# Structured & Unit Independent Search for DBRepo

Martin Weise, Sotirios Tsepelakis, Nikola Lukic, Max Spannring, Gökay Güçlü, Geoffrey Karnbach

# Motivation for a Database Repository

**Databases as important resources for research & industry**

1. Database paradigm is well-understood
2. Cost-efficient storage systems for data in use
3. Repositories as established systems to make research data FAIR

Devise a system that combines technological infrastructure with repository work-processes to provide machine-understandable data in databases.

CAVEAT: RDA WGDC recommendations https://doi.org/10.1162/99608f92.be565013 & technical knowledge gaps
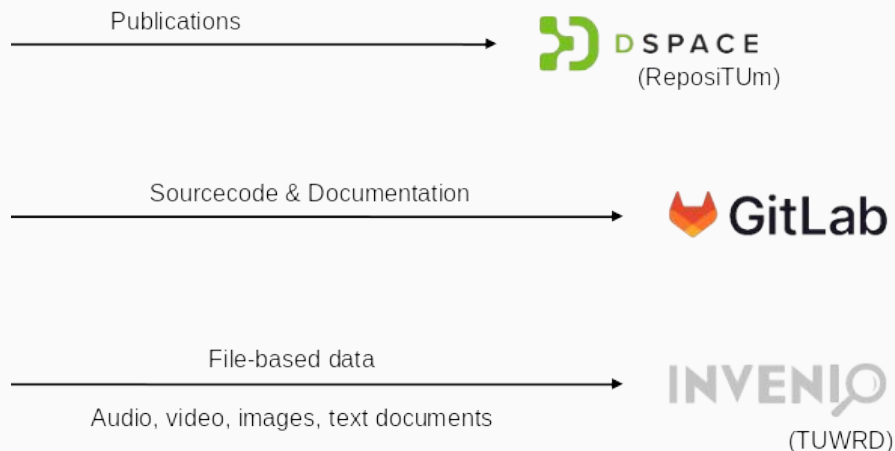
# Motivation for a Database Repository

TUWRD can handle collection of files

How about relational data in databases?

- Releasing a **data dump** every x amount of time?
- Adding **continuous data** streams, e.g. IoT?
- How to update / correct data in those databases?
- Allow **reproduction** of any subset?

Publications &rarr; D SPACE
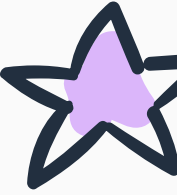(ReposiTUm)

Sourcecode & Documentation &rarr; GitLab

File-based data &rarr; INVENIO
Audio, video, images, text documents
(TUWRD)

**2020**

FAIR Data Austria
started

**2022**

FAIR Data Austria
ended

**2023**

.dcall
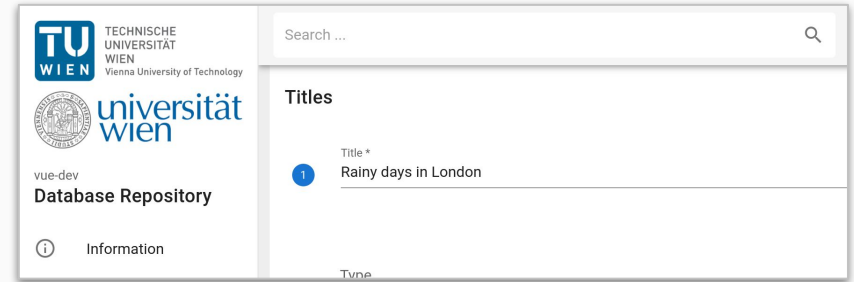Shared RDM started

1 → 2 → 3

# PREVIOUS ISSUES

ElasticSearch license change 2021

Wildcard search not accurate enough

Authentication directly to the database in the UI (unsafe!)

Does not scale

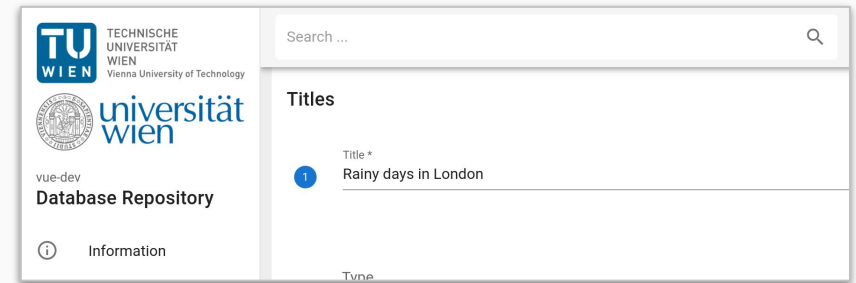Not enough data for meaningful queries

Old



... among many others

# WP1 - Extension of Indexed Metadata

## GOALS

Extend the indexed metadata in the search service to cover semantic concepts and units of measurement of columns for tables

Allow structured search through facets that assist users in filtering results based on semantic concept and/or unit of measurement

### Old

# WP1 - Extension of Indexed Metadata

### DELIVERIES

### Old

1. Faceted browsing based on semantic concepts
2. Faceted browsing based on units of measurement
3. Increase supported ontologies, e.g. I-ADOPT

/database
/table
/column
/user
/view
/identifier
/concept
/user

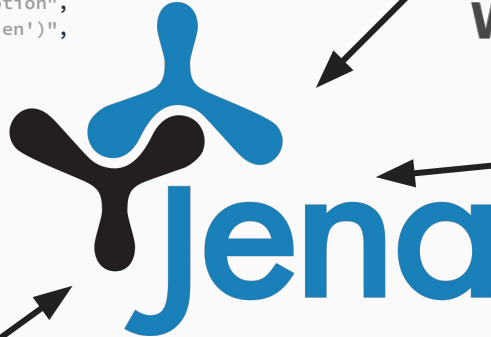OpenSearch

Refactor ElasticSearch
components

Migrate queries

Martin Weise

```java
default String ontologyToFindByLabelQuery(List<Ontology> ontologies, Ontology ontology, String label, Integer limit) {
    if (ontology.getSparqlEndpoint() != null) {
        /* prefer SPARQL endpoint over rdf */
        return String.join("\n",
            defaultNamespaces(ontologies),
            "SELECT * {",
            "  SERVICE <" + ontology.getSparqlEndpoint() + "> {",
            "    SELECT ?o ?label ?description {",
            "      ?o rdfs:label \"" + label.replace("\"", "") + "\"@en .",
            "      ?o rdfs:label ?label .",
            "      FILTER (LANG(?label) = 'en')",
            "      OPTIONAL {",
            "        ?o schema:description ?description",
            "        FILTER (LANG(?description) = 'en')",
            "      }",
            "    } LIMIT " + limit,
            "  }",
            "}");
    }
    ...
}
```



WIKIDATA

Jena

DBpedia

OM2

+ Music Ontology <http://purl.org/ontology/mo/>
+ PROV Ontology <http://www.w3.org/ns/prov/>

Martin Weise

NEW UI

NEW SERVICE



| 3 results | | | + DATABASE |
| Column | | | |

The following fields are AND connected and depend on the type above.

| ID | Name | Internal Name |

Column Type
| DECIMAL(size, d) | Is Null Allowed | Is Primary Key |

If you select a concept and unit , you can search across columns regardless of their unit of measurement.

Concept
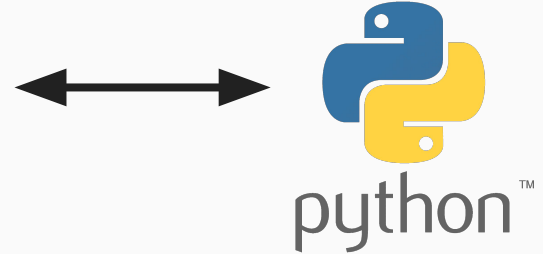| temperature | Unit | Start Value | End Value |

SEARCH

Sotirios Tsepelakis

Geoffrey Karnbach

OpenSearch client for
structured search

Authentication hidden
from UI

Enough data for
meaningful queries

Scales well

Scales well

Indexed metadata

Geoffrey Karnbach

Faceted browsing

Sotirios Tsepelakis

Faceted browsing

Sotirios Tsepelakis

Faceted browsing

Sotirios Tsepelakis

Faceted browsing

Sotirios Tsepelakis

Resolve external identifiers to increase metadata quality

Martin Weise

# WP2 - Conversion between Units of Measurement

## GOALS

Extend the metadata stored for each column that contains measurements to also allow the collection of metadata to enhance the conversation between units of measurement

Extend the search further to allow unit-independent search within the Ontology of units of measurements

## Old

# WP2 - Conversion between Units of Measurement

## DELIVERIES

Old

Conversation between Units of Measurement possible

Collect semantic metadata

http://www.wikidata.org/entity/Q11466

http://www.ontology-of-units-of-measure.org/resource/om-2/degreeFahrenheit

Martin Weise

Recommend semantic metadata based on column name

Martin Weise

Use external metadata to increase quality

Martin Weise

MariaDB

OpenSearch

```
min:     0
max:    12
mean:    4.5
median: 5
stdDev: 3.1
```

```
min:     0
max:    12
mean:    4.5
median: 5
stdDev: 3.1
```

python™

Collect statistical properties for each column

Nikola Lukic

```python
for unit_uri in unit_uris:
    gte = t1
    lte = t2
    if unit_uri != field_value_pairs["unit.uri"]:
        target_unit = unit_uri_to_unit(unit_uri)
        if not Unit.can_convert(base_unit, target_unit):
            continue
        gte = om(t1, base_unit).convert(target_unit)
        lte = om(t2, base_unit).convert(target_unit)
        searches.append({"index": "column"})
        searches.append({
            "query": {
                "bool": {
                    "must": [
                        {
                            "match": {
                                "concept.uri": {
                                    "query": field_value_pairs["concept.uri"]
                                }
                            }
                        },
                        {
                            "range": {
                                "val_min": {
                                    "gte": gte
                                }
                            }
                        },
                        {
                            "range": {
                                "val_max": {
                                    "lte": lte
                                }
                            }
                        },
                        {
                            "match": {
                                "unit.uri": {
                                    "query": unit_uri
                                }
                            }
                        }
                    ]
                }
            }
        })
```

Convert statistical properties not in the target unit

Search unit-independent between [t1, t2]

omlib.py



Max Spannring

# EXAMPLE

| | id | date | cloud_cover | sunshine | global_radiation | max_temp | mean_temp | min_temp | precipitation | pressure | snow_depth |
|---|---|---|---|---|---|---|---|---|---|---|---|
| ☐ | 1 | 19790101 | 2 | 7 | 52 | 2.3 | -4.1 | -7.5 | 0.4 | 101900 | 9 |
| ☐ | 2 | 19790102 | 6 | 1.7 | 27 | 1.6 | -2.6 | -7.5 | 0 | 102530 | 8 |
| ☐ | 3 | 19790103 | 5 | 0 | 13 | 1.3 | -2.8 | -7.2 | 0 | 102050 | 4 |
| ☐ | 4 | 19790104 | 8 | 0 | 13 | -0.3 | -2.6 | -6.5 | 0 | 100840 | 2 |
| ☐ | 5 | 19790105 | 6 | 2 | 29 | 5.6 | -0.8 | -1.4 | 0 | 102250 | 1 |
| ☐ | 6 | 19790106 | 5 | 3.8 | 39 | 8.3 | -0.5 | -6.6 | 0.7 | 102780 | 1 |
| ☐ | 7 | 19790107 | 8 | 0 | 13 | 8.5 | 1.5 | -5.3 | 5.2 | 102520 | 0 |
| ☐ | 8 | | | | | 5.8 | 6.9 | 5.3 | 0.8 | 101870 | 0 |
| ☐ | 9 | | | | | 5.2 | 3.7 | 1.6 | 7.2 | 101170 | 0 |
| ☐ | 10 | | | | | 4.9 | 3.3 | 1.4 | 2.1 | 98700 | 0 |

```
min:    -5.2
max:     12
mean:    4.5
median:  5
stdDev:  3.1
```

Rows per page: 10 ▾   1-10 of 15341   < >

## Temperature / °C

| | id | date | cloud_cover | sunshine | global_radiation | max_temp | mean_temp | min_temp | precipitation | pressure | snow_depth |
|---|---|---|---|---|---|---|---|---|---|---|---|
| ☐ | 1 | 19790101 | 2 | 7 | 52 | 2.3 | -4.1 | -7.5 | 0.4 | 101900 | 9 |
| ☐ | 2 | 19790102 | 6 | 1.7 | 27 | 1.6 | -2.6 | -7.5 | 0 | 102530 | 8 |
| ☐ | 3 | 19790103 | 5 | 0 | 13 | 1.3 | -2.8 | -7.2 | 0 | 102050 | 4 |
| ☐ | 4 | 19790104 | 8 | 0 | 13 | -0.3 | -2.6 | -6.5 | 0 | 100840 | 2 |
| ☐ | 5 | 19790105 | 6 | 2 | 29 | 5.6 | -0.8 | -1.4 | 0 | 102250 | 1 |
| ☐ | 6 | 19790106 | 5 | 3.8 | 39 | 8.3 | -0.5 | -6.6 | 0.7 | 102780 | 1 |
| ☐ | 7 | 19790107 | 8 | 0 | 13 | 8.5 | 1.5 | -5.3 | 5.2 | 102520 | 0 |
| ☐ | 8 | | | | | 5.8 | 6.9 | 5.3 | 0.8 | 101870 | 0 |
| ☐ | 9 | | | | | 5.2 | 3.7 | 1.6 | 7.2 | 101170 | 0 |
| ☐ | 10 | | | | | 4.9 | 3.3 | 1.4 | 2.1 | 98700 | 0 |

```
min:     36
max:     50
mean:    39.5
median:  38
stdDev:  2.8
```

Rows per page: 10 ▾   1-10 of 15341   < >

## Temperature / °F

"Give me tables with concept
Temperature and °C between [0, 10]"

Max Spannring

# EXAMPLE

omlib.py

"convert °F from (°C × 9/5) + 32"

"Do you contain values 0–10?"

```
min:    -5.2
max:    12
mean:    4.5
median: 5
stdDev: 3.1
```

Temperature / °C

"Do you contain values 32–50?"

```
min:    31
max:    50
mean:    39.5
median: 38
stdDev: 2.8
```

Temperature / °F

"Give me tables with concept
Temperature and °C between [0, 10]"

Max Spannring

# Helm charts for Kubernetes deployments

Generic open-source cloud deployment (for any cloud)

`oci://dbrepo.azurecr.io/helm/dbrepo-core`

TU Wien flavored cloud deployment (for any cloud)

`oci://dbrepo.azurecr.io/helm/dbrepo-tuwien`

+ SSO proxy
+ Prometheus monitoring
+ Grafana dashboard

Martin Weise

## VISION

Collaboration of TUWRD and DBRepo across VREs:

1. Get database snapshots from TUWRD (or other file-based repository)
2. Add semantic context for machine-understandability and explore tabular data in VRE
3. Seamlessly store finished research artifacts in TUWRD (e.g. plots) and data that produced these plots as queries/views in DBRepo