
Counter RNAseq Window Documentation

Release 1.0.0

Bertrand Néron

Apr 05, 2019

CONTENTS

1	User Guide	1
1.1	overview	1
1.2	Installation	2
1.3	Quick start	5
1.4	Inputs / Outputs	12
2	Developer Guide	17
2.1	Overview	17
2.2	reference API	17
2.3	Release notes	37
3	Indices and tables	39
	Python Module Index	41
	Index	43

1.1 overview

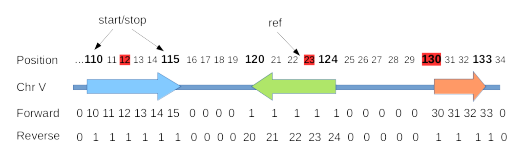
Counter RNA seq Window is a package which aim to compute and visualize the coverage of RNA seq experiment.

The *craw* package contains two scripts *craw_coverage* and *craw_htmp*. *craw_coverage* compute the coverage, whereas *craw_htmp* allow to represent graphically the results of *craw_coverage* with a heat map.

1.1.1 *craw_coverage*

craw_coverage take as input a bam file or wig file and an annotation file. **The annotation file** describe on which gene the *craw_coverage* must compute the coverage. The script compute a coverage for each position of this gene on a specified window around a position of reference on both sense and put the results on a matrix. The region of interest can be fixed for all genes (specified by the command line) or variable. In the this case the annotation file must contains two columns to specify beginning and the end of the region to take in account. The results in the matrix are centered on the position of reference of each gene. In the case of variable length of window the results are padded on left and right if necessary with *None* value. The results is saved in a file as a tabulated separated value by default with the same name as the bam file with the *.cov* extension (see [Outputs](#) for more details).

Below an example to illustrate how *craw_coverage* work. If we consider the following genome and we want to analyze 3 gene *foo*, *bar*, *buz*



On the figure above

- The first line represent the positions on the genome (1-based)
 - The bold position indicate the boundaries of region we want to analyse.
 - the red highlighted positions indicate, for each region, the position of reference.
- the second line represent the genes and their respective sense.
- the 2 last lines the coverage at each position of the genome for each strand

So to analyse these genes, we create an annotation file like following.

gene	Chr	strand	start	stop	ref
foo	V	+	110	115	112
bar	V	+	130	133	130
buz	V	-	120	124	123

the run a command line like:

```
craw_coverage --bam mygenome.bam --annot my_annotation --ref-col ref --start-col _  
↪start --stop-col stop
```

will produce the following coverage matrix

Gene	sense	-2	-1	0	1	2	3
foo	S	10	11	12	13	14	15
foo	AS	1	1	1	1	1	1
bar	S	None	None	30	31	32	33
bar	AS	None	None	1	1	1	1
buz	AS	None	1	1	1	1	1
buz	S	None	24	23	22	21	20

1.1.2 `craw_http`

`craw_http` read coverage file produced by `craw_coverage` and generate a graphical representation. It can produce either a file or an interactive graphic. The look and feel of the graphic and the format of supported outputs vary in function of the backend of matplotlib used (see [matplotlib configuration](#)). It can also produce raw images using pillow where 1 nucleotide is represent by 1 pixel.

1.1.3 Licensing

All files belonging to the Counter RNAseqWindow (`craw`) package. are distributed under the GPLv3 licensing.

You should have received a copy of the GNU General Public License along with the package (see COPYING file). If not, see <http://www.gnu.org/licenses/>.

`craw` is free software: you can redistribute it and/or modify it under the terms of the GNU General Public License as published by the Free Software Foundation, either version 3 of the License, or (at your option) any later version.

`craw` is distributed in the hope that it will be useful, but WITHOUT ANY WARRANTY; without even the implied warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the GNU General Public License for more details.

Authors: Bertrand Neron Copyright © 2017-2019 Institut Pasteur (Paris). see COPYRIGHT file for details.

1.2 Installation

There are 3 ways to use `craw`: * by install the standalone python scripts * by using docker image * by using singularity image

1.2.1 CRAW installation

Requirements

For `craw_coverage`

- python >= 3.5
- pysam >= 0.15.2

For `craw_http`

- python \geq 3.5
- pysam == 0.15.2
- pandas \geq 0.24
- numpy \geq 1.16
- matplotlib \geq 3.0
- pillow \geq 5.4
- scipy \geq 0.16.1

Installation

Installation from package

Using pip

```
pip install craw
```

Do not forget to configure the *matplotlib* backend, specially if you use virtualenv. Otherwise on some platform there won't any output. See [matplotlib configuration](#) for more explanation.

Note: On MacOS install python > 3 from image on <http://python.org> . Then install craw using pip

```
pip3 install craw
```

`craw` will be installed in `/Library/Framework/Python.Framework/Version/3.6/` So if you want to use directly `craw_coverage` and `craw_http` just create a symbolic link like this:

```
ln -s /Library/Framework/Python.Framework/Version/3.6/bin/craw_coverage /usr/local/  
↪bin/craw_coverage  
ln -s /Library/Framework/Python.Framework/Version/3.6/bin/craw_http /usr/local/bin/  
↪craw_http
```

The documentation (html and pdf) is located in `/Library/Framework/Python.Framework/Version/3.6/share/craw/`

Installation from repository

To get the last version, clone the project and install with the `setup.py`

```
git clone https://gitlab.pasteur.fr/bneron/craw.git  
cd craw  
pip install .
```

Note: Instead of installing `craw` you can directly use the scripts from the repository. You can also use the package without installing it. To do this, you have to export the **CRAW_HOME** environment variable. *CRAW_HOME* must point to the *src* directory of the project. Then you can use *craw_coverage* and *craw_http* scripts located in *bin* directory.

This project is documented using [sphinx](#). So if you use a clone, you have to generate the documentation from the source.

The project come from with some unit and functional tests. to test if everything work fine.

```
cd $CRAW_HOME python3 tests/run_tests.py -vvv
```

matplotlib configuration

matplotlib is a python library to create graphics. *craw_hmp* use this library to generate heat map. The two parameters to configure for *craw* is:

- the backend
- figure.dpi

backend

matplotlib lay on low level graphic library of your computer. This library will determine the graphical formats manage by *matplotlib* and then by *craw_hmp*. Most of backend handle *png* but some library like *Qt* can handle 'jpeg', 'eps', *pdf* ...

In your *matplotlibrc* file you must define the backend for instance to use *Qt5*

```
backend: qt5agg
```

An example of *matplotlibrc* file and all supported backend is available here: <http://matplotlib.org/users/customizing.html#a-sample-matplotlibrc-file>

figure.dpi

It's not an essential option but *matplotlib* and *craw_hmp* will produce better graphic (on screen) if you configure *matplotlib* to the native resolution of your screen. To know the resolution of your screen you can visit the following page <https://www.infobyip.com/detectmonitordpi.php> and report the resolution (for 1 inch) in *matplotlibrc* file like:

```
figure.dpi: 96
```

For full explanation on how to configure *matplotlib* read <http://matplotlib.org/users/customizing.html#the-matplotlibrc-file>.

1.2.2 Using Docker Image

Docker images are available. The two scripts are accessible through the sub-command *coverage* or *hmp*. For instance to use the latest version of *craw_hmp*:

```
docker pull c3bi/craw
docker run -v$PWD:/root -it c3bi/craw coverage --bam foo.bam --annot foo.annot --ref-
→col 'Position' --before 3 --after 5 --out foo.cov
docker run -v$PWD:/root -it c3bi/craw hmp --size raw --out foo.png foo.cov
```

Note: In docker the interactive *hmp* output is not available. So you must specify the *-out* option

1.2.3 Using Singularity Image

Singularity images are available. The two scripts are accessible through the sub-command *coverage* or *htmp*. For instance to use the latest version of *craw_htmp*:

```
singularity pull --name craw shub://c3bi/craw
./craw coverage --bam foo.bam --annot foo.annot --ref-col 'Position' --before 3 --
↪after 5 --out foo.cov
./craw htmp --size raw --out foo.png foo.cov
```

Note: instead of Docker images, in Singularity images the interactive output is available.

1.3 Quick start

1.3.1 *craw_coverage*

craw_coverage need a file bam or wig to compute coverage and an annotation file to specify on which regions to compute these coverages.

- the *-b* or *--bam* allow to specify the path to the bam file.
- or alternatively the
 - *-w*, *--wig* option to specify the path to the wig file if the both strand are encode in same file (negative value are on reverse strand)
 - *--wig-for* and *--wig-rev* to specify the paths to the wig files for the forward and reverse strand respectively
- the *-a* *--annot* allow to specify the path to the annotation file.

The *--bam* and *--wig* options are mutually exclusive but one of these option is required. *--wig* and *--wig-for* or *--wig-rev* are also mutually exclusive. the *--annot* option is required.

```
craw_coverage --bam ../WTE1.bam --annot ../annotations.txt --ref-col Position --
↪before 100 --after 500
craw_coverage --wig ../WTE1.wig --annot ../annotations.txt --ref-col Position --
↪before 100 --after 500
craw_coverage --wig-for ../WTE1_forward_strand.wig --wig-rev ../WTE1_reverse_strand.
↪wig --annot ../annotations.txt \
--ref-col Position --before 100 --after 500
```

Warning: At the same place of *bam* file, there must be the corresponding index file (the *bam.bai* file). To generate the *.bai* file you have to use *samtools* program:

```
samtools index file.bam
```

see <http://www.htslib.org/doc/> for more explanation.

with fix window

To compute the coverage on a fix window: we need to specify which column name in the annotation file define the reference position. The window will computed using this reference position.

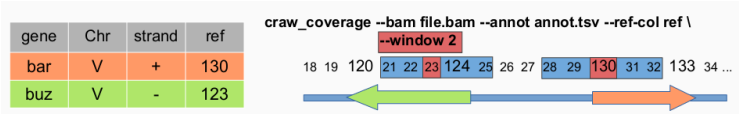
- `--ref-col`

Note: if `--ref-col` is omitted `craw_coverage` will use the column position. If there not “position” column an error will occur.

two ways to determine the window:

with `--window` option for a window centered on the reference position.

- `--window` define the number of nucleotide to take in account before and after the reference position.

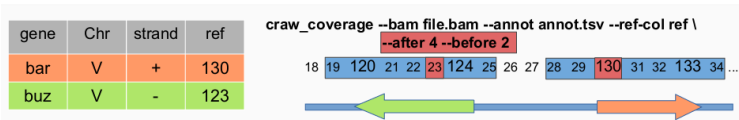


```
craw_coverage --bam ../WTE1.bam --annot ../annotations.txt --ref-col Position --  
--window 100
```

This command will compute coverage using `WTE1.bam` and with `annotations.txt` file the column used to compute the window is ‘Position’ and the window length will be 100 nucleotide before the reference position and 100 nucleotides after (201 nucleotides length).

With an non centered window we have to specify two options `--before` and `--after`

- `--before BEFORE` define the number of nucleotide to take in account before the reference position.
- `--after AFTER` define the number of nucleotide to take in account after the reference position.



```
craw_coverage --bam ../WTE1.bam --annot ../annotations.txt --ref-col Position --  
--before 100 --after 500
```

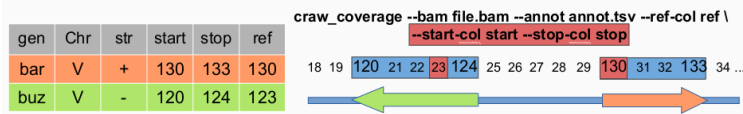
This command will compute coverage using `WTE1.bam` and with `annotations.txt` file the column used to compute the window is ‘Position’ and the window length will be 100 nucleotide before the reference position and 500 nucleotides after (201 nucleotides length).

Note: `--after` and `--before` options must be set together and are incompatible with `--window` option.

with variable window

The regions must be specified in the annotation file.

- `--start-col COL` define the name of the column in annotation file which define the start position of the region to compute.
- `--stop-col COL` define the name of the column in annotation file which define the stop position of the region to compute.



```
crawl_coverage --bam ../WTE1.bam --annot ../annotations.txt --ref-col annotation_start_
↪--start-col annotation_start --stop-col annotation_end
```

This command will compute coverage using WTE1.bam and with annotations.txt file.

- The reference position will define by the *annotation_start* column
- The first nucleotide of the window will be define by *annotation_start* column.
- The last nucleotide of the window will be define by *annotation_end* column.

By default the coverages are aligned using the column col-ref as reference, if each window have not the same size (see above) the None value is used to padded the coverage matrix. There is 2 other options which alter this behavior

justify option

When this option is set. *crawl_coverage* get the coverages between the position start and stop and then generate new values by linear interpolation to have the same number of coverage measure for all genes. If this option is used with fix window and the window is out the gene (before the first position) only the position from 0 will take in account to generates new values.

```
crawl_coverage --bam ../WTE1.bam --annot ../annotations.txt --ref-col annotation_start_
↪--start-col annotation_start --stop-col annotation_end --justify 1000
```

This command will compute coverage using WTE1.bam and with annotations.txt file. For each gene, 1000 values will be estimated by linear interpolation from the raw coverage values.

sum option

When this option is set, instead of recording the coverage for each nucleotides, *crawl_coverage* compute the sum of the coverages of the window and report it on the coverage file.:

```
crawl_coverage --bam ../WTE1.bam --annot ../annotations.txt --ref-col annotation_start_
↪--start-col annotation_start --stop-col annotation_end --sum
```

This command will compute coverage using WTE1.bam and with annotations.txt file. For each gene, the sum of all coverages is computed and reported on the coverage file.

Other options

The following option are not mandatory:

- **-q QUAL_THR, -qual-thr QUAL_THR** The minimal quality of read mapping to take it in account. (default=15)
- **-s SUFFIX, -suffix SUFFIX** The name of the suffix to use for the output file. (default= .cov)
- **-o OUTPUT, -output OUTPUT** The path of the output (default= base name of annotation file with -suffix)
- **-version** display version information and quit.
- **-verbose, -v** increase the verbosity of the output (this option can be repeat several times as -vv).

- **-quiet** decrease verbosity of the output. By default `craw_coverage` is slightly verbose and display a progress bar. This option can be useful to disable any progression information on batch run.
- **-h -help** display the inline help and exit.

Warning: by default `craw_coverage` use a quality threshold of 15 (like `pysam`)

Note: strand column must named *strand* and can take *1/-1* or *+/- for/rev* as value for forward/reverse strands.

Warning: the coverage file can be huge depending on the number of gene to compute the coverage and the size of the window for instance for 6000 genes with a window of 15000 nt the cov file will weight almost 900Mb.

1.3.2 `craw_http`

Compute a figure from a file of coverage generated by `craw_coverage`. By default, display a figure with two heatmap one for the sense the other for the antisense. But it work also if the coverage file contains *sense* or *anti sense* data only.

Mandatory arguments

- **cov_file** The path to the coverage file (the output of).

Data options

- **-crop CROP CROP:** Crop the matrix. This option need two values the name of the first and last column to keep [start col, stop col] eg `-crop -10 1000`

```
craw_http --crop 0 2000 WTE1_var_window.cov
```

This command will display only column '0' to '2000', included, of the matrix generated by `craw_coverage`.

- **-sort-using-col COL** sort the data using the column name 'COL' (descending).
- **-sort-using-file SORT_USING_FILE** sort the rows using a file. The file must have on the first line the name of the column to use for sorting and each line must match to a value contained in the matrix.
- **-sort-by-gene-size [start_col,stop_col [start_col,stop_col ...]]** The rows will be sorted by gene size using *start_col* and *stop_col* to compute length. *start_col* and *stop_col* must be a string separated by comma. If *start_col* and *stop_col* are not specify *annotation_start,annotation_end* will be used.
- **-sense-only** Display only sense matrix (default is display both).
- **-antisense-only** Display only anti sense matrix (default is display both).

Warning: Don't put the **-sort-by-gene-size** option without value as last option just before the coverage file. In this case the `craw_http` will don't work. If you want to use only this option, use the **-v** option after **-sort-by-gene-size**

```
craw_http --sort-by-gene-size -v WTE1_0_2000.cov
```

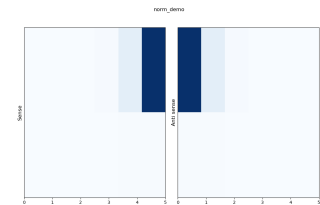
Figure options

Normalization options

craw_htmp provide several methods to normalize data before to display them: below the different figures illustrate the result of each normalization methods on the following matrix x .

Pos	0	1	2	3	4	5
S	0	1	10	100	1000	10000
AS	10000	1000	100	10	1	0
S	1	10	20	30	40	50
AS	50	40	30	20	10	1

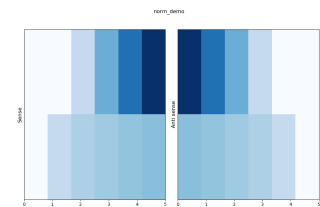
linear normalisation



–norm lin: a linear normalization is applied on the whole matrix.

- x the original matrix
- x_i the value of a cell in the original matrix
- z_i the value of a cell in the normalized matrix
- $z_i = \frac{x_i - \min(x)}{\max(x) - \min(x)}$

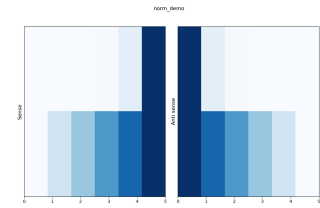
logarithmic normalisation



–norm log: a 10 base logarithm will be applied on the data before matrix normalization.

1. replace all 0 values by 1
2. $z_i = \log_{10}(x_i)$
3. $w_i = \frac{z_i - \min(z)}{\max(z) - \min(z)}$

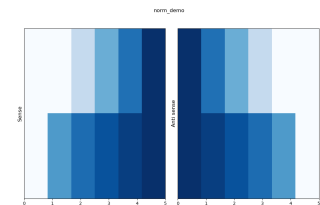
linear normalisation row by row



–norm row: mean that a linear normalisation is compute row by row.

- x the original matrix
- x_{ij} the value of a cell in the original matrix with I rows and J columns
- x_i the values of the i row
- $z_{ij} = \frac{x_{ij} - \min(x_i)}{\max(x_i) - \min(x_i)}$

logarithmic normalisation row by row



–norm log+row mean a 10 base logarithm will be applied before a normalisation row by row.

- x the original matrix
 - x_{ij} the value of a cell in the original matrix with I rows and J columns
 - x_i the values of the i row
1. replace all 0 values by 1
 2. $z_{ij} = \log_{10}(x_{ij})$
 3. $w_{ij} = \frac{z_{ij} - \min(z_i)}{\max(z_i) - \min(z_i)}$

Note:

- ‘row+log’ is an alias for ‘log+row’
 - The default normalisation is **lin**
-

Other figure options

- **–cmap CMAP** The color map used to display data. The allowed values are defined in http://matplotlib.org/examples/color/colormaps_reference.html eg: Blues, BuGn, Greens, GnBu, ... (default: Blues).

- **-title TITLE** The figure title. It will display on the top of the figure. (default: the name of the coverage file without extension).
- **-dpi DPI** The resolution of the output (default=100).

This option work only if **-out** option is specified. To set the right dpi for screen displaying use the [matplotlib configuration](#) file.

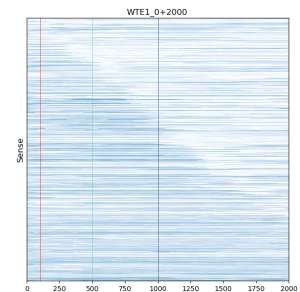
- **-size SIZE** Specify the figure size

The value must be `widexheight[unit]` or `'raw'`. If value is `'raw'` it will be produce two image files (for sense and antisense) with one pixel correspond to one coverage value. Otherwise, `'wide'` and `'height'` must be positive integers By default *unit* is in inches. eg:

- 7x10 or 7x10in for 7 inches wide by 10 inches height.
- 70x100mm for 70 mm by 100 mm.

default=7x10 or 10x7 depending of the figure orientation (see layout).

- **-mark POS <COLOR>** will draw a vertical line at the position POS with the color <COLOR>



COLOR can be the name of the most common html color red, yellow, ... or a value of a RGB in hexadecimal format like `#rgb` or `#rrggbb` for instance `#ff0000` represent pure red. (don't forget to surround the hexadecimal color with quotes on commandline)

If COLOR is omitted the color of the highest value of the color map used for the drawing will be used (The default color map is Blues).

```
craw_http --norm log --mark 1000 --mark 500 MediumAquamarine --mark 100 "#ff0000"
↪--sense-only WTE1_0+2000.cov
```

for list of HTML colors:

- https://en.wikipedia.org/wiki/Web_colors
- https://www.w3schools.com/colors/colors_names.asp

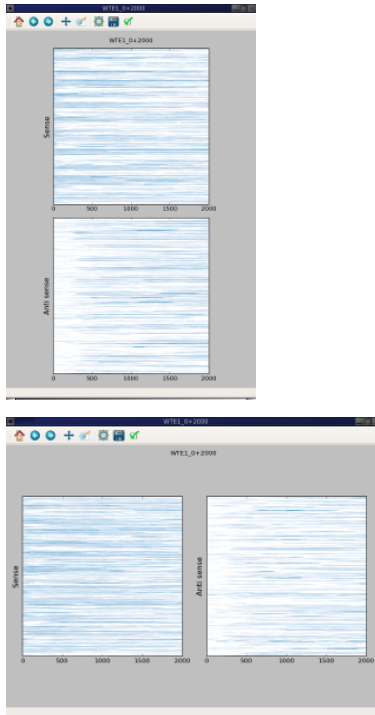
Warning: The **-mark** option must not be the last option on the command line (just before the coverage file), otherwise an error will occurred.:

```
craw_http --out my_fig.png --mark 10 red --mark 0 WTE1_0_2000.cov => raise an error
craw_http --mark 10 red --mark 0 --out my_fig.png WTE1_0_2000.cov => work
```

Layout options

- **-sense-on-left** Where to display the sense matrix relative to antisense matrix.

- **–sense-on-right** Where to display the sense matrix relative to antisense matrix.
- **–sense-on-top** Where to display the sense matrix relative to antisense matrix.
- **–sense-on-bottom** Where to display the sense matrix relative to antisense matrix.



The first screen capture uses *–sense-on-top* whereas the second capture used *–sense-on-left option*.

Note: default is top.

Other options

- **-h, –help** Display the help message and exit
- **–out OUT** The name of the file (the format will be based on the extension) to save the figure. Instead of displaying the figure on the screen, save it directly in this file.
- **-v, –verbose** Increase output verbosity. By default `craw_http` is relatively quiet (display only warning and error), if you want to display also the processing step just add `-v` on the commandline (or `-vv` to display also the debugging message).
- **–version** Display version information and quit.

1.4 Inputs / Outputs

1.4.1 `craw_coverage`

Inputs

`craw_coverage*` need a file bam or wig to compute coverage and an annotation file to specify on which regions to compute these coverages.

bam file

`craw_coverage` can use a file of alignment reads called bam file. a bam file is a short DNA sequence read alignments in the Binary Alignment/Map format (.bam). `craw_coverage` needs also the corresponding index file (bai). The index file must be located beside the bam file with the same name instead to have the .bam extension it end by .bai extension. If you have not the index file you have to create it.

To index a bam file you need samtools. The command line is

```
samtools index file.bam
```

For more explanation see <http://www.htslib.org/doc/>.

wig file

`craw_coverage` can compute coverage also from wig file see <https://wiki.nci.nih.gov/display/tcga/wiggle+format+specification> and <http://genome.ucsc.edu/goldenPath/help/wiggle.html> . for format specifications. Compared to these specifications `craw` support coverages on both strands. the positive coverages scores are on the forward strand whereas the negative ones are on the reverse strand.

```
track type=wiggle_0 name="demo" color=96,144,246 altColor=96,144,246 autoScale=on
→graphType=bar
variableStep chrom=chrI span=1
72      12.0000
73      35.0000
74      70.0000
75      127.0000
...
72      -88.0000
73      -42.0000
74      -12.0000
75      -1.0000
```

In the example above the coverage on the Chromosome I for the positions 72, 73, 74, 75 are 12, 35, 70, 127 on the forward strand and 88, 42, 12, 1 on the reverse strand.

annotation file

The annotation file is a *tsv* file by default. It's mean that it is a text file with value separated by tabulation (not spaces) or commas. But if a separator is specified (`-sep`) it can be a csv file or any columns file.

The first line of the file must be the name of the columns the other lines the values. Each line represent a row.

name	gene	chromosome	strand	Position
YEL072W	RMD6	chrV	+	14415
YEL071W	DLD3	chrV	+	17845
YEL070W	DSF1	chrV	+	21097
YEL066W	HPA3	chrV	+	27206

(continues on next page)

(continued from previous page)

YEL065W	SIT1	chrV	+	29543
YEL062W	NPR2	chrV	+	36254
YEL058W	PCM1	chrV	+	44925
YEL056W	HAT2	chrV	+	48373

All lines starting with '#' character will be ignored.

```
# This is the annotation file for Wild type
# bla bla ...
name    gene    chromosome    strand    Position
YEL072W RMD6    chrV         +         14415
YEL071W DLD3    chrV         +         17845
YEL070W DSF1    chrV         +         21097
YEL066W HPA3    chrV         +         27206
YEL065W SIT1    chrV         +         29543
YEL062W NPR2    chrV         +         36254
YEL058W PCM1    chrV         +         44925
YEL056W HAT2    chrV         +         48373
```

mandatory columns

There is 3 mandatory columns in the annotation file.

columns with fixed name

two with a fixed name:

- **strand** indicate on which strand is located the region of interest. The authorized values for this columns are +/- , 1/-1 or for/rev.
- **chromosome** the chromosome name where is located the region of interest.

columns with variable name

In addition of these two columns the column to define the position of reference is mandatory too, but the name of this column can be specified by the user. If it's not `craw_coverage` will use a column name 'position'.

If we want to compute coverage on variable window size, 2 extra columns whose name must be specified by the user by the following option:

- `--start-col` to define the beginning of the window (this position is included in the window)
- `--stop-col` to define the end of the window (this position is included in the window)

name	gene	type	chromosome	strand	annotation_start	annotation_end	transcription_start	transcription_end
YEL072W	RMD6	gene	chrV	1	13720	14415	1	14745
YEL071W	DLD3	gene	chrV	1	16355	17845	1	17881
YEL070W	DSF1	gene	chrV	1	19589	21097	1	21197
YEL066W	HPA3	gene	chrV	1	26721	27206	1	27625
YEL065W	SIT1	gene	chrV	1	27657	29543	1	29601
YEL062W	NPR2	gene	chrV	1	34407	36254	1	36401
YEL058W	PCM1	gene	chrV	1	43252	44925	1	44993

(continues on next page)

(continued from previous page)

YEL056W	HAT2	gene	chrV	1	47168	48373	1	48457	47105
YEL052W	AFG1	gene	chrV	1	56571	58100	1	58105	56537

```
craw_coverage --wig file.wig --annot annot.txt --ref-col annotation_start --start-col
↪annotation_start --stop-col annotation_end
```

The position of reference must be between start and end. The authorized values are positive integers.

Note: the position of reference can be used to define the reference and the start of the end of the window.

```
craw_coverage --bam file.bam --annot annot.txt --ref-col annotation_start --start-col
↪annotation_start --stop-col annotation_end
```

All other columns are not necessary but will be reported as is in the coverage file.

Outputs

coverage_file

It's a *tsv* file with all columns found in annotation file plus the result of coverage position by position centered on the reference position define for each line. for instance

```
craw_coverage --wig=../data/small.wig --annot=../data/annotations.txt
--ref-col=annotation_start --before=0 --after=2000
```

In the command line above, the column '0' correspond to the annotation_start position the column '1' to annotation_start + 1 on so on until '2000' (here we display only the first 3 columns of the coverage).

```
# Running Counter RnAseq Window craw_coverage
# Version: craw NOT packaged, it should be a development version | Python 3.4
# Using: pysam 0.9.1.4 (samtools 1.3.1)
#
# craw_coverage run with the following arguments:
# --after=3
# --annot=../data/annotation_wo_start.txt
# --before=5
# --chr-col=chromosome
# --output=small_wig.cov
# --qual-thr=0
# --quiet=1
# --ref-col=Position
# --sense=mixed
# --sep=
# --strand-col=strand
# --suffix=cov
# --verbose=0
# --wig=../data/small.wig
sense  name      gene      type      chromosome      strand  annotation_start  ↪
↪annotation_end  has_transcript  transcription_end  transcription_start  0  ↪
↪  1          2
S      YEL072W  RMD6      gene      chrV      +      13720  14415  1      14745  13569 ↪
↪  7          7          7
```

(continues on next page)

(continued from previous page)

AS	YEL072W	RMD6	gene	chrV	+	13720	14415	1	14745	13569
→ 0	0	0								
S	YEL071W	DLD3	gene	chrV	+	16355	17845	1	17881	16177
→ 31	33	33								

The line starting with '#' are comments and will be ignored for further processing. But in traceability/reproducibility concern, in the comments *craw_coverage* indicate the version of the program and the arguments used for this experiment.

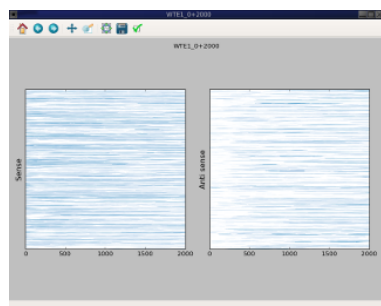
1.4.2 *craw_http*

Inputs

see *cov_out*

Outputs

The default output of *craw_http* (if *--out* is omitted) is graphical window on the screen. The figure display on the screen can be saved using the window menu.



It is also possible to generate directly a image file in various format by specifying the *--out* option. The output format will be deduced from the filename extension provide to *--out* option.

```
--out foo.jpeg for jpeg image or --out foo.png for png image
```

The supported format vary in function of the matplotlib backend used (see [matplotlib configuration](#)).

If *--size raw* is used 2 files will be generated one for the sense and the other for the antisense. If *--out* is not specified it will be the name of the coverage file without extension and the format will be png.

```
craw_http foo_bar.cov --size raw
```

will produce *foo_bar.sense.png* and *foo_bar.antisense.png*

```
craw_http foo_bar.cov --size raw --out Xyzzy.jpeg
```

will produce *Xyzzy.sense.jpeg* and *Xyzzy.antisense.jpeg*

DEVELOPER GUIDE

2.1 Overview

Scripts are located in *bin* directory, and use some modules located in *craw* directory.

- *craw_coverage* use module *craw.annotation* to handle annotation file and module *craw.coverage* to compute coverage this module rely on *pysam*.
- *craw_htmp* read coverage file generate by *craw_coverage* and produce graphical representation of data. This script use functions in module *craw.heatmap* in the form of heatmap. The module *craw.heatmap* have some capabilities to sort, crop, normalize data before represent them. this module rely on *numpy*, *pandas* to manipulate data (*craw.heatmap.sort*, *craw.heatmap.crop_matrix*, *craw.heatmap.lin_norm*, ...) and *matplotlib* and/or *pillow* to generate images (*craw.heatmap.draw_heatmap*, *craw.heatmap.draw_raw_image*)

2.2 reference API

2.2.1 *craw*

The *_get_version_message* is a private function that provide a human readable version of *craw* package and *python* which is common to all scripts. Each script have a public function *get_version_message* that call this function for the common part and add the version of all dependencies need for the script.

```
craw.get_version_message()
```

Returns A human readable version of the *craw* package version

Return type string

```
craw.init_logger(log_level, out=True)
```

Initiate the “root” logger for *craw* library all logger create in *craw* package inherits from this root logger This logger write logs on *sys.stderr*

Parameters *log_level* (*int*) – the level of the logger

2.2.2 *annotation*

The *annotation* module contains everything that is needed to parse annotation file and handle it.

AnnotationParser

The entry point to parse an annotation file is the `craw.annotation.AnnotationParser`. An annotation parser have two methods:

- `craw.annotation.AnnotationParser.get_annotations()` create a new type of `Entry` and iterate over the annotation file and for each line return a new instance of the newly `craw.annotation.Entry` class it just create on the fly.
- the other more technique give the maximum of nucleotides before and after the reference. It is needed to compute the size of the resulting matrix.

The force of this approach is to generate a new type of entry for each parsing. So it's very flexible and allow to fit with most of annotation file. But for one file, all the parsing use the same `Entry` class so it ensure the coherence in data.

new_entry_type

Is a factory which generate a new subclass of `craw.annotation.Entry` given the fields gather form the annotation file header (first line non starting with #) and the columns semantic given by the user. The first role of this factory is to check if all parameter given by user correspond ot header and do some coherence checking. If everything seems Ok it generate on the fly a new subclass of `craw.annotation.Entry`.

Entry Class

An Entry correspond to one line of the annotation file.

The Entry convert values if necessary (*strand* in a internal representation +/-, *position* in integer ...). It also expose a generic api to access some fields whatever the named of the columns.

annotation API reference

```
class craw.annotation.AnnotationParser (path,          ref_col,          strand_col='strand',
                                         chr_col='chromosome',      start_col=None,
                                         stop_col=None, sep='t')
```

Parse the annotation file

- create new type of `Entry` according to the header
- create one `Entry` object for each line of the file

```
__init__ (path,    ref_col,    strand_col='strand',    chr_col='chromosome',    start_col=None,
          stop_col=None, sep='t')
```

Parameters

- **path** (*string*) – the path to the annotation file to parse.
- **ref_col** (*string*) – the name of the column for the reference position
- **chr_col** (*string*) – the name of the column for the chromosome
- **strand_col** (*string*) – the name of the column for the strand
- **start_col** (*string*) – the name of the column for start position
- **stop_col** (*string*) – the name of the column for the stop position
- **sep** (*string*) – The separator tu use to split fields

__weakref__

list of weak references to the object (if defined)

get_annotations()

Parse an annotation file and yield a *Entry* for each line of the file.

Returns a generator on a annotation file.

max()

Returns the maximum of bases to take in count before and after the reference position.

Return type tuple of 2 int

class `craw.annotation.Entry(values)`

Handle one entry (One line) of annotation file

__eq__(*other*)

Return self==value.

__init__(*values*)

Parameters **values** (*list of string*) – the values parsed from one line of the annotation file

__str__()

Return str(self).

__weakref__

list of weak references to the object (if defined)

_convert(*field, value*)

Convert field parsed from annotation file in Entry internal value

Parameters

- **field** (*string*) – the field name associated to the value.
- **value** (*string*) – the value to convert

Returns the converted value

Return type any

Raise RuntimeError or value Error if a value cannot be converted

_switch_start_stop()

Switch start and stop value if self.start > self.stop This situation can occur if annotation regards the reverse strand

chromosome

The name of the Chromosome

header

The header of the annotation file

ref

The position of reference

start

The Position to start the coverage computation

stop

The position to end the coverage computation (included)

strand
the strand +/-

class `craw.annotation.Idx` (*col_name*, *idx*)

__getnewargs__ ()
Return self as a plain tuple. Used by copy and pickle.

static **__new__** (*_cls*, *col_name*, *idx*)
Create new instance of `Idx`(*col_name*, *idx*)

__repr__ ()
Return a nicely formatted representation string

_asdict ()
Return a new `OrderedDict` which maps field names to their values.

classmethod **_make** (*iterable*)
Make a new `Idx` object from a sequence or iterable

_replace (***kws*)
Return a new `Idx` object replacing specified fields with new values

col_name
Alias for field number 0

idx
Alias for field number 1

`craw.annotation.new_entry_type` (*name*, *fields*, *ref_col*, *strand_col*=*'strand'*,
chr_col=*'chromosome'*, *start_col*=*None*, *stop_col*=*None*)
From the header of the annotation line create a new `Entry` Class inherited from `Entry` Class

Parameters

- **name** (*str*) – The name of the new class of entry.
- **fields** (*list of string*) – The fields constituting the new type of entry.
- **ref_col** (*string*) – The name of the column representing the position of reference (default is *'position'*).
- **strand_col** (*string*) – The name of the column representing the strand (default is *'strand'*).
- **chr_col** (*string*) – The name of the column representing the name of chromosome (default is *'chromosome'*).
- **start_col** (*string*) – The name of the column representing the position of the first base to compute the coverage (inclusive).
- **stop_col** (*string*) – The name of the column representing the position of the last base to compute the coverage (inclusive).

Returns a new class child of `Entry` which is able to store information corresponding to the header.

2.2.3 wig

This module allow to parse wig files (wig file specifications are available here: <https://wiki.nci.nih.gov/display/tcga/wiggle+format+specification>, <http://genome.ucsc.edu/goldenPath/help/wiggle.html>). The wig file handle by this modules slightly differ from the canonic specifications as it allow to specify coverage on forward and reverse strand. If the coverage score is positive that mean that it's on the forward strand if it's negative, it's on the reverse strand.

The WigParser and helpers

The `craw.wig.WigParser` allow to parse the wig file. It read the file line by line, test the category of the line `trackLine`, `declarationLine` or `dataLine` and call the right method to parse the line and build the genome object.

The classes `craw.wig.VariableChunk` and `craw.wig.FixedChunk` are not keep in the final data model, they are just used to parsed the data lines and convert the wig file information (step, span) in coverages for each positions.

The data model to handle the wig information

The `craw.wig.Genome` objects contains `craw.wig.Chromosome` (each chromosomes are unique and the names of chromosomes are unique). Each chromosomes contains the coverage for the both strands. To get the coverage for region or a position just access it with indices or slices as traditional python list, tuple, on so on. The slicing return two lists. The first list correspond to the coverage on this particular region for the forward strand, the second element for the reverse strand. By default the chromosomes are initialized with 0.0 as coverage for all positions.

All information specified in the track line are stored in the `infos` attribute of `craw.wig.Genome` as a dict.

wig API reference

```
class craw.wig.Chromosome (name, size=1000000)
```

Handle chromosomes. A chromosome as a name and contains `Chunk` objects (forward and reverse)

```
__getitem__ (pos)
```

Parameters `pos` – a position or a slice (0 based) if `pos` is a slice the left indice is excluded

Returns the coverage at this position or corresponding to this slice.

Return type a list of 2 list of float `[[float,...],[float, ..]]`

Raises `IndexError` – if `pos` is not in coverage or one bound of slice is out the coverage

```
__init__ (name, size=1000000)
```

Parameters

- **name** (`str`) –
- **size** (the default size of the chromosome. Each time we try to set a value greater than the chromosome size is doubled. This is to protect the machine against memory swapping if the user provide a wig file with very big chromosomes.) –

```
__len__ ()
```

Returns the actual length of the chromosome

Return type `int`

```
__setitem__ (pos, value)
```

Parameters

- **pos** (`int` or slice object) – the position (0-based) to set value
- **value** (`float` or iterable of float) – value to assign

Raises

- **ValueError** – when `pos` is a slice and `value` have not the same length of the slice
- **TypeError** – when `pos` is a slice and `value` is not iterable
- **IndexError** – if `pos` is not in coverage or one bound of slice is out the coverage

__weakref__
list of weak references to the object (if defined)

_estimate_memory (*col_nb, mem_per_col*)
Parameters
• **col_nb** (*int*) – the number of column of the new array or the extension
• **mem_per_col** (*int*) – the memory needed to create or extend an array with one col and 2 rows fill with 0.0
Returns the estimation of free memory available after creating or extending chromosome
Return type *int*

_extend (*size=1000000, fill=0.0*)
Extend this chromosome of the size *size* and fill with *fill*. :param *size*: the size (in bp) we want to increase the chromosome. :type *size*: *int* :param *fill*: the default value to fill the chromosome. :type *fill*: *float* or *nan* :raise *MemoryError*: if the chromosome extension could overcome the free memory.

class *craw.wig.Chunk* (***kwargs*)
Represent the data following a declaration line. The a *Chunk* contains sparse data on coverage on a region of one chromosomes on both strand plus data contains on the declaration line.

__init__ (***kwargs*)
Parameters **kwargs** (*dictionary*) – the key,values pairs found on a Declaration line

__weakref__
list of weak references to the object (if defined)

is_fixed_step ()
This is an abstract methods, must be implemented in inherited class :return: True if i's a fixed chunk of data, False othewise :rtype: *boolean*

parse_data_line (*line, chrom, strand_type*)
parse a line of data and append the results in the corresponding strand This is an abstract methods, must be implemented in inherited class.
Parameters
• **line** (*string*) – line of data to parse (the white spaces at the end must be strip)
• **chrom** (*Chromosome* object.) – the chromosome to add coverage data
• **strand_type** (*string* '+' , '-' , 'mixed') – which kind of wig is parsing: forward, reverse, or mixed strand

class *craw.wig.FixedChunk* (***kwargs*)
The *FixedChunk* objects handle data of 'fixedStep' declaration line and it's coverage data

__init__ (***kwargs*)
Parameters **kwargs** (*dictionary*) – the key,values pairs found on a Declaration line

is_fixed_step ()
Returns *True*
Return type *boolean*

parse_data_line (*line, chrom, strand_type*)
parse line of data following a fixedStep Declaration. add the result on the corresponding strand (forward if coverage value is positive, reverse otherwise) :param *line*: line of data to parse (the white spaces at the end must be strip) :type *line*: *string* :param *chrom*: the chromosome to add coverage data :type *chrom*: *Chromosome* object. :param *strand_type*: which kind of wig is parsing: forward, reverse, or mixed strand :type *strand_type*: *string* '+' , '-' , 'mixed'

```
class craw.wig.Genome
```

A genome is made of chromosomes and some metadata, called infos

```
__delitem__ (name)
```

remove a chromosome from this genome

Parameters *name* (*string*) – the name of the chromosome to remove

Returns None

```
__getitem__ (name)
```

Parameters *name* (*string*) – the name of the chromosome to retrieve

Returns the chromosome corresponding to the name.

Return type *Chromosome* object.

```
__init__ ()
```

Initialize self. See help(type(self)) for accurate signature.

```
__weakref__
```

list of weak references to the object (if defined)

```
add (chrom)
```

add a chromosome in to a genome. if a chromosome with the same name already exist the previous one is replaced silently by this one.

Parameters *chrom* (*Chromosome* object.) – a chromosome to ad to this genome

Raise *TypeError* if *chrom* is not a *Chromosome* object.

```
class craw.wig.VariableChunk (**kwargs)
```

The Variable Chunk objects handle data of ‘variableStep’ declaration line and it’s coverage data

If in data there is negative values this indicate that the coverage match on the reverse strand. the chunk start with the smallest position and end to the highest position whatever on wich strand are these position. This mean that when the chunk will be convert in Coverage, the lacking positions will be filled with 0.0.

for instance:

```
variableStep chrom=chr3 span=2 10 11 20 22 20 -30 25 -50
```

will give coverages starting at position 10 and ending at 26 for both strands and with the following coverages values

```
for = [11.0, 11.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 22.0, 22.0, 0.0, 0.0, 0.0, 0.0]
```

```
rev = [0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 30.0, 30.0, 0.0, 0.0, 0.0, 50.0, 50.0]
```

```
is_fixed_step ()
```

Returns False

Return type boolean

```
parse_data_line (line, chrom, strand_type)
```

Parse line of data following a variableStep Declaration. Add the result on the corresponding strand (forward if coverage value is positive, reverse otherwise)

Parameters

• *line* (*string*) – line of data to parse (the white spaces at the end must be strip)

• *chrom* (*Chromosome* object.) – the chromosome to add coverage data

• *strand_type* (*string* ‘+’ , ‘-’, ‘mixed’) – which kind of wig is parsing: forward, reverse, or mixed strand

Raises *ValueError* – if *strand_type* is different than ‘mixed’, ‘-’, ‘+’

exception `craw.wig.WigError`

Handle error related to wig parsing

__weakref__

list of weak references to the object (if defined)

class `craw.wig.WigParser` (*mixed_wig*="", *for_wig*="", *rev_wig*="")

class to parse file in wig format. at the end of parsing it returns a *Genome* object.

__init__ (*mixed_wig*="", *for_wig*="", *rev_wig*="")

Parameters

- **mixed_wig** (*string*) – The path of the wig file to parse. The wig file code for the 2 strands:
 - The positive coverage values for the forward strand
 - The negative coverage values for the reverse strandThis parameter is incompatible with *for_wig* and *rev_wig* parameter.
- **for_wig** (*string*) – The path of the wig file to parse. The wig file code for forward strand only. This parameter is incompatible with *mixed_wig* parameter.
- **rev_wig** (*string*) – The path of the wig file to parse. The wig file code for reverse strand only. This parameter is incompatible with *mixed_wig* parameter.

__weakref__

list of weak references to the object (if defined)

static is_comment_line (*line*)

Parameters **line** (*string*) – line to parse.

Returns True if line is a comment line. False otherwise.

Return type boolean

is_data_line (*line*)

Parameters **line** – line to parse.

Returns True if it's a data line, False otherwise

is_declaration_line (*line*)

A single line, beginning with one of the identifiers *variableStep* or *fixedStep*, followed by attribute/value pairs for instance:

```
fixedStep chrom=chrI start=1 step=10 span=5
```

Parameters **line** (*string*) – line to parse.

Returns True if line is a declaration line. False otherwise.

Return type boolean

static is_track_line (*line*)

A track line begins with the identifier *track* and followed by attribute/value pairs for instance:

```
track type=wiggle_0 name="fixedStep" description="fixedStep format"  
↳visibility=full autoScale=off
```

Parameters **line** (*string*) – line to parse.

Returns True if line is a track line. False otherwise.

Return type boolean

parse()

Open a wig file and parse it. read wig file line by line check the type of line and call the corresponding method accordingly the type of the line: - comment - track - declaration - data see - <https://wiki.nci.nih.gov/display/tcga/wiggle+format+specification> - <http://genome.ucsc.edu/goldenPath/help/wiggle.html> for wig specifications. This parser does not fully follow these specification. When a score is negative, it means that the coverage is on the reverse strand. So some positions can appear twice in one block of declaration (what I call a chunk).

Returns a Genome coverage corresponding to the wig files (mixed strand on one wig or two separate wig)

Return type *Genome* object

parse_data_line(line, strand_type)

Parameters **line** (*string*) – line to parse. It must not a comment_line, neither a track line nor a declaration line.

Raises **ValueError** – if strand_type is different than ‘mixed’, ‘-’, ‘+’

parse_declaration_line(line)

Get the corresponding chromosome create one if necessary, and set the current_chunk and current_chromosome.

Parameters **line** – line to parse. The method *is_declaration_line()* must return True with this line.

parse_track_line(line, strand_type=)

fill the genome infos with the information found on the track.

Parameters **line** – line to parse. The method *is_track_line()* must return True with this line.

2.2.4 coverage

coverage module contain several functions which allow to get the coverages from data input, a *craw.wig.Genome* object or a *pysam.AlignmentFile*.

There is 2 kind of functions:

- the functions to get coverage from input.
- the functions to process the coverage.

Functions to get coverage

These low level functions are not aimed to be called directly. They are called inside function which process the coverages.

get_raw_bam_coverage

Get coverage from *pysam* for reference (*chromosome*) for an interval of positions, a quality on both strand. and convert the coverage return by *pysam*. A score on each position for each base (ACGT)) in a global coverage for this position.

This function is called for each entry of the annotation file.

get_raw_wig_coverage

Get coverage from `craw.wig.Genome` instance for reference (*chromosome*) for an interval of positions, on both strand. The quality parameter is here just to have the same signature as `get_bam_coverage` but will be ignores .

This function is called for each entry of the annotation file.

get_raw_coverage_function

Allow to choose the right `get_raw(*)_coverage` in function of the data input type (`craw.wig.Genome`, `pysam.AlignmentFile`)

```
craw.coverage.get_raw_wig_coverage (genome, annot_entry, start, stop, qual_thr=None)
```

Parameters

- **genome** (`craw.wig.Genome` object) – The genome which store all coverages.
- **annot_entry** (`annotation.Entry` object) – an entry of the annotation file
- **start** (`int`) – The position to start to compute the coverage(coordinates are 0-based, start position is included).
- **stop** (`int`) – The position to stop to compute the coverage (coordinates are 0-based, stop position is excluded).
- **qual_thr** (`None`) – this parameter is not used, It's here to have the same api as `get_bam_coverage`.

Returns the coverage (all bases)

Return type tuple of 2 list containing int or float

```
craw.coverage.get_raw_bam_coverage (sam_file, annot_entry, start, stop, qual_thr=15)
```

Compute the coverage for a region position by position on each strand

Parameters

- **sam_file** (`pysam.AlignmentFile` object.) – the samfile openend with pysam
- **annot_entry** (`annotation.Entry` object) – an entry of the annotation file
- **start** (*positive int*) – The position to start to compute the coverage(coordinates are 0-based, start position is included).
- **stop** (*positive int*) – The position to start to compute the coverage (coordinates are 0-based, stop position is excluded).
- **qual_thr** (`int`) – The quality threshold

Returns the coverage (all bases)

Return type tuple of 2 list containing int

```
craw.coverage.get_raw_coverage_function (input)
```

Parameters **input** (`wig.Genome` or `pysam.calignmentfile.AlignmentFile` object) – the input either a samfile (see `pysam` library) or a genome build from a wig file (see `wig` module)

Returns `get_wig_coverage` or `get_bam_coverage` according the type of input

Return type function

Raises RuntimeError – when input is not instance of `pysam.AlignmentFile` or `wig.Genome`

Functions to process coverages

These functions guess the right `get_raw_(*)_coverage` in function of the data input and pass it to a post processing function.

all functions returned have the same API

3 parameters as input

- **annot_entry**: an entry of the annotation file.
- **start**: The position to start to compute the coverage (coordinates are 0-based, start position is included).
- **stop**: The position to stop to compute the coverage (coordinates are 0-based, stop position is excluded).

and

- **return**: a tuple of two list or tuple containing in this order the coverages on the forward strand then the coverages on the reverse strand.

These architecture allow to combine easily the different `get_raw_coverage` function with the different post-processing. For instance:

```
bam = pysam.AlignmentFile(bam_file, "rb")
get_coverage = get_padded_coverage(bam, max_left, max_right, qual_thr=15)
forward, reverse = get_coverage(annot_entry, 10 200)
```

or

```
wig = wig_parser.parse()
get_coverage = get_resized_coverage(wig, 200)
forward, reverse = get_coverage(annot_entry, 10 200)
```

`craw.coverage.padded_coverage_maker(input_data, max_left, max_right, qual_thr=None)`

Parameters

- **input_data** (`wig.Genome` or `pysam.AlignmentFile` object) – the input either a samfile (see `pysam` library) or a genome build from a wig file (see `wig` module)
- **max_left** (`int`) – The highest number of base before the reference position to take in account.
- **max_right** (`int`) – The highest number of base after the reference position to take in account.
- **qual_thr** (`int`) – The quality threshold if input data come from wig this parameter is not used,

Returns

a function `get_padded_coverage()`, a function which compute the coverage for a gene on each strand between position `[start, stop]`

The coverage values are centered on the `annot_entry.ref` position, the matrix is padded by `None` value.:

```
[.....[ coverage ref.pos ] .....]  
[...[coverage      ref.pos ] .....]  
[.....[ cov ref.pos      ] ]
```

This function take 3 parameters:

- **annot_entry**: an entry of the annotation file.
- **start**: The position to start to compute the coverage(coordinates are 0-based, start position is included).
- **stop**: The position to stop to compute the coverage (coordinates are 0-based, stop position is excluded).

and

- **return**: a tuple with 2 tuple of float or int representing the coverage on strand forward and reverse.

Return type function

```
craw.coverage.resized_coverage_maker (input_data, new_size, qual_thr=None)
```

Parameters

- **input_data** (*craw.wig.Genome* or *pysam.AlignmentFile* object) – the input either a samfile (see *pysam* library) or a genome build from a wig file (see *wig* module)
- **new_size** (*postive int*) – the number of values in the coverage vector.
- **qual_thr** (*int*) – The quality threshold if input data come from wig this parameter is not used,

Returns

a function `get_resized_coverage()`, a function which compute the coverage for a gene on each strand between position [start, stop[This function take 3 parameters: the coverage values are generate by linear interpolation from raw values between [start, stop[using the *scipy*. see <https://docs.scipy.org/doc/scipy-0.19.0/reference/generated/scipy.interpolate.interp1d.html>

- **annot_entry**: an entry of the annotation file.
- **start**: The position to start to compute the coverage(coordinates are 0-based, start position is included).
- **stop**: The position to stop to compute the coverage (coordinates are 0-based, stop position is excluded).

and

- **return**: a tuple with 2 tuple of float or int representing the coverage on strand forward and reverse.

Return type function

```
craw.coverage.sum_coverage_maker (input_data, qual_thr=None)
```

This function return a new function `get_sum_coverage()`

Parameters

- **input_data** (*craw.wig.Genome* or *pysam.AlignmentFile* object) – the input either a samfile (see *pysam* library) or a genome build from a wig file (see *wig* module)

- **qual_thr** (*int*) – The quality threshold if input data come from wig this parameter is not used,

Returns

`get_sum_coverage()`, a function which compute the sum of coverage for a gene on each strand between position [start, stop] This function take 3 parameters:

- **annot_entry**: an entry of the annotation file.
- **start**: The position to start to compute the coverage(coordinates are 0-based, start position is included).
- **stop**: The position to stop to compute the coverage (coordinates are 0-based, stop position is excluded).

and

- **return**: a tuple with 2 tuple of float or int representing the coverage on strand forward and reverse.

Return type function

2.2.5 heatmap

`craw.heatmap.split_data()` and `craw.heatmap.sort()` work on data freshly parsed from coverage file. That mean that the data contain the metadata (all columns which are not coverage scores like chromosome, position strand , on so on)

The other functions sort normalization function work on pandas 2D DataFrame or numpy arrays containing only scores of coverage. That means all metadata was removed (`craw.heatmap.remove_metadata()`).

sort

The is one public sort function which act as proxy for several private sorting function.

normalisation

Several functions to normalize data.

The data can be normalize using min max of the whole data. Or the min max is recalculated for each row.

in both case the formula is

$$z_i = x_i - \min(x) / \max(x) - \min(x)$$

where $x=(x_1, \dots, x_n)$ and z_i is now your with normalized data. in first case x is the whole matrix in 2nd is the row.

Normalization can be precede by 10 base log transformation.

Note: In this case all 0 values are replace by 1 (10 base log is not define)

drawing heatmap

There are 2 way to generates figures, the first one is to generate a figures containing 2 heatmap for sense or antisense with axis, legend on so on. But in this representation it's not possible to display a figure with no scaling out/in. So the information of one pixel is not accessible. This representation is generate by `craw.heatmap.draw_heatmap()` and use matplotlib.

The second representation is to produce raw image where one nucleotide (one position for one gene) is represent by one pixel without any scale in/out. In this representation there si not axis legend on so on it's only a raw image.

heatmap API reference

class `craw.heatmap.Mark` (*pos, data, color_map, color=None*)

A mark is a position and a color tight together. It is used to draw a colored vertical line at the given position on the heatmap

`__init__` (*pos, data, color_map, color=None*)

Parameters

- **pos** (*int*) – The position where to draw a mark, the position is relative to the reference position (0)
- **data** (`pandas.DataFrame` object) – the coverage matrix
- **color_map** (*(class`matplotlib.pyplot.ColorMap` object)*) – the color map used to draw the heatmap
- **color** – the color of the line, the supported formats are - hexadecimal values as `#rgb` or `#rrggbb`, for instance `#ff0000` is pure red. - common html color names

`__weakref__`

list of weak references to the object (if defined)

`_color_converter` (*color, data*)

Parameters

- **color** (*string*) – the color of the line, the supported formats are - hexadecimal values as `#rgb` or `#rrggbb`, for instance `#ff0000` is pure red. - common html color names
- **data** (`pandas.DataFrame` object) – the matrix coverage

Returns `rgb` color

Return type tuple with 3 int between 0 and 255

`_get_matrix_bound` (*data*)

Parameters **data** (`pandas.DataFrame` object) – the matrix coverage

Returns the most right and left position of the coverage

Return type tuple of 2 int

`to_px` ()

translate the position of the mark relative to the reference in pixel. :return: the position of the mark in pixel. :rtype: positive int

```
craw.heatmap._sort_by_gene_size(data, start_col=None, stop_col=None, ascending=True)
```

Sort the matrix in function of the gene size.

Parameters

- **data** (`pandas.DataFrame.`) – the data to sort.
- **start_col** (`string.`) – the name of the column representing the beginning of the gene.
- **stop_col** (`string`) – the name of the column representing the end of the gene.

Returns sorted data.

Return type a `pandas.DataFrame` object.

```
craw.heatmap._sort_using_col(data, col=None, ascending=True)
```

Sort the matrix in function of the column col

Parameters

- **data** (`pandas.DataFrame.`) – the data to sort.
- **col** (`string.`) – the name of the column to use for sorting the data.

Returns sorted data.

Return type a `pandas.DataFrame` object.

```
craw.heatmap._sort_using_file(data, file=None)
```

Sort the matrix in function of file. The file must have the following structure the first line must be the name of the column the following lines must be the values, one per line each line starting by '#' will be ignore.

Parameters

- **data** (`pandas.DataFrame.`) – the data to sort.
- **file** (*a file like object.*) – The file to use as guide to sort the data.

Returns sorted data.

Return type a `pandas.DataFrame` object.

```
craw.heatmap.crop_matrix(data, start_col, stop_col)
```

Crop matrix (remove columns). The resulting matrix will be [start_col, stop_col]

Parameters

- **data** (a 2D `pandas.DataFrame` object.) – the data to sort.
- **start_col** (`string.`) – The name of the first column to keep.
- **stop_col** (`string.`) – The name of the last column to keep.

Returns sorted data.

Return type a 2D `pandas.DataFrame` object or None if data is None.

```
craw.heatmap.draw_heatmap(sense, antisense, color_map=<matplotlib.colors.LinearSegmentedColormap
                        object>, title="", sense_on='top', size=None,
                        marks=None)
```

Create a figure with subplot to represent the data as heat map.

Parameters

- **sense** (a `pandas.DataFrame` object.) – the data normalized (xi in [0,1]) representing coverage on sense.
- **antisense** – the data normalized (xi in [0,1]) representing coverage on anti sense.
- **color_map** (a `matplotlib.pyplot.cm` object.) – the color map to use to represent the data.
- **title** (*string*.) – the figure title (by default the same as the coverage file).
- **sense_on** (*string*.) – specify the lay out. Where to place the heat map representing the sense data. the available values are: 'left', 'right', 'top', 'bottom' (default = 'top').
- **size** (*tuple of 2 float*.) – the size of the figure in inches (wide, height).
- **marks** (list of *Mark* object) – list of vertical marks

Returns The figure.

Return type a `matplotlib.pyplot.Figure` object.

`craw.heatmap.draw_one_matrix(mat, ax, cmap=<matplotlib.colors.LinearSegmentedColormap object>, y_label=None, marks=None)`

Draw a matrix using matplotlib imshow object

Parameters

- **mat** (a `pandas.DataFrame` object.) – the data to represent graphically.
- **ax** (a `matplotlib.axis` object) – the axis where to represent the data
- **cmap** (a `matplotlib.pyplot.cm` object.) – the color map to use to represent the data.
- **y_label** (*string*) – the label for the data draw on y-axis.
- **marks** (list of *Mark* object) – list of vertical marks

Returns the mtp image corresponding to data

Return type a `matplotlib.image` object.

`craw.heatmap.draw_raw_image(data, out_name, color_map=<matplotlib.colors.LinearSegmentedColormap object>, format='PNG', marks=None)`

Generate an image file with one pixel for each values of the data matrix. the data can be either the coverage on sense or on antisense.

Parameters

- **data** (2D `pandas.DataFrame` or `numpy.array` object) – a **Normalized** (where all values are between 0 and 1) matrix.
- **out_name** (*string*) – The name of the generated graphic file.
- **color_map** –
- **format** (*string*) – the format of the result png, jpeg, ... (see pillow supported formats)
- **marks** (a sequence (list, tuple or set) of *Mark* objects) – the marks (vertical rule) to draw on the resulting heat map

Raise `RuntimeError` if data are not normalized.

```
craw.heatmap.get_data(coverage_file)
```

Parameters `coverage_file` (*str*) – the path of the coverage file to parse.

Returns the data as 2 dimension dataframe

Return type a `pandas.DataFrame` object

```
craw.heatmap.lin_norm(data)
```

Normalize data with linear algorithm. The formula applied to obtain the results is:

$$z_i = x_i - \min(x) / \max(x) - \min(x)$$

where $x=(x_1, \dots, x_n)$ and z_i is now your with normalized data. Ensure that the resulting values are comprise between 0 and 1. return `None` if data is `None`, return empty `pd.DataFrame` object if data is empty.

Parameters `data` (a 2D `pandas.DataFrame` object.) – the data to normalize, this 2D matrix must contains only coverage scores (no more metadata).

Returns a normalize matrix, where $0 \leq z_i \leq 1$ where $z=(z_1, \dots, z_n)$

Return type a 2D `pandas.DataFrame` object or `None` if data is `None`.

```
craw.heatmap.lin_norm_row_by_row(data)
```

Normalize data with linear algorithm but instead to normalize all the matrix, the normalization formula (see `normalize()`) is applied row by row. It ensure that all values are between 0 and 1.

Parameters `data` (a 2D `pandas.DataFrame` object.) – the data to normalize, this 2D matrix must contains only coverage scores (no more metadata).

Returns a normalize matrix, where $0 \leq z_i \leq 1$ where $z=(z_1, \dots, z_n)$

Return type a 2D `pandas.DataFrame` object or `None` if data is `None`.

```
craw.heatmap.log_norm(data)
```

The base 10 logarithm is compute for all values before a normalization (see `normalize()`) to ensure that all values are comprise between 0 and 1 .

Note: coverage scores are integers ≥ 0 . $\log_{10}(0) = -\infty$ or warning in macos prior to normalize data the 0 values are replace by 1.

Parameters `data` (a 2D `pandas.DataFrame` object.) – the data to normalize, this 2D matrix must contains only coverage scores (no more metadata).

Returns a normalize matrix, where $0 \leq z_i \leq 1$ where $z=(z_1, \dots, z_n)$

Return type a 2D `pandas.DataFrame` object or `None` if data is `None`.

```
craw.heatmap.log_norm_row_by_row(data)
```

as `normalize_row_by_row()` but prior normalisation a 10 base logarithm is applied.

Note: coverage scores are integers ≥ 0 . $\log_{10}(0) = -\infty$ to normalize data the $-\infty$ value are change in 0.

Parameters `data` (a 2D `pandas.DataFrame` object.) – the data to normalize, this 2D matrix must contains only coverage scores (no more metadata).

Returns a normalize matrix, where $0 \leq z_i \leq 1$ where $z=(z_1, \dots, z_n)$

Return type a 2D `pandas.DataFrame` object or None if data is None.

`craw.heatmap.remove_metadata(data)`

Remove all information which is not coverage value (as chromosome, strand, name, ...)

Parameters `data` (`pandas.DataFrame`.) – the data coming from a coverage file parsing containing coverage information and metadata chromosome, gene name , ...

Returns sorted data.

Return type a 2D `pandas.DataFrame` object or None if data is None.

`craw.heatmap.sort(data, criteria, **kwargs)`

Sort the matrix in function of criteria. This function act as proxy for several specific sorting functions

Parameters

- **data** (`pandas.DataFrame`.) – the data to sort.
- **criteria** (`string`.) – which criteria to use to sort the data (by_gene_size, using_col, using_file).
- **kwargs** – depending of the criteria - start_col, stop_col for sort_by_gene_size - col for using_col - file for using file

Returns sorted data.

Return type a `pandas.DataFrame` object.

`craw.heatmap.split_data(data)`

Split the matrix in 2 matrices one for sense the other for antisense.

Parameters `data` (a 2 dimension `pandas.DataFrame` object) – the coverage data to split

Returns two matrix

Return type tuple of two `pandas.DataFrame` object (sense `pandas.DataFrame`, antisense `pandas.DataFrame`)

2.2.6 argparse_util

Some utilities to improve the `argparse` module.

```
class craw.argparse_util.VersionAction(option_strings, version=None,
                                       dest='==SUPPRESS==', de-
                                       fault='==SUPPRESS==', help="show program's
                                       version number and exit")
```

Class to allow `argparse` to handel more complex version output

`__call__` (`parser`, `namespace`, `values`, `option_string=None`)

Override the `argparse._VersionAction.__call__()` to use a `RawTextHelpFormatter` only for version action whatever the class_formatter specified for the `argparse.ArgumentParser` object.

2.2.7 craw_coverage

`craw_coverage.py` is the entry point script to compute coverage.

```
craw.scripts.craw_coverage.get_result_header(annot_parser, parsed_args)
```

Compute the header for the results. the first lines start with # they contains some general information about the craw (version) and options used (for tracibility) the last line is the header of columns separated by `--sep` option and can be used as header with pandas

Parameters

- **annot_parser** (`annotation.AnnotationParser` object) – the annotation parser
- **parsed_args** (`argparse.Namespace`) – the command line argument parsed with `argparse`

Returns The header of the result file

Return type `str`

```
craw.scripts.craw_coverage.get_results_file(sense_opt, basename, suffix)
```

Parameters

- **sense_opt** (`str`) – how to managed the sense and antisense results
 - **mixed**: sense and antisense are interleaved in same file
 - **split**: sense and antisense are in separated files
 - **S**: only sense results are write down
 - **AS**: only antisense are write down
- **basename** (`str`) – the basename of the results file
- **suffix** (`str`) – the suffix of the results file

Returns the file objects where to write sense and antisense results

Return type `tuple (file object sense, file object antisense)`

```
craw.scripts.craw_coverage.get_version_message()
```

Returns The human readable CRAW version.

Return type `str`

```
craw.scripts.craw_coverage.main(args=None, log_level=None)
```

The entrypoint for `craw_coverage` script It will generate a coverage matrix around the position of interest and write the results in files

Parameters

- **args** (*list of string as given by `sys.argv` without the program name*) – the arguments and options given on the command line
- **log_level** (*positive int or logging flag `logging.DEBUG`, `logging.INFO`, `logging.ERROR`, `logging.CRITICAL`*) – the level of logger

```
craw.scripts.craw_coverage.parse_args(args)
```

Parameters **args** (*list of string*) – The options set on the command line (without the program name)

Returns

```
craw.scripts.craw_coverage.positive_int(value)
```

Parse value given by the parser

Parameters **value** (*string*) – the value given by the parser

Returns the integer corresponding to the value

Return type int

Raise `argparse.ArgumentTypeError`

`craw.scripts.craw_coverage.quality_checker` (*value*)

Parse value given by the parser

Parameters **value** (*string*) – the value given by the parser

Returns the integer ≥ 0 and ≤ 42 corresponding to the value

Return type int

Raise `argparse.ArgumentTypeError` if value does not represent a integer ≥ 0 and ≤ 42

2.2.8 `craw_http`

`craw_http.py` is the entry point script to compute heatmap.

`craw.scripts.craw_http._file_readable` (*value*)

check value given by the parser

Parameters **value** (*str*) – the value given by the parser for `--sort-using-file` option

Returns the normpath of the value

Raises `ArgumentError` if the file does not exists, or is not a file, or not readable

`craw.scripts.craw_http._gene_size_parser` (*value*)

Parse value given by the parser

Parameters **value** (*string*) – the value given by the parser for `--sort-by-gene-size` option

Returns name of column representing the start of gene, name of the column representing the end of gene

Return type tuple of 2 string

Raise `argparse.ArgumentError` object if value cannot be parsed

`craw.scripts.craw_http._size_fig_parser` (*value*)

Parse value given by the parser for `--size` option the value must follow the syntax `widexheight[unit]` if the unit is omitted unit is inch otherwise unit must be

- 'mm' for millimeters
- 'cm' for centimeters
- 'in' for inches
- 'px' for pixels

wide and height must be positive integers.

Parameters **value** (*string*) – the size of the figure

Returns the size in inch ask by the user

Return type tuple of float

Raise `argparse.ArgumentError` object

`craw.scripts.craw_http.get_version_message` ()

Returns a human readable of `craw_htmp` version and it's main dependencies.

Return type `str`

```
craw.scripts.craw_htmp.main(args=None, log_level=None, logger_out=True)
```

The entrypoint for `craw_html` script

It will generate a heatmap representing the coverage matrix around the position of interest it can display the results on the screen or write it on file depending the options

Parameters

- **args** (*list of string*) – The arguments and option representing the command line
- **log_level** – the level of verbosity
- **logger_out** – True if you want to display logs on stdout, False otherwise

Returns

```
craw.scripts.craw_htmp.parse_args(args)
```

Parameters **args** – the arguments and option as provided by `sys.argv` without the program name

Returns the argument parsed

Return type `argparse.Namespace` object

2.3 Release notes

2.3.1 Release notes

Release 0.9.0

BUG FIX

- **craw_coverage** When feature is located on reverse strand the coverage must be reverse to be read in 5'→3' ([issue 27](#))

New Feature

- **craw_coverage** add support of 2 wig files one for the forward strand the other for the reverse. ([issue 26](#))
- **craw_htmp** allow to draw a vertical line in a given color at a given position. see `–mark` option ([issue 25](#))

Release 0.8.0

This release implements the support of wig files.

- Support of wig file as input data. The user can not only provide bam file but also wig file with `–wig` option.
- The annotation file support all columns file, not only tabular file (tsv). The tabular format is still the default, but by specifying the columns separator with the `–sep` option the annotation file can be in any format for csv file use `–sep ‘,’`

- the control of verbosity has been extend with -v to increase it and -q to decrease it. These options are cumulative so it possible to use -vv or -qq
- improve traceability by displaying `craw` version and all python dependencies version `pysam`, `pandas`, `numpy`, `matplotlib`, ...

INDICES AND TABLES

- `genindex`
- `modindex`
- `search`

PYTHON MODULE INDEX

C

- `craw`, [17](#)
- `craw.annotation`, [18](#)
- `craw.argparse_util`, [34](#)
- `craw.coverage`, [27](#)
- `craw.heatmap`, [30](#)
- `craw.scripts.craw_coverage`, [34](#)
- `craw.scripts.craw_http`, [36](#)
- `craw.wig`, [21](#)

Symbols

- `__call__()` (*craw.argparse_util.VersionAction method*), 34
- `__delitem__()` (*craw.wig.Genome method*), 23
- `__eq__()` (*craw.annotation.Entry method*), 19
- `__getitem__()` (*craw.wig.Chromosome method*), 21
- `__getitem__()` (*craw.wig.Genome method*), 23
- `__getnewargs__()` (*craw.annotation.Idx method*), 20
- `__init__()` (*craw.annotation.AnnotationParser method*), 18
- `__init__()` (*craw.annotation.Entry method*), 19
- `__init__()` (*craw.heatmap.Mark method*), 30
- `__init__()` (*craw.wig.Chromosome method*), 21
- `__init__()` (*craw.wig.Chunk method*), 22
- `__init__()` (*craw.wig.FixedChunk method*), 22
- `__init__()` (*craw.wig.Genome method*), 23
- `__init__()` (*craw.wig.WigParser method*), 24
- `__len__()` (*craw.wig.Chromosome method*), 21
- `__new__()` (*craw.annotation.Idx static method*), 20
- `__repr__()` (*craw.annotation.Idx method*), 20
- `__setitem__()` (*craw.wig.Chromosome method*), 21
- `__str__()` (*craw.annotation.Entry method*), 19
- `__weakref__` (*craw.annotation.AnnotationParser attribute*), 18
- `__weakref__` (*craw.annotation.Entry attribute*), 19
- `__weakref__` (*craw.heatmap.Mark attribute*), 30
- `__weakref__` (*craw.wig.Chromosome attribute*), 21
- `__weakref__` (*craw.wig.Chunk attribute*), 22
- `__weakref__` (*craw.wig.Genome attribute*), 23
- `__weakref__` (*craw.wig.WigError attribute*), 24
- `__weakref__` (*craw.wig.WigParser attribute*), 24
- `_asdict()` (*craw.annotation.Idx method*), 20
- `_color_converter()` (*craw.heatmap.Mark method*), 30
- `_convert()` (*craw.annotation.Entry method*), 19
- `_estimate_memory()` (*craw.wig.Chromosome method*), 22
- `_extend()` (*craw.wig.Chromosome method*), 22
- `_file_readable()` (*in module craw.scripts.craw_htmp*), 36
- `_gene_size_parser()` (*in module craw.scripts.craw_htmp*), 36
- `_get_matrix_bound()` (*craw.heatmap.Mark method*), 30
- `_make()` (*craw.annotation.Idx class method*), 20
- `_replace()` (*craw.annotation.Idx method*), 20
- `_size_fig_parser()` (*in module craw.scripts.craw_htmp*), 36
- `_sort_by_gene_size()` (*in module craw.heatmap*), 30
- `_sort_using_col()` (*in module craw.heatmap*), 31
- `_sort_using_file()` (*in module craw.heatmap*), 31
- `_switch_start_stop()` (*craw.annotation.Entry method*), 19

A

- `add()` (*craw.wig.Genome method*), 23
- `AnnotationParser` (*class in craw.annotation*), 18

C

- `Chromosome` (*class in craw.wig*), 21
- `chromosome` (*craw.annotation.Entry attribute*), 19
- `Chunk` (*class in craw.wig*), 22
- `col_name` (*craw.annotation.Idx attribute*), 20
- `craw` (*module*), 17
- `craw.annotation` (*module*), 18
- `craw.argparse_util` (*module*), 34
- `craw.coverage` (*module*), 26, 27
- `craw.heatmap` (*module*), 30
- `craw.scripts.craw_coverage` (*module*), 34
- `craw.scripts.craw_htmp` (*module*), 36
- `craw.wig` (*module*), 21
- `crop_matrix()` (*in module craw.heatmap*), 31

D

- `draw_heatmap()` (*in module craw.heatmap*), 31
- `draw_one_matrix()` (*in module craw.heatmap*), 32
- `draw_raw_image()` (*in module craw.heatmap*), 32

E

- `Entry` (*class in craw.annotation*), 19

F

FixedChunk (class in *craw.wig*), 22

G

Genome (class in *craw.wig*), 22

get_annotations() (craw.annotation.AnnotationParser method), 19

get_data() (in module *craw.heatmap*), 32

get_raw_bam_coverage() (in module *craw.coverage*), 26

get_raw_coverage_function() (in module *craw.coverage*), 26

get_raw_wig_coverage() (in module *craw.coverage*), 26

get_result_header() (in module *craw.scripts.craw_coverage*), 34

get_results_file() (in module *craw.scripts.craw_coverage*), 35

get_version_message() (in module *craw*), 17

get_version_message() (in module *craw.scripts.craw_coverage*), 35

get_version_message() (in module *craw.scripts.craw_htmp*), 36

H

header (craw.annotation.Entry attribute), 19

I

Idx (class in *craw.annotation*), 20

idx (craw.annotation.Idx attribute), 20

init_logger() (in module *craw*), 17

is_comment_line() (craw.wig.WigParser static method), 24

is_data_line() (craw.wig.WigParser method), 24

is_declaration_line() (craw.wig.WigParser method), 24

is_fixed_step() (craw.wig.Chunk method), 22

is_fixed_step() (craw.wig.FixedChunk method), 22

is_fixed_step() (craw.wig.VariableChunk method), 23

is_track_line() (craw.wig.WigParser static method), 24

L

lin_norm() (in module *craw.heatmap*), 33

lin_norm_row_by_row() (in module *craw.heatmap*), 33

log_norm() (in module *craw.heatmap*), 33

log_norm_row_by_row() (in module *craw.heatmap*), 33

M

main() (in module *craw.scripts.craw_coverage*), 35

main() (in module *craw.scripts.craw_htmp*), 37

Mark (class in *craw.heatmap*), 30

max() (craw.annotation.AnnotationParser method), 19

N

new_entry_type() (in module *craw.annotation*), 20

P

padded_coverage_maker() (in module *craw.coverage*), 27

parse() (craw.wig.WigParser method), 24

parse_args() (in module *craw.scripts.craw_coverage*), 35

parse_args() (in module *craw.scripts.craw_htmp*), 37

parse_data_line() (craw.wig.Chunk method), 22

parse_data_line() (craw.wig.FixedChunk method), 22

parse_data_line() (craw.wig.VariableChunk method), 23

parse_data_line() (craw.wig.WigParser method), 25

parse_declaration_line() (craw.wig.WigParser method), 25

parse_track_line() (craw.wig.WigParser method), 25

positive_int() (in module *craw.scripts.craw_coverage*), 35

Q

quality_checker() (in module *craw.scripts.craw_coverage*), 36

R

ref (craw.annotation.Entry attribute), 19

remove_metadata() (in module *craw.heatmap*), 34

resized_coverage_maker() (in module *craw.coverage*), 28

S

sort() (in module *craw.heatmap*), 34

split_data() (in module *craw.heatmap*), 34

start (craw.annotation.Entry attribute), 19

stop (craw.annotation.Entry attribute), 19

strand (craw.annotation.Entry attribute), 19

sum_coverage_maker() (in module *craw.coverage*), 28

T

to_px() (craw.heatmap.Mark method), 30

V

VariableChunk (*class in* *craw.wig*), [23](#)

VersionAction (*class in* *craw.argparse_util*), [34](#)

W

WigError, [23](#)

WigParser (*class in* *craw.wig*), [24](#)