

CgiHunter

CAT CGTGGACG AT

User Manual

This documentation introduces the basic functions of the CgiHunterLight software tool for CpG island annotation of DNA sequences.

Table of content

1. Installation & Setup	2
2. First Steps – The Test Run.....	2
3. Command line.....	3
4. XML files.....	4

1. Installation & Setup

CgiHunterLight has been developed in the programming language python version 2.4. To use the software python 2.4 has to be installed on your computer. A python interpreter can be obtained at <http://www.python.org/download/> for Unix and Windows systems.

To install the CgiHunterLight software, download the archive and uncompress it.

Optionally, on 32-bit systems the python package 'psyco' can be installed for a performance improvement of the computation. It is available from <http://psyco.sourceforge.net/download.html>.

2. First Steps – The Test Run

The CgiHunterLight archive includes the sequences of the human chromosome 21 and 22 (assembly hg19). To generate CpG island annotations for them you can open a command line shell and change to the CgiHunterLight directory. Depending on your operating system you can start the annotation process by typing *python CgiHunterLight.py -x TestRun.xml* or *CgiHunterLight.py -x TestRun.xml*. This tells the python compiler to start the program *CgiHunterLight.py* and use the configurations in the *TestRun.xml* file to guide the annotation process. Three subfolders will be generated one for each of the chromosomes that are used to store intermediate files and the *Results* folder, which is used to store the final output.

If the program terminated successfully the *Results* folder will contain four files. *CGI_Shadowmap_hg19_chr21.bed* and *CGI_Shadowmap_hg19_chr22.bed* contain the CpG island annotations and can be directly displayed in a genome browser. *CGI_Shadowmap_hg19_chr21.txt* and *CGI_Shadowmap_hg19_chr22.txt* are files to confirm the validity of the annotation. They contain at a nucleotide resolution the information of how many genome regions overlap an individual position that fulfill the given CpG island definition.

3. Command line

Direct calls to CgiHunterLight.py are possible, but it is strongly encouraged to only use them to execute predefined XML files (e.g. `python CgiHunterLight.py -x <name of xml file>`). If direct configuration of CgiHunter via command line is demanded, use the `-h` or `--help` option to find the following overview of all command line options.

Option	Type	Description
<code>-x</code> or <code>--xml</code>	string	The name of a xml file specifying a CGIH_Task. Overrides all other parameters
<code>-X</code> or <code>--tasklist</code>	string	Relative or absolute path of a task list. Each line of this file should contain the path to a xml file that will be executed as if called by <code>-x</code> .
<code>-F</code> or <code>--force</code>	flag	The execution of an xml file will be forced, even if it is marked as finished
<code>-a</code> or <code>--assembly</code>	string	The genome assembly/species the DNA sequence belongs to
<code>-c</code> or <code>--chrom</code>	string	The chromosome the DNA sequence belongs to
<code>-d</code>	int	Debug level – the higher the number the more status messages will be shown
<code>-f</code> or <code>--comb</code>	string	Filename of the prefilterstep. This option can be used, if the prefilterstep has already been performed
<code>-m</code> or <code>--mask</code>	string	If set to True, lower cased letters in the DNA sequence are handled as repeats and are always considered as non-G and non-C. Default is False.
<code>-o</code> or <code>--offset</code>	int	This offset is added to all coordinates in the outputfiles. It should be used if fasta files are applied which start at the beginning of the chromosome.
<code>-p</code> or <code>--path</code>	string	The path to the directory where all output will be placed
<code>-s</code> or <code>--source</code>	string	Absolute or relative path of the fasta file that contains the DNA sequence
<code>--gc</code>	float	Minimal content of cytosine and guanine a region have to possess to qualify as a CpG island Only values between 50 and 100 are valid. Example: <code>--gc 50</code>
<code>-r</code> or <code>--ratio</code>	float	Minimal ratio of observed CpGs over expected CpGs a region have to possess to qualify as a CpG island. Only values between 50 and 100 are valid. Example: <code>-r 60</code>
<code>--cw</code> or <code>--combwidth</code>	float	Combwidth for the combing/filter step. Must be greater than 1. Default 1.2.
<code>-w</code> or <code>--min_win</code>	int	Minimal length threshold of CpG island
<code>-W</code> or <code>--max_win</code>	int	Optional upper bound to CpG island length. Per default set to a million base-pairs
<code>-g</code> or <code>--getxml</code>	string	If configuration is performed by command line this option can be used to retrieve the xml file that contains all options.

4. XML files

Each CgiHunter analyses can be completely specified as XML document. These documents can be either generated by using the CgiHunterGUI, by modifying the template document included in the CgiHunter distribution, by using the `-g` option of the command line or by your own custom techniques.

Each XML file has a status option that is initially set to 'unfished'. Upon completion of a CgiHunter call this mark will be set to 'finished'. CgiHunter will automatically ignore tasks that are marked as 'finished' to ensure that already processed files are not reanalyzed. To reset the status you can either load the XML file in the CgiHunterGUI and save it again under the same name or directly modify the status in the file.

If the computation of a file is interrupted, the status will contain the respective error message.

XML files can be executed directly from the command line (e.g. `python CgiHunter.py -x <name of xml file>`).