

# **SCTB review: Computing and data analysis**

Nicolas Soler - 09/11/2021  
*Computing Division*

# The Computing division

## Computing & Controls



Oscar Matilla  
Computing Division Head



Concepción Girbau  
Secretary

## Controls and DAQ



Guifré Cuni  
Controls Section Head

## Electronics



José Ávila  
Electronics Section Head

## IT Systems



Toni Pérez  
IT Systems Section Head

## MIS



MIS Section Head

## SDM

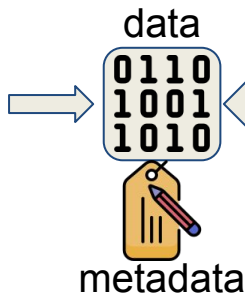
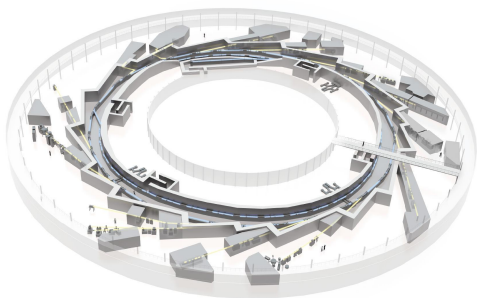


Nicolas Soler  
SDM Section Head

New  
Feb 2021

# 2021: Scientific Data Management

*A new section inside the Computing & Controls division*

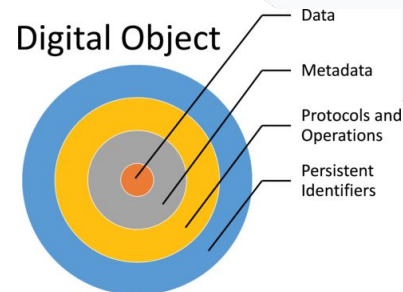


Data processing

Programming  
Optimization Java HPC  
Data C/C++ Analysis  
Pipelines Processing Visualization  
**Scientific Computing**  
Machine Learning Web scrapping  
Matlab Python Libraries  
API data reduction

Data reusability

- Metadata ingestion
- Provenance
- Persistent identifiers
- Catalogue
- DAaaS



**4 people** hired with hybrid science / computing profiles (+1 more in the beginning of 2022)

- Beamlines support
- Occasional accelerators support



<https://www.nist.gov/programs-projects/facilitating-adoption-fair-digital-object-framework-material-science>

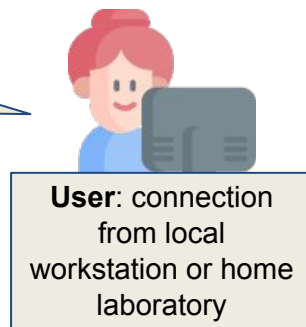
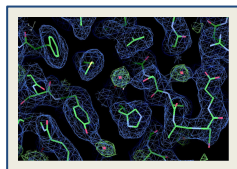
## Data availability

The data supporting this study can be made available from the corresponding author upon request.

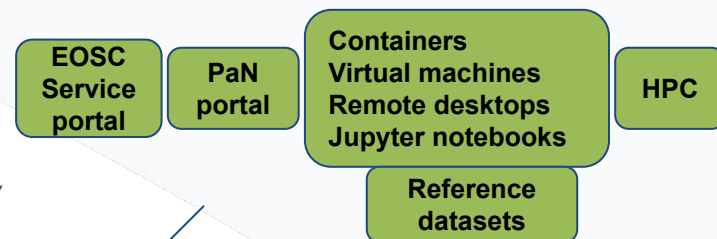
# Towards FAIR synchrotron data

panosc

ExPaNDS



## Data analysis as a service (WP4)



Use re-use

Process raw data

Annotate & store processed data + provenance

## FAIR-ready data (WP2)

## Data catalogues (WP3)

## Data reduction / compression

- Metrics
- New algorithms
- Software vs hardware
- Technique-specific
- Lossy vs non lossy
- Meta-compressors (ML)

LEAPS  
INNOVATION

(WP7)

**Certified**  
(meta)data repository

Standard metadata framework

FAIR Data Management Plan (**DMP**) for each technique

Persistent Identifier (**PID**) for data, instrument, software, sample)

Common **data policy** framework

API / UI

**EOSC hub**

**Metadata ontologies**  
(shared vocabulary)

**PaNOSC federated catalogue**

**Metadata formats**  
(NeXus)

**Catalogues**  
(eg iCAT, SciCat)

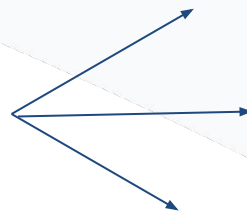
# Implementation of a data catalogue at ALBA



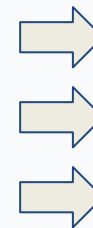
EOSC  
portal



PaNOSC  
Federated  
Metadata  
catalogue



iCAT  
iCAT  
SciCAT  
.....



Experimental  
data, metadata,  
derived data and  
related info (i.e  
PIDs)

European  
facilities  
data catalogues



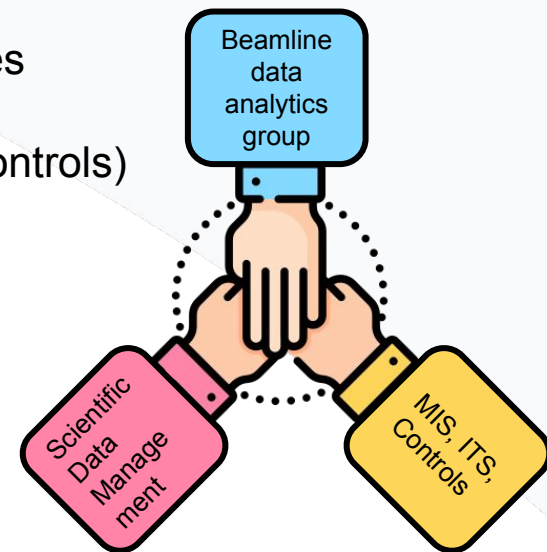
Search						
<input type="checkbox"/>	Date	Sample	Dataset	Definition	Files	Size
<input type="checkbox"/>	21:11 2 Jun 2021	PPDL30CIN2	bl11_u2020084418_9909		1232	2.8 GB
<a href="#">Download</a>						
Summary Instrument Files 1232 Metadata List						
Search						
Name	Value					
__elapsedTime	49824					
__fileCount	1232					
__volume	3055682784					
beamlineID	BL11					
datasetName	bl11_u2020084418_9909					
endDate	Thu Jun 3 11:01:25 2021					

- Data have an embargo period of 3 years
- iCAT already in use for NCD-SWEET
- Next: Implementation in MIRAS and MSPD
- Open externally by the end of the year

# 2021: Scientific Data Management

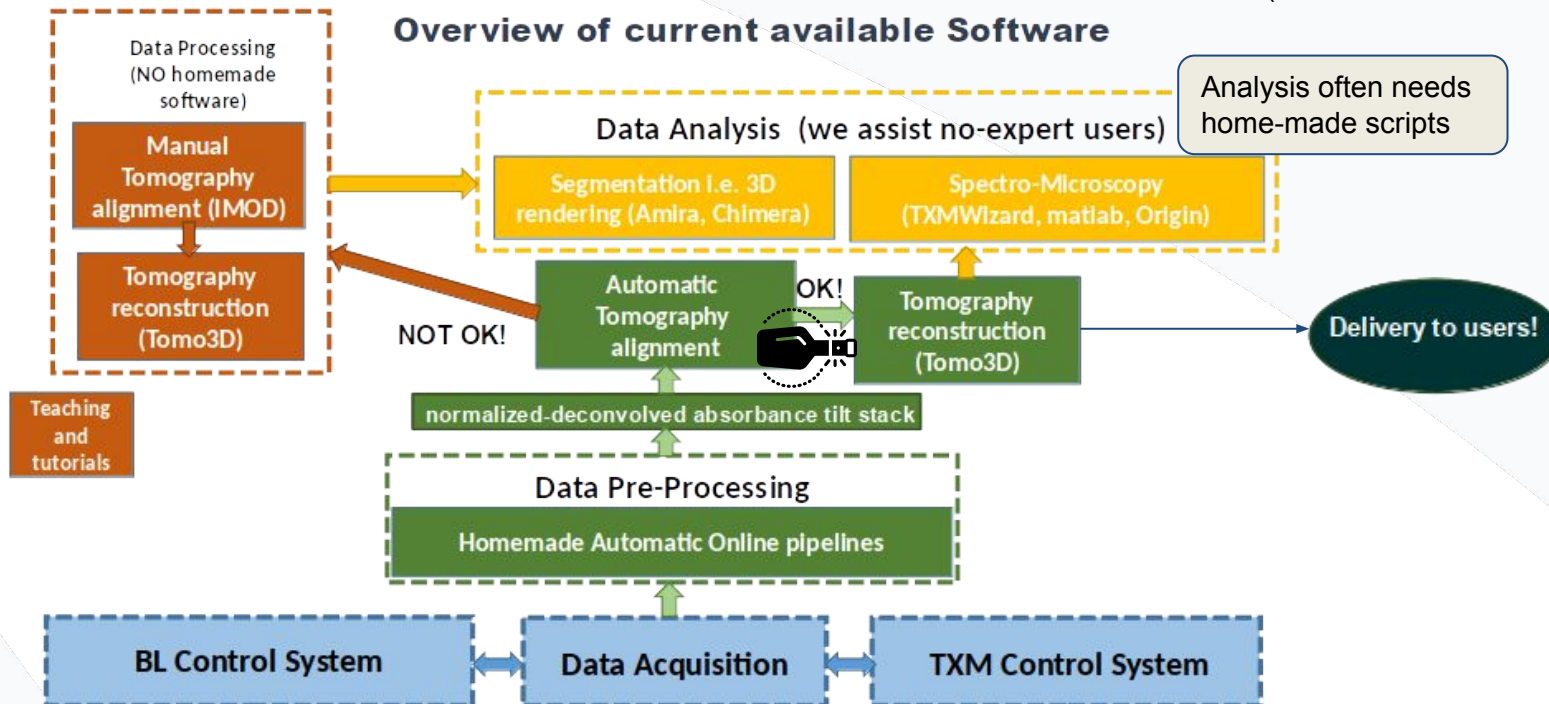
## *BUILDING A SYNERGIC SOFTWARE DEVELOPMENT PROCESS BETWEEN COMPUTING ENGINEERS AND METHOD DEVELOPMENT SCIENTISTS AT THE BEAMLINE*

- Building downstream data processing and analysis pipelines
- Real-time data processing triggering and feedback (with Controls)
- Data visualization tools and GUI
- Databases and APIs
- Data reusability ( ExPaNDS, with MIS and Controls)
- **Collaborating in data analysis methods development**



# MISTRAL data flow (cryo-SXT)

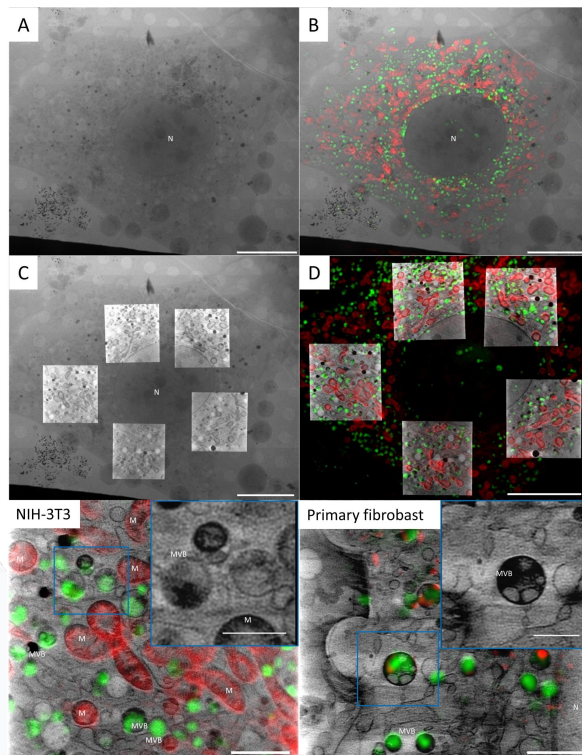
(from A.Sorrentino' slide)





# Example of data analysis

## Effect of an antifibrotic novel therapeutic agent at the cellular level

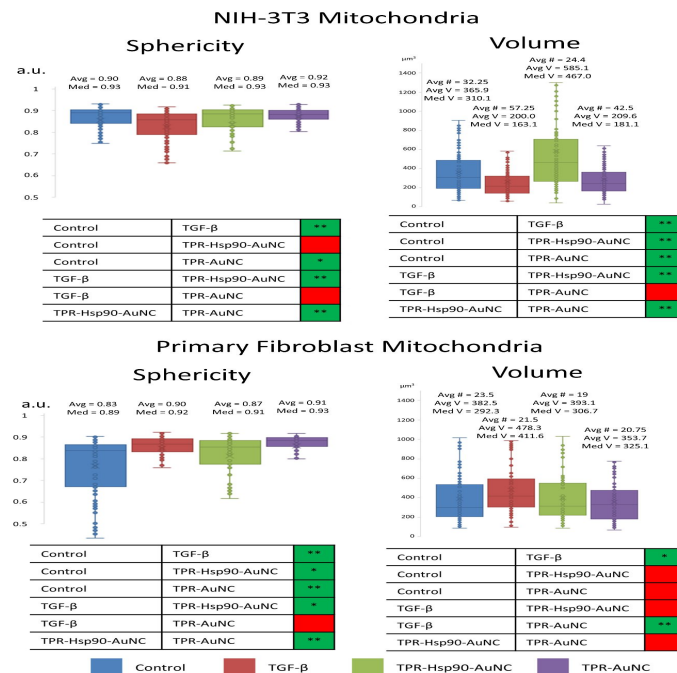


## EDGE ARTICLE

## Correlative 3D cryo X-ray imaging reveals intracellular location and effect of designed antifibrotic protein–nanomaterial hybrids†

Cite this: DOI: 10.1039/d1sc04183e

† All publication charges for this article have been paid for by the Royal Society of Chemistry

J. Groen,<sup>a</sup> A. Palanca,<sup>b,c</sup> A. Aires,<sup>d</sup> J. J. Conesa,<sup>a,c</sup> D. Maestro,<sup>b</sup> S. Rehbein,<sup>f</sup> M. Harkiolaki,<sup>g</sup> A. Villar,<sup>h,i</sup> A. L. Cortajarena,<sup>g,di</sup> and E. Pereiro<sup>g,sa</sup>

$t \leq 0.05$ : (\*) significant difference  
 $t \leq 0.01$ : (\*\*) extremely significant difference



# MISTRAL DP/DA: GAP analysis

DP feature	Current setup	Problem / disadvantage	Mitigation / future setup
<b>User interface</b>	user interface	Scripts-based	GUI / workflow manager
<b>Preprocessing</b>	Custom pipelines (single/multiple images per angle for tomos, spectroscopy 2D & 3D)	has to be launched manually after data acquisition	Automatic triggering to be installed
<b>Alignment (tomo &amp; spectra)</b>	<a href="#">IMOD</a> , EFTEM- <a href="#">TOMOJ</a> & 2D correlation Image J or python (2D spectroscopy)	presence of fiducials required but hardly controllable	Optical flow, <a href="#">TomoPy</a>
<b>Reconstruction</b>	SIRT ( <a href="#">Tomo3D</a> ), ART of <a href="#">TOMOJ</a>	ok but could be GUI-monitored	followed on workflow manager
<b>Segmentation &amp; analysis</b>	<a href="#">AMIRA</a> , <a href="#">Chimera</a> , <a href="#">SuRVoS</a>	Not enough automation?	<a href="#">ML approaches</a>
<b>Correlative microscopy 3d</b>	in progress	to test	<a href="#">Icy- eC-CLEM</a> ?
<b>Data access, reprocessing</b>	via sFTP, HDD	preservation, provenance	Data catalogue, DAaaS

# MISTRAL data processing & analysis



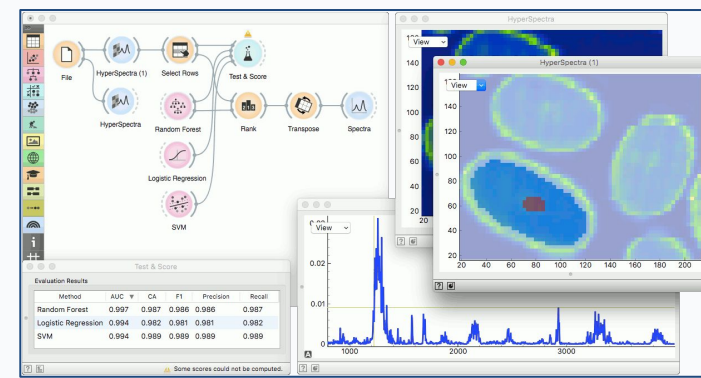
- Bring more **automation** in the current data processing pipeline. Bring graphical interfaces and data workflows where needed. Improve current pipelines.
- Install a robust, parallelized xy **fiducialless alignment** correction program (e.g. Tomopy) that could be used by default (whether or not fiducials are present).
- Provide software support for **correlative microscopy** (cryoET and cryo3DSIM volume superposition) and **segmentation methods**
- Together with on-site experts, assist in the development of **analysis** software.

- Workflow manager and previews (**GUI**)
- **preprocessing**: from raw data to phase retrieval, (normalisation etc).
- computer reconstruction (aiming at **real-time**, parallelized on the cluster)
- image **registration** (re-aligning / shifting)
  - on angular projection
  - on reconstructed volume
- **Machine learning** for low dose, segmentation, denoising, missing angles etc
  - on angular projection
  - on reconstructed volume
- Managing and visualizing **big data volumes** (size: few GB to several TB)
  - distributed memory visualization (NVIDIA index, paraview plugin)
  - Data selection (automated?), Database of scans (iCAT)
  - Data compression (probably lossy)

# Cell and tissue biology beamlines:

*How can SDM help?*

## MIRAS IR spectroscopy & microscopy



- Move from proprietary to **open source solutions** (e.g Orange, Quasar). Join the IR quasar network (SOLEIL, Elettra, Sirius, Canadian light source etc), especially bringing **data workflows**, **machine learning** capabilities and nice feature to explore **imaging**/hyperspectral data.
- Improve the data analysis as a service by switching to these solution.
- Provide advice for **statistical methods** (eg t-SNE instead of PCA)
- **Multimodal** experiments: correlation with other types of spectroscopies



orange

# Summary



- SDM has been created 9 months ago to support Experiments with **data processing, analysis and visualization needs** as well as with future data reusability (**FAIRness**)
- The first 2 engineers started in September, 1 in November, 1 more to come until in Jan. 2022, 1 to be hired next year.
- SDM will coordinates the **metadata ingestion scheme** and **catalogue**, as well as the transition towards **data analysis as a service** for remote users.
- Computing engineers with scientific background will work hand in hand with data analytics **beamline method developers** to:
  - produce and maintain **data processing & analysis pipelines**, bringing expertise in **software development** and other relevant technologies (i.e machine learning).
  - Provide support for **integration of heterogeneous data**.