# Digital Science

## Reproducibility and Visibility in Astronomy

José Enrique Ruiz on behalf of the Wf4Ever Team

**SCIOPS 2013**
**ESAC, FRIDAY 13th SEPTEMBER 2013**

# Wf4Ever
# Advanced Workflow Preservation Technologies for Enhanced Science
## 2011 - 2013

1. Intelligent Software Components (ISOCO, Spain)
2. University of Manchester (UNIMAN, UK)
3. Universidad Politécnica de Madrid (UPM, Spain)
4. Poznan Supercomputing and Networking Centre (Poland)
5. University of Oxford and OeRC (OXF, UK)

6. Instituto Astrofísica Andalucía (IAA-CSIC, Sp...
7. Leiden University Medical Centre (LUMC...

Reproducible Science

# Astronomy Research Lifecycle

Astronomy research lifecycle is **entirely digital**

- » Observation proposals
- » Data reduction pipelines
- » Analysis of science ready data
- » Catalogs of objects and data archives
- » Publish process
  - › Final data results
  - › Experiment in DL
    ADS/arXiv

**Reproducible research is still not possible in a digital world**
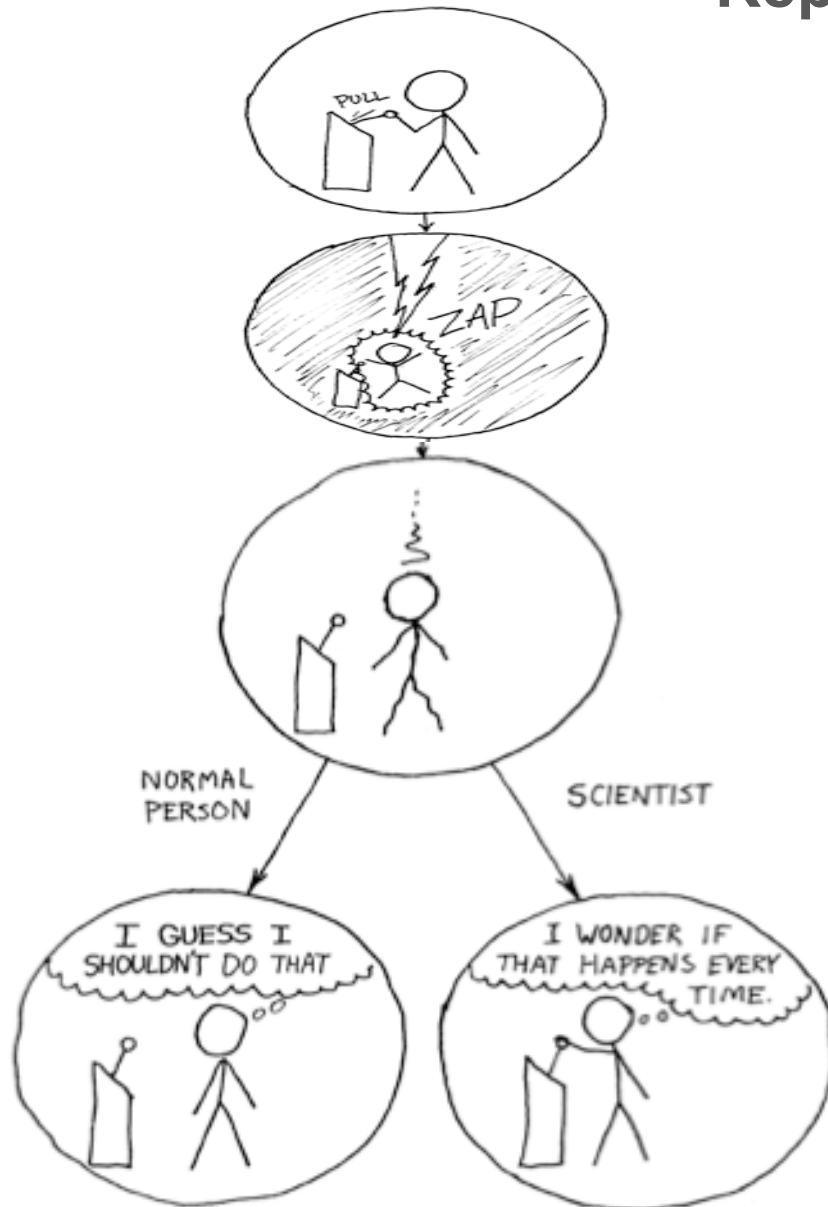
**A rich infrastructure of data is not efficiently used**

**A normalized preservation of methodology is needed**

Tools

## Reproducibility and The Scientific Method



**Benefits**

» Publishing knowledge, not advertising

» The author, the referee, the re-user

» Reputation, prestige and respect

» Higher quality of publications

  › Authors will be more careful

  › Many eyes to check results

**Challenges**

» Hard and time consuming

» Need incentives – not rewarded now

http://xkcd.com/242/

# Digital Science - Reproducibility and Visibility in Astronomy
## Visibility, Efficiency and Reuse

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 contribs are encoded in software**

More data in archives does not imply more knowledge

» Expose **complete scientific record**, not the story

» Allow easy **discovery** of methods and tools

# Paper discovery: the social dimension

## Time has come to go beyond the PDF

# Digital Astronomy in the Local Desktop

Capture
Actions, Tasks, Dependencies, Provenance

Improve
Clarity and Reproducibility

Living Tutorials
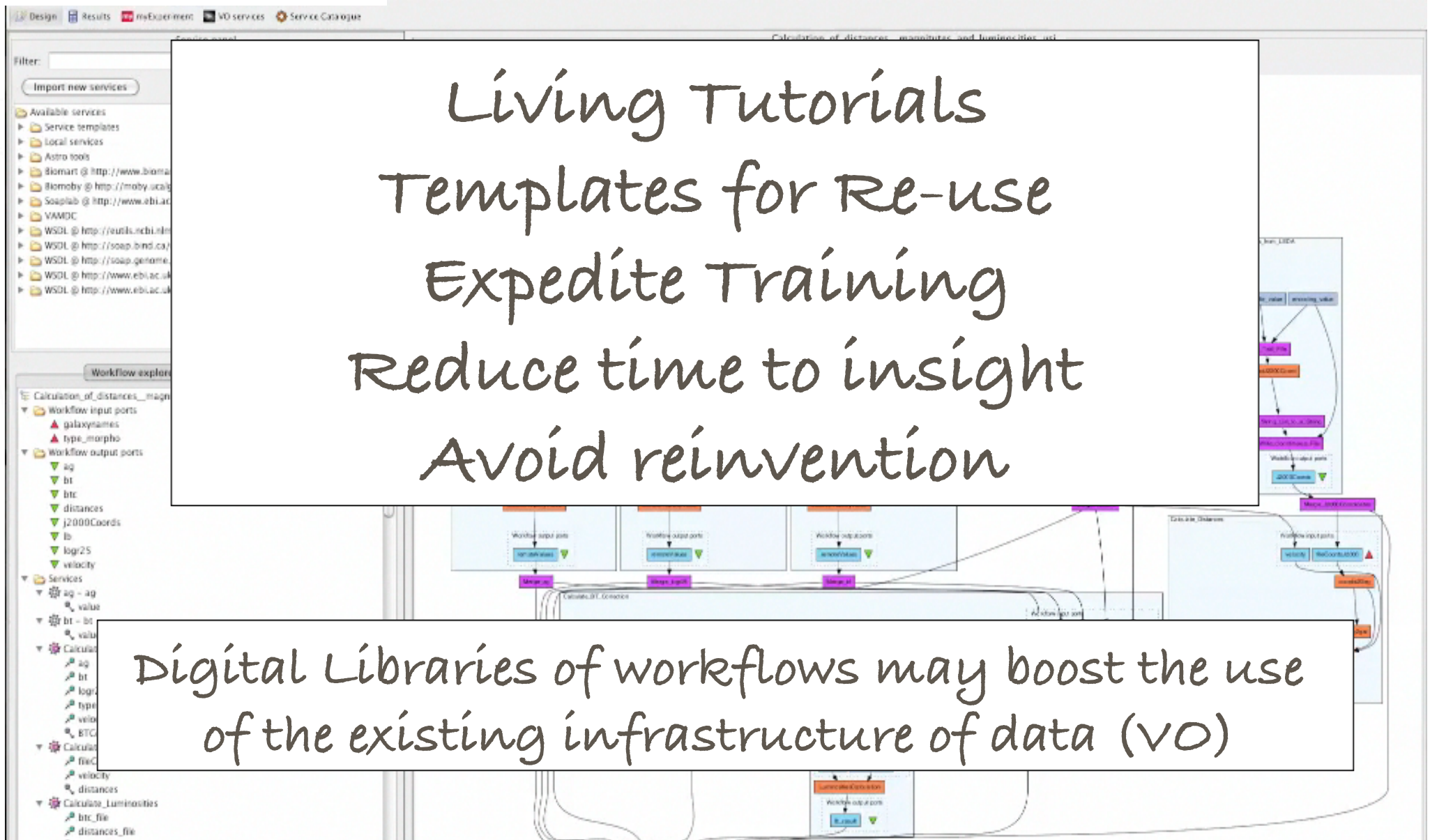Templates for Re-use
Expedite Training
Reduce time to insight
Avoid reinvention

Digital Libraries of workflows may boost the use
of the existing infrastructure of data (VO)

# Scientific Workflows

## Related Initiatives

- › ER-Flow
- › VAMDC
- › **HELIO**
- › Cyber-SKA
- › IceCore
- › Montage
- › **Astro-WISE**
- › AstroGrid

## Software

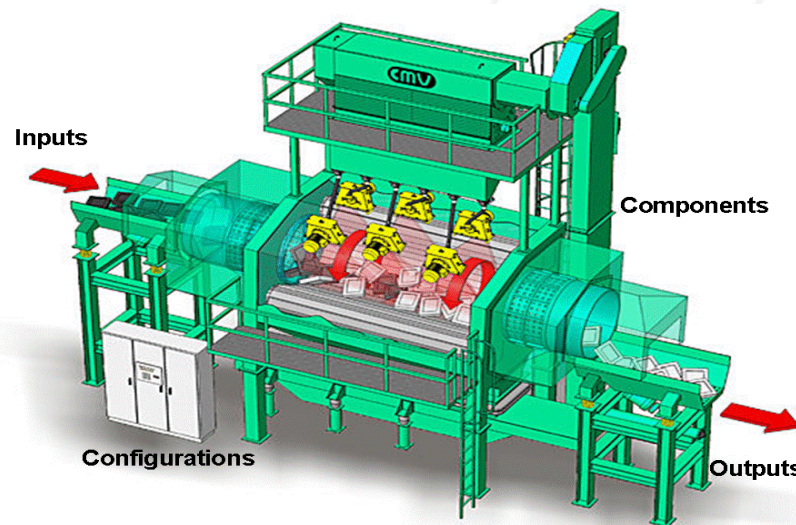- › Taverna
- › Kepler
- › Pegasus
- › Triana
- › **ESO Reflex**

## IVOA

- › AstroGrid
- › Grid&WS WG
- › VO France Wf WG
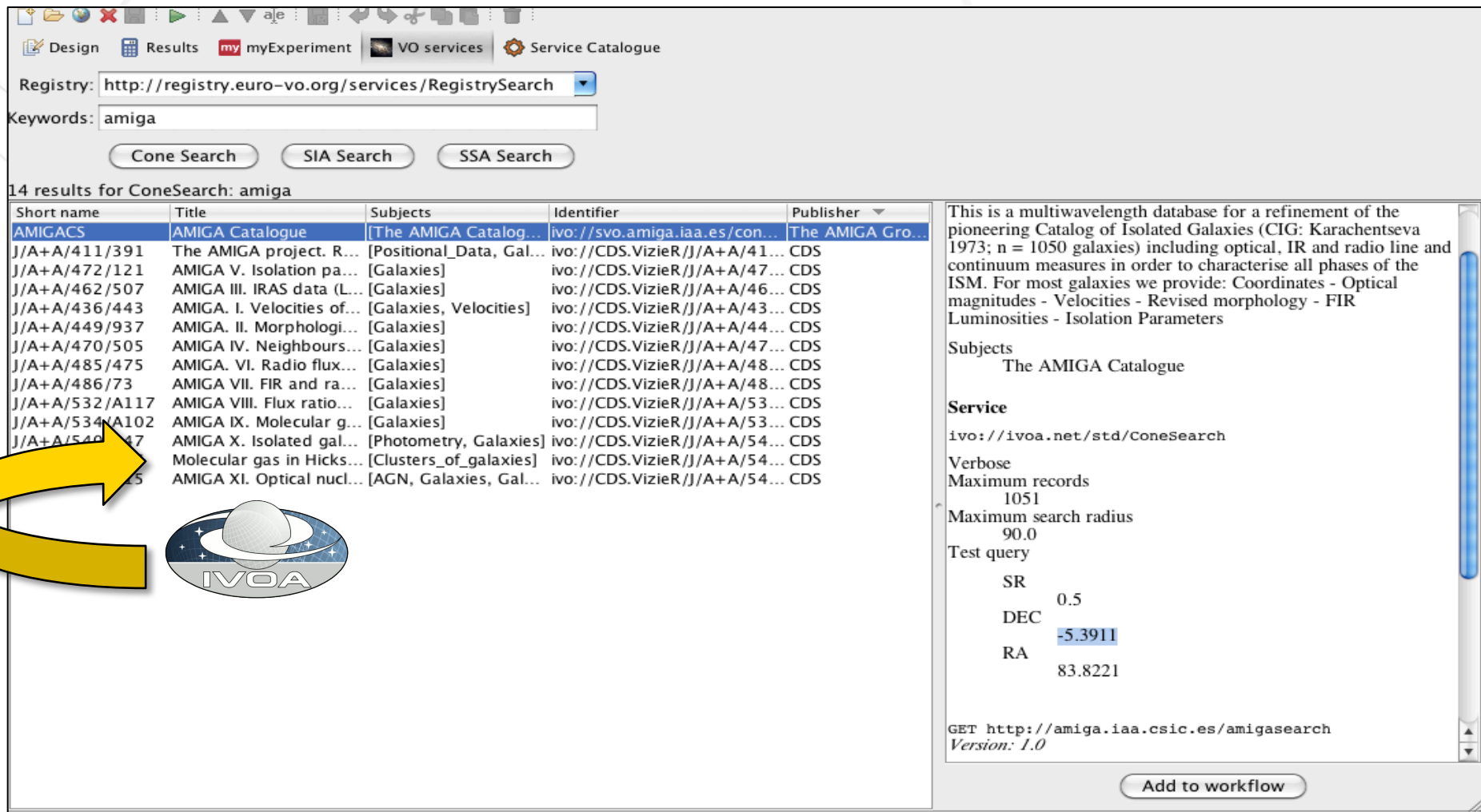
## Self descriptive WS

- › PDL
- › SimDAL, S3

Inputs

Components

Configurations

Outputs

*Interoperability Standards*

## Astronomical Research Objects in Action

## AstroTaverna: Create, annotate and run a workflow



http://amiga.iaa.es/p/290-astrotaverna.htm

# AstroTaverna: Create, annotate and run a workflow



http://amiga.iaa.es/p/290-astrotaverna.htm

## ASKAP Datacubes

| | Low Res | | High Res | | Extreme Res | |
|---|---|---|---|---|---|---|
| Number | 4 Bytes | 4B | 4 Bytes | 4B | 4 Bytes | 4B |
| Resolution | 2,048 x 2,048 | 16MB | 8,192 x 8,192 | 268MB | 12,288 x 12,288 | 603MB |
| Channels | 16,384 | 0.27TB | 16,384 | 4.39TB | 16,384 | 9.8TB |
| Stokes & Weighting | 1 | 0.27TB | 1 | 4.39TB | 4 + 1 | 49.5TB |

Prof. Kevin Vinsen

## SKA Datacubes

### Spectral Line Datacube

- Dish
    - Assume 30,000 channels
    - 27,000 x 27,000 x 30,000 x 4
    - ≈80TB
- AA
    - Assume 40,000 channels
    - 28,000 x 28,000 x 40,000 x 4
    - ≈125TB
- Stokes parameters and Weighting Map
    - Multiple by 5
    - Dish ≈ 400TB
    - AA ≈ 625TB

Prof.  Kevin Vinsen

# Much wider FoV and spectral coverage

» Large volumes for a single observed dataset

# Automated surveys

» Huge amounts of tabular data

Extraction of scientifically relevant info from a multiD param. space

» Exploration services

» Anomaly detection

» Cross-matching data

» Dimensionality reduction

Detailed inspection and subset

» Filtering

» Extraction

» Re-Projection

» Analysis services

**We are moving into a world where**

» computing and storage are cheap

» **data movement is death**

17

## The *move computing to data* paradigm

» A cloud of Web Services

Archives should evolve from ____ ____ into

» Virtual Data providers

» Software Tasks p____

» Archives speaking W____

Astronomy ____ s/facilities/wavelength

Interc____ ____roperable archives

» ____ ____bservatory

» ____ ____sks

*Web Services based Scientific Workflows*

**Preservation**

**Process should benefit of the same privileges acquired by data**

Preserving the method ensures replication of final results at any moment

**Expose experimental context in a structured way in order to be understood**



Distributed

Technical Objects

Social Objects

## IPython Notebook solutions

» Web-browser as the working desktop

» Python code, plots and data, living with rich-text documentation

» Cloud-based adaptive scalable computing environment

» Fully shareable, re-usable and executable wikis

» Social platform and Git versioning

*Similar Initiative to ESO Telbib*

## ADSLabs

## 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



Incentives

http://labs.adsabs.harvard.edu/

# The Incentive

Papers with data links are cited more than those without



1995 - 2000

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

## The Incentive

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.

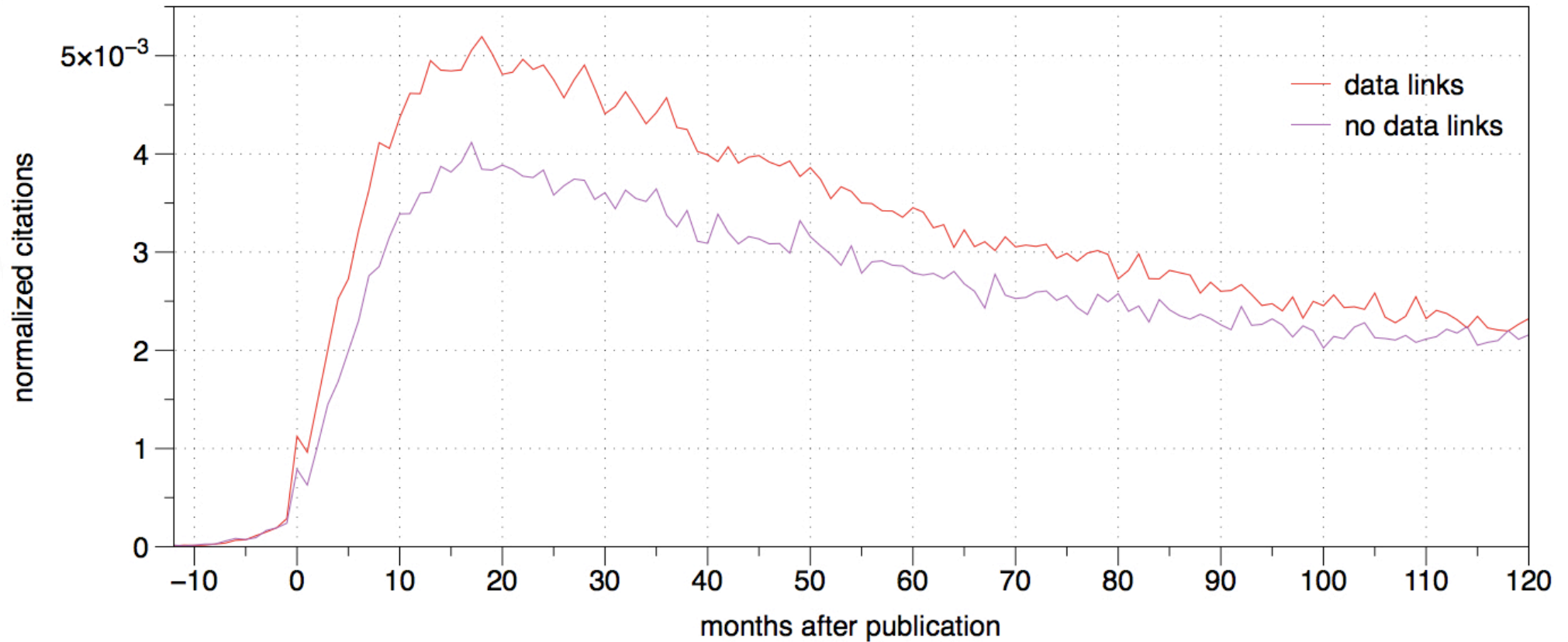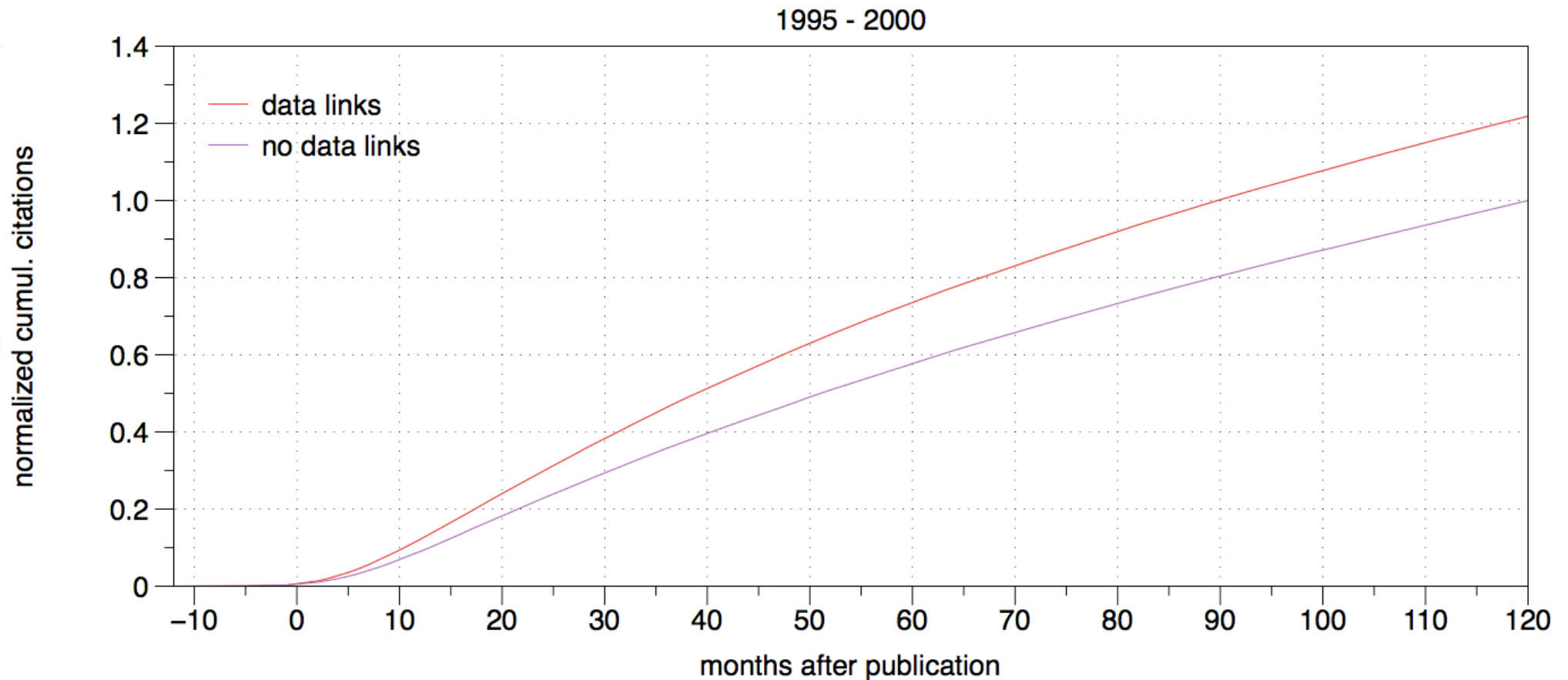» Reproducibility is at the very heart of the scientific method

» Improving visibility is key in order to avoid reinvention

» Social dimension of science stressed in the discovery process

» Highly specialized science needs re-use to achieve efficiency

» In a digital world, publish decomposable executable papers

» Capture provenance and structure in the local desktop

» Scientific workflows go beyond automation: provide clarity and structure

» Transfer rate is more than an issue for next generation of archives

» The move computing to data paradigm -> back to old terminals

» Process should benefit of the same privileges acquired by data

» Digital libraries of web-services-based workflows

» The distributed digital workflow-centric Research Object

» Preserving knowledge - not only data or advertising

jer@iaa.es