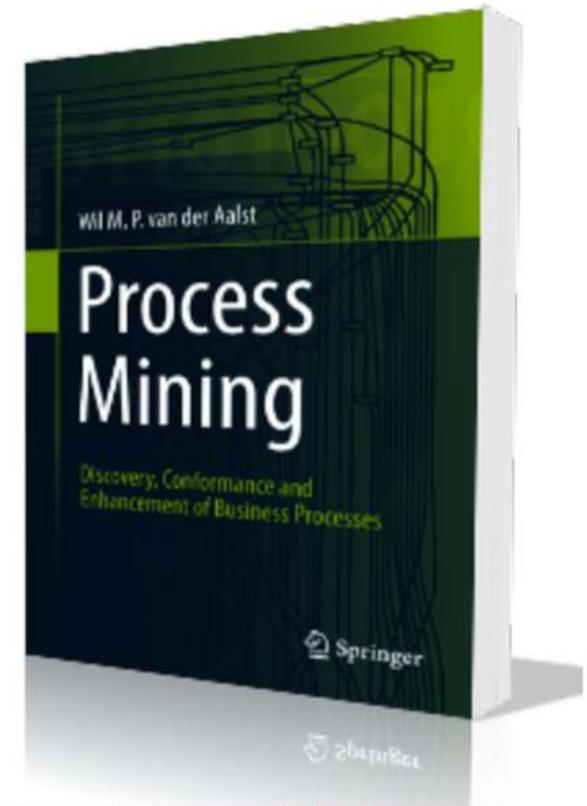
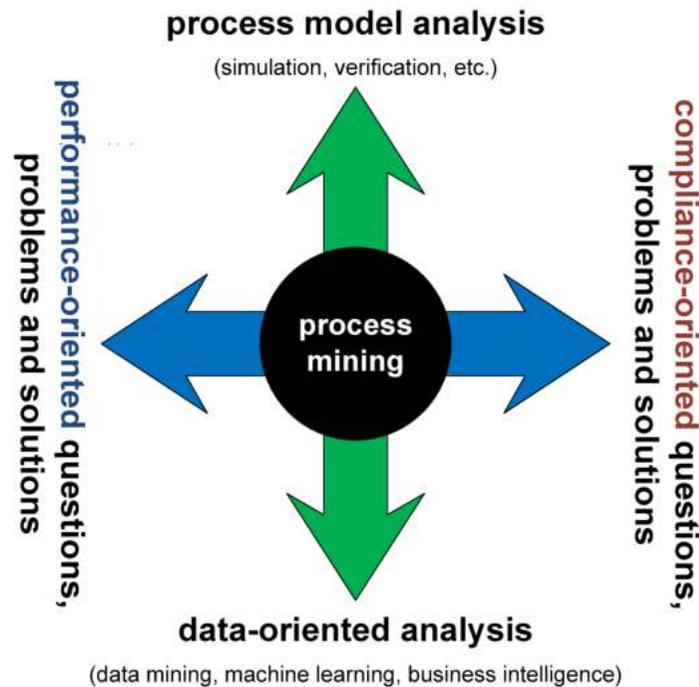


# Turning Event Data into Real Value



## Data Scientist: The Sexiest Job of the 21st Century

by Thomas H. Davenport and D.J. Patil



[processmining.org](http://processmining.org)



<http://www.win.tue.nl/ieeetfpm/>