

Table of Contents

I. Welcome to Next-Generation Sequencing	3
a. The Evolution of Genomic Science	3
b. The Basics of NGS Chemistry	4
c. Advances in Sequencing Technology	5
Paired-End Sequencing	5
Tunable Coverage and Unlimited Dynamic Range	6
Multiplexing	6
Advances in Library Preparation	7
Flexible, Scalable Instrumentation	7
II. NGS Methods	8
a. Genomics	8
Whole-Genome Sequencing	8
Exome Sequencing	9
<i>De Novo</i> Sequencing	9
Targeted Sequencing	10
b. Transcriptomics	11
Total RNA and mRNA Sequencing	11
Targeted RNA Sequencing	12
Small RNA and Noncoding RNA Sequencing	12
c. Epigenomics	12
Methylation Sequencing	12
ChIP Sequencing	12
Ribosome Profiling	12
III. Illumina DNA-to-Data NGS Solutions	13
a. The Illumina NGS Workflow	13
b. Integrated Data Analysis	13
IV. Glossary	14
V. References	15

I. Welcome to Next-Generation Sequencing

a. The Evolution of Genomic Science

DNA sequencing has come a long way since the days of two-dimensional chromatography in the 1970s. With the advent of the Sanger chain termination method¹ in 1977, scientists gained the ability to sequence DNA in a reliable, reproducible manner. A decade later, Applied Biosystems introduced the first automated, capillary electrophoresis (CE) based sequencing instruments—the AB370 in 1987 and the AB3730xl in 1998—instruments that became the primary workhorses for the NIH-led and Celera-led Human Genome Projects.² While these “first-generation” instruments were considered high throughput for their time, the Genome Analyzer emerged in 2005 and took sequencing runs from 84 kilobase (kb) per run to 1 gigabase (Gb) per run.³ The short read, massively parallel sequencing technique was a fundamentally different approach that revolutionized sequencing capabilities and launched the “next-generation” in genomic science. From that point forward, the data output of next-generation sequencing (NGS) has outpaced Moore’s law—more than doubling each year (Figure 1).

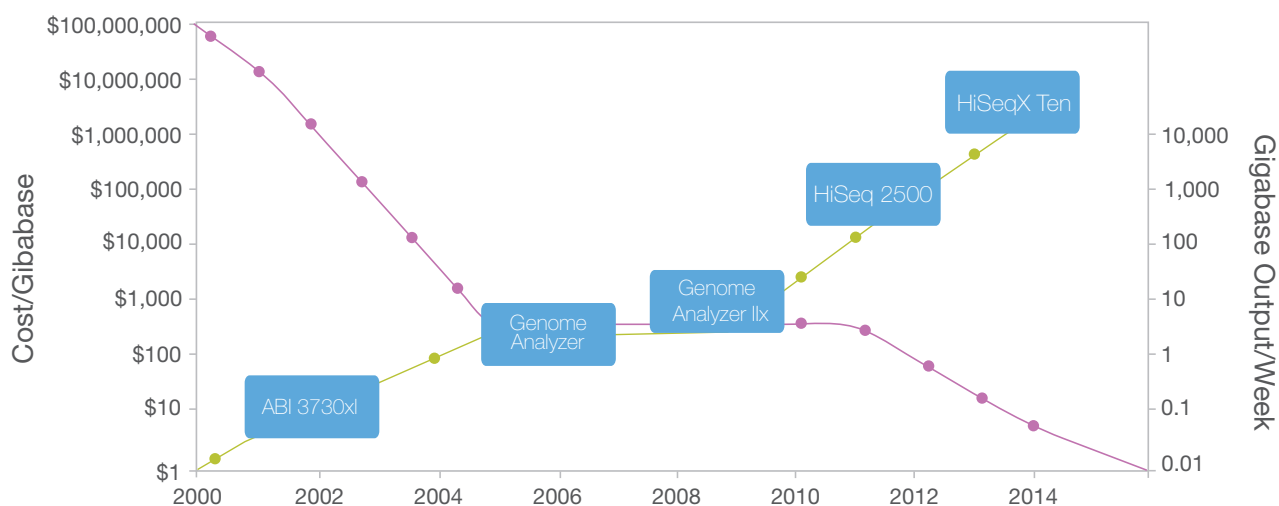


Figure 1: Sequencing Cost and Data Output Since 2000—The dramatic rise of data output and concurrent falling cost of sequencing since 2000. The Y-axes on both sides of the graph are logarithmic.

In 2005, with the Genome Analyzer, a single sequencing run could produce roughly one gigabase of data. By 2014, the rate climbed to a 1.8 terabases of data in a single sequencing run—an astounding 1000x increase. It is remarkable to reflect on the fact that the first human genome, famously copublished in *Science* and *Nature* in 2001, required 15 years to sequence and cost nearly 3 billion dollars. In contrast, the HiSeqX™ Ten, released in 2014, can sequence over 45 human genomes in a single day for approximately \$1000 each (Figure 2).⁴

Beyond the massive increase in data output, the introduction of NGS technology has transformed the way scientists think about genetic information. The \$1000 dollar genome enables population-scale sequencing and establishes the foundation for personalized genomic medicine as part of standard medical care. Researchers can now analyze thousands to tens of thousands of samples in a single year. As Eric Lander, founding director of the Broad Institute of MIT and Harvard and principle leader of the Human Genome Project, states, “The rate of progress is stunning. As costs continue to come down, we are entering a period where we are going to be able to get the complete catalog of disease genes. This will allow us to look at thousands of people and see the differences among them, to discover critical genes that cause cancer, autism, heart disease, or schizophrenia.”⁵

Human Genomes Sequenced Annually

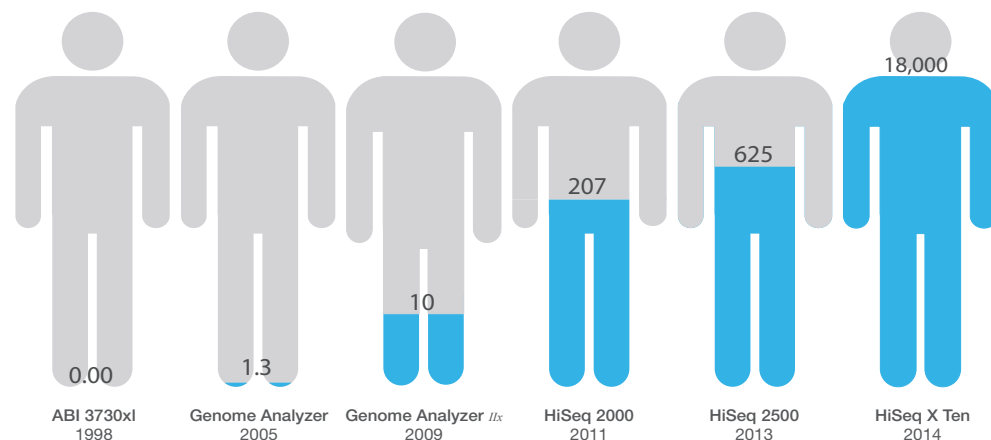


Figure 2: Human Genome Sequencing Over the Decades—The capacity to sequence all 3.2 billion bases of the human genome (at 30× coverage) has increased exponentially since the 1990s. In 2005, with the introduction of the Illumina Genome Analyzer System, 1.3 human genomes could be sequenced annually. Nearly 10 years later, with the Illumina HiSeq X Ten fleet of sequencing systems, the number has climbed to 18,000 human genomes a year.

b. The Basics of NGS Chemistry

In principle, the concept behind NGS technology is similar to CE sequencing—DNA polymerase catalyzes the incorporation of fluorescently labeled deoxyribonucleotide triphosphates (dNTPs) into a DNA template strand during sequential cycles of DNA synthesis. During each cycle, at the point of incorporation, the nucleotides are identified by fluorophore excitation. The critical difference is that, instead of sequencing a single DNA fragment, NGS extends this process across millions of fragments in a massively parallel fashion. Illumina sequencing by synthesis (SBS) chemistry is the most widely adopted chemistry in the industry and delivers the highest accuracy, the highest yield of error-free reads, and the highest percentage of base calls above Q30.⁶⁻⁸ The Illumina NGS workflows include 4 basic steps (Figure 3):

- 1. Library Preparation**—The sequencing library is prepared by random fragmentation of the DNA or cDNA sample, followed by 5' and 3' adapter ligation. Alternatively, “tagmentation” combines the fragmentation and ligation reactions into a single step that greatly increases the efficiency of the library preparation process.⁹ Adapter-ligated fragments are then PCR amplified and gel purified.
- 2. Cluster Generation**—For cluster generation, the library is loaded into a flow cell where fragments are captured on a lawn of surface-bound oligos complementary to the library adapters. Each fragment is then amplified into distinct, clonal clusters through bridge amplification. When cluster generation is complete, the templates are ready for sequencing.
- 3. Sequencing**—Illumina SBS technology utilizes a proprietary reversible terminator-based method that detects single bases as they are incorporated into DNA template strands. As all 4 reversible terminator-bound dNTPs are present during each sequencing cycle, natural competition minimizes incorporation bias and greatly reduces raw error rates compared to other technologies.^{6,7} The result is highly accurate base-by-base sequencing that virtually eliminates sequence-context-specific errors, even within repetitive sequence regions and homopolymers.
- 4. Data Analysis**—During data analysis and alignment, the newly identified sequence reads are then aligned to a reference genome. Following alignment, many variations of analysis are possible such as single nucleotide polymorphism (SNP) or insertion-deletion (indel) identification, read counting for RNA methods, phylogenetic or metagenomic analysis, and more.

A detailed animation of SBS sequencing is available at www.illumina.com/SBSvideo.

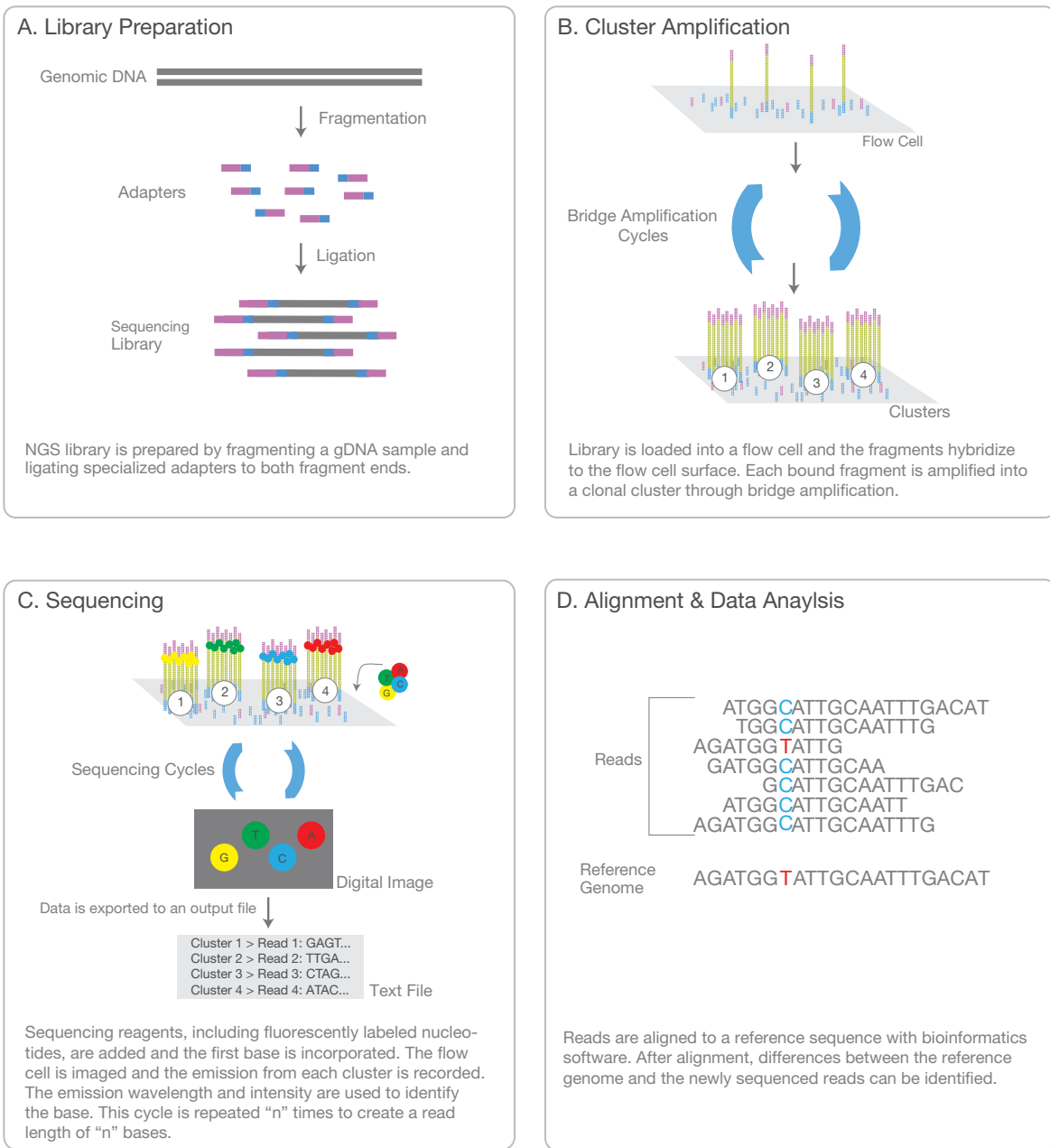


Figure 3: Next-Generation Sequencing Chemistry Overview.

c. Advances in Sequencing Technology

Paired-End Sequencing

A major advance in NGS technology occurred with the development of paired-end (PE) sequencing (Figure 4). PE sequencing involves sequencing both ends of the DNA fragments in a sequencing library and aligning the forward and reverse reads as read pairs. In addition to producing twice the number of reads for the same time and effort in library preparation, sequences aligned as read pairs enable more accurate read alignment and the ability to detect indels, which is simply not possible with single-read data.⁸ Analysis of differential read-pair spacing also allows removal of PCR duplicates, a common artifact resulting from PCR amplification during library preparation.

coverage of traditionally challenging areas such as high AT/GC-rich regions, promoters, and homopolymeric regions.¹¹ To see a complete list of Illumina library preparation kits, visit support.illumina.com/sequencing/kits.html.

To advance the process even further, Illumina has combined the precision of digital microfluidics with its ease-of-use principles to create NeoPrep™ Library Prep System—a complete, fully automated library preparation instrument. Automation of library preparation will reduce opportunities for error, increase reproducibility, and reduce the amount of hands-on time required for a process that is often a bottleneck in the sequencing workflow. For more information on library prep automation developments, visit www.illumina.com/systems.html.

Multiplexing

In addition to the rise of data output per run, the sample throughput per run in NGS has also increased over time. Multiplexing allows large numbers of libraries to be pooled and sequenced simultaneously during a single sequencing run (Figure 5). With multiplexed libraries, unique index sequences are added to each DNA fragment during library preparation so that each read can be identified and sorted before final data analysis. With PE sequencing and multiplexing, NGS has dramatically reduced the time to data for multi-sample studies and enabled researchers to go from experiment to data faster and easier than ever before.

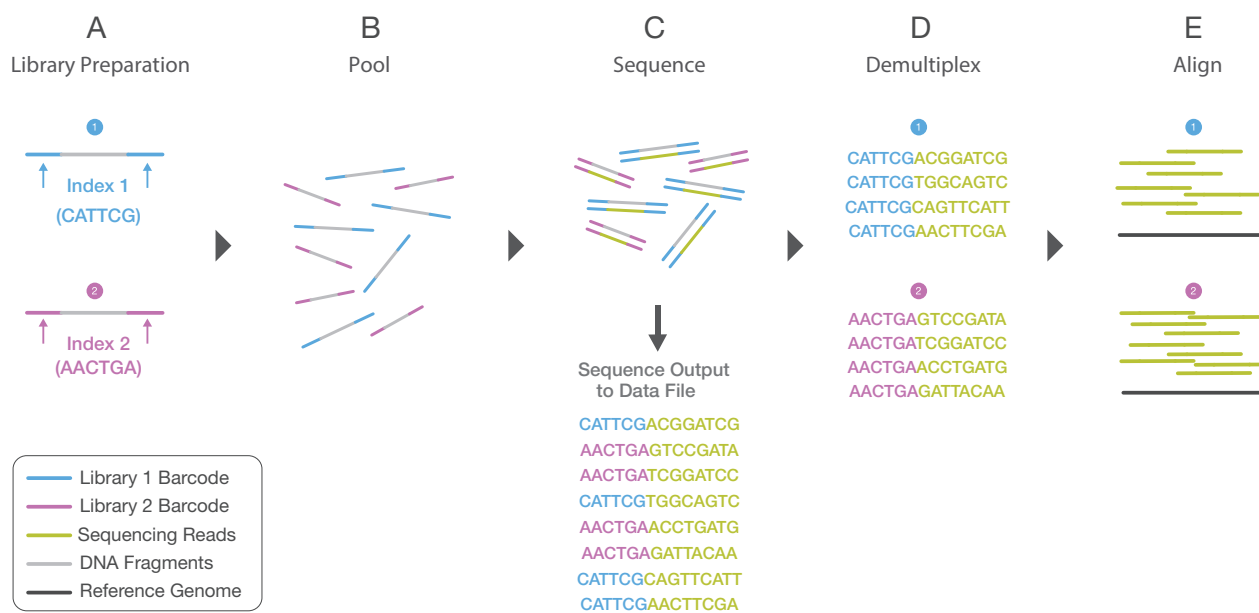


Figure 5: Library Multiplexing Overview.

- Two distinct libraries are attached to unique index sequences. Index sequences are attached during library preparation.
- Libraries are pooled together and loaded into the same flow cell lane.
- Libraries are sequenced together during a single instrument run. All sequences are exported to a single output file.
- A demultiplexing algorithm sorts the reads into different files according to their indexes.
- Each set of reads is aligned to the appropriate reference sequence.

Flexible, Scalable Instrumentation

While the latest NGS platforms can produce massive data output, NGS technology is also highly flexible and scalable. Sequencing systems are available for every method and scale of study, from small laboratories to large genome centers (Figure 6). Illumina NGS instruments range from the benchtop MiniSeq® System, with output ranging from 1.8–7.5 Gb for targeted sequencing studies, to the colossal HiSeq X Ten fleet, which can generate an impressive 16–18 Tb per run* for population-scale studies.

* With the full suite of 10 HiSeq X Systems.

Flexible run configurations are also engineered into the design of Illumina NGS sequencers. For example, the HiSeq® 2500 System offers 2 run modes and single or dual flow cell sequencing while the NextSeq® Series offers 2 flow cell types to accommodate different throughput requirements. The HiSeq 3000/4000 Series uses the same patterned flow cell technology as the HiSeq X instruments for cost-effective production-scale sequencing. This flexibility allows researchers to configure runs tailored to their specific study requirements, with the instrument of their choice. For an in-depth comparison of Illumina platforms, visit www.illumina.com/systems/sequencing.html.

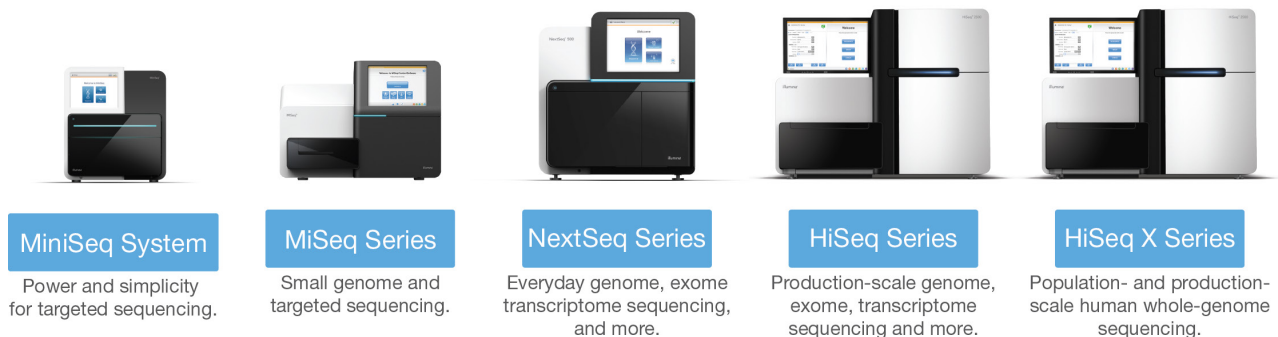


Figure 6: Sequencing Systems for Every Scale.

II. NGS Methods

Next-generation sequencing platforms enable a wide variety of methods, allowing researchers to ask virtually any question related to the genome, transcriptome, or epigenome of any organism. Sequencing methods differ primarily by how the DNA or RNA samples are obtained (eg, organism, tissue type, normal vs. affected, experimental conditions) and by the data analysis options used. After the sequencing libraries are prepared, the actual sequencing stage remains fundamentally the same regardless of the method. There are a number of standard library preparation kits that offer protocols for whole-genome sequencing, mRNA-Seq, targeted sequencing (such as exome sequencing or 16S sequencing), custom-selected regions, protein-binding regions, and more. Although the number of NGS methods is constantly growing, a brief overview of the most common methods is presented here.

a. Genomics

Whole-Genome Sequencing

Microarray-based, genome-wide association studies (GWAS) have been the most common approach for identifying disease associations across the whole genome. While GWAS microarrays can interrogate over 4 million markers per sample, the most comprehensive method of interrogating the 3.2 billion bases of the human genome is with whole-genome sequencing (WGS). The rapid drop in sequencing cost and the ability of WGS to rapidly produce large volumes of data make it a powerful tool for genomics research. While WGS is commonly associated with sequencing human genomes, the scalable, flexible nature of the technology makes it equally useful for sequencing any species, such as agriculturally important livestock, plant genomes, or disease-related microbial genomes. This broad utility was demonstrated during the recent *E. coli* outbreak in Europe in 2011, which prompted a rapid scientific response. Using the latest NGS systems, researchers quickly sequenced the bacterial strain, enabling them to better track the origins and transmission of the outbreak as well as identify genetic mutations conferring the increased virulence.¹²

Exome Sequencing

Perhaps the most widely used targeted sequencing method is exome sequencing. The exome represents less than 2% of the human genome, but contains a majority of known disease-causing variants, making whole-exome sequencing a cost-effective alternative to whole-genome sequencing. With exome sequencing, the protein-coding portion of the genome is selectively captured and sequenced. It can efficiently identify variants across a wide range of applications, including population genetics, genetic disease, and cancer studies.

De Novo Sequencing

De novo sequencing refers to sequencing a novel genome where there is no reference sequence available for alignment. Sequence reads are assembled as contigs and the coverage quality of *de novo* sequence data depends on the size and continuity of the contigs (ie, the number of gaps in the data). Another important factor in generating high-quality *de novo* sequences is the diversity of insert sizes included in the library. Combining short-insert paired-end and long-insert mate pair sequences is the most powerful approach for maximal coverage across the genome (Figure 7). The combination of insert sizes enables detection of the widest range of structural variant types and is essential for accurately identifying more complex rearrangements. 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. Therefore, using a combined approach results in higher-quality assemblies.

In parallel with NGS technology improvements, many algorithmic advances have emerged in sequence assemblers for short-read data. Researchers can perform high-quality *de novo* assembly using NGS reads and publicly available short-read assembly tools. In many instances, existing computer resources in the laboratory are enough to perform *de novo* assemblies. For example, 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.

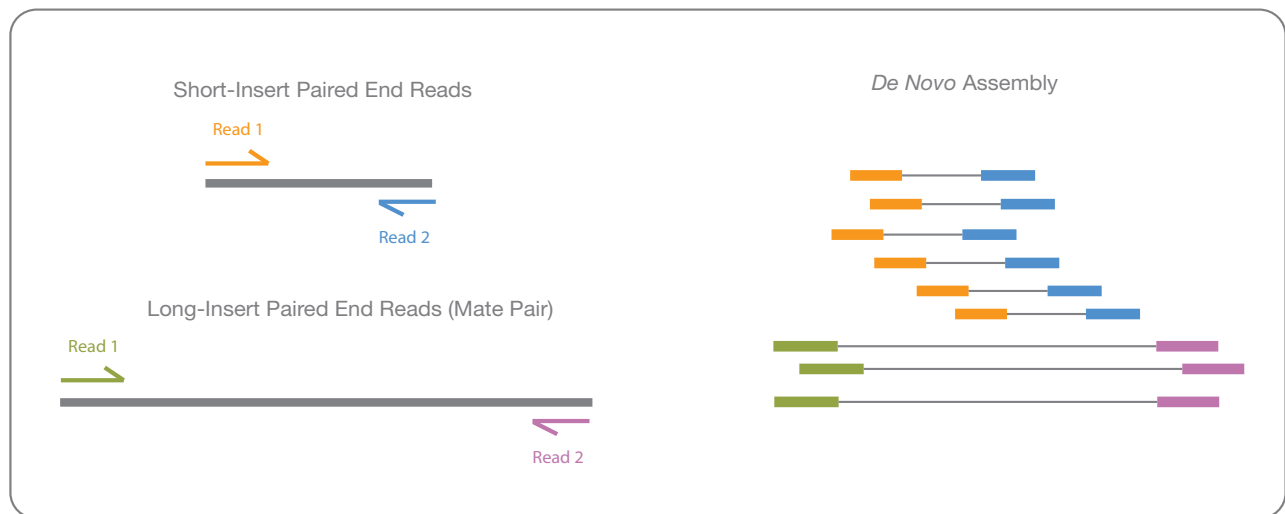


Figure 7: Mate Pairs and *De Novo* Assembly—Using a combination of short and long insert sizes with paired-end sequencing results in maximal coverage of the genome for *de novo* assembly. Because larger inserts can pair reads across greater distances, they provide a better ability to read through highly 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, a subset of genes or regions of the genome are isolated and sequenced. Targeted sequencing allows researchers to focus time, expenses, and data analysis on specific areas of interest and enables sequencing at much higher coverage levels. For example, a typical WGS study achieves coverage levels of 30x–50x per genome, while a targeted resequencing project can easily cover the target region at 500x–1000x or higher. This higher coverage allows researchers to identify rare variants—variants that would be too rare and too expensive to identify with WGS or CE-based sequencing.

Targeted sequencing panels can be purchased with fixed, preselected content or can be custom designed. A wide variety of targeted sequencing library prep kits are available, including kits with probe sets focused on specific areas of interest such as cancer, cardiomyopathy, autism, or custom probe sets (Table 2). With custom designs, researchers can target regions of the genome relevant to their specific research interests. Custom targeted sequencing is ideal for examining genes in specific pathways, or for follow-up studies from GWAS or WGS. Illumina currently supports 2 methods for targeted sequencing—target enrichment and amplicon generation (Figure 8).

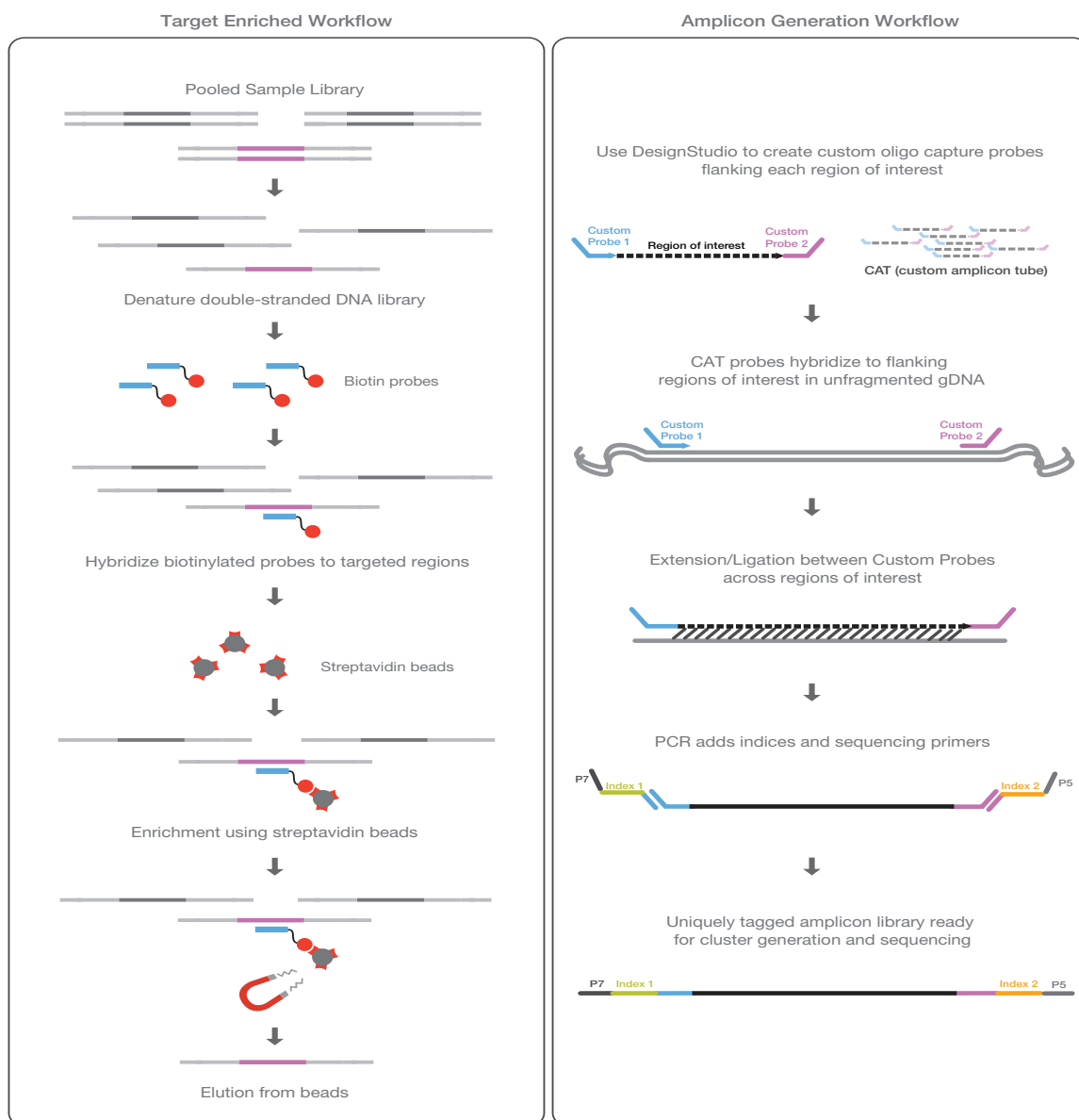


Figure 8: Target Enrichment and Amplicon Generation Workflows.

With target enrichment, specific regions of interest are captured by hybridization to biotinylated probes, then isolated by magnetic pulldown. Target enrichment captures between 10 kb–62 Mb regions depending on the library prep kit parameters. The second method, amplicon sequencing, involves the amplification and purification of regions of interest using highly multiplexed PCR oligo sets. Amplicon sequencing allows researchers to sequence 16–1536 targets at a time, spanning 2.4–652.8 kb of total content, depending on the library prep kit used. This highly multiplexed approach enables a wide range of applications for the discovery, validation, or screening of genetic variants. Amplicon sequencing is particularly useful for the discovery of rare somatic mutations in complex samples (eg, cancerous tumors mixed with germline DNA).^{13,14} Another common amplicon application is sequencing the bacterial 16S rRNA gene across multiple species, a widely used method for phylogeny and taxonomy studies, particularly in diverse metagenomic samples.¹⁵

For more information on Illumina targeted, WGS, exome, or *de novo* sequencing solutions, visit www.illumina.com/applications/sequencing/dna_sequencing.html.

b. Transcriptomics

Library preparation methods for RNA sequencing (RNA-Seq) typically begin with total RNA sample preparation followed by a ribosome removal step. The total RNA sample is then converted to cDNA before standard NGS library preparation. RNA-Seq focused on mRNA, small RNA, noncoding RNA, or microRNAs can be achieved by including additional isolation or enrichment steps before cDNA synthesis (Figure 9).

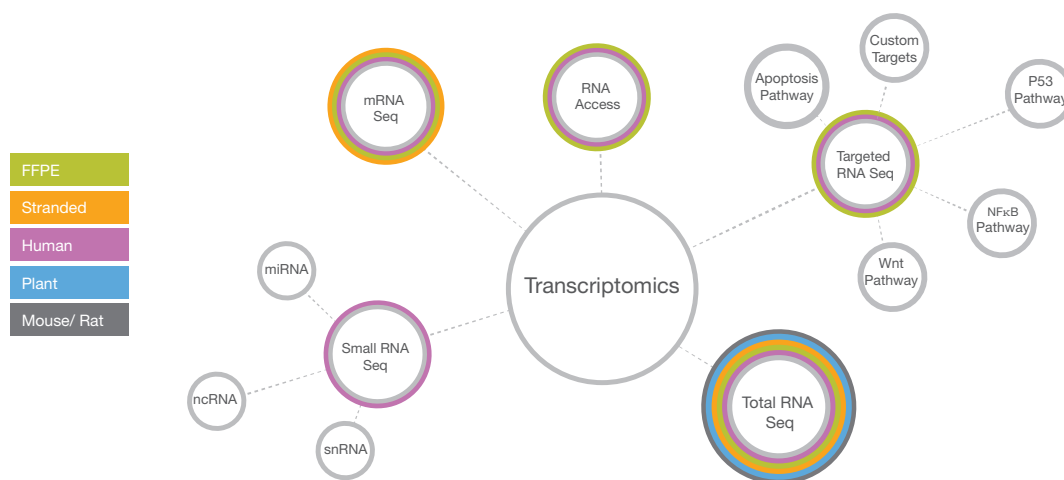


Figure 9: A Complete View of Transcriptomics with NGS—A broad range of methods for transcriptomics with NGS have emerged over the past 10 years including total RNA-Seq, mRNA-Seq, small RNA-Seq, and targeted RNA-Seq.

Total RNA and mRNA Sequencing

Transcriptome sequencing is a major advance in the study of gene expression because it allows a snapshot of the whole transcriptome rather than a predetermined subset of genes. Whole-transcriptome sequencing provides a comprehensive view of a cellular transcriptional profile at a given biological moment and greatly enhances the power of RNA discovery methods. As with any sequencing method, an almost unlimited dynamic range allows identification and quantification of both common and rare transcripts. Additional capabilities include aligning sequencing reads across splice junctions, as well as detection of isoforms, novel transcripts, and gene fusions. Library preparation kits that support precise detection of strand orientation are available for both total RNA-Seq and mRNA-Seq methods.

GAATGATAACAGTAAACACACTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACGTACCGTAAACGATCAATTGAGACTAAATATTAACGTACCATTAAAGAGCTACCGTCTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACG
TCAACGTACCGTAAACGAAACGATCATTAAGATTACTTGATCCACTGATTCAACGTACCGTAAACGAAACGATCAATTGAGACTAAATATTAACGTACCATTAAAGAGCTACCGTCTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACG
CGACGAAAGAAATGATAACAGTAAACACACTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACGTACCGTAAACGATCAATTGAGACTAAATATTAACGTACCATTAAAGAGCTACCGTCTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACG
ACGTACCATTAAGAGCTACCGTCAACAGTAAACACACTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACGTACCGTAAACGATCAATTGAGACTAAATATTAACGTACCATTAAAGAGCTACCGTCTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACG
GAAATGATAACAGTAAACACACTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACGTACCGTAAACGATCAATTGAGACTAAATATTAAGAGCTACCGTCTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACG
GATTACTTGATCCACTGATTCAACGTTAAGATTACTTGATCCACTGATTCAACGTACCGTAAACGAAACGATCAATTGAGACTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACGTACCGTAAACGAAACGATCAATTGAGACTAAATATTAACGTACCATTAAAGAGCTACCGTCTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACG
CGTATCAATTGAGACTAAATATTAACGTACCATTAAAGAGCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACGTACCGTAAACGAAACGATCAATTGAGACTAAATATTAACGTACCATTAAAGAGCTACCGTCTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACG

IV. Glossary

adapters: The oligos bound to the 5' and 3' end of each DNA fragment in a sequencing library. The adapters are complementary to the lawn of oligos present on the surface of Illumina sequencing flow cells.

bridge amplification: An amplification reaction that occurs on the surface of an Illumina flow cell. During flow cell manufacturing, the surface is coated with a lawn of 2 distinct oligonucleotides often referred to as “p5” and “p7.” In the first step of bridge amplification, a single-stranded sequencing library (with complementary adapter ends) is loaded into the flow cell. Individual molecules in the library bind to complementary oligos as they “flow” across the oligo lawn. Priming occurs as the opposite end of a ligated fragment bends over and “bridges” to another complementary oligo on the surface. Repeated denaturation and extension cycles (similar to PCR) results in localized amplification of single molecules into millions of unique, clonal clusters across the flow cell. This process, also known as “clustering,” occurs in an automated, flow cell instrument called a cBot™ or in an onboard cluster module within an NGS instrument.

clusters: A clonal grouping of template DNA bound to the surface of a flow cell. Each cluster is seeded by a single, template DNA strand and is clonally amplified through bridge amplification until the cluster has roughly 1000 copies. Each cluster on the flow cell produces a single sequencing read. For example, 10,000 clusters on the flow cell would produce 10,000 single reads and 20,000 paired-end reads.

contigs: A stretch of continuous sequence, *in silico*, generated by aligning overlapping sequencing reads.

coverage level: The average number of sequenced bases that align to each base of the reference DNA. For example, a whole genome sequenced at 30x coverage means that, on average, each base in the genome was sequenced 30 times.

digital microfluidics: Precise manipulation of droplets on a solid surface through applied voltages within a sealed microfluidic cartridge. For more information on digital microfluidics, see www.liquid-logic.com/technology.

flow cell: A glass slide with 1, 2, or 8 physically separated lanes, depending on instrument platform. Each lane is coated with a lawn of surface bound, adapter-complimentary oligos. A single library or a pool of up to 96 multiplexed libraries can be run per lane depending on application parameters.

indexes/barcodes/tags: A unique DNA sequence ligated to fragments within a sequencing library for downstream, *in silico* sorting and identification. Indexes are typically a component of adapters or PCR primers and are ligated to the library fragments during the sequencing library preparation stage. Illumina indexes are typically between 8–12 bp. Libraries with unique indexes can be pooled together, loaded into one lane of a sequencing flow cell, and sequenced in the same run. Reads are later identified and sorted via bioinformatic software. All together, this process is known as “multiplexing.”

insert: During the library preparation stage, the sample DNA is fragmented, and the fragments of a specific size (typically 200–500 bp, but can be larger) are ligated or “inserted” in between 2 oligo adapters. The original sample DNA fragments are also referred to as “inserts.”

mate pair library: A sequencing library with long inserts ranging in size from 2–5 kb typically run as paired-end libraries. The long gap length in between the sequence pairs is useful for building contigs in *de novo* sequencing, identification of indels, and other methods.

multiplexing: See “indexes/barcodes/tags.”

read: The process of next-generation DNA sequencing involves using sophisticated instruments to determine the sequence of a DNA or RNA sample. In general terms, a sequence “read” refers to the data string of A, T, C, and G bases corresponding to the sample DNA. With Illumina technology, millions of reads are generated in a single

