

1 **BugSeq 16S: NanoCLUST with Improved Consensus Sequence Classification**

2 Ana Jung¹, Samuel D Chorlton^{1,2}

3

4 1. BugSeq Bioinformatics Inc, Vancouver, Canada

5 2. Department of Pathology & Laboratory Medicine, University of British Columbia, Vancouver,

6 Canada

7

8 Corresponding Author:

9 Samuel D Chorlton

10 Email: sam@bugseq.com

11 **Abstract**

12 NanoCLUST has enabled species-level taxonomic classification from noisy nanopore 16S
13 sequencing data for BugSeq's users and the broader nanopore sequencing community. We
14 noticed a high misclassification rate of NanoCLUST-derived consensus 16S sequences due to
15 its use of BLAST top hit taxonomy assignment. We replaced the consensus sequence classifier
16 of NanoCLUST with QIIME2's VSEARCH-based classifier to enable greater accuracy. We use
17 mock microbial community and clinical 16S sequencing data to show that this replacement
18 results in significantly improved nanopore 16S accuracy (over 5% recall and 19% precision),
19 and make this new tool (BugSeq 16S) freely available for academic use at [BugSeq.com/free](https://bugseq.com/free).

20 We read with great interest the recent publication of NanoCLUST, a species-level analysis of
21 16S rRNA nanopore sequencing data (Rodríguez-Pérez *et al.*, 2020). NanoCLUST brings a
22 novel, unsupervised read clustering approach to identify groups of similar reads, and corrects
23 the errors within each cluster to identify a highly accurate amplicon sequence variant (ASV).
24 This approach reduces the typically high (>5%) nanopore sequencing error rate, enabling
25 species-level identification from nanopore 16S sequencing data.

26

27 After its release, we rapidly integrated NanoCLUST into BugSeq, our democratized online
28 bioinformatics platform, to provide leading accuracy in 16S nanopore analysis for the academic
29 community, free of charge (Fan *et al.*, 2020). While BugSeq users have reaped significant
30 benefits from the development and online deployment of NanoCLUST, we have also noted
31 significant differences from its published accuracy on real world data. Here we report on a
32 frequent discrepancy between NanoCLUST and ground truth data, our solution, and a
33 benchmark of our new classification accuracy.

34

35 When designing NanoCLUST, Rodríguez-Pérez *et al.* used BLAST to classify ASVs (Camacho
36 *et al.*, 2009). Specifically, corrected consensus sequences from clusters are BLASTed against
37 the NCBI 16S BLAST database, and the sequence is assigned to the taxonomy of the top hit.
38 This approach leads to a high misclassification rate for several reasons. First, some bacteria
39 have nearly identical 16S sequences, such that only a 100% sequence identity between query
40 and reference sequence is appropriate to assign a sequence to the reference's species (Edgar,
41 2018). No such safeguards are taken with the NanoCLUST approach - an ASV may have 80%
42 sequence identity to the top BLAST hit yet be assigned to that species. Second, NanoCLUST
43 accepts a BLAST hit of any length, as long as the expectation value of the alignment is less
44 than 11; that is, for any ASV classification, there may be up to 11 hits found just by chance.

45

46 We sought to replace the ASV classification mechanism of NanoCLUST with a more accurate
47 16S classifier. We tried IDTAXA and a naive Bayes classifier, but settled on the QIIME2
48 VSEARCH consensus classifier as the former classifiers did not provide species-level
49 classification or could not be tuned to NanoCLUST data qualities (Murali *et al.*, 2018; Bokulich
50 *et al.*, 2018). The VSEARCH classifier sets default identity (80%) and alignment length (80%)
51 thresholds, ensuring that only high quality alignments are retained. Additionally, it collapses the
52 top search hits by lowest common ancestor to ensure precision when closely related sequences
53 exist in the reference database. We deviate in one other way from NanoCLUST: the minimum
54 cluster size for HDBSCAN was recently decreased from 200 to 50 in NanoCLUST (as of
55 February 15, 2021); BugSeq has yet to follow suit (as of March 5, 2021).

56

57 We evaluated our novel 16S analysis pipeline, BugSeq 16S, with publicly available full-length
58 16S sequencing data of known mock microbial communities. Cusco *et al.* sequenced the
59 ZymoBIOMICS mock microbial community, containing eight bacteria in even abundance, in
60 duplicate (Cuscó *et al.*, 2019). Using NanoCLUST, there were 6 and 5 species correctly
61 identified in the samples, with 4 and 2 false positive species, respectively (average precision
62 66%, average recall 69%). Using BugSeq 16S, we identify 7 bacterial species correctly in each
63 sample, with 1 false positive in both (average precision 88%, average recall 88%). Specifically,
64 NanoCLUST classified the *Listeria monocytogenes* sequences as *Listeria innocua*, the *Bacillus*
65 *subtilis* sequences as *Bacillus halotolerans* and *Bacillus tequilensis*, and the *Escherichia coli*
66 sequences as *Shigella flexneri*. BugSeq 16S identified the *Bacillus subtilis* sequences to the
67 correct genus (without classifying to species level), and falsely identified an uncultured
68 bacterium from the *Escherichia/Shigella* genus, in addition to the present *E. coli*. Revising the
69 NanoCLUST minimum cluster size to 200 reads did not affect conclusions: there was still
70 misidentification of the *B. subtilis* and *L. monocytogenes* reads, with *E. coli* no longer detected.

71

72 On a more complex mock microbial community containing 20 bacteria spanning a 1000-fold
73 concentration range (BEI HM-783D), NanoCLUST detected 5 species correctly with 8 false
74 positive species; BugSeq 16S detected 6 species correctly with 1 false positive species (Cuscó
75 *et al.*, 2019). Recall was therefore 30% for BugSeq as compared to 25% for NanoCLUST, with
76 precision of 85.7% for BugSeq and 38% for NanoCLUST. Full results are available in Table 1.
77

78 We next compared BugSeq 16S with NanoCLUST on publicly available full-length 16S
79 sequencing data from patients with ventilator associated pneumonia (Maes *et al.*, 2021).
80 Twenty-nine bronchoalveolar lavage samples from 24 patients underwent microbial culture,
81 multi-pathogen TaqMan array and nanopore 16S sequencing as previously described (Maes *et*
82 *al.*, 2021). We used a combination of culture and TaqMan array results as the gold standard. As
83 TaqMan and culture results were censored by the original authors at a cycle threshold value of
84 32 and 10^4 colony forming units per milliliter, respectively, only sensitivity of 16S sequencing
85 could be calculated. Full data is available in the supplementary material.

86
87 At the species level, BugSeq 16S achieved better sensitivity ($n=28/35$, 80%) as compared to
88 NanoCLUST ($n=26/35$, 74%). Specifically, patient 16 had an *E. coli* detected by culture and
89 BugSeq 16S but NanoCLUST detected a mix of *Escherichia fergusonii*, *Shigella flexneri* and
90 *Shigella sonnei*. Patient 17 has a *Serratia marcescens* detected by culture, TaqMan and
91 BugSeq 16S, while NanoCLUST detected a *Serratia nematodiphila* and *Serratia*
92 *surfactantfaciens*. BugSeq 16S achieved this greater species-level sensitivity while predicting
93 fewer total species present in each sample. BugSeq 16S predicted a median of 13 (IQR=8)
94 fewer species present in each sample as compared with NanoCLUST.

95
96 As BugSeq 16S collapsed results based on sequence similarity, it also correctly detected a
97 *Stenotrophomonas spp.* in patient 1 (sample 2) which NanoCLUST mislabelled as

98 *Stenotrophomonas pavanii* (ground truth: *Stenotrophomonas maltophilia*) and a *Haemophilus*
99 *spp.* in patient 21 which NanoCLUST mislabelled as *Haemophilus parahaemolyticus* (ground
100 truth: *Haemophils influenzae*). Accounting for these confidence-aware classifications increases
101 BugSeq 16S's sensitivity to 85.7%.

102

103 In conclusion, BugSeq 16S builds on NanoCLUST to improve nanopore 16S sequencing
104 analysis using a popular, highly accurate sequence classifier. This change results in significantly
105 improved analysis accuracy on both mock microbial communities and real patient samples, and
106 is free for academic use at bugseq.com/free.

- Bokulich, N.A. *et al.* (2018) Optimizing taxonomic classification of marker-gene amplicon sequences with QIIME 2's q2-feature-classifier plugin. *Microbiome*, **6**, 90.
- Camacho, C. *et al.* (2009) BLAST+: architecture and applications. *BMC Bioinformatics*, **10**, 421.
- Cuscó, A. *et al.* (2019) Microbiota profiling with long amplicons using Nanopore sequencing: full-length 16S rRNA gene and the 16S-ITS-23S of the *rrn* operon. *F1000Res*, **7**, 1755.
- Edgar, R.C. (2018) Updating the 97% identity threshold for 16S ribosomal RNA OTUs. *Bioinformatics*, **34**, 2371–2375.
- Fan, J. *et al.* (2020) BugSeq: a highly accurate cloud platform for long-read metagenomic analyses. *bioRxiv*, 2020.10.08.329920.
- Maes, M. *et al.* (2021) Ventilator-associated pneumonia in critically ill patients with COVID-19. *Critical Care*, **25**, 25.
- Murali, A. *et al.* (2018) IDTAXA: a novel approach for accurate taxonomic classification of microbiome sequences. *Microbiome*, **6**, 140.
- Rodríguez-Pérez, H. *et al.* (2020) NanoCLUST: a species-level analysis of 16S rRNA nanopore sequencing data. *Bioinformatics*.

Species	Theoretical Operon Copies ($\times 10^3$) in the Community	BugSeq 16S Abundance	NanoCLUST Abundance
<i>Streptococcus mutans</i>	1000	42.34%	51.76%
<i>Escherichia coli</i>	1000	19.41%	0.00%
<i>Staphylococcus epidermidis</i>	1000	1.88%	22.78%
<i>Rhodobacter sphaeroides</i>	1000	1.10%	0.37%
<i>Streptococcus agalactiae</i>	100	1.59%	2.51%
<i>Pseudomonas aeruginosa</i>	100	0.78%	0.57%
<i>Bacillus cereus</i>	100	0.00%	0.00%
<i>Clostridium beijerinckii</i>	100	0.00%	0.00%
<i>Staphylococcus aureus</i>	100	0.00%	0.00%
<i>Acinetobacter baumannii</i>	10	0.00%	0.00%
<i>Helicobacter pylori</i>	10	0.00%	0.00%
<i>Lactobacillus gasseri</i>	10	0.00%	0.00%
<i>Listeria monocytogenes</i>	10	0.00%	0.00%
<i>Neisseria meningitidis</i>	10	0.00%	0.00%
<i>Propionibacterium acnes</i>	10	0.00%	0.00%
<i>Actinomyces odontolyticus</i>	1	0.00%	0.00%
<i>Bacteroides vulgatus</i>	1	0.00%	0.00%
<i>Deinococcus radiodurans</i>	1	0.00%	0.00%
<i>Enterococcus faecalis</i>	1	0.00%	0.00%
<i>Streptococcus pneumoniae</i>	1	0.00%	0.00%
Uncultured Escherichia/Shigella	0	6.82%	0.00%
<i>Escherichia fergusonii</i>	0	0.00%	10.11%
<i>Shigella sonnei</i>	0	0.00%	4.01%
<i>Citrobacter rodentium</i>	0	0.00%	2.25%
<i>Staphylococcus caprae</i>	0	0.00%	1.73%
<i>Bacillus tropicus</i>	0	0.00%	1.68%
<i>Shigella flexneri</i>	0	0.00%	1.08%
<i>Bacillus badius</i>	0	0.00%	0.58%
<i>Shigella boydii</i>	0	0.00%	0.57%

107 **Table 1:** Species-level abundance in the HM-783D mock microbial community using BugSeq
108 16S and NanoCLUST. Green rows reflect species actually present in the sample, red rows
109 reflect false positive species.