



WIRTSCHAFTS  
UNIVERSITÄT  
WIEN VIENNA  
UNIVERSITY OF  
ECONOMICS  
AND BUSINESS



# Open City Data Pipeline

Axel Polleres

*formerly: Siemens AG Österreich, now WU Wien, Austria*

# City Data – Important for Infrastructure Providers & for City Decision Makers

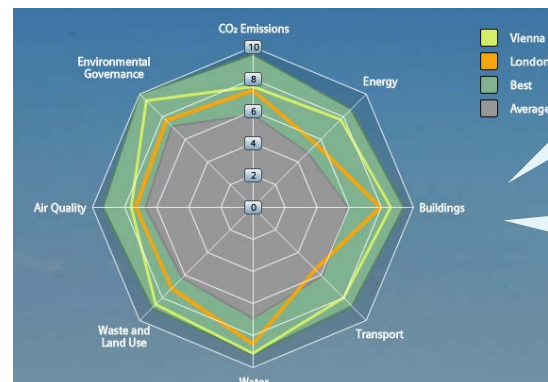
- City Assessment and Sustainability reports
- Tailored offerings by Infrastructure Providers



... however, these are often **outdated** before even published!

→ Needs **up-to-date City Data** and **calculates City KPIs** in a way that allows to display the current state and run scenarios of different product applications.

e.g. towards a “Dynamic” Green City Index:



Goal (short term):

- Leverage Open Data for calculating a city' performance from public sources on the Web **automatically**

Goal (long term):

- Define and Refine KPI models to assess specific impact of infrastructural investments and gather/check input **automatically**

# Current State of Data and Benchmarking System

- **Collecting Data for City Assessment and Benchmarking:**
- Data collection for various studies within Siemens is a manual, time-intensive process, using statistical data as well as questionnaires to city stakeholders.
- **Much of this data is available online:**
- City Open Data Initiatives publish more and more data with frequent updates



- City data format standards and regulations are developing:



e.g. EU INSPIRE directive, or



eurostat's UrbanAudit Collection of city indicators

- **Benchmarking Systems and KPIs** | The Green City Index **MOST LIVEABLE CITIES**
- City Indexes benchmark cities based on top down data (example: tax income from petrol, tax/L → Car CO<sub>2</sub> emissions). Bottom up approaches only if no top down data available, for approximation (example: No. cars, average distance driven per year, Emission factor → Car CO<sub>2</sub> emissions)
- Approaches allowing scenario comparison and calculations of system impacts only on a city specific case study base.

# Leveraging Open Data: Openly available urban indicators frameworks (already available in Linked Data!)

- Data example: Urban Audit
- (ca. 300 indicators for European
- 330 cities, maintained by Eurostat)



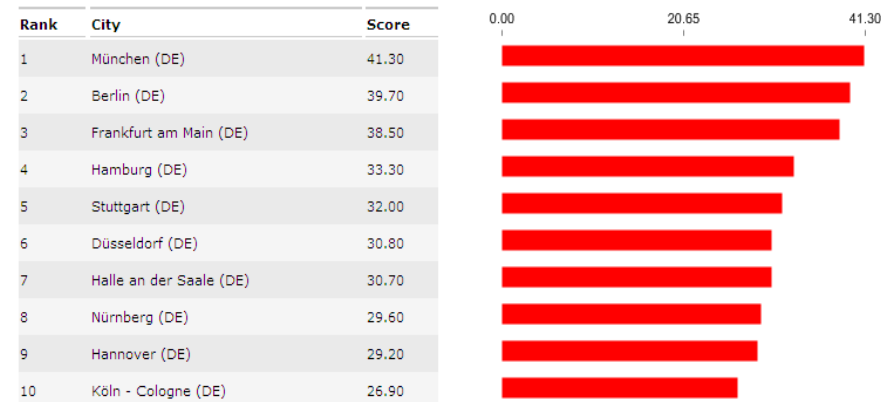
Cost of a monthly ticket for public transport (for 5-10 km)

You are on page 1 of 4 (37 records)



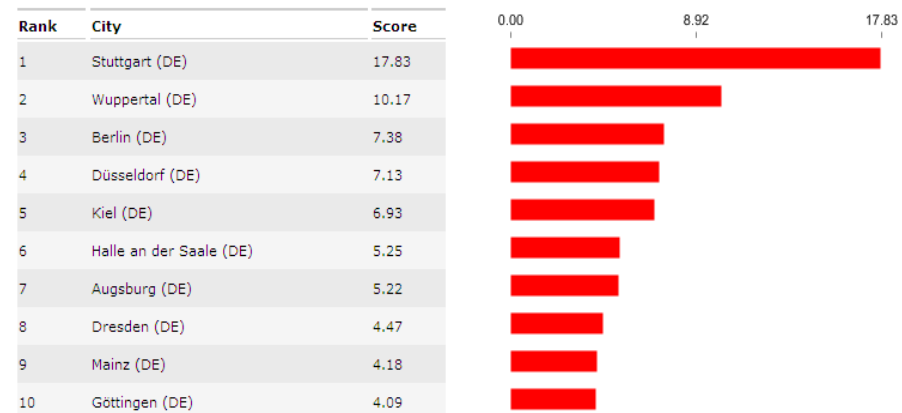
Proportion of journeys to work by public transport (rail, metro, bus, tram)

You are on page 1 of 4 (40 records)



Number of stops of public transport per km2

You are on page 1 of 4 (38 records)



# Leveraging Open Data: Other Open Data Sources

Vienna - Wikipedia, the free encyclopedia - Mozilla Firefox

en.wikipedia.org/wiki/Wienna

**Vienna**  
Coat of arms

Country: Austria  
State: Wien

Government  
 ▪ Mayor: Michael Häupl (SPÖ)  
 ▪ Vice-Mayor: Maria Vassilakou (ÖVP)

Area  
 ▪ City: 414.65 km<sup>2</sup> (160.1 mi<sup>2</sup>)  
 ▪ Land: 395.26 km<sup>2</sup> (152.8 mi<sup>2</sup>)  
 ▪ Water: 19.39 km<sup>2</sup> (7.5 sq mi)

Elevation: 151 (Lobau) – 542

Population (2011)  
 ▪ City: 1,714,142  
 ▪ Density: 4,002.2/km<sup>2</sup> (10,380/sq mi)  
 ▪ Urban: 1,983,836  
 ▪ Metro: ca. 2,419,000

OpenStreetMap - Siemens AG

http://www.openstreetmap.org/index.html?lat=48.21&lon=16.37&zoom=11

Suchen: We bin ich?

Kraftfahrzeugbestand in Wien seit 2001 - Siemens AG

http://www.wien.gv.at/mfz/statistik/wien

Kraftfahrzeugbestand in Wien seit 2001 \*

Jahr	Insgesamt	Pkw	Omnibusse	Lkw	Zug- maschinen **	Sonstige Kraft- fahrzeuge ***	Kraftträder ****
2001	782.510	646.283	3.725	58.968	3.182	5.001	65.351
2002	784.865	647.382	3.641	58.132	3.212	4.993	67.505
2003	790.963	652.418	3.602	58.396	3.282	4.907	68.358
2004	794.109	655.172	3.678	58.322	3.348	4.789	68.800
2005	795.480	655.806	3.535	58.506	3.411	4.794	69.428
2006	799.748	658.081	3.546	58.742	3.417	4.766	71.196
2007	802.209	657.426	3.604	59.619	3.487	4.737	73.336
2008	805.539	657.192	3.607	60.628	3.546	4.747	75.819
2009	814.624	663.926	3.726	60.796	3.573	4.645	77.958
2010	821.999	669.279	3.716	61.185	3.601	4.652	79.566

Quelle: Statistik Austria - Kfz-Bestand.  
 \* Stichtag 31.12.  
 \*\* Sattelfahrzeuge, Motorkarren sowie land- und forstwirtschaftliche Zugmaschinen.  
 \*\*\* Selbsterfindende Arbeitsmaschinen (einschließlich sonstige Kfz)

Structured information on most cities in the world (base indicators: population, economy, climate, ...)

Free GIS data for most cities in the world (base information: area, land-use, administrative districts, ...)

Open Government Data

# City Data Pipeline: Overview



Dynamic Calculation of KPIs at variable Granularity (City, District, Neighbourhood, Building)

*Extensible  
CityData  
Model*

Cities:

Berlin, Vienna, London, ...



+

Open Data.

DBpedia

WORLD BANK

Open Data

UNdata  
A world of information



# City Data Pipeline: Overview



Dynamic Calculation of KPIs at variable Granularity (City, District, Neighbourhood, Building)

*Extensible  
CityData  
Model*

1. Periodic Data Gathering of registered sources ("Focused Crawler"):  
Various Formats (CSV, HTML, XML ...) & Granularity (monthly, annual, daily)

Cities:

Berlin, Vienna, London, ...

+

Open Data. WORLD BANK

DBpedia Open Data UNdata  
A world of information

# City Data Pipeline: Overview



Dynamic Calculation of KPIs at variable Granularity (City, District, Neighbourhood, Building)

*Extensible  
CityData  
Model*

2. Semantic Integration: Unified Data Model, Data Consolidation

1. Periodic Data Gathering of registered sources ("Focused Crawler"):  
Various Formats (CSV, HTML, XML ...) & Granularity (monthly, annual, daily)

Cities:

Berlin, Vienna, London, ...

+

Open Data.

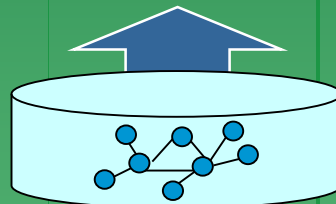
WORLD BANK  
DBpedia Open Data UNdata  
A world of information



# City Data Pipeline: Overview



Dynamic Calculation of KPIs at variable Granularity (City, District, Neighbourhood, Building)



*Extensible  
CityData  
Model*

2. Semantic Integration: Unified Data Model, Data Consolidation

1. Periodic Data Gathering of registered sources ("Focused Crawler"):  
Various Formats (CSV, HTML, XML ...) & Granularity (monthly, annual, daily)

Cities:

Berlin, Vienna, London, ...

+

Open Data.

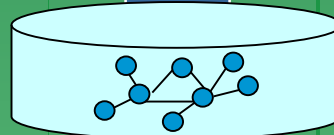
WORLD BANK  
DBpedia Open Data UNdata  
A world of information

# City Data Pipeline: Overview



Dynamic Calculation of KPIs at variable Granularity (City, District, Neighbourhood, Building)

3. Analysis/Statistical Correlation/Aggregation:  
Statistical Methods, Semantic Technologies, Constraints



*Extensible  
CityData  
Model*

2. Semantic Integration: Unified Data Model, Data Consolidation

1. Periodic Data Gathering of registered sources ("Focused Crawler"):  
Various Formats (CSV, HTML, XML ...) & Granularity (monthly, annual, daily)

Cities:

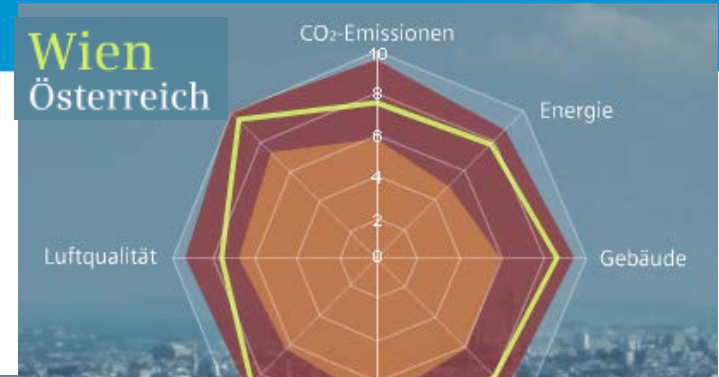
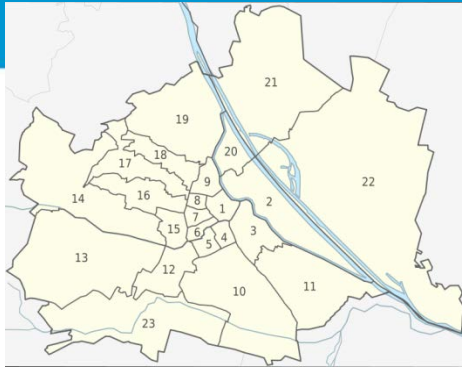
Berlin, Vienna, London, ...

+

Open Data. WORLD BANK

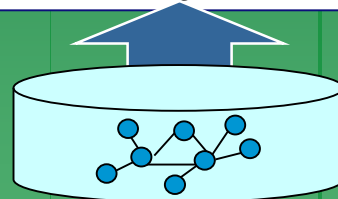
DBpedia Open Data UNdata  
A world of information

# City Data Pipeline: Overview



Dynamic Calculation of KPIs at variable Granularity (City, District, Neighbourhood, Building)

3. Analysis/Statistical Correlation/Aggregation:  
Statistical Methods, Semantic Technologies, Constraints



*Extensible  
CityData  
Model*

2. Semantic Integration: Unified Data Model, Data Consolidation

1. Periodic Data Gathering of registered sources ("Focused Crawler"):  
Various Formats (CSV, HTML, XML ... ) & Granularity (monthly, annual, daily)

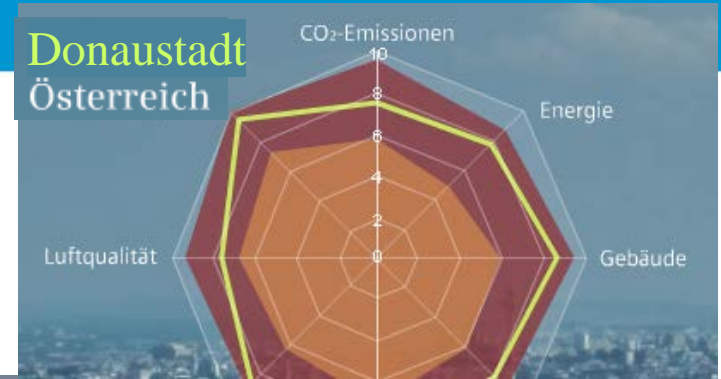
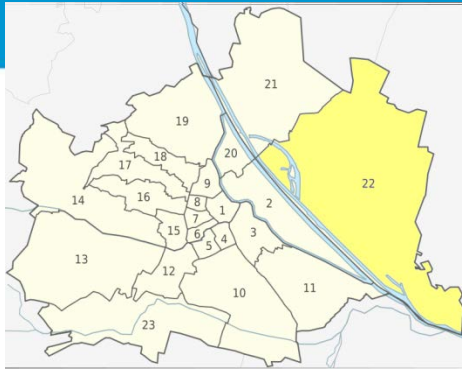
Cities:

Berlin, Vienna, London, ...

+

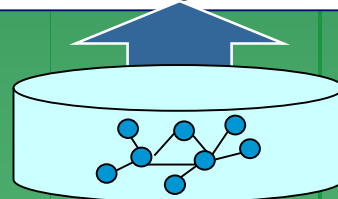
Open Data  
DBpedia Open Data UNdata  
A world of information

# City Data Pipeline: Overview



Dynamic Calculation of KPIs at variable Granularity (City, District, Neighbourhood, Building)

3. Analysis/Statistical Correlation/Aggregation:  
Statistical Methods, Semantic Technologies, Constraints



*Extensible  
CityData  
Model*

2. Semantic Integration: Unified Data Model, Data Consolidation

1. Periodic Data Gathering of registered sources ("Focused Crawler"):  
Various Formats (CSV, HTML, XML ... ) & Granularity (monthly, annual, daily)

Cities:

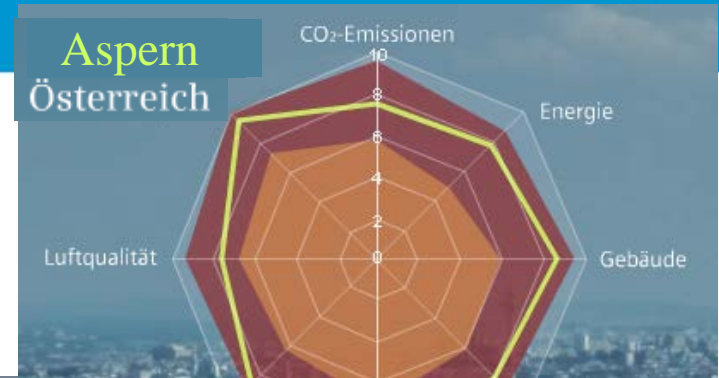
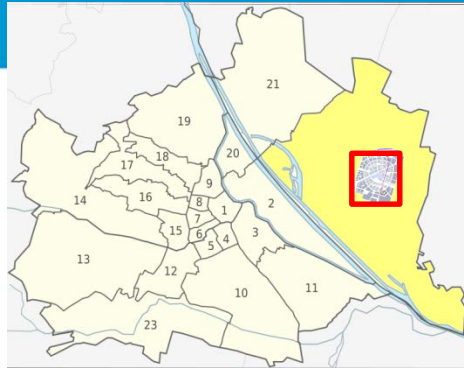
Berlin, Vienna, London, ...

+

Open Data

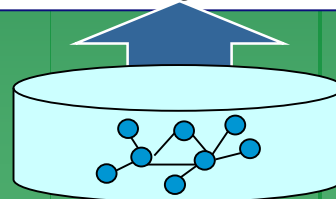
DBpedia Open Data UNdata

# City Data Pipeline: Overview



Dynamic Calculation of KPIs at variable Granularity (City, District, Neighbourhood, Building)

3. Analysis/Statistical Correlation/Aggregation:  
Statistical Methods, Semantic Technologies, Constraints



*Extensible  
CityData  
Model*

2. Semantic Integration: Unified Data Model, Data Consolidation

1. Periodic Data Gathering of registered sources ("Focused Crawler"):  
Various Formats (CSV, HTML, XML ...) & Granularity (monthly, annual, daily)

Cities:

Berlin, Vienna, London, ...

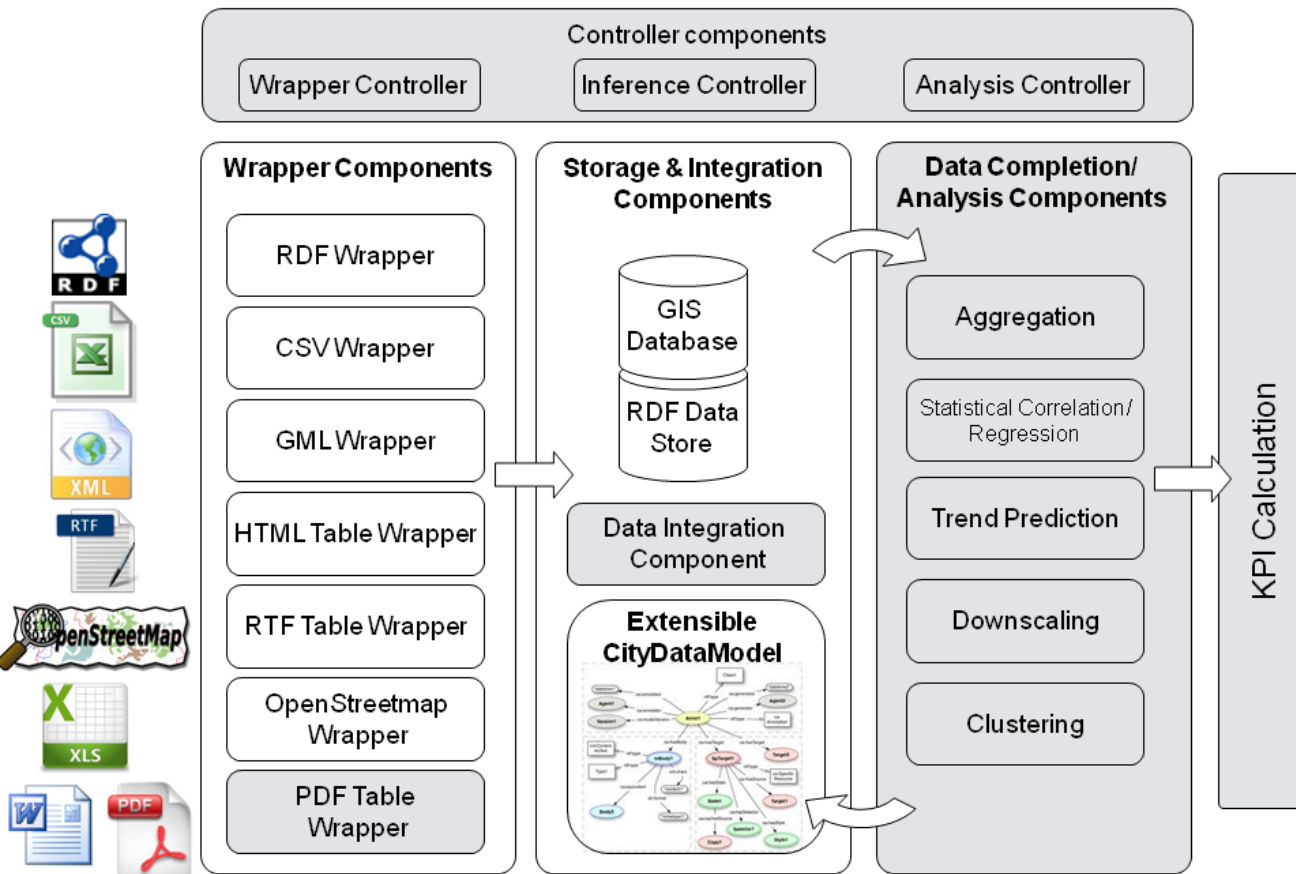
+

Open Data  
DBpedia Open Data UNdata  
A world of information

# City Data Pipeline: Architecture

We have developed a data pipeline to

- (1) (semi-)automatically collect and integrate various Open Data Sources in different formats
- (2) compose and calculate complex city KPIs from the collected data



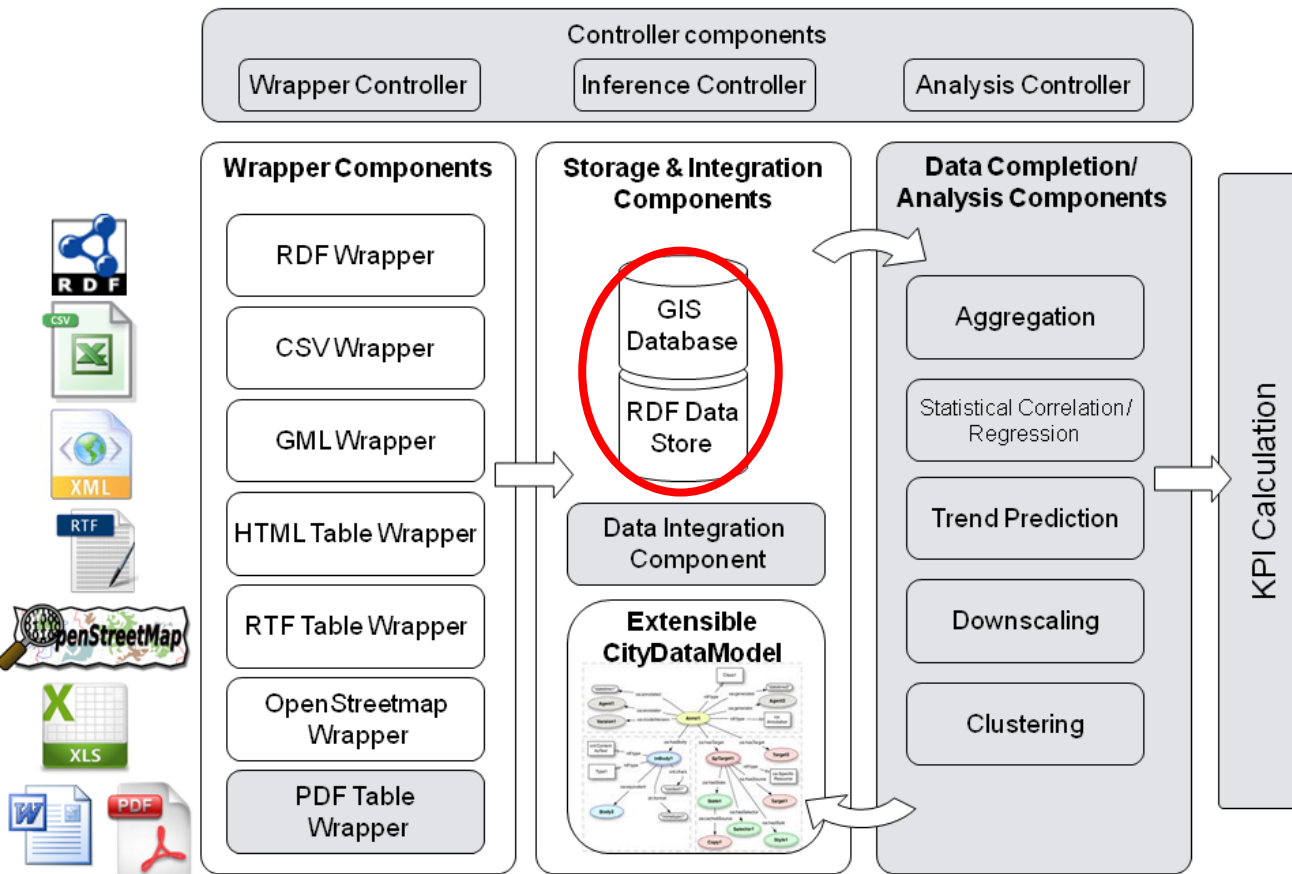


# City Data Pipeline: Architecture

We have developed a data pipeline to

- (1) (semi-)automatically collect and integrate various Open Data Sources in different formats
- (2) compose and calculate complex city KPIs from the collected data

**Semantic Integration  
 Technologies build the  
 core of our Data-Pipeline  
 (RDF Data Store,  
 Ontology-based)**

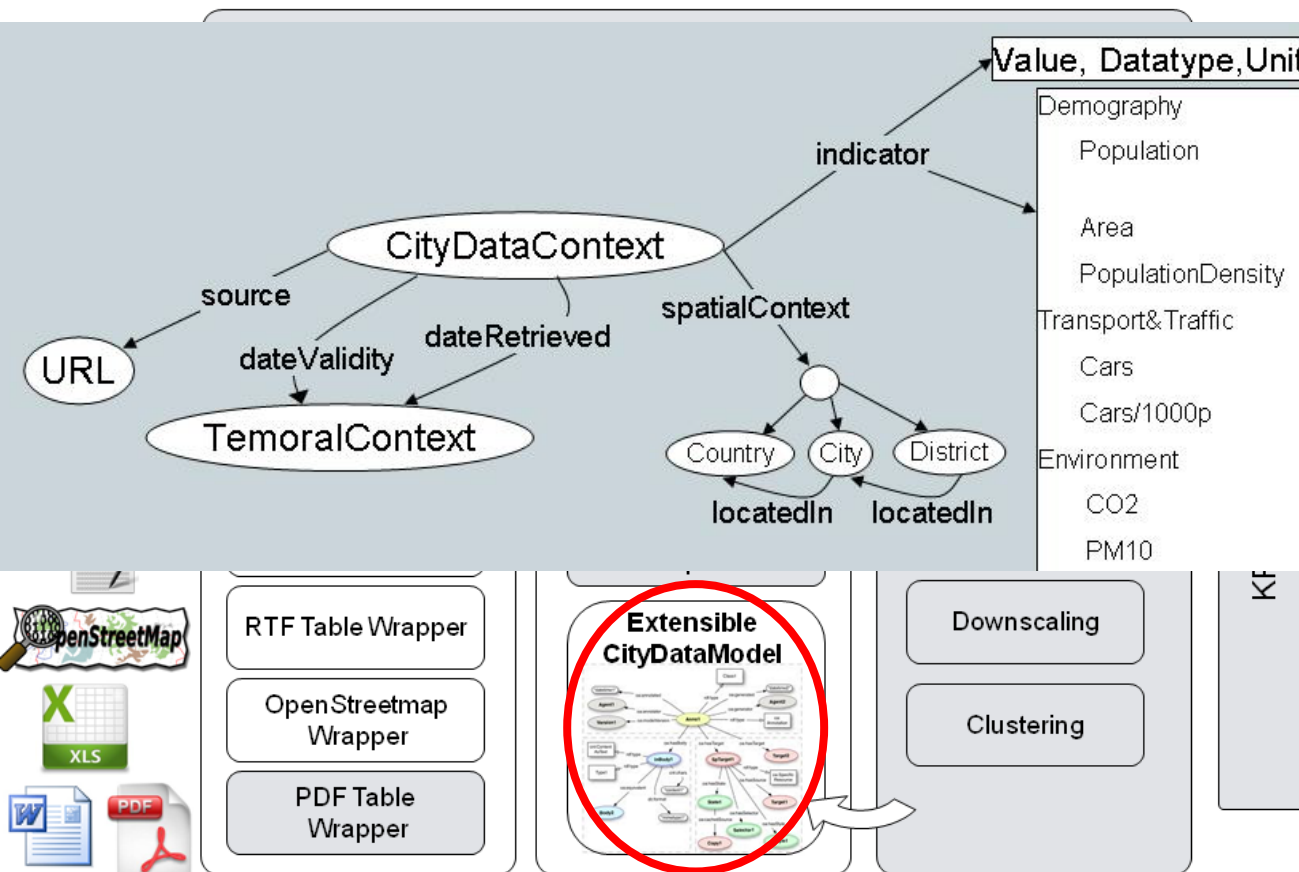


# City Data Pipeline: Architecture

We have developed a data pipeline to

- (1) (semi-)automatically collect and integrate various Open Data Sources in different formats
- (2) compose and calculate complex city KPIs from the collected data

**Semantic Integration Technologies** build the core of **our Data-Pipeline (RDF Data Store, Ontology-based)**

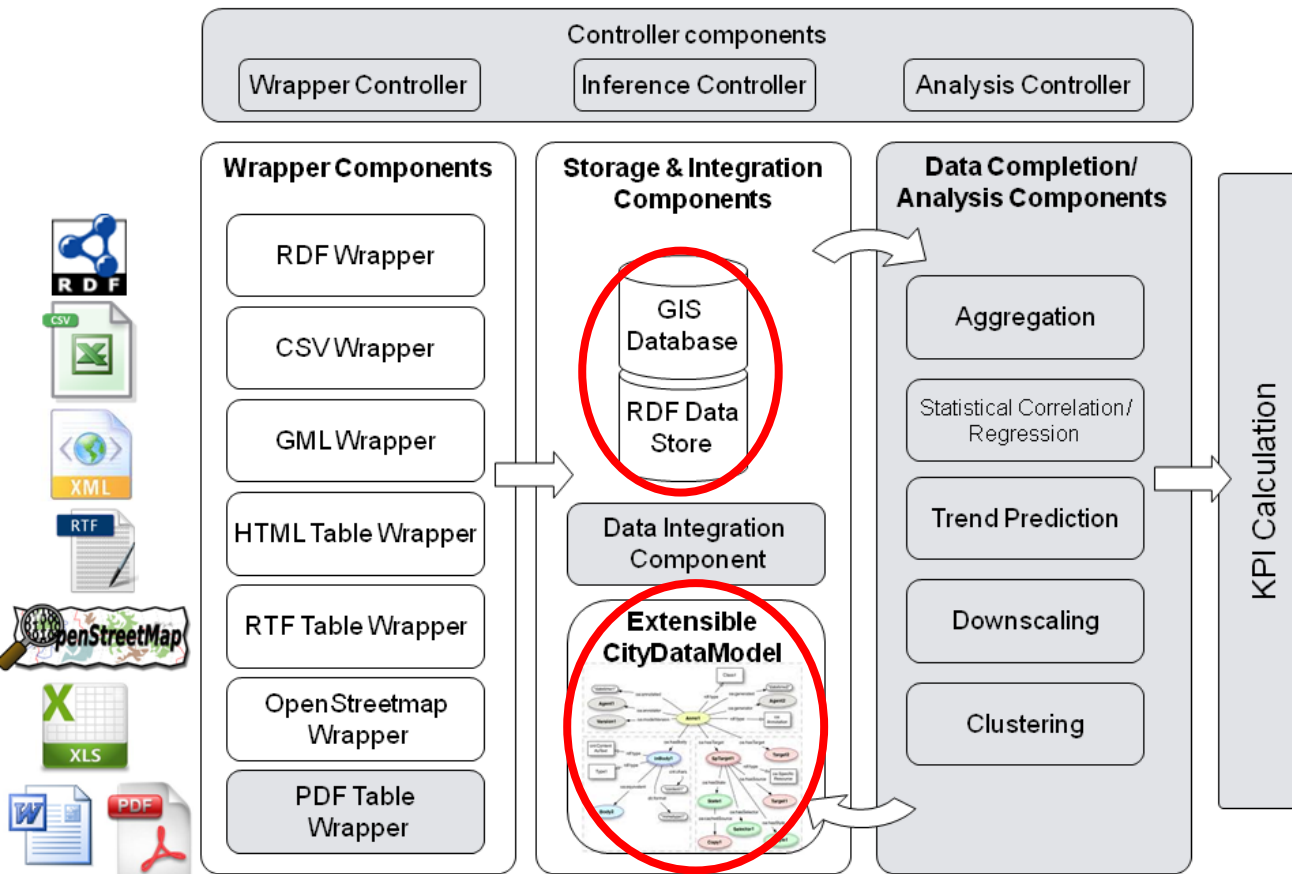


# City Data Pipeline: Architecture

We have developed a data pipeline to

- (1) (semi-)automatically collect and integrate various Open Data Sources in different formats
- (2) compose and calculate complex city KPIs from the collected data

**Semantic Integration Technologies** build the core of **our Data-Pipeline (RDF Data Store, Ontology-based)**

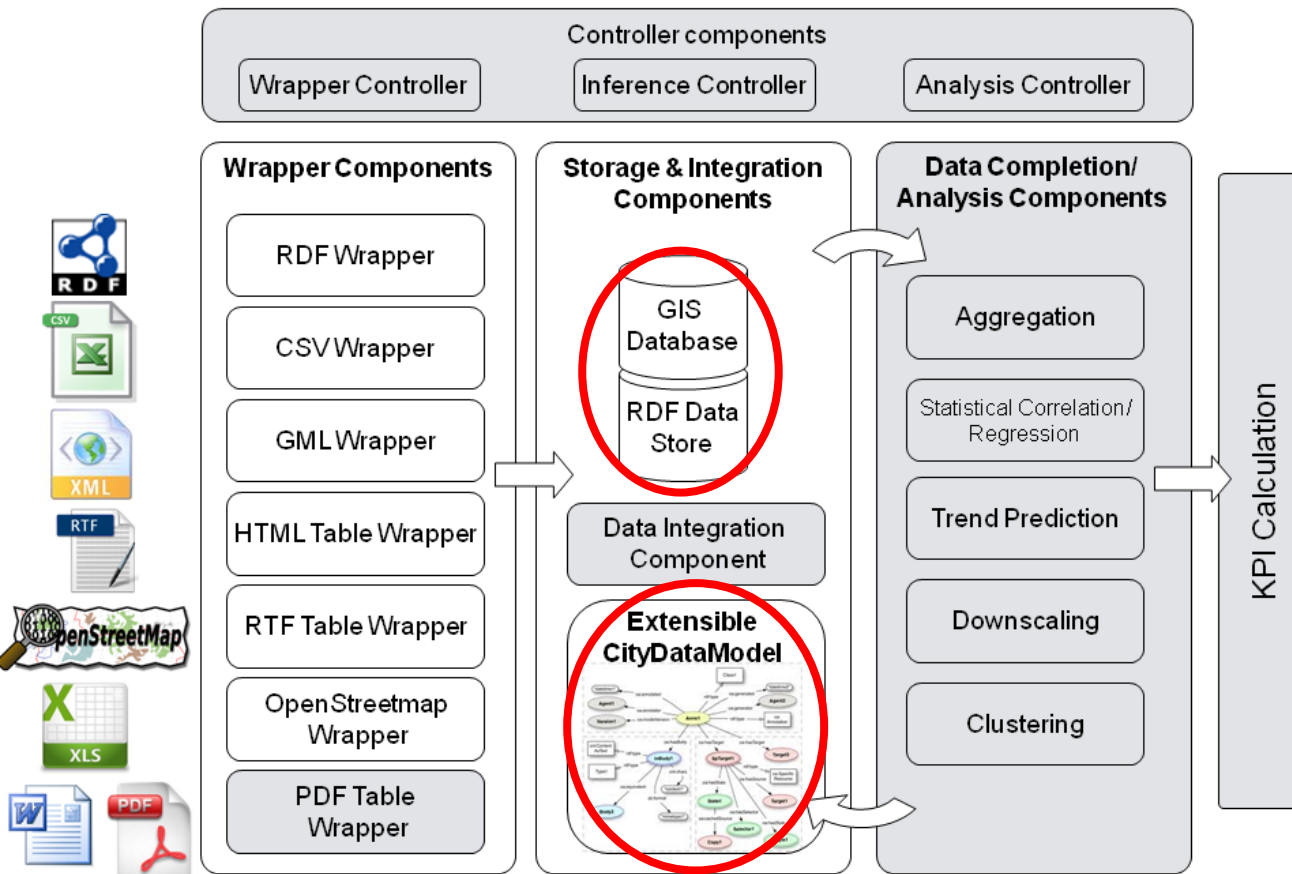


# City Data Pipeline: Architecture

We have developed a data pipeline to

- (1) (semi-)automatically collect and integrate various Open Data Sources in different formats
- (2) compose and calculate complex city KPIs from the collected data

**Semantic Integration Technologies** build the core of **our Data-Pipeline (RDF Data Store, Ontology-based)**



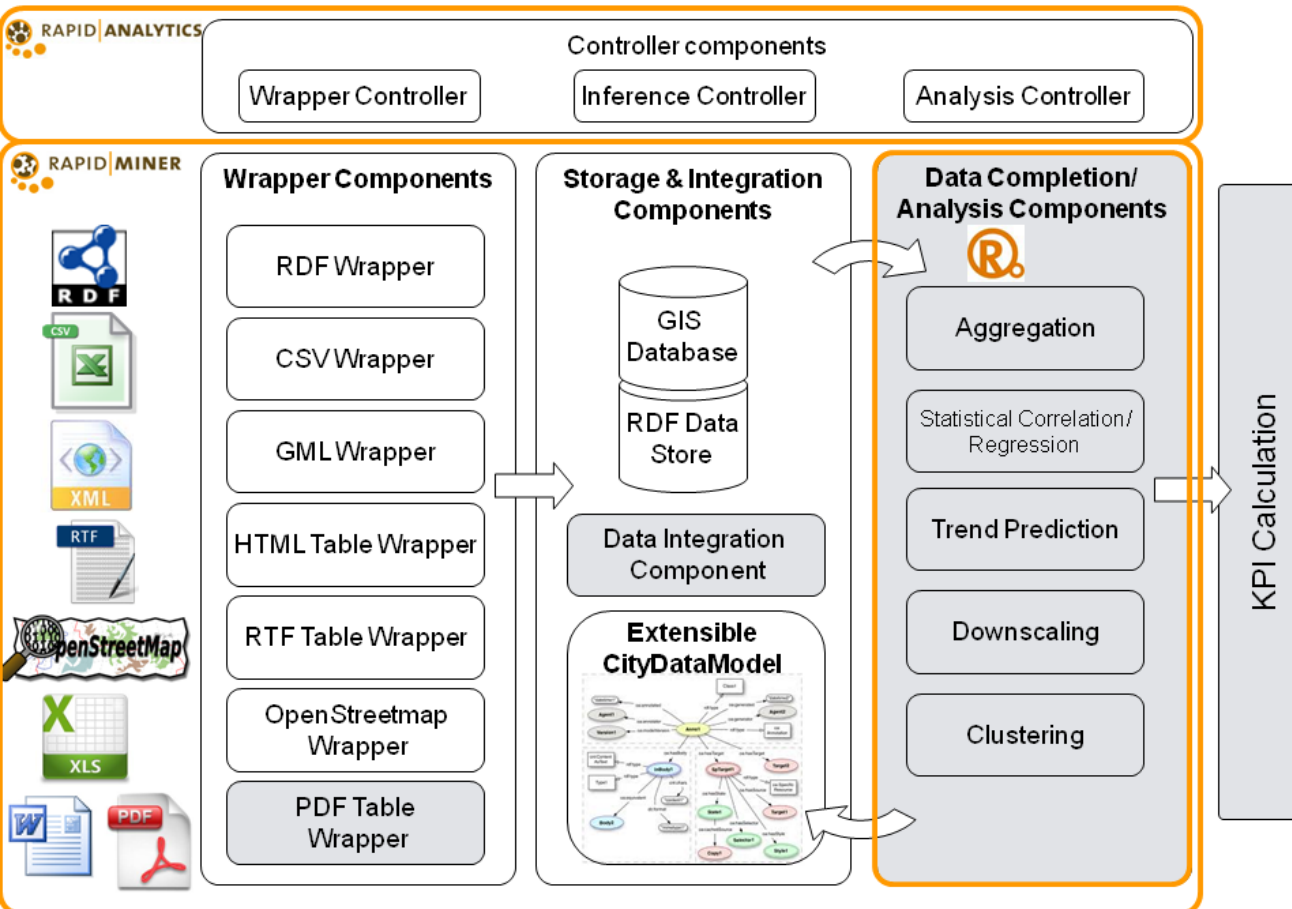
Deploying Semantic Web Standards:



# City Data Pipeline: Architecture

We have developed a data pipeline to

- (1) (semi-)automatically collect and integrate various Open Data Sources in different formats
- (2) compose and calculate complex city KPIs from the collected data



**Semantic Integration Technologies** build the core of **our Data-Pipeline (RDF Data Store, Ontology-based)**

Deploying Semantic Web Standards:



# City Data Pipeline: Current Data - Summary



# City Data Pipeline: Current Data - Summary

- Ca. **475** different indicators
  - *Categories: Demography, Geography, Social Aspects, Economy, Environment, etc.*
- from **32** sources (html, CSV, RDF, ...)
  - *Wikipedia, urbanaudit.org, Statistics from City homepages, country Statistics, iea.org*
- Covering **350+** cities in **28** European countries
  - District Data for selected cities (Vienna, Berlin)
  - Mostly snapshots, Partially covering timelines
  - On average ca. **285** facts per city.
- Examples of sources:
  - UrbanAudit (from <http://eurostat.linked-statistics.org>)
  - <http://geonames.org> (population, georeference, elevation)

Further data sources on target to integrate include:

- <http://data.worldbank.org/> WorldBank (mostly data at country level)
- [www.eea.europa.eu/](http://www.eea.europa.eu/) European Environmental Agency (weather/climate data)

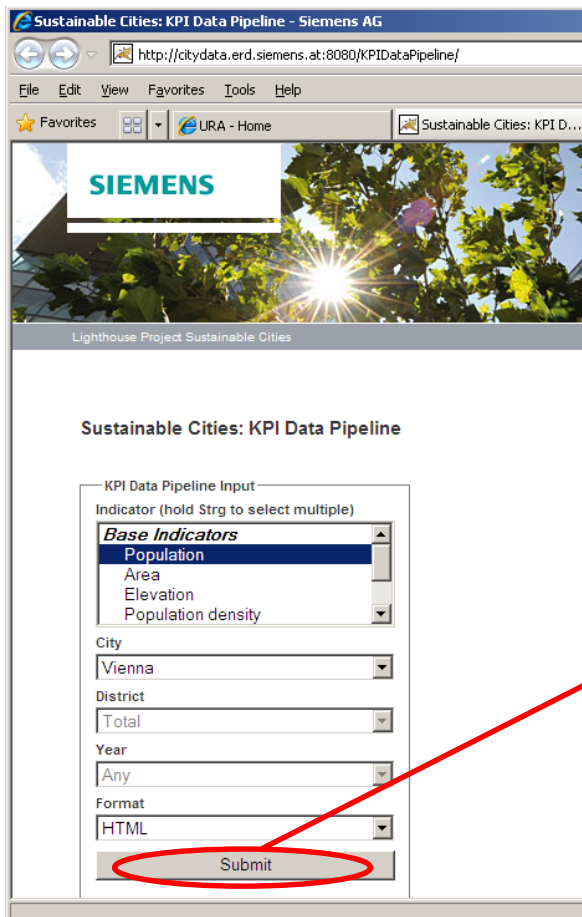
# City Data Pipeline: Web Interface

Our Web interface allows to browse data and download complex composed KPIs as Excel sheets (e.g. “Transport related CO2 emissions for Berlin”):

**1** Choose a city and set of base indicators for this city!

# City Data Pipeline: Web Interface

Our Web interface allows to browse data and download complex composed KPIs as Excel sheets (e.g. “Transport related CO2 emissions for Berlin”):



Sustainable Cities: Vienna - Population - Mozilla Firefox

citydata.erd.siemens.at:8080/KPIDataPipeline/KPIDispatcher?indicator=Population&city=http%3A%2F%2Fdbpedia.org%2...

SIEMENS

Lighthouse Project Sustainable Cities

**Vienna: Population**

**Population 2010**  
1712903 (Source: <http://dbpedia.org/>)

**Population 1996**  
1595402 (Source: <http://www.urbandaudit.org/>)

**Population 2001**  
1550123 (Source: <http://www.urbandaudit.org/>)

**Population 1991**  
1539848 (Source: <http://www.urbandaudit.org/>)

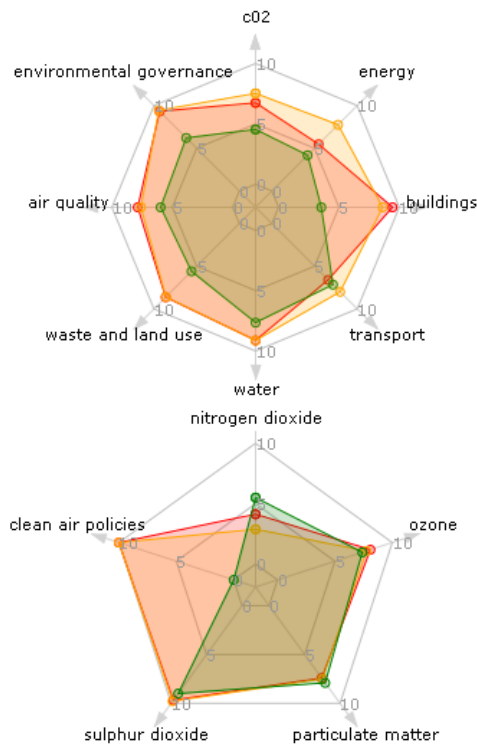
**Population 2004**  
1598626 (Source: <http://www.urbandaudit.org/>)

2 Browse available Open Data sources that contain the requested indicators

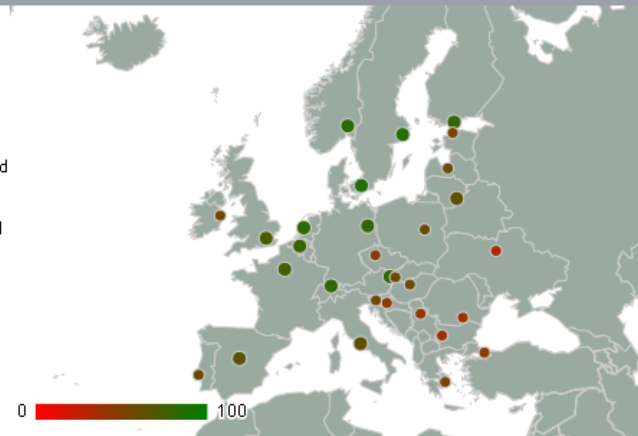
# City Data Pipeline: Web Interface



■ Vienna ■ Bratislava ■ Berlin

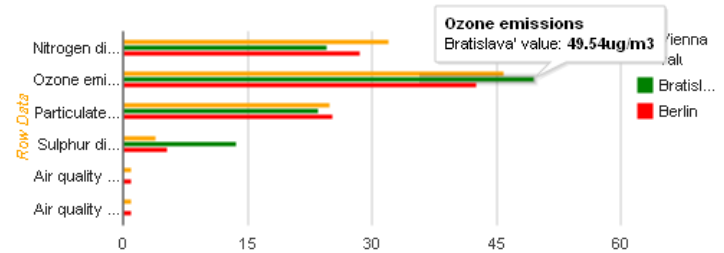


- c02
- energy
- buildings
- transport
- water
- waste and land use
- air quality
- environmental governance



3

Also available:  
Graphical user  
interface to visually  
compare base  
indicators for  
different cities.



# Applications & Challenges

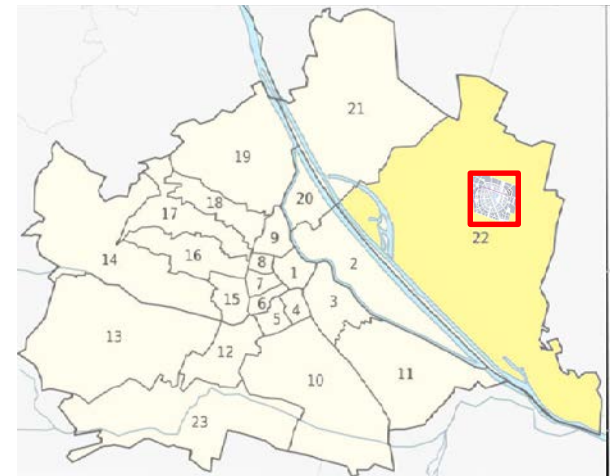
1. **Application: trend prediction - Example “Seestadt Aspern”**
2. Integration & enrichment of “Green City Index” Data
3. Challenges with Open Data Experienced

# Data Prediction/Quality: Statistical methods

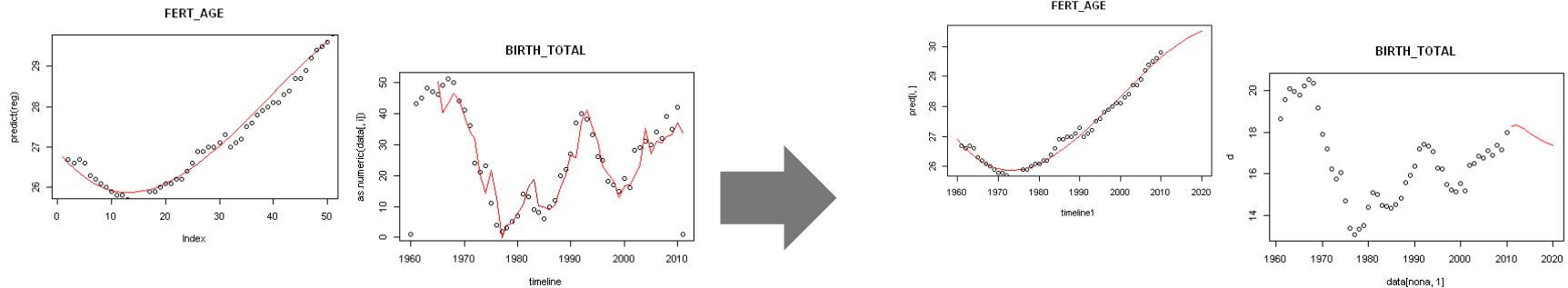
- Showcase: Estimate how Seestadt Aspern will be developing in comparison to the rest of Vienna
  - Semantic integration of open data in our data pipeline
  - Prediction of indicators for Aspern and Vienna
  - Graphical representation of the results

- **Questions:**

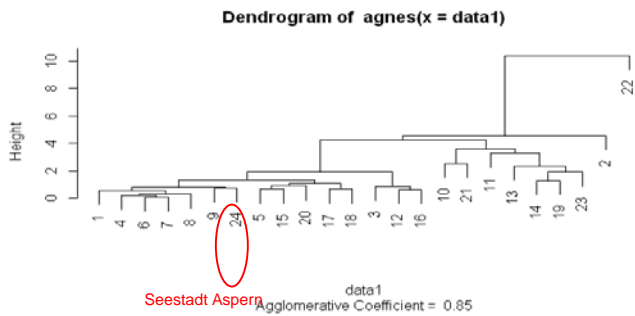
- *How will Seestadt Aspern perform in comparison to the other districts of Vienna?*
- *How do the goal indicators from the Aspern Masterplan compare with the “typical district behavior” of Vienna?*



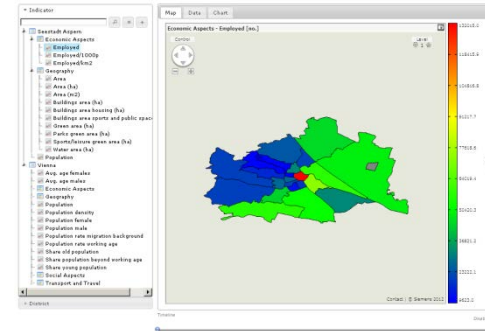




Completion & Prediction by statistical methods



finding patterns in the data



Vienna: districts and Seestadt Aspern



Berlin



Budapest



Vienna

Aim: Monitor development and compare to other cities/districts in order to take most effective infrastructural measures.

# Applications & Challenges

1. Application: trend prediction - Example “Seestadt Aspern”
- 2. Integration & enrichment of “Green City Index” Data**
3. Challenges with Open Data Experienced

# Collected Data vs. Green City Index

## Data: Overlaps

Together with colleagues from CC, we identified 20 quantitative raw data indicators that are overlapping between the Siemens' "Green City Index" and our current Data sources. The picture below visualizes the availability of data for these indicators for the cities of the European GCI:



European Green City Index

City / Raw Indicator	Gri	P	6	Ar	Lai	GE	GD	Oz	NC	PM	Ho	Joi	Le	Joi	Joi	rec	Le	Le	Le	Avc
Amsterdam	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Antwerp	X	X		X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Athens	X	X		X				X	X	X			X							X
Belgrade	X	X	X																	
Berlin	X	X		X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Bratislava	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Bremen	X	X		X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Brussels	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Bucharest	X	X	X	X	X	X	X	X	X	X										
Budapest	X	X		X	X	X	X	X	X	X			X			X	X	X	X	X
Cologne	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Copenhagen	X	X		X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Dublin	X	X	X		X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Essen	X	X		X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Frankfurt	X	X		X																X
Gothenburg	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Hamburg	X	X		X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Hanover	X	X		X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Helsinki	X	X		X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Istanbul	X																			
Kiev	X	X	X																	
Leipzig	X	X		X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Lisbon	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Ljubljana	X	X	X	X				X	X	X	X	X	X	X	X	X	X	X	X	X
London	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Luxembourg_(city)	X	X		X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Madrid	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Malm%C3%B6	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Munich	X	X		X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Nuremberg	X	X		X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Oslo	X	X																		X
Paris	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Prague	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Riga	X	X		X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Rome	X	X		X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Rotterdam	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Sofia	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Stockholm	X	X		X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Stuttgart	X	X		X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Tallinn	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
The_Hague	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Vienna	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Vilnius	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Warsaw	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Zagreb	X																			
Zurich	X	X	X																	

>65% of raw data could be covered by publically available data that we have collected automatically

- Data quality?**
- Not all indicators are 100% comparable (different scales, units, etc., sources of different quality)
  - for some indicators (e.g. Population) already less than 2% median error.
  - The more data we collect, the better the quality!

# Applications & Challenges

1. Application: trend prediction - Example “Seestadt Aspern”
2. Integration & enrichment of “Green City Index” Data
- 3. Challenges with Open Data Experienced**

# Challenges & Lessons Learnt – Is Open Data fit for industry?

Base assumption (for our use case):

Added value comes from **comparable** Open  
datasets being **combined**

# Challenges & Lessons Learnt – Is Open Data fit for industry?

- **Incomplete** Data: can be partially overcome
  - By ontological reasoning (RDF & OWL), by aggregation, or by rules & equations, e.g.
    - `:populationDensity = :population / :area`, cf. [ESWC2013]
  - *by statistical methods or Multi-dimensional Matrix Decomposition:*

$$\begin{array}{c} W \\ \boxed{\begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array}} \times \begin{array}{c} H \\ \boxed{\begin{array}{|c|c|c|c|c|} \hline & & & & \\ \hline \end{array}} \approx \begin{array}{c} V \\ \boxed{\begin{array}{|c|c|c|c|} \hline & & & \\ \hline & & & \\ \hline & & & \\ \hline \end{array}}
 \end{array}$$

unfortunately only partially successful, because these algorithms assume normally-distributed data.

- **Incomparable** Data:

dbpedia:populationTotal

dbpedia:populationCensus

- **Heterogeneity** across Open Government Data efforts:

- Different **Indicators**, Different Temporal and Spatial **Granularity**
- Different **Licenses** of Open Data: e.g. CC-BY, OGL (UK), etc.
- Heterogeneous **Formats** (CSV != CSV) ... Maybe the W3C CSV on the Web WG will solve this issue)

*→ Open Data needs stronger standards to be useful*