

Generierung angepasster RDF-Dumps von Wikidata

Bachelorverteidigung

Benno Fünfstück

15. Januar 2020

Betreuer: Prof. Dr. Markus Krötzsch
Wissensbasierte Systeme
TU Dresden



WIKIDATA

„Die freie Wissensdatenbank mit 74.018.298 Datensätzen, die jeder bearbeiten kann.“

label

item identifier

Douglas Adams (Q42)

britischer Schriftsteller (1952-2001) ----- description

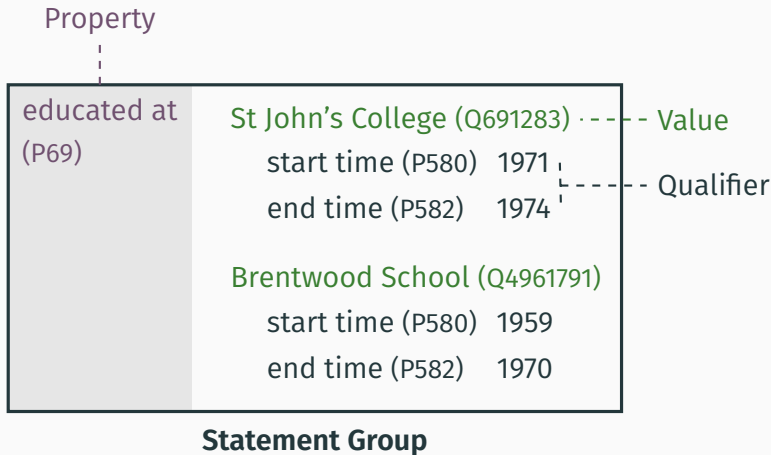
Douglas Noël Adams | Douglas Noel Adams ----- aliases

► [In weiteren Sprachen](#)

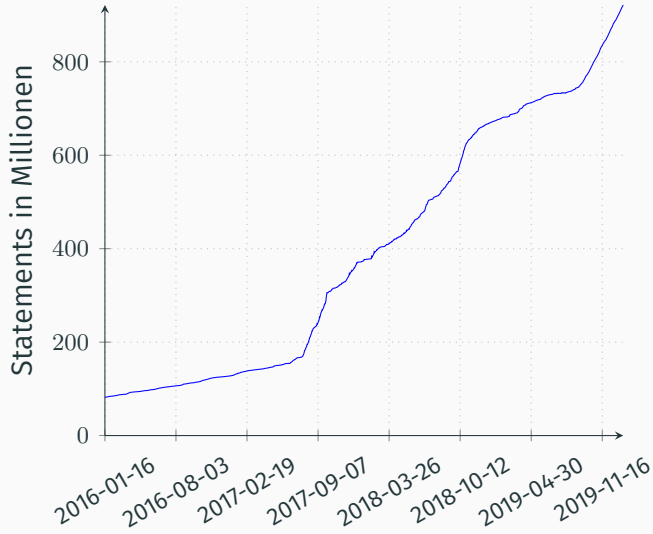
Statements

Sitelinks

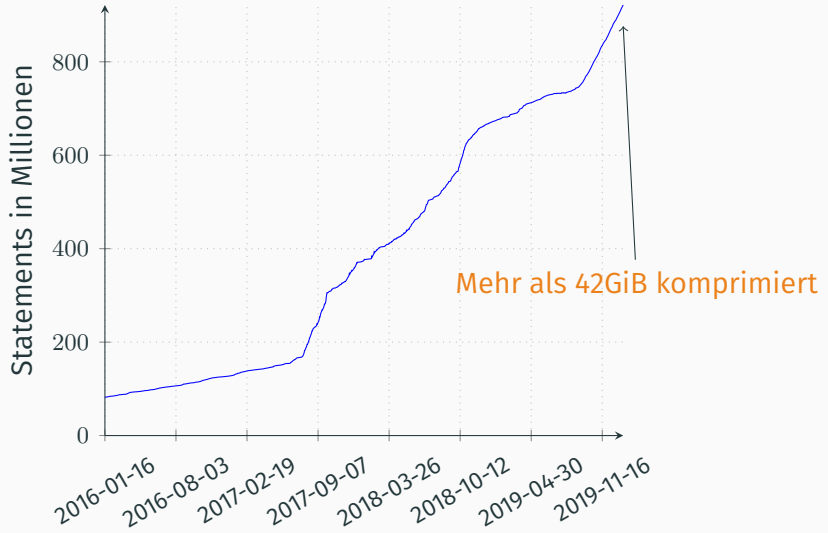
Wikidata: Statements



Wikidata: Größe

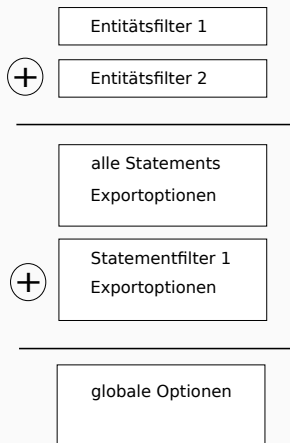


Wikidata: Größe



Idee: Tool zum Filtern der Daten

- Korrektheit, Vollständigkeit
- Filterung, Archivierung, Suche
- Statistiken, Fortschritt, Parallelverarbeitung
- Nachvollziehbarkeit, Aktualität der Daten



Format der Dumps: RDF

Wikidata als RDF (Resource Description Framework)

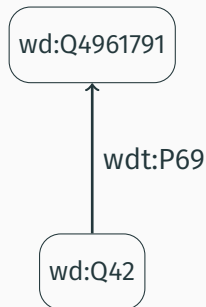
Beispiel

```
wd:Q42 rdfs:label "Douglas Adams"@en .  
wd:Q42 rdfs:label "Douglas Adams"@de .
```

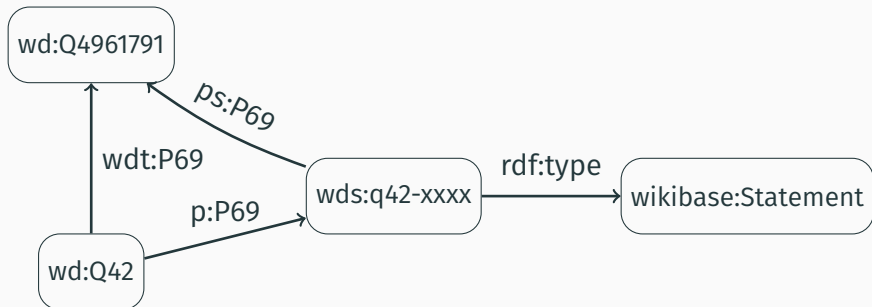
einfache Darstellung: ein Tripel für jedes Statement

```
wd:Q42 wdt:P69 wd:Q691283 .  
wd:Q42 wdt:P69 wd:Q4961791 .
```

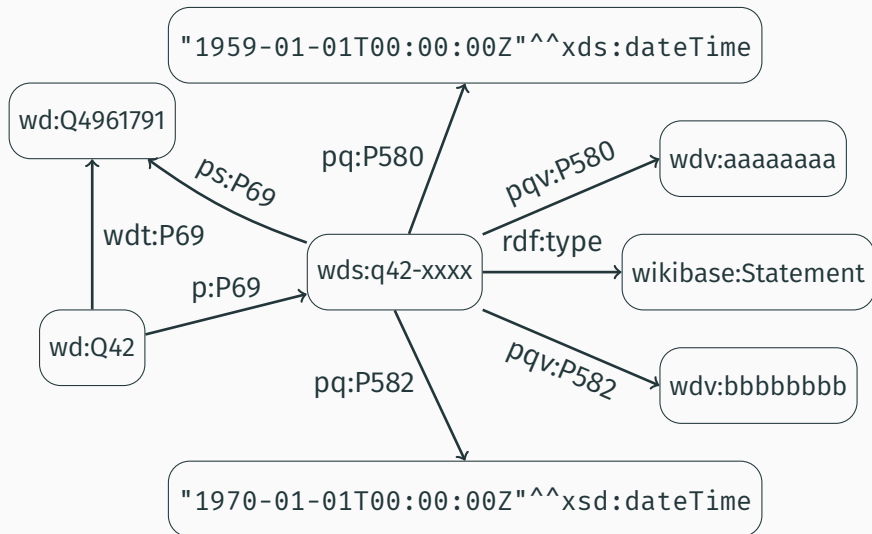
Aber: Ranks, Qualifier, komplexe Werte?



Wikidata als RDF: Reifikation



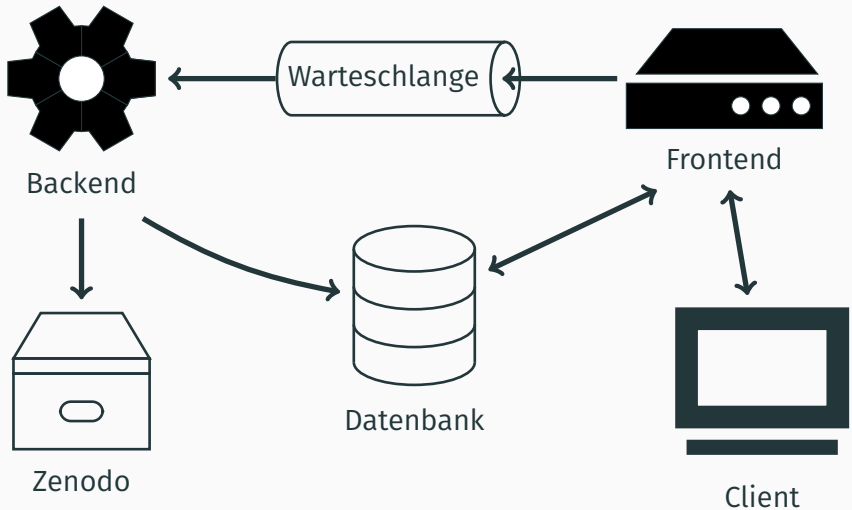
Wikidata als RDF: Reifikation

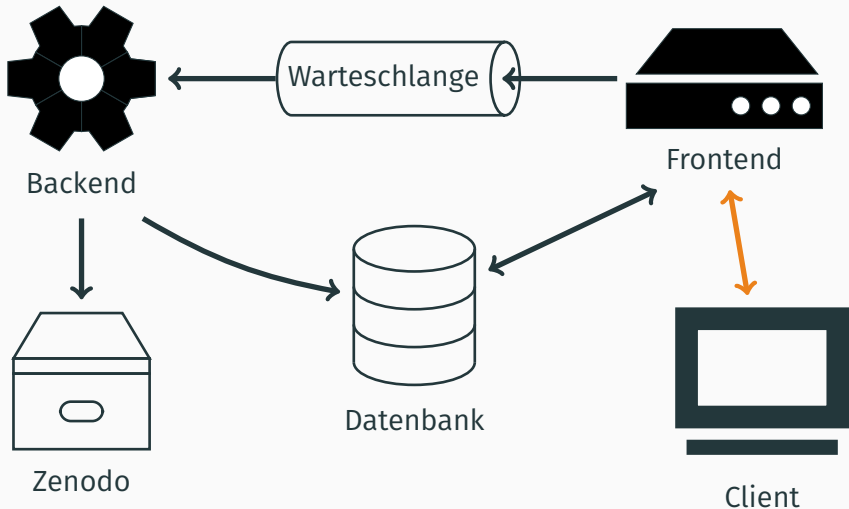


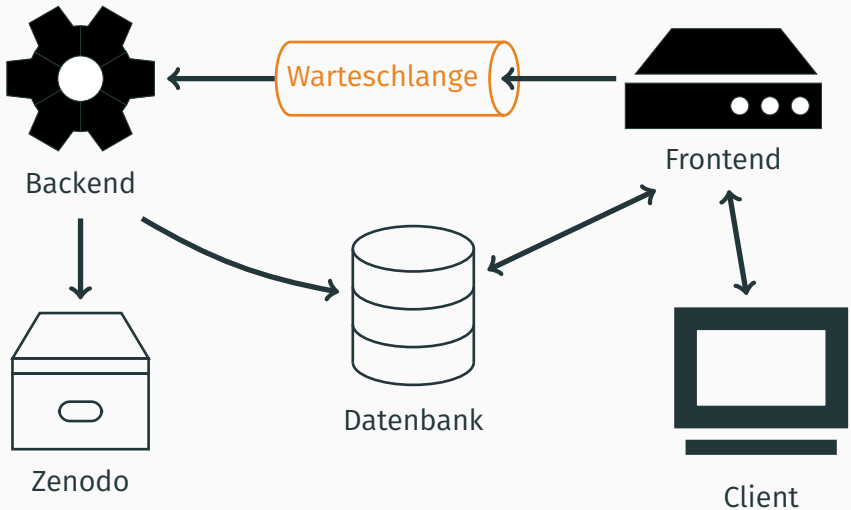
Umsetzung

	Eigener Index	SPARQL basiert	Batchverarbeitung
Geschwindigkeit	schnell	schnell	langsam
Flexibilität	gering	mittel	hoch
Implementierung	aufwändig	mittel	einfach

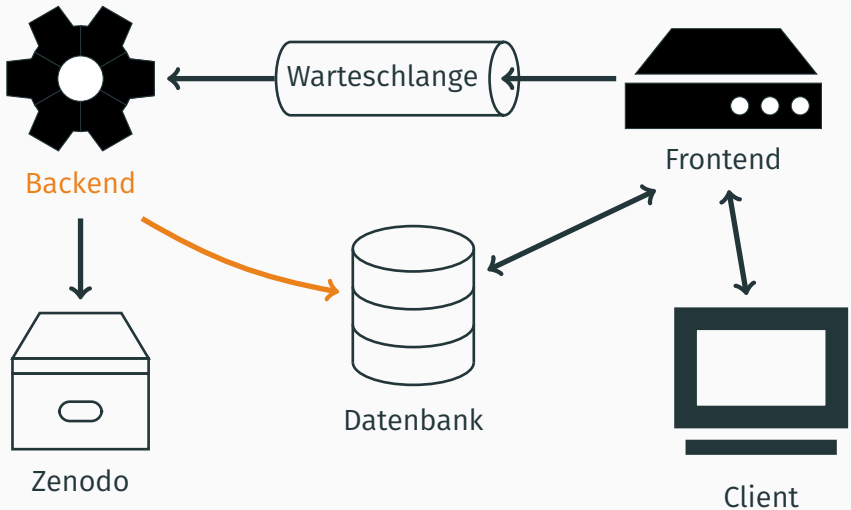
System

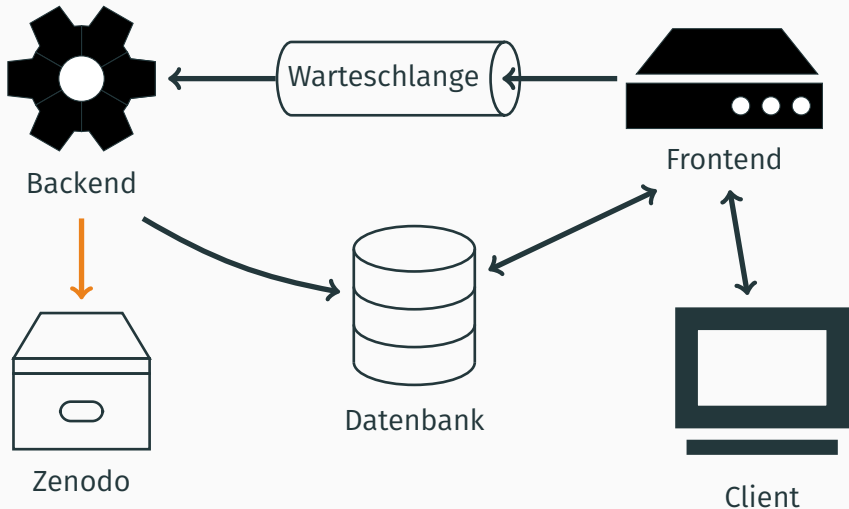






System





Evaluation

Filter entities

Choose entities to include in the dump. An entity is included if it matches at least one filter.

No filters added. All entities are included.

+ [Add basic filter](#) + [Add SPARQL query](#)

Filter statements

Choose how statements are exported. These rules are applied to all matched entities.

Default rule

This rule is applied to all matched entities

Parts to export

simple statement ☒

full statement ☐

references ☐

qualifiers ☐

+ [Add custom rule](#)

Additional settings

labels ☒

descriptions ☒

aliases ☒

sitelinks ☒

filter languages ☐

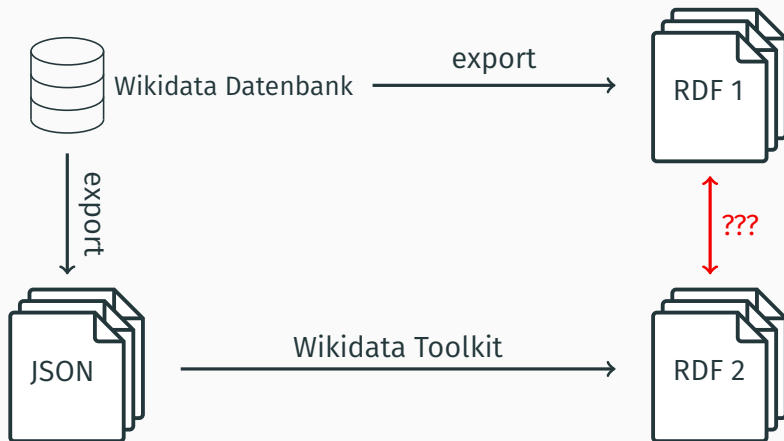
Dump metadata

Dump title

Create dump

Screenshot der
Anwendung

Die Anwendung verwendet Wikidata Toolkit zur Generierung der RDF-Daten



Viele unkritische Abweichungen:

- andere Schreibweisen
- minimal andere Struktur
- besonders bei komplexen Typen (Zeitpunkte, Orte)

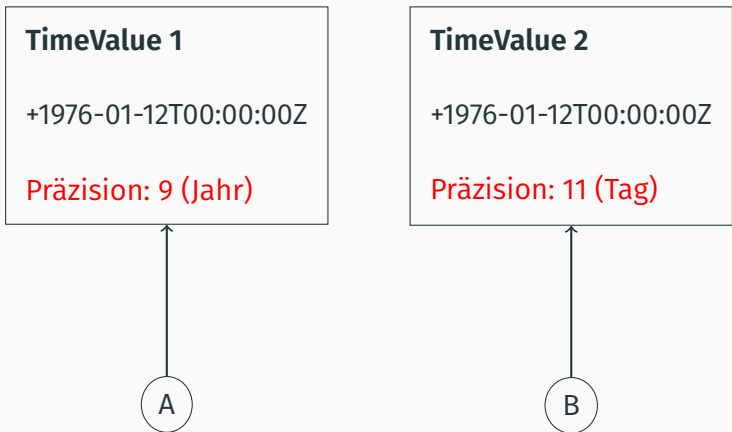
Ein paar tatsächliche Fehler:

- Vertauschung von Koordinaten bei Orten
- falsche Deduplizierung bei Value-Nodes

Einige fehlende Features in Wikidata Toolkit

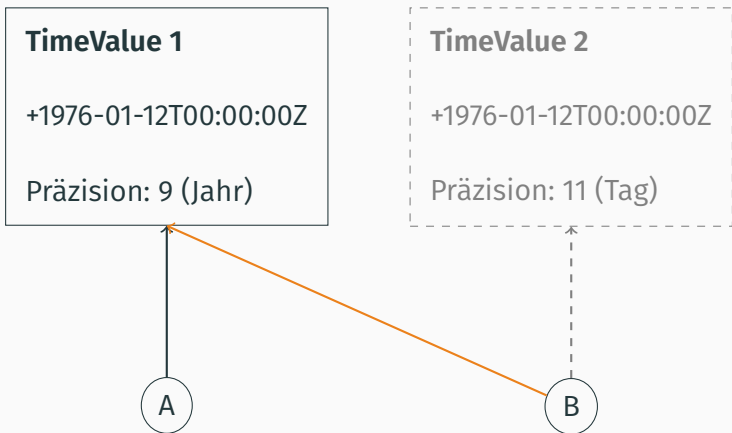
Ergebnis: falsche Deduplizierung bei Werten mit Präzisionsangabe

Für Zeit- und Ortsangaben kann die Präzision gespeichert werden



Ergebnis: falsche Deduplizierung bei Werten mit Präzisionsangabe

Für Zeit- und Ortsangaben kann die Präzision gespeichert werden



- ① erste Version der Anwendung implementiert
- ② RDF-Export von Wikidata Toolkit verbessert
- ③ viele weitere Features denkbar: Lexeme, Parallelisierung, mehr Filter, ...

- ① erste Version der Anwendung implementiert
- ② RDF-Export von Wikidata Toolkit verbessert
- ③ viele weitere Features denkbar: Lexeme, Parallelisierung, mehr Filter, ...

Fragen?

