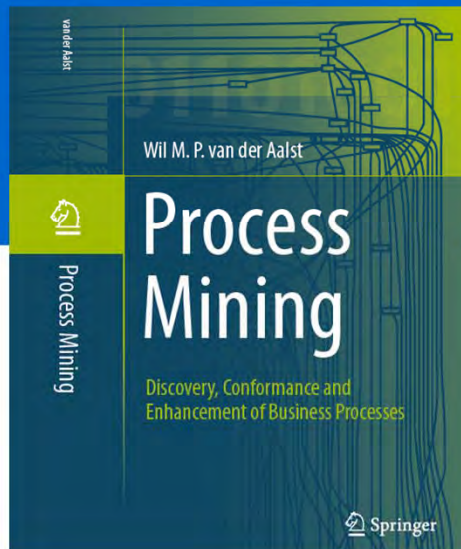


Distributed Process Discovery and Conformance Checking

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

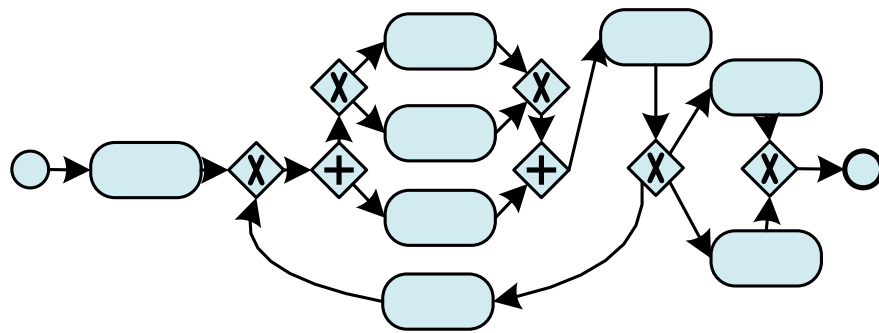


TU/e Technische Universiteit
Eindhoven
University of Technology

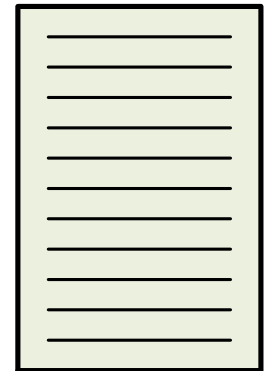
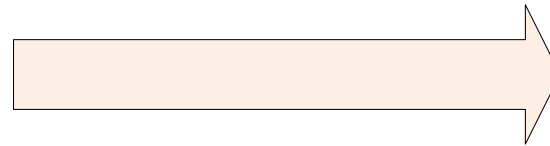
Where innovation starts

On the different roles of (process) models ...

Play-Out

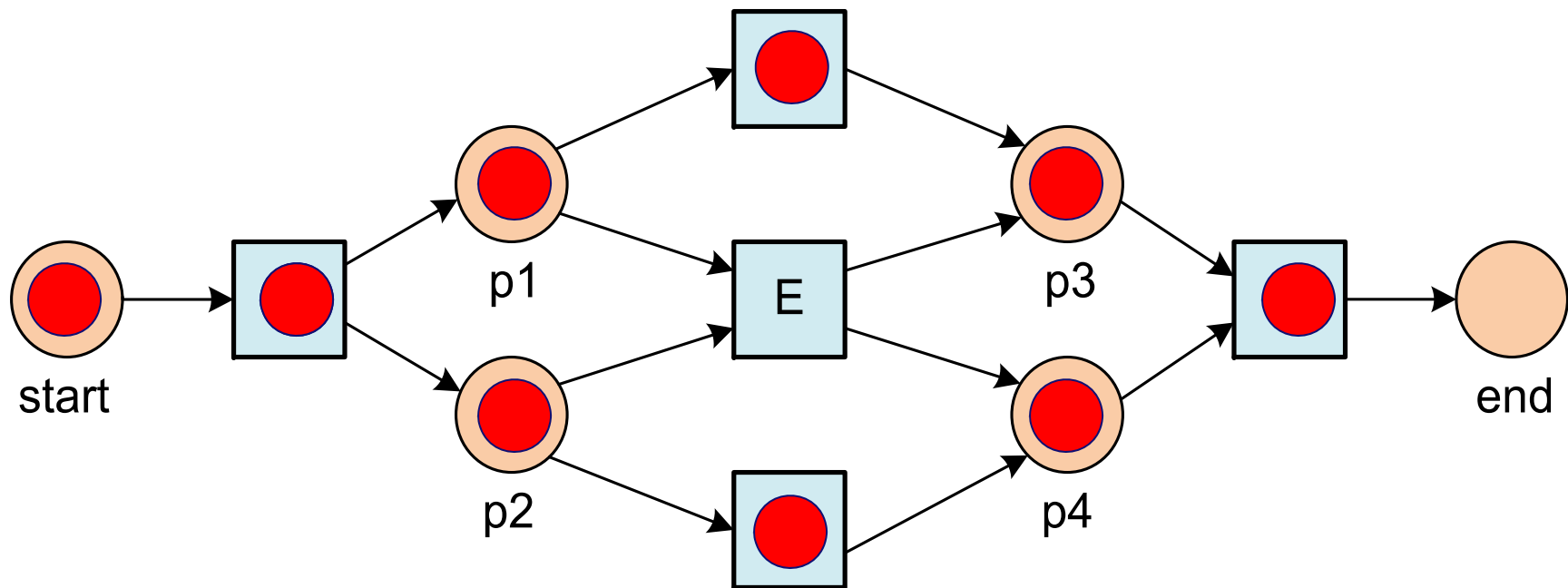


process model



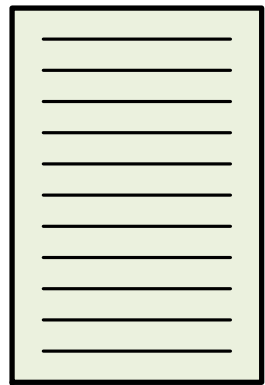
event log

Play-Out (Classical use of models)

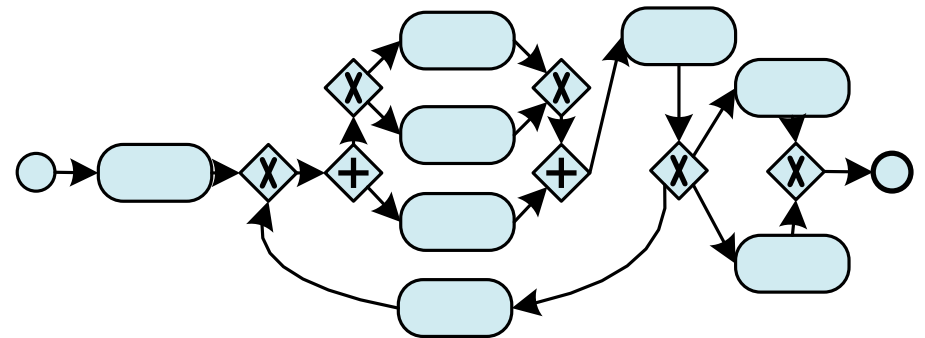


A B C D A E D A E D
A C B D A B C D A C B D
A C B D A E D A C B D

Play-In



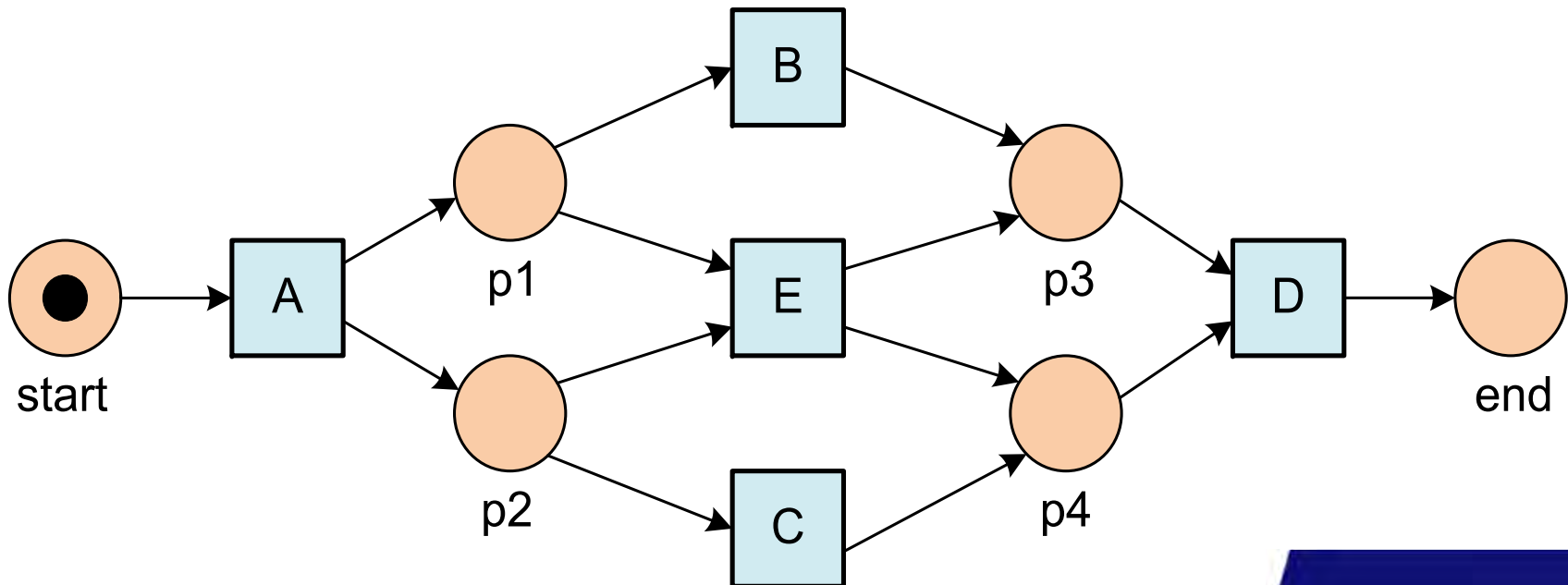
event log



process model

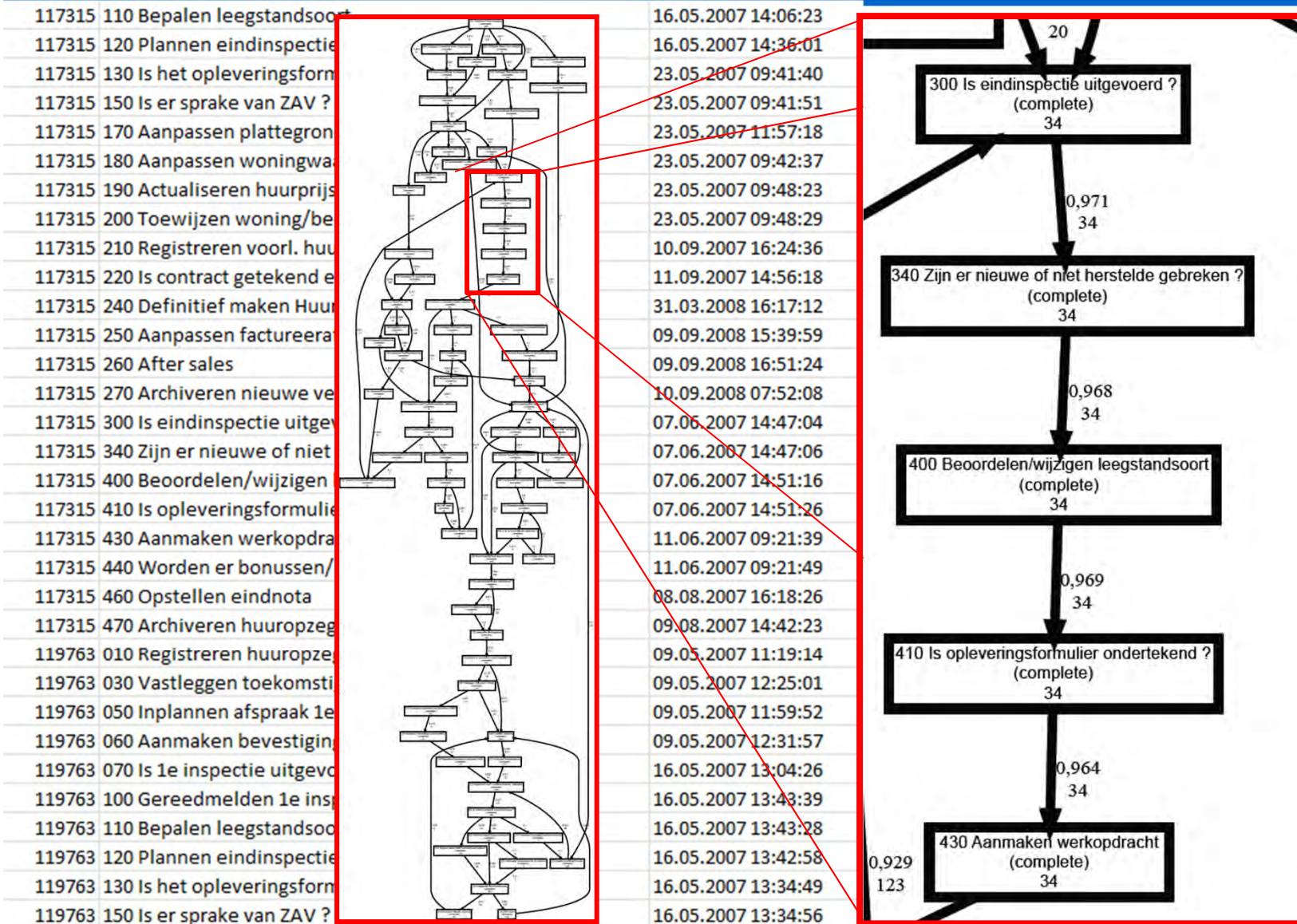
Play-In

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



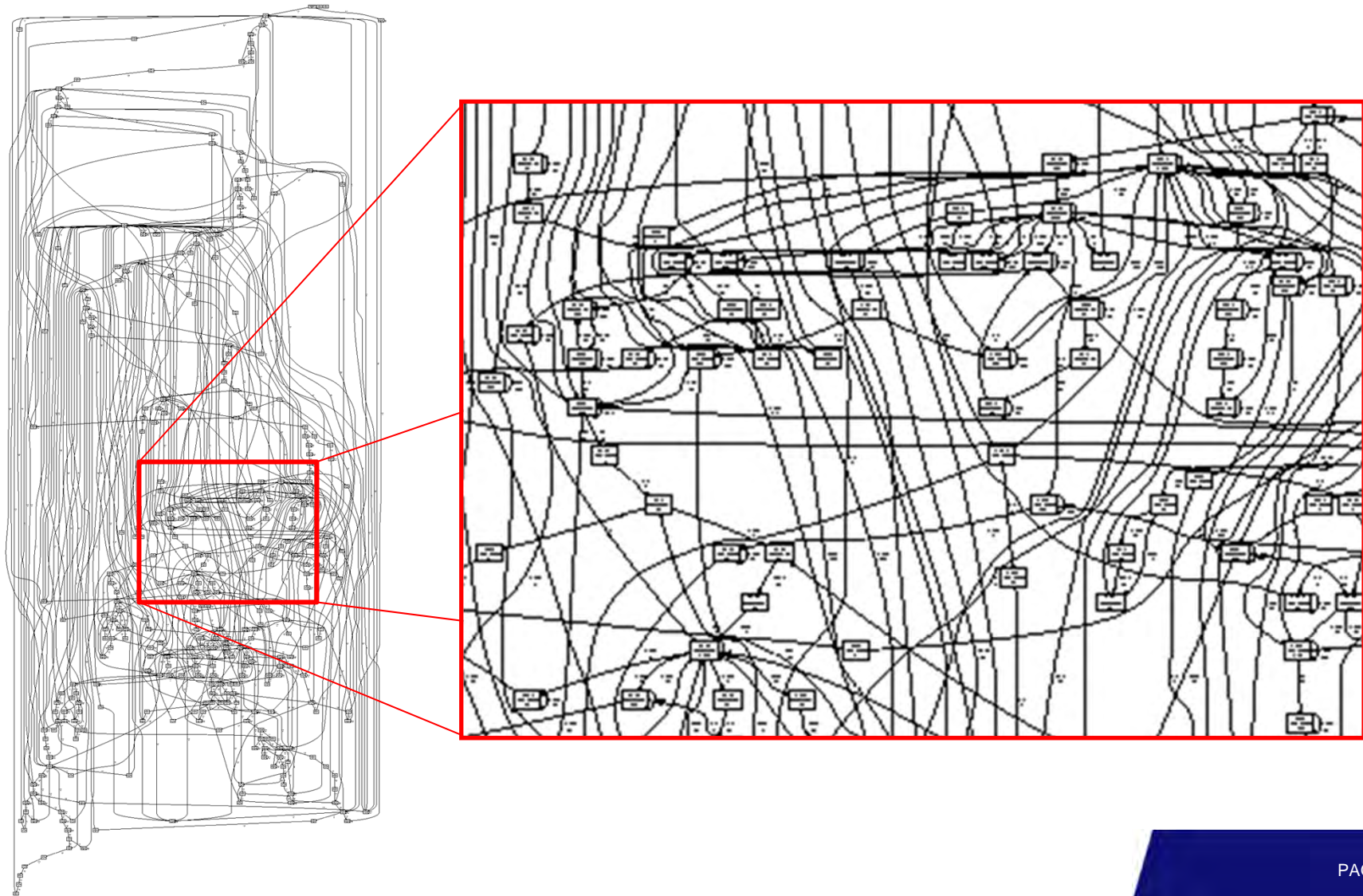
Example Process Discovery

(Vestia, Dutch housing agency, 208 cases, 5987 events)



(ASML, test process lithography systems, 154966 events)

(ASML, test process lithography systems, 154966 events)

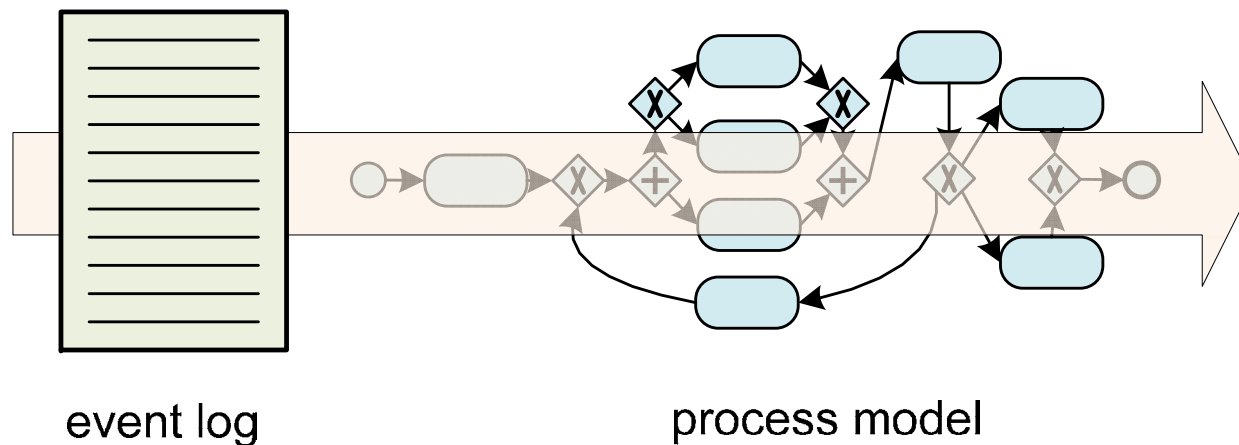


Example Process Discovery

(AMC, 627 gynecological oncology patients, 24331 events)



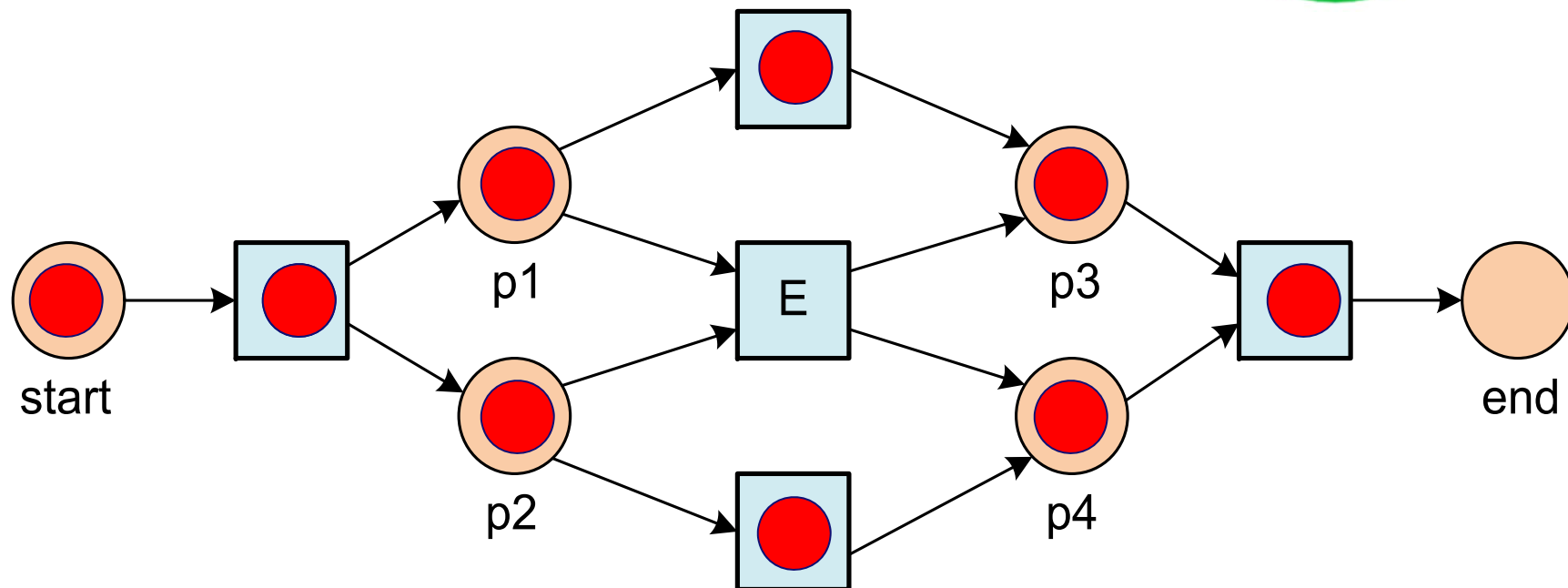
Replay



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

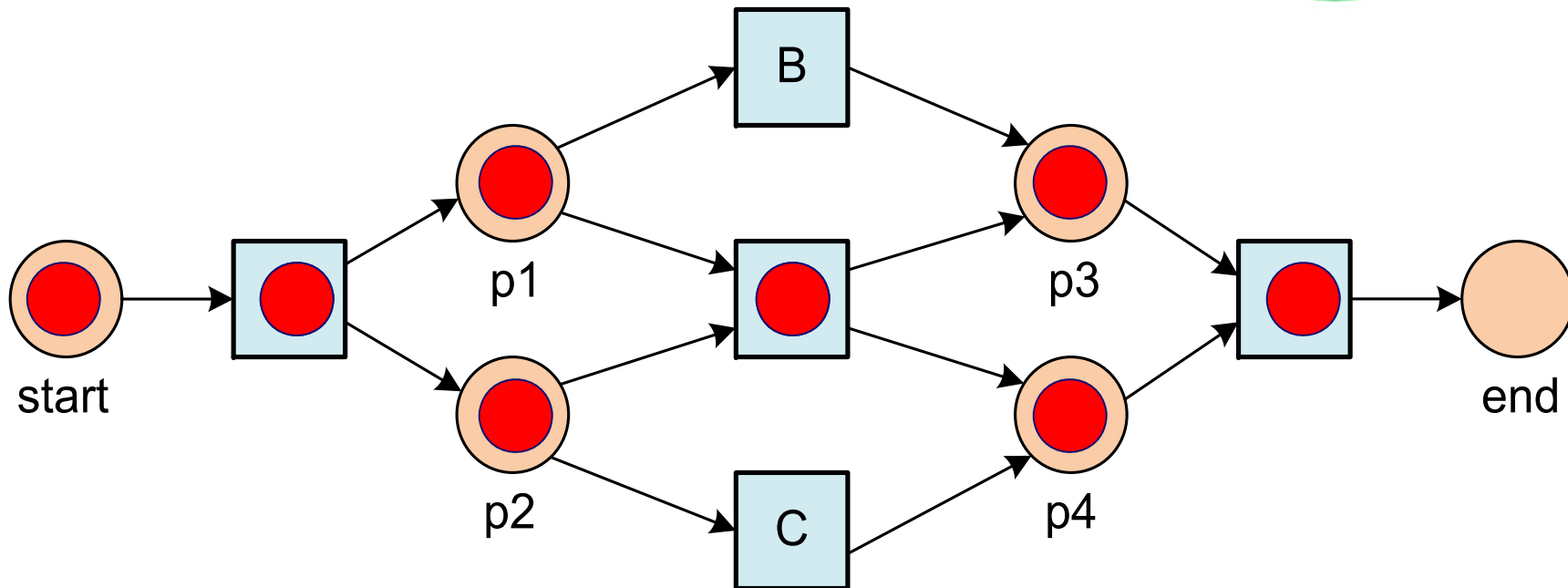
Replay

A B C D



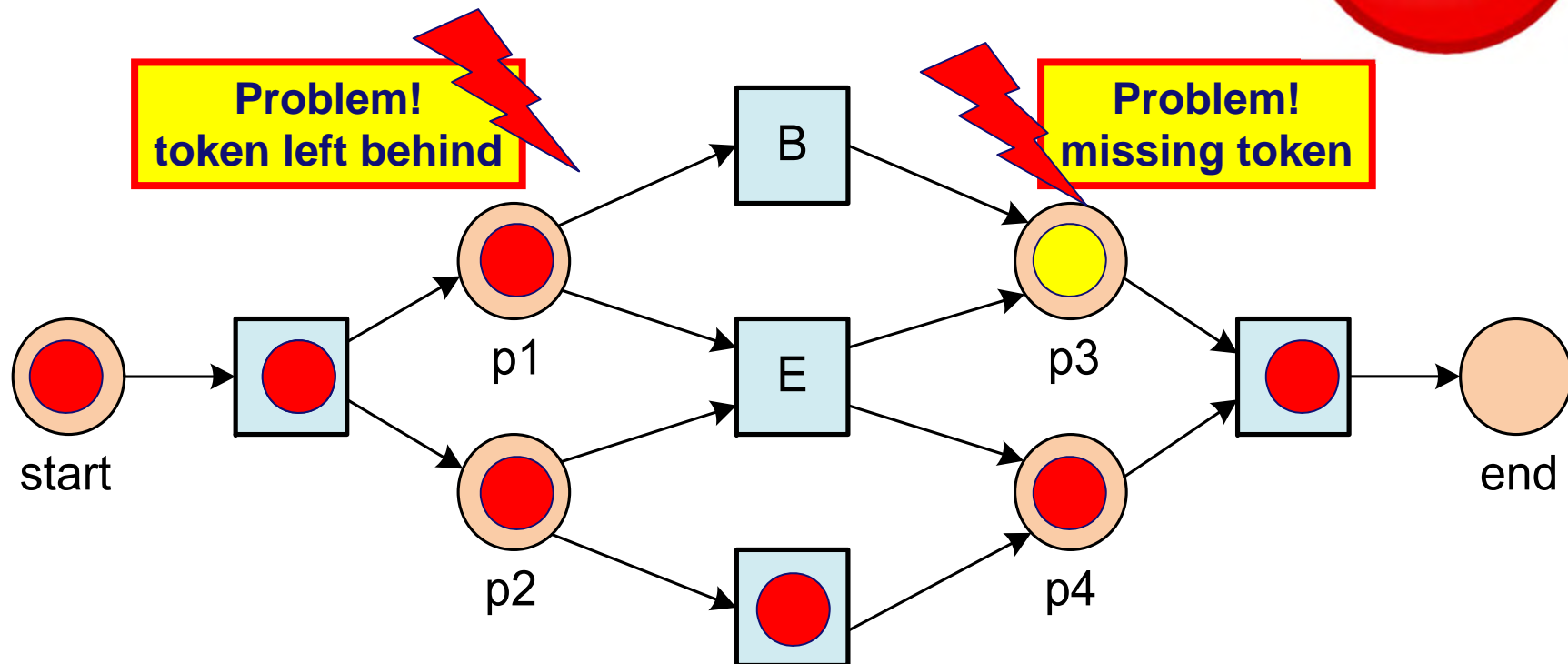
Replay

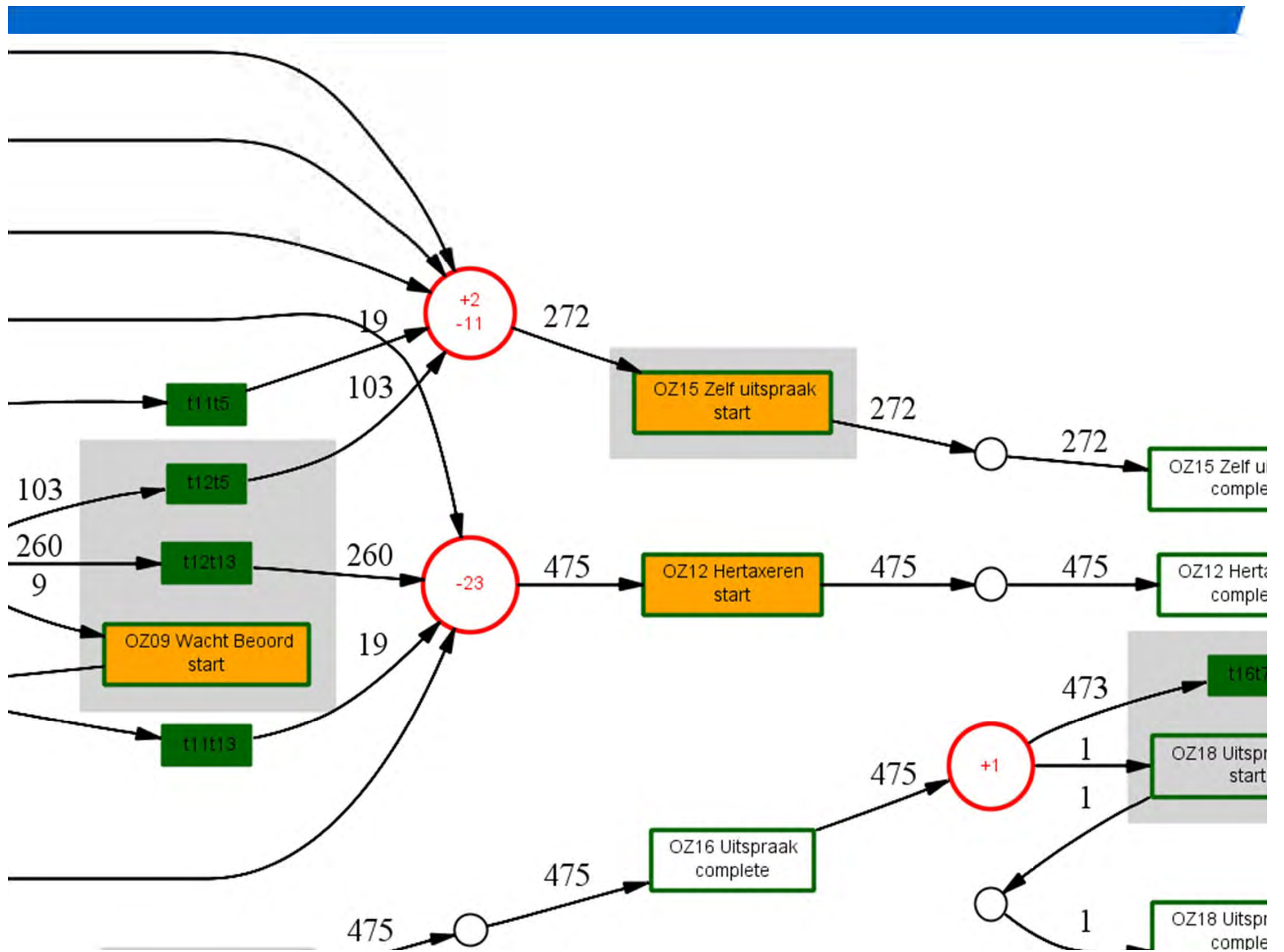
A E D



Replay can detect problems

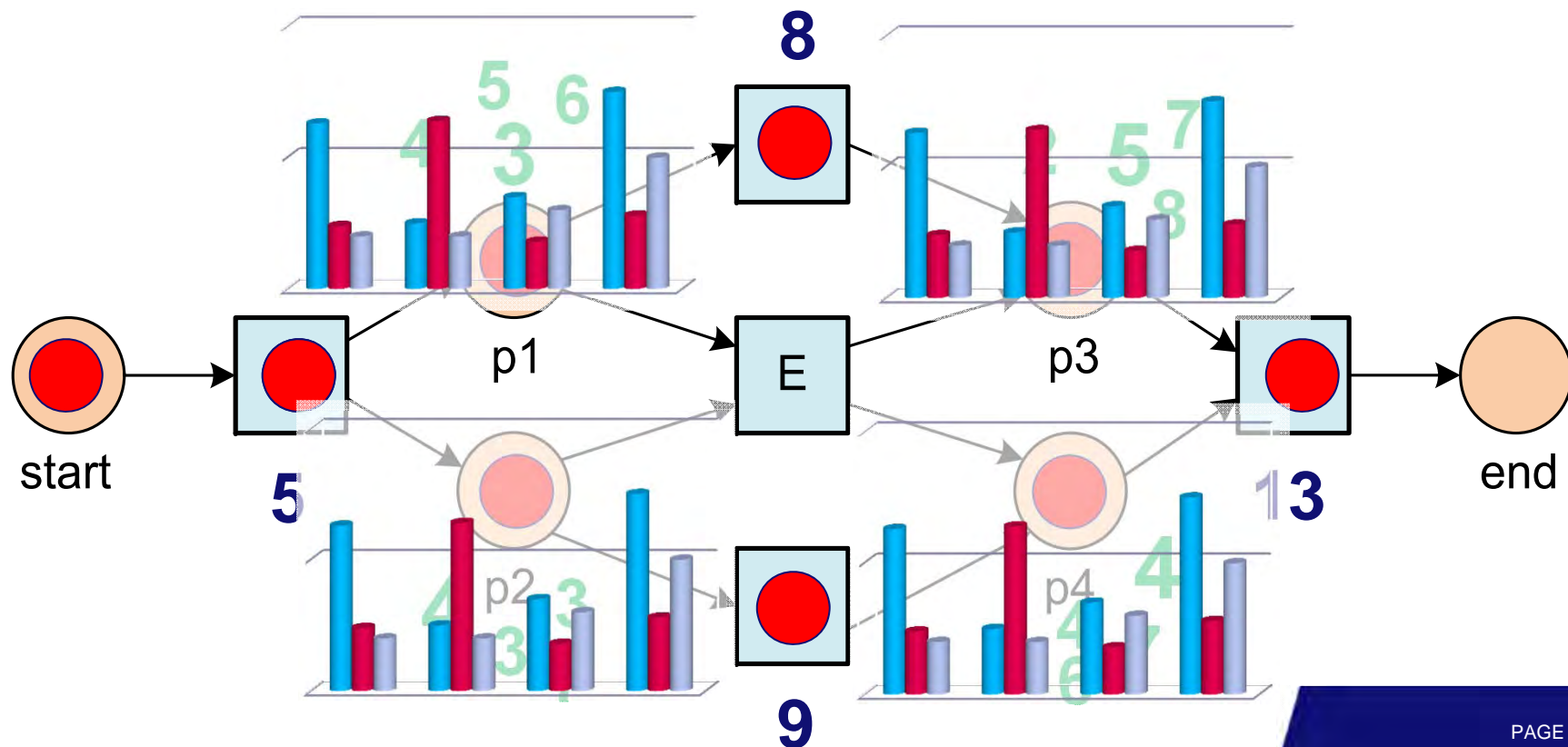
A C D





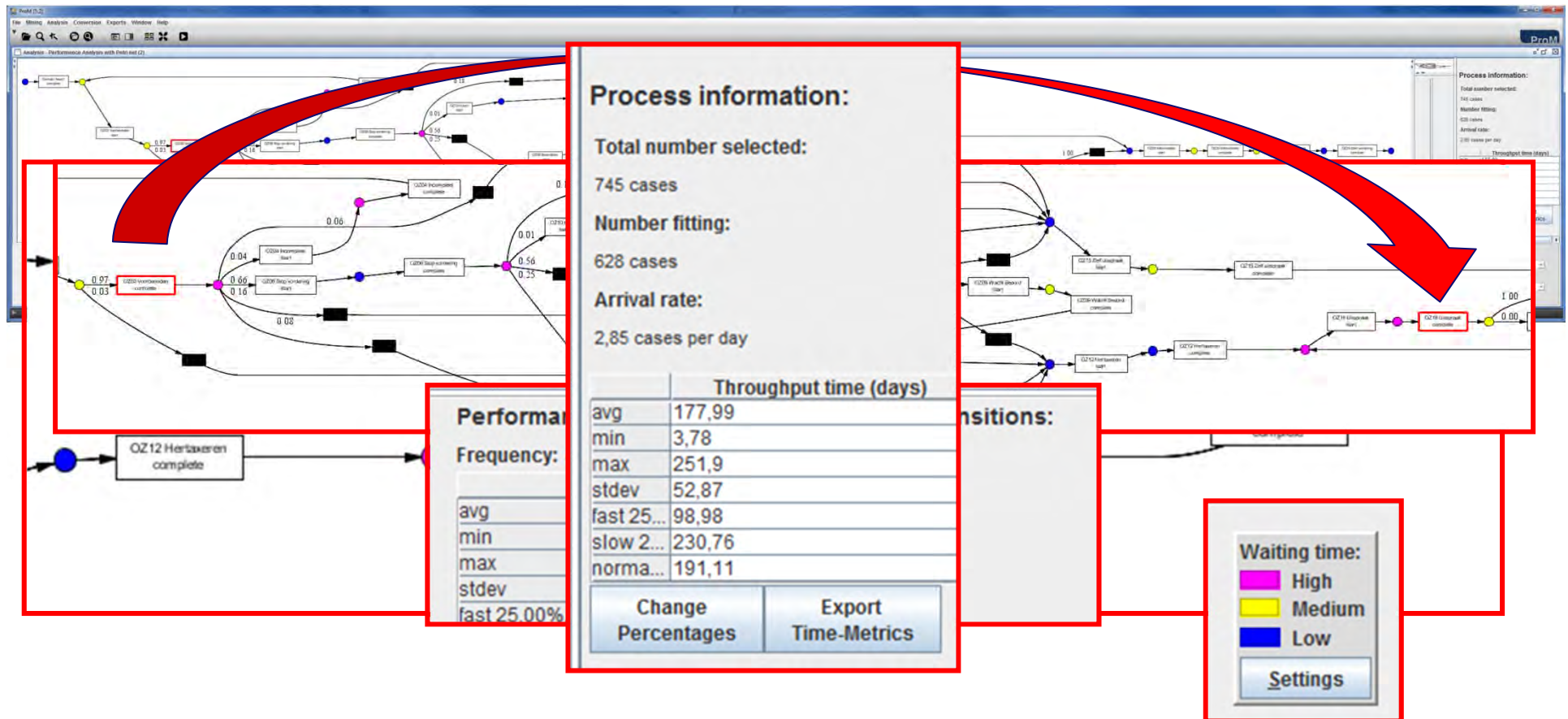
Replay can extract timing information

A⁵B⁸C⁹D¹³



Performance Analysis Using Replay

(WOZ objections Dutch municipality, 745 objections, 9583 event, $f = 0.988$)



Big Data

Big Data

“Enterprises globally stored more than 7 exabytes of new data on disk drives in 2010, while consumers stored more than 6 exabytes of new data on devices such as PCs and notebooks.”

“All of the world's music can be stored on a \$600 disk drive.”

“Indeed, we are generating so much data today that it is physically impossible to store it all. Health care providers, for instance, discard 90 percent of the data that they generate.”

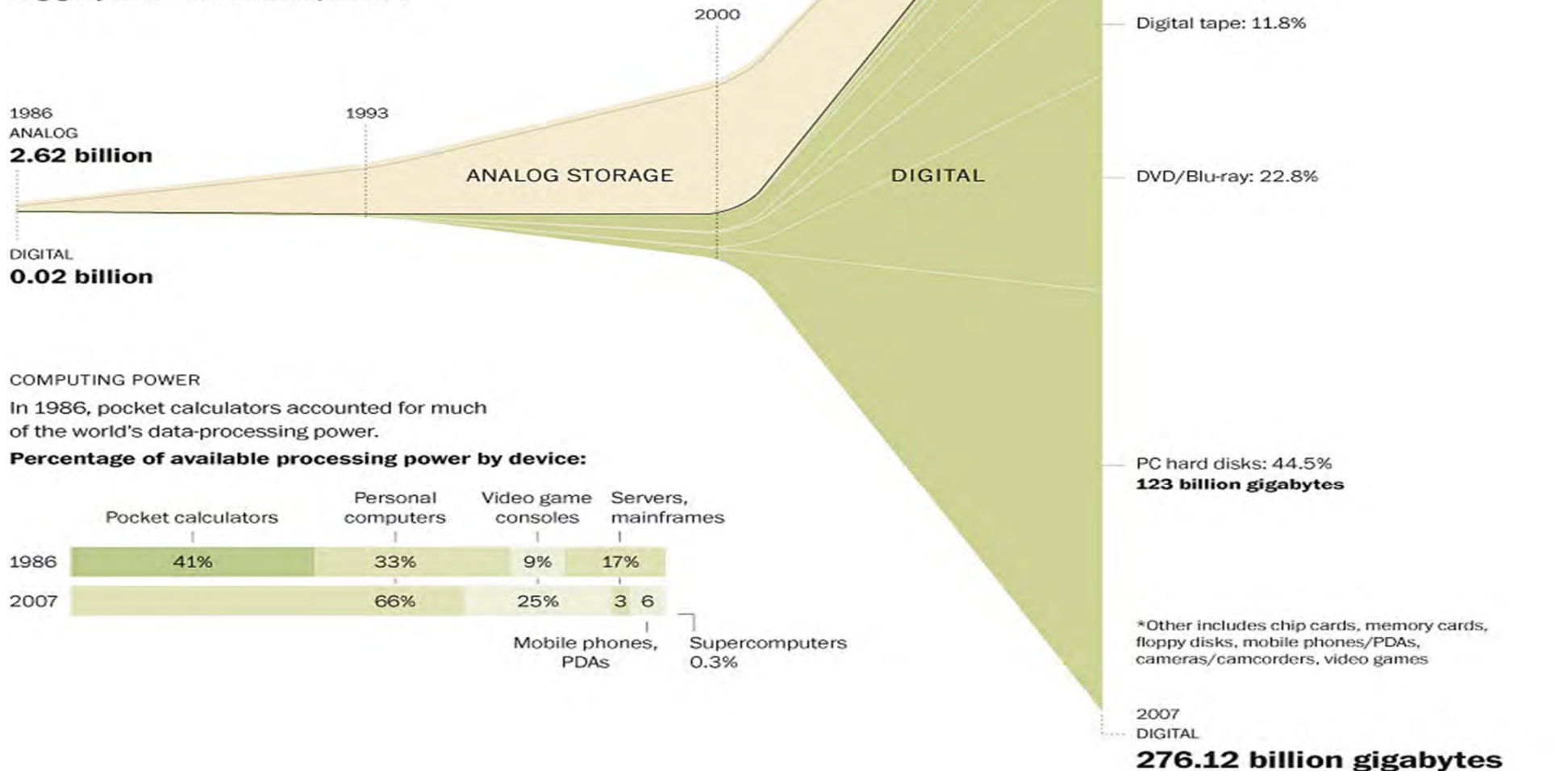
Source: “Big Data: The Next Frontier for Innovation, Competition, and Productivity” McKinsey Global Institute, 2011.

Hilbert and Lopez. The World's Technological Capacity to Store, Communicate, and Compute Information. Science, 332(6025):60-65, 2011.

THE WORLD'S CAPACITY TO STORE INFORMATION

This chart shows the world's growth in storage capacity for both analog data (books, newspapers, videotapes, etc.) and digital (CDs, DVDs, computer hard drives, smartphone drives, etc.)

In gigabytes or estimated equivalent

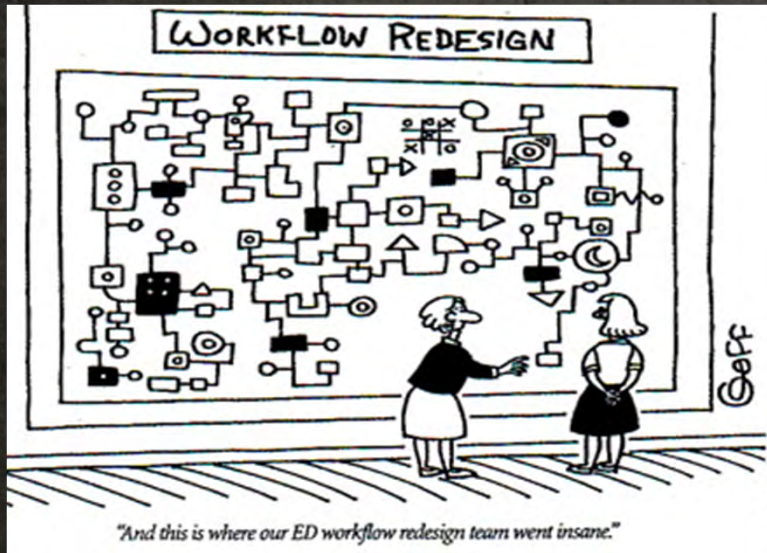








Evidence-Based Computer Science



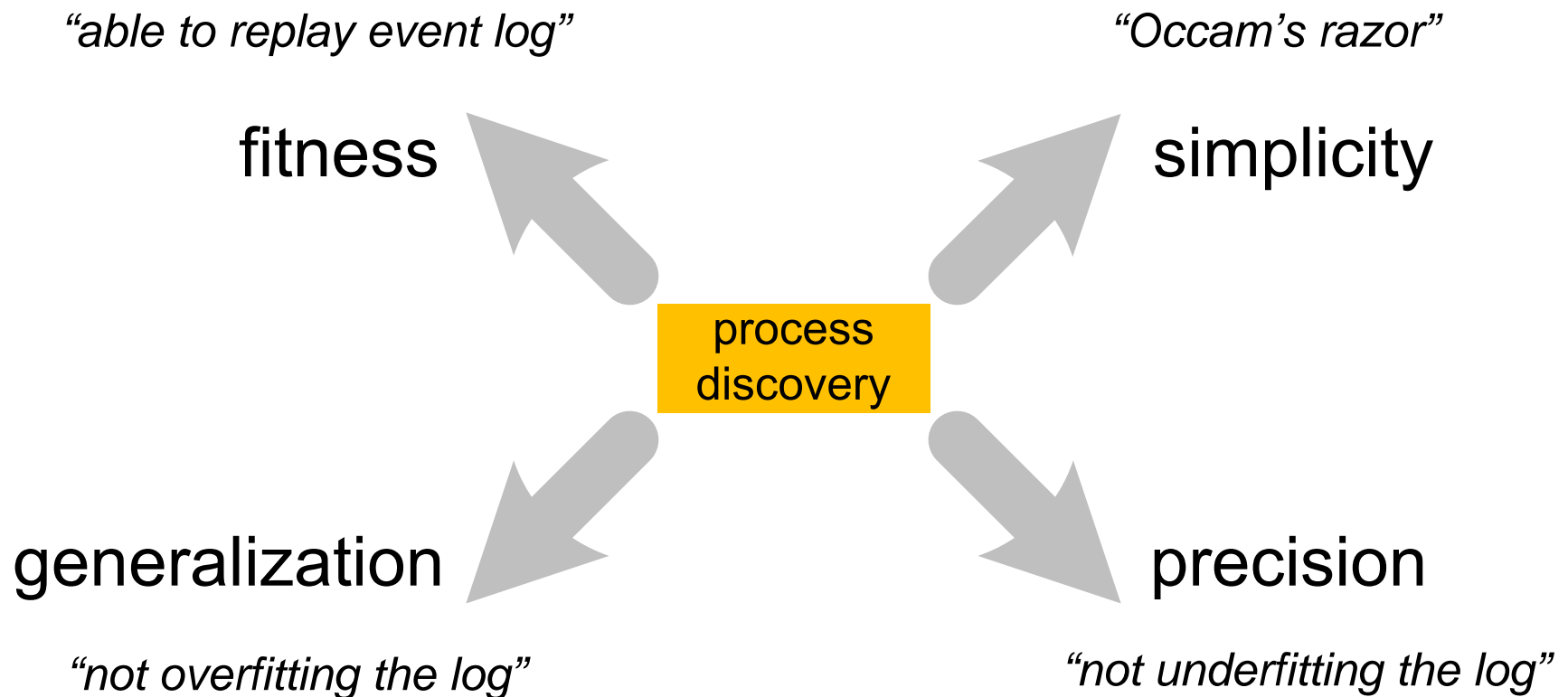
CRIME SCENE DO NOT CROSS

CRIME SCENE DO NOT CROSS

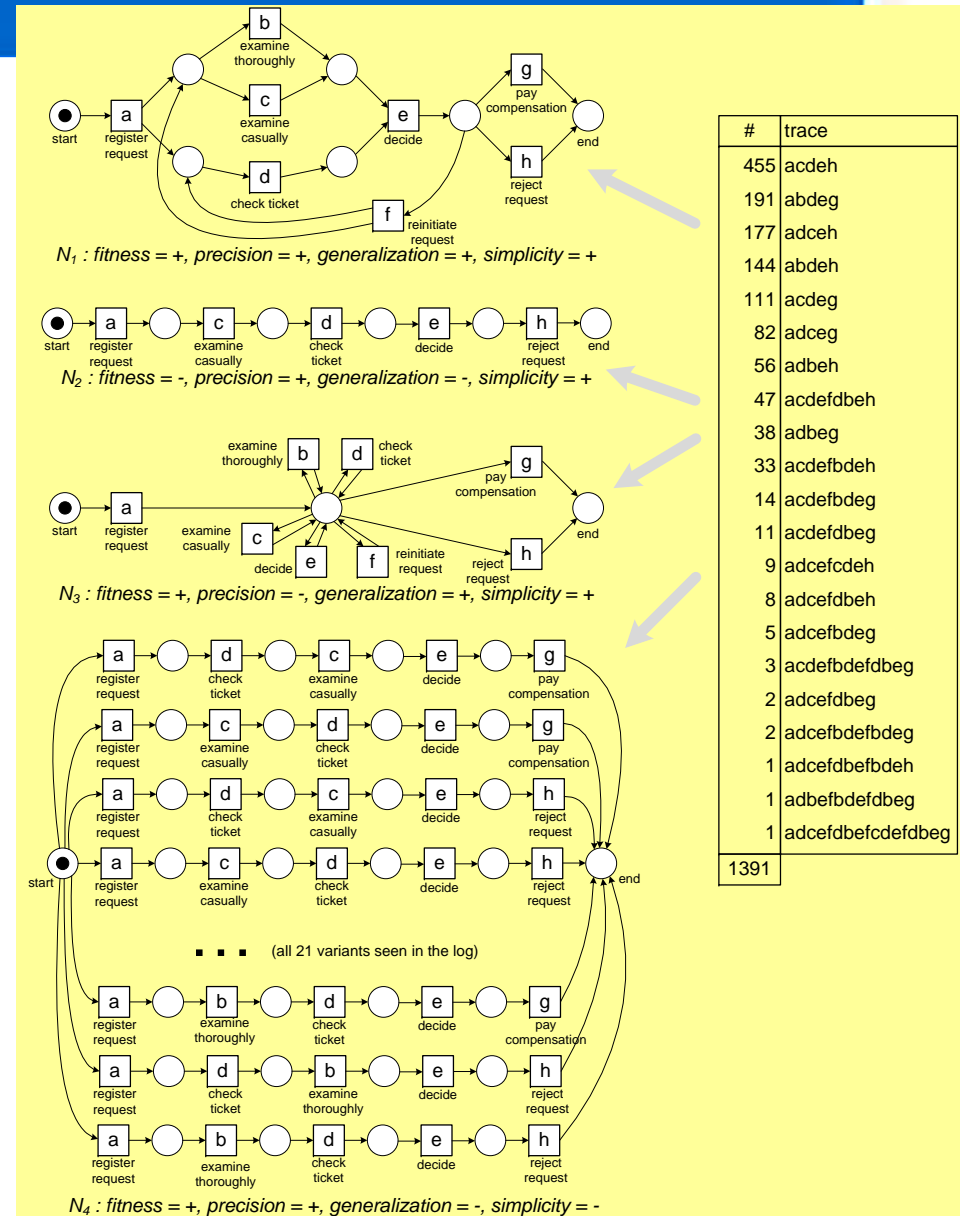
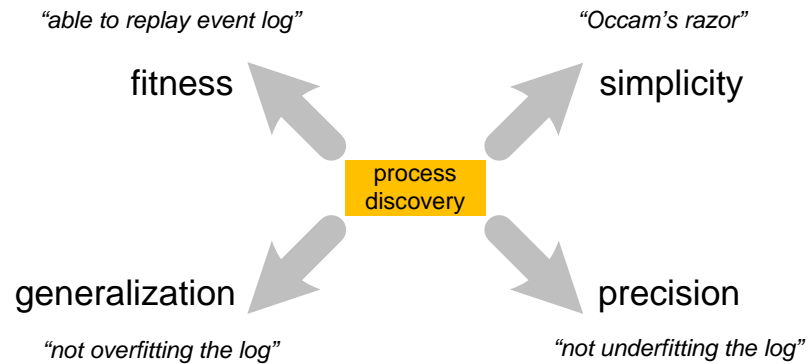
CRIME SCENE DO NOT CROSS

How good is my model?

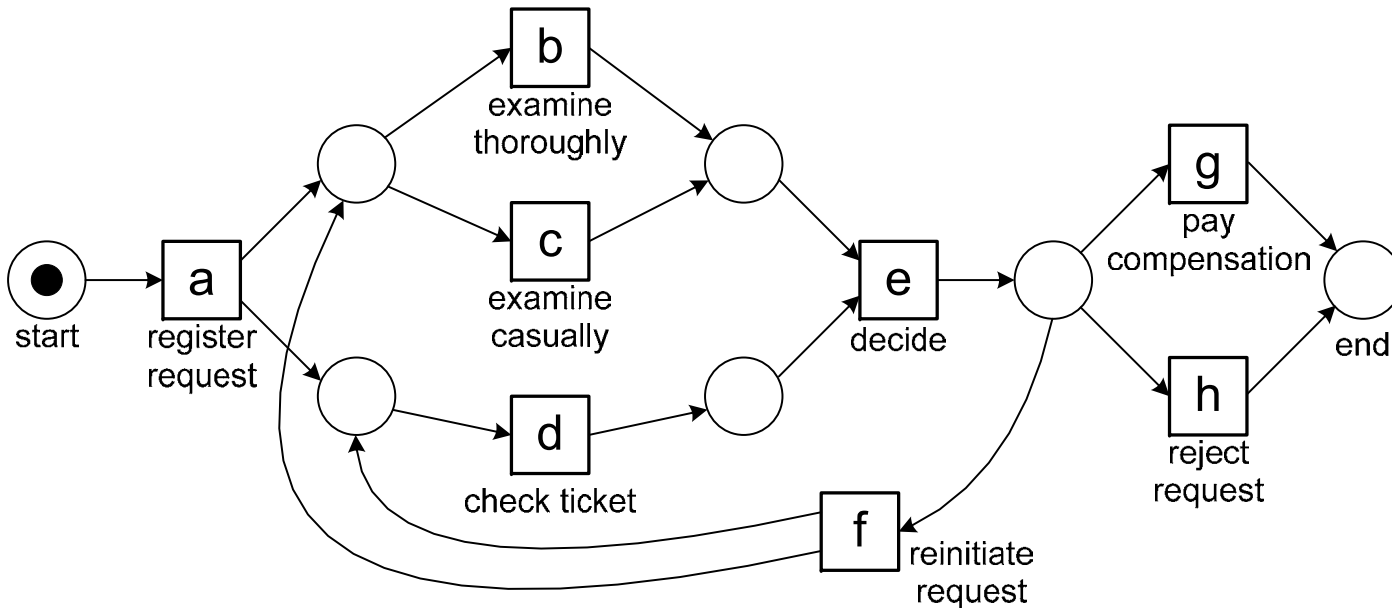
Four Competing Quality Criteria



Example: one log four models



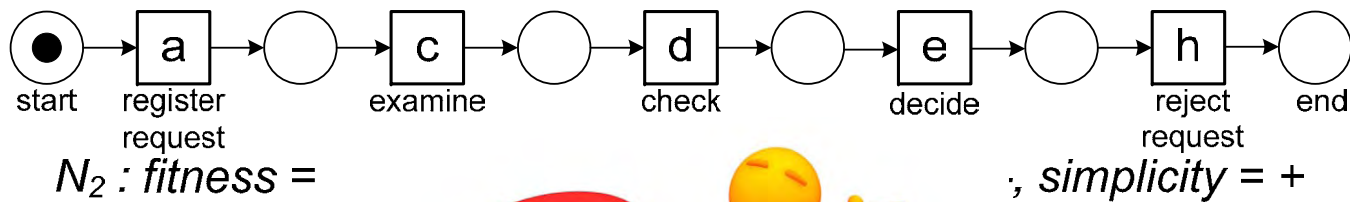
Model N₁



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

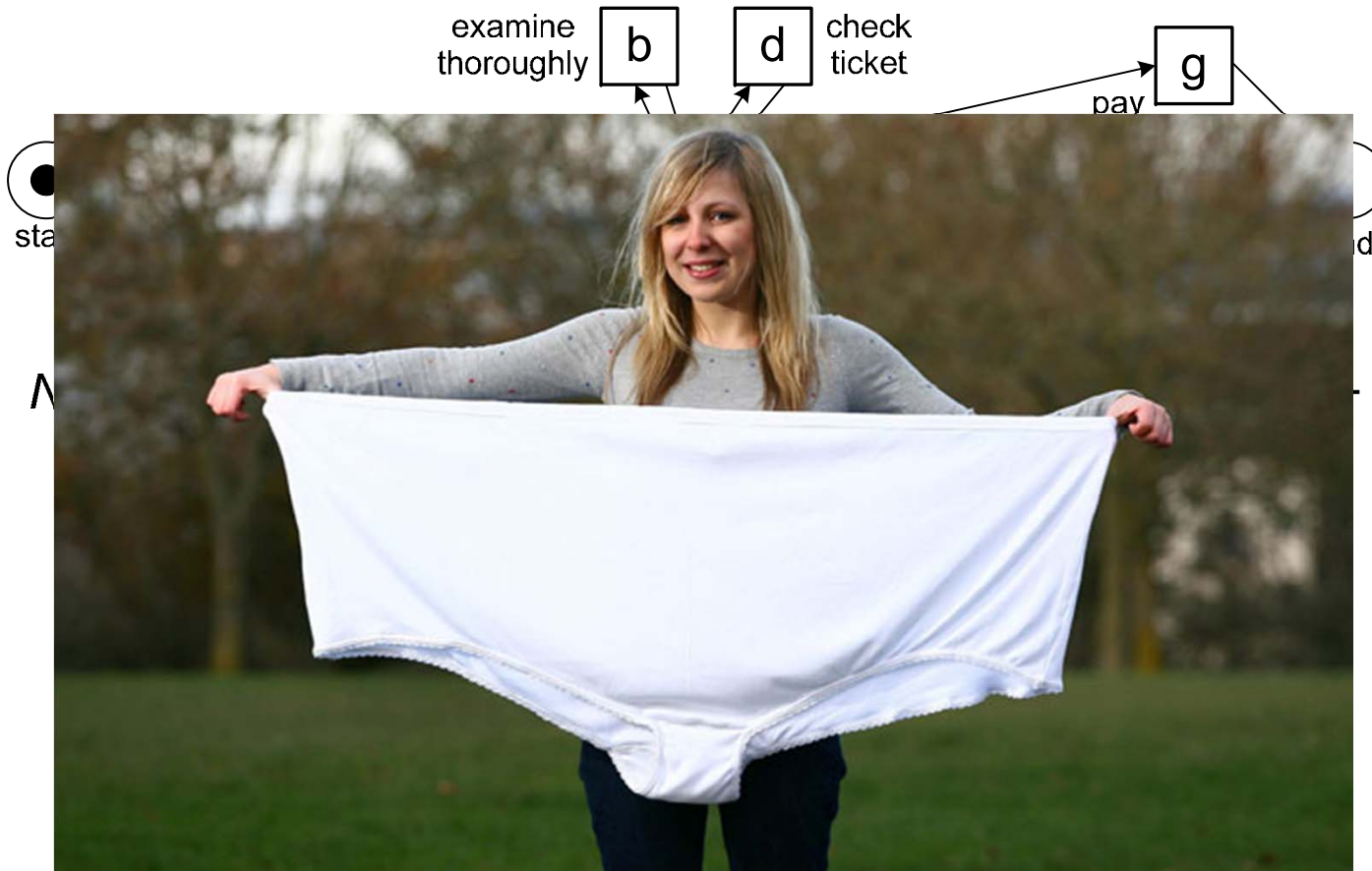
#	trace
455	acdeh
191	abdeg
177	adceh
144	abdeh
111	acdeg
82	adceg
56	adbeh
47	acdefdbeh
38	adbeg
33	acdefbdeh
14	acdefbdeg
11	acdefdbeg
9	adcefcdeh
8	adcefdbeh
5	adcefbdeg
3	acdefbdefdbeg
2	adcefdbeg
2	adcefbdefbdeg
1	adcefdbefbdeh
1	adbefbdefdbeg
1	adcefdbefcdefdbeg
1391	

Model N₂



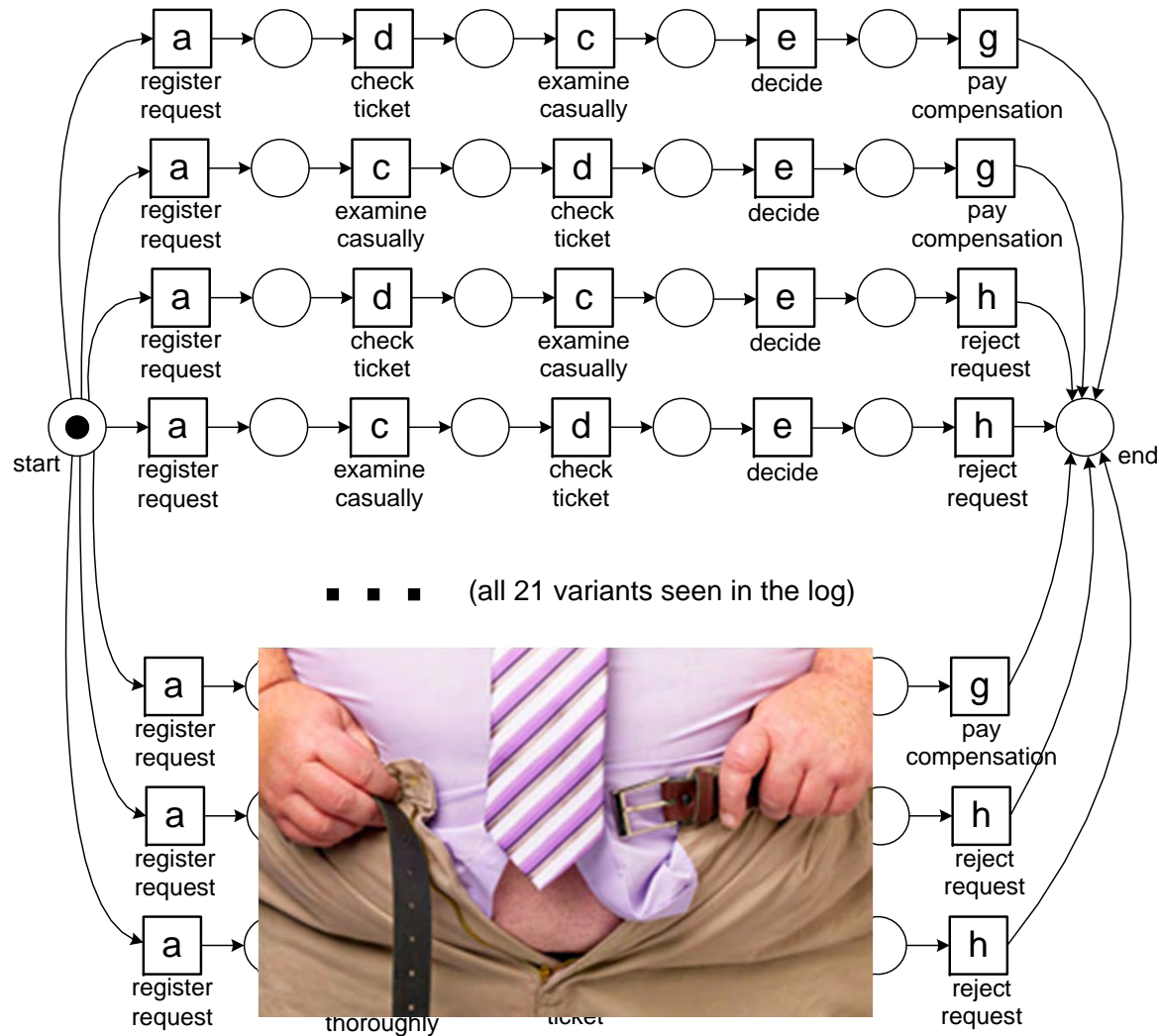
#	trace
455	acdeh
191	abdeg
177	adceh
144	abdeh
111	acdeg
82	adceg
56	adbeh
47	acdefdbeh
38	adbeg
33	acdefbdeh
14	acdefdbdeg
11	acdefdbeg
9	adcefcdeh
8	adcefdbeh
5	adcefbdeg
3	acdefbdefdbeg
2	adcefdbeg
2	adcefbdefbdeg
1	adcefdbefbdeh
1	adbefbdefdbeg
1	adcefdbefcdefdbeg
1391	

Model N₃



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

Model N₄



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

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

Process Discovery

Process Discovery (small selection)

automata-based learning

distributed genetic mining

heuristic mining

language-based regions

genetic mining

state-based regions

stochastic task graphs

LTL mining

fuzzy mining

neural networks

mining block structures

hidden Markov models

α algorithm

multi-phase mining

conformal process graph

$\alpha\#$ algorithm

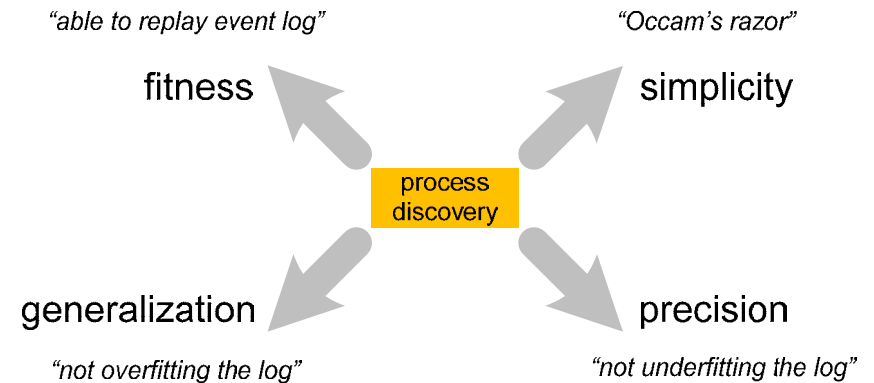
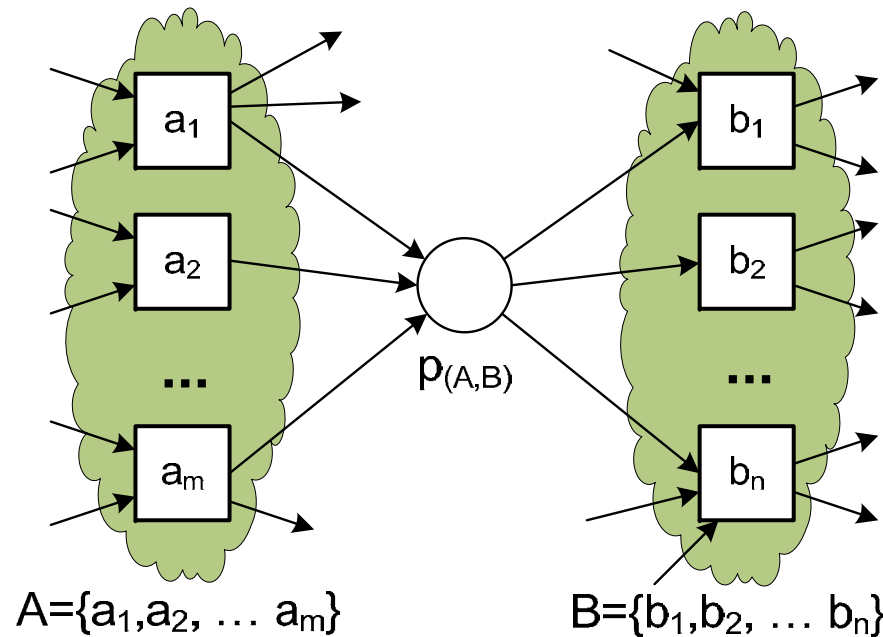
partial-order based mining

ILP mining

$\alpha++$ algorithm



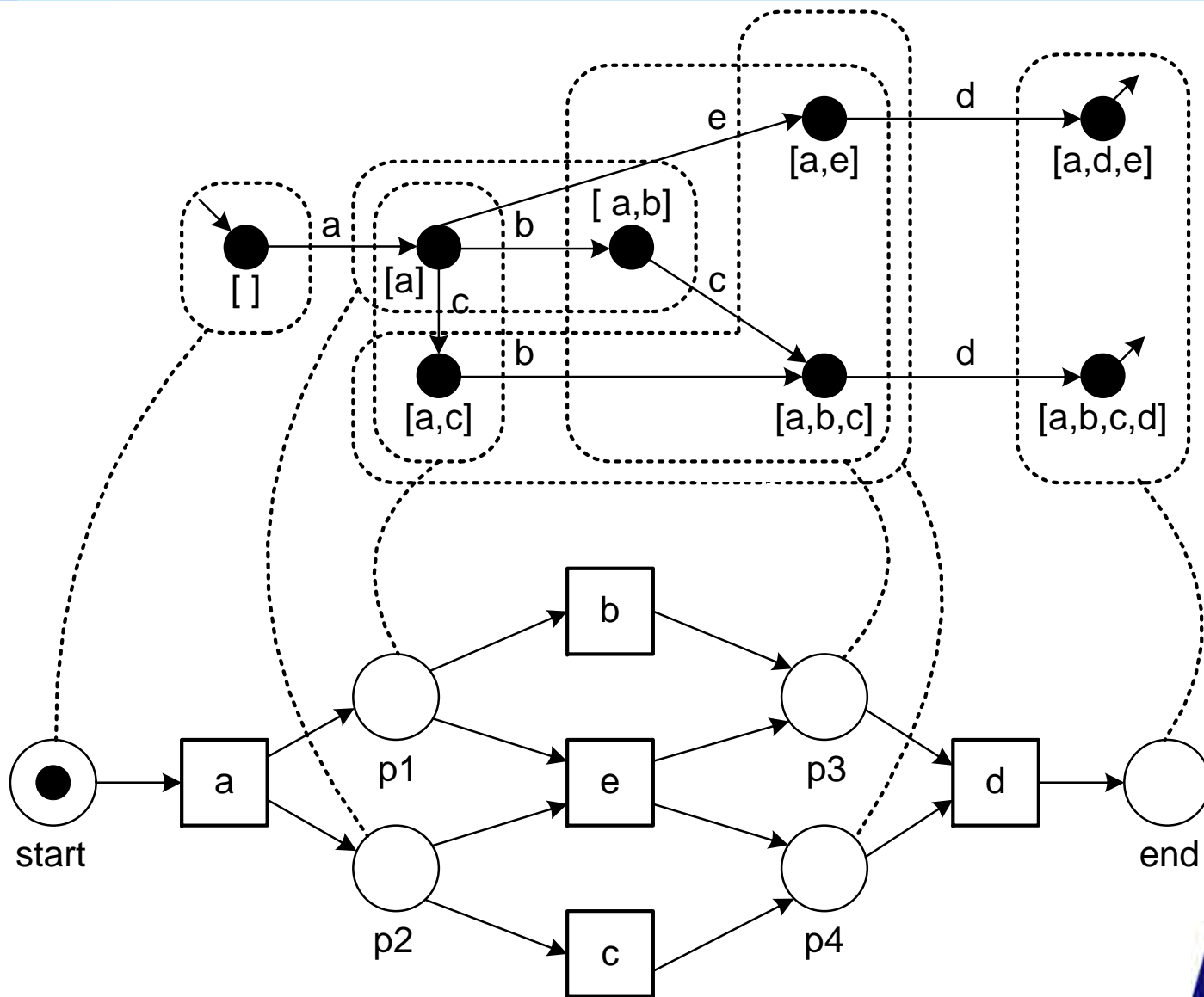
Petri net view: Just discover the places ...



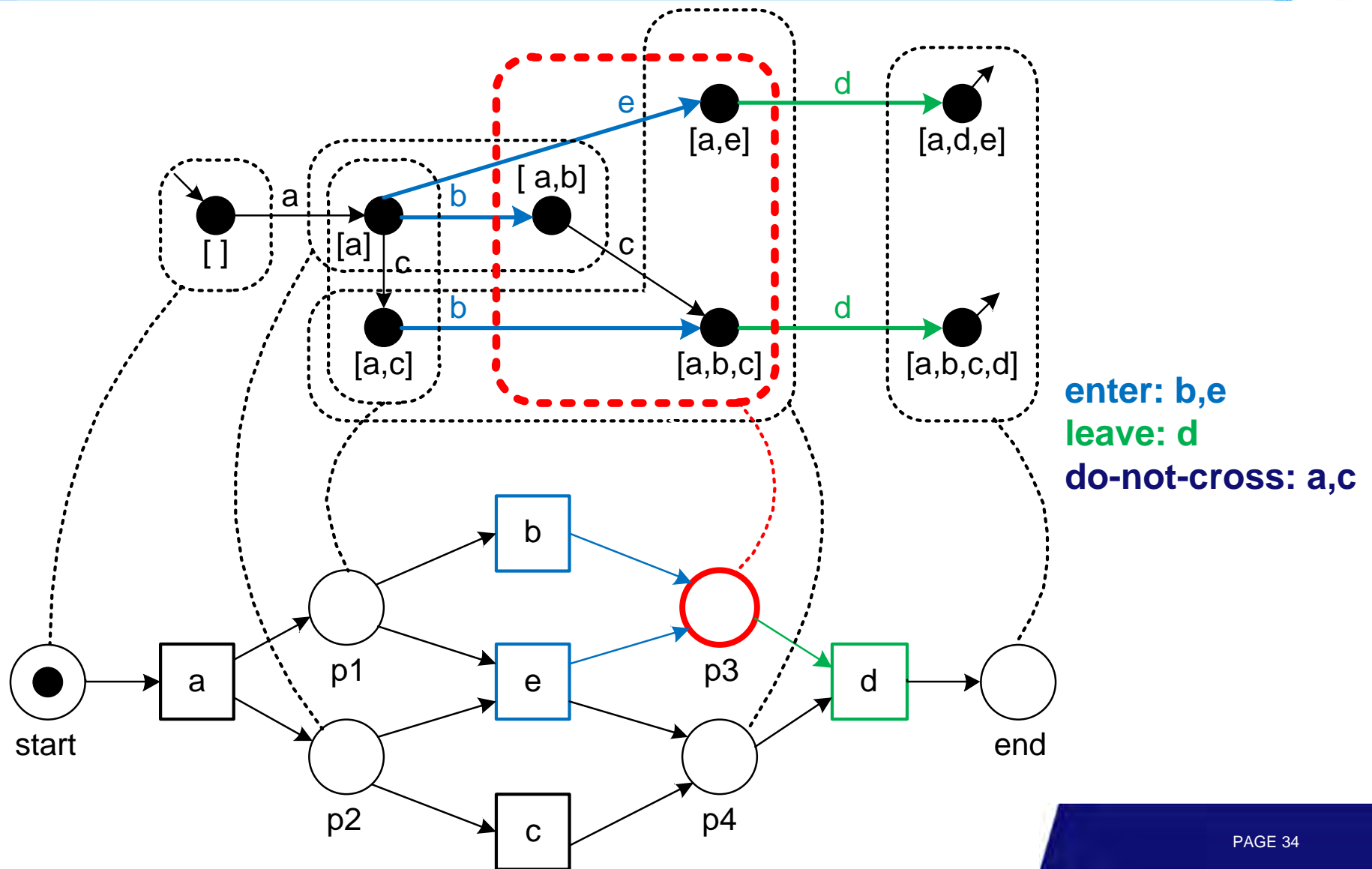
Adding a place limits behavior:

- overfitting \approx adding too many places
- underfitting \approx adding too few places

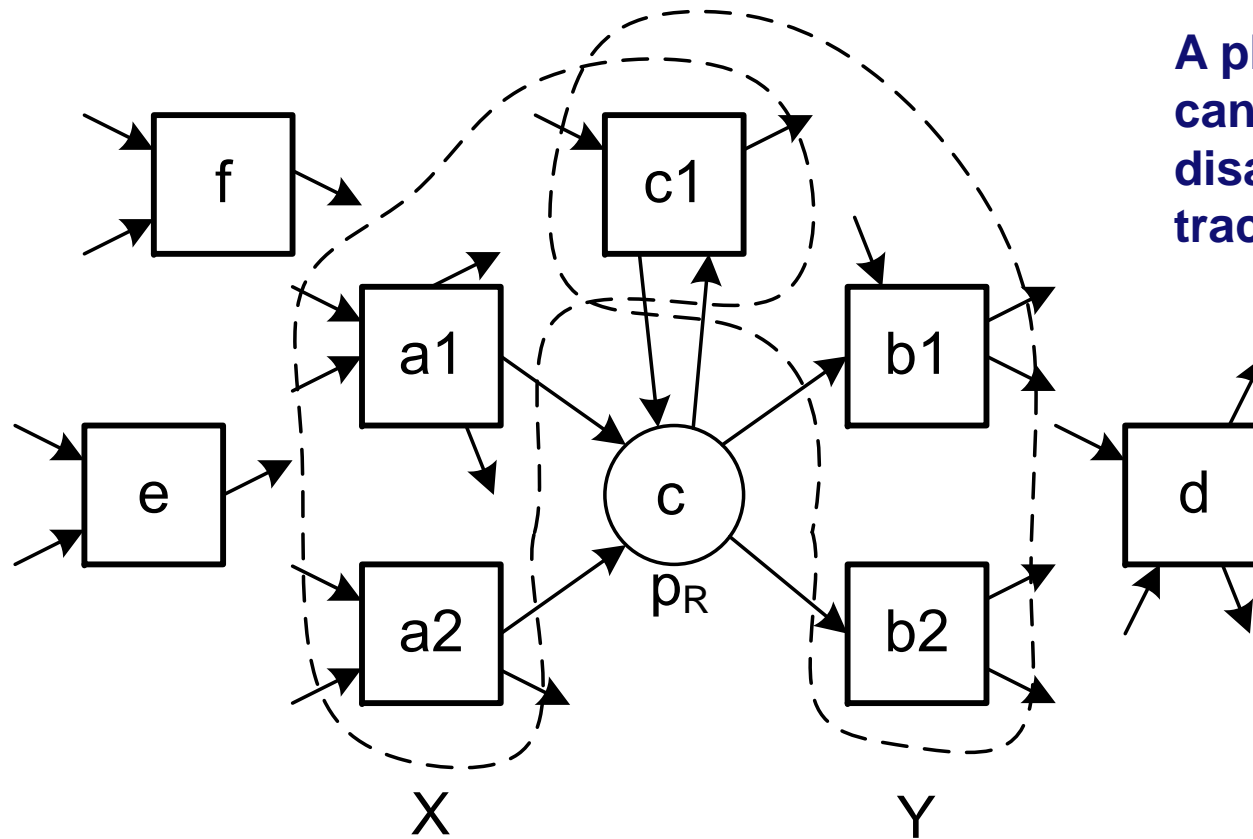
Example: Process Discovery Using State-Based Regions



Example of Region



Example: Process Discovery Using Language-Based Regions

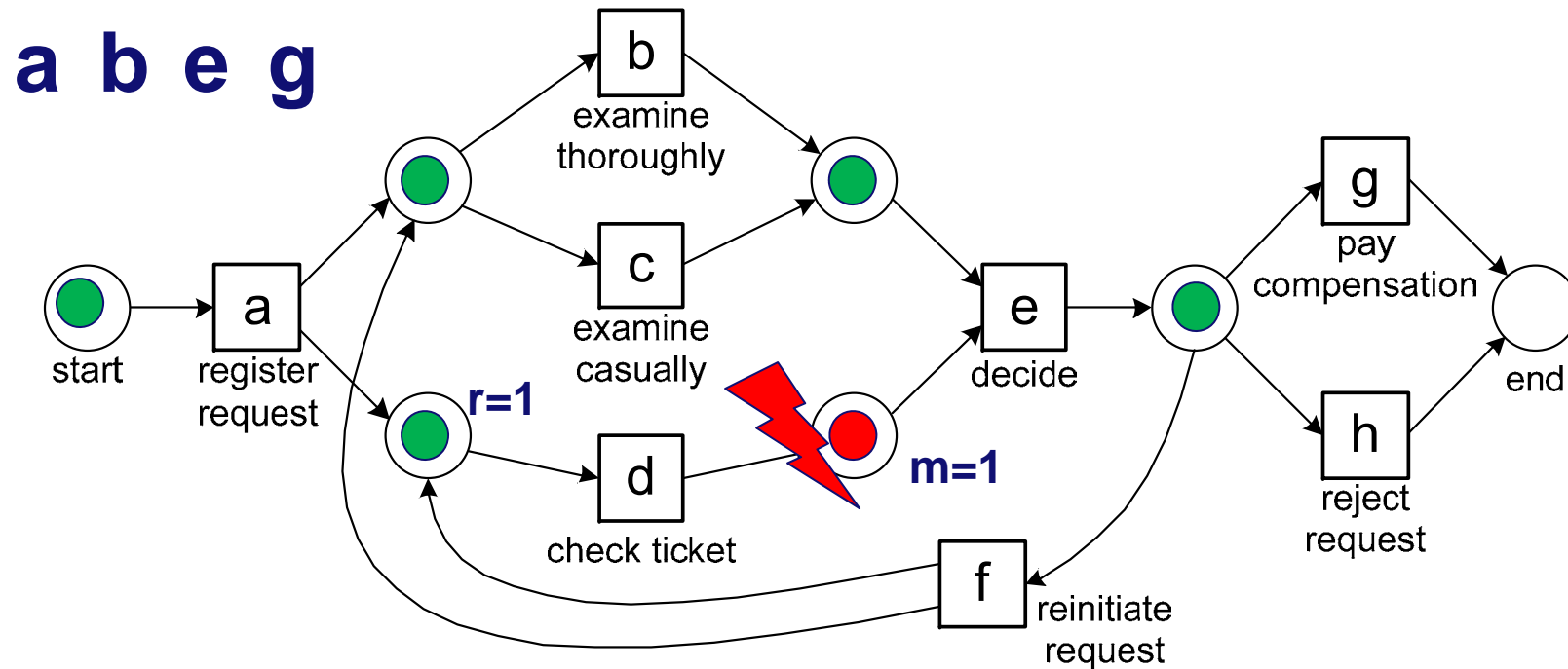


for any $\sigma \in L$, $k \in \{1, \dots, |\sigma|\}$, $\sigma_1 = hd^{k-1}(\sigma)$, $a = \sigma(k)$, $\sigma_2 = hd^k(\sigma) = \sigma_1 \oplus a$:

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

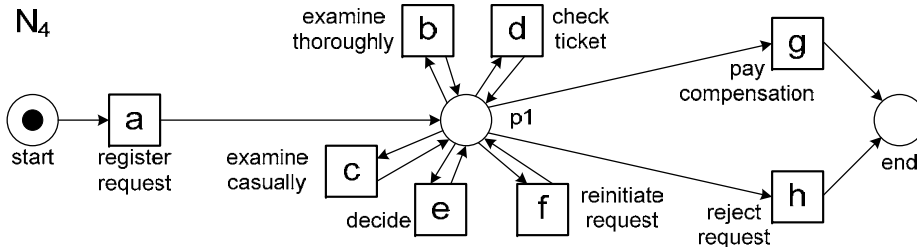
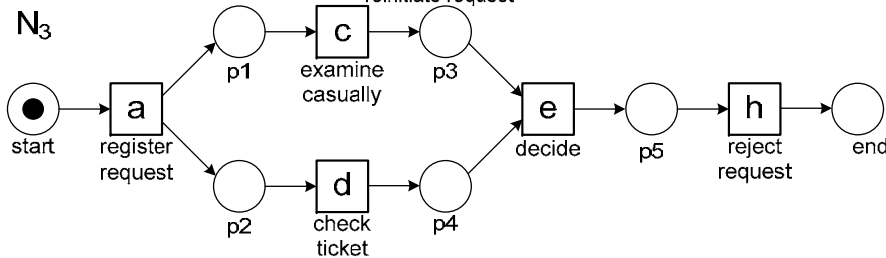
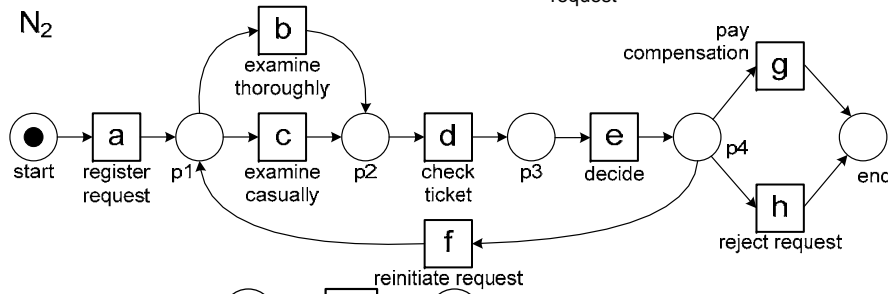
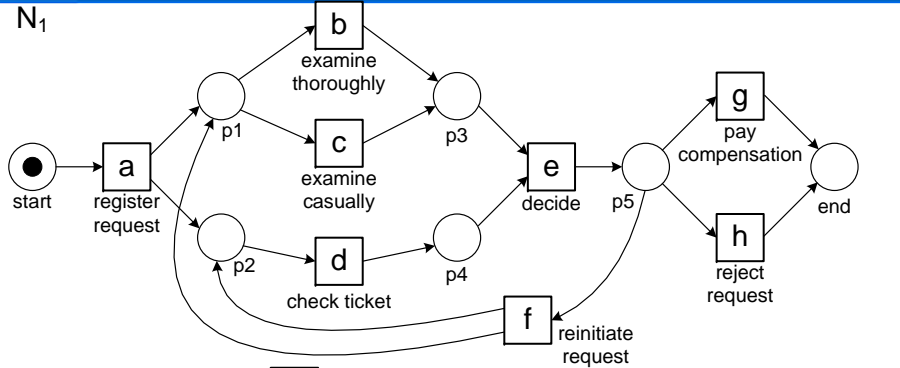
Conformance Checking

Replaying trace “abeg”



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

Can be lifted to log level



#	trace
455	acdeh
191	abdeg
177	adceh
144	abdeh
111	acdeg
82	adceg
56	adbeh

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

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

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

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

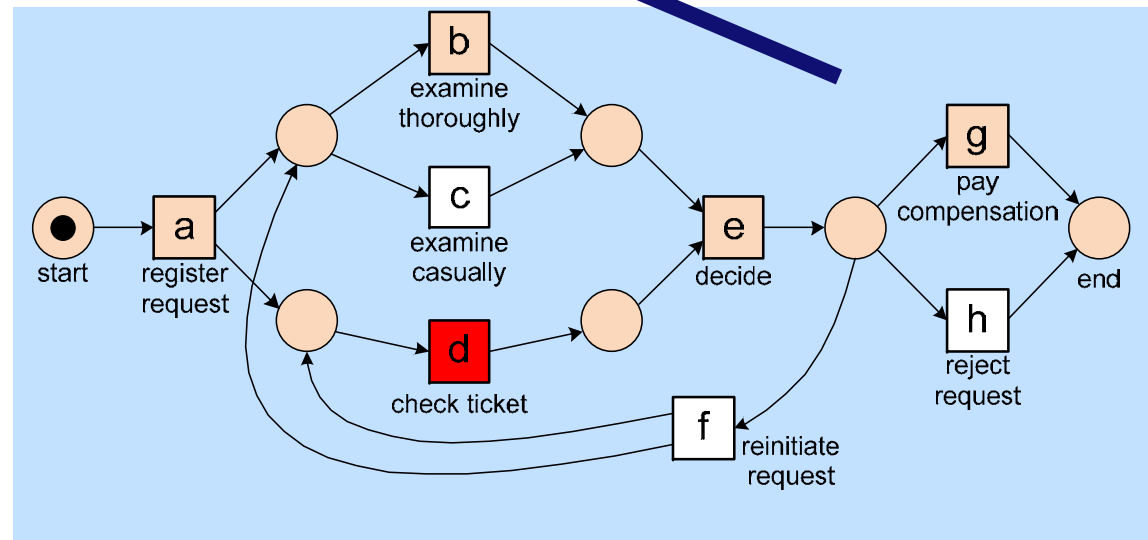
5	adcefbdeg
3	acdefbdefdbeg
2	adcefbdeg
2	adcefbdefdbeg
1	adcefdbefbdeh
1	adbefbdefdbeg
1	adcefdbefcdefdbeg
1391	

From “playing the token game” to optimal alignments ...

observed trace: “abeg”

a	b	»	e	g
a	b	d	e	g

move in
model only

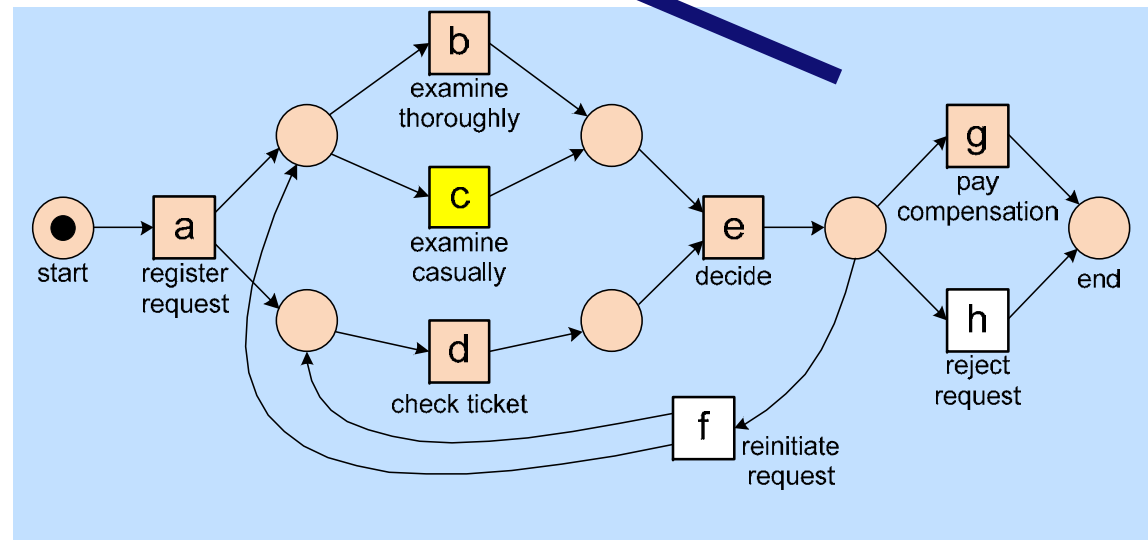


Another alignment

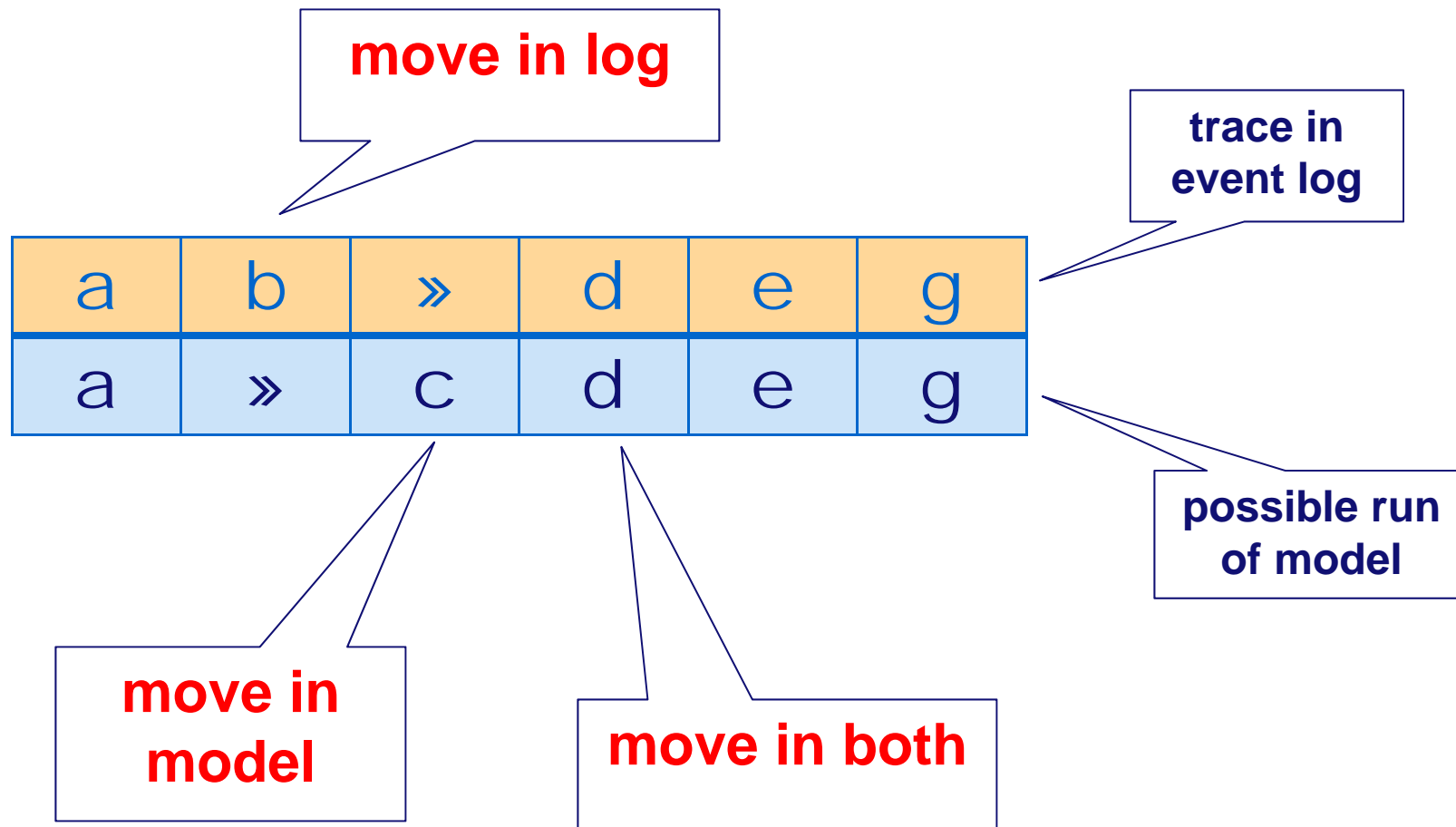
observed trace: “abcdeg”

a	b	c	d	e	g
a	b	»	d	e	g

move in
log only



Moves in an alignment



Optimal alignment describes modeled behavior closest to observed behavior

Moves have costs

...	a	...
...	»	...

...	»	...
...	a	...

...	a	...
...	a	...

...	a	...
...	b	...

- **Standard cost function:**

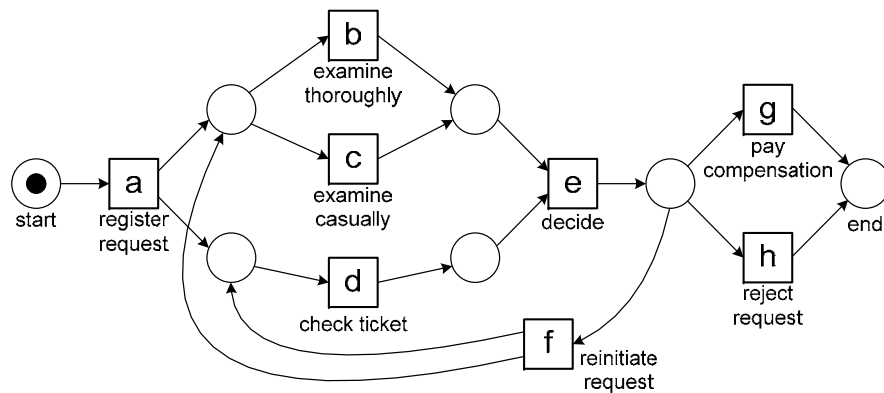
- $c(x, \text{»}) = 1$

- $c(\text{»}, y) = 1$

- $c(x, y) = 0$, if $x=y$

- $c(x, y) = \infty$, if $x \neq y$

Non-fitting trace: abefdeg



abefdeg

a	b	»	e	f	d	»	e	g
a	b	d	e	f	d	b	e	g

2

a	b	e	f	d	e	g
a	b	»	»	d	e	g

2

Any cost structure is possible

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

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

Fitness

1.0

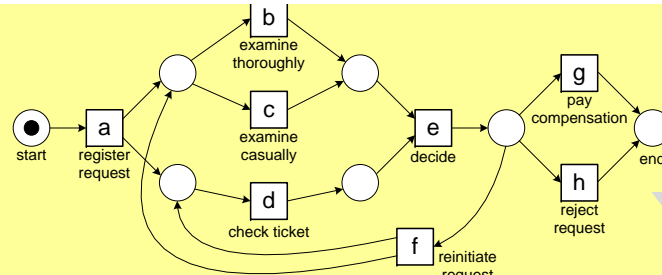
Our A* algorithm exploits the Petri net marking equation and uses other “tricks” to prune the search space.

0.8

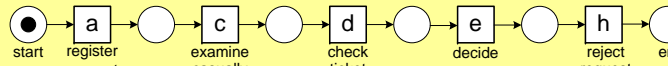
1.0

1.0

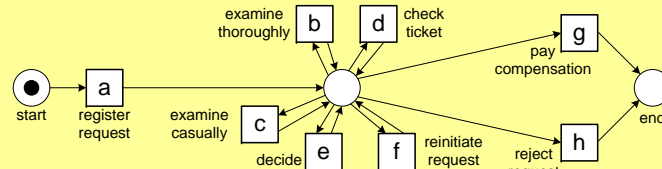
Aligned event log is starting point for other types of analysis.



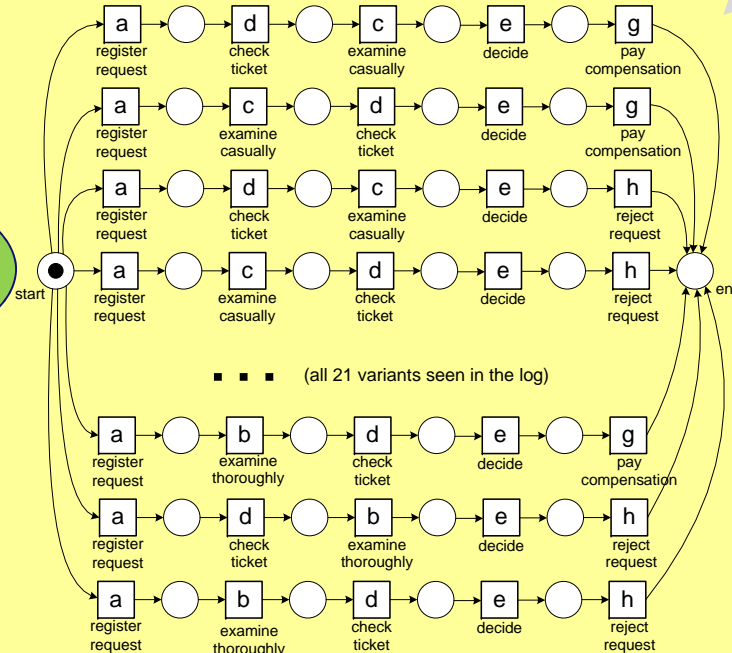
N_1 : fitness = +, precision = +, generalization = +, simplicity = +



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



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

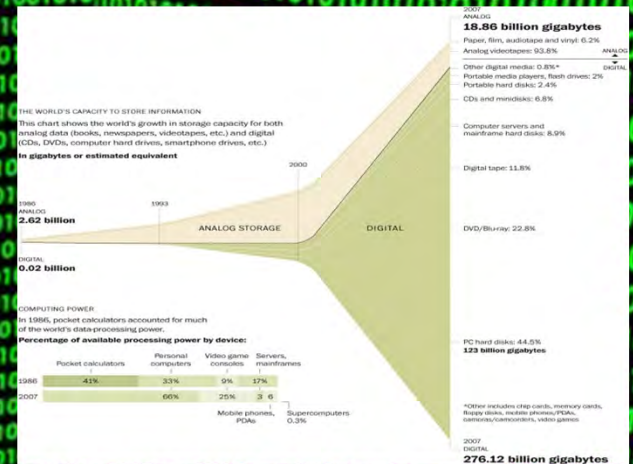
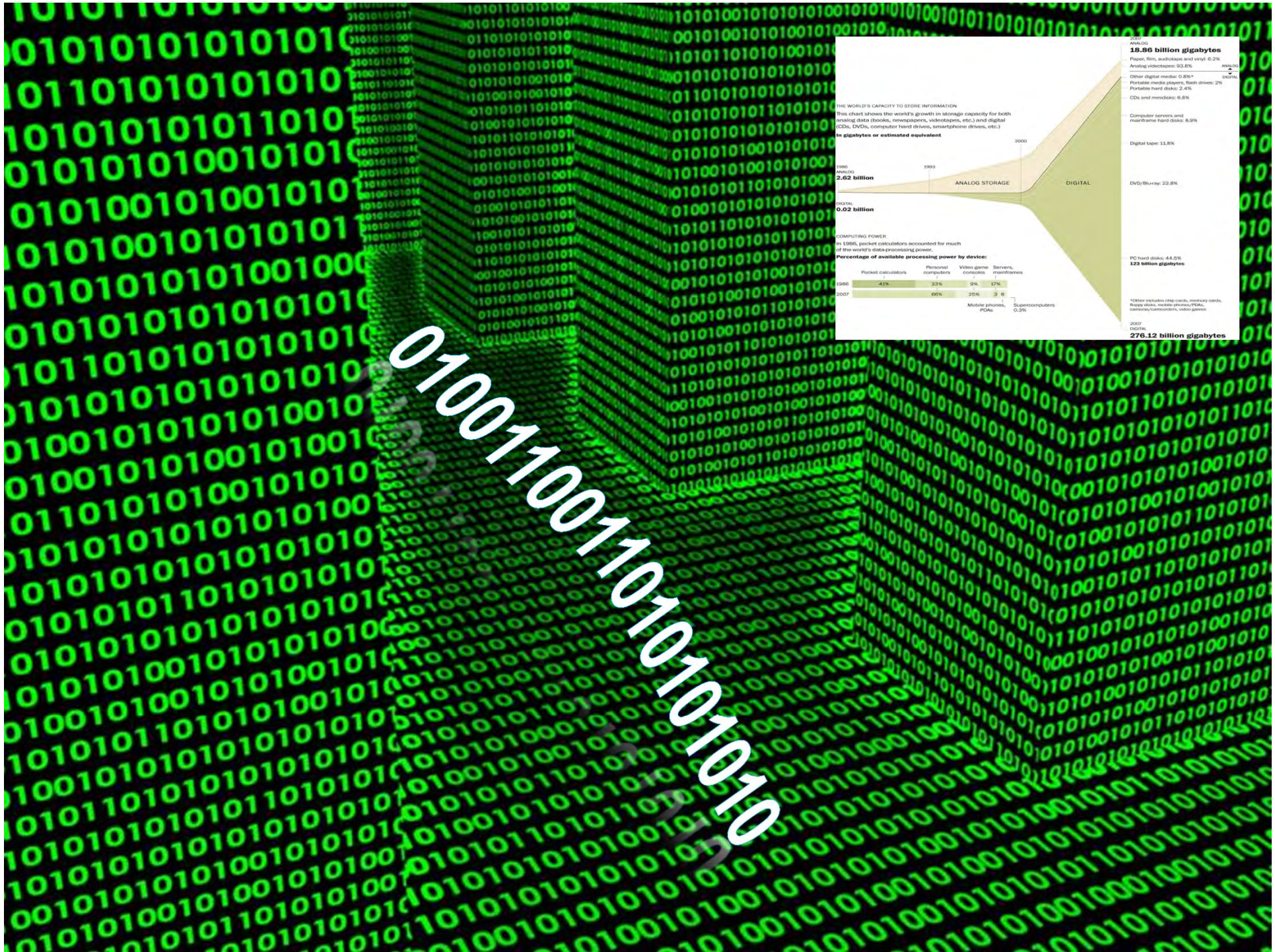


■ ■ ■ (all 21 variants seen in the log)

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

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

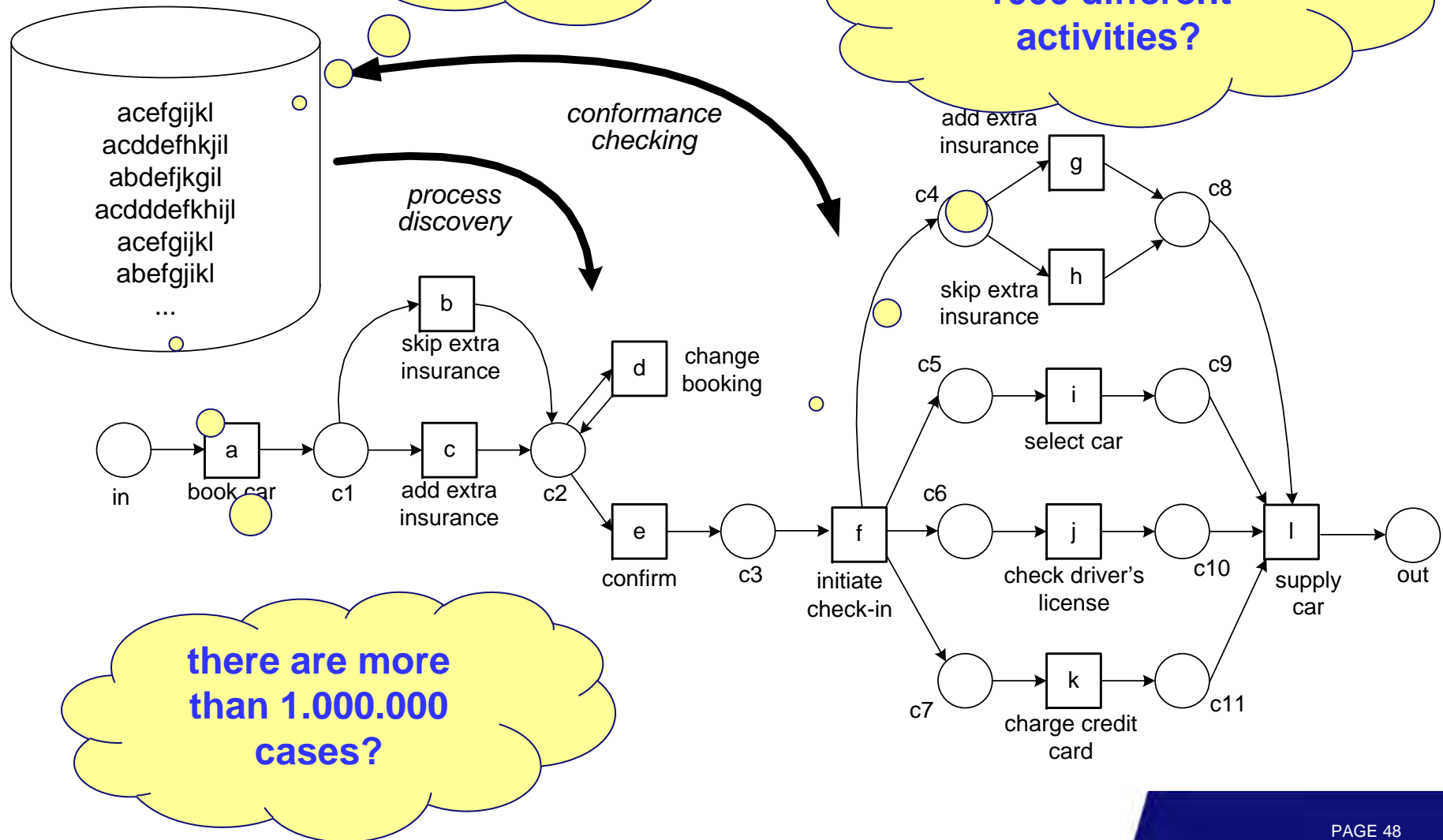
Distributing/Decomposing “Big Data” Process Mining Problems



What if?

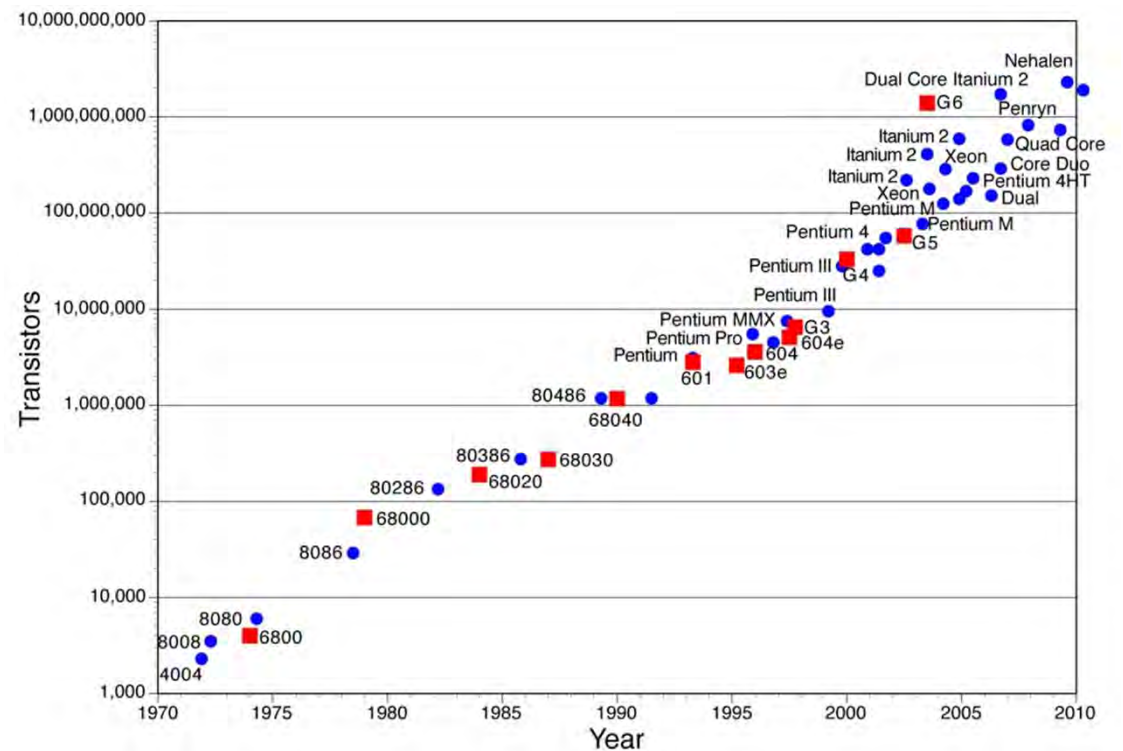
there are more
than 100.000.000
events?

there are more than
1000 different
activities?

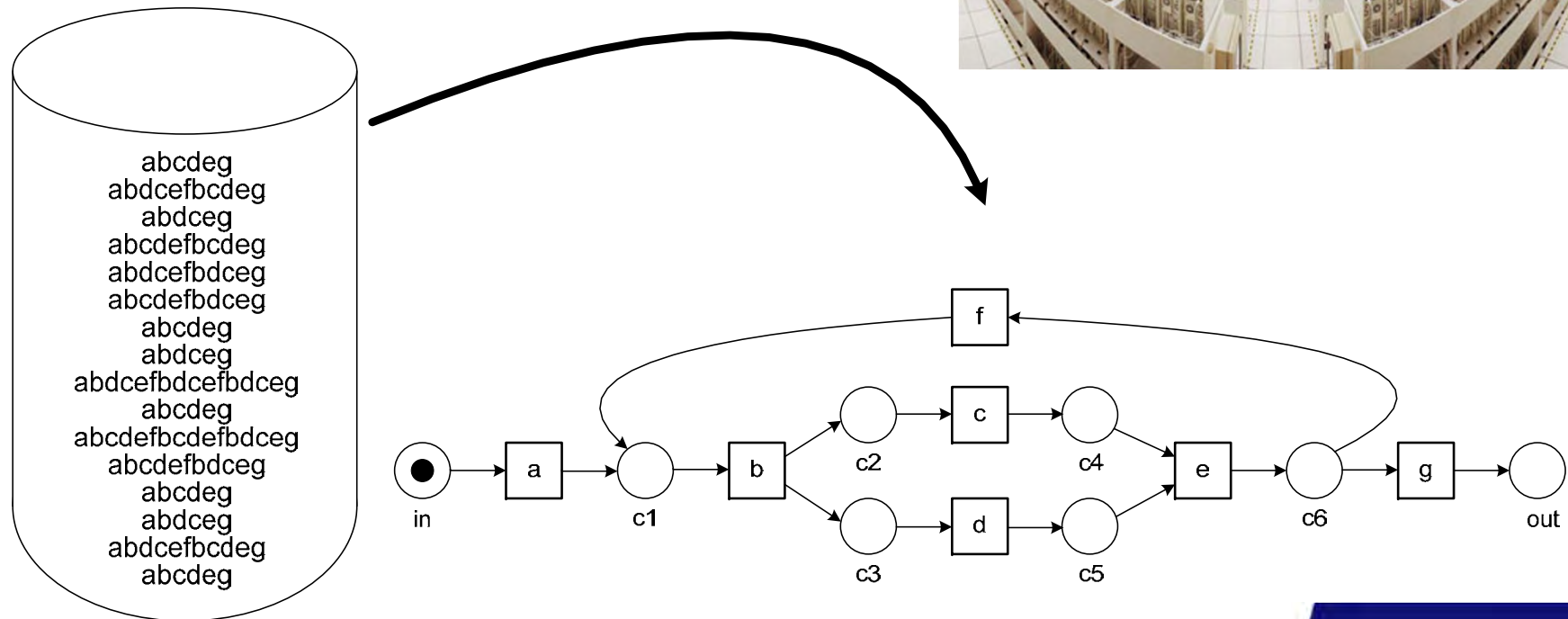
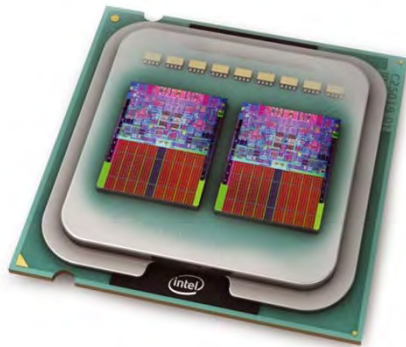


Distributed computing

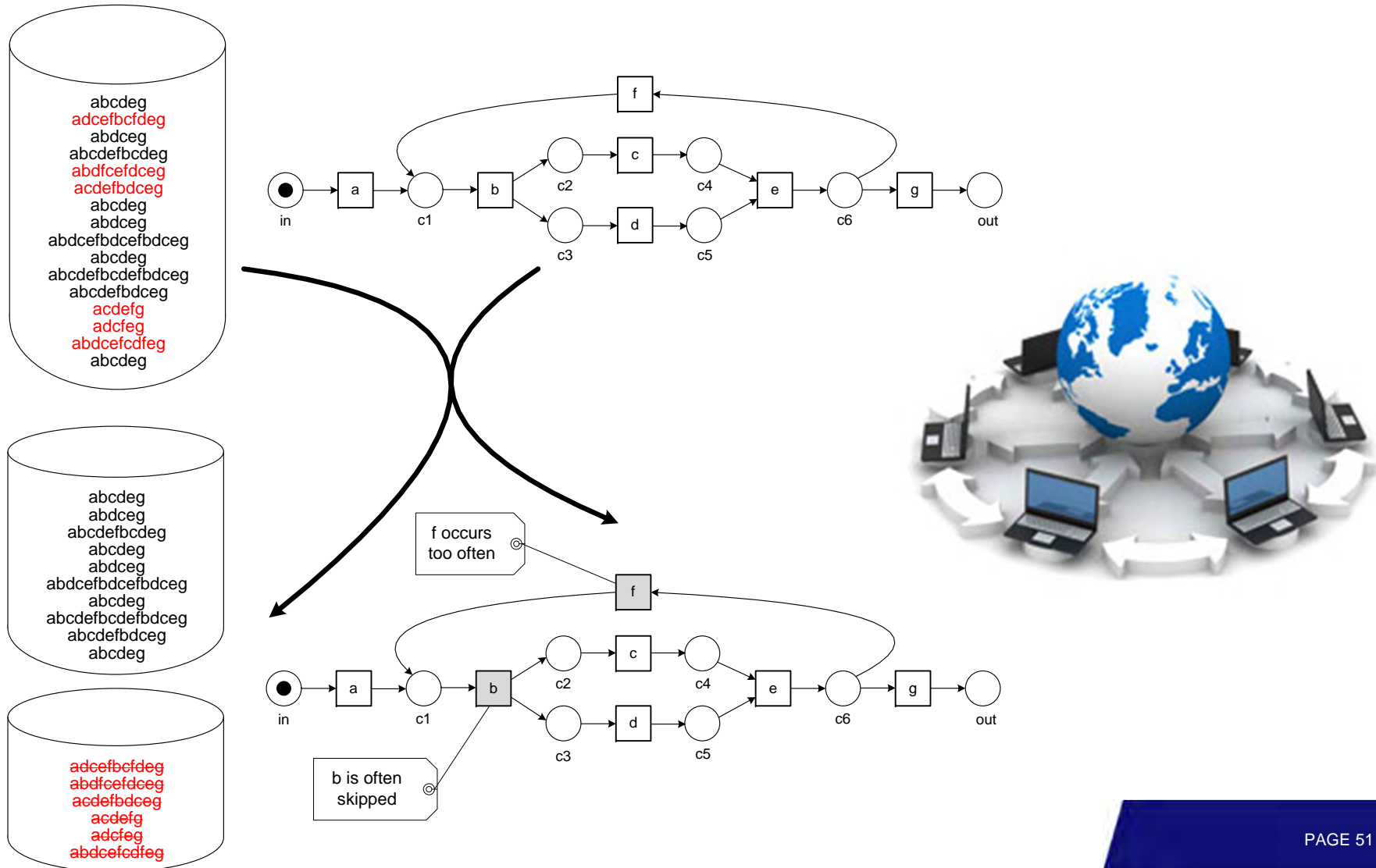
- multicore CPU
- manycore GPU
- cluster computing
- grid computing
- cloud computing
- ...



How to distribute process discovery?



How to distribute conformance checking?



Classification based on partitioning of event log: vertical and horizontal

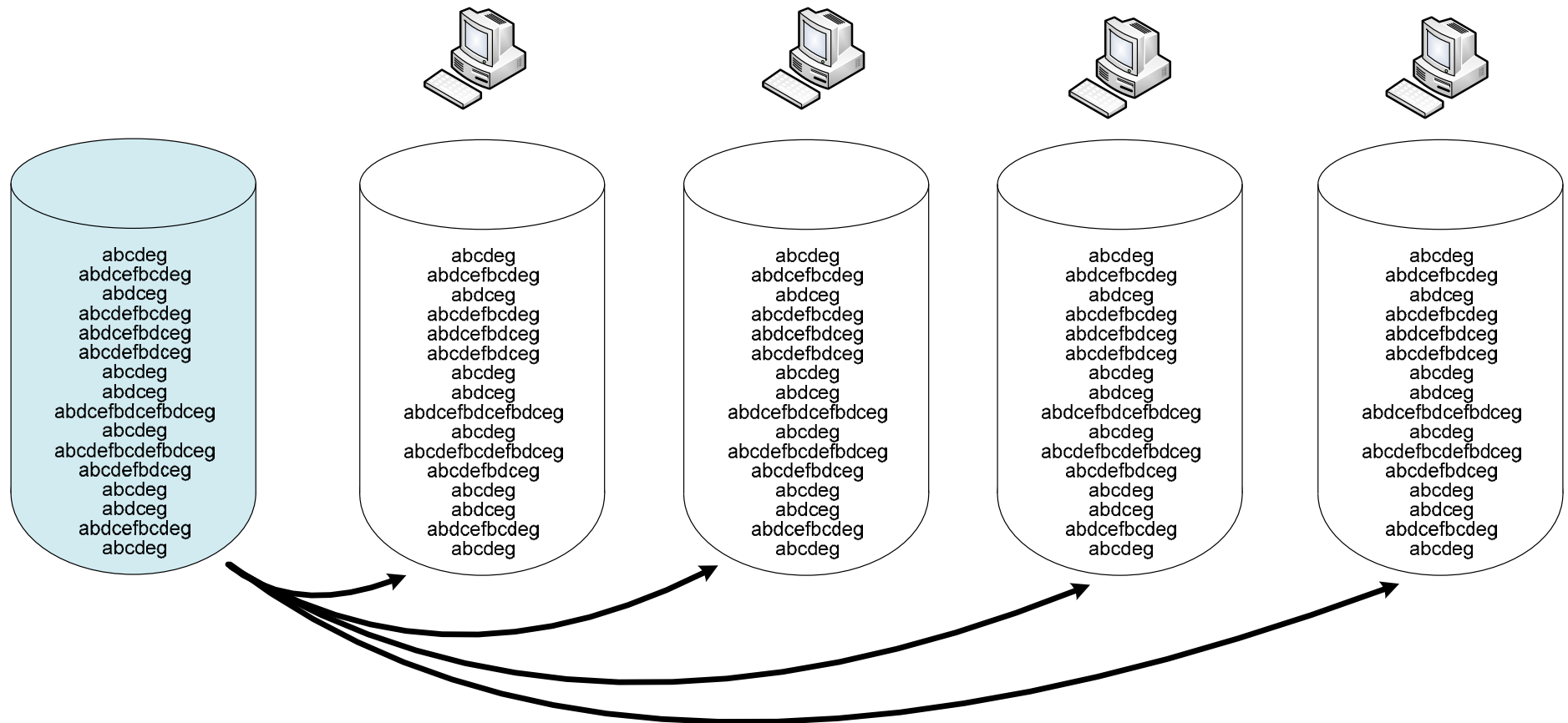


**sets of
cases**

**sets of
activities**



Replication: Same event log on all computing nodes

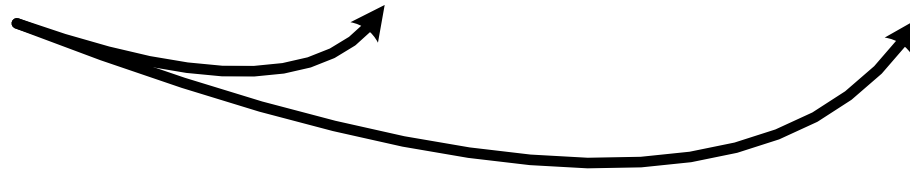
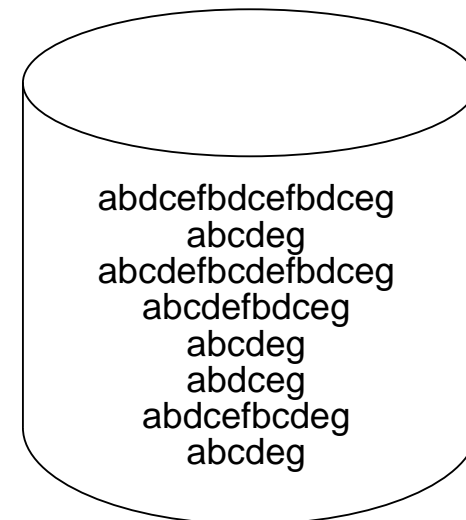
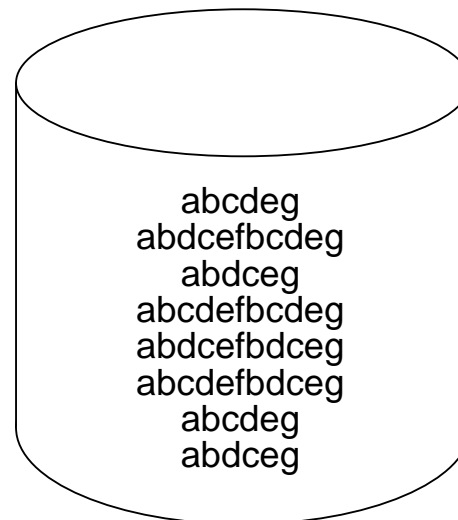
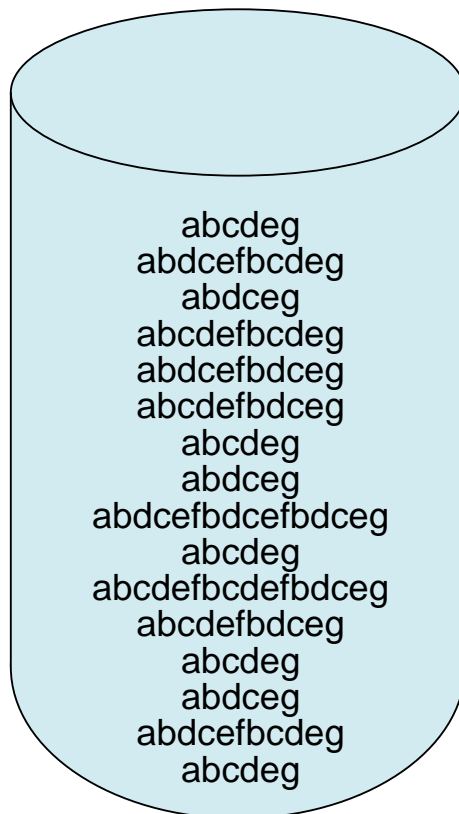
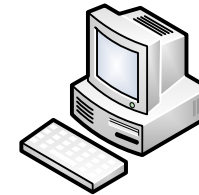
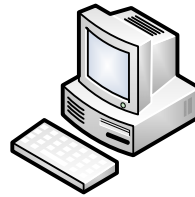


Only makes sense if random elements,
e.g., genetic process mining.

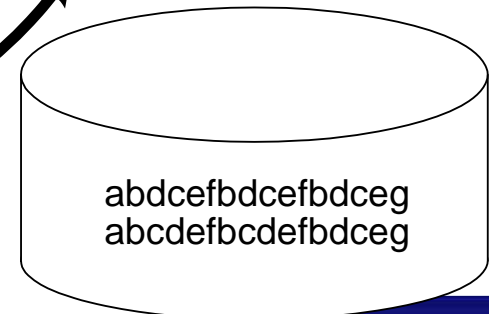
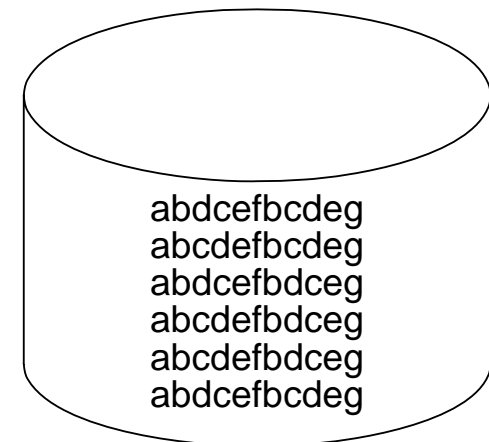
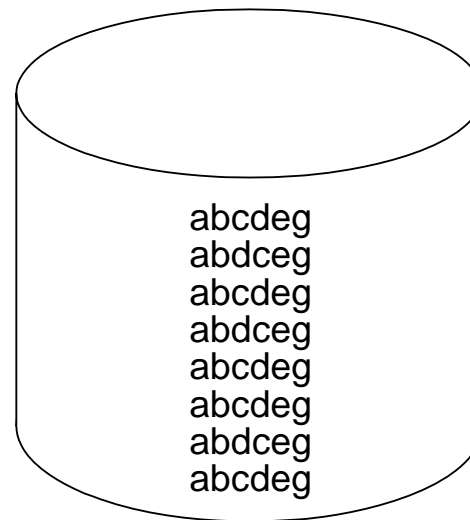
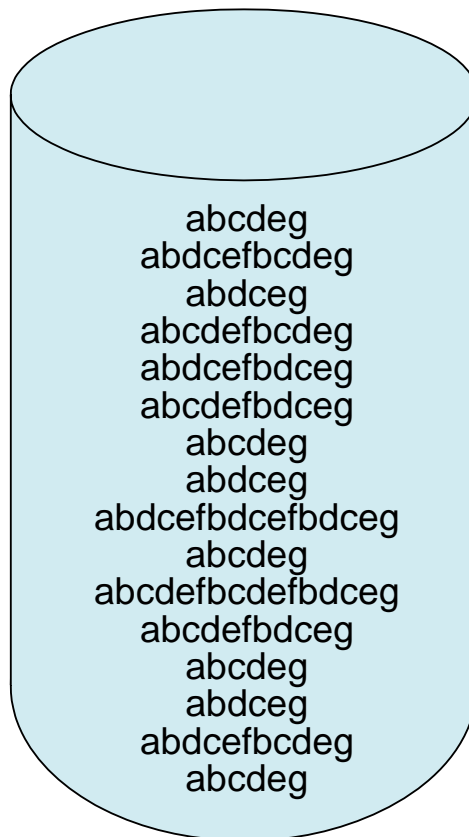
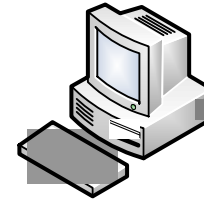
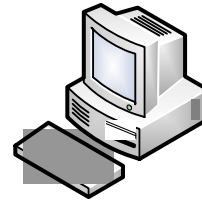
Vertical distribution I: Split cases arbitrarily



sets of
cases

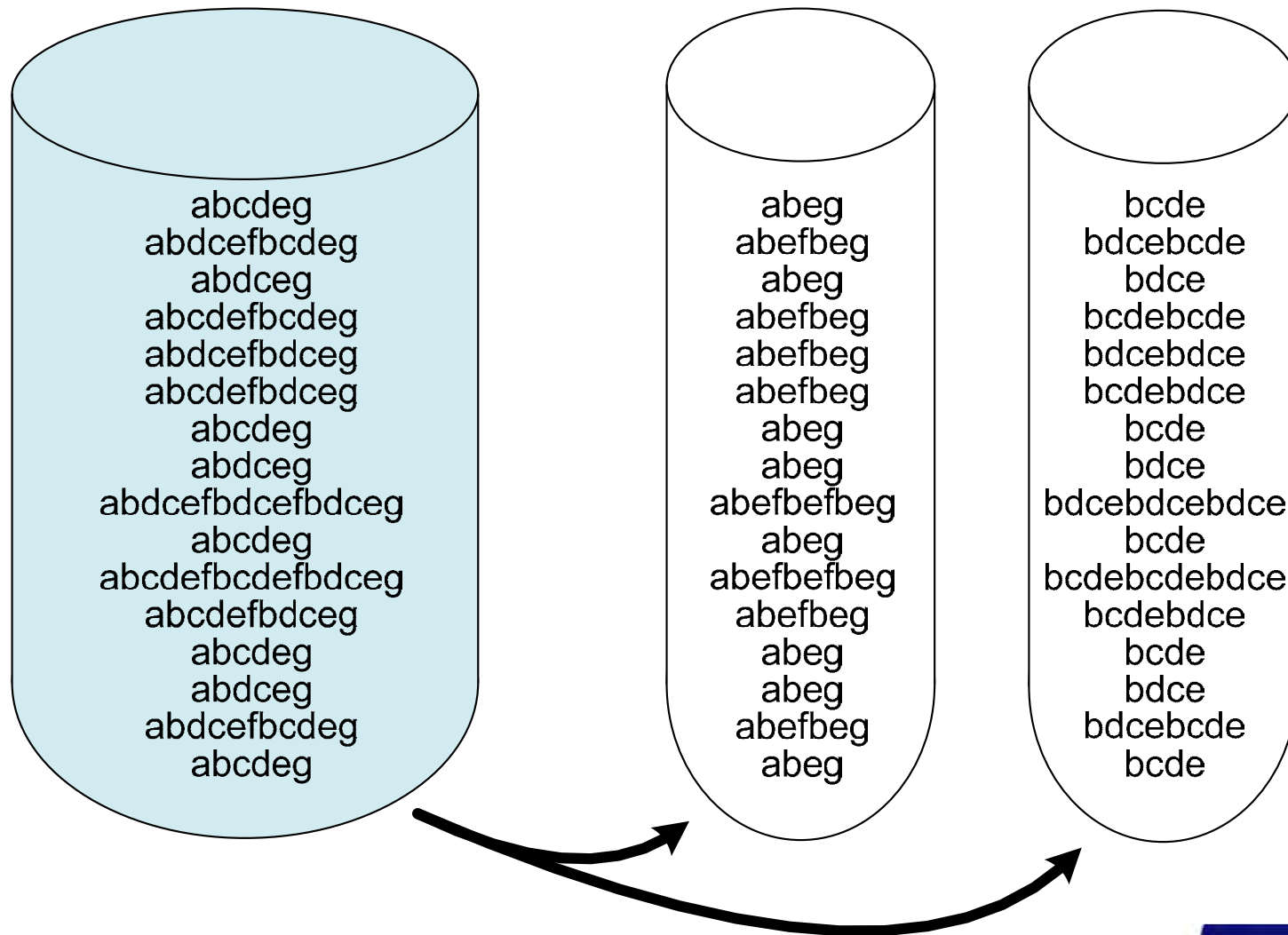


Vertical distribution II: Split cases based on a specific feature

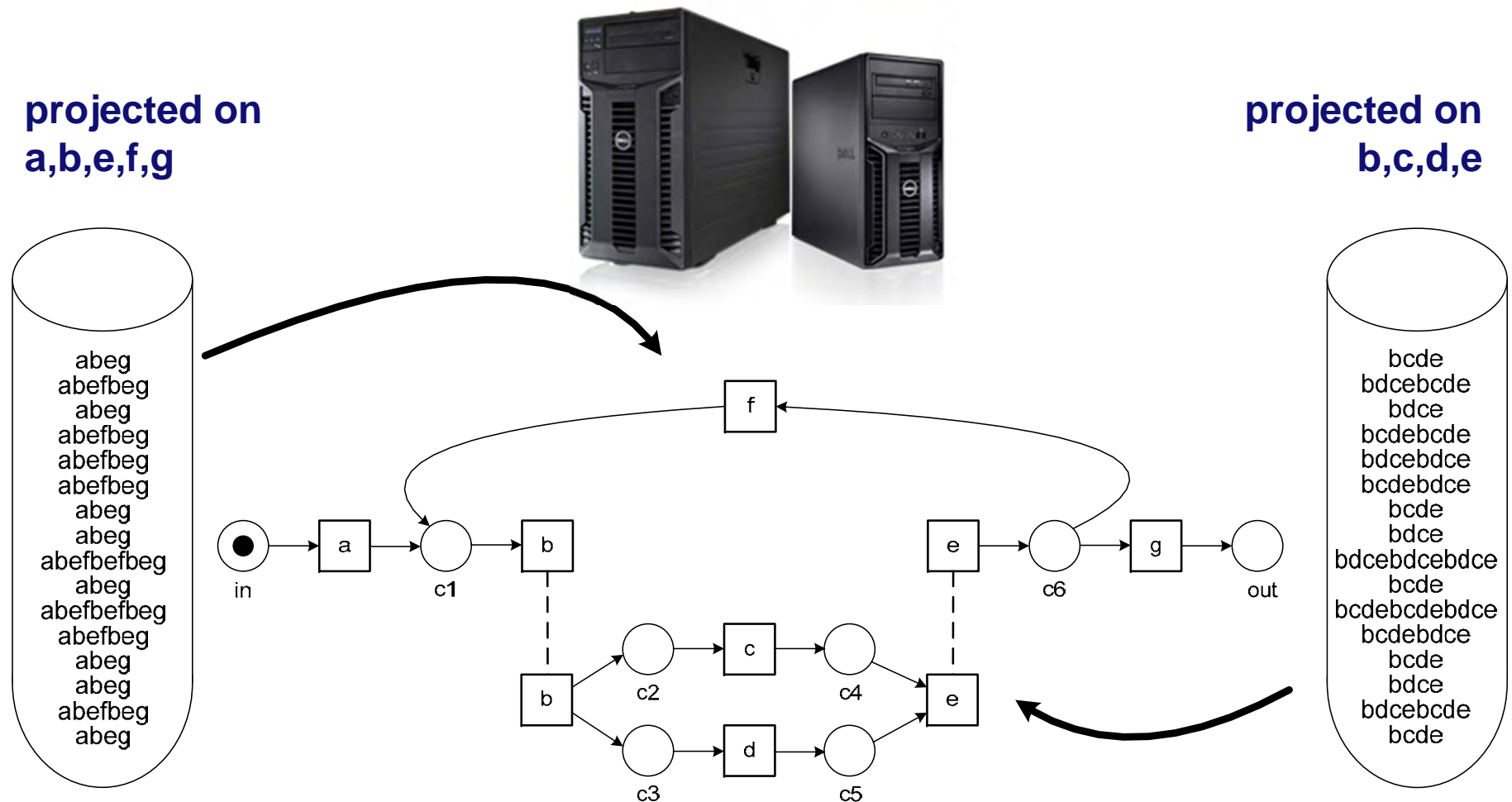


Horizontal distribution

sets of
activities



Horizontal distribution: The key idea



Passages for Horizontal Distribution

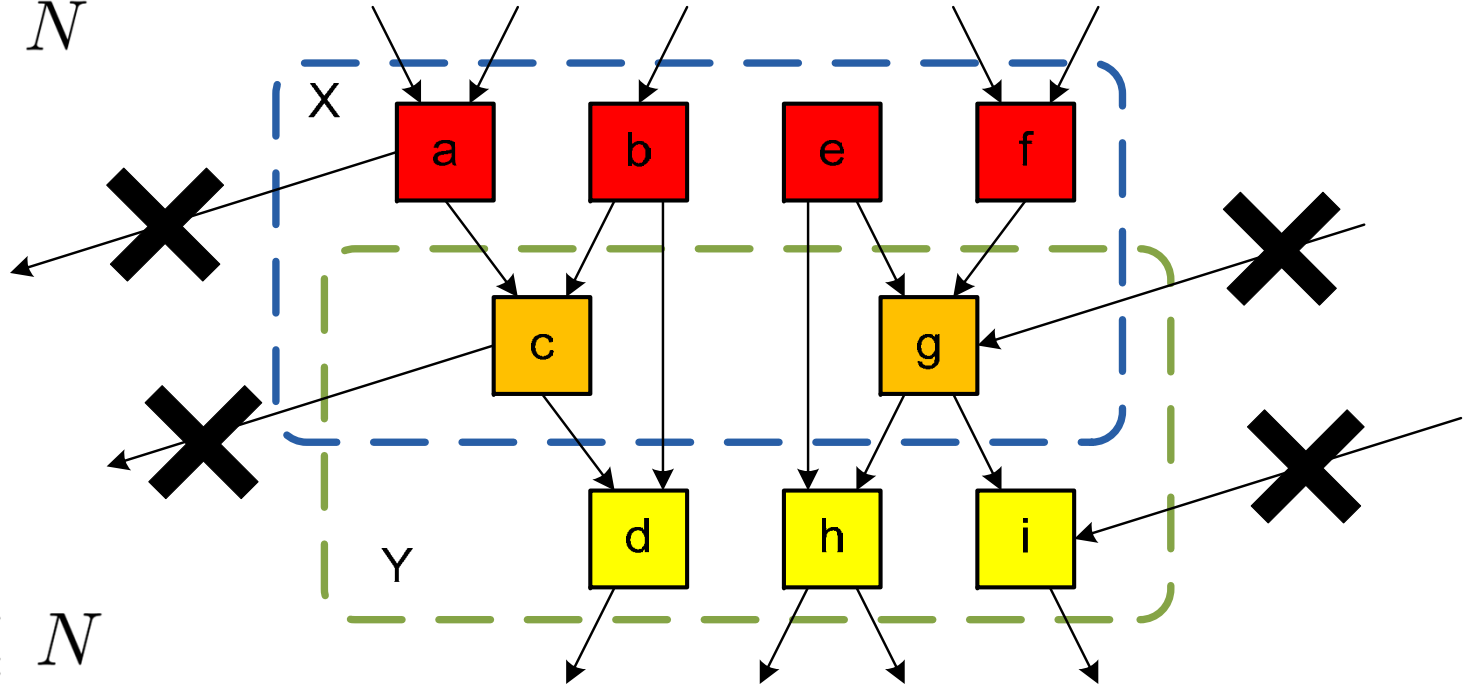
Passages



Passage $P=(X,Y)$

causal dependency:
may trigger or enable

$$\emptyset \neq X \subseteq N$$

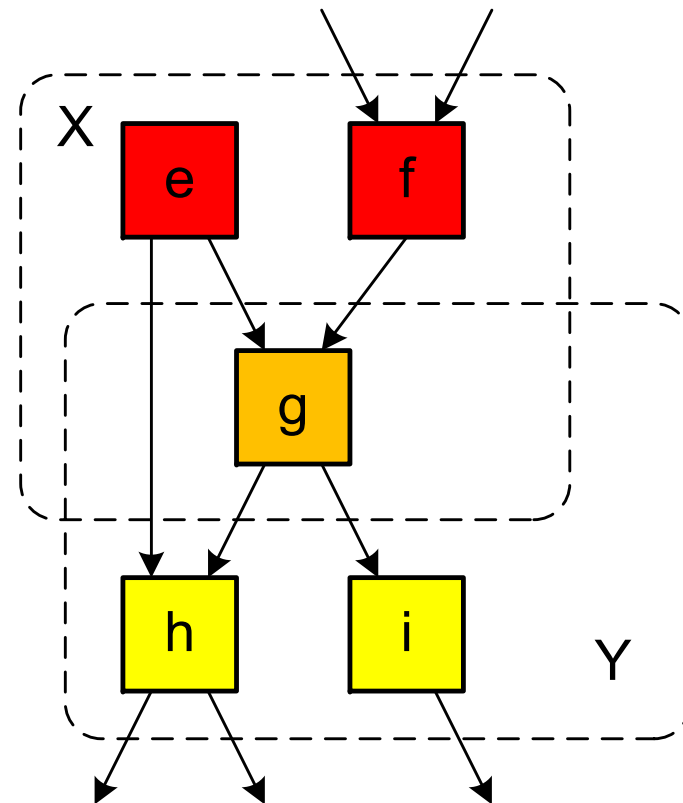
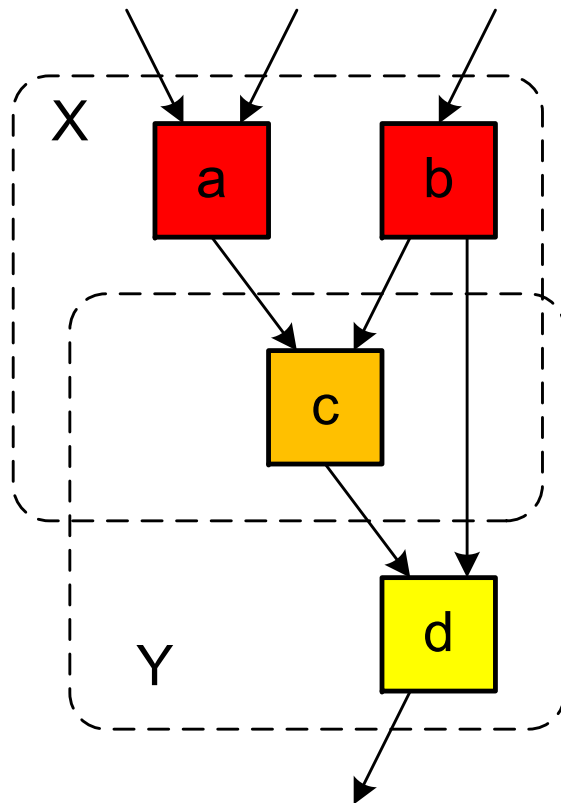


$$\emptyset \neq Y \subseteq N$$

$$X \overset{G}{\bullet} = Y$$

$$X = \overset{G}{\bullet} Y$$

Minimal passages

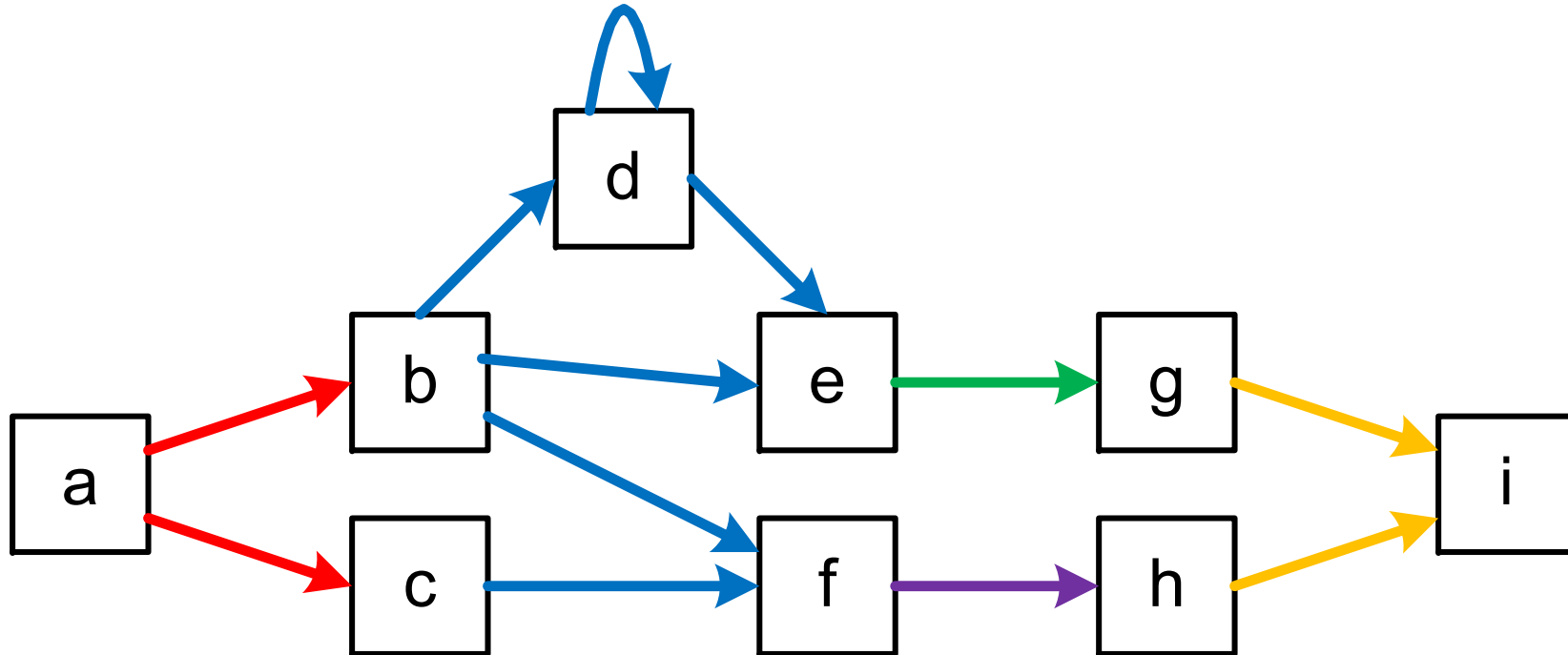


$$X \overset{G}{\bullet} = Y$$

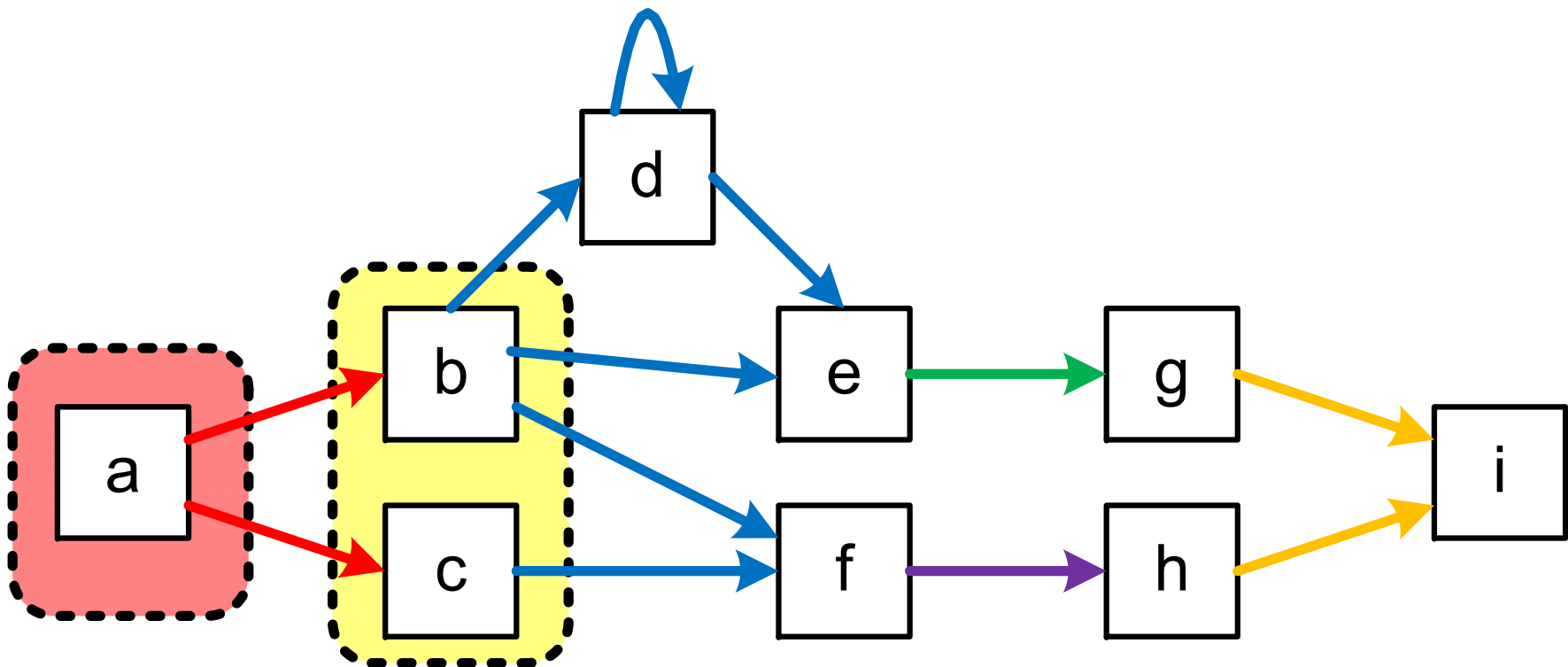
$$X = \overset{G}{\bullet} Y$$

a passage is minimal if it does not contain smaller passages

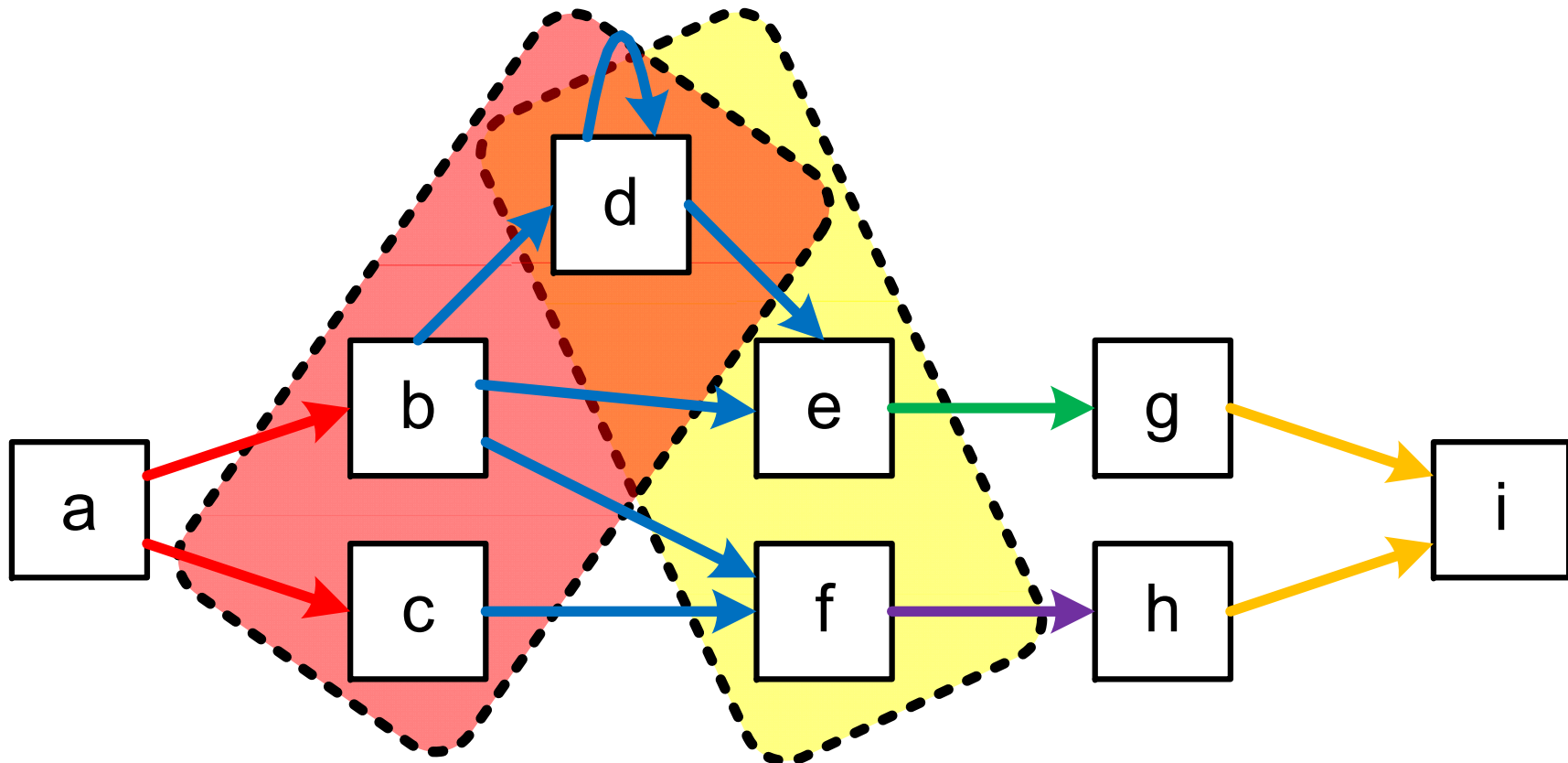
Passages define an equivalence relation on the edges in the graph



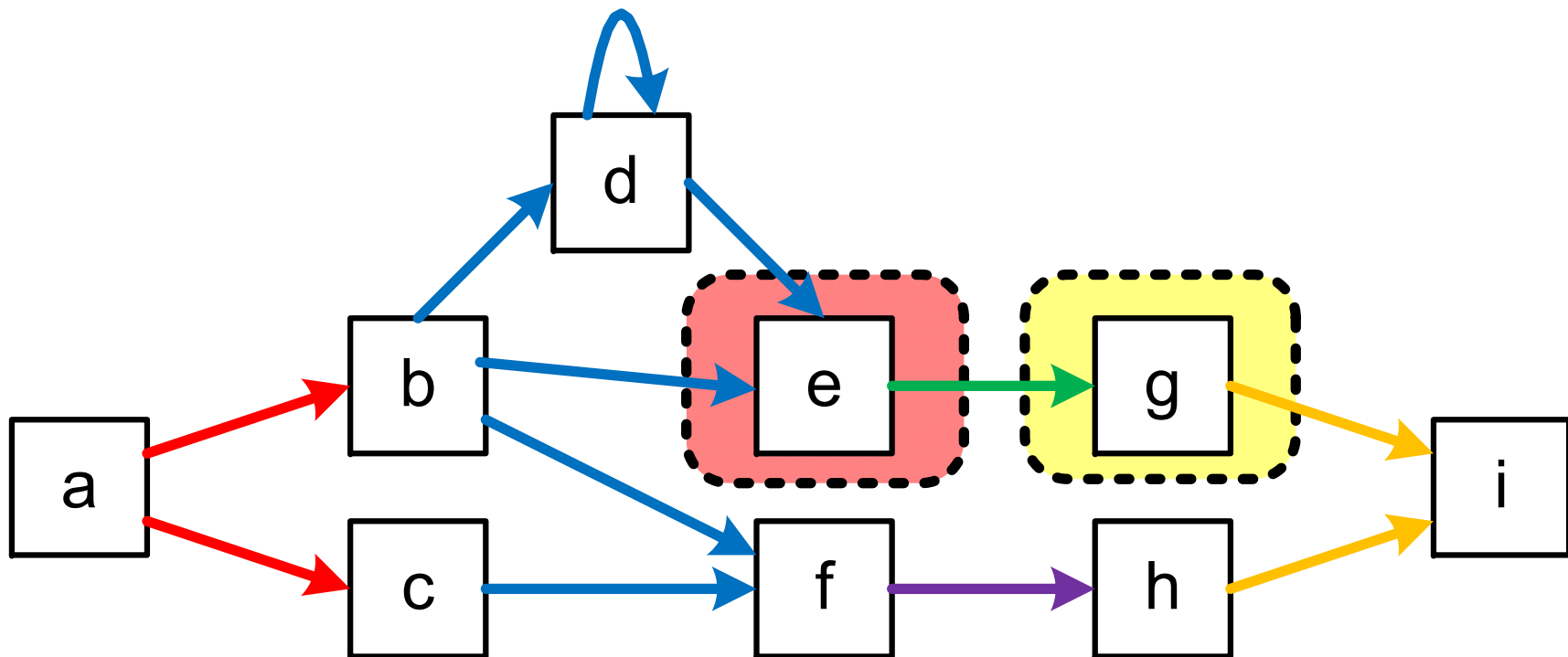
Minimal passage 1: ($\{a\}, \{b,c\}$)



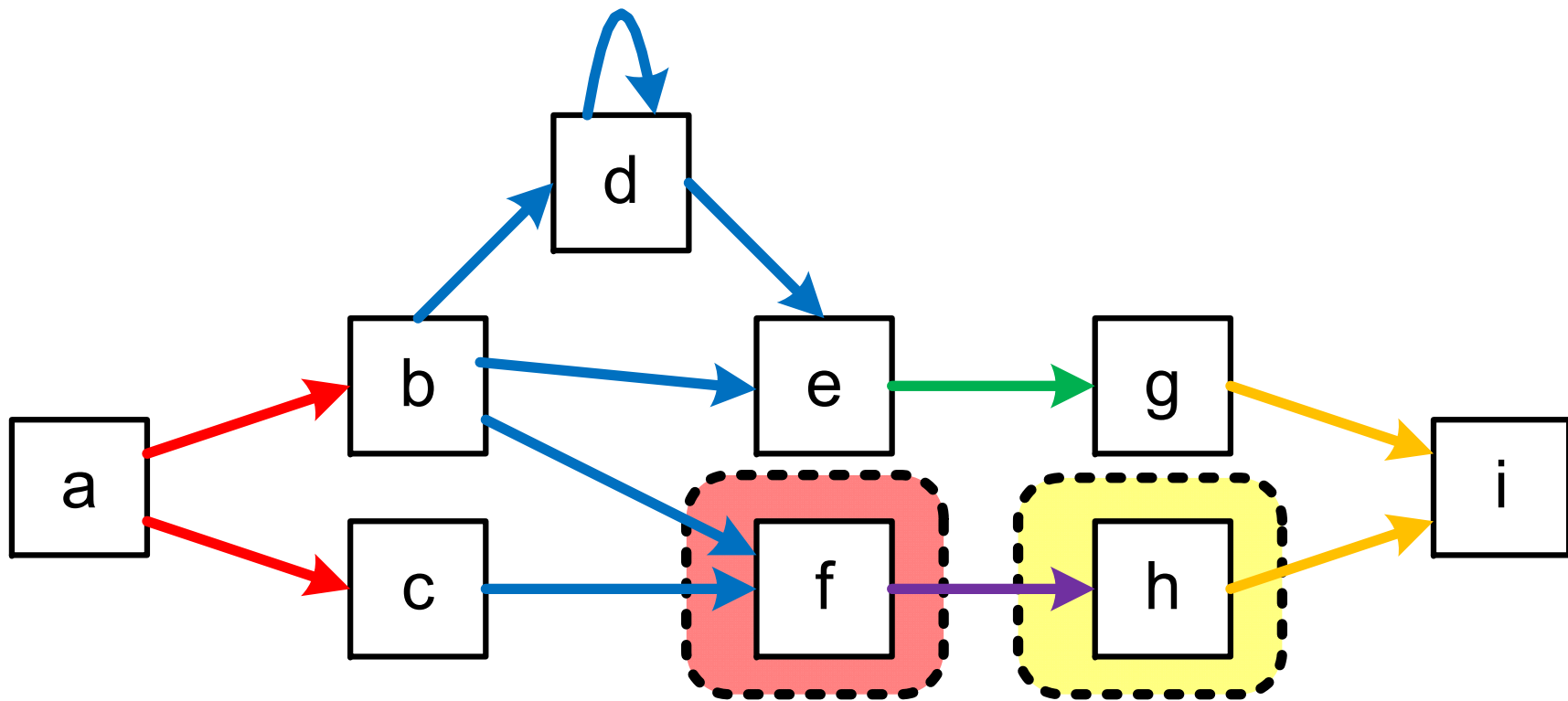
Minimal passage 2: ($\{b,c,d\},\{d,e,f\}$)



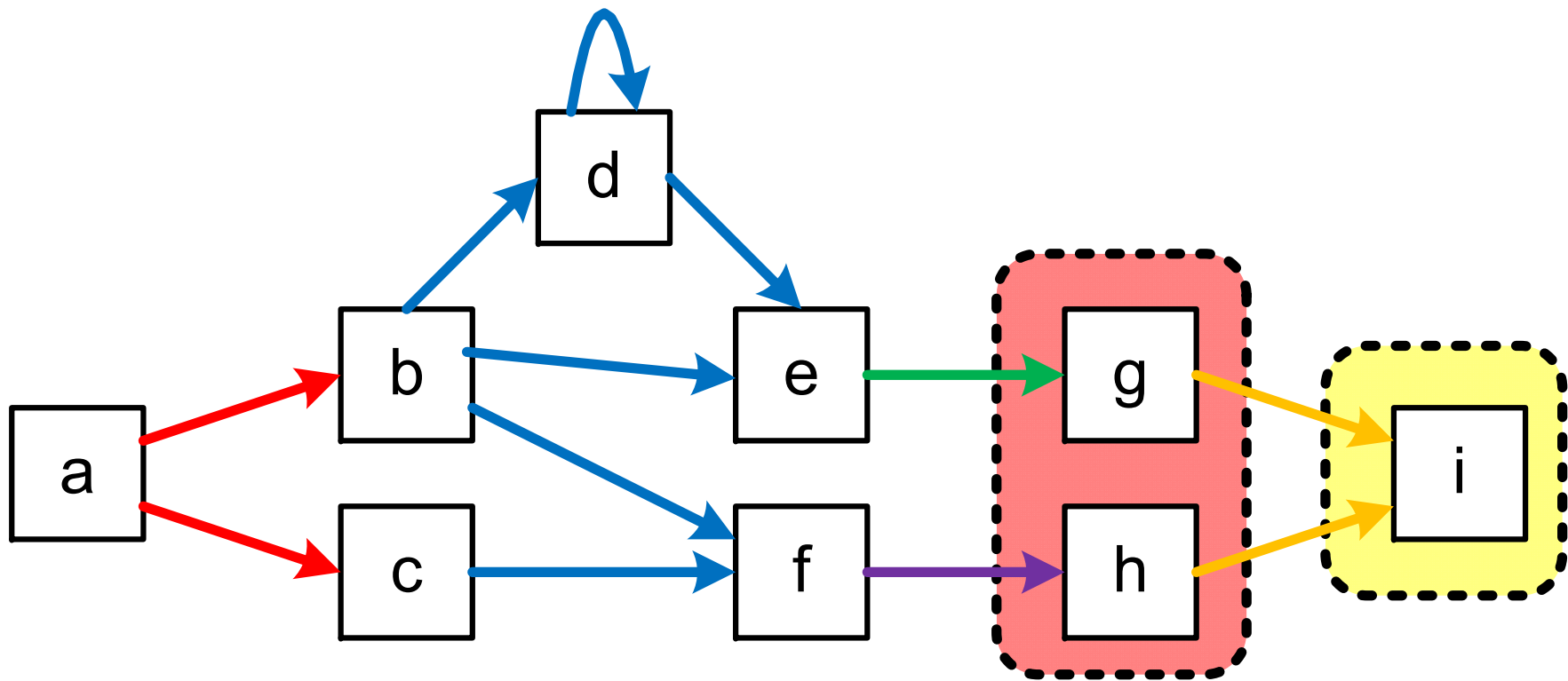
Minimal passage 3: ($\{e\}, \{g\}$)



Minimal passage 4: ($\{f\}, \{h\}$)

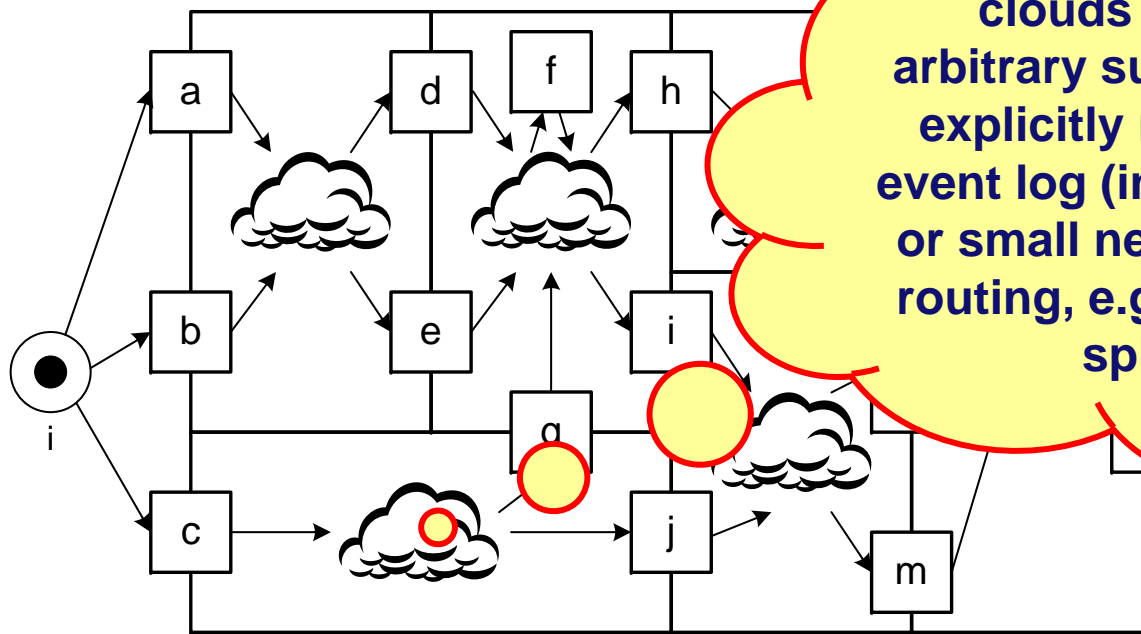


Minimal passage 5: ($\{g,h\},\{i\}$)



So What?

- Any process model can be partitioned in minimal passages.
- Discovery and conformance checking can be done per passage!



clouds may contain arbitrary subprocesses not explicitly recorded in the event log (invisible activities or small networks used for routing, e.g. XOR/AND/OR-split/joins)

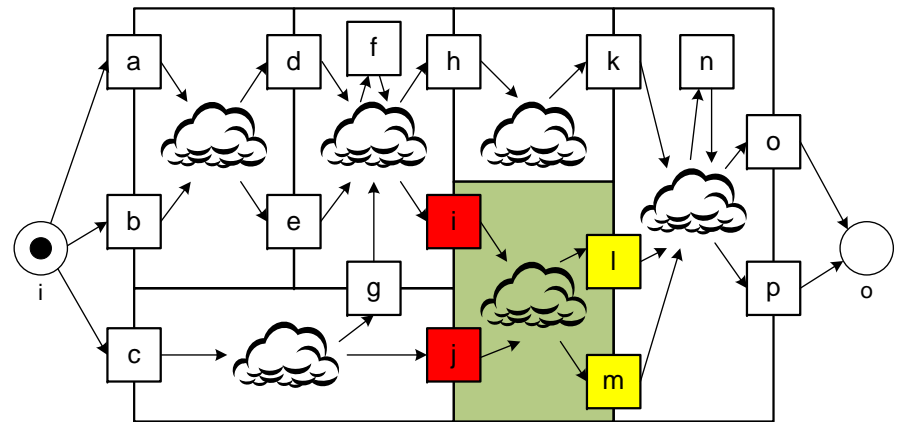
Example result for Petri nets

Theorem 1 (Main Theorem). Let $L \in \mathcal{B}(A^*)$ be an event log and let $WF = (PN, in, T_i, out, T_o)$ be a WF-net with $PN = (P, T, F, T_v)$.

L is perfectly fitting system net $SN = (PN, [in], [out])$ if and only if

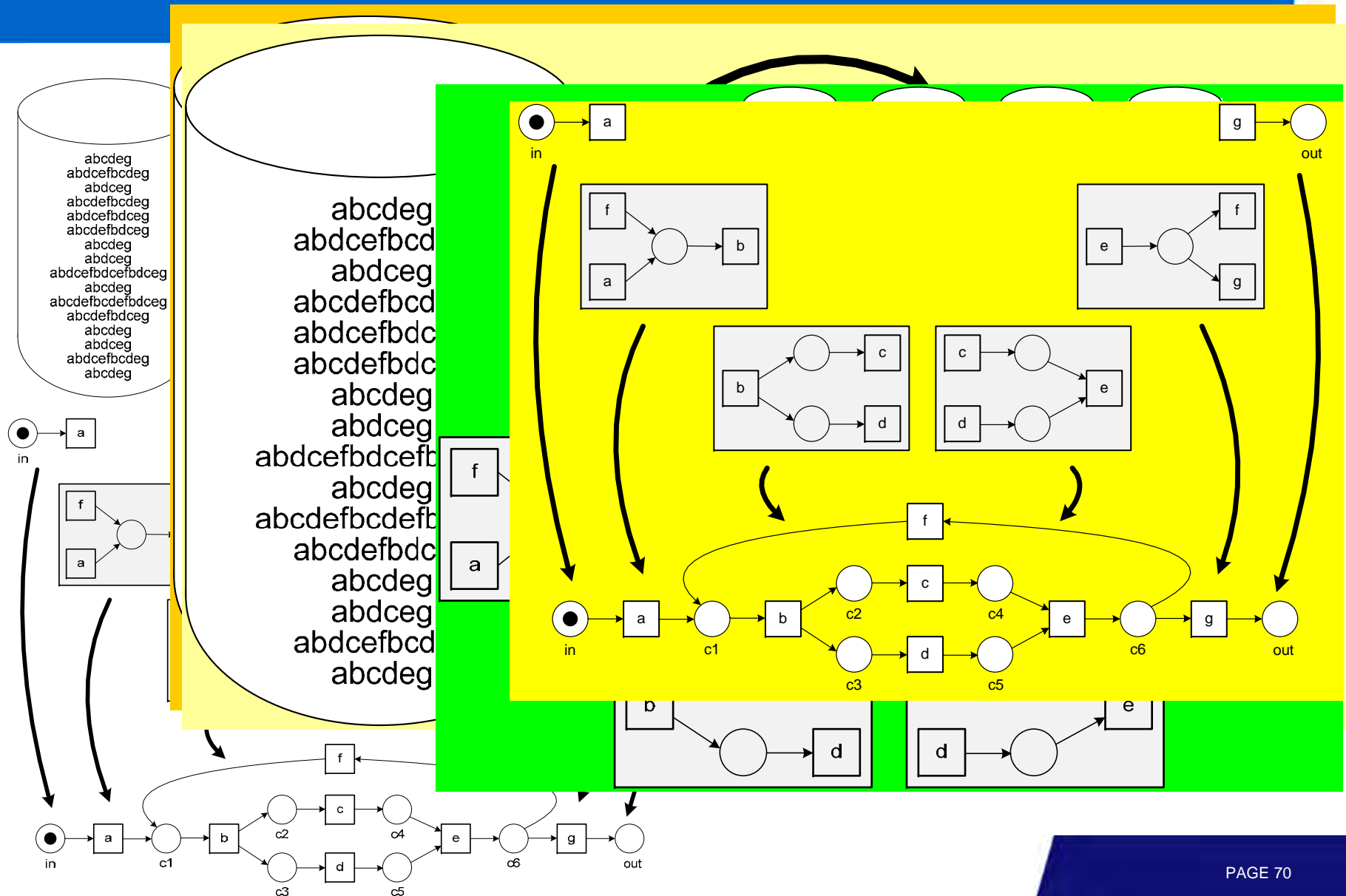
- for any $\langle a_1, a_2, \dots, a_k \rangle \in L$: $a_1 \in T_i$ and $a_k \in T_o$, and
- for any $(X, Y) \in pas_{min}(skel(PN))$: $L \upharpoonright_{X \cup Y}$ is perfectly fitting $SN^{(X,Y)} = (PN^{(X,Y)}, [], [])$.

“The event log fits all passages if and only if the event log fits the whole model.”

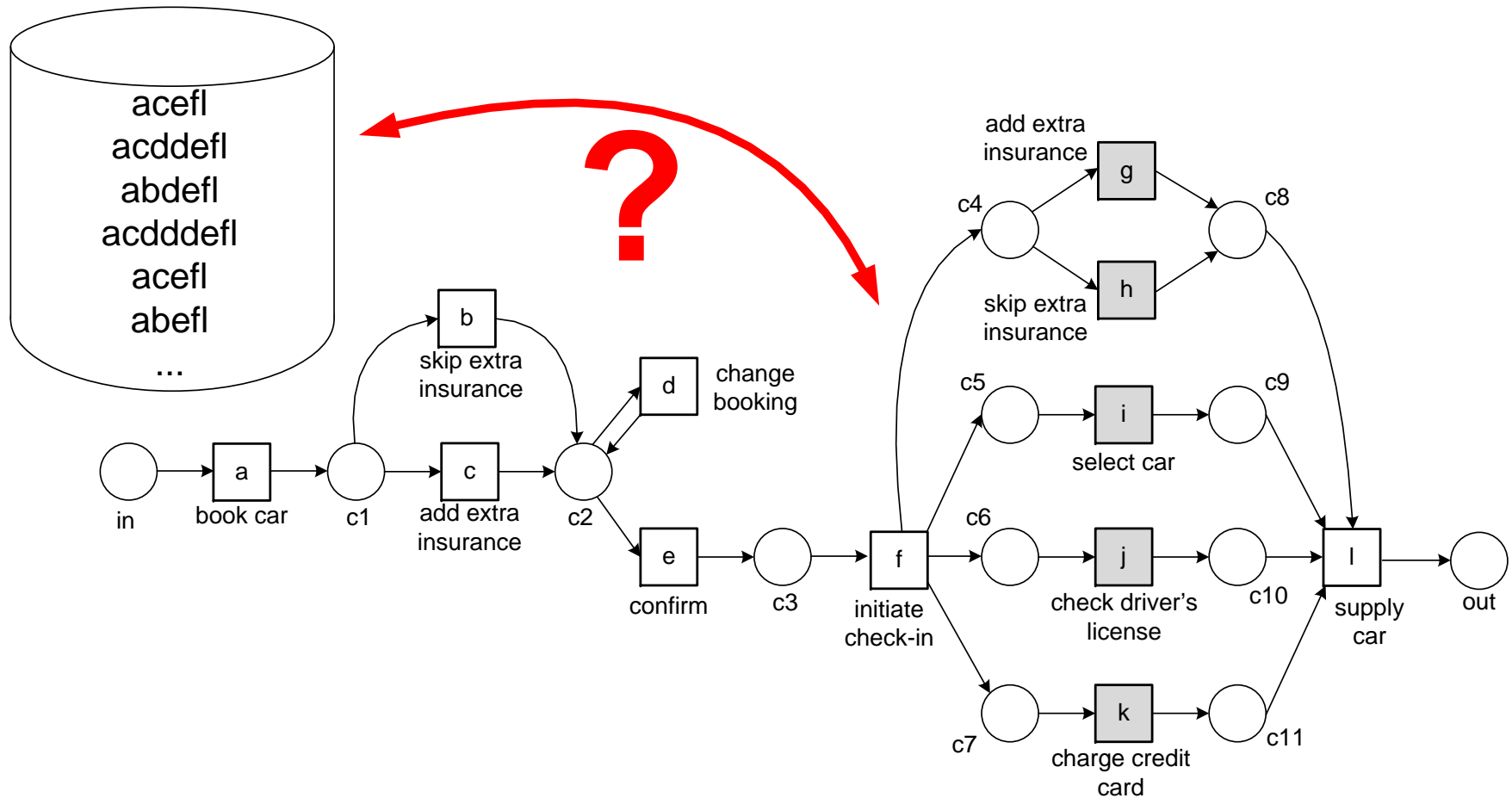


Key insight: interface transitions controlled by event log

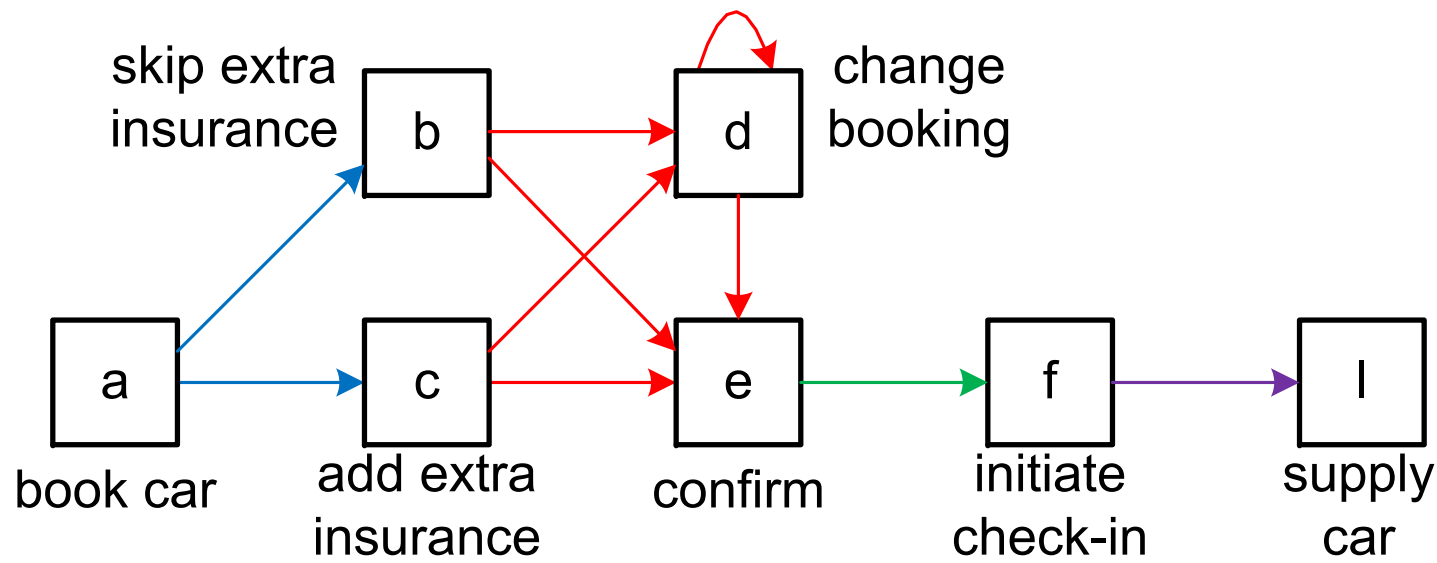
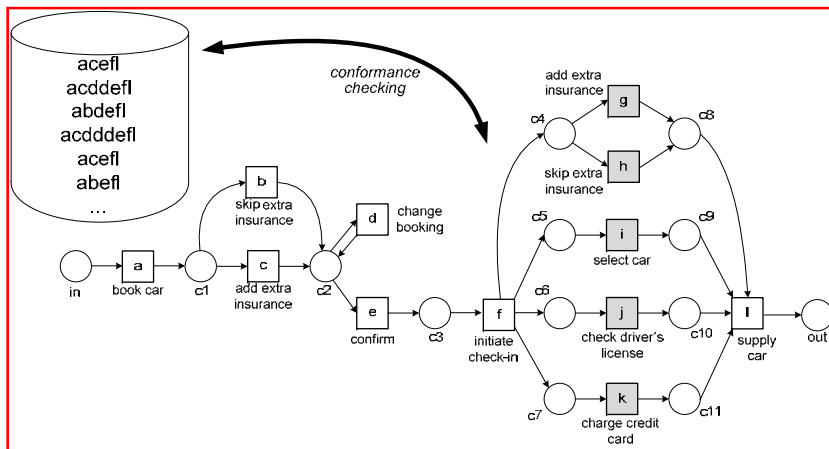
Discovery example



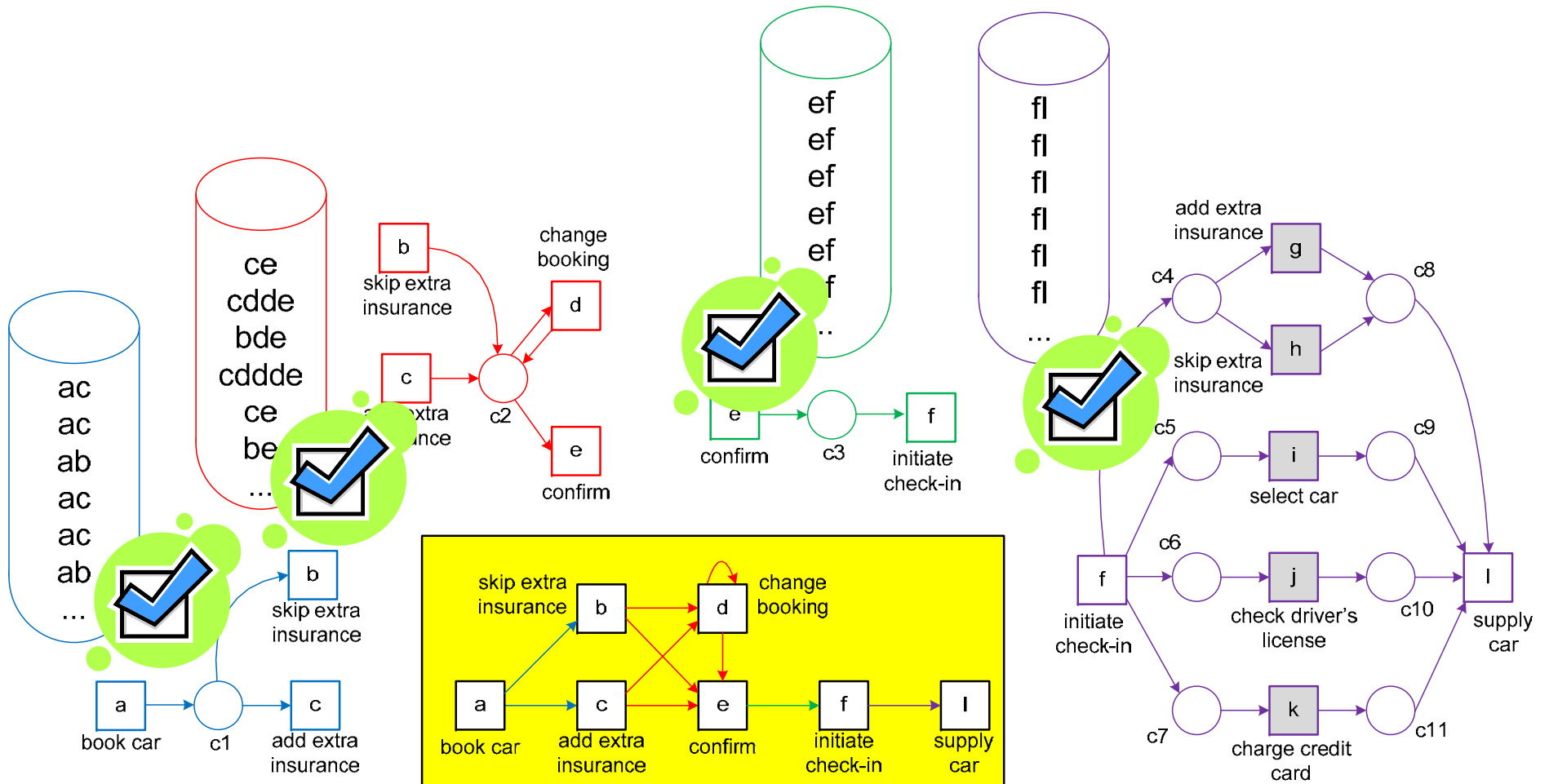
Conformance checking



Create Skeleton



Net fragments per passage



Initial implementation in ProM

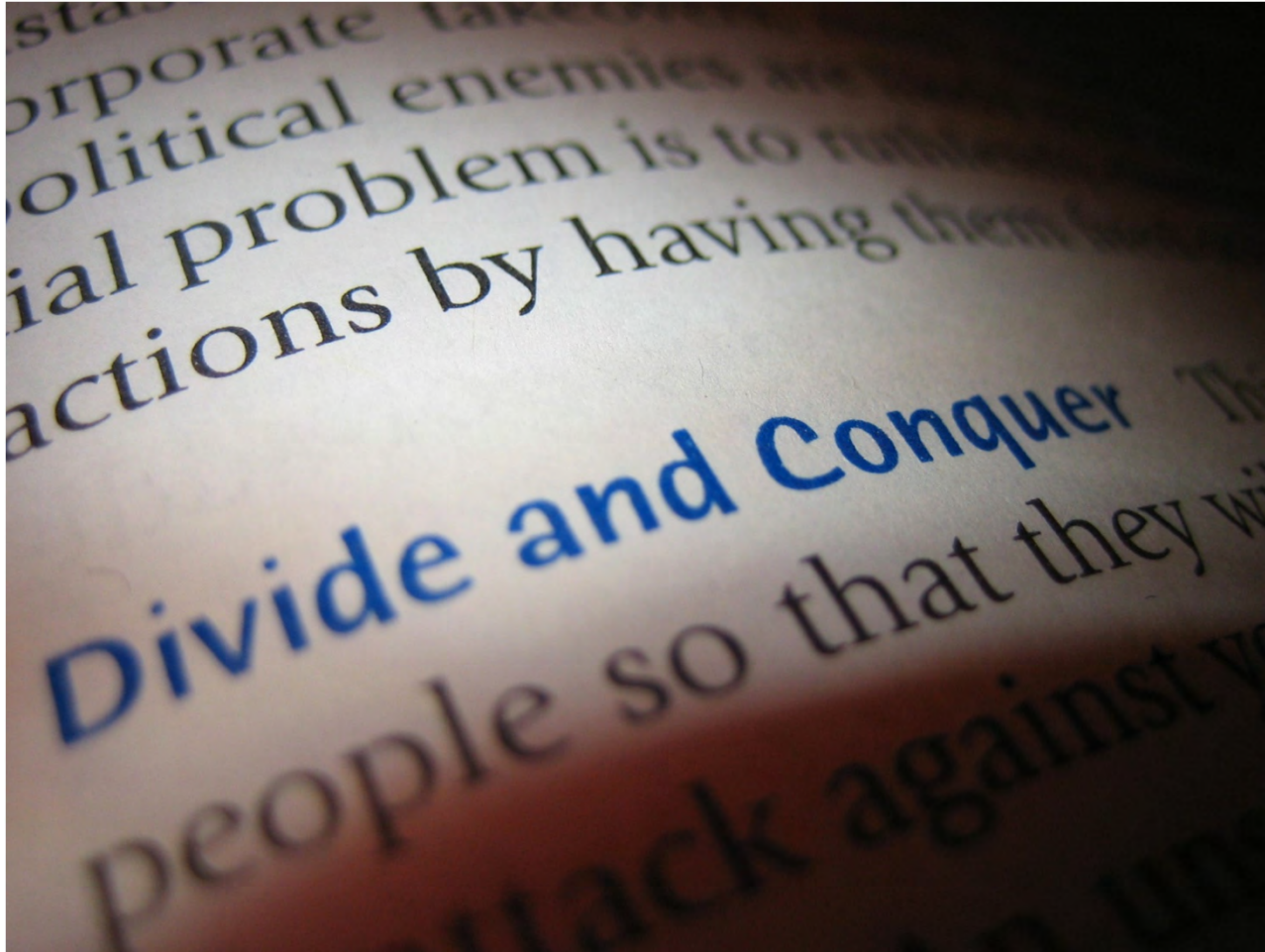
The screenshot displays the ProM 6.10.0 process mining workbench interface. The main window is titled "ProM 6.10.0" and features a toolbar with icons for file operations, execution, and visualization. The "Actions" panel on the left lists various analysis tasks, including:

- Add Artificial Events
- Add Missing Events
- Analyze Transition System
- Animate Event Log in Fuzzy Instance
- Calculate Log Meta Data
- Check Conformance using ETCConformance
- Compute Fitness
- Concept Drift
- D_IOMA
- Declare Miner
- Declare Miner Default
- Filter Log using Simple Heuristics

The "Output" panel on the right shows a dashed box with the text "Click to add output object". The main workspace displays a complex process flow diagram with nodes labeled "complete" and "Case13". A red arrow indicates a path through the flow. The bottom status bar shows the time "13:07:28".

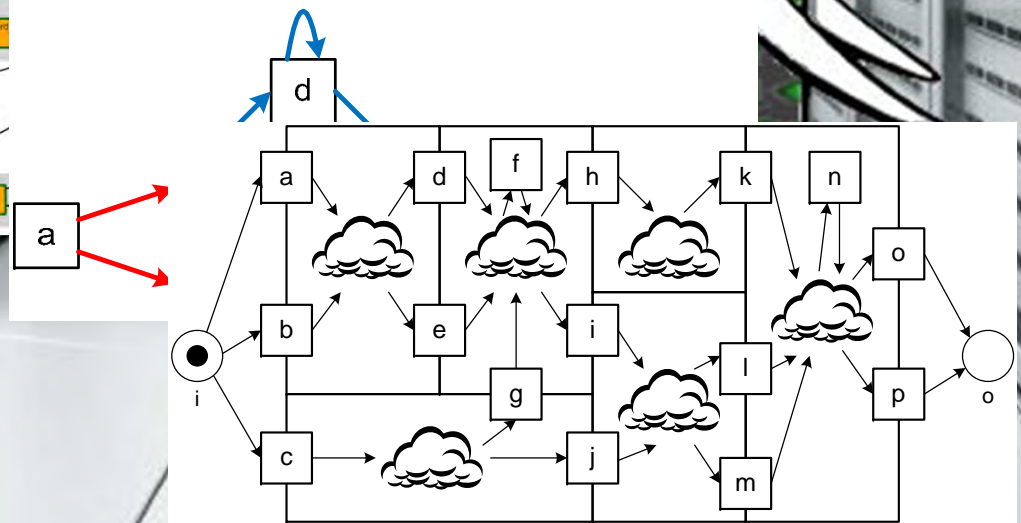
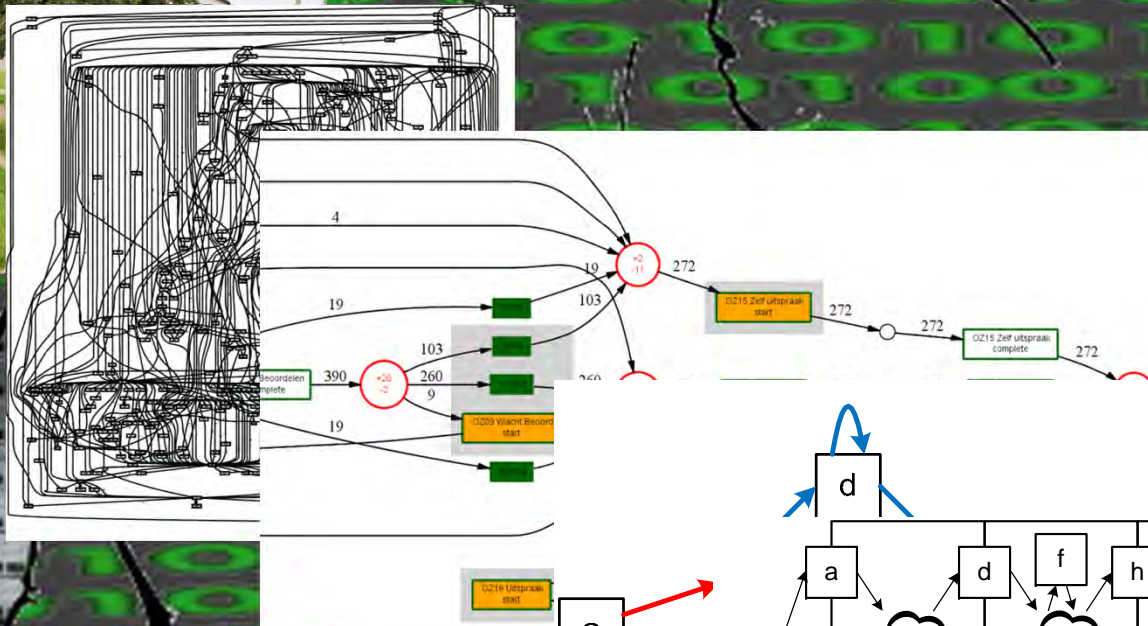
ProM
process mining workbench

Super linear speedups possible (even when using a single computer decomposition helps)



Conclusion

Conclusion



“Big Data”

Wil M. P. van der Aalst
Process Mining

Discovery, Conformance and Enhancement of Business Processes

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

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

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

Features and Benefits:

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

Computer Science

ISBN 978-3-642-19344-6



► springer.com

van der Aalst



Process Mining

Wil M. P. van der Aalst

Process Mining

Discovery, Conformance and
Enhancement of Business Processes

www.processmining.org

www.win.tue.nl/ieeetfpm/

 Springer