

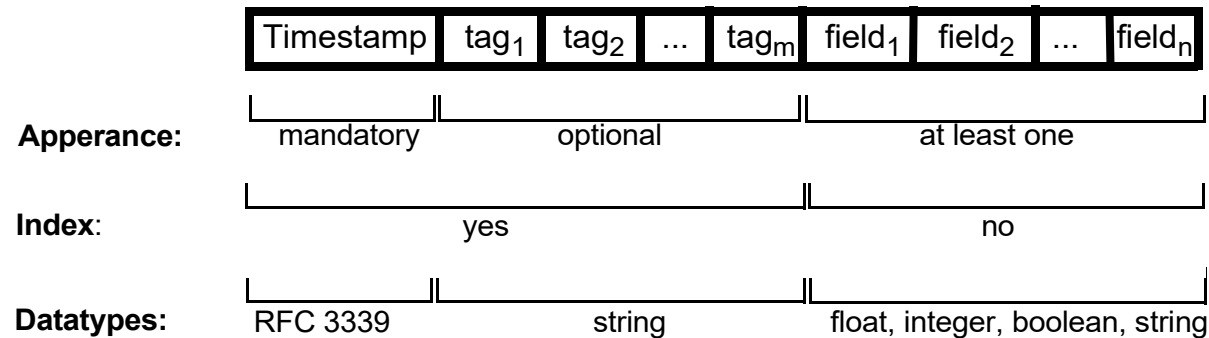
Metadata Extraction from a Huge Time Series Database

Andreas Schmidt, Mohamed Anis Koubaa, Nan Liu, Karl-Uwe Stucky, and Wolfgang Süß

Institute for Automation and Applied Informatics (IAI)

- Background
 - Influx Database at KIT
 - Zeitgeist
- Concept of MetadataExtractor
- Architecture
- First results
- Planned conceptual extensions
- Summary and Outlook

- Schema-free database, optimized to store time based datasets (precision: ns)
- Standalone/Cluster (sharding/replication)
- Up to 700.000 datasets/s for a single node
- data organized in *measurements* (comparable to a table)



- Retention policy, Continuous queries
- Query languages:
 - InfluxQL (SQL-like, limited range of functions)
 - Flux (stream-based)
- Programmatic interface: REST-API as well as multiple language bindings

- Easy publication of time series data in RO-Crate Format (data & metadata)
- Metadata comes from Influx Metadata API & manually specified (ORCID, ROR, licence, QUDT.org, description)
- Interactive specification of data to export/publish
 - restrictions on time range
 - restrictions on tag values
- Screenshots or live demo ...



Demo ...

- Actual limitation: Only one time series per export
- Goal: Bundle of multiple semantically related time series in one RO-Crate
- Question: How can we support the user to find related time series?

Zeitgeist - Export Configurator

Datasource:

Configuration: (force to read metadata from Influx DB[?])

database: fm_efficio_mirror
Measurement: efficio_raw
Time-range: [2020-12-31T23:15:00Z .. 2024-02-20T18:03:16Z]
Duration (days): 1145
Fields: displayDelta, displayEnd, displayStart, value

Export Definition:

Metadata:

ORCID:

Wolfgang Süß
KIT

License:

Server side parameter specification[?]:

Output format:

- RO-Crate 1.1
- CSV

Description :

Filters:

Time-Interval:

Start: (UTC)

End : (UTC)

Tag-Filter:

- SubCounter
- building

-
-
-
-

- deactivated
- energyType

-
-
-
-

- factor
- id
- installationPlace
- measurementUnit
- measuringValueType
- meteringType
- name
- parentName
- pointType
- raster
- unit

Export Definition:

Metadata:

ORCID:

[Wolfgang Süß](#)
[KIT](#)

License:

Server side parameter specification[?]:

Output format:

- RO-Crate 1.1
- CSV

Description [edited]:

bla bla bla ... metadata ohah !!!!

Filters:

Time-Interval:

Start: (UTC)

End : (UTC)

Tag-Filter:

- SubCounter
- building
 -
 -
 -
 -
- deactivated
- energyType
 -
 -
 -
 -
- factor
- id
- installationPlace
- measurementUnit
- measuringValueType
- meteringType
- name
- parentName
- pointType
- raster
- unit

Download URL:

https://localhost/zeitgeist/export.php/fm_efficio_mirror.fficio_raw.65d51da560add.zip

[\(config file\)](#)

InfluxQueryString:

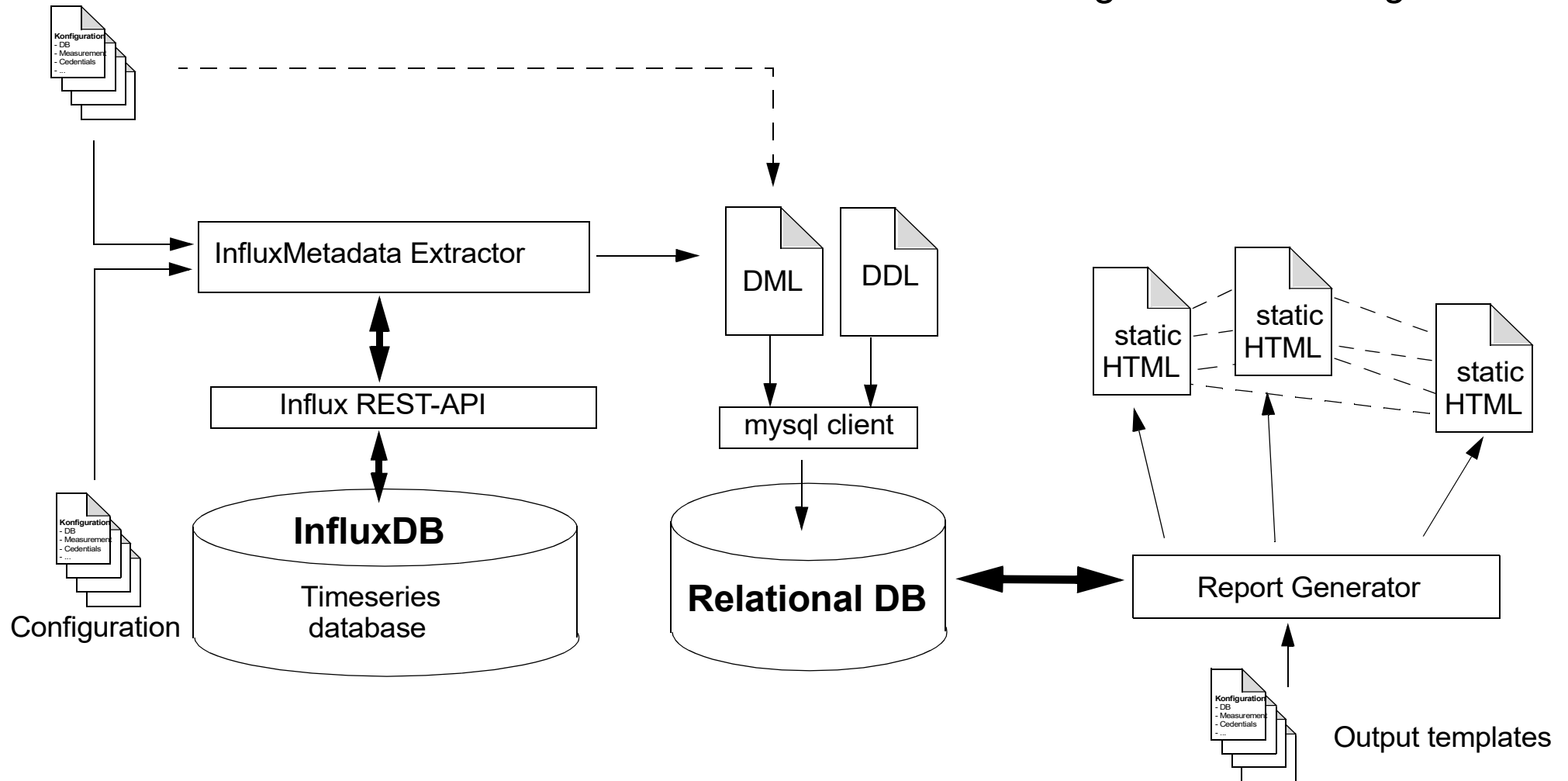
```
select *
from "efficio_raw"
where time >= '2024-02-20T17:03:16Z'
and time <= '2024-02-20T18:03:16Z'
and ("building" = '445' or "building" = '449')
and ("energyType" = 'Druckluft')
order by time
```

↓ ↑ 📁 > Dieser PC > Downloads > fm_efficio_mirror-efficio_raw-2024-02-20T17 03 16Z-2024-02-20T18 03 16Z.zip

Name	Typ	Komprimie...	Kennw...	Größe	Verhältnis	Änderungsdatum
 data.csv	Microsoft Excel-CSV-Datei	1 KB	Nein	2 KB	84%	20.02.2024 22:54
 ro-crate-metadata.json	JSON file	2 KB	Nein	9 KB	87%	20.02.2024 22:54
 ro-crate-preview.html	Firefox HTML Document	3 KB	Nein	16 KB	85%	20.02.2024 22:54

- For human inspection only
- No integration with *Zeitgeist*

Output templates



- 20 different databases
- temporal scope between 0 and 15 years
- Large number of modeling variants
 - up to 5000 measurements per database
 - up to 180 different fields per measurement
 - up to 26 tags per measurement

- Extreme examples:
 - measurement name combined with UUID
 - one measurement for a whole working group
 - for a single dataset: only one (of many fields) contains a value

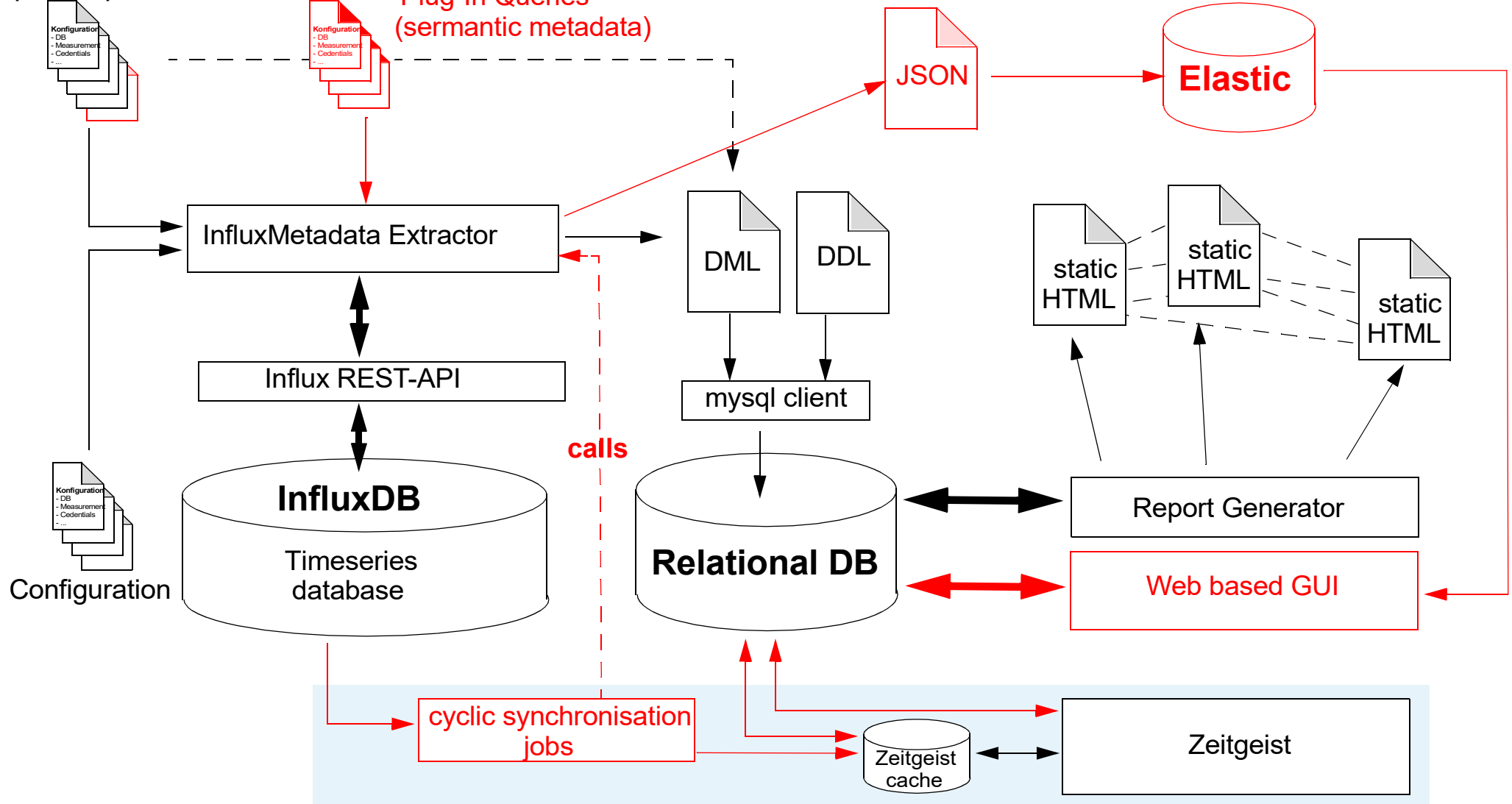
- 20 different databases
- temporal scope between 0 and 15 years
- Large number of modeling variants
 - up to 5000 measurements per database
 - up to 180 different fields per measurement
 - up to 26 tags per measurement
- Examples:
 - measurement name combined with UUID
 - one measurement for a whole working group
 - for a single dataset: only one (of many fields) contains a value

**direct consequence of schema free database
nevertheless, there are recommendations on
how to build schemas**

- Searching for relationships between data records using the data structure of a measurement is difficult
 - several unrelated time series in one measurement
 - related time series in multiple measurements
 - probably evolution of data structures over time
- More sophisticated algorithms needed:
 - search for tags, fields that occur together in a dataset
 - examination of the recorded interval
 - Look for similar values of tags (i.e. a sensor name, location)
 - Integrate the user experience (domain expert)
 - ...

Question: can a „relatedness measure“ be calculated?

Output templates



- Influx Meta API
 - Metadata calls are typically fast (< 0.05 sec.), but there can be exceptions:
 - Exceptions:
 - Extraction of field names can be slow (up to 30 s) - depending on number of shards
 - Extraction of aggregates (min/max/count) can be expensive (up to 5 s) on many (~ 100) fields
 - But: data structures typically change slowly
 - Cyclic job that queries the Influx DB and updates metainformation in relational database (+ Zeitgeist metadata cache)

- Tool for extraction of meta information from a schema-free time series database
- Primary goal: search for related/similar data (based on structure) for our *Zeitgeist* tool
- The results are currently presented as static web pages in which you can navigate
- Difficult to search for relationships between data sets based on the structure of a measurement
- In order to be used meaningfully in the context of *Zeitgeist*, more in-depth analyses than the analysis of Influx Meta-API results are necessary
- Secondary goal: Stand-alone tool to analyze the structure of a schema free InfluxDB

- Next version of MetadataExtractor will cover:
 - Data set related structural analysis (incl. structure migration)
 - Considering the recording times
 - Integration of domain expert knowledge
 - Interactive frontend with search capabilities
 - Near realtime statistics (cyclic analysis)
- Further ideas
 - calculate measure about modelling quality