

Software

Open Access

SciDBMaker: new software for computer-aided design of specialized biological databases

Riadh Hammami¹, Abdelmajid Zouhir¹, Karim Naghmouchi², Jeannette Ben Hamida¹ and Ismail Fliss*³

Address: ¹Unité de Protéomique Fonctionnelle & Biopréservation Alimentaire, Institut Supérieur des Sciences Biologiques Appliquées de Tunis, Université El Manar, Tunisie, ²Agriculture and Agri-Food Canada, Lethbridge Research Centre, Lethbridge, Alberta, T1J 4B1 Canada and ³Institut des Nutraceutiques et des Aliments Fonctionnels (INAF), Université Laval, Québec, Canada

Email: Riadh Hammami - Riadh.hammami@fst.rnu.tn; Abdelmajid Zouhir - azouheirb10@yahoo.fr; Karim Naghmouchi - naghmouchik@agr.gc.ca; Jeannette Ben Hamida - benhamida_jeannette@yahoo.fr; Ismail Fliss* - ismail.fliss@aln.ulaval.ca

* Corresponding author

Published: 25 February 2008

Received: 6 December 2007

BMC Bioinformatics 2008, 9:121 doi:10.1186/1471-2105-9-121

Accepted: 25 February 2008

This article is available from: <http://www.biomedcentral.com/1471-2105/9/121>

© 2008 Hammami et al; licensee BioMed Central Ltd.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/2.0>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Abstract

Background: The exponential growth of research in molecular biology has brought concomitant proliferation of databases for stocking its findings. A variety of protein sequence databases exist. While all of these strive for completeness, the range of user interests is often beyond their scope. Large databases covering a broad range of domains tend to offer less detailed information than smaller, more specialized resources, often creating a need to combine data from many sources in order to obtain a complete picture. Scientific researchers are continually developing new specific databases to enhance their understanding of biological processes.

Description: In this article, we present the implementation of a new tool for protein data analysis. With its easy-to-use user interface, this software provides the opportunity to build more specialized protein databases from a universal protein sequence database such as Swiss-Prot. A family of proteins known as bacteriocins is analyzed as 'proof of concept'.

Conclusion: SciDBMaker is stand-alone software that allows the extraction of protein data from the Swiss-Prot database, sequence analysis comprising physicochemical profile calculations, homologous sequences search, multiple sequence alignments and the building of new and more specialized databases. It compiles information with relative ease, updates and compares various data relevant to a given protein family and could solve the problem of dispersed biological search results.

Background

The exponential growth of molecular biology research in recent decades has brought concomitant growth in the number and size of databases used to interpret experimental findings. For example, UniProtKB/Swiss-Prot

release 53.2, dated 26-06-07, contains 272,212 sequence entries comprising 99,940,143 amino acids, abstracted from 157,086 references [1]. A variety of protein sequence databases exist, ranging from simple sequence repositories to expertly curated universal databases that cover all

species and in which the original sequence data are enhanced by manual addition of further information in each sequence record [2]. While all of these strive for completeness, the range of user interests is often beyond their scope. This may reflect the user's wish to combine different types of information or the inability of a single resource to contain the complete details of every relevant experiment. In addition, large databases with broad domains tend to offer less detailed information than smaller, more specialized, resources, with the result that data from many resources may need to be combined to provide a complete picture. There is a clear need to gather, filter and critically evaluate this mass of information so that it can be used with greater efficiency. Since scientists are continually developing new specific databases to enhance their understanding of biological processes, we created SciDBMaker to provide a tool for easy building of new specialized protein knowledge bases. This paper describes the development of new stand-alone software, *Scientific DataBase Maker*, for protein data analysis with online and/or off-line access. The software interface allows successive steps for sequence manipulation, starting from user sequence search and homologous sequence retrieval from the SwissProt databank, followed by physicochemical profile calculations, multiple sequence alignments, phylogenetic tree visualization and culminating in database export/building. All steps are performed in an interactive manner. Physical and chemical parameters, rarely found in public databases, provide a helpful tool for the analysis of a set of proteins and their calculation is achieved in a direct and interactive manner, with off-line access. SciDBMaker also processes a great number of sequences simultaneously.

Implementation

Swiss-Prot format

The Swiss-Prot format has been described previously in reference [3].

Physicochemical profiles

Protein families may be analyzed with the help of physicochemical profiles such as amino acid composition (acidic, basic, hydrophobic, polar, absent and common amino acids), atomic composition, molecular weight [4], theoretical pI [4,5], extinction coefficient [6], absorbance at 280 nm, estimated half-life in mammalian cells, yeast and *E. coli* [7,8], instability index [9], aliphatic index [10], grand average of hydropathicity (GRAVY) [11] and protein-binding potential (Boman index) [12].

Integrated tools

The European Bioinformatics Institute provides the Dbfetch tool for easy retrieval of entries from various databases [13]. Entries may be imported online into SciDBMaker from the SwissProt database using Dbfetch. To find similar sequences, the containers can be queried with either proteins from the SwissProt database or user-imported sequences, using the BLAST algorithm [14]. Multiple sequence alignments (MSA) are an essential tool for predicting protein structure and function prediction, phylogenetic inference and other common tasks in sequence analysis. To date, CLUSTALW is still the most popular alignment tool. Since it is the method of choice for biologists, CLUSTALW [15] was included in SciDBMaker for multiple sequence alignments. Generated trees may be easily viewed using phylogenetic tree visualization software such as TREEVIEW [16].

Hardware and software specifications

The executable version of the SciDBMaker software can be installed and run on a standard PC platform with a Windows operating system. The software development was done using Windows XP and tested with success on all platforms, including Win 98, Win XP and Win vista. The source code was written in Microsoft Visual Basic .NET (2005). The environment is based upon the .NET Framework library v2.0.

Results & discussion

Program description

The workflow diagram shown in Figure 1 and the following discussion illustrate how the tool works. Figure 2 illustrates a typical user interface of the program. Users may open files in Fasta or Swiss-Prot format, or import sequence entries from the Swiss-Prot database. Users may also use their own sequences, search for homologous sequences entries in Swiss-Prot database using BLAST algorithm and load selected entries into SciDBMaker (Fig. 3). The program will automatically extract available information in Swiss-Prot entries and calculate physicochemical profiles for loaded proteins. Users may also choose the information to be extracted and the properties to be calculated, as shown in Figure 4. The interface allows users to filter, search, add, remove and update data rows as required. An intuitive interface allows BLAST selection of all user sequences. Similarly, sequences may be aligned using the multiple alignment program ClustalW. Resulting trees may be shown using the phylogenetic tree visualization software TREEVIEW, as proposed by SciDBMaker. As a final step, data may be printed or saved in various file formats. Sequences may be extracted into a Fasta format file. The resulting data grid may be saved as an MS Excel data sheet, as well as database files (XML, MS Access, MySQL).

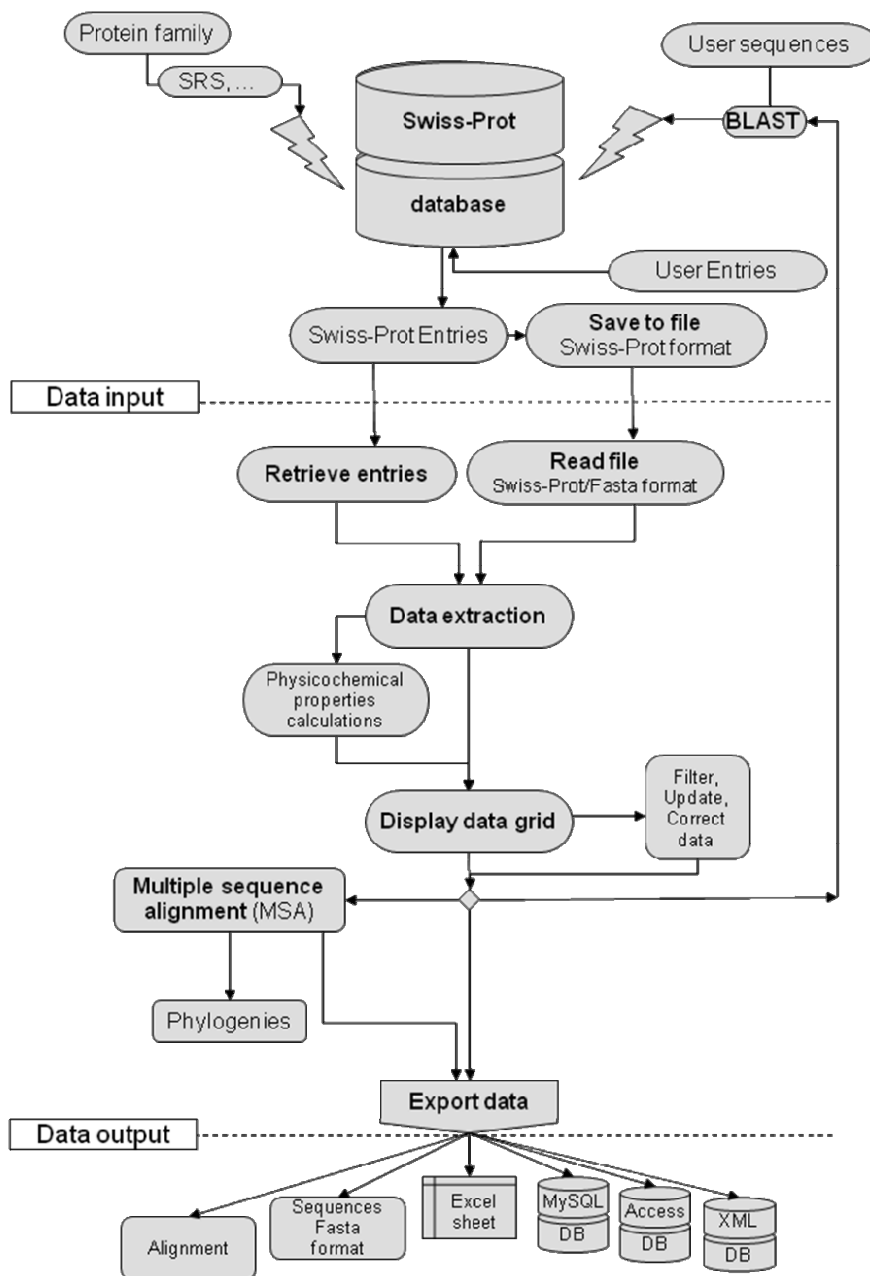


Figure 1
Workflow diagram.

Program runs

A database was developed using SciDBMaker. Named BACTIBASE, this database has been previously described in reference [17].

Conclusion

The stand-alone software SciDBMaker allows the extraction of protein data from the Swiss-Prot database,

sequence analysis comprising physicochemical profile calculations, homologous sequence searches, multiple sequence alignments and the building of new and more specialized databases. Programs of this type compile information with relative ease, update and compare various data relevant to a given protein family and could solve the problem of dispersed biological search results. Collection of a multitude of information regarding a given pro-

N°	Peptide	Sequence	Length	Gene	Swiss-Prot ID	Swiss Entry	PDB Entry	Other Database	Taxonomy	Producer organism	Organism Host	References	Feature Table
1	Leukocyte surface antigen C...	MWFLAALLL	303	Cd47	CD47_RAT	P97829; Q352...		EMBL: D8765	Eukaryota; Met...	Rattus norvegi...		[1]NUCLEOTID...	SIGNAL 1
2	Integrin alpha ina-1 precursor	MRECIISWTL...	1139	ina-1; ORFNa...	INA1_CAEL	Q03600;		EMBL: Z19155	Eukaryota; Met...	Caenorhabditis...		[1]NUCLEOTID...	SIGNAL 1
3	Integrin alpha-10 precursor	MELPFVTHLF...	1167	ITGA10; ORFN...	ITA10_HUMAN	O75578; Q6UX...		EMBL: AF0740	Eukaryota; Met...	Homo sapiens ...		[1]NUCLEOTID...	SIGNAL 1
4	Integrin alpha-11 precursor	MDLPRGLVA...	1189	ITGA11; ORFN...	ITA11_HUMAN	Q9UK05; Q8W...		EMBL: AF1096	Eukaryota; Met...	Homo sapiens ...		[1]NUCLEOTID...	SIGNAL 1
5	Integrin alpha-11 precursor	MDFPRQLLVA...	1188	Itga11	ITA11_MOUSE	P61622;		EMBL: BC058	Eukaryota; Met...	Mus musculus ...		[1]NUCLEOTID...	SIGNAL 1
6	Integrin alpha-1 (Laminin a...	ENMTFGTTLV...	285	ITGA1	ITA1_CHICK	Q90615;		EMBL: U1011...	Eukaryota; Met...	Gallus gallus (...)		[1]NUCLEOTID...	CHAIN &t...
7	Integrin alpha-PS1 precursor...	MLELPFTTIRP...	1146	mew; ORFNa...	ITA1_DROME	Q24247; Q8SY...		EMBL: X73975	Eukaryota; Met...	Drosophila me...		[1]NUCLEOTID...	SIGNAL 1
8	Integrin alpha-1 precursor (...)	MAPRRPFRAP...	1179	ITGA1	ITA1_HUMAN	P56199;	1PT6;X-ray1Q...	EMBL: AC0273	Eukaryota; Met...	Homo sapiens ...		[1]NUCLEOTID...	SIGNAL 1
9	Integrin alpha-1 precursor (...)	MVPRRPASLE...	1179	Itga1	ITA1_MOUSE	Q3V3R4;		EMBL: AK0358	Eukaryota; Met...	Mus musculus ...		[1]NUCLEOTID...	SIGNAL 1
10	Integrin alpha-1 precursor (...)	MVPRRPASLE...	1180	Itga1	ITA1_RAT	P18614;	1CK4;X-ray1M...	EMBL: X52140	Eukaryota; Met...	Rattus norvegi...		[1]NUCLEOTID...	SIGNAL 1
11	Integrin alpha-1Ib precursor ...	MARALCPLQA...	1039	ITGA2B; Synon...	ITA2B_HUMAN	P08514; Q953...	1DPK;NMR1D...	EMBL: J02764...	Eukaryota; Met...	Homo sapiens ...		[1]NUCLEOTID...	SIGNAL 1
12	Integrin alpha-1Ib precursor ...	MARASCAWH...	1033	Itga2b	ITA2B_MOUSE	Q9QUM0; Q64...		EMBL: AF1698	Eukaryota; Met...	Mus musculus ...		[1]NUCLEOTID...	SIGNAL 1
13	Integrin alpha-1Ib (Platelet ...)	QVLDSPPFTG...	604	ITGA2B	ITA2B_PAPCY	P53711;		EMBL: L12233	Eukaryota; Met...	Papio cynocep...		[1]NUCLEOTID...	CHAIN &t...
14	Integrin alpha-2 precursor (...)	PLQLVLVFSQ...	1170	ITGA2	ITA2_BOVIN	P53710;		EMBL: L25886	Eukaryota; Met...	Bos taurus (Bo...		[1]NUCLEOTID...	SIGNAL &t...
15	Integrin alpha-PS2 precursor...	MSGDSIHRR...	1396	#; ORFNAMES...	ITA2_DROME	P12080; Q9VX...		EMBL: M1905	Eukaryota; Met...	Drosophila me...		[1]NUCLEOTID...	SIGNAL 1
16	Integrin alpha-2 precursor (...)	MGPERTGAAP...	1181	ITGA2; Synony...	ITA2_HUMAN	P17301; Q145...	1A0XX-ray1D...	EMBL: X17033	Eukaryota; Met...	Homo sapiens ...		[1]NUCLEOTID...	SIGNAL 1
17	Integrin alpha-2 precursor (...)	MGPGQAGGA...	1178	Itga2	ITA2_MOUSE	Q62469; Q621...		EMBL: Z29987	Eukaryota; Met...	Mus musculus ...		[1]NUCLEOTID...	SIGNAL 1
18	Integrin alpha-3 precursor (...)	MGPGRPCAP...	1066	ITGA3	ITA3_CRIGR	P17852;		EMBL: J05281...	Eukaryota; Met...	Cricetulus gris...		[1]NUCLEOTID...	SIGNAL 1
19	Integrin alpha-PS3 precursor...	MNAESTMFP...	1115	scb; Synonyme...	ITA3_DROME	O44386; Q443...		EMBL: AF0341...	Eukaryota; Met...	Drosophila me...		[1]NUCLEOTID...	SIGNAL 1
20	Integrin alpha-3 precursor (...)	MGPQPSRAP...	1066	ITGA3; Synony...	ITA3_HUMAN	P26006;		EMBL: M5991	Eukaryota; Met...	Homo sapiens ...		[1]NUCLEOTID...	SIGNAL 1
21	Integrin alpha-3 precursor (...)	MGPQPCRPV...	1053	Itga3	ITA3_MOUSE	Q62470; Q084...		EMBL: D1386...	Eukaryota; Met...	Mus musculus ...		[1]NUCLEOTID...	SIGNAL 1
22	Integrin alpha-PS4 precursor...	MVAAPRANS...	1015	alpha-PS4; OR...	ITA4_DROME	Q9V744;		EMBL: AE0135	Eukaryota; Met...	Drosophila me...		[1]NUCLEOTID...	SIGNAL 1
23	Integrin alpha-4 precursor (...)	MFPTESAWLG...	1038	ITGA4; Synony...	ITA4_HUMAN	P13612;		EMBL: X16983	Eukaryota; Met...	Homo sapiens ...		[1]NUCLEOTID...	SIGNAL 1
24	Integrin alpha-4 precursor (...)	MFSTKSAWLR...	1039	Itga4	ITA4_MOUSE	Q00651;		EMBL: X53176	Eukaryota; Met...	Mus musculus ...		[1]NUCLEOTID...	SIGNAL 1
25	Integrin alpha-4 precursor (...)	MIRDLGKVGK...	1032	Itga4	ITA4_XENLA	Q91687; Q062...		EMBL: U5449...	Eukaryota; Met...	Xenopus laevis...		[1]NUCLEOTID...	SIGNAL 1
26	Integrin alpha-5 (Fibronecti...	CGEENICVPD...	385	ITGA5	ITA5_BOVIN	Q27977;		EMBL: U1086...	Eukaryota; Met...	Bos taurus (Bo...		[1]NUCLEOTID...	CHAIN &t...
27	Integrin alpha-PS5 precursor...	MNFSPLPNRV...	1000	alpha-PS5; OR...	ITA5_DROME	Q9W1M8;		EMBL: AE0135	Eukaryota; Met...	Drosophila me...		[1]NUCLEOTID...	SIGNAL 1
28	Integrin alpha-5 precursor (...)	MGRSPTRESPL...	1049	ITGA5; Synony...	ITA5_HUMAN	P08648; Q96H...		EMBL: X06256	Eukaryota; Met...	Homo sapiens ...		[1]NUCLEOTID...	SIGNAL 1
29	Integrin alpha-5 precursor (...)	MGSWTRPSP...	1053	Itga5	ITA5_MOUSE	P11688;		EMBL: X79003	Eukaryota; Met...	Mus musculus ...		[1]NUCLEOTID...	SIGNAL 1
30	Integrin alpha-5 precursor (...)	MQLPRGSRV...	1050	Itga5	ITA5_XENLA	Q06274;		EMBL: U1268...	Eukaryota; Met...	Xenopus laevis...		[1]NUCLEOTID...	SIGNAL 1
31	Integrin alpha-6 precursor (...)	MAAALLLYLPL...	1072	ITGA6	ITA6_CHICK	P26007;		EMBL: X56559	Eukaryota; Met...	Gallus gallus (...)		[1]NUCLEOTID...	SIGNAL 1
32	Integrin alpha-6 precursor (...)	MAAAGQLCLL...	1130	ITGA6	ITA6_HUMAN	P23229; Q084...		EMBL: X53586	Eukaryota; Met...	Homo sapiens ...		[1]NUCLEOTID...	SIGNAL 1
33	Integrin alpha-6 precursor (...)	MAVAGQLCLL...	1091	Itga6	ITA6_MOUSE	Q61739;		EMBL: X69902	Eukaryota; Met...	Mus musculus ...		[1]NUCLEOTID...	SIGNAL 1

Figure 2 User interface.

tein family allows the development of more coherent and focused approaches to structure-function relationships, thereby enhancing the development of theoretical concepts in biological sciences.

Availability and requirements

The program runs on the PC platform with a Windows operating system. The graphical environment needs the .NET Framework library v2.0. This complement is available for free download at the Microsoft website and comes pre-installed in the majority of recent computers. An installation package for SciDBMaker may be obtained from the authors free of charge upon request. The SciDBMaker website is hosted by the Centre de Calcul El Kharizmi CCK (Tunisia) and is available at <http://scidbmaker.pfba-lab.org>. The SciDBMaker software is provided 'as is' with no guarantee or warranty of any kind and is available for all non-commercial use. Any other use of the software requires special permission from the primary author.

Authors' contributions

RH programmed the software interface, performed the implementation of physicochemical parameters and drafted the manuscript. AZ participated in the design of the study, interacted with RH to carry out the physicochemical data calculation and tested the program. KN tested the program and contributed to the manuscript. JBH oversaw the project and helped define user requirements. IF conceived the study, participated in its design and coordination and helped to draft the manuscript. All authors read and approved the final manuscript.

Acknowledgements

The authors thank Dr. Stephen Davids for proofreading the manuscript. This research was supported by Ministry of Higher Education, Scientific Research and Technology, Republic of Tunisia.

References

1. **The Universal Protein Resource (UniProt).** *Nucleic Acids Res* 2007, **35**:D193-D197.
2. Apweiler R, Bairoch A, Wu HC: **Protein sequence databases.** *Curr Opin Chem Biol* 2004, **8**:76-80.
3. Boeckmann B, Bairoch A, Apweiler R, Blatter MC, Estreicher A, Gasteiger E, Martin MJ, Michoud K, O'Donovan C, Phan I, Pilbout S, Schneider M: **The SWISS-PROT protein knowledgebase and its supplement TrEMBL in 2003.** *Nucleic Acids Res* 2003, **31**:365-370.
4. Patrickios CS, Yamasaki EN: **Polypeptide amino acid composition and isoelectric point.** *Anal Biochem* 1995, **231**:82-91.
5. Bjellqvist B, Basse B, Olsen E, Celis JE: **Reference points for comparisons of two-dimensional maps of proteins from different human cell types defined in a pH scale where isoelectric points correlate with polypeptide compositions.** *Electrophoresis* 1994, **15**:529-539.
6. Henryk M, Russell MC, Randolph VL: **Statistical determination of the average values of the extinction coefficients of tryptophan and tyrosine in native proteins.** *Anal Biochem* 1992, **200**:74-80.
7. Bachmair A, Finley D, Varshavsky A: **In vivo half-life of a protein is a function of its amino-terminal residue.** *Science* 1986, **234**:179-186.
8. Gonda DK, Bachmair A, Wunning I, Tobias JW, Lane WS, Varshavsky AJ: **Universality and structure of the N-end rule.** *J Biol Chem* 1989, **264**:16700-16712.
9. Guruprasad K, Reddy BVB, Pandit MW: **Correlation between stability of a protein and its dipeptide composition: a novel approach for predicting in vivo stability of a protein from its primary sequence.** *Protein Eng* 1990, **4**:155-161.
10. Ikai AJ: **Thermostability and aliphatic index of globular proteins.** *J Biochem* 1980, **88**:1895-1898.
11. Jack K, Russell FD: **A simple method for displaying the hydrophobic character of a protein.** *J Mol Biol* 1982, **157**:105-132.
12. Radzeka A, Wolfenden R: **Comparing the polarities of amino acids: side-chain distribution coefficients between vapor phase, cyclohexane, 1-octanol and neutral aqueous solution.** *Biochemistry* 1988, **27**:1664-1670.
13. Labarga A, Valentin F, Andersson M, Lopez R: **Web Services at the European Bioinformatics Institute.** *Nucleic Acids Res* 2007:W6-11.
14. Altschul SF, Madden TL, Schaffer AA, Zhang J, Zhang Z, Miller W, Lipman DJ: **Gapped BLAST and PSI-BLAST: a new generation of protein database search programs.** *Nucleic Acids Res* 1997, **25**:3389-3402.
15. Larkin MA, Blackshields G, Brown NP, Chenna R, McGettigan PA, McWilliam H, Valentin F, Wallace IM, Wilm A, Lopez R, Thompson JD, Gibson TJ, Higgins DG: **ClustalW and ClustalX version 2.0.** *Bioinformatics* 2007, **23**(21):2947-2948.
16. Page RDM: **TREEVIEW: An application to display phylogenetic trees on personal computers.** *CABIOS* 1996, **12**:357-358.
17. Hammami R, Zouhir A, Ben Hamida J, Fliss I: **BACTIBASE: a new web-accessible database for bacteriocin characterization.** *BMC Microbiol* 2007, **7**:89.

Publish with **BioMed Central** and every scientist can read your work free of charge

"BioMed Central will be the most significant development for disseminating the results of biomedical research in our lifetime."

Sir Paul Nurse, Cancer Research UK

Your research papers will be:

- available free of charge to the entire biomedical community
- peer reviewed and published immediately upon acceptance
- cited in PubMed and archived on PubMed Central
- yours — you keep the copyright

Submit your manuscript here:
http://www.biomedcentral.com/info/publishing_adv.asp

