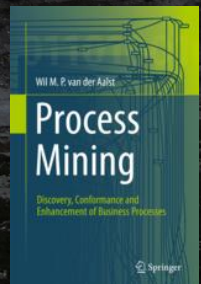


Mediating Between Modeled and Observed Behavior

The Quest for the "Right" Process

ROIS
2013

prof.dr.ir. Wil van der Aalst



Why is process discovery difficult?

How about precision and recall?

What is process mining?

What are the main research challenges?

How to measure the quality of a process model?

The future is bright, but how to get started?

What are the main pitfalls of process modeling?



Why is process discovery difficult?



How about precision and recall?



What are the main research challenges?



How to measure the quality of a process model?



The future is bright, but how to get started?

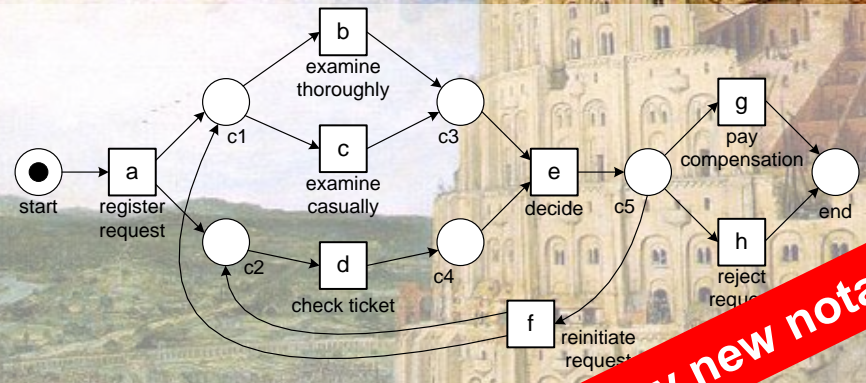
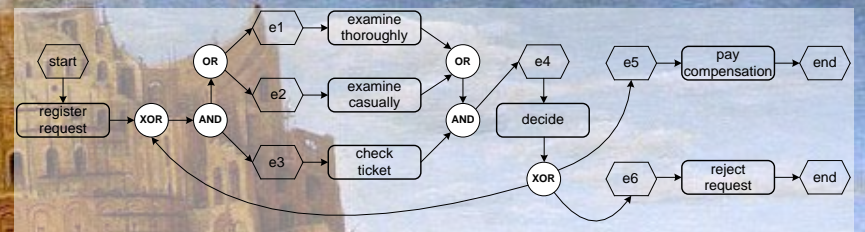


What is process mining?

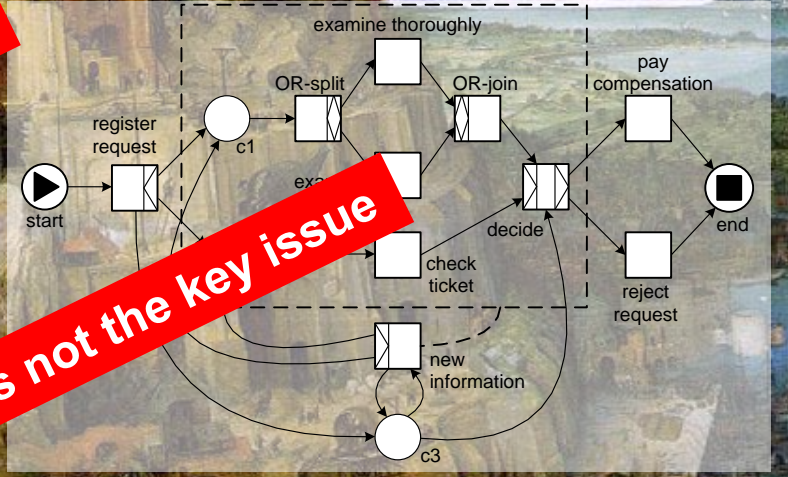


What are the main pitfalls of process modeling?

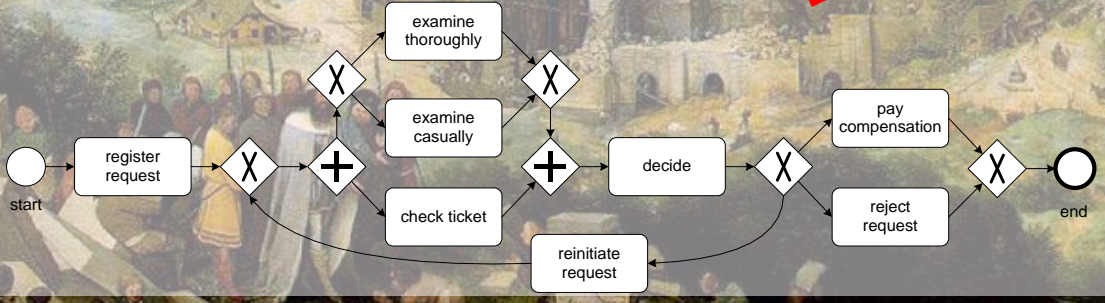




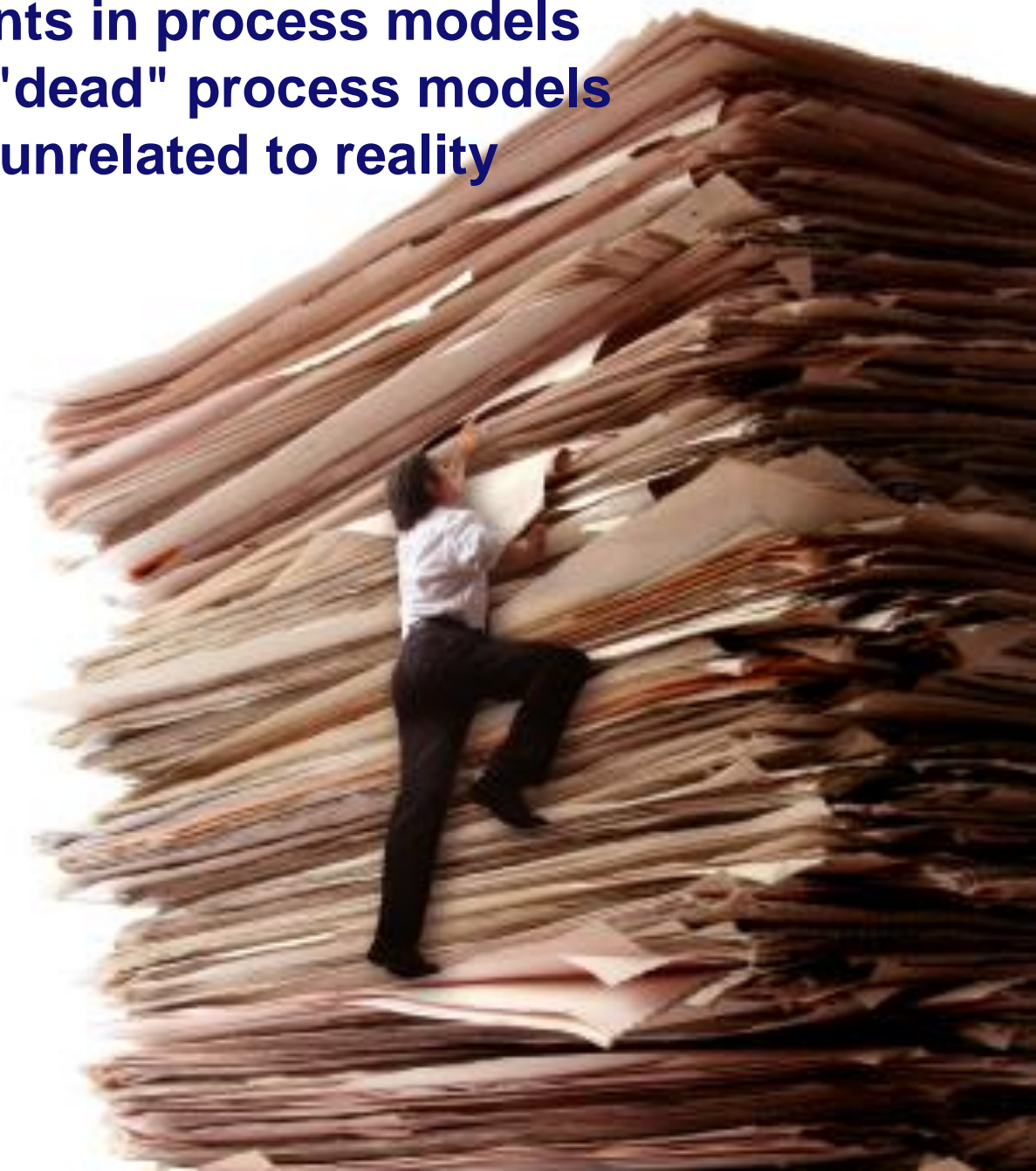
not solved by new notations



notation is not the key issue



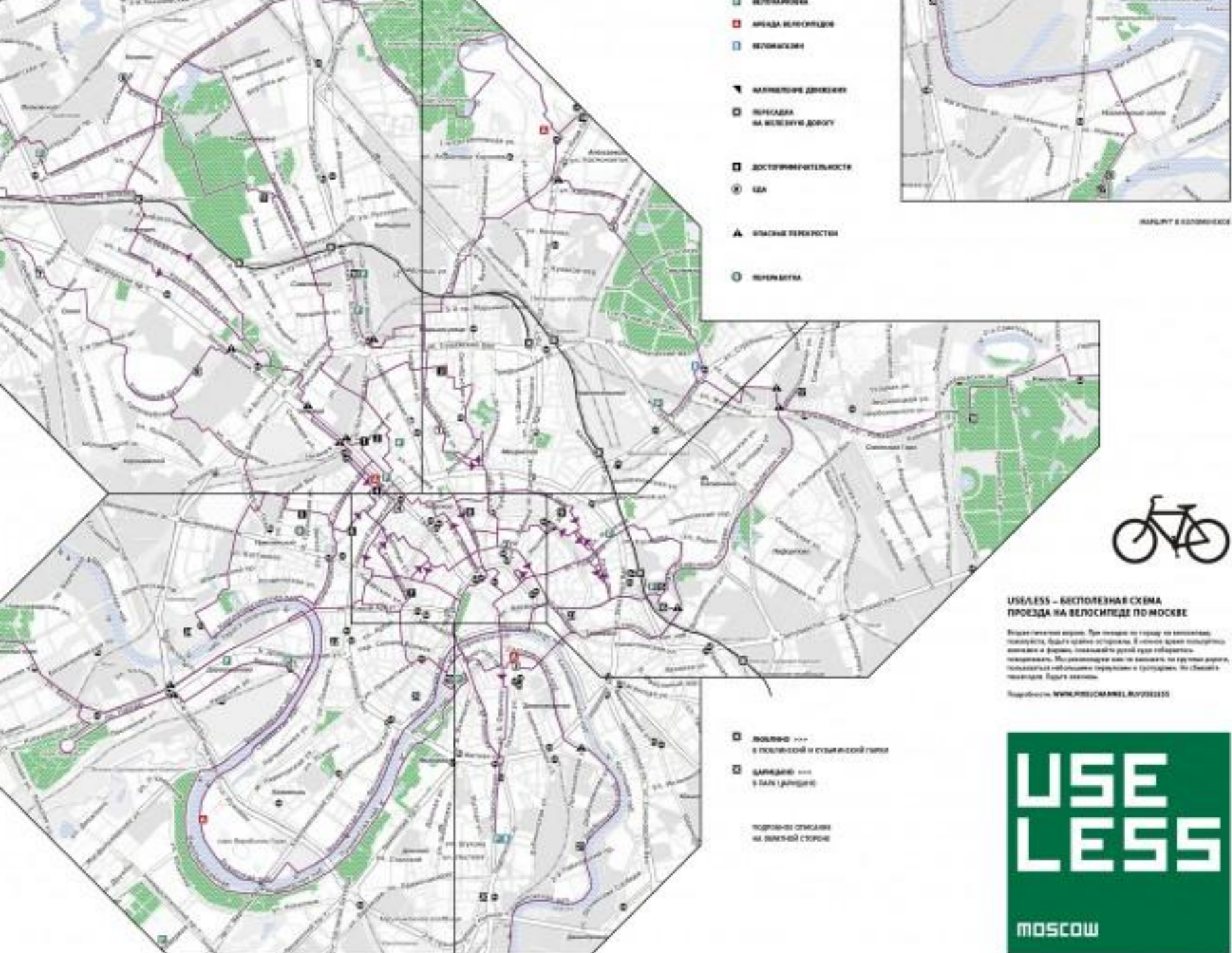
- **enormous investments in process models**
- **large collections of "dead" process models**
- **not taken seriously, unrelated to reality**





problem #1

Aiming for one model that suits all purposes



- ВЕЛОПАРКИ
- АРКАДА ВЕЛОСИДЕ
- ВЕЛОМОСКИ
- ▼ МАГНИТНЫЕ ДВИЖУКИ
- ПЕРЕСАДА НА ВЕЛОСИДУ ДОРОГУ
- ДОСТОЙНИМЫЕ
- ЦБА
- ▲ ПЛОСКИЕ ПЕРЕСЕЧЕНИЯ
- ПЕРЕСЕЧЕНИЯ



КАРТА МОСКОВСКОЙ ОБЛАСТИ



USELESS – ВЕЛОСИДНАЯ СХЕМА ПЕРЕСАДА НА ВЕЛОСИДЕ ПО МОСКВЕ

Визуализация схемы. При выборе по городу по велосипеду, маршруты, чтобы избежать столкновений, в течение дня поделены на зоны и формируются в виде карты. Визуализация схемы по городу по велосипеду, маршруты, чтобы избежать столкновений, в течение дня поделены на зоны и формируются в виде карты.

Подробнее: WWW.USELESS.MOSCOW

- ПОЛОСЫ — в пешеходной и спальной зонах
- БУЛЬВАРЫ — в парках (дети)

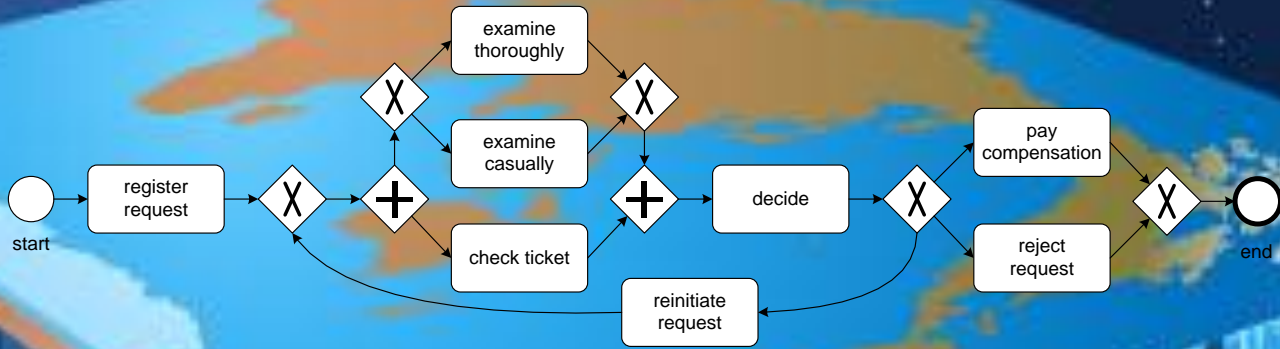
подробные описания на обратной стороне



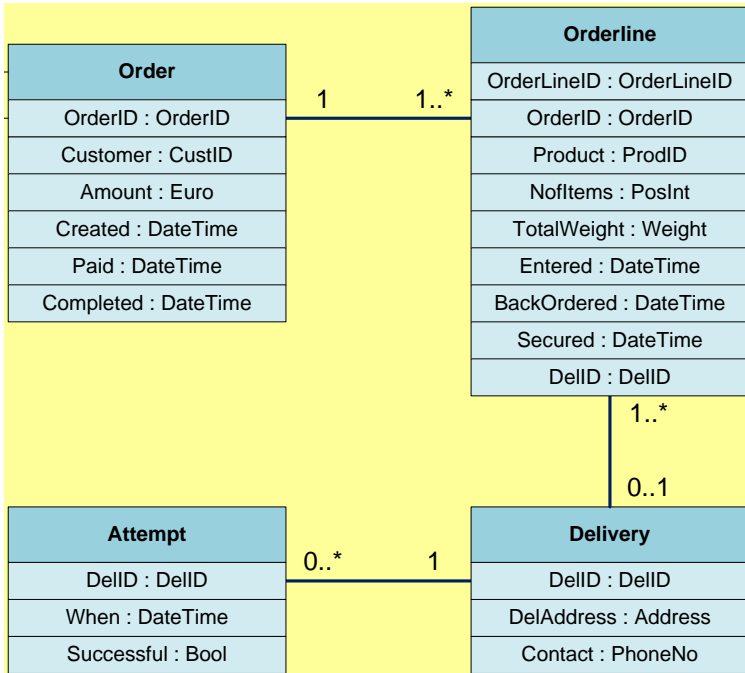


problem #2

**Straightjacketing smaller interacting processes
into one monolithic model**



What is the process instance?

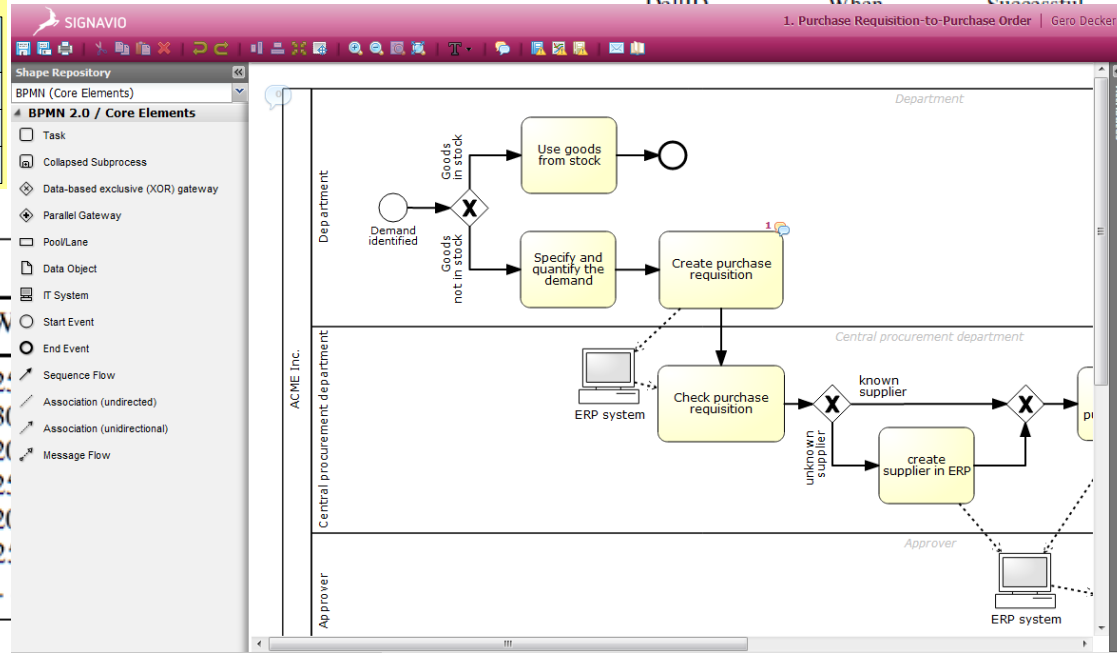


Order					
OrderID	Customer	Amount	Created	Paid	Completed
91245	John	100	28-11-2011:08.12	02-12-2011:13.45	05-12-2011:11.33
91561	Mike	530	28-11-2011:12.22	03-12-2011:14.34	05-12-2011:09.32
91812	Mary	234	29-11-2011:09.45	02-12-2011:09.44	04-12-2011:13.33
92233	Sue	110	29-11-2011:10.12	null	null
92345	Kirsten	195	29-11-2011:14.45	02-12-2011:13.45	null
92355	Pete	320	29-11-2011:16.32	null	null
...

OrderLineID	OrderID	Product	NofItems	TotalW
112345	91245	iPhone 4G	1	0.2
112346	91245	iPod nano	2	0.3
112347	91245	iPod classic	1	0.2
112448	91561	iPhone 4G	1	0.2
112449	91561	iPod classic	1	0.2
112452	91812	iPhone 4G	5	1.2
...

Delivery

Attempt





problem #3

**Using a static hierarchical decomposition
as the only abstraction mechanism**



**most process modeling notations assume a fixed hierarchy
no seamless zoom-in and zoom out!**

**traditional hierarchy concepts
don't support "Google Maps" abstraction**

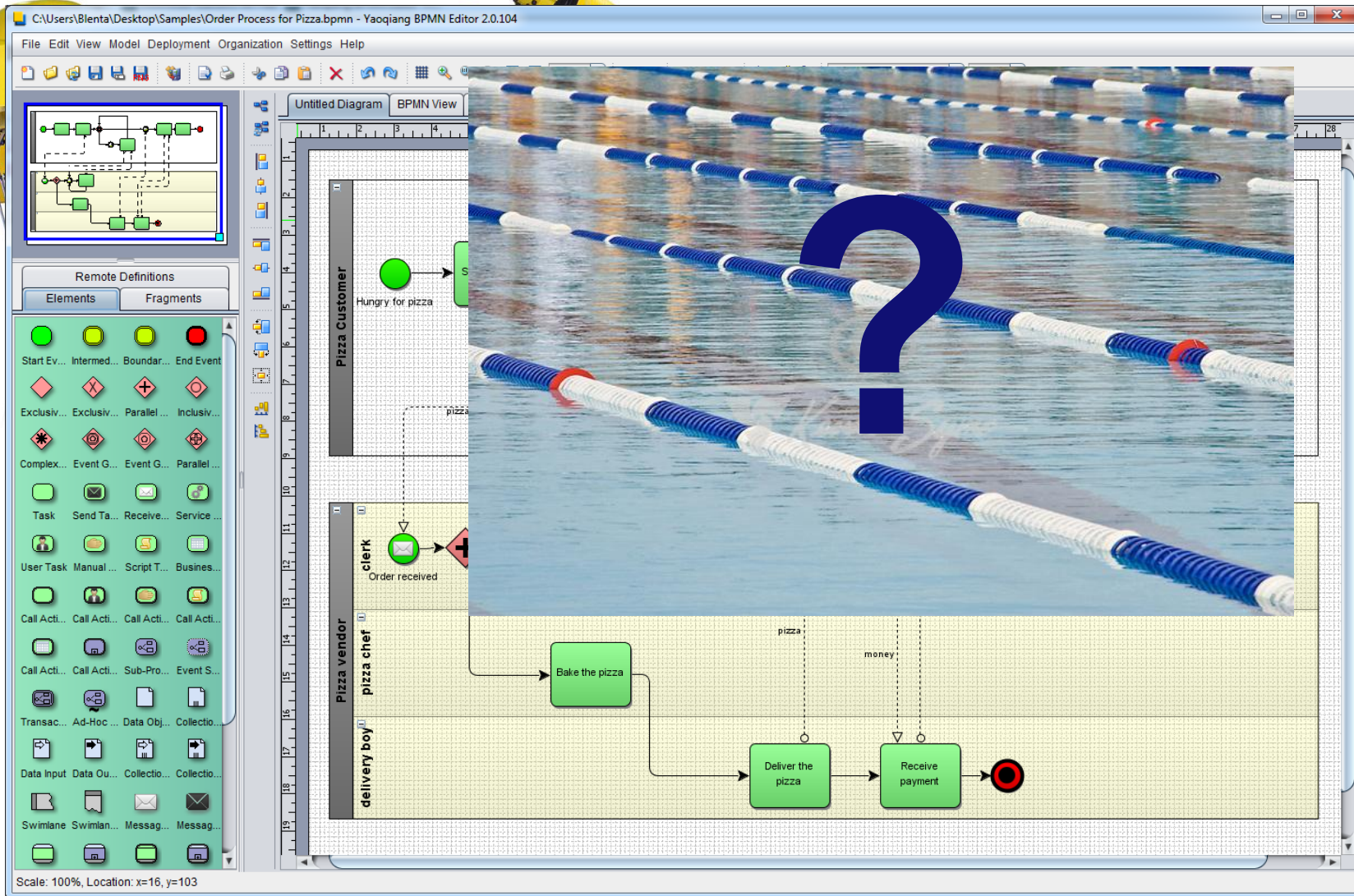


problem #4

**Modeling humans as if they are machines
doing a single task**

**"My processes are unique,
my people are artists!"**



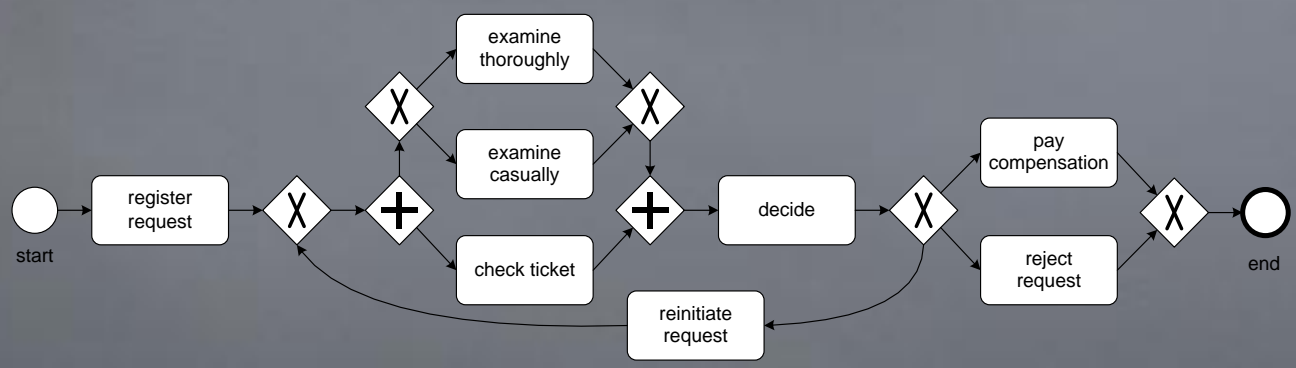




problem #5

Being vague about vagueness

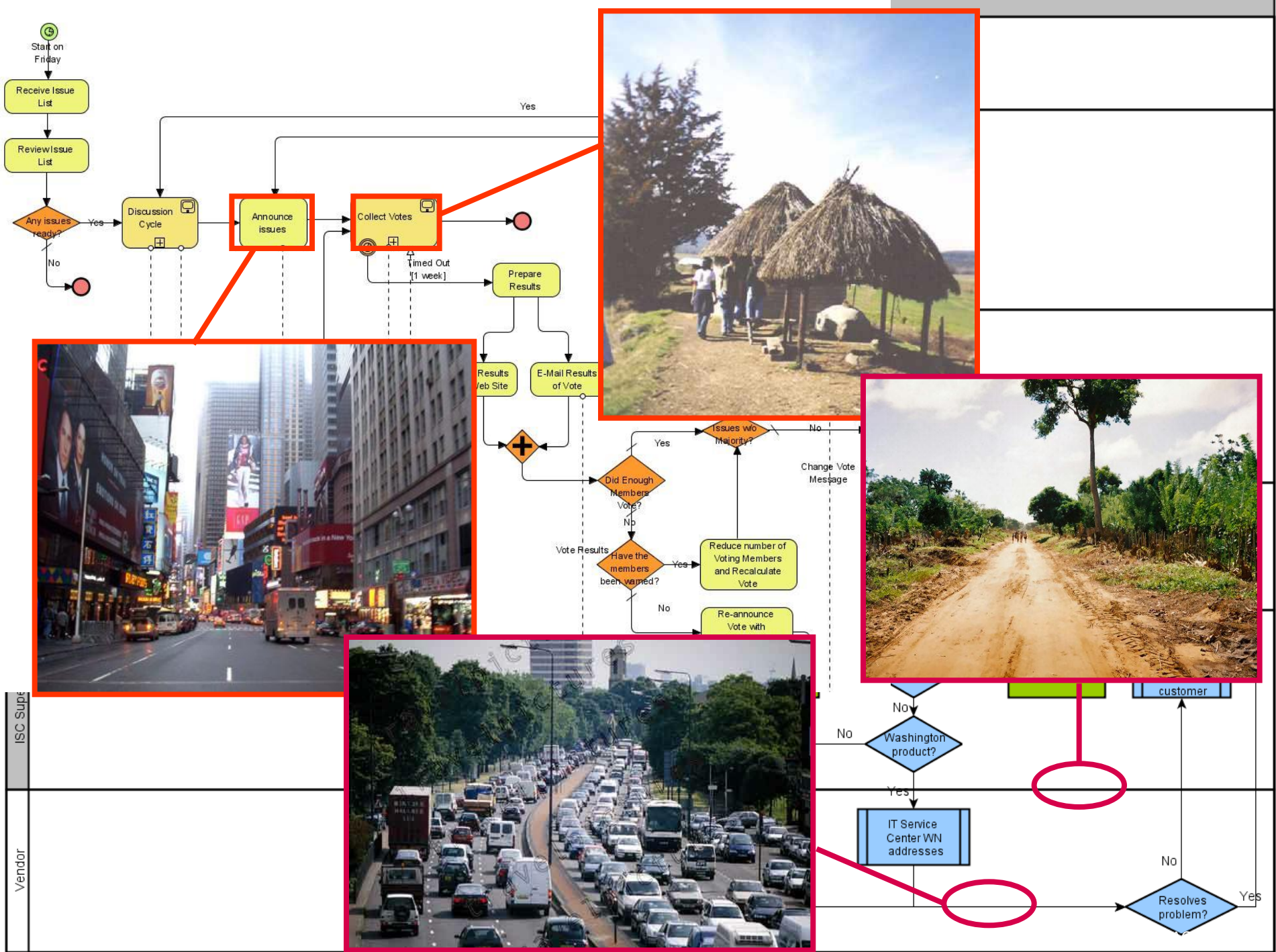
Make vagueness explicit!
(rather than having multiple notations for the same thing)





problem #6

Abstracting from the things that matter



What is process mining?



Why is process discovery difficult?



How about precision and recall?



What are the main research challenges?



How to measure the quality of a process model?



The future is bright, but how to get started?



What are the main pitfalls of process modeling?



Positioning Process Mining

process model analysis

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



performance-oriented
questions,
problems and
solutions



compliance-oriented
questions,
problems and
solutions



data-oriented analysis

(data mining, machine learning, business intelligence)

0100110011010101010

010011010101010



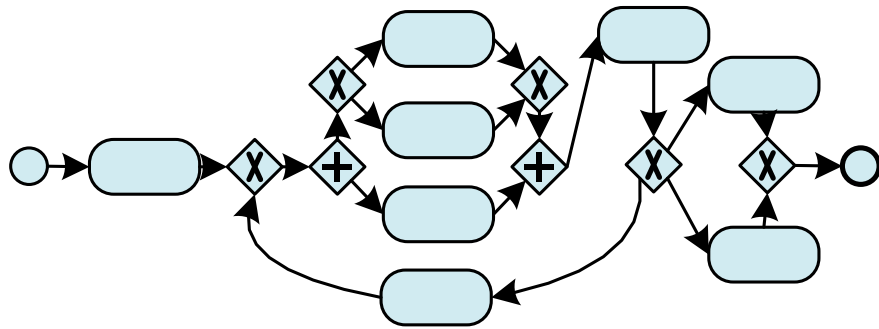
007001101010101010



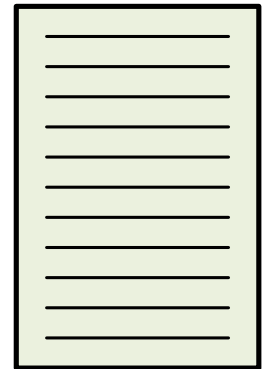


Let us take a step back and see how models and behavior relate: Let's play!

Play-Out

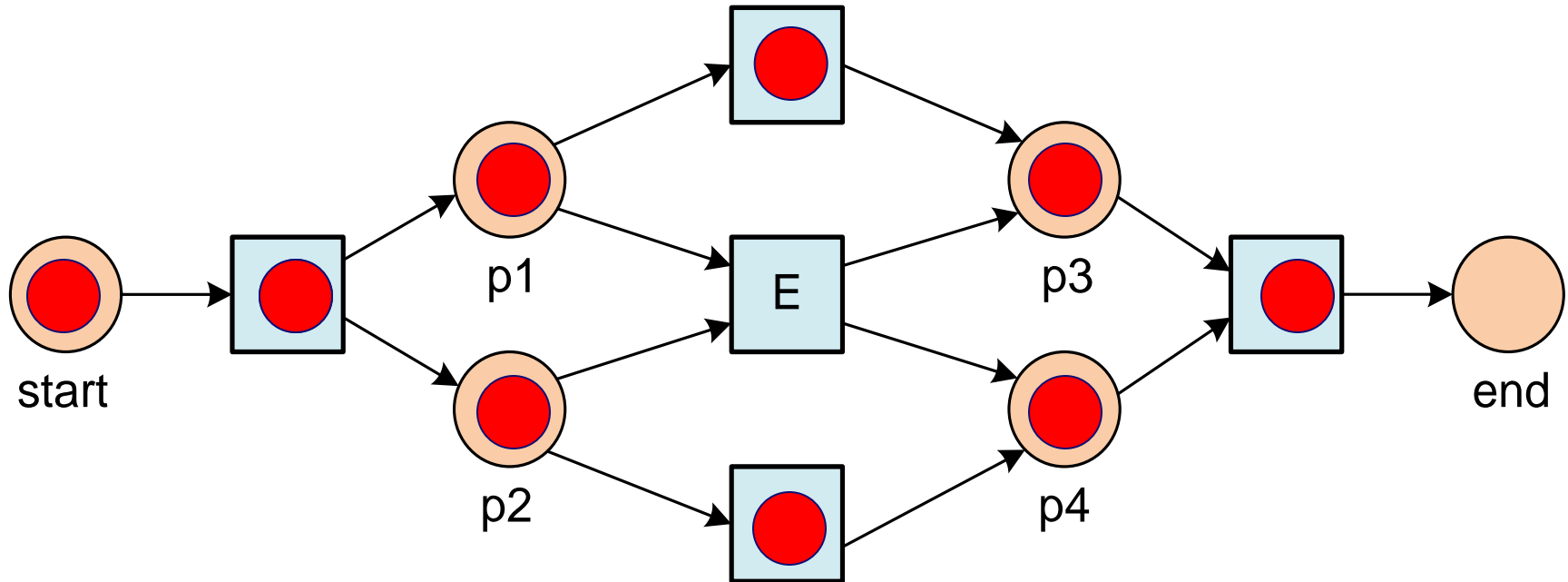


process model



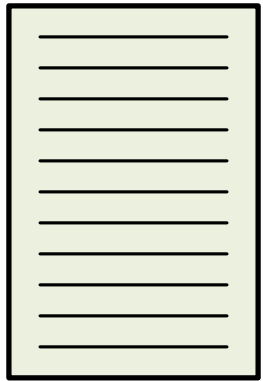
event log

Play-Out (Classical use of models)

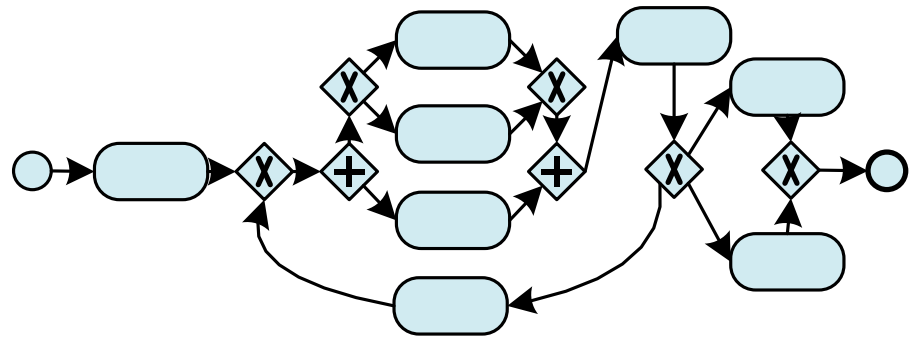
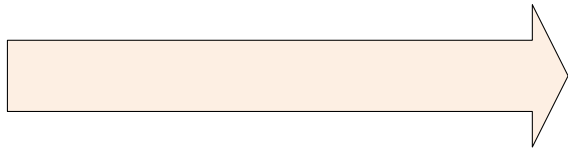


A B C D **A E D** **A E D**
A C B D **A B C D** **A C B D**
A C B D **A E D** **A C B D**

Play-In



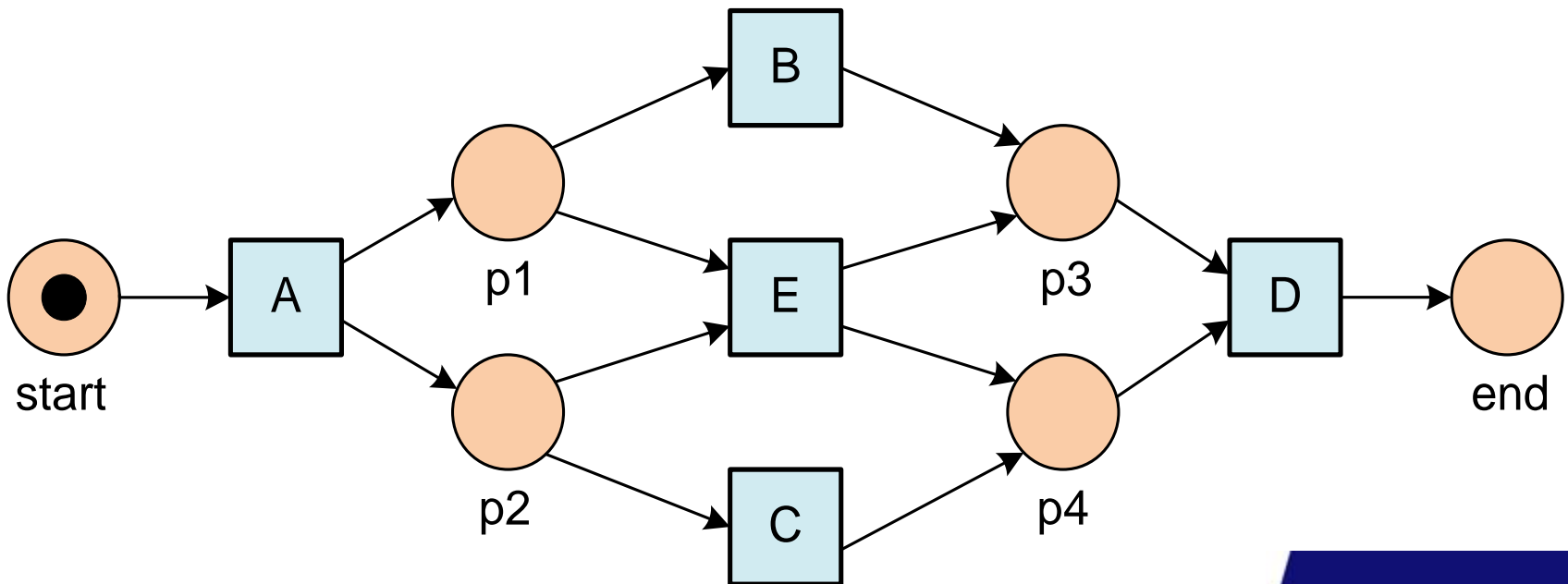
event log



process model

Play-In

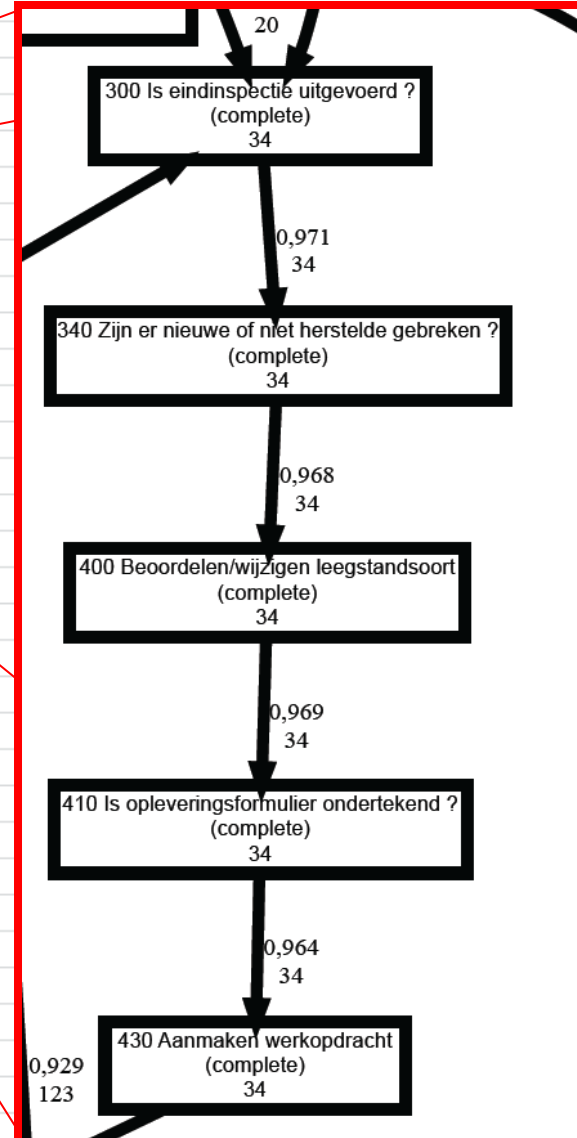
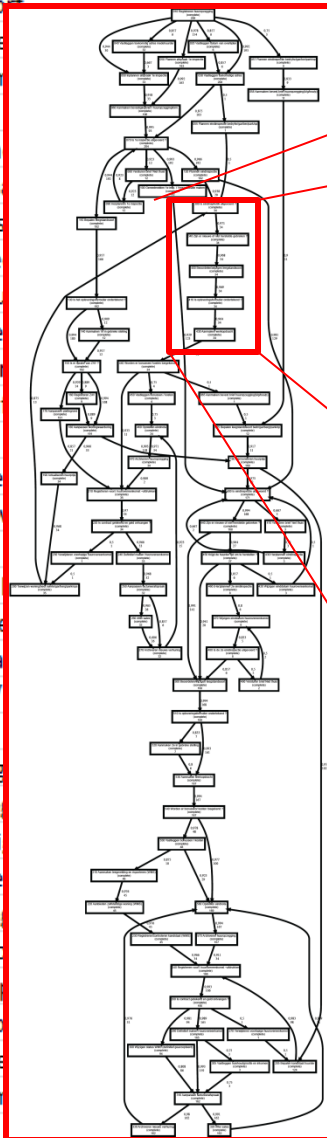
A B C D A E D A E D
A C B D A B C D A C B D
A C B D A E D A C B D



Example Process Discovery

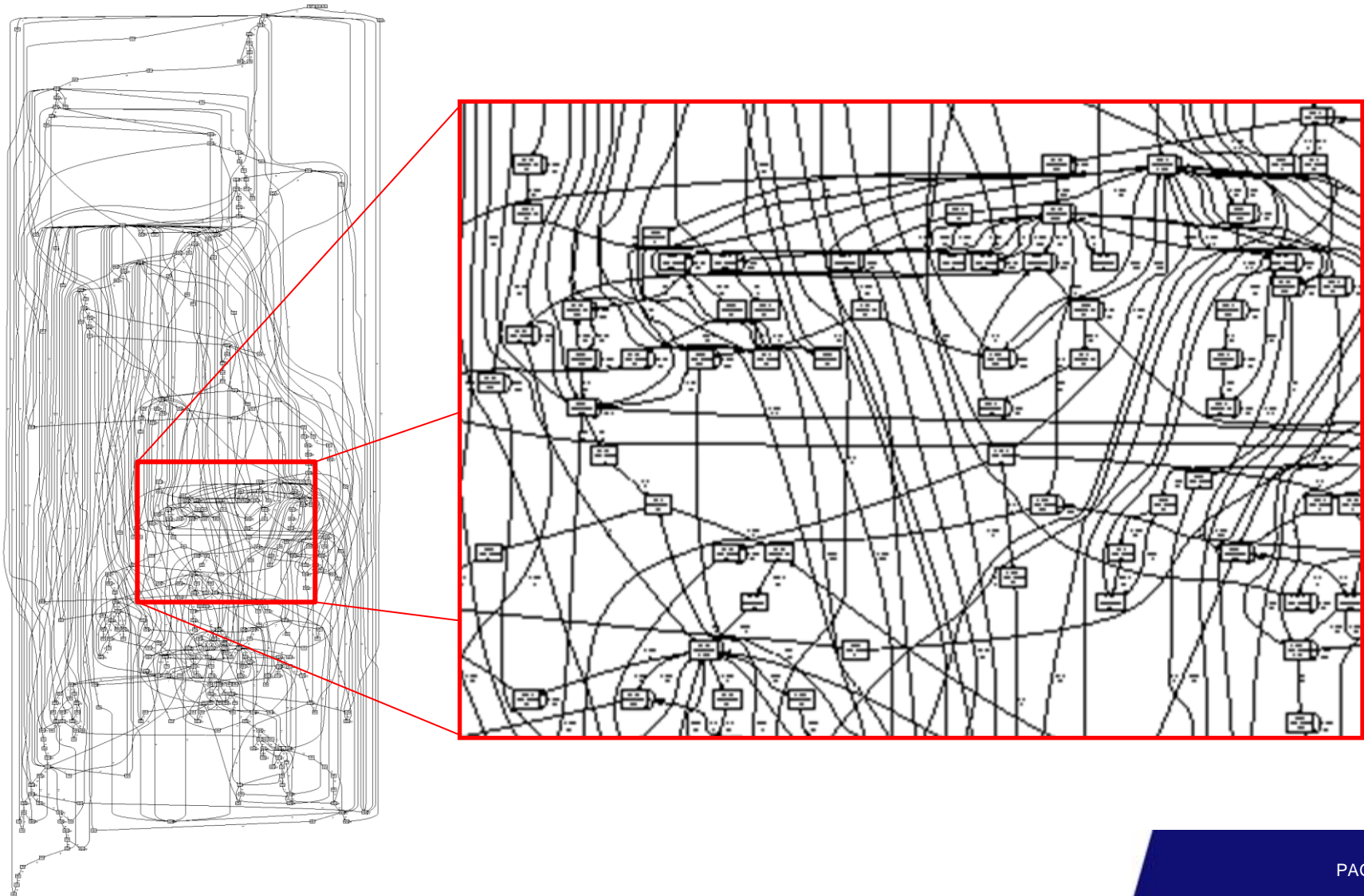
(Vestia, Dutch housing agency, 208 cases, 5987 events)

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 uitgev	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



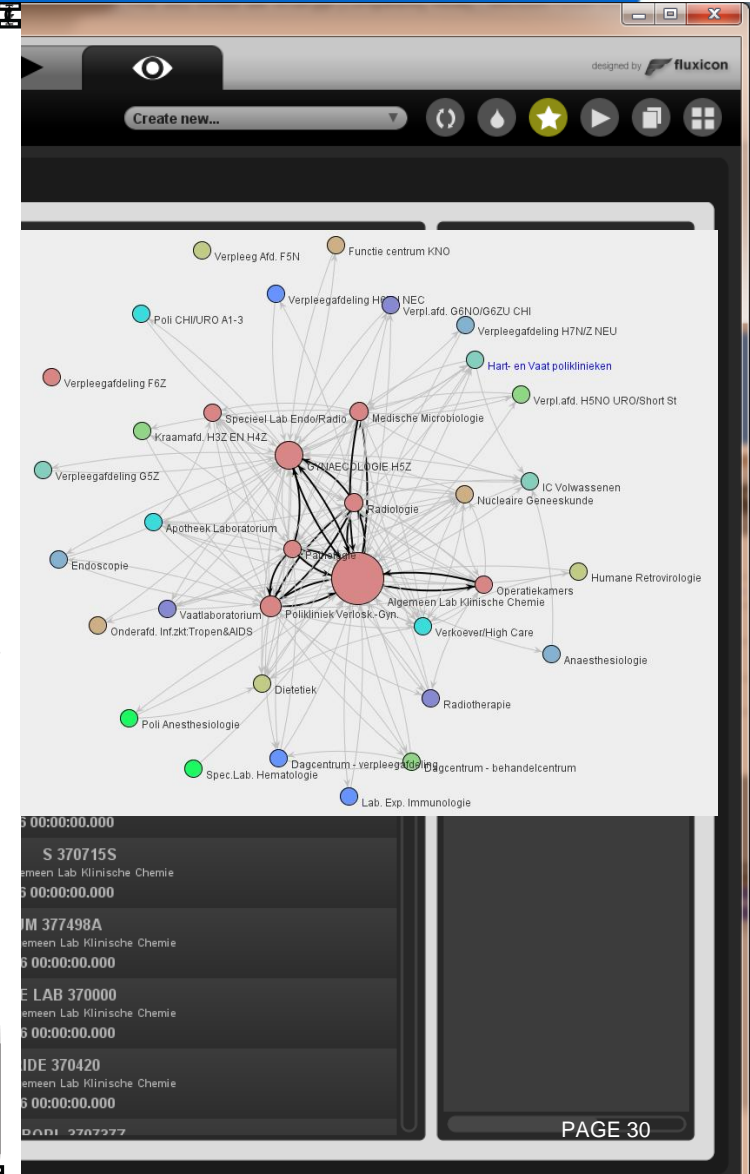
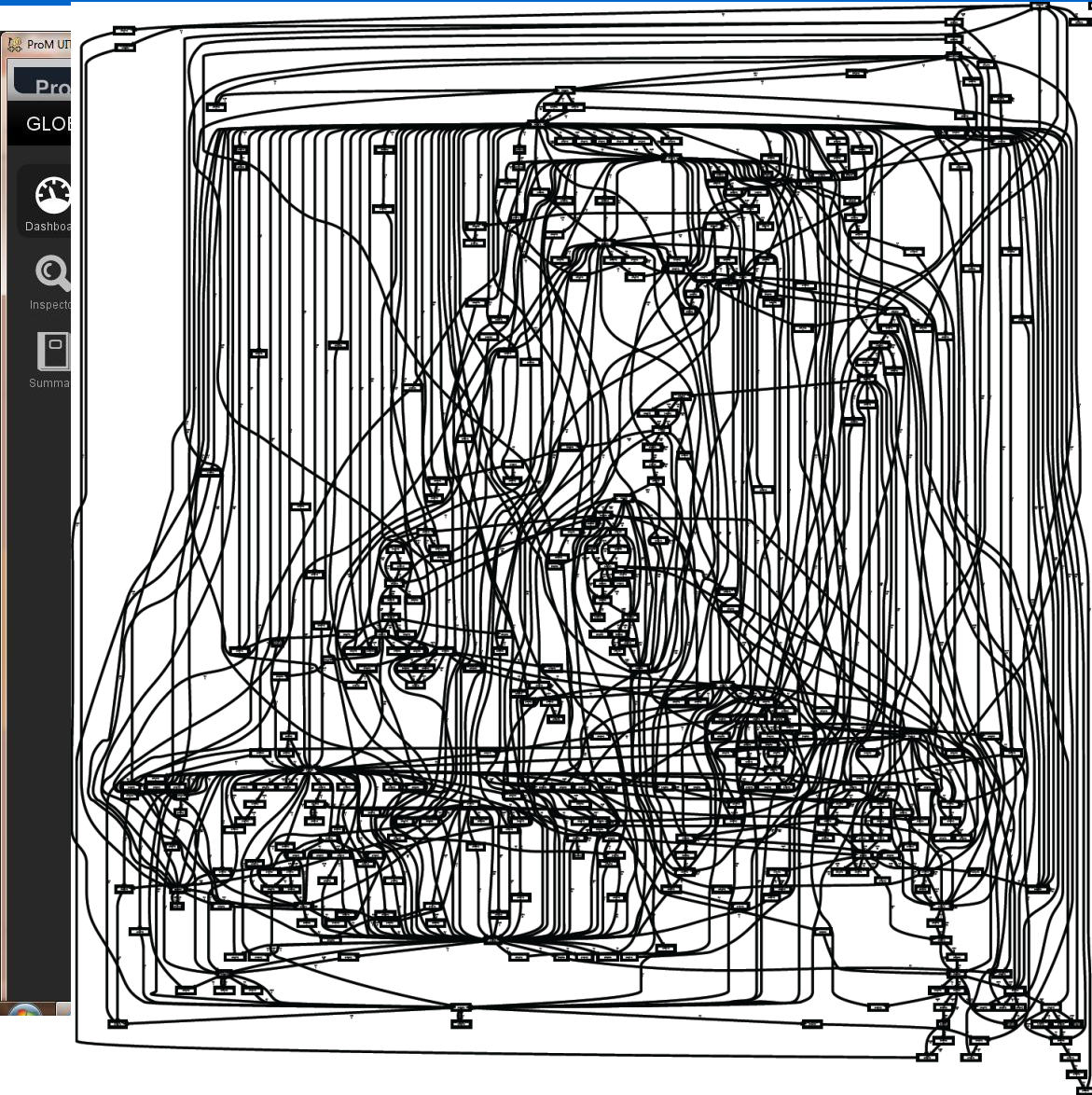
Example Process Discovery

(ASML, test process lithography systems, 154966 events)

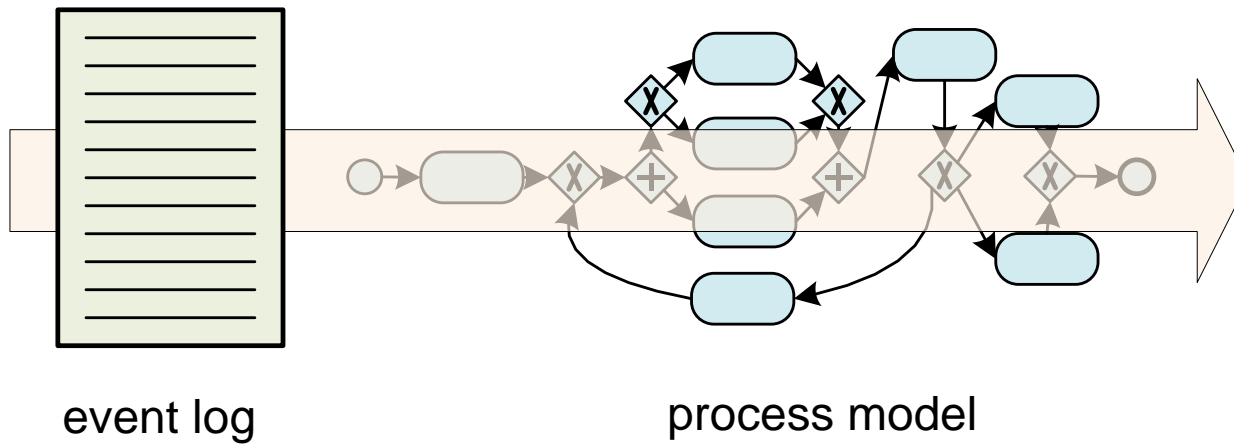


Example Process Discovery

(AMC, 627 gynecological oncology patients, 24331 events)



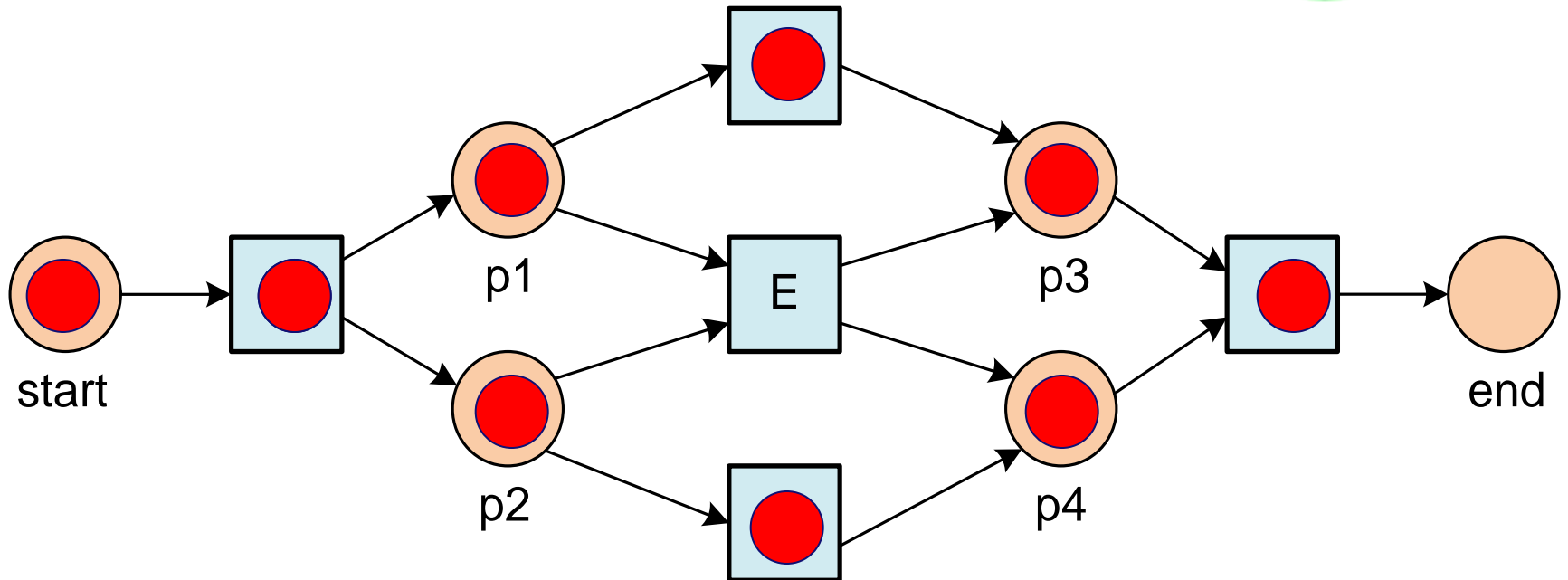
Replay



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

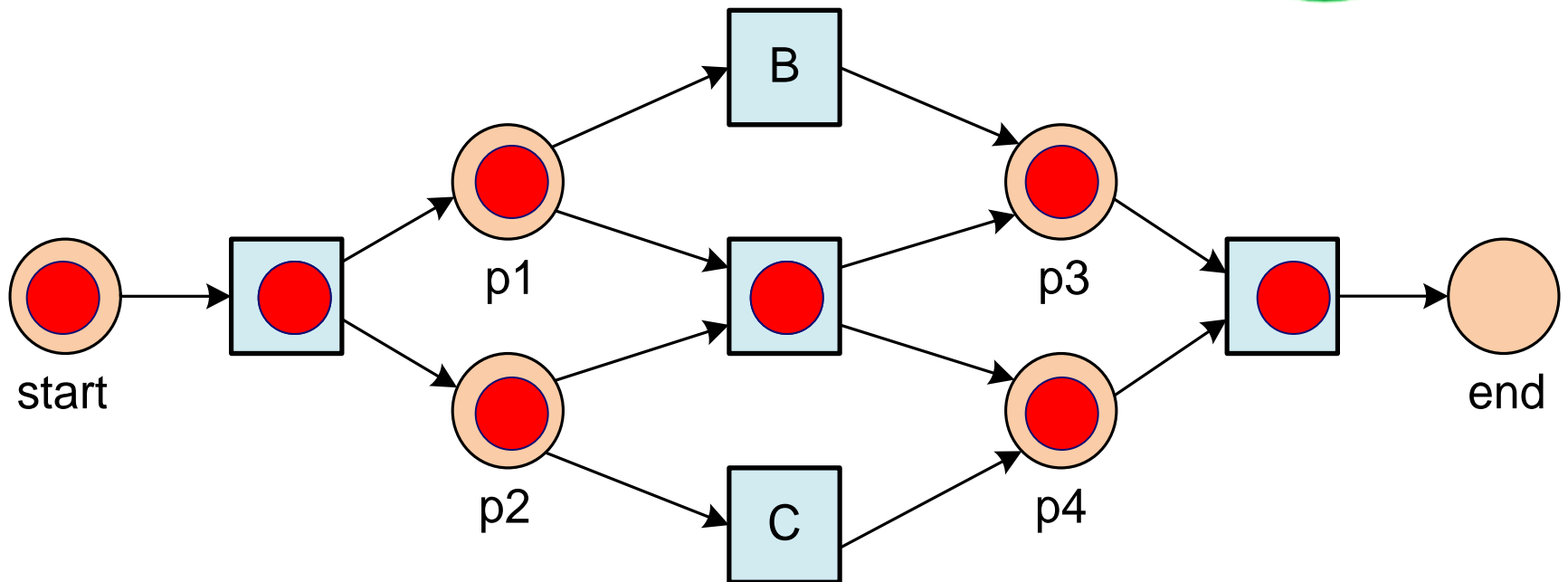
Replay

A B C D



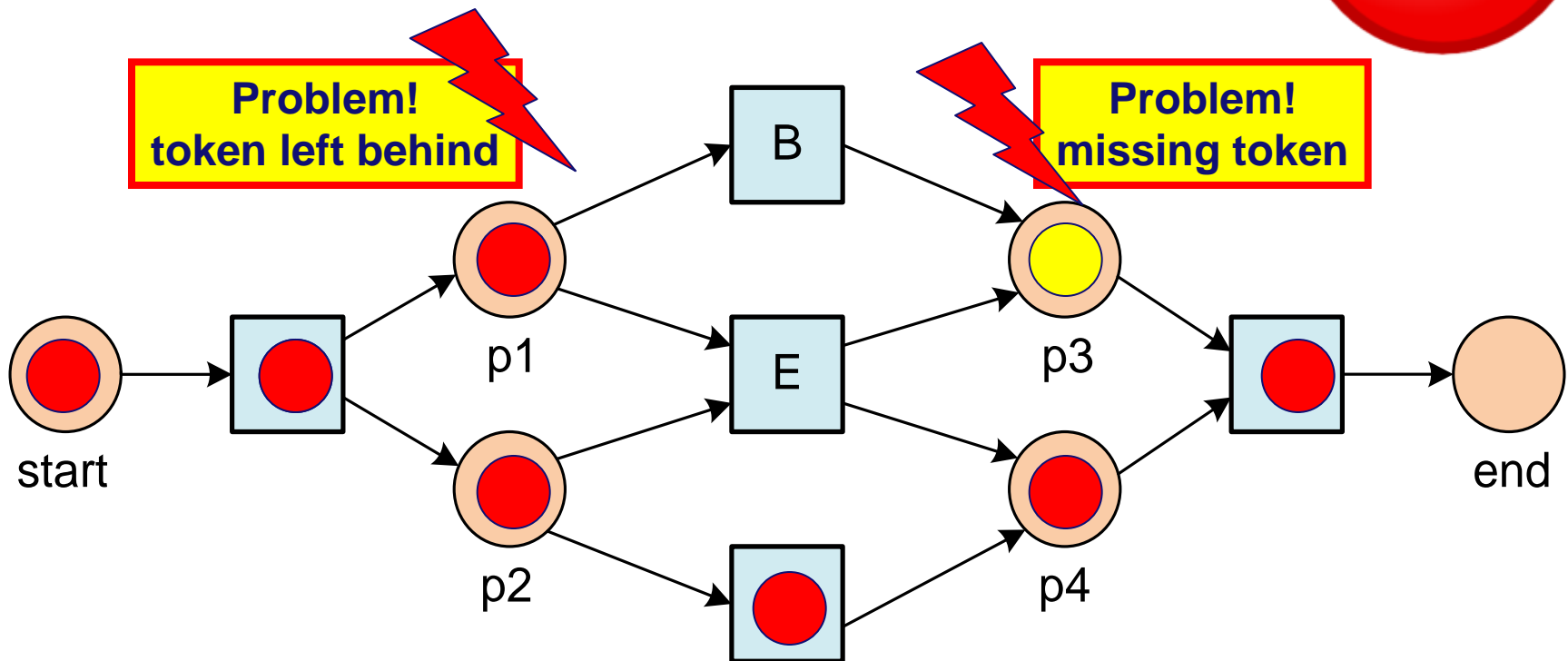
Replay

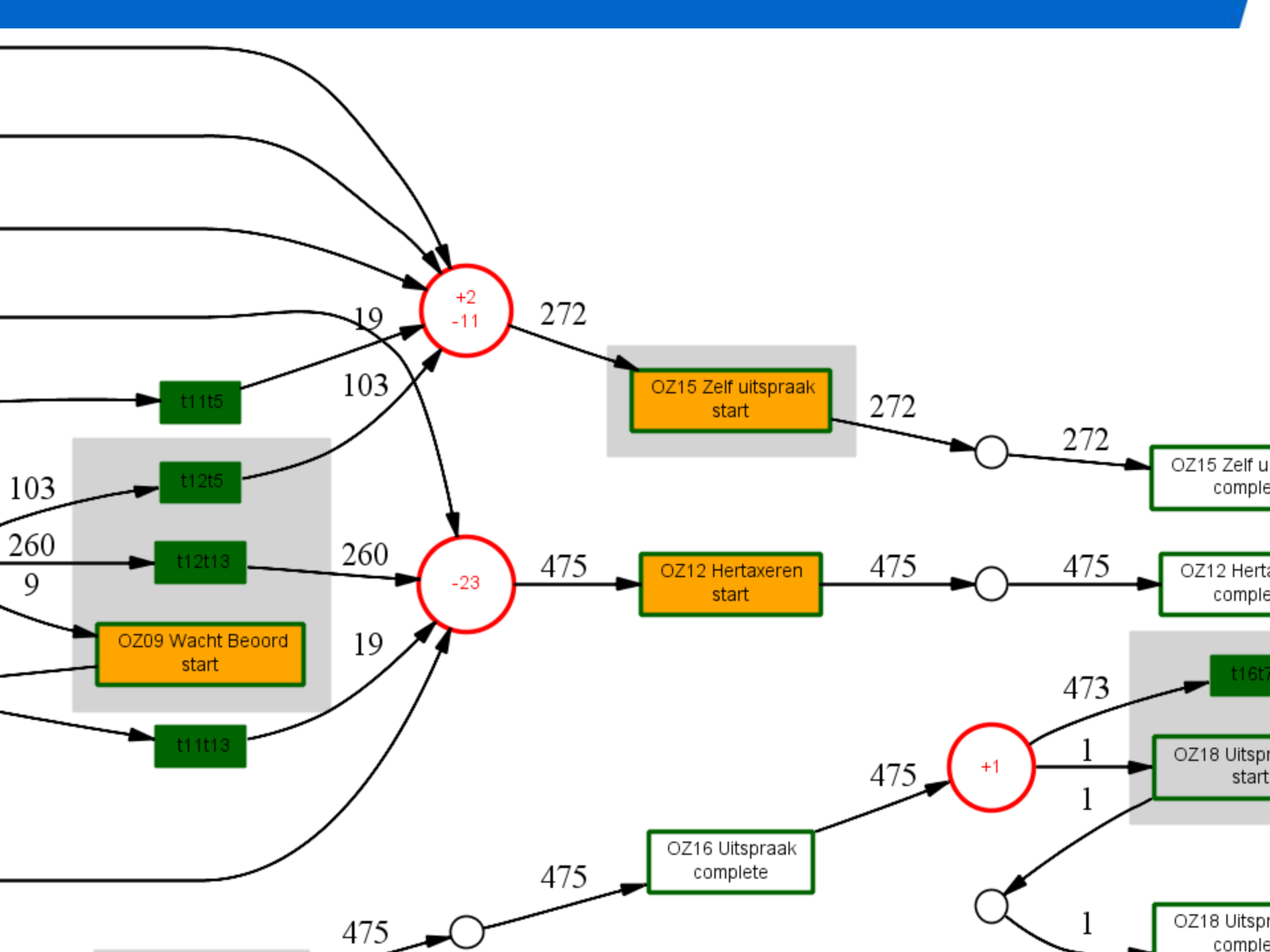
A E D



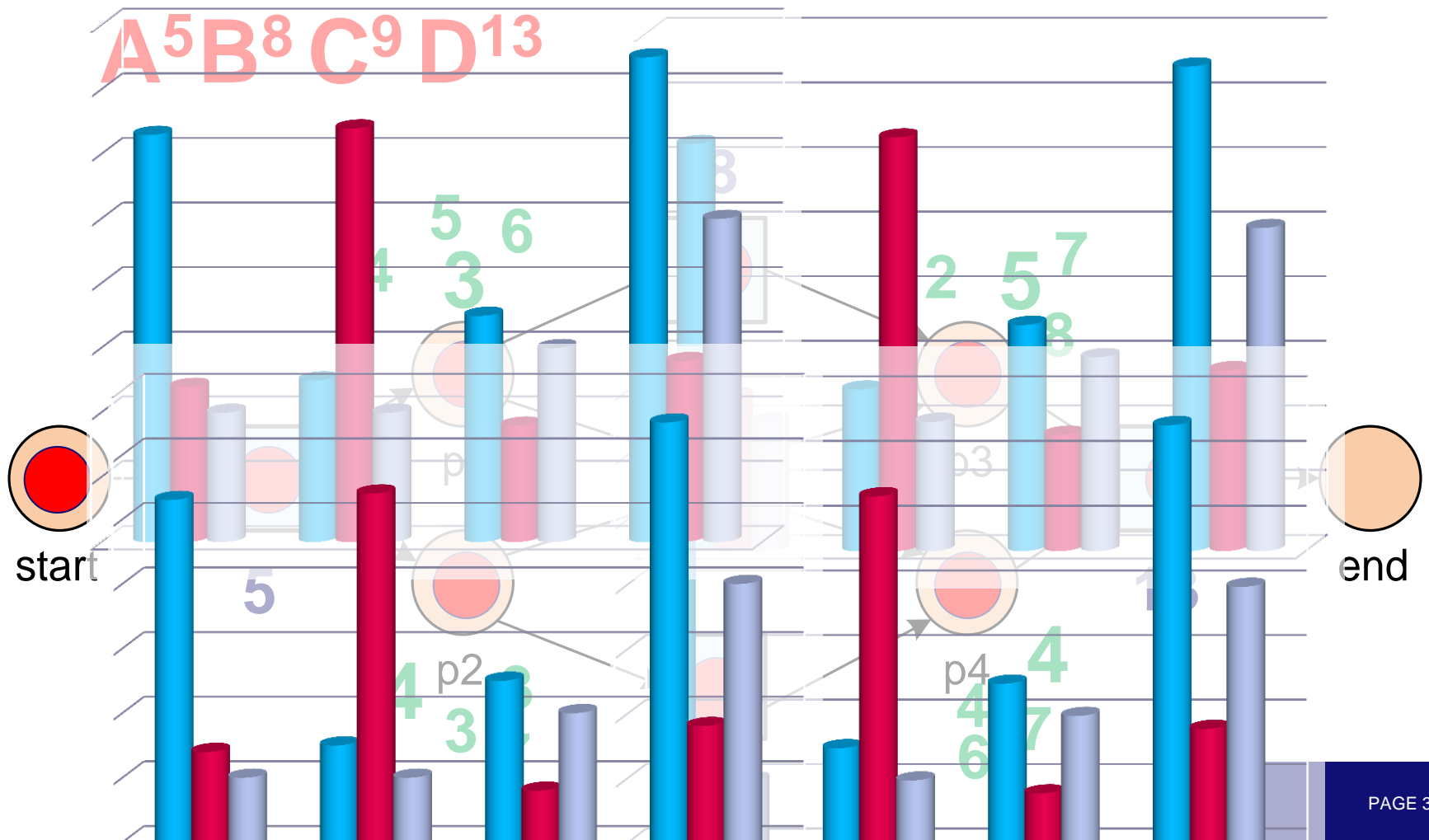
Replay can detect problems

ACD



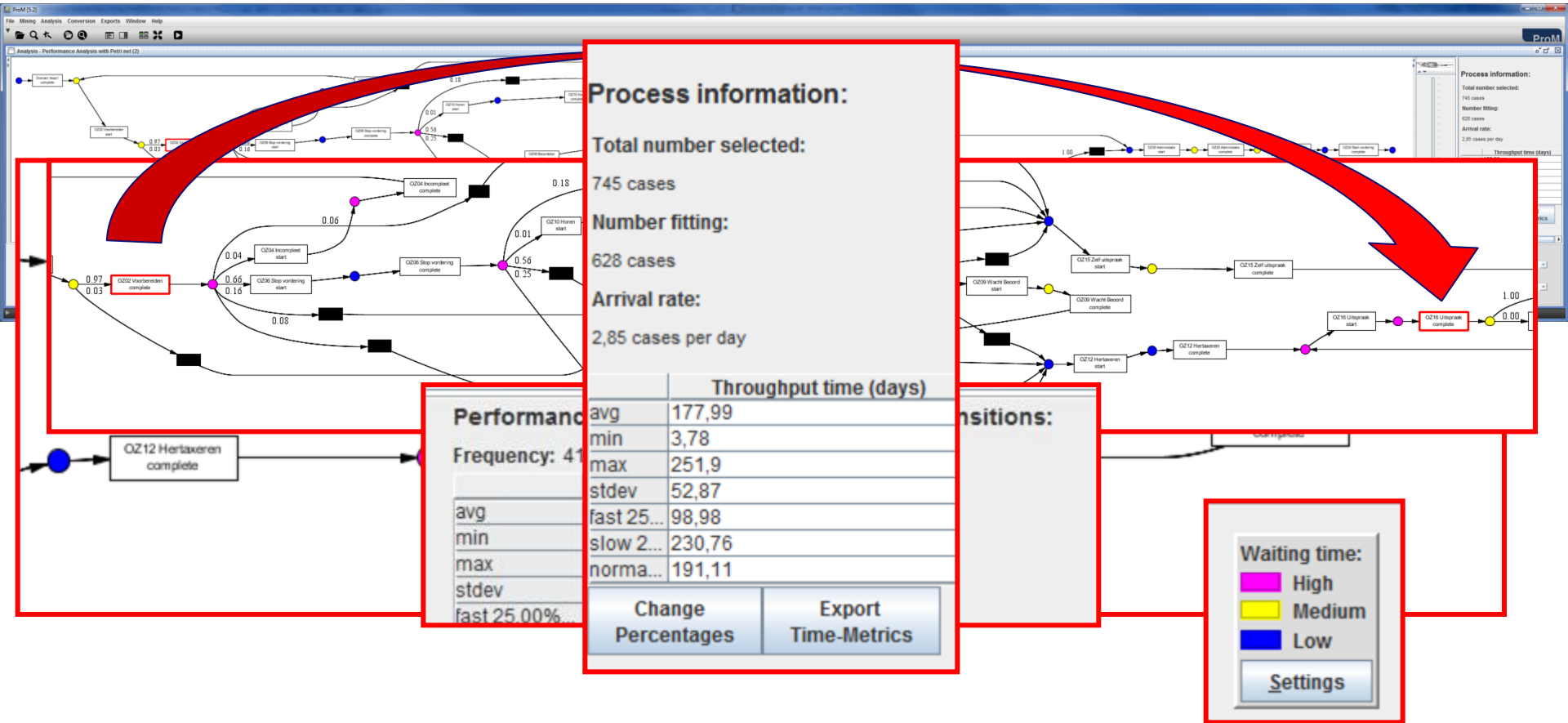


Replay can extract timing information



Performance Analysis Using Replay

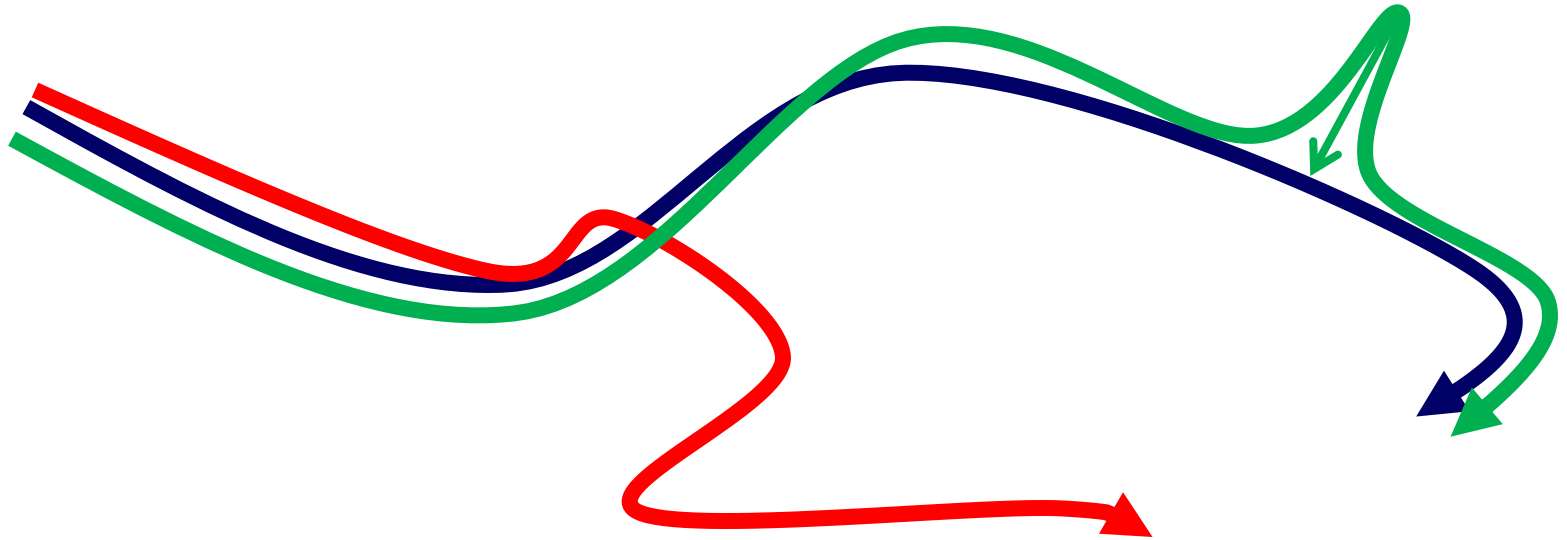
(WOZ objections Dutch municipality, 745 objections, 9583 event, f= 0.988)



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

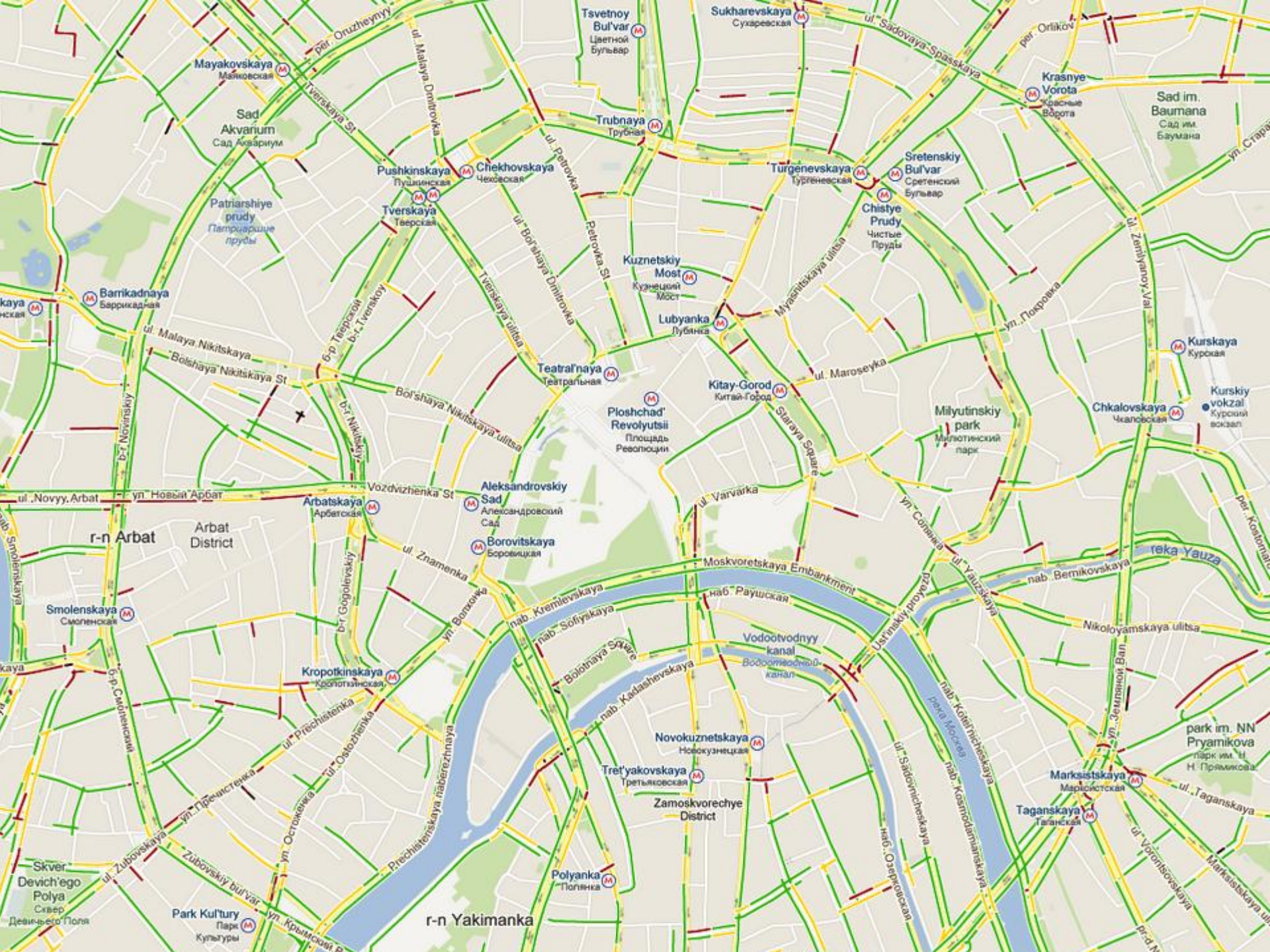


Alignments are essential!



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

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



Why is process discovery difficult?



How about precision and recall?



What are the main research challenges?



How to measure the quality of a process model?



The future is bright, but how to get started?

What is process mining?



What are the main pitfalls of process modeling?



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 (before 1998) did not consider concurrency and definitely not end-to-end business process models.
- ... with or ... examples, ... etc.

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



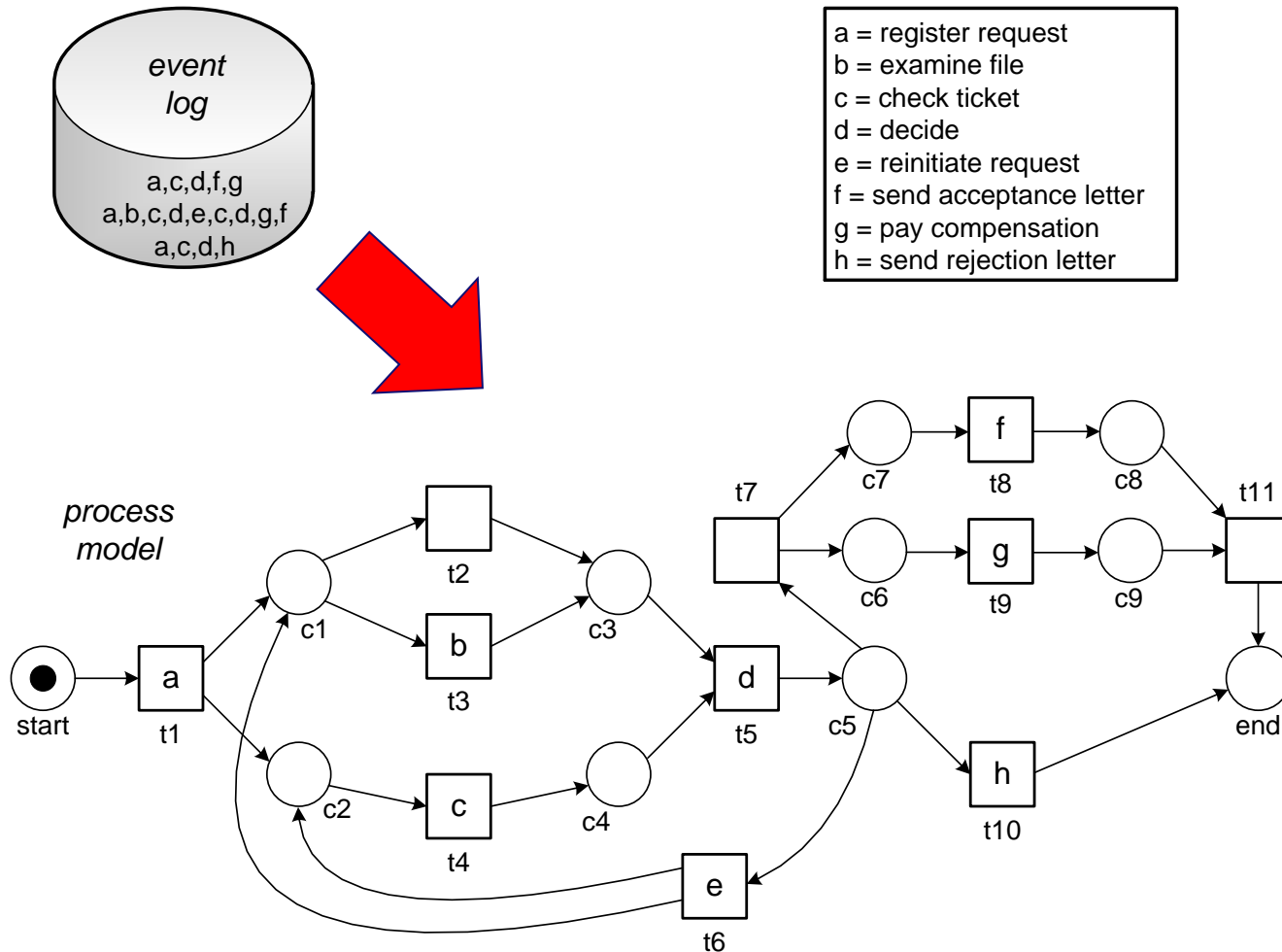
**at the start of the century,
process mining emerged as a
new research topic**

**remarkable progress over a
relatively short period**

See my keynote at <http://fluxicon.com/camp/2013/> earlier this week.

Process discovery challenge

(oversimplified no resources, data, etc.)



Process discovery algorithms (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

conformal process graph

multi-phase mining

partial-order based mining

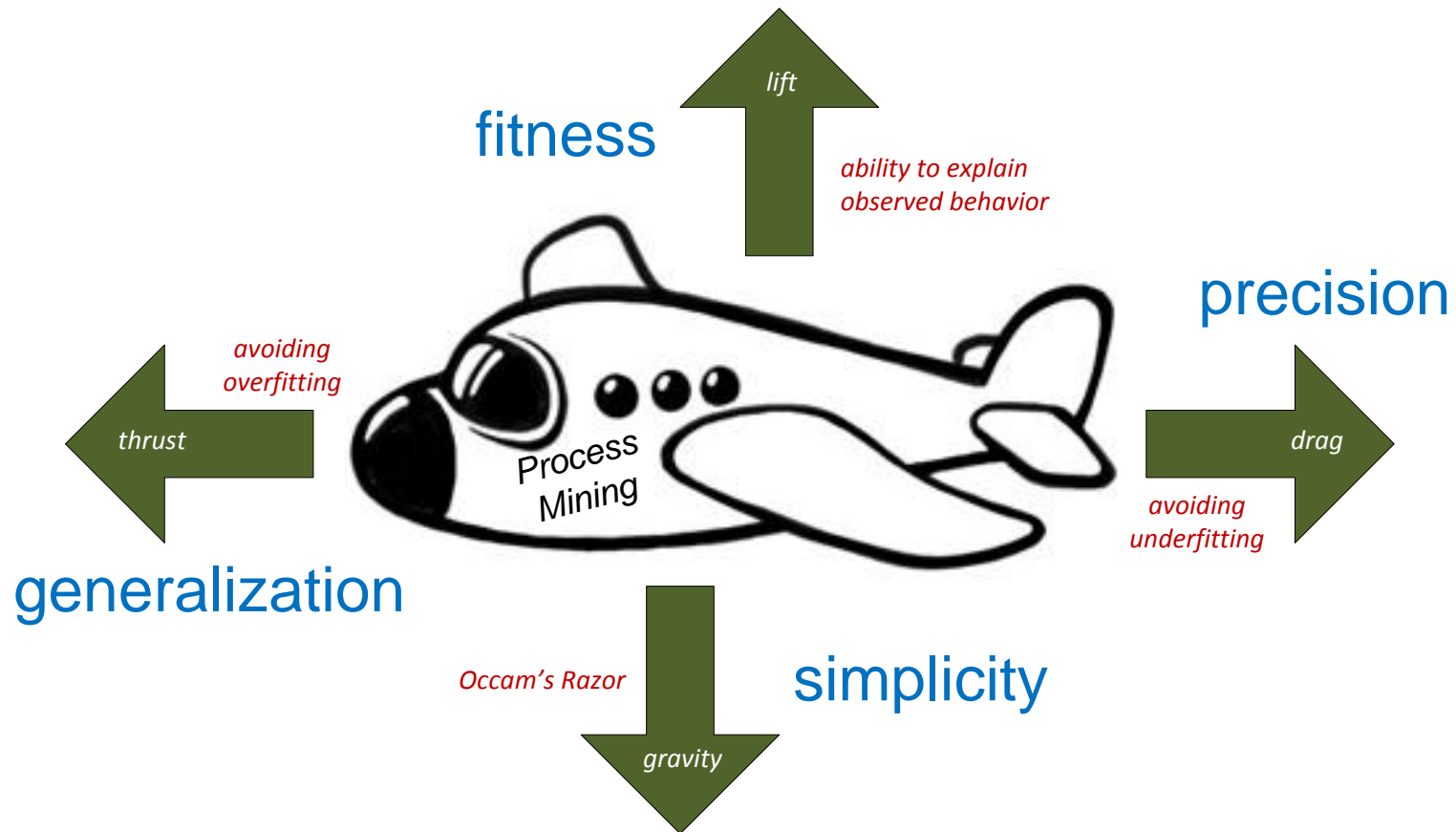
$\alpha\#$ algorithm

ILP mining

$\alpha++$ algorithm

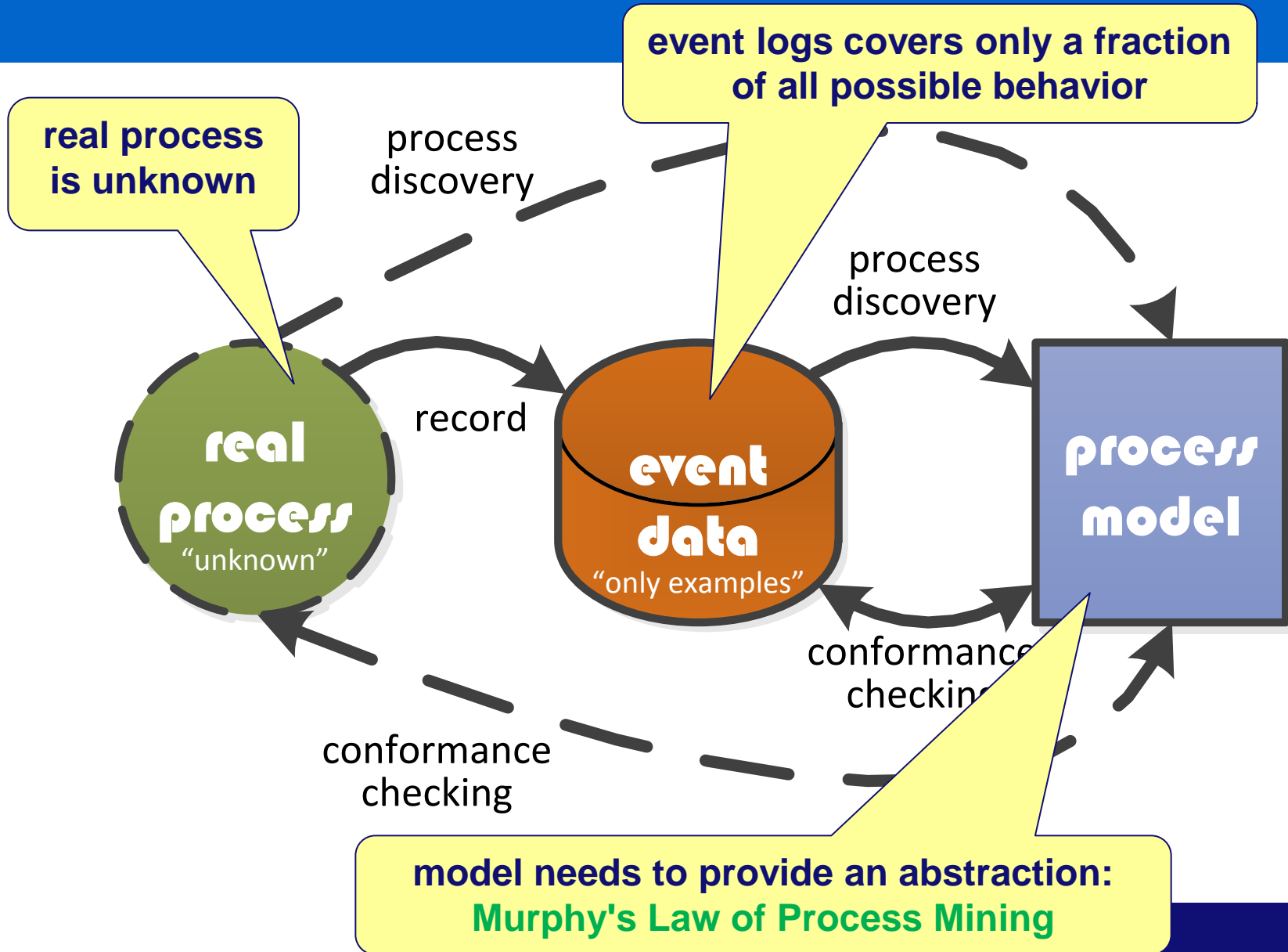


How good is my model: Four forces



Leaving out one of these dimensions during discovery will lead to degenerate cases!

Problem





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)**

Why is process discovery difficult?



How about precision and recall?



What are the main research challenges?



How to measure the quality of a process model?



The future is bright, but how to get started?



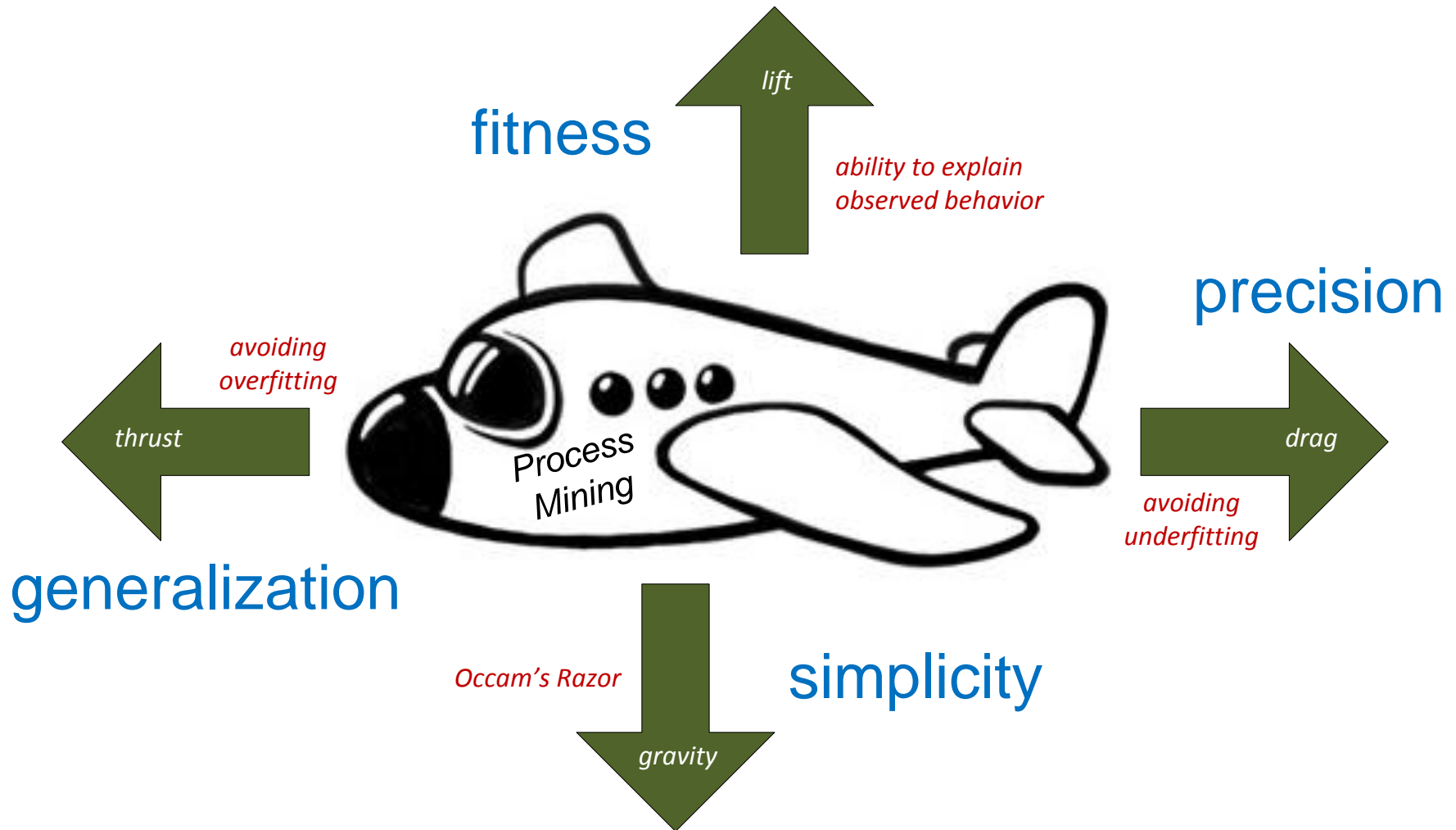
What is process mining?



What are the main pitfalls of process modeling?

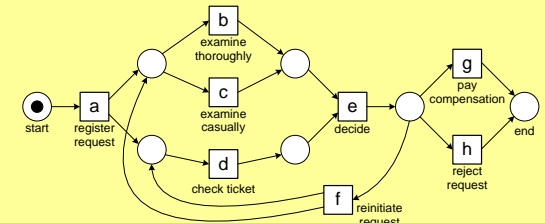
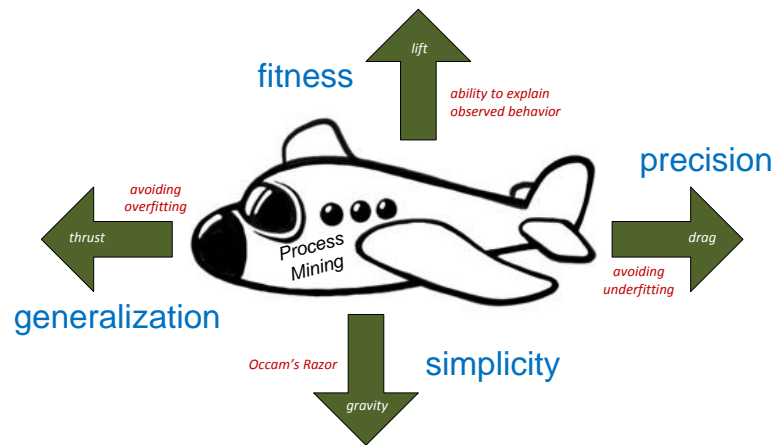


Remember the four forces

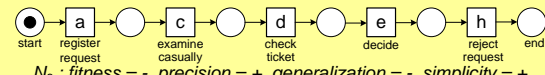


Example: one log four models

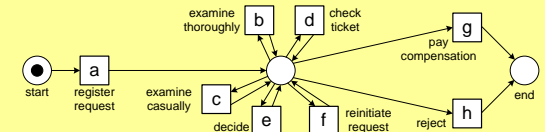
#	trace
455	acdeh
191	abdeg
177	adceh
144	abdeh
111	acdeg
82	adceg
56	adbeh
47	acdefdbeh
38	adbeg
33	acdefbdeh
14	acdefbdeg
11	acdefdbeg
9	adcefcdeh
8	adcefdbeh
5	adcefbdeg
3	acdefbdefdbeg
2	adcefbdeg
2	adcefbdefdbeg
1	adcefdbefbdeh
1	adbefbdefdbeg
1	adcefdbefcdefdbeg
1391	



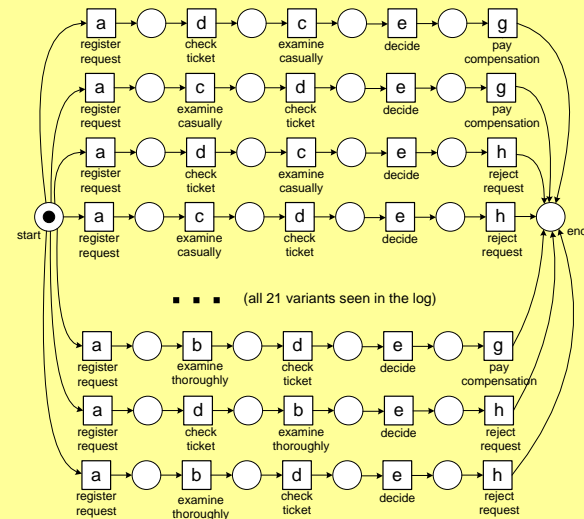
N_1 : fitness = +, precision = +, generalization = +, simplicity = +



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

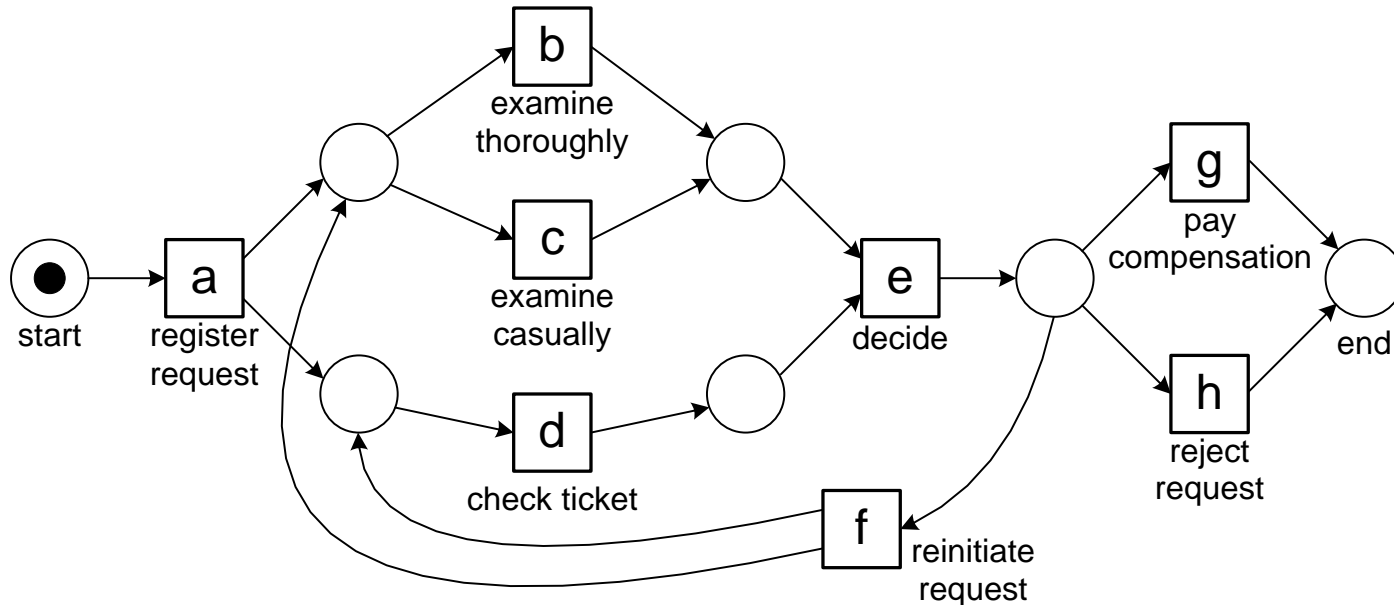


N_3 : fitness = +, precision = -, generalization = +, simplicity = +



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

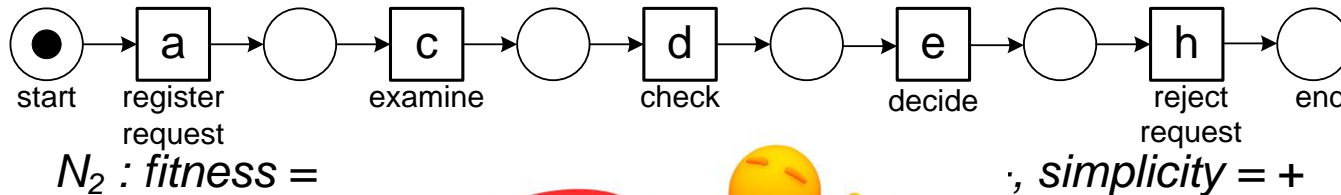
Model N_1



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	acdefbddeg
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	acdefbdeg
11	acdefdbeg
9	adcefcdeh
8	adcefdbeh
5	adcefbdeg
3	acdefbdefdbeg
2	adcefdbeg
2	adcefbdefbdeg
1	adcefdbefbdeh
1	adbefbdefdbeg
1	adcefdbefcdefdbeg
1391	

Model N₃

examine
thoroughly

b

d

check
ticket

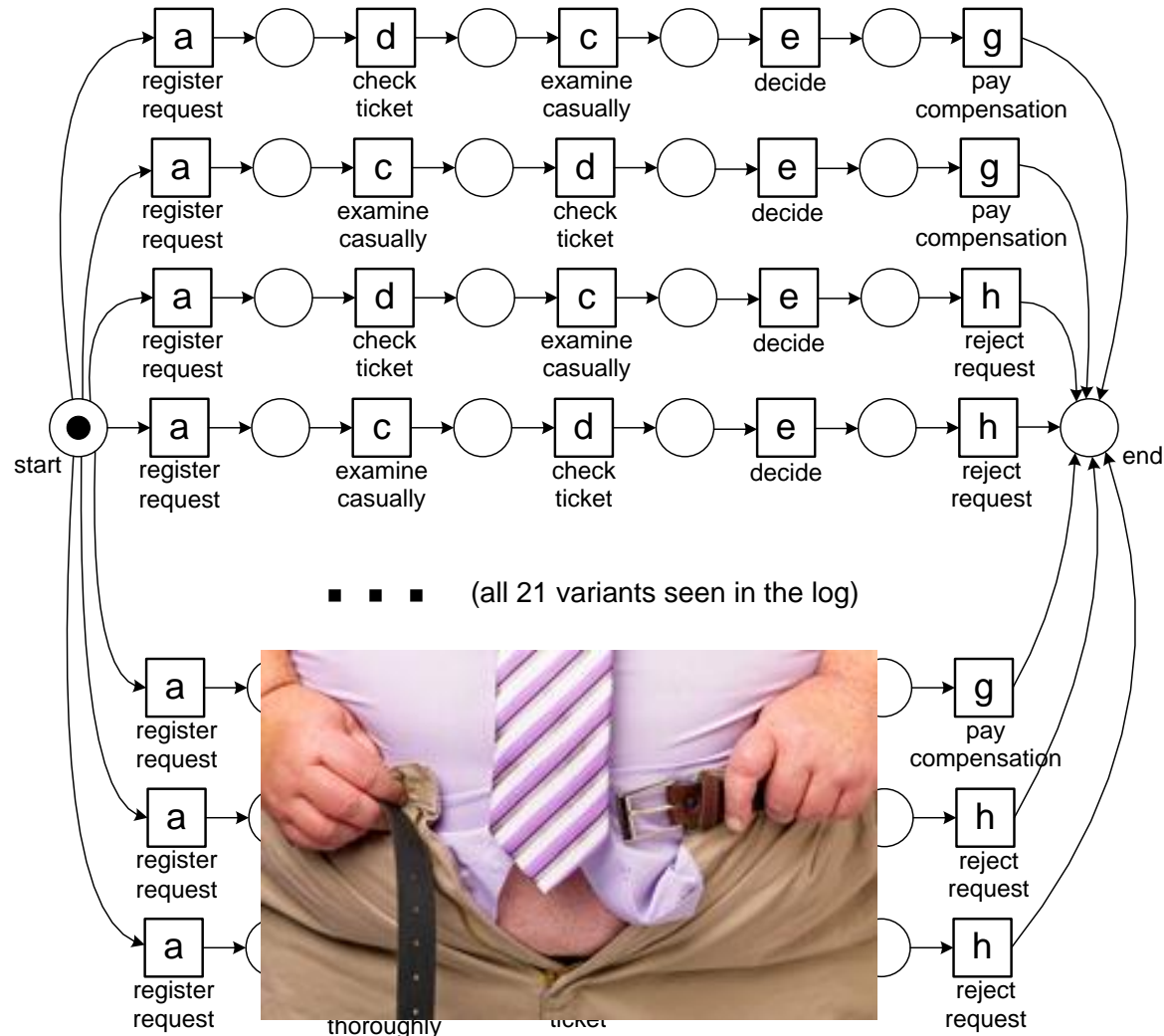
pay

g



#	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	

Why is process discovery difficult?



How about precision and recall?



What are the main research challenges?



How to measure the quality of a process model?



The future is bright, but how to get started?



What is process mining?



What are the main pitfalls of process modeling?

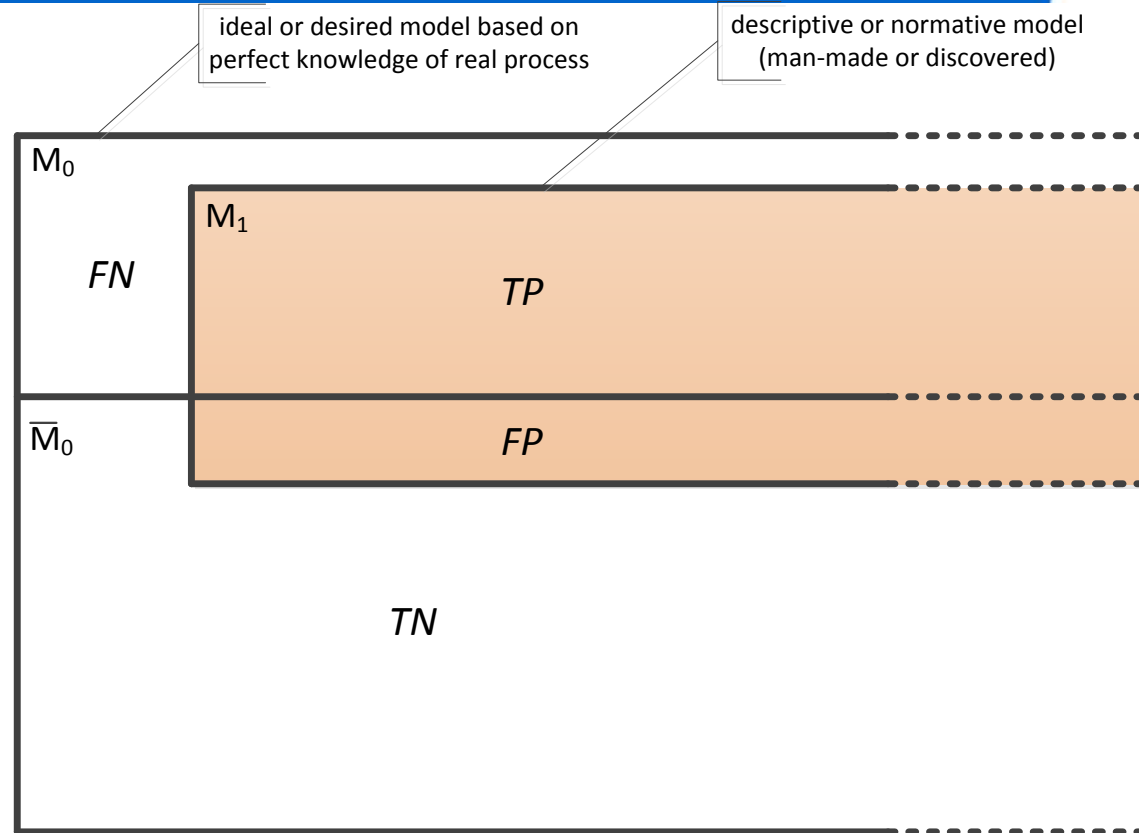
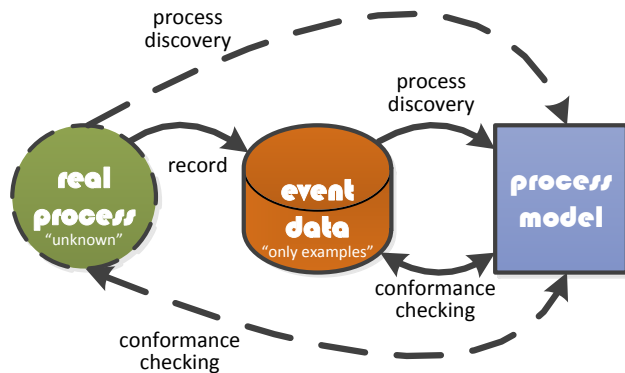


Several well-defined metrics for fitness, precision, generalization and simplicity exist ...



Why not use precision and recall as "objective" measures?

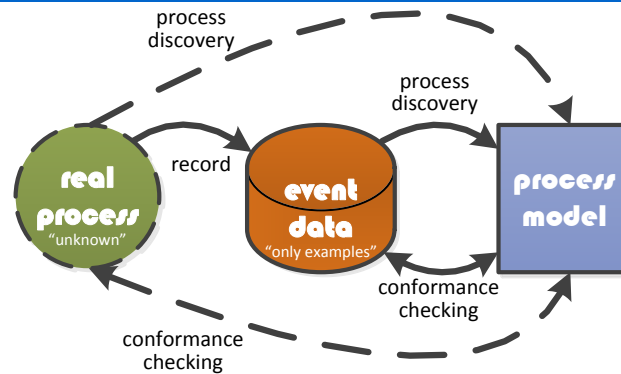
Suppose we know the "real model"



$$precision = \frac{TP}{M_1} = \frac{TP}{TP + FP}$$

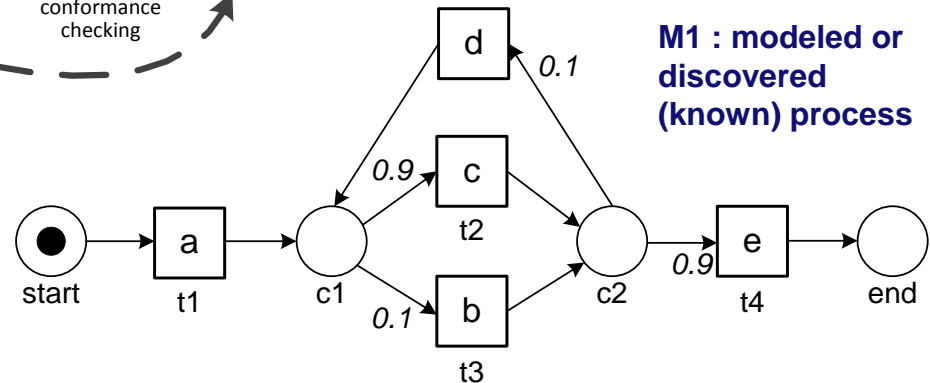
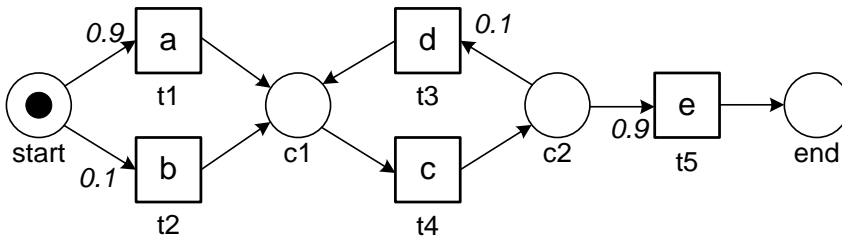
$$recall = \frac{TP}{M_0} = \frac{TP}{TP + FN}$$

Let us try ...



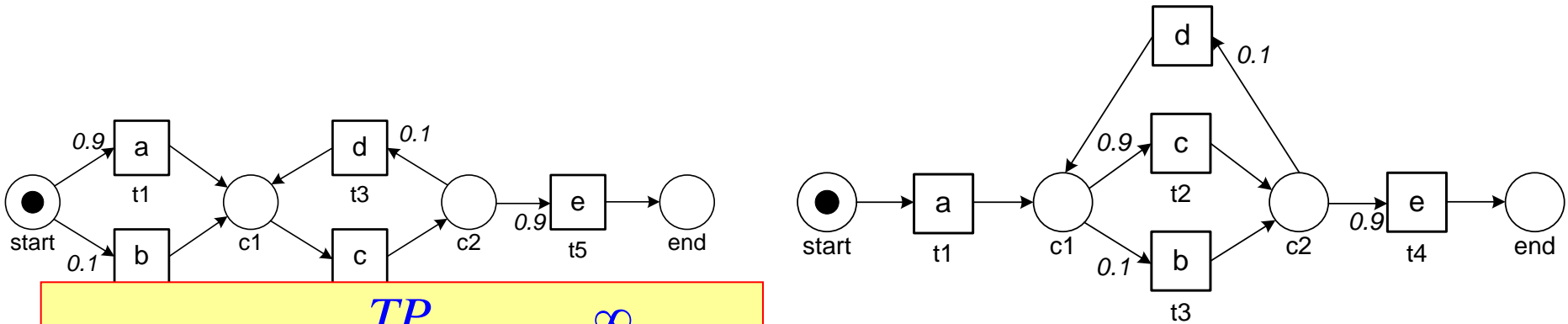
M₀ : unknown real process

M₁ : modeled or discovered (known) process



M_0		M_1	
FN	bce bcdce bcdcdce ...	TP	ace acdce acdcdce ...
\bar{M}_0		FP	
TN	aaaa abababab ccdddd ...	FP	abe acdbe abdbe ...

Oops ...



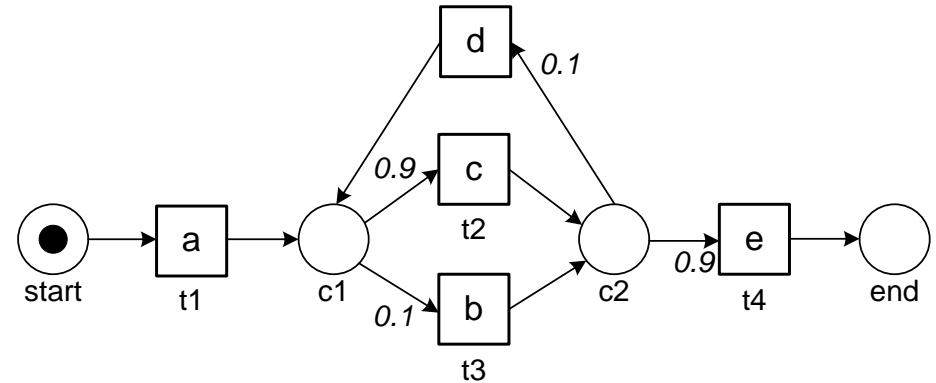
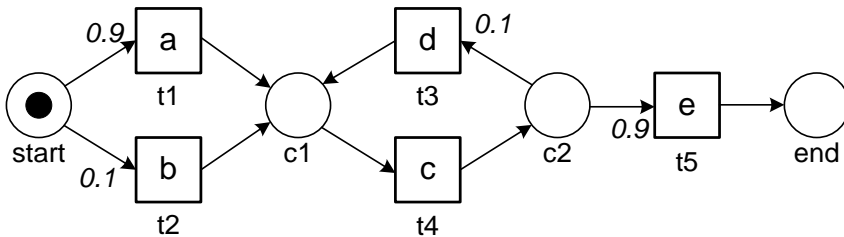
$$precision = \frac{TP}{TP + FP} = \frac{\infty}{\infty + \infty} = ?$$

Need to add probabilities / frequencies !

FN	bce bcdce bcdcdce ...	TP	ace acdce acd...
FN	aaaa abababab cccddd ...	TP	abe acdbe ...

$$recall = \frac{TP}{TP + FN} = \frac{\infty}{\infty + \infty} = ?$$

Let's compute precision and recall taking a specific viewpoint (M_0 or M_1)



~~$$precision_0 = \frac{TP}{TP + FP} = \frac{0.9}{0.9 + 0.0} = 1$$~~

ace	0.81
bce	0.09
acdce	0.081

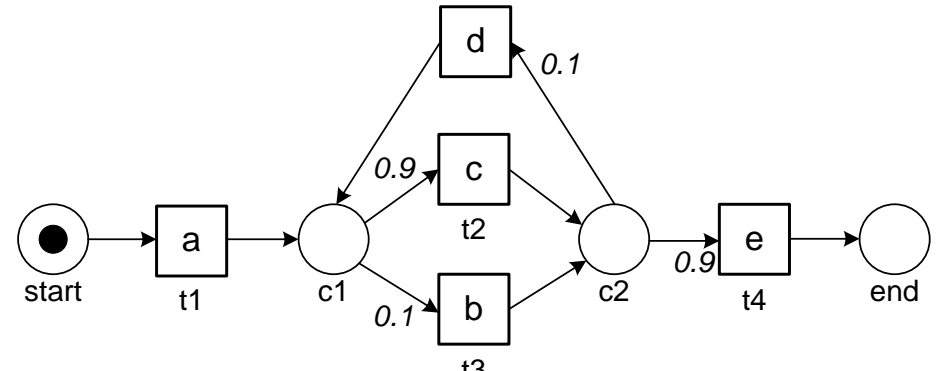
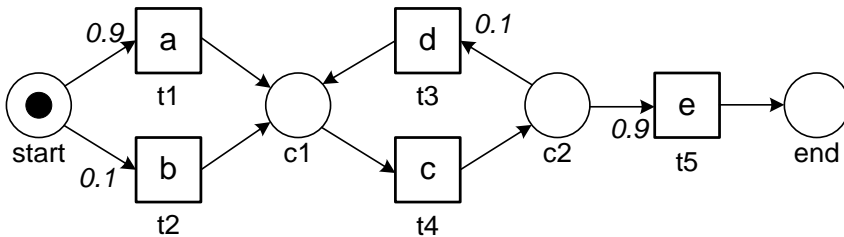
$$precision_1 = \frac{TP}{TP + FP} = \frac{0.888}{0.888 + 0.111} = 0.888$$

trace	probability in M_0
ace	0.81
bce	0.09
acdce	0.081
...	...
bce	0.0

$$recall_0 = \frac{TP}{TP + FN} = \frac{0.9}{0.9 + 0.1} = 0.9$$

~~$$recall_1 = \frac{TP}{TP + FN} = \frac{0.888}{0.888 + 0.0} = 1$$~~

Only two dimensions rather than four?



small if there are many likely traces in the real model that do not fit the discovered model

small if there are many likely traces in the discovered model that are impossible in the real model

$$recall_0 = \frac{TP}{TP + FN} = \frac{0.9}{0.9 + 0.1} = 0.9$$

$$precision_1 = \frac{TP}{TP + FP} = \frac{0.888}{0.888 + 0.111} = 0.888$$

fitness



ability to explain observed behavior

precision



avoiding underfitting

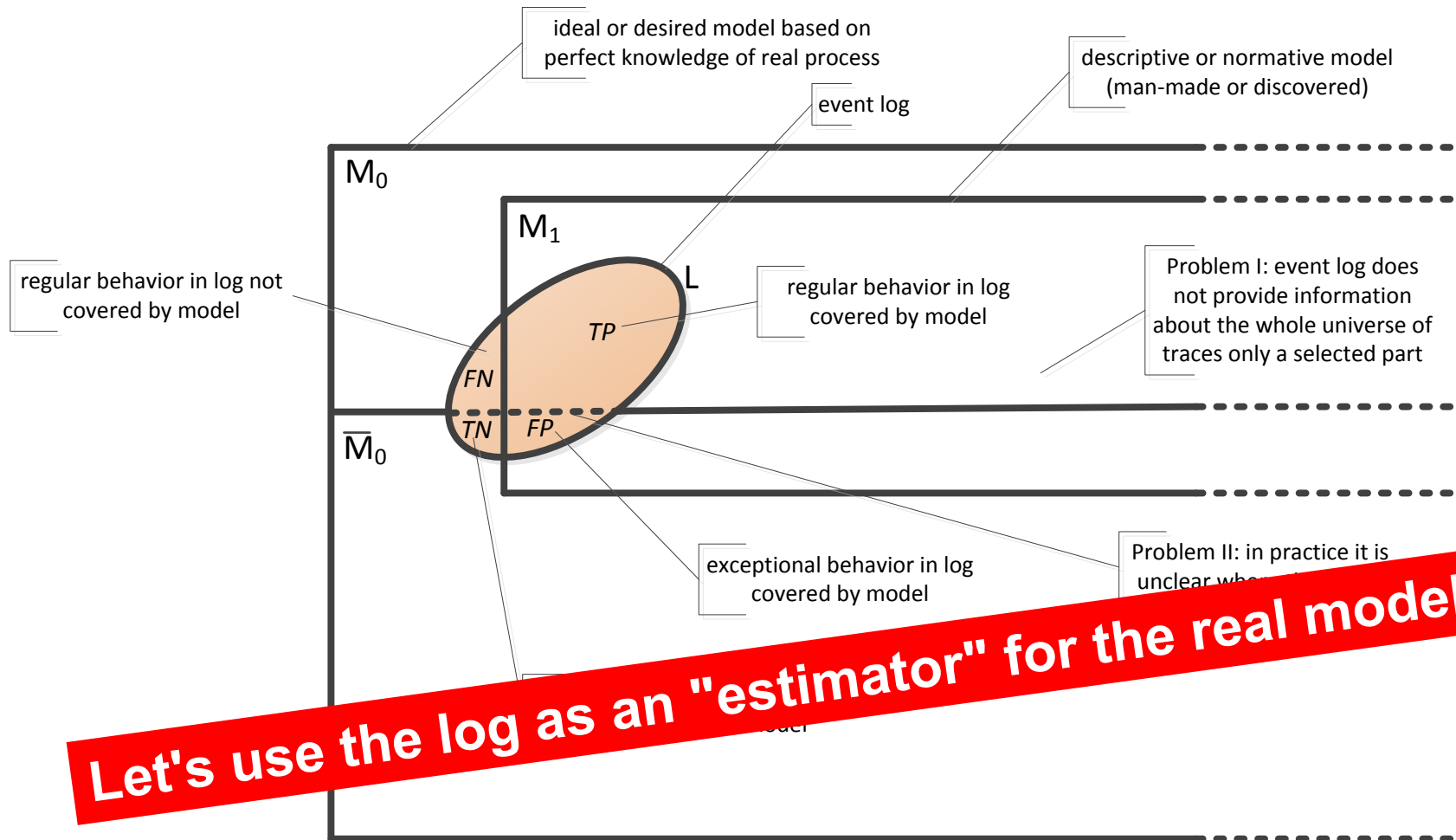
seems superfluous

Occam's razor

generalization

gravity

But we have event logs and do not know the real process ...

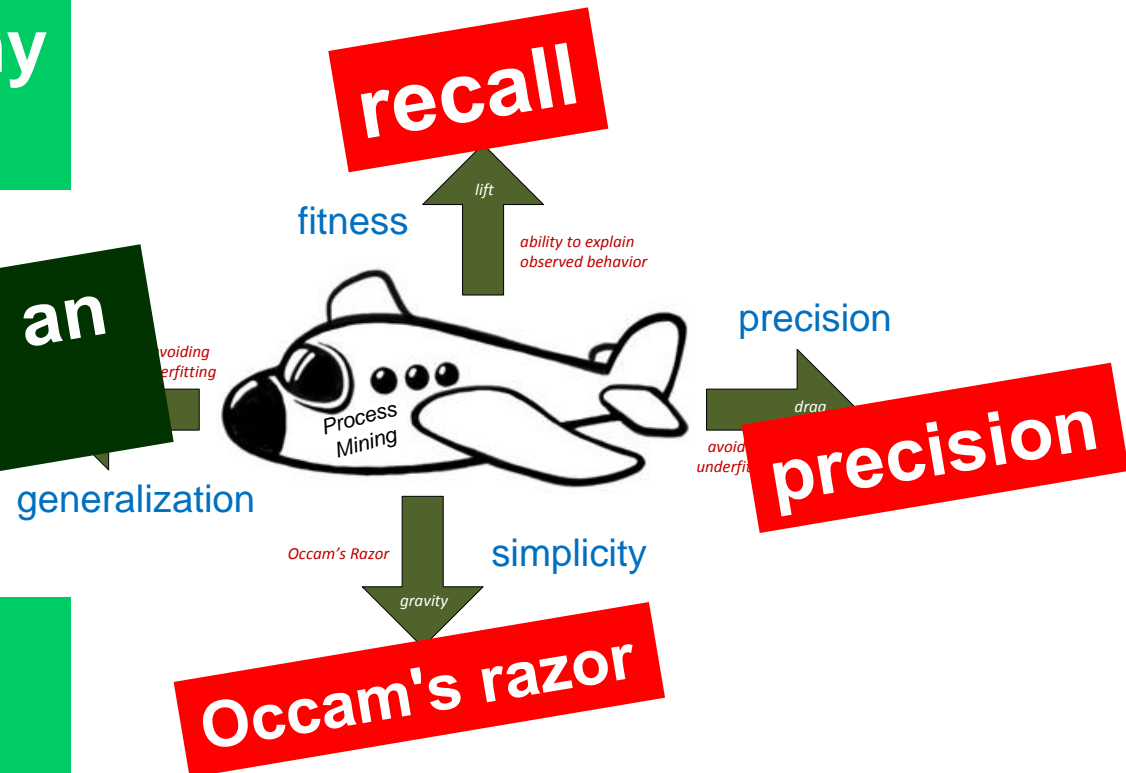


Considering event logs rather than the unknown model ...

How complete is my event log?

quality of log as an estimator

Have we seen all characteristic behaviors?



Why is process discovery difficult?



How about precision and recall?



What are the main research challenges in PM?



How to measure the quality of a process model?



The future is bright, but how to get started?



What is process mining?




What are the main pitfalls of process modeling?





Finding sheep with five legs

we are getting close...



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

On-the-fly
process mining



Operational
support

Concept drift





**cross-organizational /
comparative process mining**



**context aware
process mining**

Supporting the process of process mining



Why is process discovery difficult?



How about precision and recall?



What are the main research challenges?



How to measure the quality of a process model?



The future is bright, but how to get started?

What is process mining?



What are the main pitfalls of process modeling?



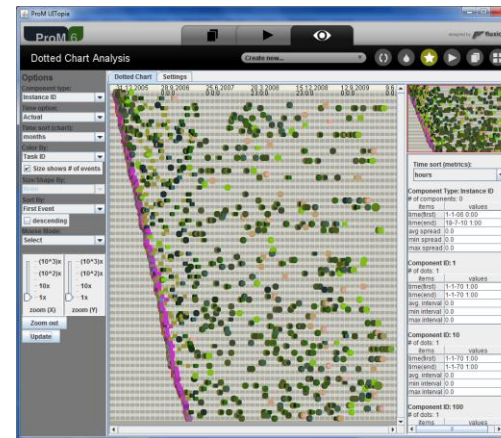
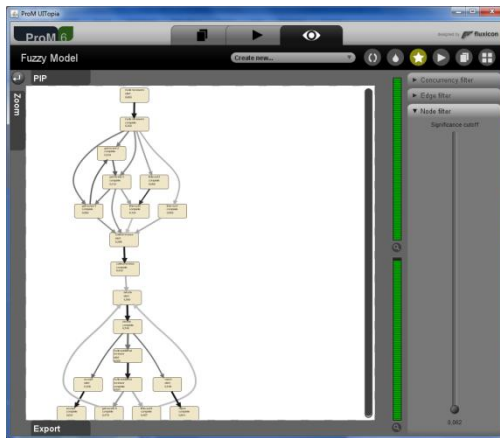
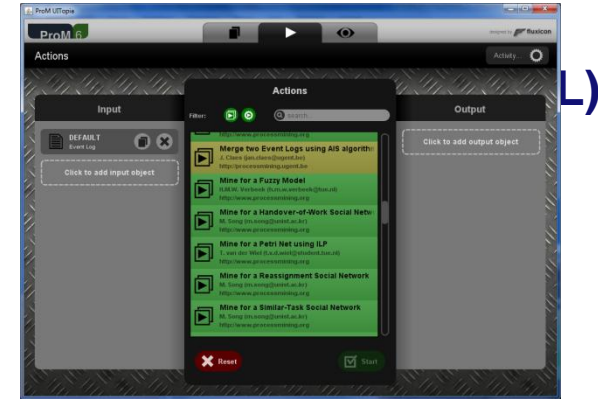
A yellow diamond-shaped sign with a black border and two mounting holes, one at the top and one at the bottom. The sign is mounted on a grey post. The background is a clear blue sky with a light gradient.

**BRIGHT
FUTURE
AHEAD**

How to get started?



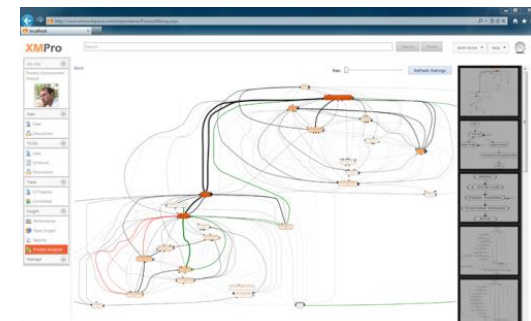
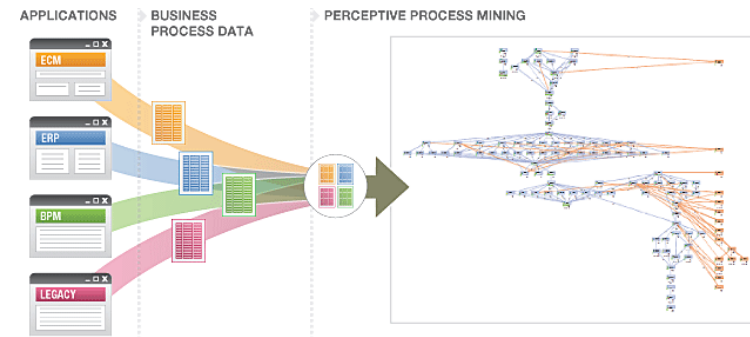
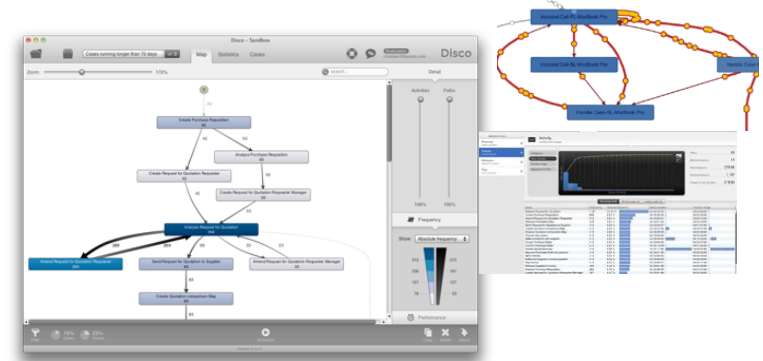
600+ plug-ins available covering the whole process mining spectrum



Download from: www.processmining.org

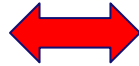
Commercial Alternatives

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



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: The Quest for the "Right" Process

