

DATA NOTE

Open Access

A reference bacterial genome dataset generated on the MinION™ portable single-molecule nanopore sequencer

Joshua Quick^{1,2}, Aaron R Quinlan³ and Nicholas J Loman^{1*}

Abstract

Background: The MinION™ is a new, portable single-molecule sequencer developed by Oxford Nanopore Technologies. It measures four inches in length and is powered from the USB 3.0 port of a laptop computer. The MinION™ measures the change in current resulting from DNA strands interacting with a charged protein nanopore. These measurements can then be used to deduce the underlying nucleotide sequence.

Findings: We present a read dataset from whole-genome shotgun sequencing of the model organism *Escherichia coli* K-12 substr. MG1655 generated on a MinION™ device during the early-access MinION™ Access Program (MAP). Sequencing runs of the MinION™ are presented, one generated using R7 chemistry (released in July 2014) and one using R7.3 (released in September 2014).

Conclusions: Base-called sequence data are provided to demonstrate the nature of data produced by the MinION™ platform and to encourage the development of customised methods for alignment, consensus and variant calling, *de novo* assembly and scaffolding. FAST5 files containing event data within the HDF5 container format are provided to assist with the development of improved base-calling methods.

Keywords: Genomics, Nanopore sequencing

Data description

Single molecule sequencing using biological nanopores was proposed nearly 20 years ago, but formidable technical challenges needed to be overcome before nucleotide sequences could be reliably read [1-5]. In Spring 2014, Oxford Nanopore Technologies released the first commercially available nanopore sequencer to early-access customers. The MinION™ is no larger than a typical smartphone and can connect to and draw power from a laptop computer via its USB 3.0 interface. Sequence data is streamed as DNA fragments translocate through the pore, permitting real-time analysis on an Internet-connected laptop. Portable sequencing may uncover new potential applications, for example near-patient testing and continuous environmental monitoring. We present the first bacterial genome data of the model organism

Escherichia coli K-12 substr. MG1655 sequenced on the MinION™ during the MinION™ Access Program (MAP). Two flowcell chemistries, R7 (released July 2014) and R7.3 (released September 2014) were used. We anticipate this dataset will serve as a useful reference for the community to develop novel bioinformatics methods for this platform [6].

DNA extraction

E. coli K-12 substr. MG1655 was streaked onto plate count agar and incubated for 48 hours at room temperature. Organisms were harvested using a L-shaped spreader and resuspended in 100 µl phosphate buffered saline (PBS). DNA extraction was performed using the Invisorb Spin Cell Mini Kit (Invitex, Birkenfeld, Germany) using the manufacturer's protocol for extraction from serum or plasma.

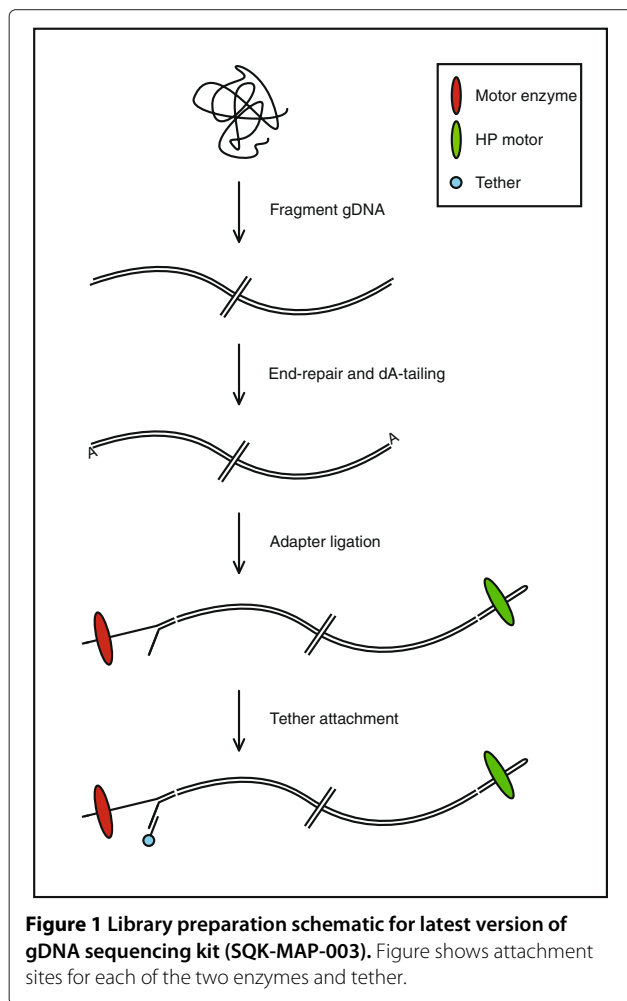
Sequencing library preparation

DNA was quantified using a Qubit fluorometer (Life Technologies, Paisley, UK) and diluted to 23.5 ng/µl.

*Correspondence: n.j.loman@bham.ac.uk

¹ Institute of Microbiology and Infection, University of Birmingham, Birmingham B15 2TT, UK

Full list of author information is available at the end of the article



85 μ l was loaded into a G-tube (Covaris, Brighton, UK) and centrifuged at 5000 rpm for 1 minute before inverting the tube then centrifuging again for 1 minute. The fragmented DNA was end-repaired in a total volume of 100 μ l using the NEBNext End-Repair module (NEB, Hitchin, UK). End-repair was performed as per the manufacturer's instructions except that the incubation time was reduced to 15 minutes. The resulting blunt-ended DNA was cleaned-up using 1.0 \times by volume AMPure XP beads (Beckman Coulter, High Wycombe, UK) according to the manufacturer's instructions with the exception that 80% ethanol was used instead of 70%, and eluted in 25 μ l molecular grade water. A-tailing was performed using the NEBNext dA-tailing module (NEB) in a total volume of

30 μ l according to the manufacturer's instructions with the exception that the incubation time was reduced to 15 minutes. We had concluded from previous experience that incubation times specified are unnecessarily long, and it is likely these incubation times could be further shortened.

Sequencing library preparation: R7-specific

For the R7 chemistry run the Genomic DNA Sequencing Kit (SQK-MAP-002) (Oxford Nanopore Technologies, Oxford, UK) was used to generate a MinION™ sequencing library. To the 30 μ l dA-tailed DNA, 50 μ l Blunt/TA ligase master mix (NEB) was added in addition to 10 μ l each of *Adapter mix* and *HP adapter*. The reaction was left to proceed at room temperature for 10 minutes. The sample was cleaned-up using 0.4 \times by volume AMPure XP beads according to the manufacturer's instructions with the exceptions that the kit supplied wash and elution buffers were used, and only a single wash was carried out. The sample was eluted in 25 μ l of elution buffer. 10 μ l of *Tether* was added and incubated for 10 minutes at room temperature. Lastly, 15 μ l of the hairpin motor (*HP motor*) was added and incubated for 30 minutes at room temperature, giving a total volume of 50 μ l library.

Library preparation: R7.3-specific

In September 2014, an updated Genomic DNA Sequencing Kit (SQK-MAP-003) was released at the same time as a new set of flow cells termed R7.3. In this kit the *HP motor* is prebound to the hairpin adapter, eliminating the incubation step. To the 30 μ l dA-tailed DNA, 50 μ l Blunt/TA ligase master mix (NEB) was added in addition to 10 μ l each of *Adapter Mix* and *HP adapter*, the reaction was left to proceed at room temperature for 10 minutes. The sample was cleaned-up using 0.4 \times by volume AMPure XP beads according to the manufacturer's instructions with the exceptions that the kit-supplied wash and elution buffers were used, and only a single wash was used. The sample was eluted in 25 μ l of elution buffer, leaving to incubate for 10 minutes before pelleting and removal of the library.

Flowcell preparation

For each run a new flowcell was removed from storage at 4°C and the protective packaging removed. The flowcell was fitted to the MinION™ device and held in place with supplied plastic screws to ensure a good thermal contact. 150 μ l *EP buffer* was loaded into the sample loading port using a P1000 pipette and left for 10 minutes to prime

Table 1 Yields for each nanopore run in reads

Run	Filename	Files	Template	Complement	All 2D	Full 2D	
1	R7	Ecoli_NONI.tgz	47195	43656	23338	20087	1598
2	R7.3	Ecoli_R73.tgz	42316	39819	18889	11823	9563

Table 2 Yields for each nanopore run in megabases

	File	Template	Complement	All 2D	Full 2D
1	R7	272.07	125.03	131.44	9.56
2	R7.3	162.68	84.35	64.53	55.68

the flowcell. This priming process was repeated a second time.

Sample loading

Each library was quantified using the Qubit fluorometer. 100 ng (R7) or 350 ng (R7.3) of library was diluted into 146 μ l using *EP buffer* and 4 μ l *Fuel mix* was added. The diluted library was loaded into the sample loading port of the flowcell using a P1000 pipette.

Sequencing

A 72-hour (R7) or 48-hour (R7.3) sequencing protocol was initiated using the MinION™ control software,

MinKNOW™ version 0.45.2.6 (R7) or 0.46.1.9 (R7.3). Read event data were base-called by the software Metrichor™ agent (version 0.16.37960) using workflow 1.0.3 (R7) or 1.2.2 rev 1.5 (R7.3). For the R7 run the flowcell was ‘topped-up’ with a freshly diluted aliquot of library every 12 hours for the first 48 hours.

Data analysis

Read data was extracted from the native HDF5 format into FASTA using poretools [7]. Histograms of read length and collector’s curves of reads were generated using the poretools *hist* and *yield_plot* functions. Alignments were performed against the *E. coli* K-12 MG1655 reference sequence (accession U00096) using the LAST [8] aligner (version 475) with two sets of parameters, each of which were determined to give high mapping rates. Both LAST alignment settings use a match score of 1 (-r1), gap opening penalty of 1 (option -a1) and a gap extension penalty of

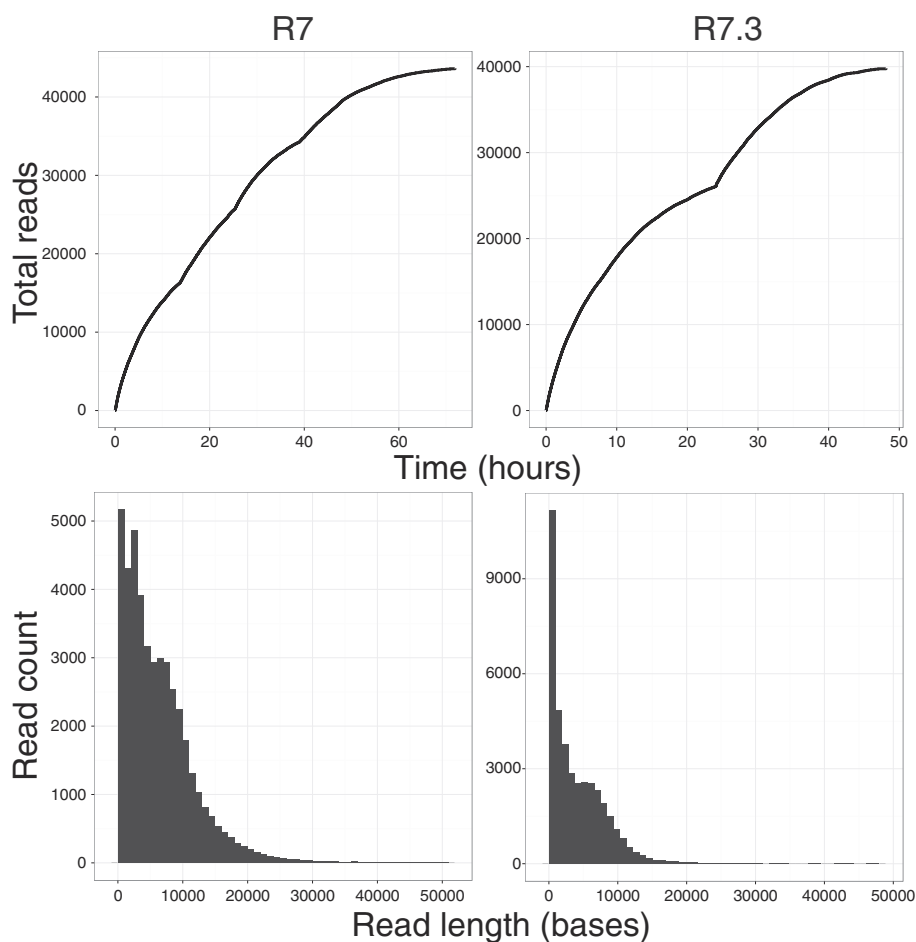


Figure 2 The top row plots demonstrate collector’s curves of sequence reads over time measured in hours for the R7 (left) and R7.3 (right) run. Visible on the R7 run is a change in rate of read acquisition associated with topping-up the flowcell with additional library at hours 12, 24, 36 and 48. For the R7.3 run, the updated MinKNOW software reselecs sequencing wells after 24 hours. The bottom row shows the histogram of read counts for reads shorter than 50,000 bp in length for the R7 (left) and R7.3 (right) run.

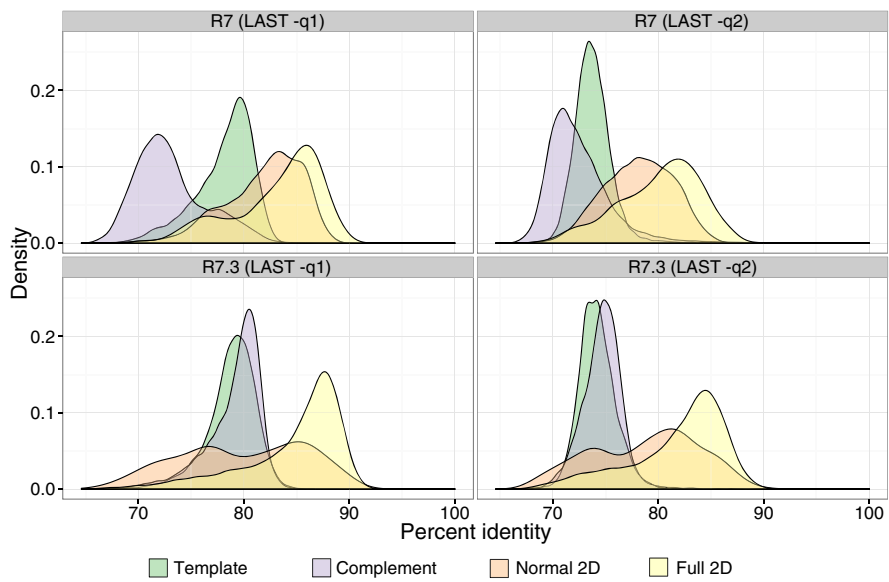


Figure 3 Kernel density plots showing accuracy for R7 and R7.3 chemistries with two different values for the LAST substitution penalty score.

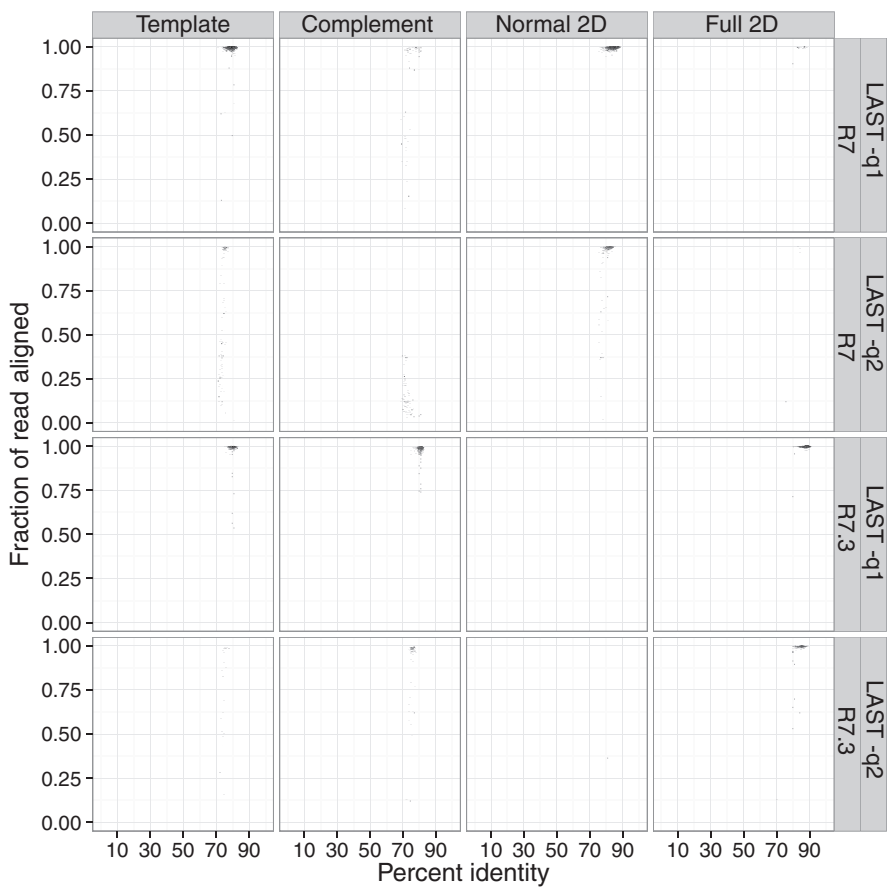


Figure 4 Alignment identity and completeness. Each plot reflects the alignment identity and the proportion of the read aligned for all 2D reads, as well as the underlying template and complement sequences. The top two panels reflect the alignment results for normal and full 2D reads from the R7 flowcell, and the bottom two panels reflect the R7.3 flowcell. Left panels employ a mismatch penalty of 1 and right panels reflect a mismatch penalty of 2. Overall, the lower mismatch penalty increases the identity and fraction of the read that aligned and this effect is greatest for full 2D reads.

Table 3 2D read alignment statistics for each nanopore run

	Run	Norm. 2D	Aligned (q1)	Aligned (q2)	Full 2D	Aligned (q1)	Aligned (q2)
1	R7	18489	85.7%	82.7%	1598	75.5%	69.7%
2	R7.3	2260	26.5%	22.4%	9563	86.9%	66.2%

1 (option -b1). Values of 1 (-q1) and 2 (-q2) were tried for the nucleotide substitution penalty parameters. Read percentage identity is defined as $100 * \text{matches}/(\text{matches} + \text{deletions} + \text{insertions} + \text{mismatches})$. Fraction of read aligned is defined as $(\text{alignment length} + \text{insertions} - \text{deletions})/(\text{alignment length} + \text{unaligned length} - \text{deletions} + \text{insertions})$. Scripts used to generate alignments and plots are available in Github [9].

Results

MinION™ reads can be classified into three types: template, complement and two-direction (2D). Template reads result from the first of two strands presented to the pore. Template strands are slowed down by a proprietary processive motor enzyme which is ligated to the leader adapter. Complement reads may be present if a hairpin has been successfully ligated. The shift from template to complement is recognised by an abasic site on the hairpin which produces a characteristic signal when in contact with the pore. The complement strand is slowed down by a second enzyme termed the *HP motor*. For optimal results, each molecule presented to the MinION™ should have a motor enzyme, tether, hairpin and hairpin motor ligated successfully (Figure 1).

The first run (R7) produced 43,656 template reads (272 Mb) and 23,338 complement reads (125 Mb). Of these, a total of 20,087 were converted into 2D reads (131 Mb) of which 8% (10 Mb) were classified as “full 2D” (explained in detail below). The mean fragment length for 2D reads was 6,543. The second run (R7.3) produced 39,819 template reads (163 Mb) and 18,889 complement reads (84 Mb). Of these, 11,823 were converted into 2D reads (64.53 Mb) of which 86% (55.68 Mb) were full 2D (Table 1 and Table 2). The mean fragment length for 2D reads was 5,458 (Figure 2).

When both template and complement strands are sequenced these are combined by the base-calling algorithm to produce 2D reads, which may be divided into two types. *Normal 2D* are defined by having fewer events detected in the complement strand than in the template strand. This suggests that there was no *HP motor* bound to the hairpin to retard its process through the pore. *Full 2D* reads are defined as having more or equal complement events than template events. These are the optimal type of reads for analysis as they are of the highest quality (Figure 3).

Since 2D reads contain information from both template and complement strands, a non-redundant dataset

would consist of 2D reads plus any remaining template or complement reads which did not get turned into 2D reads.

We compared the alignment characteristics of full 2D sequences, as well their underlying template and complementary strand sequences, using the LAST aligner with two different parameter settings (see Methods). We observe a marked increase in the identity and the proportion of the read aligned when employing equivalent mismatch, gap introduction, and gap extension penalties (Figure 4). The effect is particularly pronounced for the more accurate full 2D reads, and this effect is consistent for both R7 and R7.3. Notably, the proportion of Full 2D reads was substantially higher (22.6% vs. 3.9%) in the R7.3 run, suggesting the possibility that future improvements to the chemistry will increase the yield of the highest quality and most biologically informative reads (the alignment performance of Normal and Full 2D reads is summarized in Table 3).

Discussion

We show that the MinION™ is able to sequence entire bacterial genomes in a single run. Further work is required to determine appropriate algorithms for common secondary analysis tasks such as variant calling and *de novo* assembly. We anticipate and hope this dataset will help stimulate the development of novel methods for handling Oxford Nanopore data.

Availability of supporting data

The datasets supporting the results of this article are available in the GigaDB repository, [6] and the European Nucleotide Archive under accession number ERP007108. This DOI also contains two further MinION™ runs (filename *Ecoli_R7_NONI.tgz*) using R7 chemistry, in which the HP motor incubation was increased from 30 minutes to overnight which was found to increase the percentage of full 2D reads (data not presented in this manuscript).

Abbreviations

MAP: MinION™ Access Programme; PBS: Phosphate buffered saline.

Competing interests

NJL and AQ are part of the Oxford Nanopore MinION™ Access Programme. (MAP). Oxford Nanopore have contributed free of charge early-access reagents in support of the data presented in this manuscript.

Authors' contributions

JQ performed the DNA extraction, library preparation and sequencing and assisted in the data analysis. AQ performed the read alignment analysis. NJL performed the data analysis. NJL and JQ drafted the manuscript. All authors approved the final manuscript.

Acknowledgements

NJL is funded by a Medical Research Council Special Training Fellowship in Biomedical Informatics. JQ is funded by the National Institute for Health research (NIHR) Surgical Reconstruction and Microbiology Research Centre (partnership between University Hospitals Birmingham NHS Foundation Trust, the University of Birmingham and the Royal Centre for Defence Medicine). The views expressed are those of the author(s) and not necessarily those of the NHS, the NIHR or the Department of Health. ARQ was supported by the NIH (NGHRI; 1R01HG006693-01). Funders had no role in the design or carrying out of this work. The Medical Research Council Cloud Infrastructure for Microbial Genomics (CLIMB) platform was used for data analysis. We are grateful to Keith Robison and Minh Duc Cao for their suggestions to help improve the manuscript during the GigaScience open peer review process. The authors would like to thank Damon Huber for providing the K-12 strain used in this study.

Author details

¹Institute of Microbiology and Infection, University of Birmingham, Birmingham B15 2TT, UK. ²NIHR Surgical Reconstruction and Microbiology Research Centre, University of Birmingham, Birmingham B15 2TT, UK. ³Center for Public Health Genomics, University of Virginia, Charlottesville VA 22908, USA.

Received: 26 September 2014 Accepted: 14 October 2014
Published: 20 October 2014

References

1. Baldarelli R, Branton D, Church G, Deamer DW, Kasianowicz J: **Characterization of individual polymer molecules based on monomer-interface interactions.** Google Patents. US Patent 5,795,782 1998
2. Kasianowicz J. J, Brandin E, Branton D, Deamer DW: **Characterization of individual polynucleotide molecules using a membrane channel.** *Proc Natl Acad Sci* 1996, **93**(24):13770–13773.
3. Branton D, Deamer DW, Marziali A, Bayley H, Benner SA, Butler T, Di Ventra M, Garaj S, Hibbs A, Huang X, Jovanovich SB, Krstic PS, Lindsay S, Ling XS, Mastrangelo CH, Meller A, Oliver JS, Pershin YV, Ramsey JM, Riehn R, Soni GV, Tabard-Cossa V, Wanunu M, Wiggin M, Schloss JA: **The potential and challenges of nanopore sequencing.** *Nat Biotechnol* 2008, **26**(10):1146–1153.
4. Stoddart D, Heron AJ, Mikhailova E, Maglia G, Bayley H: **Single-nucleotide discrimination in immobilized DNA oligonucleotides with a biological nanopore.** *Proc Natl Acad Sci* 2009, **106**(19):7702–7707.
5. Cherf GM, Lieberman KR, Rashid H, Lam CE, Karplus K, Akeson M: **Automated forward and reverse ratcheting of DNA in a nanopore at 5-a precision.** *Nat Biotechnol* 2012, **30**(4):344–348.
6. Quick J, Loman NJ: **Bacterial whole-genome read data from the Oxford Nanopore Technologies MinION™ nanopore sequencer.** *GigaScience Database* 2014, [http://dx.doi.org/10.5524/100102]
7. Loman NJ, Quinlan AR: **Poretools: a toolkit for analyzing nanopore sequence data.** *Bioinformatics* 2014, doi:10.1093/bioinformatics/btu555
8. Frith MC, Hamada M, Paul H: **Parameters for accurate genome alignment.** *BMC Bioinformatics* 2010, **11**(80): doi:10.1186/1471-2105-11-80
9. Quinlan AR: **GitHub.** 2014, [https://github.com/arq5x/nanopore-scripts]

doi:10.1186/2047-217X-3-22

Cite this article as: Quick *et al.*: A reference bacterial genome dataset generated on the MinION™ portable single-molecule nanopore sequencer. *GigaScience* 2014 **3**:22.

Submit your next manuscript to BioMed Central and take full advantage of:

- Convenient online submission
- Thorough peer review
- No space constraints or color figure charges
- Immediate publication on acceptance
- Inclusion in PubMed, CAS, Scopus and Google Scholar
- Research which is freely available for redistribution

Submit your manuscript at
www.biomedcentral.com/submit

