

The Evolution of Sequencing Technology

As the basic forms of life began to evolve on earth, so too did a means to carry the essential information required for the replication and reproduction of these increasingly complex organisms. DNA, besides acting as this blue-print, can also provide us with information about our phylogeny, ancestry, environment, susceptibility to disease, and—some claim—our character and possibly our fate. It is no wonder that we as humans find the process and challenge of decoding our genetic information such an irresistible endeavor.

The first serious attempt to decode the human genome was the Human Genome Sequencing Project, started in 1990. Although it had its initial detractors, this effort turned out to be one of the most successful collaborative ventures in scientific research to date. While cooperation and hard work undoubtedly played a critical role, the core engine that drove this project was a string of engineering breakthroughs that allowed for the collection and collation of data at an unprecedented rate. Today, the capacity of those approaches is being far surpassed by technologies that cut sequencing times and costs by several orders of magnitude. Superior detection methods, massive multiplexing, and much-reduced sample size can yield complete microbial genomes in a day and human genomes in only weeks.

It is the scientific foundation that has enabled this extraordinarily rapid progress that we are attempting to capture in this poster. As our knowledge of DNA—both chemical and functional—has grown, so have sequencing technologies evolved and improved, each discovery building upon the previous as we proceed along the uncoiling DNA strand. This poster provides a snapshot of the current state of the art, as well as giving the reader a broad—and necessarily abbreviated—overview of the development of sequencing technologies over the last century and a half. If the speed of advancement in this area continues apace, even some of the more formidable challenges, such as single strand sequencing and the \$1,000 genome, might be overcome—not just in our lifetimes, but within the foreseeable future.

Writer: John Hodgson
Design: Lewis Long
Editor: Sean Sanders
Commercial Editor, *Science*

Sponsored by:

454 LIFE SCIENCES



Sponsored by:



The Evolution of Sequencing Technology

Encoding of Sequence Information



Protein-coding Sequence
DNA that is transcribed and spliced to remove introns, leaving only exonic sequences, and is then translated to form functional proteins.

Pseudogenes
Sequences that have lost protein-coding ability or are otherwise no longer expressed in the cell because of mutation.

Promoters
DNA needed for transcription: promoters include an upstream binding site for RNA polymerase, a transcription start site, and regulatory elements.

Enhancers
200-1000 bp elements upstream or downstream of promoters that influence the transcription rates by binding activator/repressor proteins. Can be located many kilobases from the promoter they influence.

Sites of Methylation/CpG Islands
300-3000 bp-sequences in which the frequency of the dinucleotide CG is over 5% rather than the 1% found in the rest of the genome. Methylation of cytosine moiety in CpG islands near promoters is associated with gene repression, and demethylation with activation.

Origins of Replication
Sites (many per chromosome) at which DNA replication is initiated; often associated with genome regions of high (>70%) A+T content.

Nuclease Hypersensitive Sites
Regions exposed to nuclease activity because they are accessible to other DNA-binding proteins such as transcription factors or RNA polymerase.

Minimal Eukaryotic Promoters

Class or Promoter	Transcription Factor Binding Site	Transcription Start Point
TATA-box	TATAAAA	30 bp downstream of TATA-box
TATA-less	PyPyANTATPyPy	Controlled by initiator
TATA-less with downstream promoter element (DPE)		Around 30 bp upstream of DPE

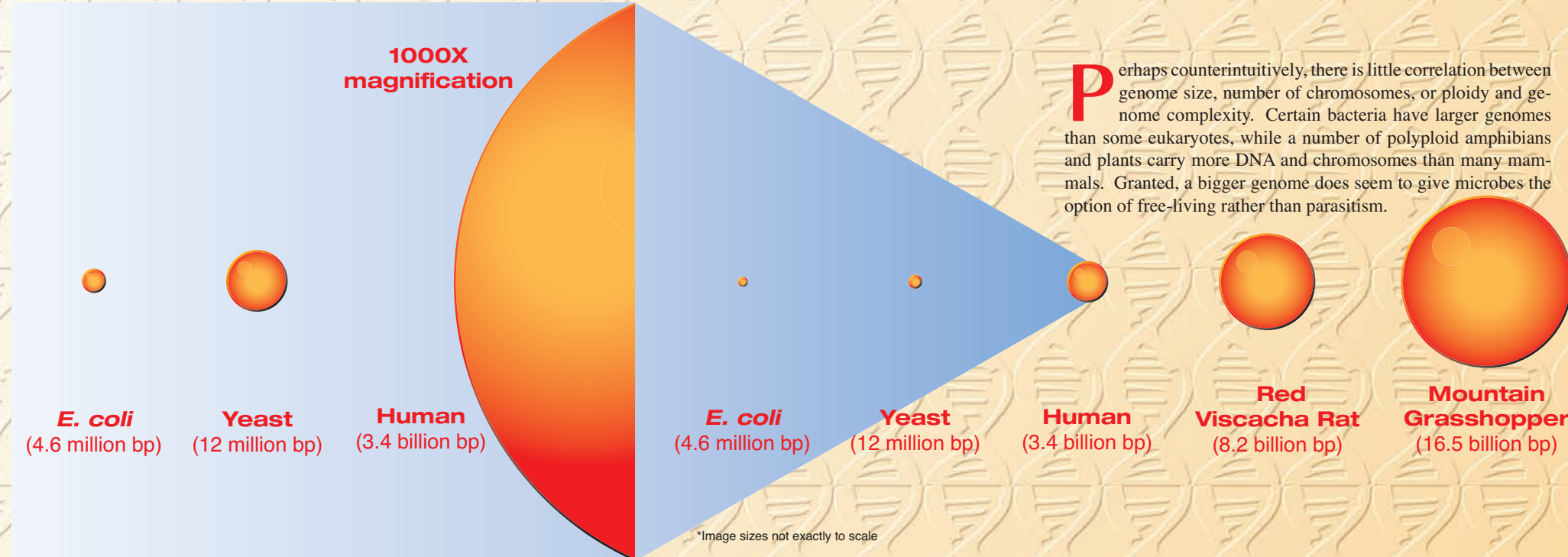
Scaffold/Matrix Attachment Sites
Generally AT-rich regions with long poly(A) sequences (over 100 bp), that create a rigid narrow minor groove structure bound by enzymes such as DNA topoisomerase II.

Histone Modifications
Acetylation, methylation, phosphorylation, and/or ubiquitination of DNA-packing histones represents a potential epigenetic code.

Histone Modifications

Modification	Enzyme	Illustrative Effects of Modification	Histones Involved
Acetylation	Histone acetyltransferases	Maintenance of active chromatin; acceleration of histone modification; transcriptional activation and silencing; dosage compensation; cell cycle progression; chromosome translocation	H1B, H2A, H2B, H3, H4
Deacetylation	Histone deacetylases	Transcription repression; trigger cell differentiation, growth arrest, apoptosis	
Methylation	Arginine methyltransferase Lysine methylase	Transcription activation or repression; cell growth and proliferation; activation/repression of methylation, acetylation, and ubiquitination of other histones Transcription activation and repression; heterochromatin formation; chromosome loss; marks euchromatin for cytosine methylation (DNA imprinting)	H3, H4 H1B, H2B, H3, H4
Demethylation	Demethylase; lysine demethylase	Transcription activation or repression	H1B, H2B, H3, H4
Phosphorylation	Histone kinases	Initiates chromosome condensation; activates transcription by promoting acetylation; response to DNA damage	H1B, H2A, H2B, H2AX, H3, H4
Ubiquitination	Ubiquitin ligases	Transcription activation; stimulation of transcription-activating methylation of H3	H2A, H2B

Genome Sizes



Perhaps counterintuitively, there is little correlation between genome size, number of chromosomes, or ploidy and genome complexity. Certain bacteria have larger genomes than some eukaryotes, while a number of polyploid amphibians and plants carry more DNA and chromosomes than many mammals. Granted, a bigger genome does seem to give microbes the option of free-living rather than parasitism.

Species	Common Name	Genome Size (billion bp)	Chromosome Number	Ploidy	Notes
<i>Encephalitozoon intestinalis</i>	Parasitic microsporidium	0.00225	10	1-2	Smallest eukaryal genome
<i>Escherichia coli</i>	-	0.0046	-	1	Typical bacterial genome size
<i>Pneumocystis carinii</i> (human)	-	0.0075	16	1	Smallest fungal genome
<i>Saccharomyces cerevisiae</i>	Baker's yeast	0.012	16	1-2	Common laboratory organism
<i>Trichoplax adhaerens</i>	-	0.054	6	1-2	Smallest animal genome
<i>Caenorhabditis elegans</i>	Worm	0.097	12	2	Common laboratory model
<i>Fragaria viridis</i>	Green strawberry	0.098	14	2	Smallest plant genome
<i>Caenocholax foveyi texensis</i>	Twisted-wing parasite	0.108	24	2	Smallest insect genome
<i>Arabidopsis thaliana</i>	Mustard weed	0.125	5	1	Common laboratory model
<i>Drosophila melanogaster</i>	Fruit fly	0.18	8	2	Common laboratory model
<i>Anopheles gambiae</i>	Malaria mosquito	0.278	6	2	Familiar insect
<i>Xenopus tropicalis</i>	Western clawed frog	1.7	20	2	Smallest amphibian chromosome number
<i>Danio rerio</i>	Zebrafish	1.7	50	2	Common laboratory model
<i>Canis familiaris</i>	Domestic dog	2.5	78	2	Familiar animal
<i>Zea mays</i>	Maize	2.5	20	4	Familiar plant
<i>Mus musculus</i>	Mouse	2.7	40	2	Familiar animal
<i>Xenopus laevis</i>	South African clawed frog	3.0	36	4	Common laboratory model
<i>Ornithorhynchus anatinus</i>	Duck-billed platypus	3.0	54	2	Familiar animal
<i>Homo sapiens</i>	Human	3.4	46	2	Familiar animal
<i>Bos taurus</i>	Domestic cow	3.7	50	2	Familiar animal
<i>Pan troglodytes</i>	Chimpanzee	3.8	48	2	Familiar animal
<i>Xenopus ravnoserratus</i>	Uganda clawed frog	7.8*	108	12	Largest amphibian chromosome number
<i>Tympanoctonus barrerae</i>	Red viscacha rat	8.2	102	4	Largest mammalian genome
<i>Podisma pedestris</i>	Mountain grasshopper	16.5	22-24	2	Largest insect genome
<i>Ambystoma mexicanum</i>	Axototl or Mexican salamander	21.9-48	28	2	Demonstrates regrowth of body parts
<i>Ophioglossum petiolatum</i>	Stalked adder's tongue	64	1020*	32-34*	Highest chromosome number and ploidy
<i>Necturus levisi</i>	Neuse River waterdog	118	38	2	Largest amphibian genome
<i>Fritillaria assyriaca</i>	Assyrian fritillary	125	48	4	Largest plant genome
<i>Protoperus aethiopicus</i>	Marbled lungfish	130	Unknown	2	Largest animal genome

Many familiar mammalian genomes are in the same 3-4 billion base pairs range as human, but a number of other mammals—bats and minkajak-like deer, for instance—have genomes under 2 billion bp. At the other extreme, the red viscacha rat (*Tympanoctonus barrerae*)—the only known mammalian tetraploid—has over twice the DNA of *Homo sapiens*, and more than double the chromosome number. The extremes of vertebrate genome sizes belong to fish: *Tetraodon nigroviridis*, the spotted green pufferfish, has a tiny 0.35 billion bp genome while the genome of *Protopterus aethiopicus*, the marbled lungfish, is approximately 130 billion bp.

A small proportion of mammals, birds, amphibians, and plants have more than 100 chromosomes. No insects (so far) have such high numbers. The top of the tree from the perspective of genome organization is the stalked adder's tongue, a 32- to 34-plant with an estimated 1,020 chromosomes.

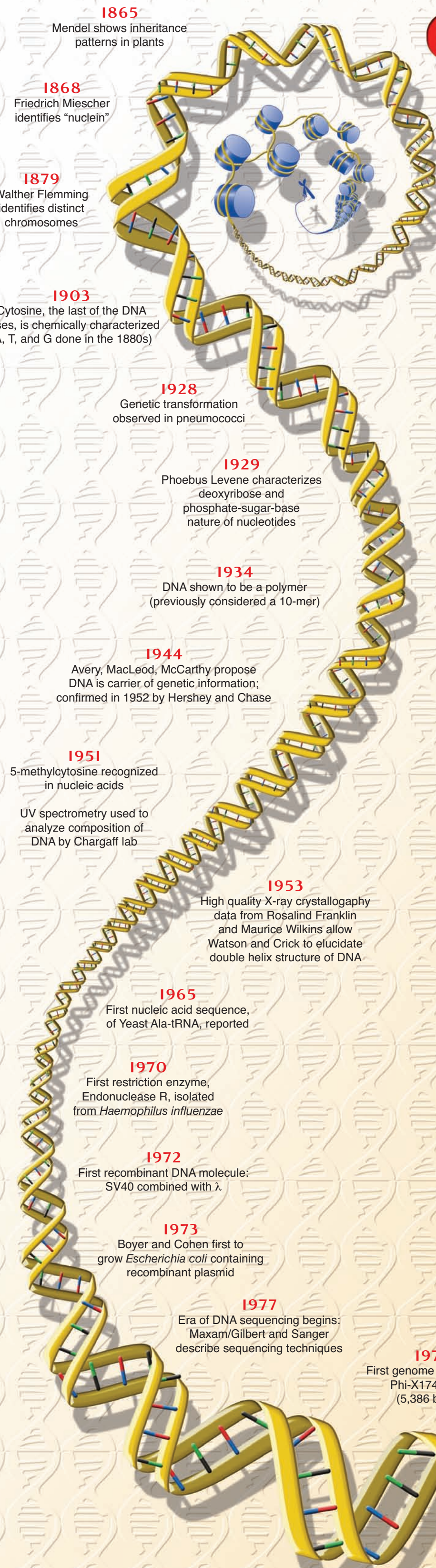
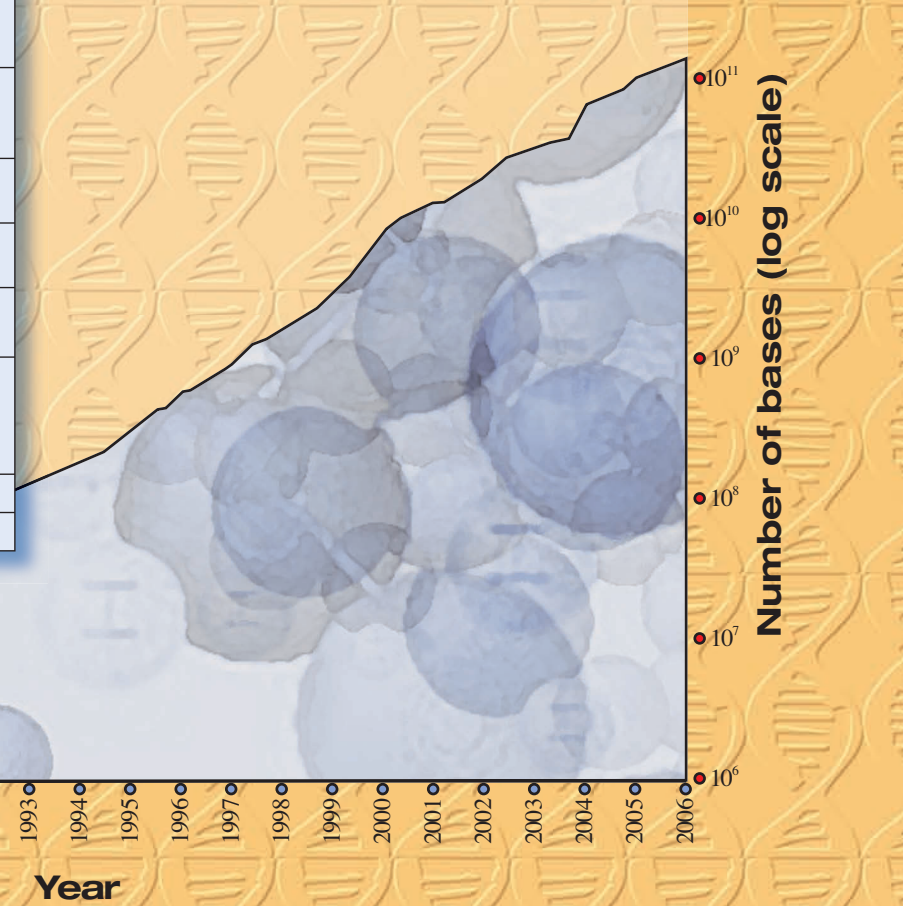
For more information, go online to www.genomesonline.org/gold.cgi and egg.ebi.ac.uk/services/cogent/

Sequencing Technologies

Sequencing principles - now and the future

Class	Subclass	Description	Optimum Sequencing Performance	Companies (and Commercially Available Systems)
Synthetic chain-terminator chemistry (Sanger method)		1977: DNA polymerase synthesizes a set of DNA fragments each one base longer than the other. Size-specific separation of the fragments (by electrophoresis, for instance) and base-specific tagging of each fragment (by fluorescence, for instance) allows sequence to be deduced.	67,000 - 96,000 bp/run (1-3h); 700-1000 bp read	Applied Biosystems (3730; 3730x1)
Sequencing-by-hybridization		1987: target sequence is annealed to a fixed array of oligonucleotide probes (8-10 bases). Sequence is deduced from hybridization pattern.	25 bp read (probe size)	Affymetrix/Perlegen (GeneChip CustomSeq Sequencing Arrays); Illumina (Beadchip); Premier Biosoft (AlleleID)
Cyclic sequencing on amplified DNA		Enzymatic methods are multiplexed in systems with a large number of addressable locations.		
Available next generation technologies	Pyrosequencing	1996: parallel sequencing-by-synthesis. Amplified tethered target sequences fixed in distinct physical locations (e.g., wells). Cycle adds each deoxynucleotide in turn with an enzyme cocktail and wells "light up" when the correct deoxynucleotide is present.	100 Mb/run (7-8h); ~250 bp read; 2400,000 reads/run	Roche Applied Science/454 Life Sciences (Genome Sequencer 20; Genome Sequencer FLX)
	Clonal single molecule array	2001: parallel sequencing-by-synthesis. Random fragments of target DNA tethered to a flow cell surface are amplified in situ. Sequencing cycle adds labeled, reversible chain terminators and interrogates all amplified clusters with a laser.	1000 Mb/run (2-3d); 25 bp read; >10 ⁷ reads/run	Illumina/Solexa (1G Genome Analyzer)
Single molecule sequencing-by-synthesis (still in development)	Sequencing-by-ligation	2007: unlabeled, tethered target DNA region defined by an "anchor primer" is probed by fluorescently labeled oligomers, the "query primers." Ligase extends the anchor primer, allowing subsequent bases to be determined in later cycles.	2000 Mb/run (>3d); 50 bp paired end read; >10 ⁷ reads/run	Applied Biosystems/Agencourt (SOLID)
	Fluorescence at just the active site of a single immobilized DNA polymerase enzyme measured using zero mode waveguide.		25 bp read; >10 ⁷ reads/run	Genovox (reagents and surface only)
Direct "reading" of single molecule (largely in early development)	Fluorescent resonance energy transfer (FRET) to channel energy via GFP-DNA polymerase to bound nucleotides.		~1000 Mb/run (~1d); ~10 ⁶ reads/run	Helicos BioSciences
	Conductance in nanopores. Differing physicochemical properties of nucleotide residues alter electric field as DNA is drawn through a pore. Protein pores, engineering nanochannels, and nanopore array systems are in development.		A DNA molecule is passed through field at rate of over 1000 bases per second	Agilent; LingVita AS
	Real-time reading of the DNA polymerase reaction using FRET. Force spectroscopy to follow the molecular mechanics of DNA synthesis.		~10 ⁶ Mb/run	Visigen; Li-Cor

Data Submitted to GenBank



Sponsored by:

454 LIFE SCIENCES

