

Thesis of the course of Biomedical Informatics
A.A. 2013-2014

Deoxyribonucleic acid
 ten-us-Deoxyribonucleic_acid.ogg
 /ˌdeɪ.əˈrɪb.əˌnjuːkleɪk ˈæ.sɪd/ (help·info)
 (DNA) is a nucleic acid that contains the genetic
 instructions used in the development and functioning
 of all known living organisms, and some viruses. The
 main role of DNA molecules is the long-term storage of
 information. DNA is often compared to a set of
 blueprints or a recipe, or a code, since it contains the
 instructions needed to construct other components of
 cells, such as proteins and RNA molecules. The
 segments that carry this genetic information are called
 genes, but other DNA sequences have structural
 purposes, or are involved in regulating the use of the
 genetic information.

Chemically, DNA consists of two long polymers of simple units called nucleotides, with backbones made of sugars and phosphate groups joined by ester bonds. These two strands run in opposite directions to each other and are therefore anti-parallel. Attached to each sugar is one of four types of molecules called bases. It is the sequence of these four bases along the backbone that encodes information. This information is read using the genetic code, which specifies the sequence of the amino acids within proteins. The code is read by copying stretches of DNA into the related messenger acid RNA, in a process called transcription.

random][plasmid

Within cells, DNA is organized into long structures called chromosomes. These chromosomes are duplicated before cells divide, in a process called DNA replication. Eukaryotic organisms (animals, plants, fungi, and protists) store most of their DNA inside the cell nucleus and some of their DNA in organelles, such as mitochondria or chloroplasts. In contrast, prokaryotes (bacteria and archaea) store their DNA only in the cytoplasm. Within the chromosomes, chromatin proteins such as histones compact and organize DNA. These compact structures guide the interactions between DNA and other proteins, helping control which parts of the DNA are transcribed.

The α -amino acid residues and α -hydroxy acids are the most abundant amino acid residues [10].

[illegible]

DNA exists in many forms, which include A-DNA, B-DNA, and Z-DNA. Only B-DNA and Z-DNA are functional organisms. The conformational change of DNA depends on the amount and type of metal ion binding to the DNA solution.[29]

The first published patterns—and also Patterson transforms—of a large amount of structural DNA [30][31]. And by Wilkins et al. [32] diffraction/scattering fibers in terms of a helical model. In the same journal, Watson and Crick's molecular model of DNA diffraction pattern for the helical [7].

Although the 'B' conformations found in the literature are in different conformations but all levels present in the diffraction and scattering molecular packing disorders.[35][36]

Compared to B-DNA, the right-handed spiral has a narrower, closer pitch and occurs under nonhydrated sample conditions produced in hybrid as well as in enzymatic segments of DNA. Chemically modified DNA shows a larger change in

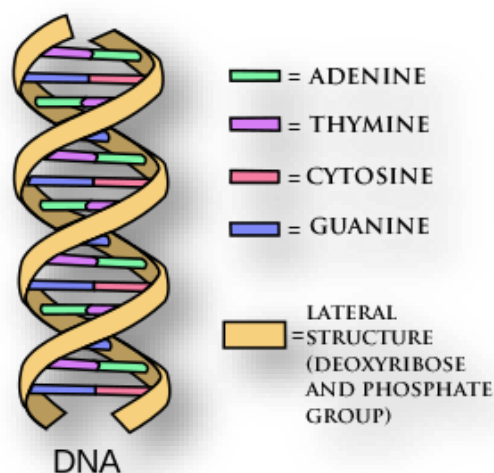
INTRODUCTION

In 1870, the Swiss chemist *Miescher* discovered inside the nucleus of a cell a giant molecule: **deoxyribonucleic acid**.

In 1953, two biochemists, the American *James Watson* and the English *Francis Crick* show that the structure of the DNA molecule is comparable to that of a spiral staircase; a sort of spiral-shaped double helix.

Deoxyribonucleic acid (DNA) is a molecule that encodes the genetic instructions used in the development and functioning of all known living organisms and many viruses. DNA is a nucleic acid; together with proteins and carbohydrates, nucleic acids compose the three major macromolecules essential for all known forms of life.

Most DNA molecules consist of *two biopolymer strands coiled around each other to form a double helix*. The two DNA strands are known as polynucleotides since they are composed of simpler units called *nucleotides*. Each nucleotide is composed of a **nitrogen-containing nucleobase**—either **guanine (G)**, **adenine (A)**, **thymine (T)**, or **cytosine (C)**—as well as a monosaccharide sugar called **deoxyribose** and a **phosphate group**. The nucleotides are joined to one another in a chain by *covalent bonds between the sugar of one nucleotide and the phosphate of the next*, resulting in an *alternating sugar-phosphate backbone*. According to base pairing rules (A with T and C with G), hydrogen bonds bind the nitrogenous bases of the two separate polynucleotide strands to make double-stranded DNA.



Picture 1 – DNA structure

DNA is well-suited for biological information storage. The DNA backbone is resistant to cleavage, and both strands of the double-stranded structure store the same biological information. Biological information is **replicated** as the two strands are separated. A significant portion of DNA (more than 98% for humans) is non-coding, meaning that these sections do not serve a function of encoding proteins.

That said, it is easy to understand how DNA is important for life. For this reason, even a small mutation (a change of the nucleotide sequence of the genome of an organism) can be decisive and cause diseases.

In this essay we will discuss a particular case of genomic mutation, the Single Nucleotide Polymorphism.

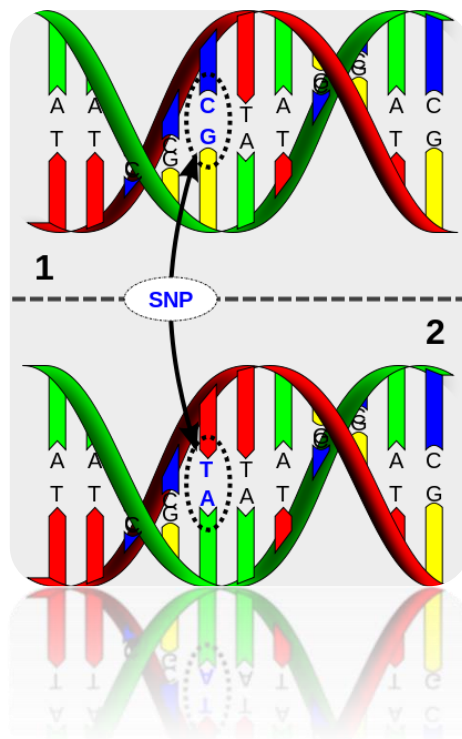
Single Nucleotide Polymorphism: what is it?

A **Single Nucleotide Polymorphism (SNP)** is a DNA sequence variation occurring commonly within a population (e.g. 1%) in which a Single Nucleotide — A, T, C or G — in the genome differs between members of a biological species or paired chromosomes.

For example, if we have two sequenced DNA fragments from different individuals (see *Picture 2*):

- AAGCCTA
- AAGCTTA

The second one contain a difference in a single nucleotide. In this case we say that there are two alleles. Almost all common SNPs have only two alleles.



Picture 2 – SNP example

The genomic distribution of SNPs is not homogenous; SNPs occur in non-coding regions more frequently than in coding regions.

The main causes of a SNP are:

1. natural selection, acting and fixating the allele of the SNP that constitutes the most favorable genetic adaptation
2. like genetic recombination
3. mutation rate

The different possible types of SNPs

As previously seen, Single Nucleotide Polymorphisms may fall within *coding* sequences of genes, *non-coding* regions of genes, as well as in the *intergenic* regions (regions between genes).

What is a coding?

To understand the difference between SNPs' types, we have to see what a coding is.

The main concept to analyse is the **Genetic Code**: it is the *set of rules* by which information encoded within genetic material (DNA or even mRNA sequences) is *translated* into proteins by living cells.

During the translation, the sequence of nitrogenous bases is treated in groups of three at a time; a group of three nitrogenous bases is called a **codon**. The code defines how codons specify which amino acid will be added next during protein synthesis. Generally, three-nucleotide codon in a nucleic acid sequence specifies a single amino acid. On the other hand, **a single amino acid can be specified by more than one codon**: this is the key concept that we will need in the following.

To understand better, here is the table that shows, for each amino acid (20 in total + START and STOP), the sequences that can generate it:

Amino acid	Codons
Ala/A	GCT, GCC, GCA, GCG
Arg/R	CGT, CGC, CGA, CGG, AGA, AGG

Asn/N	AAT, AAC
Asp/D	GAT, GAC
Cys/C	TGT, TGC
Gln/Q	CAA, CAG
Glu/E	GAA, GAG
Gly/G	GGT, GGC, GGA, GGG
His/H	CAT, CAC
Ile/I	ATT, ATC, ATA
Leu/L	TTA, TTG, CTT, CTC, CTA, CTG
Lys/K	AAA, AAG
Met/M	ATG
Phe/F	TTT, TTC
Pro/P	CCT, CCC, CCA, CCG
Ser/S	TCT, TCC, TCA, TCG, AGT, AGC
Thr/T	ACT, ACC, ACA, ACG
Trp/W	TGG
Tyr/Y	TAT, TAC
Val/V	GTT, GTC, GTA, GTG
START	ATG
STOP	TAA, TGA, TAG

Table 1 – Inverse genetic code

SNPs in the coding sequences

SNPs that fall in this category can be divided into two subcategories:

1. **Synonymous**
2. **Nonsynonymous**
 - a. **Missense**
 - b. **Nonsense**

First ones does not result in a change in the protein sequence, because the “original” sequence and the real sequence of bases both code the same amino acid.

Second ones, instead, change the amino acid sequence of protein. In their turn, they can be of two types: *Missense*, in which a single nucleotide change results in a codon that codes for a different amino acid (that can render the resulting protein non-functional), and *Nonsense*, that results in a premature stop codon, or a nonsense codon and then in a truncated, incomplete, and usually non-functional protein product.

Let us see an example of *Missense mutation*:

Original DNA code for the amino acid sequence:

C	A	T	C	A	T	C	A	T	C	A	T	C	A	T	C	A	T
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

Resulting amino acids:

His	His	His	His	His	His	His
-----	-----	-----	-----	-----	-----	-----

If we had, for example, a replacement of the eleventh nucleotide:

C	A	T	C	A	T	C	C	T	C	A	T	C	A	T	C	A	T
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

Resulting amino acids will be:

His	His	His	Pro	His	His	His
-----	-----	-----	-----	-----	-----	-----

This is, instead, an example of *Nonsense mutation*:

Original DNA code for the amino acid sequence:

A	T	G	A	C	T	C	A	C	C	G	A	G	C	G	C	G	A	A	G	C
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

Resulting amino acids:

Met	Thr	His	Arg	Ala	Arg	Ser
-----	-----	-----	-----	-----	-----	-----

If we had, for example, a replacement of the tenth nucleotide:

A	T	G	A	C	T	C	A	C	T	G	A	G	C	G	C	G	A	A	G	C
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

Resulting amino acids will be:

Met	Thr	His	Stop			
-----	-----	-----	------	--	--	--

Nonsense mutations are jointly responsible for many diseases; they can cause a genetic disease by damaging a gene responsible for a specific protein (for example *dystrophin* in *Duchenne muscular dystrophy*).

Examples of diseases in which nonsense mutations are known to be among the causes include:

- **Cystic fibrosis**
- **Duchenne muscular dystrophy** (dystrophin)
- **Beta thalassaemia** (β -globin)
- **Hurler syndrome**

On the other hand, cancer associated Missense mutations can lead to drastic destabilisation of the resulting protein.

SNPs not in coding regions

SNPs that are not in protein-coding regions may still affect:

1. gene splicing
2. transcription factor binding
3. messenger RNA degradation
4. ...

Gene expression affected by this type of SNP is referred to as an **eSNP** (*expression SNP*).

How to found SNPs: DNA Sequencing

In order to understand, according to what already said, what a SNP can cause, we must first of all *identify the SNP in the DNA* of the subject under consideration.

There are many possible types of analysis that can be performed (DNA sequencing, capillary electrophoresis, mass spectrometry, electrochemical analysis, ...); in this essay we will look at the most common one: **DNA Sequencing**.

“DNA Sequencing is the process of determining the precise order of nucleotides within a DNA molecule. It includes any method or technology that is used to determine the order of the four bases — adenine, guanine, cytosine, and thymine — in a strand of DNA.”

History of DNA Sequencing

Over the years, since the discovery of DNA by *Miescher* in 1870, the problem of DNA sequencing has been addressed in an increasingly thorough (see *Picture 3* in the next page). This also because, in 1940, *Avery* realized, by means of an experiment, that the so-called *transforming principle* (the carrier of genetic information) discovered in 1928 by *Griffith* was DNA.

Avery's experiment

In short, the *Avery* experiment was based on the *Griffith* experiment. *Griffith* used in his studies the *Streptococcus pneumoniae*. In particular, two of its strains:

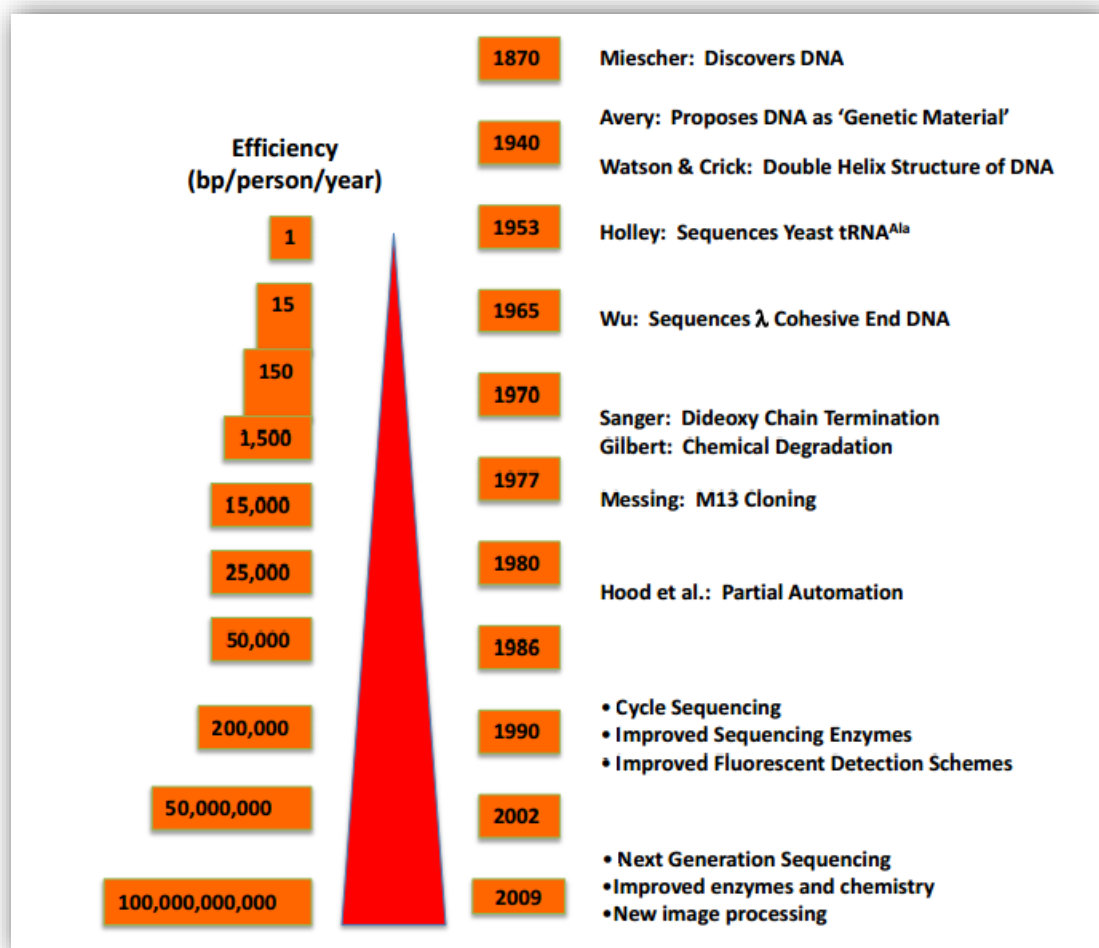
- the S strain, which can cause pneumonia in guinea pigs (virulent strain).
- the R strain is not able to cause pneumonia in guinea pigs (avirulent strain).

The main result was this:

injection in mouse of type S bacteria, killed after thermal treatment, and type R live bacteria was able to cause disease and death of the animal. From the tissues of mouse could isolate live bacteria of the S strain.

So, he verified and demonstrated that in a mixture containing either S dead bacteria and R alive bacteria, were to be happened the exchange of some substance (genetic material) that would confer virulence to bacteria R (which were then transformed into S).

The experiment of Avery aimed to determine what this substance was and it was discovered that it would necessarily be DNA.



Picture 3 – History of DNA Sequencing

Progress over the years



Picture 4 – Robert W. Holley

The first sequencing occurred in 1953 by Holley.

Since then, the efficiency of sequencing increased exponentially over the years: if in 1953 a person in a year could sequence only one *bp* (base pair), in the seventies we get to more than 1,500, in the nineties to more than 200000 and few years ago, in 2009, to more than 100 billion bps!

The first full DNA genome to be sequenced was that of bacteriophage ϕ X174 in 1977. Medical Research Council scientists deciphered the complete DNA sequence of the Epstein-Barr virus in 1984, finding it to be 170 thousand base-pairs long.

Sequencing methods

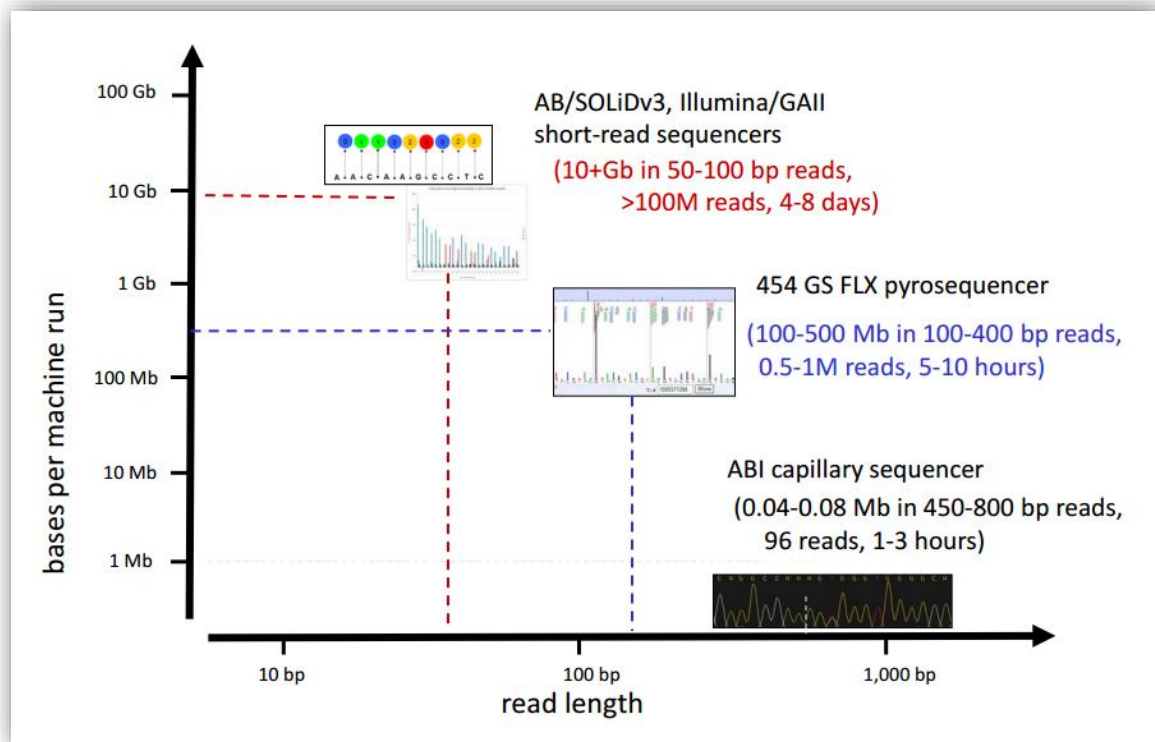
Over the years, many methods have been developed for sequencing the DNA. It goes from **basic methods**, such as the *Maxam-Gilbert sequencing* and *Chain-termination methods*, to **advanced methods** such as the *Shotgun sequencing* or *PCR Bridge*, to get to the **next-generation methods** (*Massively Parallel Signature sequencing (MPSS)*, *Po-lony sequencing*, *454 pyrosequencing*, *Illumina (Solexa) sequencing*, *SOLiD sequencing*, *Single Molecule Real Time (SMRT) sequencing*, ...).

Next-Generation Sequencing

Nowadays, thanks to technological progress we pushed even further forward. As can be seen from the following chart, it is possible to sequence **more than 100 million base pairs in about a week** (generating a very high amount of data). This is called the **Next-Generation Sequencing**.

However, the higher the speed of sequencing, the more there is a problem: **interpretation**. It often represents a real bottleneck; a single computer is not able to interpret a sequencing at the same speed of which it is presented to him.

For this reason, usually *cloud computing services* are used. They allow to take advantage of the computing power of multiple computers at the same time, parallelizing the work and thereby reducing the overall time.



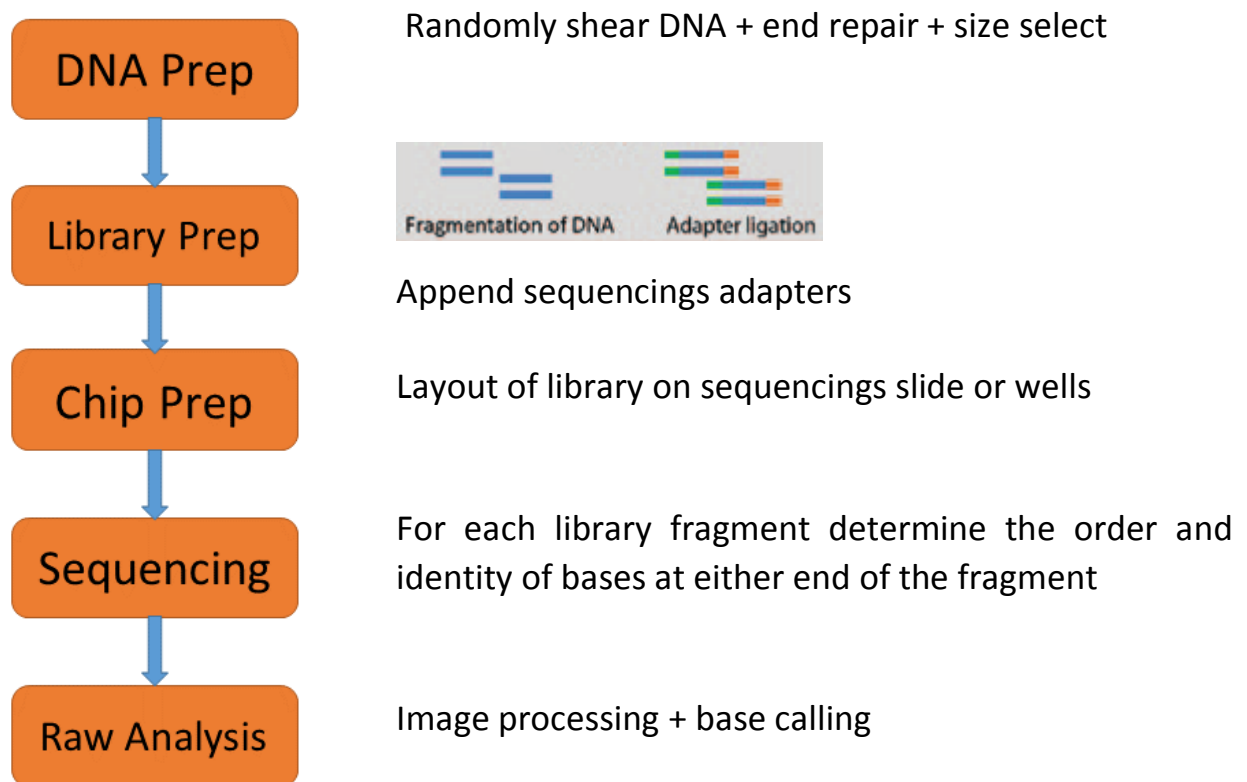
Picture 4 – Nowadays Sequencing

The high demand for low-cost sequencing has also driven the development of high-throughput sequencing (or next-generation sequencing) technologies that parallelize the sequencing process, producing thousands or millions of sequences concurrently. High-throughput sequencing technologies are intended to *lower the cost of DNA sequencing* beyond what is possible with standard methods.

In ultra-high-throughput sequencing as many as 500,000 sequencing-by-synthesis operations may be run in parallel.

Although each next-generation sequencing platform is unique in how sequencing is accomplished, there is a similar base methodology that includes preparation, sequencing, and data analysis. Within each generalized step, the individual platforms have unique aspects.

The common work-flow is the following:



DNA Sequencing Data format

Text

SNPs databases

Because SNPs are expected to facilitate large-scale association genetics studies, there has recently been great interest in SNP discovery and detection. For this reason databases can serve as a central repository. Once discovered, polymorphisms could be used by additional laboratories, using the sequence information around the polymorphism and the specific experimental conditions.

There are several databases that, nowadays, are used. The most important are:

1. **dbSNP** is a SNP database from the *National Center for Biotechnology Information (NCBI)*
2. **SNPedia** is a wiki-style database supporting personal genome annotation, interpretation and analysis.

Furthermore, there are various support databases that allow, for example, to bind a SNP to the disease that causes:

1. **OMIM** database describes the association between polymorphisms and diseases (e.g., gives diseases in text form)
2. **Human Gene Mutation Database** provides gene mutations causing or associated with human inherited diseases and functional SNPs
3. **GWAS Central** allows users to visually interrogate the actual summary-level association data in one or more genome-wide association studies
4. ...