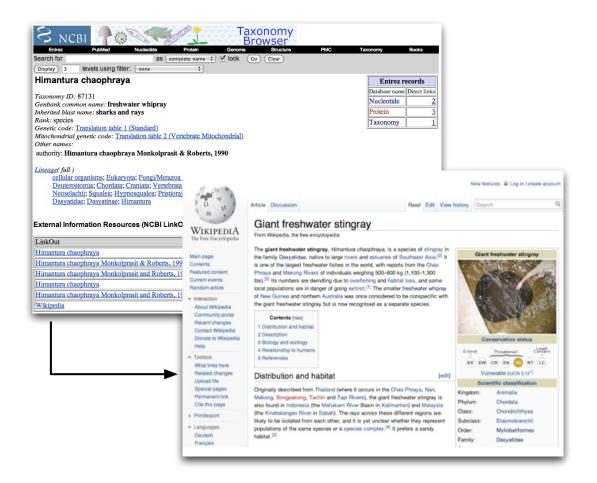# Phyloinformatics in the age of Wikipedia

Roderic D M Page
DEEB, FBLS
University of Glasgow, Glasgow G12 8QQ, UK
r.page@bio.gla.ac.uk

One of the great challenges of phyloinformatics is linking together information on phylogenies, taxa, genomes, specimens, and publications. One approach to linking disparate data is to use shared identifiers. For example, if a bibliographic database and a nomenclature database both use the same identifier for a publication (such as a DOI), then we can easily link the two pieces of information together using that identifier. An obstacle to this approach is the lack identifiers, or failure to reuse existing identifiers (Page, 2008). For example, sequences in GenBank may lack bibliographic identifiers, even if the paper in which the sequences were published has an identifier (Miller et al., 2009). This low "link density" is a major obstacle to linked data approaches to integrating biodiversity data. Link density can be increased either by reusing existing identifiers, or by creating maps between existing identifiers.

This talk describes a mapping between the NCBI taxonomy database and Wikipedia. These two databases were chosen because the NCBI taxonomy contains all the taxa for which sequences are publicly available, and for many taxa Wikipedia is the first site returned in a Google search on that taxon's scientific name (Page, 2010). The bulk of the mapping was created automatically by matching strings in the two databases, either directly, or via synonyms (NCBI) or redirect pages (Wikipedia). To support manual editing and correction the mapping was loaded into a wiki available at http://iphylo.org/linkout. Implemented using semantic wiki (http://semantic-mediawiki.org/), the wiki provides RDF export of the mapping. NCBI taxonomy ids are linked to the equivalent URI in the Uniprot database (http://www.uniprot.org), and Wikipedia pages are mapped to the equivalent URI in Dbpedia (http://dbpedia.org). Both Uniprot and Dbpedia provide data in RDF, which when combined with the mapping from http://iphylo.org/Linkout enable a wide range of queries. For example, we could query Wikipedia for images of taxa in GenBank, or query GenBank for sequences from taxa in the International Union for Conservation of Nature Red List of Threatened Species (http://www.iucnredlist.org/).

To date a total of 52,956 NCBI taxa have been mapped to the corresponding taxa in Wikipedia. The mapping for any individual NCBI taxon can be found by appending the taxonomy id to the URL "http://iphylo.org/linkout/Ncbi:", for example http://iphylo.org/linkout/Ncbi:87131. The mapping has also been uploaded NCBI's Linkout, so that the link to Wikipedia is shown on the NCBI page for the taxon (Fig. 1).

**Fig. 1.** The NCBI page for *Himantura chaophraya* (taxonomy id 87131) linked to the Wikipedia page for the Giant freshwater stingray.

## Acknowledgements

## References

Miller, Holly; Catherine Norton; Indra Neil Sarkar. (2009). GenBank and PubMed: How connected are they? BMC Research Notes 2009, 2:101 doi:10.1186/1756-0500-2-101

Page, R. D. M. (2008). Biodiversity informatics: the challenge of linking data and the role of shared identifiers. Briefings in Bioinformatics 9:345-354. doi:10.1093/bib/bbn022

Page, R. D. M. (2010). Wikipedia as an encyclopaedia of life. Nature Precedings hdl:10101/npre.2010.4242.1