# Significant Properties of Spreadsheets
## An Update On The Work Of The Open Preservation Foundation's Archives Interest Group

Remco van Veenendaal
National Archives of the Netherlands

Frederik Holmelund Kjærskov
Danish National Archives

Kati Sein
National Archives of Estonia

Jack O'Sullivan
Preservica

Anders Bo Nielsen
Danish National Archives

Phillip Mike Tømmerholt
Danish National Archives

Jacob Takema
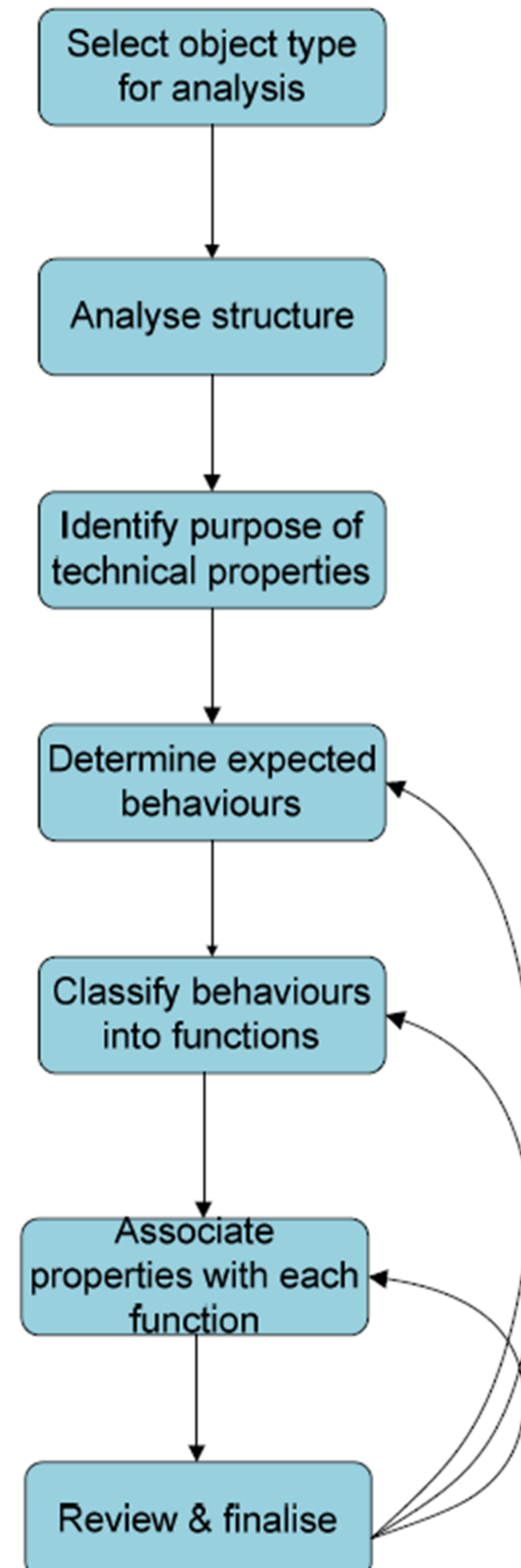National Archives of the Netherlands

## Introduction

- The Open Preservation Foundation's Archives Interest Group reports on our ongoing investigation of significant properties of spreadsheets.
- Because we were faced with the **problem of ensuring long-term accessibility of deposited spreadsheets** while preserving their significant properties
- We therefore wanted to **get hands-on experience** in investigating significant properties to understand the original deposited object, and how to preserve it

## Method

- We compared methods for identifying significant properties, and chose the methodology for assessing significant properties from the **InSpect Framework Report**.
- It is a formalised, open, standards-based model (based on Function-Behaviour-Structure design method, adapted to this new area of work)
- Used in various (digital preservation) projects, resulting in a Significant Properties Testing Report lore (for Digital Audio Recordings, Structured Text, Raster Images, etc.)

## Cooperation

- Monthly calls
- Mailing list
- Google Drive

## (Flowchart)

- Select object type for analysis
- Analyse structure
- Identify purpose of technical properties
- Determine expected behaviours
- Classify behaviours into functions
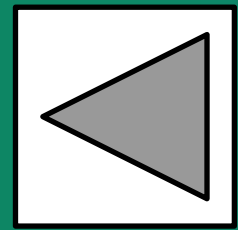- Associate properties with each function
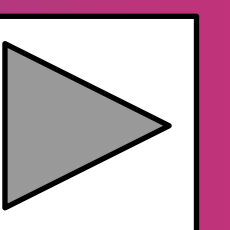- Review & finalise

## Prerequisites

- Collected a **large set of spreadsheets** from public sources and our collections
- Studied spreadsheet **format specification documents** (VisiCalc, Lotus 1-2-3, Quattro Pro, Microsoft Excel versions, ODS versions, OOXML, Numbers, Gnumeric, …)
- Tested **characterisation tools** (FITS, Apache Tika, DROID, Exiftool, File, JHOVE, NLNZ Metadata Extractor, FIDO, Siegfried, Lingfo, Dependency Discovery Tool, Officeparser.py, Python oletools, …)

## Stages

- **Object analysis** (almost done)
  Identify functional, behavioural and structural properties
- **Stakeholder analysis** (to do)
  Analyse functions relevant for a particular user group
- **Reformulation** (where applicable)
  Redevelop object to perform a set of stakeholder functions

1

## Object type

- Spreadsheets, with 2 tentative subtypes:
1. **'Simple/static' spreadsheets** used for (human) visualisation, containing static data values organised into tabular format. Can possibly be migrated to non-spreadsheet file formats.
2. **'Complex/dynamic' spreadsheets** contain formulae, notes, macros, dates, links to external data sources or other functions or behaviour. Migrating to non-spreadsheet file formats would cause severe information loss.

## A complex/dynamic spreadsheet



## Properties



## Expected behaviours



## FBS diagram



Flowchart (center):
- Select object type for analysis
- Analyse structure
- Identify purpose of technical properties
- Determine expected behaviours
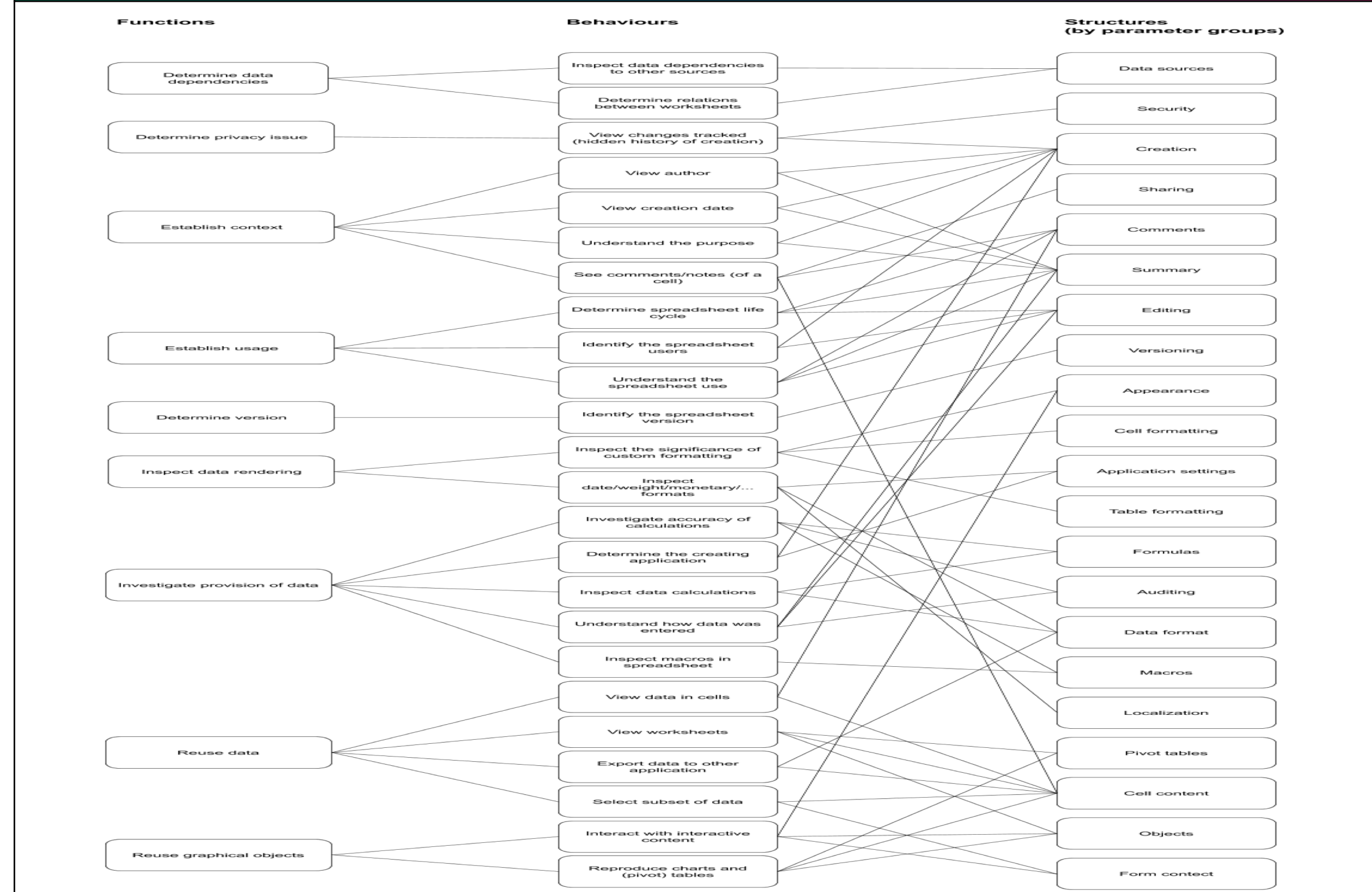- Classify behaviours into functions
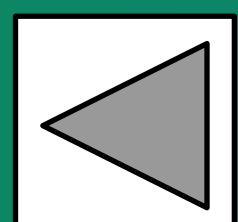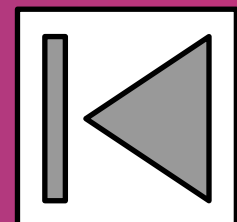- Associate properties with each function
- Review & finalise

2

# Spreadsheet Complexity Analyser and conclusion

## Spreadsheet Complexity Analyser

- Prototype (CC0) **open source tool**
- Extracts **spreadsheet-specific properties**
  **Workbook**: worksheets, fonts, defined names, cell styles, external links and revision history
  **Sheet** (totaled up): formulas, hyperlinks, comments, shapes, dates, cells used
  **VBA**: nonzero indicates possible vba macros
  **File**: file size, creation date/time, last accessed, last modified
- Assesses 'complexity' (using default or user thresholds)

```
F:\AIG_SCA>java -jar SpreadsheetComplexityAnalyser.jar . -v
Spreadsheet complexity analyser results:
File: F:\AIG_SCA\.\test.xls
        size:                   6691.0 kB
        created:                2018-03-01T12:01:15.33Z
        last accessed:          2019-09-05T22:00:00Z
        last modified:          2018-03-02T10:36:40Z
        worksheets:             3
        fonts:                  27
        defined names:          2
        cell styles:            73
        formulas:               1409
        hyperlinks:             2
        comments:               0
        vba macros:             1
        shapes:                 3
        dates:                  0
        cells used:             127404
        external links:         -1
        revision history:       -1
        tentative assessment:   complex/dynamic
Legend:
        -1 = not supported (e.g. external links extraction for XLS).
        0 or more = number of occurrences.
        At macros and revision history, nonzero means they are present.
```

```
F:\AIG_SCA>java -jar SpreadsheetComplexityAnalyser.jar
Error: please provide exactly one input DIRectory.
usage: java -jar SpreadsheetComplexityAnalyser.jar DIR [-c] [-h] [-r] [-v]
                  [-x]
 -c,--config      config file: read complexity assessment threshold values
                  from SpreadsheetComplexityAnalyser.config file
 -h,--help        help: show SpreadsheetComplexityAnalyser help
                  information (and exit)
 -r,--recursive   recurse into subdirectories
 -v,--verbose     verbose output: show number of occurrences of properties
                  in text form
 -x,--xml         xml output: show number of occurrences of properties in
                  xml form (suppresses verbose output)
 DIR              directory with *.xl[st][xm] and *.xl[akms] files to process.
```

## Acknowledgements

- Thank you to Becky McGuinness, Charlotte Armstrong and Carl Wilson of OPF for supporting our work

## Conclusion

- OPF AIG is using the InSpect methodology to investigate the significant properties of spreadsheets. We have almost finished the Object analysis stage.
- A Spreadsheet Testing Report will be added to the InSpect Testing Report lore.
- The Spreadsheet Complexity Analyser tool extracts spreadsheet properties and assesses complexity.

- Preliminary conclusions support earlier findings of significant property studies:
- The complexity of and degree of freedom inherent in spreadsheets makes creating an **exhaustive list of significant spreadsheet properties practically impossible**. But a list of (technical) significant properties does help choose suitable file formats
- To obtain such lists, one needs to **identify relevant stakeholders**, i.e. complete stakeholder analysis

## References

1. Bechhofer, S., Sierman, B., Jones, C., Elstrøm, G., Kulovits, H., Becker, C.: Final version of policy specification model. http://www.scape-project.eu/deliverable/d13-2-catalogue-of-preservation-policy-elements (2014). Accessed March 15, 2019
2. Knight, G.: InSPECT Framework Report. https://web.archive.org/web/20160520083956/http://www.significantproperties.org.uk/inspect-framework.html (2009). Accessed March 15, 2019
3. Lucker, P., Sijtsma, C., van Veenendaal, R.: Significant Significant Properties. https://osf.io/rtjw3 (2018). Accessed March 15, 2019