

 Open access • Posted Content • DOI:10.1101/628073

## **DEN-IM: Dengue Virus identification from shotgun and targeted metagenomics**

— [Source link](#) 

Catarina I. Mendes, Catarina I. Mendes, Erley Lizarazo, Miguel P. Machado ...+6 more authors

**Institutions:** University Medical Center Groningen, Instituto de Medicina Molecular

**Published on:** 06 May 2019 - bioRxiv (Cold Spring Harbor Laboratory)

**Topics:** Metagenomics

Related papers:

- [A Nanopore-based method for generating complete coding region sequences of dengue virus in resource-limited settings](#)
- [drVM: a new tool for efficient genome assembly of known eukaryotic viruses from metagenomes.](#)
- [ONTdeCIPHER: An amplicon-based nanopore sequencing pipeline for tracking pathogen variants](#)
- [Comparison of third-generation sequencing approaches to identify viral pathogens under public health emergency conditions.](#)
- [Computational Framework for Next-Generation Sequencing of Heterogeneous Viral Populations using](#)

Share this paper:    

View more about this paper here: <https://typeset.io/papers/den-im-dengue-virus-identification-from-shotgun-and-targeted-4d81tobbx3>

# DEN-IM: Dengue Virus identification from shotgun and targeted metagenomics

C I Mendes<sup>1,2,\*</sup>, †, E Lizarazo<sup>2, †</sup>, M P Machado<sup>1</sup>, D N Silva<sup>1</sup>, A Tami<sup>2</sup>, M Ramirez<sup>1</sup>, N Couto<sup>2</sup>, J W A Rossen<sup>2</sup>, J A Carriço<sup>1</sup>

**1 Instituto de Microbiologia, Instituto de Medicina Molecular, Faculdade de Medicina, Universidade de Lisboa, Lisboa, Portugal**

**2 University of Groningen, University Medical Center Groningen, Department of Medical Microbiology and Infection Prevention, Groningen, The Netherlands**

\* [cimendes@medicina.ulisboa.pt](mailto:cimendes@medicina.ulisboa.pt) † Contributed equally

## Abstract

Dengue virus (DENV) represents a public health and economic burden in affected countries. The availability of genomic data is key to understand viral evolution and dynamics, supporting improved control strategies. Currently, the use of High Throughput Sequencing (HTS) technologies, which can be applied both directly to patient samples (shotgun metagenomics) and to PCR amplified viral sequences (targeted metagenomics), is the most informative approach to monitor the viral dissemination and genetic diversity.

Despite many advantages, these technologies require bioinformatics expertise and appropriate infrastructure for the analysis and interpretation of the resulting data. In addition, the many software solutions available can hamper reproducibility and comparison of results.

Here we present DEN-IM, a one-stop, user-friendly, containerised and reproducible workflow for the analysis of DENV sequencing data, both from shotgun and targeted metagenomics approaches. It is able to infer DENV coding sequence (CDS), identify serotype and genotype, and generate a phylogenetic tree. It can easily be run on any UNIX-like system, from local machines to high-performance computing clusters, performing a comprehensive analysis without the requirement of extensive bioinformatics expertise.

Using DEN-IM, we successfully analysed two DENV datasets. The first comprised 25 shotgun metagenomic sequencing samples of varying serotype and genotype, including a spiked sample containing the existing four serotypes. The second dataset consisted of 106 targeted metagenomics samples of DENV 3 genotype III where DEN-IM allowed detection of the intra-genotype diversity.

The DEN-IM workflow, parameters and execution configuration files, and documentation are freely available at <https://github.com/B-UMMI/DEN-IM>.

**Keywords:** Dengue virus; Surveillance; Metagenomics; Reproducibility; Workflow; Containerization; Scalability

## 1 Key points

- Understanding DENV transmission in populations where the infection is endemic through the use of metagenomics is a most promising strategy to get insight into the dissemination of the virus and to support measures to prevent or decrease further spread.
- So far, the analysis, interpretation and dissemination of results encounters computational challenges, and requires appropriate expertise and infrastructure, which poses a challenge for the implementation of metagenomic approaches.
- We present DEN-IM, a reproducible, containerised and user-friendly workflow for the identification and characterisation of DENV from shotgun and targeted metagenomics data.
- The HTML reports obtained with DEN-IM can be easily shared, facilitating comparison of results from across the globe.

## 2 Background

The Dengue virus (DENV), a single-stranded positive-sense RNA virus belonging to the *Flavivirus* genus, is one of the most prevalent arboviruses and is mainly concentrated in tropical and subtropical regions. Infection with DENV results in symptoms ranging from mild fever to haemorrhagic fever and shock syndrome [1]. Transmission to humans occurs through the bite of *Aedes* mosquitoes namely *Aedes aegypti* and *Aedes albopictus* [2]. In 2010, it was predicted that the burden of dengue disease reached 390 million cases per year worldwide [3]. The high morbidity and mortality of dengue makes it the arbovirus with the highest clinical significance [4].

The viral genome of ~11,000 nucleotides, consists of a Coding Sequence (CDS) of approximately 10.2 Kb that is translated into a single polyprotein encoding three structural proteins (capsid - C, premembrane - prM, envelope - E) and seven non-structural proteins (NS1, NS2A, NS2B, NS3, NS4A, NS4B and NS5). Additionally, the genome contains two Non-Coding Regions (NCRs) at their 5' and 3' ends [5].

DENV can be classified into four serotypes (1, 2, 3 and 4), differing from each other by 25% to 40% at the amino acid level. They are further classified into genotypes that vary by up to ~3% at the amino acid level [2]. The DENV-1 serotype comprises five genotypes (I-V), DENV-2 groups six (I-VI, also named American, Cosmopolitan, Asian-American, Asian II, Asian I and Sylvatic), DENV-3 four (I-III and V), and DENV-4 also four (I-IV).

DENV is a significant public health challenge in countries where the infection is endemic due to the high health and economic burden. Despite the emergence of novel therapies and ecological strategies to control the mosquito vector, there are still important knowledge gaps in the virus biology and its epidemiology [2].

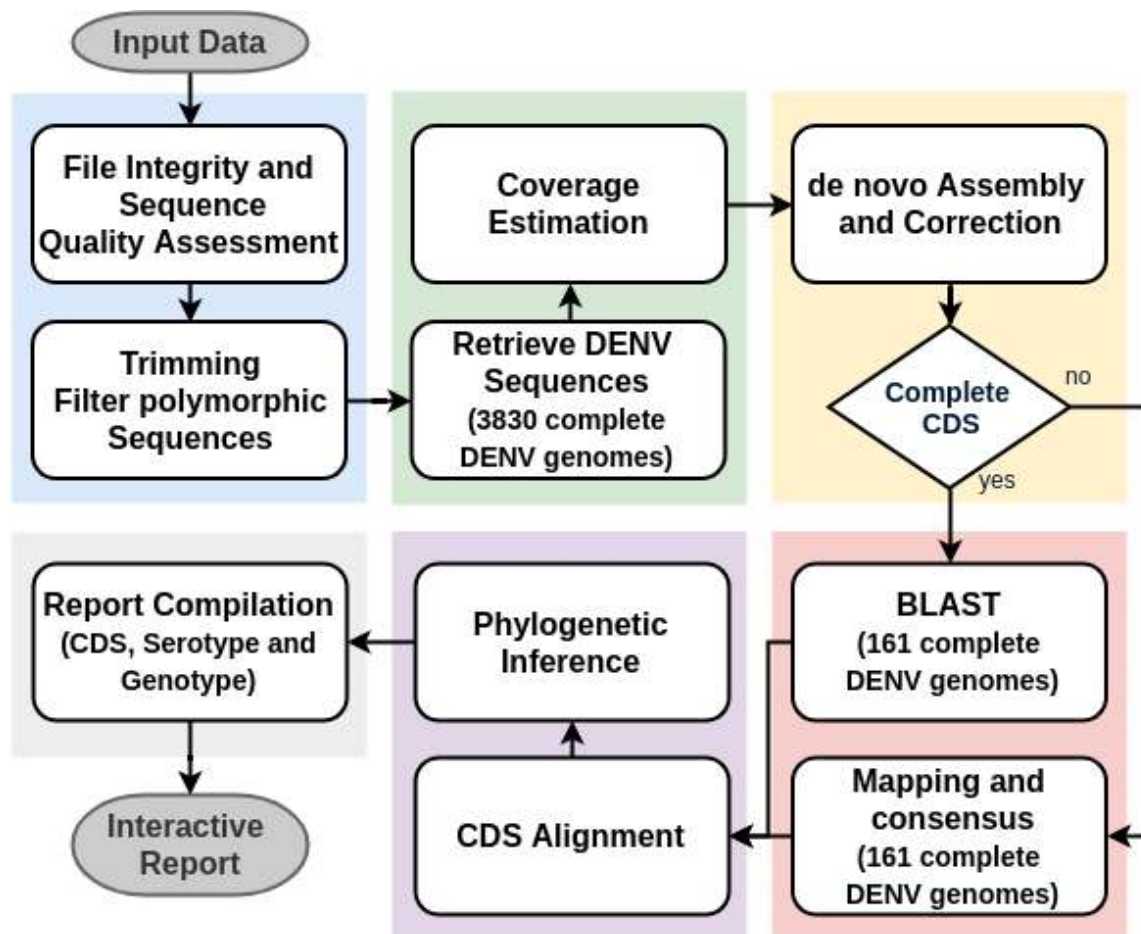
The implementation of a surveillance system relying on High Throughput Sequencing (HTS) technologies allows the simultaneous identification and surveillance of DENV cases. Due to the high sensitivity of these technologies, previous studies showed that viral sequences can be directly obtained from patient sera using a shotgun metagenomics approach [6]. Alternatively, HTS can be used in a targeted metagenomics approach in which a PCR step is used to pre-amplify viral sequences before sequencing. In recent years, HTS has been successfully used as a tool for identification of DENV directly from clinical samples [6, 7]. This also allows the rapid identification of the serotype and genotype important for disease management as the genotype may be associated with disease outcome [8].

Several initiatives aim to facilitate the identification of the DENV serotype and genotype from HTS data. The Genome Detective project (<https://www.genomedetective.com/>) offers an online Dengue Typing Tool (<https://www.genomedetective.com/app/typingtool/dengue/>) relying on BLAST and phylogenetic methods in order to identify the closest serotype and genotype, but it requires as input assembled genomes in FASTA format. Alternatively, the same project offers a Genome Detective Typing Tool (<https://www.genomedetective.com/app/typingtool/virus/>) identifying viruses present in a sample.

39 We developed DEN-IM as a ready-to-use, one-stop, reproducible bioinformatic analysis workflow  
40 for the processing and phylogenetic analysis of DENV using paired-end raw HTS data. DEN-IM is  
41 implemented in Nextflow [9], a workflow manager software that uses Docker (<https://www.docker.com>)  
42 containers with pre-installed software for all the workflow tools. The DEN-IM workflow, as well  
43 as parameters and documentation, are available at <https://github.com/B-UMMI/DEN-IM>.

### 44 3 The DEN-IM Workflow

45 DEN-IM is a user-friendly automated workflow allowing the analysis of shotgun or targeted metage-  
46 nomics data for the identification, serotyping, genotyping, and phylogenetic analysis of DENV. It is  
47 implemented in Nextflow, a workflow management system that allows the effortless deployment and



**Figure 1.** The DEN-IM workflow separated into five different components. The raw sequencing reads are provided as input to the first block (in blue), responsible for quality control and elimination of low-quality reads and sequences. After successful pre-processing of the reads, these enter the second block (green) for retrieval of the DENV reads using the mapping database of 3830 complete DENV genomes as reference. This block also provides an initial estimate of the sequencing depth. After the *de novo* assembly and assembly correction block (yellow), the coding sequences (CDSs) are retrieved and are then classified with the reduced complexity DENV typing database containing 161 sequences representing the known diversity of DENV serotypes and genotypes (red). If a complete CDS fails to be assembled, the reads are mapped against the DENV typing database and a consensus sequence is obtained for classification and phylogenetic inference. All CDSs are aligned and compared in a phylogenetic analysis (purple). Lastly, a report is compiled (gray) with the results of all the blocks of the workflow.

48 execution of complex distributed computational workflows.

49 The workflow is composed of five blocks: 1. Quality Control and Trimming, 2. Retrieval of DENV  
50 sequences, 3. Assembly, 4. *in silico* Typing, and 5. Phylogeny, described in more detail in Figure 1.  
51 DEN-IM accepts as input raw paired-end sequencing data (FASTQ files), and informs the user with  
52 an interactive HTML report with information on the quality control, mapping, assembly typing and  
53 phylogenetic analysis, as well as all the output files of the whole pipeline.

### 54 **3.1 1. Quality Control and Trimming**

55 The Quality Control (QC) and Trimming block starts with a process to verify the integrity of the  
56 input data. If the sequencing files are corrupted, the execution of the analysis of that sample is  
57 terminated.

58 The sequences are then processed by FastQC ([https://www.bioinformatics.babraham.ac.uk/  
59 projects/fastqc/](https://www.bioinformatics.babraham.ac.uk/projects/fastqc/), version 0.11.7) to determine the quality of the individual base pairs of the  
60 raw data files. The low quality bases and adapter sequences are trimmed by Trimmomatic [10]  
61 (version 0.36). By default, Trimmomatic removes the Illumina Nextera, TruSeq2 and TruSeq3 adapter  
62 sequences that are detected, and relies on Phred scores [11] and FastQC analysis data to determine  
63 to which extend the 5' and 3' ends of the reads need to be trimmed. The crop and headcrop setting  
64 is calculated per sample, based on the GC ratio of each half of the read. By default, the window size  
65 is set to 5 nucleotides with a minimum average quality of 20. In addition, paired-end reads with a  
66 read length shorter than 55 nucleotides after trimming are removed from further analyses.

67 Lastly, the low complexity sequences, containing over 50% of poly-A, poly-N or poly-T nucleotides,  
68 are filtered out of the raw data using PrinSeq [12] (version 0.10.4).

### 69 **3.2 2. Retrieval of DENV Sequences**

70 In the second step, DENV sequences are selected from the sample using Bowtie2 [13] (version  
71 2.2.9) and Samtools [14] (version 1.4.1). As a reference we provide the DENV mapping database, a  
72 curated DENV database composed of 3830 complete DENV genomes (see Methods, DENV Reference  
73 Database). A permissive approach is followed by allowing for mates to be kept in the sample even  
74 when only one read maps to the database in order to keep as many DENV derived reads as possible.  
75 The output of this step is a set of processed reads of putative DENV origin.

### 76 **3.3 3. Assembly**

77 DEN-IM applies a two assembler approach to generate assemblies of the DENV CDS. To obtain a  
78 high confidence assembly, the processed reads are first *de novo* assembled with SPAdes [15] (version  
79 3.12.0). If the full CDS fails to be assembled into a single contig, the data is re-assembled with  
80 MEGAHIT assembler [16] (version 1.1.3), a more permissive assembler developed to retrieve longer  
81 sequences from metagenomics data. The resulting assemblies are corrected with Pilon [17] (version  
82 1.22) after mapping the processed reads to the assemblies with Bowtie2 [13].

83 If more than one complete CDS is present in a sample, each of the sequences will follow the rest  
84 of the DEN-IM workflow independently. If no full CDS is assembled neither with SPAdes nor with  
85 MEGAHIT, the processed reads are passed on to the next step for consensus generation by mapping,  
86 effectively constituting DEN-IM's two pronged approach using both assemblers and mapping.

### 87 **3.4 4. Typing**

88 For each DENV complete CDS, the serotype and genotype is determined with the Seq\_Typing tool  
89 ([https://github.com/B-UMMI/seq\\_typing](https://github.com/B-UMMI/seq_typing), version 2.0) [18] using BLAST [19] and the custom  
90 Typing database of DENV containing 161 complete sequences (see Methods, DENV Reference  
91 Database). The BLAST results are first cleaned to get the best hit for each sequence in the database,  
92 based on alignment length, similarity, e-value and number of gaps. The tool determines which  
93 reference sequence is more closely related to the query based on the identity and length of the  
94 sequence covered, returning the serotype and genotype of the reference sequence for the purpose of  
95 classifying the query.

96 If a complete CDS fails to be obtained through the assembly process, the processed reads are  
97 mapped against the DENV typing database, with Bowtie2 [13], using Seq-Typing tool, with similar  
98 criteria for coverage and identity to those used with the BLAST approach. If a type is determined,  
99 the consensus sequence obtained follows through to the next step in the workflow. Otherwise, the  
100 sample is classified as Non-Typable and its process terminated.

### 101 3.5 5. Phylogeny

102 All DENV complete CDSs and consensus sequences analysed in a workflow execution are aligned  
103 with MAFFT [20] (version 7.402), in auto mode and with orientation being automatically determined  
104 and adjusted. With the resulting alignment, a Maximum Likelihood tree is inferred with RaXML [21]  
105 (version 8.2.11), using as default the GTR- $\Gamma$  substitution model and a 500 times bootstrap.

106 Optionally, the closest reference sequence in the DENV typing database to each analysed sample  
107 can be retrieved and included in the phylogenetic analysis.

## 108 4 Workflow Execution

109 The DEN-IM workflow can be executed in any UNIX-based system, from local machines to high-  
110 performance computing clusters (HPC) facilities with Nextflow and a container engine installation,  
111 such as Docker (<https://www.docker.com/>), Shifter [22] or Singularity [23].

112 Due to its Nextflow implementation, DEN-IM allows for out-of-the-box high-level parallelization  
113 and offers direct support for distributed computational environments. DEN-IM integrates Docker  
114 containerised images for all the tools necessary for its execution, ensuring reproducibility and the  
115 tracking of both software code and version, regardless of the operating system used. The Docker  
116 images provided are also compatible with other container engines.

117 Users can customise the workflow execution either by using command line options or by modifying  
118 a simple plain-text configuration file (params.config). The version of each of the tools used in  
119 DEN-IM can be changed by providing new container tags in the appropriate configuration file  
120 (containers.config). The resources for each process can also be changed (resources.config). To make  
121 the execution of the workflow as simple as possible, a set of default parameters and directives is  
122 provided.

123 The local installation of the DEN-IM workflow, including the docker containers with all the tools  
124 needed and the curated DENV database, requires 15 Gigabytes of free disk space. The minimum  
125 requirements to execute the workflow are at least 5 Gigabytes of memory and 4 CPUs, although 7  
126 Gigabytes of memory is advised. The disk space required for execution depends greatly on the size  
127 of the input data, but for the datasets used in this article DEN-IM generates approximately 20 Gb  
128 data per Gb of input data.

## 129 5 Output and Report

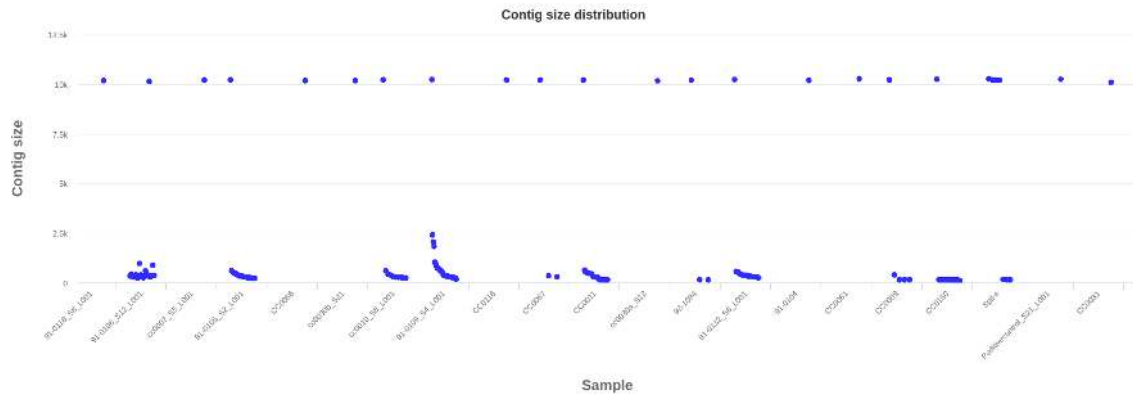
130 The output files of all tools are stored in the 'results' folder in the directory of the DEN-IM execution.  
131 The execution log file for each component, as well as a log file for the execution of the workflow, are  
132 also available.

133 DEN-IM creates an interactive HTML report (Figure S1), stored in the 'pipeline\_results' directory,  
134 containing all the information in the results divided into four sections: report overview, tables, charts  
135 and phylogenetic tree. The report can be easily exchanged between collaborators by compressing  
136 and sharing the "pipeline\_report" folder.

137 The report overview contains information about the number of samples in the analysis. It allows  
138 for selection, filtering and highlighting of particular samples and tools in the workflow.

139 The table section contains the results and statistics the quality control, assembly, read mapping  
140 and *in silico* typing results. The ***in silico* typing table** contains the results of the serotype and  
141 genotype of each CDS analysed, as well as identity and coverage and GenBank ID for the closest  
142 reference in the DENV typing database.





**Figure 2.** Contig size distribution for the shotgun metagenomics sequencing dataset. Each dot depicts an assembled DENV contig. Above the 10Kb are full CDS of DENV.

143 The **quality control table** shows information regarding the number of raw base pairs and  
144 number of reads in the raw input files and the percentage of trimmed reads. The **mapping table**  
145 includes the results for the mapping of the trimmed reads to the DENV mapping database, including  
146 the overall alignment rate, and an estimation of the sequence depth including only the DENV  
147 reads. For the **assembly statistics table**, the number of CDSs in each sample is included, and the  
148 number of contigs and the number of assembled base pairs generated by either SPAdes or MEGAHIT  
149 assemblers. The number of contigs and assembled base pairs after correction with Pilon is also  
150 presented in the table. Warning and fail messages are included in each table, as well as the ranking  
151 of the value in the cell in relation to other values in the same column, which are shown with a grey  
152 bar (Figure S1).

153 The **assembled contig size distribution scatter plot** is available in the chart section, showing  
154 the contig size distribution for the Pilon corrected assembled CDSs.

155 Lastly, a **phylogenetic tree** is included, rooted at midpoint for visualisation purposes, and  
156 with each tip coloured according to the genotyping results. If the option to retrieve the closest  
157 typing reference is selected, these sequences are also included in the tree with respective typing  
158 metadata. The tree can be displayed in several conformations provided by PhyloCanvas JavaScript  
159 library (<http://phylocanvas.net>, version 2.8.1) and it is possible to zoom in or collapse selected  
160 branches. The support bootstrap values of the branches can be displayed and the tree can be exported  
161 as a Newick tree file or as a PNG image.

## 162 6 Results

163 To evaluate the DEN-IM workflow performance, we analysed two datasets, one containing shotgun  
164 metagenomics sequencing data of patient samples and another with targeted metagenomics sequencing  
165 data obtained from Parameswaran *et al* [24].

### 166 6.1 The Shotgun Metagenomics Dataset

167 We analysed 22 shotgun metagenomics paired-end short-read Illumina sequencing datasets of clinical  
168 positive dengue cases (Table 1)(identified as such by positive qRT-PCR/RT-PCR test and serocon-  
169 version test), positive control (purified from DENV culture) and one negative control (blank), and a  
170 spiked sample containing the 4 DENV serotypes.

171 The workflow was executed using the default parameters and directives for resources, with the  
172 option to include the closest typing references in the final tree.

173 The negative control and the 92-1001 sample has no reads after trimming and filtering of low  
174 complexity reads, therefore they were removed from further analysis.

175 When mapping to the DENV mapping database, the percentage of DENV reads in the 21 clinical  
176 samples, positive control and spiked sample passing QC ranged from 0.01% (sample UCUG0186)  
177 to 85.38% (sample Positive Control). After coverage depth estimation, the analysis of the samples

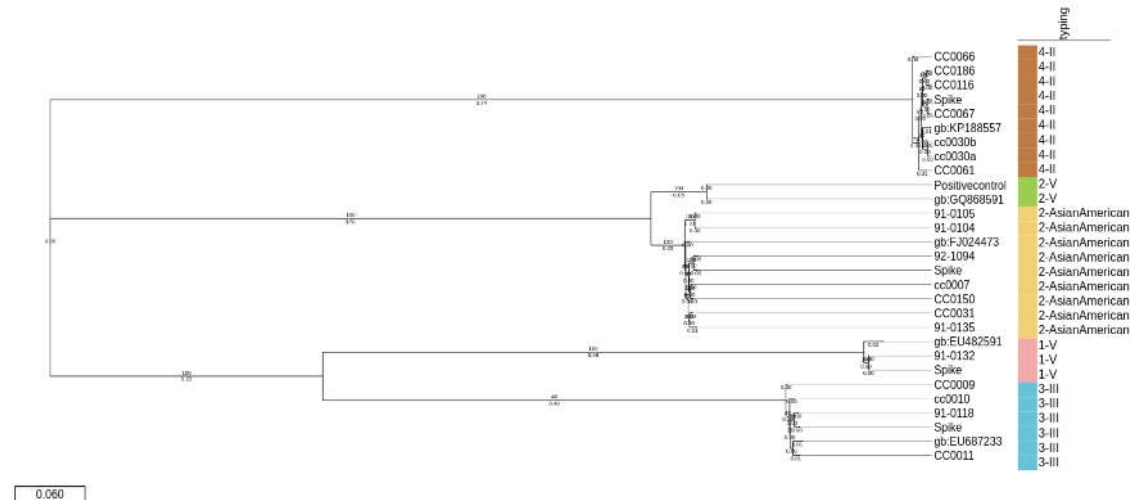
**Table 1.** Number of raw base pairs, overall alignment rate against the DENV mapping database, estimated coverage depths and serotype and genotype for 25 shotgun metagenomics sequencing samples.

| Sample                        | Raw Basepairs<br>(in megabases) | % DENV<br>Reads | Estimated Coverage<br>Depth (times) | Serotype | Genotype            |
|-------------------------------|---------------------------------|-----------------|-------------------------------------|----------|---------------------|
| 91-0104                       | 2193.71                         | 12.46           | 5944.67                             | 2        | III (AsianAmerican) |
| 91-0105                       | 191.37                          | 4.01            | 495.97                              | 2        | III (AsianAmerican) |
| 91-0115 <sup>b</sup>          | 179.24                          | 0.05            | 3.74                                | -        | -                   |
| 91-0118                       | 195.27                          | 1.69            | 86.53                               | 3        | III                 |
| 91-0132                       | 378.21                          | 20.02           | 4698.12                             | 1        | V                   |
| 91-0135                       | 91.71                           | 21.45           | 1287.52                             | 2        | III (AsianAmerican) |
| 92-1001 <sup>a</sup>          | 163.44                          | -               | -                                   | -        | -                   |
| 92-1094                       | 1197.92                         | 8.48            | 4032.21                             | 2        | III (AsianAmerican) |
| CC0007                        | 252.97                          | 3.79            | 383.77                              | 2        | III (AsianAmerican) |
| CC0009                        | 2055.13                         | 9.48            | 8226.27                             | 3        | III                 |
| CC0010                        | 368.64                          | 5.68            | 1197.58                             | 3        | III                 |
| CC0011                        | 924.69                          | 8.38            | 3016.17                             | 3        | III                 |
| CC0030a                       | 261.12                          | 52.52           | 2914.87                             | 4        | II                  |
| CC0030b                       | 399.04                          | 10.51           | 677.96                              | 4        | II                  |
| CC0031                        | 1572.1                          | 68.91           | 52318.33                            | 2        | III (AsianAmerican) |
| CC0061                        | 1262.83                         | 8.97            | 5120.4                              | 4        | II                  |
| CC0066                        | 1087.45                         | 2.8             | 569.7                               | 4        | II                  |
| CC0067                        | 1022.06                         | 5.55            | 2548.84                             | 4        | II                  |
| CC0116                        | 773.31                          | 6.72            | 2313.99                             | 4        | II                  |
| CC0150                        | 1403.69                         | 17.41           | 12065.81                            | 2        | III (AsianAmerican) |
| CC0186                        | 671.78                          | 0.03            | 14.71                               | 4        | II                  |
| UCUG0186 <sup>b</sup>         | 1116.67                         | 0.01            | 5.65                                | -        | -                   |
| Negative Control <sup>a</sup> | 163.67                          | -               | -                                   | -        | -                   |
| Positive Control              | 443.93                          | 85.38           | 19362.07                            | 2        | V (Asian I)         |
| Spike                         | 1518.93                         | 41.7            | 22289.98                            | 3        | III                 |
|                               |                                 |                 |                                     | 1        | V                   |
|                               |                                 |                 |                                     | 2        | III (AsianAmerican) |
|                               |                                 |                 |                                     | 4        | II                  |

a) Failed quality control - No sequence data after filtering of polymorphic sequences.

b) Failed quality control - Low sequence depth (<10x).





**Figure 3.** Maximum Likelihood tree in the DEN-IM report for the 24 complete CDSs (n=21 samples) obtained with the metagenomics dataset and the respective closest references in the typing database (identified by their GenBank ID). The tree is midpoint rooted for visualisation purposes. The colours depicts the DENV genotyping results.

178 91-0115 and UBUG0186 was terminated since they did not meet the threshold criterion of having an  
 179 estimated depth of coverage of 10x.

180 In the assembly module, the remaining 19 clinical samples, the spiked sample and the positive  
 181 control were assembled with DEN-IM's two assembler approach. Twenty-four full CDS were assembled  
 182 (Figure 2), despite originally having DENV reads content as low as 0.03% of the total number of  
 183 reads. Sixteen samples, including the spiked sample and the positive control, were assembled in the  
 184 first step with the SPAdes assembler, and five in the second with the MEGAHIT assembler. In the  
 185 spiked sample, all four CDSs were successfully assembled and recovered.

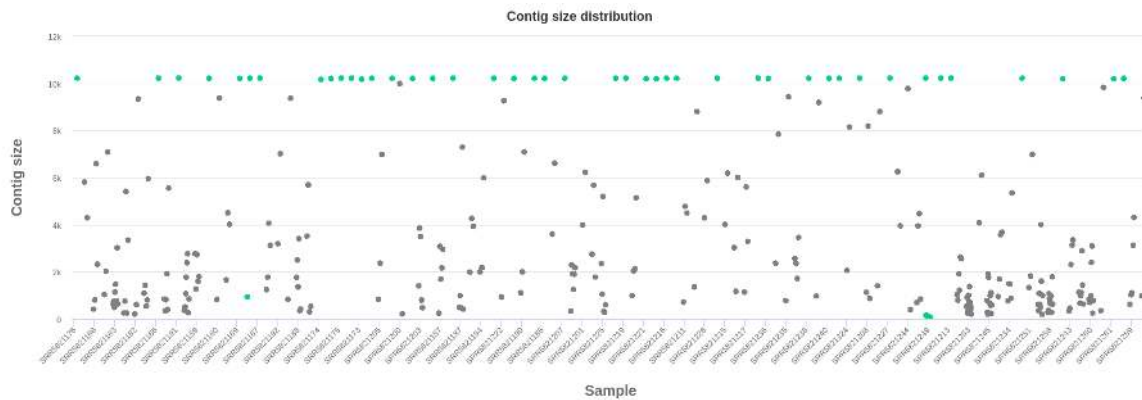
186 The serotype and genotype was successfully determined for the 24 DENV CDSs by BLAST  
 187 (Figure S1, Table 1). The most common were serotype 2 genotype III (AsianAmerican) and serotype  
 188 4 genotype II, with 8 samples each (33.33%), followed by serotype 3 genotype III (n=5, 20.83%),  
 189 serotype 1 genotype V (n=2, 8.33%) and serotype 2 genotype V (Asian I) (n=1, 4.17%). All CDSs  
 190 recovered and the respective closest reference genome in the typing database were aligned and a  
 191 maximum likelihood phylogenetic tree was obtained to visualise the relationship between the samples  
 192 (Figure 3). There was a perfect concordance between the results of serotyping and genotyping and  
 193 the major groups in the tree.

## 194 6.2 The Targeted Metagenomics Dataset

195 To validate DEN-IM's performance in a targeted metagenomics approach, 106 HTS datasets of PCR  
 196 products using primers targeting DENV-3 [24] were analysed. As with the shotgun metagenomics  
 197 dataset, the workflow was executed using the default parameters and directives for resources. This  
 198 time not including the closest reference genome present in the typing database in the multiple  
 199 sequence alignment due to the large number of input samples.

200 No samples failed the quality control block. The proportion of DENV reads ranged from 24.72%  
 201 (SRR5821236) to 99.81% (SRR5821202) of the total processed reads. The samples with less than 70%  
 202 DENV DNA were taxonomic profiled with the Kraken2 [25] with the minikraken2\_v2 database ([ftp://ftp.ccb.jhu.edu/pub/data/kraken2\\_dbs/minikraken2\\_v2\\_8GB\\_201904\\_UPDATE.tgz](ftp://ftp.ccb.jhu.edu/pub/data/kraken2_dbs/minikraken2_v2_8GB_201904_UPDATE.tgz)) and the  
 203 source of the contamination was determined to have come largely from Human DNA (Table S2).

204 Of the 106 samples, 43 (40.60%) managed to assemble a complete CDS sequence (Table S1)  
 205 whereas a mapping approach was used for the remaining 63 samples (59.90%) and a consensus CDS  
 206 was generated. For the assembled CDSs, all but one were assembled with MEGAHIT after not  
 207 producing a full CDS with SPAdes. Moreover, pronounced variation on the size of the assembled  
 208



**Figure 4.** Contig size distribution of the targeted metagenomics dataset. Each dot depicts an assembled DENV contig. Above the 10Kb are full CDS of DENV. Contigs belonging from samples that assembled a complete DENV CDS are highlighted in green, whereas the remaining are coloured in grey.

209 contigs is evident in the contig size distribution plot (Figure 4).

210 All 106 CDSs recovered belonged to serotype 3 genotype III. Despite the same classification,  
211 the maximum likelihood tree indicates that there is detectable genetic diversity within the samples  
212 (Figure 5).

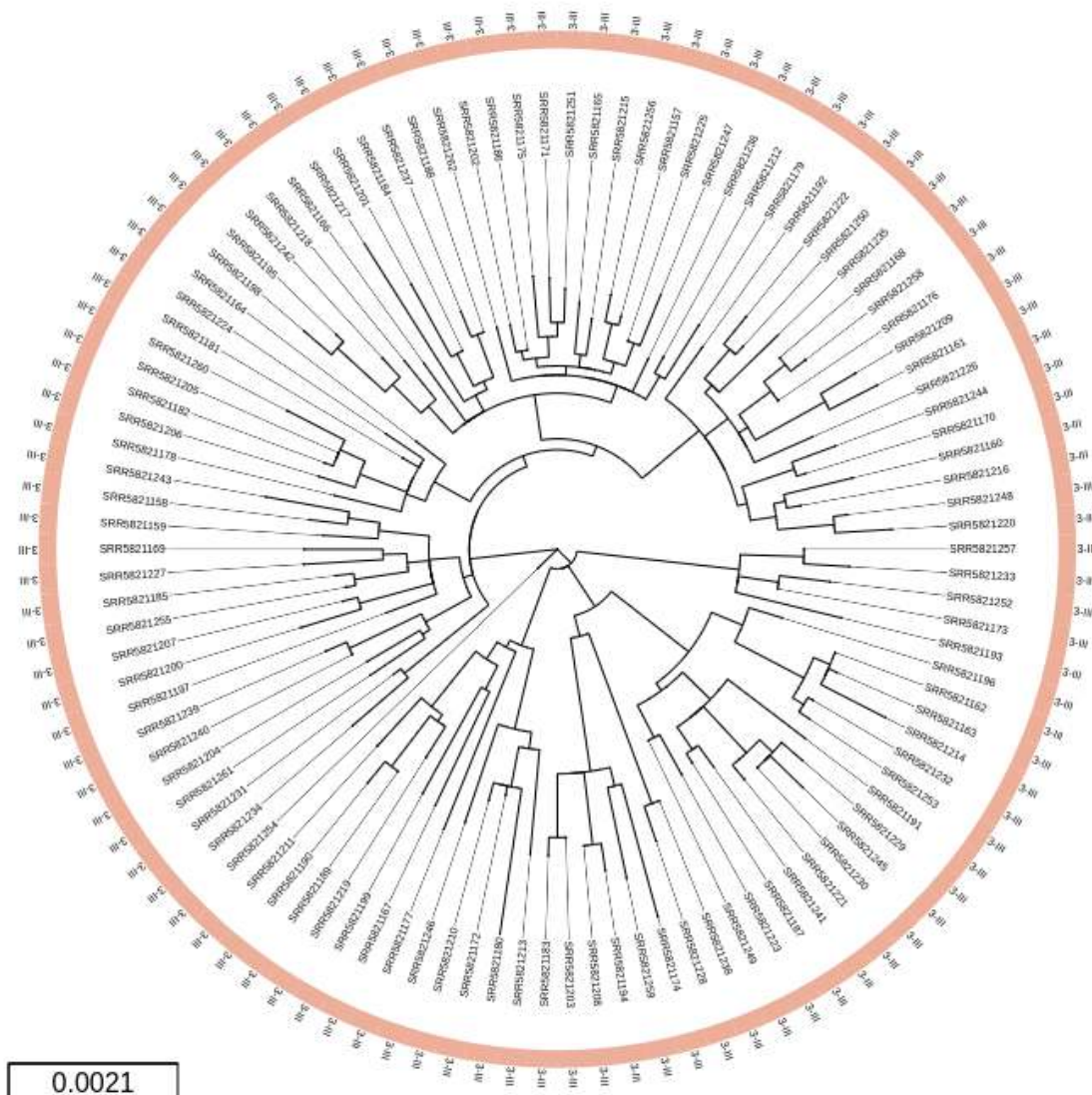
## 213 7 Discussion

214 We have successfully analysed two DENV datasets, one comprising 25 shotgun metagenomics  
215 sequencing data and 106 targeted metagenomics data.

216 In the first dataset, we recovered 24 CDSs from 19 clinical samples, a spiked sample and a positive  
217 control that were correctly serotyped and genotyped. Besides the negative control, 3 samples did not  
218 return typing information due to failing quality checks. In one case (92-1001), no DENV reads after  
219 quality control processing were detected as all the reads contained highly repetitive sequences (AAA;  
220 TTTT) and were filtered out. The two others (91-0115 and UCUG0186) had a low proportion of  
221 DENV reads (0.05% and 0.01%) and an estimated depth of coverage lower than the 10x threshold  
222 criterion (3.17x and 5.65x, respectively). Sequence data of sample CC0186 contained only 960 DENV  
223 reads (0.03%), but these were successfully assembled into a CDS with an estimated depth of coverage  
224 of 14.71x.

225 The proportion of DENV reads was the metagenomic samples is very variable. This may reflect  
226 the viral load in patients in which DENV was detected by PCR. In the spiked sample, containing  
227 4 distinct DENV serotypes, all four were correctly detected despite not being present in equal  
228 concentrations(see Methods, Targeted Metagenomics Sequencing Data). This resulted in different  
229 coverage of each serotype CDS (2032.31 times coverage for DENV-2, 229.02 times coverage for  
230 DENV-1, 76.47 times coverage for DENV-3 and 29.78 times coverage for DENV-4), in accordance  
231 with the ranking order of RT-PCR results. It highlights the potential of the DEN-IM workflow to  
232 accurately detect and recover multiple DENV genomes from samples with DENV co-infection, even  
233 if the serotypes are present in low abundance. Indeed, recent studies from areas of high endemicity  
234 suggest that co-infection with multiple DENV serotypes may frequently occur [26] [27] and the  
235 co-circulation of different DENV strains of the same serotype, but distinct genotypes, in these  
236 areas [26] raises the possibility of simultaneous infection with more than one genotype.

237 When analysing the targeted metagenomics dataset, only 43 CDS out of 106 samples were *de*  
238 *novo* assembled. For the remaining 63 samples, consensus sequences were obtained through mapping.  
239 In all samples DENV 3-III was correctly identified, demonstrating the success of DEN-IM's two  
240 pronged approach of combining assemblers and mapping. We suggest that the lower assembly success  
241 of the targeted metagenomics data may be related to errors during the amplification process resulting  
242 in low quality reads ends which are then trimmed by the quality control block potentially affecting



**Figure 5.** Maximum likelihood circular tree in the DEN-IM report for the 106 complete CDSs obtained with the targeted metagenomics dataset (n=106). All samples belong to serotype 3 genotype III.

243 the assembly process as the overlapping regions are diminished.

244 DEN-IM is built with modularity and containerisation as keystones, leveraging the parallelization  
245 of processes and guaranteeing reproducible analyses across platforms. The modular design allows for  
246 new modules to be easily added and tools that become outdated can be easily updated, ensuring  
247 DEN-IM's sustainability. The software versions are described in the Nextflow script and configuration  
248 files, and in the dockerfiles for each container, allowing traceability of each step of data processing.

249 Being developed in Nextflow, DEN-IM runs on any UNIX-like system and provides out-of-the-box  
250 support for several job schedulers (e.g., PBS, SGE, SLURM) and integration with containerised  
251 software like Docker or Singularity. While it has been developed to be ready to use by non-experts,  
252 not requiring any software installation or parameter tuning, it can still be easily customised through  
253 the configuration files.

254 The interactive HTML reports (Figure S1) provide an intuitive platform for data exploration,  
255 allowing the user to highlight specific samples, filter and re-order the data tables, and export the  
256 plots as needed. Together with the workflow and software containers, a database containing 3830  
257 complete DENV genomes for DENV sequence retrieval and a subset database with 161 curated  
258 DENV genomes for serotyping and genotyping are provided. While constructing these databases, the  
259 obstacles reported by Cuyppers *et al* [28] were apparent, namely the lack of formal definition of a  
260 DENV genotype and the lack of a standardised classification procedure that could assign sequences  
261 to a previously defined genotypic/sub-genotypic clade [28]. Discrepancies between the phylogenetic  
262 relationship and the genotype assignment were frequent and, throughout this study, the classification  
263 of some strains within the ViPR database [29] was updated.

264 As suggested previously [28], further evaluation of the DENV classification will benefit future  
265 research and investigation into the population dynamics of this virus. Our typing approach was  
266 designed to use the currently accepted DENV classification. However, DEN-IM can be easily modified  
267 if a new DENV classification system is to be established in the future.

268 In conclusion, we provide a user-friendly workflow that makes it possible to analyse paired-end raw  
269 sequencing data from shotgun or targeted metagenomics for the presence, typing and phylogenetic  
270 analysis of DENV. The use of containerised workflows, together with shareable reports, will allow an  
271 easier comparison of results globally, promoting collaborations that can benefit the populations where  
272 DENV is endemic. The DEN-IM source code is freely available in the DEN-IM GitHub repository  
273 (<https://github.com/B-UMMI/DEN-IM>), which includes a wiki with full documentation and easy to  
274 follow instructions.

## 275 8 Potential implications

276 The burden of DENV disease is already large, but is still increasing as the risk of exposure to the  
277 virus is increasing, not only through travel to endemic areas but also due to the expansion of the  
278 geographic areas of the mosquito vectors and the disease [30].

279 The decreasing costs and wider availability of HTS makes it an ideal technology to monitor  
280 DENV transmission, including the direct processing of patient samples, either through the use of a  
281 more affordable targeted metagenomics approach or through shotgun metagenomics. Either way, a  
282 ready to use bioinformatics workflow that enables the reproducible analysis of DENV is particularly  
283 relevant.

284 DEN-IM was designed to perform a comprehensive analysis without the requirement of extensive  
285 bioinformatics expertise in order to generate either assemblies or consensus of full DENV CDSs and  
286 to identify the serotype and genotype of the DENV present in the sample. Although we did not  
287 exhaustively test the capacity of DEN-IM to detect co-infection with multiple DENV serotypes and  
288 genotypes, all four genotypes present in the spiked sample were accurately detected. This raises  
289 the possibility that DEN-IM can play a role in the identification of these cases whose prevalence  
290 is increasingly appreciated in highly endemic areas. Moreover, although being ready-to-use, the  
291 DEN-IM workflow can be easily customised to optimise the data analysis.

292 DEN-IM enables reproducible and collaborative research, benefiting a wide group of researchers  
293 regardless of their computational expertise and resources available.



## 9 Methods

### 9.1 DENV Reference Database

We have compiled a database of 3830 complete DENV genomes obtained from the NIAID Virus Pathogen Database and Analysis Resource (ViPR) in January 2019 [29] (<http://www.viprbrc.org/>). The sequences were distributed unevenly throughout the four DENV serotypes, with DENV-1 being the most represented with 1636 sequences (42.72%), followed by DENV-2 with 1067 sequences (27.86%), DENV-3 with 807 sequences (21.07%), and DENV-4 with 320 sequences (8.36%). The selection criteria for the search were as follows: a) complete genome sequence only, b) human host only, c) collection year (1950-2018). Data available from all countries was included and duplicated sequences were removed and only the sequences with sub-type data were kept. A representative of DENV serotype 1 genotype III was introduced (EF457905, recovered from monkey) as no representatives were available with the search criteria used. This genotype is Sylvatic and considered extinct [31] [32]. Additionally, any sample with IUPAC codes in the sequence provided were excluded.

In order to recover the maximum number of DENV reads from the input HTS data in the first mapping step (Figure 1), we maintained the database with the 3830 complete DENV genomes to retain as much diversity as possible. This database is referred as “DENV mapping database”.

For typing purposes, overly similar sequences in the collection were removed from the database by clustering the sequences in each serotype at 98% nucleotide similarity with CD-HIT [33], leaving 161 representative sequences of all described DENV serotypes and genotypes, with 46 DENV-1 sequences (Table S4), 63 DENV-2 (Table S5), 25 DENV-3 (Table S6) and 27 DENV-4 (Table S7). This database is referred as “DENV typing database”. This step is necessary to speed up the classification step for genotyping.

Phylogenetic analysis of typing collection was performed by aligning the reference genomes with MAFFT [20], in auto mode and with automatic sequence orientation adjustment. A phylogenetic tree was inferred with RAxML (version 8.12.11) [21] using the GTR- $\Gamma$  substitution model and 500 times bootstrap. The resulting trees are available as supplemental material (Figures S2 to S5).

The sequence JF459993 from in the DENV-1 collection, as of April 2019, was annotated in ViPR as belonging to genotype IV, but in our analysis it clustered within genotype I clade (Figure S2). The classification of DENV-1 I was also obtained from GenomeDetective Dengue Typing Tool (<https://www.genomedetective.com/app/typingtool/dengue/>), so we proceeded to alter the annotation of this particular sample (Table S4). In order to harmonise dengue nomenclature, the system adopted uses Roman-numeric labels to identify the genotype, with the exception of Serotype 2 (Table S5), which used both Roman-numeric and geographic origin due to the widespread adoption of the latter.

### 9.2 Workflow Parameters

The short-read paired-end data is passed as input through the “-fastq” parameter, that by default is set to match all files in the “fastq” folder that match the pattern “\*\_R1,2\*”. In the process to verify the integrity of the paired-end raw sequencing data, the integrity of the input files is assessed by attempting to decompress and read the files. An estimation of the depth of coverage is also performed. By default, the input size (“-genomeSize”) is set to 0.012 Mb and the minimum coverage depth (“-minCoverage”) is set to 10. If any input file is found to be corrupt, its progression in the workflow is aborted.

In the FastQC and Trimmomatic module, FastQC (<https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>) is run with the parameters “-extract -nogroup -format fastq”. FastQC will inform Trimmomatic [10] on how many bases to trim from the 3' and 5' ends of the raw reads. By default, Trimmomatic uses the default set of Illumina adapters provided with the workflow but this behaviour can be overwritten with the “-adapters” parameter. The additional Trimmomatic parameters “-trimSlidingWindow”, “-trimLeading”, “-trimTrailing” and “-trimMinLength” can all be set to different values.

The removal of low complexity sequences is done with PrinSeq [12] using a custom parameter (“-pattern”), which by default is set to the value “A 50%; T 50%; N 50%”, removing sequences whose content is at least half composed of a polymeric sequence (A, T or N).

346 To retrieve the reads that map to the DENV reference database, Bowtie2 [13] is run with default  
347 parameters with the DENV mapping database as a reference. The reads and their mates that map  
348 to the reference are retrieved with "*samtools view -buh -F 12*" and "*samtools fastq*" commands.  
349 The DENV mapping database can be altered with the "*-reference*" parameter, or alternatively, a  
350 Bowtie2 index can be provided with the "*-index*" parameter. This allows for the workflow to work  
351 with other databases obtained through public and owned DENV genomes. The coverage estimation  
352 step is performed on the retrieved DENV reads with the same parameters are the first estimation  
353 ("*--genomeSize=0.012*" and "*--minCoverage=10*").

354 In the assembly process, the retrieved DENV reads are firstly assembled with SPAdes Genome  
355 Assembler [15] with the options "*-careful -only-assembler -cov-cutoff*". The coverage cutoff if  
356 dictated by the "*-spadesMinCoverage*" and "*-spadesMinKmerCoverage*" parameters, set to 2 by  
357 default. If the assembly with SPAdes fails to produce a contig equal or greater than the value  
358 defined in the "*-minimumContigSize*" parameter (default of 10000), the data is re-assembled with the  
359 MEGAHIT assembler [16] with default parameters. By default the k-mers to be used in the assembly  
360 in both tools ("*-spadesKmers*" and "*-megahitKmers*") are automatically determined depending on  
361 the read size. If the maximum read length is equal or greater than 175 nucleotides, the assembly is  
362 done with the k-mers "55, 77, 99, 113, 127", otherwise the k-mers "21, 33, 55, 67, 77" are used.

363 To correct the assemblies produced, the Pilon tool [17] is run after mapping the QC'ed reads  
364 back to the assembly with Bowtie2 and "*samtools sort*". This process also verifies the coverage and  
365 the number of contigs produced in the assembly. The behaviour can be altered with the parameters  
366 "*-minAssemblyCoverage*", "*-AMaxContigs*" and "*-genomeSize*", set to "auto", 1000 and 0.01 Mb by  
367 default. The first parameter, when set to 'auto', the minimum assembly coverage for each contig  
368 required is set to the 1/3 of the assembly mean coverage or to a minimum of 10x. The ratio of contig  
369 number per genome MB is calculated based on the genome size estimation for the samples.

370 The contigs larger than the value defined in the "*-size*" parameter (default of 10000 nucleotides)  
371 are considered to be complete CDSs and follow the rest to the workflow independently. If no complete  
372 CDS is recovered, the QC'ed read data is passed to the mapping to module that does the DENV  
373 typing database and consensus generation.

374 The serotyping and genotyping is performed with the Seq\_Typing tool [18] with the command  
375 "*seq\_typing.py assembly*" or "*seq\_typing.py reads*", using as reference the provided curated DENV  
376 typing database. It is possible to retrieve the genomes of the closest references and include them in  
377 the downstream analysis by changing the "*-get.reference*" option to "true". By default this is not  
378 included in the analysis.

379 The CDSs, and the reference sequences if requested, are aligned with the MAFFT tool [20] with  
380 the options "*-adjustdirection -auto*" and a maximum likelihood phylogenetic tree is obtained with  
381 the RAxML tool [21] with the options "*-p 12345 -f -a*". Additionally and by default, the substitution  
382 model ("*-substitutionModel*") is set to "GTRGAMMA", the bootstrap is set to 500 ("*-bootstrap*")  
383 and the seed to "12345" ("*-seedNumber*").

### 384 9.3 Shotgun Metagenomics Sequencing Data

385 Samples of plasma (n=9) and serum samples (n=13) from confirmed dengue symptomatic patients  
386 were collected in Venezuela between 2010-2015 (Table S3) (see Availability of supporting data and  
387 materials). DENV positivity was confirmed by either RT-qPCR [34] or nested RT-PCR [35].

388 As a positive control sample, the supernatant of a viral culture containing DENV-2 strain 16681  
389 was used. The negative control sample consisted of DNA- and RNA-free water (Sigma-Aldrich, St.  
390 Louis, MO, USA).

391 A spiked sample was produced consisting of a mixture of four 5  $\mu$ l of cDNA isolated from clinical  
392 samples including all DENV serotypes (DENV-1 to -4). The viral cDNA for these samples was not in  
393 equal concentration and the viral copy number in the clinical samples was assessed by RT-PCR [35].  
394 The results were as follow: DENV-2 with 1070000 copies/ $\mu$ l, DENV-1 with 117830 copies/ $\mu$ l, DENV-3  
395 with 44300 copies/ $\mu$ l and DENV-4 with 6600 copies/ $\mu$ l.

396 The cDNA libraries were generated using either the NEBNext® RNA First and Second strand  
397 modules and the Nextera XT DNA library preparation kit (NXT), or the TruSeq RNA V2 library  
398 preparation kit (TS). The libraries were sequenced in MiSeq and NextSeq instruments using 300-cycles



399 v2 paired-end cartridges.

400 The DEN-IM workflow was executed with the raw sequencing data using the default parameters  
401 and resources in an HPC cluster with 300 Cores/600 Threads of Processing Power and 3 TB RAM  
402 divided through 15 computational nodes, 9 with 254 GB Ram and 6 with 126GB RAM.

## 403 9.4 Targeted Metagenomics Sequencing Data

404 The accession numbers for the 106 DENV-3 amplicon sequencing paired-end short-read datasets  
405 are available under BioProject PRJNA394021. The list of Run Accession IDs were obtained  
406 with NCBI's RunSelector and the raw data was downloaded with the GetSeqENA tool (<https://github.com/B-UMMI/getSeqENA>).

407  
408 The DEN-IM workflow was executed with the raw sequencing data with default parameters and  
409 resources in the same HPC cluster as the shotgun metagenomics dataset.

## 410 10 Availability of source code and requirements

411 Lists the following:

- 412 • Project name: DEN-IM
- 413 • Project home page: <https://github.com/B-UMMI/DEN-IM>
- 414 • Operating system(s): UNIX-like systems.
- 415 • Programming language: Nextflow, Python, Bash
- 416 • Other requirements: Java version 8 or highest. Docker/Singularity/Shifter
- 417 • License: GNU GPL v3
- 418 • Documentation and tutorials: <https://github.com/B-UMMI/DEN-IM/wiki>

## 419 11 Availability of supporting data and materials

420 The 106 DENV-3 targeted metagenomics sequencing paired-end short-read datasets are available  
421 under BioProject PRJNA394021. The 25 shotgun metagenomics dataset is available under Bioproject  
422 PRJNA474413 The accession number for all the samples in the shotgun metagenomics dataset are  
423 available in the supplemental material (Table S3).

## 424 12 Declarations

### 425 12.1 List of abbreviations

426 DENV: Dengue Virus; CDS: Coding Sequence; ORF: Open Reading Frame; NCR: Non-Coding  
427 Region; HPC: High-Performance Computing; HTS: High Throughput Sequencing; QC: Quality  
428 Control

### 429 12.2 Ethical Approval (optional)

430 This study followed international standards for the ethical conduct of research involving human  
431 subjects. Data and sample collection was carried out within the DENVEN and IDAMS (International  
432 Research Consortium on Dengue Risk Assessment, Management and Surveillance) projects. The  
433 study was approved by the Ethics Review Committee of the Biomedical Research Institute, Carabobo  
434 University (Aval Bioetico #CBIIB(UC)-014 and CBIIB-(UC)-2013-1), Maracay, Venezuela; the  
435 Ethics, Bioethics and Biodiversity Committee (CEBioBio) of the National Foundation for Science,  
436 Technology and Innovation (FONACIT) of the Ministry of Science, Technology and Innovation,  
437 Caracas, Venezuela; the regional Health authorities of Aragua state (CORPOSALUD Aragua) and

438 Carabobo State (INSALUD); and by the Ethics Committee of the Medical Faculty of Heidelberg  
439 University and the Oxford University Tropical Research Ethics Committee.

### 440 **12.3 Consent for publication**

441 Not applicable.

### 442 **12.4 Competing Interests**

443 The authors declare that they have no competing interests.

### 444 **12.5 Funding**

445 C.I.M. was supported by the Fundação para a Ciência e Tecnologia (grant SFRH/BD/129483/2017).  
446 Erley Lizarazo received the Abel Tasman Talent Program grant from the UMCG, University of  
447 Groningen, Groningen, The Netherlands. This work was partly supported by the ONEIDA project  
448 (LISBOA-01-0145-FEDER-016417) co-funded by FEEL–Fundos Europeus Estruturais e de Investi-  
449 mento from Programa Operacional Regional Lisboa 2020 and by national funds from FCT–Fundação  
450 para a Ciência e a Tecnologia and UID/BIM/50005/2019, and by UID/BIM/50005/2019, project  
451 funded by Fundação para a Ciência e a Tecnologia (FCT)/ Ministério da Ciência, Tecnologia e Ensino  
452 Superior (MCTES) through Fundos do Orçamento de Estado.

### 453 **12.6 Author’s Contributions**

454 C.I.M., E.L., N.C., M.R., J.A.C. and J.W.A.R. designed the workflow. C.I.M implemented and  
455 optimised the workflow, created the Docker containers, and wrote the manuscript. M.P.M. imple-  
456 mented the DENV genotyping module in the workflow and D.N.S. contributed to the development of  
457 DEN-IM’s HTML report. E.L., A. T., and N.C. provided the shotgun metagenomics data used to  
458 test and validate the workflow and wrote the manuscript. A.T., N.C., M.R., J.A.C. and J.W.A.R.  
459 critically revised the article. All authors read, commented on, and approved the final manuscript.

## 460 **13 Acknowledgements**

461 The authors would like to thank Tiago F. Jesus and Bruno Ribeiro-Gonçalves for their invaluable  
462 help with the Nextflow implementation. We would also like to thank Erwin C. Raangs from the  
463 UMCG for his assistance in the sequencing of the shotgun metagenomics dataset. Additionally, the  
464 authors thank Lize Cuyppers, Krystof Theys, Pieter Libin and Gilberto Santiago for their assistance  
465 in DENV nomenclature and classification. This work was done in collaboration with the ESCMID  
466 Study Group on Molecular and Genomic Diagnostics (ESGMD), Basel, Switzerland.

## 467 **References**

- 468 1. World Health Organization, “Dengue: guidelines for diagnosis, treatment, prevention, and  
469 control,” *Special Programme for Research and Training in Tropical Diseases*, pp. x, 147, 2009.  
470 [Online]. Available: [http://whqlibdoc.who.int/publications/2009/9789241547871\\_{\\_}eng.pdf](http://whqlibdoc.who.int/publications/2009/9789241547871_{_}eng.pdf)
- 471 2. M. S. Diamond and T. C. Pierson, “Molecular Insight into Dengue Virus Pathogenesis and Its  
472 Implications for Disease Control,” *Cell*, vol. 162, no. 3, pp. 488–492, 2015. [Online]. Available:  
473 <http://dx.doi.org/10.1016/j.cell.2015.07.005>
- 474 3. S. Bhatt, P. W. Gething, O. J. Brady, J. P. Messina, A. W. Farlow, C. L. Moyes, J. M.  
475 Drake, J. S. Brownstein, A. G. Hoen, O. Sankoh, M. F. Myers, D. B. George, T. Jaenisch,  
476 G. R. William Wint, C. P. Simmons, T. W. Scott, J. J. Farrar, and S. I. Hay, “The global  
477 distribution and burden of dengue,” *Nature*, vol. 496, no. 7446, pp. 504–507, 2013. [Online].  
478 Available: <http://dx.doi.org/10.1038/nature12060>

- 479 4. J. Lourenço, W. Tennant, N. R. Faria, A. Walker, S. Gupta, and M. Recker, “Challenges in  
480 dengue research: A computational perspective,” *Evolutionary Applications*, vol. 11, no. 4, pp.  
481 516–533, apr 2018. [Online]. Available: <http://doi.wiley.com/10.1111/eva.12554>
- 482 5. K. C. Leitmeyer, D. W. Vaughn, D. M. Watts, R. Salas, I. Villalobos, de Chacon,  
483 C. Ramos, and R. Rico-Hesse, “Dengue virus structural differences that correlate with  
484 pathogenesis.” *Journal of virology*, vol. 73, no. 6, pp. 4738–47, jun 1999. [Online]. Available:  
485 <http://www.ncbi.nlm.nih.gov/pubmed/10233934>  
486 <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC112516>  
487 <http://www.ncbi.nlm.nih.gov/pubmed/10233934>  
<http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC112516>
- 488 6. N. L. Yozwiak, P. Skewes-Cox, M. D. Stenglein, A. Balmaseda, E. Harris, and J. L. DeRisi,  
489 “Virus identification in unknown tropical febrile illness cases using deep sequencing,” *PLoS*  
490 *Neglected Tropical Diseases*, vol. 6, no. 2, 2012.
- 491 7. C. K. Lee, C. W. Chua, L. Chiu, and E. S.-C. Koay, “Clinical use of targeted high-throughput  
492 whole-genome sequencing for a dengue virus variant,” *Clinical Chemistry and Laboratory*  
493 *Medicine (CCLM)*, vol. 55, no. 9, p. e209, jan 2017. [Online]. Available: [https://](https://www.degruyter.com/view/j/cclm.2017.55.issue-9/cclm-2016-0660/cclm-2016-0660.xml)  
494 [www.degruyter.com/view/j/cclm.2017.55.issue-9/cclm-2016-0660/cclm-2016-0660.xml](https://www.degruyter.com/view/j/cclm.2017.55.issue-9/cclm-2016-0660/cclm-2016-0660.xml)  
495 <http://www.degruyter.com/view/j/cclm.2017.55.issue-9/cclm-2016-0660/cclm-2016-0660.xml>
- 496 8. Z. Fatima, M. Idrees, M. A. Bajwa, Z. Tahir, O. Ullah, M. Q. Zia, A. Hussain, M. Akram,  
497 B. Khubai, S. Afzal, S. Munir, S. Saleem, B. Rauff, S. Badar, M. Naudhani, S. Butt,  
498 M. Aftab, L. Ali, and M. Ali, “Serotype and genotype analysis of dengue virus by  
499 sequencing followed by phylogenetic analysis using samples from three mini outbreaks-  
500 2007-2009 in Pakistan,” *BMC Microbiology*, vol. 11, no. 1, p. 200, 2011. [Online]. Available:  
501 <http://www.biomedcentral.com/1471-2180/11/200>
- 502 9. P. Di Tommaso, M. Chatzou, E. W. Floden, P. P. Barja, E. Palumbo, and C. Notredame,  
503 “Nextflow enables reproducible computational workflows,” *Nature Biotechnology*, vol. 35, p.  
504 316, apr 2017. [Online]. Available: <http://dx.doi.org/10.1038/nbt.3820>  
505 <http://10.0.4.14/nbt.3820>  
<https://www.nature.com/articles/nbt.3820#supplementary-information>
- 506 10. A. M. Bolger, M. Lohse, and B. Usadel, “Trimmomatic: A flexible trimmer for Illumina  
507 sequence data,” *Bioinformatics*, vol. 30, no. 15, pp. 2114–2120, 2014.
- 508 11. B. Ewing, L. Hillier, M. C. Wendl, and P. Green, “Base-Calling of Automated Sequencer  
509 Traces Using Phred. I. Accuracy Assessment,” *Genome Research*, vol. 8, no. 3, pp. 175–185,  
510 mar 1998. [Online]. Available: <http://genome.cshlp.org/lookup/doi/10.1101/gr.8.3.175>
- 511 12. R. Schmieder and R. Edwards, “Quality control and preprocessing of metagenomic datasets,”  
512 *Bioinformatics*, vol. 27, no. 6, pp. 863–864, 2011.
- 513 13. B. Langmead and S. L. Salzberg, “Fast gapped-read alignment with Bowtie  
514 2,” *Nature Methods*, vol. 9, no. 4, pp. 357–359, apr 2012. [Online]. Available:  
515 <http://www.nature.com/articles/nmeth.1923>
- 516 14. H. Li, “A statistical framework for SNP calling, mutation discovery, association mapping and  
517 population genetical parameter estimation from sequencing data,” *Bioinformatics*, vol. 27,  
518 no. 21, pp. 2987–2993, 2011.
- 519 15. A. Bankevich, S. Nurk, D. Antipov, A. A. Gurevich, M. Dvorkin, A. S. Kulikov,  
520 V. M. Lesin, S. I. Nikolenko, S. Pham, A. D. Prjibelski, A. V. Pyshkin, A. V.  
521 Sirotkin, N. Vyahhi, G. Tesler, M. A. Alekseyev, and P. A. Pevzner, “SPAdes: A New  
522 Genome Assembly Algorithm and Its Applications to Single-Cell Sequencing,” *Journal*  
523 *of Computational Biology*, vol. 19, no. 5, pp. 455–477, may 2012. [Online]. Available:  
524 <http://online.liebertpub.com/doi/abs/10.1089/cmb.2012.0021>

- 525 16. D. Li, C. M. Liu, R. Luo, K. Sadakane, and T. W. Lam, “MEGAHIT: An ultra-fast single-  
526 node solution for large and complex metagenomics assembly via succinct de Bruijn graph,”  
527 *Bioinformatics*, vol. 31, no. 10, pp. 1674–1676, 2015.
- 528 17. B. J. Walker, T. Abeel, T. Shea, M. Priest, A. Abouelliel, S. Sakthikumar, C. A. Cuomo,  
529 Q. Zeng, J. Wortman, S. K. Young, and A. M. Earl, “Pilon: An integrated tool for compre-  
530 hensive microbial variant detection and genome assembly improvement,” *PLoS ONE*, vol. 9,  
531 no. 11, 2014.
- 532 18. M. P. Machado, B. Ribeiro-Gonçalves, M. Silva, M. Ramirez, and J. A. Carriço,  
533 “Epidemiological Surveillance and Typing Methods to Track Antibiotic Resistant Strains Using  
534 High Throughput Sequencing.” *Methods in molecular biology (Clifton, N.J.)*, vol. 1520, pp.  
535 331–356, 2017. [Online]. Available: <http://www.ncbi.nlm.nih.gov/pubmed/27873262>
- 536 19. S. F. Altschul, T. L. Madden, A. A. Schäffer, J. Zhang, Z. Zhang, W. Miller, and D. J. Lipman,  
537 “Gapped BLAST and PSI-BLAST: A new generation of protein database search programs,”  
538 *Nucleic Acids Research*, vol. 25, no. 17, pp. 3389–3402, 1997.
- 539 20. T. Nakamura, K. D. Yamada, K. Tomii, and K. Katoh, “Parallelization of MAFFT for  
540 large-scale multiple sequence alignments,” *Bioinformatics*, vol. 34, no. March, pp. 2490–2492,  
541 2018. [Online]. Available: [https://academic.oup.com/bioinformatics/advance-article/doi/10.](https://academic.oup.com/bioinformatics/advance-article/doi/10.1093/bioinformatics/bty121/4916099)  
542 [1093/bioinformatics/bty121/4916099](https://academic.oup.com/bioinformatics/advance-article/doi/10.1093/bioinformatics/bty121/4916099)
- 543 21. A. Stamatakis, “RAxML version 8: A tool for phylogenetic analysis and post-analysis of large  
544 phylogenies,” *Bioinformatics*, vol. 30, no. 9, pp. 1312–1313, 2014.
- 545 22. L. Gerhardt, W. Bhimji, S. Canon, M. Fasel, D. Jacobsen, M. Mustafa, J. Porter, and  
546 V. Tsulaia, “Shifter: Containers for HPC,” *Journal of Physics: Conference Series*, vol. 898, p.  
547 082021, oct 2017. [Online]. Available: [http://stacks.iop.org/1742-6596/898/i=8/a=082021?](http://stacks.iop.org/1742-6596/898/i=8/a=082021?key=crossref.b7268cc937fc3b29093062a6749fbbbf)  
548 [key=crossref.b7268cc937fc3b29093062a6749fbbbf](http://stacks.iop.org/1742-6596/898/i=8/a=082021?key=crossref.b7268cc937fc3b29093062a6749fbbbf)
- 549 23. G. M. Kurtzer, V. Sochat, and M. W. Bauer, “Singularity: Scientific containers for mobility of  
550 compute,” *PLoS ONE*, vol. 12, no. 5, pp. 1–20, 2017.
- 551 24. P. Parameswaran, C. Wang, S. B. Trivedi, M. Eswarappa, M. Montoya, A. Balmaseda, and  
552 E. Harris, “Intrahost Selection Pressures Drive Rapid Dengue Virus Microevolution in Acute  
553 Human Infections,” *Cell Host & Microbe*, vol. 22, no. 3, pp. 400–410.e5, sep 2017. [Online].  
554 Available: <https://doi.org/10.1016/j.chom.2017.08.003>
- 555 25. D. E. Wood and S. L. Salzberg, “Kraken: ultrafast metagenomic sequence classification  
556 using exact alignments,” *Genome Biology*, vol. 15, no. 3, p. R46, mar 2014. [Online].  
557 Available: <https://doi.org/10.1186/gb-2014-15-3-r46>[http://genomebiology.biomedcentral.com/](http://genomebiology.biomedcentral.com/articles/10.1186/gb-2014-15-3-r46)  
558 [articles/10.1186/gb-2014-15-3-r46](http://genomebiology.biomedcentral.com/articles/10.1186/gb-2014-15-3-r46)
- 559 26. P. E. S. Marinho, D. B. De Oliveira, T. M. S. Candiani, A. P. C. Crispim, P. P. M. Alvarenga,  
560 F. C. d. S. Castro, J. S. Abrahão, M. Rios, R. S. Coimbra, and E. G. Kroon, “Meningitis  
561 associated with simultaneous infection by multiple dengue virus serotypes in children, Brazil,”  
562 *Emerging Infectious Diseases*, vol. 23, no. 1, pp. 115–118, 2017.
- 563 27. M. N. Reddy, R. Dungdung, L. Valliyott, and R. Pilankatta, “Occurrence of concurrent  
564 infections with multiple serotypes of dengue viruses during 2013–2015 in northern Kerala,  
565 India,” *PeerJ*, vol. 5, p. e2970, 2017.
- 566 28. L. Cuypers, P. Libin, P. Simmonds, A. Nowé, J. Muñoz-Jordán, L. Alcantara, A.-M.  
567 Vandamme, G. Santiago, and K. Theys, “Time to Harmonize Dengue Nomenclature  
568 and Classification,” *Viruses*, vol. 10, no. 10, p. 569, oct 2018. [Online]. Available:  
569 <http://www.mdpi.com/1999-4915/10/10/569>

- 570 29. B. E. Pickett, D. S. Greer, Y. Zhang, L. Stewart, L. Zhou, G. Sun, Z. Gu, S. Kumar, S. Zaremba,  
571 C. N. Larsen, W. Jen, E. B. Klem, and R. H. Scheuermann, “Virus pathogen Database and  
572 Analysis Resource (ViPR): A comprehensive bioinformatics Database and Analysis Resource  
573 for the Coronavirus research community,” *Viruses*, vol. 4, no. 11, pp. 3209–3226, 2012.
- 574 30. T. Pang, T. K. Mak, and D. J. Gubler, “Prevention and control of dengue—the light at the  
575 end of the tunnel,” *The Lancet Infectious Diseases*, vol. 17, no. 3, pp. e79–e87, 2017.
- 576 31. C. J. Villabona-Arenas and P. M. d. A. Zanutto, “Worldwide Spread of Dengue Virus Type 1,”  
577 *PLoS One*, vol. 8, no. 5, p. e62649, 2013.
- 578 32. N. Vasilakis and S. C. Weaver, “Chapter 1 The History and Evolution of  
579 Human Dengue Emergence,” in *Advances in Virus Research*, ser. Advances in Virus  
580 Research. Academic Press, 2008, vol. 72, pp. 1–76. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0065352708004016>  
581
- 582 33. W. Li and A. Godzik, “Cd-hit : a fast program for clustering and comparing large sets of  
583 protein or nucleotide sequences,” *Bioinformatics*, vol. 22, no. 13, pp. 1658–1659, 2006. [Online].  
584 Available: <https://academic.oup.com/bioinformatics/article/22/13/1658/194225>
- 585 34. CDC, “CDC DENV-1-4 Real-Time RT-PCR Assay for Detection and Serotype Identification  
586 of Dengue Virus,” pp. 1–53, 2013. [Online]. Available: [http://www.cdc.gov/dengue/resources/  
587 rt{-}pcr/CDCPackageInsert.pdf](http://www.cdc.gov/dengue/resources/rt{-}pcr/CDCPackageInsert.pdf).
- 588 35. R. Lanciotti, C. Calisher, D. Gubler, G. Chang, and A. Vorndam, “Rapid detection and  
589 typing of dengue viruses from clinical samples by using reverse transcriptase-polymerase  
590 chain reaction.” *Journal of Clinical Microbiology*, vol. 30, no. 3, pp. 545–551, 1992. [Online].  
591 Available: <http://jcm.asm.org/content/30/3/545.short>

a)

Quality control

Search ID column

| ID              | Raw BP<br>integrity coverage 1 1 | Reads<br>integrity coverage 1 1 | Coverage<br>integrity coverage 1 1 | Trimmed (%)<br>trimmomatic 1 2 | Coverage check<br>coverage 1 6 |
|-----------------|----------------------------------|---------------------------------|------------------------------------|--------------------------------|--------------------------------|
| cc030b_S21      | 399040452                        | 2642652                         | 33253.37                           | 60.28                          | 677.96                         |
| UCUG0186        | 7630478                          | 64442.24                        | 19.72                              | 2313.39                        | 5.65                           |
| 91-0115_S7_L001 | 179244760                        | 1333220                         | 14937.06                           | 32.58                          | 3.74                           |
| 91-0109_S4_L001 | 91710149                         | 656462                          | 7642.51                            | 4.23                           | 1287.52                        |
| CC0066          | 1087654460                       | 13700000                        | 90521.21                           | 47.98                          | 569.7                          |
| CC0067          | 1022064484                       | 10464336                        | 85172.04                           | 19.27                          | 2548.84                        |
| CC0061          | 1262837603                       | 12935424                        | 105236.47                          | 19.15                          | 5120.4                         |
| 91-0118_S8_L001 | 195287140                        | 1423414                         | 16272.26                           | 53.42                          | 66.53                          |

Previous Page 2 of 3 10 rows Next

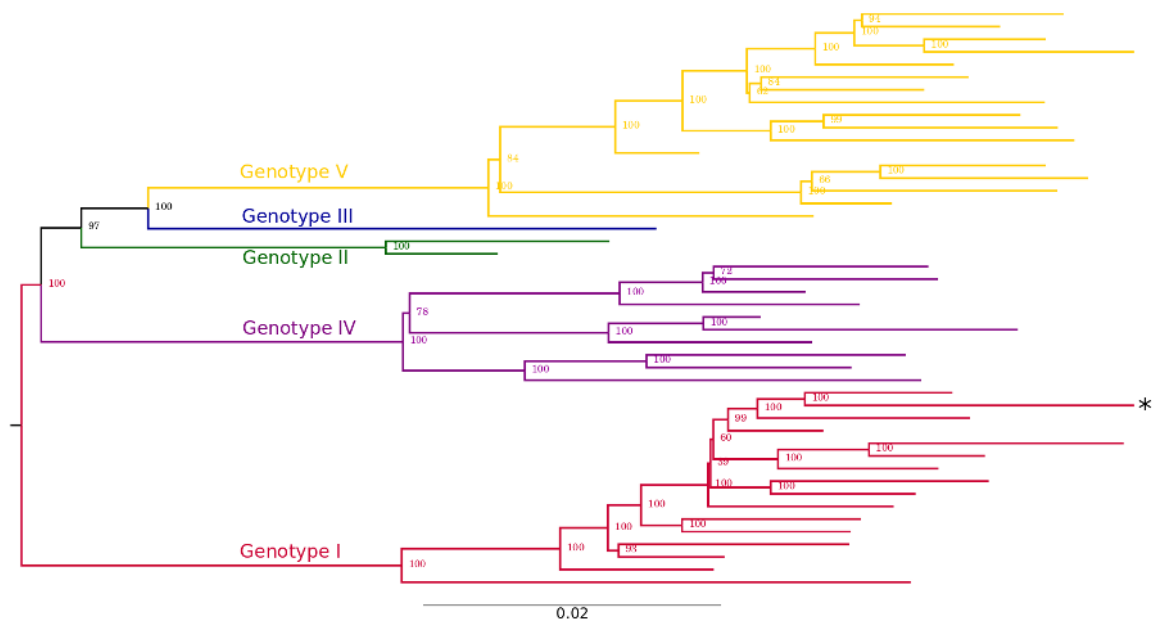
Current selection: 0

b)

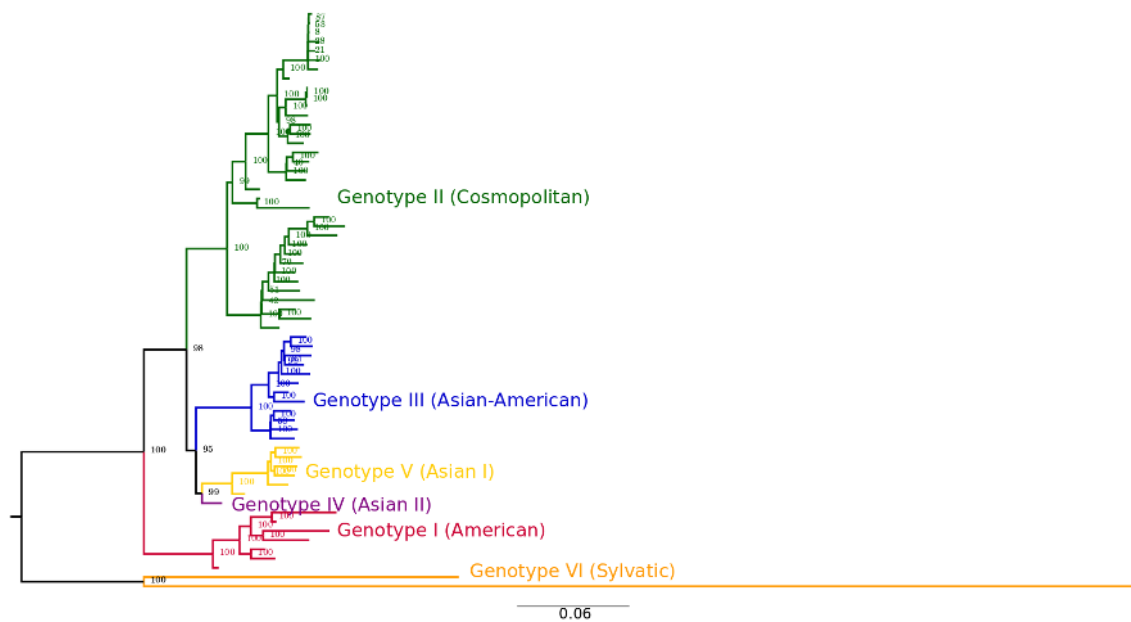
| ID   | serotyping<br>denque_typing_assembly_1_11 | identity<br>denque_typing_assembly_1_1 | Coverage<br>denque_typing_assembly_1_11 | Reference<br>denque_typing_assembly_1_11 |
|--|---|--|---|--|
| Spike_NODE_3_length_10196_cov_229.022822_pilon                         | 1-V                                       | 98.03                                  | 100                                     | gb:EU482561                              |
| 91-0132_S6_L001_NODE_1_length_10217_cov_2041.464103_pilon              | 1-V                                       | 98.03                                  | 100                                     | gb:EU482561                              |
| CC0031_k77_16_flag_0_multi_50891.9804_len_10085_pilon                  | 2-III(AsianAmerican)                      | 99.21                                  | 98.95                                   | gb:FJ024473                              |
| cc0070_S8_L001_NODE_1_length_10200_cov_119.535810_pilon                | 2-III(AsianAmerican)                      | 99.22                                  | 100                                     | gb:FJ024473                              |
| 91-0106_S2_L001_NODE_1_length_10207_cov_218.928825_pilon               | 2-III(AsianAmerican)                      | 98.72                                  | 100                                     | gb:FJ024473                              |
| Spike_NODE_4_length_10192_cov_76.477014_pilon                          | 2-III(AsianAmerican)                      | 98.65                                  | 100                                     | gb:FJ024473                              |
| CC0150_NODE_1_length_10242_cov_3878.632858_pilon                       | 2-III(AsianAmerican)                      | 99.13                                  | 100                                     | gb:FJ024473                              |
| 91-0109_S4_L001_NODE_1_length_10219_cov_652.125222_pilon               | 2-III(AsianAmerican)                      | 98.88                                  | 100                                     | gb:FJ024473                              |
| 91-0104_NODE_1_length_10181_cov_306.327573_pilon                       | 2-III(AsianAmerican)                      | 98.72                                  | 100                                     | gb:FJ024473                              |
| 92-1084_NODE_1_length_10194_cov_816.395572_pilon                       | 2-III(AsianAmerican)                      | 98.67                                  | 100                                     | gb:FJ024473                              |
| Positivecontrol_S21_L001_k77_1_flag_1_multi_19626.0647_len_10227_pilon | 2-V(Asian)                                | 100                                    | 100                                     | gb:GQ888591                              |
| CC0011_NODE_1_length_10201_cov_807.828724_pilon                        | 3-III                                     | 98.7                                   | 100                                     | gb:EU687233                              |
| Spike_NODE_1_length_10266_cov_2032.312101_pilon                        | 3-III                                     | 98.38                                  | 100                                     | gb:EU687233                              |
| CC0009_NODE_1_length_10208_cov_2013.867437_pilon                       | 3-III                                     | 98.61                                  | 99.97                                   | gb:EU687233                              |
| 91-0118_S8_L001_NODE_1_length_10178_cov_13.815371_pilon                | 3-III                                     | 98.44                                  | 99.99                                   | gb:EU687233                              |
| cc0010_S8_L001_NODE_1_length_10206_cov_450.729095_pilon                | 3-III                                     | 98.86                                  | 100                                     | gb:EU687233                              |
| CC0051_k77_1_flag_1_multi_6641.2458_len_10057_pilon                    | 4-II                                      | 98.51                                  | 100                                     | gb:KP188557                              |
| CC0067_NODE_1_length_10197_cov_734.756522_pilon                        | 4-II                                      | 98.78                                  | 100                                     | gb:KP188557                              |
| cc0030a_S12_k77_1_flag_1_multi_2805.5225_len_10163_pilon               | 4-II                                      | 98.92                                  | 99.92                                   | gb:KP188557                              |
| cc0030b_S21_NODE_1_length_10173_cov_54.900771_pilon                    | 4-II                                      | 98.92                                  | 100                                     | gb:KP188557                              |
| CC0116_k77_2_flag_1_multi_2097.0000_len_10197_pilon                    | 4-II                                      | 98.67                                  | 100                                     | gb:KP188557                              |
| Spike_NODE_2_length_10203_cov_29.787675_pilon                          | 4-II                                      | 98.75                                  | 99.95                                   | gb:KP188557                              |
| CC0066_NODE_1_length_10174_cov_40.432750_pilon                         | 4-II                                      | 98.5                                   | 100                                     | gb:KP188557                              |
| 91-0106_S12_L001_k77_17_flag_1_multi_13.3022_len_10127_pilon           | 4-II                                      | 98.72                                  | 99.67                                   | gb:KP188557                              |

**Figure S1.** a) DEN-IM's quality control report containing information of the number of basepairs and the number of reads for the analysed samples, the estimated coverage depth before and after mapping, and the percentage of reads in the input data that were trimmed. b) DEN-IM's typing report for 24 CDSs recovered from the metagenomic dataset. The ID contains the CDS contig name, the typing result for serotype-genotype, the values for identity and coverage, and the GenBank ID of the closest reference in the Typing Database containing 161 complete DENV genomes.

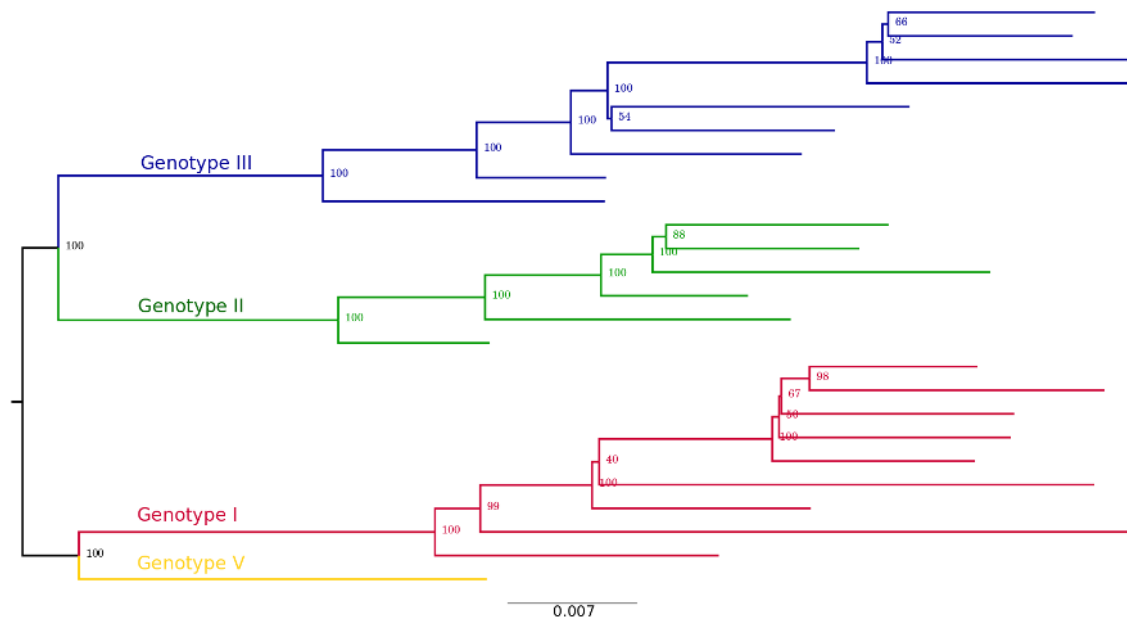




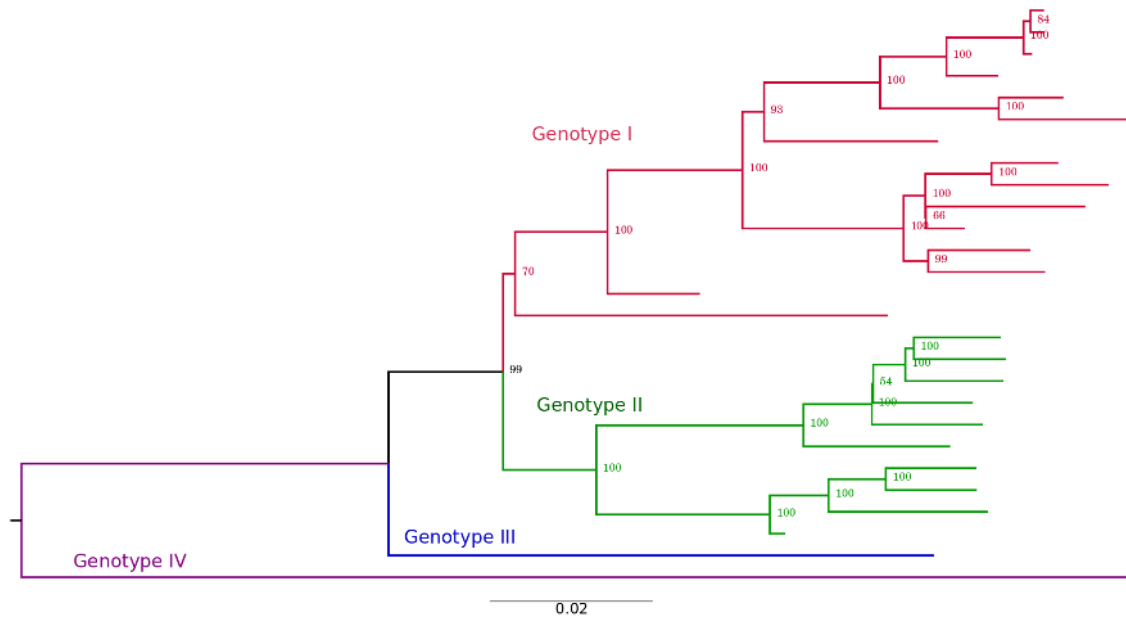
**Figure S2.** Maximum Likelihood inference of the multiple sequence alignment of the 46 DENV-1 complete genomes in the typing dataset. 1635 complete DENV-1 genomes were clustered at 98% nucleotide identity and the representative genomes were aligned with mafft. A maximum likelihood tree was inferred with RAxML. The tree is coloured according to genotype (red: genotype I; green: genotype II; blue: genotype III; purple: genotype IV). The sample JF459993, marked with a star, is currently annotated in ViPR as belonging to genotype IV but, given to the good phylogenetic support, it was re-classified as belonging to the genotype I.



**Figure S3.** Maximum Likelihood inference of the multiple sequence alignment of the 63 DENV-2 complete genomes in the typing dataset. 1067 complete DENV-1 genomes were clustered at 98% nucleotide identity and the representative genomes were aligned with mafft. A maximum likelihood tree was inferred with RAxML. The tree is coloured according to genotype (red: genotype I; green: genotype II; blue: genotype III; purple: genotype IV).



**Figure S4.** Maximum Likelihood inference of the multiple sequence alignment of the 25 DENV-3 complete genomes in the typing dataset. 807 complete DENV-3 genomes were clustered at 98% nucleotide identity and the representative genomes were aligned with mafft. A maximum likelihood tree was inferred with RAxML. The tree is coloured according to genotype (red: genotype I; green: genotype II; blue: genotype III; purple: genotype IV).



**Figure S5.** Maximum Likelihood inference of the multiple sequence alignment of the 27 DENV-4 complete genomes in the typing dataset. 320 complete DENV-4 genomes were clustered at 98% nucleotide identity and the representative genomes were aligned with mafft. A maximum likelihood tree was inferred with RAxML. The tree is coloured according to genotype (red: genotype I; green: genotype II; blue: genotype III; purple: genotype IV).