

# Documentation of S-MART

Matthias Zytnicki

May 7, 2010

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Installation and requirements</b>	<b>2</b>
2.1	For Windows . . . . .	3
2.2	For Linux or Mac . . . . .	4
2.3	Test the configuration . . . . .	6
<b>3</b>	<b>General description</b>	<b>6</b>
<b>4</b>	<b>Which tool for your need?</b>	<b>8</b>
4.1	Mappings . . . . .	8
4.1.1	Data comparison . . . . .	9
4.1.2	Sliding Windows . . . . .	12
4.1.3	Merging data . . . . .	12
4.1.4	Select part of the data . . . . .	14
4.2	Working on Sequences . . . . .	14
4.3	Modify your data . . . . .	14
4.4	Visualizing you data . . . . .	15
4.5	Conversion tools . . . . .	17
4.6	Other tasks . . . . .	19
<b>5</b>	<b>Possible pipe-lines</b>	<b>19</b>
5.1	Use two mappers . . . . .	20
5.2	Find piRNA clusters . . . . .	20
5.3	Get the letter distribution of the beginning of the data . . . . .	20
5.4	Compare two sets of reads with sliding windows . . . . .	21
5.5	Compare RNA-Seq with tiling arrays using sliding windows . . . . .	22
5.6	Compute differential expression . . . . .	23
<b>6</b>	<b>More about S-MART</b>	<b>25</b>
6.1	Data structures . . . . .	25
6.2	Tags . . . . .	25
6.3	How mySQL is used? . . . . .	26
6.4	Contribute to S-MART! . . . . .	26

<b>7</b>	<b>Contact</b>	<b>26</b>
<b>A</b>	<b>Load data on your genome browser</b>	<b>27</b>
<b>B</b>	<b>Get other data</b>	<b>27</b>
<b>C</b>	<b>Troubleshooting</b>	<b>27</b>

## 1 Introduction

S-MART should be pronounced “ess-mart” (better pronounced with a slight Spanish accent) and stands for “Short reads MART”. It could have many other meanings. The one you choose is the best.

It provides a set of Python scripts which transform your short reads which have been mapped to a genome. So, it supposes that you already have mapped your data to a reference genome. For that, you can use Maq, ZOOM, Mosaik or any other tool. S-MART supports many formats.

You can also compare S-MART with other data, such as RefSeq sequences or any kind of annotation, as long as their formats are supported by S-MART (which is usually the case). However, S-MART does not include these data, simply because they are too many of them, and too many organisms. If you want to know where I get my data, read Appendix B.

## 2 Installation and requirements

Depending on the system you are using installation can be different. However, in both cases, you will need Python, MySQL, R and Java.

**Python** Python is the language used to code the algorithms. Why Python? Well, why not?

**MySQL** MySQL is a database management system which is used to handle efficiently the reads (remember that there can be several millions reads, so that every algorithmic improvement is highly important). Normally, you should not even see that S-MART uses some databases, since it reads flat files (like GFF files), and outputs flat files. But, internally, it uses databases. If you want to know how MySQL is used, you can read Section 6.3.

**R** R is a statistical computing tool which can do many things, but, here, it is only used to plot the data.

**Java** Java is used here simply because it contains a good graphical user interface. So, it is used for the GUI in S-MART.

## 2.1 For Windows

Setting up your system for S-MART should not take more than 15 minutes.

**Python** You should have downloaded and extracted a bundle which contains all the Python scripts files. First check that Python version 2.5 is installed. You should get it from their Web site.

Each Python script uses many other scripts (basically, the structure of the files is organised by classes). So, you should add to your PATH variable the directory where you have installed the scripts. To do so, click on **My Computer**, then **Control Panel**, **System**, **Advanced**, **Environment variables**. Click on **New** and add the following variables:

name: PATH, value: %PATH%;*where\_you\_installed\_Python* (probably  
C:\Python25)

name: PYTHONPATH, value: *where\_you\_installed\_S-MART*

Now, download a Python extension for MySQL and install it.

**mySQL** You can download mySQL Server 5 *via* their Web site. Install it. You might need to reboot in order to start the mySQL daemon (yes, there will be a kind of demon in your computer... the kind which makes Word buggy, I guess) and create the sockets (yup, your computer needs sockets, otherwise it will catch a cold).

Install it, choose the **Standard Configuration**, then default options. At some point, you will have to create some user accounts. Create root account and keep the password.

Check that one user has read and write rights granted. You can do it by starting mySQL (you should find mySQL Command Line Client in the **Programs** menu) and write:

```
CREATE USER 'smart'@'localhost';
GRANT ALL PRIVILEGES ON *.* TO 'smart'@'localhost'
WITH GRANT OPTION;
```

Create also a database that S-MART will be able to use.

```
CREATE DATABASE smart_data;
GRANT ALL ON smart_data.* TO 'smart'@'localhost';
```

Close the window.

If you decided to choose a password for the mySQL, or a different login, you will have to modify consequently the file called `.pythonConnection.txt` which is the directory where you install S-MART. By default, the file contains:

```
user = smart
host = localhost
password =
database = smart
```

Write the corresponding login, password or database in the appropriate field.

**R** Again, download R *using* CRAN mirror download page.

Two packages are furthermore needed: RColorBrewer, to have a good palette of colors (I am color blind, so it is important to me), and Hmisc, which does Spearman correlations (among others). There are many ways to do so. The simplest one is to start R (find it in you **Programs** menu), and write:

```
install.packages("RColorBrewer", dependencies = TRUE)
install.packages("Hmisc", dependencies = TRUE)
```

Quit with `q()`. Done.

**Java** In the very unlikely case you had not Java (it is now a standard tool used with you Web browser), you can download it and install it afterwards.

## 2.2 For Linux or Mac

This has been tested on a Debian 2.22.3 and a MacOS 10.6.3.

**Python** You should have downloaded and extracted a bundle which contains all the Python scripts and the Java files. First check that Python is installed. I can guarantee that it works for version 2.5.2. Versions earlier than 2.4 will not work. Later versions may work. Or not.

Each Python script uses many other scripts (basically, the structure of the files is organised by classes). So, you should add to your PATH and PYTHONPATH variables the directory where you have installed the scripts. To do so, supposing you use the standard Bash shell, open your `.bashrc` file on your root directory and add the following line at the end of the file:

```
export PATH=$PATH:the_directory_where_you_installed_the_files
export PYTHONPATH=$PYTHONPATH:the_same_directory
```

Some nice guy has made a mySQL API for Python, which is useful for S-MART. You could probably download it from your package repository (look for `python-mysqldb`).

Alternatively, you can download the package from Sourceforge and install it with:

```
python setup.py build
sudo python setup.py install
```

**mySQL** You can download mySQL *via* your usual package manager or using their Web site. The Mac users can use the installation help page if they experience some problems. Just in case, reboot to be sure you have the daemon started and the sockets created. Check that one user has read and write rights granted. You can do it by starting mySQL in a console with super-user rights (type `sudo mysql` in a console for instance) and write:

```
CREATE USER 'user_name'@'localhost';
GRANT ALL PRIVILEGES ON *.* TO 'user_name'@'localhost'
WITH GRANT OPTION;
```

Create also a database that S-MART will be able to use.

```
CREATE DATABASE your_database;
GRANT ALL ON your_database.* TO 'user_name'@'localhost';
GRANT FILE ON *.* TO 'user_name'@'localhost';
```

You can quit mySQL with Ctrl-D.

Modify in your parent directory the flat file called `.pythonConnection.txt`. It should have the following content:

```
user = your_user_name
host = localhost (most probably)
password = (probably empty)
database = the_database_you_created
```

**R** Again, download R *via* your package manager or using the CRAN mirror download page.

One package is furthermore needed: RColorBrewer, to have a good palette of colors, and Hmisc, for the Spearman correlations. There are many ways to do so. The simplest one is to start R (type R in a console), then write:

```
install.packages("RColorBrewer", dependencies = TRUE)
install.packages("Hmisc", dependencies = TRUE)
```

Quit with `q()`. Done.

In case it does not work, you will have to download the RColorBrewer and Hmisc packages, then install them with:

```
R CMD INSTALL -l where/you/downloaded/the/packages RColorBrewer_xxx.tar.gz
R CMD INSTALL -l where/you/downloaded/the/packages Hmisc_xxx.tar.gz
```

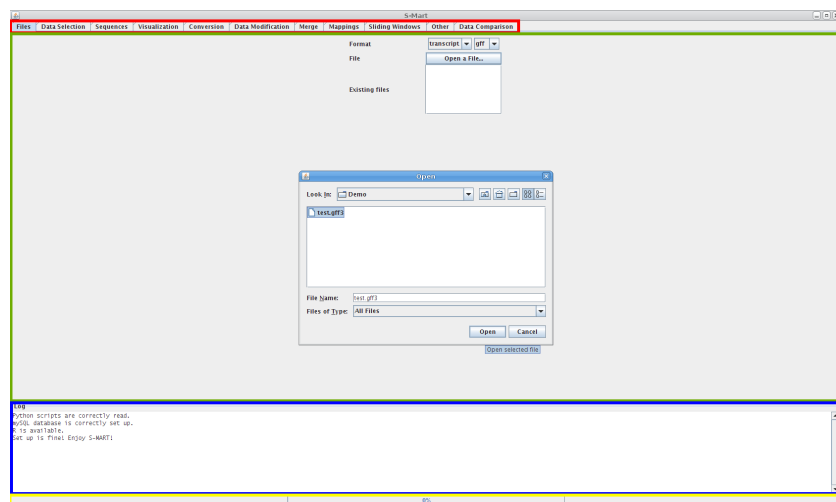


Figure 1: The “file” panel of the GUI.

**Java** I guess you already have a Java somewhere, have not you?

## 2.3 Test the configuration

To check if everything is correctly set up, double-click on the **Smart.jar** icon in the directory where you installed S-MART. The first you start it, S-MART might ask you where the executable files for Python and R are (for the Windows users, they might be at `C:\Python25\Python.exe` and `C:\Program Files\R\R-2.10.0\bin\R.exe` respectively, depending on the versions of you executables).

Is it OK? Cool! You can now start with your S-MART experience. Otherwise, you may drop me email and I might help you with the configuration.

## 3 General description

**The GUI** I have developed a graphical user interface so that every tool can be started easily.

On Windows, start it by double clicking on the **Smart.jar** file. On Linux, type `java -jar Smart.jar`.

You have to set the files that you will use (together with their formats) in the first panel. In the other panels, you will be able to start all the tools S-MART contains.

Figure 1 shows the GUI on the “file” panel. The GUI is divided into four regions, from top to bottom. On the top region (in the red rectangle), you can select the type of action you want to perform by selecting the right panel. The second area (the green one) lets you perform the task. The third region

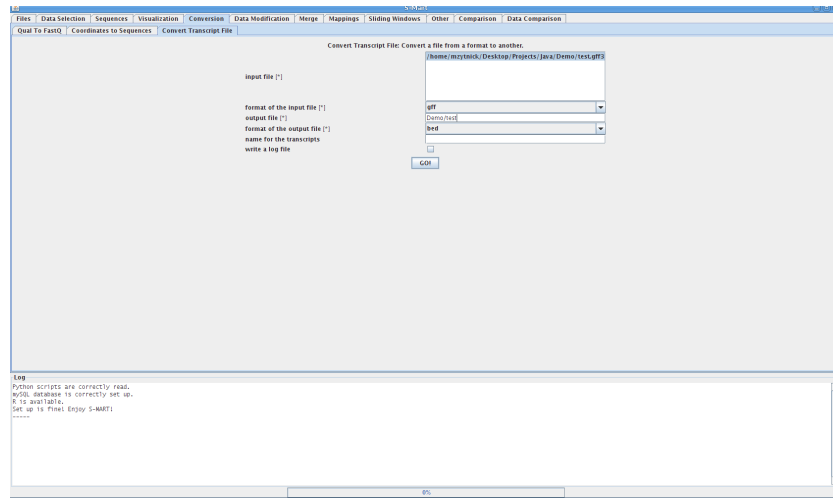


Figure 2: The “Convert Transcript File” program panel of the GUI.

(the blue one) is the log area, where you can interactively read the output of the programs, or any relevant information. The last area (the yellow one) is a progress bar and shows you how much time the program will run to perform the task.

On the “file” panel, you can submit the files that you are going to use, together with their format. In the example, we enter the file `test.gff3`, which is a transcript list in GFF3 format. First select the type of data: mapping data (coming from your mapper), transcripts and other files. Then, select the right format. As you can see, S-MART supports many formats. Finally, click on the button **Open a File** to browse your hard disk and select the right file.

You can then use any tool of the toolbox by changing the panel. Figure 2 shows a conversion utility tool. Then, we select the file that we have mentioned (`test.gff3`) to convert it into a BED file. We specify that the input file is in GFF3 and that the output file is in BED format. We also specify the output file name (do not write the extension: S-MART adds it by itself). We click on the button, the program starts, as visible in the log area. Finally, the output file appears. We can open it in a file browser (see Figure 3)

**The command line** For the real hackers, every tool can be used in command line. All the scripts are in the `Python` directory and you can start them there. They all have several parameters that you can adjust depending on what you want to do, so a typical command would be:

```
python mapperAnalyzer.py -i mappedData.psl -f psl
-q rawData.fasta -o coordinates -n 1 -s 100 -m 0
-p 0 -e -x -r -b -B -g -G -u -U -2 -y -c green
-t shortReads -v 50 -l
```

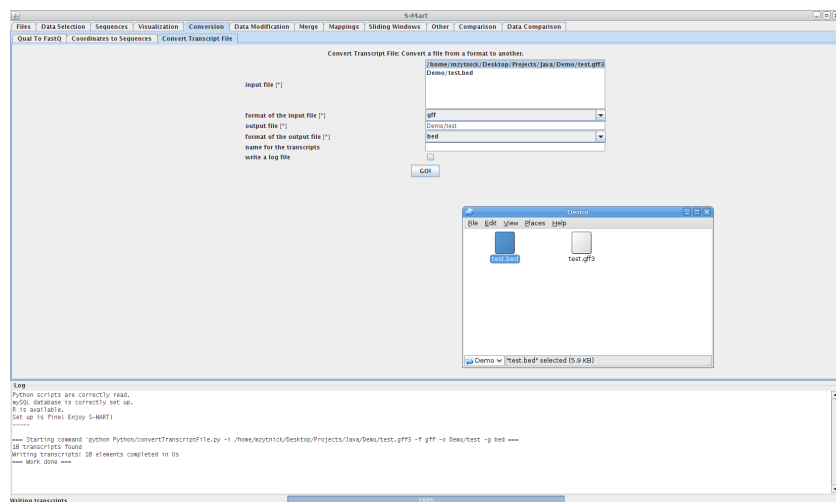


Figure 3: The result of the tool.

(It should be only one line but it does not fit in the page. Do not be afraid, commands usually are shorter.)

An important parameter is the `-v number` option, which gives the verbosity level (highest is most verbose). Another general option is the `-l` option, which writes a log file (usually utterly verbose, but sometimes useful). A last useful option is `-h`, which displays a notice, a comment for each option and exits. Moreover, the `-y`, which is sometimes available, keeps the output file in the internal representation of S-MART (which is in a MySQL database). If you want to know why it could be useful to use this option, read Section 6.3.

## 4 Which tool for your need?

This section presents the scripts that you may want to use for a particular task.

### 4.1 Mappings

**mapperAnalyzer.py** The first program you may use is **mapperAnalyzer.py**. It reads a set of mapping given by the tool you have used to map your data on the reference genome and translates it to a set of genomic coordinates. You also have the possibility to extract only those that you are interested in (few matches in the genome, few errors in the mapping, etc.). You can also select those reads which map less than a given number of times in the genome. Moreover, you can output the data in various different formats, which you can use to visualize them *via* UCSC genome browser or GBrowse<sup>1</sup>. Unmatched reads are re-written

<sup>1</sup>Look at Appendix A to know more about it.



in other file, in case you would like to try to map them with another tool (may sometimes work!).

The script can parse data given by the following programs (the corresponding option is given in parenthesis):

- Blast (use `-m 8` format for Blast and `-f blast`)
- Blat (`-f ps1`)
- Exonerate<sup>2</sup> (`-f exo`)
- Maq (`-f maq`)
- Mosaik (output in axt format for Mosaik and use `-f axt`)
- Nucmer (`-f nucmer`)
- Rmap (`-f blast`)
- Seqmap (`-f seqmap`)
- Shrimp (`-f shrimp`)
- Soap (`-f soap`)
- and more...

You can filter your data according to:

- number of errors in the mapping
- number of occurrences of the mapping in the genome
- size of the read mapped
- number of gaps in the mapping

`mappingToCoordinates.py` If you just want to convert your mapping data to genomic coordinates, without any filtering, you can use `mappingToCoordinates.py`.

#### 4.1.1 Data comparison

`compareOverlapping.py` This script may be the most important one. It basically compares two sets of transcripts and keeps those from the first set which overlap with the second one. The first set is considered as the query set (basically, your data) and the second one is the reference set (RefSeq data, for example, see Figure 4). Various modifiers are available:

---

<sup>2</sup>Exonerate can display its results in many formats. Currently, S-MART only support the following output format:

```
--ryo "%S %em %V\n" --showvulgar FALSE --showalignment FALSE
```

Please add these parameters to your command line while using Exonerate!

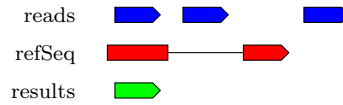


Figure 4: Simple comparison between your reads and RefSeq data, for example, using `compareOverlapping.py`.

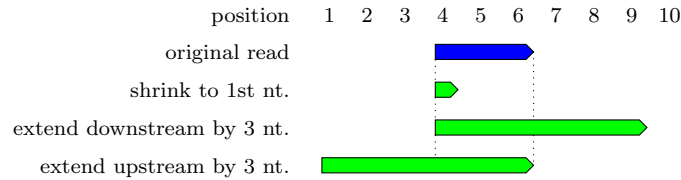


Figure 5: Shrinking and extending your data before comparison with `compareOverlapping.py`.

- Restrict query / reference set to the first nucleotide. Useful to check if the TSS of one set overlap with the other one.
- Extend query / reference set on the 5' / 3' direction. Useful to check if one set is located upstream / downstream the other one.
- Include introns in the comparison.
- Invert selection (report those which do not overlap).
- Keep colinear / anti-sense overlapping data.
- Keep the query data even if they do not strictly overlap with the reference data, but are located not further away than  $n$  nucleotide from some reference data.

The mechanism of shrinking and extending is also useful to make a fine grain comparison. For example, if you want to keep those such that the TSS is overlapping the reference set, you just shrink the query set to 1 nucleotide (see Figure 5). Now, if you want to keep those which are overlapping you data or located 2kb downstream of it, just extend the query data in the downstream direction, and you will have what you want. You can also extend in the opposite direction to get the possible transcript factor sites which are upstream.

Some option reverses the selection. Put in other words, it performs the comparison as usual, and outputs all those query data which do not overlap.

**getDifferentialExpression.py** This tool compares two sets of data and find the differential expression. One very important component of the tool is the reference set. Actually, to use the tool, you need the two input sets of data, of course, and the reference set. The reference set is a set of genomic coordinates

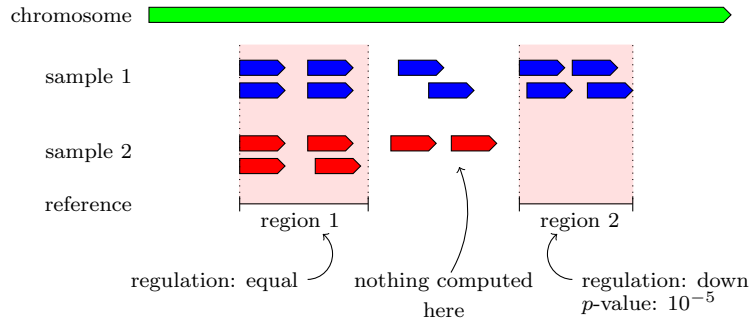


Figure 6: Differential expression computed on two reference intervals.

and, for each interval, it will count the number of feature on each sample and compute the differential expression. For each reference interval, it will output the direction of the regulation (up or down, with respect to the first input set), and a  $p$ -value from a Fisher exact test (see figure 6).

This reference set seems boring. Why not computing the differential expression without this set? The answer is: the differential expression of what? I cannot guess it. Actually, you might want to compare the expression of genes, of small RNAs, of transposable elements, of anything... So the reference set can be a list of genes, and in this case, you can compute the differential expression of genes. But you can also compute many other things.

Suppose that you clusterize the data of your two input samples (you can do it with the `clusterize` and the `mergeTranscriptLists` tools). You now have a list of all the regions which are transcribed in at least one of the input samples. This can be your reference set. This reference set is interesting since you can detect the differential expression of data which is outside any annotation.

Suppose now that you clusterize using a sliding window the two input samples (you can do it with the `clusterizeBySlidingWindows` and the `mergeSlidingWindowsClusters` tools). You can now select all the regions of a given size which contain at least one read in one of the two input samples (do it with `selectByTag` and the tag `nbElements`). Again, this can be an other interesting reference set.

In most cases, the sizes of the two input samples will be different, so you should probably normalize the data, which is an available option. The —rather crude— normalization increases the number of data in the least populated sample and decreases the number of data in the most populated sample to the average number of data.

You can also plot the differential expression. A point  $(x, y)$  refers to a reference interval which contains  $x$  data in the first sample and  $y$  data in the second sample. If you normalized the data, then the plot reports the normalized figures.

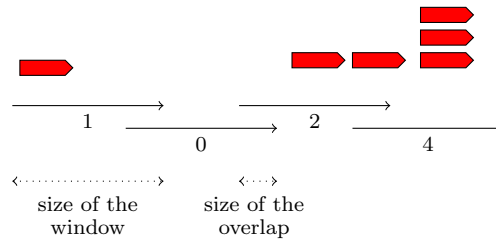


Figure 7: Sliding windows: counting the number of reads.

#### 4.1.2 Sliding Windows

Sliding windows are a convenient ways to clusterize data mapped on the genome. There are two important parameters of a sliding window: the size of the window and the size of the overlap. In Figure 7, a sliding window counts the number of reads.

**clusterizeBySlidingWindows.py** By default, sliding windows count the number of reads. However, you can basically merge any information which is contained in the tags (look at Section 6.2 if you want to know more about tags). You can compute the average, sum, median, max or min of the tags for each window. For instance, every window can contain the average cluster size, if you merge clusters instead of reads.

The output file is a GFF3 file, where each element is a window. There is a special tag for each window, whose name is `nbElements` if you counted the number of transcripts per sliding window. However, if you performed a “min” (resp. “max”, “sum”, “median”, “average”) operation on the tags `value` of the transcripts, then the tag of the window will be `minValue` (resp. `maxValue`, `sumValue`, `medValue`, `avgValue`).

You also have different option, which can select the  $n\%$  highest regions, or the regions with at least  $n$  features in it, or even the regions with at least  $n$  unique features. This last option is useful when you want to cluster the reads which have mapped only once, for instance.

**mergeSlidingWindowsClusters.py** Sliding windows are also useful to compare two (or more!) sets of data. This can be very valuable when you want to compare differential expression in two different conditions. When you have two different sliding windows sets, this function merges them into one, where each window contains the two pieces of information. You may want to plot the data afterwards using the `plot.py` function.

#### 4.1.3 Merging data

**mergeTranscriptLists.py** The script is similar to `compareOverlapping.py`, except that when data of two different sets overlap, they are merged. You can

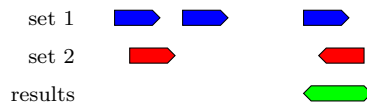


Figure 8: Finding transcription on both strands using `mergeTranscriptLists.py`.

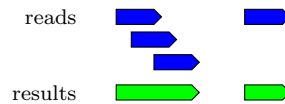


Figure 9: Clustering using `clusterize.py`.

use the same parameters as `compareOverlapping.py` and use them to look for transcription on both strands, for example (see Figure 8).

This script can also be used with one input data set. In this case, its behaviour is similar to `clusterize.py`.

**clusterize.py** The script clusterizes the reads. Two reads are clusterized when their genomic intervals overlap (see Figure 9). The output is a GFF3 file, where each element is a cluster. The number of elements in the cluster is given by the tag `nbElements`.

Alternatively, some options may clusterize the features which are closer than a given threshold.

By default, the tool clusterizes all features which overlap (or nearly overlap), even if they are on different strands. If you want to clusterize the features which are on the same strand only, you can specify it.

**findTss.py** This script is specially useful when you have 5' capped reads, that is to say, when the reads that you have mark the beginning of the transcripts. This script find all the TSS that are found by your data (see Figure 10).

In some —most, actually— cases, there is no clear TSS, but a stretch of possible TSSs. So you can choose the maximal distance between two reads for them to mark the same transcription start (for example, two reads that are distant by 20 nt. can mark 1 or 2 TSS, depending on the value of a parameter).

You can plot the distribution of the number of reads per TSS: a point  $(x, y)$  tells you that  $y$  transcripts starts are marked by  $x$  reads. The plot has sometimes

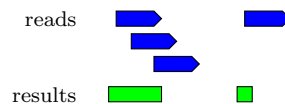


Figure 10: Finding TSS using `findTss.py`.

a long tail towards the high values in the  $x$ -axis, so you can zoom to plot only the first points on this axis by using some parameters.

#### 4.1.4 Select part of the data

This set of scripts reads of list of genomic coordinates or sequences and select those with some given properties.

`restrictSequenceList.py` Read a list of sequences and a list of sequence names, and select those sequences of which the name is in the sequence name.

`restrictFromNucleotides.py` In a list of sequences, select those with no ambiguous nucleotide, i.e. N, but also R, Y, B, etc..

`restrictGenomicCoordinates.py` In a list of genomic coordinates, keep those which are on a chromosome and / or between two positions.

`restrictFromSize.py` In a list of sequences or genomic coordinates, select those which are longer and / or shorter than a give size.

## 4.2 Working on Sequences

These tools will help you selecting and organizing your sequences, given by a FASTA or FASTQ file (usually, your reads).

`getReadDistribution.py` This function analyzes your reads, count the number of occurrences and plots this distribution. You can also select the  $n$  most sequenced reads (or the  $x\%$  highest).

`trimAdaptator.py` This function removes the adaptor from the 3' end of your reads. It can recognize partial adaptators (as it is the case when short RNAs are sequenced).

`getSequence.py` Get a sequence from you FASTA or FASTQ file, given the name of the sequence.

## 4.3 Modify your data

These tools do the “dirty job” that is sometimes useful to do: shrink or extend some genomic coordinates, take the 20 first nucleotides of your reads, etc.

`modifySequenceList.py` It keeps the  $n$  first or last nucleotides from a FASTA or FASTQ file.

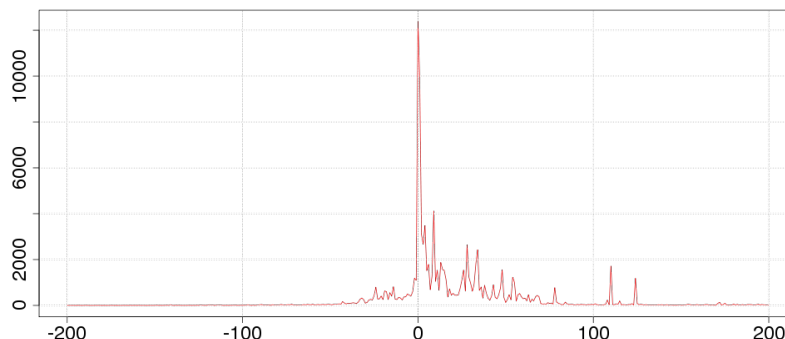


Figure 11: Distance between some reads and RefSeq genes in *Drosophila melanogaster*, using `getDistance.py`.

`modifyGenomicCoordinates.py` It extends or shrinks a set of genomic coordinates, much like what is done in Figure 5.

`changeTagName.py` It changes the name of a tag in a transcript list (see Section 6.2 to know more about tags). This may be useful to change the name of a tag which have been automatically addressed, like `nbElements` while clustering, to a more precise name (such as `nbReadsInRoot` for instance).

## 4.4 Visualizing you data

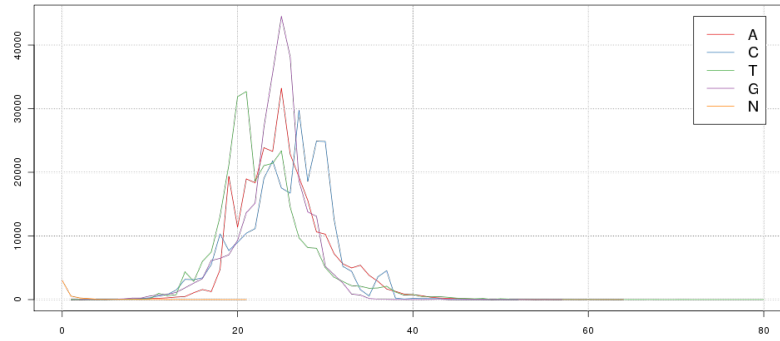
This set of tools do not modify your data, but simply output graphs or Excel files.

`getDistance.py` Give the minimum distance between every data from the first input set and the data from the second input set. Outputs the size profile (see Figure 11).

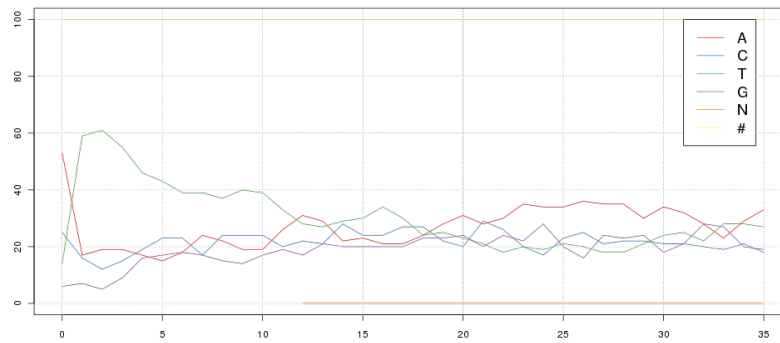
You can also shrink a extend you set of genomic coordinates, as in `compareOverlapping.py`, so that you can get the distance between starts of reads and starts or genes, for instance.

`getLetterDistribution.py` Get the nucleotide distribution the input sequence list. Outputs nucleotide distribution profile of the data (see Figure 12(a)), and the average nucleotide distribution for each position of the read (useful to detect a bias in the first nucleotides, for example, see Figure 12(b)).

`getNb.py` Get the number of occurrences of a mapping, or the number of exons of each mapping, see Figure 13.



(a) Nucleotide profile for the whole distribution. The  $x$ -axis is the proportion each nucleotide and the  $y$ -axis is the number of reads with the corresponding distribution. There is, for instance, more that 40,000 reads with around 25% of G.



(b) Nucleotide by nucleotide. The  $x$ -axis is the index of the nucleotides (start at 0),  $y$ -axis is the percentage of nucleotides. The # curve give the percentage of reads with at least the given size. In this example, all reads have exactly 36 nucleotides.

Figure 12: Nucleotide distributions using `getLetterDistribution.py`.



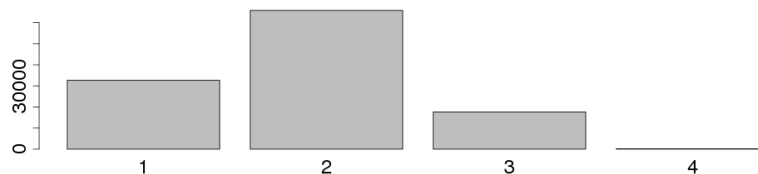


Figure 13: Number of exons per transcript of some 454 mapped reads, using `getNb.py`.

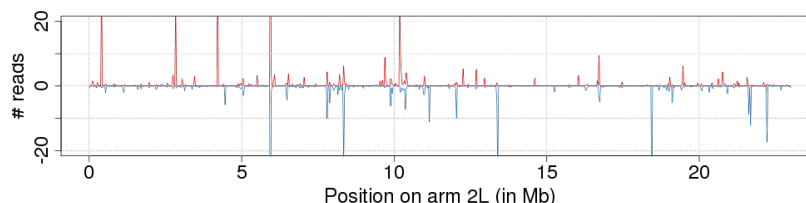


Figure 14: Density profile of some reads, using `getRepartition.py`.  $x$ -axis is the genomic coordinates of the chromosome 2L of *Drosophila melanogaster*, release 5.  $y$ -axis is the number of reads.

`getRepartition.py` Print a density profile of the data for each chromosome, see Figure 14.

`getSizes.py` Get the read size distribution, see Figure 15.

`plotTranscriptList.py` Plot the data attached as tags in a transcript list. See Section 6.2 if you want to know more about tags. This is especially interesting for displaying the comparison of different sets of sliding windows (see Figure 16).

You can display curves, barplots, points clouds, heat points clouds (where the color of the points depend on the third variable) for the distribution your tags. You can also use the log scale.

`plotRepartition.py` Plot the data attached as tags in a transcript list along the genome. See Section 6.2 if you want to know more about tags. This is especially interesting for displaying the regions where different sets of sliding windows differ (see Figure 17).

## 4.5 Conversion tools

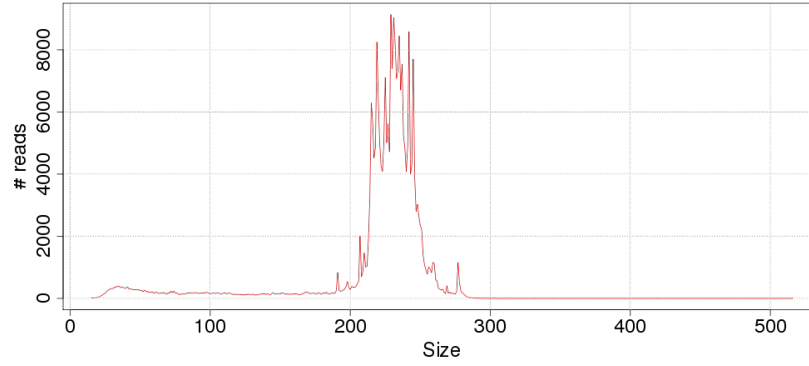
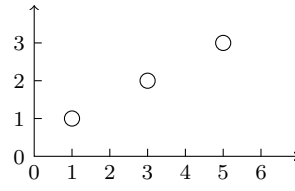


Figure 15: Size distribution of some 454 reads, using `getSizes.py`.  $x$ -axis is the size,  $y$ -axis is the number of reads.

```
chr1 S-MART transcript 100 200 . + . ID=region1;nbReadsRoot=1;nbReadsLeaf=1
chr1 S-MART transcript 200 300 . + . ID=region2;nbReadsRoot=5;nbReadsLeaf=3
chr1 S-MART transcript 300 400 . + . ID=region3;nbReadsRoot=3;nbReadsLeaf=2
```

(a) A short GFF file, with the results of sliding windows counting the number of reads in two conditions (root and leaf).



(b) The corresponding plot.

Figure 16: A plot using `plotTranscriptList.py`.

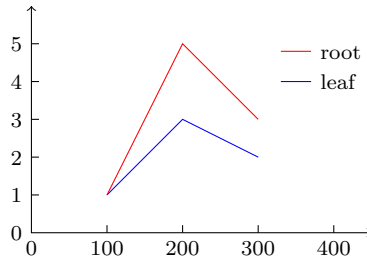


Figure 17: A plot using the data in 16(a) and `plotRepartition.py`.

**convertTranscriptFile.py** Convert some data from one format to another format. The input format has been previously listed. They can be transcript format (BED, GFF) or mapping format (SAM, Blast `-m 8`). The output formats are:

- BED format
- GFF2 format
- GFF3 format
- the format used by GBrowse
- the format used by UCSC genome browser
- the intern mySQL format

**coordinatesToSequence.py** Provide a list of genomic coordinates, a reference genome, and it will output the sequence corresponding to the genomic coordinates, in a multi-FASTA file.

**qualToFastQ.py** Convert a QUAL files (together with the corresponding FASTA file) to a FASTQ file, that S-MART can use.

## 4.6 Other tasks

**cleanGff.py** This tool tries to “clean” a GFF by removing all the irrelevant lines and possibly altering others. This is specially useful when you download a genome wide annotation file from NCBI, for instance. Look at Appendix C to see the kind of problems it can solve.

**getRandomRegions.py** Generates a set of random regions in a reference genome. Useful to compare your data with random data.

**removeAllTmpTables.py** The toolbox may generate some mySQL that you may want to drop. Use this script to do so. Notice that if you interrupt the execution of a program while it is running, the database it uses cannot be removed and you have to do it by yourself.

## 5 Possible pipe-lines

I will describe here a couple of “pipe-lines” that I have found useful.

## 5.1 Use two mappers

Up to now, there is no clear consensus about which mapper gives best results (there is probably none, though). So I used two of them. In this example, I mapped with Blat and Exonerate.

```
python mapperAnalyzer.py -i blatOutput.psl -f psl
-q reference.fasta -o mappingWithBlat
python mapperAnalyzer.py -i exonerateOutput.exo -f exo
-q reference.fasta -a mappingWithBlat -r
-o mappingWithBlatAndExo
```

## 5.2 Find piRNA clusters

Some papers show an interesting way to find potential clusters of piRNAs. They:

- sequence the transcriptome with RNA-Seq,
- map the data,
- keep the data which sizes between 25 and 31 nucleotides,
- exclude everything which overlaps with RefSeq data,
- merge the mappings into clusters (20 kb),
- keep those clusters which have at least ten elements.

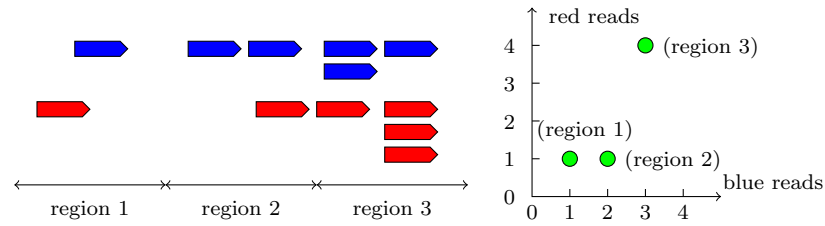
Here is how it can be done.

```
python restrictFromSize.py -i mappedData.gff3 -f gff
-m 25 -M 31 -o goodSize
python compareOverlapping.py -i goodSize.gff3 -f gff
-j refSeqGenes.bed -g bed -c -x -o noGene
python clusterize.py -i noGene.gff3 -f gff
-d 20000 -o clustered
python selectByTag.py -i clustered.gff3 -f gff
-g nbElements -m 10 -o bigClusters
```

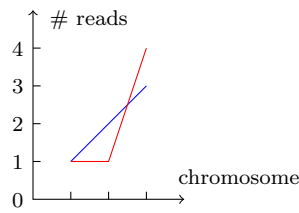
## 5.3 Get the letter distribution of the beginning of the data

It might be interesting to see if you have a bias in the distribution of the nucleotides at the beginning of the data that you have selected. For example, piRNA usually start with U, RefSeq data, with ATT, and so on. To visualize it, you can:

- select the first 30 nucleotides of your data,



(a) Example of two sets of reads (one blue, one red), clustered by sliding windows into three regions. (b) The corresponding plot. For each dot, the corresponding region is given in parenthesis.



(c) The compared distribution on the genome.

Figure 18: Comparing two sets of reads with sliding windows.

- get back the sequences from your genomic coordinates,
- get the nucleotide distribution.

It can be done this way.

```
python modifyGenomicCoordinates.py -i data.gff3 -f gff -s 30
-o cutData
python coordinatesToSequence.py -i cutData.gff3 -f gff
-s reference.fasta -o sequences.fasta
python getLetterDistribution.py -i sequences.fasta
-o nucleotideDistribution
```

## 5.4 Compare two sets of reads with sliding windows

Suppose you have two sets of reads (for instance, from two conditions) and you want to have a broad picture of a possible correlation between the two sets. What you can do is to use a sliding window to cluster the two sets of reads (see Figure 18(a)). Then, you can plot the results, where each point  $(x, y)$  comes from a sliding window which contains  $x$  reads from the first set and  $y$  reads from the second set (see Figure 18(b)).

Here are the steps you have to perform to do so:

- Use a sliding window to cluster each set (here the size of the window is 1kb, and the overlap is 500 nt.). You may or you may not consider that two reads on opposite strands can be in the same window. Use option `-2` to consider the strands independantly. You now have two GFF3 files, and the tag `nbElements` counts the number of reads for each window.
- Change the name of the tags `nbElements` to two different names (here: `set1` and `set2`). You will have some problems otherwise in the next step.
- Merge the two files.
- Plot the distribution. Here, we will use log bases for  $x$  and  $y$ -axes (`-l` option). We also have to specify the default values for missing data with the `-X` and `-Y` options. We can specify a label for the axes with the `-m` and `-n` options.

Supposing that the reads are in the files `set1.gff3` and `set2.gff3`, you can use the following commands to perform the steps:

```
python clusterizeBySlidingWindows.py -i set1.gff3 -f gff
-s 1000 -e 500 -o set1Clustered -2
python clusterizeBySlidingWindows.py -i set2.gff3 -f gff
-s 1000 -e 500 -o set2Clustered -2
python changeTagName.py -i set1Clustered -f gff -t nbElements
-n set1 -o set1ClusteredGoodTag
python changeTagName.py -i set2Clustered -f gff -t nbElements
-n set2 -o set2ClusteredGoodTag
python mergeSlidingWindowsClusters.py -i set1ClusteredGoodTag.gff3
-f gff -j set2ClusteredGoodTag.gff3 -g gff -o set12Clustered
python plotTranscriptList.py -i set12Clustered.gff3 -f gff -x set1
-y set2 -X 0 -Y 0 -l xy -n set1 -m set2 -s points
-o set12ClusteredPlot
```

And, yes, you get the Spearman rho!

Moreover, if you want to see the distribution of the two sets of reads on the chromosome (see Figure 18(c)), you can do it with:

```
python plotRepartition.py -i set12Clustered.gff3 -n set1,set2
-c blue,red -r -o set12ClusteredGenome
```

## 5.5 Compare RNA-Seq with tiling arrays using sliding windows

Let us suppose that you have some tiling array you want to compare with your sequencing. Let us also suppose you have analyzed your array data and you now have some  $p$ -value,  $t$ -value or anything which is a proxy to gene intensity.

You may want to use sliding windows to compare with your reads. Of course, for a fair comparison, you want to compare the regions where you have at least 1 chip (you may have not covered the whole genome).

The steps are the following:

- Cluster the array data into sliding windows, with a window size of 1kb, and an overlap of 500 nt., just to count the number of chips per window.
- Cluster the array data into sliding windows by using the average intensity value in the window (but we can use min, max or median value). Here we will suppose you have a GFF3 file with a `intensity` tag.
- Merge the two previous files.
- Keep the regions where you have at least 1 chip.
- Cluster the reads as done in section 5.4.
- Merge the array file and the reads file.
- Plot the distribution with log base on *y*-axis (for the reads).

```
python clusterizeBySlidingWindows.py -i array.gff3 -f gff
-s 1000 -e 500 -o arrayNb -2
python clusterizeBySlidingWindows.py -i array.gff3 -f gff
-s 1000 -e 500 -o arrayIntensity -g intensity -r avg -2
python mergeSlidingWindowsClusters.py -i arrayNb.gff3
-f gff -j arrayIntensity.gff3 -g gff -o arrayNbIntensity
python selectByTag.py -i arrayNbIntensity.gff3 -f gff
-g nbElements -m 1 -o arrayNbIntensity1chip
python clusterizeBySlidingWindows.py -i reads.gff3 -f gff
-s 1000 -e 500 -o readsClustered -2
python changeTagName.py -i readsClustered -f gff -t nbElements
-n nbReads -o readsClusteredGoodTag
python mergeSlidingWindowsClusters.py
-i arrayNbIntensity1chip.gff3 -f gff
-j readsClusteredGoodTag.gff3 -g gff -o arrayReadsClustered
python plotTranscriptList.py -i arrayReadsClustered.gff3 -f gff
-x avgIntensity -y nbReads -Y 0 -l y -n array -m reads
-s points -o arrayReadsPlot
```

## 5.6 Compute differential expression

Suppose you have two sets of reads, from two different conditions and you want to compare them. The main difficulty here is to decide where you are going to compare. You could decide to slice the genome into sliding windows and perform the comparison on the sliding windows, but that is probably not the

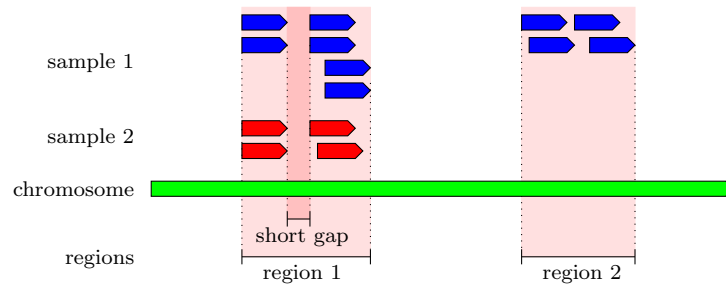


Figure 19: Finding the regions where to compare the two clusters.

best thing to do (although you can do it with S-MART). What you should do is to find the clusters of reads and compare the clusters.

For instance, look at figure 19, which represents two sets of reads mapped on a genome. Obviously, we have two clusters, although the red reads are not present in the second cluster (which is a particularly good case of differential expression). Basically, we have to find these two regions from the two sets of data. Notice that in the first region, there is a small gap which is not covered by any read. So, we have to accept some small gaps to merge the reads into clusters.

Then, it could be nice to keep the regions where there is a clear differential expression and see where they are on the genome.

From the previous reasoning, we can conclude that we have to follow the following steps:

- Cluster the reads from sample 1 with some gap allowed (here, 10 nt.). Only cluster the reads which are on the same strand (`-c` option).
- Do the same for the sample 2.
- Merge the two sets of clusters. Again, the clusters on different strands should not be merged.
- Compute the differential expression from sample 1 and sample 2 using the regions previously found. We now have a set of regions with a  $p$ -value associated to differential expression.
- Select the regions with a very low  $p$ -value (here,  $10^{-100}$ ).
- Plot these regions onto the genome (you will need the reference genome `genome.fasta` for that).

Expressed in the S-MART language, this becomes:

```
python clusterize.py -i sample1.gff3 -f gff -c -d 10
-o sample1clustered
python clusterize.py -i sample2.gff3 -f gff -c -d 10
```



```
chr1 S-MART transcript 100 400 . + . ID=read1;nbMismatches=2
chr1 S-MART exon 100 200 . + . ID=read1-1;Parent=read1
chr1 S-MART exon 300 400 . + . ID=read1-2;Parent=read1
```

Figure 20: A short GFF3 file, containing a transcript with two exons. The tags are in the last field of each line.

```
-o sample2clustered
python mergeTranscriptLists.py -i sample1clustered.gff3
  -f gff -j sample2clustered -k -c -o sample12clustered
python getDifferentialExpression.py -i sample1.gff3 -f gff
  -j sample2.gff3 -g gff -k sample12clustered -l gff
  -o sample12differential
python selectByTag.py -i sample12differential -f gff
  -g pValue -M 10e-100 -o sample12differentialLowPValue
python getRepartition.py -i sample12differentialLowPValue
  -f gff -r genome.fasta -2 -o sample12differentialLowPValuePlot
```

## 6 More about S-MART

I will give you here some details about the internal representations of the data in S-MART. Just in case you would like to know how it works...

### 6.1 Data structures

S-MART mainly use a data structure that models a transcript. A transcript can be decomposed as a set of exons, which basically are genomic intervals. So, a transcript is a set of genomic intervals. So is a mapped read; short reads usually have only one exon, but longer ones often have several.

A cluster is modeled the same way. While clusterizing, S-MART merges the exons of the transcripts one by one (if they overlap), thus forming a new set of new exons.

### 6.2 Tags

S-MART for each “transcript”, S-MART attaches some information, called *tags*. The information might be the number of mismatches of a mapped read, or the number of elements in a cluster.

S-MART automatically load the tags from a GFF3 file (see Figure 20 for an example of a GFF3 file with tags). It updates the tags while mapping the reads and clusterizing them.

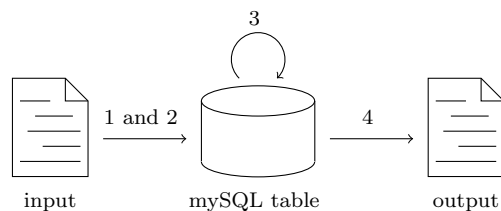


Figure 21: Use of mySQL tables in S-MART, in 4 steps.

### 6.3 How mySQL is used?

S-MART stores all the data in a mySQL table. That is to say, S-MART (see Figure 21):

1. reads your data,
2. stores them in a database,
3. processes them (it basically uses an B-tree on nested bins to compute the overlaps faster) and
4. outputs them into an output file.

So, in some cases, you can skip the steps 1 and 2. Suppose that you use the tool `compareOverlapping.py` to keep only those reads which overlap with RefSeq data, then `clusterize.py` to merge them into clusters. In `compareOverlapping.py`, you can use the `-y` option to keep into the database the data. Then, if you mention that the input format of `clusterize.py` is `-f sql`, S-MART directly uses the data which is available in your mySQL database. So, you can skip the steps 1 and 2, and save some time! Be aware, however, that it takes some place on your hard disk.

### 6.4 Contribute to S-MART!

If you want to add your own tool to S-MART, please do not hesitate. You can develop it using the API I have made. Look at the format I have used to generate the help from my Python scripts (using `OptionParser`) in the `Python` directory and add your own Python file in the same directory. That is it! It will be automatically included in the GUI and you will be able to start your own tool from there.

## 7 Contact

For any comment, suggestion, remark, do not hesitate to contact me.

## A Load data on your genome browser

**UCSC genome browser** Go to the UCSC genome browser. Select your species (if available). Click on the **add custom track** button. Upload your data in a “.ucsc” or “.bed” file.

**GBrowse** Go to your favorite GBrowse, for example the fly GBrowse. Go down, to the **Add your own tracks** section. Select your data in a “.gb” file after clicking on the **Browse...** button. Now watch amazed your reads.

## B Get other data

**UCSC** As long as your organism is covered by UCSC (including Mammals, other Vertebrates, Insects, Nematodes and a few other species), this may be the right place to get your data. Go to their database to retrieve the annotation. Select the right organism, select the group, track and table where your data are (it may take some time to find the right combination), select BED as the output format, and mention an output file. This is it!

**RefSeq** NCBI provides a list of RefSeq sequences on their public FTP. You can browse it with your normal browser. Choose the right organism, then you should download the data chromosome by chromosome (which is pain, I reckon). Be sure that you get the GFF files. Afterwards, you should concatenate the files.

**FlyBase** If you work on an insect, you can go to the FlyBase Web site to download your annotation. Select your organism, then the release (latest is best), then the GFF folder. You can then download the whole annotation in a file whose name contains the word **all**. Be careful, you will download the whole annotation of your insect! It may not exactly be what you need!

**Other data** Many Web site are dedicated to a specific organism or some specific data (miRNAs, transcription factor binding sites, etc.). I cannot cover them all, but you can probably find there the data that you need. It is a matter of patience. Just make sure that the data is available in the usual formats (GFF or BED, for instance).

## C Troubleshooting

**My GFF file is not parsed!** A problem with a GFF file may have different causes.

First, it seems that there is no gold standard for GFF files. Current files are usually produced in GFF3 format, which includes an **ID** and possibly a **Parent** fields. However, this convention is sometimes not followed. If they are

not present, it is hard for S-MART to link a line which contains a transcript annotation to the lines which contains the annotations of its exons.

Second, please make sure that your files contains *only* the data you are interested in. For instance, when you download some genome wide annotation, you have information of the size of the chromosomes, annotation of the start and stop codon, of the 5' and 3' UTR. This is usually something you are not interested in and S-MART can be confused. Basically, you just want the annotation of the transcripts and their exons (see Figure 20 for an exemple of a very simple GFF3 file).

Third, different sources may name the chromosomes differently. Compare `chr_I` with `Chr1` (or even worse: `gi|157069709|gb|AABX02000103.1|`). S-MART cannot see they actually are the same chromosome (actually, UN-S-MART could have been a better name). So you have to change the names of the chromosomes accordingly.

In all cases, you can use the tool “Clean GFF” (see Section 4.6), which tries to produce a GFF3 file which can be understood by S-MART. But please check the output file to make sure the output is correct!

## Index

changeTagName.py, 12  
cleanGff.py, 17  
clusterize.py, 11  
clusterizeBySlidingWindows.py, 10  
compareOverlapping.py, 9  
convertTranscriptFile.py, 16  
coordinatesToSequence.py, 17  
  
findTss.py, 11  
  
getDistance.py, 13  
getLetterDistribution.py, 13  
getNb.py, 13  
getRandomRegions.py, 17  
getReadDistribution.py, 12  
getRepartition.py, 13  
getSequence.py, 12  
getDistance.py, 13  
  
mapperAnalyzer.py, 8  
mappingToCoordinates.py, 9  
mergeSlidingWindowsClusters.py, 10  
mergeTranscriptLists.py, 10  
modifyGenomicCoordinates.py, 12  
modifySequenceList.py, 12  
  
plotRepartition.py, 16  
plotTranscriptList.py, 13  
  
qualToFastQ.py, 17  
  
removeAllTmpTables.py, 17  
restrictFromNucleotides.py, 12  
restrictFromSize.py, 12  
restrictGenomicCoordinates.py, 12  
restrictSequenceList.py, 11  
  
trimAdaptator.py, 12