



# *Evaluation of data catalogue systems*

ELI ALPS-PaNOSC WP3 Catalogue Integration Best Practices Meeting 05/19/2021

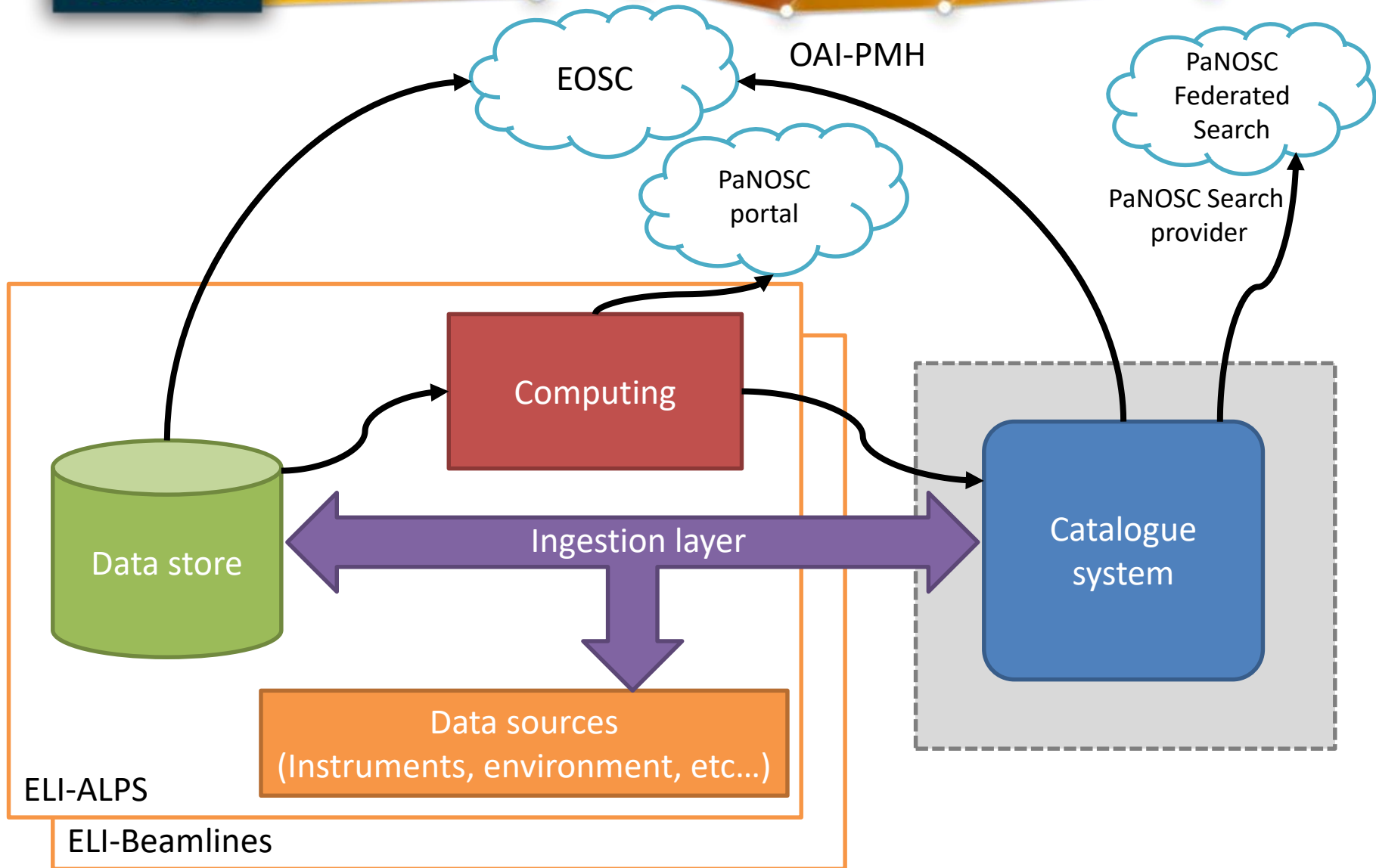
- **ELI-ERIC**
- **Planned data management architecture**
- **Evaluated catalogue systems**
- **The main aspects for the evaluation**
- **Aspects in details**
- **Summary**

- **Distributed research infrastructure**
  - ELI-ALPS, ELI-BL, ELI-NP will join in the future
- **New facilities**
- **Strategies**
  - Shared / common
    - policies
    - services
    - solutions
    - development infrastructure

KPIs - Before PaNOSC - 2018						
	ILL	ESRF	CERIC	XFEL	ELI	ESS
Data / yr	200 TB	8 PB	1 PB	3PB	< 1 PB	0
Data Policy	2011	2016	2014 (3/8)	2017	In progress	2017
Metadata catalogue	Local	Icat	Local	myMdC	No	SciCat
Metadata definitions	Nexus	Nexus	custom	myMdC	?	Nexus
DOI	2012	2018	No	2018	No	2018
Open Data	Yes	Yes	No	Yes	No	No
Data Services	Pilot	In progress	Remote	In progress	?	In progress
Common data API	No	No	No	No	No	No
User training	No	No	No	No	No	No



# Planned data management architecture



- **iCat**
- **SciCat**
- **Invenio RDM**
  
- **Other data catalogue systems:**
  - IBM Cloud Pak for Data
  - Oracle Cloud Infrastructure Data Catalog
  - Google Data Catalog
  - Apache Atlas
  - Omero



# Main aspects for evaluation

- **Data model**
- **Authentication and authorization**
- **Ingestion**
- **Search**
- **GUI**
- **Deployment and operation**
- **Documentation**

Sub-aspect \ System	iCat	SciCat	Invenio RDM
<b>Schema</b>	Based on the Core scientific data model	An inhouse developed schema	Extension of the DataCite schema
<b>Support for configurable schema</b>	Some element can have "parameters"	It supports parameters	Yes, with tutorial (it is recommended to use built-in schema)
<b>Investigation life-cycle support</b>	The model has Investigation and Publication, but the Proposal is not a part of the schema.	The model has support for Proposal and Published data and has a DatasetLifecycle element, which contains the executed operations on a Dataset.	Investigation life-cycle could be supported through the resource type property and the referenced identifiers.
<b>Interoperability with other schemas</b>	Dublin Core, DataCite (via OAI-PMH, configuration is facility dependent)	Dublin Core, DataCite (via OAI-PMH)	Dublin Core, DataCite (derived from DataCite)



# Authentication and authorization

Sub-aspect \ System	iCat	SciCat	Invenio RDM
<b>Authentication</b>	Simple, LDAP, OpenID Connect	Simple, SAML, LDAP, OAuth/OpenID Connect	Simple, SAML, LDAP, OAuth/OpenID Connect
<b>Authorization</b>	Rule based access control	Role based access control	Role based access control
<b>User and group management</b>	Yes	Yes	Yes



Sub-aspect \ System	iCat	SciCat	Invenio RDM
<b>API</b>	RESTful API (Java, Python packages)	RESTful API (openAPI/Swagger)	RESTful API (Python, CLI)
<b>HDF5 and NeXus support</b>	Yes	Yes	Yes
<b>Documentation</b>	Python tutorial	Tutorial (PSI, ESS)	Tutorial for Python and CLI

Sub-aspect \ System	iCat	SciCat	Invenio RDM
<b>Query system</b>	JPQL based (Jakarta/Java Persistence Query Language)	Loopback query (filtering)	Elasticsearch Query DSL (Domain Specific Language)
<b>OAI-PMH support</b>	Yes	Yes	Yes (provider and harvester)
<b>PaNOSC search provider</b>	Yes	Yes	N/A

Sub-aspect \ System	iCat	SciCat	Invenio RDM
<b>Web based GUI</b>	Topcat	Catanie	Invenio framework
<b>Technologies</b>	AngularJS – EOL Dec. 2021 (a React based solution in development)	Angular	React
<b>Customization</b>	A cookbook	N/A	Tutorials

# Deployment and operation

Sub-aspect \ System	iCat	SciCat	Invenio RDM
<b>High availability and scalability</b>	Yes*	Yes (Kubernetes)	Yes (Openshift, Kubernetes)
<b>Containerization</b>	It does not have official docker support.	Docker images	Docker images
<b>Technologies and languages</b>	MySQL, Oracle, MariaDB Glassfish, Payara, Java, Python(2!) EOL 2020	MongoDB, NodeJS, Rabbit MQ, Kafka, Loopback, Javascript	HAProxy, Nginx, Apache UWSGI, Gunicorn, NodeJS, Celery, PostgreSQL, Elasticsearch, RabbitMQ, Redis, Memcache Python
<b>Operation model</b>	Self service Multi-site support	Self service	Self service or commercial model: <a href="https://tind.io/">https://tind.io/</a> Multi-site support
<b>Documentation</b>	Installation process is not clear.	Docker based installation with CLI based on make	Installation is based on a simple python based CLI tool.

Sub-aspect \ System	iCat	SciCat	Invenio RDM
<b>General</b>	The documentation is fragmented.	A well written documentation, the development section builds on the documentation of the Loopback framework	Detailed documentation for both the framework and the RDM application as well.
<b>Tutorials</b>	Developer page: TODOs and “Page not found” errors	Developer guide: Contains tutorials	Developer page: Plenty of tutorials and even more for the main framework
<b>Online support</b>	Google group, Issue tracking on Github	Issue tracking on Github	Discord channel, Issue tracking on Github, LTS, Commercial support

	ICat	SciCat	Invenio RDM
<b>Pros</b>	<ul style="list-style-type: none"> <li>• used in PaN community</li> <li>• “PaNOSC ready”</li> <li>• „ERIC” support</li> </ul>	<ul style="list-style-type: none"> <li>• used in PaN community</li> <li>• “PaNOSC ready”</li> <li>• Compatible with modern solutions (containerization, kubernetes)</li> </ul>	<ul style="list-style-type: none"> <li>• Framework based</li> <li>• A customizable turn-key solution</li> <li>• Based on modern and exchangeable technologies/solutions</li> <li>• Good support</li> <li>• Commercial option</li> </ul>
<b>Cons</b>	<ul style="list-style-type: none"> <li>• EOL technologies (Angular JS, Python2)</li> <li>• Small developer community</li> <li>• Fragmented documentation</li> </ul>	<ul style="list-style-type: none"> <li>• Small developer community</li> <li>• „ERIC” support is not clear</li> </ul>	<ul style="list-style-type: none"> <li>• Under development – LTS in July 2021</li> <li>• “PaNOSC unready”</li> </ul>



*Thank you*