

# Hands-on data management with open-source software: CaosDB

Florian Spreckelsen, Daniel Hornung  
*IndiScale GmbH*

2023-03-23



# CaosDB

## History

- CaosDB started at MPI-DS (Göttingen) around 2011
- Running stable since ca. 2016, [released as open-source \(AGPLv3\)](#)<sup>[1]</sup> in 2018
- Increasing user base since 2020
- Commercial support by IndiScale GmbH
  - distribution branded as LinkAhead



## History

- CaosDB started at MPI-DS (Göttingen) around 2011
- Running stable since ca. 2016, [released as open-source \(AGPLv3\)](#)<sup>[1]</sup> in 2018
- Increasing user base since 2020
- Commercial support by IndiScale GmbH
  - distribution branded as LinkAhead



[1] <https://gitlab.com/caosdb>

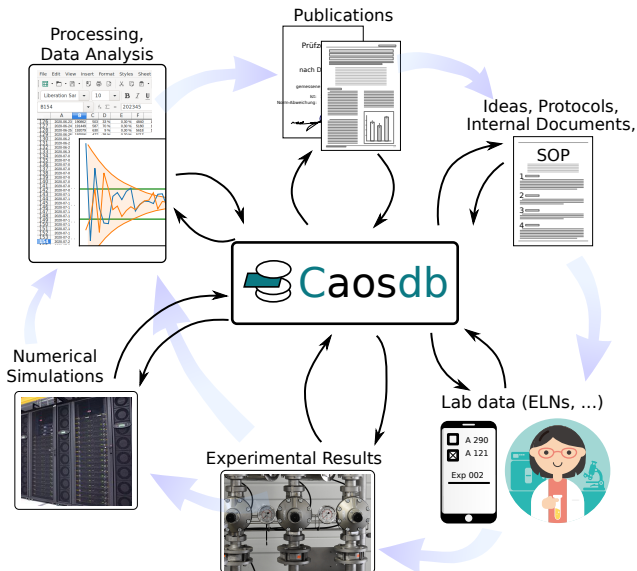
## History

- CaosDB started at MPI-DS (Göttingen) around 2011
- Running stable since ca. 2016, [released as open-source \(AGPLv3\)](#)<sup>[1]</sup> in 2018
- Increasing user base since 2020
- Commercial support by IndiScale GmbH
  - distribution branded as LinkAhead

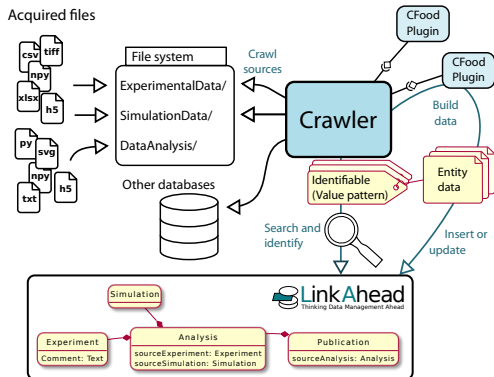
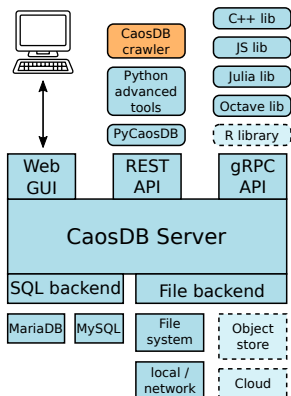


[1] <https://gitlab.com/caosdb>

# Data Life Cycle



# Architecture & Crawler LinkAhead



## Useful Links

- Documentation: <https://docs.indiscale.com/>
- Gitlab repository: <https://gitlab.com/caosdb/>

## Live-Demo

<https://demo.indiscale.com>

The screenshot shows the LinkAhead web application interface. At the top, there is a navigation bar with the LinkAhead logo and links for Entities, File System, Query, Map, and Tour. On the right side of the navigation bar, there are links for Bookmarks and Login. The main content area has a dark blue background with a geometric pattern. A white box in the center contains the text "Welcome to LinkAhead" and "Thinking Data Management Ahead". Below this, there is a paragraph explaining the flexibility of LinkAhead and a note about the demo server. Two buttons, "Start a Tour" and "Contact", are provided. To the right of the text, there is a diagram showing a central person icon connected to five circular icons representing different data management concepts: a gear, a balance scale, a person, a bar chart, and a factory. The footer of the page contains the LinkAhead logo and tagline, the IndiScale logo and tagline, and a row of links: Contact, Imprint/Impressum, Data Policy, License (AGPL-v3), and Sources.

LinkAhead Entities File System Query Map Tour Bookmarks Login

### Welcome to LinkAhead

#### Thinking Data Management Ahead

With LinkAhead, you can manage your data in a more flexible way than ever. Experience what it's like to think data management ahead!

You can explore LinkAhead yourself or take the interactive tour. Login as **admin** with password **caosdb**.

[Start a Tour](#) [Contact](#)

Note that any data you enter on this demo server is public and that changes are reset daily.

LinkAhead – Thinking data management ahead. IndiScale – We make individual data management scalable.

[Contact](#) • [Imprint/Impressum](#) • [Data Policy](#) • [License \(AGPL-v3\)](#) • [Sources](#)





## Presentation

- Main Menu



## Task

- Try the different buttons of the “Entities” Menu.

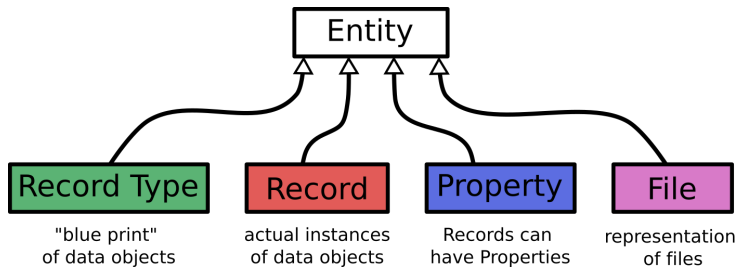


## Note

- By default only 10 Entities are shown on one page. You can get to other pages with the “Next Page” and “Previous Page” buttons.

⇒ What are the differences between the options of the “Entities” menu?

# The Data Model



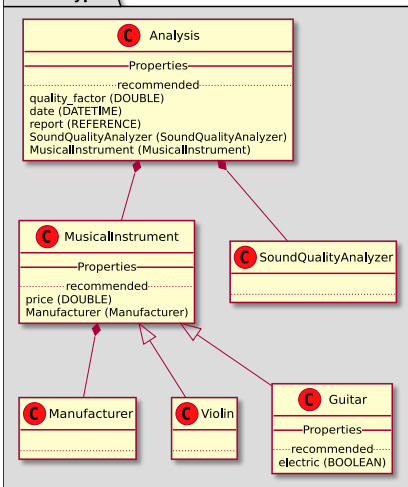
## 💡 Note

RecordTypes and Properties define the ontology. Records store the actual data and represent things, e.g., a particular measurement.

For more details, read on at [docs.indiscale.com/caosdb-server/Data-Model.html](https://docs.indiscale.com/caosdb-server/Data-Model.html).

# The Data Model

## RecordTypes



References in LinkAhead are directed.  
If Record **A** references Record **B**, then:

- The referencing Record **A** has a corresponding Property.
- The referenced Record **B** does not.

Finding “back-references” in the web interface:

← References

- The semantic data model allows efficient data access
- ChaosDB Query Language (CQL) is used to search data

## Presentation

- Simple Queries

## Note

- Try to guess what the result will be, before actually executing the query.

Documentation on the query language:

<https://docs.indiscale.com/caosdb-server/CaosDB-Query-Language.html>

# Presented Queries

- FIND RECORD MusicalInstrument
- FIND Guitar
- FIND RecordType Guitar
- FIND RECORD 'Nice Guitar'
- FIND RECORD Nice\*
- FIND Manufacturer
- FIND 'Budget Bikes'

# Searching Data: Using Properties

The power of CQL (CaosDB Query Language): using properties and references to refine search results

Find Queries:

- `FIND <Role> <Name> [Property Filter]`



## Presentation

- FIND Queries with Property filter



## Note

- Typically, the `Property Filter` has the form  
`<Property> <Operator> <Value>`, for example  
`length >= 0.7mm`.

See the [CQL-tutorial](#) for details on the possible filters.

# Presented Queries

- FIND RECORD MusicalInstrument WITH Manufacturer
- FIND RECORD MusicalInstrument WITH price=48
- FIND RECORD WITH price>48
- FIND RECORD WHICH HAS A PROPERTY temperature
- FIND RECORD WITH Manufacturer LIKE Budget\*
- FIND RECORD WITH date in 2019
- FIND RECORD WHICH REFERENCES A Guitar

# Restricting the details of results

- Using `COUNT` instead of `FIND` will only return the number of entities in the result set

## Note

This is often useful when experimenting with queries.

- Using `SELECT [Property Name], ... FROM` instead of `FIND` returns specific information in a table
- A comma-separated list of Property names can be provided behind the `SELECT` keyword.

## Presentation

### SELECT Queries



# Presented Queries

- `SELECT price, electric FROM Guitar`
- `SELECT name,  
quality_factor, report, date, Manufacturer.latitude  
FROM MusicalAnalysis WHICH REFERENCES A  
MusicalInstrument WITH electric=TRUE`

For more, read on at

<https://docs.indiscale.com/caosdb-server/CaosDB-Query-Language.html> and  
<https://docs.indiscale.com/caosdb-webui/tutorials/query.html>.

# Automating LinkAhead: The Python library

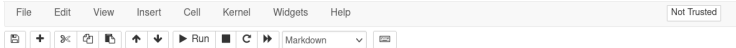
# This Hacky Hour's Jupyter Notebook

Find the Jupyter Notebook for this hacky hour on  
<https://gitlab.indiscale.com/caosdb/src/notebooks-hacky-hour-dd23>

# Connect to existing Data Analysis Software

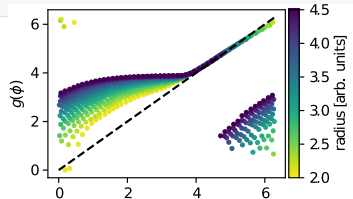
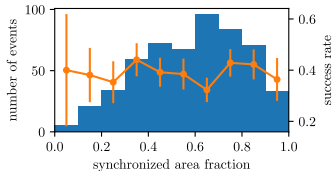


jupyter Introduction to the CaosDB-Python Client (autosaved)



## Data Analysis with Jupyter

```
In [ ]: import caosdb as db
data = db.execute_query("SELECT quality_factor FROM RECORD Analysis with quality_factor" )
table = to_table(data)
plt.plot(table.x, table.y)
```





- Online demo:
  - <https://demo.indiscale.com>
- Source code and development: <https://gitlab.com/caosdb>
- CaosDB community chat: [#caosdb:matrix.org](https://matrix.org/#caosdb)
- Documentation: <https://docs.indiscale.com>
- Scientific article, published in **Data**:  
<https://doi.org/10.3390/data4020083> 



- Fully open-source software (AGPLv3)



- Fully open-source software (AGPLv3)
- Powerful search options



- Fully open-source software (AGPLv3)
- Powerful search options
- CaosDB can be easily extended or modified





- Fully open-source software (AGPLv3)
- Powerful search options
- CaosDB can be easily extended or modified
- Integrates into existing workflows



- Fully open-source software (AGPLv3)
- Powerful search options
- CaosDB can be easily extended or modified
- Integrates into existing workflows
- Stores semantic links together with the data

Thank you for your interest!

- Comparison to SPARQL (RDF query language for e.g. WikiData)
- Data model evolution (at the AWI)

# Comparison to SPARQL

## SPARQL

```
SELECT DISTINCT ?item ?itemLabel ?givenName ?familyName WHERE {  
    ?item wdt:P31 wd:Q5; # Any instance of a human.  
    wdt:P27 wd:Q145; # United Kingdom  
    wdt:P21 wd:Q6581072; # female  
    wdt:P106 wd:Q36180; # writer  
    wdt:P569 ?birthday;  
    wdt:P570 ?diedon;  
    wdt:P734 [rdfs:label ?familyName];  
    wdt:P735 [rdfs:label ?givenName].  
    FILTER(?birthday > "1870-01-01"^^xsd:dateTime  
        && ?diedon < "1950-01-01"^^xsd:dateTime)  
    FILTER(regex(?givenName, "M.*") || regex(?familyName, "M.*"))  
    SERVICE wikibase:label { bd:serviceParam wikibase:language "en" }  
}
```

---

# Comparison to SPARQL

## SPARQL

```
SELECT DISTINCT ?item ?itemLabel ?givenName ?familyName WHERE {  
  ?item wdt:P31 wd:Q5; # Any instance of a human.  
  wdt:P27 wd:Q145; # United Kingdom  
  wdt:P21 wd:Q6581072; # female  
  wdt:P106 wd:Q36180; # writer  
  wdt:P569 ?birthday;  
  wdt:P570 ?diedon;  
  wdt:P734 [rdfs:label ?familyName];  
  wdt:P735 [rdfs:label ?givenName].  
  FILTER(?birthday > "1870-01-01"^^xsd:dateTime  
    && ?diedon < "1950-01-01"^^xsd:dateTime)  
  FILTER(regex(?givenName, "M.*") || regex(?familyName, "M.*"))  
  SERVICE wikibase:label { bd:serviceParam wikibase:language "en" }  
}
```

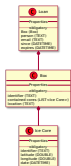
---

## CaosDB query language

```
SELECT given_name, family_name FROM Writer  
WITH gender=f AND country=UK AND birthday > 1870 AND death < 1950  
AND (given_name LIKE "M*" OR family_name LIKE "M*")
```

# Data model evolution

Data model  
ca. 2018

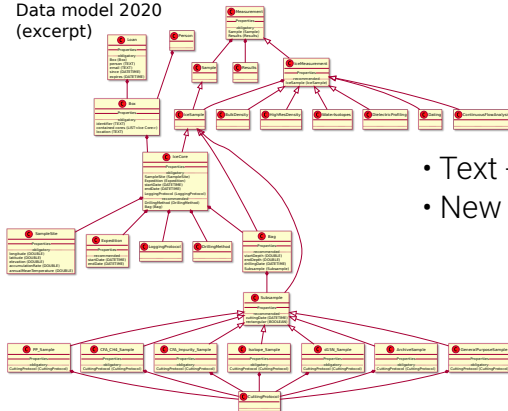


# Data model evolution

Data model  
ca. 2018



Data model 2020  
(excerpt)



- Text -> object entities
- New RecordTypes