

B2FIND records in neo4j

Vasily Bunakov
vasily.bunakov@stfc.ac.uk

EUDAT developers meeting in SURFSara,
Amsterdam, 6-8 October 2015

Data migration summary

- ~ 800 B2FIND data records imported in neo4j
- New information entities / graph nodes produced for Authors (491), Publishers(136), Tags(214) and Languages(2)
- **~ 34,000 relationships produced between B2FIND records and newly created entities**
- 10 ORCID records (for a selected Author) imported and mixed up with B2FIND records

Plain/raw B2FIND record uploaded in neo4j.
Take a note of multiple Authors squeezed in one field,
and the same applies to Publishers

The screenshot shows the Neo4j 2.2.3 web interface. The left sidebar contains navigation options: Node labels (DataRecord), Relationship types, Property keys (Author, Discipline, Format, ID, Language, Publisher, Title, Year), and Database (Location: /home/vb/neo4j/neo4j-community-2.2.3/data/graph.db, Size: 3.89 MiB).

The main area displays a Cypher query: `$ MATCH (n:DataRecord) RETURN n LIMIT 25`. The results are shown in a table format with columns for Discipline, Year, Language, Author, and Publisher. The Author and Publisher fields contain multiple values separated by semicolons.

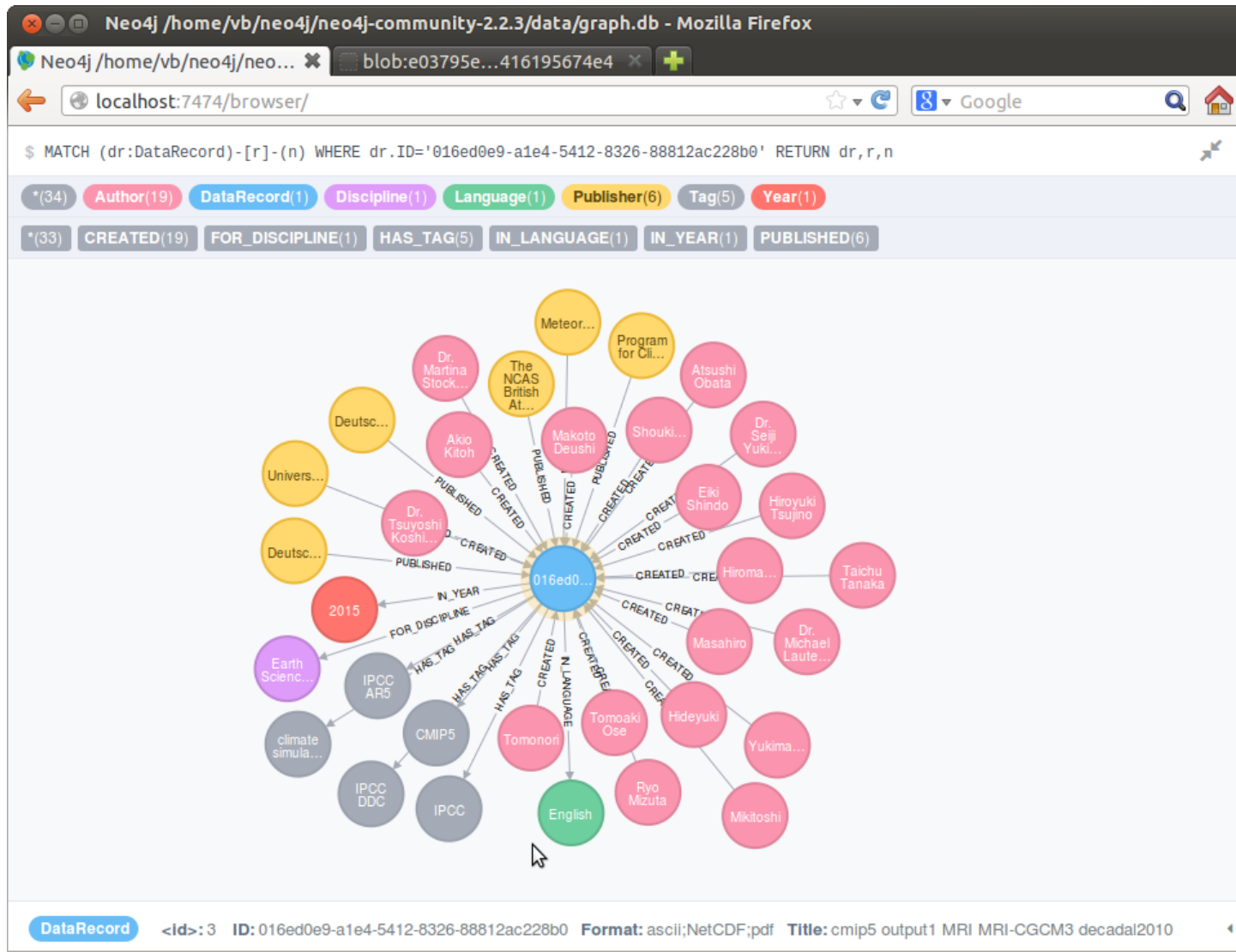
Discipline	Year	Language	Author	Publisher
Earth Sciences	2015	English	Jerry Meehl; Dr. Michael Lautenschlager; Dr. Martina Stockhause	National Center for Atmospheric Research;Deutsches Klimarechenzentrum;The NCAS British Atmospheric Data Centre;Program for Climate Model Diagnosis and Intercomparison;National Center for Atmospheric Research;National Center for Atmospheric Research;Deutsches Klimarechenzentrum (DM);Deutsches Klimarechenzentrum (DM);National Center for Atmospheric Research;Deutsches Klimarechenzentrum;The NCAS British Atmospheric Data Centre;Program for Climate Model Diagnosis and Intercomparison;National

Returned 25 rows in 355 ms.

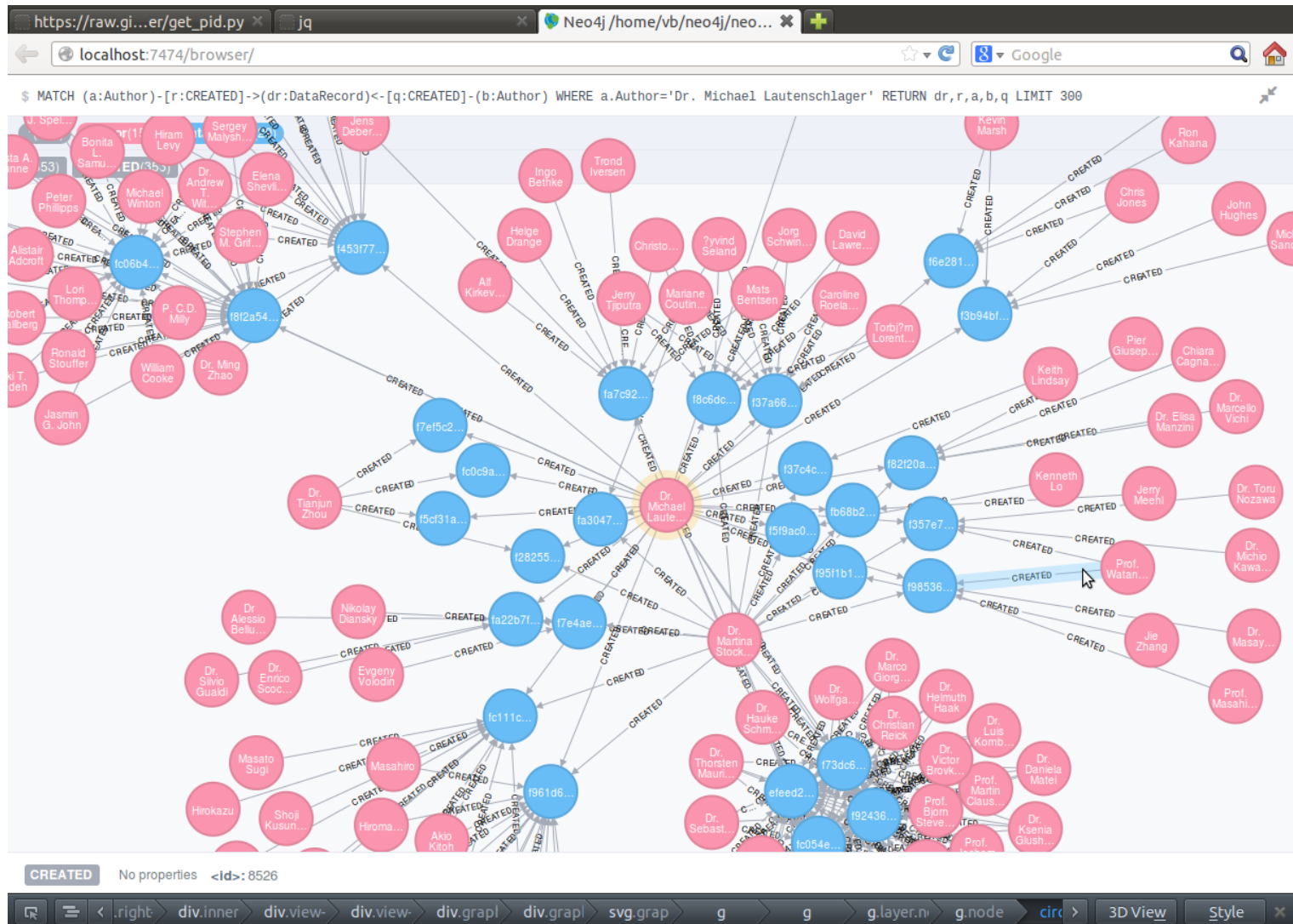
A second query is shown below: `$ MATCH (n) WHERE has(n.Language) RETURN DISTINCT "node" as element, n.Language...`. Its results are shown in a table with columns for element and Language.

element	Language
node	English
node	Multiple languages

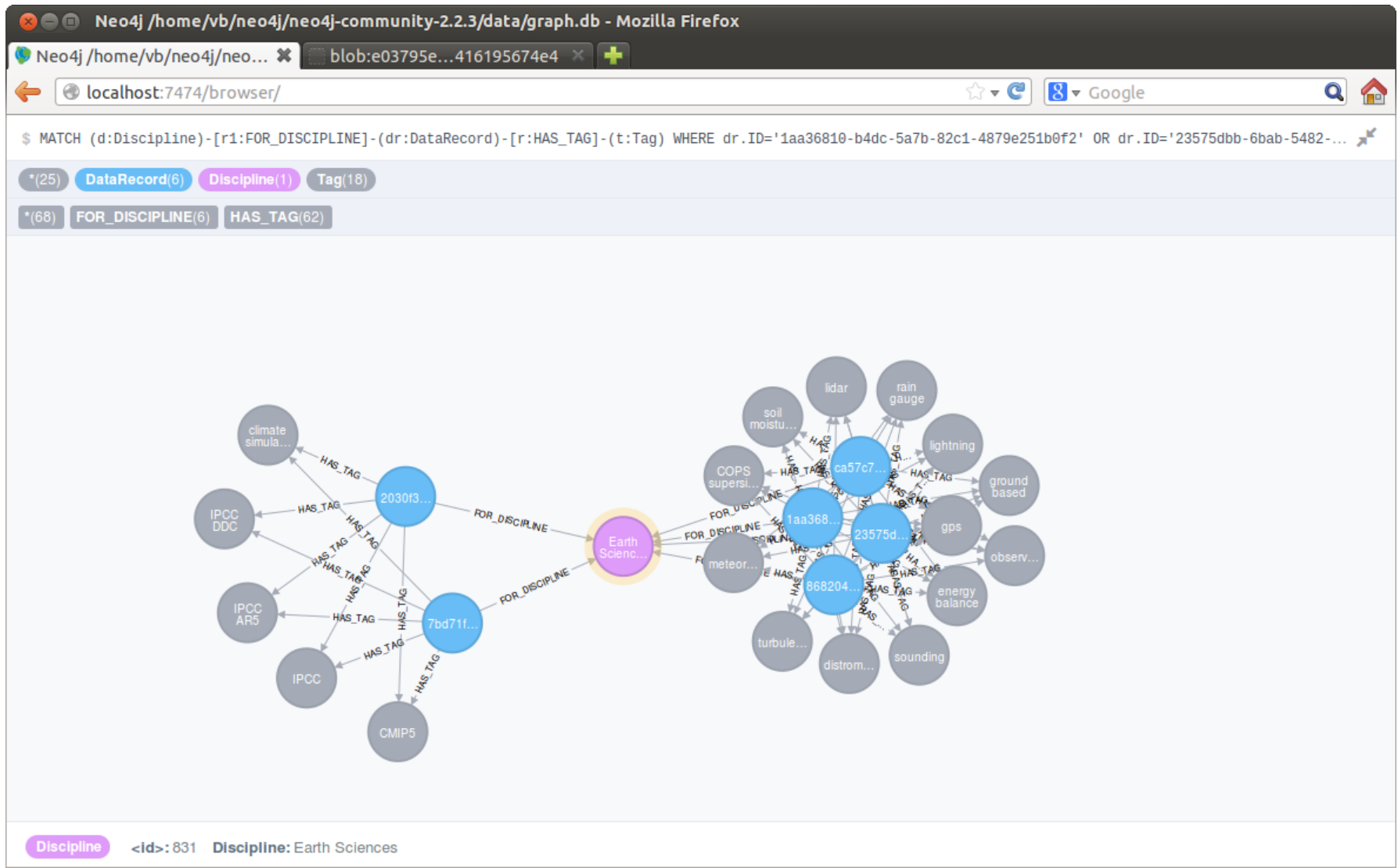
Many benefits should come from conversion of a raw B2FIND record into the graph model with new nodes created for Authors, Publishers, Languages and Tags



Authors (data creators) become new nodes.
This allows insights in data co-authorship and
contributes to finding related data



Tags become nodes. This allows looking for topical clusters and can support new flavours of data search, like “find heavily (or less heavily) related data”



Why (a small scale) experiment with ORCID records can be of interest

- ORCID records have differently structured metadata yet can be successfully mixed up with B2FIND records in one graph, then queried as necessary
- It turned out that none of data records obtained from ORCID were harvested by B2FIND, so ORCID can complement current B2FIND data sources
- Beyond ORCID, there are other perspective data sources to explore, e.g. Figshare
- In short, B2FIND now covers institutional repositories (and does this fairly well) yet ORCID and Figshare are examples of repositories popular with individual researchers; they are the “long tail of data records” that B2FIND may want to cover, too, to become a truly universal research data catalogue

