# Static Spreadsheet Analysis

Partick W. Koch, <u>Birgit Hofer</u>, and Franz Wotawa

**#VALUE! error**

Genomics papers with spreadsheet errors in supplementary files, 2005-15, %
By publishing journal

Bar chart (values read along axis 0–40):

| Journal | % |
| --- | --- |
| Nature | ~31 |
| Genes & Development | ~29 |
| Genome Research | ~27 |
| Genome Biology Evolution | ~25 |
| Nature Genetics | ~24.5 |
| Nucleic Acids Research | ~21 |
| BMC Genomics | ~20 |
| Overall Average | ~20 |
| RNA | ~19.5 |
| PLoS Computational Biology | ~19 |
| PLoS Biology | ~17.5 |
| PLoS One | ~17.5 |
| Human Molecular Genetics | ~16 |
| Science | ~16 |
| BMC Bioinformatics | ~14 |
| Bioinformatics | ~9 |
| Genome Biology Evolution | ~8 |
| DNA Research | ~6.5 |
| Molecular Biology and Evolution | ~5.5 |

Supplementary files with gene-name errors (2005–15), scale 0–200

Source: "Gene name errors are now widespread in the scientific literature", Ziemann, Eren and El-Osta, 2016

The Economist

2

Economist.com

"…the number of genomics papers packaged with error-ridden spreadsheets is **increasing by 15%** a year, far above the 4% annual growth rate in the number of genomics papers published. If we extrapolate current trends …, then **by 2025** every spreadsheet attached to a genetics paper will have an error—unless, of course, there is an error in the spreadsheet we used for this calculation"

The Economist Daily Chart, Sept 7th 2016

Partick W. Koch, Birgit Hofer, and Franz Wotawa
*Static Spreadsheet Analysis*

# Overview

- Running Example
  - Structures that should be detected

- Static Analysis Approach

- Evaluation

- Application
  - Spreadsheet smells

Partick W. Koch, Birgit Hofer, and Franz Wotawa
*Static Spreadsheet Analysis*

# We want to detect …



| | A | B | C | D | E | F |
|---|---|---|---|---|---|---|
| | **Meta-Headers** | Europe | **Headers** | | | |
| 2 | Models | 2012 | 2013 | 2014 | 2015 | Total |
| 3 | Honda | 30 | 27 **Input** | 28 | 32 | =SUM(B3:E3) |
| 4 | Mazda | 10 | 12 **Groups** | 9 | 7 | =SUM(B4:E4) |
| 5 | Fiat | 9 | 12 | 13 | 15 | =SUM(B5:E5) |
| 6 | Total | =SUM(B3:B5) | =SUM(C3:C5) | =SUM(D3:D5) | =SUM(E3:E5) | =SUM(F3:F5) |

**Formula Groups**

Partick W. Koch, Birgit Hofer, and Franz Wotawa
*Static Spreadsheet Analysis*

# Challenges

|   | A | B | C | D | E | F |
|---|---|---|---|---|---|---|
| 1 |  | Europe |  |  |  |  |
| 2 | Models | 2012 | 2013 | 2014 | 2015 | Total |
| 3 | Honda | 30 | 27 | 28 | 32 | =SUM(B3:E3) |
| 4 | Mazda | 10 | 12 | 9 | 7 | =SUM(B4:E4) |
| 5 | Fiat | 9 | 12 | 13 | 15 | =SUM(B5:E5) |
| 6 | Total | =SUM(B3:B5) | =SUM(C3:C5) | =SUM(D3:D5) | =SUM(E3:E5) | =SUM(F3:F5) |

Partick W. Koch, Birgit Hofer, and Franz Wotawa
*Static Spreadsheet Analysis*

# Structural analysis process

Spreadsheet → Grouping → Blocking → Header Assignment

Partick W. Koch, Birgit Hofer, and Franz Wotawa
*Static Spreadsheet Analysis*

# Structural analysis process

Spreadsheet → Grouping → Blocking → Header Assignment

Grouping:
- Type-based
- Partitioned
- Ref.-based

**Type-based**
- Cell type
- Neighbors
- Same formula

**Partitioned**
- 1-dimensional formula groups
- Few groups

**Ref.-based**
- Cells referenced by partitioned group
- Merge overlapping groups with same orientation

# Structural analysis process



| | A | B | C | D | E | F |
|---|---|---|---|---|---|---|
| 1 | | Europe | | | | |
| 2 | Models | 2012 | 2013 | 2014 | 2015 | Total |
| 3 | Honda | 30 | 27 | 28 | 32 | =SUM(B3:E3) |
| 4 | Mazda | 10 | 12 | 9 | 7 | =SUM(B4:E4) |
| 5 | Fiat | 9 | 12 | 13 | 15 | =SUM(B5:E5) |
| 6 | Total | =SUM(B3:B5) | =SUM(C3:C5) | =SUM(D3:D5) | =SUM(E3:E5) | =SUM(F3:F5) |

Partick W. Koch, Birgit Hofer, and Franz Wotawa
*Static Spreadsheet Analysis*

# Structural analysis process

|   | A | B | C | D | E | F |
|---|---|---|---|---|---|---|
| 1 |  | Europe |  |  |  |  |
| 2 | Models | 2012 | 2013 | 2014 | 2015 | Total |
| 3 | Honda | 30 | 27 | 28 | 32 | =SUM(B3:E3) |
| 4 | Mazda | 10 | 12 | 9 | 7 | =SUM(B4:E4) |
| 5 | Fiat | 9 | 12 | 13 | 15 | =SUM(B5:E5) |
| 6 | Total | =SUM(B3:B5) | =SUM(C3:C5) | =SUM(D3:D5) | =SUM(E3:E5) | =SUM(F3:F5) |

Partick W. Koch, Birgit Hofer, and Franz Wotawa
*Static Spreadsheet Analysis*

# Structural analysis process

Partick W. Koch, Birgit Hofer, and Franz Wotawa
*Static Spreadsheet Analysis*

# Structural analysis process

Spreadsheet → Grouping → **Blocking** → Header Assignment
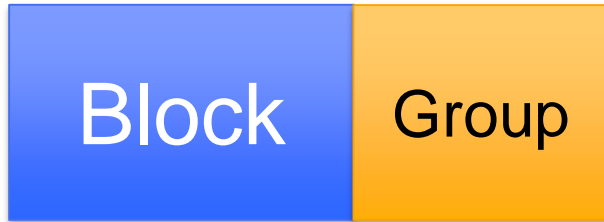
## What is a block?

- Rectangular area
- Aggregate input and formula cells, NOT header cells

## How to create a block?

1. Pick formula-group or referenced group
2. Expand

# Structural Analysis Process

Block | Group → Block

Block | Group → Block

Block / Group → Block

Partick W. Koch, Birgit Hofer, and Franz Wotawa
*Static Spreadsheet Analysis*

# Structural analysis process

Partick W. Koch, Birgit Hofer, and Franz Wotawa
*Static Spreadsheet Analysis*

# Structural analysis process

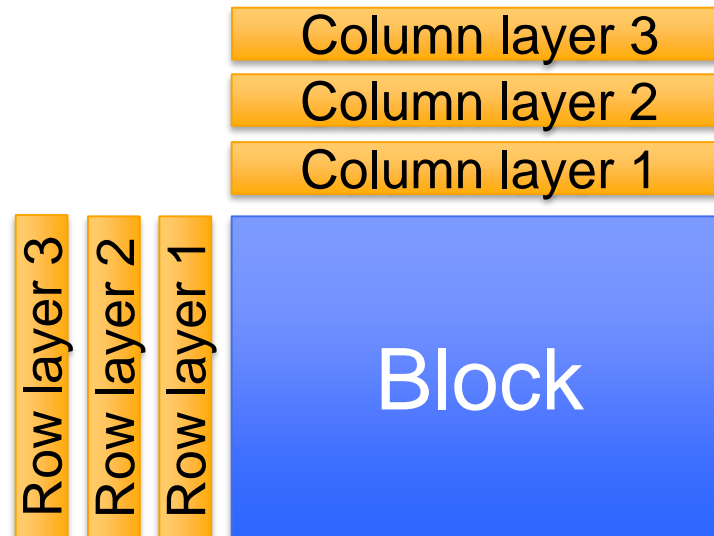| Spreadsheet | Grouping | Blocking | **Header Assignment** |

**Position of headers**

- Above blocks
- LTR vs. RTL systems
- Column vs. row headers
- Headers of headers

**Identification**

- Identify header layers

Column layer 3
Column layer 2
Column layer 1

Row layer 3  Row layer 2  Row layer 1

Block

# Structural analysis process

|  | A | B | C | D | E | F |
|---|---|---|---|---|---|---|
| 1 |  | Europe | | | | |
| 2 | Models | 2012 | 2013 | 2014 | 2015 | Total |
| 3 | Honda | 30 | 27 | 28 | 32 | =SUM(B3:E3) |
| 4 | Mazda | 10 | 12 | 9 | 7 | =SUM(B4:E4) |
| 5 | Fiat | 9 | 12 | 13 | 15 | =SUM(B5:E5) |
| 6 | Total | =SUM(B3:B5) | =SUM(C3:C5) | =SUM(D3:D5) | =SUM(E3:E5) | =SUM(F3:F5) |

## Next step

- Identify remaining meta-header

## Future work

- Better mapping approach for meta headers

Partick W. Koch, Birgit Hofer, and Franz Wotawa
*Static Spreadsheet Analysis*

# Evaluation



|  | Corpus | | POI readable | | Application readable | | Fit for evaluation |
|---|---|---|---|---|---|---|---|
| EUSES | 4,495 | | 4,172 | | 4,017 | | 1,659 |
| ENRON | 15,929 | | 15,929 | | 13,402 | | 6,691 |

Partick W. Koch, Birgit Hofer, and Franz Wotawa
*Static Spreadsheet Analysis*

# Evaluation Results



Legend:
- detected
- complete
- correct

EUSES Blocks, EUSES Column Headers, EUSES Row Headers, ENRON Blocks, ENRON Column Headers, ENRON Row Headers

Values: 99, 89, 69, 96, 89, 77, 84, 94, 76, 55, 83, 99, 82, 86, 60, 76, 72, 63

Partick W. Koch, Birgit Hofer, and Franz Wotawa
*Static Spreadsheet Analysis*

# Open Problems

| | A | B | C | D | E | F | G | H | I | J |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | | | | | | | | | | |
| 2 | | | | | | | **HANSON PERMANENTE CEMENT** | | | |
| 3 | | | | | | | Power Usage Forecast | | | |
| 4 | | CALIFORNIA SCHEDULING COORDINATION AND | | | | | | | | |
| 5 | | ENERGY PURCHASES AND SALES AGREEMENT | | | | | | | | |
| 6 | | | | | | | | | | |
| 7 | | Day of Scheduling | | | | | | | | |
| 8 | | Tuesday: Wednesday Scheduling ▼ | | | | | | | | |
| 9 | | | | | | | | | | |
| 10 | | 06/26/02 | HE | HE | HE | HE | HE | HE | HE | HE |
| 11 | | | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| 12 | | Long Term Purchase Quantity for Wednesday | 16.368 | 20.184 | 24.000 | 24.360 | 25.992 | 26.376 | 26.484 | 26.568 |
| 13 | | | | | | | | | | |
| 14 | | Expected Usage for Wednesday | 26.510 | 26.510 | 26.510 | 26.510 | 27.260 | 27.260 | 27.460 | 27.460 |
| 15 | | | | | | | | | | |
| 16 | | Preschedule Quantity for Wednesday | 26.368 | 26.184 | 27.000 | 26.360 | 26.992 | 27.376 | 27.484 | 27.568 |
| 17 | | | | | | | | | | |
| 18 | | Total Incremental Quantity (Hanson needs to purchase) for Wednesday | 10 | 6 | 3 | 2 | 1 | 1 | 1 | 1 |
| 19 | | | | | | | | | | |
| 20 | | Total Decremental Quantity (Hanson needs to sell) for Wednesday | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 21 | | | | | | | | | | |
| 22 | | Day Ahead Incremental Quantity (Hanson needs to purchase) | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 23 | | Day Ahead Decremental Quantity (Hanson needs to sell) | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 24 | | Real Time Incremental Quantity (Hanson needs to purchase) for Wednesday | 10 | 6 | 3 | 2 | 1 | 1 | 1 | 1 |
| 25 | | | | | | | | | | |
| 26 | | Real Time Decremental Quantity (Hanson needs to sell) for Wednesday | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 27 | | | | | | | | | | |

Partick W. Koch, Birgit Hofer, and Franz Wotawa
*Static Spreadsheet Analysis*

# Open Problems

Partick W. Koch, Birgit Hofer, and Franz Wotawa
*Static Spreadsheet Analysis*

# Application: Spreadsheet Smells

- Origin: **Code Smells**

- Current state of spreadsheet smells

    - Too many smells reported

    - Analysis time

→ Use structure Information

- To improve existing smells
- To identify new smells

Partick W. Koch, Birgit Hofer, and Franz Wotawa
*Static Spreadsheet Analysis*

# Improve existing smells

| Sliding Window Smells | Similarity-Based Smells | Formula-Based Smells | Long Calculation Chain | Inter-Worksheet Smells |
|---|---|---|---|---|

## Update opportunities

- Focus analysis methods on spreadsheet structures

- Analyse group formulas instead of cell fromulas

- Analyse group references instead of cell references

Partick W. Koch, Birgit Hofer, and Franz Wotawa
*Static Spreadsheet Analysis*

Motivation  Structural Analysis Process  Evaluation  Structure-based Smells  Future Work

# Example
## Sliding Window Smell

- Detects anomalies in sliding windows

- Update: limit windows to structures

|   | A | B | C |
|---|---|---|---|
| 1 |   | Europe |   |
| 2 | Models | 2012 | 2013 |
| 3 | Honda | 30 | 27 |
| 4 | Mazda | 1000 | 12 |
| 5 | Fiat | 9 | 12 |
| 6 | Total | =SUM(B3:B5) | =SUM(C3:C5) |

Partick W. Koch, Birgit Hofer, and Franz Wotawa
*Static Spreadsheet Analysis*

Motivation · Structural Analysis Process · Evaluation · Structure-based Smells · Future Work

# Structure-based Smells: Novel Smells

Duplicated Formula Groups

Formula Group Distance

Unrelated Neighbours

Inconsistent Reference Dimensions

Inconsistent Formula Group Reference

Missing Header

Overburdened Worksheet

Partick W. Koch, Birgit Hofer, and Franz Wotawa
*Static Spreadsheet Analysis*

Motivation — Structural Analysis Process — Evaluation — Structure-based Smells — Future Work

# Example
# Inconsistent Formula Group Reference

- Size mismatch between groups and group references

| | A | B | C | D | E | F |
|---|---|---|---|---|---|---|
| 1 | | Europe | | | | |
| 2 | Models | 2012 | 2013 | 2014 | 2015 | Total |
| 3 | Honda | 30 | 27 | 28 | 32 | =SUM(B3:E3) |
| 4 | Mazda | 1000 | 12 | 9 | 7 | =SUM(B4:E4) |
| 5 | Fiat | 9 | 12 | 13 | 15 | =SUM(B5:E5) |
| 6 | Total | =SUM(B3:B4) | =SUM(C3:C4) | =SUM(D3:D4) | =SUM(E3:E4) | =SUM(F3:F4) |

Partick W. Koch, Birgit Hofer, and Franz Wotawa
*Static Spreadsheet Analysis*

Motivation    Structural Analysis Process    Evaluation    Structure-based Smells    Future Work
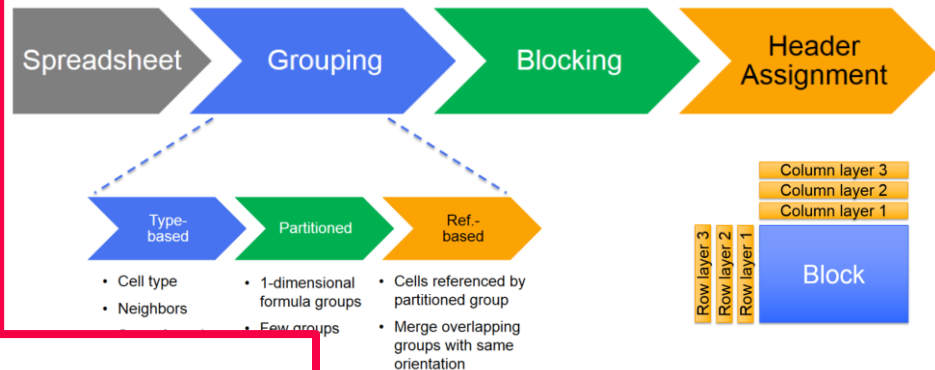
# Summary

## #VALUE! error
Genomics papers with spreadsheet errors in supplementary files, 2005-15, %
By publishing journal

*The Economist*

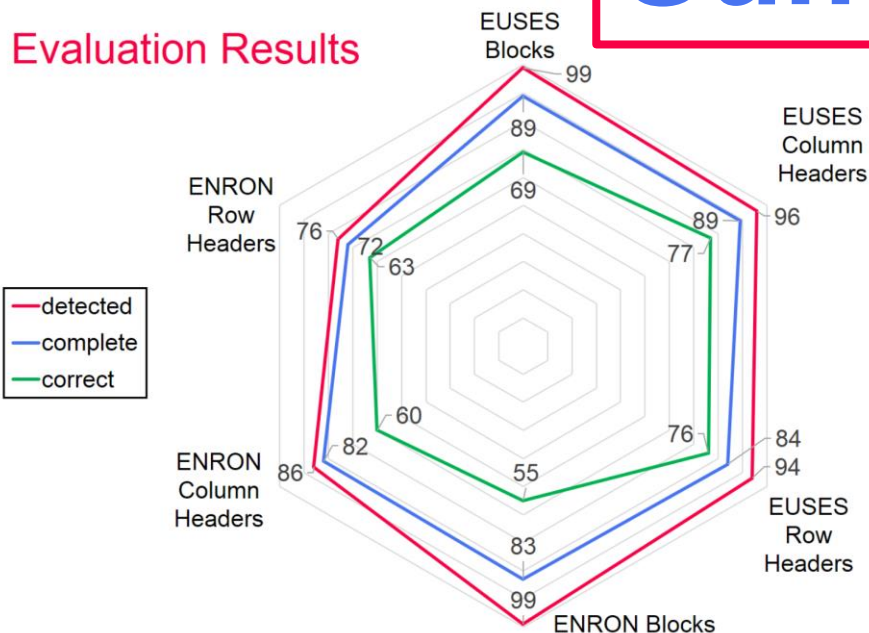## Structural analysis process

## Evaluation Results

- detected
- complete
- correct

## Application: Spreadsheet Smells

- Orign: **Code Smells**

- Current state of spreadsheet smells

  - Too many smells reported

  - Analysis time

Partick W. Koch, Birgit Hofer, and Franz Wotawa
*Static Spreadsheet Analysis*