

Combining file-based with RDMS-based scientific workflows using the LinkAhead-crawler

Henrik tom Wördén, Florian Spreckelsen, Stefan Luther,
Ulrich Parlitz, Alexander Schlemmer

2024-03-06

Mapping Hierarchical File Structures to Semantic Data Models for Efficient Data Integration into Research Data Management Systems

by  **Henrik tom Wörden**  ,  **Florian Spreckelsen**  ,  **Stefan Luther**  ,
 **Ulrich Parlitz**   and  **Alexander Schlemmer**  

1 Indiscale GmbH, 37083 Göttingen, Germany

2 Max Planck Institute for Dynamics and Self-Organization, 37077 Göttingen, Germany

3 Institute for the Dynamics of Complex Systems, Georg-August-Universität, 37077 Göttingen, Germany

4 German Center for Cardiovascular Research (DZHK), Partner Site Göttingen, 37075 Göttingen, Germany

5 Institute of Pharmacology and Toxicology, University Medical Center Göttingen, 37075 Göttingen, Germany

* Author to whom correspondence should be addressed.

Data 2024, 9(2), 24; <https://doi.org/10.3390/data9020024>

Submission received: 15 August 2023 / Revised: 8 December 2023 / Accepted: 18 December 2023 /

Published: 26 January 2024

(This article belongs to the Section Information Systems and Data Management)

tom Wörden et. al., *Data* 2024, 9, 24.

<https://doi.org/10.3390/data9020024>

Research Question

- ▶ Data integration (e.g. getting data from the disk into databases / RDMS) time-consuming, especially with heterogeneous file layouts and data formats. How can we simplify that?

Research Question

- ▶ Data integration (e.g. getting data from the disk into databases / RDMS) time-consuming, especially with heterogeneous file layouts and data formats. How can we simplify that?
- ▶ How can we use the RDMS simultaneously to traditional file systems (for interoperability)?

LinkAhead (formerly: CaosDB)

The screenshot displays the LinkAhead application interface, which is a web-based tool for managing experimental data and workflows. It features two main sections: 'DataAnalysis' and 'Run'.

DataAnalysis Section:

- date:** 2023-02-17
- identifier:** BlenderExample
- project:** 190
- runs:** 190, 194, 193, 192
- environment:** blender.sif, blender_vis.blend
- input_data:** prepared.hdfs

Run Section:

- Project:** 190
- date:** 2023-02-17T13:16:19+0000
- commandline:** commandline.sh
- output:** slurm-14873163.out, wave_video_singularity0001-0512.avi
- environment:** blender.sif
- input_data:** blender_vis.blend, prepared.hdfs
- input_parameters:** parameters.yaml
- script:** script.py
- hashsums:** sha256sums.txt

At the bottom right of the interface, there are navigation icons for back, forward, search, and refresh.

LinkAhead (formerly: CaosDB)

The screenshot displays the LinkAhead application interface, which has been renamed from CaosDB. The interface consists of two main panels, each with a header bar and a configuration table.

Top Panel (Header Bar):

- Icon: Disk icon.
- Title: LinkAhead (formerly: CaosDB).
- Buttons: R (red), DataAnalysis, References, Bookmarks, Help.

Configuration Table (Top Panel):

date	2023-02-17
identifier	BlenderExample
project	190
runs	194 193 192
environment	blender.sif blender_vis.blend
input_data	prepared.hdfs

Bottom Panel (Header Bar):

- Icon: Run icon.
- Buttons: R (red), Run, References, Bookmarks, Help.

Configuration Table (Bottom Panel):

Project	190
date	2023-02-17T13:16:19+0000
commandline	commandline.sh
output	slurm-14873163.out wave_video_singularity0001-0512.avi
environment	blender.sif
input_data	blender_vis.blend prepared.hdfs
input_parameters	parameters.yaml
script	script.py
hashsums	sha256sums.txt

Bottom Right Icons:

- Back, Forward, Home, Search, Refresh.

LinkAhead (formerly: CaosDB)

LinkAhead
~~CaosDB~~

FIND DataAnalysis with date in 2023 and identifier like *Blender*

DataAnalysis

<input checked="" type="checkbox"/> date	2023-02-17
<input checked="" type="checkbox"/> identifier	BlenderExample
<input checked="" type="checkbox"/> project	190
<input checked="" type="checkbox"/> runs	190 194 193 192
<input checked="" type="checkbox"/> environment	blender.sif blender_vis.blend
<input checked="" type="checkbox"/> input_data	prepared.hdfs

Run

<input checked="" type="checkbox"/> Project	190
<input checked="" type="checkbox"/> date	2023-02-17T13:16:19+0000
<input checked="" type="checkbox"/> commandline	commandline.sh
<input checked="" type="checkbox"/> output	slurm-14873163.out wave_video_singularity0001-0512.avi
<input checked="" type="checkbox"/> environment	blender.sif
<input checked="" type="checkbox"/> input_data	blender_vis.blend prepared.hdfs
<input checked="" type="checkbox"/> input_parameters	parameters.yaml
<input checked="" type="checkbox"/> script	script.py
<input checked="" type="checkbox"/> hashsums	sha256sums.txt

References

◀ ▶ ⏪ ⏩ 🔍

LinkAhead (formerly: CaosDB)

LinkAhead
~~CaosDB~~

FIND DataAnalysis with date in 2023 and identifier like *Blender*

R DataAnalysis

date 2023-02-17
identifier BlenderExample
project 190
runs 190, 194, 193, 192
environment blender.sif, blender_vis.blend
input_data prepared.hdfs

R Run

Project 190
date 2023-02-17T13:16:19+0000
commandline commandline.sh
output slurm-14873163.out, wave_video_singularity0001-0512.avi
environment blender.sif
input_data blender_vis.blend, prepared.hdfs
input_parameters parameters.yaml
script script.py
hashsums sha256sums.txt

LinkAhead (formerly: CaosDB)

LinkAhead
~~CaosDB~~

FIND Run which is referenced by
DataAnalysis with date in 2023 and identifier like *Blender*

The screenshot displays the LinkAhead interface, which has replaced the CaosDB system. It consists of two main sections: a search interface at the top and a detailed run view at the bottom.

Search Interface (Top):

- Search bar: FIND Run which is referenced by DataAnalysis with date in 2023 and identifier like *Blender*
- Filter buttons: R (Run) and DataAnalysis
- Results table:
 - date: 2023-02-17
 - identifier: BlenderExample
 - project: 190
 - runs: 194, 93, 192
 - environment: blender.sif, blender_vis.blend
 - input_data: prepared.hdfs

Detailed Run View (Bottom):

- Filter buttons: R (Run) and Project
- Results table:
 - Project: 190
 - date: 2023-02-17T13:16:19+0000
 - commandline: commandline.sh
 - output: slurm-14873163.out, wave_video_singularity0001-0512.avi
 - environment: blender.sif
 - input_data: blender_vis.blend, prepared.hdfs
 - input_parameters: parameters.yaml
 - script: script.py
 - hashsums: sha256sums.txt

LinkAhead (formerly: CaosDB)

- ▶ Semantic Research Data Management System (RDMS)
Fitschen et.al., *Data* 2019, 10.3390/data4020083
- ▶ Developed at the Max Planck Institute for Dynamics and Self-Organization (Göttingen) since 2010
- ▶ Open Source project since 2018: gitlab.com/linkahead
- ▶ Commercial support available by IndiScale GmbH¹ (since 2019)
- ▶ Currently approximately 20-30 instances in very different scientific disciplines
- ▶ Main features (very briefly): Flexible semantic data model, highly modular data crawler, semantic search

¹A.S. is a co-founder of IndiScale.

Data Integration

```
example/
└ runs/
    └ 2023-02-17T13_16/
        ├── prepared.hdf5
        ├── sha256sums.txt
        ├── parameters.yaml
        ├── commandline.sh
        └── slurm-14873163.out
└ script.py
└ blender.sif
```

Data Integration

```
example/
  + runs/
    + 2023-02-17T13_16/
      + prepared.hdf5
      + sha256sums.txt
      + parameters.yaml
      + commandline.sh
      + slurm-14873163.out
  + script.py
  + blender.sif
```

The screenshot shows the CaosDB application interface with two main sections: 'DataAnalytics' and 'Run'.

DataAnalytics:

- date: 2023-02-17
- identifier: BlenderExample
- project: T90
- runs: 194, 193, 192
- environment: blender.sif, blender_vs_blend
- input_data: prepared.hdf5

Run:

- Project: T90
- date: 2023-02-17T13:16:19+0000
- commandline: commandline.sh
- output: slurm-14873163.out, wave_video_singularity0001-0512.avi
- environment: blender.sif
- input_data: blender_vs_blend, prepared.hdf5
- input_parameters: parameters.yaml
- script: script.py
- hashsums: sha256sums.txt

Data Integration

```
example/
  + runs/
    + 2023-02-17T13_16/
      + prepared.hdf5
      + sha256sums.txt
      + parameters.yaml
      + commandline.sh
      + slurm-14873163.out
  + script.py
  + blender.sif
```



The image shows a screenshot of a software interface titled "CaosDB". On the left, there is a tree view of a directory structure. On the right, there are two tables showing data entries.

Table 1 (Top):

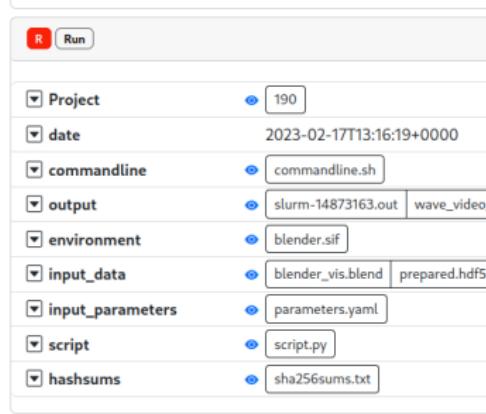
<input checked="" type="checkbox"/> date	2023-02-17
<input checked="" type="checkbox"/> identifier	BlenderExample
<input checked="" type="checkbox"/> project	190
<input checked="" type="checkbox"/> runs	194, 193, 192
<input checked="" type="checkbox"/> environment	blender.sif blender_vs_blend
<input checked="" type="checkbox"/> input_data	prepared.hdf5

Table 2 (Bottom):

<input checked="" type="checkbox"/> Project	190
<input checked="" type="checkbox"/> date	2023-02-17T13:16:19+0000
<input checked="" type="checkbox"/> commandline	commandline.sh
<input checked="" type="checkbox"/> output	slurm-14873163.out wave_video_singularity0001-0512.avi
<input checked="" type="checkbox"/> environment	blender.sif
<input checked="" type="checkbox"/> input_data	blender_vs_blend prepared.hdf5
<input checked="" type="checkbox"/> input_parameters	parameters.yaml
<input checked="" type="checkbox"/> script	script.py
<input checked="" type="checkbox"/> hashsums	sha256sums.txt

Data Integration

```
example/
└ runs/
  └ 2023-02-17T13_16/
    ├── prepared.hdf5
    ├── sha256sums.txt
    ├── parameters.yaml
    ├── commandline.sh
    └── slurm-14873163.out
  script.py
  blender.sif
```

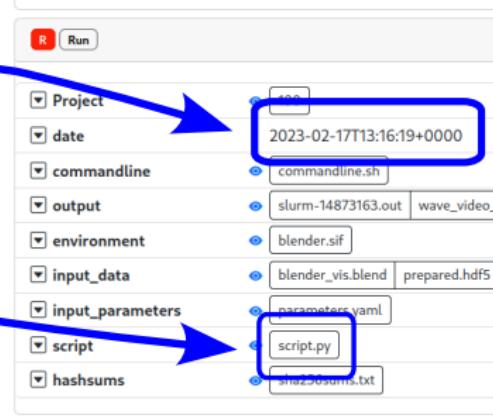


The screenshot shows a software interface with a toolbar at the top containing a red 'Run' button and a grey 'Run' button. Below the toolbar is a table with the following data:

Project	190
date	2023-02-17T13:16:19+0000
commandline	commandline.sh
output	slurm-14873163.out wave_video.
environment	blender.sif
input_data	blender_vis.blend prepared.hdf5
input_parameters	parameters.yaml
script	script.py
hashsums	sha256sums.txt

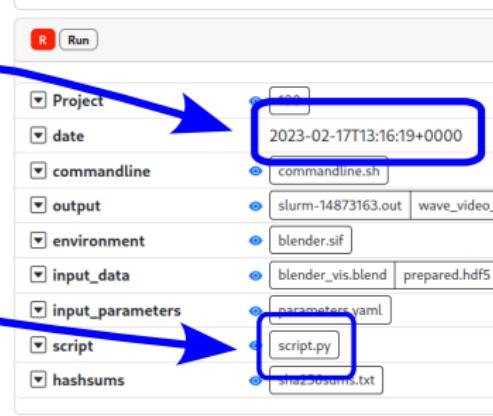
Data Integration

```
example/
  runs/
    2023-02-17T13_16/
      prepared.ndts
      sha256sums.txt
      parameters.yaml
      commandline.sh
      slurm-14873163.out
    script.py
    blender.sif
```



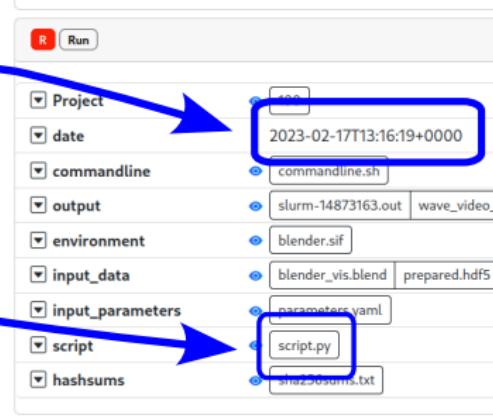
Data Integration

```
example/
  runs/
    2023-02-17T13_16/
      prepared.ndts
      sha256sums.txt
      parameters.yaml
      commandline.sh
      slurm-14873163.out
    script.py
    blender.sif
```



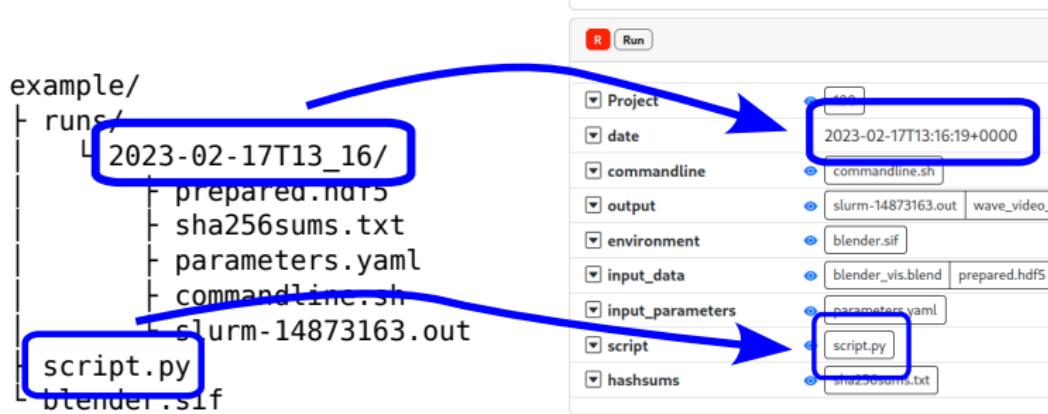
Data Integration

```
example/
  runs/
    2023-02-17T13_16/
      prepared.ndts
      sha256sums.txt
      parameters.yaml
      commandline.sh
      slurm-14873163.out
    script.py
    blender.sif
```



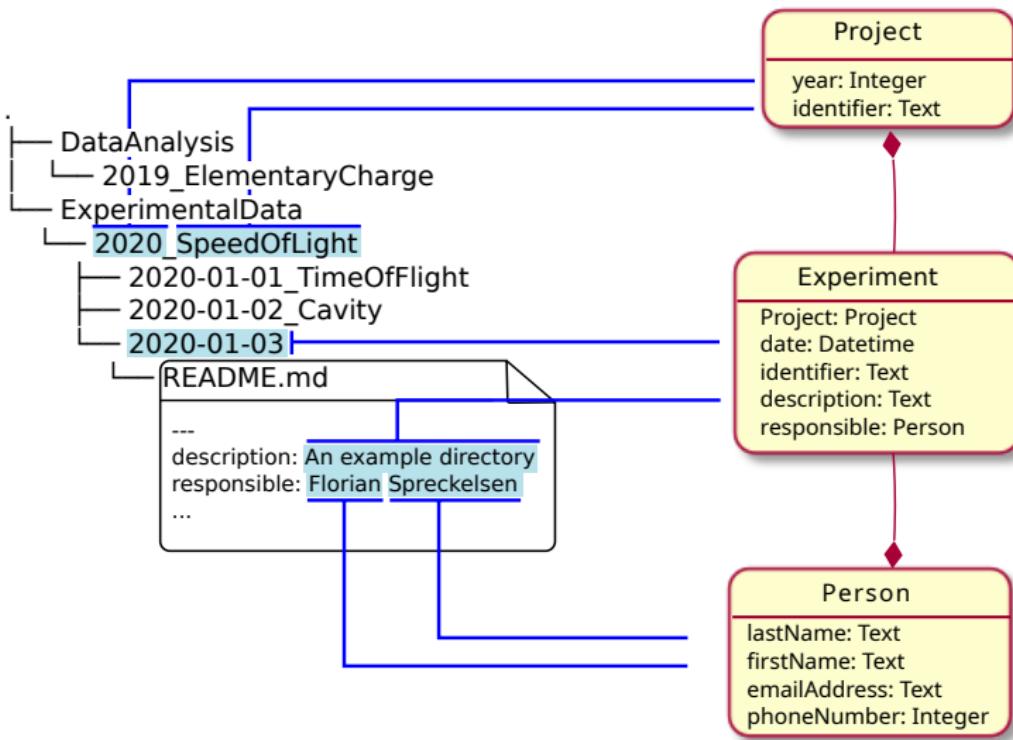
- ▶ Extract (meta) data from file names, directory names, and file contents

Data Integration



- ▶ Extract (meta) data from file names, directory names, and file contents
- ▶ Map these data onto a semantic data model (and store it in CaosDB/LinkAhead)

Example 2



Crawler Design

- ▶ Modular file crawler (Python, AGPLv3)

Crawler Design

- ▶ Modular file crawler (Python, GPLv3)
- ▶ Adaptable to file structure using YAML specification

Crawler Design

- ▶ Modular file crawler (Python, GPLv3)
- ▶ Adaptable to file structure using YAML specification
- ▶ Extendable with Python modules (custom data formats, ...)

Crawler Design

- ▶ Modular file crawler (Python, GPLv3)
- ▶ Adaptable to file structure using YAML specification
- ▶ Extendable with Python modules (custom data formats, ...)
- ▶ Synchronization into RDMS (instead of data ingest): Files and RDMS can be used simultaneously!

YAML specification

```
ExperimentalData_Dir:  
    type: Directory  
    match: "ExperimentalData"  
    subtree:  
        Project_Dir:  
            type: Directory  
            match: (?P<year>[0-9]{4,4})_(?P<name>.*)  
            records:  
                Project:  
                    year: $year  
                    name: $name  
# (...)
```

YAML specification

```
ExperimentalData_Dir:  
    type: Directory  
    match: "ExperimentalData"  
    subtree:  
        Project_Dir:  
            type: Directory  
            match: (?P<year>[0-9]{4,4})_(?P<name>.*)  
            records:  
                Project:  
                    year: $year  
                    name: $name  
# (...)
```

Converter matching folders with name
ExperimentalData

YAML specification

```
ExperimentalData_Dir:  
  type: Directory  
  match: "ExperimentalData"  
  subtree:  
    Project_Dir:  
      type: Directory  
      match: (?P<year>[0-9]{4,4})_(?P<name>.*)  
      records:  
        Project:  
          year: $year  
          name: $name  
# (...)
```

Converter for subfolders
using a regexp, e.g.: **2024_deRSE**

YAML specification

```
ExperimentalData_Dir:  
  type: Directory  
  match: "ExperimentalData"  
  subtree:  
    Project_Dir:  
      type: Directory  
      match: (?P<year>[0-9]{4,4})_(?P<name>.*)  
      records:  
        Project:  
          year: $year  
          name: $name
```



```
# (...)
```

For each subfolder, a record of type **Project** is created, setting 2 properties **year** and **name** from the regexp.

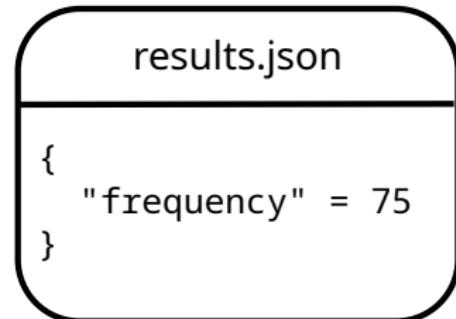
Extract contents from files

```
# (...)

DataFile:
    type: JSONFile
    match: results\.json
    subtree:
        frequency:
            type: DictTextElement
            match_name: frequency
            match_value: (?P<frequency>[0-9]+)
            records:
                DataAnalysis:
                    frequency: $frequency
# (...)
```

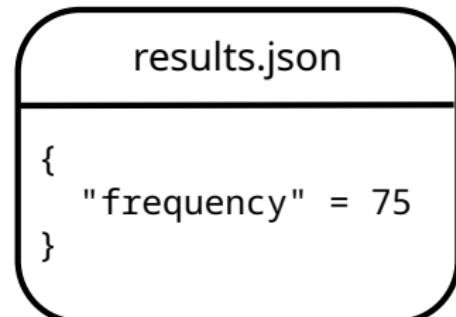
Extract contents from files

```
# (...)  
DataFile:  
    type: JSONFile  
    match: results\.json  
    subtree:  
        frequency:  
            type: DictTextElement  
            match_name: frequency  
            match_value: (?P<frequency>[0-9]+)  
            records:  
                DataAnalysis:  
                    frequency: $frequency  
# (...)
```



Extract contents from files

```
# (...)  
DataFile:  
    type: JSONFile  
    match: results\.json  
    subtree:  
        frequency:  
            type: DictTextElement  
            match_name: frequency  
            match_value: (?P<frequency>[0-9]+)  
            records:  
                DataAnalysis:  
                    frequency: $frequency  
# (...)
```



Create a record of type **DataAnalysis**
with property **frequency** = 75

Thank You! - Summary

- ▶ Crawler that continuously synchronizes (currently one-way) information from file system to semantic RDMS
- ▶ YAML based, extendable specification for adapting crawler to different file system layouts and data formats
- ▶ tom Wörden et. al., *Data* 2024, 9, 24.
<https://doi.org/10.3390/data9020024>
- ▶ gitlab.com/linkahead AND
gitlab.com/linkahead/linkahead-crawler