

BIOPERL: DEVELOPING OPEN SOURCE SOFTWARE

JASON STAJICH
DUKE UNIVERSITY

WHY TALK ABOUT BIOPERL?

- Successful open-source project
- Bioinformatics is a difficult field to straddle
- Toolkit still has (many) flaws
- An evolving project

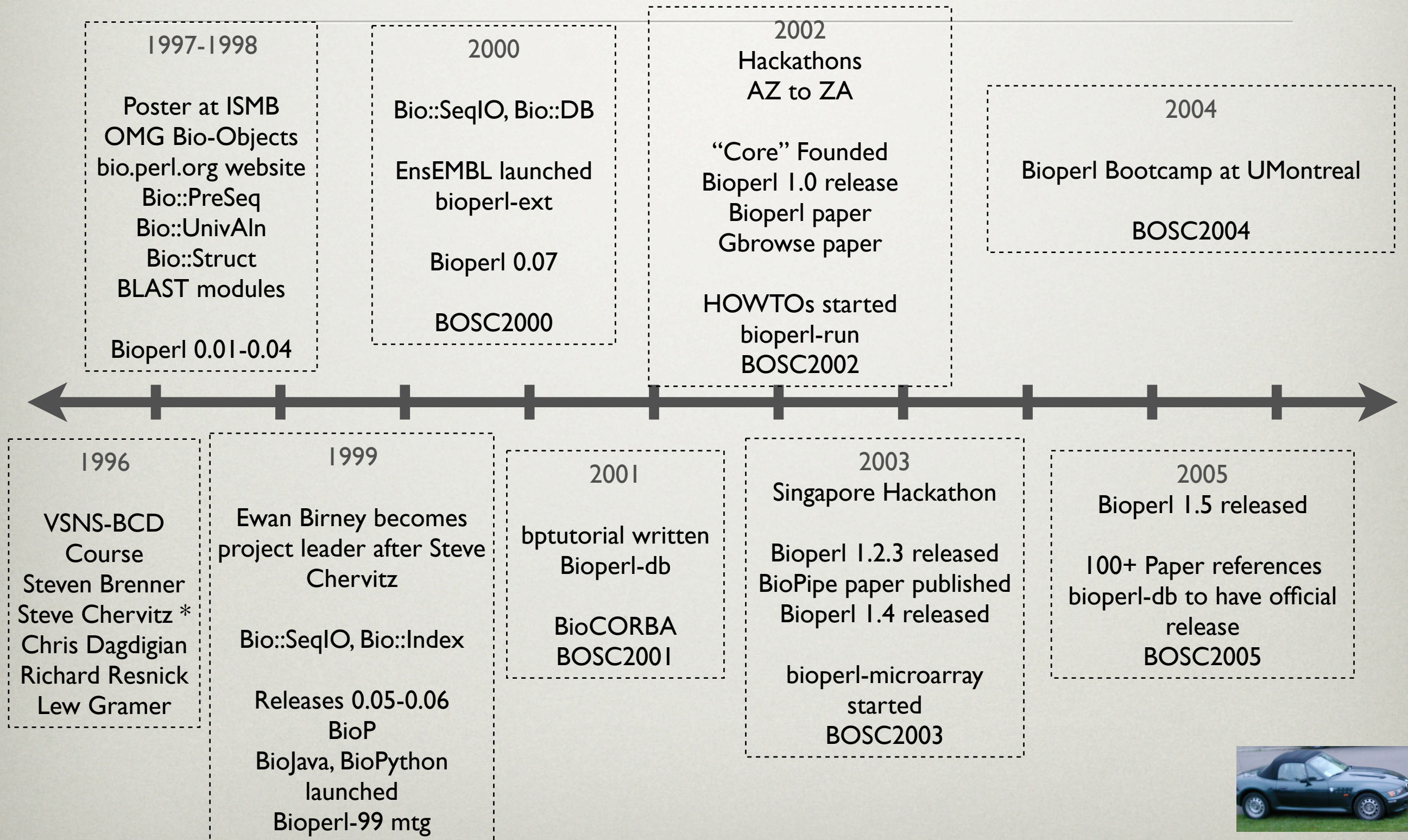
HOW DID WE GET HERE?

- A Brief History
 - Development Stats
 - People
- Bioperl in Action

BIOPERL DEVELOPMENT GOALS

- Provide a useable Perl toolkit for life science data
- Solve problems that developers need
- Avoid the 1-off “everyone write a BLAST | GenBank | FASTA parser”
- Continuity of solutions

A BRIEF HISTORY



BOSC 2005

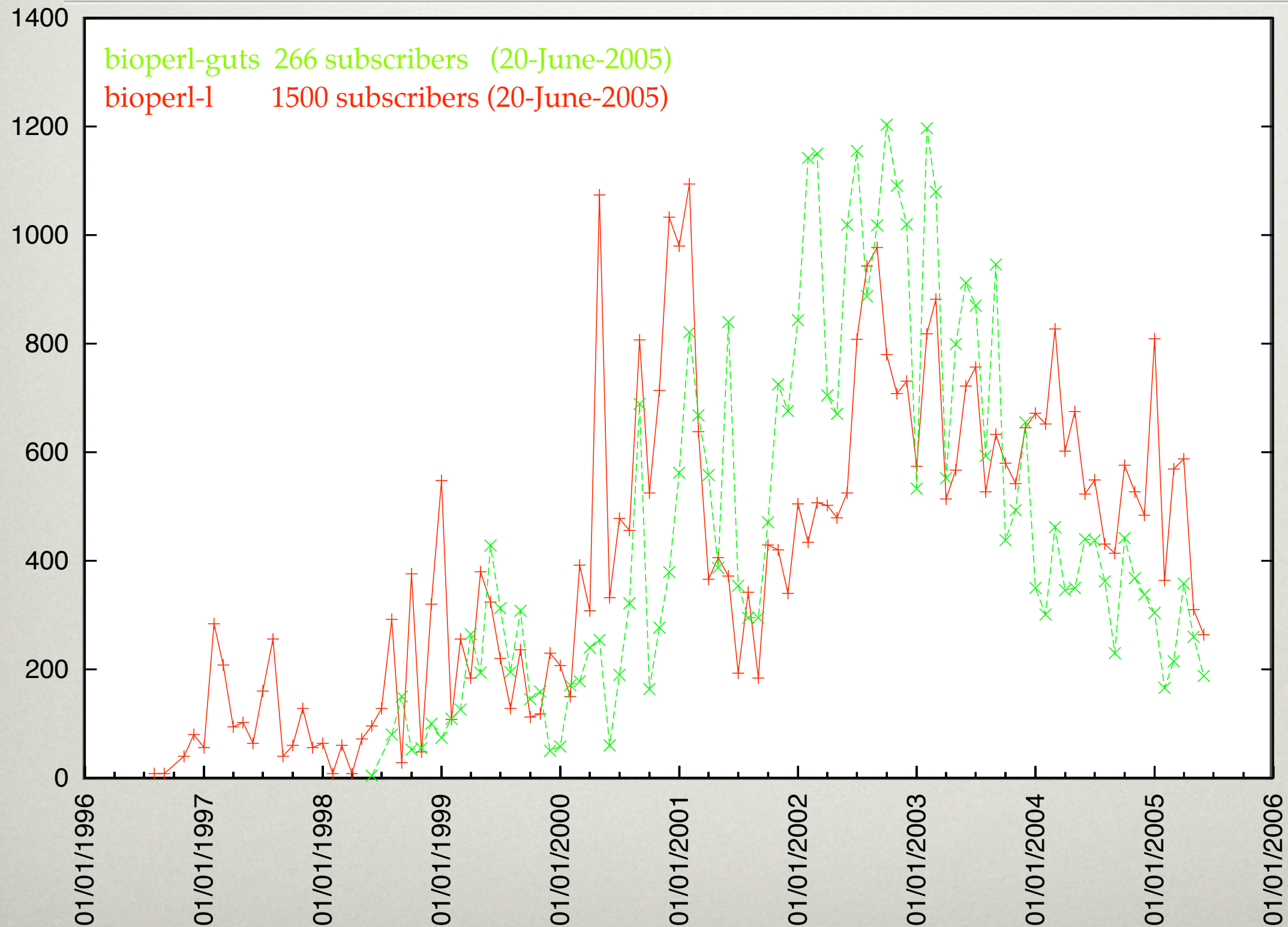
DEFINING POINTS

- Vision of module(s) set by code owner
- She / He who writes it, “wins” argument
- Project leader / core developers
 - can remove code before a release to insure a working release
 - Enforce code and module continuity
- Preserve API as much as possible

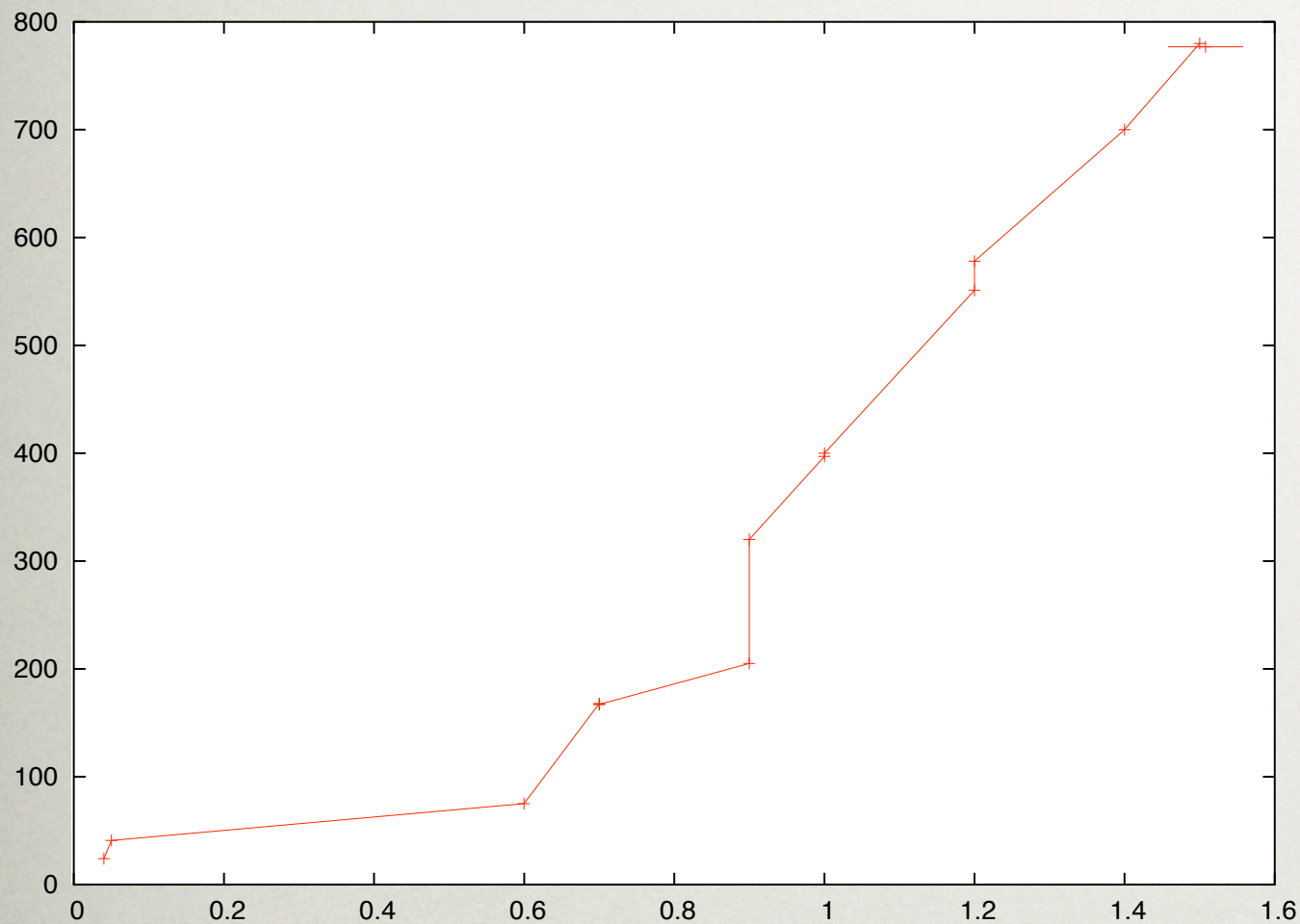
PROJECT MECHANICS

- Code is on a shared server (O | B | F)
- Publicly accessible
- Write-access to repository granted by main developers
- Discussion of code, ideas on-list

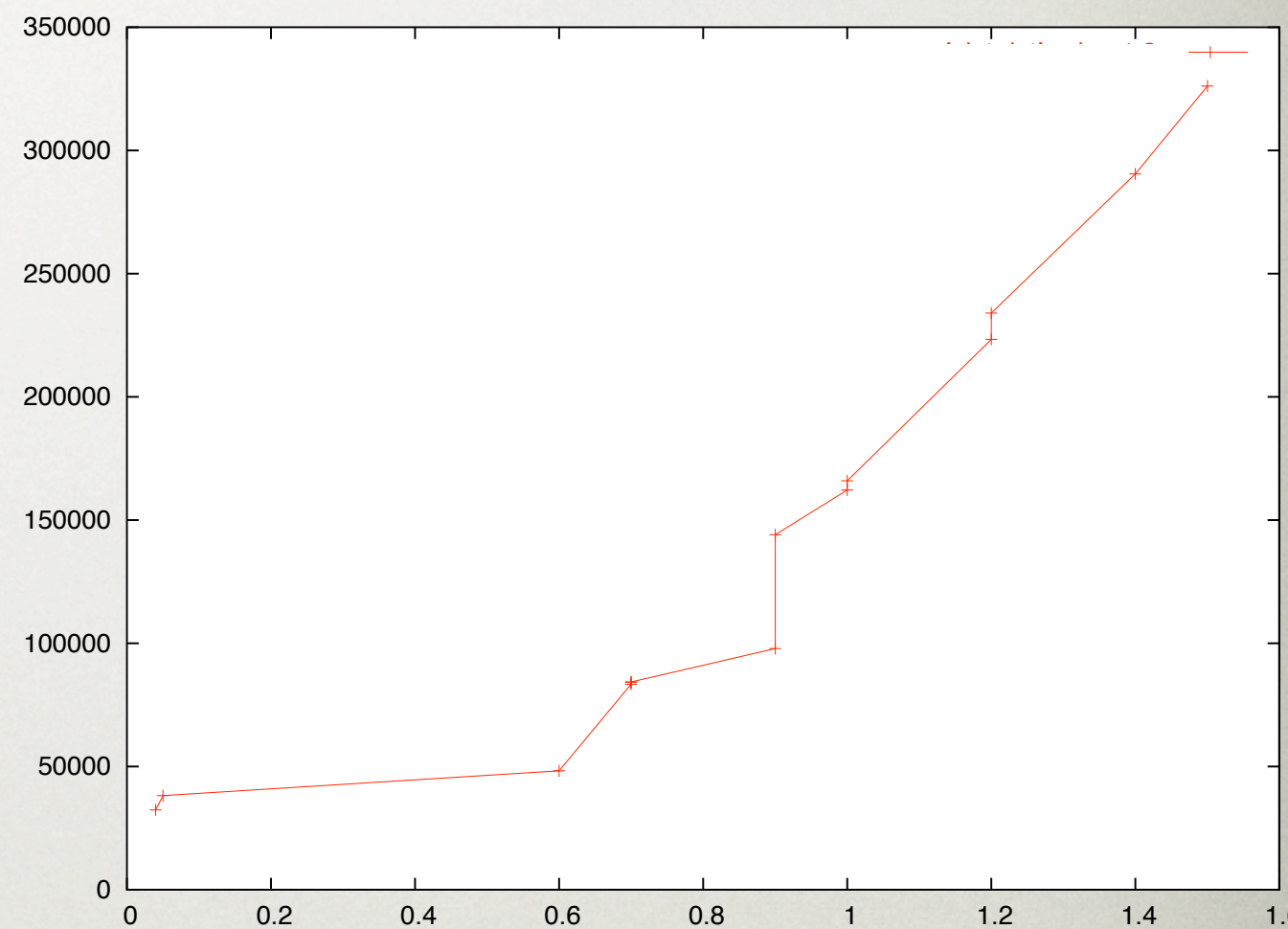
BIOPERL MAILING LIST TRAFFIC



CODE GROWTH

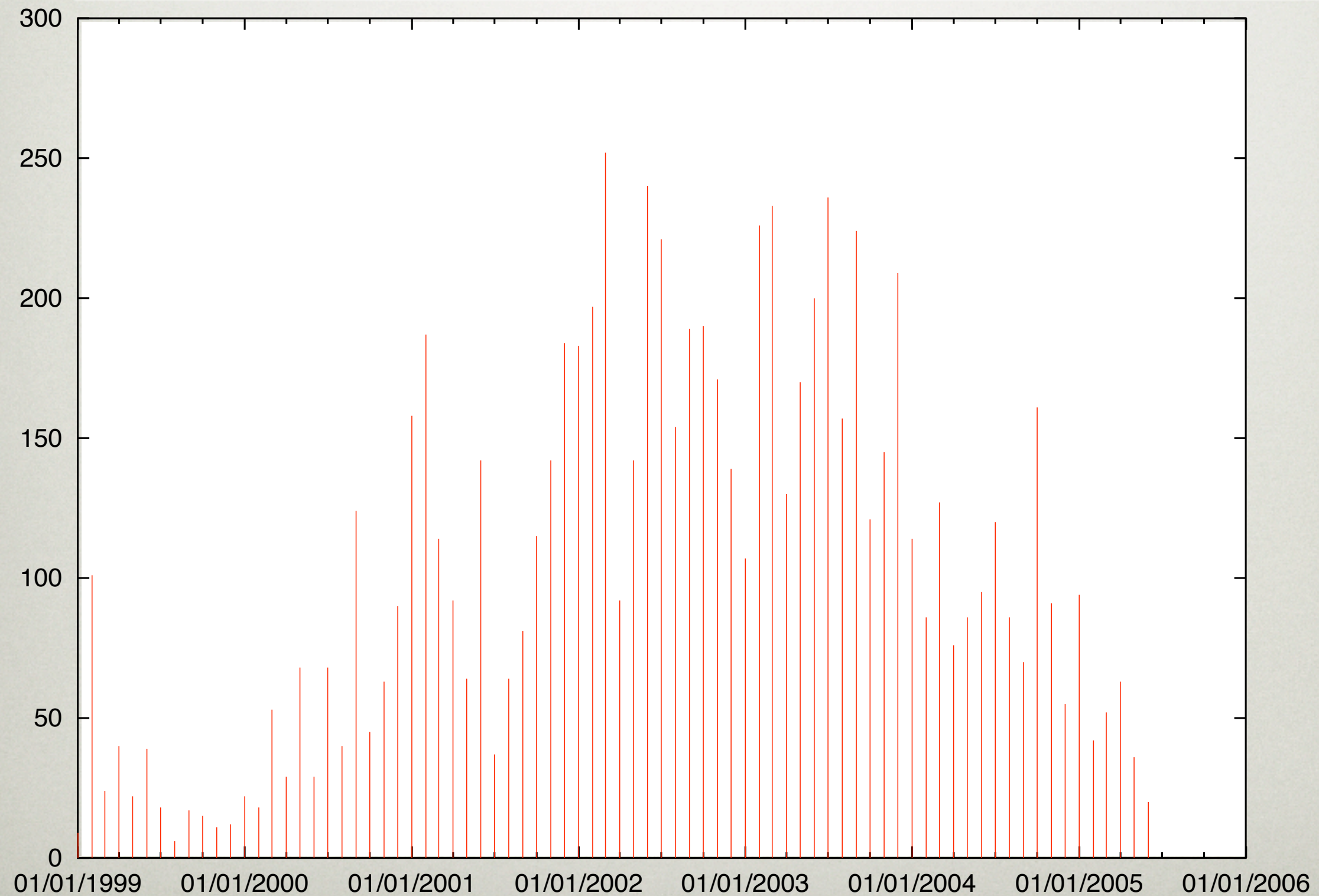


Modules per release



Lines per release

CVS COMMITS



BIOPERL IN ACTION

```
# Convert sequence formats
use Bio::SeqIO;
my $in = Bio::SeqIO->new(-format => 'genbank', -file=> 'file.gbk');
my $out = Bio::SeqIO->new(-format => 'fasta', -file => '>file.fas');

while( my $seq = $in->next_seq ) {
    $out->write_seq($seq);
}

-----

use Bio::SearchIO;
my $in = Bio::SearchIO->new(-format => 'blast',
                           -file=> 'result.blastp');
while( my $r = $in->next_result ) {
    while( my $h = $r->next_hit ) {
        while( my $hsp = $h->next_hsp ) {
            print $hsp->query->start, "..", $hsp->query->end, "\n";
        }
    }
}
```


BIOPERL IN ACTION

```
# Convert sequence formats
use Bio::SeqIO;
my $in = Bio::SeqIO->new(-format => 'embl', -file=> 'file.embl');
my $out = Bio::SeqIO->new(-format => 'gcg', -file => '>file.gcg');

while( my $seq = $in->next_seq ) {
    $out->write_seq($seq);
}
```



```
-----

use Bio::SearchIO;
my $in = Bio::SearchIO->new(-format => 'fasta',
                           -file    => 'result.fastp');
while( my $r = $in->next_result ) {
    while( my $h = $r->next_hit ) {
        while( my $hsp = $h->next_hsp ) {
            print $hsp->query->start, "...", $hsp->query->end, "\n";
        }
    }
}
```



MEASURING SUCCESS

- 2002 Paper has 100+ Citations
- Bits and pieces used in many informatics tools
 - Bio::SearchIO (Rfam, In-Paranoid)
 - Generic Genome Browser (Stein et al, 2002)
 - Comparative Genomics Library (Yandell, Mungall, et al)
 - [EnsEMBL]

BIOPERL TAUGHT AT TUTORIALS AND COURSES



Bootcamp 2004
(Montreal)



CSHL Bioinformatics
Software Courses

Various mini courses
MIT, Duke, EBI, Pasteur

KNOWLEDGE OF BIOPERL IS A MARKETABLE SKILL

<http://www.ce.com/education/Bioinformatics-Certificate-Program-10082109.htm>

Track 1: Suggested electives for computer scientists and IT professionals

- * Advanced Sequence Analysis in Bioinformatics (2 units)
- * Gene Expression and Pathways (2 units)
- * Protein Structure Analysis in Bioinformatics (2 units)
- * Design and Implementation of Bioinformatics Infrastructures (3 units)
- * DNA Microarrays - Principles, Applications and Data Analysis (2 Units)
- * Parallel and Distributed Computing for Bioinformatics (2 units)

Track 2: Suggested electives for molecular biologists and other scientists

- * Introduction to Programming for Bioinformatics II*** (3 units)
- * Introduction to Programming for Bioinformatics III (3 units)
- * **BioPerl for Bioinformatics** (2 units)
- * Design and Implementation of Bioinformatics Infrastructures (3 units)
- * Gene Expression and Pathways (2 units)
- * Protein Structure Analysis in Bioinformatics (2 units)
- * DNA Microarrays - Principles, Applications and Data Analysis (2 units)

<http://www.foothill.fhda.edu/bio/programs/bioinfo/curric.shtml>

CAREER CERTIFICATE REQUIREMENTS (49 units)*

Biotechnology Core Courses (14 units)
BTEC 51A Cell Biology for Biotechnology (3 units)
BTEC 52A Molecular Biology for Biotechnology (3 units)
BTEC 65 DNA Electrophoretic Systems (1 unit)
BTEC 68 Polymerase Chain Reaction (1 unit)
BTEC 71 DNA Sequencing & Bioinformatics (1 unit)
BTEC 76 Introduction to Microarray Data Analysis (2 unit)
BTEC 64 Protein Electrophoretic Systems (1 unit)
BTEC 66 HPLC (2 units)

Computer Science Core Courses (30 units)
CIS 52A Introduction to Data Management Systems (5 units)
CIS 52B2 Introduction to Oracle SQL (5 units)
CIS 68A Introduction to UNIX (5 units)
CIS 68E Introduction to PERL (5 units)
CIS 68H **Introduction to BioPerl** (5 units)
COIN81 Bioinformatics Tools & Databases (5 units)

http://bioag.byu.edu/botany/homepage/botweb/jobs_Ph.D.htm

Academic Facilities Coordinator II (Facilities)
Bioinformatics Position at the UCR Genomics Institute
UC, Riverside

The Center for Plant Cell Biology (CEPCEB) in the Genomics Institute of the University of California, Riverside, invites applicants for an Academic Facilities Coordinator II position, an academic-track 11-month appointment. Salary for the position is commensurate with education and experience.

The successful applicant will be expected to organize a small bioinformatics team to provide support to The Center for Plant Cell Biology. This team will implement currently available bioinformatics tools including relational database support and will develop user-specific data-mining tools. The appointee will be expected to develop research collaborations with the faculty and teach or organize short courses that will inform that will inform the local community about the available bioinformatics resources. Applicants must have a Ph.D. in the Biological Sciences (Plant Biology is preferred). Applicants with experience in leading a bioinformatics group will be given preference. Additionally, the applicant must be proficient in one or more programming languages (PERL, PYTHON, JAVA, C++) and have a good understanding of database design and implementation. In addition, applicants should have experience using **one of the open source bioinformatics frameworks such as BIOJAVA or BIOPERL**. The applicant should have experience with software collaboration tools such as CVS. The applicant will be expected to oversee the purchase, installation, and management of the necessary computer hardware and software required to provide bioinformatics support to several users. A good understanding of the UNIX operating system and systems administration is also an important qualification.

<http://careers.psgs.com/CareerOppLocation2.asp?location=BETHESDA>

CF-046: Bioinformatics Specialist

The NIAID Office of Technology and Information Systems (OTIS)/Bioinformatics and Scientific IT Program (BSIP) is seeking a bioinformatics specialist. The position includes managing our bioinformatics services and applications, training NIH scientists, creating and modifying simple bioinformatics software, and collaborating with NIH scientists on specific projects. [snip]

The qualified candidate must hold a Master's Degree (or equivalent) and three years of experience or a Ph. D in life science or computer science. The candidate must have strong interpersonal, written and oral communication skills and be a lateral thinker. Must be able to communicate current bioinformatics technology in a clear and precise manner and to discuss projects with scientists and advise what relevant tools may be used or implemented, have the ability to locate relevant data/information and put this into the context of projects they work on. Candidate must have expert knowledge of UNIX (Mac OS X Darwin a plus), with the ability to install and configure command-line-based applications and services. These include: UNIX

windowing systems (X11, FreeX86), UNIX-based open source scientific applications, **Perl and BioPerl scripts**, HTML, XML. Must have experience with one or more of the following: · Web services (WSDL, SOAP) · Genomics and proteomics · Protein structure prediction/visualization. · Sequence analysis, alignment and database searches. · Regular expressions and relational database queries. · Data mining, text mining · Biostatistics

WHY HAS THE PROJECT WORKED?

- A critical mass of individuals who wanted to solve a common problem
- Infrastructure that works
- Sense of community
- Sharing information, tutorials

SOME KEY DEVELOPERS

Not Pictured

Richard Adams

Scott Cain

Sean Davis

Stefan Kirov

Nathan Haigh

Marc Logghe



Chad



Heikki



Aaron



Steve



Jason



Shawn



Allen



Ewan



Chris M



Hilmar



Lincoln



Brian



Chris D
BOSC 2005

INFRASTRUCTURE

- Easily accessible CVS, web service
- Reliable
- Machines that we own
 - Additional services can be added
- Moving to faster machines this summer

COLLABORATIVE NATURE OF THE PROJECT

- Folks working on similar problems
- Pooled development resources
- Friendly & interactive way to do science
- Shared user support load...

OPEN SOURCE IS GOOD!

- Fix things faster
 - (If you really care about the fix)
- Everyone can see solution, audit code
- Contribution can be modular
- “Give a little, get a lot”
- More developers for less resources

OPEN SOURCE DOESN'T SOLVE EVERYTHING

- Best laid plans
 - Decision by consensus ...
 - ... Free-for-all committing
- Things don't "just happen"
 - Need strong leadership, vision, and commitment
- Documentation tends to lag
- Interesting, needed code trumps boring



GREAT POWER == GREAT RESPONSIBILITY

- Making releases
 - Point releases maintain API
 - Attempt to maintain backwards compatibility between stable releases
- Many people depend on the code
- Additional modules should maintain continuity

WE KNOW HOW TO HAVE FUN!



BIOPERL: PRESENT DAY

- Parsers for
 - Sequence & alignment parsing
 - Gene Prediction output
 - Phylogenetic trees
 - Weight Matrices (TFBS)
 - Ontologies
 - Structure data
 - Various CompBio, MolEvol apps

BIOPERL: PRESENT DAY

- Tools for
 - Sequence manipulation (I/O, manip)
 - Graphical sequence rendering
 - Sequence & alignment stats, popgen tests, molevo tests
- Access to Local & Remote Sequence Databases

PRESENT DAY STATS

- 780+ modules
- 336,000 lines of code
- ~9,000 tests in unit test 't' dir
- 57 utility scripts, 60 example scripts
- 11 HOWTO docs
- 44 Q&A in FAQ

BIOPERL MAGIC

- Using interfaces to hide specifics
 - Bio::DB::GenBank & Bio::Index::GenBank
- Objects can masquerade like other objects
 - Bio::SeqFeature::Generic, Bio::Location, DB::GFF::RelSegment, DB::GFF::Feature

BIOPERL MAGIC

```
use Bio::DB::GenBank;
```

```
use Bio::DB::Fasta;
```

```
my $db = Bio::DB::GenBank->new;
```

```
my $idx = Bio::DB::Fasta->new($faste_dir);
```

```
my $seq1 = $db->get_seq_by_acc($acc);
```

```
my $seq2 = $idx->get_seq_by_acc($acc);
```

```
use Bio::DB::FileCache;
```

```
my $cachedb = Bio::DB::FileCache->new($db);
```

```
my $acc = $cachedb->get_seq_by_acc($acc);
```


...THE PROBLEMS

- Interfaces system is cluttered
 - Learning curve is steep
 - Nonsensical to some people
- OO Code is too slow
- Developer community needs growth

WHY IS BIOPERL SLOW (IN SOME PLACES)?

1. Object-Oriented Perl is a Hack!
 1. `$class->SUPER::new(@_);`
 2. `$class->SUPER::_rearrange(@_);`
 3. Object overhead can be high.

RESEARCH WITH BIOPERL

- Genome annotation pipeline
- Genome Browsing
- Comparative genomics
- Phylogenomics

- Existing Genome Annotation
- New Annotation

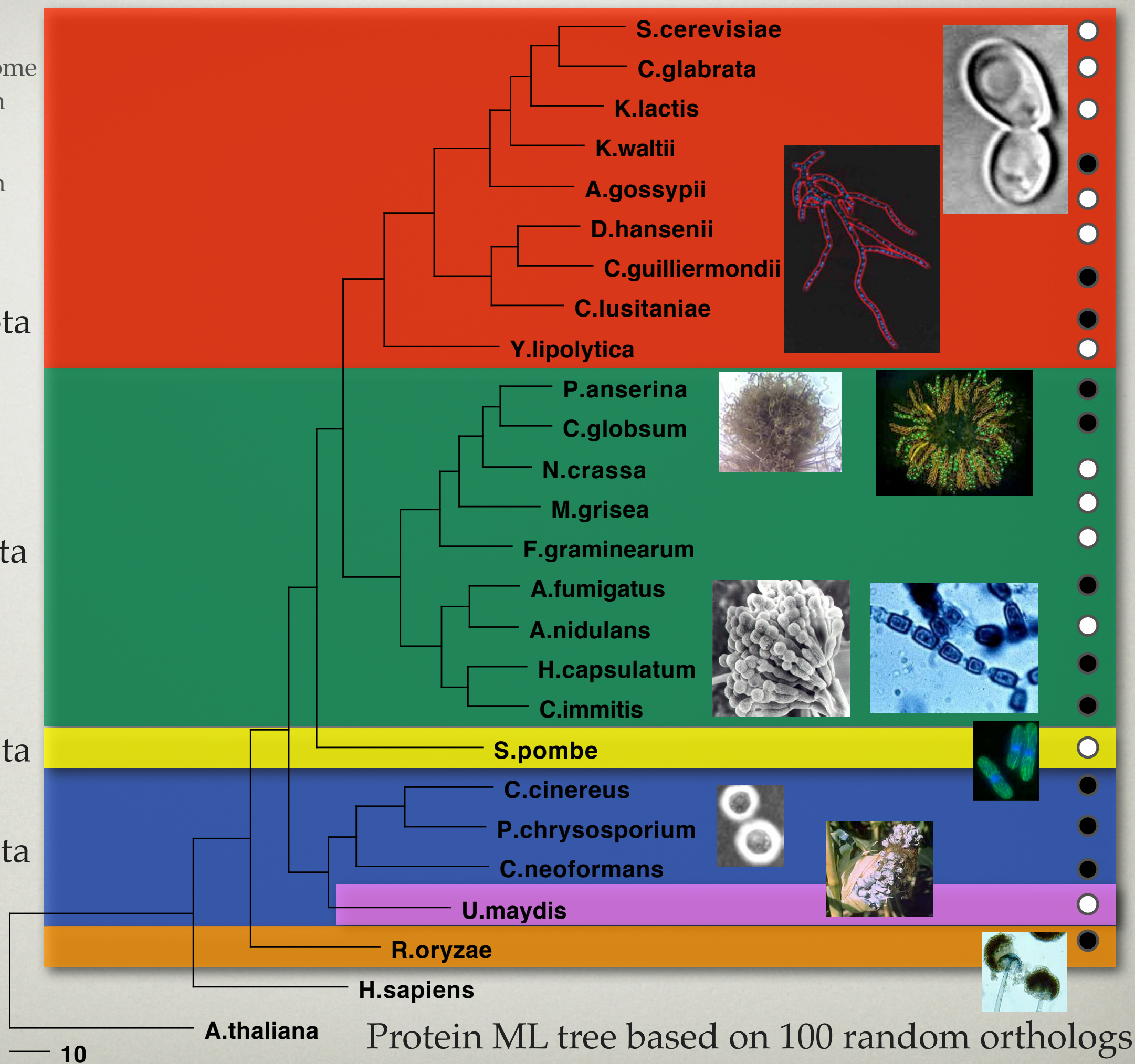
Hemiascomycota

Euascomycota

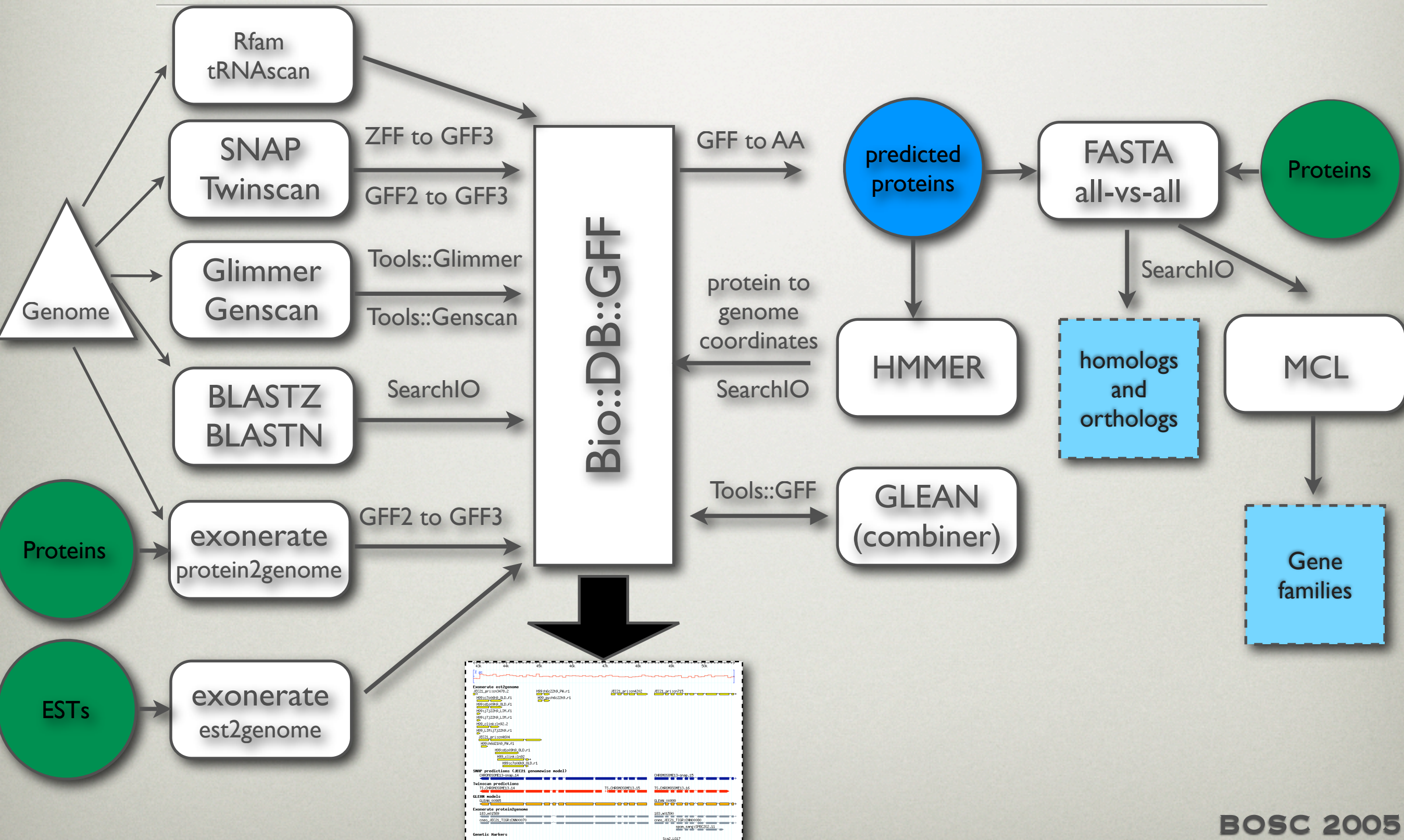
Archiascomycota

Basidiomycota

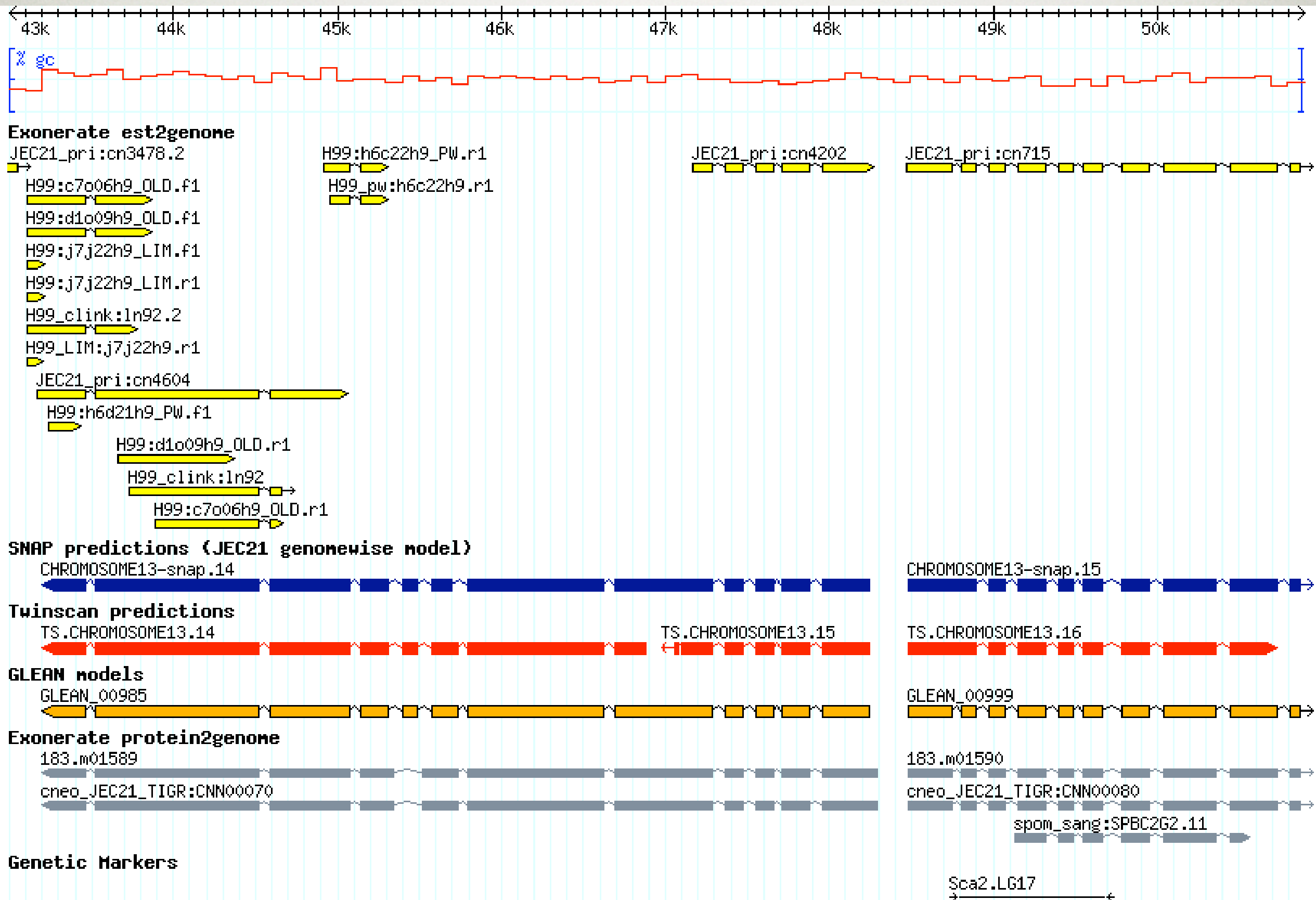
Zygomycota



GENOME ANNOTATION PIPELINE



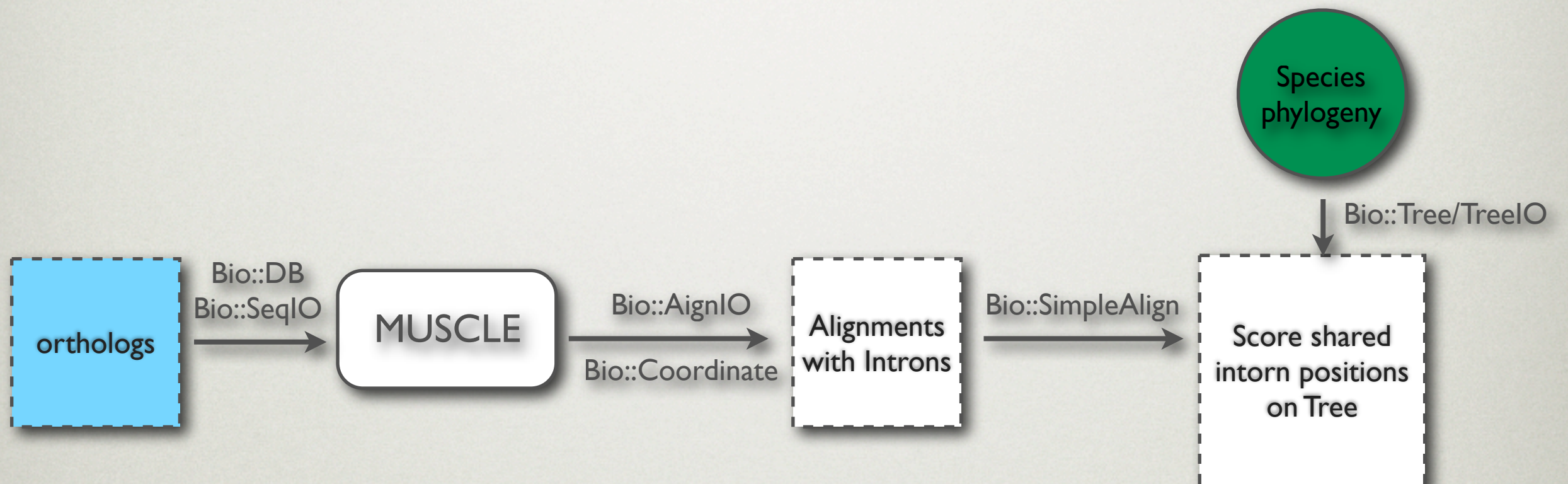
GENOME BROWSING



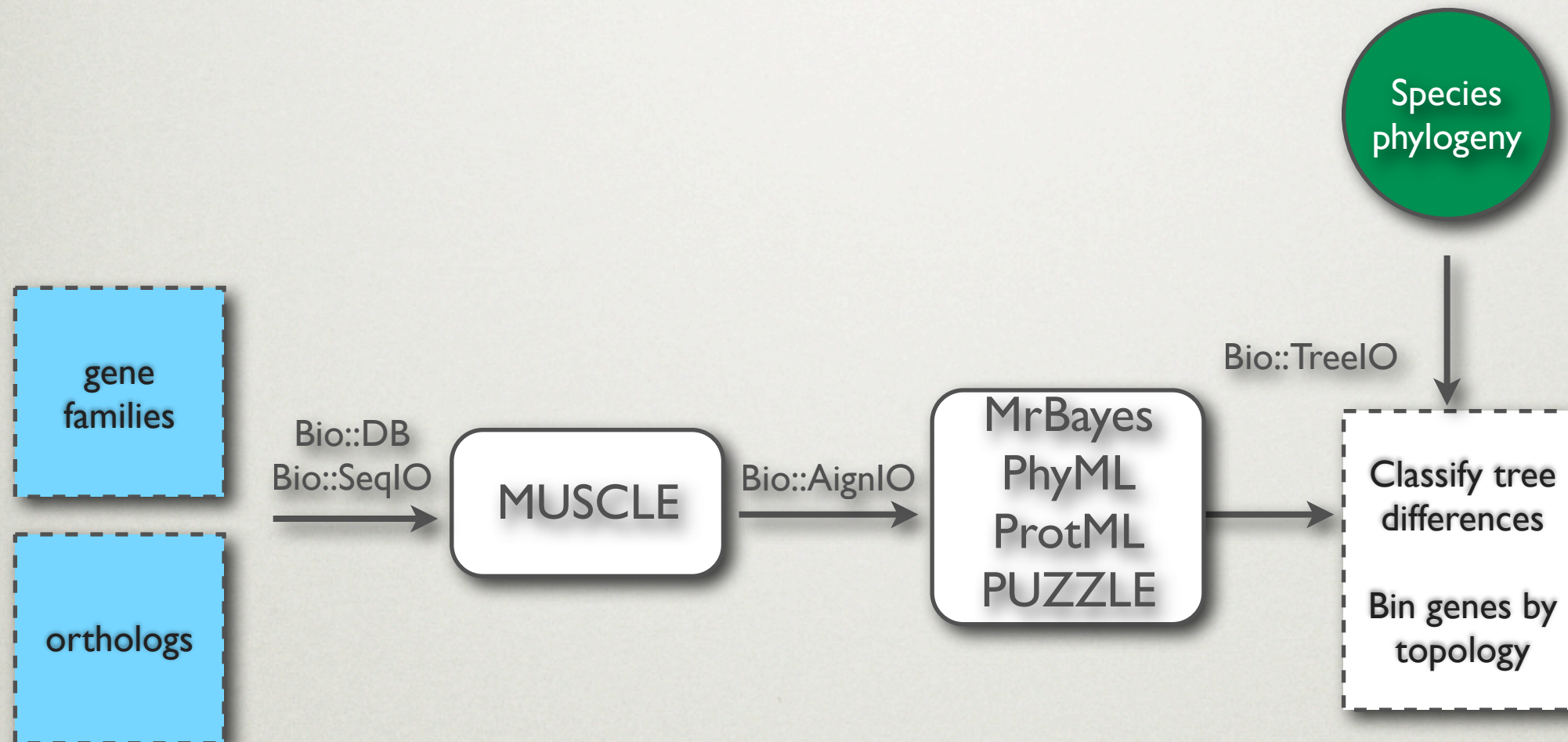
COMPARATIVE GENOMICS

- Find orthologs, align, build trees
 - Assess gene structure among orthologs
- Find gene families
 - Assess gene structure
 - Copy number among genomes

COMPARATIVE GENOMICS: GENE STRUCTURE



PHYLOGENOMICS



BIOPERL - THE FUTURE

WHERE SHOULD BIOPERL BE HEADED?

- Concise set of functionality that just works?
- Supporting leading edge of research?
- Beginner level tools?
- Advanced tools?

CHALLENGES

- Logistics
 - Building and maintaining large projects takes a lot of time
 - Project leaders with vision (and time)
 - Face-to-face time invaluable for coordination
 - Need more dedicated developers

STARTING OVER?

- Years of organic development don't make a wholly consistent toolkit
- Keeping POD up to date vs API maintenance
 - Autodoc-ing?
- Perl6 rewrite?
- Do we start a Bioperl 2, from scratch?

SUPPORTING THE TOOLKIT

- Consider applying for funding to support a developer?
- Finding companies to provide more dedicated user support?
- Gatherings (hackathons)
 - DocFests to increase documentation
- Better incentives for bug fixes

MY PRIORITY QUEUE

- Simplified API, hidden interfaces
- Overview docs point functionality to modules
- Speed
- Extensibility

A PLAN

- Continue Bioperl 1.x
 - Try and fix interface overload
 - Work out speed issues (maybe)
- Encourage dabbling into a 2.0 codebase
 - Outline some key principals in new codebase.
 - Consider prototypes & Perl6

WE'VE GROWN UP

- Over the course of the project we have
 - Gotten married
 - Had kids
 - Changed continents
 - Changed jobs



THANKS



- Developers - past and present
- Ewan Birney
- Lincoln Stein
- Chris Dagdigan
- Brian Osborne