# Data mining, management and visualization in large scientific corpus

HUI WEI

# Data collection

Some digital libraries did not supply APIs
We use raw PDF docs as input

# Data collection

1. to extract basic information of a paper such as authors, title, abstract sentences, doi

2. to extract references

3. to extract standard keywords and their frequency from each paper.

# Text mining

1. Use Jape rules to define "Macros" to find important markers,
   such as"DOI", "year", "abstract" tags.

2. Use Annie NE Transducer and Gazetteer look up
   person names like "author".

1. Use Gate ontology Gazetteer and Jape rules look up
   Computer Graphic terms in the content.

| ! | Name | Type |
|---|------|------|
| ● | Document Reset PR | Document Reset PR |
| ● | ANNIE English Tokeniser | ANNIE English Tokeniser |
| ● | ANNIE Sentence Splitter | ANNIE Sentence Splitter |
| ● | ANNIE POS Tagger | ANNIE POS Tagger |
| ● | morph | GATE Morphological analyser |
| ● | DivRootG | Onto Root Gazetteer |
| ● | DivFlexG | Flexible Gazetteer |
| ● | ANNIE Gazetteer | ANNIE Gazetteer |
| ● | Annie NE | ANNIE NE Transducer |
| ● | JapeCasa | JAPE Transducer |

Selected Processing resources

# Text mining



Figure 1. Metadata Extraction

# Keywords onto



Figure 2. Graph Model

| Subject | Predicate | Object |
|---------|-----------|--------|
| <http://www.DIV.org/Divdom/1#CG> | rdf:type | owl:Class |
| <http://www.DIV.org/Divdom/1#CG> | rdf:type | rdfs:Class |
| <http://www.DIV.org/Divdom/1#CG> | rdf:type | rdfs:Resource |
| <http://www.DIV.org/Divdom/1#CG> | rdfs:subClassOf | <http://www.DIV.org/Divdom/1#CG> |
| <http://www.DIV.org/Divdom/1#CG> | rdfs:subClassOf | rdfs:Resource |
| <http://www.DIV.org/Divdom/1#CG> | rdfs:subClassOf | owl:Thing |
| <http://www.DIV.org/Divdom/1#a_priori_estimate> | rdf:type | <http://www.DIV.org/Divdom/1#CG> |
| <http://www.DIV.org/Divdom/1#Abiotic_Factors> | rdf:type | <http://www.DIV.org/Divdom/1#CG> |
| <http://www.DIV.org/Divdom/1#Absorption_Coefficient> | rdf:type | <http://www.DIV.org/Divdom/1#CG> |
| <http://www.DIV.org/Divdom/1#Academic_Libraries> | rdf:type | <http://www.DIV.org/Divdom/1#CG> |
| <http://www.DIV.org/Divdom/1#Acceleration_of_Particles> | rdf:type | <http://www.DIV.org/Divdom/1#CG> |
| <http://www.DIV.org/Divdom/1#Accounting_Standards> | rdf:type | <http://www.DIV.org/Divdom/1#CG> |
| <http://www.DIV.org/Divdom/1#Accretion_Disk> | rdf:type | <http://www.DIV.org/Divdom/1#CG> |
| <http://www.DIV.org/Divdom/1#Accretive_Operator> | rdf:type | <http://www.DIV.org/Divdom/1#CG> |
| <http://www.DIV.org/Divdom/1#Acoustic_Communication> | rdf:type | <http://www.DIV.org/Divdom/1#CG> |

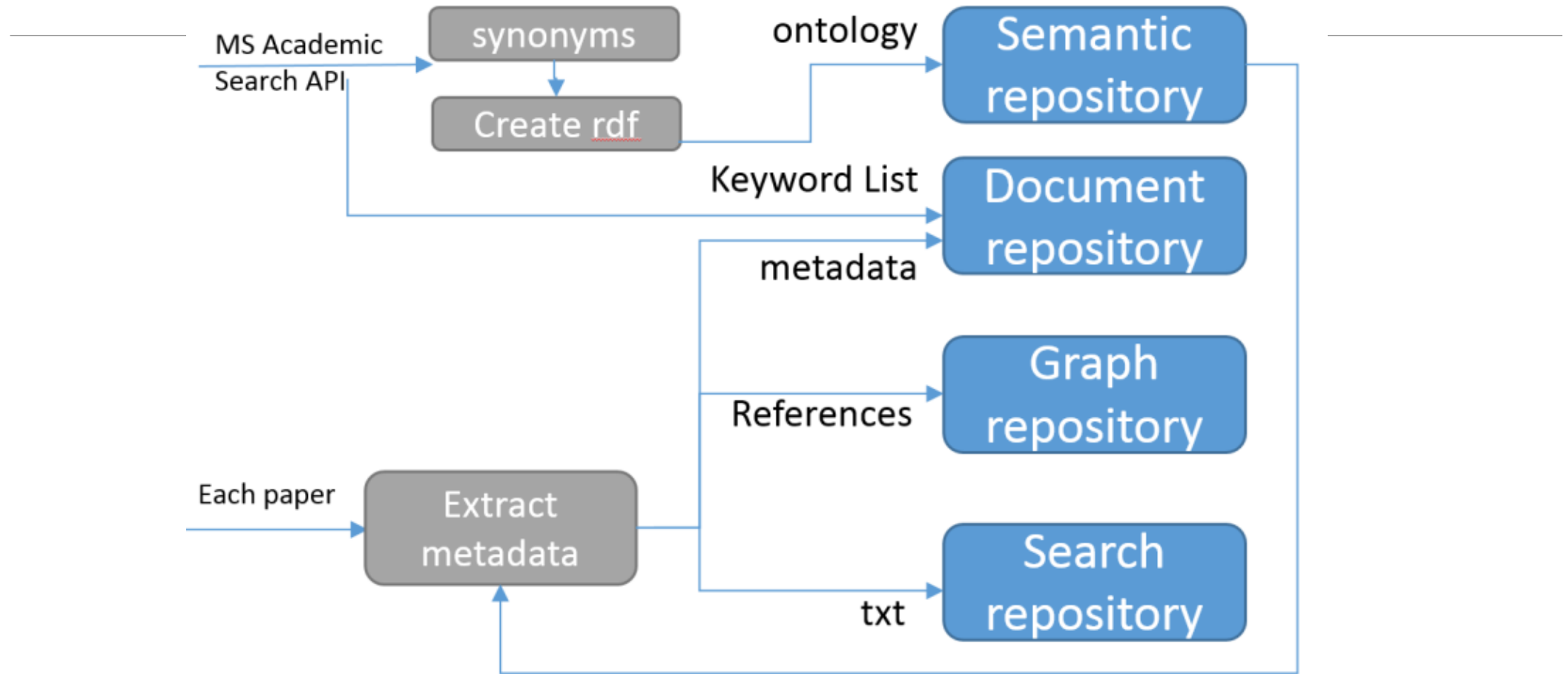# Data repositories



Graph repository

# Data repositories

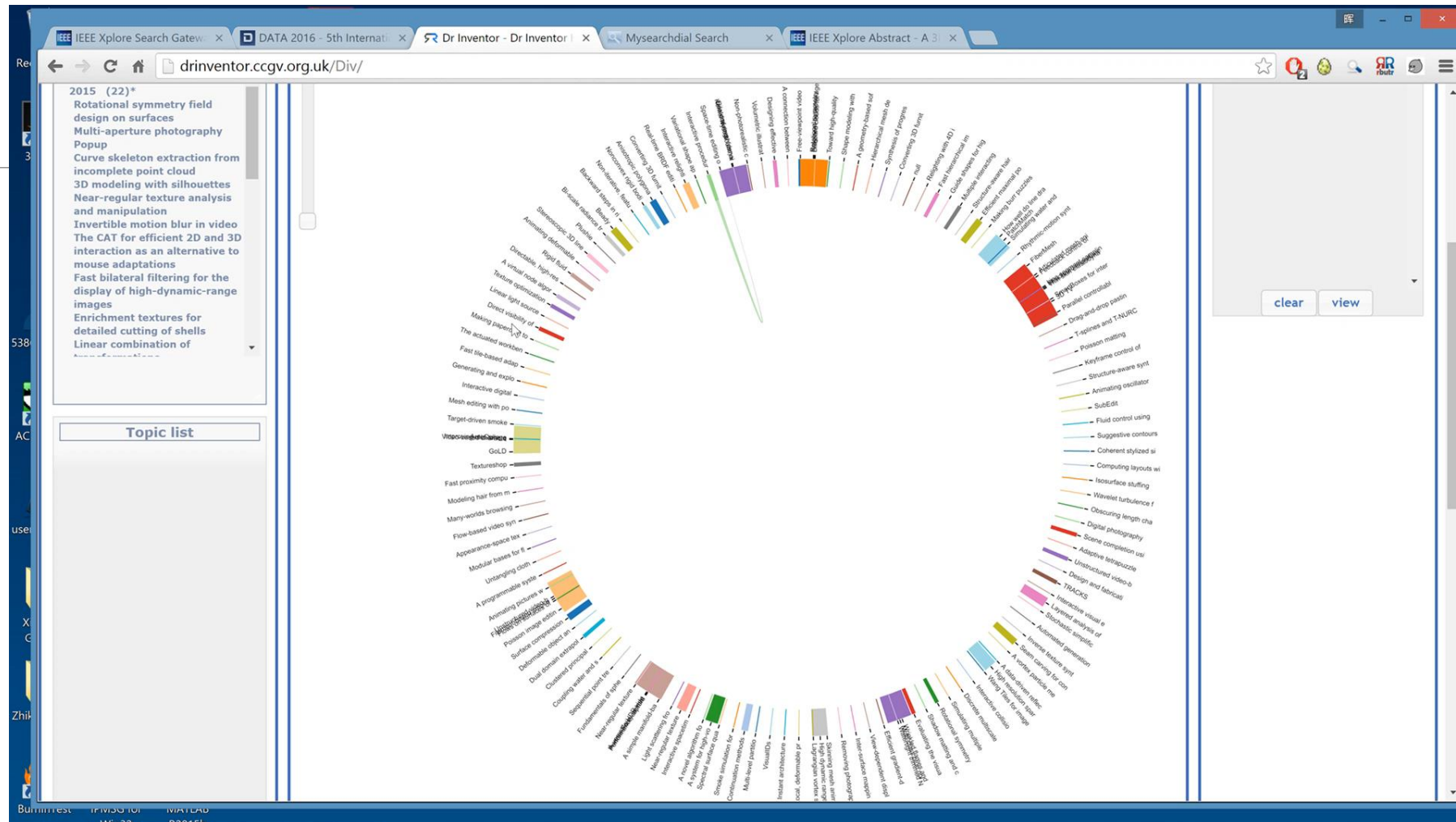

Data is managed in 4 NoSql repositories
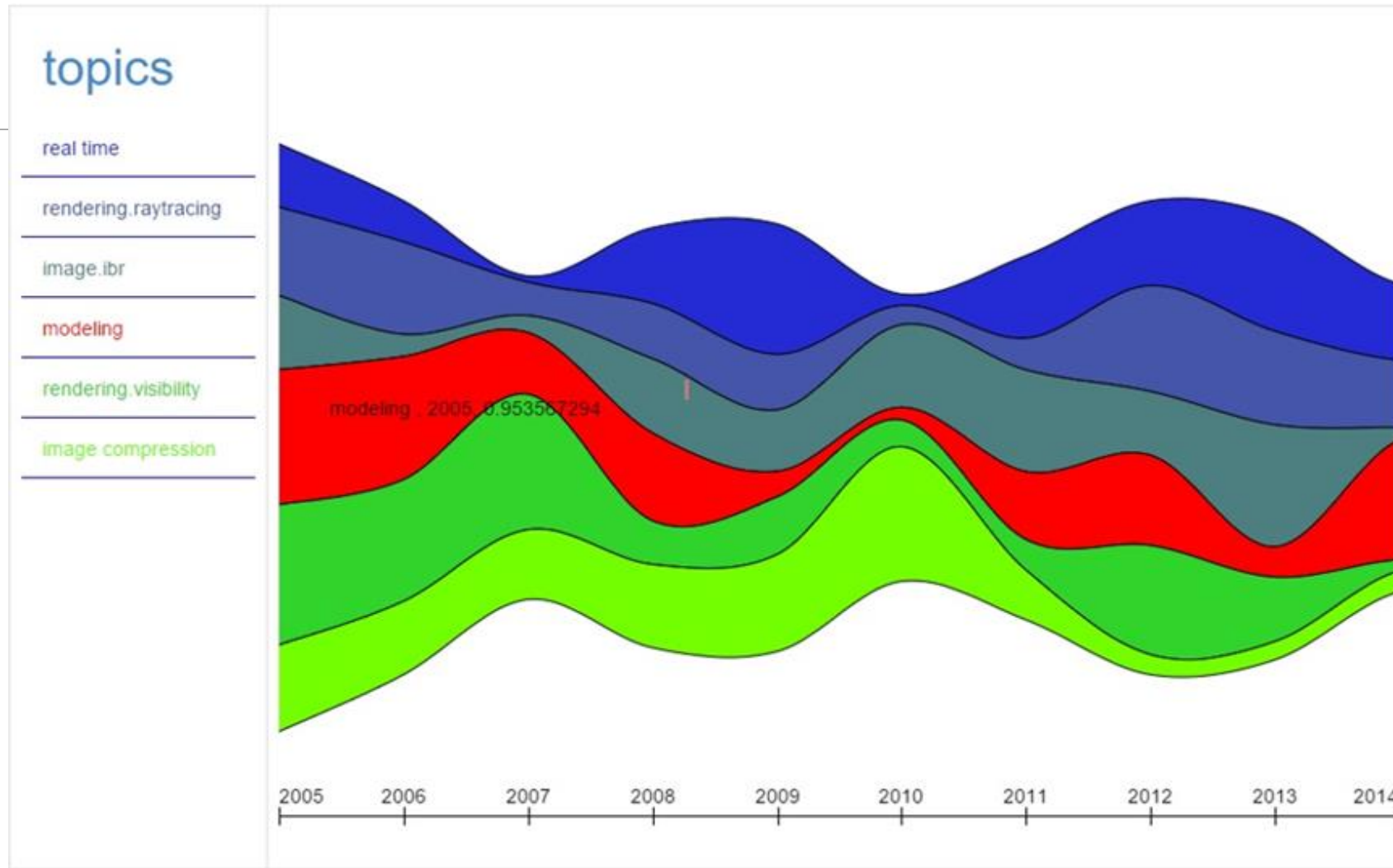
# Data repositories

Data distribution and system workflow

# Data visualization

# Topic river visualization

# Thanks

hui.wei@beds.ac.uk