

New approaches for assembly of short-read metagenomic data

Martin Ayling¹, Matthew D Clark², Richard M. Leggett^{1*}

1. Earlham Institute, Norwich Research Park, Norwich, UK.

2. Natural History Museum, London, UK

* Corresponding author: richard.leggett@earlham.ac.uk

Abstract

In recent years, the use of longer-range read data combined with advances in assembly algorithms has stimulated big improvements in the contiguity and quality of genome assemblies. However, these advances have not directly transferred to metagenomic datasets, as assumptions made by the single genome assembly algorithms do not apply when assembling multiple genomes at varying levels of abundance. The development of dedicated assemblers for metagenomic data was a relatively late innovation and for many years, researchers had to make do using tools designed for single genomes. This has changed in the last few years and we have seen the emergence of a new type of tool built using different principles. In this review, we describe the challenges inherent in metagenomic assemblies and compare the different approaches taken by these novel assembly tools.

Keywords

Metagenomics, assembly, algorithms, sequencing

1. Introduction

The increasing rate of metagenomic data generation represents a serious challenge to bioinformaticians tasked with analysing and understanding it. The EBI metagenomics portal alone has seen a tenfold increase in the number of processed samples over the last two years [1]. The emergence of next generation sequencing technologies brought with it the possibility to use shotgun sequencing to identify the composition of species within a heterogeneous sample. As high throughput sequencing technologies become more economical and widespread, the opportunity to sequence previously uncharacterised organisms directly from their environmental niche allows for a more complete view of the microbial world. Metagenomics, or envirogenomics, is a branch of genomics which seeks to explore the composition of complex communities of organisms. As many organisms can not be grown in the laboratory (e.g. as many as 99% of bacteria are considered unculturable [2]), sequencing may be the only high throughput method to measure species diversity in many niches. Metagenomics has been applied to understanding a diverse range of environments including the human microbiome [3], the New York Subway [4], the virome of bats [5], the oceans [6,7], the crop rhizosphere [8] and communities of extreme microbes living in geysers and hot springs [9]. Through metagenomic approaches, we can gain valuable insights into changing community profiles resulting from environmental changes. In traditional genomic studies, a single species is isolated and then cultured to produce a DNA-rich sample to assess. However, many species (particularly viruses) are impossible to study in this way, and have thus remained unsequenced. Similarly, many studies have used 16S (prokaryotic) or 18S (eukaryotic) rRNA marker gene protocols, methods which while simple and cheap are unsuitable for detecting viral species and have limits in their discriminatory power for other classes of organism. By using whole genome shotgun sequencing (WGS), it is possible to more completely explore environmental samples without prior isolation and culturing, enabling

previously unsequenced species to be present in the resulting datasets. These opportunities pose new problems in data analysis, as metagenomic samples are inherently heterogeneous communities, sometimes containing tens of thousands of species [10, 11]. By considering a single bacterial genome in isolation, the problem of sequence assembly is relatively simple. However, the assumptions and simplifications which can be made in the case of a single genome are not applicable to heterogeneous environmental datasets. This presents both computational and conceptual obstacles necessitating the consideration of the more complicated problem of extricating as many genomes as possible from a complex mixture.

Though it is possible to analyse sequence data without assembly, most analyses can be improved by constructing longer more contiguous sequences (contigs). Next generation (Illumina) sequencing is comparatively cheap, but the short read length limits the information within a single read. The identification of structures within a genome longer than a read, e.g. genes or operons, only becomes possible if reads are first assembled into longer sequence stretches. Whilst metagenomic assembly is a research field still in its relative infancy compared to genomic assembly, the past few years have seen an increasing interest in its potential and a subsequent deluge of new software.

In this review, we begin with a brief discussion of the genome assembly problem and then describe the specific challenges posed by metagenomic data. We describe the approaches taken by the main metagenomic assembly tools, drawing out common themes and identifying unique traits. We discuss approaches to simplifying huge datasets prior to assembly, as well as pipelines that contain assembly as just one step within a more involved analysis. Finally, we discuss what makes a 'good' assembly.

2. (Single) genome assembly

Assembling short reads into contigs has many advantages. Longer stretches of sequence are more informative, allowing the researcher to consider whole genes or even gene clusters within a genome, and to understand larger genetic variants and repeats. Additionally, it has the effect of removing most sequencing errors, though this can be at the expense of new assembly errors. First generation (Sanger) sequencing technology produced far fewer reads than second generation (or 'next generation') sequencing technology, but individual reads were significantly longer (500-1000bp). Assembly of Sanger data used overlap-layout consensus (OLC) approaches (Figure 1a), in which overlaps are computed by comparing all reads to all other reads, overlaps are grouped together to form contigs (layout) and finally a consensus contiguous sequence, or contig, is determined by picking the most likely nucleotides from the overlapping reads (e.g. Celera [12]). With the advent of second generation technologies, the number of reads increased exponentially, but the average length of a read shortened significantly. This has enabled much reduced cost per gigabase of sequence, but the computational requirements of an OLC strategy become impractical due to the need to compare all reads with every other read in the dataset (millions or even billions of reads).

To overcome this computational hurdle, de Bruijn graph based assembly strategies were introduced [13] and have become widespread in the field. A de Bruijn graph (dBg) is a mathematical construct, where each vertex (or node) in a graph represents a *kmer* (a string of nucleotides of length *k*). Nodes in the graph are connected where they differ by all but one base and this base labels the edges. A graph is built by first decomposing each read into individual overlapping kmers, with edges formed by considering each kmer in turn and its overlaps with

existing nodes (Figure 1b). In an ideal case, the de Bruijn graph would form a single line, in which each node, apart from the two edges, is connected to one other node in the forward orientation and one in the reverse orientation. Converting such a graph into a hypothetical contig would be a trivial case of starting at one edge node and following labels to reach the second edge. Of course real datasets never result in such simple graphs and complex branching structures form as a result of errors, coverage differentials, heterozygosity, repeats and other structural variants. Thus much of the innovation in genome assembly algorithms has come from developing heuristics to simplify and navigate complex graphs consisting of millions of kmers to output contigs, as well as developing approaches to link contigs together in 'scaffolds'. The principle advantage of the de Bruijn graph approach is that, unlike OLC, there is no need to consider all input reads in a pairwise manner, making the problem instantly more tractable. Additionally the repetition inherent in so many reads may be compressed, reducing memory requirements. The main drawback of the dBg is the loss of context that results from breaking reads into smaller kmers which can result in the joining up of disparate parts of a genome containing the same kmers - e.g. repeats. As a dBg becomes complex, ascertaining the most likely paths through it (the predicted genomic sequence) can be difficult.

Most short-read assemblers which have been produced in the past decade utilise de Bruijn graphs for these reasons. These include popular genomic assemblers such as Velvet [14], ABySS [15] and SOAPdenovo [16]. There are limitations inherent in de Bruijn graph assembly however, particularly the initial choice of the size of kmer with which to build the graph. Choosing an inappropriate kmer size when building a graph may dramatically affect the quality of an assembly. Smaller kmers lead to more connected graphs; larger ones provide more specificity and fewer loops, but are more disconnected as the result of gaps or errors within the

read data or lack of coverage of the genome. Some assemblers have begun to use a variety of kmer sizes during assembly to mitigate this issue. IDBA builds graphs iteratively, starting with a small k-mer size, and using the predicted contigs as hypothetical reads with which to build the next graph with a longer kmer size [17]. SPAdes [18] employs analogous techniques, moving from graphs built using smaller kmer sizes to maximise connectivity combined with larger kmer sizes for simplicity. RAMPART [19] runs a range of different assemblers with an option to produce multiple assemblies with different kmers, and a report summarising the statistics of each one. These procedures have become feasible as a result of increased computational power, but also as a result of increased sequence throughput.

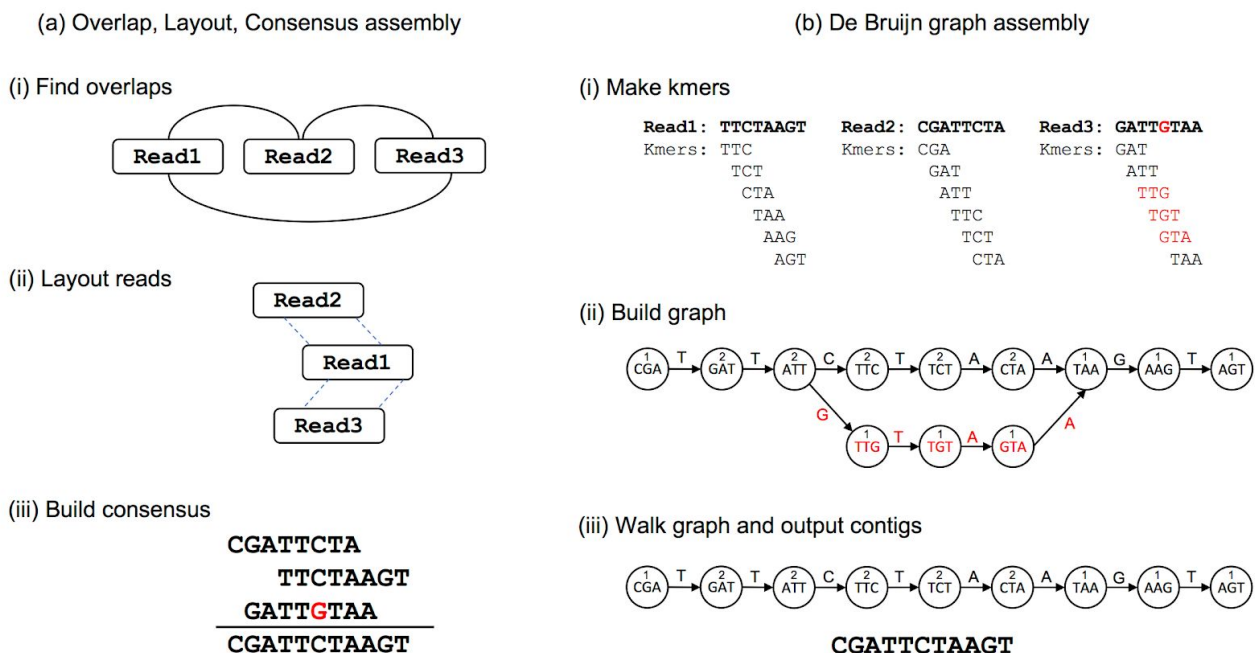


Figure 1: Two different approaches to genome assembly: (a) in Overlap, Layout, Consensus assembly, (i) overlaps are found between reads and an overlap graph constructed. (ii) Reads are laid out into contigs based on the overlaps (iii) The most likely sequence is chosen to construct consensus sequence. (b) In de Bruijn graph assembly, (i) reads are decomposed into kmers by sliding a window of size k across the reads. (ii) The kmers become vertices in the de Bruijn graph, with edges connecting overlapping kmers. Polymorphisms (red) form branches in the graph. A count is kept of how many times a kmer is seen, shown here as numbers above kmers. (iii) Contigs are built by walking the graph from edge nodes. A variety of heuristics handle branches in the graphs - for example, low coverage paths, as shown here, may be ignored.

3. The challenge of metagenomic assembly

3.1. Unknown abundance and diversity

In genomic assembly, there is an expectation that a sample contains a single species (apart from any contamination, which can be screened for prior to assembly) which allows assembly tools to make certain assumptions. The expected coverage of the target genome can be predicted from the total size of the dataset (the reads and their length) and the estimated size of the genome. In turn, it is assumed that nodes or edges in a graph occurring with very low coverage compared to the expected coverage are likely the result of sequencing errors or low level contamination, and the graph is simplified considerably by removing such nodes or paths. Similarly, nodes with much higher than average coverage can be assumed to be part of repeat structures. The typical optimal sequence coverage for a single genome assembler is in the 20-200x range, with a common “sweet spot” of ~50x [20]. However, in metagenomic datasets this assumption and simplifications cannot be made. Lower coverage nodes may originate from genomes with a lower abundance, not from errors, and so should not be discarded out of hand. Compounding this problem, the number of species within a sample, and the distribution of abundances of species is unknown. Abundance in heterogeneous samples often follows a power law [21], which means that many species will occur with similarly low abundances making the problem of distinguishing one from another problematic. The low coverage of most species means de novo assembly is unlikely unless the genome in question is relatively small, and instead we are reliant on reference genomes or gene prediction models/homology for evidence of the presence of species. It is imperative that a metagenomic assembly tool conserve as much of the less abundant species sequence as possible.

3.2. Related species

In a genomic study, it may be assumed that all sequence reads derive from the same original genome. In metagenomic studies, this is emphatically not the case, with a potentially huge diversity of species to consider. However, while distinguishing divergent species is already a difficult problem, an even more challenging problem is that of identifying closely related species or even strains within species. Often in metagenomic samples, a number of sub-species or strains are present, and this is particularly evident with viral communities which typically contain an abundance of haplotypes. Related species or subspecies introduce extensive overlaps in a kmer set, and therefore create assembly graphs which are considerably more complex as multiple genomes occupy much of the same kmer space. Branches or loops between these homologous regions make traversing the graph more complex, and if either species occurs with a low abundance, then identifying the presence of separate species will be difficult and deconvolving the graph is extremely complex. Mistakes at this point can lead to chimeric contigs containing sequence from more than one (sub-)species and a failure to capture the true diversity of the sample.

3.3. Memory and processing challenges

Metagenomic assembly brings additional processing challenges when compared to genomic assembly of similar sized organisms. Obtaining sufficient sequence coverage depth of a complex and diverse metagenomic sample can result in many times the volume of data than for a single organism, but excessive coverage can also result in observing more sequencing errors in the higher abundance genomes. As a consequence, a de Bruijn graph built to represent a metagenomic sample can require greater amounts of memory than one built to represent a similarly sized genomic sample, and will be more fragmented making graph traversal more

complicated. These challenges are mitigated directly by some assemblers, but can also be tackled using preprocessing methods before assembly, such as screening out likely host reads, subsampling and digital normalisation [22] to limit the effect of sequencing errors, or filtering/binning of reads and graph partitioning.

3.4. Filtering

An initial classification of reads within a dataset into likely taxonomic bins can be an effective first step in a metagenomic study. Removing reads from the dataset which are readily identifiable as derived from an available reference genome can streamline the process of assembling complex samples. MEGAHIT [23] stratifies the dataset into reads that are well supported by coverage, and those that are not; less well supported reads can still be included in an assembly as long as they extend well supported contigs. MetaCRAM [24], a metagenomic pipeline tool, uses Kraken [25], a kmer based taxonomic identification tool, to initially align reads to reference genomes and remove any known sequences from a dataset prior to assembly.

3.5. Graph partitioning

As metagenomic samples often contain multiple closely related species, the assembly graphs are often linked by shared sets of kmers, but this can also occur purely as a result of sequencing errors. This presents a challenge to assembler heuristics and can lead to the production of chimeric contigs which contain a mixture of sequence from multiple genomes. This would be problematic for two reasons: firstly, the chimeric contig in question may imply detection of a new species, a real possibility given the nature of metagenomics; and secondly, it would also reduce the likelihood of the contig mapping to either of the genomes underlying the chimera, thus affecting their abundance levels or even obscuring their presence. As a consequence, it can be

preferable to pre-partition the graph into sub-graphs, sacrificing contiguity for a more accurate community profile overall.

In practice, assemblers partition graphs by attempting to identify the nodes that join what appear to be distinct subgraphs, and removing them from the dataset. These nodes are often assessed as belonging in separate graphs as a result of differences in coverage, taking advantage of the notion that different species will be present in the original sample at different abundances. This difference in coverage is typically exploited on a localised basis to avoid the problems associated with the global coverage assumptions of genomic assemblers. IDBA-UD [26] and MEGAHIT both remove a node from within a graph if its coverage is significantly different to its neighbouring nodes. As a result of the iterative nature of both of these assemblers, any read which is not present in the graph as a result of this node removal will still be considered in the next round of graph building and so might be placed elsewhere in the graph (i.e. it is not discarded at this point). Metavelvet considers the overall distribution of kmers within the dataset, aiming to identify distinct peaks in a kmer distribution and using this as an indication of separate genomes with those approximate abundances. The graph is then split along these lines, maintaining local coverage within a subgraph. The critical joining node is not removed from the graph, but rather corresponding nodes are introduced into each subgraph to maintain connectivity.

Omega, although utilising an overlap based approach rather than a dBg based one, performs a similar step to graph partitioning [27]. During the process of building contigs, a contig will be split apart if the coverage of a constituent read is "significantly" below that of the rest of that contig's coverage.

Ray Meta [28] and MetaVelvet-SL [29] do not split graphs strictly based on coverage.

MetaVelvet-SL trains a support vector machine (SVM) to recognise probable chimeric nodes and remove them. Ray Meta operates a heuristic graph traversal procedure, which is based upon the minimum and peak coverages for each given read path through a graph.

3.6. Read pair information

Read pairs - both short-range paired end and long-range mate pair - are invaluable in *de novo* assembly of single genomes, providing links between disconnected contigs, scaffolding contigs and spanning areas of repeats. However, the benefits in metagenomic assembly are less clear cut, with paired information often lending support to more than one route through the graph.

Some assemblers still attempt the same scaffolding process used in genomic assembly, but others (MEGAHIT [23], Omega [27], PRICE [30], SPAdes [18], BIGMAC [31]) instead used paired reads to detect and resolve chimeric contigs produced from the misassembly of different genomes. Given the uneven coverage and low abundance of many of the species in most metagenomic samples, this produces more useful assemblies.

4. Approaches taken by metagenomic assemblers

Though initially most researchers used the common genome assembly tools to assemble metagenomic sequence data, the last few years have seen the emergence of a series of dedicated metagenomic assemblers. Here, we summarise the approaches taken, starting with the smaller group of assemblers based on overlap strategies and then considering those utilising de Bruijn graphs. Table 1 summarises the tools discussed and provides a brief comparison of their features.

Tool	Method	Key concepts	Reference
BBAP	OLC	Blast based overlap assembly, with optional intermediary assembly stage	Lin et al. 2017 [43]
Genovo	OLC	Generative probabilistic model; applies a series of hill-climbing steps iteratively until convergence; randomly (CRP prior) picks a contig to align read 'i' to. breaks up chimeric contigs by taking the edge reads off of contigs every ~5 iterations.	Laserson et al. 2011 [32] Afiahayati et al. 2013 [35]
IDBA-UD	dBg	Build graph; remove dead ends (<2k-1); merge bubbles; break graph on progressive (local) depth; error correction in reads (map reads to confident contigs; reads which match in all but a few bases can be 'corrected' to map perfectly); use mate pair info to build a 'local' assembly, avoid repeats and chimeras; hold trivial contigs, remove reads; make next graph; after k_max, partitions graph, clips tips, based on progressive (local) depth; PE requires long contigs to be effective.	Peng et al. 2012 [26]
IVA (Iterative Virus Assembler)	OLC	Aimed at viruses. Greedy kmer based extension. The most abundant kmer in the set is used as a seed, and this seed is grown out using a read which perfectly maps to it. A new kmer is drawn from the prefix of this read, which must be much more abundant than any other of the same size, and occur more than ten times in the dataset.	Hunt et al. 2015 [36]
MAP	OLC	Reads are filtered before overlap (reduce pairwise alignments made), simple paths found first, mate pair support used to simplify paths, edges removed with contradictory/insufficient mate pair support.	Lai et al. 2012 [41]
MEGAHIT	dBg	Solid kmers (occur more than a set threshold) and mercy kmers (remainder); mercy kmers that occur between two solid kmers in a read are kept; build a succinct dBG (dBG with BWT); remove tips, bubbles, progressively remove low local coverage edges; increasing kmer size, extract kmers from contigs and reads, build next graph.	Li et al. 2015 [23] Li et al. 2016 [68]
Metavelvet	dBg	dBG is first built with Velvet; population structure estimated from coverage of nodes (poisson distributions); dBG is partitioned into hypothetical	Namiki et al. 2012 [46]

		subgraphs (possibly different species) using these peaks as a guide; only nodes from primary distribution are considered - chimeric and repeat contigs are identified and split by PE info and coverage differences. Assembly produced for primary distribution; procedure repeated for next.	
MetaVelvet-SL	dBg	Similar to metavelvet - but the decision for identifying chimeric contigs is done using an SVM trained on (PE, coverage, contig lengths) for each dinucleotide (AA, AT...GG); a training set is generated from a similar population, the SVM is trained on this, then passed over the de bruijn graph for decomposition.	Afiahayati 2015 [29]
Omega	OLC	Read prefix/suffix (+/-) are stored in hashes; graph is built of V(r); simple paths (1 in, 1 out) are contracted, and transitive edges are reduced; tips removed (<10r) and bubbles are removed (hold edges with more r); minimum cost flow analysis for short (<1000bp) contigs; MP inserts are estimated from the assembly now, used to support contigs; scaffolding with LMP; remaining unresolved contigs are merged on similar coverage	Haider et al. 2014 [27]
PRICE	OLC	Reads are 'collapsed' if identical, then if near identical; then (single strand) dbg used to assemble (essentially) - greedy walking, start at highest coverage; identical contigs collapsed, then near identical contigs (ungapped) and finally gapped.	Ruby et al. 2013 [30]
Ray Meta	dBg	Extension of Ray – no graph partitioning performed, doesn't use a single peak for kmer coverage, min and peak coverage are specific for each read path; heuristics-based graph traversal; graph is coloured according to an expected taxonomic profile.	Boisvert et al. 2012 [28]
SAVAGE	OLC	Aimed at viral quasispecies recovery. Strict overlap conditions reproduces quasi-species assembly with minimal misassemblies.	Baaijens et al. 2017 [42]
SPAdes and metaSPAdes	dBg	SPAdes started out as a tool aiming to resolve uneven coverage in single cell genome data; metaSPAdes builds specific metagenomic pipeline on top of SPAdes. Multiple kmer sizes of dBG, starting with lowest kmer size and adding hypothetical kmers of (pref smallest useful size) to connect	Bankevich et al. 2012 [18] Nurk et al. 2017 [47]

		graph.	
VICUNA	Overlap	A min hash algorithm based on pairwise genetic distance threshold, inexact matching first (reads with similar or identical hash are merged) and then string matching of prefix/suffix of hashes is matched; (optional) target like reads are kept first (similar reads binned, similarity of bin is used), everything else removed.	Yang et al. 2012 [69]

Table 1: Metagenomic assembly tools: key concepts and references to papers.

4.1. Overlap based assemblers

Genovo [32] was one of the first metagenomic assemblers and is built using a generative probabilistic model that applies a series of hill-climbing steps iteratively. At each step, Genovo considers the position of every read and attempts to assign it to a new contig; upon finding a sufficiently good alignment it is added to that contig, otherwise a new contig is created. The assembly of chimeric contigs is prevented by removing the edge reads from all contigs every 5 iterations; should those reads have been correctly placed originally, they will be placed there again in the following steps. Genovo has been used in the reconstruction of bacterial and viral genomes from metagenomic samples [33, 34], and an extension to the assembler which made use of paired end read information was released later [35].

IVA [36] was developed for use with RNA virus populations, making it one of the few assemblers (along with VICUNA and SAVAGE) which specifically aims at viral rather than bacterial or eukaryotic samples. It performs greedy extensions within the dataset, starting with the most abundant kmer. This kmer is used as a seed, and this seed is grown outwards using a read which perfectly maps to it. A new kmer is drawn from the prefix of this read, and this kmer must also be common to the whole set; it must be much more abundant than any other kmer of the

same size, and occur more than ten times in the dataset. Though this assembler is not designed primarily as a metagenomic assembler, the authors assert that it is capable of performing well with samples of uneven coverage, a problem encountered when assembling environmental or heterogeneous samples. The software has been used in the assembly of viruses such as Zika virus and H1N1 influenza [37, 38]. Following a similar strategy, PRICE [30] also builds out an assembly using greedy paired end extension. However, it requires an initial assembly produced by a different assembler to start from, and then extends starting from the reads with the highest observed coverage, collapsing identical (and near-identical) reads to simplify the problem. Unlike other assemblers, it only functions in a single stranded orientation. The assembler has been used with Bunyavirus [39] and multiple water sample based metagenomic studies [40].

MAP [41] uses paired end information and specifically aims to break apart chimeric contigs in the assembly. Reads are filtered before the overlap stage to reduce the pairwise alignments required by the process, and simple paths joining reads are discovered first. Paired end reads are then used to support and simplify paths, with edges removed that are insufficiently supported in the dataset.

Omega [27] addresses the computational difficulties of OLC based assembly with a hash function built of the prefix and suffix of each read in the dataset which it uses to compute overlaps. A bi-directed graph is built up by matching reads to one another, and this is simplified by removing transitive edges (reads which are completely contained within a larger contiguous structure). Minimum cost flow analysis is performed on the basis of string copy number, to simplify the graph further, and long mate pair information is used to scaffold the contigs. There is

no explicit stage for resolving chimeric contigs; it is assumed that the nature of an OLC approach will hinder their formation.

SAVAGE [42] is an overlap based assembler of viral quasi-species, which reconstructs individual haplotypes in the final assembly by conservatively building overlap graphs (with strict minimal overlap length and of sequence similarity requirements). BBAP (the BLAST-based assembly pipeline [43]) creates a partial intermediary assembly which acts as a pseudo-reference for the remainder of the assembly process.

4.2. De Bruijn Graph based assemblers

In general when assembling using de Bruijn graph based tools, an a priori decision must be made about the size of the kmer in the underlying graph. This decision can greatly affect the resulting assembly - if the kmer size is too large, the resulting graph structure may be too disconnected, but if kmer size is too small, the graph may become overly connected making it harder to navigate paths through it. The IDBA family of assemblers (e.g. IDBA-UD) attempts to solve this problem by iterating through increasing kmer sizes, pruning the graph and merging bubbles (loops) along the way. The graph is broken up at points of significantly differing coverage, with information from paired end data included (although this is less informative in metagenomic rather than genomic cases). IDBA-UD has been used for assembly of a diverse range of bacterial and viral metagenomes (e.g. [44, 45]). More recently, MEGAHIT has used the process of increasing kmer size in assembly, but coupled it with succinct de Bruijn graphs which are more efficient in computational terms. This assembler is several orders of magnitude faster and requires significantly less memory in the process of assembly, and shows further performance increases when run on a GPU. This attention to computational performance as well

as to assembly completeness has made MEGAHIT one of the most popular of the current crop of metagenomic assemblers.

Velvet, a popular genome assembler, has received two updates aimed at metagenomic assembly in the form of MetaVelvet [46] and MetaVelvet-SL [29]. In MetaVelvet, a dBg graph is built using Velvet and the population structure is estimated from the coverage of nodes (modelled as Poisson distributions). The graph is then partitioned into subgraphs (each a hypothetical different species) using these coverage peaks as a guide. Chimeric and repeat contigs are identified and split using paired end information and local differences in coverage. This assumes that genomes are distinct mostly on coverage information (which will be relative to abundance), which may not be the case with low abundance genomes that are more susceptible to stochastic noise. MetaVelvet-SL is an extension of MetaVelvet that improves upon the decision making process for identifying chimeric contigs. An SVM (support vector machine) is trained on multiple criteria (paired end information, coverage, contig lengths) for each dinucleotide pairing (AA, AT...GG); a training set is generated from a similar population to the sample, the SVM is trained on this, and then passed over the sample graph for decomposition.

Ray is another commonly used genomic assembly to have received a metagenomic adaption in the form of RayMeta [28]. This is an extension of Ray, where no graph partitioning is performed, but unlike Ray (where a single peak coverage is expected for the whole graph and kmers with a significantly lower coverage are excluded), in RayMeta a localised coverage distribution is generated for each read path. These graphs are then walked using heuristic methods. A big emphasis for RayMeta is on computational efficiency and a lot of effort has been focused on scalability and distributability across standard clusters. This enables complex datasets to be

processed across a networked cluster of low memory machines, avoiding the need for expensive, large memory architectures.

Although not specifically a metagenomic assembler, SPAdes [18] is aimed at genome assembly from single cell data, but its core assumptions of uneven coverage also make it suitable for metagenomic assembly. It builds multiple dBgs with differing kmer sizes, and adds hypothetical kmers to ensure a connected graph. Chimeric contigs which are produced by these hypothetical kmers are then identified and split in a later stage. metaSPAdes [47] incorporates SPAdes into a metagenomic assembly pipeline and introduces new heuristics for differentiating intergenomic repeats between species.

Finally, VICUNA is one of the few assemblers released which focuses on reconstructing viral genomes. It uses neither a dBg nor an OLC approach. Rather it clusters similar reads together first, by generating a hash value for each read. Reads which are identical or similar will share the same hash value. These reads are then used to construct contigs, based on shared kmers, and reads which appear in multiple hashes can enable contigs to be merged. This is not guaranteed to detect all good suffix/prefix matches however, so a further seed based extension is performed on the now greatly reduced dataset. The authors propose this for populations of diverse but non-repetitive genomes, with high but variable coverage.

5. Assembly pipelines

A number of software pipelines are available that combine read pre-processing, metagenomic assemblers and post-assembly analysis. Perhaps the most comprehensive example is MetAMOS [48], which, at the time of writing, supports almost 20 genomic and metagenomic

assemblers, along with a wide range of pre-processing, filtering, validation and annotation tools. Users can create workflows containing combinations of the tools that are suited to their datasets.

InteMAP [49] integrates output from two DBG assemblers (ABYSS, IDBA-UD) and one OLC assembler (Celera) by separately merging low and high coverage contigs from pairs of assemblers. The authors of EnsembleAssembler also argue that merging the output from DBG and OLC assemblies can produce improved results [50]. MetaCRAM [24] is focussed on efficient storage via compression of metagenomic datasets. It taxonomically classifies reads and then assembles unclassified reads using IDBA-UD. Both the aligned reads and the unaligned read assemblies are then compressed for storage. MetaCompass [51] first maps reads against reference datasets, then generates reference-guided contigs, polishes them with Pilon [52] and finally combines unmapped reads with the polished contigs using MEGAHIT.

6. Assessing assembly quality

With an ever increasing range of metagenomic assemblers available, how can researchers choose the tool for their application? The N50 is an oft quoted statistic that is casually used to imply the quality of an assembly. If all contigs in an assembly are ordered by length, the N50 is the minimum size of contigs that contains 50% of the assembled bases. For example an N50 of 10,000 bp means that 50% of the assembled bases are contained in contigs of at least 10,000 bp. This statistic only indicates the contiguity of the assembled bases, is easy to manipulate (e.g. tools make different decisions on removal of small contigs which they consider noise or chaff), and gives no measure of assembly accuracy. A new assembler could generate long strings of random As, Cs, Gs and Ts and achieve high N50 but with no accuracy to the underlying genome, indeed the N50 could even be larger than the biological genome. Thus while it is the

most used assembly statistic, it must be used cautiously and its significance understood. For example well established assembly tools designed for single genomes may produce assemblies of metagenomic datasets with high N50 values. However, this may have been achieved by removing kmers representing lower coverage species or collapsing inter-strain variation e.g. sacrificing complexity for contiguity.

Assembly contiguity is important - after all, the whole point of assembling a metagenomic dataset is to obtain longer sequences for downstream analysis. However, the ability to capture the metagenomic diversity of a sample - including the lower abundance species and strains - may be equally important. Thus there is a compromise between the desire for long contiguous sequence and the desire for an accurate representation of community composition, possibly down to the strain level. The aim of the project should lead to a choice of assembler and assembly parameters - particularly kmer size - that moves the emphasis one way or another.

A number of tools exist for assessing metagenome assembly quality. MetaQUAST [53] performs a BLAST search of contigs against a database of 16S rRNA genes and will automatically download the top 50 references. It then performs a reference-based quality assessment of contigs that align to these references. Such an approach is limited only to bacterial sequences. BUSCO (Benchmarking Universal Single-Copy Orthologs) uses gene content to assess assembly quality and completeness [54]. It comes with a database of single-copy vertebrate, arthropod, metazoan, fungi and eukaryotic genes, as well as a smaller set of prokaryotic universal marker genes. CheckM also uses the presence of marker genes to assess assembly quality, but incorporates information about the position of a genome within a reference genome tree and collocation of genes in order to improve accuracy [55]. The Assembly Likelihood

Framework (ALE) evaluates genomic and metagenomic assemblies with a reference-free approach that incorporates read quality, mate pair orientation, read pair insert length, sequencing coverage, read alignment and k-mer frequency [56].

In the field of single genome assembly, contests have been used in an attempt to compare the performance of different algorithms. In metagenomic assembly, the Critical Assessment of Metagenome Interpretation (CAMI) set out to develop an “independent, comprehensive and bias-free evaluation” of both binning and assembly methods [57]. Its success relied on developers of tools and pipelines being willing to submit answers to a set of challenges and the organisers received six entries to the metagenomic assembly contest. Contestants were required to submit reproducible assemblies of three simulated metagenomic communities which were created from real sequencing data of newly sequenced viruses, bacteria and their plasmids. The results demonstrated substantial differences between the assemblies produced by the six teams - for example total assembly size ranged from 12.32 Mb to 1.97 Gb for a dataset with an expected assembly size of 2.80 Gb. Results also varied substantially according to the parameter settings chosen for each tool. Notably assemblers using multiple kmers performed better than those using a single kmer size. All tools struggled with assembly of closely related genomes and the authors describe this as an “unsolved problem”. Overall, there were three assembly tools that performed better - MEGAHIT, Meraga (MEGAHIT combined with Meraculous [58]) and Minia [59] - but it’s not clear that this will necessarily be the case with all datasets, or in the hands of all researchers.

7. Conclusion

Assembling genomes out of heterogeneous samples is an extremely challenging problem and one that remains unsolved. The first specialised metagenomic assembly tool was released comparatively recently, in 2011, and the intervening years have seen the introduction of a wealth of new tools. Picking the right tool and then picking the right parameters for a specific dataset are not straightforward tasks. Projects like the CAMI competition can contribute to the understanding of the strengths and weaknesses of different approaches, but researchers will benefit from trying a range of tools and parameters. As such, there is really no substitute for dedicated post-assembly analysis using both automated tools such as MetaQUAST and manual analysis by the researchers themselves.

The focus of this article has been on assembly tools for short-read metagenomics, as Illumina remains the dominant platform for metagenomics [1] due to the lowest cost per Gbp of sequence and the need for high depth of sequencing of metagenomics samples. New library methods for Illumina sequencers e.g. Illumina synthetic long reads [60], Dovetail in vitro HiC [61], and 10x Genomics microfluidics created read clouds [62] allow more contiguous assemblies but require longer DNA (10kb for synthetic long reads, and over 50kb for Dovetail and 10x Genomics) which may be hard to extract from all samples, especially without introducing bias. In vivo HiC cross-links DNA within live cells, allowing scaffolding similar to Dovetail, but uniquely it also allows grouping of chromosomes and plasmids in the same original cells [63].

Researchers are increasingly attracted to long-read technologies e.g. from established Pacific Biosciences [64] or the cheap, portable Oxford Nanopore Technologies MinION [65]. Both may simplify the need for assembly (with individual reads spanning multiple genes) or allow for

generation of much longer contiguous sequence. Assembly of reads from these third generation platforms abandons de Bruijn graph approaches and returns to the Overlap/Layout/Consensus models used in the earlier days of Sanger sequencing. As yet, there are no published tools dedicated solely to assembly of metagenomes from third generation platforms, but impressive results are possible using genome assembly tools such as Canu [66] or the very computationally efficient Minimap [67]. As the cost comes down and the accuracy and yields improve, these new technologies are likely to seem increasingly attractive platforms for metagenomic experiments.

References

1. Mitchell AL, Scheremetjew M, Denise H et al. (2017). EBI Metagenomics in 2017: enriching the analysis of microbial communities, from sequence reads to assemblies. *Nucl. Acids Res.* 46(D1):D726-D735.
2. Ling LL, Schneider T, Peoples AJ et al. (2015). A new antibiotic kills pathogens without detectable resistance. *Nature* 517:455-459.
3. The Human Microbiome Project Consortium (2012). Structure, function and diversity of the healthy human microbiome. *Nature* 486(7402):207–214.
4. Afshinnekoo E, Meydan C, Chowdhury S et al. (2015). Geospatial Resolution of Human and Bacterial Diversity with City-Scale Metagenomics. *Cell Systems* 29;1(1):72-87.
5. Baker KS, Leggett RM, Bexfield NH et al. (2013) Metagenomic study of the viruses of African straw-coloured fruit bats: Detection of a chiropteran poxvirus and isolation of a novel adenovirus. *Virology* 441(2):95–106.
6. Venter JC, Remington K, Heidelberg JF et al. (2004). Environmental Genome Shotgun Sequencing of the Sargasso Sea. *Science* 304(5667):66-74.

7. Sunagawa S, Coelho LP, Chaffron S et al. (2015). Structure and function of the global ocean microbiome. *Science* 348(6237).
8. Turner TR, Ramakrishnan K, Walshaw J et al. (2013). Comparative metatranscriptomics reveals kingdom level changes in the rhizosphere microbiome of plants. *The ISME Journal* (2013) 7, 2248–2258.
9. Strazzulli A, Fusco S, Cobucci-Ponzano B et al. (2017). Metagenomics of microbial and viral life in terrestrial geothermal environments. *Reviews in Environmental Science and Bio/Technology* 16(3):425-454.
10. Daniel R. (2005). The metagenomics of soil. *Nature Reviews Microbiology* 3:470-478.
11. Nesme J, Achouak W, Agathos SN et al. (2016). Back to the Future of Soil Metagenomics. *Frontiers in Microbiology* 7:73. Doi: doi: 10.3389/fmicb.2016.00073.
12. Myers EW, Sutton GG, Delcher AL et al. (2000). A whole-genome assembly of *Drosophila*. *Science* 287(5461):2196-204.
13. Pevzner PA, Tang H, Waterman MS (2001). An Eulerian path approach to DNA fragment assembly. *Proc Natl Acad Sci USA* 98(17): 9748–9753.
14. DR Zerbino and E Birney (2008). Velvet: Algorithms for de novo short read assembly using de Bruijn graphs. *Genome Research* 18(5): 821–829.
15. Simpson JT, Wong K, Jackman SD et al. (2009). ABySS: A parallel assembler for short read sequence data. *Genome Research* 19(6):1117–1123.
16. Li R, Zhu H, Ruan J et al. (2010). De novo assembly of human genomes with massively parallel short read sequencing. *Genome Research* 20: 265-272.
17. Peng Y, Leung HCM, Yiu SM et al. (2010). IDBA – a practical iterative de Bruijn graph de novo assembler. In: Berger B (eds) *Research in Computational Molecular Biology. RECOMB 2010. Lecture Notes in Computer Science*, vol 6044. Springer, Berlin, Heidelberg.

18. Bankevich A, Nurk S, Antipov D et al. (2012). SPAdes: A New Genome Assembly Algorithm and Its Applications to Single-Cell Sequencing. *Journal of Computational Biology* 19(5) (2012), 455-477.
19. Mapleson D, Drou N, Swarbreck D (2015). RAMPART: a workflow management system for de novo genome assembly. *Bioinformatics* 31(11):1824-6.
20. Desai A, Marwah VS, Yadav A et al. (2013). Identification of Optimum Sequencing Depth Especially for De Novo Genome Assembly of Small Genomes Using Next Generation Sequencing Data. *PLOS One* 12;8(4):e60204. doi:10.1371/journal.pone.0060204.
21. Matthews TJ, Whittaker RJ (2014). On the species abundance distribution in applied ecology and biodiversity management. *Journal of Applied Ecology* 52(2):443-454.
22. Howe AC, Jansson JK, Malfatti SA et al. (2014). Tackling soil diversity with the assembly of large, complex metagenomes. *PNAS* 111(13):4904–4909.
23. Li D, Liu CM, Luo R et al. (2015). MEGAHIT: an ultra-fast single-node solution for large and complex metagenomics assembly via succinct de Bruijn graph. *Bioinformatics* 31(10):1674-6.
24. Kim M, Zhang X, Ligo JG et al. (2016). MetaCRAM: an integrated pipeline for metagenomic taxonomy identification and compression. *BMC Bioinformatics* 17:94.
25. Wood DE, Salzberg SL (2014). Kraken: ultrafast metagenomic sequence classification using exact alignments. *Genome Biology* 15:R46.
26. Peng Y, Leung HC, Yiu SM et al. (2012). IDBA-UD: a de novo assembler for single-cell and metagenomic sequencing data with highly uneven depth. *Bioinformatics* 1;28(11):1420-8.
27. Haider B, Ahn TH, Bushnell B et al. (2014). Omega: an overlap-graph de novo assembler for metagenomics. *Bioinformatics* 30(19):2717-22.
28. Boisvert B, Raymond F, Godzaridis E et al. (2012). Ray Meta: scalable de novo metagenome assembly and profiling. *Genome Biology* 13:R122.

29. Afiahayati, Sato K, Sakakibara Y (2015). MetaVelvet-SL: an extension of the Velvet assembler to a de novo metagenomic assembler utilizing supervised learning, *DNA Research* 22(1):69–77.
30. Ruby JG, Bellare P, Derisi JL (2013). PRICE: software for the targeted assembly of components of (Meta) genomic sequence data. *G3* 20;3(5):865-80.
31. Lam K, Hall R, Clum A et al. (2016). BIGMAC : breaking inaccurate genomes and merging assembled contigs for long read metagenomic assembly. *BMC Bioinformatics* 17:435.
32. Laserson J, Jojic V, Koller D (2011). Genovo: de novo assembly for metagenomes. *Journal of Computational Biology* 18(3):429-43.
33. Gupta A, Kumar S, Prasoodanan VPK et al. (2016). Reconstruction of Bacterial and Viral Genomes from Multiple Metagenomes. *Frontiers in Microbiology* 7:469.
34. Vázquez-Castellanos JF, García-López R, Pérez-Brocal V et al. (2014). Comparison of different assembly and annotation tools on analysis of simulated viral metagenomic communities in the gut. *BMC Genomics* 15:37.
35. Afiahayati, Sato K, Sakakibara Y (2013). An extended genovo metagenomic assembler by incorporating paired-end information. *PeerJ* 1:e196.
36. Hunt M, Gall A, Ong SH et al. (2015). IVA: accurate de novo assembly of RNA virus genomes. *Bioinformatics* 31(14):2374-6.
37. Lahon A, Arya RP, Kneubehl AR et al. (2016). Characterization of a Zika Virus Isolate from Colombia. *PLoS Neglected Tropical Diseases* 10(9):e0005019.
38. Watson SJ, Langat P, Reid SM et al. (2015). Molecular Epidemiology and Evolution of Influenza Viruses Circulating within European Swine between 2009 and 2013. *Journal of Virology* 89(19):9920–9931.

39. Chandler JA, Thongsripong P, Green A et al. (2014). Metagenomic shotgun sequencing of a Bunyavirus in wild-caught *Aedes aegypti* from Thailand informs the evolutionary and genomic history of the Phleboviruses. *Virology* 464:312-319.
40. Ross DE, Gulliver D (2016). Reconstruction of a Nearly Complete *Pseudomonas* Draft Genome Sequence from a Coalbed Methane-Produced Water Metagenome. *Genome Announcements* 4(5):e01024-16.
41. Lai B, Ding R, Li Y et al. (2012). A de novo metagenomic assembly program for shotgun DNA reads. *Bioinformatics* 28(11):1455-62.
42. Baaijens JA, El Aabidine AZ, Rivals E et al. (2017). De novo assembly of viral quasispecies using overlap graphs. *Genome Research* 27:835–848.
43. Lin Y, Hsieh C, Chen J et al. (2017). De novo assembly of highly polymorphic metagenomic data using in situ generated reference sequences and a novel BLAST-based assembly pipeline. *BMC Genomics* 18:223.
44. Norman JM, Handley SA, Baldrige MT et al. (2015). Disease-Specific Alterations in the Enteric Virome in Inflammatory Bowel Disease. *Cell* 160(3):447-460.
45. Di Rienzi SC, Sharon I, Wrighton KC et al. (2013). The human gut and groundwater harbor non-photosynthetic bacteria belonging to a new candidate phylum sibling to Cyanobacteria. *eLife* 2:e01102.
46. Namiki T, Hachiya T, Tanaka H et al. (2012). Metavelvet: an extension of velvet assembler to de novo metagenome assembly from short sequence reads. *Nucleic Acids Research* 40:e155.
47. Nurk S, Meleshko D, Korobeynikov A et al. (2017). metaSPAdes: a new versatile metagenomic assembler. *Genome Research* 27: 824-834.
48. Treangen TJ, Koren S, Sommer DD et al. (2013). MetAMOS: a modular and open source metagenomic assembly and analysis pipeline. *Genome Biology* 14:R2.

49. Lai B, Wang F, Wang X et al. (2015). InteMAP: Integrated metagenomic assembly pipeline for NGS short reads. *BMC Bioinformatics* 16:244.
50. Deng X, Naccache SN, Ng T et al. (2015). An ensemble strategy that significantly improves de novo assembly of microbial genomes from metagenomic next-generation sequencing data. *Nucleic Acids Research* 43(7):e46.
51. Cepeda V, Liu B, Almeida M et al. (2017). MetaCompass: Reference-guided Assembly of Metagenomes. bioRxiv. doi:10.1101/212506.
52. Walker BJ, Abeel T, Shea T et al. (2014). Pilon: An Integrated Tool for Comprehensive Microbial Variant Detection and Genome Assembly Improvement. *PLoS One* 9:e112963.
53. Mikheenko A, Saveliev V, Gurevich A (2016). MetaQUAST: evaluation of metagenome assemblies. *Bioinformatics* 32(7):1088-1090.
54. Simão FA, Waterhouse RM, Ioannidis P et al. (2015). BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics* 31(19):3210–3212.
55. Parks DH, Imelfort M, Skennerton CT et al. (2015). CheckM: assessing the quality of microbial genomes recovered from isolates, single cells, and metagenomes. *Genome Research* 25:1043-1055.
56. Clark SC, Egan R, Frazier PI et al. (2013). ALE: a generic assembly likelihood evaluation framework for assessing the accuracy of genome and metagenome assemblies. *Bioinformatics* 29(4):435-443.
57. Sczyrba A, Hofmann P, Belmann P et al. (2017). Critical Assessment of Metagenome Interpretation - a benchmark of metagenomics software. *Nature Methods* 14, 1063–1071.
58. Chapman J, Ho I, Sunkara S et al. (2011). Meraculous: de novo genome assembly with short paired-end reads. *PLoS One* 6, e23501.

59. Chikhi R, Rizk G (2013). Space-efficient and exact de Bruijn graph representation based on a Bloom filter. *Algorithms for Molecular Biology* 8:22.
60. McCoy RC, Taylor RW, Blauwkamp TA et al. (2014). Illumina TruSeq Synthetic Long-Reads Empower De Novo Assembly and Resolve Complex, Highly-Repetitive Transposable Elements. *PLOS One* 27(5):757-767.
61. Putnam NH, O'Connell BL, Stites JC et al. (2016). Chromosome-scale shotgun assembly using an in vitro method for long-range linkage. *Genome Research* 26(3):342-50.
62. Weisenfeld NI, Kumar V, Shah P et al. (2017). Direct determination of diploid genome sequences. *Genome Research* 27(5):757-767.
63. Stewart RD, Auffret MD, Warr A et al. (2018). Assembly of 913 microbial genomes from metagenomic sequencing of the cow rumen. *Nature Communications* 9:870.
64. Frank JA, Pan Y, Tooming-Klunderud A et al. (2016). Improved metagenome assemblies and taxonomic binning using long-read circular consensus sequence data. *Scientific Reports* 6:25373.
65. Leggett RM, Clark MD (2017). A world of opportunities with nanopore sequencing. *Journal of Experimental Botany* 68(20):5419–5429.
66. Koren S, Walenz BP, Berlin K et al. (2017). Canu: scalable and accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome Research* 27:722-736.
67. Li H. (2016). Minimap and miniasm: fast mapping and de novo assembly for noisy long sequences. *Bioinformatics* 32(14):2103-2110.
68. Li D, Luo R, Liu CM et al. (2016). MEGAHIT v1.0: A fast and scalable metagenome assembler driven by advanced methodologies and community practices. *Methods* 102:3-11.
69. Yang X, Charlebois P, Gnerre S et al. (2012). De novo assembly of highly diverse viral populations. *BMC Genomics* 13:475.

Funding

This work was supported by a Biotechnology and Biological Sciences Research Council (BBSRC) grant to RML (BB/M004805/1), a Core Strategic Programme Grant to Earlham Institute (BB/J004669/1) and by the Natural History Museum.