12 July 2019

EUROPEAN
COURT
OF AUDITORS

# Process Mining in the Wild

Zsolt VARGA, ECALab
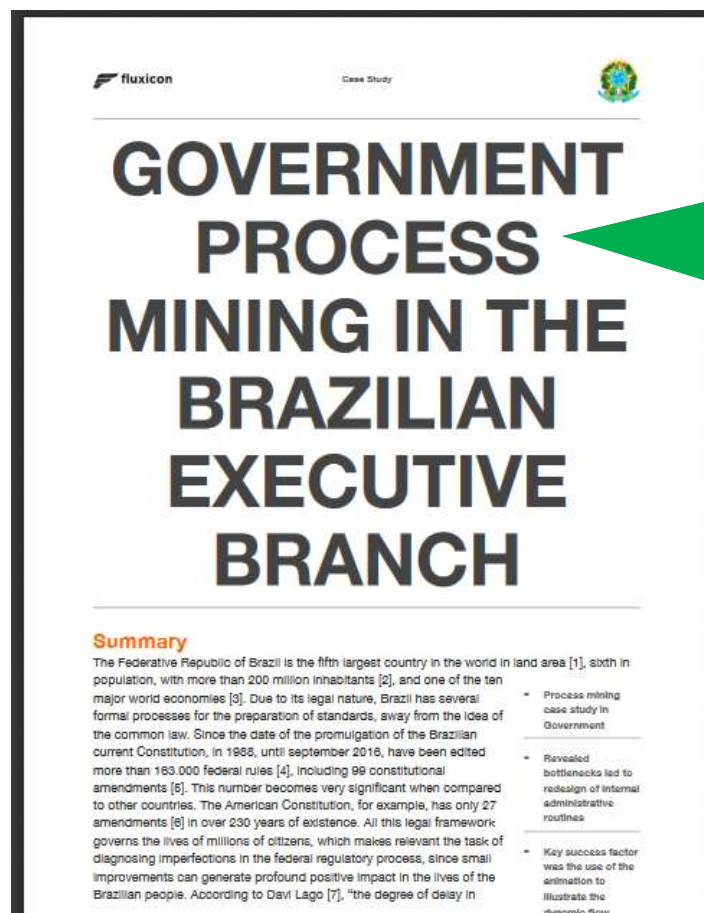
zsolt.varga@eca.europa.eu

# Process Mining Safari

- Case study 1: process mining as a new form of data representation
  - EU public consultations, a relatively simple case

- Case study 2: process mining for conformance checking
  - Audit of EU agencies, a bit more complicated
    - Getting the data, extracting from pdf files, data protection
    - Adding calculated fields in the event log
    - Fluxicon disco vs. Prom6

- Case study 3: exploratory process mining in „big data"
  - CA Service Desk manager database tables
    - Text mining methods for making sense of process data

EUROPEAN
COURT
OF AUDITORS

# Process mining in the government sector

http://romualdoalves.com/
decretos-2015-a-2018-
poc-text-mining/



Inspiration for using process mining in performance audit. It made me try PM on the audit on public consultations.

Process data is everywhere!

# Process Mining

# for performance audit

# Case study 1:

# European Public Consultations

# Process Mining Safari Part 1: playing it safe

# Audit case

- Audit question: How are public consultations implemented by the Commission?

- Dataset: Excel file about the various steps of public consultations and their start dates

- Business need: visualise the relative durations of PCs compared to each other and show deviations from the „ideal" process path

EUROPEAN
COURT
OF AUDITORS

# Input data for process mining

| Cons_ID | Event | Start | End_theoretical |
|---------|-------|-------|-----------------|
| PC-10 | Survey (National Authorities) | 2014-09-01 | 2015-02-01 |
| PC-10 | Conferences | 2014-11-01 | 2016-03-01 |
| PC-10 | Expert/focus groups | 2015-03-01 | 2016-02-01 |
| PC-1 | Roadmap | 2015-05-07 | 2015-06-04 |
| PC-18 | Survey (other) | 2015-06-26 | 2015-09-01 |
| PC-11 | Roadmap | 2015-07-23 | 2015-08-20 |
| PC-20 | Roadmap | 2015-09-28 | 2015-10-26 |
| PC-11 | Expert/focus groups | 2015-10-01 | 2016-05-01 |
| PC-10 | Workshops | 2015-10-01 | 2015-10-31 |
| PC-5 | Roadmap | 2015-11-30 | 2015-12-28 |
| PC-12 | Roadmap | 2015-12-08 | 2016-01-05 |
| PC-19 | Roadmap | 2015-12-08 | 2016-01-05 |
| PC-10 | OPC | 2016-01-12 | 2016-04-0 |
| PC-8 | Roadmap | 2016-01-18 | 20 |
| PC-9 | Conferences | 2016-02-01 | |
| PC-11 | OPC | 2016-02- | 31 |
| PC-11 | Survey (other) | | 2016-05-31 |
| PC-1 | OPC | 2016-03-01 | 2016-06-01 |
| PC-20 | Interviews | 2016-03-15 | 2016-08-24 |
| PC-14 | Expert/focus groups | 2016-04-01 | 2017-04-01 |

Open public consultations: events and start/completion dates for each event within an OPC.

Process data is everywhere!

EUROPEAN COURT OF AUDITORS

# Tableau viz, pretty close, but static, concurrencies shown as overlaps



Timeline of Events (normalized to OPC end date) (no labels)

Event type
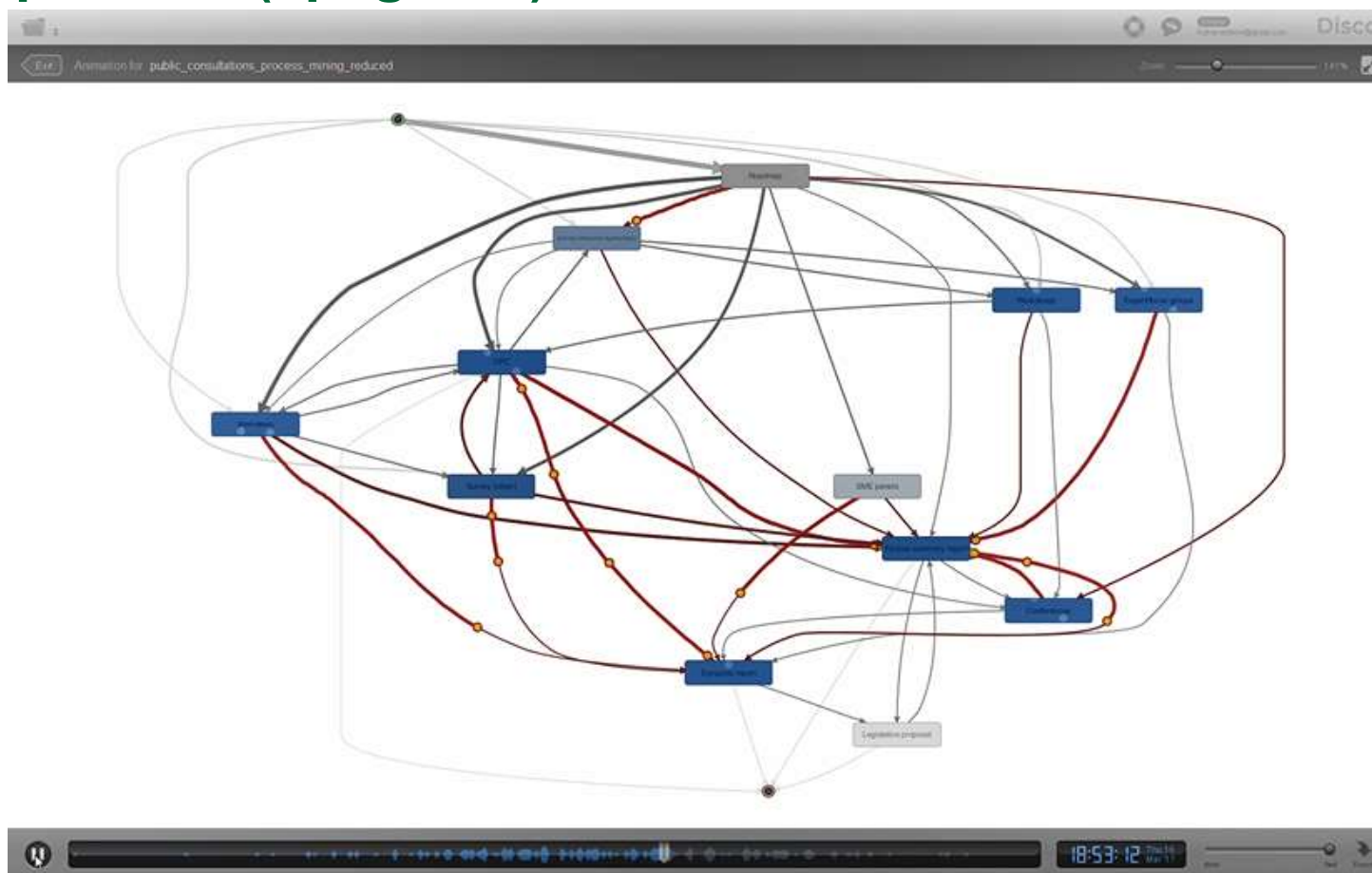- (All)
- Conferences
- Expert/focus groups
- ✓ Factual summary re...
- Interviews
- ✓ Legislative proposal
- ✓ Public consultation
- ✓ Roadmap
- SME panels
- Survey (National A...
- Survey (other)
- ✓ Synposis report
- Workshops

Event type
- Factual summary rep...
- Legislative proposal
- Public consultation
- Roadmap
- Synposis report

Caption

The reference point in the middle is the END date of the public consultation. The number of days since the start of the OPC is shown in the tooltip (positive on the right and negative on...) itself. I had to remove the bottom axis label, because a GANTT-chart can only contain real dates and the normalized date values would have been misleading.

Timeline of Events (normalized t...    Timeline of Events (normalized t...    Timeline of Events (normalized t...    **Timeline of Events (normalized...**

Data protection lesson 1: DO NOT show anything that is considered to be non-public or proprietary information

EUROPEAN COURT OF AUDITORS

# Full OPC process (spaghetti)

# Simplified OPC process (layered, „lasagne" process)

# Process conformance checking on automatically inferred model

Inductive Visual Miner plugin in ProM

# Process conformance checking on automatically inferred model



3 OPCs didn't have a synopsis report

2 OPCs had additional steps after the legislative proposal

# Lessons learned

- There's process data everywhere

- ProM/Disco is suitable for visualising performance audit results

- Replaying the log on the process model adds the time dimension by animating the path of the traces in the model

- A great way to communicate complex ideas quickly. When I demonstrated it to a visiting delegation, they immediately spotted the fact that two consultations didn't finish at the legal proposal (no need for prior briefing).

EUROPEAN COURT OF AUDITORS

# Audit of EU agencies

# Visualising and conformance checking of agency payments

# Process Mining Safari Part 2: a closer look at the beast

# Audit case

- Audit question: How can we automate the audit of EU agencies' payment processes?

- Dataset: the EC's accounting system has all the required data for certain types of agencies, but there are obstacles to directly getting the data from the database tables

- Business need: visual representation of the process workflow highlighting process deviation and non-compliance with duty segregation and deadlines

EUROPEAN
COURT
OF AUDITORS

# Small detour (between rock and a hard place)



Business area / audit chambers

IT people and database experts

ECALab / data scientist

# Small detour (between rock and a hard place)



Data protection officer

Business area / audit chambers

ECALab / data scientist

IT staff and database experts

# Access to data + data protection concerns

- Read-only access to the raw database of the Commission
    - Not much use without a data dictionary

# Relational database model

# What you are likely to get for process mining

# What you are likely to get for process mining



Data protection lesson 2: Stay away from personal data as much as possible!!

# What you are likely to get for process mining

# What you are likely to get for process mining

# Access to data + data protection concerns

- Read-only access to the raw database of the Commission
    - Not much use without a data dictionary

- Access to data warehouse
    - DW reports (that I've seen) do not contain workflow information

# Access to data + data protection concerns

- Read-only access to the raw database of the Commission
    - Not much use without a data dictionary
- Access to data warehouse
    - DW reports (that I've seen) do not contain workflow information
- Data may only be accessible through pdf exports
    - No meaningful database connection or custom legacy systems

# Access to data + data protection concerns

- Read-only access to the raw database of the Commission
    - Not much use without a data dictionary
- Access to data warehouse
    - DW reports (that I've seen) do not contain workflow information
- Data may only be accessible through pdf exports
    - No meaningful database connection or custom legacy systems
- Dataset contains user IDs and agency names
    - We are allowed to process such information for audit purposes, but I cannot show you real data in my presentation
    - I created a synthetic dataset based on the original, fully anonymised, with fictional agency names, and added some fictional anomalies

**Audit of EU agencies**

**Information extraction from pdfs**

# Why?

- PDF is evil. Although it is called a PDF "document", it's nothing like Word or HTML document.

- PDF is more like a graphic representation. PDF contents are just a bunch of instructions that tell how to place the stuff at each exact position on a display or paper.

- In most cases, it has no logical structure such as sentences or paragraphs and it cannot even identify a text box properly.

- PDFs require a disproportionate amount of resources to extract data in a structured format.

- PDFs are not meant for electronic processing. They are meant for printing!

EUROPEAN
COURT
OF AUDITORS

# Information extraction pipeline

```
┌─────────────────┐      ┌─────────────────┐      ┌─────────────────┐
│ ┌─────────────┐ │      │ ┌───────────┐ ┌─┤ │ ┌─────────────────┐ │
│ │ Pre-process │ │  →   │ │  Perform  │ │ E│ │ │ Consolidate in a│ │
│ │ documents   │ │      │ │  pattern  │ │→x│ │ │    machine      │ │
│ │             │ │      │ │  matching │ │ t│ │ │ readable format │ │
│ └─────────────┘ │      │ └───────────┘ └─r┤ │ └─────────────────┘ │
└─────────────────┘      └─────────────────┘ └─────────────────────┘
```

Pre-process documents → Perform pattern matching → Extract info from matches → Consolidate in a machine readable format

- Regular expressions
- Linguistic features
- Always corpus specific

EUROPEAN COURT OF AUDITORS

# Agency payment transactions

**Document Location**

| Local Key | Doc. List | Date | Doc. Location | Doc |
|-----------|-----------|------|---------------|-----|
| | AUT | 18/12/2017 | | Note |
| | | | | Budg |
| | AUT | 11/01/2018 | | COM |
| | | | | Mon |
| | AUT | 18/06/2018 | | COM |
| | | | | Mon |
| | AUT | 26/10/2018 | | COM |
| | | | | Mon |

**Potential Abnormal RAL**

Old Commitment (Y/N): N

Sleeping Commitment (Y/N): N

**Workflow History**

| Date | Action Taken | Step Desc. |
|------|--------------|------------|
| • 29/10/2018 09:56:23 | ACCEPT UNATTENDED-AUTOMATIC - AU | SAP R/3 - ACCOUNTANT |

Model : *** Standard WF ***
Workflow Org :
Workflow Center
Comments: Attach file from SI2 SB

| • 29/10/2018 09:56:25 | TECHNICAL ACCEPTANCE - TA | SAP R/3 - ACCOUNTANT | SAPR3 |

Model : *** Standard WF ***
Workflow Org :
Workflow Center
Comments: - Document has been accepted with visa FA .

**European Commission**
**Directorate General Budget**

**Commitment Level 2**

**Summary Information**

ABAC Internal Key
Applicable Regulation
Central Key
Type
Payment Class
User Reference
Old Responsible User
File Reference
Project
Prop. FDC ILC
Proposed FDI
Contract Sign. Date
Exception to Default FDC ILC
Justification for PP/RM : Exp
Reason updating expired FDI respecting Art. 114 FR 2018 :
EDES Justification :NA : Not applicable

Workflow :          Status: FIN          Level  100

Remarks (Y/N):  Y

cal

# A possible pdf data extraction workflow

- Camelot + pdfMiner libraries for Python
  - only works with text-based PDFs and not scanned documents
  - each table is a pandas DataFrame, which seamlessly integrates into ETL and data analysis workflows.

- My recommendation after much sweat and blood:
  - use pdfMiner to find and identify the part of the document containing the table you need
  - Use Camelot-py with the parameters you got from pdfMiner to extract the table
  - each dataset requires fine-tuning, exception handling and post-processing

EUROPEAN
COURT
OF AUDITORS

# Python code for data extraction



```python
def ExtractTablefromPDF (pdf_file, table_title):
    #coord = getCoordinatesofText ('Workflow History',        .pdf')
    coord = getCoordinatesofText (table_title, pdf_file)
    start_page = coord[0][0]
    totalpages = coord[0][1]
    y_coord = coord[0][2]
    table_area_list = []
    df = pd.DataFrame(columns=[0, 1, 2, 3, 4])

    for index in range(start_page, totalpages+1):
        #first let's try to get the first part of the table
        if index == start_page:
            table_area = '10,'+str(int(y_coord+40))+',800,70'
            table_area_list.append(table_area)
            table_first_part = camelot.read_pdf(pdf_file, str(index), flavor='stream', table_areas = table_       lumns=['72,250,35
            df = df.append(table_first_part[0].df)
#and try to get the full page for the rest of the pages
        else:
            #print (index)
            table_rest = camelot.read_pdf(pdf_file, str(index), flavor='stream', columns=['72,250,350,450'
            df = df.append(table_rest[0].df)


    columnlist = list(df.iloc[0]) # take the original columns
    columnlist.append('Model')
    columnlist.append('Workflow Org')
    columnlist.append('Workflow Center')
    columnlist.append('Comment')
    columnlist.append('Time')
    columnlist.append('Case_ID')

    df_out = pd.DataFrame(columns=columnlist)
    df = df.reset_index()
    df = df.drop('index', axis = 1)
```

A bit of hacking is needed to customise the extraction to the specificities of the dataset

## DISCLAIMER:

- The data, all agency names, anomalies and incidents portrayed in this dataset and process mining experiment are fictitious. No identification with actual persons (living or deceased), positions, organisations and events is intended or should be inferred.

- Even though the structure of the data closely follows that of the European Commission's accounting system, this presentation does not imply that the EC's control systems allow for any of the shown anomalies to occur.

# ABAC workflow data

| | CaseID | Time | Action Taken | Step Desc. | Person | Type of assignment | Signed As Other Agent | Same agents | Model |
|---|---|---|---|---|---|---|---|---|---|
| 1 | CaseID | Time | Action Taken | Step Desc. | Person | Type of assignment | Signed As Other Agent | Same agents | Model |
| 2 | GCB.2871 | 07/06/2017 11:53 | ACCORD - AC | INITIATING AGENT (first) | AGENT09 | FA | AGENT49 | No | Light version + 2A Level Ex Ante |
| 3 | GCB.2871 | 07/06/2017 11:53 | TECHNICAL ACCEPTANCE - TA | SAP R/3 | SAPR3 | | | No | Light version + 2A Level Ex Ante |
| 4 | GCB.2871 | 13/06/2017 11:27 | ACCORD - AC | VERIFYING AGENT (first) | AGENT15 | FA | AGENT01 | No | Light version + 2A Level Ex Ante |
| 5 | GCB.2871 | 13/06/2017 17:13 | Accepted without verification L2 ex ante | VERIFYING AGENT (EX-ANTE) | AGENT30 | | | No | Light version + 2A Level Ex Ante |
| 6 | GCB.2871 | 15/06/2017 14:20 | ACCORD - AC | AUTHORISING OFFICER | AGENT01 | | | No | Light version + 2A Level Ex Ante |
| 7 | GCB.2871 | 15/06/2017 14:20 | ACCEPT UNATTENDED-AUTOMATIC - AU | SAP R/3 - ACCOUNTANT | WF-BATCH | | | No | Light version + 2A Level Ex Ante |
| 8 | GCB.2871 | 15/06/2017 14:20 | TECHNICAL ACCEPTANCE - TA | SAP R/3 - ACCOUNTANT | SAPR3 | | | No | Light version + 2A Level Ex Ante |
| 9 | GCB.2919 | 13/07/2017 15:23 | ACCORD - AC | INITIATING AGENT (first) | AGENT44 | FA | AGENT44 | Yes | *** Standard WF *** |
| 10 | GCB.2919 | 13/07/2017 15:23 | TECHNICAL ACCEPTANCE - TA | SAP R/3 | SAPR3 | | | No | *** Standard WF *** |
| 11 | GCB.2919 | 13/07/2017 18:54 | ACCORD - AC | VERIFYING AGENT (first) | AGENT41 | FA | AGENT42 | No | *** Standard WF *** |
| 12 | GCB.2919 | 13/07/2017 18:59 | ACCORD - AC | AUTHORISING OFFICER | AGENT24 | | | No | *** Standard WF *** |
| 13 | GCB.2919 | 13/07/2017 19:00 | ACCEPT UNATTENDED-AUTOMATIC - AU | SAP R/3 - ACCOUNTANT | WF-BATCH | | | No | *** Standard WF *** |
| 14 | GCB.2919 | 13/07/2017 19:00 | TECHNICAL ACCEPTANCE - TA | SAP R/3 - ACCOUNTANT | SAPR3 | | | No | *** Standard WF *** |
| 15 | GCB.2919 | 05/04/2018 09:53 | ACCORD - AC | INITIATING AGENT (first) | AGENT43 | FA | AGENT43 | Yes | *** Standard WF *** |
| 16 | GCB.2919 | 05/04/2018 09:53 | TECHNICAL ACCEPTANCE - TA | SAP R/3 | SAPR3 | | | | |
| 17 | GCB.2919 | 05/04/2018 10:53 | ACCORD - AC | VERIFYING AGENT (first) | AGENT41 | FA | AGENT42 | | |
| 18 | GCB.2919 | 05/04/2018 11:01 | ACCORD - AC | AUTHORISING OFFICER | AGENT06 | | | | |
| 19 | GCB.2919 | 05/04/2018 11:01 | ACCEPT UNATTENDED-AUTOMATIC - AU | SAP R/3 - ACCOUNTANT | WF-BATCH | | | | |
| 20 | GCB.2919 | 05/04/2018 11:01 | TECHNICAL ACCEPTANCE - TA | SAP R/3 - ACCOUNTANT | SAPR3 | | | | |
| 21 | GCB.2940 | 26/10/2017 16:22 | ACCORD - AC | INITIATING AGENT (first) | AGENT39 | FA | AGENT39 | | |
| 22 | GCB.2940 | 26/10/2017 16:22 | TECHNICAL ACC | | | | | | |
| 23 | GCB.2940 | 31/10/2017 10:57 | REFUS POUR CC | | | | | | |
| 24 | GCB.2940 | 07/11/2017 09:51 | REFUS POUR CC | | | | | No | Light version + 2A Level Ex Ante |
| 25 | GCB.2940 | 07/11/2017 09:51 | REFUS POUR CC | | | | | No | Light version + 2A Level Ex Ante |

Calculated field, shows if the same agent used another function

Process log created on the basis of data extracted from 20 pdf files, semi-randomly chosen by auditors among various agencies

EUROPEAN COURT OF AUDITORS

# ABAC workflow data

Frequency based display, the darker the more frequent

Full process model based on 20 transactions. Looks good!

EUROPEAN COURT OF AUDITORS

# ABAC workflow data



Displaying mean duration, with minimum duration as a secondary metric. Outliers are interesting!

Durations: first step of performance checking

EUROPEAN COURT OF AUDITORS

# Compliance check: segregation of duties

# Compliance check: segregation of duties



Violation of segregations of duties occurs twice in the population, and at the same agency, and for cases having a total duration over a year

Agent 5 acted as verifying AND authorising officer. Maybe he just wanted to speed up the closure of long cases?

# Compliance check: suspiciously short approval times

Filters by subsequences

Filter by: Activity ▼     Reference event must be    eventually followed ▼    by a follower event.

Reference event values: (1 of 13 selected)

| ✓ | ACCORD - AC+VERIFYING AGENT (first) |
| ✗ | Accepted without verification L2 ex ante+VERIFYING AGENT (EX-ANT |
| ✗ | REFUS POUR CORRECTION - SC+INITIATING AGENT (After REFUS |
| ✗ | REFUS POUR CORRECTION - SC+INITIATING AGENT (first) |
| ✗ | REFUS POUR CORRECTION - SC+SAP R/3 |
| ✗ | REFUS POUR CORRECTION - SC+SAP R/3 - ACCOUNTANT |
| ✗ | REFUS POUR CORRECTION - SC+VERIFYING AGENT (first) |
| ✗ | TECHNICAL ACCEPTANCE - TA+SAP R/3 |

Follower event values: (1 of 13 selected)

| ✓ | ACCORD - AC+AUTHORISING OFFICER |
| ✗ | ACCORD - AC+INITIATING AGENT (first) |
| ✗ | ACCORD - AC+VERIFYING AGENT (first) |
| ✗ | Accepted |
| ✗ | REFUS P |
| ✗ | REFUS P |
| ✗ | REFUS P |
| ✗ | REFUS P |

☐ Require  the same value ▼  of  Resource ▼  for each pair of events matched above.

☑ Time between matching events must be  shorter ▼  than  2 hours ▼

> Use the Follower filter and set a time limit. You can stack filters on top of each others (filter by variant, attributes, duration, you name it) to drill into the population by criteria

# Compliance check: suspiciously short approval times

# Compliance check: suspiciously short approval times



Now, this is interesting: VA signed in as the AO followed by AO approval within an hour?

Agent27 works for the imaginary agency JFG.

# Creating a social network based on the event log



ProM is a framework with a large number of plugins. A social network graph is just another representation of the same data. In a conformance checking context, it is very useful to examine links within „populations of interest" that you isolate from the full population

# Import/conversion settings are important!!!



You need to specify the column containing person names or users at conversion time by cicking on Show Expert Configuration and choose "org:resource"

In ProM you have to import the .csv file first, and then convert it into an event log file!

# Creating a social network based on the event log



Our friend, Agent27 from the previous exercise

Working together social network in ProM, using the same agency dataset

These might be authorising officers, who act as gateways between normal employees and the accounting system (the transactions go through them for processing, but they don't otherwise interact with other agents)

# Replaying full event log on simplified process model



First create a fuzzy model based on the log and export it in ProM (and maybe add artificial start and end points)

# Replaying full event log on simplified process model



Replaying the full event log on a previously agreed model provides useful conformance and performance information (red, green, blue counters)

# Conformance checking in ProM

- There is a separate conformance checking plug-in in ProM that I have not covered in this presentation due to lack of time
- It works a bit differently than what I've shown in the previous slides
- Once a process model is agreed between auditor and auditee, the conformance checking results could be considered as audit evidence
- The ProM framework also allows you to write your own plugin in Java, so if you need checks that the system cannot currently do for you, you can implement them yourself in a plugin

# Lessons learned

- In the worst case scenario, we can still extract useful process data from pdf files for prototyping (still needs manual touch-up, not suitable for production)

- We don't need much data to build the process model if the process is well-structured, 20 transactions gave us a useable model

- ProM/Disco is suitable for conformance checking, with certain limitations

- Some pre-processing is required to add event features (e.g. same agent with different logins/functions)

- No interactive display of non-conformances in Disco, need to use filters, Celonis might be a better option for this purpose (under investigation)

EUROPEAN
COURT
OF AUDITORS

# Mining Service Desk log files for process improvement

# Process Mining Safari Part 3: being run over by an elephant

# Business case

- Audit question: Explorative analysis of business processes with a view to potential process improvements

- Dataset: MS SQL server export from CA Service Desk Manager, containing only the keys, not the resolution for persons names, requestors etc.

- Business need: to gain experience with data extraction and process mining software, determine what can be done within PM software, what pre-processing steps may be required for client systems, etc.

# My expectations

- Easy win, we just retired the system, accessing historical data on a non-live system should not be complicated

- It's a service desk management system, processes should be well-defined

- I could present something useful relatively quickly to convince people about the advantages of process mining

# Reality

# The reality

- According to the DPO everything had to be anonymised, I did not extract ANY personal data so that we wouldn't have to do extra administration
- Even the case IDs had to be encrypted, so they could not be traced back to actual persons
- Almost as many variants in the original as many cases
- French AND English event descriptions which had to be merged
- Emails had to be deleted using regexp, proper names had to be removed using Spacy EN/FR language models (NER)
- Explorative analysis had to determine the best splits of the data

# Pre-processing in Python, merging similar events

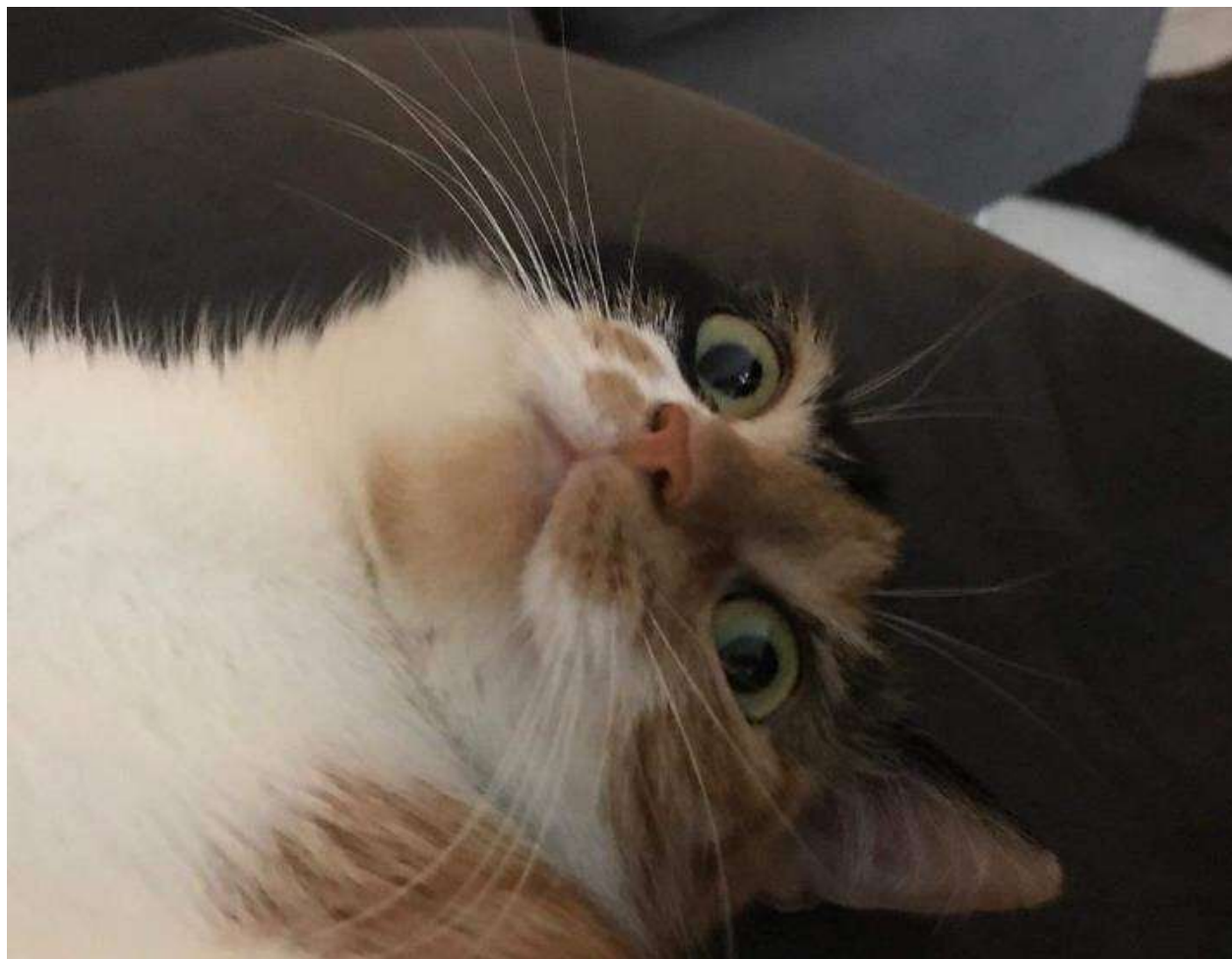# My face when I first saw the resulting process graph



#màzlithecat

# My face when I first saw the resulting process graph

# My face when I first saw the resulting process graph

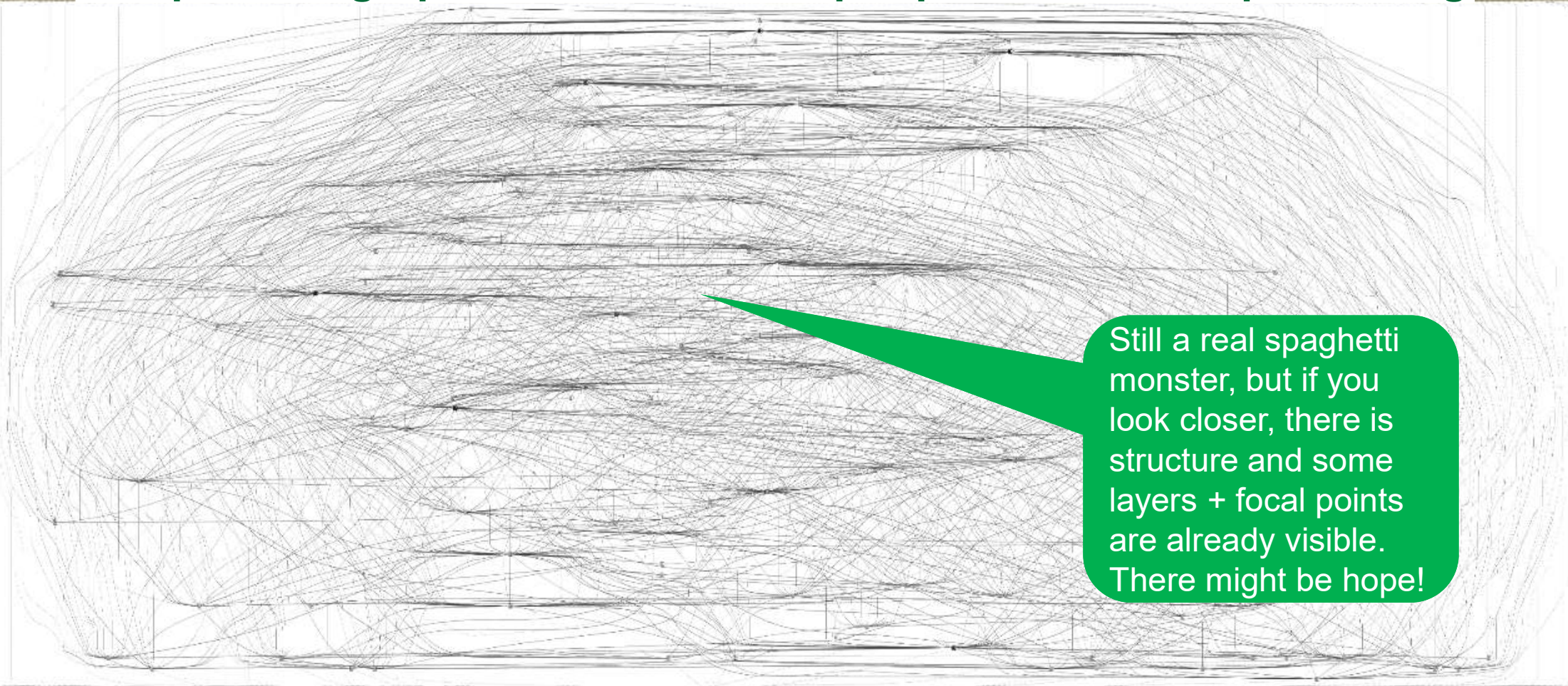# My face when I first saw the resulting process graph

# My face when I first saw the resulting process graph

# My face when I first saw the resulting process graph

# The process graph derived from the pre-processed „simplified" log
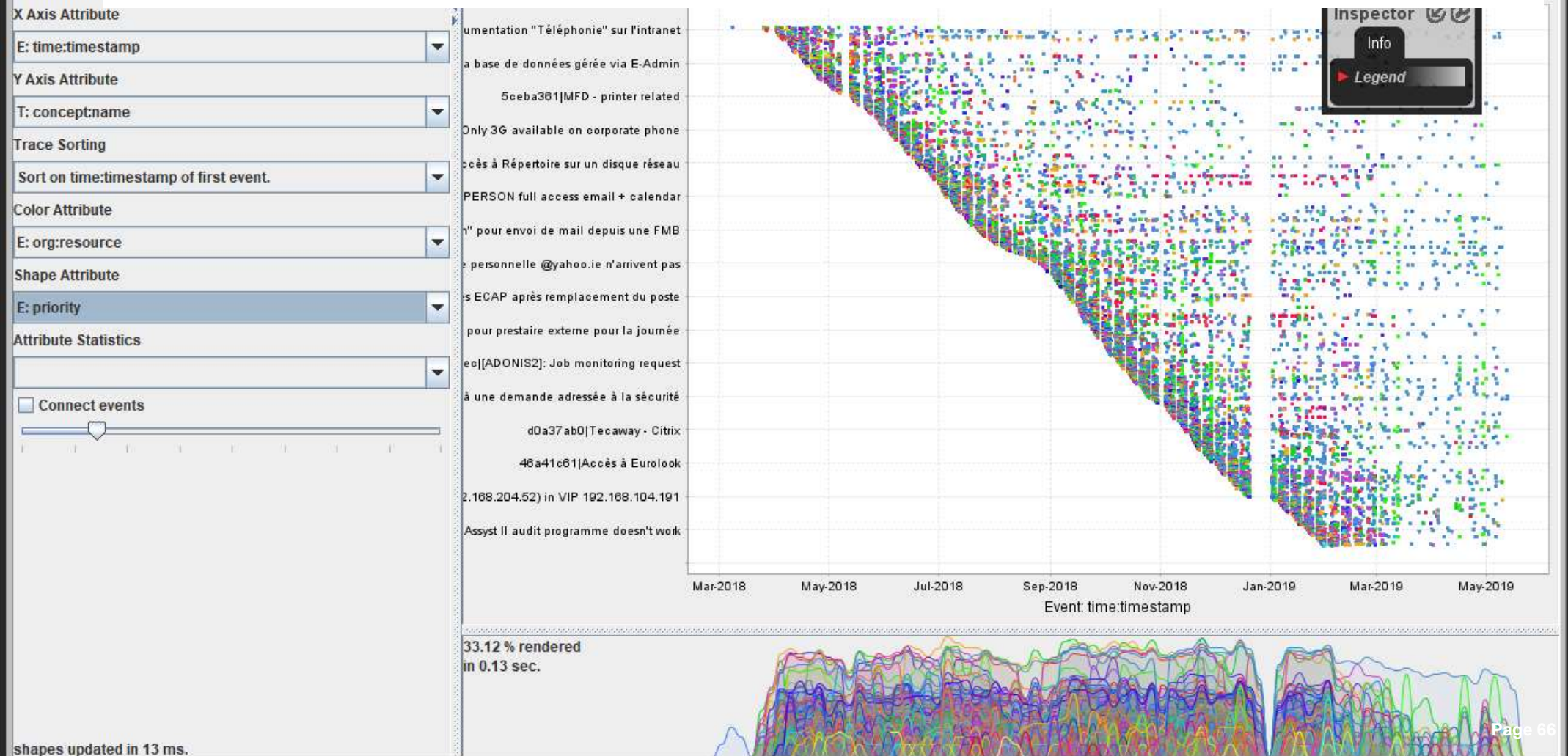


Still a real spaghetti monster, but if you look closer, there is structure and some layers + focal points are already visible. There might be hope!

# How to make lasagne out of spaghetti

- It might as well be as difficult as turning water into wine

- Slice and dice! Explore the dotted chart, correlations, histograms, patterns.

- It might be the case that the log contains not one, but several processes, so you need to separate them!

- Talk to the business area and develop an intuition for recognising recurring patterns.
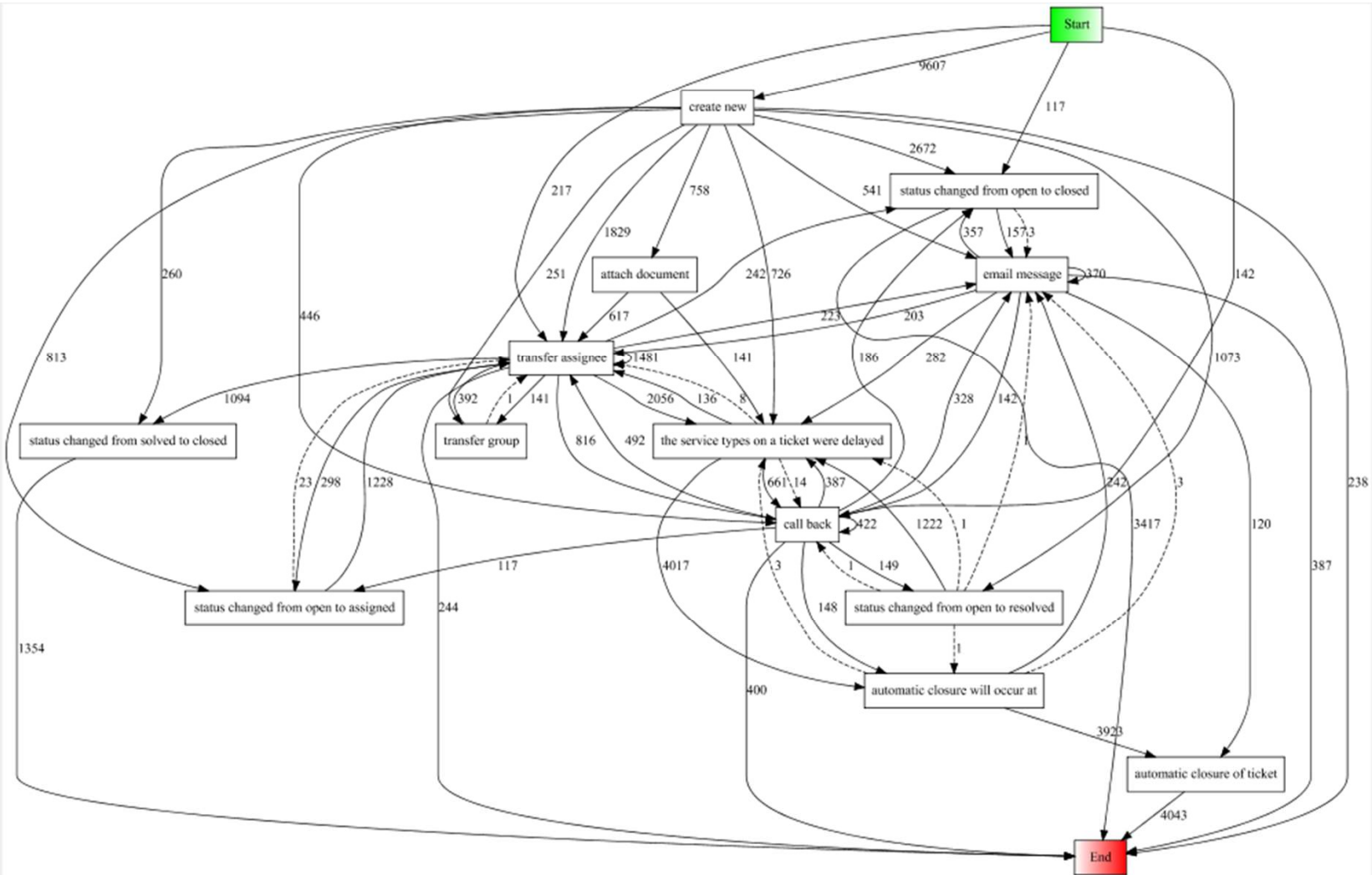
## Good old Pareto principle
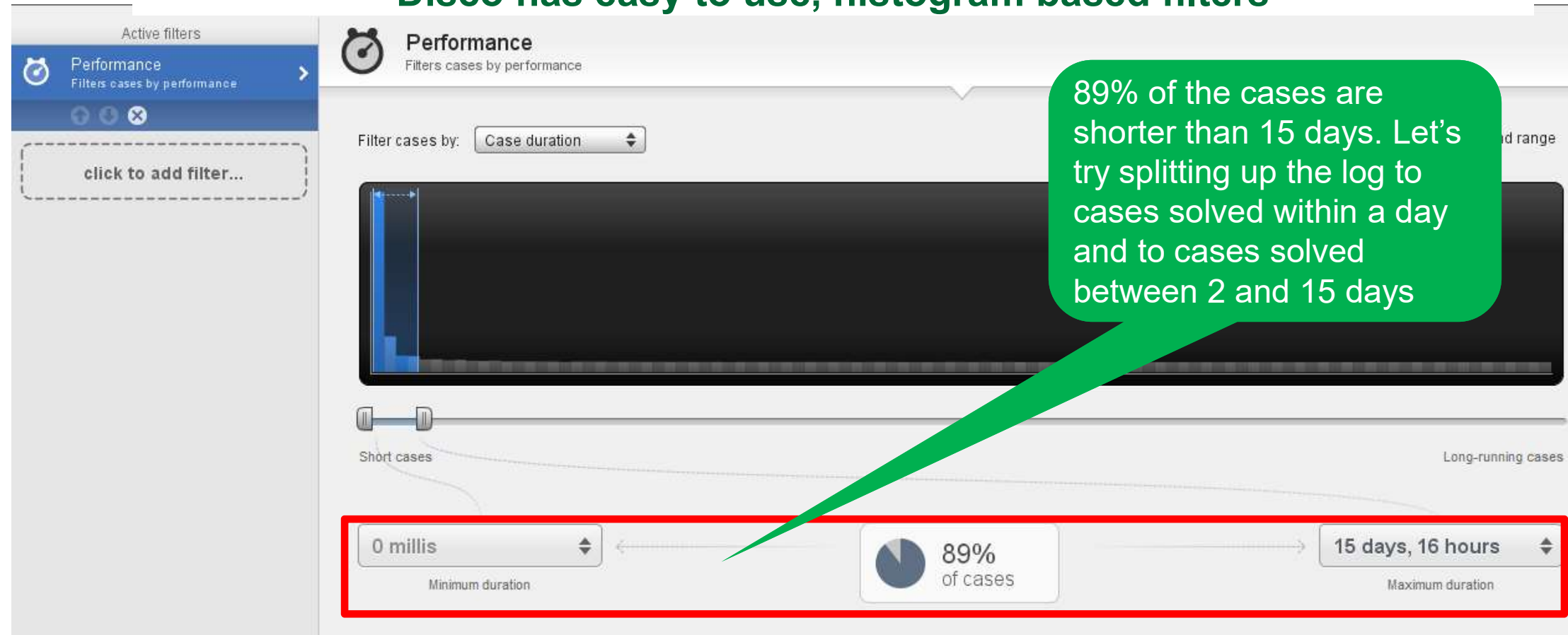
# 80/20 rule

20% of the cases/events cause 80% of the variance
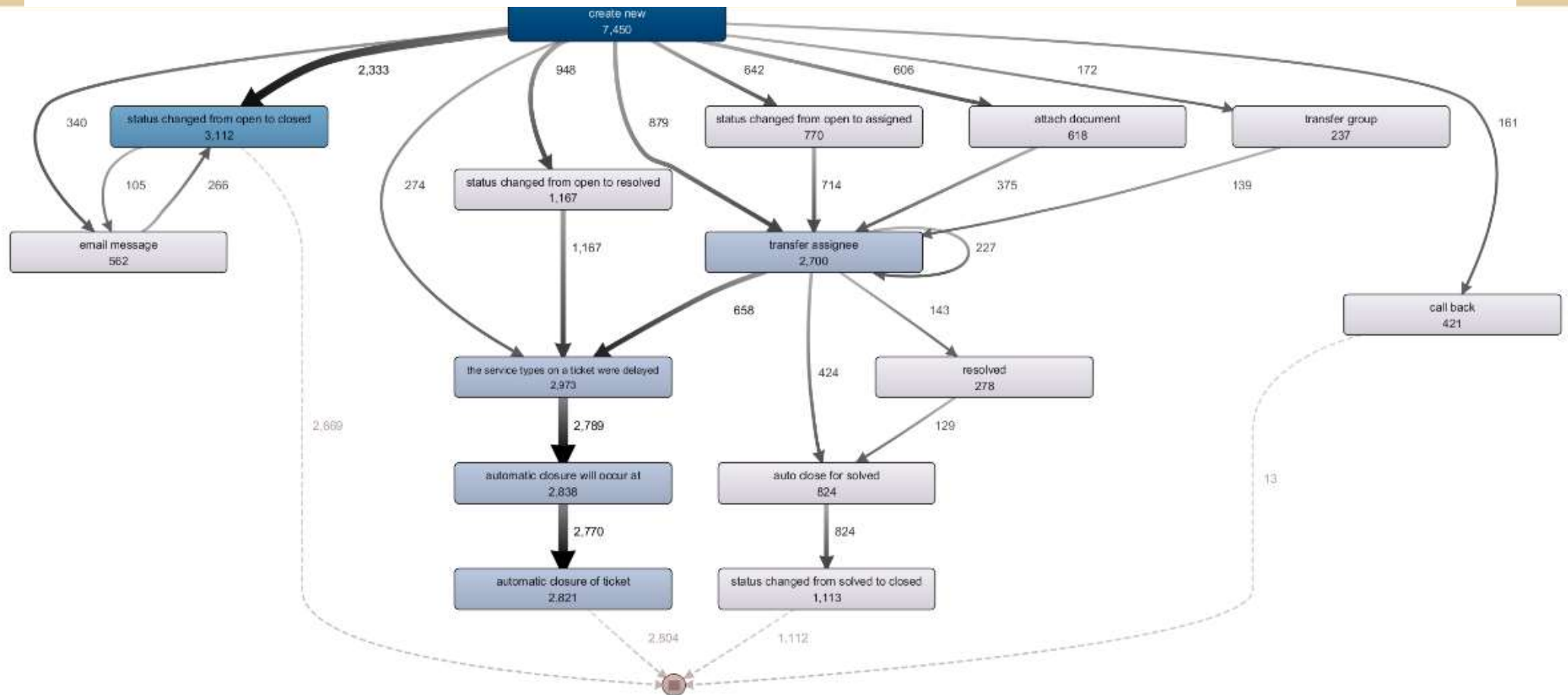
# First look with ProM and the dotted chart view

# Directly follows graph

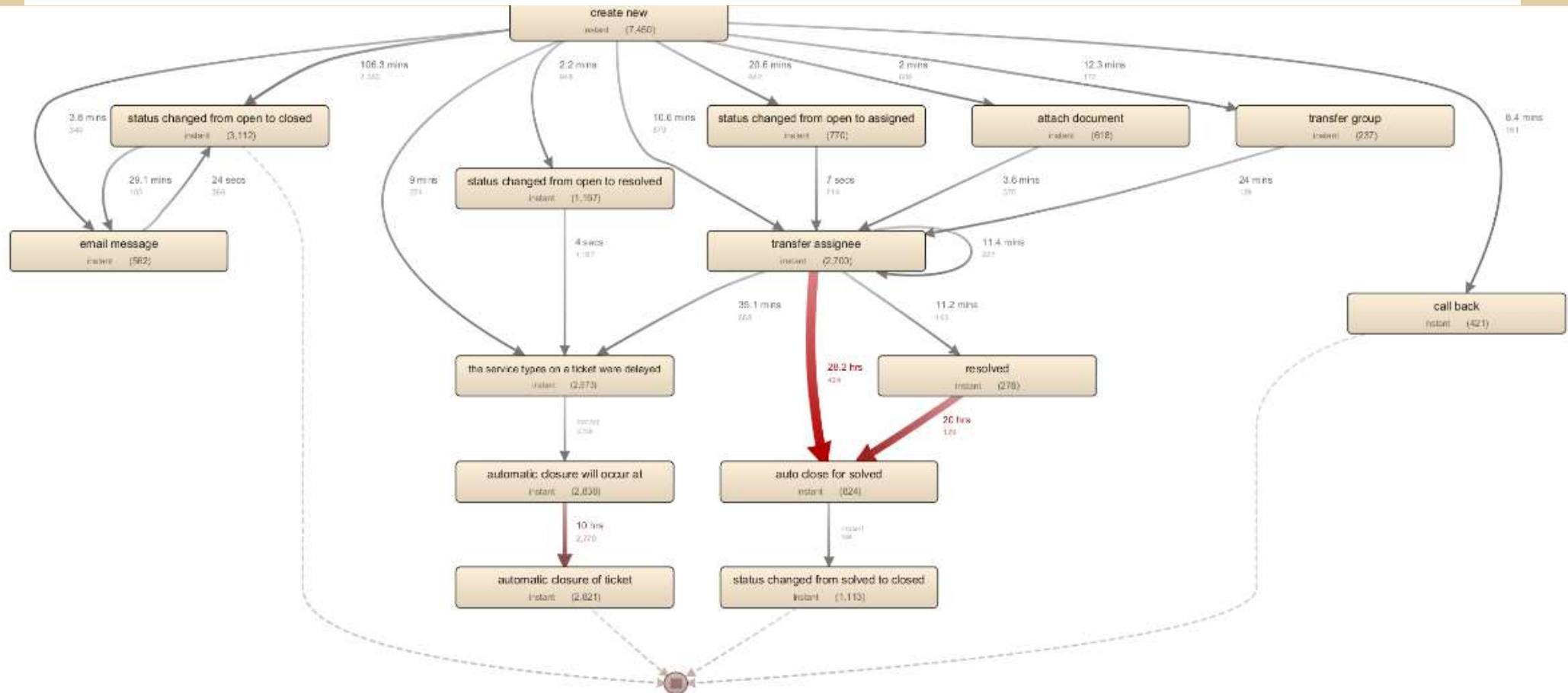# Disco has easy to use, histogram based filters



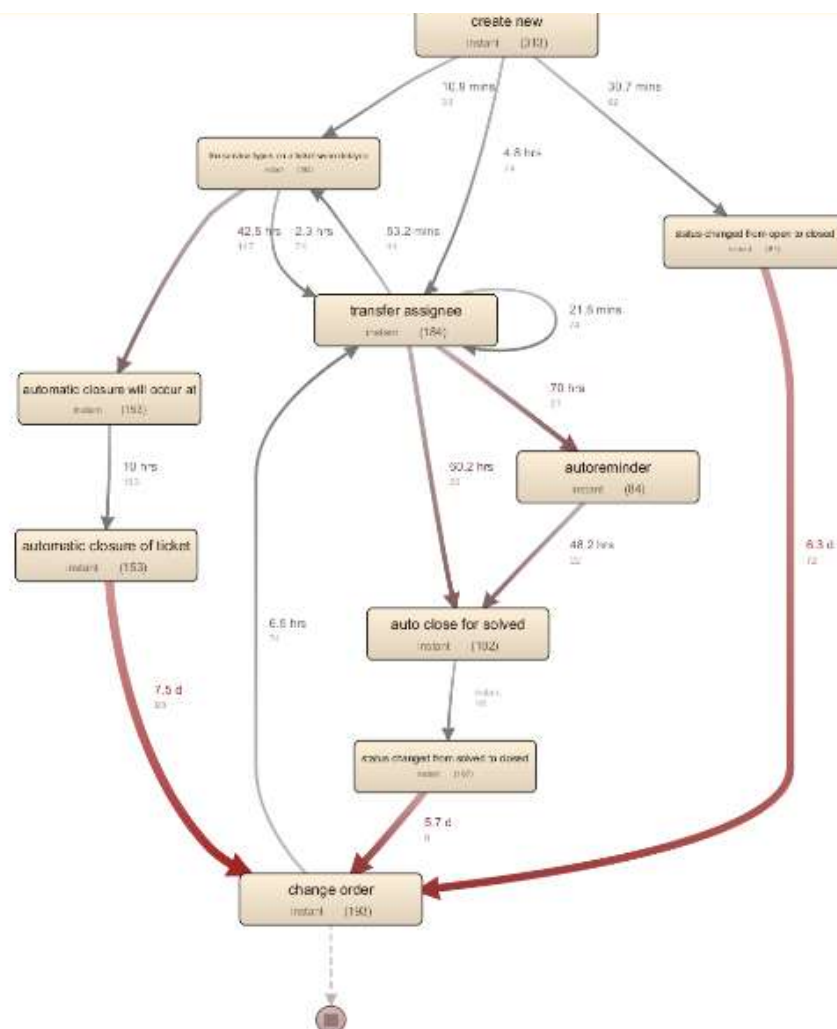89% of the cases are shorter than 15 days. Let's try splitting up the log to cases solved within a day and to cases solved between 2 and 15 days

# Split by duration (within a day)

# Split by duration (within a day)

# Agents vs time of activity

# Duration, category vs time of the week

# Lessons learned

- DO NOT assume that an event log will be easy to understand

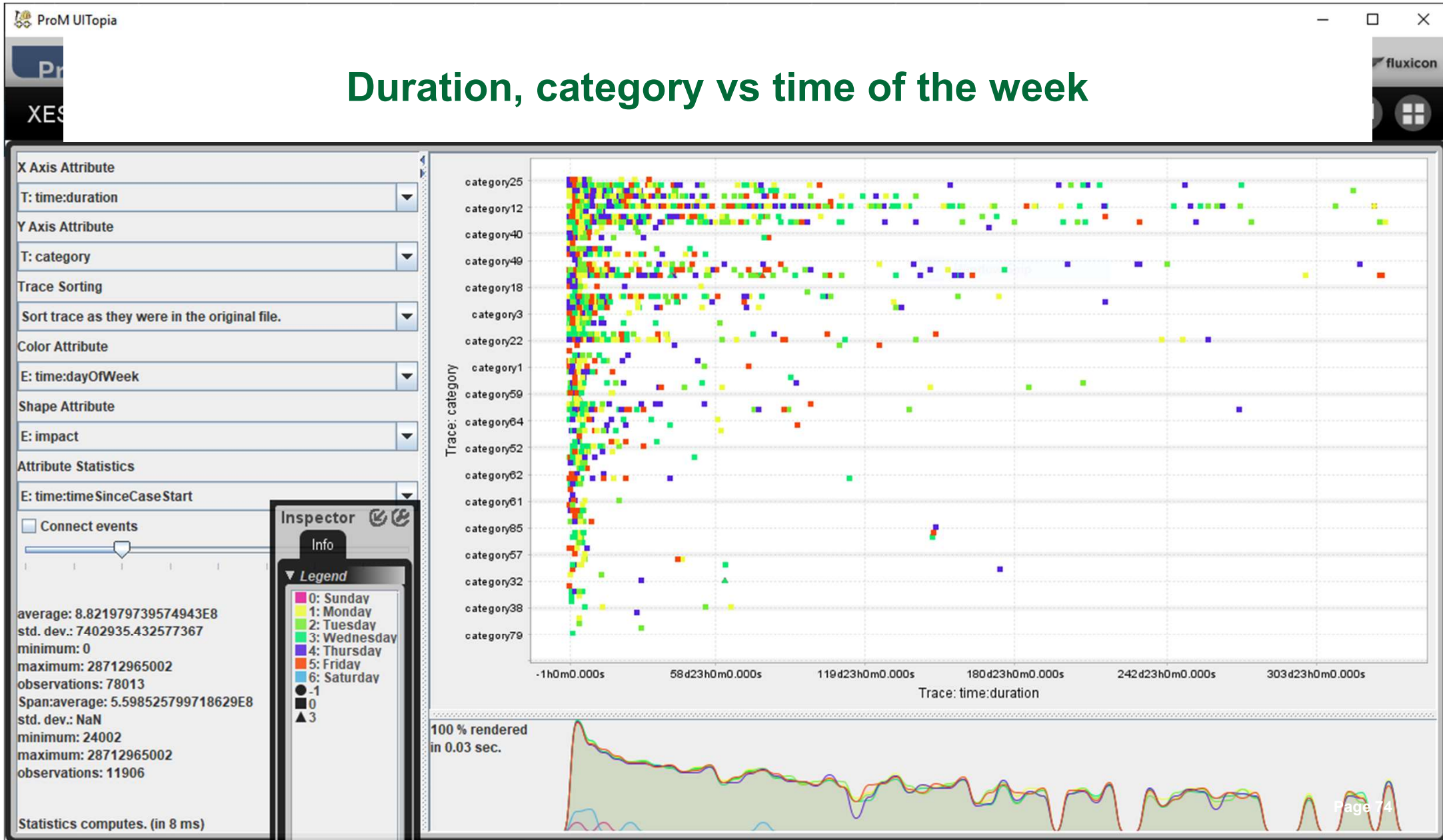- Use text mining tools to anonymise data, merge duplicate event names and quantify some of the unstructured textual data contained in description fields

- Be careful! Log content may contain personal data in fields where you are not expecting it! (e.g.: email text copied into description fields)

- Use your intuition and explore correlations and histograms

- The dotted chart view of ProM is very useful to get a „feel" for the data, and Disco has very nice histograms to help you slicing up the log data

# Recommendations

- Start small: download a copy of ProM or Fluxicon Disco (demo) or register for Celonis cloud or contact Minit

- Start small: find an existing dataset that may be a good fit for process mining explorations

- Watch the Process Mining course on Coursera (https://www.coursera.org/learn/process-mining) and the ProM course on FutureLearn (https://www.futurelearn.com/courses/process-mining).

- Read Wil's book (https://www.springer.com/gp/book/9783662498507) and website (https://processmining.org)

- Read Fluxicon's online book (http://processminingbook.com/index.html)

- Secret tip: The Flux Capacitor blog (https://fluxicon.com/blog) especially the older entries between 2011-2014 contain a lot of tips and tricks

# Recommendations

- Think big: if you are passionate you CAN make a difference
- Think big: this is a unique opportunity where public sector audit can be ahead of the private sector

- Act local: find out how automation/process mining can help your colleagues
- Act local: find out how your organisation may benefit from this

- Be realistic: it may or may not work out for you
- Be realistic: the technology might not be mature enough for your use case, but it might be overnight. Anticipate the change!

EUROPEAN
COURT
OF AUDITORS

# QUESTIONS?

Email: zsolt.varga@eca.europa.eu

LinkedIn: www.linkedin.com/in/zsvarga/

GitHub: github.com/zseebrz

Tableau Public: public.tableau.com/profile/zsolt.varga#!/