**A peer-reviewed version of this preprint was published in PeerJ on 7 February 2018.**

# DNA-spiking in viral metagenome sequencing: A new method with low bias

**Geert Cremers** [1] , **Lavinia Gambelli** [1] , **Theo van Alen** [1] , **Laura van Niftrik** [1] , **Huub Op den Camp** [Corresp. 1]

[1] Microbiology, Radboud University Nijmegen

Corresponding Author: Huub Op den Camp
Email address: h.opdencamp@science.ru.nl

With the emergence of Next Generation Sequencing, major advances were made with regard to identifying viruses in natural environments. However, bioinformatical research on viruses is still limited because of the low amounts of viral DNA that can be obtained for analysis. To overcome this limitation, DNA is often amplified with multiple displacement amplification (MDA), which causes an unavoidable bias. Here, we describe a DNA-spiking method to avoid the bias that is created when using amplification of DNA before metagenome sequencing. To obtain sufficient DNA for sequencing, a bacterial 16S rRNA gene was amplified and the obtained DNA was spiked to a DNA sample containing DNA from a bacteriophage population before sequencing using Ion Torrent technology. After sequencing, the 16S rRNA gene reads DNA was removed by mapping to the Silva database. The new DNA-spiking method was compared with the MDA technique. When MDA was applied, the overall GC content of the reads showed a bias towards lower GC%, indicating a change in composition of the DNA sample. Assemblies using all available reads from both MDA and the DNA-spiked samples resulted in six viral genomes. All six genomes could be almost completely retrieved (97.9%-100%) when mapping the reads from the DNA-spiked sample to those six genomes . In contrast, 6.3%-77.7% of three viral genomes were covered by reads obtained using the MDA amplification method and only three were nearly fully covered (97.4%-100%). The new method provides a simple and inexpensive protocol with very low bias in sequencing of metagenomes for which low amounts of DNA are available.

# 1 DNA-spiking in viral metagenome sequencing: A new
# 2 method with low bias

3 Geert Cremers, Lavinia Gambelli, Theo van Alen, Laura van Niftrik & Huub J.M.
4 Op den Camp*

5 Department of Microbiology, Institute for Water and Wetland Research, Faculty of Science,
6 Radboud University, Heyendaalseweg 135, NL 6525 AJ Nijmegen, the Netherlands.

7 *Corresponding author:

8 h.opdencamp@science.ru.nl

9 Heyendaalseweg 135, NL 6525 AJ Nijmegen, The Netherlands

10    Key words: virus, bacteriophage, metagenome, metavirome, DNA spiking, multiple
11    displacement amplification

## Abstract

13    With the emergence of Next Generation Sequencing, major advances were made with regard to
14    identifying viruses in natural environments. However, bioinformatical research on viruses is still
15    limited because of the low amounts of viral DNA that can be obtained for analysis. To overcome
16    this limitation, DNA is often amplified with multiple displacement amplification (MDA), which
17    causes an unavoidable bias. Here, we describe a DNA-spiking method to avoid the bias that is
18    created when using amplification of DNA before metagenome sequencing. To obtain sufficient
19    DNA for sequencing, a bacterial 16S rRNA gene was amplified and the obtained DNA was
20    spiked to a DNA sample containing DNA from a bacteriophage population before sequencing
21    using Ion Torrent technology. After sequencing, the 16S rRNA gene reads DNA was removed by
22    mapping to the Silva database. The new DNA-spiking method was compared with the MDA
23    technique. When MDA was applied, the overall GC content of the reads showed a bias towards
24    lower GC%, indicating a change in composition of the DNA sample. Assemblies using all
25    available reads from both MDA and the DNA-spiked samples resulted in six viral genomes. All
26    six genomes could be almost completely retrieved (97.9%-100%) when mapping the reads from
27    the DNA-spiked sample to those six genomes. In contrast, 6.3%-77.7% of three viral genomes
28    were covered by reads obtained using the MDA amplification method and only three were nearly
29    fully covered (97.4%-100%). The new method provides a simple and inexpensive protocol with
30    very low bias in sequencing of metagenomes for which low amounts of DNA are available.

## Introduction

32    Microbial research has been mainly culture-based since the work of Pasteur and Koch. This has
33    led to great improvement of our knowledge of the microbial and viral world. However, our
34    knowledge is probably still only the tip of the iceberg, as most of the microorganisms cannot be
35    cultured (Rosario & Breitbart, 2011). In recent years, there has been a greater focus on the
36    hidden bacterial and viral 'black matter' since the development of next generation sequencing
37    (NGS) techniques which allow the determination of the microbial community without the need
38    for cultivation. Without the necessity of cultivation prior to sequencing, organisms that cannot be
39    cultured under artificial conditions are now being sequenced in increasing numbers. This is
40    especially true for bacteria; sequencing viral black matter from environmental samples is still
41    hampered by a variety of factors. Besides the obvious problem that not all viruses are DNA
42    viruses (Steward et al., 2013), there is also the challenge of the low quantity of DNA that can be
43    retrieved from viruses. Although viruses outnumber bacteria 5 to 25 times in numbers (Fuhrman,
44    1999; Clokie et al., 2011), the fact that viruses have on average a significantly smaller genome
45    size means that the viral DNA yield from any given sample is significantly lower compared to
46    bacterial DNA (Brum & Sullivan, 2015). Because of this, DNA is often extracted using DNA

47  extraction methods optimized for virions gathered from a large amount (20-200 L) of sample
48  (Breitbart et al., 2002;Thurber et al., 2009;  Duhaime, 2012; Steward et al., 2013).

49  Since sampling large amounts is not always possible and because there is loss of DNA in every
50  step of DNA isolation protocols, the final yield of viral DNA remains low. Therefore, the low
51  amount of DNA must be amplified before sequencing. Several methods are available (Duhaime
52  et al., 2012; Brum & Sullivan, 2015), with the Multiple Displacement Amplification (MDA)
53  method being the most widely used. However, MDA has a few drawbacks. Because it is
54  amplification, it unavoidably introduces a bias (Kim & Bae, 2011; Marine et al., 2014) and
55  furthermore information of relative abundance in the original sample is lost (Yilmaz et al., 2010).
56  In practical terms, there is a high potential of cross contamination throughout the lab, which is
57  particularly problematic in a laboratory where viral work is the main area of research.

58  In view of the growing interest in the impact of viruses on the ecology and the natural world,
59  there is a growing need for an easy method or technology that is bias-limited and preferably
60  easy, inexpensive and readily available. However, the recurring problem on most sequencing
61  platforms is the relative large amount of DNA needed for the preparation of DNA fragments for
62  sequencing compared to the typically low yield of DNA from viral population samples. As
63  discussed above, two ways to obtain more DNA is to collect more samples or apply
64  amplification. A theoretical third option for viruses would be to artificially raise the DNA
65  concentration by adding DNA that is naturally non-occurring in any virus. Prime candidate
66  would be the bacterial 16S ribosomal gene as it has not been found in any known virus up to
67  date. In this report we describe a method for metagenome sequencing with the Ion Torrent
68  Personal Genome Machine in which we spiked low amounts ($\approx$ 0.1 ng) of viral DNA from a
69  bacteriophage population with 16S ribosomal DNA and compare this with traditional MDA
70  amplification.

## Material and methods

### Sample collection and DNA extraction

73  The bacteriophage population used for sequencing was obtained from a *Methylomirabilis oxyfera*
74  bioreactor enrichment culture (Gambelli et al., 2016). Bioreactor material was collected over a
75  period of about three months, stored at 4°C and viral particles were obtained as described before
76  (Gambelli et al., 2016). Briefly, the aggregated microbial biomass was disrupted to free the
77  bacteriophages and viral particles were precipitated using PEG8000 (Guo et al., 2012). Free
78  bacteriophages present in the bioreactor supernatant medium and not within the bacterial
79  aggregates were precipitated by iron chloride flocculation (Cunningham et al., 2015).
80  The two samples obtained by iron chloride flocculation and PEG 8000 precipitation were pooled
81  together and bacteriophages were concentrated by ultracentrifugation (Optima XE90, Beckman-
82  Coulter; Rotor: Type 90 Ti, Beckman-Coulter) at 77,000x g at 4°C for 1h. The pellet was
83  resuspended in 1 ml of supernatant and the total DNA was extracted according to the protocol

84 published by Thurber et al. (2009). Using the Qubit dsDNA HS assay kit (Thermo Scientific,
85 Waltham, USA), the extracted DNA was quantified at 0.2 ng DNA.

## Library preparation and sequencing

### MDA method

88 For the MDA method, 0.1 ng of viral metagenome DNA was amplified using the Illustra
89 GenomePhi HY DNA amplification kit (GE Healthcare, Piscataway, NJ, USA) as per
90 manufacturer's protocol. The DNA was cleaned using GeneJET plasmid Miniprep kit
91 (Fermentas, Amherst, USA) according to manufacturer's protocol, except for step one, in which
92 200 μl of DEPC was used instead of lysis-buffer. This first amplification round yielded 15 ng of
93 DNA (referred as 1 x MDA sample), after which 10 ng of DNA was used for a second
94 amplification round (referred as 2 x MDA sample) which resulted in a final yield of 5.4 μg DNA.
95 Sequence library preparation was started with 5 ng of 1 x MDA and with 100 ng 2 x MDA
96 amplified DNA. Both samples were sheared for 6 cycles (1 min on, 1 min off) on the Bioruptor®
97 Standard (Diagenode Liege, Belgium). After shearing, the samples were cleaned using a 1:1
98 volume ratio with AMPure XP beads (Beckman Coultier, High Wycombe, UK) and further
99 prepared for sequencing as per manufacturer's protocol (IonXpress Plus gDNA fragment library
100 preparation Rev C.0, Life, Carlsbad, USA).

### DNA-spiking method

102 For the DNA-spiking method, 0.1 ng of viral metagenome DNA was spiked with 43.2 ng of
103 amplified 16S rRNA DNA from "*Candidatus* Kuenenia stuttgartiensis" (GenBank CT573071)
104 (referred as DNA-spiked  sample) and sheared using the Bioruptor for 6 cycles (1 min on, 1 min
105 off) and prepared according to manufacturer's protocol (IonXpress Plus gDNA fragment library
106 preparation Rev C.0, Life).

107 The amplicons of the 16S rRNA gene were obtained by PCR of an in-house 16S rRNA gene
108 clone of "*Ca.* K. stuttgartiensis*"* using primers pla46 (5'-GGATTAGGCATGCAAGTC-'3) and
109 630R (5'-CAKAAAGGAGGTGATCC-'3) with the following settings: 5 min at 94˚C, followed
110 by 35 cycles of 40 s at 96˚C, 40 s at 49˚C and 1 min at 72˚C and finalised with an elongation step
111 of 5 min at 72˚C. After amplification, the sample was purified from non 16S ribosomal DNA by
112 excision and re-extraction of the DNA from a 0.9% gel (v/w) using the GeneJET gel extraction
113 kit (Thermo Scientific, Waltham, USA) according to manufacturer's protocol. The final
114 concentration (10.8 ng/μl) was measured using the Qubit dsDNA HS assay kit (Thermo
115 Scientific, Waltham, MA, USA).

116 Both MDA and DNA-spiking method samples were sequenced using the Personal Genome
117 Machine Ion Torrent (Thermo Scientific, Waltham, MA, USA) as per manufacturer's protocol. 1
118 x MDA was sequenced twice, once on a 314v2 chip and once on a 318v2 chip, and reads were
119 combined. Sample 2 x MDA and the DNA-spiked sample were run on a 318v2 chip. All samples
120 were constructed using the Ion PGM™ Sequencing 400 Kit and Ion PGM™ Template OT2 400

121    kit and sequenced with 850 flow cycles. The raw sequence data were submitted to the European

122    Nucleotide Archive and received accession number PRJEB20134

123    (http://www.ebi.ac.uk/ena/data/view/PRJEB20134).

## Bioinformatics

### Trimming

126    After sequencing, reads from the DNA-spiked sample were trimmed with default quality

127    settings, size 25 to 375 bp. The 1 x MDA sample was sequenced twice and reads from the first

128    run were trimmed with default quality settings, size 25 to 325 bp. Reads from the second 1 x

129    MDA sample run were trimmed with default quality settings, size 25 to 400 bp and 15 bp on 5'

130    end. Reads from the 2 x MDA sample were trimmed with default quality settings, size 25 to 375

131    bp, using CLC genomics workbench v. 8 (CLCbio, Aarhus, Denmark).

### Recovery of viral contigs

133    To remove genomic DNA from the most abundant microorganism in the bioreactor, the trimmed

134    reads from the DNA-spiked sample were mapped against the genome of "*Candidatus*

135    Methylomirabilis oxyfera" (Ettwig et al., 2010) (length 0.5, similarity 0.85) and the unmapped

136    reads were assembled (length 0.5, similarity 0.95, word size 22, bubble size 276), using CLC

137    genomics workbench v. 8 (CLCbio, Aarhus, Denmark).

138    The obtained contigs were subsequently mapped against the SILVA database 16S rRNA v119

139    (Yarza et al., 2008) (length 0.5, similarity 0.7%) and contigs that mapped to "*Ca.*

140    K. stuttgartiensis*"* were removed from the database, resulting in 4088 remaining contigs. These

141    contigs were checked with ESOM (Ultsch & Moerchen, 2005) (default settings) and from this,

142    seven clustering contigs with a high coverage were obtained and reassembled with SPAdes

143    (v.3.5.0) (Bankevich et al., 2012) using the 'trusted-contigs' and 'careful' settings for those seven

144    contigs.

145    The reads that were used in the SPAdes re-assembly, were obtained by mapping the spiked DNA

146    to the SILVA database 16S rRNA v119 (length 0.5, similarity 0.7%). The reads that did not map

147    were used.

148    Reassembly with SPAdes created 2094 contigs, 14 of which were bigger than 5000 bp. From this

149    set of contigs five putative viral genomes could be extracted (197 kbp, 86 kbp, 71 kbp, 41 kbp

150    and 17 kbp).

151    Assembly of the combined MDA data (Length fraction = 0.5, similarity fraction = 0.9, minimum

152    contig length = 1,500) resulted in 689 contigs, ranging from 130,897 to 1,504 nucleotides. With

153    the use of ESOM, one more putative viral genome was identified (42 kbp).

154 **Differential coverage**

155 For differential coverage, reads from the DNA-spiked sample were mapped against the SILVA
156 database 16S rRNA v119 (length 0.5, similarity 0.7%). Unmapped reads were combined with the
157 reads from the second 1 x MDA run and assembled (length= 0.5, similarity 0.9, word=35, bubble
158 size 271) with. Subsequent mapping of each read set (length 0.5, similarity 0.8%) was performed
159 against the assembled contigs. The depth of both sets was plotted against one another.

160 **Horizontal coverage**

161 To assess how much of each virus was present in each set, the trimmed reads from the three sets
162 were mapped against the six putative viral genomes (length 0.5, similarity 0.95%). The number
163 of mapped reads and total length was normalised to the size of the dataset and the length of virus,
164 respectively.

## Results

166 Viral DNA extracted from bioreactor biomass was sequenced following two different
167 approaches: MDA amplification (two samples) and non-amplified DNA spiked with bacterial
168 16S rRNA gene DNA. This resulted in three datasets comprising of a total of 774,366 trimmed
169 reads for 1 x MDA and 187,178 trimmed reads for 2 x MDA. After the reads from the non-
170 amplified spiked DNA were trimmed and mapped against the SILVA database, the final number
171 of reads left was 529,481. After trimming, GC graphs were created showing the GC distribution
172 of each dataset (Fig. 1). Clearly visible is the shift in GC content from non-amplified spiked
173 DNA to amplified DNA (1 x MDA, 2 x MDA). The non-amplified set, DNA-spiking method
174 (Fig. 1A), shows a large peak at 63-64% GC content and a shoulder at around 42%. For the
175 MDA method, after the first run of amplification, a large increase of the shoulder at 42% GC
176 content is visible (Fig. 1B) which becomes even more pronounced after another round of
177 amplification (Fig. 1C). The large peak at 63-64% GC content shifts to 57% after one round of
178 amplification. The second round of MDA does not lower the GC content of this peak, but the
179 relative amount of sequences is lowered within the sample.

180 A total of six different viral sequences (putative genomes) could be assembled using the reads
181 from all three datasets by a combination of various methods (see Materials and methods). The
182 DNA-spiking method resulted in five viral genomes while the MDA set only resulted in one. The
183 length of the viral genomes ranged from 197 kbp to 17 kbp with GC contents from 67% to 35%.
184 Five of the viral genomes (197 kbp, 86 kbp, 71 kbp, 42 kbp and 17 kbp) contained the same
185 sequence on each end of the contig indicating a full circular genome. Fig. 2 shows the
186 comparison of the total amount of virus genomes that could be obtained with the individual
187 datasets. Fig. 2a shows that all six virus genomes could be recovered in nearly complete length
188 from the DNA-spiked dataset. In contrast, MDA amplification clearly lowers the percentage of
189 viral genomes that can be recovered. Fig. 2b shows the percentage of reads from each viral
190 metagenome sample mapping to the six virus genomes. With the MDA amplified samples, the

191   percentages of reads mapping is low with a comparable relative abundance. The non-amplified
192   sample not only shows a lot more difference in relative abundance, but the total amount of
193   mapped reads comprises over half of the original dataset. When the MDA method is used, the
194   number of mapped reads drops to lower than 3%. Fig. 2c shows the average coverage of each
195   viral genome after mapping the reads of the individual datasets. The data are comparable to Fig.
196   2b. Using the DNA-spiked dataset, high coverage (> 35) is found for five out of six viral
197   genomes. The MDA datasets show very low coverage and only give better results for the 42 kbp
198   virus.

199   Assembly of the reads from the DNA-spiked sample and the 1 x MDA sample resulted in a total
200   of 1644 contigs. Differential coverage of these contigs, comparing the two datasets is shown in
201   Fig. 3. From the figure it is clear that the depth for each set of contigs differs completely since
202   two similar sets would result in contigs mapping on a diagonal straight line. The two datasets
203   show an almost perpendicular distribution. When looking at the location of contigs with a high
204   and low GC content it becomes clear that high GC% contigs (green) are present within the DNA-
205   spiked dataset but are low in the amplified sample while low GC% contigs (pink) are much more
206   abundant in the 1 x MDA dataset. The figure also demonstrates the wide range in sequencing
207   depth when no amplification was applied (maximum around 380) in clear contrast to
208   amplification, as the sample coverage was lowered tenfold, with a maximum around 38.

## Discussion

210   Here we describe a new method for sequencing low amount of DNA (viral DNA) and compare
211   this to a traditional MDA amplification method. Like in previous reports (Kim & Bae, 2011;
212   Marine et al., 2014), our experiments show the bias that MDA amplification of metavirome or
213   metagenome DNA introduces into the dataset. In the extreme case of applying MDA twice, the
214   dataset is completely changed with major consequences for results obtained, which as such may
215   not reflect the sampled ecosystem. This is emphasising once more the caution that has to be
216   taken when utilizing the MDA technique for environmental DNA samples. Where annotation of
217   the contigs obtained from the unamplified DNA-spiked sample showed a majority of virus
218   related genes, the MDA datasets showed a general shift towards bacterial genes (data not
219   shown). In the MDA datasets, the genes match most closely to WS6-like bacteria, a group of
220   unusually small bacteria with small genomes (~0.1 Mbp) and (in this case) low GC content
221   (Speth et al, 2016). Considering the small size of these bacteria it is very likely they were not
222   removed during the filtration process used for selecting viral particles.

223   Our results illustrate that specifically DNA with a low GC content is favoured in MDA, which
224   can lead to a considerable shift of GC content causing a severe underestimation of the quantity of
225   viruses with a higher GC percentage in the sample. Viruses with higher GC content can be easily
226   overlooked.

227  Differential binning of spiked DNA and the 1 x MDA not only shows a remarkable difference
228  between both samples, but also demonstrates the loss of information about the abundance, as the
229  range of this abundance of the different viruses shifts from 380 to 38. This shift in depth means
230  that using MDA as an abundance indicator would have been impossible in this example and as a
231  consequence we used the abundance of the spiked DNA in recovery of the viral genomes.

232  A nice side effect of the spiking method is the lowering of contamination risk with DNA from
233  the lab. By choosing the 16S rRNA gene of a microorganism which is very uncommon to the
234  laboratory for spiking, contamination can easily be recognized and filtered out. One could even
235  use eukaryotic DNA encoding the 18S rRNA gene in a prokaryotic-based lab and vice versa.

236  In this report we described an extremely low input method (only 0.1 ng of DNA needed) for viral
237  metagenome sequencing that is unbiased, inexpensive, easy and readily available for any lab
238  with sequence facilities and that can possibly be extended for other non-16S/18S containing
239  DNA like plasmids, or other Next Generation Sequencing platforms like MinIon or Illumina.


## Conclusions

241  When dealing with low quantities of DNA for Next Generation Sequencing, multiple
242  displacement amplification (MDA) of the DNA from the sample might not be the method of
243  choice to obtain enough starting material. We propose to use the DNA-spiking method, where
244  one adds DNA with a known sequence, as a more valid alternative. The reads resembling the
245  added DNA can be easily discarded afterwards. This new method is not only technically feasible,
246  but also results in a very low bias in the dataset compared to the MDA method. The obvious bias
247  of the MDA method has to be taken into consideration since it may cause major shifts in the
248  DNA profile of an ecosystem.


## Acknowledgements

## References

255  Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, Lesin VM, Nikolenko
256  SI, Pham S, Prjibelski AD, Pyshkin AV, Sirotkin AV, Vyahhi N, Tesler G, Alekseyev MA, and
257  Pevzner PA. 2012. SPAdes: A new genome assembly algorithm and its applications to single-
258  cell sequencing. Journal of Computational Biology 19:455-477. 10.1089/cmb.2012.0021

259  Breitbart M, Salamon P, Andresen B, Mahaffy JM, Segall AM, Mead D, Azam F, and Rohwer F.
260  2002. Genomic analysis of uncultured marine viral communities. Proceedings of the National
261  Academy of Sciences of the United States of America 99:14250-14255.
262  10.1073/pnas.202488399

263  Brum JR, and Sullivan MB. 2015. Rising to the challenge: accelerated pace of discovery
264  transforms marine virology. Nature Reviews Microbiology 13:147-159. 10.1038/nrmicro3404

265  Clokie MR, Millard AD, Letarov AV, and Heaphy S. 2011. Phages in nature. Bacteriophage
266  1:31-45. 10.4161/bact.1.1.14942

267  Cunningham BR, Brum JR, Schwenck SM, Sullivan MB, and John SG. 2015. An inexpensive,
268  accurate, and precise wet-mount method for enumerating aquatic viruses. Applied and
269  Environmental Microbiology 81:2995-3000. 10.1128/Aem.03642-14

270  Duhaime MB, Deng L, Poulos BT, and Sullivan MB. 2012. Towards quantitative metagenomics
271  of wild viruses and other ultra-low concentration DNA samples: a rigorous assessment and
272  optimization of the linker amplification method. Environmental Microbiology 14:2526-2537.
273  10.1111/j.1462-2920.2012.02791.x

274  Ettwig KF, Butler MK, Le Paslier D, Pelletier E, Mangenot S, Kuypers MM, Schreiber F, Dutilh
275  BE, Zedelius J, de Beer D, Gloerich J, Wessels HJ, van Alen T, Luesken F, Wu ML, van de Pas-
276  Schoonen KT, Op den Camp HJM, Janssen-Megens EM, Francoijs KJ, Stunnenberg H,
277  Weissenbach J, Jetten MSM, and Strous M. 2010. Nitrite-driven anaerobic methane oxidation by
278  oxygenic bacteria. Nature 464:543-548. 10.1038/nature08883

279  Fuhrman JA. 1999. Marine viruses and their biogeochemical and ecological effects. Nature
280  399:541-548. Doi 10.1038/21119

281  Gambelli L, Cremers G, Mesman R, Guerrero S, Dutilh BE, Jetten MSM, Op den Camp HJM,
282  and van Niftrik L. 2016. Ultrastructure and viral metagenome of bacteriophages from an
283  anaerobic methane oxidizing Methylomirabilis bioreactor enrichment culture. Frontiers in
284  Microbiology 7:1740. 10.3389/fmicb.2016.01740

285  Guo P, El-Gohary Y, Prasadan K, Shiota C, Xiao XW, Wiersch J, Paredes J, Tulachan S, and
286  Gittes GK. 2012. Rapid and simplified purification of recombinant adeno-associated virus.
287  Journal of Virological Methods 183:139-146. 10.1016/j.jviromet.2012.04.004

288  Kim KH, and Bae JW. 2011. Amplification methods bias metagenomic Libraries of uncultured
289  single-stranded and double-stranded DNA viruses. Applied and Environmental Microbiology
290  77:7663-7668. 10.1128/Aem.00289-11

291  Marine R, McCarren C, Vorrasane V, Nasko D, Crowgey E, Polson SW, and Wommack KE.
292  2014. Caught in the middle with multiple displacement amplification: the myth of pooling for

293  avoiding multiple displacement amplification bias in a metagenome. Microbiome 2. Artn
294  310.1186/2049-2618-2-3

295  Rosario K, and Breitbart M. 2011. Exploring the viral world through metagenomics. Current
296  Opinion in Virology 1:289-297. 10.1016/j.coviro.2011.06.004

297  Speth DR, In 't Zandt MH, Guerrero-Cruz S, Dutilh BE, and Jetten MS. 2016. Genome-based
298  microbial ecology of anammox granules in a full-scale wastewater treatment system. Nature
299  Communications 7:11172. 10.1038/ncomms11172

300  Steward GF, Culley AI, Mueller JA, Wood-Charlson EM, Belcaid M, and Poisson G. 2013. Are
301  we missing half of the viruses in the ocean? ISME Journal 7:672-679. 10.1038/ismej.2012.121

302  Thurber RV, Haynes M, Breitbart M, Wegley L, and Rohwer F. 2009. Laboratory procedures to
303  generate viral metagenomes. Nature Protocols 4:470-483. 10.1038/nprot.2009.10

304  Ultsch A, Moerchen F 2005. ESOM-Maps: tools for clustering, visualization, and classification
305  with Emergent SOM, Technical Report Dept. of Mathematics and Computer Science, University
306  of Marburg, Germany, No. 46

307  Yarza P, Richter M, Peplies J, Euzeby J, Amann R, Schleifer KH, Ludwig W, Glockner FO, and
308  Rossello-Mora R. 2008. The All-Species Living Tree project: a 16S rRNA-based phylogenetic
309  tree of all sequenced type strains. Systematic and Applied Microbiology. 31:241-250.
310  10.1016/j.syapm.2008.07.001

311  Yilmaz S, Allgaier M, and Hugenholtz P. 2010. Multiple displacement amplification
312  compromises quantitative analysis of metagenomes. Nature Methods 7:943-944.
313  10.1038/nmeth1210-943

315 **Legends to the figures.**

316

317 Figure 1: Distribution of reads obtained from Ion Torrent sequencing using three different
318 sample preparation methods based on their individual GC content in % of the total number of
319 reads in one sample. **A.** DNA-spiking method; **B.** 1 x MDA method; **C.** 2 x MDA method. GC
320 graphs were created using CLCgenomics workbench.

321 Figure 2:  Comparison of the viral reads from the individual datasets mapping to the six
322 assembled virus genomes (Green, DNA-spiked sample; Blue, 1 x MDA; Red, 2 x MDA). **A.**
323 Horizontal coverage of the viral genomes with reads from the individual datasets. **B.** Distribution
324 of the reads from the individual datasets over the assembled viral genomes. **C.** Depth (vertical
325 coverage) of the viral genomes with reads from the individual datasets as a measure of
326 abundance.

327 Figure 3: Differential coverage of the viral contigs assembled using a combination of the DNA-
328 spiked sample and the 1 x MDA sample with each individual read sets. Each circle represents a
329 contig present after assembly and the placement in the plot shows the abundance of the contig for
330 each read set. Two similar read sets would result in a diagonal straight line. GC content of the
331 different contigs is indicated as depicted in the colour scale. Three outliers caused by the MDA
332 amplification method are not shown in the plot.

## Figure 1(on next page)

Figure 1: Distributionof reads obtained from Ion Torrent sequencing using three different samplepreparation methods based on their individual GC content in % of the totalnumber of reads in one sample.

A. DNA-spiking method;  **B.**  1 x MDA method;  **C.**  2 x MDA method. GC graphs were created using CLCgenomics workbench.
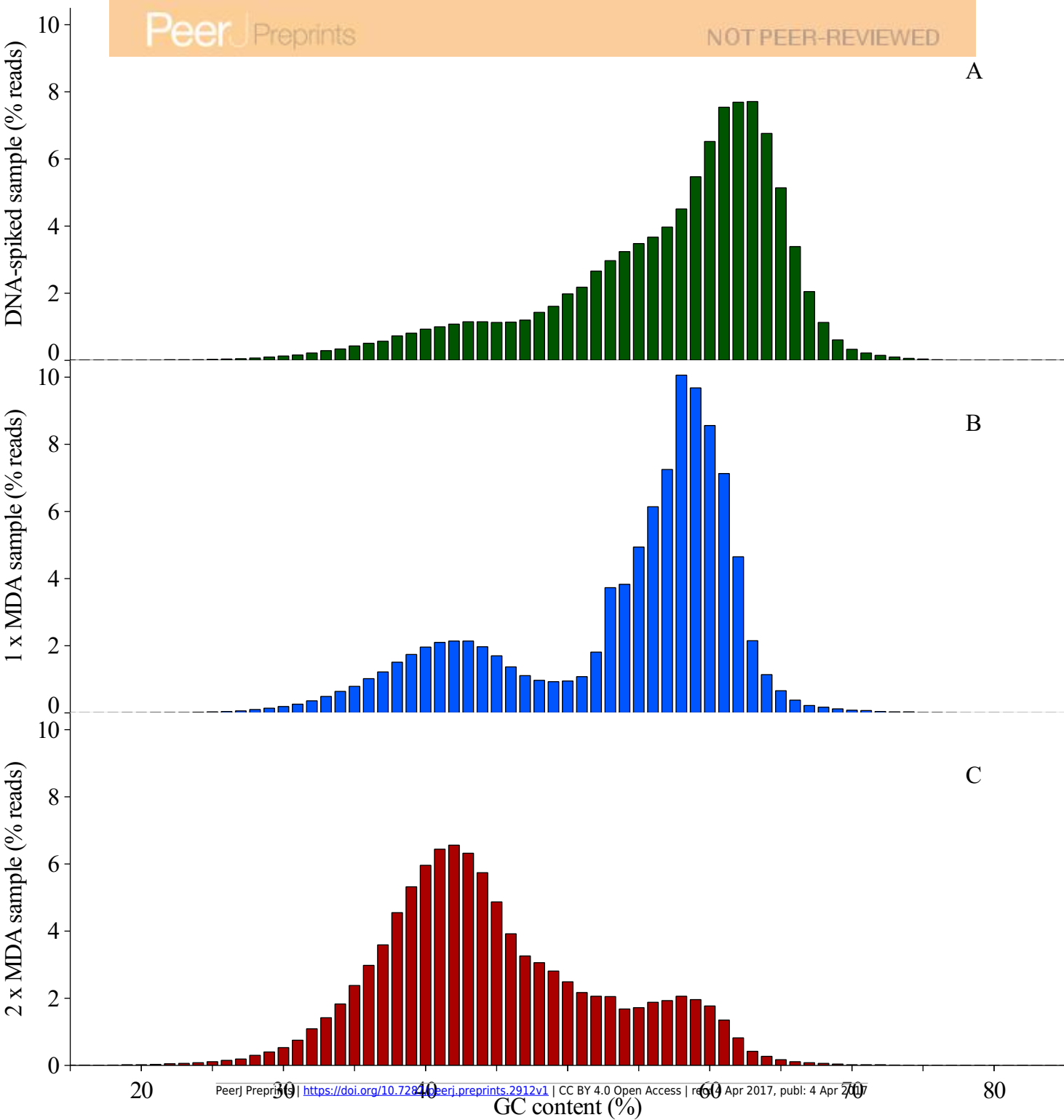
# Figure 2(on next page)

Figure 2: Comparisonof the viral reads from the individual datasets mapping to the six assembled virusgenomes.

(Green, DNA-spiked sample; Blue, 1 x MDA; Red, 2 x MDA).  **A.**  Horizontal coverage of the viral genomes with reads from the individual datasets.  **B.**  Distribution of the reads from the individual datasets over the assembled viral genomes.  **C.**  Depth (vertical coverage) of the viral genomes with reads from the individual datasets as a measure of abundance.
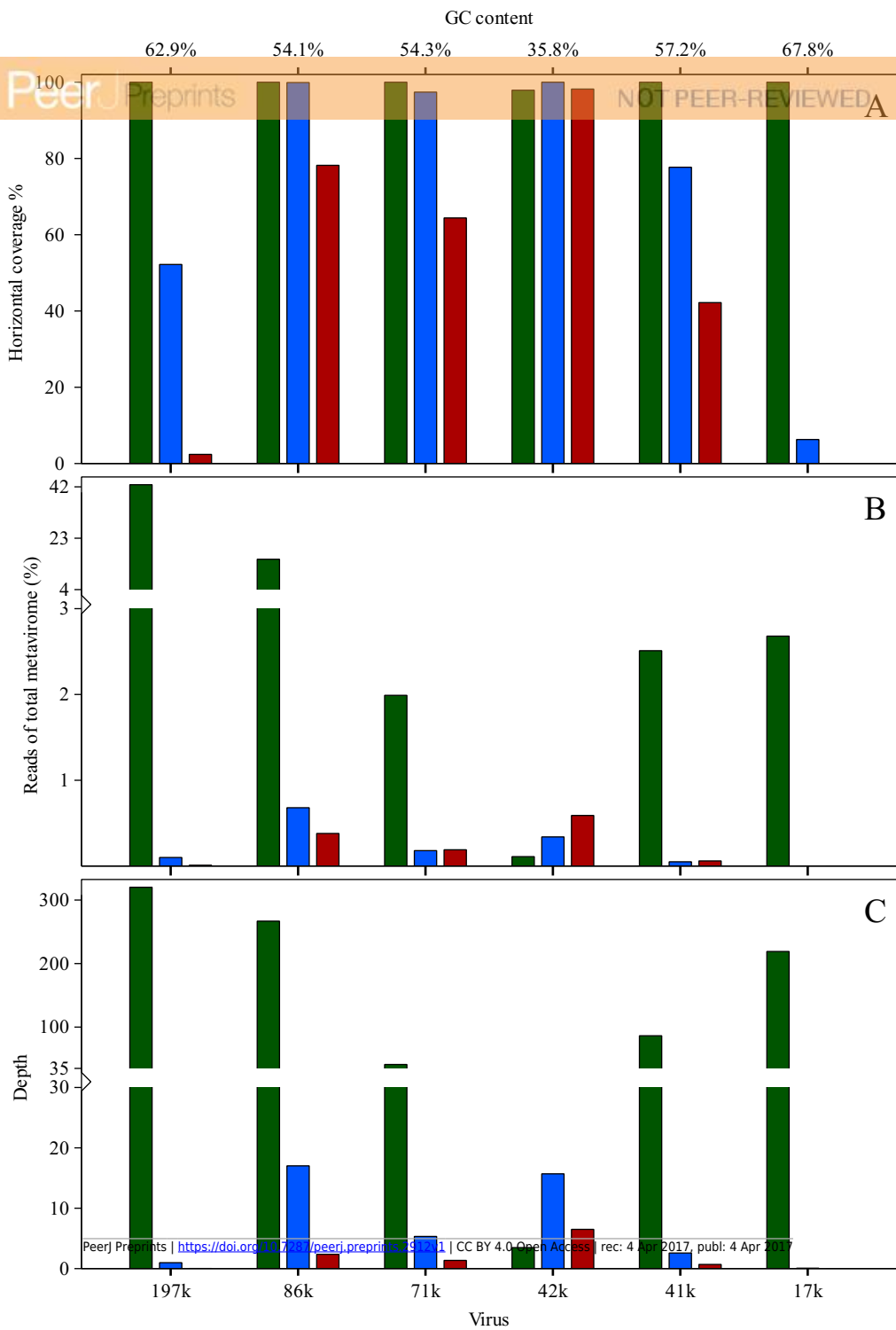
# Figure 3(on next page)

Figure3: Differential coverage of the viral contigs assembled using a combination of theDNA-spiked sample and the 1 x MDA sample with each individual read sets.

Each circle represents a contig present after assembly and the placement in the plot shows the abundance of the contig for each read set. Two similar read sets would result in a diagonal straight line. GC content of the different contigs is indicated as depicted in the colour scale. Three outliers caused by the MDA amplification method are not shown in the plot.