# Open (and Big) Data – the next challenge

**Beyond dead trees: are publishers the problem or solution?**

## Scott Edmunds

**OASPA Asia, 2nd June 2013**
**@gigascience**

# Harnessing Data-Driven Intelligence

**Enables:**

Using networking power of the internet to tackle problems

Can ask new questions & find hidden patterns & connections

Build on each others efforts quicker & more efficiently

More collaborations across more disciplines

Harness wisdom of the crowds: crowdsourcing, citizen science, crowdfunding

**Enabled by:**

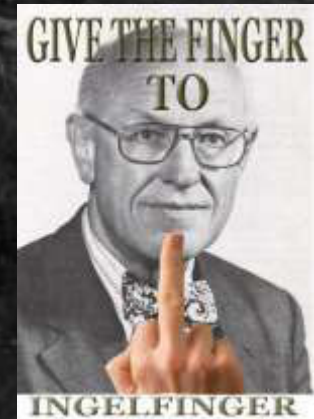Removing silos, <u>open</u> licenses, transparency, immediacy

# Dead trees not fit for purpose



1665                    1812                    1869
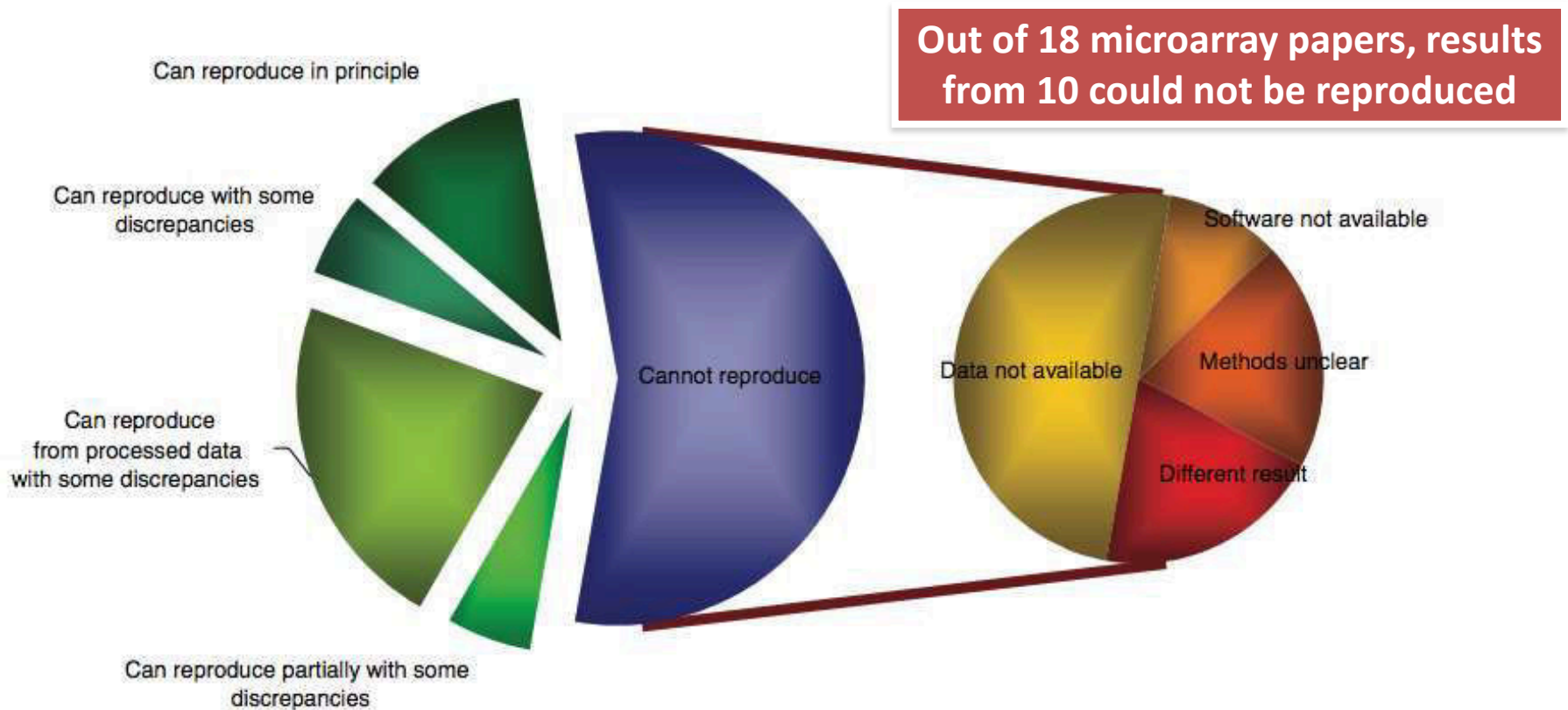
# The problems with publishing

- Scholarly articles are merely advertisement of scholarship . The actual scholarly artefacts, i.e. the data and computational methods, which support the scholarship, remain largely inaccessible --- *Jon B. Buckheit and David L. Donoho, WaveLab and reproducible research, 1995*

- Lack of transparency, lack of credit for anything other than "regular" dead tree publication.

- If there is interest in data, only to monetise & re-silo

- Traditional publishing policies and practices a hindrance

# Things holding us back:

- Disincentives to share or communicate:
  - Ingelfinger*! Embargoes, anti preprint & early data release policies
  - Page/method/citation limits

- Disincentives to remix
  - Open source approaches = plagiarism?

- Disincentives to release more quickly/more granularly
  - "Salami Slicing"

- First 2 years of citation data the only currency
  - "Faddism" v long term use or reproducibility. Publication bias.

# The consequences: growing replication gap

**Out of 18 microarray papers, results from 10 could not be reproduced**

Can reproduce in principle

Can reproduce with some discrepancies

Can reproduce from processed data with some discrepancies

Can reproduce partially with some discrepancies

Cannot reproduce

Software not available

Data not available

Methods unclear

Different result

## Essay

## Why Most Published Research Findings Are False

John P. A. Ioannidis

PLOS | MEDICINE

1.  Ioannidis et al., (2009). Repeatability of published microarray gene expression analyses. *Nature Genetics* 41: 14
2.  Ioannidis JPA (2005) Why Most Published Research Findings Are False. *PLoS Med* 2(8)

# Consequences: increasing number of retractions



## >15X increase in last decade

Strong correlation of "retraction index" with higher impact factor

1. **Science publishing: The trouble with retractions** http://www.nature.com/news/2011/111005/full/478026a.html
2. **Retracted Science and the Retraction Index** ᵛ http://iai.asm.org/content/79/10/3855.abstract?

# Consequences: growing replication gap



RESEARCH ARTICLE

## A Bacterium That Can Grow by Using Arsenic Instead of Phosphorus

**Insufficient methods**

## Retractions On the Rise

A study of the PubMed database found that the number of articles retracted from scientific journals increased substantially between 2000 and 2009.

- Fraud or fabrication 196 total
- Scientific mistake 235 total
- Other 311 total

**More retractions:**
**>15X increase in last decade**
**At current % > by 2045 as many papers published as retracted**

1. Ioannidis et al., 2009. Repeatability of published microarray gene expression analyses. *Nature Genetics* 41
2. Science publishing: The trouble with retractions http://www.nature.com/news/2011/111005/full/478026a.html
3. Bjorn Brembs: Open Access and the looming crisis in science https://theconversation.com/open-access-and-the-looming

The New York Times

# Global perceptions of Chinese Research
## Million RMB rewards for high IF publications = ?



**NewScientist**

*"Faked research is endemic in China"*

Projected growth in citations in scientific literature

*Source: Royal Society*

Fraud fighter: 'Faked research is endemic in China'
Shi-min Fang tells us how risking his life and libel writs to expose scientific misconduct in his native China has just won him the inaugural Maddox prize

## China's Publication Bazaar
A Science investigation has uncovered a smorgasbord of questionable practices including paying for author's slots on papers written by other scientists and buying papers from online brokers

**nature**
International weekly journal of science

**475**, 267 (2011)

ALTHOUGH CHINA RANKS SECOND IN TERMS OF **PUBLICATION** OUTPUT, IT RANKS ONLY NINTH IN **CITATION** NUMBERS.

## Focus on quality, not just quantity
China publishes huge amounts of scientific research. Now it must make more of it worth reading, says **Changhui Peng**.

FRAUD
## Brawl in Beijing
Critics of Chinese researchers targeted in physical attacks.

*New Scientist*, 17th Nov 2012: http://www.newscientist.com/article/mg21628910.300-fraud-fighter-faked-research-is-endemic-in-china.html
*Nature*, 29th September 2010: http://www.nature.com/news/2010/100929/full/467511a.html
*Science*, 29th November 2013: http://www.sciencemag.org/content/342/6162/1035.full
*Nature* 20th July 2011: http://www.nature.com/news/2011/110720/full/475267a.html

# Global perceptions of Chinese Research

## Million RMB rewards for high IF publications = ?

"Wide distribution of information is key to scientific progress, yet traditionally, Chinese scientists have not systematically released data or research findings, even after publication."

"There have been widespread complaints from scientists inside and outside China about this lack of transparency."

"Usually incomplete and unsystematic, [what little supporting data released] are of little value to researchers and there is evidence that this drives down a paper's citation numbers."

*"Faked research is endemic in China"*

**nature** **475**, 267 (2011)
International weekly journal of science

## Focus on quality, not just quantity

*China publishes huge amounts of scientific research. Now it must mal more of it worth reading, says* **Changhui Peng.**

ALTHOUGH CHINA RANKS SECOND IN TERMS OF **PUBLICATION** OUTPUT, IT RANKS ONLY NINTH IN **CITATION** NUMBERS.

**FRAUD**
**Brawl in Beijing**
*Critics of Chinese researchers targeted in physical attacks.*

Projected growth in citations in scientific literatur
2013 China to overtake US

*New Scientist*, 17th Nov 2012: http://www.newscientist.com/article/mg21628910.300-fraud-fighter-faked-research-is-endemic-in-china.html
*Nature*, 29th September 2010: http://www.nature.com/news/2010/100929/full/467511a.html
*Science*, 29th November 2013: http://www.sciencemag.org/content/342/6162/1035.full
*Nature* 20th July 2011: http://www.nature.com/news/2011/110720/full/475267a.html

# Issues not just in China…

**Need:**
**…to publish protocols BEFORE analysis**
**…better access to supporting data**
**…more transparent & accountable review**

protocol exchange

Home    Browse    Share protocol    Lab groups    About    Contact

PROTOCOL EXCHANGE | COMMUNITY CONTRIBUTED

Essential technical tips for STAP cell conversion culture from somatic cells

Haruko Obokata, Yoshiki Sasai & Hitoshi Niwa

STAP Group RIKEN CDB

Protocol Exchange (2014) | doi:10.1038/protex.2014.008
Published online 5 March 2014

**…to publish replication studies**

## Knoepfler Lab Stem Cell Blog
Building stem cell bridges

## Nature Rejects Publication of Paper Reporting that STAP Does Not Work

Posted on **March 24, 2014**

F1000Research

RESEARCH ARTICLE

CrossMark

Transient acid treatment cannot induce neonatal somatic cells to become pluripotent stem cells [v1; ref status: indexed, http://f1000r.es/3dq]

Mei Kuen Tang[1], Lok Man Lo[1], Wen Ting Shi[1], Yao Yao[1], Henry Siu Sum Lee[2], Kenneth Ka Ho Lee[1]

[1]Key Laboratory for Regeneration Medicine, School of Biomedical Sciences, Chinese University of Hong Kong, Shatin, Hong Kong
[2]Faculty of Life Sciences, University of Manchester, Manchester, M13 9PL, UK

# New incentives/credit

## Credit where credit is overdue:

"One option would be to provide researchers who release data to public repositories with a means of accreditation."

"An ability to search the literature for all online papers that used a particular data set would enable appropriate attribution for those who share. "

*Nature Biotechnology* **27**, 579 (2009)



- **Data**
- **Software**
- **Review**
- **Re-use…**

**= Credit**

# GigaSolution: deconstructing the paper

Combines and integrates:



(GIGA)$^n$ SCIENCE

(GIGA)$^n$ DB

(GIGA)$^n$ Galaxy by CBIIT

Open-access journal

Data Publishing Platform

Data Analysis Platform

OPEN ACCESS

OD

open source

**Utilizes big-data infrastructure and expertise from:**



华大基因 BGI

www.gigadb.org
www.gigasciencejournal.com

**BioMed** Central
The Open Access Publisher

# Rewarding open data

# $(GIGA)^n_{DB}$ Submission Workflow

Submitter logs in to GigaDB website and uploads Excel submission

Fail – submitter is provided error report

Curator contacts submitter with DOI citation and to arrange file transfer (and resolve any other questions/issues).

Curator Review

Excel submission file

Validation checks

DOI assigned

Files

Pass – dataset is uploaded to GigaDB.

GigaDB

Submitter provides files by ftp or Aspera

XML is generated and registered with DataCite

DataCite XML file

Curator makes dataset public (can be set as future date if required)

Public GigaDB dataset
DOI 10.5524/100003
Genomic data from the crab-eating macaque/cynomolgus monkey (*Macaca fascicularis*) (2011)

See: http://database.oxfordjournals.org/content/2014/bau018.abstract

- 10-100x faster download than FTP
- Provide curation & integration with other DBs

# Beneficiaries of this open data?

# Beneficiaries of this open data?

## Rice 3K project: 3,000 rice genomes, 13.4TB public data

# New Article types v Species Description <2012

# Collaborations with Pensoft & PLOS
# Cyber-centipedes & virtual worms

Taxonomic paper

*Eupolybothrus cavernicolus* Komerički & Stoev sp. n. (Chilopoda: Lithobiomorpha: Lithobiidae): the first eukaryotic species description combining transcriptomic, DNA barcoding and micro-CT imaging data

Pavel Stoev[1,2], Ana Komerički[3], Nesrine Akkari[4], Shanlin Liu[5], Xin Zhou[5], Alexander M. Weigand[6], Jeroen Hostens[7], Christopher I. Hunter[8], Scott C. Edmunds[8], David Porco[9], Marzio Zapparoli[10], Teodor Georgiev[1], Daniel Mietchen[11], David Roberts[12], Sarah Faulwetter[13], Vincent Smith[14], Lyubomir Penev[1,15]

EDITORIAL — Open Access

## Biodiversity research in the "big data" era: *GigaScience* and Pensoft work together to publish the most data-rich species description

Scott C Edmunds[1]*, Chris I Hunter[1], Vincent Smith[2], Pavel Stoev[3,4] and Lyubomir Penev[4,5]

### Abstract
With the publication of the first eukaryotic species description, combining transcriptomic, DNA barcoding, and micro-CT imaging data, GigaScience and Pensoft demonstrate how classical taxonomic description of a new species can be enhanced by applying new generation molecular methods, and novel computing and imaging technologies. This 'holistic' approach in taxonomic description of a new species of cave-dwelling centipede is published in the Biodiversity Data Journal (BDJ), with coordinated data release in the GigaScience GigaDB database.

## *Sine Systemate Chaos?* A Versatile Tool for Earthworm Taxonomy: Non-Destructive Imaging of Freshly Fixed and Museum Specimens Using Micro-Computed Tomography

Rosa Fernández[1]*, Sebastian Kvist[1], Jennifer Lenihan[1], Gonzalo Giribet[1], Alexander Ziegler[2]

1 Museum of Comparative Zoology, Department of Organismic and Evolutionary Biology, Harvard University, Cambridge, Massachusetts, United States of America, 2 Ziegler Biosolutions, Waldshut-Tiengen, Germany

### Abstract
In spite of the high relevance of lumbricid earthworms ('Oligochaeta': Lumbricidae) for soil structure and functioning, the taxonomy of this group of terrestrial invertebrates remains in a quasi-chaotic state. Earthworm taxonomy traditionally relies on the interpretation of external and internal morphological characters, but the acquisition of these data is often hampered by tedious dissections or restricted access to valuable and rare museum specimens. The present state of affairs, in conjunction with the difficulty of establishing primary homologies for multiple morphological features, has led to an almost unrivaled instability in the taxonomy and systematics of certain earthworm groups, including Lumbricidae. As a potential remedy, we apply for the first time a non-destructive imaging technique to lumbricids and explore the future application of this approach to earthworm taxonomy. High-resolution micro-computed tomography (μCT) scanning of freshly fixed and museum specimens was carried out using two cosmopolitan species, *Aporrectodea caliginosa* and *A. trapezoides*. By combining two-dimensional and three-dimensional dataset visualization techniques, we demonstrate that the morphological features commonly used in earthworm taxonomy can now be analyzed without the need for dissection, whether freshly fixed or museum specimens collected more than 60 years ago are studied. Our analyses show that μCT in combination with soft tissue staining can be successfully applied to lumbricid earthworms. An extension of the approach to

DATA NOTE — Open Access

## A dataset comprising four micro-computed tomography scans of freshly fixed and museum earthworm specimens

Jennifer Lenihan[1], Sebastian Kvist[1], Rosa Fernández[1], Gonzalo Giribet[1] and Alexander Ziegler[2]*

### Abstract
Background: Although molecular tools are increasingly employed to decipher invertebrate systematics, earthworm (Annelida: Clitellata: 'Oligochaeta') taxonomy is still largely based on conventional dissection, resulting in data that are mostly unsuitable for dissemination through online databases. In order to evaluate if micro-computed tomography (μCT) in combination with soft tissue staining techniques could be used to expand the existing set of tools available for studying internal and external structures of earthworms, μCT scans of freshly fixed and museum specimens were gathered.

# New & more transparent peer-review: open review

**EDITORIAL**      **Open Access**

# Peering into peer-review at *GigaScience*

Scott C Edmunds

**Abstract**

Fostering and promoting more open and transparent science is one of the goals of *GigaScience*. One of the ways we have been doing this is by throwing light on the peer-review process and carrying out open peer-review as standard. In this editorial, we provide our rationale for undertaking this policy, give examples of our positive experiences to date, and encourage others to open up the normally opaque publication process.

**BioData Mining**

**BMC Series Medical Journals**

**BMJ**

# **Reward open & transparent review**



End reviewer 3 Downfall parody videos, now!

# New & more transparent peer-review: pre-prints

# Reward open & transparent review

## Real-time open-review = paper in arXiv + blogged reviews



## THE ASSEMBLATHON

### Feedback and analysis of the Assemblathon 2 pre-print

There has already been some discussion of the pre-print of the Assemblathon 2 manuscript. Although a pre-print is not the same thing as a peer-reviewed, accepted paper — I don't want us to get too ahead of ourselves! — I thought it useful to start collecting together some of the online commentaries:

- Homologus blog post 1: highlights a few conclusions from the paper
- Homologus blog post 2: delves into the results, and attempts to estimate some of the costs of genome assembly. Assemblathon co-author Sébastien Boisvert adds some useful comments.
- Haldane's Sieve post: an invited blog post by lead author Keith Bradnam, that summarizes what the Assemblathons are all about by way of a pizza-themed analogy
- Reevaluating Assembly Evaluations with Feature Response Curves: GAGE and Assemblathons: this is not a blog post, but a recently published paper that evaluates some of the Assemblathon 2 data
- Thoughts on the Assemblathon 2 paper: by C. Titus Brown (a reviewer of the manuscript)
- Homologus blog post 3: reactions to the previous post by C. Titus Brown
- Assemblathon 2 review, round 1, parts thereof: a concise version of C. Titus Brown's formal manuscript review (minus the specific suggestions)
- On assembly uncertainty (inspired by the Assemblathon 2 debate): blog post by Lex Nederbragt in response to post by C. Titus Brown

www.gigasciencejournal.com/content/2/1/10          http://tmblr.co/ZzXdssfOMJfy

# Reward open & transparent review

## Real-time open-review = paper in arXiv + blogged reviews

# Readers are interested in open review

**SOAPdenovo2**
Large-size genome de-novo assembler

**THE ASSEMBLATHON**



Next step to link to ORCID

casrai
Connecting Research

# Reward better handling of metadata…

Novel tools/formats for data interoperability/handling.



**isa**infrastructure



## nature genetics

nature.com ▸ journal home ▸ archive ▸ issue ▸ commentary ▸ full text

NATURE GENETICS | COMMENTARY   OPEN

### Toward interoperable bioscience data

Susanna-Assunta Sansone, Philippe Rocca-Serra, Dawn Field, Eamonn Maguire, Chris Taylor, Oliver Hofmann, Hong Fang, Steffen Neumann, Weida Tong, Linda Amaral-Zettler, Kimberly Begley, Tim Booth, Lydie Bougueleret, Gully Burns, Brad Chapman, Tim Clark, Lee-Ann Coleman, Jay Copeland, Sudeshna Das, Antoine de Daruvar, Paula de Matos, Ian Dix, Scott Edmunds, Chris T Evelo, Mark J Forster + et al.

Affiliations | Corresponding author

Nature Genetics **44**, 121–126 (2012) | doi:10.1038/ng.1054
Published online 27 January 2012





**isa**commons
isacommons.org

**isa**creator
configurator

npg
nature publishing group
SCIENTIFIC DATA

SCIENTIFIC DATA

# Rewarding and aiding reproducibility



**OMERO: providing access to imaging data…**

# Rewarding and aiding reproducibility

## Implement workflows in a community-accepted format



Open source

Over 36,000 main Galaxy server users

Over 1,000 papers citing Galaxy use

Over 55 Galaxy servers deployed

http://galaxyproject.org

galaxy.cbiit.cuhk.edu.hk

Visualizations & DOIs for workflows

GIGA*n*
SCIENCE

# SOAPdenovo2: an empirically improved memory-efficient short-read *de novo* assembler

Ruibang Luo[1,2], Binghang Liu[1,2], Yinlong Xie[1,2,3], Zhenyu Li[1,2], Weihua Huang[1], Jianying Yuan[1], Guangzhu He[3], Yanxiang Chen[1], Qi Pan[1], Yunjie Liu[1], Jingbo Tang[3], Gengxiong Wu[1], Hao Zhang[1], Yujian Shi[1], Yong Liu[1], Chang Yu[1], Bo Wang[1], Yao Lu[1], Changlei Han[1], David W Cheung[2], Siu-Ming Yiu[2], Shaoliang Peng[4], Zhu Xiaoqian[4], Guangming Liu[4], Xiangke Liao[4], Yingrui Li[1,2], Huanming Yang[1], Jian Wang[1], Tak-Wah Lam[2*] and Jun Wang[1*]

## Abstract

**Background:** There is a rapidly increasing amount of *de novo* genome assembly using next-generation sequencing (NGS) short reads; however, several big challenges remain to be overcome in order for this to be efficient and accurate. SOAPdenovo has been successfully applied to assemble many published genomes, but it still needs improvement in continuity, accuracy and coverage, especially in repeat regions.

**Findings:** To overcome these challenges, we have developed its successor, SOAPdenovo2, which has the advantage of a new algorithm design that reduces memory consumption in graph construction, resolves more repeat regions in contig assembly, increases coverage and length in scaffold construction, improves gap closing, and optimizes for large genome.

**Conclusions:** Benchmark using the Assemblathon1 and GAGE datasets showed that SOAPdenovo2 greatly surpasses its predecessor SOAPdenovo and is competitive to other assemblers on both assembly length and accuracy. We also provide an updated assembly version of the 2008 Asian (YH) genome using SOAPdenovo2. Here, the contig and scaffold N50 of the YH genome were ~20.9 kbp and ~22 Mbp, respectively, which is 3-fold and 50-fold longer than the first published version. The genome coverage increased from 81.16% to 93.91%, and memory consumption was ~2/3 lower during the point of largest memory consumption.

**Keywords:** Genome, Assembly, Contig, Scaffold, Error correction, Gap-filling

# How are we supporting data reproducibility?

**Open-Paper**

$(GIGA)^n_{DB}$

**Data sets**

**Open-Data**

DOI:10.5524/100038

78GB CC0 data

*Linked to*

*DOI*

(GIGA)$^n$ SCIENCE

DOI:10.1186/2047-217X-1-18

**Highly accessed** >23,000 accesses

*DOI*

*Linked to*

$(GIGA)^n$ Galaxy

by CBIIT

**Analyses**

**Open-Pipelines**

**Open-Workflows**

DOI:10.5524/100044

**Open-Review**

7 reviewers tested data in ftp server & named reports published

Enabled code to being picked apart by bloggers in wiki
http://homolog.us/wiki/index.php?title=SOAPdenovo2

**Open-Code**

**source forge** open source

>20,000 downloads

Code in sourceforge under GPLv3:
http://soapdenovo2.sourceforge.net/

# Reward open & transparent review

## 7 referees downloaded & tested data, then signed reports



Pre-publication history

**Highly accessed** · **Open Access**

### SOAPdenovo2: an empirically improved memory-efficient short-read *de novo* assembler

Ruibang Luo[†], Binghang Liu[†], Yinlong Xie[†], Zhenyu Li, Weihua Huang, Jianying Yuan, Guangzhu He, Yanxiang Chen, Qi Pan, Yunjie Liu, Jingbo Tang, Gengxiong Wu, Hao Zhang, Yujian Shi, Yong Liu, Chang Yu, Bo Wang, Yao Lu, Changlei Han, David W Cheung, Siu-Ming Yiu, Shaoliang Peng, Zhu Xiaoqian, Guangming Liu, Xiangke Liao, Yingrui Li, Huanming Yang, Jian Wang, Tak-Wah Lam[*] and Jun Wang[*]

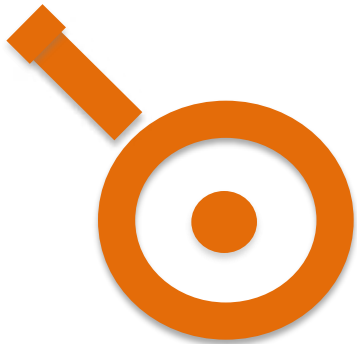* Corresponding authors: Tak-Wah Lam twlam@cs.hku.hk - Jun Wang wangj@genomics.org.cn

† Equal contributors

*GigaScience* 2012, **1**:18    doi:10.1186/2047-217X-1-18

## Pre-publication versions of this article and reviewers' reports

| | | | |
|---|---|---|---|
| Original Submission - Version 1 | Manuscript | | 24 Jul 2012 |
| Reviewer's Report | Alexander J. Nederbragt | | 15 Aug 2012 |
| Reviewer's Report | Aleksey Zimin | | 22 Aug 2012 |
| Reviewer's Report | Mario Caccamo | | 28 Aug 2012 |
| Resubmission - Version 2 | Manuscript | Author's comment | 04 Nov 2012 |
| Reviewer's Report | Alexander J. Nederbragt | | 30 Nov 2012 |
| Reviewer's Report | Aleksey Zimin | | 03 Dec 2012 |
| Resubmission - Version 3 | Manuscript | Author's comment | 04 Dec 2012 |
| Resubmission - Version 6 | Manuscript | | 10 Dec 2012 |
| Editorial acceptance | | | 10 Dec 2012 |

# Reward open & transparent review

**Post publication: bloggers pull apart code/reviews in blogs + wiki:**

## Main Program

### main.c (557 lines)

- Main program. It processes input options, and invokes various other functions (call_pregraph, call_heavygraph, call_align, call_map2contig, and call_scaffold) to assemble a genome.

## Step 1. (constructing pregraph)

### pregraph.c (229 lines )

- Constructs pregraph. From BGI's description -

```
    The main function for pregraph step. its processes are as below:
    1. Builds the kmer hash sets and remove the low coverage kmers.
    2. Removes the tips which length are no greater than 2*K.
    3. Builds edges by combining linear kmers.
    4. Maps the reads back to edges and build preArcs (the connection between edges).
```
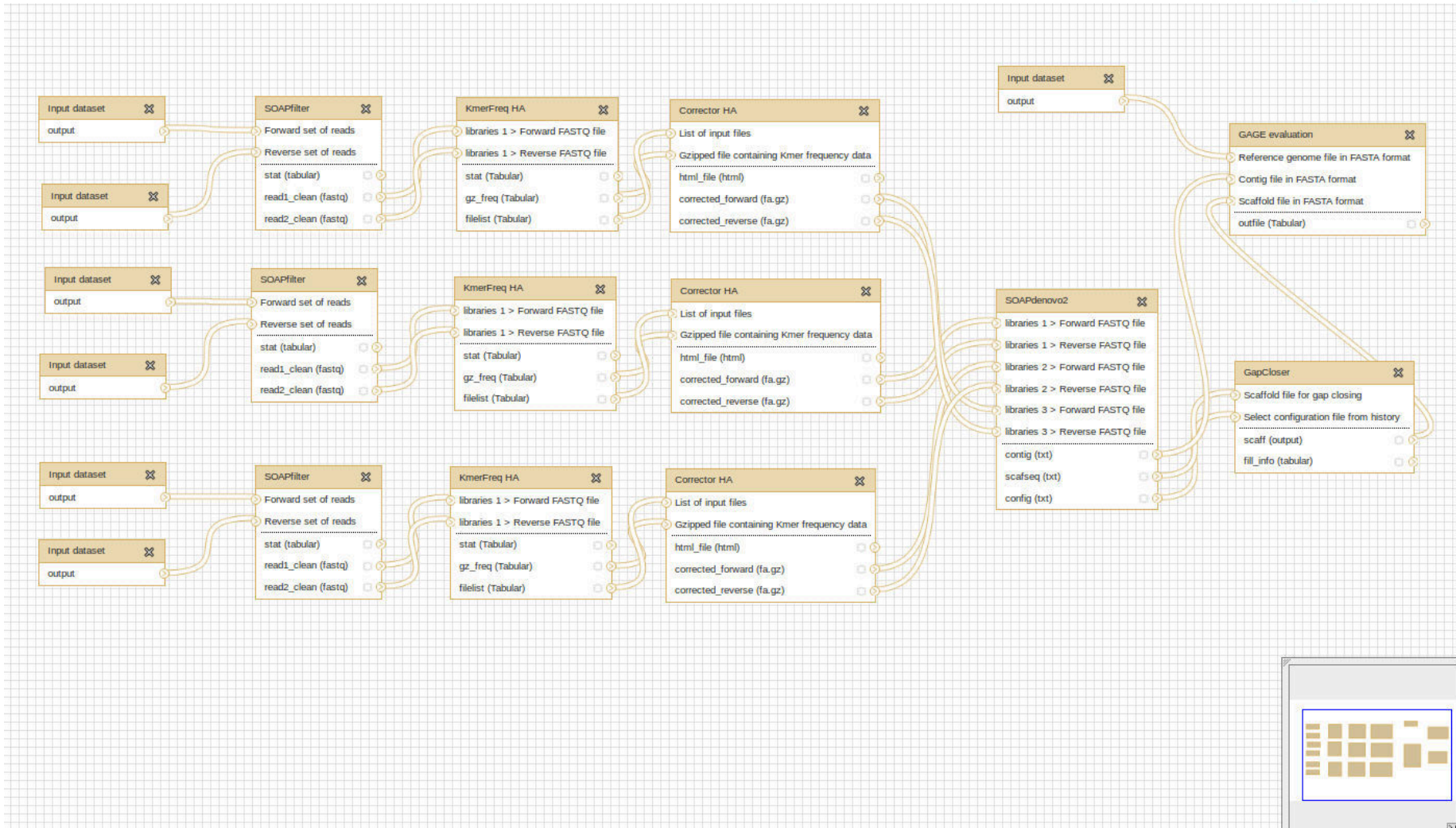
Related:

### cutTipPreGraph.c (639 lines )

### output_pregraph.c (112 lines )

## Step 2. (building contigs)

SOAPdenov2 wiki: http://homolog.us/wiki1/index.php?title=SOAPdenovo2
Homologus blogs: http://www.homolog.us/blogs/category/soapdenovo/

# SOAPdenovo2 workflows implemented in (GIGA)$^n$ Galaxy by CBIIT



galaxy.cbiit.cuhk.edu.hk

# SOAPdenovo2 workflows implemented in (GIGA)ⁿ Galaxy by CBIIT



Implemented entire workflow in our Galaxy server, inc.:

- 3 pre-processing steps

- 4 SOAPdenovo modules

- 1 post processing steps

- Evaluation and visualization tools

Also will be available to download by >36K Galaxy users in Galaxy Tool Shed

galaxy.cbiit.cuhk.edu.hk

SOAPdenovo2 *S. aureus* pipeline

**Table 2 Assemblies of *S. aureus* and *R. sphaeroides***

| Species | Version | Contigs | | | | Scaffolds | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | Number | N50 (kb) | Errors | N50 corrected(kb) | Number | N50 (kb) | Errors | N50 corrected (kb) |
| S. aureus | SOAPdenovo1 | 79 | 148.6 | 156 | 23 | 49 | 342 | 0 | 342 |
| | SOAPdenovo2 | 80 | 98.6 | 25 | 71.5 | 38 | 1,086 | 2 | 1,078 |
| | ALLPATHS-LG* | 37 | 149.7 | 13 | 117.6 | 10 | 1,477 | 1 | 1,093 |
| R. sphaeroides | SOAPdenovo1 | 2,242 | 3.5 | 392 | 2.8 | 956 | 105 | 18 | 70 |
| | SOAPdenovo2 | 721 | 18 | 106 | 14.1 | 333 | 2,549 | 4 | 2,540 |
| | ALLPATHS-LG* | 190 | 41.9 | 31 | 36.7 | 32 | 3,191 | 0 | 3,310 |

All datasets were downloaded from http://gage.cbcb.umd.edu/data/

增值機

Add Value Machine

# Taking a microscope to peer review

Essay

## Why Most Published Research Findings Are False

John P. A. Ioannidis

PLOS | MEDICINE

# The SOAPdenovo2 Case study
## Subject to and test with 3 models:



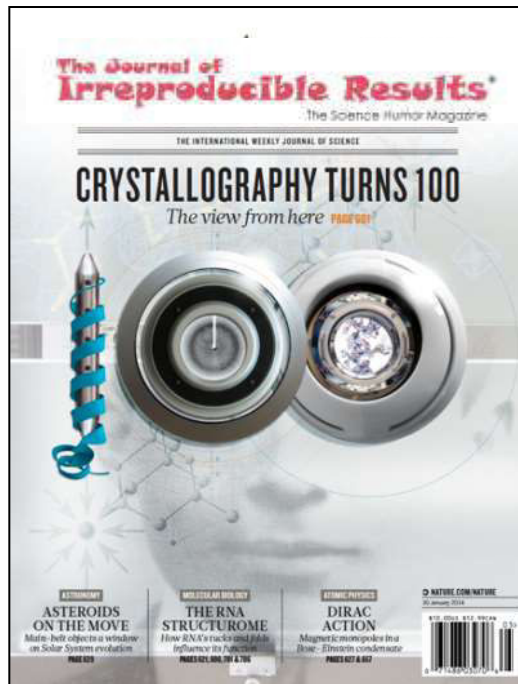| Types of resources in an RO | Models to describe each resource type |
|---|---|
| Data | ISA-TAB/ISA2OWL |
| Method/Experimental protocol | Wfdesc/ISA-TAB/ISA2OWL |
| Findings | Nanopublication |

# Lessons learned:

- Most published research findings are false. Or at least have errors.

- On a semantic level (via nanopublications) discovered 4 minor errors in text (interpretation not data)

- Is possible to push button(s) & recreate a result from a paper

- Reproducibility is COSTLY. How much are you willing to spend?

- Much easier to do this before rather than after publication

# "Regular" Journal
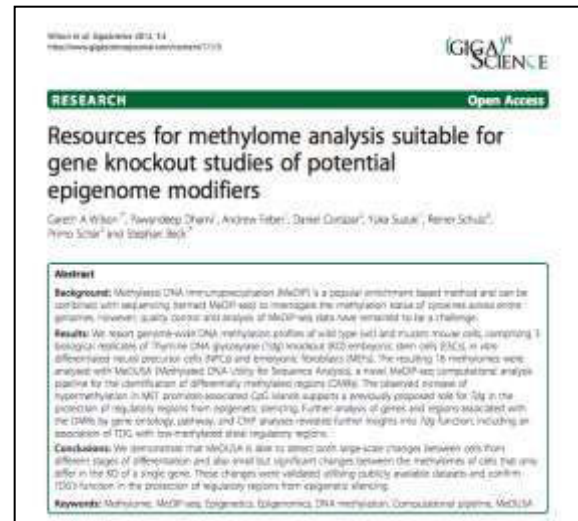
# "Conscientious" Online Journal

# "Deconstructed" Journal
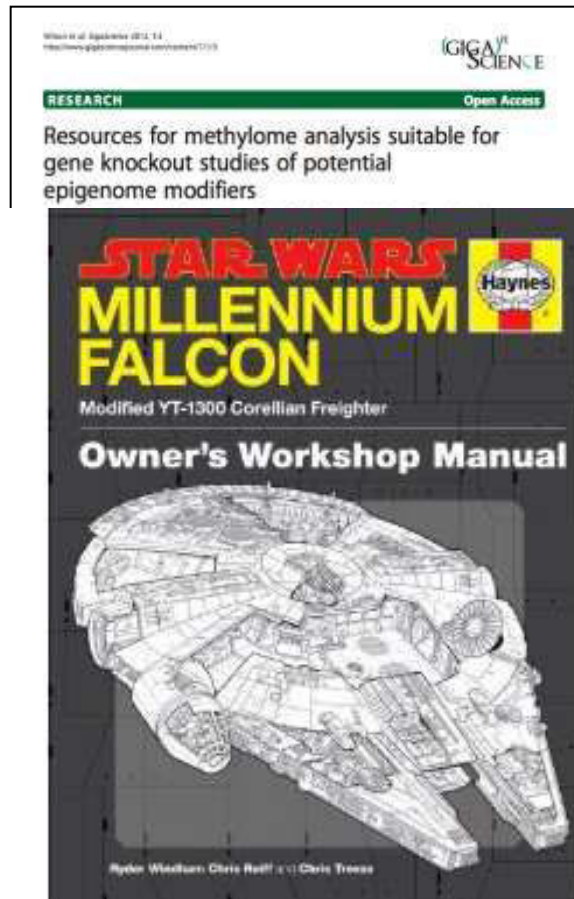
# "Regular" Journal

# "Conscientious" Online Journal

# "Deconstructed" Journal

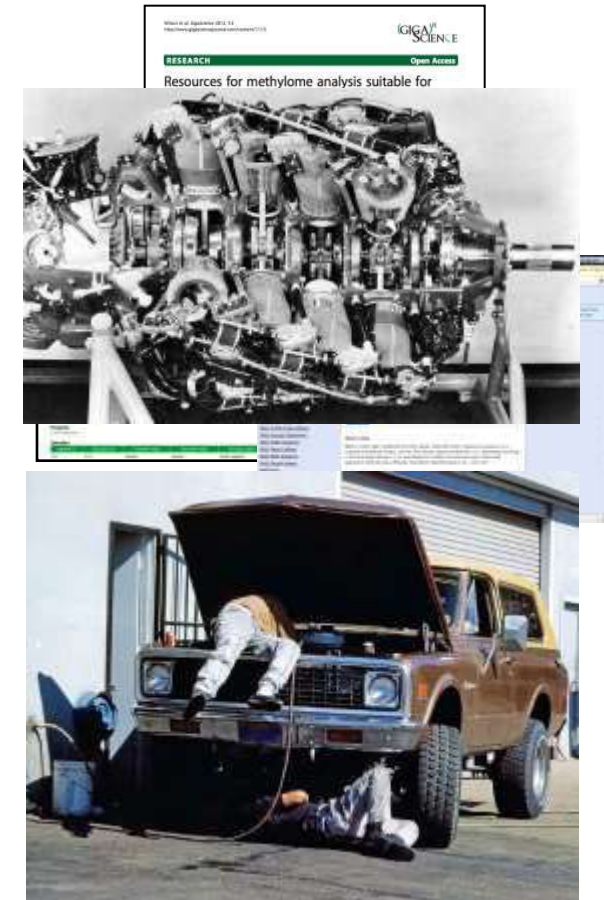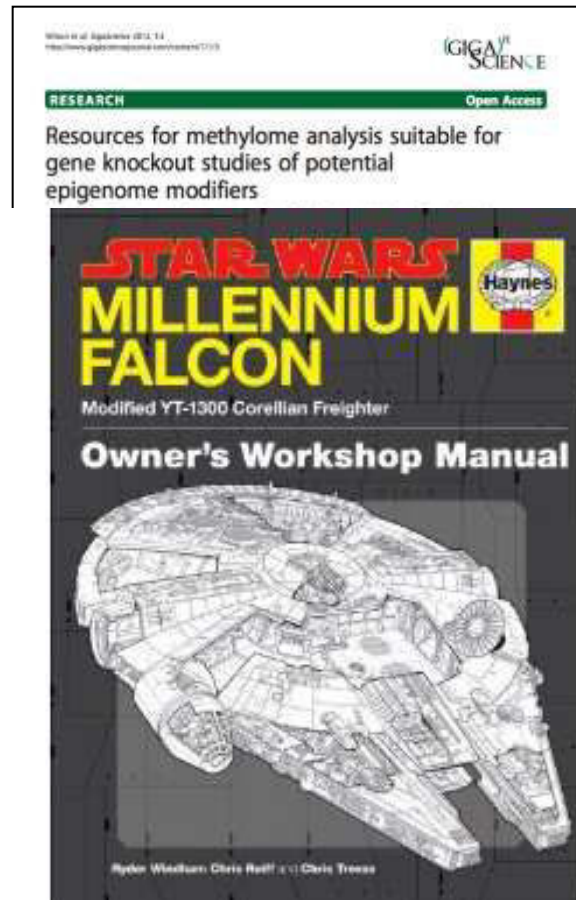# "Regular" Journal

# "Conscientious" Online Journal

# "Deconstructed" Journal

# "Regular" Journal

# "Conscientious" Online Journal

# "Deconstructed" Journal

# Thanks to:

**GIGA$^n$SCIENCE team:**

**Peter Li**
Huayan Gao
Chris Hunter
Jesse Si Zhe
Nicole Nogoy
Laurie Goodman
Amye Kenall (BMC)

**Our collaborators:**

Ruibang Luo (BGI/HKU)
Shaoguang Liang (BGI-SZ)
Tin-Lap Lee (CUHK)
Qiong Luo (HKUST)
Senghong Wang (HKUST)
Yan Zhou (HKUST)

**Case study:**

Marco Roos (LUMC)
Mark Thompson (LUMC)
Jun Zhao (Lancaster)
Susanna Sansone (Oxford)
Philippe Rocca-Serra (Oxford)
Alejandra Gonzalez-Beltran (Oxford)

@gigascience
**facebook.com/GigaScience**
**blogs.biomedcentral.com/gigablog/**

www.gigadb.org
galaxy.cbiit.cuhk.edu.hk
www.gigasciencejournal.com

GIGA$^n$SCIENCE

**BioMed** Central
The Open Access Publisher