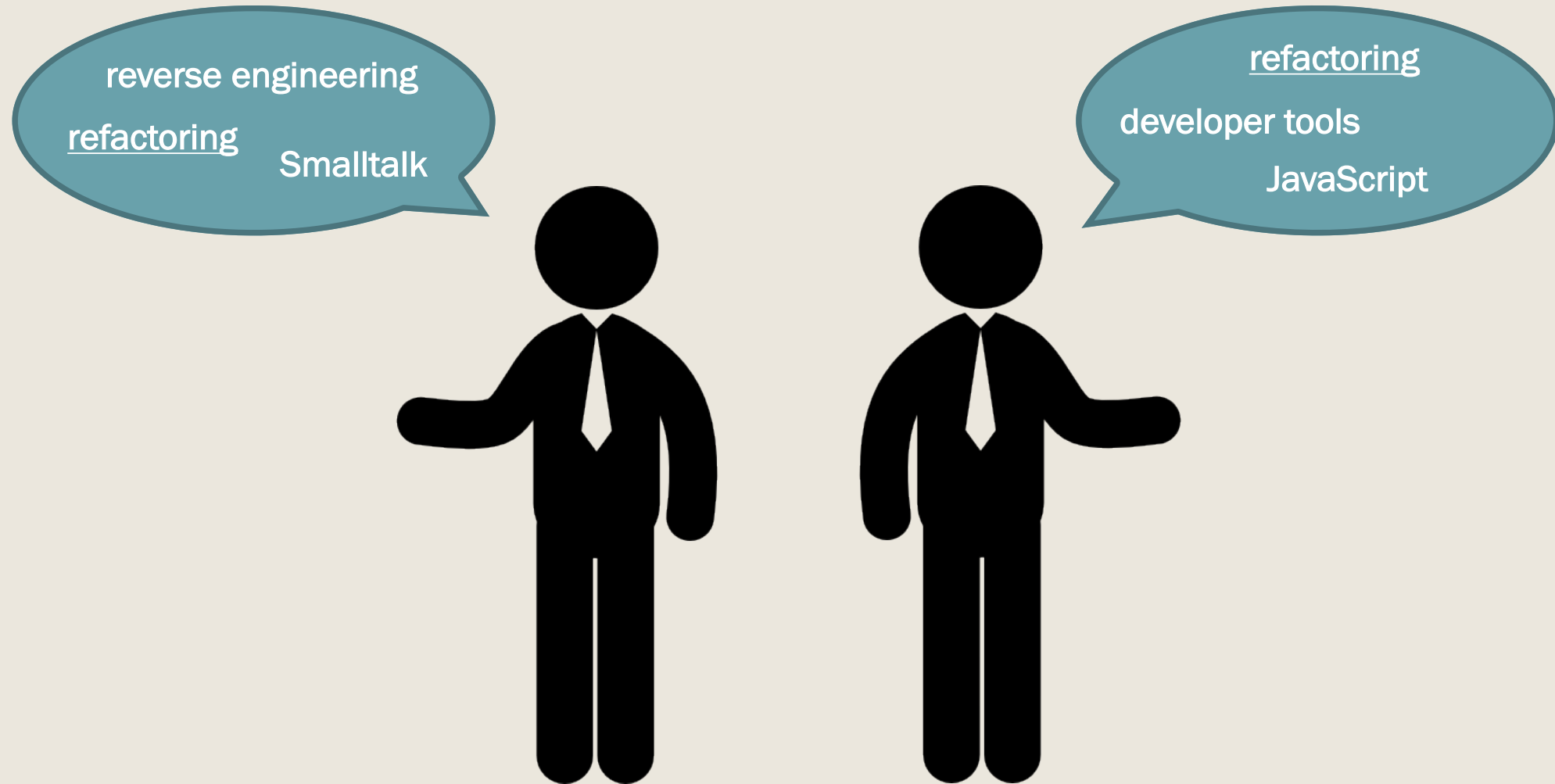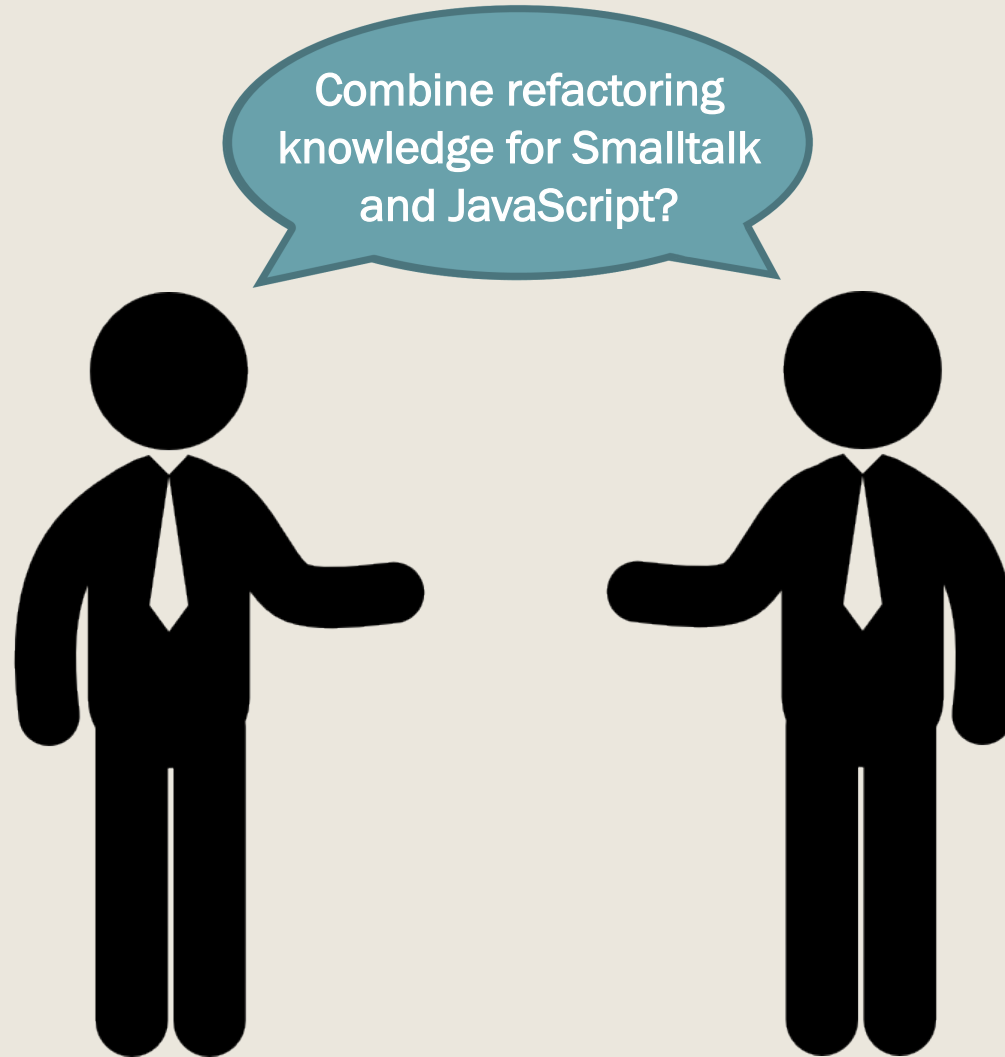# VISUALLY MINING SCIENTIFIC COMMUNITIES

Bachelor Thesis by Silas Berger (silas.berger@students.unibe.ch), supervised by Leonel Merino
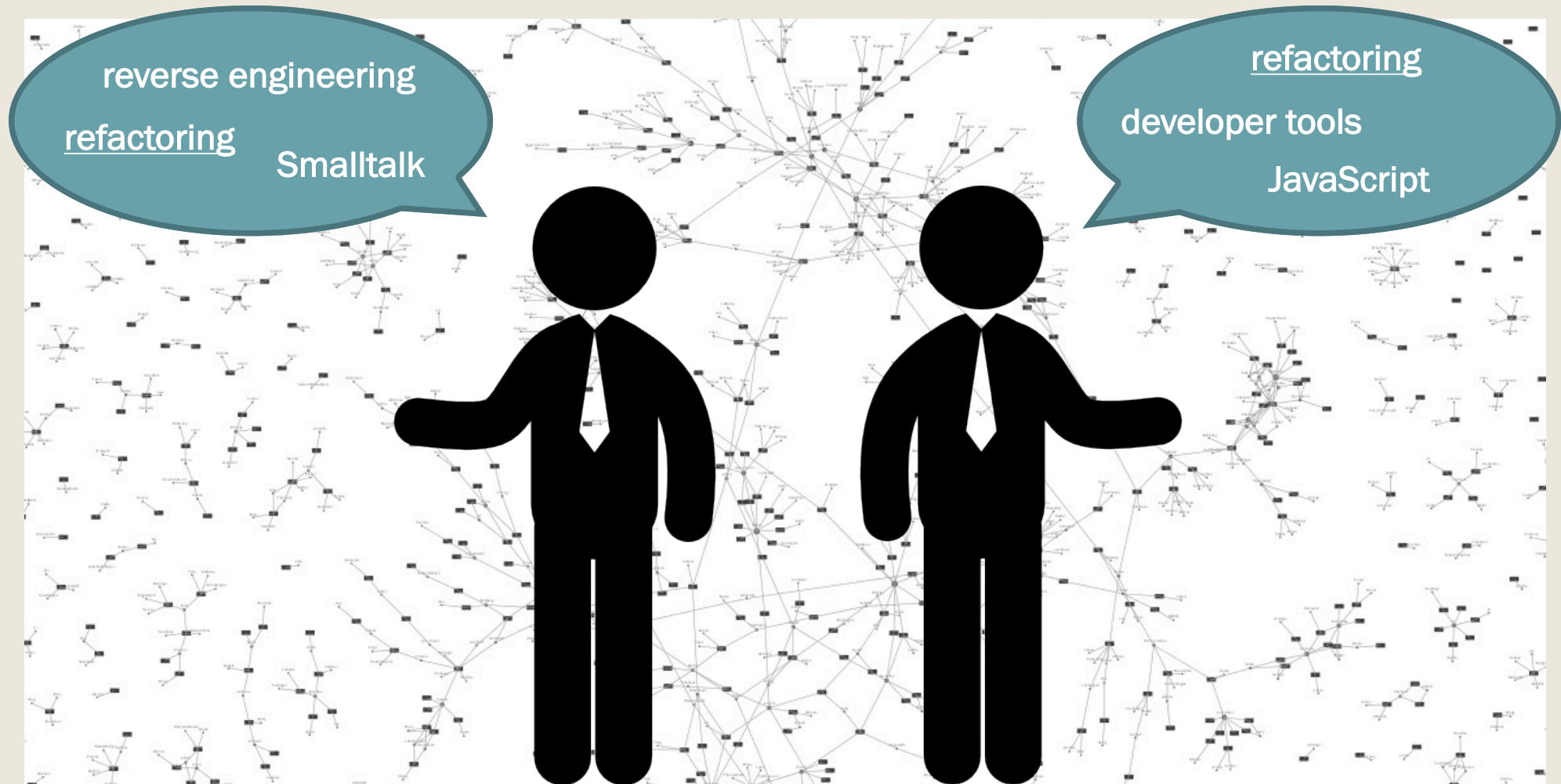
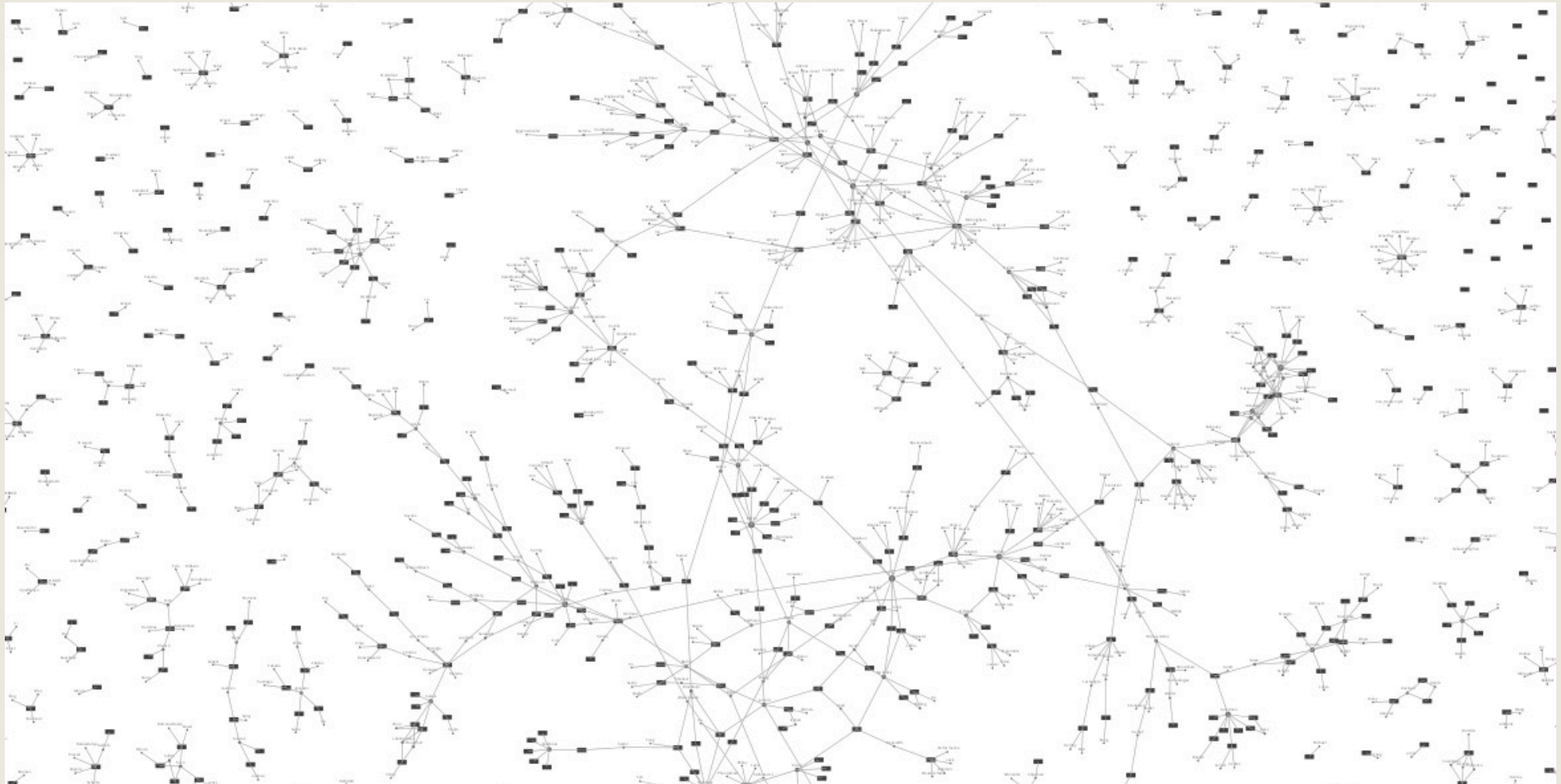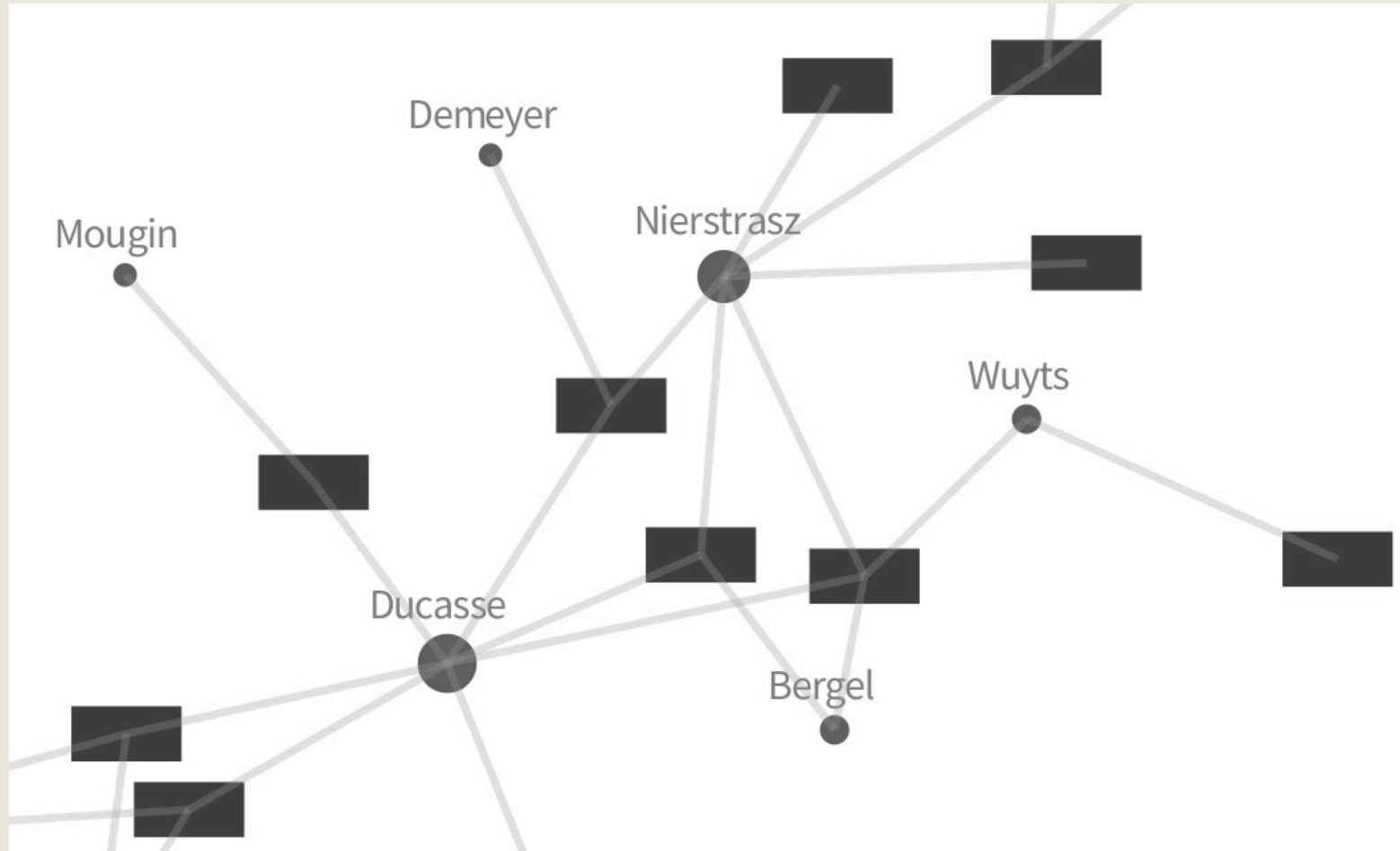# Who could I join forces with?

# Who could I join forces with?

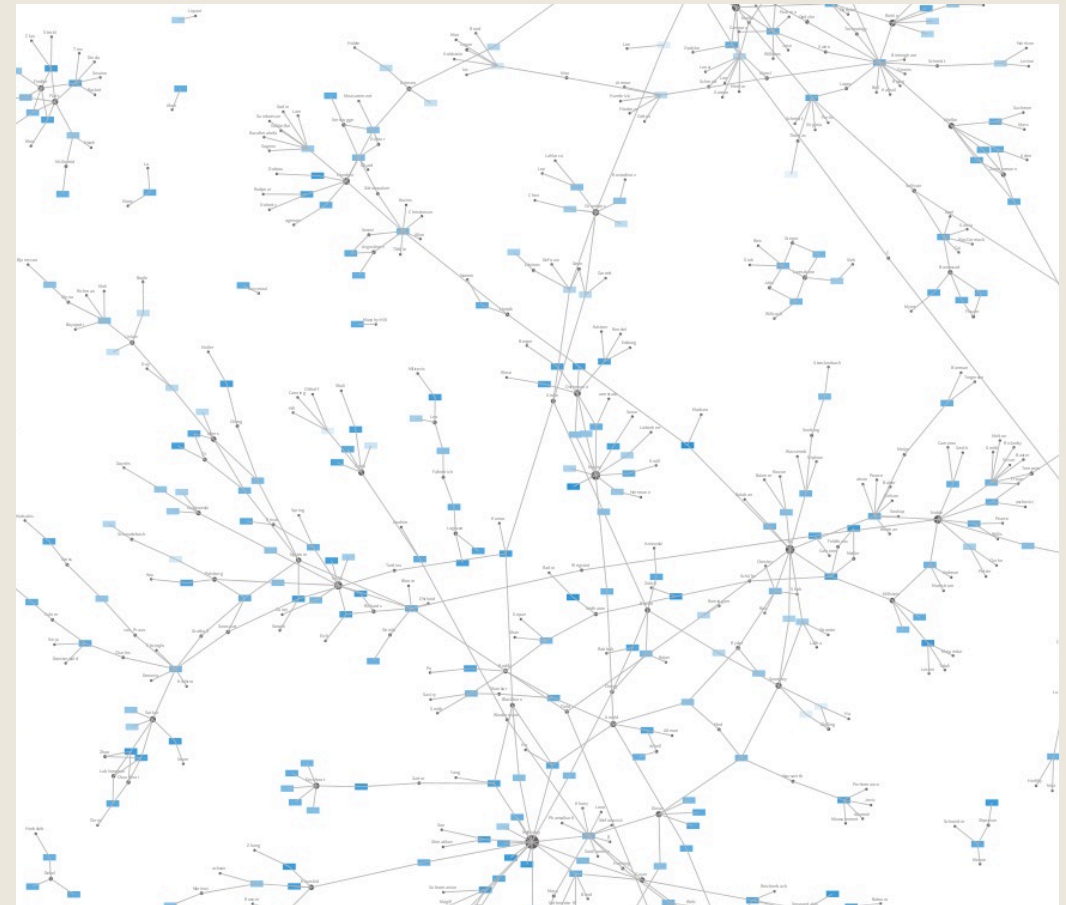# Not an easy question...

# A possible solution

# A possible solution

# The visualization

- Explorable graph

- Visualizes a scientific community in terms of papers, authors and authorship

- Additional visualization: word cloud for paper bodies

-> Live demo (corpus: 1'100 papers, published at OOPSLA, 1986-2015)

# This was just one possible query...

- How active was a certain field within the last couple of years?

- Is there a field that lately hasn't been covered anymore?

- Has a specific field recently gained more interest?

- What else can we find out about this community?

- ...

# Information is hidden in corpuses of papers

**Why Smalltalk Wins the Host Languages Shootout**

Lukas Renggli
renggli@iam.unibe.ch

Tudor Gîrba
girba@iam.unibe.ch

Software Composition Group, University of Bern, Switzerland
http://scg.unibe.ch/

**Runtime bytecode transformation for Smalltalk☆**

Marcus Denker[a,*], Stéphane Ducasse[a,b], Éric Tanter[c]

[a] Software Composition Group, IAM, Universität Bern, Switzerland
[b] Language and Software Evolution Group, LISTIC, Université de Savoie, France
[c] Center for Web Research, DCC, University of Chile, Santiago, Chile

**On the Integration of Smalltalk and Java**

Practical Experience with STX:LIBJAVA

Marcel Hlopko
Czech Technical University in Prague
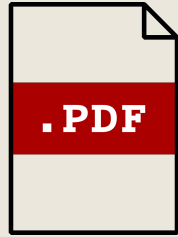
marcel.hlopko@fit.cvut.cz

Jan Kurš
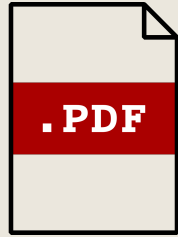Software Composition Group, University of Bern

kurs@iam.unibe.ch

Jan Vraný
Czech Technical University in Prague,
eXept Software AG

jan.vrany@fit.cvut.cz

Claus Gittinger
eXept Software AG

cg@exept.de

# What we need

# What we need



```
<algorithm name="ParsHed" version="110505">
<variant no="0" confidence="0.060039">
<title confidence="0.99942">Why Smalltalk Wins the Host Lan
<author confidence="0.999001">Lukas Renggli</author>
<email confidence="0.938906">renggli@iam.unibe.ch</email>
<author confidence="0.973453">Tudor Gîrba</author>
<email confidence="0.771593">girba@iam.unibe.ch</email>
<address confidence="0.362874">Software Composition Group,
<web confidence="0.99638">http://scg.unibe.ch/</web>
<abstract confidence="0.999470666666667">Integration of mul
```
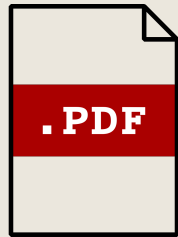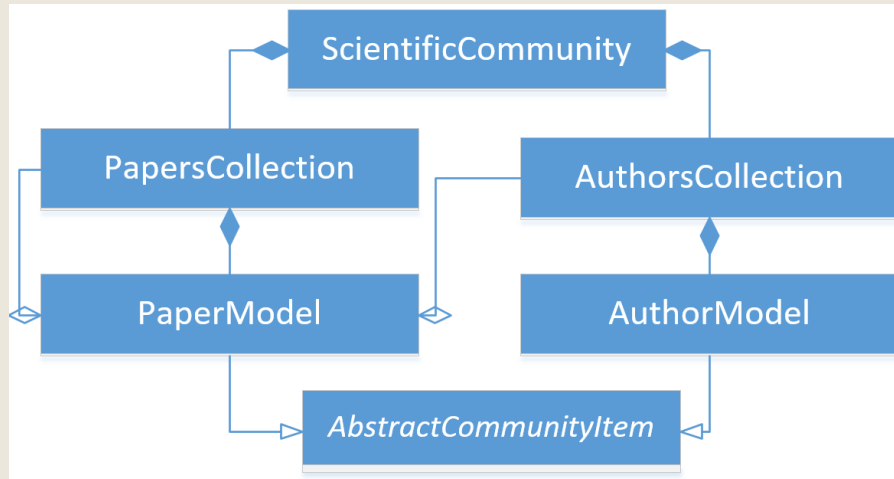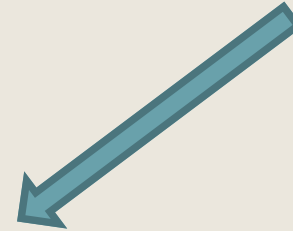
# What we need



```
<algorithm name="ParsHed" version="110505">
<variant no="0" confidence="0.060039">
<title confidence="0.99942">Why Smalltalk Wins the Host Lan
<author confidence="0.999001">Lukas Renggli</author>
<email confidence="0.938906">renggli@iam.unibe.ch</email>
<author confidence="0.973453">Tudor Gîrba</author>
<email confidence="0.771593">girba@iam.unibe.ch</email>
<address confidence="0.362874">Software Composition Group,
<web confidence="0.99638">http://scg.unibe.ch/</web>
<abstract confidence="0.999470666666667">Integration of mul
```
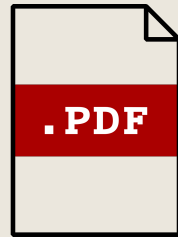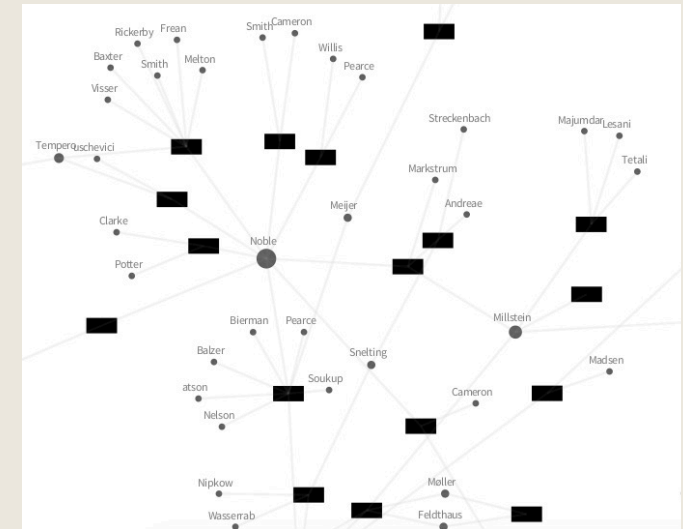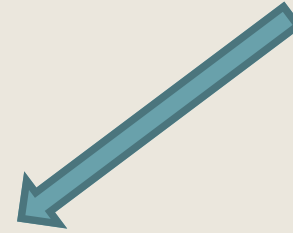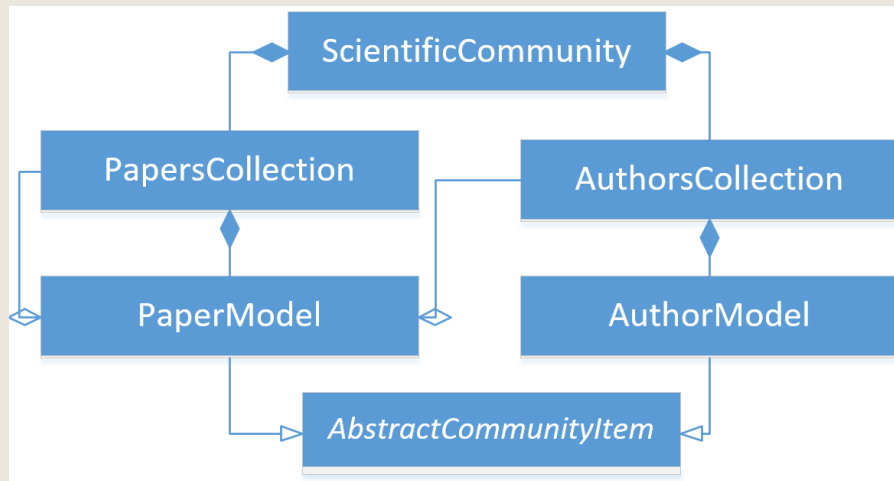
# What we need



```
<algorithm name="ParsHed" version="110505">
<variant no="0" confidence="0.060039">
<title confidence="0.99942">Why Smalltalk Wins the Host Lan
<author confidence="0.999001">Lukas Renggli</author>
<email confidence="0.938906">renggli@iam.unibe.ch</email>
<author confidence="0.973453">Tudor Gîrba</author>
<email confidence="0.771593">girba@iam.unibe.ch</email>
<address confidence="0.362874">Software Composition Group,
<web confidence="0.99638">http://scg.unibe.ch/</web>
<abstract confidence="0.999470666666667">Integration of mul
```
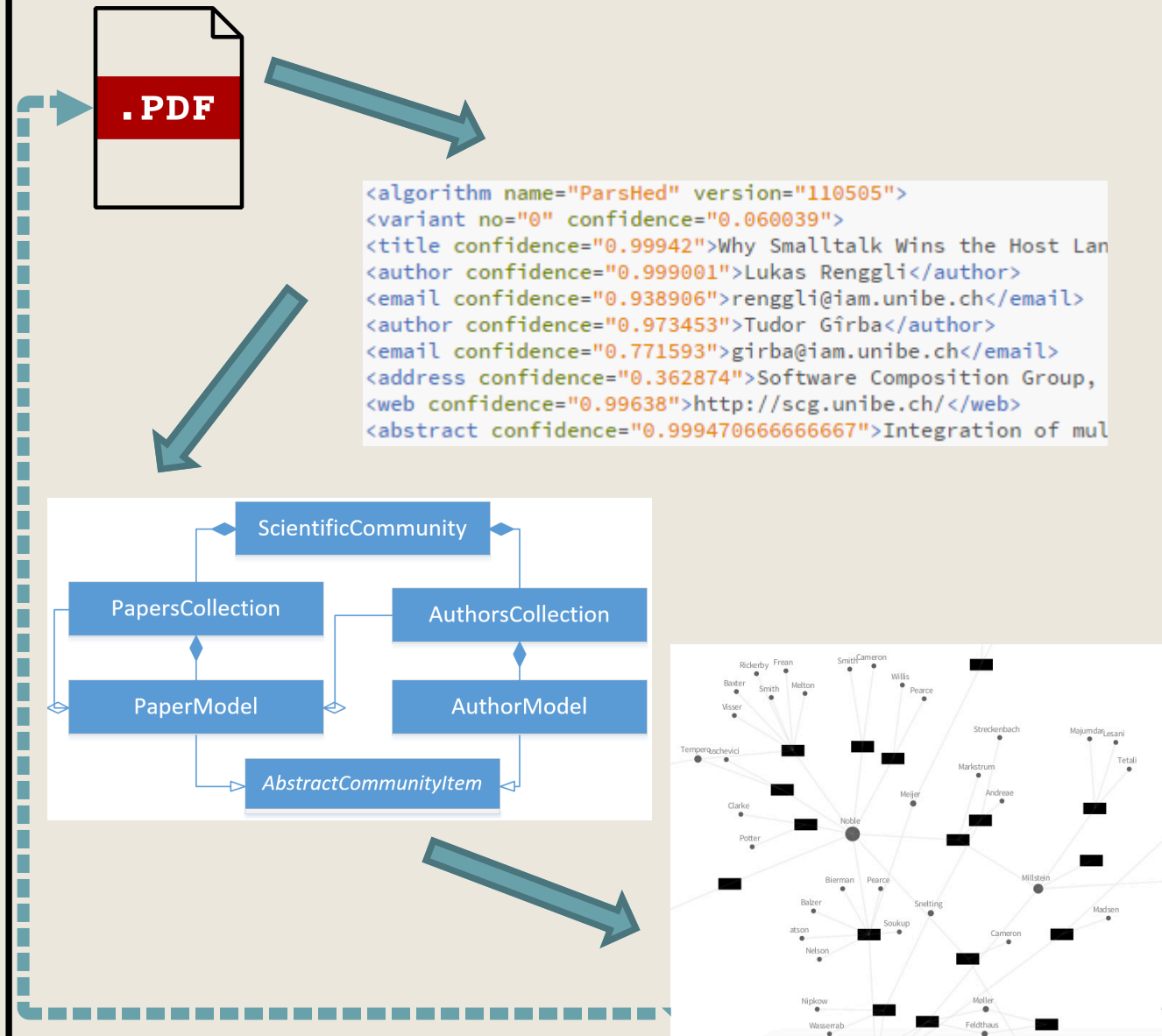


ScientificCommunity

PapersCollection

AuthorsCollection

PaperModel

AuthorModel

*AbstractCommunityItem*

# What we built

- A pipeline to convert PDFs into metadata-XML

- A query-able model of the underlying scientific community

- A graph and word cloud to visualize the model

- Builds on "EggShell", by Dominik Seliner

ExtendedEggShell



```
<algorithm name="ParsHed" version="110505">
<variant no="0" confidence="0.060039">
<title confidence="0.99942">Why Smalltalk Wins the Host Lan
<author confidence="0.999001">Lukas Renggli</author>
<email confidence="0.938906">renggli@iam.unibe.ch</email>
<author confidence="0.973453">Tudor Gîrba</author>
<email confidence="0.771593">girba@iam.unibe.ch</email>
<address confidence="0.362874">Software Composition Group,
<web confidence="0.99638">http://scg.unibe.ch/</web>
<abstract confidence="0.999470666666667">Integration of mul
```

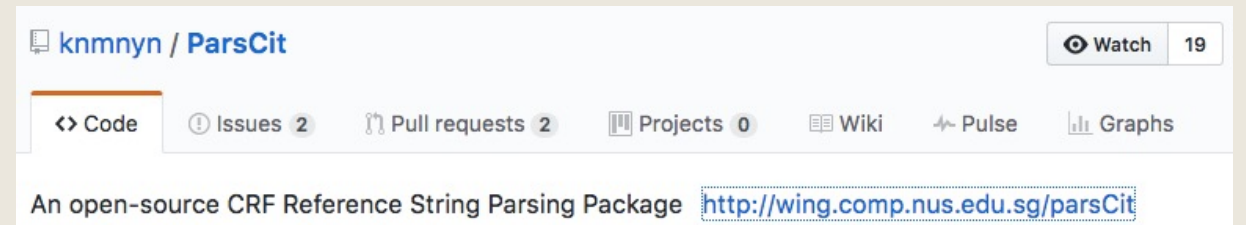*Image source: https://img.clipartfest.com*

# The metadata extraction

- PDFBox: PDF to text

- ParsCit: text to metadata-XML

■ Use third-party command line tools, used through controllers in ExtendedEggShell

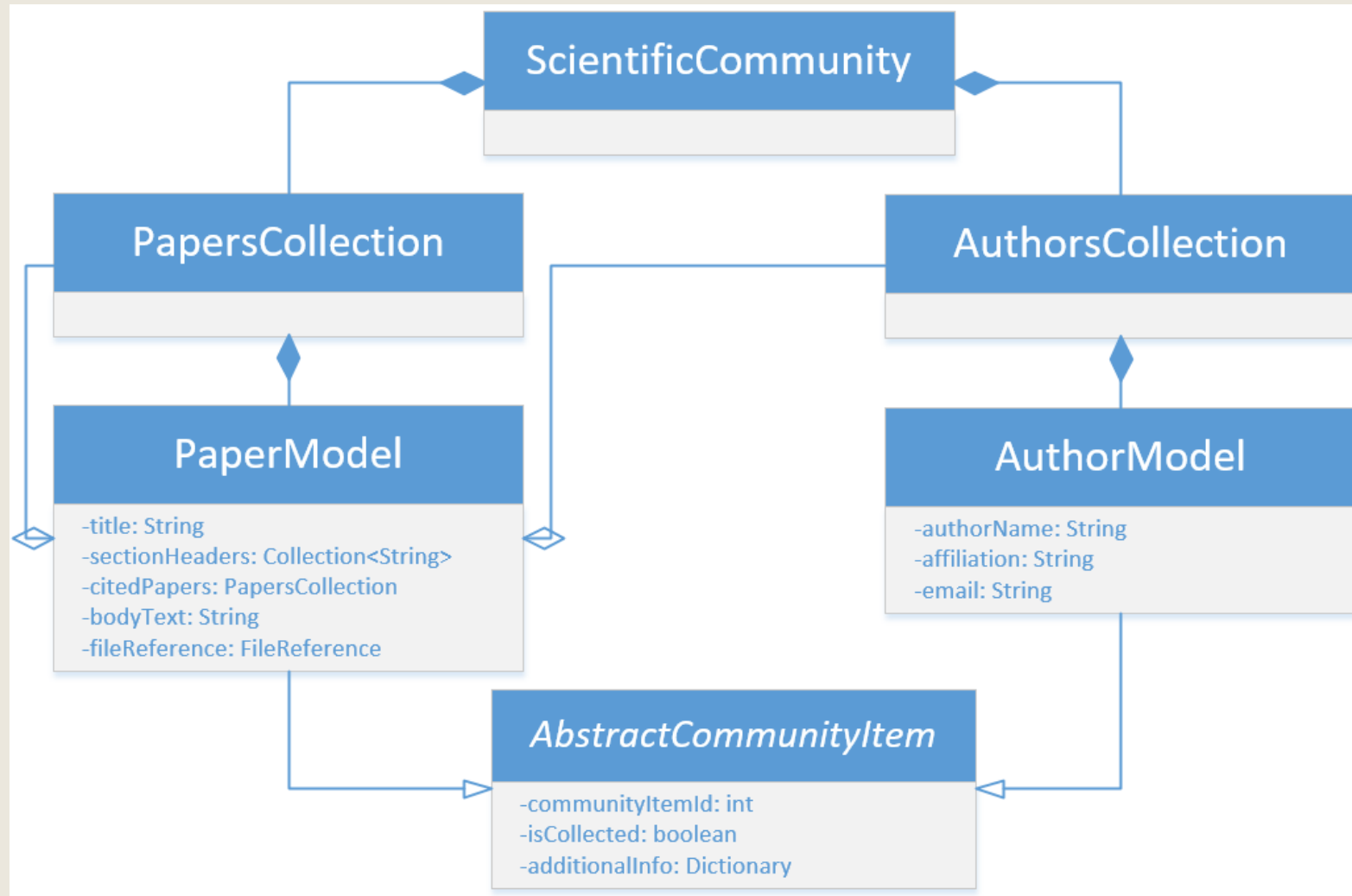■ First PDF to text, then metadata extraction from text

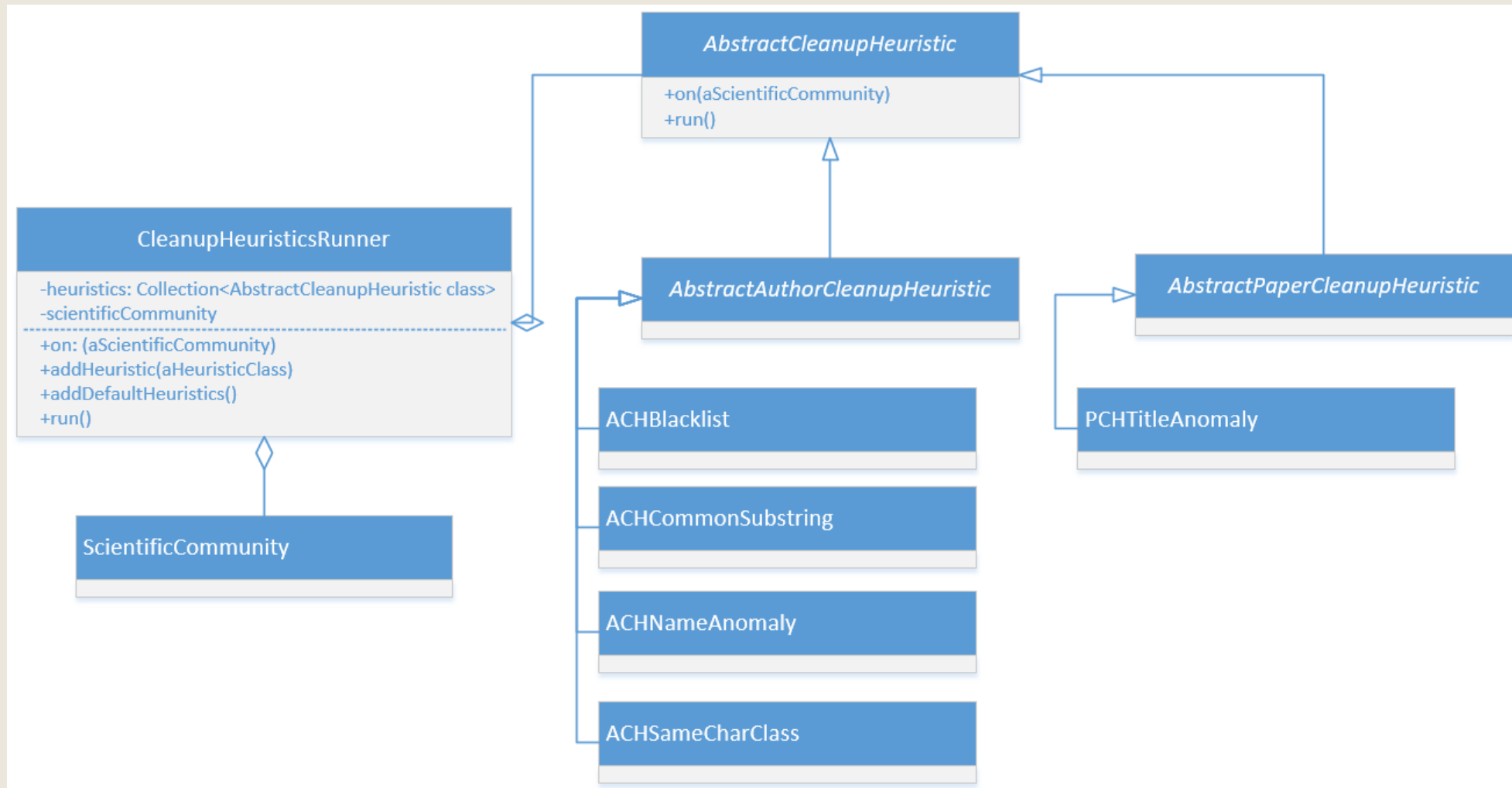# The Scientific Community model

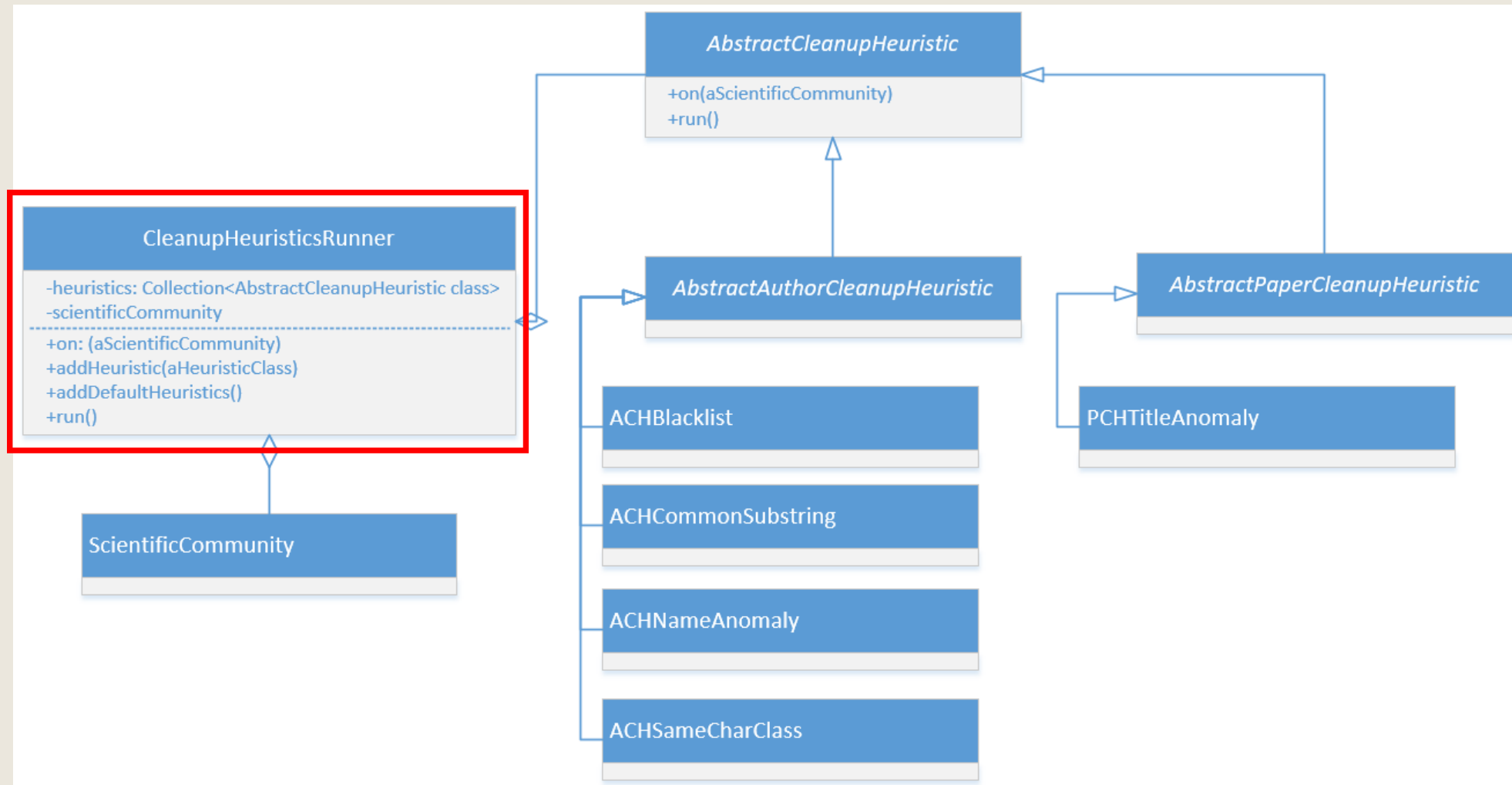# Initial model isn't perfect

- Non-alphanumeric characters in author names

- Multiple spellings of the same name

- Bad name extraction

# Model clean-up: heuristics framework

# Model clean-up: heuristics framework

# Model clean-up: heuristics runner

## CleanupHeuristicsRunner

-heuristics: Collection<AbstractCleanupHeuristic class>
-scientificCommunity

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

+on: (aScientificCommunity)
+addHeuristic(aHeuristicClass)
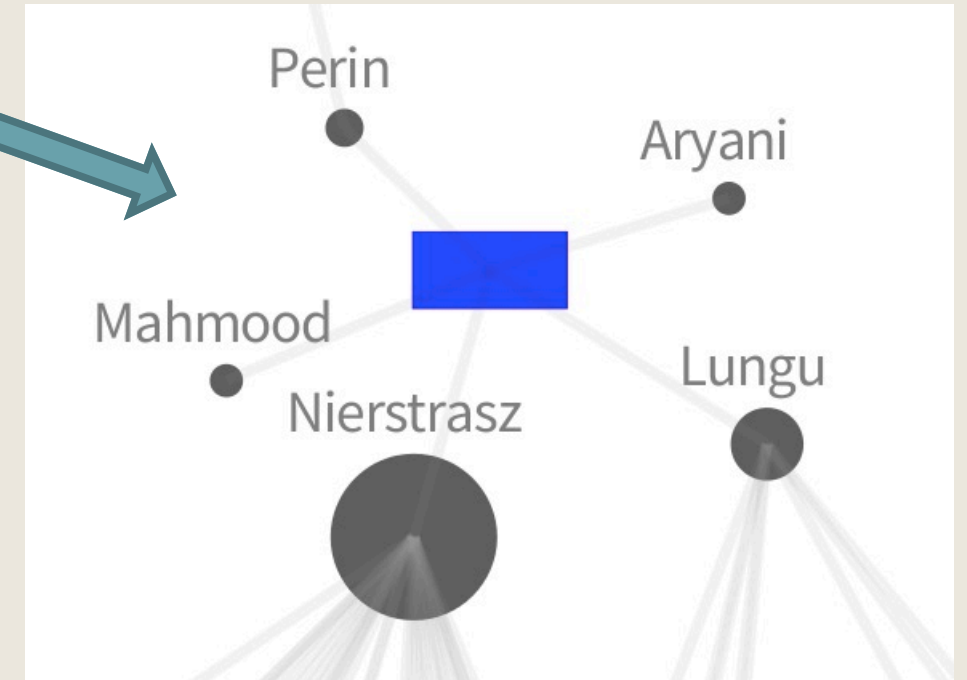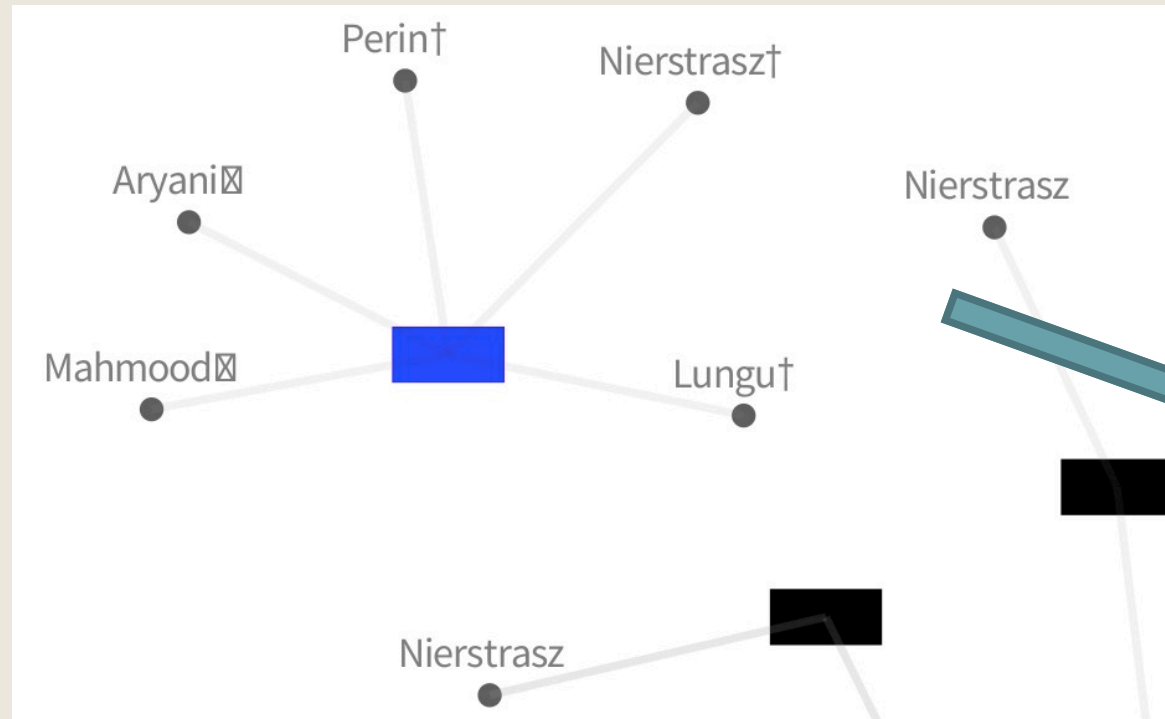+addDefaultHeuristics()
+run()

# Model clean-up: heuristics framework

# Model clean-up: example heuristics

# Model clean-up: before and after

# Conclusion

- Started with a complex, large dataset and built an explorable model and visualization, that allows for extracting insights from and about the dataset

- Users can create custom queries and have them answered visually

- Query-able model allows for easy creation of custom visualizations

- Scalability is limited

# Future work

- Further improve the model

- Web crawler for fetching cited papers

- Suggest related papers for some paper or author

# Some handy model query methods

## PaperModel

+numberOfBodyOccurrences(query)
+bodyContains(query)
+bodyContainsAll(query)
+bodyContainsAny(query)

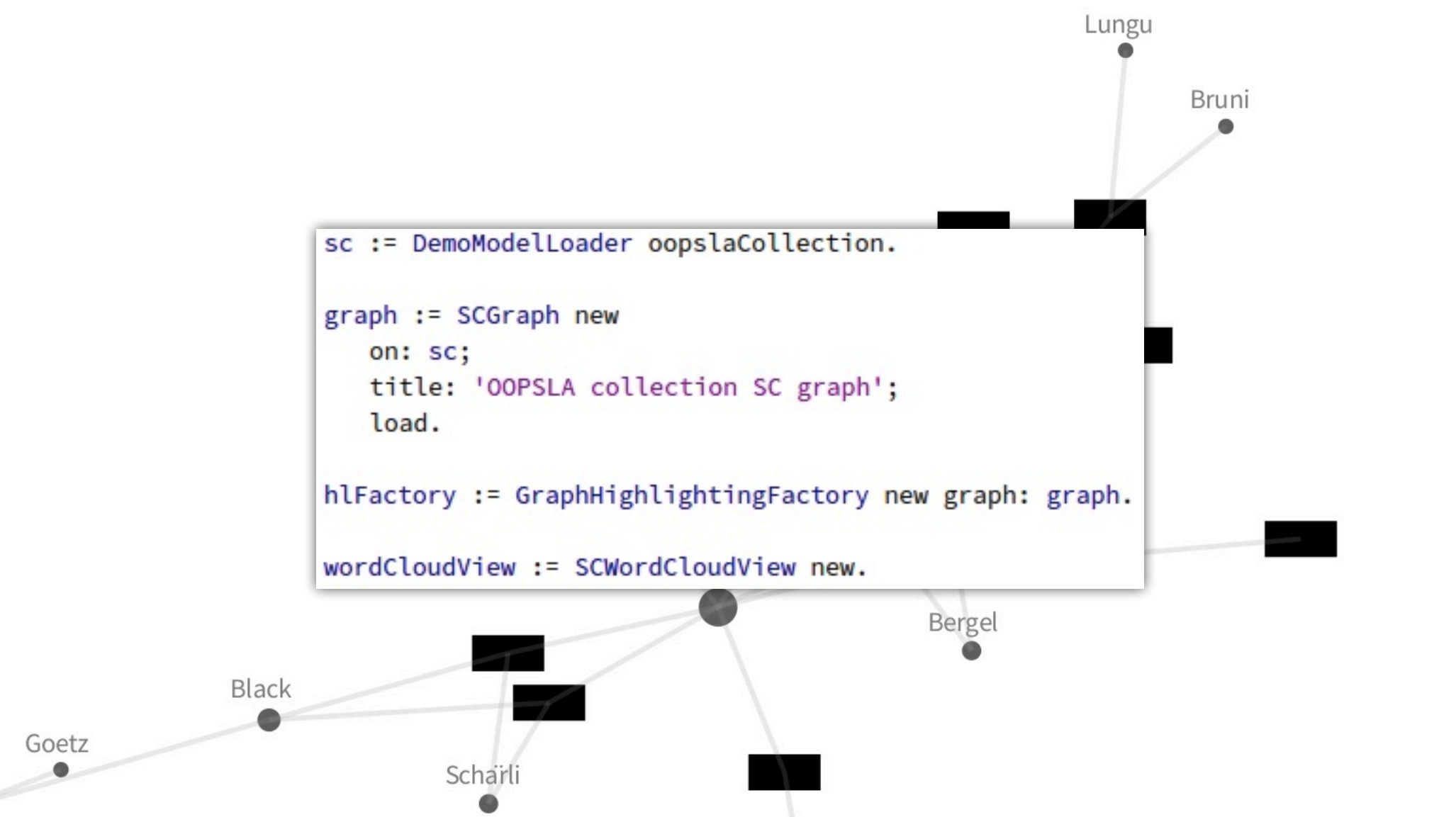## PapersCollection

+atTitle(query)
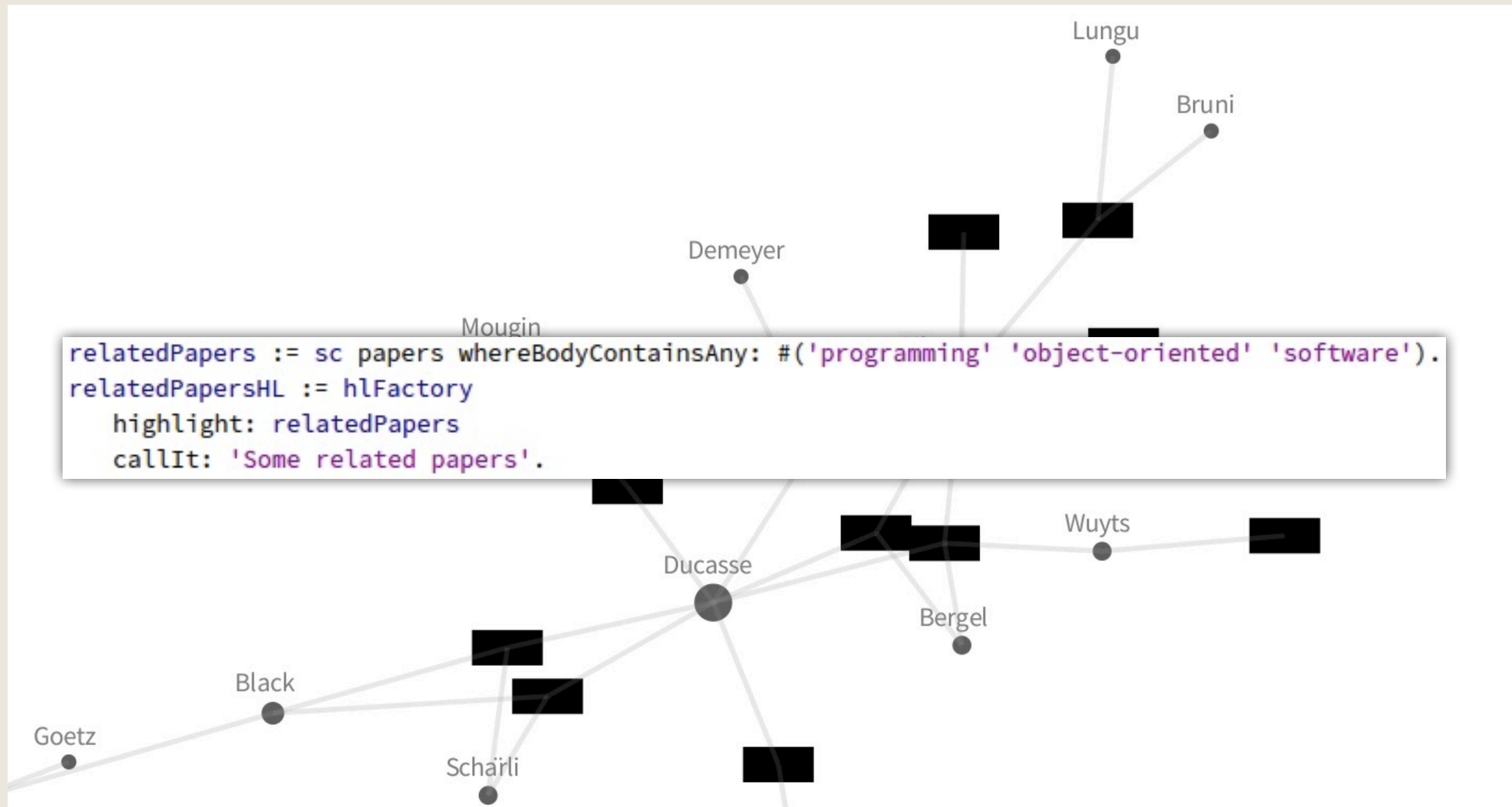+atTitleLike(query)
+atAuthorNameSubstring(query)

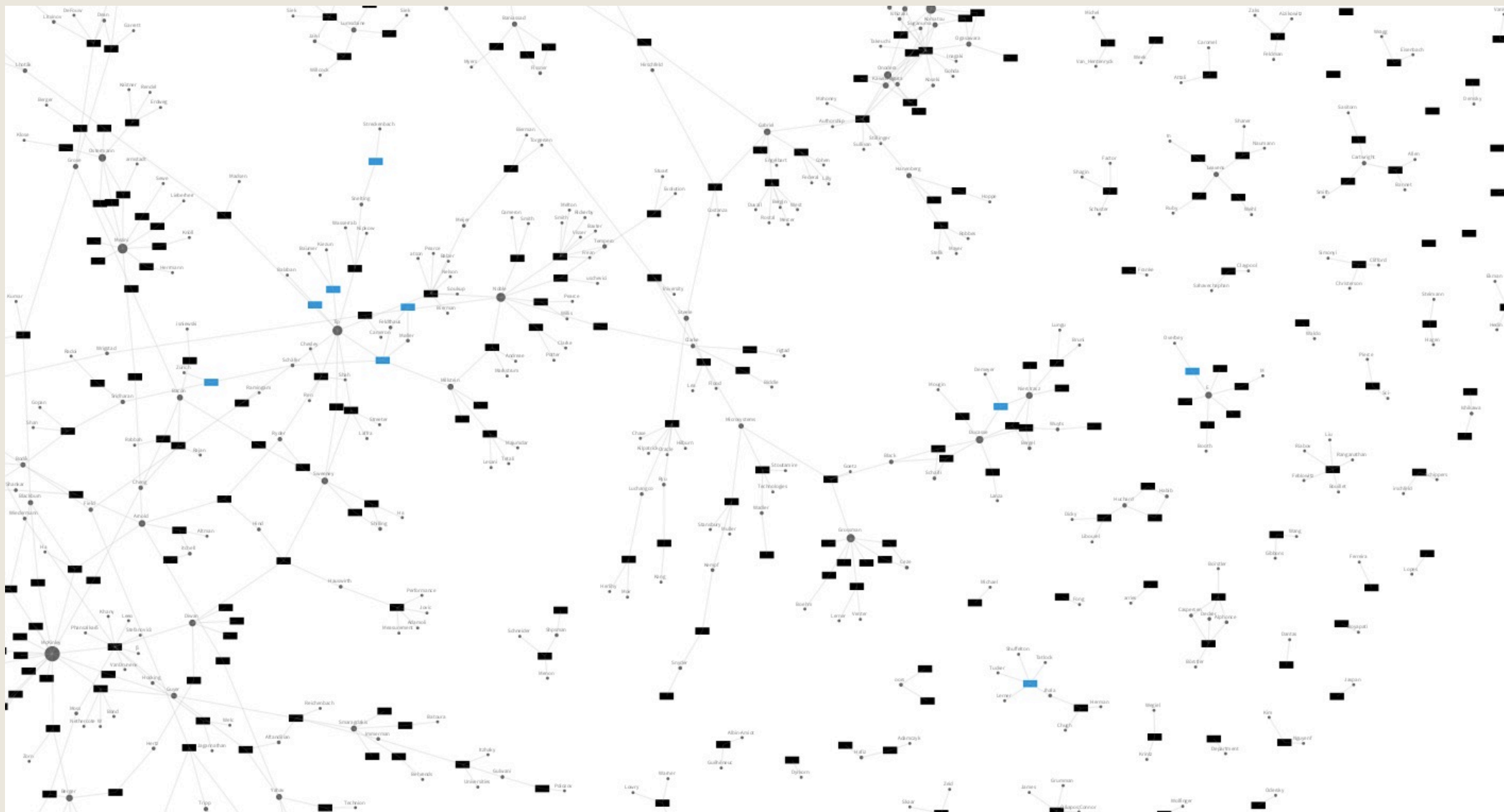# Some handy model query methods

**AuthorsCollection**

+atAuthorNameLike(query)
+atAuthorNameLike(query)
+atAuthorNameSubstring(query)
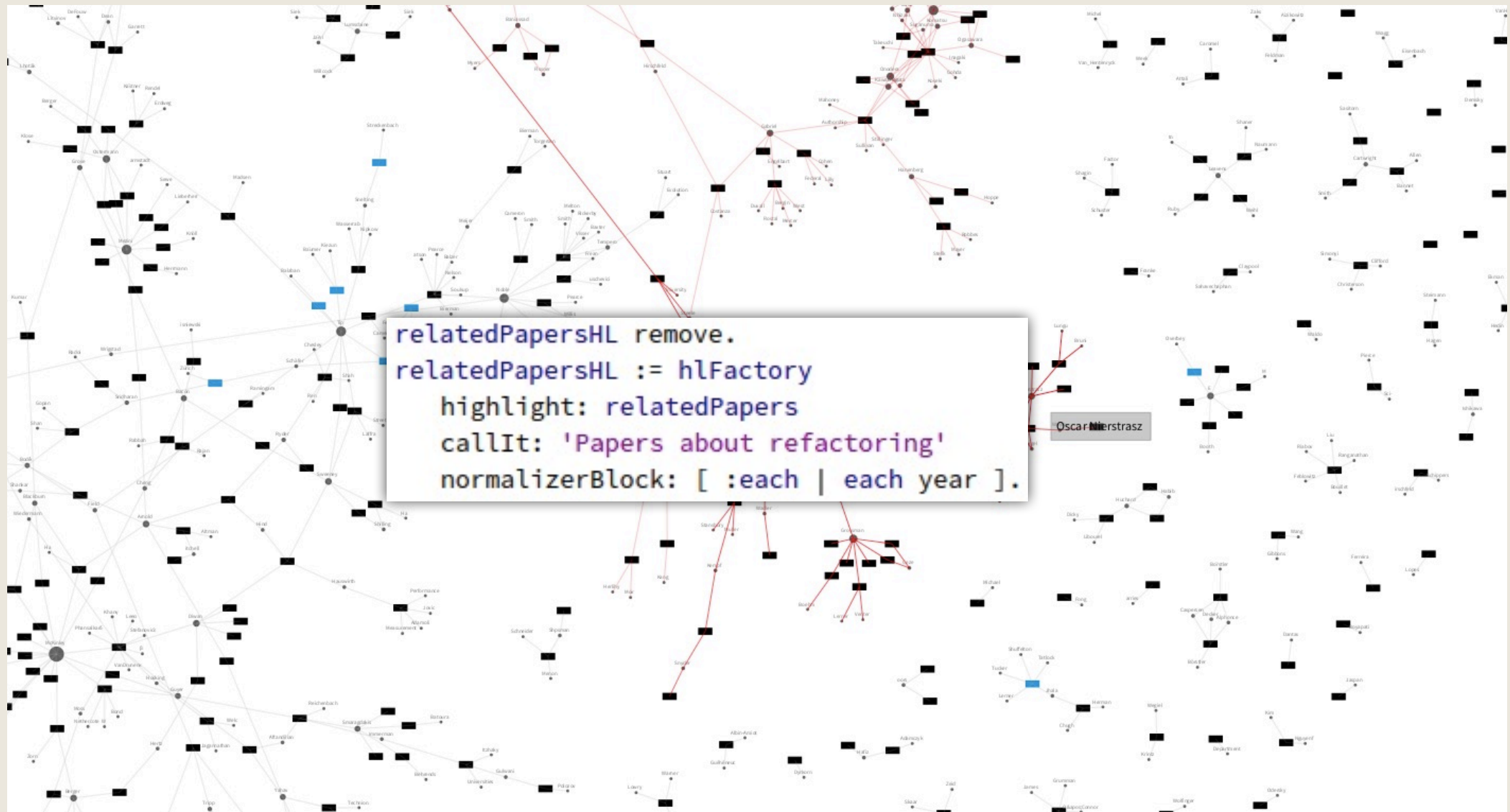
Bruni

```
sc := DemoModelLoader oopslaCollection.

graph := SCGraph new
    on: sc;
    title: 'OOPSLA collection SC graph';
    load.

hlFactory := GraphHighlightingFactory new graph: graph.

wordCloudView := SCWordCloudView new.
```

Bergel

Black

Goetz

Schärli

```
relatedPapers := sc papers whereBodyContainsAny: #('programming' 'object-oriented' 'software').
relatedPapersHL := hlFactory
    highlight: relatedPapers
    callIt: 'Some related papers'.
```

```
relatedPapersHL remove.

relatedPapers := sc papers atTitleSubstring: 'refactoring'.
relatedPapersHL := hlFactory
    highlight: relatedPapers
    callIt: 'Papers about refactoring'.
```

```
relatedPapersHL remove.
relatedPapersHL := hlFactory
    highlight: relatedPapers
    callIt: 'Papers about refactoring'
    normalizerBlock: [ :each | each year ].
```

```
interestingAuthors :=  {
    sc authors atAuthorName: 'Anders Møller'.
    sc authors atAuthorName: 'Asger Feldthaus'. }.

interstingAuthorsHL := hlFactory
    highlight: interestingAuthors
    callIt: 'Some interesting authors'.
```

```
papersMollerFeldthaus := sc papersByAny: #('Anders Møller', 'Asger Feldthaus').
wordCloudView buildCloudFor: papersMollerFeldthaus.
```

Papers about refactoring (18)

Some interesting authors (2)