



www.software.ac.uk

Referencing software to make it sustainable: what, why and how?

SciencePAD Persistent Identifiers workshop

30 January 2013, CERN

Neil Chue Hong (@npch)

N.ChueHong@software.ac.uk

Unless otherwise indicated
slides licensed under



Why is this important?
An animated explanation
Featuring a frustrated panda



<http://tinyurl.com/datasharingpanda>



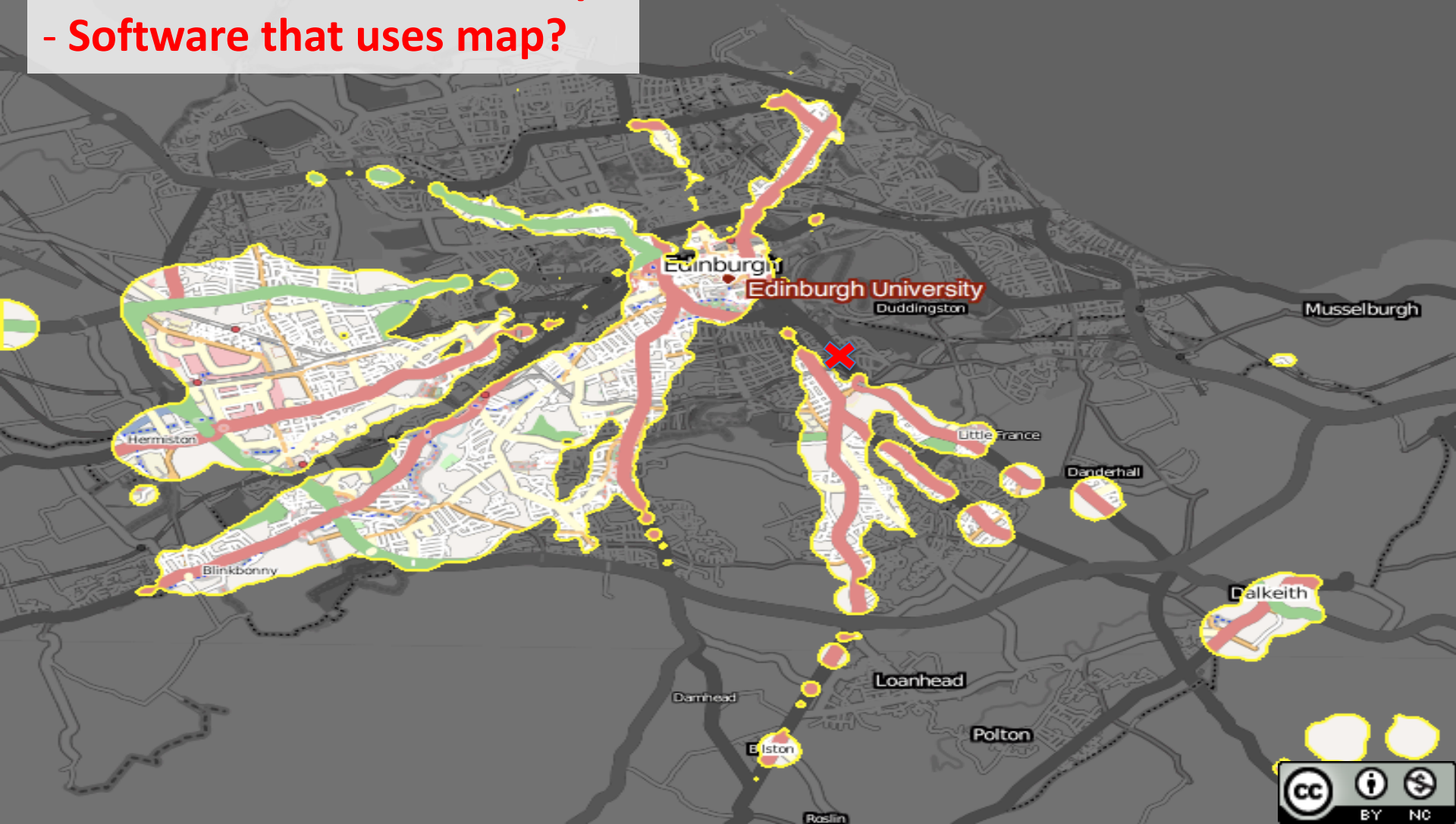
Software is no
longer easy to
define, let alone
sustain

The light, coloured areas of the map represent places where it's faster to use public transport than to drive if you want to get to work in central Edinburgh by 9AM (centred on postcode: EH1 2QL)

What do we sustain:

- Map?
- Software that creates map?
- Software that uses map?

Novel reuse of public sector data
<http://www.mysociety.org>



Boundary

Home Users Groups

AIDA | alignment | **benchmarks** | bio2rdf | BioAID | bioassist_nl | bioinform

| data integration | dbfetch | demo | design pattern | disease | e-science | ebi |

localworker | microarray | multiple sequence alignment | mygrid | ondex

annotation | pubmed | semantic_web | sequence | sequence alignment | sequence s

mining | text_mining | tutorial | uniprot | VL

Latest Last Updated Most Viewed Most Downloaded Most Fav

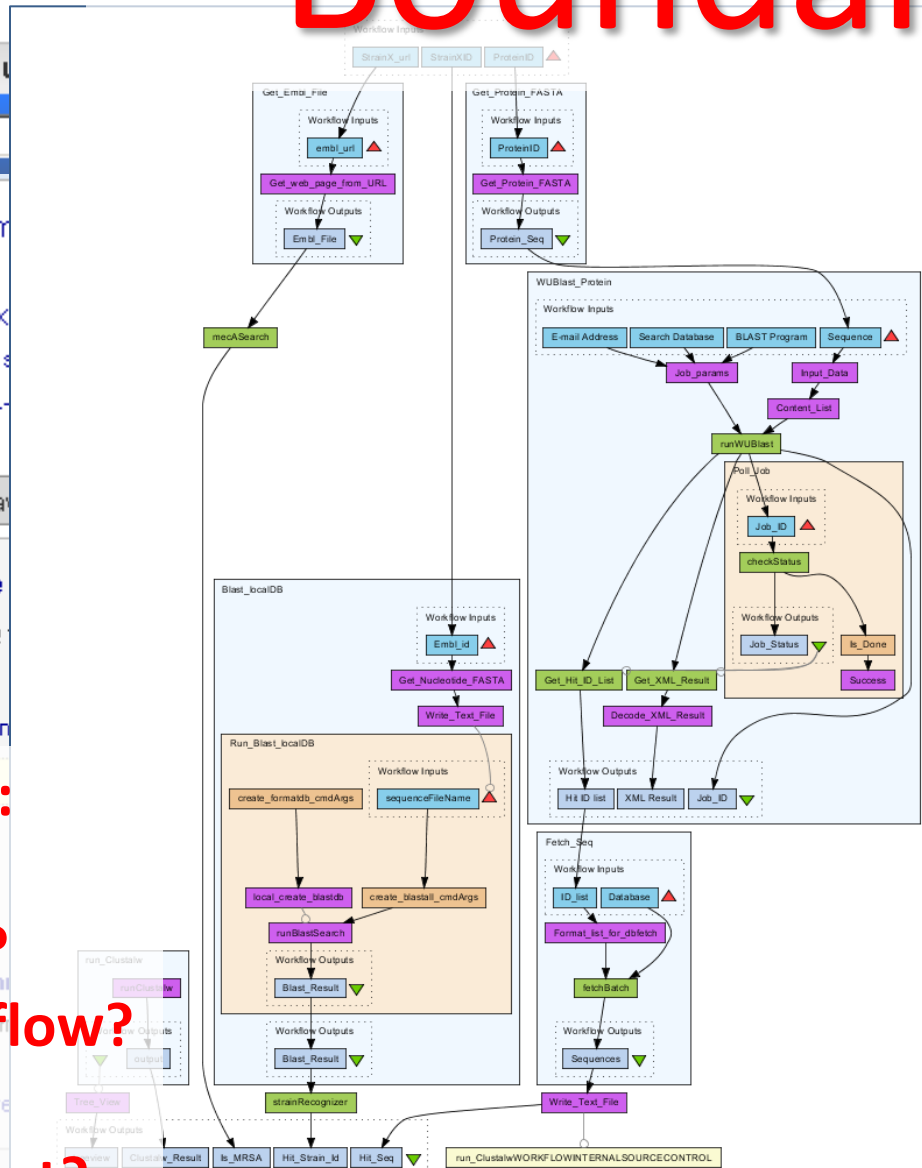
Taverna 1 **Identify_strain_and_phylogenetic_tree**

Created: 20/03/09 @ 15:20:36 | Last updated: 20/03/09 @

Credits: Aailyso Hamish McWilliam AID

License: Creative Commons Attribution-Share Alike 3.0 Licen

Original Uploader



- What do we choose to identify:**
- Workflow?
 - Software that runs workflow?
 - Software referenced by workflow?
 - Software dependencies?
- What's the minimum citable part?**

Granularity

Library / Suite / Package

Program

```
01 #include <stdio.h>
02 #include <stdlib.h>
03
04 void swap(int *x,int *y)
05 {
06     int temp;
07     temp = *x;
08     *x = *y;
09     *y = temp;
10 }
11
12 int choose_pivot(int i,int j )
13 {
14     return((i+j) /2);
15 }
16
17 void quicksort(int list[],int m,int n)
18 {
19     int key,i,j,k;
20     if( m < n)
21     {
22         k = choose_pivot(m,n);
23         swap(&list[m],&list[k]);
24         key = list[m];
25         i = m+1;
26         j = n;
27         while(i <= j)
28         {
29             while((i <= n) && (list[i] <= key))
30                 i++;
31             while((j >= m) && (list[j] > key))
32                 j--;
33             if( i < j)
34                 swap(&list[i],&list[j]);
35         }
36         // swap two elements
37         swap(&list[m],&list[j]);
38         // recursively sort the lesser list
39         quicksort(list,m,j-1);
40         quicksort(list,j+1,n);
41     }
42 }
43 void printlist(int list[],int n)
44 {
45     int i;
46     for(i=0;i<n;i++)
47         printf("%d\t",list[i]);
48 }
49
50 void main()
51 {
52     const int MAX_ELEMENTS = 10;
53     int list[MAX_ELEMENTS];
54
55     int i = 0;
56
57     // generate random numbers and fill them to the list
58     for(i = 0; i < MAX_ELEMENTS; i++) {
59         list[i] = rand();
60     }
61     printf("The list before sorting is:\n");
62     printlist(list,MAX_ELEMENTS);
63
64     // sort the list using quicksort
65     quicksort(list,0,MAX_ELEMENTS-1);
66
67     // print the result
68     printf("The list after sorting using quicksort algorithm:\n");
69     printlist(list,MAX_ELEMENTS);
70 }
```

Function

Algorithm

```
01 #include <stdio.h>
02 #include <stdlib.h>
03
04 void swap(int *x,int *y)
05 {
06     int temp;
07     temp = *x;
08     *x = *y;
09     *y = temp;
10 }
11
12 int choose_pivot(int i,int j )
13 {
14     return((i+j) /2);
15 }
16
17 void quicksort(int list[],int m,int n)
18 {
19     int key,i,j,k;
20     if( m < n)
21     {
22         k = choose_pivot(m,n);
23         swap(&list[m],&list[k]);
24         key = list[m];
25         i = m+1;
26         j = n;
27         while(i <= j)
28         {
29             while((i <= n) && (list[i] <= key))
30                 i++;
31             while((j >= m) && (list[j] > key))
32                 j--;
33             if( i < j)
34                 swap(&list[i],&list[j]);
35         }
36         // swap two elements
37         swap(&list[m],&list[j]);
38         // recursively sort the lesser list
39         quicksort(list,m,j-1);
40         quicksort(list,j+1,n);
41     }
42 }
43 void printlist(int list[],int n)
44 {
45     int i;
46     for(i=0;i<n;i++)
47         printf("%d\t",list[i]);
48 }
49
50 void main()
51 {
52     const int MAX_ELEMENTS = 10;
53     int list[MAX_ELEMENTS];
54
55     int i = 0;
56
57     // generate random numbers and fill them to the list
58     for(i = 0; i < MAX_ELEMENTS; i++) {
59         list[i] = rand();
60     }
61     printf("The list before sorting is:\n");
62     printlist(list,MAX_ELEMENTS);
63
64     // sort the list using quicksort
65     quicksort(list,0,MAX_ELEMENTS-1);
66
67     // print the result
68     printf("The list after sorting using quicksort algorithm:\n");
69     printlist(list,MAX_ELEMENTS);
70 }
```

...

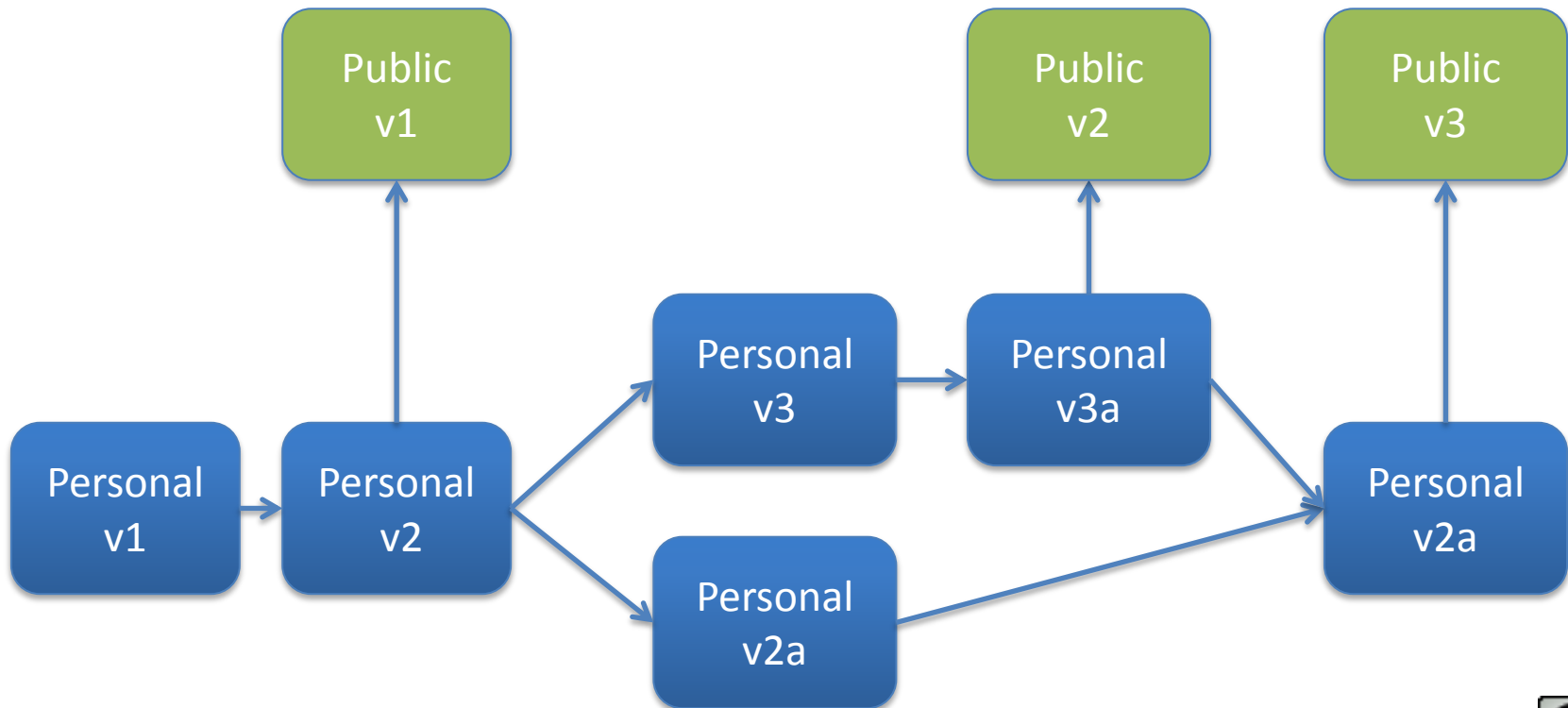
```
01 #include <stdio.h>
02 #include <stdlib.h>
03
04 void swap(int *x,int *y)
05 {
06     int temp;
07     temp = *x;
08     *x = *y;
09     *y = temp;
10 }
11
12 void bubblesort(int list[], int n)
13 {
14     int i,j;
15     for(i=0;i<n-1;i++)
16         for(j=i+1;j<n;j++)
17             if(list[i]>list[j])
18                 swap(&list[i],&list[j]);
19 }
20
21 void printlist(int list[],int n)
22 {
23     int i;
24     for(i=0;i<n;i++)
25         printf("%d\t",list[i]);
26 }
27
28 void main()
29 {
30     const int MAX_ELEMENTS = 10;
31     int list[MAX_ELEMENTS];
32
33     int i = 0;
34
35     // generate random numbers and fill them to the list
36     for(i = 0; i < MAX_ELEMENTS; i++) {
37         list[i] = rand();
38     }
39     printf("The list before sorting is:\n");
40     printlist(list,MAX_ELEMENTS);
41
42     // sort the list
43     bubblesort(list,MAX_ELEMENTS);
44
45     // print the result
46     printf("The list after sorting using bubble sorting algorithm:\n");
47     printlist(list,MAX_ELEMENTS);
48 }
```



Versioning

Why do we version?

- To indicate a change
- To allow sharing
- To confer special status



Authorship

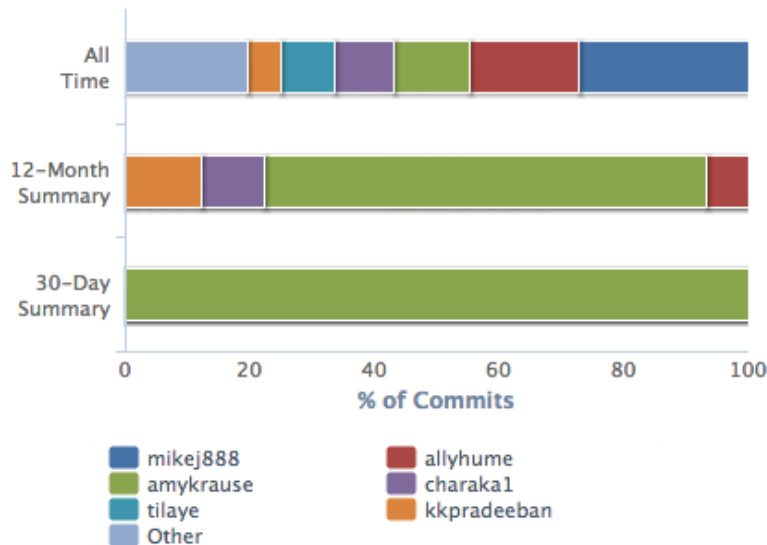
- Which authors have had what impact on each version of the software?
- Who had the largest contribution to the scientific results in a paper?

<http://beyond-impact.org/?p=175>

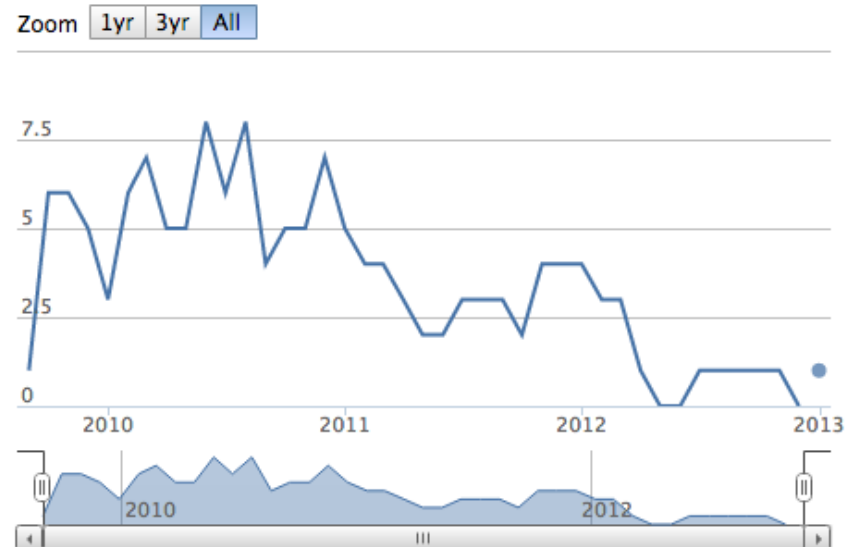
Contributors

Commits by Top Contributors

amykrause generated more than 50% of all commits during the past 12 months.



Number of Contributors



OGSA-DAI projects statistics
from Ohloh

Software

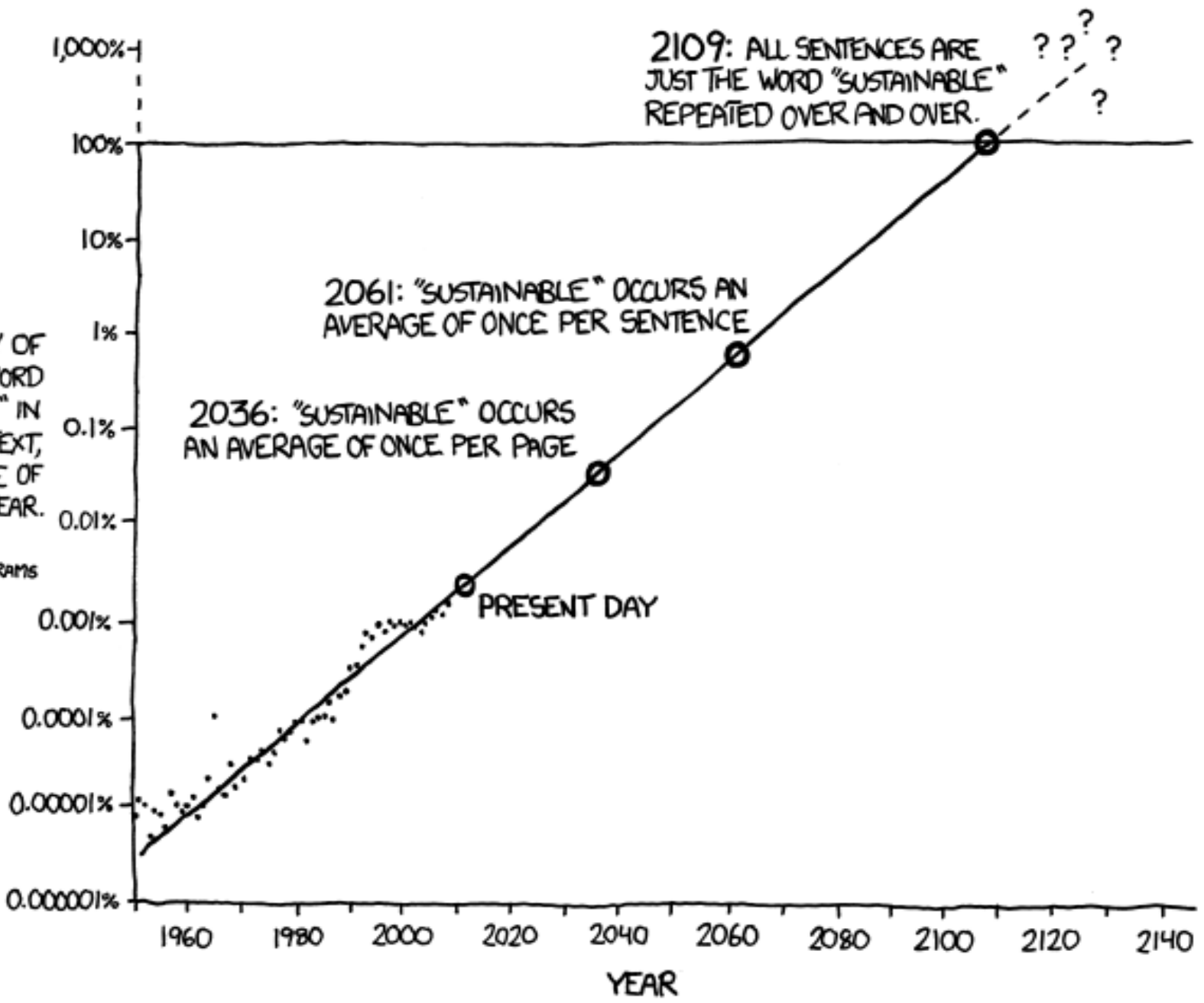
sustainability is the
ability to continue
to use, support,
maintain, and
evolve software



www.software.ac.uk

FREQUENCY OF USE OF THE WORD "SUSTAINABLE" IN US ENGLISH TEXT, AS A PERCENTAGE OF ALL WORDS, BY YEAR.

SOURCE: GOOGLE NGRAMS



THE WORD "SUSTAINABLE" IS UNSUSTAINABLE.

From xkcd.com

RCUK Gateway to Research



www.software.ac.uk



Gateway to research BETA

software

Did you mean? [asoftware](#) [nsoftware](#) [softwares](#) [oftware](#) [softare](#)

Research Funded (1486) Research Publications (75) People (2461) Organisations (21)

Relevance Start Date Grant Value 3 4 5 6 7 8 9 10 11 12 >>> 25 50 100

Project Status
Active (874) REMOVE
Funded Amount
Above 10M (0)
1M to 10M (93)
£100K to £1M (582)
Up to £100K (199)
Grant Category
Research Grant (794)
Fellowship (53)
Training Grant (13)
Intramural (8)
Other Grant (5)
P&Cs (1)
Funder
EPSRC (443)
BBSRC (151)
STFC (87)
NERC (68)
AHRC (60)
MRC (45)
ESRC (20)
Start Year
2010 (186)
2009 (158)
2011 (155)

£9,112,338 Feb 07 - Oct 12	Centre for Plant Integrative Biology BBSRC award to University of Nottingham and Charlie Hodgman
£9,106,003 Jan 07 - Dec 12	Centre for Systems Biology at Edinburgh BBSRC award to University of Edinburgh and Andrew Millar
£8,834,903 May 12 - May 18	DAASE: Dynamic Adaptive Automated Software Engineering EPSRC award to University College London and Mark Harman
£8,778,373 Jun 06 - Dec 10	OMII-UK Centre and Managed Programme EPSRC award to University of Southampton and David De Roure
£8,236,104 Apr 11 - Apr 16	Natural Speech Technology EPSRC award to University of Edinburgh and Steve Renals
£5,767,968 Jul 11 - Jul 16	EPSRC Centre for Innovative Manufacturing in Through-life Engineering Services EPSRC award to Granfield University and Rajkumar Roy
£5,492,969 Jan 11 - Dec 15	HUMAN-AGENT COLLECTIVES: FROM FOUNDATIONS TO APPLICATIONS [ORCHID] EPSRC award to University of Southampton and Nick Jennings

£5,055,294 Jun 07 - Jun 13	LSCITS-RPv2: Large-Scale Complex IT Systems Initiative - Research Programme v2 EPSRC award to University of Bristol and Dave Cliff	2012 (142) 2008 (125) 2007 (67) 2006 (36) 2000 (2) 1998 (1) 2005 (1) 1992 (1)
£4,923,928 Mar 06 -	Epidemiological Laboratory and IT Development (Area C) MRC award to Medical Research Council and Rory Collins	
£4,780,286 Aug 11 - Aug 16	EPSRC Centre for Innovative Manufacturing in Advanced Metrology EPSRC award to University of Huddersfield and Xiangqian Jiang	
£4,569,566 Mar 09 - Feb 14	Centre for Secure Information Technologies (CSIT) EPSRC award to Queen's University of Belfast and John McCanny	
£4,550,814 Jul 09 - Jul 14	Numerical Algorithms and Intelligent Software for the Evolving HPC Platform EPSRC award to University of Edinburgh and Nigel Brown	
£4,364,751 May 10 - May 15	SSI: The UK Software Sustainability Institute EPSRC award to University of Edinburgh and Neil Chue Hong	
£4,340,016 Oct 12 - Oct 16	Experimental Particle Physics Consolidated Grant (2012-2016) STFC award to University College London and David Waters	
£4,131,937 Mar 05 -	Analysis and prediction of protein structure and sequence MRC award to MRC National Inst for Medical Research and Willie Taylor	
£3,989,307 May 11 - May 16	A new approach to Science at the Life Sciences Interface EPSRC award to University of Oxford and David Gavaohan	

The Software Sustainability Institute



www.software.ac.uk

A national facility for cultivating world-class research through software

- Better software enables better research
- Software reaches boundaries in its development cycle that prevent improvement, growth and adoption
- Providing the expertise and services needed to negotiate to the next stage
- Developing the policy and tools to support the community developing and using research software



Supported by EPSRC
Grant EP/H043160/1

SSI: Long Term Goals



www.software.ac.uk

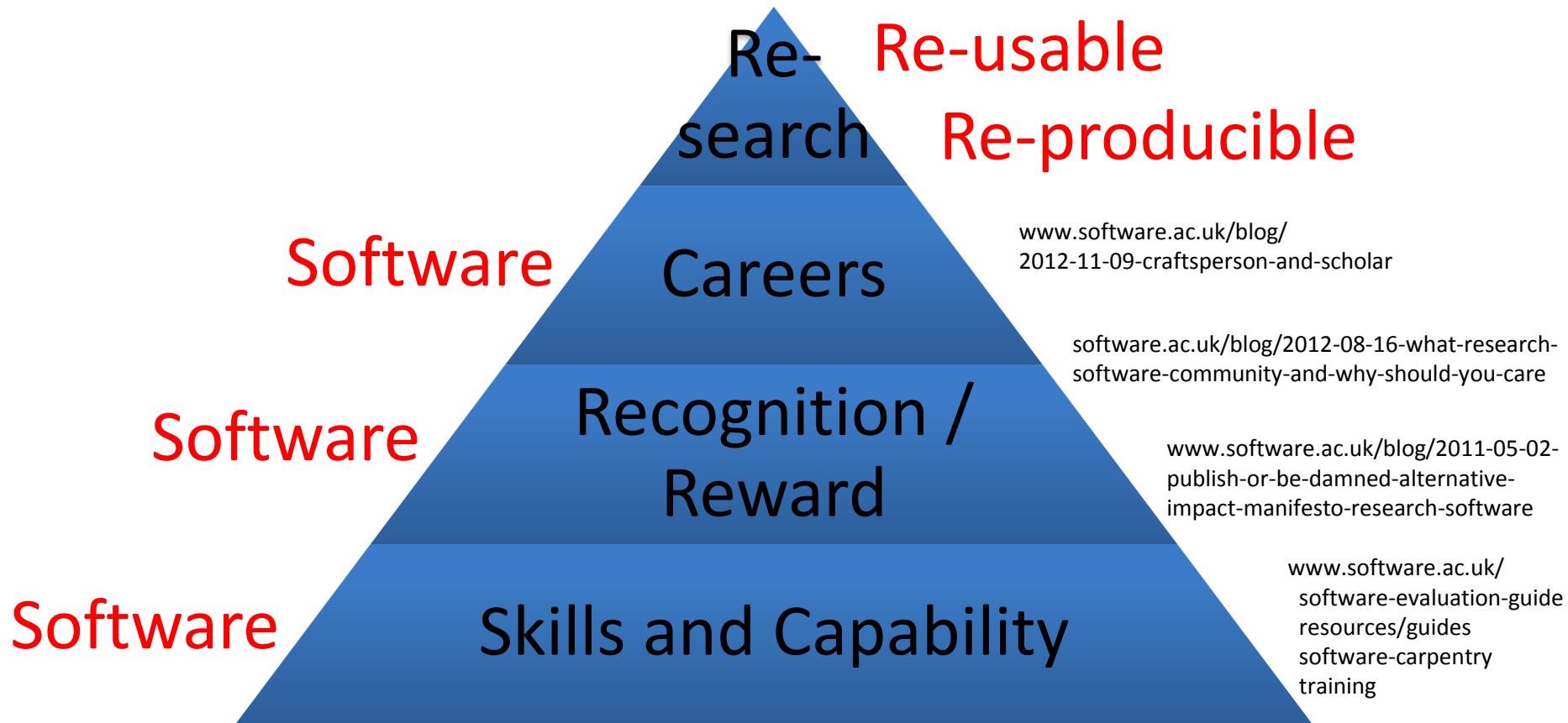
- Provision of useful, effective services for research software community
- Development and sharing of research community intelligence and interactions
- Promotion of research software best practice
- Mantra:
 - Keep the software in its respective community
 - Work with the community, to increase ability
 - Don't introduce dependency on SSI as the developer
 - Expand and exploit networks and opportunities



The Foundations of Digital Research



www.software.ac.uk



Prlić A, Procter JB (2012) Ten Simple Rules for the Open Development of Scientific Software
PLoS Comput Biol 8(12): e1002802. doi:10.1371/journal.pcbi.1002802

Wilson G, et al. (2013) Best Practices for Scientific Computing
Submitted to PNAS. <http://arxiv.org/abs/1210.0530>

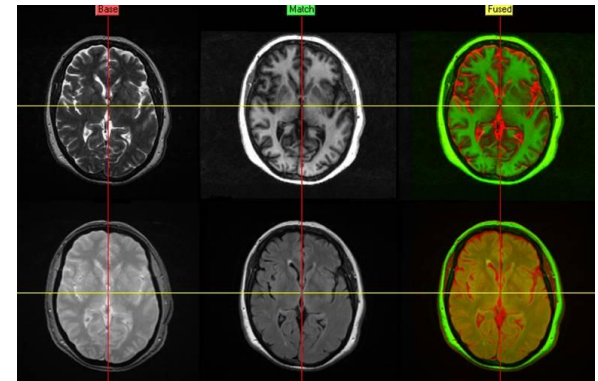


Case Study: Brain Imaging



www.software.ac.uk

- Brain Research Imaging Centre, Edinburgh
 - Develop PrivacyGuard software, a DICOM image deidentification toolkit
 - Created software to support new multispectral colouring modulation and variance identification technique (“MCMxxxVI”) to identify white matter lesions that are indicative of declining cognitive ability
 - BRIC are not principally software developers, but do provide software to other researchers
- SSI’s work means the software has been reviewed and refactored
 - Looked at exploitation
 - Usability review, Naming/trademark review
 - Made it easier for BRIC staff to maintain and develop
 - Move to standard repositories, testing and documentation processes
 - Examination of licencing for MCMxxxVI
 - Extraction and refactoring to create standalone tools
- <http://www.software.ac.uk/who-do-we-work/brain-research-imaging-centre-edinburgh>
- <http://www.bric.ed.ac.uk/>





www.software.ac.uk

Software
sustainability
requires thriving
communities of users
and developers:
how do we do this?

5 Stars of Scientific Software



www.software.ac.uk

- We need a 5 stars for software!
 - Existence – there is accurate metadata that defines the software
 - Availability – you can access and run the software
 - Openness – the software has an open permissible license
 - Linked – related data, dependencies and papers are referenced
 - Assured – the software provides ways of demonstrating “correctness”



c.f.
5 Stars of Linked Data
(Berners-Lee)
5 Stars of Online Journals
(Shotton)

Discoverable Software?



www.software.ac.uk

- To grow a community around software, first it must be discoverable
 - For users, wanting to find a solution
 - For developers, wanting to reuse or extend
 - For funders, wanting to promote or feature
- For sustainability
 - Provide useful information
 - Make it easier to attract and add contributors
 - Enable dormant projects to re-activate?



Software Hub Prototyping



www.software.ac.uk

Jisc Software Hub How do people search for software?

Xerte

Posted January 22, 2013 by admin & filed under Authoring, Learning and Teaching, Open Educational Resources.



Xerte Online Toolkits is a server-based suite of tools for content authors. Elearning materials can be authored quickly and easily using browser-based tools, with no programming required. Xerte Online

What can be imported from other sites?
What metadata must be collected to produce this information?
What is the user benefit?

Search [Search] Search

Openness Rating

legal	77%
governance	38%
standards	67%
knowledge	72%
market	74%

about the Openness Rating

Programming Factoids

- Increasing year-over-year development activity
- Large, active development team
- Well-established codebase
- Average number of code comments

What information is useful?
Do both provider and user benefit?

Levels of Showcasing



www.software.ac.uk

- Level 1: internal
 - Has had support from Jisc
 - Has produced a software output
 - Metadata is incomplete
- Level 2: awaiting approval
 - enough metadata to publish it externally
 - perhaps not all quality criteria met
- Level 3: published
 - meets quality criteria
 - enough information to allow comparison
- Level 4: featured
 - seen as particularly useful, exciting, best of breed etc.
 - associated screencasts, tutorials to show off
- Offer incentives to move up the levels



Collecting Software Metadata



www.software.ac.uk

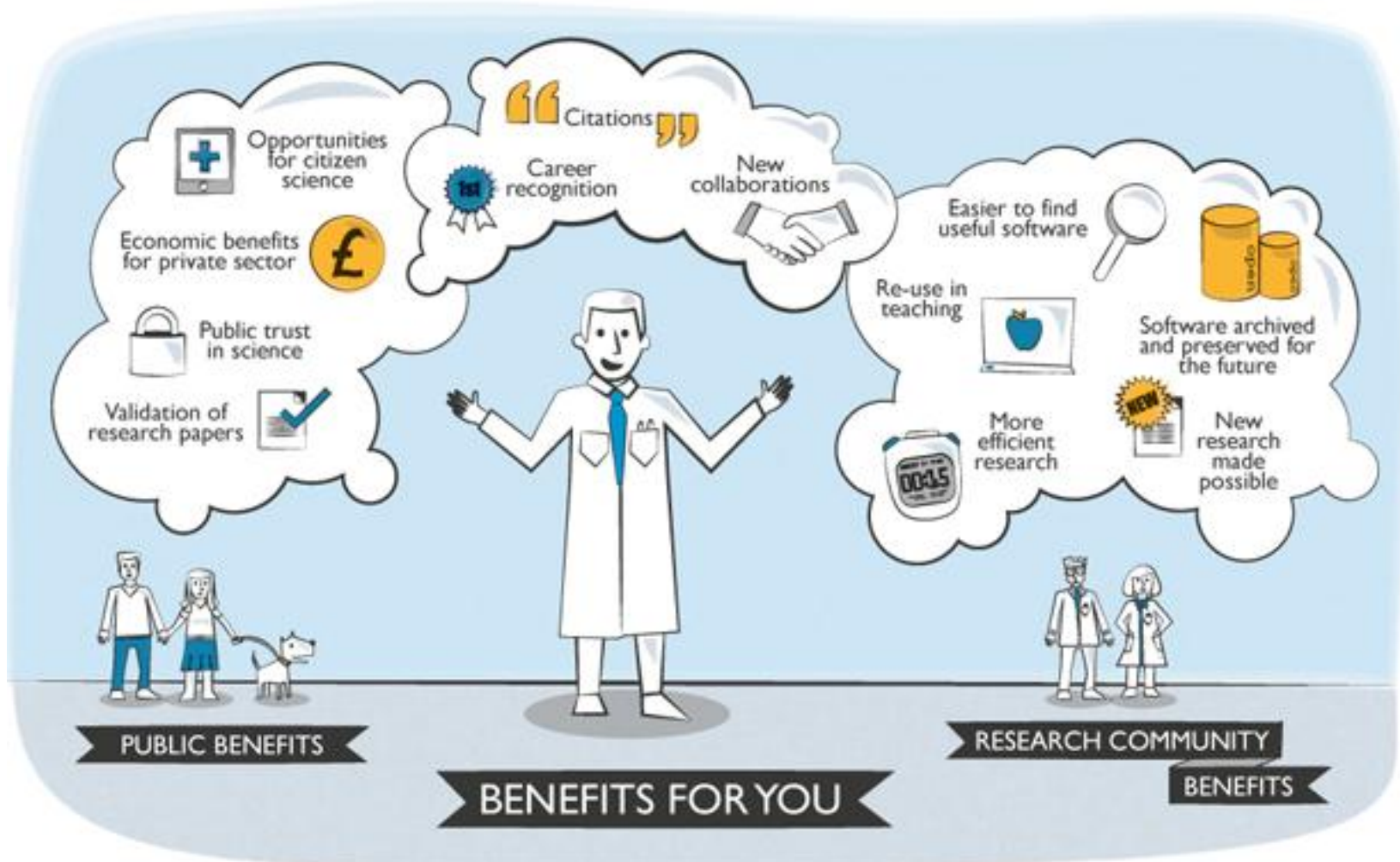
- Can we make the software metadata collection process work?
 - What are the benefits to provider and user?
 - Distinction between Project information and Product Information
 - Difference between information that enables discovery and choice, and the metadata that allows this information to be displayed
 - E.g. “vitality” of project different for developer vs user



Citable Software?



www.software.ac.uk

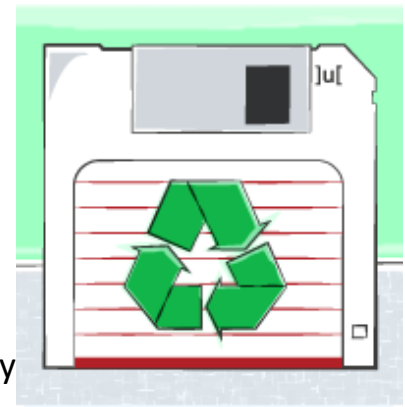


Software Metapapers



www.software.ac.uk

- Create a complete scholarly record including “standard” publication, method, dataset and models, and software
 - e.g. modelling and simulation, statistical analysis
 - Enable replay, reproduction and reuse
- Pragmatic approach is to create a metadata record for the software, and link it to a copy of the software in some storage infrastructure
 - This is a software metapaper
 - Peer-review the metadata, not the software
- Journal of Open Research Software:
 - <http://openresearchsoftware.metajnl.com/>



See: <http://openresearchsoftware.metajnl.com/faq/>

and the work by B. Matthews et al: The Significant Properties of Software: A Study





- Does the DataCite approach work with software?
 - What is the cost of minting a DOI?
 - What level do you mint DOIs for software?
 - What is the cost of storing the metadata associated with a software asset?
 - What is the cost of a software asset associated with a DOI disappearing?

Alternative Impact Stories



www.software.ac.uk

ImpactStory.

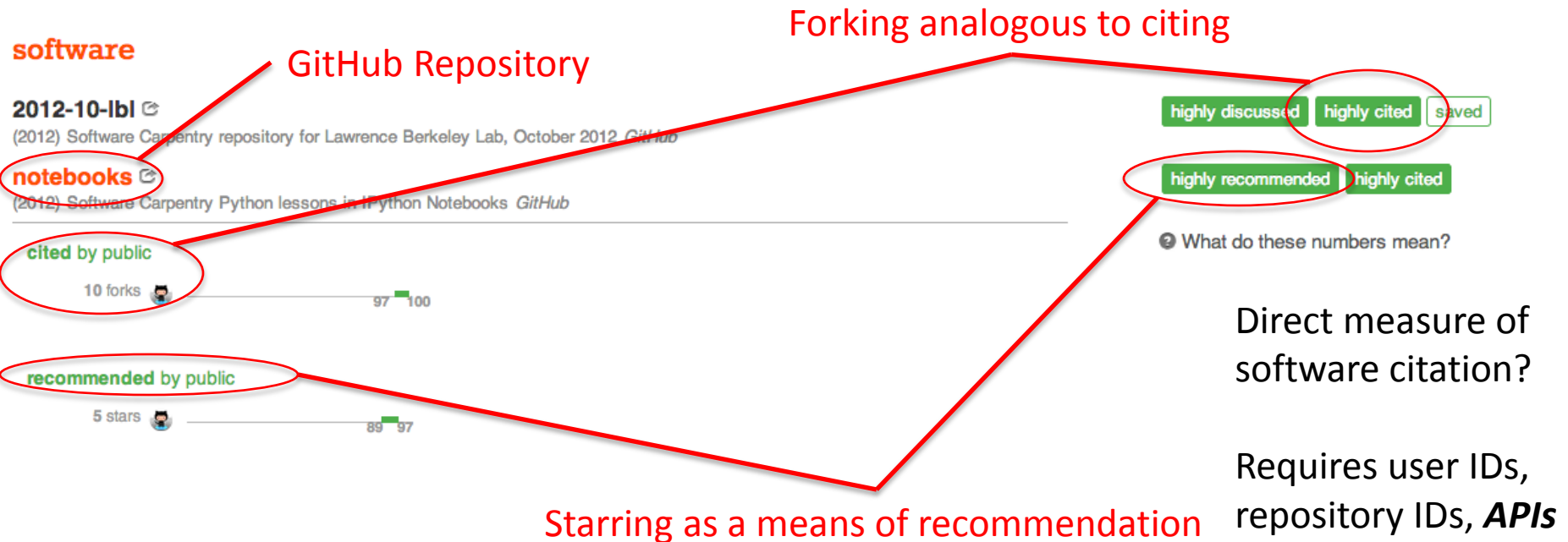
embed create follow about

hi, N.ChueHong@software.ac.uk! (logout)

BETA
Send us your feedback!

Software Carpentry

19 items (expand all) update json csv Tweet 0





www.software.ac.uk

Software
sustainability
requires identifiers
(and open APIs)
for discovery, growth,
attribution + reward

Why is this important?
An animated explanation
Featuring a frustrated panda



<http://tinyurl.com/datasharingpanda>