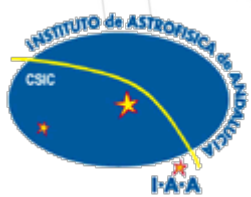


Digital Science

Going beyond automation

José Enrique Ruiz, Lourdes Verdes-Montenegro, Susana Sánchez,
Julian Garrido, Juan de Dios Santander and the Wf4Ever Team



SESIONES CCD
GRANADA, SEPTEMBER 26th 2012



Digital Science - Reproducibility and Visibility in Astronomy

Astronomy Research Lifecycle

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 rich infrastructure of data (VO) is not efficiently used



A normalized preservation of methodology is needed

Tools

Digital Science - Reproducibility and Visibility in Astronomy

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 contributions are encoded in softw**

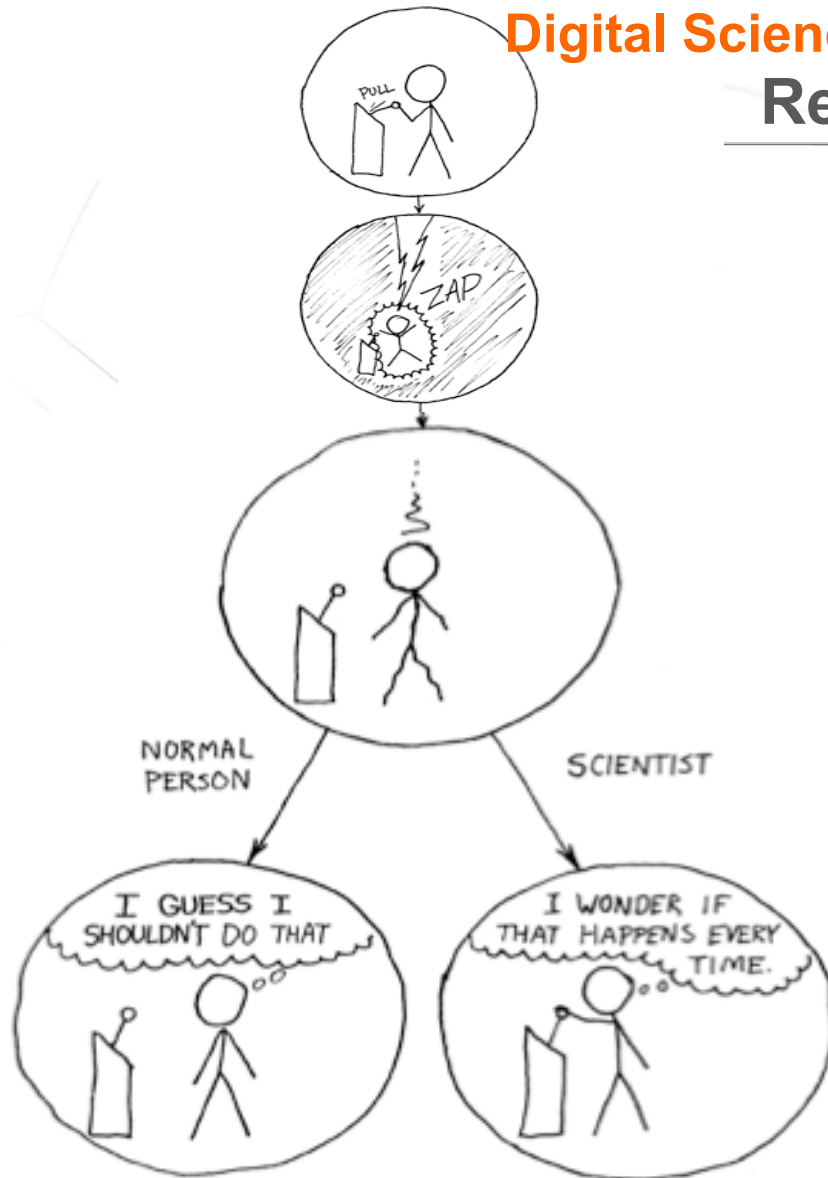
More data in archives does not imply more knowledge

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



Digital Science - Reproducibility and Visibility in Astronomy

Reproducibility and The Scientific Method



<http://xkcd.com/242/>

Benefits

- » Publishing knowledge, **not advertising**
- » The author, the referee and 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

Initiatives

- » **Elsevier Executable Papers Challenge**
- » Open Data / Open Science

Barriers to Data and Code Sharing in Computational Science

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

Code		Data
77%	Time to document and clean up	54%
52%	Dealing with questions from users	34%
44%	Not receiving attribution	42%
40%	I don't know how	-
34%	Legal barriers (i.e. 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%

Digital Science - Reproducibility and Visibility in Astronomy Publishing: Discovery, Visibility and Credit

nature International weekly journal of science


Search [Advanced search](#)

[Home](#) | [News & Comment](#) | [Research](#) | [Careers & Jobs](#) | [Current Issue](#) | [Archive](#) | [Audio & Video](#) | [For Authors](#)

[Archive](#) > [Specials & supplements archive](#) > [Science Metrics](#)

SPECIALS

[▶ See all specials](#)

[Journal home](#) | [Subscribe](#) | 

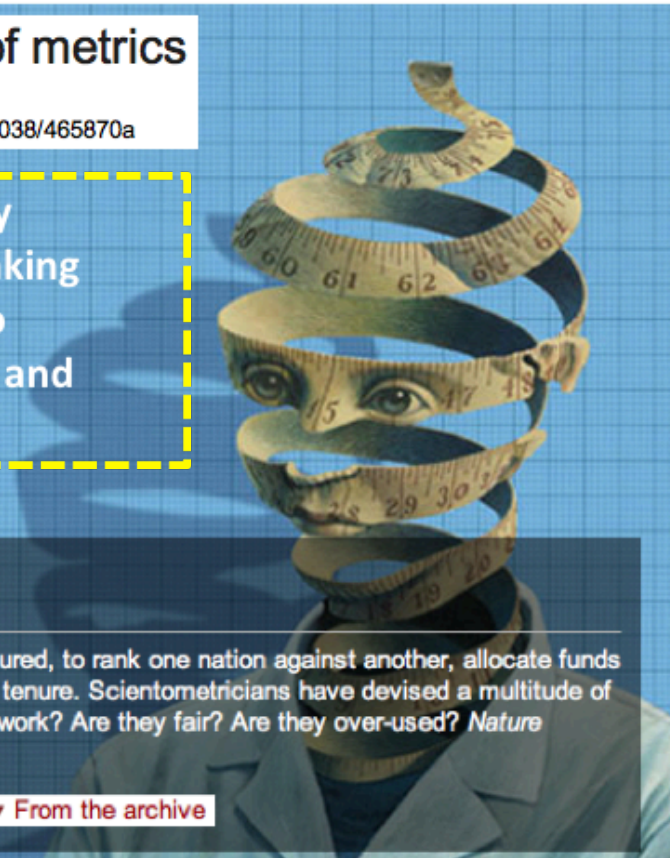
[Current issue](#) | [E-alert sign up](#)

[For authors](#) | [RSS feed](#)

How to improve the use of metrics

Nature 465, 870–872 (17 June 2010) | doi:10.1038/465870a

... “Science is being killed by numerical ranking,” [...] Ranking systems lures scientists into pursuing high rankings first and good science second.



SCIENCE METRICS

The value of scientific output is often measured, to rank one nation against another, allocate funds between universities, or even grant or deny tenure. Scientometricians have devised a multitude of 'metrics' to help in these rankings. Do they work? Are they fair? Are they over-used? *Nature* investigates.

[▼ Editorial](#) | [▼ Features](#) | [▼ Opinion](#) | [▼ From the archive](#)

Top content

Emailed	Downloaded	Blogged
1. Public health: The toxic truth about sugar <i>Nature</i> 01 February 2012		
2. The case for open computer programs <i>Nature</i> 22 February 2012		
3. The great beyond <i>Nature</i> 29 February 2012		
4. The darker side of stem cells <i>Nature</i> 29 February 2012		
5. Cancer: Solving an age-old problem <i>Nature</i> 29 February 2012		
View all ▶		

Digital Science - Reproducibility and Visibility in Astronomy Publishing: Discovery, Visibility and Credit

nature International weekly journal of science


Search [Advanced search](#)

[Home](#) | [News & Comment](#) | [Research](#) | [Careers & Jobs](#) | [Current Issue](#) | [Archive](#) | [Audio & Video](#) | [For Authors](#)

[Archive](#) > [Specials & supplements archive](#) > [Science Metrics](#)

SPECIALS

[▶ See all specials](#)

[Journal home](#) | [Subscribe](#) | 

[Current issue](#) | [E-alert sign up](#)

[For authors](#) | [RSS feed](#)

How to improve the use of metrics

Nature 465, 870–872 (17 June 2010) | doi:10.1038/465870a

Research reverts to a kind of 'academic prostitution', in which work is done to please editors and referees rather than to further knowledge.



SCIENCE METRICS

The value of scientific output is often measured, to rank one nation against another, allocate funds between universities, or even grant or deny tenure. Scientometricians have devised a multitude of 'metrics' to help in these rankings. Do they work? Are they fair? Are they over-used? *Nature* investigates.

- [▼ Editorial](#)
- [▼ Features](#)
- [▼ Opinion](#)
- [▼ From the archive](#)

Top content

Emailed	Downloaded	Blogged
1. Public health: The toxic truth about sugar <i>Nature</i> 01 February 2012		
2. The case for open computer programs <i>Nature</i> 22 February 2012		
3. The great beyond <i>Nature</i> 29 February 2012		
4. The darker side of stem cells <i>Nature</i> 29 February 2012		
5. Cancer: Solving an age-old problem <i>Nature</i> 29 February 2012		

[View all ▶](#)

Digital Science - Reproducibility and Visibility in Astronomy Publishing: Discovery, Visibility and Credit

nature International weekly journal of science


Search [Advanced search](#)

[Home](#) | [News & Comment](#) | [Research](#) | [Careers & Jobs](#) | [Current Issue](#) | [Archive](#) | [Audio & Video](#) | [For Authors](#)

[Archive](#) > [Specials & supplements archive](#) > [Science Metrics](#)

SPECIALS

[▶ See all specials](#)

[Journal home](#) | [Subscribe](#) | 

[Current issue](#) | [E-alert sign up](#)

[For authors](#) | [RSS feed](#)

How to improve the use of metrics

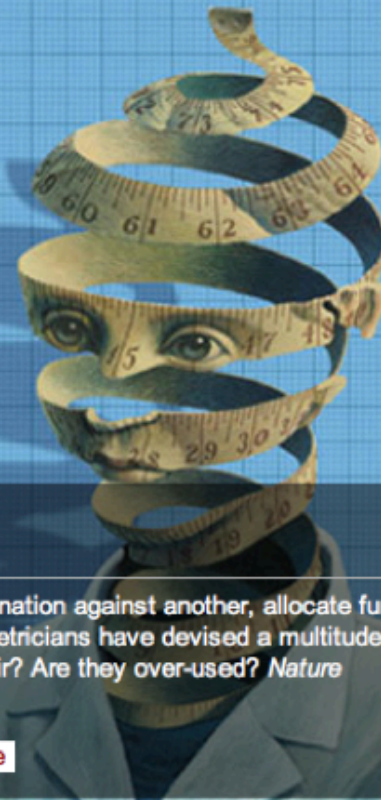
Nature 465, 870–872 (17 June 2010) | doi:10.1038/465870a

... an author's h-index can reflect longevity as much as quality — and can never go down with age, even if a researcher drops out of science altogether.

SCIENCE METRICS

The value of scientific output is often measured, to rank one nation against another, allocate funds between universities, or even grant or deny tenure. Scientometricians have devised a multitude of 'metrics' to help in these rankings. Do they work? Are they fair? Are they over-used? *Nature* investigates.

[▼ Editorial](#) | [▼ Features](#) | [▼ Opinion](#) | [▼ From the archive](#)

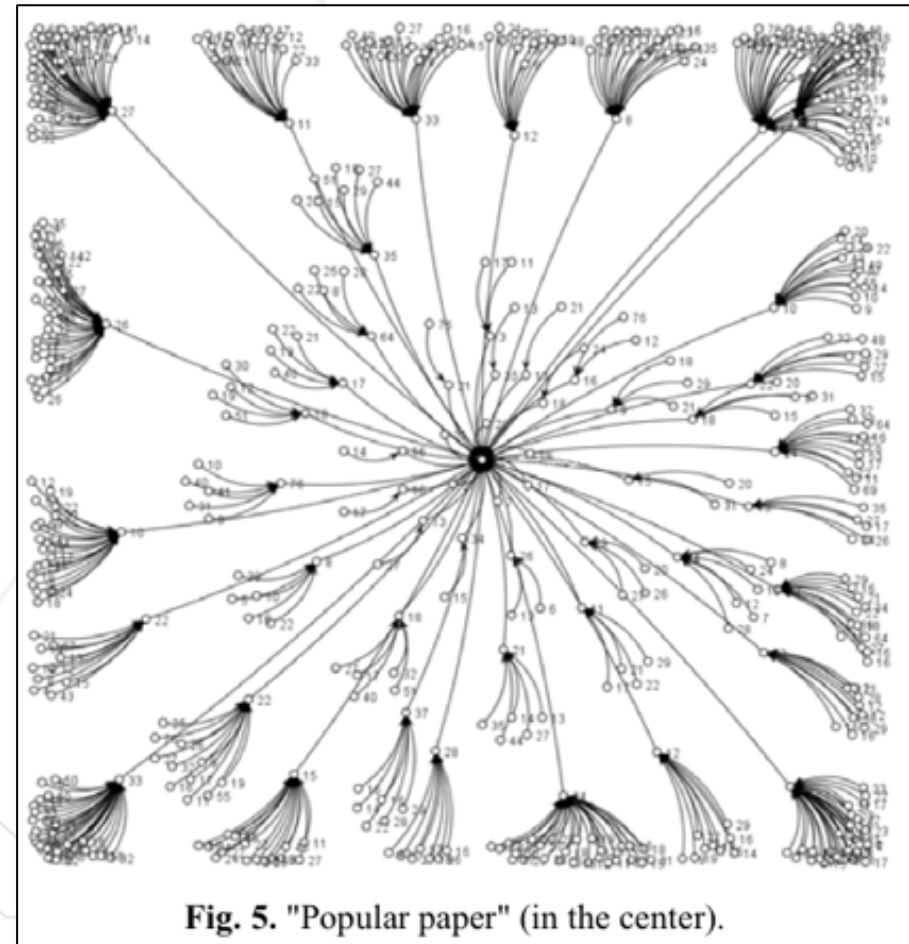
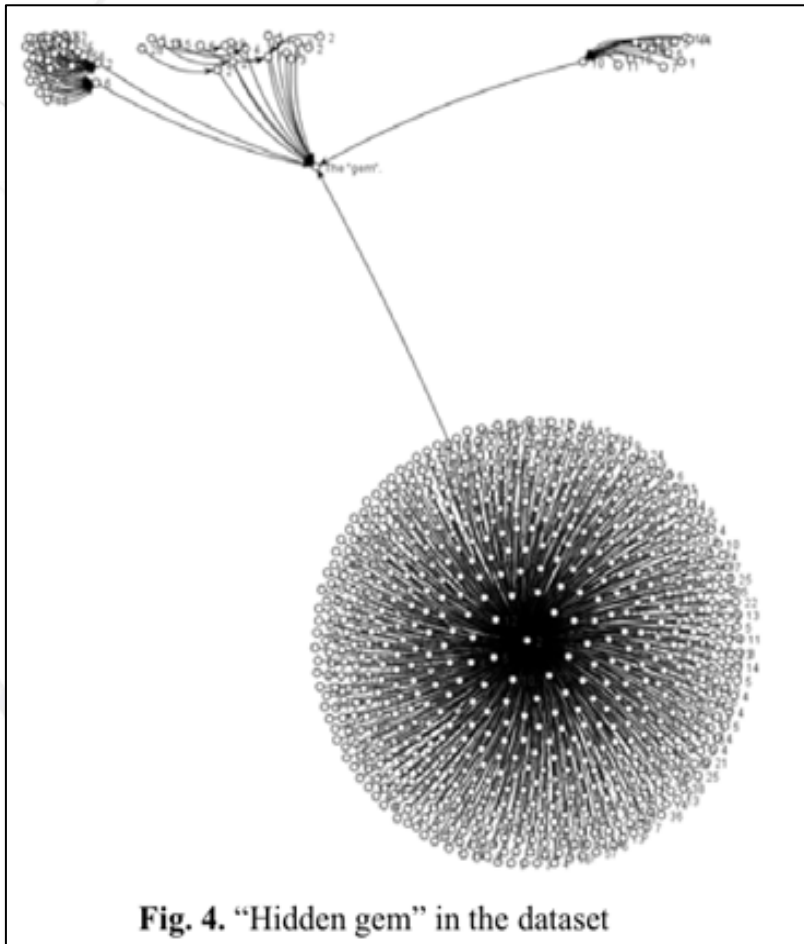


Top content

Emailed	Downloaded	Blogged
1. Public health: The toxic truth about sugar <i>Nature</i> 01 February 2012		
2. The case for open computer programs <i>Nature</i> 22 February 2012		
3. The great beyond <i>Nature</i> 29 February 2012		
4. The darker side of stem cells <i>Nature</i> 29 February 2012		
5. Cancer: Solving an age-old problem <i>Nature</i> 29 February 2012		
View all ▶		

Digital Science - Reproducibility and Visibility in Astronomy Publishing: Discovery, Visibility and Credit

Exploring and understanding scientific metrics in citation



Digital Science - Reproducibility and Visibility in Astronomy

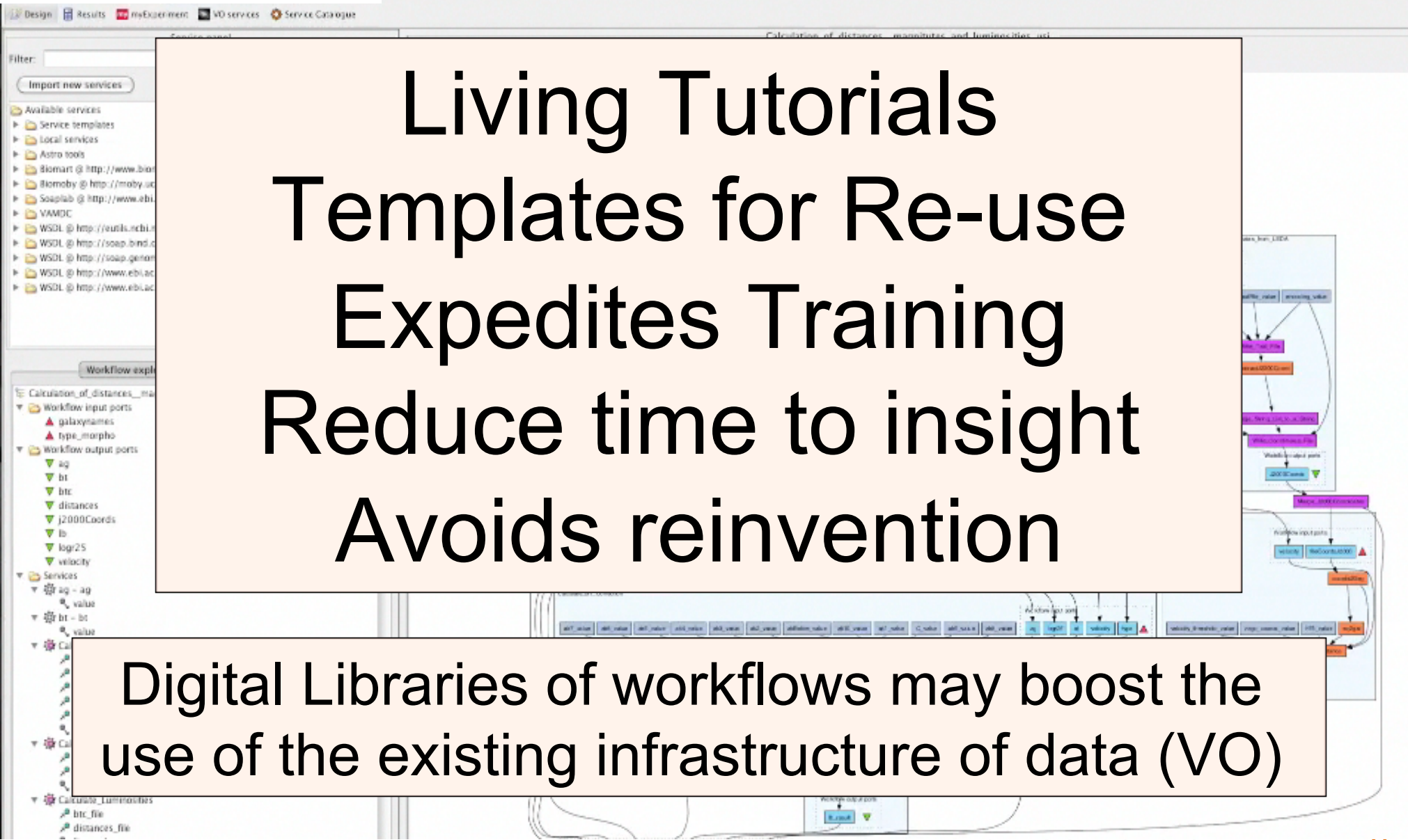
Publishing: Discovery, Visibility and Credit

Paper discovery: the social dimension

The image is a collage of various digital science and social networking tools. At the top left is the 'peerefevaluation' logo with the tagline 'empowering scholars'. Below it is the 'citeulike' logo. To the right is the 'MENDELEY' logo with a red molecular structure icon and a user profile snippet for 'Lourdes Verdes-Montenegro'. Further right is the 'YouTube' logo. Below these are 'BibSonomy' and 'ResearchGate' logos. In the center is the 'klænk' logo with the tagline 'Spread your research results'. To the right is the 'twitter' logo. Below 'klænk' is the 'delicious social bookmarking' logo. At the bottom left is the 'slideshare BETA' logo. In the center bottom is a screenshot of the 'AstroBetter' website, which features a search bar and navigation links like 'Blog', 'About', 'Archives', 'Support', and 'Wiki'. The main content of the screenshot is an article titled 'Learning Python - The Interactive Way' by Jessica, dated June 18, 2012. To the right of the AstroBetter screenshot is the text 'Collabgraph!' in large green letters, followed by a paragraph: 'Collaborating in your field of research. Just [green]ary or upload a bibtex file, containing your graph will create a fancy graph showing'. At the bottom right is the 'zotero' logo.

Digital Science - Reproducibility and Visibility in Astronomy

Going beyond automation: Scientific Workflows



The background image is a screenshot of a scientific workflow management system. On the left, there is a 'Service Catalogue' with a search filter and a list of 'Available services' including 'Service templates', 'Local services', and various 'WSO2' services. Below this is a 'Workflow explorer' showing a tree view of a workflow named 'Calculation_of_distances_ma'. The main area displays a complex workflow diagram with nodes and arrows, representing the process of calculating distances, magnitudes, and luminosities. The text is overlaid on this interface.

Living Tutorials
Templates for Re-use
Expedites Training
Reduce time to insight
Avoids reinvention

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

Digital Science - Reproducibility and Visibility in Astronomy

Going beyond automation: Scientific Workflows

A STORY TOLD IN FILE NAMES:

Location: C:\user\research\data

Filename	Date Modified	Size	Type
data_2010.05.28_test.dat	3:37 PM 5/28/2010	420 KB	DAT file
data_2010.05.28_re-test.dat	4:29 PM 5/28/2010	421 KB	DAT file
data_2010.05.28_WTF.dat	9:58 PM 5/28/2010	30 KB	DAT file
data_2010.05.29_aaarrrob.dat	12:37 AM 5/29/2010	30 KB	DAT file
data_2010.05.29_USETHRISONE.dat	5:08 AM 5/29/2010	2,894 KB	DAT file
analysis_graphs.xls	7:13 AM 5/29/2010	455 KB	XLS file
ThesisOutline!.doc	7:26 AM 5/29/2010	38 KB	DOC file
Notes_Meeting_with_ProfSmith.txt	11:38 AM 5/29/2010	1,673 KB	TXT file
JUNK...	2:45 PM 5/29/2010		Folder
data_2010.05.30_startingover.dat	8:37 AM 5/30/2010	420 KB	DAT file

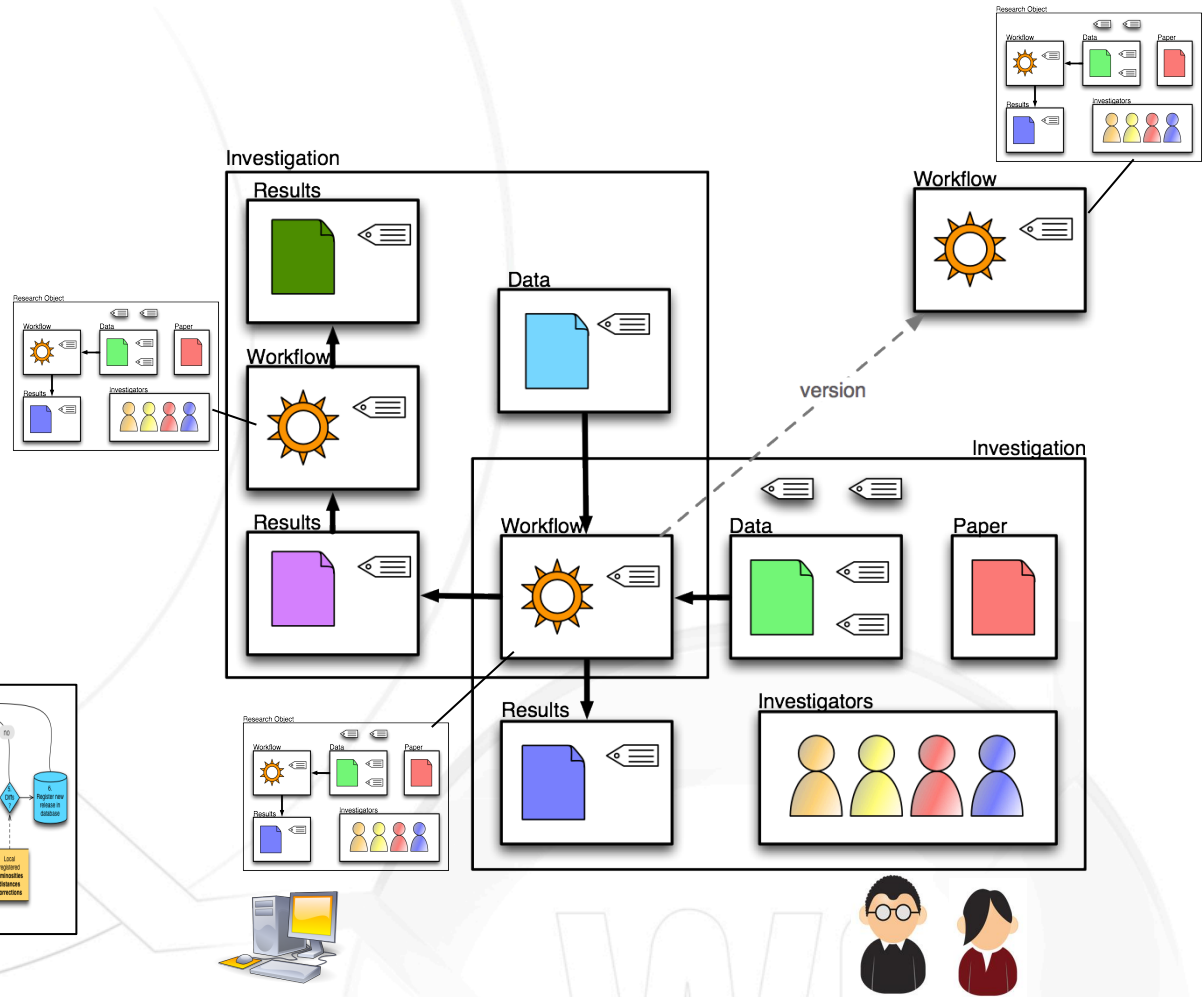
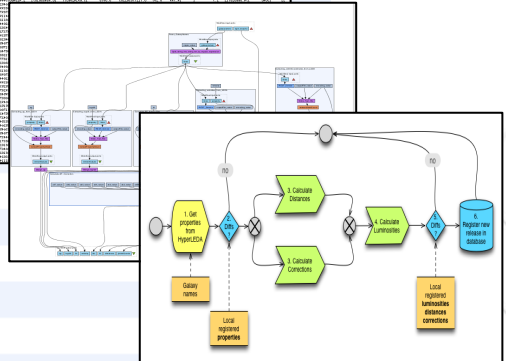
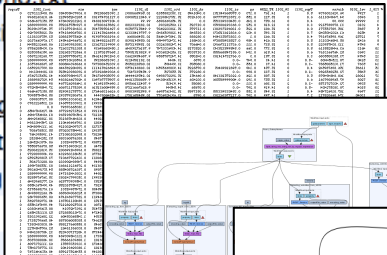
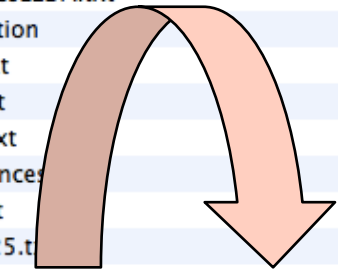
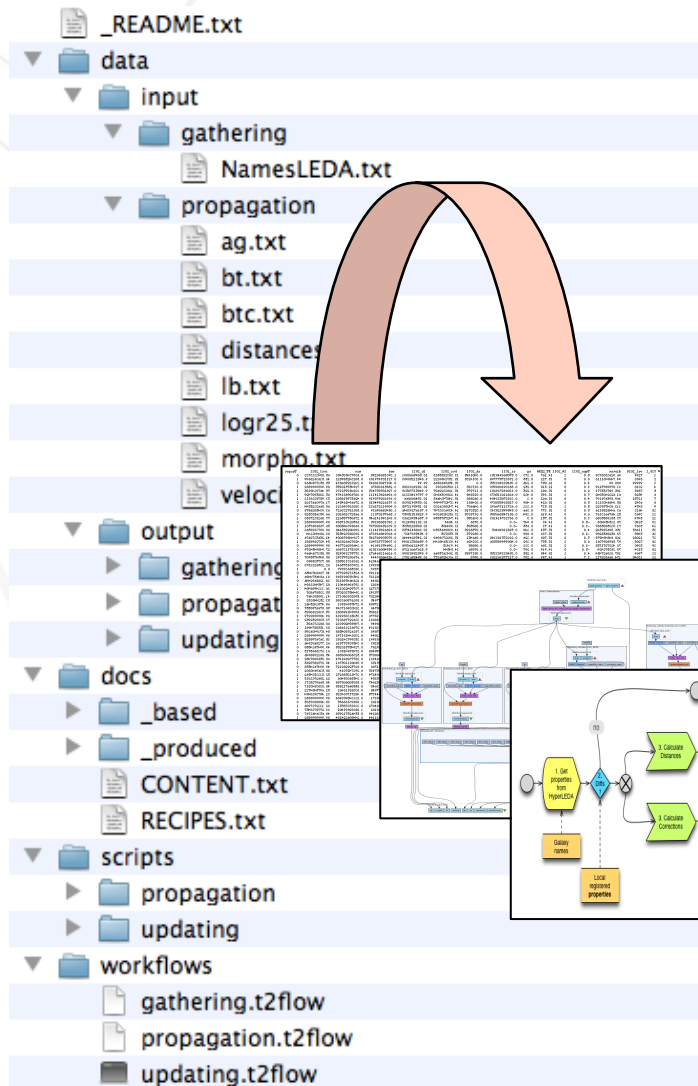
Automation does not imply organization

Assistive building
Completeness evaluation

Type: Ph.D Thesis Modified: too many times Copyright: Jorge Cham www.phdcomics.com

Digital Science - Reproducibility and Visibility in Astronomy Research Objects

Expose experiment in a structured way in order to be **understood**



Technical Objects
Distributed

Social Objects

Similar initiatives in Astronomy

- » **Semantic curation** of digital objects
 - › CDS Centre Données Strasbourg
 - › US Virtual Astronomical Observatory
 - › SAO/NASA ADSLabs
- » **Workflow users platforms**
 - › Cyber-SKA
 - › IceCore
 - › Montage
 - › Astro-WISE
 - › Helio-VO
- » **Semantically auto descriptive WS**
 - › Workflows VO-France

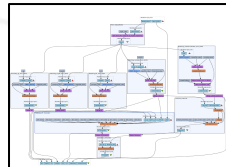
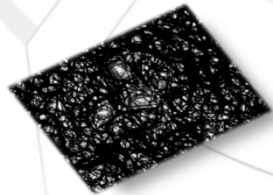


Digital Science - Reproducibility and Visibility in Astronomy Research Objects

ADSLabs Initiative

ADO Linked Components

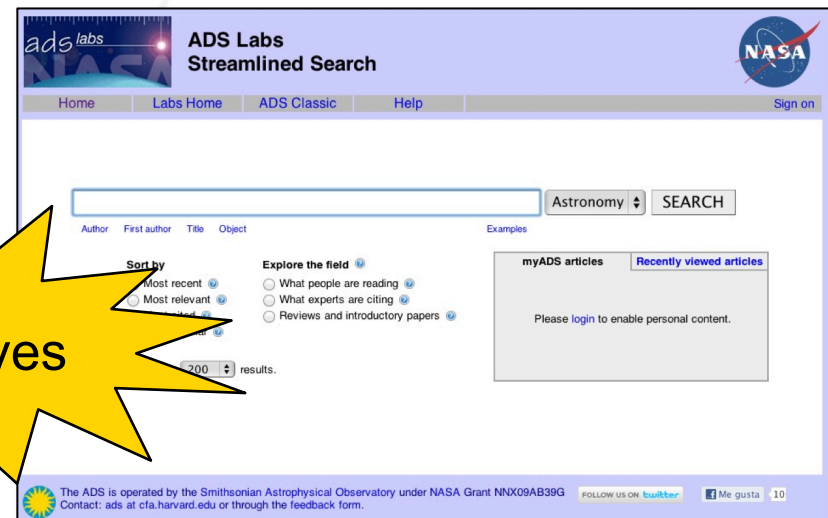
- » Authors
- » Publications
- » Journals
- » Objects SIMBAD
- » Tabular data behind the plots CDS
- » ASCL reference of used software
- » Observing time Proposals
- » Used facilities, surveys or missions



Incentives



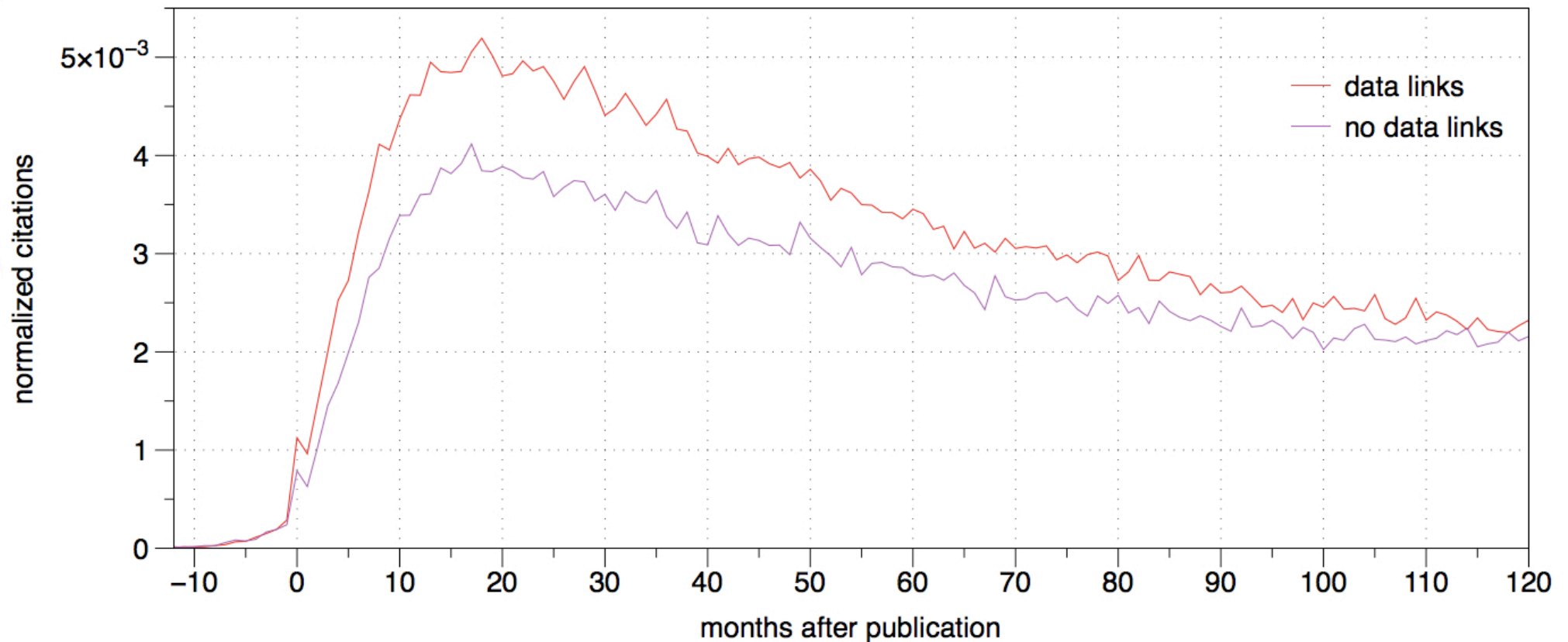
<http://labs.adsabs.harvard.edu/>



The Incentive

Papers with data links are cited more than those without

1995 - 2000

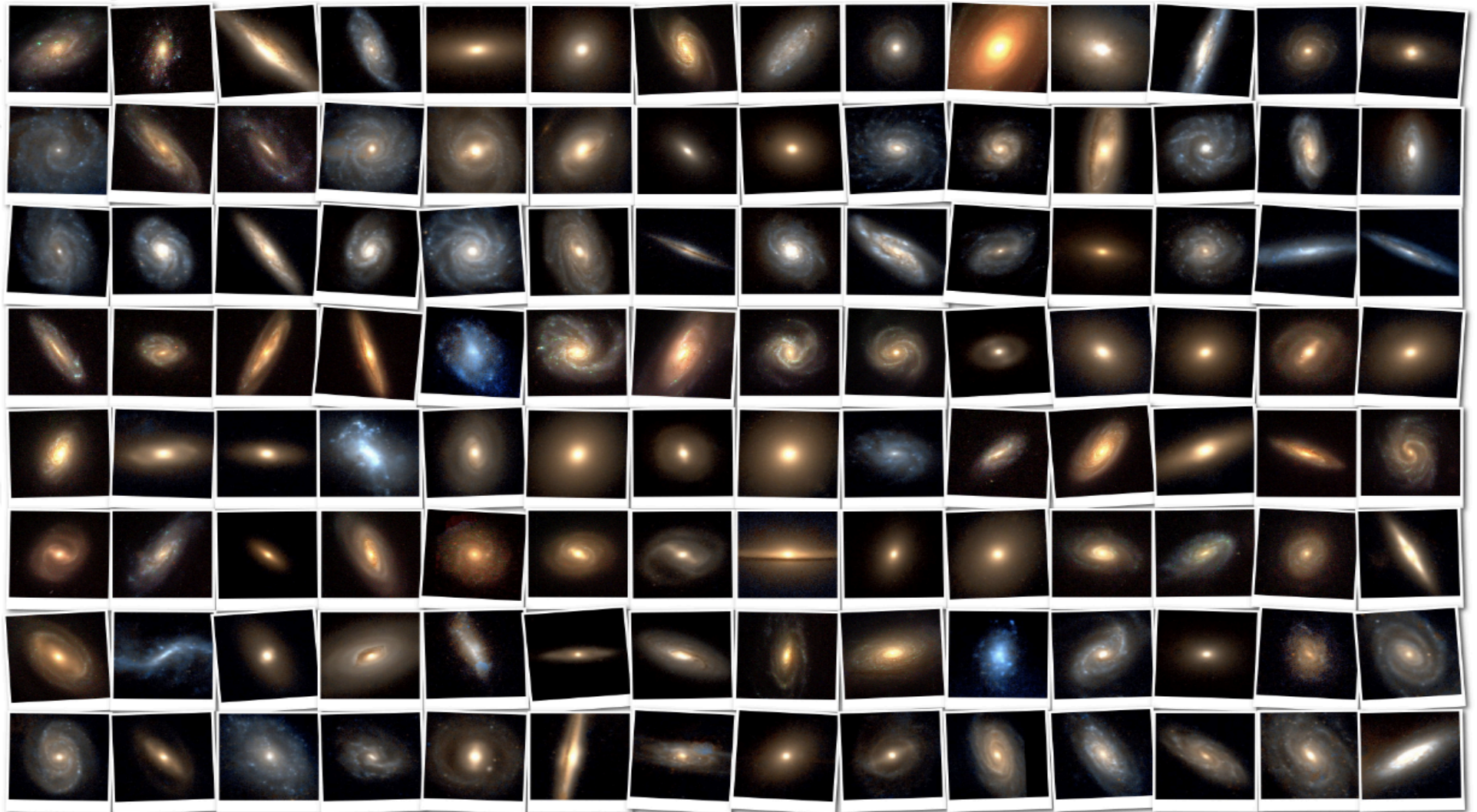


Effect of E-printing on Citation Rates in Astronomy and Physics
2006. Edwin A. Henneken et al.

Digital Science - Reproducibility and Visibility in Astronomy

Astronomical Research Objects in Action

Curation by inspecting propagation of changes in quantities

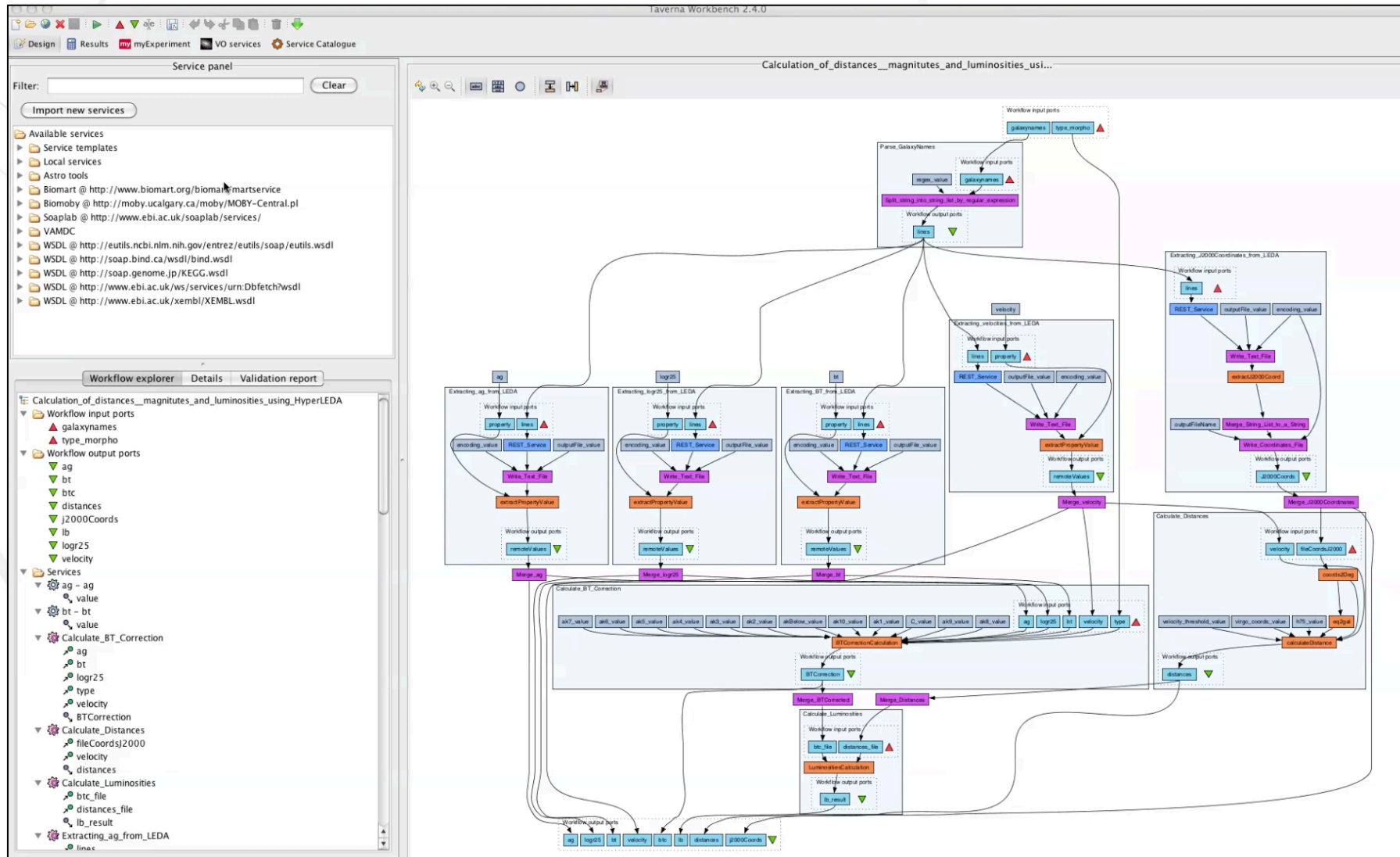


Credit: Zsolt Frei and James E. Gunn. The Galaxy Catalog

Digital Science - Reproducibility and Visibility in Astronomy

Astronomical Research Objects in Action

Create, annotate and run a workflow



Digital Science - Reproducibility and Visibility in Astronomy

Astronomical Research Objects in Action

Populate the Research Object and annotate

Wf4Ever - RO Annotator MOCKUP

Research Object: Distance Estimation

- Datasets
 - Galaxy_Names.csv
 - Apparent_Magnitudes.csv
- Scripts
- Web Services
- Workflows

Annotating "Galaxy_Names.csv"

Type: Comma-separated-value
Keywords: src; meta.name, galaxies, ...
Description: Names of galaxies whose
Role: Input file
Created At: 2011-09-06 16:32:18

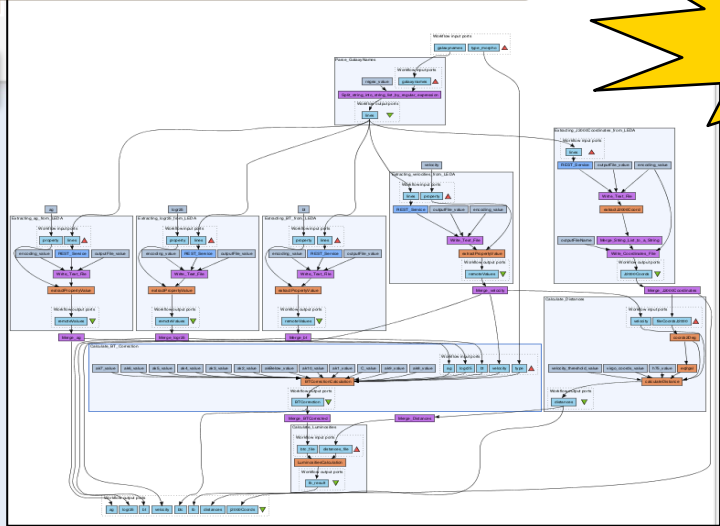


What kind of annotation is this?
Description:

Value for the annotation

Names of galaxies whose distance is to be estimated. Each line represents a different galaxy. Optional information on the galaxy is added as comma-separated values, in this

- Galaxy name
- Morphology type (NED)
- NED distance

Save Changes Cancel



Digital Science - Reproducibility and Visibility in Astronomy

Astronomical Research Objects in Action

Add documents and references

Wf4Ever - RO Annotator MOCKUP

Research Object: Distance Estimation

- Datasets
 - Galaxy_Names.csv
 - Apparent_Magnitudes.csv
- Scripts
- Web Services
- Workflows
- Docs
 - 2012A&A...536A.108V.AMIGA XIII: Workflow-based distance assessment...

Annotating "Galaxy_Names.csv"

Type: Comma-separated-value

Keywords: src; meta.name, galaxies, ...

Description: Names of galaxies whose

Role: Input file

Created At: 2011-09-06 16:32:18

Modified At: 2012-02-07 08:44:32

What kind of annotation is this?

Description:

Value for the annotation

Names of galaxies whose distance is to be estimated. Each line represents a different galaxy. Optional information on the galaxy is added as comma-separated values, in this order:

- Galaxy name
- Morphology type (NED)
- NED distance
- Estimation Method

Save Changes Cancel

User-agent: *
Allow: /

```
# October 11 2010 Whyte & Mackay whisky are running a promotion where 250 bottles of 30 year old whisky, each worth £150, are hidden in bottles of Whyte & Mackay Special whisky (learn more at http://bit.ly/whiskyhuntvideo )  
  
# The bottles are hidden in stores across the UK but we wanted to hide one in our new-look website as well - so if you're reading this congratulations on being a winner  
  
# Drop an email to richard at themasterblender dot com with the subject I Read Robots.txt Files and if you're one of the first to reply, are of a legal drinking age and your local licensing laws allow it we'll send you a bottle of 30 year old Whyte & Mackay. If we can't send you that, we'll send you something else.  
  
# And of course make sure you are following us on Twitter and Facebook - @the_nose @whyteandmackay facebook.com/whyteandmackay see what we do next. We do like our stunts!  
  
# If you weren't first, there's still 200 bottles hidden in stores across the UK  
October 11 2010
```

PDF

Digital Science - Reproducibility and Visibility in Astronomy

Astronomical Research Objects in Action

Add schema of the experiment

Wf4Ever - RO Annotator MOCKUP

Research Object: Distance Estimation

- Datasets
 - Galaxy_Names.csv
 - Apparent_Magnitudes.csv
- Scripts
- Web Services
- Workflows
- Docs

Annotating "Galaxy_Names.csv"

Type: Comma-separated-value

Keywords: src; meta.name, galaxies, ...

Description: Names of galaxies whose

Role: Input file

Created At: 2011-09-06 16:32:18

Modified At: 2012-02-07 08:44:32

What kind of annotation is this?

Description:

Value for the annotation

Names of galaxies whose distance is to be estimated. Each line represents a different galaxy. Optional information on the galaxy is added as comma-separated values, in this order:

- Galaxy name
- Morphology type (NED)
- NED distance
- Estimation Method

```
graph TD; Start(( )) --> T1{{1. Get properties from HyperLEDA}}; T1 --> D1{2. Diff?}; D1 -- no --> T1; D1 --> T2(( )); T2 --> T3[3. Calculate Distances]; T2 --> T4[3. Calculate Corrections]; T3 --> T5(( )); T4 --> T5; T5 --> T6[4. Calculate Luminosities]; T6 --> D2{5. Diff?}; D2 -- no --> T1; D2 --> T7[(6. Register new release in database)]; T7 --> End(( ))
```

The workflow diagram illustrates the process of distance estimation. It begins with a start node leading to a task '1. Get properties from HyperLEDA', which takes 'Galaxy names' as input. This leads to a decision diamond '2. Diff?'. If 'no', it loops back to '1. Get properties from HyperLEDA'. If 'yes', it proceeds to a merge node, then splits into two parallel tasks: '3. Calculate Distances' and '3. Calculate Corrections'. Both lead to another merge node, followed by '4. Calculate Luminosities'. This leads to a second decision diamond '5. Diff?'. If 'no', it loops back to '2. Diff?'. If 'yes', it leads to a task '6. Register new release in database', which takes 'Local registered luminosities distances corrections' as input. The process ends at a final node.

Digital Science - Reproducibility and Visibility in Astronomy

Astronomical Research Objects in Action

Publication for later discovery

Home / Research Object: <http://sandbox.wf4ever-project.org/rosrs5/ROs/HyperLEDA%20Luminosities/>

Interactive Conceptual Physical

HyperLEDA Luminosities/

- Web Services
- Datasets
 - agNew.txt
 - lbOld.txt
 - j2000Coords.txt
 - lbNew.txt
 - diff_lb.txt
 - lb.sql
 - NamesLEDA.txt
 - logr25New.txt
 - velocitiesNew.txt
 - distancesNew.txt
 - morphoNew.txt
 - btcNew.txt
 - btNew.txt
- Scripts
- Workflows
 - comparison_and_update_values_475535.
 - calculating_the_total_luminosity_of_a_galaxy_using_properties_from_text_1
 - gathering_galaxy_properties_using_hyperleda_129473.
- Workflow Runs
- Documents
 - GoldenTrace.txt

Item info


Created by: Jose Enrique Ruiz


Created on: 2012.01.08 17:09:14 CET

File size: --

Number of annotations: 1

Keywords [galaxies][catalogs]

Integrity  50%

Rating 



Downloads 36

Citations 1

Re-used 4

Comments 2

<< Previous version | Next version >>



Digital Science - Reproducibility and Visibility in Astronomy

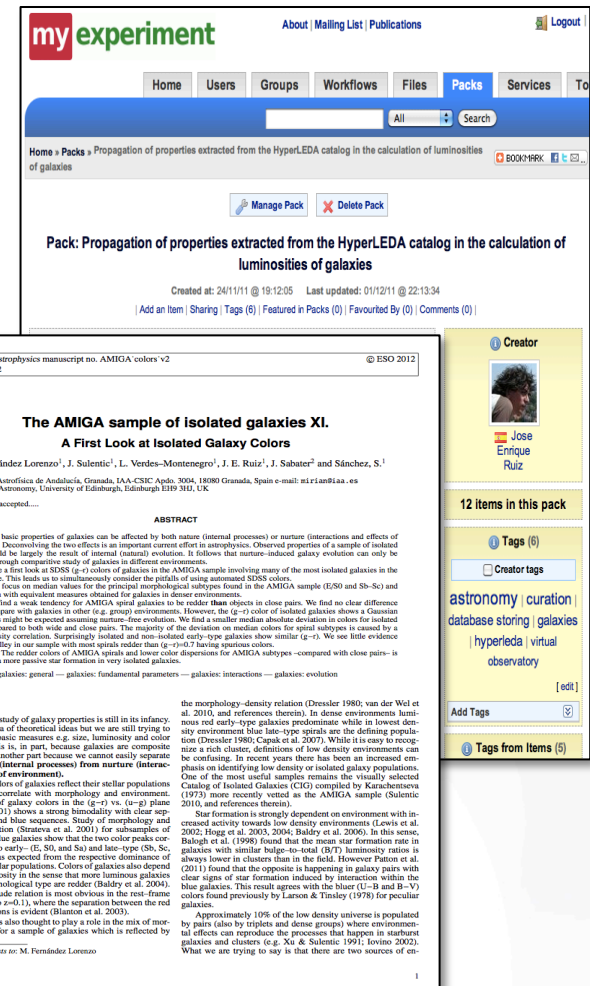
Astronomical Research Objects in Action

Curation by inspecting propagation of changes in quantities

- » Taverna 2.3 
- » MyExperiment Pack
 - » <http://www.myexperiment.org/packs/231>

Related Publication

The AMIGA sample of isolated galaxies XI. A First Look at Isolated Galaxy Colors 2012 A&A 540, A.47



The screenshot shows the MyExperiment website interface. At the top, there's a navigation bar with 'Home', 'Users', 'Groups', 'Workflows', 'Files', 'Packs', 'Services', and 'Logout'. Below this is a search bar and a breadcrumb trail: 'Home > Packs > Propagation of properties extracted from the HyperLEDA catalog in the calculation of luminosities of galaxies'. The main content area displays the pack title, creation and update dates, and a list of actions like 'Manage Pack' and 'Delete Pack'. On the right side, there's a 'Creator' section for Jose Enrique Ruiz, showing 12 items in the pack and a list of tags including 'astronomy', 'curation', 'database storing', 'galaxies', 'hyperleda', and 'virtual observatory'. The bottom part of the screenshot shows the abstract of the publication 'The AMIGA sample of isolated galaxies XI. A First Look at Isolated Galaxy Colors' by M. Fernández Lorenzo et al., published in A&A 540, A.47 (2012).

SOFTWARE

- <http://wf4ever.github.com/astrotaverna/>
- <http://www.taverna.org.uk/>



Hands on
AstroTaverna !

VIDEO TUTORIALS

- <http://amiga.iaa.es/files/tavernavideos/AstroTavernaIntro.m4v>
- <http://amiga.iaa.es/files/tavernavideos/NEDImages.mov>

How NOT to be a good e-astronomer

- » Write a **obscure paper**, do not say clearly how to reproduce the results
- » Do things **quickly** and forget about them once you've submitted the paper
- » Be untidy, **spread your code and data** in a variety of formats, folders and disks
- » Practise the "**data mine-ing**" – data are mine
- » Practise the "**data flirting**" – call me if you would like to see more
- » Do not provide data results, **including the plots is just fine**
- » Always **cite the same** authors and papers or those that cite you
- » Do not cite other resources than papers, **neither provide their URL links**
- » Do not search info on **Internet** with other tools than ADS or arXiv
- » **Work alone** and email/phone one friend if you have any doubt

 <http://amiga.iaa.es/p/212-workflows.htm>

 <http://www.wf4ever-project.org>

 jer@iaa.es