

Jumpstart your genomics pipelines with genomepy

This manuscript ([permalink](#)) was automatically generated from [vanheeringen-lab/genomepy_manuscript@63dc557](#) on March 11, 2021.

Authors

- **Siebre Frolich**

 [0000-0001-6925-8446](#) ·  [siebrenf](#)

Department of Molecular Developmental Biology, Radboud University

- **Maarten van der Sande**

 [0000-0001-7803-1526](#) ·  [Maarten-vd-Sande](#) ·  [MaartenvdSande](#)

Department of Molecular Developmental Biology, Radboud University

- **Simon van Heeringen**

 [0000-0002-0411-3219](#) ·  [simonvh](#) ·  [svheeringen](#)

Department of Molecular Developmental Biology, Radboud University

Abstract

Summary

Analyzing genomics data, such as ATAC, ChIP and RNA-sequencing, requires multiple types of genomic data, such as genome sequence and gene annotations. These resources can generally be retrieved from multiple organizations, where they exist at multiple versions, and may have been generated with varying methods. Which set of genomic data to use depends on the context of the research, such as collaboration partners, data reuse or the quality of the genomic data. Many of the bioinformatic workflows and pipelines available to date require the user to make this informed decision and supply the genomic data manually. Obtaining this data can be a tedious and error-prone process and does not allow for full computational reproducibility.

Here we present genomepy, a quality-of-life enhancement tool, that can navigate the genome databases of Ensembl, UCSC and NCBI. Genomepy can search and install genome sequences and gene annotation data from these providers in a consistent, reproducible and documented manner. The search function retrieves genomes related to the search term, and can do so for one or all providers to allow the user to make an informed decision. The install function downloads a specified genome with sensible defaults, while providing full control to advanced features. Additionally, gene annotations can be downloaded and converted to commonly used formats, with built-in checks for compatibility with the genome. Genomepy can optionally create genome indexes for commonly used aligners, including splice-aware aligners utilizing both genome and gene annotations. Genomes and gene annotations not available on supported databases can be processed by genomepy as well, providing a consistent workflow with any genome.

Genomepy provides this functionality and more via command line interface and Python application programming interface, aimed at easy of use and integration in automated pipelines.

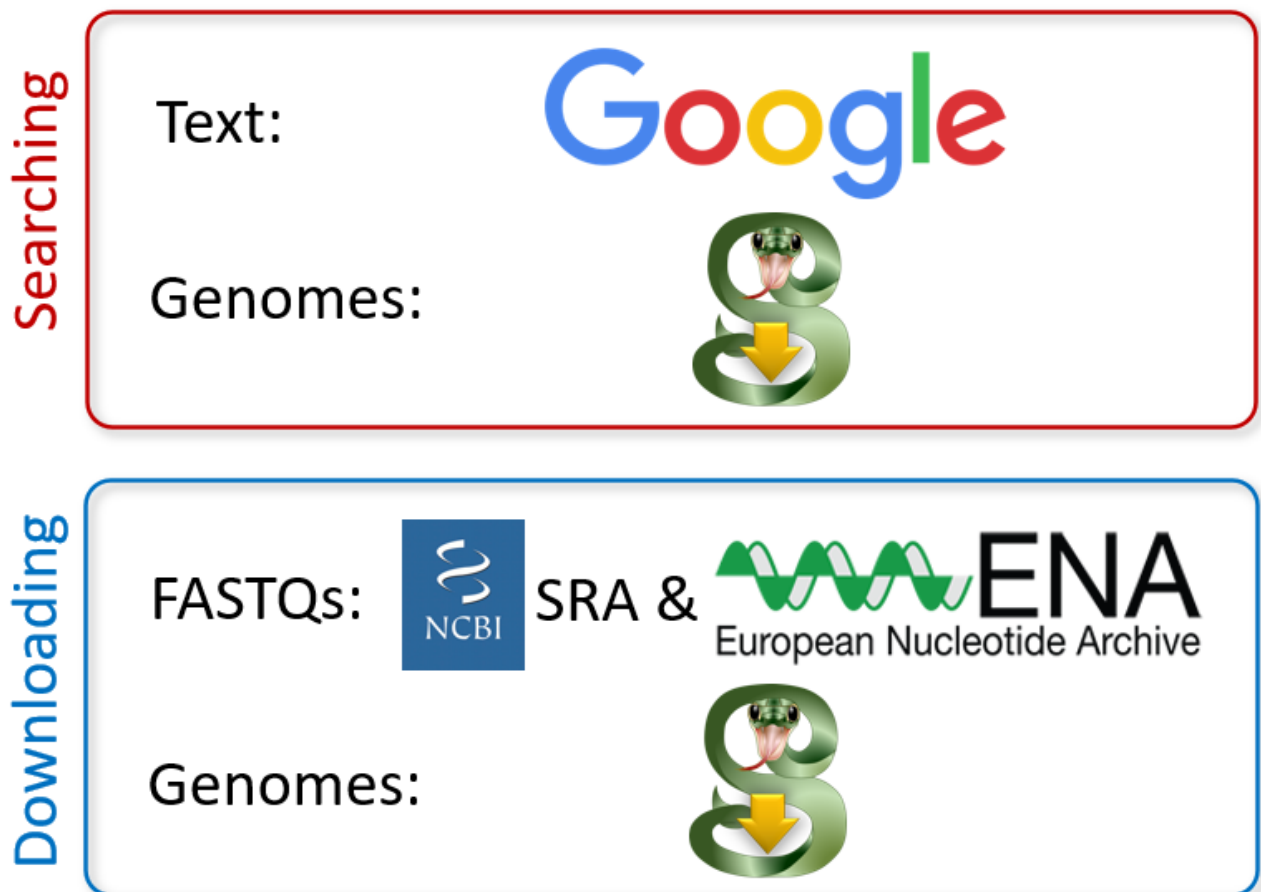


Figure 1: genomepy in a nutshell

Availability and implementation

Genomepy can be installed using [Bioconda](#) and [Pip](#), and the code is available at <https://github.com/vanheeringen-lab/genomepy>.

Introduction

As high-throughput sequencing matured over the past decade and a half, the size and amount of sequencing data has exploded. Over the past five years, the number of datasets published on the NCBI GEO database increased by an average of 2000 per year, while the number of samples increased by 100.000 per year [1](#). This explosion of data highlights the need for scalable, robust and automatable methods for data processing and analysis.

A large amount of bioinformatics tools have been created to facilitate this process, including (pre)processing tools, analysis tools and workflow managers to link these tools together (such as the snakemake wrapper [2](#)). Together, this allows the human component to shift towards management, reducing both human labor and error. However, not all input data can be obtained fully automatic. Notably the genomic data, which includes the genome assembly, gene annotation and derived files. These data can be obtained from a variety of different providers. These include three major genome providers, Ensembl [3](#), UCSC [4](#) and NCBI [5](#), and many niche providers, such as flybase [6](#), wormbase [7](#) or xenbase [8](#). Each provider has a different method of generating genome assemblies and gene annotations, which can affect available data formats, naming schema, information density, as well as availability, accessibility and relevance. Take the major providers: Ensembl hosts a large database of

curated genomes, most of which are actively maintained. They provide gene annotation in commonly used data formats (GTF, GFF3 and BED12), and their GTF files contain the most metadata of the major providers. However, Ensembl uses a chromosome naming scheme that clashes with several bioinformatics tools (e.g. "1" for chromosome 1). NCBI hosts the largest database by accepting submissions from the Genome Reference Consortium as well as from individuals. As a result, the database provides a trove of data with varying quality and completeness. Finally, UCSC hosts and maintains a modest set of reference genomes, with gene annotations derived from multiple pipelines for most. However, their gene annotations have inconsistent data formats, which also clashes with bioinformatics tools (e.g. the gene names and gene identifiers are mixed). These differences impact the compatibility of the reference data with research tools [9](#), other reference databases and other research. Therefore, the choice of provider and reference data is of significant importance.

Making an informed decision requires an overview of available options. This is achieved by determining which providers host the genomic data of interest, then which genomic data is most suitable, then downloading and processing the genomic data and finally starting your analysis. Performing these steps manually creates room for human error in the locating, processing and reporting of these steps as well as the reasoning behind them. For the sake of sanity and reproducibility, it would be better if this could be done in a standardized system.

In order to assist in searching through genome providers for, and the processing of genomic data, we developed genomepy. Using the search function one can get an overview of all genomes containing the search term in their name, description or accession identifier, as well as all genomes matching a taxonomy identifier. The search function can provide this overview for each of the major providers, Ensembl, UCSC, NCBI or all three. Once a selection is made the genome and gene annotations can be downloaded and prepared for use with the install function (figure [2](#)). The install function automatically performs processing of the genomic data, including genome indexing with pyfaidx [10](#), generating support files (chromosome sizes and sequencing gaps), matching chromosome names between genome and gene annotation, and optional aligner index generation). Using the command line interface will automatically log which of these features is used. Once downloaded genomepy installs the path to genome FASTA file into command line variables, for quick reference in downstream tools. Using the Python application programming interface the genomepy Genome object can be used to similarly access the metadata of the installed genomic data, such as the location and content of additional data files (chromosome sizes, gene annotation files). Both interfaces are developed for incorporation in workflows to automate the downloading and processing steps.

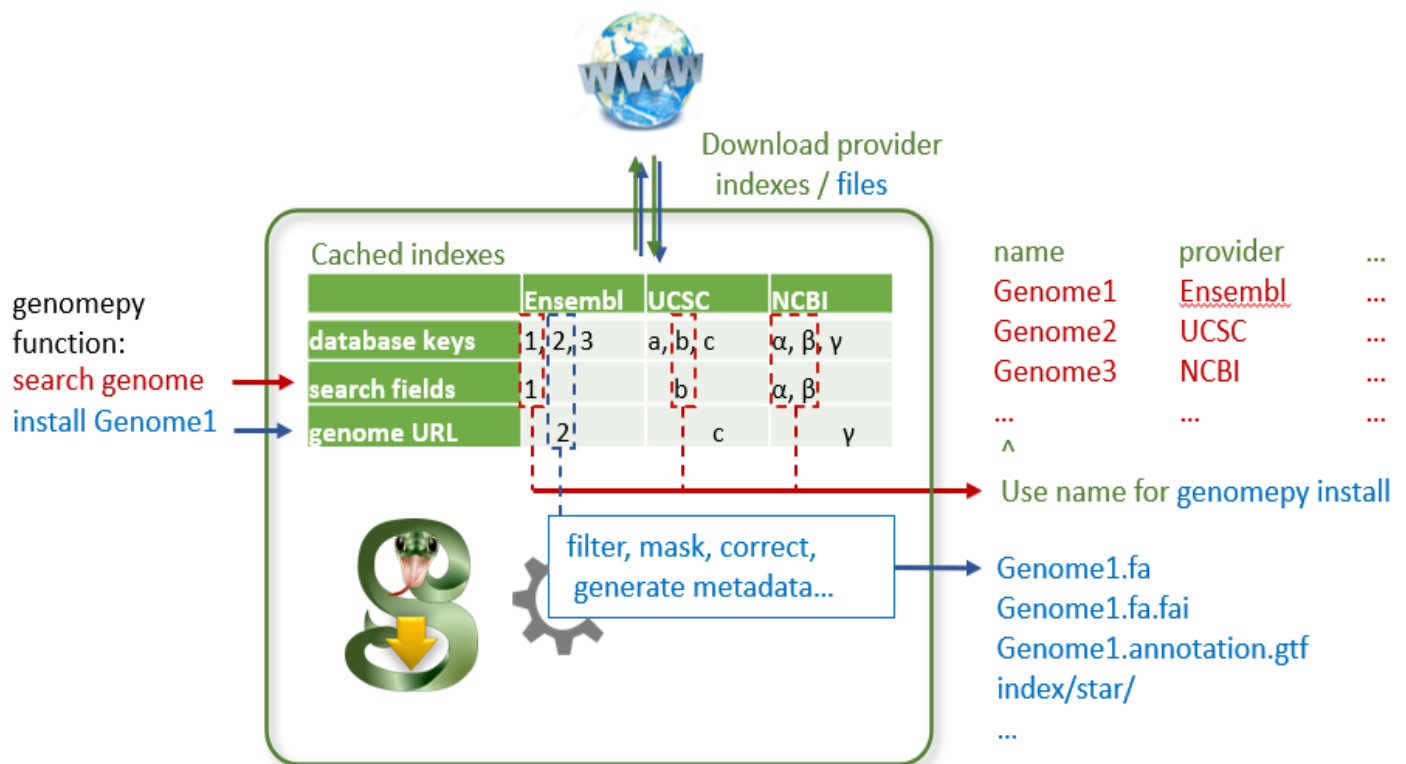


Figure 2: genomepy workflow

```
user@comp:~$ genomepy install dm3 --annotation --genomes_dir ./genomes

user@comp:~$ ls ./genomes/dm3 -l
README.txt
dm3.annotation.bed.gz
dm3.annotation.gtf.gz
dm3.blacklist.bed
dm3.fa
dm3.fa.fai
dm3.fa.sizes
dm3.gaps.bed
index
```

Figure 3: executive overview

Related Work

Ensembl, UCSC and NCBI all support downloading from their individual databases via accessible FTP archives, web portals, and REST APIs. To access these databases programmatically, there exists several external tools, such as the ncbi-genome-download tool [11](#) and ucsc-genomes-downloader [12](#). However, to our knowledge no tool exists that can consistently search or download from all three major genome providers.

There are several existing tools for reproducibly sharing reference data between projects. These data management tools accept reference data and derived assets such as aligner indexes from any source, such as iGenomes [13](#), refGenie [14](#) and Go Get Data [15](#). These tools excel in their ability to reproducibly share data, a feature which is not present in genomepy, and can be used to obtain and manage previously generated data with ease. However, these tools require the user to supply the reference data to any new assembly (such as non-model organisms), new assembly version (such as the latest patch of the human genome) or in some cases assets (such as an aligner index not present in the hosted data).

We conclude that there is a need for a tool that can assist in obtaining and preparing genomic data for downstream analysis. This includes providing an overview of available choices, downloading the data, error checking (and if possible correcting) genome and gene annotations, and outputting data types commonly used in bioinformatics tools. Genomepy was created to fit this need for both automated and human-operated workflows. Additionally, genomepy provides a framework that facilitates the downstream use of the data.

genomepy

The core functionality of genomepy is to search, download and prepare genomes and gene annotations. These functions are split over two command line functions: `genomepy search` and `genomepy install`.

Search

Search will query the three major providers for a given search term, with the option query only one specified provider. Genomepy can search for text terms present in genome names or descriptions, taxonomy identifiers and accession numbers, and will automatically detect which is used in the search term. The search results are returned with available metadata for review.

Install

When an assembly has been selected, the name (as returned from the search function) can be passed to the install function. The files will be downloaded from the major provider that hosts it. Users can optionally specify a provider if the data is available on multiple. Alternatively, if the assembly originates from another source, the url to the genome can be passed. Next, the genome assembly is downloaded with the desired sequence masking level [[16](#),[17](#)]. By default soft masked genomes (repetitive sequences written in lower case) are downloaded, but unmasked or hard masked (repetitive sequences written as Ns) can be obtained as well. If the provider does not have the genome at this masking level, genomepy will edit the FASTA to match.

Reference assemblies often contain alternate sequences to reflect biological diversity. For the purpose of sequence alignment however, the best results are obtained if there is one reference per nucleotide. Therefore genomepy filters out alternative regions, unless specified otherwise. Additionally, regex filters may be passed to either include or exclude contigs (chromosomes, scaffolds, etc.) by name. For instance to filter out (or filter for) chromosomes, unplaced-, unknown- or mitochondrial sequences.

Once processing is performed, genomepy generates several commonly used support files. The genome is indexed using pyfaidx [10](#), and contig sizes and contig gap sizes are stored in separate files.

If specified, genomepy will attempt to download a gene annotation: genomepy will search the database for a GFF, GTF, BED or (for UCSC only) text format gene annotation. The annotation is then

processed to the commonly used GTF and BED output formats using publicly available conversion tools [18](#). Afterward, contig names between the genome and gene annotation are checked for compatibility. Should these mismatch, genomepy will attempt to match the names in the annotations to those in the genome.

For data provenance and reproducibility, a README file is kept with the output which logs the URLs to the source files, the steps performed, and contigs filtered out.

Plugins

Genomepy facilitates several optional processing steps via plugins. Using the command line, these options can be inspected and toggled with the `genomepy plugin` function. Enabled plugins will execute upon (re)running the install function.

The blacklist plugin downloads blacklists by the Kundaje lab [19](#) for the supported UCSC genomes (and GRCh38). If no blacklist is available the program will proceed after giving a warning.

Currently, other plugins support the generation of several popular genome aligners. These include DNA aligner indexes for Bowtie2 [20](#), BWA [21](#), GMAP [22](#) or Minimap2 [23](#), and splice-aware aligners such as STAR [24](#) and HISAT2 [25](#). The splice-aware aligners function both with and without gene annotations. Should no gene annotation be downloaded, the program will issue a warning and proceed to generate a splice-unaware index.

Conclusion

Obtaining suitable genomic data is a principal step in any genomics project. Genomepy was developed to provide a consistent overview of genomes on the three major providers, and reproducibly download and process genomic data ready for downstream use. Using genomepy, a project can utilize genomes from any provider and expect consistent output with data from a major provider or otherwise. This allows genomepy to automate a step in genomic data preprocessing that was performed by hand prior. Additionally, it facilitates downstream analysis, by setting up paths to the genome FASTA file, providing genome metadata within the Python Genome object and generating support files. Combined, these features make genomepy ideal for integration in automated sequencing workflows, as demonstrated in seq2science [26](#), paving the way for robust and reproducible analysis.

Future prospects

Further improvements to openness and reproducibility could be achieved by combining or coupling of genomepy with data management tools such as Go Get Data and RefGenie. Additionally, further improvements could be made in the ability to specify the assembly version, in order to allow annotations of different versions with an assembly.

Acknowledgements

We thank the Department of Molecular (Developmental) Biology, our github [contributors](#), and issue posters for their patience, feedback and insight. We thank black, pytest, CodeCoverage and TravisCI for enduring our abuse and teaching us patience. And finally, we thank Manubot [27](#) for assisting with this manuscript.

Code availability

Genomepy can be installed using [Bioconda](#) and [Pip](#). The code is available at <https://github.com/vanheeringen-lab/genomepy>.

Supplementary Information

The full genomepy documentation including examples can be viewed [here](#)

References

1. NCBI GEO: archive for functional genomics data sets-update

Tanya Barrett, Stephen E Wilhite, Pierre Ledoux, Carlos Evangelista, Irene F Kim, Maxim Tomashevsky, Kimberly A Marshall, Katherine H Phillippy, Patti M Sherman, Michelle Holko, ... Alexandra Soboleva

Nucleic acids research (2013-01) <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3531084/>
DOI: [10.1093/nar/gks1193](https://doi.org/10.1093/nar/gks1193) · PMID: [23193258](https://pubmed.ncbi.nlm.nih.gov/23193258/) · PMCID: [PMC3531084](https://pubmed.ncbi.nlm.nih.gov/PMC3531084/)

2. GENOMEPTY — Snakemake Wrappers tags/0.72.0 documentation <https://snakemake-wrappers.readthedocs.io/en/stable/wrappers/genomepty.html>

3. Ensembl 2020

Andrew D. Yates, Premanand Achuthan, Wasiu Akanni, James Allen, Jamie Allen, Jorge Alvarez-Jarreta, M. Ridwan Amode, Irina M. Armean, Andrey G. Azov, Ruth Bennett, ... Paul Flicek

Nucleic Acids Research (2020-01-08) <https://academic.oup.com/nar/article/48/D1/D682/5613682>
DOI: [10.1093/nar/gkz966](https://doi.org/10.1093/nar/gkz966)

4. The Human Genome Browser at UCSC

W. J. Kent, C. W. Sugnet, T. S. Furey, K. M. Roskin, T. H. Pringle, A. M. Zahler, a. D. Haussler

Genome Research (2002-05-16) <https://doi.org/fpf5rm>
DOI: [10.1101/gr.229102](https://doi.org/10.1101/gr.229102) · PMID: [12045153](https://pubmed.ncbi.nlm.nih.gov/12045153/) · PMCID: [PMC186604](https://pubmed.ncbi.nlm.nih.gov/PMC186604/)

5. The UCSC Table Browser data retrieval tool

Donna Karolchik, Angela S Hinrichs, Terrence S Furey, Krishna M Roskin, Charles W Sugnet, David Haussler, W James Kent

Nucleic acids research (2004-01-01) <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC308837/>
DOI: [10.1093/nar/gkh103](https://doi.org/10.1093/nar/gkh103) · PMID: [14681465](https://pubmed.ncbi.nlm.nih.gov/14681465/) · PMCID: [PMC308837](https://pubmed.ncbi.nlm.nih.gov/PMC308837/)

6. FlyBase 2.0: the next generation

Jim Thurmond, Joshua L. Goodman, Victor B. Strelets, Helen Attrill, L. Sian Gramates, Steven J.

Marygold, Beverley B. Matthews, Gillian Millburn, Giulia Antonazzo, Vitor Trovisco, ... Phillip Baker

Nucleic Acids Research (2019-01-08) <https://academic.oup.com/nar/article/47/D1/D759/5144957>
DOI: [10.1093/nar/gky1003](https://doi.org/10.1093/nar/gky1003)

7. WormBase: a modern Model Organism Information Resource

Todd W. Harris, Valerio Arnaboldi, Scott Cain, Juancarlos Chan, Wen J. Chen, Jaehyoung Cho, Paul Davis, Sibyl Gao, Christian A. Grove, Ranjana Kishore, ... Paul W. Sternberg

Nucleic Acids Research (2020-01-08) <https://academic.oup.com/nar/article/48/D1/D762/5603222>
DOI: [10.1093/nar/gkz920](https://doi.org/10.1093/nar/gkz920)

8. Xenbase: a genomic, epigenomic and transcriptomic model organism database

Kamran Karimi, Joshua D Fortriede, Vaneet S Lotay, Kevin A Burns, Dong Zhou Wang, Malcom E Fisher, Troy J Pells, Christina James-Zorn, Ying Wang, V G Ponferrada, ... Peter D Vize

Nucleic Acids Research (2018-01-04) <https://doi.org/ghk99d>
DOI: [10.1093/nar/gkx936](https://doi.org/10.1093/nar/gkx936) · PMID: [29059324](https://pubmed.ncbi.nlm.nih.gov/29059324/) · PMCID: [PMC5753396](https://pubmed.ncbi.nlm.nih.gov/PMC5753396/)

9. A comprehensive evaluation of ensembl, RefSeq, and UCSC annotations in the context of RNA-seq read mapping and gene quantification

Shanrong Zhao, Baohong Zhang

BMC Genomics (2015-02-18) <https://doi.org/10.1186/s12864-015-1308-8>
DOI: [10.1186/s12864-015-1308-8](https://doi.org/10.1186/s12864-015-1308-8)

10. **Efficient “pythonic” access to FASTA files using pyfaidx**
Matthew D. Shirley, Zhaorong Ma, Brent S. Pedersen, Sarah J. Wheelan
PeerJ PrePrints (2015-04-08) <https://peerj.com/preprints/970>
11. **kblin/ncbi-genome-download**
Kai Blin
(2021-03-05) <https://github.com/kblin/ncbi-genome-download>
12. **ucsc-genomes-downloader: Python package to quickly download genomes from the UCSC.**
Luca Cappelletti
https://github.com/LucaCappelletti94/ucsc_genomes_downloader
13. **iGenomes** https://support.illumina.com/sequencing/sequencing_software/igenome.html
14. **Refgenie: a reference genome resource manager**
Michał Stolarczyk, Vincent P. Reuter, Jason P. Smith, Neal E. Magee, Nathan C. Sheffield
GigaScience (2020-02-01) <https://academic.oup.com/gigascience/article/9/2/giz149/5717403>
DOI: [10.1093/gigascience/giz149](https://doi.org/10.1093/gigascience/giz149)
15. **Go Get Data (GGD): simple, reproducible access to scientific data**
Michael J. Cormier, Jonathan R. Belyeu, Brent S. Pedersen, Joseph Brown, Johannes Koster, Aaron R. Quinlan
bioRxiv (2020-09-11) <https://www.biorxiv.org/content/10.1101/2020.09.10.291377v2>
DOI: [10.1101/2020.09.10.291377](https://doi.org/10.1101/2020.09.10.291377)
16. **RepeatMasker Home Page** <http://repeatmasker.org/>
17. **A Fast and Symmetric DUST Implementation to Mask Low-Complexity DNA Sequences**
Aleksandr Morgulis, E. Michael Gertz, Alejandro A. Schäffer, Richa Agarwala
Journal of Computational Biology (2006-06) <https://doi.org/fqs85g>
DOI: [10.1089/cmb.2006.13.1028](https://doi.org/10.1089/cmb.2006.13.1028) · PMID: [16796549](https://pubmed.ncbi.nlm.nih.gov/16796549/)
18. **Index of /admin/exe** <http://hgdownload.cse.ucsc.edu/admin/exe/>
19. **The ENCODE Blacklist: Identification of Problematic Regions of the Genome**
Haley M. Amemiya, Anshul Kundaje, Alan P. Boyle
Scientific Reports (2019-06-27) <https://www.nature.com/articles/s41598-019-45839-z>
DOI: [10.1038/s41598-019-45839-z](https://doi.org/10.1038/s41598-019-45839-z)
20. **Fast gapped-read alignment with Bowtie 2**
Ben Langmead, Steven L. Salzberg
Nature Methods (2012-04) <https://www.nature.com/articles/nmeth.1923>
DOI: [10.1038/nmeth.1923](https://doi.org/10.1038/nmeth.1923)
21. **Fast and accurate short read alignment with Burrows–Wheeler transform**
Heng Li, Richard Durbin
Bioinformatics (2009-07-15) <https://academic.oup.com/bioinformatics/article/25/14/1754/225615>
DOI: [10.1093/bioinformatics/btp324](https://doi.org/10.1093/bioinformatics/btp324)

22. **GMAP: a genomic mapping and alignment program for mRNA and EST sequences**
Thomas D. Wu, Colin K. Watanabe
Bioinformatics (2005-05-01) <https://academic.oup.com/bioinformatics/article/21/9/1859/409207>
DOI: [10.1093/bioinformatics/bti310](https://doi.org/10.1093/bioinformatics/bti310)
23. **Minimap2: pairwise alignment for nucleotide sequences**
Heng Li
Bioinformatics (2018-09-15) <https://academic.oup.com/bioinformatics/article/34/18/3094/4994778>
DOI: [10.1093/bioinformatics/bty191](https://doi.org/10.1093/bioinformatics/bty191)
24. **STAR: ultrafast universal RNA-seq aligner**
Alexander Dobin, Carrie A. Davis, Felix Schlesinger, Jorg Drenkow, Chris Zaleski, Sonali Jha, Philippe Batut, Mark Chaisson, Thomas R. Gingeras
Bioinformatics (2013-01-01) <https://academic.oup.com/bioinformatics/article/29/1/15/272537>
DOI: [10.1093/bioinformatics/bts635](https://doi.org/10.1093/bioinformatics/bts635)
25. **HISAT: a fast spliced aligner with low memory requirements**
Daehwan Kim, Ben Langmead, Steven L. Salzberg
Nature Methods (2015-04) <https://www.nature.com/articles/nmeth.3317>
DOI: [10.1038/nmeth.3317](https://doi.org/10.1038/nmeth.3317)
26. **seq2science**
Maarten Van Der Sande, Siebren Frölich, Jos Smits, Tilman Schäfers, Rebecca Snabel, Simon Van Heeringen
Zenodo (2021-03-03) <https://doi.org/ghktgZ>
DOI: [10.5281/zenodo.3921913](https://doi.org/10.5281/zenodo.3921913)
27. **Open collaborative writing with Manubot**
Daniel S. Himmelstein, Vincent Rubinetti, David R. Slochower, Dongbo Hu, Venkat S. Malladi, Casey S. Greene, Anthony Gitter
PLOS Computational Biology (2019-06-24) <https://doi.org/c7np>
DOI: [10.1371/journal.pcbi.1007128](https://doi.org/10.1371/journal.pcbi.1007128) · PMID: [31233491](https://pubmed.ncbi.nlm.nih.gov/31233491/) · PMCID: [PMC6611653](https://pubmed.ncbi.nlm.nih.gov/PMC6611653/)