

# SCIENTIFIC REPORTS



OPEN

## Improved metagenome assemblies and taxonomic binning using long-read circular consensus sequence data

Received: 06 October 2015

Accepted: 12 April 2016

Published: 09 May 2016

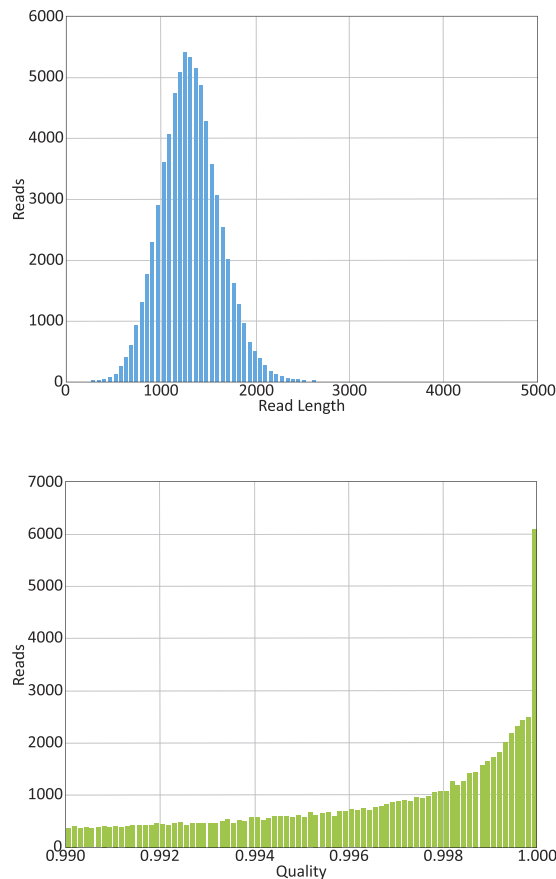
J. A. Frank<sup>1</sup>, Y. Pan<sup>2</sup>, A. Tooming-Klunderud<sup>3</sup>, V. G. H. Eijsink<sup>1</sup>, A. C. McHardy<sup>2</sup>,  
A. J. Nederbragt<sup>3</sup> & P. B. Pope<sup>1</sup>

DNA assembly is a core methodological step in metagenomic pipelines used to study the structure and function within microbial communities. Here we investigate the utility of Pacific Biosciences long and high accuracy circular consensus sequencing (CCS) reads for metagenomic projects. We compared the application and performance of both PacBio CCS and Illumina HiSeq data with assembly and taxonomic binning algorithms using metagenomic samples representing a complex microbial community. Eight SMRT cells produced approximately 94 Mb of CCS reads from a biogas reactor microbiome sample that averaged 1319 nt in length and 99.7% accuracy. CCS data assembly generated a comparative number of large contigs greater than 1 kb, to those assembled from a ~190x larger HiSeq dataset (~18 Gb) produced from the same sample (i.e. approximately 62% of total contigs). Hybrid assemblies using PacBio CCS and HiSeq contigs produced improvements in assembly statistics, including an increase in the average contig length and number of large contigs. The incorporation of CCS data produced significant enhancements in taxonomic binning and genome reconstruction of two dominant phylotypes, which assembled and binned poorly using HiSeq data alone. Collectively these results illustrate the value of PacBio CCS reads in certain metagenomics applications.

Metagenome assembly is a key methodological stage in all environmental sequencing projects, which has significant repercussions on all down-stream analyses such as taxonomic classification, genome reconstruction, and functional gene annotation. It is commonly a very complex process, with many sequencing platform-specific issues such as read length and number. Similarly, there are also many sample-specific issues such as the numbers, frequencies, types and sizes of microbial genomes present in highly diverse communities. The goal of metagenomic assemblies is relatively straightforward: obtain large contig sizes coupled with the fewest possible misassemblies. However, metagenomic assemblies often consist of a fragmented collection of short contigs, which are difficult to taxonomically and functionally assign accurately. There are at least two current approaches to metagenomic assembly: (i) assembly of all data<sup>1</sup>, which is typically computationally demanding, or (ii) using binning or normalization methods to select subsets of reads that are then assembled separately<sup>2,3</sup>. Methods that use data from multiple sequencing platforms are still infrequent, despite indications that combined approaches yield improvements in contig length and integrity<sup>4</sup>.

Current sequencing technologies offer a range of read lengths. Methods that produce short reads (<250 nucleotides (nt)) such as Illumina can generate high sequencing depth with minimal costs, however when used for analyzing complex communities data assembly typically requires massive computational resources and the resulting contigs remain relatively short<sup>1</sup>. In theory, longer read sequencing technologies can overcome many of the known assembly problems associated with short reads, however these technologies have traditionally been accompanied with one or more inherent shortcomings, such as lower sequencing depth, higher costs and higher error rates. Several technologies exist that can produce longer reads. For example, Ion Torrent and Roche 454 offer read

<sup>1</sup>Department of Chemistry, Biotechnology and Food Science, Norwegian University of Life Sciences, Ås, 1432 Norway. <sup>2</sup>Computational Biology of Infection Research, Helmholtz Centre for Infection Research, Inhoffenstraße 7, 38124 Braunschweig, Germany. <sup>3</sup>University of Oslo, Department of Biosciences, Centre for Ecological and Evolutionary Synthesis, Blindern, 0316 Norway. Correspondence and requests for materials should be addressed to P.B.P. (email: phil.pope@nmbu.no)



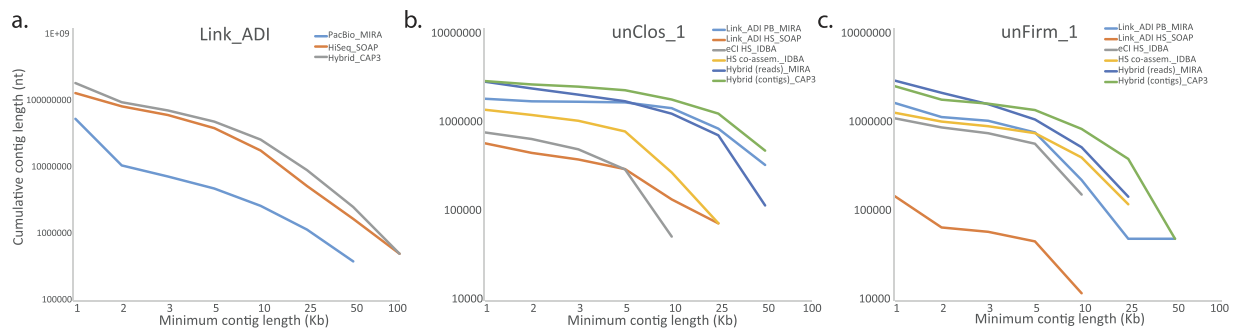
**Figure 1. Read length and quality distribution of PacBio “Circular Consensus Sequence” (CCS) reads produced from a Link\_ADI-derived shotgun library (~1.5 kb inserts) sequenced on a PacBio RS II instrument using P4-C2 chemistry.** In total, eight SMRT cell were used for sequencing. (a) Read length distribution of PacBio CCS reads that passed a 0.99 quality score for which an average of 10 insert passes was required. (b) Quality distribution of the 71,254 PacBio CCS reads that passed the 0.99 cutoff using the SMRT portal (average 99.7%).

lengths of up to 400 nt and 1000 nt, respectively, but these technologies are more costly per base pair and are vulnerable to generating homopolymer (single-nucleotide repeats) sequencing errors. Pacific Biosciences (PacBio) has designed a sequencing technology based on single-molecule, real-time (SMRT) detection that can provide much greater read lengths, with ~50% of reads in a single run exceeding 14 kb and 5% exceeding 30 kb<sup>5</sup>. High error rates, reported as high as 15% in individual reads, have previously prevented the use of raw PacBio reads in metagenomics<sup>6,7</sup>. Interestingly, the error rates may be reduced by using circular consensus sequencing (CCS) that entails the repeated sequencing of a circular template, and subsequent generation of a consensus of individual DNA inserts. Consensus quality increases with each sequencing pass, and this approach can ultimately result in high-quality sequences of about 500 to ~2,500 nt in length with greater than 99% accuracy (Q20 or better)<sup>8,9</sup>.

Here, we present various applications of PacBio CCS data in a metagenomic analysis of the complex microbial community in a commercial biogas reactor. We compare individual assemblies of short read HiSeq2000 and PacBio CCS data as well as hybrid assemblies of subsets from both platforms. PacBio CCS data provides a dramatic improvement in the assembly of universal marker genes in comparison to HiSeq2000 data, allowing for custom training data for phylogenomic binning algorithms and accurate taxonomic binning of assembled contigs from both data types. Subsequently this enabled enhancements in genome reconstructions of uncultured microorganisms that inhabit complex communities.

## Results

**PacBio CCS reads improve assembly statistics.** For the purpose of this study we analyzed and compared two sequence datasets generated from the same biological sample, a methanogenic biogas reactor microbiome containing an estimated 480 individual phylotypes, hereafter referred to as Link\_ADI (Table S1). These datasets comprised approximately one lane of HiSeq sequence data and data from eight PacBio SMRT cells, respectively. HiSeq sequencing entailed 175 nt library construction and generation of  $2 \times 100$  nt paired end sequence data, totaling approximately 149 million read pairs (18.5 Gb). For PacBio, a library was constructed with inserts of approximately 1.5 kb, which were sequenced using a RS II instrument and P4-C2 chemistry. A total of 522,695 PacBio reads were generated with a mean accuracy of 86%, totaling approximately 3.3 Gb. Of these reads, 71,254 were CCS that averaged 99.7% accuracy and 1,319 nt in length (totaling 95.4 Mb) (Fig. 1). Given the two



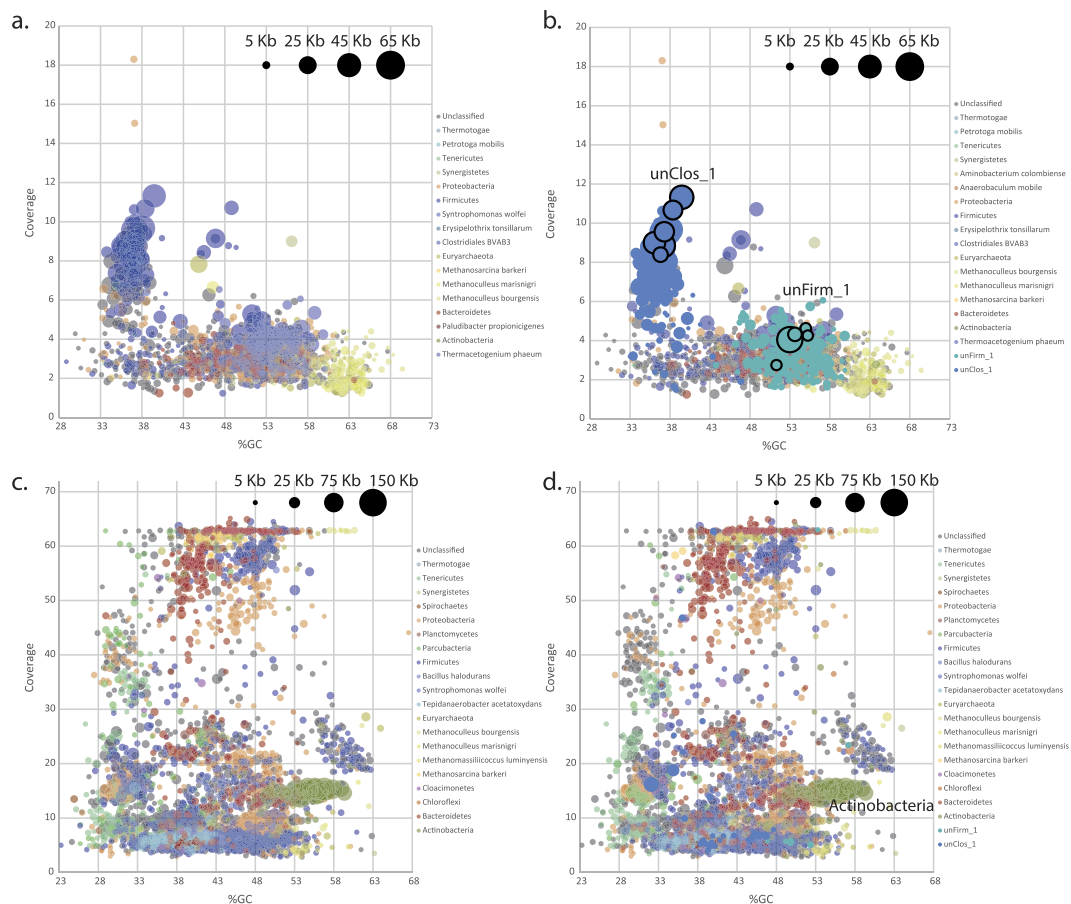
**Figure 2.** Cumulative number of assembled nucleotides in contigs of different minimum lengths for (a) Link\_ADI, (b) unClos\_1, and (c) unFirm\_1. Each line corresponds to a different sample (Link\_ADI or eCI, where noted), sequencing method (HiSeq [HS] or PacBio [PB]), different assembly method (co-assembly across samples Link\_ADI and eCI, hybrid using mapped reads from HiSeq and PacBio, or hybrid using contigs from HiSeq and PacBio), or assembly program (CAP3, IDBA\_UD, MIRA, or SOAPdenovo).

different sequencing platforms, multiple assembly algorithms were used. MIRA 4.0<sup>10</sup> was used to assemble the PacBio CCS reads, which resulted in approximately 46% of the CCS reads assembling into 2,181 contigs averaging 4,459 nt with the max contig length of 65,165 nt (Table S2). SOAPdenovo2<sup>11</sup> was used to assemble approximately 35.6% of the HiSeq reads generated for Link\_ADI (~96 million read pairs unassembled), which produced 55,633 contigs greater than 1 kb (average contig length: 2411 nt) with a maximum length of 148,797 nt.

Comparing the assembly statistics from the two assemblies showed that, despite the much smaller size of the raw PacBio CCS dataset (around 190-fold less sequence), the total length of large contigs produced from the MIRA assembly was in the range of those produced from the HiSeq assembly (Fig. 2 and Table S2). The MIRA assembly produced 34,513 contigs and unassembled reads that were greater than 1 kb in length, which totaled approximately 54.9 Mb (Table S2). In contrast, the HiSeq assembly generated 55,633 contigs greater than 1 kb (134.2 Mb). The total size of the 100 biggest MIRA contigs totaled 52% of the equivalent HiSeq subset. Attempts to perform hybrid assemblies using raw HiSeq and PacBio CCS reads were ultimately unsuccessful, presumably due to the large number of sequencing reads and a paucity of algorithms customized for this particular hybrid input (to our knowledge). Therefore, as an alternative we used a downstream approach that was more amenable to our datasets and available assemblers. Both subsets of assembled HiSeq and CCS contigs greater than 1 kb (including unassembled CCS reads > 1 kb) were further assembled using the “Sanger”-era program CAP3<sup>12</sup>, which was designed for use with long sequencing reads. In total, 21.31% and 10.98% of PacBio and HiSeq contigs greater than 1 Kb (respectively) assembled into a dataset that included only 4767 hybrid contigs, with the remaining 90,183 contigs not assembling. Despite the modest incorporation rate, the assembly provided an increase in cumulative nucleotides from contigs larger than 10 kb (PacBio + HiSeq: 21.01 Mb, Hybrid: 26.8 Mb) and 25 kb (PacBio + HiSeq: 6.5 Mb, Hybrid: 9.3 Mb) (Fig. 2 and Table S2).

**PacBio CCS reads improve genome binning of difficult to assemble phylotypes.** Community characterization of Link\_ADI using short subunit (SSU) rRNA gene amplicon analysis identified approximately 480 individual phylotypes, of which two exhibited high relative abundance and no close taxonomic relationship to cultivated bacterial species (Table S1). Phylotype unClos\_1 is an as-yet uncultured bacterium affiliated to the Clostridiales family and was estimated to represent ~36% of the total microbiome, whereas unFirm\_1 is a deeply-branched uncultured representative affiliated to the Firmicutes, accounting for ~5%. In order to functionally characterize both phylotypes and determine their contribution to the microbiomes metabolic network, we sought to reconstruct and annotate their genomes. Given the high levels of relative abundance, both organisms were anticipated to be represented by high DNA levels within the metagenomic datasets, and thus conducive to greater assembly in terms of coverage and contig length. First pass comparisons of the assembled HiSeq contigs focusing on contig coverage, size and GC %, gave no clear patterns that are indicative of several numerically dominating organisms (i.e. a cluster of large high-coverage contigs within a narrow GC % range, Fig. 3c). In contrast, coverage vs GC % comparisons of assembled PacBio CCS contigs revealed one clear cluster of higher coverage contigs that were large and within a narrow GC % range (Fig. 3a).

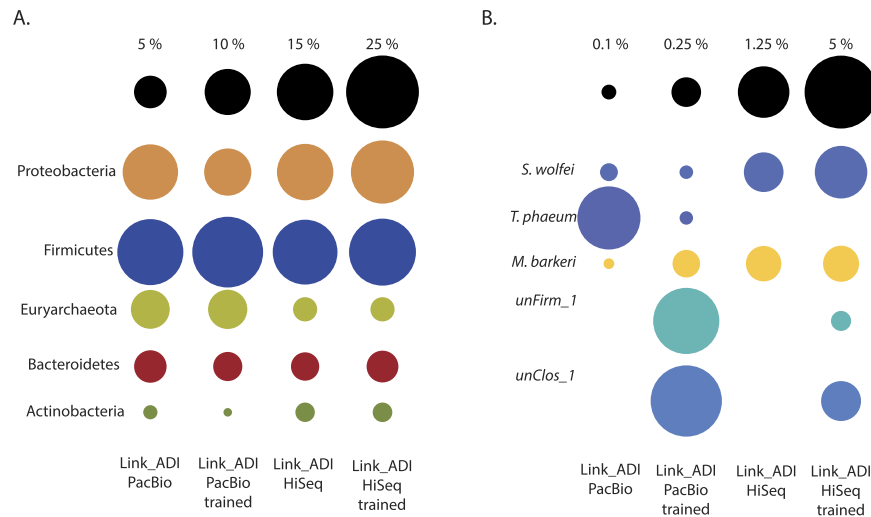
Phylogenomic binning methods were subsequently used in attempts to recover genome sequence information for unClos\_1 and unFirm\_1 and for as many other phylotypes as possible. The presence of only one biological sample and DNA extraction, pre-determined the use of sequence compositional binning algorithms and prevented the use of temporal and/or multi-sample binning methods that have been recently shown to produce accurate genomes from metagenomic datasets<sup>13,14</sup>. PhyloPythiaS+<sup>15</sup> was initially used to assign taxonomy to PacBio CCS and HiSeq contigs (greater than 1 kb), which produced very few taxonomic assignments to a strain or species level (Table S3). Instead, the vast majority of contigs were binned to higher-ranking taxa at a phylum or order level, implying that the data provides limited functional and structural insights into the individual organisms making up the microbial community. This result was not unexpected as the SSU rRNA gene analyses indicated that the Link\_ADI microbiome is composed of uncharacterized species (Table S1) that are distantly related to the available prokaryotic genomes in NCBI used to train PhyloPythiaS+.



**Figure 3.** Visualization of GC %, coverage and size of assembled contigs generated from PacBio CCS (a,b) and HiSeq data (c,d) from a biogas reactor microbiome (Link\_ADI). Contigs are coloured based on taxonomic binning that was performed using PhyloPythiaS+ under default settings (a,c) and after including custom phylotype-specific training data (b,d). Contig lengths are indicated by circle sizes. PacBio CCS contigs that contain marker genes and were used as training data for phylotype unClos\_1 and unFirm\_1 are outlined in black. For the purposes of clarity, only HiSeq contigs greater than 5 kb are represented (c,d).

In cases where PhyloPythiaS and its predecessors have had phylotype-specific training data (at least 100 kb) from a given metagenome, the binning and genome reconstruction of the target phylotype has proven to be highly accurate<sup>16,17</sup>. Therefore, to improve the resolution of PhyloPythiaS+ we compiled as much phylotype-specific training data as possible. All contigs were evaluated for coverage vs. GC % metrics and the presence of taxonomically informative marker genes<sup>18</sup>, with the aim of identifying contigs that correspond to the abundant phylotypes identified in our samples and can therefore be used as training data. The complexity and fragmented nature of the HiSeq assembly (Fig. 3c) made identification of species-specific genome information problematic. This had direct implications on the ability to obtain the ~100 kb high-confidence assemblages of training data that are required for accurate species level binning<sup>17</sup>. However, the increased length and improved clustering of the assembled PacBio CCS contigs provided large and accurate training data collections for unClos\_1 and unFirm\_1 in particular. We pooled together six contigs totaling 200 kb for unClos\_1 and seven contigs totaling 107 kb for unFirm\_1 (Highlighted in Fig. 3b). Interestingly this included large contigs that encoded complete SSU rRNA operons, which are notoriously difficult to assemble using short-read NGS data, such as reads obtained using HiSeq. In total, we identified 17 SSU rRNA gene fragments in the PacBio CCS contigs and 86 when including unassembled reads (compared to six in the HiSeq contigs greater than 1 kb). For unClos\_1, we identified 16S rRNA genes in three contigs that totaled 96 kb in length.

Both the total collection of HiSeq contigs greater than 1 kb and the PacBio CCS contigs, including unassembled reads, were binned with the custom training model for PhyloPythiaS+, that includes all the available prokaryotic genomes in NCBI and the two phylotype-specific contig subsets described above. The output produced a greatly improved recovery of phylotype-level binning for both unClos\_1 and unFirm\_1 in both HiSeq and PacBio CCS contigs from Link\_ADI (Fig. 4). For unClos\_1, 189 PacBio sequences (PacBio contigs and unassembled CCS reads, totaling 1,913,759 nt) and 182 HiSeq contigs (600,903 nt) were assigned to the phylotype (Table S2). 576 PacBio sequences (1,710,231 nt) and 77 HiSeq contigs (151,790 nt) were binned to unFirm\_1. The binning of unClos\_1 and unFirm\_1 contigs also revealed patterns that indicate assembly differences between PacBio CCS and HiSeq. Despite the indications from the SSU rRNA gene amplicon analyses that phylotypes unClos\_1 and



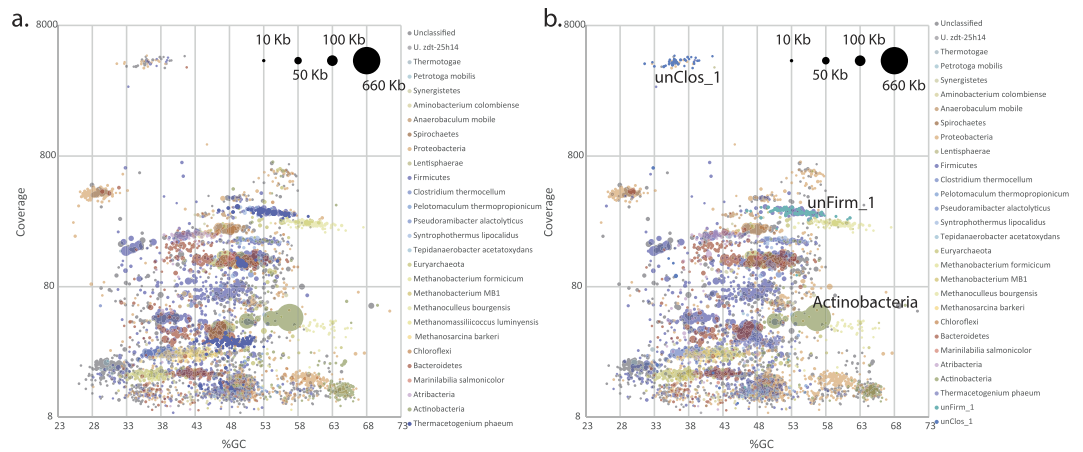
**Figure 4. Selected taxonomic bins generated via PhyloPythiaS+ binning using default settings with and without use of custom training data.** Circle size indicates relative bin size; for complete binning information see Table S3. The proportion of total DNA binned in the major phyla (A) represented in the Link\_ADI microbiome was similar for both PacBio CCS and HiSeq contigs regardless of the use of training data. However, use of training data enhanced the recovery of unClos\_1 and unFirm\_1 (B) in both the PacBio and HiSeq assemblies. Differences between the sequencing methods were also evident at a species level where some abundant species assembled and binned better with PacBio (*Thermacetogenium phaeum*, unClos\_1, and unFirm\_1), whereas others produced better results with HiSeq data (*Syntrophomonas wolfei* and *Methanosarcina barkeri*).

unFirm\_1 were the most abundant in Link\_ADI, neither phylotype were attributed to the longest HiSeq contigs (Fig. 3d). Nine of the ten largest HiSeq contigs from Link\_ADI binned to the Order Actinomycetales (Fig. 3c), totaling around 2.2 Mb over 203 contigs (Table S2). Only one phylotype affiliated to the Actinomycetales was identified in SSU rRNA gene amplicon analysis, which was ranked 61<sup>st</sup> most abundant (Table S1). In addition, the coverage for each of the Actinomycetales-affiliated HiSeq contigs was on average approximately two-fold higher than the contigs binning as unClos\_1 (Fig. 3d). In contrast, the Actinomycetales-affiliated PacBio CCS contigs were much shorter and exhibited lower coverage than unClos\_1 (Figs 3 and 4).

The integrity of the Link\_ADI HiSeq and PacBio CCS binning results for both unClos\_1 and unFirm\_1 phylotypes were evaluated by determining completeness and redundancy in each respective reconstructed genome. Genomes were evaluated in terms of conserved single-copy gene (CSCG) coverage<sup>19</sup>, which was significantly lower for unClos\_1 HiSeq contigs (24 of 107, 4 redundancies), compared to PacBio CCS contigs (96 of 107, 10 redundancies). For unFirm\_1, only four CSCGs (no redundancies) were identified in binned HiSeq contigs, whereas 60 (6 redundancies) were found in the corresponding PacBio CCS dataset. Overall, these results illustrate that for both these numerically abundant phylotypes, PacBio CCS data produced greater assembly and binning results, subsequently generating more-competent reconstructed genomes.

The custom trained PhyloPythiaS+ with training data obtained from the PacBio CCS contigs also showed enhanced binning when used for other biological samples and metagenomics datasets where unClos\_1 and unFirm\_1 were found (Fig. 5). An independently created cellulose enrichment (eCI) was inoculated from Link\_ADI and exhibited reduced species complexity, with both unClos\_1 and unFirm\_1 demonstrating comparable numerical dominance (~48% and ~7%, respectively) (Table S4). Similar to the Link\_ADI HiSeq dataset that was assembled with SOAPdenovo, assembly of eCI with IBDA\_UD<sup>20</sup> did not generate long marker-gene encoding contigs representative of unClos\_1 and unFirm\_1, and phylotype-specific binning was not possible using this dataset alone (Fig. 5a). Therefore, training data generated from the Link\_ADI PacBio CCS dataset was used to taxonomically bin the eCI HiSeq dataset (Fig. 5b). The binning produced after training improved cluster visualization, and binning assignments were concurrent with coverage vs GC % comparisons, which indicated explicit clusters for each phylotype (Fig. 5b). Subsequently, the recovery of genomic information linked to the unClos\_1 and unFirm\_1 phylotypes was substantially larger (Table S3). Similar to the Link\_ADI SOAP\_denovo assembly, discrepancies were also noted in the eCI IDBA\_UD assembly, where the most abundant organisms (unClos\_1 and unFirm\_1) did not assemble into the largest contigs. Instead, Actinobacteria-affiliated DNA from one predicted phylotype (0.45% relative abundance: Table S4) formed the largest contigs (Fig. 5).

**Hybrid assembly of genome bins improves overall genomic reconstruction.** In an effort to reconstruct improved genomes for both unClos\_1 and unFirm\_1, we used a two-step hybrid assembly approach that was refined to include only PacBio and HiSeq data that binned to either phylotype. With the intention of generating as complete as possible genomes, we used all genomic material that was available for both phylotypes from both the Link\_ADI and eCI samples. Binned HiSeq contigs from Link\_ADI and the cellulose enrichment eCI datasets were first deconstructed into individual reads and then pooled into one file prior to assembly using



**Figure 5. Visualization of GC %, coverage and size of assembled contigs generated from eCI HiSeq data.** Sample eCI originated from a lab-scale enrichment grown on cellulose that was inoculated from Link\_ADI. Contig lengths are indicated by circle sizes. Contigs are coloured based on phylogenetic binning that was performed using PhyloPythiaS+ under default settings (a) and PacBio-derived custom phylotype-specific training data (b). For the purposes of clarity, only HiSeq contigs greater than 5 kb are represented.

IBDA\_UD. This initial “phylotype-specific” HiSeq co-assembly improved the cumulative nucleotides from larger contigs (Table S2), which were then assembled together with Pacbio CCS contigs and unassembled reads binned to the same phylotype. This phylotype-specific hybrid approach improved genome reconstruction in terms of total genome size as well as improved average contig length and large contig assembly (Fig. 2b,c and Table S2). For unClos\_1, a total of 811 out of 1178 sequences (PacBio contigs, unincorporated PacBio reads, and co-assembled Link\_ADI and eCI HiSeq contigs; 3,350,596 nt in length) assembled into 71 hybrid contigs (909,704 nt) and 367 unincorporated sequences (2,120,602 nt), totaling 3,030,306 nt. Average contig lengths increased 127%, whereas improvements in large contig assembly was particularly evident for contigs greater than 25 Kb (136%) and 50 Kb (145%) (Table S2). For unFirm\_1, 520 out of 1,212 sequences (3,037,687 nt) were assembled into 123 hybrid contigs (1,086,984 nt) and 692 unincorporated sequences (1,563,729), totaling 2,650,713 nt. Average contig lengths increased 118% and large contig assembly was significantly improved for contigs greater than 10 Kb (217%) and 25 Kb (796%) (Table S2). Hybrid MIRA assemblies that used the individual sequencing reads (that formed the original contigs) instead of a two-step approach using CAP3, resulted in contigs that were on average smaller for both unClos\_1 and unFirm\_1 (Fig. 2b,c and Table S2).

## Discussion

Many of the commonly used second generation sequencing methods in (meta) genome sequencing provide gigabases of data. While this provides high levels of sequencing depth per sample, the short read lengths can restrict the ability to assemble longer contigs, particularly when evaluating complex microbial communities. Specific exemplary problems include the presence of genes with low evolutionary divergence between organisms or repetitive genomic regions that are larger than a sequencing read (e.g., rRNA operons). One way of circumventing this is by combining multiple sequencing technologies that can overcome each other’s limitations. For example, Illumina HiSeq provides high sequencing depth, but with low sequencing breadth; in other words this technique has a high ability to sample across multiple genomes with the drawback that individual reads sample a very small proportion of each genome. This can be complemented by additional PacBio sequencing, which has high breadth (providing at least 10–30-fold more data per read), but a lot lower depth. By combining the two methods, one has a higher probability of covering regions problematic for short read sequencing methods. Several studies have illustrated this convincingly for bacterial genomes, where a hybrid Illumina-PacBio approach has enabled near-complete chromosome closure with no necessary secondary sequencing or primer-walking methods<sup>21</sup>. Previously, the high error rate of PacBio reads (~10–15%) has prevented their use in metagenomic analysis of complex communities, where the coverage required to compensate the erroneous reads was not financially or technically feasible. However, use of the CCS provides high quality long reads that are suitable for metagenomic applications. Here we illustrate the features that PacBio CCS data may bring to a metagenomics project, with respect to increased contig lengths, assembly of problematic genomic regions, improved phylogenomic binning, and genome reconstruction of the uncultured phylotypes that dominate microbial communities.

Specific benefits of the PacBio CCS contigs for Link\_ADI were the considerably larger average contig sizes as well as the number of large contigs, with the later being comparable to the HiSeq assembly that was generated from 190-fold more data. In metagenomic analyses, larger contigs are key to producing higher quality output that is needed for downstream applications such as taxonomic assignments<sup>17</sup>, gene calling, and annotation of operons that often exceed 10 kb in length<sup>16</sup>. The assembly output from both platforms varied considerably in both contig size and distribution (Figs 2 and 4 and Table S2). In particular, numerically dominating organisms did not necessarily assemble into the largest HiSeq contigs (Figs 3b,d and 5), irrespective of species diversity or the assembly algorithms used (Link\_ADI: SOAP\_denovo, eCI: IDBA\_UD), which in contrast transpired for PacBio CCS contigs (Fig. 3a,b). Despite the similar size of the PacBio CCS and HiSeq > 1 kb contig datasets available for

binning, the size of the unClos\_1 and unFirm\_1 genomic bins obtained from the PacBio CCS data were, on average, ~3x and ~6x larger, respectively (Fig. 4 and Table S2). Another observation was the examples of PacBio CCS contigs containing difficult to assemble regions such as SSU rDNA. On average, PacBio CCS contigs that contained relevant SSU rDNA data were 15-fold larger than the SSU rDNA containing HiSeq contigs. Conventional composition-based binning was shown to be substantially improved with the addition of PacBio-derived custom training data that contained genomic information specific for unClos\_1 and unFirm\_1 (Fig. 4 and Table S3). The collection of these phylotype-specific training subsets was only possible in the PacBio CCS contig dataset, since neither phylotype produced contigs of sufficient length in HiSeq datasets. Hence, this approach presents an alternative means to reconstruct genomes in instances where phylotypes are not conducive to HiSeq assembly and experimental design that will not allow multiple sample timepoints or several differential DNA extractions, which are necessary for accurate binning algorithms that use differential coverage of populations<sup>13,14</sup>.

Whilst this study shows the potential value PacBio CCS reads can exert upon a metagenomics study, there is certainly room for improvement. The comparative high costs of PacBio data (approximately 8 times the cost per Gb of data, this study), can restrict the depth of raw data used. Moreover, one of the key concerns with the use of PacBio CCS reads is data wastage with respect to the number of reads generated and the number that pass CCS quality cutoffs. Upcoming PacBio upgrades will increase read lengths and produce a higher amount of high-quality CCS reads per SMRT cell, which will generate greater assemblies and thus less wastage. Notably, closer examination reveals that read wastage is also applicable for the use of Illumina in metagenomic applications. For example, in the present study only 35.6% of the paired-end HiSeq reads assembled into contigs greater than 1,000 nt, an arbitrary cutoff that is used in many metagenomic analyses. The improved PacBio CCS assembly statistics for the two dominant phylotypes, also suggests that greater depth of PacBio CCS data will increase read incorporation rates and average contig lengths in assemblies of lesser abundant phylotypes within complex communities. Whilst the hybrid assemblies of the PacBio CCS and HiSeq contigs from the large Link\_ADI metagenomes improved assembly statistics, they only produced modest incorporation rates, presumably due to low levels of overlap between the two datasets. The observed low level of hybrid assembly overlap is possibly attributed to the relatively low amount of raw data used and the high species complexity of the Link\_ADI sample, which contains approximately 480 phylotypes. As expected, hybrid assembly overlap was improved for contigs that were taxonomically assigned to a phylotype prior to hybrid assembly (Table S2).

Despite the relative small size of the PacBio CCS dataset and high species complexity of the sample, hybrid assemblies for both the total community dataset and phylotype-specific bins produced improvements (Fig. 2 and Table S2), and this represents just a start. In the future, there will be access to better long read data and it is anticipated that further improvement of assembly algorithms customized to incorporate multiple sequencing technology inputs will improve hybrid assembly performance. Regardless, these aspects need further attention in moving forward, so that the full potential of longer read technology can be exploited to deepen our insight into complex microbial communities. This study also shows that as long reads become more common, they will make further software extensions of binning algorithms such as PhyloPythiaS+ very valuable and will allow automatic assignment of training contigs to novel phylotypes and not just the higher ranking assignments. Increased capabilities to reconstruct accurate genomes representative of uncultured microorganisms are of major importance since they allow accurate mapping of community metabolism and are a prerequisite for meaningful “meta-omic” studies that may reveal genes and/or proteins with novel functions that cannot be recognized by bioinformatics alone.

## Methods

**Samples.** Sample Link\_ADI was obtained from a commercial biogas reactor in Linköping, Sweden, fed on a mixture of slaughterhouse waste, food waste, and plant biomass (Reactor I)<sup>22</sup>. Sample eCI was taken from a batch enrichment using the same commercial biogas plant as inoculum source and cellulose as substrate<sup>23</sup>.

**DNA extraction and sequencing.** Total genomic DNA was prepared using the FastDNA Spin Kit for Soil (MP Biomedicals, Santa Ana, CA, USA). For both Link\_ADI and eCI, an aliquot of 200 µl was used for DNA extraction following the manufacturer’s protocol. For SSU rRNA gene sequencing, library preparation was performed as per manufacturers recommendations (Illumina, 2013). V3 and V4 regions of bacterial SSU rRNA genes were amplified using the 341F (5’-TCGTCGGCAGCGTCAGATGTGTATAAGAGACAGCCTACGGGNGGCWGCAG-3’) and 785R (5’-GTCTCGTGGGCTCGGAGATGTGTATAAGAGACAGGGACTACHVGGGTATCTAATCC-3’) modified primer set<sup>24</sup>, where the underlined sequence corresponds to the Illumina adaptor. The amplicon PCR reaction mixture (25 µl) consisted of 12.5 ng microbial gDNA, 12.5 µl iProof HF DNA polymerase mix (BioRad) and 0.2 µM of each primer. The PCR reaction was performed with an initial denaturation step at 98 °C for 30 s, followed by 25 cycles of denaturation at 98 °C for 30 s, annealing at 55 °C for 30 s, and extension at 72 °C for 30 s followed by a final elongation at 72 °C for 5 min. A new PCR reaction was carried out to attach unique 6 nt indices (Nextera XT Index Kit) to the Illumina sequencing adaptors to allow multiplexing of samples. The PCR conditions were as follows: 98 °C for 3 min., 8 cycles of 95 °C for 30s., 55 °C for 30s., and 72 °C for 30 s., followed by a final elongation step at 72 °C for 5 min. AMPure XP beads were used to purify the resulting 16S rRNA amplicons. The 16S rRNA amplicons were quantified (Quant-IT™ dsDNA HS Assay Kit and Qubit™ fluorometer, Invitrogen, Carlsbad, CA, USA), normalized and then pooled in equimolar concentrations. The multiplexed library pool was then spiked with 25% PhiX control to improve base calling during sequencing. A final concentration of 8 pM denatured DNA was sequenced on an Illumina MiSeq instrument using the MiSeq reagent v3 kit chemistry with paired end, 2 × 300 bp cycle run. HiSeq Shotgun sequencing runs were performed on libraries (175 nt, to ensure overlap and allow for merging of the paired-ends) prepared from Link\_ADI and enrichment eCI DNA using TruSeq PE Cluster Kit v3-cBot-HS sequencing kit (Illumina Inc.). In addition, libraries prepared from Link\_ADI DNA were shotgun sequenced using the PacBio RS II Single Molecule, Real-Time (SMRT®) DNA Sequencing System. Library. The library was prepared using the PacBio 2 kb library preparation

protocol and sequenced on 8 SMRT cells using P4-C2 chemistry. To increase PacBio read quality, raw reads were filtered using RS\_ReadsOfInsert.1 pipeline of SMRT Analysis (smrtanalysis\_2.3.0.140936.p4.150482) with minimum predicted accuracy of 99.

**SSU rRNA gene amplicon analysis.** Paired end reads were joined using the QIIME v1.8.0 toolkit included python script join\_paired\_ends.py (with the default method fastq-join) and quality filtered (at Phred  $\geq$  Q20) before proceeding with downstream analysis<sup>25</sup>. USEARCH61 was used for detection of chimeric sequences followed by clustering (at 97% sequence similarity) of non-chimera sequences and denovo picking of OTUs<sup>26,27</sup>. Joined reads were assigned to OTUs using the QIIME v1.8.0 toolkit<sup>25</sup>, where uclust<sup>28</sup> was applied to search sequences against a subset of the Greengenes database<sup>29</sup> filtered at 97% identity. Sequences were assigned to OTUs based on their best hit to the Greengenes database, with a cut-off at 97% sequence identity. Taxonomy was assigned to each sequence by accepting the Greengenes taxonomy string of the best matching Greengenes sequence. filter\_otus\_from\_otu\_table.py (included with QIIME) was used to filter out OTUs making up less than 0.005% of the total using default parameters and `-min_count_fraction` set to 0.00005 as previously reported<sup>30</sup>.

**Raw data assembly.** Due to the species complexity of sample Link\_ADI (~480 phylotypes) and the computational resources available to this study, HiSeq data was assembled using SOAPdenovo-63mer (SOAPdenovo2 <http://soap.genomics.org.cn/soapdenovo.html>) using the following parameters: `-K 51 -p 40 setting max_rd_len = 125, avg_ins = 100, reverse_seq = 0, and asm_flags = 1`. For all other HiSeq assemblies with samples exhibiting lower species complexity (enrichment eCI and phylotype-specific co-assembly, see below), IDBA\_UD was used. Sequence data from enrichment eCI was trimmed using sickle pe (version 0.940 <https://github.com/najoshi/sickle>) with default parameters, converted to an interleaved FASTA using the program fq2fa (bundled with IDBA\_UD) with the parameters `-merge -filter`, and assembled with IDBA\_UD v1.1.1 ([http://i.cs.hku.hk/~alse/hkubrg/projects/idba\\_ud/index.html](http://i.cs.hku.hk/~alse/hkubrg/projects/idba_ud/index.html)) using the parameters `-pre_correction -num_threads 15 -maxk 60`. PacBio reads for Link\_ADI were filtered using the SMRT portal, with only those CCS reads that produced a minimum accuracy of 0.99 (average 10 passes) being considered for further analysis (ranging from one to three kb in length). PacBio CCS reads were assembled using slightly modified parameters in MIRA 4.0 (<http://sourceforge.net/p/mira-assembler/wiki/Home/>): `COMMON_SETTINGS -DI:trt = ./ -NW:cmrl = warn \ PCBIOHQ_SETTINGS -CL:pec = yes`.

**Identification of marker genes in contigs.** For the identification of protein coding marker genes, open reading frame calling was first performed using MetaGeneMark<sup>31</sup> version 1 metagenome ORF calling model (`gmhmp -m MetaGeneMark_v1.mod -f G -a -d`). Output was subsequently converted into a multiple FASTA using the included `aa_from_gff.pl` script. The resulting proteins sequences were compared against the 31 AMPHORA marker gene HMMs using HMMSCAN (part of HMMER version 3.0<sup>32</sup>), that form the basis of an automated phylogenomic inference pipeline for bacterial sequences<sup>18</sup>. The marker genes used are: *dnaG*, *frz*, *infC*, *nusA*, *pgk*, *pyrG*, *rplA*, *rplB*, *rplC*, *rplD*, *rplE*, *rplF*, *rplK*, *rplL*, *rplM*, *rplN*, *rplP*, *rplS*, *rplT*, *rpmA*, *rpoB*, *rpsB*, *rpsC*, *rpsE*, *rpsI*, *rpsJ*, *rpsK*, *rpsM*, *rpsS*, *smpB* and *tsf*. Matches with  $e$ -values of  $< 1.e^{-5}$  were considered legitimate. SSU rDNA searches were conducted using BLASTN (`-e 1e-20 -r 1 -q -1 -v 5 -b 5 -F F`) against a database of phylogenetically diverse representative sequences from sequenced genomes<sup>33</sup>.

The completeness of the unClos\_1 and unFirm\_1 HiSeq and PacBio contig sets was calculated using the identification of 107 CSCGs<sup>19</sup>. Using a provided HMM file<sup>19</sup>, the protein lists from each contig set were compared to the essential protein list using default hmmscan parameters with the `-tblout` flag. To determine if a match is legitimate, we relied on the trusted cutoff values for each protein's HMM (provided in the file) and excluded all matches with scores below the trusted cutoff for each essential protein.

**Genomic binning.** The GC % was calculated for each contig and the coverage values for each were provided by each assembler (IDBA\_UD provides a single coverage value, MIRA provides average coverage, and SOAPdenovo provides k-mer coverage). From this, we created a table of GC % versus coverage for each contig, allowing us to visualize clustering of contigs. Using contig clustering and marker gene analysis of our PacBio contigs (because they are on average longer and contain greater marker gene representation including SSU rDNA fragments), we were able to generate phylotype-specific training data for the two most abundant organisms (unClos\_1 and unFirm\_1). The taxonomic congruence of marker genes in the training data contigs was verified by BLASTP<sup>34</sup> and CLUSTALX alignments. Each marker gene was compared against NCBI nr database using default parameters except for a more stringent  $e$ -value cut off of  $1e^{-5}$ . Representative sequences of the best matches, corresponding to sequences from *Mageeibacillus indolicus* (NC\_013895), *Ruminoclostridium thermocellum* (NC\_009012), *Caldicellulosiruptor kristjanssonii* (NC\_014721), and *Peptoniphilus sp.* 1-1 (NZ\_LM997412) for unClos\_1 and *Dethiobacter alkaliphilus* (NZ\_ACJM01000000), *Desulfotomaculum kuznetsovii* (NC\_015573), *Ruminoclostridium thermocellum* (NC\_009012), and *Mahella australiensis* (NC\_015520) for unFirm\_1. We also included sequences from *Escherichia coli* (NC\_000913) and *Thermotoga maritima* (NC\_000853) as outliers. For both sets of marker genes, each sequence was trimmed by the BLASTP alignment with representatives from the above mentioned organisms to relevant amino acids and compared using CLUSTALX 2.1<sup>35</sup> with default parameters except the BLOSUM series protein weight matrix was used instead of Gonnet. Bootstrapped, mid-point rooted, neighbor-joining trees were generated from the alignments using the random number generator seed of 111 and 1000 trials. The resulting phylip trees were visualized in phenogram form using the Phylodendron website ([iubio.bio.indiana.edu/treeapp/treeprint-form.html](http://iubio.bio.indiana.edu/treeapp/treeprint-form.html)).

These subsets consisted of contigs totaling more than 100 kb, the minimum necessary for custom binning using PhyloPythiaS+<sup>15</sup>. Contigs that met the criteria for phylotype-specific training data were larger than 7 kb, exhibited consist coverage ( $\pm 2x$ ) and GC % ( $\pm 3\%$ ) values and encoded a SSU rRNA gene or marker gene



that demonstrated phylogenomic grouping with the representative OTU sequence identified via 16S rRNA gene amplicon analysis. Binning was performed using PhyloPythiaS+ using both default settings, against a database consisting of all publically available prokaryotic genomes in NCBI, and with our custom training data.

**Co- and hybrid assembly.** Various merged assemblies were performed in an attempt to improve assembly statistics of the Link\_ADI community metagenome and the genome reconstructions of dominate phylotypes (unClos\_1 and unFirm\_1). Hybrid assemblies of whole community contigs (>1 kb) from both the HiSeq and PacBio CCS contig subsets were performed using CAP3<sup>12</sup> (version date 12/21/07) with default parameters except a minimum overlap percent identity (-p) of 0.95.

In order to reconstruct as large as possible genomes for unClos\_1 and unFirm\_1, we performed hybrid assemblies of binned contigs for each phylotype from all of our samples including the PacBio and HiSeq data from Link\_ADI and the HiSeq data from enrichment eCI. This was carried out in two stages. The first stage consisted of mapping HiSeq reads to their corresponding phylotype contigs using BWA mem<sup>36</sup> (version 0.7.8-r455) with default parameters. The reads that mapped from each sample (Link\_ADI and eCI) were identified by parsing the resulting SAM files, pooled together for each phylotype, and co-assembled with IDBA\_UD using the same workflow as eCI above into cross-sample HiSeq contigs. The second stage consisted of pooling together the cross-sample HiSeq contigs with the phylotype-specific PacBio contigs, which were hybrid assembled using CAP3, with the same parameters as above. The unincorporated contigs from the hybrid assemblies (contigs that went into the assembly but were not incorporated into hybrid contigs) were also included in the final reconstructed genomes used in this study.

A hybrid assembly of raw sequences between both platforms was also performed using MIRA 4.0. The cross-sample HiSeq reads used above in each co-assembly were used as input along with PacBio reads that mapped to each species-specific bin (identified through the MIRA supplied CAF result file). MIRA 4.0 was run using the following parameters: COMMON\_SETTINGS -SK:mmhr = 1 -NW:cac = warn -NW:cdrn = no -NW:cmrl = warn \ PCBIOHQ\_SETTINGS -CL:pec = yes \ SOLEXA\_SETTINGS -CL:pec = yes. For the HiSeq readgroup, the following information was supplied: template\_size = 100 400 and segmet\_naming = solexa.

## References

- Hess, M. *et al.* Metagenomic discovery of biomass-degrading genes and genomes from cow rumen. *Science* **331**, 463–467 (2011).
- Albertsen, M. *et al.* Genome sequences of rare, uncultured bacteria obtained by differential coverage binning of multiple metagenomes. *Nat. Biotechnol.* **31**, 533–538 (2013).
- Brady, A. & Salzberg, S. L. Phymm and PhymmBL: metagenomic phylogenetic classification with interpolated Markov models. *Nat. Meth.* **6**, 673–676 (2009).
- Scholz, M., Lo, C. C. & Chain, P. S. Improved assemblies using a source-agnostic pipeline for MetaGenomic Assembly by Merging (MeGAMerge) of contigs. *Sci Rep.* **4**, e6480 (2014).
- Lee, H. *et al.* Error correction and assembly complexity of single molecule sequencing reads. *bioRxiv*. 10.1101/006395 (2014).
- English, A. C. *et al.* Mind the gap: Upgrading genomes with pacific biosciences RS long-read sequencing technology. *PLoS One* **11**, e47768 (2012).
- Koren, S. *et al.* Hybrid error correction and *de novo* assembly of single-molecule sequencing reads. *Nat. Biotechnol.* **30**, 693–700 (2012).
- Eid, J. *et al.* Real-time DNA sequencing from single polymerase molecules. *Science* **323**, 133–138 (2009).
- Travers, K. J., Chin, C.-S., Rank, D. R., Eid, J. S. & Turner, S. W. A flexible and efficient template format for circular consensus sequencing and SNP detection. *Nucleic Acids Res.* **38**, e159 (2010).
- Chevreaux, B., Wetter, T. & Suhai, S. Genome Sequence Assembly Using Trace Signals and Additional Sequence Information. *Computer Science and Biology: Proceedings of the German Conference on Bioinformatics (GCB)* **99**, 45–46 (1999).
- Luo, R. *et al.* SOAPdenovo2: an empirically improved memory-efficient short-read *de novo* assembler. *Gigascience* **1**, 18 (2012).
- Huang, X. & Madan, A. CAP3: A DNA sequence assembly program. *Genome Res.* **9**, 868–877 (1999).
- Imelfort, M. *et al.* GroopM: an automated tool for the recovery of population genomes from related metagenomes. *PeerJ*. **2**, e603 (2014).
- Alneberg, J. *et al.* Binning metagenomic contigs by coverage and composition. *Nat. Methods*. **11**, 1144–1146 (2014).
- Gregor, I., Dröge, J., Schirmer, M., Quince, C. & McHardy, A. C. PhyloPythiaS+: A self-training method for the rapid reconstruction of low-ranking taxonomic bins from metagenomes. *arXiv.org q-bio.QM*, arXiv:1406.7123 (2014).
- Pope, P. B. *et al.* Adaptation to herbivory by the Tamar wallaby includes bacterial and glycoside hydrolase profiles different to other herbivores. *Proc. Natl Acad. Sci. USA* **107**, 14793–14798 (2010).
- Patil, K. R. *et al.* Taxonomic metagenome sequence assignment with structured output models. *Nat. Meth.* **8**, 191–192 (2011).
- Wu, M. & Eisen, J. A. A simple, fast, and accurate method of phylogenomic inference. *Genome Biol.* **9**, R151 (2008).
- Dupont, C. L. *et al.* Genomic insights to SAR86, an abundant and uncultivated marine bacterial lineage. *ISME J.* **6**, 1186–1199 (2012).
- Peng, Y., Leung, H. C., Yiu, S. M. & Chin, F. Y. IDBA-UD: a *de novo* assembler for single-cell and metagenomic sequencing data with highly uneven depth. *Bioinformatics* **28**, 1420–1428 (2012).
- Koren, S. & Phillippy, A. M. One chromosome, one contig: complete microbial genomes from long-read sequencing and assembly. *Curr. Opin. Microbiol.* **23**, 110–120 (2015).
- Sun, L., Müller, B., Westerholm, M. & Schnürer, A. Syntrophic acetate oxidation in industrial CSTR biogas digesters. *J. Biotechnol.* **171**, 39–44 (2014).
- Sun, L. Biogas production from lignocellulosic materials. *Degree of Doctor in Philosophy* (Swedish University of Agricultural Sciences, Uppsala) (2015).
- Klindworth, A. *et al.* Evaluation of general 16S ribosomal RNA gene PCR primers for classical and next-generation sequencing-based diversity studies. *Nucleic acids research* **41**, e1 (2013).
- Caporaso, J. G. *et al.* QIIME allows analysis of high-throughput community sequencing data. *Nat Methods* **7**, 335–336 (2010).
- Edgar, R. C. Search and clustering orders of magnitude faster than BLAST. *Bioinformatics* **26**, 2460–2461 (2010).
- Edgar, R. C., Haas, B. J., Clemente, J. C., Quince, C. & Knight, R. UCHIME improves sensitivity and speed of chimera detection. *Bioinformatics* **27**, 2194–2200 (2011).
- Edgar, R. C. Search and clustering orders of magnitude faster than BLAST. *Bioinformatics* **26**, 2460–2461 (2010).
- DeSantis, T. Z. *et al.* Greengenes, a chimera-checked 16S rRNA gene database and workbench compatible with ARB. *Appl Environ Microbiol* **72**, 5069–5072 (2006).

30. Bokulich, N. A. *et al.* Quality-filtering vastly improves diversity estimates from Illumina amplicon sequencing. *Nat. Meth.* **10**, 57–59 (2013).
31. Zhu, W., Lomsadze, A. & Borodovsky, M. *Ab initio* gene identification in metagenomic sequences. *Nucleic Acids Res.* **38**, e132 (2010).
32. Finn, R. D., Clements, J. & Eddy, S. R. HMMER web server: interactive sequence similarity searching. *Nucleic Acids Res.* **39**, W29–W37 (2011).
33. Frank, J. A. *et al.* Critical evaluation of two primers commonly used for amplification of bacterial 16S rRNA genes. *Appl. Environ. Microb.* **74**, 2461–2470 (2008).
34. Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local alignment search tool. *J. Mol. Biol.* **215**, 403–410 (1990).
35. Larkin, M. A. *et al.* Clustal W and Clustal X version 2.0. *Bioinformatics* **23**, 2947–2948 (2007).
36. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler Transform. *Bioinformatics* **25**, 1754–1760 (2009).

## Acknowledgements

J.A.F. and P.B.P. are supported by a grant from the European Research Council (336355-MicroDE). The sequencing service was provided by the Norwegian Sequencing Centre ([www.sequencing.uio.no](http://www.sequencing.uio.no)), a national technology platform hosted by the University of Oslo and supported by the “Functional Genomics” and “Infrastructure” programs of the Research Council of Norway and the Southeastern Regional Health Authorities. DNA preparations from Biogas reactor samples were supplied by Professor Anna Schnürer and Li Sun from the Department of Microbiology, Swedish University of Agricultural Science, Uppsala, Sweden. SSU rDNA analyses were performed by Live H. Hagen from Department of Chemistry, Biotechnology and Food Science, Norwegian University of Life Sciences, Ås, Norway. We thank Professor Abigail A. Salyers from the University of Illinois for her helpful advice and correspondence.

## Author Contributions

P.B.P., A.J.N. and V.G.H.E. proposed this project. J.A.F., A.J.N., A.C.M. and P.B.P. designed the experiments and supervised the project. J.A.F., A.T.-K. and Y.P. did the experiments. J.A.F., A.T.-K., Y.P. and P.B.P. analyzed the data. J.A.F., A.J.N., V.G.H.E. and P.B.P. contributed to analysis of the results and paper writing.

## Additional Information

**Accession codes:** Datasets are available at the NCBI Sequence Read Archive under the BioProject PRJNA294734 with the accession numbers SRR212703, SRR2420276, and SRR2420280.

**Supplementary information** accompanies this paper at <http://www.nature.com/srep>

**Competing financial interests:** The authors declare no competing financial interests.

**How to cite this article:** Frank, J. A. *et al.* Improved metagenome assemblies and taxonomic binning using long-read circular consensus sequence data. *Sci. Rep.* **6**, 25373; doi: 10.1038/srep25373 (2016).



This work is licensed under a Creative Commons Attribution 4.0 International License. The images or other third party material in this article are included in the article’s Creative Commons license, unless indicated otherwise in the credit line; if the material is not included under the Creative Commons license, users will need to obtain permission from the license holder to reproduce the material. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>