

Conformance Checking (1/2)

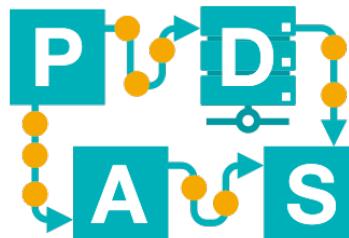
Lecture 12

prof.dr.ir. Wil van der Aalst

www.vdaalst.com @wvdaalst

www.pads.rwth-aachen.de

BPI-L12

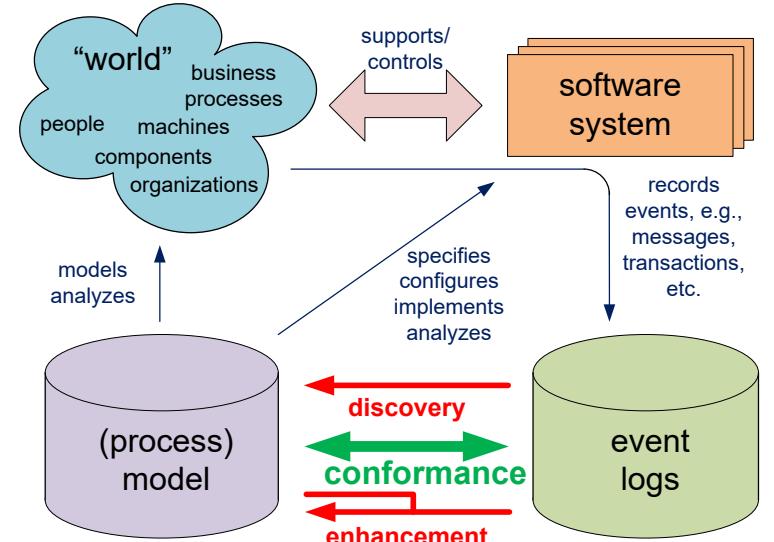


Chair of Process
and Data Science

RWTHAACHEN
UNIVERSITY

Conformance checking approaches

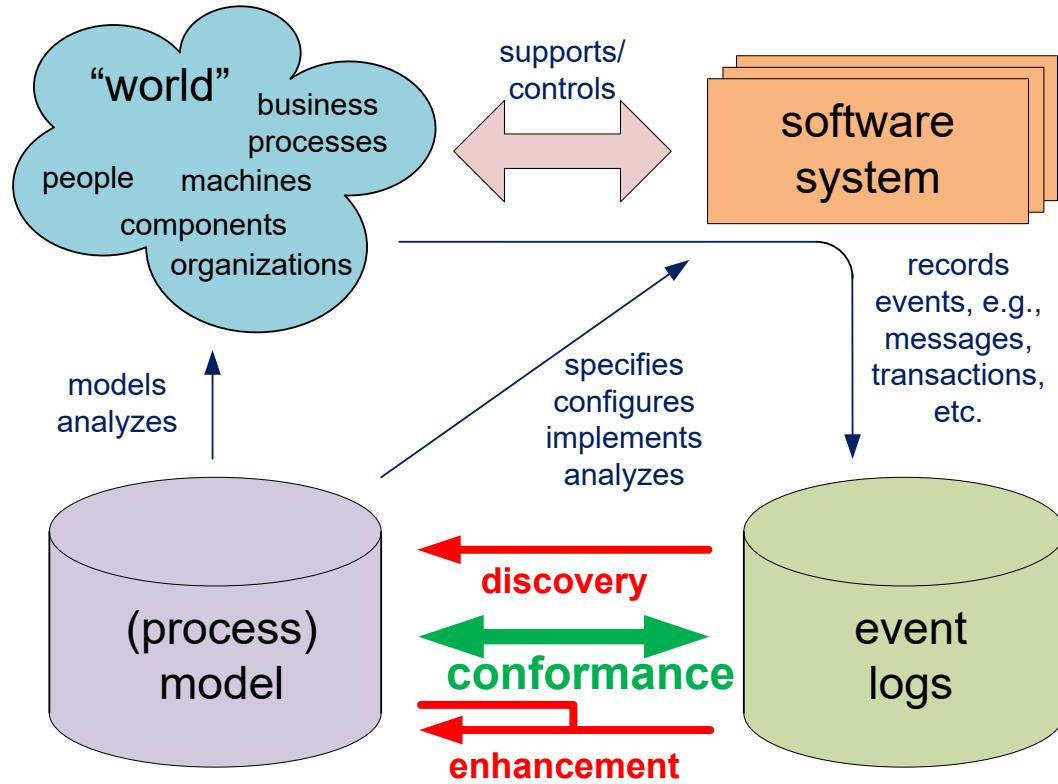
1. Conformance checking using **causal footprints**.
2. Conformance checking based on **token-based replay**.
3. **Alignment-based conformance checking.**



Conformance Checking



Conformance checking



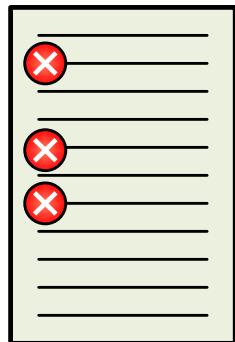
Three main use cases

- **Compliance checking (for auditing, fraud detection, etc.)**
 - Audits are performed to ascertain the validity and reliability of information about organizations and their associated processes.
 - This is done to check whether business processes are executed within certain boundaries set by managers, governments, and other stakeholders.
- **Evaluating process discovery results / algorithms**
 - Comparing discovered process models with the data used to learn the model or with unseen test data.
 - Evaluating a model or an algorithm (see k-fold cross validation).
- **Conformance to specification (software, services, etc.).**



Positive or negative deviants?

“Breaking the glass” may save lives!



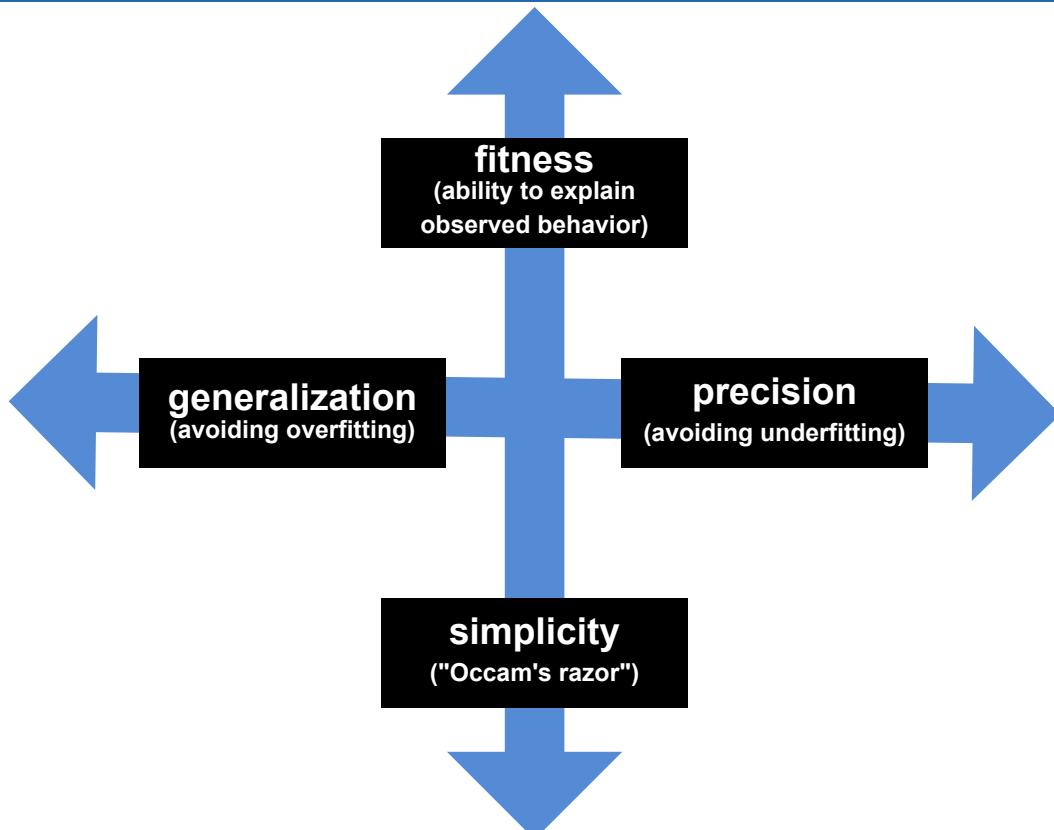
event log

Is the model wrong or is t

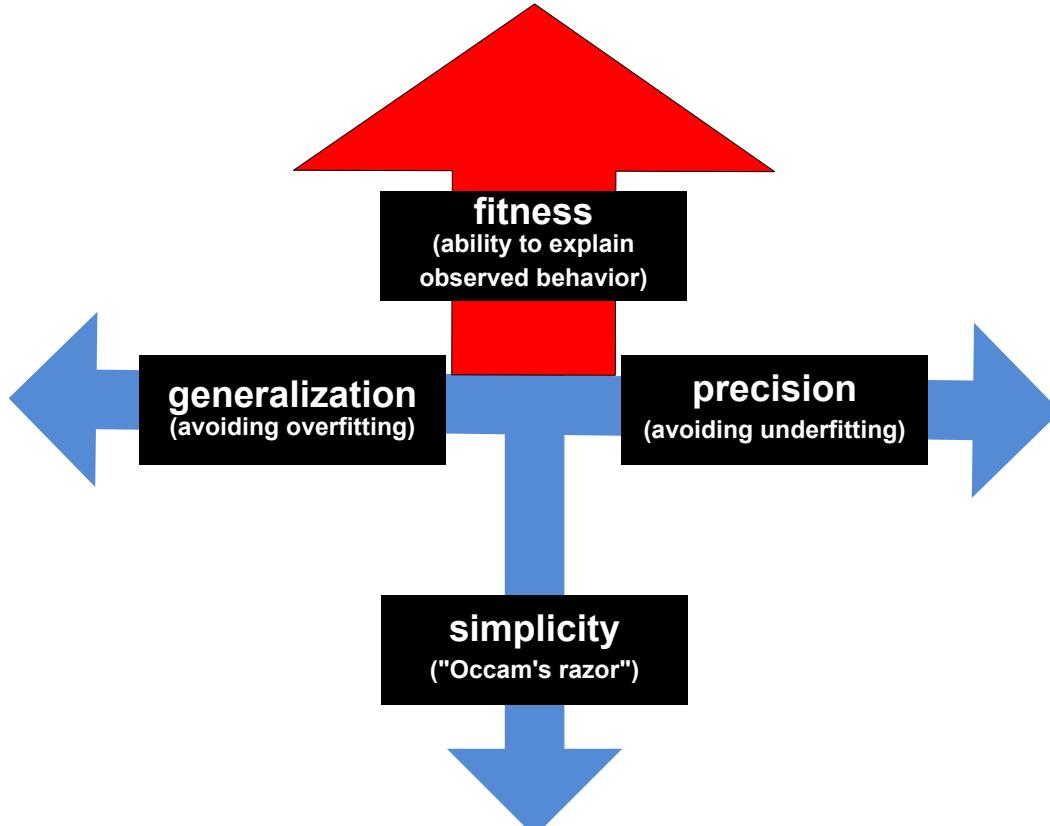


learning from
positive deviants

Four dimensions to compare log and model



Replay fitness is dominant



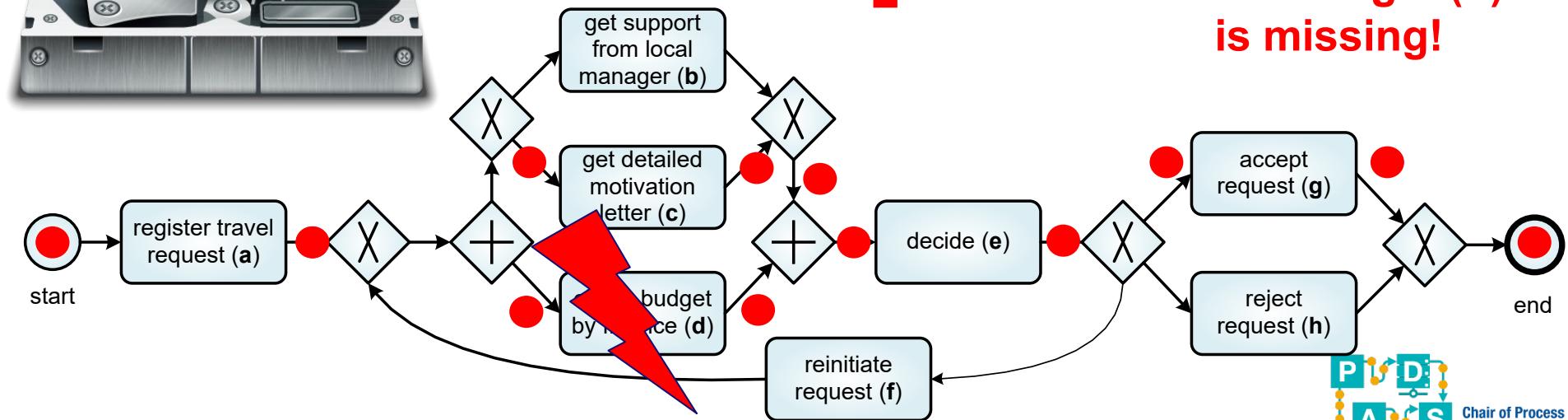
Replay example used before



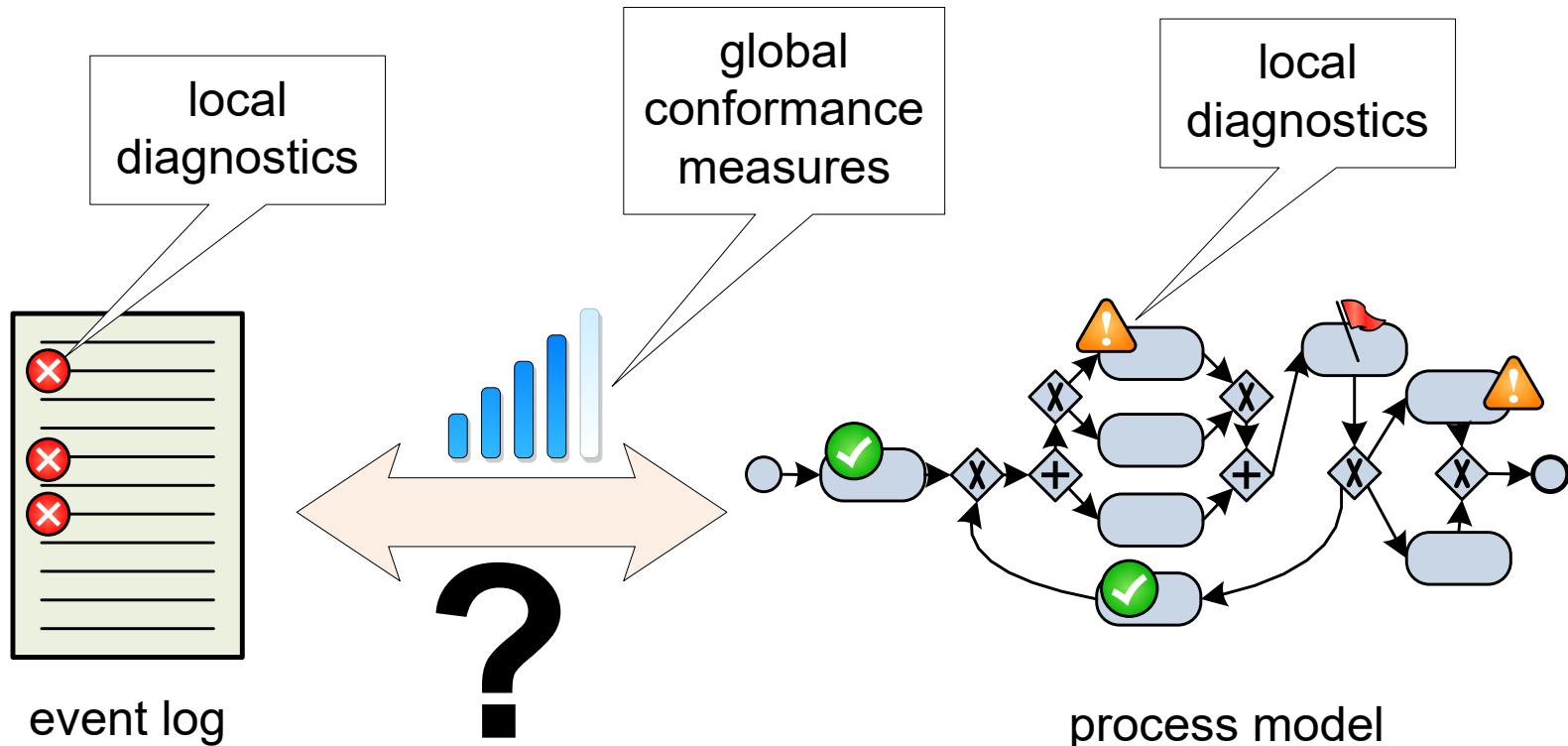
a c e g
?



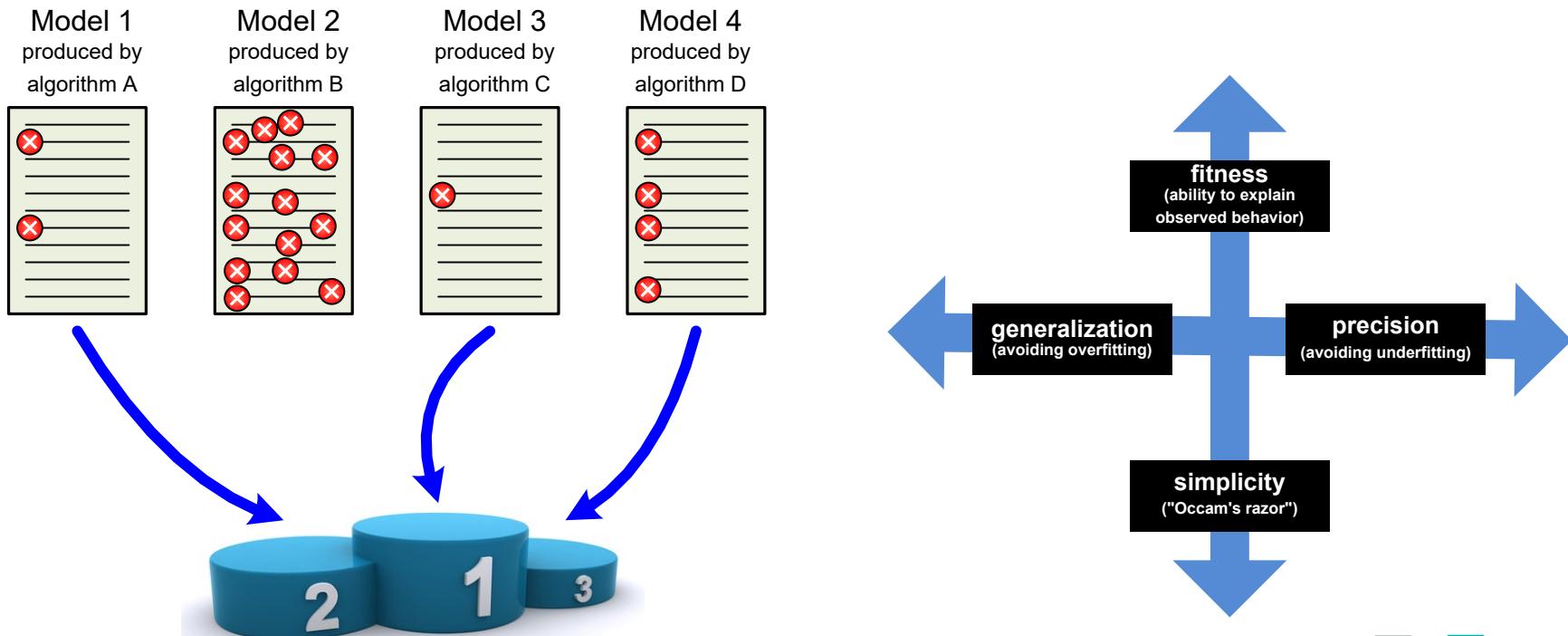
check budget (d)
is missing!



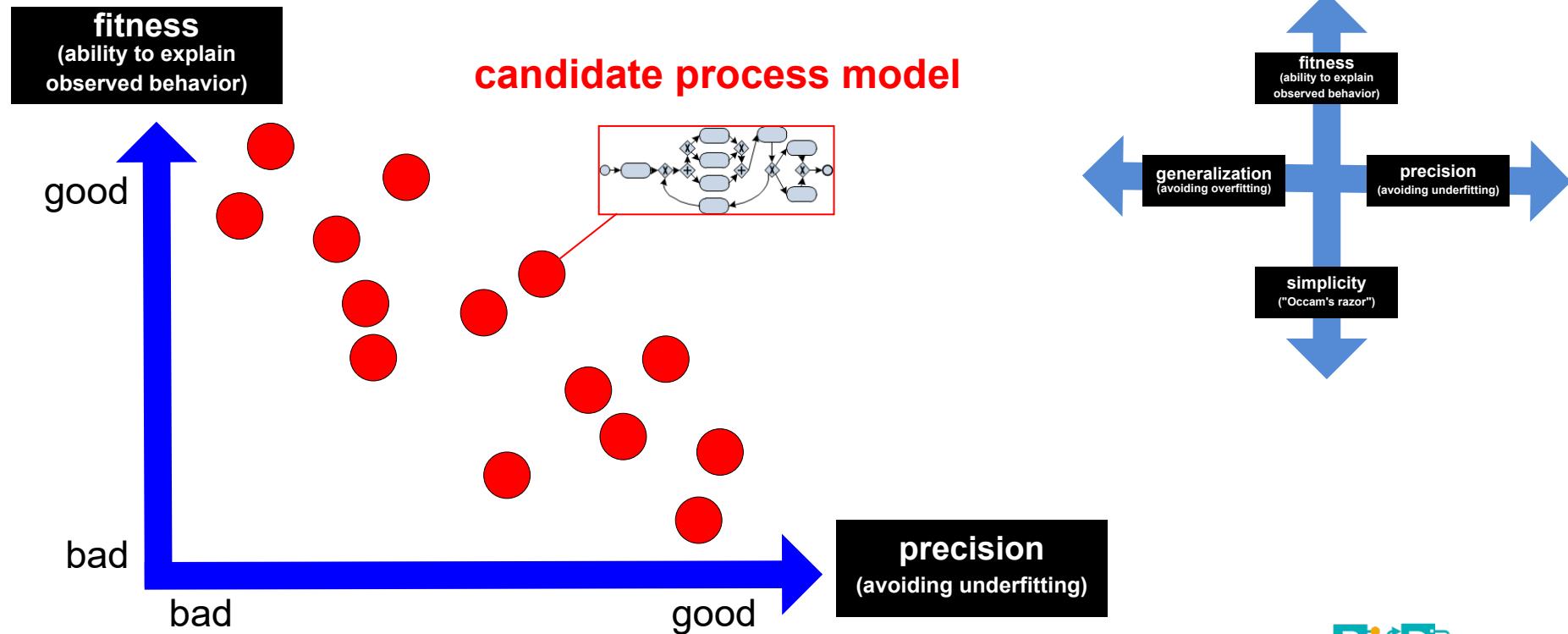
Conformance diagnostics and measures



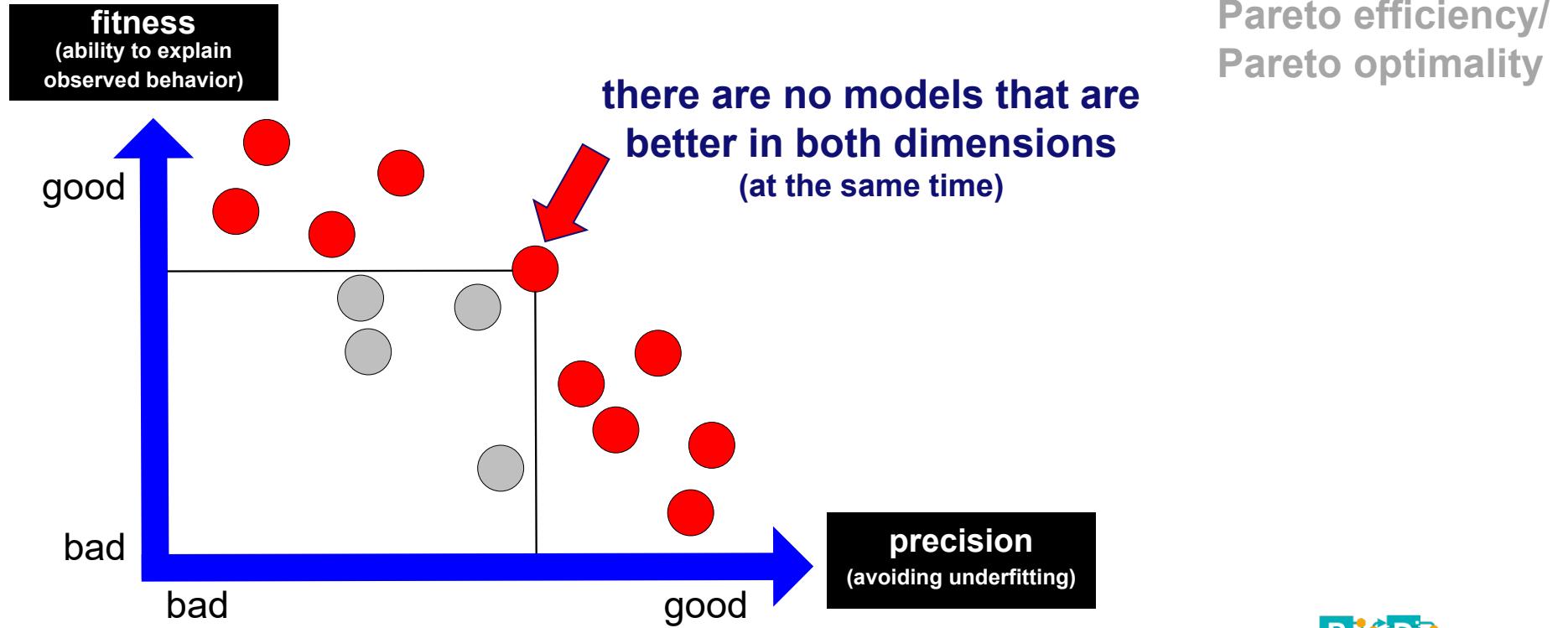
Evaluating process discovery algorithms



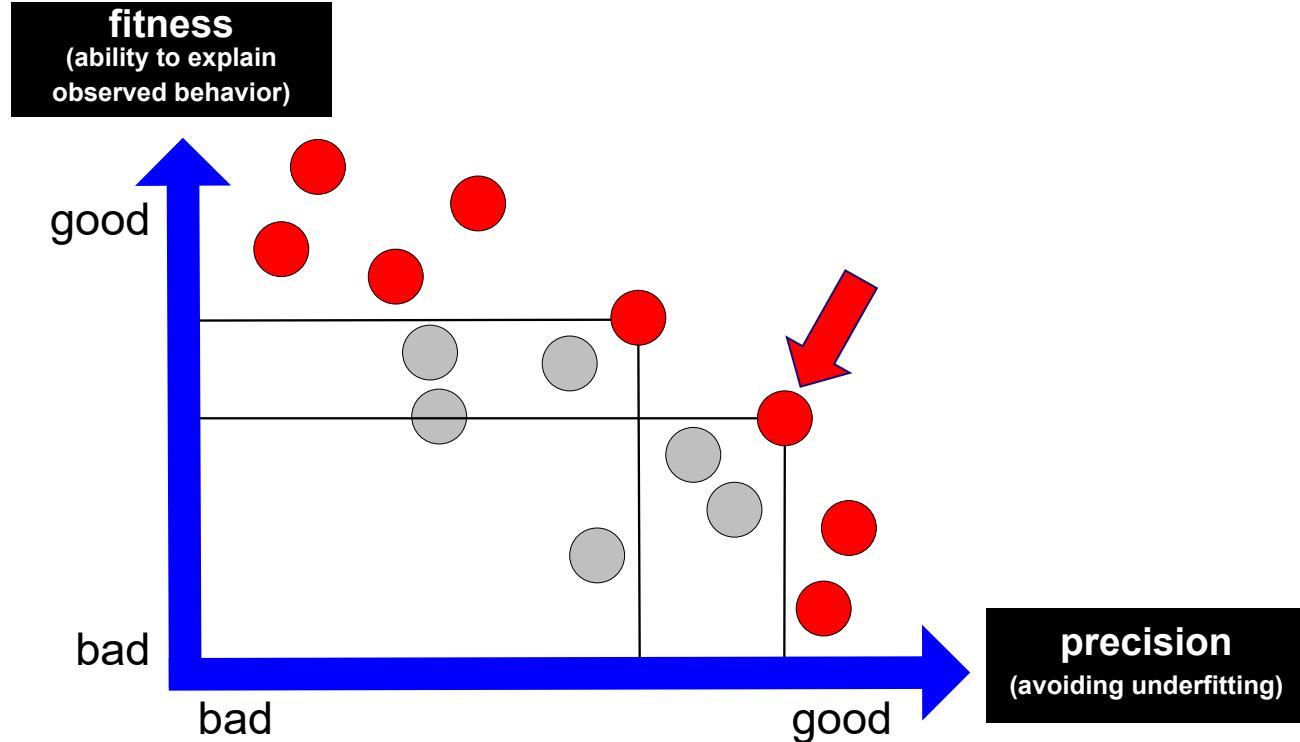
Pareto front considering two dimensions



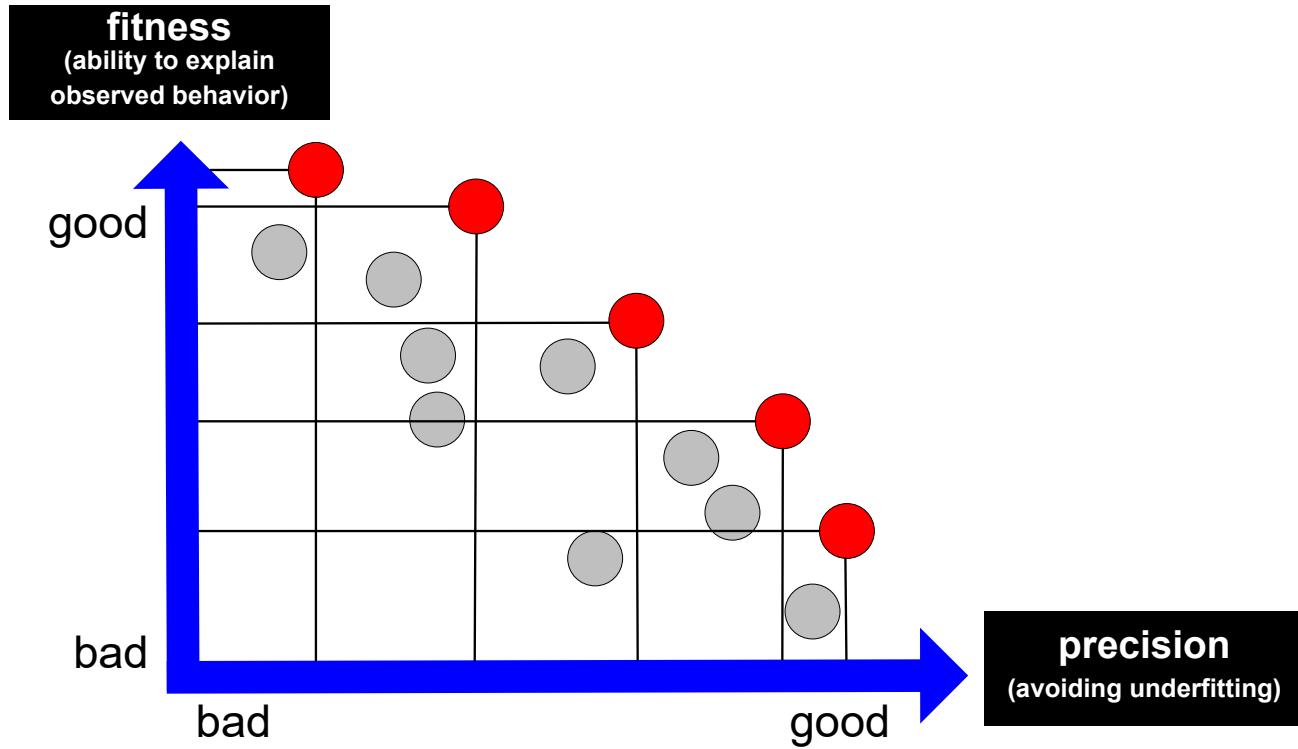
Pareto front considering two dimensions



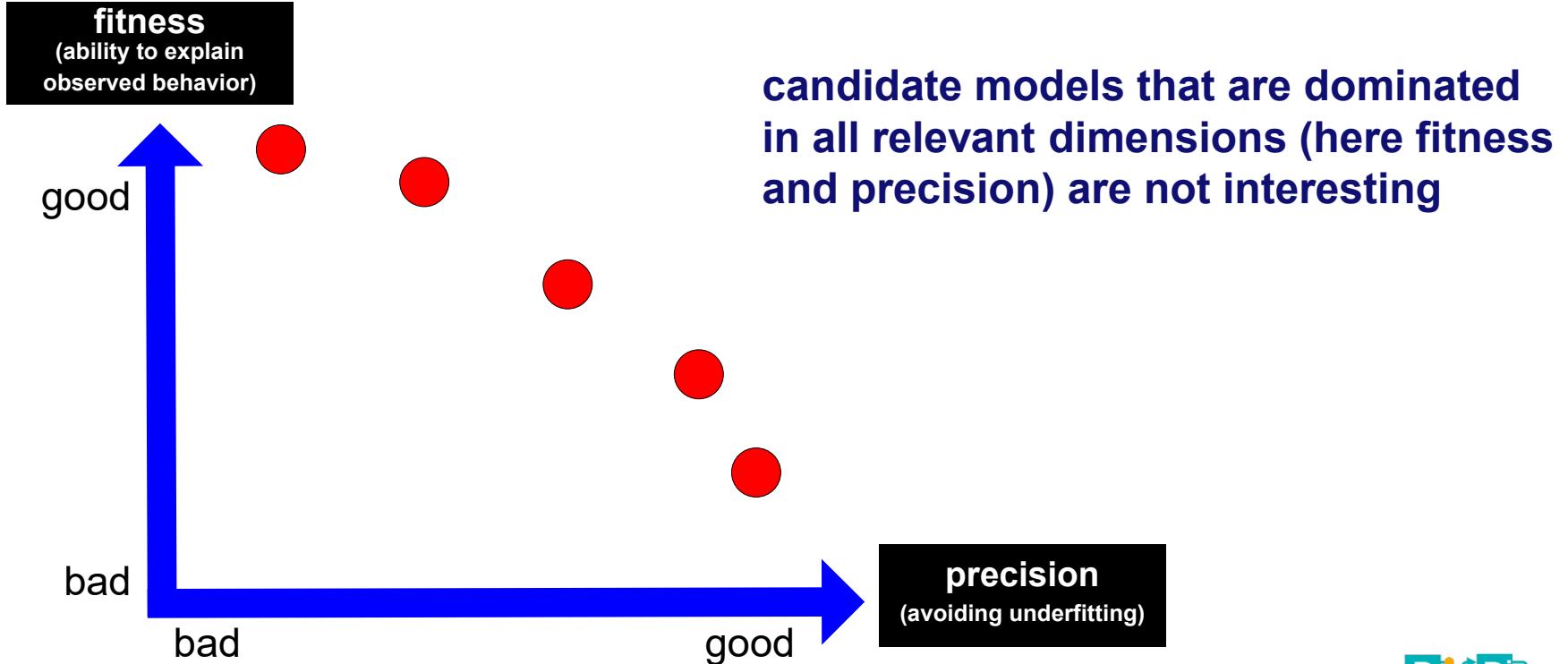
Pareto front considering two dimensions



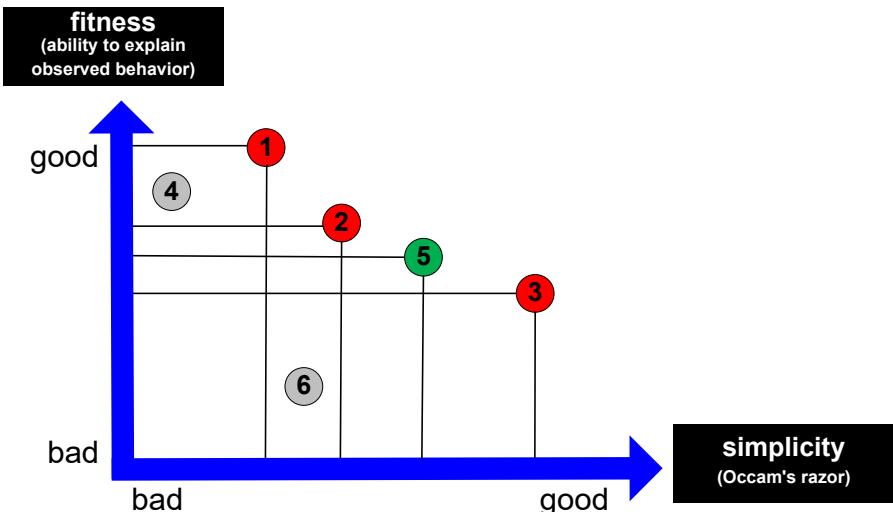
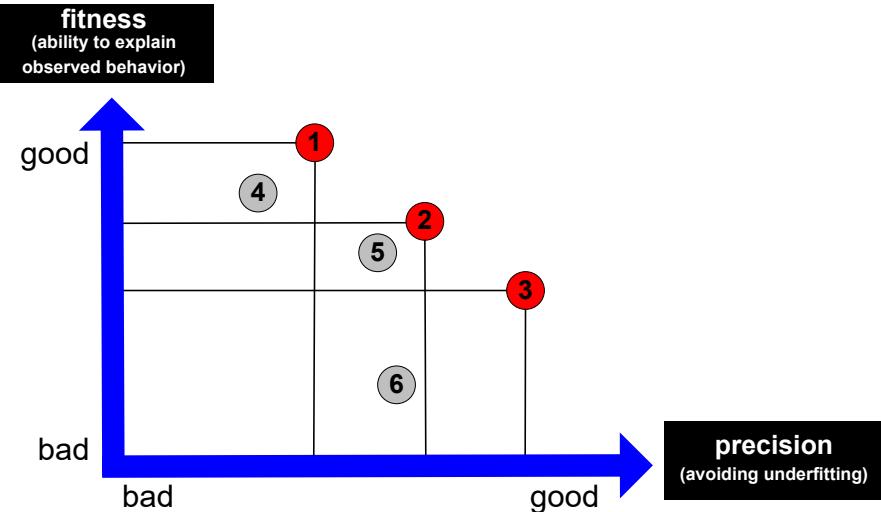
Pareto front considering two dimensions



Pareto front considering two dimensions

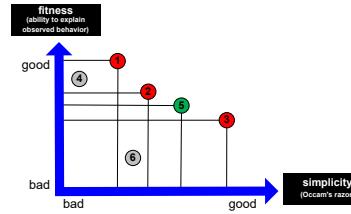
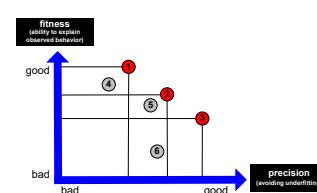
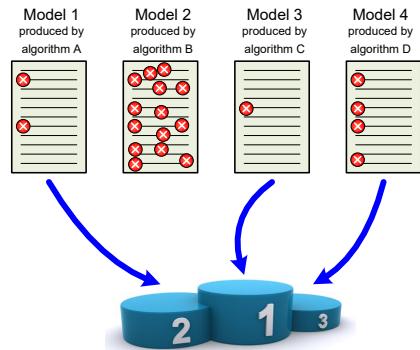


Adding another dimension: Simplicity



There is no model that (at the same time) has a better fitness, precision, and simplicity than model 5, i.e., it is not dominated.

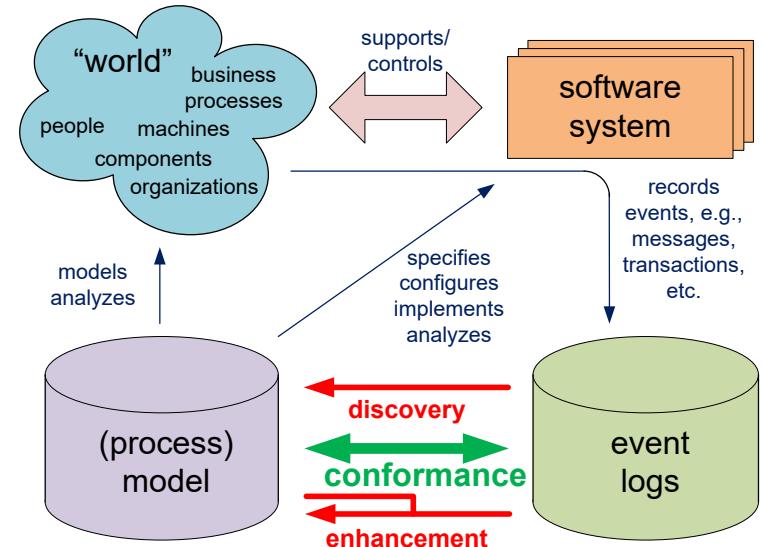
Comparing discovered models is not easy



- There is no such thing as "the best model".
 - Moreover, models may serve different purposes (see maps).
 - Initially, we focus on quantifying (replay) fitness (a.k.a. recall).

Conformance checking approaches

1. Conformance checking using **causal footprints**.
2. Conformance checking based on **token-based replay**.
3. **Alignment-based conformance checking.**



Focusing just on control-flow

- An event log is a finite multiset of traces:
$$L = [\langle a, b, c, d \rangle^6, \langle a, b, c, d \rangle^4]$$
- A model describes a (possibly infinite) set of traces:
$$M = \{ \langle a, b, c, d \rangle, \langle a, b, c, d \rangle \}$$
- As indicated before: Simple precision and recall measures do not work:
 - Loops lead to a recall of 0
 - Traces may be almost fitting
 - Log is only a sample and typically very incomplete
- Idea: Create a common finite abstraction, e.g., a DFG or causal footprint (next), dealing with incompleteness and loops.



Causal footprints



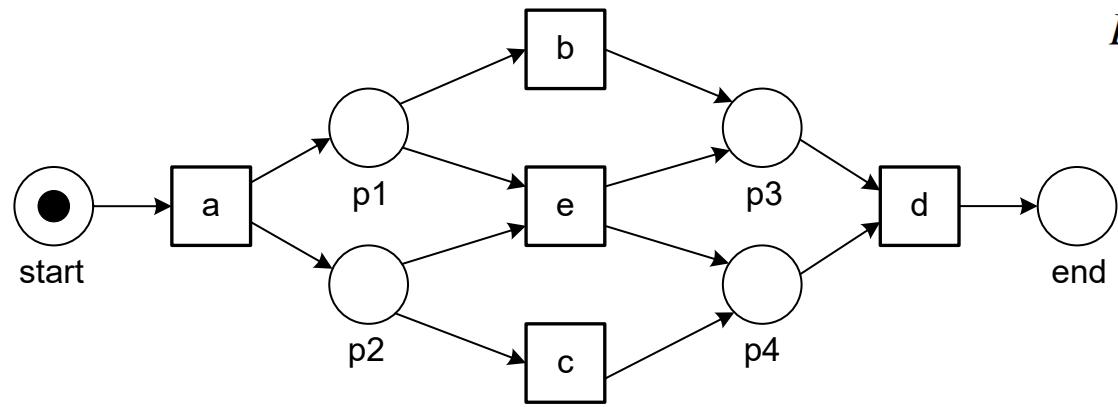
Footprint of L_1

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

	a	b	c	d	e
a	$\#_{L_1}$	\rightarrow_{L_1}	\rightarrow_{L_1}	$\#_{L_1}$	\rightarrow_{L_1}
b	\leftarrow_{L_1}	$\#_{L_1}$	\parallel_{L_1}	\rightarrow_{L_1}	$\#_{L_1}$
c	\leftarrow_{L_1}	\parallel_{L_1}	$\#_{L_1}$	\rightarrow_{L_1}	$\#_{L_1}$
d		\leftarrow_{L_1}	\leftarrow_{L_1}	$\#_{L_1}$	\leftarrow_{L_1}
e		$\#_{L_1}$	$\#_{L_1}$	\rightarrow_{L_1}	$\#_{L_1}$

- Direct succession: **x>y iff for some case x is directly followed by y.**
- Causality: **x→y iff x>y and not y>x.**
- Parallel: **x||y iff x>y and y>x**
- Choice: **x#y iff not x>y and not y>x.**

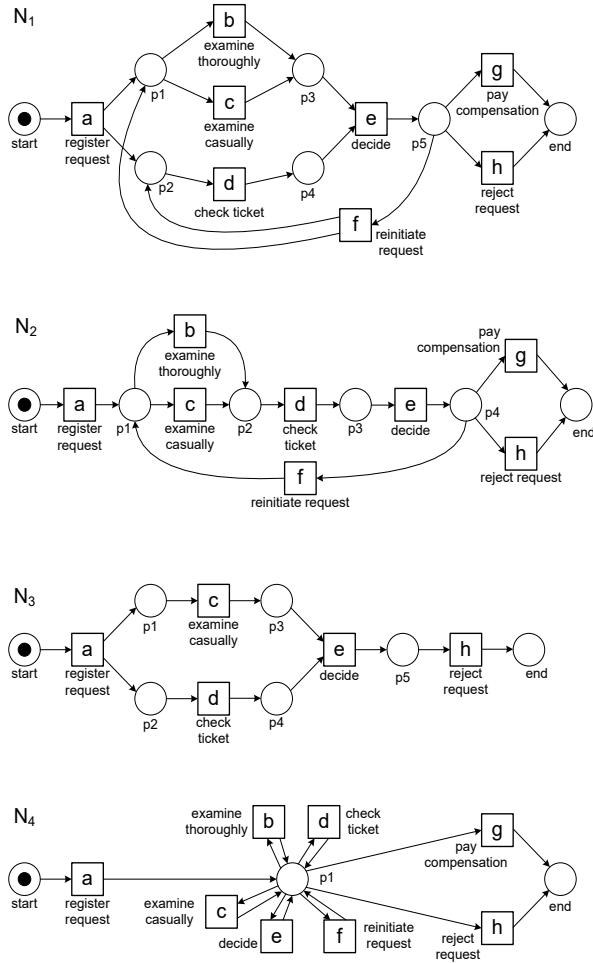
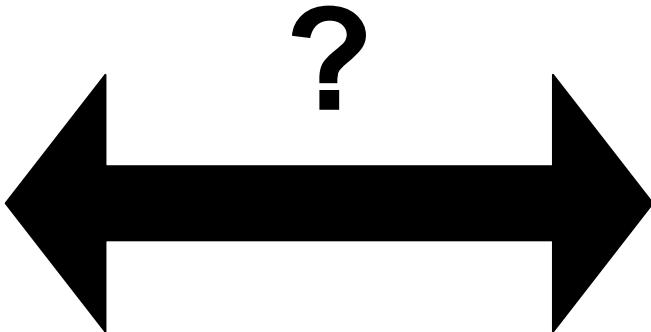
Discovered model has the same footprint



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

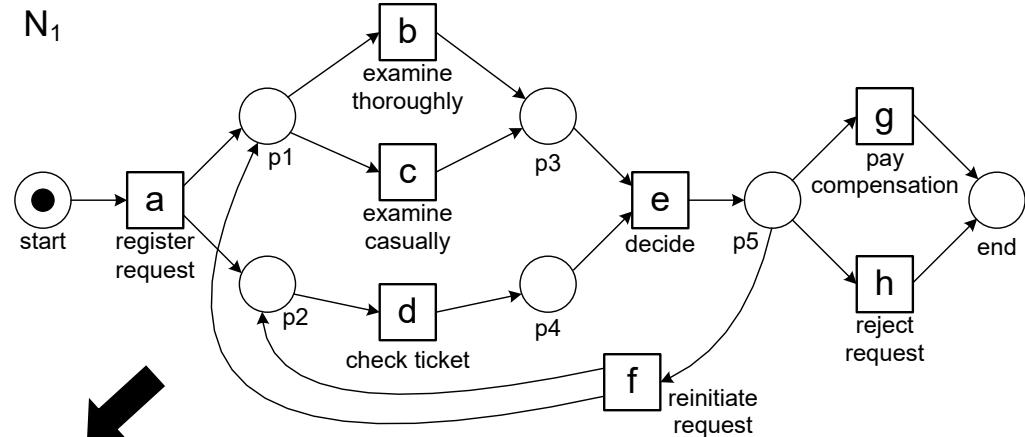
	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>
<i>a</i>	$\#_{L_1}$	\rightarrow_{L_1}	\rightarrow_{L_1}	$\#_{L_1}$	\rightarrow_{L_1}
<i>b</i>	\leftarrow_{L_1}	$\#_{L_1}$	\parallel_{L_1}	\rightarrow_{L_1}	$\#_{L_1}$
<i>c</i>	\leftarrow_{L_1}	\parallel_{L_1}	$\#_{L_1}$	\rightarrow_{L_1}	$\#_{L_1}$
<i>d</i>	$\#_{L_1}$	\leftarrow_{L_1}	\leftarrow_{L_1}	$\#_{L_1}$	\leftarrow_{L_1}
<i>e</i>	\leftarrow_{L_1}	$\#_{L_1}$	$\#_{L_1}$	\rightarrow_{L_1}	$\#_{L_1}$

#	trace
455	acdeh
191	abdeg
177	adceh
144	abdeh
111	acdeg
82	adceg
56	adbeh
47	acdefdfbeh
38	adbeg
33	acdefbdeh
14	acdefbddeg
11	acdefdbeg
9	adcefcdbeh
8	adcefdbeh
5	adcefbdeg
3	acdefbdefdbeg
2	adcefdbeg
2	adcefbddefbddeg
1	adcefdbefbdbeh
1	adbefbddefdbeg
1	adcefdbefcdefdbeg
1391	



#	trace
455	acdeh
191	abdeg
177	adceh
144	abdeh
111	acdeg
82	adceg
56	adbhe
47	acdefdbeh
38	adbeg
33	acdefbdeh
14	acdefbdbeh
11	acdefbdieg
9	adcefdeh
8	adcefdeh
5	adcfbdieg
3	acdefbdiegfbeg
2	adcfbdieg
2	adcfbdiegfbeg
1	adcfbdiegfbdeh
1	adbfbdiegfbeg
1	adcfdbefcdefdbeg
1	adcfdbefcdefdbeg

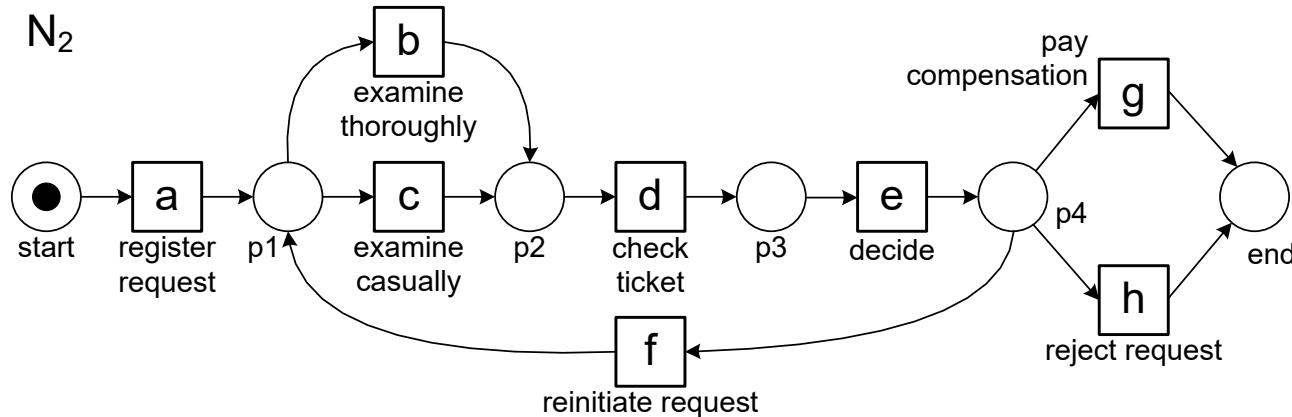
L_{full} and N_1



	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>	<i>g</i>	<i>h</i>
<i>a</i>	#	→						#
<i>b</i>	←	#						#
<i>c</i>	←	#						#
<i>d</i>	←				"	"	"	#
<i>e</i>	#	←	←	←	#	→	→	→
<i>f</i>	#	→	→	→	→	#	#	#
<i>g</i>								#
<i>h</i>	#	#	#	#	←	#	#	#

footprint-based conformance = 1
 (1 = perfect match)
 (0 = worst match possible)

footprints of log and model coincide



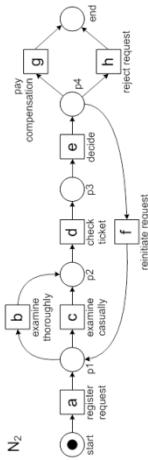
	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>	<i>g</i>	<i>h</i>
<i>a</i>	#	→	→	#	#	#	#	#
<i>b</i>	←	#	#	→	#	←	#	#
<i>c</i>	←	#	#	→	#	←	#	#
<i>d</i>	#	←	←	#	→	#	#	#
<i>e</i>	#	#	#	←	#	→	→	→
<i>f</i>	#	→	→	#	←	#	#	#
<i>g</i>	#	#	#	#	←	#	#	#
<i>h</i>	#	#	#	#	←	#	#	#

L_{full}

#	trace
455	acdeh
191	abdeg
177	adceh
144	abdeh
111	acdeg
82	adceg
56	adbeh
47	acdefdbeh
38	adbeg
33	acdefbdeh
14	acdefbdg
11	acdefdbeg
9	adcefdeh
8	adcefbeh
5	adcefbdg
3	adcefbdedefbeg
2	adcefdbeg
2	adcefbddefbdg
1	adcefdbefbdbeh
1	adbefdbefdfbeg
1	adcefbcfdefdbeg
1391	

	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>	<i>g</i>	<i>h</i>
<i>a</i>	#	\rightarrow	\rightarrow	\rightarrow	#	#	#	#
<i>b</i>	\uparrow	#	#	\uparrow	\uparrow	\uparrow	#	#
<i>c</i>	\uparrow	#	#	\uparrow	\uparrow	\uparrow	#	#
<i>d</i>	\uparrow	\uparrow	\uparrow	#	\uparrow	\uparrow	#	#
<i>e</i>	#	\uparrow	\uparrow	#	#	\uparrow	\rightarrow	\rightarrow
<i>f</i>	#	\rightarrow	\rightarrow	\rightarrow	\uparrow	\uparrow	#	#
<i>g</i>	#	#	#	#	\uparrow	\uparrow	#	#
<i>h</i>	#	#	#	#	\uparrow	$\#$	#	#

1391



N_2

	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>	<i>g</i>	<i>h</i>
<i>a</i>	#	\rightarrow	\rightarrow	#	#	#	#	#
<i>b</i>	\uparrow	#	#	\uparrow	$\#$	\uparrow	#	#
<i>c</i>	\uparrow	#	#	\uparrow	$\#$	\uparrow	#	#
<i>d</i>	#	\uparrow	\uparrow	#	\uparrow	$\#$	#	#
<i>e</i>	#	#	#	#	#	\uparrow	\rightarrow	\rightarrow
<i>f</i>	#	\rightarrow	\rightarrow	#	\uparrow	\uparrow	#	#
<i>g</i>	#	#	#	#	\uparrow	\uparrow	#	#
<i>h</i>	#	#	#	#	\uparrow	$\#$	#	#

Quantifying the differences

	a	b	c	d	e	f	g	h
a				→: #				
b				:→	→: #			
c				:→	→: #			
d	←: #	:←	:←				←: #	
e		←: #	←: #					
f				→: #				
g								
h								

(x:y where x is in log and y in N₂)

footprint-based conformance



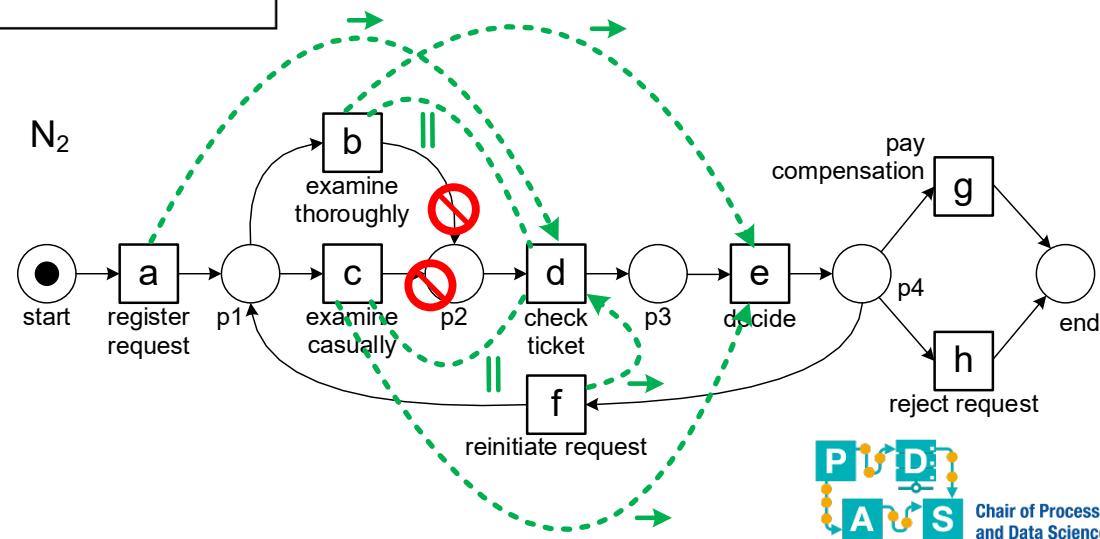
Chair of Process
and Data Science

$$1 - \frac{12}{64} = 0.8125$$

Diagnostics

	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>	<i>g</i>	<i>h</i>
<i>a</i>					$\rightarrow: \#$			
<i>b</i>					$\parallel : \rightarrow$	$\rightarrow: \#$		
<i>c</i>					$\parallel : \rightarrow$	$\rightarrow: \#$		
<i>d</i>	$\leftarrow: \#$	$\parallel : \leftarrow$	$\parallel : \leftarrow$				$\leftarrow: \#$	
<i>e</i>		$\leftarrow: \#$	$\leftarrow: \#$					
<i>f</i>							$\rightarrow: \#$	
<i>g</i>								
<i>h</i>								

(*x:y* where *x* is in log and *y* in N_2)

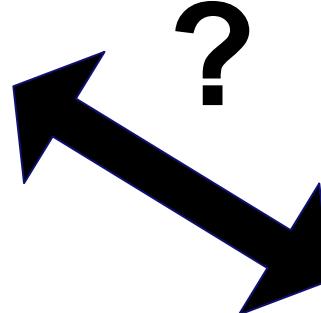


Question

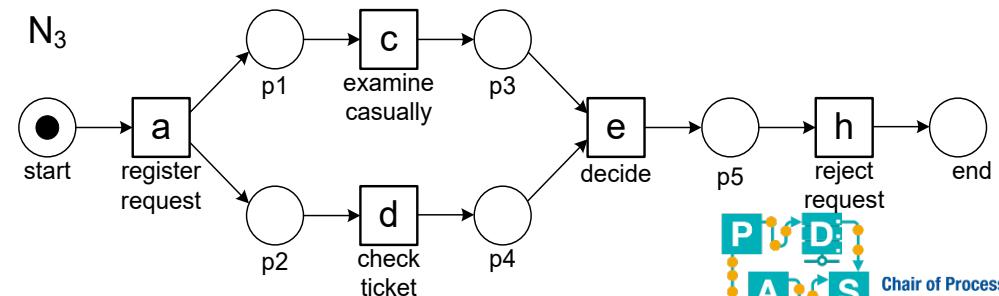
Estimate footprint-based conformance

#	trace
455	acdeh
191	abdeg
177	adceh
144	abdeh
111	acdeg
82	adceg
56	adbeh
47	acdefdbeh
38	adbeg
33	acdeffbdeh
14	acdefbddeg
11	acdefdbeg
9	adcefcdbeh
8	adcefdbbeh
5	adcefbddeg
3	acdefbfdfdbeg
2	adcefdbeg
2	adcefbddefbddeg
1	adcefdbefbdbeh
1	adbefbfdfdfbeg
1	adcefdbefcdfdfbeg
1391	

	a	b	c	d	e	f	g	h
a	#	→	→	→	#	#	#	#
b	←	#	#		→	←	#	#
c	←	#	#		→	←	#	#
d	←			#	→	←	#	#
e	#	←	←	←	#	→	→	→
f	#	→	→	→	←	#	#	#
g	#	#	#	#	←	#	#	#
h	#	#	#	#	←	#	#	#



Estimate the fraction of matching cells in footprint matrices

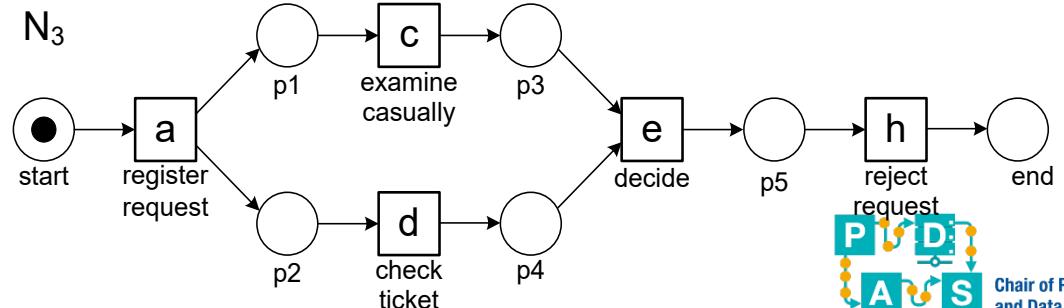


Answer

	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>	<i>g</i>	<i>h</i>
<i>a</i>	#	#	→	→	#	#	#	#
<i>b</i>	#	#	#	#	#	#	#	#
<i>c</i>	←	#	#		→	#	#	#
<i>d</i>	←	#		#	→	#	#	#
<i>e</i>	#	#	←	←	#	#	#	→
<i>f</i>	#	#	#	#	#	#	#	#
<i>g</i>	#	#	#	#	#	#	#	#
<i>h</i>	#	#	#	#	←	#	#	#

$$1 - \frac{16}{64} = 0.75$$

footprint-based conformance



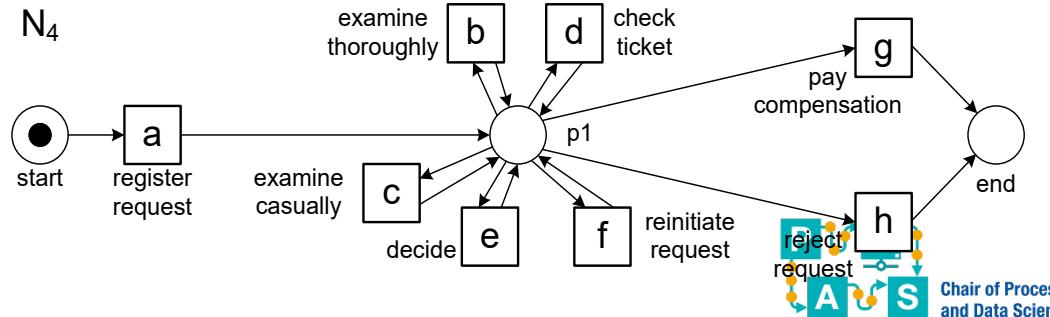
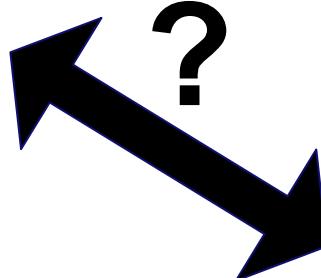
Question

Estimate footprint-based conformance

#	trace
455	acdeh
191	abdeg
177	adceh
144	abdeh
111	acdeg
82	adceg
56	adbeh
47	acdefdbeh
38	adbeg
33	acdefbdbeh
14	acdefbddeg
11	acdefdbeg
9	adcefcdbeh
8	adcefdbbeh
5	adcefbddeg
3	acdefbdfdfbeg
2	adcefdbeg
2	adcefbddefbddeg
1	adcefdbefbdbeh
1	adbefbddefdbeg
1	adcefdbeccdefdbeg
1391	

	a	b	c	d	e	f	g	h
a	#	→	→	→	#	#	#	#
b	←	#	#		→	←	#	#
c	←	#	#		→	←	#	#
d	←			#	→	←	#	#
e	#	←	←	←	#	→	→	→
f	#	→	→	→	←	#	#	#
g	#	#	#	#	←	#	#	#
h	#	#	#	#	←	#	#	#

Estimate the fraction of matching cells in footprint matrices



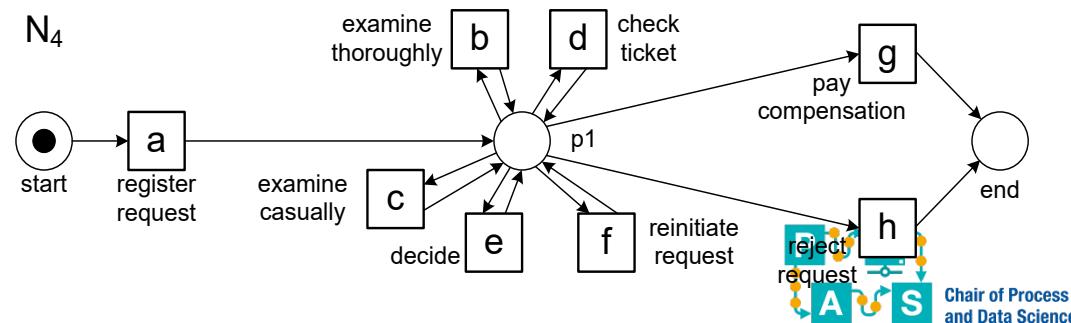
Answer

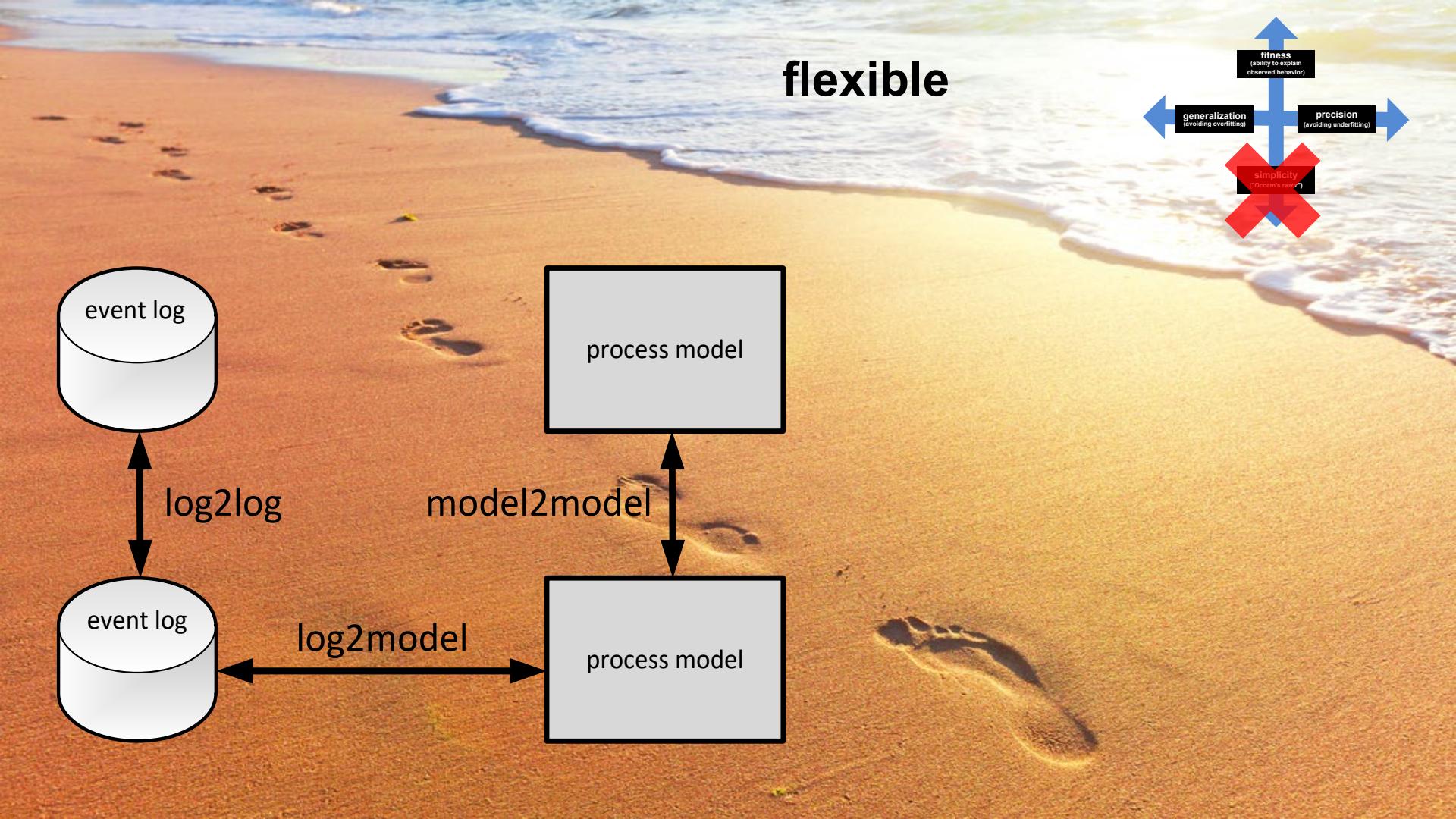
	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>	<i>g</i>	<i>h</i>
<i>a</i>	#	→	→	→	#	#	#	#
<i>b</i>	←					←	#	#
<i>c</i>	←					←	#	#
<i>d</i>	←					←	#	#
<i>e</i>	#	←	←	←			→	→
<i>f</i>	#						#	#
<i>g</i>	#	←	←	←	←	←	#	#
<i>h</i>	#	←	←	←	←	←	#	#

Color
 • Log
 • Model

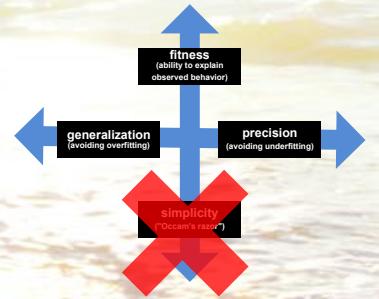
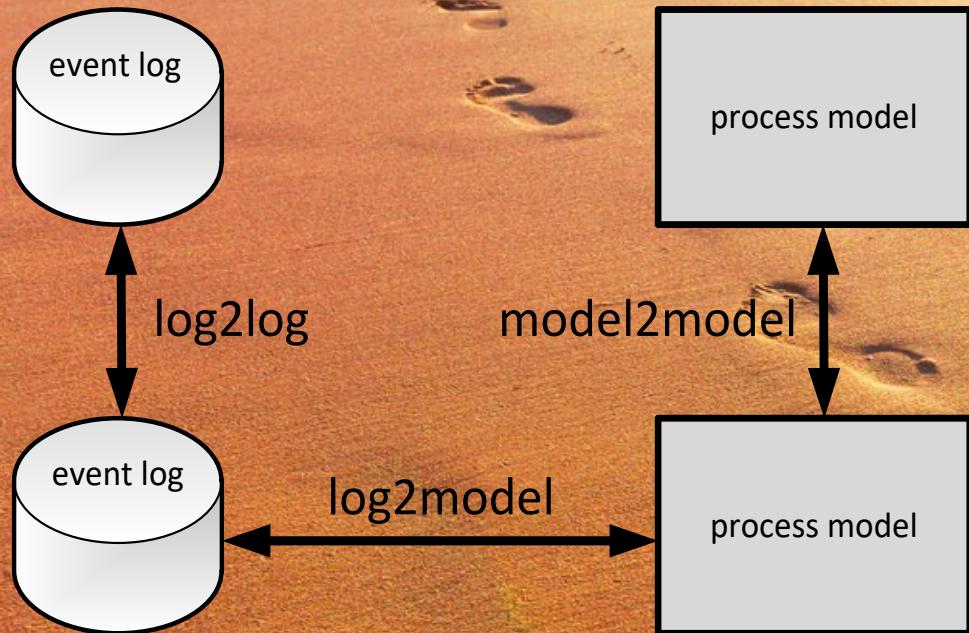
$$1 - \frac{45}{64} = 0.296875$$

footprint-based conformance



The background of the diagram is a photograph of a sandy beach meeting the ocean at the water's edge. Several sets of footprints are visible in the sand, leading towards the water.

flexible



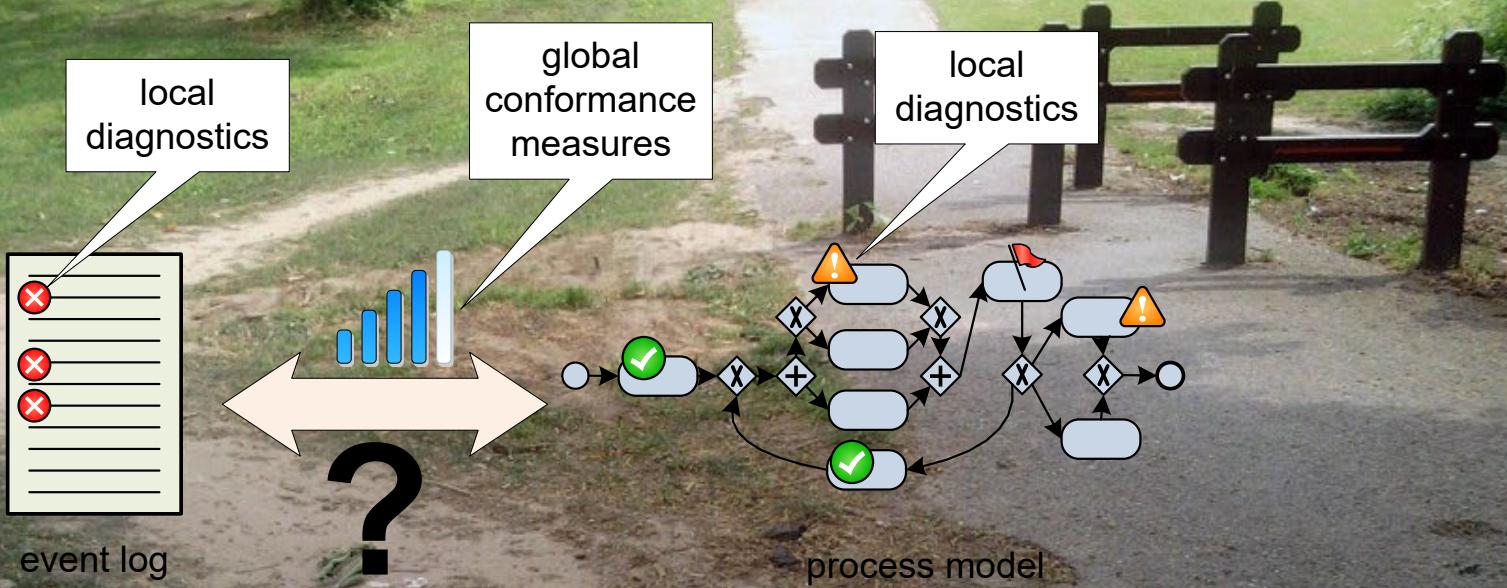
Limitations

- Frequencies are not used.
- Behavior is only considered indirectly (directly follows relation).
- Aims to capture fitness, precision and generalization in a single metric.

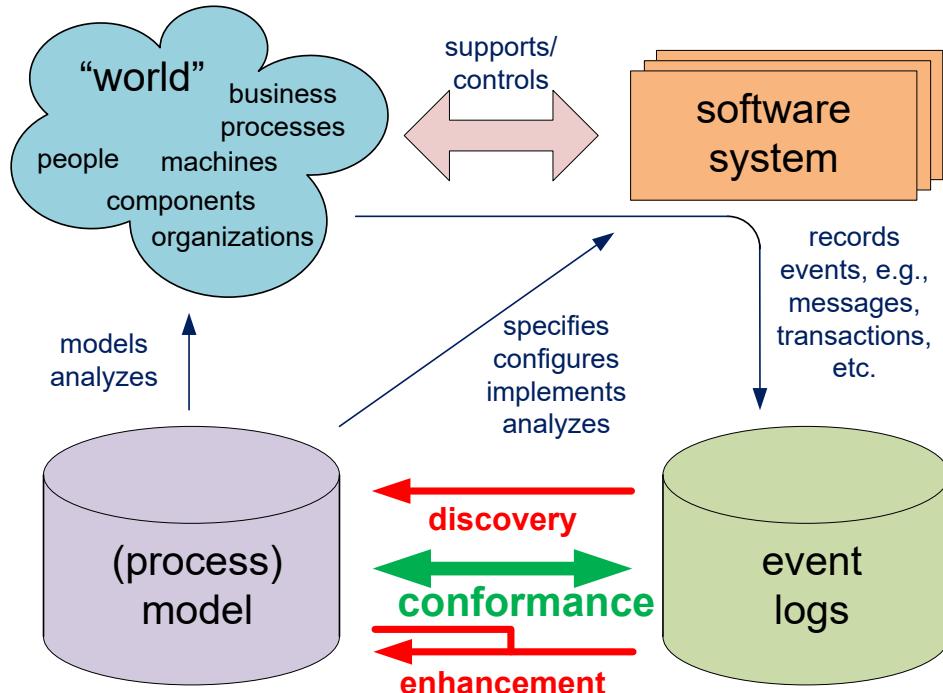
Next: conformance checking using token-based replay

Conformance checking using token-based replay



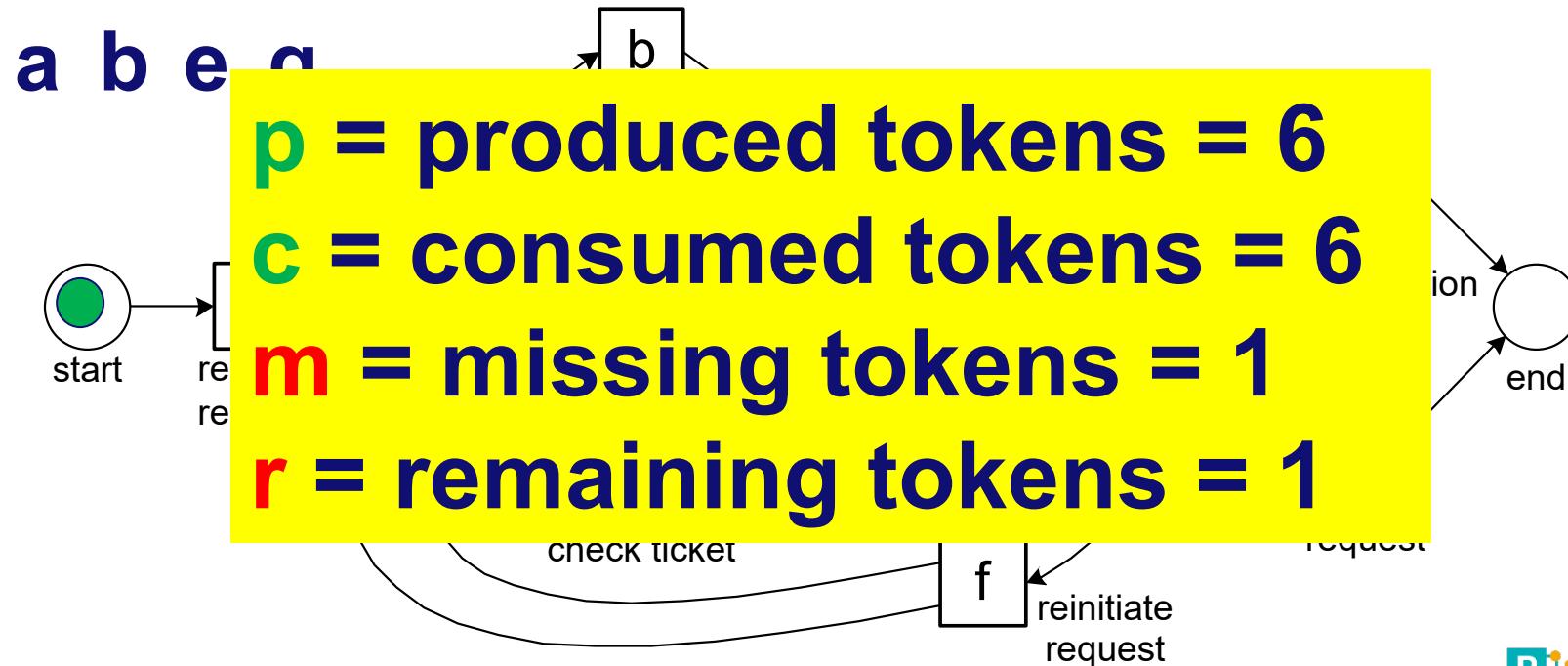


Conformance checking



1. Conformance checking using causal footprints.
2. Conformance checking based on **token-based replay**.
3. Alignment-based conformance checking.

Counting tokens while replaying



Quantifying fitness at the trace level

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

p = produced tokens = 6

c = consumed tokens = 6

m = missing tokens = 1

r = remaining tokens = 1

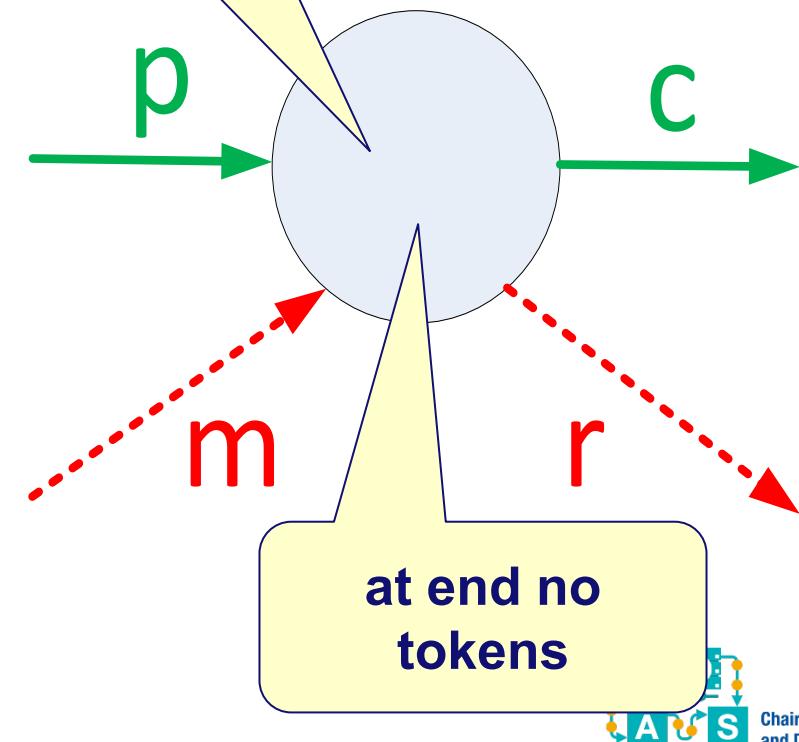


Approach (1/3)

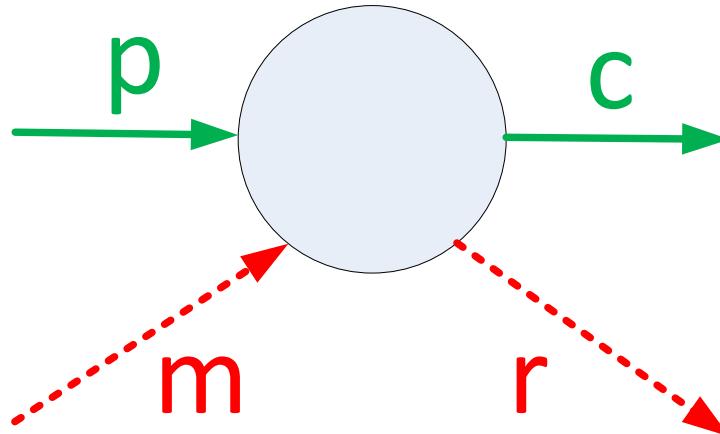
while running
 $p+m-c$ tokens

Use four counters:

- **p = produced tokens**
- **c = consumed tokens**
- **m = missing tokens**
(consumed while not there)
- **r = remaining tokens**
(produced but not consumed)

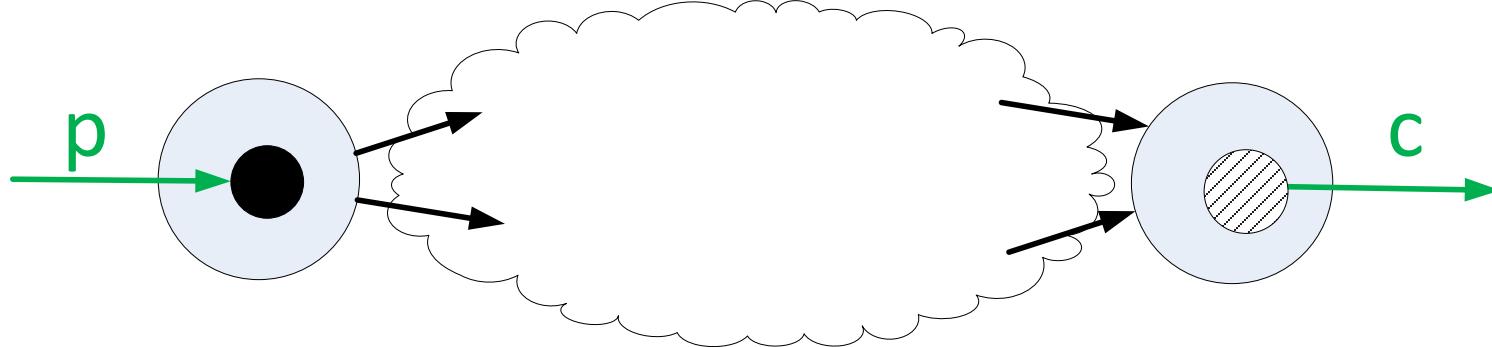


Approach (2/3)



- Invariants
 - At any time: $p+m \geq c \geq m$ (also per place)
 - At the end: $r = p + m - c$ (also per place)

Approach (3/3)

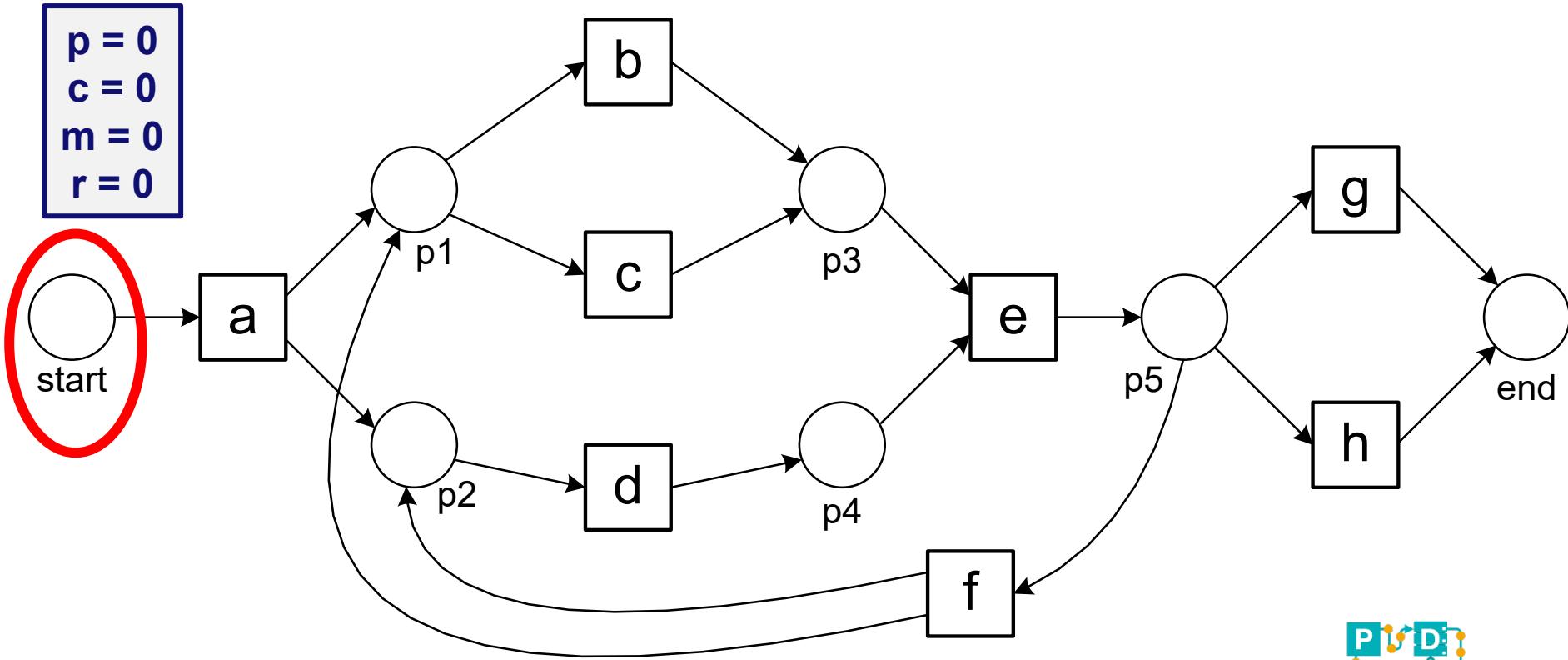


Initialization and finalization:

- In the beginning a token is **produced** for the source place: $p = 1$.
- At the end a token is **consumed** from the sink place (also if not there): $c' = c + 1$.

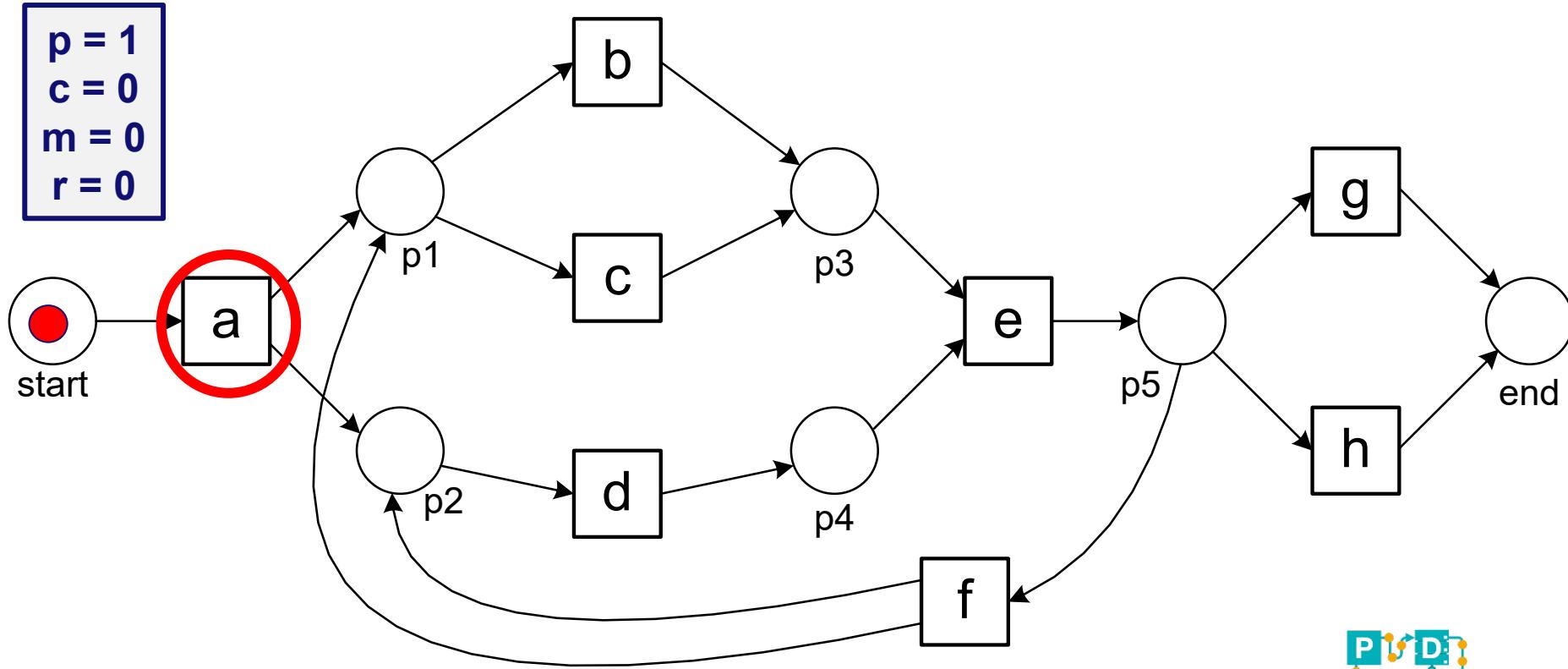
Replaying

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



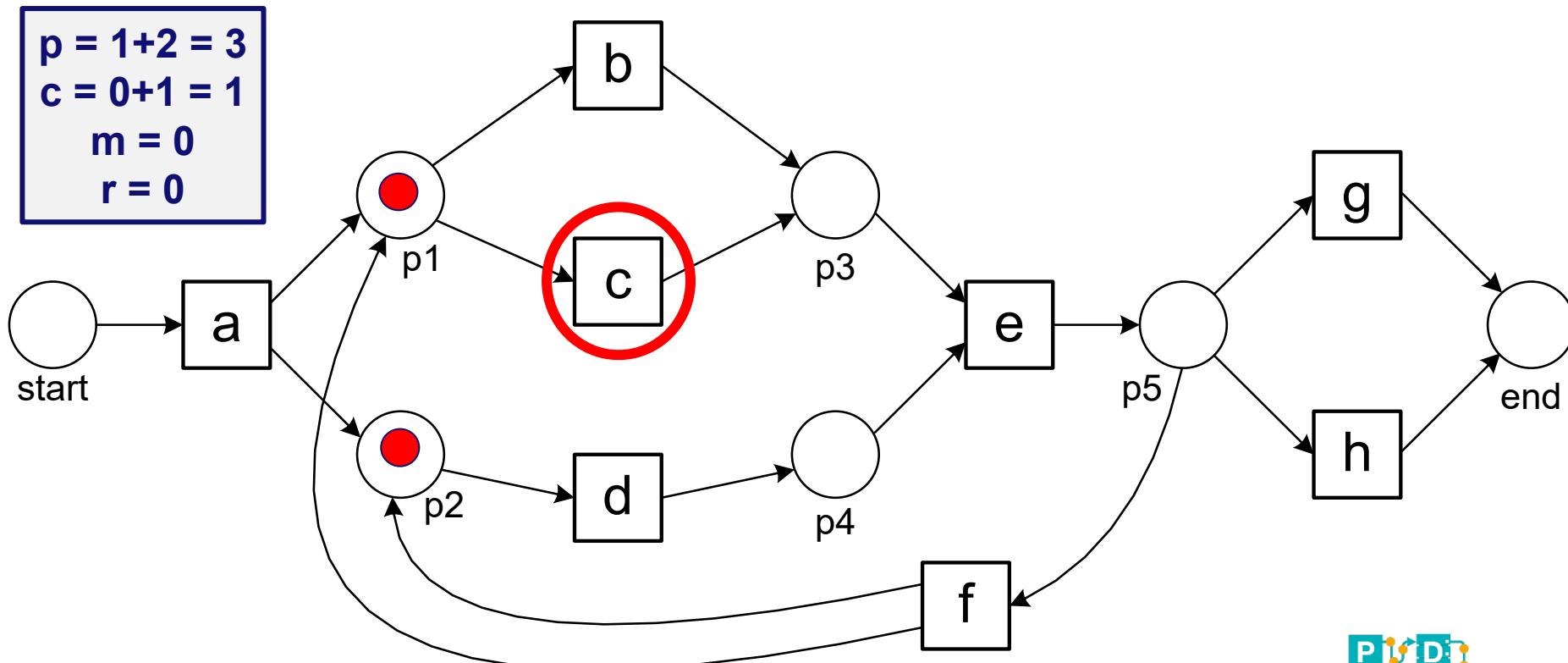
Replaying

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



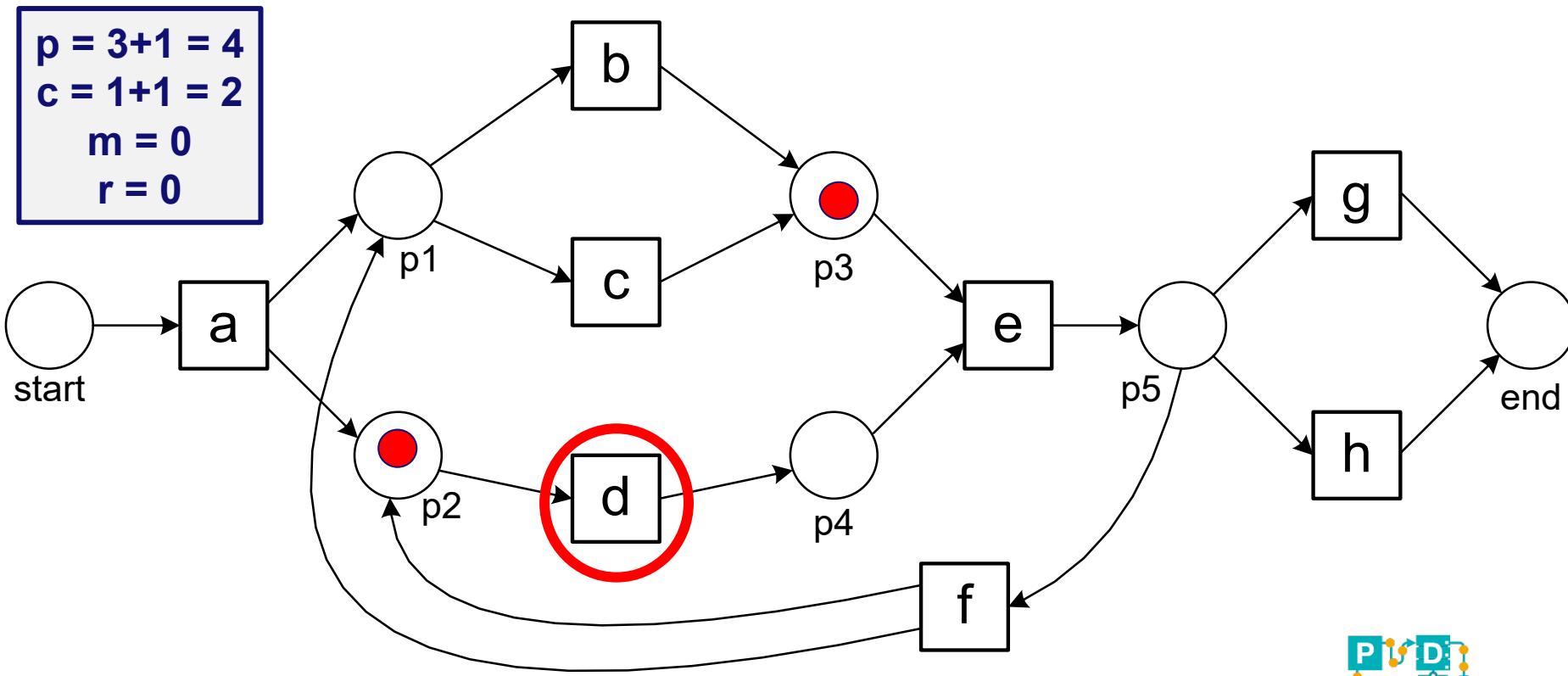
Replaying

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



Replaying

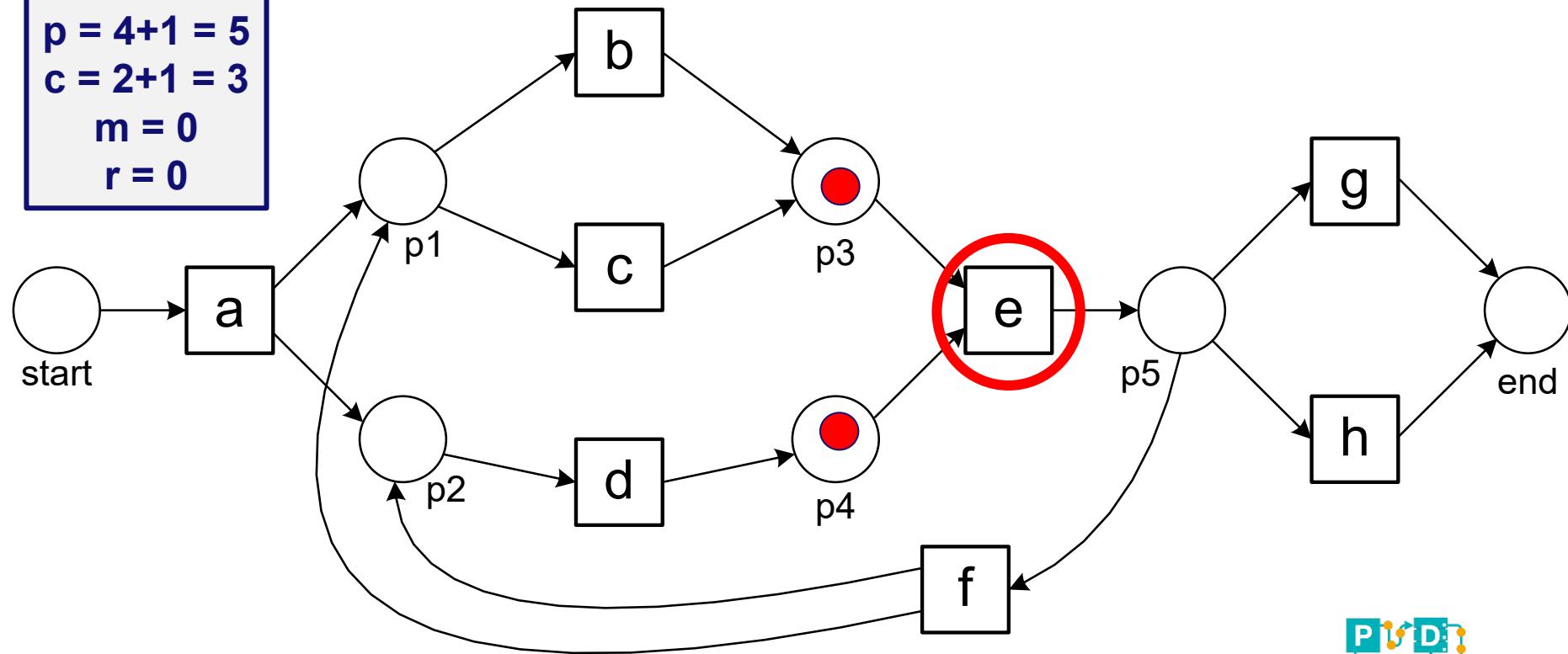
$$\sigma_1 = \langle a, c, d, e, h \rangle$$



Replaying

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

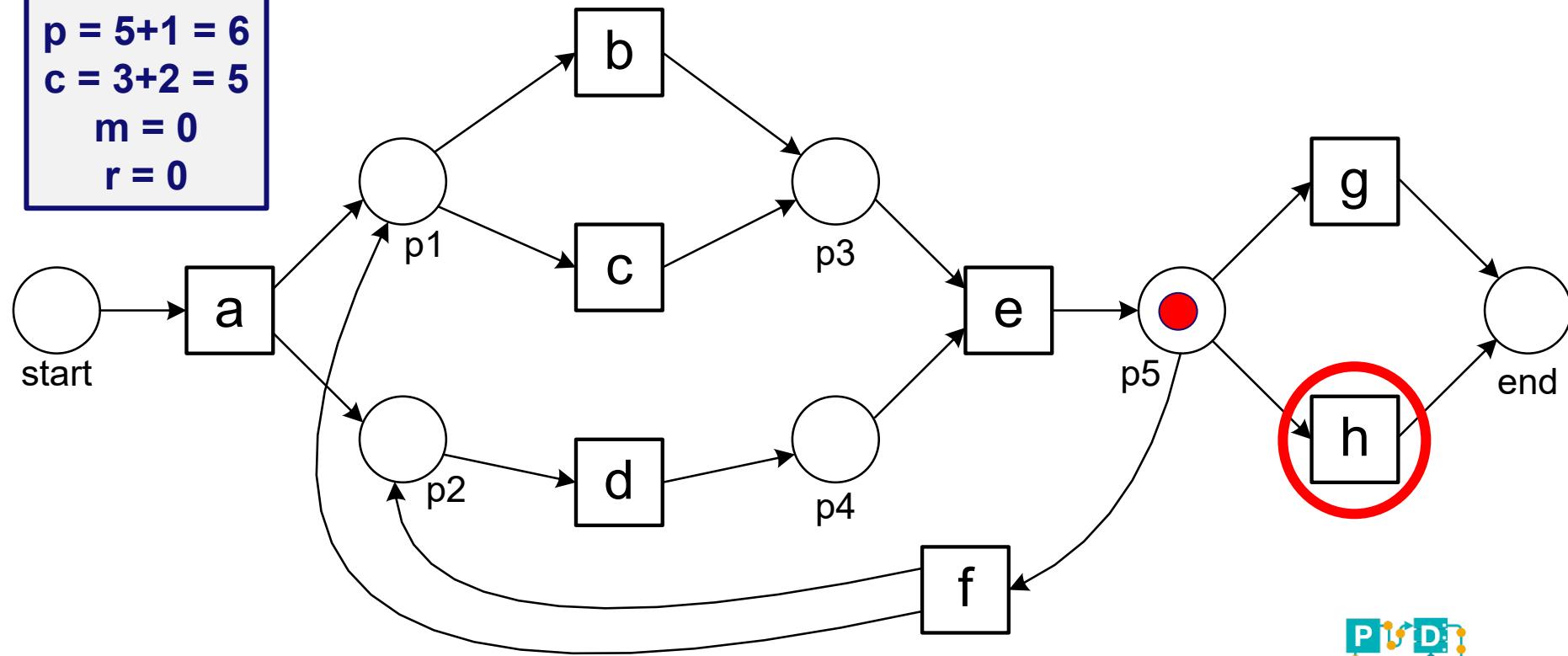
$p = 4+1 = 5$
 $c = 2+1 = 3$
 $m = 0$
 $r = 0$



Replaying

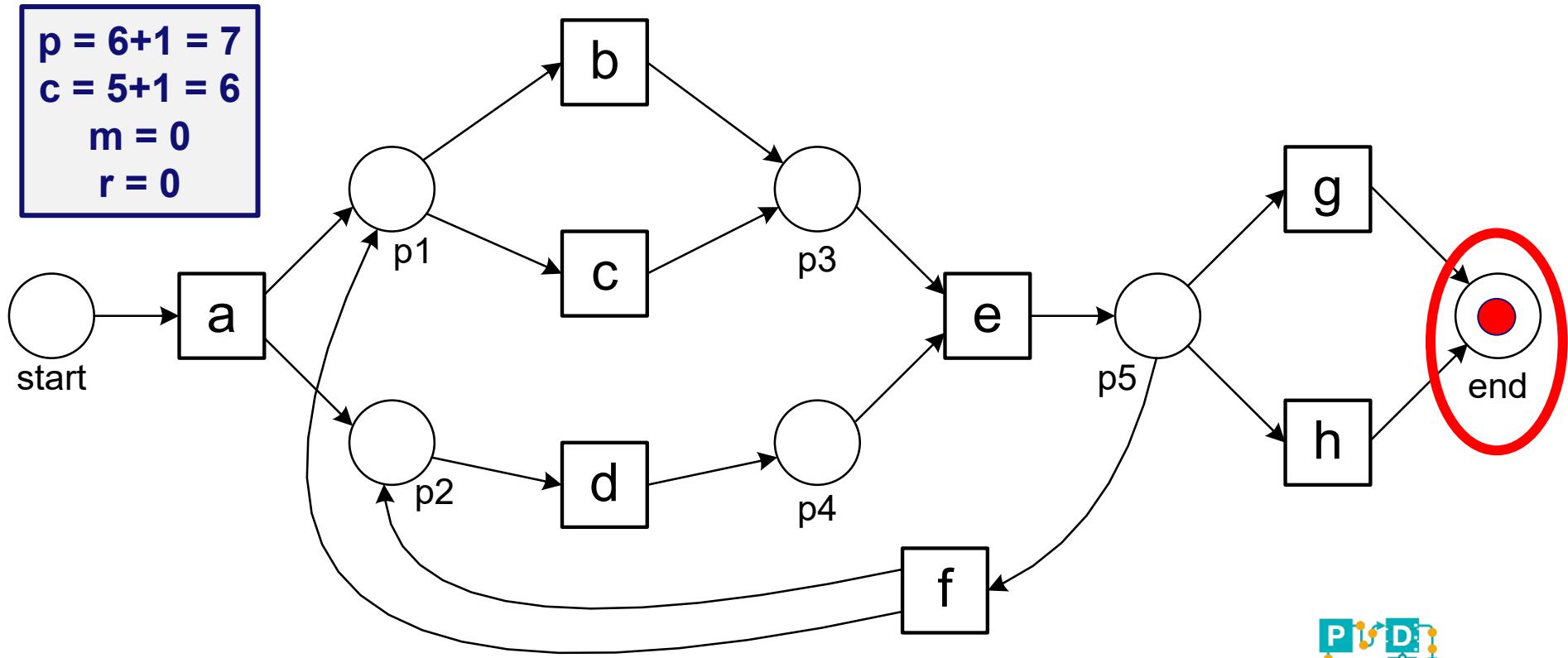
$$\sigma_1 = \langle a, c, d, e, h \rangle$$

$p = 5+1 = 6$
 $c = 3+2 = 5$
 $m = 0$
 $r = 0$



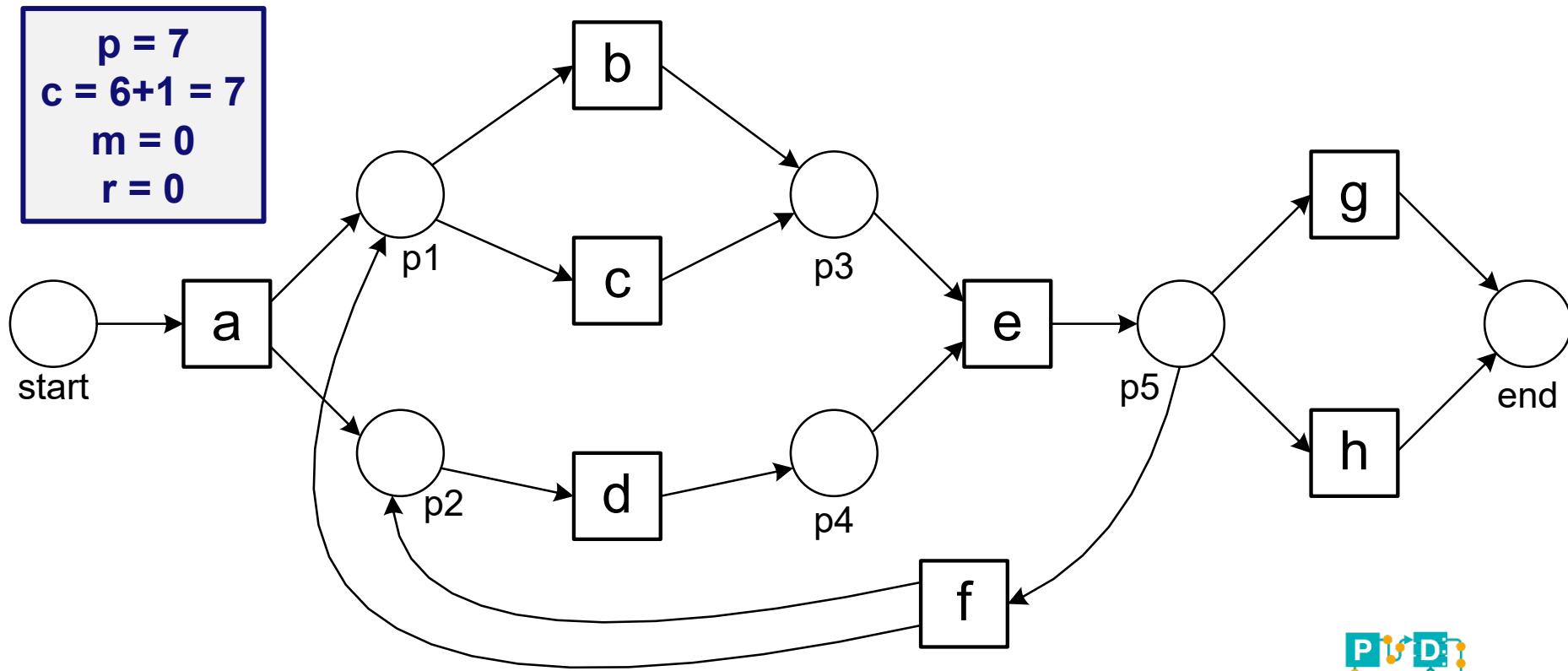
Replaying

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



Replaying

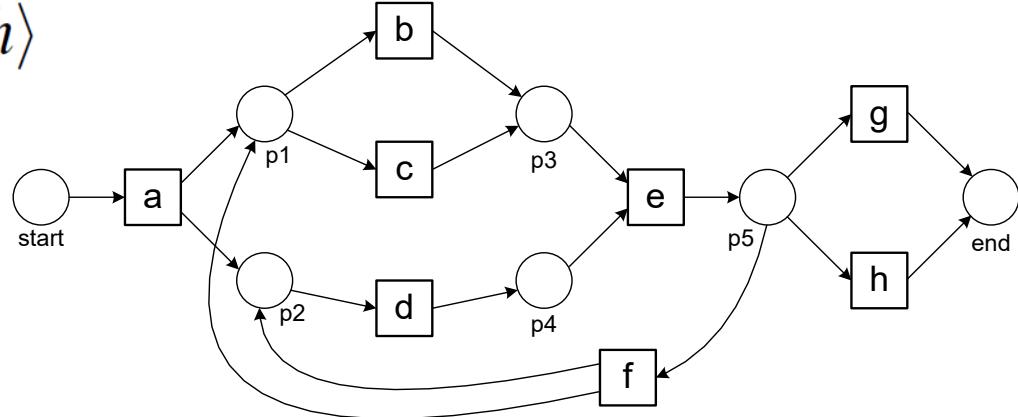
$$\sigma_1 = \langle a, c, d, e, h \rangle$$



Quantifying fitness at the trace level

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

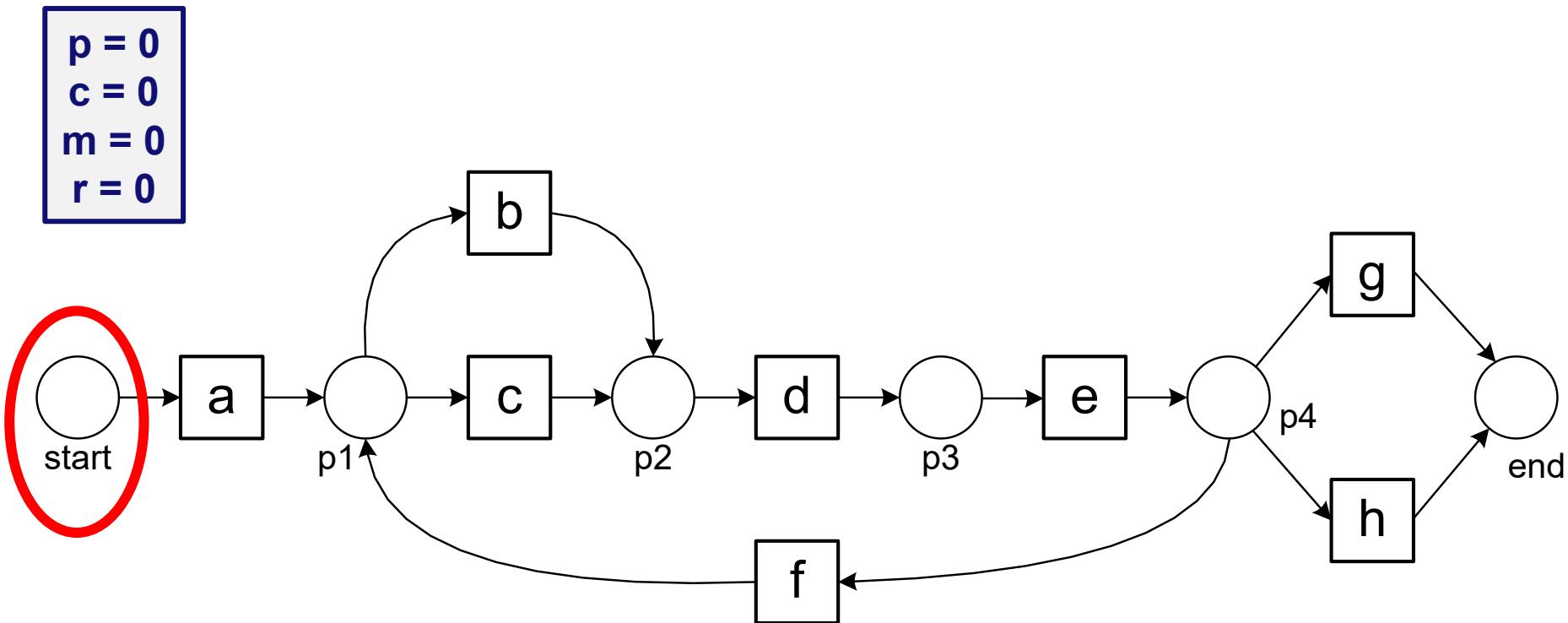
$$\sigma_1 = \langle a, c, d, e, h \rangle$$



$$fitness(\sigma, N) = \frac{1}{2} \left(1 - \frac{0}{7} \right) + \frac{1}{2} \left(1 - \frac{0}{7} \right) = 1$$

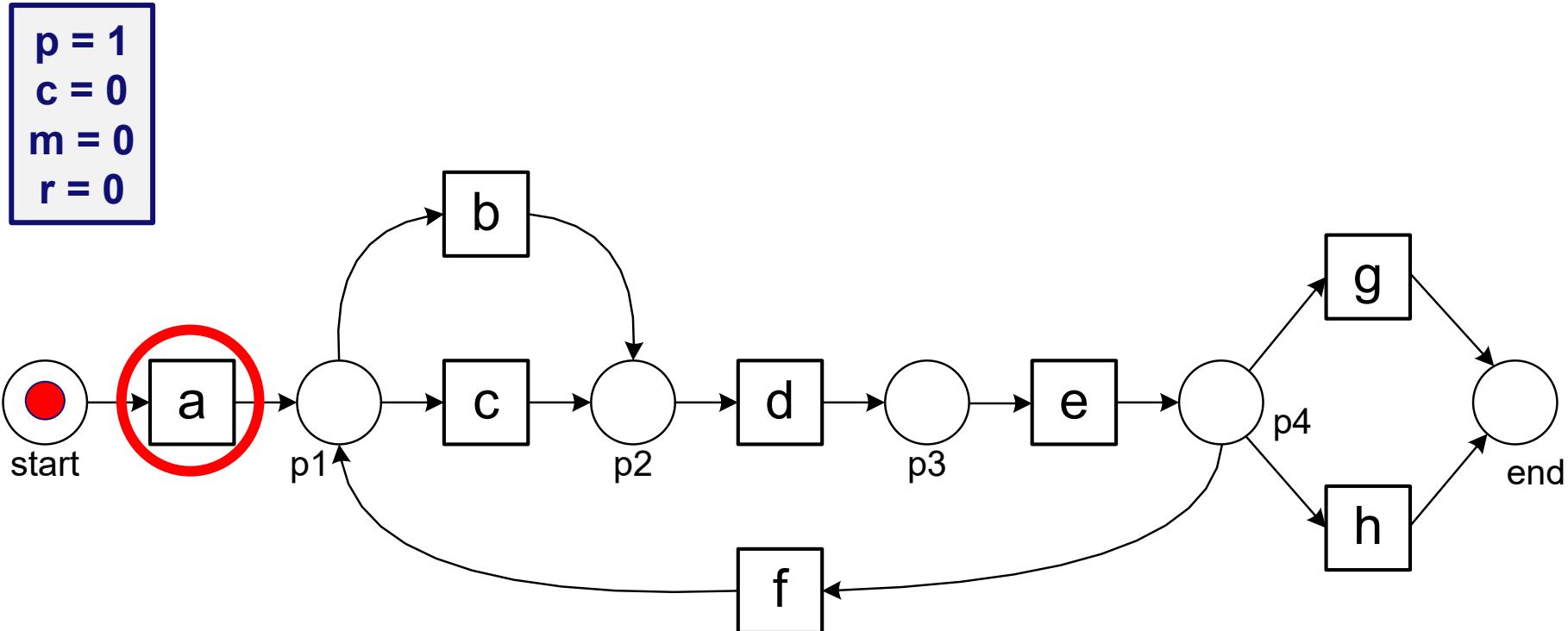
Replaying

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



Replaying

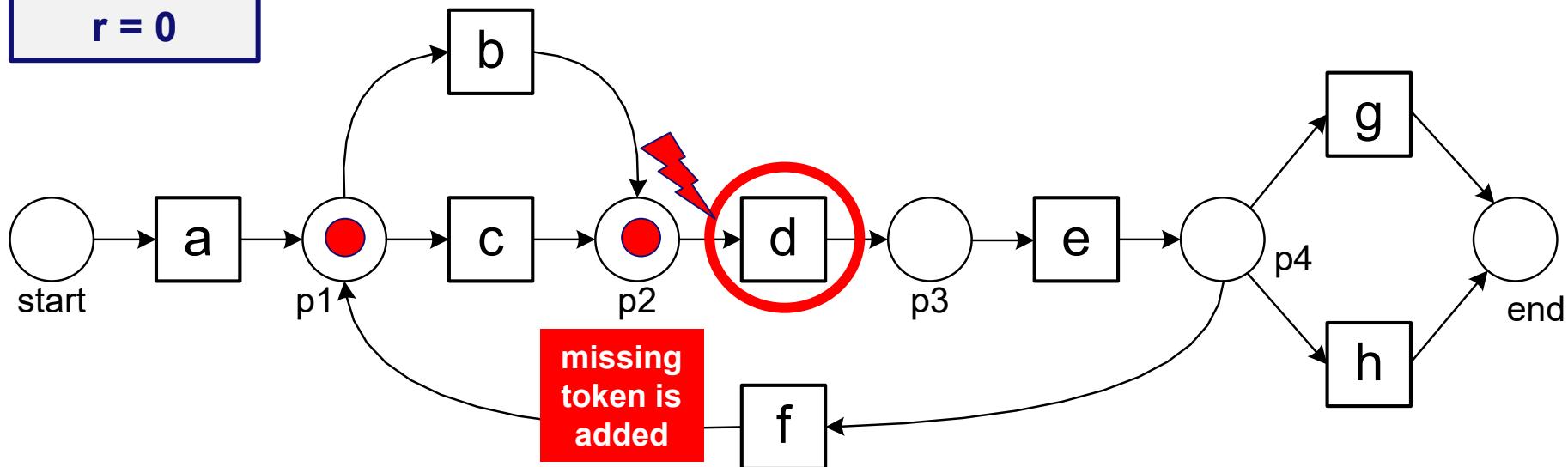
$$\sigma_3 = \langle a, d, c, e, h \rangle$$



Replaying

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

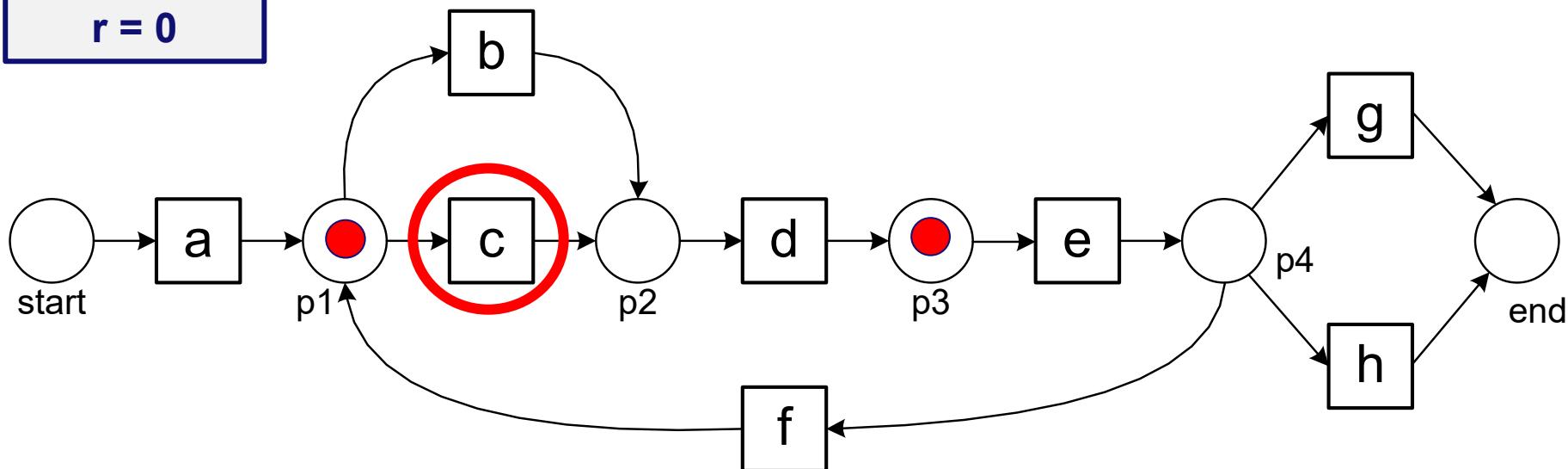
$p = 1+1 = 2$
 $c = 0+1 = 1$
 $m = 0$
 $r = 0$



Replaying

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

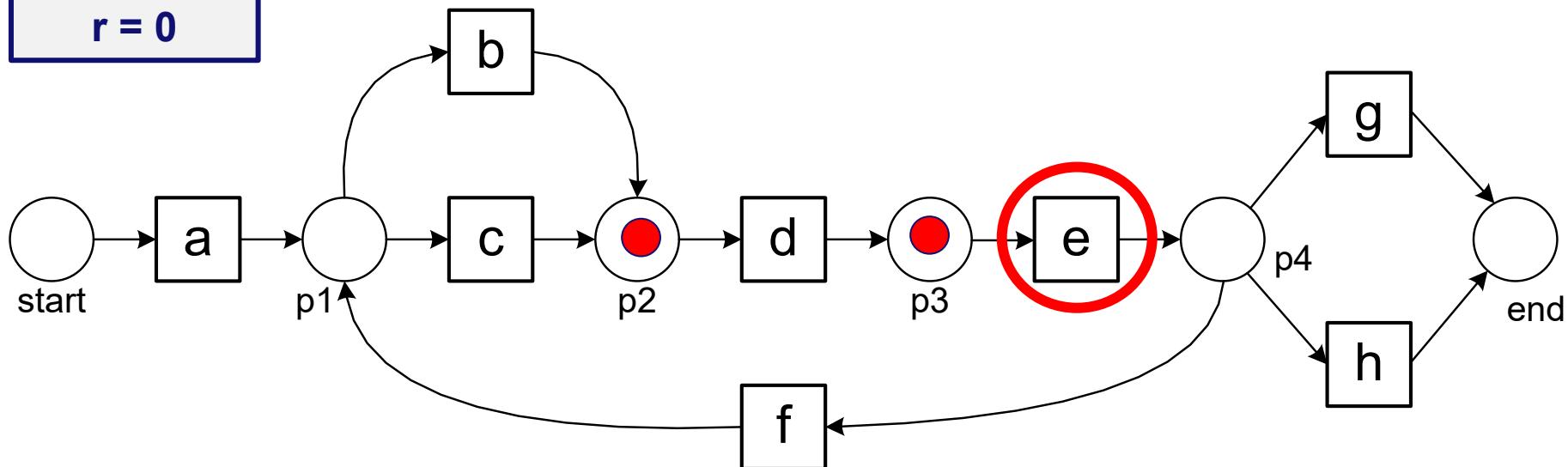
$$\begin{aligned} p &= 2+1 = 3 \\ c &= 1+1 = 2 \\ m &= 0+1 = 1 \\ r &= 0 \end{aligned}$$



Replaying

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

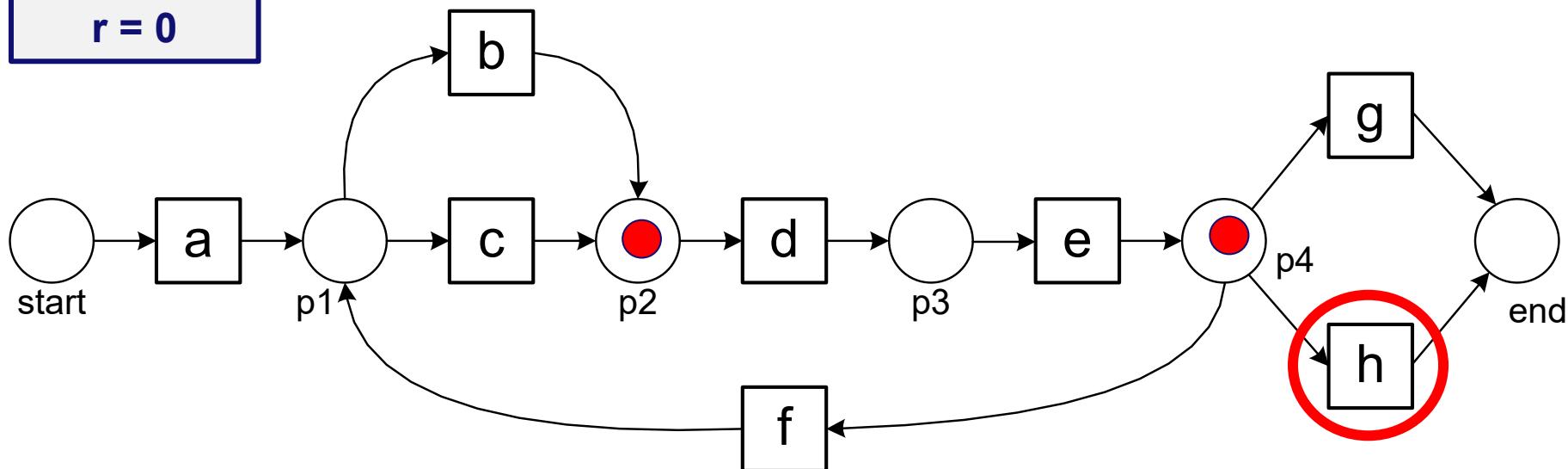
$$\begin{aligned} p &= 3+1 = 4 \\ c &= 2+1 = 3 \\ m &= 1 \\ r &= 0 \end{aligned}$$



Replaying

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

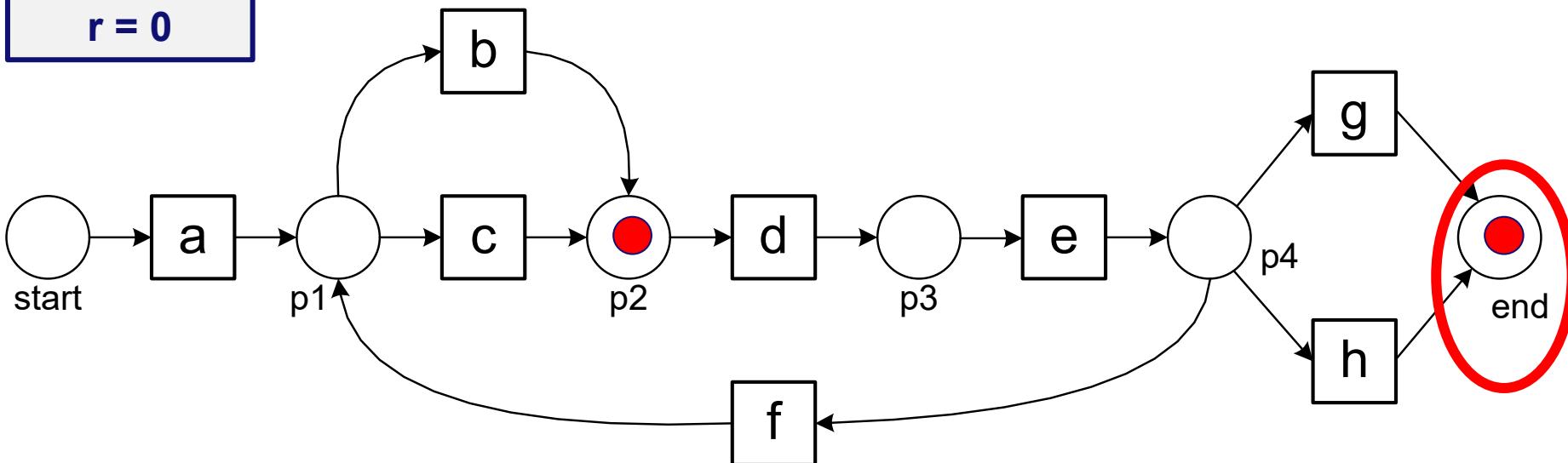
$$\begin{aligned} p &= 4+1 = 5 \\ c &= 3+1 = 4 \\ m &= 1 \\ r &= 0 \end{aligned}$$



Replaying

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

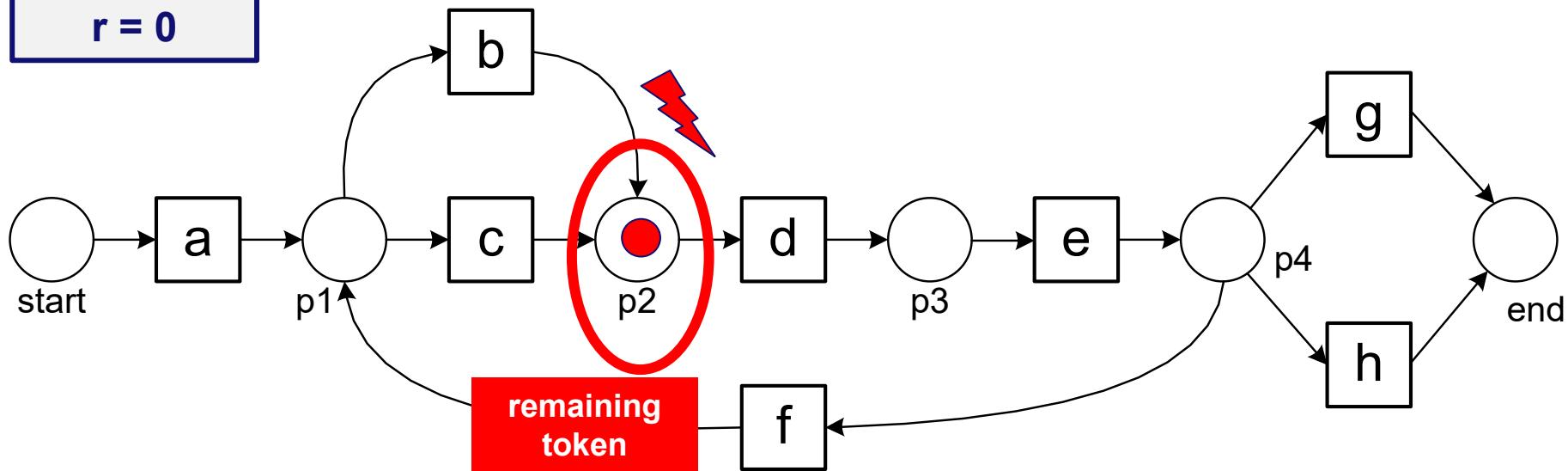
$$\begin{aligned} p &= 5+1 = 6 \\ c &= 4+1 = 5 \\ m &= 1 \\ r &= 0 \end{aligned}$$



Replaying

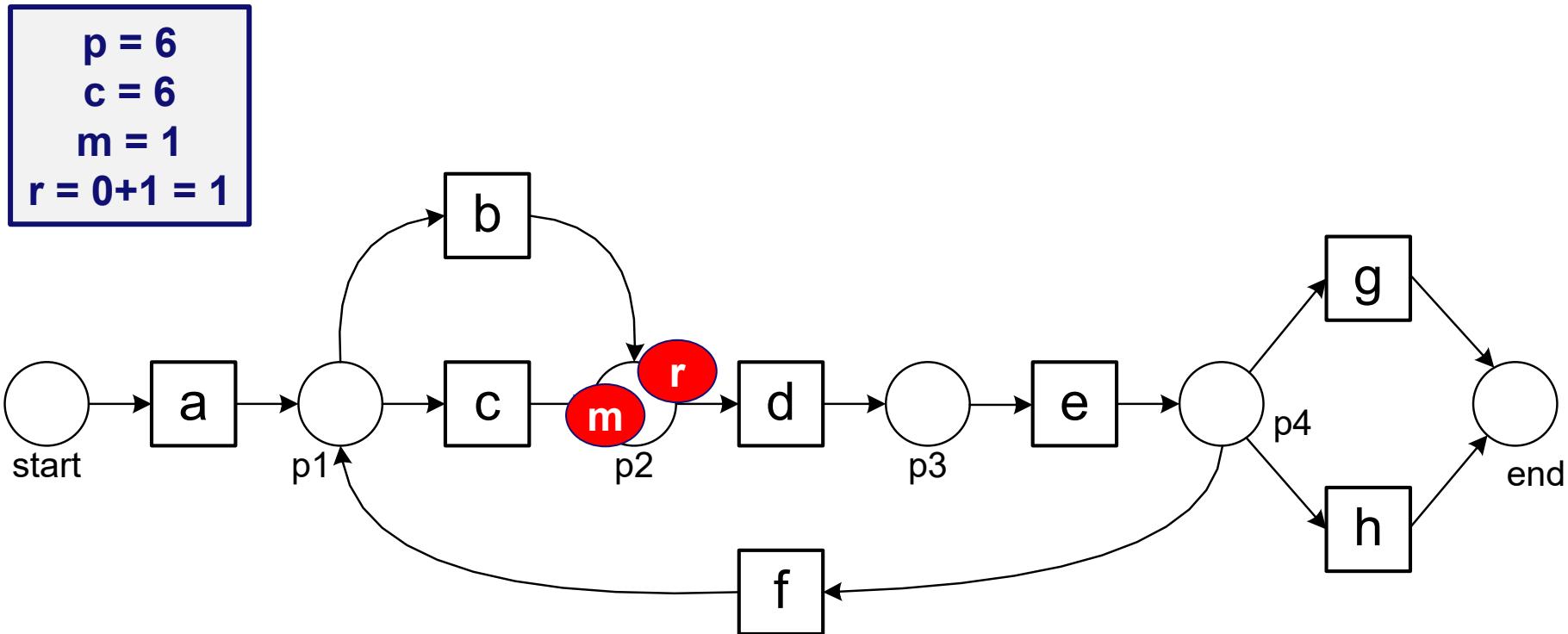
$$\sigma_3 = \langle a, d, c, e, h \rangle$$

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



Replaying

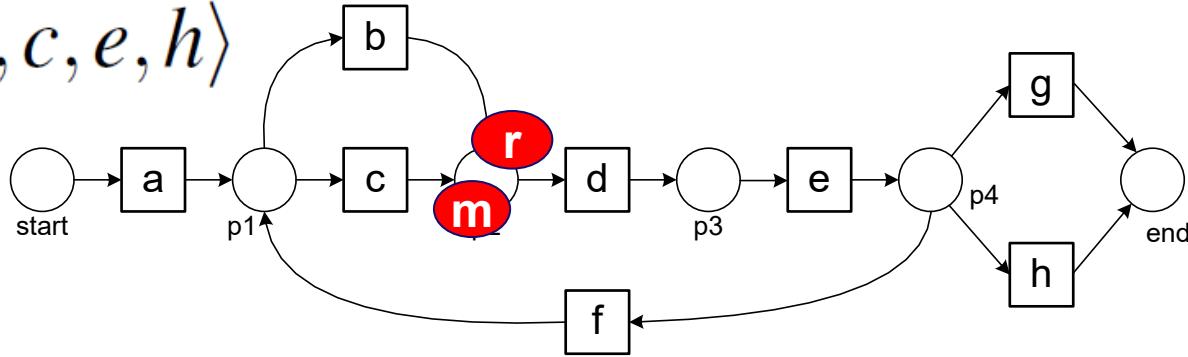
$$\sigma_3 = \langle a, d, c, e, h \rangle$$



Quantifying fitness at the trace level

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

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



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

Fitness at the log level

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

missing tokens

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

consumed tokens

remaining tokens

produced tokens

Looks scar
just needs t
sums of p, c, m, and r
over the multiset of
traces in de

#	trace
455	acdeh
191	abdeg
177	adceh
144	abdeh
111	acdeg
82	adceg
56	adbeh
47	acdefdbeh
38	adbeg
33	acdefbdbeh
14	acdefbddeg
11	acdefdbeg
9	adcefcdbeh
8	adcefdbeh
5	adcefbdeg
3	acdefbdefdbeg
2	adcefdbebeg
2	adcefbdefbdeg
1	adcefdbefbdeh
1	adbefbdefdbeg
1	adcefdbefcdefdbeg
1391	



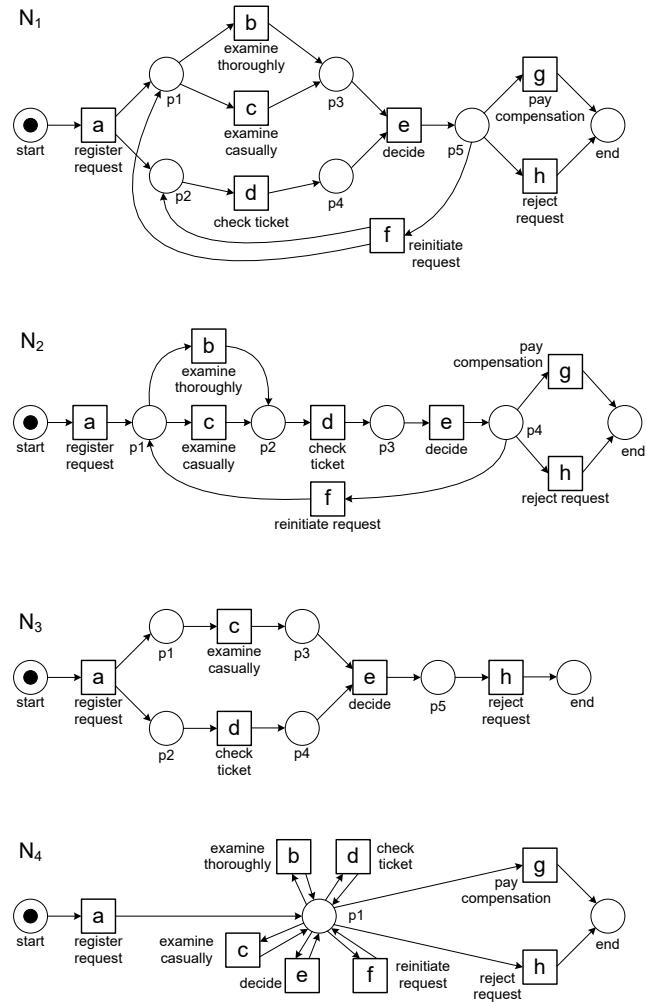
$$fitness(L, N) = \frac{1}{2} \left(1 - \frac{\sum_{\sigma \in L} L(\sigma) \times m_{N,\sigma}}{\sum_{\sigma \in L} L(\sigma) \times c_{N,\sigma}} \right) + \frac{1}{2} \left(1 - \frac{\sum_{\sigma \in L} L(\sigma) \times r_{N,\sigma}}{\sum_{\sigma \in L} L(\sigma) \times p_{N,\sigma}} \right)$$

$$fitness(L_{full}, N_1) = 1$$

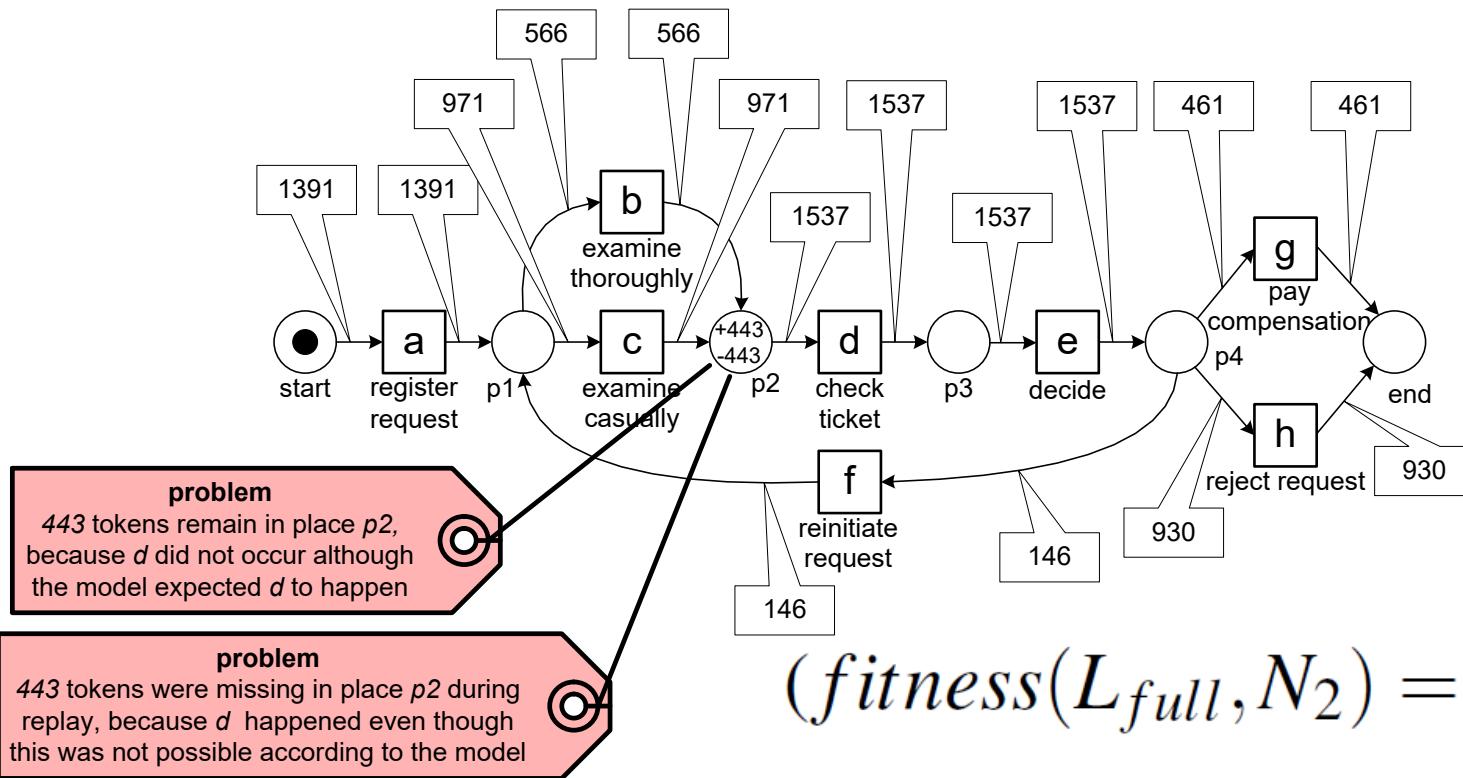
$$fitness(L_{full}, N_2) = 0.9504$$

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

$$fitness(L_{full}, N_4) = 1$$



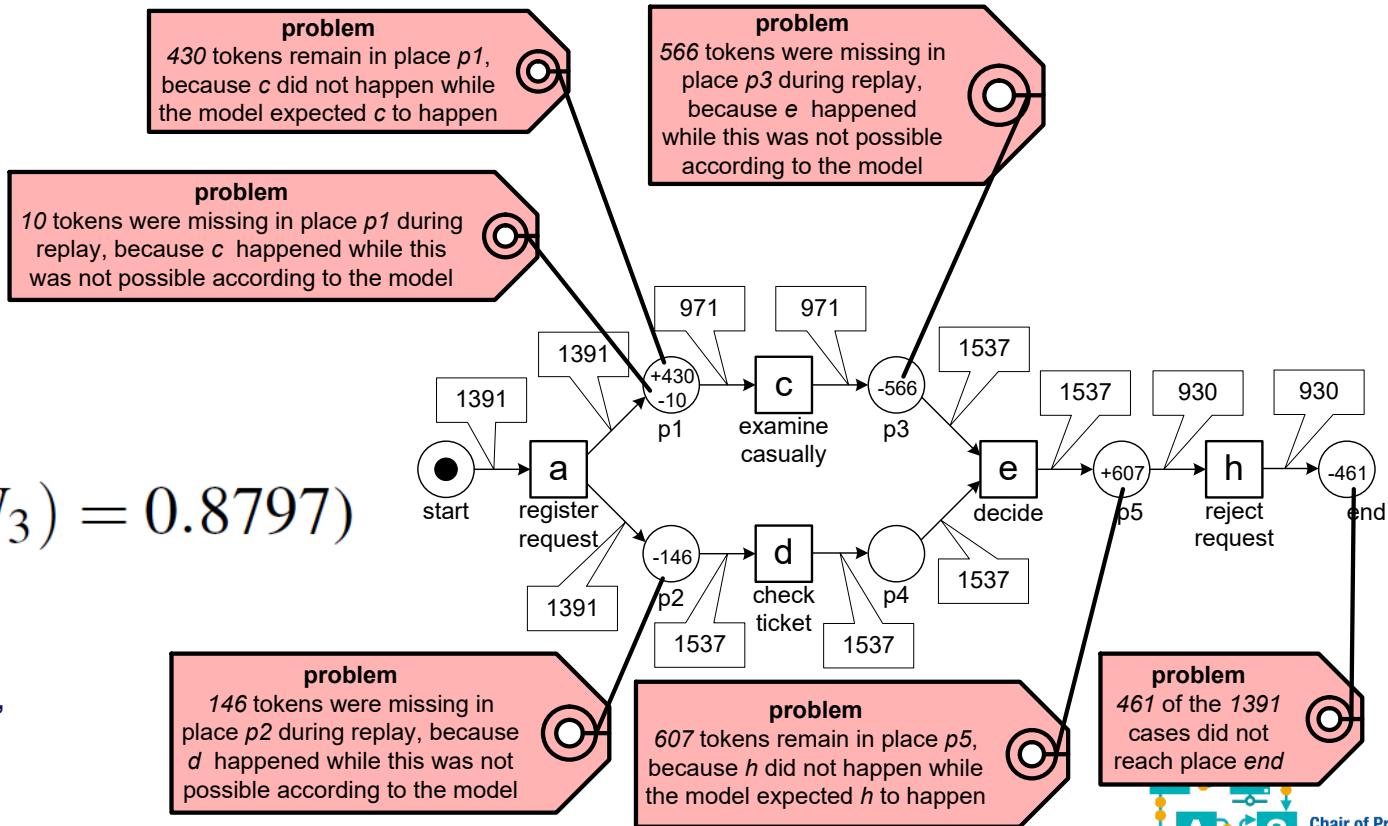
Diagnostics



Diagnostics

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

Remark: If event log and model consider different sets of activities, this should be addressed first.
Activities in log, but not in model, are simply ignored.

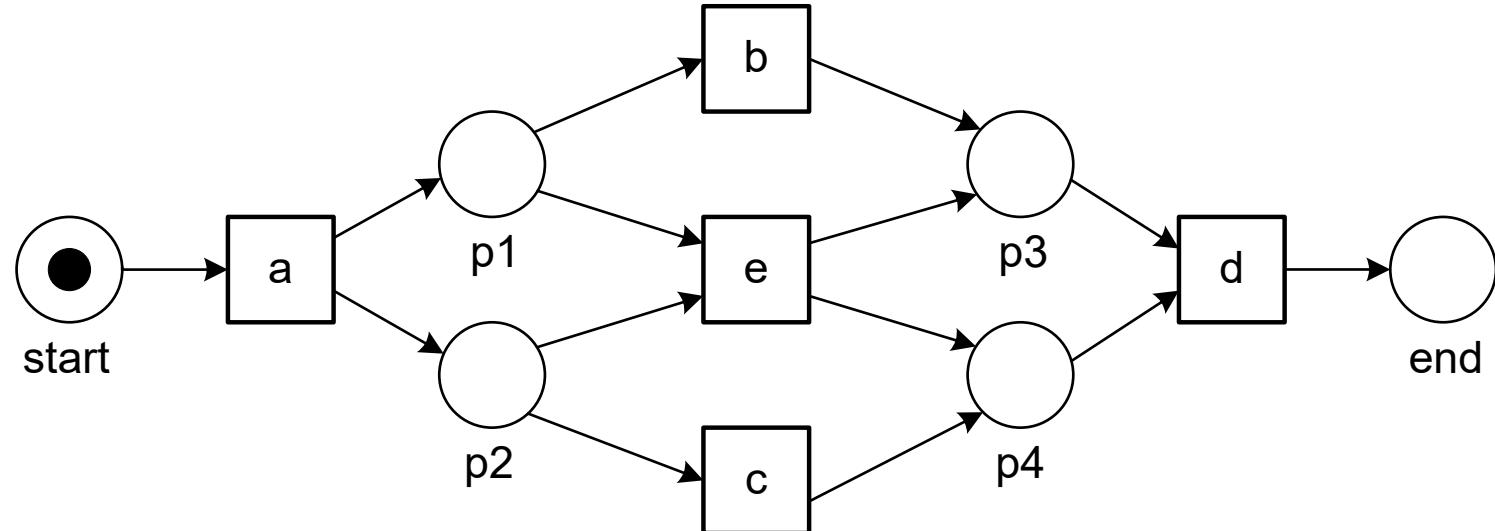


Examples

Question (may take some time)

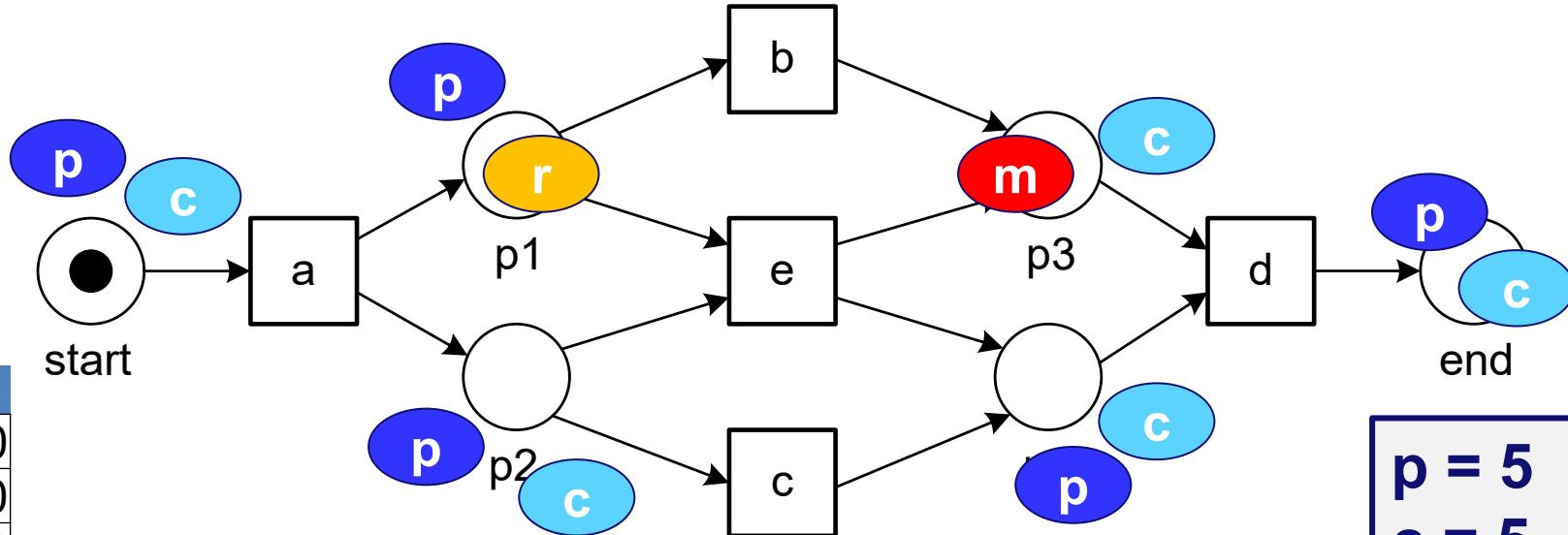
Compute fitness using missing and remaining tokens

trace	frequency
abcd	10
acbd	10
aed	10
abd	2
acd	1
ad	1
abbd	1



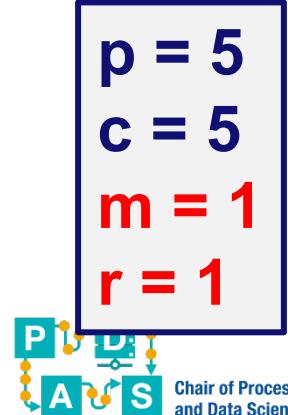
- Consider the event log containing 35 cases.
- What is the fitness?

Let us pick one trace: acd



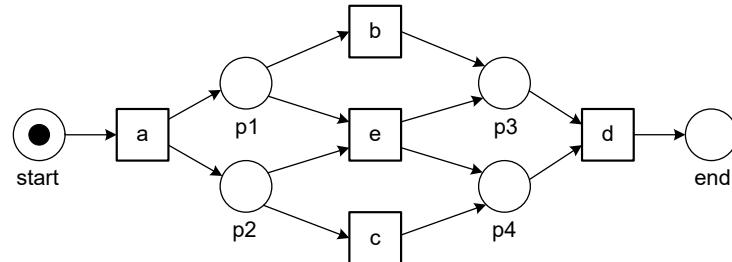
trace	frequency
abcd	10
acbd	10
aed	10
abd	2
acd	1
ad	1
abbd	1

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



Fitness = 0.9658

trace	frequency	produced tokens (p)	remaining tokens (r)	consumed tokens (c)	missing tokens (m)	produced tokens (pII)	remaining tokens (rII)	consumed tokens (cII)	missing tokens (mII)
abcd	10	6	0	6	0	60	0	60	0
acbd	10	6	0	6	0	60	0	60	0
aed	10	6	0	6	0	60	0	60	0
abd	2	5	1	5	1	10	2	10	2
acd	1	5	1	5	1	5	1	5	1
ad	1	4	2	4	2	4	2	4	2
abbd	1	6	2	6	2	6	2	6	2



205	7	205	7
sum p	sum r	sum c	sum m

fitness	0.965853659
---------	-------------

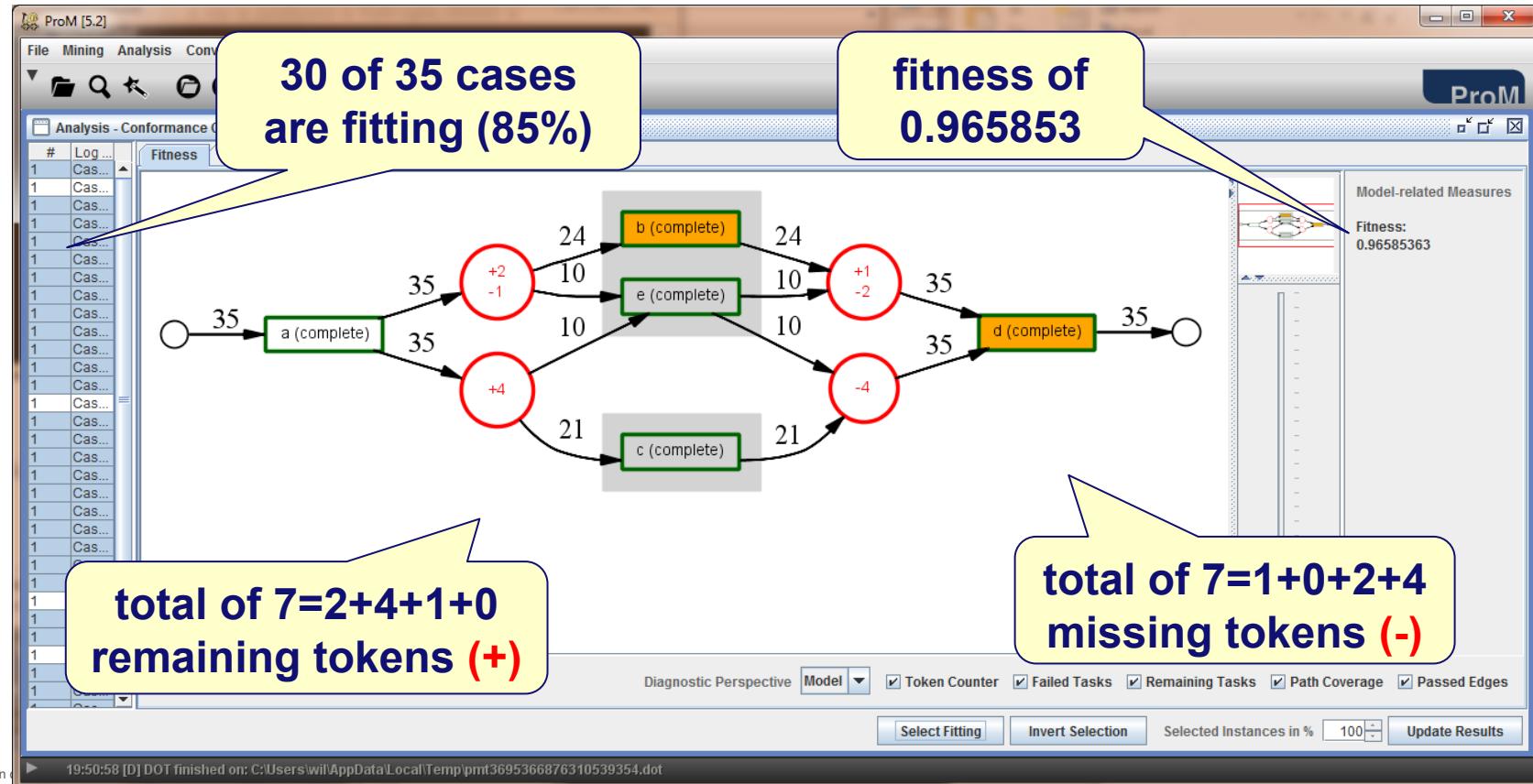
$$fitness(L, N) = \frac{1}{2} \left(1 - \frac{\sum_{\sigma \in L} L(\sigma) \times m_{N,\sigma}}{\sum_{\sigma \in L} L(\sigma) \times c_{N,\sigma}} \right) + \frac{1}{2} \left(1 - \frac{\sum_{\sigma \in L} L(\sigma) \times r_{N,\sigma}}{\sum_{\sigma \in L} L(\sigma) \times p_{N,\sigma}} \right)$$



ProM 5.2 output

Note that PM4Py and Celonis support variants of token-based replay.

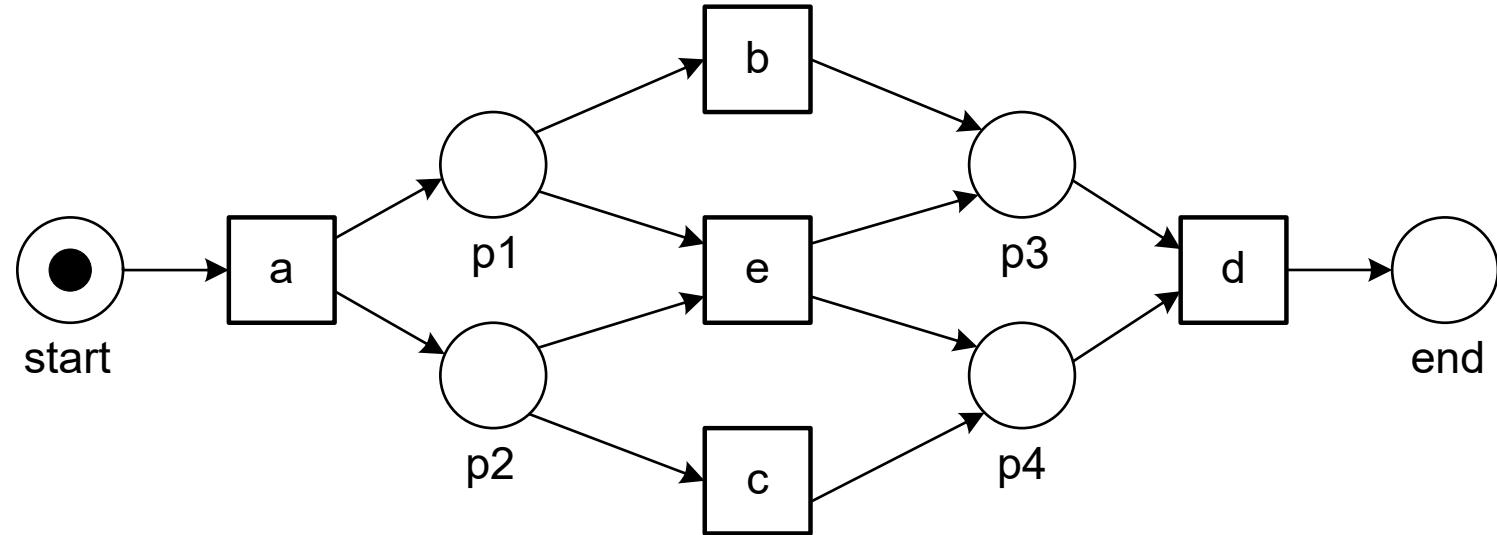
(ProM 6 only supports more advanced conformance checking techniques like alignments)



Question

Compute fitness using missing and remaining tokens

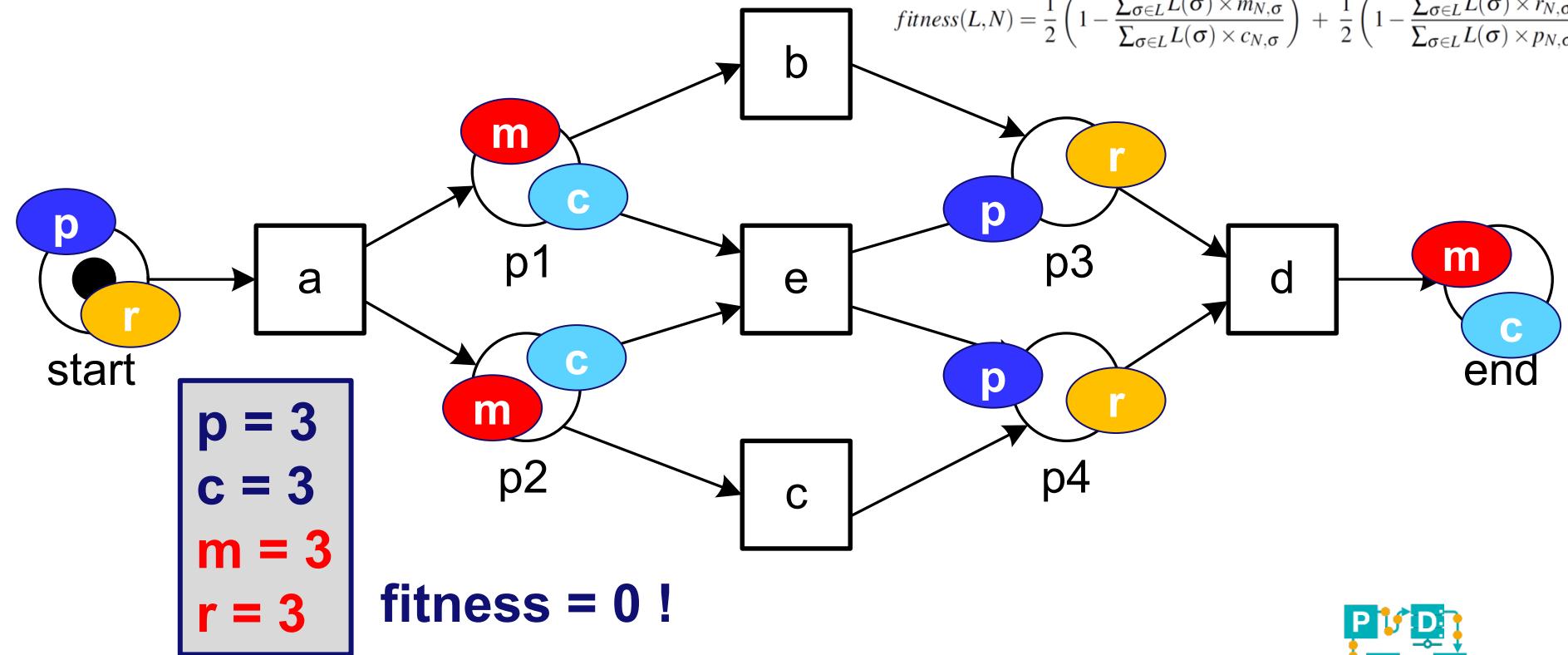
trace	frequency
e	1



- Consider the event log containing just one case: $L = [\langle e \rangle]$.
- What is the fitness (using token-based replay)?

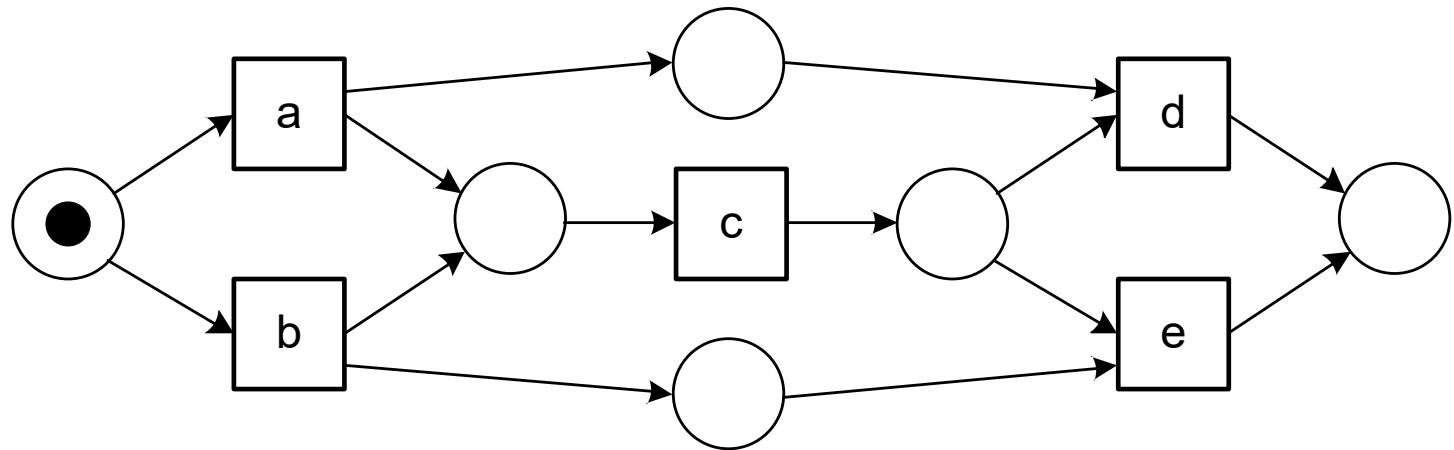
Answer obtained by replaying $\langle e \rangle$

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



Another example

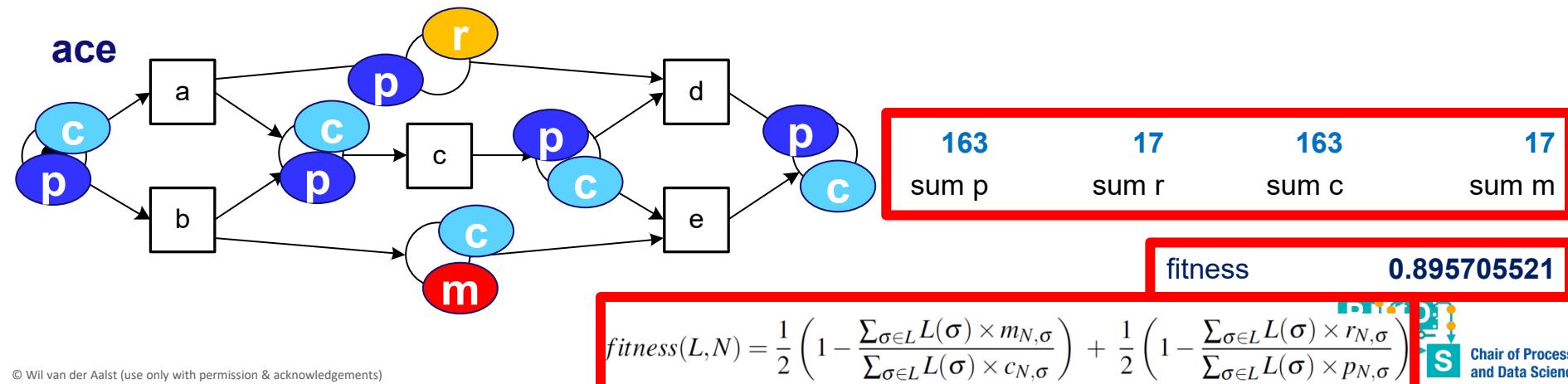
trace	frequency
acd	10
bce	10
ace	5
bcd	5
dca	1
abd	1
d	1



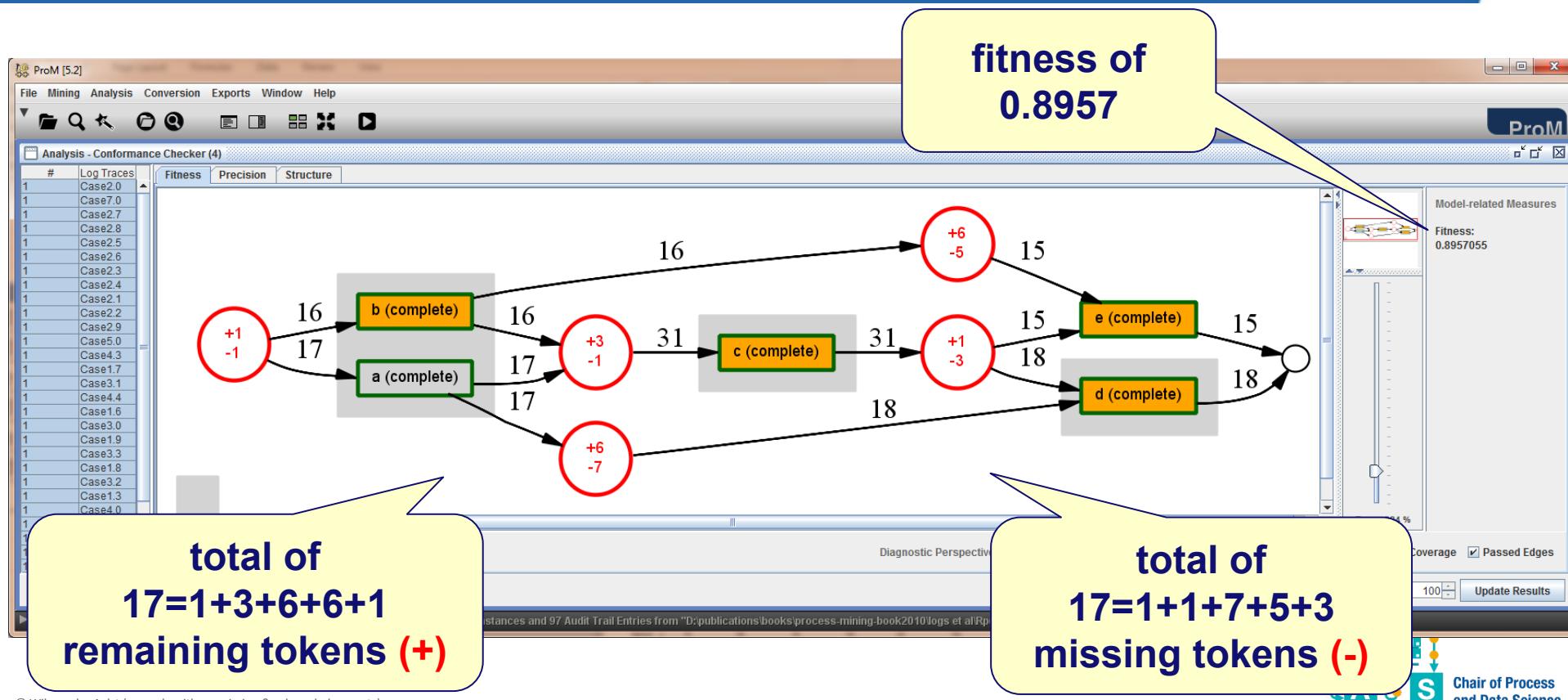
- Consider the event log containing 33 cases.
- What is the fitness?

Fitness = 0.895705521

trace	frequency	produced tokens (p)	remaining tokens (r)	consumed tokens (c)	missing tokens (m)	produced tokens (all)	remaining tokens (all)	consumed tokens (all)	missing tokens (all)
acd	10	5	0	5	0	50	0	50	0
bce	10	5	0	5	0	50	0	50	0
ace	5	5	1	5	1	25	5	25	5
bcd	5	5	1	5	1	25	5	25	5
dca	1	5	3	5	3	5	3	5	3
abd	1	6	3	5	2	6	3	5	2
d	1	2	1	3	2	2	1	3	2



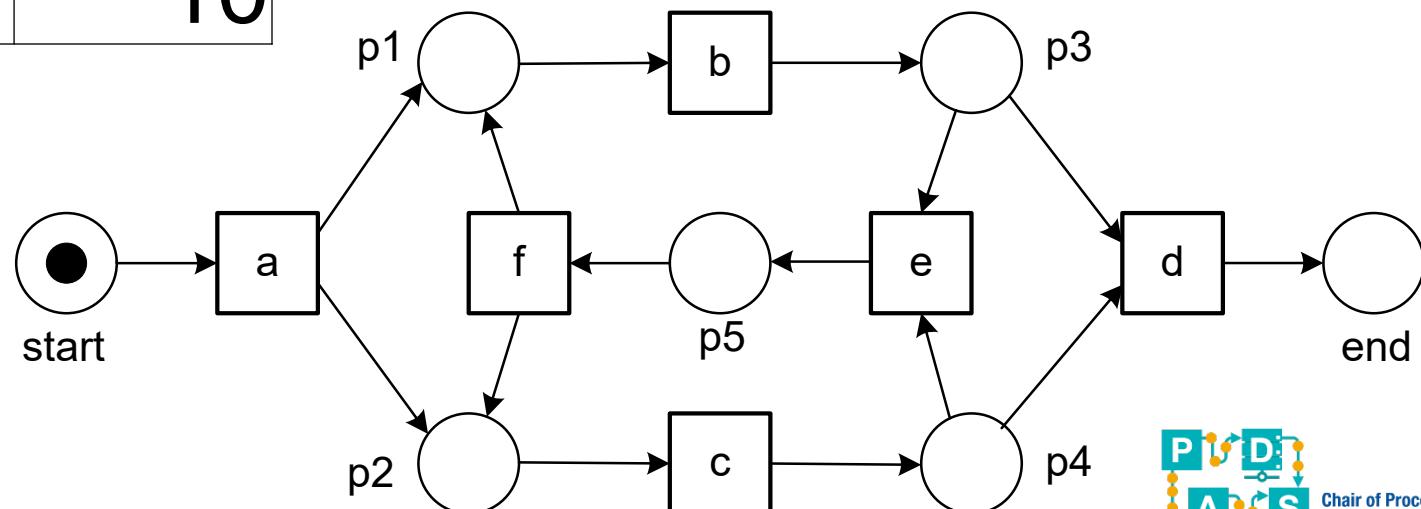
ProM 5.2 diagnostics



Another example

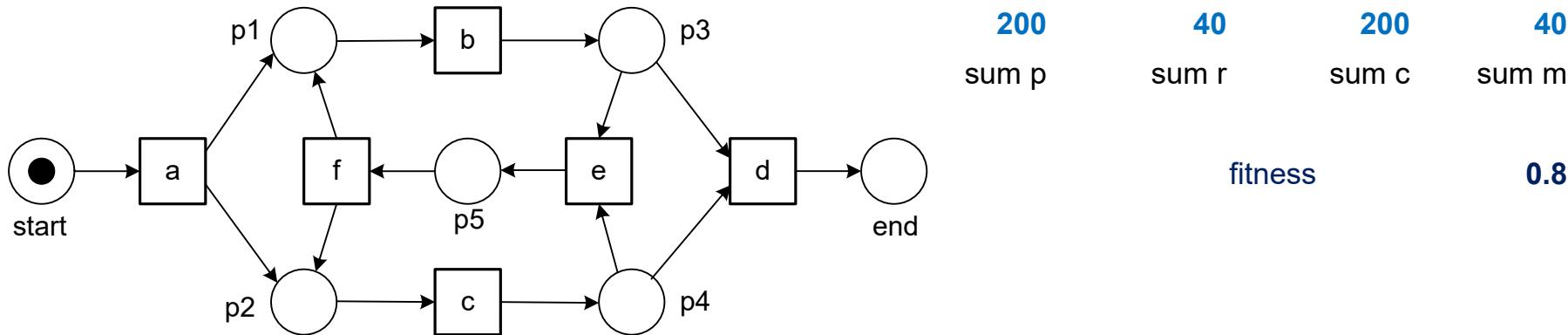
trace	frequency
abefcd	10
abbefcccd	10

- Consider the event log containing 20 cases.
- What is the fitness?



Fitness = 0.8

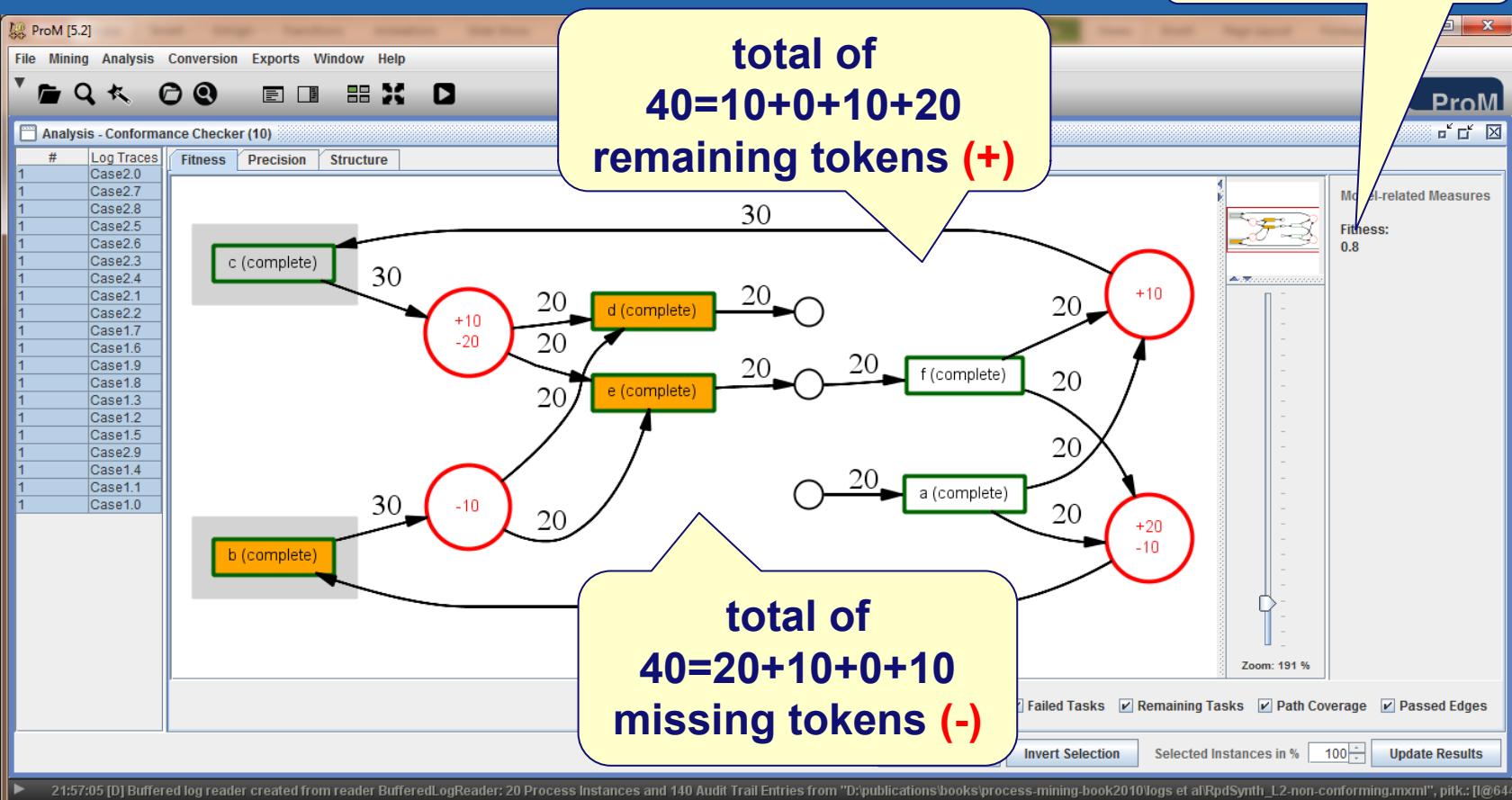
trace	frequency	produced tokens (p)	remaining tokens (r)	consumed tokens (c)	missing tokens (m)	produced tokens (all)	remaining tokens (all)	consumed tokens (all)	missing tokens (all)
abefcd	10	9	2	9	2	90	20	90	20
abbefcccd	10	11	2	11	2	110	20	110	20



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

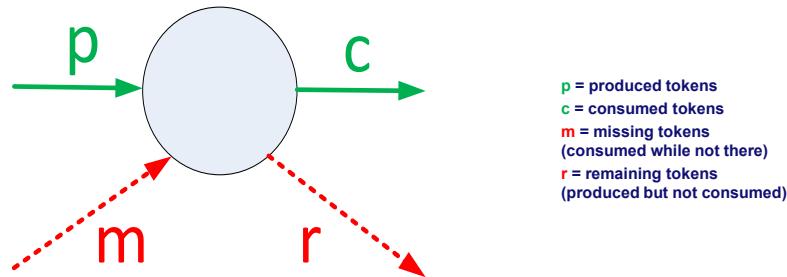
ProM 5.2 diagnostics

fitness of 0.8



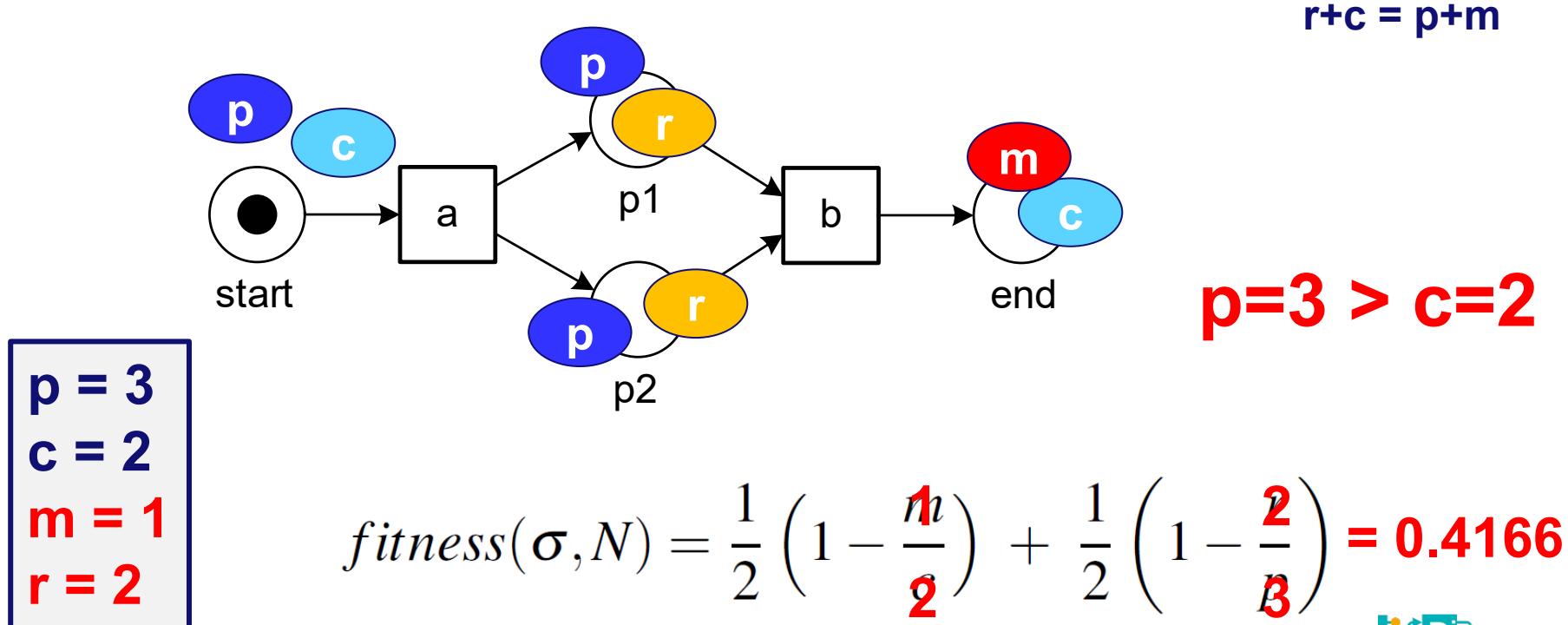
Thus far $c=p$ and $m=r$. Always?

- Provide, if possible, a log and model such that $c \neq p$ and $m \neq r$ at end.

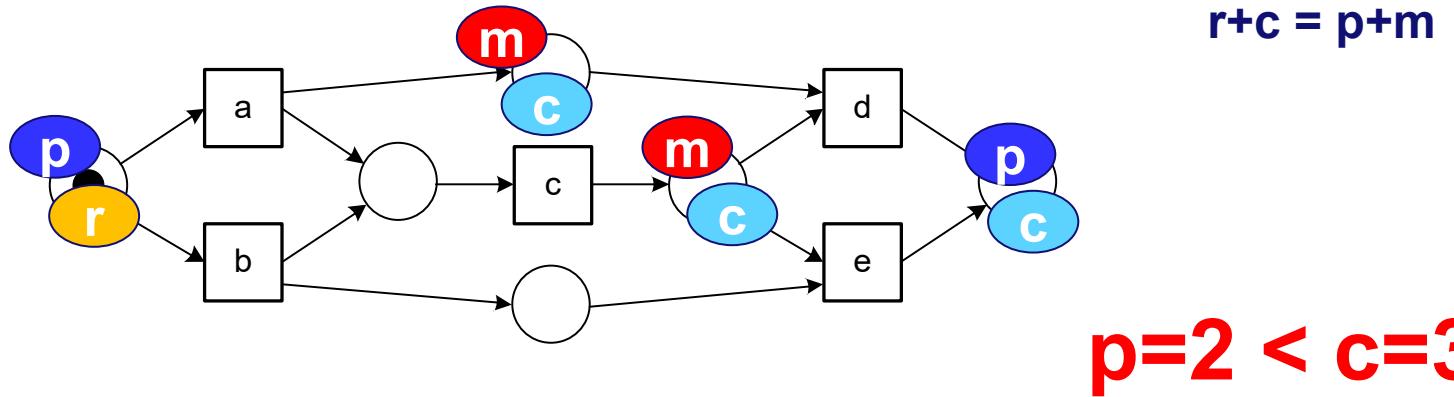


Hint: Recall that $r = p+m-c$ (i.e., $p+m=c+r$) at end.

Example: Model below and event log [$\langle a \rangle$]



Example: Model below and event log [$\langle d \rangle$]



$$\begin{aligned} p &= 2 \\ c &= 3 \\ m &= 2 \\ r &= 1 \end{aligned}$$

$$fitness(\sigma, N) = \frac{1}{2} \left(1 - \frac{2}{3} \right) + \frac{1}{2} \left(1 - \frac{1}{2} \right) = 0.4166$$

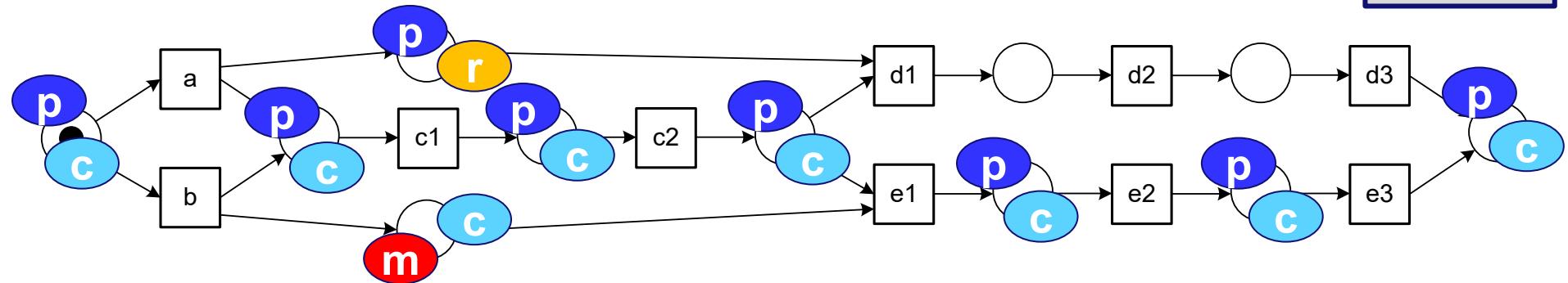
Limitations

- Basic replay approach assumes **visible & uniquely labeled transitions**.
- ProM implementation uses **heuristics** to deal with silent transitions and multiple transitions having the same label.
- Conformance values sometimes **too optimistic** due to "token flooding".
- Local decision making may lead to misleading results.

Local decision making is not enough ...

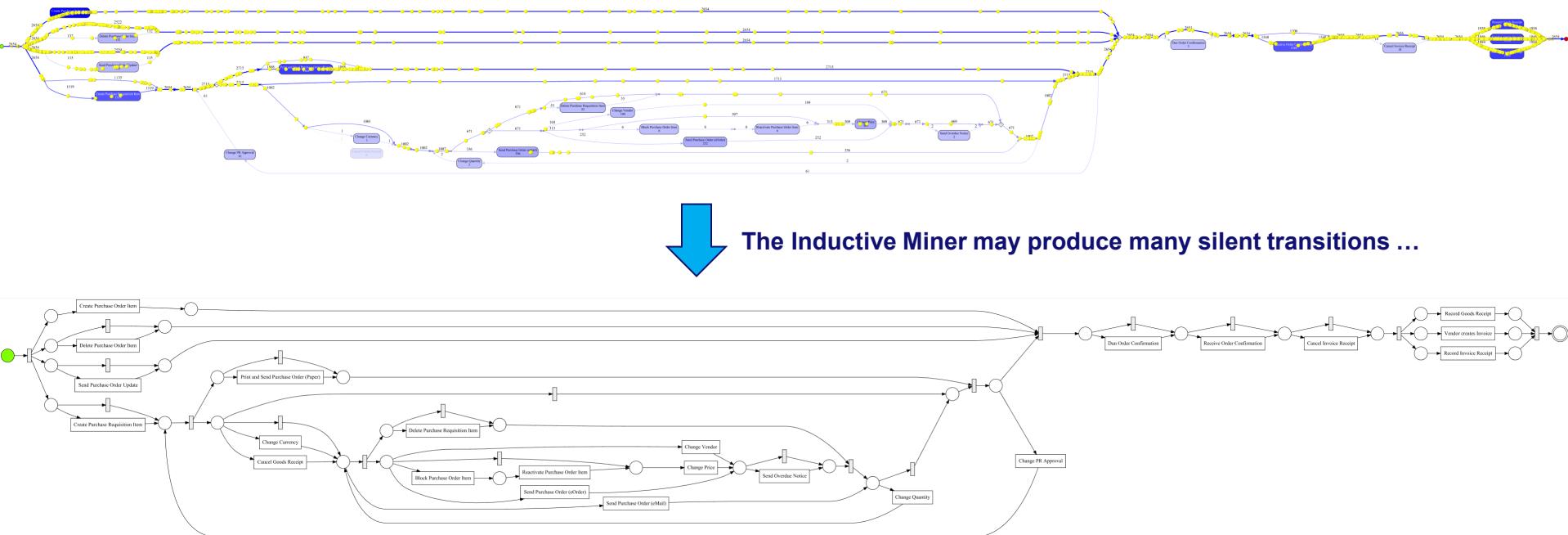
$p = 8$
 $c = 8$
 $m = 1$
 $r = 1$
 $f = 0.875$

$\langle a, c1, c2, e1, e2, e3 \rangle$



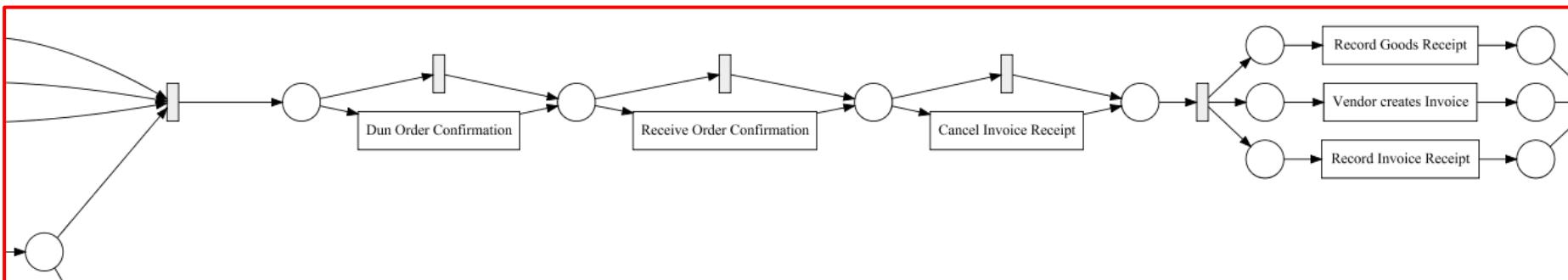
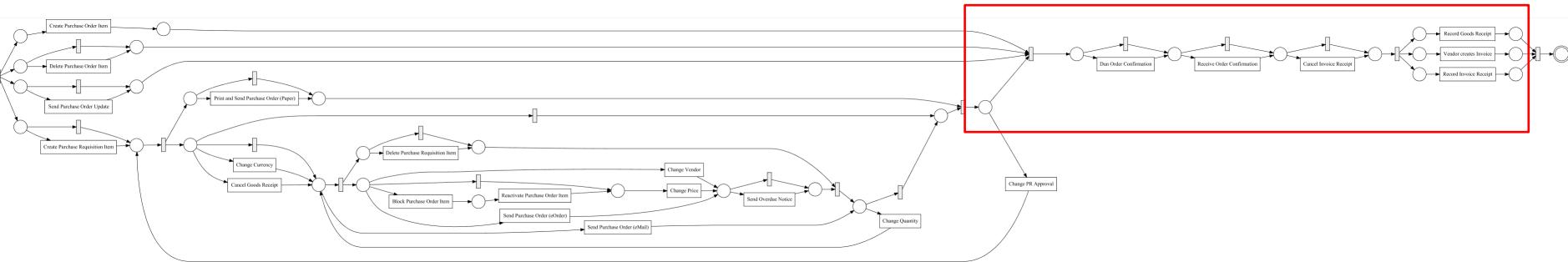
- Replay technique does **not** provide a corresponding path through the model (vital for conformance/performance analysis and other diagnostics).
- We would like to see the "closest path", i.e., $\langle b, c1, c2, e1, e2, e3 \rangle$.

Challenge: Silent and duplicate transitions



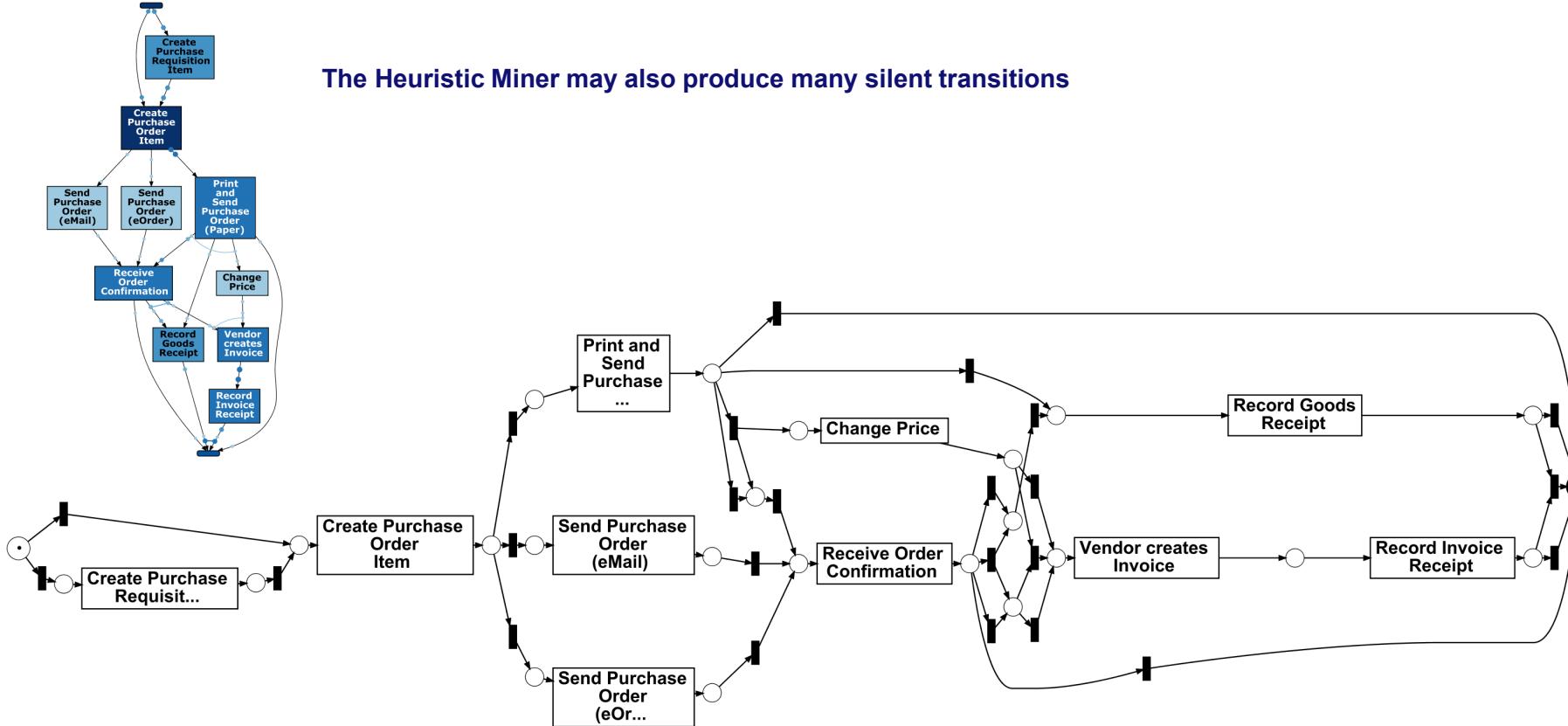
Chair of Process
and Data Science

Challenge: Silent and duplicate transitions



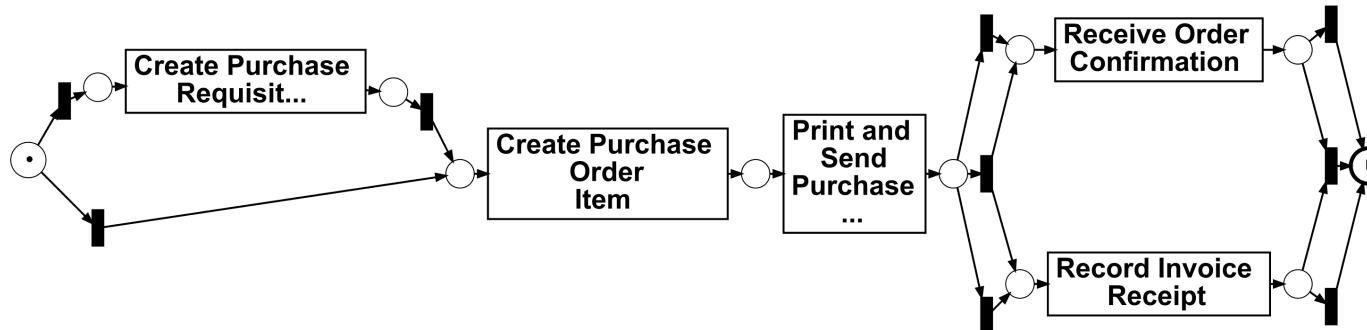
Challenge: Silent and duplicate transitions

The Heuristic Miner may also produce many silent transitions

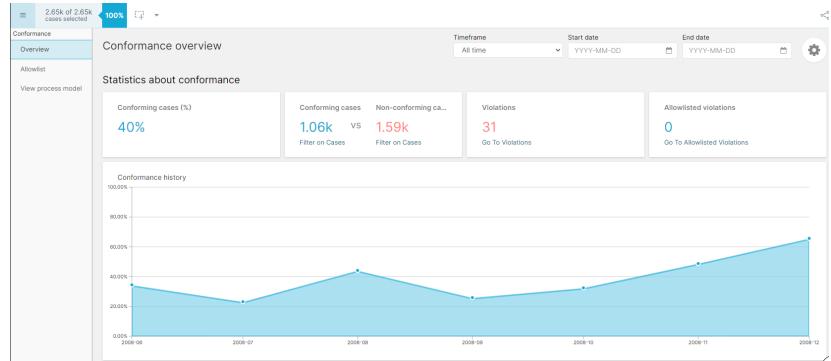
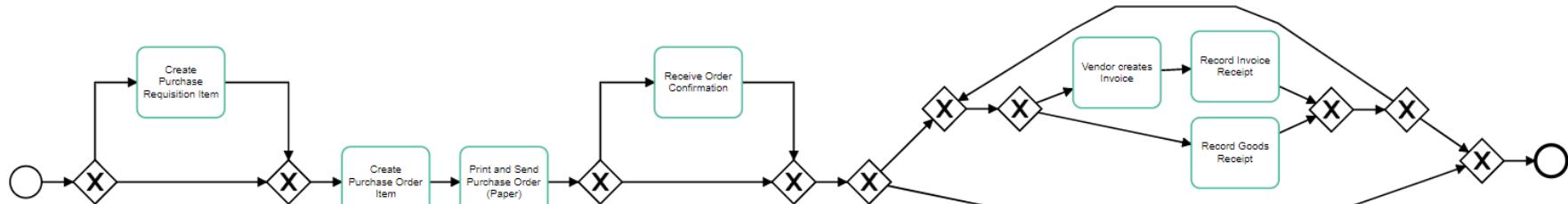


Challenge: Silent and duplicate transitions

- Some form of **state-space exploration** is needed to decide which silent transitions to fire.
- If a trace is **non-fitting**, the set of possibilities grows further.



Celonis's classical conformance checker uses a variant of token-based replay



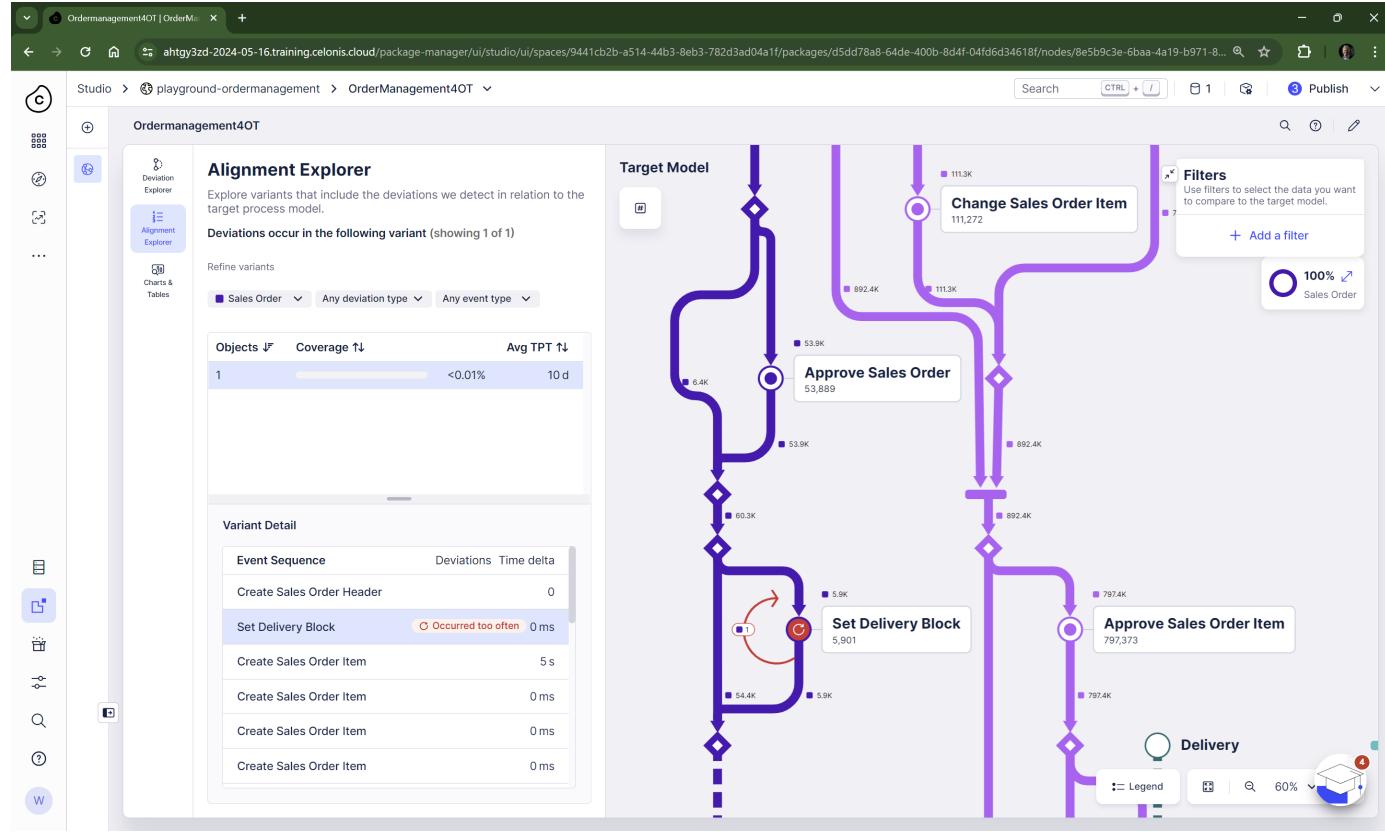
Violations

23% of cases	Change Price is an undesired activity View cases in... Effect on throughput time: 16 Days longer Effect on steps per case: + 2.6 Steps per case
14% of cases	Create Purchase Order Item is followed by Receive Order Confirmation View cases in... Effect on throughput time: 9 Days longer Effect on steps per case: + 1.6 Steps per case
13% of cases	Send Purchase Order (eMail) is an undesired activity View cases in... Effect on throughput time: 4 Days longer Effect on steps per case: + 0.8 Steps per case



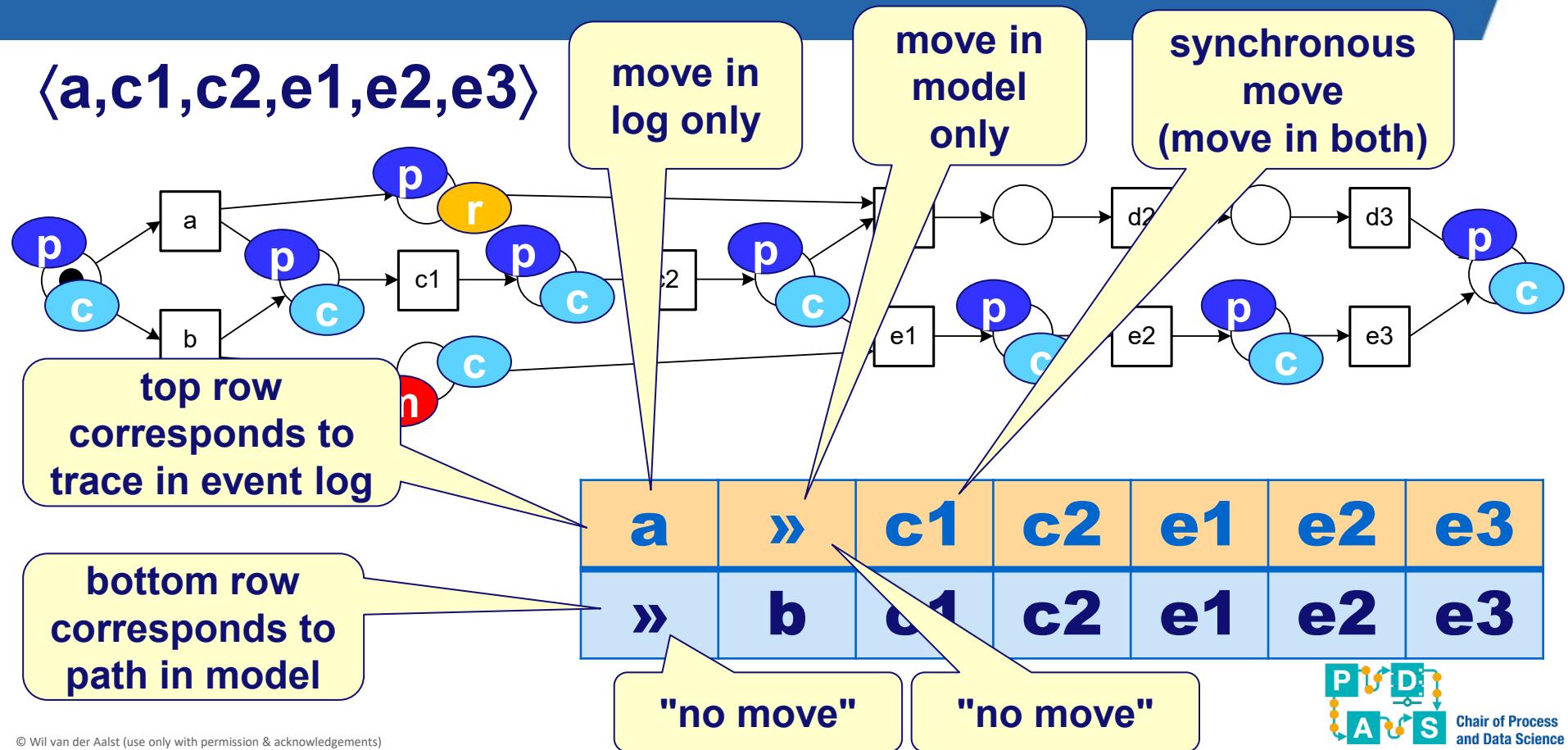
Chair of Process
and Data Science

Newer capabilities of Celonis (e.g. PAM) use alignments



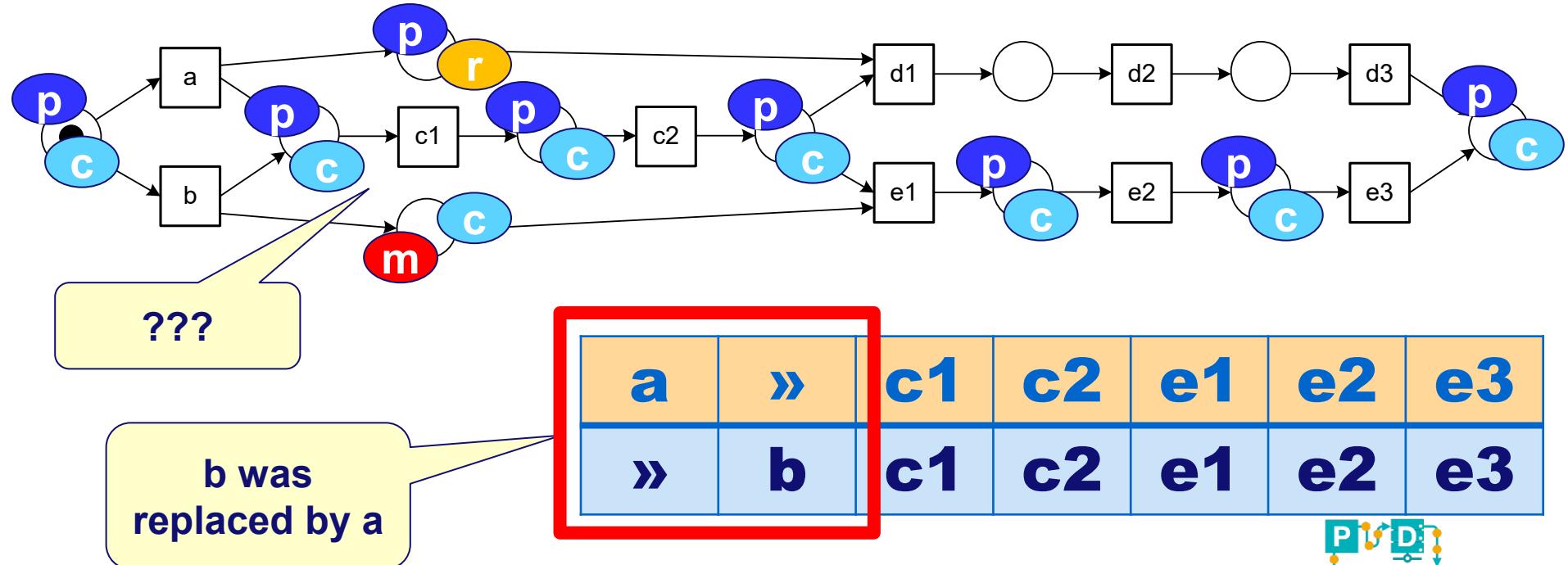
Next: alignments

$\langle a, c1, c2, e1, e2, e3 \rangle$

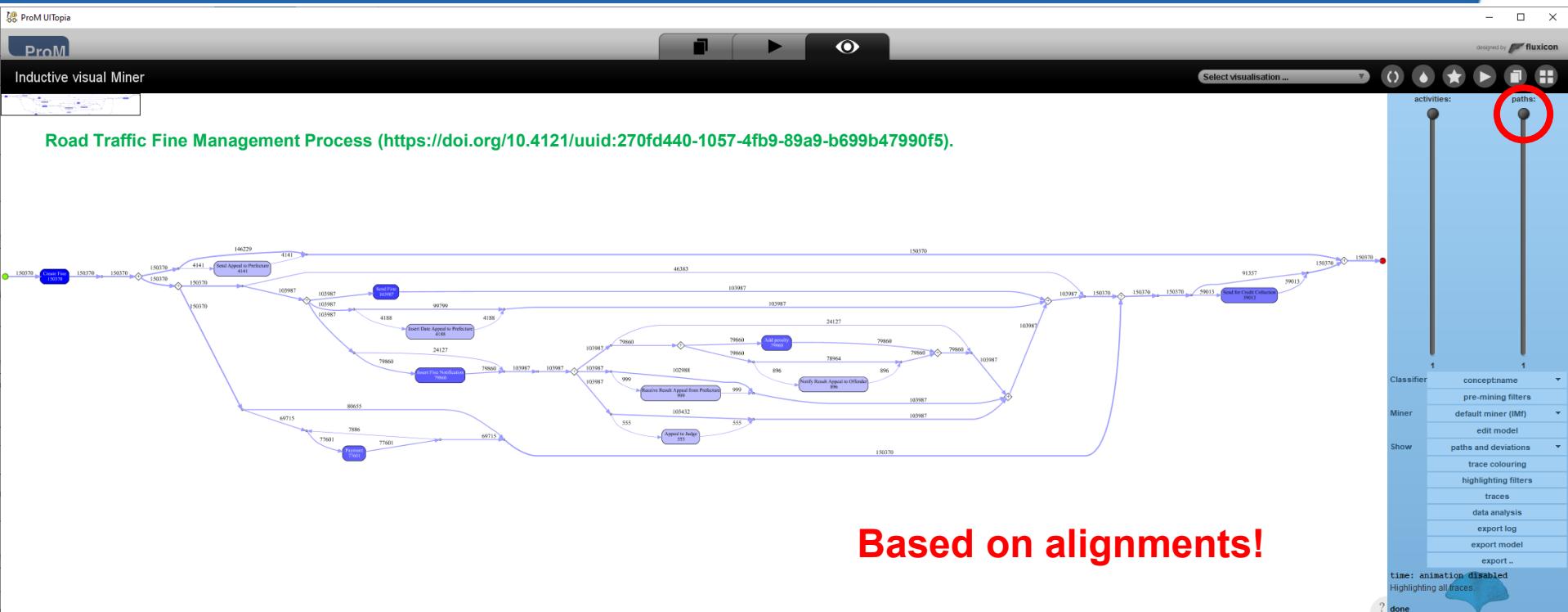


Alignments provide better diagnostics

$\langle a, c1, c2, e1, e2, e3 \rangle$

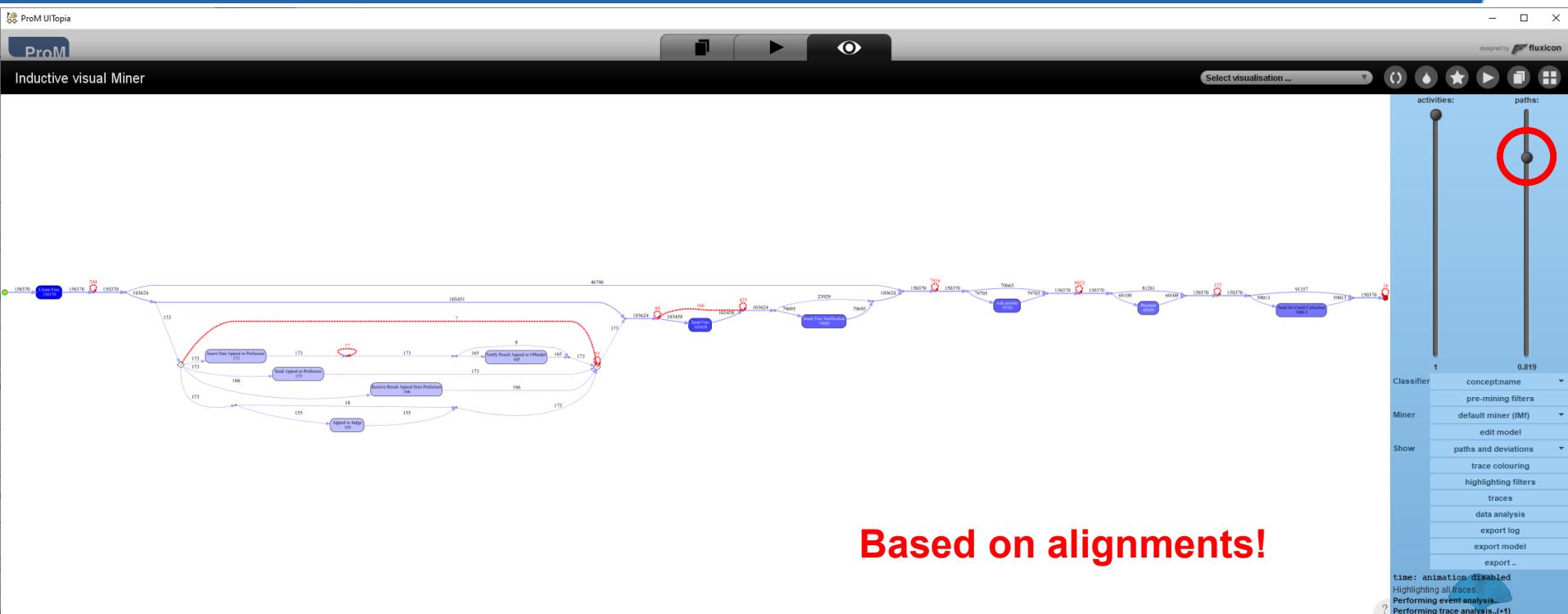


Example: Inductive Visual Miner



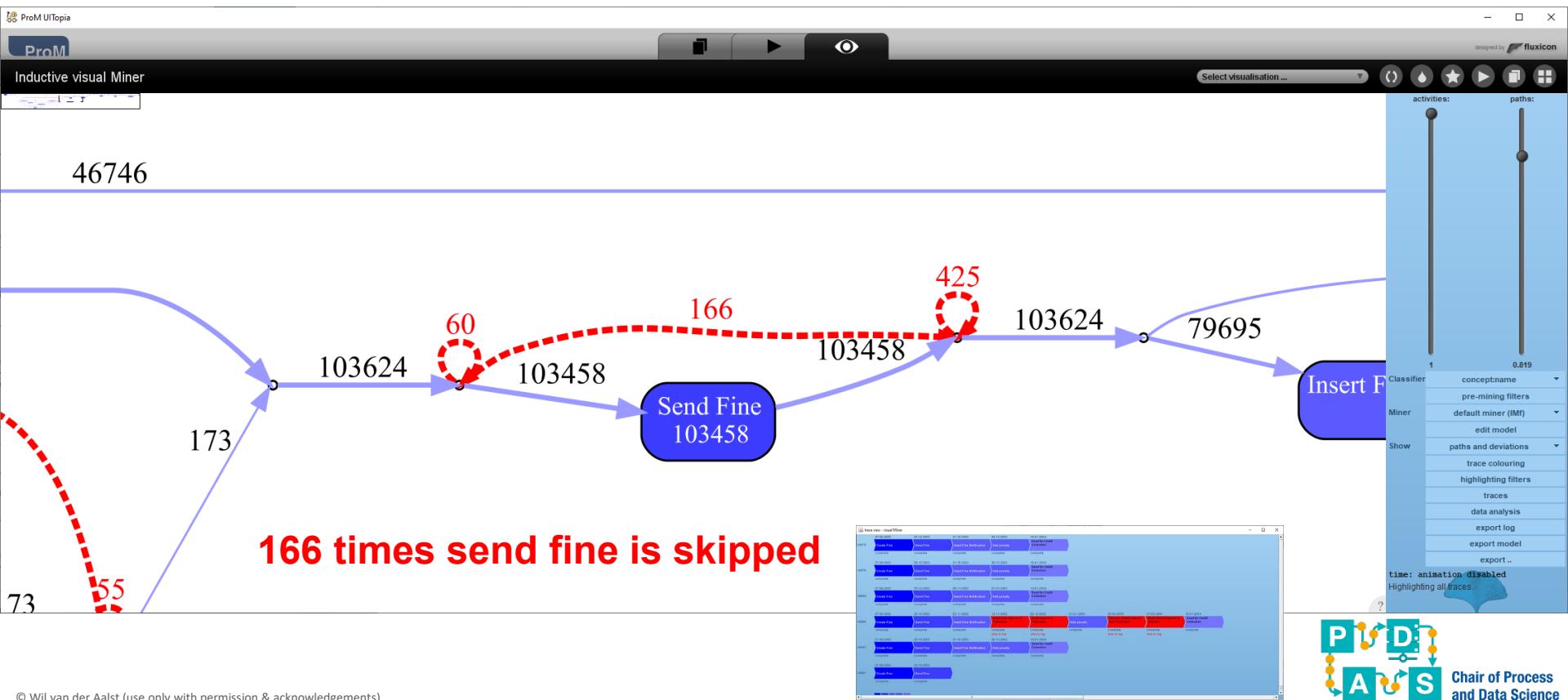
Chair of Process
and Data Science

Example: Inductive Visual Miner

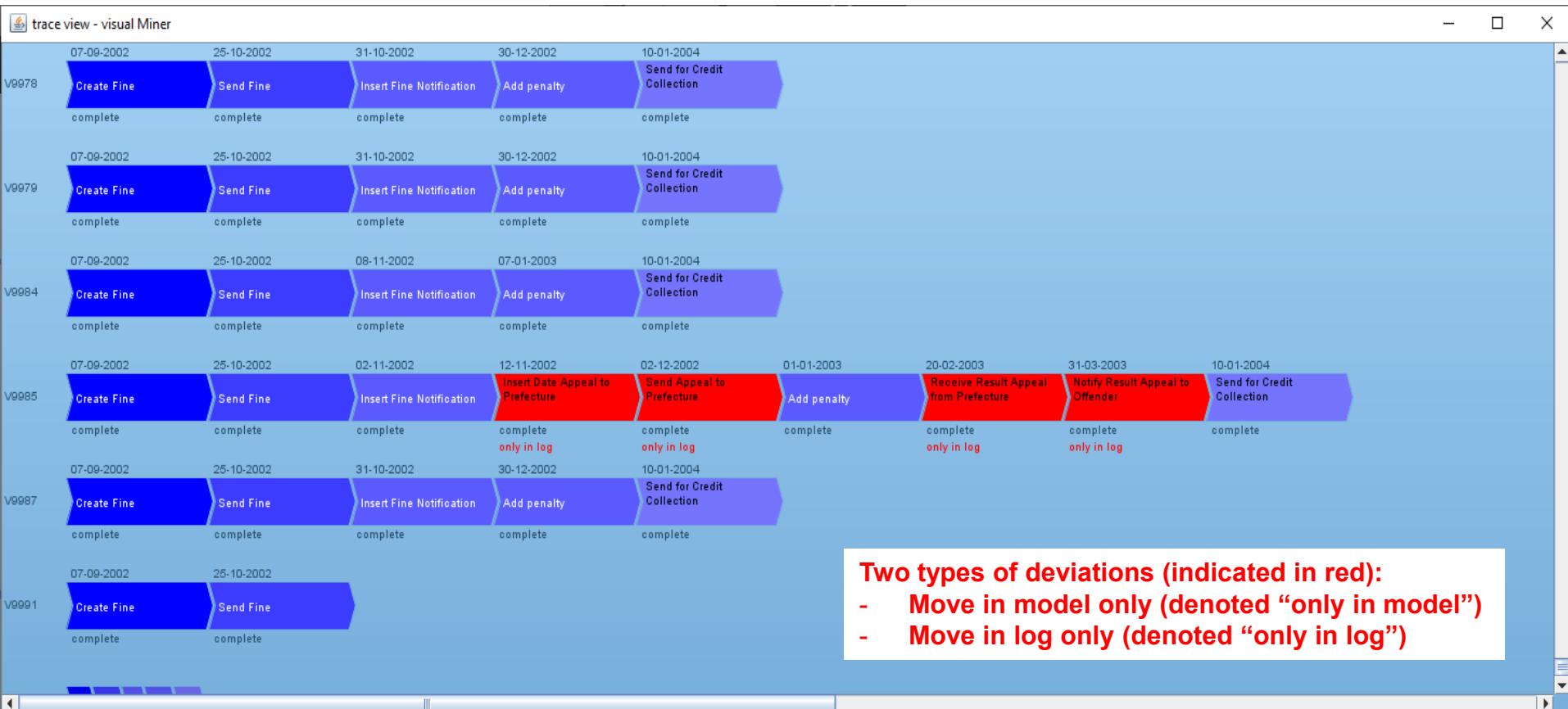


Based on alignments!

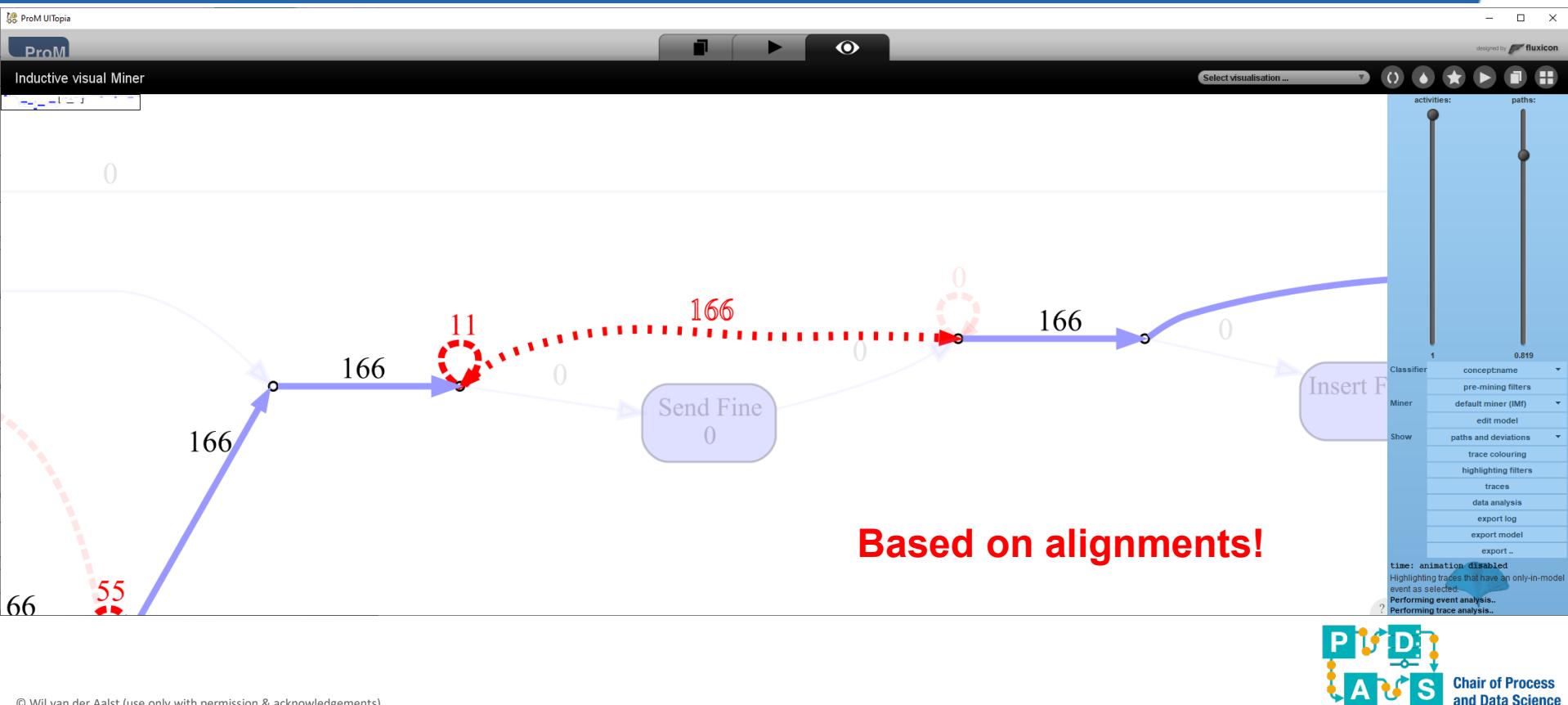
Example: Inductive Visual Miner



Example: Inductive Visual Miner



Example: Inductive Visual Miner



Example: Inductive Visual Miner

ProM UItopia designed by fluxicon

Inductive visual Miner

Select visualisation ...

activities: paths:

Classifier conceptname

Miner pre-mining filters

Show default miner (IMF) edit model paths and deviations trace colouring highlighting filters traces data analysis export log export model export ..

time: animation disabled

Highlighting traces that have an only-in-model event as selected.

Performing event analysis..

Performing trace analysis..

The screenshot shows the ProM UItopia interface with the 'Inductive visual Miner' tab selected. On the left, there is a process flow diagram with various nodes and transitions, some of which are highlighted in red. In the center, a table displays 166 traces, each consisting of a sequence of events with their start and end times. The traces are color-coded, and some events are highlighted in red. On the right, there are several filter and configuration options, including dropdown menus for 'activities' and 'paths', and sections for 'Classifier', 'Miner', and 'Show'. At the bottom, there are buttons for 'Event analysis' and 'Trace analysis'.

166 traces that have this deviation

Example: Inductive Visual Miner



Part I: Introduction

Chapter 1
Data Science in Action

Chapter 2
Process Mining:
The Missing Link

Part II: Preliminaries

Chapter 3
Process Modeling
and Analysis

Chapter 4
Data Mining

Part III: From Event Logs to Process Models

Chapter 5
Getting the Data

Chapter 6
Process Discovery:
An Introduction

Chapter 7
Advanced Process
Discovery Techniques

Part IV: Beyond Process Discovery

Chapter 8
Conformance
Checking

Chapter 9
Mining Additional
Perspectives

Chapter 10
Operational Support

Part V: Putting Process Mining to Work

Chapter 11
Process Mining
Software

Chapter 12
Process Mining in the
Large

Chapter 13
Analyzing “Lasagna
Processes”

Chapter 14
Analyzing “Spaghetti
Processes”

Part VI: Reflection

Chapter 15
Cartography and
Navigation

Chapter 16
Epilogue



ID	Topic	Date	Date	Place
	Lecture 1 Introduction to Process Mining	08.04.24	Monday	AH V
	Lecture 2 Data Science: Supervised Learning	09.04.24	Tuesday	AH V
	<i>Exercise 1 Tool Introduction</i>	09.04.24	Tuesday	AH III
	Lecture 3 Data Science: Unsupervised Learning and Evaluation	15.04.24	Monday	AH V
	Lecture 4 Introduction to Process Discovery	16.04.24	Tuesday	AH V
	<i>Exercise 2 Data Mining</i>	16.04.24	Tuesday	AH III
	Lecture 5 Alpha Algorithm 1	22.04.24	Monday	AH V
	Lecture 6 Alpha Algorithm 2	23.04.24	Tuesday	AH V
	<i>Exercise 3 Petri Nets</i>	23.04.24	Tuesday	AH III
	Lecture 7 Model Quality Representation	29.04.24	Monday	AH V
	Lecture 8 Heuristic Mining	30.04.24	Tuesday	AH V
	<i>Exercise 4 Alpha Miner</i>	30.04.24	Tuesday	AH III
	Lecture 9 Region-Based Mining	06.05.24	Monday	AH V
	<i>Exercise 5 Heuristic Mining and Region-Based Mining</i>	07.05.24	Tuesday	AH III
	Lecture 10 Inductive Mining	13.05.24	Monday	AH V
	Lecture 11 Event Data and Exploration	14.05.24	Tuesday	AH V
	<i>Exercise 6 Inductive Mining</i>	14.05.24	Tuesday	AH III
	Lecture 12 Conformance Checking 1	27.05.24	Monday	AH V
	Lecture 13 Conformance Checking 2	28.05.24	Tuesday	AH V
	<i>Q&A Session Assignment Part I</i>	28.05.24	Tuesday	AH III
	Deadline Assignment Part I	02.06.24	Sunday	
	<i>Exercise 7 Footprint and Token-Based Replay (Exercise)</i>	03.06.24	Monday	AH V
	<i>Exercise 8 Alignments (Exercise)</i>	04.06.24	Tuesday	AH V
	Lecture 14 Decision Mining	10.06.24	Monday	AH V
	<i>Lecture 15 Celonis Guest Lecture</i>	11.06.24	Tuesday	AH V
	<i>Exercise 9 Decision Mining</i>	11.06.24	Tuesday	AH III
	Lecture 16 Performance Analysis and Organizational Mining	17.06.24	Monday	AH V
	<i>Exercise 10 Performance Analysis (Exercise)</i>	18.06.24	Tuesday	AH V
	<i>Exercise 11 Organizational Mining</i>	18.06.24	Tuesday	AH III
	<i>Exercise 12 Celonis Case Study</i>	24.06.24	Monday	AH V
	Lecture 17 Operational Support and Process Mining Applications	01.07.24	Monday	AH V
	Lecture 18 Distributed, Streaming, and Comparative Process Mining	02.07.24	Tuesday	AH V
	<i>Exercise 13 Operational Process Mining</i>	02.07.24	Tuesday	AH III
	Lecture 19 Closing	08.07.24	Monday	AH V
	<i>Q&A Session Assignment Part II</i>	09.07.24	Tuesday	AH III
	Deadline Assignment Part II	14.07.24	Sunday	
	<i>Q&A Session Exam</i>	16.07.24	Tuesday	AH III



Conformance Checking (2/2)

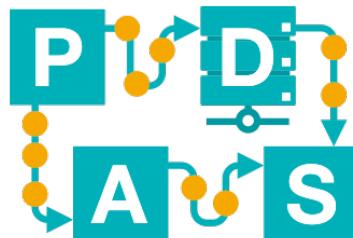
Lecture 13

prof.dr.ir. Wil van der Aalst

www.vdaalst.com @wvdaalst

www.pads.rwth-aachen.de

BPI-L13



Chair of Process
and Data Science

RWTHAACHEN
UNIVERSITY

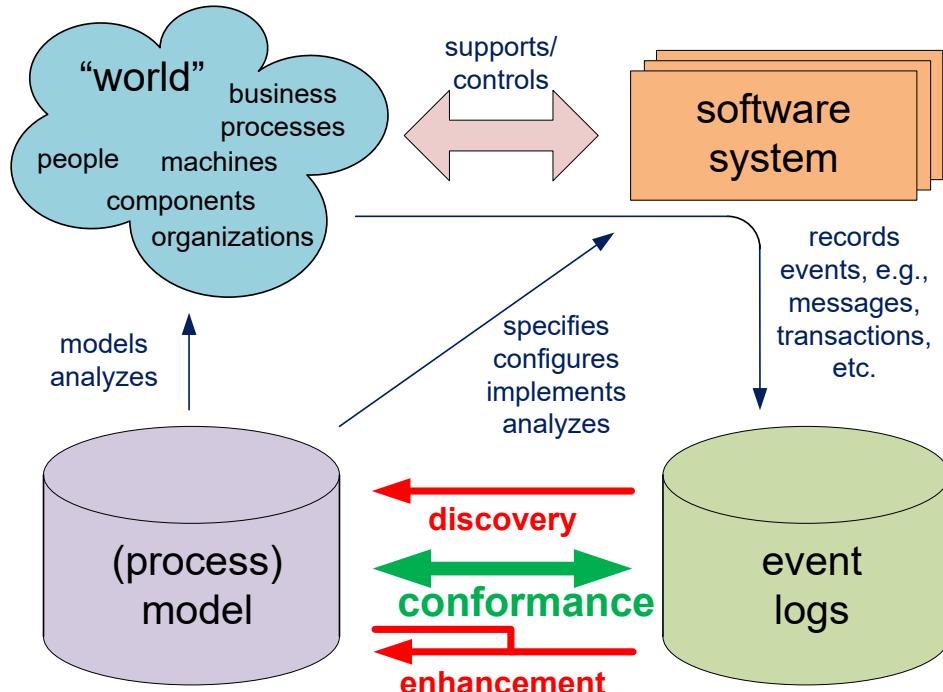
Outline

- A few end-to-end examples using token-based replay.
- Alignment-based conformance checking.
- Tool support for conformance checking.
- Applications of process mining.

Token-based replay revisited



Conformance checking



1. Conformance checking using causal footprints.
2. Conformance checking based on **token-based replay**.
3. Alignment-based conformance checking.

Last lecture: Token-based replay

#	trace
455	acdeh
191	abdeg
177	adceh
144	abdeh
111	acdeg
82	adceg
56	adbeh
47	acdefdfbeh
38	adbeg
33	acdefbdeh
14	acdefbddeg
11	acdefdbeg
9	adcefcdgeh
8	adcfdbeh
5	adcefbddeg
3	adcfbddefdbeg
2	adcefdbeg
2	adcefbddefbddeg
1	adcefdbefbdbeh
1	adbefbddefdbeg
1	adcfdbefcdefdbeg
1391	

missing tokens

?

consumed tokens

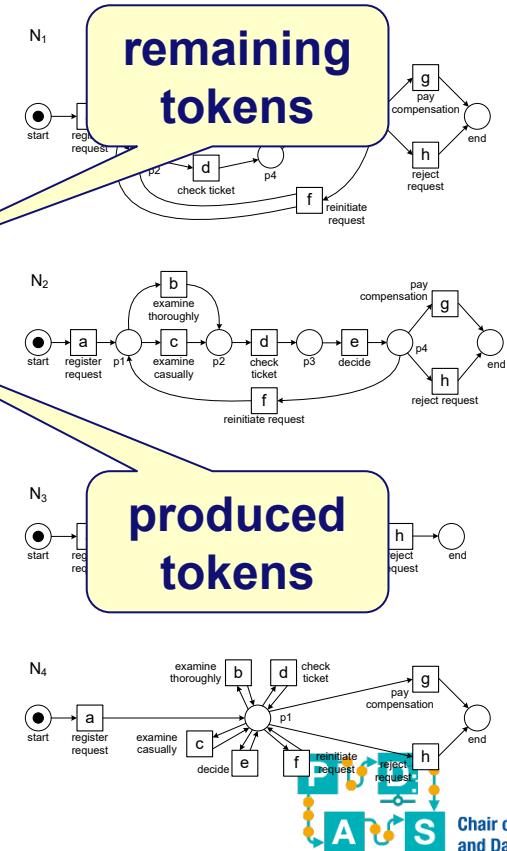
$$fitness(L, N) = \frac{1}{2} \left(1 - \frac{\sum_{\sigma \in L} L(\sigma) \times m_{N,\sigma}}{\sum_{\sigma \in L} L(\sigma) \times c_{N,\sigma}} \right) + \frac{1}{2} \left(1 - \frac{\sum_{\sigma \in L} L(\sigma) \times r_{N,\sigma}}{\sum_{\sigma \in L} L(\sigma) \times p_{N,\sigma}} \right)$$

$$ss(L_{full}, N_1) = 1$$

$$ss(L_{full}, N_2) = 0.9504$$

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

$$fitness(L_{full}, N_4) = 1$$



End-to-end examples

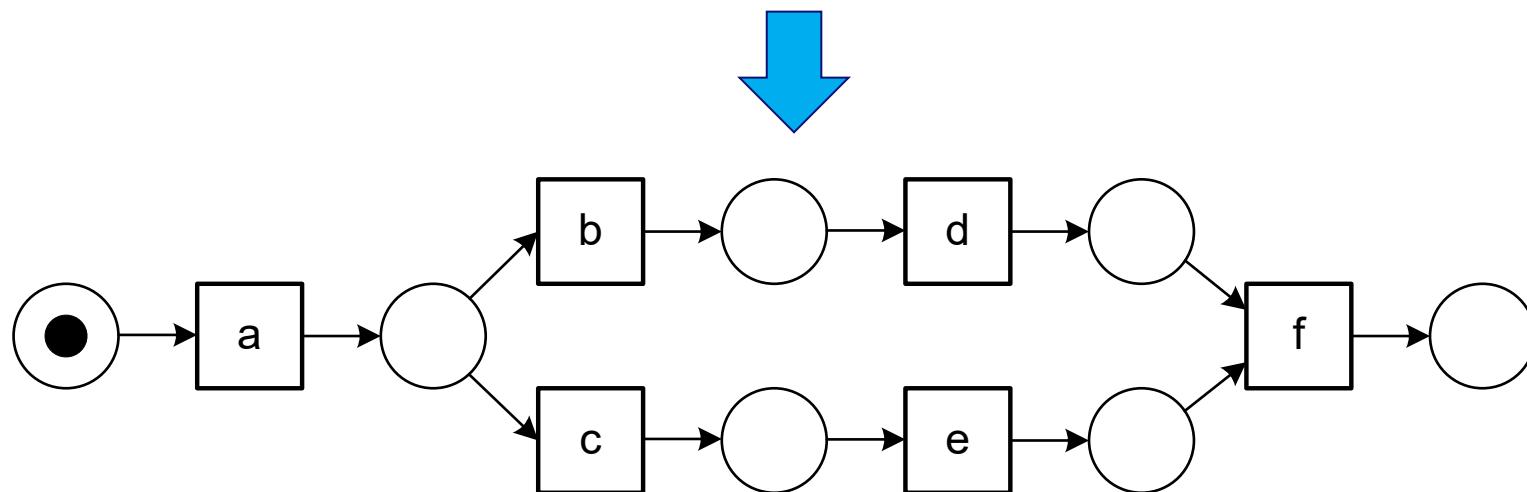
Question

$$L = [\langle a, b, d, e, f \rangle^{10}, \langle a, c, e, d, f \rangle^{10}]$$

- Consider the above event log.
- Give the model that the Alpha algorithm generates.
- Compute fitness using missing and remaining tokens.
- Comment on the findings.

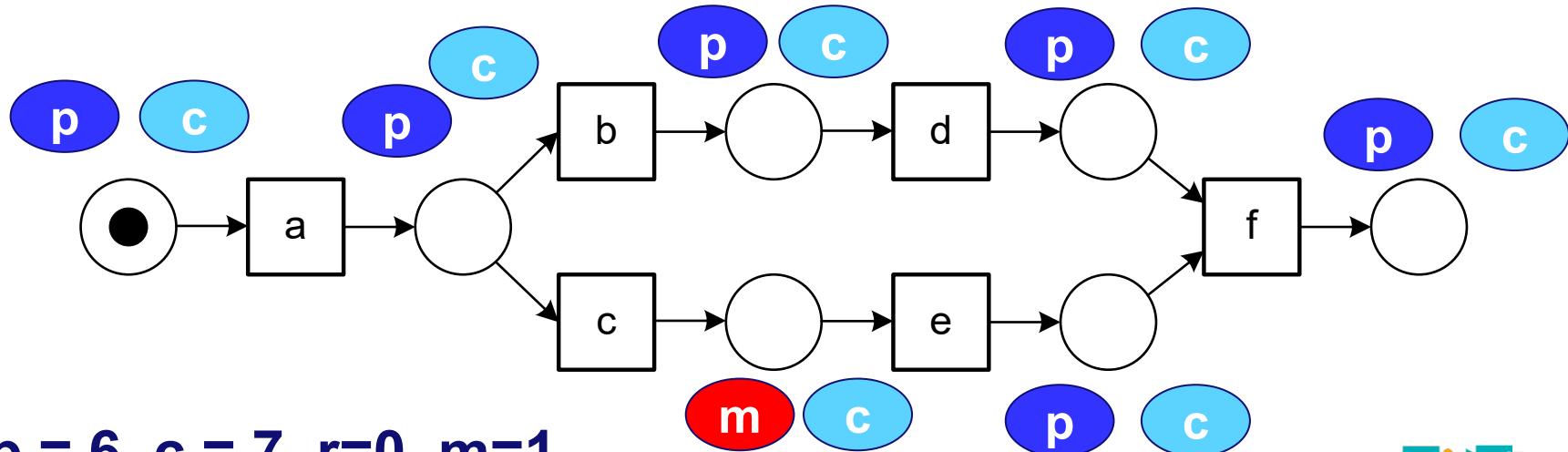
Model generated by the Alpha Algorithm

$$L = [\langle a, b, d, e, f \rangle^{10}, \langle a, c, e, d, f \rangle^{10}]$$



Trace abdef

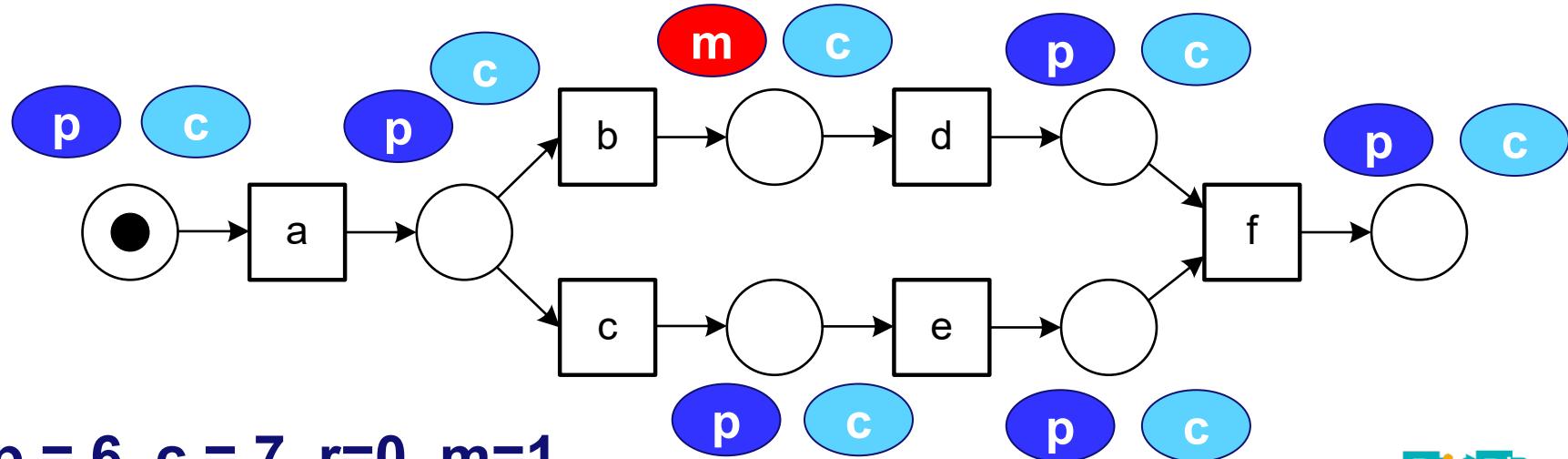
$$L = [\langle a, b, d, e, f \rangle^{10}, \langle a, c, e, d, f \rangle^{10}]$$



$p = 6, c = 7, r=0, m=1$

Trace acedf

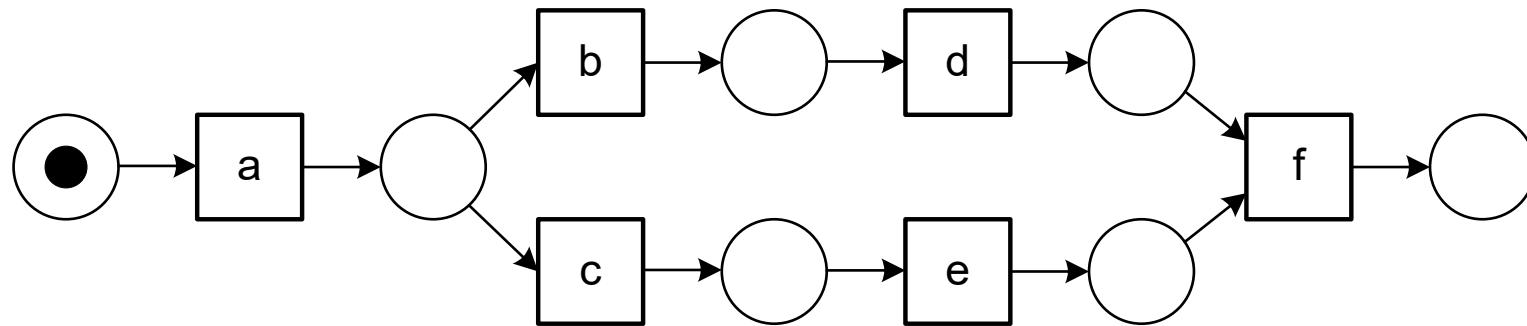
$$L = [\langle a, b, d, e, f \rangle^{10}, \boxed{\langle a, c, e, d, f \rangle^{10}}]$$



Overall fitness

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

$$L = [\langle a, b, d, e, f \rangle^{10}, \langle a, c, e, d, f \rangle^{10}]$$

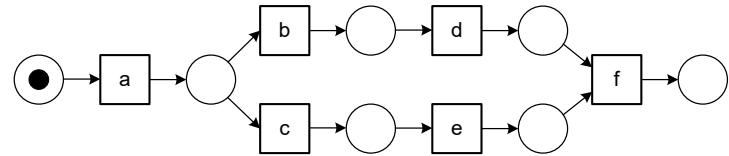


$$p = 2 * 10 * 6 = 120, c = 2 * 10 * 7 = 140, r = 2 * 10 * 0 = 0, m = 2 * 10 * 1 = 20$$

$$\frac{1}{2} \left(1 - \frac{20}{140} \right) + \frac{1}{2} \left(1 - \frac{0}{120} \right) = \frac{13}{14} \approx 0.93$$

Model is not sound!

$$L = [\langle a, b, d, e, f \rangle^{10}, \langle a, c, e, d, f \rangle^{10}]$$



$$\frac{1}{2} \left(1 - \frac{20}{140} \right) + \frac{1}{2} \left(1 - \frac{0}{120} \right) = \frac{13}{14} \approx 0.93$$

- The model is not sound. Actually there is no firing sequence leading to the target marking!
- How to interpret the result? Therefore, we typically require “relaxed soundness”, i.e., there is at least one firing sequence leading to the target marking.

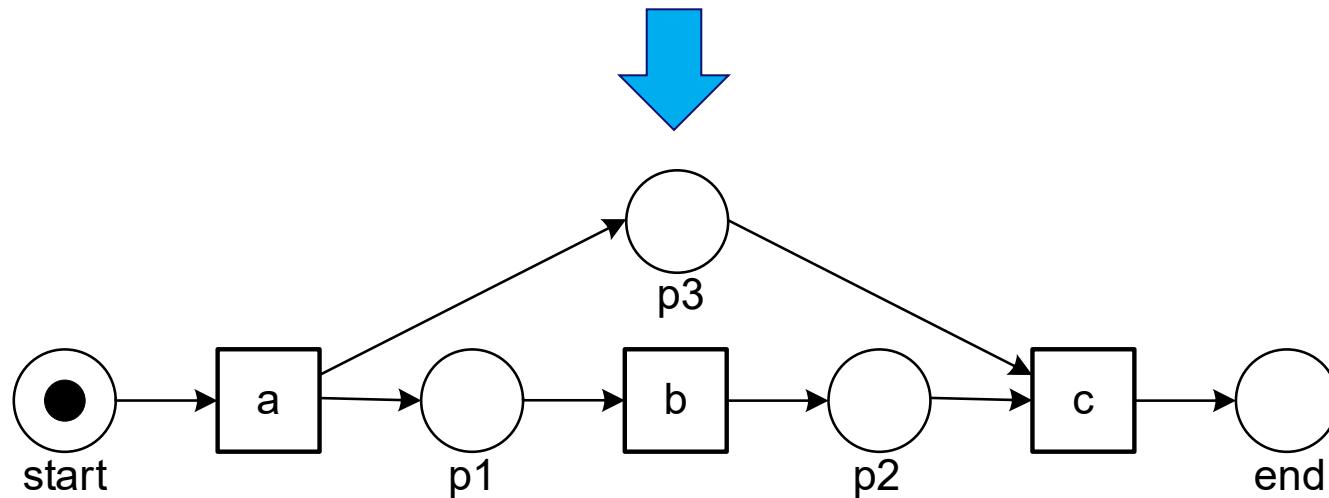
Question

$$L_{11} = [\langle a, b, c \rangle^{20}, \langle a, c \rangle^{30}]$$

- Consider the above event log.
- Give the model that the Alpha algorithm generates
- Compute fitness using missing and remaining tokens.
- Comment on the findings.

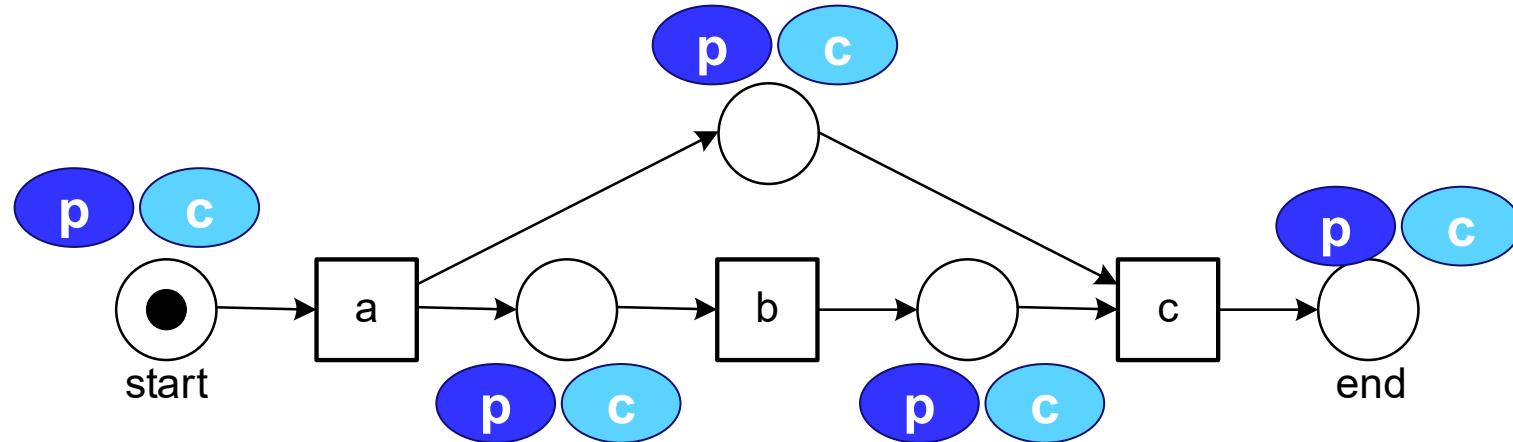
Model generated by the Alpha Algorithm

$$L_{11} = [\langle a, b, c \rangle^{20}, \langle a, c \rangle^{30}]$$



Trace abc

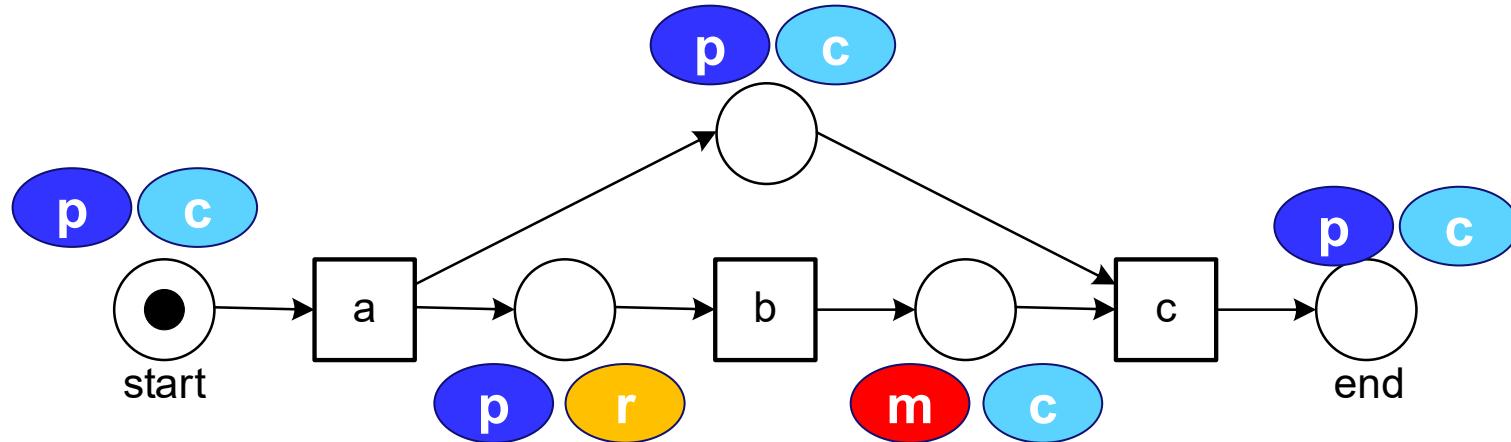
$$L_{11} = [\langle a, b, c \rangle^{20}, \langle a, c \rangle^{30}]$$



$p = 5, c = 5, r=0, m=0$

Trace ac

$$L_{11} = [\langle a, b, c \rangle^{20}, \boxed{\langle a, c \rangle^{30}}]$$



$p = 4, c = 4, r=1, m=1$

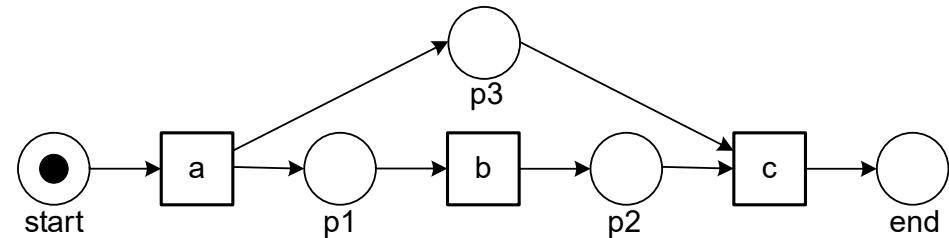
Overall fitness

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

$$L_{11} = [\langle a, b, c \rangle^{20}, \langle a, c \rangle^{30}]$$

p = 5, c = 5, r=0, m=0

p = 4, c = 4, r=1, m=1

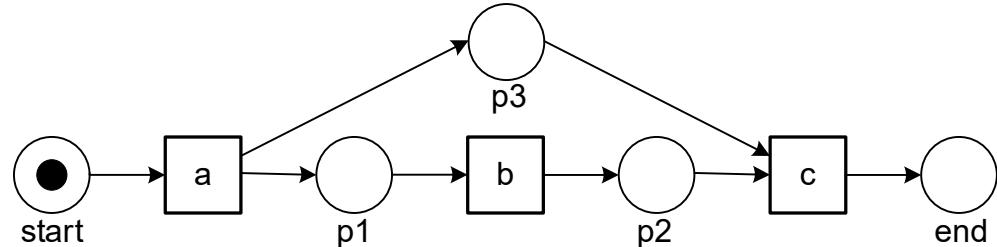


$$\begin{aligned} p &= 20*5+30*4 = 220, \quad c = 20*5+30*4=220, \\ r &= 20*0+30*1=30, \quad m=20*0+30*1=30 \end{aligned}$$

$$\frac{1}{2} \left(1 - \frac{30}{220} \right) + \frac{1}{2} \left(1 - \frac{30}{220} \right) = \frac{19}{22} \approx 0.86$$

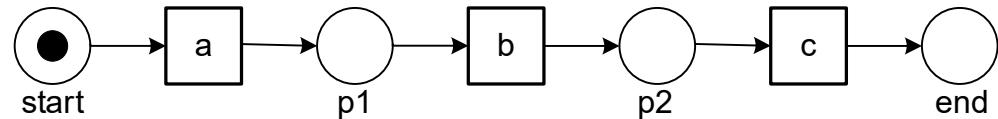
Redundant places impact fitness (1/2)

$$L_{11} = [\langle a, b, c \rangle^{20}, \langle a, c \rangle^{30}]$$



$$\frac{1}{2} \left(1 - \frac{30}{220} \right) + \frac{1}{2} \left(1 - \frac{30}{220} \right) = \frac{19}{22} \approx 0.86$$

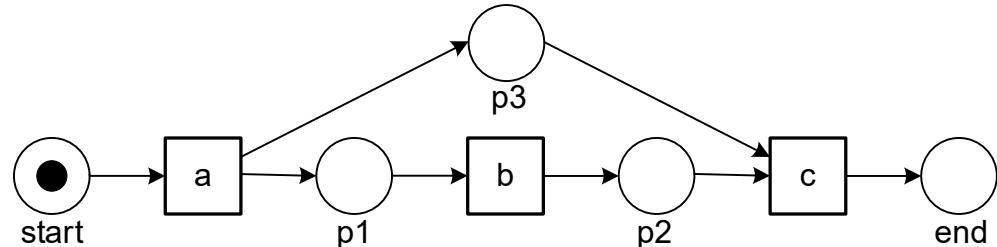
$$L_{11} = [\langle a, b, c \rangle^{20}, \langle a, c \rangle^{30}]$$



$$\frac{1}{2} \left(1 - \frac{30}{170} \right) + \frac{1}{2} \left(1 - \frac{30}{170} \right) = \frac{14}{17} \approx 0.82$$

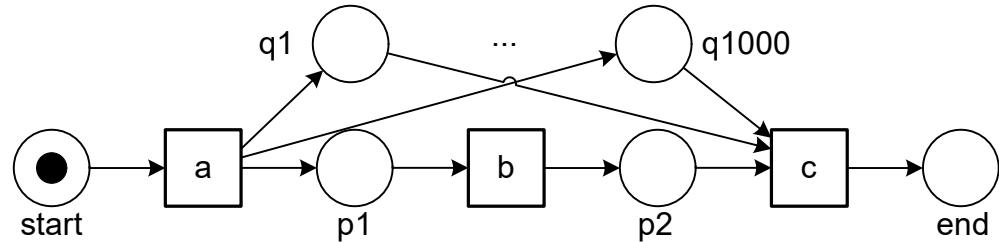
Redundant places impact fitness (2/2)

$$L_{11} = [\langle a, b, c \rangle^{20}, \langle a, c \rangle^{30}]$$



$$\frac{1}{2} \left(1 - \frac{30}{220} \right) + \frac{1}{2} \left(1 - \frac{30}{220} \right) = \frac{19}{22} \approx 0.86$$

$$L_{11} = [\langle a, b, c \rangle^{20}, \langle a, c \rangle^{30}]$$

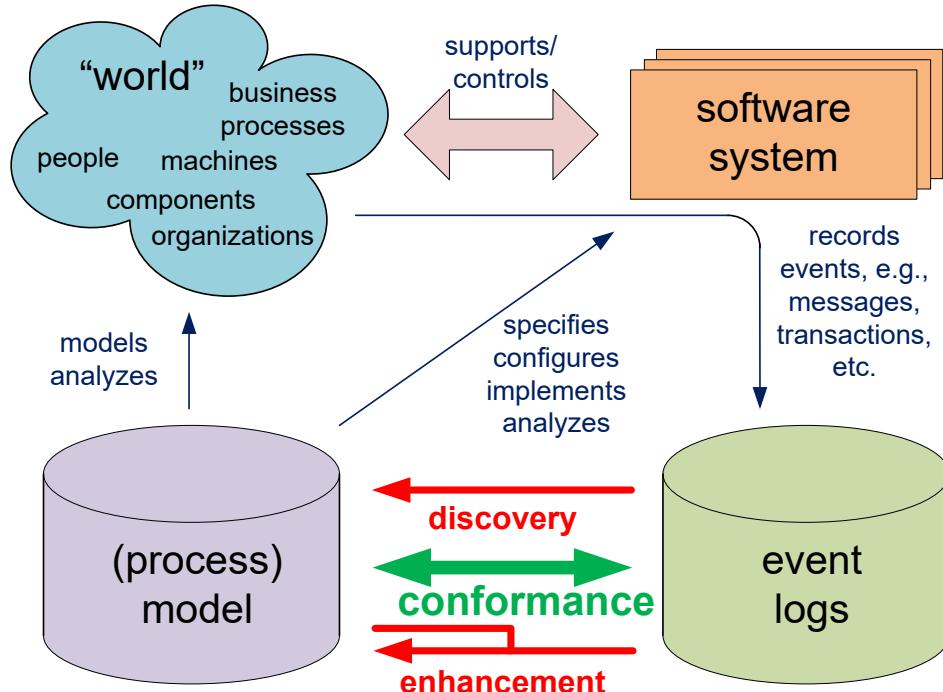


$$\frac{1}{2} \left(1 - \frac{30}{50170} \right) + \frac{1}{2} \left(1 - \frac{30}{50170} \right) = \frac{5014}{5017} \approx 0.999$$

Aligning Observed and Modeled Behavior



Conformance checking

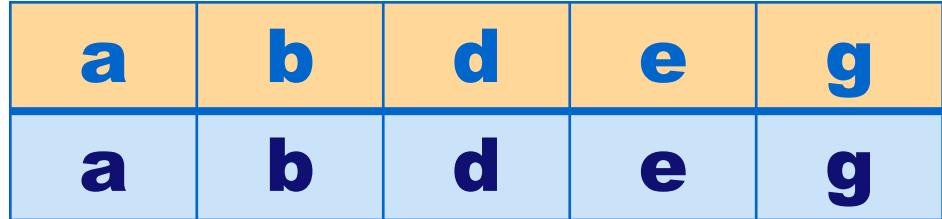


1. Conformance checking using causal footprints.
2. Conformance checking based on token-based replay.
3. Alignment-based conformance checking.

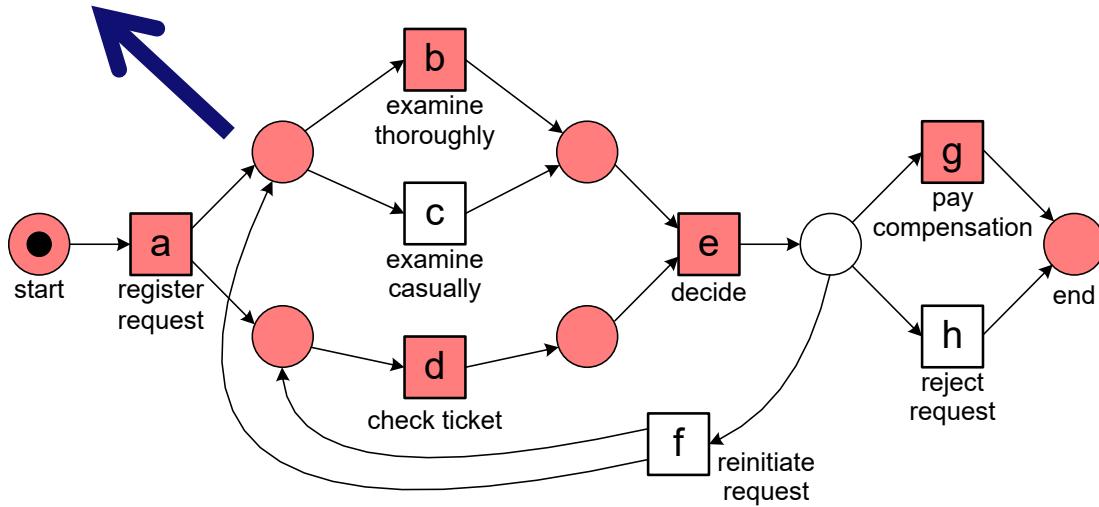
Requirements

- Conformance checking should **not** impose restrictions on the process notation (e.g., silent transitions and two transitions with the same label should be possible).
- Two **semantically equivalent** models should have the same conformance value.
- Should provide a "**closest matching path**" through the process model for any trace in the event log.
 - Also required for **performance analysis!**
 - **Beyond** the analysis of replay fitness (advanced diagnostics, precision, generalization, etc.).

Alignments



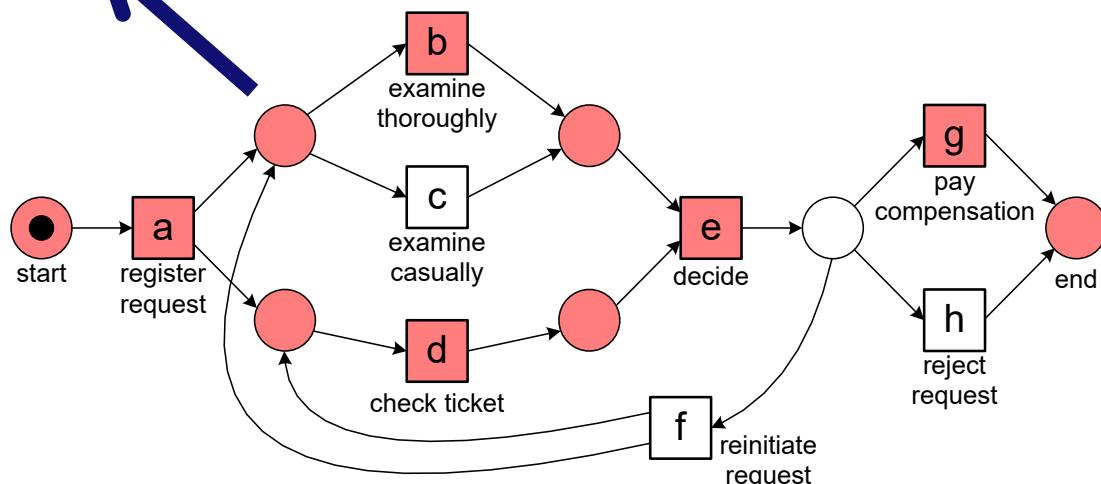
$\langle a, b, d, e, g \rangle$



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

Alignments

a	»	d	e	g	h
a	b	d	e	g	»



Terminology

alignment
(sequence of moves)

move in log only

move in model

move in both

a			d	e	g
	c	d	e	g	

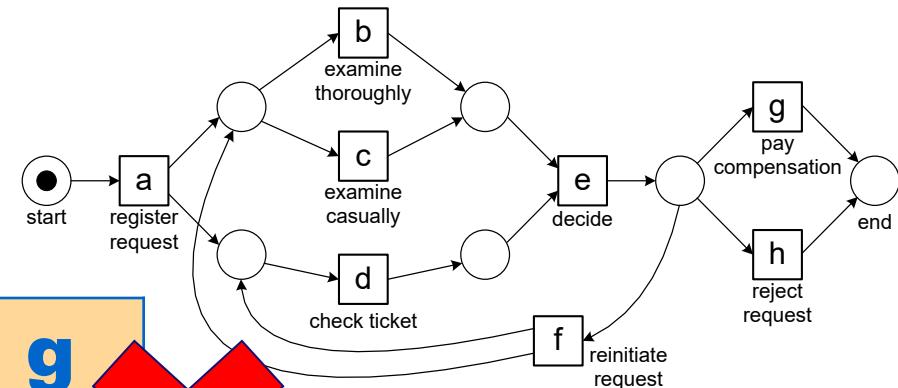
- Independent of process model notation!
- Projection on top row (remove "no moves") corresponds to the **run** in the event log.
 - Projection on bottom row (remove "no moves") corresponds to a **run of the model**.

Optimal alignment for $\langle a,b,d,e,g \rangle$

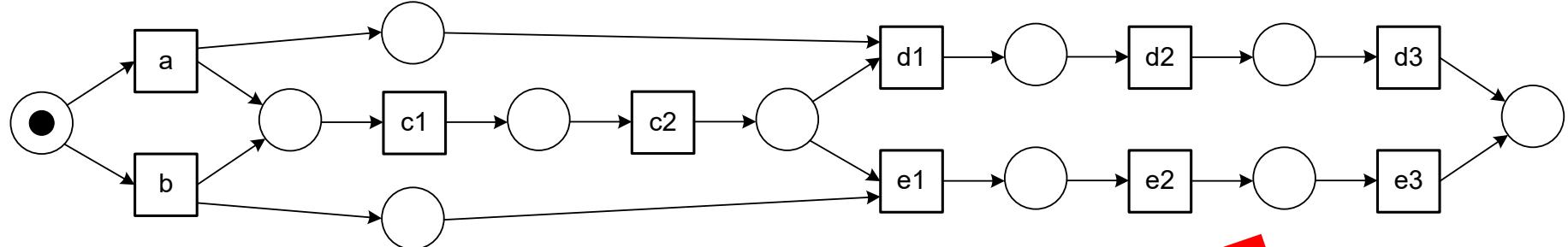
a	b	d	e	g
a	b	d	e	g

a	b	»	d	e	g
a	»	c	d	e	g

a	b	d	e	g	»	»	»	»	»	»
»	»	»	»	»	a	c	d	e	g	»



Optimal alignment for $\langle a, c_1, c_2, e_1, e_2, e_3 \rangle$?

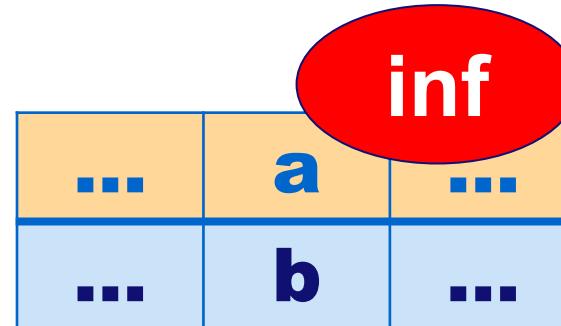
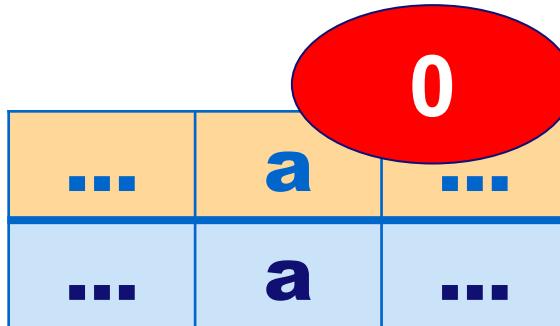
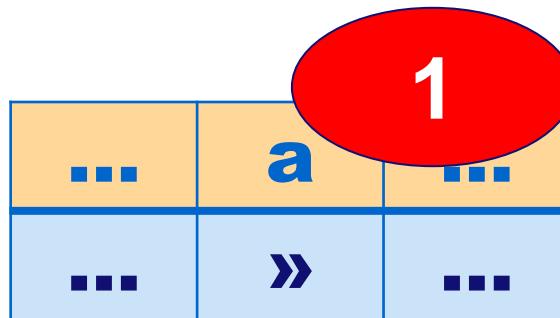


a	»	c1	c2	e1	e2	e3
»	b	c1	c2	e1		

a	c1	»	»	e1	e2	e3
a	c	d1	d2	d3	»	»

Depends on cost function!

Standard cost function count »'s in alignment



Using the standard cost function

optimal

there is no other alignment
that has lower costs

a	b	d	e	g
a	b	d	e	g

0

0

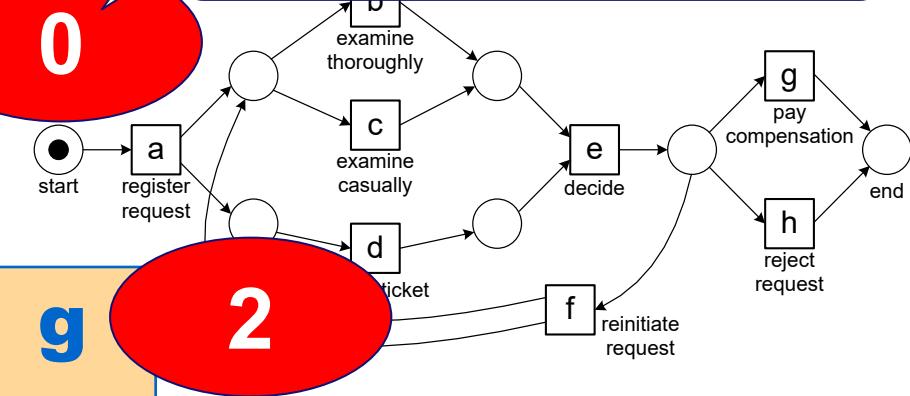
0

0

0

a	b	»	d	e	g
a	»	c	d	e	g

2



a	b	d	e	g	»	»	»	»	»	»
»	»	»	»	»	a	c	d	e	g	D

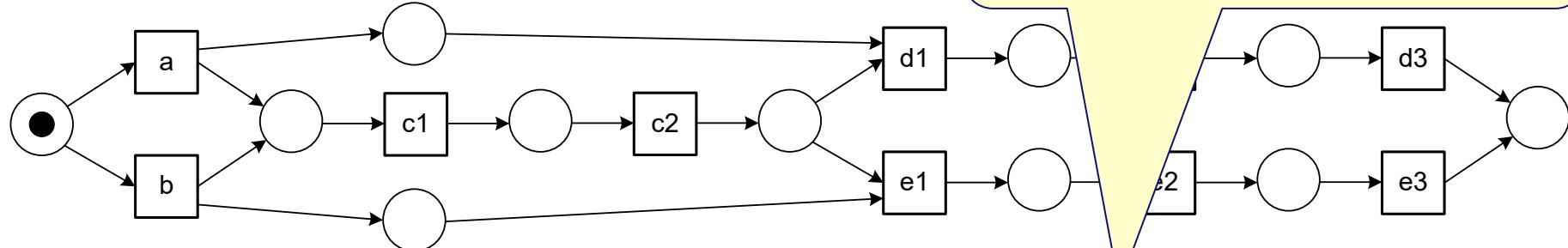
10

Using the standard cost function

trace in log: $\langle a, c1, c2, e1, e2, e3 \rangle$

optimal

there is no other alignment
that has lower costs



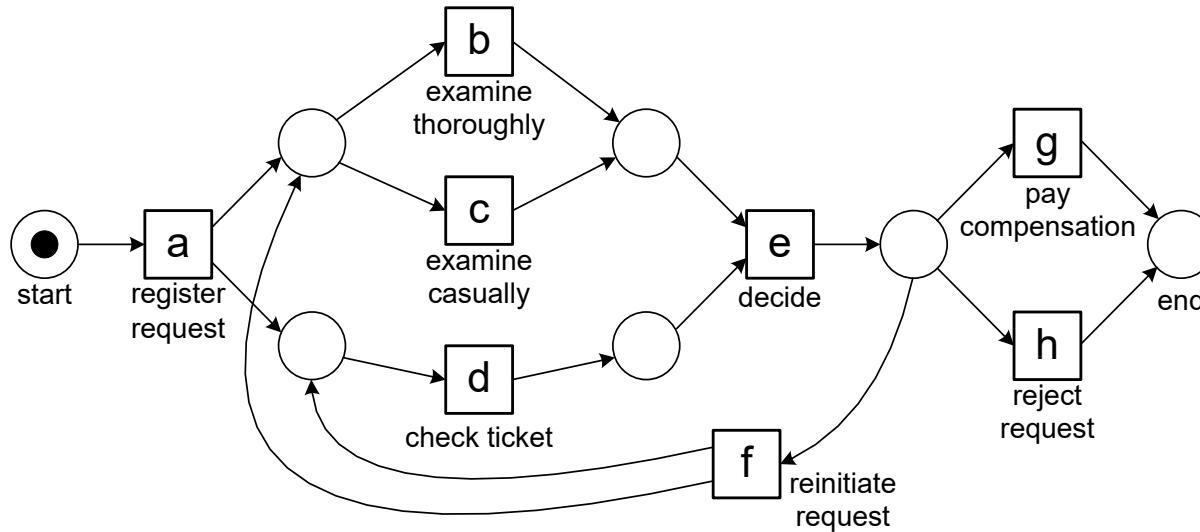
a	»	c1	c2	e1	e2	e3
»	b	c1	c2	e1	e2	e3

a	c1	c2	»	»	»	e1	e2	e3
a	c1	c2	d1	d2	d3	»	»	»

6

Optimal alignment for $\langle a,b,e,f,d,e,g \rangle$ (1/2)

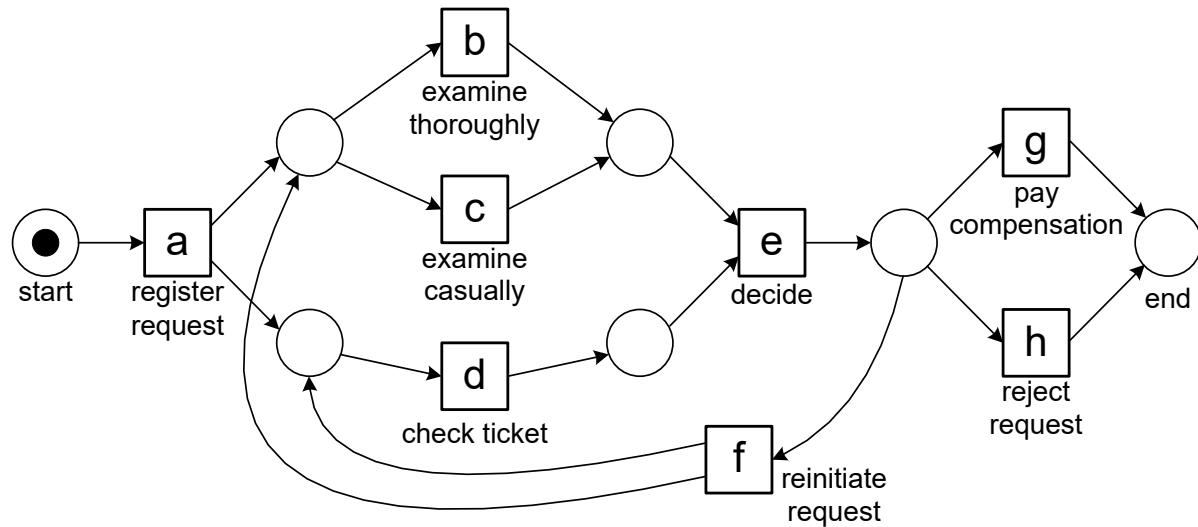
a	b	e	f	d	e	g	
a	b	»	»	d	e	g	2



loop is not taken: e and f in event log are discarded

Optimal alignment for $\langle a,b,e,f,d,e,g \rangle$ (2/2)

a	b	»	e	f	d	»	e	g	2
a	b	d	e	f	d	b	e	g	



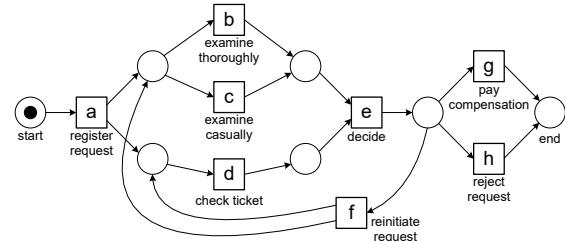
loop is taken:
d and b are
missing in
event log

Not one unique optimal alignment

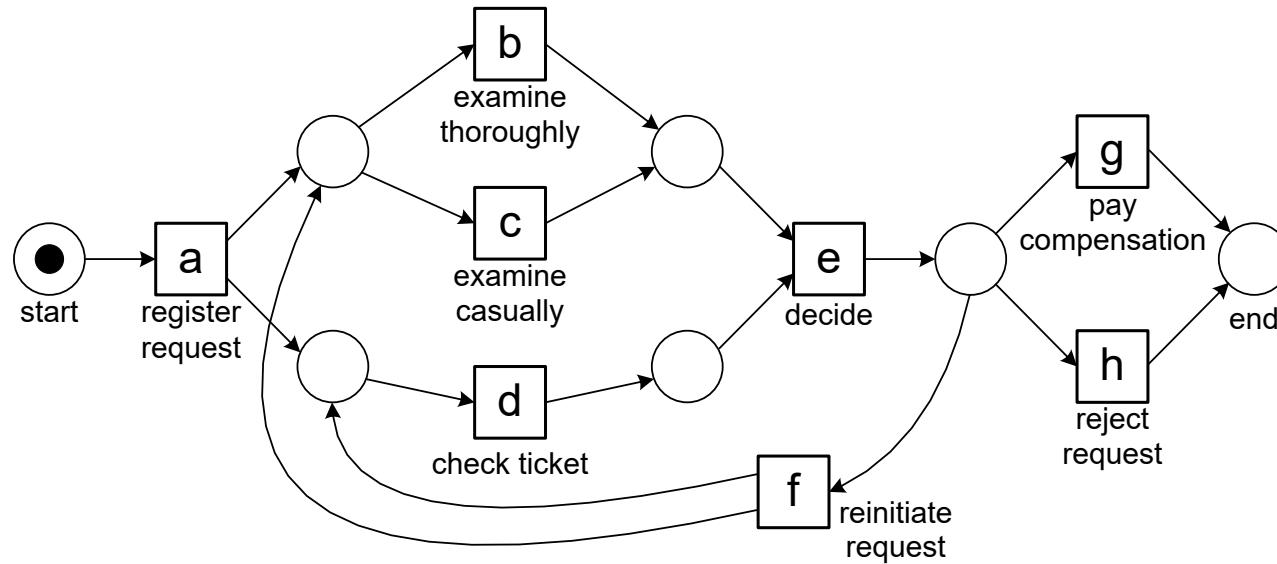
a	b	»	e	f	d	»	e	g	2
a	b	d	e	f	d	b	e	g	

a	b	e	f	d	e	g	2
a	b	»	»	d	e	g	

...



Question: How many optimal alignments are there for $\langle a,b,e,f,d,e,g \rangle$?



$\langle a,b,e,f,d,e,g \rangle$

Answer: 9

$1 + (2 \times 2) + (2 \times 2) = 9$ optimal alignments having cost 2

a	b	e	f	d	e	g
a	b	»	»	d	e	g

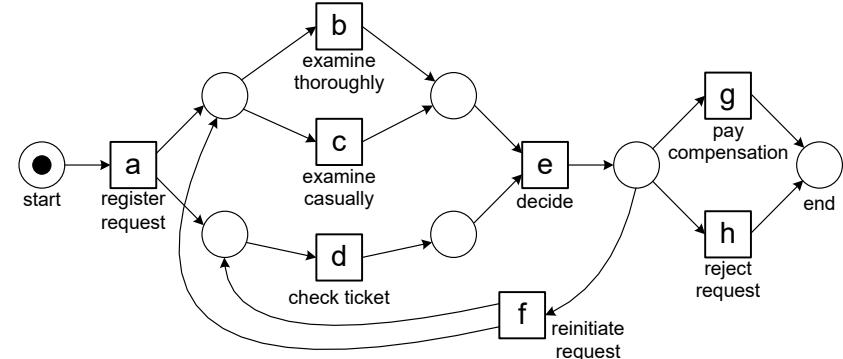
1x

move in model can be reordered in concurrent part



4x

a	b	»	e	f	d	»	e	g
a	b	d	e	f	d	b	e	g



4x

a	b	»	e	f	d	»	e	g
a	b	d	e	f	d	c	e	g

Any cost structure is possible

...	send-letter(John,4 weeks, \$400)	...
...	send-email(Sue,3 weeks,\$500)	...

Any cost structure is possible

...	send-letter(John,4 weeks, \$400)	...
...	send-email(Sue,3 weeks,\$500)	...

similar activities (lower costs for related activities)

Any cost structure is possible

...	send-letter(John,4 weeks, \$400)	...
...	send-email(Sue,3 weeks,\$50)	...

**resource-related conformance costs
(done by someone that does or does not have
the specified role)**



Any cost structure is possible

...	send-letter(John,4 weeks, \$400)	...
...	send-email(Sue,3 weeks,\$500)	...

**time-related conformance costs
(activity should happen within a preset deadline)**



Any cost structure is possible

...	send-letter(John,4 weeks, \$400)	...
...	send-email(Sue,3 weeks,\$500)	...

**data-related conformance costs
(routing condition is violated, e.g., path
only for more valuable orders)**

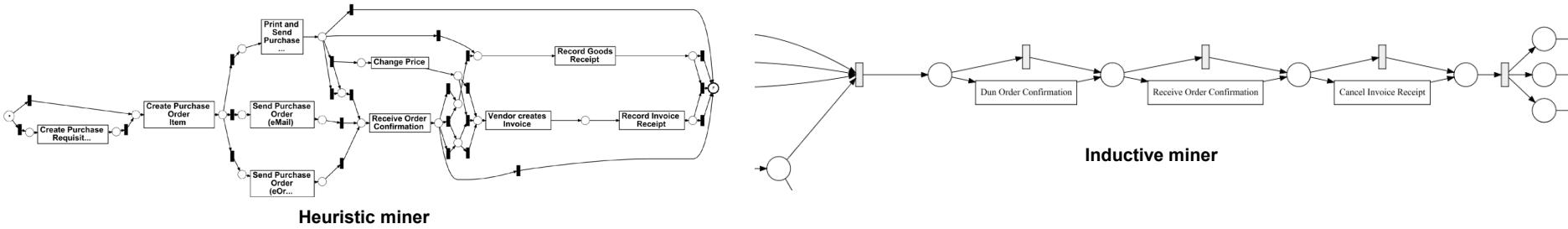


Any cost structure is possible

...	send-letter(John,4 weeks, \$400)	...
...	send-email(Sue,3 weeks,\$500)	...

risk-related conformance costs, context-dependent conformance costs, ...

Side note: Silent transitions

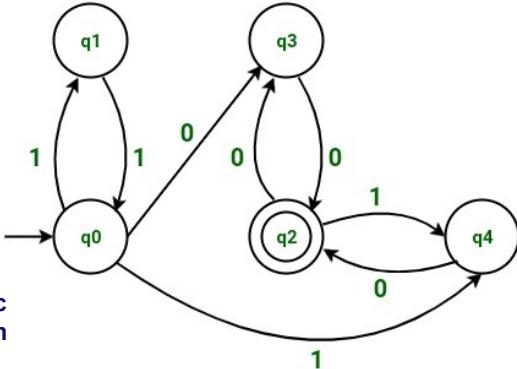
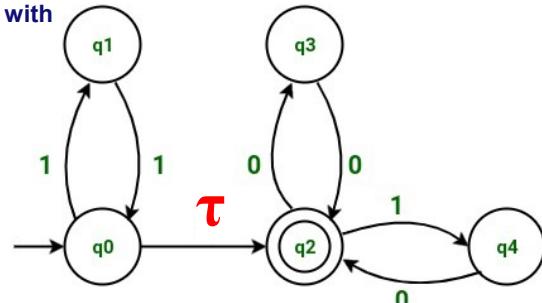


- Two options:
 - Remove these transitions when considering the model's behavior (Model reduction: any Nondeterministic Finite Automaton (NFA) with silent activities, can be translated into an NFA without silent activities, which, in turn, can be translated into a Deterministic Finite Automaton (DFA).)
 - The corresponding moves on model have cost 0 (or a very small epsilon in case of loops)

Example Transformations

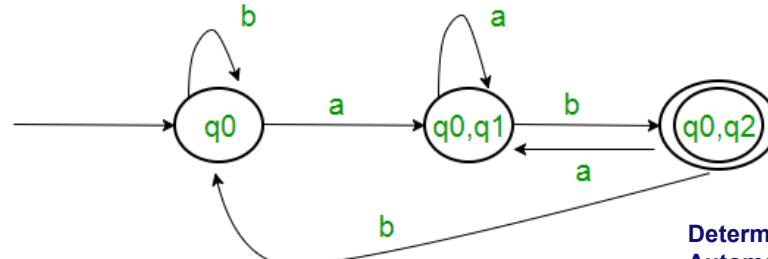
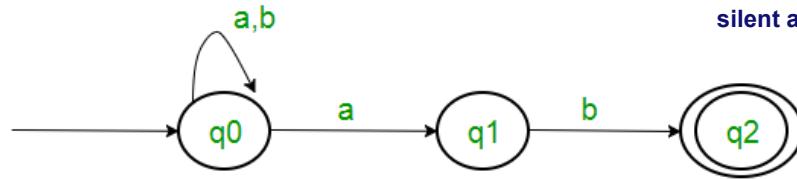
(outside scope course)

Nondeterministic Finite Automaton (NFA) with silent activities



Nondeterministic Finite Automaton NFA without silent activities

Nondeterministic Finite Automaton NFA without silent activities



Deterministic Finite Automaton (DFA)



Chair of Process and Data Science

Examples taken from <https://www.geeksforgeeks.org/>

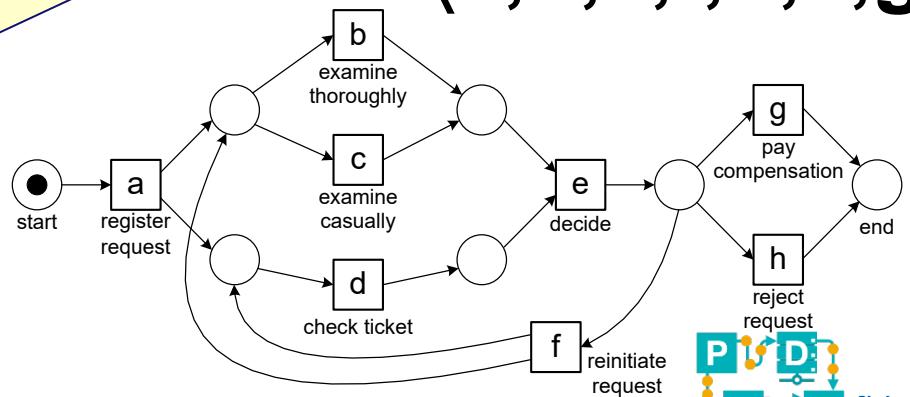
Computing fitness

Computing fitness
all events cause a move in log only
model.

$$1 - \frac{2}{7+5} = 0.833$$

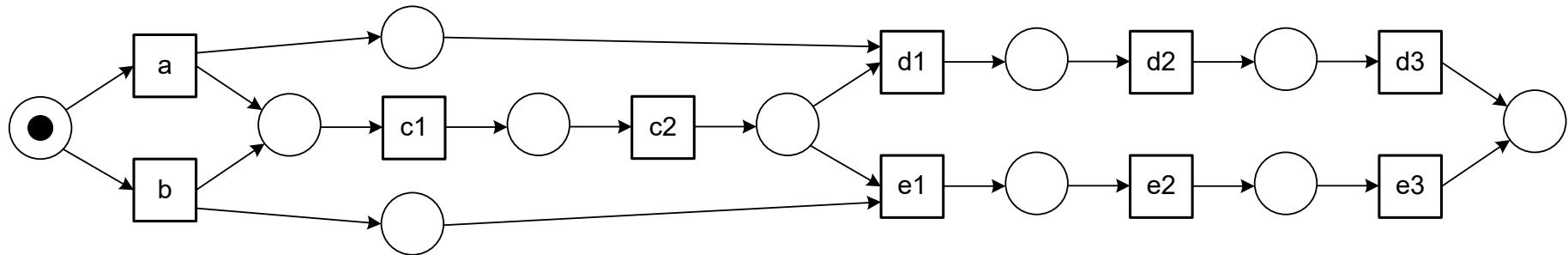
an optimal alignment with a
shortest path from initial state to final state

$\langle a,b,e,f,d,e,g \rangle$



Question: Compute alignment-based fitness

$\langle a, c1, c2, e1, e2, e3 \rangle$

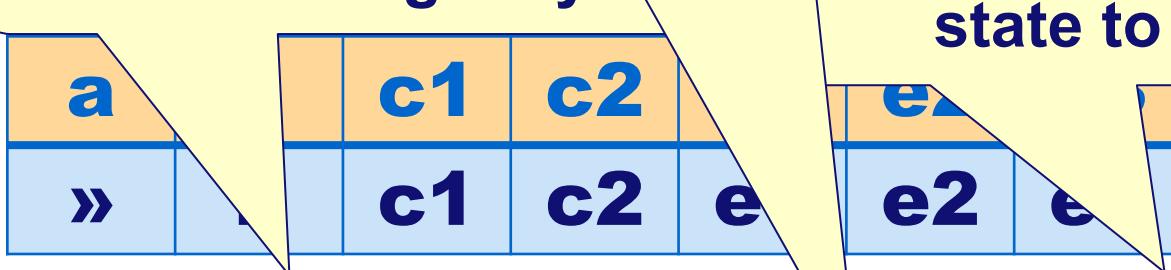


Answer

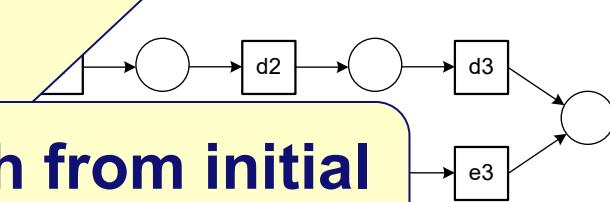
cost of optimal alignment = 2

$\langle a, c1, c2, e1, e2, e3 \rangle$

all even worst-case scenario
moves bring only



shortest path from initial state to final state

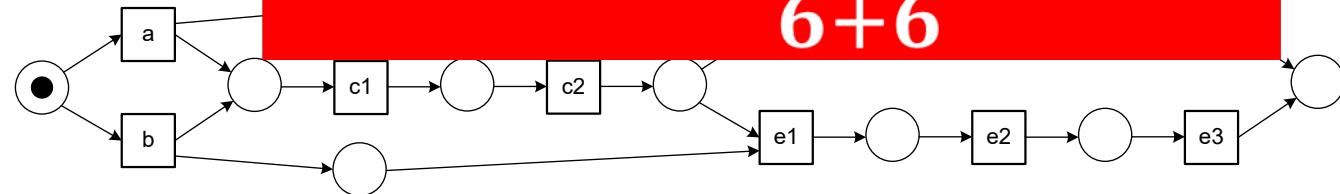


a	c1	c2	e1	e2	e3	»	»	»	»	»	»
»	»	»	»	»	»	a	c1	c2	d1	d2	d3

Answer

$\langle a, c1, c2, e1, e2, e3 \rangle$

$$\text{fitness} = 1 - \frac{2}{6+6} = 0.833$$

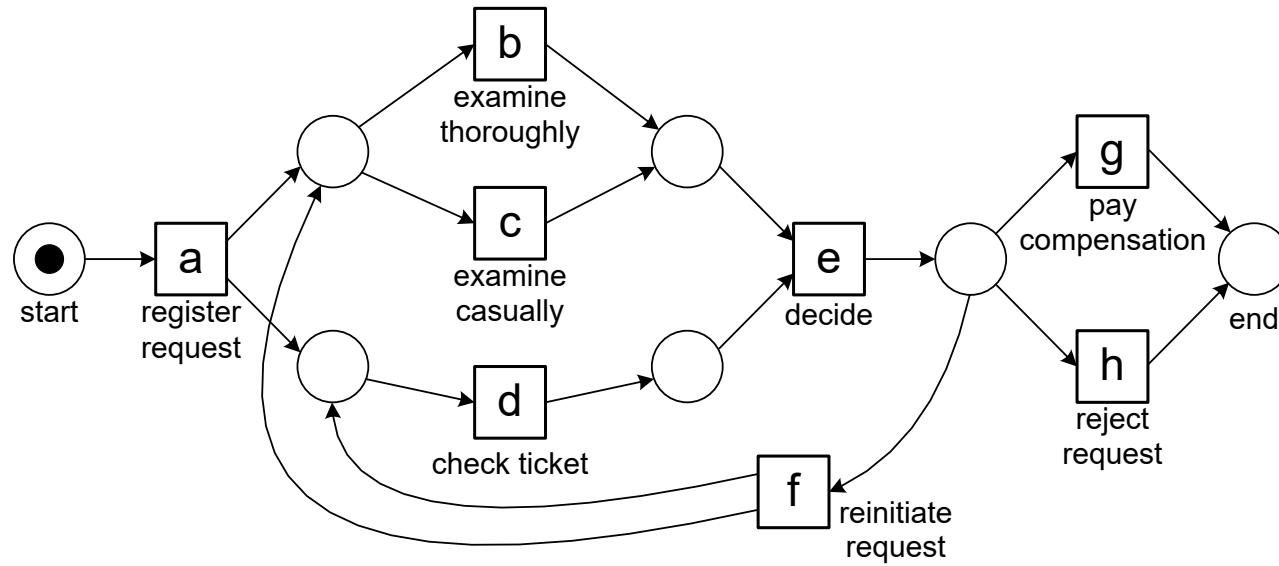


a	»	c1	c2	e1	e2	e3
»	b	c1	c2	e1	e2	e3

a	c1	c2	e1	e2	e3	»	»	»	»	»	»
»	»	»	»	»	»	a	c1	c2	d1	d2	d3

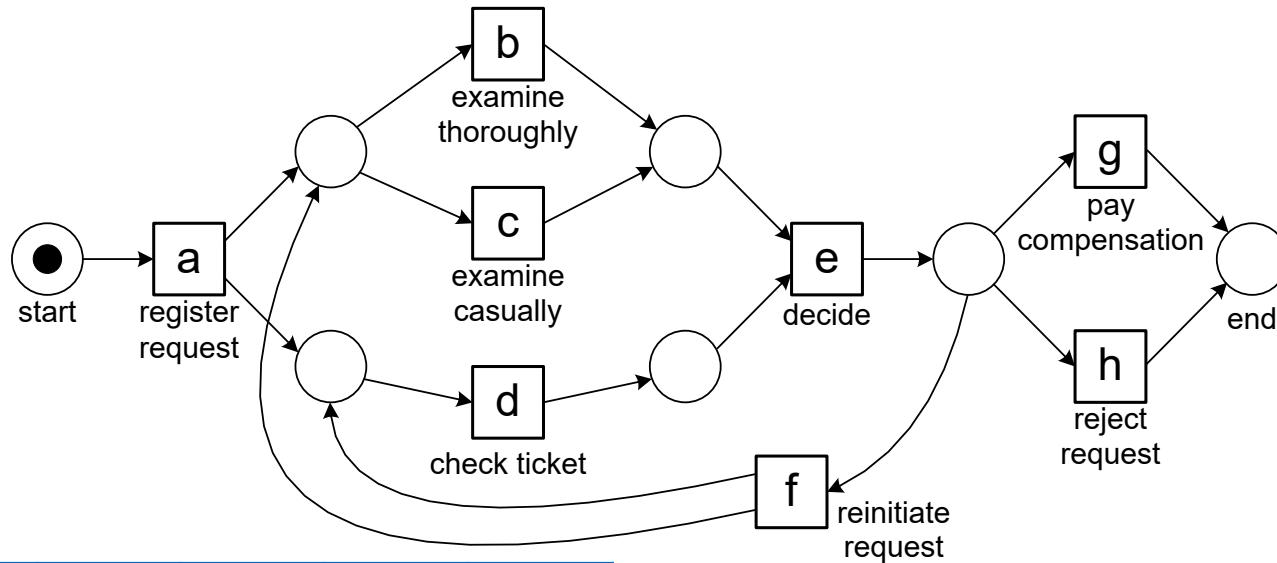
6

Another example: What is the alignment-based fitness of the trace $\langle b, c \rangle$?



$\langle b, c \rangle$

What is the alignment-based fitness of the trace $\langle b, c \rangle$?

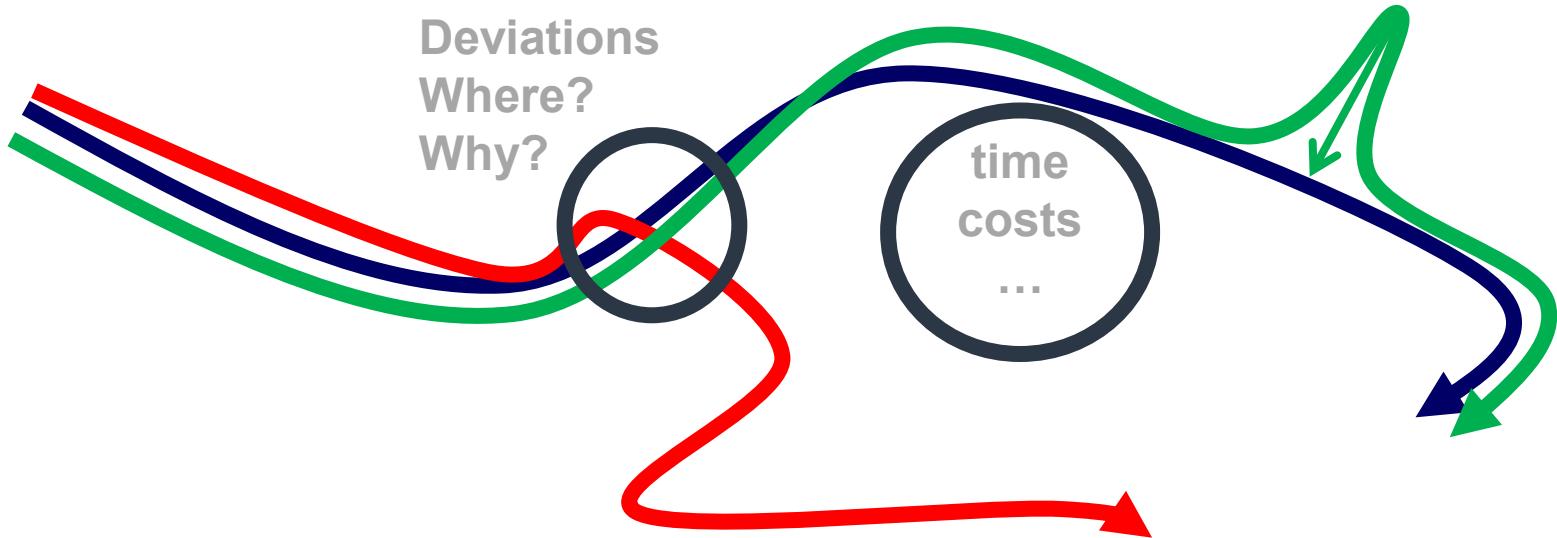


»	b	c	»	»	»
a	b	»	d	e	h

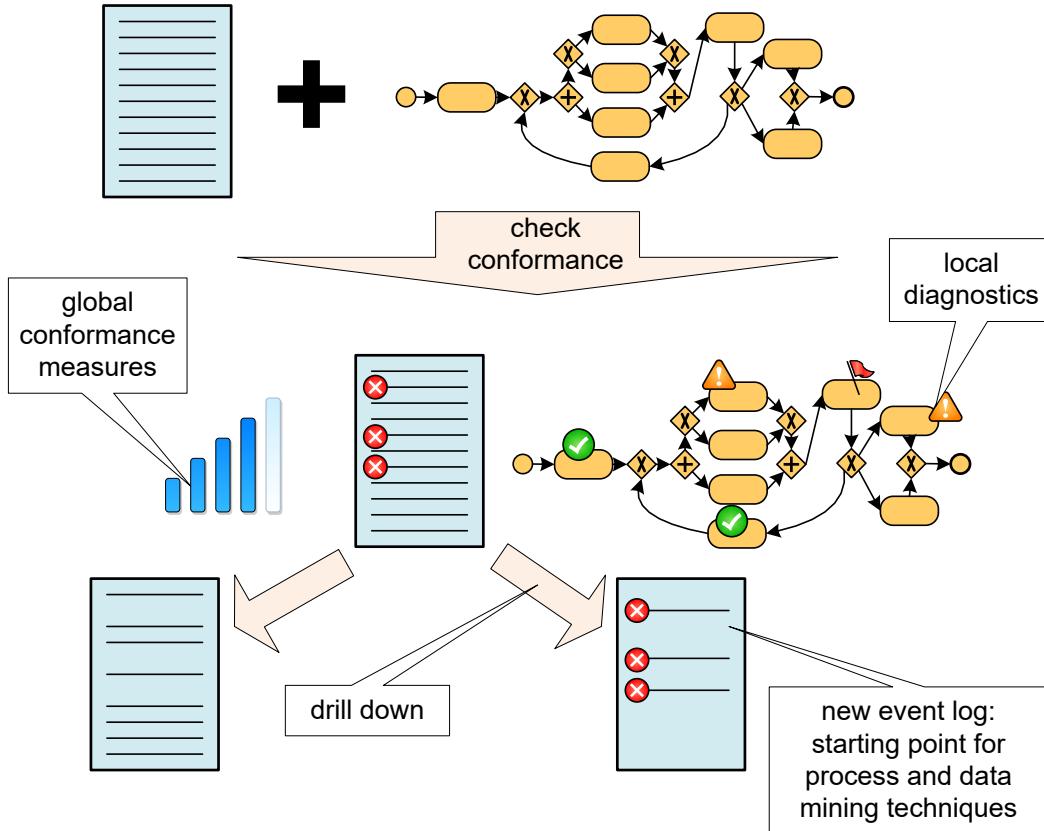
$$1 - \frac{5}{2 + 5} = \frac{2}{7} = 0.286$$

Advantages of aligning log and model

- Observed behavior is directly related to modeled behavior.
- Very flexible (any cost structure).
- Detailed diagnostics.
- After aligning log and model, other quality dimensions can be investigated (separation of concerns).



Drilling down



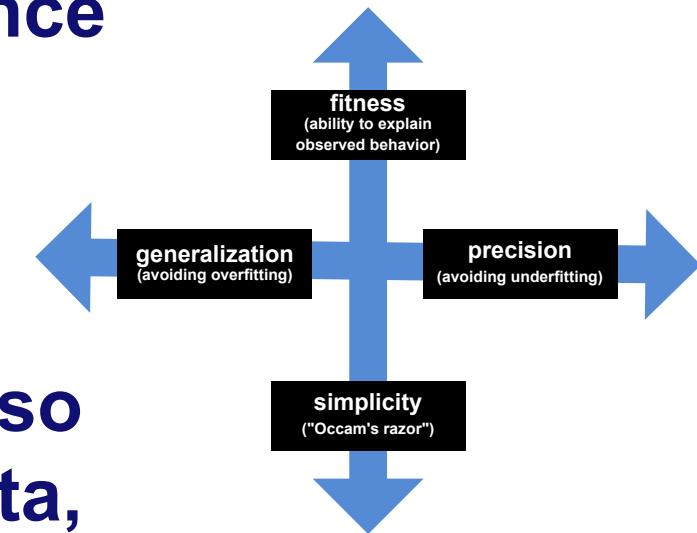
Example approach

- Create event log containing deviating (or non-deviating) cases.
- Apply process discovery to new log.
- Compare process models.

Later more on comparative process mining.

Beyond fitness and control-flow

- There are also solid conformance measures for **precision**, **generalization**, and **simplicity**.
- Multiple definitions possible.
- Conformance checking may also include **other perspectives** (data, resources, time, cost, etc.).
- Example: **data-aware alignments**.



Example: Precision (1/2)

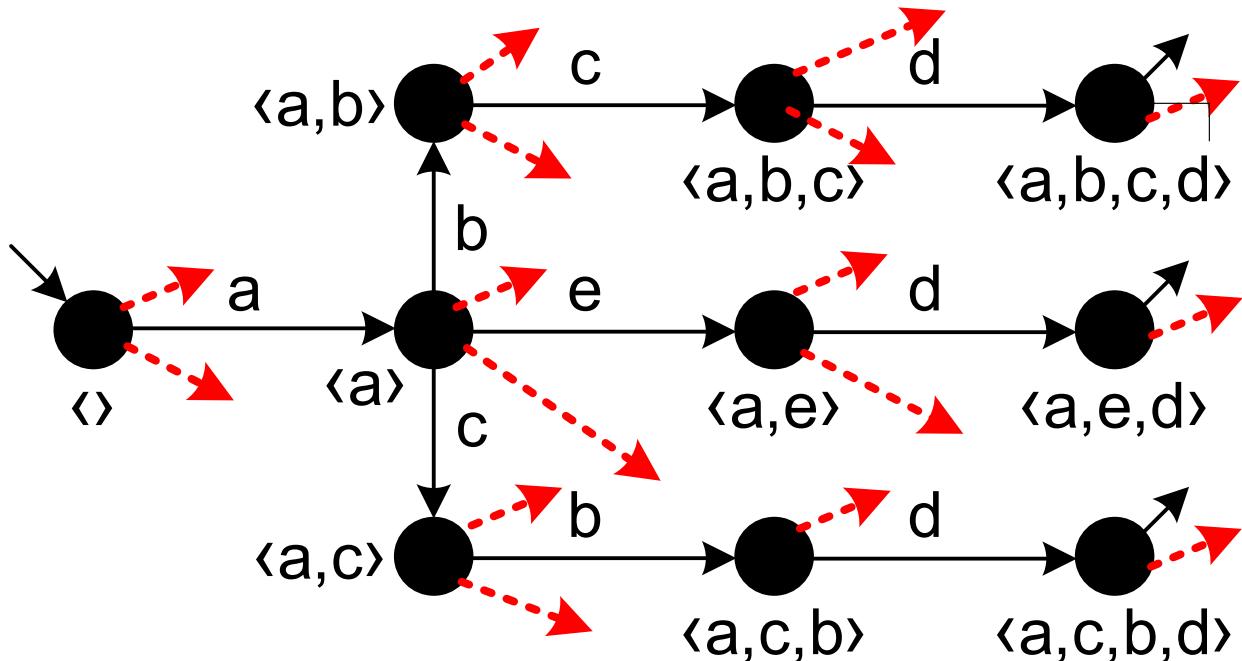
- **Naïve approach:** Compute the fraction of modeled behavior actually observed in the event log.
- **Problems:**
 - If “behavior” = “trace”, then the model with loops has a precision of 0 (because there are ∞ -many traces).
 - Sensitive to the size of the event log. Precision gets better when a longer period is taken.
 - Recall that an event log contains example behavior and often many unique traces.



Example: Precision (2/2)

- Many smarter approaches, e.g., using **escaping edges**.
- Assume the event log has been aligned, i.e., the remaining events are synchronous moves or model moves.
- Build a **prefix automaton** (see lecture on region-based mining) based on the aligned event log.
- Extend the prefix automaton with **escaping edges**, i.e., situation where the model allows for more behavior.
- Quantify such escaping edges.

Prefix Automaton with escaping edges



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

W.M.P. van der Aalst, A. Adriansyah, and B. van Dongen. Replay History on Process Models for Conformance Checking and Performance Analysis. WIREs Data Mining and Knowledge Discovery, 2(2):182-192, 2012.
A. Adriansyah, J. Munoz-Gama, J. Carmona, B.F. van Dongen, and W.M.P. van der Aalst. Measuring Precision of Modeled Behavior. Information Systems and e-Business Management, 13(1):37-67, 2015.

Red arcs indicate that while replaying the model could do more than what was observed in the event log. It is possible to qualify precision by taking into account how often a node is visited and what the fraction of escaping edges is.



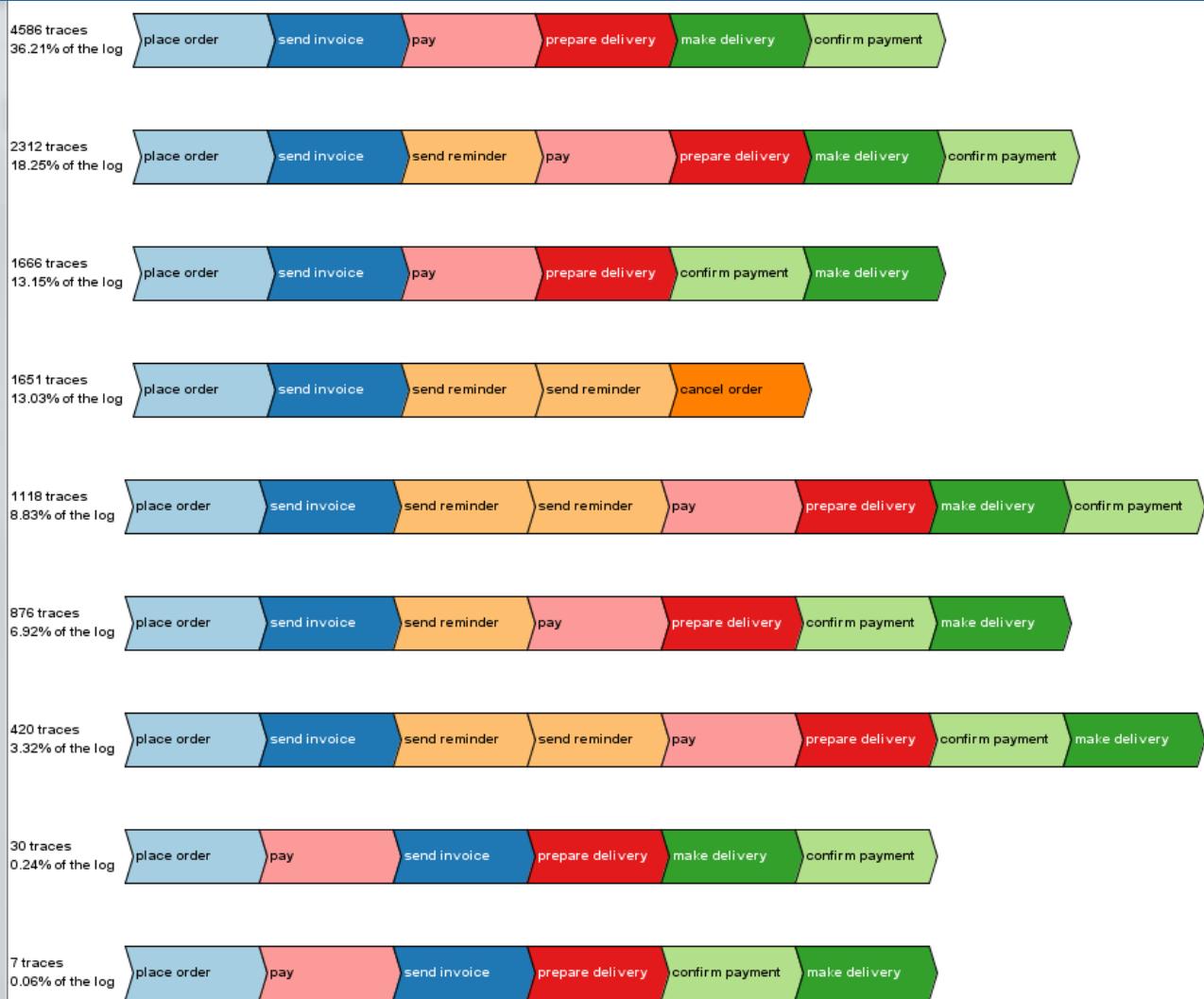
Tooling



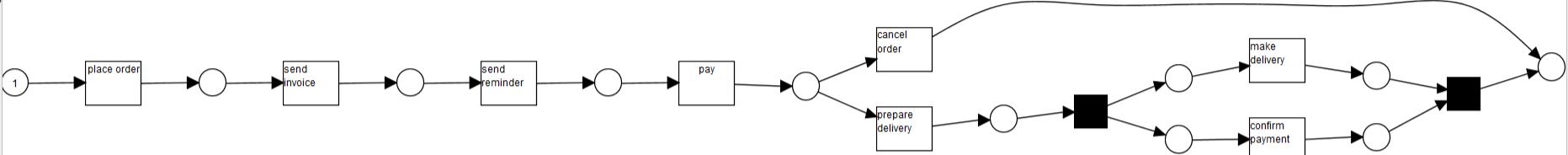
ProM: Load Event Log



Event log



- **12,666 cases**
- **80,609 events**
- **8 unique activities**



No send reminder



No send reminder



Two send reminders



Two send reminders



Two send reminders

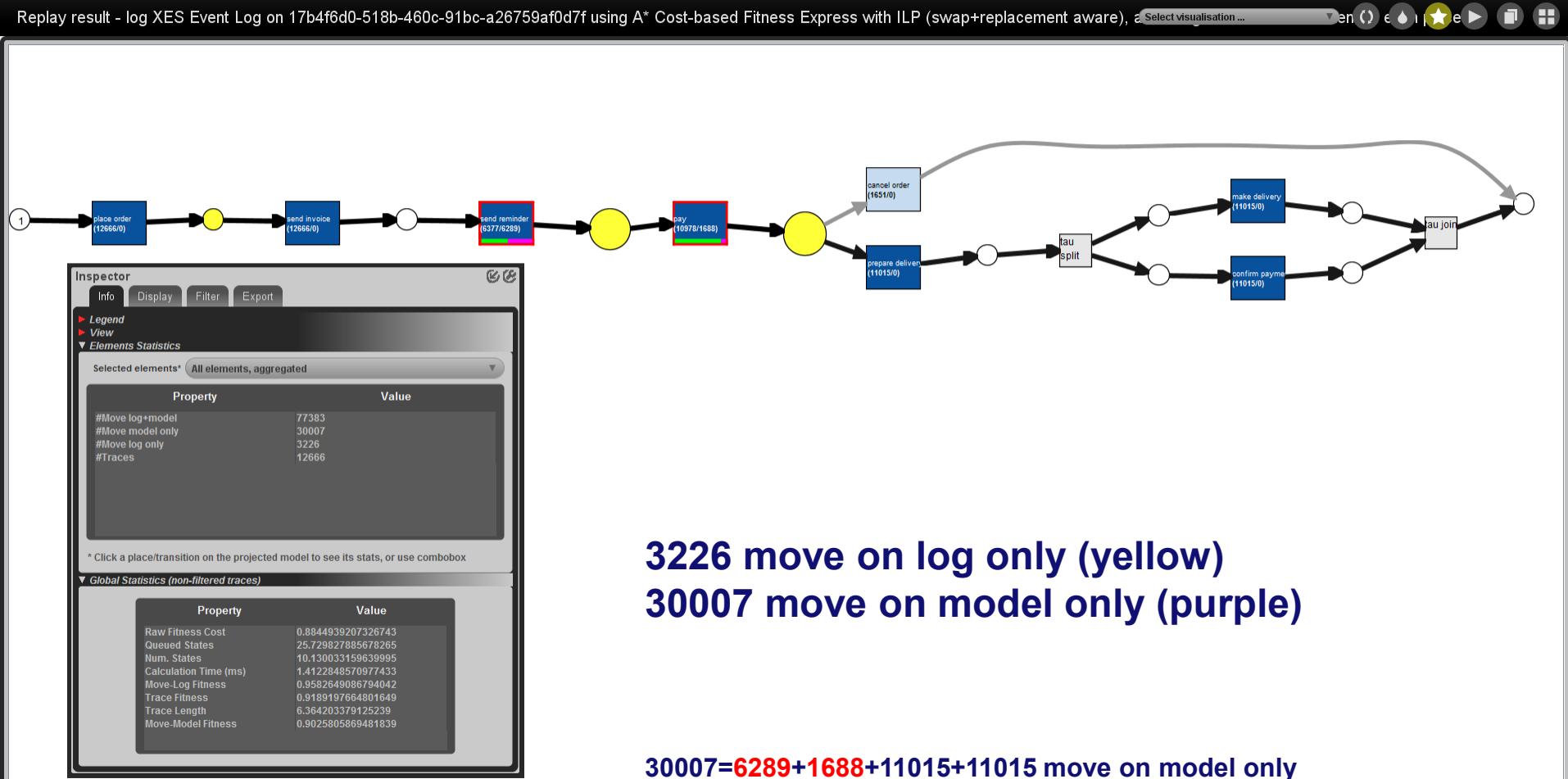


No send reminder and pay before send invoice

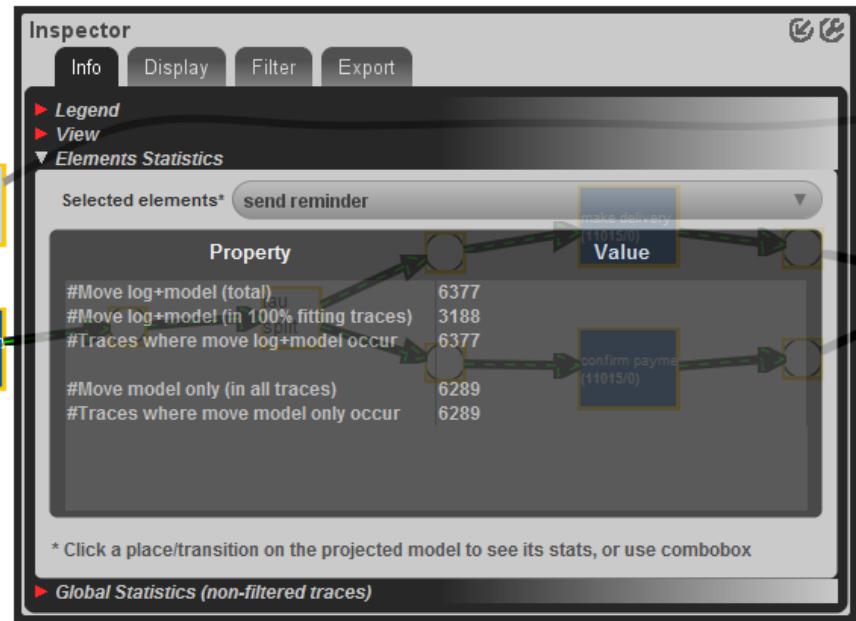
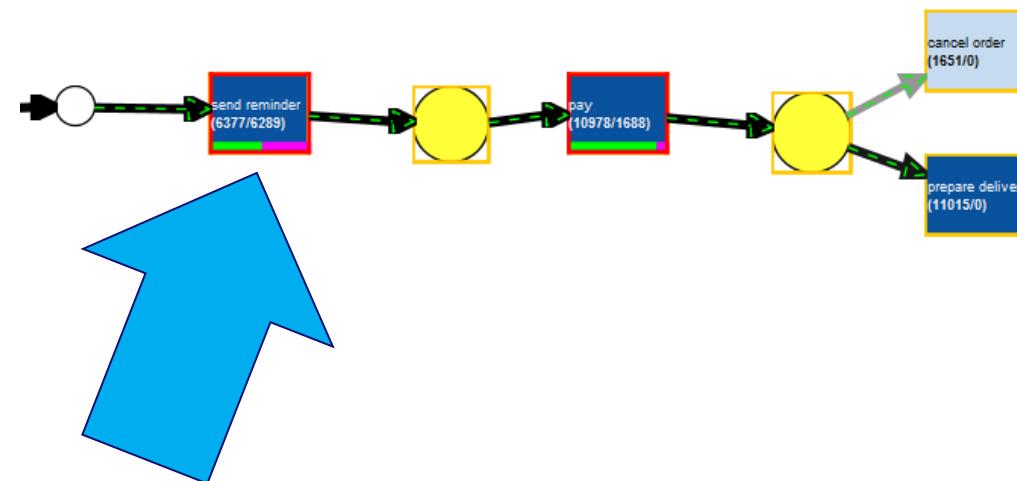


No send reminder and pay before send invoice

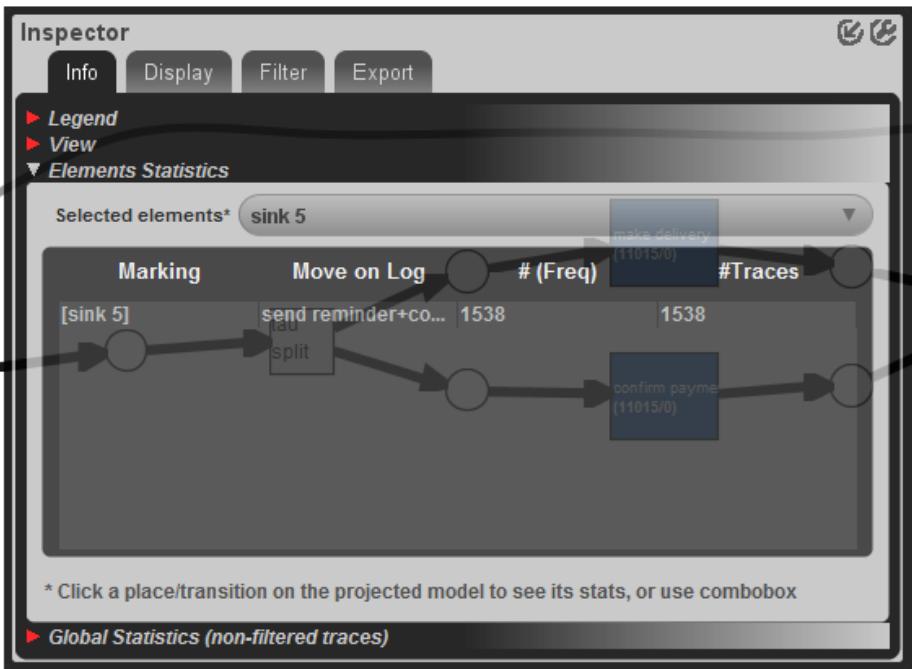
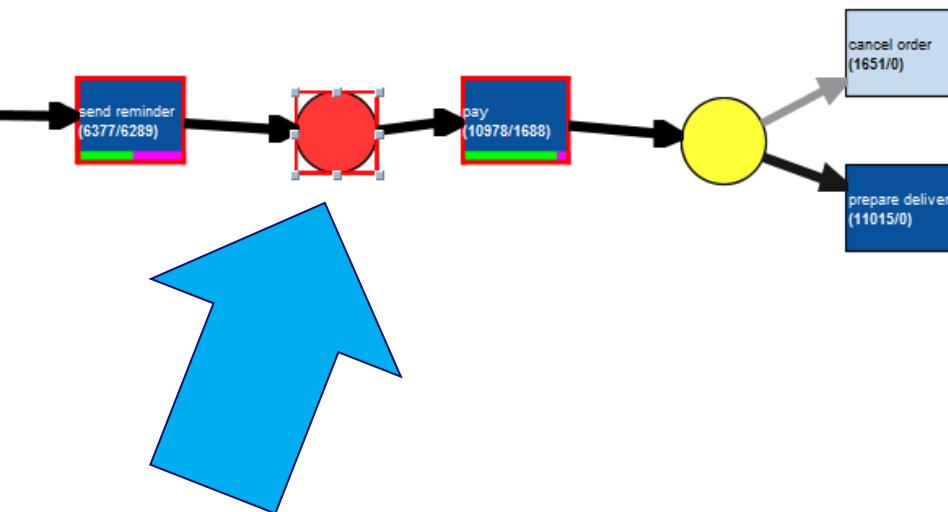




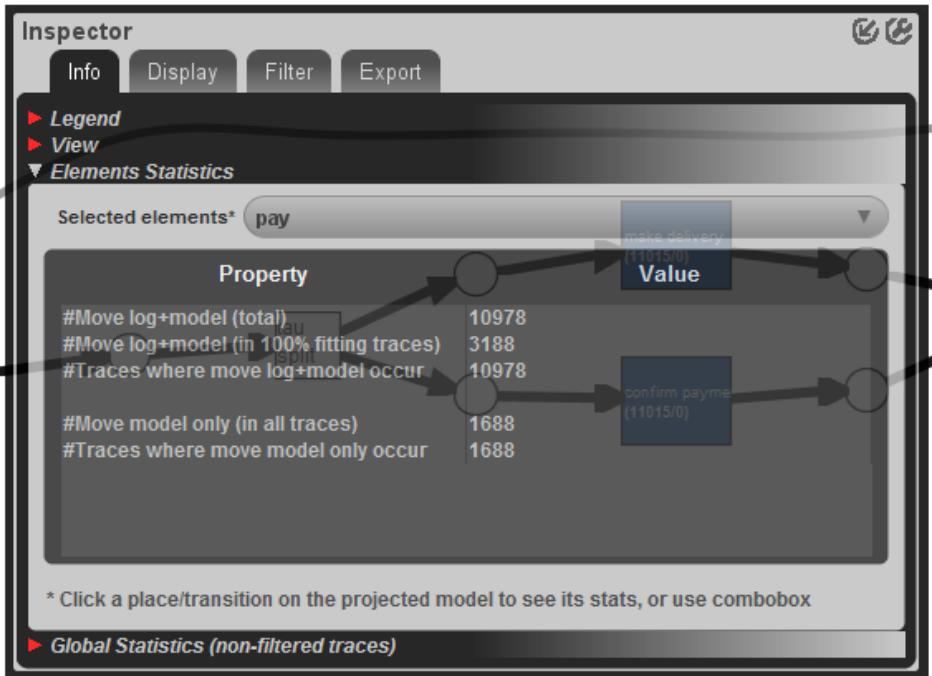
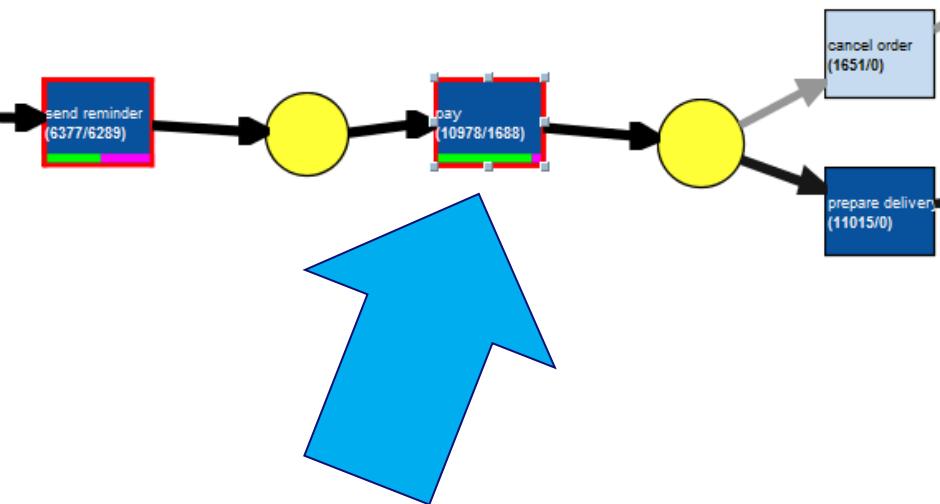
Send reminder is often skipped (6289 times)



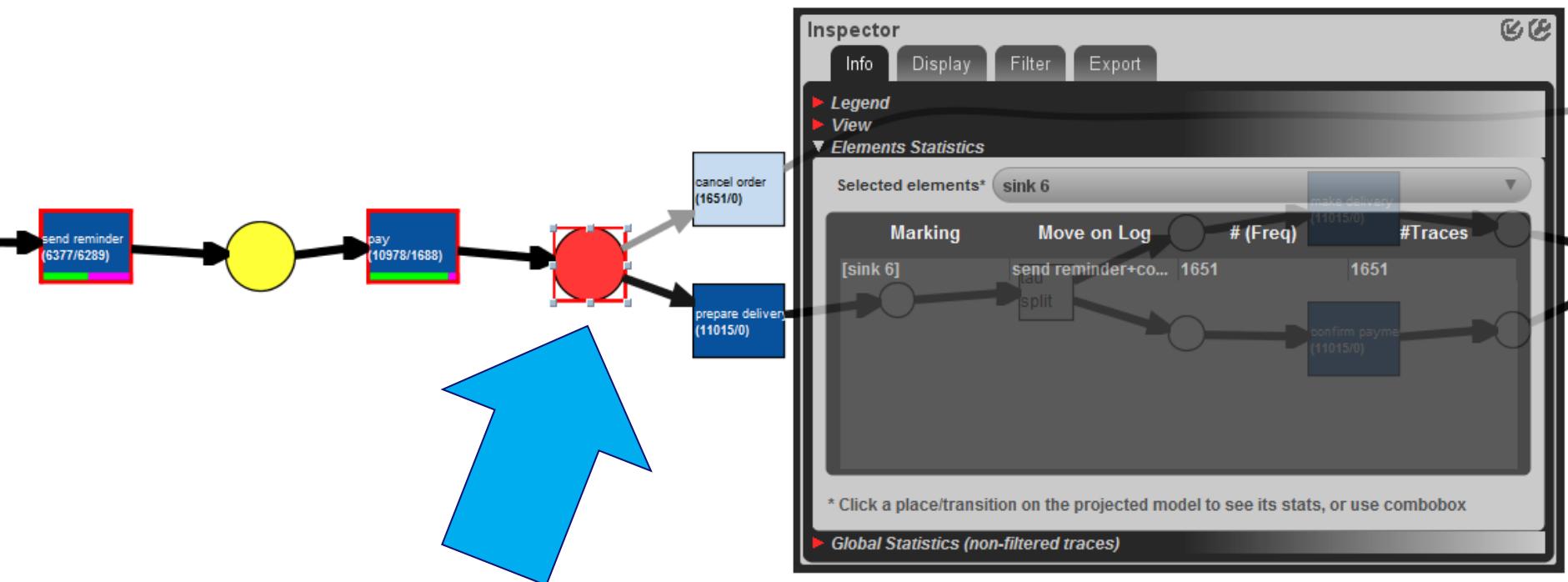
Additional send reminder (1538 times in this place)



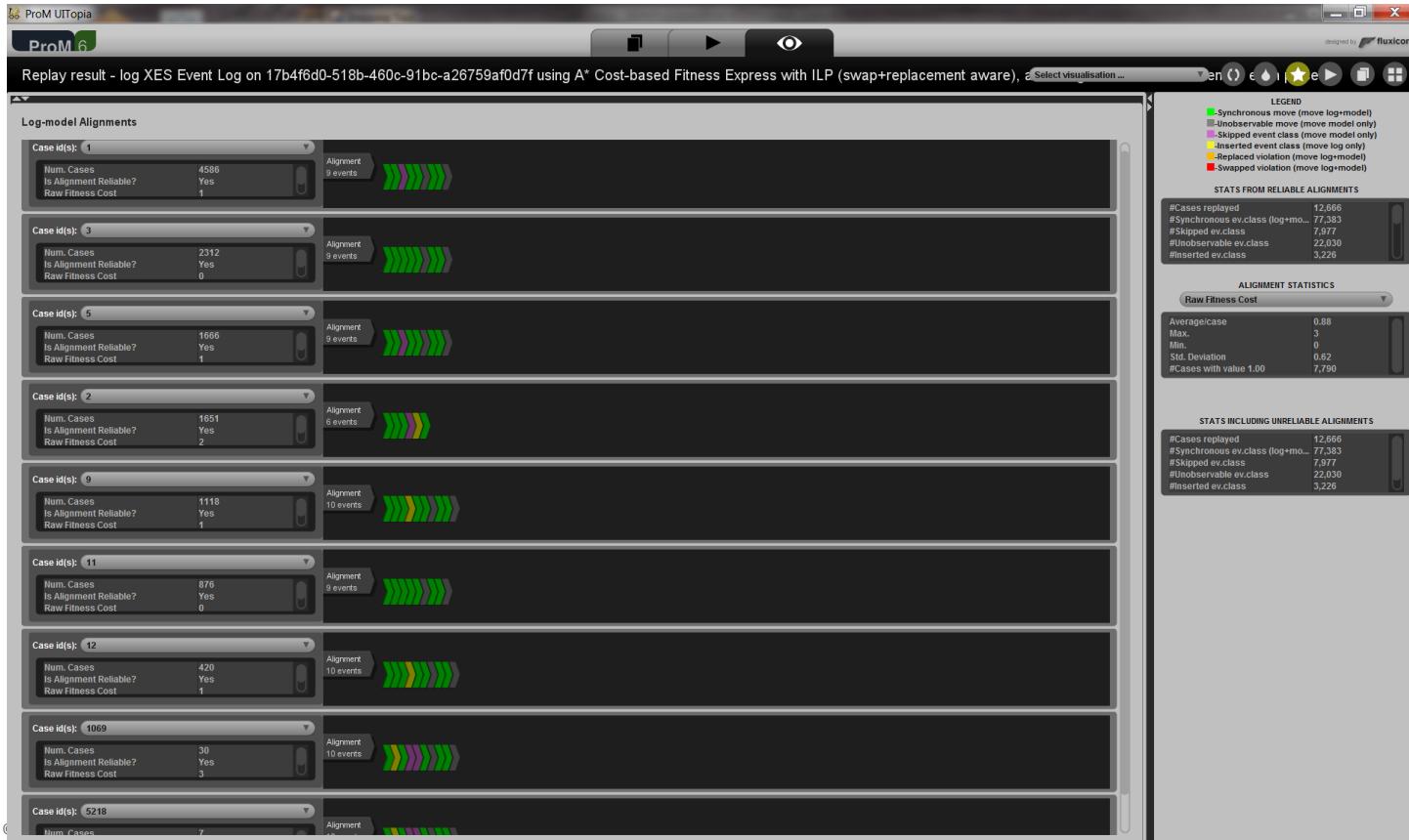
Pay is sometimes skipped (1688 times)



Additional send reminder (1651 times in this place)



Log view



synchronous move

move on model
(required activity was skipped in event log)

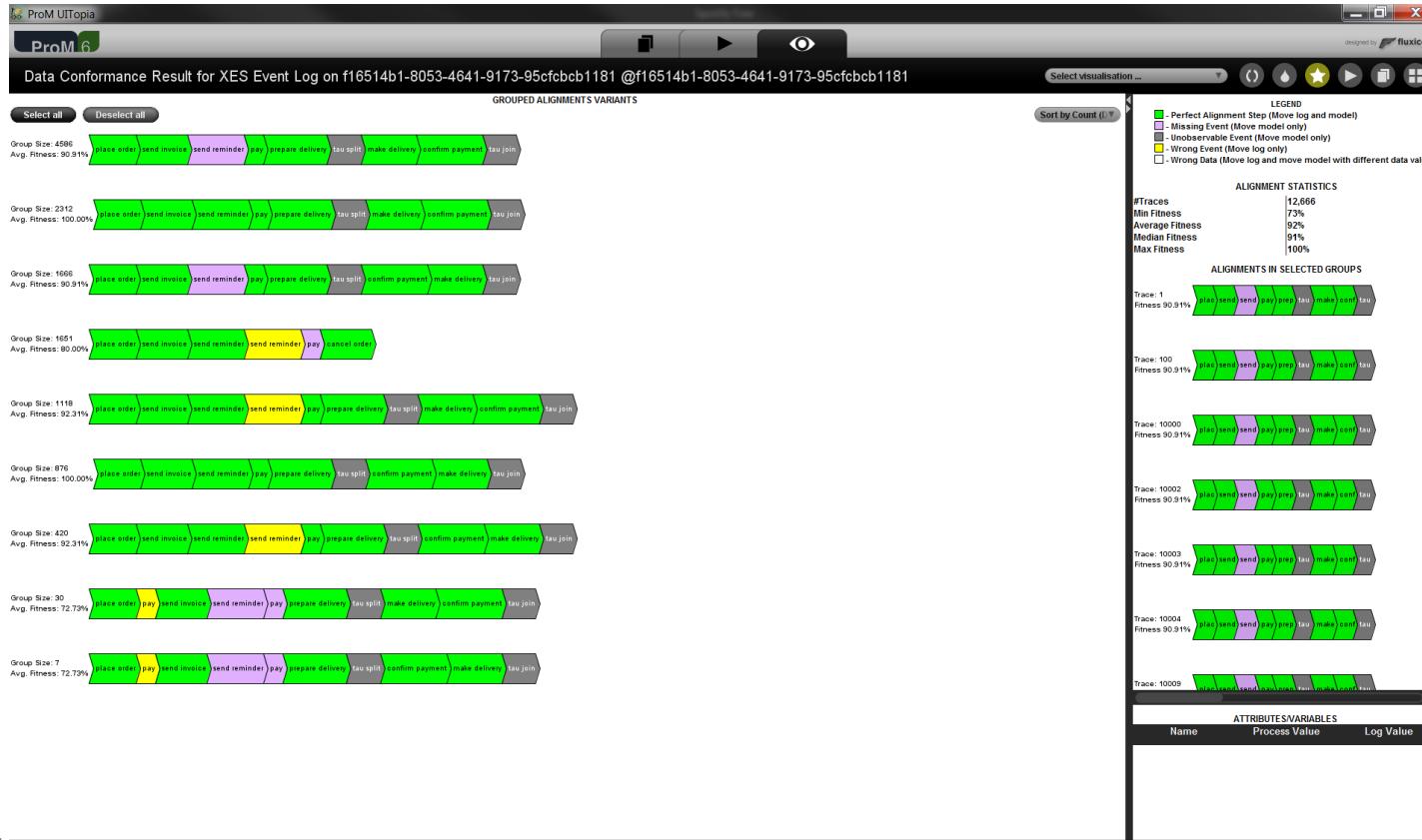
move on model
(silent transition in model was fired)

move on log
(activity in log was not possible in model)



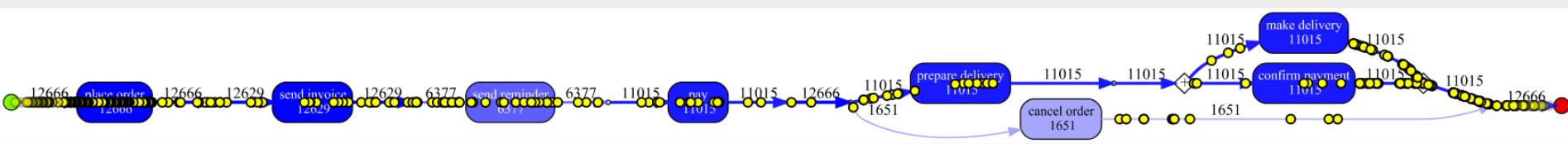
Chair of Process
and Data Science

Many more plug-ins exploring conformance (based on alignments or not)



Chair of Process
and Data Science

Inductive miner



No send reminder



No send reminder



Two send reminders



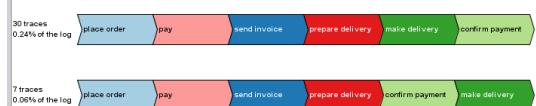
Two send reminders



Two send reminders

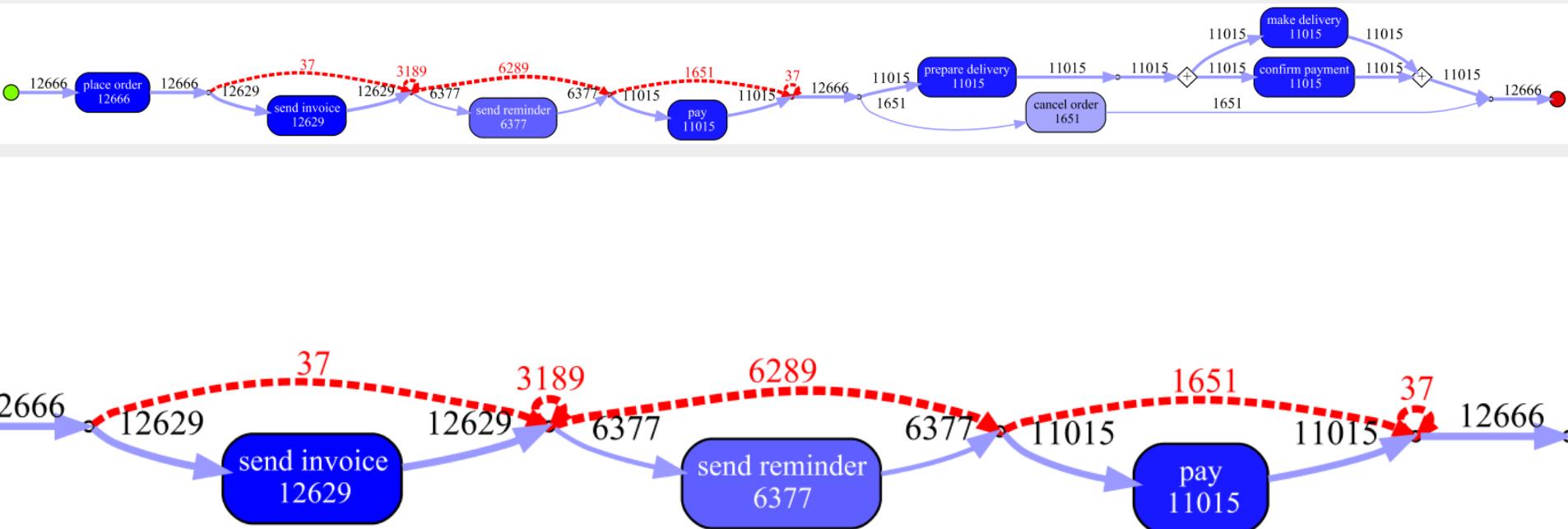


No send reminder and pay before send invoice



No send reminder and pay before send invoice

Deviations in inductive miner



Note that some numbers are different than before because of different alignments.

Drill down on deviations

ProM UItopia

Inductive visual Miner

ProM 6

10001 pla sen sen sen pay can

10005 pla sen sen sen pay can

10014 pla sen sen sen pay can

10019 pla sen sen sen pay can

10020 pla sen sen sen pay can

10027 pla sen sen sen pay can

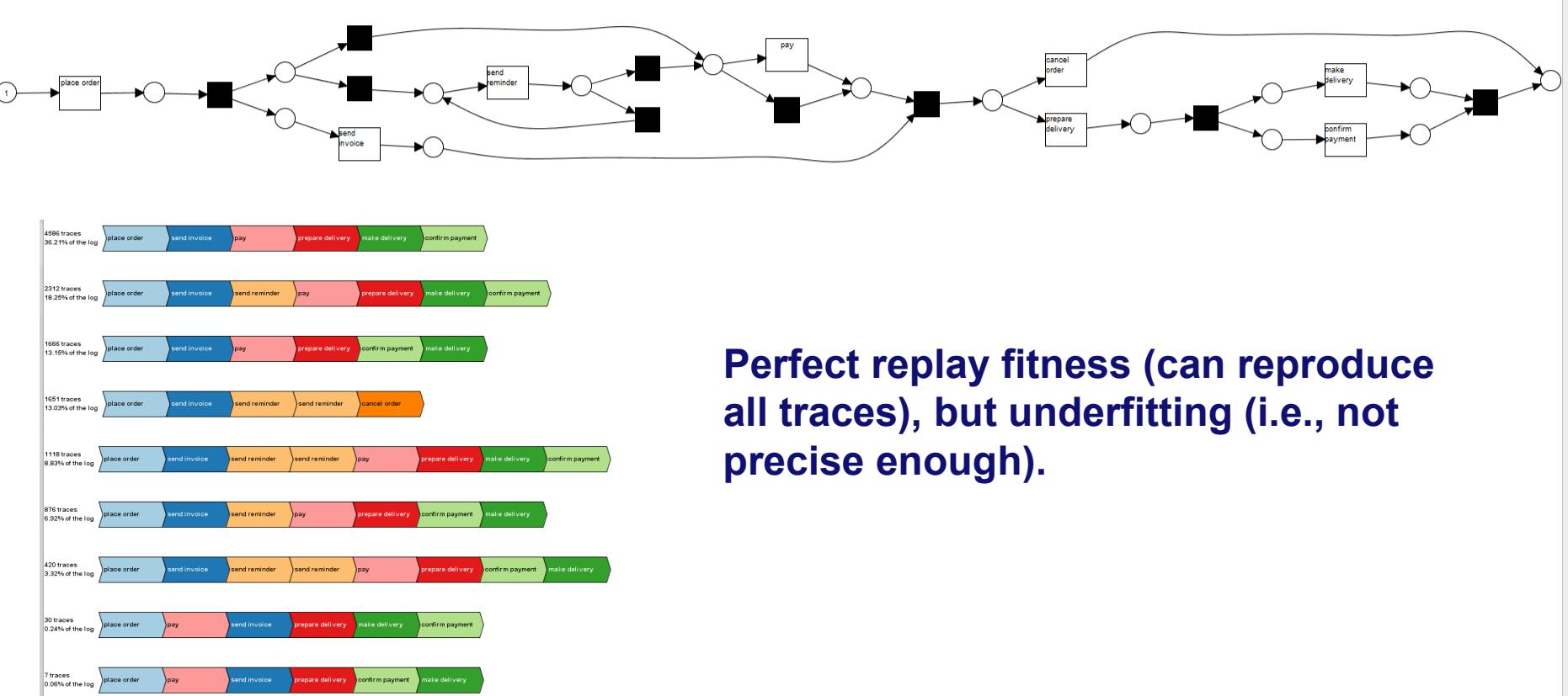
Pay is skipped 1651 times and all of these cases had an extra send reminder (these are the cases that were cancelled).

synchronous move

move on model
(required activity was skipped in event log)

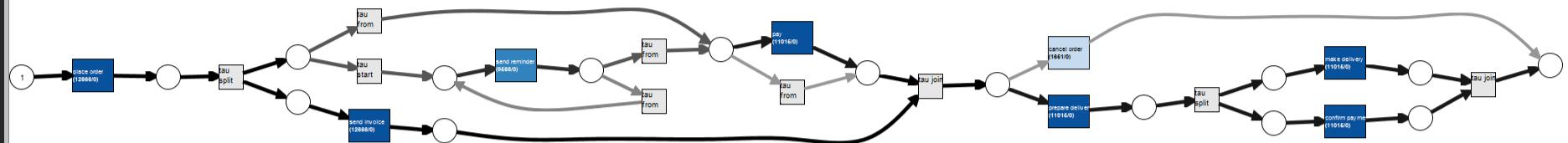
move on log
(activity in log was not possible in model)

Same event log and different (underfitting) model



Perfect replay fitness (can reproduce all traces), but underfitting (i.e., not precise enough).

Replay result - log XES Event Log on ad14ea2b-9e58-4d4c-b840-02be6fb3dff3 using A* Cost-based Fitness Express with ILP (swap+replacement aware), [Select visualisation ...](#)



Inspector

- Info
- Display
- Filter
- Export

Legend

View

Elements Statistics

Property	Value
#Move log+model	80609
#Move model only	71245
#Move log only	0
#Traces	12666

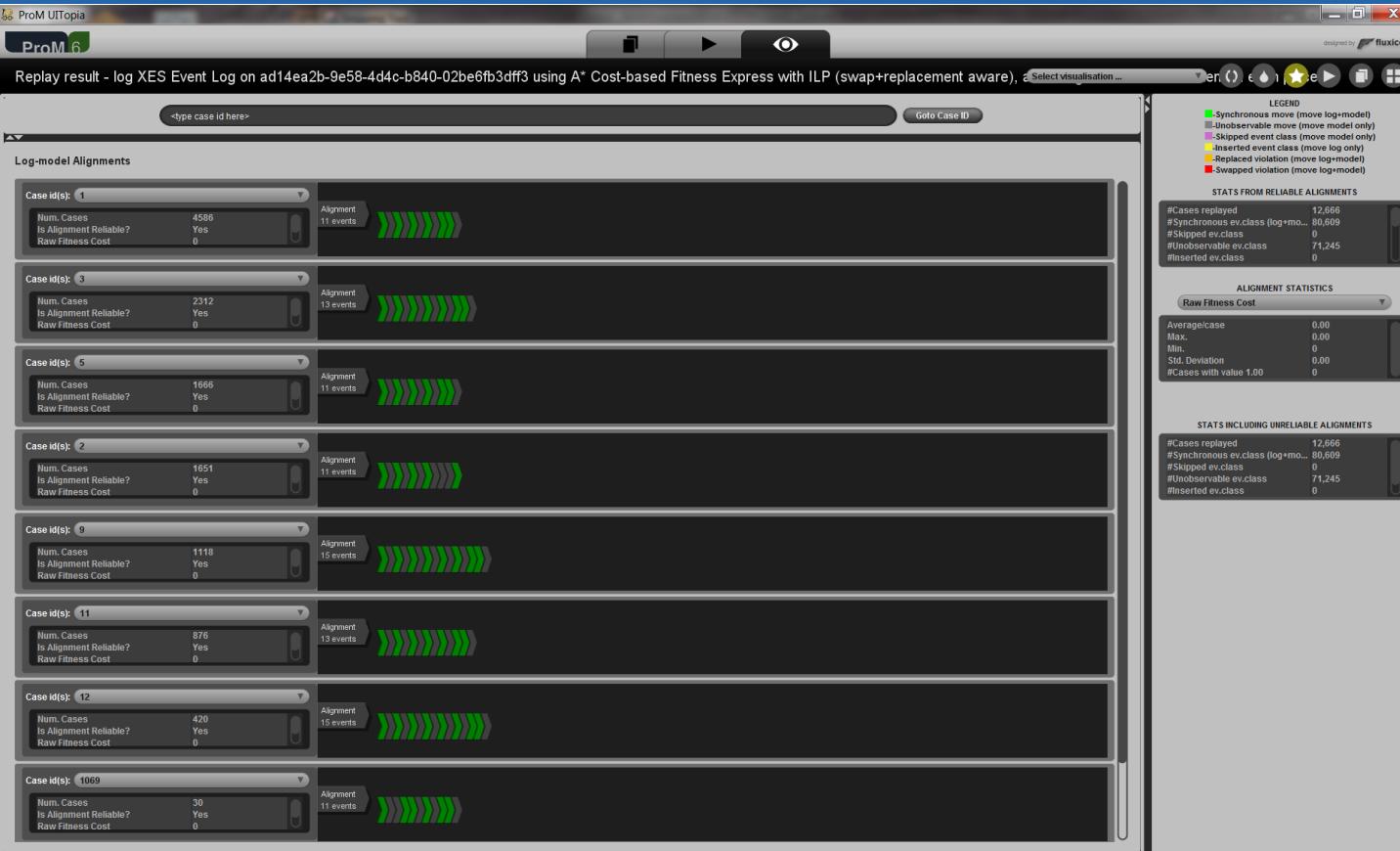
* Click a place/transition on the projected model to see its stats, or use combobox

Global Statistics (non-filtered traces)

Property	Value
Raw Fitness Cost	0.0
Queued States	45.91386388757287
Num. States	15.896652455392376
Calculation Time (ms)	8.713721774830239
Move-Log Fitness	1.0
Trace Fitness	1.0
Trace Length	6.364203379125239
Move-Model Fitness	1.0

No problems, only moves on model for silent steps.

Log view

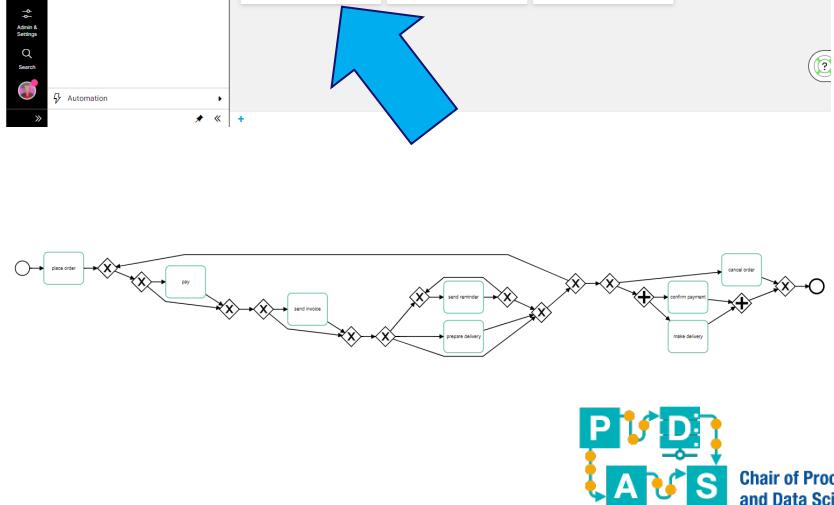


Conformance checking in Celonis

- Input = BPMN model
- Output = List of violations
- Internally: A variant of token-based replay

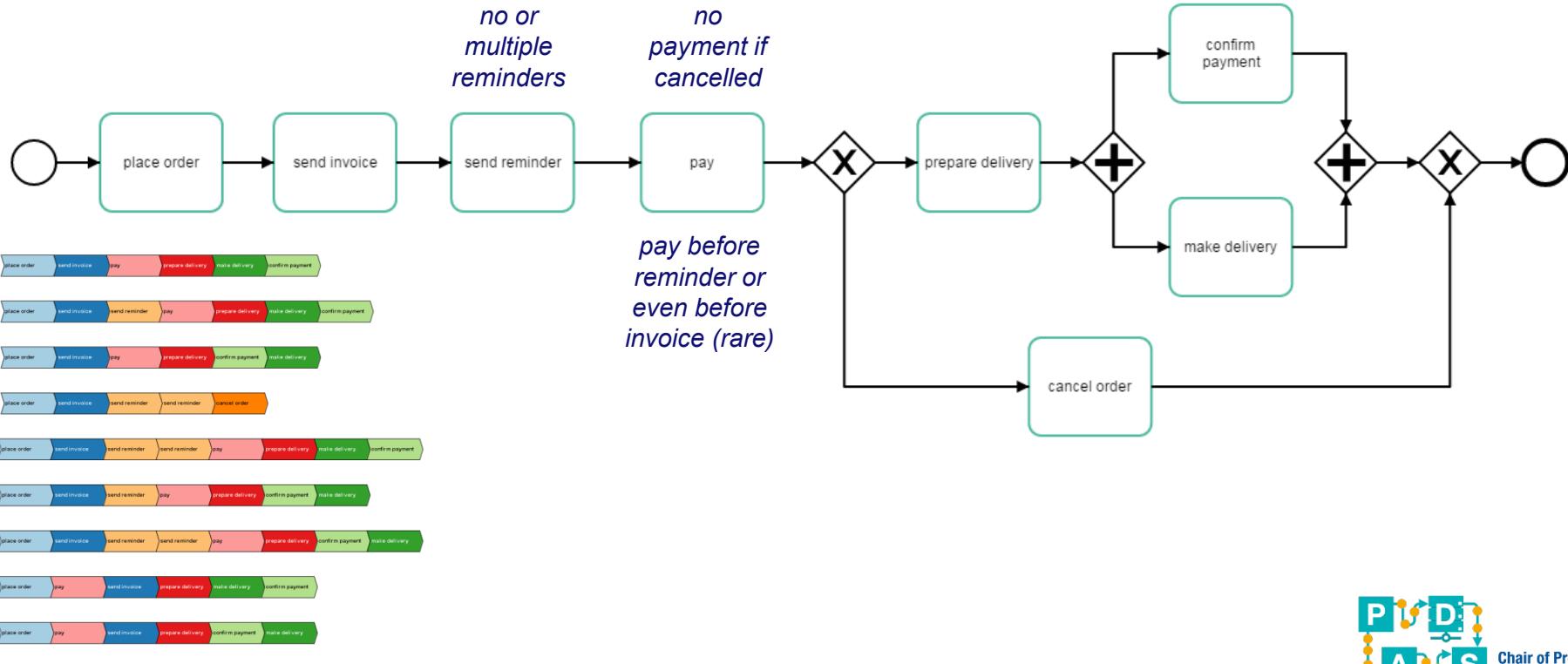
The screenshot shows the Celonis Platform interface. On the left, a sidebar lists various analysis packages: BPI 2022, Create Package, Order Handling 12666, counter example, Pizza, XXXXX, Performance Problems, Purchase-to-Pay, Purchase-to-Pay-Analysis, Vesta-Cancel-Apartment, Vesta-Cancel-Apartment-Analysis, and Order Handling Analysis. The main area displays the 'Order Handling 12666 (Draft)' package. It shows a progress bar at 100% completion with 12.7k of 12.7k cases selected. Below the progress bar are four cards: 'New Sheet' (A new sheet waiting to be built.), 'Process AI' (Detect and analyze deviations from the most common path.), 'Process Overview' (Get the main insights on your process.), and 'Process Explorer' (Analyze and understand your process.). A large blue arrow points upwards towards the 'Conformance' card, which is described as 'Compare the real process to your target process.'.

This screenshot shows the Celonis Platform interface under the 'Automation' tab. It displays three cards: 'Mine the target process' (Upload process model, Select file), 'New process model' (Pull from process repository, Select file), and 'Conformance' (100% completion). The sidebar on the left is identical to the one in the top screenshot, listing the same analysis packages.



Chair of Process
and Data Science

BPMN model (hand-made)



Conformance overview

Timeframe

All time

From

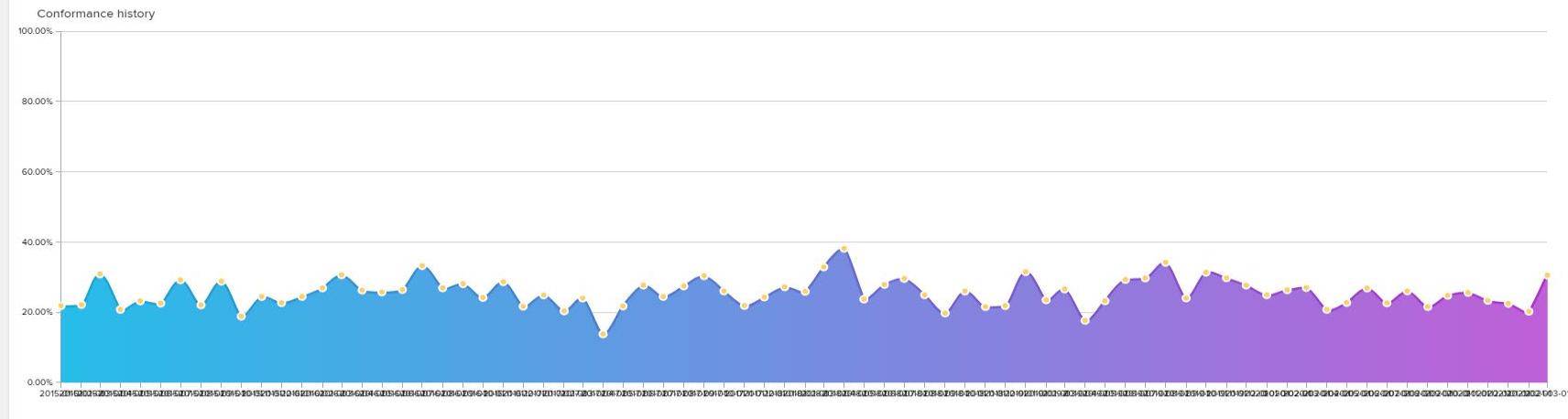
2015-01-04

To

2021-04-27



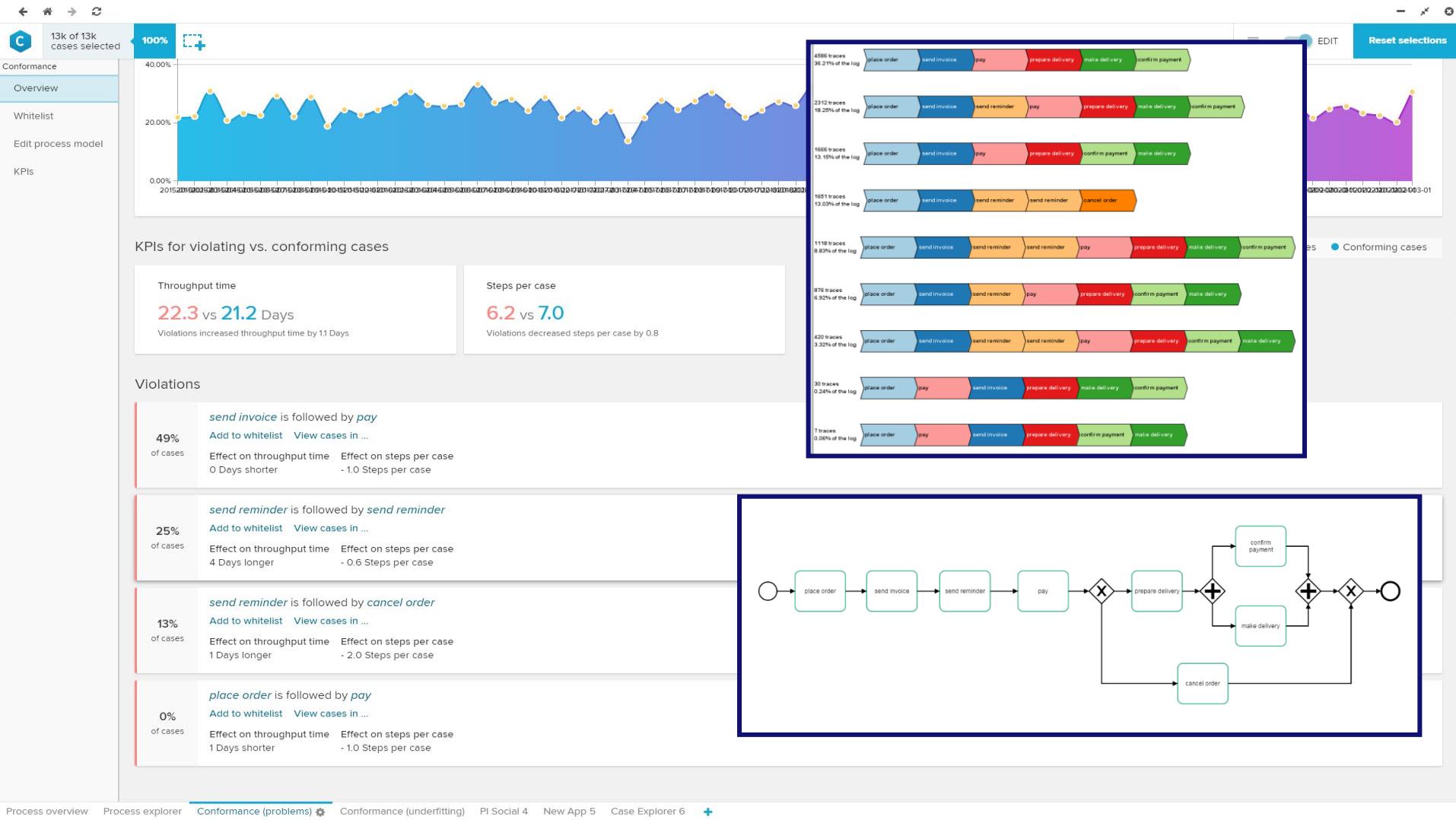
Statistics about conformance

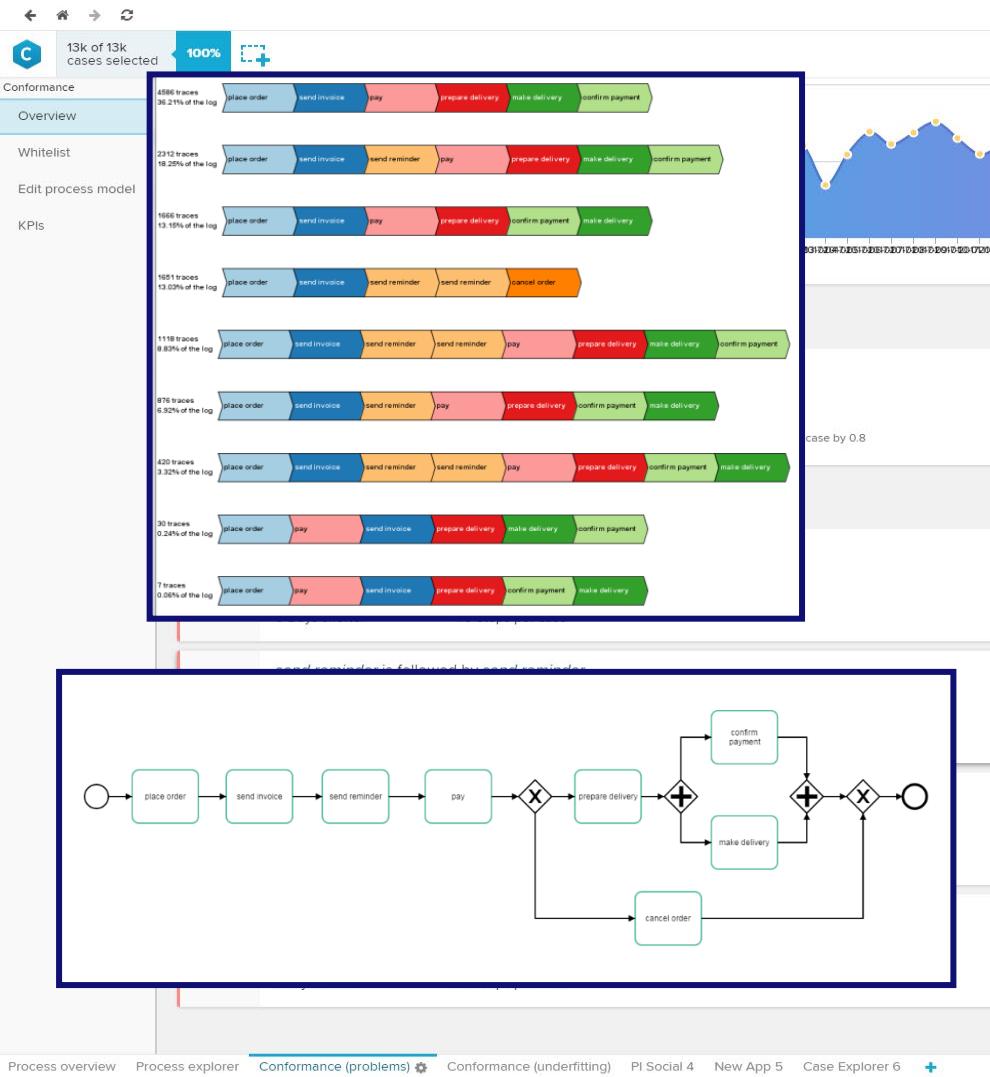


KPIs for violating vs. conforming cases

● Violating cases ● Conforming cases

Violations





Violations

send invoice is followed by *pay*

49%
of cases

Effect on throughput time
0 Days shorter

Effect on steps per case
- 1.0 Steps per case

send reminder is followed by *send reminder*

25%
of cases

Effect on throughput time
4 Days longer

Effect on steps per case
- 0.6 Steps per case

send reminder is followed by *cancel order*

13%
of cases

Effect on throughput time
1 Days longer

Effect on steps per case
- 2.0 Steps per case

place order is followed by *pay*

0%
of cases

Effect on throughput time
1 Days shorter

Effect on steps per case
- 1.0 Steps per case

Drill-down on most rare violation (37 cases)

37 of 13k cases selected 0% Violation: place order → pay 1

EDIT Reset selections

Activities

Connections

```
graph TD; Start((Process Start)) -- 37 --> PlaceOrder[place order]; PlaceOrder -- 37 --> Pay[pay]; Pay -- 37 --> SendInvoice[send invoice]; SendInvoice -- 37 --> PrepareDelivery[prepare delivery]; PrepareDelivery -- 31 --> MakeDelivery[make delivery]; MakeDelivery -- 6 --> ProcessEnd((Process End)); PlaceOrder -- 37 --> Pay; PlaceOrder -.-> SendInvoice; PlaceOrder -.-> PrepareDelivery; Pay -.-> SendInvoice; Pay -.-> PrepareDelivery; SendInvoice -.-> PrepareDelivery; PrepareDelivery -.-> MakeDelivery; MakeDelivery -.-> ProcessEnd;
```

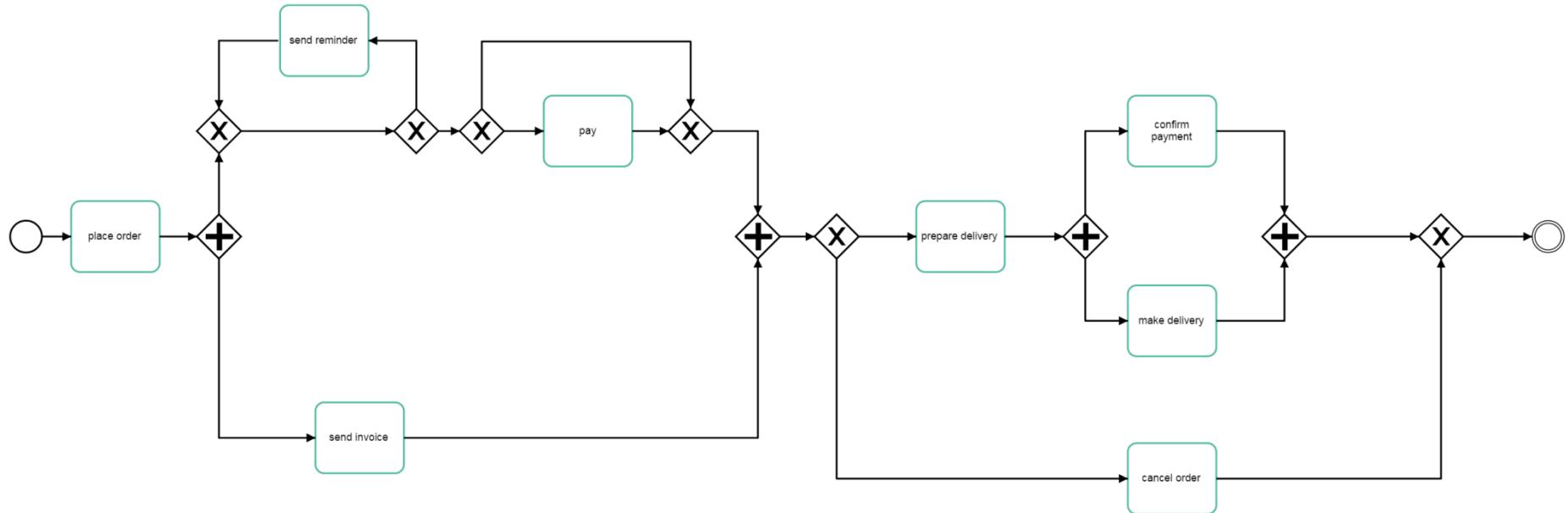
place order is followed by pay
Add to whitelist View cases in ...

Effect on throughput time 1 Days shorter

Effect on steps per case - 1.0 Steps per case

“Almost” perfectly fitting model

(not expecting any deviations)



Conformance

Overviews

Whitelists

Edit process

KPIs

place order ▾
Tue, Aug 4, 2015 12:15 PM -14d

pay ▾
Mon, Aug 17, 2015 9:31 AM -1d

prepare delivery ▾
Tue, Aug 18, 2015 4:33 PM 0

case	1133
end time	Tue, Aug 18, 2015 4:41 PM
resource	Sophia
product	SAMSUNG Galaxy S4
prod-price	3290
quantity	5
address	NL-7943MC-4

send invoice ▾
Tue, Aug 18, 2015 4:33 PM 0

case	1133
end time	Tue, Aug 18, 2015 4:41 PM
resource	Lily
product	SAMSUNG Galaxy S4
prod-price	3290
quantity	5
address	NL-7943MC-4

make delivery ▾
Wed, Aug 19, 2015 4:41 PM +1d

confirm payment ▾
Thu, Aug 20, 2015 3:47 PM +2d

Timeframe
All time From 2015-01-04 To 2021-04-27

g cases /S 4.00

Violations 1 found in process model

Whitelisted violations 0 configured in whitelist

4 cases where prepare delivery and send invoice have exactly the same timestamp

Violations

0.0

100% of cases

misleading diagnostics (send invoice is put after prepare delivery)

pay is followed by **prepare delivery**

Add to whitelist [View cases in ...](#)

Effect on throughput time 18 Days longer

Effect on steps per case + 6.0 Steps per case

pay is followed by prepare delivery

Process overview Process explorer Conformance (problems) Conformance (underfitting) PI Social 4 New App 5 Case Exp

Process Adherence Management (PAM): Discover Object-Centric BPMN

The screenshot shows the Celonis Model Miner interface within a web browser window titled "junk | OrderManagement4OT.pml". The main area displays a BPMN-like process flow for Order Management, featuring various objects like Sales Order, Sales Order Item, Delivery Item, and Customer. The flow consists of several parallel tracks and decision points, primarily colored purple and teal. A large blue banner at the top right of the process diagram states "Uses the Inductive Mining algorithm". On the left side of the interface, there is a sidebar with various icons and a central panel for configuring mining parameters such as "Select object types" (set to 10), "Set the data scope", "Exclude incomplete cases", and "Filter object attributes". The bottom right corner of the interface features the Chair of Process and Data Science logo.

Studio > playground-ordermanagement > OrderManagement4OT

Search CTRL + / 1 Publish

junk Model Miner

Mine a model that is closest to your target process

Build the process model by selecting variants that are close to your target process. In the next step you will be able to edit this model to create your final model.

Select object types Select variants

10 (of 418) most frequent variants selected

Custom

Set the data scope

Exclude incomplete cases

Focus on complete cases to eliminate deviations caused by flows that are still in progress. Define complete cases by selecting valid start and end events for object types.

Filter object attributes

If your process model has a specific scope (e.g. region or product), please set your filters here.

Sales Order

Create Sales Order Header

Create Sales Order Item

Approve Sales Order

Set Delivery Block

Approve Sales Order Item

Delivery

Create Delivery Header

Create Delivery Item

Post Goods Issue

Sign Proof Of Delivery

Create Customer

Change Sales Order Item

Delivery Item

Legend

Save & Continue Cancel

Uses the Inductive Mining algorithm

Zoom to fit 30% 1

P D A S Chair of Process and Data Science

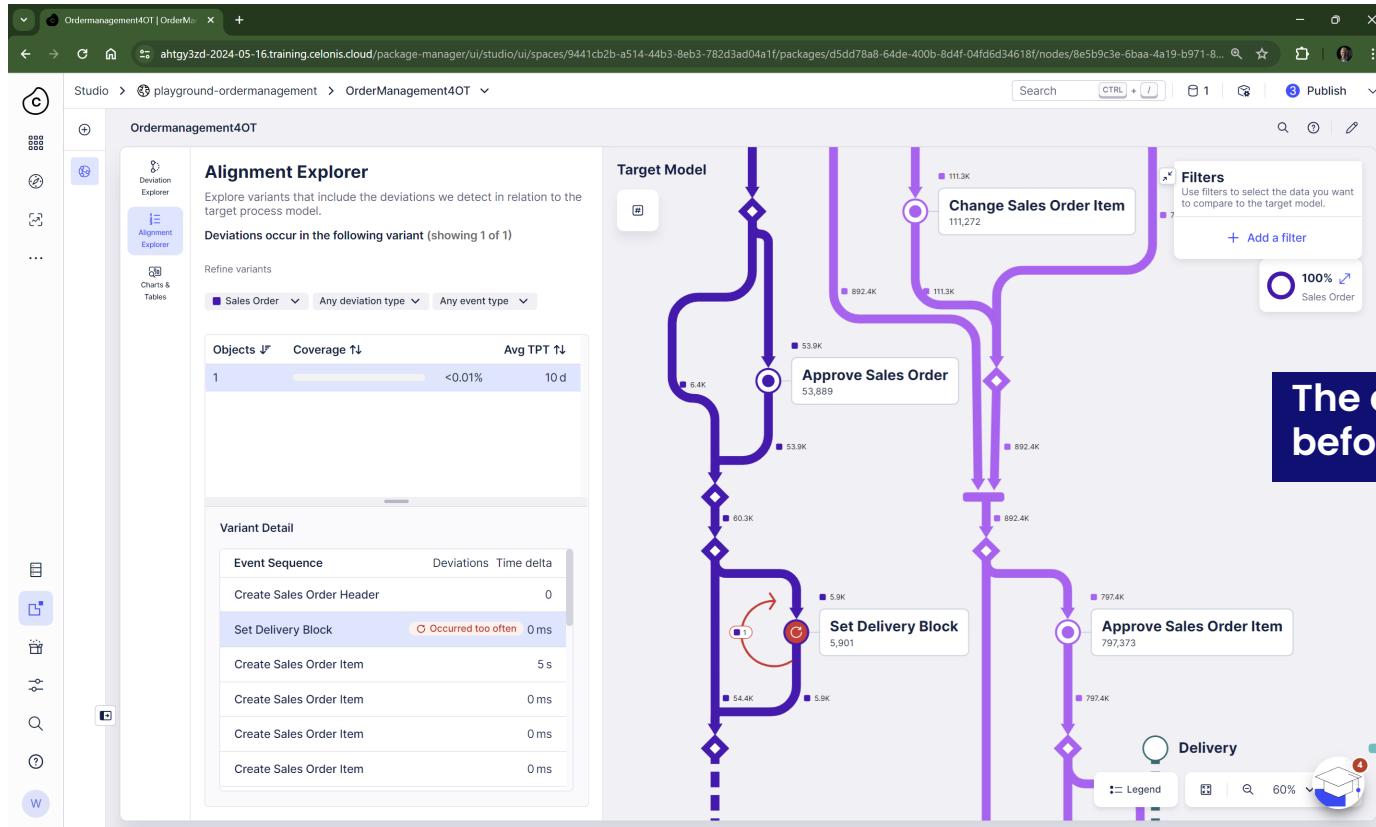
Process Adherence Management (PAM): Check Conformance

The screenshot displays the Celonis Platform interface for the Ordermanagement4OT process. The main view shows a complex process flow diagram with various steps and their associated data. On the left, the 'Deviation Explorer' section provides insights into process deviations, including Conformance Rate (99.99% for Sales Order, 99.57% for Delivery), Actual behavior (Violated exclusive gateway for Create Customer Invoice and Post Goods Issue), and Occurred too often for Set Delivery Block. The 'Target Model' section on the right illustrates the ideal process flow with nodes like Sales Order, Create Sales Order Header, Create Sales Order Item, Change Sales Order Item, Approve Sales Order, and Sales Order Item, each with specific counts (e.g., 60.3K, 892.4K). A 'Filters' panel allows users to select data to compare against the target model. The bottom right corner features a large 'Celonis' logo.

Computes alignments



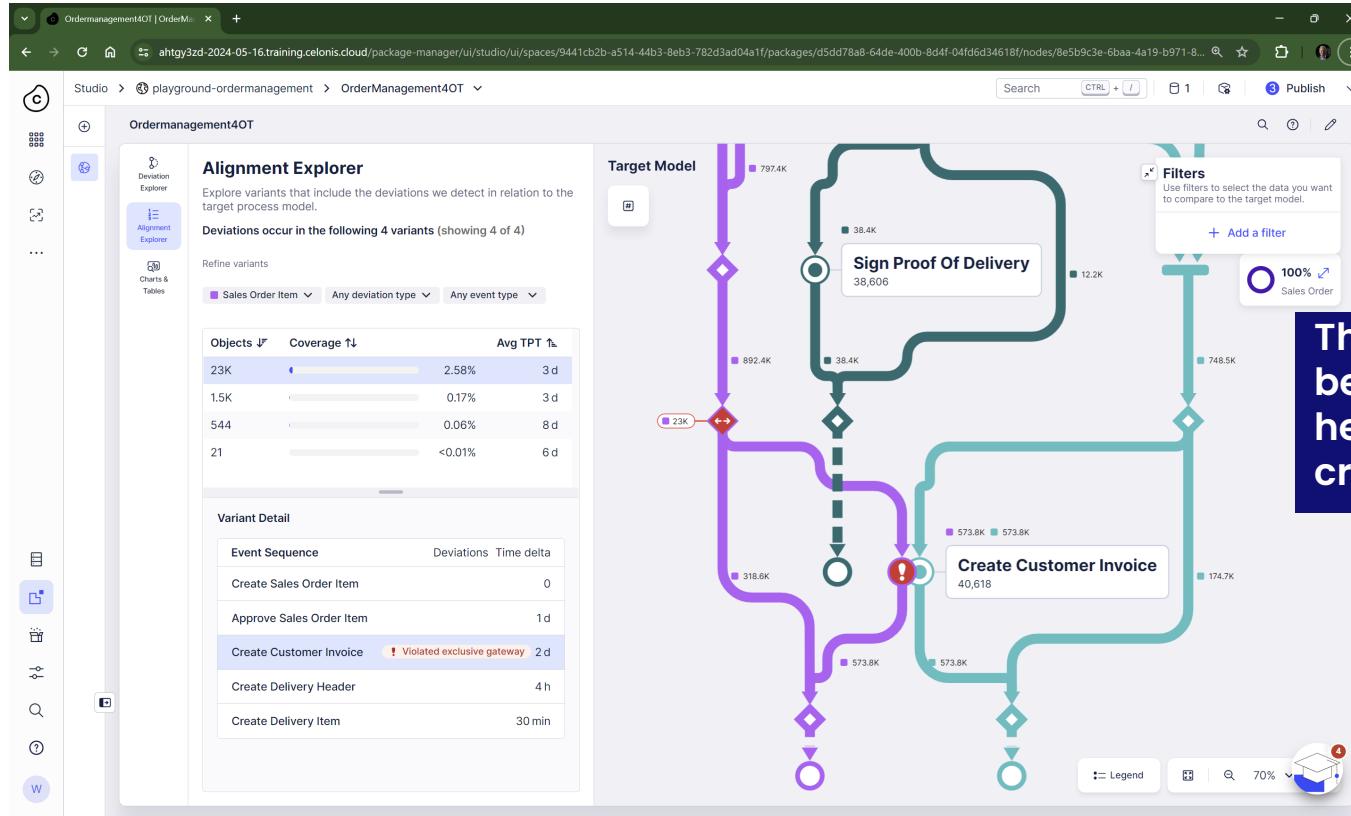
Process Adherence Management (PAM): Show Alignments



The delivery block is set before items are created



Process Adherence Management (PAM): Show Alignments



The invoice is created before the delivery header and item are created.

Process Adherence Management (PAM): Analyze Performance

The screenshot shows the OrderManagement4OT studio interface with the following components:

- Deviation Explorer:** Displays conformance rates for Sales Order (99.99%) and Delivery (99.57%). It also lists actual behaviors: "Violated exclusive gateway Create Customer Invoice" (2.81% objects) and "Violated exclusive gateway Post Goods Issue" (1.55% objects).
- Target Model:** A process flow diagram showing the target model for the order management process.
- Actual behavior:** A process flow diagram showing the actual behavior of the system, overlaid on the target model. Colored lines indicate deviations from the target: purple for Sales Order items and teal for Delivery items.
- Throughput Time:** A detailed analysis of throughput time for specific steps:
 - Create Sales Order Header:** Sales Order, First. Filter: Sales Order.
 - Create Customer Invoice:** Delivery Item, Last. Filter: Delivery Item.Statistics shown: avg. 14.62 Days, med. 12 Days. A histogram shows the distribution of throughput times from 2 to 37 days.

Compute the time between the creation of the order and the delivery of all items in the order



Demo



Applications of Process Mining



Process Mining Is Used Everywhere

Technology



Financial Services & Insurance



Life Sciences & Chemicals



Consumer & Retail



Manufacturing



Telecommunications & Media



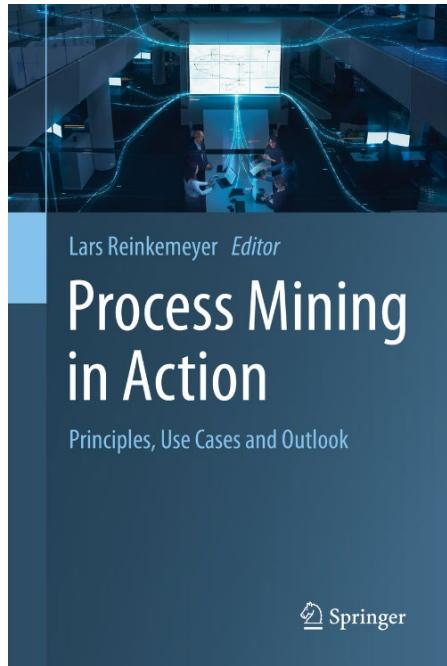
Energy & Utilities



Oil & Gas



Some Case Studies



Part II Best Practice Use Cases

- 9 **Siemens: Driving Global Change with the Digital Fit Rate in Order2Cash** 49
Gia-Thi Nguyen
- 10 **Uber: Process Mining to Optimize Customer Experience and Business Performance** 59
Martin Rowson
- 11 **BMW: Process Mining @ Production** 65
Patrick Lechner
- 12 **Siemens: Process Mining for Operational Efficiency in Purchase2Pay** 75
Khaled El-Wafi
- 13 **athenahealth: Process Mining for Service Integrity in Healthcare** 97
Corey Balint, Zach Taylor, and Emily James
- 14 **EDP Comercial: Sales and Service Digitization** 109
Ricardo Henrique
- 15 **ABB: From Mining Processes Towards Driving Processes** 119
Heymen Jansen
- 16 **Bosch: Process Mining—A Corporate Consulting Perspective** 129
Christian Buhrmann
- 17 **Schukat: Process Mining Enables Schukat Electronic to Reinvent Itself** 135
Georg Schukat
- 18 **Siemens Healthineers: Process Mining as an Innovation Driver in Product Management** 143
Jutta Reindler
- 19 **Bayer: Process Mining Supports Digital Transformation in Internal Audit** 159
Arno Boenner
- 20 **Tekom: Process Mining in Shared Services** 169
Gerrit Lillig



Chair of Process
and Data Science

Example: Siemens Order-to-Cash (O2C)

Siemens: Driving Global Change with the Digital Fit Rate in Order2Cash

Gia-Thi Nguyen

“Using out-of-the-box reports from the Process Mining tool, the results included an increase in automation by 24% and a reduction in manual rework by 11%.”

“Over 70 million sales order items across data from 90 countries with over 1.5 million process variants.”

# Sales Order Items	70,286,004	Statistical Transactional Value in EUR	232,797,668,141
# Activities	411,462,971	# Process Variants	1,511,644
Digital FIT Rate	2.18	SAP Systems	28
Automation Rate	63%	# Countries	90
Rework Rate	36%	# AREs	255
eBiz Rate All-in	64%	# Customers	257,236
Total Cycle Time (average)	48 Days	# Materials	1,728,677



Chair of Process
and Data Science

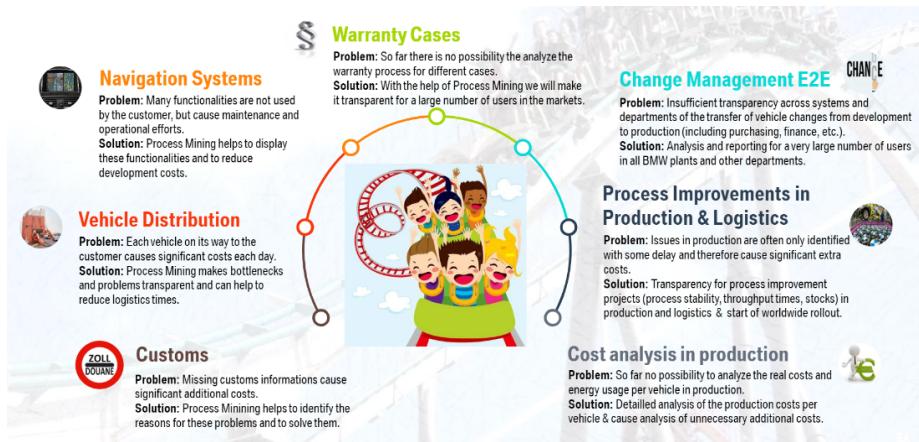
Example: BMW

BMW: Process Mining @ Production

11

Bringing Innovation to Production Processes and Beyond

Patrick Lechner



More than **850 registered users** on BMW Celonis Process Mining Infrastructure

49 of BMW's 50 data models are in use

More than **10 terabyte of raw relational data** is handled on the Process Mining Databases

More than **6.4 million process variants** are currently being analyzed within the BMW Process Mining Infrastructure

More than **500 million events** are analyzed on the BMW Celonis Process Mining Infrastructure

More than **30 million cases** are analyzed on the BMW Celonis Process Mining Infrastructure

Between **400 – 600 analytical views** are provided **each day** for the Process Mining Users of BMW

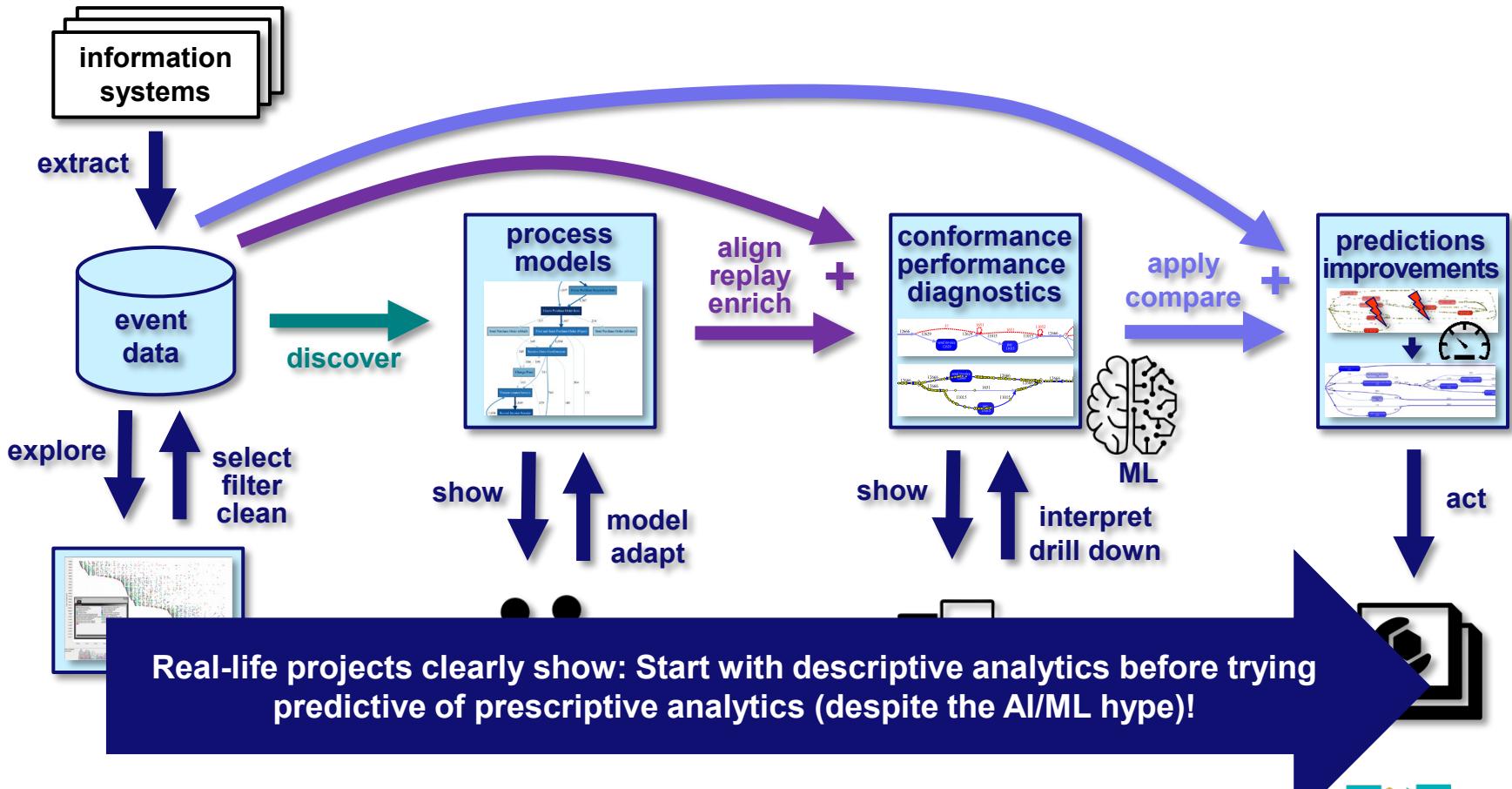
600 unique analyses are currently available for the 49 data models



Chair of Process
and Data Science

Example: Lufthansa





Part I: Introduction

Chapter 1
Data Science in Action

Chapter 2
Process Mining:
The Missing Link

Part II: Preliminaries

Chapter 3
Process Modeling
and Analysis

Chapter 4
Data Mining

Part III: From Event Logs to Process Models

Chapter 5
Getting the Data

Chapter 6
Process Discovery:
An Introduction

Chapter 7
Advanced Process
Discovery Techniques

Part IV: Beyond Process Discovery

Chapter 8
Conformance
Checking

Chapter 9
Mining Additional
Perspectives

Chapter 10
Operational Support

Part V: Putting Process Mining to Work

Chapter 11
Process Mining
Software

Chapter 12
Process Mining in the
Large

Chapter 13
Analyzing “Lasagna
Processes”

Chapter 14
Analyzing “Spaghetti
Processes”

Part VI: Reflection

Chapter 15
Cartography and
Navigation

Chapter 16
Epilogue



ID	Topic	Date	Date	Place
	Lecture 1 Introduction to Process Mining	08.04.24	Monday	AH V
	Lecture 2 Data Science: Supervised Learning	09.04.24	Tuesday	AH V
	<i>Exercise 1 Tool Introduction</i>	09.04.24	Tuesday	AH III
	Lecture 3 Data Science: Unsupervised Learning and Evaluation	15.04.24	Monday	AH V
	Lecture 4 Introduction to Process Discovery	16.04.24	Tuesday	AH V
	<i>Exercise 2 Data Mining</i>	16.04.24	Tuesday	AH III
	Lecture 5 Alpha Algorithm 1	22.04.24	Monday	AH V
	Lecture 6 Alpha Algorithm 2	23.04.24	Tuesday	AH V
	<i>Exercise 3 Petri Nets</i>	23.04.24	Tuesday	AH III
	Lecture 7 Model Quality Representation	29.04.24	Monday	AH V
	Lecture 8 Heuristic Mining	30.04.24	Tuesday	AH V
	<i>Exercise 4 Alpha Miner</i>	30.04.24	Tuesday	AH III
	Lecture 9 Region-Based Mining	06.05.24	Monday	AH V
	<i>Exercise 5 Heuristic Mining and Region-Based Mining</i>	07.05.24	Tuesday	AH III
	Lecture 10 Inductive Mining	13.05.24	Monday	AH V
	Lecture 11 Event Data and Exploration	14.05.24	Tuesday	AH V
	<i>Exercise 6 Inductive Mining</i>	14.05.24	Tuesday	AH III
	Lecture 12 Conformance Checking 1	27.05.24	Monday	AH V
	Lecture 13 Conformance Checking 2	28.05.24	Tuesday	AH V
	<i>Q&A Session Assignment Part I</i>	28.05.24	Tuesday	AH III
	Deadline Assignment Part I	02.06.24	Sunday	
	<i>Exercise 7 Footprint and Token-Based Replay (Exercise)</i>	03.06.24	Monday	AH V
	<i>Exercise 8 Alignments (Exercise)</i>	04.06.24	Tuesday	AH V
	Lecture 14 Decision Mining	10.06.24	Monday	AH V
	<i>Lecture 15 Celonis Guest Lecture</i>	11.06.24	Tuesday	AH V
	<i>Exercise 9 Decision Mining</i>	11.06.24	Tuesday	AH III
	Lecture 16 Performance Analysis and Organizational Mining	17.06.24	Monday	AH V
	<i>Exercise 10 Performance Analysis (Exercise)</i>	18.06.24	Tuesday	AH V
	<i>Exercise 11 Organizational Mining</i>	18.06.24	Tuesday	AH III
	<i>Exercise 12 Celonis Case Study</i>	24.06.24	Monday	AH V
	Lecture 17 Operational Support and Process Mining Applications	01.07.24	Monday	AH V
	Lecture 18 Distributed, Streaming, and Comparative Process Mining	02.07.24	Tuesday	AH V
	<i>Exercise 13 Operational Process Mining</i>	02.07.24	Tuesday	AH III
	Lecture 19 Closing	08.07.24	Monday	AH V
	<i>Q&A Session Assignment Part II</i>	09.07.24	Tuesday	AH III
	Deadline Assignment Part II	14.07.24	Sunday	
	<i>Q&A Session Exam</i>	16.07.24	Tuesday	AH III

