# Mediating Between Modeled and Observed Behavior: The Quest for the "Right" Process

**prof.dr.ir. Wil van der Aalst**

PAIS lab
Лаборатория ПОИС

NATIONAL RESEARCH UNIVERSITY

TU/e Technische Universiteit **Eindhoven** University of Technology

**Where innovation starts**

# Outline

Introduction to Process Mining (short)

The 4+ dimensions of conformance

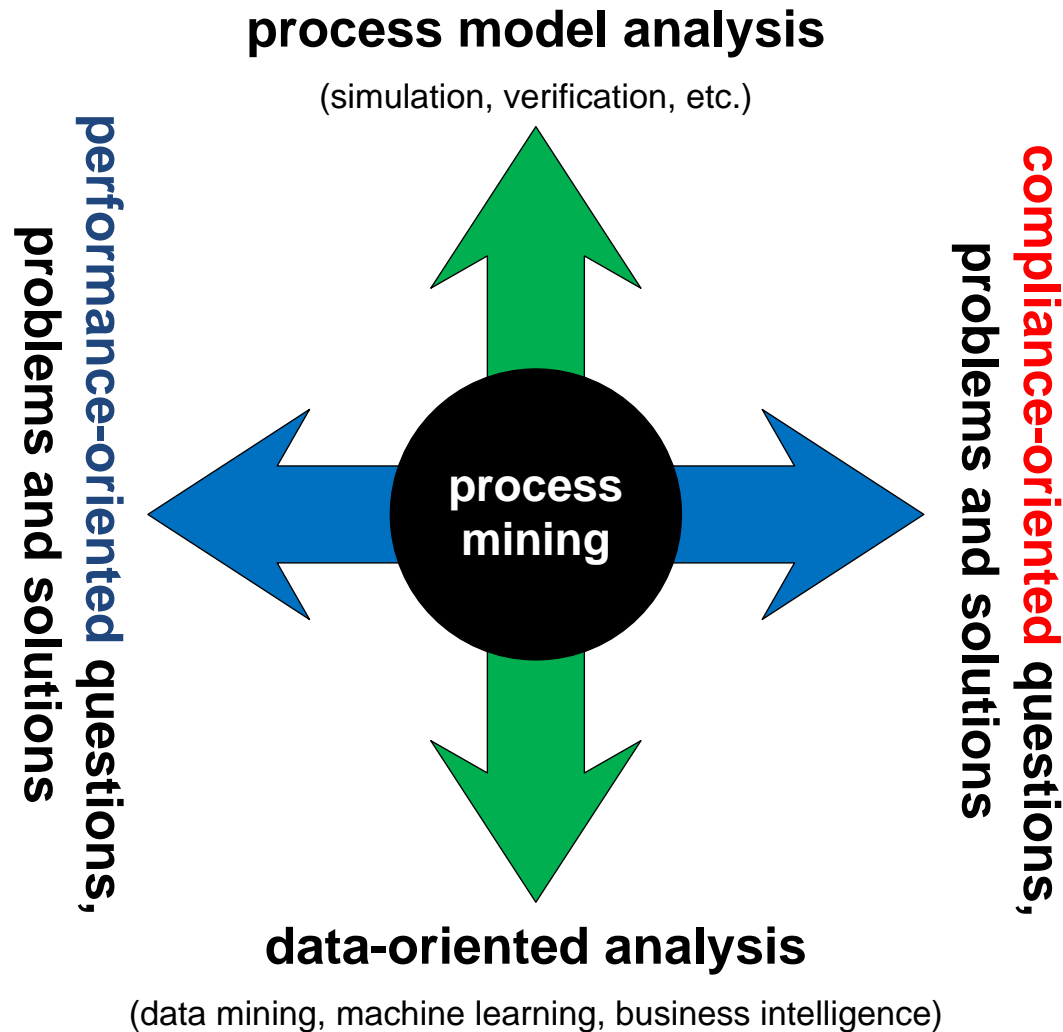Importance of alignments to relate observed and modeled behavior

Representational bias

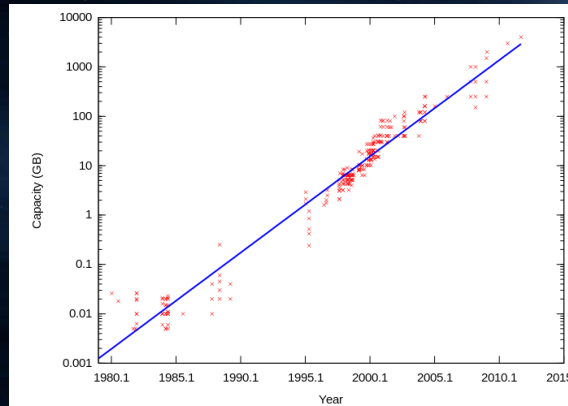Mediating between a reference model and observed behavior (model repair)

Discovering configurable process models

Decomposing process mining problems to deal with Big Data

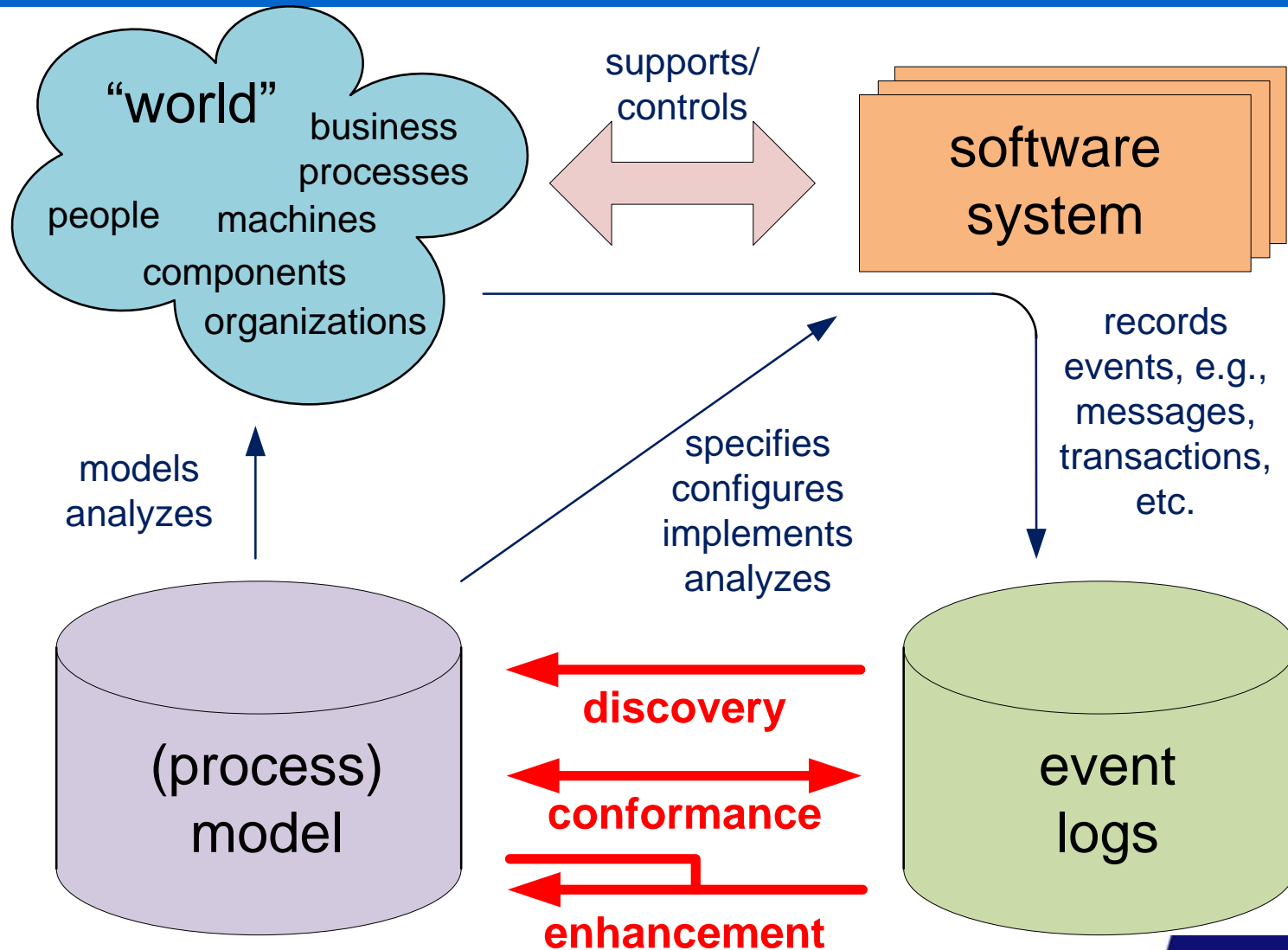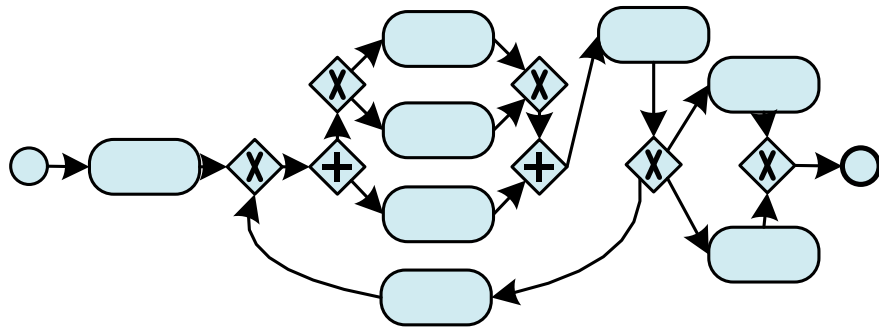# Positioning Process Mining



**process model analysis**

(simulation, verification, etc.)

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

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

**process mining**

**data-oriented analysis**

(data mining, machine learning, business intelligence)

# Moore's Law

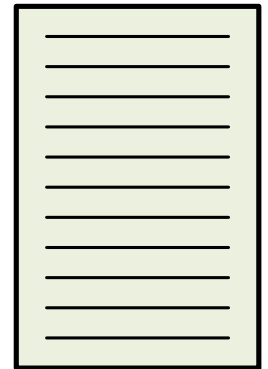Microprocessor Transistor Counts 1971-2011 & Moore's Law



2060

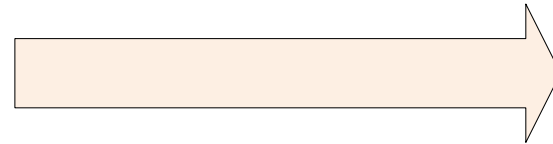# Three Types of Process Mining

# Play-Out



process model

event log

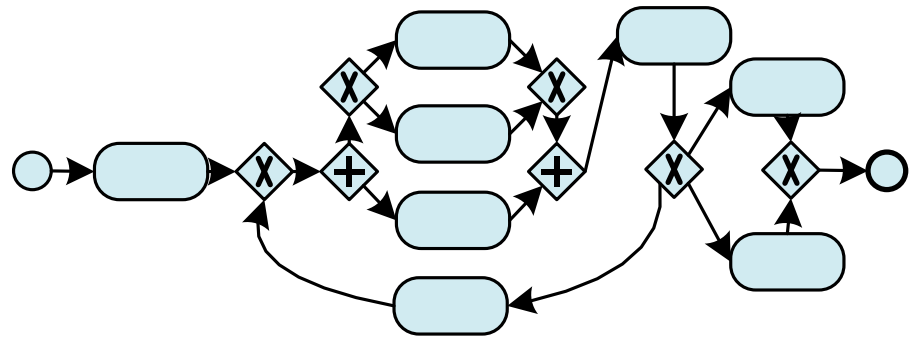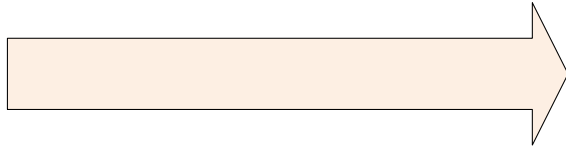# Play-Out (Classical use of models)
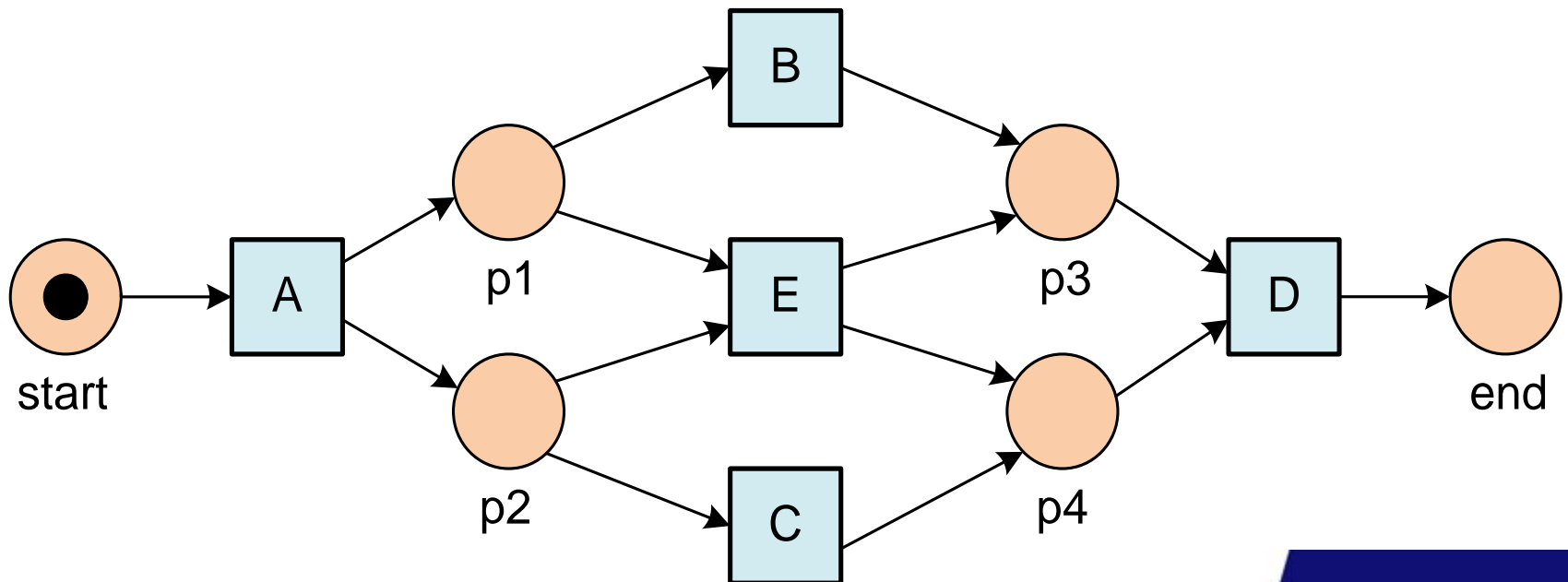
# Play-In



event log

process model

# Play-In

A B C D   A E D   A E D

A E D

A B C D   A C B D

A C B D

A E D   A C B D

# Example Process Discovery
## (AMC, 627 gynecological oncology patients, 24331 events)

# Replay



event log

process model

- extended model showing times, frequencies, etc.
- diagnostics
- predictions
- recommendations

# Replay

**A B C D**

# Replay

**A E D**

# Replay can detect problems

**A C D**



Petri net replay diagram. Red "X" error icon in top right. Yellow boxes labeled "Problem! token left behind" and "Problem! missing token" with red lightning bolts. Net flows: start → transition → p1/p2 → B, E → p3/p4 → transition → end. p3 shown with yellow token (missing), others with red tokens.

# Replay can extract timing information

# Hundreds of plug-ins available covering the whole process mining spectrum



**open-source (L-GPL)**

**Download from: www.processmining.org**

# Commercial Alternatives

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

# Three Key Observations

# #1 Alignments are essential!

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

move on log

move on model

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

move on model (harmless)

move on model

**#2 Models are like the glasses required to see and understand event data!**

Mayakovskaya
Маяковская

Sad Akvarium
Сад Акариум

Patriarshiye prudy
Патриаршие пруды

Barrikadnaya
Баррикадная

ul. Malaya Nikitskaya

"Bolshaya" Nikitskaya St

Pushkinskaya
Пушкинская

Tverskaya
Тверская

Chekhovskaya
Чеховская

б-р Тверской
б-т Tverskoy

Tverskaya ulitsa

ul. Bol'shaya Dmitrovka

Petrovka St

ul. Petrovka

Trubnaya
Трубная

Turgenevskaya
Тургеневская

Krasnye Vorota
Красные Ворота

Sad im. Baumana
Сад им. Баумана

Sretenskiy Bul'var
Сретенский Бульвар

Chistye Prudy
Чистые Пруды

Kuznetskiy Most
Кузнецкий Мост

Myasnitskaya ulitsa

ул. Покровка

Lubyanka
Лубянка

Teatral'naya
Театральная

Bol'shaya Nikitskaya ulitsa

б-т Nikitskiy

Ploshchad' Revolyutsii
Площадь Революции

Kitay-Gorod
Китай-Город

ul. Maroseyka

Milyutinskiy park
Милютинский парк

Kurskaya
Курская

Kurskiy vokzal
Курский вокзал

Chkalovskaya
Чкаловская

Starya Square

Vozdvizhenka St

Aleksandrovskiy Sad
Александровский Сад

Arbatskaya
Арбатская

r-n Arbat

Arbat District

ul. Znamenka

Borovitskaya
Боровицкая

ul. Varvarka

Ул. Солянка

Ул. Яузская

nab. Bernikovskaya

reka Yauza

ul. Novyy Arbat
ул. Новый Арбат

б-р Гоголевский

Smolenskaya
Смоленская

Moskvoretskaya Embankment

nab. Kremlevskaya

nab. Sofiyskaya

Ust'inskiy proyezd

Nikoloyamskaya ulitsa

Kropotkinskaya
Кропоткинская

ul. Volkhonka

Bolotnaya Square

Vodootvodnyy kanal
Водоотводный канал

ul. Prechistenka

ul. Ostozhenka

nab. Kadashevskaya

Novokuznetskaya
Новокузнецкая

Tret'yakovskaya
Третьяковская

Zamoskvorechye District

park im. NN Pryamikova
Парк им. Н. Пряникова

Marksistskaya
Марксистская

Taganskaya
Таганская

ул. Taganskaya

Skver Devich'ego Polya

Zubovskaya
Зубовская

Zubovskiy bul'var

ul. Prechistenka

Polyanka
Полянка

Zamoskvorechye District

Park Kul'tury
Парк Культуры

Крымский Вал

r-n Yakimanka

r-n Zamoskvorechye

# Alignments

**Joint work with Arya Adriansyah, Boudewijn van Dongen, Elham Ramezani, Dirk Fahland, Massimiliano de Leoni, et al.**

**a b e g**

r=1

m=1

$$fitness(\sigma, N) = \frac{1}{2}\left(1 - \frac{1}{6}\right) + \frac{1}{2}\left(1 - \frac{1}{6}\right) = 0.83333$$

# From "playing the token game" to optimal alignments …

observed trace: "abeg"

| a | b | » | e | g |
|---|---|---|---|---|
| a | b | d | e | g |

move in model only

# Another alignment

observed trace: "abcdeg"



| a | b | c | d | e | g |
|---|---|---|---|---|---|
| a | b | » | d | e | g |

**move in log only**

# Moves have costs

| ... | a | ... |
|-----|---|-----|
| ... | » | ... |

| ... | » | ... |
|-----|---|-----|
| ... | a | ... |

| ... | a | ... |
|-----|---|-----|
| ... | a | ... |

| ... | a | ... |
|-----|---|-----|
| ... | b | ... |

- **Standard cost function:**
  - $c(x,») = 1$
  - $c(»,y) = 1$
  - $c(x,y) = 0$, if $x=y$
  - $c(x,y) = \infty$, if $x \neq y$

# Any cost structure is possible

| ... | send-letter(John,2 weeks, $400) | ... |
|-----|----------------------------------|-----|
| ... | send-email(Sue,3 weeks,$500) | ... |

- **Similar activities** (more similarity implies lower costs).
- **Resource conformance** (done by someone that does not have the specified role).
- **Data conformance** (path is not possible for this customer).
- **Time conformance** (missed the legal deadline)

# Using Alignments

- An **optimal alignment** has the lowest possible costs.
- If multiple alignments are optimal, pick one or use all.
- Like an "**oracle**" revealing paths in the model.
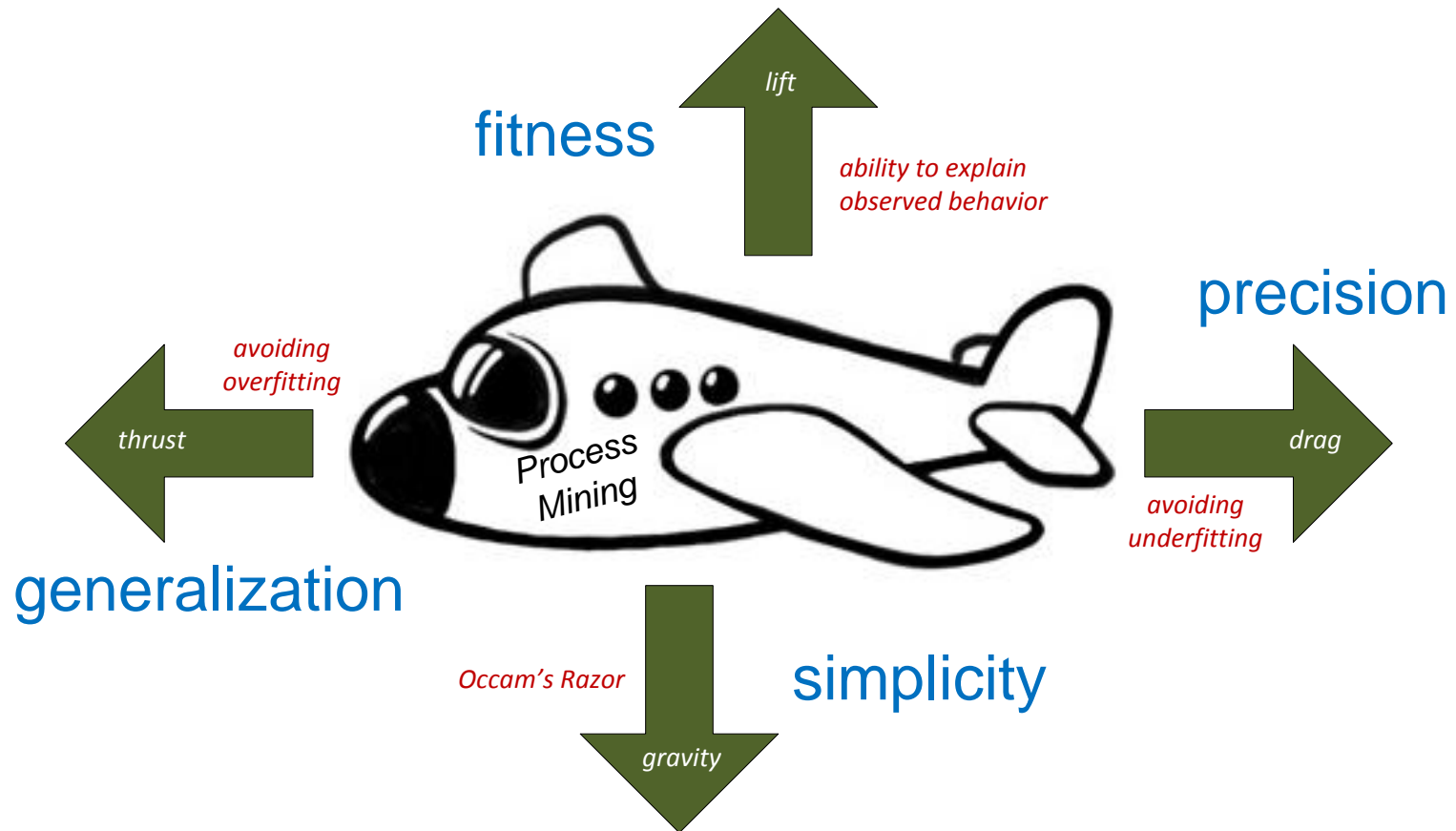- These paths can be used for further analysis!



- Can be used for **quantifying various types conformance** (possibly using a different cost function).

# Some pointers

- **Wil M. P. van der Aalst, Arya Adriansyah, Boudewijn F. van Dongen: Replaying history on process models for conformance checking and performance analysis. Wiley Interdisc. Rew.: Data Mining and Knowledge Discovery 2(2): 182-192 (2012)**

- **Arya Adriansyah, Jorge Munoz-Gama, Josep Carmona, Boudewijn F. van Dongen, Wil M. P. van der Aalst: Alignment Based Precision Checking. Business Process Management Workshops 2012: 137-149**

- **Arya Adriansyah, Boudewijn F. van Dongen, Wil M. P. van der Aalst: Conformance Checking Using Cost-Based Fitness Analysis. EDOC 2011: 55-64**

- **Massimiliano de Leoni, Wil M. P. van der Aalst, Boudewijn F. van Dongen: Data- and Resource-Aware Conformance Checking of Business Processes. BIS 2012: 48-59**

- **Joos C. A. M. Buijs, Boudewijn F. van Dongen, Wil M. P. van der Aalst: On the Role of Fitness, Precision, Generalization and Simplicity in Process Discovery. OTM Conferences (1) 2012: 305-322**

- **Elham Ramezani, Dirk Fahland, Wil M. P. van der Aalst: Where Did I Misbehave? Diagnostic Information in Compliance Checking. BPM 2012: 262-278**

- **A. Adriansyah, B.F. van Dongen, W.M.P. van der Aalst. Memory-Efficient Alignment of Observed and Modeled Behavior. BPM Center Report BPM-13-03, BPMcenter.org, 2013**
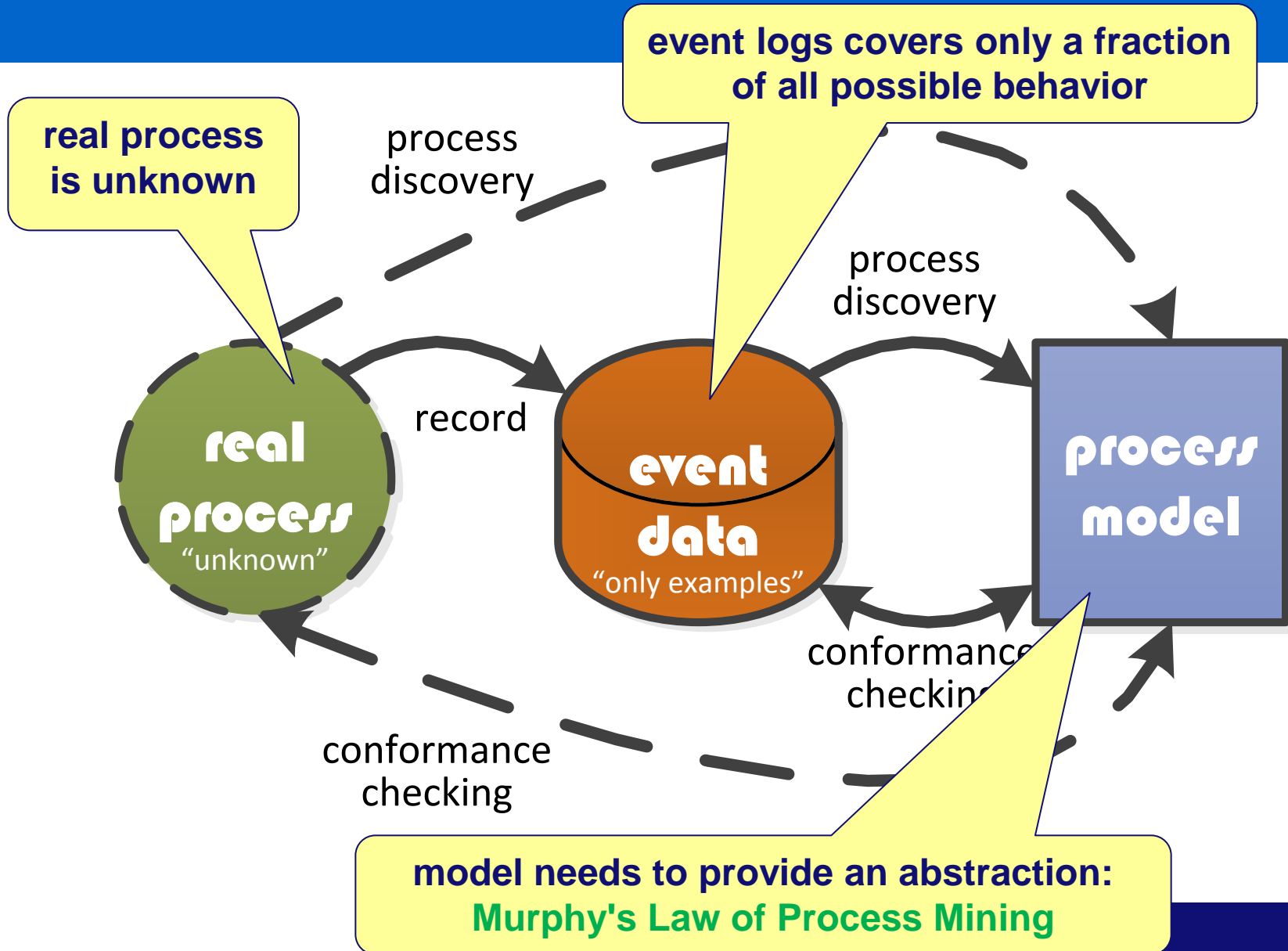
# Conformance: Taking a Step Back
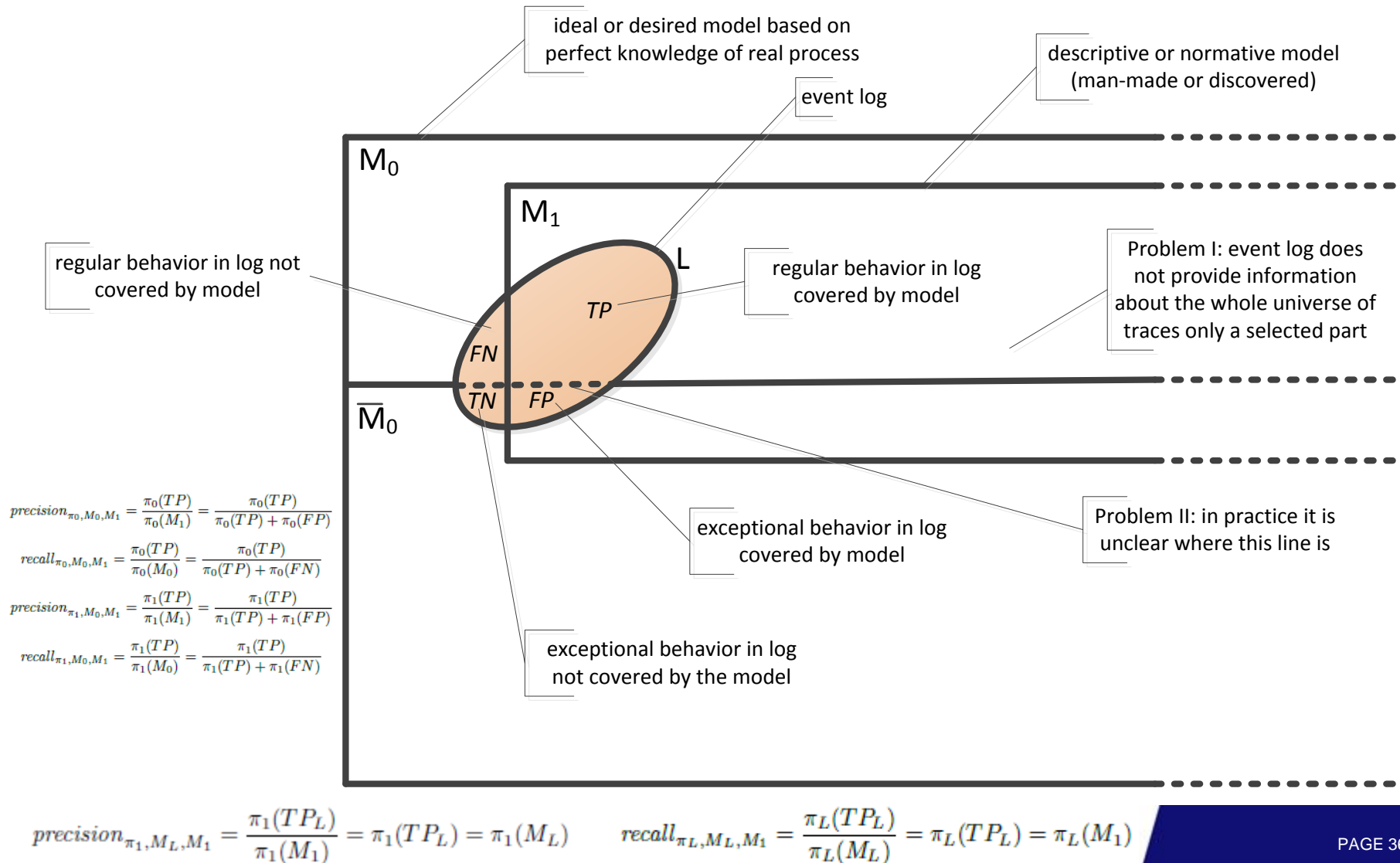
# Conventional Conformance Notions



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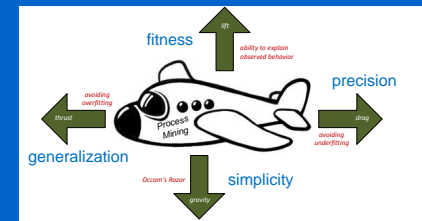
# Problem



event logs covers only a fraction of all possible behavior

real process is unknown

process discovery

process discovery

record

**real process**
"unknown"

**event data**
"only examples"

**process model**

conformance checking

conformance checking

model needs to provide an abstraction:
Murphy's Law of Process Mining

# Traditional Notions such as Precision and Recall do NOT Apply



ideal or desired model based on perfect knowledge of real process

descriptive or normative model (man-made or discovered)

event log

$M_0$

$M_1$

L

regular behavior in log not covered by model

regular behavior in log covered by model

Problem I: event log does not provide information about the whole universe of traces only a selected part

TP

$\overline{M}_0$

FN

TN    FP

Problem II: in practice it is unclear where this line is

exceptional behavior in log covered by model

exceptional behavior in log not covered by the model

$$precision_{\pi_0, M_0, M_1} = \frac{\pi_0(TP)}{\pi_0(M_1)} = \frac{\pi_0(TP)}{\pi_0(TP) + \pi_0(FP)}$$

$$recall_{\pi_0, M_0, M_1} = \frac{\pi_0(TP)}{\pi_0(M_0)} = \frac{\pi_0(TP)}{\pi_0(TP) + \pi_0(FN)}$$

$$precision_{\pi_1, M_0, M_1} = \frac{\pi_1(TP)}{\pi_1(M_1)} = \frac{\pi_1(TP)}{\pi_1(TP) + \pi_1(FP)}$$

$$recall_{\pi_1, M_0, M_1} = \frac{\pi_1(TP)}{\pi_1(M_0)} = \frac{\pi_1(TP)}{\pi_1(TP) + \pi_1(FN)}$$

$$precision_{\pi_1, M_L, M_1} = \frac{\pi_1(TP_L)}{\pi_1(M_1)} = \pi_1(TP_L) = \pi_1(M_L) \qquad recall_{\pi_L, M_L, M_1} = \frac{\pi_L(TP_L)}{\pi_L(M_L)} = \pi_L(TP_L) = \pi_L(M_1)$$

# Operationalizing the Four Conformance dimensions



- **Fitness** (fraction of observed behavior possible according to the model).
  - Measure at the case or event level.
  - How to continue after a deviation (where to put the "blame"), cf. duplicate and silent activities.
- **Precision** (avoiding underfitting; fraction of allowed behavior never observed).
  - Log only contains examples.
  - Metrics are e.g. based on escaping edges.
- **Generalization** (avoiding overfitting; probability that the next unseen case will not fit).
  - Reasoning about unseen behavior, strongly related to log completeness.
- **Simplicity** (Occam's Razor: the simplest of two or more competing theories is preferable)
  - Easy to operationalize (e.g., number of nodes or arcs).
  - Often subjective (valuation of AND/XOR/OR-split/joins).

# Some pointers

- **Wil M. P. van der Aalst: Mediating Between Modeled and Observed**
- **Behavior: The Quest for the "Right" Process. Seventh IEEE International Conference on Research Challenges in Information Science (RCIS 2013), (2013)**
- **Wil M. P. van der Aalst, Arya Adriansyah, Boudewijn F. van Dongen: Replaying history on process models for conformance checking and performance analysis. Wiley Interdisc. Rew.: Data Mining and Knowledge Discovery 2(2): 182-192 (2012)**
- **Joos C. A. M. Buijs, Boudewijn F. van Dongen, Wil M. P. van der Aalst: On the Role of Fitness, Precision, Generalization and Simplicity in Process Discovery. OTM Conferences (1) 2012: 305-322**
- **Wil M. P. van der Aalst: Process Mining - Discovery, Conformance and Enhancement of Business Processes. Springer 2011, isbn 978-3-642-19344-6, pp. I-XVI, 1-352**

# Representational Bias
# in Process Mining
**(not about visualization)**

**Joint work with Joos Buijs, Sander Leemans, and Boudewijn van Dongen.**

# Typical Representational Bias

- **(Labeled) Petri Nets, WF-nets, etc.**
- **Subsets of**
  - **BPMN diagrams,**
  - **UML Activity Diagrams,**
  - **Event-Driven Process Chains (EPCs),**
  - **YAWL,**
  - **etc.**
- **Transition Systems**
- **(Hidden) Markov Models**
- **…**

# Huge Search Space When Discovering a Petri Net, BPMN model, and the like …

# … with just a few interesting candidates

# Alternative Representational Bias

1. **C-nets** (XOR/AND/OR-split/join graphs; more likely to be sound due to declarative semantics).

2. **Declare models** (constraint based, grounded in LTL; anything is possible unless forbidden)

3. **Process Trees** (similar to subsets of various process algebras; sound by structure)



today's focus

# Another Representational Bias: Process Trees



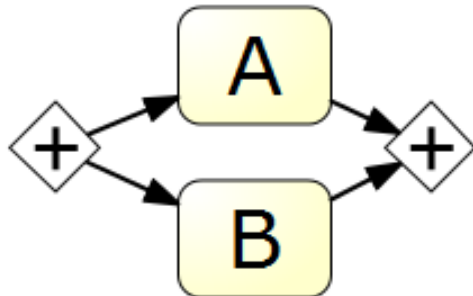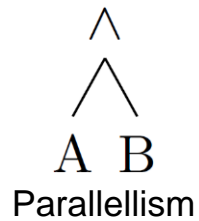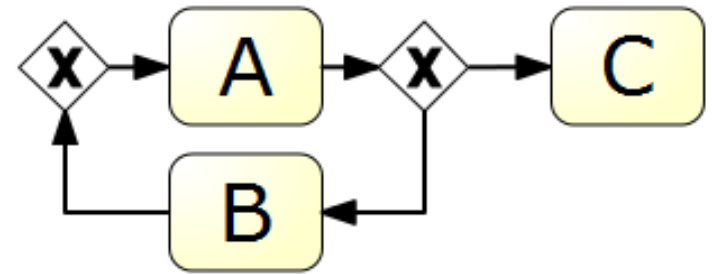- **Always sound because of the block structure**
- **Also Loop and OR operator**

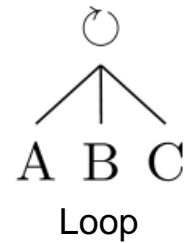A (BA)*

# Petri Net Semantics
**(used for comparison and conformance checking only)**



Sequence

Exclusive Choice

Loop

Parallellism

Or Choice

# … and BPMN.
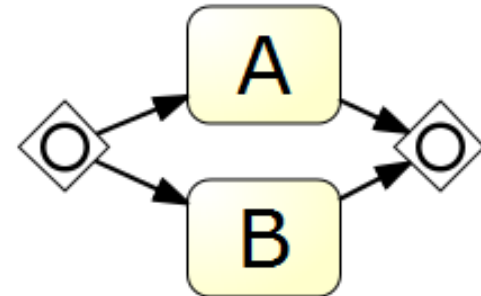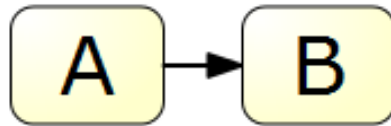


Sequence

Exclusive Choice

Parallellism

Or Choice

Loop

# A Discovery Algorithm Using Process Trees: Evolutionary Tree Miner (ETM)

- **Process trees** as representation (= limit search space to "good" models).
- **Genetic** approach (= very flexible)
- Fitness function uses all **four criteria** (= seamlessly balance the different "forces")

# Population Change

# Example

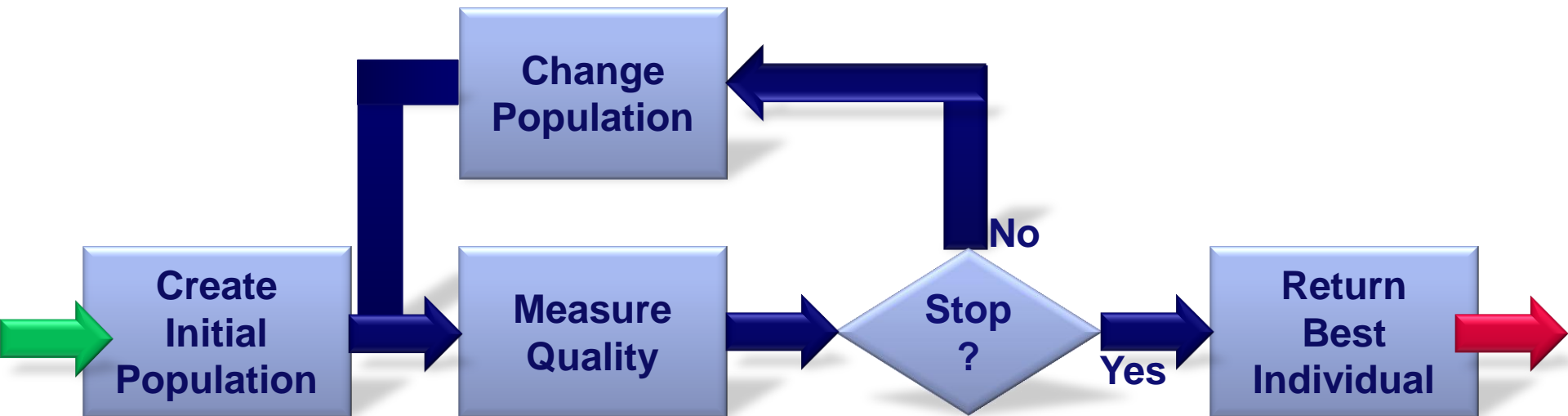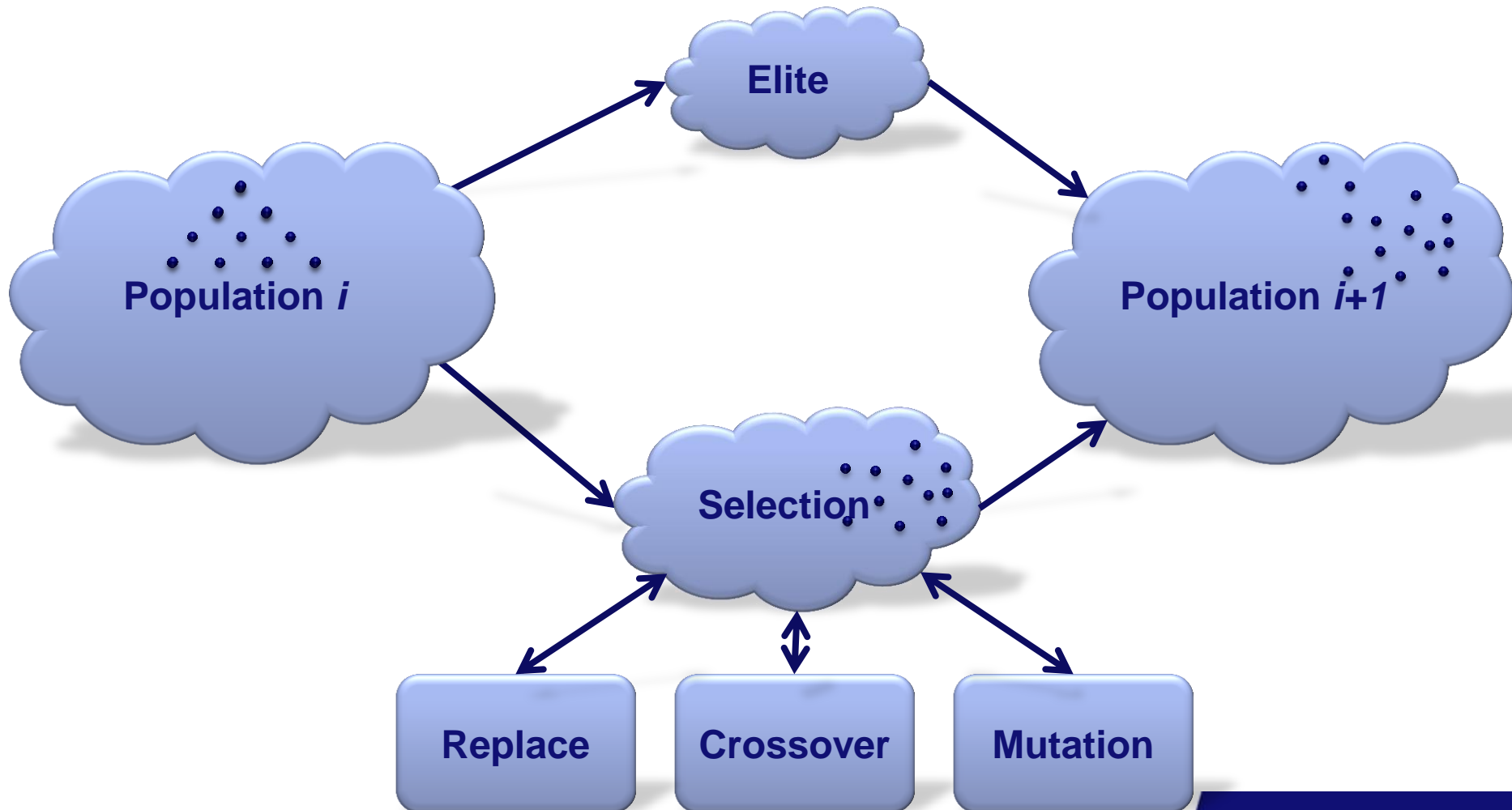| Trace | # |
|---|---|
| A B C D E G | 6 |
| A B C D F G | 38 |
| A B D C E G | 12 |
| A B D C F G | 26 |
| A B C F G | 8 |
| A C B E G | 1 |
| A D B C F G | 1 |
| A D B C E G | 1 |
| A D C B F G | 4 |
| A C D B F G | 2 |
| A C B F G | 1 |



A = send e-mail, B = check credit,
C = calculate capacity, D = check system,
E = accept, F = reject, G = send e-mail

**alpha miner**



*sound*

*... lucky*

| f: 0,992 | p: 0,995 |
|----------|----------|
| s: 1,000 | g: 0,889 |

**low fitness**

**ILP miner**



*sound*

| f: 1,000 | p: 0,784 |
|----------|----------|
| s: 0,933 | g: 0,830 |

**low precision**

**language-based region miner**



*sound*

| f: 0,992 | p: 0,957 |
|----------|----------|
| s: 1,000 | g: 0,889 |

**low fitness**

# Conventional Algorithms (2/3)

**heuristic miner**



**multi-phase miner**

# Conventional Algorithms (3/3)

**genetic miner**



unsound

| f: 1,000 | p: 0,922 |
| --- | --- |
| s: 0,737 | g: 0,790 |

**state-based region miner**



sound

| f: 1,000 | p: 0,893 |
| --- | --- |
| s: 0,933 | g: 0,830 |

# Genetic Mining (ETM) While Considering Only One Criterion



(a) Only replay fitness

| f: 1,000 | p: 0,341 |
| s: 0,737 | g: 0,681 |

(b) Only Precision

| f: 0,449 | p: 1,000 |
| s: 0,400 | g: 0,797 |

(c) Only Simplicity

| f: 0,504 | p: 0,587 |
| s: 1,000 | g: 0,661 |

(d) Only Generalization

| f: 0,961 | p: 0,394 |
| s: 0,923 | g: 0,916 |

best value possible for this log

**ETM with weight zero to three out of four perspectives.**

# Considering Replay Fitness and One Other Criterion



(a) Replay Fitness and Precision

(b) Replay Fitness and Simplicity

(c) Replay Fitness and Generalization

(a) no precision

# Considering All Four Criteria with Emphasis on Fitness

| f: 1,000 | p: 0,923 |
|----------|----------|
| s: 1,000 | g: 0,889 |



**fitness has weight 10**

# Initial Model Versus Discovered Model

**Better than existing algorithms (but patience is needed)!**

| Trace | # |
|---|---|
| A B C D E G | 6 |
| A B C D F G | 38 |
| A B D C E G | 12 |
| A B D C F G | 26 |
| A B C F G | 8 |
| A C B E G | 1 |
| A D B C F G | 1 |
| A D B C E G | 1 |
| A D C B F G | 4 |
| A C D B F G | 2 |
| A C B F G | 1 |

**Discovered model outperforms initial model with respect too all criteria!**

**simulated**

| f: 1,000 | p: 0,893 |
|---|---|
| s: 0,933 | g: 0,830 |

**discovered by ETM**

| f: 1,000 | p: 0,923 |
|---|---|
| s: 1,000 | g: 0,889 |

**1) Carefully choose your representational bias during discovery: Unrelated to presentation/visualization!**

Lessons Learned

**2) Consider all conformance dimensions (replay fitness, precision, generalization, and simplicity)!**

# Some pointers

- **Wil M. P. van der Aalst, Arya Adriansyah, Boudewijn F. van Dongen: Causal Nets: A Modeling Language Tailored towards Process Discovery. CONCUR 2011: 28-42**

- **Joos C. A. M. Buijs, Boudewijn F. van Dongen, Wil M. P. van der Aalst: On the Role of Fitness, Precision, Generalization and Simplicity in Process Discovery. OTM Conferences (1) 2012: 305-322**

- **Joos C. A. M. Buijs, Boudewijn F. van Dongen, Wil M. P. van der Aalst: A genetic algorithm for discovering process trees. IEEE Congress on Evolutionary Computation 2012: 1-8**

- **S.J.J. Leemans, D. Fahland, W.M.P. van der Aalst. Discovering Block-Structured Process Models From Event Logs – A Constructive Approach. BPM Center Report BPM-13-06, BPMcenter.org, 2013**

# Mediating Between a Reference Model and Real Observed Behavior

**Joint work with Joos Buijs, Boudewijn van Dongen, and Dirk Fahland.**

# Compromise Based on Two Main Forces



**Various techniques to compare graphs, e.g., edit distance notions (add, remove, replace).**

# Force Between Reference Model and Candidate Model Can be Viewed as a 5th Conformance Notion

# Example From CoSeLoG Project



(d) Similarity x1

| sim: 1.000 | 0 edits |
|---|---|
| f: 0.744 | p: 0.785 |

| sim: 0.650 | 42 edits |
|---|---|
| f: 0.974 | p: 0.933 |
| s: 0.613 | g: 0.747 |

| sim: 0.650 | 42 edits |
|---|---|
| f: 0.974 | p: 0.933 |
| s: 0.613 | g: 0.747 |

# Some pointers

- J.C.A.M. Buijs ,M. La Rosa, H.A. Reijers, B.F. van Dongen, and W.M.P. van der Aalst: Improving Business Process Models using Observed Behavior, SIMPDA 2012 post-proceedings, Lecture Notes in Business Information Processing, 2013.

- Joos C. A. M. Buijs, Boudewijn F. van Dongen, Wil M. P. van der Aalst: On the Role of Fitness, Precision, Generalization and Simplicity in Process Discovery. OTM Conferences (1) 2012: 305-322

- Dirk Fahland, Wil M. P. van der Aalst: Repairing Process Models to Reflect Reality. BPM 2012: 229-245.

- Dirk Fahland, Wil M. P. van der Aalst: Simplifying discovered process models in a controlled manner. Inf. Syst. 38(4): 585-605 (2013).

# Mining Configurable Process Models

**Joint work with Joos Buijs, Boudewijn van Dongen, and Florian Gottschalk.**

# Two variants of the same process …

# Configurable Process Model

# Variants of the same process

# Approach 1: Merge Separately Mined Models

| | Overall | Fitness | Precision | Simplicity | Generalization | Size | #C.P. | Similarity |
|---|---|---|---|---|---|---|---|---|
| Combined | 0.989 | 0.999 | 0.999 | 0.981 | 0.220 | 53 | 4 | - |
| Variant 0 | 0.986 | 0.995 | 0.995 | 0.981 | 0.235 | 14 | 3 | 0.418 |
| Variant 1 | 0.989 | 1.000 | 1.000 | 0.981 | 0.263 | 16 | 3 | 0.464 |
| Variant 2 | 0.989 | 1.000 | 1.000 | 0.981 | 0.174 | 10 | 3 | 0.317 |
| Variant 3 | 0.989 | 1.000 | 1.000 | 0.981 | 0.264 | 16 | 3 | 0.464 |

# Approach 2

# Results Approach 2 for running example



(a) Process model discovered from combined event log

(b) Process model individualized for event log 1 (c) Process model individualized for event log 2

(d) Process model individualized for event log 3 (e) Process model individualized for event log 4

(f) Configurable process model of ... er merging models (b) throug ...

| | Overall | Fitness | Precision | Simplicity | Generalization | Size | #C.P. | Similarity |
|---|---|---|---|---|---|---|---|---|
| Combined | 0.958 | 0.974 | 0.921 | 0.968 | 0.212 | 46 | 4 | - |
| Variant 0 | 0.981 | 0.995 | 0.995 | 0.968 | 0.232 | 12 | 3 | 0.414 |
| Variant 1 | 0.984 | 1.000 | 1.000 | 0.968 | 0.246 | 13 | 3 | 0.441 |
| Variant 2 | 0.984 | 1.000 | 1.000 | 0.968 | 0.180 | 10 | 3 | 0.357 |
| Variant 3 | 0.869 | 0.886 | 0.649 | 0.968 | 0.232 | 14 | 3 | 0.467 |

(g) Quality statistics of the configurable process model of (f)

# Approach 3

(a) Configurable process model discovered when first discovering the process model and then the configurations

| | Overall | Fitness | Precision | Simplicity | Generalization | Size | #C.P. | Similarity |
|---|---|---|---|---|---|---|---|---|
| Combined | 0.988 | 1.000 | 0.981 | 0.986 | 0.374 | 42 | 11 | - |
| Variant 0 | 0.990 | 1.000 | 0.990 | 0.986 | 0.400 | 20 | 6 | 0.645 |
| Variant 1 | 0.992 | 1.000 | 1.000 | 0.986 | 0.408 | 20 | 7 | 0.645 |
| Variant 2 | 0.992 | 1.000 | 1.000 | 0.986 | 0.285 | 13 | 8 | 0.473 |
| Variant 3 | 0.977 | 1.000 | 0.922 | 0.986 | 0.496 | 24 | 6 | 0.727 |

(b) Quality statistics of the configurable process model of (a)

# Approach 4

# Results Approach 4 for running example



(a) Configurable ... odel discov ... using the integrate ... y approach

| | Overall | Fitness | Precision | Simplicity | Generalization | Size | #C.P. | Similarity |
|---|---|---|---|---|---|---|---|---|
| Combined | 0.983 | 0.962 | 0.999 | 1.000 | 0.684 | 12 | 6 | - |
| Variant 0 | 0.996 | 0.995 | 0.994 | 1.000 | 0.738 | 10 | 2 | 0.909 |
| Variant 1 | 0.957 | 0.894 | 1.000 | 1.000 | 0.723 | 10 | 3 | 0.909 |
| Variant 2 | 0.998 | 1.000 | 1.000 | 1.000 | 0.614 | 8 | 5 | 0.800 |
| Variant 3 | 0.961 | 0.905 | 1.000 | 1.000 | 0.741 | 10 | 3 | 0.909 |

(b) Quality statistics of the configurable process model of (a)

# Genetic algorithms are versatile and can consider different forces at the same time

# Some pointers

- **J.C.A.M. Buijs, B.F. van Dongen, and W.M.P. van der Aalst: Mining Configurable Process Models from Collections of Event Logs, under review, 2013.**

- **Florian Gottschalk, Teun A. C. Wagemakers, Monique H. Jansen-Vullers, Wil M. P. van der Aalst, Marcello La Rosa: Configurable Process Models: Experiences from a Municipality Case Study. CAiSE 2009: 486-500**

- **Florian Gottschalk, Wil M. P. van der Aalst, Monique H. Jansen-Vullers: Merging Event-Driven Process Chains. OTM Conferences (1) 2008: 418-426**

- **Wil M. P. van der Aalst: Business Process Configuration in the Cloud: How to Support and Analyze Multi-tenant Processes? ECOWS 2011: 3-10**

# Decomposing Process Mining Problems

**Joint work with Eric Verbeek, Jorge Munoz , and Josep Carmona.**

# Big Data: Opportunities and Challenges

# System Net



a = register request
b = examine file
c = check ticket
d = decide
e = reinitiate request
f = send acceptance letter
g = pay compensation
h = send rejection letter

- **Labeled Petri net (P,T,F,l)**
- **Silent transitions and visible transitions (unique or not),**
- **One initial marking $M_{init}$, one final marking $M_{final}$**

# Traces and Logs

- **A system net SN has a set of <span style="color:red">possible visible traces</span> φ(SN) starting in $M_{init}$ and ending $M_{final}$ only showing the visible steps.**

- **An <span style="color:red">event log L</span> is a multiset of traces.**

- **Two main process mining problems:**

  1. **Conformance checking: Given L and SN, evaluate the "conformance" (e.g., fitness, precision, generalization, etc.) of L and φ(SN)**

  2. **Process discovery: Given L, create SN such that the conformance of L and φ(SN) is "as good as possible"**

# Valid Decomposition



- **Union of subnets is original net**
- **No shared places**
- **Shared transitions are visible and have unique label**

# Another Decomposition



**Requirement**

- **Union of subnets is original net**
- **No shared places**
- **Shared transitions are visible and unique**

# Maximal Valid Decomposition

# Maximal Decomposition



- **Construction: group arcs iteratively**
- **Maximal decomposition is unique**
- **There is always a valid decomposition**

# Non-unique visible labels



- **Union of subnets is original net**
- **No shared places**
- **Shared transitions are visible and unique**

# Conformance checking can be decomposed !!!

- **Let L be an event log, SN a system net, and $D=\{SN^1, SN^2, \ldots SN^n\}$ a valid decomposition**

- **$L^i$ is the sublog of $SN^i$ (L projected onto visible transitions of $SN^i$)**



**L is perfectly fitting SN**
**if and only if**
**each projected log is $L^i$ is perfectly fitting $SN^i$**

**a,b,c,d,e,c,d,g,f**



**Etc.**

$$\gamma_3 = \begin{array}{|c|c|c|c|c|c|c|c|c|c|c|c|} \hline 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 \\ \hline a & b & c & d & e & c & \gg & d & \gg & g & f & \gg \\ \hline a & b & c & d & e & c & \tau & d & \tau & g & f & \tau \\ \hline t1 & t3 & t4 & t5 & t6 & t4 & t2 & t5 & t7 & t9 & t8 & t11 \\ \hline \end{array}$$

$$\gamma_3^1 = \begin{array}{|c|} \hline 1 \\ \hline a \\ \hline a \\ \hline t1 \\ \hline \end{array} \qquad \gamma_3^2 = \begin{array}{|c|c|c|c|c|c|} \hline 1 & 2 & 4 & 5 & 7 & 8 \\ \hline a & b & d & e & \gg & d \\ \hline a & b & d & e & \tau & d \\ \hline t1 & t3 & t5 & t6 & t2 & t5 \\ \hline \end{array} \qquad \gamma_3^3 = \begin{array}{|c|c|c|c|} \hline 1 & 3 & 5 & 6 \\ \hline a & c & e & c \\ \hline a & c & e & c \\ \hline t1 & t4 & t6 & t4 \\ \hline \end{array}$$

$$\gamma_3^4 = \begin{array}{|c|c|c|c|} \hline 3 & 4 & 6 & 8 \\ \hline c & d & c & d \\ \hline c & d & c & d \\ \hline t4 & t5 & t4 & t5 \\ \hline \end{array} \qquad \gamma_3^5 = \begin{array}{|c|c|c|c|c|c|} \hline 4 & 5 & 8 & 9 & 10 & 11 \\ \hline d & e & d & \gg & g & f \\ \hline d & e & d & \tau & g & f \\ \hline t5 & t6 & t5 & t7 & t9 & t8 \\ \hline \end{array} \qquad \gamma_3^6 = \begin{array}{|c|c|c|} \hline 10 & 11 & 12 \\ \hline g & f & \gg \\ \hline g & f & \tau \\ \hline t9 & t8 & t11 \\ \hline \end{array}$$

# So What?
# Quantifying Conformance

- **Exact result: Fraction of cases perfectly fitting SN equals the fraction of cases for which each projected log is $L^i$ is perfectly fitting $SN^i$.**

- **Bound: For fitness at the event level (costs based alignment) it is possible to compute an optimistic value.**

# Discovery can also be distributed!

1. **Assume the set of activities is split in overlapping sets.**

2. **Split log L in sublogs $L^i$ based on these sets.**

3. **Discover a model $SN^i$ per sublog.**

4. **Merge the models $SN^i$ into SN and return result.**

All the earlier guarantees hold, e.g., L is perfectly fitting SN if and only if each projected log is $L^i$ is perfectly fitting $SN^i$

**Das Orakel**

$$L = [\langle a,b,c,d,h \rangle^{30}, \langle a,c,b,d,e,b,c,d,j,g,f,k \rangle^{28}, \langle a,c,b,d,e,c,b,d,h \rangle^{28}, \langle a,c,b,d,e,c,b,d,j,f,g,k \rangle^{27}, \langle a,b,c,d,e,c,i,d,j,g,f,k \rangle^{26}, \langle a,b,c,d,e,i,c,d,h \rangle^{26}, \langle a,i,c,d,e,c,b,d,j,g,f,k \rangle^{25}, \langle a,i,c,d,j,f,g,k \rangle^{24}, \langle a,c,b,d,e,b,c,d,h \rangle^{24}, \langle a,b,c,d,e,b,c,d,j,g,f,k \rangle^{24}, \langle a,c,b,d,e,c,i,d,h \rangle^{22}, \langle a,b,c,d,e,c,i,d,h \rangle^{22}, \langle a,c,i,d,j,f,g,k \rangle^{21}, \langle a,c,i,d,j,g,f,k \rangle^{20}, \langle a,c,i,d,e,b,c,d,h \rangle^{18}, \langle a,c,i,d,e,c,i,d,h \rangle^{17}, \langle a,i,c,d,h \rangle^{17}, \langle a,c,b,d,j,f,g,k \rangle^{17}, \langle a,c,b,d,h \rangle^{15}, \langle a,i,c,d,e,i,c,d,j,f,g,k \rangle^{14}, \langle a,c,i,d,e,c,b,d,h \rangle^{14}, \langle a,c,i,d,e,c,b,d,j,f,g,k \rangle^{12}, \langle a,b,c,d,e,i,c,d,j,f,g,k \rangle^{12}, \langle a,b,c,d,e,c,b,d,h \rangle^{11}, \langle a,c,b,d,j,f,g,k \rangle^{10}, \langle a,i,c,d,j,g,f,k \rangle^{9}, \langle a,i,c,d,e,b,c,d,h \rangle^{9}, \langle a,i,c,d,e,c,b,d,h \rangle^{9}, \langle a,i,c,d,e,c,b,d,j,f,g,k \rangle^{8}, \langle a,c,b,d,e,b,c,d,j,f,g,k \rangle^{8}, \langle a,b,c,d,e,b,c,d,h \rangle^{7}, \langle a,c,b,d,e,i,c,d,j,g,f,k \rangle^{7}, \langle a,i,c,d,e,c,i,d,j,g,f,k \rangle^{6}, \langle a,c,i,d,h \rangle^{6}, \langle a,c,b,d,e,i,c,d,h \rangle^{6}, \langle a,i,c,d,e,b,c,d,j,g,f,k \rangle^{6}, \langle a,b,c,d,j,f,g,k \rangle^{5}, \langle a,i,c,d,e,i,c,d,h \rangle^{5}, \langle a,i,c,d,e,c,i,d,h \rangle^{4}, \langle a,c,i,d,e,i,c,d,h \rangle^{4}, \langle a,b,c,d,e,c,i,d,j,f,g,k \rangle^{4}, \langle a,b,c,d,e,c,b,d,j,g,f,k \rangle^{4}, \langle a,b,c,d,j,g,f,k \rangle^{3}, \langle a,c,i,d,e,b,c,d,j,f,g,k \rangle^{3}, \langle a,b,c,d,e,b,c,d,j,f,g,k \rangle^{3}, \langle a,i,c,d,e,b,c,d,j,f,g,k \rangle^{3}, \langle a,b,c,d,e,c,b,d,j,f,g,k \rangle^{3}, \langle a,c,i,d,e,i,c,d,j,f,g,k \rangle^{2}, \langle a,c,i,d,e,i,c,d,e,b,c,d,e,b,c,d,j,g,g,f,k \rangle^{2}, \langle a,b,c,d,e,i,c,d,j,g,f,k \rangle^{2}, \langle a,c,b,d,e,c,i,d,j,g,f,k \rangle^{2}, \langle a,c,i,d,e,c,i,d,j,f,g,k \rangle^{2}, \langle a,c,b,d,e,b,c,d,e,c,b,d,j,g,f,k \rangle^{2}, \langle a,c,i,d,e,c,b,d,j,g,f,k \rangle^{2}, \langle a,c,b,d,e,i,c,d,j,f,g,k \rangle^{2}, \langle a,c,i,d,e,i,c,d,e,i,c,d,j,g,f,k \rangle^{1}, \langle a,c,b,d,e,c,i,d,j,f,g,k \rangle^{1}, \langle a,i,c,d,e,c,i,d,e,c,i,d,e,i,c,d,e,i,c,d,j,g,f,k \rangle^{1}, \langle a,i,c,d,e,c,i,d,j,f,g,k \rangle^{1}, \langle a,c,i,d,e,c,i,d,j,g,f,k \rangle^{1}]$$

$$A^1 = \{a,b,d,e,i\}$$

$$A^2 = \{a,c,d,e\}$$

$$A^3 = \{d,e,h,j\}$$

$$A^4 = \{f,g,h,j,k\}$$

- **Let's use the Alpha algorithm**
- **Alpha algorithm is not very suitable for discovering transition bordered subnets.**
- **By adding start and end activities we get (in this case) perfectly fitting subnets.**

$$L' = [\langle \top \rangle \cdot \sigma \cdot \langle \bot \rangle \mid \sigma \in L_o] = [\langle \top, a,b,c,d,h, \bot \rangle^{30}, \langle \top, a,c,b,d,e,b,c,d,j,g,f,k, \bot \rangle^{28}, \langle \top, a,c,b,d,e,c,b,d,h, \bot \rangle^{28}, \langle \top, a,c,b,d,e,c,b,d,j,f,g,k, \bot \rangle^{27}, \langle \top, a,b,c,d,e,c,i,d,j,g,f,k, \bot \rangle^{26}, \ldots].$$

# Discover model per sublog (Alpha algorithm)



$$L^1 = [\langle \top, a, b, d, \bot \rangle^{30}, \langle \top, a, b, d, e, b, d, \bot \rangle^{28}, \langle \top, a, b, d, e, b, d, \bot \rangle^{28}, \langle \top, a, b, d, e, b, d, \bot \rangle^{27}, \langle \top, a, b, d, e, i, d, \bot \rangle^{26}, \ldots]$$

$$L^2 = [\langle \top, a, c, d, \bot \rangle^{30}, \langle \top, a, c, d, e, c, d, \bot \rangle^{28}, \langle \top, a, c, d, e, c, d, \bot \rangle^{28}, \langle \top, a, c, d, e, c, d, \bot \rangle^{27}, \langle \top, a, c, d, e, c, d, \bot \rangle^{26}, \ldots]$$

$$L_3 = [\langle \top, d, h, \bot \rangle^{30}, \langle \top, d, e, d, j, \bot \rangle^{28}, \langle \top, d, e, d, h, \bot \rangle^{28}, \langle \top, d, e, d, j, \bot \rangle^{27}, \langle \top, d, e, d, j, \bot \rangle^{26}, \ldots]$$

$$L_4 = [\langle \top, h, \bot \rangle^{30}, \langle \top, j, g, f, k, \bot \rangle^{28}, \langle \top, h, \bot \rangle^{28}, \langle \top, j, f, g, k, \bot \rangle^{27}, \langle \top, j, g, f, k, \bot \rangle^{26}, \ldots]$$

# Merge models



**L is perfectly fitting SN because each projected log is L$^i$ is perfectly fitting SN$^i$**

# Simplify by removing redundant places and initialization



**Log is perfectly fitting this model !**

- **Few large process mining tasks or many smaller process mining tasks.**

- **Choice of level is driven by computation time and desired diagnostics!**

- **All implemented in ProM.**

ProM
process mining workbench

# Diagnosis



**Diagnose Subprocess**
Detect an unfitting subprocess,
analyze it in isolation,
and diagnose the cause of the problems.

**Diagnose Non Fitting Net**
Detect all unfitting subprocesses,
compose the net that contains all them,
and diagnose the cause of the problems.

# Example: Conformance Checking based on SESEs

k=inf (no decompostion)

some overhead for easy problems

intractable problems become tractable

| | | | [9,11] | | k = 50 | | | | k = | | | | k = 200 | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $P$ | $T$ | $f$ | $t$ | $S/B$ | >5 | nf | $t$ | $S/B$ | >5 | nf | $t$ | $S/B$ | >5 | nf | $t$ |
| prAm6 | 347 | 363 | 0.92 | 75 | 129/57 | 29 | 7(3%) | 423 | 62/27 | 14 | 1(9%) | 323 | 27/12 | 7 | 1(10%) | 180 |
| prBm6 | 317 | 317 | 1 | 88 | 93/38 | 22 | 0(0%) | 608 | 66/29 | 14 | 0(0%) | 318 | 36/16 | 8 | 0(0%) | 114 |
| prCm6 | 317 | 317 | 0.57 | 2743 | 93/38 | 22 | 58(92%) | 189 | 66/29 | 14 | 41(94%) | 185 | 36/16 | 8 | 22(96%) | 502 |
| prDm6 | 529 | 429 | - | - | 105/34 | 33 | 5(8%) | 1386 | 60/23 | 18 | 4(14%) | 986 | 33/15 | 9 | 4(23%) | 1284 |
| prEm6 | 277 | 275 | 0.97 | 3566 | 82/35 | 20 | 2(5%) | 529 | 35/13 | 11 | 2(5%) | 343 | 15/7 | 5 | 2(6%) | 211 |
| prFm6 | 362 | 299 | - | - | 108/43 | 28 | 2(6%) | 1667 | 57/23 | 15 | 2(21%) | 813 | 21/9 | 5 | 1(23%) | 562 |

stopped after 10 hours

problems are often local

suitable k-value depends on model/log

**k = maximal number of arcs in one SESE, P = # places, T = # transitions, f = fitness, t = time in seconds, S = # transition bounded SESEs, B = # bridges, >5 = # more than 5 arcs, nf = # non-fitting parts.**

**ProM**
process mining workbench

# Some pointers

- **Wil M.P. van der Aalst. Decomposing Petri Nets for Process Mining: A Generic Approach. BPM Center Report BPM-12-20, BPMcenter.org, 2012**

- **Wil M. P. van der Aalst: Distributed Process Discovery and Conformance Checking. FASE 2012: 1-25**

- **Wil M. P. van der Aalst: Decomposing Process Mining Problems Using Passages. Petri Nets 2012: 72-91**

- **H. M. W. (Eric) Verbeek, Wil M. P. van der Aalst: An Experimental Evaluation of Passage-Based Process Discovery. Business Process Management Workshops 2012: 205-210**

- **Wil M.P. van der Aalst and Eric Verbeek. Process Discovery and Conformance Checking Using Passages. BPM Center Report BPM-12-21, BPMcenter.org, 2012**

- **J. Munoz-Gama, J. Carmona, W. van der Aalst: Hierarchical Conformance Checking of Process Models Based on Event Logs. In: Applications and Theory of Petri Nets, 2013**

# Conclusion

# Conclusion (1/2)

- **Alignments are essential for relating observed and modeled behavior!**
- **Conformance has (at least) four dimensions!**
- **Representational bias is important (and should not be confused with visualization)!**
- **New questions are emerging:**
  - **mediating between a reference model and observed behavior**
  - **discovering configurable process models**
- **Decomposing process mining problems to deal with Big Data.**

# Conclusion (2/2)

**Still many challenging and highly relevant open problems in process mining!**



**process model analysis**
(simulation, verification, etc.)

performance-oriented questions, problems and solutions

**process mining**

compliance-oriented questions, problems and solutions

**data-oriented analysis**
(data mining, machine learning, business intelligence)

**Harvard Business Review**

**Data Scientist: The Sexiest Job of the 21st Century**
by Thomas H. Davenport and D.J. Patil

Wil M. P. van der Aalst

**Process Mining**

Discovery, Conformance and Enhancement of Business Processes

Springer

**ProM**
process mining workbench

**processmining.org**