

SNAPPy: A snakemake pipeline for scalable HIV-1 subtyping by phylogenetic pairing

Pedro M.M. Araújo,^{1,2} Joana S. Martins,^{1,2} and Nuno S. Osório^{1,2,*},†

¹Life and Health Sciences Research institute (ICVS), School of Medicine, University of Minho, Braga, Portugal and ²ICVS/3B's - PT Government Associate Laboratory, Braga, Guimarães, Portugal

*Corresponding author: E-mail: nosorio@med.uminho.pt

† <http://orcid.org/0000-0003-0949-5399>

Abstract

Human immunodeficiency virus 1 (HIV-1) genome sequencing is routinely done for drug resistance monitoring in hospitals worldwide. Subtyping these extensive datasets of HIV-1 sequences is a critical first step in molecular epidemiology and evolution studies. The clinical relevance of HIV-1 subtypes is increasingly recognized. Several studies suggest subtype-related differences in disease progression, transmission route efficiency, immune evasion, and even therapeutic outcomes. HIV-1 subtyping is mainly done using web-servers. These tools have limitations in scalability and potential noncompliance with data protection legislation. Thus, the aim of this work was to develop an efficient method for large-scale local HIV-1 subtyping. We designed SNAPPy: a snakemake pipeline for scalable HIV-1 subtyping by phylogenetic pairing. It contains several tasks of phylogenetic inference and BLAST queries, which can be executed sequentially or in parallel, taking advantage of multiple-core processing units. Although it was built for subtyping, SNAPPy is also useful to perform extensive HIV-1 alignments. This tool facilitates large-scale sequence-based HIV-1 research by providing a local, resource efficient and scalable alternative for HIV-1 subtyping. It is capable of analyzing full-length genomes or partial HIV-1 genomic regions (GAG, POL, and ENV) and recognizes more than ninety circulating recombinant forms. SNAPPy is freely available at: <https://github.com/PMMARAujo/snappy/releases>.

Key words: HIV-1; subtyping; genetic diversity; scalability; phylogeny.

1. Introduction

The number of HIV-1 partial or complete genomic sequences in databases largely increased along the years after a noteworthy surge of almost tenfold in the 2000s. HIV-1 genomic data are extremely valuable for fundamental research, translating into several epidemiological applications such as antiretroviral resistance surveillance or transmission history reconstruction (Abecasis et al. 2013; Yebra et al. 2015; Araújo et al. 2019). Subtyping is a primary analysis done on HIV-1 sequences to allow further investigation.

HIV-1 was consensually divided in four groups (M, N, O, and P), a consequence of multiple cross-species transmission events from non-human primates to humans. M group is the only with

a worldwide dispersion, and due to genetic differences and divergent evolutionary stories viruses from this group were divided into nine Subtypes (A to D, F to H, J, and K) and Sub-subtypes (e.g. A1, A2, F1, and F2). HIV-1 genomes composed of parts of different subtypes are known as recombinant forms, which can be named circulating recombinant forms (CRFs) if several cases are detected or unique recombinant forms (URFs) for more sporadic cases (Robertson et al. 2000; Hemelaar 2013). It has been reported that different HIV-1 subtypes may be better adapted for specific transmission routes (Renjifo et al. 2004; John-Stewart et al. 2005), contain higher prevalence of polymorphisms known to influence immune systems (Bartolo et al. 2011; Serwanga et al. 2015) or antiretroviral treatment evasion (Brenner et al. 2003; Abecasis et al. 2005; Camacho and

Vandamme 2007), and lead to differences in the disease progression rate (Baeten et al. 2007; Kiwanuka et al. 2008; Easterbrook et al. 2010; Araujo et al. 2019).

There are three main classes of approaches to perform HIV-1 subtyping: similarity-based (e.g. Stanford (Liu and Shafer 2006) and NCBI subtyping tool (Rozañov et al. 2004)); statistical-based (e.g. COMET (Struck et al. 2014), jpHMM (Schultz et al. 2009), and STAR (Myers et al. 2005)); and phylogenetic-based (e.g. REGA (Pineda-Peña et al. 2013) and SCUEAL (Kosakovsky Pond et al. 2009)). Phylogenetic-based tools are considered the most sensitive and specific but also more time and computational resource consuming (Pineda-Peña et al. 2013; Fabeni et al. 2017). Most of the currently available tools are made available in the form of web-servers, making them easy to access and use. Nevertheless, this distribution mode raises issues regarding the scalability of the implementation; it is unreasonable to provide a web-server without limitations in the input size or number of jobs. Making large-scale analysis like multicenter molecular epidemiology studies, systematic reviews or databases curations practically impossible. Moreover, the HIV-1 genomic material corresponds to a clinical result and is often under data protection legislation as such, requiring in many cases an ethic approval for data sharing or submission in external servers.

Despite the large interest in using phylogeny for HIV-1 subtyping, existing tools have failed to address scalability and privacy issues. To answer these limitations, we used the Snakemake workflow management system (Koster and Rahmann 2012) to create a reproducible and scalable HIV-1 subtyping method based on phylogenetic pairing (SNAPPy). By combining established tools with an innovative approach, this pipeline is capable of scaling according to the available computational resources, allowing the local analysis of large amounts (tens of thousands) of HIV-1 genomes. SNAPPy was built on top of the assumption that the phylogenetic relationship provides the best possible identification of the HIV-1 subtype (Pineda-Peña et al. 2013; Fabeni et al. 2017). However, recombination events represent exceptions to the assumption of a common ancestor (coalescent) (Pérez-Losada et al. 2015). Therefore, we complemented the phylogenetic inference with the similarity search method BLAST (Camacho et al. 2009). Reproducibility and efficient transmission of protocols are current challenges in bioinformatics, critical to share domain-specific knowledge (Koster and Rahmann 2012; Di Tommaso et al. 2017). Therefore, one of the focus in SNAPPy is to give the user access to all the relevant intermediate files created and how the final subtyping decision was performed.

Overall, we present a problem-solving pipeline to allow local large-scale HIV-1 subtyping, based on phylogenetic inference and complemented with similarity search tasks.

2. Implementation

2.1 SNAPPy architecture

The SNAPPy pipeline was built on the Snakemake workflow management system (Koster and Rahmann 2012). Several tools/software were used to perform different tasks within this pipeline: MAFFT v7.245 (Katoh and Standley 2013) for multiple sequence alignment (MSA); the Biopython v1.72 (Cock et al. 2009) module SeqIO and the ETE Toolkit (ete3) v3.1.1 (Huerta-Cepas, Serra and Bork 2016) for data parsing and manipulation; BLAST v2.7.1 (Camacho et al. 2009) for local database search; IQ-TREE v1.6.9 (Nguyen et al. 2015) for phylogenetic inference. Other Python v3.6 (Van Rossum and Drake 2009) packages were also

used to create tests, pytest (Krekel et al. 2019), and data manipulation, numpy (Oliphant 2006) and pandas (McKinney 2010). For package management and to create a contained environment for SNAPPy, we recommend Conda (Anaconda Software Distribution 2019), and provided a ready to use file to this end ('environment.yml') as well as instructions on how to install and utilize it in SNAPPy's documentation page (Araújo, Martins and Osório 2019).

The term 'target' used in this manuscript refers to the file currently being processed by SNAPPy. When used for subtyping, SNAPPy performs the following tasks to a given target sequence: 1) split the input in multiple single FASTA files; 2) alignment to the reference genome; 3) BLAST against a set of HIV-1 reference sequences; 4) perform phylogenetic inferences using the BLAST top hits, the target, and an outgroup sequence; 5) sliding window BLAST against a database of HIV-1 reference sequences; 6) concatenation and analysis of the results obtained in the previous tasks and creation of the output results. Fig. 1 is a schematic representation of this pipeline. At the end of the subtyping task, SNAPPy produces two files the 'subtype_results.csv' and the 'report_subtype_results.csv', corresponding to a simplified version of the subtyping result and a more extensive report of all the outputs created by SNAPPy.

Alternatively, as any other Snakemake (Koster and Rahmann 2012) pipeline, intermediate tasks can be performed without the execution of the entire pipeline, making SNAPPy extremely useful for HIV-1 MSA (Fig. 1, Point 7). To match the names of the intermediate files created by SNAPPy and the header of the target sequences, a file named 'keys_and_ids.csv' is created. An in-depth description of each of these general steps can be seen in the following sections.

2.1.1 Reference sequences

In all instances of SNAPPy, the HXB2 (GenBank: K03455) reference genome was used as a genomic position reference. The outgroup sequence used in the phylogenetic analysis corresponds to the CONSENSUS_CPZ sequence from the HIV sequence database (Theoretical Biology and Biophysics Group 2019) (Alignment type: Consensus/Ancestral, Year: 2002, Organism: Other SIV, Alignment ID: S02CG1). The creation of a comprehensive set of HIV-1 reference sequences proved to be a challenge. We based our dataset creation on the previously curated 'HIV sequence database 2010 subtype reference genomes sequences compendium' (Kuiken et al. 2010) and the 'HIV Drug Resistance Database reference sequences' (Shafer 2006). The final subtype reference dataset for SNAPPy consisted of 491 genomic sequences. It included references for groups N, O, P, and within group M for Subtypes B, C, D, G, H, J, and K, Sub-subtypes A1, A2, F1, and F2 and CRFs until the number 99. Please notice that some CRFs are not represented due to lack of at least two high-quality genomes available (CRFs numbered 30, 66, 75, 76, 84, 89, 91, 95, 97, and 98). The full reference dataset (491 sequences) is used for the BLAST task (see Section 2.1.3) and a subset only containing groups, subtypes, Sub-subtypes and CRFs 1 and 2 references (56 sequences) is used in the sliding window BLAST (see Section 2.1.5). For more information on these reference datasets please consult the Supplementary Table S1.

2.1.2 Alignment to reference

After splitting the MSA into several single sequence FASTA files, each of them is aligned to the HIV-1 reference genome (HXB2). The module SeqIO from Biopython (Cock et al. 2009) is used to parse and manipulate the FASTA files in SNAPPy. The alignment is done using MAFFT (Katoh and Standley 2013). The

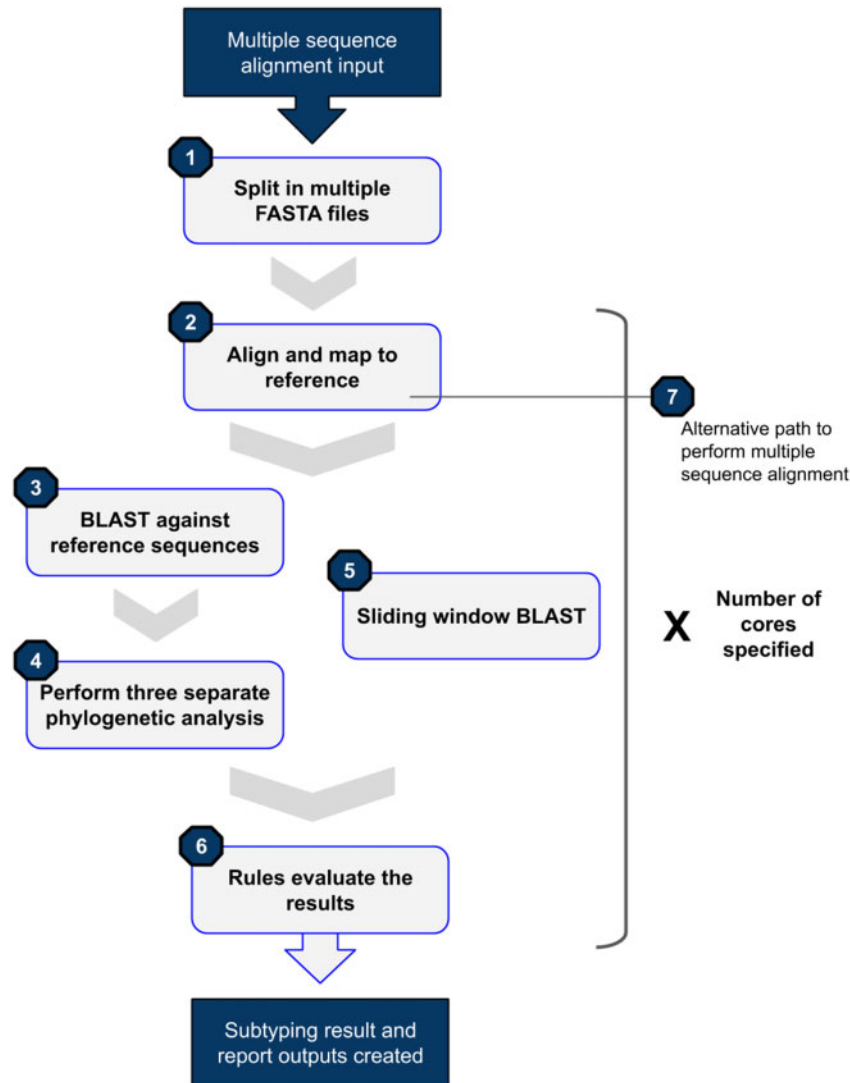


Figure 1. SNAPPy workflow diagram.

alignment method used does not allow the insertion of gaps in the reference sequence. After the alignment is performed, the target sequence is trimmed to only contain the genomic region specified by the user in the 'config.yaml' file. Being the currently available options 'GAG', 'PR', 'RT', 'PR-RT', 'INT', 'POL', 'ENV', and 'GAG-POL-ENV', which correspond to the HIV-1 genomic regions with the same names in the HXB2 reference genome. The resulting files are then written to the 'aligned' folder.

2.1.3 BLAST

The obtained alignments are then BLASTed against a local database of 491 HIV-1 reference sequences (see Section 2.1.1). For this task, BLAST (Camacho et al. 2009) is used. The results were sorted by bitscore (considering higher is better) and the best scoring result is outputted in the 'report_subtype_results.csv' file in the column 'closer_ref'. The BLAST results are also used to make three groups of reference sequences: containing the first forty-eight results; containing the first forty-eight results of only subtype references; containing the first forty-eight results of only CRF references. These three groups of reference sequences are then used in the phylogenetic analysis. The selected number of sequences (forty-eight) showed good

compromise between analysis time and reproducibility, as discussed in Section 2.1.4. Since this BLAST task is done as preparation step to the phylogenetic analysis, it must have high sensitivity, so no related reference sequences are missed. To achieve this, the word size parameter was set to 10. We also restricted the cutoff E-value to $1.0e-10$ to avoid the creation of large output files, without restricting the results. The intermediate files of the BLAST analysis are outputted to the 'blast' folder, being available for further consulting. For the split in subtype and CRF references in this step of SNAPPy, CRFs 1 and 2 are treated as subtypes, not CRFs. This decision was made based on the high prevalence of these CRFs and their ambiguous origin (Gao et al. 1996; Abecasis et al. 2007).

2.1.4 Phylogenetic inference

To the three previously selected groups of forty-eight references (see Section 2.1.3), the target sequence and a non-HIV-1 sequence for rooting (see Section 2.1.1) are added. Obtaining three sets of fifty sequences that will serve as inputs for the phylogenetic analysis. Groups of fifty sequences showed to be contained and yet a comprehensive set of sequences to perform the phylogenetic inference. To perform the phylogenetic analysis,

IQ-TREE (Nguyen et al. 2015) was used with the general time reversible (GTR) nucleotide substitution model, empirical base frequencies (+F), and a discrete Gamma model with four rate categories (+G4), with the fast tree search mode, and zero as seed number. The ETE toolkit (Huerta-Cepas, Serra and Bork 2016) was applied to parse and manipulate the phylogenetic trees created within SNAPPy. After rooting on the outgroup; it is inferred if the target sequence belongs to a monophyletic clade with sequences of one, and only one, subtype or CRF. If this happened, we consider that there is phylogenetic evidence of the relationship between the target sequence and a reference subtype/CRF. The result of this inference, together with the support values for that node (Shimodaira-Hasegawa like approximate likelihood ratio test with 1,000 replicates test, as implemented in IQ-TREE), are then outputted to the 'report_subtype_results.csv' file. Resulting in six output columns: 'node_all_refs', 's_node_all_refs', 'node_pure_refs', 's_node_pure_refs', 'node_recomb_refs', and 's_node_recomb_refs'. The intermediate files of the phylogenetic analysis are outputted to the 'trees' folder. The notation 'all', 'pure', and 'recomb' refers to the set of references used for that phylogenetic reconstruction.

2.1.5 Sliding window BLAST

The sliding window BLAST can be performed in parallel with the tasks described in Sections 2.1.3 and 2.1.4, being only dependent on the outputs from the Alignment to the reference task (Section 2.1.2). Initially, the positions in the target sequence corresponding to gaps ('-') are excluded. For this task, BLAST (Camacho et al. 2009) is used. The length of the sliding window was set to 400 nucleotides, a size previously reported to allow phylogenetic inference in HIV-1 (Pineda-Peña et al. 2013). Smaller fragments/sequences are not processed by this method. The step size used is fifty nucleotides, creating eight bins for each window. The result for each BLAST window, and consequently its eight bins, is the subtype of the top result (bitscore) reference sequence. If more than one sequence of different subtypes has the same top score, the output for all bins of that window is null ('-'). If the method fails to produce an output, the result for all bins of that window is null. After all possible sliding windows have been BLASTed, several bins will have multiple outputs. Then, a majority rule is applied to decide the final subtype for that bin. In case of a tie, the result for that bin is null. Fig. 2 contains a schematic representation of this process. The database used to BLAST against in this task only contains the group, subtype, sub-subtype and CRF1 1 and 2 references, as described in the references section (fifty-six sequences). The word size parameter applied was 30, with the purpose of obtaining high specificity. Values higher than this showed to cause instability in our tests (low reproducibility). An E-value cutoff of 1.e-50 was used to ensure the generated outputs were not too large, without losing real BLAST hits. The outputs of this inference are written to the 'report_subtype_results.csv' file in the column 'recomb_result' and the resulting files from these tasks are outputted to the 'blast' folder.

2.1.6 Decision rules

The results generated by the full sequence BLAST, phylogenetic analysis, and the sliding window BLAST are then processed using a set of rules to produce the final output. These rules are executed in order, meaning that the second rule is only applied if the first rule criteria were not met, and so forth. The list of these rules can be consulted in Supplementary Table S2. At the end of the process, two final outputs are created in the snappy folder:

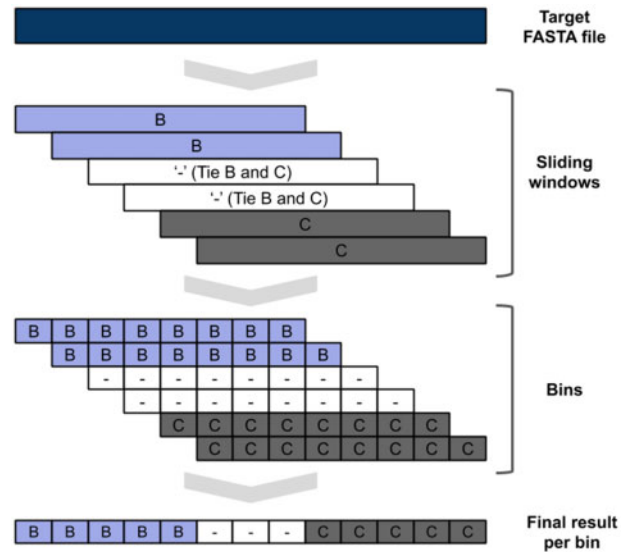


Figure 2. Sliding window BLAST schematic representation.

'subtype_results.csv' and 'report_subtype_results.csv', as mentioned earlier.

2.1.7 System management

SNAPPy is a pipeline that performs several tasks, some of them generate a relatively large number of outputs. Therefore, in order to avoid wasting unnecessary disk space and simplifying the user experience, some of the intermediate files produced by SNAPPy are deleted before the end of the process. However, all the relevant files for the subtyping decision are kept and available for consulting after the pipeline finishes. At the end of each SNAPPy run, a snakemake hidden folder named 'snakemake' is deleted because it occupies a substantial amount of space. However, this folder contains all the logs about the tasks performed and may be useful for debugging.

SNAPPy is distributed with a series of built-in tests created using pytest (Krekel et al. 2019). After installing SNAPPy or after making alterations in SNAPPy's folder is recommended to run the tests to infer if the pipeline is behaving as expected.

SNAPPy's documentation (Araújo, Martins and Osório 2019) includes detailed instructions on: pipeline installation; tutorials on how to use it; a list of available commands; an in-depth description of each pipeline step; FAQ; and how to cite and user license information.

2.2 Pipeline evaluation

2.2.1 Reproducibility

One of the bases of SNAPPy is phylogenetic inference, which has some stochasticity involved. As implemented in IQ-TREE (Nguyen et al. 2015), the initial topologies are constructed based on heuristic methods, which afterwards is optimized with maximum-likelihood rearrangements. In theory, this could lead to variance in the output of the subtyping pipeline. Furthermore, for the evaluation of the branch support, we were faced with the possibility of using statistic-based (e.g. Shimodaira-Hasegawa test) or sampling-based (e.g. bootstrapping) methods. In our tests, we found that the usage of bootstrapping approaches leads to some lack of reproducibility in the pipeline outputs, even at one thousand replicates, as previously reported (Pineda-Peña et al. 2013). Which may be

explained by the low percentage of informative sites found in the tested HIV-1 sequences, which are extremely similar among each other. The increment of the bootstrapping replicates would lead to an exponential increment of the phylogenetic inference step computational time, making the pipeline much slower. Therefore, we decided to use a statistic-based branch support inference method, the Shimodaira-Hasegawa test, for phylogenetic inference in SNAPPy. As stated in the Sections 2.1.3 and 2.1.5, when using BLAST, the word size and cutoff parameters were selected to achieve the desired objectives (sensitivity or specificity) and ensure the stability of the analysis. After the pipeline was constructed, we performed 6 sets of 3 independent SNAPPy runs with a test set of 5,285 sequences (see Section 2.2.2) and compared the outputs of each independent run in terms of reproducibility. The obtained result was 100% reproducibility, meaning that the output file 'subtyping_results.csv' for each of the independent runs were exactly the same.

2.2.2 Scalability

SNAPPy was built for large-scale analysis, taking advantage of modern multi core/thread CPUs. The usage of Snakemake (Koster and Rahmann 2012) as a workflow manager allows the construction of a directed acyclic graph of jobs, inferring which tasks need to be performed sequential and which can run in parallel. To infer the overall scalability of SNAPPy regarding multithreading, we performed the subtyping of a test set of 5,285 sequences with the following number of CPU threads: 1, 2, 4, 8, 16, and 32. The selection of these numbers was made having as objective the comparison of the computation time reduction by half (halving) when doubling the amount of computational resources. In Fig 3, there is a comparison of the real time that SNAPPy took to subtype the test set versus the expected halving time. This expected time reduction is purely theoretical and constructed based on the time SNAPPy took to subtype the test set with one core and subsequence duplication of the number of computational resources used. The performance regarding smaller test sets and for specific genomic regions was also evaluated and can be consulted in Supplementary Table S3. These tests were performed in a server with double Xeon E5-2680 2.50 GHz CPU (twelve cores/twenty-four threads), 128 GB of ram 2,133 MHz, in a SATA III SSD hard drive.

SNAPPy manages the generated files in order to give the user all the information needed to understand the results and decisions made. This feature is a tradeoff purposely made to give users the maximum amount of information without wasting disk space. Nevertheless, when used for large-scale analysis (tens of thousands of sequences), SNAPPy will create a large number of small files that together occupy a considerable amount of disk space. As an indicator, a SNAPPy run of 50,000 HIV-1 sequences occupied at the peak 59 GB and <4 GB after the depletion of the snakemake hidden logs folder.

2.2.3 Subtyping methods comparison

The division of HIV-1 in groups, subtypes, and sub-subtypes is extremely valuable for epidemiological inferences (Abecasis et al. 2013; Yebra et al. 2015; Araújo et al. 2019). However, this division is a man-made construction that only makes sense in the eyes of a phylogenetic or epidemic reconstruction (Hemelaar 2013; Araújo et al. 2019). Therefore, it makes sense to argue that phylogenetic reconstruction is the gold standard for HIV-1 subtyping (Pineda-Peña et al. 2013; Fabeni et al. 2017), with the exception of recombination events that represent a deviation

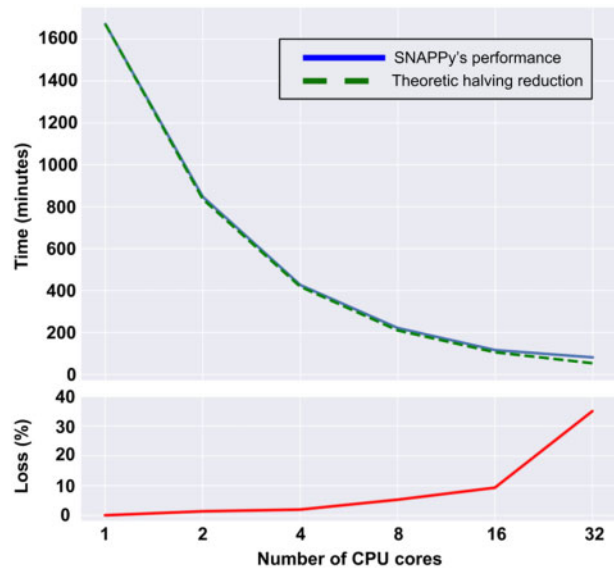


Figure 3. Multiple thread CPU time performance and associated percentage loss.

from the coalescent assumption (Pérez-Losada et al. 2015). Since SNAPPy is based on phylogenetic inference, using phylogeny to evaluate SNAPPy's performance would be poorly informative. Given this limitation, we decided to test SNAPPy against a set of other HIV-1 subtyping methods, evaluating their convergence and divergence. The selected HIV-1 subtyping methods were REGA v3.0 (Pineda-Peña et al. 2013), COMET v2.3 (Struck et al. 2014), and SCUEAL (Kosakovsky Pond et al. 2009). This selection was made based on our best knowledge of available tools including statistical and phylogeny-based methods.

For this comparison, a test set of 5,285 sequences was created (the ids of the test set sequences can be consulted in Supplementary Table S4). From all the available complete HIV-1 genomes (>9,100) in the HIV-1 sequence database (Theoretical Biology and Biophysics Group 2019) 10% of sequences for each of the subtypes present (according to the database) were selected at random, comprising a total of 1,057 genomes. Those genomes were then trimmed for four genomic regions: ENV, GAG, POL, and PR-RT. Together, the full genome sequences and the trimmed replicates compose the full test set (5,285 sequences). This test set was designed to explore the capabilities of each subtyping method in an extensive variety of subtypes while using different HIV-1 genomic regions as input. However, there are differences in the methods implementations; SCUEAL only subtypes sequences of the POL region, therefore the comparison dataset for this tool is smaller (1,057 sequences \times 3 regions = 3,171 sequences), and is capable of recognizing CRFs until the number 43; REGA is able to recognize CRFs until the number 47; COMET is capable of recognizing CRFs until the number 96. The outputs for each subtyping tool required some manipulation in order to achieve a 'common language', making it possible to compare all the tools outputs. Regarding REGA outputs, the terms 'like' and 'potential recombinant' were excluded, maintaining the assigned subtype; the outputs with 'recombination of' were named URFs of the described subtypes or URF_CPX if there was an indication of more than two recombining subtypes. REGA results marked only with the information of 'Recombinant' were transformed into URF_CPX; and when the outputs were 'Check the report' no subtyping result was assigned. For the COMET outputs the only transformation

was to change ‘unassigned’ to URFs of the subtypes present or URF_CPX if more than two subtypes were reported. Concerning the SCUEAL outputs, the words ‘ancestral’ and ‘like’ were excluded; ‘Complex’ was converted to URF_CPX and ‘AE’ to CRF_01. The SCUEAL results with ‘recombinant’ and more than two subtypes were converted to URF_CPX and those with less than two subtypes converted to URFs; for outputs with ‘U’ and ‘FAILED’ no subtyping result was assigned. The comparison between the four subtyping methods results can be seen in Fig. 4. The highest level of agreement was observed between SNAPPy and COMET (83%), whereas SCUEAL and REGA had the lowest concordance (61%). The remaining pairs showed results in the range between 72 and 78%.

We also calculated the precision, recall, and F1 scores (balance of precision and recall) for the three subtyping methods tested (REGA, SCUEAL, and COMET) versus SNAPPy (Supplementary Tables S5–S7). This analysis was performed to give an indication for users comparing results obtained with other tools and SNAPPy, for each HIV-1 group, subtype, sub-subtype, CRF, and URF. The precision and recall metrics can be seen as indicators if SNAPPy is classifying a given subtype in the test set more or less often, respectively, than the subtyping method it is being compared with. Without surprise, in the test set the results for Subtypes B and C (the most abundant) and non-M HIV-1 groups (N, O, and P) showed the highest F1 scores. The results for URFs and CRFs showed great variability.

3. Discussion and conclusions

The quantity of available HIV-1 genomes is ever increasing; the manual handling of such large amounts of data is impractical, leading to the need of creating analysis pipelines. Such pipelines targeting specific challenges are a practical and effective way of disseminating domain knowledge and increasing reproducibility (Koster and Rahmann 2012; Di Tommaso et al. 2017). The test set used for the different metrics evaluated here is composed of 1,057 sequences of HIV-1 genomes and the same sequences trimmed for the genomic regions ENV, GAG, POL, and PR-RT corresponding to a total of 5,285 sequences.

There is some stochasticity involved in the phylogenetic inference process as established in IQ-TREE (Nguyen et al. 2015). Nevertheless, the parameters selected for the branching support evaluation and the number of samples per tree allowed 100% reproducibility among independent SNAPPy runs results in an extremely diverse test set. This outcome highlights the versatility and reliability of this pipeline.

Together with the increment in the amount of data available, there have been hardware improvements, particularly in recent years CPUs with a high number of cores/threads (≥ 8) have reached the mainstream segment of the market. Therefore, tool building should be done to take maximum advantage of these resources. Snakemake (Koster and Rahmann 2012), the pipeline workflow management systems that SNAPPy is built upon, allows an almost linear scaling in the ratio of computational time/number of CPU cores used. The expected reduction by half of the computational time by doubling the number of threads used was observed in SNAPPy runs of the test set with minor percentage lost when using 2, 4, 8, or 16 CPU threads (1, 2, 5, and 9%, respectively). However, for the runs with thirty-two threads, the drop from the expected halving time was almost 35%. This drop may be a consequence of the hardware used, since it is composed of twenty-four independent CPU cores (two CPUs \times twelve cores) passing that number may cause single thread performance loss. Therefore, we do not advice the

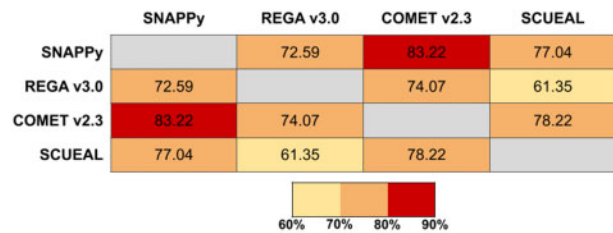


Figure 4. Percentage of agreement observed among the HIV-1 subtyping tools tested.

usage of SNAPPy multithreading capabilities in more instances than the number of physical CPU cores of the machine used.

The classification of different HIV-1 sequences in groups, subtypes, sub-subtypes, or recombinant forms is a challenging and sometimes ambiguous process. Sequence-based phylogenetic reconstruction of the evolutionary history assuming a common ancestry of the viral samples is consensually identified as the best approach for HIV-1 subtyping (Pineda-Peña et al. 2013; Fabeni et al. 2017). However, the coalescent assumption is not fulfilled in the cases of recombination (Pérez-Losada et al. 2015). Therefore, two complementary approaches were used in SNAPPy, one based on BLAST (Camacho et al. 2009) similarity searches and another based on the phylogenetic inference (IQ-TREE, Nguyen et al. 2015).

The side-by-side comparison of different HIV-1 subtyping methodologies is complex and sometimes impossible. Here, we compared SNAPPy, REGA (Pineda-Peña et al. 2013), COMET (Struck et al. 2014), and SCUEAL (Kosakovsky Pond et al. 2009) outputs for a test set of 5,285 sequences (Fig. 4). As described in the methods section, the regions of the HIV-1 genome these tools are capable of subtyping and the number of CRFs they are able to identify varies. Without surprise, COMET and SNAPPy have the highest value of accordance ($>80\%$) among the tested tools. This outcome is highly influenced by the fact that these tools are prepared to identify a large number of CRFs (>90) in comparison with the remaining two tools in this test. The lowest accordance was observed between the pair REGA and SCUEAL (61%). The lowest pairing for SNAPPy was REGA with 73% followed by SCUEAL with 77%.

In Supplementary Tables S5–S7, we show the precision, recall, and F1 scores resulting from the comparison of three subtyping methods (REGA, COMET and SCUEAL) with SNAPPy, for each HIV-1 group, subtype, sub-subtype, CRF, and URF. The overall F1 scores for REGA and SCUEAL suffer from the fact that these tools identify a narrower range of CRFs than SNAPPy. Moreover, SCUEAL only subtypes sequences for the POL HIV-1 genomic region, being sequences from the remaining regions treated as missing data, and therefore driving the F1 score further down. Nevertheless, is it expected a great reproducibility among these two tools and SNAPPy for Subtypes B and C, group N and several CRFs (1, 5, 13, 27, 35, 40, 42, and 47 for REGA and 5, 17, 18, 19, 24, 31, and 33 for SCUEAL). On the other hand, COMET is capable of identifying a wide range of CRFs, similarly to SNAPPy, which is observed in the overall F1 score (0.83), precision (0.87), and recall (0.83) results. Regarding the test set, these two tools identified CRFs 17, 27, 34, 40, 47, 68, and 74 in exactly the same cases and Subtypes B and C, and Groups N, P, and O with a high similarity rate (F1 scores, respectively: 0.99, 0.94, 1.0, 1.0, and 0.96). The results for subtypes that showed less reproducibility among the three tested methods and SNAPPy, were Subtypes H and K and sub-Subtype A2 and F2.

These results highlight the variability observed among HIV-1 subtyping tools, which is expected (Gifford et al. 2006; Pineda-Peña et al. 2013; Fabeni et al. 2017) and should not be seen as a drawback but instead as an interval of possibilities around a subtyping result. Moreover, these results demonstrate that when several HIV-1 tools agree in one result, there is a high degree of confidence in that outcome. SNAPPy is not strictly better or worse than the other tools regarding the final result, but it is a needed addition to this space, allowing local large-scale HIV-1 subtyping while being versatile, reliable and cable of scaling. The results reported here emphasized the importance of SNAPPy to facilitate the subtype annotation of large datasets of HIV-1 genomic sequences. This work represents a novel approach for HIV-1 subtyping that can contribute significantly towards a better understanding of the relevant roles and traits of the different HIV-1 subtypes.

Key points

- The amount of available HIV-1 genomic information is increasing, therefore there is a need to create tools to perform scalable analysis of HIV-1 subtypes.
- HIV-1 subtyping methods have considerable differences in their implementations and, consequently, their outputs.
- SNAPPy is capable of local large-scale HIV-1 subtyping with great reproducibility while being able to scale according to the computational resources available.

Data availability

SNAPPy source code is freely available via GitHub at: <https://github.com/PMMAraujo/snappy/releases>. SNAPPy documentation can be consulted at: <https://snappy-hiv1-subtyping.readthedocs.io/>.

Supplementary data

Supplementary data are available at *Virus Evolution* online.

Funding

This work was supported by FEDER, COMPETE, and FCT by the projects NORTE-01-0145-FEDER-000013, POCI-01-0145-FEDER-007038, and IF/00474/2014; FCT PhD scholarship PDE/BDE/113599/2015; FCT contract IF/00474/2014.

Conflict of interest: None declared.

References

Abecasis, A. B. et al. (2005) 'Protease Mutation M89I/V Is Linked to Therapy Failure in Patients Infected with the HIV-1 non-B Subtypes C, F or G', *AIDS*, 19: 1799–806.

— et al. (2007) 'Recombination Confounds the Early Evolutionary History of Human Immunodeficiency Virus Type 1: Subtype G Is a Circulating Recombinant Form', *Journal of Virology*, 81: 8543–51.

— et al. (2013) 'HIV-1 Subtype Distribution and Its Demographic Determinants in Newly Diagnosed Patients in

Europe Suggest Highly Compartmentalized Epidemics', *Retrovirology*, 10: 7.

Anaconda Software Distribution. Computer Software. Vers. 3-4.6.14. Miniconda, April 2019. <<https://anaconda.com>>.

Araujo, P. M. M. et al. (2019) 'Characterization of a Large Cluster of HIV-1 A1 Infections Detected in Portugal and Connected to Several Western European Countries', *Scientific Report*, 9: 7223.

Araújo, P. M. M., Martins, J. S., and Osório, N. S. SNAPPy documentation. September 2019. <<https://snappy-hiv1-subtyping.readthedocs.io/en/latest/>>.

Baeten, J. M. et al. (2007) 'HIV-1 Subtype D infection is associated with Faster Disease Progression than Subtype a in Spite of Similar Plasma HIV-1 Loads', *The Journal of Infectious Diseases*, 195: 1177–80.

Bartolo, I. et al. (2011) 'Origin and Epidemiological History of HIV-1 CRF14_BG', *PLoS One*, 6: e24130.

Brenner, B. et al. (2003) 'A V106M Mutation in HIV-1 Clade C Viruses Exposed to Efavirenz Confers Cross-Resistance to Non-Nucleoside Reverse Transcriptase Inhibitors', *AIDS*, 17: F1–5.

Camacho, C. et al. (2009) 'BLAST+: Architecture and Applications', *BMC Bioinformatics*, 10: 421.

Camacho, R. J., and Vandamme, A.-M. (2007) 'Antiretroviral Resistance in Different HIV-1 Subtypes: Impact on Therapy Outcomes and Resistance Testing Interpretation', *Current Opinion in HIV and AIDS*, 2: 123–9.

Cock, P. J. A. et al. (2009) 'Biopython: Freely Available Python Tools for Computational Molecular Biology and Bioinformatics', *Bioinformatics*, 25: 1422–3.

Di Tommaso, P. et al. (2017) 'Nextflow Enables Reproducible Computational Workflows', *Nature Biotechnology*, 35: 316–9.

Easterbrook, P. J. et al. (2010) 'Impact of HIV-1 Viral Subtype on Disease Progression and Response to Antiretroviral Therapy', *Journal of the International Aids Society*, 13: 4.

Fabeni, L. et al. (2017) 'Comparative Evaluation of Subtyping Tools for Surveillance of Newly Emerging HIV-1 Strains', *Journal of Clinical Microbiology*, 55: 2827–37.

Gao, F. et al. (1996) 'The Heterosexual Human Immunodeficiency Virus Type 1 Epidemic in Thailand Is Caused by an Intersubtype (a/E) Recombinant of African Origin', *Journal of Virology*, 70: 7013–29.

Gifford, R. et al. (2006) 'Assessment of Automated Genotyping Protocols as Tools for Surveillance of HIV-1 Genetic Diversity', *AIDS*, 20: 1521–9.

Hemelaar, J. (2013) 'Implications of HIV Diversity for the HIV-1 Pandemic', *Journal of Infection*, 66: 391–400.

Huerta-Cepas, J., Serra, F., and Bork, P. (2016) 'ETE 3: Reconstruction, Analysis, and Visualization of Phylogenomic Data', *Molecular Biology and Evolution*, 33: 1635–8.

John-Stewart, G. C. et al. (2005) 'Subtype C Is Associated with Increased Vaginal Shedding of HIV-1', *The Journal of Infectious Diseases*, 192: 492–6.

Katoh, K., and Standley, D. M. (2013) 'MAFFT Multiple Sequence Alignment Software Version 7: Improvements in Performance and Usability', *Molecular Biology and Evolution*, 30: 772–80.

Kiwanuka, N. et al. (2008) 'Effect of Human Immunodeficiency Virus Type 1 (HIV-1) Subtype on Disease Progression in Persons from Rakai, Uganda, with Incident HIV-1 Infection', *The Journal of Infectious Diseases*, 197: 707–13.

Kosakovsky Pond, S. L. et al. (2009) 'An Evolutionary Model-Based Algorithm for Accurate Phylogenetic Breakpoint Mapping and Subtype Prediction in HIV-1', *PLoS Computational Biology*, 5: e1000581.

- Koster, J., and Rahmann, S. (2012) 'Snakemake—a Scalable Bioinformatics Workflow Engine', *Bioinformatics*, 28: 2520–2.
- Krekel, H. et al. Computer software. Vers. 4.5.0. Pytest, April 2019. <<https://docs.pytest.org/en/latest/>>.
- Kuiken, C. et al., eds. (2010). *HIV Sequence Compendium*. NM: Theoretical Biology and Biophysics Group, Los Alamos National Laboratory, LA-UR 10-03684.
- Liu, T. F., and Shafer, R. W. (2006) 'Web Resources for HIV Type 1 Genotypic-Resistance Test Interpretation', *Clinical Infectious Diseases : An Official Publication of the Infectious Diseases Society of America*, 42: 1608–18.
- McKinney, W., van der Walt, S. and Millman, J. (eds) (2010) 'Data Structures for Statistical Computing in Python'. in van der Walt S. and Millman, J. (eds) *Proceedings of the 9th Python in Science Conference*, pp. 51–6. <http://conference.scipy.org/proceedings/scipy2010/mckinney.html>.
- Myers, R. E. et al. (2005) 'A Statistical Model for HIV-1 Sequence Classification Using the Subtype Analyser (STAR)', *Bioinformatics*, 21: 3535–40.
- Nguyen, L.-T. et al. (2015) 'IQ-TREE: A Fast and Effective Stochastic Algorithm for Estimating Maximum-Likelihood Phylogenies', *Molecular Biology and Evolution*, 32: 268–74.
- Oliphant, T. E. (2006) *A Guide to NumPy*. USA: Trelgol Publishing.
- Pérez-Losada, M. et al. (2015) 'Recombination in Viruses: Mechanisms, Methods of Study, and Evolutionary Consequences', *Infection, Genetics and Evolution*, 30: 296–307.
- Pineda-Peña, A.-C. et al. (2013) 'Automated Subtyping of HIV-1 Genetic Sequences for Clinical and Surveillance Purposes: Performance Evaluation of the New REGA Version 3 and Seven Other Tools', *Infection, Genetics and Evolution*, 19: 337–48.
- Renjifo, B. et al. (2004) 'Preferential in-Utero Transmission of HIV-1 Subtype C as Compared to HIV-1 Subtype a or D', *AIDS*, 18: 1629–36.
- Robertson, D. L. et al. (2000) 'HIV-1 Nomenclature Proposal', *Science*, 288: 55–6.
- Rozanov, M. et al. (2004) 'A Web-Based Genotyping Resource for Viral Sequences', *Nucleic Acids Research*, 32: W654–9.
- Schultz, A.-K. et al. (2009) 'jpHMM: Improving the Reliability of Recombination Prediction in HIV-1', *Nucleic Acids Research*, 37: W647–51.
- Serwanga, J. et al. (2015) 'Frequencies of Gag-Restricted T-Cell Escape 'Footprints' Differ across HIV-1 Clades A1 and D Chronically Infected Ugandans Irrespective of Host HLA B Alleles', *Vaccine*, 33: 1664–72.
- Shafer, R. W. (2006) 'Rationale and Uses of a Public HIV Drug-Resistance Database', *The Journal of Infectious Diseases*, 194(Suppl): S51–8.
- Struck, D. et al. (2014) 'COMET: Adaptive Context-Based Modeling for Ultrafast HIV-1 Subtype Identification', *Nucleic Acids Research*, 42: e144.
- Theoretical Biology and Biophysics Group, Los Alamos National Laboratory, HIV sequence database, April 2019. <<http://www.hiv.lanl.gov/>>.
- Van Rossum, G., and Drake, F. L. (2009) *Python 3 Reference Manual*. CA: CreateSpace.
- Yebra, G. et al. (2015) 'Analysis of the History and Spread of HIV-1 in Uganda Using Phylodynamics', *Journal of General Virology*, 96: 1890–8.