

# BioRuby : Open-source Bioinformatics Library



Naohisa Goto<sup>1</sup>, Mitsuteru C. Nakao<sup>2</sup>, Shuichi Kawashima<sup>3</sup>, Toshiaki Katayama<sup>2</sup>, Minoru Kanehisa<sup>3</sup>

<sup>1</sup> Genome Information Research Center, Osaka University    <sup>2</sup> Human Genome Center, Institute of Medical Science, University of Tokyo

<sup>3</sup> Bioinformatics Center, Institute for Chemical Research, Kyoto University    + various contributors on the Internet

E-mail: <staff@bioruby.org>

## Summary

**BioRuby** is an open-source project which aims to provide a reusable library for biological tasks for the Ruby language. Ruby is an interpreted object-oriented scripting language with a simple and powerful syntax and native object-oriented programming support. Ruby was started by a Japanese author and is now accepted not only by Japanese but also by many professional programmers around the world as a highly productive language.

Ruby has many advantageous features to process text files and for system management tasks, which are frequently needed for bioinformatics tools. Compared to other languages, it has native support for object-oriented programming with a simple but powerful syntax, with which we can easily describe and manipulate complicated biological data structures efficiently. These are the main reason why we decided to implement a bioinformatics library in Ruby, even though BioPerl, BioJava, and BioPython were developed previously.

BioRuby project was started in late 2000, and is still in progress. Currently, in version 0.5.3 (latest release), there are over 80 files and 15,000 lines (except comment-only lines). BioRuby is available as free software and is licensed under the LGPL.

## History

(1995/12/21 Ruby language (0.95) was opened to public (fj.sources))

2000/11/21 **BioRuby** project started

2001/03/18 mailing list started

2001/06/21 bioruby-0.1

2001/07/19 ISMB/BOSC2001 lightning talks

2001/10/24 bioruby-0.3 w/ CVS repository

2002/01-02 BioHackathon w/ BioFetch server

2002/12/16 GIW2002 software demonstration

2003/01/28 bioruby-0.4.0 released

2003/02/17 BioHackathon 2003

2003/06/25 bioruby-0.5.0 released

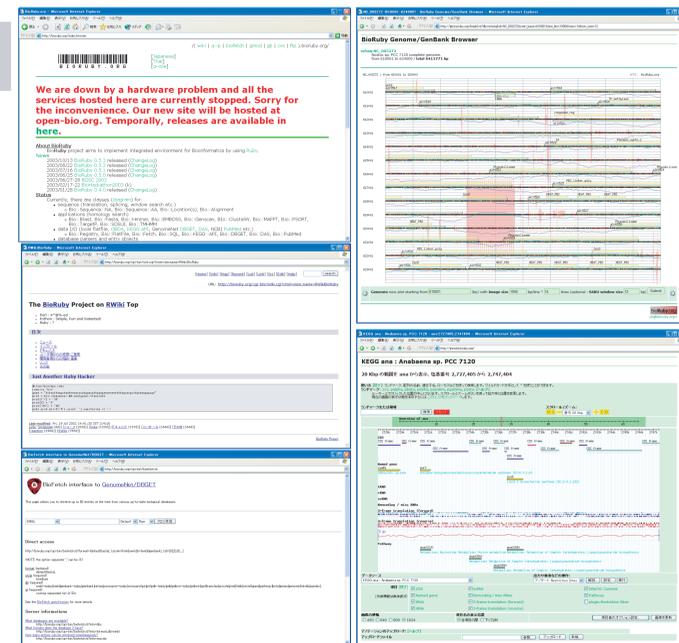
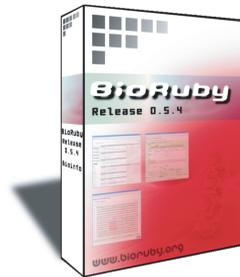
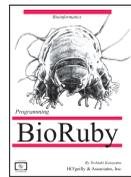
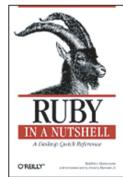
2003/06/27 ISMB/BOSC2003 talks

2003/07/16 bioruby-0.5.1 released

2003/08/22 bioruby-0.5.2 released

2003/10/13 bioruby-0.5.3 released

<http://bioruby.org/>



## Classes in BioRuby

### Basic data structures

<code>Bio::Sequence::NA</code>	Nucleic acid sequences
<code>Bio::Sequence::AA</code>	Amino acid sequences
<code>Bio::Locations</code>	Locations
<code>Bio::Features</code>	Annotations
<code>Bio::Reference</code>	Literatures
<code>Bio::Relation</code>	Graphs / Relations
<code>Bio::Pathway</code>	Pathways
<b>NEW</b> <code>Bio::Alignment</code>	Alignments

### Wrappers and parsers for bioinformatics tools

<b>updated</b> <code>Bio::Blast</code>	BLAST (similarity search)
<code>Bio::Fasta</code>	FASTA (similarity search)
<code>Bio::HMMER</code>	HMMER (similarity search)
<b>NEW</b> <code>Bio::ClustalW</code>	CLUSTAL W (multiple alignment)
<b>NEW</b> <code>Bio::MAFFT</code>	MAFFT (multiple alignment)
<b>NEW</b> <code>Bio::PSORT</code>	PSORT (protein subcellular localization)
<b>NEW</b> <code>Bio::TargetP</code>	TargetP (protein subcellular localization)
<b>NEW</b> <code>Bio::SOSUI</code>	SOSUI (transmembrane helix prediction)
<b>NEW</b> <code>Bio::TMHMM</code>	TMHMM (transmembrane helix prediction)
<b>NEW</b> <code>Bio::GenScan</code>	GenScan (gene finding)
<code>Bio::EMBOSS</code>	EMBOSS (analysis package)

### Databases and sequence file formats

<code>Bio::FastaFormat</code>	FASTA format
<code>Bio::GenBank</code>	GenBank / DDBJ
<code>Bio::EMBL</code>	EMBL
<code>Bio::SPTR</code>	SwissProt and TrEMBL
<b>NEW</b> <code>Bio::NBRF</code>	PIR
<b>NEW</b> <code>Bio::PDB</code>	Protein Data Bank
<code>Bio::PROSITE</code>	PROSITE motifs
<code>Bio::AAindex</code>	AAindex
<b>NEW</b> <code>Bio::GO</code>	Gene Ontology
<b>NEW</b> <code>Bio::GFF</code>	General Feature Format
<code>Bio::KEGG::GENES</code>	KEGG GENES
<code>Bio::KEGG::GENOME</code>	KEGG Genomes
<b>NEW</b> <code>Bio::KEGG::KO</code>	KEGG Orthology
<code>Bio::KEGG::ENZYME</code>	KEGG enzyme
<code>Bio::KEGG::COMPOUND</code>	KEGG compound
<code>Bio::KEGG::CELL</code>	KEGG CELL
<code>Bio::KEGG::Microarrays</code>	KEGG microarrays
<code>Bio::KEGG::BRITE</code>	Biomolecular Reactions
<code>Bio::LITDB</code>	Protein/peptide literature database
<code>Bio::TRANSFAC</code>	The transcription factor database
<b>NEW</b> <code>Bio::FANTOM</code>	Functional annotation of mouse
<code>Bio::MEDLINE</code>	MEDLINE bibliographic database

### File, network, and database I/O

<code>Bio::Registry</code>	OBDA Registry service
<code>Bio::SQL</code>	OBDA BioSQL RDB schema
<code>Bio::Fetch</code>	OBDA BioFetch via HTTP
<code>Bio::FlatFileIndex</code>	OBDA flat file indexing system
<code>Bio::FlatFile</code>	Flat file reader with data □□□ format autodetection
<code>Bio::PubMed</code>	NCBI PubMed service
<b>NEW</b> <code>Bio::DAS</code>	Distributed Annotation System
<b>NEW</b> <code>Bio::KEGG::API</code>	SOAP/WSDL interface in KEGG
<b>NEW</b> <code>Bio::DDBJ::XML</code>	DDBJ web services

### Command-line applications

<b>NEW</b> <code>biogetseq</code>	OBDA Registry sequence retrieval
<code>biofetch</code>	OBDA BioFetch sequence retrieval
<code>bioflat</code>	Creates/searches OBDA flat file index

### Sample / Miscellaneous

<b>NEW</b> <code>goslim.rb</code>	GO slim histogram
<b>NEW</b> <code>tdiary.rb</code>	Plug-in for tDiary

**NEW** : newly added in 2003  
**updated** : significant improvement  
(Some classes are not listed here.)

## Sample codes

### Sequence manipulation

```
#!/usr/bin/env ruby
require 'bio'
gene = Bio::Sequence::NA.new("ggatctg")
puts gene.complement[""] # -> "ccagatacc"
puts gene.composition[""] # -> {"a"=>1,"c"=>1, ...}
puts gene.gc[""] # -> 55.6
prot = gene.translate[""] # -> Bio::Sequence::AA
puts prot[""] # -> "GIW"
puts prot.molecular_weight[""] # -> 374.45
```

### Flat file database access

```
#!/usr/bin/env ruby
require 'bio'
ff = Bio::FlatFile.auto("gbbct1.seq")
ff.each do |gb|
  gb.each_gene do |f|
    pos = f.position
    title = gb.entry_id + pos
    puts gb.seq.splicing(pos).to_fasta(title)
  end
end
```

### Entry fetch by OBDA

```
#!/usr/bin/env ruby
require 'bio'
registry = Bio::Registry.new
db = registry.db("swissprot")
puts db.fetch("TETW_BUTFI")
```



### 'bioflat' command

```
% bioflat -create -location ./ \
  -dbname gbphg \
  gbphg.seq
```

### 'biofetch' command

```
% biofetch eco b0001
```

### 'biogetseq' command

```
% biogetseq -dbname eco \
  [""] b0001 b0002
```