# Passages in Big Data

Partitioning Event Logs and Process Models to
Speed Up Process Mining Algorithms

**Wil van der Aalst**

*www.processmining.org*
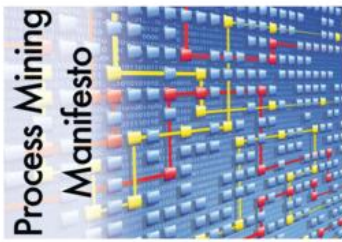
**TU/e** Technische Universiteit
**Eindhoven**
University of Technology

**Where innovation starts**

# Advances in Process Mining

- **Many process discovery and conformance checking algorithms and tools are available (cf. the various ProM packages).**

- **Also commercial software based on these ideas:**
  **Disco (Fluxicon), Reflect (Futura), BPMOne (Pallas Athena/Perceptive), ARIS Process Performance Manager (Software AG), Futura Reflect (Futura Technology), Interstage Automated Process Discovery (Fujitsu), QPR ProcessAnalyzer/Analysis (QPR Software), flow (fourspark), Discovery Analyst (StereoLOGIC), etc.**

- **We applied process mining in over 100 organizations.**

**More than 75 people involving more than 50 organizations created the Process Mining Manifesto in the context of the IEEE Task Force on Process Mining.**
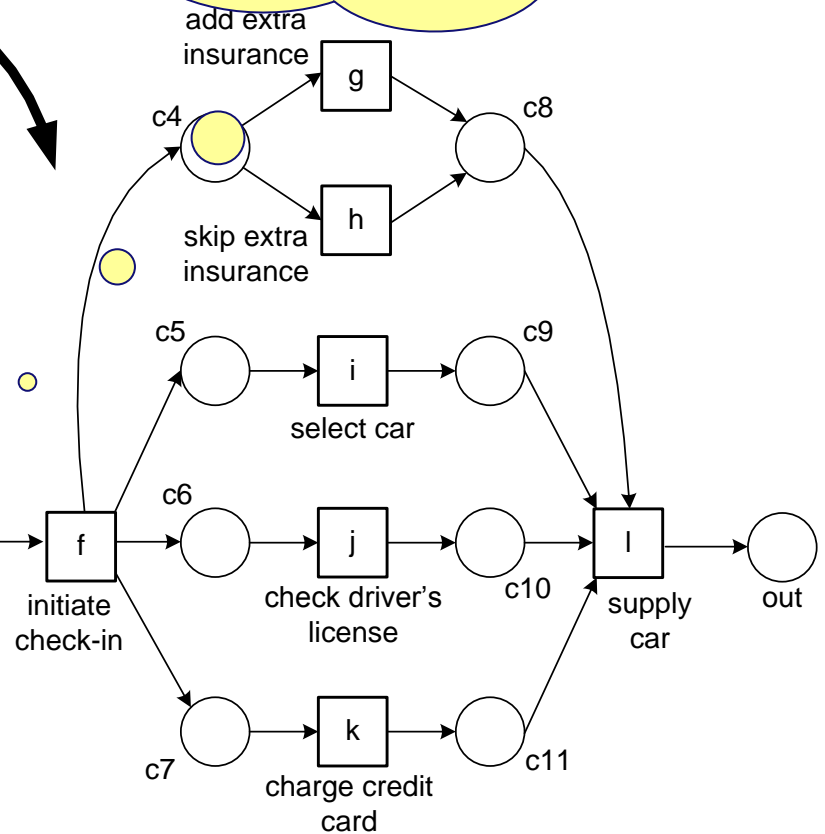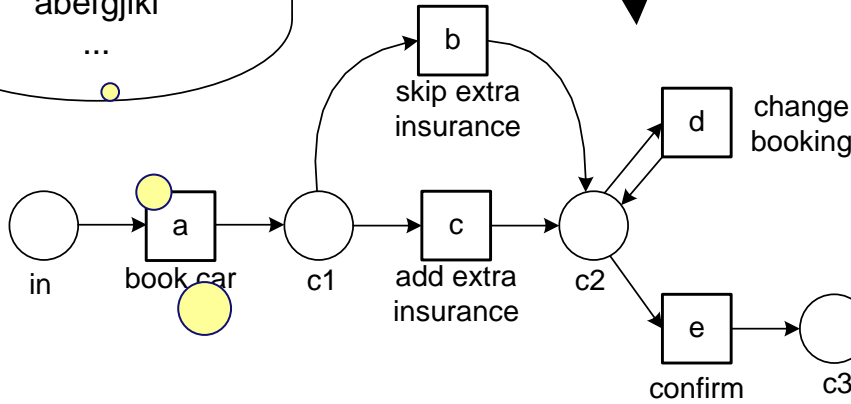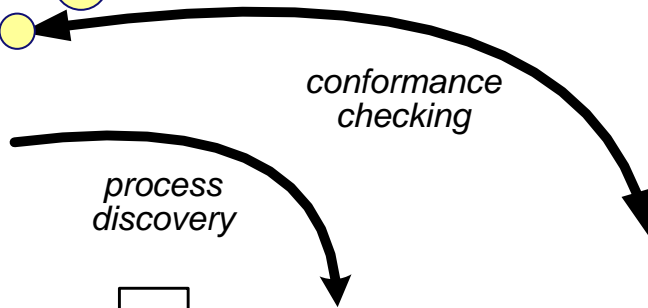
**Available in 13 languages**

killer app for Petri nets!

# Big Data: Opportunities and Challenges

# Distributed Computing

- **multicore CPU**
- **manycore GPU**
- **cluster computing**
- **grid computing**
- **cloud computing**
- **…**

# How to distribute process discovery?



abcdeg
abdcefbcdeg
abdceg
abcdefbcdeg
abdcefbdceg
abcdefbdceg
abcdeg
abdceg
abdcefbdcefbdceg
abcdeg
abcdefbcdefbdceg
abcdefbdceg
abcdeg
abdceg
abdcefbcdeg
abcdeg

**?**

# How to distribute process discovery?

# How to distribute conformance checking?

# How to distribute conformance checking?

# Replication: Same event log on all computing nodes



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

# Classification based on partitioning of event log: vertical and horizontal



sets of cases

sets of activities

# Vertical distribution I: Split cases arbitrarily



sets of cases

abcdeg
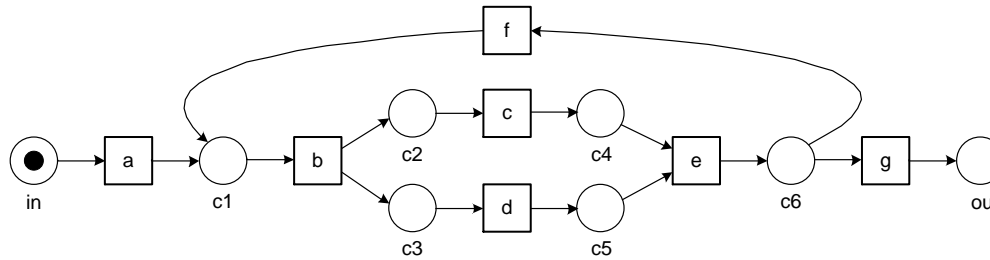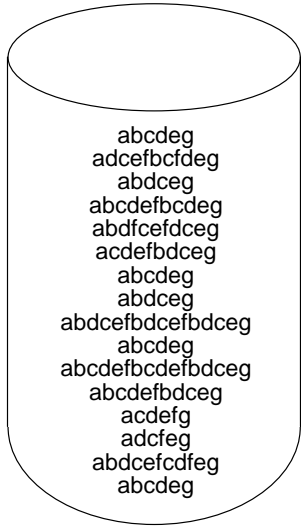abdcefbcdeg
abdceg
abcdefbcdeg
abdcefbdceg
abcdefbdceg
abcdeg
abdceg
abdcefbdcefbdceg
abcdeg
abcdefbcdefbdceg
abcdefbdceg
abcdeg
abdceg
abdcefbcdeg
abcdeg

abcdeg
abdcefbcdeg
abdceg
abcdefbcdeg
abdcefbdceg
abcdefbdceg
abcdeg
abdceg
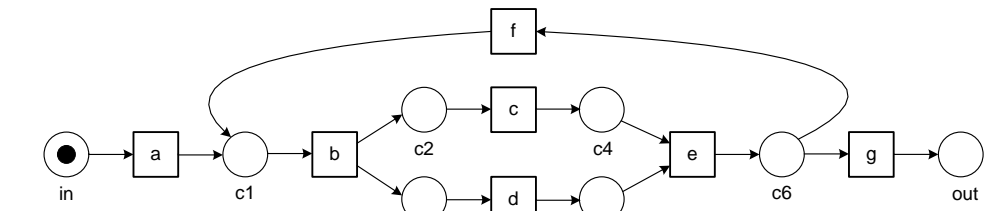
abdcefbdcefbdceg
abcdeg
abcdefbcdefbdceg
abcdefbdceg
abcdeg
abdceg
abdcefbcdeg
abcdeg

# Vertical distribution II:
# Split cases based on a specific feature

# Horizontal distribution

| abcdeg | abeg | bcde |
| abdcefbcdeg | abefbeg | bdcebcde |
| abdceg | abeg | bdce |
| abcdefbcdeg | abefbeg | bcdebcde |
| abdcefbdceg | abefbeg | bdcebdce |
| abcdefbdceg | abefbeg | bdcebdce |
| abcdeg | abeg | bcde |
| | | bdce |
| | | ebdcebdce |
| | | bcde |
| | | ebcdebdce |
| | | cdebdce |
| | | bcde |
| | | bdce |
| | | dcebcde |
| | | bcde |

$$abcdeg = abeg + bcde$$

# Horizontal distribution: The key idea



**projected on a,b,e,f,g**

abeg
abefbeg
abeg
abefbeg
abefbeg
abefbeg
abeg
abeg
abefbefbeg
abeg
abefbefbeg
abefbeg
abeg
abeg
abefbeg
abeg

**projected on b,c,d,e**

bcde
bdcebcde
bdce
bcdebcde
bdcebdce
bcdebdce
bcde
bdce
bdcebdcebdce
bcde
bcdebcdebdce
bcdebdce
bcde
bdce
bdcebcde
bcde

Passages

Moscow GUM

# Passage *P=(X,Y)*

causal dependency: may trigger or enable

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

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

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

$$X = \overset{G}{\bullet} Y$$
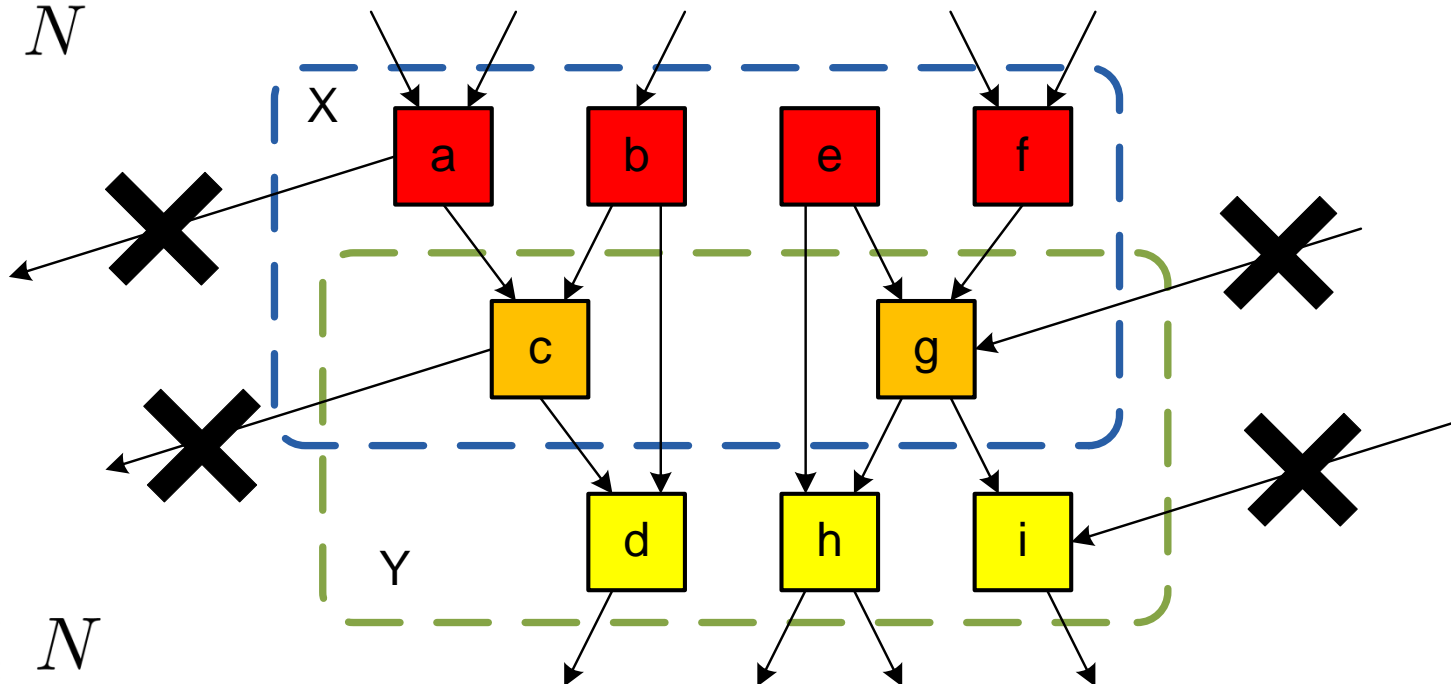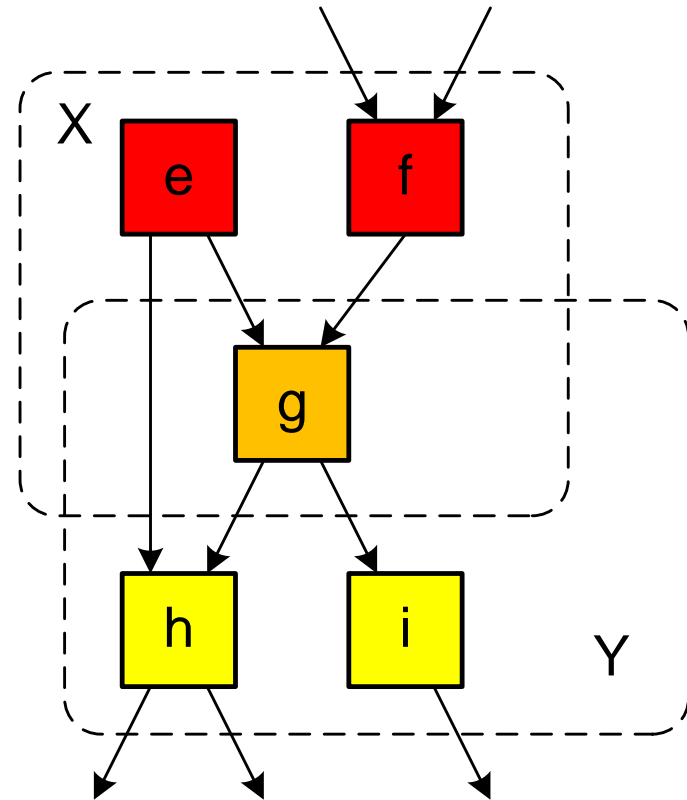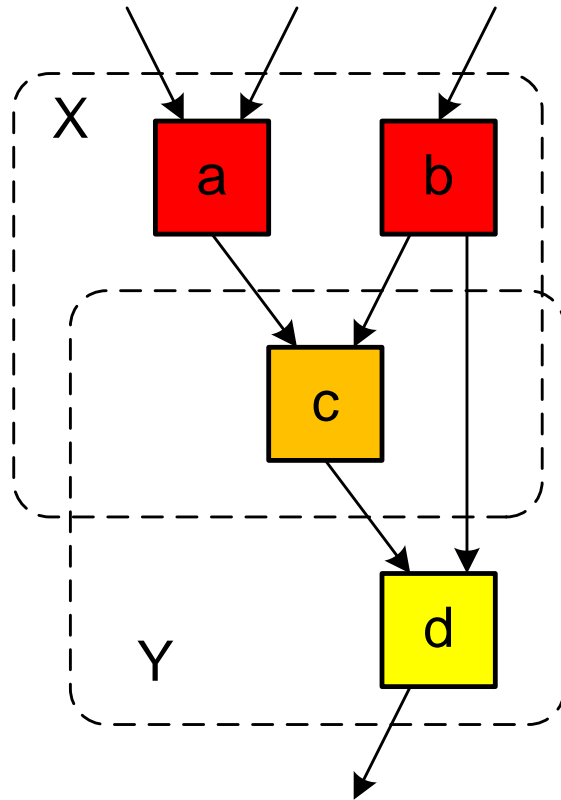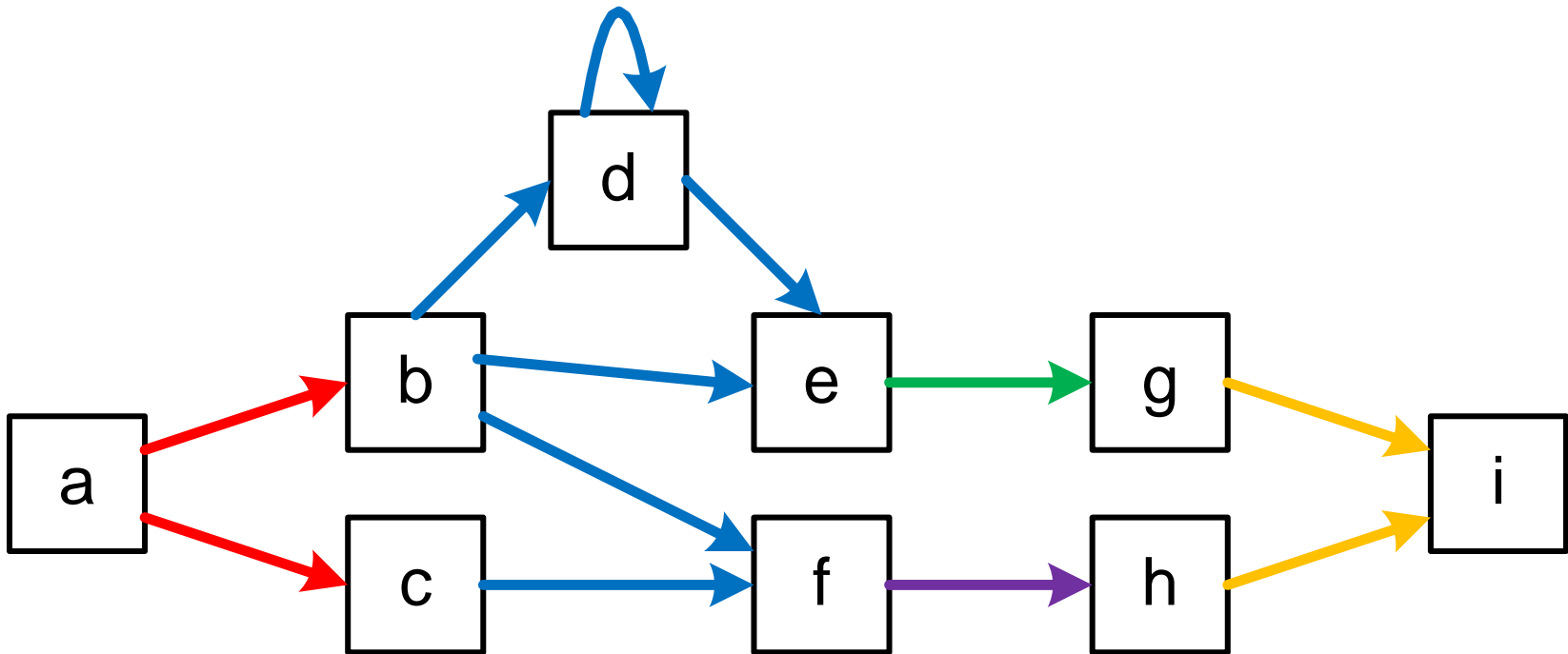
# Minimal passages
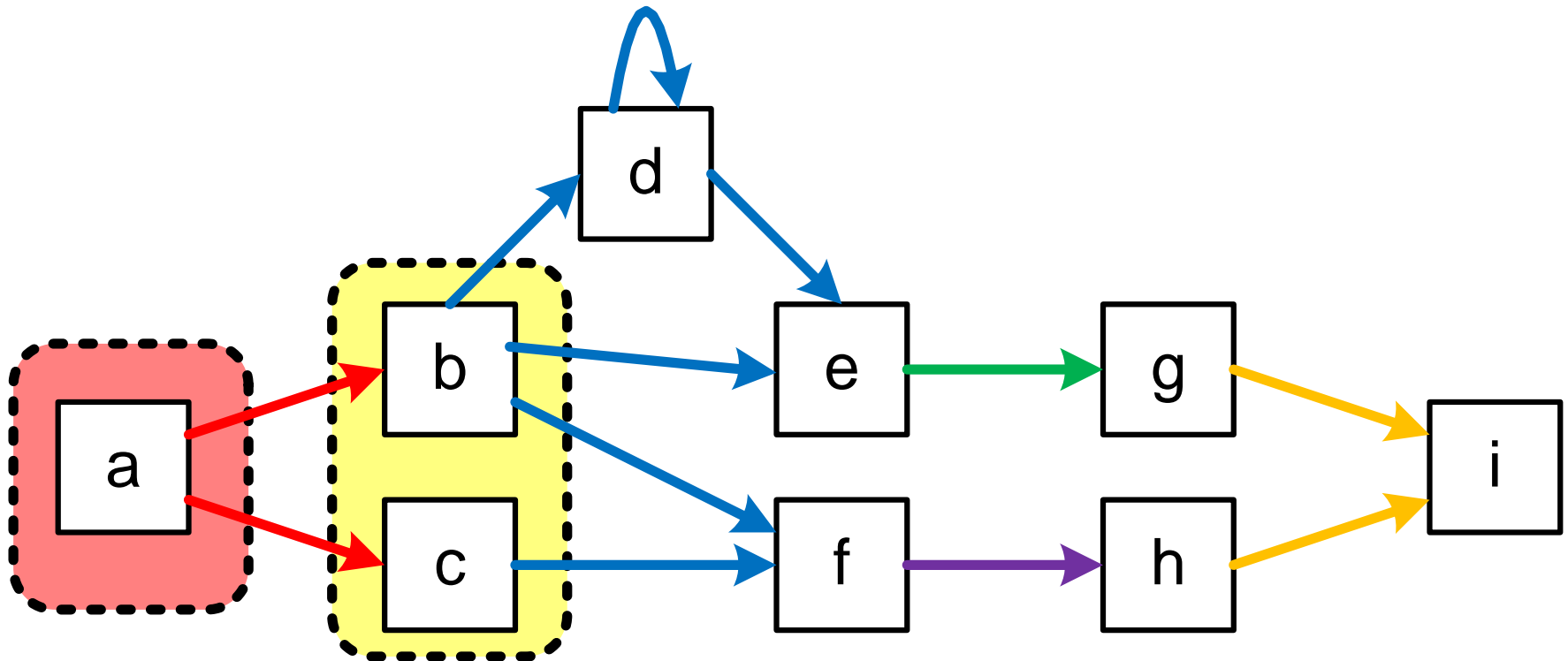


$$X \overset{G}{\bullet} = Y$$

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

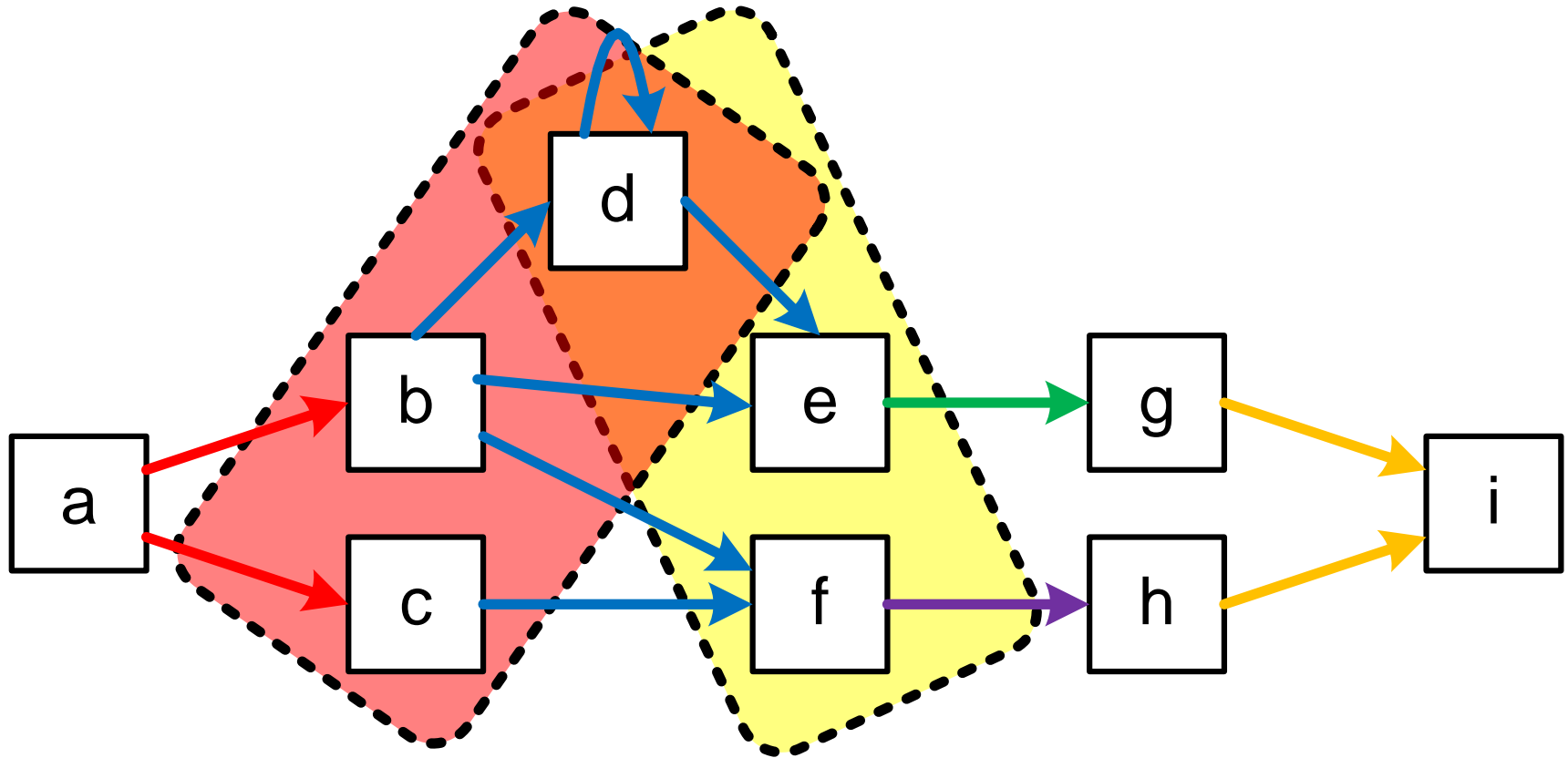**a passage is minimal if it does not contain smaller passages**

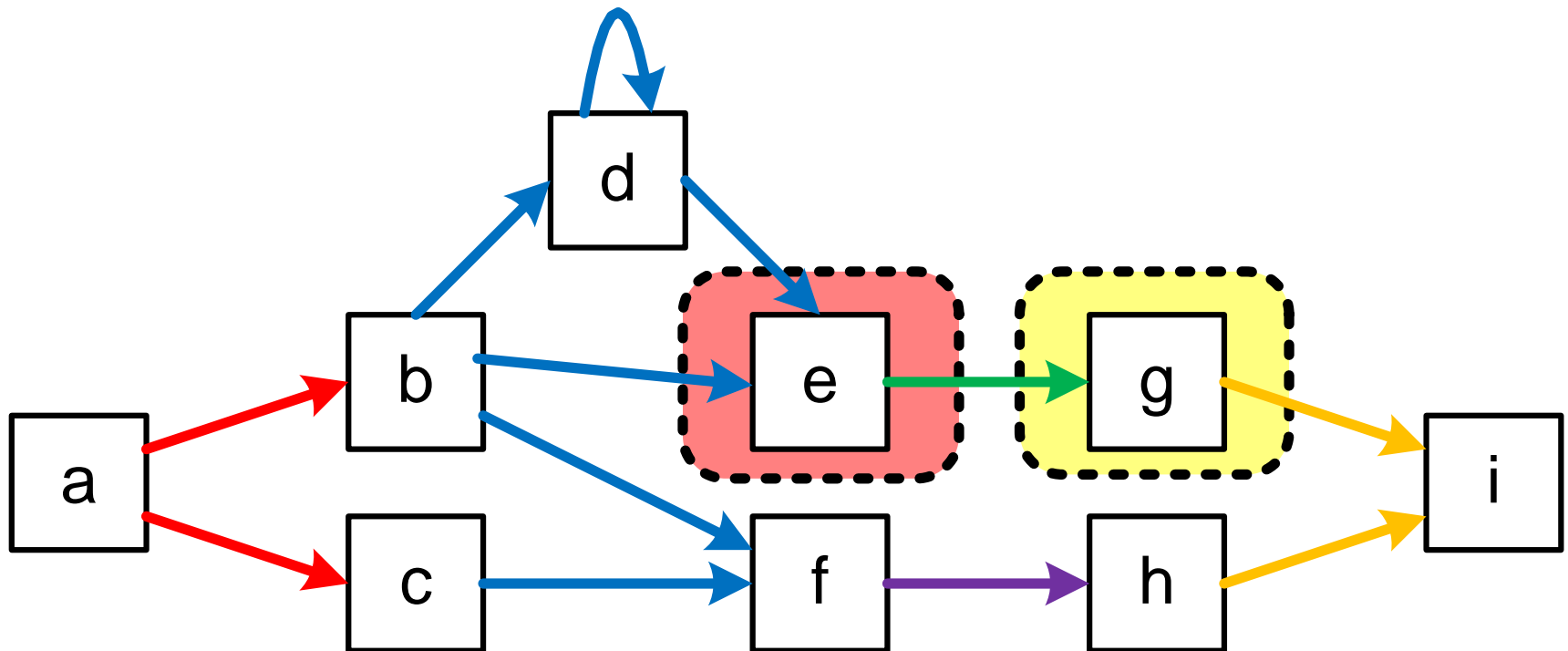# Passages define an equivalence relation on the edges in the graph

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

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

# Minimal passage 4: ({f},{h})

# So What?

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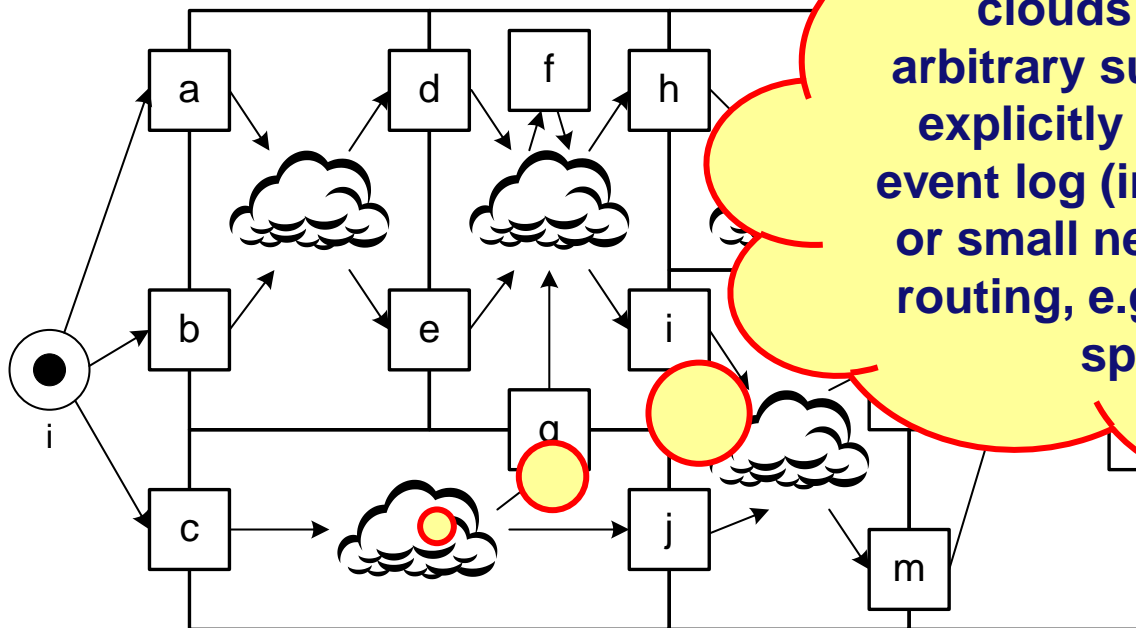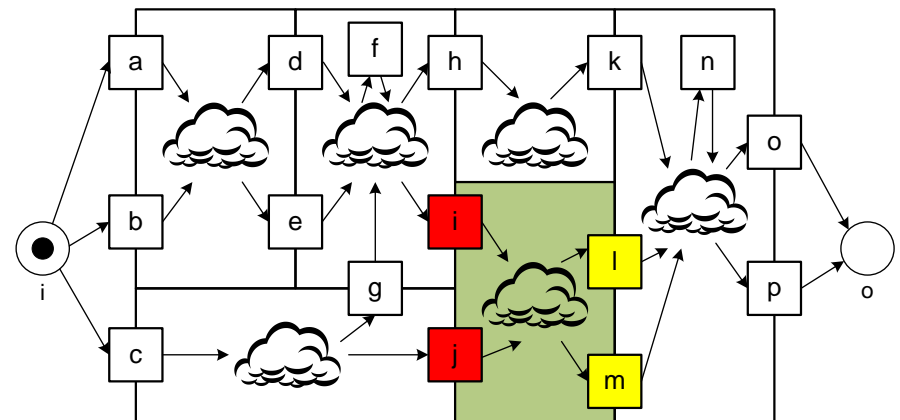
**Theorem 1** (**Main Theorem**). *Let* $L \in \mathcal{B}(A^*)$ *be an event log and let* $WF = (PN, in, T_i, out, T_o)$ *be a WF-net with* $PN = (P, T, F, T_v)$.

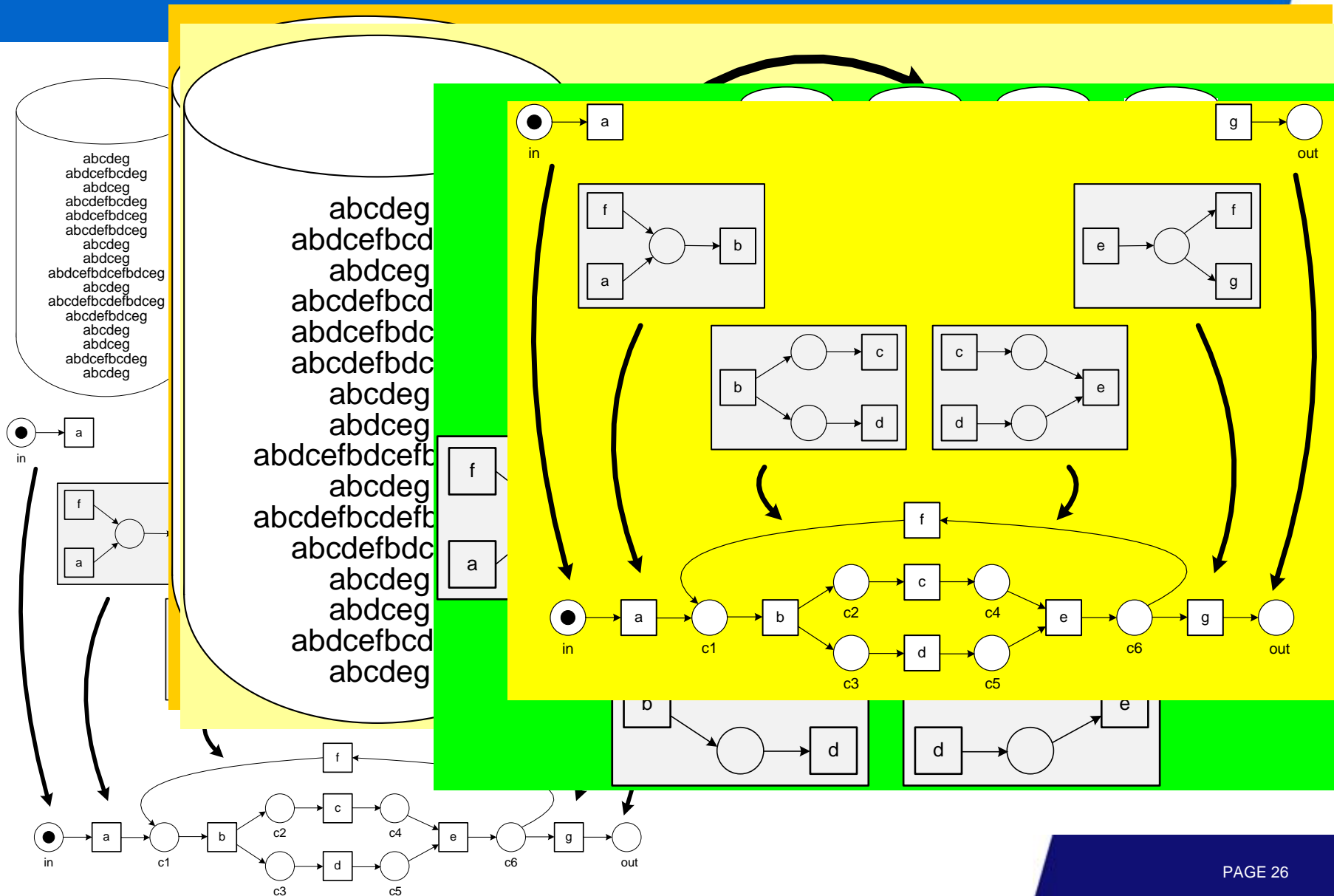*$L$ is perfectly fitting system net* $SN = (PN, [in], [out])$ *if and only if*

- *for any* $\langle a_1, a_2, \ldots a_k \rangle \in L$: $a_1 \in T_i$ *and* $a_k \in T_o$, *and*
- *for any* $(X, Y) \in pas_{min}(skel(PN))$: $L \restriction_{X \cup Y}$ *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."**
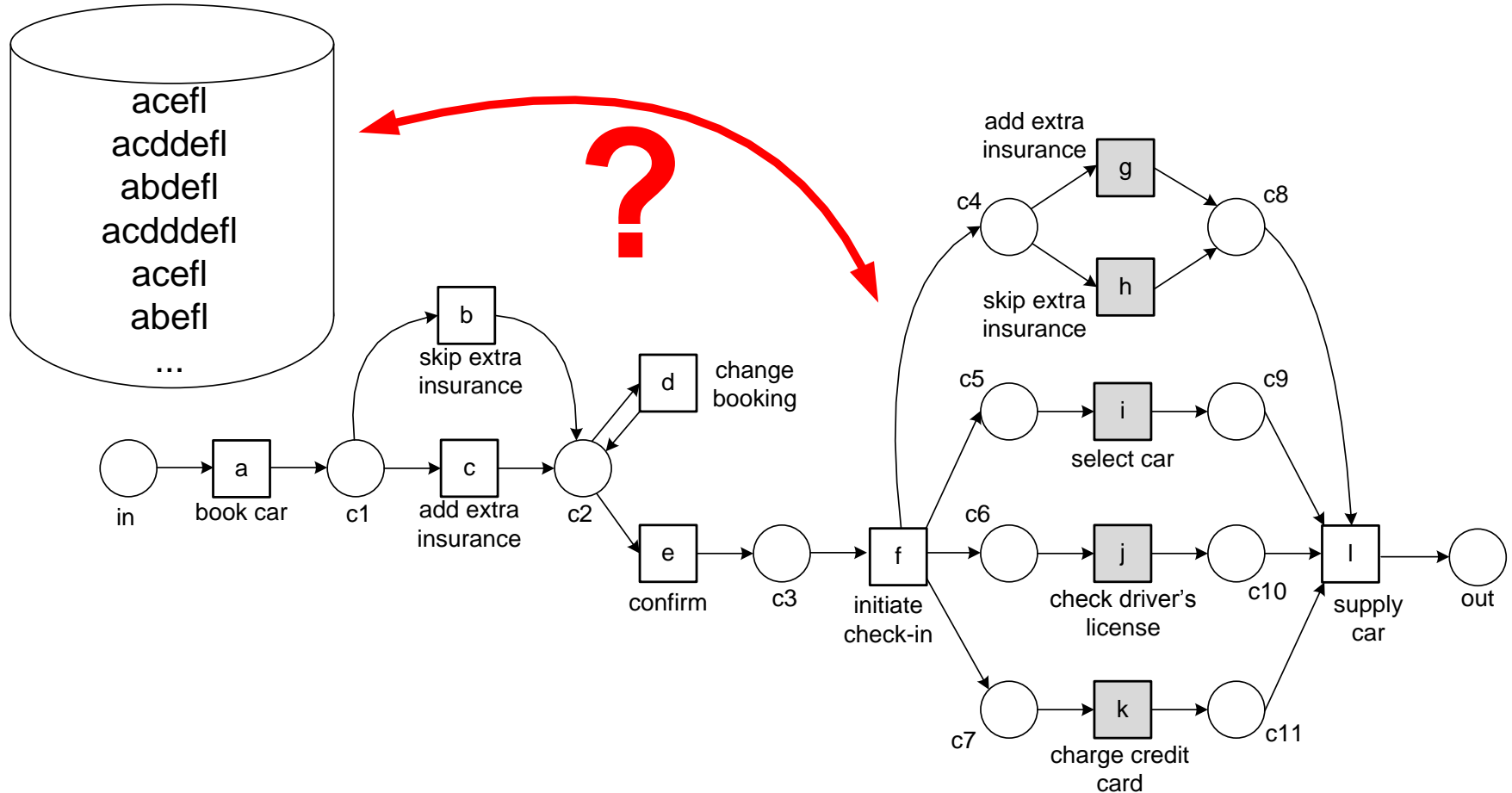


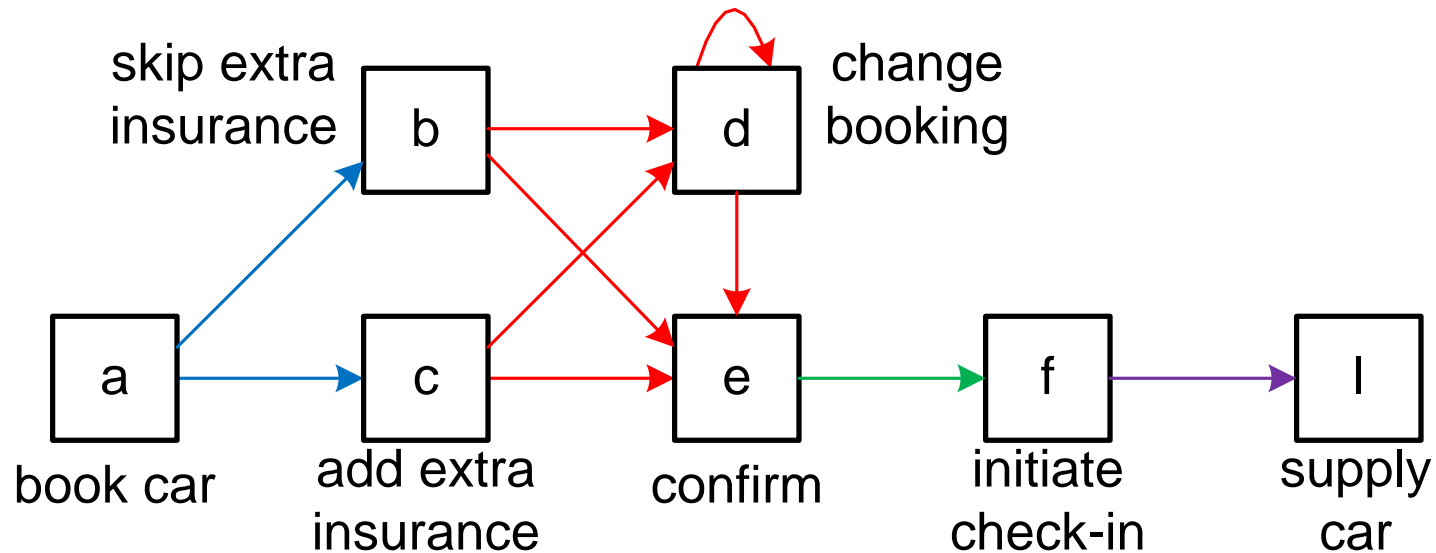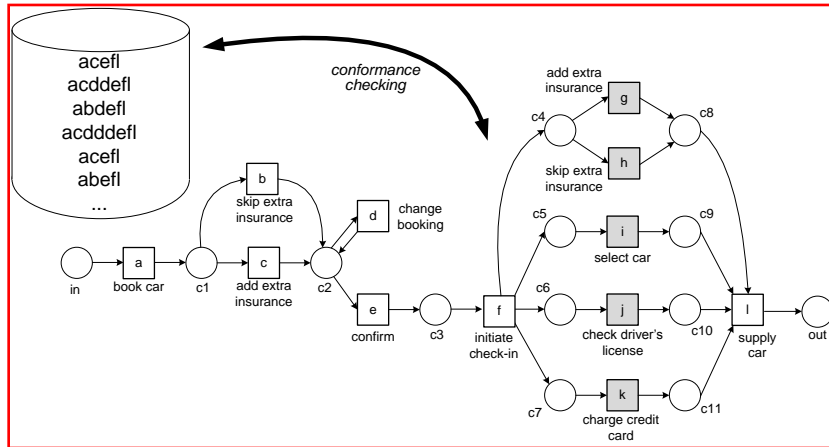**Key insight: interface transitions controlled by event log**
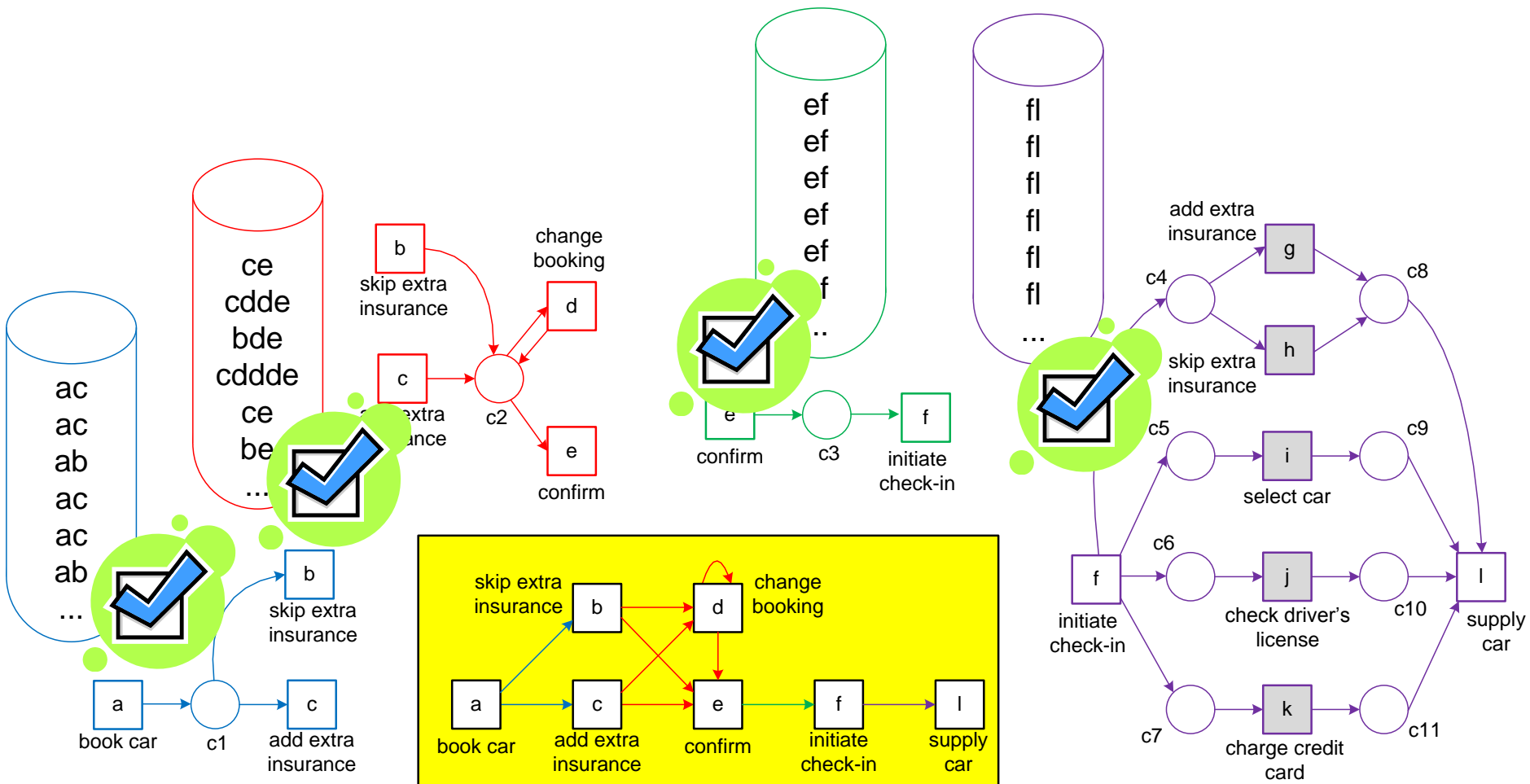
# Discovery example
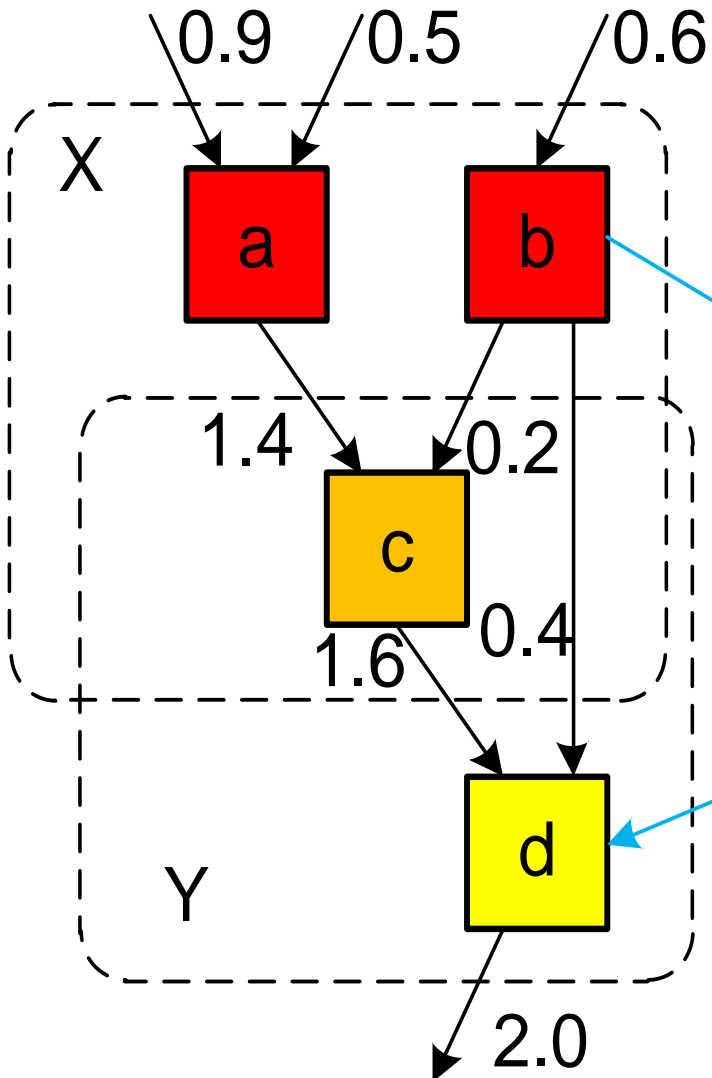
# Conformance checking

# Create Skeleton

# Limitations

- **Need to discover causal dependencies first (only issue for discovery, use fuzzy/heuristic rules).**
- **"Interface transitions" need to have a unique label.**
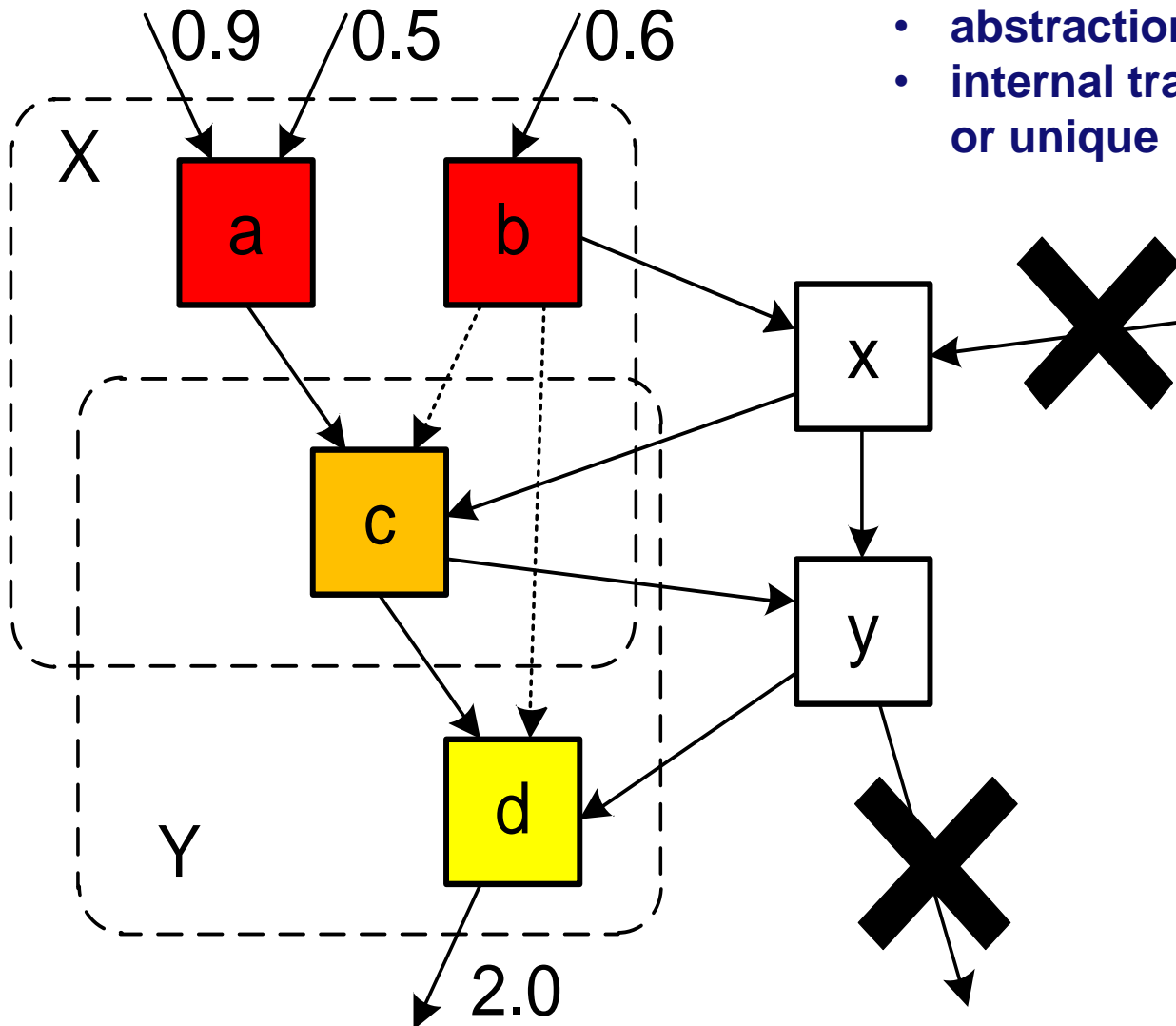- **Minimal passages may be large in dense graphs.**

# "Almost passages"



- **balance between size and quality**
- **discard same arcs in all passages (use an iterative procedure)**
- **also adding arcs? (don't think so)**

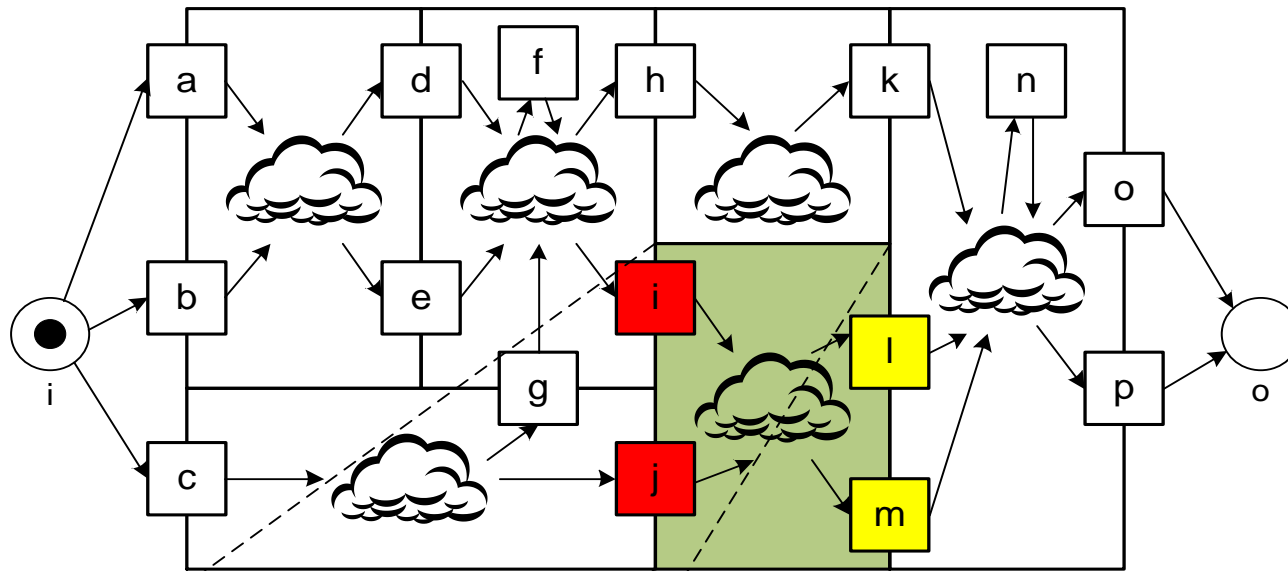$$1 - \frac{0.1 + 0.1}{0.9 + 0.5 + 0.6 + 0.1 + 1.4 + 0.2 + 1.6 + 0.4 + 0.1 + 2.0}$$

**Goal: small passages with only low-frequent arcs violating the rule.**

# Extended passages



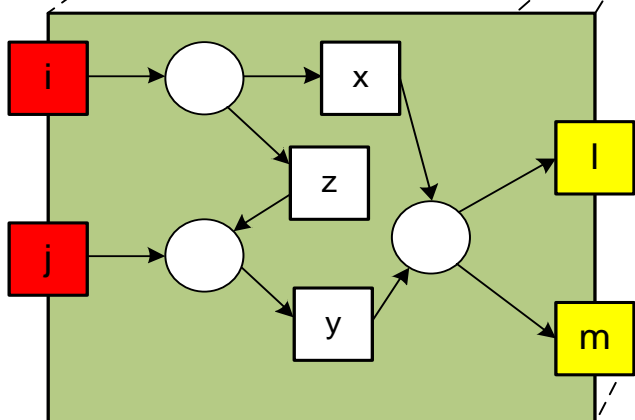- **abstraction yields passage**
- **internal transitions may be silent or unique**
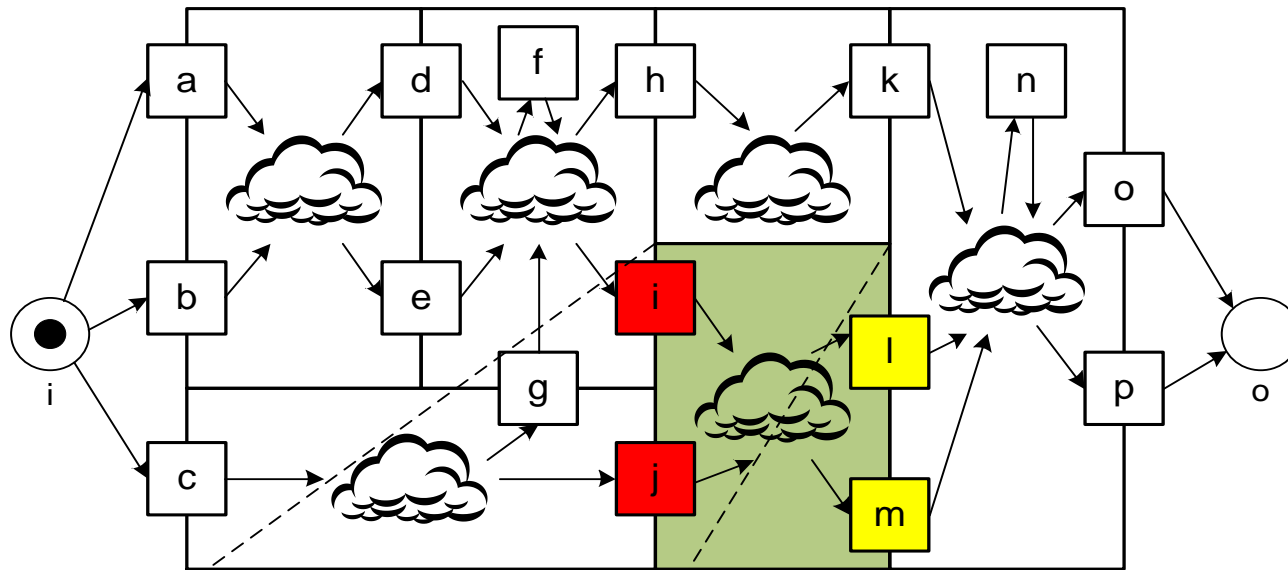
# Conformance checking (with silent steps inside passages)

# Aggregating conformance metrics
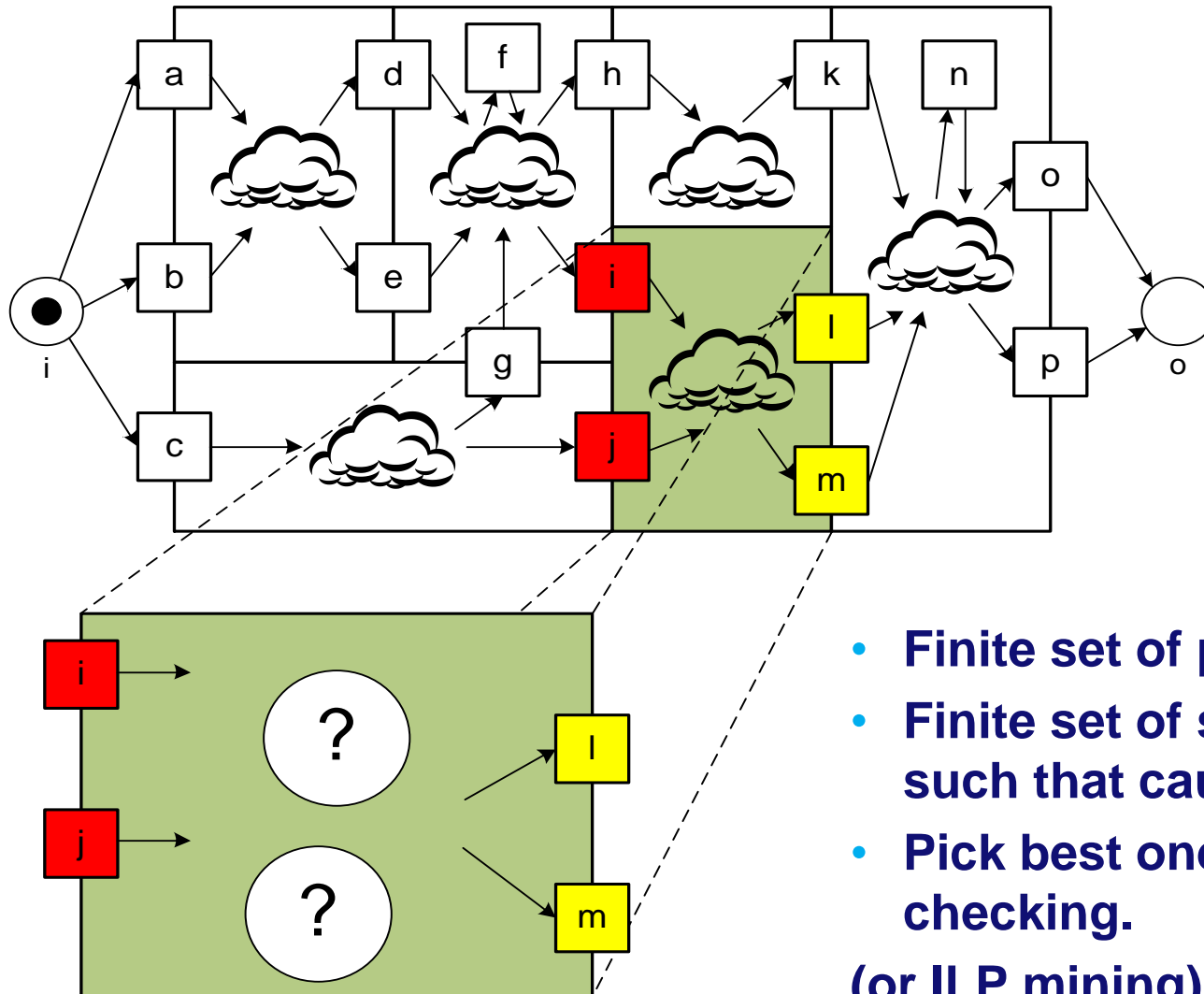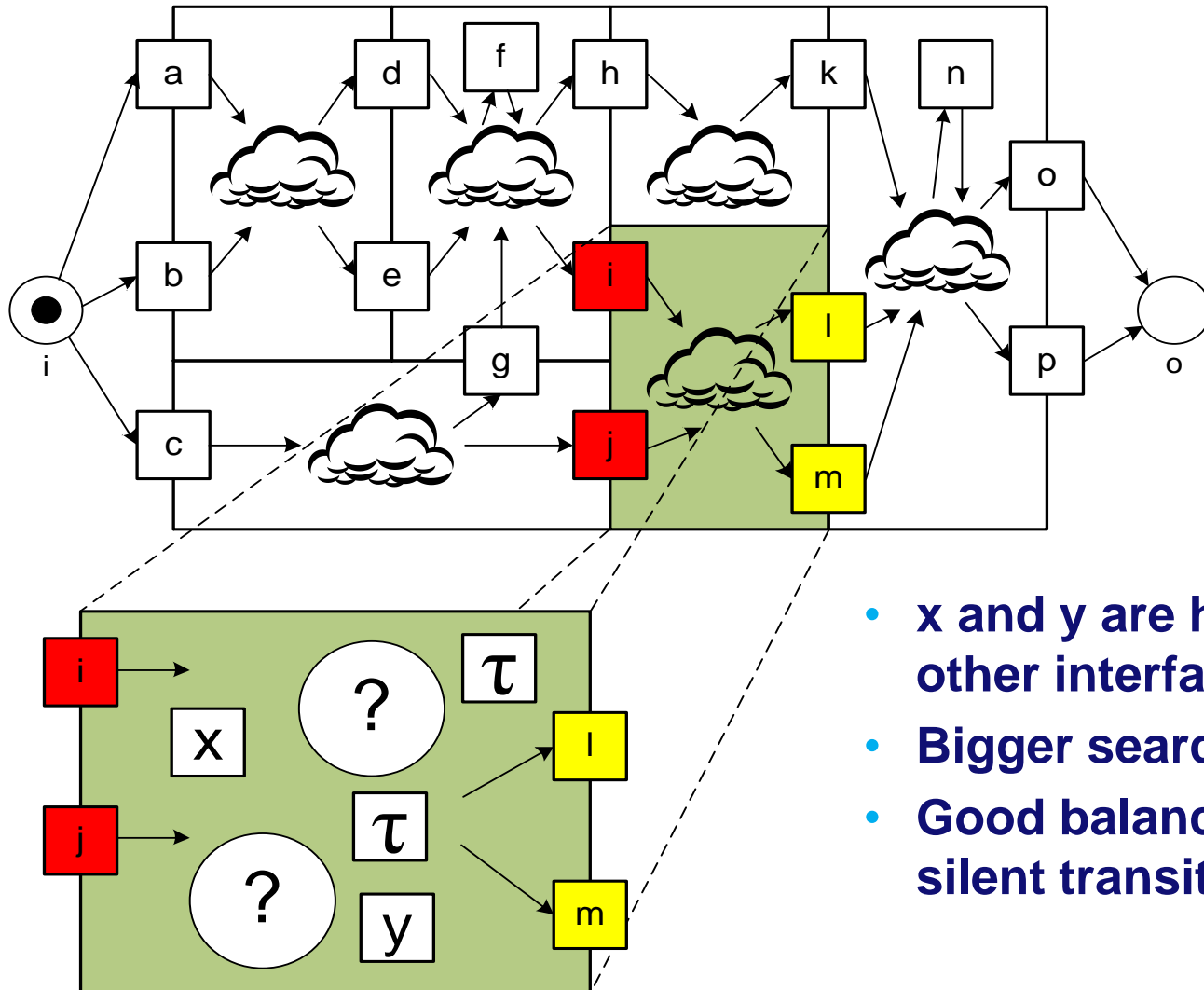


- **For all four conformance dimensions (fitness, simplicity, generalization, and precision)?**
- **Overall metrics: aggregated local values should be close to global values (e.g., computed fitness value is a lower bound).**
- **Local diagnostics (problem spots).**

# Discovery (no silent/internal transitions)



- **Finite set of possible places.**
- **Finite set of subsets of places such that causalities hold.**
- **Pick best one using conformance checking.**

**(or ILP mining)**

# Discovery (with silent or internal transitions)



- **x and y are handled like the other interface transitions.**
- **Bigger search space.**
- **Good balance: one layer of silent transitions.**

# Tool Support in ProM
## (implemented by Eric Verbeek)

# Passage-Based Conformance Checking

# Process Model with 11 Passages

# Conformance Checking per Passage

# Discovery (no initial model, just events)

Process Mining in the Large

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.

van der Aalst

Wil M. P. van der Aalst

Process Mining

# Process Mining

Discovery, Conformance and Enhancement of Business Processes

**www.processmining.org**

**www.win.tue.nl/ieeetfpm/**

Springer