

In contrast, CE-based Sanger sequencing requires genomic DNA to be fragmented first and cloned into either bacterial artificial chromosomes (BACs) or yeast artificial chromosomes (YACs). Then, each BAC/YAC must be further subcloned into a sequencing vector and transformed into the appropriate microbial host. Template DNA is then purified from individual colonies or plaques prior to sequencing. This process can take days or even weeks to complete, depending on the size of the genome.

Data Analysis Algorithms

Along with sample preparation, data analysis is an important factor to consider for sequencing applications. A range of data analysis algorithms are available that perform specific tasks related to a given application. Some applications, such as *de novo* sequencing, require specialized assembly of sequencing reads. Other applications, such as RNA-Seq, require algorithms that quantify read counts to provide information about gene expression levels. A number of algorithms exist to address the needs of each application. While some of these are commercially available from software vendors, many are freely available open-source algorithms from academic institutions. Illumina collaborates closely with commercial and academic software developers to create an ecosystem of data analysis tools that address the needs of various research objectives¹.

Whole-Genome Sequencing

Until recently, sequencing an entire genome was a major endeavor. Even for fairly compact viral genomes with overlapping genes and few to no repetitive regions, whole-genome sequencing using CE-based Sanger technology requires a significant commitment of time and resources. For example, *de novo* whole-genome sequencing of vaccinia virus—a large and complex DNA virus (~200 kb)—using CE-based methods would involve roughly 4,000 sequencing reactions (assuming 10x coverage and 500 bp read lengths), each conducted in separate tubes or wells. Using NGS technology, the same sequencing project, including library preparation, can be completed in just a few days with one sequencing run at 30x coverage or greater.

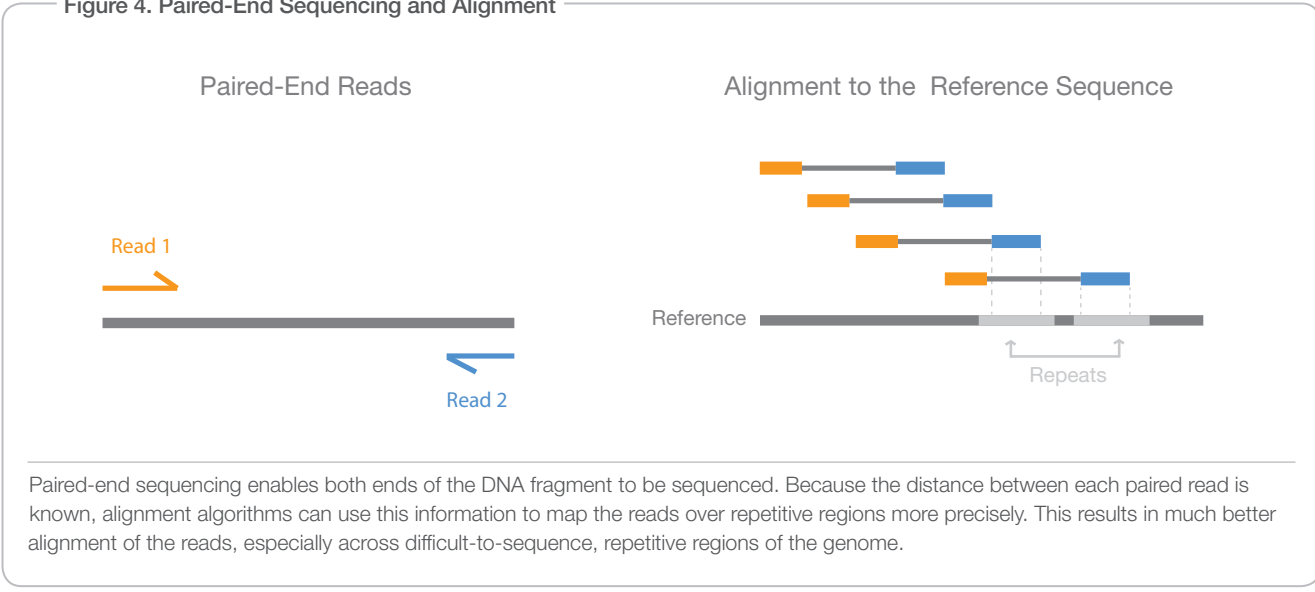
Sequencing Small Genomes

The ability of NGS platforms to produce a large volume of data in a short period of time makes it a powerful tool for whole-genome sequencing in a laboratory environment. While NGS is commonly associated with sequencing large genomes, especially human genomes, the scalability of the technology makes it just as useful for small viral or bacterial genomes. This was demonstrated during the recent *E. coli* (EHEC) outbreak in Europe, which prompted a rapid scientific response. Using the latest NGS systems, researchers were able to quickly generate a high-quality whole-genome sequence of the bacterial strain², enabling them to better understand the genetic mutations conferring the increased virulence.

One challenge associated with sequencing small genomes is the lack of reference genomes available for most species. This means that whole-genome sequencing must often be done *de novo*, where the reads are assembled without aligning to a reference sequence. The coverage quality of a *de novo* sequencing data set depends upon the quality of the contigs, or continuous sequences generated by aligning overlapping sequencing reads. The size and continuity of the contigs will affect the number of

gaps present in the data. A problem for *de novo* sequencing is that the short read lengths generated by NGS can lead to a higher number of gaps, regions where no reads align, resulting in greater fragmentation and smaller contigs—poorer data quality. This is especially true for regions of the genome containing repetitive sequence elements. To overcome this challenge, some NGS platforms offer paired-end (PE) sequencing protocols (Figure 4), where both ends of a DNA fragment are sequenced, as opposed to single-read sequencing where only one end is sequenced. Paired-end reads result in superior alignment across regions containing repetitive sequences and produce longer contigs for *de novo* sequencing by filling gaps in the consensus sequence, resulting in complete overall coverage.

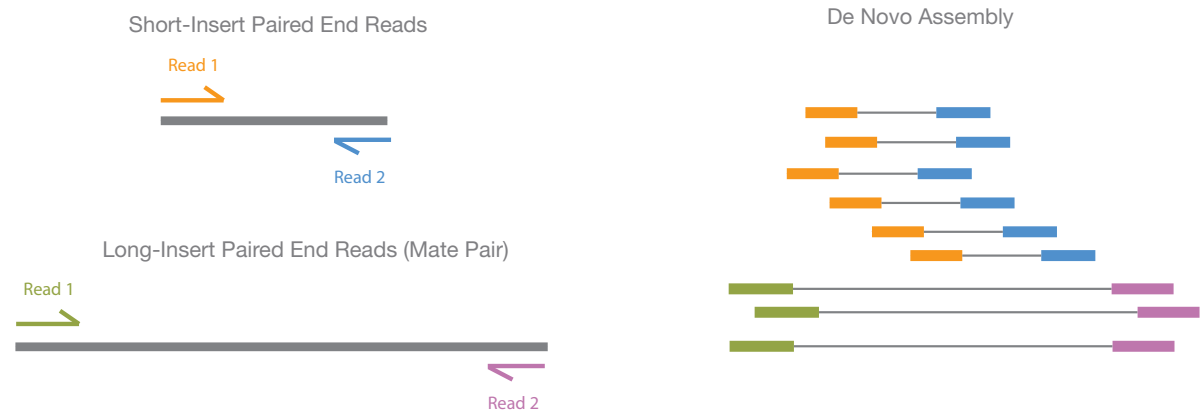
Figure 4. Paired-End Sequencing and Alignment



Another important factor in generating high quality *de novo* sequences is the diversity of insert (DNA fragment) sizes in the library. Using longer inserts provides the highest fragment diversity relative to starting input material, yielding more uniform sequencing coverage. When long inserts are prepared for pair-end sequencing, a mate pair library is generated. These can include insert sizes ranging from 2 to 5 kb, optimal for *de novo* assembly applications, including both genome scaffold generation and genome finishing. In general, libraries with larger insert sizes will result in less fragmented assemblies and larger contigs. Combining short-insert paired-end and long-insert mate pair sequencing is the most powerful approach for maximal coverage across the genome (Figure 5). The combination of insert sizes enables detection of the widest range of structural variant types and is essential for accurately identifying more complex rearrangements, which results in a higher quality assembly. The short-insert reads sequenced at higher depths can fill in gaps not covered by the long inserts, which are often sequenced at lower read depths.

In parallel with NGS technological improvements, many algorithmic advances have been made in *de novo* sequence assemblers for short-read data. Researchers can perform high-quality *de novo* assembly using NGS reads and publicly available short-read assemblers. In many instances, existing computer resources in the laboratory are enough to perform *de novo* assemblies. The *E. coli* genome can be assembled in as little as 15 minutes using a 32-bit Windows desktop computer with 32 GB of RAM.

Figure 5. De Novo Assembly with Mate Pairs



Using a combination of short and long insert sizes with paired-end sequencing results in maximal coverage of the genome for de novo assembly. Because larger inserts can pair reads across greater distances, they provide a better ability to read through highly repetitive sequences and regions where large structural rearrangements have occurred. Shorter inserts sequenced at higher depths can fill in gaps missed by larger inserts sequenced at lower depths. Thus a diverse library of short and long inserts results in better de novo assembly, leading to fewer gaps, larger contigs, and greater accuracy of the final consensus sequence.

Targeted Sequencing

With targeted sequencing, only a subset of genes or defined regions in a genome are sequenced, allowing researchers to focus time, expenses, and data storage on the regions of the genome in which they are most interested. This approach is typically used to sequence many individuals to discover, screen, or validate genetic variation within a population. The ability to pool samples and obtain high sequence coverage during a single run allows NGS to identify rarer variants that are missed, or too expensive to identify, using CE-based sequencing approaches. There are two different methods for making libraries for targeted sequencing and resequencing projects—target enrichment and amplicon generation (Table 3).

Table 3: Targeted Sequencing Sample Prep at a Glance

CE-based Sanger Sequencing	Next-Generation Sequencing
Library preparation more involved—each sample must contain a single template, either from a single PCR purified from single bacterial colonies	Library preparation more streamlined—each sample can be a population and does not require clonal purification
Suitable for sequencing amplicons and clone checking	Suitable for sequencing amplicons and clone checking
Complete within days to weeks, depending upon the size of the genome being sequenced	Complete within hours

GGAATGATAACAGTAAACACACTTCTGTAAACCTTAAGATTACTTGATCCACTGATCAACGTACCCTAACGAAACGTATCAATTGAGACTAAATATTAACGTACCATTAAGAGCTACCGTCTCTGTAAACCTTAAGATTACTTGATCCACTGATCAACG
TCAACGTACCCTAACGAAACGTATCAATTGAGATTACTTGATCCACTGATCAACGTACCCTAACGAAACGTATCAATTGAGACTAAATATTAACGTACCATTAAGAGCTACCGTCTCTGTAAACCTTAAGATTACTTGATCCACTGATCAACG
GACGAAAAGAATGATAACAGTAAACACACTTCTGTAAACCTTAAGATTACTTGATCCACTGATCAACGTACCCTAACGAAACGTATCAATTGAGACTAAATATTAACGTACCATTAAGAGCTACCGTCTCTGTAAACCTTAAGATTACTTGAT
ACCTACCATTAAGAGCTACCGTCTCTGTAAACCTTAAGATTACTTGATCCACTGATCAACGTACCCTAACGAAACGTATCAATTGAGACTAAATATTAACGTACCATTAAGAGCTACCGTCTCTGTAAACCTTAAGATTACTTGAT
AGAATTGATAACAGTAAACACACTTCTGTAAACCTTAAGATTACTTGATCCACTGATCAACGTACCCTAACGAAACGTATCAATTGAGACTAAATATTAACGTACCATTAAGAGCTACCGTCTCTGTAAACCTTAAGATTACTTGATCAACG
GATTACTTGATCCACTGATCAACGTACCCTAACGAAACGTATCAATTGAGACTAAATATTAACGTACCATTAAGAGCTACCGTCTCTGTAAACCTTAAGATTACTTGATCCACTGATCAACGTACCCTAACGAAACGTATCAATTGAGACTAAATATTAACG
CGTATCAATTGAGACTAAATATTAACGTACCATTAAGAGCTGTTAAACCTTAAGATTACTTGATCCACTGATCAACGTACCCTAACGAAACGTATCAATTGAGACTAAATATTAACGTACCATTAAGAGCTACCGTCTCTGTAAACCTTAAGATTACTTGAT

From Innovation to Publication

The advent of NGS has enabled researchers to study biological systems at a level never before possible. As the technology has evolved, an increasing number of innovative sample preparation methods and data analysis algorithms have enabled a broad range of scientific applications. Researchers are making fascinating discoveries in a number of biological fields, unlocking answers never before possible. As a result, there has been an explosion in the number of scientific publications. Illumina sequencing alone has resulted in thousands of peer-reviewed publications. Selected recent examples are listed below.

Whole-Genome Sequencing

- Srivatsan A, et al. (2008) High-precision, whole-genome sequencing of laboratory strains facilitates genetic studies. *PLoS Genet* 4: e1000139.
- Rasmussen M, et al. (2010) Ancient human genome sequence of an extinct Palaeo-Eskimo. *Nature* 463:757-762.
- Li R, et al. (2010) The sequence and *de novo* assembly of the giant panda genome. *Nature* 463:311-317.
- Pelak K, et al. (2010) The characterization of twenty sequenced human genomes. *PLoS Genet* 6:e1001111.

Targeted Resequencing

- Ram JL, Karim AS, Sendler ED, Kato (2011) Strategy for microbiome analysis using 16S rRNA gene sequence analysis on the Illumina sequencing platform. *Syst Biol Reprod Med.* 57(3):117-8.
- McEllistrem M.C. (2009) Genetic diversity of the pneumococcal capsule: implications for molecular-based serotyping. *Future Microbiol* 4:857-865.
- Lo YMD, Chiu RWK. (2009) Next-generation sequencing of plasma/serum DNA: an emerging research and molecular diagnostic tool. *Clin. Chem* 55:607-608.
- Robinson PN (2010) Whole-exome sequencing for finding *de novo* mutations in sporadic mental retardation. *Genome Biol* 11:144.
- Araya CL, Fowler DM (2011) Deep mutational scanning: assessing protein function on a massive scale. *Trends Biotechnol.* doi:10.1016/j.tibtech.2011.04.003.

References

1. www.illumina.com/software/illumina_connect.ilmn
2. www.illumina.com/Documents/products/appnotes/appnote_miseq_ecoli.pdf
3. www.illumina.com/Documents/products/appnotes/appnote_miseq_denovo.pdf
4. Bomar L, Maltz M, Colston S, and Graf J. (2011) Directed culturing of microorganisms using metatranscriptomics. *mBio.* 2:e00012-11

FOR RESEARCH USE ONLY

© 2011-2012 Illumina, Inc. All rights reserved.
Illumina, illuminaDx, BaseSpace, BeadArray, BeadXpress, cBot, CSpPro, DASL, DesignStudio, Eco, GAlIx, Genetic Energy, Genome Analyzer, GenomeStudio, GoldenGate, HiScan, HiSeq, Infinium, iSelect, MiSeq, Nextera, Sentrix, SeqMonitor, Solexa, TruSeq, VeraCode, the pumpkin orange color, and the Genetic Energy streaming bases design are trademarks or registered trademarks of Illumina, Inc. All other brands and names contained herein are the property of their respective owners.
Pub No. 770-2012-008 Current as of 28 February 2012

