

# Mapping Spreadsheets to RDF – Supporting Excel in RML

Markus Schröder, Christian Jilek, Andreas Dengel



Second International Workshop on  
Knowledge Graph Construction

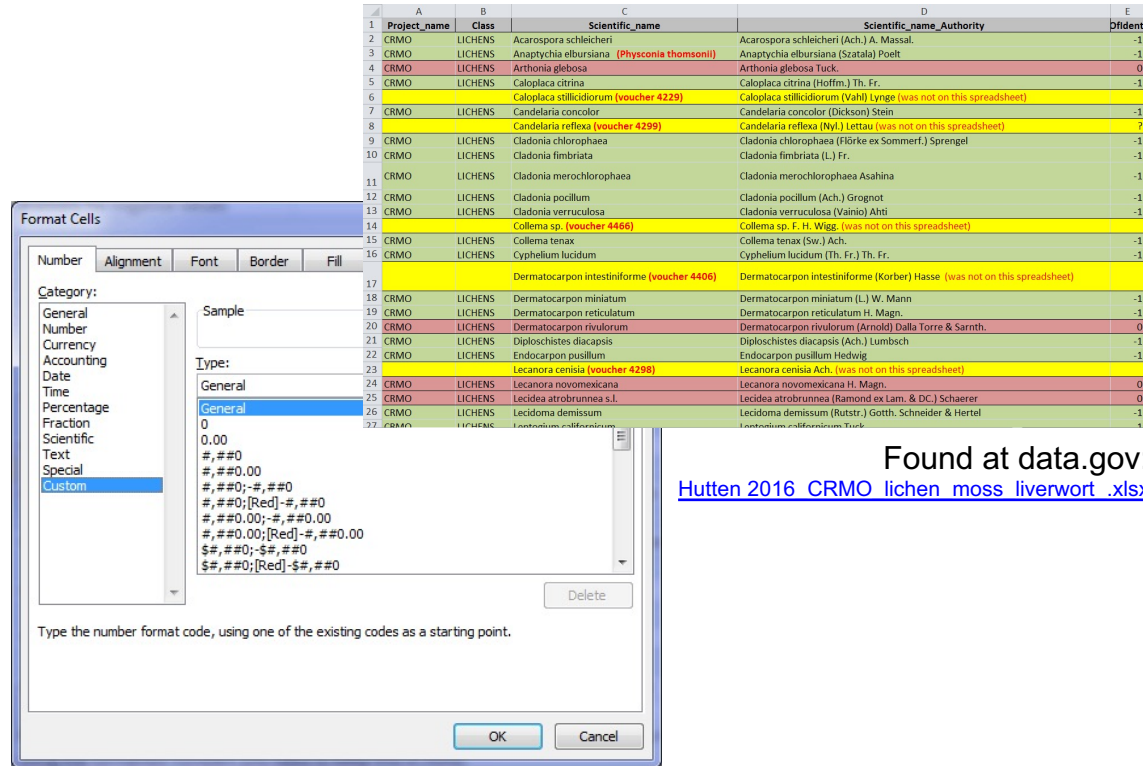
06.06.2021

# A Spreadsheet

(or sometimes a canvas painted by a data artist)

## ◆ Spreadsheets

- well understood
- easy and fast possibility to enter data
- complex workbooks
- multiple sheets
- cells having rich meta data
  - formats, colors, fonts, styles, borders, etc.
  - arbitrarily arranged
  - can lead to inconsistent and unstructured content



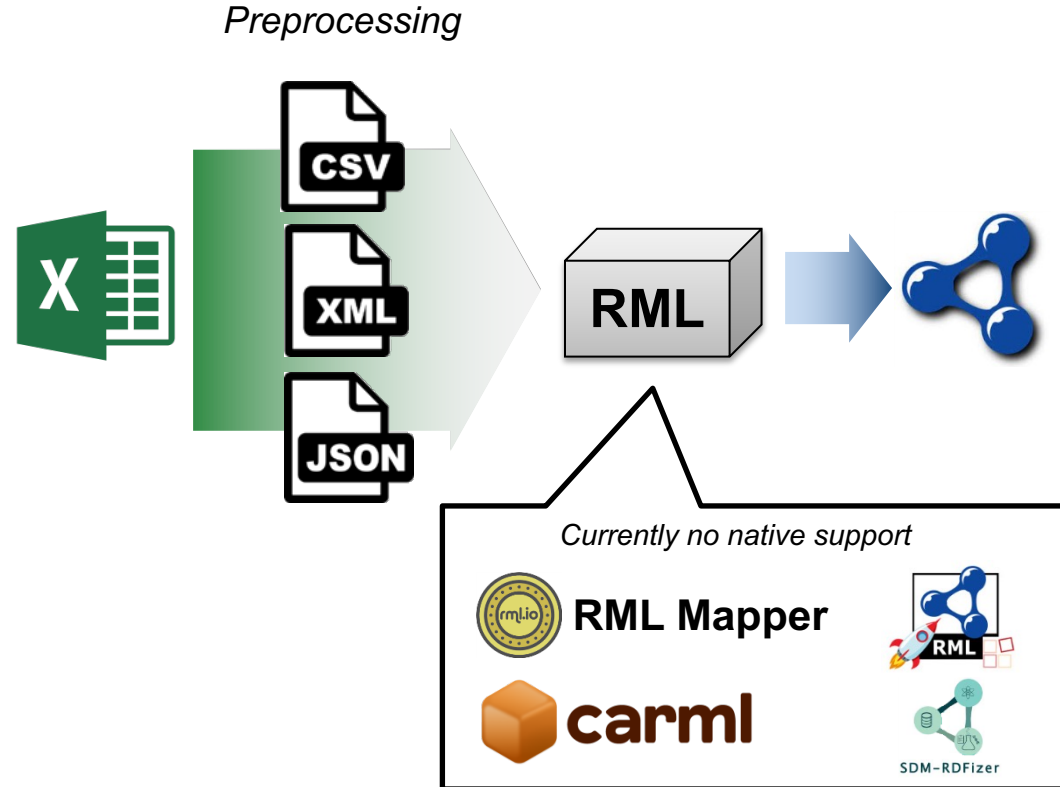
The image shows a screenshot of an Excel spreadsheet with a table of lichen data. The table has columns for Project\_name, Class, Scientific\_name, Scientific\_name\_Authority, and Dfidnet. The data includes various lichen species like Acarospora schleicheri, Anaptychia elburiana, Arthonia glebosa, Caloplaca citrina, Caloplaca sticticidiorum, Candelaria concolor, Candelaria reflexa, Cladonia chlorophaea, Cladonia fimbriata, Cladonia merchlorophaea, Cladonia pocillum, Cladonia verruculosa, Collema sp., Collema tenax, Cyphellum lucidum, Dermatocarpon intestinforme, Dermatocarpon miniatum, Dermatocarpon reticulatum, Dermatocarpon rivulorum, Diploschistes diacapsis, Endocarpon pusillum, Lecanora censis, Lecanora novomexicana, Lecidea atrobrunnea, and Lecidoma demissum. Some cells are highlighted in yellow, and some are in red. A 'Format Cells' dialog box is open over the spreadsheet, showing the 'Number' tab with the 'Custom' category selected. The 'Sample' text box is empty, and the 'Type' list shows various number formats. The 'OK' and 'Cancel' buttons are at the bottom of the dialog box.

Project_name	Class	Scientific_name	Scientific_name_Authority	Dfidnet
CRMO	LICHENS	Acarospora schleicheri	Acarospora schleicheri (Ach.) A. Massal.	-1
CRMO	LICHENS	Anaptychia elburiana (Physconia thomsonii)	Anaptychia elburiana (Szatala) Poelt	-1
CRMO	LICHENS	Arthonia glebosa	Arthonia glebosa Tuck.	0
CRMO	LICHENS	Caloplaca citrina	Caloplaca citrina (Hoffm.) Th. Fr.	-1
		Caloplaca sticticidiorum (voucher 4229)	Caloplaca sticticidiorum (Vahl) Lynge (was not on this spreadsheet)	
		Candelaria concolor	Candelaria concolor (Dickson) Stein	-1
		Candelaria reflexa (voucher 4299)	Candelaria reflexa (Nyl.) Lettau (was not on this spreadsheet)	?
CRMO	LICHENS	Cladonia chlorophaea	Cladonia chlorophaea (Förke ex Sommerf.) Sprengel	-1
CRMO	LICHENS	Cladonia fimbriata	Cladonia fimbriata (L.) Fr.	-1
CRMO	LICHENS	Cladonia merchlorophaea	Cladonia merchlorophaea Asahina	-1
		Cladonia pocillum	Cladonia pocillum (Ach.) Grognot	-1
		Cladonia verruculosa	Cladonia verruculosa (Vainio) Abti	-1
		Collema sp. (voucher 4466)	Collema sp. F. H. Wigg. (was not on this spreadsheet)	
		Collema tenax	Collema tenax (Sw.) Ach.	-1
		Cyphellum lucidum	Cyphellum lucidum (Th. Fr.) Th. Fr.	-1
		Dermatocarpon intestinforme (voucher 4406)	Dermatocarpon intestinforme (Korber) Hasse (was not on this spreadsheet)	
CRMO	LICHENS	Dermatocarpon miniatum	Dermatocarpon miniatum (L.) W. Mann	-1
CRMO	LICHENS	Dermatocarpon reticulatum	Dermatocarpon reticulatum H. Magn.	-1
		Dermatocarpon rivulorum	Dermatocarpon rivulorum (Arnold) Dalla Torre & Sarnth.	0
CRMO	LICHENS	Diploschistes diacapsis	Diploschistes diacapsis (Ach.) Lumbsch	-1
CRMO	LICHENS	Endocarpon pusillum	Endocarpon pusillum Hedwig	-1
		Lecanora censis (voucher 4298)	Lecanora censis Ach. (was not on this spreadsheet)	
		Lecanora novomexicana	Lecanora novomexicana H. Magn.	0
CRMO	LICHENS	Lecidea atrobrunnea s.l.	Lecidea atrobrunnea (Ramond ex Lam. & DC.) Schaerer	0
CRMO	LICHENS	Lecidoma demissum	Lecidoma demissum (Rutstr.) Gotth. Schneider & Hertel	-1
		Lecidoma demissum	Lecidoma demissum Tuck.	

Found at data.gov:  
[Hutten 2016 CRMO lichen moss liverwort .xlsx](#)

# From Spreadsheet to RDF

- ◆ Use input format that is supported
- ◆ Advantages of native support
  - Eliminates extraneous conversion efforts:  
No preprocessing and transformation needed
  - All aspects of a spreadsheet can be exploited
  - Eases mapping rule communication
  - For RML practitioners no extra language to learn



# Related Work

## Spread2RDF (Ruby syntax)

```
worksheet 'MaterialelementeKlassen',
  name: :MaterialElementClasses,
  start: :B5,
  subject: { uri: { namespace: PSM.MaterialElement },
    type: RDF::RDFS.Class,
    sub_class_of: PSM.MaterialElement
  } do
  column :name, predicate: RDFS.label
  column :uri

  column :sub_class_of, predicate: RDFS.subClassOf,
    object: { from: :MaterialElementClasses }
  column_block :parameter, subject: { uri: :bnode, type: PSM.Parameter },
    predicate: PSM.materialParameter,
    statement: :restriction do
    column :name, predicate: PSM.parameterName
    column :description, predicate: PSM.parameterDescription

    column :min, predicate: PSM.parameterMinQuantity,
      object: { uri: :bnode, type: QUOT.QuantityValue },
      &quantity_mapping
  end
end
```

## XLWrap (TriG syntax)

```
{ [ ] a xl:Mapping ;
  ...
  xl:templateGraph :ANamedGraph
  ...
}

:ANamedGraph {
  [ xl:uri "'http://example.org/' & URLENCODE(A2 & B2)"^xl:Expr ] a foaf:Person ;
  foaf:name "A2 & ' ' & B2"^xl:Expr ;
  foaf:mbox_sha1sum "SHA(C2)"^xl:Expr ;
}
```

## Mapping Master (Manchester syntax)

```
Class: @A*(rdfs:label 'analyte assay')
EquivalentTo:
(achieves_planned_objective some 'analyte measurement objective') and
(realizes some ('evaluant role' and (role_of some
  @D*(material_entity)))) and
(realizes some ('analyte role' and
  (role_of some ('scattered molecular aggregate' and
    ('has grain' only
      @B*('molecular entity'))))))

SubClassOf:
has specified output some
  ('scalar measurement datum' and
    ('is quality measurement of' some 'molecular concentration')) and
  ('has measurement unit label' some
    @F*('measurement unit label'))
```

## Sheet2RDF (PEARL syntax)

```
prefix : <http://baseuri.org#>
prefix coda: <http://art.uniroma2.it/coda/contracts>
...

rule it.uniroma2.art.Sheet2RDFAnnotation id:row {
  nodes = {
    subject uri col0_skos_Concept/value
    col1_skos_broadener_value uri col1_skos_broadener/value
    col2_skos_prefLabel_value literal@en col2_skos_prefLabel/value
  }
  graph = {
    $subject rdf:type <http://www.w3.org/2004/02/skos/core#Concept> .
    $subject skos:inScheme <http://sheet2rdf#main> .
    OPTIONAL { $subject skos:broadener $col1_skos_broadener_value . }
    OPTIONAL { $subject skos:prefLabel $col2_skos_prefLabel_value . }
  }
}
```

◆ Various languages with different features

◆ RML + Excel

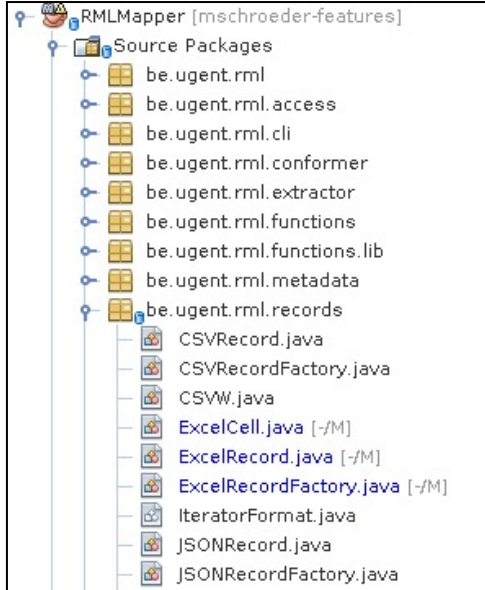
- Familiar syntax and concepts
- Utilize RML features

# Approach: Technical Integration in RMLMapper

## ◆ Extension to the RML Mapper tool



[GitHub fork](#)



```
/** This class creates records based on RML rules ...3 lines */
public class RecordsFactory {

    private Map<Access, Map<String, Map<String, List<Record>>>> recordCache;
    private AccessFactory accessFactory;
    private Map<String, ReferenceFormulationRecordFactory> referenceFormulationRecordFactoryMap;

    public RecordsFactory(String basePath) {
        accessFactory = new AccessFactory(basePath);
        recordCache = new HashMap<>();

        referenceFormulationRecordFactoryMap = new HashMap<>();
        referenceFormulationRecordFactoryMap.put(NAMESPACES.QL + "XPath", new XMLRecordFactory());
        referenceFormulationRecordFactoryMap.put(NAMESPACES.QL + "JSONPath", new JSONRecordFactory());
        referenceFormulationRecordFactoryMap.put(NAMESPACES.QL + "CSV", new CSVRecordFactory());
        referenceFormulationRecordFactoryMap.put(NAMESPACES.QL + "Spreadsheet", new ExcelRecordFactory());
    }
}
```

```
:ls1
a rml:LogicalSource ;
rml:referenceFormulation ql:Spreadsheet ;
rml:source [
    a ss:Workbook;
    ss:url "workbook.xlsx" ;
    ss:sheetName "Papers" ;
    ss:range "A2:A5" ;
] .
```

## ◆ Record Factory

- Use parameters to return records



## ◆ Record

- Interpret references

# Conceptual Integration in RML

## Cell Location Reference

- ◆ Cell iteration using ranges (e.g. B2:B5)
  - Instead of row iteration (like in CSV)
- ◆ Arbitrarily structured tables without any anchor points like column names
  - Relative reference to neighboring cells with parenthesis notation
  - Absolute reference to cells with square brackets

```
rr:predicateObjectMap [
  a rr:PredicateObjectMap ;
  rr:predicateMap [
    a rr:PredicateMap ;
    rr:constant ex:numberOfPages
  ] ;
  rr:objectMap [
    a rr:ObjectMap ;
    rml:reference "(1,0).valueInt" ;
    rr:datatype xsd:integer
  ]
] .
```

x-offset: -1    0    +1

y-offset: -1    0    +1

	A	B	C
1	Title	Pages	Price
2	Deducing Suppres	4	\$10,40
3	Structuring Fair-M	6	\$5,23
4	Utilizing Symbolic	8	\$2,50
5	Colonizing Loyal	3	\$11,23

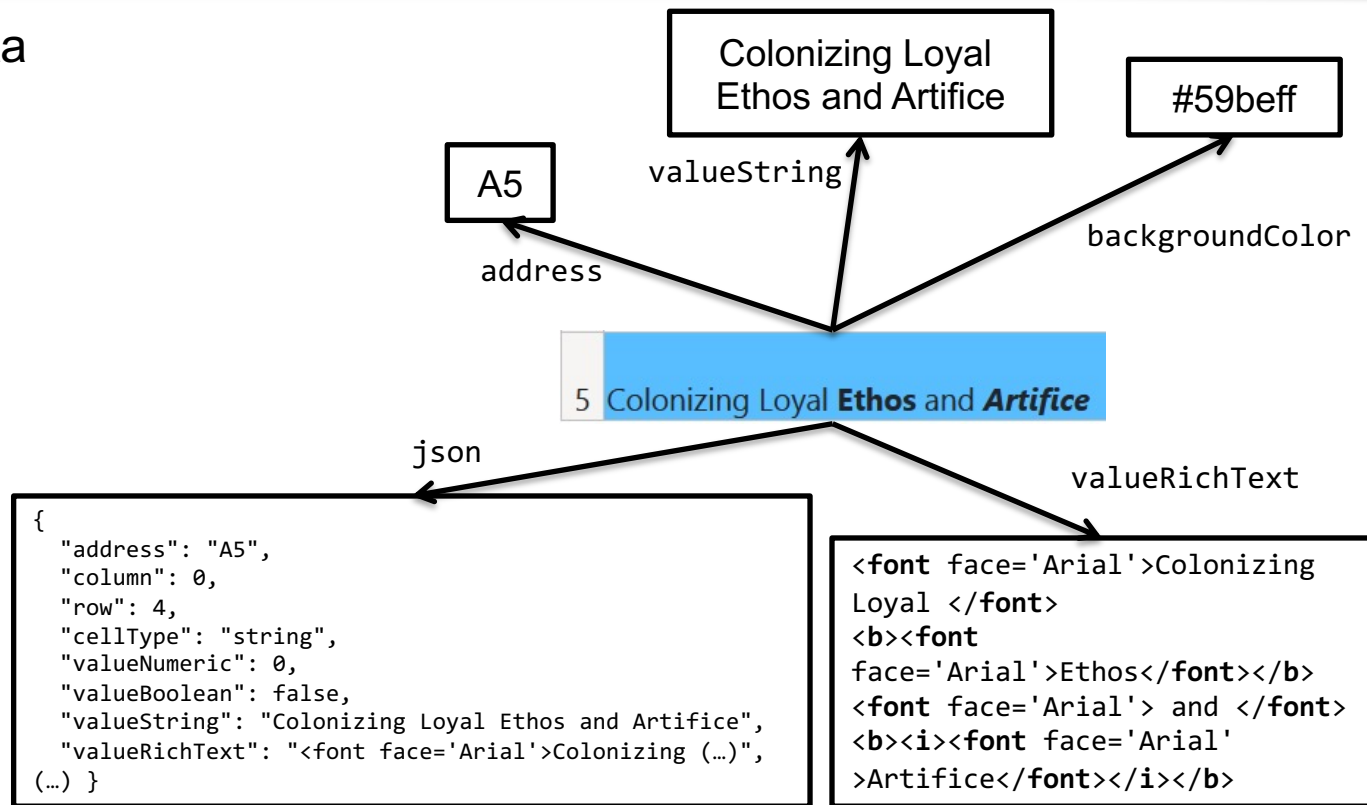
(1,0)

# Conceptual Integration in RML

## Meta Data References

### ◆ Access cell meta data

- address
- column
- row
- backgroundColor
- foregroundColor
- fontColor
- fontName
- fontSize
- valueNumeric
- valueInt
- valueBoolean
- valueFormula
- valueError
- valueString
- valueRichText
- json
- value





# Experimental Features

## Multiple Different Properties in a Cell

- ◆ Each piece of information corresponds to a different property
- ◆ Zip together
  - properties
  - returned objects from FnO function
- ◆ Just a shortcut
  - Usually requires separate predicate object maps

```
rr:predicateObjectMap [  
  a rr:PredicateObjectMap ;  
  ss:zip true ;  
  rr:predicateMap [  
    a rr:PredicateMap ;  
    rr:constant ( ex:numberOfPages ex:price )  
  ] ;  
  rr:objectMap [  
    a rr:ObjectMap , fnml:FunctionMap ;  
    fnml:functionValue [  
      rr:predicateObjectMap [  
        rr:predicate fno:executes ;  
        rr:object <java:ifRegexReturnGroup>  
      ] ;  
      # String value  
      rr:predicateObjectMap [  
        rr:predicate <java:parameter.predicate.string.0> ;  
        rr:objectMap [ rml:reference "(5,0).valueString" ]  
      ] ;  
      # String pattern  
      rr:predicateObjectMap [  
        rr:predicate <java:parameter.predicate.string.1> ;  
        rr:objectMap [ rr:constant "\\d+\\.?\\d*" ]  
      ] ;  
    ]  
  ] ;
```

F	
Pages + Price	
4, \$10.4	
6 \$5.23	
8 – 2.50\$	
3, \$11.23	

page information

price information

```
<http://example.org/A2> a <http://example.org/Paper>;  
<http://example.org/numberOfPages> 4.0;  
<http://example.org/price> 10.4 .  
  
<http://example.org/A3> a <http://example.org/Paper>;  
<http://example.org/numberOfPages> 6.0;  
<http://example.org/price> 5.23 .  
  
<http://example.org/A4> a <http://example.org/Paper>;  
<http://example.org/numberOfPages> 8.0;  
<http://example.org/price> 2.5 .  
  
<http://example.org/A5> a <http://example.org/Paper>;  
<http://example.org/numberOfPages> 3.0;  
<http://example.org/price> 11.23 .
```

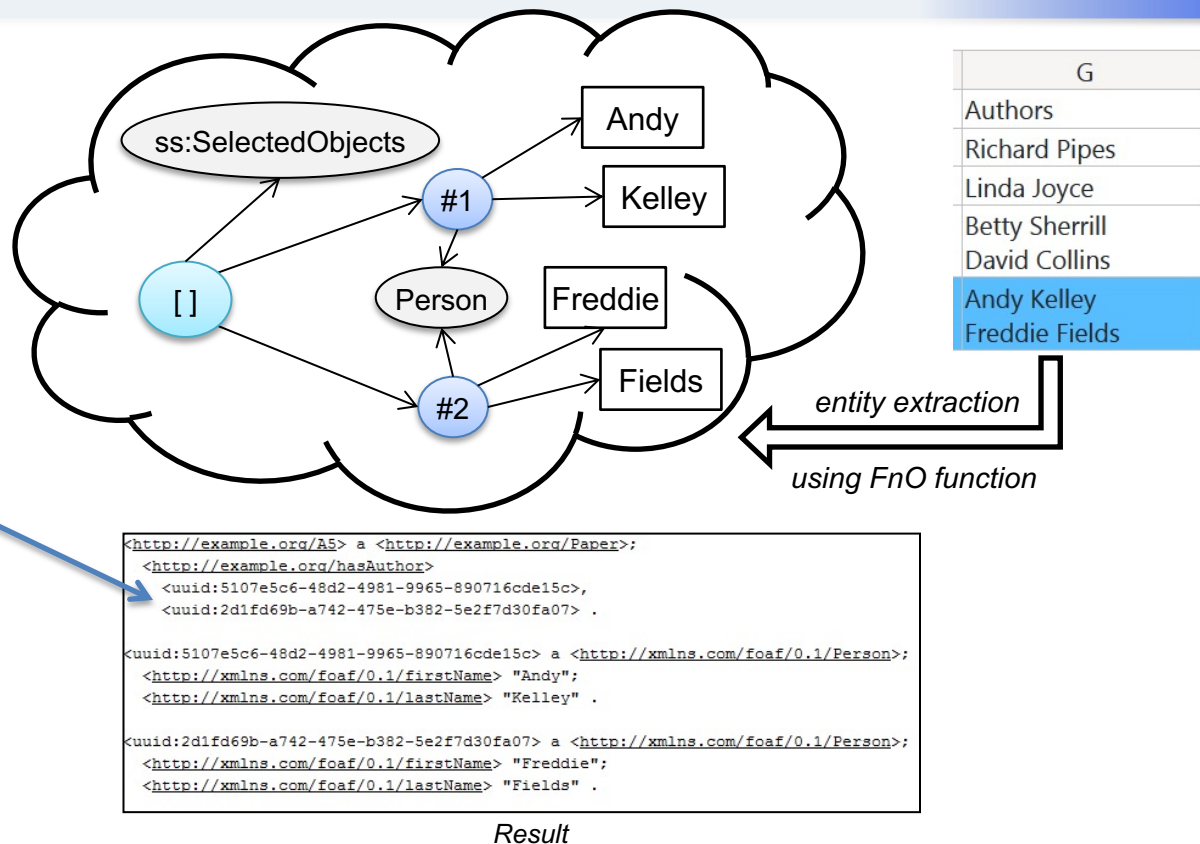
Result



# Experimental Features

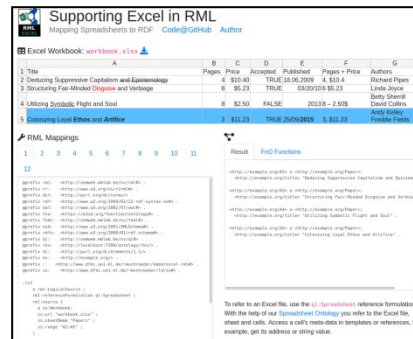
## Multiple Complex Entities in a Cell

- ◆ For example: list of persons having first and last names
- ◆ RDF graph in turtle syntax
- ◆ New term type `ss:Graph`
- ◆ Graph is added to result
- ◆ Selected ones are mapped using `ss:SelectedObjects`



# Conclusion

- ◆ Implemented Excel support in RML Mapper
- ◆ Cell iteration, location, meta data access
- ◆ Try it out on our [demo page](#)
  - [GitHub code](#)



- ◆ You may also be interested in
  - [Data Sprout](#) – Dataset generation for evaluating KGC
- ◆ More related research in our project
  - [SensAI](#)



Thank you for your attention.  
Questions?