

**Figure 1.** Analysis of links between the Entrez Gene and PubMed databases. **(A)** Examining the degree of gene annotation from the perspective of Entrez Gene, we found that while a few genes are very well annotated with links to PubMed references, the vast majority of genes have few or no linked references. **(B)** Examining links from the perspective of PubMed, we found that only a small fraction of published articles are linked to human genes. Taken together, these findings suggest that the traditional model of centralized curation is not scaling well with the rate of scientific research, and that complementary approaches based on community intelligence may be worth exploring.

has been very successful in its goal to systematically advance gene annotation, creating essential tools and ontologies in the process.

However, this model alone may not be sufficient to efficiently and systematically annotate gene function. Many leading voices in the gene annotation and model organism communities recently wrote a feature article in *Nature* describing the current state and future of biocuration (8). They noted the immense challenge to the curator community (typically numbering in tens to hundreds of people) to keep pace with the biomedical literature (currently 18 million articles in PubMed, roughly 750 000 new articles per year). Specifically, these curation experts suggest that merely preserving the existing models of gene annotation will lead to an increasing lag between curated data and biological knowledge, and that 'sooner or later, the research community will need to be

involved in the annotation effort to scale up to the rate of data generation' (8).

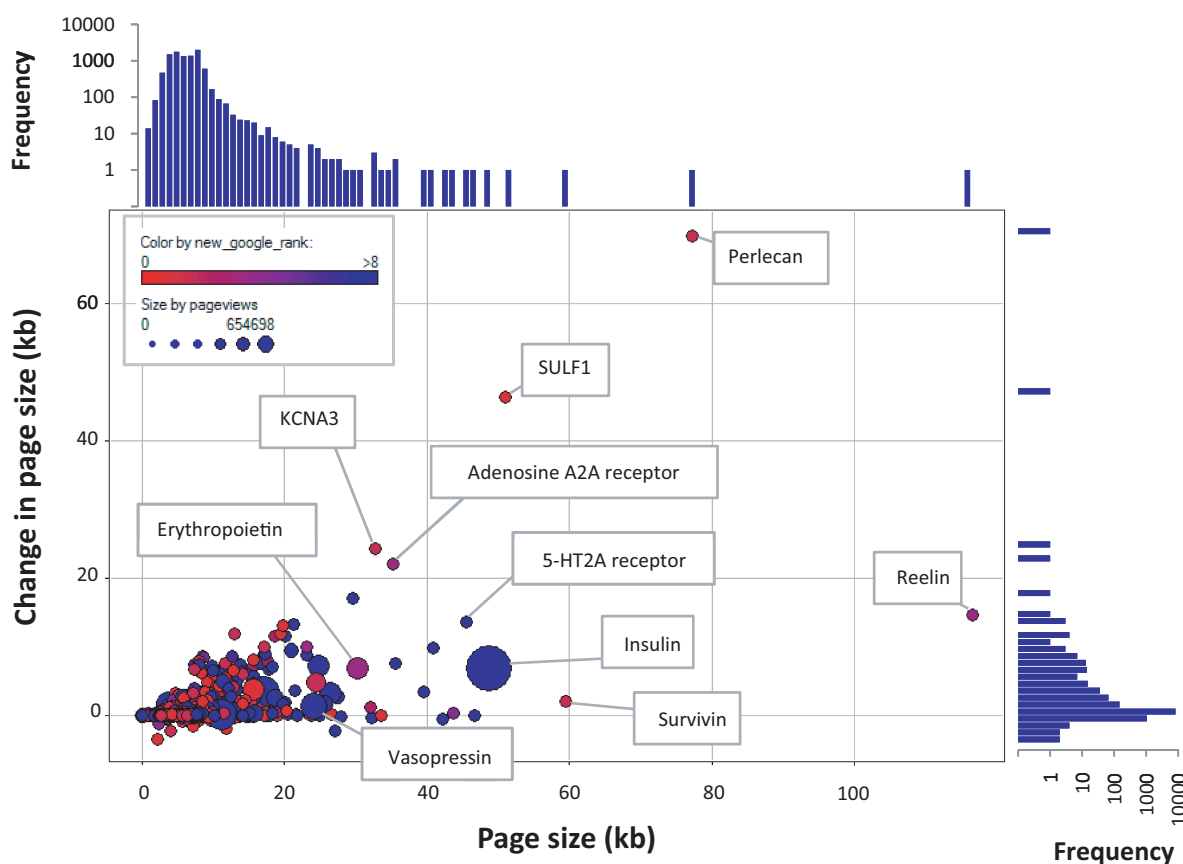
Thus, although leaders in the curation community have successfully set up a robust pipeline and infrastructure, and although the individuals in the curation community are clearly skilled in the annotation process, the amount of resources devoted to this important task may be simply insufficient relative to the volume of biomedical data being generated.

Recently, several efforts have been published, which attempt to harness the principle of 'community intelligence' (9–15). In particular, we introduced the Gene Wiki (11), an effort to systematically annotate articles in the online encyclopedia, Wikipedia, for approximately 9000 human genes. Articles were created or amended with content mined from structured gene annotation databases, including Entrez Gene, Ensembl, UniProt and









**Figure 4.** Analysis of page size versus change in page size. The change in size of each Gene Wiki article during the first 6 months of 2009 is plotted as a function of the page size at the end of June 2009. Larger markers indicate more page views, and marker color indicates the Google rank of the page when searching by gene symbol. Labels are shown for a few representative pages.

When examining the usage of statistics, we noticed spikes in the viewing of certain genes, especially those mentioned in the popular press. To explore this observation, we identified the 771 Gene Wiki pages with the most recent variability in monthly page views. Of these, 69 had been searched often enough to have data in Google Trends (<http://www.google.com/trends>), a service that quantifies how many Google searches have been done for a particular term over time relative to the total number of Google searches. The correlation between Gene Wiki page views and Google Trends over time is readily apparent, with 43% of examined pages having significant correlation ( $R > 0.3$ ;  $P < 0.01$ ).

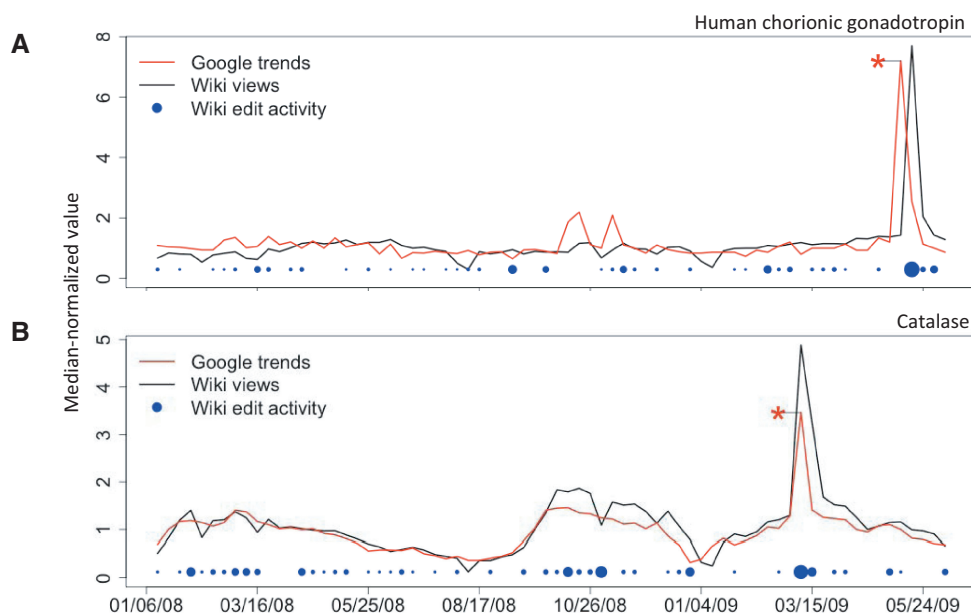
In many cases, the strong relationship between page views and Google Trends was driven by articles in the popular press (Figure 5). For example, the Wikipedia article for human chorionic gonadotropin (HCG) is one of the most frequently viewed articles in the Gene Wiki, presumably for its common usage in pregnancy tests. In May 2009, the Wikipedia article for this gene experienced a sharp spike in views (and edits) when Manny Ramirez was suspended for using HCG as a performance-enhancing drug. Similarly, catalase is frequently viewed article for its relevance to many areas of biology including aging and cancer. However, following a scientific report linking catalase function to premature gray hair in

February 2009 (20), a prominent spike occurred in the viewing and editing of its Gene Wiki entry. Taken in sum, these data show a dynamic relationship between scientific publications, reports on this science in the popular press and usage of the Gene Wiki. These observations also underscore the potential opportunity and effectiveness of using the Gene Wiki for public outreach and scientific education.

## DISCUSSION AND FUTURE DIRECTIONS

With the explosion in biological wikis, it is clear that the community intelligence model resonates with the biology and scientific community (9–15). Despite the enthusiasm in the potential of this model, it is also clear that realizing this potential is not trivial. Many of these biological wikis appear to suffer from a lack of participation. Establishing a critical mass of users and useful content appears to be the most common obstacle in these efforts.

By integrating directly with Wikipedia, establishing critical mass has not been an issue for the Gene Wiki. Clearly, Wikipedia already had a critical mass of users and articles, and the Gene Wiki has been able to effectively leverage those resources as demonstrated by the usage and editing metrics presented above. Moreover, within the last year, the American Society for Cell Biology, the Society



**Figure 5.** Timelines of Wikipedia views and Google Trends information for HCG (A) and catalase (B). Blue dots represent editing events, with bigger size corresponding to more editing events. The asterisks mark major events in the popular press for HCG and catalase, where Manny Ramirez was linked to performance-enhancing drugs and catalase was linked to premature gray hair, respectively.

for Neuroscience and the National Institutes of Health have all held workshops or initiated efforts focused on science articles in Wikipedia. However, the Gene Wiki inherited a completely different set of challenges. First and most notably, Wikipedia allows users to remain completely anonymous, which often leads to fears of inaccuracy and bias. And second, Wikipedia is primarily focused on building unstructured articles (free text, images, diagrams, etc.) with relatively little attention to how contributed knowledge can be structured for downstream analyses in the way that Gene Ontology annotations, for example, can be utilized (21).

We intend to focus on these issues in future developments of the Gene Wiki. Although previous studies have suggested that Wikipedia is of comparable accuracy to traditionally curated works (22), other efforts have been developed to explicitly account for trustworthiness of content based on historical editing patterns of each user (23). Moreover, while we still believe that a completely unstructured Gene Wiki article is useful to the community (similarly to a gene-specific review article), we are also investigating methods to integrate community intelligence with data structure using novel technical solutions [e.g. Semantic MediaWiki (24)] and biomedical ontologies (25).

It is essential to emphasize that community intelligence efforts are not a replacement for traditionally curated gene annotation authorities (16,26–28). In contrast, we believe that community intelligence resources are complementary to existing databases and offer a different set of strengths and weaknesses. Certainly, the data generation model is very different, and users of the Gene Wiki need to recognize that the Gene Wiki, like Wikipedia itself, should be treated differently than the primary literature and expert-curated databases.

Ultimately, we believe that a variety of solutions in the area of community intelligence are worth exploring. Future Gene Wiki development will focus on addressing the challenges described above, and we are also very enthusiastic about complementary efforts as they work to build critical mass and encourage participation. Regardless, the usage metrics presented above demonstrate that the Gene Wiki is relevant right now, certainly to the general public and also to a growing number of scientists. We hope that the scientific community embraces this opportunity both to collaboratively annotate gene function and to directly communicate with the public in science education and outreach.

## ACCESSIBILITY

Wikipedia is freely available for viewing at <http://wikipedia.org>, and the Gene Wiki Portal page can be accessed at [http://en.wikipedia.org/wiki/Portal:Gene\\_Wiki](http://en.wikipedia.org/wiki/Portal:Gene_Wiki). All text is licensed under the Creative Commons Attribution/Share-Alike License 3.0 (Unported).

## SUPPLEMENTARY DATA

Supplementary Data are available at NAR Online.

## ACKNOWLEDGEMENTS

The authors acknowledge Konrad F. Koehler for helpful suggestions and enthusiastic editing, Jeff Janes and Julia Turner for technical assistance, as well as the entire community of Wikipedia editors and the Molecular and Cellular Biology WikiProject (<http://en.wikipedia.org/wiki/WP:MCB>) for contributions and feedback.

## FUNDING

Funding for this work and for the open access charge was provided by the Novartis Research Foundation and the National Institutes of Health [Grant Number 1R01GM083924 to A.S.].

*Conflict of interest statement.* None declared.

## REFERENCES

1. Su, A.I. and Hogenesch, J.B. (2007) Power-law-like distributions in biomedical publications and research funding. *Genome Biol.*, **8**, 404.
2. Blake, J.A., Eppig, J.T., Bult, C.J., Kadin, J.A. and Richardson, J.E. (2006) The Mouse Genome Database (MGD): updates and enhancements. *Nucleic Acids Res.*, **34**, D562–D567.
3. Cherry, J.M., Ball, C., Weng, S., Juvik, G., Schmidt, R., Adler, C., Dunn, B., Dwight, S., Riles, L., Mortimer, R.K. *et al.* (1997) Genetic and physical maps of *Saccharomyces cerevisiae*. *Nature*, **387**, 67–73.
4. Grumbling, G. and Strelets, V. (2006) FlyBase: anatomical data, images and queries. *Nucleic Acids Res.*, **34**, D484–D488.
5. Camon, E., Magrane, M., Barrell, D., Lee, V., Dimmer, E., Maslen, J., Binns, D., Harte, N., Lopez, R. and Apweiler, R. (2004) The Gene Ontology Annotation (GOA) Database: sharing knowledge in Uniprot with Gene Ontology. *Nucleic Acids Res.*, **32**, D262–D266.
6. Ashburner, M., Ball, C.A., Blake, J.A., Botstein, D., Butler, H., Cherry, J.M., Davis, A.P., Dolinski, K., Dwight, S.S., Eppig, J.T. *et al.* (2000) Gene Ontology: tool for the unification of biology. The Gene Ontology Consortium. *Nat. Genet.*, **25**, 25–29.
7. Hill, D.P., Smith, B., McAndrews-Hill, M.S. and Blake, J.A. (2008) Gene Ontology annotations: what they mean and where they come from. *BMC Bioinformatics*, **9**(Suppl. 5), S2.
8. Howe, D., Costanzo, M., Fey, P., Gojobori, T., Hannick, L., Hide, W., Hill, D.P., Kania, R., Schaeffer, M., St Pierre, S. *et al.* (2008) Big data: the future of biocuration. *Nature*, **455**, 47–50.
9. Daub, J., Gardner, P.P., Tate, J., Ramsköld, D., Manske, M., Scott, W.G., Weinberg, Z., Griffiths-Jones, S. and Bateman, A. (2008) The RNA WikiProject: community annotation of RNA families. *RNA*, **14**, 2462–2464.
10. Hoffmann, R. (2008) A wiki for the life sciences where authorship matters. *Nat. Genet.*, **40**, 1047–1051.
11. Huss, J.W., Orozco, C., Goodale, J., Wu, C., Batalov, S., Vickers, T.J., Valafar, F. and Su, A.I. (2008) A gene wiki for community annotation of gene function. *PLoS Biol.*, **6**, e175.
12. Mons, B., Ashburner, M., Chichester, C., van Mulligen, E., Weeber, M., den Dunnen, J., van Ommen, G., Musen, M., Cockerill, M., Hermjakob, H. *et al.* (2008) Calling on a million minds for community annotation in WikiProteins. *Genome Biol.*, **9**, R89.
13. Pico, A.R., Kelder, T., van Iersel, M.P., Hanspers, K., Conklin, B.R. and Evelo, C. (2008) WikiPathways: pathway editing for the people. *PLoS Biol.*, **6**, e184.
14. Stokes, T.H., Torrance, J.T., Li, H. and Wang, M.D. (2008) ArrayWiki: an enabling technology for sharing public microarray data repositories and meta-analyses. *BMC Bioinformatics*, **9**, S18 Suppl. 6.
15. Hodis, E., Prilusky, J., Martz, E., Silman, I., Moulton, J. and Sussman, J.L. (2008) Proteopedia - a scientific 'wiki' bridging the rift between three-dimensional structure and function of biomacromolecules. *Genome Biol.*, **9**, R121.
16. Maglott, D., Ostell, J., Pruitt, K.D. and Tatusova, T. (2007) Entrez Gene: gene-centered information at NCBI. *Nucleic Acids Res.*, **35**, D26–D31.
17. Hughes-Davies, L., Huntsman, D., Ruas, M., Fuks, F., Bye, J., Chin, S., Milner, J., Brown, L.A., Hsu, F., Gilks, B. *et al.* (2003) EMSY links the BRCA2 pathway to sporadic breast and ovarian cancer. *Cell*, **115**, 523–535.
18. Berman, H.M., Westbrook, J., Feng, Z., Gilliland, G., Bhat, T.N., Weissig, H., Shindyalov, I.N. and Bourne, P.E. (2000) The Protein Data Bank. *Nucleic Acids Res.*, **28**, 235–242.
19. Andreeva, A., Howorth, D., Chandonia, J., Brenner, S.E., Hubbard, T.J., Chothia, C. and Murzin, A.G. (2008) Data growth and its impact on the SCOP database: new developments. *Nucleic Acids Res.*, **36**, D419–D425.
20. Wood, J.M., Decker, H., Hartmann, H., Chavan, B., Rokos, H., Spencer, J.D., Hasse, S., Thornton, M.J., Shalhaf, M., Paus, R. *et al.* (2009) Senile hair graying: H<sub>2</sub>O<sub>2</sub>-mediated oxidative stress affects human hair color by blunting methionine sulfoxide repair. *FASEB J.*, **23**, 2065–2075.
21. Subramanian, A., Tamayo, P., Mootha, V.K., Mukherjee, S., Ebert, B.L., Gillette, M.A., Paulovich, A., Pomeroy, S.L., Golub, T.R., Lander, E.S. *et al.* (2005) Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proc. Natl Acad. Sci. USA*, **102**, 15545–15550.
22. Giles, J. (2005) Internet encyclopaedias go head to head. *Nature*, **438**, 900–901.
23. Adler, B.T., Chatterjee, K., de Alfaro, L., Faella, M. and Pye, I. (2008) Assigning trust to Wikipedia content. In WikiSym 08: Proceedings of the International Symposium on Wikis, Porto, Portugal; 150.
24. Krötzsch, M., Vrandečić, D. and Völkel, M. (2006) *The Semantic Web – ISWC 2006*. Springer, Heidelberg, Berlin.
25. Noy, N.F., Shah, N.H., Whetzel, P.L., Dai, B., Dorf, M., Griffith, N., Jonquet, C., Rubin, D.L., Storey, M., Chute, C.G. *et al.* (2009) BioPortal: ontologies and integrated data resources at the click of a mouse. *Nucleic Acids Res.*, **37**, W170–W173.
26. Flicek, P., Aken, B.L., Beal, K., Ballester, B., Caccamo, M., Chen, Y., Clarke, L., Coates, G., Cunningham, F., Cutts, T. *et al.* (2008) Ensembl 2008. *Nucleic Acids Res.*, **36**, D707–D714.
27. Amberger, J., Bocchini, C.A., Scott, A.F. and Hamosh, A. (2009) McKusick's Online Mendelian Inheritance in Man (OMIM). *Nucleic Acids Res.*, **37**, D793–D796.
28. Bairoch, A., Bougueleret, L., Altairac, S., Amendolia, V., Auchincloss, A., Argoud-Puy, G., Axelsen, K., Baratin, D., Blatter, M., Boeckmann, B. *et al.* (2009) The Universal Protein Resource (UniProt) 2009. *Nucleic Acids Res.*, **37**, D169–D174.