

Essential information for synthetic DNA sequences

To the Editor:

Following a discussion by the workgroup for Data Standards in Synthetic Biology, which met in June 2010 during the Second Workshop on Biodesign Automation in Anaheim, California, we wish to highlight a problem relating to the reproducibility of the synthetic biology literature. In particular, we have noted the very small number of articles reporting synthetic gene networks that disclose the complete sequence of all the constructs they describe.

To our knowledge, there are only a few examples where full sequences have been released. In 2005, a patent application¹ disclosed the sequences of the toggle switches published four years earlier in a paper by Gardner *et al.*². The same year, Basu *et al.*³ deposited their construct sequences for programmed pattern formation into GenBank³. Examples of synthetic DNA sequences derived from standardized parts that have been made available in GenBank include the refactored genome of the bacteriophage phage T7 (ref. 4) and a BioBrick-based plasmid⁵. More recently, the full genome sequence of synthetic *Mycoplasma mycoides* JCVI-syn1.0 clone sMmYCp235-1 also has been made available in GenBank (accession no. CP002027)⁶.

In contrast, most publications provide a variety of methods, information and/or partial sequences to explain the constructs used in a paper; for the research community, piecing together the full sequences of constructs is thus laborious, error-prone and sometimes impossible. A paper from your journal provides a recent example; although Kemmer *et al.*⁷ provided admirably detailed Supplementary Information on the construction methods for their plasmids, they failed to provide access to the final sequences. Indeed, the

gaps between key components are almost never reported, presumably because they are not considered crucial to the report. Yet, synthetic biology relies on the premise that synthetic DNA can be engineered with base-level precision.

Missing sequence information in papers hurts reproducibility, limits reuse of past work and incorrectly assumes that we know fully which sequence segments are important. For example, many synthetic biologists are currently realizing that translation initiation rates are dependent on more than the Shine-Dalgarno sequence⁸. Sequences upstream of the

start codon are crucial for translation rates, yet are underreported. Similarly, it has been demonstrated that intron length can affect the dynamics of genetic oscillators⁹. Many more such examples are likely to emerge.

Because full sequence disclosure is critical, we wonder why the common requirement by many journals to provide GenBank entries

for genomes and natural sequences has not been enforced for synthetic DNA and engineered genetic constructs. In an environment where word count is a constant battle, replacing plasmid construction method sections with references to annotated GenBank entries would be a welcome change. We therefore feel that including a completely annotated sequence of the construct would greatly contribute to the development of our discipline. We hope that in the future you will encourage the authors you publish to submit this information to GenBank or other appropriate databases. In the long term, we hope to establish a minimal information guideline around the Minimal Information about a Biomedical or Biological Investigation (MIBBI; http://mibbi.org/index.php/Main_Page) project

and welcome contributions from the greater community.

COMPETING FINANCIAL INTERESTS

The authors declare no competing financial interests.

Jean Peccoud¹, J Christopher Anderson², Deepak Chandran³, Douglas Densmore⁴, Michal Galdzicki⁵, Matthew W Lux¹, Cesar A Rodriguez⁶, Guy-Bart Stan⁷ & Herbert M Sauro³

¹Virginia Bioinformatics Institute, Virginia Tech, Blacksburg, Virginia, USA. ²Department of Bioengineering, QB3: California Institute for Quantitative Biological Research, University of California, Berkeley, California, USA.

³Department of Bioengineering, University of Washington, Seattle, Washington, USA.

⁴Department of Electrical and Computer Engineering, Boston University, Boston, Massachusetts, USA. ⁵Biomedical and Health Informatics, University of Washington, Seattle, Washington, USA. ⁶BIOFAB, Emeryville, California, USA. ⁷Department of Bioengineering and Centre for Synthetic Biology and Innovation, Imperial College London, London, UK.

e-mail: peccoud@vt.edu

- Gardner, T.S. & Collins, J.J. US patent 6,841,376 (2005).
- Gardner, T.S., Cantor, C.R. & Collins, J.J. *Nature* **403**, 339–342 (2000).
- Basu, S., Gerchman, Y., Collins, C.H., Arnold, F.H. & Weiss, R. *Nature* **434**, 1130–1134 (2005).
- Chan, L.Y., Kosuri, S. & Endy, D. *Mol. Syst. Biol.* **1**, 2005.0018 (2005).
- Shetty, R.P., Endy, D. & Knight, T.F. Jr. *J. Biol. Eng.* **2**, 5 (2008).
- Gibson, D.G. *et al. Science* **329**, 52–56 (2010).
- Kemmer, C. *et al. Nat. Biotechnol.* **28**, 355–360 (2010).
- Salis, H.M., Mirsky, E.A. & Voigt, C.A. *Nat. Biotechnol.* **27**, 946–950 (2009).
- Swinburne, I.A., Miguez, D.G., Landgraf, D. & Silver, P.A. *Genes Dev.* **22**, 2342–2346 (2008).

Nature Biotechnology replies:

Kemmer *et al.*¹ have now lodged the sequences of the constructs used in their paper with GenBank HQ644133, HQ644134, HQ644135, HQ644136 and HQ644137. *Nature Biotechnology* and other Nature research journals currently require disclosure only of the sequences of genomes, deep sequencing and short-read data, short stretches of novel

