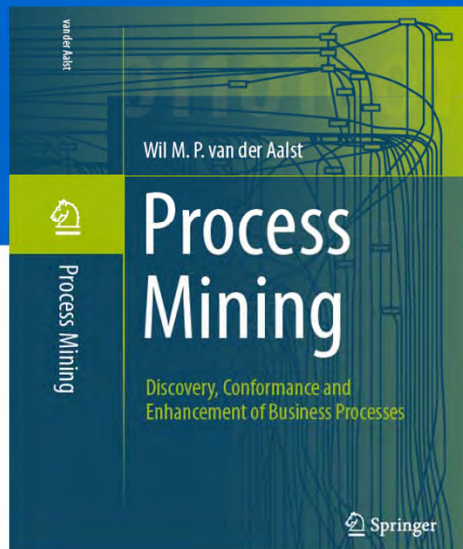# Discovering Concurrency
## Learning (Business) Process Models from Examples

**Invited Talk CONCUR 2011, 8-9-2011, Aachen.**

**prof.dr.ir. Wil van der Aalst**
**www.processmining.org**

van der Aalst

Wil M. P. van der Aalst

Process Mining

## Process Mining
Discovery, Conformance and
Enhancement of Business Processes

Springer

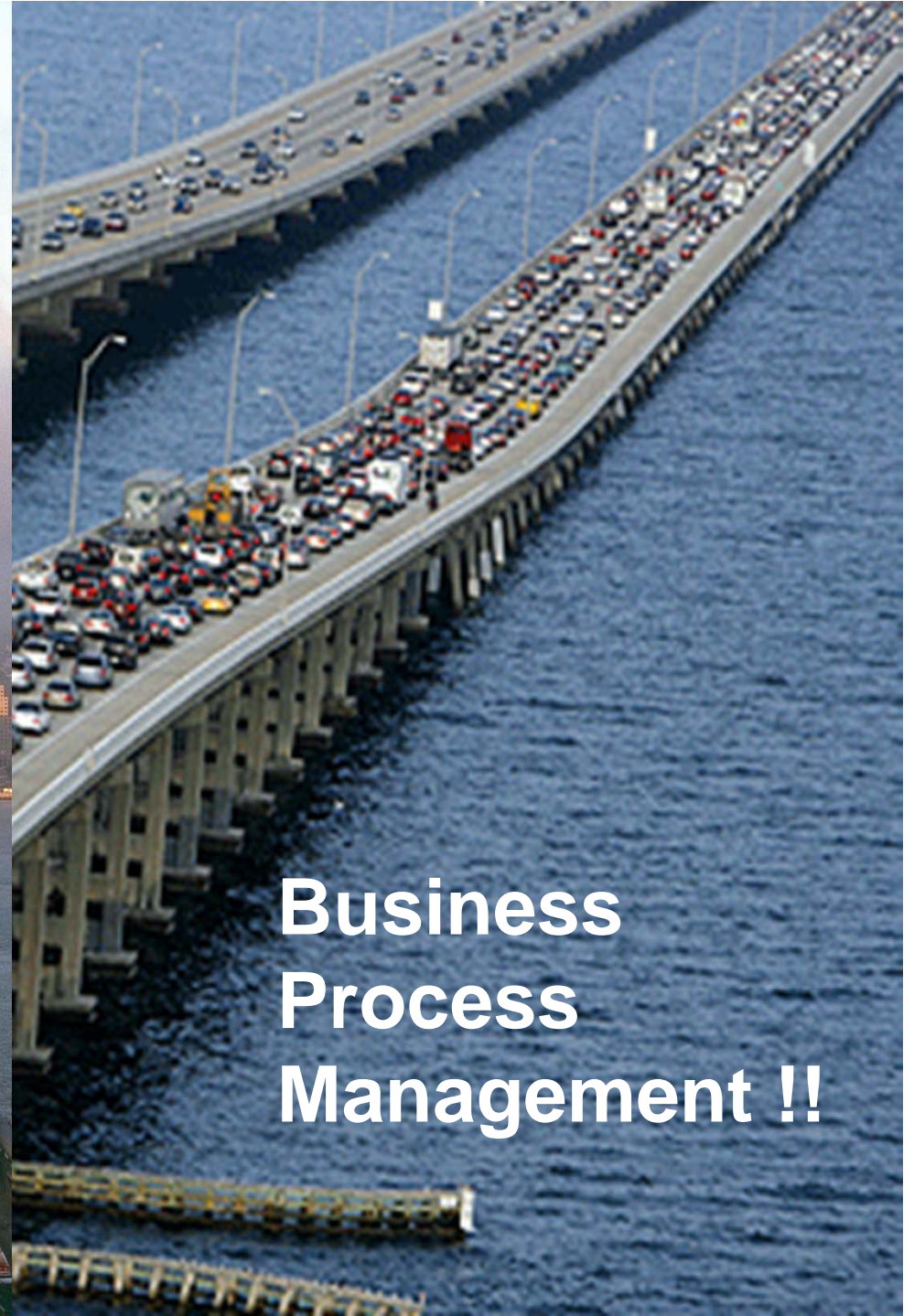TU/e
Technische Universiteit
**Eindhoven**
University of Technology
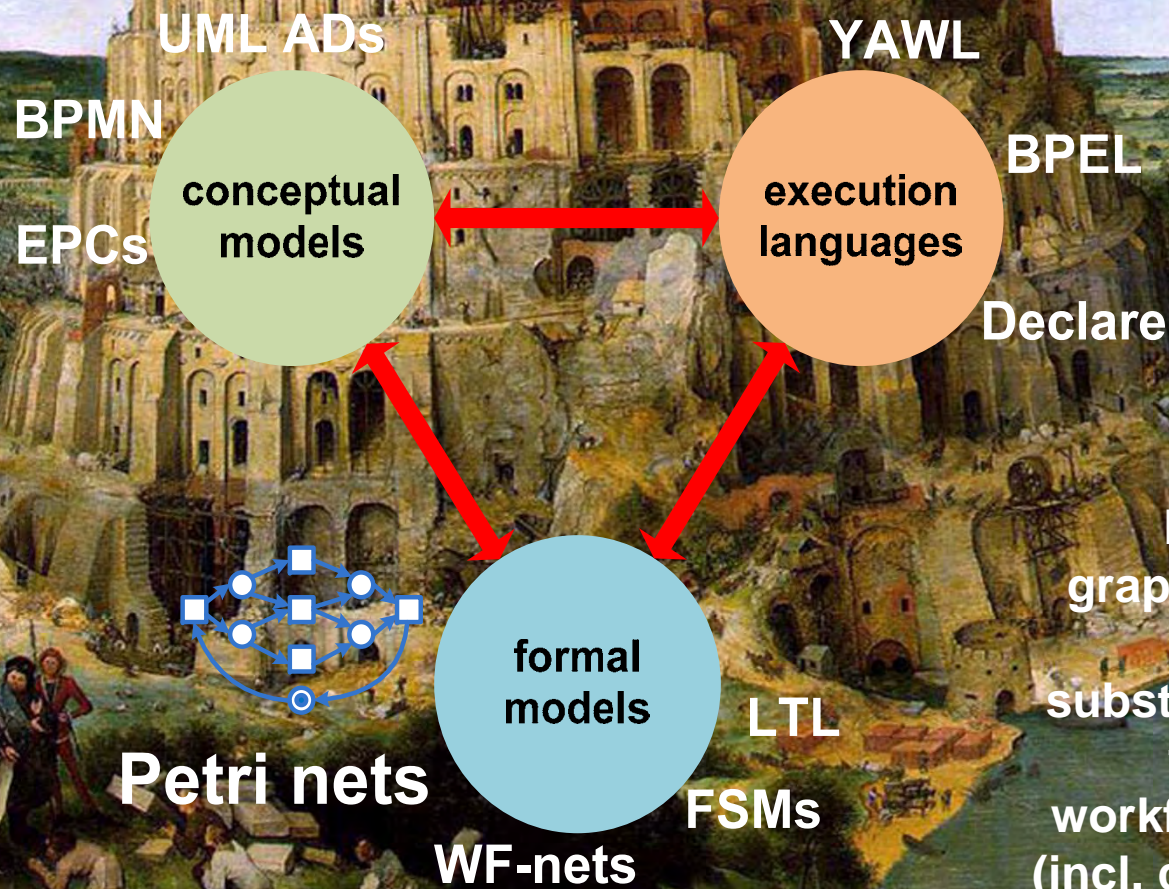
**Where innovation starts**

Business Process Management ?

**Business Process Management !!**

# Types of Process Models Used in BPM

UML ADs

YAWL

BPMN

BPEL

**conceptual models**

**execution languages**

EPCs

Declare

**Petri nets**

**formal models**

LTL

FSMs

WF-nets

Emphasis on graphical models supporting a substantial part of the so-called workflow patterns (incl. concurrency)

# Classical Challenges in BPM

- **Verification (cf. soundness problem in WF-nets)**
- **Performance analysis (e.g., simulation)**
- **Converting models into running systems**
- **Providing flexibility without loosing control**
- **…**

**Many problems have been "solved": the real problem is adoption in practice!**

# A more interesting challenge: Dealing with variability



- **Variants of the same process exist in various domains, e.g., Dutch Municipalities, Hertz, Suncorp, Salesforce, Easychair, etc.**
- **Configurable process models to generate concrete processes, cf. C-YAWL.**
- **Merging process models is a challenging problem.**
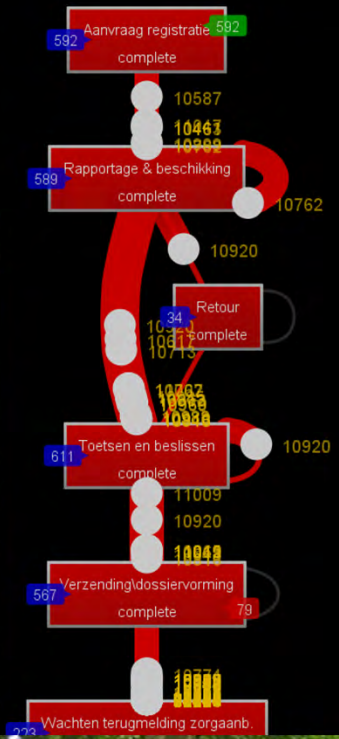- **Model similarity (rather than equivalence).**

So we are interested in processes …

… but most of us do not study them!

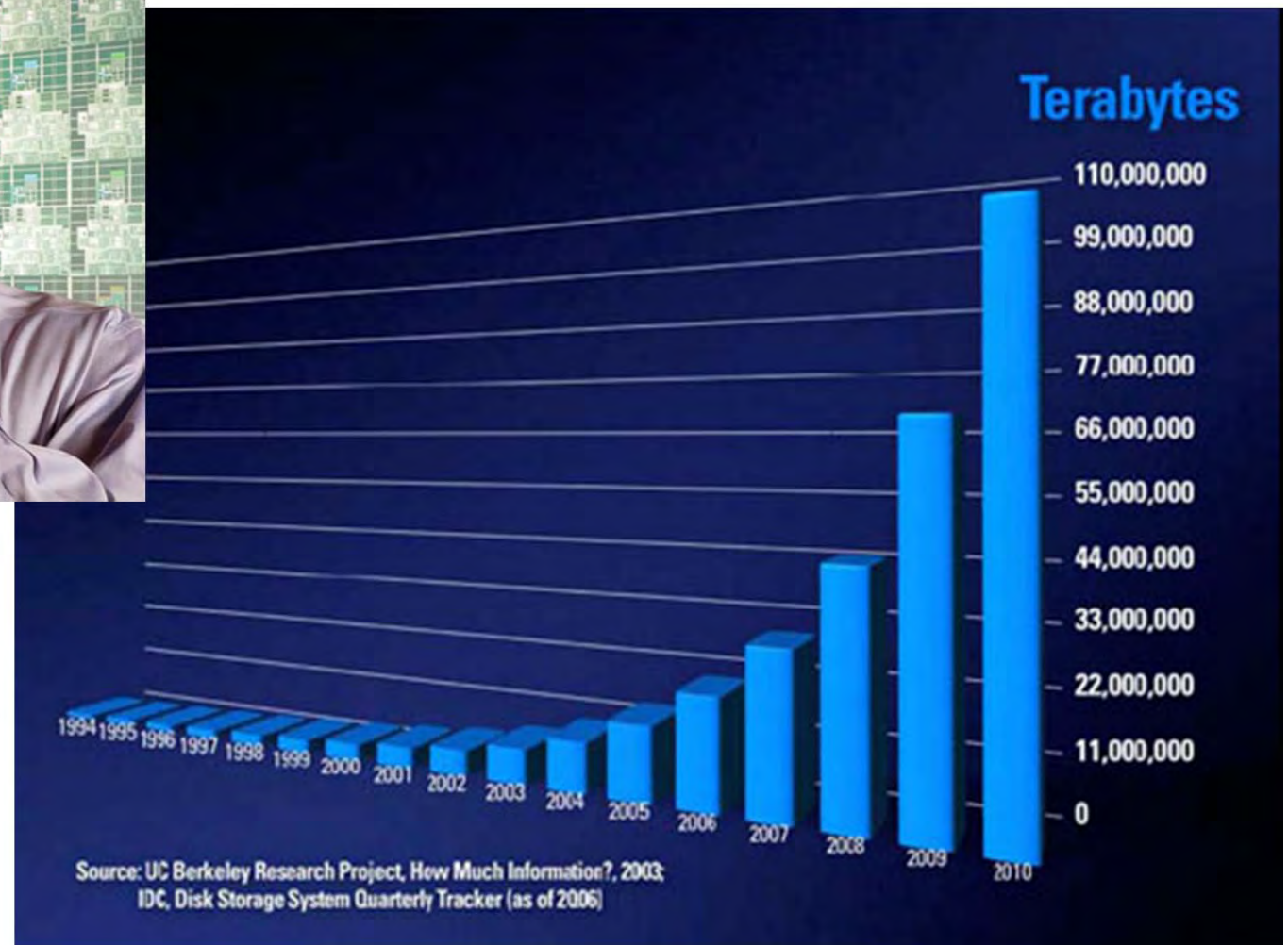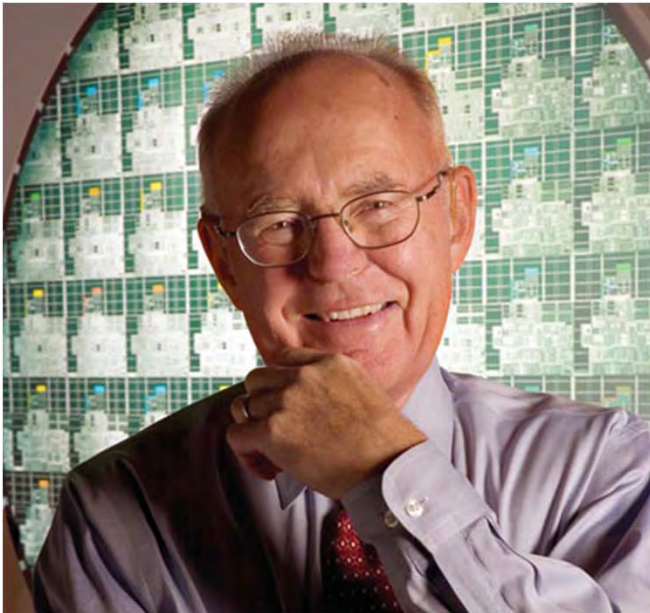# Desire lines in process models

# Data explosion



Terabytes

110,000,000
99,000,000
88,000,000
77,000,000
66,000,000
55,000,000
44,000,000
33,000,000
22,000,000
11,000,000
0

1994 1995 1996 1997 1998 1999 2000 2001 2002 2003 2004 2005 2006 2007 2008 2009 2010

Source: UC Berkeley Research Project, How Much Information?, 2003;
IDC, Disk Storage System Quarterly Tracker (as of 2006)

**Process Mining =**

**Event Data + Processes**

**Data Mining + Process Analysis**

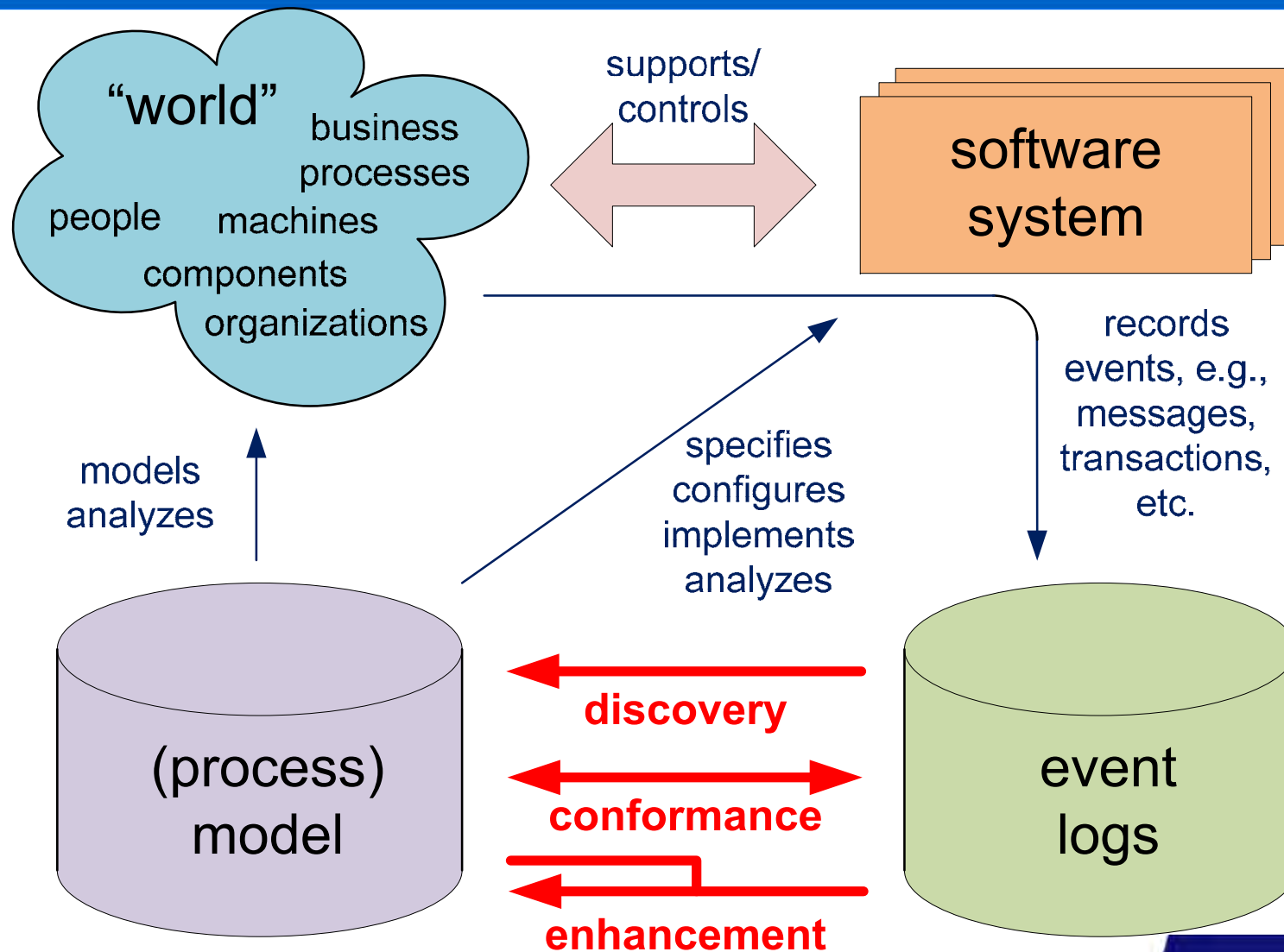**Machine Learning + Formal Methods**

# Process Mining



- **Process discovery: "What is really happening?"**
- **Conformance checking: "Do we do what was agreed upon?"**
- **Performance analysis: "Where are the bottlenecks?"**
- **Process prediction: "Will this case be late?"**
- **Process improvement: "How to redesign this process?"**
- **Etc.**

# Process Mining

# Starting point: event log

| 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 | ... |

| case id | event id | properties | | | | |
|---|---|---|---|---|---|---|
| | | timestamp | activity | resource | cost | ... |
| 3 | 35654521 | 30-12-2010:14.32 | register request | | | |
| | 35654522 | 30-12-2010:15.06 | examine casually | | | |
| | 35654524 | 30-12-2010:16.34 | check ticket | | | |
| | 35654525 | 06-01-2011:09.18 | decide | | | |
| | 35654526 | 06-01-2011:12.18 | reinitiate request | | | |
| | 35654527 | 06-01-2011:13.06 | examine thoroughly | | | |
| | 35654530 | 08-01-2011:11.43 | check ticket | | | |
| | 35654531 | 09-01-2011:09.55 | decide | | | |
| | 35654533 | 15-01-2011:10.45 | pay compensation | | | |
| 4 | 35654641 | 06-01-2011:15.02 | register request | | | |
| | 35654643 | 07-01-2011:12.06 | check ticket | | | |
| | 35654644 | 08-01-2011:14.43 | examine thoroughly | | | |
| | 35654645 | 09-01-2011:12.02 | decide | | | |
| | 35654647 | 12-01-2011:15.44 | reject request | | | |
| 5 | 35654711 | 06-01-2011:09.02 | register request | | | |
| | 35654712 | 07-01-2011:10.16 | examine casually | | | |
| | 35654714 | 08-01-2011:11.22 | check ticket | | | |
| | 35654715 | 10-01-2011:13.28 | decide | | | |
| | 35654716 | 11-01-2011:16.18 | reinitiate request | | | |
| | 35654718 | 14-01-2011:14.33 | check ticket | | | |
| | 35654719 | 16-01-2011:15.50 | examine casually | | | |
| | 35654720 | 19-01-2011:11.18 | decide | Sara | 200 | ... |
| | 35654721 | 20-01-2011:12.48 | reinitiate request | Sara | 200 | ... |
| | 35654722 | 21-01-2011:09.06 | examine casually | Sue | 400 | ... |
| | 35654724 | 21-01-2011:11.34 | check ticket | Pete | 100 | ... |
| | 35654725 | 23-01-2011:13.12 | decide | Sara | 200 | ... |
| | 35654726 | 24-01-2011:14.56 | reject request | Mike | 200 | ... |
| 6 | 35654871 | 06-01-2011:15.02 | register request | Mike | 50 | ... |
| | 35654873 | 06-01-2011:16.06 | examine casually | Ellen | 400 | ... |
| | 35654874 | 07-01-2011:16.22 | check ticket | Mike | 100 | ... |
| | 35654875 | 07-01-2011:16.52 | decide | Sara | 200 | ... |
| | 35654877 | 16-01-2011:11.47 | pay compensation | Mike | 200 | ... |
| ... | ... | ... | ... | ... | ... | ... |

**XES, MXML, SA-MXML, CSV, etc.**

# Simplified event log

| case id | event id | timestamp | activity | resource | |
|---------|----------|-----------|----------|----------|---|
| 1 | 35654423 | 30-12-2010:11.02 | register request | Pete | |
| | 35654424 | 31-12-2010:10.06 | examine thoroughly | Sue | |
| | 35654425 | 05-01-2011:15.12 | check ticket | Mike | |
| | 35654426 | 06-01-2011:11.18 | decide | Sara | |
| | 35654427 | 07-01-2011:14.24 | reject request | Pete | |
| 2 | 35654483 | 30-12-2010:11.32 | register request | Mike | |
| | 35654485 | 30-12-2010:12.12 | check ticket | Mike | |
| | 35654487 | 30-12-2010:14.16 | examine casually | Pete | |
| | 35654488 | 05-01-2011:11.22 | decide | Sara | |
| | 35654489 | 08-01-2011:12.05 | pay compensation | Ellen | |
| 3 | 35654521 | 30-12-2010:14.32 | register request | Pete | |
| | 35654522 | 30-12-2010:15.06 | examine casually | Mike | |
| | 35654524 | 30-12-2010:16.34 | check ticket | Ellen | |
| | 35654525 | 06-01-2011:09.18 | decide | Sara | |
| | 35654526 | 06-01-2011:12.18 | reinitiate request | Sara | |
| | 35654527 | 06-01-2011:13.06 | examine thoroughly | Sean | |
| | 35654530 | 08-01-2011:11.43 | check ticket | Pete | |
| | 35654531 | 09-01-2011:09.55 | decide | Sara | |
| | 35654533 | 15-01-2011:10.45 | pay compensation | Ellen | |
| 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 ... |
| 5 | 35654711 | 06-01-2011:09.02 | register request | Ellen | 50 ... |
| | 35654712 | 07-01-2011:10.16 | examine casually | Mike | 400 ... |
| | 35654714 | 08-01-2011:11.22 | check ticket | Pete | 100 ... |
| | 35654715 | 10-01-2011:13.28 | decide | Sara | 200 ... |
| | 35654716 | 11-01-2011:16.18 | reinitiate request | Sara | 200 ... |
| | 35654718 | 14-01-2011:14.33 | check ticket | Ellen | 100 ... |
| | 35654719 | 16-01-2011:15.50 | examine casually | Mike | 400 ... |
| | 35654720 | 19-01-2011:11.18 | decide | Sara | 200 ... |
| | 35654721 | 20-01-2011:12.48 | reinitiate request | Sara | 200 ... |
| | 35654722 | 21-01-2011:09.06 | examine casually | Sue | 400 ... |
| | 35654724 | 21-01-2011:11.34 | check ticket | Pete | 100 ... |
| | 35654725 | 23-01-2011:13.12 | decide | Sara | 200 ... |
| | 35654726 | 24-01-2011:14.56 | reject request | Mike | 200 ... |
| 6 | 35654871 | 06-01-2011:15.02 | register request | Mike | 50 ... |
| | 35654873 | 06-01-2011:16.06 | examine casually | Ellen | 400 ... |
| | 35654874 | 07-01-2011:16.22 | check ticket | Mike | 100 ... |
| | 35654875 | 07-01-2011:16.52 | decide | Sara | 200 ... |
| | 35654877 | 16-01-2011:11.47 | pay compensation | Mike | 200 ... |
| ... | ... | ... | ... | ... | ... |

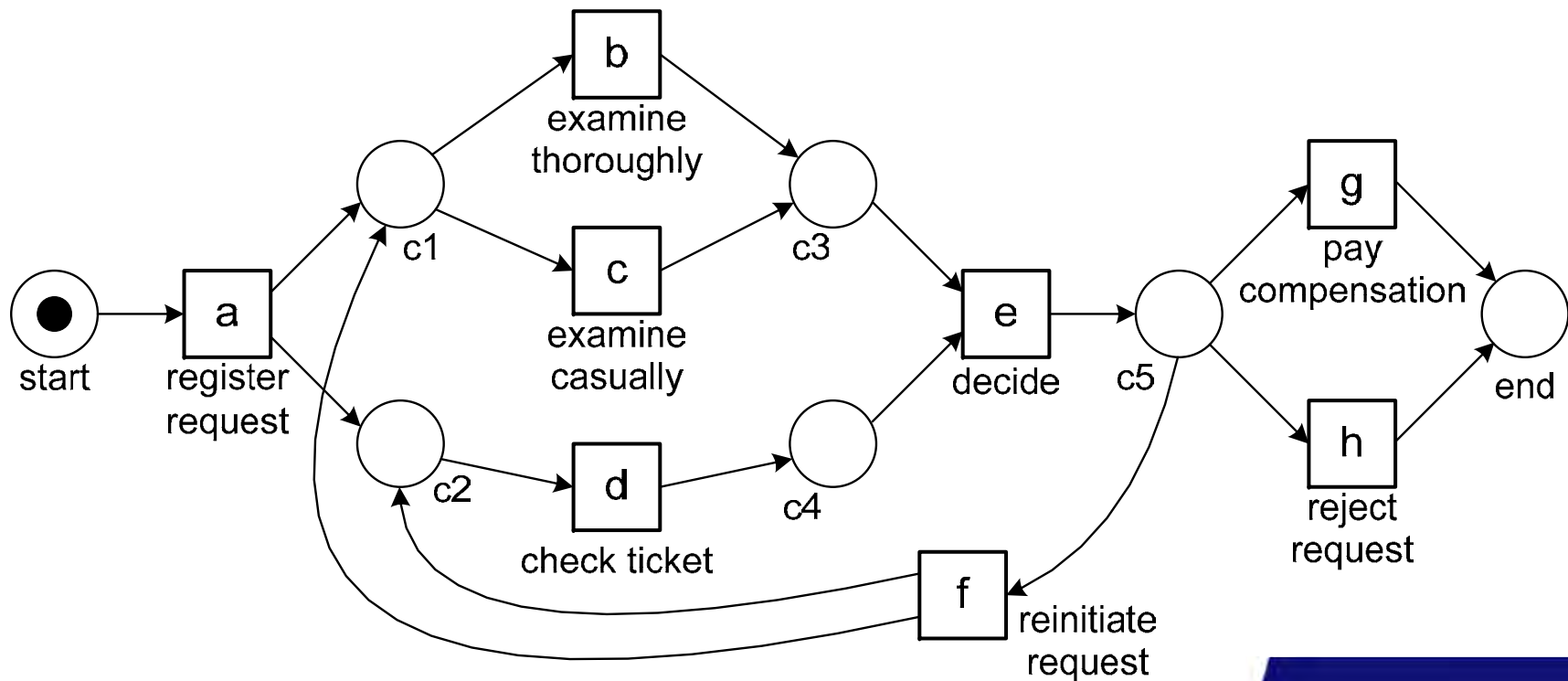| case id | trace |
|---------|-------|
| 1 | $\langle a,b,d,e,h \rangle$ |
| 2 | $\langle a,d,c,e,g \rangle$ |
| 3 | $\langle a,c,d,e,f,b,d,e,g \rangle$ |
| 4 | $\langle a,d,b,e,h \rangle$ |
| 5 | $\langle a,c,d,e,f,d,c,e,f,c,d,e,h \rangle$ |
| 6 | $\langle a,c,d,e,g \rangle$ |
| ... | ... |

a = register request,
b = examine thoroughly,
c = examine casually,
d = check ticket,
e = decide,
f = reinitiate request,
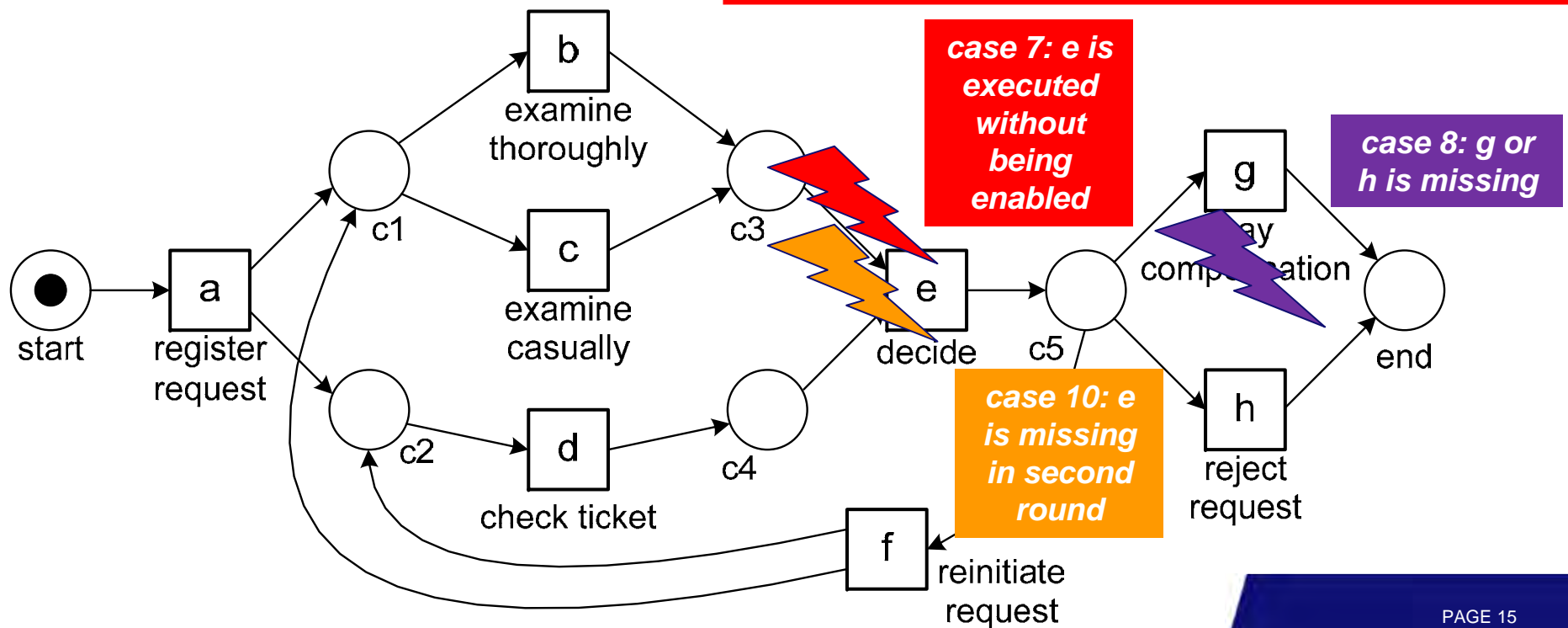g = pay compensation,
and h = reject request

# Process discovery

| case id | trace |
|---------|-------|
| 1 | $\langle a,b,d,e,h \rangle$ |
| 2 | $\langle a,d,c,e,g \rangle$ |
| 3 | $\langle a,c,d,e,f,b,d,e,g \rangle$ |
| 4 | $\langle a,d,b,e,h \rangle$ |
| 5 | $\langle a,c,d,e,f,d,c,e,f,c,d,e,h \rangle$ |
| 6 | $\langle a,c,d,e,g \rangle$ |
| … | … |

# Conformance checking

| case id | trace |
|---------|-------|
| 1 | $\langle a,b,d,e,h \rangle$ |
| 2 | $\langle a,d,c,e,g \rangle$ |
| 3 | $\langle a,c,d,e,f,b,d,e,g \rangle$ |
| 4 | $\langle a,d,b,e,h \rangle$ |
| 5 | $\langle a,c,d,e,f,d,c,e,f,c,d,e,h \rangle$ |
| 6 | $\langle a,c,d,e,g \rangle$ |
| 7 | $\langle \mathbf{a,b,e,g} \rangle$ |
| 8 | $\langle \mathbf{a,b,d,e} \rangle$ |
| 9 | $\langle a,d,c,e,f,d,c,e,f,b,d,e,h \rangle$ |
| 10 | $\langle \mathbf{a,c,d,e,f,b,d,g} \rangle$ |



**case 7: e is executed without being enabled**

**case 8: g or h is missing**

**case 10: e is missing in second round**

# Extension: Adding perspectives to model based on event log



The event log can be used to discover roles in the organization (e.g., groups of people with similar work patterns). These roles can be used to relate individuals and activities.

Performance information (e.g., the average time between two subsequent activities) can be extracted from the event log and visualized on top of the model.

Decision rules (e.g., a decision tree based on data known at the time a particular choice was made) can be learned from the event log and used to annotated decisions.

# We applied ProM in >100 organizations

- **Municipalities** (e.g., Alkmaar, Heusden, Harderwijk, etc.)
- **Government agencies** (e.g., Rijkswaterstaat, Centraal Justitieel Incasso Bureau, Justice department)
- **Insurance related agencies** (e.g., UWV)
- **Banks** (e.g., ING Bank)
- **Hospitals** (e.g., AMC hospital, Catharina hospital)
- **Multinationals** (e.g., DSM, Deloitte)
- **High-tech system manufacturers and their customers** (e.g., Philips Healthcare, ASML, Ricoh, Thales)
- **Media companies** (e.g. Winkwaves)
- **...**

# All supported by …



- **Open-source (L-GPL), cf. www.processmining.org**
- **Plug-in architecture**
- **Plug-ins cover the whole process mining spectrum and also support classical forms of process analysis**

# Process discovery



"world"
business
processes
people    machines
components
organizations

supports/
controls

software
system

models
analyzes

specifies
configures
implements
analyzes

records
events, e.g.,
messages,
transactions,
etc.

(process)
model

discovery

conformance

enhancement

event
logs

# Process Discovery Techniques (small selection)

distributed genetic mining

automata-based learning

language-based regions

heuristic mining

state-based regions

partial-order based mining

genetic mining

LTL mining

pattern-based mining

neural networks

stochastic task graphs

fuzzy mining

mining block structures

hidden Markov models

α algorithm

multi-phase mining

conformal process graph

α# algorithm

ILP mining

α++ algorithm

# Why is process discovery such a difficult problem?

- **There are no negative examples (i.e., a log shows what has happened but does not show what could not happen).**

- **Due to concurrency, loops, and choices the search space has a complex structure and the log typically contains only a fraction of all possible behaviors.**

- **There is no clear relation between the size of a model and its behavior (i.e., a smaller model may generate more or less behavior although classical analysis and evaluation methods typically assume some monotonicity property).**

# Challenge: four competing quality criteria

*"able to replay event log"*

*"Occam's razor"*

fitness

simplicity

process discovery

generalization

precision

*"not overfitting the log"*

*"not underfitting the log"*

# Example: one log four models



*"able to replay event log"*

fitness

*"Occam's razor"*

simplicity

process discovery

generalization

*"not overfitting the log"*

precision

*"not underfitting the log"*

$N_1$ : fitness = +, precision = +, generalization = +, simplicity = +

$N_2$ : fitness = -, precision = +, generalization = -, simplicity = +

$N_3$ : fitness = +, precision = -, generalization = +, simplicity = +

(all 21 variants seen in the log)

$N_4$ : fitness = +, precision = +, generalization = -, simplicity = -

| #    | trace             |
|------|-------------------|
| 455  | acdeh             |
| 191  | abdeg             |
| 177  | adceh             |
| 144  | abdeh             |
| 111  | acdeg             |
| 82   | adceg             |
| 56   | adbeh             |
| 47   | acdefdbeh         |
| 38   | adbeg             |
| 33   | acdefbdeh         |
| 14   | acdefbdeg         |
| 11   | acdefdbeg         |
| 9    | acdefcdeh         |
| 8    | adcefdbeh         |
| 5    | adcefbdeg         |
| 3    | acdefbdefdbeg     |
| 2    | adcefdbeg         |
| 2    | adcefbdefbdeg     |
| 1    | adcefdbefbdeh     |
| 1    | adbefbdefdbeg     |
| 1    | adcefdbefcdefdbeg |
| 1391 |                   |

# Model N₁



$N_1$ : fitness = +, precision = +, generalization = +, simplicity = +

| # | trace |
|---|---|
| 455 | acdeh |
| 191 | abdeg |
| 177 | adceh |
| 144 | abdeh |
| 111 | acdeg |
| 82 | adceg |
| 56 | adbeh |
| 47 | acdefdbeh |
| 38 | adbeg |
| 33 | acdefbdeh |
| 14 | acdefbdeg |
| 11 | acdefdbeg |
| 9 | adcefcdeh |
| 8 | adcefdbeh |
| 5 | adcefbdeg |
| 3 | acdefbdefdbeg |
| 2 | adcefdbeg |
| 2 | adcefbdefbdeg |
| 1 | adcefdbefbdeh |
| 1 | adbefbdefdbeg |
| 1 | adcefdbefcdefdbeg |
| 1391 | |

# Model N<sub>2</sub>



$N_2$ : fitness = -, precision = +, generalization = -, simplicity = +

| # | trace |
|---|---|
| 455 | acdeh |
| 191 | abdeg |
| 177 | adceh |
| 144 | abdeh |
| 111 | acdeg |
| 82 | adceg |
| 56 | adbeh |
| 47 | acdefdbeh |
| 38 | adbeg |
| 33 | acdefbdeh |
| 14 | acdefbdeg |
| 11 | acdefdbeg |
| 9 | adcefcdeh |
| 8 | adcefdbeh |
| 5 | adcefbdeg |
| 3 | acdefbdefdbeg |
| 2 | adcefdbeg |
| 2 | adcefbdefbdeg |
| 1 | adcefdbefbdeh |
| 1 | adbefbdefdbeg |
| 1 | adcefdbefcdefdbeg |
| 1391 | |

# Model N₃



$N_3$ : fitness = +, precision = -, generalization = +, simplicity = +

| # | trace |
|---|---|
| 455 | acdeh |
| 191 | abdeg |
| 177 | adceh |
| 144 | abdeh |
| 111 | acdeg |
| 82 | adceg |
| 56 | adbeh |
| 47 | acdefdbeh |
| 38 | adbeg |
| 33 | acdefbdeh |
| 14 | acdefbdeg |
| 11 | acdefdbeg |
| 9 | adcefcdeh |
| 8 | adcefdbeh |
| 5 | adcefbdeg |
| 3 | acdefbdefdbeg |
| 2 | adcefdbeg |
| 2 | adcefbdefbdeg |
| 1 | adcefdbefbdeh |
| 1 | adbefbdefdbeg |
| 1 | adcefdbefcdefdbeg |
| 1391 | |

# Model N₄



| # | trace |
|---|---|
| 455 | acdeh |
| 191 | abdeg |
| 177 | adceh |
| 144 | abdeh |
| 111 | acdeg |
| 82 | adceg |
| 56 | adbeh |
| 47 | acdefdbeh |
| 38 | adbeg |
| 33 | acdefbdeh |
| 14 | acdefbdeg |
| 11 | acdefdbeg |
| 9 | adcefcdeh |
| 8 | adcefdbeh |
| 5 | adcefbdeg |
| 3 | acdefbdefdbeg |
| 2 | adcefdbeg |
| 2 | adcefbdefbdeg |
| 1 | adcefdbefbdeh |
| 1 | adbefbdefdbeg |
| 1 | adcefdbefcdefdbeg |
| 1391 | |

N₄ : fitness = +, precision = +, generalization = -, simplicity = -

# Example of a process discovery technique: Genetic Mining



- **Characteristics**
  - requires a lot of computing power, but can be distributed easily,
  - can deal with noise, infrequent behavior, duplicate tasks, invisible tasks,
  - allows for incremental improvement and combinations with other approaches (heuristics post-optimization, etc.).

# Genetic process mining: Overview

# Example: crossover

# Example: mutation

# Example of a process discovery technique: Theory of Regions

- **Two types of regions theory:**
  - **State-based regions**
  - **Language-based regions**



$A=\{a_1, a_2, \ldots a_m\}$

$B=\{b_1, b_2, \ldots b_n\}$

$p_{(A,B)}$

**All about discovering places!**

# State-based regions: Two step approach

$$L_1 = [\langle a,b,c,d \rangle^3, \langle a,c,b,d \rangle^2, \langle a,e,d \rangle]$$

# Example region
# a enters, b and e exit, c and d do not cross

# Language-based regions



**A place is feasible if it can be added without disabling any of the traces in the event log.**

for any $\sigma \in L, k \in \{1, \ldots, |\sigma|\}, \sigma_1 = hd^{k-1}(\sigma), a = \sigma(k), \sigma_2 = hd^k(\sigma) = \sigma_1 \oplus a$:

$$c + \sum_{t \in X} \partial_{multiset}(\sigma_1)(t) - \sum_{t \in Y} \partial_{multiset}(\sigma_2)(t) \geq 0.$$

# Conformance checking

| frequency | reference | trace |
|---|---|---|
| 455 | $\sigma_1$ | $\langle a,c,d,e,h \rangle$ |
| 191 | $\sigma_2$ | $\langle a,b,d,e,g \rangle$ |
| 177 | $\sigma_3$ | $\langle a,d,c,e,h \rangle$ |
| 144 | $\sigma_4$ | $\langle a,b,d,e,h \rangle$ |
| 111 | $\sigma_5$ | $\langle a,c,d,e,g \rangle$ |
| 82 | $\sigma_6$ | $\langle a,d,c,e,g \rangle$ |
| 56 | $\sigma_7$ | $\langle a,d,b,e,h \rangle$ |
| 47 | $\sigma_8$ | $\langle a,c,d,e,f,d,b,e,h \rangle$ |
| 38 | $\sigma_9$ | $\langle a,d,b,e,g \rangle$ |
| 33 | $\sigma_{10}$ | $\langle a,c,d,e,f,b,d,e,h \rangle$ |
| 14 | $\sigma_{11}$ | $\langle a,c,d,e,f,b,d,e,g \rangle$ |
| 11 | $\sigma_{12}$ | $\langle a,c,d,e,f,d,b,e,g \rangle$ |
| 9 | $\sigma_{13}$ | $\langle a,d,c,e,f,c,d,e,h \rangle$ |
| 8 | $\sigma_{14}$ | $\langle a,d,c,e,f,d,b,e,h \rangle$ |
| 5 | $\sigma_{15}$ | $\langle a,d,c,e,f,b,d,e,g \rangle$ |
| 3 | $\sigma_{16}$ | $\langle a,c,d,e,f,b,d,e,f,d,b,e,g \rangle$ |
| 2 | $\sigma_{17}$ | $\langle a,d,c,e,f,d,b,e,g \rangle$ |
| 2 | $\sigma_{18}$ | $\langle a,d,c,e,f,b,d,e,f,b,d,e,g \rangle$ |
| 1 | $\sigma_{19}$ | $\langle a,d,c,e,f,d,b,e,f,b,d,e,h \rangle$ |
| 1 | $\sigma_{20}$ | $\langle a,d,b,e,f,b,d,e,f,d,b,e,g \rangle$ |
| 1 | $\sigma_{21}$ | $\langle a,d,c,e,f,d,b,e,f,c,d,e,f,d,b,e,g \rangle$ |

| frequency | reference | trace |
|---|---|---|
| 455 | $\sigma_1$ | $\langle a,c,d,e,h \rangle$ |
| 191 | $\sigma_2$ | $\langle a,b,d,e,g \rangle$ |
| 177 | $\sigma_3$ | $\langle a,d,c,e,h \rangle$ |
| 144 | $\sigma_4$ | $\langle a,b,d,e,h \rangle$ |
| 111 | $\sigma_5$ | $\langle a,c,d,e,g \rangle$ |
| 82 | $\sigma_6$ | $\langle a,d,c,e,g \rangle$ |
| 56 | $\sigma_7$ | $\langle a,d,b,e,h \rangle$ |
| 47 | $\sigma_8$ | $\langle a,c,d,e,f,d,b,e,h \rangle$ |
| 38 | $\sigma_9$ | $\langle a,d,b,e,g \rangle$ |
| 33 | $\sigma_{10}$ | $\langle a,c,d,e,f,b,d,e,h \rangle$ |
| 14 | $\sigma_{11}$ | $\langle a,c,d,e,f,b,d,e,g \rangle$ |
| 11 | $\sigma_{12}$ | $\langle a,c,d,e,f,d,b,e,g \rangle$ |
| 9 | $\sigma_{13}$ | $\langle a,d,c,e,f,c,d,e,h \rangle$ |
| 8 | $\sigma_{14}$ | $\langle a,d,c,e,f,d,b,e,h \rangle$ |

# Replaying (2/7) $\sigma_1$ on $N_1$

$$\sigma_1 = \langle a, c, d, e, h \rangle$$



```
p=1
c=0
m=0
r=0
```

```
p=3
c=1
m=0
r=0
```

p' = p+2
c' = c+1

# Replaying (3/7) $\sigma_1$ on $N_1$

$$\sigma_1 = \langle a, c, d, e, h \rangle$$

p=3
c=1
m=0
r=0

p=4
c=2
m=0
r=0

p' = p+1
c' = c+1

$$\sigma_1 = \langle a, c, d, e, h \rangle$$



p=4
c=2
m=0
r=0

p=5
c=3
m=0
r=0

**p' = p+1**
**c' = c+1**

# Replaying (5/7) $\sigma_1$ on $N_1$

$$\sigma_1 = \langle a, c, d, e, h \rangle$$

p=5
c=3
m=0
r=0

p=6
c=5
m=0
r=0

p' = p+1
c' = c+2

p=6
c=5
m=0
r=0

p=7
c=6
m=0
r=0

p' = p+1
c' = c+1

# Replaying (7/7) $\sigma_1$ on $N_1$

$$\sigma_1 = \langle a, c, d, e, h \rangle$$

p=7
c=6
m=0
r=0

p=7
c=7
m=0
r=0

start

a

b

c

d

e

f

g

h

p1

p2

p3

p4

p5

end

$c' = c+1$

m = 0
r = 0
no problems encountered

$$\sigma_3 = \langle a, d, c, e, h \rangle$$



| p=0 | p=1 |
|-----|-----|
| c=0 | c=0 |
| m=0 | m=0 |
| r=0 | r=0 |

**p' = p+1**

**p = produced**
**c = consumed**
**m = missing ≤ c**
**r = remaining ≤ p**
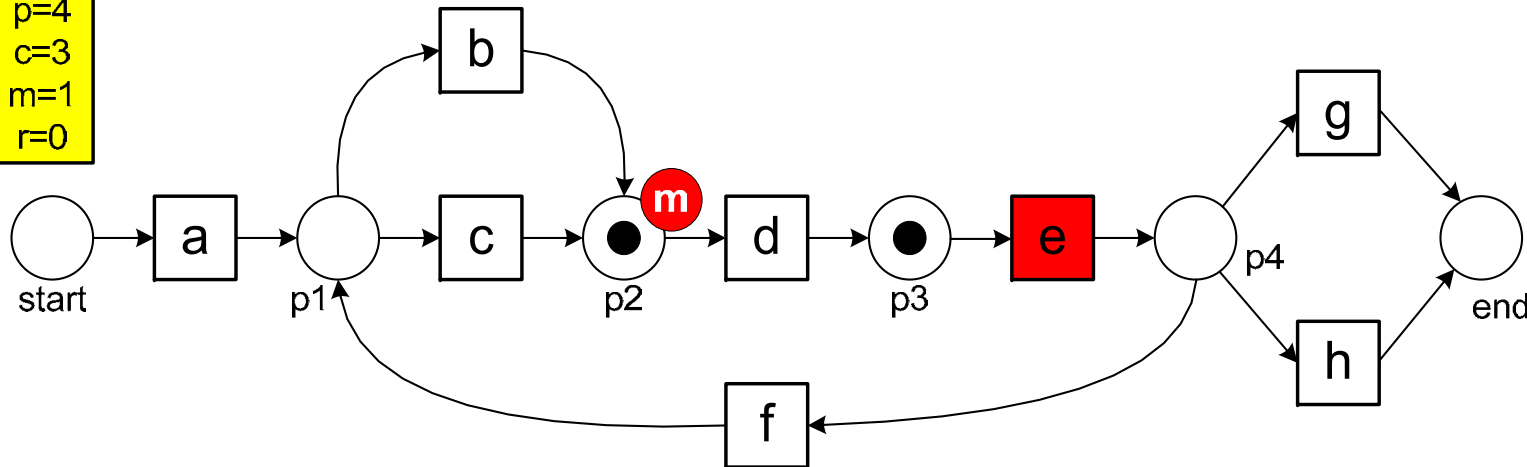
# Replaying (2/7) $\sigma_3$ on $N_2$

$$\sigma_3 = \langle a, d, c, e, h \rangle$$



p' = p+1
c' = c+1

# Replaying (3/7) $\sigma_3$ on $N_2$
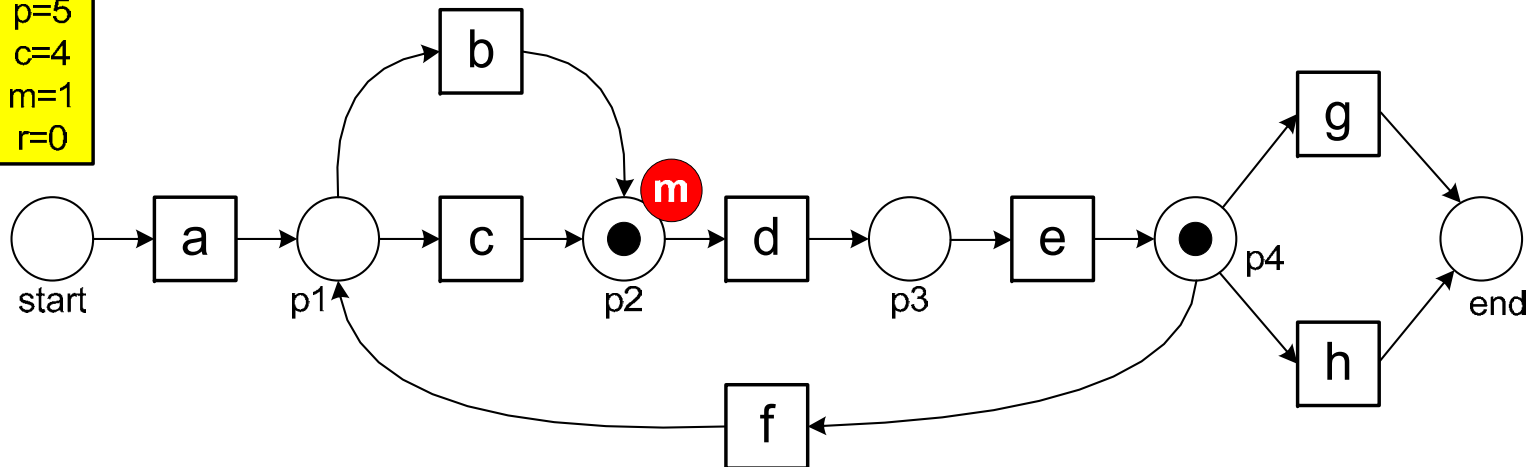
$$\sigma_3 = \langle a, d, c, e, h \rangle$$

p=2
c=1
m=0
r=0

p=3
c=2
m=1
r=0

p' = p+1
c' = c+1
m' = m+1

p=3
c=2
m=1
r=0

start  a  p1  b  c  m  d  p2  p3  e  p4  g  h  end  f

p' = p+1
c' = c+1

p=4
c=3
m=1
r=0
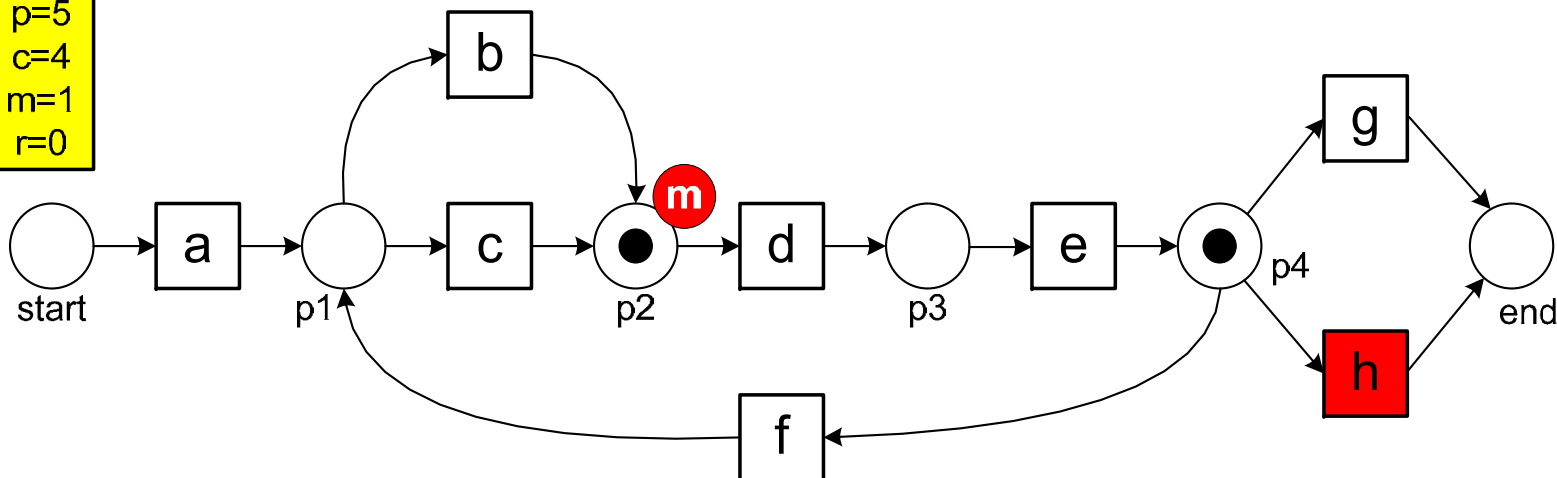
start  a  p1  b  c  m  d  p2  p3  e  p4  g  h  end  f

# Replaying (5/7) $\sigma_3$ on $N_2$

$$\sigma_3 = \langle a, d, c, e, h \rangle$$



p=4
c=3
m=1
r=0

p=5
c=4
m=1
r=0

p' = p+1
c' = c+1

# Replaying (6/7) σ₃ on N₂
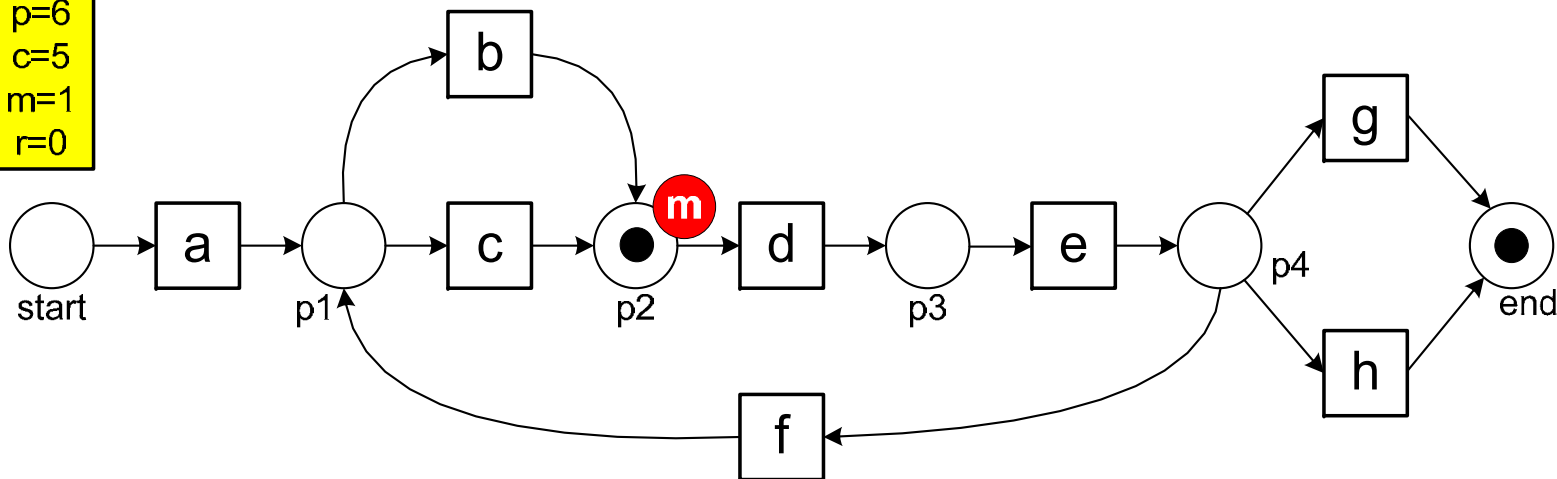
$$\sigma_3 = \langle a, d, c, e, h \rangle$$
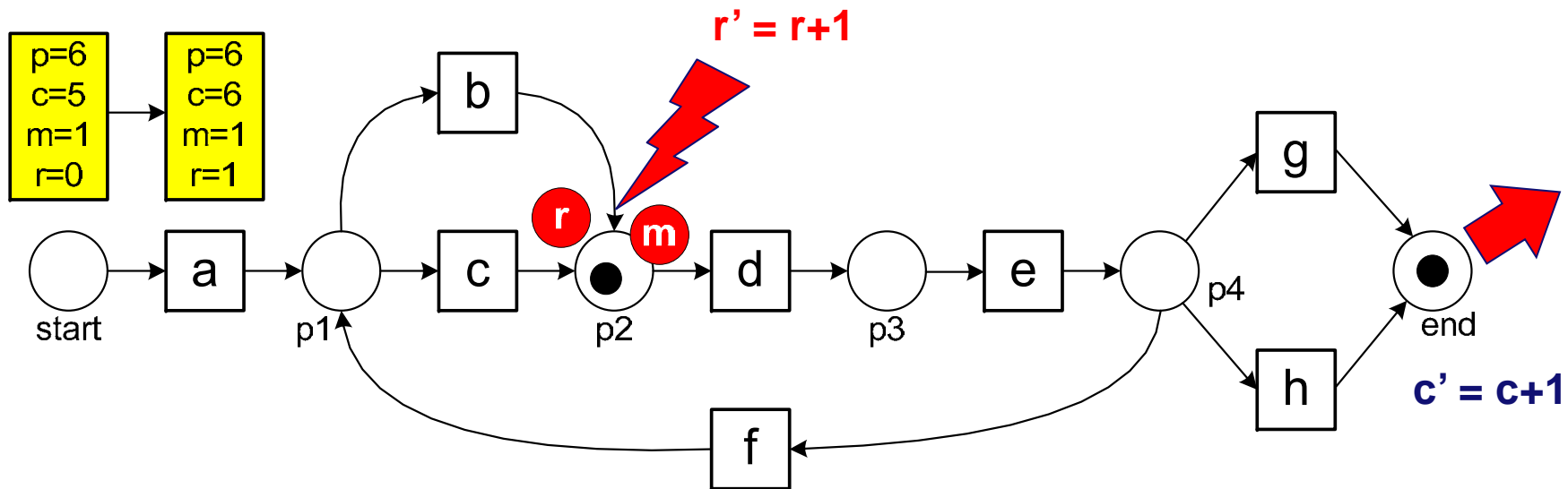
p=5
c=4
m=1
r=0

p=6
c=5
m=1
r=0

p' = p+1
c' = c+1

# Replaying (7/7) $\sigma_3$ on $N_2$
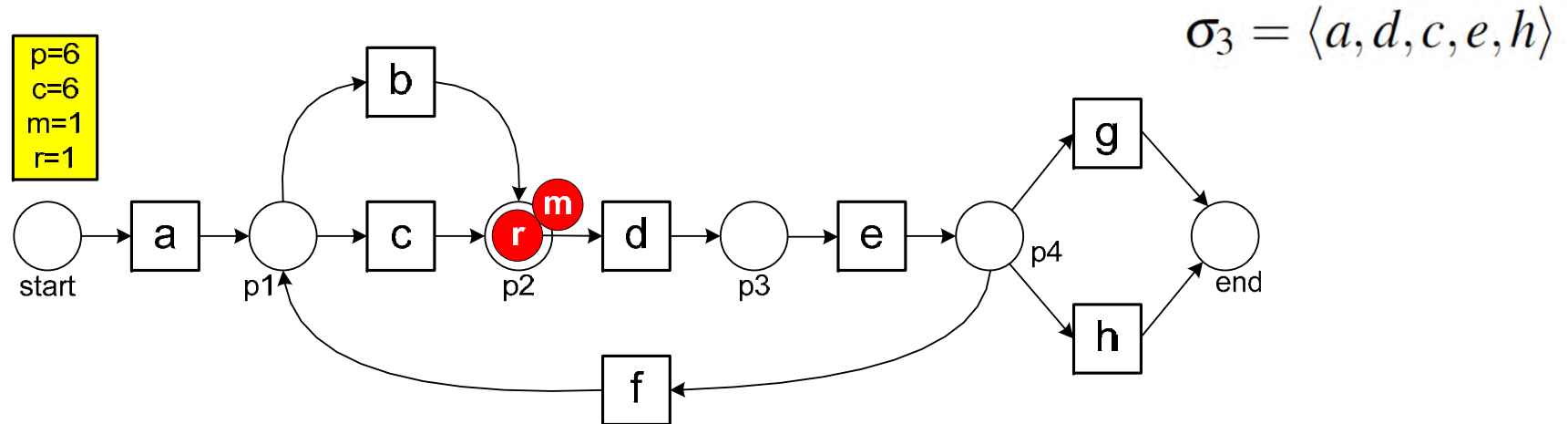
$$\sigma_3 = \langle a, d, c, e, h \rangle$$



**Problems:**
- m = 1 : d was forced to occur without being enabled
- r = 1 : output of c was not used by d

# Computing fitness at trace level



$$\sigma_3 = \langle a, d, c, e, h \rangle$$

p=6
c=6
m=1
r=1

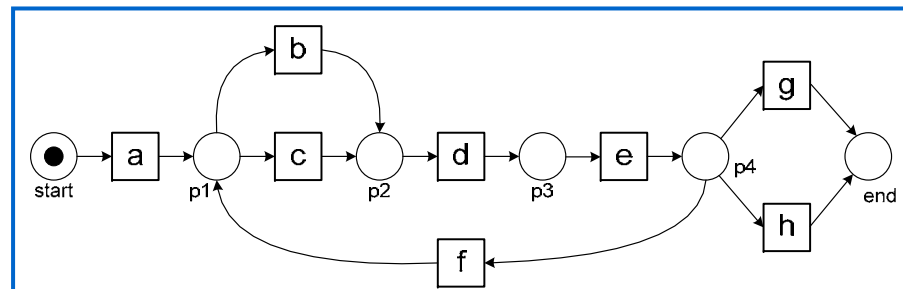$$fitness(\sigma, N) = \frac{1}{2}\left(1 - \frac{m}{c}\right) + \frac{1}{2}\left(1 - \frac{r}{p}\right)$$

$$fitness(\sigma_3, N_2) = \frac{1}{2}\left(1 - \frac{1}{6}\right) + \frac{1}{2}\left(1 - \frac{1}{6}\right) = 0.8333$$

$$fitness(L,N) = \frac{1}{2}\left(1 - \frac{\sum_{\sigma \in L} L(\sigma) \times m_{N,\sigma}}{\sum_{\sigma \in L} L(\sigma) \times c_{N,\sigma}}\right) +$$

$$\frac{1}{2}\left(1 - \frac{\sum_{\sigma \in L} L(\sigma) \times r_{N,\sigma}}{\sum_{\sigma \in L} L(\sigma) \times p_{N,\sigma}}\right)$$

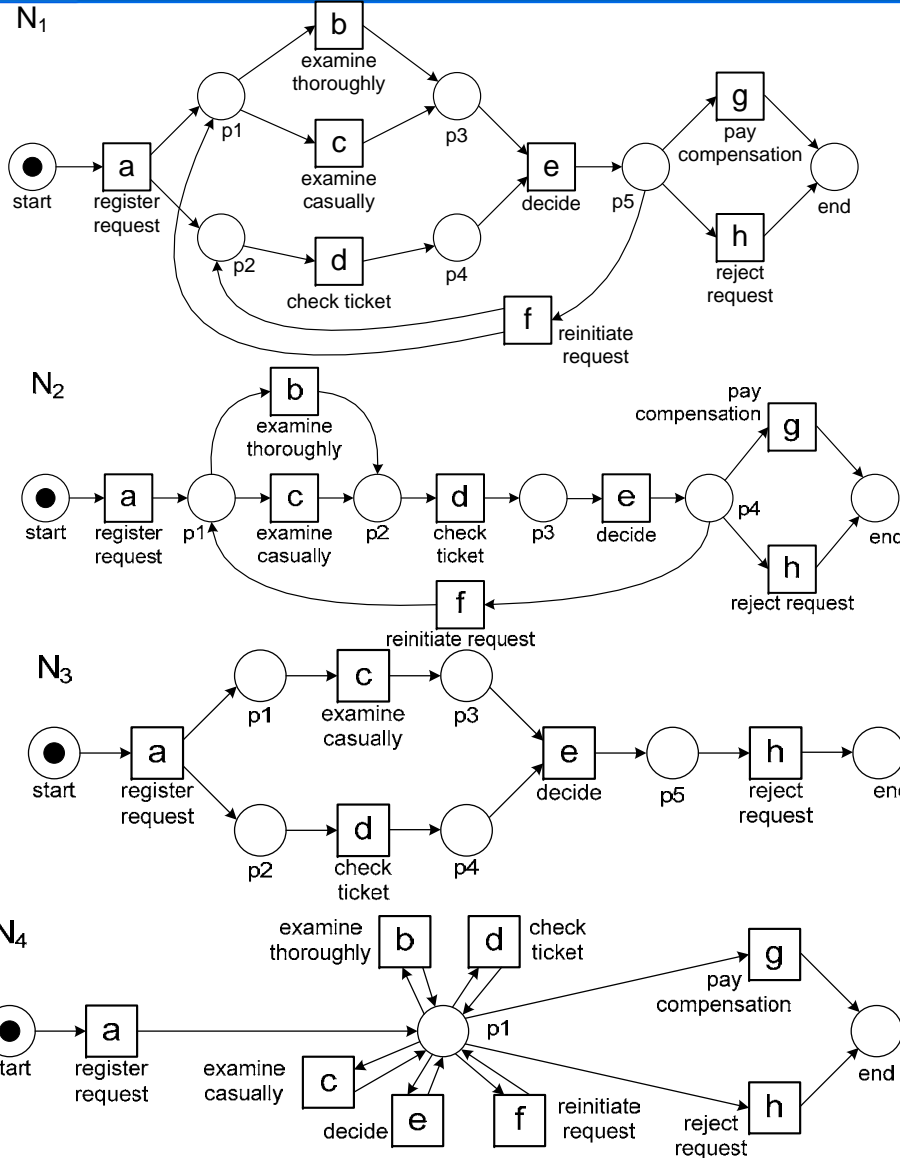| frequency | reference | trace |
|---|---|---|
| 455 | $\sigma_1$ | $\langle a,c,d,e,h \rangle$ |
| 191 | $\sigma_2$ | $\langle a,b,d,e,g \rangle$ |
| 177 | $\sigma_3$ | $\langle a,d,c,e,h \rangle$ |
| 144 | $\sigma_4$ | $\langle a,b,d,e,h \rangle$ |
| 111 | $\sigma_5$ | $\langle a,c,d,e,g \rangle$ |
| 82 | $\sigma_6$ | $\langle a,d,c,e,g \rangle$ |
| 56 | $\sigma_7$ | $\langle a,d,b,e,h \rangle$ |
| 47 | $\sigma_8$ | $\langle a,c,d,e,f,d,b,e,h \rangle$ |
| 38 | $\sigma_9$ | $\langle a,d,b,e,g \rangle$ |
| 33 | $\sigma_{10}$ | $\langle a,c,d,e,f,b,d,e,h \rangle$ |
| 14 | $\sigma_{11}$ | $\langle a,c,d,e,f,b,d,e,g \rangle$ |
| 11 | $\sigma_{12}$ | $\langle a,c,d,e,f,d,b,e,g \rangle$ |
| 9 | $\sigma_{13}$ | $\langle a,d,c,e,f,c,d,e,h \rangle$ |
| 8 | $\sigma_{14}$ | $\langle a,d,c,e,f,d,b,e,h \rangle$ |
| 5 | $\sigma_{15}$ | $\langle a,d,c,e,f,b,d,e,g \rangle$ |
| 3 | $\sigma_{16}$ | $\langle a,c,d,e,f,b,d,e,f,d,b,e,g \rangle$ |
| 2 | $\sigma_{17}$ | $\langle a,d,c,e,f,d,b,e,g \rangle$ |
| 2 | $\sigma_{18}$ | $\langle a,d,c,e,f,b,d,e,f,b,d,e,g \rangle$ |
| 1 | $\sigma_{19}$ | $\langle a,d,c,e,f,d,b,e,f,b,d,e,h \rangle$ |
| 1 | $\sigma_{20}$ | $\langle a,d,b,e,f,b,d,e,f,d,b,e,g \rangle$ |
| 1 | $\sigma_{21}$ | $\langle a,d,c,e,f,d,b,e,f,c,d,e,f,d,b,e,g \rangle$ |

# Example values

# Diagnostics

$$(fitness(L_{full}, N_3) = 0.8797)$$



**problem**
*430* tokens remain in place *p1*, because *c* did not happen while the model expected *c* to happen

**problem**
*566* tokens were missing in place *p3* during replay, because *e* happened while this was not possible according to the model

**problem**
*10* tokens were missing in place *p1* during replay, because *c* happened while this was not possible according to the model

**problem**
*146* tokens were missing in place *p2* during replay, because *d* happened while this was not possible according to the model

**problem**
*607* tokens remain in place *p5*, because *h* did not happen while the model expected *h* to happen

**problem**
*461* of the *1391* cases did not reach place *end*

# Challenges related to conformance checking

- **Not as simple as it seems!**

- **In case of duplicate tasks (two transition with the same label) or silent tasks ($\tau$ labeled transitions), multiple paths need to be considered (state space analysis, heuristics, or optimization).**

- **More general formulation of the problem with costs associated to skipping/inserting particular tasks, see ProM latest conformance checker (A\* algorithm).**

- **Computing the most likely alignment is needed for other types of process mining (time analysis, measuring precision, social network analysis, etc.).**

# How can process mining help?

- **Detect bottlenecks**
- **Detect deviations**
- **Performance measurement**
- **Suggest improvements**
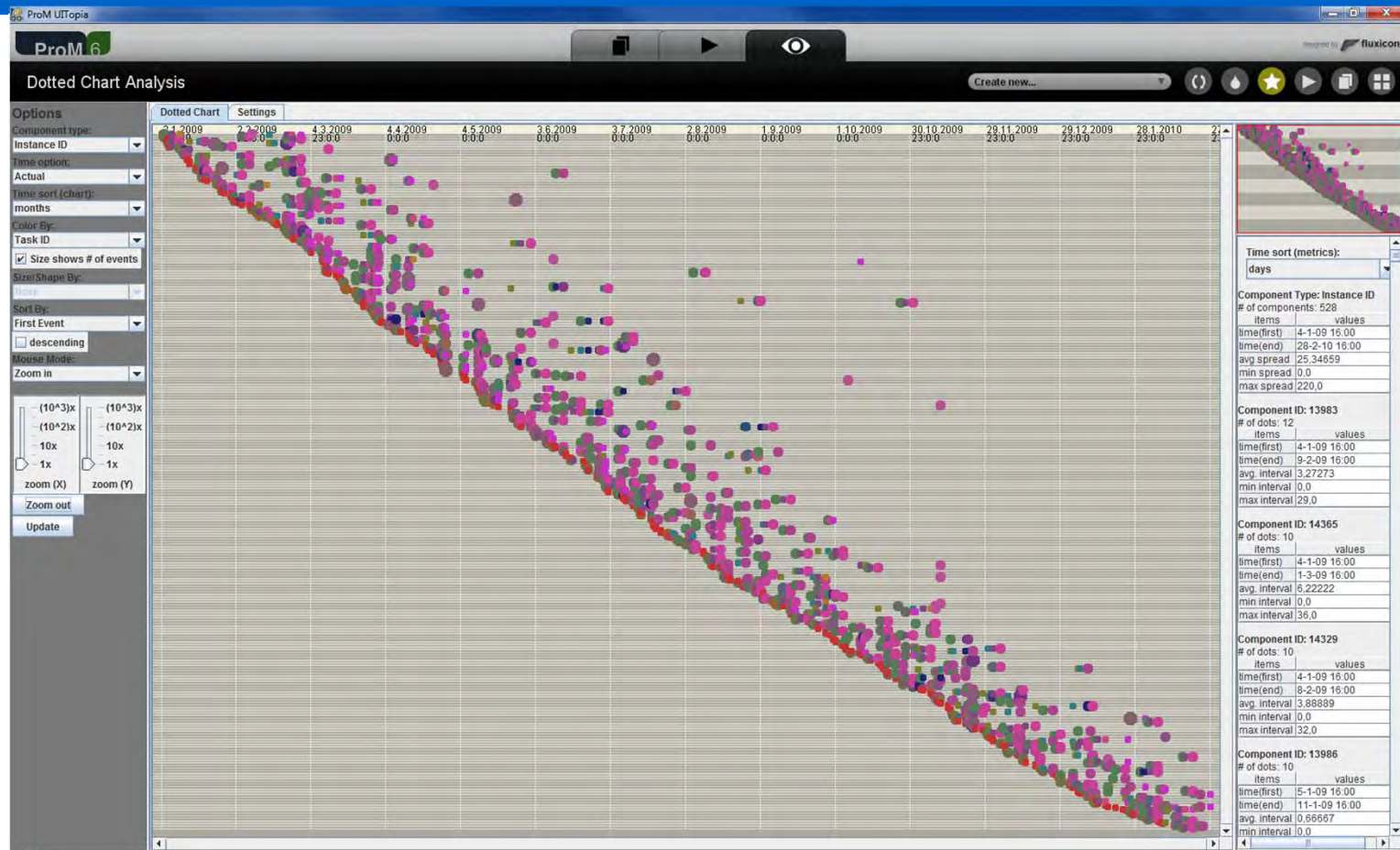- **Decision support (e.g., recommendation and prediction)**

- **Provide mirror**
- **Highlight important problems**
- **Avoid ICT failures**
- **Avoid management by PowerPoint**
- **From "politics" to "analytics"**

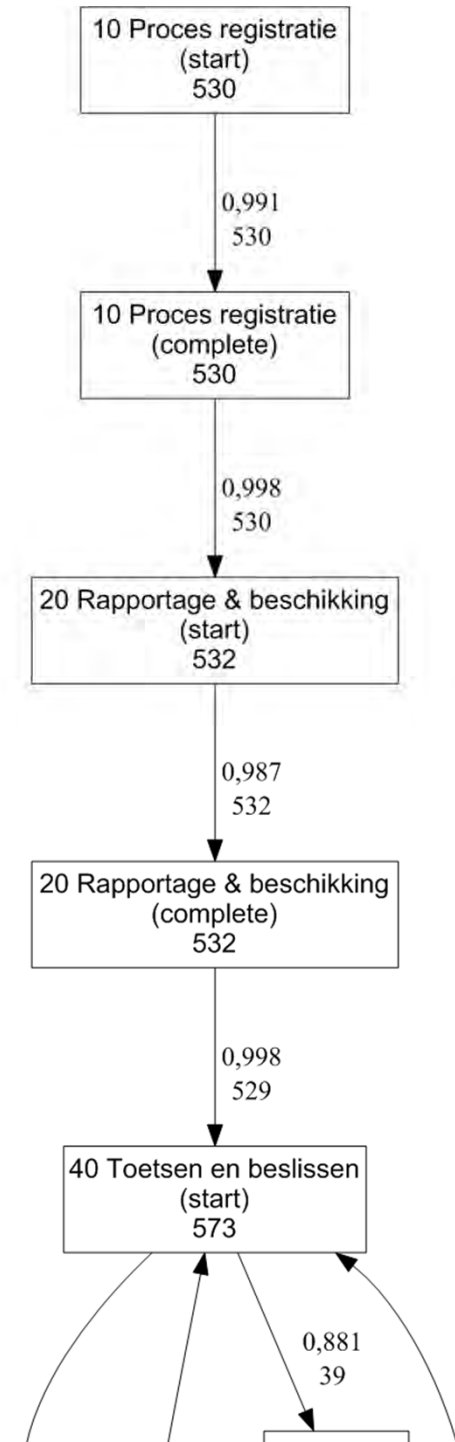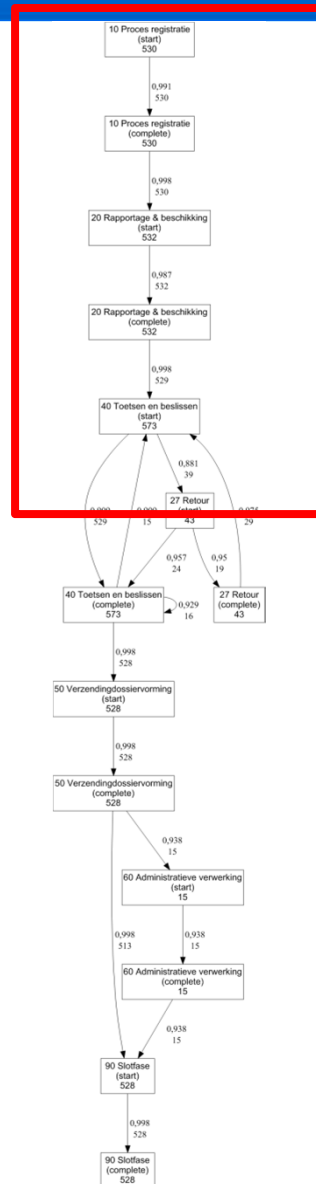# Example of a Lasagna process: WMO process of a Dutch municipality



**Each line corresponds to one of the 528 requests that were handled in the period from 4-1-2009 until 28-2-2010. In total there are 5498 events represented as dots. The mean time needed to handled a case is approximately 25 days.**
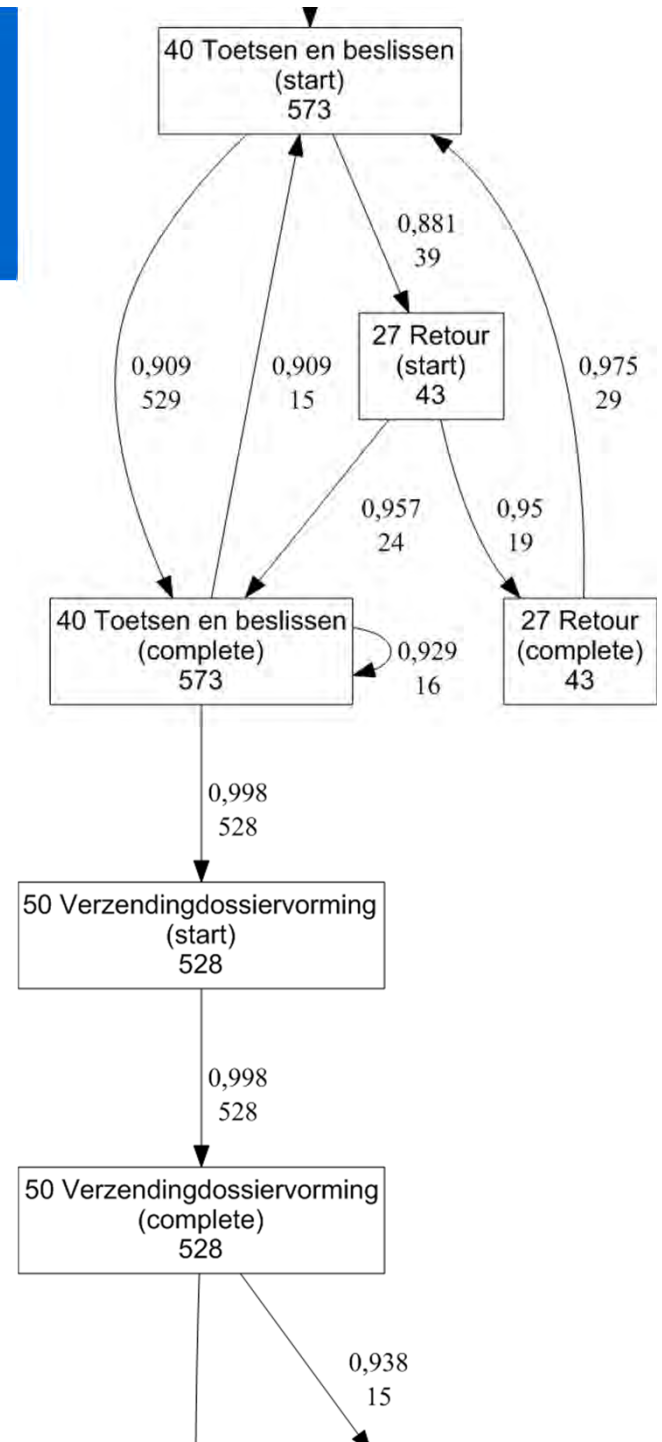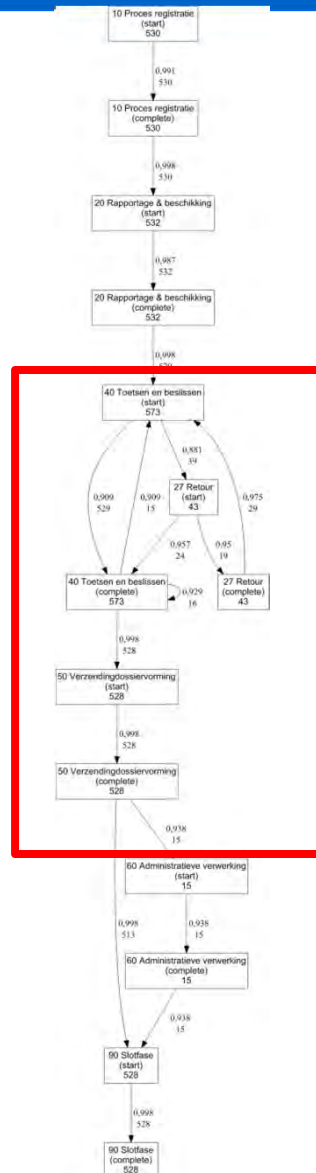
# WMO process
## (Wet Maatschappelijke Ondersteuning)

- **WMO refers to the social support act that came into force in The Netherlands on January 1st, 2007.**

- **The aim of this act is to assist people with disabilities and impairments. Under the act, local authorities are required to give support to those who need it, e.g., household help, providing wheelchairs and scootmobiles, and adaptations to homes.**

- **There are different processes for the different kinds of help. We focus on the process for handling requests for household help.**

- **In a period of about one year, 528 requests for household WMO support were received.**

- **These 528 requests generated 5498 events.**

# C-net discovered using heuristic miner (1/3)
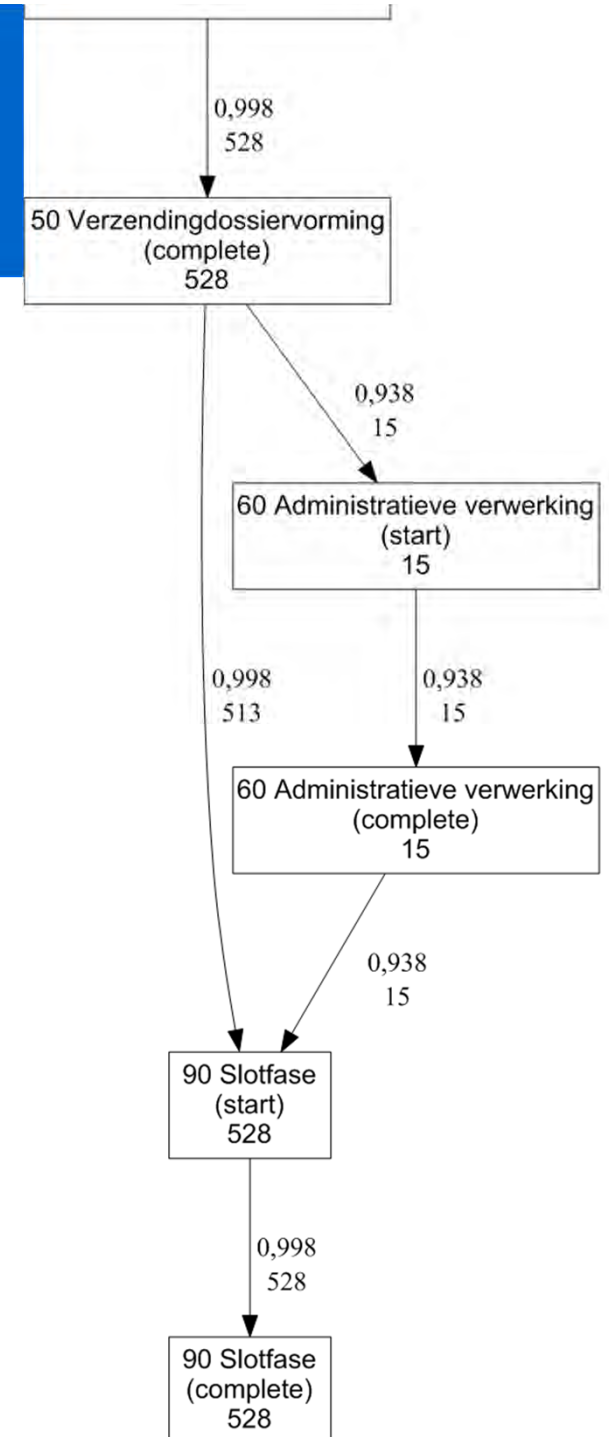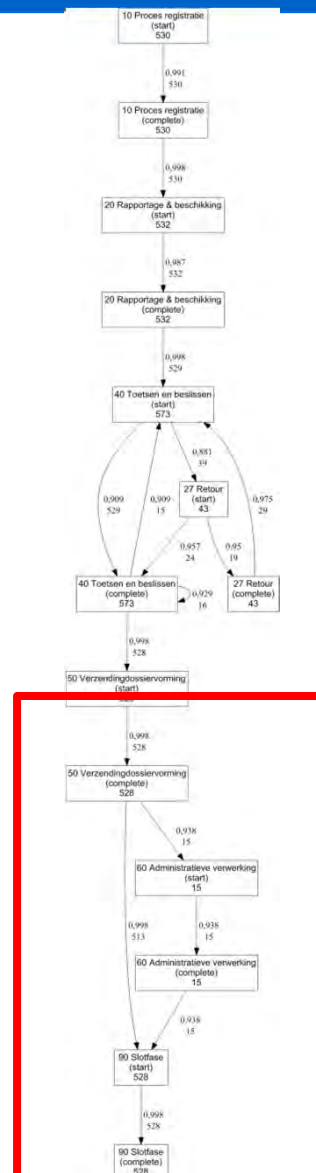
# C-net discovered using heuristic miner (2/3)

# C-net discovered using heuristic miner (3/3)

# Conformance check WMO process (3/3)



**The fitness of the discovered process is 0.99521667. Of the 528 cases, 496 cases fit perfectly whereas for 32 cases there are missing or remaining tokens.**

# Bottleneck analysis WMO process (1/3)

# Bottleneck analysis WMO process (2/3)

# Bottleneck analysis WMO process (3/3)



**flow time of approx. 25 days with a standard deviation of approx. 28**

## Process information:

**Total number selected:**

528 cases

**Number fitting:**

496 cases

**Arrival rate:**

1,21 cases per day

| | Throughput time (days) |
|---|---|
| avg | 24,66 |
| min | 0,0 |
| max | 220,0 |
| stdev | 27,86 |
| fast 25... | 3,54 |
| slow 2... | 60,27 |
| norma... | 17,42 |

# Two additional Lasagna processes

**RWS
("Rijkswaterstaat")
process**

**WOZ ("Waardering
Onroerende Zaken")
process**

# RWS Process

- **The Dutch national public works department, called "Rijkswaterstaat" (RWS), has twelve provincial offices. We analyzed the handling of invoices in one of these offices.**

- **The office employs about 1,000 civil servants and is primarily responsible for the construction and maintenance of the road and water infrastructure in its province.**

- **To perform its functions, the RWS office subcontracts various parties such as road construction companies, cleaning companies, and environmental bureaus. Also, it purchases services and products to support its construction, maintenance, and administrative activities.**

# C-net discovered using heuristic miner

# Social network constructed based on handovers of work



Each of the 271 nodes corresponds to a civil servant. Two civil servants are connected if one executed an activity causally following an activity executed by the other civil servant

# Social network consisting of civil servants that executed more than 2000 activities in a 9 month period.



The darker arcs indicate the strongest relationships in the social network. Nodes having the same color belong to the same clique.

# WOZ process

- **Event log containing information about 745 objections against the so-called WOZ ("Waardering Onroerende Zaken") valuation.**

- **Dutch municipalities need to estimate the value of houses and apartments. The WOZ value is used as a basis for determining the real-estate property tax.**

- **The higher the WOZ value, the more tax the owner needs to pay. Therefore, there are many objections (i.e., appeals) of citizens that assert that the WOZ value is too high.**

- **"WOZ process" discovered for another municipality (i.e., different from the one for which we analyzed the WMO process).**

# Discovered process model



**The log contains events related to 745 objections against the so-called WOZ valuation. These 745 objections generated 9583 events. There are 13 activities. For 12 of these activities both start and complete events are recorded. Hence, the WF-net has 25 transitions.**

# Conformance checker:
## (fitness is 0.98876214)

# Performance analysis



Waiting time:
- **High** (magenta)
- **Medium** (yellow)
- **Low** (blue)

bottleneck detection: places are colored based on average durations

time required to move from one activity to another

information on total flow time

## Performance information of the selected transitions:

Frequency: 416 cases

| | Time in between (days) |
|---|---|
| avg | 202,73 |
| min | 126,89 |
| max | 245,98 |
| stdev | 19,74 |
| fast 25.00%... | 177,2 |

## Arrival rate:

2,85 cases per day

| | Throughput time (days) |
|---|---|
| avg | 177,99 |
| min | 3,78 |
| max | 251,9 |
| stdev | 52,87 |
| fast 25... | 98,98 |
| slow 2... | 230,76 |
| norma... | 191,11 |

# Resource-activity matrix (four groups discovered)

| user | $a_1$ | $a_2$ | $a_3$ | $a_4$ | $a_5$ | $a_6$ | $a_7$ | $a_8$ | $a_9$ | $a_{10}$ | $a_{11}$ | $a_{12}$ | $a_{13}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| user 1 | 0 | 0 | 51 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| user 2 | 1 | 2 | 0 | 0 | 2 | 0 | 0 | 0 | 0 | 38 | 0 | 69 | 0 |
| user 3 | 0 | 9 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| user 4 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| user 5 | 117 | 0 | 4 | 0 | 3 | 0 | 0 | 0 | 0 | 1 | 0 | 20 | 6 |
| user 6 | 172 | 6 | 14 | 0 | 7 | 3 | 0 | 0 | 1 | 2 | 0 | 48 | 53 |
| user 7 | 1 | 41 | 8 | 14 | 275 | 8 | 8 | 865 | 55 | 180 | 0 | 128 | 5 |
| user 8 | 2 | 868 | 7 | 6 | 105 | 0 | 0 | 79 | 266 | 441 | 0 | 844 | 3 |
| user 9 | 90 | 0 | 2 | 0 | 1 | 2 | 0 | 0 | 1 | 2 | 0 | 27 | 28 |
| user 10 | 0 | 0 | 0 | 899 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1019 |
| user 11 | 336 | 1 | 3 | 1 | 4 | 2 | 0 | 0 | 0 | 1 | 0 | 18 | 23 |
| user 12 | 1 | 645 | 13 | 21 | 419 | 3 | 0 | 3 | 217 | 281 | 1 | 334 | 9 |
| user 13 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| user 14 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| user 15 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 2 | 0 | 0 | 2 | 0 |
| user 16 | 1 | 3 | 3 | 2 | 1 | 0 | 0 | 1 | 2 | 3 | 1 | 0 | 0 |
| user 17 | 0 | 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| user 18 | 9 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| user 19 | 13 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 4 | 0 |
| user 20 | 0 | 0 | 0 | 21 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 258 |

# Example of a Spaghetti process



**Spaghetti process describing the diagnosis and treatment of 2765 patients in a Dutch hospital. The process model was constructed based on an event log containing 114,592 events. There are 619 different activities (taking event types into account) executed by 266 different individuals (doctors, nurses, etc.).**

# Fragment
## 18 activities of the 619 activities (2.9%)

# Another example
## (event log of Dutch housing agency)



The event log contains 208 cases that generated 5987 events. There are 74 different activities.

010 Registreren huuropzegging
(complete)
208

0,857
6

0,978
154

0,857
6

0,944
30

040 Vastleggen toekomstig adres medehuurder
(complete)
32

020 Vastleggen datum van overlijden
(complete)
6

0,992
193

057 Plannen eindinspectie bedryfsr/gar/ber/park/op
(complete)
9

050 Plannen afspraak 1e inspectie
(complete)
163

0,667
3

0,857
6

0,833
9

050 Inplannen afspraak 1e inspectie
(complete)
33

0,992
163

030 Vastleggen toekomstige adres
(complete)
208

058 Aanmaken bevest.brief huuropzegging(b/g/bso/p)
(complete)
11

0,958
33

0,5
1

060 Aanmaken bevestigingbrief / huuropzeggingform.
(complete)
196

0,875
103

055 Plannen eindinspectie bedryfsr/gar/ber/park/op
(complete)
1

0,995
93

070 Is 1e inspectie uitgevoerd ?
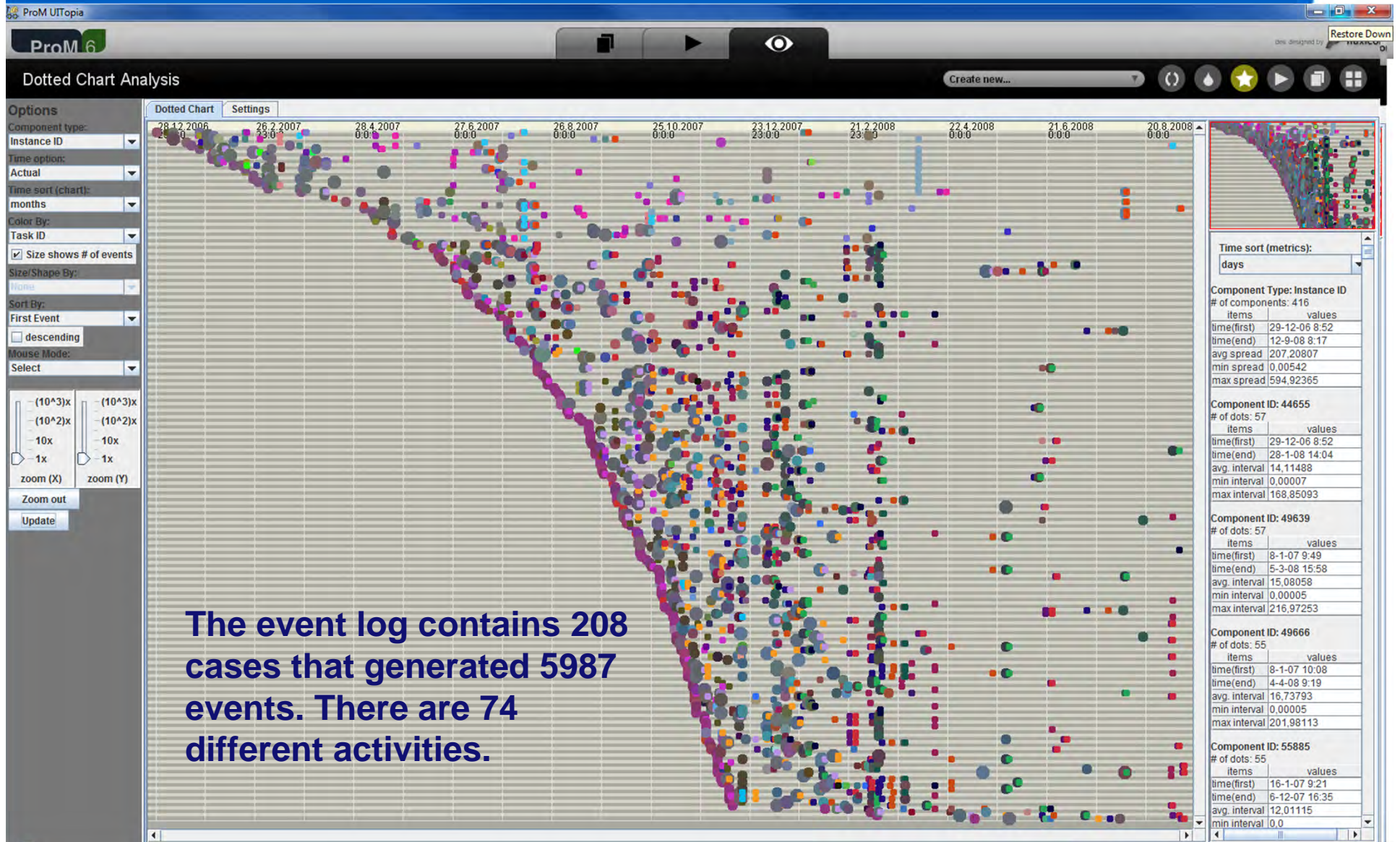(complete)
204

0,923
12

0,993
192

0,966
192

0,5
1

0,944
192

0,923
8

080 Versturen brief 'Niet thuis'
(complete)
12

120 Plannen eindinspectie
(complete)
192

100 Gereedmelden 1e insp. / Voorcalculatie maken
(complete)
192

0,938
20

0,923
12

090 Herplannen 1e inspectie
(complete)
12

300 Is eindinspectie uitgevoerd ?
(complete)
34

110 Bepalen leegstandsoort
(complete)
192

0,971
34

340 Zijn er nieuwe of niet herstelde gebreken ?
(complete)

# Conclusion

- **Many concurrency-related BPM challenges.**
- **Process leave traces in event logs. So, if you are interested in processes, use them!**
- **Process mining: challenging and highly relevant.**
- **Process discovery challenge**
  - **balancing between different objectives**
  - **only example behavior**
- **Conformance checking challenge**
  - **finding the most likely trace**
  - **dealing with silent/duplicate steps**
- **Eldorado for exciting concurrency research!**

More and more information about business processes is recorded by information systems in the form of so-called "event logs". Despite the omnipresence of such data, most organizations diagnose problems based on fiction rather than facts. Process mining is an emerging discipline based on process model-driven approaches and data mining. It not only allows organizations to fully benefit from the information stored in their systems, but it can also be used to check the conformance of processes, detect bottlenecks, and predict execution problems.

Wil van der Aalst delivers the first book on process mining. It aims to be self-contained while covering the entire process mining spectrum from process discovery to operational support. In Part I, the author provides the basics of business process modeling and data mining necessary to understand the remainder of the book. Part II focuses on process discovery as the most important process mining task. Part III moves beyond discovering the control flow of processes and highlights conformance checking, and organizational and time perspectives. Part IV guides the reader in successfully applying process mining in practice, including an introduction to the widely used open-source tool ProM. Finally, Part V takes a step back, reflecting on the material presented and the key open challenges.

Overall, this book provides a comprehensive overview of the state of the art in process mining. It is intended for business process analysts, business consultants, process managers, graduate students, and BPM researchers.

**Features and Benefits:**

- First book on process mining, bridging the gap between business process modeling and business intelligence.
- Written by one of the most influential and most-cited computer scientists and the best-known BPM researcher.
- Self-contained and comprehensive overview for a broad audience in academia and industry.
- The reader can put process mining into practice immediately due to the applicability of the techniques and the availability of the open-source process mining software ProM.

van der Aalst

Process Mining

Wil M. P. van der Aalst

# Process Mining

Discovery, Conformance and Enhancement of Business Processes

**www.processmining.org**

**www.win.tue.nl/ieeetfpm/**

Springer