



Polish Roadmap toward
Domain-Specific Infrastructure
for Supporting Computational Science
in European Research Area

Reproducible e-science with GridSpace2 platform

Eryk Ciepiela, Bartosz Wilk, Daniel Hareźlak,
Marek Kasztelnik, Maciej Pawlik, Jan Meizner,
Marian Bubak Academic Computer Center CYFRONET
AGH University of Science and Technology

EGI Technical Forum, Madrid, 17 Sep 2013



GridSpace2
GridSpace2

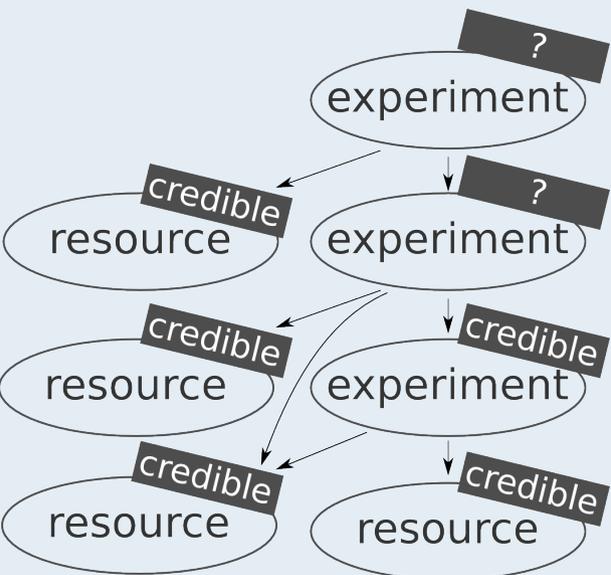
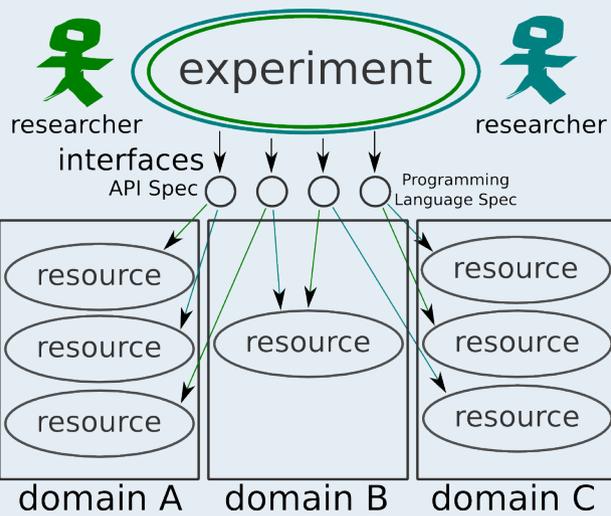
- Reproducibility problem
- Methodology proposed
- Technology proposed (GridSpace2)
- Case study (Collage)
- Challenges and future prospects
- Conclusions

Reproducibility Motivation



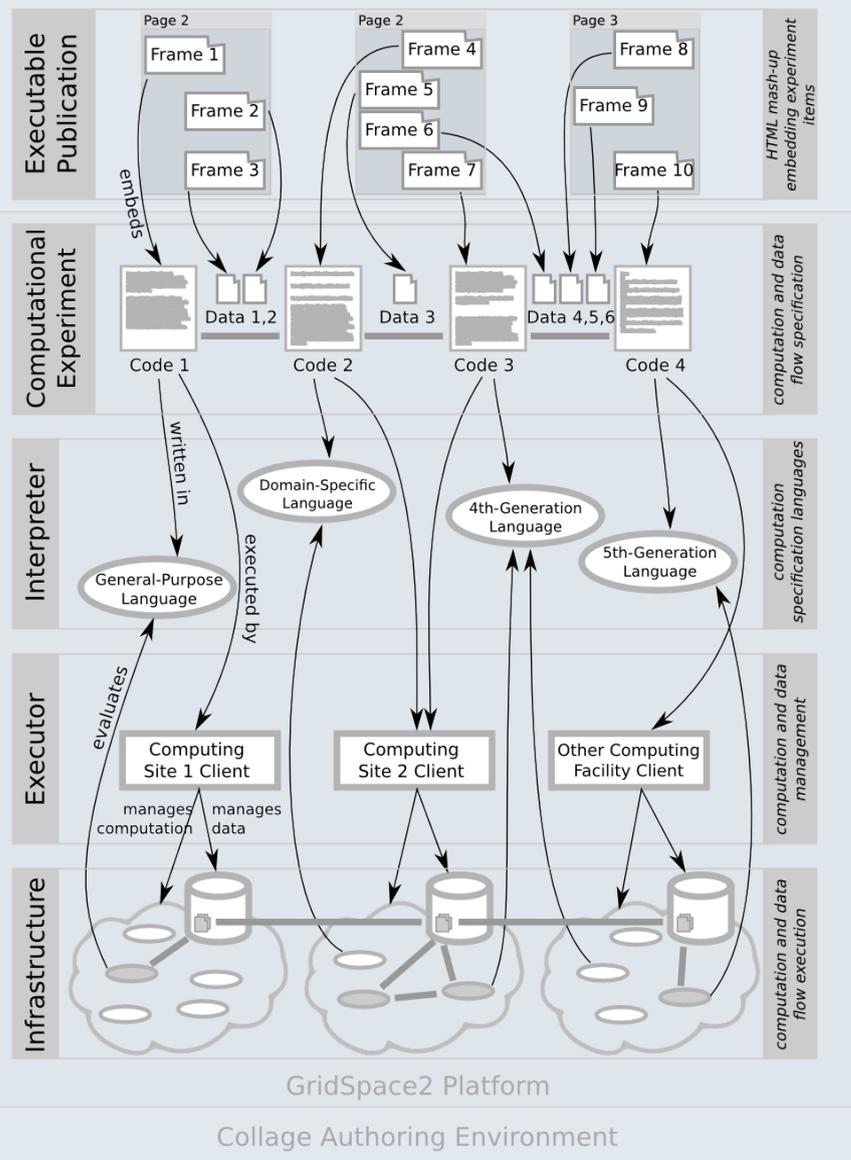
- **152 entries** for **retracted** scientific papers due to **reproducibility issues** since August 2010 reported by Retracting Watch (<https://retractionwatch.wordpress.com>)
- **10-100k** of articles published in 1950-2004 **ought to be retracted** (Cokol et al. *How many scientific papers should be retracted?*, EMBO Reports, 2007)
 - *Nature* – retracted 30 out of 45–67 articles that should have been retracted in 1999-2004
 - Lower IF journals – lower rate by the orders of magnitude
- Publishers share concern of **credibility crisis**
 - *Article of the Future* umbrella project (Elsevier)
 - Executable papers: publish reproducible experiments along with articles
 - *Executable Paper Grand Challenge* (2011, Elsevier, won by our solution: *Collage Authoring Environment*)
- Programmers critical about scientific computing
- Science as an open enterprise (Royal Society)
- R-dimension of e-research (David de Roure)

Reproducibility Problem Statement



- *Reproducibility – the ability of an entire experiment or study to be reproduced, either by the researcher or by someone else working independently. It is one of the main principles of the scientific method (...) [Wikipedia]*
- *Reproducibility – the ability to carry out the same experiment or study the other time and using the same or other resources within given time and cost constraints*
- *...other time...*
 - preservable
 - durable
 - susceptible to decay in ever evolving environment
- *...using the same other resources...*
 - coupled through interfaces
 - based on standards/specifications
 - accross administrative domains
- *...within given time and cost constraints*
 - the more straightforward the better
- *Credible – proved, considered or evaluated as correct or derived from the other credible*
 - reproducibility gives a way to evaluate correctness thus bring credibility
 - re-use of reproducible experiments and credible resources efficient also in terms of credibility evaluation: credibility chain
- **How to bring reproducibility to e-science experiments?**

GridSpace2 Model Methodology



- Experiments are workflows consisting of **code**, **data** and **requisite** items
- Experiments and their items are web-enabled: URL-accessible, embeddable on web sites e.g. on-line publication
- Code items are written in general-purpose, domain-specific, 4th or 5th generation programming languages that we call **interpreters**
- Data items are simply files or directories
- For executing code and storing data items the underlying e-infrastructures are used that we refer to as **executors**
- Experiments use e-infrastructures through abstract Executor API that constitutes an interoperability layer between various types and instances of e-infrastructures
- Experiments use interpreters without being bound to any specific interpreter implementation
- Concrete executors and interpreters' implementations are of user's choice
- GridSpace2 – web-oriented distributed computing platform
- Collage – GridSpace2 + executable publications support

Collage Authoring Workbench

You are logged in to **collage-exp host.elsevier.com** as **eciepiela**

collage-exp host.elsevier.com: collatz/collatz | collage-exp host.elsevier.com: hello

The Collatz Conjecture

Computing sequences for Ruby 1.8.7 with collage-exp host.elsevier.com

Generating plots GnuPlot 4.2.6 with collage-exp host.elsevier.com

```
set output 'collatz/hwm.png'
plot 'collatz/results_raw.txt' using 1:3 title "The biggest number (high water mark) reached when iterating"

set output 'collatz/last.png'
plot 'collatz/results_raw.txt' using 1:4 title "Value in last iteration"
```

1 input data defined for this code. 3 output data defined for this code.

Output

```
gnuplot> set output 'collatz/hwm.png'
gnuplot> plot 'collatz/results_raw.txt' using 1:3 title "The biggest number (high h water mark) reached when iterating"
```

Files: arguments.txt, collatz.exp.xml, hwm.png, iters.png, last.png, results_raw.txt

Filter by file name

Path: eciepiela/collatz/

Releases

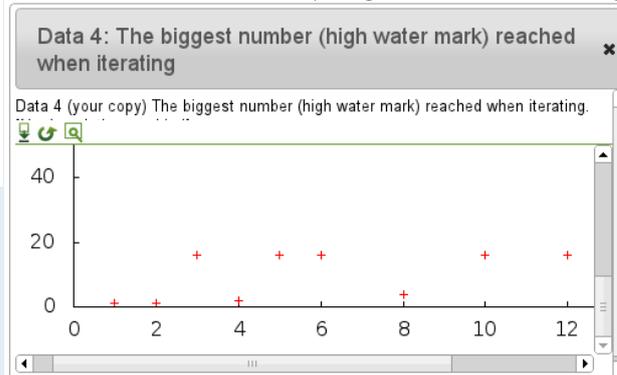
© 2012 Spectrum - IT Division Please report any problems related to this portal to our Issue Tracker Acknowledgements

The Collatz Conjecture

eciepiela, DOI: 10.0000/1358511059290

The experiment was released by **eciepiela** on **Fri Jan 18 13:10:59 CET 2013** in the **private** scope. No, below is not an article. It's only generated text with injected labels that navigate to particular experiment items to show you how in-text links work.

Euismod in pellentesque massa placerat dui ultricies lacus. Penatibus et magnis dis parturient. Ipsum consequat nisl vel pretium lectus. Interdum velit euismod in pellentesque. Nascetur ridiculus mus mauris vitae ultricies leo integer malesuada nunc. Sapien nec sagittis aliquam malesuada. Donec ultrices tincidunt arcu non sodales neque sodales. Sed velit dignissim sodales ut eu. [Data 1](#) Ac tortor vitae purus faucibus ornare. Mattis enim ut tellus elementum sagittis vitae et. Vitae semper quis lectus nulla at. Vitae purus faucibus ornare suspendisse sed nisi lacus sed viverra. Amet tellus cras adipiscing enim eu. [Code 1](#) Tellus integer



Data 1: Arguments

Code 1: Computing sequences for given arguments

Code 1 (your copy) Computing sequences for

Source: Ruby 1.8.7 Output

```
step+=1
end
results_raw.puts("#(argument)\t\t#"
```

Data 2: Raw results to be visualized afterwards

Data 2 (original) Raw results to be visualized

1	0	1	1
2	1	1	1
3	7	16	1
4	2	2	1
5	5	16	1
6	8	16	1
7	10	52	10

Code 2: Generating plots

Data 3: Number of iterations

Data 4: The biggest number (high water mark) reached when iterating

Data 5: Value in last iteration

Save experiment Reset

Experiment authoring and running through a dedicated web application: **Experiment Workbench** (up)

Experiment publishing on arbitrary web site, e.g. Elsevier ScienceDirect portal (right)

Collage Executable Papers

Case Study



7

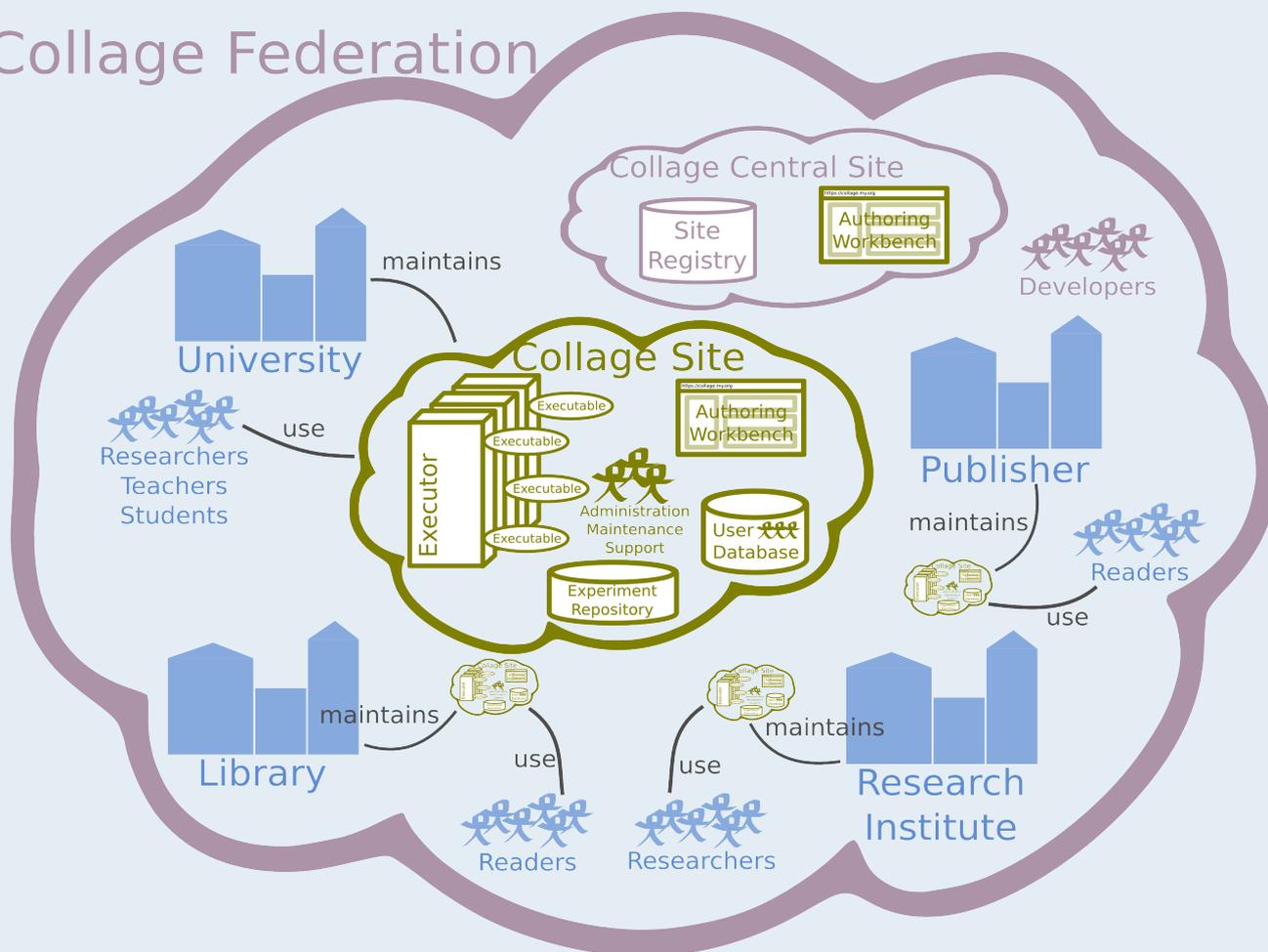


- GridSpace2-based Collage Authoring Environment won Executable Paper Grand Challenge in 2011
- Authoring, reviewing, publishing and executing of experiments offered by Collage proved suitable for scientific publishing in its idea
- Collage was integrated with Elsevier ScienceDirect portal so papers can be linked and presented with corresponding computational experiments
<https://collage.elsevier.com>
- Special Issue of *Computers & Graphics* journal featuring 7 executable papers was released in May 2013
<http://www.journals.elsevier.com/computers-and-graphics/news/special-issue-on-3d-object-retrieval/>
(video)

Federated Collage

Challenges and Future Prospects

Collage Federation



- Collage Site
 - University, Library, Research Institute, Computer Center etc.
 - User database
 - Experiment repository
 - Executors
 - Maintenance, support
- Central Collage Site
 - Registry of interpreters and federated Collage Sites
 - Experiments portable between Collage Sites
- “Publish globally – compute locally” model
- Collage/GridSpace2 software as public open-source project
 - Further development voluntary
- Collage/GridSpace2 software shipped to Collage Sites as virtual machine appliance

- Reproducibility is an inherent part of scientific method
- E-science suffers from lack of reproducibility what causes credibility crisis in scholarly papers
- Poor quality of scientific programs leads to retracted papers so they also must be subject to review
- Credibility of e-science experiments must be efficiently evaluated
- New methodology is needed to embrace above mentioned concerns – GridSpace2 model is the example
- New tools are needed to aid such methodology – like GridSpace2 platform
- GridSpace2/Collage platform is scalable and their instances can be federated
- With federated instances of GridSpace2/Collage platform we can achieve reproducibility across administrative domains and scientific communities

GridSpace2
GridSpace2