# Research Objects and WS Characterization

**Jose Enrique Ruiz**
**jer@iaa.es**

**November 29th 2012**
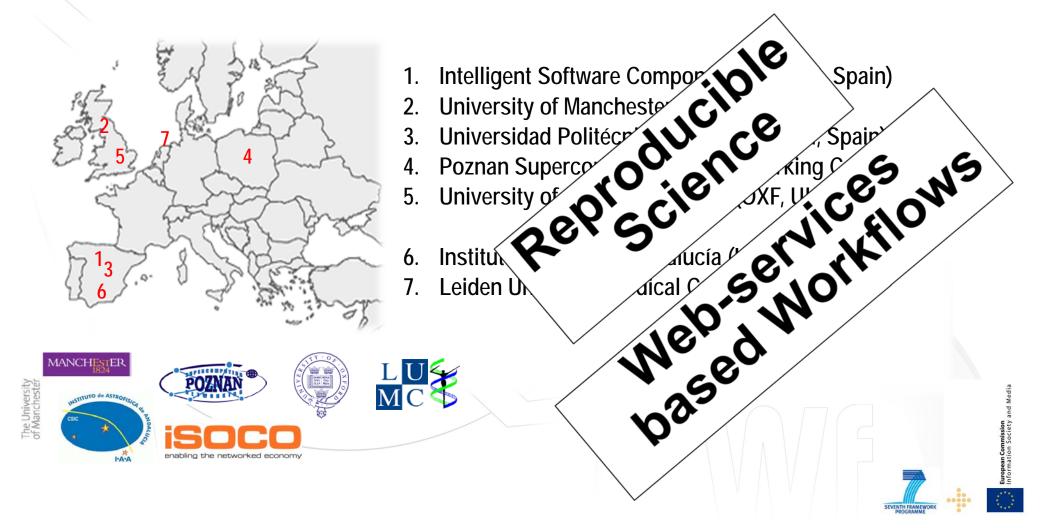**Workflows Group ASOV France**

# Wf4Ever
# Advanced Workflow Preservation Technologies for Enhanced Science

1. Intelligent Software Compon... ...Spain)
2. University of Manchester...
3. Universidad Politécn... ...Spain)
4. Poznan Superco... ...king C...
5. University of... ...OXF, U...

6. Institut... ...lucía (...
7. Leiden U... ...dical C...

**Reproducible Science**

**Web-services based Workflows**

Astronomy research lifecycle is **entirely digital**

» Observation proposals



» Data reduction pipelines



» Analysis of science ready data

» Catalogs of objects and data

» Publish process



› Final data results

› Experiment in DL

ADS/arXiv



**Reproducible research is still not possible in a digital world**

**A normalized preservation of methodology is needed**

**Efficient use of rich data infrastructure (VO) may be improved**

Tools
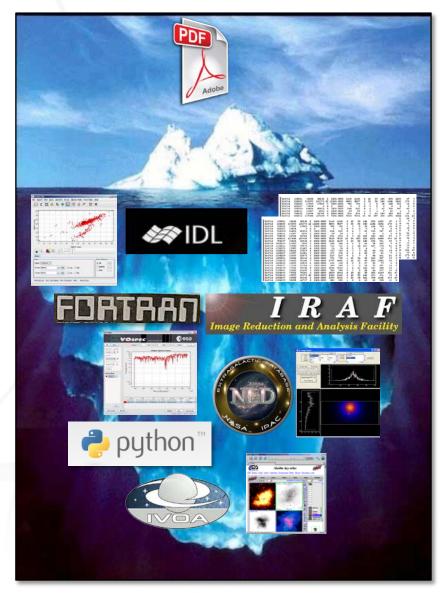
Optimize return on investments made on big facilities

- » Avoid duplication of efforts and reinvention
- » How to discover and not duplicate ?
- » How to re-use and not duplicate ?
- » How to make use of best practices ?
- » How to use the rich infrastructure of data ?
- » **Intellectual contributions are encoded in soft**

More data in archives does not imply more knowledge

- » Time has come to go beyond the PDF
- » Expose complete scientific record, not the story
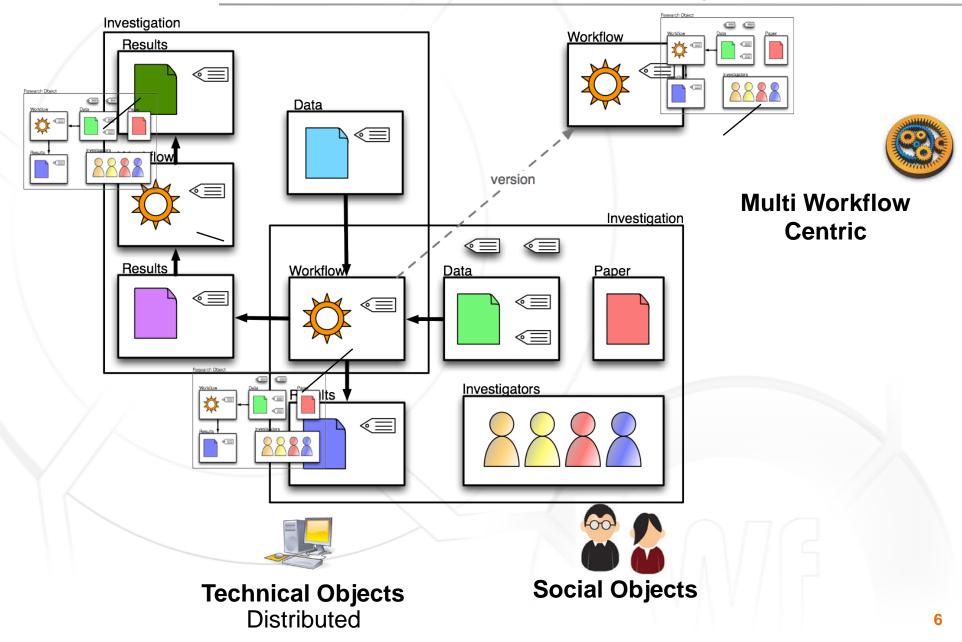- » Allow easy discovery of methods and tools

# Barriers to Data and Code Sharing in Computational Science

Survey of Machine Learning Community, NIPS (Stodden, 2010):

I don't know how

Tools

| Code | | Data |
|---|---|---|
| 77% | Time to document and clean | 54% |
| 52% | Dealing with questions from users | 34% |
| 44% | Not receiving attribution | 42% |
| 40% | Possibility of patents | - |
| 34% | Legal Barriers (ie. copyright) | 41% |
| - | Time to verify release with admin | 38% |
| 30% | Potential loss of future publications | 35% |
| 30% | Competitors may get an advantage | 33% |
| 20% | Web/disk space limitations | 29% |

**Multi Workflow Centric**

**Technical Objects**
Distributed

**Social Objects**

# RO Content

› Process (workflows), data, external resources and bibliography
› Execution environment set-up and local software dependencies
› Experimental protocol followed
› Roles, types and relationships among all digital components
› Provenance of intermediate and final results
› Decomposable attribution and authoring
› Fine-grained access control and permissions
› Example datasets for demonstration, reproducibility, monitoring, etc

# RO Template

› Placeholders to ease the aggregation process
› Completeness checking/quality assessment

# Semantic Annotations

- » Author of an **annotation**
- » **Author and co-authors** of a **workflow**; reference link to a re-used workflow and its author
- » Who has performed the **execution** of a workflow leading to the results provided in the RO
- » Computing execution environment of the RO and local software **dependencies**
- » Special **access requirements** to web services
- » Datasets **provider**: person, webpage, survey, data release, etc.
- » How much **time** does it take to run a workflow using the full data and the provided subsample
- » The number of **elements** of the sample dataset where one workflow and/or RO iterates
- » Previous and subsequent workflows to be executed, as in the experimental **protocol**
- » Research institution, country, and scientific domain of the RO
- » The actual **size** of the RO and/or a folder
- » The **version** of a workflow

*VOTable DataLink*

# Research Object Digital Library Architecture

**Workflow 4Ever**

User Clients

Extension Services

Foundation Services

**Access & Usage Functionalities**

| Use | Edit | Create | Annotate | ... |
|-----|------|--------|----------|-----|

**Data Management & Analysis Functionalities**

| Stability Evaluation | Completeness Evaluation | Recommenda-tions | Visualization | Collaboration | ... |
|---|---|---|---|---|---|

**Storage Functionalities**

| Storage | Retrieval | Maintenance | ... |
|---|---|---|---|

**Lifecycle Functionalities**

| Execution | Publication | Archival | ... |
|---|---|---|---|

# Research Object Digital Library Architecture

# Luminosity Profiles RO

LuminosityProfiles
  biblio
  config
  CONTENT.txt
  data
  LICENCE.txt
  process
  README.txt
  workflows

1010 Files, 200 MB
External Sources ~ 8 GB

5 Main Workflows, 14 Nested Workflows, 25 Scripts, 11 Configuration files
10 Software dependencies, 1 Web Service

Dataset: 90 galaxies observed in 3 bands

## Reproducibility

**When organization is better than automation**

# Credit and attribution
## Papers with data links are cited more than those without



Effect of E-printing on Citation Rates in Astronomy and Physics
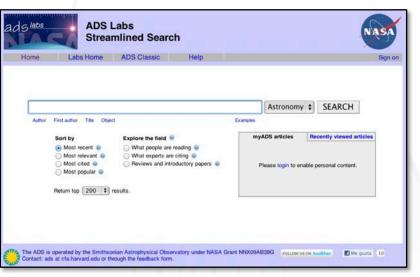2006. Edwin A. Henneken et al.

## ADSLabs Research Objects

## ADO Linked Components

» Authors

» Publications

» Journals

» Objects SIMBAD

» Tabular data behind the plots CDS

» ASCL reference of used software

» Observing time Proposals

» Used facilities, surveys or missions

# The next generation of archives

Much wider FoV and spectral coverage
- Huge sized datasets (~ tens TB)
- Big Data science highly dependent on I/O data rates
- Subproducts as virtual data generated on-the-fly

**We are moving into a world where**
- **computing and storage are cheap**
- **data movement is death**

# The next generation of archives

Much wider FoV and spectral coverage
- Huge sized datasets (~ tens TB)
- Big Data science highly dependent on I/O data rates
- Subproducts as virtual data generated on-the-fly

**The *move computing to data* paradigm**

Archives should evolve from data providers into **services providers**, where web services may help to solve bandwidth issues.

# The next generation of archives

Much wider FoV and spectral coverage
- Huge sized datasets (~ tens TB)
- Big Data science highly dependent on I/O data rates
- Subproducts as virtual data generated on-the-fly
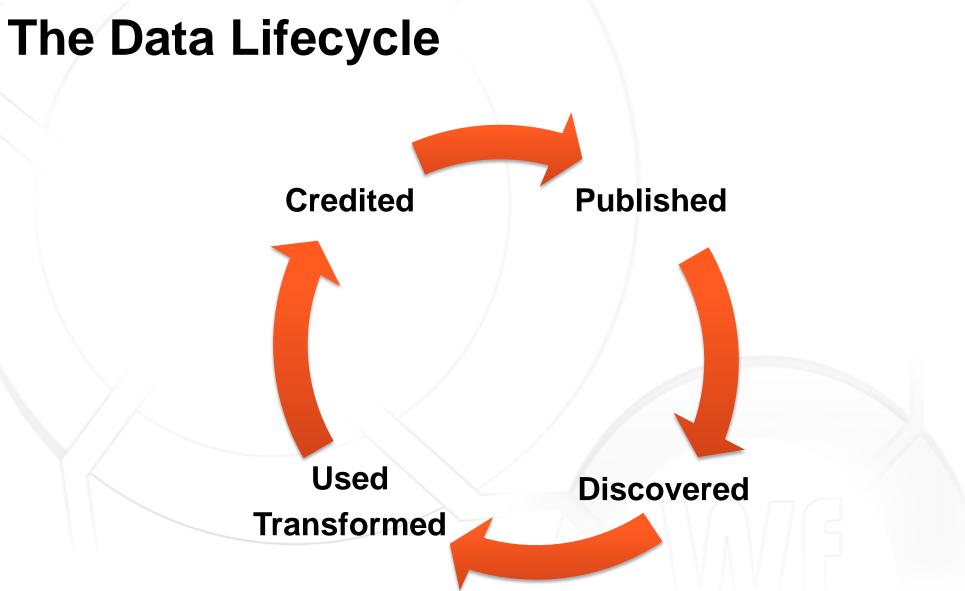
**Data Discovery**
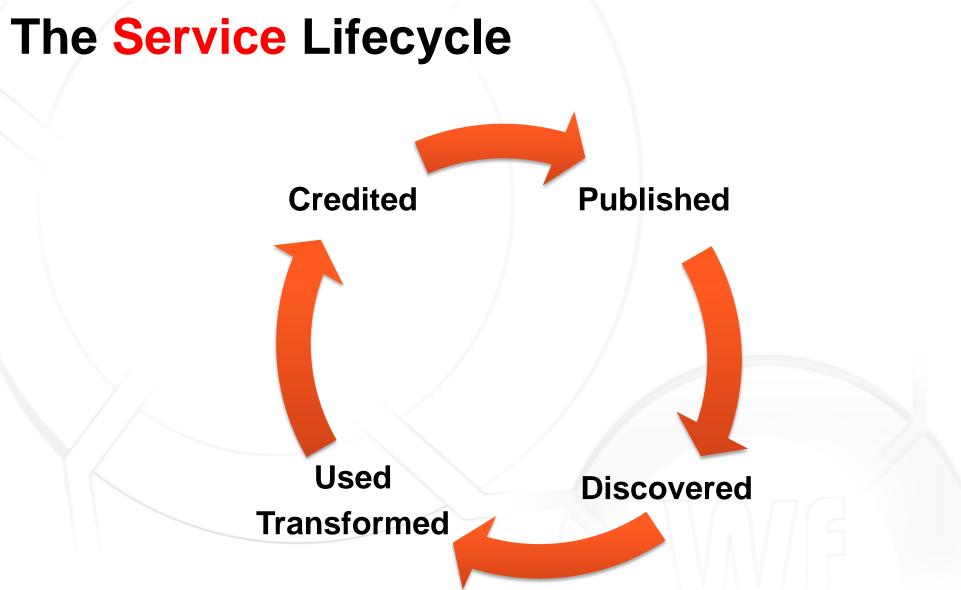**Data Access**
**Data Management**
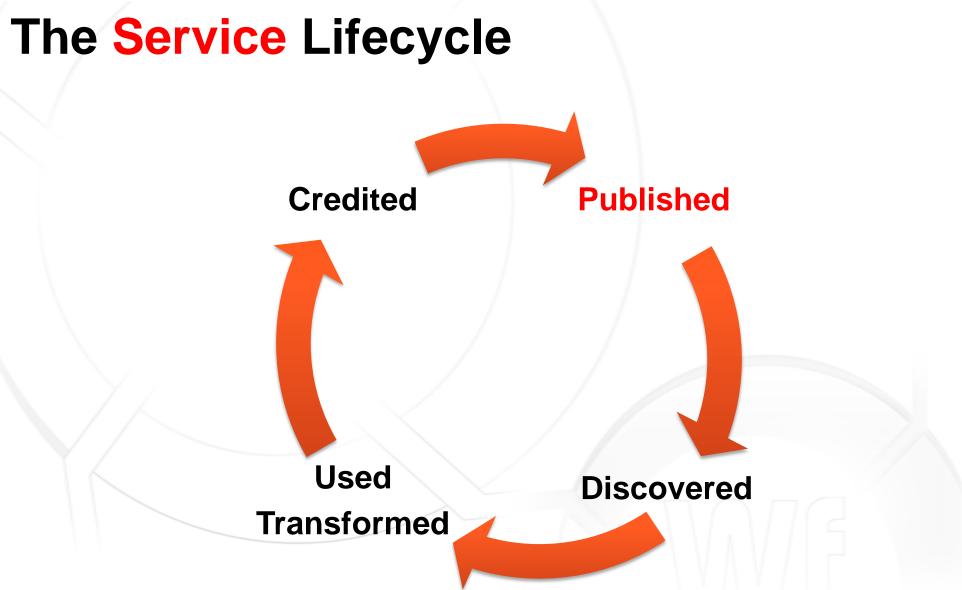
# The next generation of archives

Much wider FoV and spectral coverage
- Huge sized datasets (~ tens TB)
- Big Data science highly dependent on I/O data rates
- Subproducts as virtual data generated on-the-fly

**Web Services Discovery**
**Web Services Access**
**Web Services Management**

The Web-
Services
Deluge

# The Data Lifecycle



**Credited**

**Published**

**Discovered**

**Used
Transformed**

# The **Service** Lifecycle

# The Service Lifecycle

**Credited**

**Published**
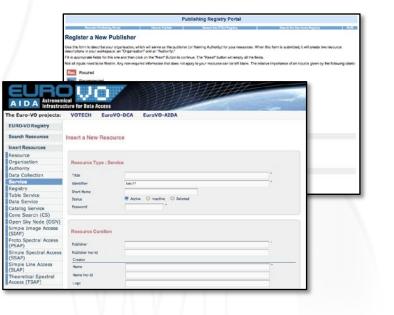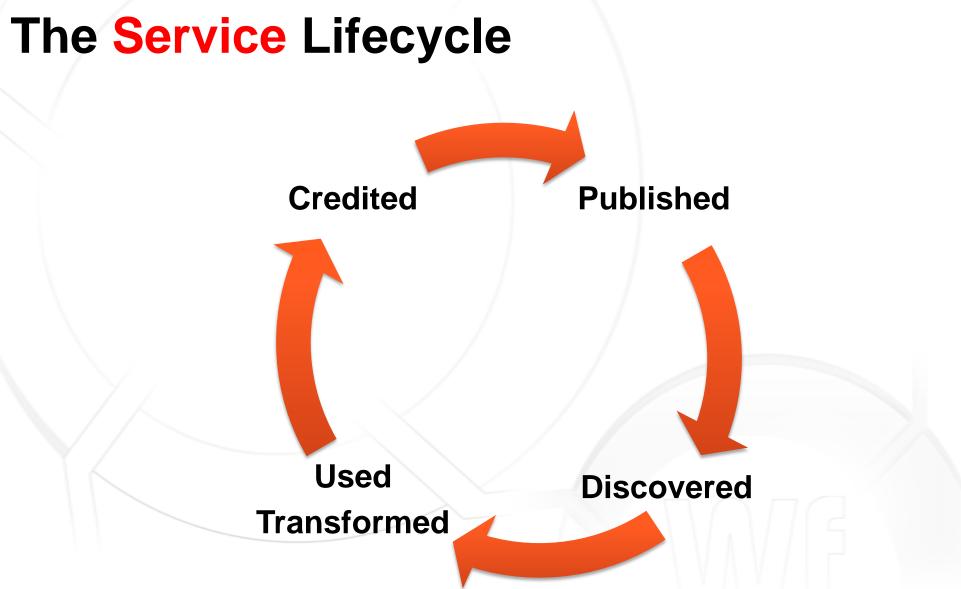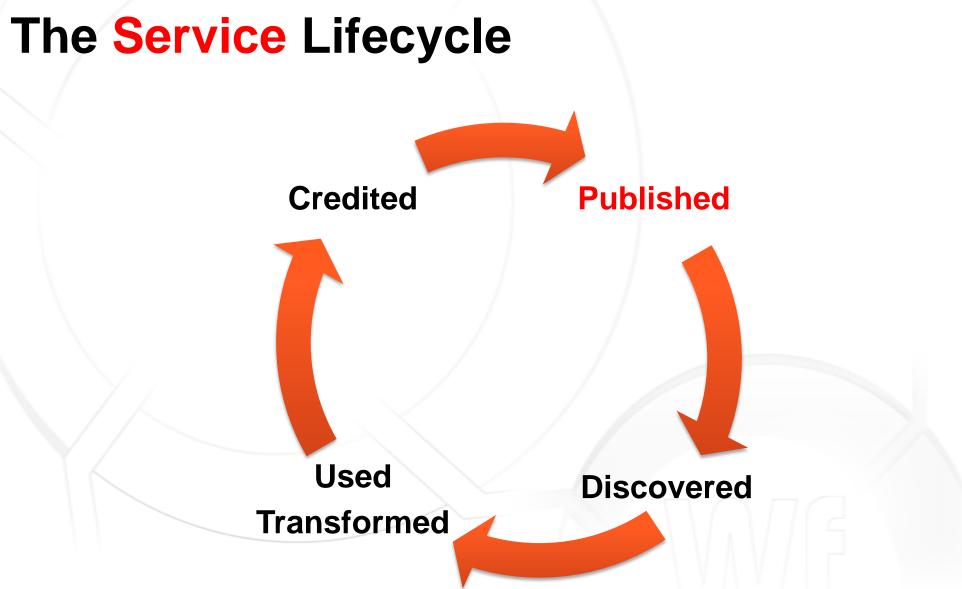
**Discovered**

**Used Transformed**

# Published

- The VO Registry
- Easier to publish services than datasets in the VO ?
- WS are not exclusive property of big data archives

- Publication is not Preservation
- Backup strategies
- Replication/Mirrors
- Versioning
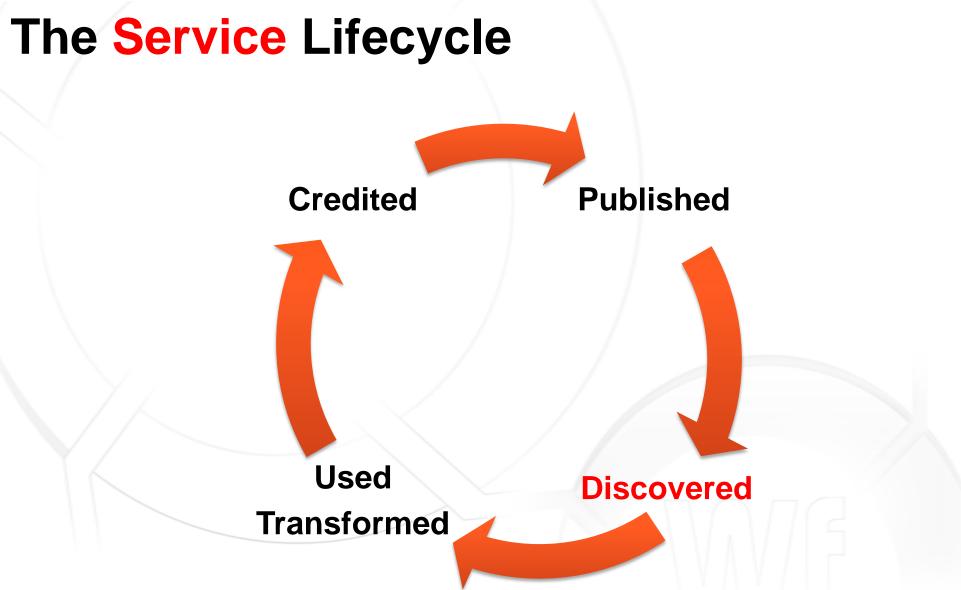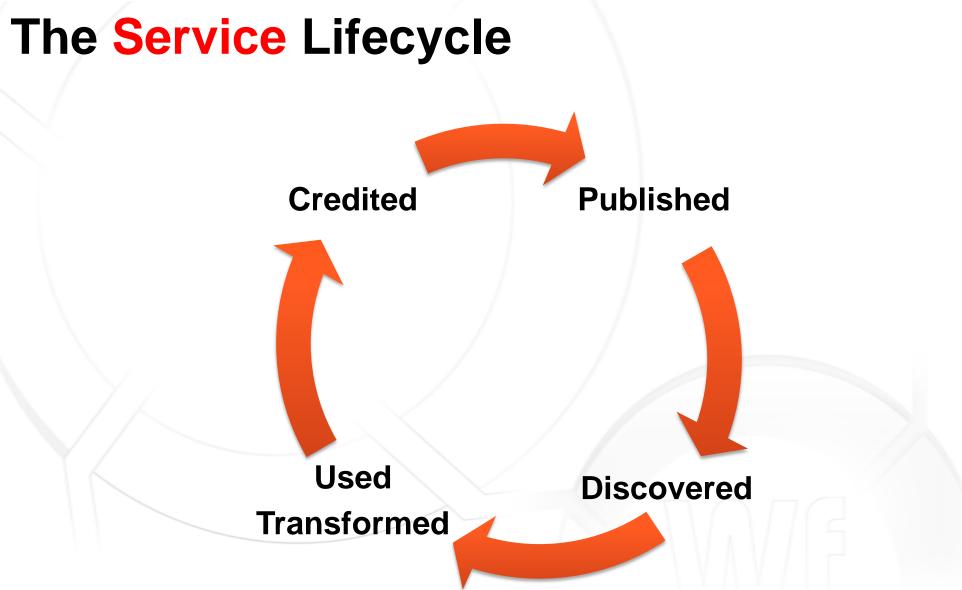
- Software Publishing Platforms

# The Service Lifecycle

**Credited**

**Published**

**Discovered**

**Used Transformed**

# The Service Lifecycle



Credited

Published

Used
Transformed

Discovered

# The Service Lifecycle



Credited

Published

Used
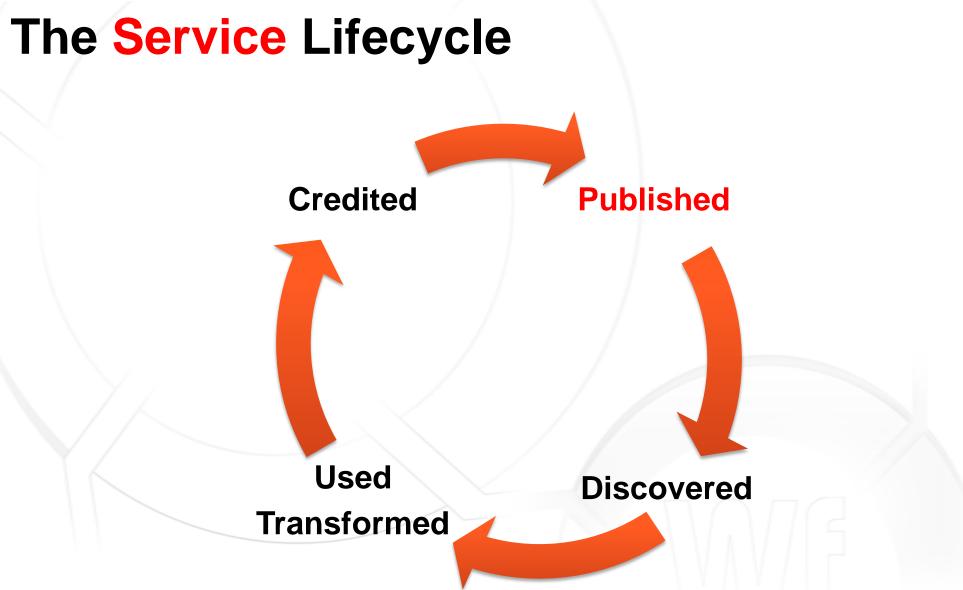Transformed

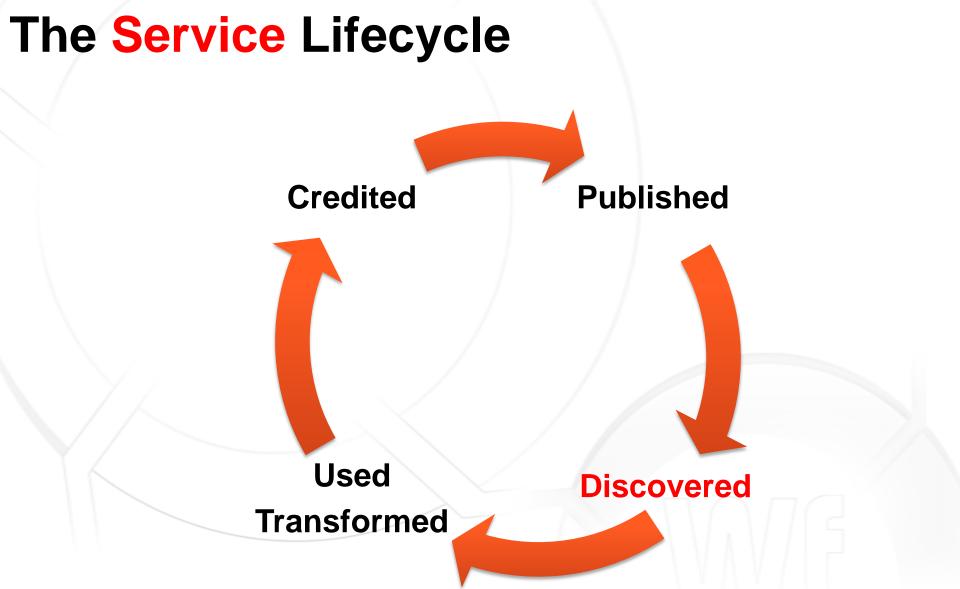Discovered

# Discovered

- **Search Criteria**
  - Relevant Keywords (Semantics)
  - Authoring - Institution, Archive
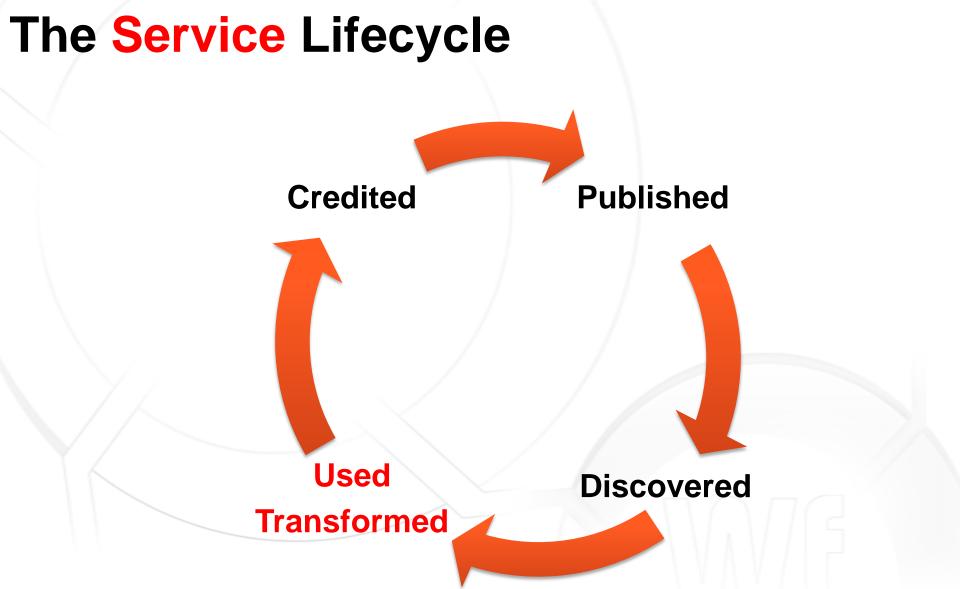  - Waveband, Science
  - Function-based
    - VO Services mainly focused on Data Discovery and Access (DAL)
    - Wrapped Legacy Apps and Data Processing (SIAv2, Theory IG)
    - KDD IG
  - Input/Output Data (TAP, UTypes, VOSI #tables)
  - Access Policy (Authentication – SSO, OAuth)
  - A-Synchrony (SOAP, REST) and Stage Data (VOSpace)
  - Allocation of CPU/Storage, Estimated Computing Time

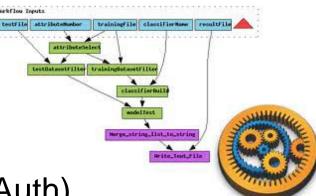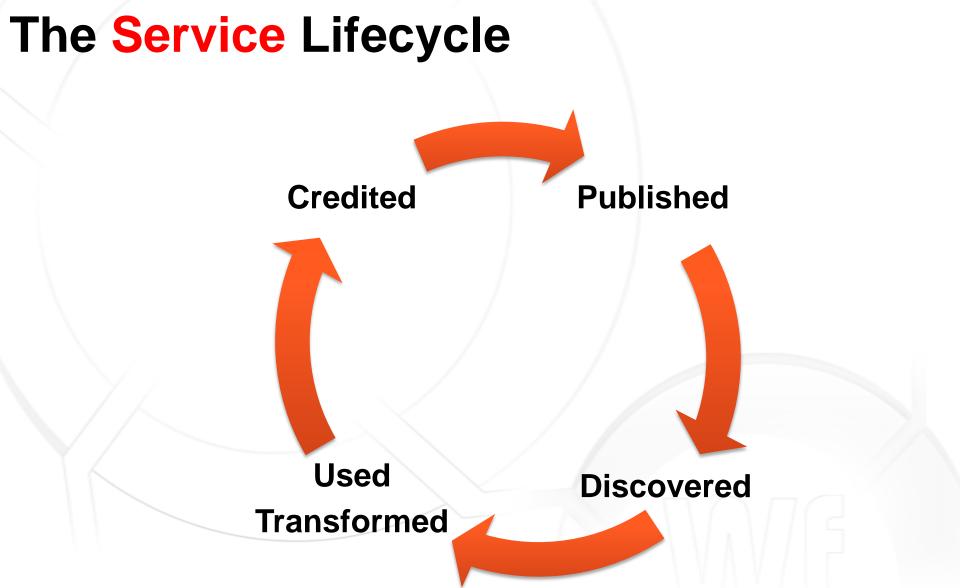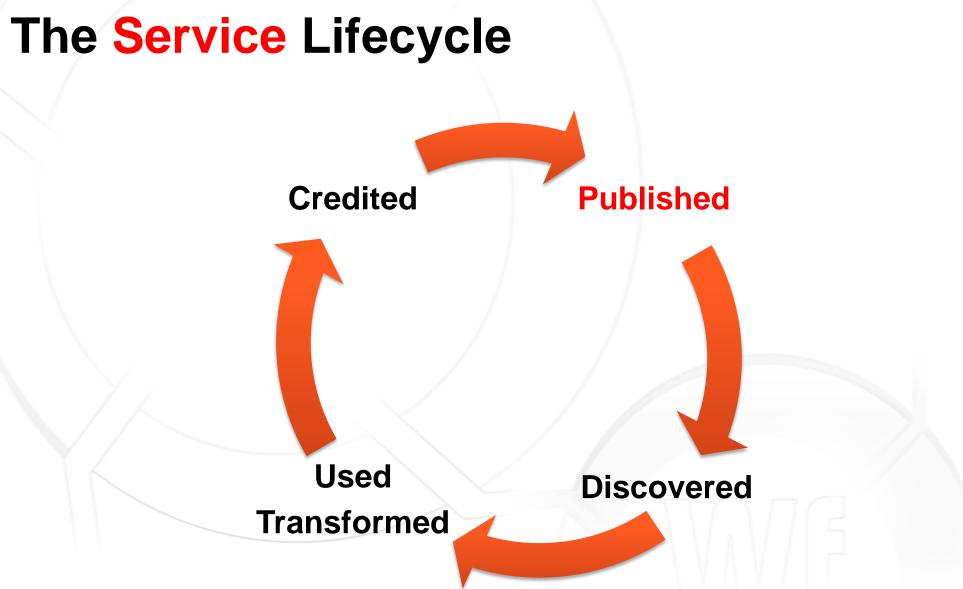- **Proposition of alternatives and similars**

# The Service Lifecycle



Credited → Published → Discovered → Used Transformed → Credited

# The Service Lifecycle



Credited

Published

Used Transformed

Discovered

# The Service Lifecycle

# The <span style="color:red">Service</span> Lifecycle



**Credited**

**Published**

**Used Transformed**

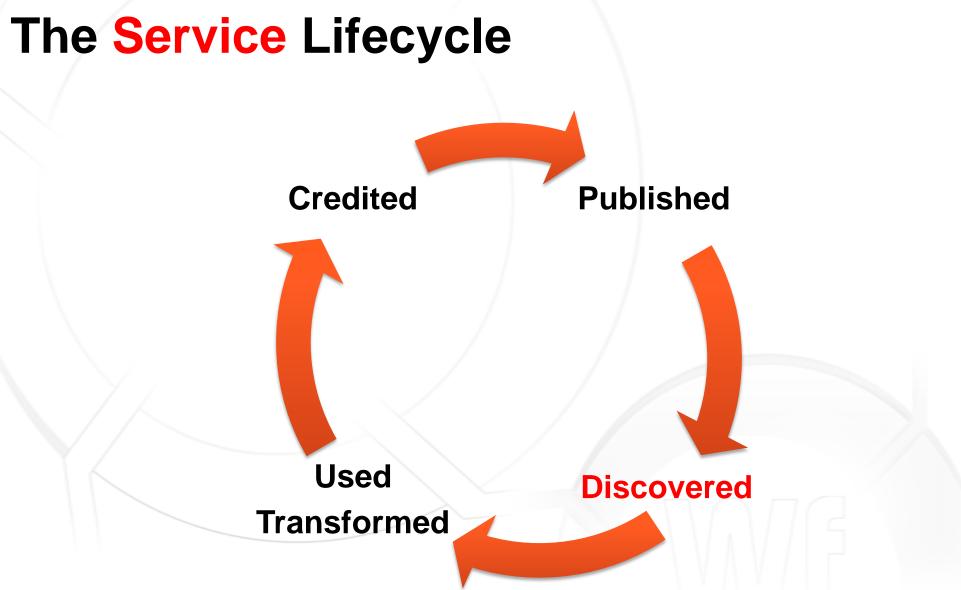**Discovered**

# Used and Transformed

- **How to use them ?** (WADL, WSDL – VOSI #capabilities)
  - Input Data -> Parameters needed and formats
  - Self-described WS (PDL, S3, SimDAL, SimDB)
  - Output Data -> Response format - TAP
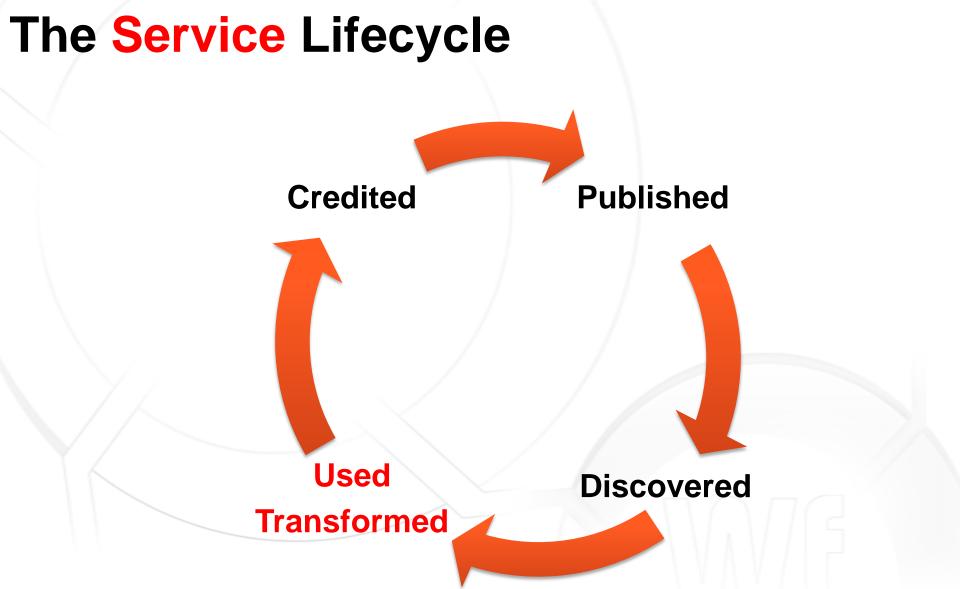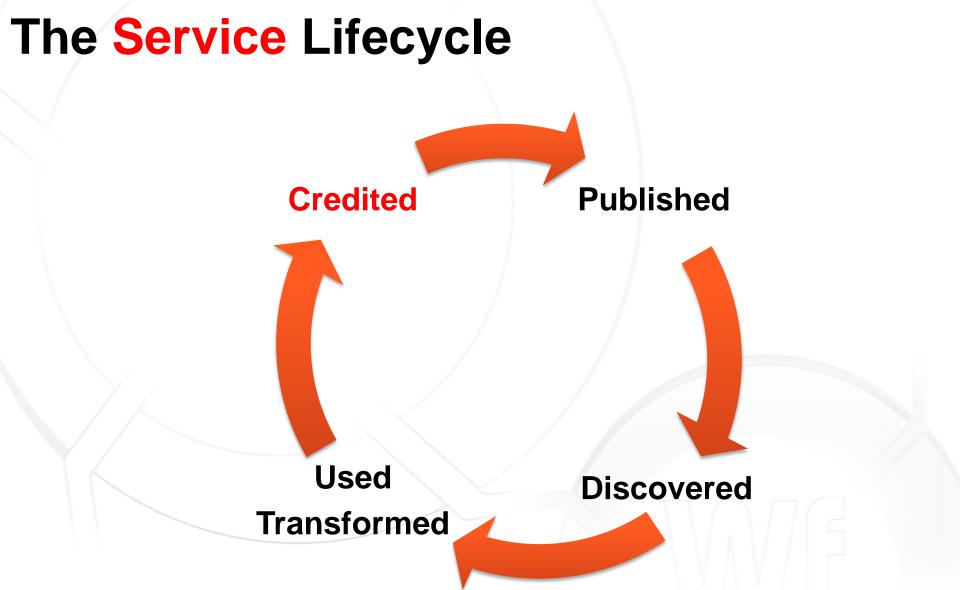  - Example Data, Self-Consistency Checking



- Access Policy (**Authentication** – SSO, OAuth)
- **WS orchestration in Workflows** (Data-flow vs. Control-flow)
- How the **community** uses WS ?
- Propositions based on patterns of statistical use or popularity
- **Provenance** of the methods is Wf-evolution by re-use
- Consumed by Humans and Machines - **Interoperable** (WS-I)

# The Service Lifecycle



Credited

Published

Discovered

Used
Transformed

# The Service Lifecycle



Credited

Published

Used Transformed

Discovered

# The Service Lifecycle

**Credited**

**Published**

**Discovered**

**Used Transformed**

# The **Service** Lifecycle



**Credited**

**Published**

**Discovered**

**Used Transformed**

# The Service Lifecycle



Credited

Published

Used
Transformed

Discovered

# Credited



- **Linked to related Artefacts**
  - Data Facilities and Archives
  - Authors, ASCL Software, Wfs
- **Quality Assessment**
  - Technical and scientific
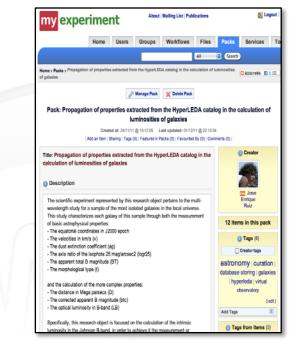  - Penalize abandoned and award the maintained
- **Automate Monitoring (VOSI #availability)**
  - Decay
  - Performance, WS Analytics
  - Modifs. on interfaces, permissions, etc.
- **Community Curation**
  - Blogging
  - Recommendation
  - Folksonomy

In a cloud of web services and data..
**Web Services should benefit of the *same privileges* acquired by Data until now.**

Start thinking on how to provide
- **Detailed curation**
- **Thorough characterization**