

Distributed
Archiving & Preservation System
Système de Préservation
et d'Archivage Réparti (SPAR)

Lessons from the implementation

Sun-PASIG Malta

Agenda

- Missions and context
- From theory to practice
 - the model
 - the expectations
 - handling “bad” data
 - lessons
- Management through risk analysis
 - Focus shift
 - Dealing with human skills

■ Missions :

- to build up the collections,
- to preserve and communicate them to the public,
- to produce a reference catalog,
- to cooperate with other institutions,
- to participate to research programs.

Legal deposit :

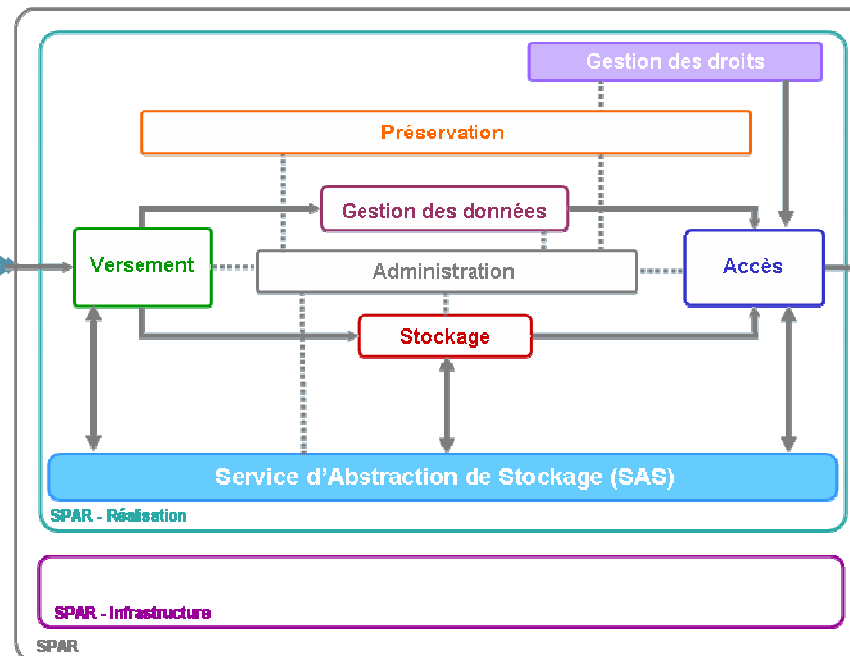
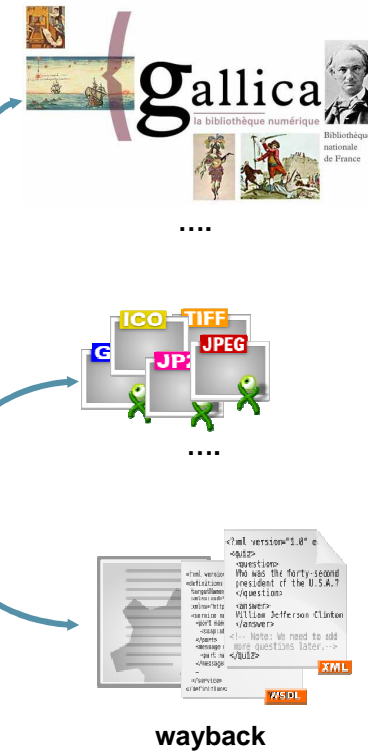
- legal deposit since 1537 for printed materials
- 1648: engravings and maps
- 1793: musical scores
- 1925: photos
- 1938: phonograms
- 1975: videograms
- 1992: electronic documents
- 2006: Web legal deposit

Digital archiving at BnF

Production applications



Dissemination applications



Decomposition in channels

- To deal with the variability and heterogeneity of the data, definition of channels
- build on the relation between the digital objects and the archival system, independently of any given organization:
 - Preservation digitization
 - Audiovisual material
 - Negotiated legal deposit (dark Web, regional press)
 - Automatic legal deposit (surface Web)
 - Administrative production
 - Deposit / Third party archiving
 - Acquisition / Donation

Agenda

- Missions and context
- From theory to practice
 - the model
 - the expectations
 - handling “bad” data
 - lessons
- Management through risk analysis
 - Focus shift
 - Dealing with human skills

- Each channel is formally defined by a reference package:
 - description of the Service Level Agreement (SLA)
 - includes human readable definition
 - machine actionable parameters
 - links to accepted formats at various levels of commitment (stored, identified, known, managed)
 - links to the used tools

Channel description (sample)

```
<sla:serviceLevelAgreement>
  <sla:header>
    <sla:channelIdentifier>FIL_NUM_CONS_A</sla:channelIdentifier>
    <sla:type>info:bnf/spar/context/channel#ingest</sla:type>
  </sla:header>
  <sla:packageAttribute>
    <sla:minSize unit="kilobyte">10</sla:minSize>
    <sla:maxSize unit="gigabyte">40</sla:maxSize>
    <sla:maximumNumberOfFiles>3000</sla:maximumNumberOfFiles>
    <sla:packageContent>
      <sla:formatCategory type="info:bnf/spar/representation#storedFormat"
        order="deny,allow">
        <sla:formatList action="deny"><format>*</format></sla:formatList>
      </sla:formatCategory>
      <sla:formatCategory type="info:bnf/spar/representation#managedFormat"
        order="deny,allow">
        <sla:formatList action="allow">
          <format type="ark">ark:/12148/ftiff_6_0w</format>
        </sla:formatList>
        <sla:formatList action="deny"><format>*</format></sla:formatList>
      </sla:formatCategory>
    </sla:packageContent>
  </sla:serviceLevelAgreement>
```


- For the digitized collection,
 - all formats are managed
- BUT we deal with 10 years of data creation
 - historical choices
 - changes in requirements
 - changes in scope
 - learning mechanisms

- Testing the formal method, we discover “bad” HTML files => not valid against the W3C validator
- BUT these files have been
 - produced for 2 years
 - displayed for 8 years
- Possible strategies:
 - correct the data before ingestion
 - accept “bad” data and plan for correction

Previous experience : SGML files

■ Previous experience :

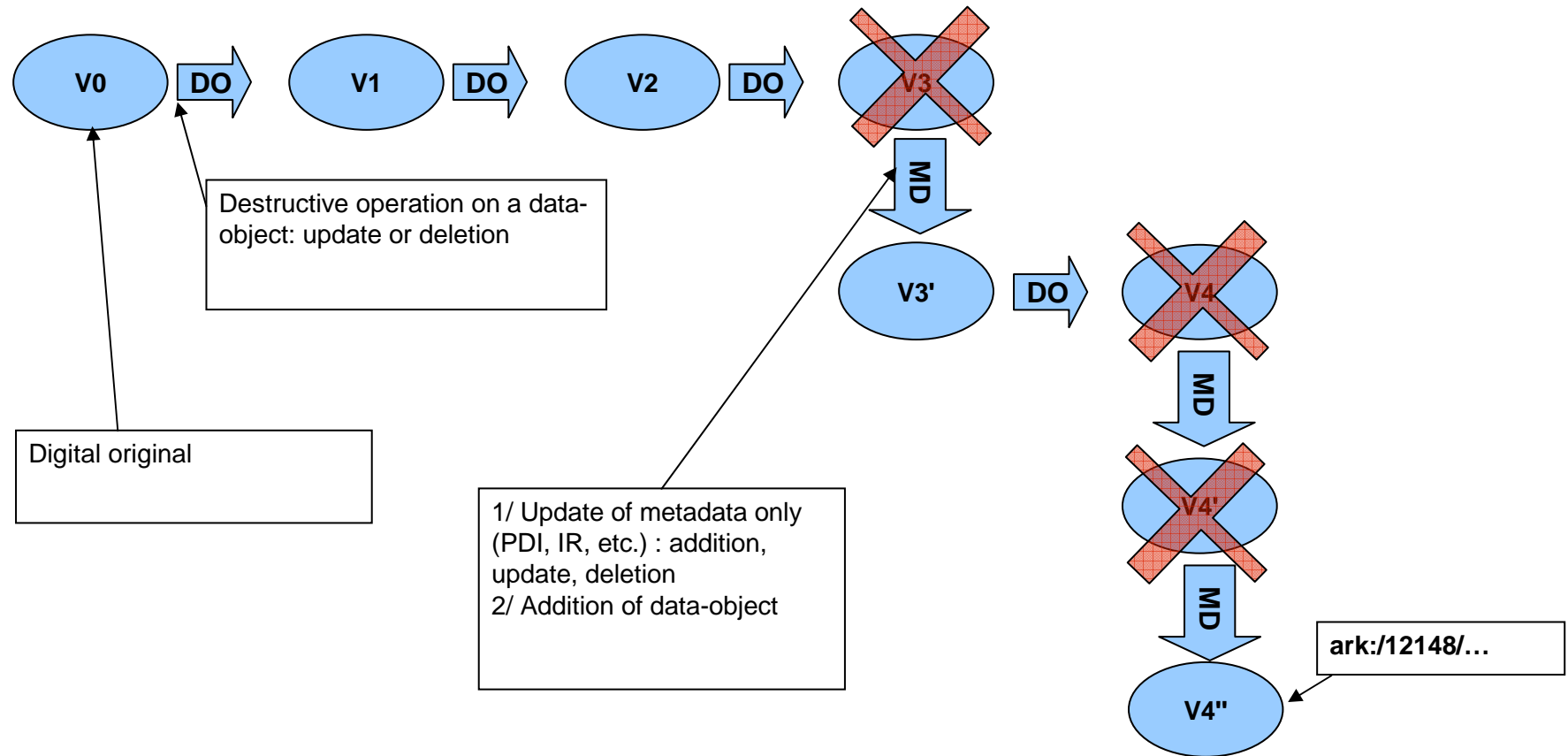
- production made on SGML format (special TEI profile)
- uncontrolled conversion made in XML (the current format at the time)
- at ingestion time, no way to display correctly the files

- need for *archeological* search, to come up with:
 - ❖ the original XML files
 - ❖ the original specifications (SGML DTD)
 - ❖ the expected XML files
 - ❖ the needed tools (SGML parsers)
 - ❖ the required skills

- The idea is to accept the files in a control way
 - Lax the requirements on the channel (some identified formats)
 - Use the ingest mechanisms to characterize the files as much as possible
 - Formally determine the set of “bad” data
 - Plan for *curation* with a well defined and documented migration plan

- Benefits
 - the history of transformation is preserved, the decision are traced and the tools used are known

Lifecycle of an Archival Package (AIP)



- Each version/edition has its own internal identifier
- The persistent identifier is unique, in ark: format

- Lesson 1:
 - never think your data is perfect
- Lesson 2:
 - the model works !!!
- Lesson 3:
 - migration is part of the preservation process
 - keep the original data and trace the operations made

Instructions to use the original tool



Agenda

- Missions and context
- From theory to practice
 - the model
 - the expectations
 - handling “bad” data
 - lessons
- Management through risk analysis
 - Focus shift
 - Dealing with human skills

- Risk management is at the core of preservation strategy
- Necessity for reassessing the risks on a regular basis
- Inclusion of the preservation system in the whole organization:
 - sustainability of the economics
 - add value to the service

- Building a smooth transition in two phases
 - ingesting and auditing
 - direct access
- BUT dissemination has become critical in our ecosystem:
 - access through Gallica
 - cooperation with Europeana
 - experimentation with publishers
- Risk to high to modify access at the moment

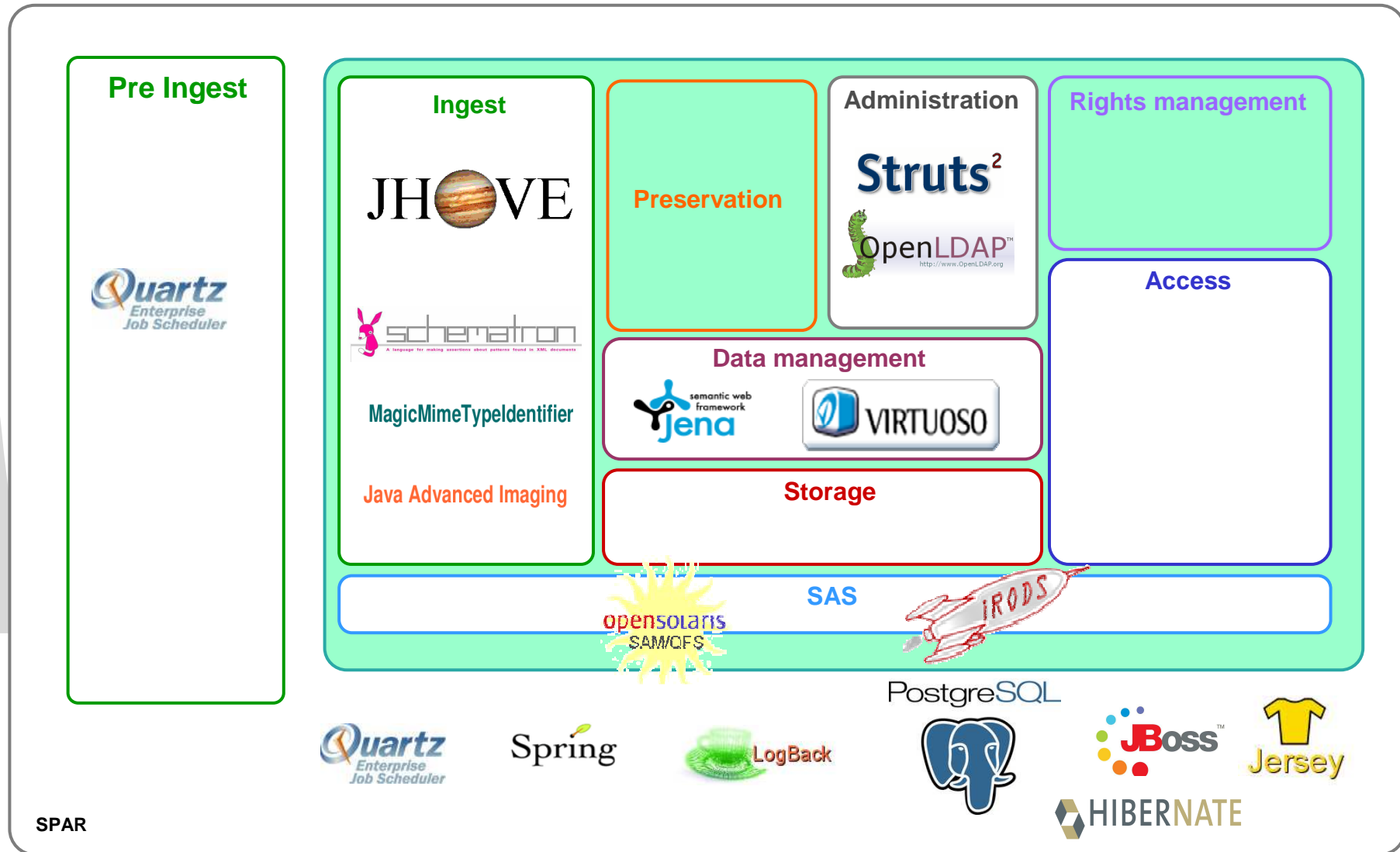


Focus shift to preservation repository

- The SPAR project focus shifts to
 - ingesting data
 - storing it safely
 - and managing it
- It acts as a digital repository for preservation

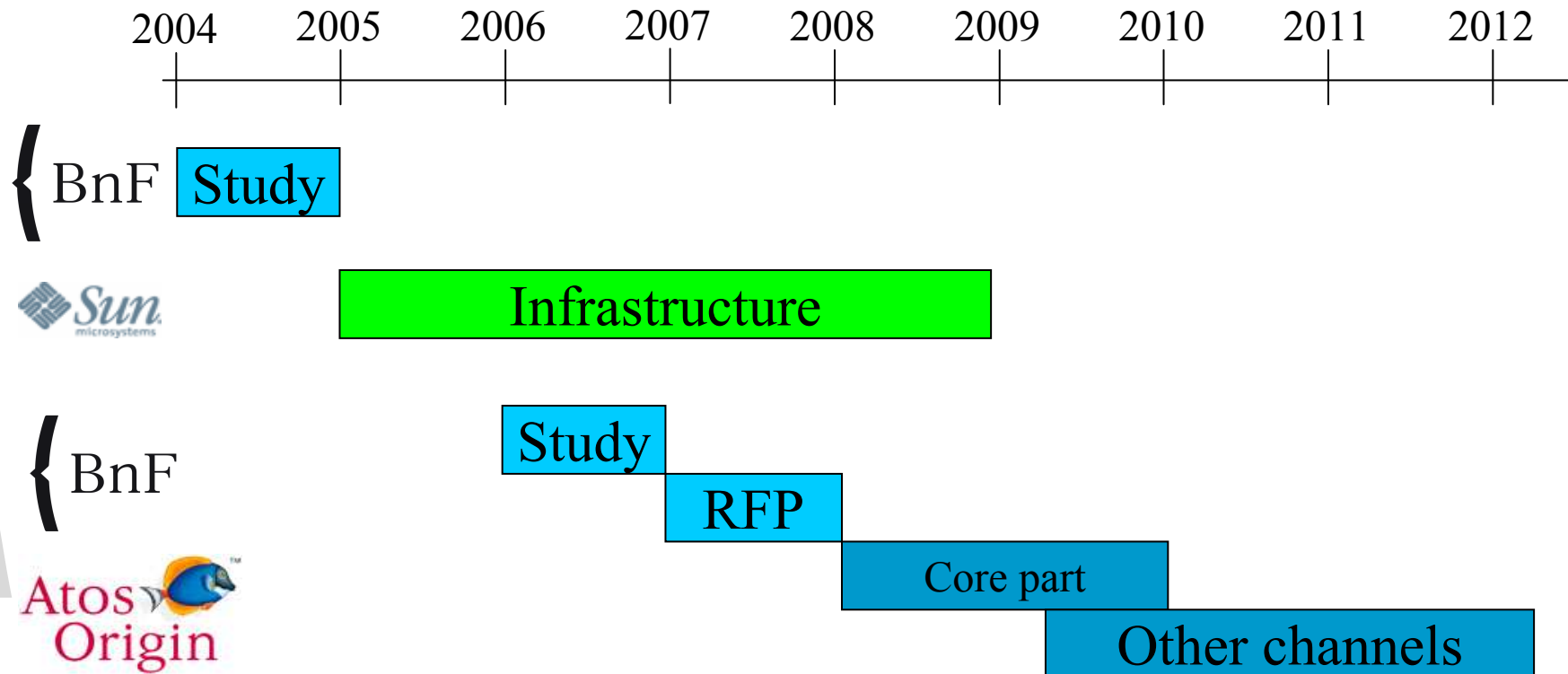
- Need for enhanced skills in the library
 - preservation expert: digital archeologist, model and formats
 - digital curator: management of digital collections
 - computer scientists: advise on tools
 - administrator: dealing with the data deluge
 - manager: understanding of digital issues, endorsement of channels

Implementation



SPAR

Planning



Thank you for your attention

Questions ?

More information : <http://bibnum.bnf.fr/spar>

Thomas Ledoux
thomas.ledoux_AT_bnf.fr