

1 **CarrierSeq: a sequence analysis workflow for low-input nanopore sequencing**

2

3 Angel Mojarro^{1,*}, Julie Hachey², Gary Ruvkun³, Maria T. Zuber¹, and Christopher E. Carr^{1,3}

4

5 ¹Department of Earth, Atmospheric and Planetary Sciences, Massachusetts Institute of
6 Technology, Cambridge, MA, USA

7 ²ReadCoor, Cambridge, MA, USA

8 ³Department of Molecular Biology, Massachusetts General Hospital, Boston, MA, USA

9

10 *Address correspondence to:

11 Angel Mojarro

12 Massachusetts Institute of Technology

13 77 Massachusetts Ave

14 Cambridge, MA 02139

15 E-mail: mojarro@mit.edu

16

17

18

19

20

21

22

23

24 **Abstract**

25 **Motivation:** Long-read nanopore sequencing technology is of particular significance for
26 taxonomic identification at or below the species level. For many environmental samples, the total
27 extractable DNA is far below the current input requirements of nanopore sequencing, preventing
28 “sample to sequence” metagenomics from low-biomass or recalcitrant samples.

29 **Results:** Here we address this problem by employing carrier sequencing, a method to sequence
30 low-input DNA by preparing the target DNA with a genomic carrier to achieve ideal library
31 preparation and sequencing stoichiometry without amplification. We then use CarrierSeq, a
32 sequence analysis workflow to identify the low-input target reads from the genomic carrier. We
33 tested CarrierSeq experimentally by sequencing from a combination of 0.2 ng *Bacillus subtilis*
34 ATCC 6633 DNA in a background of 1 µg *Enterobacteria phage λ* DNA. After filtering of carrier,
35 low quality, and low complexity reads, we detected target reads (*B. subtilis*), contamination reads,
36 and “high quality noise reads” (HQNRs) not mapping to the carrier, target or known lab
37 contaminants. These reads appear to be artifacts of the nanopore sequencing process as they are
38 associated with specific channels (pores). By treating reads as a Poisson arrival process, we
39 implement a statistical test to reject data from channels dominated by HQNRs while retaining
40 target reads.

41 **Availability:** CarrierSeq is an open-source bash script with supporting python scripts which
42 leverage a variety of bioinformatics software packages on macOS and Ubuntu. Supplemental
43 documentation is available from Github - <https://github.com/amojarro/carrierseq>. In addition, we
44 have compiled all required dependencies in a Docker image available from -
45 <https://hub.docker.com/r/mojarro/carrierseq>.

46

47 **1 Introduction**

48 Environmental metagenomic sequencing poses a number of challenges. First, complex soil
49 matrices and tough-to-lyse organisms can frustrate the extraction of deoxyribonucleic acid (DNA)
50 and ribonucleic acid (RNA) (Lever et al., 2015). Second, low-biomass samples require further
51 extraction and concentration steps which increase the likelihood of contamination (Barton et al.,
52 2006). Third, whole genome amplification may bias population results (Sabina and Leamon, 2015)
53 while targeted amplification (e.g., 16S rRNA amplicon) may decrease taxonomic resolution
54 (Poretsky et al., 2014). To address these challenges, we have developed extraction protocols
55 compatible with low-biomass recalcitrant samples and difficult to lyse organisms (Mojarro,
56 Ruvkun, et al., 2017). These protocols, developed using tough-to-lyse spores of *Bacillus subtilis*,
57 allow us to achieve at least 5% extraction yield from a 50 mg sample containing 2×10^5 cells/g of
58 soil without centrifugation (Carr et al., 2017). Furthermore, in order to avoid possible amplification
59 biases and additional points of contamination, we have experimented with utilizing a genomic
60 carrier (*Enterobacteria phage λ*) to shuttle low-input amounts of target DNA (*B. subtilis*) through
61 library preparation and sequencing with ideal stoichiometry (Mojarro, Hachey, et al., 2017). This
62 approach has allowed us to detect down to 0.2 ng of *B. subtilis* DNA prepared with 1 μ g of Lambda
63 DNA using the Oxford Nanopore Technologies (ONT) MinION sequencer (supplementary data,
64 <https://www.ncbi.nlm.nih.gov/bioproject/398368>). Here we present CarrierSeq, a sequence
65 analysis workflow developed to identify target reads from a low-input sequencing run employing
66 a genomic carrier.

67

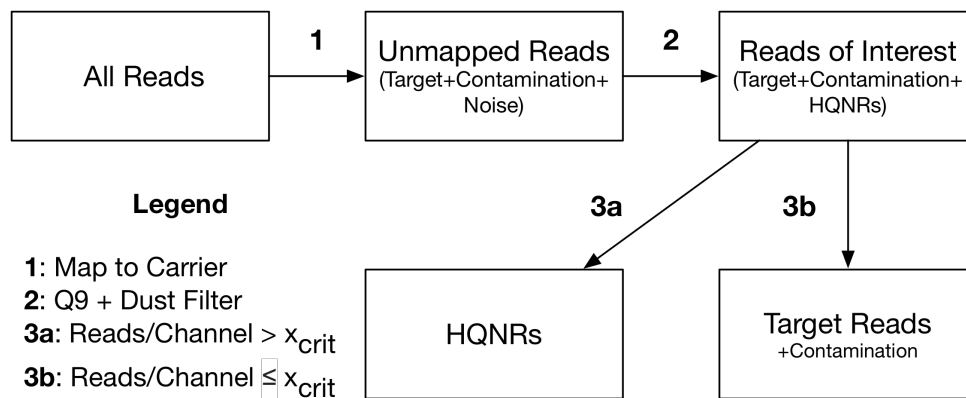
68

69

70 2 Methods

71 CarrierSeq implements `bwa-mem` (Li, 2013) to first map all reads to the genomic carrier then
72 extracts unmapped reads by using `samtools` (Li et al., 2009) and `seqtk` (Li, 2012). Thereafter,
73 the user can define a quality score threshold and CarrierSeq proceeds to discard low-complexity
74 reads (Morgulis et al., 2006) with `fqtrim` (Pertea, 2015). This set of unmapped and filtered reads
75 are labeled “reads of interest” (ROI) and should theoretically comprise target reads and likely
76 contamination. However, ROIs also include “high-quality noise reads” (HQNRs), defined as reads
77 that satisfy quality score and complexity filters yet do not match to any database and dis-
78 proportionately originate from specific channels. By treating reads as a Poisson arrival process,
79 CarrierSeq models the expected ROIs channel distribution and rejects data from channels
80 exceeding a reads/channels threshold (x_{crit}) (Figure 1).

81



82

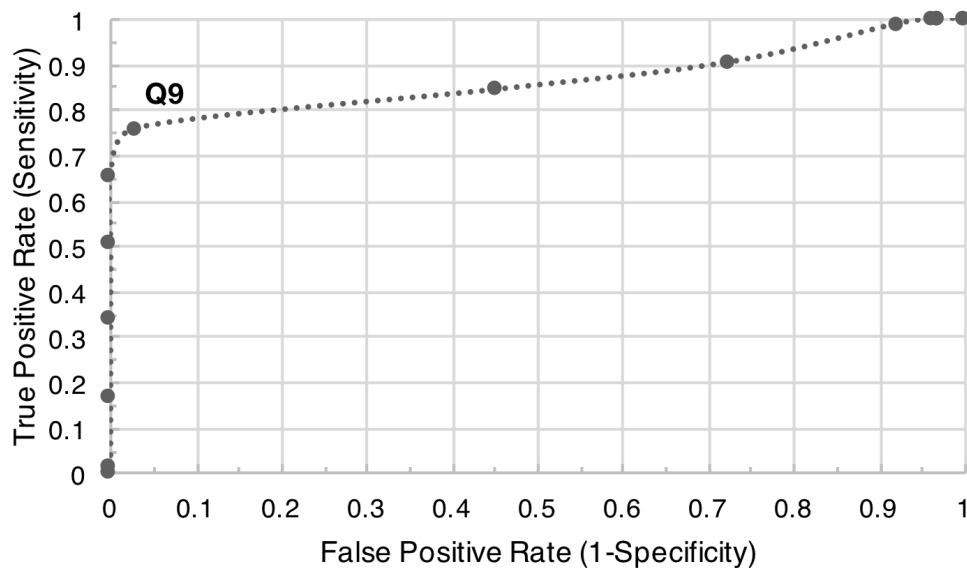
83

84 **Fig. 1. CarrierSeq workflow.** Starting from all reads, CarrierSeq identifies unmapped reads then
85 applies a quality score and complexity filter to discard low-quality reads. Afterwards, CarrierSeq
86 applies a Poisson distribution test to sort *likely* high-quality noise reads (HQNRs) from target
87 reads.

88 2.1 Quality Score Filter

89 The default per-read quality score threshold (Q9) was determined through receiver operating
90 characteristic curve (ROC) analysis (Fawcett, 2006) of carrier sequencing runs of *B. subtilis* and
91 Lambda DNA (Figure 2). This threshold is best suited for Lambda carriers that are 99% library by
92 mass and essentially function as a pseudo “lambda burn-in” experiment (Nanoporetech.com,
93 2017). Therefore, the user is encouraged to define their own threshold based on their libraries’
94 quality control metrics (e.g., carrier to target ratio, quality distribution, sequencing accuracy
95 achieved, and basecaller confidence).

96



97

98

99 **Fig. 2. Receiver operating characteristic curve.** Q9 provides a good threshold which discards
100 the majority of low-quality and noise reads (0.76 True Positive Rate and 0.03 False Positive Rate)
101 for carrier runs that are 99% Lambda DNA by mass. A perfect quality score threshold would plot
102 in the top left of the ROC curve.

103 **2.2 Poisson Distribution Sorting**

104 Assuming that sequencing is a stochastic process, CarrierSeq is able to identify channels producing
105 spurious reads by calculating the expected Poisson distribution of reads/channel. Given total ROIs
106 and number of active sequencing channels, CarrierSeq will determine the arrival rate ($\lambda = \text{reads of}$
107 $\text{interest/active channels}$). CarrierSeq then calculates an x_{crit} threshold ($x_{\text{crit}} = \text{poisson.ppf}(1 - p$
108 $\text{value}, \lambda)$) and sorts ROIs into target reads ($\text{reads/channel} \leq x_{\text{crit}}$) or HQNRs ($\text{reads/channel} > x_{\text{crit}}$)
109 (supplementary data).

110

111 **2.3 Implementation**

112 Reads to be analyzed must be compiled into a single fastq file and the carrier reference genome
113 must be in fasta format. Run CarrierSeq with:

114

```
115 ./carrierseq.sh -i <input.fastq> -r <reference.fasta> -o <output_directory>
```

116

117 **3 Results & Discussion**

118 From experimenting with low-input carrier sequencing and CarrierSeq we observed that the
119 abundance of HQNRs may vary per run, perhaps due to sub-optimal library preparation, delays in
120 initializing sequencing, or other sequencing conditions. In addition, target DNA purity and lysis
121 carryover (e.g., proteins) may conceivably contribute to HQNR abundance. Possibly due to pore
122 blockages from unknown macromolecules that result in erroneous reads. While the cause or
123 significance of HQNRs have yet to be determined, future work will focus on developing a method
124 to identify HQNRs on a per-read basis. In contrast, the current approach discards entire HQNR-
125 associated channels at the risk of discarding target reads. Moreover, some reads in non-HQNR-
126 associated channels may also be artifacts. The ability to identify HQNRs on a per-read basis is

127 especially important for metagenomic studies of novel microbial communities where HQNRs may
128 complicate the identification of an unknown organism, or in a life detection application (Carr et
129 al., 2017) where artefactual reads not mapping to known life could represent a false-positive.

130

131 **4 Summary**

132 CarrierSeq was developed to analyze low-input carrier sequencing data and identify target reads.
133 We have since deployed CarrierSeq to test the limits of detection of ONT's MinION sequencer
134 from 0.2 ng down to 2 pg of low-input carrier sequencing. CarrierSeq may be a particularly
135 valuable tool for in-situ metagenomic studies where limited sample availability (e.g., low biomass
136 environmental samples) and laboratory resources (i.e., field deployments) may benefit from
137 sequencing with a genomic carrier.

138

139 **Acknowledgements**

140 The authors would like to thank Michael Micorescu at Oxford Nanopore Technologies for
141 providing and granting us permission to utilize his fastq quality filter script.

142

143 **Funding**

144 This work has been supported by NASA MatISSE award NNX15AF85G

145

146 *Conflict of Interest:* none declared.

147

148 **References**

149 Barton, H.A. *et al.* (2006) DNA extraction from low-biomass carbonate rock: An improved method

150 with reduced contamination and the low-biomass contaminant database. *Journal of*
151 *Microbiological Methods*, **66**, 21–31.

152 Carr,C.E. *et al.* (2017) Towards in situ sequencing for life detection. 2017 *Aerospace Conference*,
153 1–18.

154 Fawcett,T. (2006) An introduction to ROC analysis. *Pattern Recognition Letters*, **27**, 861–874.

155 Lever,M.A. *et al.* (2015) A modular method for the extraction of DNA and RNA, and the
156 separation of DNA pools from diverse environmental sample types. *Front. Microbiol.*, **6**, 1–25.

157 Li,H. (2013) Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
158 *arXiv preprint arXiv1303.3997*, 1–3.

159 Li,H. (2012) seqtk Toolkit for processing sequences in FASTA/Q formats.

160 Li,H. *et al.* (2009) The Sequence Alignment/Map format and SAMtools. *Bioinformatics*, **25**, 2078–
161 2079.

162 Mojarro,A Ruvkun,G *et al.* (2017) Nucleic Acid Extraction from Synthetic Mars Analog Soils for
163 in situ Life Detection. *Astrobiology*.

164 Mojarro,A., Hachey,J., *et al.* (2017) Nucleic Acid Extraction and Sequencing from Low-Biomass
165 Synthetic Mars Analog Soils *Lunar & Planetary Science XLVIII*, 1-2.

166 Morgulis,A. *et al.* (2006) A fast and symmetric DUST implementation to mask low-complexity
167 DNA sequences. *Journal of Computational Biology*, **13**, 1028–1040.

168 Nanoporetech.com. (2017) *Getting started with MinION - what you need to know* Available at:
169 <https://nanoporetech.com/community/faqs>

170 Pertea,G. (2015) Fqtrim: v0. 9.4 release.

171 Poretsky,R. *et al.* (2014) Strengths and Limitations of 16S rRNA Gene Amplicon Sequencing in
172 Revealing Temporal Microbial Community Dynamics. *PLoS ONE*, **9**, e93827.

- 173 Sabina,J. and Leamon,J.H. (2015) Bias in Whole Genome Amplification: Causes and
174 Considerations. In, *Whole Genome Amplification*, Methods in Molecular Biology. Springer New
175 York, New York, NY, pp. 15–41.