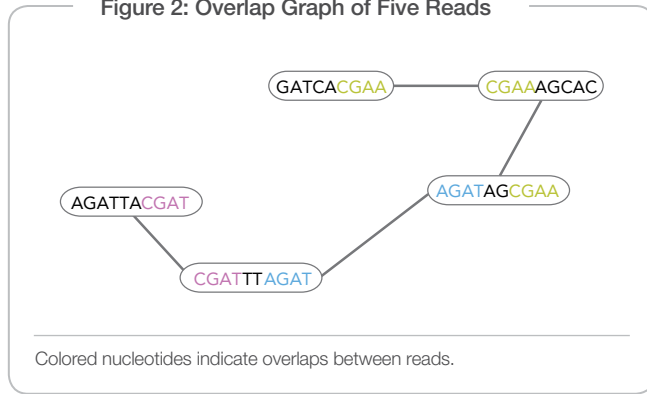


Figure 2: Overlap Graph of Five Reads

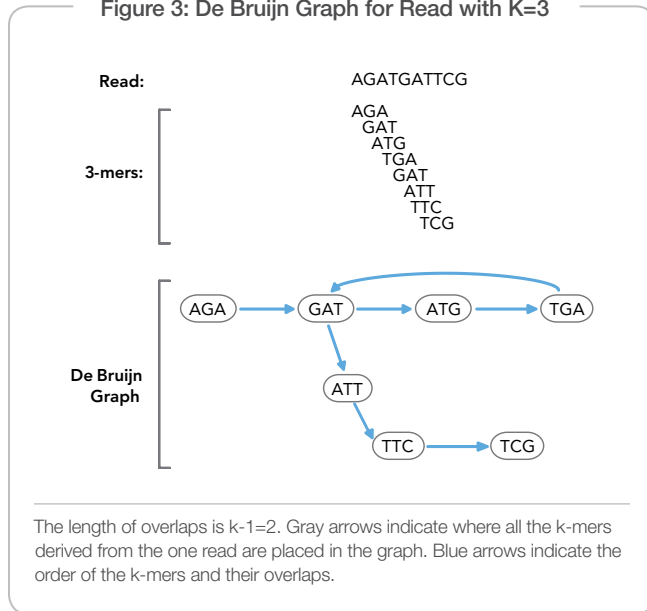


Some assemblers for next-generation sequence data use overlap graphs, but this traditional approach is computationally intensive: even a *de novo* assembly of simple organisms needs millions of reads, making the overlap graph extremely large.

De Bruijn Graphs

Because overlap graphs do not scale well with increasing numbers of reads, most assemblers for next-generation sequencing use de Bruijn graphs. De Bruijn graphs reduce the computational effort by breaking reads into smaller sequences of DNA, called k-mers, where the parameter k denotes the length in bases of these sequences. The de Bruijn graph captures overlaps of length k-1 between these k-mers and not between the actual reads (Figure 3).

Figure 3: De Bruijn Graph for Read with K=3



By reducing the entire data set down to k-mer overlaps the de Bruijn graph reduces the high redundancy in short-read data sets. The maximum efficient k-mer size for a particular assembly is determined by the read length as well as the error rate. The value of the parameter k has significant influence on the quality of the assembly. Estimates of good values can be made before the assembly, but often the optimal value

is best found by testing a small range of values. We provide more details in the "Recommendations" section.

Another attractive property of de Bruijn graphs is that repeats in the genome can be collapsed in the graph and do not lead to many spurious overlaps, although this does not mean that they can be more easily bridged or resolved. The maximum size of the de Bruijn graph is independent of sequence depth with an upper bound of 4k. Depending upon the genome being sequenced and the value of k, the de Bruijn graph may not reach the theoretical maximum, but in the presence of sequencing errors or biological variation, the memory footprint of the graph increases. Nevertheless, it has been our experience that reasonable error rates do not significantly increase the memory requirement.

A Sampling of Assemblers for Short Reads

The software package Velvet¹ was among the first assemblers for short reads and is now widely used. It implements an approach based on de Bruijn graphs, uses information from read pairs, and implements various error correction steps after building the graph. Velvet has successfully been used to assemble bacterial genomes¹.

SOAPdenovo² also implements a de Bruijn graph approach. In contrast to Velvet, error correction is performed before the actual graph is built.

The assemblers ABySS3 also uses the de Bruijn graph method. Its advantage is that it can be run in a parallel environment and thus has the potential to assemble much larger genomes. For example, Simpson et al. demonstrate the assembly of a human genome using ABySS3. SOAPdenovo also implements a parallel assembly algorithm based on de Bruijn graphs but details of this tool are not yet published.

Forge⁵ implements an overlap-layout-consensus approach with various changes to accommodate Illumina reads. It distributes the computational and memory consumption on various nodes and has therefore the potential to assemble much larger genomes, despite not being a de Bruijn graph method.

An overview of the tested assemblers is given in Table 1.

Note

The analysis presented here represents a snapshot in time of a subset of the currently available assemblers. For example, much of our analysis was performed using Velvet version 0.7.31 but several releases have occurred since we downloaded and tested this software. Assemblers evolve constantly and we anticipate that new methods will be developed to allow mammalian genomes to be more rapidly and efficiently assembled.

Comparing Assembly Outcomes

The outcome of an assembly is a set of contigs. A contig is a contiguous assembled piece of DNA sequence. Some assemblers also compute scaffolds, which is a set of contigs for which the relative orientation and distance is known. An alternative to scaffolds are supercontigs: contigs in which gaps are allowed. Gaps are usually denoted by the letter 'N' in the DNA sequence.

Table 1: Overview of Tested Assemblers

Algorithm	Description	Strength	Genomes Assembled
Velvet	De Bruijn graph based Error corrections after graph is built	Fast (~30 mins) Easy to use Larger supercontig N50	Bacterial (Ref. 1; this technical note)
SOAPdenovo	De Bruijn graph based Error correction before graph is built	Easy to use Multi-threaded mode	Panda, Bacterial (Ref. 11; this technical note)
ABYSS	De Bruijn graph based Can be run in parallel Distributed memory model (efficient)	Easy to use Largest contigs/scaffolds Best suited for large genomes	Human (Ref. 3; this technical note)
Forge	Overlap-layout-consensus method Modifications to accommodate Illumina reads	Largest contigs/supercontigs Good "long read" assembler	Bacterial (this technical note)

The following metrics are often used to compare the quality of assemblies:

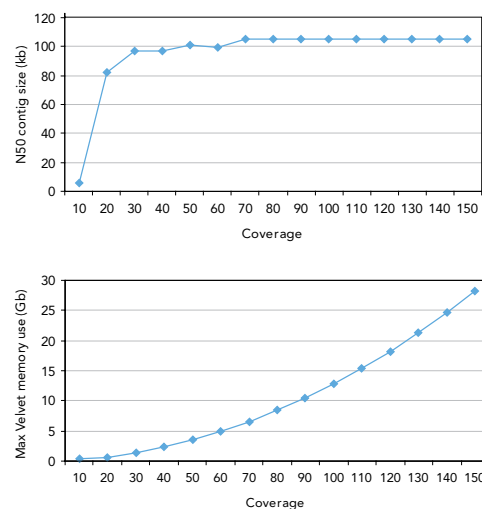
- **N50**—The contig length such that 50% of the *de novo* assembled genome lies in blocks of this size or larger. N80 or N60 are also used.
- **Genome coverage**—The percentage of bases in the reference covered by the assembled contigs. This can only be computed if a reference genome exists, and the way the assembled contigs are mapped to the reference can have a significant influence on this parameter. Popular tools for mapping include BLAST⁶, MUMmer⁷, and SSAHA⁸. Contigs aligning with large internal gaps or low confidence indicate either structural variants or misassemblies, and should be carefully monitored. If paired reads are available, mapping them back onto the contigs is another source of information. Large deviations in the mapping distances or read pairs in which only one read maps against the contig hint at misassemblies rather than structural variants.
- **Maximum/median/average contig size**—Usually calculated after removing the smallest contigs (for example removing all contigs less than 150 bp in length).

The contig sizes and their distribution are convenient and straightforward quality metrics but they do not contain all of the information needed to judge quality. For example, an assembly comprising several medium-sized contigs can be high-quality if these contigs cover the inter-repeat regions of the genome with high accuracy. On the other hand, an assembly consisting of one large contig with roughly the length of the genome being assembled is useless if the contig is incorrectly assembled.

Sequencing Parameters

The most important sequencing parameters are discussed in this section. We also provide examples of assemblies of *E. coli* data using Velvet to illustrate the influence of the parameters on the quality of the assembly. We performed several assemblies with Velvet, using different values of *k*, and determined that Velvet works best with *k*=31 for the *E. coli* data set. Note that higher values of *k* may perform better but Velvet 0.7.31 only supports values up to 31. This value of *k* was used throughout the work described here.

Figure 4: Effect of Coverage



Effect of coverage on N50 contig size and memory requirements in an *E. coli* de novo assembly.

Coverage

A high-quality de novo assembly cannot be achieved unless there is a sufficient number of error-free reads covering the entire genome. To

achieve this and thus produce a high-quality assembly, a high depth of coverage is essential. The coverage needed will depend on the organism, its genome size, and the repeat content. To give an example, the Beijing Genomics Institute sequenced the Giant Panda genome using 75 bp reads at a coverage of 50x¹¹.

To evaluate the influence of coverage on a bacterial assembly, we created simulated data sets based on the *E. coli* genome using error-free reads. We simulated 75 bp paired-end data sets with a 200 bp insert size assembled using Velvet 0.7.31 (Figure 4) and different coverages.

As can be seen from Figure 4, a coverage of more than 50x does not yield a significant improvement in terms of the contig sizes. Since Velvet's memory usage increases with sequencing depth, higher

coverage can require a significantly larger amount of memory with little or no improvement in assembly quality.

To further examine whether 50x coverage is sufficient, we used real sequencing data: a single GA lane from a 200 bp insert library of *E. coli* with 75 bp paired reads. After removing reads that did not pass the GA analysis software Failed_Chastity filter (described later) or containing Ns, the coverage is 320x. From this starting data, we randomly removed reads to generate samples with lower coverage, and assembled using Velvet 0.7.31.

Table 2 shows that the contig sizes drop as expected when the coverage drops under 50x, but contig sizes remain stable at higher cover-

Table 2: Effect of Coverage on Assembly Quality

Coverage	N50 contig size	Largest contig	Genome coverage
320x	95,313 bp	215,645 bp	99.47%
160x	95,368 bp	209,234 bp	99.72%
50x	97,333 bp	223,793 bp	99.72%
21x	35,828 bp	119,071 bp	99.38%

age. We expect similar results for genomes with sizes, base compositions, and repeat contents similar to *E. coli*. In principle, it is possible to obtain several good assemblies from a single GA lane using multiplexing to sequence several bacterial samples.

In general, the coverage threshold above which no improvement in N50 is possible depends on the size and repeat content of the genome to be sequenced. Still, even for larger genomes, it is expected that after a certain coverage is reached, adding more short-insert reads will not improve the assembly.

Read Length

The Genome Analyzer can generate paired reads with a length of 100 bp and more. There are, however, alternate technologies such as pyrosequencing that produce longer unpaired reads. In this experiment, we investigate the influence of read length on an assembly.

Chaisson et al.¹² used simulations to show that reads with more than 36 bp and 60 bp do not improve assemblies of *E. coli* and *S. cerevisiae*, respectively. They restricted their experiment to paired reads and the two abovementioned organisms with relatively simple genomes.

Table 3: Effect of Read Length

Sample	N50 contig size	Largest contig	Genome coverage
<i>E. coli</i> , 100 bp pe	132,786 bp	326,886 bp	99.87 %
<i>E. coli</i> , 400 bp sr	22,902 bp	127,976 bp	99.87 %
Chr. 20, 100 bp pe	70,744 bp	484,312 bp	92.69 %
Chr. 20, 400 bp sr	2,319 bp	22,823 bp	92.65 %

We simulated paired reads of 100 bp length (with a 400 bp fragment size) and 400 bp unpaired reads from the *E. coli* genome and human chromosome 20, respectively. All reads were simulated as error-free and at 50x base coverage. We used Velvet 0.7.31 with a k-mer size of 31 to assemble all four data sets.

Table 3 shows the contig sizes for each assembly. The assemblies using the 100 bp paired reads yield far larger contigs as compared to the 400 bp single read assemblies. This is in part due to the way a de Bruijn graph assembler such as Velvet works: the reads are split into smaller pieces, the k-mers. In this experiment, we used a k-mer size of 31 which is the maximum for Velvet 0.7.31. Therefore we compared in principle k-mers of length 31 with and without pairing information. The strong impact of read pairs as opposed to single-ended reads is evident.

To further explore the advantage of paired-end reads, we performed an assembly of the 100 bp paired reads in which we ignored the read pairing information. This means that Velvet treats the reads essentially as 100 bp single reads. Table 4 shows that the assembly quality decreases strongly when not using paired-ends. The 100 bp single-read contig sizes match surprisingly well the 400 bp assembly of both *E. coli* and chromosome 20 (Table 3). This illustrates the critical importance of read pairs for obtaining high-quality assemblies.

There are newer versions of Velvet which support, in theory, k-mers of unlimited size. In practice, the maximum k-mer size is restricted by the available RAM since longer k-mers need to be stored in a different data structure with higher memory requirements. Furthermore, for real data sets, sequencing errors and base coverage limit

the maximum k-mer size as well. In the next experiment, we repeated the 400 bp assemblies with the most recent version of Velvet, 0.7.55, and a k-mer size of 119 (Velvet allows only odd k-mer sizes). The k-mer size increased significantly the memory footprint and the RAM usage peaked at ~80 GB for the assembly of the Human chromosome. Even more RAM will be required for real (non-perfect) data and this is why k-mers of this size or larger will be of limited practical use.

Table 5 shows the results of this experiment. The N50 of resulting assemblies are still smaller than the ones obtained with 100 bp paired reads (Table 3) but are of the same order of magnitude.

These experiments show that 400 bp reads, even at high coverages, do not provide an advantage over shorter paired reads even for Human chromosomes. Longer k-mers in a de Bruijn graph assembly increase the contigs obtained from 400 bp reads but in practice there is a limit on the maximum k-mer size that can be employed due to memory requirements and sequencing errors.

Table 4: Effect of Pairing Reads

Sample (100 bp reads)	N50 contig size	Largest contig	Genome coverage
<i>E. coli</i> , paired-end	132,786 bp	326,886 bp	99.87 %
<i>E. coli</i> , single read	23,326 bp	127,976 bp	99.87 %
Chr. 20, paired-end	70,744 bp	484,312 bp	92.69 %
Chr. 20, single read	2,320 bp	22,823 bp	92.43 %

Table 8: Comparison of Contig Assembly

Software package	N50	Largest contig	Genome coverage
Velvet 0.7.31, k=31	61,802 bp	115,666 bp	99.72%
ABYSS 1.0.8, k=42	45,171 bp	140,706 bp	99.64%
Forge 1.0, k=15	70,447 bp	444,471 bp	99.4%
SOAPdenovo 1.0	3,026 bp	20,258 bp	99.51%

Error Rate

We created a series of simulated data sets based on the E. coli genome to investigate the influence of sequencing errors. We simulated different error rates in sequencing reads, and used Velvet 0.7.31 to perform an assembly at 150x coverage (Figure 5).

The results for an error rate less than ~4% match the contig sizes we obtained using real E. coli reads. There is a sharp drop in contig sizes as soon as the error rates surpass 4%. This error rate is well above the average error rate for a good GA run, indicating that sequencing error does not usually limit the assembly quality (as shown in Table 2).

Testing Assemblers

Comparison of Assemblers on a Bacterial Genome

We compared currently available assemblers for Illumina reads using a single GA lane from 200 bp insert library of E. coli with 75 bp paired reads, down-sampled to 50x coverage.

It is difficult to compare the results of assemblers directly, since they produce different outputs: ABySS computes only contigs without gaps whereas Velvet, Forge, and SOAPdenovo compute sets of contigs, “sequence-connected-supercontigs (SCSS)” in Velvet, supercontigs in Forge and SOAPdenovo. To make the results comparable, we generated two tables.

Table 8 shows a comparison based on the contig sizes, where supercontigs/scaffolds for Velvet, Forge, and SOAPdenovo were split whenever at least one gap character ('N') occurs. Table 9 shows a comparison based on the supercontig/scaffold sizes. Since ABySS does not generate supercontigs, it is omitted from this table.

Comparing supercontig/scaffolds, Velvet produced the largest N50 statistic in the E. coli assembly using short inserts, but Forge and SOAPdenovo computed assemblies of similar quality and contig size distribution (Table 9). In fact, Forge produced a much longer scaffold, but took ~50 times longer than Velvet to run (~30 minutes for Velvet versus ~24 hours for Forge). We executed Velvet on a machine with 60 GB of RAM and 16 CPUs with 2.4 GHz. Note that Velvet does not

Table 9: Comparison of Supercontig/Scaffold Assembly

Software package	N50	Largest scaffold	Genome coverage
Velvet 0.7.31, k=31	97,333 bp	223,793 bp	99.72%
Forge 1.0,k=15	82,595 bp	482,322 bp	99.4%
SOAPdenovo 1.0	95,472 bp	223,876 bp	98.61%

make use of multiple CPUs. Forge was executed in a parallel fashion on a cluster with 20 CPUs and 4 GB RAM per CPU. Most computing time was spent on building the scaffold and traversing the overlap graph.

Assembly of Larger Genomes

Since ABySS uses parallelization and de Bruijn graphs, it can be used for de novo assembly of larger genomes. The other assemblers have limitations that become prohibitive when assembling a large genome, such as a mammalian genome.

We tested ABySS using reads from a Yoruba male (child of the individual published in Bentley et al.¹⁰) with the HapMap reference number NA18506. The data set consisted of 100 bp paired reads sampled at 30x coverage with an insert size of 600 bp. We first assembled chromosomes 1 and 20, to serve as medium-sized genome test cases. After that, we assembled the whole human genome.

Medium-Sized Genome Assemblies

We aligned all reads from the Yoruba male against the NCBI human reference genome and used reads aligning to chromosome 1 and 20 to assemble both chromosomes. These chromosomes have a size of 247 Mb and 62 Mb respectively and thus fall into the gap in genome size between E. coli and mammalian genomes.

After assembly, we discarded contigs with less than 100 bp to make the results comparable to previously published data³ (Table 10).

Table 10: Assembly of Human Chromosome 1 (K=55) and Chromosome 20 (K=62) by AByss 1.0.8

Chromosome	Size (bp)	N50 contig size (bp)	Largest contig (bp)	Bases in contigs (Mb)
Chr. 20	62,435,965	4,743	48,538	64
Chr. 1	247,199,719	2,879	32,516	197

ABySS assembles both data sets into reasonably sized contigs. Contigs of this size can be useful for characterizing Single Nucleotide Polymorphisms (SNPs) and small to medium-sized structural variants. Further improvements in contig size can be obtained by adding long-insert libraries.

Whole Human Genome Assembly

We also performed a prototype assembly of the whole genome. The first stage of assembly, which was performed without the read pairing information, took ~20 hours on a cluster with 150 cores. Joining and error correcting the resulting contigs required an additional three days.

Due to the high repeat content and the small insert size, this assembly is highly fragmented. The largest contig had a size of 27,534 bp, but the N50 is much lower than the N50 that we achieved for chromosome 1. Whole-genome assembly of a mammalian genome with ABySS may therefore provide a starting point, but requires significant hands-on assembly afterwards. However, we expect that assemblies of whole mammalian genomes will improve with further improvements in algorithms and the application of long-insert libraries.

