

Bio-Mirror project for public bio-data distribution

Don Gilbert^{1*}, Yoshihiro Ugawa², Markus Buchhorn³, Tan Tin Wee⁴, Akira Mizushima⁵,
Hyunchul Kim⁶, Kilnam Chon⁶, Seyeon Weon⁷, Juncai Ma⁸, Yoshihiro Ichiyanagi⁸, Der-Ming
Liou⁹, Somnuk Keretho¹⁰, Suhaimi Napis¹¹

¹ Department of Biology, Indiana University, Bloomington, IN 47405 USA

² Miyagi University of Education, Japan

³ ANU Internet Futures Project, Australian National University, Australia

⁴ BioInformatics Centre, National University of Singapore, Singapore

⁵ National Agricultural Research Organization, Japan

⁶ Korea Advanced Institute of Science and Technology, Daejeon, South Korea

⁷ Bioinformatics Research Laboratories, Co., Taejon, South Korea

⁸ Institute of Microbiology, Chinese Academy of Sciences, Beijing, China

⁹ Bioinformatics Research Center, National Yang-Ming University, Taiwan

¹⁰ Kasetsart University, Thailand

¹¹ Universiti Putra Malaysia, Malaysia

* To whom correspondence should be addressed

Running head: Bio-Mirror data distribution

ABSTRACT

Summary: Timely worldwide distribution of biosequence and bioinformatics data depends on high performance networking and advances in Internet transport methods. The Bio-Mirror project focuses on providing up-to-date distribution of this rapidly growing and changing data. It offers FTP, Web and Rsync access to many high-volume databanks from several sites around the world. Experiments with data grids and other methods offer future improvements in biology data distribution.

Availability: <http://www.bio-mirror.net/>

Contact: bio-mirror@apan.net

INTRODUCTION

The Bio-Mirror project for distribution of public biosequence and bioinformatics data is a collaborative service of several worldwide bioinformatics centers. As multi-gigabyte public bioinformatics databanks grow and change daily, access to them is hampered by limits on Internet bandwidth. The Bio-Mirror project addresses this problem with rapid redistribution from several sites. Such mirrors reduce the burden on source providers, and mitigate Internet outages and slow distant connections.

The Bio-Mirror project was formed in 1998 as a collaboration of Asian-Pacific bioinformatics centers (APBionet) and IUBio Archive at Indiana University. It distributes databanks using Internet2 connections between continents. Many APBionet members are part of new and growing bioinformatics centers, where Bio-Mirror sites provide a short path to current data from US and European sources. Project sites serve data to a range of educational, government and industry bioinformatics groups. We also investigate new technologies in science data grid informatics to improve data distribution.

CONTENTS AND METHODS

Contents of the Bio-Mirror databank set include core databanks from the three collaborating DNA databanks GenBank, EMBL and DDBJ. SwissProt, TrEMBL, PIR, PDB, Pfam and InterPro protein databanks are included, as well as BLAST and RefSeq mixed databanks. Ensembl, NCBI Genomes, LocusLink, euGenes, Unigene, Gene Ontology, and related genome and other data are included. The current collection exceeds 140 Gigabytes (compressed, or about 500 GB uncompressed). Approximately 10% of these change daily, and most are updated in the course of 3 months. Complete databank file sets are mirrored from sources in the US, Europe, and Asia-Pacific. The European Bioinformatics Institute (EBI) and National Center for Biotechnology Information (NCBI) centers now include Bio-Mirror sites in their access information.

The Bio-Mirror project is provided as a public service by member centers with support of several organizations. Participants include APBionet and bioinformatics centers in Japan, Australia, Singapore, China, Korea, Malaysia, Taiwan, Thailand and the USA (Table 1). High performance network infrastructure and collaborative help have been essential to this project, including Trans-Pacific network (TransPAC) and Asia-Pacific Advanced Network (APAN).

Table 1. Bio-Mirror sites.

Country	Host	Web site	Bulk data access
Australia	Australian National University	--	rsync: and ftp://bio-mirror.-au.apan.net/biomirror/
China	Institute of Microbiology, Chinese Academy of Sciences	http://bio-mirror.cn.apan.net/	ftp://bio-mirror.cn.apan.net/pub/biomirror/
Japan	Computer Center for AFFRC	http://bio-mirror.jp.apan.net/	ftp://bio-mirror.jp.apan.net/pub/biomirror/
Korea	Korea Advanced Institute of Science and Technology	http://bio-mirror.kr.apan.net/	ftp://bio-mirror.kr.apan.net/pub/biomirror/
Malaysia	Universiti Putra Malaysia	http://ingene2.upm.edu.my/	ftp://ingene2.upm.edu.my/
Singapore	National University of Singapore	http://bio-mirror.sg.apan.net/	ftp://bio-mirror.sg.apan.net/biomirrors/
Taiwan	National Yang-Ming University	http://bio-mirror.ym.edu.tw/	ftp://bio-mirror.ym.edu.tw/biomirror/
Thailand	Kasetsart University	http://bio-mirror.ku.ac.th/	http://bio-mirror.ku.ac.th/biomirror/
USA	IUBio Archive, Indiana University	http://www.bio-mirror.net/	rsync: and ftp://bio-mirror.net/biomirror/

Project sites currently serve many Terabytes per month to thousands of bioinformatics centers and labs. Bulk distribution has risen from 100 Gigabytes/month in year 2001 to 3 Terabytes/month in 2003 at the US node (Figure 1). File transfer (FTP) provides the best access, as FTP servers have been tuned for large file transfer. The *rsync* protocol (Tridgell, 2002) is supported at some Bio-Mirror servers as an efficient alternative. Bio-Mirror sites also offer search and analyses services of these data with SRS (Zdobnov *et al.*, 2002) and other programs. The Perl FTP *mirror* package (McLoughlin, 1998) is used at Bio-Mirror sites to maintain daily updates, with additional Perl tools for updating local databanks available in the project collection.

[INSERT FIGURE 1 HERE]

DISCUSSION

The quantity of popular bio-data served through this project has grown from under 10 GB in 1998 to 150 GB in 2003. Much of this has come with growth in core databanks. As Human Genome and related projects matured, genome databases such as Ensembl, euGenes, and NCBI Genomes have added significant quantities of new data. European bioinformatics centers in the EMBnet group offer similar bio-data distribution, and many bioinformatics groups see need for improving data distribution. Others are encouraged to join the Bio-Mirror project.

Improving bio-data distribution. Rsync is a useful alternative to FTP, as it includes file system synchronization similar to the Perl mirror package that updates only changed files. Rsync also attempts to synchronize only changed sections within files, though for binary-compressed large databanks this may reduce efficiency. GridFTP (Allcock *et al.*, 2002), another possible improvement, supports parallel transfers and other efficiency methods. Although it can double transfer rates in a local network, our tests indicate problems such as lack of anonymous transfers, limited support for 64-bit systems, limited mirroring options, and bandwidth costs associated with parallel transfer.

Data Grids. New technology for computable access to large databanks is being developed in the context of science data grids (Avaki Corp., 2002; OSGA-DAI, 2002; iVDGL, 2002). Web Services examples in bioinformatics include bio-databank access (XEMBL, XDDBJ). Lightweight Directory Access Protocol (LDAP) directories of bio-databanks are also available (Gilbert, 2002). Open Grid Services Architecture (OGSA) has growing support in bioinformatics (EGD-WP10, 2000; MyGrid, 2002), with a database access and integration component (OSGA-DAI, 2002) that is relevant.

Object versus Bulk distribution. Bioinformatics groups should consider investing in object or record-oriented exchange. When only changed records are updated, transport time and costs are minimized. Record-oriented transport has many uses, where selection of data subsets tied to search systems is often of most interest. For grid computing, one wants to rapidly distribute task and data subsets to many compute nodes. One practical goal for the Bio-Mirror project is to develop more efficient object-level distribution using available technology. Tests are in progress with FTP, LDAP, SOAP and SRS to search and retrieve gigabytes of biosequence records (Gilbert, 2002), and suggest that LDAP matches or surpasses FTP when data selection is included, and is about 10 times faster than Web Services, due to its compact binary-encoded transport and more direct route from server storage to client applications.

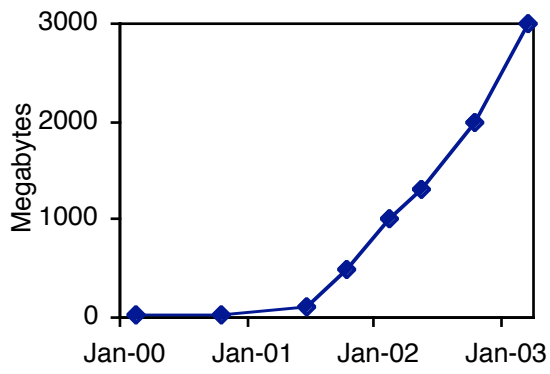


Figure 1. Monthly data transfer at US Bio-mirror.

ACKNOWLEDGEMENTS

Rick McMullen, Doug Pearson and Hisashi Eguchi (TransPac and APAN), Kazuya Tamura (Forestry and Forest Products Research Institute, Japan), Peter Stoehr and Rodrigo Lopez (EBI), Mark Cavanaugh and Scott McGinnis (NCBI) have provided valuable assistance. Bio-Mirror project has received support from Australian National University, Australia; Agriculture, Forestry and Fisheries Research Council, Japan; APBioNet of Asia-Pacific Advanced Network; Chinese Academy of Sciences, China; Korea Advanced Institute of Science and Technology, Korea; National Yang-Ming University, Taiwan; National Science Foundation and Indiana University, USA; National University of Singapore, Singapore.

REFERENCES

- Allcock, B, J. Bester, J. Bresnahan, A. L. Chervenak, I. Foster, C. Kesselman, S. Meder, V. Nefedova, D. Quesnal, S. Tuecke. (2002) Data Management and Transfer in High Performance Computational Grid Environments. *Parallel Computing Journal*, Vol. 28 (5), pp. 749-771.
- Avaki Corporation (2002) Avaki Data Grid 3.0, Java-based data grid software. URL: <http://www.avaki.com/news/release20021209.html>
- DDBJ. The Dna Databank of Japan. URL: <http://www.ddbj.nig.ac.jp/>
- EBI. The European Bioinformatics Institute. URL: <http://www.ebi.ac.uk/>
- EGD-WP10 (2000) Grid-aware biomedical applications for datagrid testbed assessment, WP10 of European Union Data Grid Project. URL: <http://edg-wp10.healthgrid.org/>, <http://marianne.in2p3.fr/datagrid/wp10/>
- EMBnet. The European Molecular Biology network. URL: <http://www.embnet.org/>
- Foster, I. and Kesselman, C., eds. (1999) *The Grid: Blueprint for a New Computing Infrastructure*. Morgan Kaufmann.
- Gilbert, D.G. (2002) Directories of Bio-data. <http://iubio.bio.indiana.edu/biogrid/directories/>
- iVDGL (2002) The International Virtual-Data Grid Laboratory. URL: <http://www.ivdgl.org/>
- McLoughlin, L. (1998) Mirror perl package. URL: <http://sunsite.org.uk/packages/mirror/>
- MyGrid (2002) The MyGrid project, Univ. Manchester, UK. URL: <http://www.mygrid.org.uk/>
- NCBI. The National Center for Biotechnology Information. URL: <http://www.ncbi.nih.gov/>
- OSGA-DAI (2002) Open Grid Services Architecture: Data Access and Integration project. URL: <http://www.ogsadai.org.uk/>
- Tridgell, A. (2002) Rsync, remote file synchronization system. URL: <http://rsync.samba.org/>
- XDDBJ. XML Central of DDBJ. URL: <http://xml.nig.ac.jp/>
- XEMBL. EMBL Nucleotide Sequence data in XML. URL: <http://www.ebi.ac.uk/xembl/>
- Zdobnov, E.M., Lopez,R., Apweiler,R., Eitzold, T. (2002) The EBI SRS server – recent developments. *Bioinformatics*, 18, 368-373.