



HAL
open science

Bio++: a set of C++ libraries for sequence analysis, phylogenetics, molecular evolution and population genetics

Julien Dutheil, Sylvain Gaillard, Eric Bazin, Sylvain Glémin, Vincent Ranwez,
Nicolas Galtier, Khalid Belkhir

► To cite this version:

Julien Dutheil, Sylvain Gaillard, Eric Bazin, Sylvain Glémin, Vincent Ranwez, et al.. Bio++: a set of C++ libraries for sequence analysis, phylogenetics, molecular evolution and population genetics. BMC Bioinformatics, 2006, 7, pp.188. 10.1186/1471-2105-7-188 . halsde-00323971

HAL Id: halsde-00323971

<https://hal.science/halsde-00323971>

Submitted on 25 May 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

Software

Open Access

Bio++: a set of C++ libraries for sequence analysis, phylogenetics, molecular evolution and population genetics

Julien Dutheil*¹, Sylvain Gaillard¹, Eric Bazin¹, Sylvain Glémin¹, Vincent Ranwez², Nicolas Galtier¹ and Khalid Belkhir¹

Address: ¹CNRS UMR 5171 – Génome, Populations, Interactions, Adaptation (GPIA), Université Montpellier 2, France and ²CNRS UMR 5554 – Institut des Sciences de l'Évolution de Montpellier (ISE-M), Université Montpellier 2, France

Email: Julien Dutheil* - Julien.Dutheil@univ-montp2.fr; Sylvain Gaillard - syl.gaillard@gmail.com; Eric Bazin - bazin@univ-montp2.fr; Sylvain Glémin - glemin@univ-montp2.fr; Vincent Ranwez - ranwez@isem.univ-montp2.fr; Nicolas Galtier - galtier@univ-montp2.fr; Khalid Belkhir - belkhir@univ-montp2.fr

* Corresponding author

Published: 04 April 2006

Received: 06 December 2005

BMC Bioinformatics 2006, 7:188 doi:10.1186/1471-2105-7-188

Accepted: 04 April 2006

This article is available from: <http://www.biomedcentral.com/1471-2105/7/188>

© 2006 Dutheil 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: A large number of bioinformatics applications in the fields of bio-sequence analysis, molecular evolution and population genetics typically share input/output methods, data storage requirements and data analysis algorithms. Such common features may be conveniently bundled into re-usable libraries, which enable the rapid development of new methods and robust applications.

Results: We present Bio++, a set of Object Oriented libraries written in C++. Available components include classes for data storage and handling (nucleotide/amino-acid/codon sequences, trees, distance matrices, population genetics datasets), various input/output formats, basic sequence manipulation (concatenation, transcription, translation, etc.), phylogenetic analysis (maximum parsimony, markov models, distance methods, likelihood computation and maximization), population genetics/genomics (diversity statistics, neutrality tests, various multi-locus analyses) and various algorithms for numerical calculus.

Conclusion: Implementation of methods aims at being both efficient and user-friendly. A special concern was given to the library design to enable easy extension and new methods development. We defined a general hierarchy of classes that allow the developer to implement its own algorithms while remaining compatible with the rest of the libraries. Bio++ source code is distributed free of charge under the CeCILL general public licence from its website <http://kimura.univ-montp2.fr/BioPP>.

Background

The design of re-usable software components into libraries has proved to be useful to the rapid development of bioinformatics applications, as witnesses the success of the ever growing open source libraries known as the Bio{*} projects. These libraries cover many fields of bio-

informatics, but rarely offer tools for intensive calculus, as required for phylogenetic analyses for instance. Such computer-expensive applications are usually written in C, a compiled low-level language. Programming time-consuming tasks often requires the implementation of specific algorithms optimized for each particular problem.

However, despite these specificities, there are some common points (e.g. likelihood calculation, minimization), and developers will appreciate to use modules pre-programmed in a flexible manner.

The use of the C++ language combines both the computer efficiency of the C language and the convenience of an Object-Oriented (OO) approach. We present Bio++, a set of C++ libraries dedicated to sequence analysis, molecular evolution and population genetics. Other libraries taking these advantages of C++ have been developed, including the libsequence library [1] for population genetics analysis (single nucleotides polymorphism and coalescent), PDBlib for manipulating tridimensional structures [2] or libcov [3] for protein phylogenetics and molecular evolution. The Bioinformatics Template Library [4] also provides algorithms for manipulating and analysing three dimensional structures, using a generic programming approach.

Implementation

We did a collaborative effort to design a hierarchy of useful classes for the rapid development of efficient applications in the fields of sequence analysis, phylogenetics, molecular evolution and population genetics. These libraries combine tools for data handling and numerical calculations, providing an easy-to-use, powerful and general development environment. The three main objectives of Bio++ are (i) to speed up the implementation of a new method in order to test it, (ii) to make easier the development of robust applications for distribution and (iii) to provide a high level programming language for biologists.

Library design

We used object-oriented programming paradigms to develop Bio++. We defined for each class an abstract basal class that contains a number of purely virtual functions (dummy functions with an empty body). Such a basal class is called an interface and can be seen as a kind of contract that a class must follow in order to be used by other methods. Such a contract only specifies what an object is at least able to do, and not how it does it. Most classes in Bio++ are derived from interfaces, and some common instantiations are proposed for each interface. The user may, however, write its own implementation of a given method, and still remain compatible with the rest of the library. For instance, one may easily design a new alphabet to handle a new kind of sequences, or write a new substitution model and perform likelihood calculation without writing any likelihood computation function. Conversely, one may implement a new likelihood computation algorithm without re-developing substitution model classes. Several interfaces and abstract classes (partial implementations) linked by inheritance specify different levels of specialization. Depending on their

objectives and programming level, users might simply use fully-specified objects, re-implement specific methods, or even design new classes.

The sequence container hierarchy is a representative example of the class hierarchies defined in Bio++ (see figure 1). The SequenceContainer basal interface only specifies that sequences can be accessed by their name. The OrderedSequenceContainer has the additional requirement that sequences can be accessed by index, which is generally the case, even if ordering has no biological meaning. The SiteContainer interface requires that sequences have the same length (i.e. are aligned) and hence may also be accessed by site (= column in the alignment). Several implementations are proposed: in the VectorSequenceContainer class, data are stored as a vector of sequences. Each sequence can hence be accessed in $O(1)$. If there are N sequences, they are hence accessed in $O(N)$. The AlignedSequenceContainer is derived from the VectorSequenceContainer class, and adds site access. Since data are stored as sequences, the access time for a site is in $O(N)$ and the complete set of L sites is accessed in $O(N \times L)$. The VectorSiteContainer proposes an alternative implementation by storing sequences as sites instead of sequence objects. The sequence access is hence achieved in $O(L)$, and the site access in $O(1)$. All methods working on containers only deal with the SequenceContainer or the OrderedSequenceContainer interface, without any assumption about the implementation. The execution time, however, may vary depending on the implementation used.

Development is facilitated by the use of the code documentation generated by the doxygen program [5]. Full class documentation can be consulted online or downloaded in navigable HTML format. A short tutorial is also available.

Library content

We split Bio++ classes into five libraries: three main biological libraries (SeqLib, PhylLib and PopGen-Lib) and two utility libraries (Utils and NumCalc), see table 1 for their content and dependencies:

- the Utils library contains core classes and utility functions;
- the NumCalc library contains classes for numerical calculus, including several optimization tools, random number generators, probability distributions and statistical functions;
- the SeqLib library is dedicated to sequence analysis: sequence storage and manipulation, input/output toward several file formats, alphabet translation, etc.

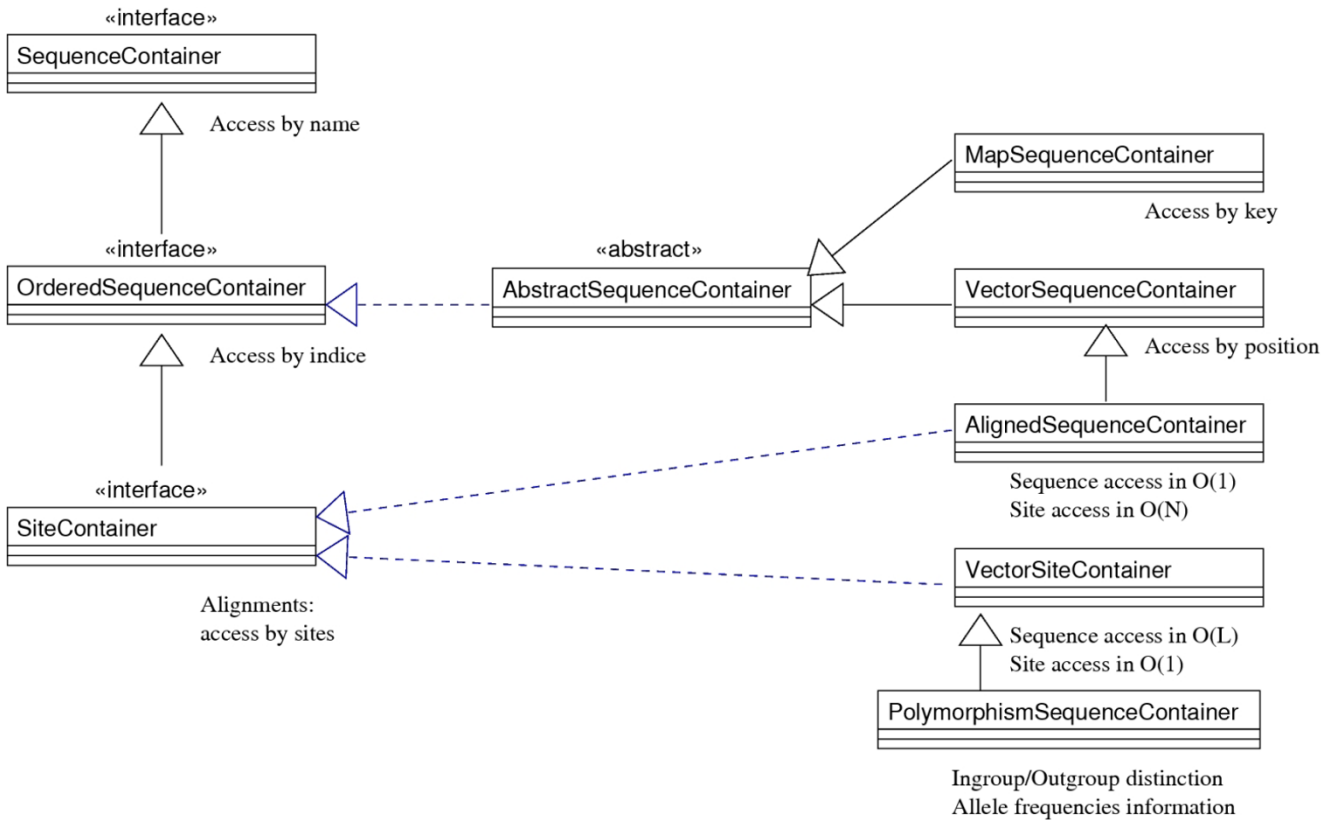


Figure 1
Class hierarchy for the sequence container family. "Interfaces" are classes with only pure virtual methods. "Abstract" classes are partial implementations of interfaces. Full arrows indicate inheritance links and dashed arrows represent implementation links.

- the PhylLib library deals with phylogenetic trees: Markov models, likelihood calculation, etc. The library is more focused on molecular evolution (model fitting, ancestral state reconstruction, across-site rate variation) than pure phylogenetic reconstruction, although several topology reconstruction methods are implemented (distance methods, nearest-neighbor interchanges (NNI) search for maximum likelihood (ML) and parsimony scores);

- the PopGenLib library is dedicated to population genetics, with particular sequence and codominant markers storage. Methods are provided for polymorphism analysis, population divergence estimation and neutrality tests.

Results

Figure 2 shows a full program example. This application builds a Neighbor-Joining tree from a sequence file in Phylip format, and re-estimates parameters (branch lengths and shape of the rate across-site distribution) using maximum likelihood. The final tree is then written to a file in Newick format. Additional output files are also created, providing detailed information about the estima-

tion procedure. The program begins with the creation of a ProteicAlphabet object (line 3):

```
const ProteicAlphabet * alphabet = new ProteicAlphabet();
```

This object specifies the type of sequence the program will use. The Alphabet object family contains classes for nucleic alphabets (DNA and RNA), proteins or codons (nuclear and mitochondrial, for vertebrates, echinoderms/nematodes and other invertebrates). A dataset is then readed from a Phylip sequence file. A sequence reader is created using

```
Phylip * seqReader = new Phylip(false, false);
```

Other file formats for reading and writing are supported, including the commonly used Fasta and Clustal. Sequences are then read and stored in a container:

```
SiteContainer * sites = seqReader->read("Myoglobin.phy", alphabet);
```

Table 1: Libraries available and dependencies.

Library	Utils	NumCalc	SeqLib	Content
Utils				Basal classes and utilities. Exceptions, text and file manipulation.
NumCalc	✓			Numerical calculus. Vector and matrix manipulation, optimization, algebra, probability distribution and statistics, random number generation.
SeqLib	✓	✓		Sequence manipulation, alphabets, chemical properties and distances, input/output in various formats (Fasta, Phylip, Clustal, etc.).
PhylLib	✓	✓	✓	Phylogenetics. Tree reconstruction by maximum parsimony, distance methods or maximum likelihood, parameter estimations, ancestral states reconstruction, etc.
PopGenLib	✓	✓	✓	Population genetics. Polymorphism statistics, linkage disequilibrium, neutrality tests and recombination.

Since sequences are aligned, they are actually stored in a SiteContainer object (see figure 1). Line 7 to 10 display on screen some properties of the container, namely the number of sequences and sites. In this example we need to restrict our analysis to sites without any gap, so at line 11 a new container called "completeSites" is created from the former:

```
SiteContainer * completeSites = SiteContainerTools::getSitesWithoutGaps(sites);
```

Several `{*}Tools` classes providing utility functions are available. For instance, the SiteContainerTools class contains static functions that work on SiteContainer objects.

A JTT proteic SubstitutionModel object is created [6,7]:

```
SubstitutionModel * model = new JTT92(alphabet);
```

Here again, several commonly used models are available. We also use a discrete gamma rate across sites distribution with 4 rate classes and a shape parameter of 0.5 [8]:

```
DiscreteDistribution * rateDist = new GammaDiscreteDistribution(4, 0.5);
```

This model is used first to estimate the distance matrix from the data. A DistanceEstimation object is then created, by giving the model and data as parameters:

```
DistanceEstimation distEstimation(model, rateDist, completeSites);
```

```
DistanceMatrix * matrix = distEstimation.getMatrix();
```

Distances will be estimated using a maximum likelihood (ML) procedure. The NumCalc library provides an object-oriented implementation of several general optimization procedures (Optimizer objects). For phylogenetic optimization, specific optimizers are also available. The program

then computes a neighbor-joining tree [9]. A NeighborJoining object is instantiated and the resulting tree stored in a variable:

```
NeighborJoining nj(* matrix, false);
```

```
Tree<Node> * tree = nj.getTree();
```

This variable is passed to a TreeLikelihood object, together with the data set, model and rate distribution. An additional parameter tells the object to display a few informations. The DRHomogeneousTreeLikelihood object is for homogeneous models, with rate across sites variations.

```
DiscreteRatesAcrossSitesTreeLikelihoodFunction = new DRHomogeneousTreeLikelihood(*tree, *completeSites, model, rateDist, true);
```

Several TreeLikelihood implementations are available. The DRHomogeneousTreeLikelihood object uses one 'view' by neighbor node for each node, *i.e.* three views by node in case of bifurcating trees [10]. This implementation saves computational time during NNI-mediated topology search, and is also more convenient to compute branch-lengths derivatives. An alternative implementation is proposed: HomogeneousTreeLikelihood, and others could be easily added if required. For parameter optimization, we used a utility function from the OptimizationTools class. This function automatically instantiates the appropriate optimizer object. It receives as argument a pointer toward the likelihood function, a tolerance number and a maximum number of function evaluations. Two optional log-files are provided to monitor the optimization process ("profiler" and "messenger", which contain several detailed informations, like for instance all parameters and function values at each step of the process), and the last parameter specifies the verbose level:

```

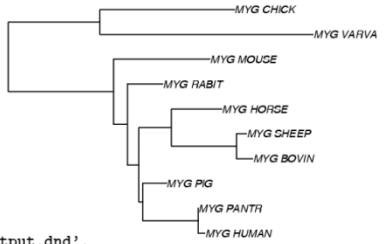
Input sequence file:
10 167
MYG.HUMAN  GLSDQEWQLVLNVWGKVEAD---IPGHQEVLRIRLFKQHPETELEKFDKFK
MYG.PANTR  GLSDQEWQLVLNVWGKVEAD---IPGHQEVLRIRLFKQHPETELEKFDKFK
MYG.VARVA  GLSDDEEKVVDIWKGVED---LPSHQEVIIRMFQNHPEIQDRFAKFK
MYG.CHICK  GLSDQEWQQLVITWGVKVEAD---IAGHGEVLRIRLFHDHPEILDRFDKFK
MYG.PIG    GLSDQEWQLVLNVWGKVEAD---VAGHQEVLRIRLFKQHPETELEKFDKFK
MYG.MOUSE  GLSDQEWQLVLNVWGKVEAD---LAGHQEVLRIRLFKTHPEILDKFDKFK
MYG.RABIT  GLSDQEWQLVLNVWGKVEAD---LAGHQEVLRIRLFHTHPEILEKFDKFK
MYG.BOVIN  GLSDQEWQLVLNVWGKVEAD---VAGHQEVLRIRLFTGHPEILEKFDKFK
MYG.SHEEP  GLSDQEWQLVLNVWGKVEAD---VAGHQEVLRIRLFTGHPEILEKFDKFK
MYG.HORSE  GLSDQEWQQLVLNVWGKVEAD---IAGHQEVLRIRLFTGHPEILEKFDKFK
...
    
```

```

Number of sequences in alignment: 10
Number of sites in alignment: 377
Number of complete sites in alignment: 153
Estimating distance matrix...
*****Done.
Building tree
Done.
Estimating parameters using ML
Double-Recursive Homogeneous Tree Likelihood
Initializing data structure..... Done.
Number of distinct sites.....: 83
    
```

```

Branch lengths:
*****
Rate distribution:
alpha:.....
Branch lengths:
****
Rate distribution:
alpha:.....
Branch lengths:
*
Rate distribution:
alpha:.....Done.
Writing results to file 'output.dnd'.
    
```



```

1 int main() {
2     try {
3         const ProteicAlphabet *alphabet = new ProteicAlphabet();
4         Phylip *seqReader = new Phylip(false, false);
5         SiteContainer *sites =
6             seqReader->read("Myoglobin.phy", alphabet);
7         cout << "Number_of_sequences_in_alignment:uuuuuu"
8             << sites->getNumberOfSequences() << endl;
9         cout << "Number_of_sites_in_alignment:uuuuuuuuuu"
10            << sites->getNumberOfSites() << endl;
11        SiteContainer * completeSites =
12            SiteContainerTools::getSitesWithoutGaps(* sites);
13        cout << "Number_of_complete_sites_in_alignment:u"
14            << completeSites->getNumberOfSites() << endl;
15        delete sites;
16
17        SubstitutionModel *model = new JTT92(alphabet);
18        DiscreteDistribution *rateDist =
19            new GammaDiscreteDistribution(4, 0.5);
20
21        cout << "Estimating_distance_matrix..." << endl;
22        DistanceEstimation distEstimation(
23            model, rateDist, completeSites);
24        DistanceMatrix *matrix = distEstimation.getMatrix();
25        cout << "Done." << endl;
26
27        cout << "Building_tree_u" << endl;
28        NeighborJoining nj(*matrix, false); // unrooted tree.
29        Tree<Node> *tree = nj.getTree();
30        cout << "Done." << endl;
31
32        cout << "Estimating_parameters_using_ML" << endl;
33        DiscreteRatesAcrossSitesTreeLikelihood *likFunction =
34            new DRHomogeneousTreeLikelihood(
35                *tree, *completeSites, model, rateDist, true);
36
37        // Save messages to separate files:
38        ofstream *profiler = new ofstream("profile.txt", ios::out);
39        ofstream *messenger = new ofstream("messages.txt", ios::out);
40        OptimizationTools::optimizeNumericalParameters(
41            likFunction, 0.000001, 1000000, messenger, profiler, 3);
42        cout << "Done." << endl;
43        cout << "Writing_results_to_file_'output.dnd'." << endl;
44        Newick newick;
45        newick.write(*tree, "output.dnd");
46
47        delete alphabet, seqReader, completeSites, model,
48            rateDist, matrix, tree, likFunction;
49    } catch(Exception & e) {
50        cout << "ERROR!u" << e.what() << endl;
51    }
52 }
    
```

Figure 2
Code example. Lines numbered in red are referenced in the text.

```

OptimizationTools:: optimizeNumericalPa
rameters(likFunction, 0.000001,
1000000messenger, profiler, 3);
    
```

Finally, the optimized tree is written to a file named output.dnd in Newick file format:

```

Newick newick;
newick.write(*tree, "output.dnd");
    
```

This program creates three files: two for the ML tree estimation and one for the final tree.

Conclusion

Bio++ is a mature project which has been used in previous works like molecular coevolution analysis [11] or codon analysis [12]. However it is an active project still receiving new methods and improvements.

Development snapshots may be accessed by anonymous CVS on the library website. The website also provides on-line documentation of classes, tutorial and several code

examples. Several programs have been developed with the Bio++ libraries, including bppML for ML likelihood tree estimation and bppSG for sequence generation by simulation under different kinds of models. Any contribution will be welcome, as specific functions or as additional libraries compatible with the present ones.

Availability and requirements

C++ libraries are not organized in a tree-like hierarchy as java packages or perl modules. They are bundled in a non-nested way, and may be compiled in two flavours, either static or shared (= dynamic). Dynamic libraries are loaded during the program execution and hence can be shared by several applications, while static libraries are hard-coded into the executable, which no longer requires the library to be installed.

Each Bio++ library is compliant with the GNU standards and uses the autotools suite for compilation and installation. Bio++ can hence be built with these GNU tools on any unix-like system. Alternatively, sources can be imported and compiled in any C++ development environment (IDE). It has been successfully installed on Linux,

MacOS X and Windows using the Cygwin port and MinGW. Bio++ is distributed freely under the CeCILL public license (the French free software license, compatible with the GNU General Public License). It is available from its website at <http://kimura.univ-montp2.fr/BioPP>.

Authors' contributions

JD developed the Utils, NumCalc, SeqLib and PhylLib libraries and drafted the manuscript. SGa developed the PopGenLib library, provided algorithms for the NumCalc library and helped with the GNU configuration tools. EB developed tools for the PopGenLib library. SG1 provided tools for the PopGenLib and SeqLib libraries. VR helped in the development of the PhylLib library. NG participated in the design of the library. KB supervised the whole project. SG1, VR, NG and KB helped to draft the manuscript. All authors read and approved the final manuscript.

Additional material

Additional File 1

Sources and data files of the example as a zipped archive.

Click here for file

[<http://www.biomedcentral.com/content/supplementary/1471-2105-7-188-S1.zip>]

Acknowledgements

JD would like to thank Tal Pupko for providing help with C++ in the early development of this project. Some code in the NumCalc library have been adapted from the SEMPHY library [13], Ziheng Yang's PAML package [14] and the C++ port of the Java Matrix (JAMA) package <http://math.nist.gov/tnt>. This work was supported by French Ministère de la Recherche ACI IMPBio.

References

- Thornton K: **Libsequence: a C++ class library for evolutionary genetic analysis.** *Bioinformatics* 2003, **19**:2325-7.
- Chang W, Shindyalov IN, Pu C, Bourne PE: **Design and application of PDBlib, a C++ macromolecular class library.** *Computer Applications In The Biosciences* 1994, **10**:575-86.
- Butt D, Roger AJ, Blouin C: **libcov: a C++ bioinformatic library to manipulate protein structures, sequence alignments and phylogeny.** *BMC Bioinformatics* 2005, **6**:138-138.
- Pitt WR, Williams MA, Steven M, Sweeney B, Bleasby AJ, Moss DS: **The Bioinformatics Template Library – generic components for biocomputing.** *Bioinformatics* 2001, **17**:729-37.
- van Heesch D: **Doxygen, a documentation system for C++.** 2005 [<http://www.stack.nl/~dimitri/doxygen/>].
- Jones DT, Taylor WR, Thornton JM: **The rapid generation of mutation data matrices from protein sequences.** *Computer Applications In The Biosciences* 1992, **8**:275-82.
- Kosiol C, Goldman N: **Different versions of the Dayhoff rate matrix.** *Molecular Biology And Evolution* 2005, **22**:193-9.
- Yang Z: **Maximum likelihood phylogenetic estimation from DNA sequences with variable rates over sites: approximate methods.** *Journal Of Molecular Evolution* 1994, **39**:306-14.
- Saitou N, Nei M: **The neighbor-joining method: a new method for reconstructing phylogenetic trees.** *Molecular Biology And Evolution* 1987, **4**:406-25.
- Felsenstein J: *Inferring Phylogenies* Sunderland, Massachusetts: Sinauer Associates; 2004.
- Dutheil J, Pupko T, Jean-Marie A, Galtier N: **A model-based approach for detecting coevolving positions in a molecule.** *Molecular Biology And Evolution* 2005, **22**:1919-28.
- Galtier N, Bazin E, Bierné N: **GC-biased segregation of non-coding polymorphisms in Drosophila.** *Genetics* 2005, **172**:221-8.
- Friedman N, Ninio M, Pe'er I, Pupko T: **A structural EM algorithm for phylogenetic inference.** *Journal Of Computational Biology* 2002, **9**:331-53.
- Yang Z: **PAML: a program package for phylogenetic analysis by maximum likelihood.** *Computer Applications In The Biosciences* 1997, **13**:555-6.

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

