



An Introduction to Next-Generation Sequencing Technology



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 capillary electrophoresis (CE)-based sequencing in 1977, scientists gained the ability to sequence the full genome of any species in a reliable, reproducible manner.¹ A decade later, Applied Biosystems introduced the first automated, 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 to sequencing 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).

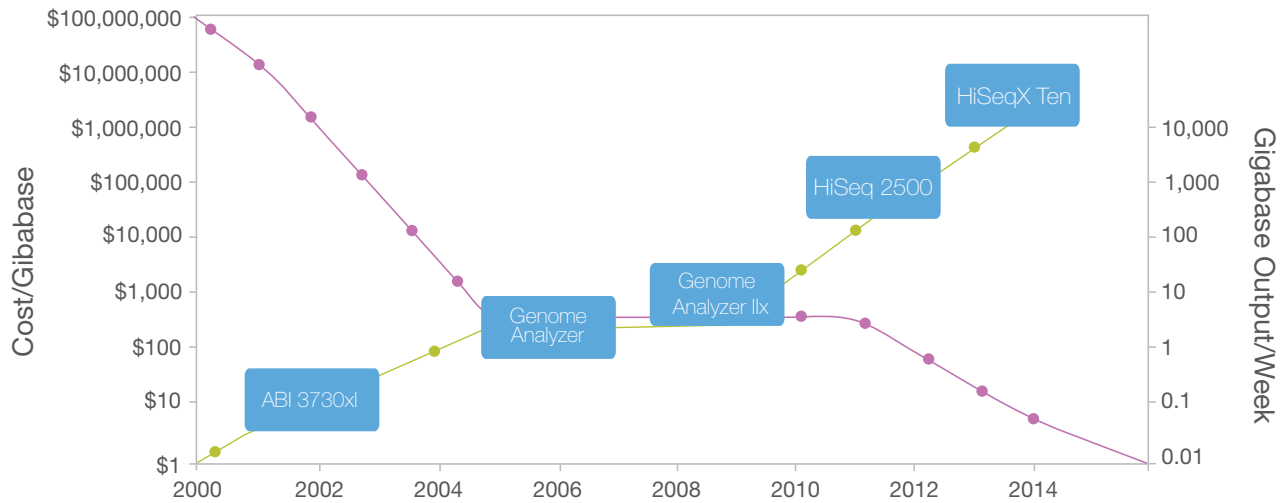


Figure 1: Sequencing Cost and Data Output Since 2000—The dramatic rise of data output and concurrent falling cost of sequencing since 2000. The Y-axes on both sides of the graph are logarithmic.

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

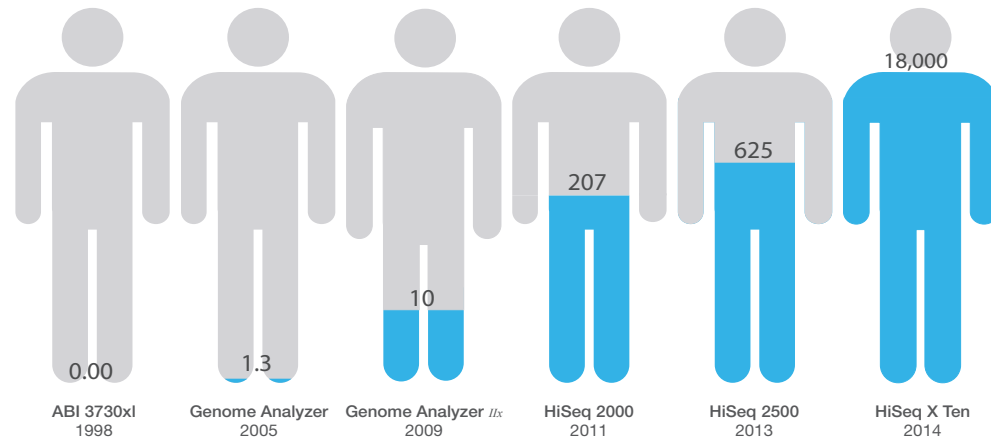


Figure 2: Human Genome Sequencing Over the Decades—The capacity to sequence all 3.2 billion bases of the human genome (at 30× coverage) has increased exponentially since the 1990s. In 2005, with the introduction of the Illumina Genome Analyzer System, 1.3 human genomes could be sequenced annually. Nearly 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.

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.

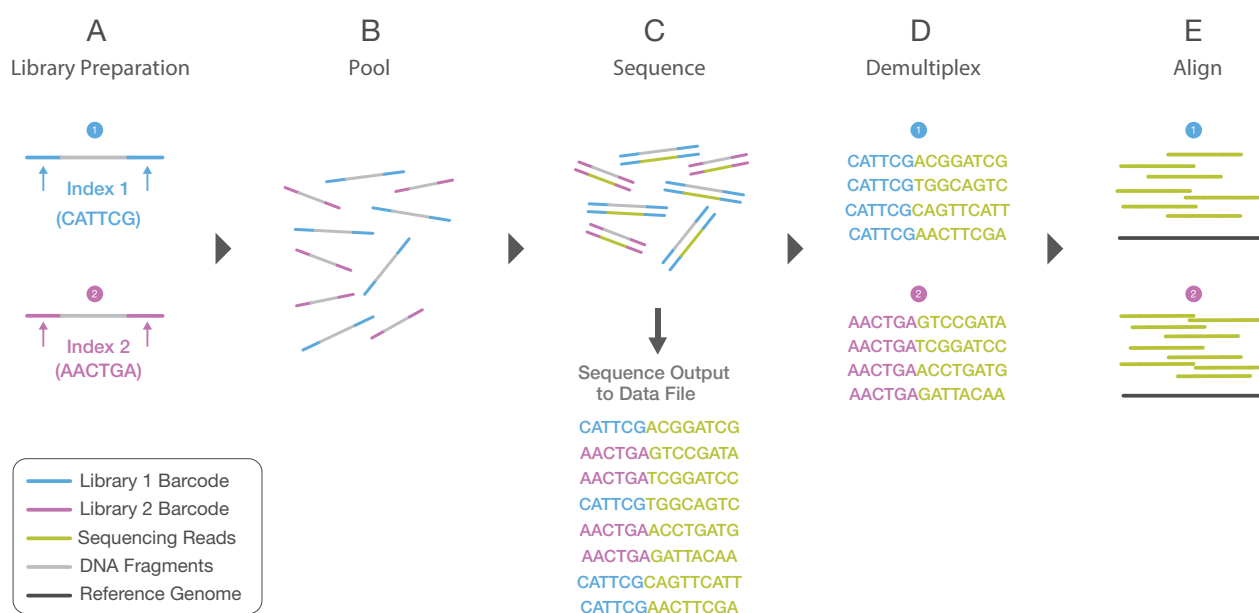


Figure 5: Library Multiplexing Overview.

- 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 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 desktop MiSeq® Series, with output ranging from 0.3–15 Gb for small genome, amplicon, or 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



With target enrichment, specific regions of interest are captured by hybridization to biotinylated probes, then isolated by magnetic pulldown. Target enrichment captures between 20 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 26–1536 targets at a time, spanning 150 bp–1.5 kb per target, 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).

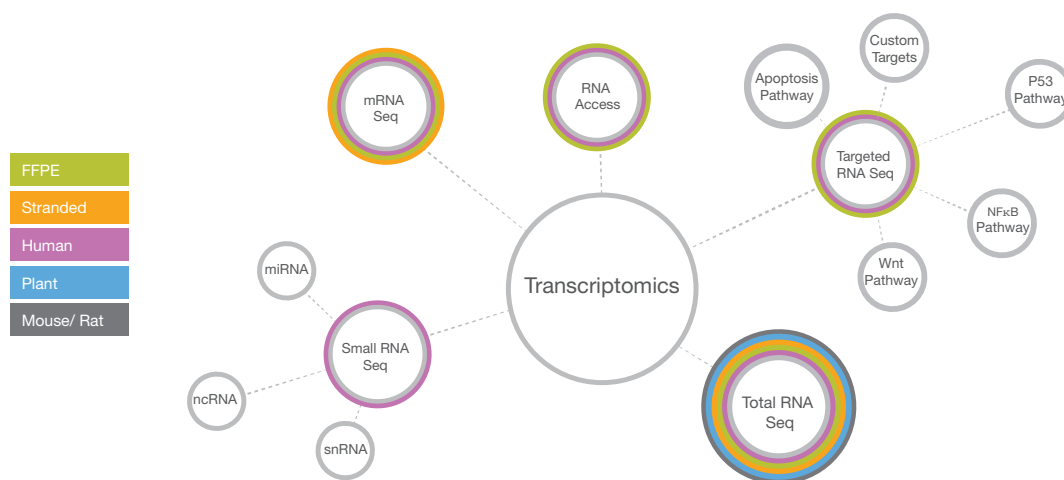


Figure 9: A Complete View of Transcriptomics with NGS—A broad range of methods for transcriptomics with NGS have emerged over the past 10 years including total RNA-Seq, mRNA-Seq, small RNA-Seq, and targeted RNA-Seq.

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 quantitation 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.

sequencing run. In more specific terms, 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.

reference genome: A reference genome is a fully sequenced and assembled genome that acts as a scaffold against which new sequence reads are aligned and compared. Typically, reads generated from a sequencing run are aligned to a reference genome as a first step in data analysis. In the absence of a reference genome, the newly sequenced reads must be constructed by contig assembly (*de novo* sequencing).

sequencing by synthesis (SBS): SBS technology uses 4 fluorescently labeled nucleotides to sequence the tens of millions of clusters on the flow cell surface in parallel. During each sequencing cycle, a single labeled dNTP is added to the nucleic acid chain. The nucleotide label serves as a “reversible terminator” for polymerization: after dNTP incorporation, the fluorescent dye is identified through laser excitation and imaging, then enzymatically cleaved to allow the next round of incorporation. As all 4 reversible terminator-bound dNTPs (A, C, T, G) are present, natural competition minimizes incorporation bias. Base calls are made directly from signal intensity measurements during each cycle, which greatly reduces raw error rates compared to other technologies. The result is highly accurate base-by-base sequencing that eliminates sequence-context-specific errors, enabling robust base calling across the genome, including repetitive sequence regions and within homopolymers.

V. References

1. Sanger F, Nicklen S, Coulson AR. DNA sequencing with chain-terminating inhibitors. *PNAS* 1977;74:5463-5467.
2. Collins FS, Morgan M, Patrinos A. The human genome project: lessons from large-scale biology. *Science*. 2003;300:286-290.
3. Davies K. (2010) 13 years ago, a beer summit in an English pub led to the birth of Solexa. *BioIT World* (www.bio-itworld.com/) 28 Sept 2010.
4. Illumina (2014) HiSeqX Ten preliminary system specification sheet. (www.illumina.com/documents/products/datasheets/datasheet-hiseq-x-ten.pdf)
5. Fallows J. (2013) When will genomics cure cancer? A conversation with Eric S. Lander. *The Atlantic* (www.theatlantic.com/) 22 Dec 2013.
6. Ross MG, Russ C, Costello M, et al. Characterizing and measuring bias in sequence data. *Gen Biol*. 2013;14:R51.
7. Bentley DR, Balasubramanian S, Swerdlow HP, et al. Accurate whole human genome sequencing using reversible terminator chemistry. *Nature*. 2008;456:53-59.
8. Nakazato T, Ohta T, Bono H. Experimental design-based functional mining and characterization of high-throughput sequencing data in the sequence read archive. *PLoS One*. 2013;22;8(10):e77910.
9. Illumina (2014) Nextera DNA Library Preparation Kits data sheet. (www.illumina.com/documents/products/datasheets/datasheet_nextera_dna_sample_prep.pdf)
10. Illumina (2014) Nextera XT DNA Library Preparation Kit data sheet. (www.illumina.com/documents/products/datasheets/datasheet_nextera_xt_dna_sample_prep.pdf)
11. Illumina (2013) TruSeq DNA PCR-Free Library Preparation Kit data sheet. (www.illumina.com/documents/products/datasheets/datasheet_truseq_dna_pcr_free_sample_prep.pdf)
12. Grad YH, Lipsitch M, Feldgarden M, et al. Genomic epidemiology of the *Escherichia coli* O104:H4 outbreaks in Europe, 2011. *PNAS*. 2012;109:3065-3070.
13. McEllistrem MC. Genetic diversity of the pneumococcal capsule: implications for molecular-based serotyping. *Future Microbiol*. 2009;4:857-865.
14. Lo YMD, Chiu RWK. Next-generation sequencing of plasma/serum DNA: an emerging research and molecular diagnostic tool. *Clin Chem*. 2009;55:607-608.

GAATGATAACAGTAAACACACTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACGTACCCTAACGAACGTATCAATTGAGACTAAATATTAACGTACCATTAAAGAGCTACCCTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACG
TCAACGTACCCTAACGAACGTATCAATTGAGATTACTTGATCCACTGATTCAACGTACCCTAACGAACGTATCAATTGAGACTAAATATTAACGTACCATTAAAGAGCTACCCTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACG
CGACGAAAGAAATGATAACAGTAAACACACTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACGTACCCTAACGAACGTATCAATTGAGACTAAATATTAACGTACCATTAAAGAGCTACCCTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACG
ACGTACCAATTAAGAGCTACCGTCAACAGTAAACACACTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACGTACCCTAACGAACGTATCAATTGAGACTAAATATTAACGTACCATTAAAGAGCTACCCTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACG
AGAATGATAACAGTAAACACACTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACGTACCCTAACGAACGTATCAATTGAGACTAAATATTAACGTACCATTAAAGAGCTACCCTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACG
GATTACTTGATCCACTGATTCAACGTAAAGATTACTTGATCCACTGATTCAACGTACCCTAACGAACGTATCAATTGAGACTAAATATTAACGTACCATTAAAGAGCTACCCTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACG
CGTATCAATTGAGACTAAATATTAACGTACCATTAAAGATTACTTGATCCACTGATTCAACGTACCCTAACGAACGTATCAATTGAGACTAAATATTAACGTACCATTAAAGAGCTACCCTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACG

15. Ram JL, Karim AS, Sandler ED, and Kato I. Strategy for microbiome analysis using 16S rRNA gene sequence analysis on the Illumina sequencing platform. *Syst Biol Reprod Med.* 2011;57:117-118.
16. Wang Y, Kim S, Kim IM. Regulation of metastasis by microRNAs in ovarian cancer. *Front Oncol.* 2014;10:143.
17. Dior Up, Kogan L, Chill HH, Eizenberg N, Simon A. Emerging roles of microRNA in the embryo-endometrium cross talk. *Semin Reprod Med.* 2014;32:402-409.

Illumina • 1.800.809.4566 toll-free (US) • +1.858.202.4566 tel • techsupport@illumina.com • www.illumina.com

For Research Use Only. Not for use in diagnostic procedures.

© 2015 Illumina, Inc. All rights reserved. Illumina, BaseSpace, cBot, HiSeq, MiSeq, NeoPrep, Nextera, NextSeq, TruSeq, and the pumpkin orange color are trademarks of Illumina, Inc. and/or its affiliate(s) in the U.S. and/or other countries.
Pub. No. 770-2012-008 Current as of 21 April 2015



GAATGATAACAGTAAACACACTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACGTACCGTAACGAACGTATCAATTGAGACTAAATATTAACGTACCATTAAAGAGCTACCGTCTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACG
TCAACGTACCGTAAACGAACGTATCAATTAAAGATTACTTGATCCACTGATTCAACGTACCGTAACGAACGTATCAATTGAGACTAAATATTAACGTACCATTAAAGAGCTACCGTCTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACG
CGACGAAAGAATGATAACAGTAAACACACTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACGTACCGTAAAGATTACTTGATCCACTGATTCAACGTACCGTAAACGAACGTATCAATTGAGACTAAATATTAACGTACCATTAAAGAGCTACCG
ACGTACCAATTAAGAGCTACCGTCAACAGTAAACACACTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACGTACCGTAAACGAACGTATCAATTGAGACTAAATATTAACGTACCATTAAAGAGCTACCGTCAACGAACGAAAGAATGATAA
AGAATGATAACAGTAAACACACTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACGTACCGTAAACGAACGTATCAATTGAGACTAAATATTAACGTACCATTAAAGAGCTACCGTCTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACG
GATTACTTGATCCACTGATTCAACGTAAAGATTACTTGATCCACTGATTCAACGTACCGTAAACGAACGTATCAATTGAGACTAAATATTAACGTACCATTAAAGAGCTACCGTCTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACG
CGTATCAATTGAGACTAAATATTAACGTACCATTAAAGACTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACGTACCGTAAACGAACGTATCAATTGAGACTAAATATTAACGTACCATTAAAGAGCTACCGTCAACGAACGAAAGAATGATAACAG