

# Graphbasierte Modellierung und Repräsentation eines Forschungsdateninventars

Daniel Jettka

INEL-Projekt (Indigenous Northern Eurasian Languages)

Akademie der Wissenschaften in Hamburg

## Grammatiken, Korpora und Sprachtechnologie für indigene nordeurasische Sprachen

Leitung: Prof. Dr. Beáta Wagner-Nagy  
 Antragssteller: Prof. Dr. Beáta Wagner-Nagy, Dr. Michael Rießler,  
 Hanna Hedeland, Timm Lehmberg  
 Wissenschaftliche Mitarbeiter: Dr. Alexandre Arkhipov (Forschungskordinator),  
 Timm Lehmberg (Technischer Koordinator),  
 Dr. Maria Brykina, Chris Lasse Däbritz, Anne Ferger,  
 Daniel Jettka, Tiina Klooster, Svetlana Orlova

**Erschließung sprachlicher Ressourcen indigener Sprachen  
 sowie ihre Bereitstellung über eine digitale Forschungsinfrastruktur**

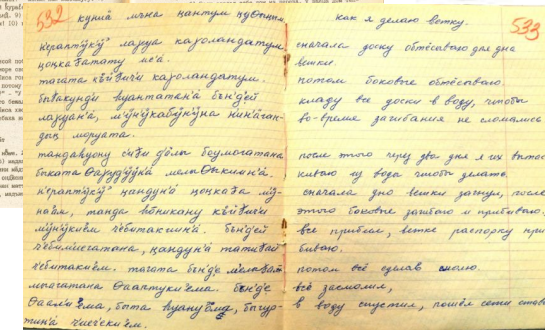
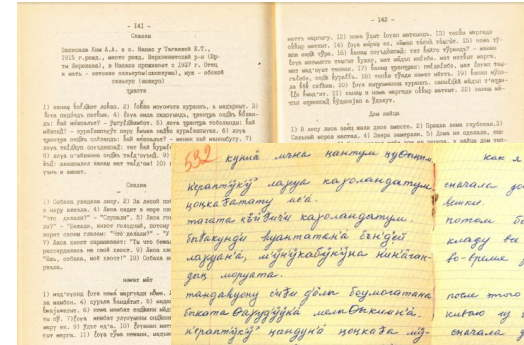


## Große Mengen von Sprachdaten

- in Archiven (nicht alles öffentlich)
- in Publikationen

## Daten/Texte mit und ohne Audio

## Archive in Tomsk, St. Petersburg, Moskau, Hamburg, Tartu, ...



# Inventur

Ziel: Inventar der vorhandenen Daten

erweiterbar

durchsuchbar

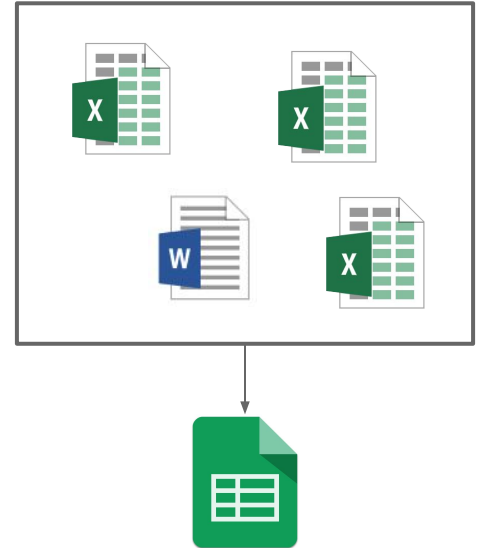
web-basiert

Problem:

Mehrere fortlaufend ergänzte Inventurlisten (versch. Struktur und Formate)

Lösung:

Inventur und Zusammenführung von Inventurlisten  
(Google Spreadsheet für Übergangszeit)





## Vorteile Google Spreadsheets



1. Einfache Kollaboration
  2. Höheres Maß an Konsistenz durch:
    - a. Auswahllisten (auf Basis vorhandener Werte/Listen)
    - b. Validierung (z.B. Dopplungen)
    - c. Dynamisch generierte Inhalte (z.B. volle Namen, IDs)
    - d. Vglw. einfache Erweiterung der Oberfläche/Funktionalität
- Erste Schritte der Datenmodellierung (Entitäten, Kategorien, Referenzen)

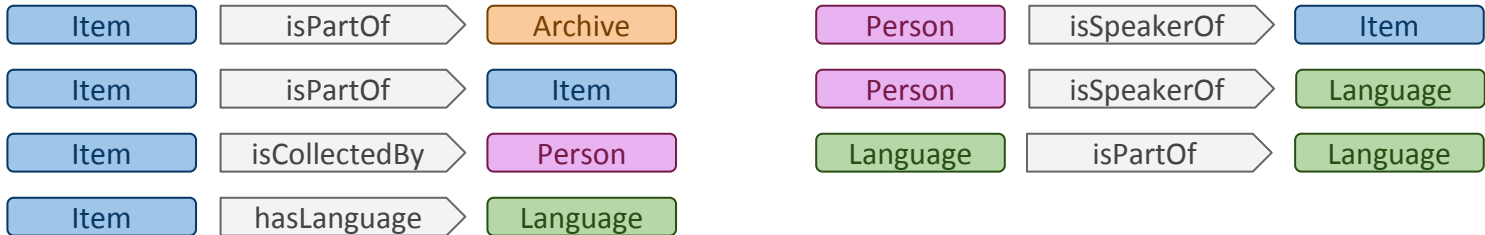
# Inventur

## Datenmodellierung:

- Identifikation/Definition von Entitäten
- Identifikation/Definition relevanter Eigenschaften der Entitäten
- Identifikation/Definition von Relationen zw. Den Entitäten
- Validierung/Kuration existierender Beschreibungen

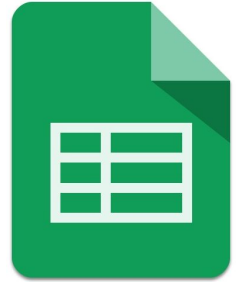


### Relationen zwischen Entitäten:



## Nachteile Google Spreadsheets

1. Beschränkungen des tabellarischen Datenmodells
2. Beschränkungen für Generierung dynamischer Inhalte
3. Datenschutz



- Versch. Datenbanktechnologien möglich
- Evaluation der Graphdatenbank Neo4j



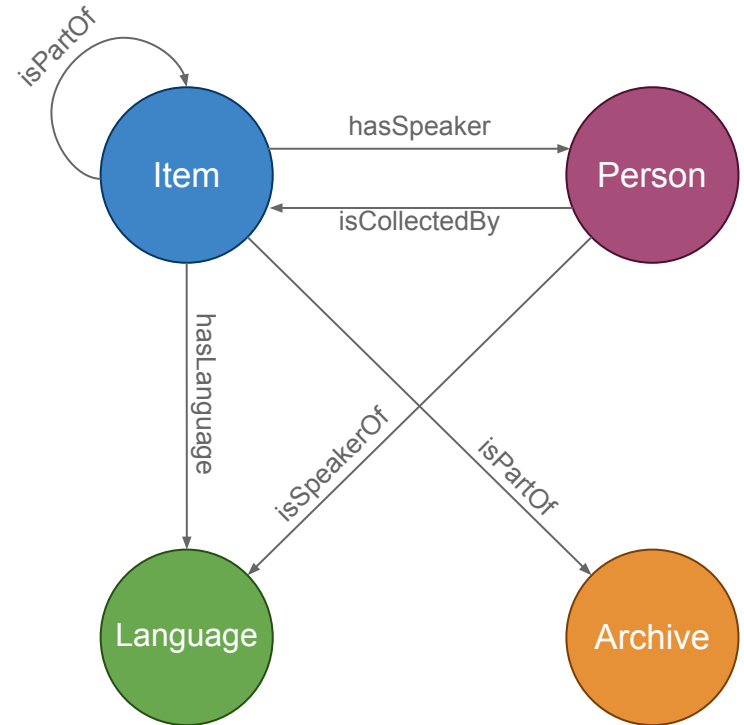


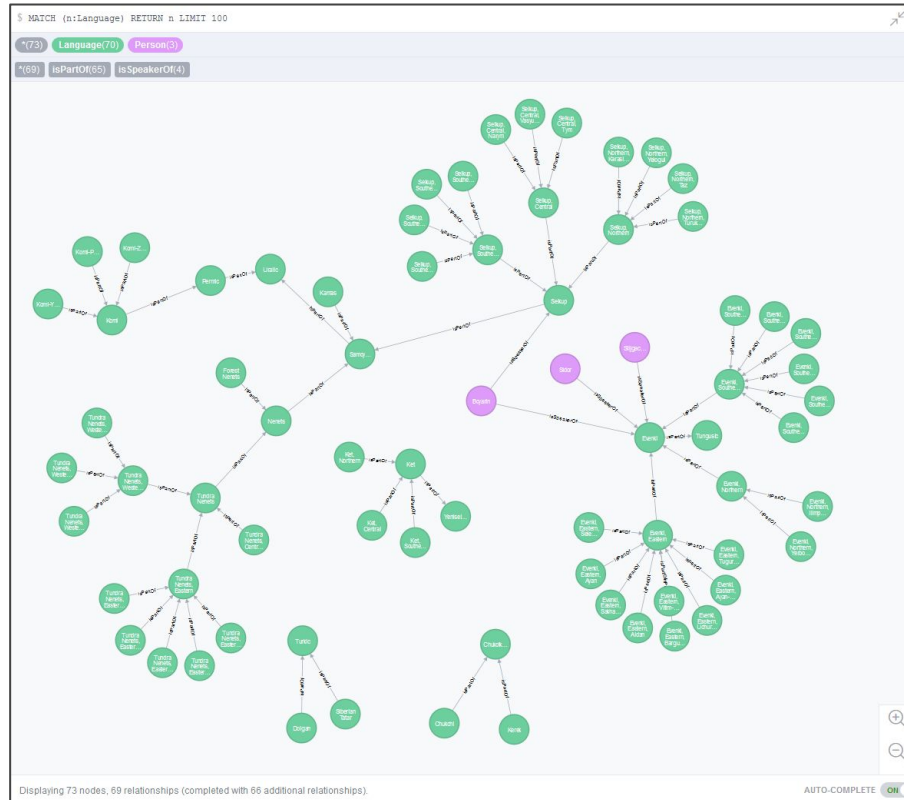
- Community Edition frei verwendbar
  - Einfache Installation
  - Schema-los; simple Datentypen
  - REST Endpoint
  - Webschnittstelle für Suche und kommando-basierte Pflege von Daten
- Transformation existierender Beschreibungen in Cypher-Ausdrücke (JSON)
- Import in Neo4j-Datenbank



## Vorteile:

- Flexible Datenmodellierung, Repräsentation, Abfrage







```

1 MATCH p = (l:Language)-[:hasLanguage]-(i:Item)-[:isPartOf*]->(a:Archive)
2 WITH a, l, collect(DISTINCT l.name_eng) AS languages
3 RETURN a.name_eng AS archive, collect(languages) AS langs, count(languages) AS count ORDER BY count desc

```

\$ MATCH p = (l:Language)-[:hasLanguage]-(i:Item)-[:isPartOf\*]->(a:Archive) WITH a, l, collect(DISTINCT l.name\_eng) AS languages RETURN a.name\_eng AS archive, collect(languages) ...

archive	langs	count
Pushkin House (IRLI)	[[Selkup], [Kerek], [Chukchi], [Komi], [Evenki], [Ket], [Nenets], [Dolgan]]	8
Dulson Archive (TSPU)	[[Evenki], [Selkup], [Dolgan]]	3

Returned 2 rows in 88 ms.

**Namen und Anzahl von Languages verbunden mit Items in versch. Archives**

```

MATCH p = (l:Language)-[:hasLanguage]-(i:Item)-[:isPartOf*]->(a:Archive)
WITH a, l, collect(DISTINCT l.name_eng) AS languages
RETURN a.name_eng AS archive, collect(languages) AS langs, count(languages) AS count ORDER BY count desc

```

\$ MATCH p = (n:Person)-->() WHERE n.lastname\_eng='Dibikova' RE...  
 \*(4) Language(1) Person(1)  
 \*(5) hasLanguage(2) isSpeakerOf(3)

\$ MATCH (n:Language) RETURN n LIMIT 100  
 \*(73) Language(70) Person(3)  
 isSpeakerOf(4)

Person <id>: 3841 lastname\_rus: Дибикова identifier: DVP  
 lastname\_eng: Dibikova patronymicname\_rus: П. firstname\_eng: V.  
 patronymicname\_eng: P. firstname\_rus: B.

Person <id>: 4245 lastname\_rus: Боярин lastname\_eng: Boyarin patronymicname\_rus: A. firstname\_eng: D. patronymicname\_eng: A. id: BDA

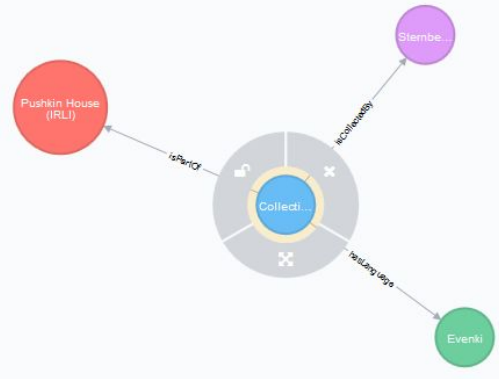
# Collections

```

    $ MATCH p = (n:Item)-->() WHERE n.title_eng='Collection 20' RE...
    
```

\*(4) Archive(1) Item(1) Language(1) Person(1)

\*(3) hasLanguage(1) isCollectedBy(1) isPartOf(1)



```

    graph LR
        A((Pushkin House (IRL))) -- isPartOf --> B((Collection 20))
        B -- isCollectedBy --> C((Sterbe...))
        B -- hasLanguage --> D((Evenki))
    
```

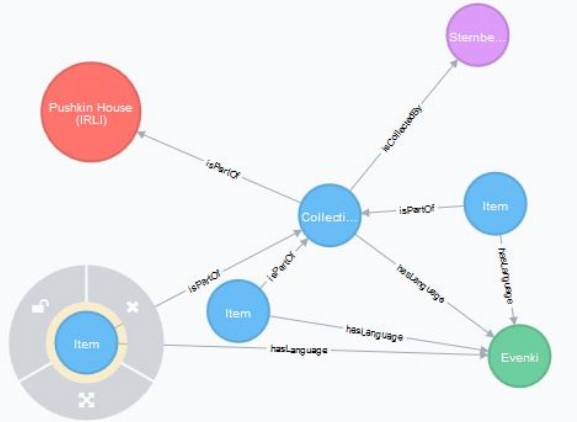
**Item** <id>: 4092 title\_eng: Collection 20 identifier: ph-coll-20  
**description\_rus:**  
 Записи выполнены летом 1910 г. на восковых валиках (фонографических цилиндрах) на территории Амурского края от нанайцев, нивхов, негидальцев, эвенков. Предварительное описание экспедиции (С.М.Широкогоров)  
**type:** Collection **collection\_ref:** Coll. 20

```

    $ MATCH p = (n:Item)-->() WHERE n.title_eng='Collection 20' RE...
    
```

\*(7) Archive(1) Item(4) Language(1) Person(1)

\*(9) hasLanguage(4) isCollectedBy(1) isPartOf(4)



```

    graph LR
        A((Pushkin House (IRL))) -- isPartOf --> B((Collection 20))
        B -- isCollectedBy --> C((Sterbe...))
        B -- isPartOf --> D((Item))
        D -- isPartOf --> E((Item))
        E -- hasLanguage --> F((Evenki))
        B -- hasLanguage --> F
        E -- hasLanguage --> F
    
```

**Item** <id>: 4095 identifier: ph-coll-20-item-3 description\_rus: Песня  
**description\_deu:** Lied **collection\_ref:** Coll. 20 **type:** Item **description\_eng:** Song



```

1 MATCH (l:Language)-[:hasLanguage]-(i:Item)-[:isCollectedBy]->(p:Person)
2 RETURN p.firstname_eng+" "+p.patronymicname_eng+" "+p.lastname_eng AS researcher, collect(DISTINCT l.name_eng) AS langs, count(i) AS
   item_count
3 ORDER BY item_count desc

```

\$ MATCH (l:Language)-[:hasLanguage]-(i:Item)-[:isCollectedBy]->(p:Person) RETURN p.firstname\_eng+" "+p.patronymicname\_eng+" "+p.lastname\_eng AS researcher, collect(DISTINCT l.name\_eng) AS langs, count(i) AS item\_count

researcher	langs	item_count
Andrej Petrovich Duljzon	[Evenki, Selkup]	253
N. P. Maksimova	[Selkup]	126
Z. P. Demjyanenko	[Dolgan, Selkup]	67
V. V. Bykonya	[Selkup]	64
E. S. Kuznecova	[Selkup]	62
N. M. Voevodina	[Selkup]	58
I. A. Iljashenko	[Selkup]	56
E. G. Bekker	[Selkup]	53
N. V. Denning	[Selkup]	32
A. A. Kim	[Selkup]	31
G. K. Verner	[Selkup]	30
Yu. A. Moreva	[Selkup]	29
N. P. Beljtyukova	[Selkup]	28
T. M. Kosheverova	[Selkup]	19
O. A. Osipova	[Selkup]	17
L. A. Alitkina	[Selkup]	13

Returned 55 rows in 253 ms.

Anzahl von Items und versch. verbundenen Languages, die von jedem Forscher erhoben wurden

```

MATCH (l:Language)-[:hasLanguage]-(i:Item)-[:isCollectedBy]->(p:Person)
RETURN p.firstname_eng+" "+p.patronymicname_eng+" "+p.lastname_eng, collect(DISTINCT l.name_eng) AS langs, count(i) AS item_count
ORDER BY item_count desc

```







## Nachteil:

- Beschränkte Zugänglichkeit für Linguisten  
(Nicht-Techniker)

## Lösung:

- Evaluation vorhandener GUIs (z.B. structr, Linkurious, popoto.js)  
→ entweder zu teuer oder nur teilweise einsetzbar
- Implementierung einer eigenen GUI



## Prototyp I: HTML-Formulare

- Formulare mit Bootstrap (HTML5, CSS, JavaScript)
  - Validierung, Auto-completion, Felder mit Mehrfachauswahl
  - Verbindung mit Neo4j REST Endpoint
- ➔ Nicht nutzerfreundlich (Tabellenfunktionen gewünscht)

The image shows a web application interface for managing research items. The main view is titled 'Manage research items' and has tabs for 'Item', 'Archive', 'Language', and 'Person'. It contains several sections:

- Related research items:** Includes an 'Archive' button and a dropdown menu showing 'Dulson Archive (TSPU) | Дульзонский...'. Below are 'Parent item/s' and 'Derivation of' dropdowns.
- Related persons:** Includes 'Collector/s' and 'Author/s' dropdowns. The 'Author/s' dropdown is open, showing a list of names with 'AEE: E. Aksenova | E. Аксенова' selected.
- Title:** Includes dropdowns for 'eng' and 'rus' with 'Enter a...' text.
- Description:** Includes dropdowns for 'eng' and 'rus' with 'Enter a description in...' text.
- Year:** Includes a 'YYYY' input field.

The modal window for editing 'Dolgan' has the following fields:

- Language name:** Three input fields for 'eng' (Dolgan), 'rus' (Enter a language name), and 'deu' (Enter a language name).
- Type:** A dropdown menu.
- Language Codes:** Two input fields for 'ISO-639-3' (dlg) and 'Glottolog' (dolg1241).
- Relation/s:** A dropdown menu showing 'hasParent' and 'Turkic', with a green '+' button.
- Comment:** A text area containing the URL 'http://glottolog.org/resource/languoid/id/dolg1241'.
- Buttons:** 'Cancel' and 'Submit' buttons.

# GUI Prototyp II

Prototyp II:



+



DataTables

- Einlesen von DB-Inhalten in DataTables-basierte HTML-Tabelle
- Anpassung der Tabellendaten auf Client-Seite
- Schreiben aktualisierter Daten in DB durch JavaScript-Trigger
- Zugriffsschutz via Shibboleth

INEL Resource Catalogue

Items Persons Languages Archives

Manage Items

Add Item

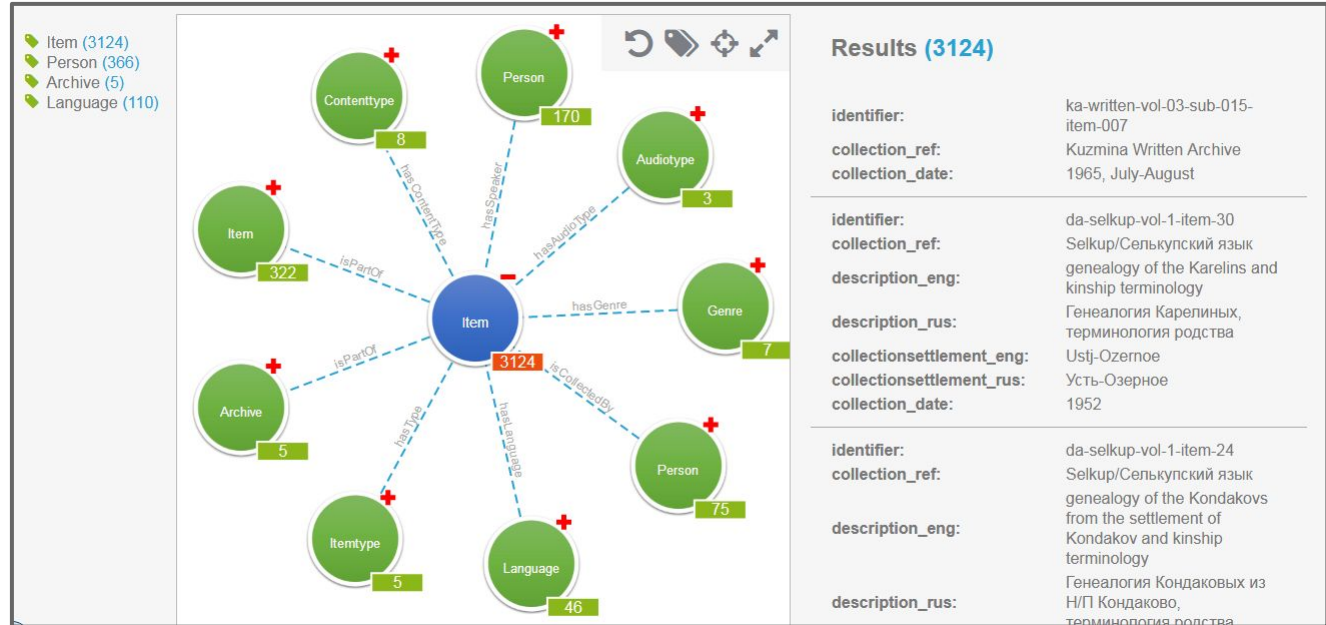
Search

Delete column (3/30)

ID	Engl Title	Engl Title	Engl Title	Type	Collection	Engl Description	Engl Description	Engl Description	Belongs to website	Contained in Item	Language(s)	Collectors	Speakers	Genre	Exact genre	Audio type	Reference audio	Audio duration
AA_1514_Brother_R	Brother	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item
AA_1514_Corpus_R	Corpus	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item
AA_1514_Girl_R	Girl	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item
AA_1514_Hair_R	Hair	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item
AA_1514_Head_R	Head	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item
AA_1514_Khan_R	Khan	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item
AA_1514_Lament_sq	Lament	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item
AA_1514_Mansar_R	Mansar	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item	Item

# Graphische Recherche in Daten

popoto.js



JavaScript-Bibliothek für Neo4j-Visualisierung

# Zusammenfassung

- Inventar von Forschungsdaten in Graphdatenbank
- Prototyp einer Webschnittstelle



- Ansatz anwendbar auf andere Daten



- Evaluation Geo-Plugin für Neo4j
- Evaluation von Triple Stores (Graphdatenbanken für RDF)
  - z.B. Fuseki, Blazegraph
- (zusätzl.) Verwendung existierender Ontologien
- Webschnittstelle mit anderen Data Stores verbinden
- Vernetzung von Datenbeständen