

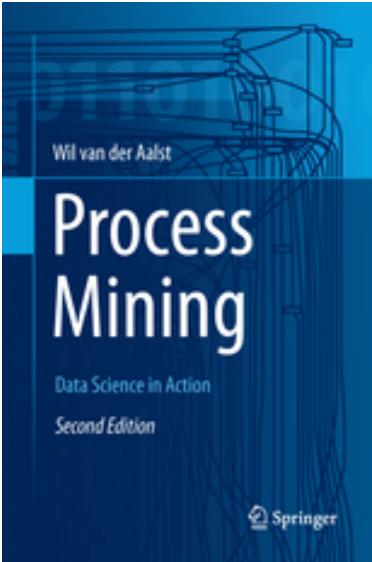
# Introduction to Process Mining

Josep Carmona

# Syllabus

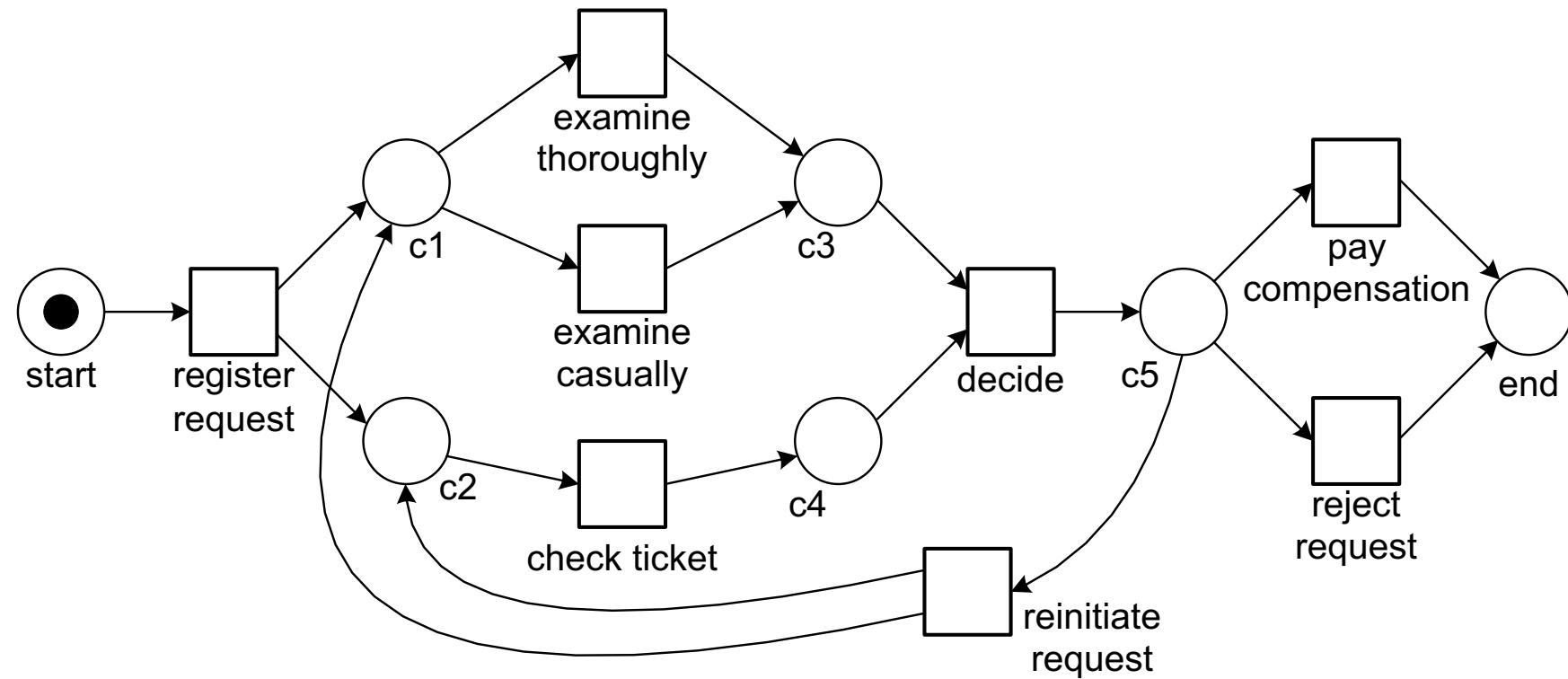
- 9/05 (TODAY): Process mining intro, Process Discovery [T]
- 13/05: DISCO + Inductive Miner [L]
- 15/05: Conformance Checking [T]
- 20/05: Process Mining in Python [L]

# Disclaimer

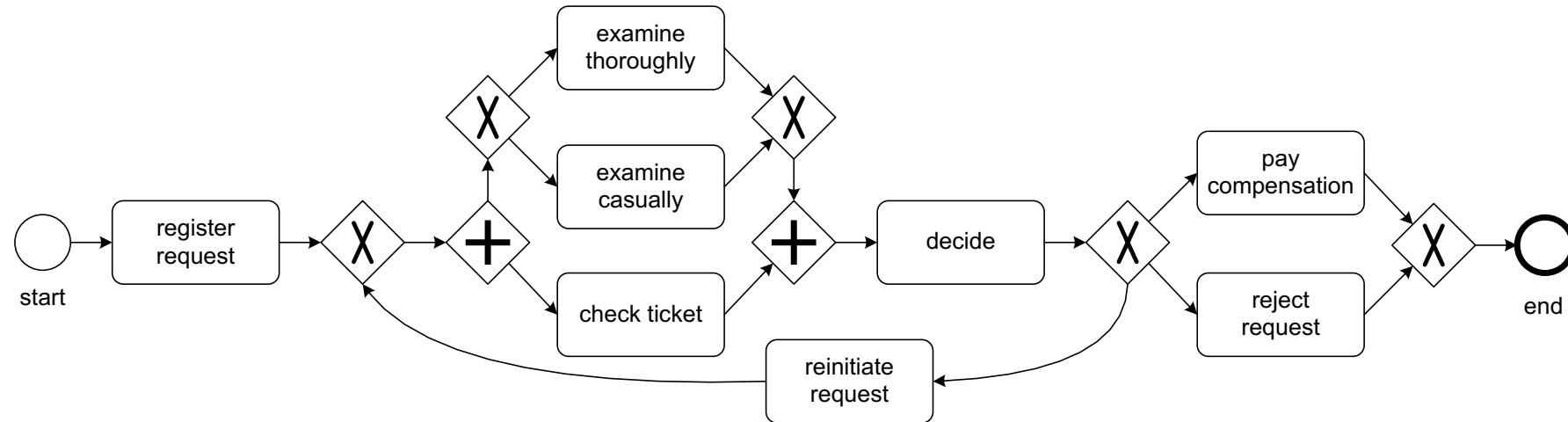


Part of this presentation is taken from the slides of the book from Prof. Wil van der Aalst:  
**"Process Mining - Discovery, Conformance and Enhancement of Business Processes.** Springer 2016, ISBN 978-3-662-49850-7"

# Example process model



# Same process in terms of BPMN rather than Petri nets



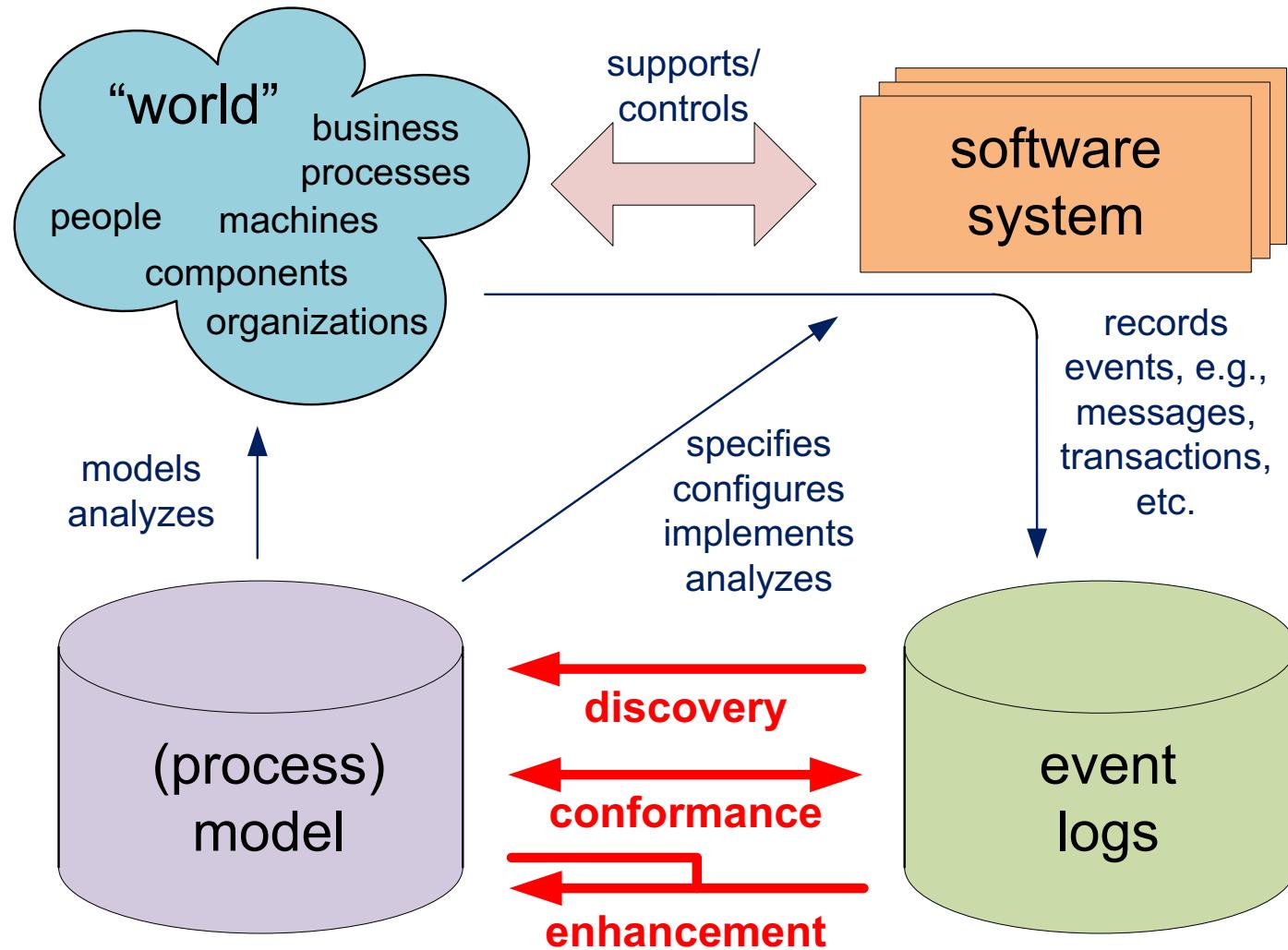
# What are process models used for?

- **insight**: while making a model, the modeler is triggered to view the process from various angles;
- **discussion**: the stakeholders use models to structure discussions;
- **documentation**: processes are documented for instructing people or certification purposes (cf. ISO 9000 quality management);
- **verification**: process models are analyzed to find errors in systems or procedures (e.g., potential deadlocks);
- **performance analysis**: techniques like simulation can be used to understand the factors influencing response times, service levels, etc.;
- **animation**: models enable end users to “play out” different scenarios and thus provide feedback to the designer;
- **specification**: models can be used to describe a PAIS before it is implemented and can hence serve as a “contract” between the developer and the end user/management; and
- **configuration**: models can be used to configure a system.

# Limitations

- Executable models may be used to force people to work in a particular manner.
- However, most models are **not well-aligned** with reality.
- Most hand-made models are **disconnected from reality** and provide only an idealized view on the processes at hand: “paper tigers”.
- Given (a) the interest in process models, (b) the abundance of event data, and (c) the limited quality of hand-made models, it seems worthwhile to relate event data to process models: **process mining!**

# The three main types of process mining: discovery, conformance, and enhancement



# Orthogonal: Perspectives

- The **control-flow perspective** focuses on the control-flow, i.e., the ordering of activities.
- The **organizational perspective** focuses on information about resources hidden in the log, i.e., which actors (e.g., people, systems, roles, and departments) are involved and how are they related.
- The **case perspective** focuses on properties of cases, e.g., cases can also be characterized by the values of the corresponding data elements.
- The **time perspective** is concerned with the timing and frequency of events.

# Starting point: event log

case id	event id	properties				
		timestamp	activity	resource	cost	...
1	35654423	30-12-2010:11.02	register request	Pete	50	...
	35654424	31-12-2010:10.06	examine thoroughly	Sue	400	...
	35654425	05-01-2011:15.12	check ticket	Mike	100	...
	35654426	06-01-2011:11.18	decide	Sara	200	...
	35654427	07-01-2011:14.24	reject request	Pete	200	...
2	35654483	30-12-2010:11.32	register request	Mike	50	...
	35654485	30-12-2010:12.12	check ticket	Mike	100	...
	35654487	30-12-2010:14.16	examine casually			
	35654488	05-01-2011:11.22	decide			
	35654489	08-01-2011:12.05	pay compensation			
3	35654521	30-12-2010:14.32	register request			
	35654522	30-12-2010:15.06	examine casually			
	35654524	30-12-2010:16.34	check ticket			
	35654525	06-01-2011:09.18	decide			
	35654526	06-01-2011:12.18	reinitiate request			
	35654527	06-01-2011:13.06	examine thoroughly			
	35654530	08-01-2011:11.43	check ticket			
	35654531	09-01-2011:09.55	decide			
	35654533	15-01-2011:10.45	pay compensation			
	35654641	06-01-2011:15.02	register request			
4	35654643	07-01-2011:12.06	check ticket			
	35654644	08-01-2011:14.43	examine thoroughly			
	35654645	09-01-2011:12.02	decide			
	35654647	12-01-2011:15.44	reject request			
	35654711	06-01-2011:09.02	register request			
5	35654712	07-01-2011:10.16	examine casually			
	35654714	08-01-2011:11.22	check ticket			
	35654715	10-01-2011:13.28	decide			
	35654716	11-01-2011:16.18	reinitiate request			
	35654718	14-01-2011:14.33	check ticket			
	35654719	16-01-2011:15.50	examine casually			
	35654720	19-01-2011:11.18	decide	Sara	200	...
	35654721	20-01-2011:12.48	reinitiate request	Sara	200	...
	35654722	21-01-2011:09.06	examine casually	Sue	400	...
	35654724	21-01-2011:11.34	check ticket	Pete	100	...
6	35654725	23-01-2011:13.12	decide	Sara	200	...
	35654726	24-01-2011:14.56	reject request	Mike	200	...
	35654871	06-01-2011:15.02	register request	Mike	50	...
	35654873	06-01-2011:16.06	examine casually	Ellen	400	...
	35654874	07-01-2011:16.22	check ticket	Mike	100	...
...	35654875	07-01-2011:16.52	decide	Sara	200	...
	35654877	16-01-2011:11.47	pay compensation	Mike	200	...
	...	...	...	...	...	...

case id	event id	properties				
		timestamp	activity	resource	cost	...
1	35654423	30-12-2010:11.02	register request	Pete	50	...
1	35654424	31-12-2010:10.06	examine thoroughly	Sue	400	...
1	35654425	05-01-2011:15.12	check ticket	Mike	100	...
1	35654426	06-01-2011:11.18	decide	Sara	200	...
1	35654427	07-01-2011:14.24	reject request	Pete	200	...
2	35654483	30-12-2010:11.32	register request	Mike	50	...
2	35654485	30-12-2010:12.12	check ticket	Mike	100	...
2	35654487	30-12-2010:14.16	examine casually	Pete	400	...
2	35654488	05-01-2011:11.22	decide	Sara	200	...
2	35654489	08-01-2011:12.05	pay compensation	Ellen	200	...

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

PAGE 10

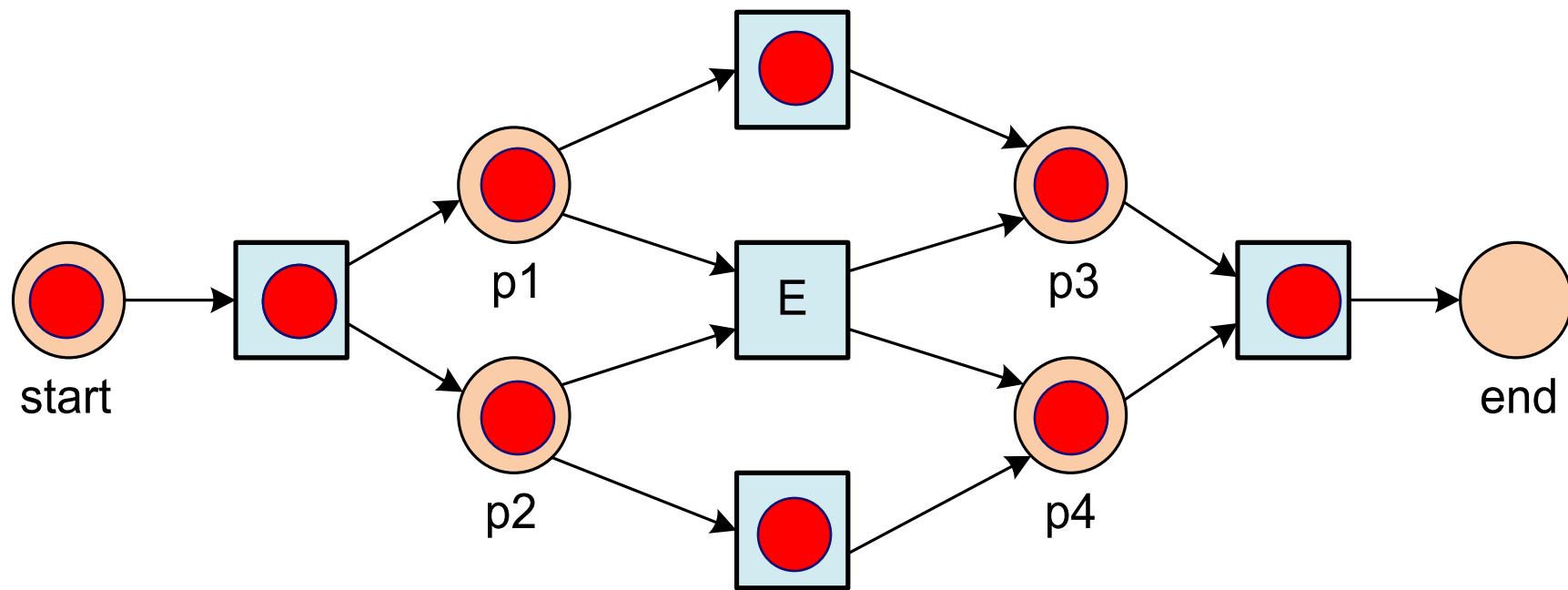
# Simplified event log

case id	event id	properties		
		timestamp	activity	resource
1	35654423	30-12-2010:11.02	register request	Pete
	35654424	31-12-2010:10.06	examine thoroughly	Sue
	35654425	05-01-2011:15.12	check ticket	Mike
	35654426	06-01-2011:11.18	decide	Sara
	35654427	07-01-2011:14.24	reject request	Pete
2	35654483	30-12-2010:11.32	register request	Mike
	35654485	30-12-2010:12.12	check ticket	Mike
	35654487	30-12-2010:14.16	examine casually	Pete
	35654488	05-01-2011:11.22	decide	Sara
	35654489	08-01-2011:12.05	pay compensation	Ellen
3	35654521	30-12-2010:14.32	register request	Pete
	35654522	30-12-2010:15.06	examine casually	Mike
	35654524	30-12-2010:16.34	check ticket	Ellen
	35654525	06-01-2011:09.18	decide	Sara
	35654526	06-01-2011:12.18	reinitiate request	Sara
	35654527	06-01-2011:13.06	examine thoroughly	Sean
	35654530	08-01-2011:11.43	check ticket	Pete
	35654531	09-01-2011:09.55	decide	Sara
	35654533	15-01-2011:10.45	pay compensation	Ellen
	35654641	06-01-2011:15.02	register request	Pete
4	35654643	07-01-2011:12.06	check ticket	Mike
	35654644	08-01-2011:14.43	examine thoroughly	Sean
	35654645	09-01-2011:12.02	decide	Sara
	35654647	12-01-2011:15.44	reject request	Ellen
	35654711	06-01-2011:09.02	register request	Ellen
5	35654712	07-01-2011:10.16	examine casually	Mike
	35654714	08-01-2011:11.22	check ticket	Pete
	35654715	10-01-2011:13.28	decide	Sara
	35654716	11-01-2011:16.18	reinitiate request	Sara
	35654718	14-01-2011:14.33	check ticket	Ellen
	35654719	16-01-2011:15.50	examine casually	Mike
	35654720	19-01-2011:11.18	decide	Sara
	35654721	20-01-2011:12.48	reinitiate request	Sara
	35654722	21-01-2011:09.06	examine casually	Sue
	35654724	21-01-2011:11.34	check ticket	Pete
	35654725	23-01-2011:13.12	decide	Sara
	35654726	24-01-2011:14.56	reject request	Mike
	35654871	06-01-2011:15.02	register request	Mike
6	35654873	06-01-2011:16.06	examine casually	Ellen
	35654874	07-01-2011:16.22	check ticket	Mike
	35654875	07-01-2011:16.52	decide	Sara
	35654877	16-01-2011:11.47	pay compensation	Mike
	...	...	...	...

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

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

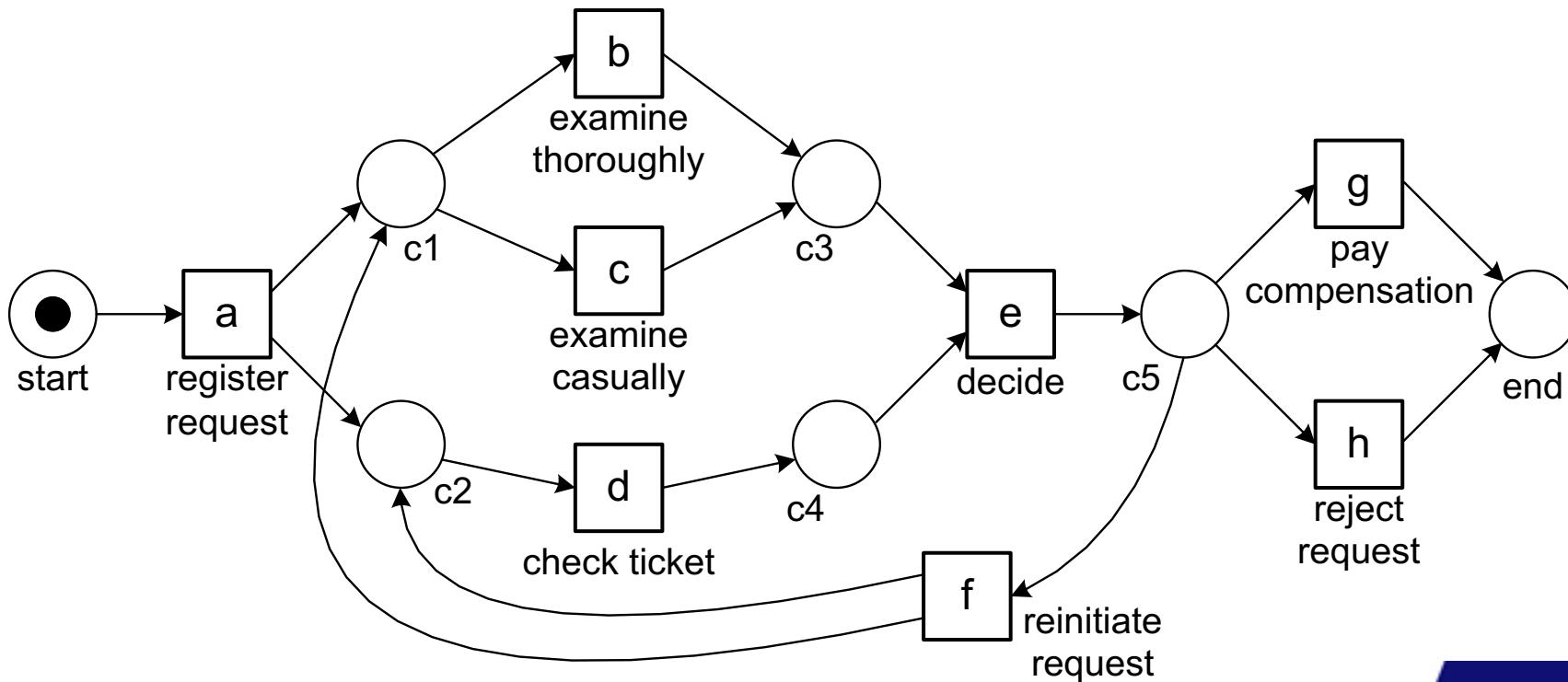
# Petri net semantics



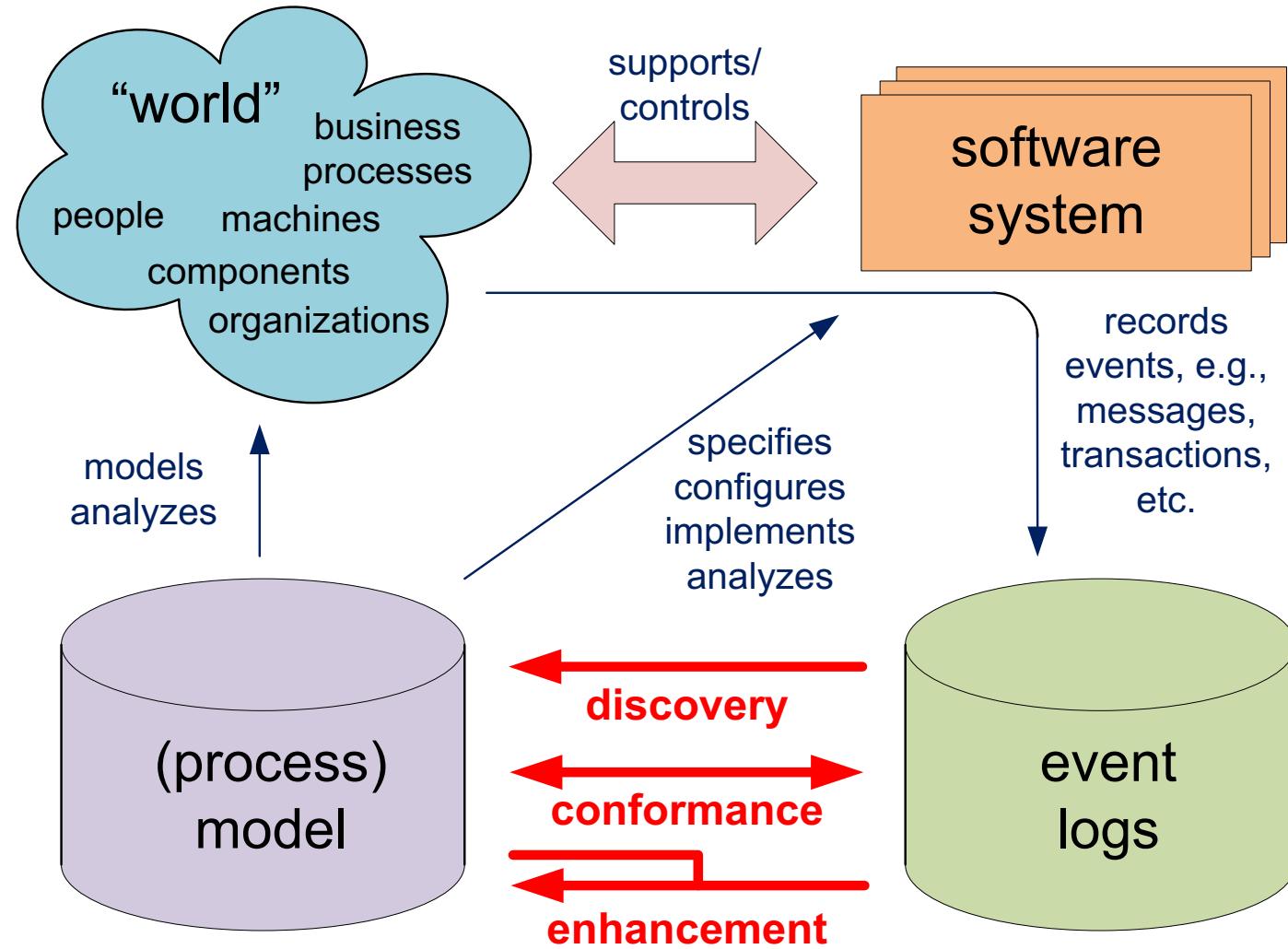
A B C D   A E D   A E D  
A C B D   A B C D   A C B D  
A C B D   A E D   A C B D

# Process discovery

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

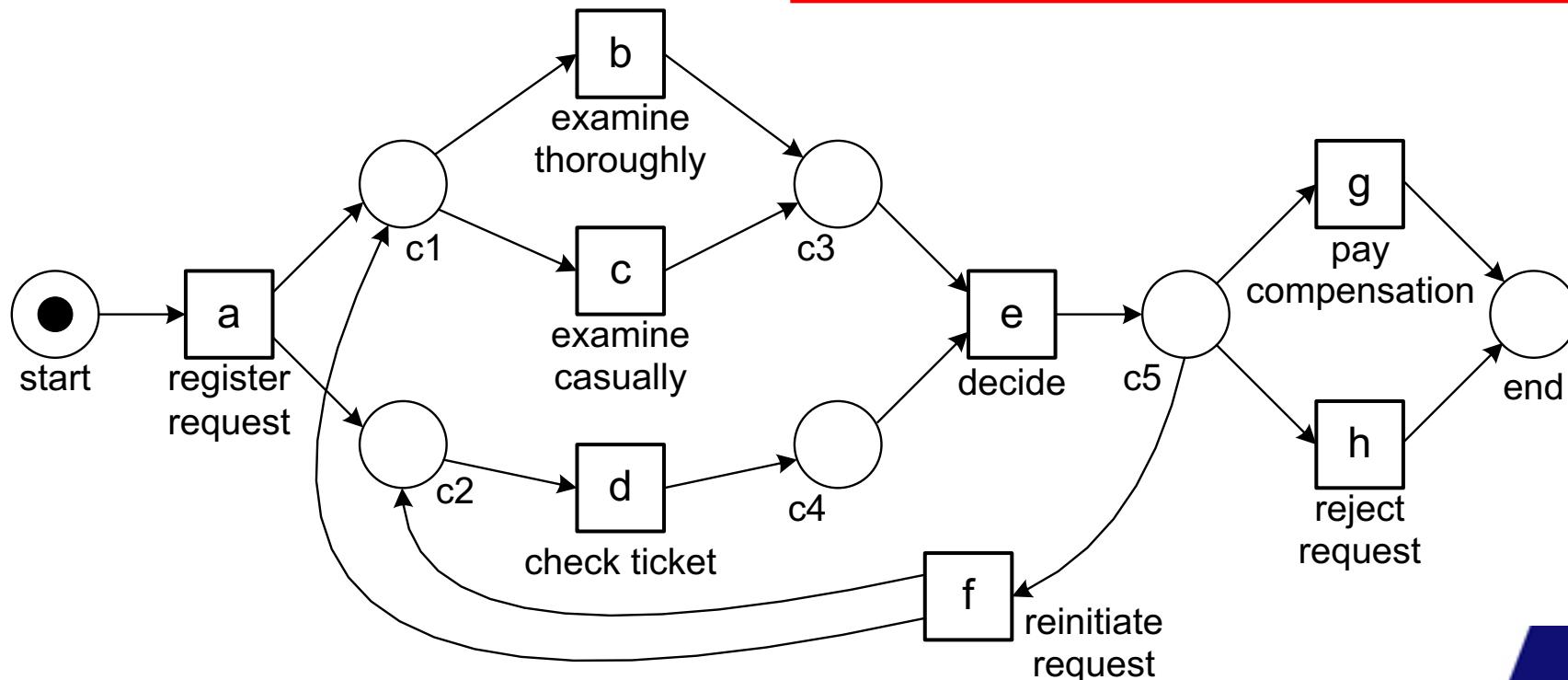


# Beyond discovery: conformance and enhancement

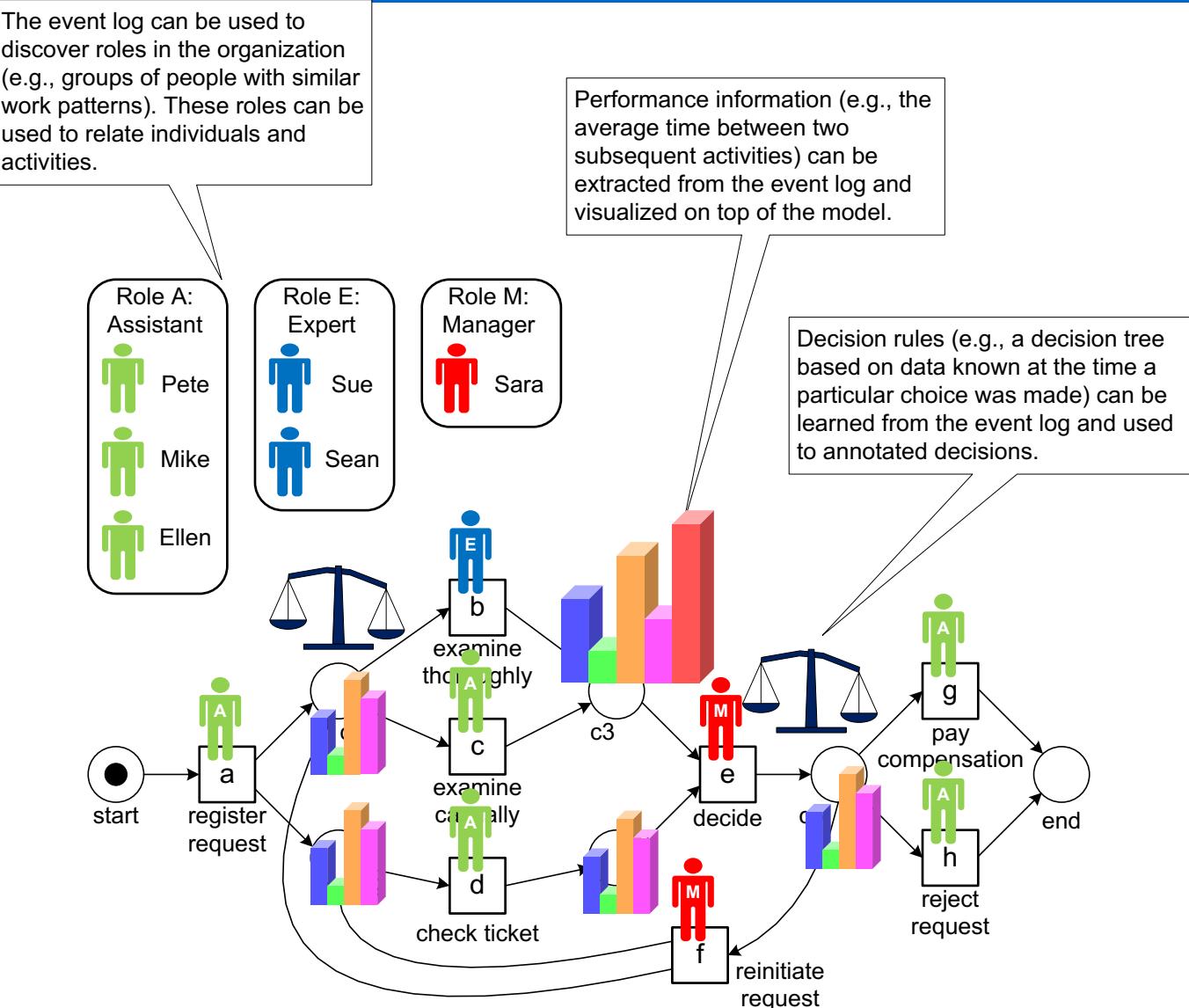


# Another event log

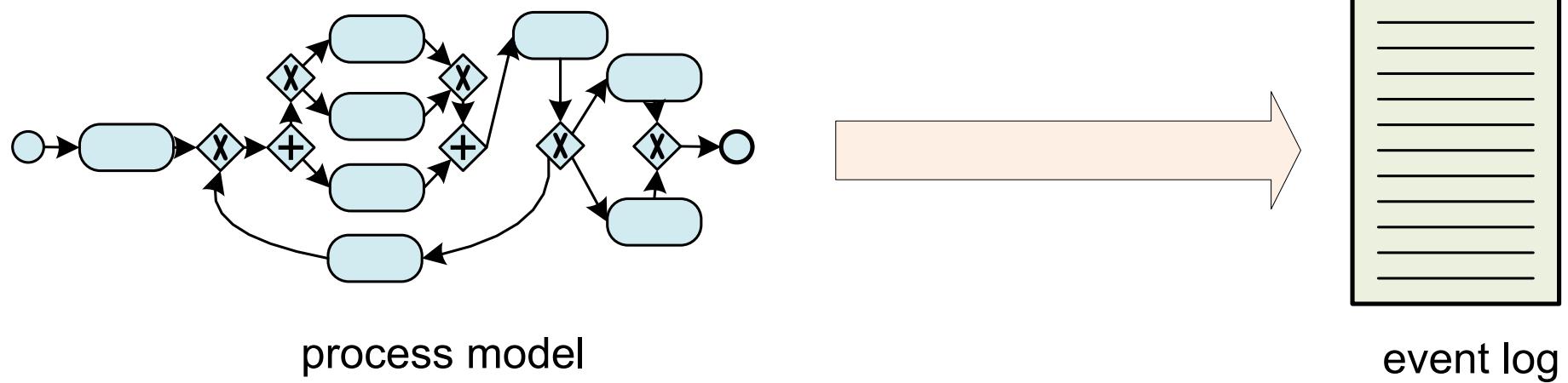
case id	trace
1	$\langle a, b, d, e, h \rangle$
2	$\langle a, d, c, e, g \rangle$
3	$\langle a, c, d, e, f, b, d, e, g \rangle$
4	$\langle a, d, b, e, h \rangle$
5	$\langle a, c, d, e, f, d, c, e, f, c, d, e, h \rangle$
6	$\langle a, c, d, e, g \rangle$
7	$\langle a, b, e, g \rangle$
8	$\langle a, b,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



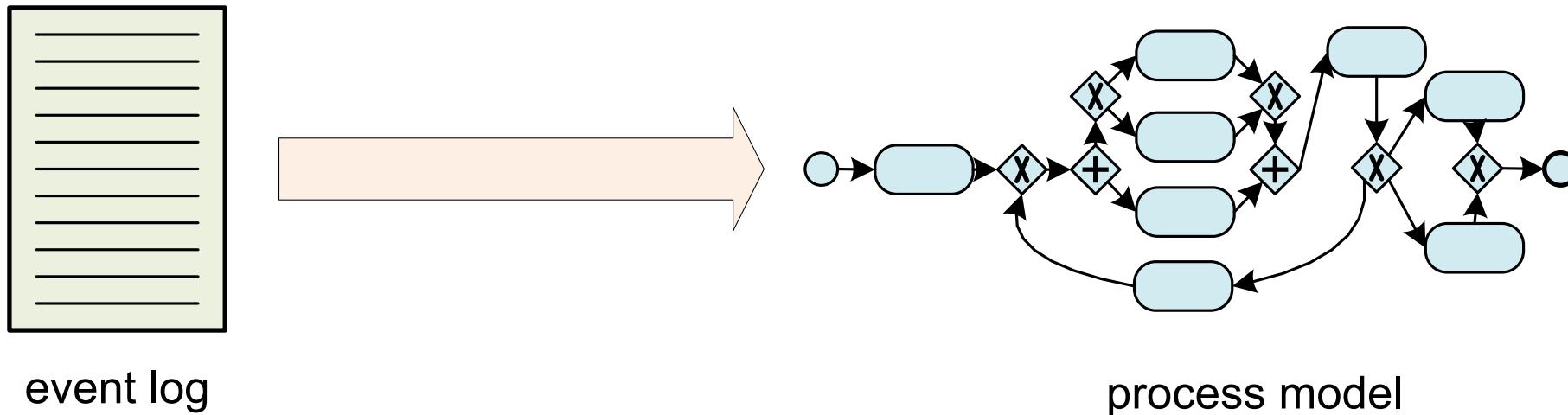
# Extension



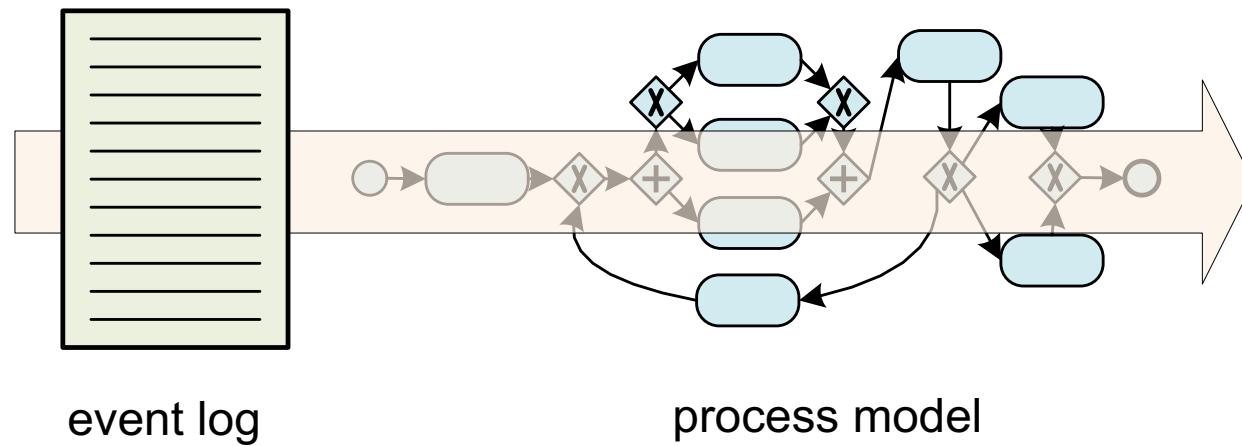
# Play-Out



# Play-In



# Replay



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

# Replay

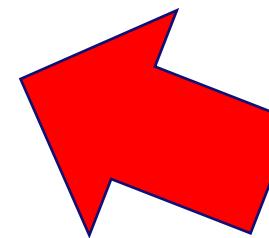
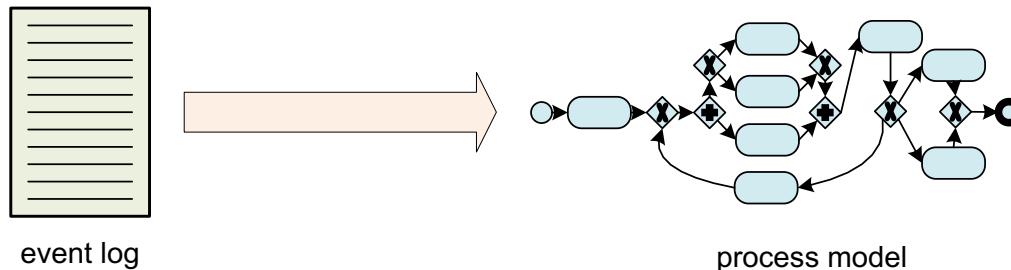
- **Connecting models to real events is crucial!**
- **Possible uses:**
  - Conformance checking
  - Repairing models
  - Extending the model with frequencies and temporal information
  - Constructing predictive models
  - Operational support (prediction, recommendation, etc.)

# Relation between data mining and process mining

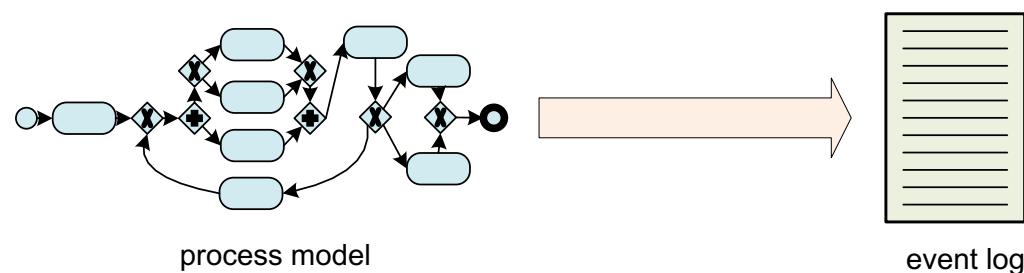
- Process mining: about end-to-end processes.
- Data mining: data-centric and not process-centric.
- Judging the quality of data mining and process mining: many similarities, but also some differences.
- Clearly, process mining techniques can benefit from experiences in the data mining field.

# Process discovery = Play-In

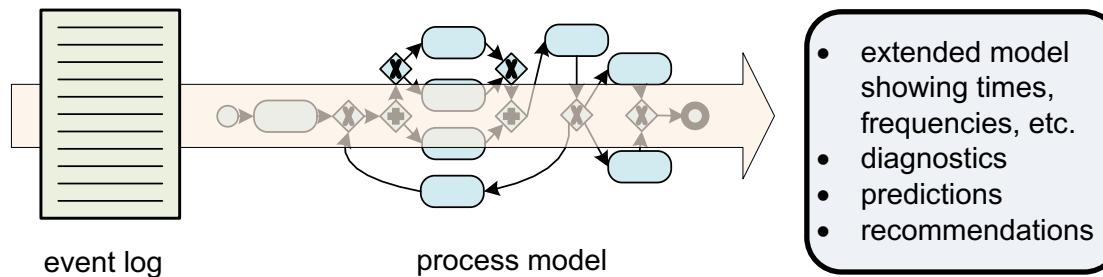
Play-In



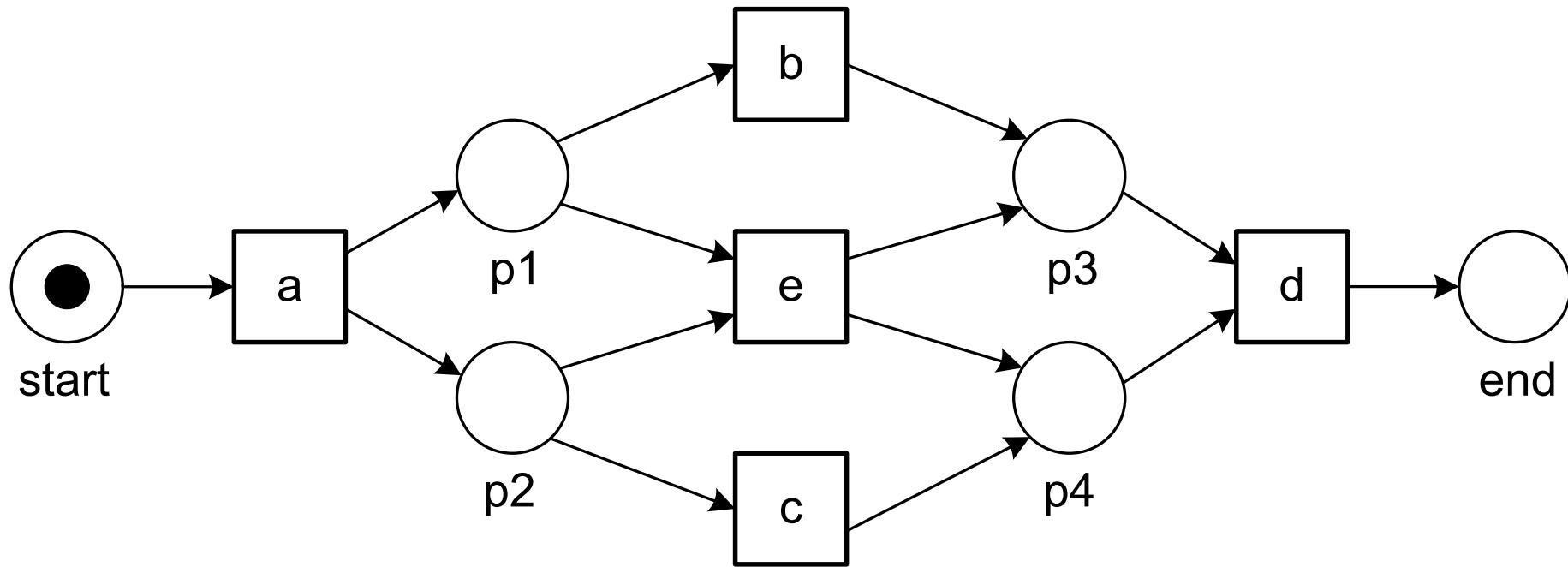
Play-Out



Replay



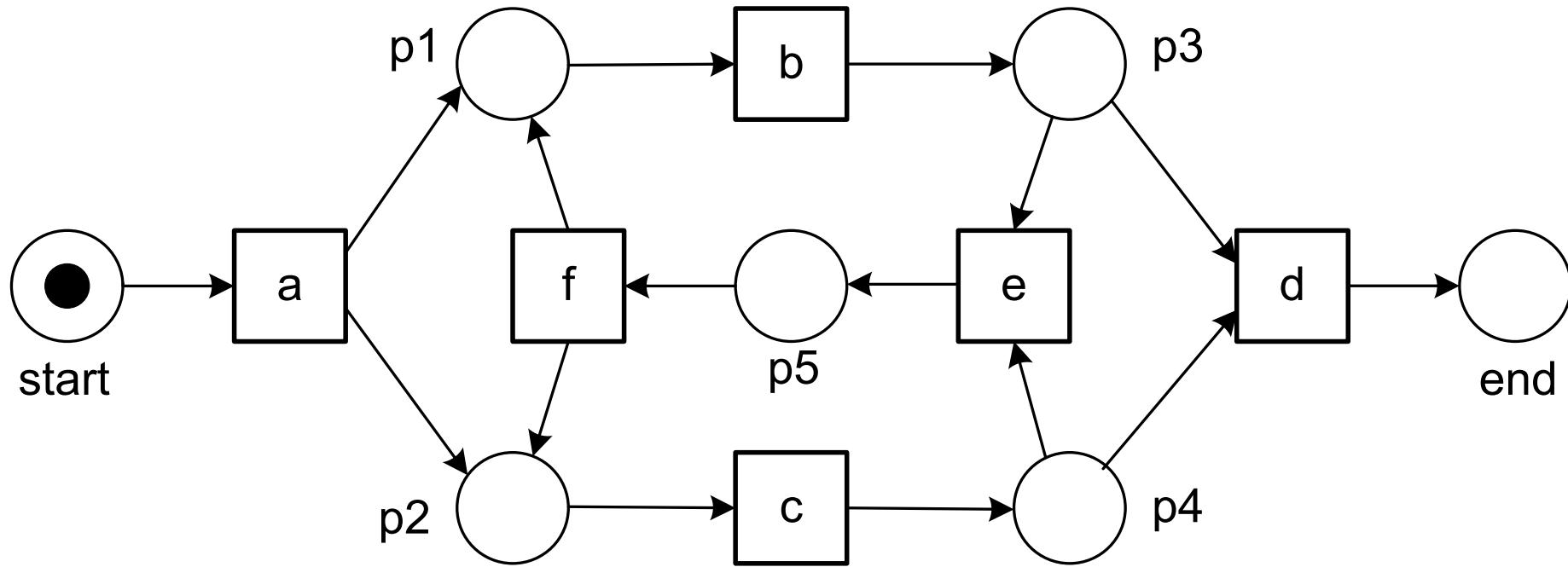
# Example



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

Event log contains all possible traces of model and vice versa.

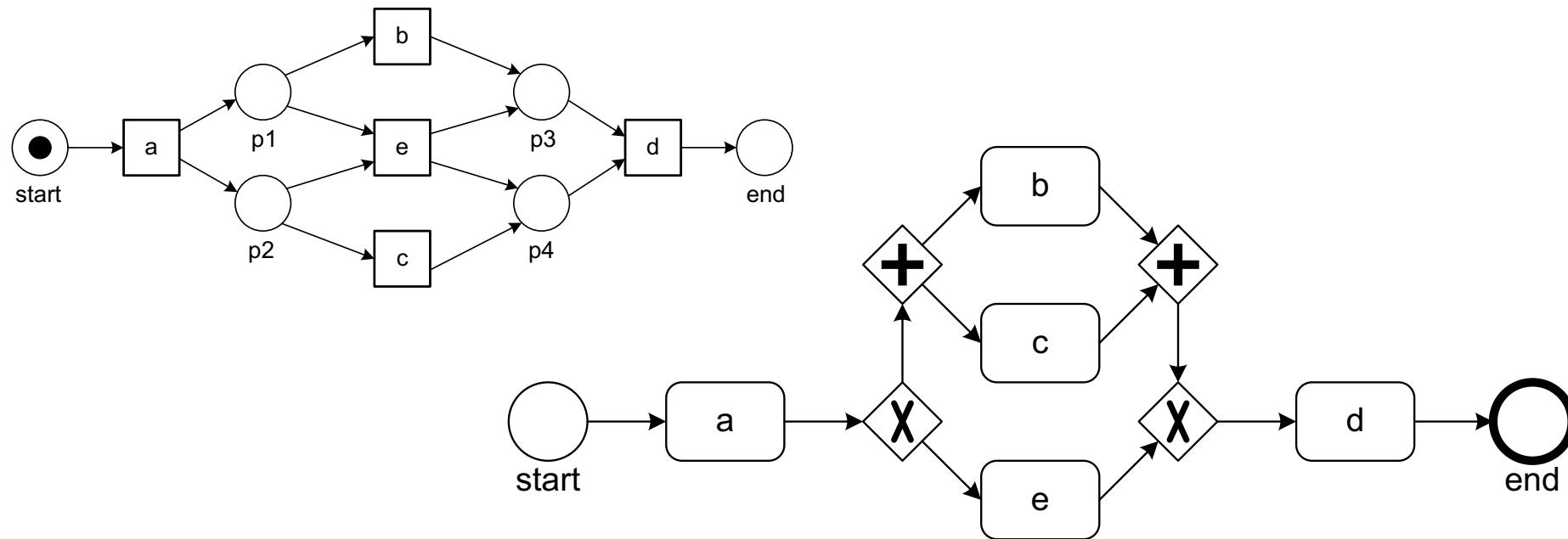
# Another example



$$L_2 = [\langle a, b, c, d \rangle^3, \langle a, c, b, d \rangle^4, \langle a, b, c, e, f, b, c, d \rangle^2, \langle a, b, c, e, f, c, b, d \rangle, \\ \langle a, c, b, e, f, b, c, d \rangle^2, \langle a, c, b, e, f, b, c, e, f, c, b, d \rangle]$$

**Generalization:** event log contains only subset of all possible traces of model.

# Notation is less relevant (e.g. BPMN)



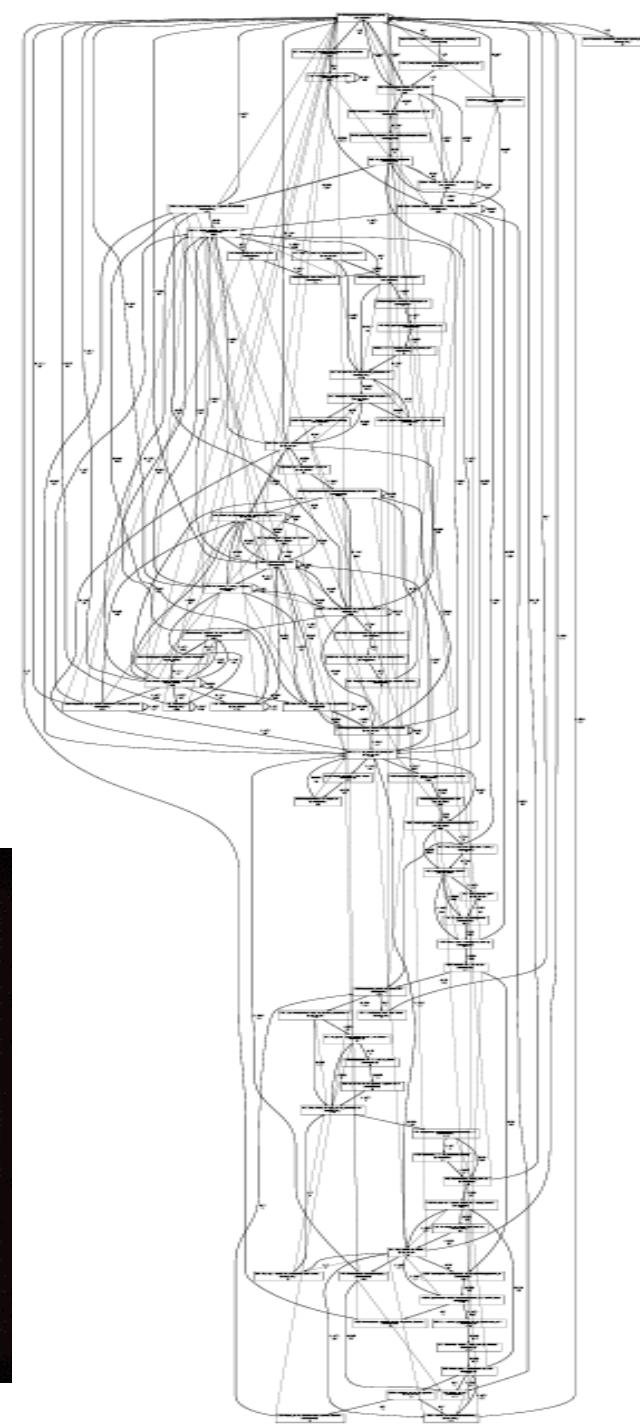
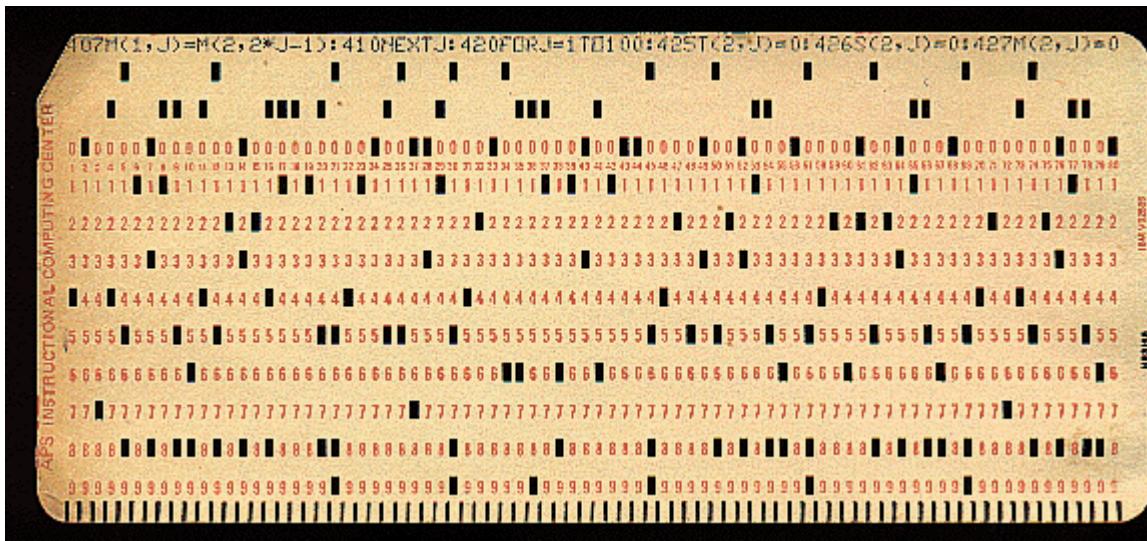
$$L_1 = [\langle a, b, c, d \rangle^3, \langle a, c, b, d \rangle^2, \langle a, e, d \rangle]$$

# Challenge

- In general, there is a trade-off between the following four quality criteria:
  1. **Fitness**: the discovered model should allow for the behavior seen in the event log.
  2. **Precision (avoid underfitting)**: the discovered model should not allow for behavior completely unrelated to what was seen in the event log.
  3. **Generalization (avoid overfitting)**: the discovered model should generalize the example behavior seen in the event log.
  4. **Simplicity**: the discovered model should be as simple as possible.

# Process Discovery: example of algorithm

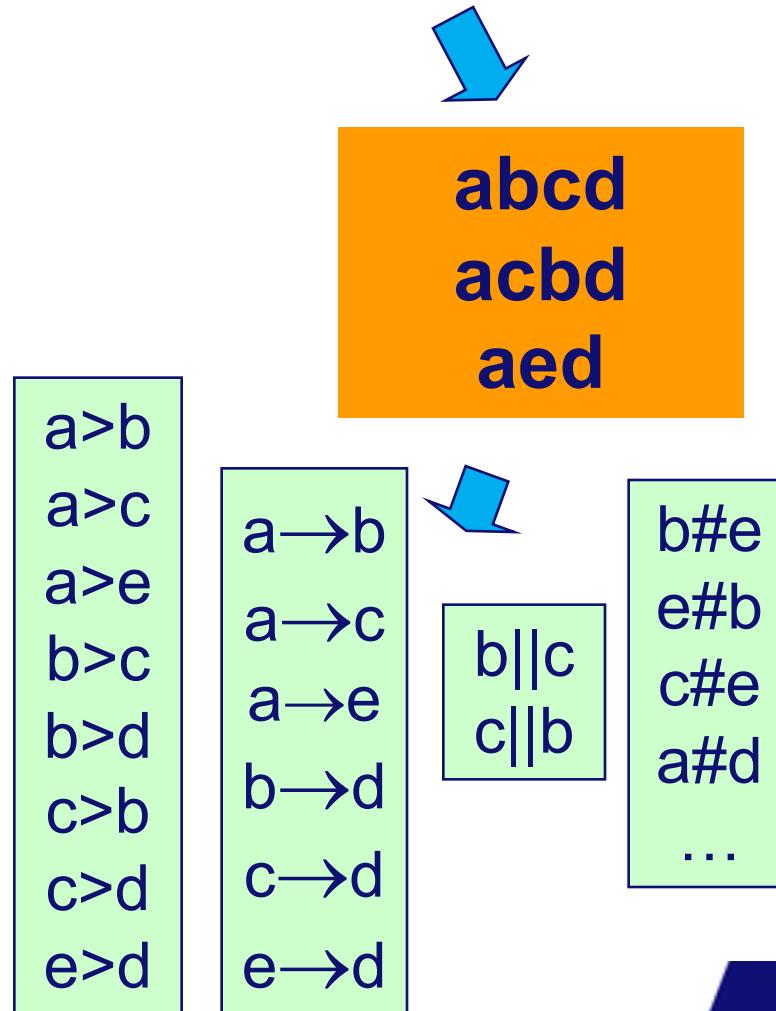
$\alpha$



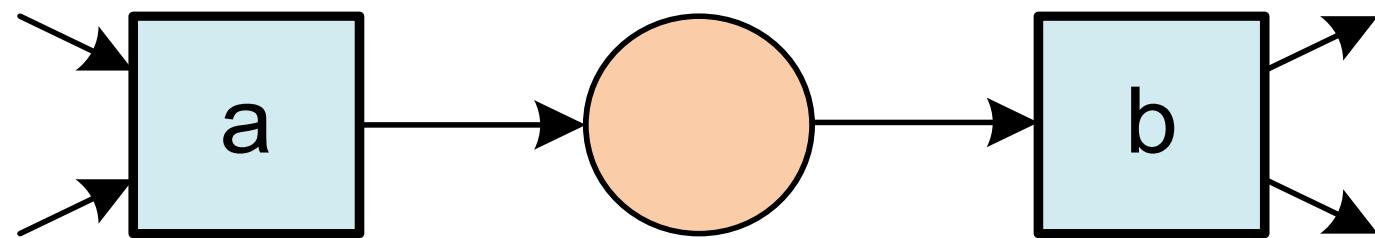
# $>$ , $\rightarrow$ , $||$ , $\#$ relations

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

- Direct succession:  $x>y$  iff for some case  $x$  is directly followed by  $y$ .
- Causality:  $x\rightarrow y$  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  $y>x$ .

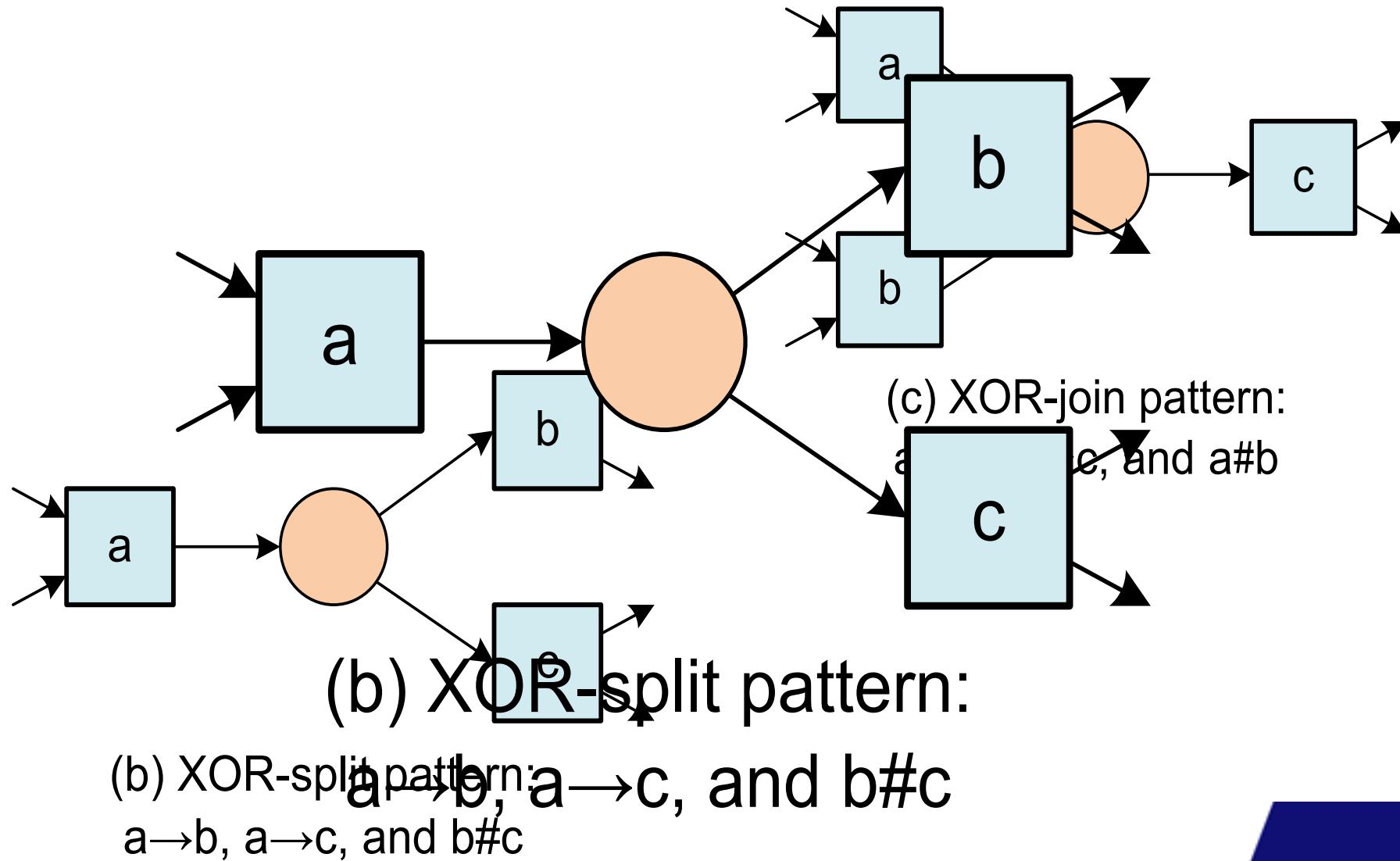


# Basic Idea Used by $\alpha$ Algorithm (1)

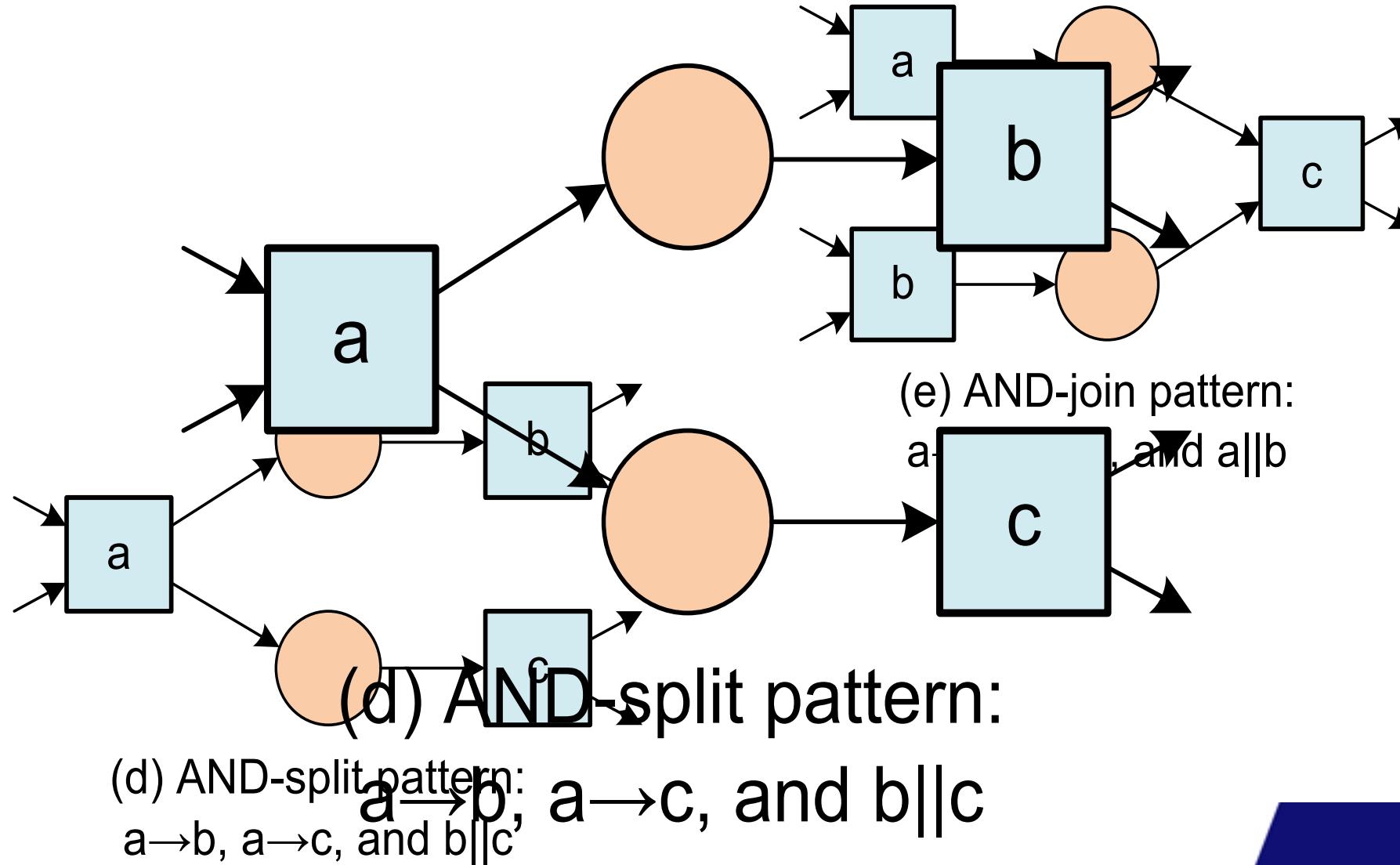


(a) sequence pattern:  $a \rightarrow b$

## Basic Idea Used by $\alpha$ Algorithm (2)

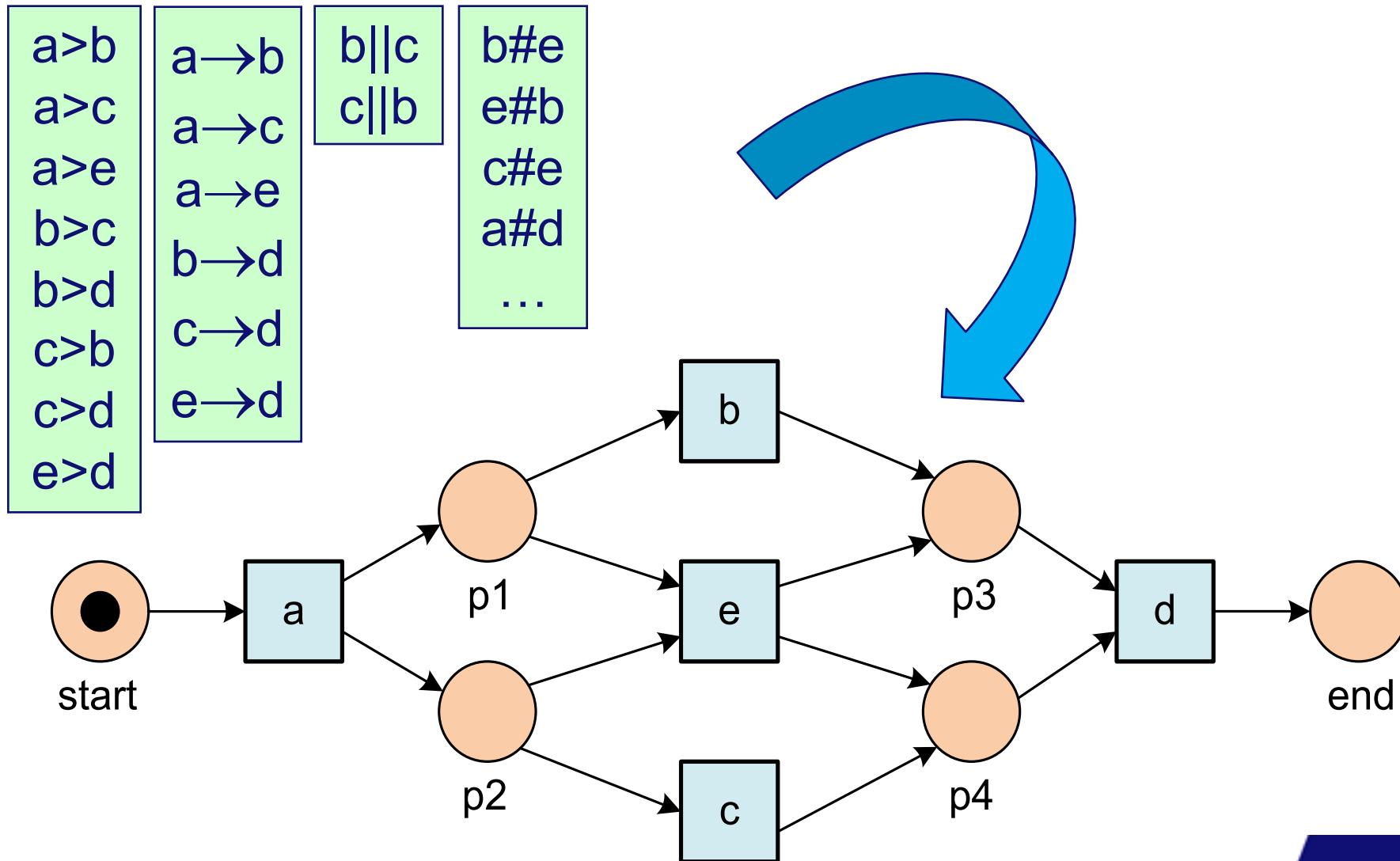


# Basic Idea Used by $\alpha$ Algorithm (3)



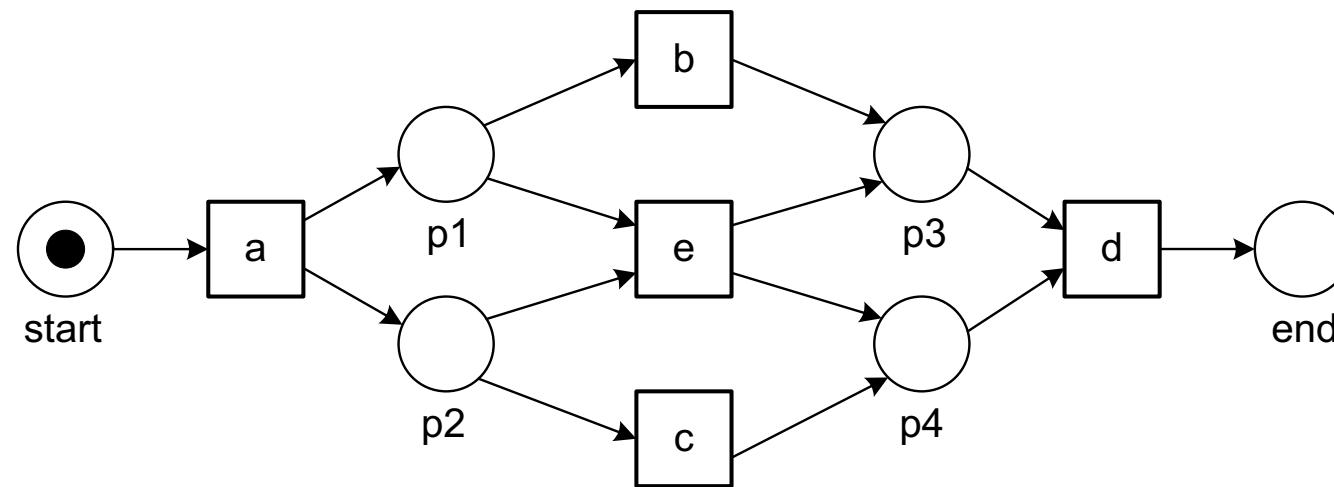
# Example Revisited

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



# Footprint of $L_1$

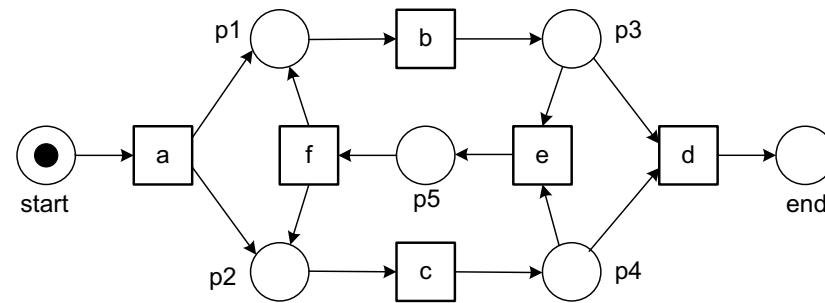
$$L_1 = [\langle a, b, c, d \rangle^3, \langle a, c, b, d \rangle^2, \langle a, e, d \rangle]$$



	<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>	$\leftarrow_{L_1}$	# $L_1$	# $L_1$	$\rightarrow_{L_1}$	# $L_1$

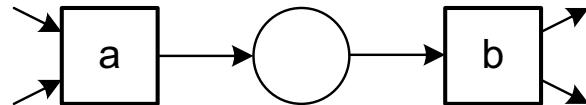
# Footprint of $L_2$

$$L_2 = [\langle a, b, c, d \rangle^3, \langle a, c, b, d \rangle^4, \langle a, b, c, e, f, b, c, d \rangle^2, \langle a, b, c, e, f, c, b, d \rangle, \\ \langle a, c, b, e, f, b, c, d \rangle^2, \langle a, c, b, e, f, b, c, e, f, c, b, d \rangle]$$

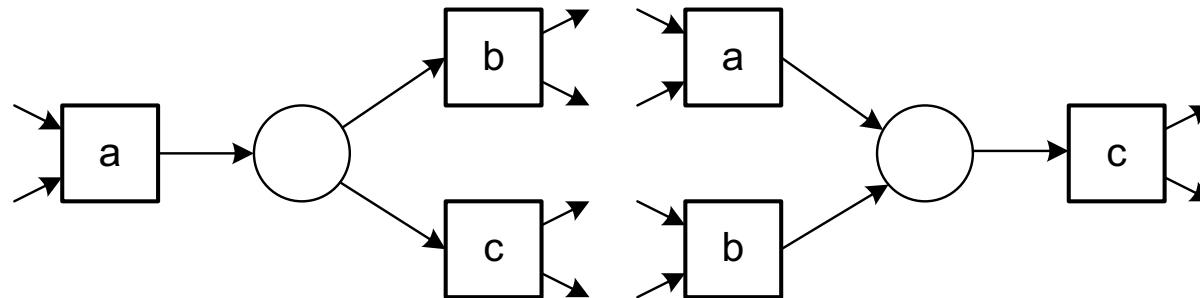


	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
<i>a</i>	#	→	→	#	#	#
<i>b</i>	←	#		→	→	←
<i>c</i>	←		#	→	→	←
<i>d</i>	#	←	←	#	#	#
<i>e</i>	#	←	←	#	#	→
<i>f</i>	#	→	→	#	←	#

# Simple patterns

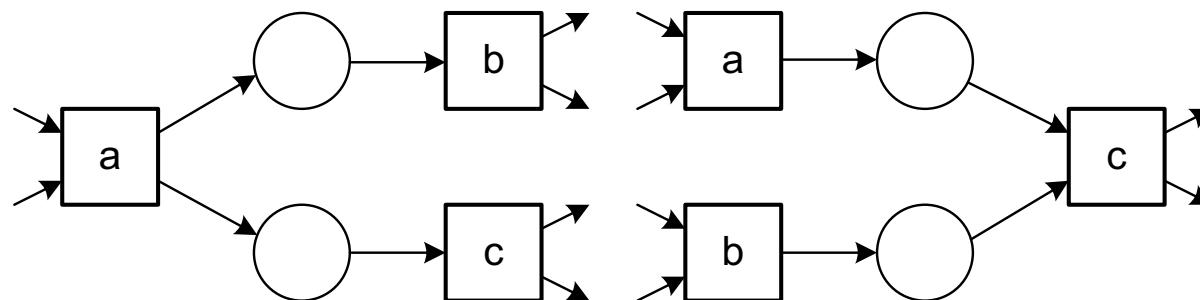


(a) sequence pattern:  $a \rightarrow b$



(b) XOR-split pattern:  
 $a \rightarrow b$ ,  $a \rightarrow c$ , and  $b \# c$

(c) XOR-join pattern:  
 $a \rightarrow c$ ,  $b \rightarrow c$ , and  $a \# b$



(d) AND-split pattern:  
 $a \rightarrow b$ ,  $a \rightarrow c$ , and  $b \parallel c$

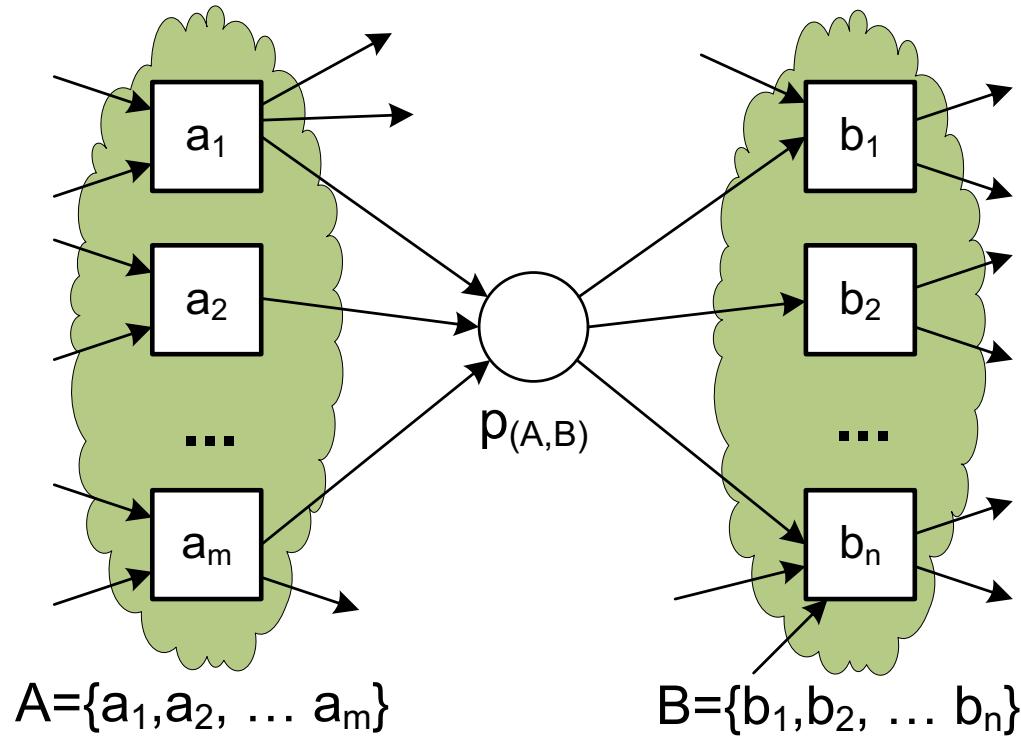
(e) AND-join pattern:  
 $a \rightarrow c$ ,  $b \rightarrow c$ , and  $a \parallel b$

# Algorithm

Let  $L$  be an event log over  $T$ .  $\alpha(L)$  is defined as follows.

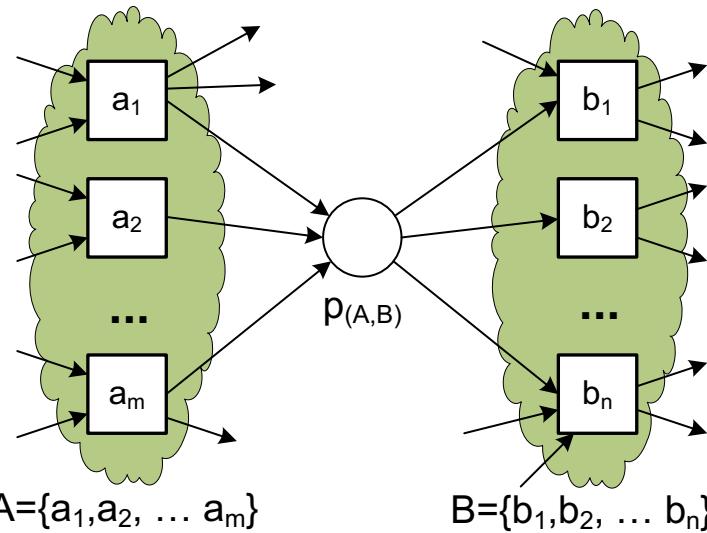
1.  $T_L = \{ t \in T \mid \exists_{\sigma \in L} t \in \sigma \},$
2.  $T_I = \{ t \in T \mid \exists_{\sigma \in L} t = \text{first}(\sigma) \},$
3.  $T_O = \{ t \in T \mid \exists_{\sigma \in L} t = \text{last}(\sigma) \},$
4.  $X_L = \{ (A, B) \mid A \subseteq T_L \wedge A \neq \emptyset \wedge B \subseteq T_L \wedge B \neq \emptyset \wedge \forall_{a \in A} \forall_{b \in B} a \rightarrow_L b \wedge \forall_{a_1, a_2 \in A} a_1 \#_L a_2 \wedge \forall_{b_1, b_2 \in B} b_1 \#_L b_2 \},$
5.  $Y_L = \{ (A, B) \in X_L \mid \forall_{(A', B') \in X_L} A \subseteq A' \wedge B \subseteq B' \Rightarrow (A, B) = (A', B') \},$
6.  $P_L = \{ p_{(A, B)} \mid (A, B) \in Y_L \} \cup \{ i_L, o_L \},$
7.  $F_L = \{ (a, p_{(A, B)}) \mid (A, B) \in Y_L \wedge a \in A \} \cup \{ (p_{(A, B)}, b) \mid (A, B) \in Y_L \wedge b \in B \} \cup \{ (i_L, t) \mid t \in T_I \} \cup \{ (t, o_L) \mid t \in T_O \},$  and
8.  $\alpha(L) = (P_L, T_L, F_L).$

# Key idea: find places



4.  $X_L = \{ (A, B) \mid A \subseteq T_L \wedge A \neq \emptyset \wedge B \subseteq T_L \wedge B \neq \emptyset \wedge \forall a \in A \forall b \in B a \rightarrow_L b \wedge \forall_{a_1, a_2 \in A} a_1 \#_L a_2 \wedge \forall_{b_1, b_2 \in B} b_1 \#_L b_2 \},$
5.  $Y_L = \{ (A, B) \in X_L \mid \forall_{(A', B') \in X_L} A \subseteq A' \wedge B \subseteq B' \Rightarrow (A, B) = (A', B') \},$

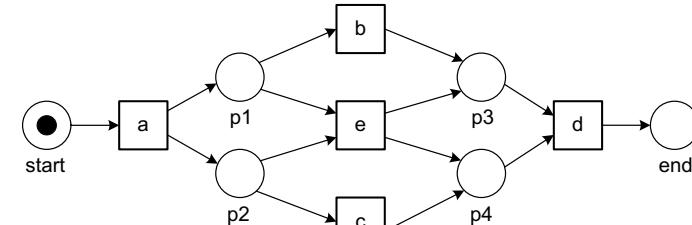
# Places as footprints



	$a_1$	$a_2$	$\dots$	$a_m$	$b_1$	$b_2$	$\dots$	$b_n$
$a_1$	#	#	$\dots$	#	→	→	$\dots$	→
$a_2$	#	#	$\dots$	#	→	→	$\dots$	→
$\dots$								
$a_m$	#	#	$\dots$	#	→	→	$\dots$	→
$b_1$	←	←	$\dots$	←	#	#	$\dots$	#
$b_2$	←	←	$\dots$	←	#	#	$\dots$	#
$\dots$								
$b_n$	←	←	$\dots$	←	#	#	$\dots$	#

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

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



$$X_{L_1} = \{(\{a\}, \{b\}), (\{a\}, \{c\}), (\{a\}, \{e\}), (\{a\}, \{b, e\}), (\{a\}, \{c, e\}), \\ (\{b\}, \{d\}), (\{c\}, \{d\}), (\{e\}, \{d\}), (\{b, e\}, \{d\}), (\{c, e\}, \{d\})\}$$

$$Y_{L_1} = \{(\{a\}, \{b, e\}), (\{a\}, \{c, e\}), (\{b, e\}, \{d\}), (\{c, e\}, \{d\})\}$$

# Another event log $L_3$

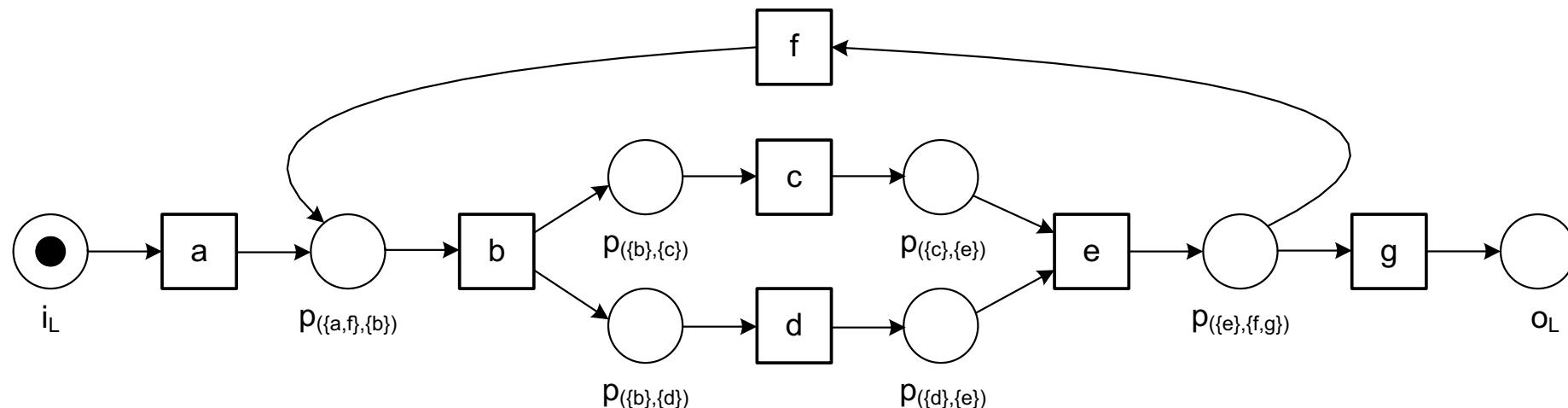
$$L_3 = [\langle a, b, c, d, e, f, b, d, c, e, g \rangle, \\ \langle a, b, d, c, e, g \rangle^2, \\ \langle a, b, c, d, e, f, b, c, d, e, f, b, d, c, e, g \rangle]$$

	$a$	$b$	$c$	$d$	$e$	$f$	$g$
$a$	#	$\rightarrow$	#	#	#	#	#
$b$	$\leftarrow$	#	$\rightarrow$	$\rightarrow$	#	$\leftarrow$	#
$c$	#	$\leftarrow$	#	$\parallel$	$\rightarrow$	#	#
$d$	#	$\leftarrow$	$\parallel$	#	$\rightarrow$	#	#
$e$	#	#	$\leftarrow$	$\leftarrow$	#	$\rightarrow$	$\rightarrow$
$f$	#	$\rightarrow$	#	#	$\leftarrow$	#	#
$g$	#	#	#	#	$\leftarrow$	#	#

# Model for $L_3$

	$a$	$b$	$c$	$d$	$e$	$f$	$g$
$a$	#	$\rightarrow$	#	#	#	#	#
$b$	$\leftarrow$	#	$\rightarrow$	$\rightarrow$	#	$\leftarrow$	#
$c$	#	$\leftarrow$	#	$\parallel$	$\rightarrow$	#	#
$d$	#	$\leftarrow$	$\parallel$	#	$\rightarrow$	#	#
$e$	#	#	$\leftarrow$	$\leftarrow$	#	$\rightarrow$	$\rightarrow$
$f$	#	$\rightarrow$	#	#	$\leftarrow$	#	#
$g$	#	#	#	#	$\leftarrow$	#	#

$$L_3 = [\langle a, b, c, d, e, f, b, d, c, e, g \rangle, \\ \langle a, b, d, c, e, g \rangle^2, \\ \langle a, b, c, d, e, f, b, c, d, e, f, b, d, c, e, g \rangle]$$



# Event log $L_5$

$$L_5 = [\langle a, b, e, f \rangle^2, \langle a, b, e, c, d, b, f \rangle^3, \langle a, b, c, e, d, b, f \rangle^2, \\ \langle a, b, c, d, e, b, f \rangle^4, \langle a, e, b, c, d, b, f \rangle^3]$$

	$a$	$b$	$c$	$d$	$e$	$f$
$a$	#	$\rightarrow$	#	#	$\rightarrow$	#
$b$	$\leftarrow$	#	$\rightarrow$	$\leftarrow$	$\parallel$	$\rightarrow$
$c$	#	$\leftarrow$	#	$\rightarrow$	$\parallel$	#
$d$	#	$\rightarrow$	$\leftarrow$	#	$\parallel$	#
$e$	$\leftarrow$	$\parallel$	$\parallel$	$\parallel$	#	$\rightarrow$
$f$	#	$\leftarrow$	#	#	$\leftarrow$	#

$$T_L = \{a, b, c, d, e, f\}$$

$$T_I = \{a\}$$

$$T_I = \{f\}$$

$$X_L = \{(\{a\}, \{b\}), (\{a\}, \{e\}), (\{b\}, \{c\}), (\{b\}, \{f\}), (\{c\}, \{d\}), \\ (\{d\}, \{b\}), (\{e\}, \{f\}), (\{a, d\}, \{b\}), (\{b\}, \{c, f\})\}$$

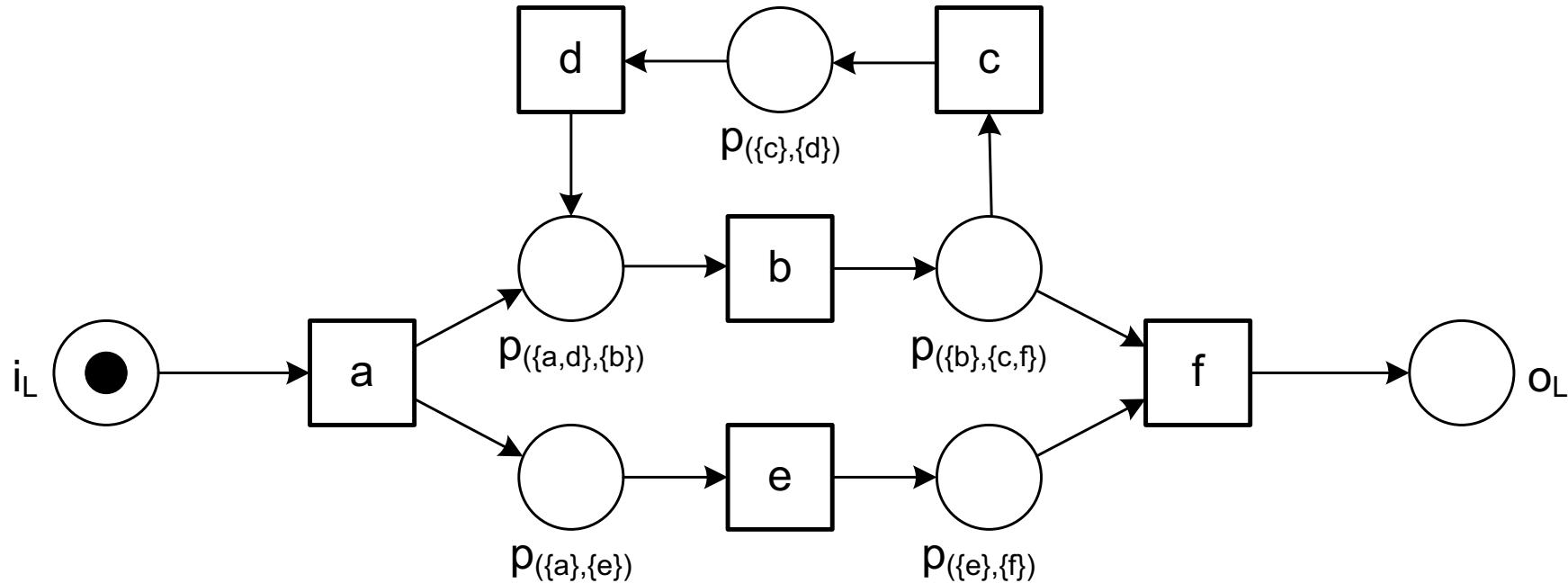
$$Y_L = \{(\{a\}, \{e\}), (\{c\}, \{d\}), (\{e\}, \{f\}), (\{a, d\}, \{b\}), (\{b\}, \{c, f\})\}$$

$$P_L = \{p_{(\{a\}, \{e\})}, p_{(\{c\}, \{d\})}, p_{(\{e\}, \{f\})}, p_{(\{a, d\}, \{b\})}, p_{(\{b\}, \{c, f\})}, i_L, o_L\}$$

$$F_L = \{(a, p_{(\{a\}, \{e\})}), (p_{(\{a\}, \{e\})}, e), (c, p_{(\{c\}, \{d\})}), (p_{(\{c\}, \{d\})}, d), \\ (e, p_{(\{e\}, \{f\})}), (p_{(\{e\}, \{f\})}, f), (a, p_{(\{a, d\}, \{b\})}), (d, p_{(\{a, d\}, \{b\})}), \\ (p_{(\{a, d\}, \{b\})}, b), (b, p_{(\{b\}, \{c, f\})}), (p_{(\{b\}, \{c, f\})}, c), (p_{(\{b\}, \{c, f\})}, f), \\ (i_L, a), (f, o_L)\}$$

$$\alpha(L) = (P_L, T_L, F_L)$$

# Discovered model

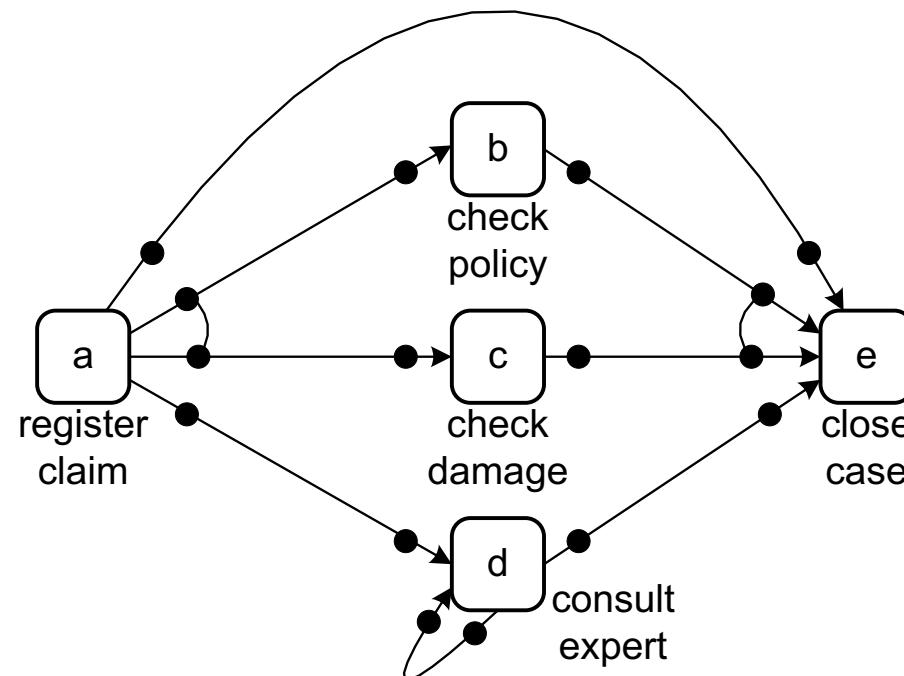


$$X_L = \{(\{a\}, \{b\}), (\{a\}, \{e\}), (\{b\}, \{c\}), (\{b\}, \{f\}), (\{c\}, \{d\}), (\{d\}, \{b\}), (\{e\}, \{f\}), (\{a, d\}, \{b\}), (\{b\}, \{c, f\})\}$$

$$Y_L = \{(\{a\}, \{e\}), (\{c\}, \{d\}), (\{e\}, \{f\}), (\{a, d\}, \{b\}), (\{b\}, \{c, f\})\}$$

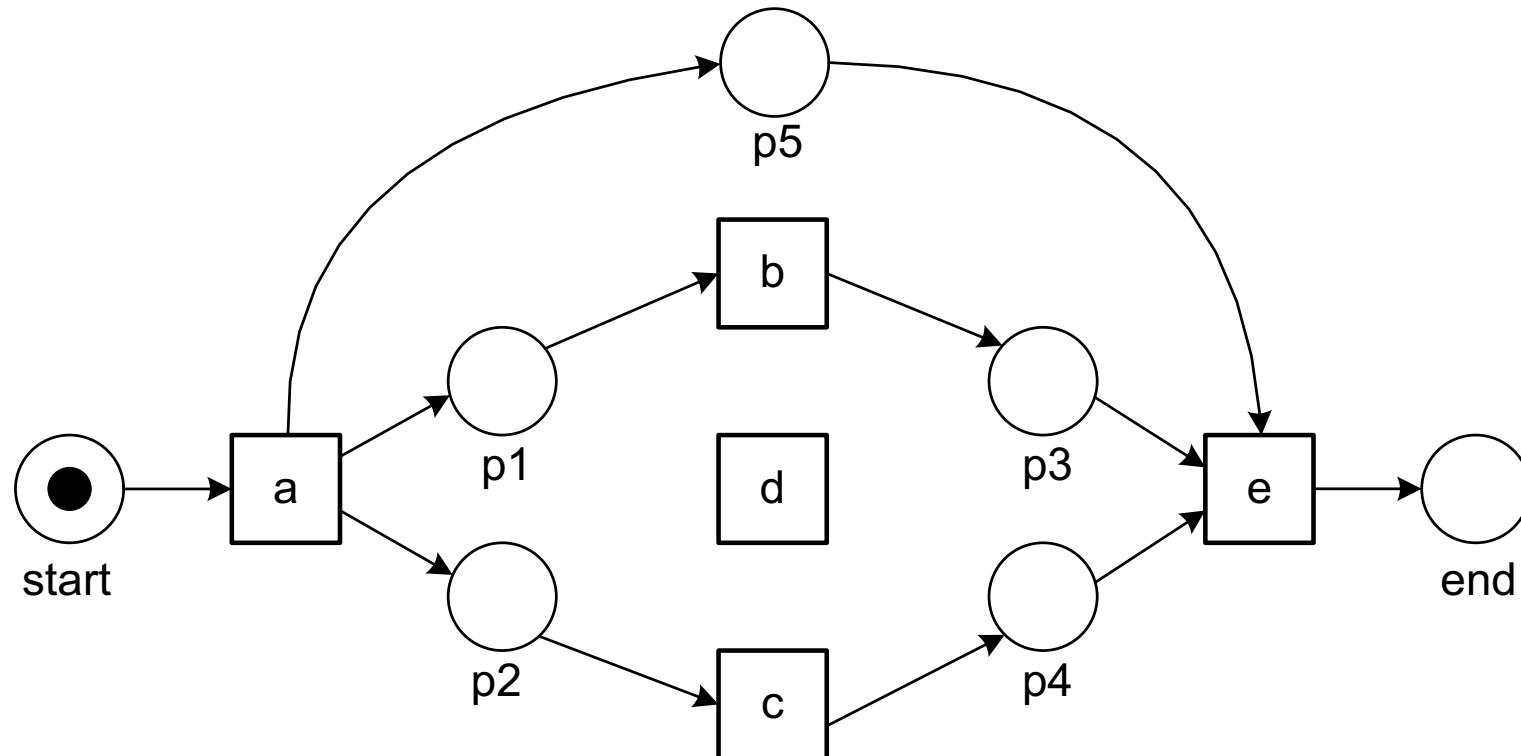
# Heuristic mining

- To deal with noise and incompleteness.
- To have a better representational bias than the  $\alpha$  algorithm (AND/XOR/OR/skip).
- Uses C-nets.



# Example log; problem α algorithm

$$L = [\langle a, e \rangle^5, \langle a, b, c, e \rangle^{10}, \langle a, c, b, e \rangle^{10}, \langle a, b, e \rangle^1, \langle a, c, e \rangle^1, \\ \langle a, d, e \rangle^{10}, \langle a, d, d, e \rangle^2, \langle a, d, d, d, e \rangle^1]$$



# Taking into account frequencies

$$L = [\langle a, e \rangle^5, \langle a, b, c, e \rangle^{10}, \langle a, c, b, e \rangle^{10}, \langle a, b, e \rangle^1, \langle a, c, e \rangle^1, \\ \langle a, d, e \rangle^{10}, \langle a, d, d, e \rangle^2, \langle a, d, d, d, e \rangle^1]$$

$$|a >_L b| = \sum_{\sigma \in L} L(\sigma) \times |\{1 \leq i < |\sigma| \mid \sigma(i) = a \wedge \sigma(i+1) = b\}|$$

$  >_L  $	$a$	$b$	$c$	$d$	$e$
$a$	0	11	11	13	5
$b$	0	0	10	0	11
$c$	0	10	0	0	11
$d$	0	0	0	4	13
$e$	0	0	0	0	0

# Dependency measure

$$|a >_L b| = \sum_{\sigma \in L} L(\sigma) \times |\{1 \leq i < |\sigma| \mid \sigma(i) = a \wedge \sigma(i+1) = b\}|$$

$|a \Rightarrow_L b|$  is the value of the dependency relation between  $a$  and  $b$ :

$$|a \Rightarrow_L b| = \begin{cases} \frac{|a >_L b| - |b >_L a|}{|a >_L b| + |b >_L a| + 1} & \text{if } a \neq b \\ \frac{|a >_L a|}{|a >_L a| + 1} & \text{if } a = b \end{cases}$$

# Example

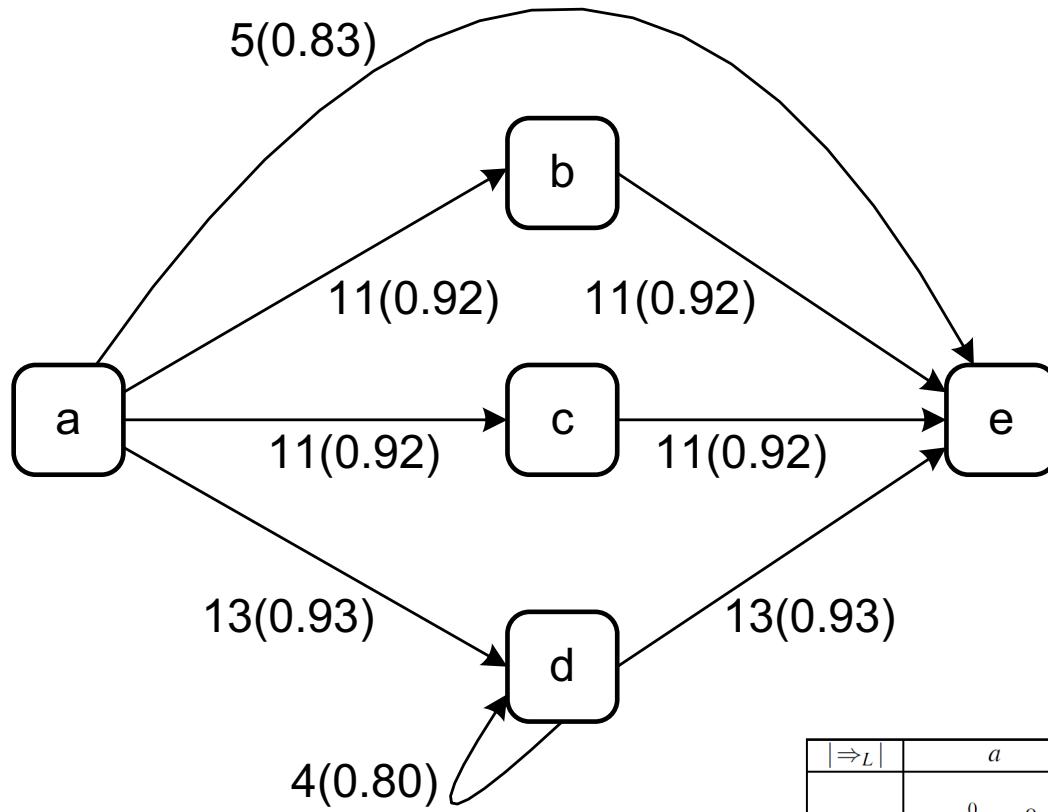
$  \Rightarrow_L  $	$a$	$b$	$c$	$d$	$e$
$a$	$\frac{0}{0+1} = 0$	$\frac{11-0}{11+0+1} = 0.92$	$\frac{11-0}{11+0+1} = 0.92$	$\frac{13-0}{13+0+1} = 0.93$	$\frac{5-0}{5+0+1} = 0.83$
$b$	$\frac{0-11}{0+11+1} = -0.92$	$\frac{0}{0+1} = 0$	$\frac{10-10}{10+10+1} = 0$	$\frac{0-0}{0+0+1} = 0$	$\frac{11-0}{11+0+1} = 0.92$
$c$	$\frac{0-11}{0+11+1} = -0.92$	$\frac{10-10}{10+10+1} = 0$	$\frac{0}{0+1} = 0$	$\frac{0-0}{0+0+1} = 0$	$\frac{11-0}{11+0+1} = 0.92$
$d$	$\frac{0-13}{0+13+1} = -0.93$	$\frac{0-0}{0+0+1} = 0$	$\frac{0-0}{0+0+1} = 0$	$\frac{4}{4+1} = 0.80$	$\frac{13-0}{13+0+1} = 0.93$
$e$	$\frac{0-5}{0+5+1} = -0.83$	$\frac{0-11}{0+11+1} = -0.92$	$\frac{0-11}{0+11+1} = -0.92$	$\frac{0-13}{0+13+1} = -0.93$	$\frac{0}{0+1} = 0$

$|a \Rightarrow_L b|$  is the value of the dependency relation between  $a$  and  $b$ :

$$|a \Rightarrow_L b| = \begin{cases} \frac{|a >_L b| - |b >_L a|}{|a >_L b| + |b >_L a| + 1} & \text{if } a \neq b \\ \frac{|a >_L a|}{|a >_L a| + 1} & \text{if } a = b \end{cases}$$

$  >_L  $	$a$	$b$	$c$	$d$	$e$
$a$	0	11	11	13	5
$b$	0	0	10	0	11
$c$	0	10	0	0	11
$d$	0	0	0	4	13
$e$	0	0	0	0	0

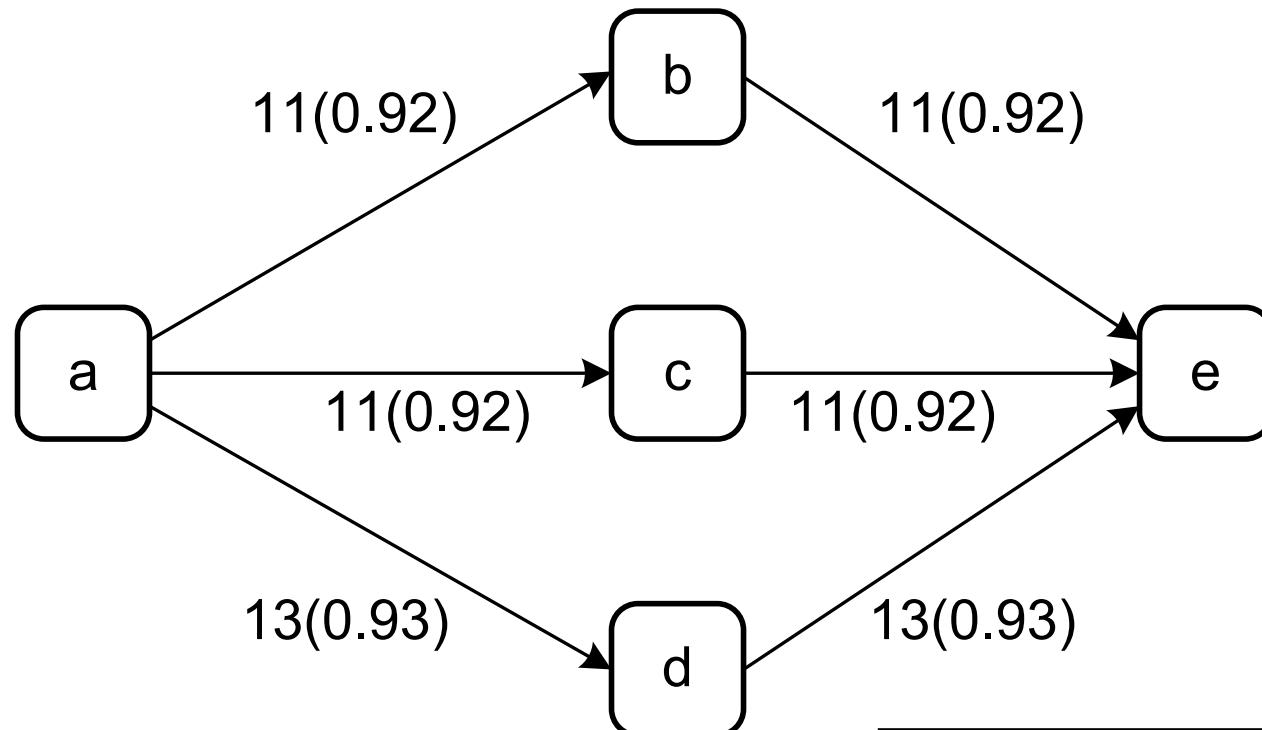
# Lower threshold (2 direct successions and a dependency of at least 0.7)



$ >_L $	a	b	c	d	e
a	0	11	11	13	5
b	0	0	10	0	11
c	0	10	0	0	11
d	0	0	0	4	13
e	0	0	0	0	0

$ \Rightarrow_L $	a	b	c	d	e
a	$\frac{0}{0+1} = 0$	$\frac{11-0}{11+0+1} = 0.92$	$\frac{11-0}{11+0+1} = 0.92$	$\frac{13-0}{13+0+1} = 0.93$	$\frac{5-0}{5+0+1} = 0.83$
b	$\frac{0-11}{0+11+1} = -0.92$	$\frac{0}{0+1} = 0$	$\frac{10-10}{10+10+1} = 0$	$\frac{0-0}{0+0+1} = 0$	$\frac{11-0}{11+0+1} = 0.92$
c	$\frac{0-11}{0+11+1} = -0.92$	$\frac{10-10}{10+10+1} = 0$	$\frac{0}{0+1} = 0$	$\frac{0-0}{0+0+1} = 0$	$\frac{11-0}{11+0+1} = 0.92$
d	$\frac{0-13}{0+13+1} = -0.93$	$\frac{0-0}{0+0+1} = 0$	$\frac{0-0}{0+0+1} = 0$	$\frac{4}{4+1} = 0.80$	$\frac{13-0}{13+0+1} = 0.93$
e	$\frac{0-5}{0+5+1} = -0.83$	$\frac{0-11}{0+11+1} = -0.92$	$\frac{0-11}{0+11+1} = -0.92$	$\frac{0-13}{0+13+1} = -0.93$	$\frac{0}{0+1} = 0$

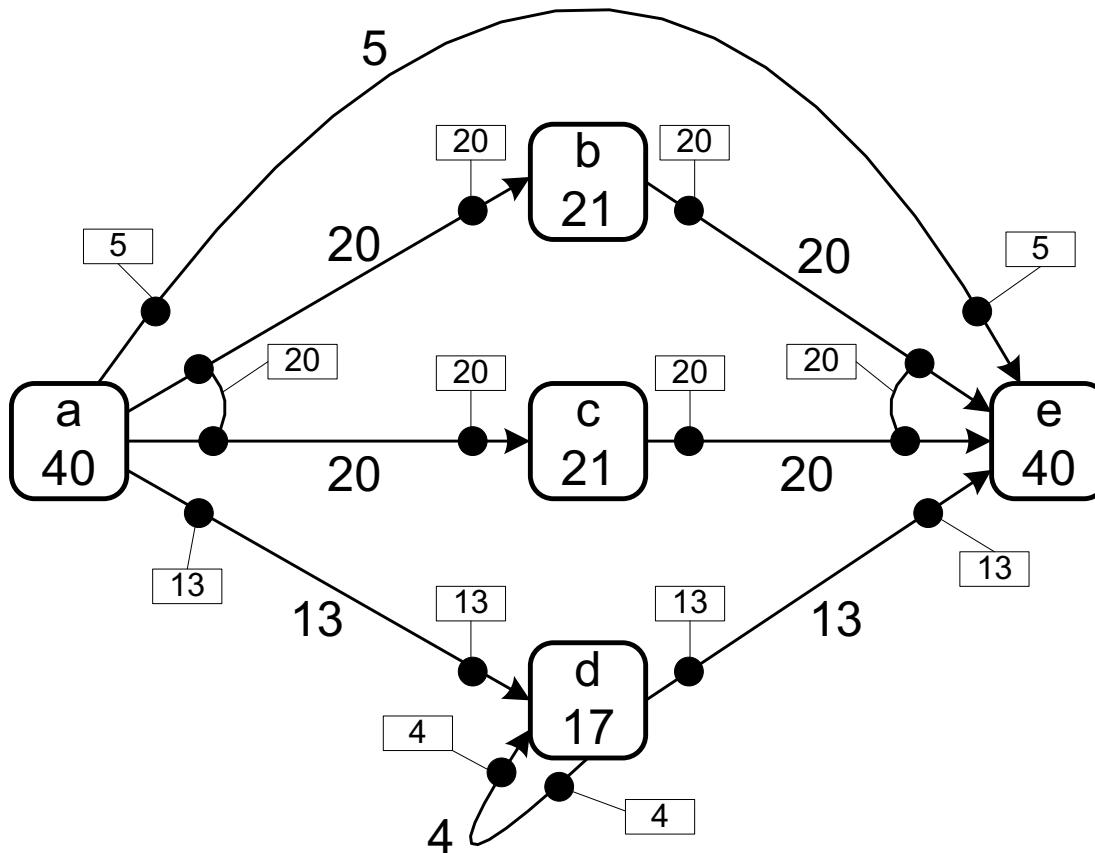
# Higher threshold (5 direct successions and a dependency of at least 0.9)



$ >_L $	$a$	$b$	$c$	$d$	$e$
$a$	0	11	11	13	5
$b$	0	0	10	0	11
$c$	0	10	0	0	11
$d$	0	0	0	4	13
$e$	0	0	0	0	0

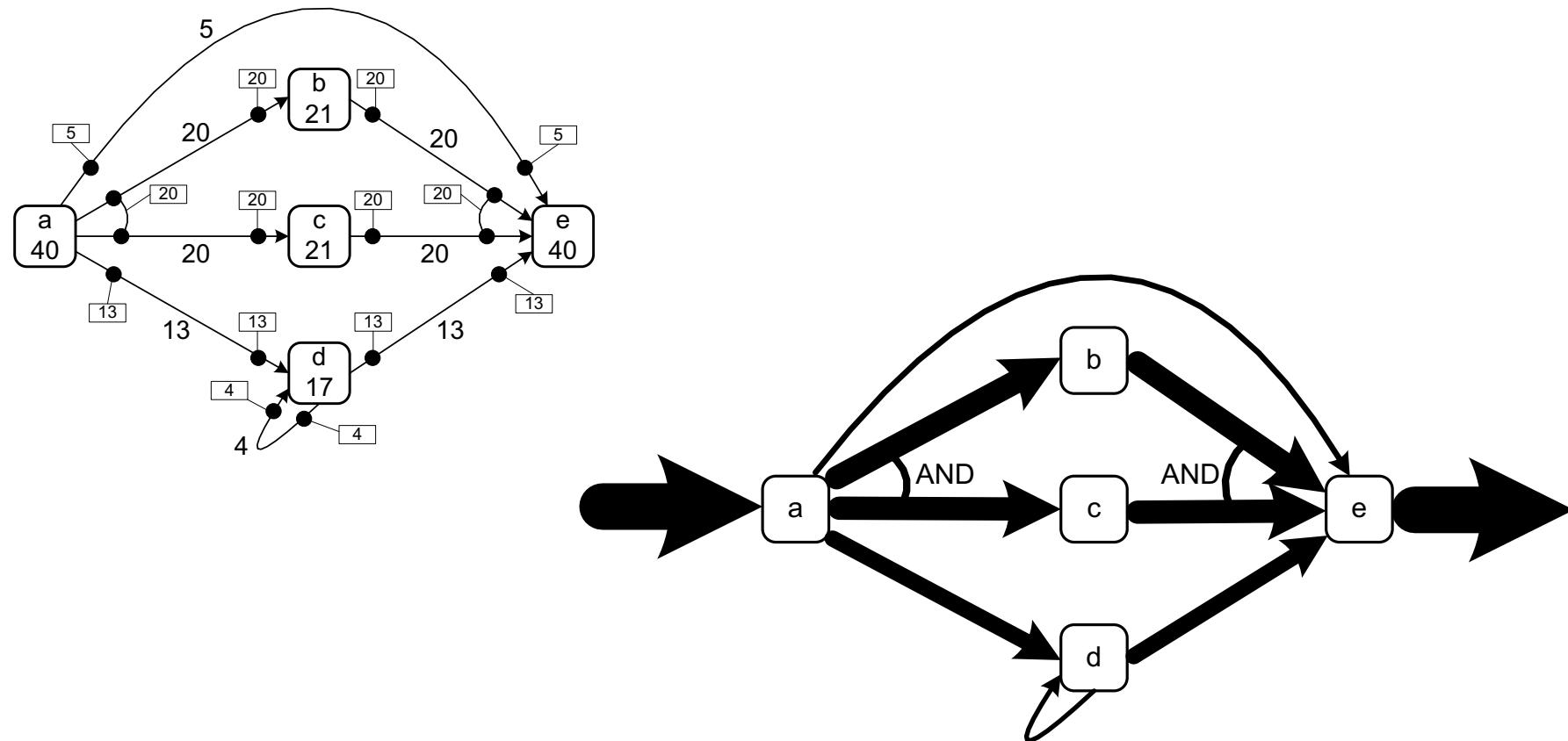
$\Rightarrow_L$	$a$	$b$	$c$	$d$	$e$
$a$	$\frac{0}{0+1} = 0$	$\frac{11-0}{11+0+1} = 0.92$	$\frac{11-0}{11+0+1} = 0.92$	$\frac{13-0}{13+0+1} = 0.93$	$\frac{5-0}{5+0+1} = 0.83$
$b$	$\frac{0-11}{0+11+1} = -0.92$	$\frac{0}{0+1} = 0$	$\frac{10-10}{10+10+1} = 0$	$\frac{0-0}{0+0+1} = 0$	$\frac{11-0}{11+0+1} = 0.92$
$c$	$\frac{0-11}{0+11+1} = -0.92$	$\frac{10-10}{10+10+1} = 0$	$\frac{0}{0+1} = 0$	$\frac{0-0}{0+0+1} = 0$	$\frac{11-0}{11+0+1} = 0.92$
$d$	$\frac{0-13}{0+13+1} = -0.93$	$\frac{0-0}{0+0+1} = 0$	$\frac{0-0}{0+0+1} = 0$	$\frac{4}{4+1} = 0.80$	$\frac{13-0}{13+0+1} = 0.93$
$e$	$\frac{0-5}{0+5+1} = -0.83$	$\frac{0-11}{0+11+1} = -0.92$	$\frac{0-11}{0+11+1} = -0.92$	$\frac{0-13}{0+13+1} = -0.93$	$\frac{0}{0+1} = 0$

# Learning splits and joins



$$L = [\langle a, e \rangle^5, \langle a, b, c, e \rangle^{10}, \langle a, c, b, e \rangle^{10}, \langle a, b, e \rangle^1, \langle a, c, e \rangle^1, \\ \langle a, d, e \rangle^{10}, \langle a, d, d, e \rangle^2, \langle a, d, d, d, e \rangle^1]$$

# Alternative visualization

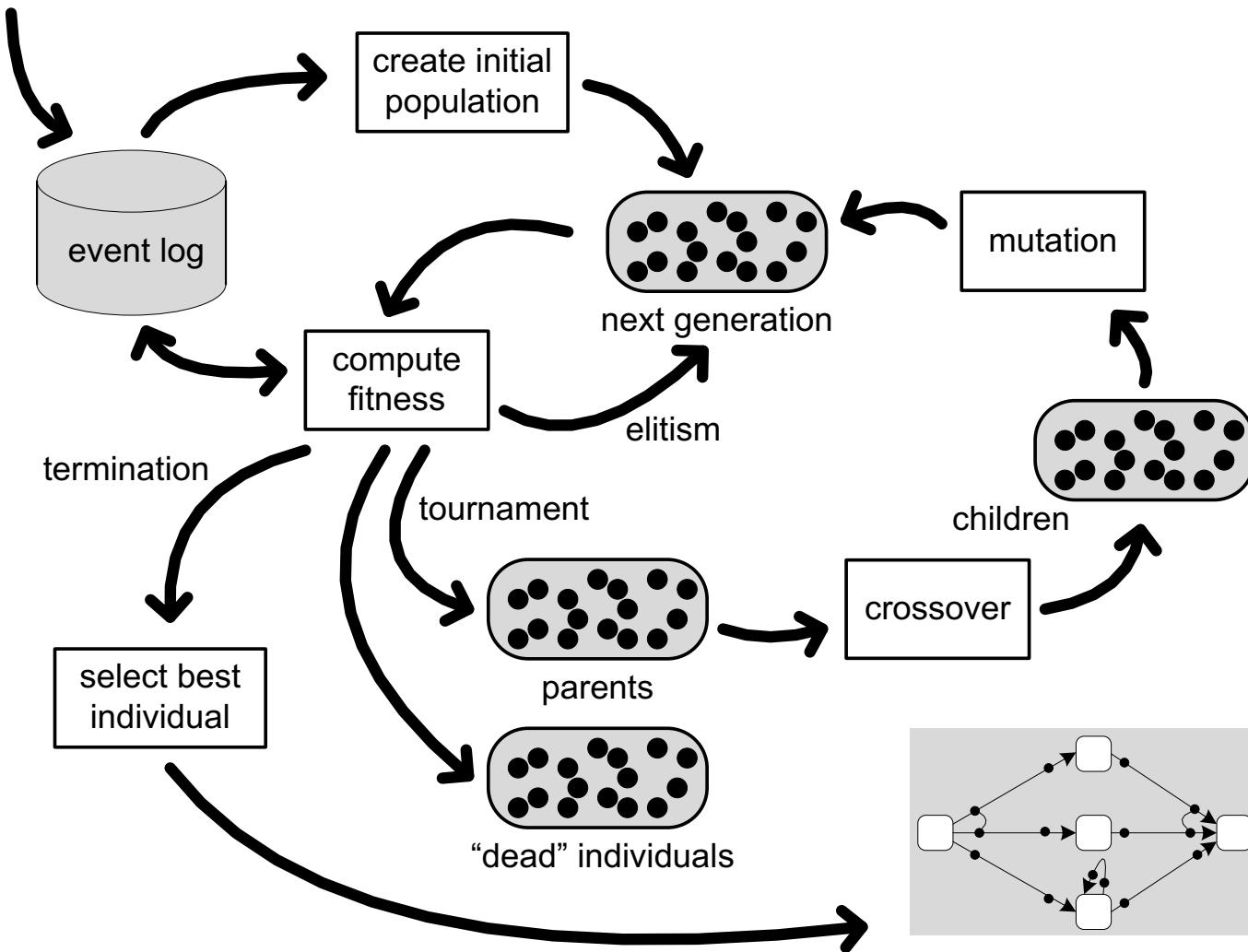


$$\begin{aligned} L = & [\langle a, e \rangle^5, \langle a, b, c, e \rangle^{10}, \langle a, c, b, e \rangle^{10}, \langle a, b, e \rangle^1, \langle a, c, e \rangle^1, \\ & \langle a, d, e \rangle^{10}, \langle a, d, d, e \rangle^2, \langle a, d, d, d, e \rangle^1] \end{aligned}$$

# Characteristics of heuristic mining

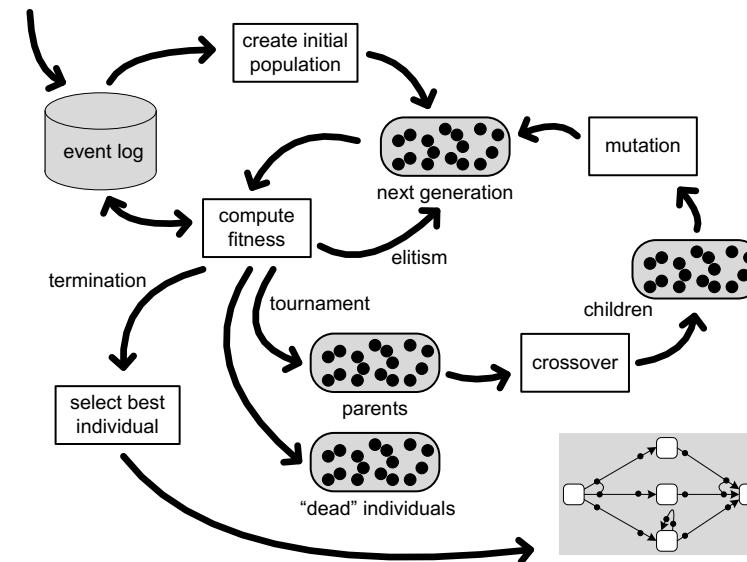
- Can deal with noise and therefore quite robust.
- Improved representational bias.
- Split and join rules are only considered locally (therefore most of the discovered model are not sound and require repair actions).

# Genetic process mining

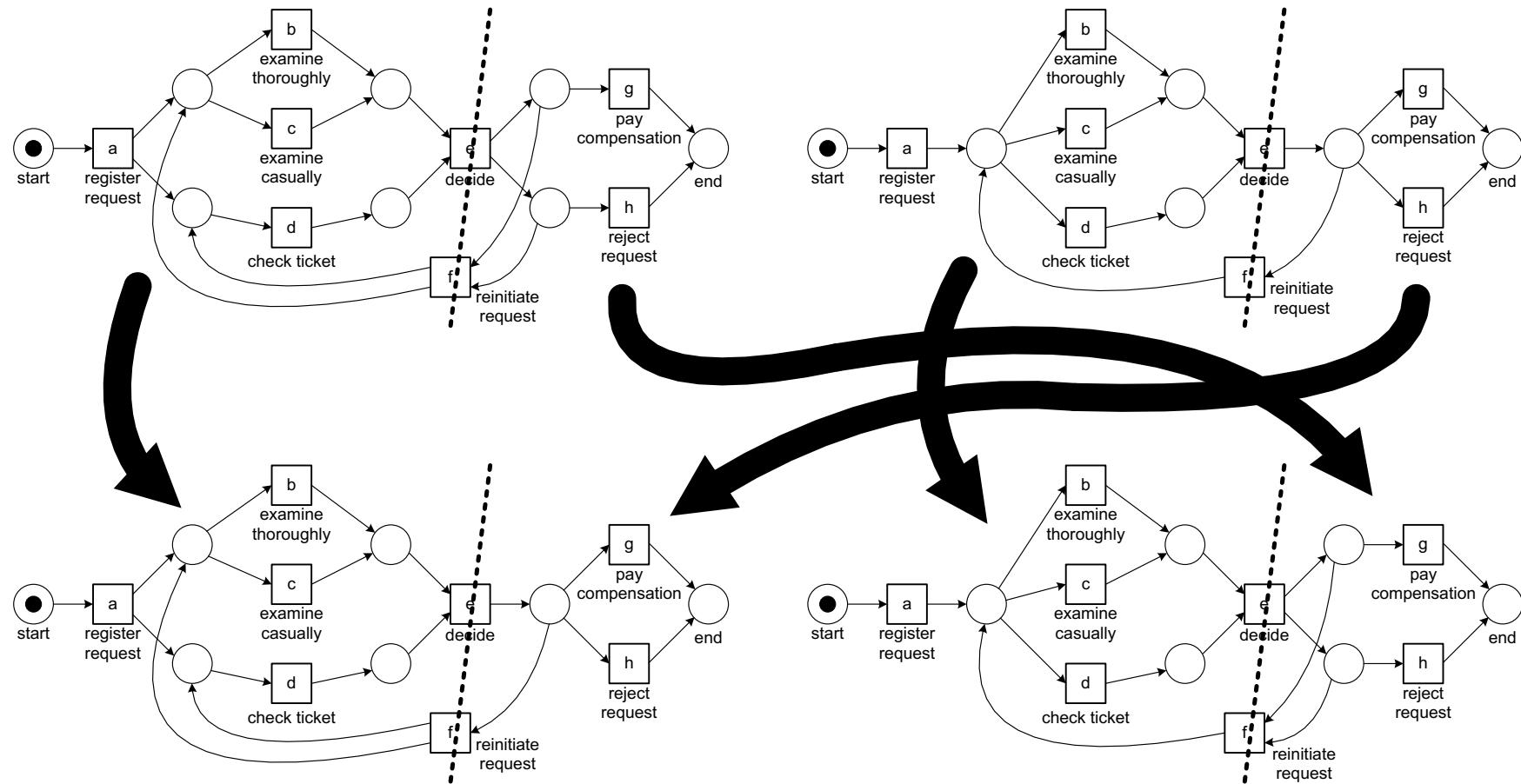


# Design decisions

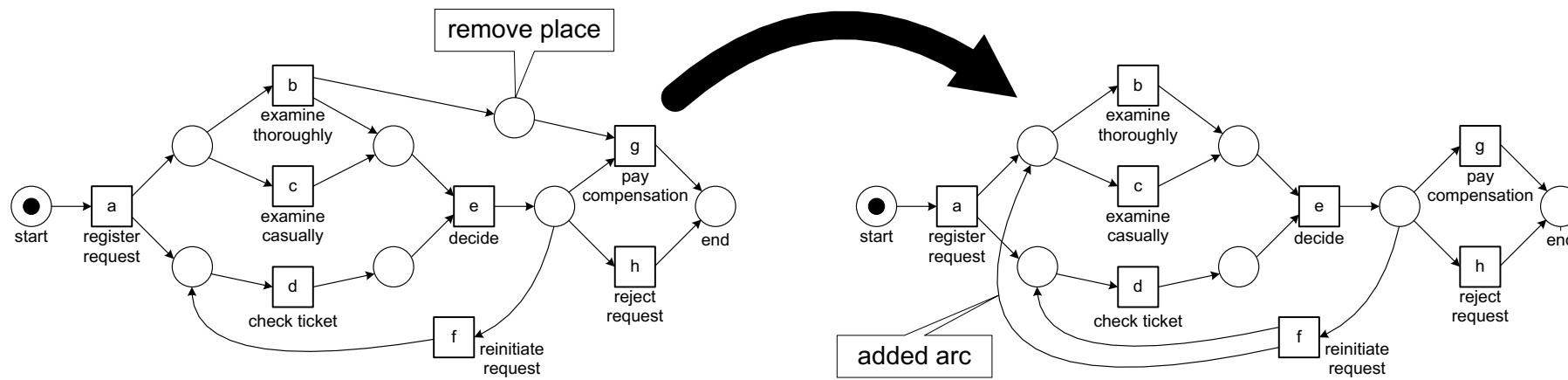
- Representation of individuals
- Initialization
- Fitness function
- Selection strategy (tournament and elitism)
- Crossover
- Mutation



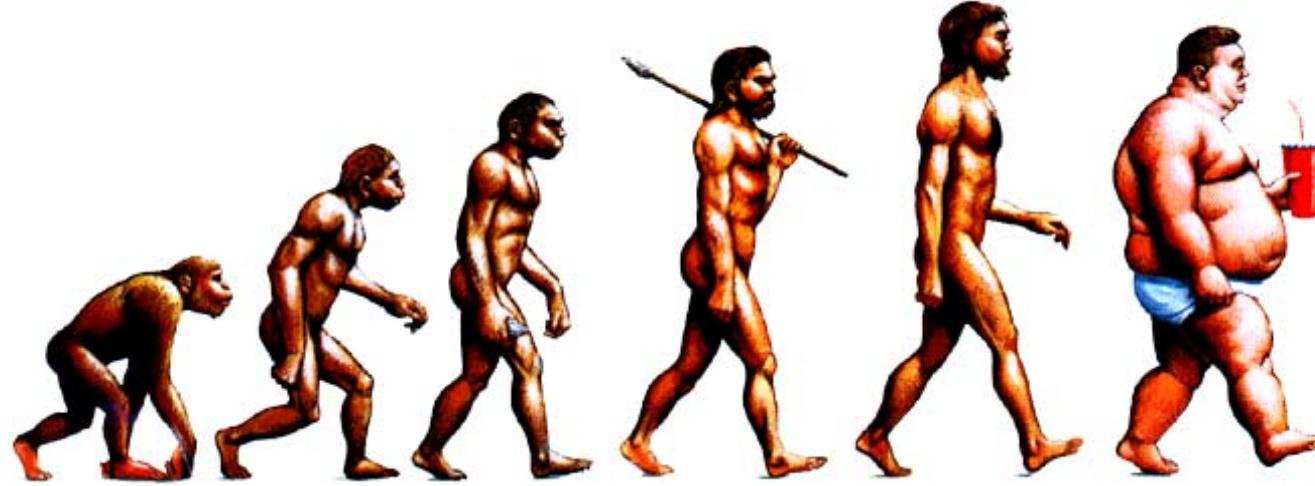
# Example: crossover



# Example: mutation



# Characteristics of genetic process mining



- Requires a lot of computing power.
- Can be distributed easily.
- Can deal with noise, infrequent behavior, duplicate tasks, invisible tasks, etc.
- Allows for incremental improvement and combinations with other approaches (heuristics post-optimization, etc.).

# Summary

- Process mining = data + processes
- Process discovery = events + learning
- Three different discovery algorithms covered
  - Alpha algorithm
  - Heuristic miner
  - Genetic miner
- But there are many many more! (Inductive Miner will be seen in the lab)
- **Next session: how good models are in describing reality ?**