

RESEARCH

Open Access



Choice of assembly software has a critical impact on virome characterisation

Thomas D. S. Sutton^{1,2†} , Adam G. Clooney^{1,2†}, Feargal J. Ryan^{1,2,3†}, R. Paul Ross^{1,2,4} and Colin Hill^{1,2*}

Abstract

Background: The viral component of microbial communities plays a vital role in driving bacterial diversity, facilitating nutrient turnover and shaping community composition. Despite their importance, the vast majority of viral sequences are poorly annotated and share little or no homology to reference databases. As a result, investigation of the viral metagenome (virome) relies heavily on de novo assembly of short sequencing reads to recover compositional and functional information. Metagenomic assembly is particularly challenging for virome data, often resulting in fragmented assemblies and poor recovery of viral community members. Despite the essential role of assembly in virome analysis and difficulties posed by these data, current assembly comparisons have been limited to subsections of virome studies or bacterial datasets.

Design: This study presents the most comprehensive virome assembly comparison to date, featuring 16 metagenomic assembly approaches which have featured in human virome studies. Assemblers were assessed using four independent virome datasets, namely, simulated reads, two mock communities, viromes spiked with a known phage and human gut viromes.

Results: Assembly performance varied significantly across all test datasets, with SPAdes (meta) performing consistently well. Performance of MIRA and VICUNA varied, highlighting the importance of using a range of datasets when comparing assembly programs. It was also found that while some assemblers addressed the challenges of virome data better than others, all assemblers had limitations. Low read coverage and genomic repeats resulted in assemblies with poor genome recovery, high degrees of fragmentation and low-accuracy contigs across all assemblers. These limitations must be considered when setting thresholds for downstream analysis and when drawing conclusions from virome data.

Keywords: Virome, Viral, Assembly, Metagenome, Benchmark, Comparison, Bacteriophage, Phage

Background

The rapid evolution of metagenomics and high throughput sequencing technologies has revolutionised the study of microbial communities, giving new insights into the role and identity of the uncultivated microbes which account for the majority of metagenomic sequences [52]. However, the majority of microbial sequencing efforts have focused on the characterisation of prokaryotic microbes. Viral metagenomes (viromes) are dominated by novel sequences, often with up to 90% of sequences sharing little to no homology to reference databases [2].

Bacteriophage, the most abundant member of viral communities, play a key role in the shaping the composition of microbial communities and facilitate horizontal gene transfer [42]. Viromes have been shown to play a role in global geochemical cycles [7] and have been studied in varied ecosystems including the ocean [21]. Viromes of the human body are of particular interest, where they have been linked to disease status [39], maintaining human health [31] and shaping the gut microbiome in early life [26, 34]. Due to the predominance of uncharacterised viral sequences “viral dark matter” [45] and the lack of a universal marker gene, virome studies rely on database-independent analysis methods and depend heavily on de novo assembly to resolve viral genomes from metagenomic sequencing reads.

* Correspondence: c.hill@ucc.ie

[†]Thomas D. S. Sutton, Adam G. Clooney and Feargal J. Ryan contributed equally to this work.

¹APC Microbiome Ireland, Cork, Ireland

²School for Microbiology, University College Cork, Cork, Ireland

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



Metagenomic assemblers typically use de Bruijn graph (DBG) approaches to address the complexity and size of metagenomic datasets in an accurate and efficient manner. Microbial metagenomes pose significant challenges to DBG assembly when compared to single genome assemblies often complicating the DBG and leading to fragmentation and/or misassembly [41]. These challenges include uneven sequencing coverage of organisms within the metagenome, the presence of conserved regions across different species, repeat regions within genomes and the introduction of false k -mers by both closely related genomes at differing abundances and sequencing errors at high read coverage. This hampers the use of coverage statistics to resolve repeat regions between and within genomes [41].

A wide array of metagenomic assembly programs have been employed, each addressing aspects of metagenomic challenges to varying degrees. However, many of these programs have been designed and optimised for bacterial metagenomes, which share many assembly challenges of viromes but to a lesser degree. Virome data is characterised by high proportions of repeat regions within viral genomes [36], hypervariable genomic regions associated with host interaction [55] and high mutation rates which lead to increased metagenomic complexity and strain variation [44]. Low DNA yields also limit read coverage and often require a multiple displacement amplification (MDA) step which has been shown to preferentially amplify small single-stranded DNA viruses [23]. Extremes in read coverage caused by MDA bias and dominant viral taxa such as crAssphage, which can make up large proportions of human gut viromes [10], sequester sequencing resources and result in insufficient coverage of low abundance viruses. These challenges result in fragmented virome assemblies [14], limiting their use in downstream analysis. Despite benchmarks of bacterial metagenomes having highlighted failings and benefits of particular assembly programs, many poorly performing assemblers have featured in virome studies [12, 17, 19].

Accurate comparison of metagenomic assemblers is complicated by the unknown composition of metagenomic datasets and the limited applicability of general assembly statistics such as N50 [9, 54]. To address this, the accuracy and efficacy of metagenomic assembly programs are often evaluated using simulated datasets and mock communities of known composition. Although these simulated datasets are undergoing constant improvements [13, 48], they have focused primarily on bacterial metagenomes and remain limited in their ability to accurately replicate the challenges of true metagenomes. While some virome-specific assembly benchmarks have been performed, many have been limited to a small number of assemblers, 454 pyrosequencing data or subsections of virome studies which have exclusively used simulated data [3, 14, 20, 44, 51, 53].

Here, we expand upon previous studies and present a detailed investigation of assembly software for virome analysis which compares all those previously used in human virome studies to date, as well as other popular or more recently published assemblers (Table 1). We compare assembly efficacy and accuracy using simulated viromes, mock viral communities and human gut viromes spiked with a known exogenous bacteriophage (Additional file 1; Additional file 2; Additional file 3; Additional file 4). Furthermore, we confirm these findings using human virome data from published datasets and assess computational parameters such as runtime and RAM usage. We also investigate in detail the impact of sequencing coverage and genomic repeats on assembly performance and highlight important considerations for future virome studies. Together, these data comprise most comprehensive virome assembly benchmark to date.

Results

Simulated virome dataset

Normalised genome abundance of 572 members of a published simulated community, (Fig. 1a) [20] the proportion of genome recovered and the degree of fragmentation were assessed by aligning the resulting contigs from each assembler to the reference genomes (Fig. 1b). MetaVelvet was not included in this analysis as it failed to reach completion after 7 days. Approximately half of the genomes in the community featured an average recovered genome fraction less than 75% and exhibited higher degrees of fragmentation (> 10 contigs per genome on average) across all assemblers. For 87 of the 572 genomes, there was an average recovered genome fraction of less than 20% across all assemblers (the low recovered genome fraction of VICUNA was excluded as an outlier). Of these genomes, 84 were present at low abundance (lowest 40% of all abundances normalised to genome length). The remaining three genomes were present at higher normalised abundances (50–80th percentile) but featured some of the highest proportions of genomic repeats (70th–90th percentile).

Normalised genome abundance within the community had a strong positive correlation with recovered genome fraction across all assemblers (Additional file 5: Table S1) and was verified using a linear model (Additional file 5: Table S2), with the exception of SOAPdenovo2, which was negative. Normalised abundance also correlated negatively with the degree of fragmentation (number of contigs) across all assemblers except Velvet which was positively correlated and Geneious which was not significant (Additional file 5: Table S1). None of the genomes in the lower 30th percentile of normalised abundance featured an average recovered genome fraction greater than 75%, further exemplifying the impact of low sequencing coverage. However, high abundance did not

Table 1 A list of assemblers used in this study

	Link	Version used	Reference
ABYSS	http://www.bcgsc.ca/downloads/abyss/	v2.0.2	[50]
CLC	https://www.qiagenbioinformatics.com/products/clc-assembly-cell/	v5.0.5	https://www.qiagenbioinformatics.com/
Geneious	https://www.geneious.com/features/assembly-mapping/	v11.0.3	[22]
IDBA UD	https://i.cs.hku.hk/~alse/hkubrg/projects/idba_ud	v1.1.1	[43]
MEGAHIT	https://github.com/voutcn/megahit	v1.1.1-2	[25]
MetaVelvet	https://metavelvet.dna.bio.keio.ac.jp/	v1.2.02	[38]
MIRA	http://www.chevreux.org/mira_downloads.html	v4.0.2	[14]
Ray Meta	http://denovoassembler.sourceforge.net/	v2.3.0	[5]
SOAPdenovo2	http://soap.genomics.org.cn/soapdenovo.html	v2.04	[29]
SPAdes	http://cab.spbu.ru/software/spades/	v3.10.0	[4]
SPAdes meta	http://cab.spbu.ru/software/spades/ (variation of SPAdes applied with flag)	v3.10.0	[40]
Velvet	https://www.ebi.ac.uk/~zerbino/velvet/	v1.2.10	[58]
VICUNA	https://github.com/broadinstitute/mvicuna	v1.3	[53]

consistently improve genome recovery, and of the 172 genomes in the top 30% of normalised abundance, 20 featured an average genome fraction below 50%. The distance of the log-transformed (due to extremes in values) normalised abundances from the mean was negatively correlated with recovered genome fraction across all assemblers (correlation coefficient -0.42 , p value $< 2.2e-16$). Of 171 genomes in the 40th–60th percentile of normalised abundance, 29 featured an average genome fraction below 50%. This indicates factors other than abundance may hamper genome recovery. MIRA and Geneious both recovered a greater fraction of low abundance genomes with fewer contigs than other assemblers. However, MIRA assemblies of 13 of the most abundant genomes in the community (highest 10%) exhibited the highest degree of fragmentation in the study, generating between 401 and 2983 contigs per genome.

The proportion of inverted repeats, palindromic repeats, tandem repeats and a total proportion of genomic repeats was calculated for each genome. The total percentage of repeat regions predicted in each genome was positively correlated with the degree of fragmentation observed in each assembly across all assemblers with the exception of Ray Meta (Additional file 5: Table S3) and negatively correlated with recovered genome fraction across all assemblers except ABySS (k -mer 63/127), Geneious and SOAPdenovo2. When this relationship between repeat regions and the recovered genome fraction was assessed using a linear model, correlations were significant for CLC, MIRA, Ray Meta, Velvet and all parameters of SPAdes (Additional file 5: Table S2). Both the proportion of repeat regions in a genome and the relative abundance of that genome contribute to the variation in recovered genome fraction, though each

explains a separate aspect of this variation. No interaction was found between these two metrics.

VICUNA, Ray Meta, SOAPdenovo2, Geneious, ABySS (both k -mer sizes) and Velvet recovered under 50% of the total genome fraction (all genomes in the community). VICUNA produced just four contigs in total with high levels of mismatches (174 per 100 kb on average) which could possibly be linked to the format of the artificial reads as this was not observed in real sequencing data. The five assemblers which recovered the highest genome fraction overall were SPAdes (default), MEGAHIT, SPAdes (single cell), SPAdes (single cell + careful) and CLC. All assemblers achieving a minimum average genome fraction of 50% were subjected to a ranking system (Additional file 5: Table S4). To compare both recovery and fragmentation, assemblers were ordered from best to worst based on genome recovery and number of aligned contigs. The average rank resulted in SPAdes (default) performing best, recovering 72.2% overall genome sequences with 8230 contigs. The remaining top five assemblers of this combined rank were SPAdes (meta) 68.2% with 7419 contigs, SPAdes (single cell) 68.9% with 9506 contigs, CLC 68.6% with 9152 contigs and MEGAHIT 69.6% with 10,083 contigs. The number of assemblies which recovered greater than 90% of the target genome in one single contig was compared (Fig. 2). SPAdes (default) performed best, recovering 210, and SPAdes (meta), SPAdes (single cell + careful), CLC and SPAdes (single cell) each recovered 179, 168, 162 and 160 genomes, respectively.

The accuracy of assemblies was assessed by calculating the average count of indels, mismatches and misassemblies per 100 kb across all genomes. These counts were normalised to the number of genomes each assembler recovered with a minimum genome fraction of 50%.

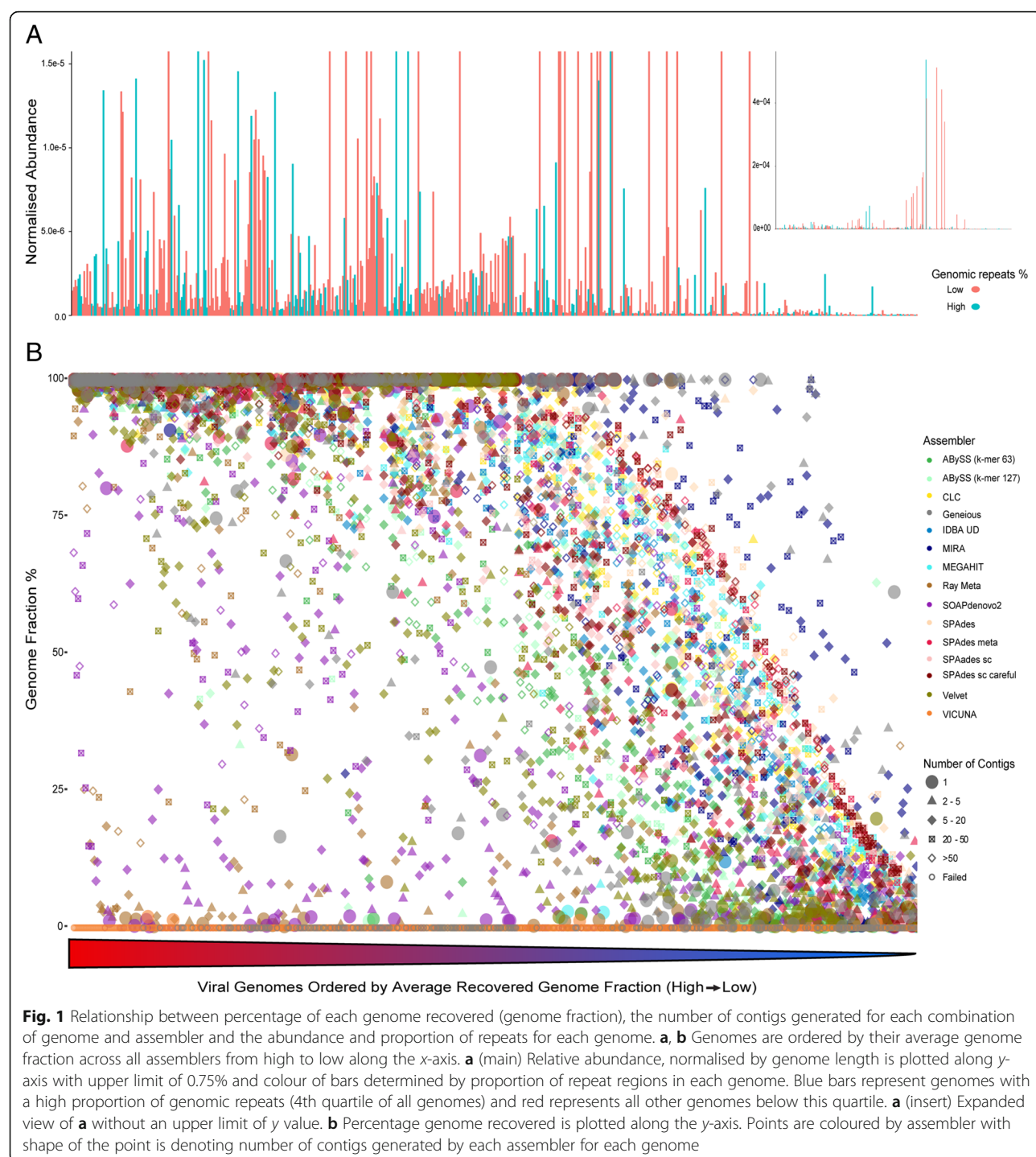
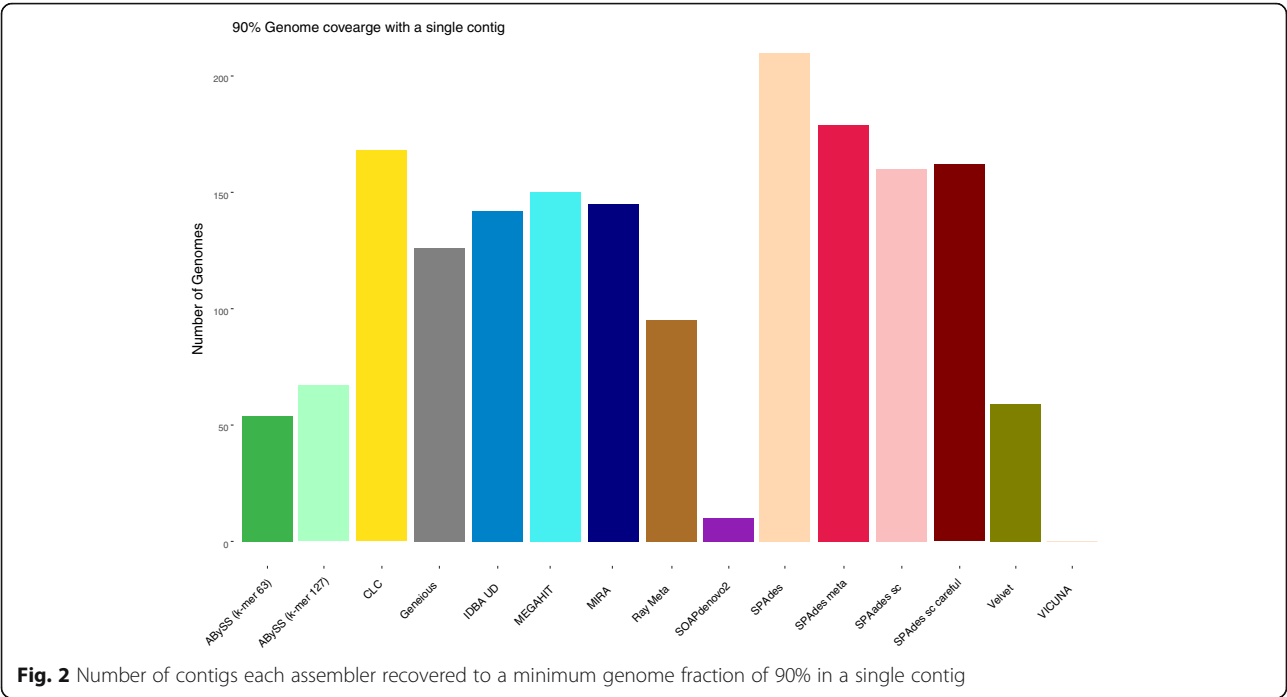


Fig. 1 Relationship between percentage of each genome recovered (genome fraction), the number of contigs generated for each combination of genome and assembler and the abundance and proportion of repeats for each genome. **a, b** Genomes are ordered by their average genome fraction across all assemblers from high to low along the x-axis. **a** (main) Relative abundance, normalised by genome length is plotted along y-axis with upper limit of 0.75% and colour of bars determined by proportion of repeat regions in each genome. Blue bars represent genomes with a high proportion of genomic repeats (4th quartile of all genomes) and red represents all other genomes below this quartile. **a** (insert) Expanded view of **a** without an upper limit of y value. **b** Percentage genome recovered is plotted along the y-axis. Points are coloured by assembler with shape of the point is denoting number of contigs generated by each assembler for each genome

These were ranked according to their performance in all three metrics (Additional file 5: Table S4), with assemblies from Velvet having the lowest overall counts followed by ABYSS, IDBA UD, MEGAHIT and Ray Meta. With the exception of Ray Meta and SOAPdenovo2, the number of mismatches per 100 kb was negatively correlated with both genome abundance and recovered genome fraction across all assemblers (Additional file 5: Table S1).

The rate of false positive (no alignment to reference genomes) and false negative (recovered genome fraction of 0%) contigs assembled allowed for the determination of sensitivity. A number of assemblers had a sensitivity greater than 97%; however, each returned greater than 7000 contigs, inferring a high degree of fragmentation (Table 2). MIRA assembled (partial or complete) 559 of the genomes with a false positive count of just four.



However, this was achieved from more than 27,000 contigs. ABySS (both *k*-mer sizes), Geneious, Ray Meta and Velvet returned very few false positives but failed to detect many of the genomes present. SPAdes (meta) performed best with 558 of the 572 genomes detected and only five false positives resulting from 7419 contigs.

Mock community dataset

Two mock viral communities were used to investigate the impact of high and low abundance ssDNA viruses on assembly performance. Mock A (Table 3a) contained 12 viral genomes, 10 of which were at equal abundance (9.82% of the total community) and 2 ssDNA genomes

Table 2 The number of false positive and false negative contigs generated by each assembler for the simulated community, together with the sensitivity rates

	False positives	False negative	True positives	No. of contigs returned ^a	Sensitivity
ABSS (k-mer 63)	0	111	461	7957	80.59
ABySS (k-mer 127)	1	123	449	7732	78.50
CLC	34	5	567	9152	99.13
Geneious	9	190	382	958	66.78
IDBA UD	25	9	563	8999	98.43
MEGAHIT	21	8	564	10,083	98.60
MetaVelvet	N/A	N/A	N/A	N/A	N/A
MIRA	4	13	559	27,600	97.73
Ray Meta	0	213	359	4224	62.76
SOAPdenovo2	536	116	456	11,548	79.72
SPAdes	29	3	569	8230	99.48
SPAdes meta	5	14	558	7419	97.55
SPAdes sc	38	7	565	9506	98.78
SPAdes sc careful	40	6	566	9724	98.95
Velvet	1	65	507	6343	88.64
VICUNA	0	558	14	4	2.45

^a572 in community

Table 3 The number of false positive and false negative contigs generated by each assembler for (a) mock community A and (b) mock community B along with the sensitivity rates for each

	False positives	False negative	True positive	No. of contigs returned ^a	Sensitivity
A					
ABYSS (<i>k</i> -mer 63)	52	4	8	61	66.67
ABYSS (<i>k</i> -mer 127)	50	6	6	56	50.00
CLC	1143	0	12	1299	100.00
Geneious	53	0	12	65	100.00
IDBA UD	0	0	12	12	100.00
MEGAHIT	0	0	12	13	100.00
MetaVelvet	0	3	9	26	75.00
MIRA	0	0	12	89	100.00
Ray Meta	0	0	12	12	100.00
SOAPdenovo2	2	0	12	23	100.00
SPAdes	0	0	12	14	100.00
SPAdes meta	0	0	12	14	100.00
SPAdes sc	1513	0	12	1527	100.00
SPAdes sc careful	0	0	12	15	100.00
Velvet	0	3	9	26	75.00
VICUNA	4969	0	12	5385	100.00
B					
ABYSS (<i>k</i> -mer 63)	60	4	8	69	66.67
ABYSS (<i>k</i> -mer 127)	132	6	6	139	50.00
CLC	450	0	12	505	100.00
Geneious	14	0	12	30	100.00
IDBA UD	0	0	12	12	100.00
MEGAHIT	0	0	12	14	100.00
MetaVelvet	0	1	11	24	91.67
MIRA	94	1	11	157	91.67
Ray Meta	0	0	12	13	100.00
SOAPdenovo2	2	2	10	27	83.33
SPAdes	0	0	12	13	100.00
SPAdes meta	0	0	12	14	100.00
SPAdes sc	593	0	12	607	100.00
SPAdes sc careful	0	0	12	14	100.00
Velvet	0	1	11	24	91.67
VICUNA	0	0	12	15	100.00

^a12 in community

(NC_001330 and NC_001422) at low abundance (0.92%). Analysis of this community showed that although some assemblers, namely CLC, Geneious, SPAdes (single cell) and VICUNA, detected all 12 genomes, this was at the expense of a large number of false positives (1143, 53, 1513 and 4969, respectively). Velvet and MetaVelvet generated no false positives but failed to assemble three genomes, while ABYSS (for both *k*-mers) generated a large number of false positives and failed to assemble four and six

genomes, respectively. IDBA UD and Ray Meta outperformed the other assemblers with an equal number of contigs to genomes (12), followed by MEGAHIT, SPAdes (default) and SPAdes (meta) with 13, 14 and 14. Mock B (Table 3b) also contained 12 genomes but with a higher abundance of ssDNA genomes NC_001330 and NC_001422 (32.47%). VICUNA assemblies of mock B improved upon those from mock A as no false positives were generated, while the false positive rate in the MIRA

assembler increased from none to 94 in mock A. IDBA UD performed best followed by SPAdes (default), Ray Meta, MEGAHIT and SPAdes (meta) based on sensitivity and number of contigs, while ABySS (both k -mer sizes) and SOAPdenovo2 had the lowest sensitivity. Despite being a relatively simple community consisting of 12 members, not all assemblers were able to recover all members (Additional file 5: Tables S5-S6). A greater number of assemblers (six) failed to assemble all members of mock B than mock A (four). ABySS(k -mer 63), ABySS(k -mer 127), Velvet and MetaVelvet failed to assemble 6, 4, 3 and 3 genomes, respectively, in mock A and 6, 4, 1 and 1 genomes, respectively, in mock B. In addition, MIRA and SOAPdenovo2 failed to assemble 1 and 2 genomes, respectively, in mock B.

All but three VICUNA assemblies in mock A exhibited a high level of fragmentation, generating 34.7 ± 35 (mean \pm standard deviation) contigs per genome. Fragmentation was also seen in MIRA assemblies to a lesser degree with 7.4 ± 10 contigs per genome on average. There was a high rate of fragmentation in CLC with one community member generating 144 contigs for genome KF302035. Average recovered genome fraction of $85.4 \pm 6.4\%$ was skewed by ABySS (k -mer 63), ABySS (k -mer 127), Velvet, MetaVelvet, SOAPdenovo2 and VICUNA which recovered on average 49.5%, 66.6%, 73.8%, 73.8%, 29.7% and 76.6%, respectively. All other assemblers recovered over 99% of each genome in the community (Additional file 6: Figure S1).

Closer inspection of the two ssDNA genomes present at lower relative abundance highlighted significant differences in the average number of indels across all assemblies of the NC_001330 and NC_001422 genomes versus other members of the community (p value = 0.037). These genomes exhibited an average of 41.7 ± 18.5 and 9.4 ± 20.4 indels per 100 kb, while all other genomes featured an average of 7.8 ± 18.9 indels per 100 kb. The low abundant ssDNA genomes NC_001330 and NC_001422 also featured the highest average mismatches per 100 kb at 148.7 ± 3 and 302.5 ± 10.7 , respectively (Additional file 6: Figure S1).

The degree of fragmentation observed by VICUNA and MIRA in mock B was lower than in mock A with a mean of 1.3 ± 0.89 and 5.3 ± 7.7 contigs per genome, respectively. CLC fragmented genome KF302035 in mock B (44 contigs), but to a lesser degree than mock A (144 contigs). MEGAHIT, which recovered at least 98% of all genomes in mock A, also recovered over 98% of all genomes in mock B except for the ssDNA genome NC_001422, of which 56.5% was recovered in two contigs. The majority of assemblies exhibited 147.9 ± 0 and 297 ± 1 mismatches per 100 kb for NC_001330 and NC_001422 (high abundance ssDNA), respectively, identical values to those measured in mock A. Velvet and MetaVelvet were exceptions with 184.2 and 860.2 for genome NC_001422 and

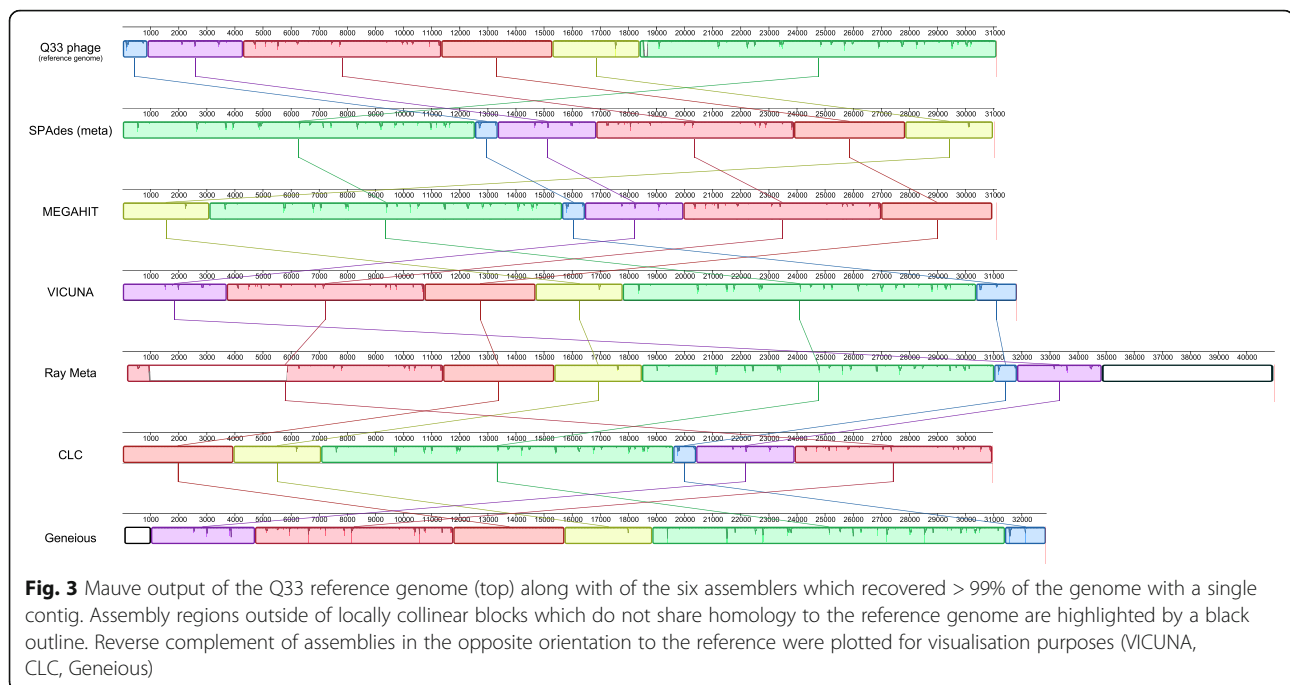
NC_001330. A similar pattern of high values across a narrow range was also observed with the number of indels, with 49.3 to 32.9 present in all assemblies NC_001330. Genome NC_001422 featured 18.57 indels across all SPAdes assemblies (all parameters) and 860.2 across both Velvet and MetaVelvet assemblies. All other assemblers which successfully recovered this genome did not feature any indels (Additional file 6: Figure S1).

Q33

Five assemblers failed to generate contigs which met alignment thresholds and were subsequently excluded from further analysis, namely ABySS (k -mer 63), ABySS (k -mer 127), SOAPdenovo2, Velvet and MetaVelvet. All remaining assemblers recovered over 90% of the spiked Q33 genome with the exception of MIRA (8.5%). Six assemblers recovered over 99% of the Q33 genome in a single contig—SPAdes (meta) 99.74%, MEGAHIT (99.6%), VICUNA (99.6%), Ray Meta (99.6%), CLC (99.5%) and Geneious (99.1) (Fig. 3). However, only MEGAHIT assembled the Q33 genome with a contig equal in length to the genome itself. SPAdes (meta) and CLC generated assemblies shorter than the reference genome by 86 and 141 bases. VICUNA (723), Geneious (1765) and Ray Meta (9884) each generated assemblies longer than the reference genome. SPAdes (default), SPAdes (single cell), IDBA UD and SPAdes (single cell + careful) each assembled Q33 in 2, 3, 4, 5 and 5 contigs, respectively. Ray Meta and VICUNA assemblies had the lowest number of mismatches and indels per 100 kb; however, Ray Meta exhibited the highest rate of misassemblies (two relocations, one inversion). All assemblers featured a minimum of one local misassembly with the exception of SPAdes (meta) which did not feature any. The six best assemblies of the Q33 genome and the genome itself are syntenic (although occasionally on the reverse strand), and the start and end point were not conserved (Fig. 3).

Read depth analysis (time and RAM)

Assemblers were compared for practicality by measuring the time to reach completion and maximum RAM usage via four published healthy human gut viromes [31] and various sequencing depths. It must be noted that all assembly tasks were allocated five threads; however, specifying the number of threads did not change the number of threads used by certain programs. MetaVelvet was not included in this analysis as it failed to reach completion after running for 7 days. CLC and Geneious were performed on a desktop computer and therefore excluded from time and RAM analysis. Runtime is dependent upon the number of reads, and this is largely linear in scale with more reads leading to an increased assembly time (Fig. 4a). MIRA and VICUNA (Fig. 4a insert) were



the slowest with MIRA taking over 15 times longer than the other software to assemble 3.5 million reads. SOAPdenovo2 had the shortest completion time followed by IDBA UD and Velvet. Most assemblers were consistent across samples (observed via error bars) with the exception of MIRA and Ray Meta. MIRA, VICUNA and Velvet (Fig. 4b insert) had the highest max RAM usage, while the lowest was Ray Meta, IDBA UD and SPAdes (meta) (Fig. 4b). The majority of assemblers observed a linear scale pattern similar to that of runtime.

Read depth analysis N50 and longest contig length

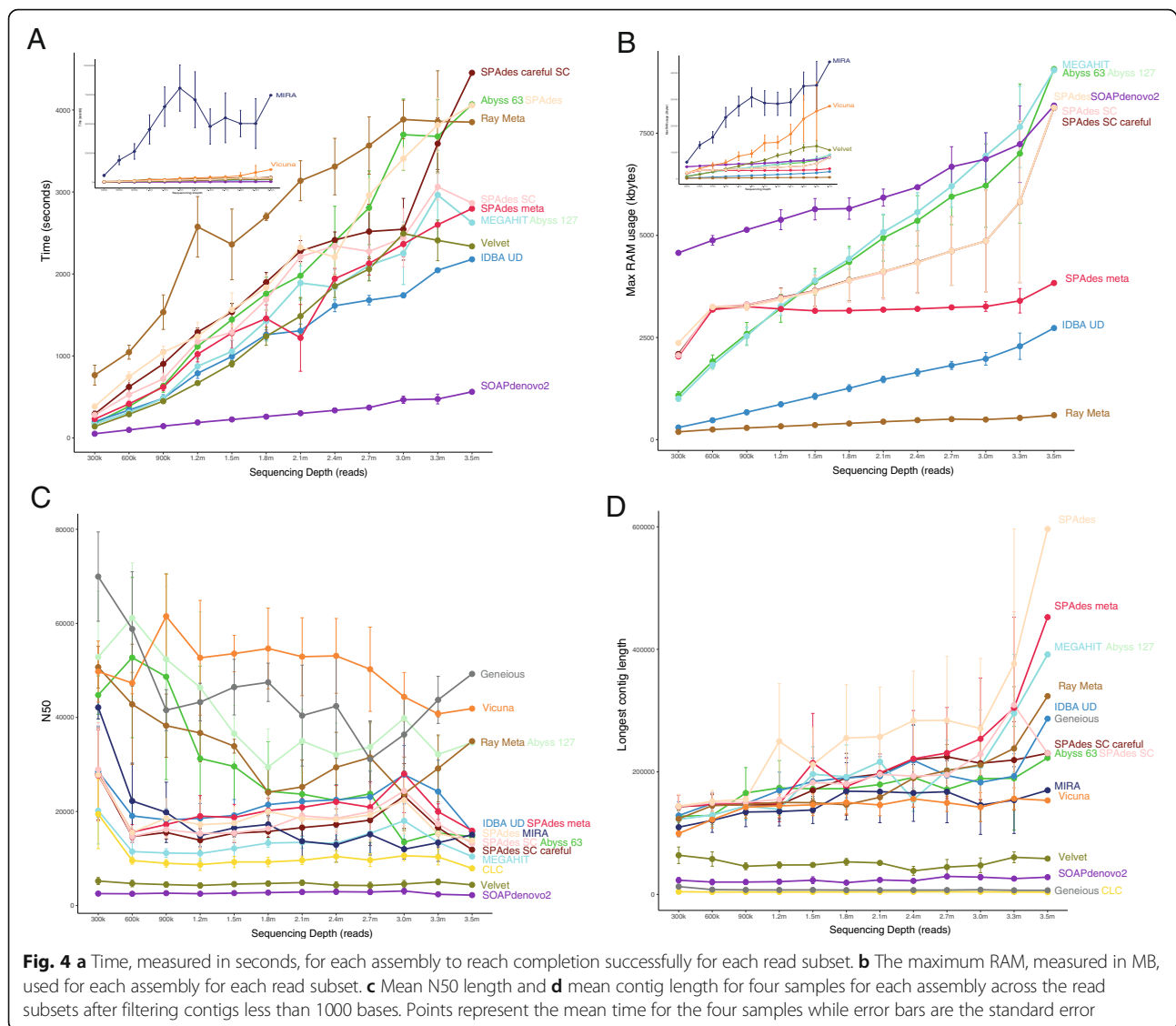
For both the N50 (Fig. 4c) and the longest contig length (Fig. 4d), there was a large amount of variation between samples for the majority of assemblers. The longest contig length showed a large increase at the final sequencing depth. Particular assemblers, namely SPAdes (default), SPAdes (meta), MEGAHIT and ABySS (k -mer 127), produced longer contigs as the sequence depth was increased.

Discussion

Many bacterial metagenomic assembly comparisons have highlighted that the choice of assembler has a significant impact on downstream analysis and the accuracy of the reconstructed metagenome [15, 28, 33, 54]. We have found this also to be true for viral metagenomes, where accurate and complete assembly are of particular importance given the lack of viral representation in reference databases. Virome studies depend heavily on the assembly step and possess many features which

are challenging to successful assembly. In this study, we compared the performance of those assemblers used to date in human viral metagenomics studies using datasets of known and unknown composition and varying complexity. These included a Q33-spiked virome, mock virome communities, a simulated virome and the “healthy human gut phageome” [31]. Each dataset provided unique attributes allowing for comparison of assembly performance on a number of levels. The combination of artificial and real viromes used in this study allows for the comparison of various aspects of assembly performance across a range of datasets rather than depending on simulated viromes alone, as is commonly carried out in assembly comparisons [13, 33].

The simulated dataset featured 572 viral genomes at various relative abundances as published by Vázquez-Castellanos and colleagues [53]. Fragmented assemblies of individual genomes within microbial communities hamper downstream analysis and limit the conclusions which can be drawn from metagenomic data such as taxonomic and functional profiles [11]. Consequently, the percentage genome recovery and degree of fragmentation were assessed across each assembler, with SPAdes (default, meta and single cell) each performing well. VICUNA performed very poorly, recovering only four contigs with high numbers of mismatches and misassemblies, despite having performed well with other datasets and being designed to address challenges of heterogeneous viral populations [57]. This failure may reflect the computational challenges relating to the format of the simulated reads, as benchmarks carried



out within the VICUNA study itself only include actual sequencing reads [57]. However, similar poor performance has been previously observed in virome assembly comparison using VICUNA and 454 reads [53]. For those assemblers which could recover greater than 90% of the reference genome in a single contig, SPAdes (default) outperformed SPAdes (meta). This may be explained by a lack of strain variants in the dataset and the fact that SPAdes (meta) was optimised to combine strain variants of each species to form consensus sequences.

A subset of genomes were poorly recovered (< 20% genome fraction) by nearly all assemblers. This observation indicates that there are challenging aspects of viral genomes and metagenomes which cannot be overcome with current assembly strategies. The strong positive correlations between the relative abundance and genome fraction

suggest that a low abundance threshold applies to virome assembly, below which assemblies will consist of small fractions of the viral genome, and in most cases be highly fragmented. This detrimental impact of low coverage has been well established in previous assembly comparison studies [13, 14, 44]. Highly abundant genomes also caused similar recovery and fragmentation issues across all assemblers, which is of particular importance due to the prevalence of extremely high abundance genomes in viral data (crAssphage, certain ssDNA viruses). As both abundance extremes are common in virome data, their impact must be considered when designing virome studies (i.e. sequencing depth). As relative abundance alone did not fully explain the variation in genome fraction recovered, the role of genomic repeats (a well-established assembly challenge [1]) was also investigated. However, genomic repeats

could explain the variation in genome fraction recovered to a lesser degree than relative abundance, suggesting other factors contribute to poor genome recovery.

Compositional differences between final assemblies and viromes themselves must be taken into account when drawing conclusions about virome composition and setting parameters for downstream analysis. Challenges such as genomic content and strain variation are not currently addressed in human virome assembly strategies and impact the reconstruction of certain members of a virome. Hybrid sequencing, which uses both long and short reads to resolve genomic regions associated with poor assembly [55], is a promising new technology which could address current virome assembly challenges. Library preparation methods which may reduce the bias introduced by MDA steps include using Swift Biosciences 1S Plus kit [46] and/or increasing overall sequencing depth or read length to improve recovery of lowly abundant viral genomes will be key. Furthermore, utilising an assembler which can robustly deal with ultra-high coverage genomes ($>1000\times$ coverage) is an important but often not appreciated aspect of virome assembly analysis. While promising, these potential solutions highlight a requirement for ongoing optimisation and extermination of virome analysis protocols.

Performance of some assemblers in this study was hampered by high coverage sequences (primarily overlap consensus assemblers). VICUNA assemblies exhibited the highest degree of fragmentation of all assemblers with mock A, despite having resolved both high abundance ssDNA genomes of mock B to a single contig. MIRA also exhibited a high degree of fragmentation with high abundance genomes in both simulated and mock datasets. However, MIRA was least affected by low abundance reads, recovering a greater genome fraction of low abundance genomes than other assemblers with fewer contigs. Performance of assemblers hampered by high coverage sequences in viromes may potentially be improved by sub-setting reads similar to the assembly approach used by SLICEMBLER [37].

Multi-assembler approaches such as the use of Geneious to generate consensus sequences from separate assemblers have been developed [9, 24, 47] but are not commonly included in human virome studies using short reads. MIRA assemblies of the Q33 genome and some low abundance genomes in the simulated dataset were improved using Geneious, resolving greater genome fractions with fewer contigs (despite Geneious recovering a lower genome fraction of the simulated dataset overall). It is possible that using these approaches could address issues facing each assembler, i.e. combine the assemblies of SPAdes (meta) which performs well across all four datasets but struggles to recover low abundant genomes, with MIRA assemblies

which are less affected by low abundance but have difficulty resolving genomes of higher abundance. Comparison of multi-assembler approaches and combinations of various assemblers were not within the scope of this study, but should not be ruled out as a potential method of improving virome assembly in cases where composition could be assessed and obvious assembly challenges were known to be present.

Across all analysis methods in this study, SPAdes (meta) performed consistently well and would be our recommendation. It performed best in the simulated data based on false positives, true positives and false negatives, best assembled the Q33 genome (recovery, fragmentation, misassemblies and genome size) and performed well with both mock communities in recovering all members accurately in one or two contigs. SPAdes (meta) RAM usage was low and did not increase to the same degree as other assemblers with increasing sequencing depth. This recommendation is in agreement with previous comparisons [54] which also suggested using SPAdes (meta) due to its ability to accurately assemble members of bacterial metagenomes. SPAdes (meta) is less able to accurately reconstruct micro-diversity as it generates a consensus assembly of “strain contigs” in a metagenome, which means it is better equipped to address the high mutation rates observed in virome data [40]. This recommendation is also concurrent with a previous study [44] which found IDBA UD, MEGAHIT and SPAdes (meta) to perform equally well when assessed using 14 simulated viromes. However, we found that SPAdes (meta) outperformed IDBA UD and MEGAHIT in the Q33-spiked dataset, RAM usage in relation to increasing sequencing depth, and in its ability to recover members of the simulated virome in a single contig. This recommendation contradicts two previous assembly comparisons which found CLC [20] and Velvet [56] to be best suited to virome data. However, SPAdes (meta) was not included in either study. Though SPAdes (meta) was outperformed by MIRA in the assembly of low abundance genomes in the Simulated dataset, MIRA has limited application to large datasets. MEGAHIT also performed well across all datasets performing well in relation to recovery, fragmentation and accuracy but encountered some recovery issues in mock datasets and minor accuracy issues with the Q33 genome.

The higher levels of accuracy (low mismatch indel and misassembly counts) of assemblers which performed poorly in other metrics, namely velvet and ABySS (k -mer 63), highlight the trade-off between accuracy and contiguity observed in previous assembly studies [16, 27]. However, both IDBA UD and MEGAHIT performed well for accuracy, genome recovery and fragmentation. These assemblers may be worth considering if strain level detail is of particular importance. The mock A and B datasets were

used to assess the impact of amplification bias on assembly performance. All ssDNA assemblies featured an equal minimum number of mismatches across both mock A and B. This may be caused by challenges in the genomes themselves hampering accurate assembly, but is more likely to reflect strain variation between genome sequence featured in the original publication and the genome of the phage used in the community itself.

The Q33-spiked virome consisted of pooled reads from three healthy human faecal samples, each of which having been spiked with 10^7 PFU ml⁻¹ of lactococcal phage Q33 prior to virome extraction. This allowed for assembly comparison of one abundant member of a challenging viral community. Despite the high relative abundance of the Q33 genome, only 6 assemblers could recover over 90% of the genome in a single contig, of these SPAdes (meta) and MEGAHIT reconstructed the Q33 genome accurately without the introduction foreign or chimeric DNA. It was also noted that the genome synteny was conserved across these six assemblies. This may reflect circularization of the linear Q33 genome during DNA extraction as the presence of cos sites has been previously predicted [30].

The longest contigs of each assembler were only detected at the highest sequencing depths and varied across assemblers, which may indicate that high coverage is necessary to recover the largest viral genomes in a community. However, it is also possible that these long contigs may reflect misassemblies and duplication events caused by read errors at high sequencing depths which must be considered when analysing high coverage data. At almost all sequencing depths, Geneious, VICUNA, Ray Meta and ABySS (k-mer 127) exhibited the highest N50 values, despite performing poorly in other metrics. This further highlights the limitation of using N50 alone as a metric of metagenomic assembly [54].

A further important consideration when performing any metagenomic assembly is practicality, size of dataset, computational resources, bioinformatic resources and how much hands-on time is required from the end user. Both CLC and Geneious are available as a GUI (albeit requiring a licence fee) which widens their audience to researchers with limited command-line experience (CLC can also be run using the windows command line). However, this limits their practicality for large scale virome studies as they are limited to the computational power of desktop computers and are not suited to the assembly of large numbers of samples. Despite the limitations of computational power, CLC performed well in all datasets in terms of genome recovery and fragmentation. Of the freely available open-source assemblers, MIRA and VICUNA are the least efficient in terms of RAM usage and assembly time, reflecting limitations of the overlap consensus approach to assembly. This limits

their applicability to large virome datasets and further increases the time required to carry out the Geneious assembly approach which requires the output of both assemblers. Despite the long runtime, VICUNA did not adhere to the number of cores specified, instead using all available cores. All other assemblers had a similar time requirements (with the exception of SOAPdenovo2 which performed poorly across all datasets). Of the assemblers which consistently performed well in terms of accuracy, genome fraction recovered and fragmentation, SPAdes (meta) was the most efficient in terms of RAM usage, which did not increase to the same degree as other assemblers with increasing sequencing depth. MIRA stood out in terms of impracticality by generating by far the largest intermediate files of any assembler, requiring several gigabytes of storage space for intermediate files.

The combination of results from four datasets facilitates accurate comparison of assemblers as the limitations of each individual dataset vary. Application of Phi29 MDA to amplify virome DNA to sufficient quantities for sequencing can introduce significant bias and skew the original composition of the virome, making quantitative viromics difficult [23, 46]. As a result, it is likely that true diversity of viral metagenomes is not being accurately captured using current virome extraction methods. However, as these procedures move away from steps known to introduce bias, greater diversity will be observed. In this sense, the level of complexity of the Q33 dataset, which pooled three independent human viromes, provides a useful testbed for metagenomic assemblers in future virome studies as extraction methods improve. Additionally, Q33 was not present in the viromes prior to spiking and assemblers were not challenged by the presence of native strain variations of Q33 genome. In this study, assemblers were compared without individual optimisation to the specific dataset. Feasibility dictates that this “straight out of the box” approach to assembly is used by almost all metagenomic assembly comparisons. Additionally, as the true composition of metagenomes is unknown, any impact of parameter optimisation must be estimated from general assembly statistics such as N50 and longest contig which have been highlighted to be of limited usefulness [3, 54]. Any parameter optimisation performed in the study (i.e. ABySS k-mer lengths, SPAdes careful vs. single cell) reflected parameters used in published virome studies and was not analysed in greater depth. While it is possible that parameter optimisation could improve individual assemblers, we believe that the differences in assembly algorithms are the primary drivers of assembly performance.

This study describes the impact of a crucial analysis step on virome characterisation and highlights the need

for a standardised analysis protocol across future virome studies. Such a protocol would allow for comparison across studies and facilitate accurate meta and cross analyses. This will be crucial should virome sequencing be utilised in diagnostic and clinical settings. However, it must be noted that any workflow will be somewhat limited and biased to the detection of particular viral taxa. Consequently, studies (e.g. identifying novel viruses) may benefit from implementing multiple assembly approaches due to the large number of factors, both technical (read length, quality, paired-end information etc.) and biological (genetic complexity, evenness etc.), which impact virome assembly.

Conclusions

Of all assembly programs used in human virome studies, SPAdes (meta) addressed the challenges of virome data most effectively. However, all assemblers have limitations and are hampered by aspects of virome data. Low read coverage and high genomic repeats lead to assemblies with low recovered genome fraction and a higher degree of fragmentation, with the assemblies themselves being less accurate. This pattern was seen across all assemblers used in this study.

As assembler choice has significant implications for virome composition and the conclusions which can be drawn from a dataset, assemblers which performed poorly in this study (i.e. low genome recovery or accuracy and high degree of fragmentation) highlight a potential untapped resource in the sequence data of previously conducted virome studies. It is highly likely that many viral sequences were poorly assembled and reanalysis using a more effective assembler may yield new insights. Researchers conducting meta-analysis of virome sequencing studies should take particular care when evaluating viral assemblies from different assembly programs. Design of future virome studies should carefully consider the impact of sequencing depth, as extremes in read coverage will prevent the assembly and detection of viral genomes at both high and low abundance.

Methods

Each assembler with the exception of Geneious and CLC was run as per manual with default parameters (unless stated) using a Lenovo x3650 M5 server with an intel Xeon processor E5-2690 v3 and 512Gb RAM running Ubuntu 14.04.5. Geneious assembly approach mirrored that used in [31] by generating consensus sequences from the assemblies of both MIRA and VICUNA. CLC and Geneious were run on a 64-bit windows 10 computer with an i5-4690 CPU and 16 GB of RAM.

Data sources

Sequencing reads from mock communities A and B featured in [46], simulated Virome dataset featured in [20],

reads used to compare the impact of sequencing depth on time and RAM usage featured in [31] and human viromes spiked with 10^7 PFU of Lactococcal phage Q33 [30] and originated from [49] .

Read pre-processing

Raw read quality was assessed with FastQC v0.11.5, and sequencing adapters were removed with cutadapt v1.9.1 [32] for the mock, spiked and healthy gut virome data sets. Trimming and filtering were carried out with Trimmomatic v0.36 [6] using parameters specific to each dataset. A sliding window size of 4 with a minimum Phred score of 30 and a minimum length of 60 bp was used with reads from both mock communities. The leading 15 bp and trailing 60 bp were removed from “healthy human gut phageome” reads, and a sliding window of 4 bp with a minimum Phred score of 20 was applied. The leading 10 bp and trailing 100 bp were removed from the Q33-spiked virome reads and a sliding window size of 4 bp with a minimum Phred score of 30. Filtered reads shorter than 60 bp were removed.

Analysis methods

Quality filtered reads from the Q33-spiked dataset consisted of three individual viromes which were pooled and subsequently assembled. Contigs were aligned to the published Q33 using Blastn with an e value cut-off of $1e^{-20}$. Top hit alignments to the Q33 genome with a minimum alignment length of 800 bases and which shared 95% identity were included in further analysis using QUAST (v. 4.4) [18] with “--unique mapping” flag. Further comparison and visualisation of Q33 assemblies were carried out using Mauve (v. 20150226, build 10) [8].

Alignment and comparison of assemblies from mock and simulated data sets to reference genomes were carried using MetaQUAST (v. 4.4) [35] with “--unique mapping” flag and default parameters (minimum contig length of 500 bp, minimum alignment length of 65 bp, minimum identity threshold of 95%). Correlations were carried out using Spearman’s method, and plots were generated using the package ggplot2 (v 3.0.0) package in R (v.3.4.3). These correlations were validated using a linear model in R base library. For data which was not normally distributed, log transformation was carried out.

Reads from the “healthy human gut phageome” were analysed to compare the overall assembler efficiency and the impact of sequencing depth. Reads were randomly subset in pairs (both the forward and reverse read of a pair were retained) to different depths using an in-house python script. Samples were subset in increments of 300,000 reads to their respective maximum depth (2.7, 3.5, 3 and 3.3 million reads). GNU time was utilised to measure the maximum RAM and length of time for

each assembly to reach completion. All assemblers were run using five threads where possible with the exception of CLC, Geneious, Ray Meta, Velvet and VICUNA. Ray Meta and Velvet were run with 10 one thread(s), respectively. Ray Meta failed to run with five while Velvet ran with one core despite five being allocated. VICUNA was also allocated five threads, however used upwards of 20. MetaVelvet was run, but after 7 days had failed to reach completion and was therefore removed from the subsequent analysis of these metrics. Contig statistics and filtering (contigs greater than 1 kb retained) were performed using the assembly-stats script (v1.0.1) from the Pathogen Informatics group at the Wellcome Sanger Institute (<https://github.com/sanger-pathogens/assembly-stats>).

Additional files

Additional file 1: Simulated virome MetaQUAST output. (HTML 6369 kb)
Additional file 2: Mock virome A MetaQUAST output. (HTML 528 kb)
Additional file 3: Mock virome B MetaQUAST output. (HTML 528 kb)
Additional file 4: Q33-spiked virome QUAST output. (HTML 360 kb)
Additional file 5: **Table S1.** Spearman correlation values from the relationships of indel, mismatch and misassembly counts, recovered genome fraction, abundance and total proportion of genomic repeats within the simulated virome. *GF–recovered genome fraction. **Table S2.** Linear modelling correlation values comparing recovered genome fraction, total proportion of genomic repeats and abundance for the Simulated virome. **Table S3.** Spearman correlation values from the relationships of inverted, tandem, palindromic and total repeats, abundance and the number of contigs generated by each assembler for the Simulated virome. **Table S4.** (A) Ranking table comparing recovered genome fraction and contig numbers for assemblers which recovered at least 50% of the total genome fraction. (B) Ranking table of indel, mismatch and misassembly counts per 100 kb, normalised to the number of genomes recovered to at least 50%. **Table S5.** Number of aligned and unaligned contigs generated by each assembler for mock community A. **Table S6.** Number of aligned and unaligned contigs generated by each assembler for mock community B. (XLSX 33 kb)
Additional file 6: **Figure S1.** Analysis of recovered genome fraction and indel/mismatch counts for mock communities A and B. Triangles represent NV A values for mismatches and indels caused by assembly failures. (PDF 293 kb)

Acknowledgements

We thank the authors of all datasets used in this study for the availability of their data.

Funding

This research was conducted with the financial support of Science Foundation Ireland (SFI) under Grant Number SFI/12/RC/2273 a Science Foundation Ireland's Spokes Programme which is co-funded under the European Regional Development Fund under Grant Number SFI/14/SP APC/B3032, and a research grant from Janssen Biotech, Inc.

Availability of data and materials

The Sequencing reads which support this study are available from the following links.

Mock communities A and B are available at: http://datacommons.cyverse.org/browse/iplant/home/shared/Virus/DNA_Viromes_library_comparison.

Simulated virome reads are available at: https://figshare.com/articles/Simulated_virome_dataset_for_assembly_and_annotation_tests/5151163

Reads used to compare the impact of sequencing depth on time and RAM usage are available from the NCBI SRA; <http://www.ncbi.nlm.nih.gov/sra> under the accession numbers SAMN04415496 to SAMN04415499.

Reads from human viromes spiked with 10^7 PFU of Lactococcal phage Q33 phage are available at <http://www.ncbi.nlm.nih.gov/sra> under the accession numbers SRX3240741, SRX3240716, SRX3240715.

Authors declare that all other data supporting the findings of this study are available within the article and its additional files.

Authors' contributions

FJR conceived the study. FJR and TDSS designed the experiments. TDSS, AGC and FJR carried out the bioinformatics analysis and drafted the manuscript. All authors approve and contributed to the final manuscript.

Ethics approval and consent to participate

Not applicable.

Consent for publication

Not applicable

Competing interests

The authors declared that they have no competing interests.

Abbreviations/glossary

The following terms, genome fraction, N50, number of contigs, misassemblies and local misassemblies, are defined by QUAST [35] Genome fraction "is the total number of aligned bases in the reference, divided by the genome size. A base in the reference genome is counted as aligned if there is at least one contig with at least one alignment to this base. Contigs from repeat regions may map to multiple places, and thus may be counted multiple times in this quantity."

N50 "is the contig length such that using longer or equal length contigs produces half (50%) of the bases of the assembly. Usually there is no value that produces exactly 50%, so the technical definition is the maximum length x such that using contigs of length at least x accounts for at least 50% of the total assembly length."

Number of contigs "is the total number of contigs in the assembly that have size greater than or equal to 0 bp."

Misassemblies "is the number of positions in the assembled contigs where the left flanking sequence aligns over 1 kbp away from the right flanking sequence on the reference (relocation) or they overlap on more than 1 kbp (relocation) or flanking sequences align on different strands (inversion) or different chromosomes (translocation)."

Local misassemblies "A local misassembly has two or more distinct alignments covering the breakpoint, the gap between left and right flanking sequences is less than 1 kbp and the left and right flanking sequences both are on the same strand of the same chromosome of the reference genome."

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Author details

¹APC Microbiome Ireland, Cork, Ireland. ²School for Microbiology, University College Cork, Cork, Ireland. ³Present Address: South Australian Health and Medical Research Institute, Adelaide, Australia. ⁴Teagasc Food Research Centre, Fermoy, Cork, Ireland.

Received: 16 October 2018 Accepted: 14 January 2019

Published online: 28 January 2019

References

1. Acuña-Amador L, Primot A, Cadieu E, Roulet A, Barloy-Hubler F. Genomic repeats, misassembly and reannotation: a case study with long-read resequencing of *Porphyromonas gingivalis* reference strains. *BMC Genomics*. 2018;19(1):54.
2. Aggarwala V, Liang G, Bushman FD. Viral communities of the human gut: metagenomic analysis of composition and dynamics. *Mob DNA*. 2017;8(1):12.
3. Aguirre de Cárcer D, Angly FE, Alcamí A. Evaluation of viral genome assembly and diversity estimation in deep metagenomes. *BMC Genomics*. 2014;15(1):989.

4. Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, Lesin VM, Nikolenko SI, Pham S, Pribelski AD. SPAdes: a new genome assembly algorithm and its applications to single-cell sequencing. *J Comput Biol.* 2012;19(5):455–77.
5. Boisvert S, Raymond F, Godzaridis É, Laviolette F, Corbeil J. Ray Meta: scalable de novo metagenome assembly and profiling. *Genome Biol.* 2012; 13(12):R122.
6. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics.* 2014;30(15):2114–20.
7. Breitbart M. Marine viruses: truth or dare; 2011.
8. Darling AE, Mau B, Perna NT. progressiveMauve: multiple genome alignment with gene gain, loss and rearrangement. *PLoS One.* 2010;5(6):e11147.
9. Deng X, Naccache SN, Ng T, Federman S, Li L, Chiu CY, Delwart EL. An ensemble strategy that significantly improves de novo assembly of microbial genomes from metagenomic next-generation sequencing data. *Nucleic Acids Res.* 2015;43(7):e46.
10. Dutilh BE, Cassman N, McNair K, Sanchez SE, Silva GG, Boling L, Barr JJ, Speth DR, Seguritan V, Aziz RK. A highly abundant bacteriophage discovered in the unknown sequences of human faecal metagenomes. *Nat Commun.* 2014;5:ncomms5498.
11. Florea L, Souvorov A, Kalbfleisch TS, Salzberg SL. Genome assembly has a major impact on gene content: a comparison of annotation in two *Bos taurus* assemblies. *PLoS One.* 2011;6(6):e21400.
12. Foulongne V, Sauvage V, Hebert C, Dereure O, Cheval J, Gouilh MA, Pariente K, Segondy M, Burguière A, Manuguerra J-C. Human skin microbiota: high diversity of DNA viruses identified on the human skin by high throughput sequencing. *PLoS One.* 2012;7(6):e38499.
13. Fritz A, Hofmann P, Majda S, Dahms E, Droegge J, Fiedler J, Lesker TR, Belmann P, DeMaere MZ, Darling AE. CAMISIM: simulating metagenomes and microbial communities, vol. bioRxiv; 2018. p. 300970.
14. García-López R, Vázquez-Castellanos JF, Moya A. Fragmentation and coverage variation in viral metagenome assemblies, and their effect in diversity calculations. *Front Bioeng Biotechnol.* 2015;3:141.
15. Greenwald WW, Klitgord N, Seguritan V, Yooseph S, Venter JC, Garner C, Nelson KE, Li W. Utilization of defined microbial communities enables effective evaluation of meta-genomic assemblies. *BMC Genomics.* 2017; 18(1):296.
16. Gritsenko AA, Nijkamp JF, Reinders MJ, Ridder D d. GRASS: a generic algorithm for scaffolding next-generation sequencing assemblies. *Bioinformatics.* 2012;28(11):1429–37.
17. Guo L, Hua X, Zhang W, Yang S, Shen Q, Hu H, Li J, Liu Z, Wang X, Wang H. Viral metagenomics analysis of feces from coronary heart disease patients reveals the genetic diversity of the Microviridae. *Virol Sin.* 2017;32(2):130–8.
18. Gurevich A, Saveliev V, Vyahhi N, Tesler G. QUAST: quality assessment tool for genome assemblies. *Bioinformatics.* 2013;29(8):1072–5.
19. Hannigan GD, Meisel JS, Tyldesley AS, Zheng Q, Hodkinson BP, SanMiguel AJ, Minot S, Bushman FD, Grice EA. The human skin double-stranded DNA virome: topographical and temporal diversity, genetic enrichment, and dynamic associations with the host microbiome. *MBio.* 2015;6(5):e01578–15.
20. Hesse U, van Heusden P, Kirby BM, Olonade I, van Zyl LJ, Trindade M. Virome assembly and annotation: a surprise in the Namib Desert. *Front Microbiol.* 2017;8:13.
21. Hurwitz BL, Sullivan MB. The Pacific Ocean Virome (POV): a marine viral metagenomic dataset and associated protein clusters for quantitative viral ecology. *PLoS One.* 2013;8(2):e57355.
22. Kearse M, Moir R, Wilson A, Stones-Havas S, Cheung M, Sturrock S, Buxton S, Cooper A, Markowitz S, Duran C. Geneious Basic: an integrated and extendable desktop software platform for the organization and analysis of sequence data. *Bioinformatics.* 2012;28(12):1647–9.
23. Kim K-H, Bae J-W. Amplification methods bias metagenomic libraries of uncultured single-stranded and double-stranded DNA viruses. *Appl Environ Microbiol.* 2011;77:00289–11.
24. Koren S, Treangen TJ, Hill CM, Pop M, Phillippy AM. Automated ensemble assembly and validation of microbial genomes. *BMC bioinformatics.* 2014; 15(1):126.
25. Li D, Luo R, Liu C-M, Leung C-M, Ting H-F, Sadakane K, Yamashita H, Lam T-W. MEGAHIT v1.0: a fast and scalable metagenome assembler driven by advanced methodologies and community practices. *Methods.* 2016;102:3–11.
26. Lim ES, Zhou Y, Zhao G, Bauer IK, Droit L, Ndao IM, Warner BB, Tarr PI, Wang D, Holtz LR. Early life dynamics of the human gut virome and bacterial microbiome in infants. *Nat Med.* 2015;21(10):1228.
27. Lin S-H, Liao Y-C. CISA: contig integrator for sequence assembly of bacterial genomes. *PLoS One.* 2013;8(3):e60843.
28. Lindgreen S, Adair KL, Gardner PP. An evaluation of the accuracy and speed of metagenome analysis tools. *Sci Rep.* 2016;6:19233.
29. Luo R, Liu B, Xie Y, Li Z, Huang W, Yuan J, He G, Chen Y, Pan Q, Liu Y. SOAPdenovo2: an empirically improved memory-efficient short-read de novo assembler. *Gigascience.* 2012;1(1):18.
30. Mahony J, Martel B, Tremblay DM, Neve H, Heller KJ, Moineau S, van Sinderen D. Identification of a new P335 subgroup through molecular analysis of lactococcal phages Q33 and BM13. *Appl Environ Microbiol.* 2013;79(14):4401–9.
31. Manrique P, Bolduc B, Walk ST, van der Oost J, de Vos WM, Young MJ. Healthy human gut phageome. *Proc Natl Acad Sci.* 2016;113(37):10400–5.
32. Martin M. Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet journal.* 2011;17(1):10–2.
33. Mavromatis K, Ivanova N, Barry K, Shapiro H, Goltsman E, McHardy AC, Rigoutsos I, Salamov A, Korzeniewski F, Land M. Use of simulated data sets to evaluate the fidelity of metagenomic processing methods. *Nat Methods.* 2007;4(6):495.
34. McCann A, Ryan FJ, Stockdale SR, Dalmasso M, Blake T, Ryan CA, Stanton C, Mills S, Ross PR, Hill C. Viromes of one year old infants reveal the impact of birth mode on microbiome diversity. *PeerJ.* 2018;6:e4694.
35. Mikheenko A, Saveliev V, Gurevich A. MetaQUAST: evaluation of metagenome assemblies. *Bioinformatics.* 2015;32(7):1088–90.
36. Minot S, Grunberg S, Wu GD, Lewis JD, Bushman FD. Hypervariable loci in the human gut virome. *Proc Natl Acad Sci.* 2012;109(10):3962–6.
37. Mirebrahim H, Close TJ, Lonardi S. De novo meta-assembly of ultra-deep sequencing data. *Bioinformatics.* 2015;31(12):i9–i16.
38. Namiki T, Hachiya T, Tanaka H, Sakakibara Y. MetaVelvet: an extension of Velvet assembler to de novo metagenome assembly from short sequence reads. *Nucleic Acids Res.* 2012;40(20):e155.
39. Norman JM, Handley SA, Baldrige MT, Droit L, Liu CY, Keller BC, Kambal A, Monaco CL, Zhao G, Fleshner P. Disease-specific alterations in the enteric virome in inflammatory bowel disease. *Cell.* 2015;160(3):447–60.
40. Nurk S, Meleshko D, Korobeynikov A, Pevzner PA. metaSPAdes: a new versatile metagenomic assembler. *Genome Res.* 2017;27(5):824–34.
41. Olson ND, Treangen TJ, Hill CM, Cepeda-Espinoza V, Ghurye J, Koren S, Pop M. Metagenomic assembly through the lens of validation: recent advances in assessing and improving the quality of genomes assembled from metagenomes. *Brief bioinform.* 2017.
42. Paul JH. Prophages in marine bacteria: dangerous molecular time bombs or the key to survival in the seas? *ISME J.* 2008;2(6):579.
43. Peng Y, Leung HC, Yiu S-M, Chin FY. IDBA-UD: a de novo assembler for single-cell and metagenomic sequencing data with highly uneven depth. *Bioinformatics.* 2012;28(11):1420–8.
44. Roux S, Emerson JB, Elie-Fadrosh EA, Sullivan MB. Benchmarking viromics: an in silico evaluation of metagenome-enabled estimates of viral community composition and diversity. *PeerJ.* 2017;5:e3817.
45. Roux S, Hallam SJ, Woyke T, Sullivan MB. Viral dark matter and virus–host interactions resolved from publicly available microbial genomes. *Elife.* 2015; 4:e08490.
46. Roux S, Solonenko NE, Dang VT, Poulos BT, Schwenck SM, Goldsmith DB, Coleman ML, Breitbart M, Sullivan MB. Towards quantitative viromics for both double-stranded and single-stranded DNA viruses. *PeerJ.* 2016;4:e2777.
47. Schürch AC, Schipper D, Bijl MA, Dau J, Beckmen KB, Schapendonk CM, Raj VS, Osterhaus AD, Haagmans BL, Tryland M. Metagenomic survey for viruses in Western Arctic caribou, Alaska, through iterative assembly of taxonomic units. *PLoS One.* 2014;9(8):e105227.
48. Szczyrba A, Hofmann P, Belmann P, Koslicki D, Janssen S, Dröge J, Gregor I, Majda S, Fiedler J, Dahms E. Critical assessment of metagenome interpretation—a benchmark of metagenomics software. *Nat Methods.* 2017;14(11):1063.
49. Shkoporov AN, Ryan FJ, Draper LA, Forde A, Stockdale SR, Daly KM, McDonnell SA, Nolan JA, Sutton TD, Dalmasso M. Reproducible protocols for metagenomic analysis of human faecal phageomes. *Microbiome.* 2018; 6(1):68.
50. Simpson JT, Wong K, Jackman SD, Schein JE, Jones SJ, Birol I. ABySS: a parallel assembler for short read sequence data. *Genome Res.* 2009;19(6):1117–23.
51. Smits SL, Bodewes R, Ruiz-Gonzalez A, Baumgärtner W, Koopmans MP, Osterhaus AD, Schürch AC. Assembly of viral genomes from metagenomes. *Front Microbiol.* 2014;5:714.

52. Solden L, Lloyd K, Wrighton K. The bright side of microbial dark matter: lessons learned from the uncultivated majority. *Curr Opin Microbiol*. 2016; 31:217–26.
53. Vázquez-Castellanos JF, García-López R, Pérez-Brocal V, Pignatelli M, Moya A. Comparison of different assembly and annotation tools on analysis of simulated viral metagenomic communities in the gut. *BMC Genomics*. 2014;15(1):37.
54. Vollmers J, Wiegand S, Kaster A-K. Comparing and evaluating metagenome assembly tools from a microbiologist's perspective-not only size matters. *PLoS One*. 2017;12(1):e0169662.
55. Warwick-Dugdale J, Solonenko N, Moore K, Chittick L, Gregory AC, Allen MJ, Sullivan MB, Temperton B. Long-read metagenomics reveals cryptic and abundant marine viruses. *bioRxiv*. 2018.
56. White DJ, Wang J, Hall RJ. Assessing the impact of assemblers on virus detection in a de novo metagenomic analysis pipeline. *J Comput Biol*. 2017; 24(9):874–881.
57. Yang X, Charlebois P, Gnerre S, Coole MG, Lennon NJ, Levin JZ, Qu J, Ryan EM, Zody MC, Henn MR. De novo assembly of highly diverse viral populations. *BMC Genomics*. 2012;13(1):475.
58. Zerbino D, Birney E. Velvet: algorithms for de novo short read assembly using de Bruijn graphs. *Genome Res*. 2008;18(5):821–9.

Ready to submit your research? Choose BMC and benefit from:

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more biomedcentral.com/submissions

