

BioSmalltalk v 1.0

Hernán Morales (hernan.morales@gmail.com)
Institute of Veterinary Genetics (IGEVET - CONICET - CCT La Plata)
Buenos Aires, Argentina

15/05/2017

Contents

Introduction	4
Notes and Disclaimers	4
Installation	5
Smalltalk Basics	6
Additional References	6
Scripting Languages versus Smalltalk	8
Executing Scripts	9
Bioinformatics with BioSmalltalk	10
Sequences	10
Mutable Sequences	12
Operations on Sequences	13
Basic Operations	13
Sequence complement	13
Reverse complement	13
Transcription and back transcription	14
Sequence translation	14
Alphabets	16
International Union of Pure and Applied Chemistry (IUPAC)	16
Sequence Records	18
Repetitive sequences	20
Interspersed repeats.	20
Tandem repeat sequences	21

Segmental duplications	23
Panels of Markers	23
Microsatellites and Minisatellites Locus	24
Wrappers	25
Parsing	25
About parsing	25
XML parsing modes	25
Parsing in BioSmalltalk	26
Parsing Sequences	27
Parsing Accession Numbers	27
Parsing FASTA files	28
Parsing XML files	29
Alignments	30
BLAST	31
Reading and Writing Sequences	33
Reading FASTA records	33
Entrez	35
Entrez Examples	35
Swiss-Prot and ExPASy	36
PDB	36
PDBePISA	36
Population Genetics	38
Organism Names	38
Population Repository Model	38
PLINK	39
ShapeIt	43
LAMP-LD	44
Wrapper execution	45
HMMER	46
Phyx	47
Phylogenetics	48
Nexus File Format	49
Newick File Format	50
Clustering	51
Graphics	52
DotPlots	52
Bar Charts	53
Histograms	53

Widgets	54
Chromosome Selection	54
CSV Explorer	56
Matrix Explorer	56
Third Party Libraries	57
Persistency	57
Mathematics	58
OS Communication	59
Developing BioSmalltalk	60
Code Organization	60
User Code Organization	60
Developer Code Organization	63
Building from Source	64
Testing	66
Executing Tests	66
Writing Tests	67
Test Data	67
SubSystems	69
The Engines SubSystem	69
The Parsing SubSystem	71
The Entrez SubSystem	74
The Files SubSystem	76
Troubleshooting	77
Community	77
Mailing lists	77

Introduction

BioSmalltalk is a library for doing Bioinformatics with the Smalltalk programming system. This guide's goal is to be useful for bioinformaticians, biologists, or life scientists, with little or no experience at all in programming or object technology. It will explain everything needed to understand how to use BioSmalltalk to reach your objectives, also contains the specific context and vocabulary to communicate properly in mailing lists when you need help. Object-Oriented Programming concepts are not required for basic tasks (like building sequences, parsing files, and exploring results) but are specially necessary to build more complex objects and/or debugging code, for example.

Notes and Disclaimers

BioSmalltalk is free of charge and far from perfect. The developer does not provide any guarantees for quality, availability, or fitness for particular purpose. BioSmalltalk is licensed with the MIT License.

If you feel that you can contribute to **BioSmalltalk** development, please do not hesitate to contact through the available communities.

Installation

BioSmalltalk is currently implemented in Pharo, a Smalltalk dialect. It is distributed in several forms depending of the target platform:

- For Microsoft Windows as a .exe installer.
- For Mac/OS X as a compressed ZIP file.
- For GNU/Linux as a compressed .tgz file.

Once installed, all distributions contains:

- A virtual image file (.image extension)
- A changes file (.changes extension)
- The Smalltalk Virtual Machine (VM) (executable file)
- A group of folders including examples and test data.

The usual way to run BioSmalltalk is to open the "Environment", a run-time Integrated Developement Environment (IDE) where all coding activities are performed. This means launching the VM and selecting a virtual image. At first you will have just one image - the one you downloaded in the BioSmalltalk distribution - but as you gain more experience with Smalltalk you will find having several images for different purposes.

You may also work by executing scripts from .st files (explained in the Executing Scripts section) however you would miss powerful Smalltalk tools like Browsers, Refactoring Tools, Code Critics, Finders, etc.

Smalltalk Basics

You may evaluate scripts or code in BioSmalltalk as in most Smalltalk environments: Selecting the code, right button, and choose one of the following operations:

- **Print it (Show it, or Display it):** To print the results as a simple String next to the selected code.
- **Inspect it:** To open an Inspector window. You may observe and send messages over the selected object in the left pane
- **Explore it:** To open an Explorer window. You may expand the tree items, and send messages over the selected object.

At least you should work with two tools. A "System Browser" and a "Workspace". The famous System Browser, Class Hierarchy Browser, or simply Browser, is the most important tool for development in Smalltalk and enables to explore and modify the categories, classes, methods in the system. Actually is so powerful that you can do more advanced actions like finding class references, perform code critics and refactorings, find method senders and implementors, execute tests, group classes, manage breakpoints, find method examples and past versions, among other actions.

Additional References

A comprehensive description of any Smalltalk system would deserve a complete book on its own. Fortunately, there are several good books devoted to describe the Smalltalk language and environment. I recommend the following ones as starting points:

- If you are very new about object technology, a **must** read and one of the best books is "Object Technology" (1997, ISBN: 0-201-30994-7) from David Taylor.
- The Smalltalk Best Practices Patterns is a good book for coding conventions and style: <http://www.amazon.com/Smalltalk-Best-Practice-Patterns-Kent/dp/013476904X>
- There are free on-line books in the Stéphane Ducasse personal web site: <http://stephane.ducasse.free.fr/FreeBooks.html>
- Pharo Books: <http://files.pharo.org/books/>

Finally, a useful advice often done by experienced Smalltalkers is to actually experiment with the environment, make mistakes, and ask in the mailing lists or forums in case of doubt.

Scripting Languages versus Smalltalk

Table 1: Comparison of coding tasks

Python/Perl/Ruby	Pharo Smalltalk
Write a script (in a file) called foo.py/foo.pl/etc	Write a script in some environment tool (Workspace, Script Manager, etc.)
Specify a shebang line	Not necessary
Code navigation	Open the System Browser
Search import modules and how to "invoke" them	Not necessary
Search for senders/callers of a function: Configure ctags for your language	In a method, press the Senders button
Place your source code in directories or/and files	Not necessary, the snapshot system will persist your source code
Execute your program : Find a way to execute the command line interpreter/compiler specifying with your script name as input	Select any piece of valid source code and press the mouse right button, then do it or print it

Usually depends of a web browser on the generated pages from cross-linking the source code and comments by an autodoc tool. This setting usually doesn't allow targeted navigation.

Executing Scripts

You can execute scripts in .st files from command-line. Just cd into the directory where the "pharo" script is found and evaluate the following test expression to check everything is working properly:

```
./pharo -nodisplay Pharo.image eval "3 + 4"
```

To run a script in a file you can pass the file name as argument to the image:

```
./pharo Pharo.image my_script.st
```

Bioinformatics with BioSmalltalk

Developing and using Object-Oriented technology in BioSmalltalk is about finding your objects of interest in the (virtual) image. Think about the virtual image, or simply the image, as a snapshot of the system current state.

If you want to experiment with sequences, your first task is then to find classes related to sequences (hint: Search for **BioSequence** in the Browser). Once found, you can easily select methods and explore the senders to check how they are used. You will see some methods require a minimum instantiation effort, like a reverse complement, some others will require you to set up a whole context of many objects to figure out how to use them.

Although the above workflow could be valid for all the Bioinformatics libraries, in BioSmalltalk you have probably the most powerful tools to explore a system through targeted navigation. This is a big change if you are too used to scripting and sweeping raw textual files.

Sequences

A biological sequence is a single, continuous molecule of nucleic acid or protein, and represents a naturally occurring, or experimentally generated, fragment of genetic or protein material or any intermediate product (like the messenger RNA). It is possibly the most important object in bioinformatics and most bioinformatics workflows involve using sequences at some point.

In BioSmalltalk sequences are represented with **BioSequence** object. We will talk about representation later. If you have previous programming experience, a sequence is essentially a **String** object with an alphabet. It is read-only, if you want to modify a sequence in place you may use **BioMutableSeq**.

Creating a DNA sequence

```
BioSequence newDNA: 'atcggtcggctta'.
```

The above one-liner creates a "BioSequence instance". Nucleotide letters are case-insensitive, this is, meaning does not change based on differing use of uppercase and lowercase letters. When the alphabet information of a sequence is missing (as in this case, or when reading sequences from FASTA files), **BioSequence** will try to guess the alphabet type. Notice you can also create a sequence

specifying the alphabet, for example:

```
BioSequence newAmbiguousDNA: 'AAGTCAGTGTACTATTAGCATGCATGTGCAACACATTAGCTG'.  
BioSequence newUnambiguousDNA: 'AAGTCAGTGTACTATTAGCATGCATGTGCAACACATTAGCTG'.
```

Similarly, to create a RNA and Protein sequence, you may use:

```
BioSequence newRNA: 'augccuacauaggc'.  
BioSequence newProtein: 'AGFAVENDSA'.
```

Please note: Next to the following expressions, between quotations, is displayed the result of the "Print it" (or "Show it" depending the Smalltalk flavor you are using) operation

The format of the **BioSequence** printString is:

```
"  
aBioSequence( sequence_size ) ([ letters_in_alphabet ] alphabet_class) [sequence_letters  
  ]  
"
```

An easier way to create sequences, without specifying the **BioSequence** class, is converting a **String** into a Sequence object. The following expressions display the contents of the #asSequence message.

```
'ATGCTAGN' asSequence. " a BioSequence(8) ([GATCN] IUPAC DNA) [ATGCTAGN] "  
  
'RSWAFFGH' asSequence. " a BioSequence(8) ([ACDEFGHIKLMNPQRSTVWY] IUPAC -> Protein) [  
  RSWAFFGH] "  
  
'AGCUCGGCU' asSequence. " a BioSequence(9) ([GAUCRYWSMKHBVDN] IUPAC -> RNA -> Ambiguous  
  ) [AGCUCGGCU] "  
  
'VPFSNATGVVRSPFEYPQYYLAE' asSequence  
" a BioSequence(23) ([ACDEFGHIKLMNPQRSTVWY] IUPAC -> Protein) [VPFSNATGVVRSPFEYPQYYLAE]  
"
```

Mutable Sequences

A biological mutable sequence is like a normal sequence, but it is editable, i.e. not read-only.

Take note that Smalltalk, as in R, index starts counting from 1. One of the benefits of 1-based system is a more human-centered thinking, rather than 0-based counting which are more machine centered. The following script modifies the first position in a RNA sequence:

```
(BioMutableSeq newRNA: 'auugccuacauaggc')  
  at: 1 put: $g;  
  yourself.  
" a BioMutableSeq(15) ([GAUC] IUPAC -> RNA -> Unambiguous) [guugccuacauaggc] "
```

You may completely replace the sequence if appropriate, this could be convenient for preserving additional information (annotations like sequence features) still valid for the new sequence.

```
(BioMutableSeq newUnambiguousDNA: 'ATGC')  
  seq: 'CGTA';  
  yourself.
```

Operations on Sequences

The following subsections describes available operations on sequences.

Basic Operations

Accessing a subsequence and display it as a **String**:

```
((BioSequence newDNA: 'atcggtcggctta') copyFrom: 3 to: 5) asString. " --> 'cgg' "
```

Count the number of bases of a sequence:

```
(BioSequence newUnambiguousDNA: 'AAGTCAGTGTAATTAGCATGGCTG') size.  
" --> 26 "
```

Get the molecular weight (MW) of a sequence:

```
'AGTACACTGGT' asSequence molecularWeightNonDegen " --> 3435.23"
```

Get the GC nucleotide content:

```
'AGTACACTGGT' asSequence gc.
```

Sequence complement

```
(BioSequence newDNA: 'AAGTCAGTGTAATTAGCATGCATGTGCAACACATTAGCTG') complement  
" a BioSequence(42) ([GATCN] IUPAC DNA) [TTCAGTCACATGATAATCGTACGTACACGTTGTGTAATCGAC] "
```

Reverse complement

This message answers a copy of the receiver's sequence complemented and with element order reversed. You may need to work with the reverse complement of a sequence if it contains an ORF on the reverse strand.

Upper and lower case are preserved as can be used to mark regions of interest. Spaces and special letters (like '-' and spaces) are left untouched.

```
(BioSequence newAmbiguousDNA: 'AcTG C-NH' ) reverseComplement.  
" a BioSequence(9) ([ACGTWSMKRYBDHVN] IUPAC DNA -> Ambiguous) [DN-G CAgT] "
```

Ambiguous and unambiguous nucleotide sequences are supported

```
(BioSequence newUnambiguousDNA: 'GGGGaaaaaaaatttatatat') reverseComplement  
" a BioSequence(21) ([GATCN] IUPAC DNA) [ATATATAAATTTTTTTCCCC] "
```

Supports BioSequence and BioMutableSeq

```
(BioMutableSeq newDNA: 'AAGTCAGTGCTACTATTAGCATGCATGTGCAACACATTAGCTG') reverseComplement  
" a BioMutableSeq(42) ([GATCN] IUPAC DNA) [CAGCTAATGTGTTGCACATGCATGCTAATAGTACACTGACTT] "
```

In the answered sequence, the alphabet is preserved.

```
| mySeq |  
mySeq := BioSequence newDNA: 'AAGTCAGTGCTACTATTAGCATGCATGTGCAACACATTAGCTG'.  
mySeq alphabet = mySeq reverseComplement alphabet
```

Transcription and back transcription

Sequence transcription is a cellular process which firstly reads a DNA gene, producing what is known as *primary transcript* represented by an equivalent RNA sequence (a mRNA precursor). As a second step, a subprocess named *splicing* or *intron splicing* removes the introns from the mRNA precursor and produces the final product: a messenger RNA or mRNA.

Sequence translation

Translation of RNA is another cellular process which could determine the codons of a mRNA segment, and could be performed in BioSmalltalk with a one-liner:

```
(BioSequence newRNA: 'AUGGCCAUUGUAAUGGGCCGCUGAAAGGGUGCCCGAUAG' ) translate.  
" a BioSequence(13) ([ACDEFGHIKLMNPQRSTVWY] IUPAC -> Protein) [MAIVMGR*KGAR*] "
```

the same could be achieved stopping at the first stop codon:

```
(BioSequence newRNA: 'AUGGCCAUUGUAAUGGGCCGCUGAAAGGGUGCCCGAUAG' ) translateToStop  
" a BioSequence(7) ([ACDEFGHIKLMNPQRSTVWY] IUPAC -> Protein) [MAIVMGR]"
```

Translate and concatenate a DNA sequence to Protein sequence:

```
| dna protein |
dna := BioSequence newDNA: 'atcggtcggctta'.
protein := BioSequence newProtein: 'AGFAVENDSA'.
protein , dna translate.
```

It is possible to specify the translation table to use alternative genetic codes, according to the [NCBI Genetic Codes](#) tables (check the transl_table identifier):

```
(BioSequence newRNA: 'AUGGCCAUUGUAAUGGGCCGCUGAAAGGGUGCCCGAUAG' ) translateWithTableId:
2.
" a BioSequence(13) ([ACDEFGHIKLMNPQRSTVWY] IUPAC -> Protein) [MAIVMGRWKGAR*]"
```

Warning: After you select another translation table, the selection remains until you otherwise specify other table. To see which table is currently selected evaluate:

```
BioCodonTables currentCodonTable
```

which in our last example should answer:

```
" a BioCodonTable [2] a Dictionary('AAA'→'K' 'AAC'→'N' 'AAG'→'K' 'AAU'→'N' 'ACA'
'→'T' 'ACC'→'T' 'ACG'→'T' 'ACU'→'T' 'AGC'→'S' 'AGU'→'S' 'AUA'→'M' 'AUC'
'→'I' 'AUG'→'M' 'AUU'→'I' 'CAA'→'Q' 'CAC'→'H' 'CAG'→'Q' 'CAU'→'H' 'CCA'
'→'P' 'CCC'→'P' 'CCG'→'P' 'CCU'→'P' 'CGA'→'R' 'CGC'→'R' 'CGG'→'R' 'CGU'
'→'R' 'CUA'→'L' 'CUC'→'L' 'CUG'→'L' 'CUU'→'L' 'GAA'→'E' 'GAC'→'D' 'GAG'
'→'E' 'GAU'→'D' 'GCA'→'A' 'GCC'→'A' 'GCG'→'A' 'GCU'→'A' 'GGA'→'G' 'GGC'
'→'G' 'GGG'→'G' 'GGU'→'G' 'GUA'→'V' 'GUC'→'V' 'GUG'→'V' 'GUU'→'V' 'UAC'
'→'Y' 'UAU'→'Y' 'UCA'→'S' 'UCC'→'S' 'UCG'→'S' 'UCU'→'S' 'UGA'→'W' 'UGC'
'→'C' 'UGG'→'W' 'UGU'→'C' 'UUA'→'L' 'UUC'→'F' 'UUG'→'L' 'UUU'→'F' ) "
```

To set the default codon table, evaluate:

```
BioCodonTables currentCodonTable: BioCodonTables defaultCodonTable
```

Alphabets

International Union of Pure and Applied Chemistry (IUPAC)

IUPAC is one of the many non-governmental organizations composing the International Council for Science. In practice, for the bioinformatician, serves as a "standard" for the naming of the chemical elements and compounds.

Amongst other nomenclatures, in bioinformatics the *IUPAC codes* usually identify amino acids and nucleotide bases. These codes can consist of either one letter code or a three letter code. For example DNA letters are A, G, C, T bases, which compose the DNA Alphabet. Alphabets are used to determine the contents of a sequence, because simple String(s) are not enough. Consider the following sequence:

```
'GTCG'
```

it could be interpreted as:

- A DNA fragment of Guanine, Thymine, Cytosine and Guanine.
- A Peptide with Glycine, Threonine, Cysteine and Glycine.

the solution is to disambiguate specifying the correct Alphabet objects when appropriate.

Alphabets classes are located under the **BioAlphabet** class hierarchy. To access the ambiguous alphabets just ask to the IUPAC alphabet abstract class:

```
BioIUPACAlphabet ambiguousAlphabets. --> an OrderedCollection(BioIUPACAmbiguousRNA  
BioIUPACAmbiguousDNA)
```

The following script displays the alphabet letters available in the IUPAC specific classes:

```
BioIUPACAmbiguousDNA codes --> ACGTWSMKRYBDHVN  
BioIUPACAmbiguousRNA codes --> GAUCRYWSMKHBVDN  
BioIUPACAmbiguousDNA ambiguityCodes --> WSMKRYBDHVN  
BioIUPACAmbiguousRNA ambiguityCodes --> RYWSMKHBVDN  
BioIUPACUnambiguousDNA codes --> GATCN  
BioIUPACUnambiguousRNA codes --> GAUC  
BioIUPACDNAExtended codes --> GATCBDSW  
BioIUPACProtein codes --> ACDEFGHIKLMNPQRSTVWY  
BioIUPACProteinExtended codes --> ACDEFGHIKLMNPQRSTVWYBXZJUO  
BioNullAlphabet codes --> -
```


Notice however the IUPAC encoding lacks of symbols for termination, gap (-), or selenocysteine.

Create and display sequence objects guessing the alphabet:

```
| sequence |  
sequence := 'atcggtcggctta' asSequence.  
sequence alphabet.  
  
sequence := 'gccuacau' asSequence.  
sequence alphabet.
```

Sequence Records

BioSeqRecord stores a sequence with additional information, which could be "annotations" or sequence features. You can think of a *Sequence Record* like a wrapper for a **BioSequence** to attach additional information to a sequence. Equivalent classes in other Bio toolkits are `Bio::Sequence` for BioRuby and `Bio.SeqRecord` for Biopython. In a sequence record usually the `#id` method is used to store an unique identifier such as the Accession Number (a **BioAccession**), and the `#name` is used to hold a **String** with a more descriptive name, such as clone name. The `#description` method is used to set up a **String** with a textual, i.e. human readable, description of the sequence.

A sequence record responds to the `#annotations` message to retrieve the sequence annotations. Such message will answer a Smalltalk **Dictionary**, however, if you need to customize the annotations you may subclass **BioSeqRecord** and override `#annotationsClass` method to answer a different kind of **Dictionary**.

Creating a **BioSeqRecord**:

```
| sequence |
sequence := BioSequence newProtein: 'mdstnvrsgmksrkkkpkttviddddcmtcsacqslklndfas'.
(BioSeqRecord new: sequence)
    id: 'P20994.1';
    name: 'P20994';
    description: 'Protein A19';
    dbxRefs: #('Pfam:PF05077' 'InterPro:IPR007769' 'DIP:2186N');
    annotationAt: #note put: 'A simple note';
    yourself.
```

Mostly you would want to iterate over a collection of sequences with metadata information, and populate Sequence Records with annotations and features. Some Bio* libraries use to parse a GenBank file (.gbk) which is a textual format intended to be read by humans (even when there exists for years programming utilities like Entrez E-Utils to retrieve information in a machine-readable format). One format which is easier to parse, is the XML format, and the following example retrieves a XML-formatted file with annotations, which are then stored in **BioSeqRecord** objects:

Work in progress

```
'seq_records.xml'
```

To access sequence features, the following methods are available:

- `sequenceFeatures`, `sequenceFeatures`: Answer or set a **Collection** the receiver's sequence **BioSequenceFeatures**.
- `addSeqFeature`: Add a **BioSequenceFeature** passed as parameter, to the receiver's sequence.
- `removeSeqFeature`: Remove a **BioSequenceFeature** passed as parameter, to the receiver's sequence.

Similarly, to access sequence annotations and per-letter annotations, the following methods are used:

- `annotations`, `annotations`: Answer or set a **Collection** the receiver's sequence **BioSeqAnnotations**.
- `annotationsAt:put`: At a key first parameter, add a **Collection** with the receiver's annotations for the sequence.
- `removeAnnotations` : Completely remove all the annotations from the receiver.
- `removeAnnotationsAt`: Remove annotation passed as parameter from the receiver.
- `letterAnnotations`, `letterAnnotations`: Answer or set a **Collection** with the receiver's per-letter annotations.
- `letterAnnotationsAt:put`: At a key first parameter, add a value **Collection** with the receiver's per-letter annotation for the sequence. The collection passed as parameter must match the length of the receiver's sequence length.
- `removeLetterAnnotations` : Completely remove all the letter annotations from the receiver.

Repetitive sequences

There are several types of [repetitive DNA sequences](#), and multiple ways of classify them, for example according to their function, their processes or structure. They are often seen as mobile genetic elements which have been replicated in a genome by inserting new copies of themselves into different positions. Some types of repetitive DNA segments, also known as "DNA fingerprints", are useful for identifying individuals in crime forensics, as genetic markers in linkage analysis and population studies.

The root class for accessing repetitive sequences is **BioRepeatSequence**, and the following subsections explore the possibilities of BioSmalltalk with repetitive sequences. In the next sections, in sequences considered to be repeated, the repeat units are called "motif".

Interspersed repeats.

Also known as Transposon-Derived Repeats, interspersed repeats can be defined as repetitive sequences dispersed throughout the genome as single units flanked by unique sequence. They could be divided between "Transposable DNA" and "Non-Transposable DNA". At the same time the Transposable DNA, could be subdivided in "DNA transposons" and "Retrotransposons" both of them are going to be described in the context of usage with the BioSmalltalk library.

DNA transposons

Also referred as "fossil" DNA elements (as far it is known, DNA transposons are not active anymore in humans), they are the simplest repetitive sequences and represent mobile DNA elements that transpose to new sites directly as DNA. They encode for an enzyme called "transposase" in a process named conservative transposition. An example of an actual DNA transposon [found](#) in *c.elegans* in 1983, is this sequence of 1610 bp:

```
BioDNATransposon new
  name: 'C.elegans Tc1 transposable element';
  motif: 'CAGTGCTGGCCAAAAGATATCCACTTTTGTTTTGTGTGTAACTTTTCTCAAGCATCCATTTGACTTGAA
TTTTCCGTGTGCATAAAGCGAAATGTTACGCAAATTTGCGGACCAACATTACATGATTATCGATTTTTCTGAATTTATTTCAA
TTTTTGATTTTTTCGTTTTTCAATTTTCATTATTTTTTTGAATTATCAATAAACGCACTCTGTTTGTGCACTGGATTTGTTT
GGTTGATAAATTATTTTAAGGTATGGTAAATCTGTTGGGTGTAAAAATCTTTCCTTGGACGTCAAGAAAGCCATTGTAGCTGGCT
TCGAACAAGGAATACCCACGAAAGCTCGCGCTGCAAATTCACGTTCTCCGTCGACTATTTGGAAGTAATCAAGAAGTACCAAAC
TGAGGTGAGTTCGAAAAATATTATTTTTTAATAATAAATGTTTAGAAATCCGTCGCTTTGAGAATCTCGCCCGGCAGGCCTCGAGTG
```

```
ACAACCCATAGGATGGATCGCAACATCCTCCGATCAGCAAGAGAAGATCCGCATAGGACCGCCACGGATATTCAAATGATTATAAGT
TCTCCAAATGAACCTGTACCAAGTAACGAACTGTTCTGTCGACGTTTACAGCAAGCAGGCTACACGGACGAAAGCCAGTCAAGAAAC
CGTTCATCAGTAAGAAAAATCGCATGGCTCGAGTTGCGTGGGCAAGCGCATCTTCGTTGGGGACGTCAGGAATGGGCTAACACATCT
GGTCTGACGAAGCAATTCATTTGTTCCGGAGTGATGGAAATTCCTGGGTACGTCGTCCTGTTGGCTCTAGGTAATCTCCAAAGTATC
AATGCCAACCGTTAAGCATGGAGGTGGGAGCGTCATGGTGTGGGGTGCTTCACCAGCACTTCCATGGGCCCACTAAGAGATCCAAA
CATTATGGATCGTTTTCAATACGAAAACATCTTTGAAACTACAATGCGACCCTGGGCACTTCAAATGTGGGCCGTGGCTTCGTGTT
TCAGCAGGATAACGATCCTAGCATACTTCTCTTCATGTGCGTTTCATGGTTTCAACGTCGTCATGTGCATTTGCTCGATTGGCCAGTC
AGTCTCCGGAATGATCCATAGAGCATTGTGGGAAGGTTGGJUAGCGTCTTGGAGGTATTCGGGCTTCAAATGCAGATGCCAAATT
CAACCAGTTGAAAAACGCTTGGAAAGCTATCCCATGTCAAGTTATTCACAAGCTGATCGACTCGATGCCACGTCGTTGTCAAGCTGT
TATTGATGCAACCGGATACGCGACAAAGTATTAAGCATAATTATGTTGTTTTAAATCCAATTGCTCATATTCGGTACTTTAATTG
TCATTTCTTGAACCTCGGTTTTTCAATATTTCTAGTTTTTCGATTTTGAATTTTCTGAAGTTTTTCAAATCTGTTGAACAT
TTTTGATGAATATTGTGTTTTTAGATTTGTGAACACTGTGGTGAAGTTTCAAAACAAATAACCCTTAGAAAAAGTTACACACAAA
AAACCAAAAGTGGATATCTTTTGCCAGCACTG';
repeats: 30 " Bristol genome "
yourself.
```

Tandem repeat sequences

If the copies of a sequence motif lie adjacent to each other in a block, or an array, we are speaking about **tandem repeats**. There are three types of Tandem Repeat DNA or "Satellite DNA":

- Common "Satellite DNA"
- Microsatellites (or STR)
- Minisatellites.

Both Minisatellites and STR's are also called VNTRs because the number of repeats in a given satellite varies among individuals.

Microsatellites

Also known as STR (Short Tandem Repeats) or SSR (Simple Sequence Repeats), microsatellites are nuclear DNA sequences with fewer than 5 base pairs. In BioSmalltalk we use the microsatellite definition at the US National Library of Medicine Medical Subject Headings (MeSH) considering a STR sequence to be from 2 to 5 base pairs long. STRs are used as genetic markers and population studies, to study gene duplication or deletion.

Microsatellites can be classified in several ways, for example by their repetition pattern (or motif type):

- Perfect: Uninterrupted repeats: ---TATATATA----

- Imperfect: Interrupted by non-repeat nucleotides: ---TATATA--TATATA---TATA
- Compound: Two or more STRs are found adjacent to one another: ---AGAGAGTAGATAGATAGA---
- Complex: ---GAGAGA---AATAATAATAAT---

And also could be classified also by the number of base pairs:

- Mono-nucleotide: AAAAAAAAA
- Di-nucleotides: CACACACA
- Tri-nucleotides: TTATTATTATTA
- Tetra-nucleotides: GATAGATAGATA
- Penta-nucleotides: GTCAAGTCAAGTCAAGTCAA
- Hexa-nucleotides: ACGGTAACGGTAACGGTA

The following expression creates a Microsatellite sequence of 2 base pairs and 5 repeat motif, and prints its representation as an expanded sequence:

```
(BioSTRSequence new
  motif: 'AT';
  repeats: 5;
  yourself) asString.
```

Setting a motif also performs a range check. If the provided motif size is fewer or larger than the allowed by the class, an exception is raised.

One can also create an instance of a Microsatellite by writing as a String:

```
'(AT)5' asBioSTRSequence
```

Minisatellites

A minisatellite is a DNA region (or locus) made up of short repeated sequences of at least 10 base pairs, according to the US National Library of Medicine Medical Subject Headings (MeSH). Minisatellites share much of the same properties as Microsatellites, and could be instantiated using a similar pattern:

```
BioMinisatellite new
  motif: 'TGCGCCATTGTAGCCA';
  repeats: 3;
  yourself.
```

Satellite DNA

Classical satellite DNA consists of repeats larger than minisatellites, with sizes between 100 and 6500 base pairs. The class to access satellite DNA is **BioSatelliteSequence**.

Segmental duplications

Segmental duplications are blocks which may consist of 1000 kb to 200,000 kb of nucleotide pairs that are also present at two or more sites in the genome. Specifically in a chromosome (intrachromosomal) or another chromosome (interchromosomal), usually in subtelomeric regions.

Work in progress

Panels of Markers

Panels or markers are used for routine use in species parentage testing and identification. BioSmalltalk provide classes for accessing information about recommended panels of markers. The root class for the hierarchy of panels is **BioPanel**. Subclasses implements specific panels including class methods for particular species. For example the class **BioFAOPanel** is used to access a panel of loci recommended by the *Food and Agriculture Organization of the United Nations (FAO)* for genetic studies of domestic animals. Also a class for the panel marker information provided by the *International Society for Animal Genetics (ISAG)* is available in the class **BioISAGPanel**.

The markers information could be accessed using the following methods:

```
BioFAOPanel markerAlleleRangesForCattle.  
BioISAGPanel markerAlleleRangesForCattle.
```

answering a **Collection** of **BioLocus**, each one representing a position in a chromosomal location.

The marker allele ranges along with its positions in both cases could be accessed as:

```
BioFAOPanel markerAlleleRangesForCattle.  
BioISAGPanel markerAlleleRangesForCattle.
```

Each panel should implement the method `#markersForSpecies`, providing a **Dictionary** of names of supported species, with its associated method. Also each class provides a method for determining which types of markers (STR, Minisatellite, etc.) are supported by the panel. For example display the markers available for the ISAG Panel answer its marker associated class:

```
BioISAGPanel supportedMarkerTypes " (#BioSTRSequence) "
```

Microsatellites and Minisatellites Locus

BioSmalltalk can use sequences with locis through the **BioLocus** class which is used to represent "static" or unvariable information about a microsatellite or minisatellite. A locus is used by **BioSTRSequence** or **BioMinisatelliteSequence** respectively. In the following examples we are using as reference the Cattle STR information accessible at <http://www.cstl.nist.gov/biotech/strbase/cattleSTRs.htm>. Specifically we are going to set up the BM1818 loci through code. A microsatellite locus can have a simple or a compound repeat structure. You can set this with the proper message sends:

```
BioLocus
  name: 'BM1818';
  chrLocation: 'D23S21';
  beSimpleRepeat;
  yourself
```


Wrappers

Issuing command line settings is a typo prone process, specially for command line programs with many options. To prevent such situation, BioSmalltalk includes a common set of classes for building wrappers for external programs, which lets use messages to configure program parameters.

Parsing

About parsing

Parsing is one of the most important part of developing software in Bioinformatics. Parsing is a process by which a developer makes sense of a sentence (acting as input data, usually a **String**), usually by breaking it down into letters or words.

A parser is an object that scans an input **String**, decomposes it into its constituent parts (which you may define), and then optionally processes the parts in some suitable manner.

Programming literature usually refer to writing a parser *by hand* when you write a parser without any support programming library. Currently taking such approach is rare and not recommended (because keeping track of contingencies caused by different inputs and input paths is extremely difficult), so developers usually use a parser generator, i.e. a software library, and write rules for such generator.

XML parsing modes

One of the major file formats in Bioinformatics is XML. XML is a widely supported machine-friendly format much easier to parse than the error-prone processing of hardly structured ASCII format (intended for human reading). Also, another huge advantage of XML is that it constitutes a stack or building block upon which other related technologies has been developed such as XSLT, XML Schema, DocBook, XQuery, XPath, XQL, XPHIL, and others.

Under the XML core building block, there are several types of parsers used against a XML formatted file:

- Tree Parser: This is what you will read everywhere as the "DOM parser"

- Event Parser: This is denoted by S*X and could be -- SAX Parser: Known as the "Push parser" -- StAX Parser: Known also too as the "Pull parser"
- VTD Parser : This is known as "Virtual Token Descriptor"

Note: as far as I know currently there is no VTD parser for Smalltalk. This is the most usual classification of XML Parsers:

There are several other classifications depending of the parser characteristics and what you want to do or how. You may be interested in making modifications or just processing, for which the following classification would make more sense:

- For modifications: The parser creates long-lived representations from the XML document (necessary for modifications): You should choose DOM or VTD -- Do you *need* to query or modify the objects (parser creates nodes): DOM -- You do not need the objects (parser creates integers and locations caches): VTD
- For processing: The parser doesn't creates long-lived objects: SAX or StAX.

or, if you are interested in the type of access to XML nodes, the following "taxonomy" could serve:

- Back-and-forth: Access the data after the parsing is complete: DOM or VTD
- Massive or very frequent access: Choose DOM
- Rare or simple access: Choose VTD
- Sequential: Access the data while you're processing the document: SAX or StAX
- Processing all tokens: SAX
- Processing interested tokens (allows skipping forward): StAX

Parsing in BioSmalltalk

BioSmalltalk uses different libraries to simplify and enable parsing of the many formats out there in the Bioinformatics arena. The libraries used are:

- Support for parsing XML format in DOM or SAX mode.
- XMLPullParser for parsing XML in StAX mode.
- PetitParser for parsing formatted or structured text.

Parsing Sequences

For example, finding if a given **String**, like 'TCGTACGA', is actually a DNA sequence you just evaluate:

```
#dnaSequence asParser matches: 'TCGTACGA'. --> true
#dnaSequence asParser matches: 'TCGTZZZACGA'. --> false
```

And for parsing:

```
#dnaSequence asParser parse: 'TCGTACGA'. --> 'TCGTACGA'
#dnaSequence asParser parse: 'TCGTA3344CGA'. --> end of input expected at 5
```

The #asParser method actually instantiates a parser, which is then sent the message #matches: to test the given input.

Another common way to use the existing parsers is through the **BioParser class**. In BioSmalltalk we have two modes of parsing: Tokenize and Parse. The tokenize messages always answer a Collection-like object containing other collections or primitive Smalltalk objects (this is, something like #('object1' 'object2') or 1234, 74.32, true, etc.), while the parse messages answer **BioObject**'s.

You may want BioObjects if you need to keep working with the resulting bio-objects from your parsing. Or to learn relationships between BioObjects.

You may work with primitive Smalltalk objects like **Integers** or **Strings** if you want "cheap" objects.

Parsing Accession Numbers

An accession number is a unique sequential number assigned to each record (sequence) in a repository. This number allows for tracking of different versions of a sequence record and the associated sequence over time. Parsing an accession number answers a very simple object which you may use to store or index separately from its sequence, thus, reducing storage space and processing time. To parse an accession number:

```
BioParser parseAccession: 'CAF97855.1' --> a BioAccession CAF97855 1
```

You may also query if the accession number is versioned:

```
(BioParser parseAccession: 'CAF97855.1') isVersioned --> true
```

And query the version number:

```
(BioParser parseAccession: 'CAF97855.1') version --> '1'
```

Another shortcut for the same message is #asAccession

```
'XP_425521.2' asAccession version --> '2'
```

Parsing FASTA files

This is how to parse a FASTA file with multiple sequences:

```
BioParser parseMultiFasta: (BioFASTAFile on: 'my_sequencesHM.fasta') contents.
```

Given a header, split its identifiers:

```
BioParser tokenizeFastaHeader: '>gi|198282148|ref|NC_011206.1| Acidithiobacillus  
ferrooxidans ATCC 53993 chromosome, complete genome'.  
  
" #('>gi' '198282148' 'ref' 'NC_011206.1' 'Acidithiobacillus ferrooxidans ATCC 53993  
chromosome, complete genome') "
```

Print in the Transcript each sequence description with its sequence size from a FASTA String:

```
| fastaString |  
  
fastaString := '>gi|2765658|emb|Z78533.1|CIZ78533 C.irapeanum 5.8S rRNA gene and ITS1  
and ITS2 DNA  
CGTAACAAGGTTTCCGTAGGTGAACCTGCGGAAGGATCATTGATGAGACCGTGGAATAAACGATCGAGTG  
AATCCGGAGGACCGGTGTAAGCTCAGCTACCGGGGGCATTGCTCCCGTGGTGACCCTGATTGTTGTTGGG  
CCGCCTCGGGAGCGTCCATGGCGGTTTGAACCTCTAGCCCGGCGCAGTTTGGGCGCCAAGCCATATGAA  
AGCATCACCGCGAATGGCATTGTCTTCCCAAAACCCGAGCGGCGGCGTGCTGTCGCGTGCCCAATGA  
ATTTTGATGACTCTCGAAACGGGAATCTTGGCTCTTGCATCGGATGGAAGGACGACGCGAAATGCGAT  
AAGTGGTGTGAATTGCAAGATCCCGTGAACCATCGAGTCTTTGAACGCAAGTTGCGCCGAGGCCATCA  
GGCTAAGGGCACGCTGCTTGGGCGTCGCGCTTCGTCTCTCTGCCAATGCTTGCCCGGCATACAGCC  
AGGCCGGCGTGGTGCAGATGTAAAGATTGGCCCTTGTGCCTAGGTGCGGCGGGTCCAAGAGCTGGTGT  
TTTGATGGCCCGGAACCCGGCAAGAGGTGGACGGATGCTGGCAGCAGCTGCCGTGCGAATCCCCATGTT  
GTCGTGCTTGTGCGACAGGCAGGAGAACCTTCCGAACCCCAATGGAGGGCGGTTGACCGCCATTGCGAT  
GTGACCCAGGTCAGGCGGGGGCACCCGCTGAGTTTACGC  
  
>gi|2765657|emb|Z78532.1|CCZ78532 C.californicum 5.8S rRNA gene and ITS1 and ITS2 DNA  
CGTAACAAGGTTTCCGTAGGTGAACCTGCGGAAGGATCATTGTTGAGACAACAGAATATATGATCGAGTG  
AATCTGGAGGACCTGTGGTAACCTCAGCTCGTCGTGGCACTGCTTTTGTGCGTGACCCTGCTTGTGTTGG
```

```
GCCTCCTCAAGAGCTTTCATGGCAGGTTTGAACCTTAGTACGGTGCAGTTTGCGCCAAGTCATATAAAGC
ATCACTGATGAATGACATTATTGTGAGAAAAATCAGAGGGGCAGTATGCTACTGAGCATGCCAGTGAAT
TTTTATGACTCTCGCAACGGATATCTTGGCTCTAACATCGATGAAGAACGCAGCTAAATGCGATAAGTGG
TGTGAATTGCAGAATCCCGTGAACCATCGAGTCTTTGAACGCAAGTTGCGCTCGAGGCCATCAGGCTAAG
GGCACGCCTGCCTGGGCGTCGTGTGTTGCGTCTCTCCTACCAATGCTTGCTTGGCATATCGCTAAGCTGG
CATTATACGGATGTGAATGATTGGCCCTTGTGCCTAGGTGCGGTGGGTCTAAGGATTGTTGCTTTGATG
GGTAGGAATGTGGCACGAGGTGGAGAATGCTAACAGTCATAAGGCTGCTATTTGAATCCCCATGTTGTT
GTATTTTTTCGAACCTACACAAGAACCTAATTGAACCCCAATGGAGCTAAAATAACCATTGGGCAGTTGA
TTTCATTGAGATGCGACCCAGGTCAGGCGGGGCCACCCGCTGAGTTGAGGC
'.
```

```
(BioParser parseMultiFasta: fastaString) do: [ : seqRecord |
    Transcript
        show: seqRecord name; cr;
        show: seqRecord size; cr; cr ]
```

Parsing XML files

Alignments

Underlying sequences in a `BioAlignment` are sequence based (contrary to String based). This could enable to access track information like quality score through the alignment itself, instead of maintaining multiple associated matrices.

Creating a Multiple Sequence Alignment

```
align := BioAlignment new.  
align  
  addSequence: 'ACTGCTAGCTAG';  
  addSequence: 'ACT-CTAGCTAG';  
  addSequence: 'ACTGCTAGATAG';  
  addSequence: 'ACTGCTTGCTAG';  
  addSequence: 'ACTGCTTGATAG'.
```

BLAST

Provided by the National Center for Biotechnology Information (NCBI), BLAST is a sequence homology search program that finds all the sequences that match a particular given sequence from a database. The NCBI-GenBank database contains over 150 billion nucleotide bases as of June 2013, and doubling its contents approximately every 18 months. It is considered the most fundamental tool in bioinformatics today, although there are many other similar programs for different cases. In BioSmalltalk you may query programmatically a BLAST at NCBI with the following parameters:

- Query the nucleotide (nr) database.
- Use an "unknown" DNA sequence.
- Get only the first 200 hits.
- Filter query sequence with low compositional complexity.
- Use the default expect value (10).
- Use the BLASTN programs to search nucleotide databases using a nucleotide query.
- Get the results in XML for further querying

which is translated as:

Make a BLAST from BioSmalltalk

```
| search |
search := BioNCBIWebBlastWrapper new nucleotide
    query: 'TCGAAATAACGCGTGTCTCAACGCGGTCGCGCAGATGCCTTTGCTCATC
AGATGCGACCGCAACCACGTCCGCGCCTTGTTGCGCGTCCCGTGCCTC
AACCACCACCACGGTGTCTCTTCCCGAACGCGTCCCGGTGAGCCAGCC
TCCACGCGCCGCGCGCGGAGTGCCATTGCGGCCGACGTGCGACGGT
GCCGCTCAGATTCTGTGTGGCAGGCGCGTGTGGAGTCTAAA';
    hitListSize: 200;
    filterLowComplexity;
    expectValue: 10;
    blastn;
    formatTypeXML;
    fetch.
```

Now the "search" object contains a BioWebBlastResult. The following line write the XML contents to a timestamped file in the directory where the .image file is located:

```
search outputToFile: 'blast-' , BioObject currentSeconds , '.xml'.
```

"Note: When the Transcript is opened (Tools --> Transcript) you may view a detailed log of message sends."

Reading and Writing Sequences

Reading FASTA records

The following script count letters for each fasta record in the 'Fasta_example_1.fa' file included in the BioSmalltalk distribution. The file contains multiple fasta records which are parsed and iterated collecting the numbers of letters for each record.

```
| mFasta fPath |  
  
fPath := BioObject testFilesFullDirectoryName , 'Fasta_example_1.fa'.  
mFasta := BioParser parseMultiFasta: (BioFASTAFile on: fPath) contents.  
mFasta collect: #occurrencesOfLetters
```

And of course the output is not nice for the human eye:

```
an OrderedCollection(a Dictionary($A->24 $C->18 $G->25 $T->23 ) a Dictionary($A->33 $  
C->30 $G->23 $T->34 ) a Dictionary($A->25 $C->30 $G->15 $T->39 ) a Dictionary($  
A->17 $C->9 $G->19 $T->15 ) a Dictionary($A->26 $C->17 $G->36 $T->41 ) a  
Dictionary($A->8 $C->4 $G->11 $T->31 ) a Dictionary($A->31 $C->16 $G->17 $T->22 )  
a Dictionary($A->16 $C->14 $G->15 $T->25 ) a Dictionary($A->14 $C->20 $G->11 $T->21  
) a Dictionary($A->20 $C->29 $G->25 $T->41 ) a Dictionary($A->12 $C->14 $G->8 $  
T->26 ) a Dictionary($A->26 $C->8 $G->11 $T->24 ) a Dictionary($A->36 $C->37 $G->16  
$T->36 ) a Dictionary($A->30 $C->30 $G->33 $T->27 ) a Dictionary($A->14 $C->11 $  
G->13 $T->20 ) a Dictionary($A->34 $C->31 $G->20 $T->27 ) a Dictionary($A->23 $  
C->31 $G->29 $T->37 ) a Dictionary($A->30 $C->42 $G->25 $T->31 ) a Dictionary($A->7  
$C->23 $G->5 $T->18 ) a Dictionary($A->15 $C->28 $G->6 $T->18 ) a Dictionary($  
A->23 $C->44 $G->12 $T->37 ))
```

Suppose we want to display the output in CSV format. Each object could be formatted for CSV by implementing #outputAsCsvTo: aStream. Then modifying the script:

```
| mFasta filePath |  
  
filePath := BioObject testFilesFullDirectoryName , 'Fasta_example_1.fa'.  
mFasta := BioParser parseMultiFasta: (BioFASTAFile on: filePath) contents.  
BioCSVFormatter new exportFrom: (mFasta collect: #occurrencesOfLetters).
```

It is possible to get a timestamped file like "CSVOutput3536647736.csv" with contents like this:

```
A;24
```

G;25
C;18
T;23

A;33
G;23
C;30
T;34

A;25
G;15
C;30
T;39

...

Entrez

The Entrez API is an interface for accessing to the Entrez Utilities provided by the NCBI.

Entrez Examples

FASTA record from Protein

- Input: Accession number
- Database: Protein
- Output: FASTA record

```
(BioEntrezClient new protein
  term: 'NP_031402.3';
  useWebEnv;
  setFasta;
  fetch) outputFile: 'NP_031402.3.fasta'
```

SeqID record from Protein

- Input: Accession number
- Database: Protein
- Output: SeqID

```
(BioEntrezClient new protein
  term: 'NP_031402.3';
  useWebEnv;
  setSeqId;
  fetch) outputFile: 'NP_031402.3.seqid'
```

GenBank records from Nuccore

- Input: List of UIDs
- Database: Nuccore
- Output: GenBank formatted file (.gb)

```
(BioEntrezClient new nuccore
```

```
uids: #(57240072 57240071 6273287 6273291 6273290 6273289 6273286 6273285
6273284);
setGb;
fetch) outputToFile: 'fetchNucore3.gb'
```

XML records from Pubmed

- Input: List of UIDs
- Database: Pubmed
- Output: XML

```
| result |
result := BioEntrezClient new pubmed
      uids: #( 11877539 11822933 );
      setModeXML;
      fetch.
result outputToFile: 'eFetch1.xml'
```

XML records from Protein

- Input: Search terms
- Database: Protein
- Output: List of UIDs in XML

```
(BioEntrezClient new protein
  term: 'insulin AND homo';
  search) outputToFile: 'insulin AND homo_protein.xml'
```

Swiss-Prot and ExPASy

PDB

PDBePISA

BioSmalltalk includes an API to access the service provided by the EBI: Protein Data Bank in Europe : Proteins, Interfaces, Structures and Assemblies tool (Krissinel and Henrick, 2007). This tool can help in the analysis of the protein-protein interfaces within structures, search for related protein structures,

verify the reliability of the models interfaces, find the interactions between the ligand and receptor proteins, and other [set of related features](#).

To search inside the PDBePISA server you should provide a list of PDB codes, each one represented by 4-letter Protein Databank (PDB) identifier, as parameter of the method #pdbCodes:.

Inside BioSmalltalk, highlight the following expression and bring the contextual menu to Inspect or Explore the results:

```
BioPDBePISAClient new interfaces
  pdbCodes: { '3gcb' . 'lsar' };
  fetch.
```

You can also download the results in the XML format provided by the PDBePISA.

```
(BioPDBePISAClient new interfaces
  pdbCodes: { '3onz' . 'ld66' };
  fetch) dumpToFileNamed: 'fetchPdbEPisa_01.xml'
```

To further filter results, you can use the XML parser provided by BioSmalltalk. The following code includes two PDB entries, which instead of downloading directly into a file, we can parse.

```
(BioPDBePISAClient new interfaces
  pdbCodes: { '3gcb' . 'lsar' };
  fetch)
```

Population Genetics

Population genetics is the study of genetic variation, specifically the allele frequency changes, in a population, as the population is subject to the four main evolutionary processes: natural selection, genetic drift, mutation and gene flow. It also takes into account the factors of recombination, population subdivision and population structure.

Organism Names

BioSmalltalk includes lists of common organism names, which can be accessed by:

```
BioSpeciesLists commonAnimalSpecies.  
BioSpeciesLists faoSpecies.
```

You may also configure a "current" species list, for your particular application.

```
BioSpeciesLists currentSpeciesList
```

Population Repository Model

BioSpecies serves as a global repository for organism populations. Each organism is represented by a **BioPopRepository**. In turn, each **BioPopRepository** stores repository versions, instances of **BioPopRepositoryVersion**, which has the actual varieties (read: breeds, models, strains, etc. depending on your organism taxonomy vocabulary) .

You may start inspecting this model by evaluating:

```
BioSpecies repositories.
```

Some repositories could be already populated:

```
BioSpecies populatedRepositories.
```

To access a specific organism repository, you may send #repositoryVersions as shown below:

```
(BioSpecies repositoryFor: 'Cattle') repositoryVersions.
```

And accessing a particular repository version for an organism (the following expression works in PhyloclassTalk, a phylogeography environment based in BioSmalltalk) :

```
(BioSpecies repositoryFor: 'Cattle') repositoryVersionAt: 'Mason'.
```

PLINK

The PLINK wrapper exists in at least two variants. One for working with the 1.07 version (class BioPLINK1Wrapper) and another one for working with the "2" version (actually 1.9, class BioPLINK2Wrapper).

Wrapper execution

File format conversions

Make a binary pedigree file (.BED):

```
| samplesDir |  
  
samplesDir := BioPLINK2Wrapper testFilesFullDirectoryName.  
BioPLINK2Wrapper new  
  file: (samplesDir / 'sample_ped') fullName;  
  out: 'sample_out';  
  makeBed;  
  execute.
```

Transpose to output .TPED and .TFAM files:

```
| samplesDir |  
  
samplesDir := BioPLINK2Wrapper testFilesFullDirectoryName.  
BioPLINK2Wrapper new  
  file: (samplesDir / 'sample_ped') fullName;  
  out: 'sample_out';  
  transpose;  
  execute.
```

To get out numeric alleles:

```
| samplesDir |
```

```

samplesDir := BioPLINK2Wrapper testFilesFullDirectoryName.
BioPLINK2Wrapper new
  file: (samplesDir / 'sample_ped') fullName;
  out: 'sample_out';
  recode12;
  execute.

```

Basic filters

To get out specific chromosome excluding individuals with genotyping error > 5% (or call rate at least of 95%):

```

| samplesDir |

samplesDir := BioPLINK2Wrapper testFilesFullDirectoryName.
BioPLINK2Wrapper new
  file: (samplesDir / 'sample_ped') fullName;
  out: 'sample_out';
  chr: 23;
  mind: 0.5;
  execute.

```

Combining filters in one pass:

```

| samplesDir |

samplesDir := BioPLINK2Wrapper testFilesFullDirectoryName.
BioPLINK2Wrapper new
  file: (samplesDir / 'sample_ped') fullName;
  out: 'sample_out';
  chr: 23;
  mind: 0.5;
  geno: 0.1;
  maf: 0.05;
  execute.

```

Advanced filters

Splitting a genome PED by chromosome. This script generates N files according to the number of chromosomes if non-human data is used, using the provided species parameter (cow, dog, horse, mouse, etc). The list of supported organisms is currently available at: https://www.cog-genomics.org/plink/1.9/input#chr_set


```
| samplesDir |

samplesDir := BioPLINK2Wrapper testFilesFullDirectoryName.
BioPLINK2Wrapper new
  file: (samplesDir / 'sample_ped') fullName;
  out: 'sample_out';
  cow;
  geno: 0.1;
  execSplitChrs
```

Splitting by family IDs requires to specify a **Collection** with valid families in your data set.

```
| samplesDir |

samplesDir := BioPLINK2Wrapper testFilesFullDirectoryName.
BioPLINK2Wrapper new
  bfile: (samplesDir / 'sample_ped') fullName;
  out: (samplesDir / 'sample_bed_fams') fullName;
  keepFams: #('4304' '7319' '6968');
  makeBed;
  execute.
```

Reading PLINK (PEDMAP) files

BioSmalltalk includes a simple utility to browse and select PED/MAP files. You can use it for your projects and launch it by evaluating:

```
BioPEDMAPFilesLoader open.
```

The window will require mandatory population name, a .PED and .MAP files to be selected.

Once opened, you can browse to your PED/MAP files and click Load to perform an action.

Which action is performed by default? If no action was specified, then the Load button will close the PED/MAP loader and return an Association with the file name of the selected .PED file as key, and the name of the .MAP file as value. The answered directory name will have a system dependent path separator (e. g. / for Unix) appended:

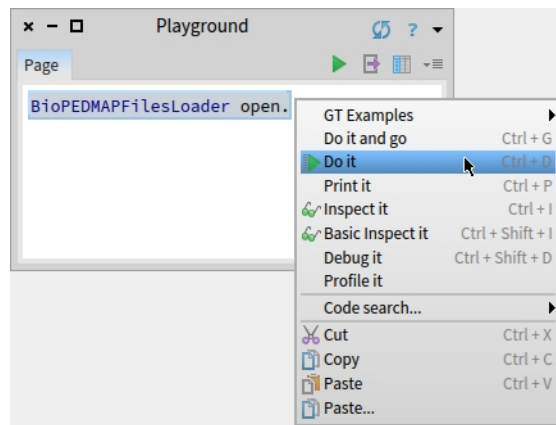


Figure 1: Opening the PEDMAP loader

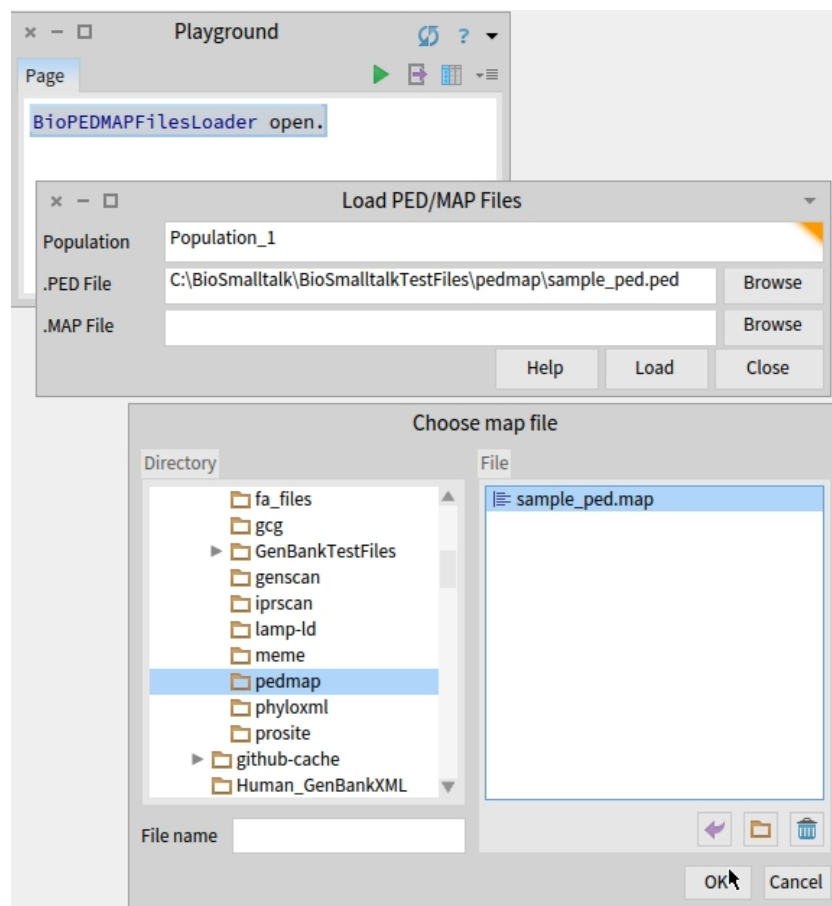


Figure 2: Selecting MAP file

```
''C:\BioSmalltalk\BioSmalltalkTestFiles\pedmap\sample_ped.ped' -> 'C:\BioSmalltalk\
BioSmalltalkTestFiles\pedmap\sample_ped.map''
```

You can also specify which action to realize after the Load button gets clicked. The action block in the following code acts as a callback, receiving two String parameters: The PED file name and the MAP file name

```
BioPEDMAPFilesLoader open
  actionBlock: [ : pedFilename : mapFilename | "..."]
```

Forcing reference alleles

A key step in preparing files for other software packages (such as LAMP-LD) is to add proper reference alleles information to data set, to ensure allele coding consistency. The inconsistency problem is relevant when merging multiple datasets from different experiments or platforms.

The following script adds a reference allele file to avoid allele coding discordance:

```
BioPLINK2Wrapper new
  inputFile: '';
  outputFilename: '';
  referenceAlleles: 'allele_A.txt';
  makeBed;
  execute.
```

ShapeIt

Input preparation

ShapeIt is a fast and accurate method for estimation of haplotypes (a.k.a. phasing) from a set of SNP genotypes (.ped format or its .bed/.bim/.fam binary version) and a genetic map (.map format), and produces as output, either a single set of estimated haplotypes, or a haplotype graph that encapsulates the uncertainty about the underlying haplotypes. The software is currently only available in Unix-like OS.

ShapeIt2 outputs a pair .haps (IMPUTE2 format) and .sample files.

Wrapper execution

To use the wrapper the program binary must be in the system PATH environment variable and all input files, being binarized PLINK (bed, bim, fam) or textual PLINK (ped, map) must have the same name. The following expression launches ShapeIt2 from BioSmalltalk, setting several parameters such as:

- The number of burn-in MCMC iterations
- The input file name (without extension),
- The output file name for the best haplotypes
- The number of threads to use the multi-threading capabilities

```
BioShapeIt2WrapperR727 new
  burn: 10;
  inputBinarized: 'input_brangus';
  outputMax: 'output_brangus';
  thread: 8;
  execute
```

You can also explicitly specify the PEDMAP files as parameters:

```
BioShapeIt2WrapperR644 new
  burn: 10;
  inputTextual: 'input_brangus.ped';
  inputMap: 'input_brangus.map';
  outputMax: 'output_brangus';
  thread: 8;
  execute
```

LAMP-LD

BioLAMPLD is a package which includes several classes for working with the LAMP-LD software for doing inference of locus-specific ancestry in recently admixed populations. LAMP-LD is currently available only for Unix/Linux systems.

Input preparation

The LAMP-LD package includes a class for writing the admixed genotype sample file (.geno or .gen) which contains genotype data for the problem population.

```
BioLAMPLDGenotypeFormatter new
```

```
bimFilePath: 'file.bim';
pedFilePath: 'file.ped';
export.
```

Wrapper execution

The main class for interacting with LAMP-LD is BioLAMPLDWrapper which **replace** the need to run the perl script run_LAMPLD.pl provided by LAMP-LD software. The wrapper can be executed for a single or a range of window size(s).

The wrapper automatically adds the positional parameters in the order expected by LAMP-LD, which enables the user to send the configuration messages in any order. Adding ancestry population files (.hap) requires to provide a positive Integer with the population order through #addPopFile:atOrder:

Usage example for 2-way populations with a single window size (50) :

```
| samplesDir |

samplesDir := BioLAMPLDWrapper testFilesFullDirectoryName / '2-way'.
BioLAMPLDWrapper new
  unolanc2way;
  windowSize: 50;
  nrFounders: 2;
  positionsFile: (samplesDir / 'chr.pos') fullName;
  addPopFile:    (samplesDir / 'pop1.hap') fullName atOrder: 1;
  addPopFile:    (samplesDir / 'pop2.hap') fullName atOrder: 2;
  genosFile:     (samplesDir / 'genofile.gen') fullName;
  outputFile:    'lamp-ld_example1.out';
  execute.
```

To execute for several window sizes, say from 30 to 150, you can specify an **Interval** instance :

```
| samplesDir |

samplesDir := BioLAMPLDWrapper testFilesFullDirectoryName / '2-way'.
BioLAMPLDWrapper new
  unolanc2way;
  windowSizes: (30 to: 150);
  nrFounders: 2;
  positionsFile: (samplesDir / 'chr.pos') fullName;
  addPopFile:    (samplesDir / 'pop1.hap') fullName atOrder: 1;
  addPopFile:    (samplesDir / 'pop2.hap') fullName atOrder: 2;
```

```
genosFile:      (samplesDir / 'genofile.gen') fullName;  
executeWS.
```

HMMER

To test if the HMMER wrapper is working properly for your platform evaluate:

```
BioHMMERWrapper new  
  help;  
  execute.
```

Successful execution should answer a BioResultObject with results:

Now this is translated code from the HMMER tutorial:

```
| samplesDir hmmFile stoFile fastaFile outputFile |  
  
samplesDir := BioObject testFilesFullDirectoryName.  
" hmmbuild globins4.hmm tutorial/globins4.sto "  
hmmFile := samplesDir , 'globins4.hmm'.  
stoFile := samplesDir , 'globins4.sto'.  
BioHMMERWrapper new  
  addParameter: hmmFile;  
  addParameter: stoFile;  
  runHmmBuild;  
  execute.  
  
" hmmsearch globins4.hmm uniprot sprot.fasta > globins4.out "  
fastaFile := samplesDir , 'globins45.fa'.  
outputFile := samplesDir , 'globins4.out'.  
BioHMMERWrapper new  
  runHmmSearch;  
  addParameter: hmmFile;  
  addParameter: fastaFile;  
  addOption: 'o' value: outputFile;  
  execute.
```

Phyx

- Work in progress *

Phylogenetics

Phylogenetics use genetic variation between species and populations (taxa) to infer evolutionary relationships.

- Work in progress *

Nexus File Format

Newick File Format

Newick is a file format used for representing phylogenetic trees in computer-readable form. In contrast to Nexus files they contain no further syntax elements or other information than the trees. representing rooted trees with labeled nodes and specified lengths between parent and child

- Trees are represented as an annotated nested list of node names
- Node names are a letter optionally followed by more letters, digits or an underline
- Node names can be followed by a length separated from the name by a colon. For example (ant:17, (bat:31, cow:22):7, dog:22, (elk:33, fox:12):40);
- Optionally one can label internal nodes. For example: ((cow:12, gnu:10)bigThings:3, (ant:23, bat:19)smallThings:5);
- Officially the tree ends in a semicolon (redundant since the parens must balance)
- The standard also does not allow length for the distance for a line leading to the root since this has no meaning in the context of phylogeny

Clustering

Graphics

Graphics are built and rendered using the [Roassal agile visualization engine](#), which is already installed in the BioPharo distribution. The Roassal engine includes a large set of examples using short Smalltalk expressions.

DotPlots

A simple visualization for the comparison of two sequences (or one sequence against itself) is [DotPlot](#), which provides an overview of the extent of similarity between sequences without diving into details. DotPlots operates by projecting two sequences to be compared along the horizontal and vertical axis of a Matrix. DotPlots requires three main parameters:

- A "Word Size", which is a **Number** representing the length of the identical polynucleotide or polypeptide that must be found in both sequences in order to generate a dot on the chart. The bigger the word size, the lower is the probability that the same segment is present in both sequences.
- A "Scoring Scheme", specified as a **Matrix**. The simplest scoring scheme is the Identity Matrix, but others can be used, for example the EMBOSS DNA scoring matrix. Other extended matrix is used when sequences contain ambiguity codes, or using proteins.
- A "Threshold" or "cut-off score".

When plotting mummer output, it is necessary to use the lengths of the input sequences to set the plot ranges, otherwise the plot will be automatically scaled around the minimum and maximum data points

Any character or symbol that does not belong to the [a-zA-Z] set is ignored.

Drawing a 2D DotPlot with two random sequences

```
BioRTDotPlot
  x: BioSequence atRandom y: BioSequence atRandom;
  window: 3;
  threshold: 2;
  open.
```

Drawing a 2D DotPlot for single analysis (identify repeats in a sequence)

```
BioRTDotPlot
  sequence: BioSequence atRandom;
  window: 3;
```

```
threshold: 2;  
open.  
open.
```

Bar Charts

Under construction, please contact mailing list to collaborate.

Histograms

Histograms are frequently used to get an overview of quantitative distributions. Examples where histograms are applied are: Distribution of codon usage, sequence lengths distribution, etc.

Drawing a histogram of a FASTA file with multiple sequences

```
(BioParser parseMultiFastaFile: BioObject testFilesDirectoryName asFileReference / '  
ls_orchid.fasta')  
plotHistogramBins: 20  
xAxisLabel: 'Sequence lengths (bp)'  
yAxisLabel: 'Count'  
color: Color red
```

Widgets

BioSmalltalk includes widgets to provide user interfaces. These utilities are small re-usable windows which could be easily integrated into your applications.

Chromosome Selection

A small utility to select organisms and their chromosomes for posterior analysis. Each organism has already configured their corresponding list of chromosomes.

[BioChrSelector](#) open.

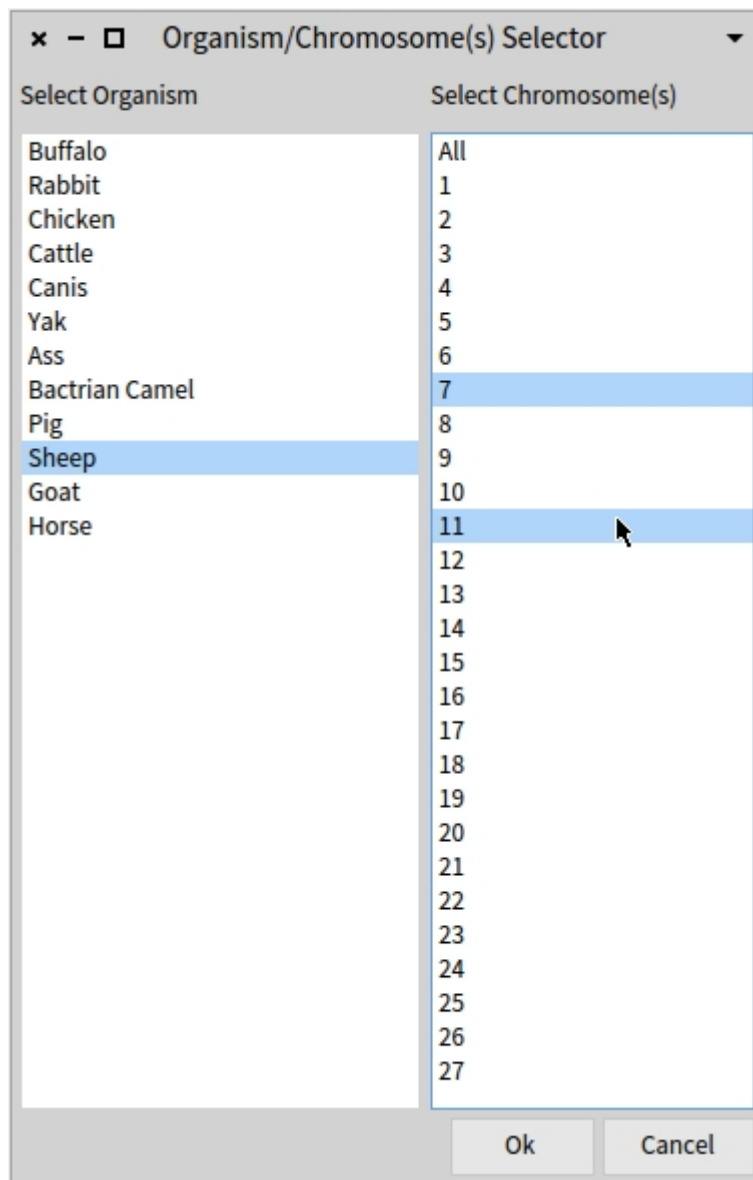


Figure 3: Selecting MAP file

CSV Explorer

The BioCSVExplorer class builds a widget that allow users to visually explore a CSV file and import its contents to the image.

Work in progress

Matrix Explorer

Work in progress

Third Party Libraries

This section includes description of libraries not specific to "Bio" tasks, but useful for development of any software utility.

Persistency

All current Pharo images includes Fuel, a fast binary serializer that defines a pickle format based on clustering similar objects. It reduces deserialization overhead compared to other approaches. It supports class shape changing (when a variable has been added, or removed, or its index changed), cyclic object graphs, serialization of several special objects like **CompiledMethod** (in a binary way, without compiling code) and Classes, etc.

Mathematics

The package *Numerical Methods* includes common numerical algorithms implemented in Smalltalk. The package library is documented in the free downloadable book [Object-Oriented Implementation of Numerical Methods. An Introduction with Smalltalk](#), which covers all the available algorithms:

- Interpolation: Lagrange, Newton, Neville, Bulirsch-Stoer, Cubic Spline.
- Iterative algorithms
- Integration methods: Trapeze, Simpson, Romberg.
- Series: Infinite Series, Continued Fractions, etc.
- Linear Algebra: Vectors, Matrices, Linear Equations, LUP decomposition, Eigenvalues and Eigenvectors.
- Distributions: Normal, Gamma, etc.
- Statistical Tests: Chisquare-test, F-test, T-test, Linear Regression, Maximum Likelihood, etc.
- Optimization: Hill Climbing, Powell, Simplex, Genetic, etc.

OS Communication

The package `OSProcess` provides access to operating system functions, including pipes and child process creation. It is implemented using pluggable primitives in a shared library for Unix or Linux, and a DLL for Windows. The Smalltalk code, including the classes which implement pluggable primitives for Unix or Win32 operating system functions, may be loaded into any Squeak image, but the primitives are only useful on Unix and Windows systems. Placeholder classes are provided for MacOS, OS/2 and RiscOS, but are not yet implemented.

`CommandShell` is a Smalltalk implementation of a command processor shell and terminal window. It is intended to behave like to a simple terminal window (like `xterm`) running a Unix command shell (like `/bin/sh`). It lacks some elements of Unix shell syntax, and does not provide terminal emulation, but it adds some nice Smalltalk enhancements such as a text editor which works in a command pipeline, and the ability to evaluate Smalltalk expressions in a command pipeline with unix commands.

Developing BioSmalltalk

Code Organization

User Code Organization

At this point, you probably have your own scripts and wonder how to maintain your source code. You have several choices, depending on what is your task on course:

The Script Manager

The Script Manager is a script browser tool developed by Torsten Bergman which enables to organize your scripts. It allows to have multiple workspaces organized by categories and within each category, to name scripts, execute, save and export them. It also supports syntax highlighting and autocompletion. Unfortunately the tool was removed from official Pharo images, but you can install it through the Catalog Browser:

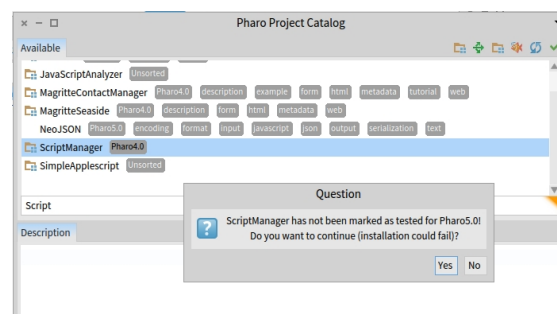


Figure 4: Installing the Script Manager

To open the tool, go to "Tools" -> "ScriptManager"

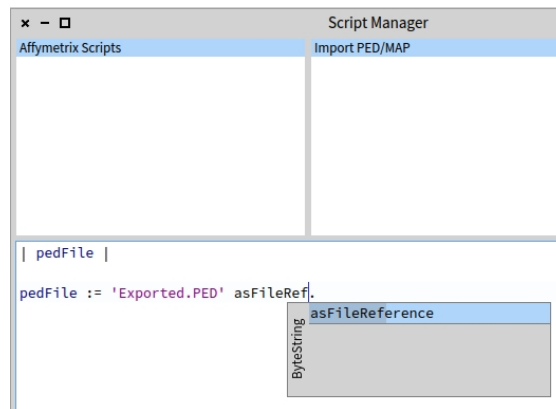


Figure 5: Using the Script Manager

The Nautilus Bio Browser

To facilitate working with Bio API, all the library classes are grouped in the System Browser:

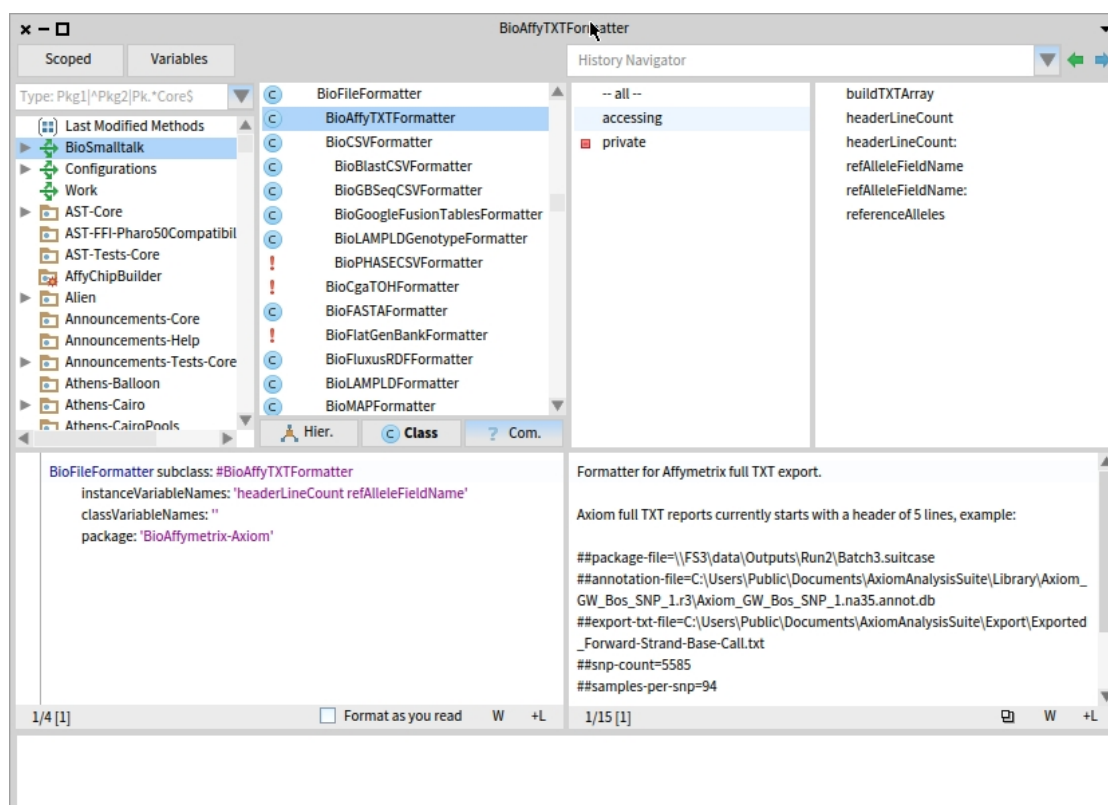


Figure 6: Browsing the BioSmalltalk Classes

Developer Code Organization

Packages are organized in the following way:

Table 2: BioSmalltalk Code Organization

Package	Description
BioBlast	Desktop and web wrappers, filter and reader classes for Blast results and software
BioClassifier	A basic rule engine system and feature classifier
BioEntrez	Access to the Entrez Utilities API
BioFormatters	Basically file format writers. Naming strategies and other utilities also included here
BioTools	Samples could be methods returning String's for other subsystems or Scripts writing output to files
BioToolsSamples	Not necessary, the snapshot system will persist your source code
BioSupport	Classes and methods working in all Smalltalk platforms
BioPharo	Code specific for the Pharo platform
BioSqueak	Code specific for the Squeak platform
BioNCBI	Code to access NCBI's QBLAST service and convenience utilities.
BioParsers	Contains all the parsers. Patches for specific platforms are in their respective packages.
BioPopulation	Classes and methods for working with population genetics.
BioProject	A basic application framework for organizing work in projects.
BioWrappers	Core for wrapping desktop and web programs + specific wrappers for common bioinformatics software.

Building from Source

To build an image from source you will need at least `wget` or `curl` to download a Pharo clean image. In all cases, the following instructions creates a BioSmalltalk directory and download a clean Pharo image. Packages are fetched from SmalltalkHub, although there are plans to migrate source code to GitHub. You can build the BioSmalltalk image in two ways: Interactive or non-interactive.

Non-interactive building in Windows

It is highly recommendable to download and install MSYS Shell from the MinGW project. Additionally, you should install the `msys-wget`, `msys-zip` and `msys-unzip` packages from the MinGW Package Manager.

```
mkdir /c/BioSmalltalk
cd /c/BioSmalltalk
wget -O- get.pharo.org/50+vm | bash -
./pharo Pharo.image config "http://smalltalkhub.com/mc/hernan/BioSmalltalk" "
ConfigurationOfBioSmalltalk" --printVersion --install=development
```

Non-interactive building in GNU/Linux

```
mkdir /usr/local/src/BioSmalltalk
cd /usr/local/src/BioSmalltalk
wget -O- get.pharo.org/50+vm | bash -
./pharo Pharo.image config "http://smalltalkhub.com/mc/hernan/BioSmalltalk" "
ConfigurationOfBioSmalltalk" --printVersion --install=development
```

Interactive Building

Interactive building enables you to view the installation progress and to live debug installation problems. You should start the image from command line, and execute a Smalltalk installer script. To start the Pharo image evaluate:

```
./pharo Pharo.image &
```

One inside Smalltalk, you can evaluate the installer using Gofer or Metacello. Gofer is an API on top of [Monticello](#), a package system of Smalltalk. Both Gofer and Monticello are preinstalled in current Pharo images. The following script uses Gofer:


```
Gofer it
smalltalkhubUser: 'hernan' project: 'BioSmalltalk';
configurationOf: 'BioSmalltalk';
loadDevelopment
```

Metacello is a newer cross-smalltalk package management system. Supports named versioning, dependencies, metadata, conditional package loading and other features.

```
Metacello new
smalltalkhubUser: 'hernan' project: 'PhyloclassTalk';
configuration: 'PhyloclassTalk';
version: #development;
load.
```

Testing

Testing is the evaluation of your developed software product. Being Smalltalk the technology where the original testing framework was invented (SUnit), writing comprehensive tests is a traditionally important part of the package validation process.

A rapid and excellent introduction to Unit Testing in Smalltalk can be found in the Stéphane Ducasse article: SUnit Explained Revisited <http://www.cs.pdx.edu/~black/OOP/papers/SUnitEnglish2.pdf>. A tutorial about testing is beyond the scope of this guide, however I would like to introduce some issues about the nature of test writing before starting to write tests:

- Assume the presence of errors, always.
- It's impossible, generally, to find all errors.
- Test to find errors, not to show that it works.
- Not finding errors doesn't mean that your software doesn't contain errors.
- A successful test run discovers errors!
- Write for both valid and invalid conditions.

More useful advices are found in the Pharo By Example book: <http://pharobyexample.org/>

Executing Tests

You can execute tests by selecting the method you want to test in the Browser:

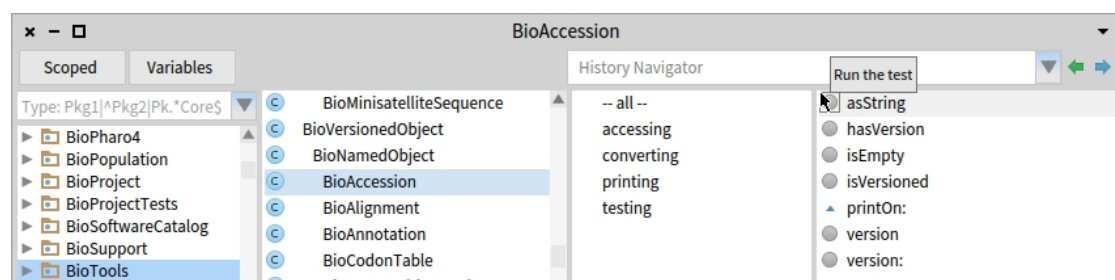


Figure 7: Selecting method to be tested

Once clicked, if the test pass the icon will be displayed in green color.

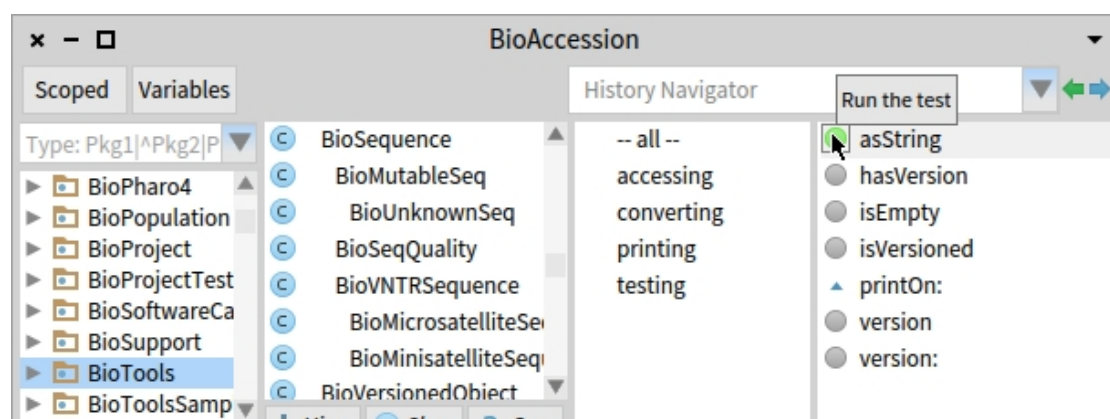


Figure 8: Executing a test on a method

Writing Tests

These are guidelines for writing tests in BioSmalltalk. Eventually they will be mandatory.

- All test classes should inherit from **BioAbstractTest** class
- Tests should be in the “testing” method category
- Tests selectors should start with the prefix “test”.

Test Data

Including your own test data set:

Each BioSmalltalk distribution must include a “BioSmalltalkTestFiles” subdirectory (from now the “test files subdirectory”) to place all the files used for the test suites.

- If you want to include data sets in your tests, then place them as files in this subdirectory.
- If your data is bigger than 8 Kb, then compress in one or several files and send appropriate messages to uncompress them.
- Do not override the contents or edit the file index.html in the test files subdirectory, it is needed for serving a “document root” when testing the local web servers.
- Testing of databases results is NOT encouraged.

Using already existing data sets in the image:

- In your BioSamples subclass, instance side, add a category named “samples-data” and add appropriate sample data. Beware that each Smalltalk flavor has their own limits in the method size and supporting huge data sets in the image is not encouraged.

```
BioYourTestCase>>testFilesDirectoryName
```

```
^ super testFilesDirectoryName , BioObject fileDirectoryClass slash , '  
  GenBankTestFiles'
```

SubSystems

The Engines SubSystem

BioSmalltalk supports interchangeable strategies by enabling the selection of the most appropriate underlying software package (called engine in our context) or to discover the best option for specific tasks. This is because in Smalltalk flavors very often you have several packages to perform the same task.

The developer should never hard-code a specific external replaceable package class name in a method, instead, he should write the corresponding adapter subclass and link it to the particular package (see below). Then the BioSmalltalk user could use a generic name for the particular tasks, instead of referring to implementations, which could fail, become unsupported or too slow. For example, to serialize your objects, use

```
BioSerializationEngine serialize: anObject.
```

(In Pharo: to materialize your serialized file just open the File Browser, select the file and Materialize from the right click menu)

To get the contents at an URL

```
(BioWebEngine for: anURL) httpGet
```

Adding an Engine

BioAbstractAdapter class defines the interface and abstract behavior for Adapters.

- What we call “Providers” refers to the external classes which are adaptee (**OSPlatformExecutor**, **KomServer**, etc). They could be missing in your platform.
- What we call “Adapters” refers to the final subclasses acting as adapters (prefixed by “Bio” for example: **BioOSPlatformExecutor**, **BioKomServer**, etc)

The information of which is the current engine in a hierarchy of engines is maintained by the AdapterClass class variable in the corresponding hierarchy. To access the current adapter for an engine evaluate:

```
BioSerializationEngine adapter
```

The information of which is the current provider in the current engine is maintained by the “provider” instance variable.

If you want to add your engine to BioSmalltalk, these are the required steps:

- 1) Check for the kind of engine you want to support at subclasses of BioAbstractAdapter, currently there are engines for:
 - Web clients (The historical HTTP client from Squeak, WebClient from Andreas Raab, ZnHttpClient from Sven Van Caekenberghe)
 - Serialization clients (Fuel, SRP, SmartReferenceStream, ReferenceStream, Ma object serialization, BOSS, SIXX/SIF)
 - OS Execution
 - .

If there is no engine type listed, you will have to check by yourself specific users of your methods and refactor code to support multiple strategies, this can be as easy as changing one method, adding a superclass and implementing some template methods, or as complex as reorganizing a whole set of BioSmalltalk hierarchies.

- 2) Create your subclass and implement these CLASS methods:

```
#clientClass
" Answer a Symbol, the main external class name for this adapter "
```

Remember to name the class as a Symbol, do not assume the client has installed your adapter. To answer the actual class use for example: `^ self classFor: #WebUtils`

```
#isPreferredAdapter
" Private – Answer <true> if the receiver is the best option for the current image and platform "
^ true
```

Don't assume the client will have ANY client actually installed, so you might want to check if `#providerIsAvailable`

The Parsing SubSystem

This intends to be a general guide for writing a parser in BioSmalltalk. Parsing is a whole area in its own, and writing a parser generally is very difficult and takes time, so be careful when choosing your learning sources. To build a framework to start learning about writing parsers, the following concepts could be useful:

- Determine that you actually need a parser. Recall: A parser is an object that scans an input String, decomposes it into its constituent parts (which you may define), and then optionally processes the parts in some suitable manner.
- How: Most tutorials calls writing a parser “by hand” when you write a parser without any support programming library, so developers use a “parser generator” and write rules for this generator. Why: Because keeping track of contingencies caused by different inputs and input paths is extremely difficult, and a generated parser code can be proved mathematically.
- **BioParser** is the root class for parsers. Parsers use several parsing libraries
- **XMLSupport** for parsing in DOM or SAX mode (see below to understand XML parsing modes)
- **XMLPullParser** for StAX parsing (see below to understand XML parsing modes)
- **PetitParser** for parsing formatted or structured text. If you don’t know anything about PetitParser a good guide to start is <http://scg.unibe.ch/research/helvetia/petitparser>

Adding a new parser

Create the access point in the BioParser class

- Instead of forcing the developers to find the appropriate parser subclass on usage, the BioParser class was implemented as a Facade pattern providing class methods which hides the complexity of the hierarchy, this is, put every access point (method) in the **BioParser class** side.
- Comment the implemented method, this is mandatory.

Create the parser class

- If you want to parse formatted/structured text, subclass from BioAbstract-TextParser or specific subclasses (read class comments to find out which superclass is appropriate for your parser)

- If you want to parse XML, regardless of the mode, subclass from BioXML-Parser

Create the main access method(s) for your parser:

- If you need to return a **Collection** object (**Array**, **OrderedCollection**, etc.), then prefix the selector with #tokenize, example: #tokenizeLocus:, #tokenizeFastaHeader:, #tokenizeSwissProtEntryName:
- If you need to return a BioObject, then prefix the selector with #parse, example: #parseFasta: will return a BioFastaRecord, #parseAmbiguousWithSeparators: will answer a BioSequence with a **BioIUPACAmbiguousDNA** as alphabet.

The BLAST XML parser

Our motivation for building the initial BLAST XML parser is the chance of having customized DOM sub-trees in memory that could be subsequently filtered through messages, not by specifying node names which could change, be sensitive to typos, and are not suitable for performing computational reflection on them. For example, given a `blastXmlTree` object, instead of filtering by

```
blastXmlTree filter: [: hit |  
  hit name = 'Hit_def and: [ hit value = 'Homo sapiens' ] ]
```

the idea is to add semantics to our tree:

```
blastXmlTree hitDefinitionsEqualTo: 'Homo sapiens'.
```

which enables to modify completely the implementation or algorithms in the method `#hitDefinitionsEqualTo`: preserving the same functionality for the programmer. The key in parsing BLAST XML contents are the `#select...` methods to narrow the parsing to specific element groups. This is, instead of building the whole XML tree into memory (which would be the DOM technique) the idea is that users will only be interested in specific DOM sub-trees.

The Entrez SubSystem

The main class of this subsystem is **BioEntrezClient**, it is the only class that should be visible for users. The database classes are located in the **BioEntrez-Database** hierarchy, below in this document there is a step guide to add a new database. The database hierarchy does not validate currently the support rettype's and retmodes as described in http://www.ncbi.nlm.nih.gov/books/NBK25499/table/chapter4.chapter4_table1/?report=objectonly

This means the developer is responsible to build correctly the request with valid parameters.

The entrez commands are supported through the **BioEntrezCommand** hierarchy, which contains code to validate command parameters and the base url of the requests. Currently the mandatory documentation for implementation is located at: <http://www.ncbi.nlm.nih.gov/books/NBK25499/>

The code samples are located at **BioESamples**, when providing new features, a new method should be created at the appropriate method category.

Steps for adding an Entrez database

- If the database supports retmode and rettype parameters, subclass from **BioERetDatabase**
- If the database supports any of the following parameters, subclass from **BioESeqDatabase**: strand, complexity, seq_start, seq_stop
- If the database is any of the Literature Databases (PubMed, PubMed Central (PMC), Journals, OMIM), subclass from **BioELitDatabase**
- Add a new method named #dbName in class side of the new recently created class, for example

```
dbName
  " See superimplementor's comment "
  ^ 'sra'
```

- Add a new method in instance side of class **BioEntrezClient**, category "accessing public – databases" for example:

```
biosample
  " Set and answer the receiver's working database "

  ^ self database: (BioEntrezSample new client: self)
```

- Add a simple test case for the database name in the #testDbClassOf method, example:

```
self assert: (BioEntrezDatabase dbClassOf: 'biosample') = BioEntrezSample.
```

The Files SubSystem

- A **BioFile** is used to represent a file, it's not intended to be used for parsing a file. This is useful when having a **BioCSVFile** which we need to set specific parsing parameters, for example, numbers of columns to read or ignore header lines.
- A **BioFormatter** is what to use when a developer has an INPUT file and possible an OUTPUT file/stream. It is not necessarily a writer.

File Formats

File formats information is available through the class **BioFileFormat** which contains a registry of file formats information. Included information doesn't mean that the format is supported for reading or writing though, it just acts as an internal small database of formats available. You may query if a format is registered by evaluating:

```
BioFileFormat isRegistered: 'CSV'.  
BioFileFormat fileExtensionFor: 'FASTA'.
```

For example to register a new format for CSV files, evaluate;

```
BioFileFormat  
  registerFormat: 'CSV'  
  extensions: #('csv')  
  description: 'General ''Comma'' Separated Values format'  
  rootClass: BioCSVFile
```

The root class is not necessary only for registration purposes, so for now you can pass nil as parameter. In the future a GUI wizard will be provided to generate code automatically for registering and unregistering file formats.

Troubleshooting

In case Browsing tools are broken, you can still obtain message senders by evaluating:

```
SystemNavigation new allLocalCallsOn: '' ofClass: ProtoObject .
```

Community

Mailing lists

A mailing list for users is available as [biosmalltalk-users](#).

A mailing list for developers is available as [biosmalltalk-developers](#).

Beware that posts are (or will be) publicly available through popular search engines and mailing lists archives. If you want to create another list for a particular package or subsystem, contact the project manager through the project page.

To ask questions please use the etiquette suggested at: <http://www.catb.org/~esr/faqs/smart-questions.html>