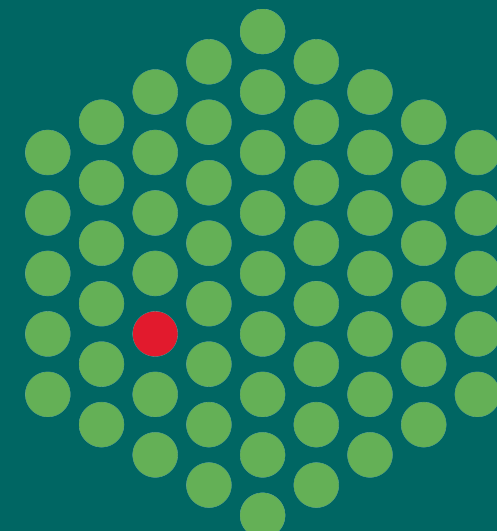


10 simple rules for making data web friendly

EMBL-EBI



Nick Juty¹, Julie McMurry^{2,3}, BioMedBridges Consortium⁴, Melissa Haendel^{2,3}, Carole Goble^{5,6}, Helen Parkinson¹

¹EMBL-EBI, Hinxton, Cambridge, UK. CB10 1SD.

²Department of Medical Informatics and Epidemiology and OHSU Library, Oregon Health & Science University, Portland, USA.

³Monarch Initiative; <http://monarchinitiative.org/>

⁴BioMedBridges Consortium; <http://www.biomedbridges.eu/partners>

⁵School of Computer Science, The University of Manchester, Manchester, UK.

⁶Open PHACTS; <http://www.openphactsfoundation.org/>



Introduction

Life Science data continues to grow, becoming increasingly available via the Web. Our handling of identifiers for this distributed data, and lack of consideration for how data evolves, has led to phenomena such as 'link rot' (dead links) and 'content drift' (evolved data underlying a stagnant identifier). In addition, poorly conceived strategies for making data available on the Web, built upon poor choices in identifier design and lack of consideration of the global data landscape, is causing downstream integration issues. As a consequence, there is often a need to provide additional and potentially costly procedures for the source-dependent processing of data, or else the provision of alternative mapping solutions. BioMedBridges, in association with an international community of partners, have begun to determine the common technical bridges required to allow data integration across the biological domain. Based on our experience we describe ten simple rules for best practice in the provision and reuse of identifiers for web-based Life Science data [1].

Re-use established identifiers

Never create new identifiers unnecessarily

- *Link using native identifiers*
- *Document relationships to existing entities*
eg. owl:sameAs, skos:broader

Follow identifier conventions

Follow conventions when minting new identifiers

- *No internal whitespace*
- *Avoid 'date' or 'exponent' misinterpretation*
- *Specify a 'pattern' for identifiers*

Assign namespaces

Provide a namespace and mapping to allow local identifiers to be used globally

- *Document the namespace used*
- *Document the mapping to a resolving location*
- *Register namespace where appropriate*

Some desirable identifier characteristics

Unambiguous – associated with a single entity

Stable – identifier to entity relationship

Versionable – to allow evolution of entity or metadata

Persistent – identifier is never deleted

Resolvable – information accessible

Defined format – adheres to documented pattern

Web ready – contains no reserved characters or special handling

Documented – identifier scheme and policy available

Reference responsibly

Use appropriate referencing and ...

- *Always provide the full URI*
- *Consider a CURIE for narrative online text*
- *Provide additional file to document mappings if necessary*
- *Attempt to ensure that all cross-references remain 'live'*

<http://identifiers.org/registry>

Provide URIs

Implement intuitive and resolvable URIs

- *Avoid format, file extension and administrative parameters*
- *Provide a landing page*
- *Consider CURIEs for display*

CURIE

Compact URI = <Prefix>:<Local Identifier>

Prefix is defined and documented

Allows deterministic expansion to resolving location

Mnemonic useful in human communication

Suitable for convenient display to users

ZFIN:ZDB-GENE-980526-166

Document your practices

Document the identifiers you create, and those you use from other providers. Include HCLS Dataset Description [2], e.g.

- *Entity scope*
- *licence*
- *Namespace and mapping*
- *Machine-readable formats*



Avoid embedded semantics

Would embedding semantic information in an identifier make it prone to obsolescence in an evolving data domain?

- *Define the entity being identified*
- *Provide metadata about the entity*
- *Semantic information acceptable in some cases (InChI)*

UniProtKB - P29358 (143B_BOVIN)



This entry is obsolete

On October 25, 2004 this entry became obsolete. It can now be found as secondary accession in P68250 and P68251. [List]
For previous versions of this entry, please look at its history.

History

Previous versions for UniProtKB entry: P12345

Download

	Versions	Sequence	Info	Releases
	Entry		Entry name	Database
<input type="checkbox"/>	100 bt	2 fasta	AATM_RABIT	Swiss-Prot
<input type="checkbox"/>	99 bt	2 fasta	AATM_RABIT	Swiss-Prot
<input type="checkbox"/>	98 bt	2 fasta	AATM_RABIT	Swiss-Prot
<input type="checkbox"/>	97 bt	2 fasta	AATM_RABIT	Swiss-Prot

For a permanent link to this page, which will not change with the next release of Ensembl, use:

http://Sep2015.archive.ensembl.org/Mus_musculus/Gene/Summary?g=ENSMUSG00000033577;r=9:80165031-80311729#

We aim to maintain all archives for at least two years; some key releases may be maintained for longer

Ensembl

Never reassign

Never re-assign an identifier to a different record. (cf. Ghost Busters 'don't cross the streams')

- *Provide a list of all obsolete identifiers*
- *Create new identifier on record merge*
- *Create new identifiers for demerged*
- *Always provide links for the above*

Use versioning

Document change history at the record level, and/or provide versioned identifiers

- *Tombstone page for obsolete*
- *Use ':' to version local identifier*
- *provide links to latest record*

Make URIs clear

Allow users to easily identify the URI for reuse

- *Advertise a permanent link*
- *Clearly label out-dated records*
- *Consider providing 'cite this' button*

Summary

A plethora of issues undermine the identification of data in the Life Sciences, hampering our ability to integrate the ever increasing amounts of data being produced. We have outlined best practices for data providers, as well as guidance for data integrators and redistributors, to begin the process of solving at least some of these issues. We hope that the information provided here will also be useful to data generators and end users, allowing them to appreciate the additional complexity of making data available in a web environment. We anticipate that, over time, improved tooling will become available, lowering the barriers to adoption. We also wish to acknowledge the work of several groups internationally[3,4], who are also converging on standards that will be applicable to identifiers.

References

- McMurry J, et al. (2015) 10 Simple rules for design, provision, and reuse of identifiers for web-based life science data: Zenodo. 10.5281/zenodo.31765
- Gray AJG, Baran J, Marshall MS, Dumontier M. (2014) Identifiers in Dataset Descriptions: HCLS Community Profile. In: HCLS Community Profile: <http://www.w3.org/2001/sw/hcls/notes/hcls-dataset/>
- Data Citation Synthesis Group (2015) Joint Declaration of Data Citation Principles: <https://www.force11.org/datacitation>
- FORCE11 (2015) The FAIR data Guiding Principles: <https://www.force11.org/group/fairgroup/fairprinciples>



EMBL-EBI
Wellcome Trust Genome Campus
Hinxton, Cambridgeshire, CB10 1SD, UK

Tel. +44 (0) 1223 494 444
juty@ebi.ac.uk
www.ebi.ac.uk