

The Path to Building Consensus: Can good minimum information checklists be created by the community?

Dawn Field

NERC Centre for Ecology and Hydrology &
the University of Oxford e-Research Centre



Steps toward a standard of Minimum Information About a Phylogenetic Analysis (MIAPA)

*Hilmar Lapp, Nico Cellinese, Jim Leebens-Mack,
Enrico Pontelli, Arlin Stoltzfus*

ABSTRACT Many phylogenetic analysis results are published in ways that present serious barriers to their reuse in numerous research applications that would stand to benefit from them. While some of these barriers are well understood, such as issues with adherence to standard exchange formats, those centering on the associated metadata necessary for researchers to evaluate or reuse a published phylogeny have only recently begun to be articulated. One of the critical next steps towards formalizing these metadata requirements as a minimum reporting standard is to convene meetings of key stakeholder communities with the goal to identify information attributes necessary and desirable for facilitating reuse, and to build consensus on their priority. To this end, we are holding a workshop at the 2011 Biodiversity Information Standards (TDWG) Conference to determine how a future reporting standard for phylogenetic analyses can best serve biodiversity science and related research applications. We invite all interested colleagues to participate.



What are Standards?

- A standard is a convention that gives uniformity to an area of research or innovation.
- Standards unite groups and enable collective change.
- Standards provide the language in which innovation is written.



standards

Principles:

Not everything should be ‘standardized’

Aggregation of data, information, and knowledge requires
standard ways of doing things

Standards provide foundations; Standards should drive innovation
(think of electrical plugs or the internet)

Pick the right concepts to standardize – at the right time, with
the right people

Requires good ‘group think’ – or ‘systems thinking’

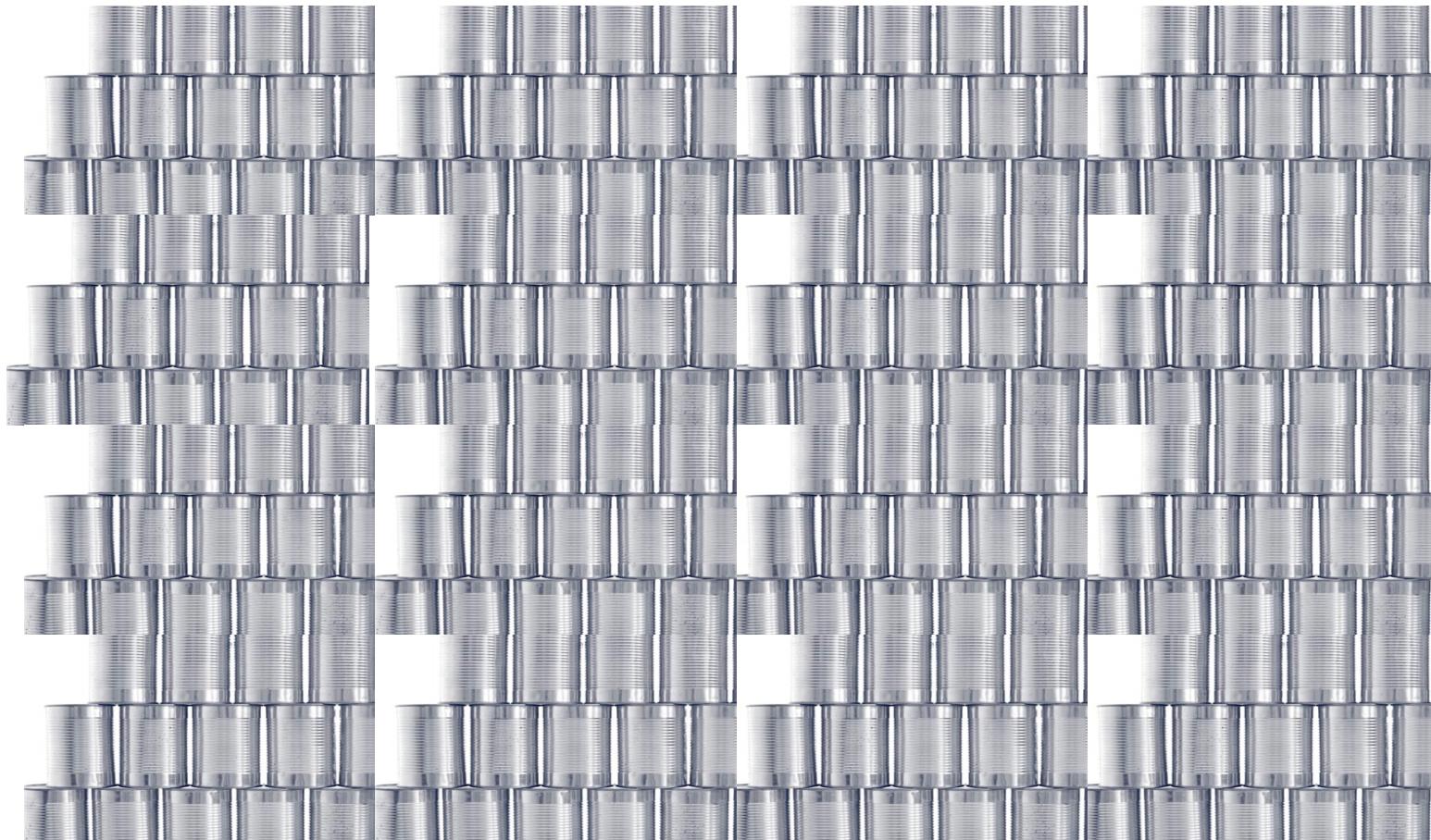


SuperMarket



DataMarket

Norman Morrison



Packaging data



Labels for data

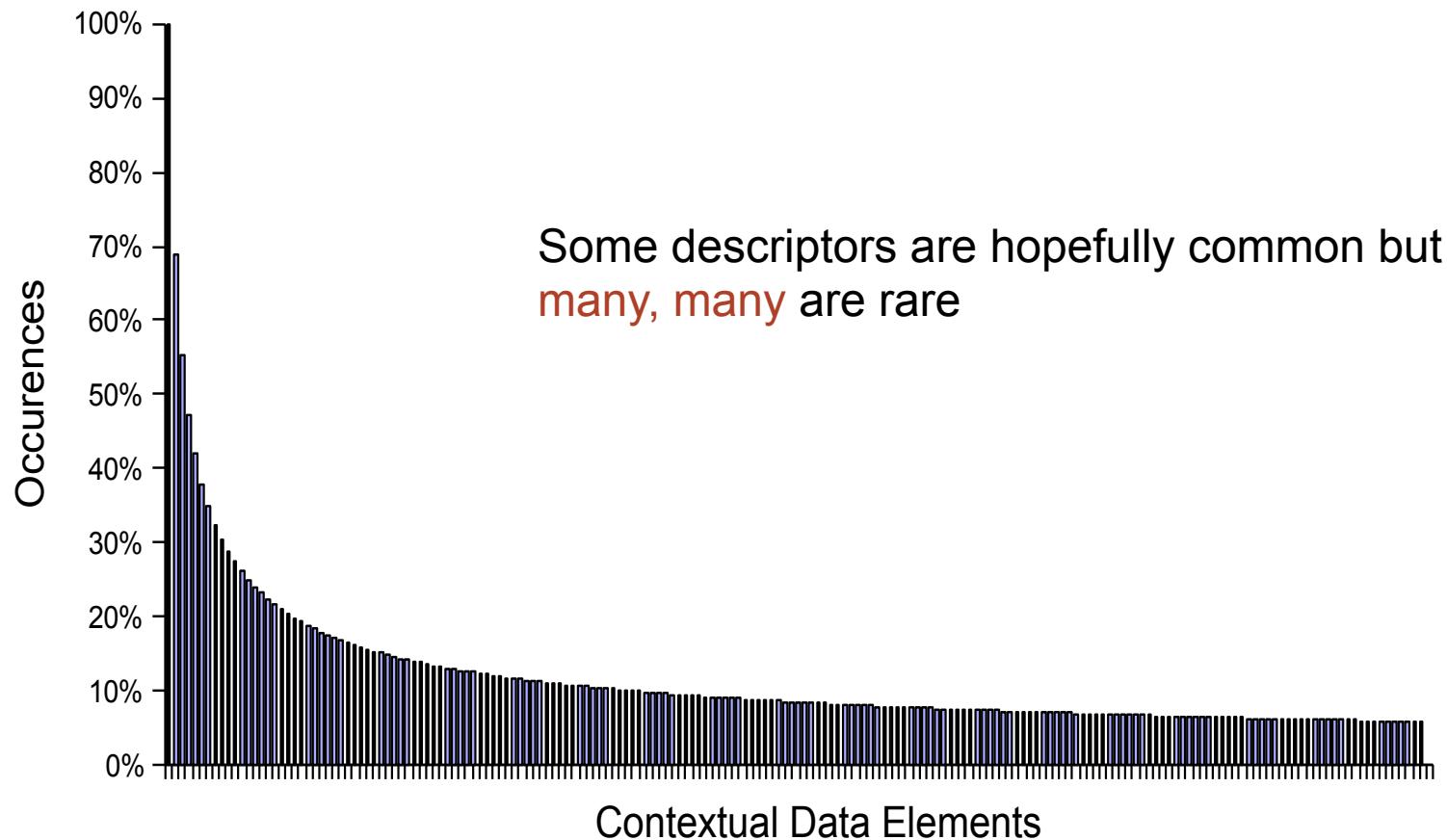


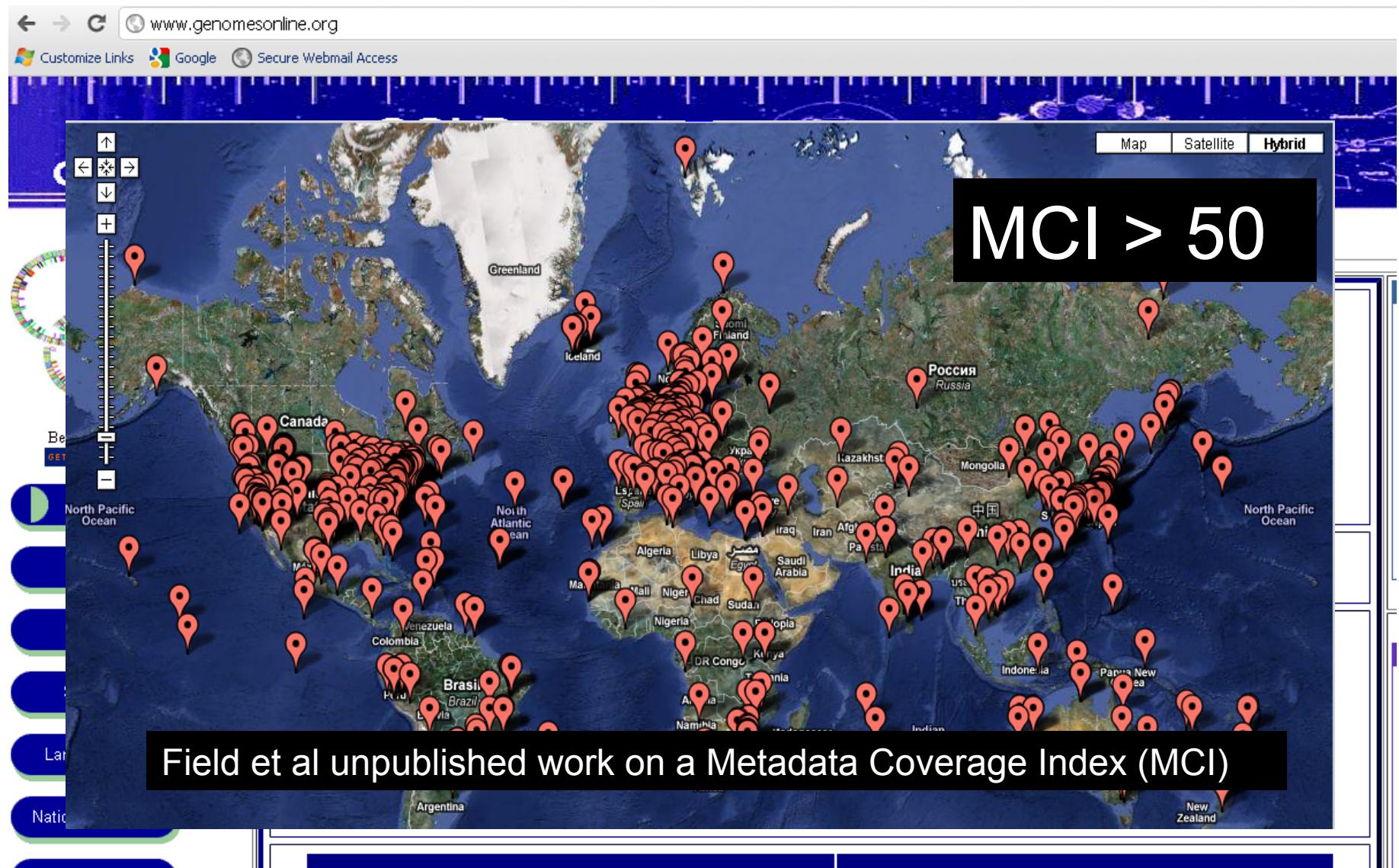
What does your metadata look like?

- How easily will you be able to build consensus?
- If you polled everyone in the room for the top 5-10 descriptors how many would be the same?
- Surveys are helpful
- When you start to reach consensus, polling new people largely yeilds the same results
- Can we make this an objective process?



Truth: Prevalence of descriptors





Community-driven solutions:

The Common Path:

- Identify the problem
- Define a community to address it
- Define scope of the solution
- Implement solution
- Gain adoption of solution

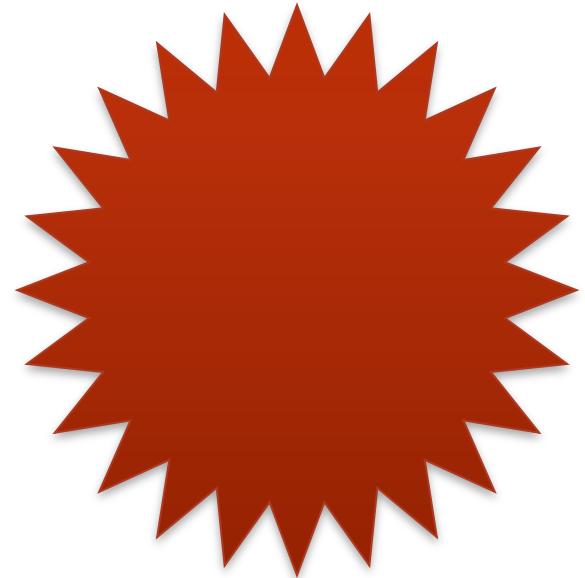
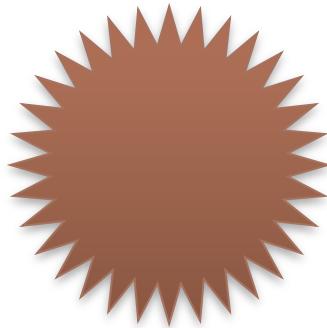


Progression of ideas

- Pre-meeting thoughts – wish lists
- Beer discussions
- **Good suggested concepts**
- Working groups – with at least one champion
- Formal launch at a meeting – ideally, recorded in meeting report publication
- Actively producing outputs (papers, standards, tools)
- Funded (workshops, development)
- **Changing the landscape of your field**
- **Changing the landscape of science**



Joining up Communities



The Genomic Standards Consortium



Innovation through Collaboration

- Established in September 2005
- Governed by a Board
- 100+ members
- Representatives from **INSDC** (EMBL/DDBJ/GenBank), EBI, Sanger, JGI, JCVI, GOLD, CAMERA, MG-Rast, RDP, Silva, Greengenes, VAMPS, and many more
- Open membership defined by participation

GSC 10
Argonne, 2010



GSC 11,
Hinxton, 2011



GSC 12
Bremen, 2011



GSC 13
BGI 2012



The minimum information about a marker gene sequence (MIMARKS) and minimum information about any (x) sequence (MIxS) specifications

http://www.nature.com/naturebiotechnology

Dawn Field^{1*}, George Gar Nicholas Thomson⁸, Mich Sandra Baldauf¹², Stuart B Claude dePamphilis¹⁸, Rol Frank Oliver Glöckner²³, I Henning Hermjakob⁶, Chi Jessie Kennedy²⁷, George I Jim Leebens-Mack³³, Suza Victor Markowitz³⁷, Jenni Julian Parkhill⁸, Lita Proct Paul Swift¹, Chris Taylor⁶, Naomi Ward⁴⁵, Trish Whe

Minimum information about a marker gene sequence (MIMARKS) and minimum information about any (x) sequence (MIxS) specifications

Pelin Yilmaz^{1,2*}, Renzo Kottmann¹, Dawn Field³, Rob Knight^{4,5}, James R Cole^{6,7}, Linda Amaral-Zettler⁸, Jack A Gilbert^{9–11}, Ilene Karsch-Mizrachi¹², Anjanette Johnston¹², Guy Cochrane¹³, Robert Vaughan¹³, Christopher Hunter¹³, Joonhong Park¹⁴, Norman Morrison^{3,15}, Philippe Rocca-Serra¹⁶, Peter Sterk³, Manimozhiyan Arumugam¹⁷, Mark Bailey³, Laura Baumgartner¹⁸, Bruce W Birren¹⁹, Martin J Blaser²⁰, Vivien Bonazzi²¹, Tim Booth³, Peer Bork¹⁷, Frederic D Bushman²², Pier Luigi Buttigieg^{1,2}, Patrick S G Chain^{7,23,24}, Emily Charlson²², Elizabeth K Costello⁴, Heather Huot-Creasy²⁵, Peter Dawyndt²⁶, Todd DeSantis²⁷, Noah Fierer²⁸, Jed A Fuhrman²⁹, Rachel E Gallery³⁰, Dirk Gevers¹⁹, Richard A Gibbs^{31,32}, Inigo San Gil³³, Antonio Gonzalez³⁴, Jeffrey I Gordon³⁵, Robert Guralnick^{28,36}, Wolfgang Hankeln^{1,2}, Sarah Highlander^{31,37}, Philip Hugenholtz³⁸, Janet Jansson^{23,39}, Andrew L Kau³⁵, Scott T Kelley⁴⁰, Jerry Kennedy⁴, Dan Knights³⁴, Omry Koren⁴¹, Justin Kuczynski¹⁸, Nikos Kyriakis²³, Robert Larsen⁴, Christian L Lauber⁴², Teresa Legg²⁸, Ruth E Ley⁴¹, Catherine A Lozupone⁴, Wolfgang Ludwig⁴³, Donna Lyons⁴², Eamonn Maguire¹⁶, Barbara A Methé⁴⁴, Folker Meyer¹⁰, Brian Muegge³⁵, Sara Nakielsky⁴, Karen E Nelson⁴⁴, Diana Nemergut⁴⁵, Josh D Neufeld⁴⁶, Lindsay K Newbold³, Anna E Oliver³, Norman R Pace¹⁸, Giriprakash Palanisamy⁴⁷, Jörg Peplies⁴⁸, Joseph Petrosino^{31,37}, Lita Proctor²¹, Elmar Pruesse^{1,2}, Christian Quast¹, Jeroen Raes⁴⁹, Sujeewan Ratnasingham⁵⁰, Jacques Ravel²⁵, David A Relman^{51,52}, Susanna Assunta-Sansone¹⁶, Patrick D Schloss⁵³, Lynn Schriml²⁵, Rohini Sinha²², Michelle I Smith³⁵, Erica Sodergren⁵⁴, Aymé Spor⁴¹, Jesse Stombaugh⁴, James M Tiedje⁷, Douglas V Ward¹⁹, George M Weinstock⁵⁴, Doug Wendl¹⁴, Owen White²⁵, Andrew Whitlow³, Andreas Wilcock¹⁰



Promoting coherent minimum reporting guidelines for biological and biomedical investigations: the MIBBI project

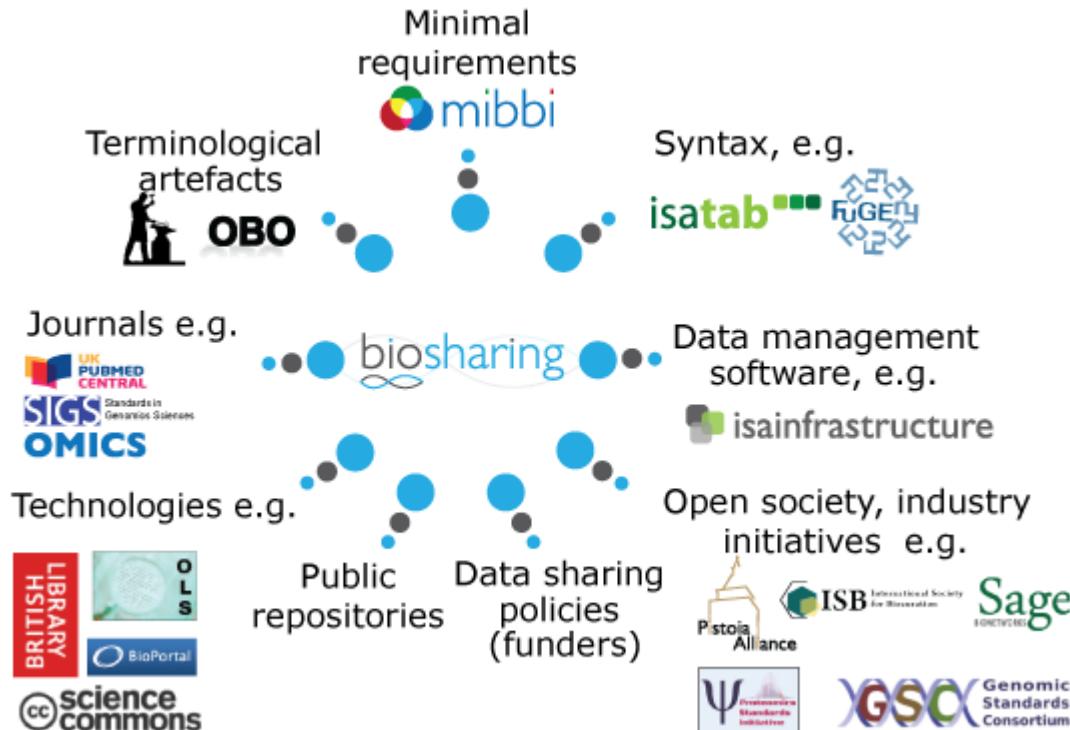
Chris F Taylor^{*1,2}, Dawn Field^{2,3}, Susanna-Assunta Sansone^{1,2}, Jan Aerts⁴, Rolf Apweiler¹, Michael Ashburner⁵, Catherine A Ball⁶, Pierre-Alain Binz^{7,8}, Molly Bogue⁹, Tim Booth², Alvis Brazma¹, Ryan R Brinkman¹⁰, Adam Michael Clark¹¹, Eric W Deutsch¹², Oliver Fiehn¹³, Jennifer Fostel¹⁴, Peter Ghazal¹⁵, Frank Gibson¹⁶, Tanya Gray^{2,3}, Graeme Grimes¹⁵, John M Hancock¹⁷, Nigel W Hardy¹⁸, Henning Hermjakob¹, Randall K Julian Jr¹⁹, Matthew Kane²⁰, Carsten Kettner²¹, Christopher Kinsinger²², Eugene Kolker^{23,24}, Martin Kuiper²⁵, Nicolas Le Novère¹, Jim Leebens-Mack²⁶, Suzanna E Lewis²⁷, Phillip Lord¹⁶, Ann-Marie Mallon¹⁷, Nishanth Marthandan²⁸, Hiroshi Masuya²⁹, Ruth McNally³⁰, Alexander Mehrle³¹, Norman Morrison^{2,32}, Sandra Orchard¹, John Quackenbush³³, James M Reecy³⁴, Donald G Robertson³⁵, Philippe Rocca-Serra^{1,36}, Henry Rodriguez²², Heiko Rosenfelder³¹, Javier Santoyo-Lopez¹⁵, Richard H Scheuermann²⁸, Daniel Schober¹, Barry Smith³⁷, Jason Snape³⁸, Christian J Stoeckert Jr³⁹, Keith Tipton⁴⁰, Peter Sterk¹, Andreas Untergasser⁴¹, Jo Vandesompele⁴² & Stefan Wiemann³¹



MEGASCIENCE

'Omics Data Sharing

Dawn Field,^{1*}†‡ Susanna-Assunta Sansone,^{1,2†} Amanda Collis,^{3†} Tim Booth,¹ Peter Dukes,⁴ Susan K. Gregurick,⁵ Karen Kennedy,⁶ Patrik Kolar,⁷ Eugene Kolker,⁸ Mary Maxon,⁹ Siân Millard,¹⁰ Alexis-Michel Mugabushaka,¹¹ Nicola Perrin,¹² Jacques E. Remacle,⁷ Karin Remington,¹³ Philippe Rocca-Serra,¹² Chris F. Taylor,¹² Mark Thorley,¹⁴ Bela Tiwari,¹ John Wilbanks¹⁵



Data sharing, and the good annotation practices it depends on, must become part of the fabric of daily research for researchers and funders.

¹U.K. Natural Environment Research Council (NERC), Environmental Bioinformatics Centre. ²European Molecular Biology Laboratory (EMBL) Outstation, The European Bioinformatics Institute (EBI). ³U.K. Biotechnology and Biological Sciences Research Council. ⁴U.K. Medical Research Council. ⁵U.S. Department of Energy. ⁶Genome Canada. ⁷Unit for Genomics and Systems Biology, European Commission. ⁸Seattle Childrens Hospital. ⁹Marine Microbiology Initiative, Gordon and Betty Moore Foundation. ¹⁰U.K. Economic and Social Research Council. ¹¹European Science Foundation. ¹²The Wellcome Trust. ¹³U.S. National Institute of General Medical Sciences, NIH. ¹⁴NERC. ¹⁵Science Commons.

* The first three authors contributed equally to this article.



“ BioSharing works to facilitate data sharing policies and the communication of high-quality sharing of high-quality data.”

We work with communities to:

1. develop catalogues - enriching the seeking process
2. moderate and promote policies

POLICIES



A catalogue of data preservation, management and sharing policies from international funding agencies and regulators.

STANDARDS



A catalogue of domain-specific data formats and terms used by organizations to describe their data.

OLS



A CATALOGUE OF STANDARDS

You can **sort** columns and **browse** the reporting guidelines content, or you can view [all the standards](#), or [terminological artifacts](#) or [exchange formats](#) only; or go back to the [catalogue main page](#).

ACRONYM	FULL NAME	TYPE▲	DOMAIN	VERSION	PUBLICATION	CONTACT
BioPAX	Biological Pathway Exchange	exchange format	biological pathway	Level 3	Demir et al; Nat Biotech; 2010	BioPAX community
CellML	Cell Markup Language	exchange format	cell modelling	v 1.1	Cue Sim	
SBML	System Biology Markup Language	exchange format	computational modelling (biochemical reaction networks)	level 3, v 1 core	Huck Bio 200	
FuGE	Functional Genomics Experiment Markup Language	exchange format	experimental description	v 1.0	Jon Bio	
ISA-Tab	Investigation/Study/Assay Tabular	exchange format	experimental description	v 1.0	Roc Bio 201	
MINiML	MIAME Notation in Markup Language	exchange format	experimental description (functional genomics)	v 1.0		
SOFT						
GCDML						
MAGE-Tab						
GelML						
mzML						
MIABIE						
MIPFE						
BioCoreDB						
GIATE	About Therapy Experiments	guideline	experiments	not specified	Eng Des Sel; 2009	Antibody Society
MIRIAM	Minimal Information Required In the Annotation of biochemical Models	reporting guideline	computational modeling	not specified	Le Novère et al; Nature Biotech; 2005	BioModels.net



SBML
System Biology Markup Language

ID
bsg-000052

TYPE
exchange format

DOMAIN(S) COVERED
computational modelling (biochemical reaction networks)

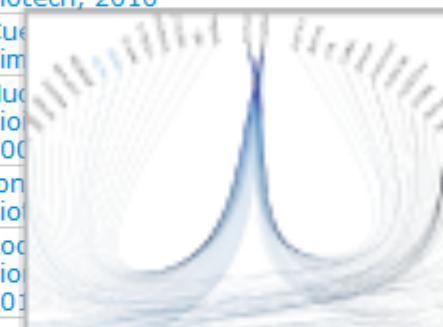
PUBLICATION(S)
Hucka et al; Bioinformatics; 2003

ORGANIZATION
SMBL community

MAIN CONTACT(S)

Like
 Tweet 0

» [Login](#) or [register](#) to post comments



show relation between
exchange formats and
reporting guidelines

Escalating number of standardization efforts in bioscience, e.g.:



formats



terminologies



guidelines



To exploit fully the promise of these data
we need both scientific innovation and
community agreement on how to provide
appropriate stewardship of these
resources for the benefit of all.

Requires the evolution of our scientific, technological and
sociological thinking....



Conclusions

- The era of real data sharing is just beginning...
- Self-organization by the scientific community can pay dividends (i.e. consensus building, large-scale co-ordination)
 - Standards are keys to unlocking data
 - Group thinking overcomes the tragedy of the commons
- Many communities and ‘solutions’ to learn from and work with

