

1 **Text S1: General pipeline used to create semi-artificial and completely artificial**
2 **datasets**

3
4 **Semi-artificial datasets**

5 ***Analysis of the real datasets with Geneious***

6 In this study, 8 real datasets were used. Some of them were spiked with artificial data in order
7 to create semi-artificial datasets. These 8 real datasets are all RNA-seq datasets from
8 virus/viroid-infected plants obtained using Illumina technology (either HiSeq, MiSeq or
9 NextSeq sequencer). They have been chosen in order to cover as much as possible host plant
10 diversity (fruit trees, vegetables and biological indicator plants), pathogen diversity (RNA and
11 DNA viruses, viroids) and sequencing options (reads length ranging from 50 to 301 bp between
12 each dataset, number of reads per dataset from 65,177 to 49,052,832 reads, and single-end
13 or paired-end reads). All real datasets are therefore composed of plant host and virus/viroid
14 reads. A more detailed description of these real datasets is available in Table S1.

15 Before the creation of the artificial reads, each real dataset was analysed in order to check that
16 we were able to detect every viruses/viroids and to assess their presumed frequencies. Note
17 that for each dataset, the presence of the viruses/viroids has also been confirmed by PCR
18 and/or ELISA. The datasets were analysed in Geneious Prime 2019.2.1 (Biomatters, New
19 Zealand). The main steps performed were:

- 20 - Pairing of the reads for paired-end datasets. For each paired-end dataset, all the reads
21 have always been paired (i.e. no unpaired reads have been used).
- 22 - Merging/assembling of the paired reads to obtain longer reads for paired-end datasets
- 23 - Trimming of the merged and unmerged reads (BBduk, minimum quality = Q25,
24 minimum length = 20 bp).
- 25 - Removing duplicates of the trimmed, merged and unmerged reads. We obtain the final
26 “filtered reads”.
- 27 - De novo assembly of the filtered reads using SPAdes for the long reads (301 or 150
28 bp) or Velvet for the short reads (50 or 75 bp).
- 29 - Blasting of the contigs against the RefSeq virus database
30 (<https://ftp.ncbi.nlm.nih.gov/refseq/release/viral/>) using tblastx from the blast+ suite
31 ([https://blast.ncbi.nlm.nih.gov/Blast.cgi?CMD=Web&PAGE_TYPE=BlastDocs&DOC](https://blast.ncbi.nlm.nih.gov/Blast.cgi?CMD=Web&PAGE_TYPE=BlastDocs&DOC_TYPE=Download)
32 [TYPE=Download](https://blast.ncbi.nlm.nih.gov/Blast.cgi?CMD=Web&PAGE_TYPE=BlastDocs&DOC_TYPE=Download)).
- 33 - Mapping of the filtered reads against either the top blast hit reference sequences or the
34 contigs (if we have a complete genome, which depends on the dataset). If a read could

35 be mapped with equal match to two or more strains, it was randomly assigned to one
36 of them.

37 - Assessment of the proportion of each virus/viroid of the dataset.

38

39 In the end, 8 real datasets were analysed. Three were already interesting and showed
40 challenges that we wanted to test. Therefore, these 3 datasets (No. 7, 8 and 9) were been
41 spiked with artificial viral reads or modified. The other 5 real datasets were used to create 7
42 semi-artificial datasets.

43 ***Creation of the semi-artificial datasets with ART and Linux tools***

44 ART (Huang *et al.*, 2012) is a set of simulation tools allowing the generation of synthetic next-
45 generation sequencing reads. It mimics the real sequencing process and allows the creation
46 of FASTQ files with single or paired-end reads of different sizes.

47 The first challenge was to create artificial reads showing the same quality score as the reads
48 from the real datasets. ART allows the creation of our own quality profile files. Therefore, for
49 each dataset, quality profiles were created using the real reads. In the case of paired-end
50 datasets, quality profiles were generated both for the reads 1 and reads 2. This was done using
51 the “art_profiler_illumina” tool.

52 Then, artificial reads showing the same size and quality score as the one of the real datasets
53 were created using the “art_illumina” tool. The following options were used:

54 - qprof1: the quality profile file generated for the reads 1.

55 - qprof2: the quality profile file generated for the reads 2.

56 - p: for a paired-end read simulation (only for the paired-end datasets).

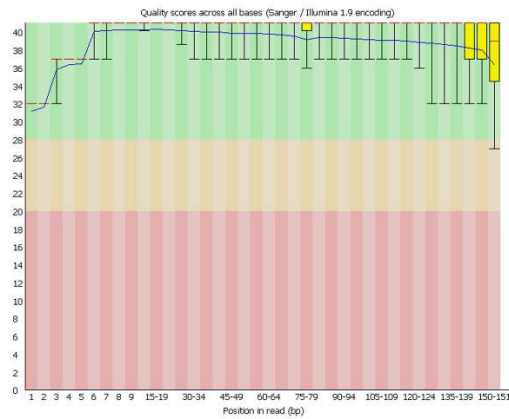
57 - l: the read length (bp), which is the same as the real dataset read length.

58 - f: the expected fold coverage. This number varies between the datasets and the
59 strains. In the GitLab page ([https://gitlab.com/ilvo/VIROMOCKchallenge/-](https://gitlab.com/ilvo/VIROMOCKchallenge/-/tree/master/Datasets)
60 [/tree/master/Datasets](https://gitlab.com/ilvo/VIROMOCKchallenge/-/tree/master/Datasets)), a table presenting the composition of each dataset is available.

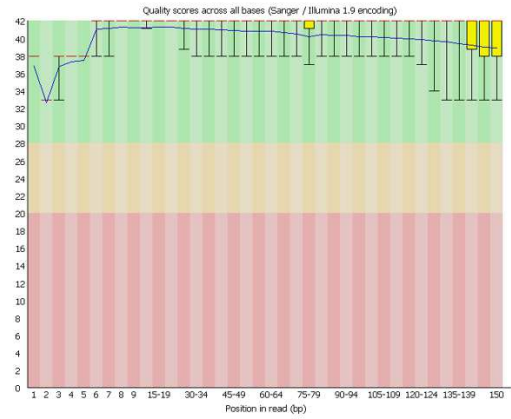
61 This number is provided in the “Expected average number of reads per position”
62 column.

63 Figure S1 shows the results of a FastQC analysis on both the real and artificial reads generated
64 for dataset 1. The quality scores are extremely similar between both type of reads, and the
65 same kind of results have been obtained for the other semi-artificial datasets.

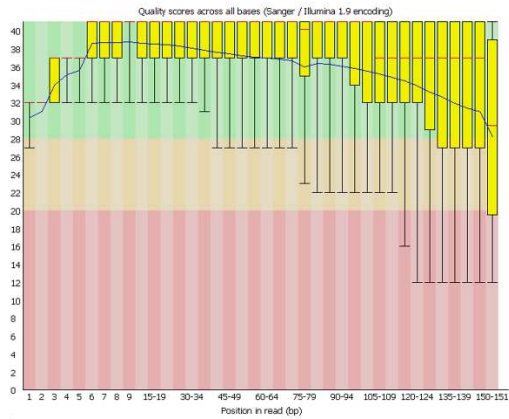
Real reads R1



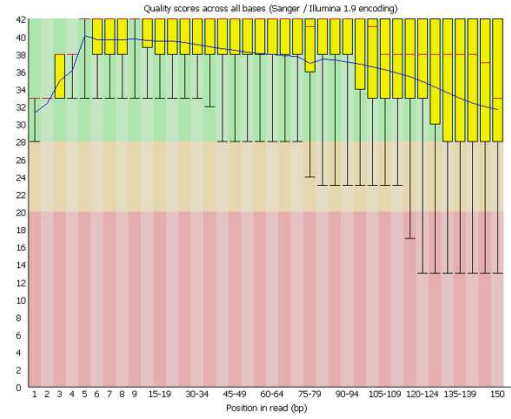
Artificial reads R1



Real reads R2



Artificial reads R2



66

67 **Figure S1:** Quality scores of real and artificial reads of the paired-end dataset 1. The artificial reads
68 correspond to the citrus tristeza virus (CTV) JQ911663 strain.

69

70 The second challenge was to merge artificial and real reads in the same file and with similar
71 headers. Indeed, ART creates its own header for each read and does not allow to change it.
72 The different steps to overcome this issue are summarized in Figure S2. These steps were
73 performed using command-line in Linux and are the following:

74

- Both artificial and real reads were merged in the same file.

75

- The old headers were removed and new headers were added. These new headers are all the same, they look like a “real” header, but they have a number in the end. If the dataset contains x reads, x different numbers have been generated and put randomly at the end of the headers. Here is an example of headers for a hypothetical paired-end dataset with 3 reads, the number added being in red:

76

77

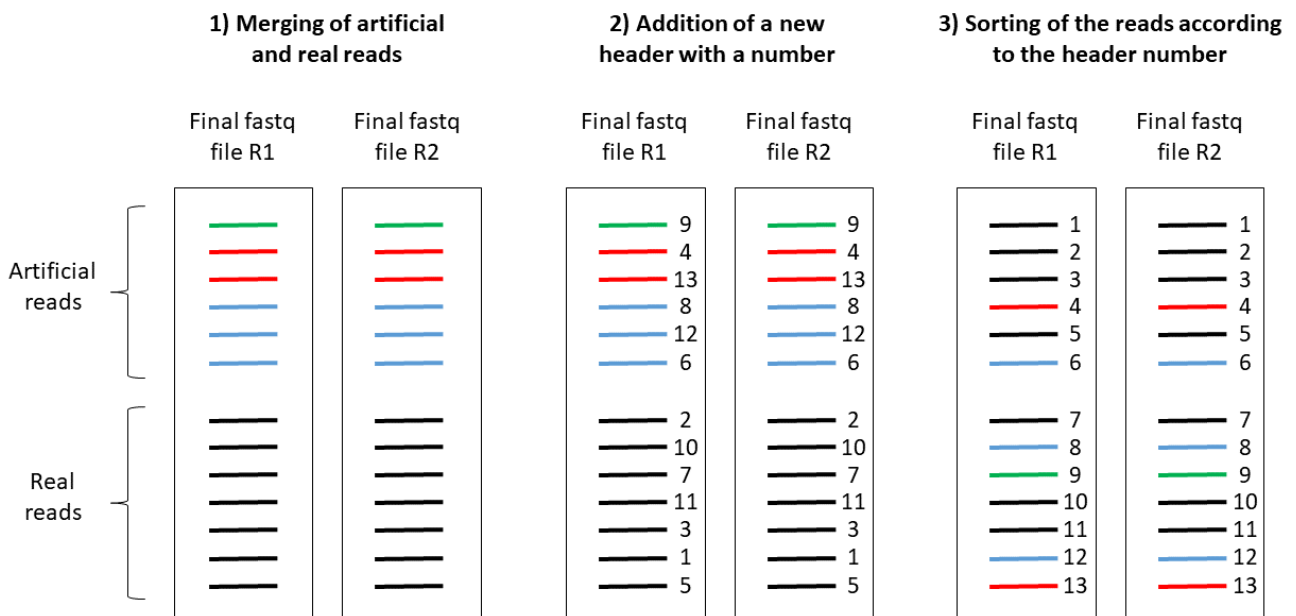
78

79

80

81
82
83 @E00526:48:HT7CTCXY:3:1101:16295:993:1 1:N:0:NGAGCTAC+NTAGCCTT Read 1
84 @E00526:48:HT7CTCXY:3:1101:16295:993:1 2:N:0:NGAGCTAC+NTAGCCTT Read 2
85
86 @E00526:48:HT7CTCXY:3:1101:16295:993:3 1:N:0:NGAGCTAC+NTAGCCTT Read 1
87 @E00526:48:HT7CTCXY:3:1101:16295:993:3 2:N:0:NGAGCTAC+NTAGCCTT Read 2
88
89 @E00526:48:HT7CTCXY:3:1101:16295:993:2 1:N:0:NGAGCTAC+NTAGCCTT Read 1
90 @E00526:48:HT7CTCXY:3:1101:16295:993:2 2:N:0:NGAGCTAC+NTAGCCTT Read 2
91

92 - The reads were sorted according to number of the header. It allowed for the mixing of
93 the real and artificial reads.



94 **Figure S2:** Final steps of the semi-artificial datasets creation.

95

96 **Completely artificial datasets**

97 Eight artificial datasets were generated (Datasets 11-18). Each artificial dataset consists of a
98 mix of several strains from the same viral species showing different frequencies. The virus
99 species were selected to be as divergent as possible. Therefore, the selected viruses have (i)
100 a DNA or RNA genome, (ii) a single or double-stranded genome, (iii) a linear, circular and/or
101 segmented genome, and (iv) show a genome length ranging from 2.8 to 17.1 kb. The virus
102 species and strains selected are presented in the GitLab page of each dataset
103 (<https://gitlab.com/ilvo/VIROMOCKchallenge/-/tree/master/Datasets>).

104

105 For each viral species, all the complete genome sequences publicly available in April 2019
106 were retrieved from NCBI GenBank. Then, a multiple alignment of all the sequences belonging
107 to the same species has been performed using the “Geneious Alignment” tool from Geneious
108 Prime (Biomatters, New Zealand) with default parameters. The obtained pairwise identity
109 matrices have been used to select between 3 and 6 isolates per species. The isolates were
110 selected to be as divergent as possible from each other. The average percentage of identity
111 among isolates ranges between 62.8 (for banana streak virus) and 89.7% (for eggplant mottled
112 dwarf virus).

113 For each isolate, artificial viral reads of 150 bp have been synthesized using the ART software
114 (Huang *et al.*, 2012) from NCBI reference genomes and no single nucleotide polymorphisms
115 (SNPs) have been added. The artificial reads were created using the “art_illumina” tool. The
116 following options have been used:

117 - p: for a paired-end read simulation.

118 - l: 150 bp.

119 - f: the expected fold coverage. This number varies between the datasets and the
120 isolates. In the GitLab page, a table presenting the composition of each dataset is
121 available. This number is provided in the “Expected average number of reads per
122 position” column.

123 - ss: HiSeq2000.

124

125 **References**

126 **Huang, W., Li, L., Myers, J.R. and Marth, G.T.** (2012) ART: a next-generation sequencing read
127 simulator. *Bioinformatics* **28**, 593–594.