

RESEARCH

Open Access



# Species-level resolution of 16S rRNA gene amplicons sequenced through the MinION™ portable nanopore sequencer

Alfonso Benítez-Páez<sup>\*</sup>, Kevin J. Portune and Yolanda Sanz

## Abstract

**Background:** The miniaturised and portable DNA sequencer MinION™ has been released to the scientific community within the framework of an early access programme to evaluate its application for a wide variety of genetic approaches. This technology has demonstrated great potential, especially in genome-wide analyses. In this study, we tested the ability of the MinION™ system to perform amplicon sequencing in order to design new approaches to study microbial diversity using nearly full-length 16S rDNA sequences.

**Results:** Using R7.3 chemistry, we generated more than 3.8 million events (nt) during a single sequencing run. These data were sufficient to reconstruct more than 90 % of the 16S rRNA gene sequences for 20 different species present in a mock reference community. After read mapping and 16S rRNA gene assembly, consensus sequences and 2d reads were recovered to assign taxonomic classification down to the species level. Additionally, we were able to measure the relative abundance of all the species present in a mock community and detected a biased species distribution originating from the PCR reaction using ‘universal’ primers.

**Conclusions:** Although nanopore-based sequencing produces reads with lower per-base accuracy compared with other platforms, the MinION™ DNA sequencer is valuable for both high taxonomic resolution and microbial diversity analysis. Improvements in nanopore chemistry, such as minimising base-calling errors and the nucleotide bias reported here for 16S amplicon sequencing, will further deliver more reliable information that is useful for the specific detection of microbial species and strains in complex ecosystems.

**Keywords:** MinION, Nanopore sequencer, 16S rDNA amplicon sequencing, Microbial diversity, Long-read sequencing

## Background

The third generation of DNA sequencers is based on single-molecule analysis technology that constantly is under development to minimise errors and produce high quality reads. Oxford Nanopore Technologies (ONT) released the first miniaturised and portable DNA sequencer to researchers in early 2014, within the framework of the MinION™ Access Programme. The MinION™ is a USB stick-sized device operated from a computer via USB 3.0. Real-time data analysis can be visualised in terms of number of reads and length distribution. Nucleotide base-calling and quality assessment of reads

require further processing, where data exchange of Hierarchical Data Format (HDF5) files, containing a large amount of numerical data, is indispensable. This data exchange is done via the Internet through the Metrichor platform; a process that can optionally be launched after the sequencing process itself. According to its theoretical capabilities, the MinION™ provides new alternatives for genomic analyses. One of the most attractive capabilities of the MinION™ platform is the sequencing of complete bacterial genomes, as demonstrated recently by Quick *et al.* [1]. Another major advantage of the MinION™ platform, compared to other popular sequencing technologies, is its performance in terms of read length. Theoretically, nanopore-sensing technology is able to generate thousands of reads that are hundreds to thousands of nucleotides in length; the only limitation being the DNA fragments generated during nucleic acid extraction

\* Correspondence: [abenitez@iata.csic.es](mailto:abenitez@iata.csic.es)

Microbial Ecology, Nutrition & Health Research Unit, Institute of Agrochemistry and Food Technology Institute, Spanish National Research Council (IATA-CSIC), C. Catedrático Agustín Escardino Benlloch, 7, 46980 Paterna, Valencia, Spain

procedures, which frequently produce fragmented DNA with an average length of 50 kb. Although short-read length sequencing approaches deliver high quality sequences, these partial genome sequences with unsolved repetitive elements make it impossible to study genetic variation or molecular evolution directly or indirectly associated to such elements. Therefore, long-read approaches offer new insights into genomic analysis, facilitating the assembly of complete genomes through hybrid strategies [2]. In addition to genome sequencing analysis, microbial diversity and taxonomic approaches are also deeply limited by short-read strategies. Early massive sequencing approaches producing 50 nt (Genome Analyzer, Solexa/Illumina) to 200 nt (454 Roche) effective reads with a modest average quality only allowed accurate exploration of diversity at the phylum level. However, thanks to improvements in the chemistry of the most common, popular sequencing platforms in recent years, it is now possible to characterise microbial communities in detail down to the family or even genus level. To date, paired-end short read approaches for massive sequencing permit the analysis of sequence information of roughly 30 % (~500 nt) of the full 16S rRNA gene, which means taxonomic assignment of reads at the species level is elusive. Therefore, implementation of long-read sequencing approaches to study 16S rRNA genes will permit the design of new studies to provide evidence for the central role of precise bacterial species/strains in a great variety of microbial consortia. As a consequence, we present a preliminary study of 16S rDNA amplicon sequencing of a mock microbial community composed of genomic DNA from 20 different bacterial species (BEI Resources) using the MinION™ sequencing platform. The aim of this study is to evaluate the application of nanopore technology in performing bacterial diversity and taxonomic analysis on nearly full-length bacterial 16S rRNA genes.

### Data description

Raw data collected in this experiment were obtained as fast5 files using MinKNOW software v0.50.1.15 (Oxford Nanopore Technologies), after conversion of electric signals into base calls via the Metrichor Agent v2.29 and the 2D Basecalling workflow v1.16. Base-called data passing quality control and filtering were downloaded and basic statistical analysis was carried out using *poretools* [3] and *poRe* [4]. Mapping statistics are depicted in Table 1. Fast5 raw data can be accessed at the European Nucleotide Archive (ENA) under the project ID PRJEB8730 (sample ERS760633). Only one data set was generated after a sequencing run of MinION™.

### Analyses

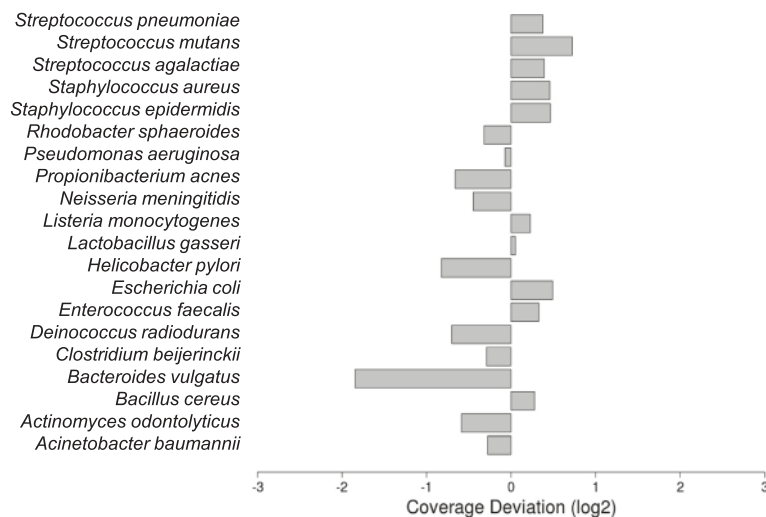
DNA reads derived from MinION sequencing can be classified into three types: ‘template’, ‘complement’, and

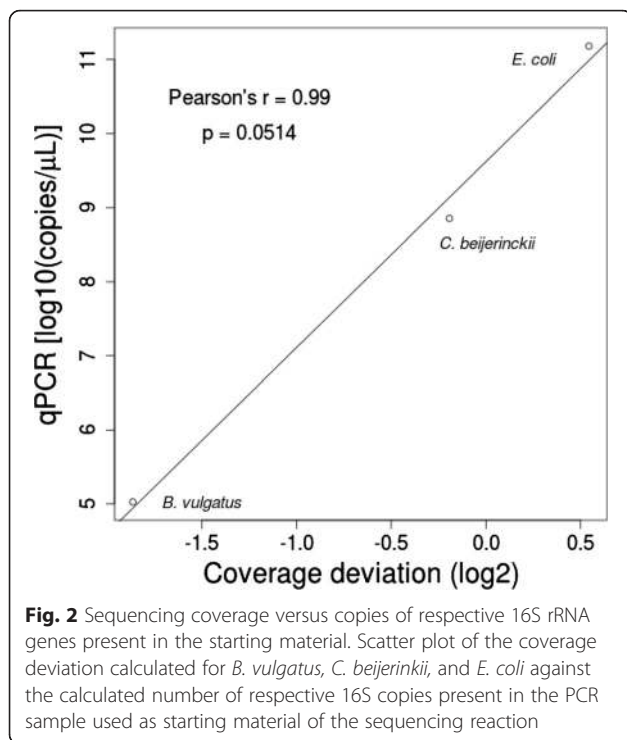
‘2d’ reads. While template reads come from DNA strands that are primed by a leader adapter and passed through the pore, the complement reads are generated only if a second adapter (hairpin adapter) is present in the same DNA fragment, thus permitting sequencing both strands of a single molecule in a concatenated manner. The 2d reads are products of aligning and merging sequences from template and complement reads generated from the same DNA fragment: these contain a lower error rate, owing to strand comparison and mismatch correction. After the sequencing process, we obtained 3404 reads, of which 58.5 % were template reads (1991), 23.8 % were complement reads (812), and 17.7 % were 2d reads (601). Read lengths had a wide distribution ranging from 12 nt to more than 50,000 nt in length, with a median of 1100 nt. We hypothesised that extremely large reads might be products of the ligation of multiple amplicons. However, when we tried to align these large reads to reference sequences, we detected no matches (data not shown). Accordingly, a filtering step was performed by retaining 97 % of the original dataset (3297 reads), with a size range between 100 and 2000 nt in length for downstream analysis.

In the first step of our analysis, we used the large set of template and complement reads to assess the global performance of the amplicon sequencing process. Consequently, we analysed basic read-mapping statistics to uncover potential pitfalls of the MinION platform and tried to reconstruct the reference sequences. We assembled more than 90 % of 16S rRNA gene sequences for all organisms included in the mock community (Table 1). We observed that even at very low coverage, such as that retrieved for the *Bacteroides vulgatus* 16S rRNA gene (Fig. 1), it is possible to reconstruct almost 93 % of the entire gene. Indeed, the maximum size of amplicons sequenced in all cases was close to the expected amplicon size according to the universal primers used in the PCR design (Table 1). In terms of coverage, we hypothesised that the lower than expected number of 16S reads from *B. vulgatus* species (Fig. 1) resulted from a bias caused during PCR amplification, despite using high coverage primers [5, 6], or as a result of the sequencing process itself. To further investigate this matter, we performed an absolute quantification of 16S rRNA genes using qPCR from three different species with a high coverage, close-to-expected coverage, and the lowest coverage, respectively, which were present in the initial PCR sample used for library construction and sequencing. A correlation between the number of molecules present in the starting material and the coverage obtained after the sequencing process (Pearson’s  $r = 0.99$ ,  $p \leq 0.0514$ ) was detected (Fig. 2), indicating that the sequencing process faithfully reproduced the proportion of amplicons present in the sample and the coverage bias

**Table 1** Statistics of the mapping process using 16S rDNA reads produced by MinION™

Organism	Mapped reads	Mapped bases	Strand mapping <sup>a</sup>	Max length <sup>b</sup>	Mean length	Variants after assembly	Consensus	rRNAgene <sup>c</sup>	Assembled 16S <sup>d</sup>
<i>Acinetobacter baumannii</i>	98	99,352	0.46:0.54	1390	1013	32	1415	1529	0.93
<i>Actinomyces odontolyticus</i>	79	73,480	0.49:0.51	1377	930	8	1407	1528	0.92
<i>Bacillus cereus</i>	144	151,668	0.46:0.54	1419	1053	31	1415	1508	0.94
<i>Bacteroides vulgatus</i>	33	29,499	0.58:0.42	1346	893	25	1403	1510	0.93
<i>Clostridium beijerinckii</i>	97	99,476	0.46:0.54	1393	1025	13	1408	1505	0.94
<i>Deinococcus radiodurans</i>	73	69,940	0.45:0.55	1390	958	8	1398	1502	0.93
<i>Enterococcus faecalis</i>	149	153,581	0.50:0.50	1398	1030	8	1444	1549	0.93
<i>Escherichia coli</i>	167	181,084	0.45:0.55	1398	1084	0	1434	1542	0.93
<i>Helicobacter pylori</i>	67	62,838	0.46:0.54	1390	937	11	1411	1498	0.94
<i>Lactobacillus gasseri</i>	123	128,120	0.51:0.49	1407	1041	0	1467	1579	0.93
<i>Listeria monocytogenes</i>	139	140,478	0.50:0.50	1343	1010	13	1374	1486	0.92
<i>Neisseria meningitidis</i>	87	86,916	0.48:0.52	1390	999	11	1433	1544	0.93
<i>Propionibacterium acnes</i>	75	70,160	0.48:0.52	1375	935	21	1401	1525	0.92
<i>Pseudomonas aeruginosa</i>	113	120,520	0.55:0.45	1398	1066	14	1425	1536	0.93
<i>Rhodobacter sphaeroides</i>	95	89,750	0.52:0.48	1416	944	5	1352	1463	0.92
<i>Staphylococcus epidermidis</i>	164	177,084	0.51:0.49	1423	1079	0	1443	1540	0.94
<i>Staphylococcus aureus</i>	163	179,477	0.51:0.49	1423	1101	1	1435	1554	0.92
<i>Streptococcus agalactiae</i>	156	166,420	0.52:0.48	1411	1066	5	1439	1551	0.93
<i>Streptococcus mutans</i>	196	221,682	0.47:0.53	1411	1131	2	1440	1552	0.93
<i>Streptococcus pneumoniae</i>	154	168,657	0.52:0.48	1411	1095	2	1442	1560	0.92

<sup>a</sup>Proportion of reads mapped against the forward and complementary strand, respectively<sup>b</sup>Maximum length of reads mapped<sup>c</sup>Length of the 16S sequence used as reference<sup>d</sup>Numbers are generated from Consensus/rRNA gene ratio**Fig. 1** Species abundance in the mock community detected by MinION™. Species coverage was calculated by obtaining the fold-change ( $\log_2$ ) of species-specific read counting against the expected average for the entire community. A coverage bias was assumed when coverage deviation was lower than  $-1$  or higher than  $1$



**Fig. 2** Sequencing coverage versus copies of respective 16S rRNA genes present in the starting material. Scatter plot of the coverage deviation calculated for *B. vulgatus*, *C. beijerinckii*, and *E. coli* against the calculated number of respective 16S copies present in the PCR sample used as starting material of the sequencing reaction

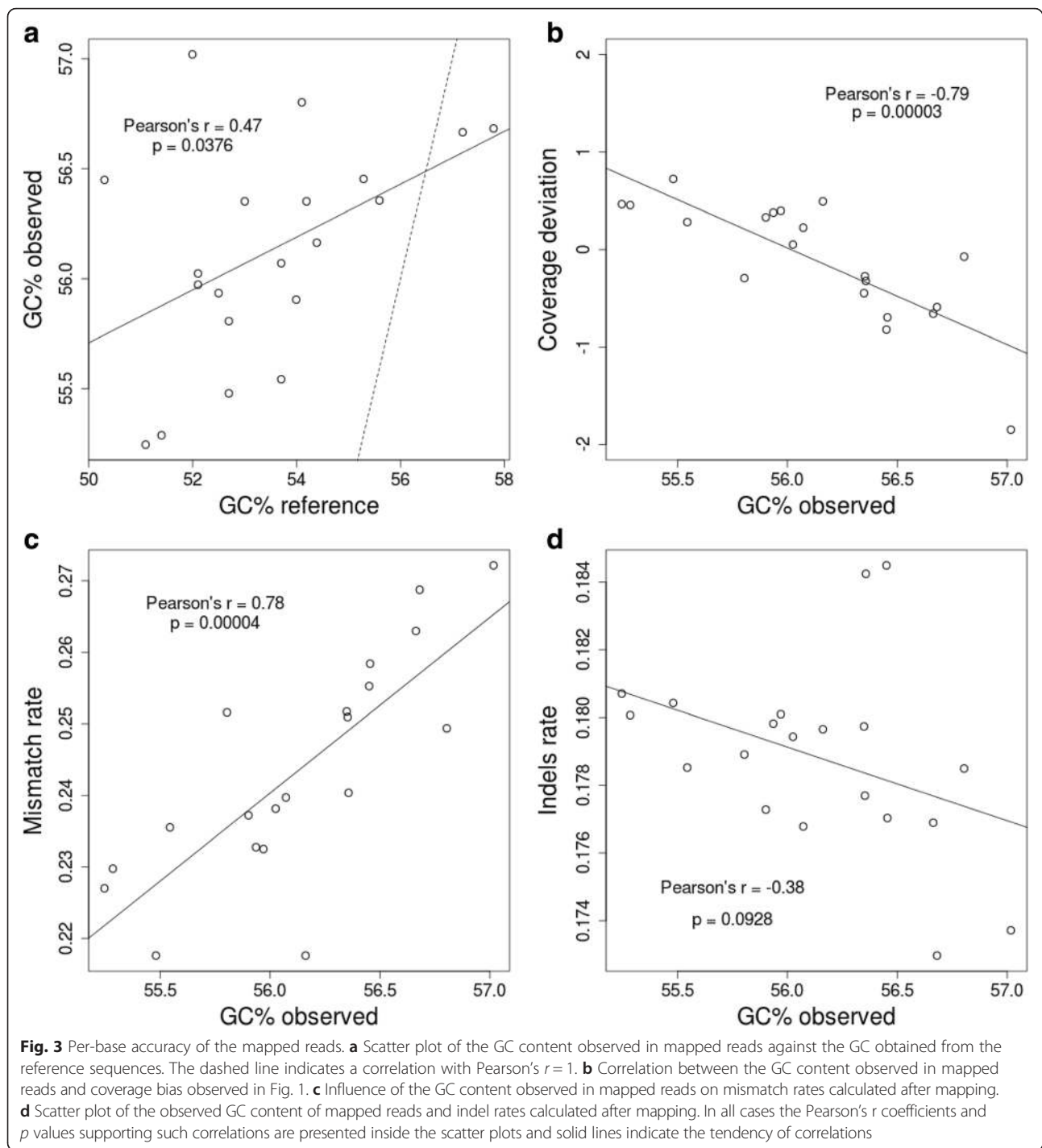
was therefore derived from the starting material generated by PCR. Despite this bias, the 16S rRNA gene from *B. vulgatus* was fully assembled with low variation ( $25/1403 = 1.78\%$ ) after DNA read alignment and pileup (Table 1).

Read-mapping statistics were analysed to further measure the performance of MinION™ sequencing in microbial diversity analysis based on 16S rDNA sequences. The GC content of reads produced by MinION™ showed an important and significant correlation (Pearson's  $r = 0.47$ ,  $p \leq 0.0376$ ) against the GC content of reference values (Fig. 3a), which indicates that the GC content of 16S rDNA sequences is fairly well replicated during sequencing. However, we found a 16S rDNA GC content bias, to some extent, in the reads obtained from MinION™, which in almost all cases exceeds the GC content of the reference (Fig. 3a). To test the probable influence of GC content bias in base-calling accuracy, we performed linear comparisons against mismatch rates, indel rates, and coverage deviation. We observed that both coverage deviation ( $p \leq 0.00003$ ) and mismatch rate ( $p \leq 0.00004$ ) are significantly influenced by read GC content (Fig. 3b and c, respectively). In the first case, the influence of GC content on coverage deviation could have a minimal effect because 95 % of species analysed show no more than a one-fold deviation. However, with GC bias detected in reads from the MinION™ sequencer, this effect could be magnified, especially in species where GC content is high. On the other hand, we found a strong correlation between the GC content of reads and the mismatch rate retrieved from

alignments, which would insinuate again that GC content is a factor that influences 16S rDNA amplicon sequencing in the MinION™ platform. Conversely, GC content did not appear to profoundly affect indel rate (Fig. 3d).

The complete assembly of the amplified 16S rRNA gene permitted the quantification of the level of sequence variants in the consensus sequence. These variants were recovered after a pileup of reads against reference sequences, and they were variable in number with a median of 8 variants per 16S rRNA gene (Table 1). This number of nucleotide substitutions means that approximately 0.5 % of the 16S rDNA sequence assembled from MinION™ reads retained unnatural genetic variants directly generated from the sequencing process itself, theoretically allowing a bona fide identification and taxonomic assignment of 16S rDNA sequences at the species level. In the worst cases, where the number of variants were meaningful ( $\sim 2.3\%$  of the full assembly), such as those observed for *Acinetobacter baumannii* and *Bacillus cereus* (Table 1), direct BLAST comparisons of these assembled 16S rDNA sequences against the NCBI 16S database only produced matches with homologous sequences belonging to the same species, respectively (data not shown).

A final step in our analysis tested whether or not the information obtained through sequencing of nearly full-size 16S rRNA genes using the MinION™ platform is useful to perform taxonomic assignment with tools commonly employed in microbial community analysis. For this aim, we used the information derived from 2d reads (601 reads), which is limited in terms of the effective number of reads but more reliable in terms of sequence identity. Using the SINA web service [7] we obtained the taxonomic assignment of 2d reads to the Silva bacterial 16S database [8]. The results of this approach are shown in Table 2. Out of the 17 different genera present in the mock community, we retrieved information for six of them, with an assignment threshold of 80 %, seven with 70 % and eight with 60 %. Using 60 % as the lowest assignment threshold, we started to retrieve unexpected genera composition in our 2d data set, indicating that reliable identifications must be set with a higher identity threshold. As expected, taxonomy assignment was limited to those species with a higher coverage during sequencing processing (Fig. 1), which is consistent with the number of 2d reads expected after aligning and merging respective template and complement reads obtained from the 16S rRNA genes of species over-represented in the starting material. We expect that the whole repertoire of species present in the sample can be detected by increasing the performance of the sequencing process. This would allow us to obtain a larger raw dataset and, particularly, more 2d reads containing more reliable information to perform taxonomic assignments and



disclose the full inventory of species present in the microbial community under study. Finally, a BLAST-based assignment against the NCBI bacterial 16S rRNA gene database retrieved the identities of 8 of the 20 species presented in the mock community analysed (Table 2). Although other species were also retrieved, they exhibited a high level of affiliation in terms of the 16S rDNA

sequence identity to the true species included in the mock community.

### Discussion

The inventory of microbial species based on 16S rDNA sequencing is frequently used in biomedical research to determine microbial organisms inhabiting the human

**Table 2** Taxonomy assignment of 2d reads derived from 16S rDNA sequencing using MinION™

Genus assignment	At 80 % identity (SINA)	At 70 % identity (SINA)	At 60 % identity (SINA)	Species assignment (Blast <sup>3</sup> )
<i>Bacillus</i>	3	3	4	<b><u>B. cereus</u></b>
<i>Escherichia</i>	61	93	103	<b><u>E. coli</u></b> , <i>E. fergusonii</i> , <i>Shigella flexneri</i>
<i>Streptococcus</i>	69	96	111	<b><u>S. mutans</u></b> , <i>S. gordonii</i> , <b><u>S. agalactiae</u></b>
<i>Lactobacillus</i>	2	2	3	<b><u>L. gasseri</u></b> , <i>L. fornicalis</i>
<i>Staphylococcus</i>	37	55	60	<b><u>S. epidermidis</u></b> , <i>S. capitis</i> , <b><u>S. aureus</u></b>
<i>Rhodobacter</i>	4	8	10	<b><u>R. sphaeroides</u></b>
<i>Clostridium</i>	0	2	2	N/A
<i>Listeria</i>	0	0	1	N/A
<i>Cronobacter</i> *	0	1	1	N/A
<i>Yersinia</i> *	0	0	1	N/A
<i>Citrobacter</i> *	0	0	1	<i>C. koseri</i> *, <i>C. muytjensii</i> *
<i>Alteromonas</i> *	0	0	1	N/A
<i>Burkholderia</i> *	0	0	1	N/A
<i>Sulfitobacter</i> *	0	0	1	N/A
Total assigned	176/601	260/601	300/601	172/601

<sup>3</sup>Blast searching against the NCBI bacterial 16S database. Valid assignments were based on E-values < 1e<sup>-03</sup> and alignments larger than 800 bp. N/A means no species was found for the respective genus. Species matching with those present in the mock community are underlined and in bold. Rows with genera, species, and read counts highlighted with asterisk (\*) indicate they are not expected to be present in the mock community

body and their relationship with disease. Identification of microbial species inhabiting different areas and cavities of the human body currently relies on the handling and processing of millions of DNA sequences obtained through the second generation of massive and parallel sequencing methods. However, these methods are still limited, mainly in terms of DNA read length. The inability to determine complete 16S rDNA sequences during massive sequencing has led to the development of multiple algorithms dedicated to theoretically discerning microbial species present in samples according to the sequence similarity degree, or Operational Taxonomic Units (OTUs). Despite high accuracy and a constant update of the methods used in OTU-based approaches, available algorithms produce no consensus outputs, leaving a high degree of uncertainty when the number of theoretical species and their abundance is the subject of study [9–12].

Thanks to the fact that they overcome DNA read limitations at the expense of decreasing throughput, a third generation of sequencing methods based on single-molecule technology offers new possibilities to study microbial diversity and taxonomic composition. MinION™ is one of these single-molecule methodologies, which has demonstrated its capacity in genome sequencing [1, 13]. Recent studies have reported the application of this technology in medical microbiology by using amplicon sequencing to determine bacterial and viral infections [14, 15]. Our results indicate that the MinION™ per-base accuracy (65–70 % for template reads, and 85 % for 2d reads) is in concordance with previous results [1, 14, 16].

We found that sequence coverage was close to expected values in most cases, with the exception of that of *B. vulgatus* (gene GC = 52 %), which was 1.84-fold less than the expected coverage. Using absolute quantification of molecules presented in the starting material, we demonstrated that such coverage bias came from the PCR process used to generate the 16S amplicons, despite using ‘universal’ primers with higher coverage among bacterial species [5]. Despite this, such coverage was enough to reconstruct 93 % of the 16S rRNA gene of *B. vulgatus* with a low proportion of unnatural variants.

We observed a general influence of GC content in the mismatch rate but not in the indel rate. This suggests that base miscalling could be associated with the amplicon GC content. Moreover, a slight correlation between the amplicon GC content observed and coverage bias was evidenced, indicating that GC content could be negatively affecting amplicon coverage to some extent. Although MinION™ was able to replicate the GC content expected for every amplicon sequenced fairly well, we observed a slight over-representation of GC in all reads obtained. This over-calling of GC bases in 16S rDNA amplicons could additionally influence the issues stated above in a negative manner.

The R7.3 chemistry used in MinION™ allowed the acquisition of reads of moderate quality, which were enough to reconstruct more than 90 % of the 16S rRNA gene in all 20 bacterial species analysed. None of the 20 16S rDNA consensus sequences assembled showed more than 3 % of sequence variation, which can be considered as a threshold for canonical species identification.

Therefore, the consensus sequence assembled was useful to obtain a reliable taxonomic identification at the species level. As expected, unnatural variants were associated with low coverage regions. Therefore, increasing the sequencing coverage will drastically reduce the ambiguities of the assembled sequences. When we tested the high quality reads (2d) in common routines for the analysis of microbial communities, the SINA web server retrieved a taxonomic assignment, indicating the presence of 7 genera out of the 17 expected for the mock community without any mismatches (using 70 % sequence identity as a threshold). Although this number of matches can be considered low, it was directly associated with the sequencing coverage, therefore, a larger 2d data set generated from a greater sequencing effort would produce enough information to identify the entire community.

### Conclusions

In terms of the study of microbial communities, results obtained using 16S rDNA amplicon sequencing through the MinION™ device are promising. Despite the observed modest per-base accuracy of this sequencing platform, we were able to reconstruct nearly full-length 16S rDNA sequences for 20 different species analysed from a mock bacterial community, and were able to obtain an acceptable taxonomy assignment for high quality 2d reads, only limited by the sequencing effort. This seems to be the major handicap of the MinION™ platform for microbial diversity analysis. To date, MinION™ and nanopore technologies have demonstrated great potential in DNA sequencing by allowing the retrieval of whole bacterial genome sequences with a minimum level of variation [1]. With the results presented here, we postulate that the MinION™ platform is a reliable methodology to study the diversity of microbial communities. It permits: i) a taxonomic identification at the species level through 16S rDNA sequence comparisons, and ii) a relative quantification to determine the species abundance. This type of analysis will likely become more accurate over time as nanopore chemistry is improved in future releases, with the concomitant increasing of the throughput, pivotal to disclose the hundreds of species present in complex microbial communities. The implementation of the “What’s In My Pot” (WIMP) Metrichor workflow, which aims to acquire real-time taxonomic sequence identification by comparing against different bacterial references databases (i.e., NCBI, SILVA [8], GreenGenes [17]), will be helpful in other types of analyses related to those presented here. Accordingly, sequence studies of the entire 16S rDNA molecule could allow OTU-based analysis to be bypassed completely, thus making it feasible to obtain a direct inventory of bacterial species and relative abundance, as well as to determine the key players at the species level in different microbial communities of interest.

### Methods

#### Bacterial DNA and 16S rDNA amplicons

Genomic DNA for the reference mock microbial community was kindly donated by BEI Resources (<http://www.beiresources.org>). This mock community (HM-782D) is composed of a genomic DNA mix from 20 bacterial strains containing equimolar ribosomal RNA operon counts (100,000 copies per organism per  $\mu\text{L}$ ), as indicated by the manufacturer. According to instructions provided by BEI Resources, 1  $\mu\text{L}$  of mock community DNA was used to amplify 16S rRNA genes. DNA was amplified by 30 PCR cycles at 95 °C for 20 s, 47 °C for 30 s, and 72 °C for 60 s. Phusion High-Fidelity Taq Polymerase (Thermo Scientific) and the primers S-D-Bact-0008-c-S-20 and S-D-Bact-1391-a-A-17, which target a wide range of bacterial 16S rRNA genes, were used during PCR [5, 6]. Amplicons consisted of ~1.5 kbp blunt-end fragments, which were purified using the Illustra GFX PCR DNA and Gel Band Purification Kit (GE Healthcare). Amplicon DNA was quantified using a Qubit 3.0 fluorometer (Life Technologies).

#### Amplicon DNA library preparation

The Genomic DNA Sequencing Kit SQK-MAP-005 was used to prepare the amplicon library to be loaded into the MinION™. Approximately 250 ng of amplicon DNA (0.25 pmol) was processed for end repair using the NEBNext End Repair Module (New England Biolabs), followed by purification using Agencourt AMPure XP beads (Beckman Coulter). Subsequently, according to the manufacturer’s instructions, we used 200 ng of the purified amplicon DNA (~0.2 pmol) to perform dA-tailing using the NEBNext dA-tailing module (New England Biolabs) with a total volume of 30  $\mu\text{L}$ , and incubated the sample at 37 °C for 15 min. Fifty  $\mu\text{L}$  of Blunt/TA ligase master mix (New England Biolabs), 10  $\mu\text{L}$  of adapter mix, and 2  $\mu\text{L}$  of HP adapter were added to the 30  $\mu\text{L}$  dA-tailed amplicon DNA. The reaction was incubated at 16 °C for 15 min. The adapter-ligated amplicon was recovered using Dynabeads® His-Tag (Life Technologies) and washed with washing buffer provided with the Genomic DNA Sequencing Kit SQK-MAP-005 (Oxford Nanopore Technologies). Finally, the sample was eluted from the Dynabeads® by adding 25  $\mu\text{L}$  of elution buffer and incubating for 10 min at room temperature before pelleting in a magnetic rack.

#### Flowcell set-up

A brand new, sealed R7.3 flowcell was stored at 4 °C before first use. It was fitted to the MinION™ with plastic screws, ensuring a good thermal contact. The R7.3 flowcell was primed twice using 71  $\mu\text{L}$  premixed nuclease-free water, 75  $\mu\text{L}$  2x running buffer, and 4  $\mu\text{L}$  fuel mix. At least 10 min were needed to equilibrate the flowcell

before each round of priming and before final DNA library loading.

#### Amplicon DNA sequencing

The sequencing mix was prepared with 63  $\mu$ L nuclease-free water, 75  $\mu$ L 2x running buffer, 8  $\mu$ L DNA library, and 4  $\mu$ L fuel mix. A standard 48-h sequencing protocol was initiated using the MinKNOW™ v0.50.1.15. Base-calling was performed through data transference using the Metrichor™ agent v2.29.1 and 2D base-calling workflow v1.16. During the sequencing run, one additional freshly diluted aliquot of DNA library was loaded after 12 h of initial input.

#### Data analysis

Quality assessment of read data and conversion to fasta format was performed using the *poretools* [1] and *poRe* [4] packages. Fasta sequences were filtered by retaining those with a length between 100 and 2000 nt. Read-mapping was performed against the 16S ribosomal RNA sequences for the species present in the mock community (see Availability of supporting data).

Read-mapping was performed using the LAST aligner v.1.89 [18] with parameters -q1 -b1 -Q0 -a1 -e45, which were configured to give the best balance between 16S rDNA assembly length and variants. LAST outputs were converted to sam files and processed with *samtools* [19] to build indexed bam files and obtain consensus sequences from alignments and variant calling. Read-mapping stats from sam files were calculated with the *ea-utils* package and its *sam-stats* function [20]. Different comparisons, GC content correlations, and plots were performed and drawn in R v3.2.0 (<https://cran.r-project.org>). Species coverage was calculated by obtaining fold-change ( $\text{Log}_2$ ) of species-specific read counting against the expected average for the entire community. A coverage bias was assumed when coverage deviation was lower than -1 or higher than 1. Taxonomy assignment of 2d reads was performed using the Silva database [8] and the SINA aligner [19]. Sequences were submitted to the SINA alignment web server using 80 %, 70 %, and 60 % of identity thresholds to ensure a reliable identification. Additional identification at the species level was done using BLAST and the reference NCBI 16S rDNA database.

For the absolute quantification of 16S amplicons we used the following primers: *Escherichia coli* GGACGGGT GAGTAATGTCCTGG and ACCTACTAGCTAATCCC ATCTG; *Clostridium beijerinckii* AGAACCTTACCTA GACTTGACATC and GCTACTAACAATAAGGGTT GCG; and *Bacteroides vulgatus* CACGGGTGAGTAA CACGTATCC and GCATCCCCATCGTCTACCGGAA. Single-stranded DNA (ssDNA), fully covering the respective 16S rDNA regions to amplify for the *E. coli*, *C. beijerinckii*, and *B. vulgatus* species, was obtained from

Isogen Life Science B.V (Utrecht, The Netherlands) where it was synthesized, PAGE-purified, and quantified and used in molecule titration for qPCR. The qPCR was performed on a LightCycler® 480 instrument (Roche Life Science) using the SYBR Green I Master Mix reagent (Roche Life Science), 0.625  $\mu$ M oligos, and 1  $\mu$ L of 1:20 diluted and purified PCR product generated with ‘universal primers’. After 35 cycles of amplification at 95 °C for 10 s, 64 °C for 20 s, and 72 °C for 15 s, absolute quantification was determined using LightCycler® 480 SW v1.5 software (Roche Life Science). Ct values were obtained from serial dilutions of respective ssDNA with known concentrations.

#### Availability of supporting data

Accessions for the 16S ribosomal RNA sequences for the species present in the mock community are available at GenBank: NC\_009085 range c3505652–3504124, NZ\_GG753639 range 96928–98455, NC\_003909 range 82453–83960, NC\_009614 range c4744649–4743140, NC\_009617 range c5775228–5773724, NC\_001263 range c2287019–2285518, NC\_017316 range 213429–214977, NC\_000913 range 4208147–4209688, NC\_000915 range c1512634–1511137, NC\_008530 range c1560731–1559153, NC\_003210 range 243556–245041, NC\_003112 range c2137452–2415909, NC\_006085 range 606163–607687, NC\_002516 range c6044743–6043208, NC\_007493 range 1–1463, NC\_010079 range c2003413–2001874, NC\_004461 range c1816154–1811601, NC\_004116 range 348575–350125, NC\_004350 range 185749–187300, and NC\_003028 range c1815064–1813505). Alternatively, a multi-fasta file containing the 16S reference sequences for the species included in the mock community is available at [https://github.com/alfbenpa/16S\\_MinION](https://github.com/alfbenpa/16S_MinION).

Further supporting data can be found in the *GigaScience* database, GigaDB [21].

#### Abbreviations

BLAST: basic local alignment tool; EC: European commission; ENA: European nucleotide archive; HDF: Hierarchical data format; NCBI: National center for biotechnology information; ONT: Oxford nanopore technologies; OTU: operational taxonomic unit; PCR: polymerase chain reaction; rDNA: DNA encoding for the Ribosomal RNA; rRNA: Ribosomal RNA; SINA: SILVA incremental aligner; USB: universal serial bus; WIMP: what’s in my pot metrichor workflow.

#### Competing interests

ABP is part of the MinION™ Access Programme supported by ONT. Sequencing kits used in this research were partially donated by ONT.

#### Authors’ contributions

ABP and YS designed the study and managed the project. ABP performed the experiments, and analysed and managed the data. ABP, KP, and YS wrote the manuscript. All authors read and approved the final manuscript.

#### Acknowledgements

Authors thank the European 7th Framework Programme for funding ABP and KP, who were supported by the EC Project no. 613979 (MyNewGut).



Received: 27 July 2015 Accepted: 12 January 2016

Published online: 28 January 2016

## References

- Quick J, Quinlan AR, Loman NJ. A reference bacterial genome dataset generated on the MinION portable single-molecule nanopore sequencer. *Gigascience*. 2014;3:22.
- Utturkar SM, Klingeman DM, Land ML, Schadt CW, Doktycz MJ, Pelletier DA, et al. Evaluation and validation of de novo and hybrid assembly techniques to derive high-quality genome sequences. *Bioinformatics*. 2014;30:2709–16.
- Loman NJ, Quinlan AR. Poretools: a toolkit for analyzing nanopore sequence data. *Bioinformatics*. 2014;30:3399–401.
- Watson M, Thomson M, Risse J, Talbot R, Santoyo-Lopez J, Gharbi K, et al. poRe: an R package for the visualization and analysis of nanopore sequencing data. *Bioinformatics*. 2014;31:114–5.
- Klindworth A, Pruesse E, Schweer T, Peplies J, Quast C, Horn M, et al. Evaluation of general 16S ribosomal RNA gene PCR primers for classical and next-generation sequencing-based diversity studies. *Nucleic Acids Res*. 2012;41:e1.
- Loy A, Maixner F, Wagner M, Horn M. probeBase—an online resource for rRNA-targeted oligonucleotide probes: new features 2007. *Nucleic Acids Res*. 2007;35:D800–4.
- Pruesse E, Peplies J, Glockner FO. SINA: accurate high-throughput multiple sequence alignment of ribosomal RNA genes. *Bioinformatics*. 2012;28:1823–9.
- Quast C, Pruesse E, Yilmaz P, Gerken J, Schweer T, Yarza P, et al. The SILVA ribosomal RNA gene database project: improved data processing and web-based tools. *Nucleic Acids Res*. 2013;41:D590–6.
- Koskinen K, Auvinen P, Björkroth KJ, Hultman J. Inconsistent Denoising and Clustering Algorithms for Amplicon Sequence Data. *J Comput Biol* 2014. doi:10.1089/cmb.2014.0268
- Schmidt TS, Matias Rodrigues JF, von Mering C. Ecological consistency of SSU rRNA-based operational taxonomic units at a global scale. *PLoS Comput Biol*. 2014;10:e1003594.
- Schmidt TS, Matias Rodrigues JF, von Mering C. Limits to robustness and reproducibility in the demarcation of operational taxonomic units. *Environ Microbiol* 2014, doi:10.1111/1462-2920.12610
- He Y, Caporaso JG, Jiang XT, Sheng HF, Huse SM, Rideout JR, et al. Stability of operational taxonomic units: an important but neglected property for analyzing microbial diversity. *Microbiome*. 2015;3:20.
- Ashton PM, Nair S, Dallman T, Rubino S, Rabsch W, Mwaigwisya S, et al. MinION nanopore sequencing identifies the position and structure of a bacterial antibiotic resistance island. *Nat Biotechnol*. 2015;33:296–300.
- Kilianski A, Haas JL, Corriveau EJ, Liem AT, Willis KL, Kadavy DR, et al. Bacterial and viral identification and differentiation by amplicon sequencing on the MinION nanopore sequencer. *Gigascience*. 2015;4:12.
- Quick J, Ashton P, Calus S, Chatt C, Gossain S, Hawker J, et al. Rapid draft sequencing and real-time nanopore sequencing in a hospital outbreak of *Salmonella*. *Genome Biol*. 2015;16:114.
- Mikheyev AS, Tin MM. A first look at the Oxford Nanopore MinION sequencer. *Mol Ecol Resour*. 2014;14:1097–102.
- DeSantis TZ, Hugenholtz P, Larsen N, Rojas M, Brodie EL, Keller K, et al. Greengenes, a chimera-checked 16S rRNA gene database and workbench compatible with ARB. *Appl Environ Microbiol*. 2006;72:5069–72.
- Frith MC, Hamada M, Horton P. Parameters for accurate genome alignment. *BMC Bioinformatics*. 2010;11:80.
- Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The sequence alignment/map format and SAMtools. *Bioinformatics*. 2009;25:2078–9.
- Aronesty E. ea-utils: Command-line tools for processing biological sequencing data. <http://code.google.com/p/ea-utils/>; 2011.
- Benitez-Paez A, Portune K, Sanz Y. Supporting information for "Species-level resolution of 16S rRNA gene amplicons sequenced through MinION portable nanopore sequencer". *GigaScience Database* 2016; <http://dx.doi.org/10.5524/100185>.

Submit your next manuscript to BioMed Central and we will help you at every step:

- We accept pre-submission inquiries
- Our selector tool helps you to find the most relevant journal
- We provide round the clock customer support
- Convenient online submission
- Thorough peer review
- Inclusion in PubMed and all major indexing services
- Maximum visibility for your research

Submit your manuscript at  
[www.biomedcentral.com/submit](http://www.biomedcentral.com/submit)

