1    Choice of assembly software has a critical impact on virome characterisation.

2

3    Thomas D.S. Sutton*[1,2,#], Adam G. Clooney*[1,2], Feargal J. Ryan*[1,2,3], R. Paul Ross[1,2,4],

4    ColinHill[1,2]

5    * Authors contributed equally

6    # Corresponding author, t.sutton@umail.ucc.ie

7    1. APC Microbiome Ireland, Cork, Ireland

8    2. School for Microbiology, University College Cork

9    3. (Current location) South Australian Health and Medical Research Institute, Adelaide,

10    Australia.

11    4. Teagasc Food Research Centre, Fermoy, Cork, Ireland

12

13    **<u>Abstract</u>**

14    *Background*

15    The viral component of microbial communities play a vital role in driving bacterial diversity,

16    facilitating nutrient turnover and shaping community composition. Despite their importance, the vast

17    majority of viral sequences are poorly annotated and share little or no homology to reference

18    databases. As a result, investigation of the viral metagenome (virome) relies heavily on *de novo*

19    assembly of short sequencing reads to recover compositional and functional information.

20    Metagenomic assembly is particularly challenging for virome data, often resulting in fragmented

21    assemblies and poor recovery of viral community members. Despite the essential role of assembly in

22    virome analysis and difficulties posed by these data, current assembly comparisons have been limited

23    to subsections of virome studies or bacterial datasets.

24 *Design*

25 This study presents the most comprehensive virome assembly comparison to date, featuring 16

26 metagenomic assembly approaches which have featured in human virome studies. Assemblers were

27 assessed using four independent virome datasets, namely; simulated reads, two mock communities,

28 viromes spiked with a known phage and human gut viromes.

29 *Results*

30 Assembly performance varied significantly across all test datasets, with SPAdes (meta) performing

31 consistently well. Performance of MIRA and VICUNA varied, highlighting the importance of using a

32 range of datasets when comparing assembly programs. It was also found that while some assemblers

33 addressed the challenges of virome data better than others, all assemblers had limitations. Low read

34 coverage and genomic repeats resulted in assemblies with poor genome recovery, high degrees of

35 fragmentation and low accuracy contigs across all assemblers. These limitations must be considered

36 when setting thresholds for downstream analysis and when drawing conclusions from virome data.

## **Keywords**

37

38 Virome, viral, assembly, metagenome, benchmark, comparison, bacteriophage, phage

39

## **Background**

41    The rapid evolution of metagenomics and high throughput sequencing technologies has revolutionised

42    the study of microbial communities, giving new insights into the role and identity of the uncultivated

43    microbes which account for the majority of metagenomic sequences (Solden, Lloyd et al. 2016).

44    However, the majority of microbial sequencing efforts have focused on the characterisation of

45    prokaryotic microbes. Viral metagenomes (viromes) are dominated by novel sequences, often with up

46    to 90% of sequences sharing little to no homology to reference databases (Aggarwala, Liang et al.

47    2017).  Bacteriophage, the most abundant members of viral communities, play a key role in the

48    shaping the composition of microbial communities and facilitate horizontal gene transfer (Paul 2008).

49    Viromes have been shown to play a role in global geochemical cycles (Breitbart 2011) and have been

50    studied in varied ecosystems including the ocean (Hurwitz and Sullivan 2013). Viromes of the human

51    body are of particular interest, where they have been linked to disease status (Norman, Handley et al.

52    2015), maintaining human health (Manrique, Bolduc et al. 2016) and shaping the gut microbiome in

53    early life (Lim, Zhou et al. 2015, McCann, Ryan et al. 2018).  Due to the predominance of

54    uncharacterised viral sequences "viral dark matter"; (Roux, Hallam et al. 2015), and the lack of a

55    universal marker gene, virome studies rely on database independent analysis methods and depend

56    heavily on *de novo* assembly to resolve viral genomes from metagenomic sequencing reads.

57            Metagenomic assemblers typically use de Bruijn graph (DBG) approaches to address the

58    complexity and size of metagenomic datasets in an accurate and efficient manner. Microbial

59    metagenomes pose significant challenges to DBG assembly when compared to single genome

60    assemblies often complicating the DBG and leading to fragmentation and/or misassembly (Olson,

61    Treangen et al. 2017). These challenges include uneven sequencing coverage of organisms within the

62    metagenome, the presence of conserved regions across different species, repeat regions within

63    genomes and the introduction of false $k$-mers by both closely related genomes at differing abundances

64    and sequencing errors at high read coverage. This hampers the use of coverage statistics to resolve

65    repeat regions between and within genomes (Olson, Treangen et al. 2017).

66    A wide array of metagenomic assembly programs have been employed, each addressing

67    aspects of metagenomic challenges to varying degrees. However, many of these programs have been

68    designed and optimised for bacterial metagenomes, which share many assembly challenges of

69    viromes but to a lesser degree. Virome data is characterised by high proportions of repeat regions

70    within viral genomes (Minot, Grunberg et al. 2012), hypervariable genomic regions associated with

71    host interaction (Warwick-Dugdale, Solonenko et al. 2018) and high mutation rates which lead to

72    increased metagenomic complexity and strain variation (Roux, Emerson et al. 2017). Low DNA

73    yields also limit read coverage and often require a multiple displacement amplification (MDA) step

74    which has been shown to preferentially amplify small single stranded DNA viruses (Kim and Bae

75    2011). Extremes in read coverage caused by MDA bias and dominant viral taxa such as crAssphage,

76    which can make up large proportions of human gut viromes (Dutilh, Cassman et al. 2014), sequester

77    sequencing resources and result in insufficient coverage of low abundance viruses. These challenges

78    result in fragmented virome assemblies (García-López, Vázquez-Castellanos et al. 2015), limiting

79    their use in downstream analysis. Despite benchmarks of bacterial metagenomes having highlighted

80    failings and benefits of particular assembly programs, many poorly performing assemblers have

81    featured in virome studies (Foulongne, Sauvage et al. 2012, Hannigan, Meisel et al. 2015, Guo, Hua et

82    al. 2017).

83    Accurate comparison of metagenomic assemblers is complicated by the unknown

84    composition of metagenomic datasets and the limited applicability of general assembly statistics such

85    as N50 (Deng, Naccache et al. 2015, Vollmers, Wiegand et al. 2017).  To address this, the accuracy

86    and efficacy of metagenomic assembly programs is often evaluated using simulated datasets and

87    mock communities of known composition. Although these simulated datasets are undergoing constant

88    improvements (Sczyrba, Hofmann et al. 2017, Fritz, Hofmann et al. 2018), they have focused

89    primarily on bacterial metagenomes and remain limited in their ability to accurately replicate the

90    challenges of true metagenomes. While some virome-specific assembly benchmarks have been

91    performed, many have been limited to a small number of assemblers, 454 data or subsections of

92    virome studies and have exclusively used simulated data (Aguirre de Cárcer, Angly et al. 2014, Smits,

93    Bodewes et al. 2014, Vázquez-Castellanos, García-López et al. 2014, García-López, Vázquez-

94    Castellanos et al. 2015, Hesse, van Heusden et al. 2017, Roux, Emerson et al. 2017).

95          Here we expand upon previous studies and present a detailed investigation of assembly

96    software for virome analysis which compares all those previously used in human virome studies to

97    date, as well as other popular or more recently published assemblers (Table 1). We compare assembly

98    efficacy and accuracy using simulated viromes, mock viral communities and human gut viromes

99    spiked with a known exogenous bacteriophage. Furthermore we confirm these findings using human

100   virome data from published datasets and assess computational parameters such as runtime and RAM

101   usage. We also investigate in detail the impact of sequencing coverage and genomic repeats on

102   assembly performance and highlight important considerations for future virome studies. Together

103   these data comprise most comprehensive virome assembly benchmark to date.

| | Link | Version used | Reference |
|---|---|---|---|
| ABySS (*k*-mer 63) | http://www.bcgsc.ca/downloads/abyss/ | v2.0.2 | (Simpson, Wong et al. 2009) |
| ABySS (*k*-mer 127) | http://www.bcgsc.ca/downloads/abyss/ | v2.0.2 | (Simpson, Wong et al. 2009) |
| CLC | https://www.qiagenbioinformatics.com/products/clc-assembly-cell/ | v5.0.5 | https://www.qiagenbioinformatics.com/ |
| Geneious | https://www.geneious.com/features/assembly-mapping/ | | (Kearse, Moir et al. 2012) |
| IDBA UD | https://i.cs.hku.hk/~alse/hkubrg/projects/idba_ud | v1.1.1 | (Peng, Leung et al. 2012) |
| MEGAHIT | https://github.com/voutcn/megahit | V1.1.1-2 | (Li, Luo et al. 2016) |
| MetaVelvet | https://metavelvet.dna.bio.keio.ac.jp/ | V1.2.02 | (Namiki, Hachiya et al. 2012) |
| MIRA | http://www.chevreux.org/mira_downloads.html | V4.0.2 | (García-López, Vázquez-Castellanos et al. 2015) |
| Ray Meta | http://denovoassembler.sourceforge.net/ | V2.3.0 | (Boisvert, Raymond et al. 2012) |
| SOAPdenovo2 | http://soap.genomics.org.cn/soapdenovo.html | v2.04 | (Luo, Liu et al. 2012) |
| SPAdes | http://cab.spbu.ru/software/spades/ | V3.10.0 | (Bankevich, Nurk et al. 2012) |
| SPAdes meta | http://cab.spbu.ru/software/spades/ (variation of SPAdes applied with flag) | V3.10.0 | (Nurk, Meleshko et al. 2017) |
| SPAdes sc | http://cab.spbu.ru/software/spades/ (variation of SPAdes applied with flag) | V3.10.0 | (Bankevich, Nurk et al. 2012) |
| SPAdes sc careful | http://cab.spbu.ru/software/spades/ (variation of SPAdes applied with flag) | V3.10.0 | (Bankevich, Nurk et al. 2012) |
| Velvet | https://www.ebi.ac.uk/~zerbino/velvet/ | V1.2.10 | (Zerbino and Birney 2008) |
| VICUNA | https://github.com/broadinstitute/mvicuna | V1.3 | (Vázquez-Castellanos, García-López et al. 2014) |

104   *Table 1: A list of assemblers used in this study*

## **Results**

*Simulated virome dataset*

The ability to accurately recover each of the 572 members of a published simulated community (Hesse, van Heusden et al. 2017) and the degree of fragmentation, was assessed by aligning the resulting contigs from each assembler to the reference genomes (Fig. 1). MetaVelvet was not included in this analysis as it failed to reach completion after seven days. Approximately half of the genomes in the community featured an average recovered genome fraction less than 75% and exhibited higher degrees of fragmentation (>10 contigs per genome on average) across all assemblers. For 87 of the 572 genomes there was an average recovered genome fraction of less than 20% across all assemblers (the low recovered genome fraction of VICUNA was excluded as an outlier). Of these genomes, 84 were present at low abundance (lowest 40% of all abundances normalised to genome length). The remaining three genomes were present at higher normalised abundances (50 – 80$^{th}$ percentile) but featured the some of the highest proportions of genomic repeats (70$^{th}$-90$^{th}$ percentile).

Normalised genome abundance within the community had a strong positive correlation with recovered genome fraction across all assemblers (Supplementary Table 1, Additional file 5) and was verified using a linear model (Supplementary Table 2, Additional file 5), with the exception of SOAPdenovo2, which was negative.  Normalised abundance also correlated negatively with the degree of fragmentation (number of contigs) across all assemblers except Velvet which was positively correlated and Geneious which was not significant (Supplementary Table1, Additional file 5). None of the genomes in the lower 30$^{th}$ percentile of normalised abundance featured an average recovered genome fraction greater than 75%, further exemplifying the impact of low sequencing coverage. However high abundance did not consistently improve genome recovery and of the 172 genomes in the top 30% of normalised abundance, 20 featured an average genome fraction below 50%. The distance of the log transformed (due to extremes in values) normalised abundances from the mean was negatively correlated with recovered genome fraction across all assemblers (correlation coefficient: -0.42, p-value $< 2.2e^{-16}$). Of 171 genomes in the 40$^{th}$ – 60$^{th}$ percentile of normalised abundance 29 featured an average genome fraction below 50%. This indicates factors other than abundance may

132     hamper genome recovery. MIRA and Geneious both recovered a greater fraction of low abundance

133     genomes with fewer contigs than other assemblers. However, MIRA assemblies of 13 of the most

134     abundant genomes in the community (highest 10%) exhibited the highest degree of fragmentation in

135     the study, generating between 401 and 2983 contigs per genome.

136        The proportion of inverted repeats, palindromic repeats, tandem repeats and a total proportion

137     of genomic repeats was calculated for each genome. The total percentage of repeat regions predicted

138     in each genome was positively correlated with the degree of fragmentation observed in each assembly

139     across all assemblers with the exception of Ray Meta (Supplementary Table 3, Additional file 5), and

140     negatively correlated with recovered genome fraction across all assemblers except ABySS ($k$-mer

141     63/127), Geneious, and SOAPdenovo2. When this relationship between repeat regions and the

142     recovered genome fraction was assessed using a linear model, correlations were significant for  CLC,

143     MIRA, Ray Meta, Velvet, and all parameters of SPAdes (Supplementary Table 2, Additional file 5).

144     Both the proportion of repeat regions in a genome and the relative abundance of that genome

145     contribute to the variation in recovered genome fraction, though each explain a separate aspect of this

146     variation. No interaction was found between these two metrics.

147        VICUNA, Ray Meta, SOAPdenovo2, Geneious, ABySS (both $k$-mer sizes) and Velvet

148     recovered under 50% of the total genome fraction (all genomes in the community). VICUNA

149     produced just four contigs in total with high levels of mismatches (174 per 100kb on average) which

150     could possibly linked to the format of the artificial reads as this was not observed in real sequencing

151     data.  The five assemblers which recovered the highest genome fraction overall were SPAdes

152     (default), MEGAHIT, SPAdes (single cell), SPAdes (single cell + careful) and CLC. All assemblers

153     achieving a minimum average genome fraction of 50% were subjected to a ranking system

154     (Supplementary Table 4, Additional file 5). To compare both recovery and fragmentation assemblers

155     were ordered from best to worst based on genome recovery and number of aligned contigs. The

156     average rank resulted in Spades (default) performing best, recovering 72.2% overall genome

157     sequences with 8230 contigs. The remaining top five assemblers of this combined rank were SPAdes

158     (meta) 68.2% with 7419 contigs, SPAdes (single cell) 68.9% with 9506 contigs, CLC 68.6% with

159    9152 contigs and MEGAHIT 69.6% with 10083 contigs. The number of assemblies which recovered

160    greater than 90% of the target genome in one single contig was compared (Fig 2). SPAdes (default)

161    performed best, recovering 210, SPAdes (meta), SPAdes (single cell + careful), CLC, and SPAdes

162    (single cell) each recovered 179, 168, 162 and 160 genomes respectively.

163        The accuracy of assemblies was assessed by calculating the average count of indels,

164    mismatches, and misassemblies per 100kb across all genomes.  These counts were normalised to the

165    number of genomes each assembler recovered with a minimum genome fraction of 50%. These were

166    ranked according to their performance in all three metrics (Supplementary Table 4, Additional file 5),

167    with assemblies from Velvet having the lowest overall counts followed by ABySS, IDBA UD,

168    MEGAHIT and Ray Meta. With the exception of Ray Meta and SOAPdenovo2, the number of

169    mismatches per 100kb was negatively correlated with both genome abundance and recovered genome

170    fraction across all assemblers (Supplementary Table 1, Additional file 5).

171        The rate of false positive (no alignment to reference genomes) and false negative (recovered

172    genome fraction of 0%) contigs assembled allowed for the determination of sensitivity. A number of

173    assemblers had a sensitivity greater than 97%, however each returned greater than 7,000 contigs,

174    inferring a high degree of fragmentation (Table 2). MIRA assembled (partial or complete) 559 of the

175    genomes with a false positive count of just four. However, this was achieved from more than 27,000

176    contigs. ABySS (both *k*-mer sizes), Geneious, Ray Meta and Velvet returned very few false positives

177    but failed to detect many of the genomes present. SPAdes (meta) performed best with 558 of the 572

178    genomes detected and only five false positives resulting from 7419 contigs.

| | False Positives | False Negative | True Positives | No. of contigs returned* | Sensitivity |
|---|---|---|---|---|---|
| **ABSS (*k*-mer 63)** | 0 | 111 | 461 | 7957 | 80.59 |
| **ABySS (*k*-mer 127)** | 1 | 123 | 449 | 7732 | 78.50 |
| **CLC** | 34 | 5 | 567 | 9152 | 99.13 |
| **Geneious** | 9 | 190 | 382 | 958 | 66.78 |
| **IDBA UD** | 25 | 9 | 563 | 8999 | 98.43 |
| **MEGAHIT** | 21 | 8 | 564 | 10083 | 98.60 |
| **MetaVelvet** | N/A | N/A | N/A | N/A | N/A |
| **MIRA** | 4 | 13 | 559 | 27600 | 97.73 |

| | | | | | |
|---|---|---|---|---|---|
| **Ray Meta** | 0 | 213 | 359 | 4224 | 62.76 |
| **SOAPdenovo2** | 536 | 116 | 456 | 11548 | 79.72 |
| **SPAdes** | 29 | 3 | 569 | 8230 | 99.48 |
| **SPAdes meta** | 5 | 14 | 558 | 7419 | 97.55 |
| **SPAdes sc** | 38 | 7 | 565 | 9506 | 98.78 |
| **SPAdes sc careful** | 40 | 6 | 566 | 9724 | 98.95 |
| **Velvet** | 1 | 65 | 507 | 6343 | 88.64 |
| **VICUNA** | 0 | 558 | 14 | 4 | 2.45 |
| | | | | *572 in community | |

179    *Table 2: The number of false positive, false negative contigs generated by each assembler for the*

180    *Simulated community, together with the sensitivity rates*

181

182    *Mock community dataset*

183    Two mock viral communities were used to investigate the impact of high and low abundance ssDNA

184    viruses on assembly performance. Mock A (Table 3a) contained 12 viral genomes, 10 of which were

185    at equal abundance (9.82% of the total community) and two ssDNA genomes (NC_001330 and

186    NC_001422) at low abundance (0.92%). Analysis of this community showed that although some

187    assemblers, namely CLC, Geneious, SPAdes (single cell) and VICUNA, detected all 12 genomes, this

188    was at the expense of a large number of false positives (1143, 53, 1513 and 4969 respectively). Velvet

189    and MetaVelvet generated no false positives, but failed to assemble three genomes, while ABySS (for

190    both *k*-mers) generated a large number of false positives and failed to assemble four and six genomes,

191    respectively. IDBA UD and Ray Meta outperformed the other assemblers with an equal number of

192    contigs to genomes (12), followed by MEGAHIT, SPAdes (default) and SPAdes (meta) with 13, 14

193    and 14. Mock B (Table 3b) also contained 12 genomes but with a higher abundance of ssDNA

194    genomes NC_001330 and NC_001422 (32.47%). VICUNA assemblies of Mock B improved upon

195    those from Mock A as no false positives were generated, while the false positive rate in the MIRA

196    assembler increased to 94 from none in Mock A. IDBA UD performed best followed by SPAdes

197    (default), Ray Meta, MEGAHT and SPAdes (meta) based on sensitivity and number of contigs, while

198    ABySS (both *k*-mer sizes) and SOAPdenovo2 had the lowest sensitivity. Despite being a relatively

199    simple community consisting of 12 members, not all assemblers were able to recover all members

200    (Supplementary Table 5-6, Additional file 5). A greater number of assemblers (six) failed to assemble

201    all members of Mock B than Mock A (four). ABySS(*k*-mer 63), ABySS(*k*-mer 127), Velvet and

202    MetaVelvet failed to assemble 6, 4, 3 and 3 genomes respectively, in Mock A and 6, 4 ,1 and 1

203    genomes, respectively in Mock B. In addition, MIRA and SOAPdenovo2 failed to assemble 1 and 2

204    genomes respectively in Mock B.

**A)**

| | False Positives | False Negative | True Positive | No. of contigs returned* | Sensitivity |
|---|---|---|---|---|---|
| **ABySS (*k*-mer 63)** | 52 | 4 | 8 | 61 | 66.67 |
| **ABySS (*k*-mer 127)** | 50 | 6 | 6 | 56 | 50.00 |
| **CLC** | 1143 | 0 | 12 | 1299 | 100.00 |
| **Geneious** | 53 | 0 | 12 | 65 | 100.00 |
| **IDBA UD** | 0 | 0 | 12 | 12 | 100.00 |
| **MEGAHIT** | 0 | 0 | 12 | 13 | 100.00 |
| **MetaVelvet** | 0 | 3 | 9 | 26 | 75.00 |
| **MIRA** | 0 | 0 | 12 | 89 | 100.00 |
| **Ray Meta** | 0 | 0 | 12 | 12 | 100.00 |
| **SOAPdenovo2** | 2 | 0 | 12 | 23 | 100.00 |
| **SPAdes** | 0 | 0 | 12 | 14 | 100.00 |
| **SPAdes meta** | 0 | 0 | 12 | 14 | 100.00 |
| **SPAdes sc** | 1513 | 0 | 12 | 1527 | 100.00 |
| **SPAdes sc careful** | 0 | 0 | 12 | 15 | 100.00 |
| **Velvet** | 0 | 3 | 9 | 26 | 75.00 |
| **VICUNA** | 4969 | 0 | 12 | 5385 | 100.00 |
| | | | | *12 in community | |

205

**B)**

| | False Positives | False Negative | True Positives | No. of contigs returned* | Sensitivity |
|---|---|---|---|---|---|
| **ABySS (*k*-mer 63)** | 60 | 4 | 8 | 69 | 66.67 |
| **ABySS (*k*-mer 127)** | 132 | 6 | 6 | 139 | 50.00 |
| **CLC** | 450 | 0 | 12 | 505 | 100.00 |
| **Geneious** | 14 | 0 | 12 | 30 | 100.00 |
| **IDBA UD** | 0 | 0 | 12 | 12 | 100.00 |
| **MEGAHIT** | 0 | 0 | 12 | 14 | 100.00 |
| **MetaVelvet** | 0 | 1 | 11 | 24 | 91.67 |
| **MIRA** | 94 | 1 | 11 | 157 | 91.67 |
| **Ray Meta** | 0 | 0 | 12 | 13 | 100.00 |
| **SOAPdenovo2** | 2 | 2 | 10 | 27 | 83.33 |

| | | | | | |
|---|---|---|---|---|---|
| **SPAdes** | 0 | 0 | 12 | 13 | 100.00 |
| **SPAdes meta** | 0 | 0 | 12 | 14 | 100.00 |
| **SPAdes sc** | 593 | 0 | 12 | 607 | 100.00 |
| **SPAdes sc careful** | 0 | 0 | 12 | 14 | 100.00 |
| **Velvet** | 0 | 1 | 11 | 24 | 91.67 |
| **VICUNA** | 0 | 0 | 12 | 15 | 100.00 |
| | | | | *12 in community | |

206 *Table 3: The number of false positive, false negative contigs generated by each assembler for a) Mock*

207 *community A and b) Mock community B) along with the sensitivity rates for each*

208       All but three VICUNA assemblies in Mock A exhibited a high level of fragmentation,

209 generating 34.7 ± 35 (mean ± standard deviation) contigs per genome. Fragmentation was also seen in

210 MIRA assemblies to a lesser degree with 7.4 ± 10 (mean ± standard deviation) contigs per genome on

211 average. There was a high rate of fragmentation in CLC with one community member generating 144

212 contigs for genome KF302035. Average recovered genome fraction of 85.4 ± 6.4 % was skewed by

213 ABySS (*k*-mer 63), ABySS (*k*-mer 127), Velvet, MetaVelvet, SOAPdenovo2, and VICUNA which

214 recovered on average 49.5%, 66.6%, 73.8%, 73.8%, 29.7% and 76.6%, respectively. All other

215 assemblers recovered over 99% of each genome in the community (Supplementary Figure 1,

216 Additional file 6).

217       Closer inspection of the two ssDNA genomes present at lower relative abundance highlighted

218 significant differences in the average number of indels across all assemblies of the NC_001330 and

219 NC_001422 genomes versus other members of the community (p-value = 0.037). These genomes

220 exhibited an average of 41.7 ± 18.5 and 9.4 ± 20.4 indels per 100kb, while all other genomes featured

221 an average of 7.8 ± 18.9 indels per 100kb. The low abundant ssDNA genomes NC_001330 and

222 NC_001422 also featured the highest average mismatches per 100kb at 148.7 ± 3 and 302.5 ± 10.7,

223 respectively (Supplementary Figure 1, Additional file 6).

224 The degree of fragmentation observed by VICUNA and MIRA in Mock B was lower than in Mock A

225 with a mean of 1.3 ± 0.89 and 5.3 ± 7.7 contigs per genome, respectively. CLC fragmented genome

226 KF302035 in Mock B (44 contigs), but to a lesser degree than Mock A (144 contigs). MEGAHIT,

227 which recovered at least 98% of all genomes in Mock A, also recovered over 98% of all genomes in

228    Mock B except for the ssDNA genome NC_001422, of which 56.5% was recovered in two contigs.

229    The majority of assemblies exhibited 147.9 ± 0 and 297 ± 1 mismatches per 100kb for NC_001330

230    and NC_001422 (high abundance ssDNA), respectively, identical values to those measured in Mock

231    A. Velvet and MetaVelvet were exceptions with 184.2 and 860.2 for genome NC_001422 and

232    NC_001330. A similar pattern of high values across a narrow range was also observed with the

233    number of indels, with 49.3 to 32.9 present in all assemblies NC_001330. Genome NC_001422

234    featured 18.57 indels across all SPAdes assemblies (all parameters) and 860.2 across both Velvet and

235    Metavelvet assemblies. All other assemblers which successfully recovered this genome did not feature

236    any indels (Supplementary Figure 1, Additional file 6).

237    *Q33*

238    Five assemblers failed to generate contigs which met alignment thresholds and were subsequently

239    excluded from further analysis - namely ABySS (*k*-mer 63), ABySS (*k*-mer 127), SOAPdenovo2,

240    Velvet and MetaVelvet. All remaining assemblers recovered over 90% of the spiked Q33 genome

241    with the exception of MIRA (8.5%). Six assemblers recovered over 99% of the Q33 genome in a

242    single contig - SPAdes (meta) 99.74%, MEGAHIT (99.6%), VICUNA (99.6%), Ray meta (99.6%),

243    CLC (99.5%) and Geneious (99.1) (Fig. 3). However, only MEGAHIT assembled the Q33 genome

244    with a contig equal in length to the genome itself. SPAdes (meta) and CLC generated assemblies

245    shorter than the reference genome by 86 and 141 bases. VICUNA (723), Geneious (1765), and Ray

246    Meta (9884) each generated assemblies longer than the reference genome. SPAdes (default) SPAdes

247    (single cell), IDBA UD and SPAdes (single cell + careful) each assembled Q33 in 2, 3, 4, 5 and 5

248    contigs, respectively. Ray Meta and VICUNA assemblies had the lowest number of mismatches and

249    indels per 100kb, however Ray Meta exhibited the highest rate of misassemblies (2 relocations, 1

250    inversion). All assemblers featured a minimum of one local misassembly with the exception of

251    SPAdes (meta) did not feature any. The six best assemblies of the Q33 genome and the genome itself

252    are syntenic (although occasionally on the reverse strand) and the start and end point were not

253    conserved (Fig .3).

254    *Read depth analysis (Time and RAM)*

255    Assemblers were compared for practicality by measuring the time to reach completion and maximum

256    RAM usage via four published healthy human gut viromes (Manrique, Bolduc et al. 2016) and various

257    sequencing depths . It must be noted that all assembly tasks were allocated five threads, however

258    specifying the number of threads did not change the number of threads used by certain programs.

259    MetaVelvet was not included in this analysis as it failed to reach completion after running for seven

260    days. CLC and Geneious were performed on a desktop computer and therefore excluded from time

261    and RAM analysis. Run time is dependent upon the number of reads and this is largely linear in scale

262    with more reads leading to an increased assembly time (Fig. 4a). MIRA and Vicuna (Fig. 4a insert)

263    were the slowest with MIRA taking over 15 times longer than the other software to assemble 3.5

264    million reads.  SOAPdenovo2 had the shortest completion time followed by IBDA UD and Velvet.

265    Most assemblers were consistent across samples (observed via error bars) with the exception of

266    MIRA and Ray Meta. MIRA, Vicuna and Velvet (Fig. 4b insert) had the highest max RAM usage

267    while the lowest was Ray Meta, IDBA UD and SPAdes (meta) (Fig. 4b). The majority of assemblers

268    observed a linear scale pattern similar to that of run time.

269    *Read depth analysis N50 and Longest contig length*

270    For both the N50 (Fig. 4c) and the longest contig length (Fig. 4d), there was a large amount of

271    variation between samples for the majority of assemblers. The longest contig length showed a large

272    increase at the final sequencing depth. Particular assemblers, namely SPAdes (default), SPAdes

273    (meta), MEGAHIT and ABySS (*k*-mer 127), produced longer contigs as the sequence depth was

274    increased.

## **Discussion**

276    Many bacterial metagenomic assembly comparisons have highlighted that the choice of assembler has

277    a significant impact on downstream analysis and the accuracy of the reconstructed metagenome

278    (Mavromatis, Ivanova et al. 2007, Lindgreen, Adair et al. 2016, Greenwald, Klitgord et al. 2017,

279    Vollmers, Wiegand et al. 2017). We have found this also to be true for viral metagenomes, where

280    accurate and complete assembly are of particular importance given the lack of viral representation in

281    reference databases. Virome studies depend heavily on the assembly step and possess many features

282    which are challenging to successful assembly. In this study we compared the performance of those

283    assemblers used to date in human viral metagenomics studies using datasets of known and unknown

284    composition and varying complexity. These included a Q33 spiked virome, mock virome

285    communities, a simulated virome and the "Healthy human gut phageome" (Manrique, Bolduc et al.

286    2016). Each dataset provided unique attributes allowing for comparison of assembly performance on a

287    number of levels. The combination of artificial and real viromes used in this study allows for the

288    comparison of various aspects of assembly performance across a range of datasets rather than

289    depending on simulated viromes alone, as is commonly carried out in assembly comparisons

290    (Mavromatis, Ivanova et al. 2007, Fritz, Hofmann et al. 2018) .

291        The Simulated dataset featured 572 viral genomes at various relative abundances as published

292    by Vázquez-Castellanos and colleagues (Vázquez-Castellanos, García-López et al. 2014). Fragmented

293    assemblies of individual genomes within microbial communities hamper downstream analysis and

294    limit the conclusions which can be drawn from metagenomic data such as taxonomic and functional

295    profiles (Florea, Souvorov et al. 2011). Consequently, the percentage genome recovery and degree of

296    fragmentation was assessed across each assembler, with SPAdes (default, meta and single cell) each

297    performing well. VICUNA performed very poorly, recovering only four contigs with high numbers of

298    mismatches and misassemblies, despite having performed well with other datasets and being designed

299    to address challenges of heterogeneous viral populations (Yang, Charlebois et al. 2012). This failure

300    may reflect the computational challenges relating to the format of the simulated reads, as benchmarks

301    carried out within the VICUNA study itself only include actual sequencing reads (Yang, Charlebois et

302    al. 2012). However, similar poor performance has been previously observed in virome assembly

303    comparison using VICUNA and 454 reads (Vázquez-Castellanos, García-López et al. 2014). For

304    those assemblers which could recover greater than 90% of the reference genome in a single contig,

305    SPAdes (default) outperformed SPAdes (meta). This may be explained by a lack of strain variants in

306    the dataset and the fact that SPAdes (meta) was optimised to combine strain variants of each species

307    to form consensus sequences.

308        A subset of genomes were poorly recovered (<20% genome fraction) by nearly all

309    assemblers. This observation indicates that there are challenging aspects of viral genomes and

310    metagenomes which cannot be overcome with current assembly strategies. The strong positive

311    correlations between the relative abundance and genome fraction suggest that a low abundance

312    threshold applies to virome assembly, below which assemblies will consist of small fractions of the

313    viral genome, and in most cases be highly fragmented.  This detrimental impact of low coverage has

314    been well established in previous assembly comparison studies (García-López, Vázquez-Castellanos

315    et al. 2015, Roux, Emerson et al. 2017, Fritz, Hofmann et al. 2018). Highly abundant genomes also

316    caused similar recovery and fragmentation issues across all assemblers, which is of particular

317    importance due to the prevalence of extremely high abundance genomes in viral data (crAssphage,

318    certain ssDNA viruses). As both abundance extremes are common in virome data, their impact must

319    be considered when designing virome studies (i.e. sequencing depth). As relative abundance alone did

320    not fully explain the variation in genome fraction recovered, the role of genomic repeats (a well-

321    established assembly challenge (Acuña-Amador, Primot et al. 2018) was also investigated. However,

322    genomic repeats could explain the variation in genome fraction recovered to a lesser degree than

323    relative abundance, suggesting other factors contribute to poor genome recovery.

324        Compositional differences between final assemblies and viromes themselves must be taken

325    into account when drawing conclusions about virome composition and setting parameters for

326    downstream analysis. Challenges such as genomic content and strain variation are not currently

327    addressed in human virome assembly strategies and impact the reconstruction of certain members of a

328    virome. Hybrid sequencing, which uses both long and short reads to resolve genomic regions

329    associated with poor assembly (Warwick-Dugdale, Solonenko et al. 2018) is a promising new

330    technology which could address current virome assembly challenges. Extraction methods which may

331    reduce the bias introduced by MDA steps include using Swift Biosciences 1S Plus kit (Roux,

332    Solonenko et al. 2016) and/or increasing overall sequencing depth to improve coverage of lowly

333    abundant viral genomes (in conjunction with an assembler which is less sensitive to high coverage

334    sequences).

335    Performance of some assemblers in this study was hampered by high coverage datasets

336    (namely overlap consensus assemblers). VICUNA assemblies exhibited the highest degree of

337    fragmentation of all assemblers with Mock A, despite having resolved both high abundance ssDNA

338    genomes of Mock B to a single contig. MIRA also exhibited a high degree of fragmentation with high

339    abundance genomes in both simulated and mock datasets. However, MIRA was least affected by low

340    abundance reads, recovering a greater genome fraction of low abundance genomes than other

341    assemblers with fewer contigs. Assembly challenges of high coverage sequences in viromes may

342    potentially be addressed by sub-setting reads similar to the assembly approach used by

343    SLICEMBLER (Mirebrahim, Close et al. 2015).

344    Multi-assembler approaches such as the use of Geneious to generate consensus sequences from

345    separate assemblers have been developed (Koren, Treangen et al. 2014, Schürch, Schipper et al. 2014,

346    Deng, Naccache et al. 2015) but have not been included in human virome studies using short reads.

347    MIRA assemblies of the Q33 genome and some low abundance genomes in the Simulated dataset

348    were improved using Geneious, resolving greater genome fractions with fewer contigs (despite

349    Geneious recovering a lower genome fraction of the Simulated dataset overall). It is possible that

350    using these approaches could address issues facing each assembler, i.e. combine the assemblies of

351    SPAdes (meta) which performs well across all 4 datasets but struggles to recover low abundant

352    genomes, with MIRA assemblies which are less affected by low abundance but has difficulty

353    resolving genomes of higher abundance. Comparison of multi-assembler approaches and

354    combinations of various assemblers was not within the scope of this study, but should not be ruled out

355    as a potential method of improving virome assembly in cases where composition could be assessed

356    and obvious assembly challenges were known to be present.

357    Across all analysis methods in this study, SPAdes (meta) performed consistently well and

358    would be our recommendation. It performed best in the Simulated data based on false positives, true

359    positives and false negatives, best assembled the Q33 genome (recovery, fragmentation,

360    misassemblies and genome size) and performed well with both mock communities in recovering all

361    members accurately in one or two contigs. SPAdes (meta) RAM usage was low and did not increase

362    to the same degree as other assemblers with increasing sequencing depth. This recommendation is in

363    agreement with previous comparisons (Vollmers, Wiegand et al. 2017) which also suggested using

364    SPAdes (meta) due to its ability to accurately assemble members of bacterial metagenomes. SPAdes

365    (meta) is less able to accurately reconstruct micro-diversity as it generates a consensus assembly of

366    "strain–contigs" in a metagenome, which means it is better equipped to address the high mutation

367    rates observed in virome data (Nurk, Meleshko et al. 2017). This recommendation is also concurrent

368    with a previous study (Roux, Emerson et al. 2017) which found IDBA UD, MEGAHIT and SPAdes

369    (meta) to perform equally well when assessed using 14 simulated viromes. However, we found that

370    SPAdes (meta) outperformed IDBA UD and MEGAHIT in the Q33 spiked dataset, RAM usage in

371    relation to increasing sequencing depth, and in its ability to recover members of the Simulated virome

372    in a single contig.  This recommendation contradicts two previous assembly comparisons which found

373    CLC (Hesse, van Heusden et al. 2017) and Velvet (White, Wang et al. 2017) to be best suited to

374    virome data. However, SPAdes (meta) was not included in either study. Though SPAdes (meta) was

375    out performed by MIRA in the assembly of low abundance genomes in the Simulated dataset, MIRA

376    has limited application to large datasets. MEGAHIT also performed well across all datasets

377    performing well in relation to recovery, fragmentation and accuracy, but encountered some recovery

378    issues in mock datasets and minor accuracy issues with the Q33 genome.

379        The higher levels of accuracy (low mismatch indel and misassembly counts) of assemblers

380    which performed poorly in other metrics namely (velvet and ABySS ($k$-mer 63), highlights the trade-

381    off between accuracy and contiguity observed in previous assembly studies (Gritsenko, Nijkamp et al.

382    2012, Lin and Liao 2013). However, both IDBA-UD and MEGAHIT performed well for accuracy,

383    genome recovery and fragmentation. These assemblers may be worth considering if strain level detail

384    is of particular importance. The mock A and B datasets were used to assess the impact of

385    amplification bias on assembly performance.  All ssDNA assemblies featured an equal minimum

386    number of mismatches across both Mock A and B. This may be caused by challenges in the genomes

387    themselves hampering accurate assembly, but is more likely to reflect strain variation between

388    genome sequence featured in the original publication and the genome of the phage used in the

389    community itself.

390         The Q33 spiked virome consisted of pooled reads from three healthy human faecal samples,

391    each of which having been spiked with $10^7$ PFU ml$^{-1}$ of lactococcal phage Q33 prior to virome

392    extraction. This allowed for assembly comparison of one abundant member of a challenging viral

393    community. Despite the high relative abundance of the Q33 genome, only 6 assemblers could recover

394    over 90% of the genome in a single contig, of these SPAdes (meta) and MEGAHIT reconstructed the

395    Q33 genome accurately without the introduction foreign or chimeric DNA. It was also noted that the

396    genome synteny was conserved across these six assemblies. This may reflect circularization of the

397    linear Q33 genome during DNA extraction as the presence of cos sites has been previously predicted

398    (Mahony, Martel et al. 2013).

399         The longest contigs of each assembler were only detected at the highest sequencing depths

400    and varied across assemblers, which may indicate that high coverage is necessary to recover the

401    largest viral genomes in a community. However, it is also possible that these long contigs may reflect

402    misassemblies and duplication events caused by read errors at high sequencing depths which must be

403    considered when analysing high coverage data. At almost all sequencing depths Geneious, Vicuna,

404    Ray Meta and ABySS (k-mer 127) exhibited the highest N50 values, despite performing poorly in

405    other metrics. This further highlights the limitation of using N50 alone as a metric of metagenomic

406    assembly (Vollmers, Wiegand et al. 2017).

407         A further important consideration when performing any metagenomic assembly is

408    practicality; size of dataset, computational resources, bioinformatic resources, and how much hands-

409    on time is required from the end user. Both CLC and Geneious are available as a GUI (albeit requiring

410    a licence fee) which widens their audience to researchers with limited command-line experience (CLC

411    can also be run using the windows command line). However, this limits their practicality for large

412    scale virome studies as they are limited to the computational power of desktop computers and are not

413    suited to the assembly of large numbers of samples. Despite the limitations of computational power,

414    CLC performed well in all datasets in terms of genome recovery and fragmentation. Of the freely

415    available open source assemblers, MIRA and VICUNA are the least efficient in terms of RAM usage

416    and assembly time, reflecting limitations of the overlap consensus approach to assembly. This limits

417    their applicability to large virome datasets, and further increases the time required to carry out the

418    Geneious assembly approach which requires the output of both assemblers. Despite the long runtime,

419    VICUNA did not adhere to the number of cores specified, instead using all available cores.  All other

420    assemblers had a similar time requirements (with the exception of SOAPdenovo2 which performed

421    poorly across all datasets). Of the assemblers which consistently performed well in terms of accuracy,

422    genome fraction recovered and fragmentation, SPAdes (meta) was most efficient in terms of RAM

423    usage, which did not increase to the same degree as other assemblers with increasing sequencing

424    depth. MIRA stood out in terms of impracticality by generating by far the largest intermediate files of

425    any assembler, requiring several gigabytes of storage space for intermediate files.

426    The combination of results from four datasets facilitates accurate comparison of assemblers as

427    the limitations of each individual dataset vary.  Application of Phi29 MDA to amplify virome DNA to

428    sufficient quantities for sequencing can introduce significant bias and skew the original composition

429    of the virome, making quantitative viromics difficult (Kim and Bae 2011, Roux, Solonenko et al.

430    2016). As a result, it is likely that true diversity of viral metagenomes is not being accurately captured

431    using current virome extraction methods. However, as these procedures move away from steps known

432    to introduce bias, greater diversity will be observed. In this sense, the level of complexity of the Q33

433    dataset, which pooled three independent human viromes, provides a useful testbed for metagenomic

434    assemblers in future virome studies as extraction methods improve. Additionally, Q33 was not present

435    in the viromes prior to spiking, assemblers were not challenged by the presence of native strain

436    variations of Q33 genome.  In this study, assemblers were compared without individual optimisation

437    to the specific dataset. Feasibility dictates that, this "straight out of the box" approach to assembly is

438    used by almost all metagenomic assembly comparisons. Additionally, as the true composition of

439    metagenomes is unknown, any impact of parameter optimisation must be estimated from general

440    assembly statistics such as N50 and longest contig which have been highlighted to be of limited

441    usefulness (Aguirre de Cárcer, Angly et al. 2014, Vollmers, Wiegand et al. 2017).

## **Conclusions**

442

443    Of all assembly programs used in human virome studies, SPAdes (meta) addressed the challenges of

444    virome data most effectively. However, all assemblers have limitations and are hampered by aspects

445    of virome data. Low read coverage and high genomic repeats lead to assemblies with low recovered

446    genome fraction and a higher degree of fragmentation, with the assemblies themselves being less

447    accurate. This pattern was seen across all assemblers used in this study.

448            As assembler choice has significant implications for virome composition and the conclusions

449    which can be drawn from a dataset, assemblers which performed poorly in this study (i.e. low genome

450    recovery or accuracy and high degree of fragmentation) highlight a potential untapped resource in the

451    sequence data of previously conducted virome studies. It is highly likely that many viral sequences

452    were poorly assembled and reanalysis using a more effective assembler may yield new insights.

453    Design of future virome studies should carefully consider the impact of sequencing depth, as extremes

454    in read coverage will prevent the assembly and detection of viral genomes at both high and low

455    abundance.

456

## **Methods**

457

458    Each assembler with the exception of Geneious and CLC was run as per manual with default

459    parameters (unless stated) using a Lenovo x3650 M5 server with an intel Xeon processor E5-2690 v3

460    and 512Gb RAM . Geneious assembly approach mirrored that used in (Manrique, Bolduc et al. 2016)

461    by generating consensus sequences from the assemblies of both MIRA and Vicuna. CLC and

462    Geneious were run on a 64-bit windows 10 computer with an i5-4690 CPU and 16 GB of RAM.

463    *Data sources*

464    Sequencing reads from mock communities A and B featured in (Roux, Solonenko et al. 2016),

465    Simulated Virome dataset featured in (Hesse, van Heusden et al. 2017),  reads used to compare the

466    impact of sequencing depth on time and RAM usage featured in (Manrique, Bolduc et al. 2016) and

467      human viromes spiked with $10^7$ PFU of Lactococcal phage Q33 (Mahony, Martel et al. 2013) and

468      originated from (Shkoporov, Ryan et al. 2018) .

469      *Read Pre-processing*

470      Raw read quality was assessed with FastQC v0.11.5 and sequencing adapters were removed with

471      cutadapt v1.9.1 (Martin 2011) for the mock, Spiked and healthy gut virome data sets. Trimming and

472      filtering was carried out with Trimmomatic v0.36 (Bolger, Lohse et al. 2014)  using parameters

473      specific to each dataset. A sliding window size of 4 with a minimum Phred score of 30 and a

474      minimum length of 60bp was used with reads from both mock communities. The leading 15bp and

475      trailing 60bp were removed from "Healthy human gut phageome" reads and a sliding window of 4bp

476      with a minimum phred score of 20 was applied. The leading 10bp and trailing 100bp were removed

477      from the Q33 spiked virome reads and a sliding window size of 4bp with a minimum Phred score of

478      30. Filtered reads were through a minimum length filter of 60bp.

479      *Analysis methods*

480      Quality filtered reads from the Q33 spiked dataset consisted of 3 individual viromes which were

481      pooled and subsequently assembled. Contigs were aligned to the published Q33 using Blastn with an

482      e-value cut-off of $1e^{-20}$. Top hit alignments to the Q33 genome with a minimum alignment length of

483      800 bases and which shared 95% identity were included in further analysis using QUAST (v. 4.4)

484      (Gurevich, Saveliev et al. 2013) with "--unique mapping" flag. Further comparison and visualisation

485      of Q33 assemblies was carried out using Mauve (v. 20150226, build 10) (Darling, Mau et al. 2010).

486          Alignment and comparison of assemblies from mock and simulated data sets to reference

487      genomes was carried using MetaQUAST (v. 4.4) (Mikheenko, Saveliev et al. 2015) with "--unique

488      mapping" flag. Correlations were carried out using Spearman method and plots were generated using

489      the package ggplot2 (v 3.0.0) package in R (v.3.4.3).  These correlations were validated using a linear

490      model in R base library. For data which was not normally distributed, log transformation was carried

491      out.

492    Reads from the "healthy human gut phageome" were analysed to compare the overall

493    assembler efficiency and the impact of sequencing depth. Reads were randomly subset in pairs (both

494    the forward and reverse read of a pair were retained) to different depths using an in-house python

495    script. Samples were subset in increments of 300,000 reads to their respective maximum depth (2.7,

496    3.5, 3 and 3.3 million reads). The shell script *time* script, location */usr/bin/time*, was utilised to

497    measure the maximum RAM and length of time for each assembly to reach completion. All

498    assemblers were run using 5 threads where possible with the exception of CLC, Geneious, Ray Meta,

499    Velvet and Vicuna. Ray Meta and Velvet were run with 10, 1 thread(s) respectively. Ray Meta failed

500    to run with 5 while Velvet ran with 1 core despite 5 being allocated. Vicuna was also allocated 5

501    threads however used upwards of 20. MetaVelvet was run, but after 7 days had failed to reach

502    completion and was therefore removed from the subsequent analysis of these metrics. Contig statistics

503    and filtering (contigs greater than 1kb retained) were performed using the assembly-stats script from

504    the Pathogen Informatics group at the Wellcome Sanger Institute (https://github.com/sanger-

505    pathogens/assembly-stats).

506    **Figure Legends**

507    Figure 1: Relationship between percentage of each genome recovered (genome fraction), the number

508    of contigs required for each combination of genome and assembler and the abundance and proportion

509    of repeats for each genome. (A and B) Genomes are ordered by their average genome fraction across

510    all assemblers from high to low along the x-axis. (A main) Relative abundance, normalized by

511    genome length is plotted along y-axis with upper limit of 0.75% and colour of bars determined by

512    proportion of repeat regions in each genome. Blue bars represent genomes with a high proportion of

513    genomic repeats (4[th] quartile of all genomes), red represents all other genomes below this quartile. (A

514    insert) Expanded view of (A) without an upper limit of y value. (B) Percentage genome recovered is

515    plotted along the y axis. Points are coloured by assembler with shape of the point is denoting number

516    of contigs generated by each assembler for each genome.

517    Figure 2: Number of contigs each assembler recovered to a minimum genome fraction of 90% in a

518    single contig.

519    Figure 3: Mauve output of the Q33 reference genome (top) along with of the six assemblers which

520    recovered >99% of the genome with a single contig. Assembly regions outside of locally collinear

521    blocks which do not share homology to the reference genome are highlighted by a black outline.

522    Reverse complement of assemblies in the opposite orientation to the reference were plotted for

523    visualisation purposes (VICUNA, CLC, Geneious)

524    Figure 4: (A) Time, measured in seconds, for each assembly to reach completion successfully for each

525    read subset,  (B) the maximum RAM, measured in MB, used for each assembly for each read subset,

526    (C) mean N50 length and (D) mean contig length for 4 samples for each assembly across the read

527    subsets after filtering contigs less than 1000 bases. Points represent the mean time for the 4 samples

528    while error bars are the standard error.

### Abbreviations/Glossary

530    The following terms; Genome fraction, N50, number of contigs, misassemblies, local misassemblies,

531    are defined by QUAST (Mikheenko, Saveliev et al. 2015)

532    Genome fraction "is the total number of aligned bases in the reference, divided by the genome size. A

533    base in the reference genome is counted as aligned if there is at least one contig with at least one

534    alignment to this base. Contigs from repeat regions may map to multiple places, and thus may be

535    counted multiple times in this quantity."

536    N50 "is the contig length such that using longer or equal length contigs produces half (50%) of the

537    bases of the assembly. Usually there is no value that produces exactly 50%, so the technical definition

538    is the maximum length x such that using contigs of length at least x accounts for at least 50% of the

539    total assembly length."

540    Number of contigs "is the total number of contigs in the assembly that have size greater than or equal

541    to 0 bp."

542    Misassemblies "is the number of positions in the assembled contigs where the left flanking sequence

543    aligns over 1 kbp away from the right flanking sequence on the reference (relocation) or they overlap

544 on more than 1 kbp (relocation) or flanking sequences align on different strands (inversion) or

545 different chromosomes (translocation)."

546 Local misassemblies "A local misassembly has two or more distinct alignments covering the

547 breakpoint, the gap between left and right flanking sequences is less than 1 kbp and the left and right

548 flanking sequences both are on the same strand of the same chromosome of the reference genome."

549 **Data sources**

550 Sequencing reads from mock communities A and B:

551 http://datacommons.cyverse.org/browse/iplant/home/shared/iVirus/DNA_Viromes_library_compariso

552 n .

553 Simulated virome reads:

554 https://figshare.com/articles/Simulated_virome_datasest_for_assembly_and_annotation_tests/515116

555 3 .

556 Reads used to compare the impact of sequencing depth on time and RAM usage from the NCBI SRA;

557 http://www.ncbi.nlm.nih.gov/sra under the accession numbers SAMN04415496 to SAMN04415499

558 Human viromes spiked with $10^7$ PFU of Lactococcal phage Q33 phage

559 http://www.ncbi.nlm.nih.gov/sra under the accession numbers SRX3240741, SRX3240716,

560 SRX3240715

561

562 **Declarations**

563 **Ethics approval and consent to participate**

564 Not applicable

565 **Consent for publication**

566 Not applicable

567 **Competing interests**

568 The authors declared that they have no competing interests.

574 **Authors' contributions**

575 TDSS, AGC, FJR, PR, and CH conceived and designed experiments. TDSS, AGC, FJR

576 carried out bioinformatics analysis and drafted the manuscript. All authors approve and contributed to

577 the final manuscript.

580

581 **References**

582

583

584

585

586 Acuña-Amador, L., A. Primot, E. Cadieu, A. Roulet and F. Barloy-Hubler (2018). "Genomic repeats,
587 misassembly and reannotation: a case study with long-read resequencing of Porphyromonas gingivalis
588 reference strains." BMC genomics **19**(1): 54.
589 Aggarwala, V., G. Liang and F. D. Bushman (2017). "Viral communities of the human gut:
590 metagenomic analysis of composition and dynamics." Mobile DNA **8**(1): 12.
591 Aguirre de Cárcer, D., F. E. Angly and A. Alcamí (2014). "Evaluation of viral genome assembly and
592 diversity estimation in deep metagenomes." BMC Genomics **15**(1): 989.
593 Bankevich, A., S. Nurk, D. Antipov, A. A. Gurevich, M. Dvorkin, A. S. Kulikov, V. M. Lesin, S. I.
594 Nikolenko, S. Pham and A. D. Prjibelski (2012). "SPAdes: a new genome assembly algorithm and its
595 applications to single-cell sequencing." Journal of computational biology **19**(5): 455-477.

596 Boisvert, S., F. Raymond, É. Godzaridis, F. Laviolette and J. Corbeil (2012). "Ray Meta: scalable de
597 novo metagenome assembly and profiling." Genome biology **13**(12): R122.
598 Bolger, A. M., M. Lohse and B. Usadel (2014). "Trimmomatic: a flexible trimmer for Illumina
599 sequence data." Bioinformatics **30**(15): 2114-2120.
600 Breitbart, M. (2011). "Marine viruses: truth or dare."
601 Darling, A. E., B. Mau and N. T. Perna (2010). "progressiveMauve: multiple genome alignment with
602 gene gain, loss and rearrangement." PloS one **5**(6): e11147.
603 Deng, X., S. N. Naccache, T. Ng, S. Federman, L. Li, C. Y. Chiu and E. L. Delwart (2015). "An
604 ensemble strategy that significantly improves de novo assembly of microbial genomes from
605 metagenomic next-generation sequencing data." Nucleic acids research **43**(7): e46-e46.
606 Dutilh, B. E., N. Cassman, K. McNair, S. E. Sanchez, G. G. Silva, L. Boling, J. J. Barr, D. R. Speth,
607 V. Seguritan and R. K. Aziz (2014). "A highly abundant bacteriophage discovered in the unknown
608 sequences of human faecal metagenomes." Nature communications **5**: ncomms5498.
609 Florea, L., A. Souvorov, T. S. Kalbfleisch and S. L. Salzberg (2011). "Genome assembly has a major
610 impact on gene content: a comparison of annotation in two Bos taurus assemblies." PLoS One **6**(6):
611 e21400.
612 Foulongne, V., V. Sauvage, C. Hebert, O. Dereure, J. Cheval, M. A. Gouilh, K. Pariente, M. Segondy,
613 A. Burguière and J.-C. Manuguerra (2012). "Human skin microbiota: high diversity of DNA viruses
614 identified on the human skin by high throughput sequencing." PloS one **7**(6): e38499.
615 Fritz, A., P. Hofmann, S. Majda, E. Dahms, J. Droege, J. Fiedler, T. R. Lesker, P. Belmann, M. Z.
616 DeMaere and A. E. Darling (2018). "CAMISIM: Simulating metagenomes and microbial
617 communities." bioRxiv: 300970.
618 García-López, R., J. F. Vázquez-Castellanos and A. Moya (2015). "Fragmentation and Coverage
619 Variation in Viral Metagenome Assemblies, and Their Effect in Diversity Calculations." Frontiers in
620 Bioengineering and Biotechnology **3**(141).
621 Greenwald, W. W., N. Klitgord, V. Seguritan, S. Yooseph, J. C. Venter, C. Garner, K. E. Nelson and
622 W. Li (2017). "Utilization of defined microbial communities enables effective evaluation of meta-
623 genomic assemblies." BMC genomics **18**(1): 296.
624 Gritsenko, A. A., J. F. Nijkamp, M. J. Reinders and D. d. Ridder (2012). "GRASS: a generic
625 algorithm for scaffolding next-generation sequencing assemblies." Bioinformatics **28**(11): 1429-1437.
626 Guo, L., X. Hua, W. Zhang, S. Yang, Q. Shen, H. Hu, J. Li, Z. Liu, X. Wang and H. Wang (2017).
627 "Viral metagenomics analysis of feces from coronary heart disease patients reveals the genetic
628 diversity of the Microviridae." Virologica Sinica **32**(2): 130-138.
629 Gurevich, A., V. Saveliev, N. Vyahhi and G. Tesler (2013). "QUAST: quality assessment tool for
630 genome assemblies." Bioinformatics **29**(8): 1072-1075.
631 Hannigan, G. D., J. S. Meisel, A. S. Tyldsley, Q. Zheng, B. P. Hodkinson, A. J. SanMiguel, S. Minot,
632 F. D. Bushman and E. A. Grice (2015). "The human skin double-stranded DNA virome:
633 topographical and temporal diversity, genetic enrichment, and dynamic associations with the host
634 microbiome." MBio **6**(5): e01578-01515.
635 Hesse, U., P. van Heusden, B. M. Kirby, I. Olonade, L. J. van Zyl and M. Trindade (2017). "Virome
636 Assembly and Annotation: A Surprise in the Namib Desert." Frontiers in Microbiology **8**(13).
637 Hurwitz, B. L. and M. B. Sullivan (2013). "The Pacific Ocean Virome (POV): a marine viral
638 metagenomic dataset and associated protein clusters for quantitative viral ecology." PloS one **8**(2):
639 e57355.
640 Kearse, M., R. Moir, A. Wilson, S. Stones-Havas, M. Cheung, S. Sturrock, S. Buxton, A. Cooper, S.
641 Markowitz and C. Duran (2012). "Geneious Basic: an integrated and extendable desktop software
642 platform for the organization and analysis of sequence data." Bioinformatics **28**(12): 1647-1649.
643 Kim, K.-H. and J.-W. Bae (2011). "Amplification methods bias metagenomic libraries of uncultured
644 single-stranded and double-stranded DNA viruses." Applied and environmental microbiology: AEM.
645 00289-00211.
646 Koren, S., T. J. Treangen, C. M. Hill, M. Pop and A. M. Phillippy (2014). "Automated ensemble
647 assembly and validation of microbial genomes." BMC bioinformatics **15**(1): 126.
648 Li, D., R. Luo, C.-M. Liu, C.-M. Leung, H.-F. Ting, K. Sadakane, H. Yamashita and T.-W. Lam
649 (2016). "MEGAHIT v1. 0: A fast and scalable metagenome assembler driven by advanced
650 methodologies and community practices." Methods **102**: 3-11.

Lim, E. S., Y. Zhou, G. Zhao, I. K. Bauer, L. Droit, I. M. Ndao, B. B. Warner, P. I. Tarr, D. Wang and L. R. Holtz (2015). "Early life dynamics of the human gut virome and bacterial microbiome in infants." Nature medicine **21**(10): 1228.

Lin, S.-H. and Y.-C. Liao (2013). "CISA: contig integrator for sequence assembly of bacterial genomes." PloS one **8**(3): e60843.

Lindgreen, S., K. L. Adair and P. P. Gardner (2016). "An evaluation of the accuracy and speed of metagenome analysis tools." Scientific reports **6**: 19233.

Luo, R., B. Liu, Y. Xie, Z. Li, W. Huang, J. Yuan, G. He, Y. Chen, Q. Pan and Y. Liu (2012). "SOAPdenovo2: an empirically improved memory-efficient short-read de novo assembler." Gigascience **1**(1): 18.

Mahony, J., B. Martel, D. M. Tremblay, H. Neve, K. J. Heller, S. Moineau and D. van Sinderen (2013). "Identification of a new P335 subgroup through molecular analysis of lactococcal phages Q33 and BM13." Applied and environmental microbiology **79**(14): 4401-4409.

Manrique, P., B. Bolduc, S. T. Walk, J. van der Oost, W. M. de Vos and M. J. Young (2016). "Healthy human gut phageome." Proceedings of the National Academy of Sciences **113**(37): 10400-10405.

Martin, M. (2011). "Cutadapt removes adapter sequences from high-throughput sequencing reads." EMBnet. journal **17**(1): pp. 10-12.

Mavromatis, K., N. Ivanova, K. Barry, H. Shapiro, E. Goltsman, A. C. McHardy, I. Rigoutsos, A. Salamov, F. Korzeniewski and M. Land (2007). "Use of simulated data sets to evaluate the fidelity of metagenomic processing methods." Nature methods **4**(6): 495.

McCann, A., F. J. Ryan, S. R. Stockdale, M. Dalmasso, T. Blake, C. A. Ryan, C. Stanton, S. Mills, P. R. Ross and C. Hill (2018). "Viromes of one year old infants reveal the impact of birth mode on microbiome diversity." PeerJ **6**: e4694.

Mikheenko, A., V. Saveliev and A. Gurevich (2015). "MetaQUAST: evaluation of metagenome assemblies." Bioinformatics **32**(7): 1088-1090.

Minot, S., S. Grunberg, G. D. Wu, J. D. Lewis and F. D. Bushman (2012). "Hypervariable loci in the human gut virome." Proceedings of the National Academy of Sciences **109**(10): 3962-3966.

Mirebrahim, H., T. J. Close and S. Lonardi (2015). "De novo meta-assembly of ultra-deep sequencing data." Bioinformatics **31**(12): i9-i16.

Namiki, T., T. Hachiya, H. Tanaka and Y. Sakakibara (2012). "MetaVelvet: an extension of Velvet assembler to de novo metagenome assembly from short sequence reads." Nucleic acids research **40**(20): e155-e155.

Norman, J. M., S. A. Handley, M. T. Baldridge, L. Droit, C. Y. Liu, B. C. Keller, A. Kambal, C. L. Monaco, G. Zhao and P. Fleshner (2015). "Disease-specific alterations in the enteric virome in inflammatory bowel disease." Cell **160**(3): 447-460.

Nurk, S., D. Meleshko, A. Korobeynikov and P. A. Pevzner (2017). "metaSPAdes: a new versatile metagenomic assembler." Genome research: gr. 213959.213116.

Olson, N. D., T. J. Treangen, C. M. Hill, V. Cepeda-Espinoza, J. Ghurye, S. Koren and M. Pop (2017). "Metagenomic assembly through the lens of validation: recent advances in assessing and improving the quality of genomes assembled from metagenomes." Briefings in bioinformatics.

Paul, J. H. (2008). "Prophages in marine bacteria: dangerous molecular time bombs or the key to survival in the seas?" The ISME journal **2**(6): 579.

Peng, Y., H. C. Leung, S.-M. Yiu and F. Y. Chin (2012). "IDBA-UD: a de novo assembler for single-cell and metagenomic sequencing data with highly uneven depth." Bioinformatics **28**(11): 1420-1428.

Roux, S., J. B. Emerson, E. A. Eloe-Fadrosh and M. B. Sullivan (2017). "Benchmarking viromics: an in silico evaluation of metagenome-enabled estimates of viral community composition and diversity." PeerJ **5**: e3817.

Roux, S., S. J. Hallam, T. Woyke and M. B. Sullivan (2015). "Viral dark matter and virus–host interactions resolved from publicly available microbial genomes." Elife **4**: e08490.

Roux, S., N. E. Solonenko, V. T. Dang, B. T. Poulos, S. M. Schwenck, D. B. Goldsmith, M. L. Coleman, M. Breitbart and M. B. Sullivan (2016). "Towards quantitative viromics for both double-stranded and single-stranded DNA viruses." PeerJ **4**: e2777.

704      Schürch, A. C., D. Schipper, M. A. Bijl, J. Dau, K. B. Beckmen, C. M. Schapendonk, V. S. Raj, A. D.
705      Osterhaus, B. L. Haagmans and M. Tryland (2014). "Metagenomic survey for viruses in Western
706      Arctic caribou, Alaska, through iterative assembly of taxonomic units." PLoS One **9**(8): e105227.
707      Sczyrba, A., P. Hofmann, P. Belmann, D. Koslicki, S. Janssen, J. Dröge, I. Gregor, S. Majda, J.
708      Fiedler and E. Dahms (2017). "Critical assessment of metagenome interpretation—a benchmark of
709      metagenomics software." Nature methods **14**(11): 1063.
710      Shkoporov, A. N., F. J. Ryan, L. A. Draper, A. Forde, S. R. Stockdale, K. M. Daly, S. A. McDonnell,
711      J. A. Nolan, T. D. Sutton and M. Dalmasso (2018). "Reproducible protocols for metagenomic analysis
712      of human faecal phageomes." Microbiome **6**(1): 68.
713      Simpson, J. T., K. Wong, S. D. Jackman, J. E. Schein, S. J. Jones and I. Birol (2009). "ABySS: a
714      parallel assembler for short read sequence data." Genome research: gr. 089532.089108.
715      Smits, S. L., R. Bodewes, A. Ruiz-Gonzalez, W. Baumgärtner, M. P. Koopmans, A. D. Osterhaus and
716      A. C. Schürch (2014). "Assembly of viral genomes from metagenomes." Frontiers in microbiology **5**:
717      714.
718      Solden, L., K. Lloyd and K. Wrighton (2016). "The bright side of microbial dark matter: lessons
719      learned from the uncultivated majority." Current opinion in microbiology **31**: 217-226.
720      Vázquez-Castellanos, J. F., R. García-López, V. Pérez-Brocal, M. Pignatelli and A. Moya (2014).
721      "Comparison of different assembly and annotation tools on analysis of simulated viral metagenomic
722      communities in the gut." BMC genomics **15**(1): 37.
723      Vollmers, J., S. Wiegand and A.-K. Kaster (2017). "Comparing and evaluating metagenome assembly
724      tools from a microbiologist's perspective-Not only size matters!" PloS one **12**(1): e0169662.
725      Warwick-Dugdale, J., N. Solonenko, K. Moore, L. Chittick, A. C. Gregory, M. J. Allen, M. B.
726      Sullivan and B. Temperton (2018). "Long-read metagenomics reveals cryptic and abundant marine
727      viruses." bioRxiv.
728      Yang, X., P. Charlebois, S. Gnerre, M. G. Coole, N. J. Lennon, J. Z. Levin, J. Qu, E. M. Ryan, M. C.
729      Zody and M. R. Henn (2012). "De novo assembly of highly diverse viral populations." BMC
730      genomics **13**(1): 475.
731      Zerbino, D. and E. Birney (2008). "Velvet: algorithms for de novo short read assembly using de
732      Bruijn graphs." Genome research: gr. 074492.074107.

733

734

## **Additional files**

736      Additional file 1. (html) Simulated Virome MetaQUAST output.

737      Additional file 2. (html) Mock virome A MetaQUAST output.

738      Additional file 3. (html) Mock virome B MetaQUAST output.

739      Additional file 4. (html) Q33 spiked virome QUAST output.

740      Additional file 5. (xls)

741      S. Table 1: Spearman correlation values from the relationships of indel, mismatch and misassembly
742      counts, recovered genome fraction, abundance and total proportion of genomic repeats within the
743      Simulated virome. *GF – Recovered genome fraction

744      Additional file 5. (xls)

745      Supplementary Table 2: Linear modelling correlation values comparing recovered genome fraction,

746      total proportion of genomic repeats and abundance for the Simulated virome.

747      Additional file 5. (xls)

748    Supplementary Table 3: Spearman correlation values from the relationships of inverted, tandem,

749    palindromic and total repeats, abundance and the number of contigs generated by each assembler for

750    the Simulated virome.

751    Additional file 5. (xls)

752    Supplementary Table 4: (A) Ranking table comparing recovered genome fraction and contig numbers

753    for assemblers which recovered at least 50% of the total genome fraction. (B) Ranking table of indel,

754    mismatch and misassembly counts per 100kb, normalised to the number of genomes recovered to at

755    least 50%.

756    Additional file 5. (xls)

757    Supplementary Table 5: Number of aligned and unaligned contigs generated by each assembler for
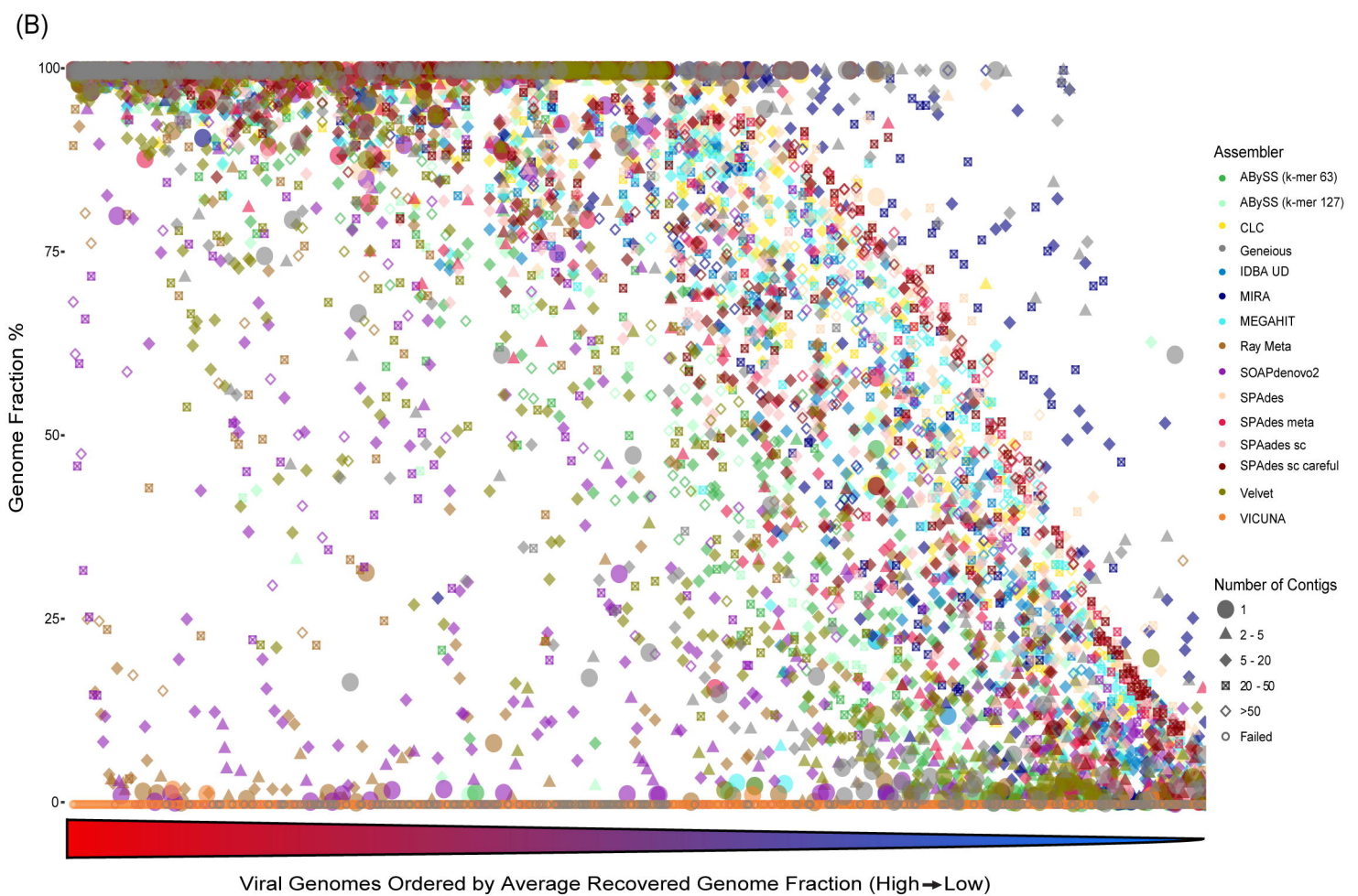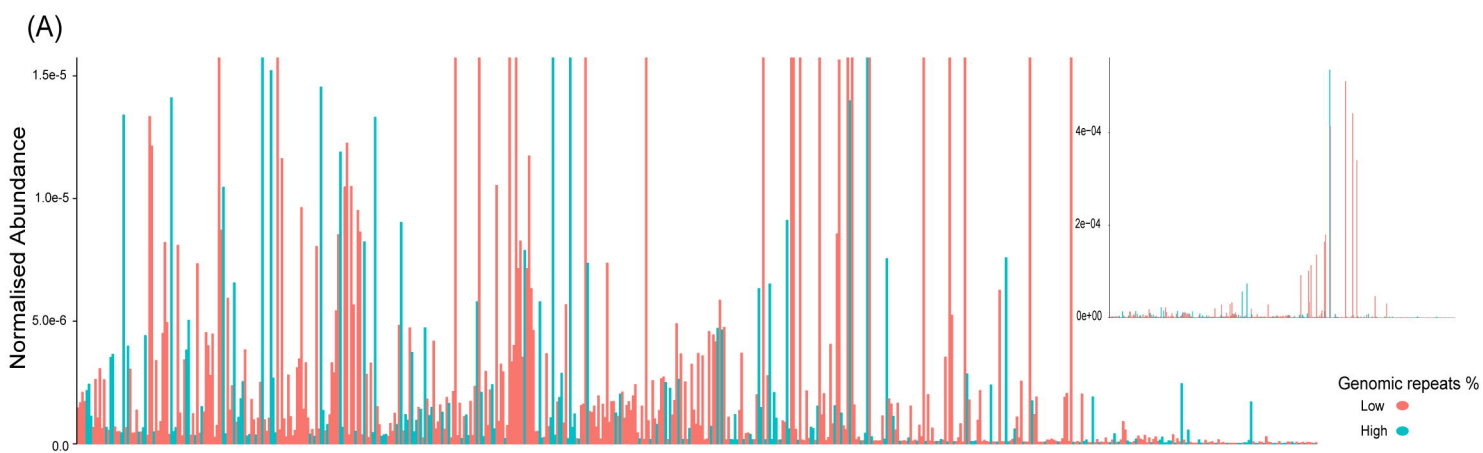
758    Mock Community A.

759    Additional file 5. (xls)

760    Supplementary Table 6: Number of aligned and unaligned contigs generated by each assembler for

761    Mock Community B.

762    Additional file 6. Supplementary Figure 1. Analysis of recovered genome fraction and indel/mismatch

763    counts for Mock communities A and B. Triangles represent N/A values for mismatches and indels

764    caused by assembly failures.

765

766

(A)

Normalised Abundance

Genomic repeats %
Low
High

(B)

Genome Fraction %

Assembler
ABySS (k-mer 63)
ABySS (k-mer 127)
CLC
Geneious
IDBA UD
MIRA
MEGAHIT
Ray Meta
SOAPdenovo2
SPAdes
SPAdes meta
SPAades sc
SPAdes sc careful
Velvet
VICUNA

Number of Contigs
1
2 - 5
5 - 20
20 - 50
>50
Failed

Viral Genomes Ordered by Average Recovered Genome Fraction (High→Low)

90% Genome covearge with a single contig

Number of Genomes

200

150

100

50

0

ABySS (k-mer 63), ABySS (k-mer 127), CLC, Geneious, IDBA UD, MEGAHIT, MIRA, Ray Meta, SOAPdenovo2, SPAdes, SPAdes meta, SPAades sc, SPAdes sc careful, Velvet, VICUNA

**A)**