

EXPERIMENTAL DATA STORAGE MANAGEMENT IN NEXUS FORMAT AT SYNCHROTRON SOLEIL

S. Poirier, C. Marechal, M. Ounsy, A. Buteau, P. Martinez, B. Gagey, P. Pierrot, M. Mederbel, JM Rochat, SOLEIL, Gif-sur-Yvette, France

Abstract

At Synchrotron SOLEIL, some 20 beamlines are already in operation and produce daily experimental data files with volumes ranging from a few MB up to 100 GB.

Almost all experimental data is stored using the NeXus data format [1]. This logical format, based on an HDF physical format, is suitable for storing any kind of scientific data produced in Neutron or Synchrotron sources. The NeXus format allows experimental data to be recorded together with all needed contextual information about the experiment, instrumentation, sample, user, and so on.

Several tools have been developed or are currently under development at SOLEIL, including a web data browser to retrieve, browse and download data, and a command line tool used to export data from NeXus to ASCII or binary data files.

The NeXus storage system is fully integrated into the Tango bus [2] through a set of dedicated devices. It is therefore possible to record any data coming from Tango devices into a NeXus file. The storage system is designed to allow recording data from various data sources using a plugin system.

In terms of hardware, a high availability system using an innovative cellular storage concept known as Active Circle [3] has been in use since December 2006.

OUR CONTEXT

Missions and Organization of the Computing Division

The SOLEIL computing division is in charge of managing the complete lifecycle of data file production for the machine and the beamlines. The division is able to fulfill this permanent and strategic mission thanks to the collaborative work between its four groups:

- Data is first acquired on dedicated electronics (*CompactPCI boards, detectors, CCDs, etc.*) which are managed by the Electronics group.
- The Software for Controls and Data Acquisition provides Tango devices to get the data from electronics, along with other software modules (*ScanServer device, Passerelle[4]*) to automate the data collection process.
- Once data is available on the Tango software bus, the Data Management group provides software applications to store these files in a coherent way and tools to easily retrieve them.
- Lastly, the System and Network Infrastructure group is in charge of providing storage infrastructure for the substantial volume of data files generated.

Management of Experimental Data Files is Strategic for SOLEIL

The production of experimental data files reflects the output of an instrument such as SOLEIL. Being able to store and retrieve this heritage during the future years of SOLEIL operation is one of the important missions of the computing division.

Facing the Diversity of Scientific Applications

SOLEIL beamlines are dedicated to numerous scientific domains using different energy ranges and many different experimental techniques. This diversity leads to a significant number of detection systems and acquisition processes. It is therefore quite challenging to provide software tools able to cover these different scientific applications in a generic way.

From the data storage point of view, this diversity means that data volumes and data file retention periods may be very different from one beamline to another.

OUR MOTIVATIONS

Decouple the Data Acquisition, Data Storage and Data Collection Processes

Our vision was that software for collecting data files on a beamline may be seen as (*at least*) a three-way process:

- The data acquisition specialist should take care of all issues regarding the interfacing of the detection system.
- The beamline scientist should focus on the definition of the experimental process needed to get a good measurement.
- The data storage specialist should take care of the details regarding the proper organization of all data describing the measurement.

This vision leads us to work on providing a high level data storage service which hides ugly details of how data files are organized from the data acquisition specialist and beamline scientist.

Organize Data Files with a Common Set of Metadata on all Beamlines

To be able to provide common tools for data retrieval and visualization, a common set of metadata associated with scientific data has been defined, and their correct definition by scientists enforced by various means (*dedicated GUI components which have been integrated in various end user applications, connection with our DUO system, etc.*).

Keep Close to the Standards of our Scientific Community

It was important for SOLEIL that the choice of data format should be coherent with our scientific community, so that in the future we can directly use data reduction and analysis applications developed in other synchrotron and neutron sources. NeXus therefore appeared as a natural choice in this respect.

Easily Manage the Lifecycle of Data Files

The retention time for data files varies depending on the beamline and whether an experiment is conducted by a SOLEIL or external user. The storage infrastructure should provide system administrators with a transparent way to manage these different life cycles.

Data Storage Infrastructure should be Independent from Hardware Solutions

Storage vendors often offer solutions based on proprietary hardware. Instead, we wanted to keep the future extensions of our system open to new technological progress in the mass data storage industry.

PRODUCTION OF NEXUS DATA FILES

A High Level Software Service on the Tango Bus in Charge of Managing Data Storage

SOLEIL data acquisition systems are based on the Tango software bus which carries all important information. To manage Nexus files, dedicated Tango devices (*DataRecorder*, *AuthServer*, *NeXusReader*) have been developed to perform the following functions:

- Create the data files with the access rights associated with the user conducting the experiment.
- Gather metadata and experimental data associated with the experiment.
- Organize the internal structure of the Nexus file.

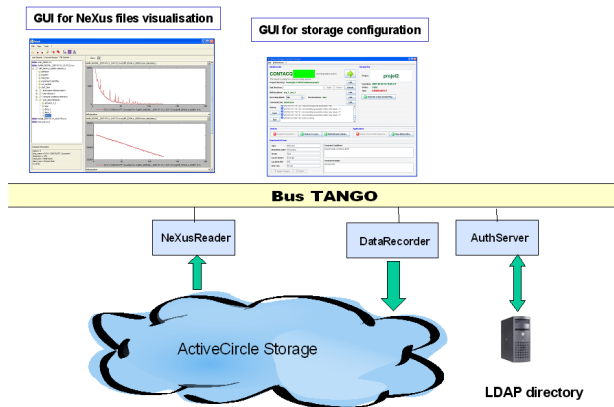


Figure 1: Software components in charge of NeXus files storage.

Use NeXus as a Container in a First Approach

To be ready before data analysis applications are adapted to the NeXus format, we decided that NeXus file would initially be used only as a data container. Thus, from a NeXus file, an extraction tool allows users to generate files in the format expected by their analysis applications. Moreover, an IGOR plugin also allows NeXus files to be directly opened with this scientific tool.

File Searching and Visualization Tools

The metadata stored in the NeXus files can be used to index them in a relational database. A dedicated graphical user interface, accessible through the web, acts as an advanced file searching engine for SOLEIL users. Coupled with a powerful data visualization application, NeXus files can then be remotely found, browsed and downloaded by SOLEIL users directly from their home institutes.

Project Code	Beamline	Type	Comment condit.	Comment sample	Beamline leader	User List	Log book ref.	Date
pt1	CONTACO	type	conditions exper.	sample conditions ligne	ba	logbookref1	10162008	14 41
pt1	CONTACO	type	conditions exper.	sample conditions ligne	ba	logbookref1	10172008	11 49
pt1	CONTACO	type	conditions exper.	sample conditions ligne	ba	logbookref1	10172008	12 39
pt1	CONTACO	type	conditions exper.	sample conditions ligne	ba	logbookref1	10242008	16 20
pt1	CONTACO	type	conditions exper.	sample conditions Feed	ba	logbookref1	10242008	16 22
pt1	CONTACO	type	conditions exper.	sample conditions Feed	ba	logbookref1	10242008	16 24
pt1	CONTACO	type	conditions exper.	sample conditions Feed	ba	logbookref1	10242008	16 30
pt1	CONTACO	type 1	conditions exper.	information sur l'	Thomson	bois ba	01172007	16 28
pt1	CONTACO	type 1	conditions exper.	information sur l'	Thomson	bois ba	01182007	11 28
pt1	CONTACO	type 1	conditions exper.	conditions Archiv.	Thomson	bois ba	01182007	11 58

Location	NeXus Entry	Date
database_2008-10-24_16-32-39.nxs	exp_1_run_0+NeXusEntry>	2008-10-24
database_2008-10-24_16-32-39.nxs	exp_1_run_0+NeXusEntry>	2008-10-24
database_2008-10-25_16-00-01.nxs	exp_1_run_0+NeXusEntry>	2008-10-25
database_2008-10-25_16-00-01.nxs	exp_1_run_0+NeXusEntry>	2008-10-25
database_2008-10-25_16-08-52.nxs	exp_1_run_0+NeXusEntry>	2008-10-25
database_2008-10-25_16-08-52.nxs	exp_1_run_0+NeXusEntry>	2008-10-25
database_2008-10-25_16-16-56.nxs	exp_1_run_0+NeXusEntry>	2008-10-25
database_2008-10-25_16-21-36.nxs	exp_1_run_0+NeXusEntry>	2008-10-25
database_2008-10-25_16-21-36.nxs	exp_1_run_0+NeXusEntry>	2008-10-25
database_2008-10-25_16-28-32.nxs	exp_1_run_0+NeXusEntry>	2008-10-25
database_2008-10-25_16-28-32.nxs	exp_1_run_0+NeXusEntry>	2008-10-25
database_2008-10-25_16-28-32.nxs	exp_1_run_0+NeXusEntry>	2008-10-25

Project Code	Beamline	Type	Comment condit.	Comment sample	Beamline leader	Feed	Local contact	Group	User List	Log book ref.	Date	
pt1	CONTACO	type	conditions exper.	conditions experimentale	ba		Dedei	group1	ba	logbookref1	110642008	16 32

Figure 2: TWIST [4].

THE STORAGE INFRASTRUCTURE FOR THESE DATA FILES

ActiveCircle Solution

The ActiveCircle product has been chosen for the SOLEIL's data storage infrastructure because it fulfilled all our requirements regarding management of data file lifecycles.

Moreover, it provides built-in mechanisms for data protection, high availability, and operations continuity planning. Data protection is assured by file versioning, which allows users to retrieve previous versions in the event of an error. High availability is delivered by the Virtual File Server cluster: if one node fails, another automatically takes its place. Finally, operations continuity is ensured through the replication of data over multiple nodes:

- On the beamline node (for continuity of operation in case of a network backbone cut)
- In each of our computing rooms located in two different buildings

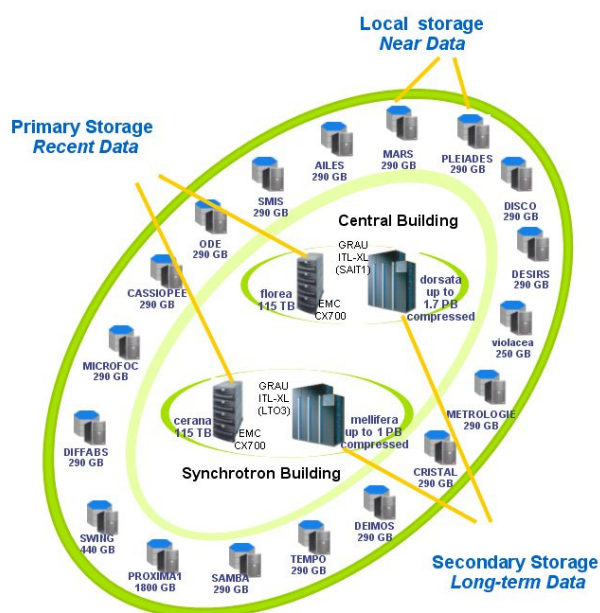


Figure 3: ActiveCircle deployment architecture.

SUCSESSES AND DIFFICULTIES

NeXus is Today Used on the Majority of the Beamlines

Table 1: Nexus Files Production Statistics

Beamline	NeXus files	Volume(GB)	Average Size (MB)
Cristal	7 733	7	0,9
Diffabs	10 340	61	5,9
Ode	45 542	11	0,3
Swing	15 116	169	11,2
Samba	91 679	205	2,2
Deimos	72	3	40,4
Proxima1	1 331	1	1,1
Cassiopee	1 402	0	0,1
Tempo	704	0	0,1
Antares	375	1	1,9
Lucia	2 784	7	2,2
Total	180 002	465	6,0

The Users did not Initially Understand our Motivations

During the beamline construction phase, the beamline scientists were focused on getting their first data files, regardless of format. Organizing data file production was not a top priority for them, and at this stage we lacked real usage scenarios to help validate and improve the technical solutions we proposed.

The Use of NeXus on Proprietary Systems has its Limitations

In some cases, data files are directly produced by turnkey proprietary acquisition systems, and access to metadata describing the acquisition is therefore often difficult. Addressing these specific problems has not been a priority for us at this time, even if some (imperfect) solutions are foreseen.

NeXus is not yet a Fully Defined Standard

The organization of experimental data and metadata within Nexus files must be defined by a long term (and not yet complete!!) standardization process. Initial exchanges of data files with ESRF have convinced us that higher level APIs should be used to decouple data analysis code development from internal file organization.

ActiveCircle: An Innovative Choice

SOLEIL has been one of ActiveCircle’s first customers. This company’s distributed storage technology provides a virtual data storage space, easily extensible to our growing needs, and which has proved to suit our requirements.

Nevertheless, as with any new product, we had to cope with a number of bugs and technical problems during the first year of operation in 2007. Thanks to the quality of the support, and after a few software upgrades, our ActiveCircle infrastructure is today very reliable.

CONCLUSION

NeXus was the Right Strategic Choice

When we made the choice of NeXus, this format was only used in a few neutron source experiments. Today, at SOLEIL, NeXus has become a *de facto* standard which federates software development for data reduction and data analysis. As an example, a fully NeXus-based SAXS data reduction application has recently been installed on the SWING beamline.

Future Steps

Of course, the most challenging goal is still to be able to transparently exchange data files between institutes. This would allow us to benefit from the data analysis tools developed in our scientific community. In this respect, we are trying to work with ESRF and DIAMOND which have also shown an interest on NeXus.

REFERENCES

- [1] <http://www.nexusformat.org>
- [2] <http://www.tango-controls.org>
- [3] <http://www.active-circle.com>
- [4] <https://twist.synchrotron-soleil.fr>