# Large Scale Data Facility:
# Design of meta data and community-specific services

**Rainer Stotzka, Jos van Wezel**

In close collaboration with:
Steinbuch Centre for Computing
Institute of Toxicology and Genetics
Institute for Applied Computer Science

Institute for Data Processing and Electronics
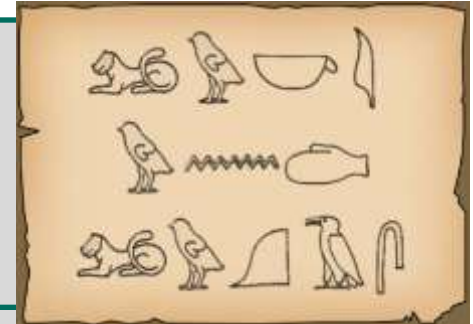
# Objectives of this talk

Title:            Large Scale Data Facility:
Design of meta data and community-specific services

- Why is the LSDF different?
- Why is meta data important?
- Data and meta data management
- Advantages for the user:
    - Long term sustainability
    - Additional services
    - High throughput data analysis
- Examples

# LSDF objectives (from the user's point of view)

**Storage**

- Dedicated for science data
- ExaByte scale data
- To archive data, long term sustainability
  (10 yrs. – ?)

**Interactivity**

- To enable scientists to gain better scientific results by providing
  - Data intensive analysis
  - Added value services for data intensive processing
- To provide high performance access, high throughput
- "Barrier free" access (easy-to-use)

Institute for Data Processing and Electronics

# Why is meta data necessary?

Meta data describe the contents of data

- Everybody uses meta data:
    - File name and extension
      (e.g. `rainer.jpg, budget.xls, Readme.doc`)
    - Location
      (e.g. `/…/EU-projects/2010/Fishy/budget.xls`)
    - Personal know-how
- → Sufficient for small file systems

Have you ever tried to locate a file or info-somewhere-in-a-file-system

- 15 years old ?
- in the file system of a colleague ?
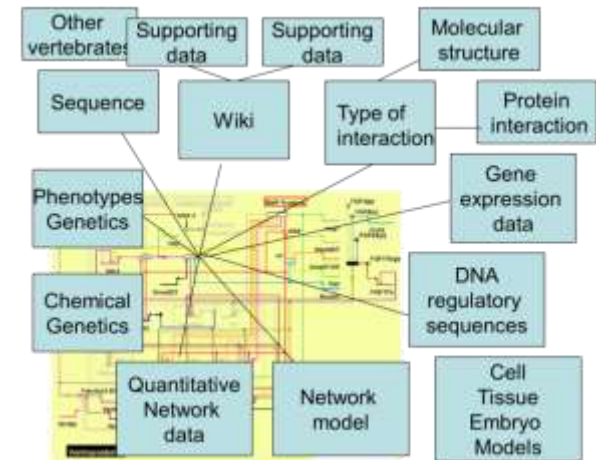- in a 100 PetaByte file system ?
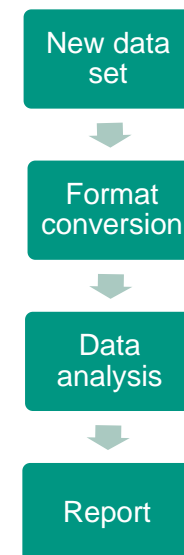
PANIC ?

# Applications requiring meta data

- Data archiving and retrieval (libraries)
- Fusion of complex data from various sources (data integration)

Community-specific services:

- Automatic processing
  (e.g. automatic analysis starts when data appears)
- Analysis chains
  (reporting analysis workflow, results and errors)
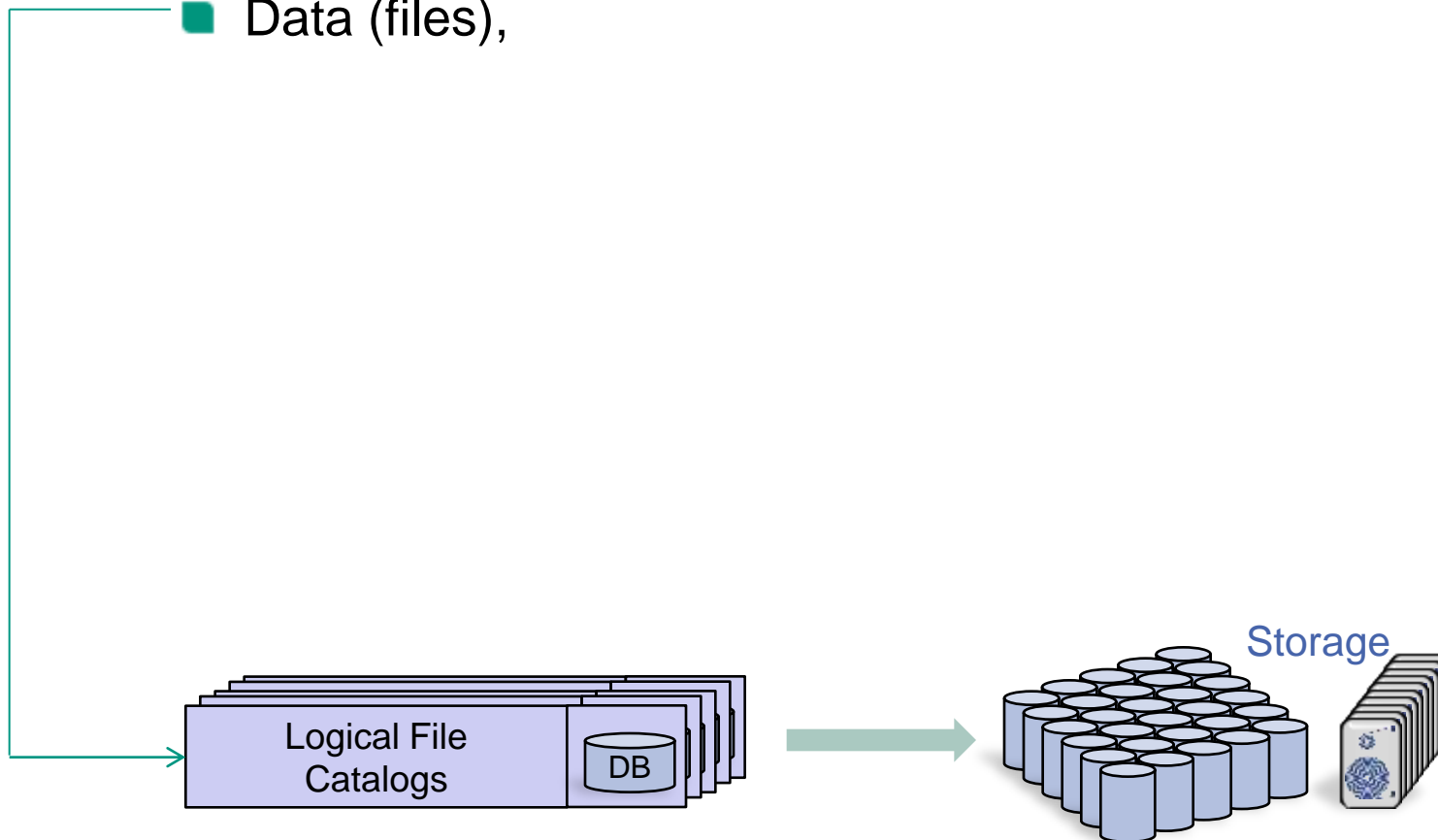- Google and Yacy
- Etc.



Source: Uwe Strähle

# Model of the LSDF meta data management
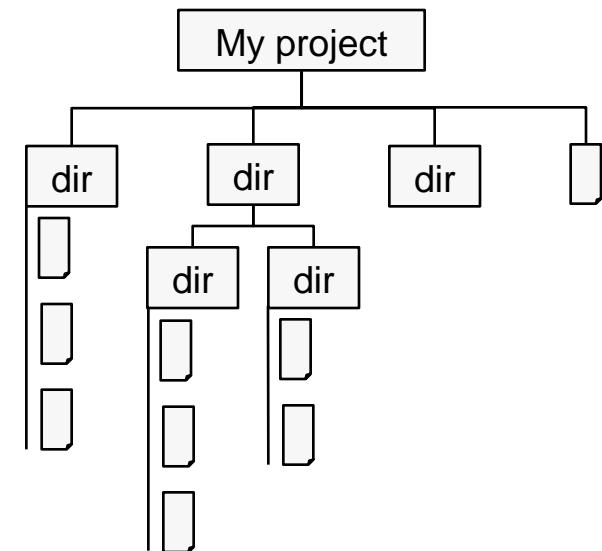
**Idea**:
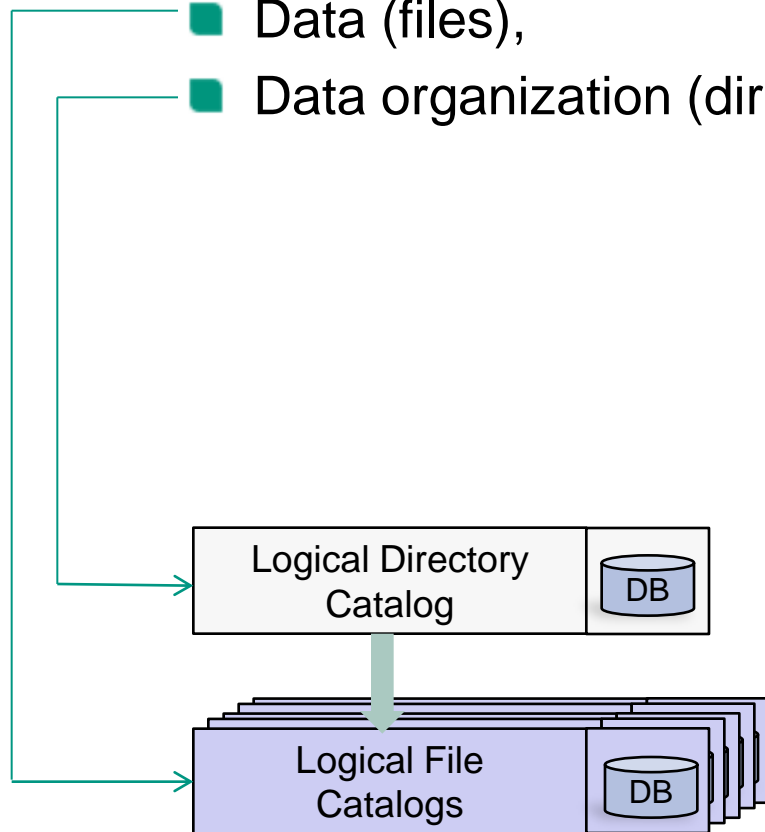
Clear separation between

- Data (files),

Storage

Logical File Catalogs | DB

# Model of the LSDF meta data management

**Idea**:

Clear separation between

- Data (files),
- Data organization (directory structure)

My project

dir    dir    dir

dir    dir

Logical Directory Catalog — DB

Logical File Catalogs — DB

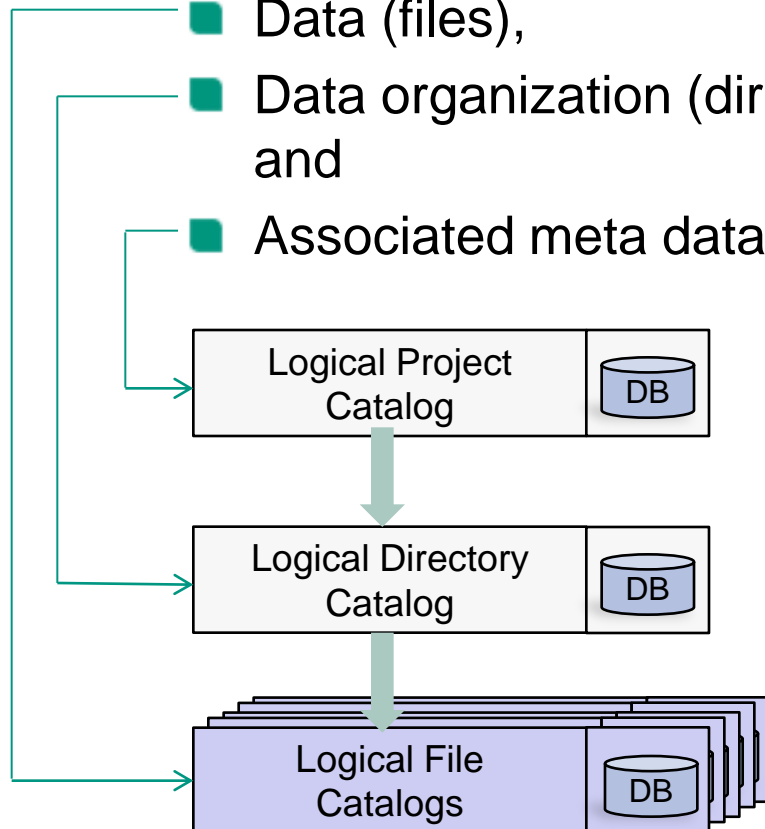Institute for Data Processing and Electronics

# Model of the LSDF meta data management

**Idea**:

Clear separation between

- Data (files),
- Data organization (directory structure) and
- Associated meta data

Logical Project Catalog — DB

Logical Directory Catalog — DB

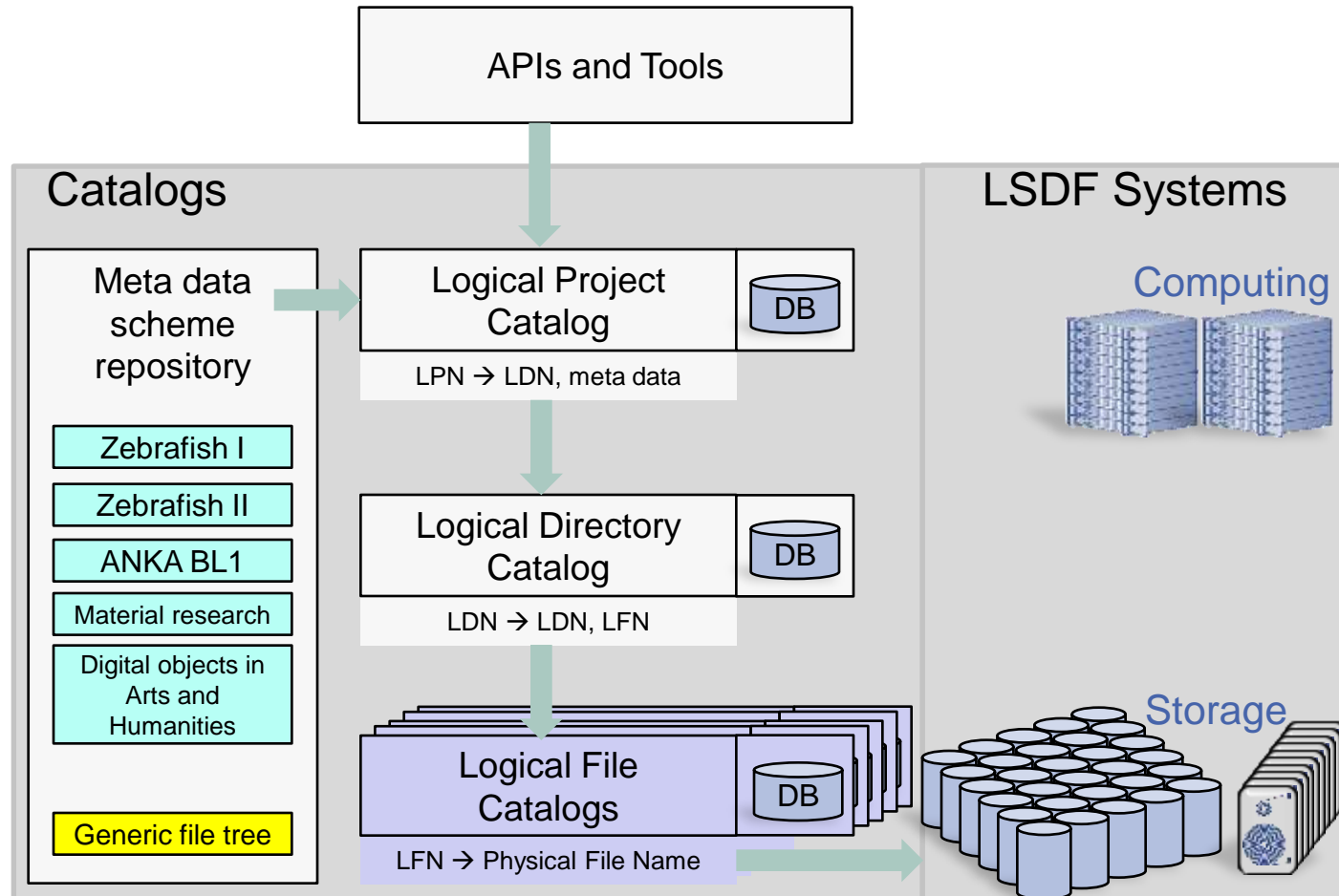Logical File Catalogs — DB

- name
- owners
- access rights
- date
  - community
  - (sub)subcommunity
    - measurement type
    - device, instrument
    - ...

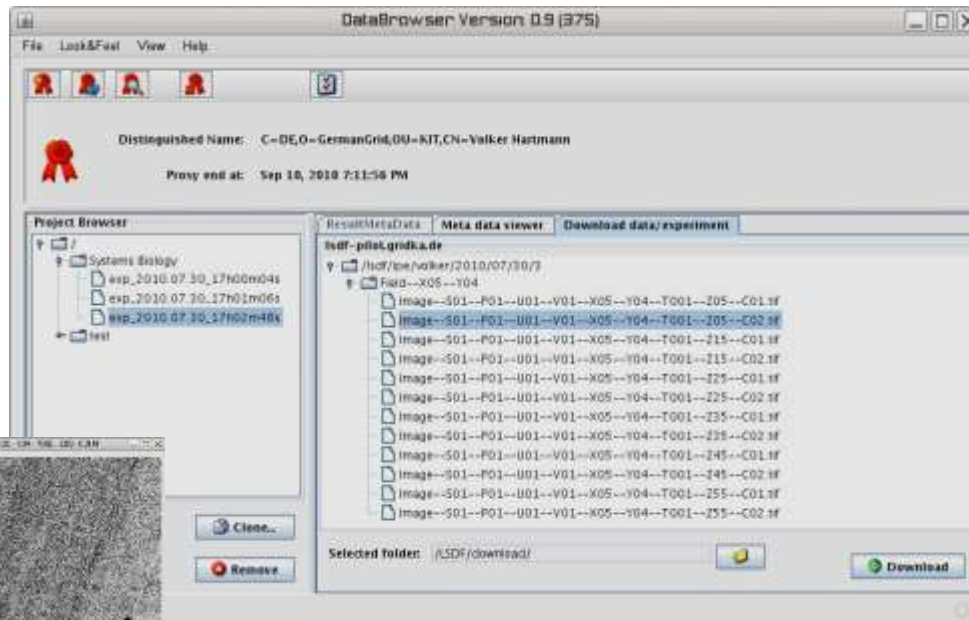Meta data structure depends on project, instruments, **time**, …

# Hierarchical Catalog System

- Sustainable
- Easily extensible
- Independent of data formats
- Enhanced performance: distribution of access
- Safety by redundancy

- Easy-to-use?

APIs and Tools

**Catalogs**

Meta data scheme repository

| Zebrafish I |
| Zebrafish II |
| ANKA BL1 |
| Material research |
| Digital objects in Arts and Humanities |

| Generic file tree |

Logical Project Catalog — DB

LPN → LDN, meta data

Logical Directory Catalog — DB

LDN → LDN, LFN

Logical File Catalogs — DB

LFN → Physical File Name

**LSDF Systems**

Computing

Storage

Institute for Data Processing and Electronics

# How to handle the complexity?

- *Apparently more complex: how do I use it?*
  - → Simple access tools, which can be easily adapted to your specific needs
    - → **LSDF DataBrowser** is a File-, Data- and Project-Explorer



DataBrowser allows:
- Authentication
- Project and file browsing
- Upload
- Download
- Edit meta data
- Data visualization
- Control data analysis

Features:
- Extensible
- Huge variety of communication protocols
- Open source

Institute for Data Processing and Electronics

# How to handle the complexity?

■ *How do I insert a new scientific project ?*

→ Data and meta data organization experts for projects with specific needs

→ Generic meta data format for simple file trees

■ *How do I transfer my data to a different location? Do I loose my meta data?*

→ Import-export to standard data and meta data formats

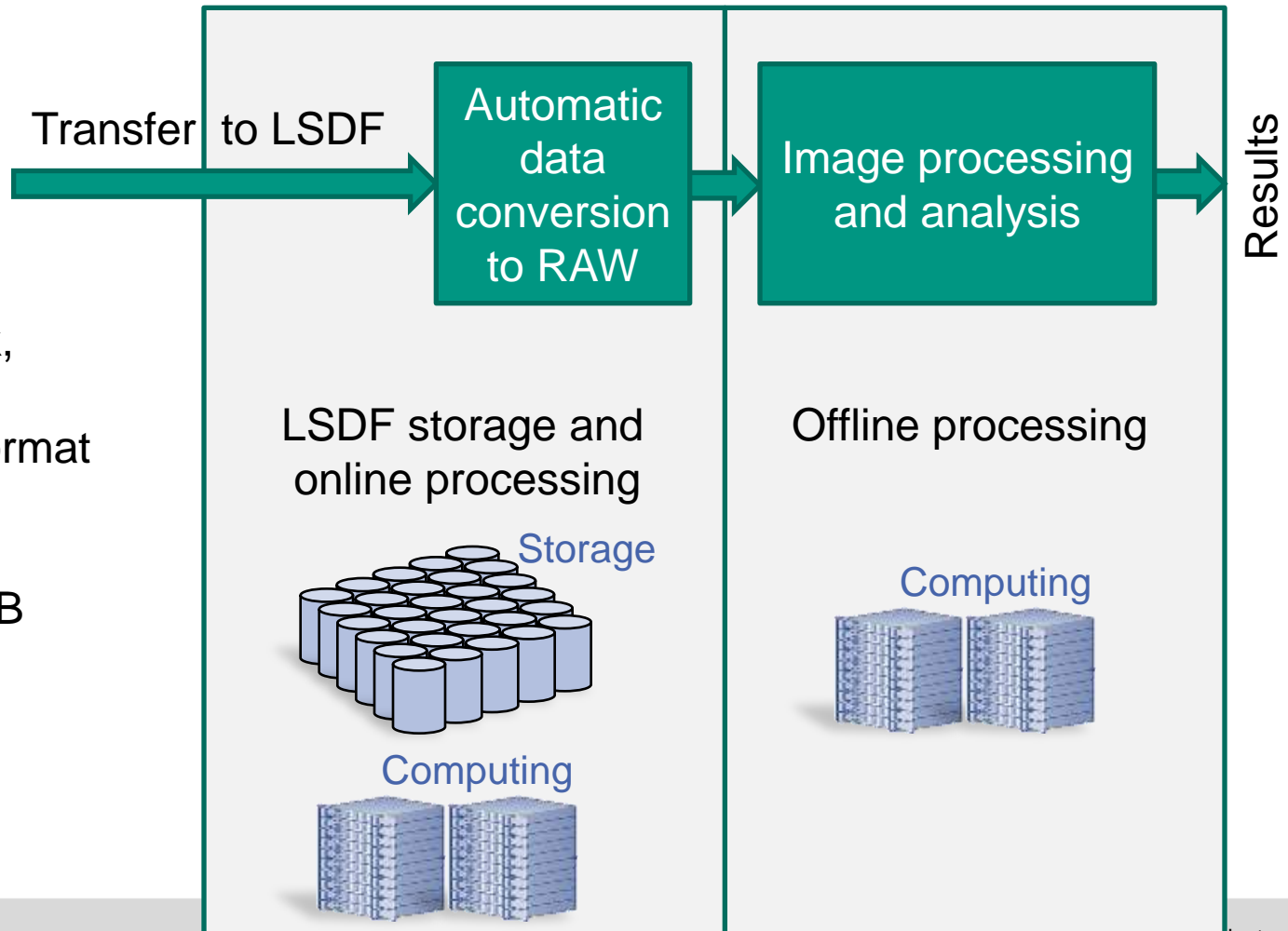→ Archive-in-a-box (Web installer or DVD, zip-archive, etc.)

Institute for Data Processing and Electronics

# Example: Toxicology in early life stages *in vivo*

■ Complex image analysis chain:



3D image stack,
time series,
Leica Image Format

data set size:
      100 GB
# data sets:
      > 100

Transfer to LSDF

**Automatic data conversion to RAW**

**Image processing and analysis**

Results

LSDF storage and online processing

Offline processing

Storage

Computing

Computing

# **Example: Toxicology in early life stages *in vivo***

- Close cooperation  ITG, IAI, SCC and IPE
  (Thanks to Jens C. Otte for the images)

Data Browser:

- Meta data organization

- Adapted Data Browser implementation

- Implementation of data conversion

- Automatic data conversion workflow at LSDF
  steered by meta data

Estimated effort:          ~ 2 PM

# Example: ITG adapted DataBrowser

Institute for Data Processing and Electronics

# Scientific communities

- Systems biology (ITG, BioQuant, Immunogenetics)
    - Vertebrate development studies and
    - Deconvolution (5000 data sets $\rightarrow$ <180 min.)
- Synchroton facilities and beamlines
    - ANKA data storage
    - HGF "High Data Rate Initiative"
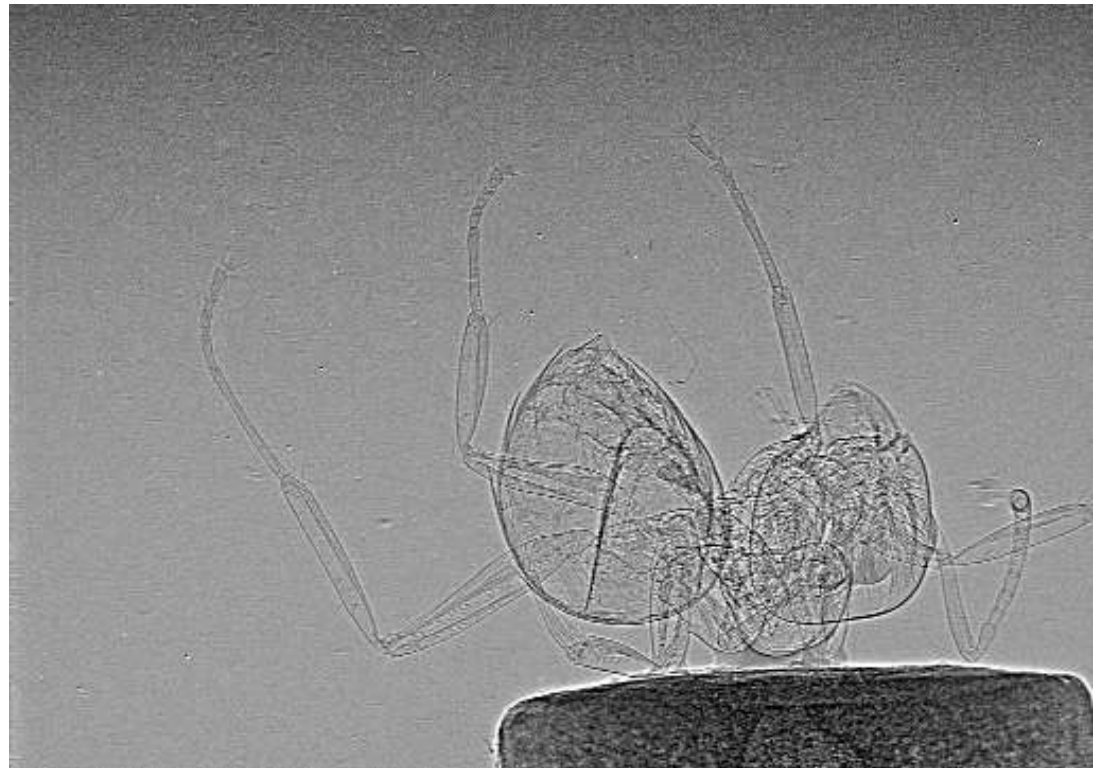- Climate research
- Material research
- Arts and humanities



*»Il Cenacolo« von Da Vinci (1494-98)*



*»L'ultima cena« von Julius Romanus (1754)*

Institute for Data Processing and Electronics

# Data intensive science

- Algorithms for data analysis
- Visualization of huge 3D data sets: online visualization of 500 GB data sets

Institute for Data Processing and Electronics

# Conclusions

- LSDF is a powerful structure
  → more than data storage and cluster computing
- Design for future requirements → ExaByte storage + interactivity
  - R&D in progress

LSDF offers

- Sustainability and safety
- Flexibility for future requirements
- Interactivity
- Community-specific services
- Support

→ To gain faster and better scientific results

# Thanks to

The team at IPE:
Volker Hartmann
Thomas Jejkal
Michael Sutter
Francesca Rindone
Michael Götter
Patrick Neuberger
Simon Ochsenreither

The team behind LSDF at SCC:
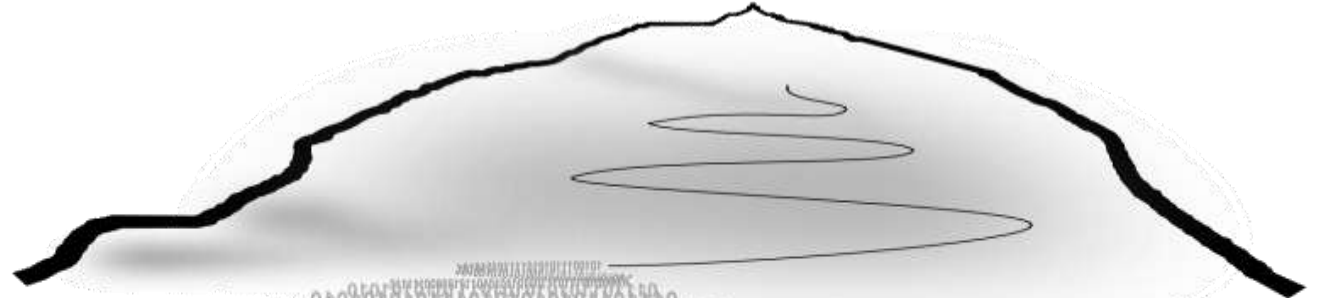Serguei Bourov
Ariel Garcia
Bruno Hoeft
Rainer Kupsch
Achim Streit
Bernhard Verstege

LSDF

01110110101010101011101010101010