

KONZEPTION UND UMSETZUNG EINER DSL ZUR INFORMATIONSFUSION AUF VERTEILTEN HETEROGENEN GRAPHEN

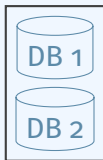
ALEXANDER KERN

BTW 2019

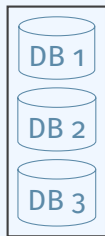
05.03.2019

INFORMATIONSFUSION

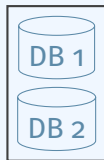
Sales



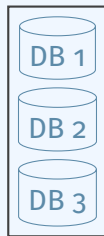
Public Relations

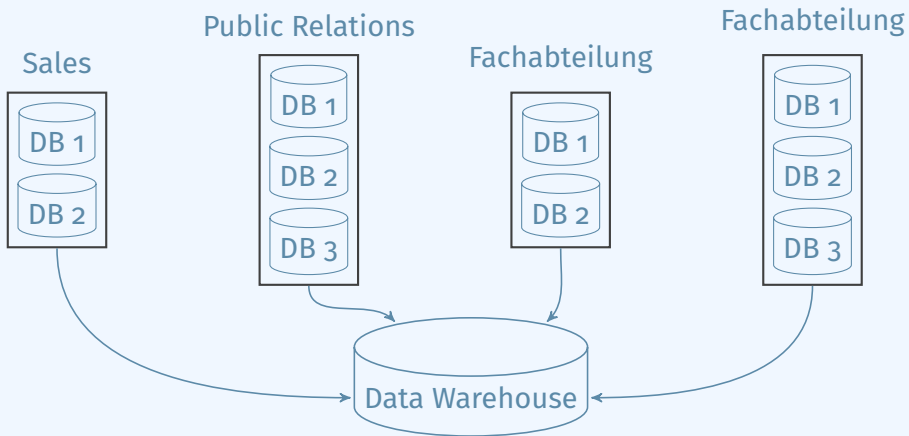


Fachabteilung



Fachabteilung





- Verknüpfung verschiedener Datenbanken
- Aggregation von Suchergebnissen
- Sensormesswerte

- Neue Zusammenhänge
- Wissen präzisieren
- Fehler finden
- Lücken füllen

- Entity Resolution
- Duplikatserkennung
- Schemaintegration
- Informationsfusion

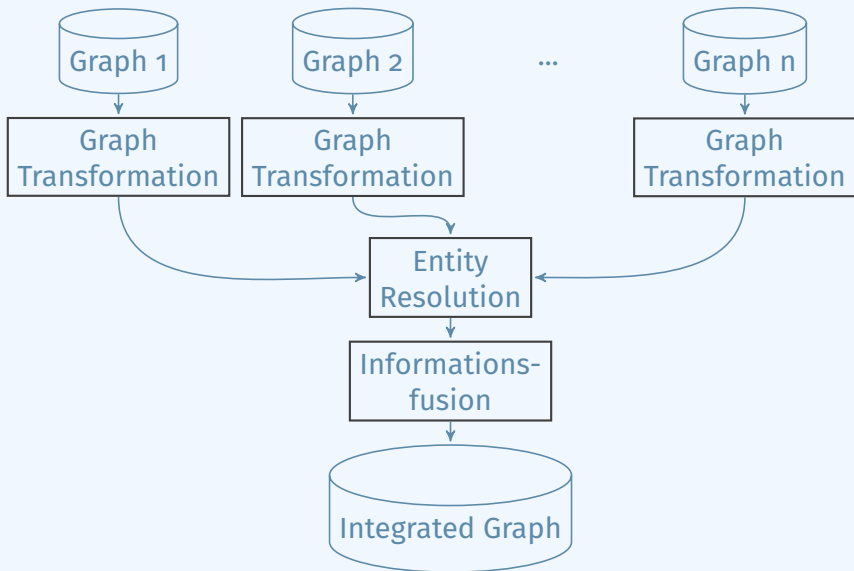
- Entity Resolution
- Deduplikation
- Schemaintegration
- Informationsfusion

Modellschema	Adresse	Straße, Nummer, PLZ, Ort
Inhaltliche Fehler	name=Alexander Augustusplatz 10 Augustusplatz 10	name=Kern Ritterstrasse 9-13 Augsutusplatz 10
Datenschema	2003-01-07	7. Jan. 2018

- Apache Flink-basiertes Framework
- Verteilte deklarative Graphanalysen
- Extended Property Graph Model



PROZESS IN GRADOOP



- Unterstützung heterogener Datenschema
- Beschreibung des Ausgabeschemas
- Zusammenhang mit Eingabedaten
- Bestimmung der Konfliktlösungsstrategien
→ Transformationsregeln

DOMÄNENSPEZIFISCHE SPRACHE

- Declarative Data Cleaning: Language, Model, and Algorithms, Helena Galhardas et. al., 2001
- SQL-ähnliche Sprache zur Beschreibung von Informationsintegration
- Datenaufbereitung, Clustering und Merging

```
CREATE MERGING MergeAuthors
USING clusterAuthors(cluster id) ca
LET name = getLongestAuthorName(
    DirtyAuthors(ca).name)
key = generateKey()
{ SELECT key AS authorKey,
    name AS name INTO Authors g }
```

- Auswahl von (heterogenen) Eingabeattributen
- Strategien mit Optionen
- Kurzformen:
 - ▶ homogene Eingabeattribute
 - ▶ Attribute ohne Regeln

- Ausführung mit Gradoop/Flink
- Deklarativ

Strategie	Beschreibung
Straight	Erster Nicht-Null-Wert nach Ordnung
Retain	Konkatenation
Priority	Erster Wert in Ordnung der Quellen
Newest	Aktuellster Zeitstempel
Majority	Meistauf tretender Wert
Source	Höchstes summiertes Gewicht
Property	Auswahl anhand der vorhandenen Attributwerte

Vergleiche [2]

```
CREATE MERGING
```

```
LET Author
```

```
  name = property(  
    (Source1.author.name,  
     Source2.author.full_name),  
    textual:longest),  
  authorKey = straight(key)
```

```
node Author {  
    property name {  
        Source1.author.name  
        Source2.author.full_name  
    } strategy property { textual:longest }  
  
    property authorKey { key } strategy straight  
}
```

nodes:

- label: Person
- transformations:
 - outputProperty: name
 - strategy:
 - type: property
 - options: [textual-longest]
 - properties:
 - Source1.author.name
 - Source2.author.full_name
- outputProperty: authorKey
- strategy:
 - type: straight
 - properties: [key]

- Vorkenntnisse, Bekanntheit
- Toolunterstützung
- Einarbeitungszeit, Komplexität
- Abbildung der Problemdomäne
- ...

- Vorkenntnisse
- Bekanntheit
- Einarbeitungszeit
 - SQL-basierte DSL

ERGEBNIS - GRADOOP SERVICE

scaDS COMPETENCE CENTER FOR SCALABLE DATA SERVICES AND SOLUTIONS

graph Graph Analytics Service

Navigation: [Home](#) [About](#) [Contact](#) [Feedback](#)

- Graph
- Filter
- Clustering
- Grouping
- Vertex Fusion
- Combine
- Exclude
- ExpandGraph
- WCC
- PatternMining
- RDFGraph
- Linking
- Sampling
- EdgeFusion
- Cypher
- Overlap
- SchemaGraph
- Schema Matching
- PageRank
- Output

Workflow: --

Save Save As Load Clear List Toggle

- amazon-gp-linking
- combinedTest1
- combinedTest2
- combinedTest3
- combinedTest4

Zoom is Out default

Run RunAllOnce RunDefault

Graph cities

- Download Graph
- Upload Graph
- Drop Inp graphs
- Load as Cluster Graph

Drawing Properties

Sampling: No Sampling

Threshold: 0.2

Drawing: Vvigraph

Done View as table

Vertex Size: [Slider]

Edge Opacity: [Slider]

Font Size: [Slider]

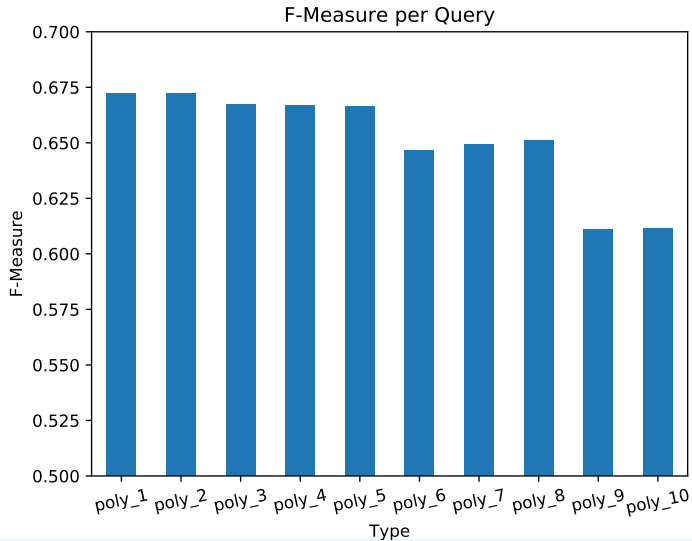
Scale Nodes: none

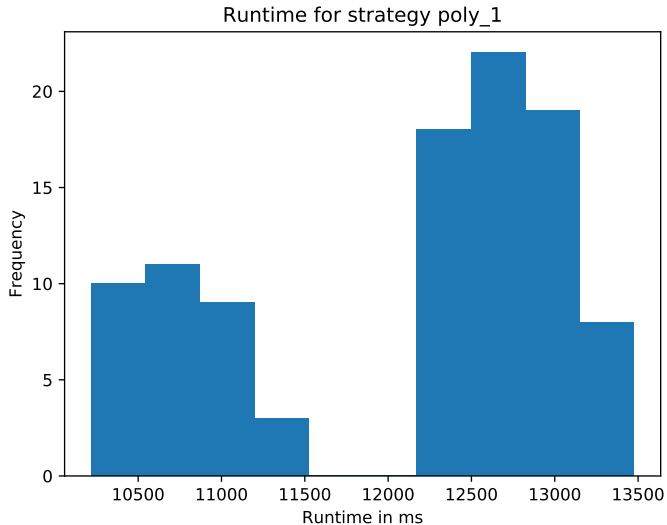
Stop Rendering

Color Chooser

EVALUATION

EVALUATION - MATCHING QUALITÄT





Parallelisierung	64	32	16	8	4	2
Mittelwert	35.29	35.02	47.12	75.49	155.7	271.06
Median	35.46	34.42	47.30	75.31	156.62	272.57
Faktor	-	0.97	1.37	1.59	2.08	1.74



FAZIT UND AUSBLICK

- Beispielhafte DSL-Implementierung
- In Code und Service nutzbar

- Evaluation mit Endnutzern
- Erweiterung für Kanten
- Weitere Konfliktlösungsstrategien
- Zusätzliche Anwendungsfälle

VIELEN DANK FÜR IHRE
AUFMERKSAMKEIT!

REFERENCES

-  HELENA GALHARDAS, DANIELA FLORESCU, AND DENNIS SHASHA.
DECLARATIVE DATA CLEANING: LANGUAGE, MODEL, AND ALGORITHMS.
In *In VLDB*, pages 371–380, 2001.
-  FELIX NAUMANN AND MATTHIAS HÄUSSLER.
DECLARATIVE DATA MERGING WITH CONFLICT RESOLUTION.
In *International Conference on Information Quality (IQ 2002)*. 2002,
pages 212–224, 2002.