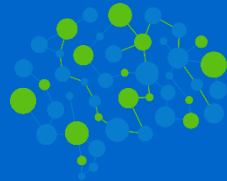# Learning Hybrid Process Models from Events

*Process Mining for the Real World*
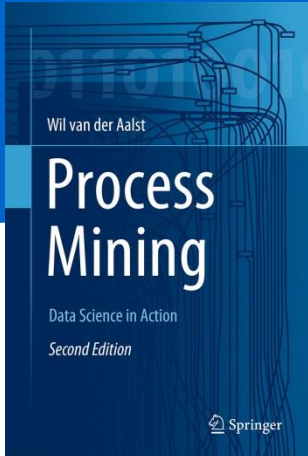
**Smart Data Analytics (SDA) research group,
University of Bonn, 29-11-2017**

SMART
DATA
ANALYTICS
FROM DATA TO KNOWLEDGE

## Wil van der Aalst

**www.vdaalst.com | @wvdaalst**

Wil van der Aalst

**Process Mining**

Data Science in Action

*Second Edition*

Springer

TU/e

TU/e
DSC/e

# Positioning PM

Launched in 2013

**Data Science as "the enabler"**

| infrastructure | analysis | effect |
|---|---|---|
| "volume and velocity" | "extracting knowledge" | "people, organizations, society" |

- instrumentation
- big data infrastructures and distributed systems
- databases and data management
- programming
- security
- ...

- statistics
- data/process mining
- machine learning/artificial intelligence
- operations research
- algorithms
- visualization
- ...

- ethics & privacy
- IT law
- human technology interaction
- operations management
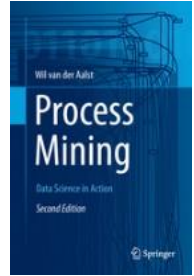- business models
- entrepreneurship
- ...

# the data science pipeline

**TU/e**

JADS

Jheronimus
Academy
of Data Science

# Uptake of process mining



Start of process mining at TU/e (1999)

Alpha miner & Heuristic miner (2000-2002)

First version of ProM (2004)

Evolved from 29 plug-ins in 2004 to 274 plug-ins 2009 (ProM 5.2) to over 1500 plug-ins today.

Conformance checking, other perspectives, prediction, etc. (2005 - )

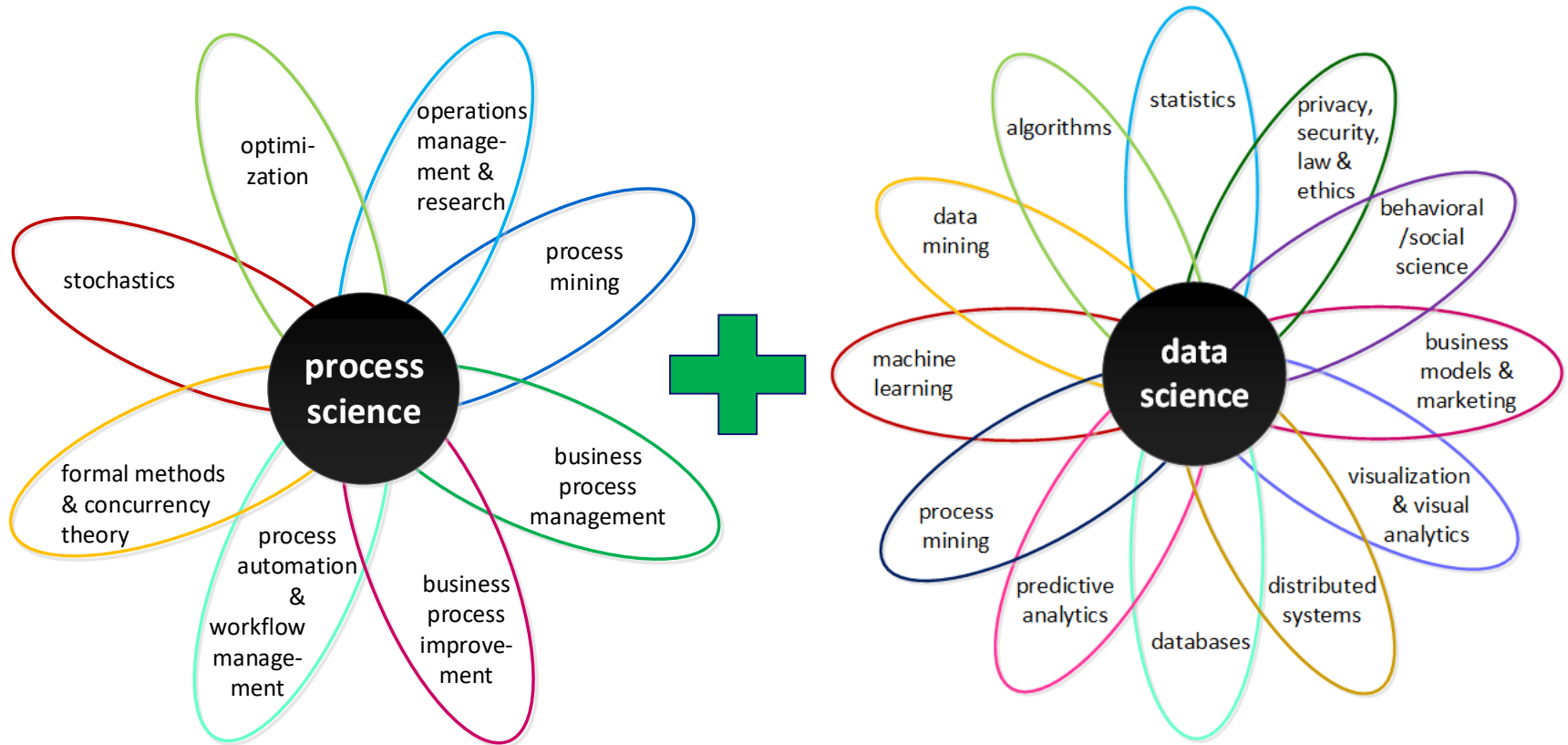Process mining book, courses, uptake commercial tools, etc. (2011 - )

publications on process mining (2016/17 incomplete) based on Scopus data

**process science**
- optimization
- operations management & research
- process mining
- stochastics
- business process management
- formal methods & concurrency theory
- process automation & workflow management
- business process improvement

**+**

**data science**
- statistics
- algorithms
- privacy, security, law & ethics
- data mining
- behavioral /social science
- machine learning
- business models & marketing
- process mining
- visualization & visual analytics
- predictive analytics
- databases
- distributed systems

TU/e

process mining as the missing link

# Taxonomy: Not just CF discovery!



Task
- Process discovery
- Conformance checking
- Extending process models
- Decision making

Perspective
- Control Flow (CF) only
- CF + time
- CF + data
- CF + resources

Type
- Offline
- Online

TU/e

# Taxonomy: Not just CF discovery!

offline control-flow discovery

**Task**
- Process discovery
- Conformance checking
- Extending process models
- Decision making

**Perspective**
- Control Flow (CF) only
- CF + time
- CF + data
- CF + resources

**Type**
- Offline
- Online

TU/e

process model discovered using the inductive miner (showing only the most frequent paths)

# using conformance checking to see all deviations

seamless abstraction:
one log many views

# Data-driven & process-centric

# Answering two types of questions



**performance questions**

*Why are these cases late?*

*Where are the real bottlenecks?*

*Which resources are overloaded?*

**data-driven**

**P**

**process-centric**

**M**

**event data**

**process models**

*How often is the four-eyes principle violated?*

*Which activities are often skipped?*

*Which resources cause deviations?*

**conformance questions**

TU/e

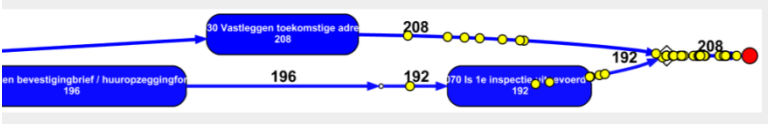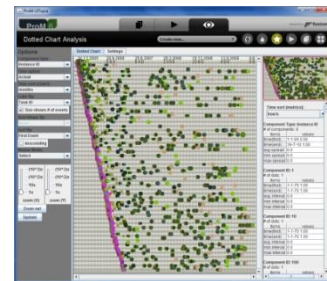Process mining results may hurt and trigger resistance, but this only supports the need for it.
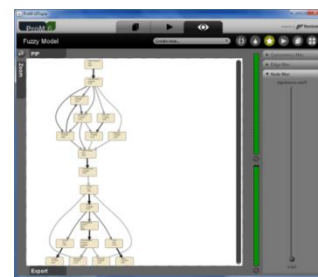
# Process Mining Software

1500+ plug-ins available covering the whole process mining spectrum

100% FREE

>130k downloads

ProM
process mining workbench
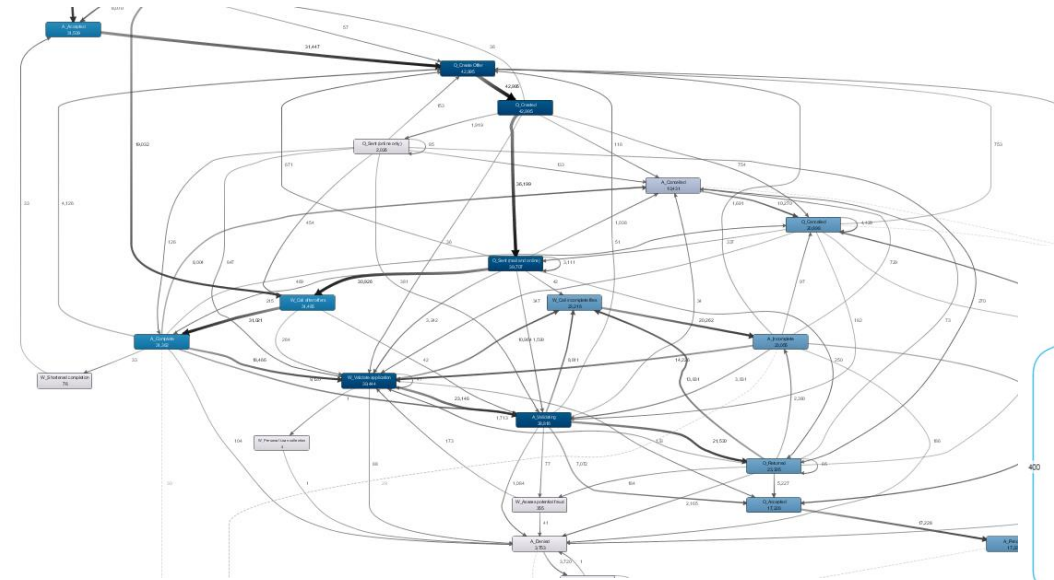
# Interaction with industry

# Job done?

# *Not really …*

# Concurrency and Semantics

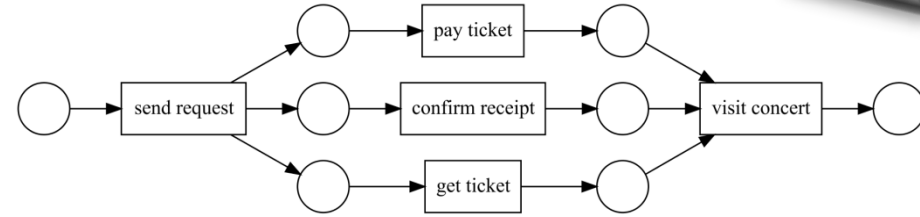# Boxes and arrows: What do they mean?
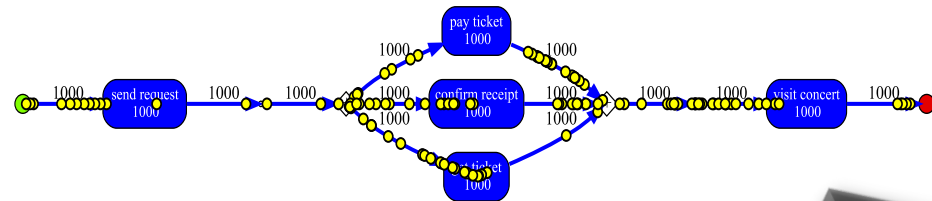
# An example log and "proper models"

402 traces
40.20% of the log
send request | pay ticket | confirm receipt | get ticket | visit concert

314 traces
31.40% of the log
send request | confirm receipt | pay ticket | get ticket | visit concert

250 traces
25.00% of the log
send request | pay ticket | get ticket | confirm receipt | visit concert

20 traces
2.00% of the log
send request | confirm receipt | get ticket | pay ticket | visit concert

10 traces
1.00% of the log
send request | get ticket | pay ticket | confirm receipt | visit concert

4 traces
0.40% of the log
send request | get ticket | confirm receipt | pay ticket | visit concert
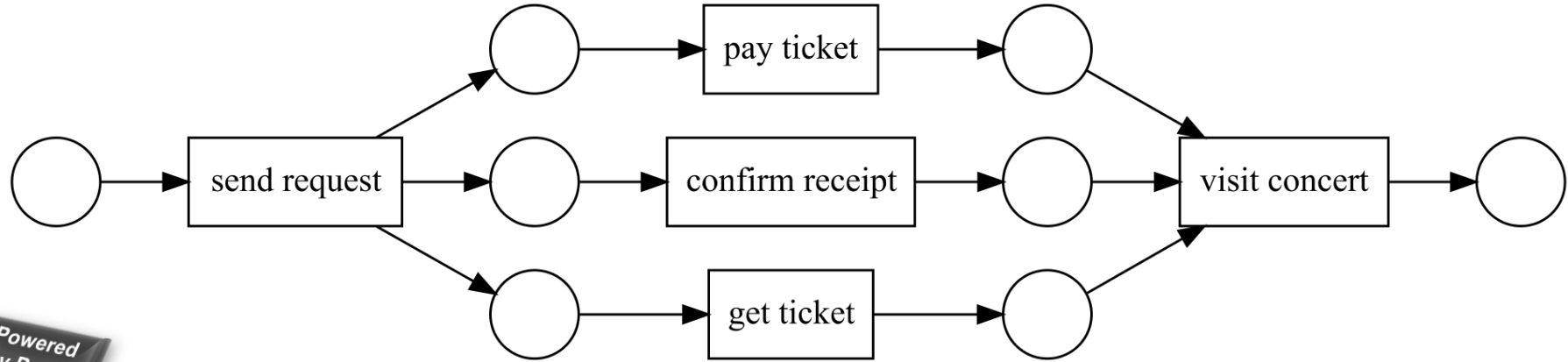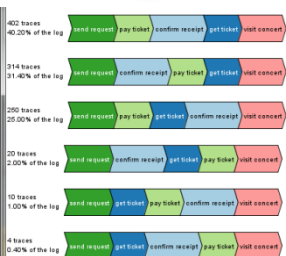
**1000 cases and 5000 events**



Powered by ProM

**Inductive miner, ILP miner, Alpha miner, etc.**
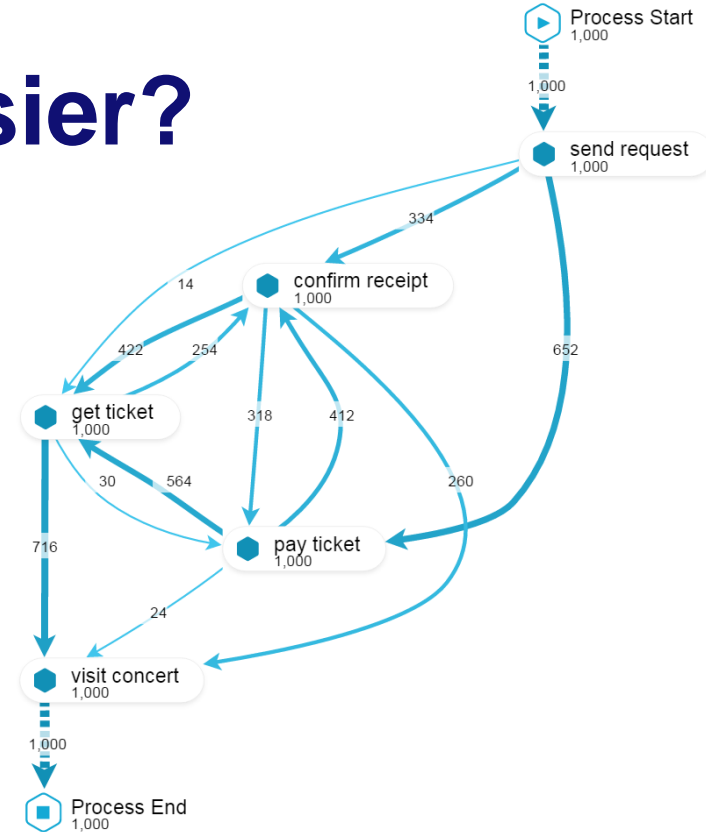
TU/e

# Model with 3 concurrent activities



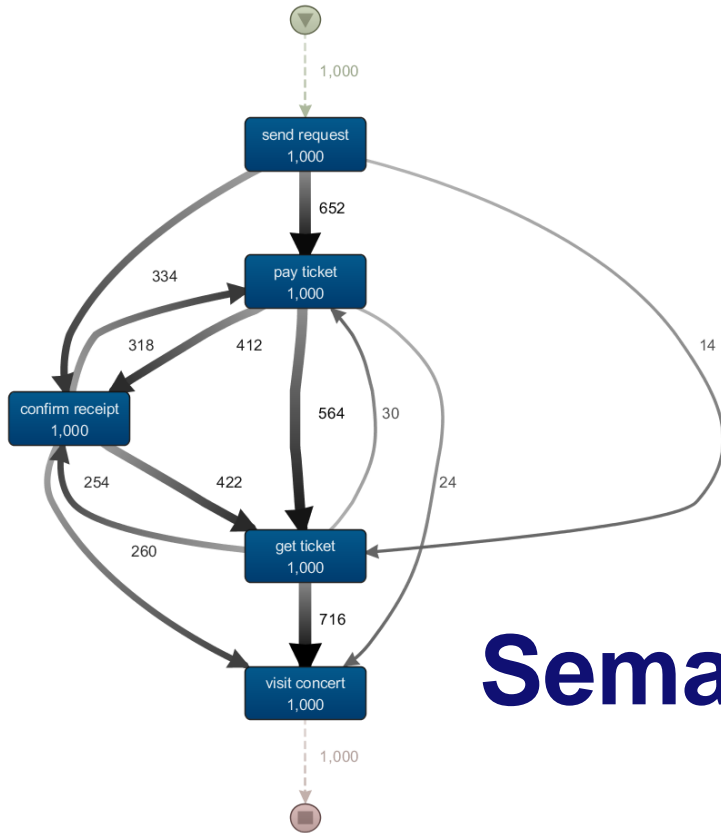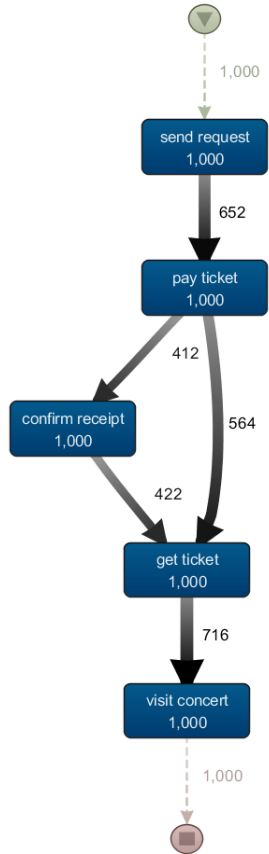# Too difficult ☺?

TU/e

# Concurrency & Semantics
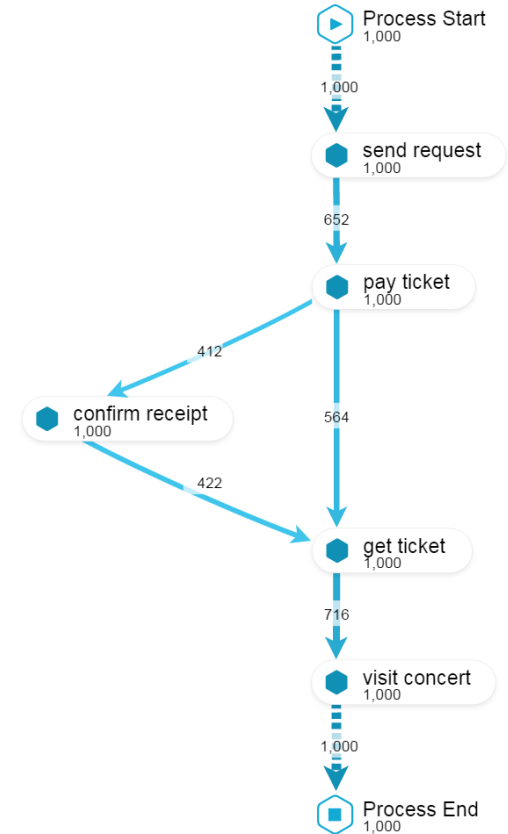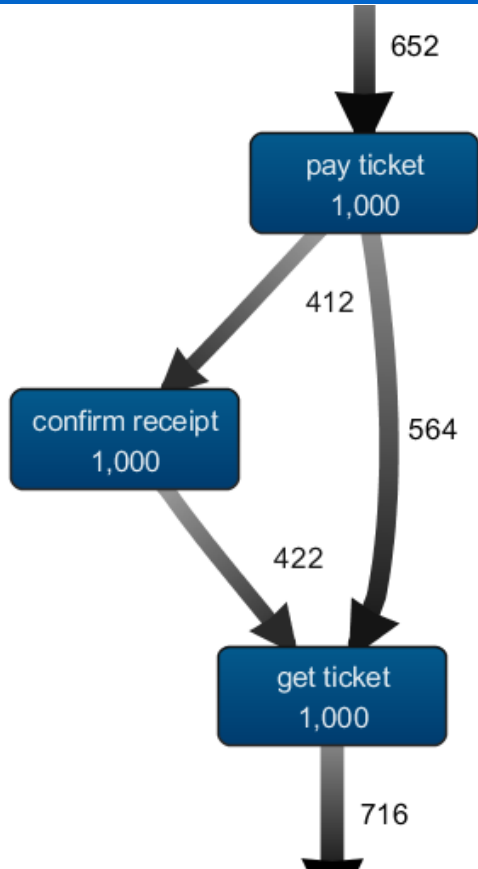


**Much easier?**

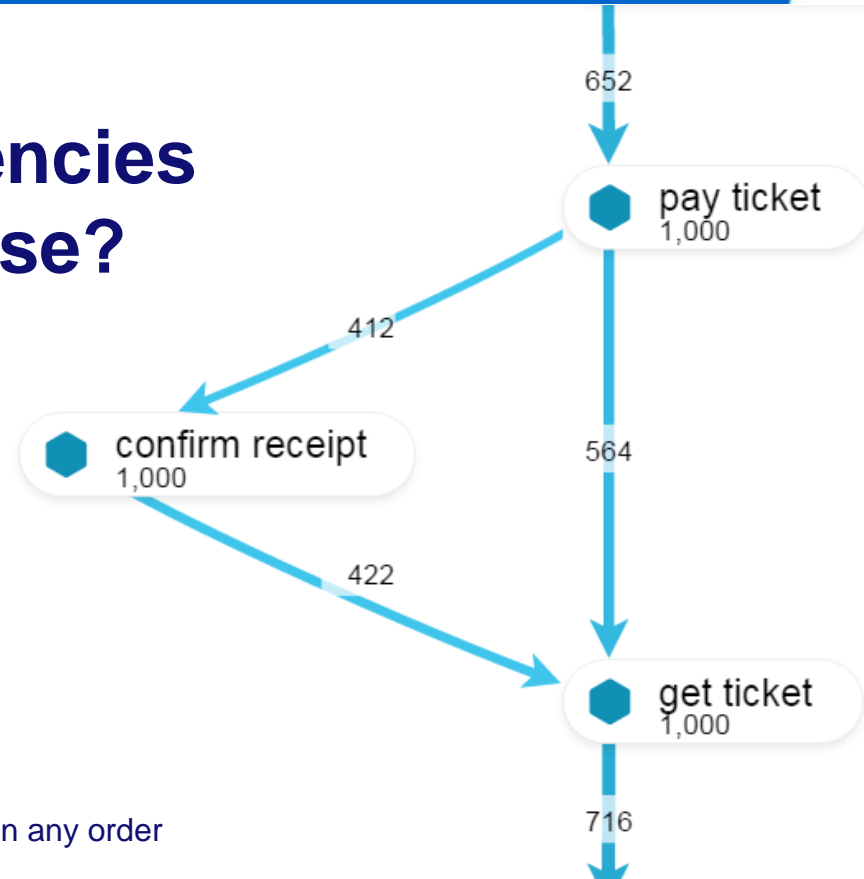**Semantics?**

# Concurrency & Semantics



**Do frequencies make sense?**
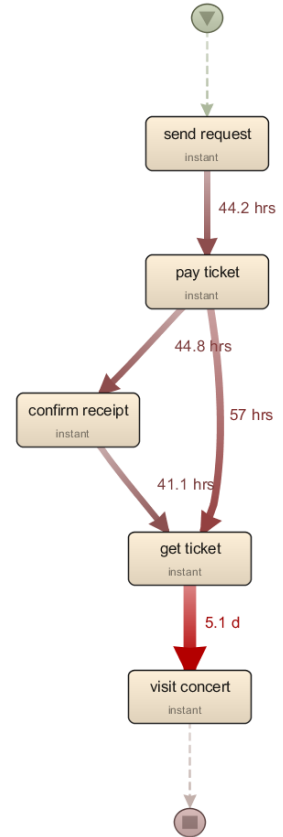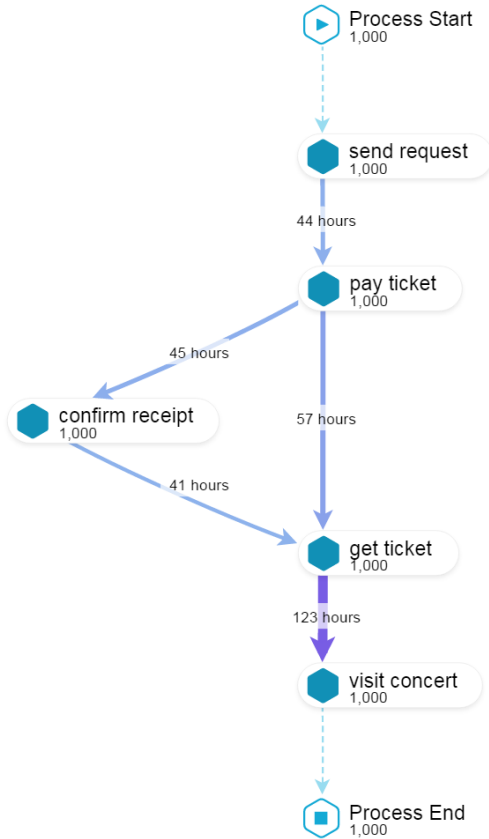
# Concurrency & Semantics



**Do frequencies make sense?**

➢ do not add up
➢ were performed in any order
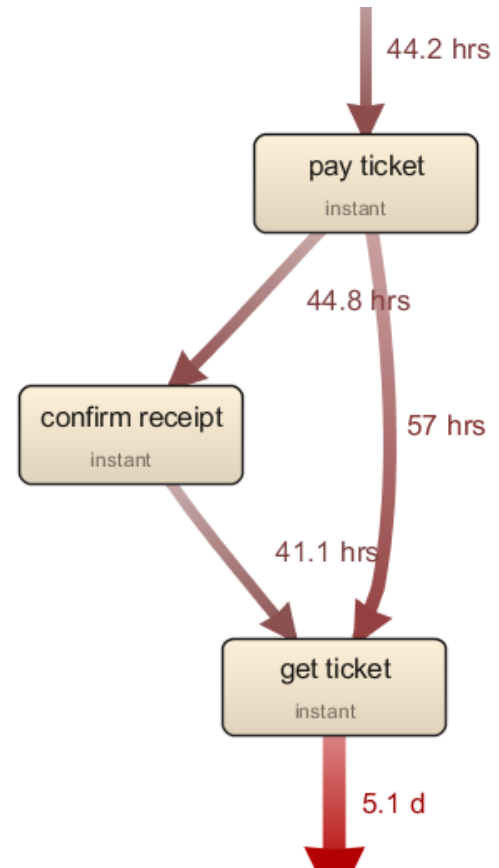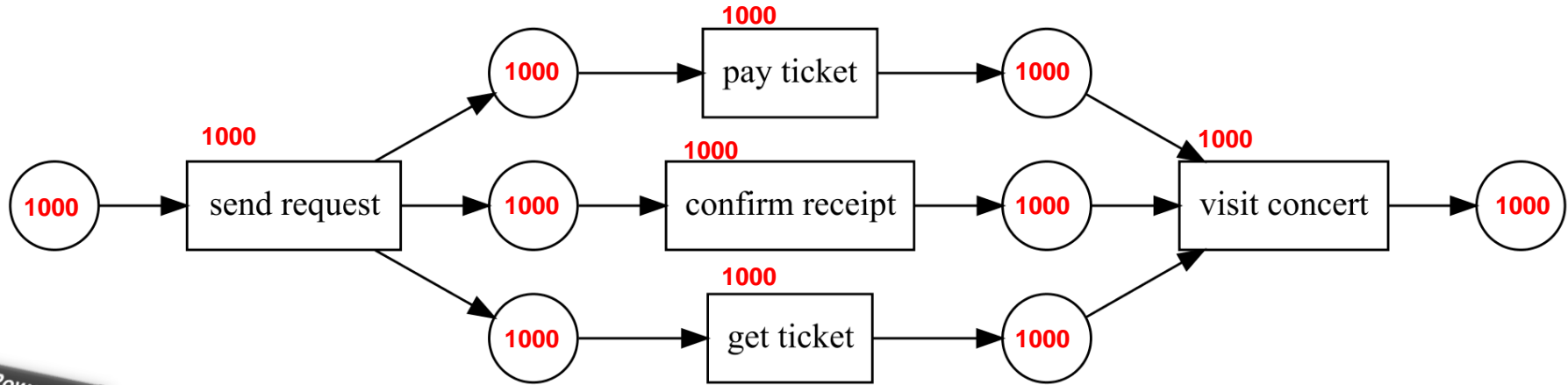
# Concurrency & Semantics



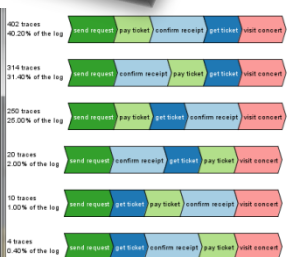**Do times make sense?**

# Concurrency & Semantics



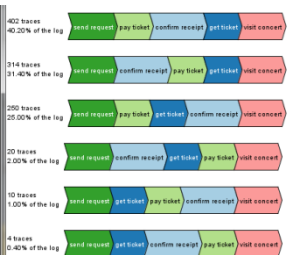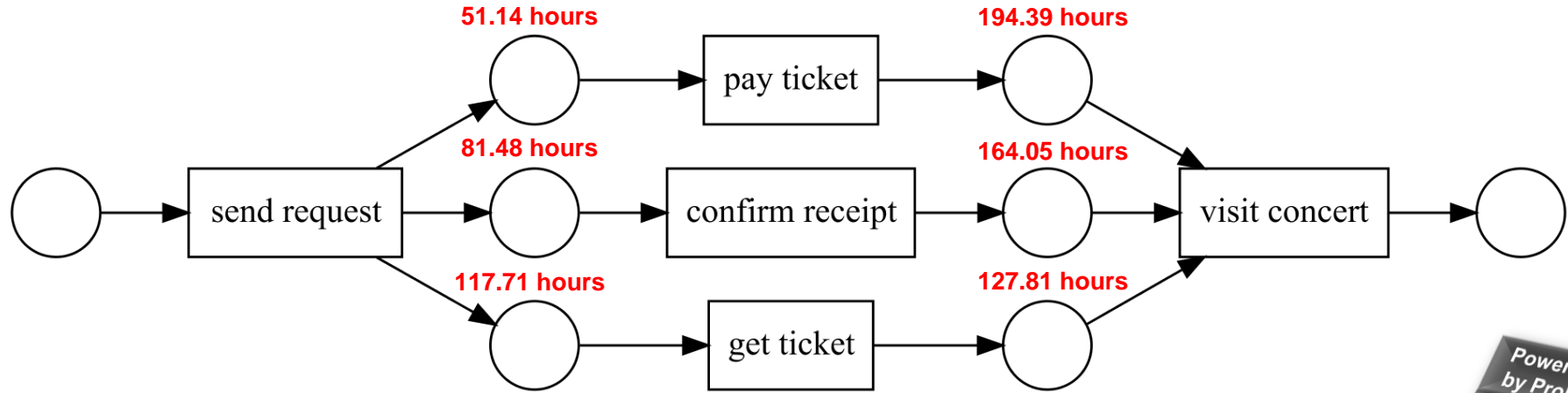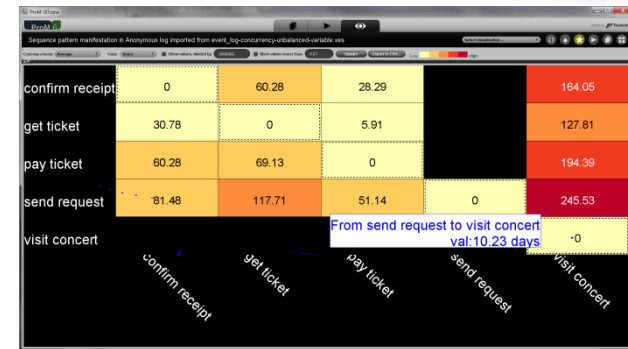**Do times make sense?**

# Concurrency & Semantics



**Correct numbers**

# Concurrency & Semantics



**51.14 hours** **194.39 hours**

pay ticket

**81.48 hours** **164.05 hours**

send request   confirm receipt   visit concert

**117.71 hours** **127.81 hours**

get ticket

Powered by ProM

## Correct times

| | confirm receipt | get ticket | pay ticket | send request | visit concert |
|---|---|---|---|---|---|
| confirm receipt | 0 | 60.28 | 28.29 | | 164.05 |
| get ticket | 30.78 | 0 | 5.91 | | 127.81 |
| pay ticket | 60.28 | 69.13 | 0 | | 194.39 |
| send request | 81.48 | 117.71 | 51.14 | 0 | 245.53 |
| visit concert | | | | | -0 |

From send request to visit concert
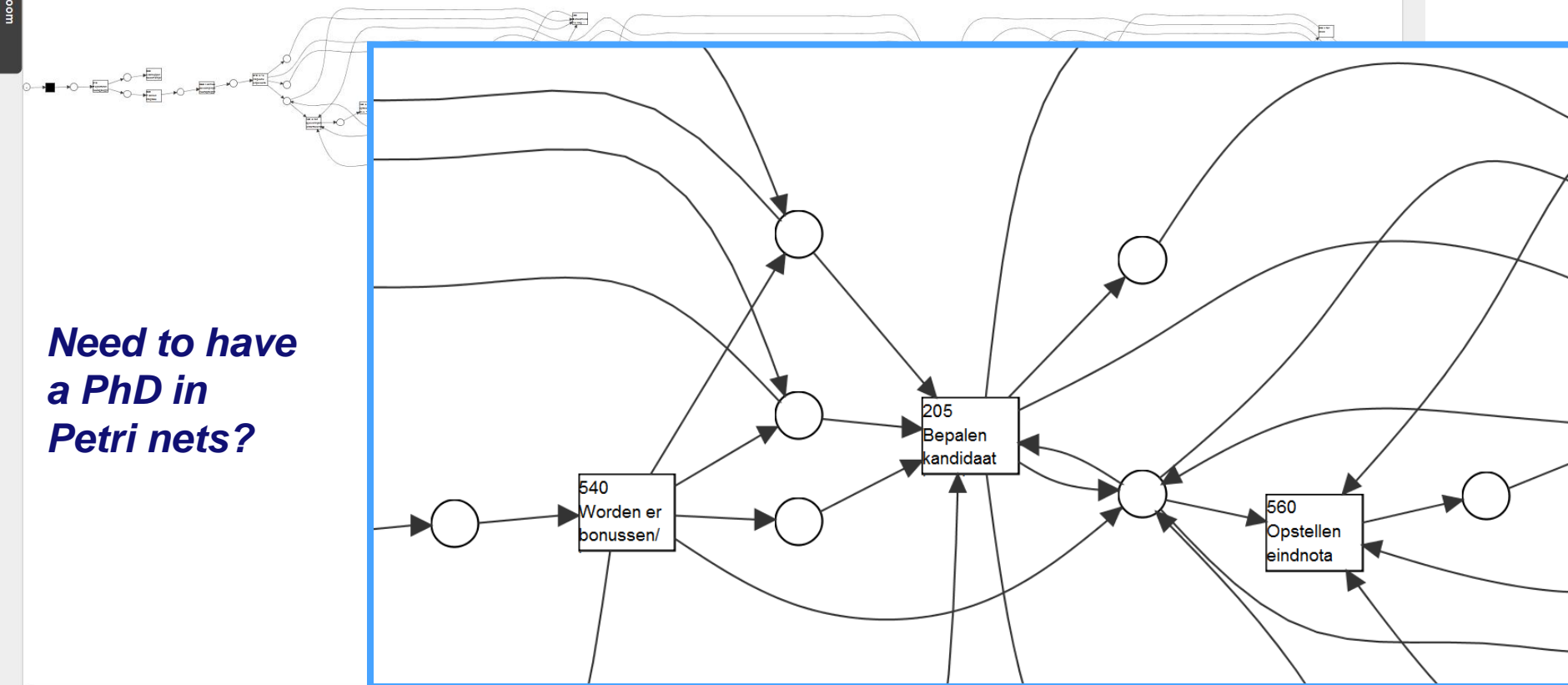val:10.23 days

But ...

faking confidence
(even with fitness = 1.0)
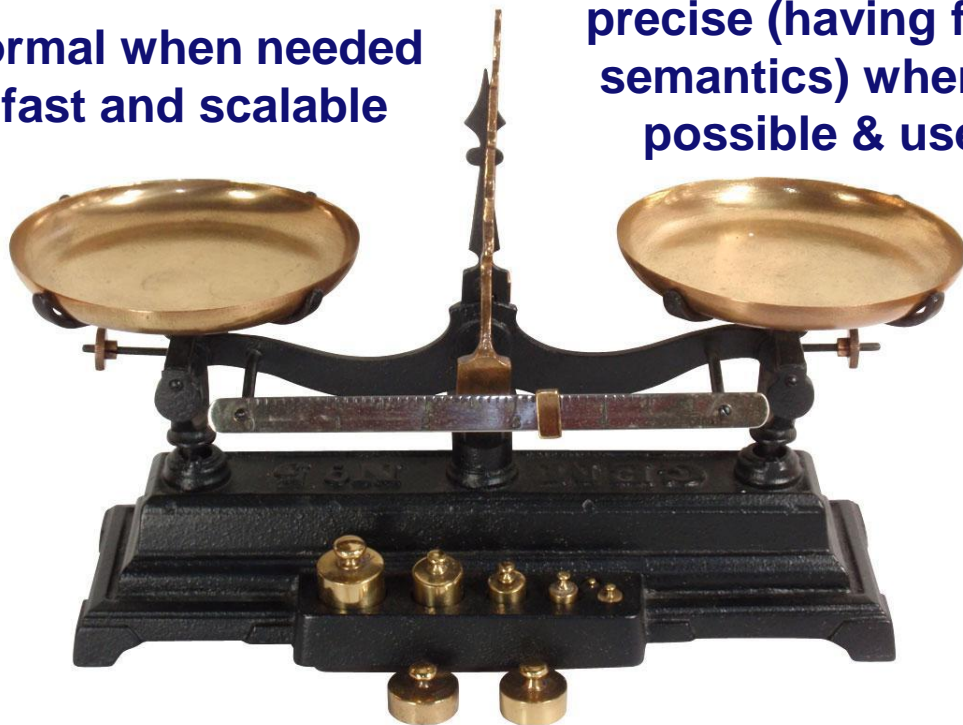
Need to have a PhD in Petri nets?

# How to combine the best of both worlds?

**informal when needed & fast and scalable**

**precise (having formal semantics) whenever possible & useful**

- **commercial tools**
- **heuristic miner**
- **fuzzy miner**
- **etc.**

- **basic inductive miner**
- **ILP miner**
- **other region-based approaches**
- **etc.**

TU/e

# Idea: Hybrid process models

Joint work with Riccardo De Masellis, Chiara Di Francescomarino, Chiara Ghidini

# Vagueness in models of socio-technical systems

THOMAS HERRMANN and KAI-UWE LOSER

Special Field of Informatics and Society, Department of Computer Science, University of Dortmund, FB Informatik, D-44221 Dortmund, Germany; e-Mail: {herrmann, loser}@iug.cs.uni-dortmund.de

**Abstract.** This article presents graphical modeling concepts, especially for the modeling of socio-technical processes. This requires the representation of those parts of knowledge which cannot be stated definitely and have to be modeled vaguely. The presented modeling concepts allow the extension of existing graphical and textual modeling methods to model facts without making unnecessary and unwelcome commitments about not already completed knowledge. In the same way it also allows the modeling of facts which cannot be modeled completely, like aspects of social systems comprising of cooperation and communication. A special modeling notation (SeeMe) is used to present the concepts. A systematic differentiation of vagueness shows the alternative ways for modelers to express vague facts. Expressing undetermined decisions is another element of vague modeling in SeeMe.

## 1. Introduction

# Semistructured models are surprisingly useful for user-centered design

Thomas Herrmann, Marcel Hoffmann, Kai-Uwe Loser, Klaus Moysich

*Informatics & Society, Dept. Computer Science, University of Dortmund, Dortmund, Germany {herrmann, hoffmann, loser, moysich}@iug.cs.uni-dortmund.de*

**Abstract.** Diagrammatic representations are commonly accepted as valuable tools in requirements engineering and systems design. However, the most prominent techniques, are not sufficient for requirements negotiation with users because they focus on the design of technical systems. In user-centered design of socio-technical systems there is a strong demand for models which integrate different viewpoints. We bel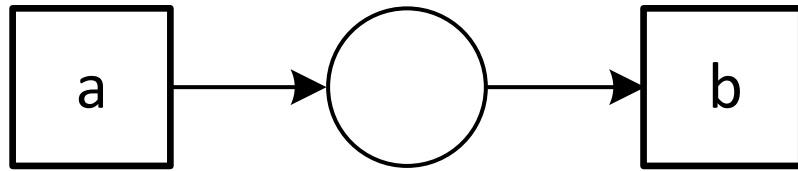ieve that appropriate semi-formal diagramming techniques can facilitate the negotiation of the design, especially when they are combined with additional representations. Therefore we have designed a notation that supports the generation of integrated models of organizational, social, and technical structures, e.g. business processes, social relations and dependencies among protagonists, resources, work-objects, and software functionality. SeeMe, the diagramming-technique for modeling semistructured socio-technical systems moreover provides special concepts for the representation of vagueness, incompleteness, and contradictions that are inherent to user requirements. In this paper we present a first evaluation of the SeeMe-diagramming technique. The results are drawn from four different case studies. We briefly introduce the main features of the SeeMe Diagramming technique and sub-

**People do not hate Petri nets (BPMN, etc.): they hate to be <u>precise</u> (when …)!**
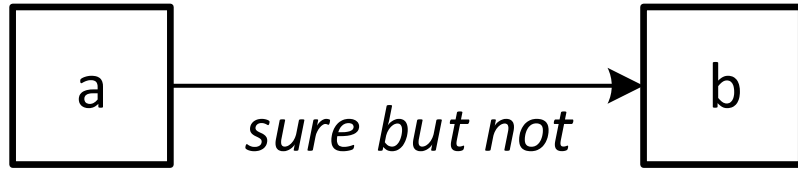
**Vagueness can be a feature!**

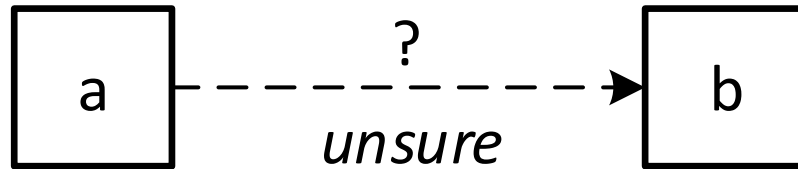**Semi-structuredness can be deliberate!**

Wil van der Aalst, Riccardo De Masellis, Chiara Di Francescomarino, Chiara Ghidini:
Learning Hybrid Process Models From Events: Process Discovery Without Faking Confidence.

July 2016

# Hybrid Petri nets have three types of arcs



*sure and precise*

*sure but not precise*

? *unsure*
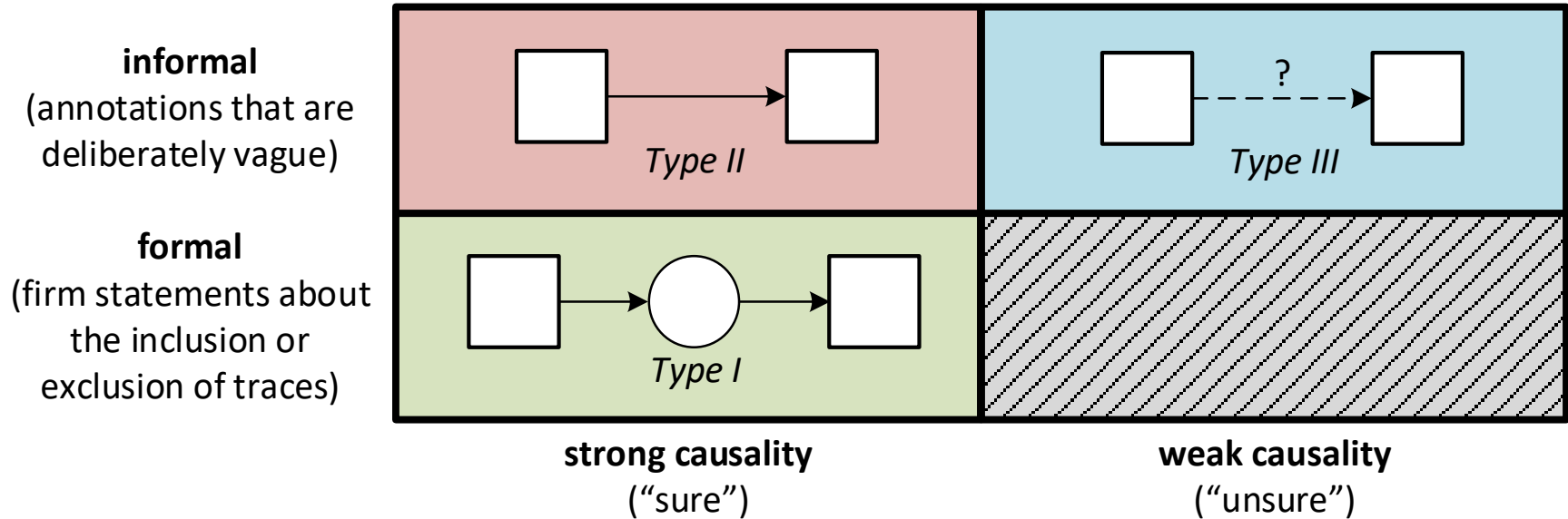
**strong causality**

determined based on thresholds

**weak causality**

determined based on thresholds

**logic can be captured***

determined based on thresholds

**logic cannot be captured***

**\* = easily / of a certain quality / within the representational bias**

# Hybrid Petri nets have three types of arcs

# Phase 0: Get data

# Phase 1: Learn a Causal Graph

# Phase 2: Learn a Hybrid Petri Net

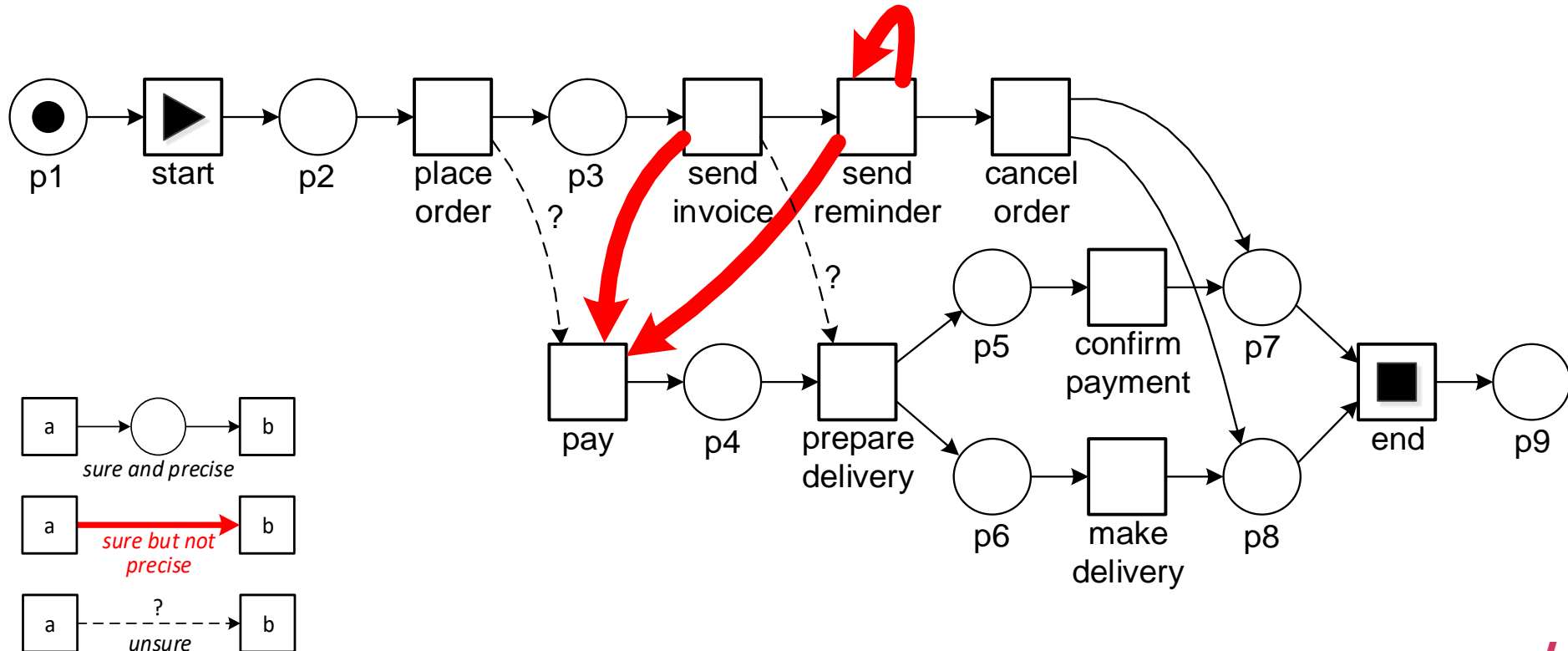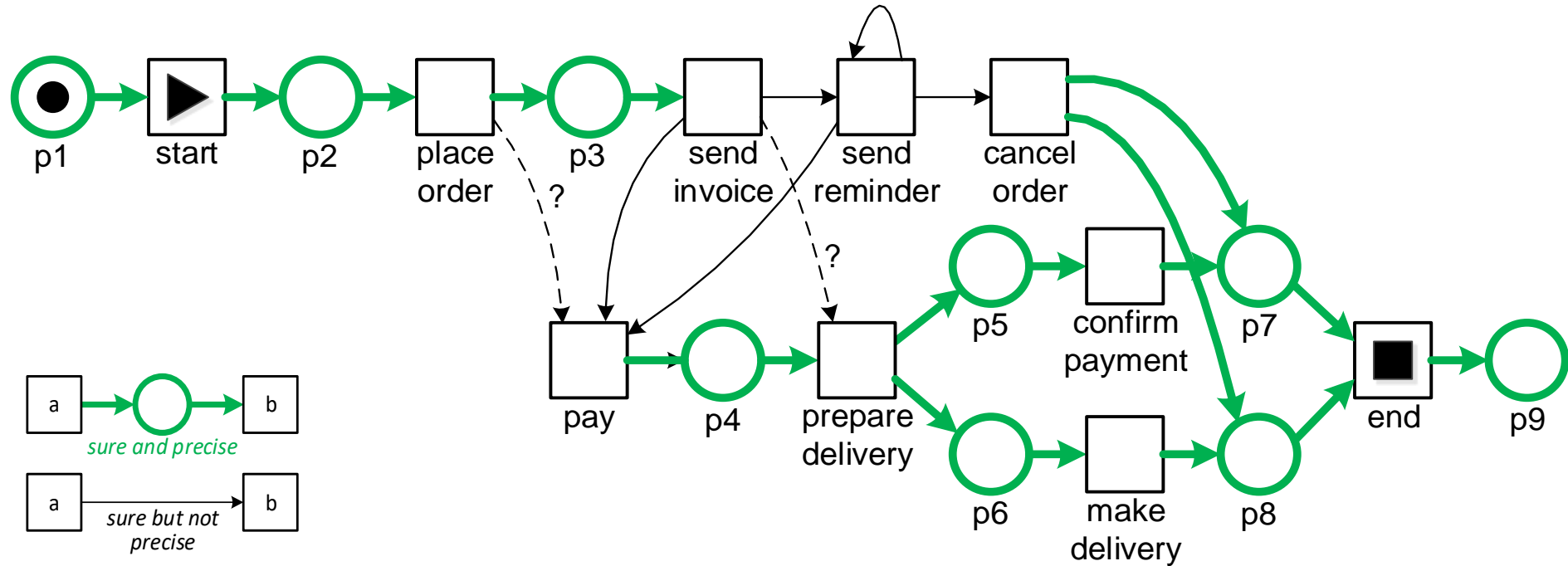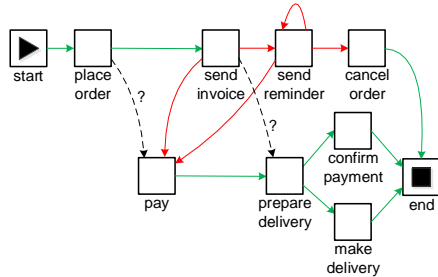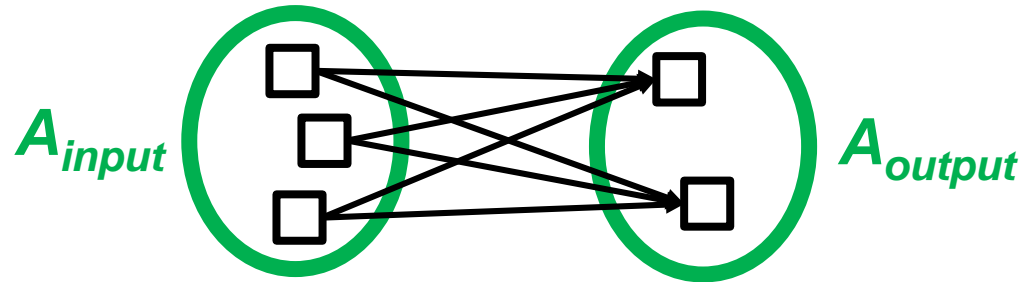# Phase 2: Learn a Hybrid Petri Net

# Phase 2: Learn a Hybrid Petri Net

# Phase 2: Learn a Hybrid Petri Net

# Phase 2: Learn a Hybrid Petri Net



- **How? ►Generate candidate places**
- **A candidate place is characterized by two sets of activities $A_{input}$ and $A_{output}$ such that sure arcs are connecting any activity in $A_{input}$ to any activity in $A_{output}$**
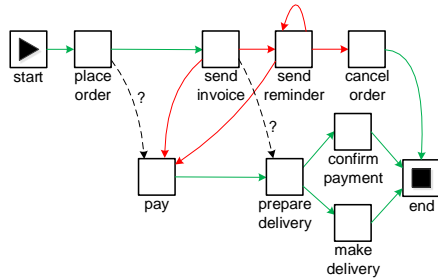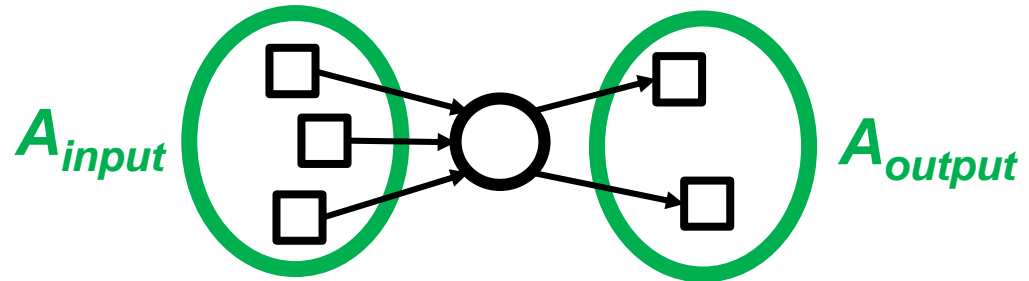
# Phase 2: Learn a Hybrid Petri Net



- **How? ►Generate candidate places**
- **A candidate place is characterized by two sets of activities $A_{input}$ and $A_{output}$ such that sure arcs are connecting any activity in $A_{input}$ to any activity in $A_{output}$**
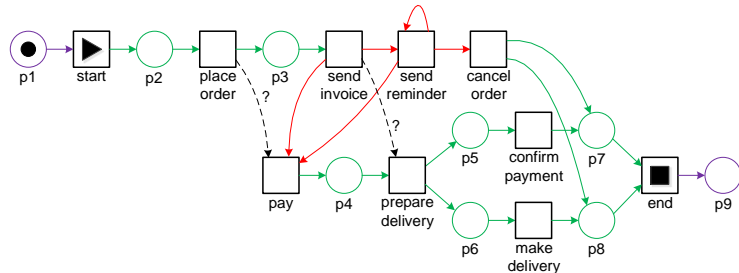
# Phase 2: Learn a Hybrid Petri Net



- **Determine the quality of each candidate place $p=(A_{input}, A_{output})$**



$A_{input}$        $A_{output}$

- **Ideally: start empty, finish empty, not negative (compare ILP miner)**

# Phase 2: Learn a Hybrid Petri Net



- **Three possible scoring functions:**
  a) **Fraction of cases perfectly fitting**
  b) **Fraction of relevant cases perfectly fitting**
  c) **Global score (extremely efficient)**

# Generate candidate places and evaluate



({co},{e})
({cp},{e})
({md},{e})
({co,cp},{e})
({co,md},{e})
({cp,md},{e})
({co,cp,md},{e})

# Generate candidate places and evaluate



({co},{e})
({cp},{e})
({md},{e})
({co,cp},{e})
({co,md},{e})
({cp,md},{e})
({co,cp,md},{e})

# ProM Implementation

# ProM Implementation



sliders

sure arc

unsure arc

fuzzy causal graph

fuzzy Petri net

# Parameters

- **Threshold for activity frequency ($t_{freq}$)**
- **Parameters used to compute strength of relations taking into account concurrency and loops ($c$ and $w$)**
- **Thresholds for strong and weak causalities ($t_{R_S}$ and $t_{R_W}$ with $t_{R_S} > t_{R_W}$)**
- **Threshold for place quality ($t_{replay}$)**

TU/e

# Evaluation (see paper and technical report for more details)

| Log | $t_{freq}$ | $t_{R_S}$ | $t_{R_W}$ | $w$ | $t_{replay}$ | $|T|$ | $|P|$ | $|\widehat{F_1}|$ | $|F_2|$ | $|F_3|$ | Fitness | Precision | Time (ms) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| BPI-2011 | 343 | 0.81 | 0.80 | 0.10 | 0.80 | 38 | 6 | 4 | 200 | 6 | 0.84 | 0.04 | 11772 |
| BPI-2012 | 3926 | 0.90 | 0.89 | 0.10 | 0.80 | 14 | 8 | 7 | 20 | 1 | 0.90 | 0.26 | 12414 |
| BPI-2014 | 13985 | 0.90 | 0.90 | 0.10 | 0.80 | 10 | 5 | 3 | 13 | 0 | 0.93 | 0.54 | 21233 |
| BPI-2015 | 360 | 0.45 | 0.40 | 0.50 | 0.80 | 59 | 26 | 24 | 145 | 75 | 0.74 | 0.05 | 7055 |
| BPI-2016 | 445 | 0.50 | 0.50 | 0.10 | 0.80 | 12 | 2 | 0 | 31 | 0 | 0.83 | 0.10 | 31428 |
| BPI-2017 | 9453 | 0.51 | 0.50 | 0.50 | 0.80 | 22 | 8 | 7 | 36 | 12 | 0.95 | 0.12 | 24772 |

- **Behaves as expected, e.g., when $t_{replay}$ goes up fitness goes up and precision goes down**

arxiv.org/abs/1703.06125

TU/e

# Performance

- **Good, but room for improvement**
- **Smartly pruning the set of candidate places (avoid conflicting or less informative places)**
- **Lazy place evaluation**
- **Distribution/decomposition using e.g. Spark (see joint work with Long Cheng and Boudewijn van Dongen in a slightly different setting)**

TU/e

# Evaluation is not easy

# But, the need is obvious

# But, the need is obvious

# Conclusion

**Learning Hybrid Process Models from Events**

**Process Discovery Without Faking Confidence**

TU/e