# CaosDB - An open scientific database

**Alexander Schlemmer**, Henrik tom Wörden, Timm Fitschen, Daniel Hornung, Ulrich Parlitz, and Stefan Luther

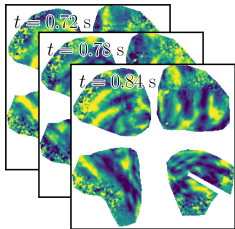Research Group Biomedical Physics, MPI for Dynamics and Self-Organization, Göttingen
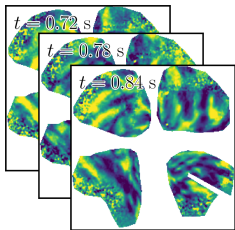
IndiScale GmbH

2019-09-23



Caosdb
an open scientific database

# Current Challenges in Scientific Data Management

**Diverse and Complex Data**

# Current Challenges in Scientific Data Management



**Diverse and Complex Data**

# Current Challenges in Scientific Data Management



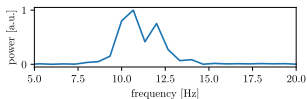**Diverse and Complex Data**

# Current Challenges in Scientific Data Management
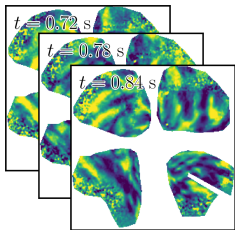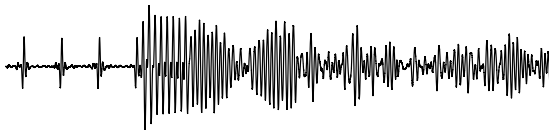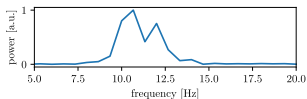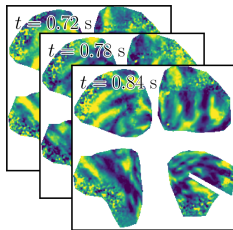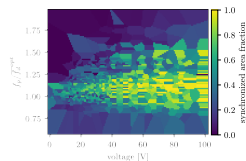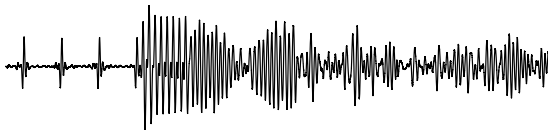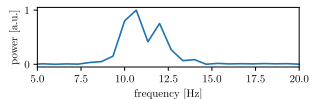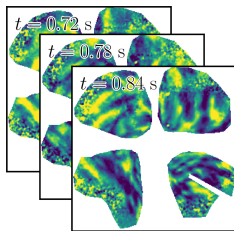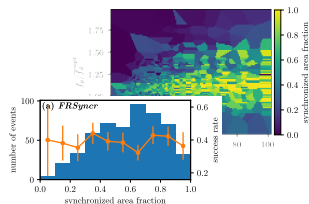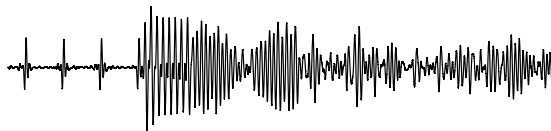


**Diverse and Complex Data**

**Diverse and Complex Data**

**Diverse and Complex Data**

**Diverse and Complex Data**

**Diverse and Complex Data**

**Diverse and Complex Data**

**Diverse and Complex Data**

# Challenges in Scientific Data Management

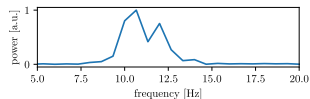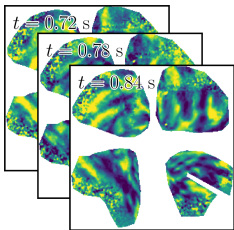- Standardization of file storage

# Challenges in Scientific Data Management

- Standardization of file storage

- Data formats

# Challenges in Scientific Data Management

- Standardization of file storage

- Data formats

- Metadata

# Challenges in Scientific Data Management

- Standardization of file storage

- Data formats

- Metadata

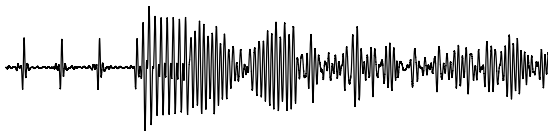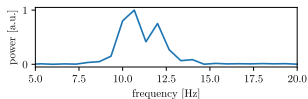- Linking of raw data, processed data, analysis results, documentation, software

# Challenges in Scientific Data Management

- Standardization of file storage

- Data formats

- Metadata

- Linking of raw data, processed data, analysis results, documentation, software

- Retrievability and searchability

# File formats, file names and folders

- Filename and folder structure conventions

  Example:
  /Photos/Holidays/Europe/2019/img1.jpg or
  /Photos/2019/Holidays/Europe/img1.jpg ?

- No central file storage

- Vendor-lock-in / Proprietary file formats / Missing APIs (application programming interfaces)

- Undocumented file formats

# Metadata

Many file formats don't provide metadata storage.

Popular workarounds:

- Store metadata in filename:
  datafile_networksimulation_20190207_a120_b17_CRPGM_debug.dat

- Non-standardized text file formats

→ How to search metadata?

  "Find all data files which were recorded by Person X."

## "Missing Link"

- I have the plot `plot_network.pdf`, where is the plotting script and the raw data?

- I have the data file `datafile_networksimulation_20190207_a120_b17_CRPGM_debug.dat`, but haven't I plotted the data before?

# Important Requirements

- Search functionality

## Important Requirements

- Search functionality

- Ability to store every data format, at any file size

# Important Requirements

- Search functionality

- Ability to store every data format, at any file size

- Possibility to store, link and retrieve raw data, processed data, analysis results and documentation

# Important Requirements

- Search functionality

- Ability to store every data format, at any file size

- Possibility to store, link and retrieve raw data, processed data, analysis results and documentation

- Support all kinds of data analysis software, from simple scripts to high-level software

## Important Requirements

- Search functionality

- Ability to store every data format, at any file size

- Possibility to store, link and retrieve raw data, processed data, analysis results and documentation

- Support all kinds of data analysis software, from simple scripts to high-level software

- Minimally invasive workflow

# Important Requirements

- Search functionality

- Ability to store every data format, at any file size

- Possibility to store, link and retrieve raw data, processed data, analysis results and documentation

- Support all kinds of data analysis software, from simple scripts to high-level software

- Minimally invasive workflow

- Scientific environments change often: Need for flexible data model

# CaosDB

# Research Data Management during Data Analysis

## Data Aquisition

Electronic Lab
Notebooks



z.B. RSpace, IDBS E-WorkBook,
Biovia

## Data Publication

Data Repositories



ePIC / PID

Metadata

Scientist picture: `201705 Scientist bench F.svg` from commons.wikimedia.org/wiki/Category:Life_science_images_from_DBCLS, CC-BY 4.0
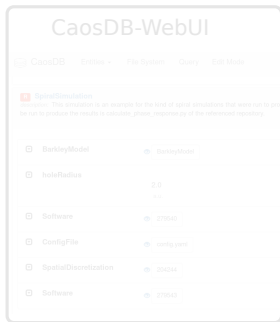
Bookshelf: https://openclipart.org/detail/289378/bookshelf-with-blue-books

# Research Data Management during Data Analysis



Data Aquisition

Electronic Lab Notebooks

z.B. RSpace, IDBS E-WorkBook, Biovia

Data Publication

Data Repositories

ePIC / PID

Metadata

Data Analysis

Caosdb
an open scientific database

Scientist picture: 201705 Scientist bench F.svg from commons.wikimedia.org/wiki/Category:Life_science_images_from__DBCLS, CC-BY 4.0

Bookshelf: https://openclipart.org/detail/289378/bookshelf-with-blue-books

# CaosDB Overview



Data Acquisition:
Use your desired workflow!

Data Files → FileSystem

CaosDB-WebUI

CaosDB-Python-Interface

```
lasse@salexan-x1 ~ % ipython
    results = db.execute_query(
```

CaosDB-Crawler

Automatic
file indexing

RESTful
XML-Protocol

Caosdb
an open scientific database

OpenSource! https://gitlab.gwdg.de/bmp-caosdb

# CaosDB Overview



Data Acquisition:
Use your desired workflow!

Data Files

FileSystem

CaosDB-WebUI

CaosDB-Python-Interface

CaosDB-Crawler

RESTful
XML-Protocol

Automatic
file indexing

Caosdb
an open scientific database

OpenSource! https://gitlab.gwdg.de/bmp-caosdb

# CaosDB Overview

# CaosDB Overview



OpenSource! https://gitlab.gwdg.de/bmp-caosdb

# State and Future of the Project

- CaosDB is a scientific project at the Research Group Biomedical Physics

- It is developed since ≈9 years and running stable since ≈2016

- CaosDB has been released as OpenSource software in 2018

- CaosDB is currently tested in several other workgroups in and outside of Göttingen

- Since May 2019: Commercial support and development by IndiScale GmbH

# Demo: Tracing a publication back to its data

$\rightarrow$ Query: FIND Publication with author with firstname = Henrik

$\rightarrow$ Query: FIND Publication with author with firstname = Henrik

# Demo: Tracing a publication back to its data

# Demo: Tracing a publication back to its data

# Demo: Tracing a publication back to its data

Caosdb — an open scientific database

- OpenSource Project: https://gitlab.com/caosdb

- Paper in Data https://doi.org/10.3390/data4020083

- Demo Instance (hosted by IndiScale) https://demo.indiscale.com (Beta test!)

# Thank You!