

# LDIF

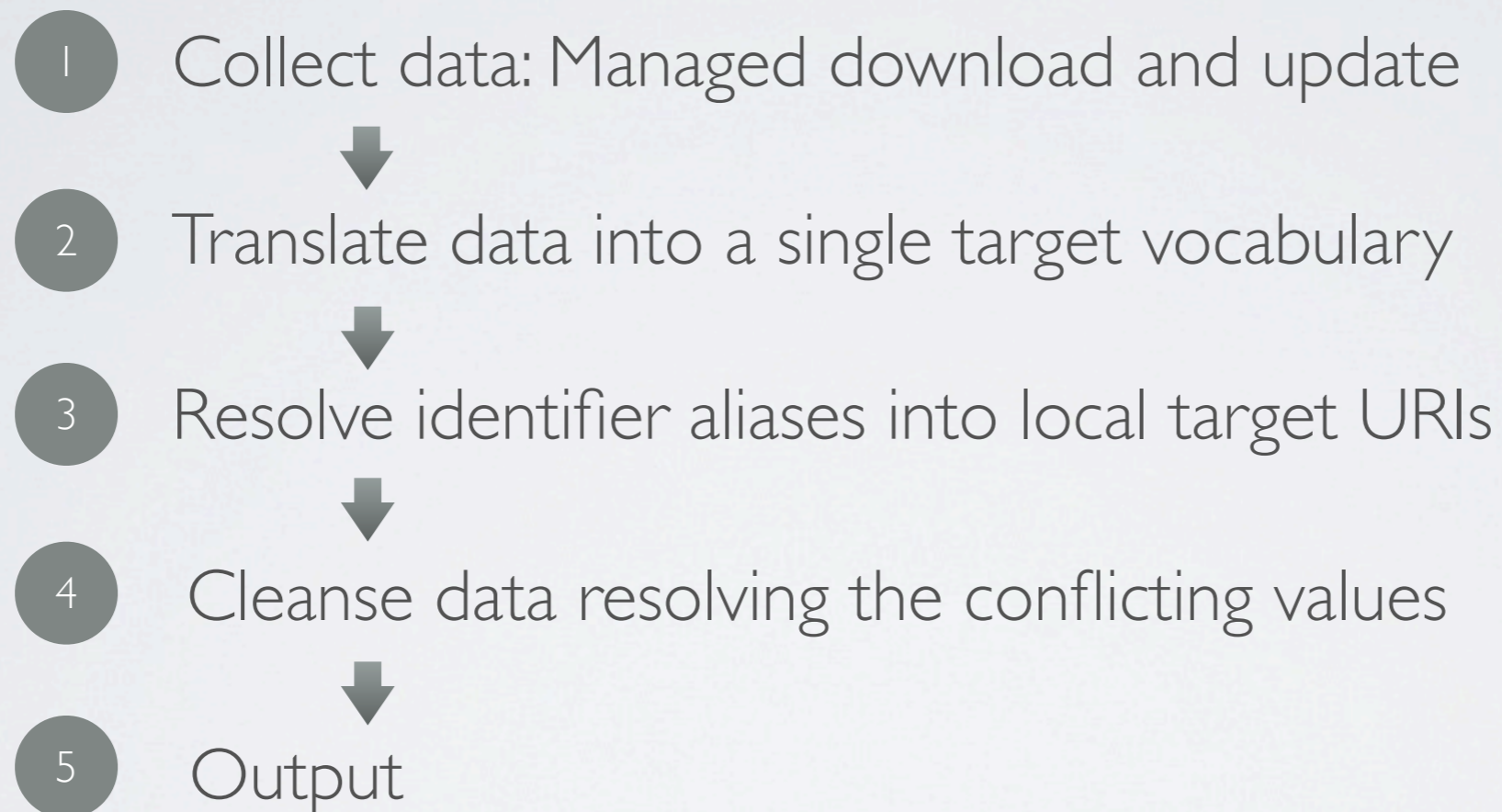
Linked Data Integration Framework

# LINKED DATA CHALLENGES

- Data sources that overlap in content may:
  - use a wide range of different RDF vocabularies
  - use different identifiers for the same real-world entity
  - provide conflicting values for the same properties
- Implications:
  - Queries are usually hand-crafted against individual sources – no different than an API
  - Improvised or manual merging of entities
- Integrating public datasets with internal databases poses the same problems

# LDIF

- LDIF homogenizes Linked Data from multiple sources into a clean, local target representation while keeping track of data provenance



- Open source (Apache License, Version 2.0)
- Collaboration between Freie Universität Berlin and mes|semantics

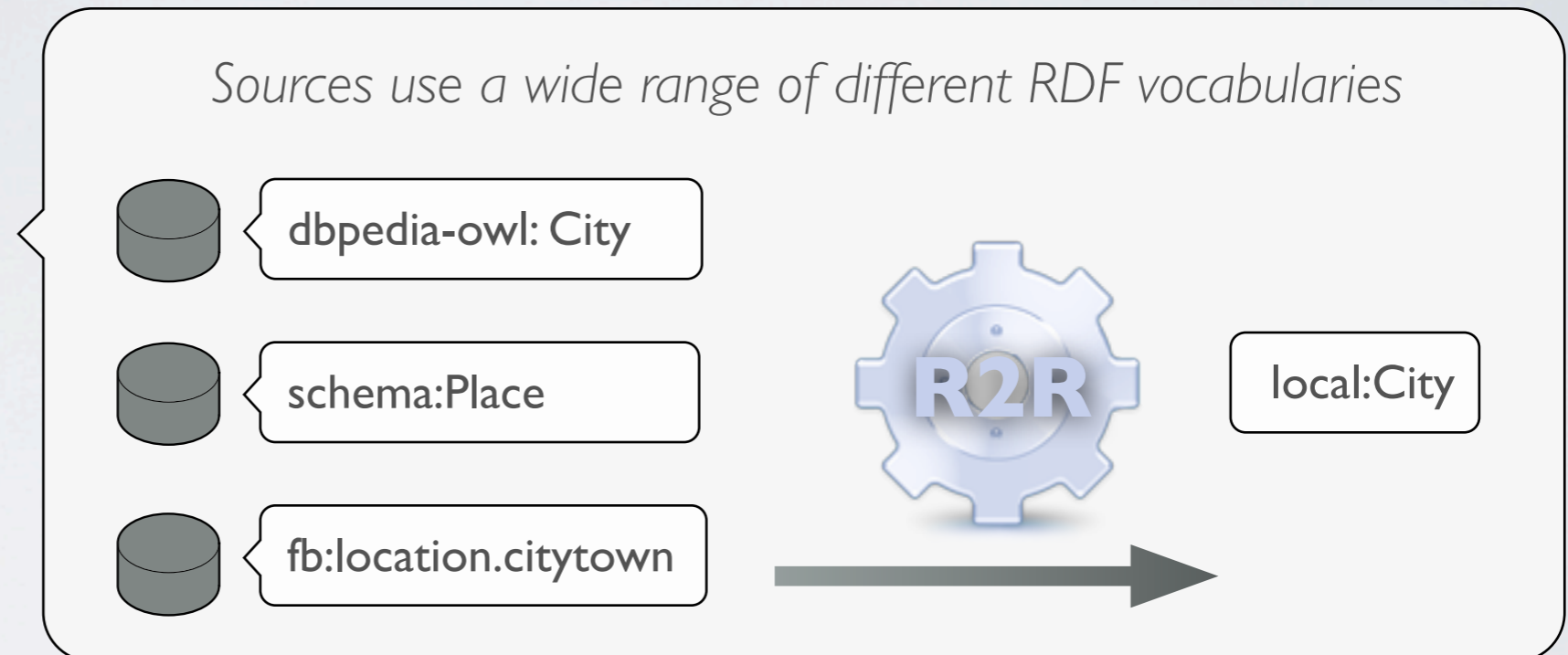
# LDIF PIPELINE



Supported data sources:

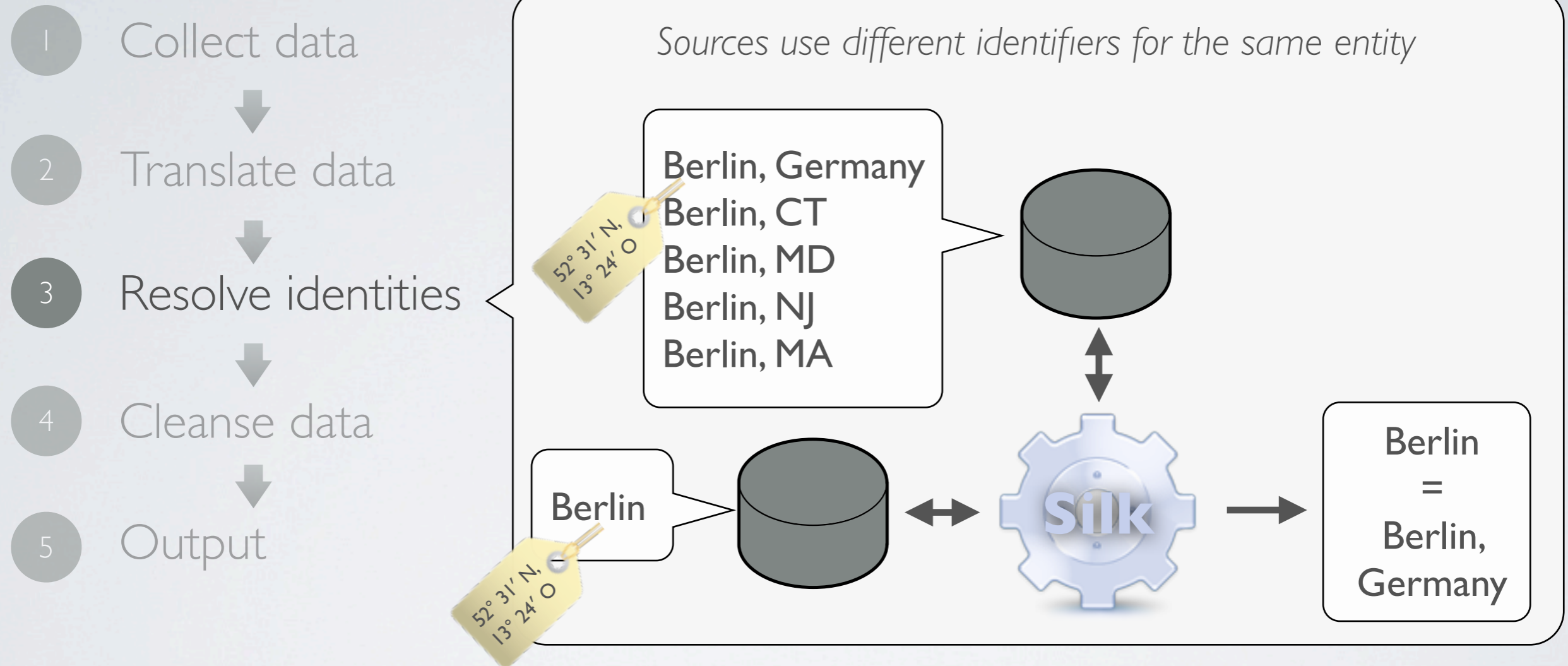
- RDF dumps (various formats)
- SPARQL Endpoints
- Crawling Linked Data

# LDIF PIPELINE



- Mappings expressed in RDF (Turtle)
- Simple mappings using OWL / RDFs statements (x rdfs:subClassOf y)
- Complex mappings with SPARQL expressivity
- Transformation functions

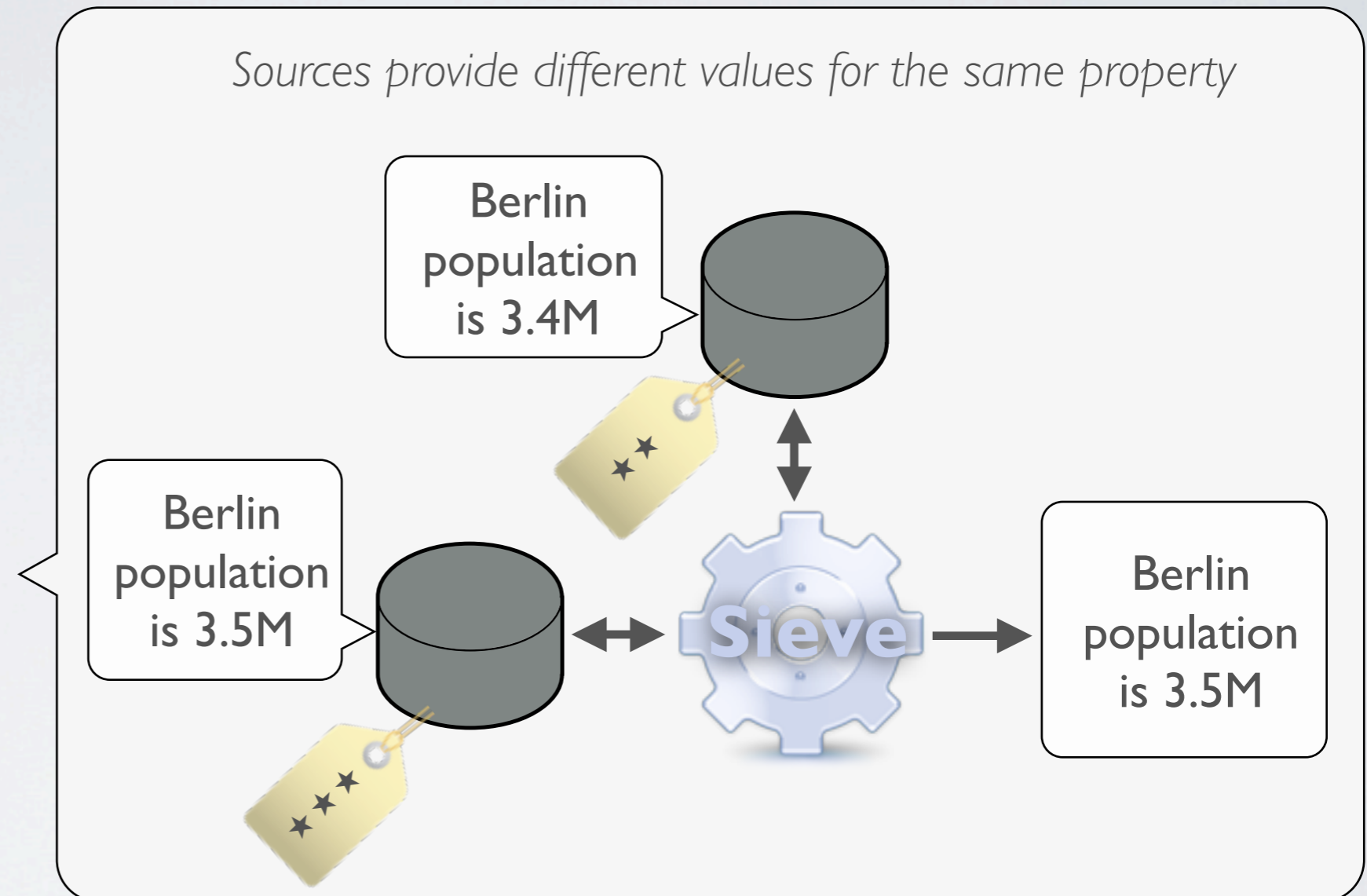
# LDIF PIPELINE



- Profiles expressed in XML
- Supports various comparators and transformations

# LDIF PIPELINE

- 1 Collect data
- 2 Translate data
- 3 Resolve identities
- 4 Cleanse data
- 5 Output



- Profiles expressed in XML
- Supports various quality assessment policies and conflict resolution methods

# LDIF PIPELINE



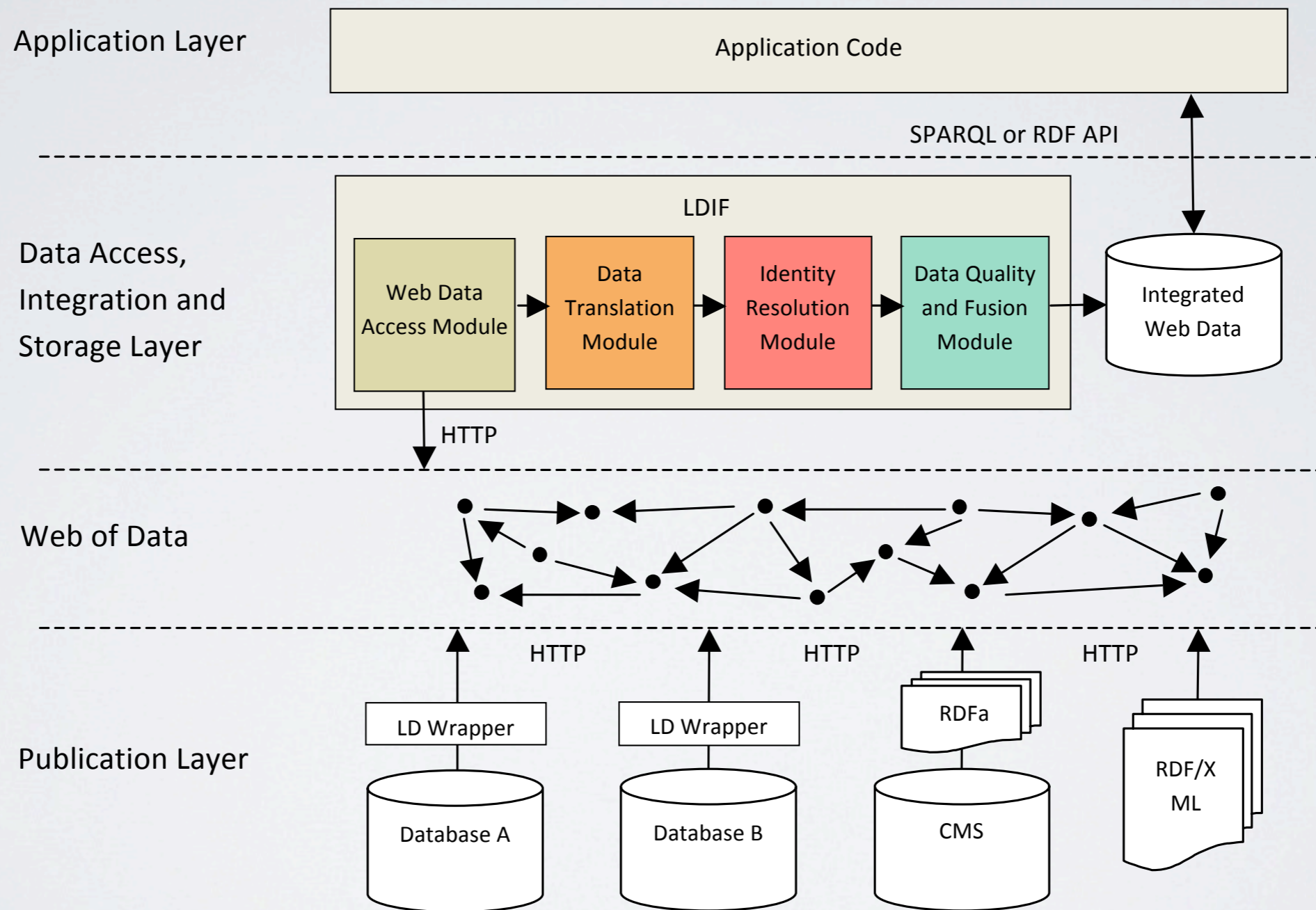
## Output options:

- N-Quads
- N-Triples
- SPARQL Update Stream
- Provenance tracking using Named Graphs





# LDIF ARCHITECTURE



# LDIF VERSIONS

- In-memory
  - keeps all intermediate results in memory
  - fast, but scalability limited by local RAM
- RDF Store (TDB)
  - stores intermediate results in a Jena TDB RDF store
  - can process more data than In-memory but doesn't scale
- Cluster (Hadoop)
  - scales by parallelizing work across multiple machines using Hadoop
  - can process a virtually unlimited amount of data

# THANK YOU

- Website: <http://ldif.wbsg.de>
- Google group: <http://bit.ly/ldifgroup>
- Supported in part by
  - Vulcan Inc. as part of its [Project Halo](#)
  - EU FP7 project [LOD2 - Creating Knowledge out of Interlinked Data](#) (Grant No. 257943)