

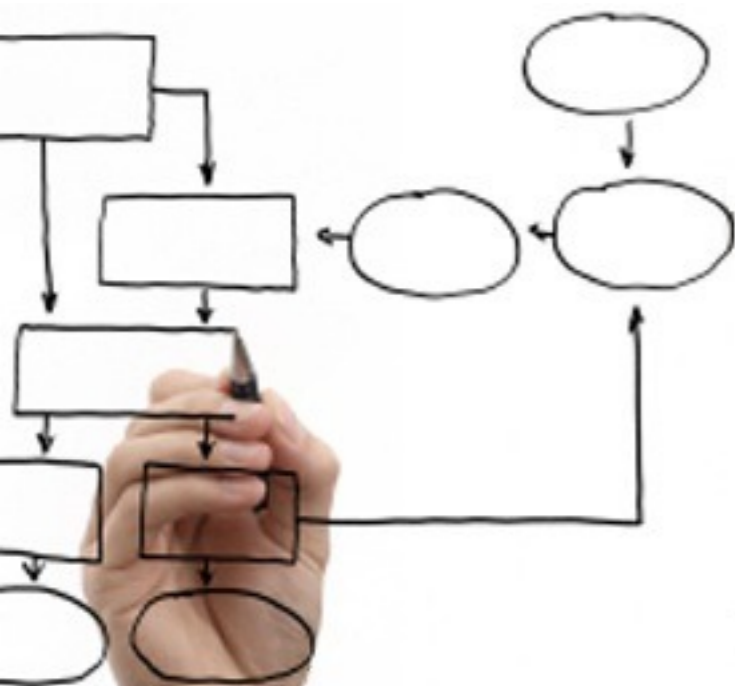
Methods for the specification and verification of business processes

MPB (6 cfu, 295AA)

Roberto Bruni

<http://www.di.unipi.it/~bruni>

24 - Process Mining



Object

We overview the key principles of
process mining

Process Mining

Process mining is a relative young research discipline that sits between

machine learning and data mining on the one hand and process modeling and analysis on the other hand.

The idea of process mining is to discover, monitor and improve real processes (i.e., not assumed processes) by extracting knowledge from event logs readily available in today's systems.

Processes, Cases, Events, Attributes

A process consists of cases.

A case consists of events such that each event relates to precisely one case.

Events within a case are ordered.

Events can have attributes.

Examples of typical attribute names are activity, time, costs, and resource.

Event Logs

Let us assume that it is possible to sequentially record events such that each event:

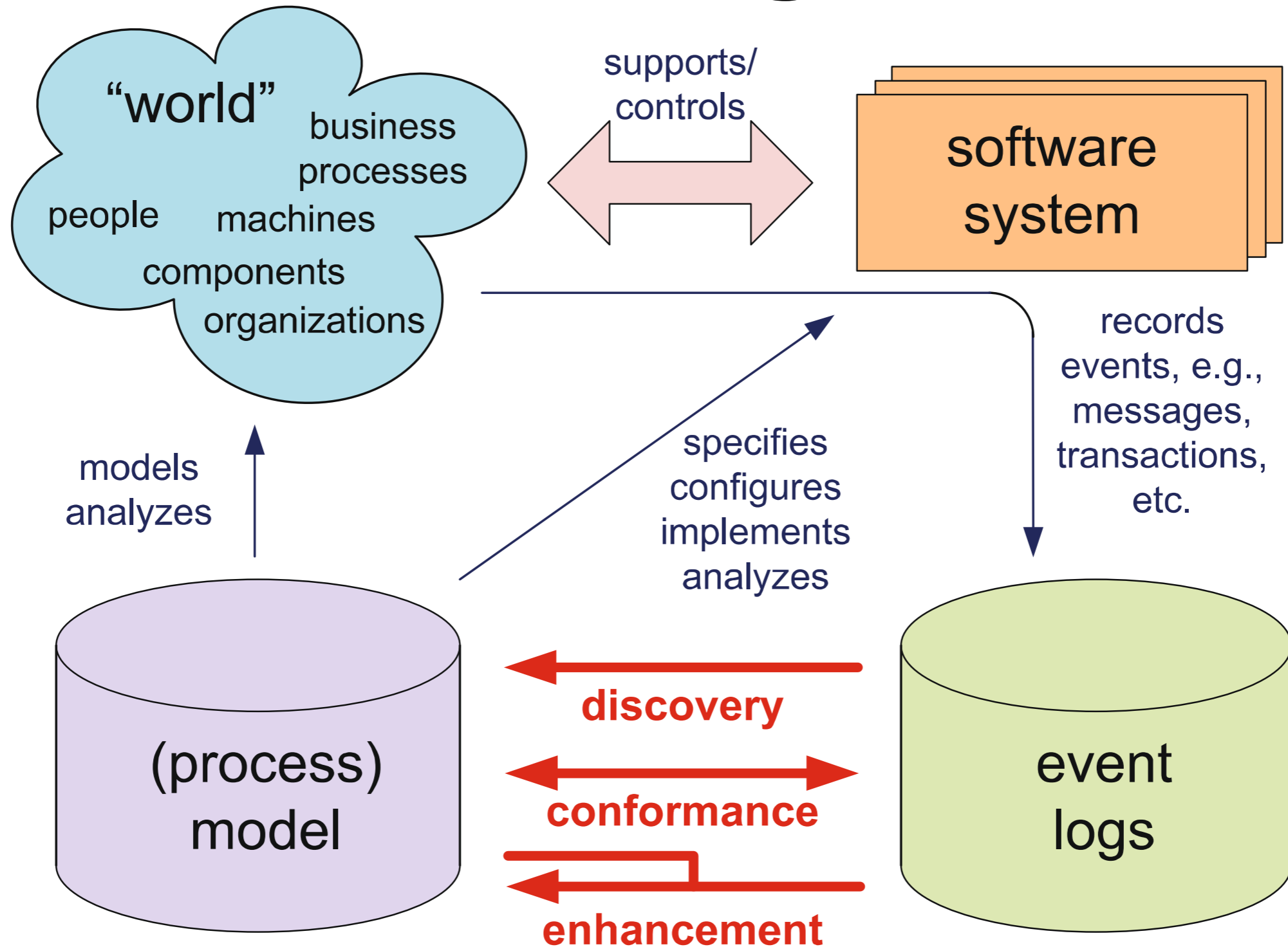
refers to an activity
(i.e., a well-defined step in the process)

and is related to a particular case
(i.e., a process instance).

Event Log Example

Case id	Event id	Properties				
		Timestamp	Activity	Resource	Cost	...
1	35654423	30-12-2010:11.02	Register request	Pete	50	...
	35654424	31-12-2010:10.06	Examine thoroughly	Sue	400	...
	35654425	05-01-2011:15.12	Check ticket	Mike	100	...
	35654426	06-01-2011:11.18	Decide	Sara	200	...
	35654427	07-01-2011:14.24	Reject request	Pete	200	...
2	35654483	30-12-2010:11.32	Register request	Mike	50	...
	35654485	30-12-2010:12.12	Check ticket	Mike	100	...
	35654487	30-12-2010:14.16	Examine casually	Pete	400	...
	35654488	05-01-2011:11.22	Decide	Sara	200	...
	35654489	08-01-2011:12.05	Pay compensation	Ellen	200	...

Process Mining Scheme



Discovery

A **discovery** technique takes an event log and produces a model without using any a-priori information.

If the event log contains information about resources, one can also discover resource-related models, e.g., a social network showing how people work together in an organization.

Conformance

An existing process model is compared with an event log of the same process.

Conformance checking can be used to check if reality, as recorded in the log, conforms to the model and vice versa.

Conformance checking may be used to detect, locate and explain deviations, and to measure the severity of these deviations.

Enhancement

The idea is to **extend/improve an existing process model** using information about the actual process recorded in some event log.

Whereas conformance checking measures the alignment between model and reality, this third type of process mining aims at changing or extending the a-priori model.

Enhancement: Repair

One type of enhancement is **repair**,
i.e., modifying the model to better reflect reality.

For example, if two activities are modeled sequentially
but in reality can happen in any order,
then the model may be corrected to reflect this.

Four Perspectives

Control-Flow Perspective

The **control-flow perspective** focuses on the control-flow, i.e., the ordering of activities.

The goal of mining this perspective is to find a good characterization of all possible paths, e.g., expressed in terms of a Petri net or some other notation (e.g., EPC, BPMN, and UML AD).

We shall focus on this perspective

Organizational Perspective

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 they are related.

The goal is to either structure the organization by classifying people in terms of roles and organizational units or to show the social network.

Case Perspective

The **case perspective** focuses on properties of cases.

Obviously, a case can be characterized by its path in the process or by the originators working on it.

However, cases can also be characterized by the values of the corresponding data elements.

For example, if a case represents a replenishment order, it may be interesting to know the supplier or the number of products ordered.

Time Perspective

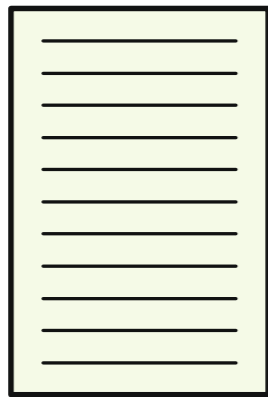
The **time perspective** is concerned with the timing and frequency of events (performance checking).

When events bear timestamps it is possible to discover bottlenecks, measure service levels, monitor the utilization of resources, and predict the remaining processing time of running cases.

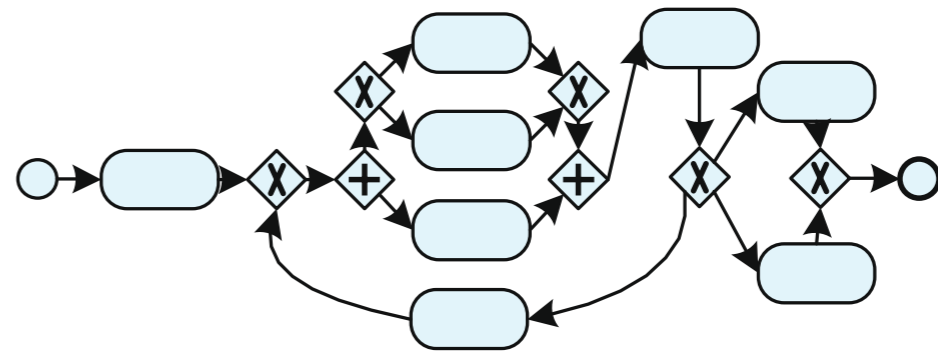
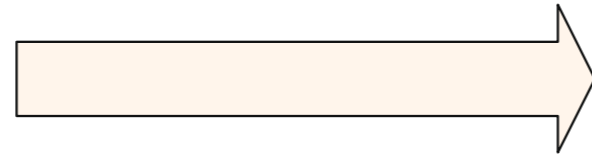
Play-in, Play-out, Replay

Play-in

Play-In



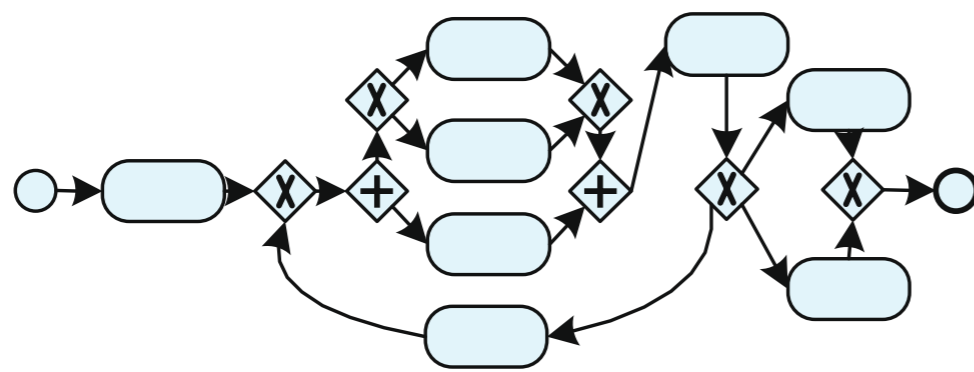
event log



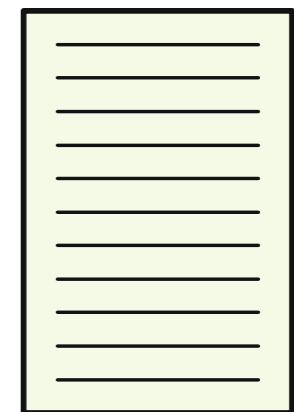
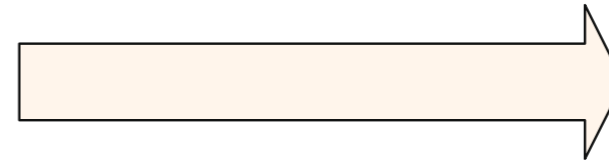
process model

Play-out

Play-Out



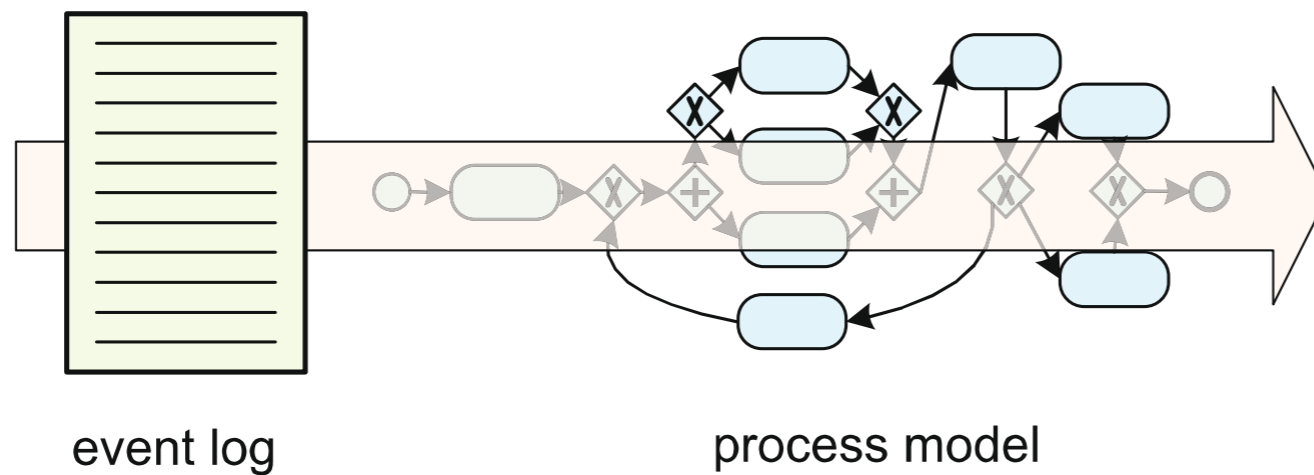
process model



event log

Replay

Replay



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

An Example

Event Log Example

Case id	Event id	Properties				
		Timestamp	Activity	Resource	Cost	...
1	35654423	30-12-2010:11.02	Register request	Pete	50	...
	35654424	31-12-2010:10.06	Examine thoroughly	Sue	400	...
	35654425	05-01-2011:15.12	Check ticket	Mike	100	...
	35654426	06-01-2011:11.18	Decide	Sara	200	...
	35654427	07-01-2011:14.24	Reject request	Pete	200	...
2	35654483	30-12-2010:11.32	Register request	Mike	50	...
	35654485	30-12-2010:12.12	Check ticket	Mike	100	...
	35654487	30-12-2010:14.16	Examine casually	Pete	400	...
	35654488	05-01-2011:11.22	Decide	Sara	200	...
	35654489	08-01-2011:12.05	Pay compensation	Ellen	200	...

Case id	Event id	Properties			
		Timestamp	Activity	Resource	Cost
1	35654423	30-12-2010:11.02	Register request	Pete	50
	35654424	31-12-2010:10.05	Examine thoroughly	Sue	400
	35654425	05-01-2011:15.12	Check ticket	Ellen	100
	35654426	06-01-2011:12.18	Decide	Sara	200
	35654427	07-01-2011:14.21	Reject request	Pete	200
2	35654483	30-12-2010:11.32	Register request	Mike	50
	35654485	30-12-2010:12.12	Check ticket	Mike	100
	35654487	30-12-2010:14.16	Examine casually	Pete	400
	35654488	05-01-2011:11.22	Decide	Sara	200
	35654489	08-01-2011:12.05	Pay compensation	Ellen	200
3	35654521	30-12-2010:14.32	Register request	Pete	50
	35654522	30-12-2010:15.06	Examine casually	Mike	400
	35654524	30-12-2010:16.34	Check ticket	Ellen	100
	35654525	06-01-2011:09.18	Decide	Sara	200
	35654526	06-01-2011:12.18	Reinitiate request	Sara	200
	35654527	06-01-2011:13.06	Examine thoroughly	Sean	400
	35654530	08-01-2011:11.43	Check ticket	Pete	100
	35654531	09-01-2011:09.55	Decide	Sara	200
	35654533	15-01-2011:10.45	Pay compensation	Ellen	200
4	35654641	06-01-2011:15.02	Register request	Pete	50
	35654643	07-01-2011:12.06	Check ticket	Mike	100
	35654644	08-01-2011:14.43	Examine thoroughly	Sean	400
	35654645	09-01-2011:12.02	Decide	Sara	200
	35654647	12-01-2011:15.44	Reject request	Ellen	200
5	35654711	06-01-2011:09.02	Register request	Ellen	50
	35654712	07-01-2011:10.16	Examine casually	Mike	400
	35654714	08-01-2011:11.22	Check ticket	Pete	100
	35654715	10-01-2011:13.28	Decide	Sara	200
	35654716	11-01-2011:16.18	Reinitiate request	Sara	200
	35654718	14-01-2011:14.33	Check ticket	Ellen	100
	35654719	16-01-2011:15.50	Examine casually	Mike	400
	35654720	19-01-2011:11.18	Decide	Sara	200
	35654721	20-01-2011:12.48	Reinitiate request	Sara	200
	35654722	21-01-2011:09.06	Examine casually	Sue	400
	35654724	21-01-2011:11.34	Check ticket	Pete	100
	35654725	23-01-2011:13.12	Decide	Sara	200
	35654726	24-01-2011:14.56	Reject request	Mike	200

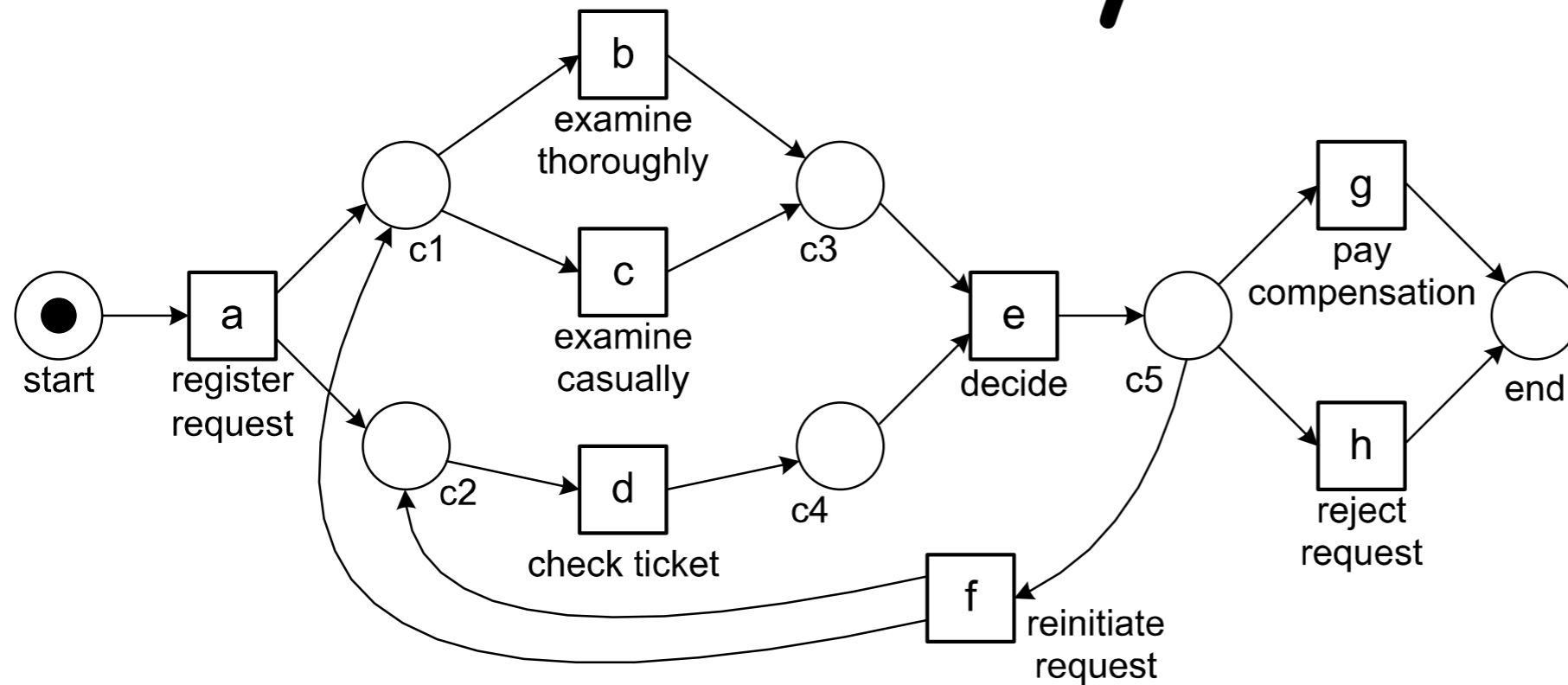
Event Log Example

Case id	Event id	Properties			
		Timestamp	Activity	Resource	Cost
6	35654871	06-01-2011:15.02	Register request	Mike	50
	35654873	06-01-2011:16.06	Examine casually	Ellen	400
	35654874	07-01-2011:16.21	Check ticket	Mike	100
	35654875	07-01-2011:15.52	Decide	Sara	200
	35654877	16-01-2011:11.47	Pay compensation	Mike	200
...

Table 1.2 A more compact representation of log shown in Table 1.1: *a* = register request, *b* = examine thoroughly, *c* = examine casually, *d* = check ticket, *e* = decide, *f* = reinitiate request, *g* = pay compensation, and *h* = reject request

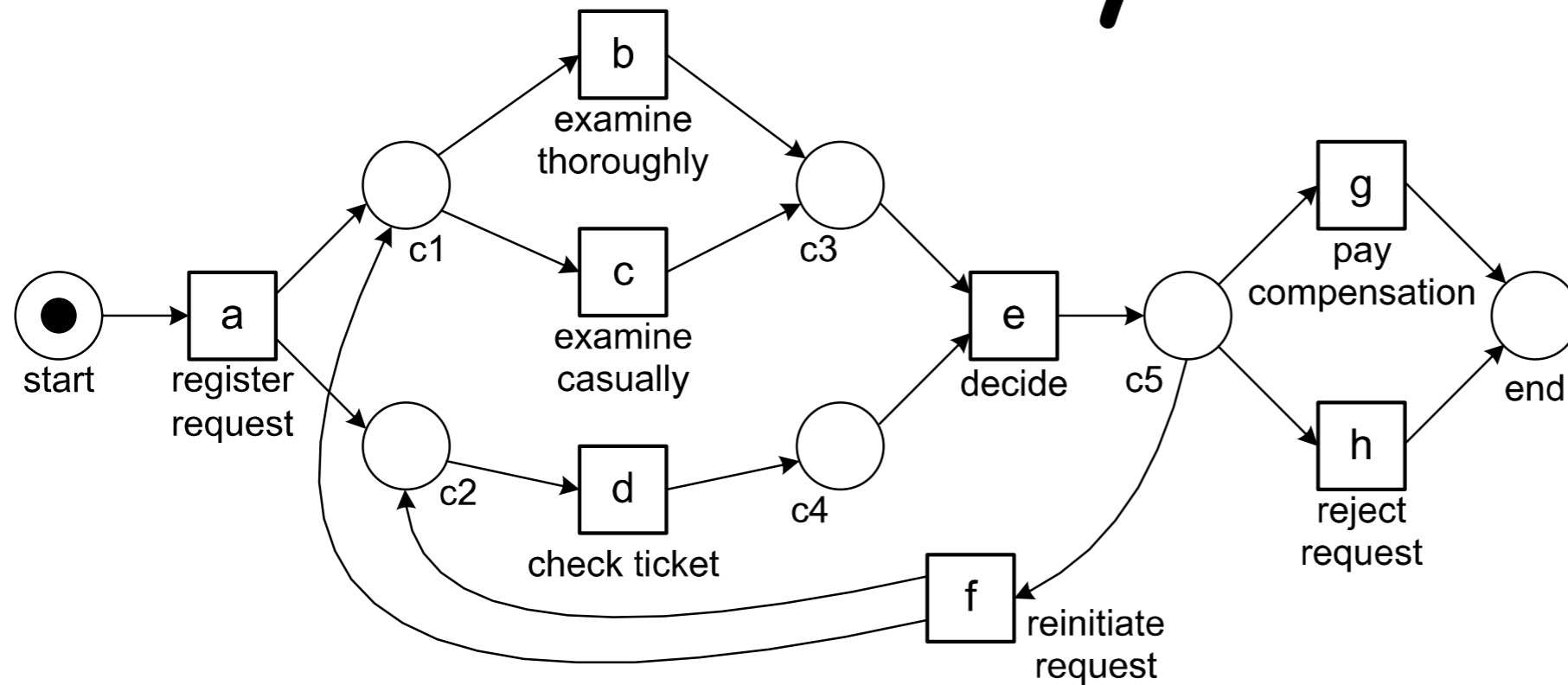
Case id	Trace
1	$\langle a, b, d, e, h \rangle$
2	$\langle a, d, c, e, g \rangle$
3	$\langle a, c, d, e, f, b, d, e, g \rangle$
4	$\langle a, d, b, e, h \rangle$
5	$\langle a, c, d, e, f, d, c, e, f, c, d, e, h \rangle$
6	$\langle a, c, d, e, g \rangle$
...	...

Discovery Example



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	$\langle a, c, d, e, g \rangle$
...	...

Discovery Example

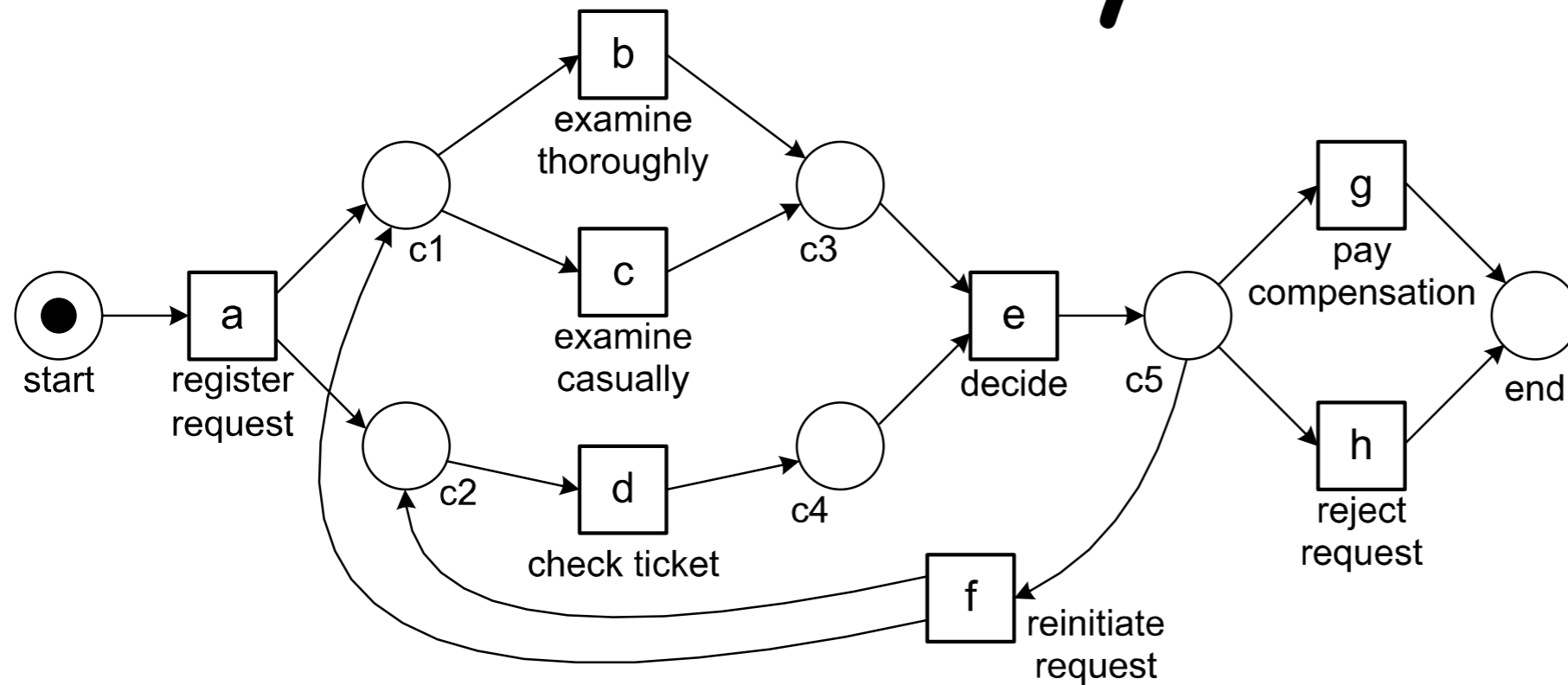


All cases start with a and end with either g or h.

Every e is preceded by d and one of the examination activities (b or c).

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	$\langle a, c, d, e, g \rangle$
...	...

Discovery Example

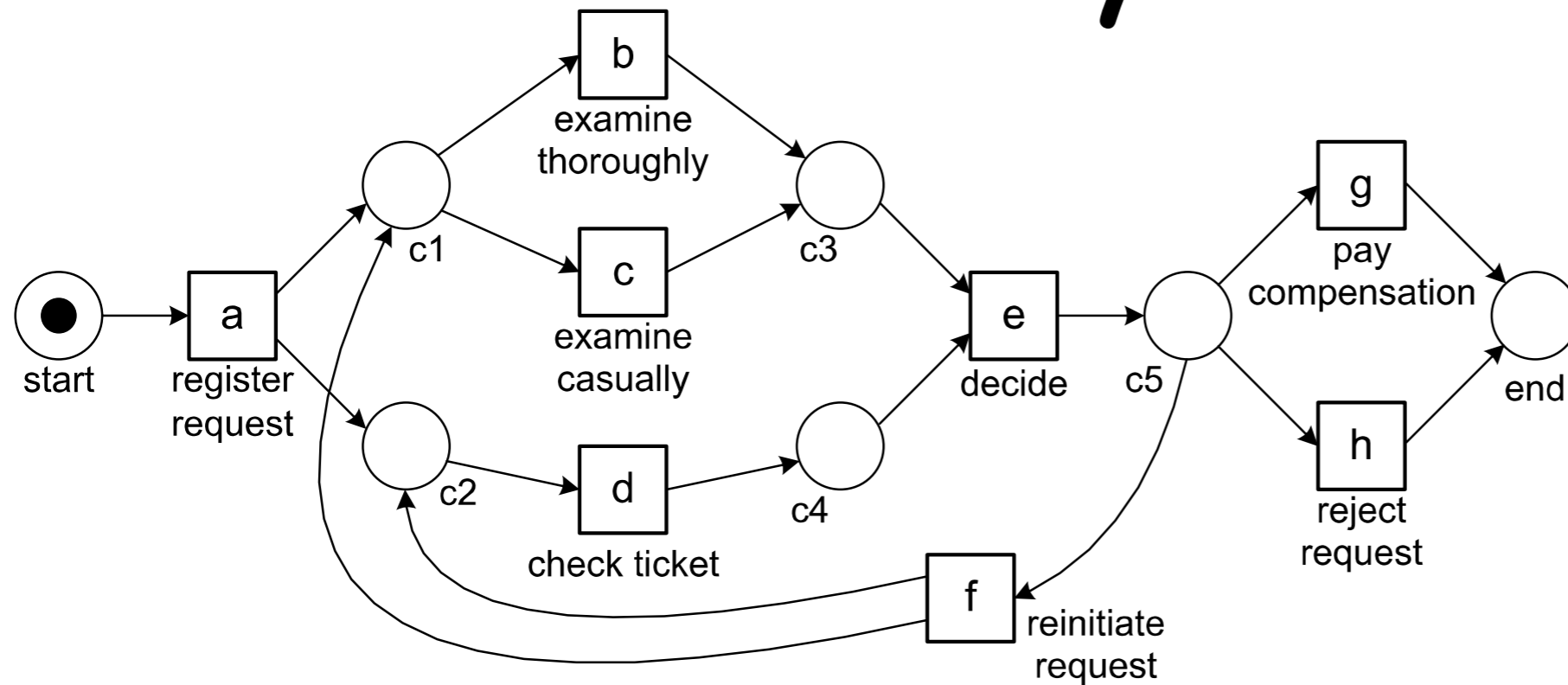


Moreover, e is followed by f , g , or h.

The repeated execution of b or c, d, and e suggests the presence of a loop.

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	$\langle a, c, d, e, g \rangle$
...	...

Discovery Example



These characteristics
are adequately captured
by the net.

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	$\langle a, c, d, e, g \rangle$
...	...

Overfitting and Underfitting

One of the challenges of process mining is to balance between

overfitting (the model is too specific and only allows for the accidental behavior observed) and

underfitting (the model is too general and allows for behavior unrelated to the behavior observed).

Discussion

The Petri net shown also allows for traces not in the log.

For example, other possible traces are

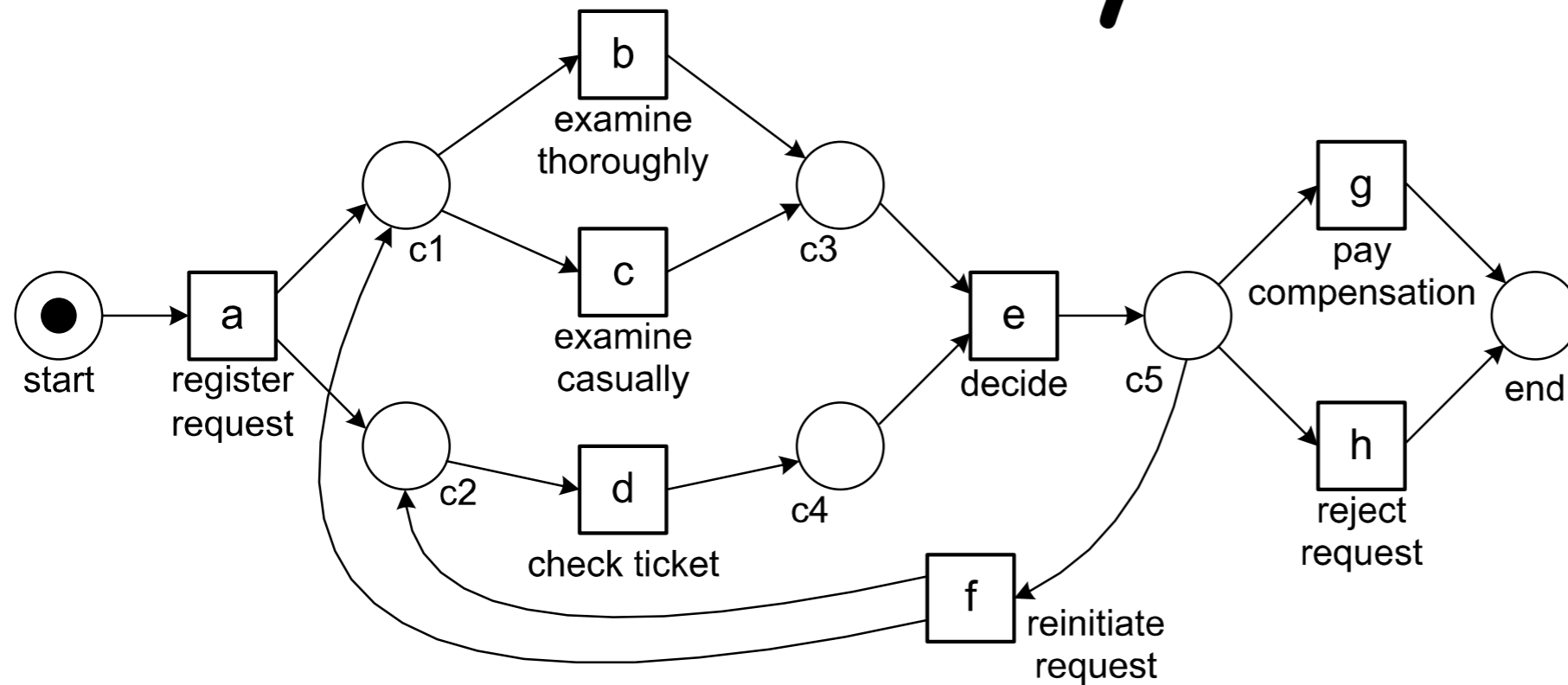
$\langle a, d, c, e, f, b, d, e, g \rangle$ and

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

This is a desired phenomenon as the goal is not to represent just the particular set of example traces in the event log.

Process mining algorithms need to generalize the behavior contained in the log to show the most likely underlying model that is not invalidated by the next set of observations

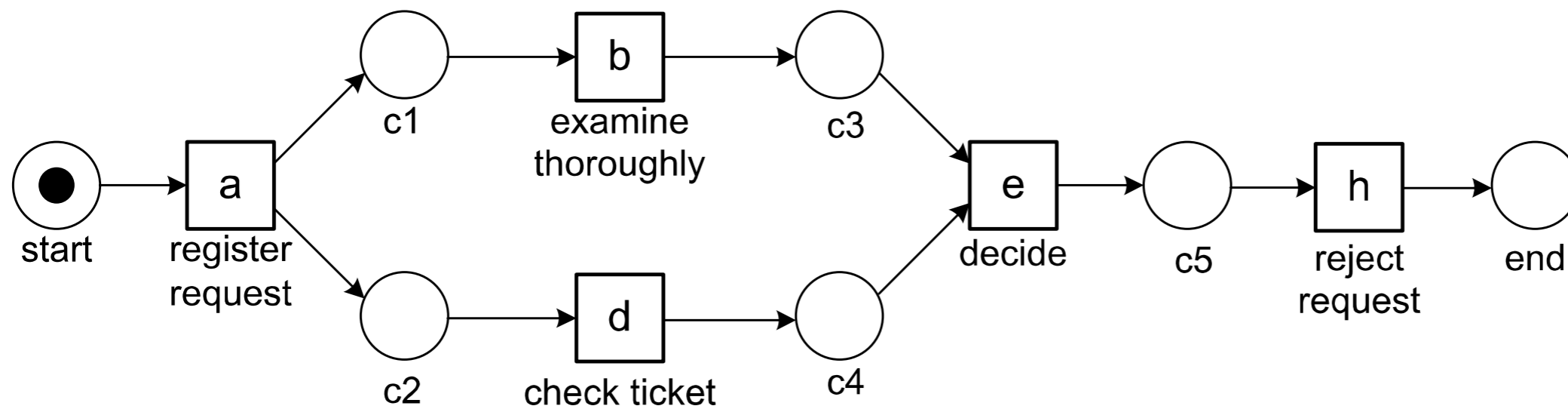
Discovery Example



When comparing the event log and the model, there seems to be a good balance between “overfitting” and “underfitting”.

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	$\langle a, c, d, e, g \rangle$
...	...

Another Discovery Example

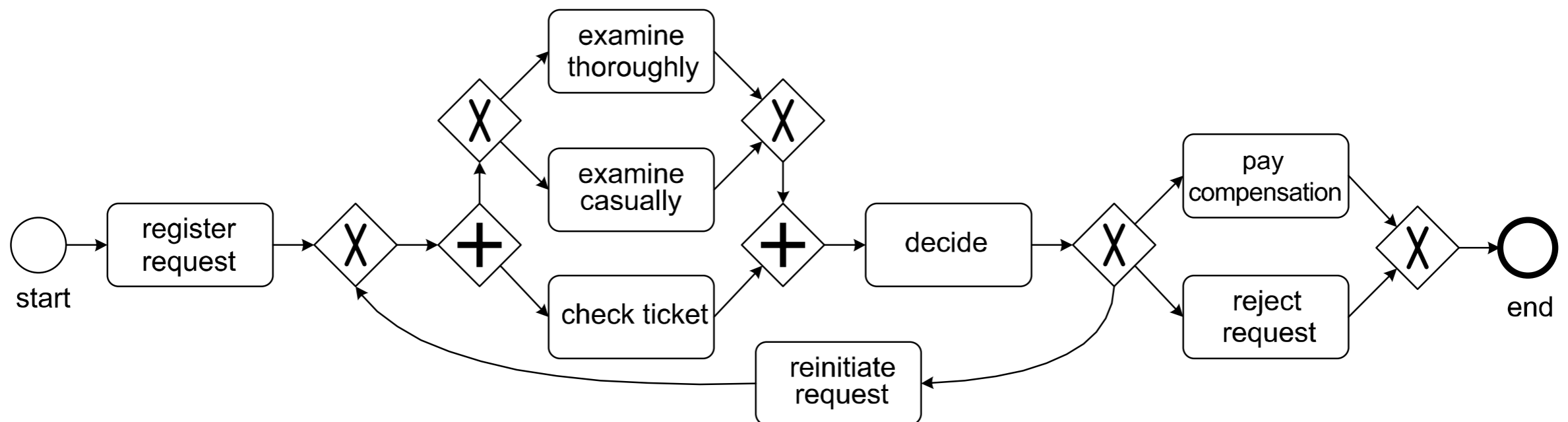


Case id	Trace
1	$\langle a, b, d, e, h \rangle$
2	$\langle a, d, e, e, g \rangle$
3	$\langle a, e, d, e, f, b, d, e, g \rangle$
4	$\langle a, d, b, e, h \rangle$
5	$\langle a, e, d, e, f, d, e, e, f, e, d, e, h \rangle$
6	$\langle a, e, d, e, g \rangle$
...	...

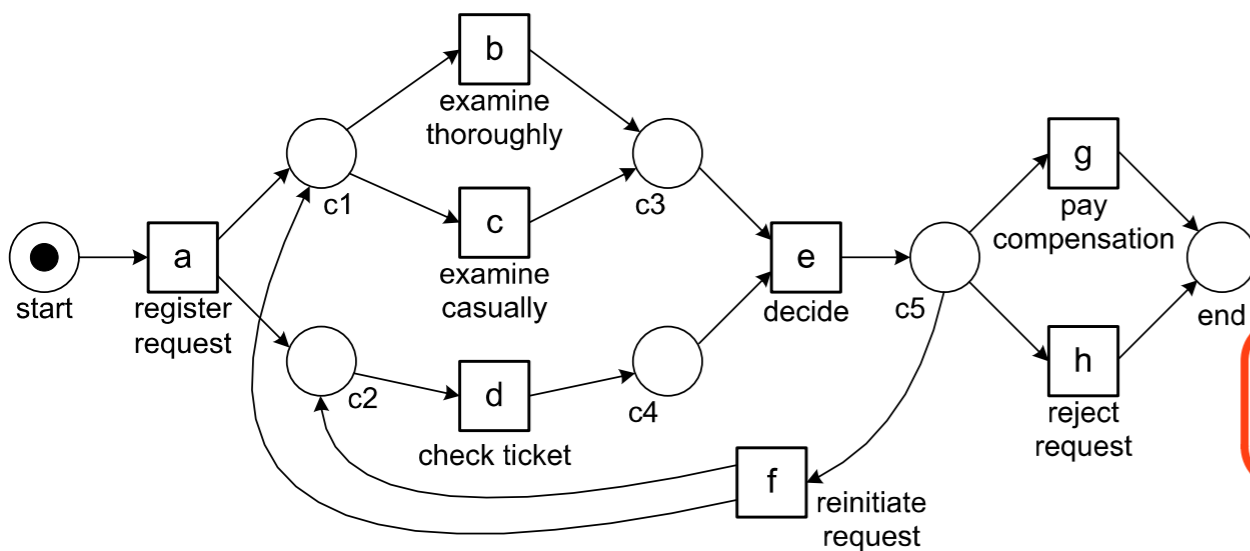
Mining Other Models

We used Petri nets to represent the discovered process models, because Petri nets are a succinct way of representing processes and have unambiguous but intuitive semantics.

However, some mining techniques are independent of the desired representation.

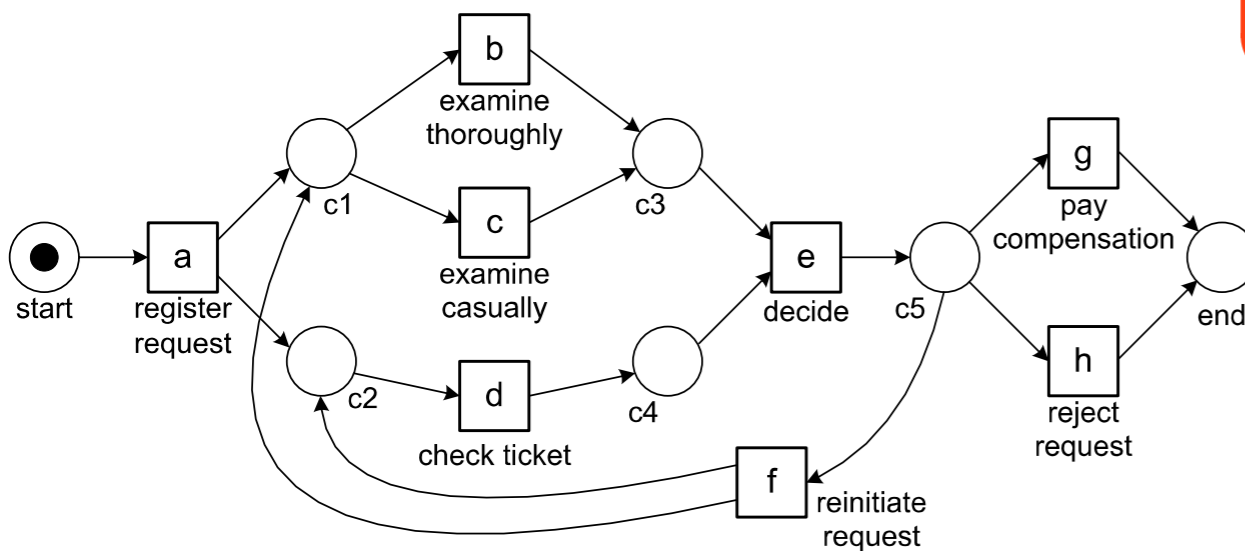
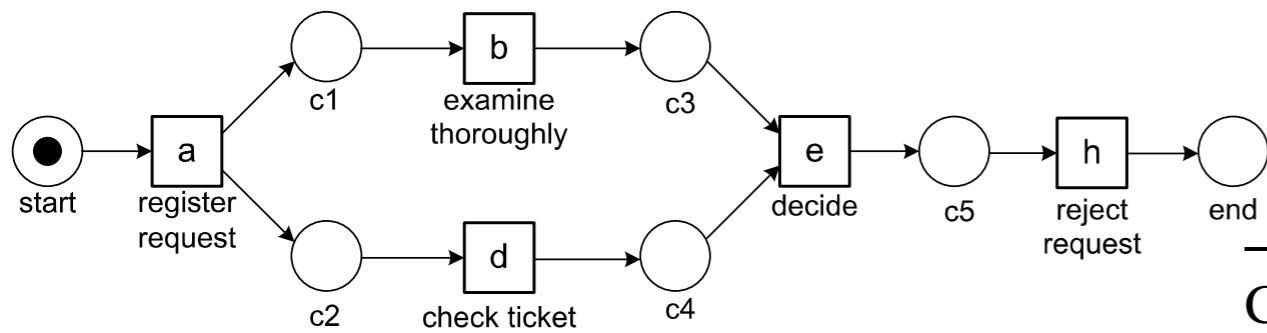


Conformance Example



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

Conformance Example



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

Process Discovery: α -Algorithm

Process Discovery

Process discovery is the activity that combines Discovery with the Control-flow Perspective.

The general problem:

A **process discovery algorithm** is a function that maps an event log L onto a process model M such that the model M is “representative” for the behavior seen in the event log L .

We focus on *simple event logs* and Petri net models (possibly sound workflow nets).

Simple Event Log

Let A be a set of activities.

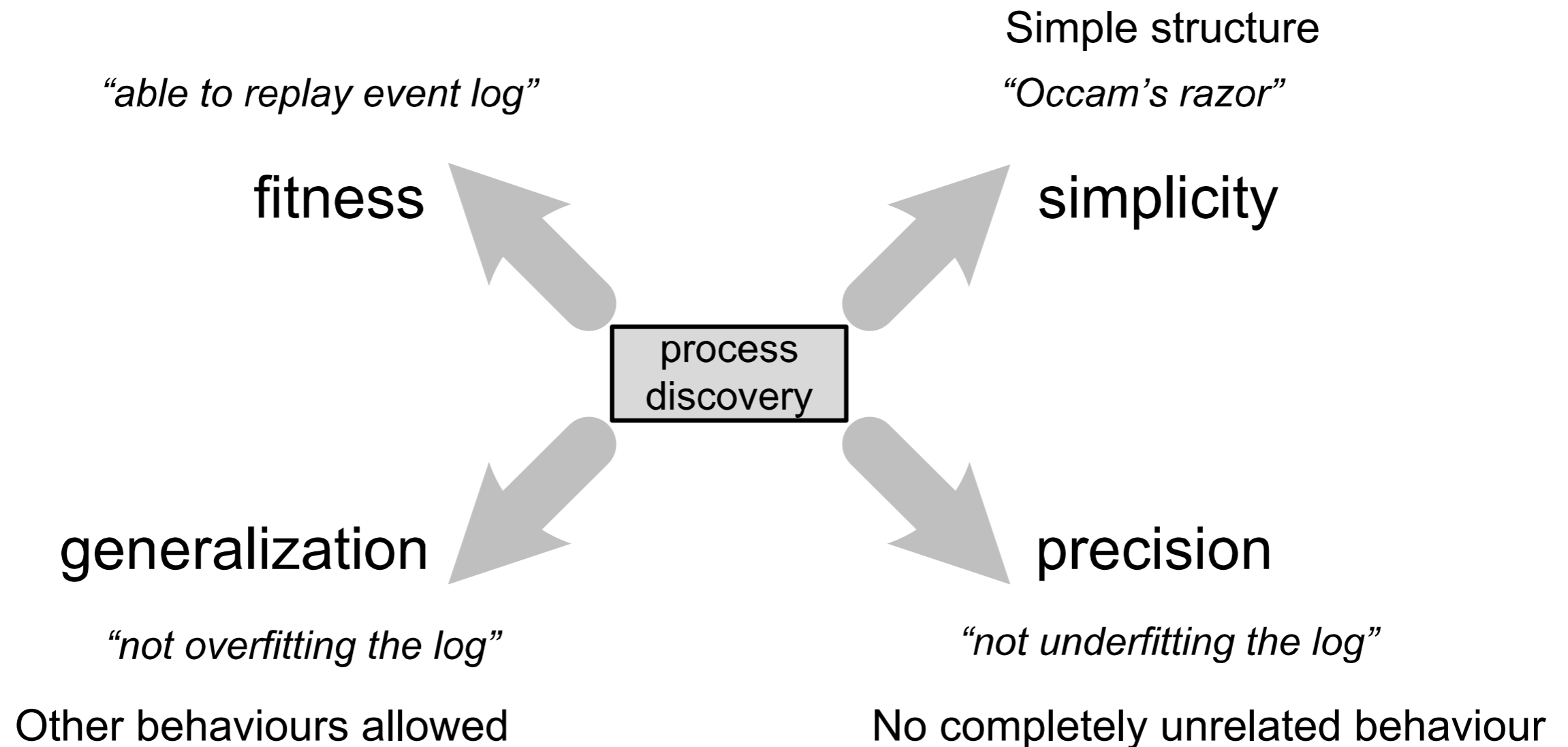
A **simple trace** over A is a finite sequence of activities.

A **simple event log** over A is a multiset of traces.

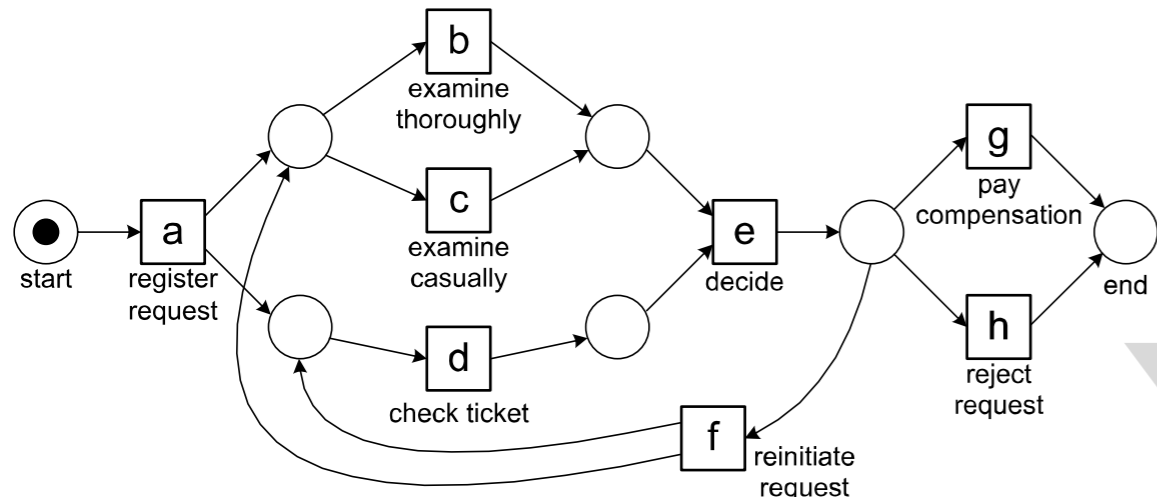
$$L_1 = [\langle a, b, c, d \rangle^3, \langle a, c, b, d \rangle^2, \langle a, e, d \rangle]$$

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

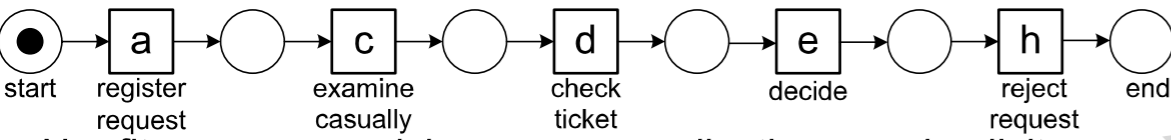
Challenges



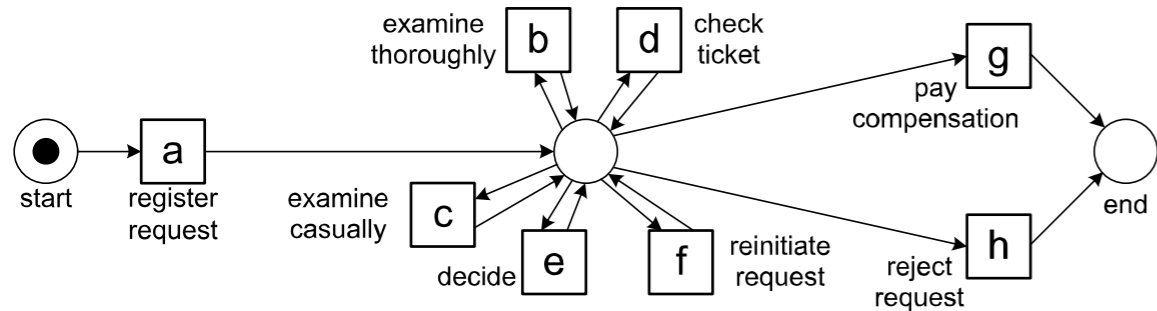
Appropriateness



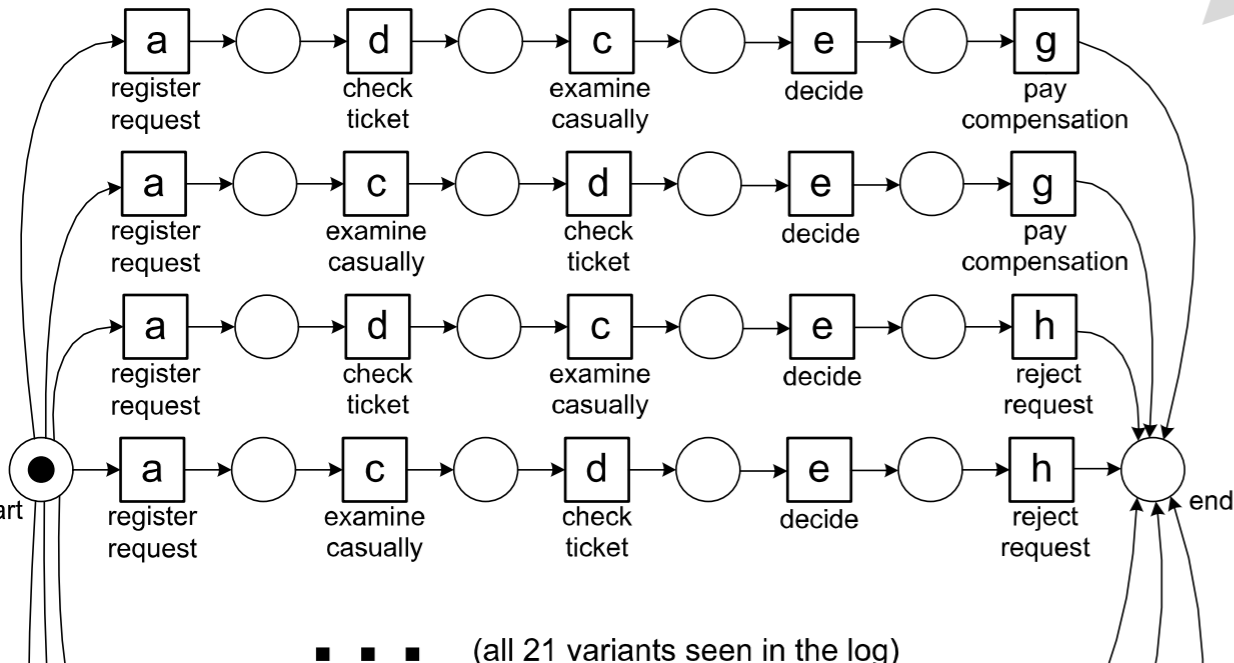
N_1 : fitness = +, precision = +, generalization = +, simplicity = +



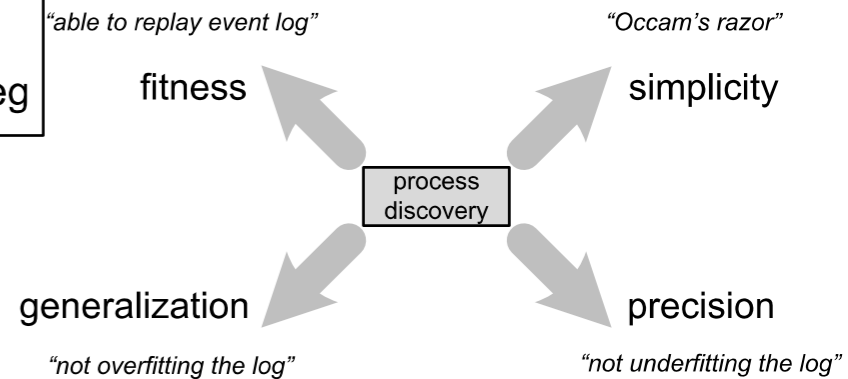
N_2 : fitness = -, precision = +, generalization = -, simplicity = +



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



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



α -Algorithm

The α -algorithm was one of the first process discovery algorithms that could adequately deal with concurrency.

It has several limitations,

but it provides a good introduction into the topic:

The α -algorithm is simple and many of its ideas have been embedded in more complex and robust techniques.

The α -algorithm scans the event log for particular patterns, called **log-based ordering relations**, to create a **footprint** of the log.

Log-based Ordering Relations

$a >_L b$ if and only if there is a trace $\sigma = \langle t_1, t_2, t_3, \dots, t_n \rangle$ and $i \in \{1, \dots, n - 1\}$ such that $\sigma \in L$ and $t_i = a$ and $t_{i+1} = b$

- $a \rightarrow_L b$ if and only if $a >_L b$ and $b \not>_L a$
- $a \#_L b$ if and only if $a \not>_L b$ and $b \not>_L a$
- $a \parallel_L b$ if and only if $a >_L b$ and $b >_L a$

$x \rightarrow_L y, y \rightarrow_L x, x \#_L y, \text{ or } x \parallel_L y$

Log-based Ordering Relations: Example

- $a \rightarrow_L b$ if and only if $a >_L b$ and $b \not>_L a$
- $a \#_L b$ if and only if $a \not>_L b$ and $b \not>_L a$
- $a \parallel_L b$ if and only if $a >_L b$ and $b >_L a$

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

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

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

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

$$\parallel_{L_1} = \{(b, c), (c, b)\}$$

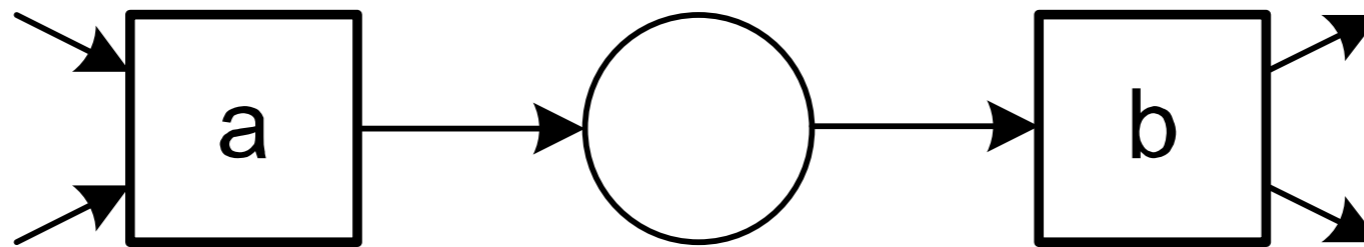
Footprint Matrix: Example

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

	<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}	$\#_{L_1}$

Patterns

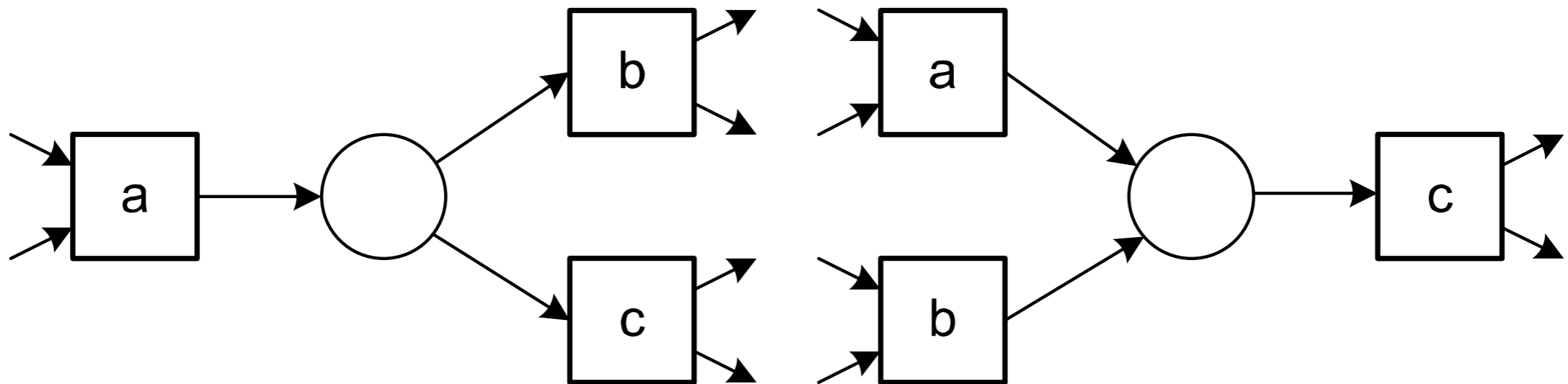
Footprints are useful to discover typical patterns of activities in the corresponding process model



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

Patterns

Footprints are useful to discover typical patterns of activities in the corresponding process model

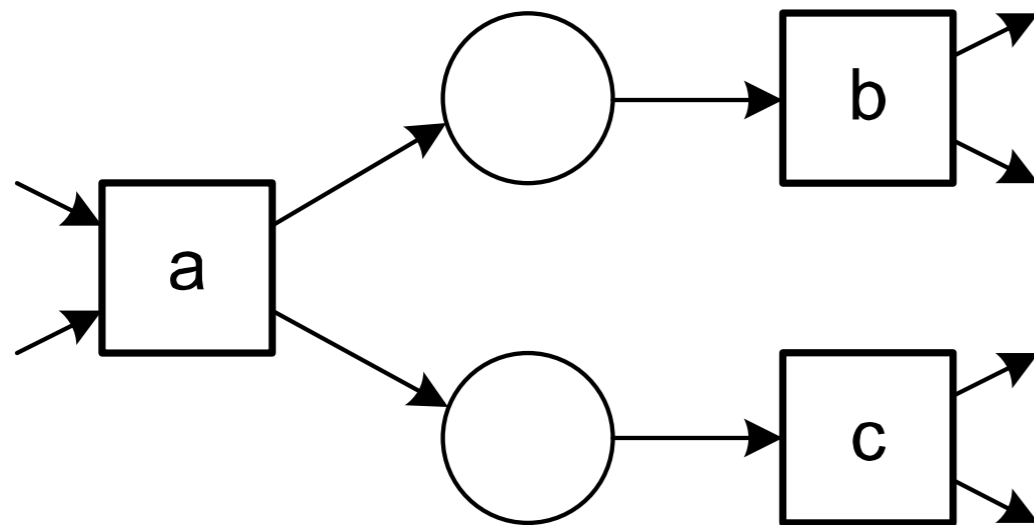


(b) XOR-split pattern:
 $a \rightarrow b$, $a \rightarrow c$, and $b \# c$

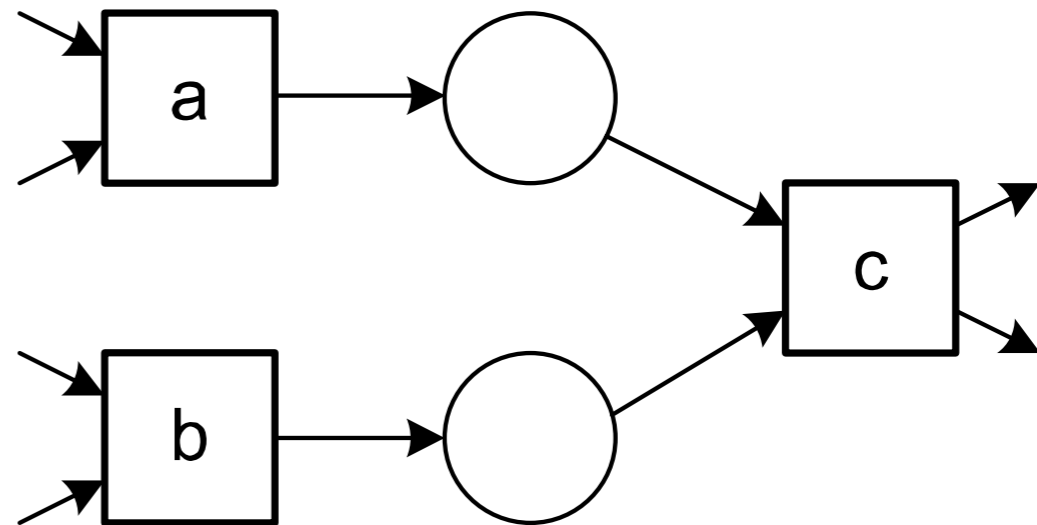
(c) XOR-join pattern:
 $a \rightarrow c$, $b \rightarrow c$, and $a \# b$

Patterns

Footprints are useful to discover typical patterns of activities in the corresponding process model



(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 $a \parallel b$

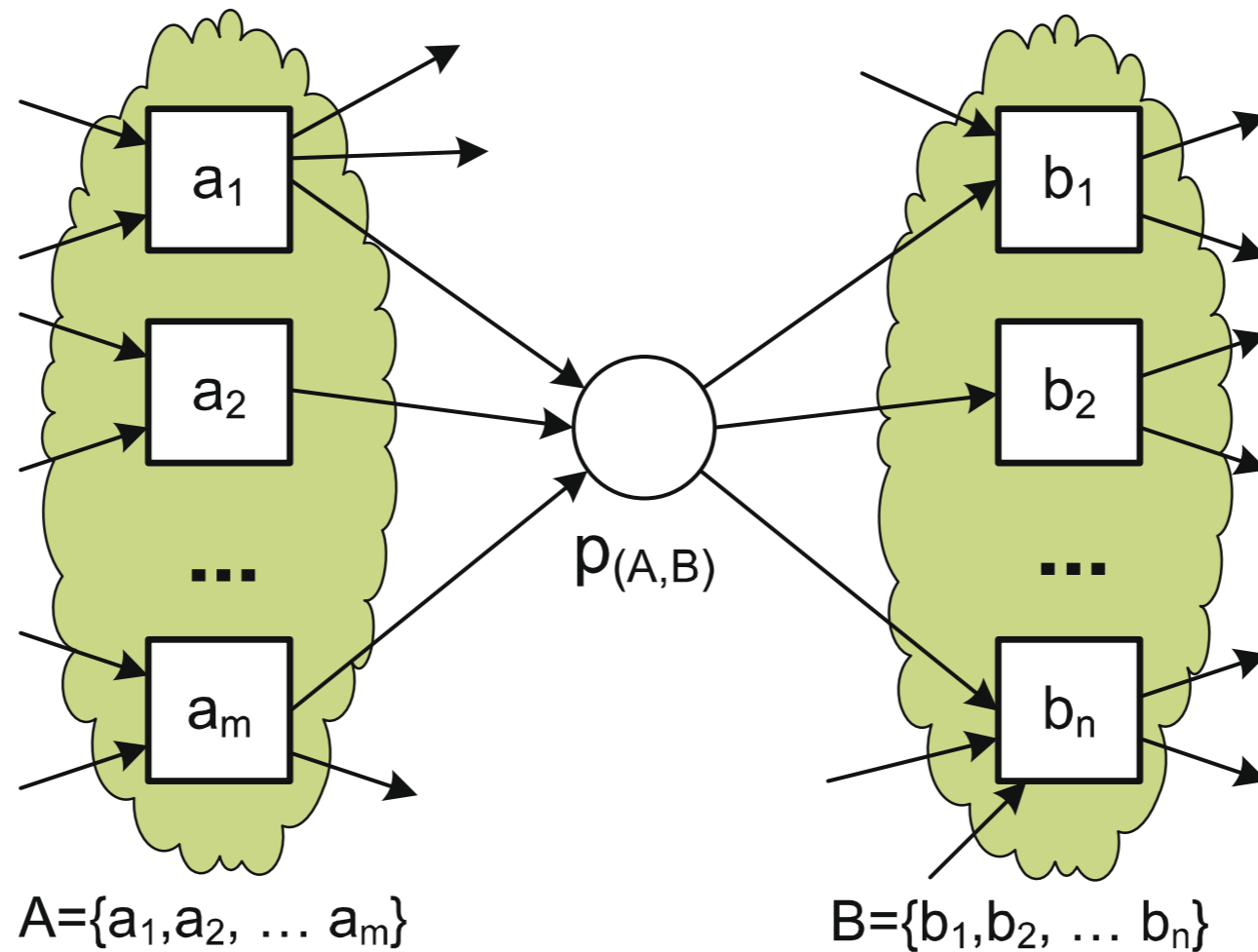
The Algorithm

- (1) $T_L = \{t \in T \mid \exists \sigma \in L \ t \in \sigma\}$ **transitions**
- (2) $T_I = \{t \in T \mid \exists \sigma \in L \ t = \text{first}(\sigma)\}$ **start event**
- (3) $T_O = \{t \in T \mid \exists \sigma \in L \ t = \text{last}(\sigma)\}$ **end event**
- (4) $X_L = \{(A, B) \mid A \subseteq T_L \wedge A \neq \emptyset \wedge B \subseteq T_L \wedge B \neq \emptyset \wedge \forall a \in A \forall b \in B \ a \rightarrow_L b \wedge \forall a_1, a_2 \in A \ a_1 \#_L a_2 \wedge \forall b_1, b_2 \in B \ b_1 \#_L b_2\}$ **decision point**
- (5) $Y_L = \{(A, B) \in X_L \mid \forall (A', B') \in X_L \ A \subseteq A' \wedge B \subseteq B' \implies (A, B) = (A', B')\}$ **max decision point**
- (6) $P_L = \{p_{(A, B)} \mid (A, B) \in Y_L\} \cup \{i_L, o_L\}$ **places**
- (7) $F_L = \{(a, p_{(A, B)}) \mid (A, B) \in Y_L \wedge a \in A\} \cup \{(p_{(A, B)}, b) \mid (A, B) \in Y_L \wedge b \in B\} \cup \{(i_L, t) \mid t \in T_I\} \cup \{(t, o_L) \mid t \in T_O\}$ **arcs**
- (8) $\alpha(L) = (P_L, T_L, F_L)$ **net**

The Core of the Algorithm: Steps 4, 5

	a_1	a_2	...	a_m	b_1	b_2	...	b_n
a_1	#	#	...	#	→	→	...	→
a_2	#	#	...	#	→	→	...	→
...
a_m	#	#	...	#	→	→	...	→
b_1	←	←	...	←	#	#	...	#
b_2	←	←	...	←	#	#	...	#
...
b_n	←	←	...	←	#	#	...	#

The Core of the Algorithm: Step 4, 5



The Algorithm: Example

	<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}	$\#_{L_1}$

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

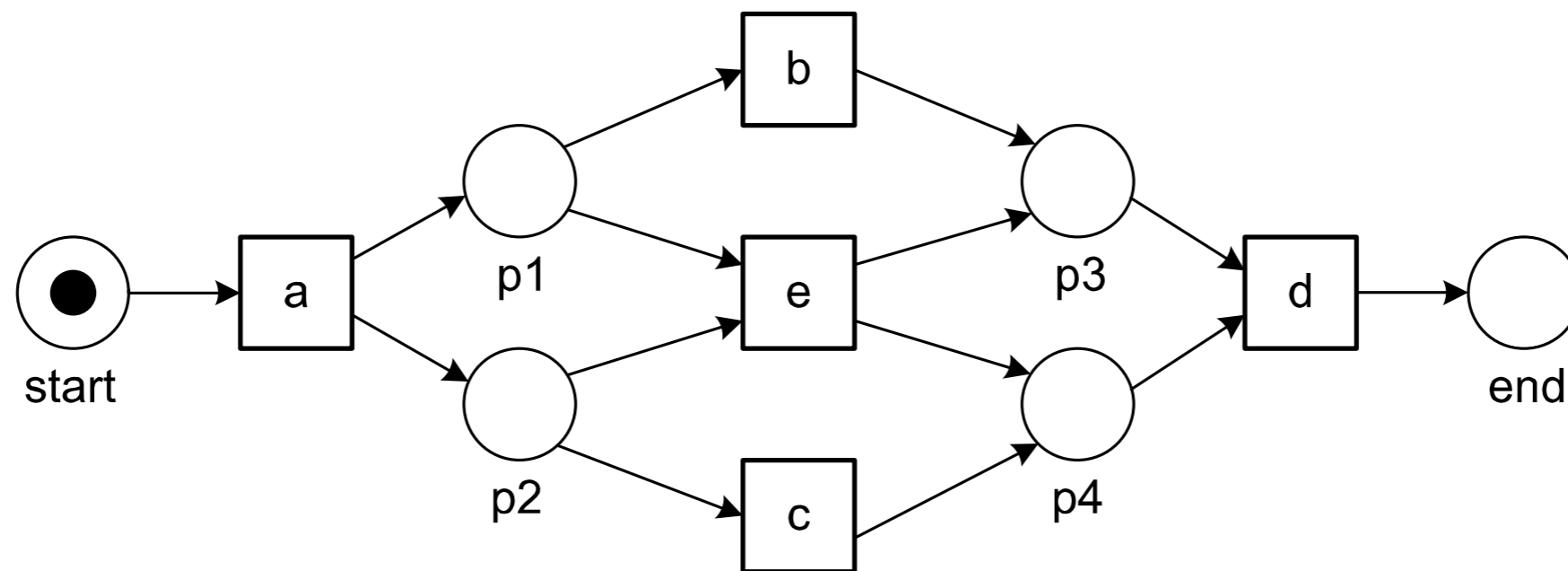
$$(\{b\}, \{d\}), (\{c\}, \{d\}), (\{e\}, \{d\}), (\{b, e\}, \{d\}), (\{c, e\}, \{d\})\}$$

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

The Algorithm: Example

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

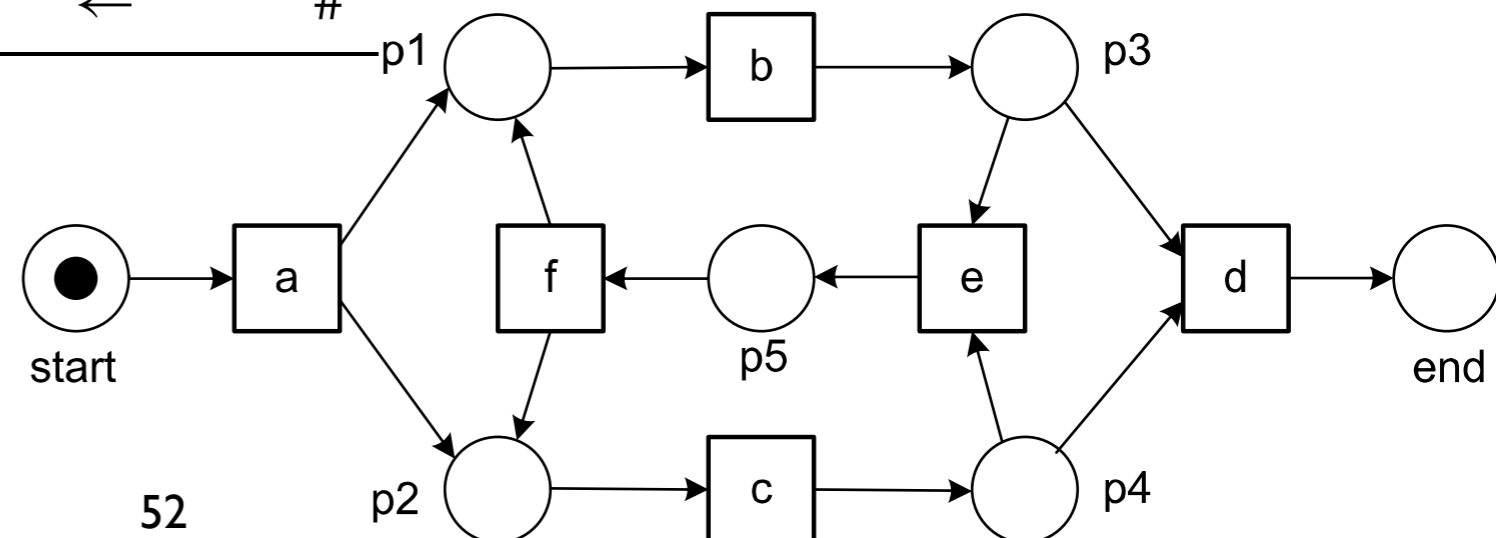
$$Y_{L_1} = \{(\{a\}, \{b, e\}), (\{a\}, \{c, e\}), (\{b, e\}, \{d\}), (\{c, e\}, \{d\})\}$$



Other Examples

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

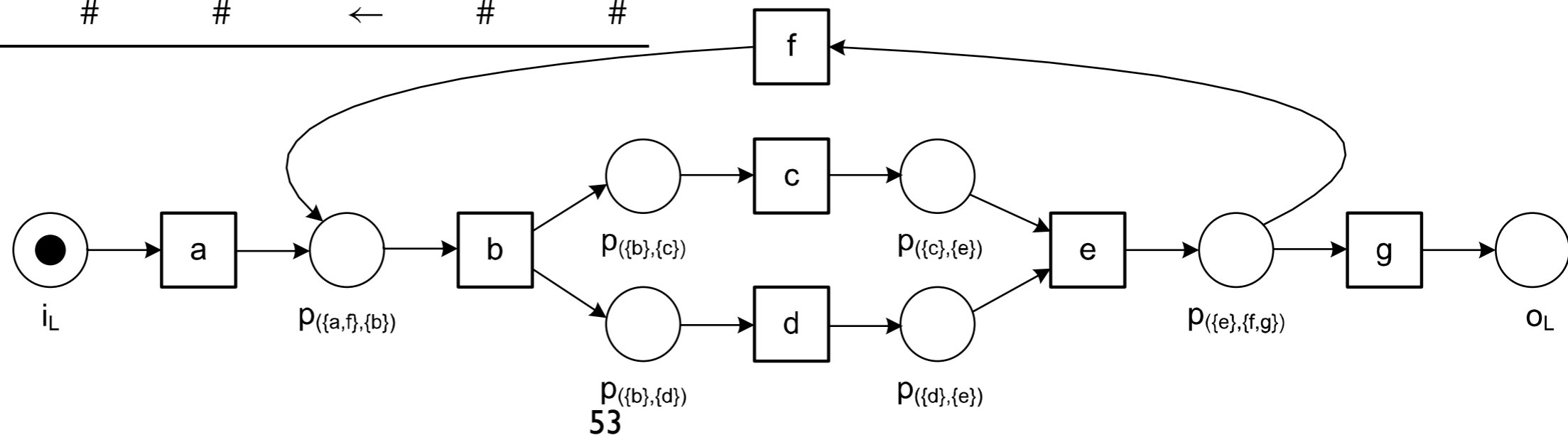
	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
<i>a</i>	#	→	→	#	#	#
<i>b</i>	←	#		→	→	←
<i>c</i>	←		#	→	→	←
<i>d</i>	#	←	←	#	#	#
<i>e</i>	#	←	←	#	#	→
<i>f</i>	#	→	→	#	←	#



Other Examples

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

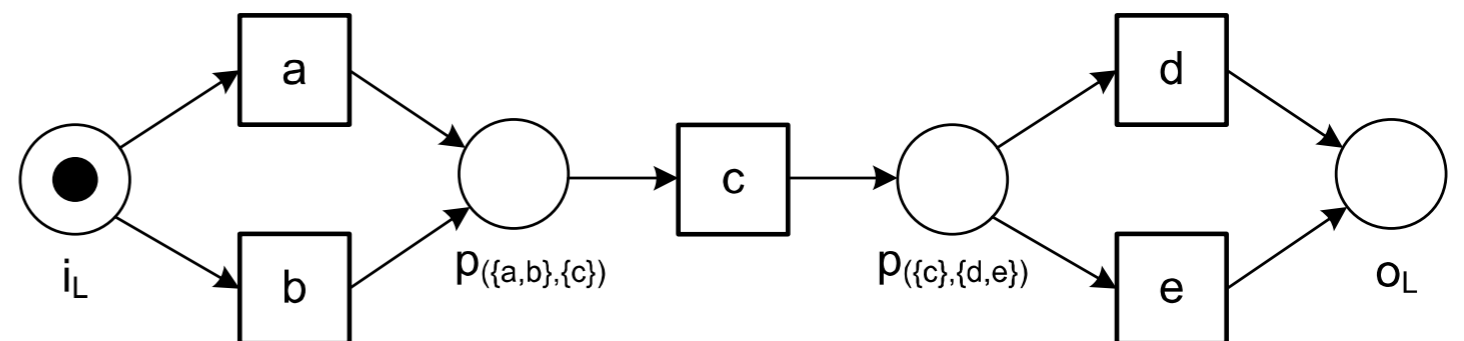
	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>	<i>g</i>
<i>a</i>	#	→	#	#	#	#	#
<i>b</i>	←	#	→	→	#	←	#
<i>c</i>	#	←	#		→	#	#
<i>d</i>	#	←		#	→	#	#
<i>e</i>	#	#	←	←	#	→	→
<i>f</i>	#	→	#	#	←	#	#
<i>g</i>	#	#	#	#	←	#	#



Other Examples

$$L_4 = [\langle a, c, d \rangle^{45}, \langle b, c, d \rangle^{42}, \langle a, c, e \rangle^{38}, \langle b, c, e \rangle^{22}]$$

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



Other Examples

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

	<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>	#	←	#	#	←	#

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

$$T_I = \{a\}$$

$$T_I = \{f\}$$

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

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

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

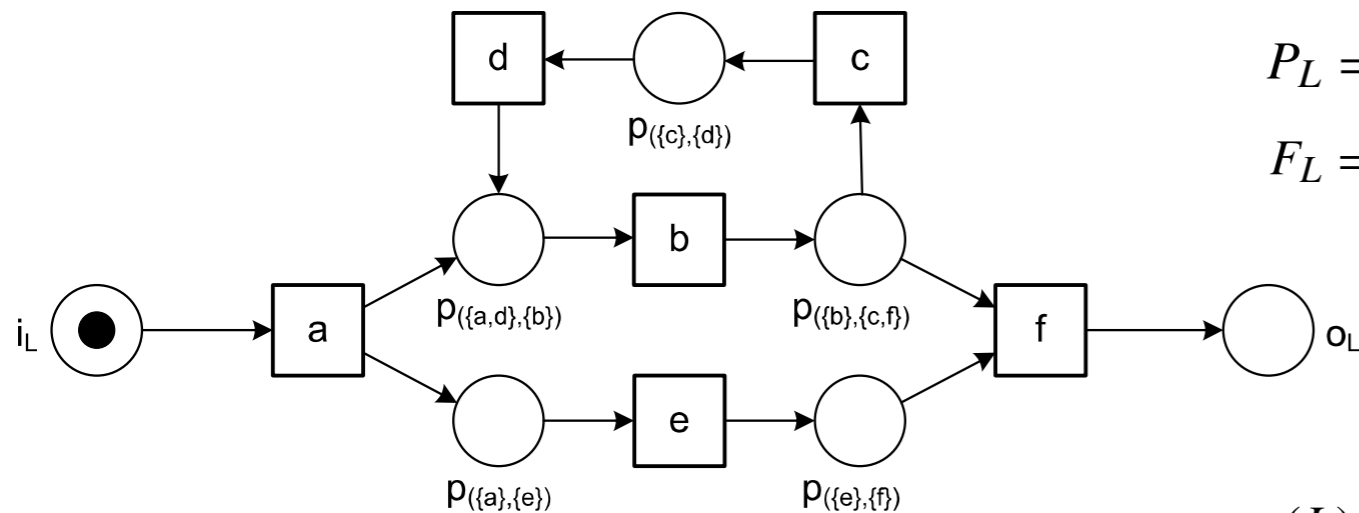
$$F_L = \{(a, p(\{a\}, \{e\})), (p(\{a\}, \{e\}), e), (c, p(\{c\}, \{d\})), (p(\{c\}, \{d\}), d),$$

$$(e, p(\{e\}, \{f\})), (p(\{e\}, \{f\}), f), (a, p(\{a, d\}, \{b\})), (d, p(\{a, d\}, \{b\})),$$

$$(p(\{a, d\}, \{b\}), b), (b, p(\{b\}, \{c, f\})), (p(\{b\}, \{c, f\}), c), (p(\{b\}, \{c, f\}), f),$$

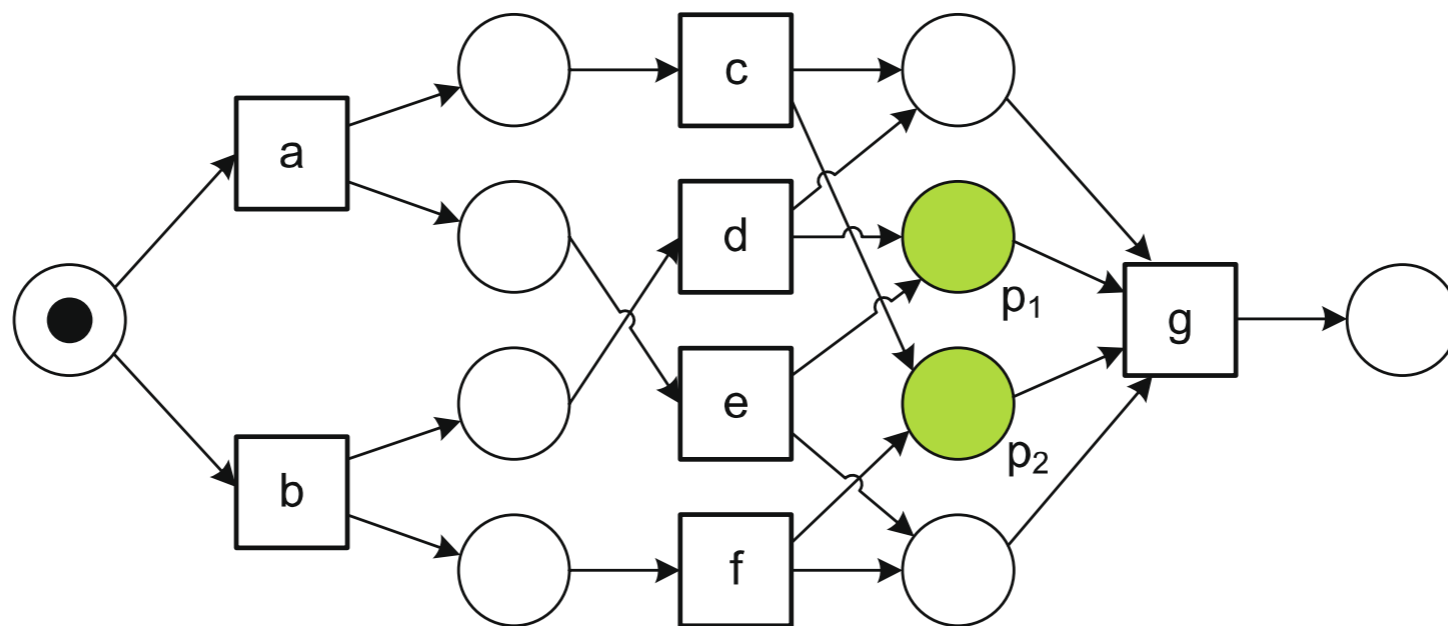
$$(i_L, a), (f, o_L)\}$$

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



Limitation: Implicit Places

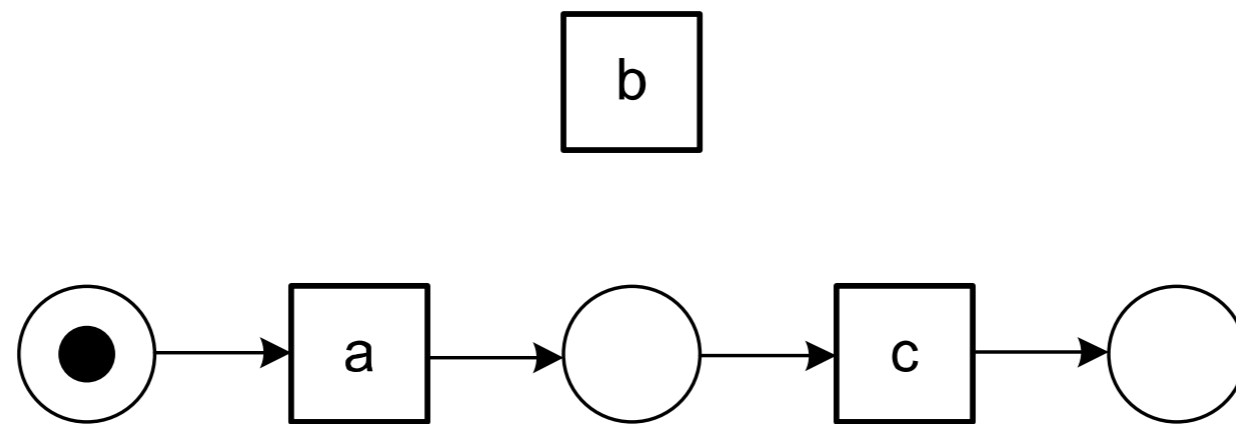
$$L_6 = [\langle a, c, e, g \rangle^2, \langle a, e, c, g \rangle^3, \langle b, d, f, g \rangle^2, \langle b, f, d, g \rangle^4]$$



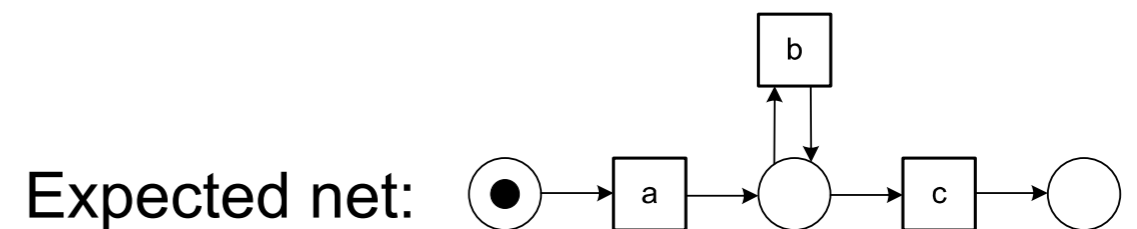
p1 and *p2* are redundant

Limitation: Short Loop

$$L_7 = [\langle a, c \rangle^2, \langle a, b, c \rangle^3, \langle a, b, b, c \rangle^2, \langle a, b, b, b, b, c \rangle^1]$$



b is disconnected from the model



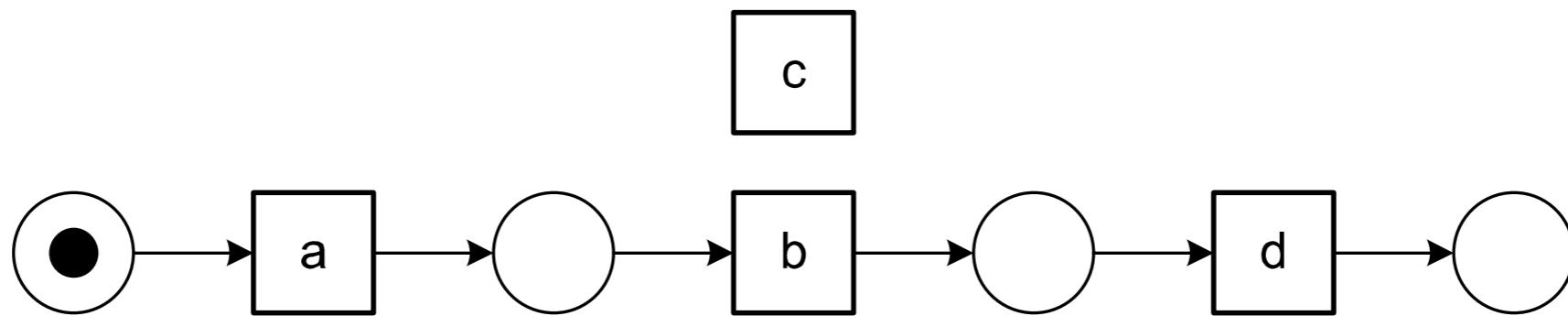
Limitation: Short Loop

$$L_8 = [\langle a, b, d \rangle^3, \langle a, b, c, b, d \rangle^2, \langle a, b, c, b, c, b, d \rangle]$$

$a \rightarrow_{L_8} b,$

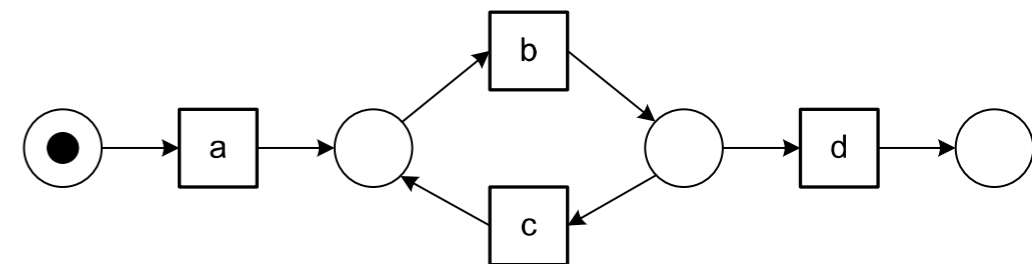
$b \rightarrow_{L_8} d,$

$b \parallel_{L_8} c.$



c is disconnected from the model

Expected net:

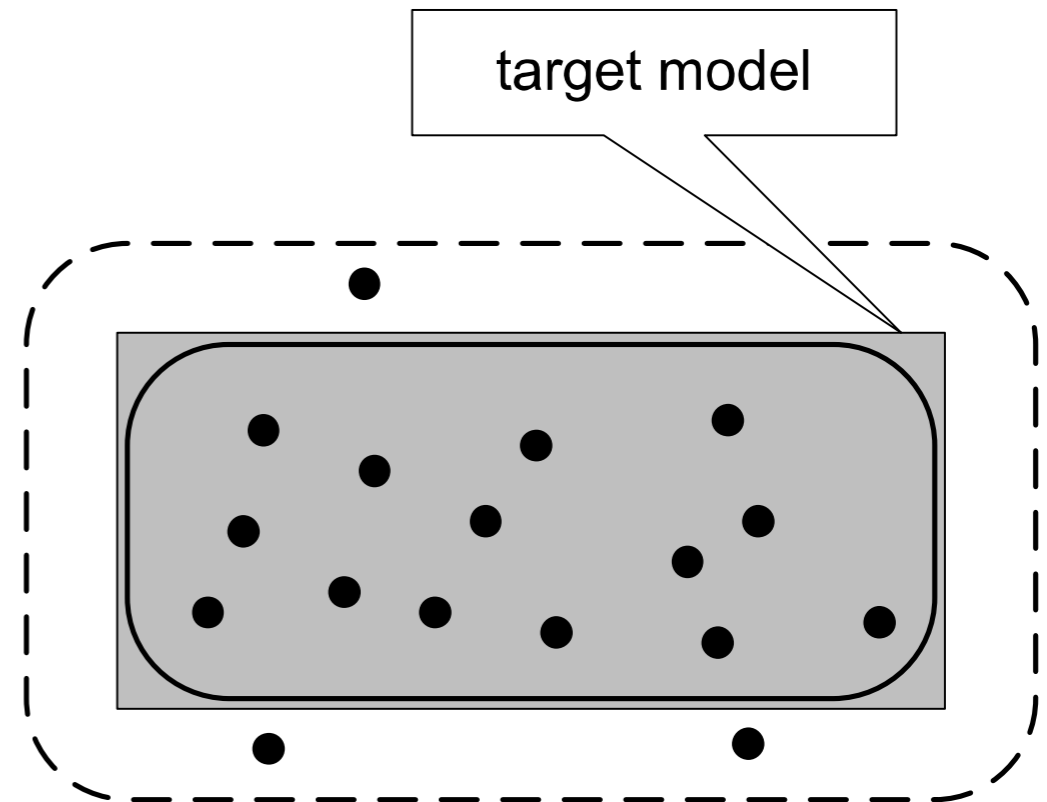
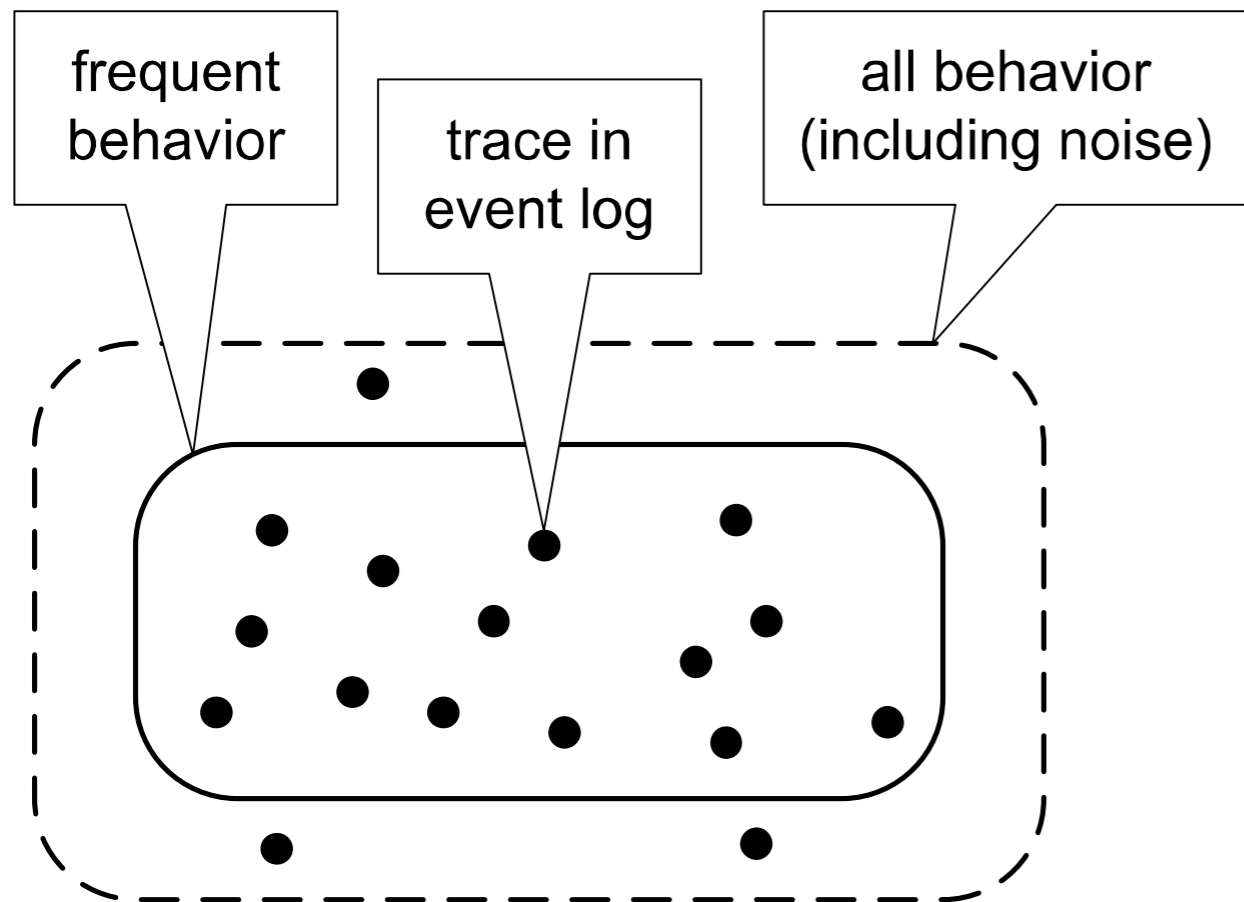


Limitation: Noise

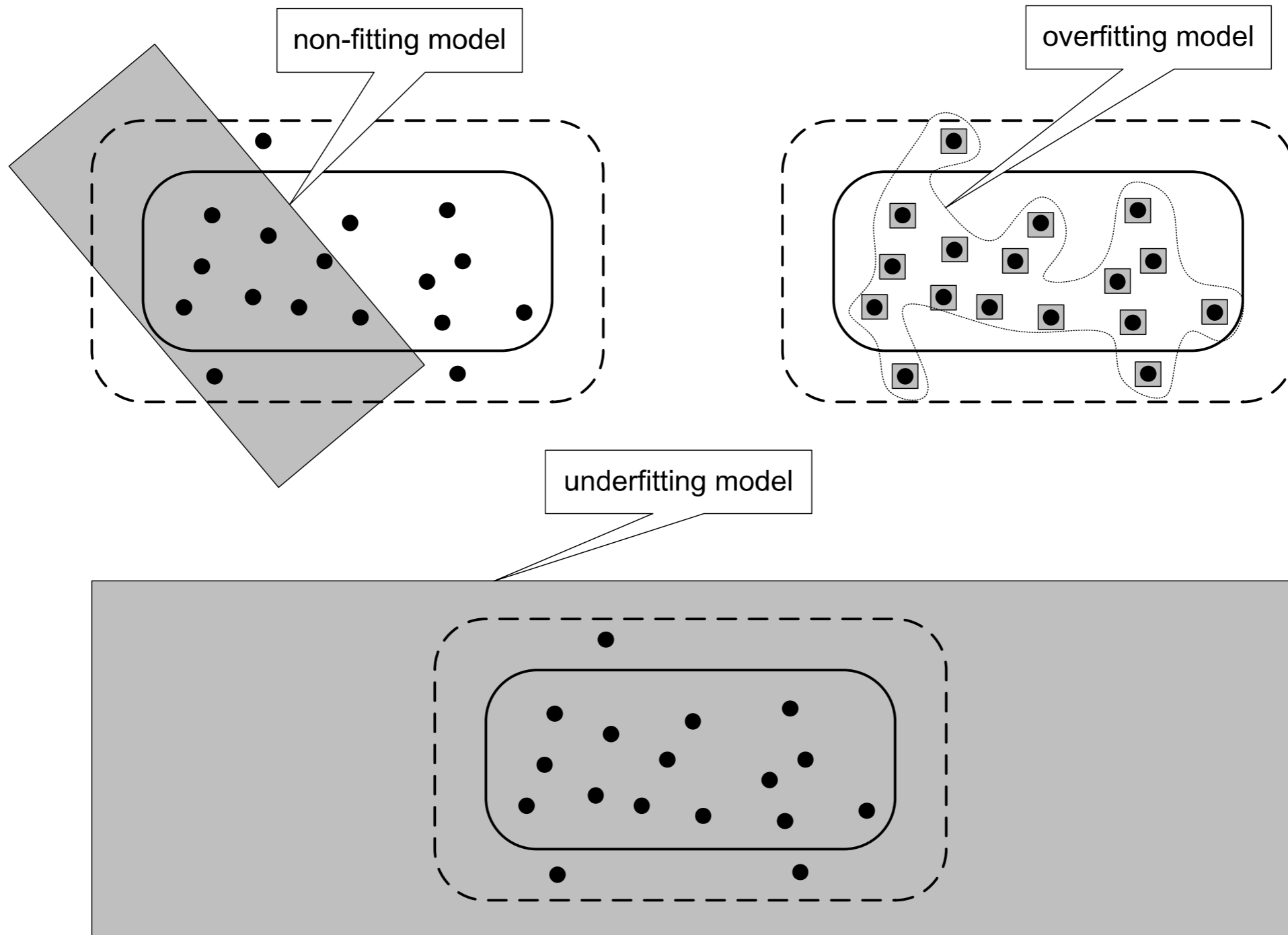
We use the term “noise” to refer to rare and infrequent behavior rather than errors related to event logging.

For example, frequencies are not taken into account by the α -algorithm (discard less frequent traces?).

Limitation: Noise



Limitation: Noise



Limitation: Incompleteness

Whereas noise refers to the problem of having “*too much data*” (describing rare behavior), **(in)completeness** refers to the problem of having “*too little data*”.

Process models typically allow for an exponential or even infinite number of different traces (in case of loops).

Moreover, some traces may have a much lower probability than others. Therefore, it is unrealistic to assume that every possible trace is present in the event log.

Limitation: Incompleteness

The α -algorithm uses a **local completeness notion**:

if there are two activities a and b ,
and a can be directly followed by b ,
then this should be observed at least once in the log.

Conformance Checking

Two Angles

Conformance check is based on the comparison between an event log and a process model.

(Un)desirable deviations can be detected.

First viewpoint (the model is supposed to be **descriptive**):
the model does not capture the real behavior
("the model is wrong, how to improve it?")

Second viewpoint (the model is **normative**):
reality deviates from the desired model
("the event log is wrong, how to impose control?").

Measures and Diagnostic

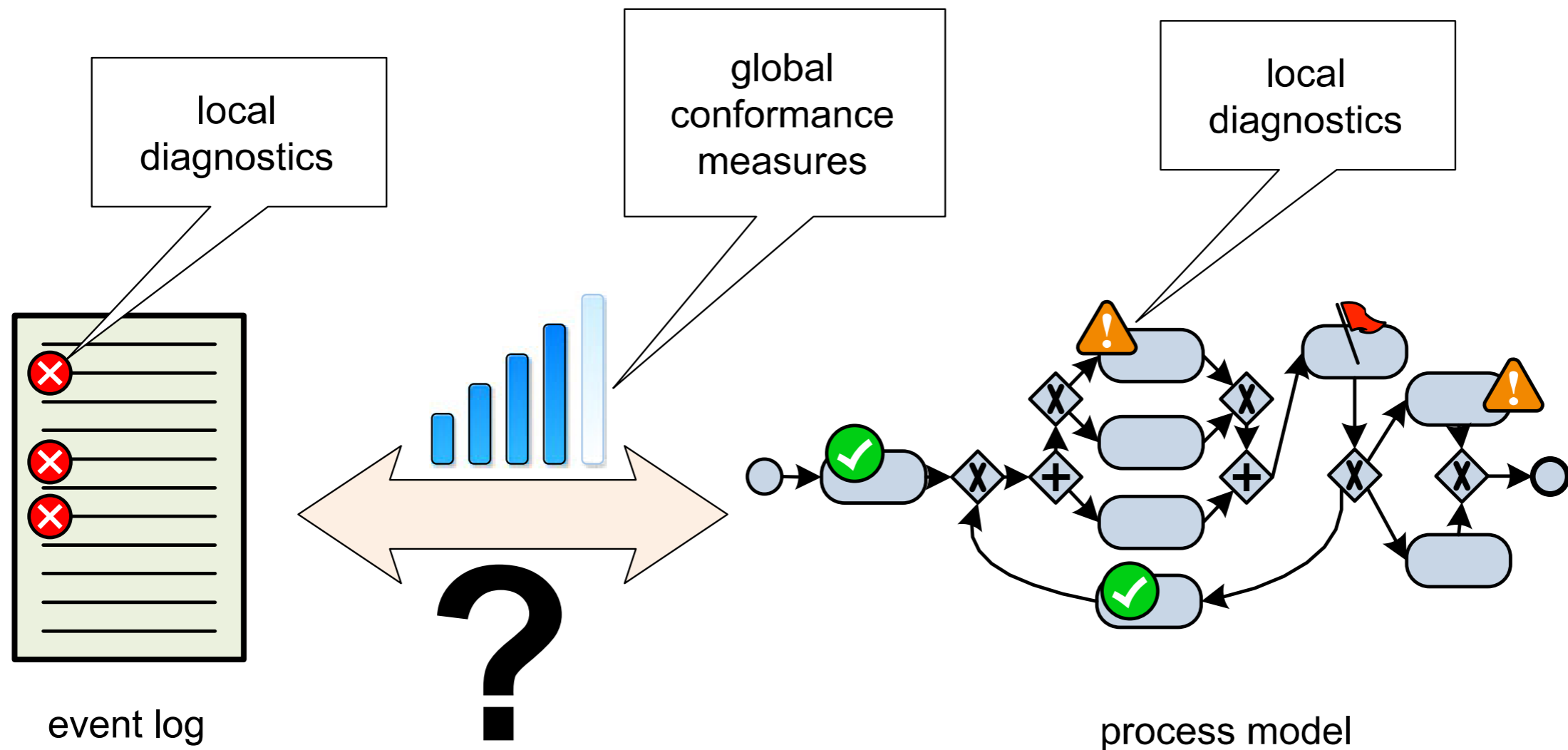


Fig. 7.1 Conformance checking: comparing observed behavior with modeled behavior. Global conformance measures quantify the overall conformance of the model and log. Local diagnostics are given by highlighting the nodes in the model where model and log disagree. Cases that do not fit are highlighted in the visualization of the log 66

Business Alignment

The goal of business alignment is to make sure that the information systems and the real business processes are well aligned.

People should be supported by the information system rather than work behind its back to get things done.

Process mining can assist in improving the alignment of information systems, business processes, and the organization.

By analyzing the real processes and diagnosing discrepancies, new insights can be gathered showing how to improve the support by information systems.

Auditing

The term auditing refers to the evaluation of organizations and their processes.

Audits are performed to ascertain the validity and reliability of information about these organizations and associated processes.

This is done to check whether business processes are executed within certain boundaries set by managers, governments, and other stakeholders.

Rules violations may indicate fraud, malpractice, risks, and inefficiencies.

New Forms of Auditing

However, today detailed information about processes is being recorded in the form of event logs, audit trails, transaction logs, databases, data warehouses, etc.

All events in a business process can be evaluated and this can be done while the process is still running.

The availability of log data and advanced process mining techniques enables new forms of auditing, and conformance checking in particular, provide the means to do so.

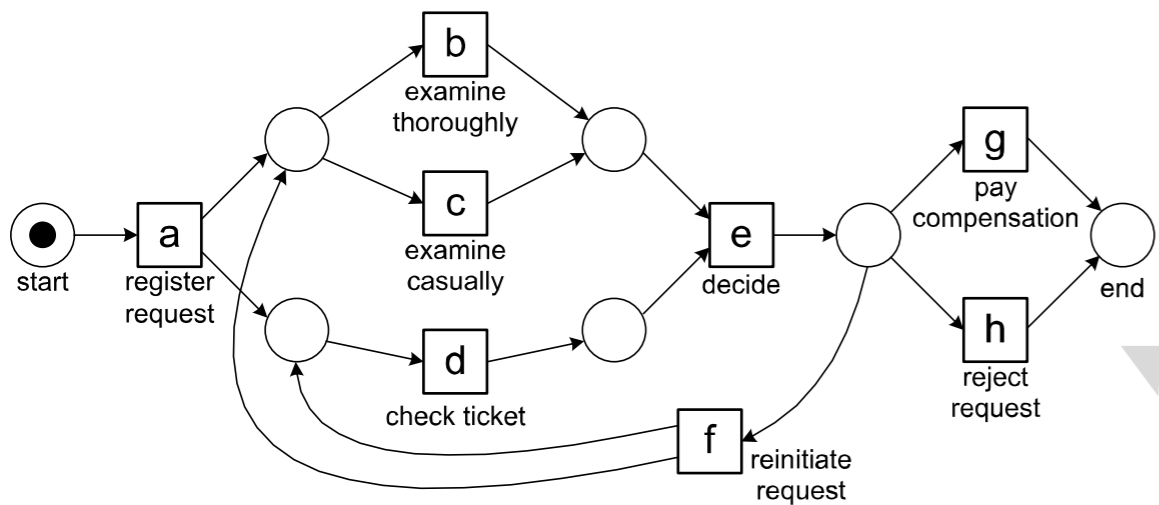
Quality Criteria

We have seen four quality criteria:
fitness, precision, generalization, and simplicity.

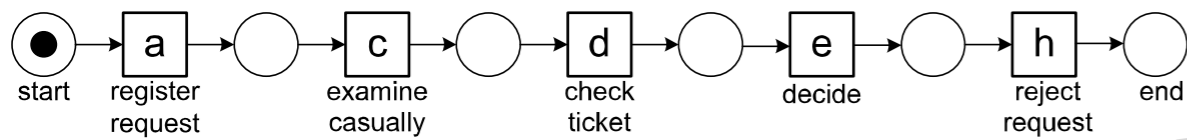
In an example shown, for each of these models, a subjective judgment is given with respect to the four quality criteria. As the models are rather extreme, the scores +/- for the various quality criteria are evident.

We discuss how the notion of fitness can be quantified.

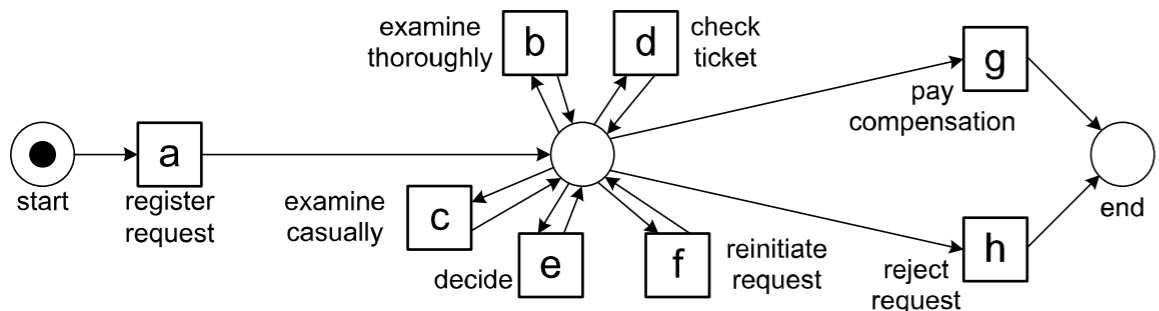
Appropriateness



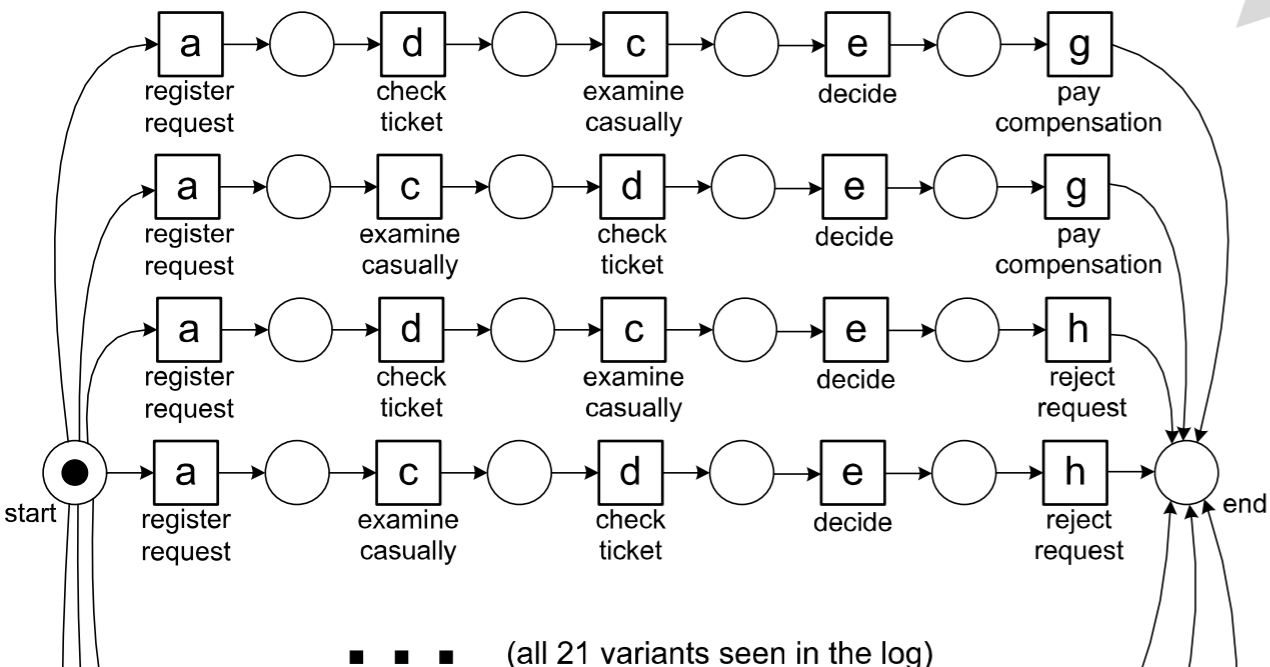
N_1 : fitness = +, precision = +, generalization = +, simplicity = +



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

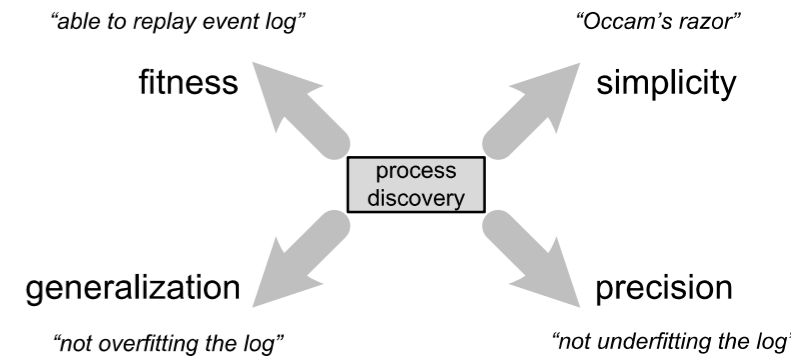


N_3 : fitness = +, precision = -, generalization = +, simplicity = +



■ ■ ■ (all 21 variants seen in the log)

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



Measuring Fitness

However, in a more realistic setting it is much more difficult to judge the quality of a model.

Fitness measures “the proportion of behavior in the event log possible according to the model”.

Of the four quality criteria,
fitness is most related to conformance.

A naïve approach toward conformance checking would be to count the fraction of cases that can be “parsed completely” (i.e., the proportion of cases corresponding to firing sequences leading from [start] to [end]).

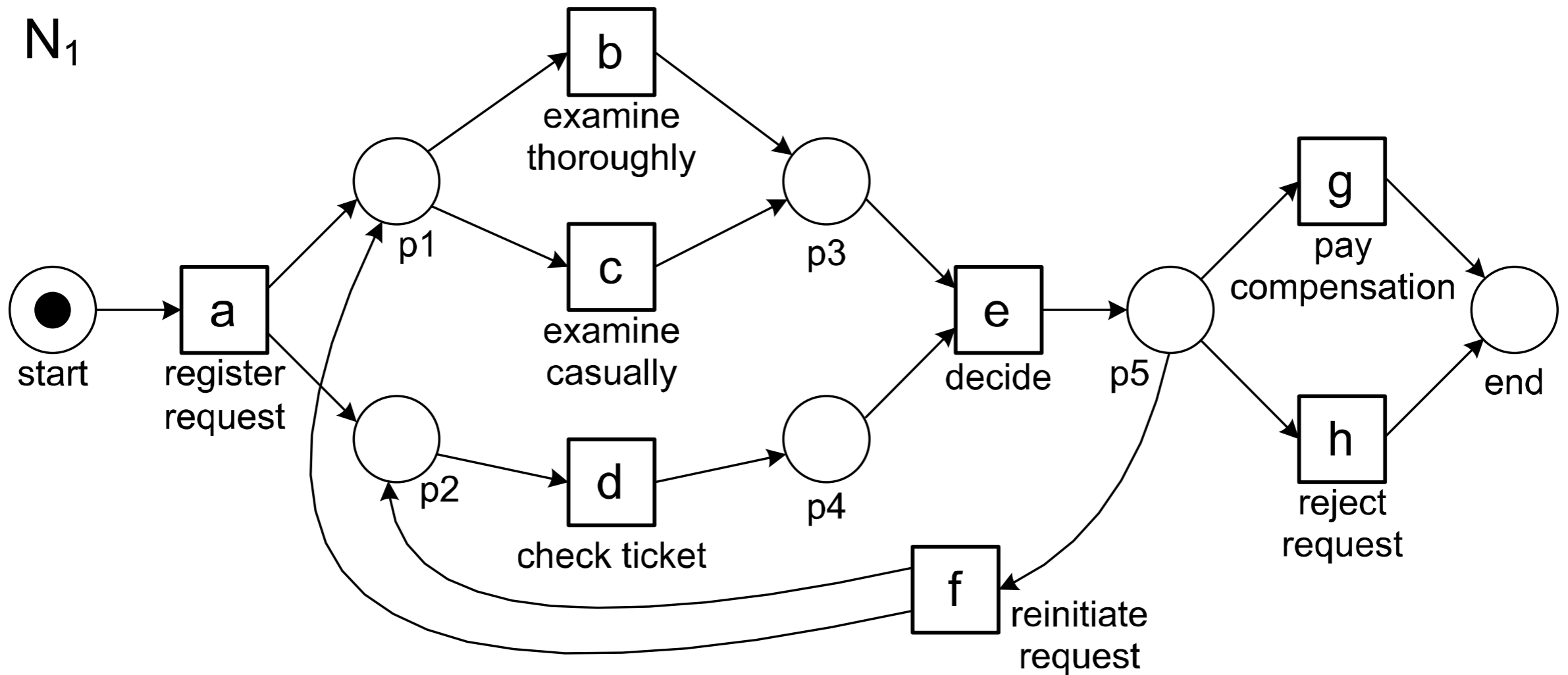
Table 7.1 Event log L_{full} : $a = register\ request$, $b = examine\ thoroughly$, $c = examine\ casually$, $d = check\ ticket$, $e = decide$, $f = reinitiate\ request$, $g = pay\ compensation$, and $h = reject\ request$

1391 cases

Frequency	Reference	Trace
455	σ_1	$\langle a, c, d, e, h \rangle$
191	σ_2	$\langle a, b, d, e, g \rangle$
177	σ_3	$\langle a, d, c, e, h \rangle$
144	σ_4	$\langle a, b, d, e, h \rangle$
111	σ_5	$\langle a, c, d, e, g \rangle$
82	σ_6	$\langle a, d, c, e, g \rangle$
56	σ_7	$\langle a, d, b, e, h \rangle$
47	σ_8	$\langle a, c, d, e, f, d, b, e, h \rangle$
38	σ_9	$\langle a, d, b, e, g \rangle$
33	σ_{10}	$\langle a, c, d, e, f, b, d, e, h \rangle$
14	σ_{11}	$\langle a, c, d, e, f, b, d, e, g \rangle$
11	σ_{12}	$\langle a, c, d, e, f, d, b, e, g \rangle$
9	σ_{13}	$\langle a, d, c, e, f, c, d, e, h \rangle$
8	σ_{14}	$\langle a, d, c, e, f, d, b, e, h \rangle$
5	σ_{15}	$\langle a, d, c, e, f, b, d, e, g \rangle$
3	σ_{16}	$\langle a, c, d, e, f, b, d, e, f, d, b, e, g \rangle$
2	σ_{17}	$\langle a, d, c, e, f, d, b, e, g \rangle$
2	σ_{18}	$\langle a, d, c, e, f, b, d, e, f, b, d, e, g \rangle$
1	σ_{19}	$\langle a, d, c, e, f, d, b, e, f, b, d, e, h \rangle$
1	σ_{20}	$\langle a, d, b, e, f, b, d, e, f, d, b, e, g \rangle$
1	σ_{21}	$\langle a, d, c, e, f, d, b, e, f, c, d, e, f, d, b, e, g \rangle$

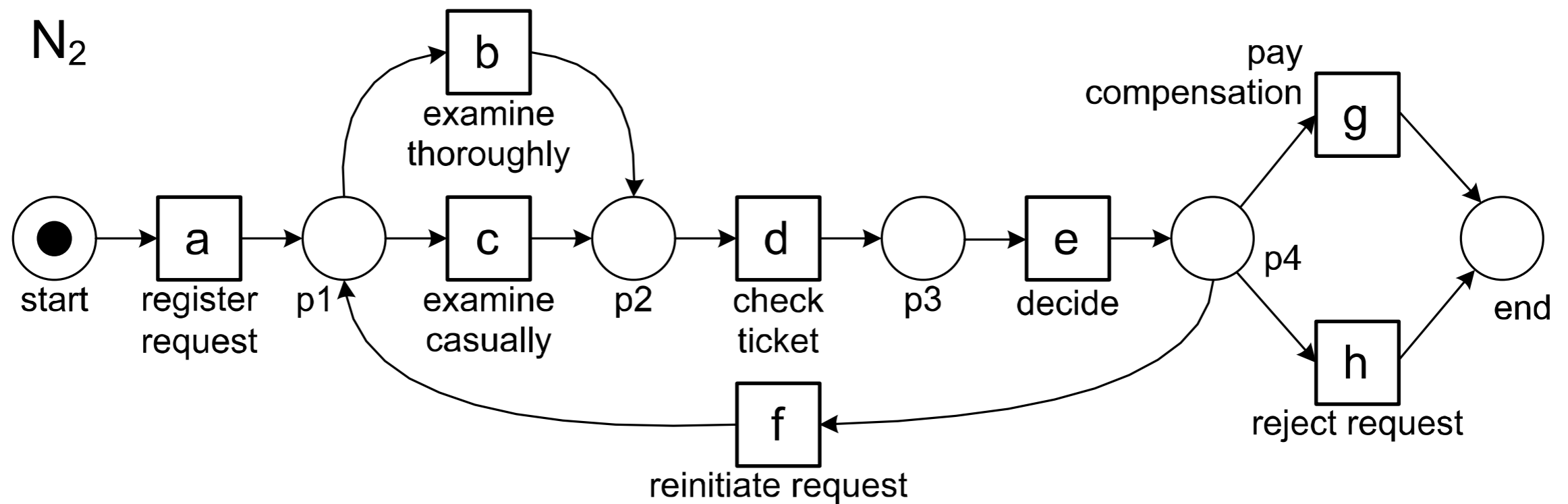
Example

Example N1



naïve fitness $\frac{1391}{1391} = 1$

Example N2



443 cases do not correspond to a firing sequence

$\langle a, d, c, e, h \rangle^{177}$

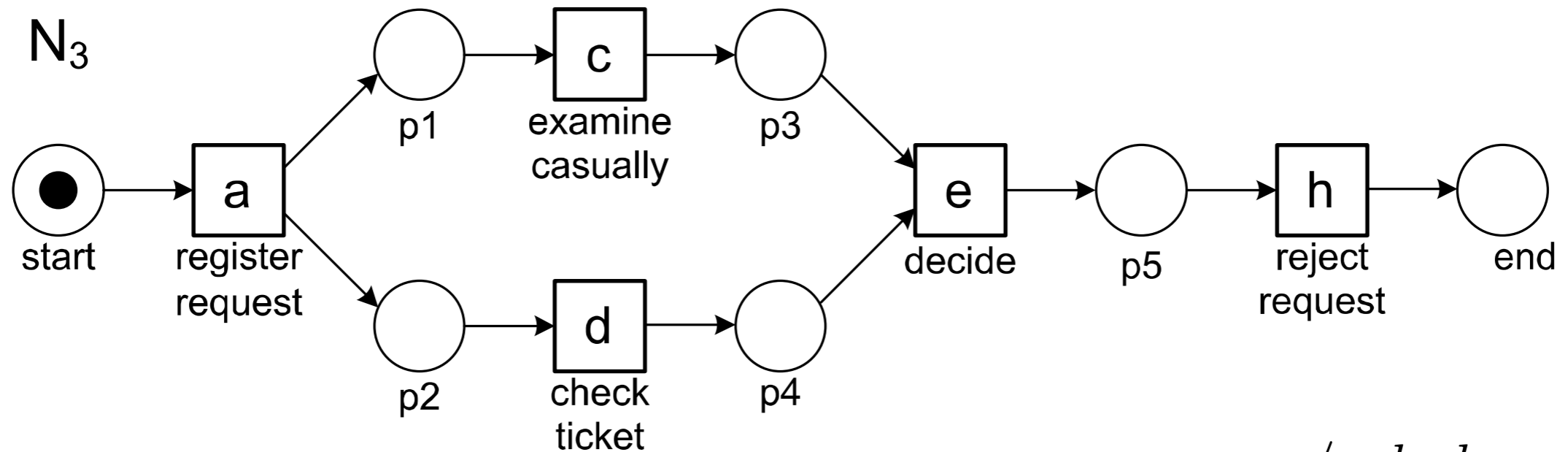
$\langle a, d, c, e, g \rangle^{82}$

$\langle a, d, b, e, h \rangle^{56}$

...

naïve fitness $\frac{948}{1391} = 0.6815$

Example N3

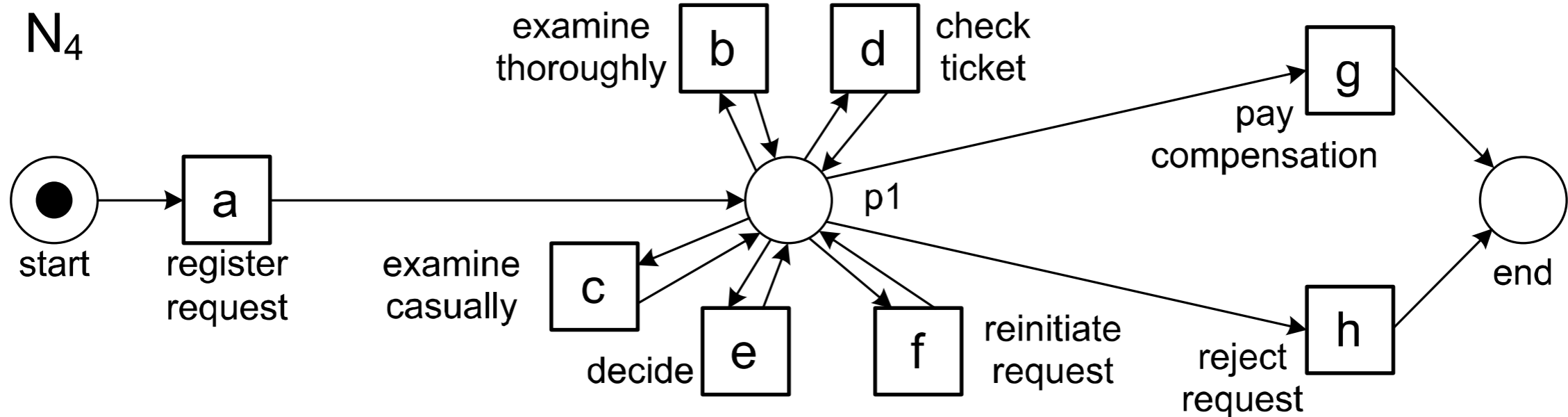


759 cases do not correspond to a firing sequence

- $\langle a, b, d, e, g \rangle^{191}$
- $\langle a, b, d, e, h \rangle^{144}$
- $\langle a, c, d, e, g \rangle^{111}$
- ...

naïve fitness $\frac{632}{1391} = 0.4543$

Example N4



“flower model” (poorly structured)

naïve fitness $\frac{1391}{1391} = 1$

Almost Fitting Traces

This naïve fitness notion seems to be too strict as traces can differ only slightly and not be counted at all.

$$\sigma = \langle a_1, a_2, \dots, a_{100} \rangle$$

Now consider a model that cannot replay σ , but that can replay 99 of the 100 events in σ .

Then, consider another model that can only replay 10 of the 100 events in σ .

Using the naïve fitness metric, the trace would simply be classified as nonfitting for both models without acknowledging that σ was almost fitting in one model and in complete disagreement with the other.

Missing and Remaining Tokens

We introduce a fitness notion defined at the level of events rather than full traces.

In the naïve fitness computation just described, we stopped replaying a trace once we encounter a problem (and mark it as nonfitting).

Let us instead just continue replaying the trace on the model but record all situations where a transition is forced to fire without being enabled, i.e., we count all missing tokens.

Moreover, we record the tokens that remain at the end.

Four Counters

p (produced tokens)

c (consumed tokens)

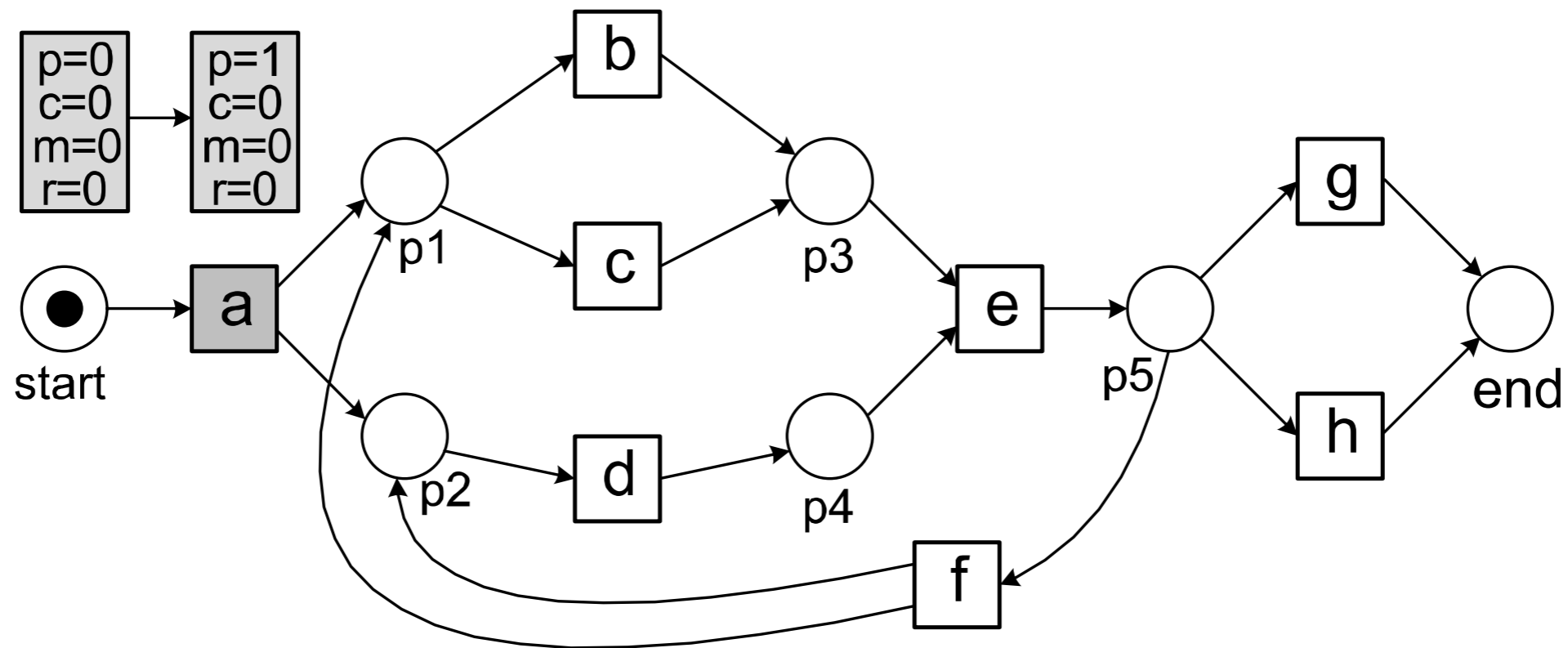
m (missing tokens)

r (remaining tokens)

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

Example

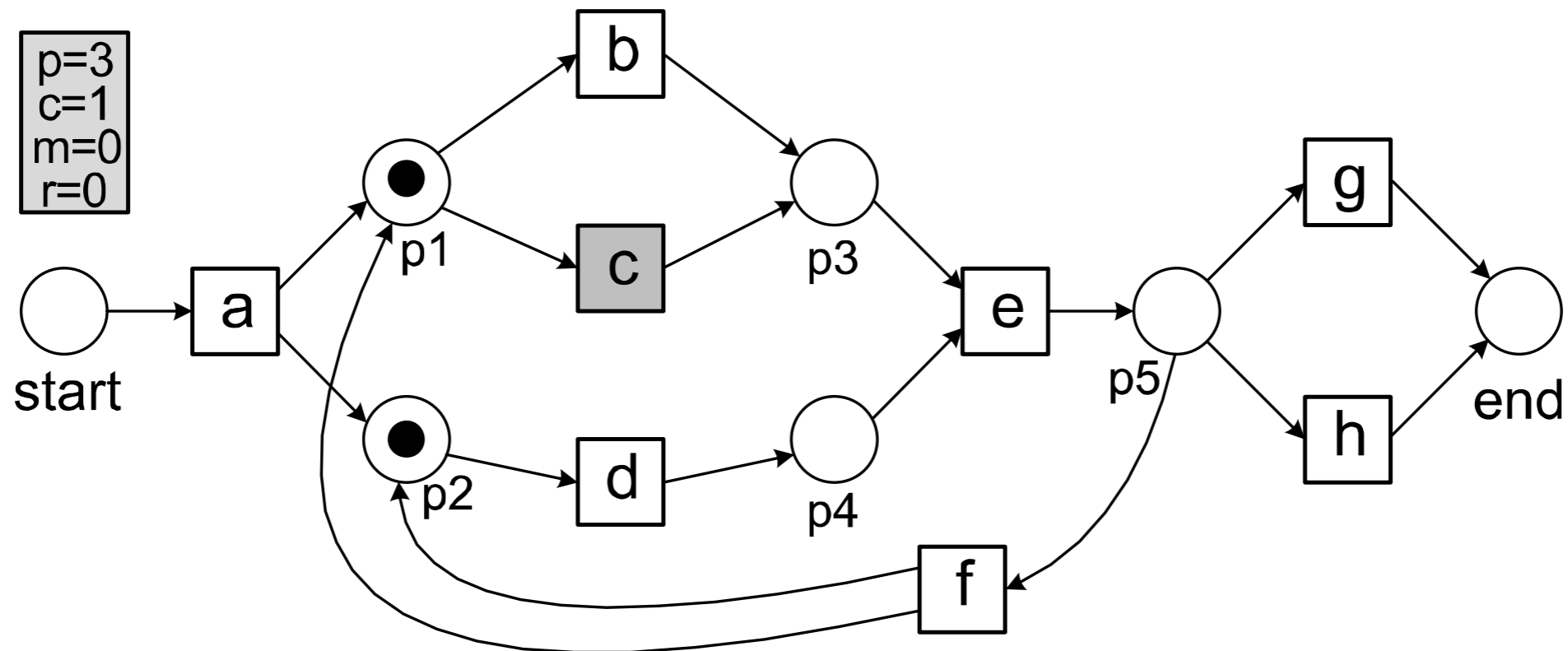
the environment produces a token for place start



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

Example

replaying a is possible
one token is consumed, two produced

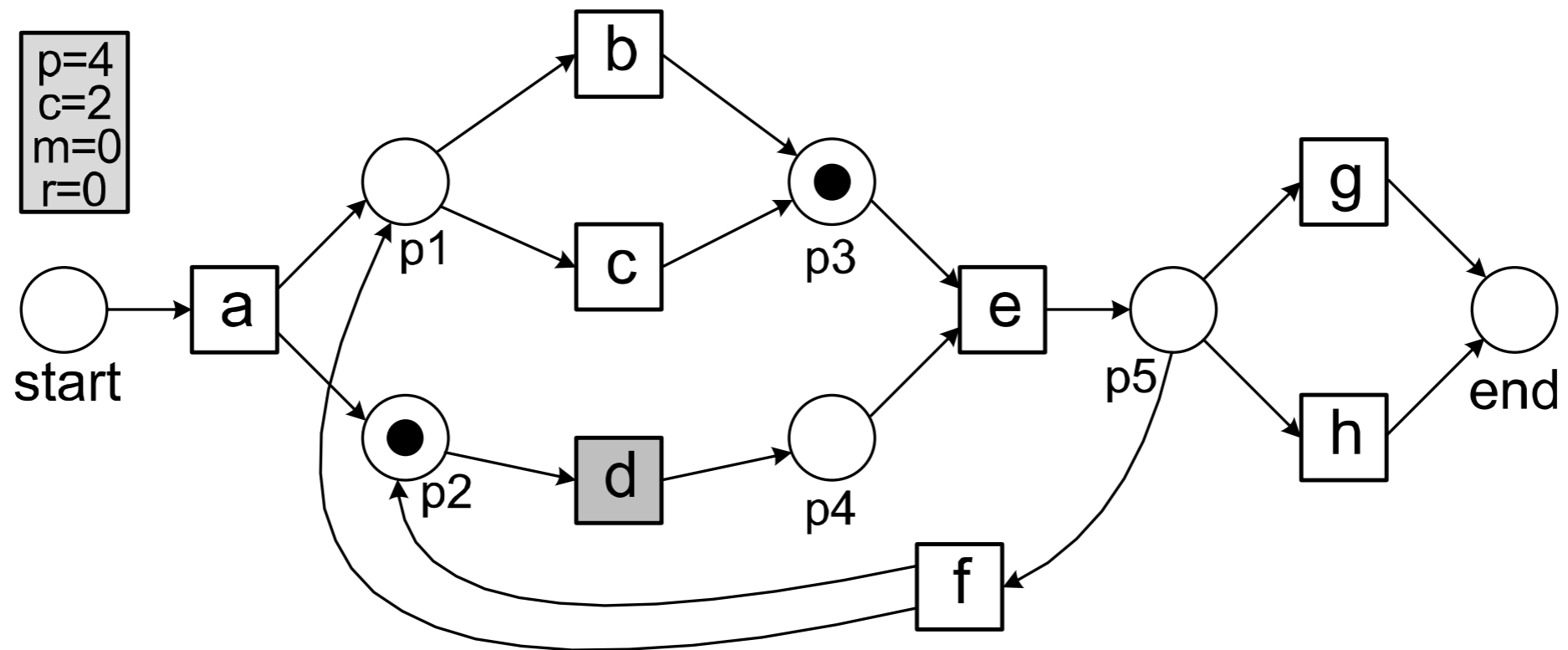


$$\sigma_1 = \langle a, c, d, e, h \rangle$$

Example

replaying c is possible

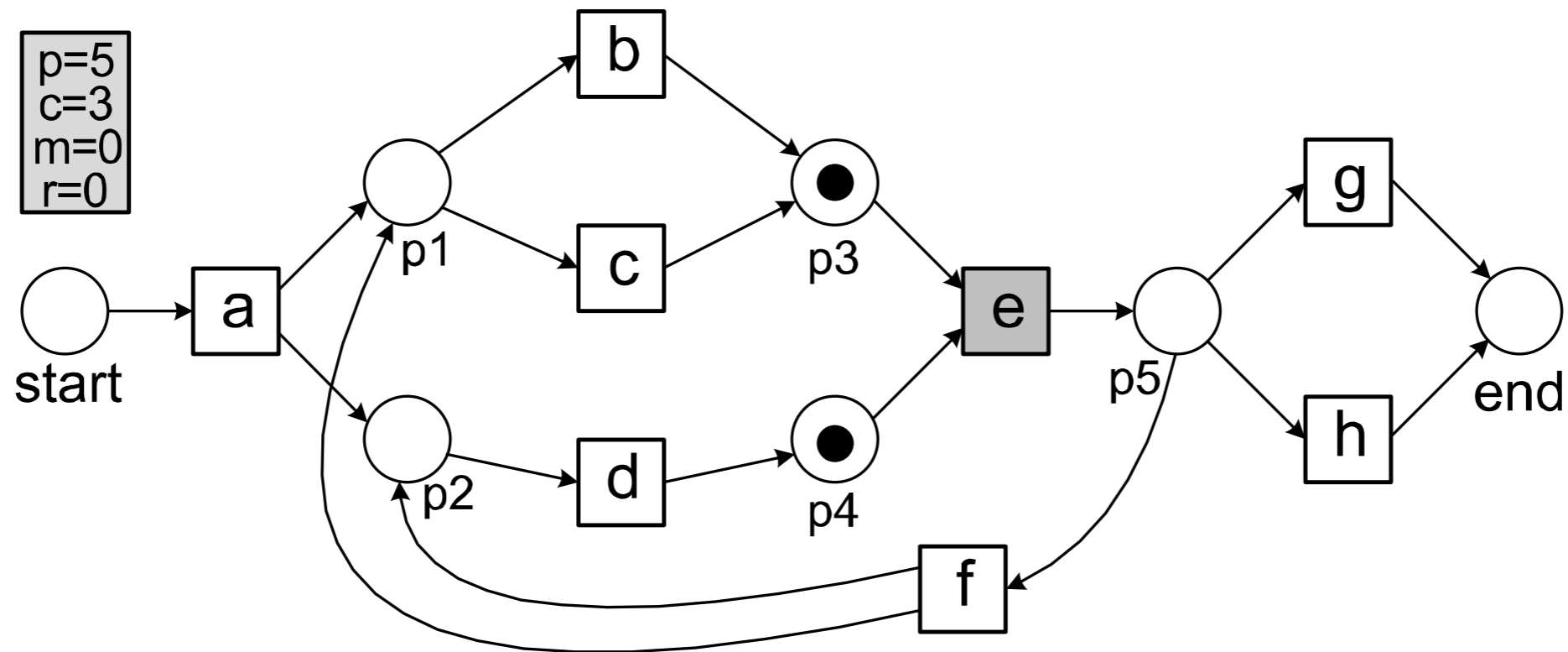
one token is consumed, one produced



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

Example

replaying d is possible
one token is consumed, one produced

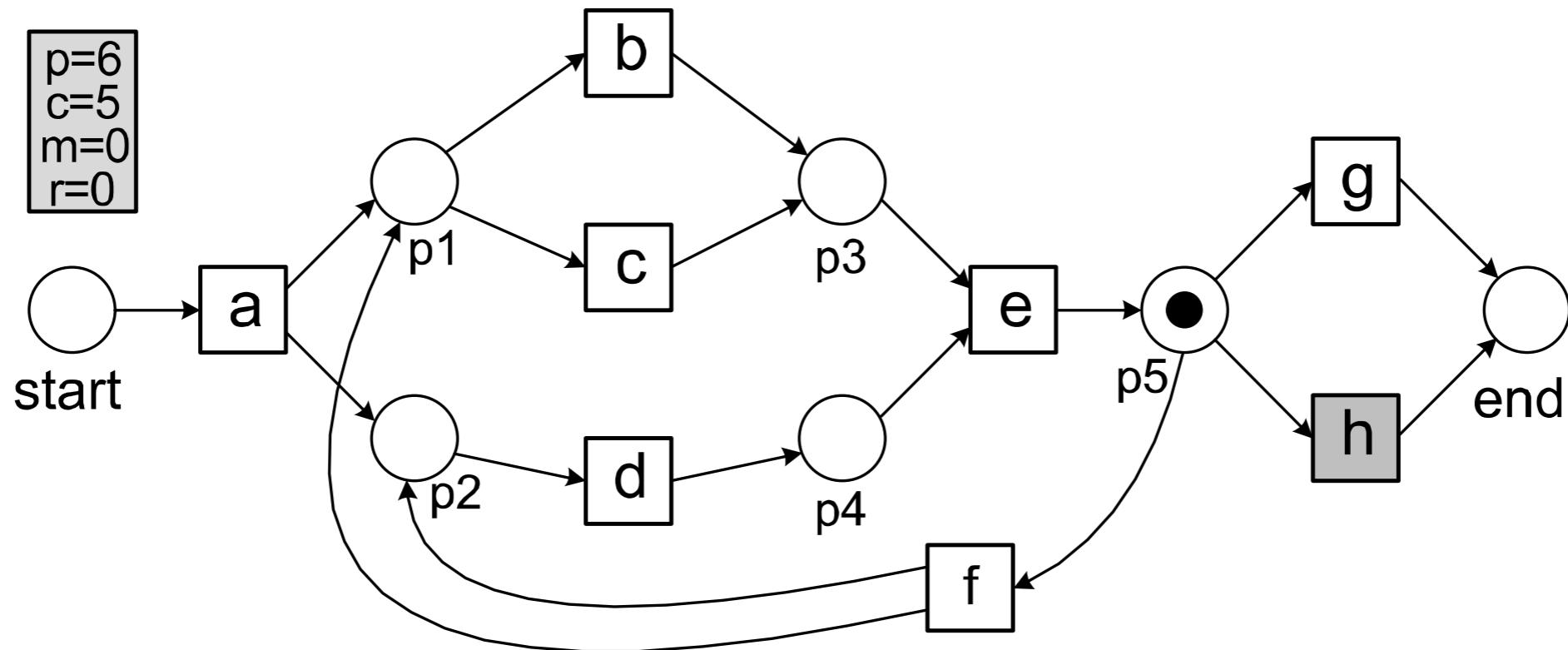


$$\sigma_1 = \langle a, c, d, e, h \rangle$$

Example

replaying e is possible

two tokens are consumed, one produced



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

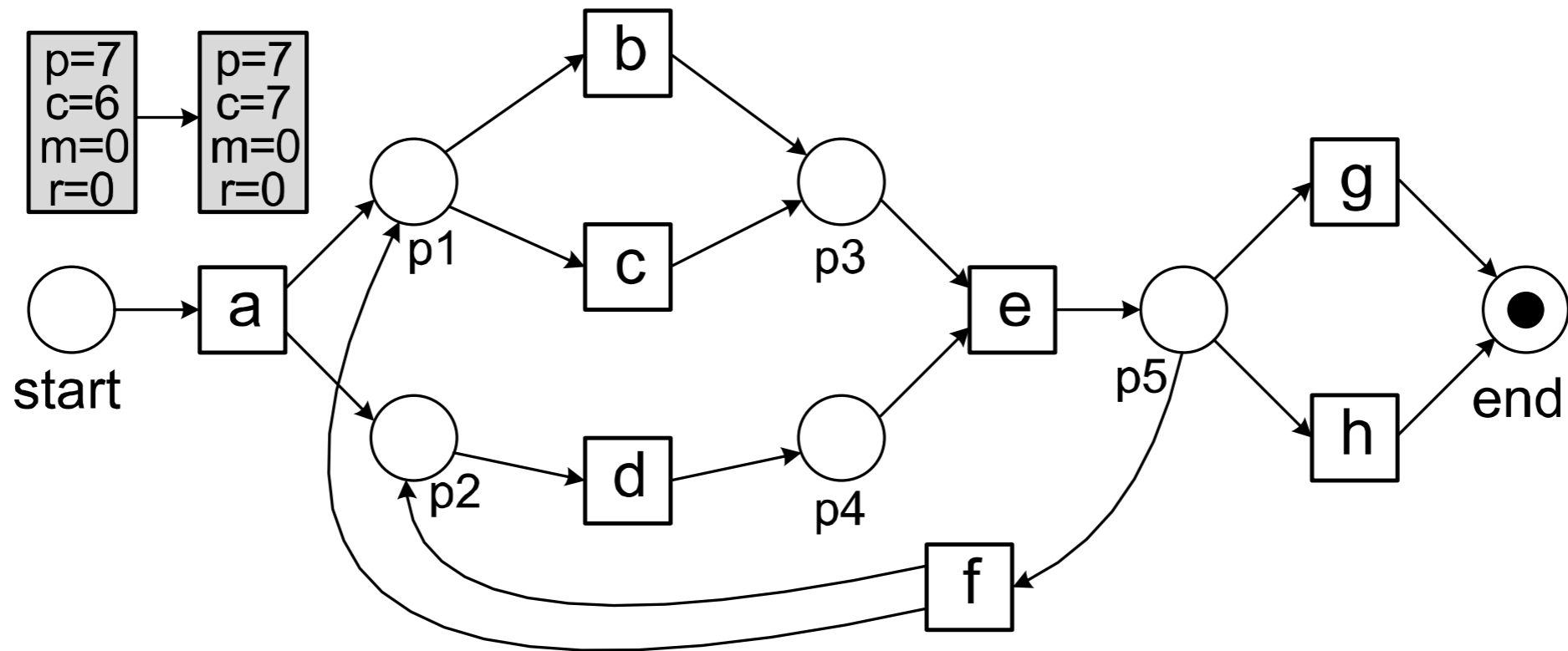
Example

replaying h is possible

one token is consumed, one produced

At the end,

the environment consumes a token from place end.

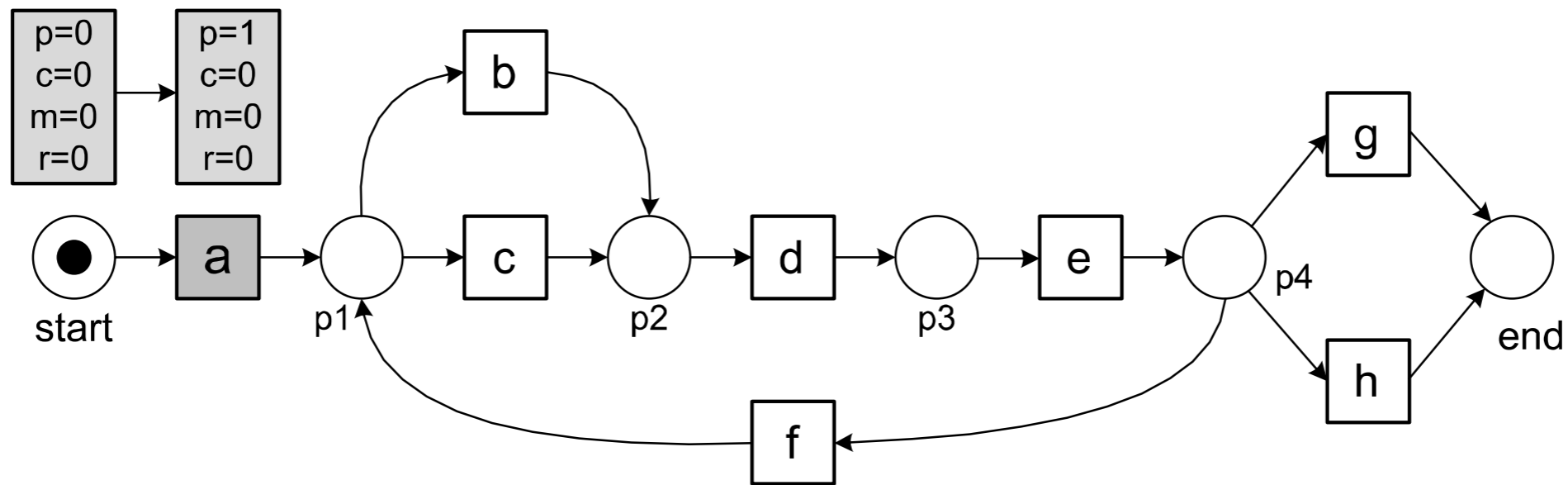


$$fitness(\sigma_1, N_1) = \frac{1}{2} \left(1 - \frac{0}{7}\right) + \frac{1}{2} \left(1 - \frac{0}{7}\right) = 1$$

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

Example: Missing Token

the environment produces a token for place start

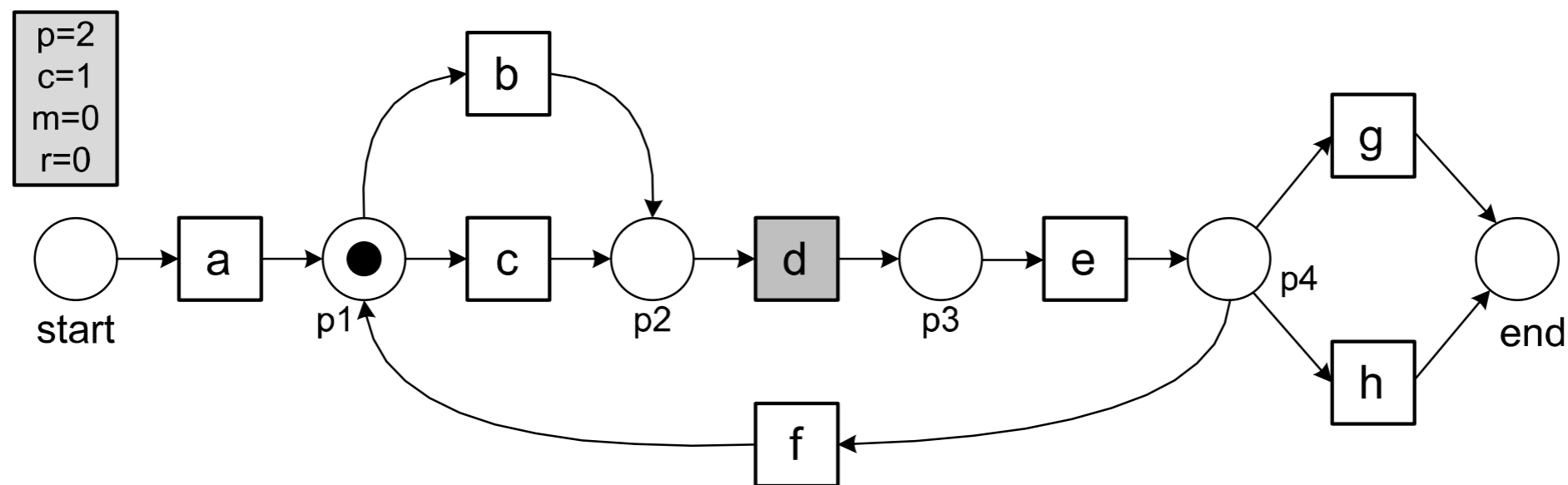


$$\sigma_3 = \langle a, d, c, e, h \rangle$$

Example: Missing Token

replaying a is possible

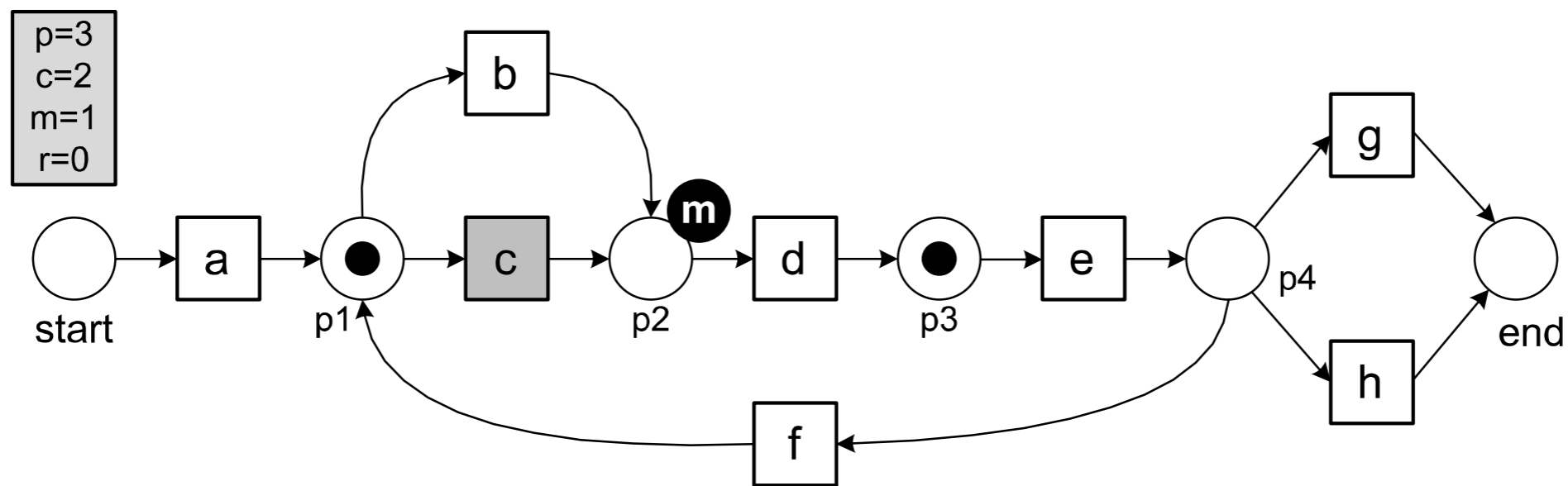
one token is consumed, one produced



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

Example: Missing Token

replaying d is NOT possible
one token is missing,
one produced, one consumed

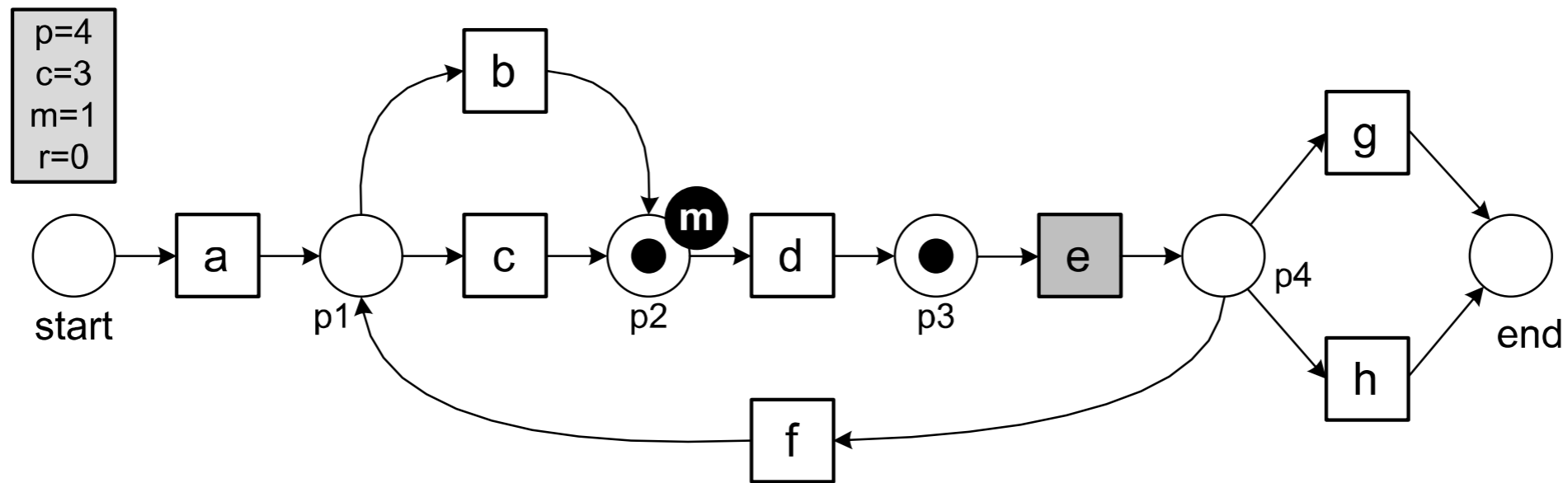


$$\sigma_3 = \langle a, d, c, e, h \rangle$$

Example: Missing Token

replaying c is possible

one token is produced, one consumed

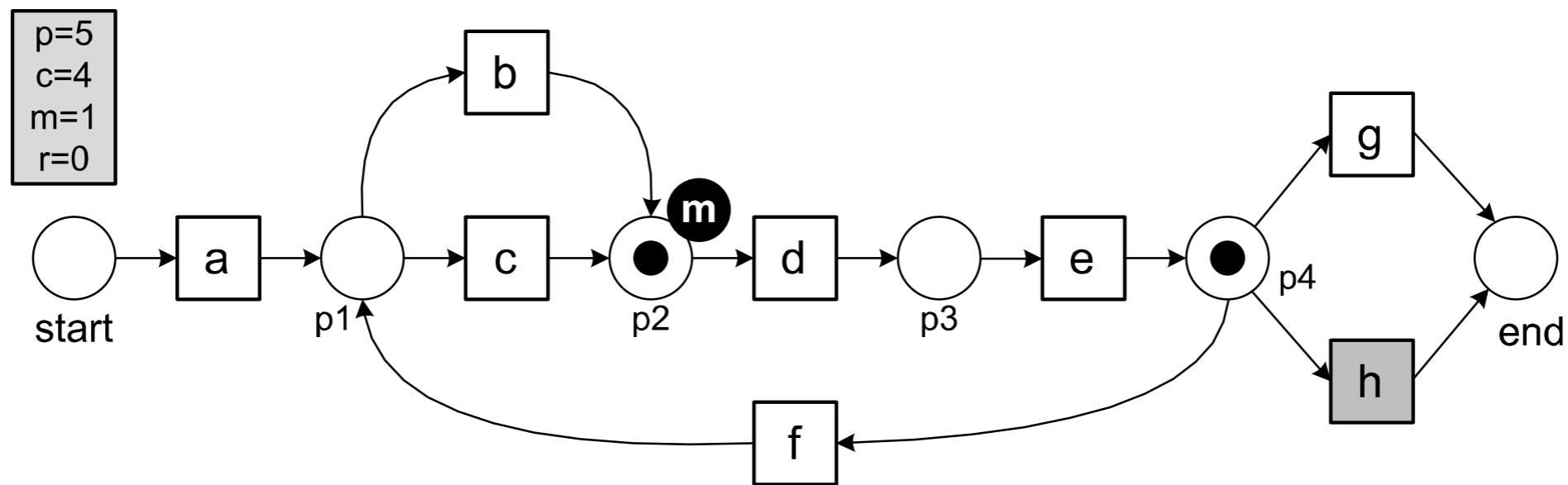


$$\sigma_3 = \langle a, d, c, e, h \rangle$$

Example: Missing Token

replaying e is possible

one token is produced, one consumed



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

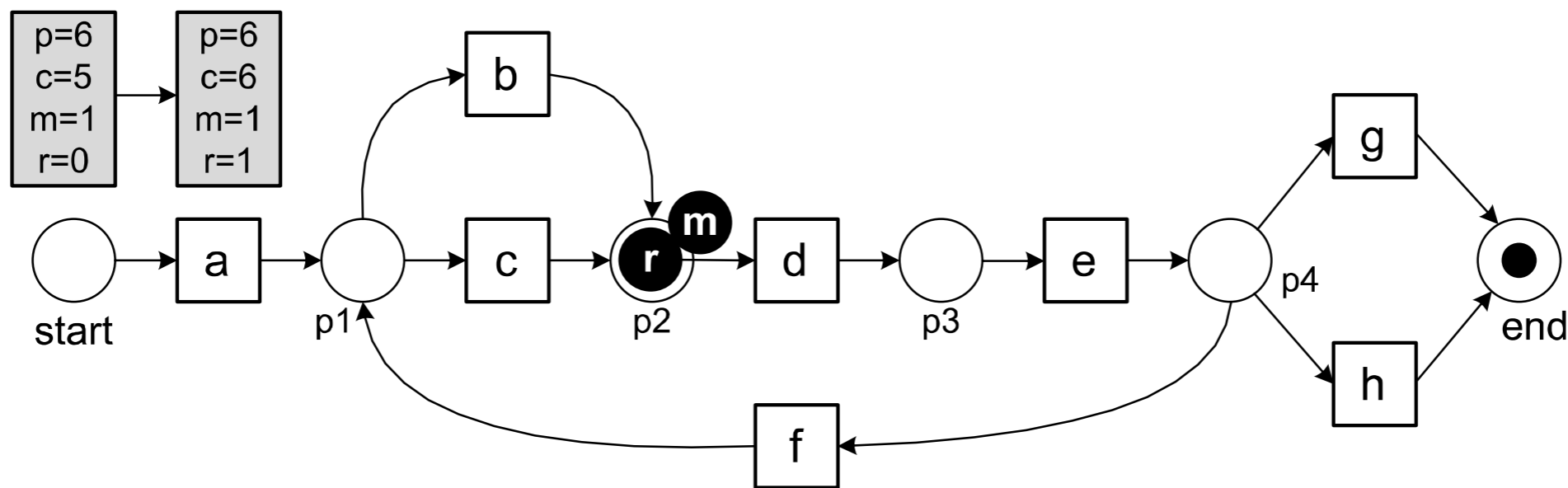
Example: Missing Token

replaying h is possible

one token is produced, one consumed

At the end,

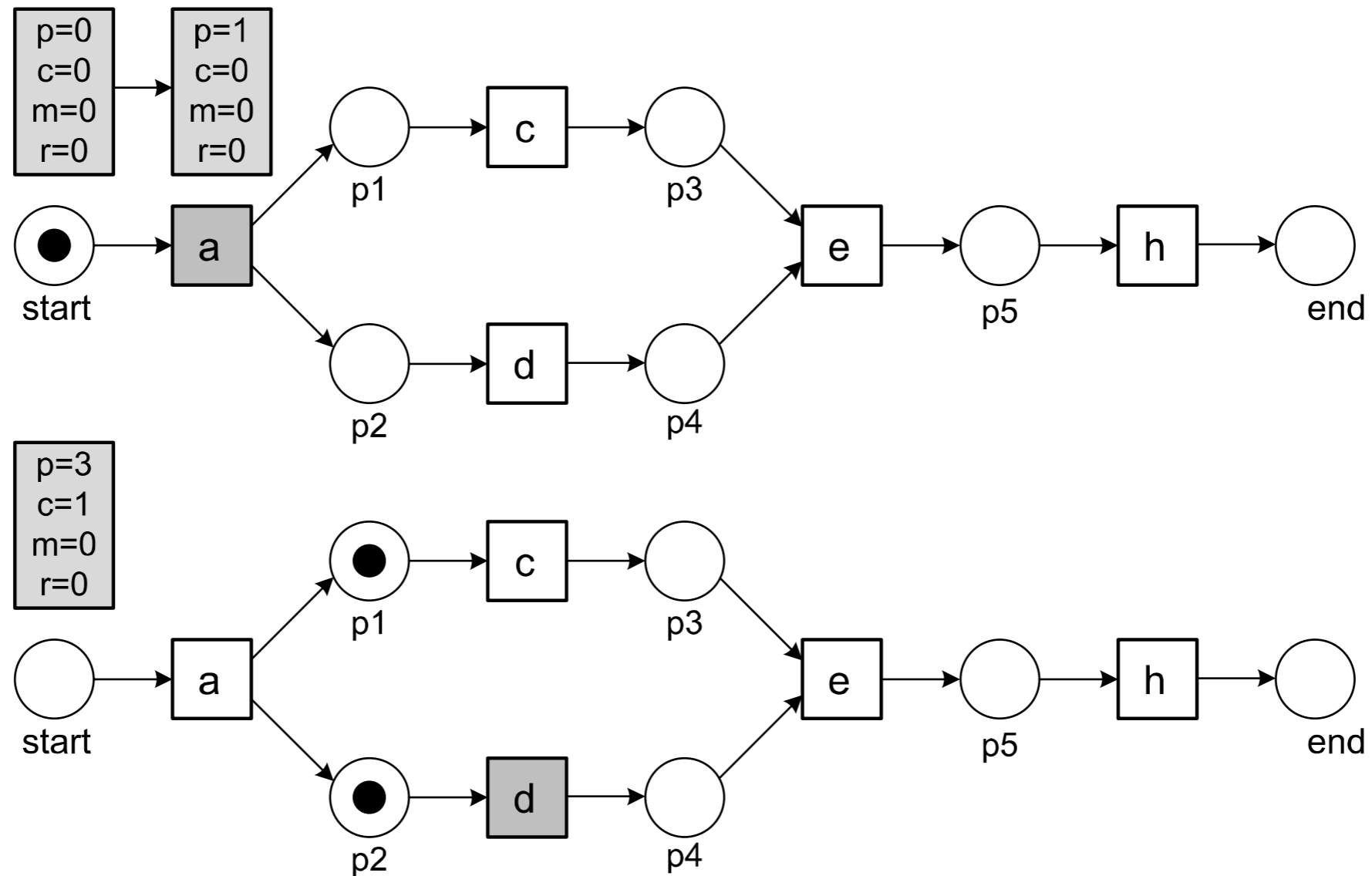
the environment consumes a token from place end.



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

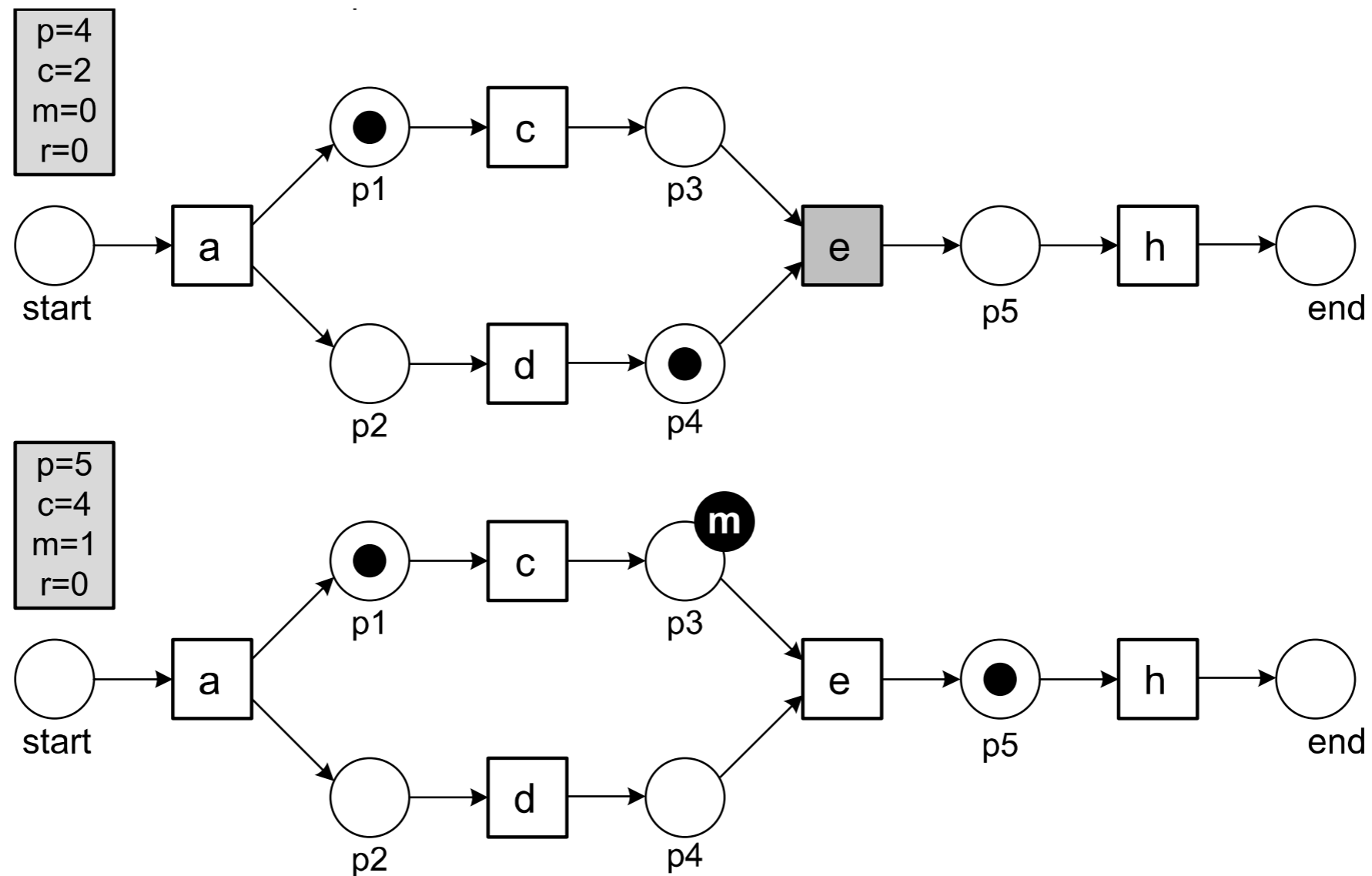
$$\sigma_3 = \langle a, d, c, e, h \rangle$$

Example: Event Removal



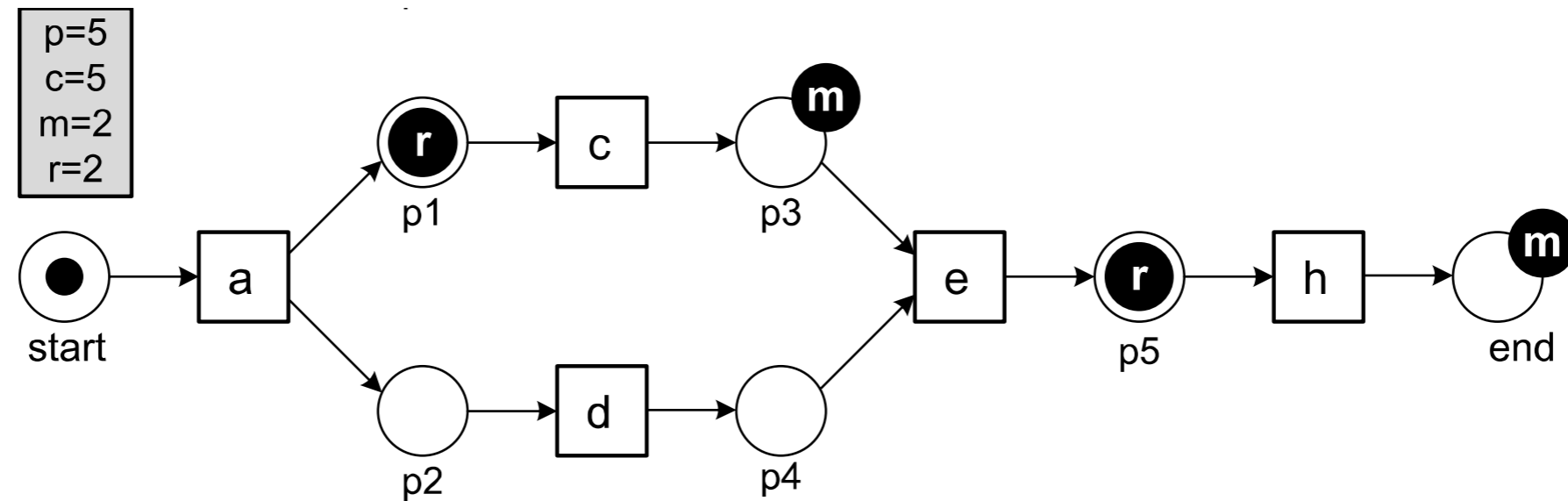
$$\sigma_2 = \langle a, b, d, e, g \rangle \quad \sigma'_2 = \langle a, d, e \rangle$$

Example: Event Removal



$$\sigma_2 = \langle a, b, d, e, g \rangle \quad \sigma'_2 = \langle a, d, e \rangle$$

Example: Event Removal



$$fitness(\sigma_2, N_3) = \frac{1}{2} \left(1 - \frac{2}{5} \right) + \frac{1}{2} \left(1 - \frac{2}{5} \right) = 0.6$$

$$\sigma_2 = \langle a, b, d, e, g \rangle \quad \sigma'_2 = \langle a, d, e \rangle$$

Fitness of a Log

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

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

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

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

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

Diagnostic Information

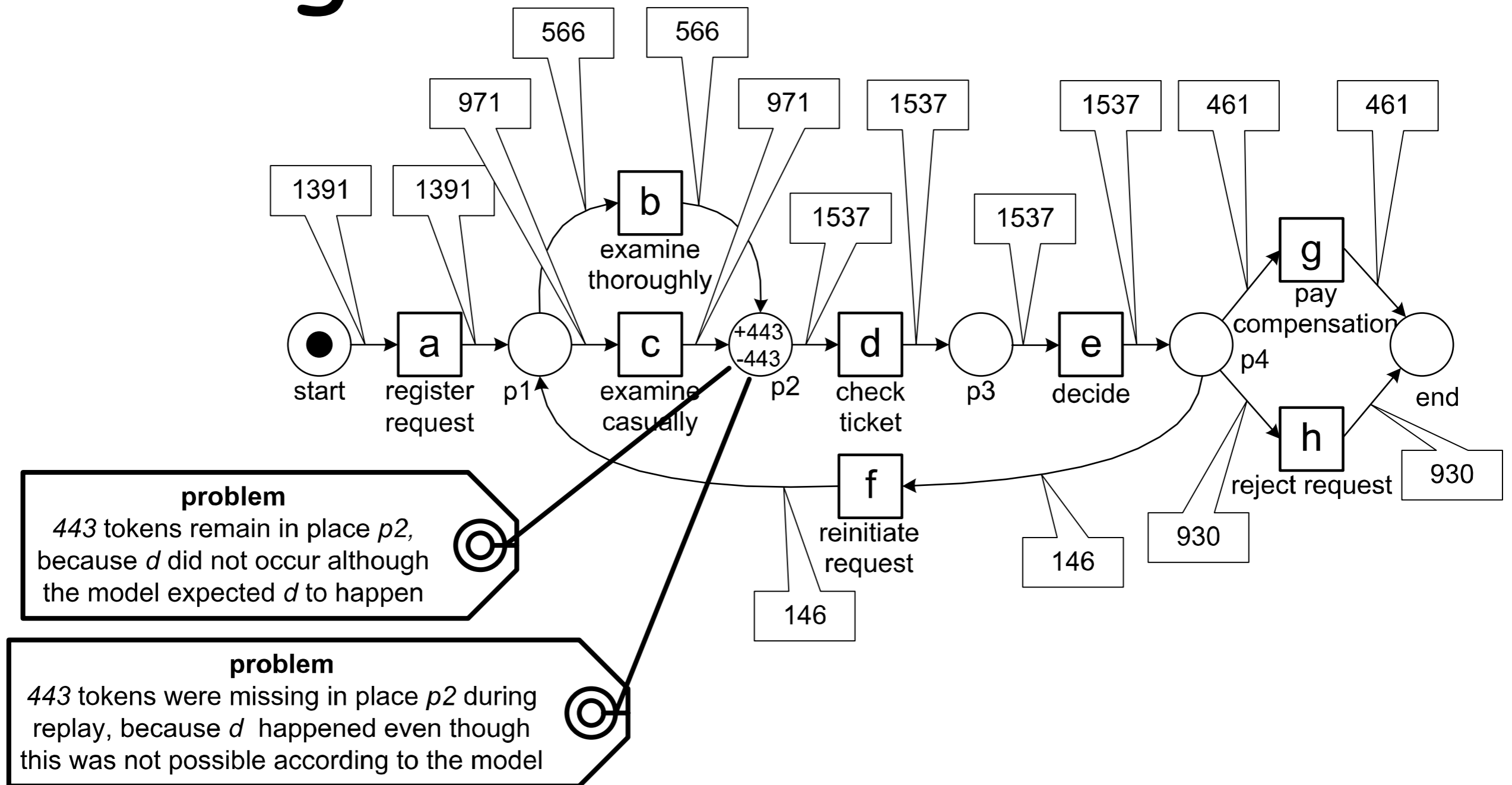


Fig. 7.6 Diagnostic information showing the deviations ($fitness(L_{full}, N_2) = 0.9504$)

Diagnostic Information

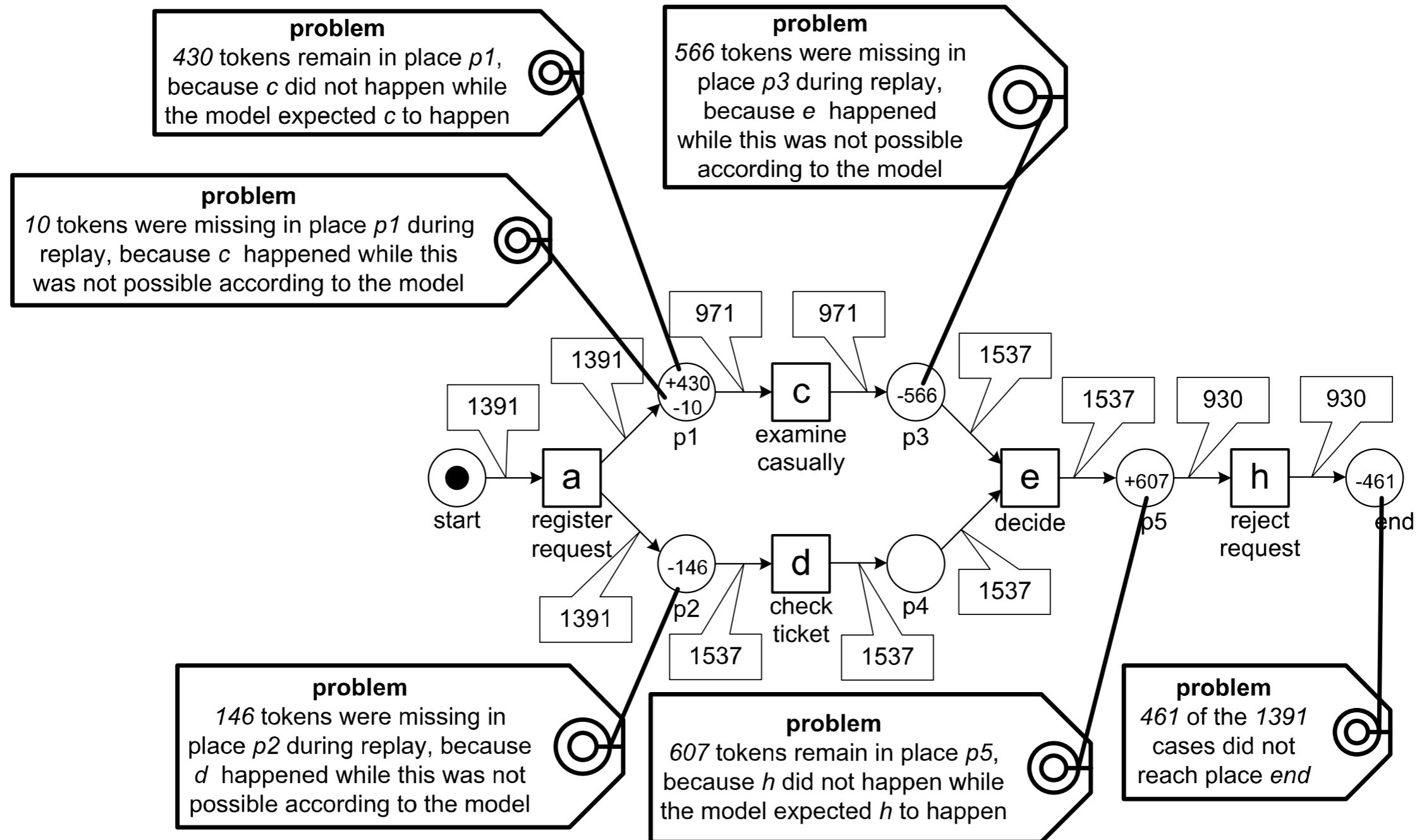


Fig. 7.7 Diagnostic information showing the deviations ($fitness(L_{full}, N_3) = 0.8797$)

Drill Down

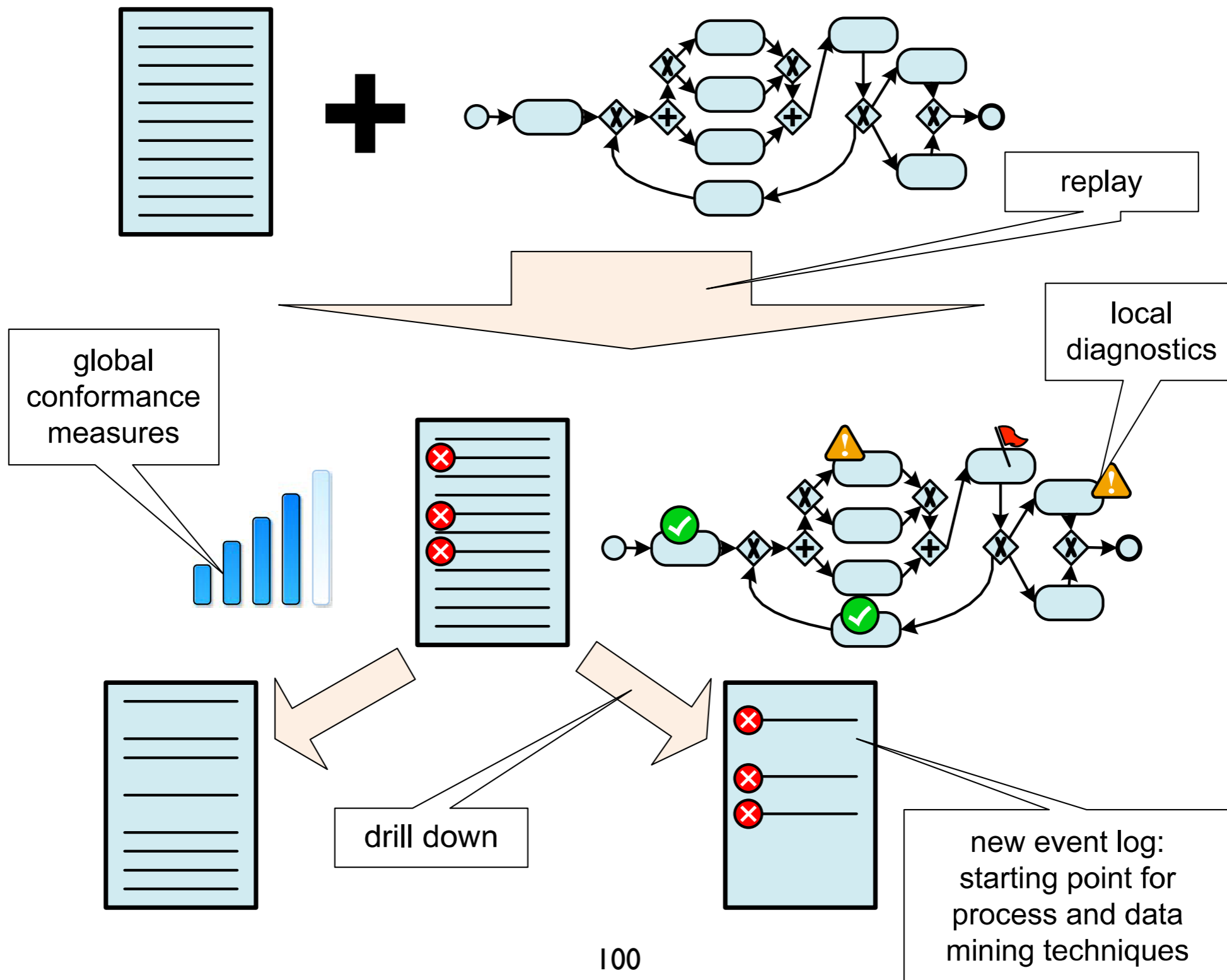
An event log can be split into two sublogs:
one event log containing only fitting cases and
one event log containing only non-fitting cases.

The second event log can be used to discover a different
process model.

Also other data and process mining techniques can be used.
For instance, it is interesting to know which people handled
the deviating cases and whether these cases took
longer or were more costly.

In case fraud is suspected, one may create a social
network based on the event log with deviating cases.

Drill Down



Comparing Footprints

Footprint from Play-out

Given a workflow net, the play-out technique can be used to extract a local complete set of traces.

If we see the set of traces as an event log (without multiplicities), then we can derive the relation $>$.

Then, we can construct the footprint (i.e. a matrix showing causal dependencies between events) of the net model based on such relation $>$.

(From the viewpoint of a footprint matrix, an event log is complete if and only if all activities that can follow one another do so at least once in the log.)

Footprint-based Conformance

Footprints are available for logs and models (nets).

This allows for:

log vs model conformance

(do the log and the model agree on the ordering of activities?)

model vs model conformance

(quantification of their similarities)

log vs log comparison

(*concept drift*: how does the work changes in sub-logs?)

Conformance based on footprints

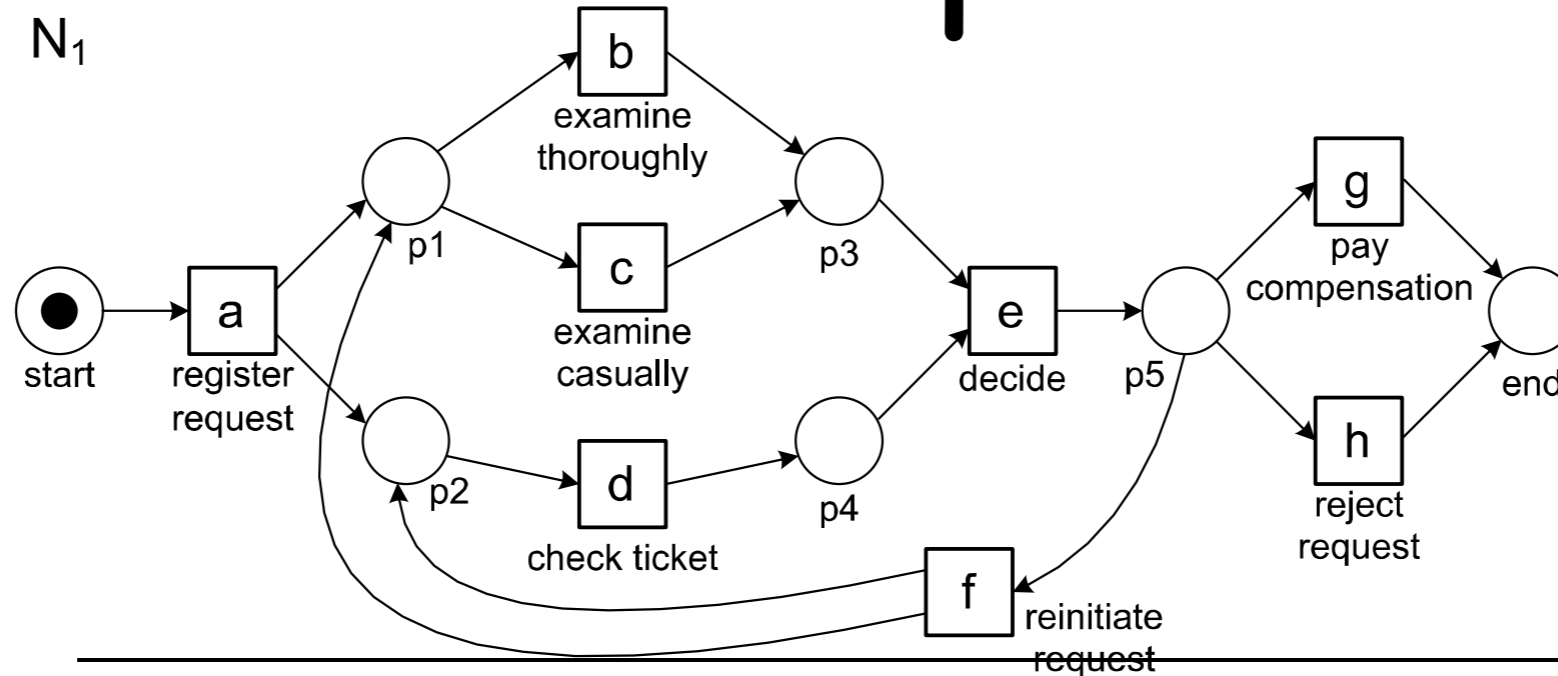
The conformance based on footprints can be computed by taking:

n : total number of cells in the footprint matrix

d : number of cells with different content between the two matrices

$$1 - \frac{d}{n}$$

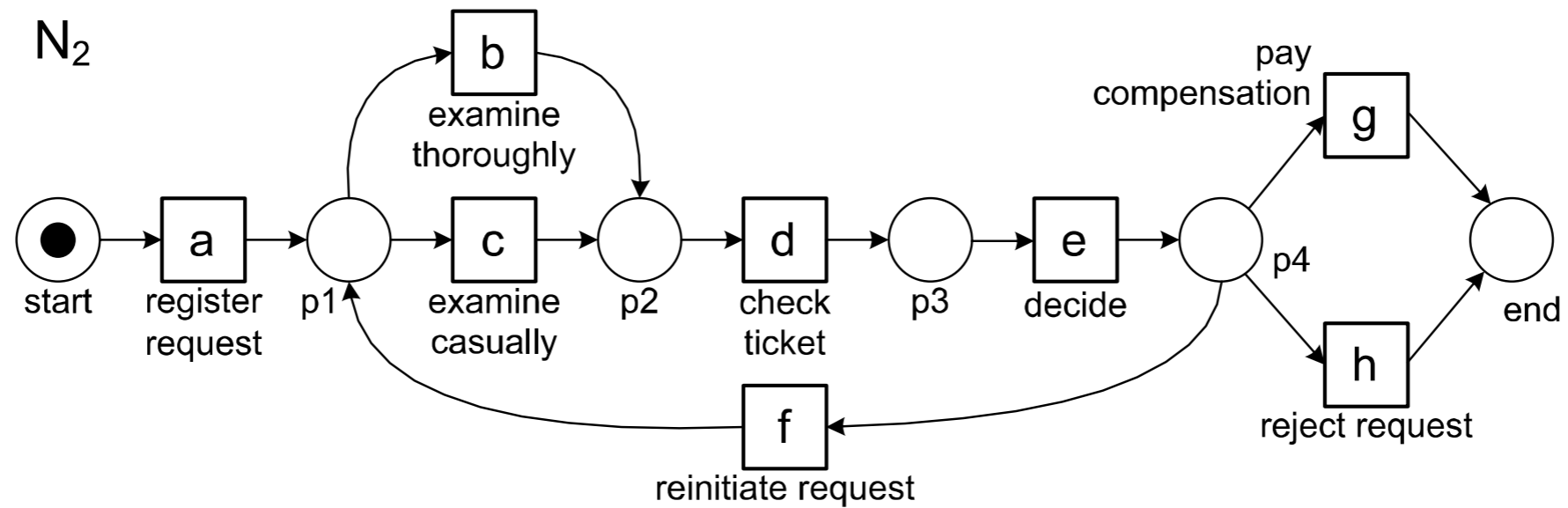
Example



	<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>	#	#	#	#	←	#	#	#

Also
Footprint of L_{full}

Example



	<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>	#	#	#	#	←	#	#	#

Example

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

Example

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

<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>							

Conclusion

Requirements gone bad



How the customer explained it

Conclusion

**We have overviewed the iceberg tip of
business process management**

**more notation, theory, technology,
tools, methodology, encoding,
validation, verification, research
lie down there,
more or less deep,
below the surface...**

...for all of us to explore