# Motivation & Aims

**Open Biodiversity Knowledge Management**

**Building a semantic knowledge graph from literature-extracted biodiversity data**

**Conceptual modelling of the biodiversity publishing domain**

# Introducing OpenBiodiv-O

OpenBiodiv-O

DarwinCore-based ontologies

SPAR ontologies

Journal of
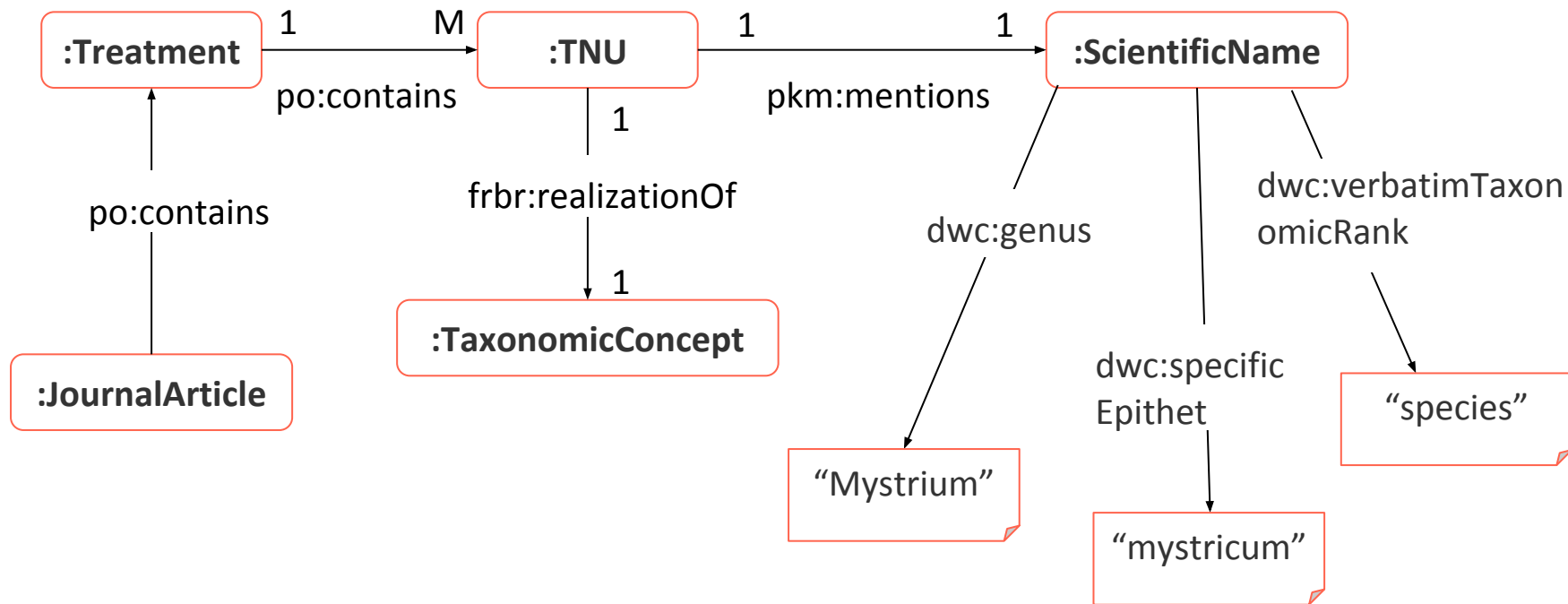Biomedical Semantics

**RESEARCH**

**Open Access**

CrossMark

# OpenBiodiv-O: ontology of the OpenBiodiv knowledge management system

Viktor Senderov[1,2]*, Kiril Simov[3], Nico Franz[4], Pavel Stoev[1,7], Terry Catapano[5], Donat Agosti[5], Guido Sautter[5], Robert A. Morris[6] and Lyubomir Penev[1,2]

# Biodiversity Publishing

- Taxonomic articles:
  - Abstract, Introduction, Materials and Methods, Results, Conclusions
  - Taxonomic treatment:
    - Nomenclature
    - Type material
    - Etymology
    - Diagnosis
    - Description
    - Distribution
- Metadata and various identifiers (DOI, ORCID, Zoobank ID)
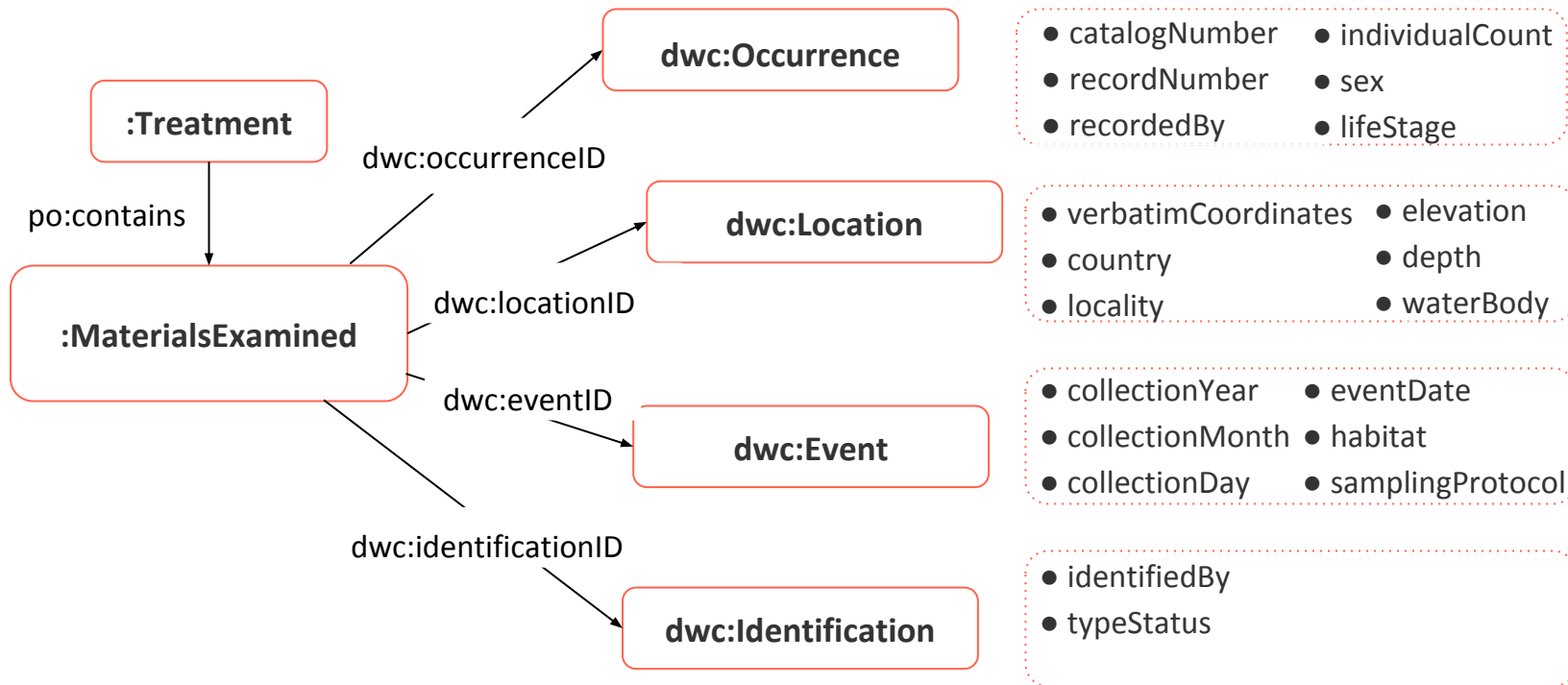- **Taxonomic names**

# Modelling Taxonomic Names

# Bottom-up approach to modifying OpenBiodiv-O

# 1. Materials Examined

# 2. Institutional identifiers

GRSciColl: The Global Registry of Scientific Collections

http://biocol.org/urn:lsid:biocol.org:col:34985

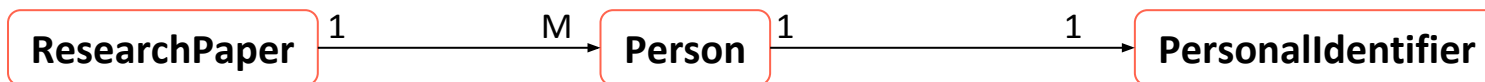# 3. Molecular resource identifiers

- OpenBiodiv-O, DataCite, Fabio

# 4. Personal identifiers

openbiodiv:8B0A6890-3094-4431-8262-23748A86B071   rdf:type   fabio:ResearchPaper ;
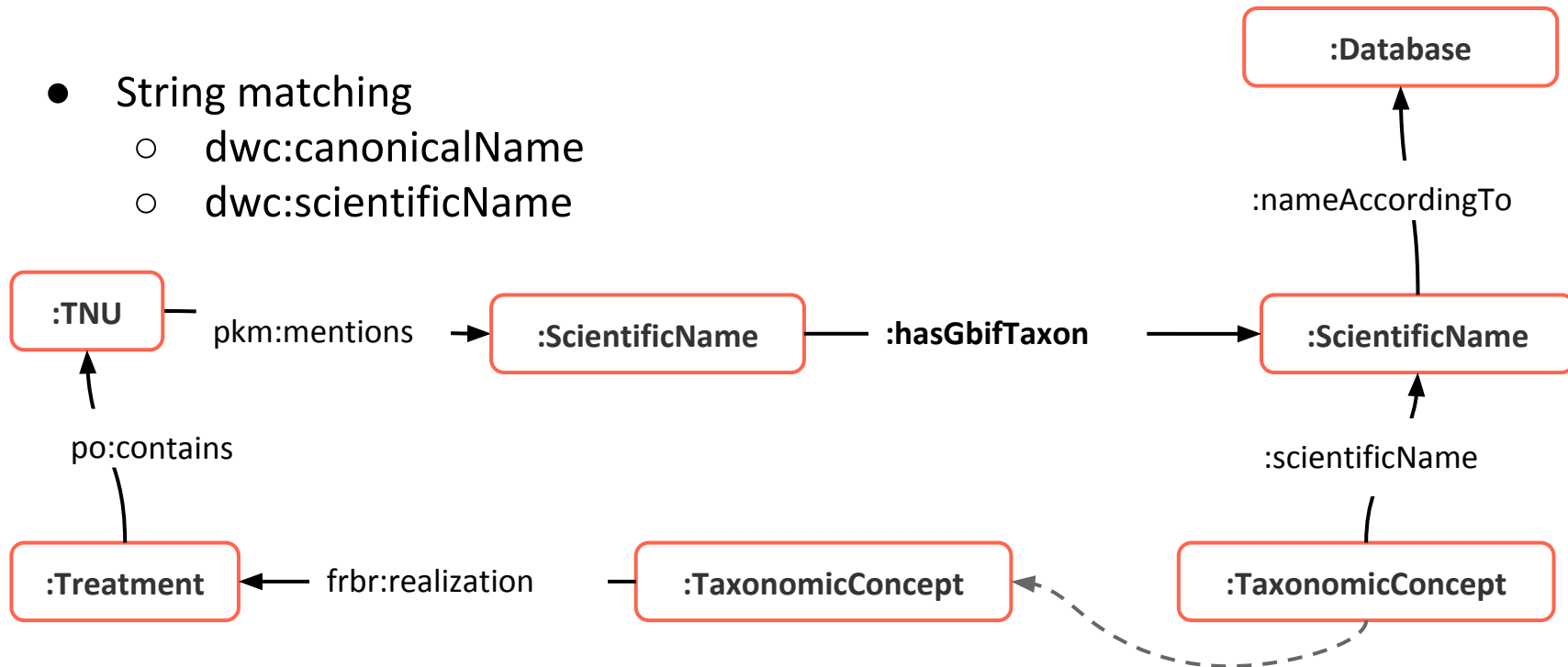        dcterms:creator   openbiodiv:7F19D49E-B4CA-4E7A-808D-A57DAA7E02A3.

openbiodiv:7F19D49E-B4CA-4E7A-808D-A57DAA7E02A3   rdf:type   foaf:Person ;
        rdfs:label   "Mengmeng Liu" ;
        openbiodiv:affiliation   "College of Ecology, Lishui University, Lishui, Zhejiang, China" ;
        datacite:hasIdentifier   orcid:0000-0002-0985-5852 .

orcid:0000-0002-0985-5852   rdf:type   datacite:PersonalIdentifier ;
        datacite:usesIdentifierScheme   datacite:orcid ;
        rdfs:label   "0000-0002-0985-5852" .

ResearchPaper  1 ──── M →  Person  1 ──── 1 →  PersonalIdentifier

# 5. Explicit mapping to GBIF's taxonomy

- String matching
  - dwc:canonicalName
  - dwc:scientificName

**Results**

## Use cases

- Modelling and linking resources across domains
    - Publishing
    - Taxonomy
    - Genomics
- Serving users from different groups
    - Taxonomists
    - Ecologists
    - Curators
    - Institutions

# Application                    PS 0101

## Next steps

- Facilitate federated SPARQL queries by mapping external ontologies to OpenBiodiv-O.
  - Wikidata
- Expanding the ontology as the number of extracted entity types increases

# Thank you!

## Questions?