

Biopython @  
EuroSciPy 2010



EuroSciPy  
Annual European Conference  
for Scientists using Python

Peter Cock, Plant Pathology, SCRI, Dundee, UK

EuroSciPy 2010, 3<sup>rd</sup> European meeting on Python in Science

Ecole Normale Supérieure, Paris, France, 10 July 2010

O|B|F

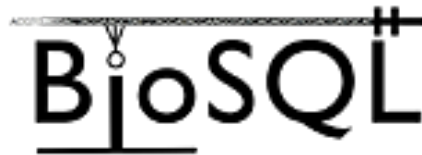


# Open Bioinformatics Foundation (OBF)



The OBF supports:

- BioPerl
- Biopython
- BioJava
- BioRuby
- BioSQL
- EMBOSS
- ...



BioPerl



BioRuby



O|B|F



- Brief introduction to Biopython & history
- Examples:
  - Sequence manipulation
  - 3D Biological structures
- Current and future projects
- Developers, git and github, ...

# Biopython



- Free, open source library for bioinformatics
- Supported by Open Bioinformatics Foundation
- Runs on Windows, Linux, Mac OS X, etc
- International team of volunteer developers
- Currently about four releases per year
- Extensive “Biopython Tutorial & Cookbook”
- See [www.biopython.org](http://www.biopython.org) for details

# Biopython's Ten Year (and a bit) History



- 1999 • Started
- 2000 • First release
- 2001 • Biopython 1.00
- ...
- 2007 • Biopython 1.43, ...
- 2008 • Biopython 1.45, ...
- 2009 • Biopython 1.50, ...
- Application note
- 2010 • Biopython 1.54, ...

## Sequence analysis

### Biopython: freely available Python tools for computational molecular biology and bioinformatics

Peter J. A. Cock<sup>1,\*</sup>, Tiago Antao<sup>2</sup>, Jeffrey T. Chang<sup>3</sup>, Brad A. Chapman<sup>4</sup>, Cymon J. Cox<sup>5</sup>, Andrew Dalke<sup>6</sup>, Iddo Friedberg<sup>7</sup>, Thomas Hamelryck<sup>8</sup>, Frank Kauff<sup>9</sup>, Bartek Wilczynski<sup>10,11</sup> and Michiel J. L. de Hoon<sup>12</sup>

<sup>1</sup>Plant Pathology, SCRI, Invergowrie, Dundee, DD2 5DA, <sup>2</sup>Liverpool School of Tropical Medicine, Liverpool, L3 5QA, UK, <sup>3</sup>Institute for Genome Sciences and Policy, Duke University Medical Center, Durham, NC, <sup>4</sup>Department of Molecular Biology, Simches Research Center, Massachusetts General Hospital, Boston, MA 02114, USA, <sup>5</sup>Centro de Ciências do Mar, Universidade do Algarve, Faro, Portugal, <sup>6</sup>Andrew Dalke Scientific, AB, Gothenburg, Sweden, <sup>7</sup>California Institute for Telecommunications and Information Technology, University of California, San Diego, 9500 Gilman Dr., La Jolla, CA 92093-0446, USA, <sup>8</sup>Bioinformatics Center, Department of Biology, University of Copenhagen, Ole Maaloes Vej 5, 2200 Copenhagen N, Denmark, <sup>9</sup>Molecular Phylogenetics, Department of Biology, TU Kaiserslautern, 67653 Kaiserslautern, UK, <sup>10</sup>EMBL Heidelberg, Meyerhofstraße 1, 69117 Heidelberg, Germany, <sup>11</sup>Institute of Informatics, University of Warsaw, Poland and <sup>12</sup>RIKEN Omics Science Center, 1-7-22 Suehiro-cho, Tsurumi-ku, Yokohama-shi, Kanagawa-ken, 230-0045, Japan

Received and revised on March 11, 2009; accepted on March 16, 2009

Advance Access publication March 20, 2009

Associate Editor: Dmitri Frishman

#### ABSTRACT

**Summary:** The Biopython project is a mature open source international collaboration of volunteer developers, providing Python libraries for a wide range of bioinformatics problems. Biopython includes modules for reading and writing different sequence file formats and multiple sequence alignments, dealing with 3D macromolecular structures, interacting with common tools such as BLAST, ClustalW and EMBOSS, accessing key online databases, as well as providing numerical methods for statistical learning.

**Availability:** Biopython is freely available, with documentation and source code at [www.biopython.org](http://www.biopython.org) under the Biopython license.

**Contact:** All queries should be directed to the Biopython mailing lists, see [www.biopython.org/wiki/Mailing\\_lists](http://www.biopython.org/wiki/Mailing_lists); [peter.cock@scri.ac.uk](mailto:peter.cock@scri.ac.uk).

#### 1 INTRODUCTION

Python ([www.python.org](http://www.python.org)) and Biopython are freely available open source tools, available for all the major operating systems. Python is a very high-level programming language, in widespread commercial and academic use. It features an easy to learn syntax, object-oriented programming capabilities and a wide array of libraries. Python can interface to optimized code written in C, C++ or even FORTRAN, and together with the Numerical Python project *numpy* (Oliphant, 2006), makes a good choice for scientific programming (Oliphant, 2007). Python has even been used in the numerically demanding field of molecular dynamics (Hinsen, 2000). There are also high-quality plotting libraries such as *matplotlib* ([matplotlib.sourceforge.net](http://matplotlib.sourceforge.net)) available.

\*To whom correspondence should be addressed.

© 2009 The Author(s)

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

Since its founding in 1999 (Chapman and Chang, 2000), Biopython has grown into a large collection of modules, described briefly below, intended for computational biology or bioinformatics programmers to use in scripts or incorporate into their own software. Our web site lists over 100 publications using or citing Biopython.

The Open Bioinformatics Foundation (OBF, [www.open-bio.org](http://www.open-bio.org)) hosts our web site, source code repository, bug tracking database and email mailing lists, and also supports the related *BioPerl* (Stajich *et al.*, 2002), *BioJava* (Holland *et al.*, 2008), *BioRuby* ([www.bioruby.org](http://www.bioruby.org)) and *BioSQL* ([www.biosql.org](http://www.biosql.org)) projects.

#### 2 BIOPYTHON FEATURES

The *Seq* object is Biopython's core sequence representation. It behaves very much like a Python string but with the addition of an alphabet (allowing explicit declaration of a protein sequence for example) and some key biologically relevant methods. For example,

```
>>> from Bio.Seq import Seq
>>> from Bio.Alphabet import generic_dna
>>> gene = Seq("ATGAAAGCAATTTTCGTACTG",
...           "AAAGGTTGGTGCCCACTTGA",
...           generic_dna)
>>> print gene.transcribe()
AUGAAGCCAAUUGUCUACUCGAAAGCUUGUUGUCCOCCACUUGA
>>> print gene.translate(table=11)
MKALFVLKGVWRT*
```

Sequence annotation is represented using *SeqRecord* objects which augment a *Seq* object with properties such as the record name, identifier and description and space for additional key/value terms. The *SeqRecord* can also hold a list of *SeqFeature*

# Examples

---



# Sequence vs Sequence



- In biology the word “sequence” generally means an ordered collection of letters representing a directed molecular chain

DNA usually A, C, G and T

RNA usually A, C, G and U

Proteins usually 20 single letter codes

- Python strings are often a good model
- Biopython has a Seq object...

# String like methods for Seq objects



- Seq has an alphabet (DNA, RNA or Protein)

```
>>> from Bio.Seq import Seq
>>> from Bio.Alphabet import generic_dna
>>> dna = Seq("GATCGATGGGCCTATATAGGATCGAAAATCGC", generic_dna)
>>> print dna, dna.alphabet
GATCGATGGGCCTATATAGGATCGAAAATCGC DNAAlphabet()
```

```
>>> len(dna)
32
>>> dna.count('C')
6
>>> dna.find("TATAT")
12
>>> print dna[:12] + "-----" + dna[17:]
GATCGATGGGCC-----AGGATCGAAAATCGC
```

```
>>> print dna.lower()
gatcgatgggcctatataggatcgaaaatcgc
```

Explicit declaration  
of the alphabet  
(sequence type).



# Biological methods for sequences



- DNA to RNA to Protein - “The Central Dogma”

```
>>> from Bio.Seq import Seq
>>> from Bio.Alphabet import generic_dna
>>> dna = Seq("GATCGATGGGCCTATATAGGATCGAAAATCGC", generic_dna)
>>> print dna, dna.alphabet
GATCGATGGGCCTATATAGGATCGAAAATCGC DNAAlphabet()
```

```
>>> print dna.complement()
CTAGCTACCCGGATATATCCTAGCTTTTAGCG
>>> print dna.reverse_complement()
GCGATTTTCGATCCTATATAGGCCCATCGATC
```

```
>>> rna = dna.transcribe()
>>> print rna, rna.alphabet
GAUCGAUGGGCCUAUAUAGGAUCGAAAUCGC RNAAlphabet()
```

```
>>> protein = rna.translate()
>>> print protein, protein.alphabet
DRWAYIGSKI ExtendedIUPACProtein()
```

# Sequence File Manipulation

---



- Manipulating nucleotide and protein sequences is a common task in Bioinformatics
- Manipulating plain text sequence files is too
- There are *lots* of different file formats ☹️
- Motivation for common object and API

# Reading a FASTA file with Bio.SeqIO



```
>FL3B07415JACDX
TTAATTTTATTTTGTCCGGCTAAAGAGATTTTGTAGCTAAACGTTCAATTGCTTTAGCTGAA
GTACGAGCAGATACTCCAATCGCAATTGTTTCTTCATTTAAAATTAGCTCGTCGCCACCT
TCAATTGGAAATTTATAATCACGATCTAACCAGATTGGTACATTATGTTTTGCAAATCTT
GGATGATATTTAATGATGTAATCCATGAATAATGATTCACGTCTACGCGCTGGTTCTCTC
ATCTTATTTATCGTTAAGCCA
>FL3B07415I7AFR
CATTAATAA...
```

```
from Bio import SeqIO
for rec in SeqIO.parse("phage.fasta", "fasta") :
    print rec.id, len(rec.seq), rec.seq[:10]+"..."
```

```
FL3B07415JACDX 261 TTAATTTTAT...
FL3B07415I7AFR 267 CATTAATAA...
FL3B07415JCAY5 136 TTTCTTTTCT...
FL3B07415JB41R 208 CTCTTTTATG...
FL3B07415I6HKB 268 GGTATTTGAA...
FL3B07415I63UC 219 AACATGTGAG...
...
```

Focus on the filename and format ("fasta")...

# Reading a FASTQ file with Bio.SeqIO



```
@FL3B07415JACDX  
TTAATTTTATTTTGTGGCTAAAGAGATTTTGTAGCTAAACGTTCAATTGCTTTAGCTGAAGTACGAGCAGATACTCCAATCGCAATTGTTTCTTC  
ATTTAAAATTAGCTCGTCGCCACCTTCAATTGGAATTTATAATCACGATCTAACCAGATTGGTACATTATGTTTTGCAAATCTTGATGATATT  
TAATGATGTACTCCATGAATAATGATTCACGTCTACGCGCTGGTTCCTCATCTTATTTATCGTTAAGCCA  
+  
BBBB2262=1111FFGGGHHHHIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII  
FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFGGGGFFFFFFFFFFFFFFFFFFFF  
BBCFFFFFFFFFFFFFFFFFFFFFFFFGGGGGGGIIIIIIIGGGIIIGGGIIGGG@AAAAA?===@@@???  
@FL3B07415I7AFR  
CATTAATA...
```

```
from Bio import SeqIO  
  
for rec in SeqIO.parse("phage.fastq", "fastq") :  
    print rec.id, len(rec.seq), rec.seq[:10]+"..."  
    print rec.letter_annotations["phred_quality"][:10], "..."
```

```
FL3B07415JACDX 261 TTAATTTTAT...  
[33, 33, 33, 33, 17, 17, 21, 17, 28, 16] ...  
FL3B07415I7AFR 267 CATTAATA...  
[37, 37, 37, 37, 37, 37, 37, 37, 38, 38] ...  
FL3B07415JCA5 136 TTTCTTTTCT...  
[37, 37, 36, 36, 29, 29, 29, 29, 36, 37] ...  
FL3B07415JB41R 208 CTCTTTTATG...  
[37, 37, 37, 38, 38, 38, 38, 38, 37, 37] ...  
FL3B07415I6HKB 268 GGTATTTGAA...  
[37, 37, 37, 37, 34, 34, 34, 37, 37, 37] ...  
FL3B07415I63UC 219 AACATGTGAG...  
[37, 37, 37, 37, 37, 37, 37, 37, 37, 37] ...  
...
```

Just filename and  
format changed  
("fasta" to "fastq")

# Sequence File Manipulation

---



- Common object for sequence file entries, SeqRecord, for Seq plus annotation like ID
- Sequence file API based on iterators
- Memory efficient!
- Scales to millions of reads as seen in current sequencing platforms (Roche, Illumina, etc)

# Trimming a FASTA file with Bio.SeqIO



```
>FL3B07415JACDX
TTAATTTTATTTTGTCCGGCTAAAGAGATTTTGTAGCTAAACGTTCAATTGCTTTAGCTGAA
GTACGAGCAGATACTCCAATCGCAATTGTTTCTTCATTTAAAATTAGCTCGTCGCCACCT
TCAATTGGAAATTTATAATCACGATCTAACCAGATTGGTACATTATGTTTTGCAAATCTT
GGATGATATTTAATGATGTAATCCATGAATAATGATTCACGTCTACGCGCTGGTTCTCTC
ATCTTATTTATCGTTAAGCCA
>FL3B07415I7AFR
CATTAATAA...
```

```
from Bio import SeqIO
```

```
recs = (r[:10] for r in SeqIO.parse("phage.fasta", "fasta"))
```

```
SeqIO.write(recs, "long.fasta", "fasta")
```

```
>FL3B07415JACDX
TTAATTTTAT
>FL3B07415I7AFR
CATTAATAA
>FL3B07415JCAY5
TTTCTTTTCT
>FL3B07415JB41R
CTCTTTTATG
...
```

Generator  
expression

# Filtering a FASTA file with Bio.SeqIO



```
>FL3B07415JACDX  
TTAATTTTATTTTGTCTGGCTAAAGAGATTTTTAGCTAAACGTTCAATTGCTTTAGCTGAA  
GTACGAGCAGATACTCCAATCGCAATTGTTTCTTCATTTAAAATTAGCTCGTCGCCACCT  
TCAATTGGAAATTTATAATCACGATCTAACCAGATTGGTACATTATGTTTTGCAAATCTT  
GGATGATATTTAATGATGTAATCATGAATAATGATTCACGTCTACGCGCTGGTTCTCTC  
ATCTTATTTATCGTTAAGCCA  
>FL3B07415I7AFR  
CATTAATAA...
```

```
from Bio import SeqIO
```

```
recs = (r for r in SeqIO.parse("phage.fasta", "fasta") if len(r)>200)
```

```
SeqIO.write(recs, "long.fasta", "fasta")
```

```
>FL3B07415JACDX  
TTAATTTTATTTTGTCTGGCTAAAGAGATTTTTAGCTAAACGTTCAATTGCTTTAGCTGAA  
GTACGAGCAGATACTCCAATCGCAATTGTTTCTTCATTTAAAATTAGCTCGTCGCCACCT  
TCAATTGGAAATTTATAATCACGATCTAACCAGATTGGTACATTATGTTTTGCAAATCTT  
GGATGATATTTAATGATGTAATCATGAATAATGATTCACGTCTACGCGCTGGTTCTCTC  
ATCTTATTTATCGTTAAGCCA  
>FL3B07415I7AFR  
CATTAATAA...
```

Generator  
expression

# Filtering and converting FASTQ to FASTA



```
@FL3B07415JACDX
TTAATTTTATTTTGTCTGGCTAAAGAGATTTTGTAGCTAAACGTTCAATTGCTTTAGCTGAAGTACGAGCAGATACTCCAATCGCAATTGTTTCTTC
ATTTAAAATTAGCTCGTCGCCACCTTCAATTGGAAATTTATAATCACGATCTAACCAGATTGGTACATTATGTTTTGCAAATCTTGGATGATATT
TAATGATGTACTCCATGAATAATGATTCACGTCTACGCGCTGGTTCTCTCATCTTATTTATCGTTAAGCCA
+
BBBB2262=1111FFGGGHHHHIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII
FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF
FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFGGGFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFGGGGFFFFFFFFFFFFFFFFFFFFFFFFGG
BBCFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFGGGGGGGIIIIIIIGGGIIIGGGIIGGG@AAAAA?===@@@???
```

```
from Bio import SeqIO

recs = (r for r in SeqIO.parse("phage.fastq", "fastq") if len(r)>200)

SeqIO.write(recs, "long.fasta", "fasta")
```

Generator expression

```
>FL3B07415JACDX
TTAATTTTATTTTGTCTGGCTAAAGAGATTTTGTAGCTAAACGTTCAATTGCTTTAGCTGAA
GTACGAGCAGATACTCCAATCGCAATTGTTTCTTCATTTAAAATTAGCTCGTCGCCACCT
TCAATTGGAAATTTATAATCACGATCTAACCAGATTGGTACATTATGTTTTGCAAATCTT
GGATGATATTTAATGATGTACTCCATGAATAATGATTCACGTCTACGCGCTGGTTCTCTC
ATCTTATTTATCGTTAAGCCA
>FL3B07415I7AFR
CATTAECTAA...
```



# General sequence file conversion



- Separate parse and write calls (as before):

```
from Bio import SeqIO
recs = SeqIO.parse("roche.sff", "sff")
SeqIO.write(recs, "reads.fastq", "fastq")
```

- Shorthand convert call (for the typical case):

```
from Bio import SeqIO
SeqIO.convert("roche.sff", "sff", "reads.fastq", "fastq")
```

- Simple to switch file formats:

```
from Bio import SeqIO
SeqIO.convert("roche.sff", "sff", "phage.fasta", "fasta")
```

- Some conversions are optimized

# Analysing a FASTA file with Bio.SeqIO



```
>>> from Bio import SeqIO
>>> sizes = [len(r) for r in SeqIO.parse("ls_orchid.fasta", "fasta")]
>>> len(sizes), min(sizes), max(sizes)
(94, 572, 789)
>>> sizes
[740, 753, 748, 744, 733, 718, 730, 704, 740, 709, 700, 726, ..., 592]
```

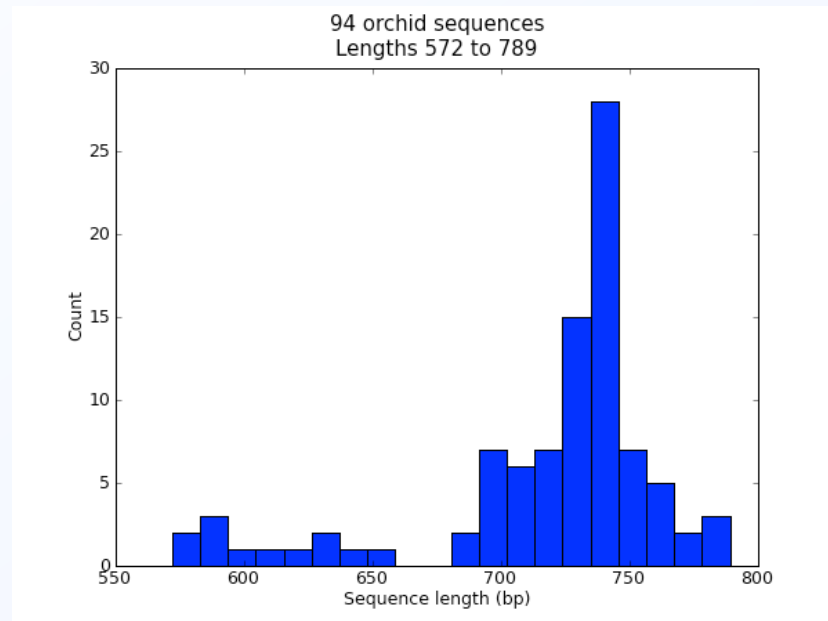
List comprehension

# Analysing a FASTA file with Bio.SeqIO



```
from Bio import SeqIO
sizes = [len(r) for r in SeqIO.parse("ls_orchid.fasta", "fasta")]

import pylab
pylab.hist(sizes, bins=20)
pylab.title("%i orchid sequences\nLengths %i to %i" \
            % (len(sizes),min(sizes),max(sizes)))
pylab.xlabel("Sequence length (bp)")
pylab.ylabel("Count")
pylab.show()
```



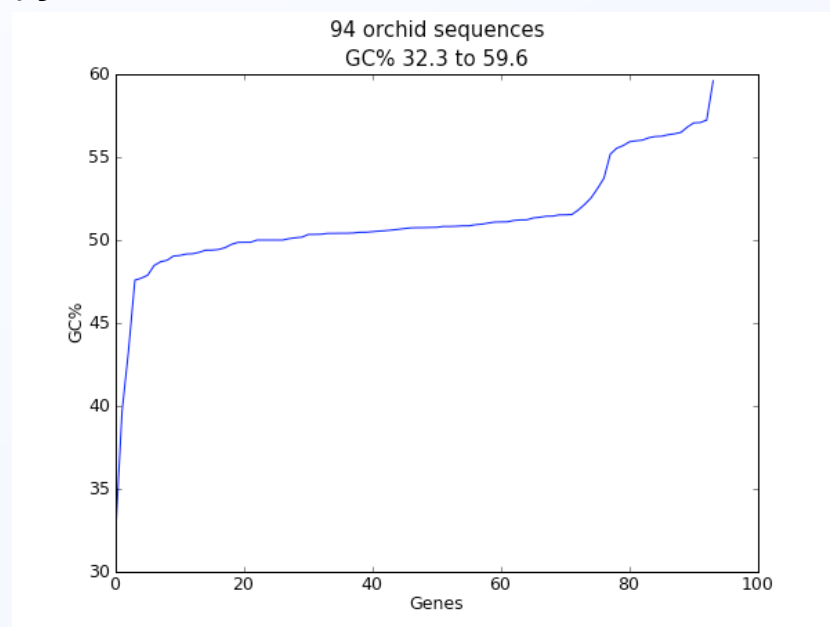
Plot with pylab  
(aka matplotlib)

# Analysing a FASTA file with Bio.SeqIO



```
from Bio import SeqIO
from Bio.SeqUtils import GC
val = sorted(GC(r.seq) for r in SeqIO.parse("ls_orchid.fasta", "fasta"))

import pylab
pylab.plot(val)
pylab.title("%i orchid sequences\nGC% %0.1f to %0.1f" \
            % (len(val), min(val), max(val)))
pylab.xlabel("Genes")
pylab.ylabel("GC%")
pylab.show()
```



Calculate percentage of DNA sequence using the letters G or C (biologically important)

# Querying online database – e.g. NCBI



```
>>> from Bio import Entrez
>>> Entrez.email = "A.N.Other@example.com"      # Tell NCBI who you are
>>> record = Entrez.read(Entrez.egquery(term="biopython"))
>>> for row in record["eGQueryResult"]:
...     print row["DbName"], row["Count"]
...
pubmed 10
pmc 109
journals 0
mesh 0
books 0
omim 0
omia 0
ncbisearch 0
nuccore 0
...
```

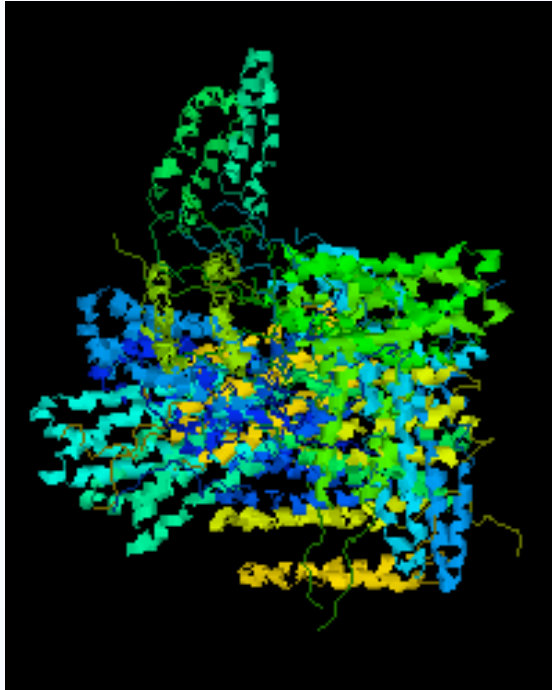
Using NCBI  
XML parser

numpy gets 74 hits in PubMedCentral,  
scipy gets 5 in PubMed and 91 in PMC

# Manipulating 3D Biological structures



- Spatial alignment (using NumPy internally)

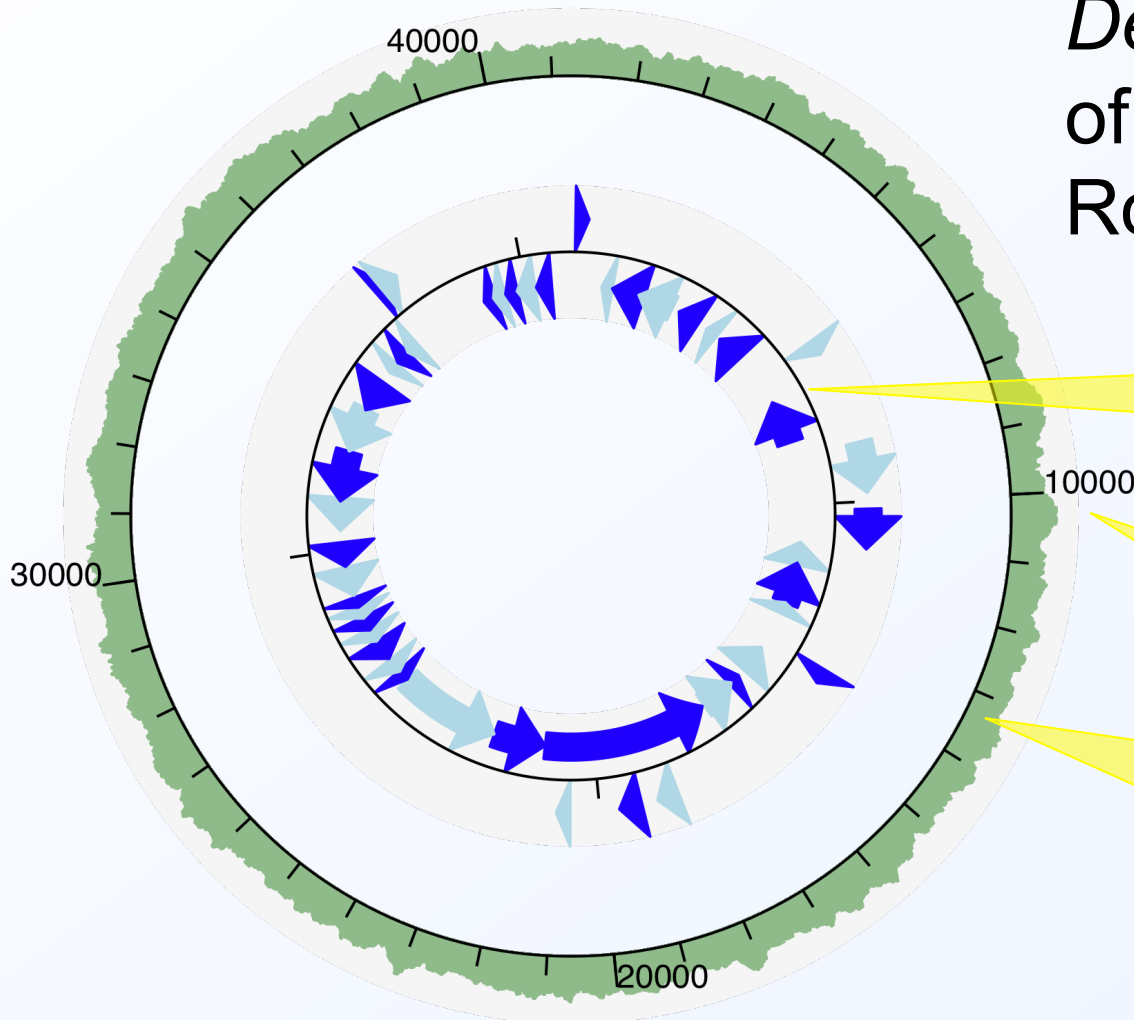


- See [http://www.warwick.ac.uk/go/peter\\_cock/python/protein\\_superposition/](http://www.warwick.ac.uk/go/peter_cock/python/protein_superposition/)
- Visualisation using OpenRasMol

# Circular GenomeDiagram



*De novo* assembly  
of 42kb phage from  
Roche 454 data



“Feature Track”  
showing ORFs

Scale tick marks

“Barchart Track” of  
read depth (~100,  
scale max 200)

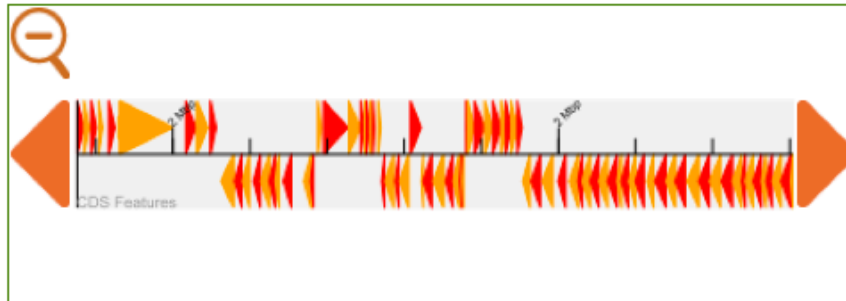
# Linear - GenomeDiagram



NC\_000913.2



Overview of region 2037589 to 2130379 (alternatively, [view in GBrowse](#)):



Escherichia coli K12, complete genome.  
Sequence length 4639675 bp.



Screenshot from an in-house web server using:

- Biopython
- BioSQL
- ReportLab
- SQLAlchemy
- Turbogears



# Other functionality not discussed

---



- Calling and parsing BLAST (local and online)
- Call command line tools (e.g. clustalw)
- Restriction enzymes
- Multiple Sequence Alignments
- Clustering (Bio.Cluster)
- Phylogenetics (Bio.Phylo, Bio.Nexus)
- BioSQL support (common schema)
- Population genetics (Bio.PopGen)
- ...

# Current and Future Work

---



- Python 3 support (now NumPy is almost there)
- Google Summer of Code 2009:
  - Nick Matzkes, Biogeography
  - (Erik Talevich , phyloXML, already merged)
- Google Summer of Code 2010:
  - João Rodrigues, extending Bio.PDB module
- Lots of other stuff!

# Development



- Moved from CVS to git a year ago
- Hosted on github.com at <http://github.com/biopython/biopython>
- Over 50 people have made a branch
- New features are now routinely developed on public branches
- Still work from a main stable branch

# What do I personally use (Bio)python for?

---



- Scripting command line tools
- Basic sequence manipulation
- Preparing input files for genome assembly
- Analysis of genome assembly coverage etc
- Working with gene annotation
- Visualising genomic information
- Calling R scripts with rpy or rpy2
- ...

# Acknowledgements



- Other Biopython contributors & developers!
- Open Bioinformatics Foundation (OBF) supports Biopython (and BioPerl etc)
- My Biopython work was/is supported by:

O|B|F

- EPSRC funded PhD (MOAC DTC, University of Warwick, UK)
- SCRI (Scottish Crop Research Institute), who also paid my conference fees and travel to be here

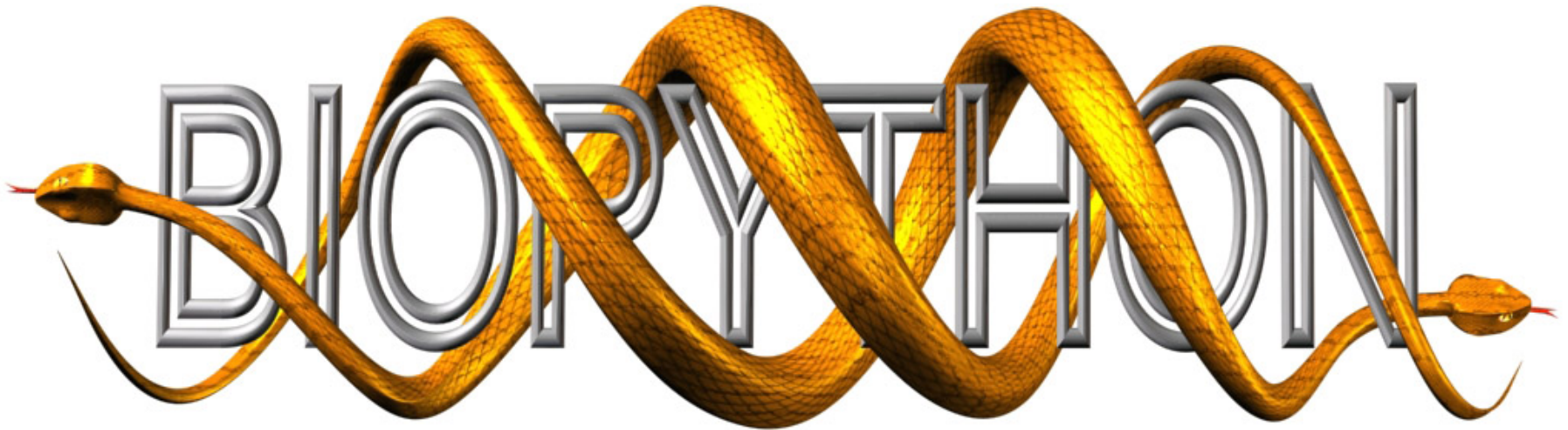
EPSRC



## What next?



- Sign up to our mailing list?



- Homepage [www.biopython.org](http://www.biopython.org)