

Semantic Data Management



TMDT



**BERGISCHE
UNIVERSITÄT
WUPPERTAL**

Different actors have different assumptions, goals and challenges when working with data



Data Provider

- Knows his data and use cases
- **Provides data** so others can work with them
- Needs a way to **share his knowledge**



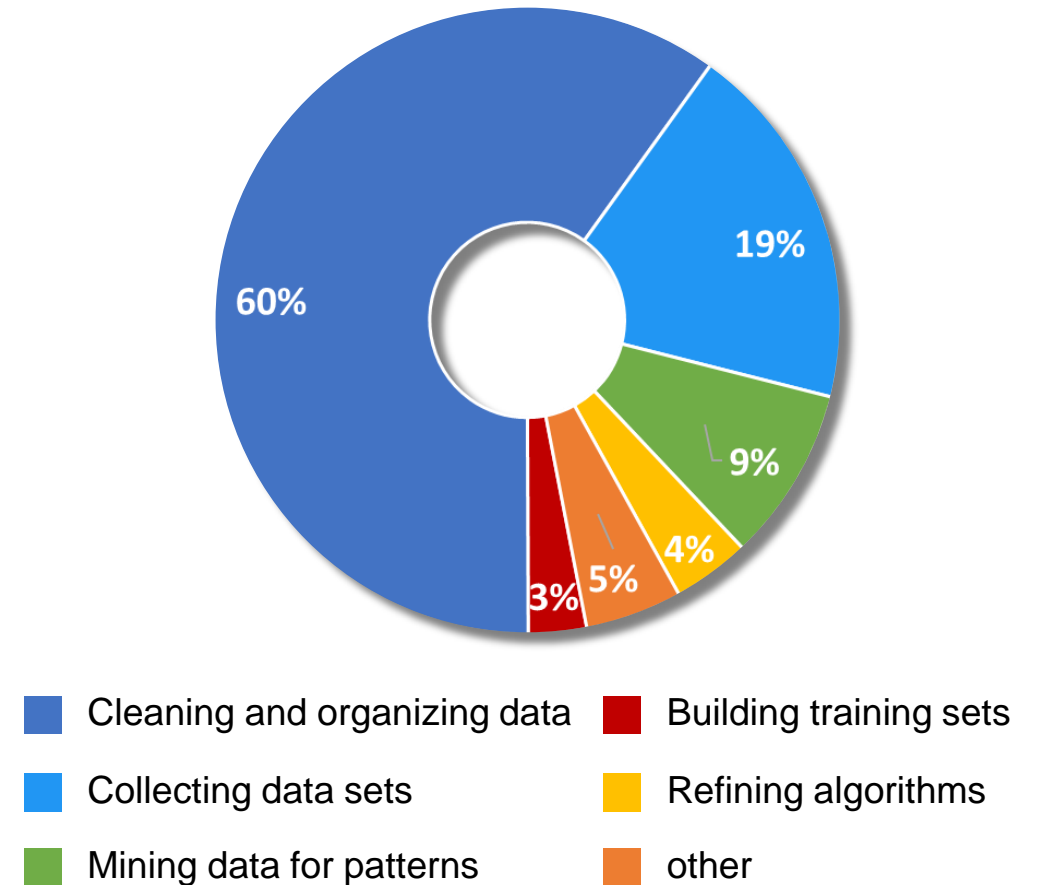
Data Consumer

- Wants to work with data
- Needs to **find the right data** for his use case
- Needs to **access the data**
- Needs to **understand the data**

Current Situation



- Data scientists spend almost **80% of their time with collecting, cleaning and organizing data**
- Common challenges:
 - **Quality, comprehensiveness, availability, retrievability, and expandability** of data
- **To make the most of your data**, they need to be available in a way that not only allows but supports further processing and future use cases
- Semantic data management effectively reduces the **time-to-application** and prepares your data for **future use cases**

What data scientists spend the most time doing





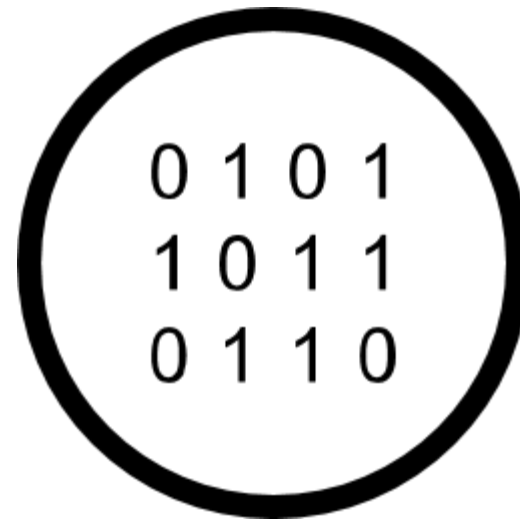
Source: Forbes, 2016




Challenges





Heterogeneous Provision  


Heterogeneous Models  Relational  Hierachy

Heterogeneous Quality  



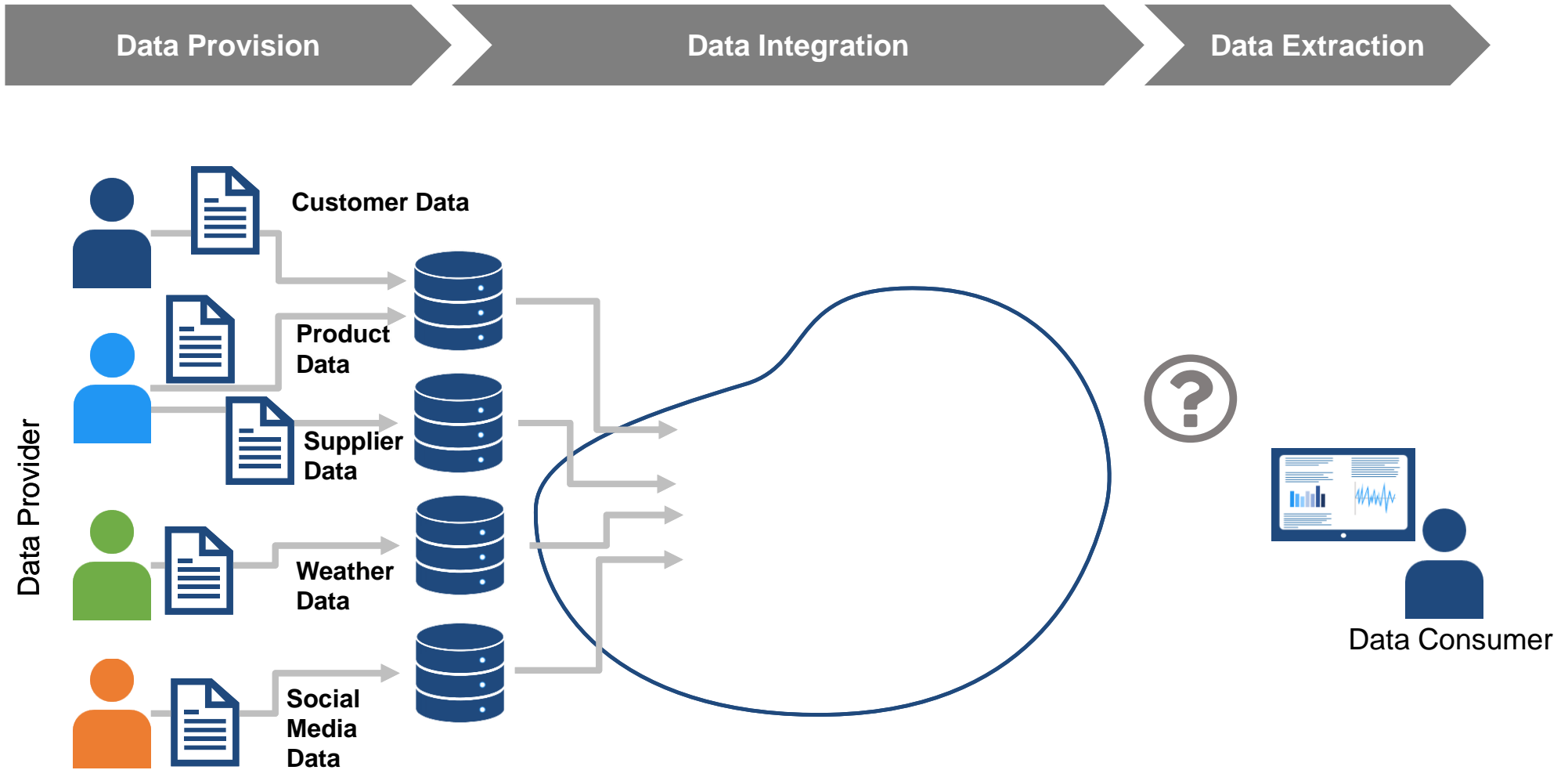
Heterogeneous Formats   

Heterogeneous Use Cases    

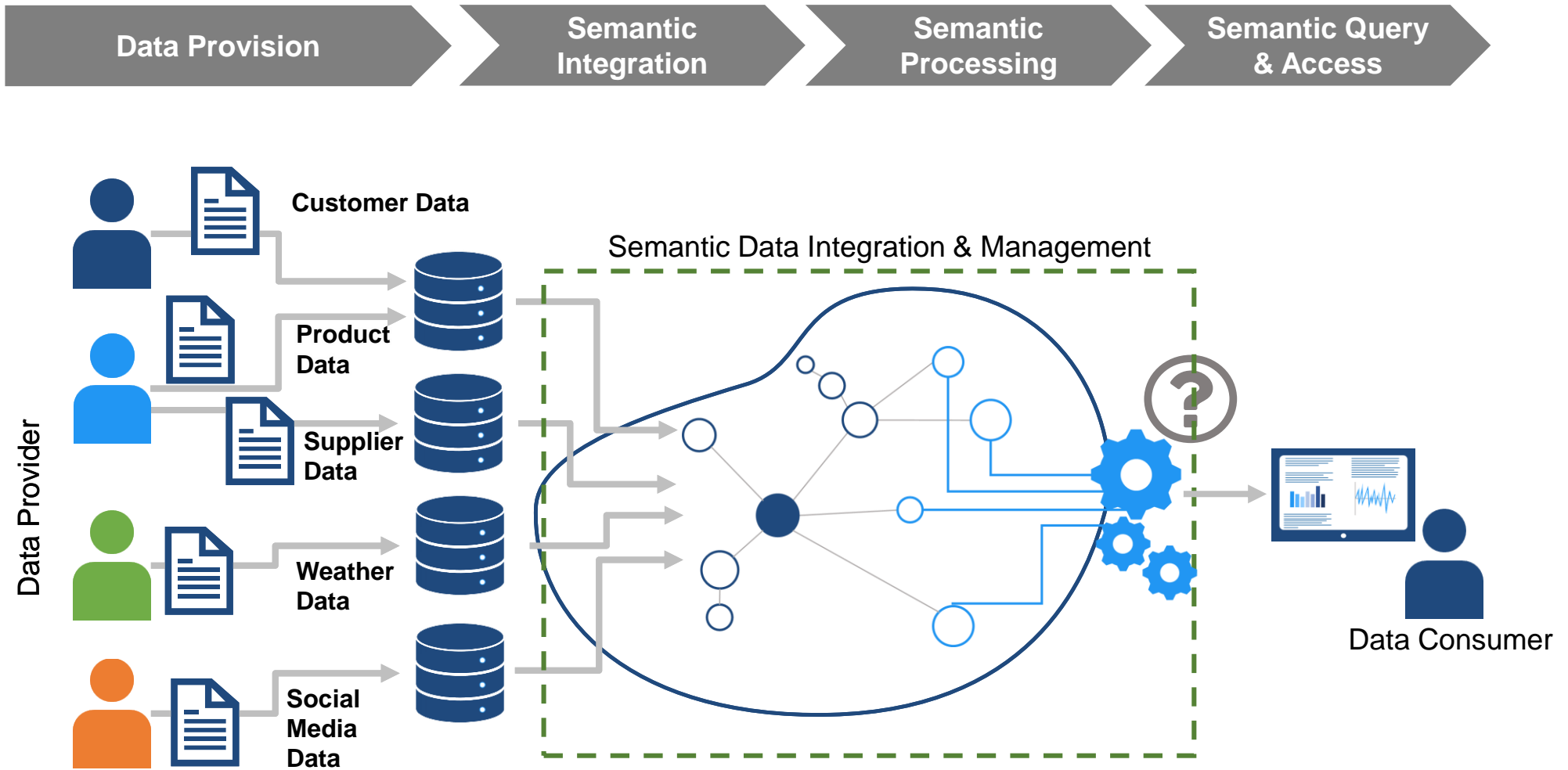
Heterogeneous Semantics  °C °F K

Semantic Data Management

Schematic Representation of Initial Situation



Schematic Representation of Semantic Data Integration & Management



Our Mission

- Google **finds** and "**understands**" **websites without the need** for people to describe them explicitly
 - How can **enterprise data** be provided and found without complex and possibly incomprehensible descriptions?
- Our aim is to reduce the **Time-to-Application** of enterprise data by supporting data scientist to answer pressing questions like:
 - Where to find the data needed for an application, analysis etc.?
 - What information is contained in the available data?
 - Are there other data sources that might be relevant to solve a challenge?

Research Goals



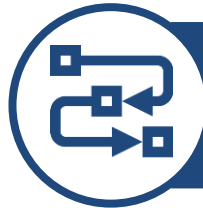
Bridge the gap between physical & virtual world by developing intuitive concepts for Big Data collection



Develop new approaches for constructing and evolving Knowledge Graphs



Develop semantic mapping approaches that leverage machine learning & external knowledge



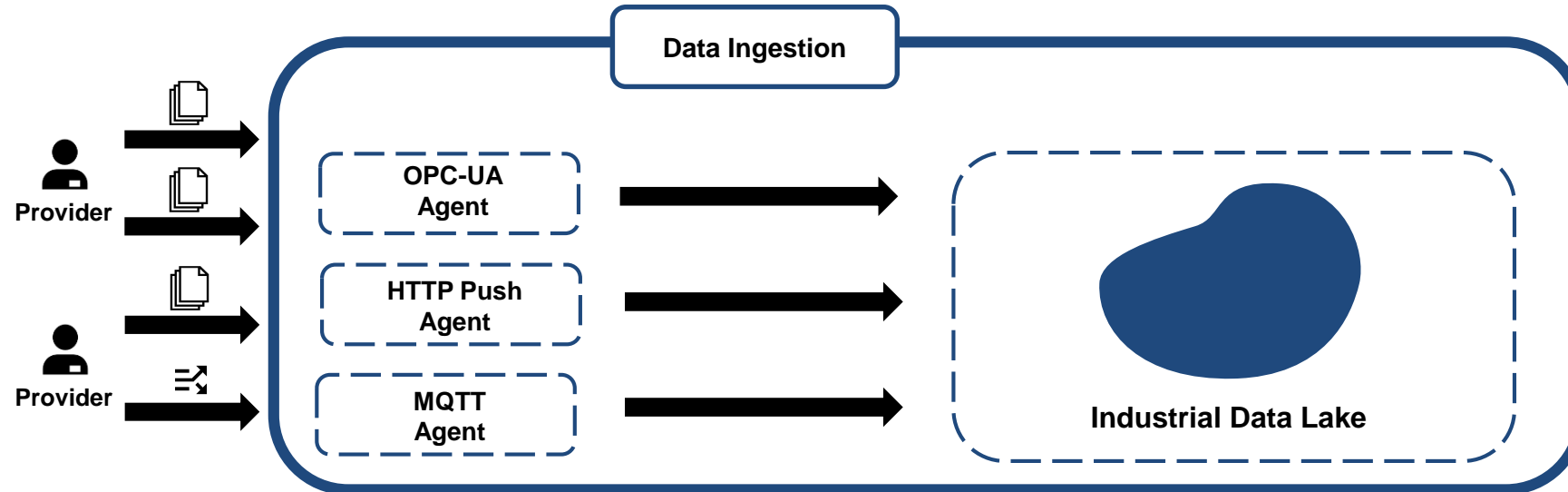
Enable semantic processing for data based on evolving Knowledge Graphs



Improve the usability & accessibility of data with the help of semantic querying



Bridge the gap between physical & virtual world by developing intuitive concepts for Big Data collection



- Plug & Ingest: Autonomous agents monitor data sources and collect the data in heterogeneous company-wide environments
- Agents proactively react to changing data sources and interfaces
- Data is ingested into a central location, e.g., into a data lake



Develop new approaches for constructing and evolving Knowledge Graphs

Knowledge graph learns continuously from interacting with its environment



Data Provider

Defines Semantic Models



Data Consumer

Queries Data Sets

01100
10110
11110

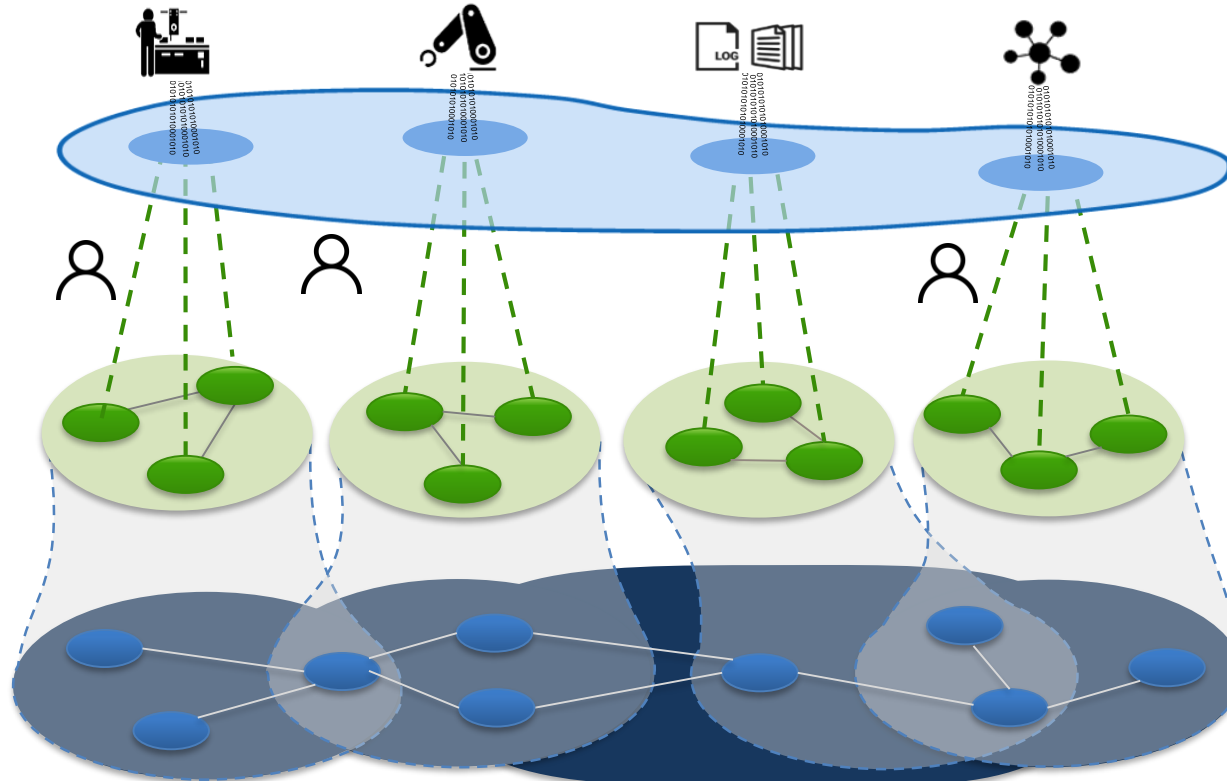
System

Autonomous Evolution

- Knowledge Graph **grows and strengthens continuously**
- Knowledge Graph allows **statements about relationships** of data sets
- **Advanced AI algorithms** monitor and support the evolution of the knowledge graph



Develop new approaches for constructing and evolving Knowledge Graphs



▶ **Data collection and access**
Construction Industrial Data Lake

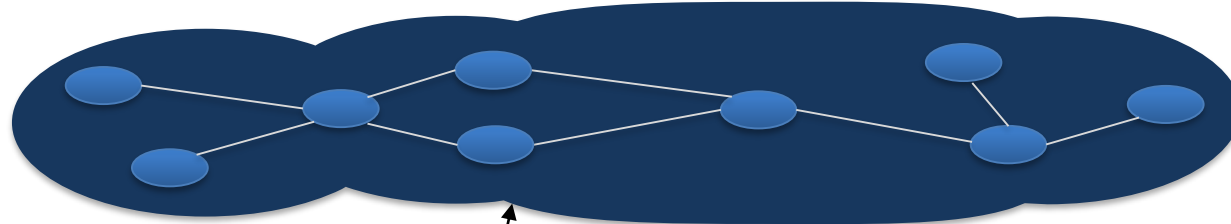
▶ **Data Description**
Local Knowledge Representation

▶ **Learning from Data**
Global Knowledge Graph




Semantic Mapping: External Knowledge Recommendation



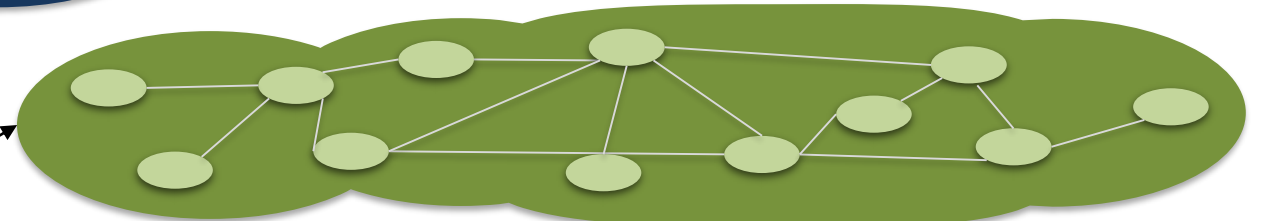
Develop semantic mapping approaches that leverage machine learning & external knowledge



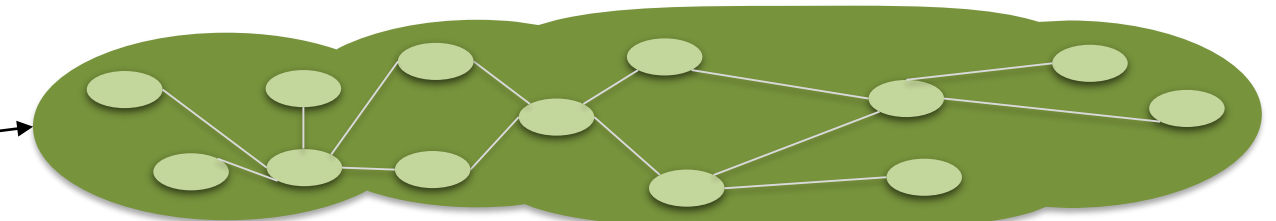
Knowledge Graph

-  Timestamp (92,2%)
-  Identifier (4,2%)
-  ...

TS	TEMP	POS_Y	POS_X	ALT
1476835200	3	50.780296	6.104511	12
1476835200	4	50.774423	6.085856	4
1476835200	10	50.776057	6.043272	11



External Knowledge 1

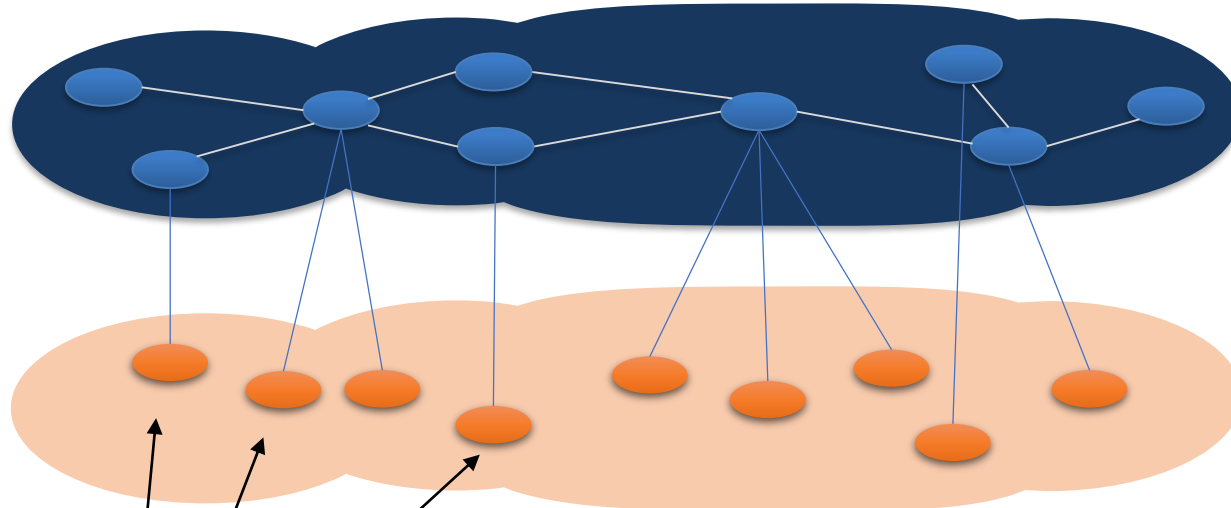


External Knowledge 2

Semantic Mapping: Machine Learning Recommendation



Develop semantic mapping approaches that leverage machine learning & external knowledge



Legend:

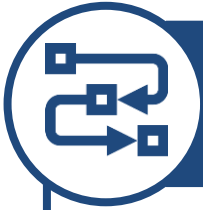
- Blue circle: Timestamp (87,2%)
- Orange circle: Number (9,1%)
- ...: ...

TS	TEMP	POS_Y	POS_X	ALT
1476835200	3	50.780296	6.104511	12
1476835200	4	50.774423	6.085856	4
1476835200	10	50.776057	6.043272	11

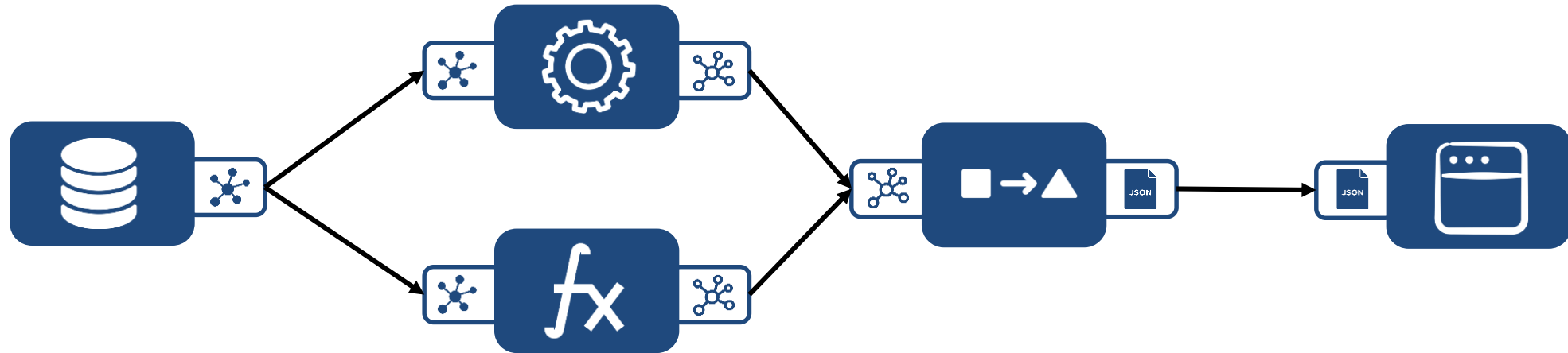
▶ Knowledge Graph

▶ Machine Learning Layer

▶ Recommendation



Enable semantic processing for data based on evolving Knowledge Graphs



- Consider the continuously changing knowledge graph and enable the processing of data based on the provided semantics
- Prepare the data so that they are directly useable for analytics or applications



Improve the usability & accessibility of data with the help of semantic querying



Natural Language &
Visual Queries



Expand queries



Personalize the search
experience

- Enable data set queries based on natural language and visual interfaces
- Leverage external knowledge databases to expand queries
- Learn from users queries in order to improve the personalized search experience



Selected Research Projects and Applications

Cross-site Enterprise Data Lake in Automotive Industry

Cross-site Enterprise Data Lake



Connection of **five plants on two continents**



Logistics
Construction
Paint Shop
Assembly
Rework



Live operation for more than 2 years
currently more than 6 TB/year



Configurable, intelligent **software connectors**



Optimization of throughput times on the conveyor belt

Goal

- **Construction** of a **cross-site enterprise data lake** that continuously collects and **integrates data from multiple factories**

Approach

- Conceptual **design of cross-site, scalable, non-invasive data acquisition and -management**
- **Integration of heterogeneous data sources** in a common long-term data storage for each site through **configurable connector modules**
- **Automatic metadata enrichment**
- **Cross-site transparent consolidation of local data lakes** based on **information models**
- Provision of consolidated data through **common methods and familiar tools** (Tableau, Power BI, RapidMiner,...)

Results

- Functional cross-site data acquisition and management infrastructure for five plants on two continents
- Over 1000 different data sources integrated

Semantic Data Management Knowledge Graph-based Data Management Platform (KGBDM Platform)

KGBDM Platform

The screenshot displays the KGBDM Platform interface. The top section shows a 'Dashboard' with '3 Data Sources' and 'Data Views'. Below this, there are sections for 'Highest Rated Data Sources' and 'Most Recent Data Sources'. The middle section shows a 'Data Inventory' with a 'Create Data Source' form and a flow diagram: 'File Collector' -> 'Live Content Reader' -> 'JSON Converter'. The bottom section shows a 'Weather Station' knowledge graph. The graph has a 'ROOT' node connected to 'Temperature # temp', 'Pressure # press', and 'Timestamp # timestamp'. 'Timestamp # timestamp' is connected to 'lat_lon', which is further connected to 'Latitude # 1st' and 'Longitude # 2nd'. A sidebar on the left lists 'Entities' and 'Relations' with a search bar.

Goal

- Reduce the **time-to-application** by handling **heterogeneous data sources** on an information level

Approach

- Provide **different connector modules** for adding data sources
- Data providers **create semantic models** for their data sets
- Semantic model creation is **supported with the help of recommendation systems and artificial intelligence**
- Semantic Data Integration** based on the provided semantic models
- Data access via the **semantic search**
- No need to define ontologies.** The KGBDM platform **learns a knowledge graph** based on the user provided semantic models and queries

Results

- Available state of the art semantic data management platform** that is **used in various research and development projects**

Your Contact Person:

André Pomp, M.Sc.

Tel: +49 (0)202 439 1153

pomp@uni-wuppertal.de

Chair for Technologies and Management of Digital Transformation

Univ. Prof. Dr. Ing. Tobias Meisen

<https://www.tmdt.uni-wuppertal.de/>

Campus Freudenberg

Rainer-Gruenter-Str. 21

D-42119 Wuppertal

Germany

University of Wuppertal

School of Electrical, Information and Media Engineering



TMDT



**BERGISCHE
UNIVERSITÄT
WUPPERTAL**