

# HTSeq - a Python framework to work with high-throughput sequencing data

Dragoş Alin Rotaru

University of Bucharest

24 April, 2015

# Overview for section 1

- 1 Introduction
- 2 Trying to solve HTSeq-data problems
- 3 Prerequisites
- 4 Investigations
- 5 misc

- Low cost DNA sequencing

- Low cost DNA sequencing
- Technologies that can take a snapshot of RNA at a moment (RNA-Seq)

- Low cost DNA sequencing
- Technologies that can take a snapshot of RNA at a moment (RNA-Seq)
- Size of human genome is about 3,234.83 Mb (Mega-basepairs)

- Low cost DNA sequencing
- Technologies that can take a snapshot of RNA at a moment (RNA-Seq)
- Size of human genome is about 3,234.83 Mb (Mega-basepairs)
- Dealing with HTS(High Throughput Sequencing) data is inevitably

# Overview for section 2

- 1 Introduction
- 2 Trying to solve HTSeq-data problems
- 3 Prerequisites
- 4 Investigations
- 5 misc

# Current problems

- Most tools are built only for specific purposes



# Current problems

- Most tools are built only for specific purposes
- Part of them: not open-source

- Most tools are built only for specific purposes
- Part of them: not open-source
- Part of the others: not suitable for high performance tasks

# What is HTSeq?

- Framework for Python
- Written in Cython
- Great and easy for writing customizable scripts
- High level scripts
- Well documented

# What is HTSeq?

HTSeq 0.6.1p2 documentation »

previous | next | index

Table Of Contents

HTSeq: Analysing high-throughput sequencing data with Python

- Paper
- Documentation overview
- Author
- License

Previous topic

HTSeq: Analysing high-throughput sequencing data with Python

Next topic

Prerequisites and installation

This Page

Show Source

Quick search

Go

Enter search terms or a module, class or function name.

HTSeq: Analysing high-throughput sequencing data with Python

HTSeq is a Python package that provides infrastructure to process data from high-throughput sequencing assays.

- Please see the chapter *A tour through HTSeq* first for an overview on the kind of analysis you can do with HTSeq and the design of the package, and then look at the reference documentation.
- While the main purpose of HTSeq is to allow you to write your own analysis scripts, customized to your needs, there are also a couple of stand-alone scripts for common tasks that can be used without any Python knowledge. See the *Scripts* section in the overview below for what is available.
- For downloads and installation instructions, see *Prerequisites and installation*.

Paper

HTSeq is described in the following publication:

Simon Anders, Paul Theodor Pyl, Wolfgang Huber  
*HTSeq — A Python framework to work with high-throughput sequencing data*  
Bioinformatics (2014), in print, online at doi:10.1093/bioinformatics/btu638

If you use HTSeq in research, please cite this paper in your publication.

Documentation overview

- Prerequisites and installation*

Download links and installation instructions can be found here

- A tour through HTSeq*

The Tour shows you how to get started. It explains how to install HTSeq, and then demonstrates typical analysis steps with explicit examples. Read this first, and then see the Reference for details.

- A detailed use case: TSS plots*

This chapter explains typical usage patterns for HTSeq by explaining in detail three different solutions to the same ~~problematic task~~

## HTSeq - Documentation Website

Dragoş Alin Rotaru (UniBuc)

University of Bucharest

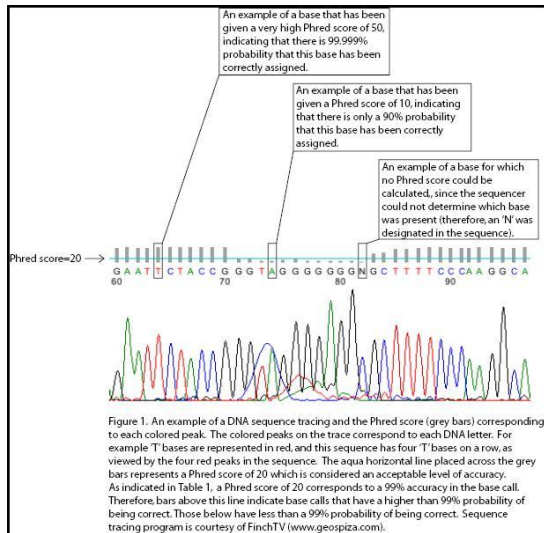
24 April, 2015

6 / 25

# Overview for section 3

- 1 Introduction
- 2 Trying to solve HTSeq-data problems
- 3 Prerequisites**
- 4 Investigations
- 5 misc

# Back to biology... - Phred Quality Scores



Phred Quality Scores - Source: Wiki

# File formats?

- FASTA files - representing nucleotide (A-C-G-T) or peptide (aminoacid) sequences
- FASTQ files - usually nucleotide sequences, this time with their scores

# File formats?

```
@HWI-EAS209_0006_FC706VJ:5:58:5894:21141#ATCACG/1
TTAATTGGTAAATAAATCTCCTAATAGCTTAGATNTTACCTTNNNNNNNNNNNTAGTTTCTTC
+HWI-EAS209_0006_FC706VJ:5:58:5894:21141#ATCACG/1
efcfffffcfeefffcfffffddf'feed]'_Ba_^__[YBBBBBBBBBBRTT\]][]dc
```



- SAM - Sequence Alignments Map

- SAM - Sequence Alignments Map
- SAM can be derived from FASTQ files by DNA aligning reads with some reference genome

- BAM

# More file formats

- BAM
- Binary version of SAM

- BAM
- Binary version of SAM
- There are different tools which does the conversion between SAM and BAM or others (samtools)

# Now HTSeq; SAM - Aligned; FASTQ - Phred score

- Fortunately, HTSeq have built-in parsers for those files!

# Now HTSeq; SAM - Aligned; FASTQ - Phred score

- Fortunately, HTSeq have built-in parsers for those files!
- Can perform quality assessments on SAM and FASTQ directly with `htseq-qc` command

# Now HTSeq - counting

- Given a list of features, genes, exons etc, how many reads map to them?



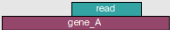
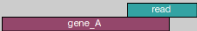



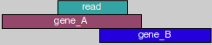
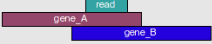
# Now HTSeq - counting

- Given a list of features, genes, exons etc, how many reads map to them?
- Can use prebuilt htseq-count command or custom scripts

# Now HTSeq - counting

- Given a list of features, genes, exons etc, how many reads map to them?
- Can use prebuilt htseq-count command or custom scripts
- Advanced data structures already implemented such as GenomicInterval or GenomicArrayOfSets

# Now HTSeq - counting

	union	intersection_strict	intersection_nonempty
	gene_A	gene_A	gene_A
	gene_A	no_feature	gene_A
	gene_A	no_feature	gene_A
	gene_A	gene_A	gene_A
	gene_A	gene_A	gene_A
	ambiguous	gene_A	gene_A
	ambiguous	ambiguous	ambiguous

HTSeq count modes - Source: HTSeq WebSite

# Now HTSeq - transcription start sites

- Compute a 'window' profile from an interval with respect to some features

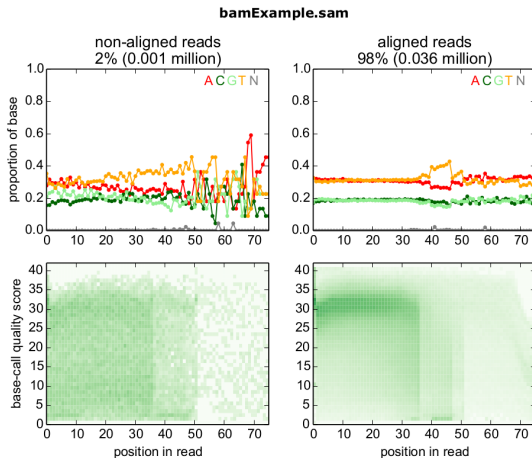
# Overview for section 4

- 1 Introduction
- 2 Trying to solve HTSeq-data problems
- 3 Prerequisites
- 4 Investigations**
- 5 misc

# Playing with HTSeq and different files

- Realize that in practice you can't find real SAM files
- BAM files have at least 2GB size
- Samtools doesn't work very well when it comes to converting BAM to SAM
- Managed to get a QA plot

# Playing with HTSeq and different files



Toy Example

Live demonstration



- HTSeq is a strong tool for manipulating High-Throughput data

# Conclusions

- HTSeq is a strong tool for manipulating High-Throughput data
- Often, theory differs from practice

# Thank you!