

An introduction to high-throughput sequencing experiments: design and bioinformatics analysis

Rachelly Normand¹ and Itai Yanai²

1. Life Sciences and Engineering, HT-Seq Unit, Technion – Israel Institute of Technology
2. Department of Biology, Technion – Israel Institute of Technology

Correspondence to: Itai Yanai, Tel: +972-4-829-3763; Fax: +972-4-822-5153;
Email: yanai@technion.ac.il

Summary

The dramatic fall in the cost of DNA sequencing has revolutionized the experiments within reach in the life sciences. Here we provide an introduction for the domains of analyses possible using high-throughput sequencing, distinguishing between “counting” and “reading” application. We discuss the steps in designing a high-throughput sequencing experiment, introduce the most widely used applications, and describe basic sequencing concepts. We review the various software programs available for many of the bioinformatics analysis required to make sense of the sequencing data. We hope that this introduction will be accessible to biologists with no previous background in bioinformatics, yet with a keen interest in applying the power of high-throughput sequencing in their research.

Introduction

High-throughput sequencing is the process of identifying the sequence of millions of short DNA fragments in parallel. In this chapter, we will discuss applications and analyses of high-throughput sequencing done on the Illumina platform. The main advantage of this technology is that it allows a very high-throughput; currently up to 1.6 billion DNA fragments can be sequenced in parallel in a single run, to produce a total of 320Gbp (HiSeq 2000, version 3 kits). One challenge with this technology, however, is that the sequenced fragments are relatively short – currently up to 150bp (MiSeq instrument) or 100bp (HiSeq 2000 instrument) – though double this can be produced using the paired-end option (see below).

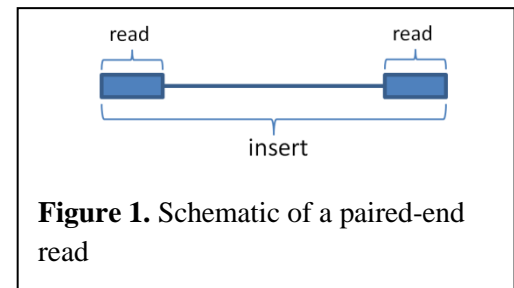
We operate a service unit in a University setting providing high-throughput sequencing (henceforth, HTS) sample preparation, sequencing and initial bioinformatics analysis. Based upon our experiences over the past two years we provide the following notes. We do not aim to provide a complete picture of all of the innumerable resources available for any one of the described applications. Rather, our goal is to provide a basic overview of the opportunities and challenges that HTS represents. The field is clearly changing rapidly and so the details are to be taken with caution as they will surely need revision as new algorithms and technology emerge.

While many applications are supported by HTS, the actual input to the instrument is the same: libraries comprised of billions of DNA strands of roughly the same length (typically 300bp) with particular sequences (linkers) on either end. “Sample preparation” is the process by which an initial sample arrives at this highly ordered state. When genomic DNA is the starting material, it is fragmented and then size-selected for the tight size distribution. If the starting material is RNA, often times it is polyA-selected to limit the sequencing to mRNA. The RNA is reverse transcribed to DNA and then also size-selected. Irrespective of the application, linker DNA molecules of particular sequences are ligated to the ends of the strands. These consist of two fragments: adaptors and indices. The adaptors hybridize the DNA fragments to the flowcell on which they are sequenced. The indices are 6-7bp sequences tagging different samples within the same library that will be sequenced together. Importantly there is a PCR amplification step in many of the sample preparation protocols which has implications for the structure of the data: identical sequences may be a result of the amplification or reflect recurrence in the original sample of DNA.

Basic concepts in high-throughput sequencing

Figure 1 indicates the anatomy of an insert. The following are additional basic definitions important for high-throughput sequencing:

- **Insert** – the DNA fragment that is used for sequencing.
- **Read** – the part of the insert that is sequenced.
- **Single Read (SR)** – a sequencing procedure by which the insert is sequenced from one end only.
- **Paired End (PE)** – a sequencing procedure by which the insert is sequenced from both ends.
- **Flowcell** – a small glass chip on which the DNA fragments are attached and sequenced. The flowcell is covered by probes that allow hybridization of the adaptors that were ligated to the DNA fragments.
- **Lane** – the flowcell consist of 8 physically separated channels called lanes. The sequencing is done in parallel on all lanes.
- **Multiplexing / Demultiplexing** – sequencing a few samples on the same lane is called multiplexing. The separation of reads that were sequenced on one lane to different samples is called demultiplexing and is done by a script that recognizes the index of each read and compares it to the known indices of each sample.
- **Pipeline** – a series of computational processes.



High-throughput sequencing applications

HTS applications can be divided into two main categories: ‘reading’ and ‘counting’. In reading applications the focus of the experiment is the sequence itself, for example for finding genomic variants or assembling the sequence of an unknown genome. Counting applications are based on the ability to count amounts of reads and compare these counts, for example to assess gene expression levels. Table 1 shows some of the main applications enabled by high-throughput sequencing. These represent but a sampling of the main HTS applications. It should be noted that

one can invoke HTS in practically any experiment that produces DNA fragments. What should be considered and planned before the sequencing however is the method by which the analysis of the sequenced fragments will be done to extract the meaning from the experiment. As an example of a unique HTS experiment, chromatin interactions can be identified by PE sequencing (1). This procedure includes capturing interacting loci in the genome by immune-precipitating cross-linked fragments of DNA and proteins from fixed cells. There are many others, published at a rate of about one per day.

Table 1. HTS applications.

	Application	Goal	Experiment details	Basic analysis summary
Reading	Re-sequencing	Find variants in a given sample relative to reference genome.	Extract DNA from the relevant cells, conduct sample preparation consisting of DNA fragmentation and sequencing.	Mapping of the sequenced fragment to the reference genome and identifying variants relative to the reference genome by summarizing the differences of the fragments from the genomic loci to which they map.
	Target-enriched sequencing	Target enrichment sequencing is a specific form of re-sequencing that is focused only on certain genomic loci. This is useful for organisms with large genomes where enrichment increases the coverage on the loci of interest thereby reducing costs	After the DNA is extracted from the cells and undergoes sample preparation, an enrichment process is done to capture the relevant loci. Target enrichment can be done on specific regions of the genome using “tailored” target-enrichment probes, or by using available kits such as exome-enrichment kits.	Same as in resequencing.
	De-novo assembly	Identify a genomic sequence without any additional reference.	Same as in re-sequencing.	The assembly process relies on overlaps of DNA fragments. These overlaps are merged into consensus sequences called contigs and scaffolds.
Counting	ChIP-Seq/ RIP-Seq	Find the binding locations of RNA- or DNA-binding proteins.	First, the ChIP/RIP experiment is done: proteins are bound to the DNA/ RNA and are cross-linked to it. The DNA/RNA is then fragmented. The proteins are pulled down by an immuno precipitation process and are then the cross-linking is reversed. The DNA/RNA fragments that are enriched in the protein binding sites locations are then sequenced.	The sequenced fragments are mapped to the genome. The enriched locations in the genome are found by detecting “peaks” of mapped fragments along the genome. These peaks should be significantly higher than the mapped fragments in the surrounding loci, and significantly higher compared to a control sample – usually the input DNA of the ChIP experiments or another sample of immuno-precipitation done by a non-specific antibody.
	RNA-Seq	Detecting and comparing gene expression levels.	Total RNA is extracted from the cells. In a sample preparation process the mRNA is pulled down and fragmented. The mRNA fragments are then reversed transcribed to cDNA. The cDNA fragments are sequenced.	The cDNA fragments are mapped to the reference genome. The fragments that map to each gene are counted and normalized to allow comparisons between different genes and different samples. Un-annotated genes and transcripts can be found in an RNA-Seq experiment by detecting bundles of fragments that are mapped to the genome in an un-annotated region.
Reading/ Counting	microRNA-Seq	Detect and count microRNAs.	Total RNA is extracted from the cells, and the microRNA is isolated by recognizing the	The sequenced fragments are mapped to the genome. The microRNA can then be detected and counted.

			natural structure common to most known microRNA molecules. The microRNA fragments are then reversed transcribed and sequenced.	
--	--	--	--	--

Sequence coverage

In reading applications, coverage corresponds to the number of reads that cover each base in the genome on average. Coverage can be calculated as:

$$\text{average coverage} = \frac{\text{read length} \cdot \text{number of reads}}{\text{genome size}}$$

Note that only the number of mapped reads should be included in the above calculation. In general, 30X coverage is considered a minimum for identifying genomic variants, while *de-novo* assembly usually requires a much higher coverage. Furthermore, the needed coverage depends on the experiment design. For example, if re-sequencing is done on a population and the sample includes pooling of heterogenic genomes, the coverage must be higher for the robust detection of rare variants.

Contaminations may not pose a great difficulty for ‘reading’ applications since they will not map to the reference genome. However contaminations “steal” coverage from the sample, and should be taken into account when estimating the expected coverage. If it is not possible to assess what percentage of contaminations the sample will contain, a pilot experiment may again prove useful: sequencing of just one or two samples in low coverage, and then assessing by mapping the percentage of contaminants. In *de-novo* assembly, contaminations may be a lot more difficult to detect and thus attempts to eliminate contamination should be made when extracting the DNA, before sequencing and analysis.

In counting applications, such as RNA-Seq, the notion of coverage is not straightforward since the number of reads along the genome is not expected to be uniform. For example most RNA-Seq reads will correspond to highly expressed transcripts, whereas lowly expressed transcripts will be less represented. This notion presents the question of how many reads are required for a particular application. In general, this is a trial and error process, and consequently we have found it useful to begin with a pilot experiment of a few samples to provide an estimate of the transcriptomic complexity.

An analysis that can help assess whether enough reads have been sequenced is a “saturation report” (Figure 2, (2)). In this “jack-knifing” method, the expression levels are determined using all of the reads. The expression levels are then compared to those recalculated using only a fraction of the reads. Examining the expression levels at each cut of the data informs at which point the expression level remain unchanged despite additional data. As expected, additional data is most helpful in resolving the expression levels of the lowly expressed genes. After deciding how many reads are required per sample, the samples are divided into lanes according to the number of sequenced reads per lane, which is a fixed amount.

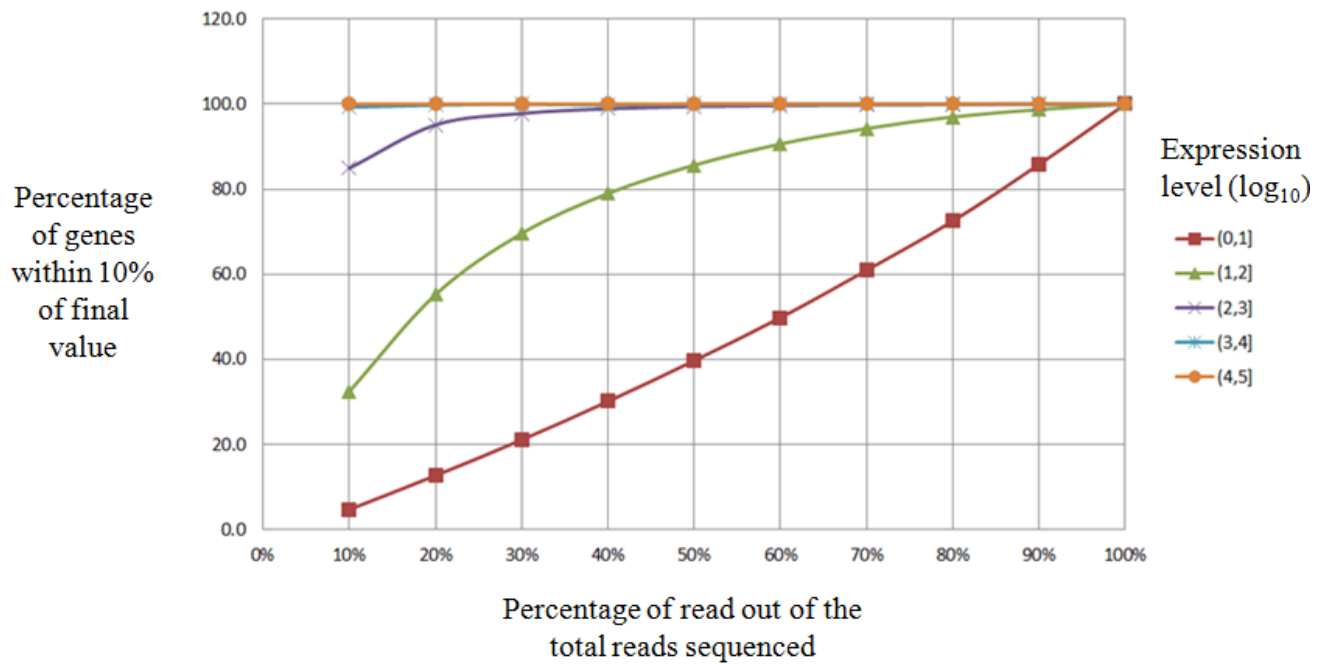


Figure 2. Saturation report. The different series are sets of genes that differ in their final expression values using the complete dataset (in this case, 32 million reads). Highly expressed genes are saturated with even 10% of the reads, whereas lowly expressed genes require a higher amount of reads, while very lowly expressed genes remain unsaturated even with the complete dataset.

Sequencing recipe – Single-read vs. Paired-end, insert size and read length

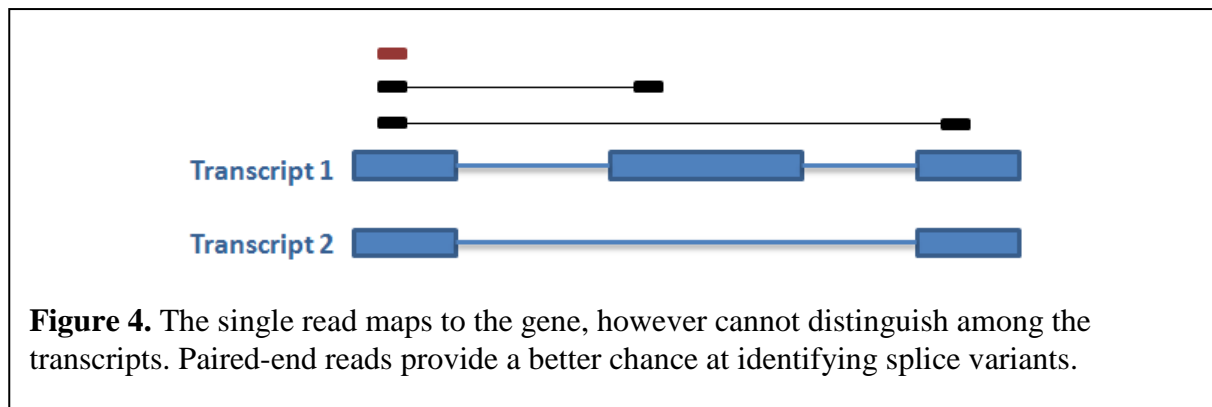
The sequencing recipe is influenced by several factors:

The repetitive nature of the genome. Human and mouse genomes have ~20% repetitive sequences (3). Consequently, to uniquely score a read mapping to a repetitive region it must be longer than the repetitive region or border the neighboring non-repetitive sequence. Longer reads or PE reads allow “rescue” of the non-unique end and also mapping to non-unique regions in the genome (Figure 3).



Figure 3. The red end would not have been uniquely mapped if sequenced as a single read as opposed to a paired-end read.

Differentially spliced variants. When assessing gene expression levels in RNA-Seq, it is potentially informative to discover the differential expression levels of different transcripts of the same gene. Reads that map to an exon shared by more than one transcript pose a difficulty in assessing the transcript of origin. PE reads may solve this problem if one end of the sequenced fragment maps to an exon that is unique to one of the transcripts. Figure 4 shows an example in which one cannot determine with certainty from which transcripts the SR originated. Sequencing it as PE resolves this problem.



Genetic distance of the sequenced sample from the reference genome. If the sequenced samples are genetically distant from the reference genome, longer reads may be required to determine the origin of each read in the genome. The mappings of each read will contain more mismatches and thus making it difficult to unambiguously determine its correct location, thereby increasing the probability that more than one location may be possible. Thus, the longer the read, the more likely a unique mapping becomes.

Finding structural variations. Structural variations in the genome, such as long insertions or deletions, inversions, and translocations, can be found using PE information. For example, if a large deletion is present in the sequenced strain, the inserts lengths will be longer than expected (Figure 5).

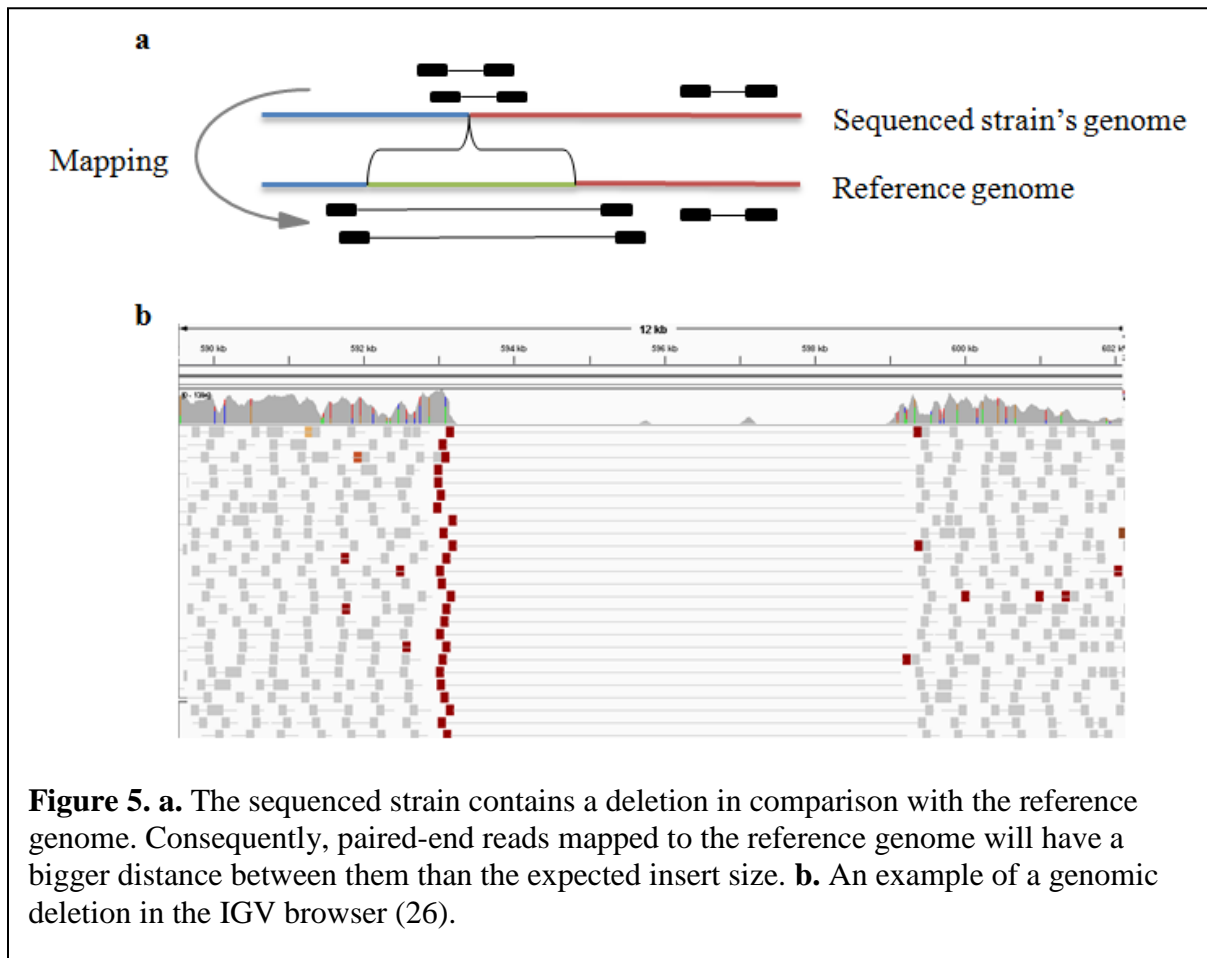


Figure 5. a. The sequenced strain contains a deletion in comparison with the reference genome. Consequently, paired-end reads mapped to the reference genome will have a bigger distance between them than the expected insert size. **b.** An example of a genomic deletion in the IGV browser (26).

De novo assembly. Assembling a new genome from short sequenced reads consist of overcoming many challenges, such sequencing errors, low complexity regions and repetitive regions among others (4,5). De novo assembly remains a notoriously difficult problem and often the genome of a metazoan remains in thousands of contigs. Obviously, longer PE reads lead to better assemblies. It has also been shown that using a few sequencing libraries with different insert length may improve the assembly process (4).

Number of samples for sequencing

Resequencing. If the reference genome to which the sequenced reads are mapped is genetically distant, sequencing the actual strain in its baseline state (before the mutagenesis, without the phenotypic change, etc.) will be beneficial for interpreting the data. This will help in distinguishing the variations that are due to evolutionary distance from those that cause the actual phenotypic trait under study.

RNA-Seq. It is highly recommended to sequence a few biological replicates to control for biological noise. Technical replicates will also inevitably show variation (6). Some gene expression software programs, such as Cufflinks (7), can use the data from different replicates and merge it into one value with a higher statistical significance.

ChIP-Seq. A ChIP-Seq experiment should include the IP DNA and one more sample that will serve as a control. The control sample may be the input DNA, before the IP process, or an IP done on the same DNA with a non-specific antibody, such as IgG (8,9). Sequencing a control sample enables detection of enriched regions that are also significantly enriched compared to the control sample, and not only enriched compared to the area surrounding them in the IP sample. This may reduce false-positive peaks detected solely because of areas in the genome that have a higher coverage due to better DNA fragmentation compared to the surrounding area.

Analysis Pipelines

Figure 6 shows the bioinformatics pipelines involved in four main applications: resequencing, de-novo assembly, RNA-Seq and ChIP-Seq. Several processes are common to all or multiple applications.

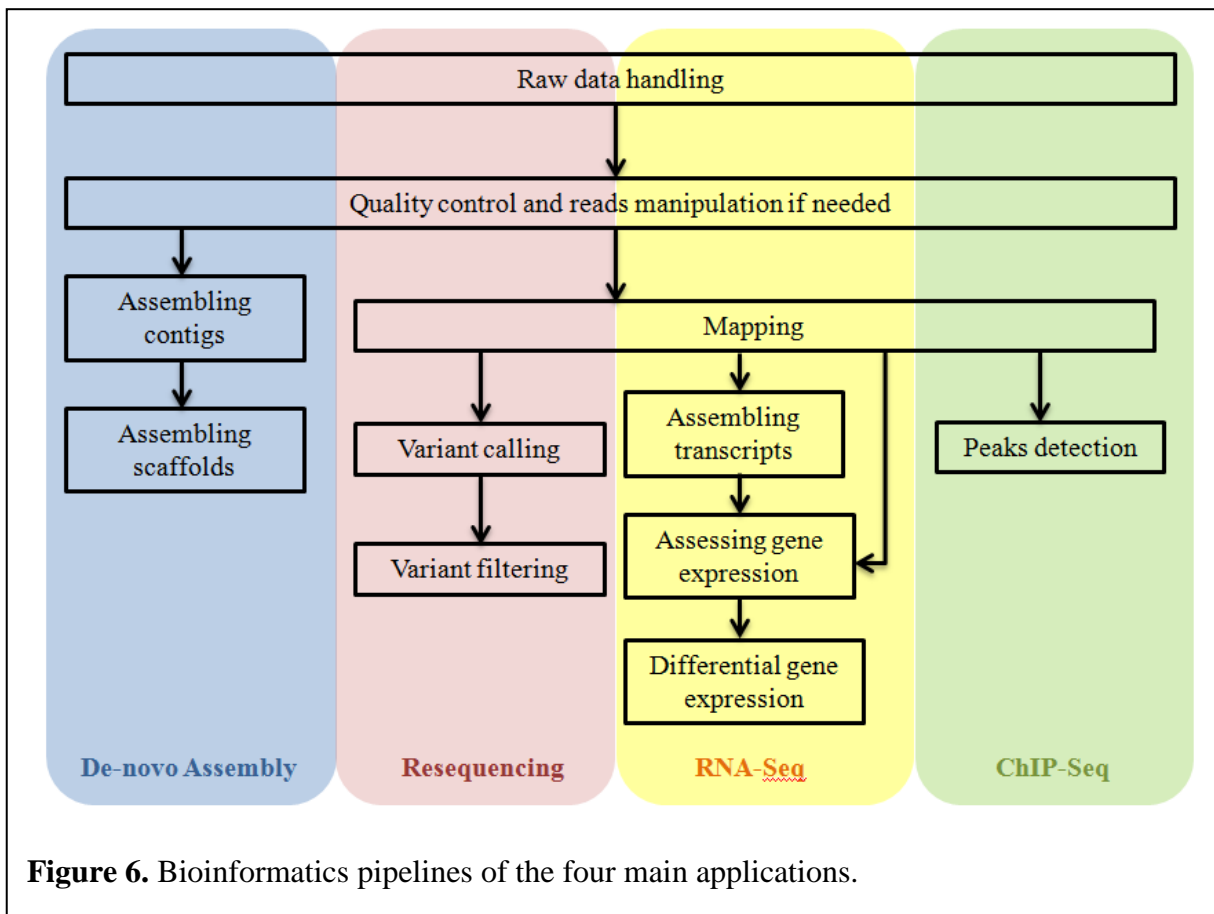


Figure 6. Bioinformatics pipelines of the four main applications.

1. Raw data handling. Available software for this step: Illumina’s CASAVA software. The Illumina run produces “base-calling” files (*.bcl) which only become useful bioinformatically when converted to the general fastq format (see below). During this file conversion, the demultiplexing process is also carried out, which is the separation of reads from different samples that were sequenced on the same lane.

2. Quality control and reads manipulation. Available software for this step: CASAVA, FastQC (Babraham Bioinformatics). After a sequencing run is completed and before starting the analysis, the run's quality should be checked for the following parameters which may be telling of the quality of the sample and run.

- a. Pass Filter (PF) reads – The number and percentage of PF reads in each lane and for each sample should match the number of expected sequenced reads. If it is dramatically lower, this might indicate a low quality run, and may reduce the expected coverage.
- b. Control reads – Apart from the DNA libraries, control DNA from the viral PhiX genome is spiked-in at 1% concentration with the sample onto each lane of the flowcell. Reads are automatically mapped by the Illumina software to the PhiX genome. The percentage of reads from each lane mapping to this genome and the amount of mismatches in the mapping are used as control values for the lane's quality. A good run typically has ~1% sequencing errors, as detected by the mismatches to the PhiX genome.
- c. Quality scores of the reads – As will also be explained in the next section (“Diving into the technical details”) each base of each sequenced read is associated with a quality score providing the confidence in the particular base. In general, the quality scores drop towards the end of the sequenced read. These confidences should be assessed to check for the overall quality of the run. The quality scores may be automatically produced by the sequencing platform, and may also be created by programs like FastQC that provide other statistics on the sequenced reads, such as overrepresented sequences, per base GC content and more (Figure 7).

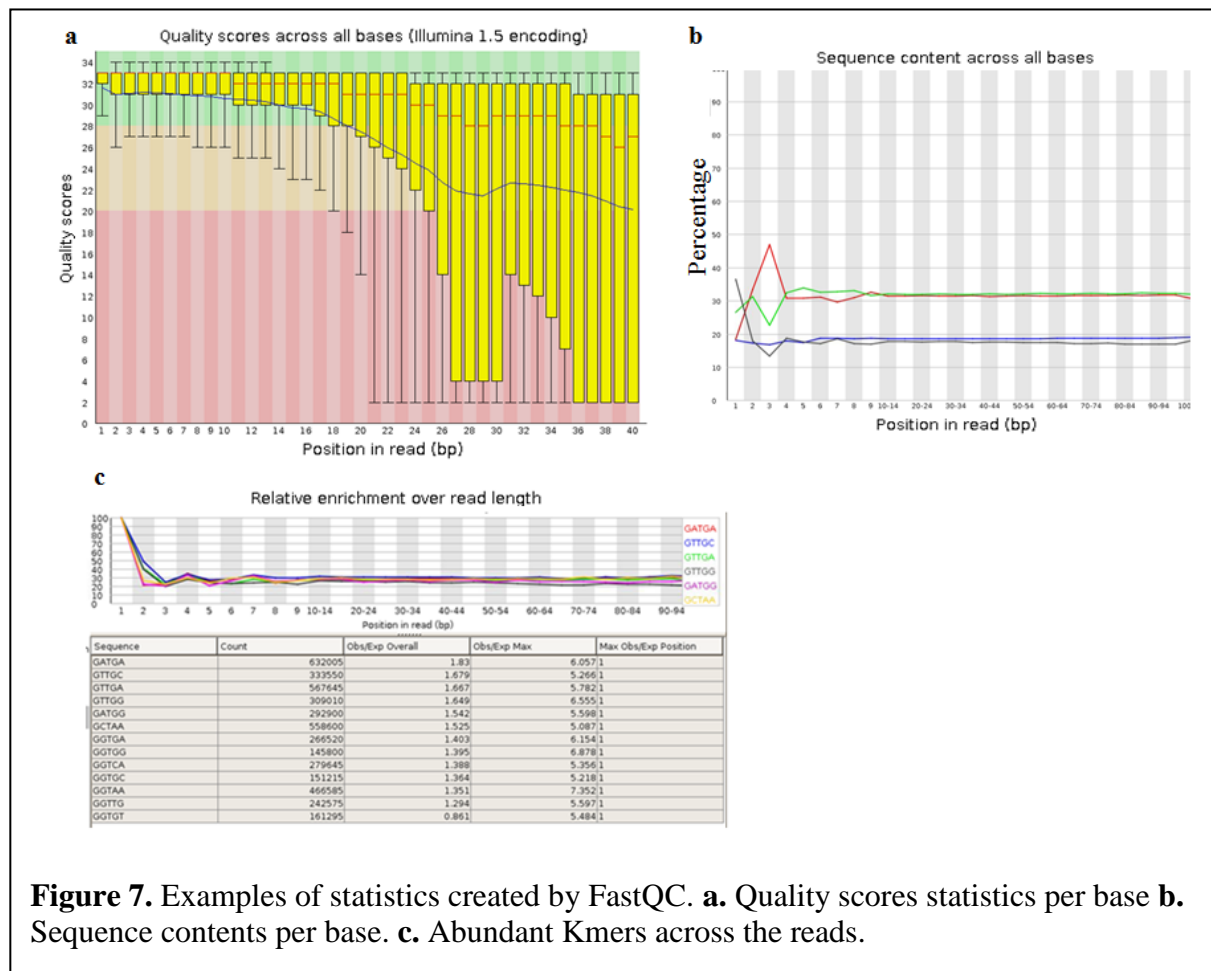


Figure 7. Examples of statistics created by FastQC. **a.** Quality scores statistics per base **b.** Sequence contents per base. **c.** Abundant Kmers across the reads.

Based upon these parameters, we found it advantageous in particular instances to further manipulate the sequences. For example, sequences may be trimmed to reduce low quality ends, filtering reads by quality, and removal of adaptors.

3. Assembling contigs and scaffolds for *de novo* assembly. Available software for this step: SOAPdenovo (10), ABySS (11), Velvet (12), ALL-PATHS (13). *De novo* assembly is the most challenging application and continues to be the subject of intense algorithmic research. The process generally consists of 3 basic steps (Figure 8):

- a. Contig-ing – The first step in the assembly consists of detecting overlaps between single reads. Bundle of overlapping reads are merged into a consensus sequence, called a contig. Repetitive or low complexity regions in the assembled genomic sequence often prevent the construction of one long sequence at this initial step. This step typically results in >10,000 contigs, depending of course on the size of the genome and the number and length of sequenced reads.
- b. Scaffolding – For *de novo* sequencing of complex genomes, it is crucial for the sequenced reads to be of paired-ends inserts. If so, the many contigs can then be merged onto longer segments called scaffolds by taking into account the paired-end information of the reads. Since the paired-end inserts contain an unknown sequence between the two reads, the scaffold may contain unknown sequence (represented as N's) of a size that can be determined by the average insert length.
- c. Gap closing – After creating the scaffolds, the sequence of any remaining gaps within the scaffolds may be resolved by mapping the original paired-end reads to the scaffolds and searching for a read that informs the gap regions. This function may be an integrated process of some assemblers or a separate function may need to be run as in SOAPdenovo.

a. Contig

```
AGGCTGACTT
AGTCTGGCTG
GCTGGCTGA
GGCTGACTTG
ACTTGGTCCA
CTGGGTGACT
```

AGTCTGGCTGACTTGGTCCA **Contig**

b. Scaffolding

```
AGGCTGACTT-----ATGCCGTGGA
AGTCTGGCTG-----GCTATGCCGT
GCTGGCTGA-----GCAAGTGGACA
GGCTGACTTG-----CGCTATGCCG
AACTGTTACT-----ACTTGGTCCA
TGTTAACTGT-----CTGGGTGACT
TGTTAACTGTTACT-----AGTCTGGCTGACTTGGTCCA-----CGCTATGCCGTGGACA Scaffold
```

c. Closing intra-scaffold gaps

```
AGGCTGACTT-----ATGCCGTGGA
AGTCTGGCTG-----GCTATGCCGT
GCTGGCTGA-----GCAAGTGGACA
GGCTGACTTG-----CGCTATGCCG
AACTGTTACT-----ACTTGGTCCA
TGTTAACTGT-----CTGGGTGACT
TGTAGTCTGG-----CCCTGAATGG
AGTCTGGCTG-----TGAATGGTCT
TGTTAACTGTTACT-----AGTCTGGCTGACTTGGTCCA-----CCCTGAATGGTCTCGCTATGCCGTGGACA Scaffold
```

Figure 8. Three basic steps of *de-novo* assembly: **a.** Aligning reads to find overlaps **b.** Connecting contigs into scaffold by using PE information **c.** Closing intra-scaffold gaps

It should be noted that *de novo* assembly projects may include a reference genome of a close strain, or sequences that are known to be included in the assembly, which may help with the assembly process. In this section we will discuss the basic *de-novo* assembly process that does not rely on additional reference sequences.

In the assembly process the identification of sequencing errors is more difficult than it is when mapping reads to a reference genome. Detection of sequencing errors in the process of finding overlaps and merging them into a consensus sequence is possible if there is enough coverage. This is one of the reasons that a higher coverage is required for *de novo* assembly compared to application that consist of a known reference genome.

4. Mapping. Available software for this step: BWA (14), Bowtie (15), TopHat (7). The process of mapping is done in any application that includes a known reference genome. Each read is mapped to the reference genome separately under the conditions of the mapping software, as defined by the input parameters. PE reads are each mapped separately and only then the distance between their mappings is measured.

The main parameters inputted for a mapping software deal with the measure of difference between the read and the reference genome. As in many other bioinformatics methods, deciding on the measure of similarity between reads and the reference genome raises the dilemma between sensitivity and specificity: Allowing too much difference may result in false positive

mappings, while allowing too little difference may lead to missing true positives. From our experience the best way to decide on the parameters is to try a few values and see how they affect the results.

There are two main methods to control the measure of dissimilarity between reads and the reference genome:

- a. Number of differences per read – Apply the mapping software with a value that defines the maximum number of allowed differences (mismatches, insertions and deletions) between the read and the reference genome.
- b. Seed mapping – In this method the software looks for a sequence of certain length inside the read that does not contain differences or contains a small amount of differences compared to the reference genome. The rest of the alignment is elongated without limiting the amount of differences. The parameters given to the software control the seed length, the amount of differences allowed in it and sometimes also the intervals in the reads in which it is searched.

In general, seed mapping is a more permissive approach and is suitable for sequences strains that are distant from the reference genomes they are mapped to. The first method is more strict and is suitable for strains that are known to be close to the reference genome and when trying to avoid false positives. It should be noted that when using the first methods and allowing many differences per read, the results become similar to those that are received in the second methods. The sensitivity and specificity can be tuned also by the parameters of each method.

It is important to remember that the way the mapping step is done affects the rest of the analysis. Allowing a low amount of mismatches may cause regions in the reference genome that contain many variations compared to the sequenced strain to have little to no coverage. Regions in the genome with little to no coverage may be caused by a few reasons. First, region is not present in the sequenced strain – the zero coverage implies a deletion compared to the reference genome. Second, the region does exist in the sequenced strain but is not represented in the sequenced library because of a bias caused in the sample preparation process (for example, because some regions in the genome that are not sonicated as well as others). Finally, the low coverage may also be caused by allowing too few differences per read to a region in the genome that contains many variations in the sequenced strain compared to the reference genome. Trying to map the reads again with a higher percentage of differences may cause these low coverage regions to “fill-up”.

After the mapping is done one can choose to use only a partial set of the mappings:

- a. Use only uniquely mapped reads: It is very common for initial analyses to use only reads that map to one unique location in the genome. Under the mapping conditions, defined by the parameters, reads may be mapped to more than one location in the genome. In this case, one cannot surely determine where the read has originated from. There are a few approaches to deal with such reads – map them randomly to one of the possible locations, map them to all locations, apply an even amount of coverage to every possible location, etc. Each of these approaches may cause a bias in the results, and can be ignored in the initial analysis by using only the uniquely mapped reads.
- b. Use mappings with a minimum mapping score: One can choose to use only mappings of higher quality in order to disregard low quality mappings that may introduce false positives.

- c. Filter mappings with certain insert sizes: PE reads are first mapped separately and only then is the distance between them measured. Long insert sizes, or reads that map to different chromosomes may imply structural variations such as large deletions, inversion and translocations (Figure 5). One can choose to use only mappings with irregular insert size to find such structural variations or use only mappings with normal insert size for initial variant analysis. BreakDancer is an example to a program that uses PE information to find structural variations (16).
- d. Removal of PCR duplicates: PCR amplification is part of the sample preparation, and may introduce bias. PCR duplicates may be identified as reads that map to the exact same location, and in PE reads have the same insert size.

5. Variant calling and filtering. Available software for this step: SAMtools (17), GATK (18), MAQ (19). Based on the mapping done in the previous step, variants can be called by finding the consensus sequence from the mapped reads. The first step in this process is to create a “pileup” file of the mapped reads. This file summarizes for each base in the reference genome the coverage of the reads that are mapped at the loci and the called bases of these reads. Depending on the software that creates the pileup file, more information can be obtained from it, such as genotype calling, mapping qualities, p-values etc. The information in the pileup file can be used to detect and filter variants. The two basic parameters that help detect variants are:

- a. Coverage at the loci – The detected variants should rely on a sufficient coverage. A minimum number of reads should be set as a threshold for initial filtering.
- b. Frequency of the allele that was sequenced – The variant should have sufficient frequency out of the total reads covering the loci. If one read out of 15 reads covering a loci shows a base different from the reference genome, it may not imply a variant but rather a sequencing error. To find a heterozygous variants the frequency should be ~0.5, for homozygous variants the frequency should be ~1, if pooling was done than the frequency should match the expected percentage in the pooling. When filtering by allele frequency taking a margin of security is recommended, especially if the coverage is low. For example, for heterozygous variants filter by a frequency of 0.4 or 0.3.

The above are two basic parameters for variant filtering, but other parameters can be used for variant filtering, for example the mapping and base qualities in the variant locations.

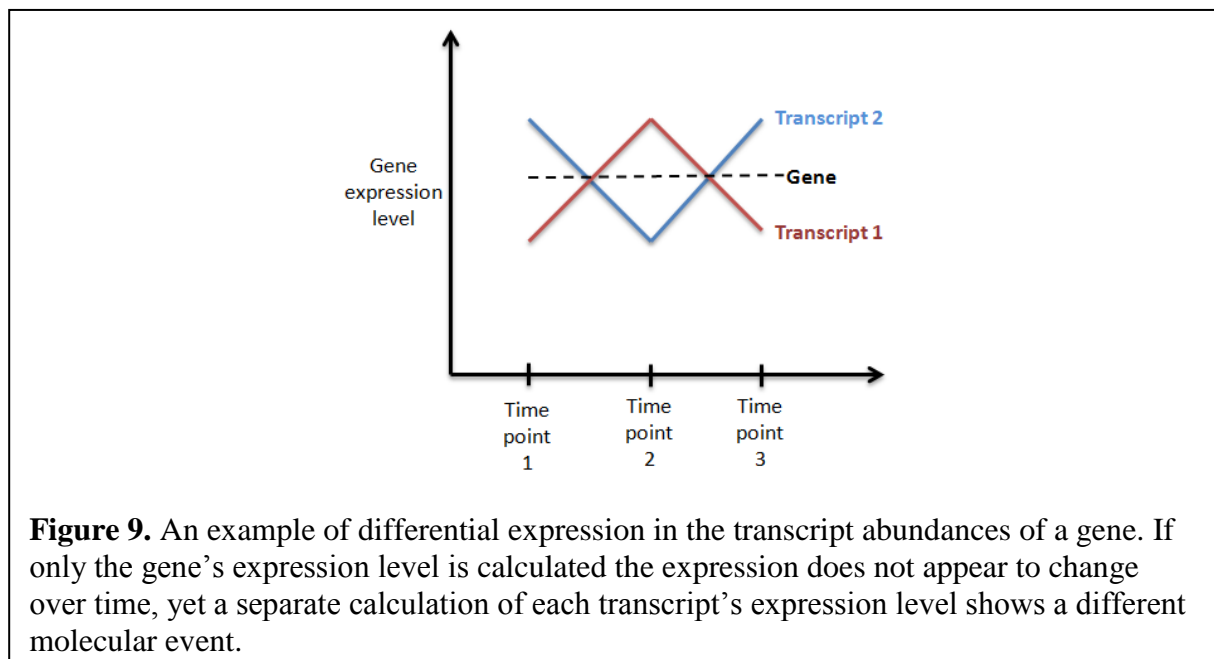
6. Assembling transcripts. In strains that do not have full or sufficient gene annotations, novel annotations can be found by HTS. The idea is to sequence mRNA, map the reads to the reference genome, and infer transcripts from the detected bundles of reads in a certain loci. Based on these annotations a gene expression analysis can then be done. In principle, one can assemble the whole genome before performing RNA-Seq experiment, or assemble the transcriptome only in an application called “de-novo RNA” (20) (or a combination of both).

7. Gene expression analysis. Available software for this step: Cufflinks (7), Myrna (21). After mapping the reads to the reference genome an assessment of their abundance can be made by the gene annotations. In general, the amount of reads that overlap each gene is counted. The raw count must be normalized for further analysis. A common normalization method is called FPKM (Fragments per Kilobase Million) and is calculated as follows:

$$FPKM = \frac{\text{raw count}}{\text{gene length} \cdot \text{number of mapped reads in millions}}$$

The normalization takes into account the gene's length, to avoid a bias toward higher expression in longer genes. FPKM also takes into account the total number of mapped reads in each sample, to avoid a bias because of difference in number of reads in each sample.

A basic approach to gene expression is to count all the reads that map to a gene's annotation, normalize them and set this value as the expression level of that gene. If a gene has more than one transcript due to alternative splicing, not separating the reads that map to it to each of the transcript can cause a great bias and change the results entirely (Figure 9). Finding the expression levels of different transcript of the same gene is challenging, since reads that map to exons that belong to more than one transcript cannot be unambiguously correlated to one transcript (7). The software Cufflinks (7) attempts to assess transcripts expression levels by using the reads that can be unambiguously correlated with certain exons to infer the expression of all the reads (Figure 10). Cufflinks' algorithm uses maximum likelihood to assess the abundances of each transcript.



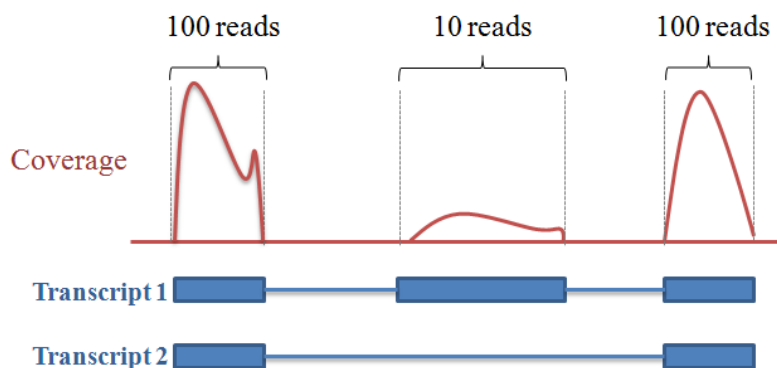


Figure 10. Assessing transcript abundance. Since 10 reads undoubtedly originated from transcript 1, it may be inferred that 90 reads from each shared exon originated from transcript 2 while 10 reads from each shared exon originated from transcript 1.

8. Peaks detection. Available software for this step: MACS (22,23), SICER (24). A ChIP-Seq experiment is done to detect enriched regions in the IP sample. These regions, called “peaks” or “islands”, should be significantly higher both from their surrounding in the IP sample and from the same loci in the control sample. The peaks are found by statistical modeling of the enriched regions compared to the control. There are two important parameters for peaks detection: the abundance in the genome and the width of the binding sites. We introduce two programs for peaks detection, each addressing binding sites with different abundance and width characteristics. MACS is more suitable for narrow peaks that represent short and specific binding sites, for example of transcription factors. SICER is more suitable for wide peaks that extends over thousands of base-pairs, these peaks are typical for histone modification experiments, in which many close binding across the genome. Their proximity to each other makes the peaks merge to wide enriched regions rather than short and sharp peaks.

Tuning up the pipelines

The pipelines detailed above are general. It is crucial to examine each project specifically and decide what pipeline is best suited for it. Tuning up the parameters of each step in the pipeline may be vital for accurate results. Tuning up the pipeline and parameters can be done by following the general pipeline presented above and conducting quality control measurements after each step, to allow identifying a phenomenon that might infer some insight or require special action in the analysis.

Quality control measurements should be done after each step in the analysis. After the raw data handling a quality control step is done as detailed above. After the mapping step the mapping statistics should be assessed. How many reads were not mapped, uniquely mapped and multi-mapped? High values in the first and third parameters may infer a problem. How does the coverage profile look like? What percentage of the genome is covered with sufficient coverage, and what is the average coverage? For exome projects, what is the coverage over exons? It is

highly recommended to look at the mappings in a genome viewer. Some phenomena can be detected easier visually (Figure 11).

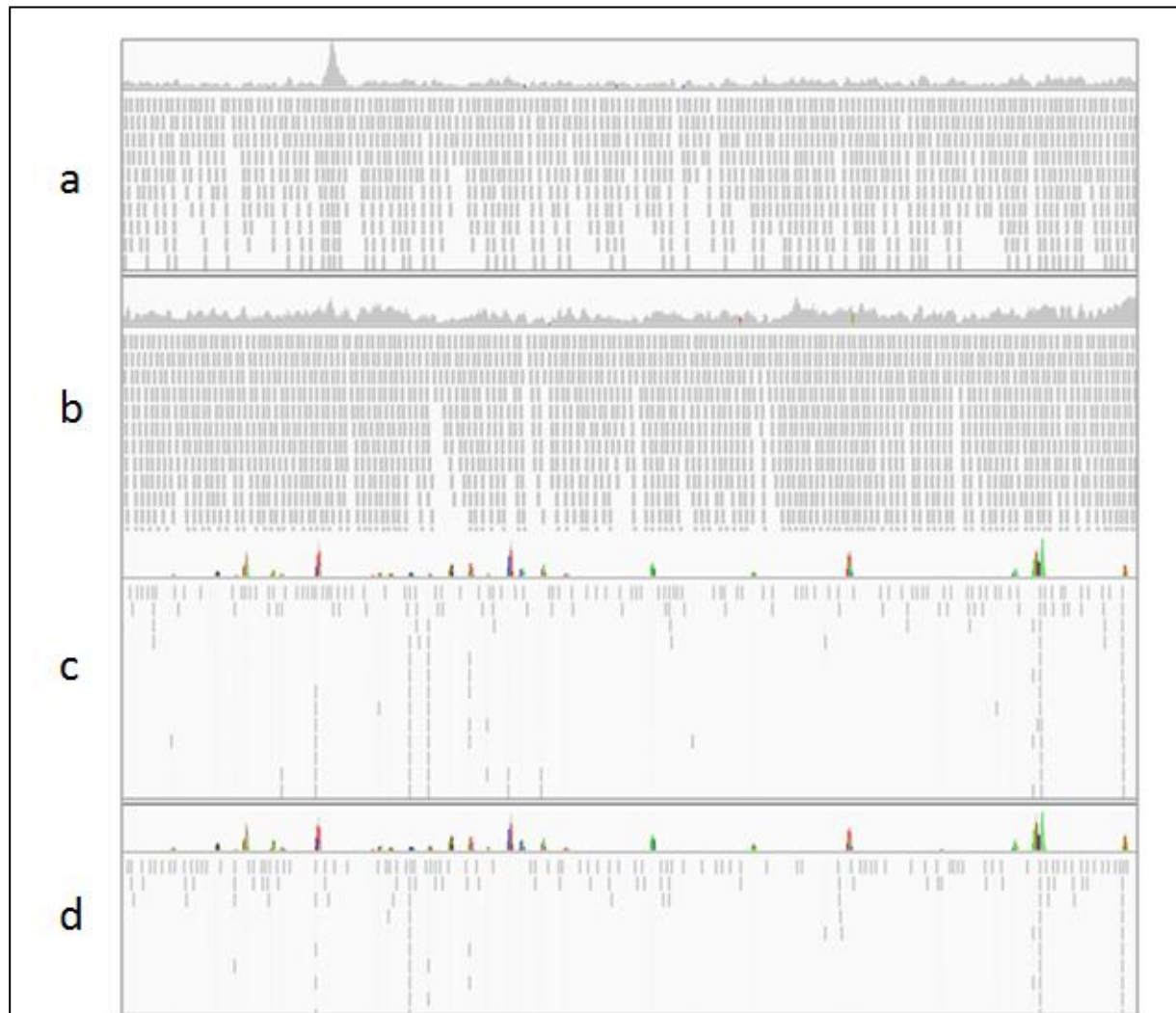


Figure 11. Four bacteria samples were sequenced and mapped to the same reference genome. The mapping statistics of all of them showed that 96-98% of the reads were unmapped. Viewing them on a genome viewer reveals the different phenomenon in each sample. **a,b.** only 2-3% of the reads are of the expected strain, while the rest are contamination. This can be seen by the high and continuous coverage and lack of variants. **c,d.** These sequenced samples seem to be evolutionary distant from the reference genome, as can be seen by the low and segmented coverage and many variants.

Tuning up the parameters in each step of the analysis allows to control the balance between sensitivity and specificity. For example, if we allow one mismatch per 50bp read in the mapping step, it will reduce the rate of incorrect mappings, but we will not be able to detect 2-base indel or areas in the genome that have more than one variant per 50 bp, the coverage in these regions

will be low or zero due to incapability of mapping. Another example from gene expression analysis: when comparing gene expression between two samples one can choose to statistically test only genes that have a minimum amount of reads mapped to them in at least one sample (7). Choosing a high threshold may cause missing interesting genes, but choosing a low threshold may include genes that their differential gene expression is not significant - a gene can be expressed in a fold change of 5 if the ration between the samples is 1 read vs. 5 reads or 1,000 reads vs. 5,000 reads.

Diving into the technical details: file formats

In this section we will overview the formats of some basic files used in HTS data analysis. Though not all useable formats are mentioned here, this section provides a general idea of how the files used in the analysis are constructed, as their structure is similar and the same concepts generally apply. All the files we will present in this section and most of the files used for HTS analysis are plain text files and usually tab delimited, which enables easy management by various tools and scripts.

1. FastQ - raw reads format (fasta + quality). A fastq file is constructed out of quadruplets lines (Figure 12), each quadruplet representing a read and contains the following information:

1) Read identifier – PE reads will have the same identifier. The read's identifier is unique and is constructed in the following way (CASAVA 1.8.2): @<instrument>:<run number>:<flowcell ID>:<lane>:<tile>:<x-pos on tile>:<y-pos on tile>

<read (1/2)>:<is filtered>:<control number>:<index sequence>

2) Sequence.

3) Read description (optional).

4) Quality score per base. Each base is associated with a quality score that defines how reliable tha base is. The score is called Phred quality score and defined as: $P(\text{base is mis called}) = \frac{1}{10^{\frac{Q}{10}}}$

where Q is the quality value. The quality values are typically between 2-50. For example, the quality scores 20, 30 and 4, refer to an error probability of 1/100, 1/1000, and 1/10000, respectively. In order to encode each quality score into one character in the fastq file the following procedure is done: A value is added to the quality score, either 33 or 64, and the new value is then encoded into a character using the ASCII table (Figure 13). The number that is usually being added is 33, while some old CASAVA versions used to add 64 instead.



Figure 12. The anatomy of a Fastq file.

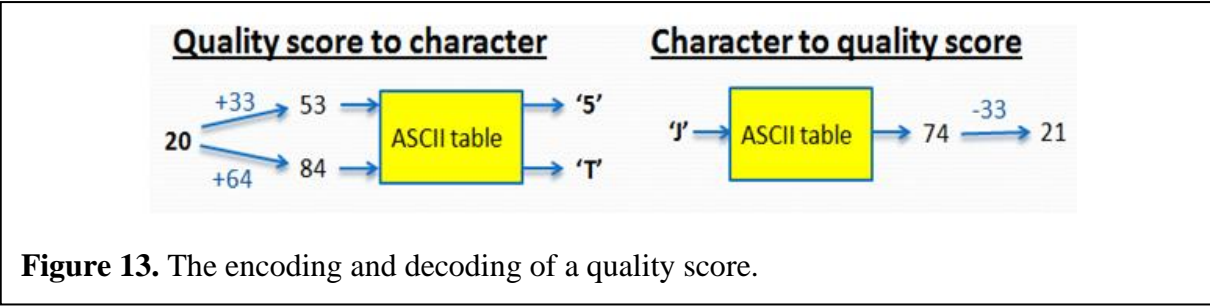


Figure 13. The encoding and decoding of a quality score.

Table 2. Description of the columns of a SAM format file.

Column	Field	Description
1	Read ID	The read's identifier as it appears in the fastq file.
2	Flag	See text.
3	Chromosome	The chromosome to which the read was mapped to. "*" if the read isn't mapped.
4	Position	The position on the chromosome to which the read was mapped to. "*" if the read isn't mapped.
5	Mapping quality	The mapping quality score that was specified by the mapping software. "255" if the mapping quality is not available.
6	CIGAR	See text
7	Mate chromosome	The chromosome to which the read's mate was mapped to. "*" if the mate isn't mapped. "=" if it identical to the chromosome of the read.
8	Mate position	The position on the chromosome to which the read's mate was mapped to. "*" if the mate isn't mapped.
9	Insert length	The distance between the mappings of the two reads, inferring the insert size.
10	Sequence	The read's sequence, as it appears in the fastq file.
11	Quality scores	The read's quality scores, as it appears in the fastq file.
12 (optional)	Program specific attributes	See text

The flag is one number that contains answers to the following 11 YES/NO questions regarding the read's mapping:

1. Is the read paired
2. Is the read mapped in proper pair (in the expected insert length)
3. Is the read unmapped
4. Is the mate unmapped
5. Is the read mapped to the reverse strand
6. Is the mate mapped to the reverse strand
7. Is the read the first in the pair
8. Is the read the second in the pair
9. Is the mapping not a primary alignment
10. Did the read fails platform/vendor quality checks
11. Is the read a PCR or an optical duplicate

To encode the answers to these questions to one number, a NO answer is encoded as "0", and a YES answer is encoded as "1". The binary number resulting from the serie of answers is then converted to a decimal number (see Figure 15 for an example).

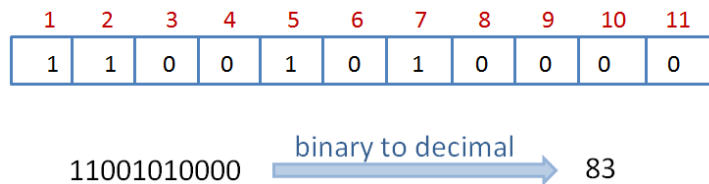


Figure 15. The flag in the SAM file is binary encoded with the following 11 bits of information: read is paired, read is mapped in a proper pair, read is on reverse strand, read is the first in the pair.

3. CIGAR (Compact Idiosyncratic Gapped Alignment Report). Details the mapping structure between the read and reference genome, according to a specific encoding. As an example, ‘M’ corresponds to a matched alignment. Thus, 101M means that there were 101 matches or mismatches without gaps opened, and 73M1I27M means that the first 73 bases were a match or mismatch compared to the reference genome, then there was one base insertion, and then another 27 matches/mismatches.

The optional attributes detail more information about the mapping. Some of the options are the edit distance in the mapping, mismatches positions, number of gap openings. These attributes are tab delimited and will be of the form <Tag>:<Type>:<Details>

- Tag – identifies what kind of information is detailed, according to the SAM specification
- Type- I for integer, Z for string
- Details – the details themselves

For example, NM:i:3 means an edit distance of 3 and MD:Z:74G26 means that there is a mismatch in the 75th position of the reads, a ‘G’ instead of the reference base.

The SAM specification defines some of these options, and reserves attribute for program specific needs, the reserved options start with X. These attribute should be defined in the mapper’s documentations. SAMtools (17) is a program that enables manipulation, conversion and data retrieval from SAM files.

4. VCF – Variant Call Format. A VCF file details information per base of the reference genome, accumulated from the mappings in a SAM file (Figure 16). A VCF file is a tab delimited file, also constructed from header lines and variant lines. The header lines begin with ‘#’ character and detail general information on the file, such as the program used for the variant calling process, and the attributes that appear in each variant line. Each line in the rest of the document contains information about a specific base in the genome. Only bases that have a coverage of at least one appear in a VCF file. A “raw” VCF file contains information about every base with coverage in the genome. It can then be filtered to contain bases that define a variant compared to the reference genome. VCFtools (25) is a program package designed to working with VCF files.

```
#CHROM POS ID REF ALT QUAL FILTER INFO FORMAT
  _genome 1 . A . 111 . DP=27;AF1=0;AC1=0;DP4=12,15,0,0;MQ=37;FQ=-108 PL 0
  _genome 2 . A . 174 . DP=48;AF1=0;AC1=0;DP4=27,21,0,0;MQ=37;FQ=-171 PL 0
  _genome 3 . T . 220 . DP=63;AF1=0;AC1=0;DP4=38,25,0,0;MQ=37;FQ=-217 PL 0
  _genome 4 . A . 283 . DP=85;AF1=0;AC1=0;DP4=49,36,0,0;MQ=37;FQ=-282 PL 0
  _genome 5 . C . 283 . DP=105;AF1=0;AC1=0;DP4=60,45,0,0;MQ=37;FQ=-282 PL 0
  _genome 6 . A . 283 . DP=122;AF1=0;AC1=0;DP4=69,53,0,0;MQ=37;FQ=-282 PL 0
  _genome 7 . A . 283 . DP=146;AF1=0;AC1=0;DP4=73,72,0,0;MQ=37;FQ=-282 PL 0
  _genome 8 . C . 283 . DP=168;AF1=0;AC1=0;DP4=82,86,0,0;MQ=37;FQ=-282 PL 0
  _genome 9 . A . 283 . DP=190;AF1=0;AC1=0;DP4=91,98,0,0;MQ=37;FQ=-282 PL 0
  _genome 10 . A . 283 . DP=212;AF1=0;AC1=0;DP4=98,114,0,0;MQ=37;FQ=-282 PL 0

##INFO=<ID=DP,Number=1,Type=Integer,Description="Raw read depth">
##INFO=<ID=DP4,Number=4,Type=Integer,Description="# high-quality
ref-forward bases, ref-reverse, alt-forward and alt-reverse bases">
```

Figure 16. The anatomy of a VCF file

5. GFF – General Feature Format. A GFF file contains details about annotations of a specific genome sequence (Figure 17). A GFF file should be of the same build or version of the genome sequence it refers to. A GFF file is constructed of header lines that begin with a “#” character and feature lines. The feature lines are tab-delimited and contain the attributes shown in Table 3.

```
##gff-version 3
##sequence-region SL2.30ch00 1 21839854
SL2.30ch00 EuGene gene 16480 17940 . + . gene_id "Solyc00g005000"; transcript_id "Solyc00g005000"
SL2.30ch00 EuGene mRNA 16480 17940 . + . gene_id "Solyc00g005000.1.1"; transcript_id "Solyc00g005000.1.1"
SL2.30ch00 EuGene exon 16480 17275 . + . gene_id "Solyc00g005000.1.1"; transcript_id "Solyc00g005000.1.1.1"
SL2.30ch00 EuGene exon 17336 17940 . + . gene_id "Solyc00g005000.1.1"; transcript_id "Solyc00g005000.1.1.2"
SL2.30ch00 EuGene CDS 16480 17275 . + 0 gene_id "Solyc00g005000.1.1"; transcript_id "Solyc00g005000.1.1.1"
SL2.30ch00 EuGene CDS 17336 17940 . + 2 gene_id "Solyc00g005000.1.1"; transcript_id "Solyc00g005000.1.1.2"
SL2.30ch00 EuGene gene 66298 67449 . + . gene_id "Solyc00g005010"; transcript_id "Solyc00g005010"
SL2.30ch00 EuGene mRNA 66298 67449 . + . gene_id "Solyc00g005010.1.1"; transcript_id "Solyc00g005010.1.1"
SL2.30ch00 EuGene exon 66298 66567 . + . gene_id "Solyc00g005010.1.1"; transcript_id "Solyc00g005010.1.1.1"
SL2.30ch00 EuGene exon 66779 66946 . + . gene_id "Solyc00g005010.1.1"; transcript_id "Solyc00g005010.1.1.2"
SL2.30ch00 EuGene exon 67033 67449 . + . gene_id "Solyc00g005010.1.1"; transcript_id "Solyc00g005010.1.1.3"
```

Figure 17. The anatomy of a GFF file

Table 3. Description of the columns of a GFF format file.

Column	Field	Description
1	Chromosome	The chromosome on which the feature is located.
2	Source	The source of this feature, usually the prediction software or a public DB.
3	Feature	The feature type.
4	Start	The start position on the chromosome on which the feature is located.
5	End	The end position on the chromosome on which the feature is located.
6	Score	A floating point value
7	Strand	The strand the feature is originated from (‘+’, ‘-’ or ‘.’ If the strand isn’t relevant).
8	Frame	The position of the feature in the ORF (‘0’, ‘1’, ‘2’ or ‘.’)
9	Attributes	More details about the feature, separated by ‘;’. For example gene ID, gene description, exon number, description.

The file formats detailed above, similarly to other files used in HTS data analysis, enable easy retrieval of information using simple scripts or public programs. Knowing how the data is stored and where enables to ask questions such as:

- Which sequences were not mapped? (look for lines with bit4 in the flag equals to 0 in a SAM file)
- What is the average coverage in a certain region in the genome (Calculate the average of the DP values in a region in a VCF file)
- What kind of annotations are known in a reference genome (find the possible options in column 3 in a GFF file)

References

1. Goh, Y., Fullwood, M.J., Poh, H.M., Peh, S.Q., Ong, C.T., Zhang, J., Ruan, X. and Ruan, Y. (2012) Chromatin Interaction Analysis with Paired-End Tag Sequencing (ChIA-PET) for mapping chromatin interactions and understanding transcription regulation. *J Vis Exp*.
2. Toung, J.M., Morley, M., Li, M. and Cheung, V.G. RNA-sequence analysis of human B-cells. *Genome Res*, **21**, 991-998.
3. Koehler, R., Issac, H., Cloonan, N. and Grimmond, S.M. The uniqueome: a mappability resource for short-tag sequencing. *Bioinformatics*, **27**, 272-274.
4. Paszkiewicz, K. and Studholme, D.J. (2010) De novo assembly of short sequence reads. *Brief Bioinform*, **11**, 457-472.
5. Schatz, M.C., Witkowski, J. and McCombie, W.R. (2012) Current challenges in de novo plant genome sequencing and assembly. *Genome Biol*, **13**, 243.
6. McIntyre, L.M., Lopiano, K.K., Morse, A.M., Amin, V., Oberg, A.L., Young, L.J. and Nuzhdin, S.V. (2011) RNA-seq: technical variability and sampling. *BMC Genomics*, **12**, 293.
7. Trapnell, C., Williams, B.A., Pertea, G., Mortazavi, A., Kwan, G., van Baren, M.J., Salzberg, S.L., Wold, B.J. and Pachter, L. (2010) Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation. *Nat Biotechnol*, **28**, 511-515.
8. Kharchenko, P.V., Tolstorukov, M.Y. and Park, P.J. (2008) Design and analysis of ChIP-seq experiments for DNA-binding proteins. *Nat Biotechnol*, **26**, 1351-1359.
9. Auerbach, R.K., Euskirchen, G., Rozowsky, J., Lamarre-Vincent, N., Moqtaderi, Z., Lefrancois, P., Struhl, K., Gerstein, M. and Snyder, M. (2009) Mapping accessible chromatin regions using Sono-Seq. *Proc Natl Acad Sci U S A*, **106**, 14926-14931.
10. Li, R., Zhu, H., Ruan, J., Qian, W., Fang, X., Shi, Z., Li, Y., Li, S., Shan, G., Kristiansen, K. *et al.* (2010) De novo assembly of human genomes with massively parallel short read sequencing. *Genome Res*, **20**, 265-272.
11. Simpson, J.T., Wong, K., Jackman, S.D., Schein, J.E., Jones, S.J. and Birol, I. (2009) ABySS: a parallel assembler for short read sequence data. *Genome Res*, **19**, 1117-1123.
12. Zerbino, D.R. and Birney, E. (2008) Velvet: algorithms for de novo short read assembly using de Bruijn graphs. *Genome Res*, **18**, 821-829.
13. Butler, J., MacCallum, I., Kleber, M., Shlyakhter, I.A., Belmonte, M.K., Lander, E.S., Nusbaum, C. and Jaffe, D.B. (2008) ALLPATHS: de novo assembly of whole-genome shotgun microreads. *Genome Res*, **18**, 810-820.
14. Li, H. and Durbin, R. (2009) Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics*, **25**, 1754-1760.
15. Langmead, B., Trapnell, C., Pop, M. and Salzberg, S.L. (2009) Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biol*, **10**, R25.
16. Chen, K., Wallis, J.W., McLellan, M.D., Larson, D.E., Kalicki, J.M., Pohl, C.S., McGrath, S.D., Wendl, M.C., Zhang, Q., Locke, D.P. *et al.* (2009) BreakDancer: an algorithm for high-resolution mapping of genomic structural variation. *Nat Methods*, **6**, 677-681.

17. Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G. and Durbin, R. (2009) The Sequence Alignment/Map format and SAMtools. *Bioinformatics*, **25**, 2078-2079.
18. McKenna, A., Hanna, M., Banks, E., Sivachenko, A., Cibulskis, K., Kernytsky, A., Garimella, K., Altshuler, D., Gabriel, S., Daly, M. *et al.* (2010) The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res*, **20**, 1297-1303.
19. Li, H., Ruan, J. and Durbin, R. (2008) Mapping short DNA sequencing reads and calling variants using mapping quality scores. *Genome Res*, **18**, 1851-1858.
20. Grabherr, M.G., Haas, B.J., Yassour, M., Levin, J.Z., Thompson, D.A., Amit, I., Adiconis, X., Fan, L., Raychowdhury, R., Zeng, Q. *et al.* (2010) Full-length transcriptome assembly from RNA-Seq data without a reference genome. *Nat Biotechnol*, **29**, 644-652.
21. Langmead, B., Hansen, K.D. and Leek, J.T. (2010) Cloud-scale RNA-sequencing differential expression analysis with Myrna. *Genome Biol*, **11**, R83.
22. Zhang, Y., Liu, T., Meyer, C.A., Eeckhoute, J., Johnson, D.S., Bernstein, B.E., Nusbaum, C., Myers, R.M., Brown, M., Li, W. *et al.* (2008) Model-based analysis of ChIP-Seq (MACS). *Genome Biol*, **9**, R137.
23. Feng, J., Liu, T. and Zhang, Y. (2011) Using MACS to identify peaks from ChIP-Seq data. *Curr Protoc Bioinformatics*, **Chapter 2**, Unit 2 14.
24. Zang, C., Schones, D.E., Zeng, C., Cui, K., Zhao, K. and Peng, W. (2009) A clustering approach for identification of enriched domains from histone modification ChIP-Seq data. *Bioinformatics*, **25**, 1952-1958.
25. Danecek, P., Auton, A., Abecasis, G., Albers, C.A., Banks, E., DePristo, M.A., Handsaker, R.E., Lunter, G., Marth, G.T., Sherry, S.T. *et al.* (2011) The variant call format and VCFtools. *Bioinformatics*, **27**, 2156-2158.
26. Robinson, J.T., Thorvaldsdottir, H., Winckler, W., Guttman, M., Lander, E.S., Getz, G. and Mesirov, J.P. (2011) Integrative genomics viewer. *Nat Biotechnol*, **29**, 24-26.