

Mine Your Own Business

Using process mining to turn big data
into better processes and systems

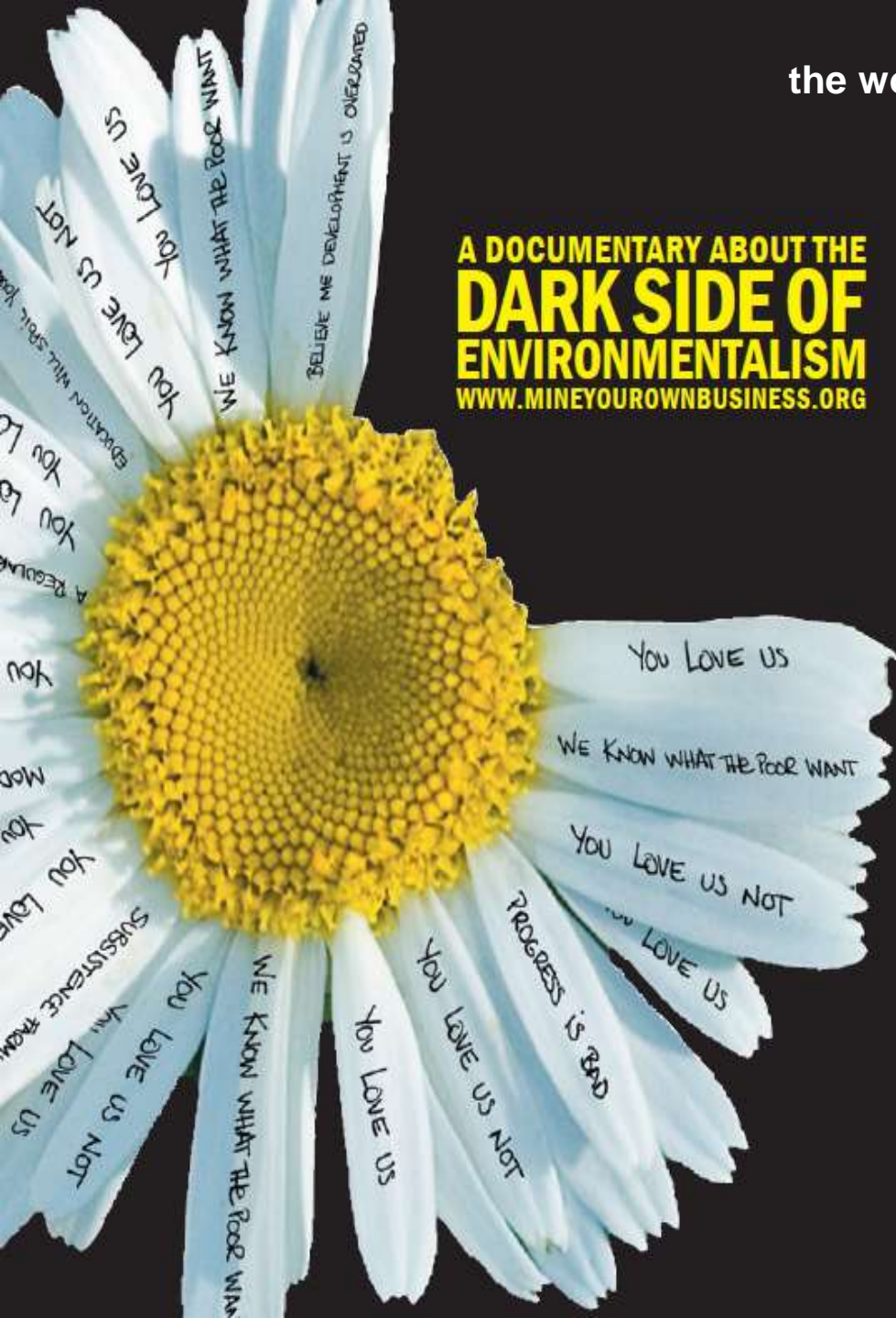
prof.dr.ir. Wil van der Aalst

SCOOBY DOO, WHERE ARE YOU!
IN:
**MINE YOUR
OWN BUSINESS**

© 1969 HANNA-BARBERA PRODUCTIONS, INC.



Season 1, Episode 4 (1969)



**A DOCUMENTARY ABOUT THE
DARK SIDE OF
ENVIRONMENTALISM**
WWW.MINEYOUROWNBUSINESS.ORG

"Mine Your Own Business" (2006)
the world's first anti-environmentalist documentary

A FILM BY PHELIM MCALEER & ANN MCELHINNEY

MINE YOUR OWN BUSINESS

WWW.MINEYOUROWNBUSINESS.ORG



**Mine your own business:
Turning big data into real value**

process mining as
the missing link



aligning model
and reality



divide and
conquer



process
discovery



Big (Event)
Data



challenges



getting
started



process mining as
the missing link



aligning model
and reality



divide and
conquer



process
discovery



Big (Event)
Data



challenges



getting
started





process model analysis

(simulation, verification, optimization, gaming, etc.)

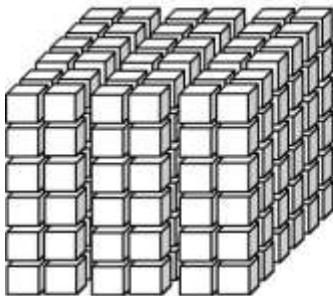


performance-oriented
questions,
problems and
solutions



compliance-oriented
questions,
problems and
solutions

**process
mining**



data-oriented analysis

(data mining, machine learning, business intelligence)

0100110011010101010

010011010101010



10070011010101010





小草对您微微笑
请您把路绕一绕

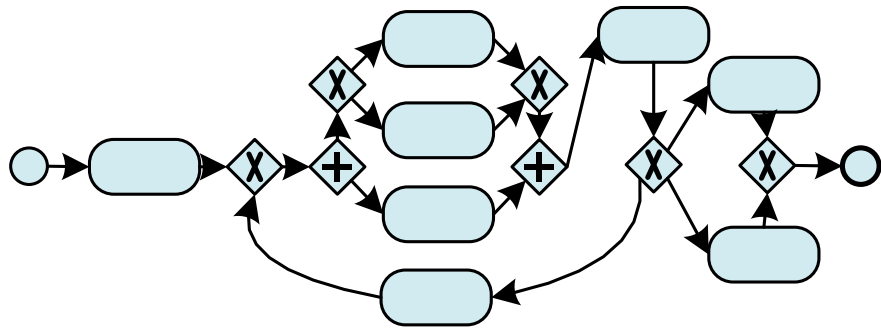
KEEP OFF GRASS

绿色大学办公室
修缮中心园林科

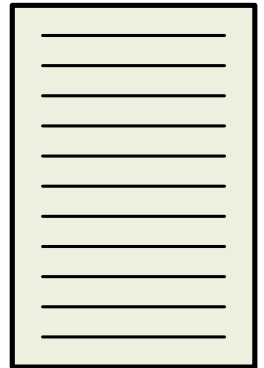


let's play

Play-Out

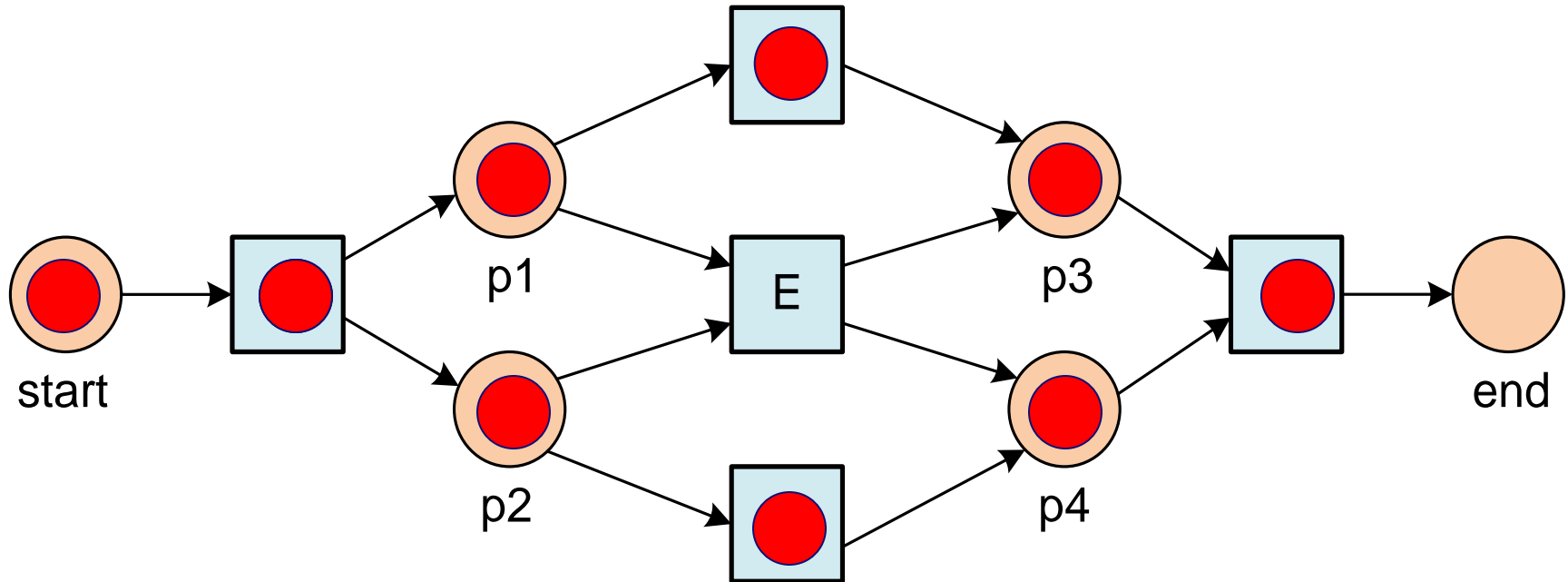


process model



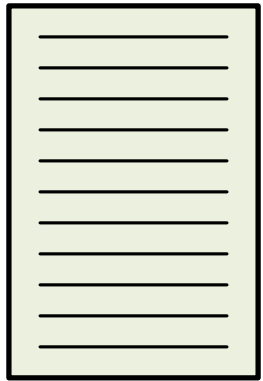
event log

Play-Out (Classical use of models)

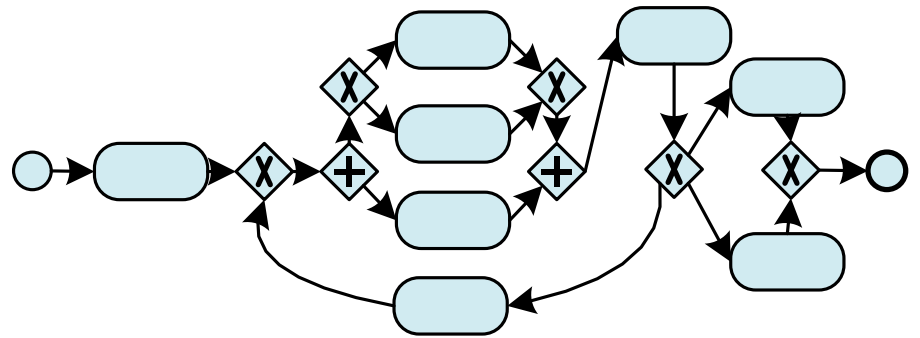
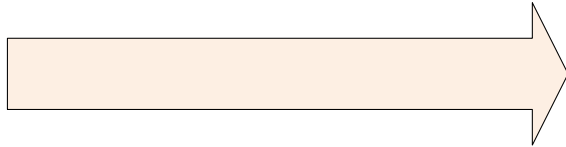


A B C D **A E D** **A E D**
A C B D **A B C D** **A C B D**
A C B D **A E D** **A C B D**

Play-In



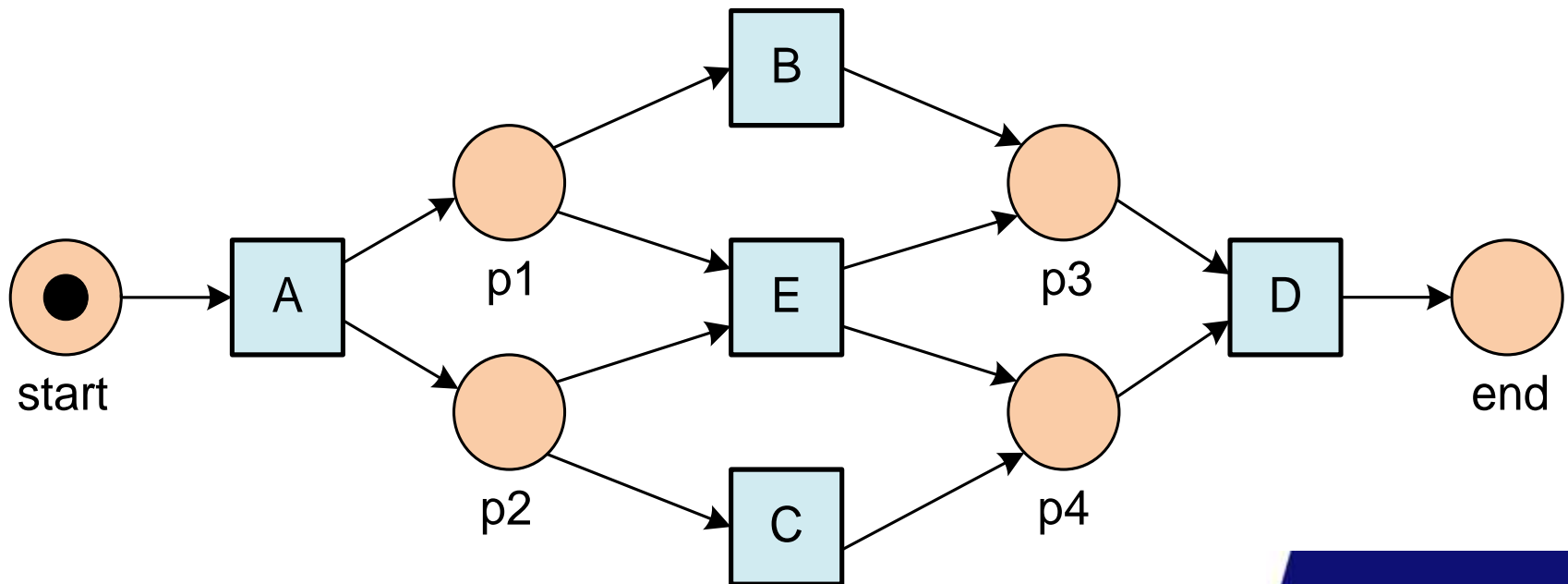
event log



process model

Play-In

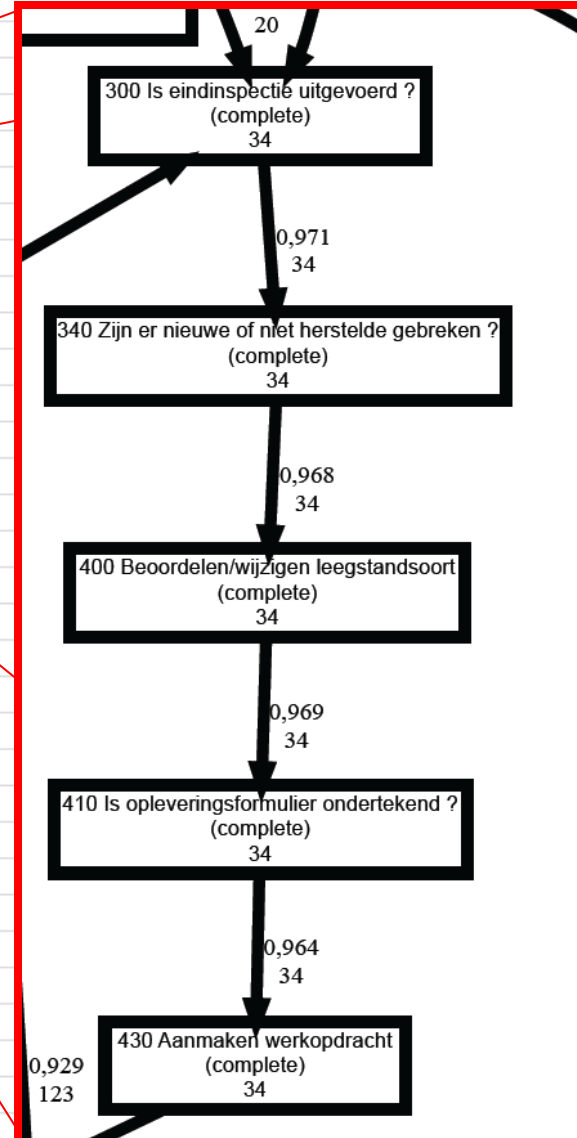
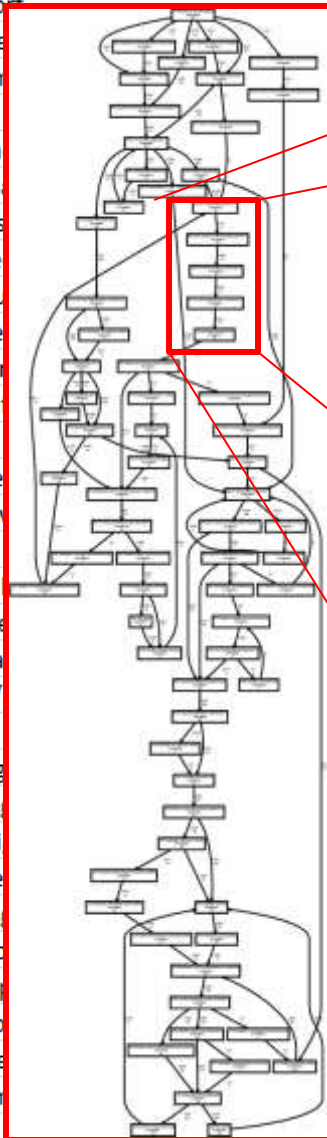
A B C D A E D A E D
A C B D A B C D A C B D
A C B D A E D A C B D



Example Process Discovery

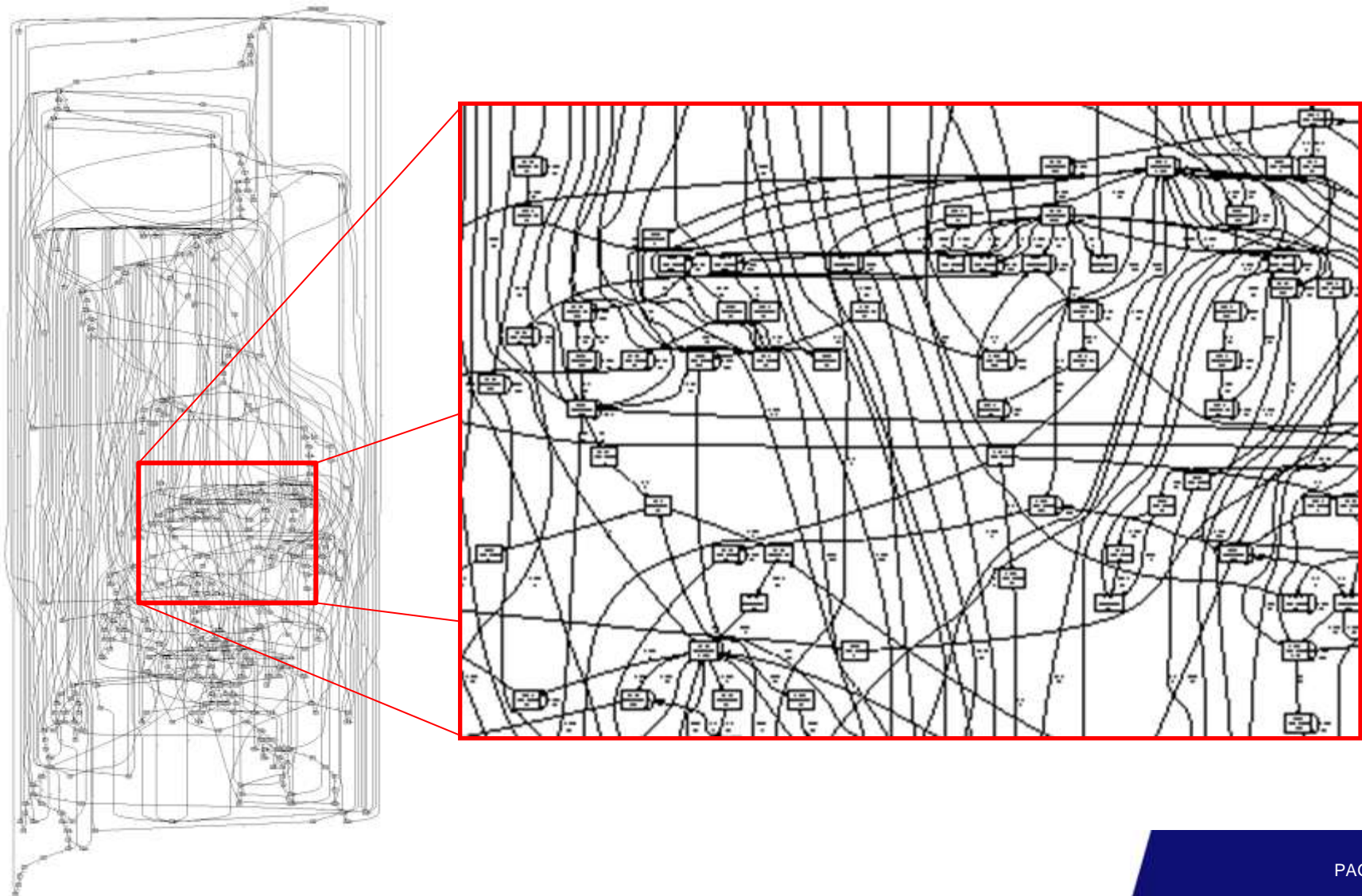
(Vestia, Dutch housing agency, 208 cases, 5987 events)

117315	110	Bepalen leegstandsoort	16.05.2007 14:06:23
117315	120	Plannen eindinspectie	16.05.2007 14:36:01
117315	130	Is het opleveringsform	23.05.2007 09:41:40
117315	150	Is er sprake van ZAV ?	23.05.2007 09:41:51
117315	170	Aanpassen plattegron	23.05.2007 11:57:18
117315	180	Aanpassen woningwa	23.05.2007 09:42:37
117315	190	Actualiseren huurprijs	23.05.2007 09:48:23
117315	200	Toewijzen woning/be	23.05.2007 09:48:29
117315	210	Registreren voorl. hu	10.09.2007 16:24:36
117315	220	Is contract getekend e	11.09.2007 14:56:18
117315	240	Definitief maken Huu	31.03.2008 16:17:12
117315	250	Aanpassen factuureera	09.09.2008 15:39:59
117315	260	After sales	09.09.2008 16:51:24
117315	270	Archiveren nieuwe ve	10.09.2008 07:52:08
117315	300	Is eindinspectie uitge	07.06.2007 14:47:04
117315	340	Zijn er nieuwe of niet	07.06.2007 14:47:06
117315	400	Beoordelen/wijzigen	07.06.2007 14:51:16
117315	410	Is opleveringsformulie	07.06.2007 14:51:26
117315	430	Aanmaken werkopdra	11.06.2007 09:21:39
117315	440	Worden er bonussen/	11.06.2007 09:21:49
117315	460	Opstellen eindnota	08.08.2007 16:18:26
117315	470	Archiveren huuropzeg	09.08.2007 14:42:23
119763	010	Registreren huuropze	09.05.2007 11:19:14
119763	030	Vastleggen toekomst	09.05.2007 12:25:01
119763	050	Inplannen afspraak 1e	09.05.2007 11:59:52
119763	060	Aanmaken bevestigin	09.05.2007 12:31:57
119763	070	Is 1e inspectie uitgev	16.05.2007 13:04:26
119763	100	Gereedmelden 1e ins	16.05.2007 13:43:39
119763	110	Bepalen leegstandsoo	16.05.2007 13:43:28
119763	120	Plannen eindinspectie	16.05.2007 13:42:58
119763	130	Is het opleveringsform	16.05.2007 13:34:49
119763	150	Is er sprake van ZAV ?	16.05.2007 13:34:56



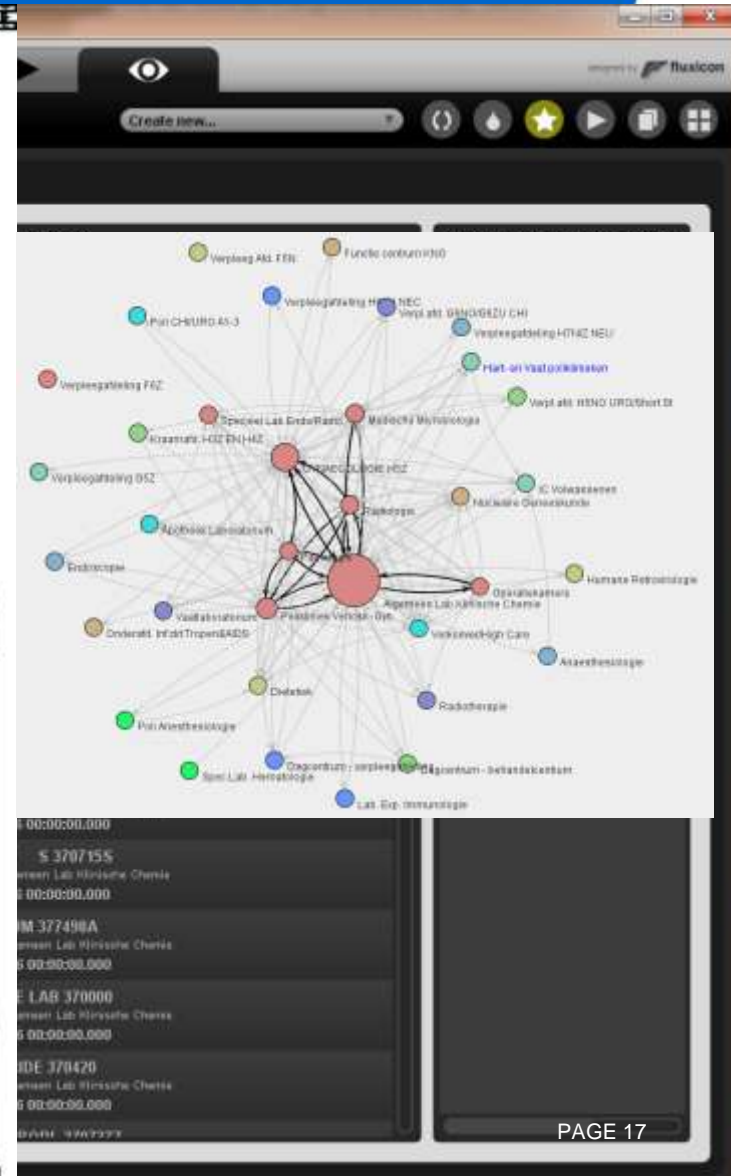
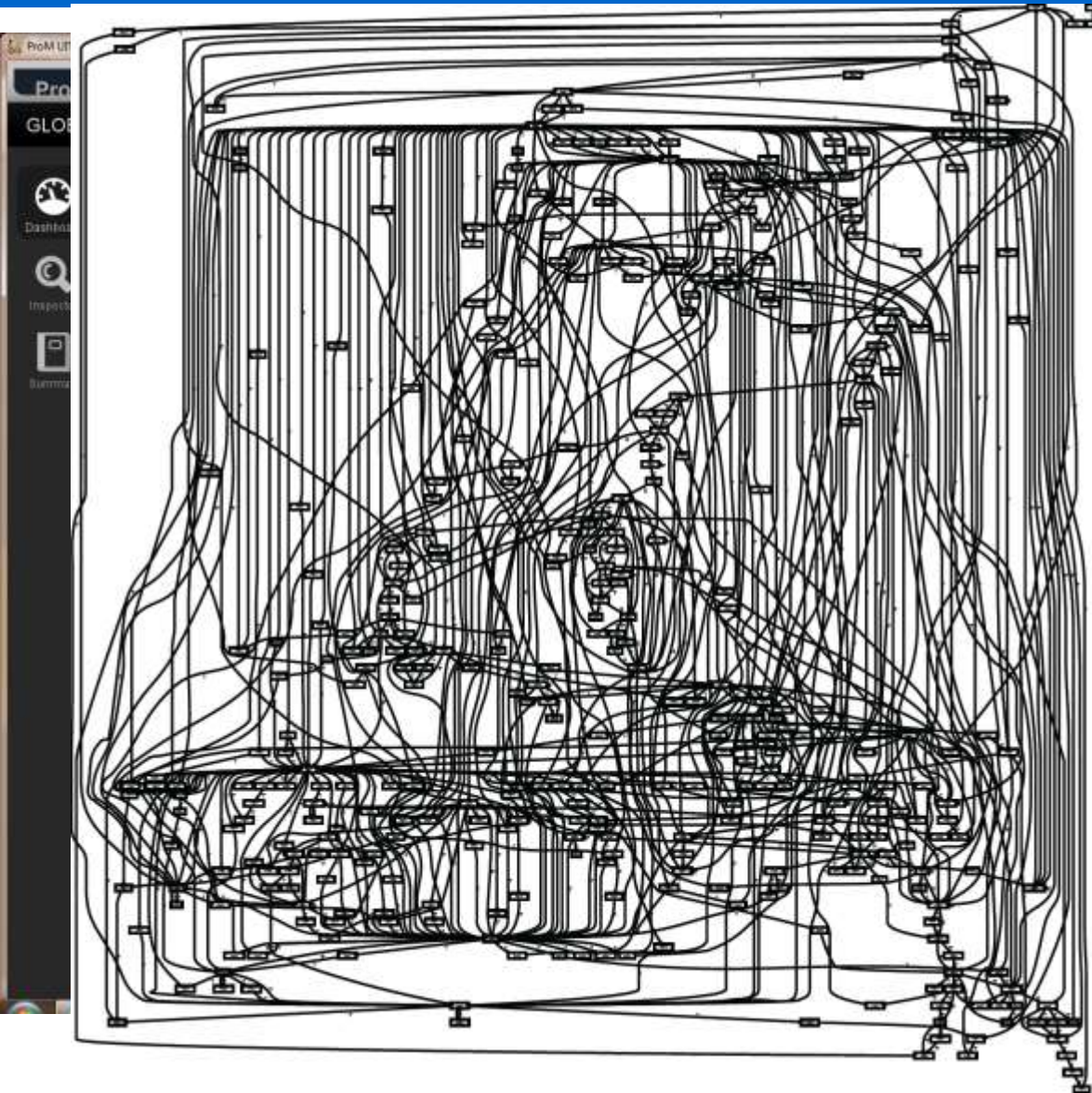
Example Process Discovery

(ASML, test process lithography systems, 154966 events)

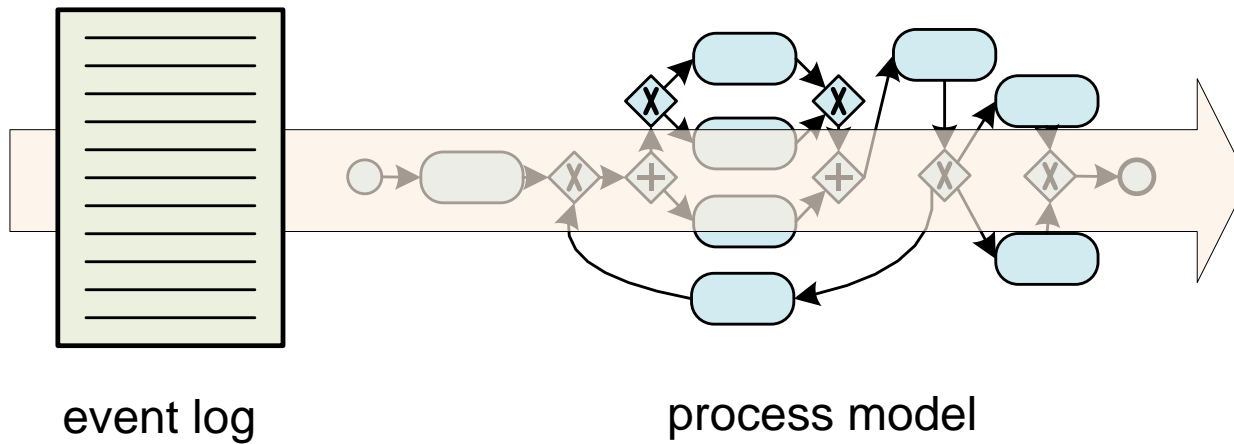


Example Process Discovery

(AMC, 627 gynecological oncology patients, 24331 events)



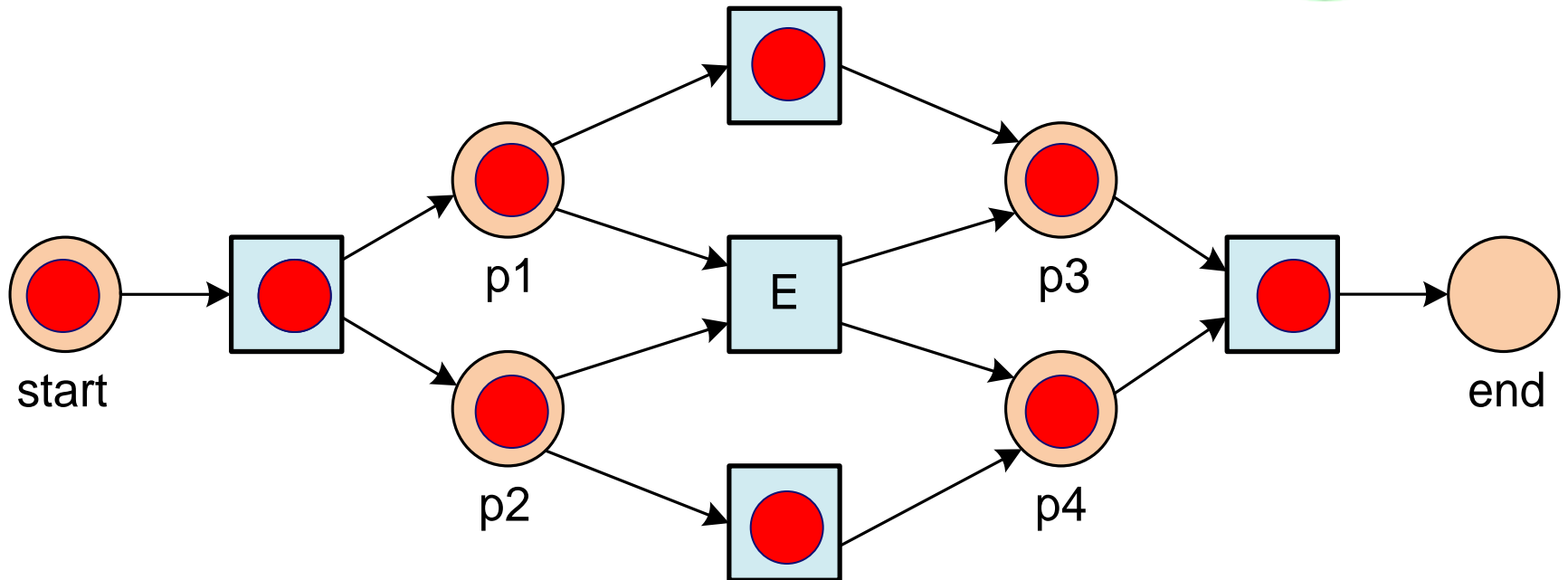
Replay



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

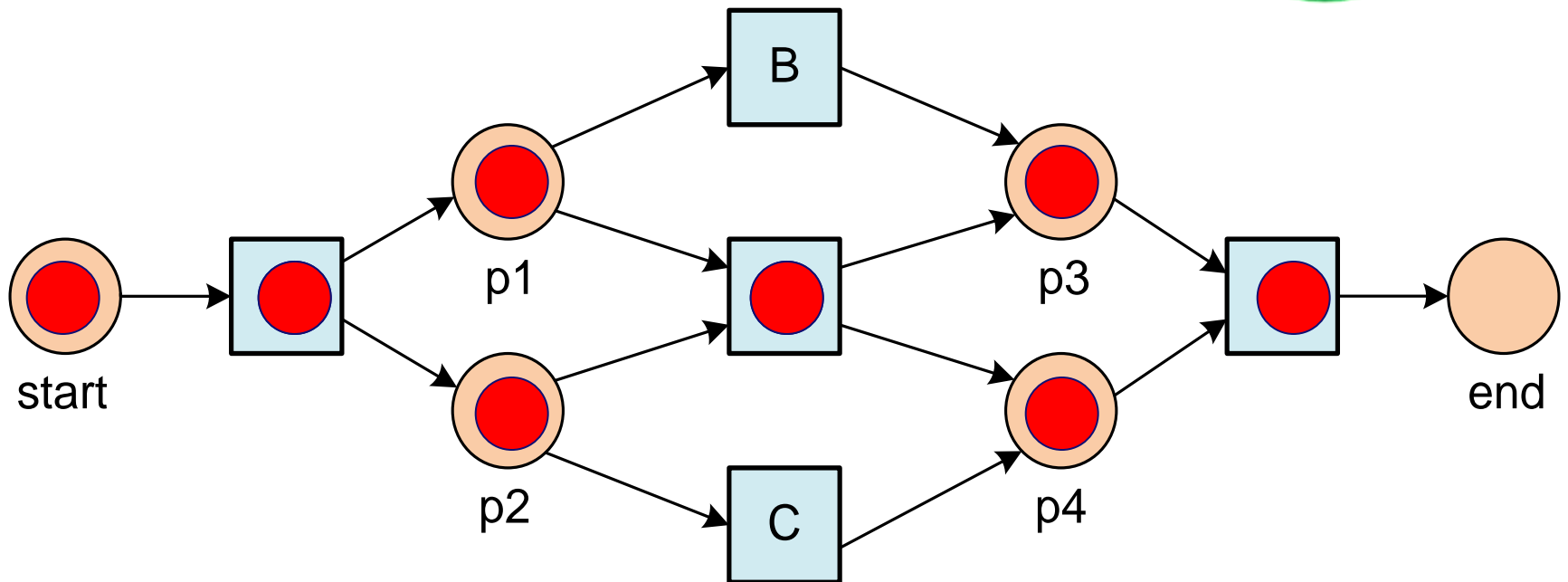
Replay

A B C D



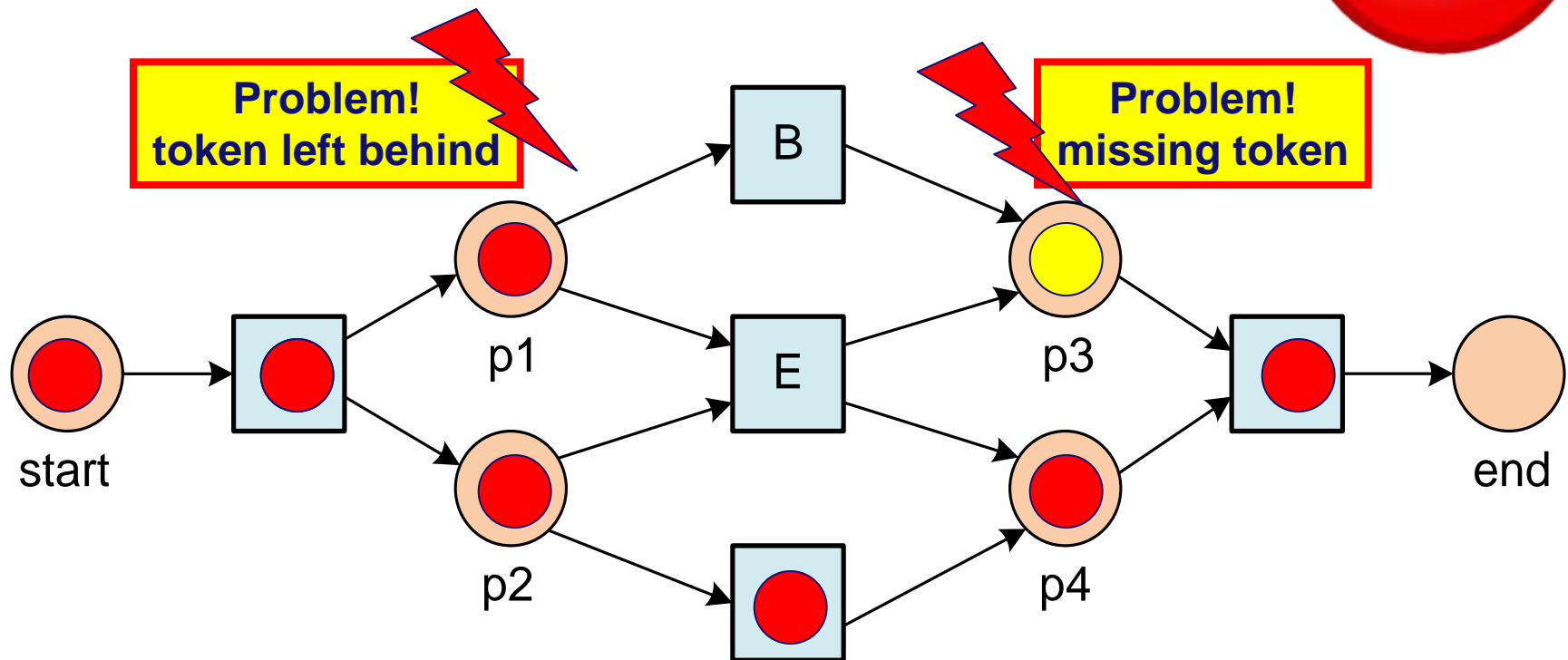
Replay

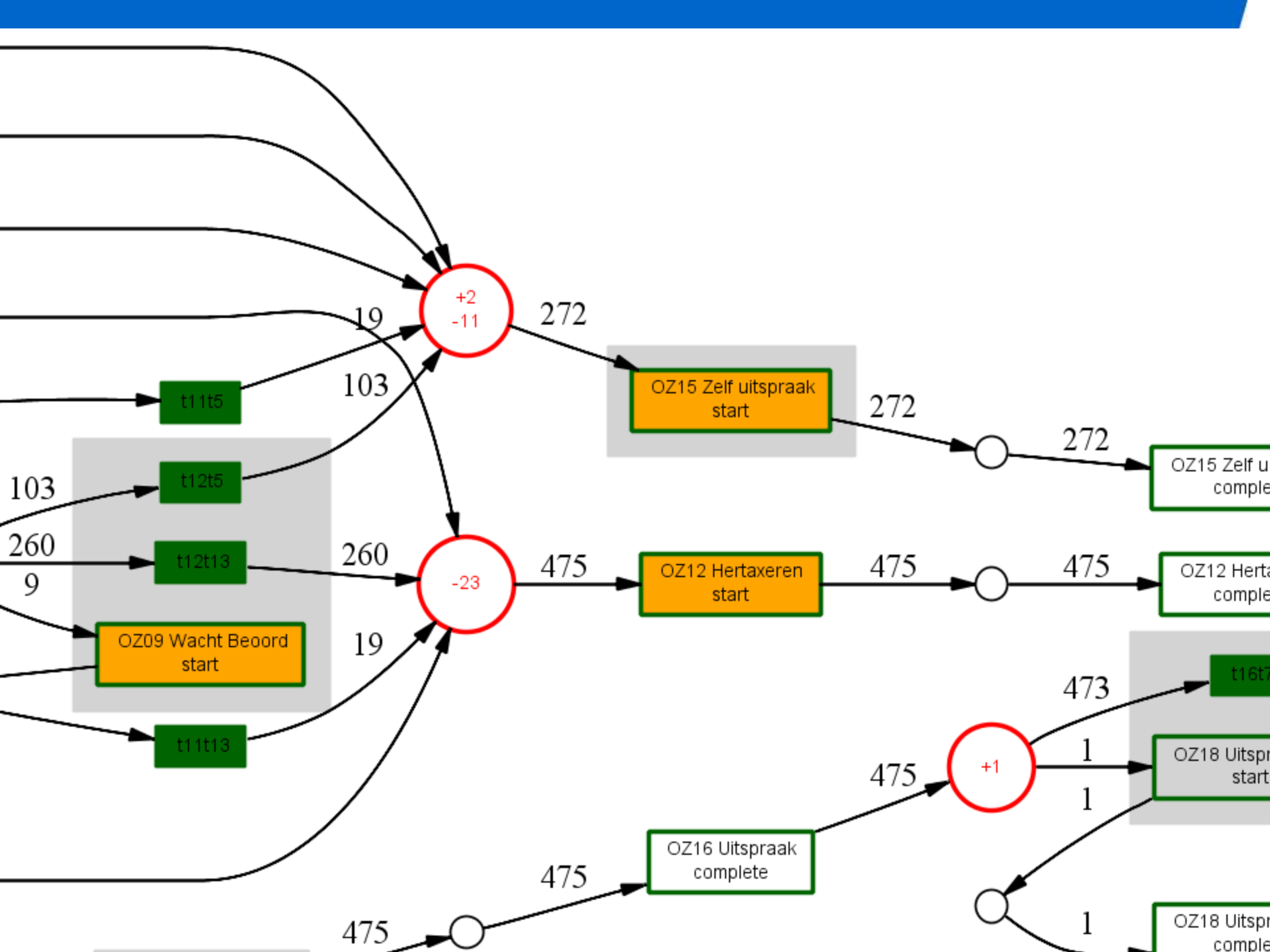
A E D



Replay can detect problems

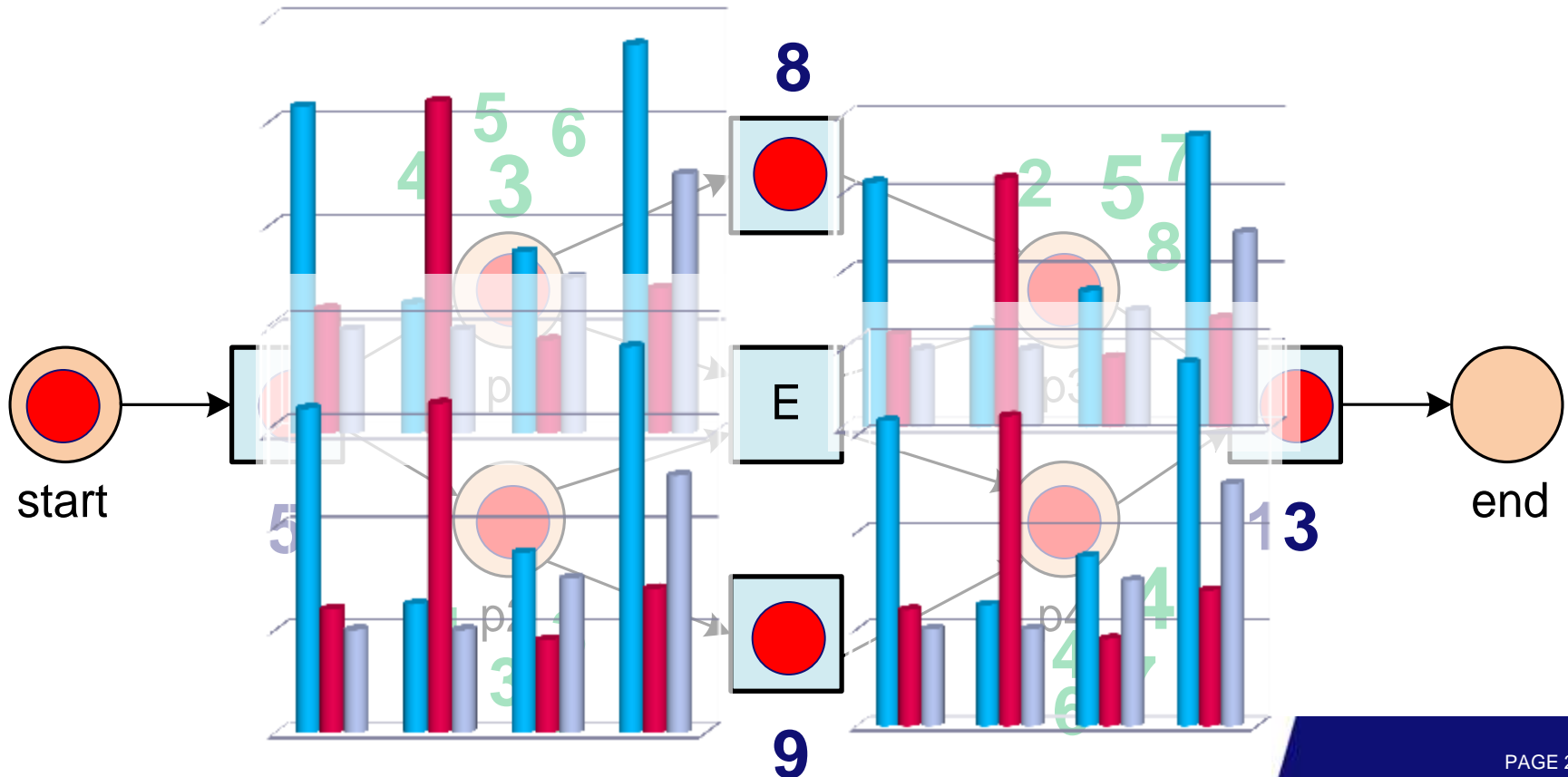
ACD





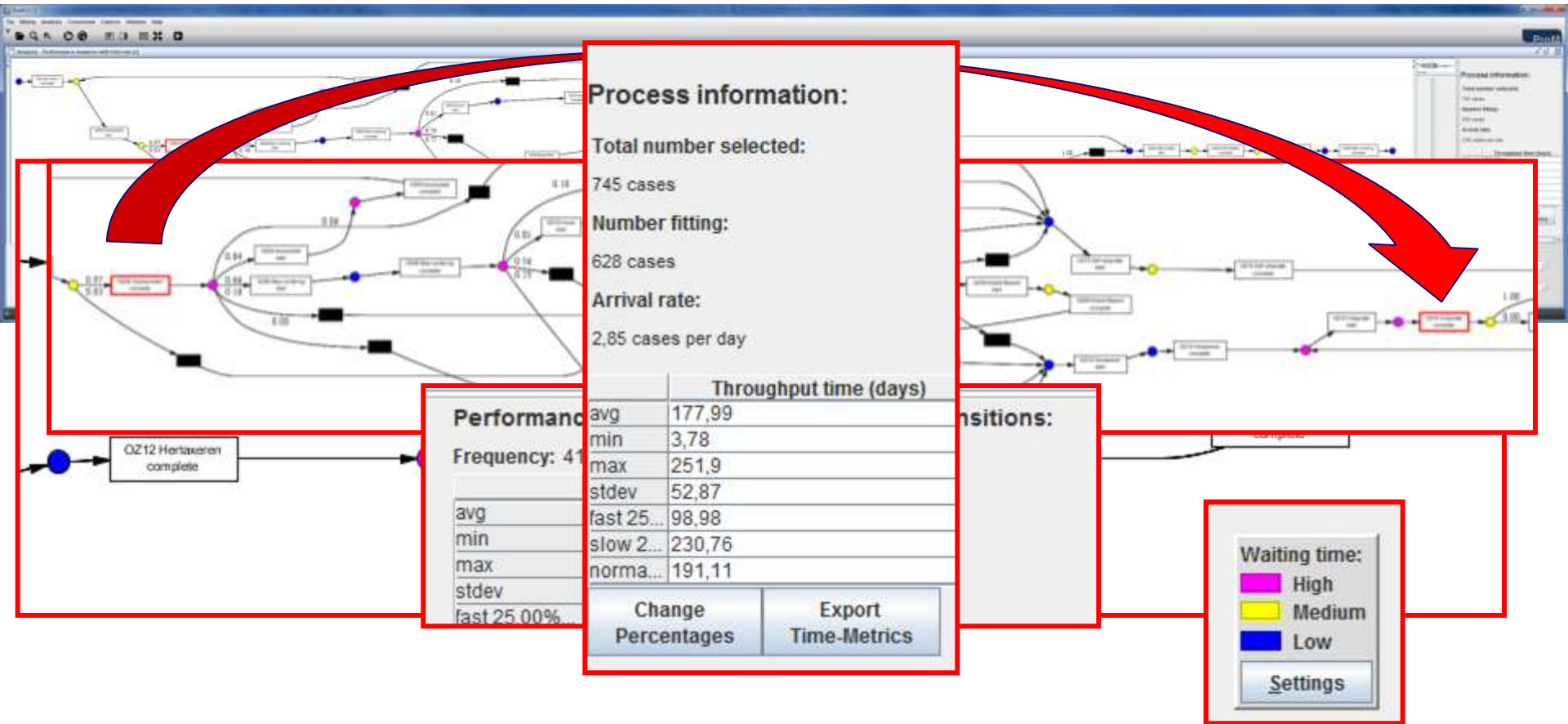
Replay can extract timing information

A⁵B⁸C⁹D¹³



Performance Analysis Using Replay

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



Models are like the glasses required to see and understand event data!



process mining as
the missing link



aligning model
and reality



divide and
conquer



process
discovery



Big (Event)
Data



challenges



getting
started



Language identification in the limit (Mark Gold 1967)



A language is **learnable in the limit** if there exists a perfect child that generates only finitely many hypotheses.

Learning is not easy ...

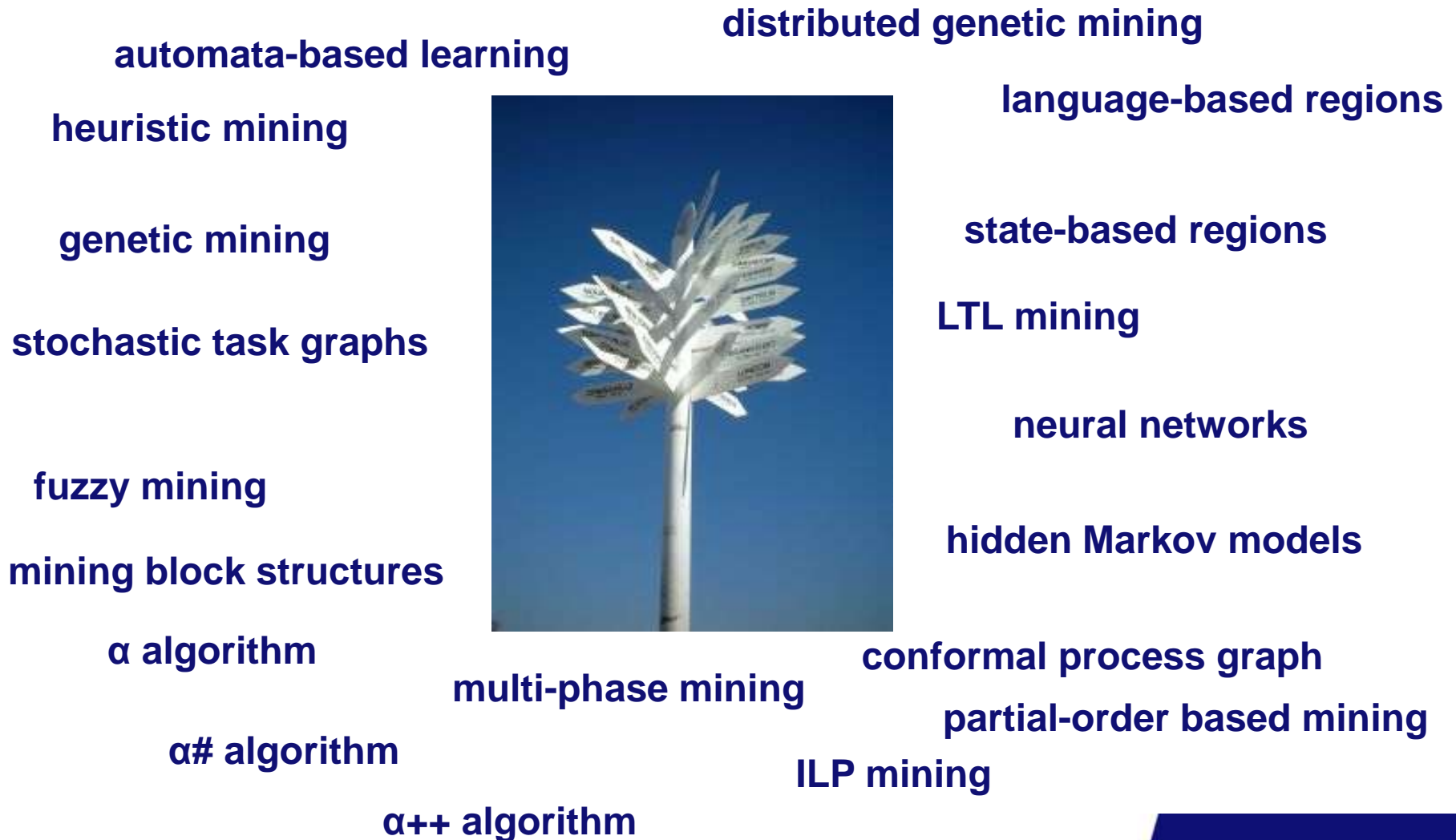


- Even simple languages (e.g. regular languages) are not learnable in general
- Most models do not consider concurrency and definitely not end-to-end business process models.

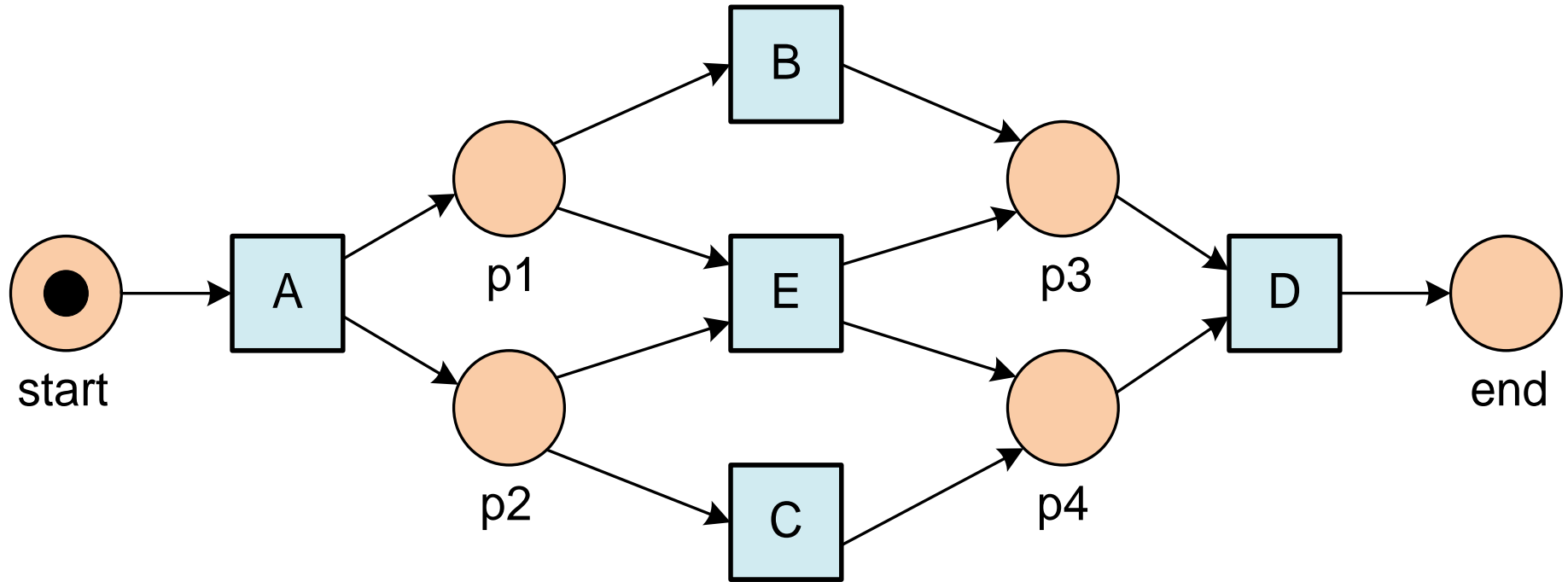
Classical approaches (before 1998) did not consider concurrency and definitely not end-to-end business process models.

reference \cong trace in event log
language \cong process model

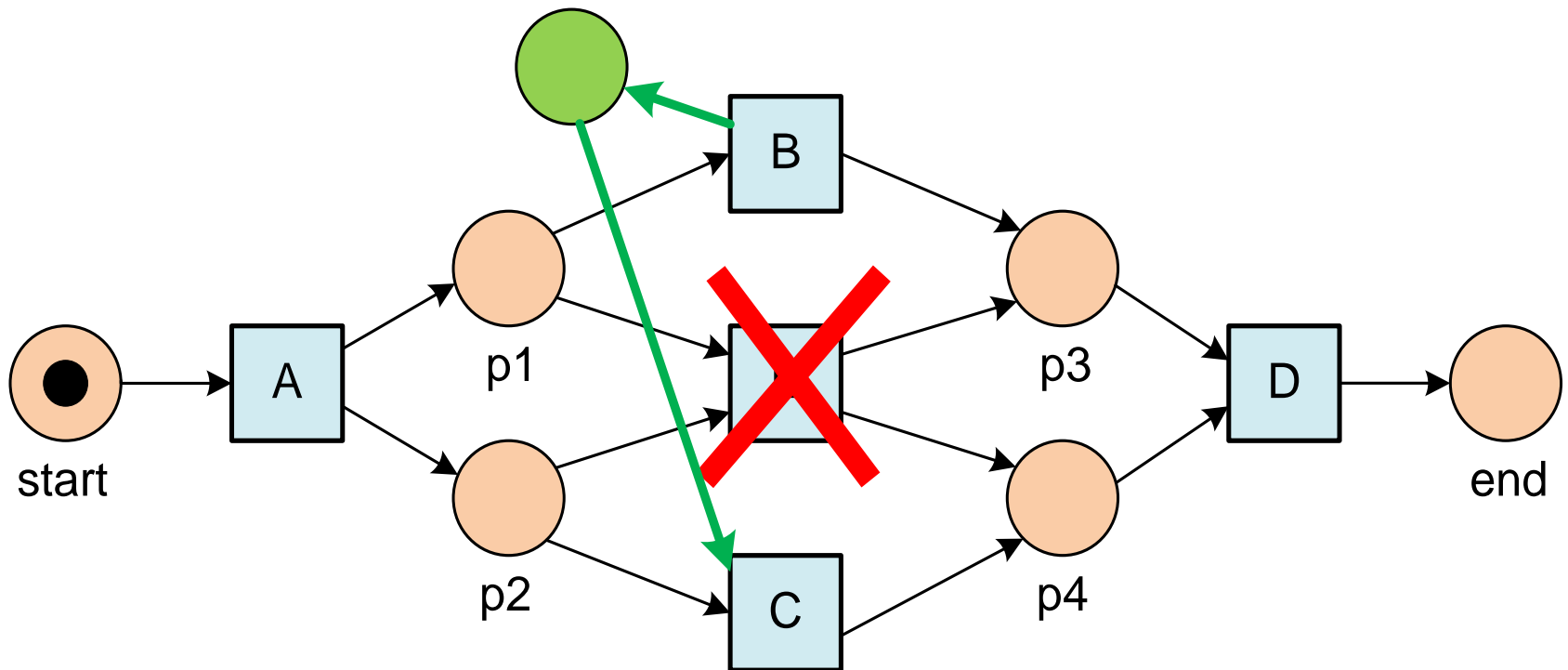
Process discovery algorithms (small selection)



Quiz Question: How to remove behavior?

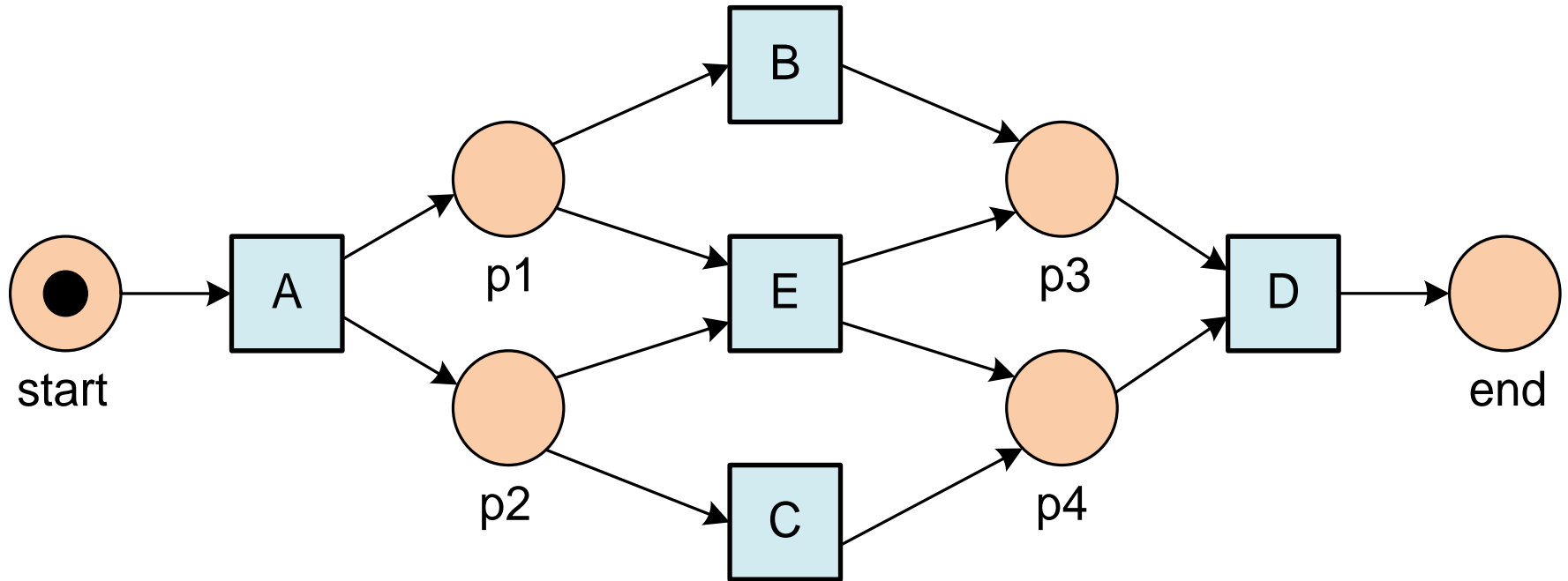


Quiz Question: How to remove behavior?

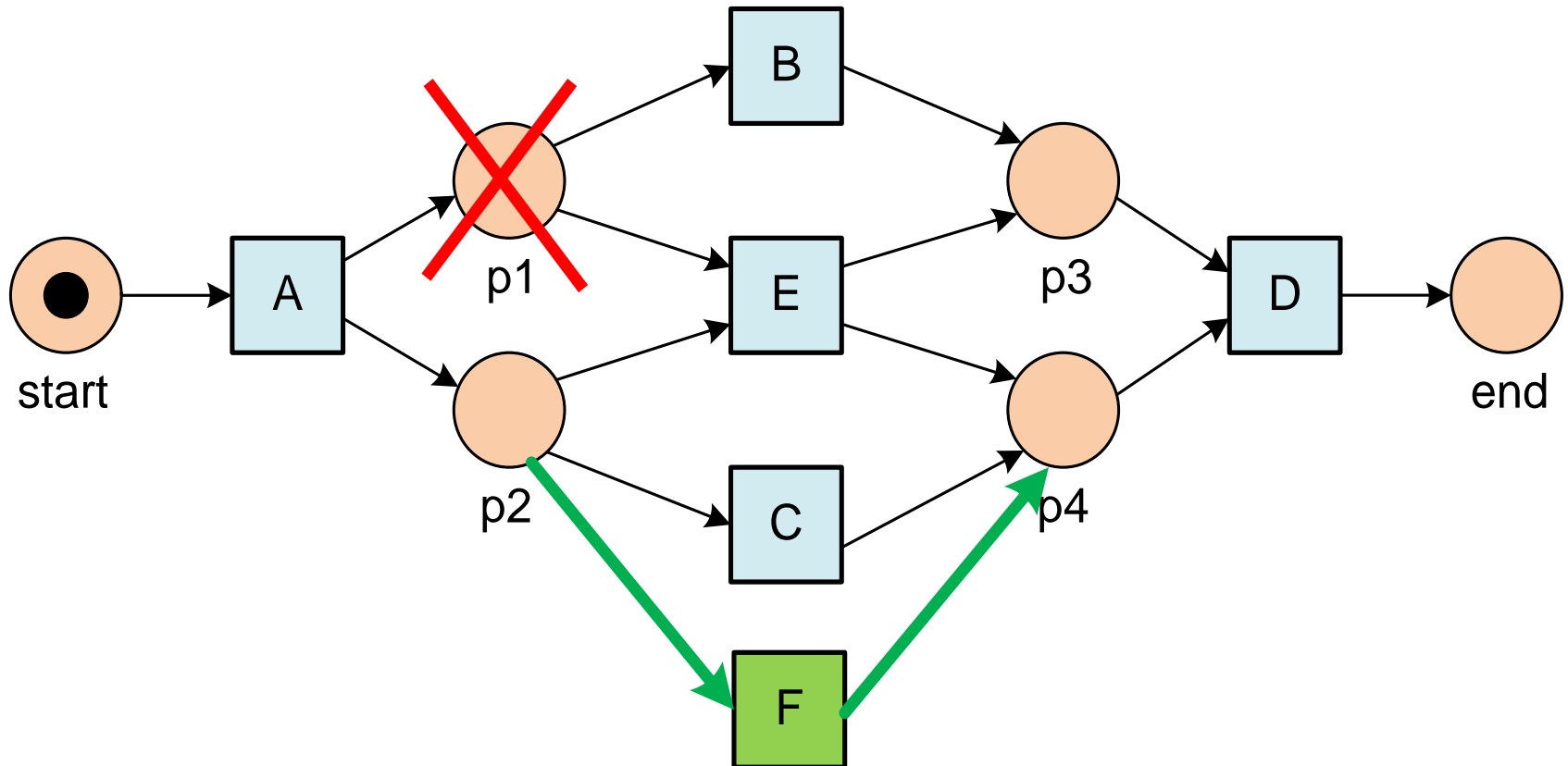


Add places or remove transitions!

Quiz Question: How to add behavior?

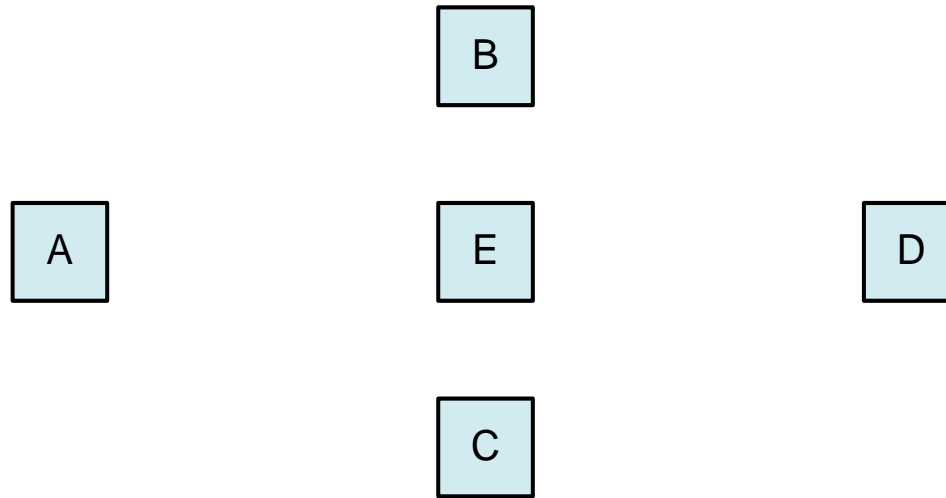


Quiz Question: How to add behavior?



Add transition or remove places!

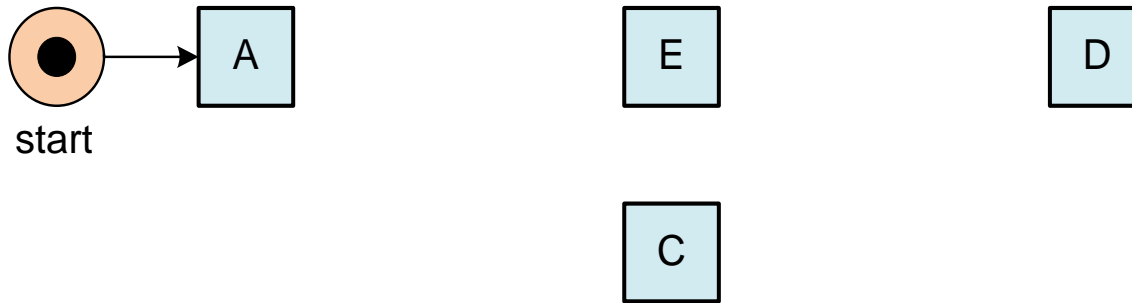
Places limit behavior



- **abcd**
- **ad**
- **abed**
- **abccd**
- **acbd**

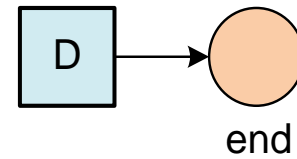
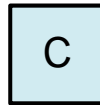
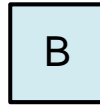
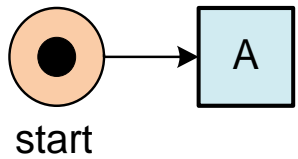
- **aebcd**
- **aed**
- **aad**
- **caed**
- **aded**

Places limit behavior



- abcd
- ad
- abed
- abccd
- acbd
- aeccd
- aed
- **aad**
- caed
- aded

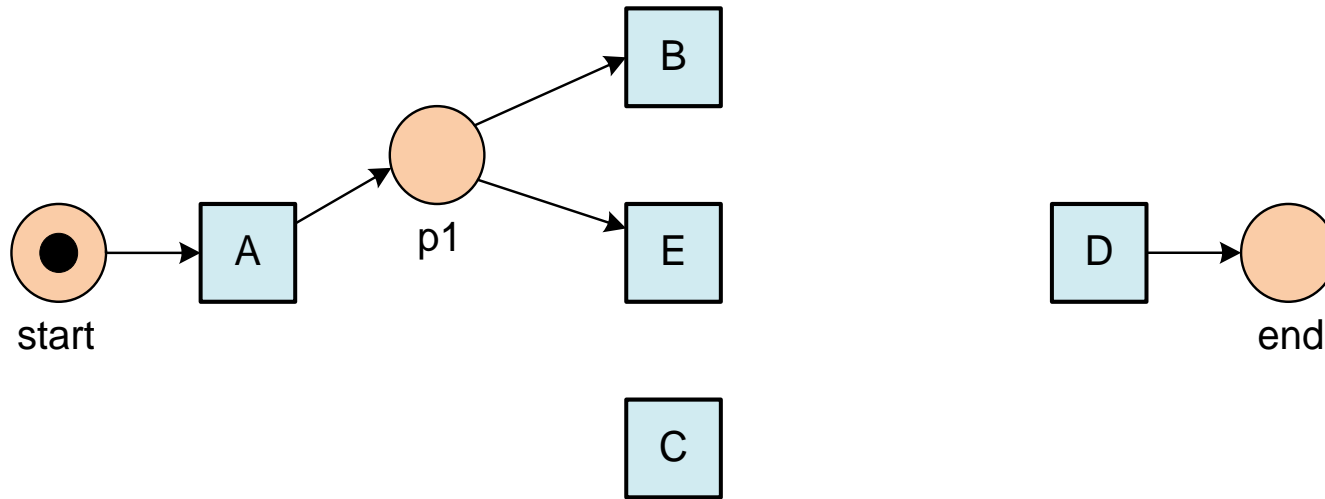
Places limit behavior



- **abcd**
- **ad**
- **abed**
- **abccd**
- **acbd**

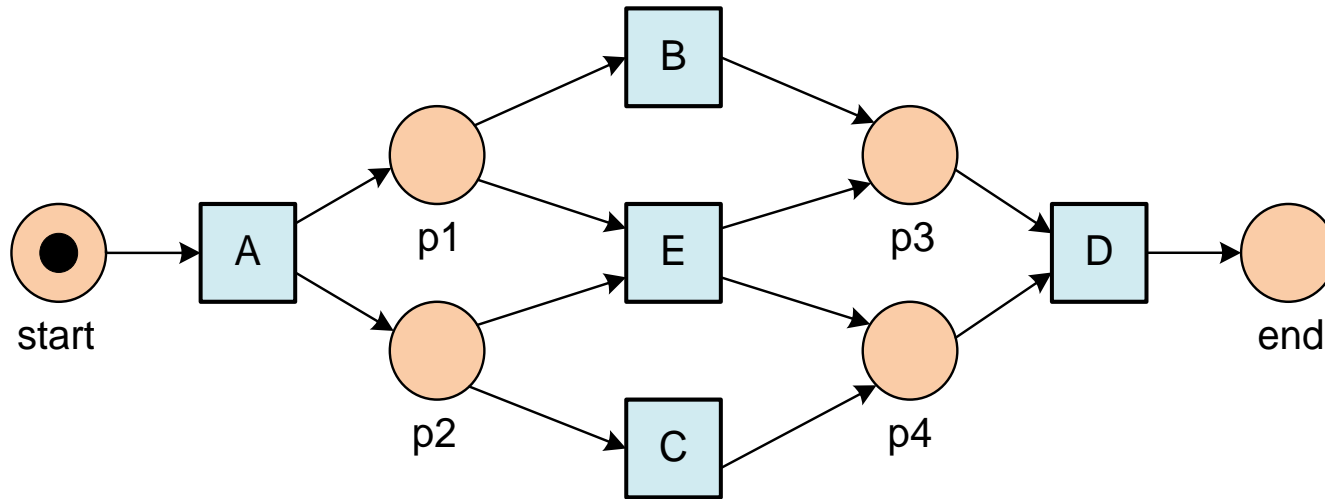
- **aebcd**
- **aed**
- **aad**
- **caed**
- **aded**

Places limit behavior



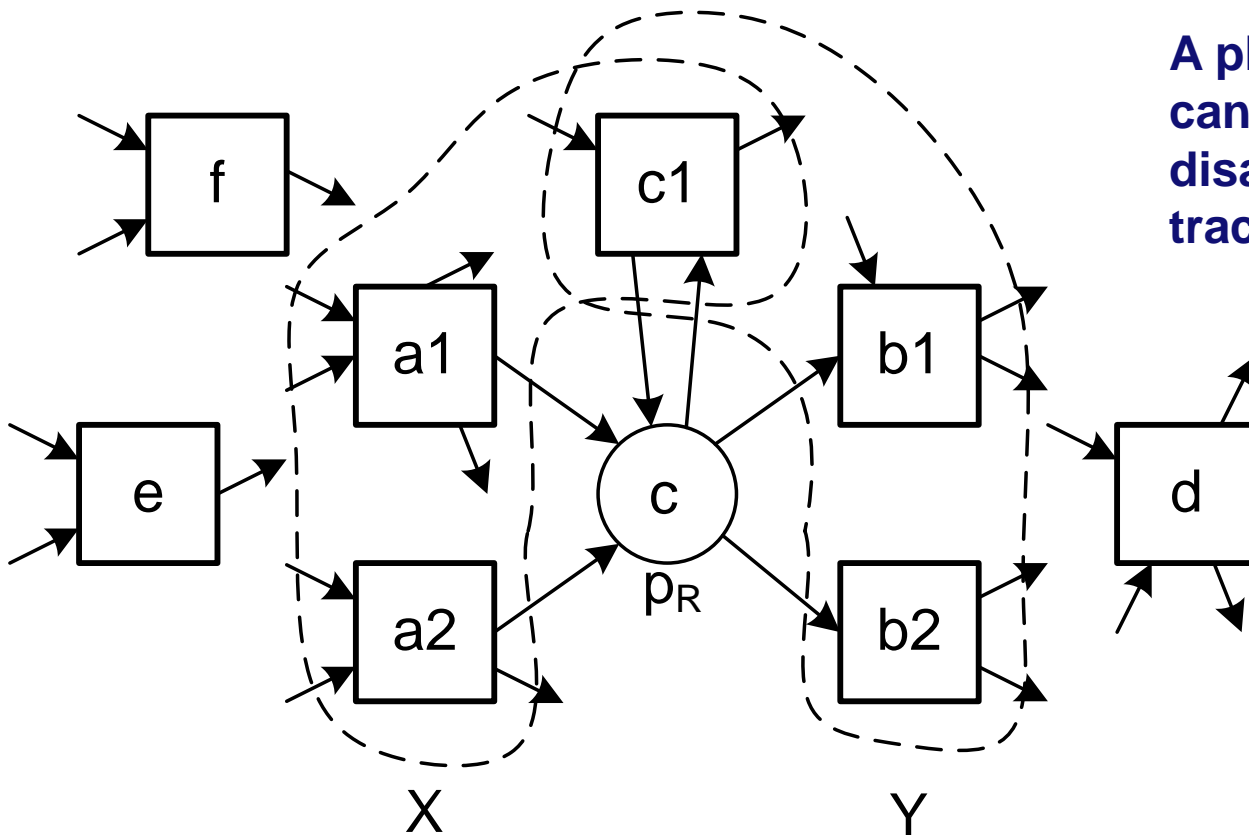
- **abcd**
- **ad**
- **abed**
- **abccd**
- **acbd**
- **aebcd**
- **aed**
- **aad**
- **caed**
- **aded**

Places limit behavior



- **abcd**
- **ad**
- **abed**
- **abccd**
- **acbd**
- **aebcd**
- **aed**
- **aad**
- **caed**
- **aded**

Example: Process Discovery Using Language-Based Regions

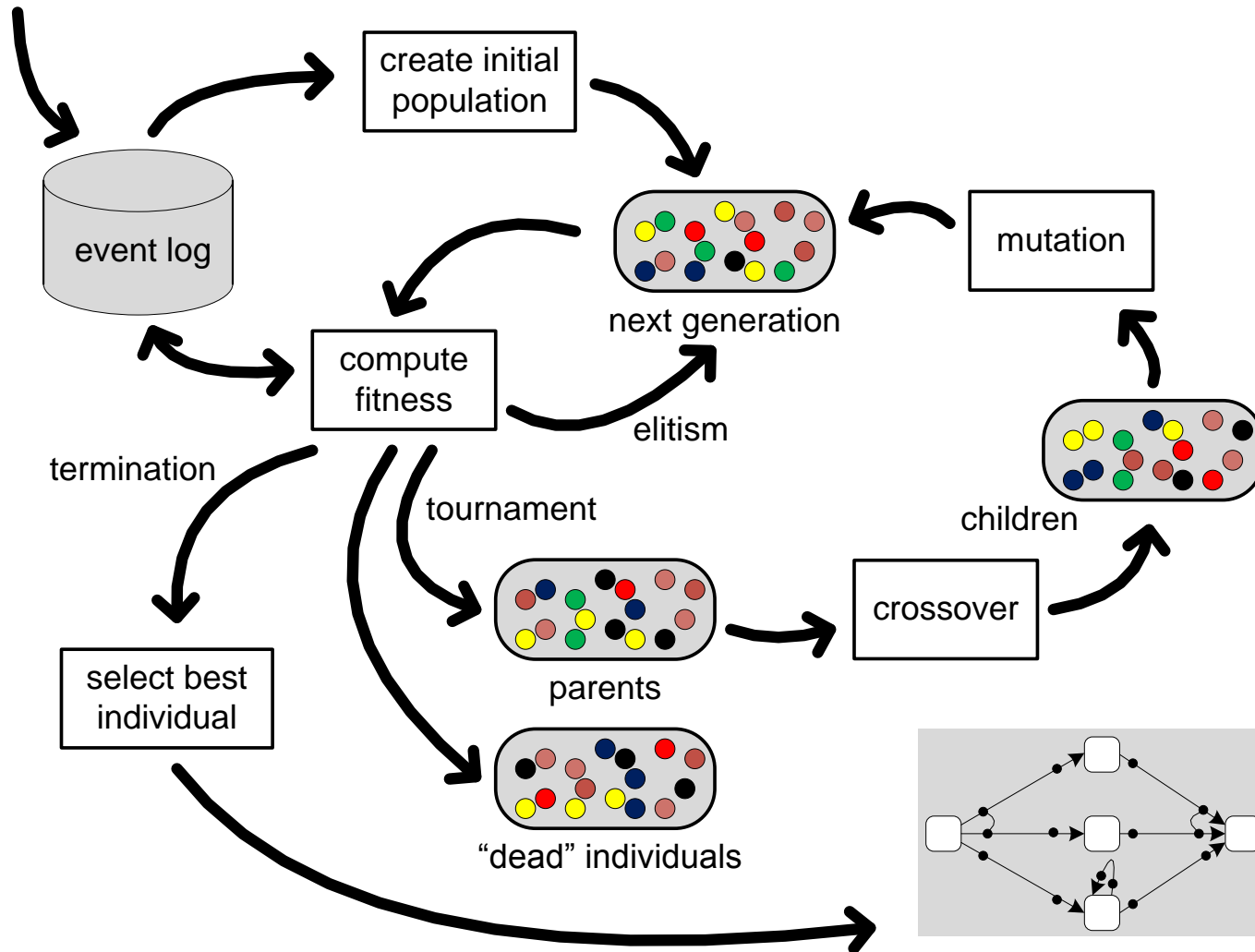


A place is **feasible** if it can be added without disabling any of the traces in the event log.

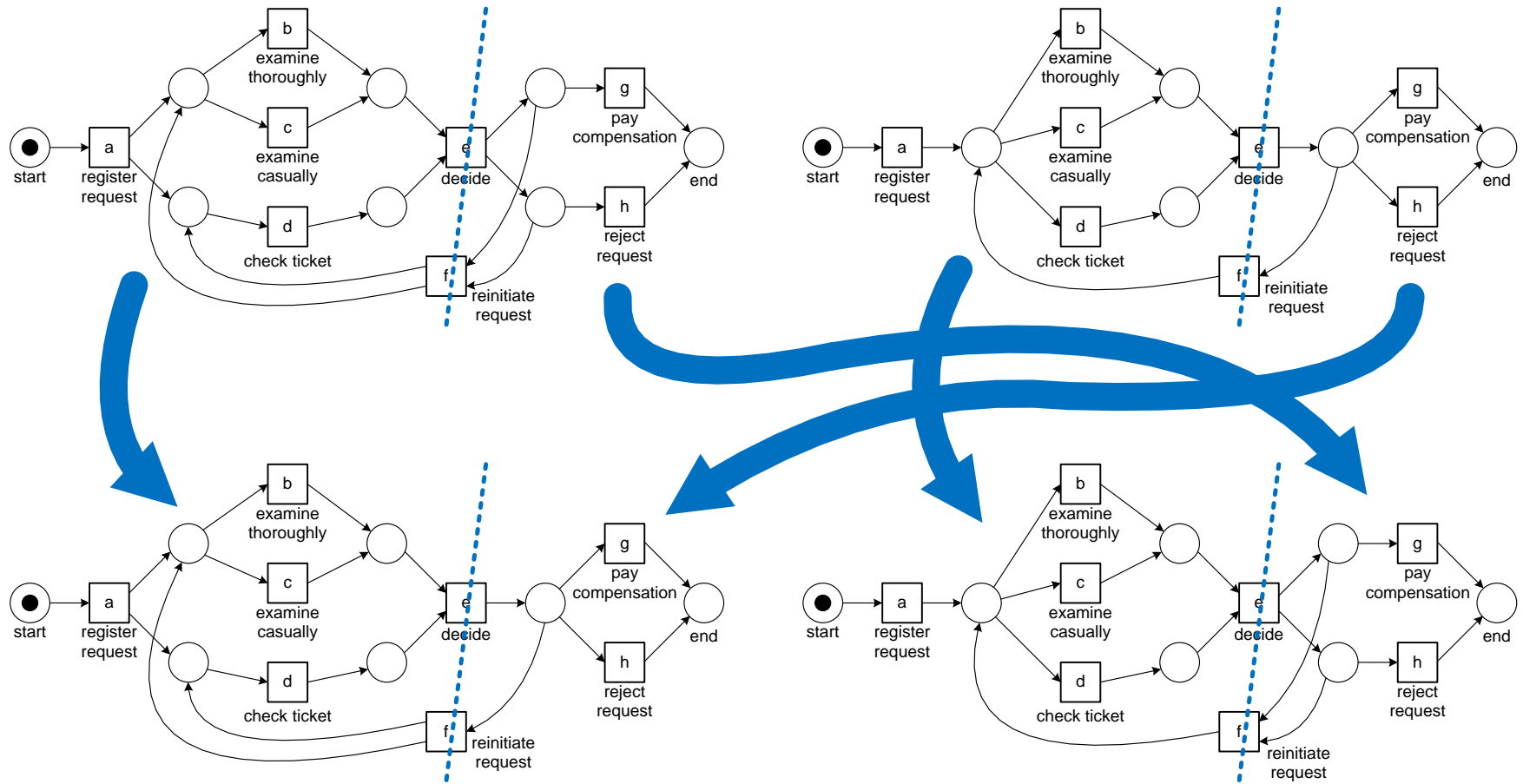
for any $\sigma \in L$, $k \in \{1, \dots, |\sigma|\}$, $\sigma_1 = hd^{k-1}(\sigma)$, $a = \sigma(k)$, $\sigma_2 = hd^k(\sigma) = \sigma_1 \oplus a$:

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

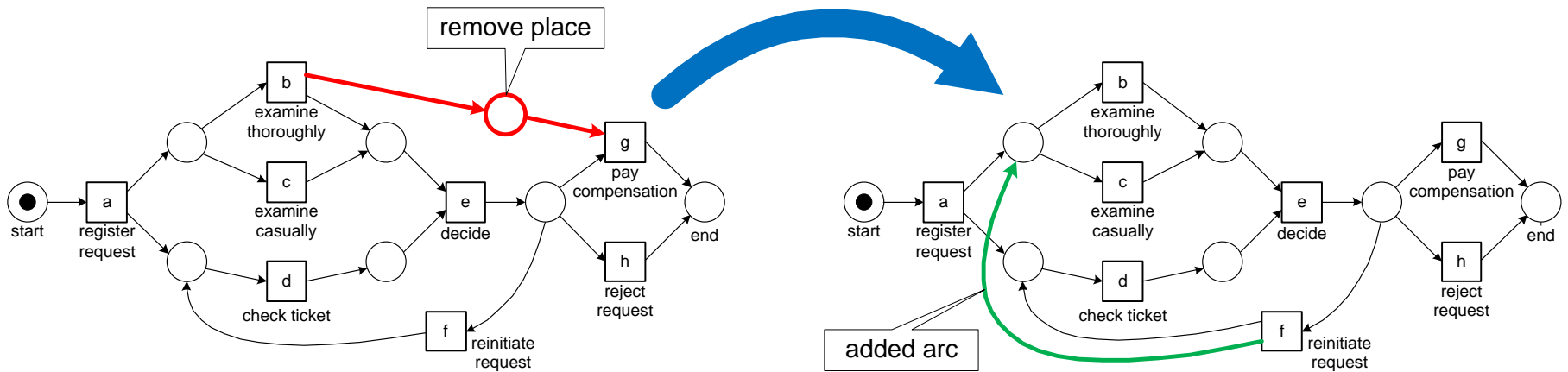
Genetic process mining: Overview

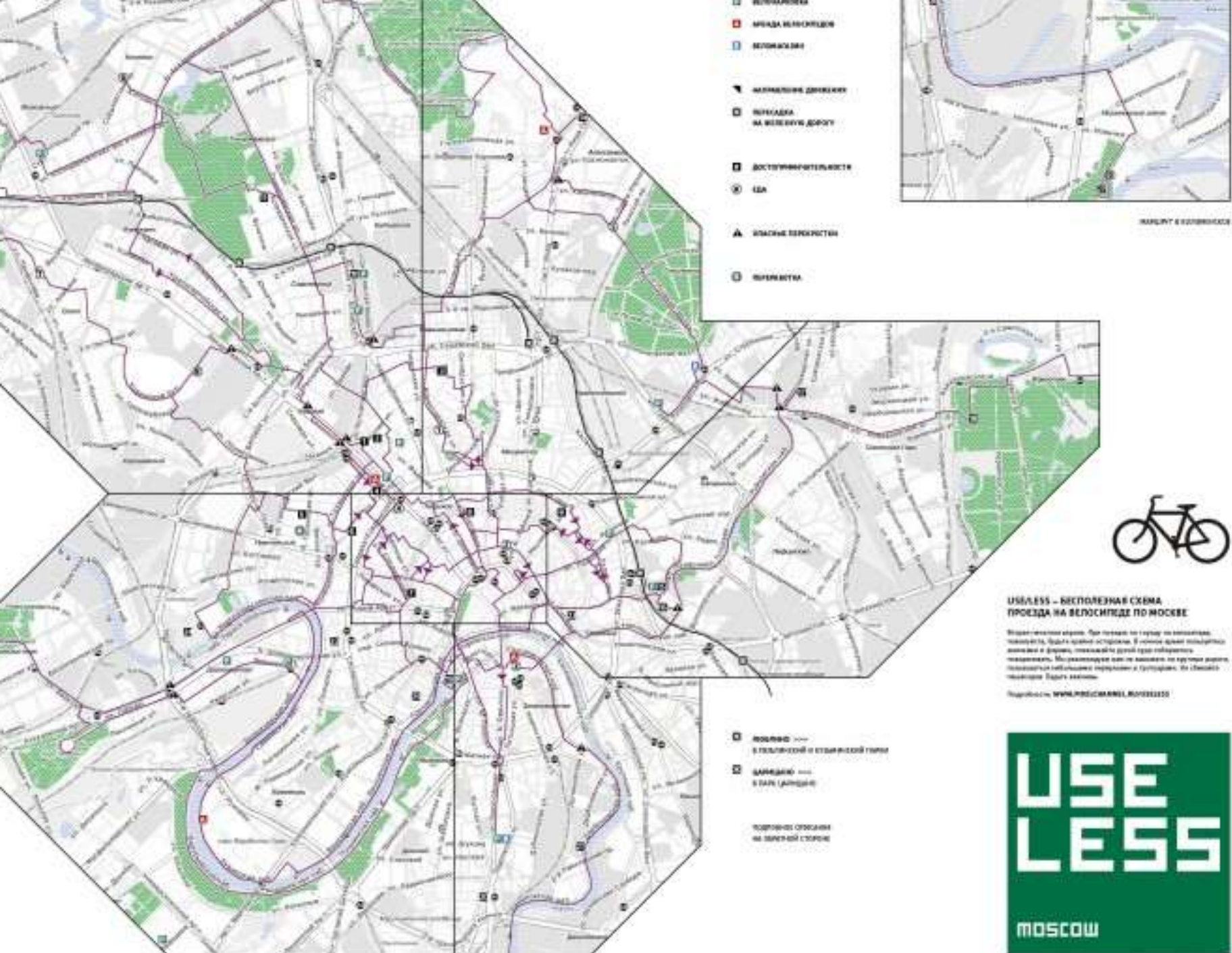


Example: crossover



Example: mutation





USELESS – БЕСПОЛЕЗНАЯ СХЕМА ПРОЕЗДА НА ВЕЛОСИКЕ ПО МОСКВЕ

Видео-информационный ресурс. Для поиска по городу по различным параметрам, чтобы найти интересные и новые для вас маршруты, увидеть и скачать спутниковый снимок и другие полезные функции. Мы работаем над тем, чтобы сделать использование нашего ресурса наиболее удобным и удобным. На сайте Москва.Район.России.

Сайт: WWW.MOSCOW.MOSKVA.RU

- ПОСЛЕД — в гольф-клубе и старинной гонимой
- БАРЬЕРЫ — в БАР (Алматы)

ПОДРОБНОЕ ОПИСАНИЕ НА ЗАКРЕПЛЕННОЙ СТОРОНЕ



process mining as
the missing link



aligning model
and reality



divide and
conquer



process
discovery



Big (Event)
Data



challenges



getting
started



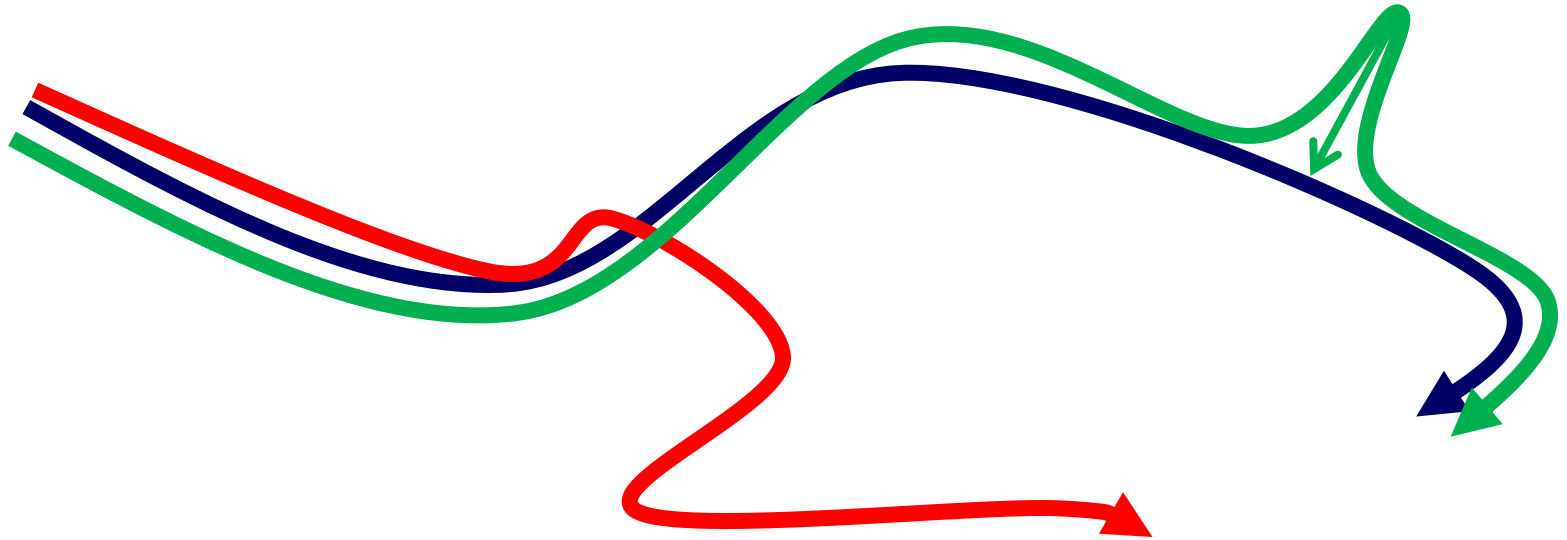
Conformance checking

The image shows a screenshot of a Microsoft Word document titled "hello world.docx". The document text is "Recent breakthroughs in process mining research make it possible to discover, analyze, and improve business processes based on event data. People, machines, and software leave trails of events such as entering a customer order into SAP, checking in for a flight, booking a room for a patient, and rejecting a building information system. Over the past few years, there has been a spectacular growth of data. Moreover, the digital universe and the physical universe has becoming more and more aligned." The text is annotated with five yellow callout boxes:

- Top-left: "an activity that should not happen happened" (points to "analyze")
- Top-right: "an activity was executed by the wrong person" (points to "machines")
- Center: "an activity was executed too late" (points to "booking a room")
- Bottom-left: "an activity that should happen did not happen" (points to "rejecting a building")
- Bottom-center: "two activities were swapped" (points to "booking a room" and "checking in for a flight")

Microsoft Word interface elements include the ribbon (File, Home, Insert, Page Layout, References, Mailings, Review, View), the ribbon tabs (Themes, Page Setup, Page Background, Paragraph), and the status bar (Page: 1 of 1, Words: 95, English (U.S.), 103%).

Alignments are essential!



- conformance checking to diagnose deviations
- squeezing reality into the model to do model-based analysis

<i>a</i>	<i>c</i>	\gg	<i>d</i>	\gg	<i>f</i>	\gg
<i>a</i>	<i>c</i>	<i>b</i>	<i>d</i>	τ	\gg	<i>h</i>
<i>t1</i>	<i>t4</i>	<i>t3</i>	<i>t5</i>	<i>t7</i>		<i>t10</i>

process
model

event log

synchronous
move

a	c	\gg	d	\gg	f	\gg
a	c	b	d	τ	\gg	h
$t1$	$t4$	$t3$	$t5$	$t7$		$t10$

move on
model only

move on log
only

Example: BPI Challenge 2012

(Dutch financial institute, doi:10.4121/uuid:3926db30-f712-4394-aebc-75976070e91f)

Loops of “W_Completeren aanvraag” and “W_Nabellen offertes” are often performed

“O_DECLINED” and “W_Wijzigen contractgegevens” are often skipped



O_DECLINED+CO
(781/8550)

W_Wijzigen
contractgegevens
(4/6739)

Loops of “W_Completeren aanvraag” and “W_Nabellen offertes” are often performed

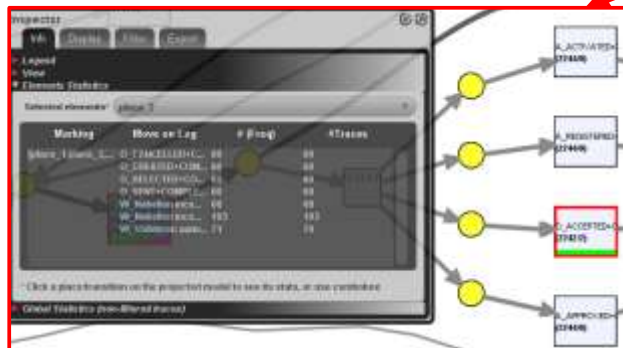
Loops of “W_Completeren aanvraag” and “W_Nabellen offertes” are often performed

“O_DECLINED” and “W_Wijzigen contractgegevens” are often skipped

“O_DECLINED” and “W_Wijzigen contractgegevens” are often skipped

“O_DECLINED” and “W_Wijzigen contractgegevens” are often skipped

Many moves on log of “O_CANCELLED”, “O_CREATED”, “O_SELECTED”, “O_SENT” occurred with the same frequency value (i.e. 60) before parallel branch



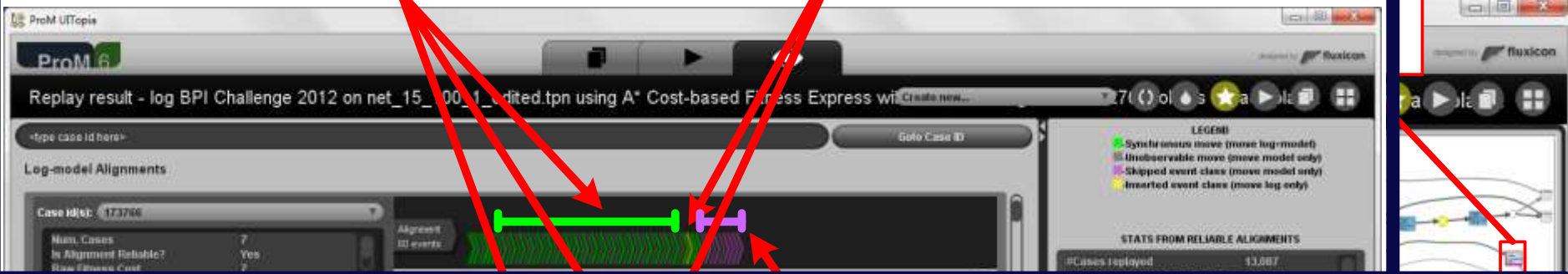
Marking	Move on Log	# (Freq)	#Traces
place_11	A_ACCEPTED+COMPLETE	39	38
place_11	A_PREACCEPTED+COMPLETE	481	381
place_11	W_Afhandelen leads+SCHEDULE	2321	2321
place_11	W_Afhandelen leads+START	2554	2321
place_11	W_Completeren aanvraag+COMPLETE	62	62
place_11	W_Completeren aanvraag+SCHEDULE	421	421
place_11	W_Wijzigen contractgegevens+START	578	421

Marking	Move on Log	# (Freq)	#Traces
place_12	A_ACCEPTED+COMPLETE	60	16
place_12	A_CANCELLED+COMPLETE	1067	1067
place_12	A_PREACCEPTED+COMPLETE	89	89
place_12	A_SELECTED+COMPLETE	156	156
place_12	O_CANCELLED+COMPLETE	524	524
place_12	O_CREATED+COMPLETE	54	54
place_12	W_Afhandelen leads+COMPLETE	2233	2225

Many moves on log of “W_Afhandelen leads” (> 2200 times) occurred in the end of traces

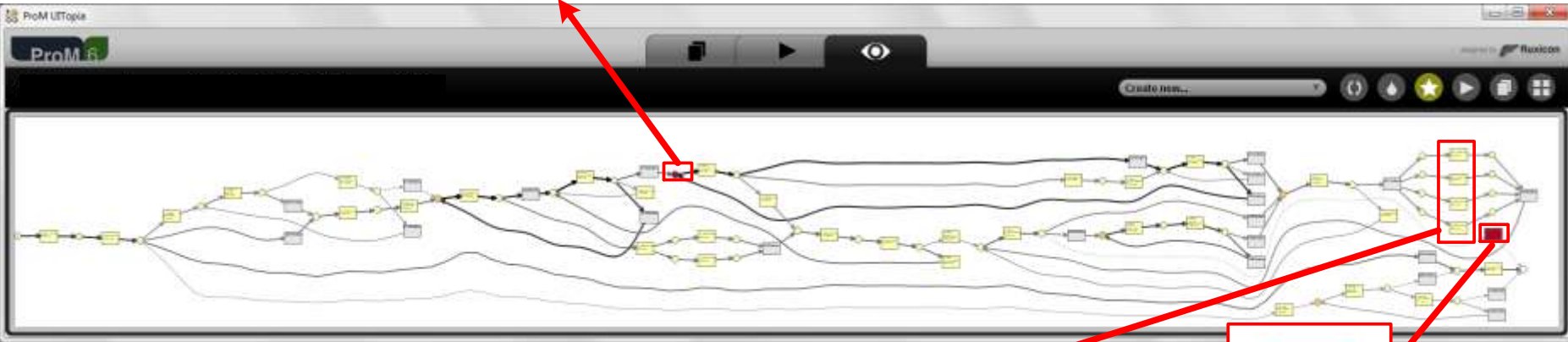
Synchronous moves of
"Completeren aanvraag"

Move on log of "Completeren aanvraag"

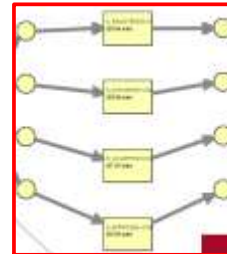


Property	Min.	Max.	Avg.	Std. Dev	Freq.
Waiting time	0.00 ms	29.78 days	2.83 days	3.30 days	24,229
Synchronization time	0.00 ms	0.00 ms	0.00 ms	0.00 ms	24,229
Sojourn time	0.00 ms	29.78 days	2.83 days	3.30 days	24,229

The average waiting time for the input place of "W_Nabellen offertes+START" is very long (2.83 days) compares to the average waiting time of other places

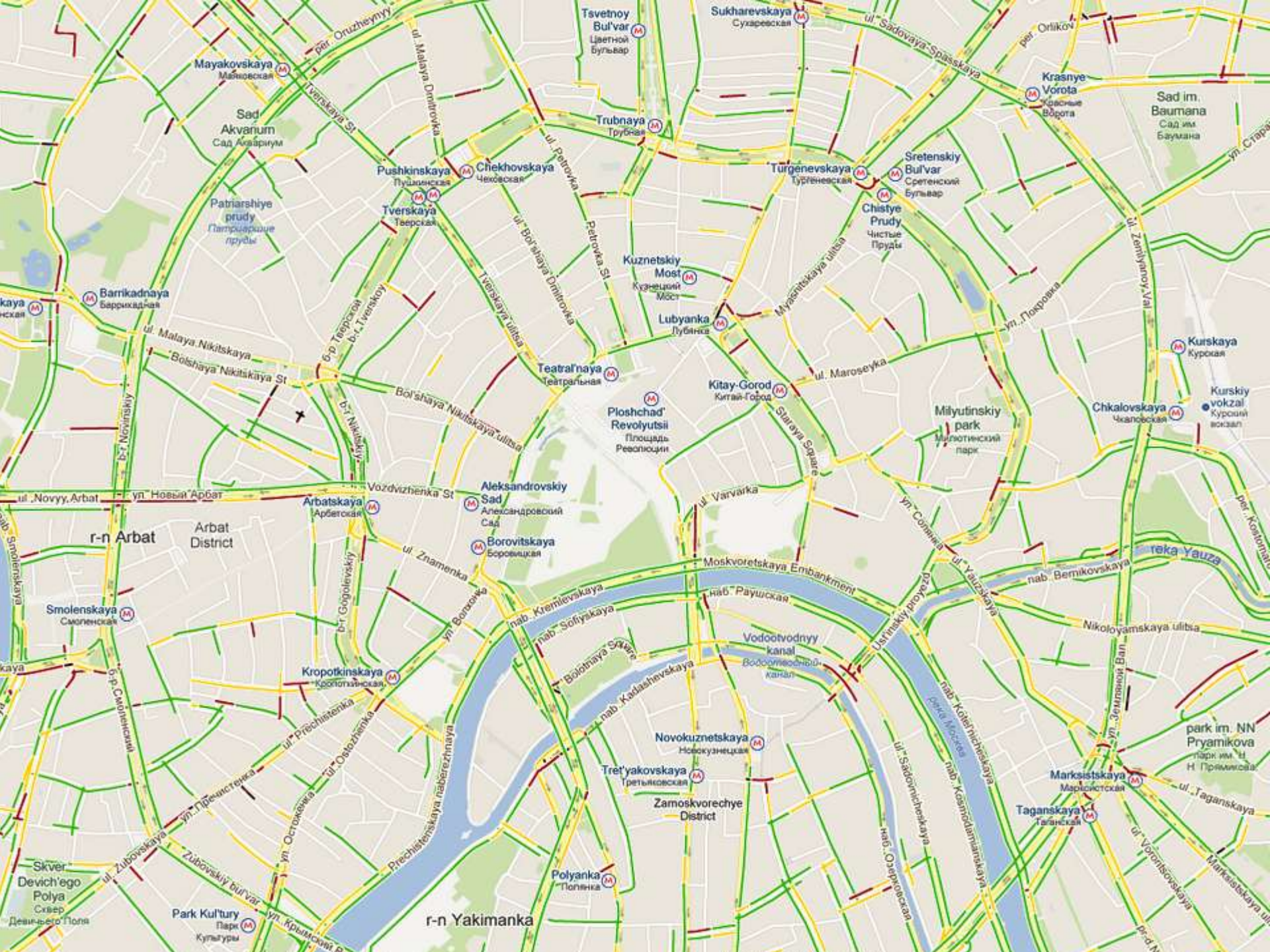


"O_ACCEPTED" has average sojourn time of 27.07 minutes, while "A_REGISTERED", "A_ACTIVATED", and "A_APPROVED" have average sojourn time of 29.56 minutes



Property	Min.	Max.	Avg.	Std. Dev	Freq.
Throughput time	0.00 ms	0.00 ms	0.00 ms	0.00 ms	4
Waiting time	1.55 hours	3.43 months	1.14 months	1.55 months	4
Sojourn time	1.55 hours	3.43 months	1.14 months	1.55 months	4
#Unique cases ...	4				

Activity "W_Wijzigen contractgegevens" is the bottleneck, but it occurred rarely (only 4 times)



process mining as
the missing link



aligning model
and reality



divide and
conquer



process
discovery



Big (Event)
Data



challenges



getting
started



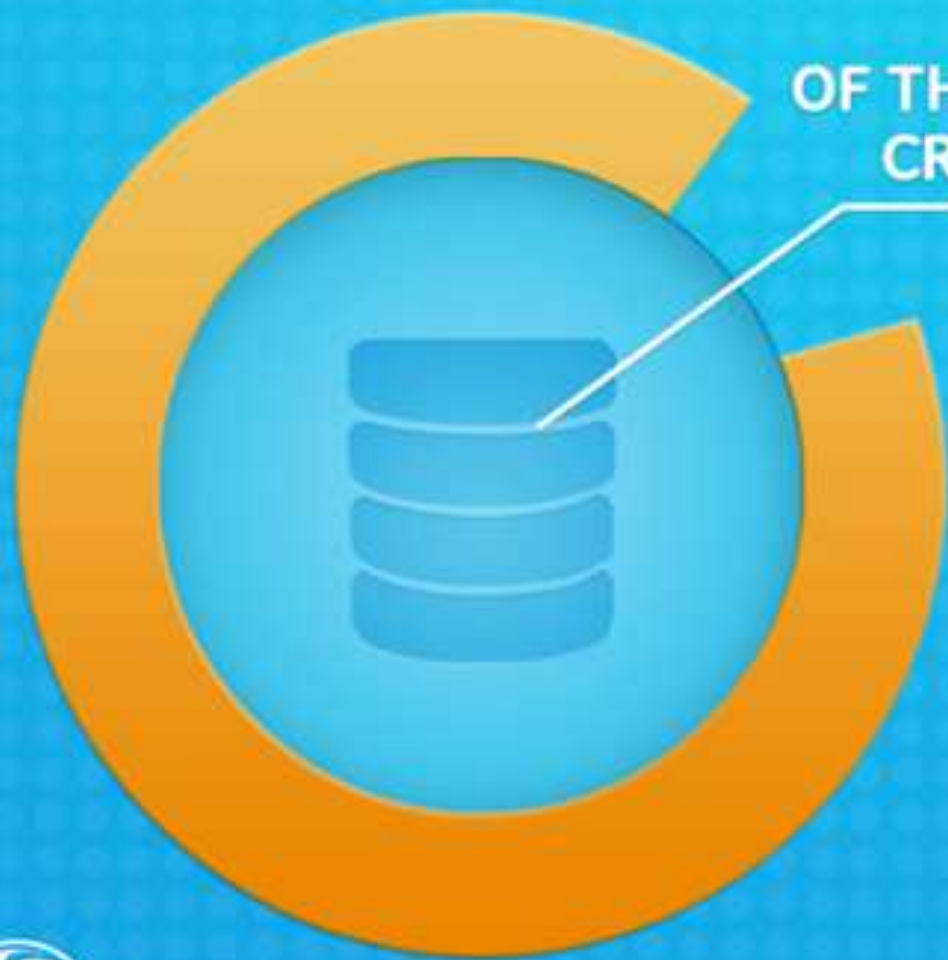
The Digital Universe: 50-fold Growth from the Beginning of 2010 to the End of 2020

In 10 years we will have 50 times as much data!

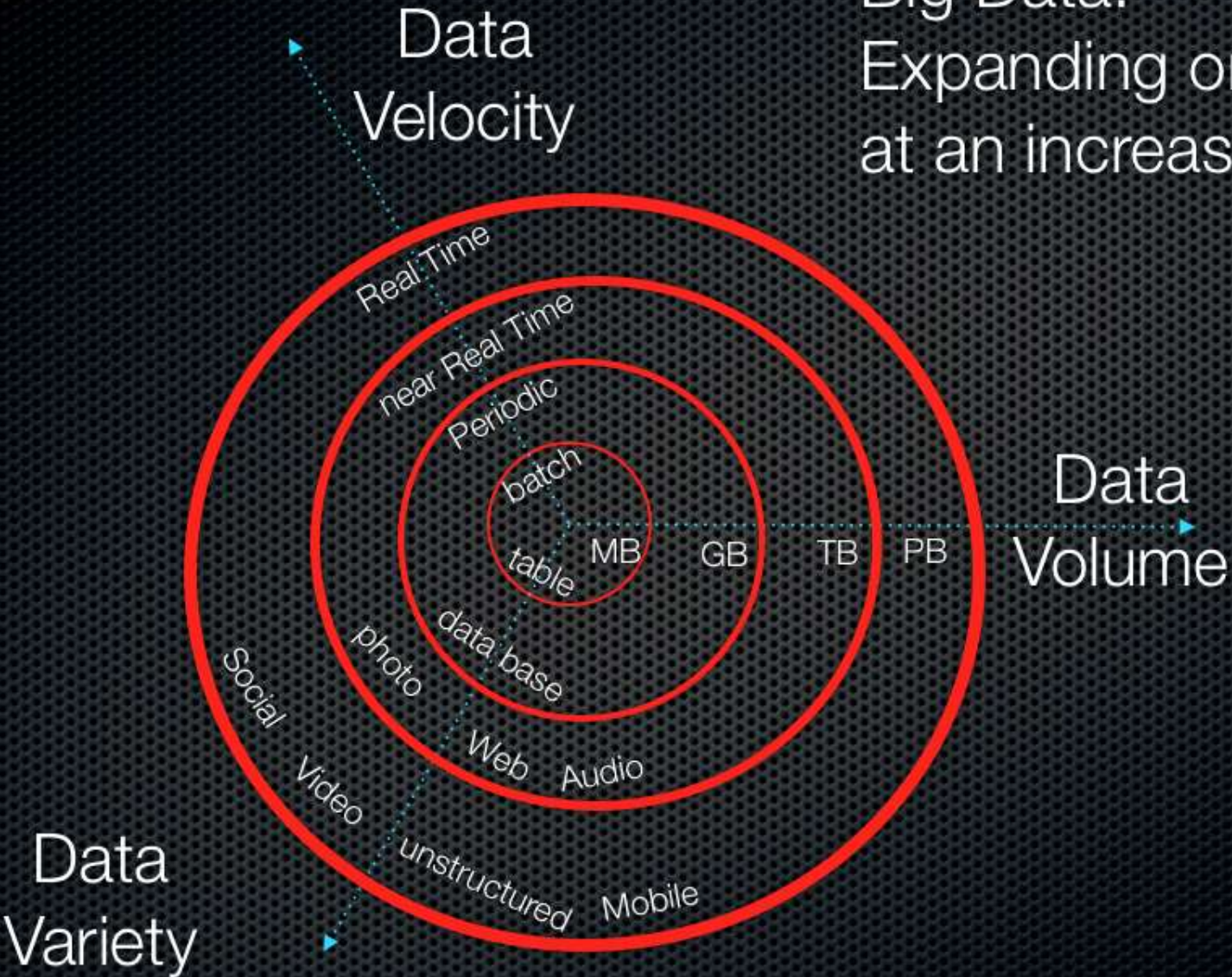


90%

OF THE WORLD'S DATA WAS
CREATED IN THE LAST
TWO YEARS



Big Data:
Expanding on 3 fronts
at an increasing rate.



Acquisition

Marshalling

Analysis

Action

Data Acquisition



VLDW and BI Appliances



Analytics



BPM & Action



Including Complex Event Processing (CEP) tools

Data Providers



And all your own data
And your partners data

No SQL



Content Management



Data Virtualization COMPOSITE SOFTWARE



BI Tools



Data Governance



Capgemini - Capping IT off
Manuel Sevilla - 2012

The Economist

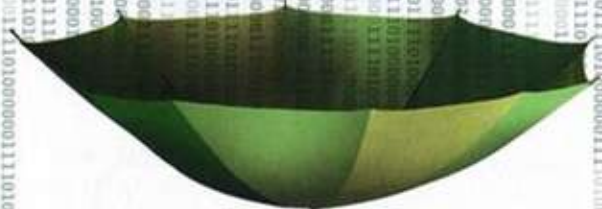
FEBRUARY 27TH - MARCH 5TH 2010

Economist.com

Obama the warrior
Misgoverning Argentina
The economic shift from West to East
Genetically modified crops blossom
The right to eat cats and dogs

The data deluge

AND HOW TO HANDLE IT: A 14-PAGE SPECIAL REPORT





Big Data ?



Big ... or fast and efficient?

process mining as
the missing link



aligning model
and reality



divide and
conquer



process
discovery



Big (Event)
Data



challenges



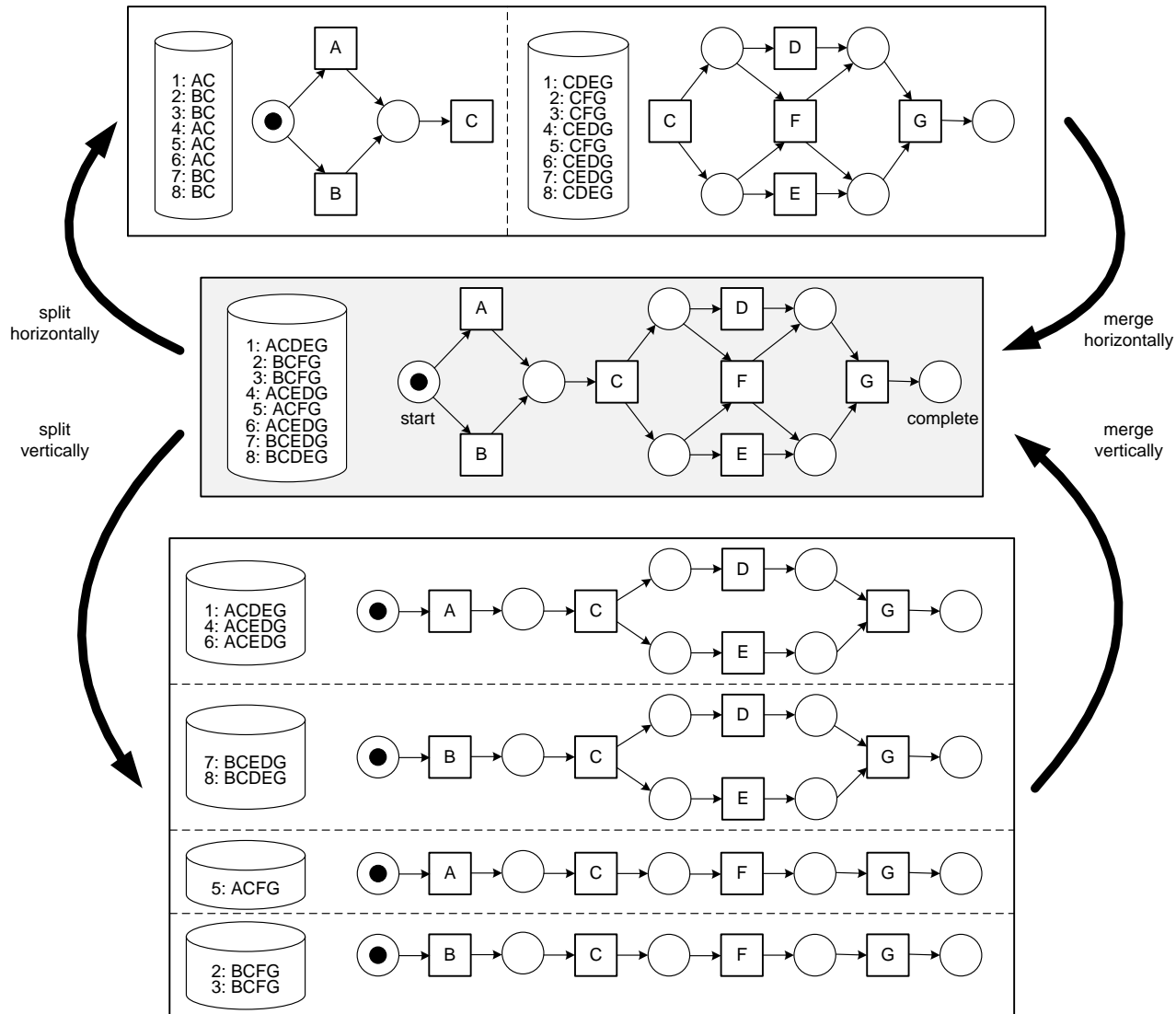
getting
started



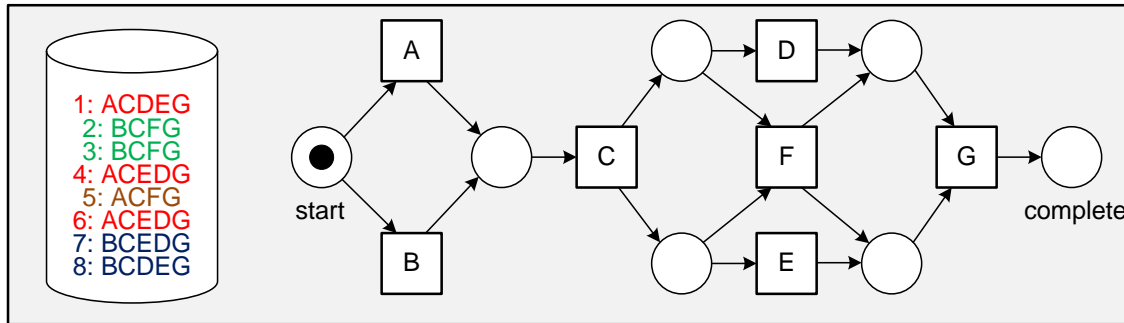
Big Data: Opportunities and Challenges



Divide and Conquer

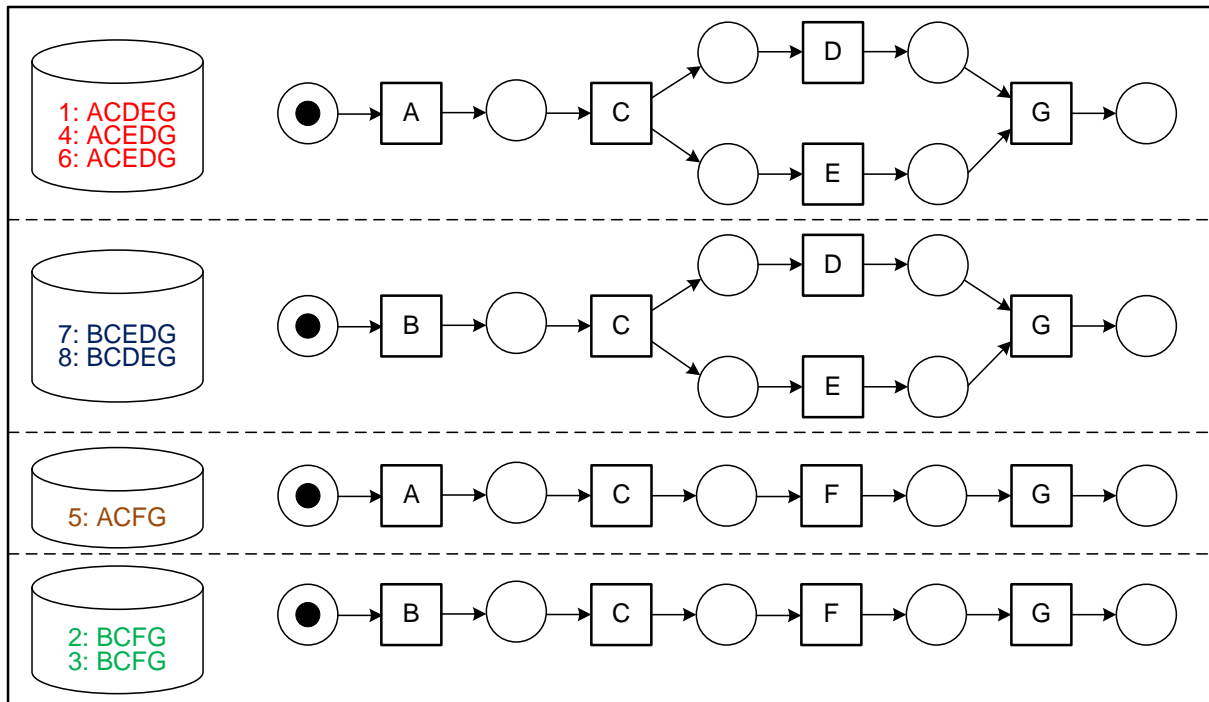


Vertical Decomposition

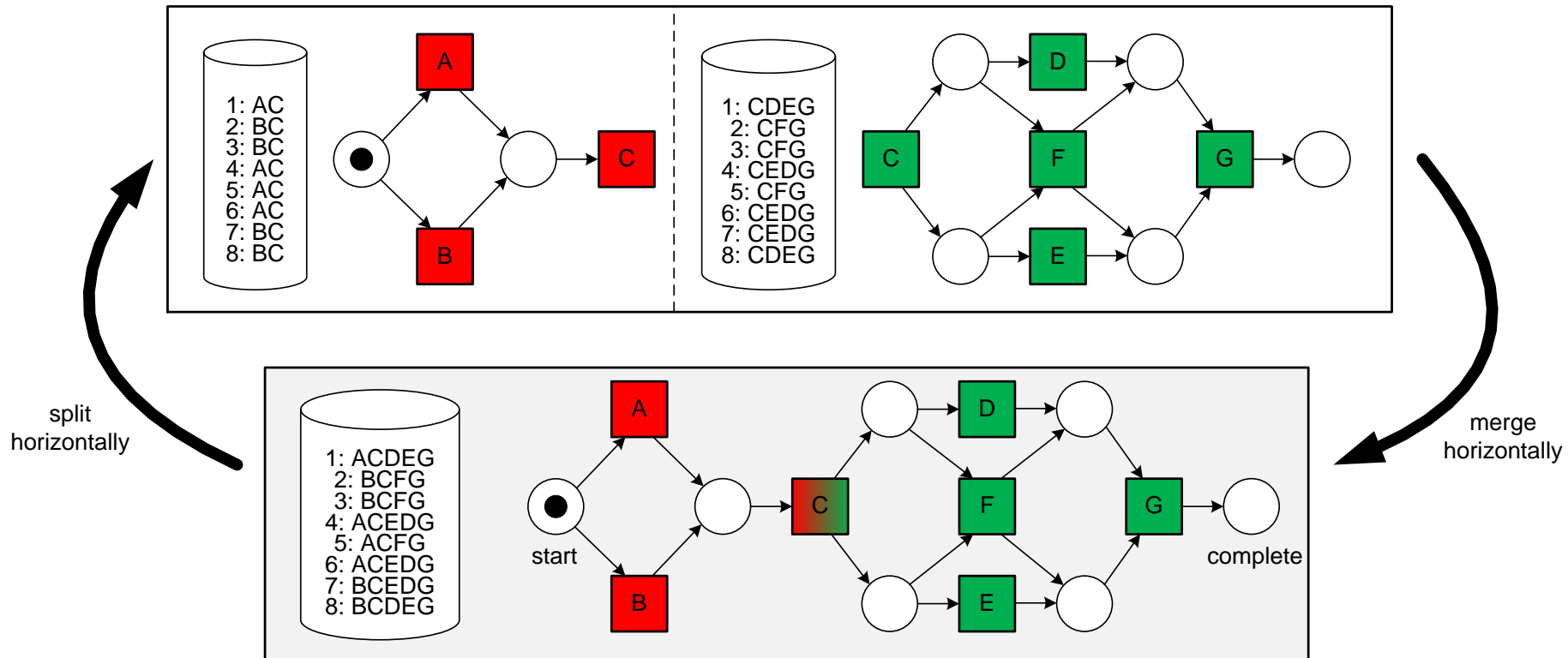


split vertically

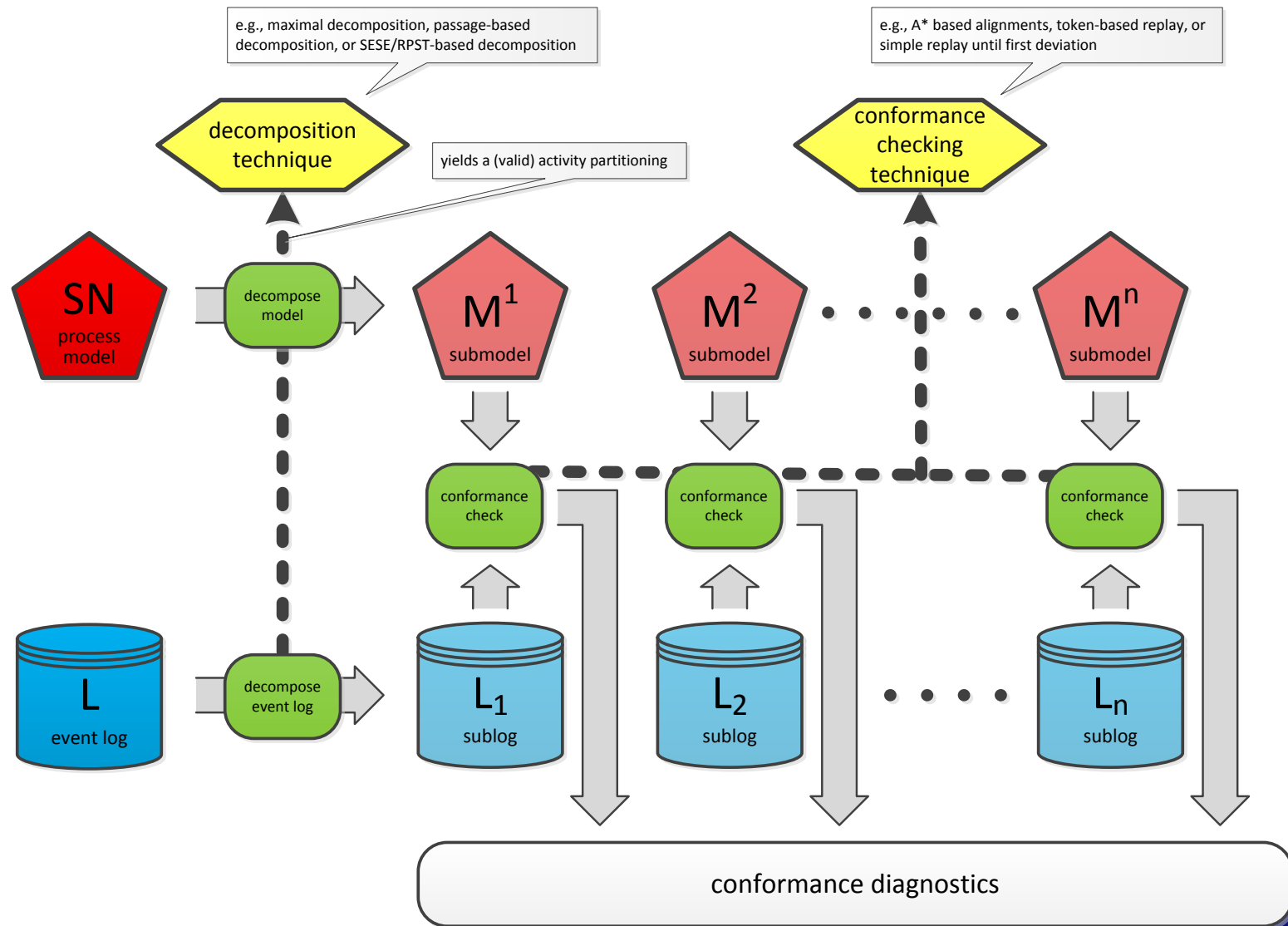
merge vertically



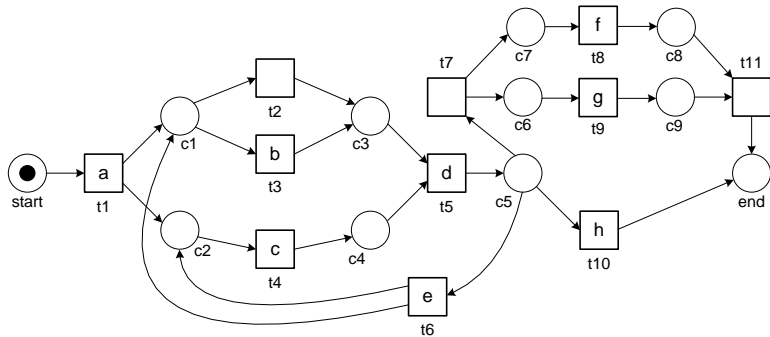
Horizontal Decomposition



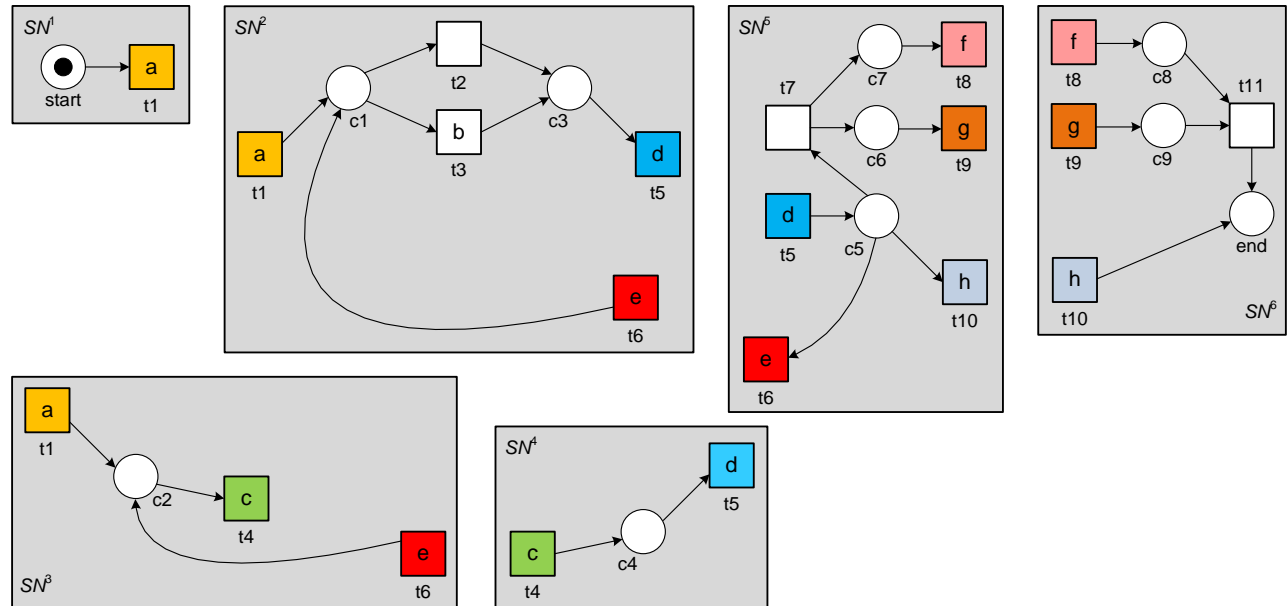
Decomposing Conformance Checking



Example of a valid decomposition

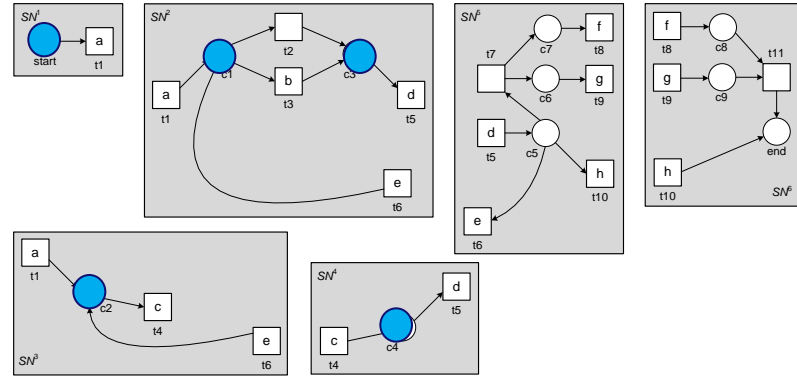
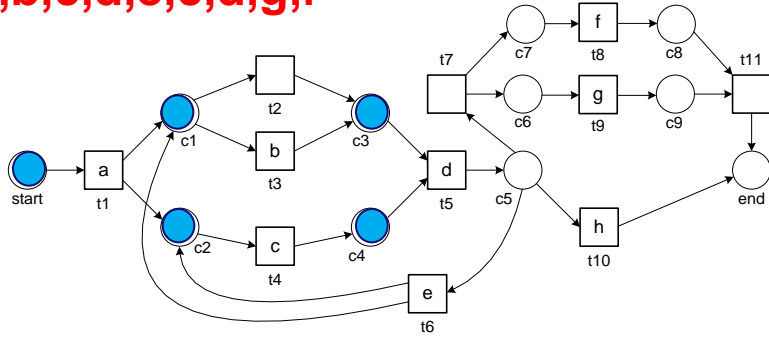


Log can be split in the same way!



Example of alignment for observed trace a,b,c,d,e,c,d,g,f

a,b,c,d,e,c,d,g,f



$\gamma_3 =$

1	2	3	4	5	6	7	8	9	10	11	12
a	b	c	d	e	c	\gg	d	\gg	g	f	\gg
a	b	c	d	e	c	τ	d	τ	g	f	τ
t1	t3	t4	t5	t6	t4	t2	t5	t7	t9	t8	t11

Etc.

$\gamma_3^1 =$

1
a
a
t1

 $\gamma_3^2 =$

1	2	4	5	7	8
a	b	d	e	\gg	d
a	b	d	e	τ	d
t1	t3	t5	t6	t2	t5

 $\gamma_3^3 =$

1	3	5	6
a	c	e	e
a	c	e	e
t1	t4	t6	t4

$\gamma_3^4 =$

3	4	6	8
c	d	c	d
c	d	c	d
t4	t5	t4	t5

 $\gamma_3^5 =$

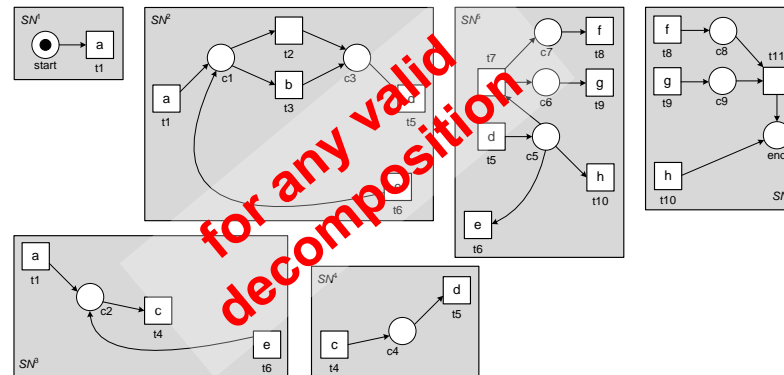
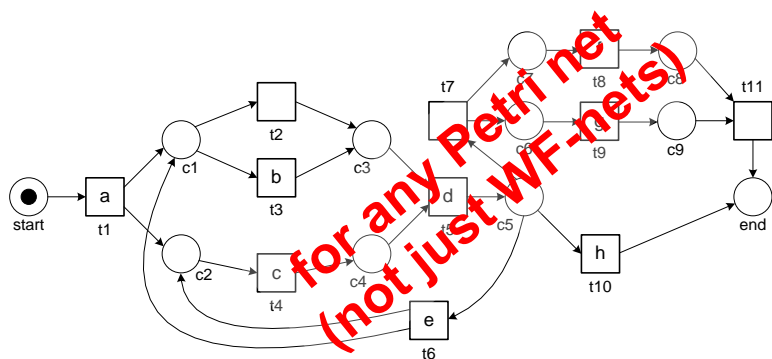
4	5	8	9	10	11
d	e	d	\gg	g	f
d	e	d	τ	g	f
t5	t6	t5	t7	t9	t8

 $\gamma_3^6 =$

10	11	12
g	f	\gg
g	f	τ
t9	t8	t11

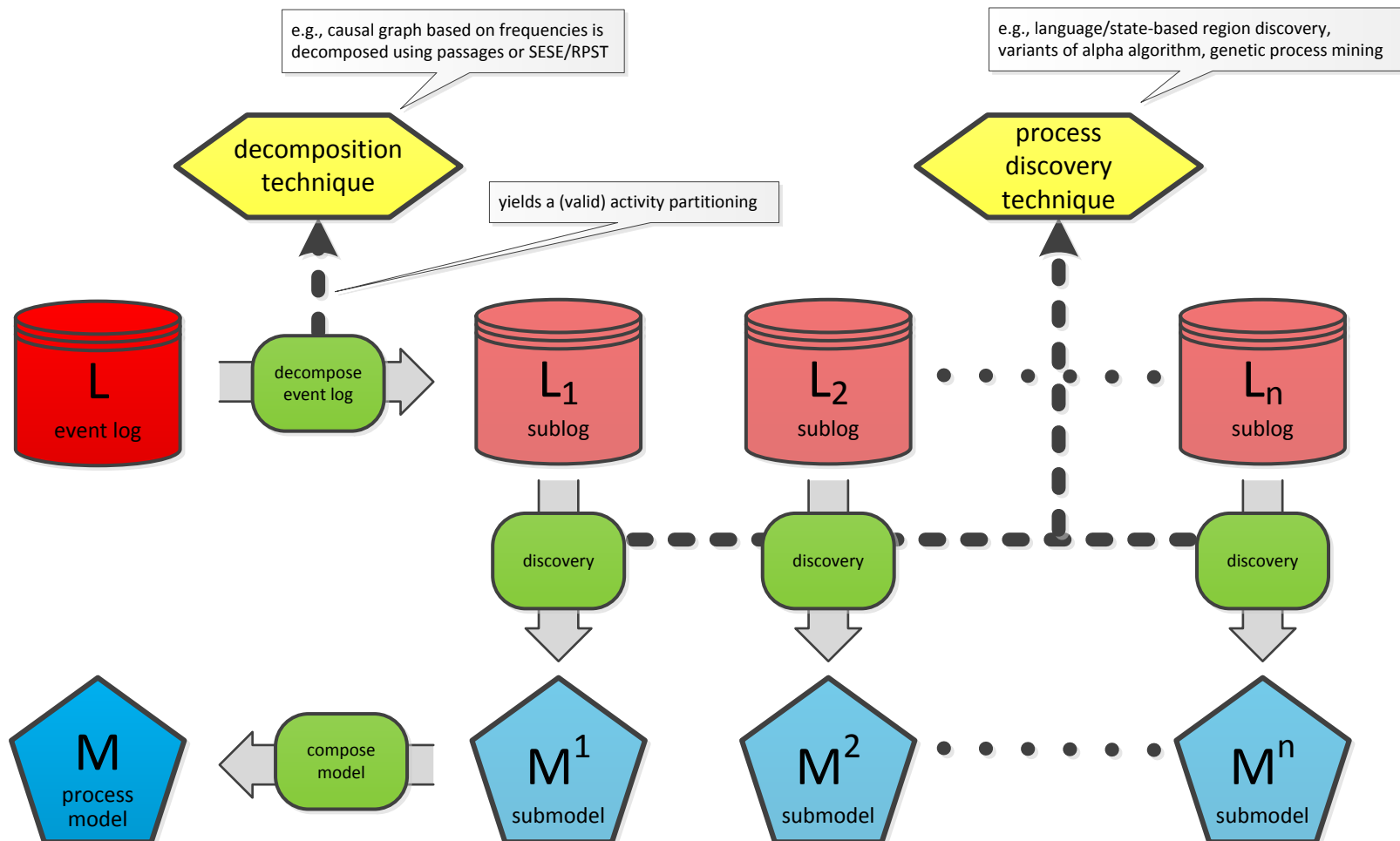
Conformance checking can be decomposed !!!

- **General result for any valid decomposition: Any event log or trace is perfectly fitting the overall model if and only if it is also fitting all the individual fragments**



Wil van der Aalst, Decomposing Petri nets for process mining: A generic approach. Distributed and Parallel Databases, Volume 31, Issue 4, pp 471-507, 2013

Decomposing Process Discovery



Learn more about decomposing process mining problems?

- **W.M.P. van der Aalst. Decomposing Petri Nets for Process Mining: A Generic Approach. *Distributed and Parallel Databases*, 31(4):471-507, 2013.**
- W.M.P. van der Aalst. A General Divide and Conquer Approach for Process Mining. In M. Ganzha, L. Maciaszek, and M. Paprzycki, editors, *Federated Conference on Computer Science and Information Systems (FedCSIS 2013)*, pages 1-10. IEEE Computer Society, 2013.
- W.M.P. van der Aalst. Decomposing Process Mining Problems Using Passages. In S. Haddad and L. Pomello, editors, *Applications and Theory of Petri Nets 2012*, volume 7347 of *Lecture Notes in Computer Science*, pages 72-91. Springer-Verlag, Berlin, 2012.
- J. Munoz-Gama, J. Carmona, and W.M.P. van der Aalst. Hierarchical Conformance Checking of Process Models Based on Event Logs. In J.M. Colom and J. Desel, editors, *Applications and Theory of Petri Nets 2013*, volume 7927 of *Lecture Notes in Computer Science*, pages 291-310. Springer-Verlag, Berlin, 2013.
- J. Munoz-Gama, J. Carmona, and W.M.P. van der Aalst. Conformance Checking in the Large: Partitioning and Topology. In F. Daniel, J. Wang, and B. Weber, editors, *International Conference on Business Process Management (BPM 2013)*, volume 8094 of *Lecture Notes in Computer Science*, pages 130-145. Springer-Verlag, Berlin, 2013.
- E. Verbeek and W.M.P. van der Aalst. Decomposing Replay Problems: A Case Study. In D. Moldt and H. Roelke, editors, *Proceedings of the International Workshop on Petri Nets in Software Engineering (PNSE 2013)*, volume 989 of *CEUR Workshop Proceedings*, pages 213-232. CEUR-WS.org, 2013.

process mining as
the missing link



aligning model
and reality



divide and
conquer



process
discovery



Big (Event)
Data

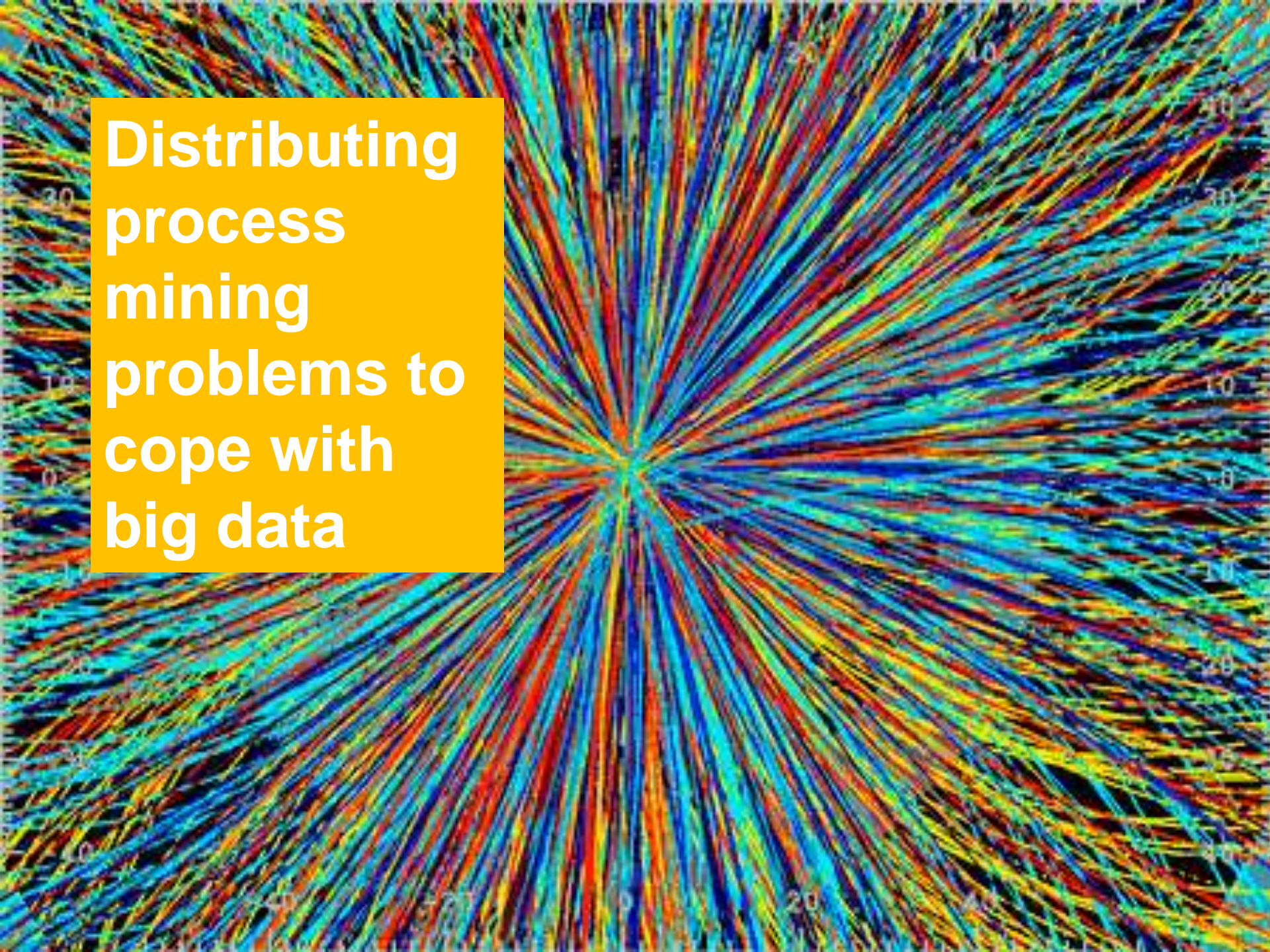


challenges




getting
started





**Distributing
process
mining
problems to
cope with
big data**

A close-up photograph of a red fire hydrant with a green top. Water is spraying out from the side outlet on the left. The hydrant has a red cap on the right side and a chain hanging from it. The background is dark and out of focus.

streaming event data

(sensors, RFID, messages, etc.)



**process discovery: finding
sheep with five or more legs**

1

**formal
(not just a
picture)**

2

**fast
(should not
take years)**

**ability to balance
all conformance
dimensions
(fitness, precision,
generalization, and
simplicity) incl.
noise**

3

4

**sound
(result should
at least be free
of deadlocks,
etc.)**

5

**provide
guarantees
(not just a best
effort)**

**On-the-fly
process mining**



**Operational
support**

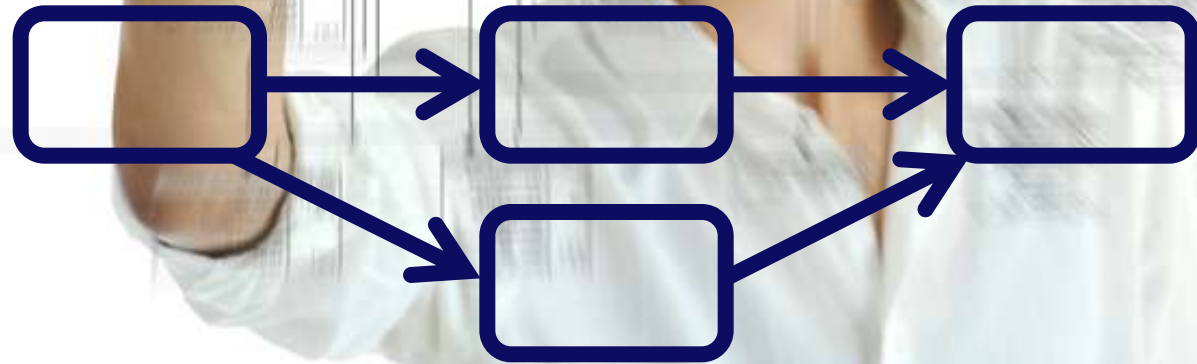
Concept drift





**cross-organizational /
comparative process mining**

Supporting the process of process mining



process mining as
the missing link



aligning model
and reality



divide and
conquer



process
discovery



Big (Event)
Data



challenges



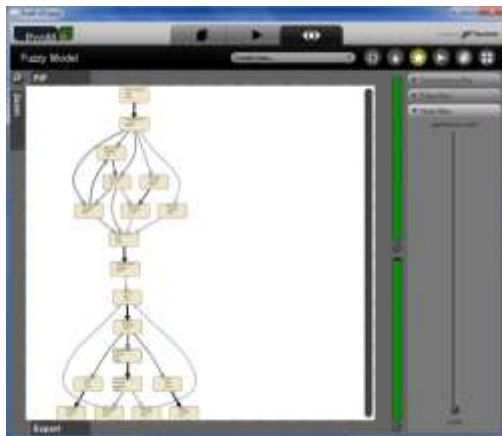
getting
started



How to get started?



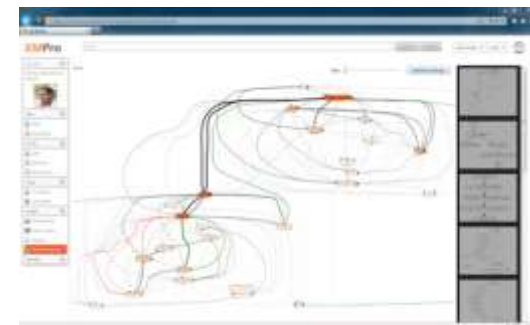
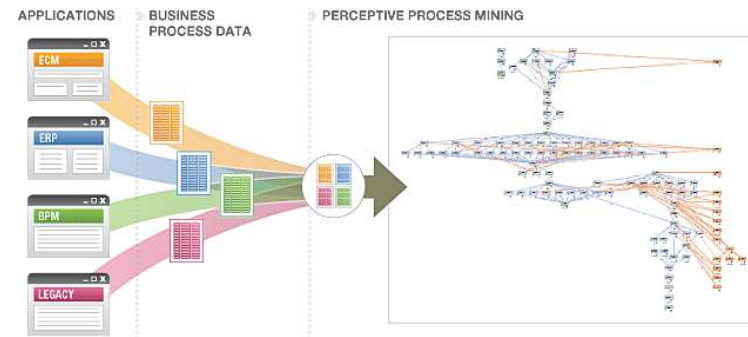
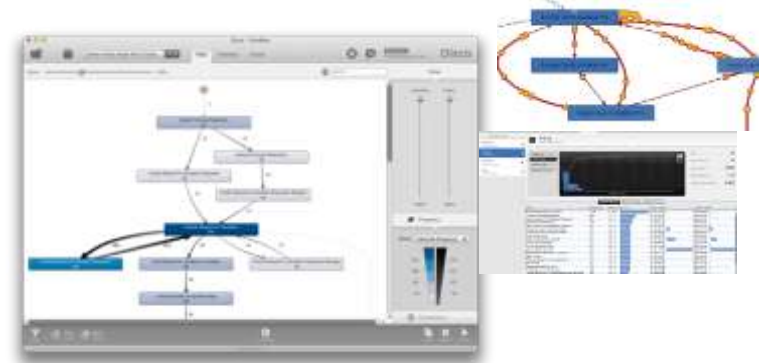
600+ plug-ins available covering the whole process mining spectrum



Download from: www.processmining.org

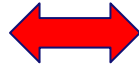
Commercial Alternatives

- **Disco (Fluxicon)**
- **Perceptive Process Mining**
(before Futura Reflect and BPM|one)
- **ARIS Process Performance Manager**
- **QPR ProcessAnalyzer**
- **Interstage Process Discovery (Fujitsu)**
- **Discovery Analyst (StereoLOGIC)**
- **XMAalyzer (XMPPro)**
- ...



How to Get Started?

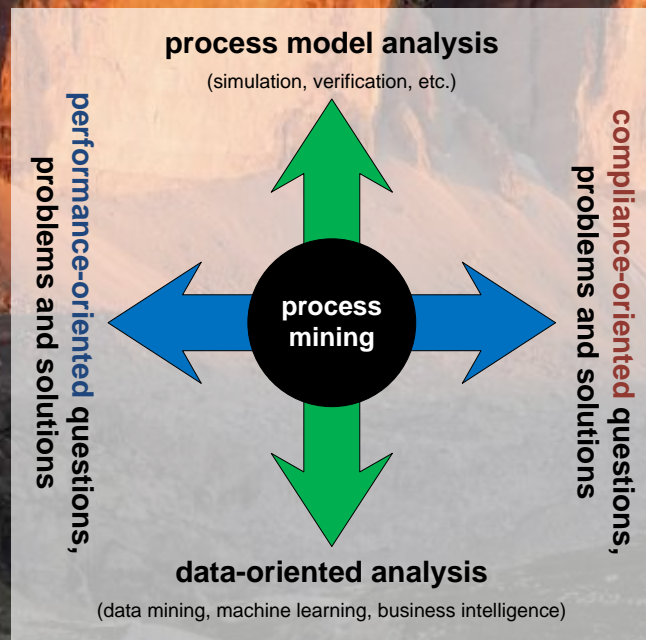
Collect event data



Collect questions

- **Minimal requirement:** events referring to an activity name and a process instance.
 - **Good to have:** timestamps, resource information, additional data elements.
 - **Challenges:** scoping and sometimes correlation.
- **What kind problems would you like to address (cost, time, risk, compliance, service, etc.)?**
 - **Related to discovery, conformance, enhancement?**
 - **Iterative process:** can be “curiosity driven” initially.

Join our expedition: Mine your processes!



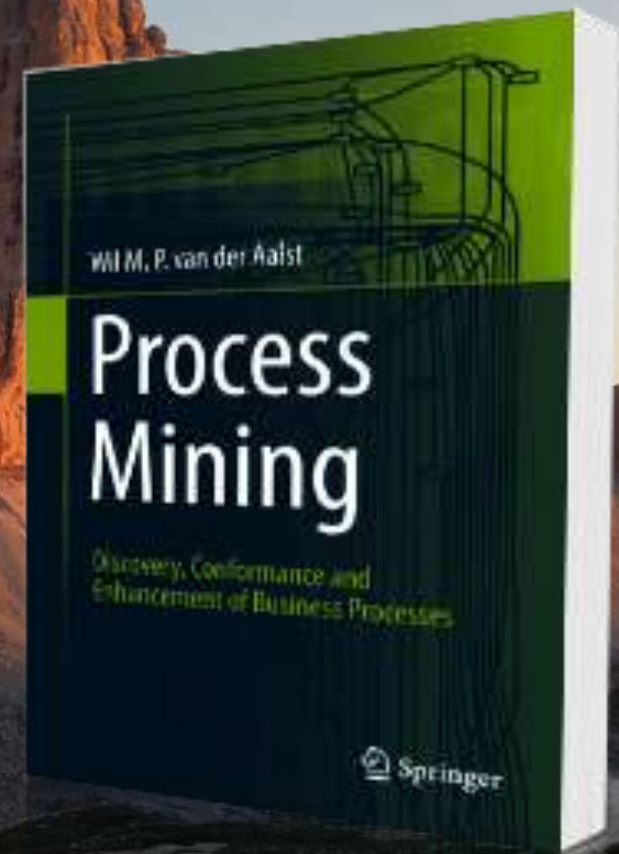
Learn more?



http://www.youtube.com/watch?v=7oat7MatU_U

<http://www.win.tue.nl/ieetfpm/>

**Informal PM Meeting (15.50 today).
Thanks to Krzysztof Kluza!
Building C2 (entrance through buildings
C1 or C3), room no. 316 (3rd floor).**



processmining.org