

Como processar os dados de RNA-Seq?

Dr. Pablo Rodrigo Sanches

O que é isso para você?



```
>Sequence_1 assembly1  
CCCTAAACCCTAAACCCTAAACCCTAAACCTCTGAATCCTTAATCCCTAAATCCCTAAAT  
CTTTAAATCCTACATCCATGAATCCCTAAATACCTAATTCCTAAACCCGAAACCGGTTT  
CTCTGGTTGAAAATCATTGTGTATATAATGATAATTTTATCGTTTTTATGTAATTGCTTA  
TTGTTGTGTGTAGATTTTTTAAAAATATCATTTGAGGTCAATACAAATCCTATTTCTTGT  
GGTTTTCTTTCCTTCACTTAGCTATGGATGGTTTATCTTCATTTGTTATATTGGATACAA  
GCTTTGCTACGATCTACATTTGGGAATGTGAGTCTCTTATTGTAACCTTAGGGTTGGTTT  
ATCTCAAGAATCTTATTAATTGTTTGGACTGTTTATGTTTGGACATTTATTGTCATTCTT
```

Uma sequência de caracteres?
Um arquivo texto?



vs.

Um Gene?
A parte de um genoma?



E isso?

```
while(my $seq = $seqio->next_seq) {  
  @idseqs = `cat $ARGV[0]`;  
  foreach $idseq (@idseqs)  
  {  
    chomp $idseq;  
    if($seq->desc =~ /$idseq/)  
    {  
      if($seq->desc =~ m/$ARGV[3]/) {  
        $pos = $-[0];  
      }  
      $desc = substr($seq->desc,$pos);  
      print INFO ">" . $idseq . " " . $desc . "\n";  
      print INFO $seq->seq . "\n";  
    }  
  }  
}
```

Um texto em uma língua desconhecida?
Palavras sem sentido organizadas?

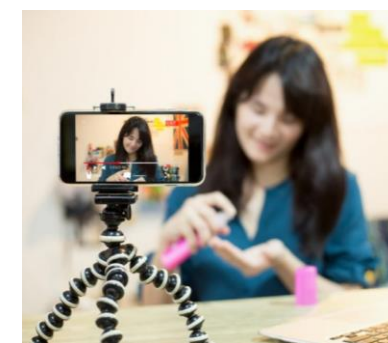


VS.

Um algoritmo?
Um código-fonte escrito em linguagem Perl?



Mercado de trabalho



Piloto de drones, Engenheiro de robôs, Youtuber, Streamer gamer, Cyber atleta, Influenciador digital, Cientista de dados, Técnico de saúde assistida por Inteligência Artificial, Walker/Talker, ...

E ainda...

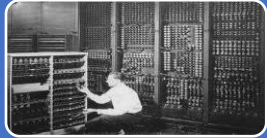
- Bioinformata → Bioinformática
 - Etimologia
 - Bio = “*bios*” (vida) + Informática = “*informatik*” (informação automática)
 - Não é exatamente nova, porém pouco conhecida... Ainda 🤔



Existe uma relação entre Biologia e Informática?



Histórico



A história começa na década de 1940 com a invenção do moderno computador digital



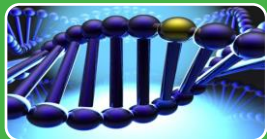
Ele se chama digital, pois os dados são armazenados com um alfabeto binário (0 e 1)



A descoberta da dupla hélice, em 1953, por Watson e Crick, mostrou que a informação genética também é armazenada de forma digital



Mas diferente do alfabeto binário dos computadores, os dados genéticos são armazenados com um alfabeto quaternário (A, C, G e T)

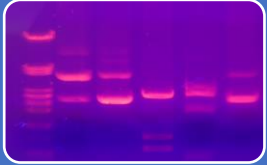


Mais tarde se descobriu que a forma dos genes operarem também é digital. Até certo ponto, os genes podem ser “ligados” ou “desligados”



Apenas estas observações já seriam suficientes para prever, na década de 1950, que um dia informática e biologia molecular iriam juntas fazer nascer uma nova área de conhecimento

Histórico



Apesar da estrutura do DNA ter sido desvendada em 1953, a informação nela contida não podia ser “lida”



Foi preciso esperar até fins da década de 1980 para que aparecesse uma “lente de aumento” suficientemente boa que permitisse a leitura do DNA em grande quantidade



Na década de 1970 a unidade básica de armazenamento de informação era o kilobyte - aproximadamente 1000 letras;

Um computador da época tinha alguns kbytes de memória;

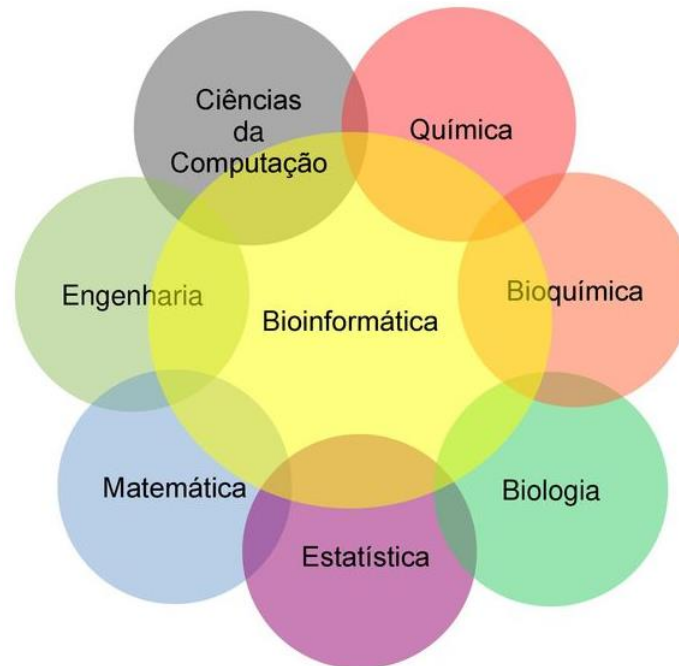
Com tal memória um computador desses não seria capaz de processar nem sequer o genoma de um vírus (20 kb), ou 20 mil letrinhas; que dirá o genoma humano, com seus 3 bilhões de letrinhas.



Foi preciso esperar alguns anos para que essas duas áreas alcançassem formas de produzir a biologia em larga-escala. Produção de dados em massa gera demanda para análises computacionais => Bioinformática.

Bioinformática

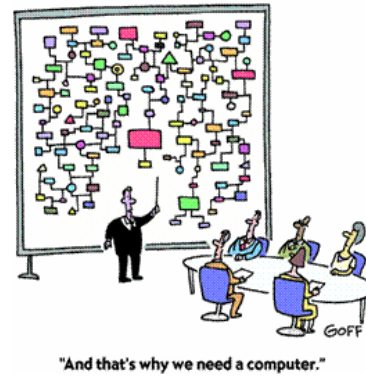
- “...aplicação das **técnicas da informática**, no sentido de **análise da informação** na área de estudo da **biologia**...”;
- “...a utilização de **técnicas computacionais e matemáticas** relacionadas ao **conhecimento químico, físico, biológico, biomédico,...**”.



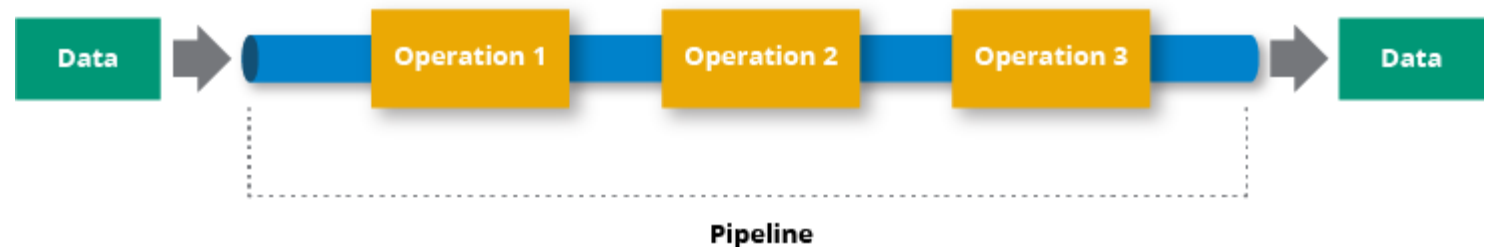
Desafios

- Sequenciamento
- *Base Calling*
- *Qualidade do sequenciamento*
- *Alinhamento*/Montagem
- Predição/Anotação
- Vias metabólicas
- *Expressão diferencial*
- *Splicing alternativo*
- Identificação de mutações
- Filogenia
- Regiões promotoras

- Regiões não-codificantes
- *RNA-Seq*, miRNA-Seq, ChiP-Seq
- Domínios de proteínas
- Bioinformática estrutural
- Bancos de dados biológicos
- *Big Data*
- *Machine Learning*
- Biologia sintética
- Medicina personalizada
- ...



Como processar os dados de RNA-Seq?



O que é RNA-Seq?

- Nome dado às novas tecnologias de sequenciamento (*Next-generation sequencing*) aplicadas aos transcriptomas, ou seja, às regiões do DNA transcritas em moléculas de RNAs.
- Podemos fazer:
 - Analisar a expressão diferencial em diferentes tecidos ou condições ambientais;
 - Analisar diferentes isoformas (*alternative splicing*);
 - Descobrir novas regiões dos genomas que são transcritas;
 - Identificar moléculas de RNAs que participam de processos regulatórios;
 - etc.

Tipos de bibliotecas de sequenciamento

- Bibliotecas de fragmentos (***single-end***)
 - Resultam no sequenciamento de apenas uma das extremidades do fragmento, sendo a metodologia mais simples e barata.
- Bibliotecas de extremidades pareadas (***paired-end***)
 - Resultam em duas leituras para cada fragmento sequenciado, uma referente à fita *forward* e outra à fita *reverse*.

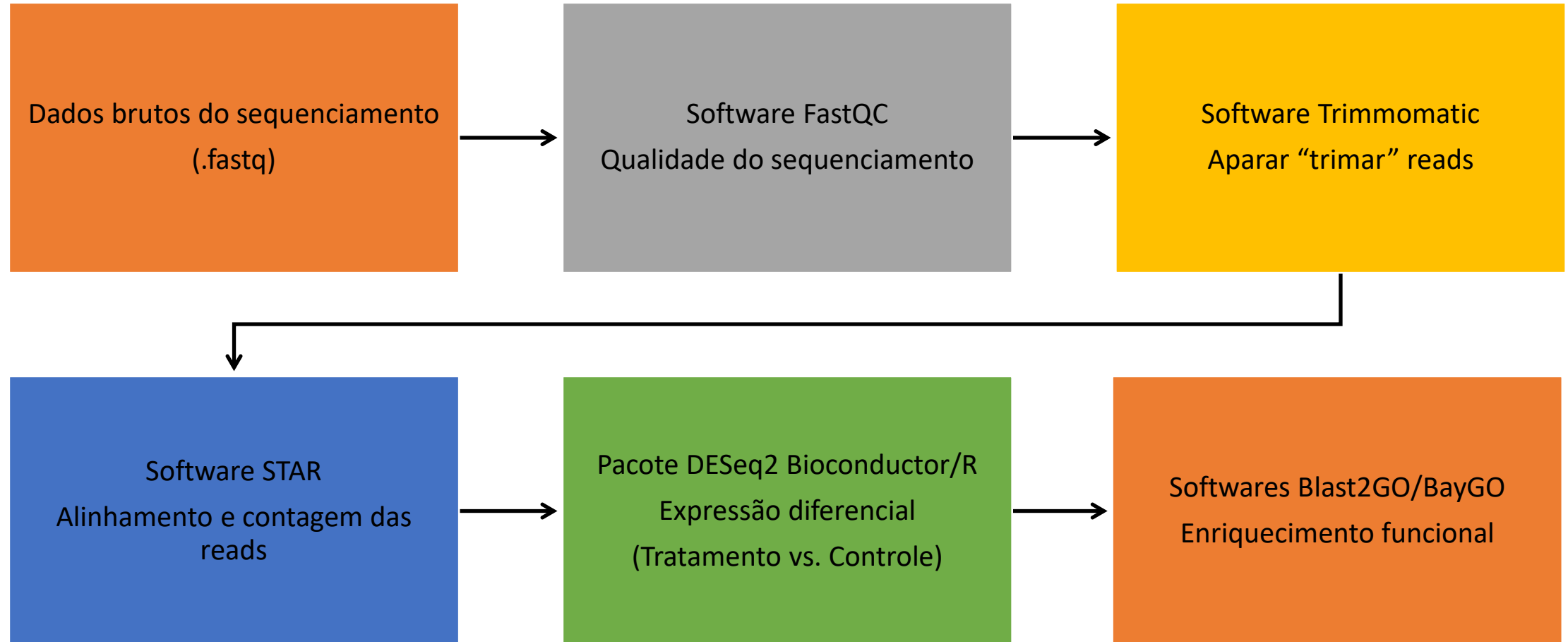
Single-end



Paired-end



Exemplo de pipeline (RNA-Seq) – 6 passos



Dados brutos do sequenciamento

- Dados gerados pelos sequenciadores automáticos;
- Exemplos de formatos de arquivos gerados:
 - FASTA
 - SFF (ROCHE 454)
 - CSFASTA (ABI SOLiD)
 - FASTQ
 - ...

Formato FASTQ

```
@SRR001666.1 071112_SLXA-EAS1_s_7:5:1:817:345 length=36
GGGTGATGGCCGCTGCCGATGGCGTCAAATCCCACC
+SRR001666.1 071112_SLXA-EAS1_s_7:5:1:817:345 length=36
IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII9IG9IC
```

Cada read é representada por 4 linhas no arquivo

```
@ read ID
Sequência
+ read ID
Qualidade
```


Dados brutos disponíveis na web

- GEO Gene Expression Omnibus
 - <http://www.ncbi.nlm.nih.gov/geo/>
- Exemplo:
 - *A transcriptome survey of Trichophyton rubrum exposed to undecanoic acid*
 - <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE102872>
- Após download converter arquivos SRA para FASTQ
 - \$ `fastq-dump -B arquivo.sra`

Shell do Linux

```
mars@marsmain ~$ pwd
/home/mars
mars@marsmain ~$ cd /usr/portage/app-shells/bash
mars@marsmain /usr/portage/app-shells/bash$ ls -al
total 130
drwxr-xr-x 3 portage portage 1024 Jul 25 10:06 .
drwxr-xr-x 33 portage portage 1024 Aug 7 22:39 ..
-rw-r--r-- 1 root root 35808 Jul 25 10:06 ChangeLog
-rw-r--r-- 1 root root 27002 Jul 25 10:06 Manifest
-rw-r--r-- 1 portage portage 4645 Mar 23 21:37 bash-3.1_p17.ebuild
-rw-r--r-- 1 portage portage 5977 Mar 23 21:37 bash-3.2_p39.ebuild
-rw-r--r-- 1 portage portage 6151 Apr 5 14:37 bash-3.2_p48-r1.ebuild
-rw-r--r-- 1 portage portage 5980 Mar 23 21:37 bash-3.2_p48.ebuild
-rw-r--r-- 1 portage portage 5643 Apr 5 14:37 bash-4.0_p10-r1.ebuild
-rw-r--r-- 1 portage portage 5640 Apr 5 14:37 bash-4.0_p10.ebuild
-rw-r--r-- 1 portage portage 5532 Apr 14 05:52 bash-4.0_p17-r1.ebuild
-rw-r--r-- 1 portage portage 5660 May 30 03:35 bash-4.0_p17.ebuild
-rw-r--r-- 1 portage portage 5660 Jul 25 09:43 bash-4.0_p24.ebuild
-rw-r--r-- 1 root root 5660 Jul 25 09:43 bash-4.0_p28.ebuild
drwxr-xr-x 2 portage portage 2048 May 30 03:35 files
-rw-r--r-- 1 portage portage 468 Feb 9 04:35 metadata.xml
mars@marsmain /usr/portage/app-shells/bash$ cat metadata.xml
<?xml version="1.0" encoding="UTF-8"?>
<!DOCTYPE pkgmetadata SYSTEM "http://www.gentoo.org/dtd/metadata.dtd">
<pkgmetadata>
  <herd>base-system</herd>
  <use>
    <flag name="bashlogger">Log ALL commands typed into bash; should ONLY be
      used in restricted environments such as honeypots</flag>
    <flag name="net">Enable /dev/tcp/host/port redirection</flag>
    <flag name="plugins">Add support for loading builtins at runtime via
      'enable' </flag>
  </use>
</pkgmetadata>
mars@marsmain /usr/portage/app-shells/bash$ sudo /etc/init.d/bluetooth status
Password:
* status: started
mars@marsmain /usr/portage/app-shells/bash$ ping -q -c1 en.wikipedia.org
PING rr.esams.wikimedia.org (91.198.174.2) 56(84) bytes of data.

--- rr.esams.wikimedia.org ping statistics ---
1 packets transmitted, 1 received, 0% packet loss, time 2ms
rtt min/avg/max/mdev = 49.820/49.820/49.820/0.000 ms
mars@marsmain /usr/portage/app-shells/bash$ grep -i /dev/sda /etc/fstab | cut --fields=3
/dev/sda1      /boot
/dev/sda2      none
/dev/sda3      /
mars@marsmain /usr/portage/app-shells/bash$ date
Sat Aug 8 02:42:24 MSD 2009
mars@marsmain /usr/portage/app-shells/bash$ lsmdu
Module      Size  Used by
rndis_wlan  23424 0
rndis_host   8696 1 rndis_wlan
cdc_ether    5672 1 rndis_host
usbnet      18688 3 rndis_wlan,rndis_host,cdc_ether
parport_pc  38424 0
fglrx       2388128 20
parport      39648 1 parport_pc
iTCO_wdt     12272 0
i2c_i801     9380 0
mars@marsmain /usr/portage/app-shells/bash$
```

- Interpretador de comandos do Linux;
- Os comandos digitados pelo usuário podem ser comandos internos (embutidos) do *shell*, mas na maioria das vezes eles são programas externos.

Qualidade do sequenciamento

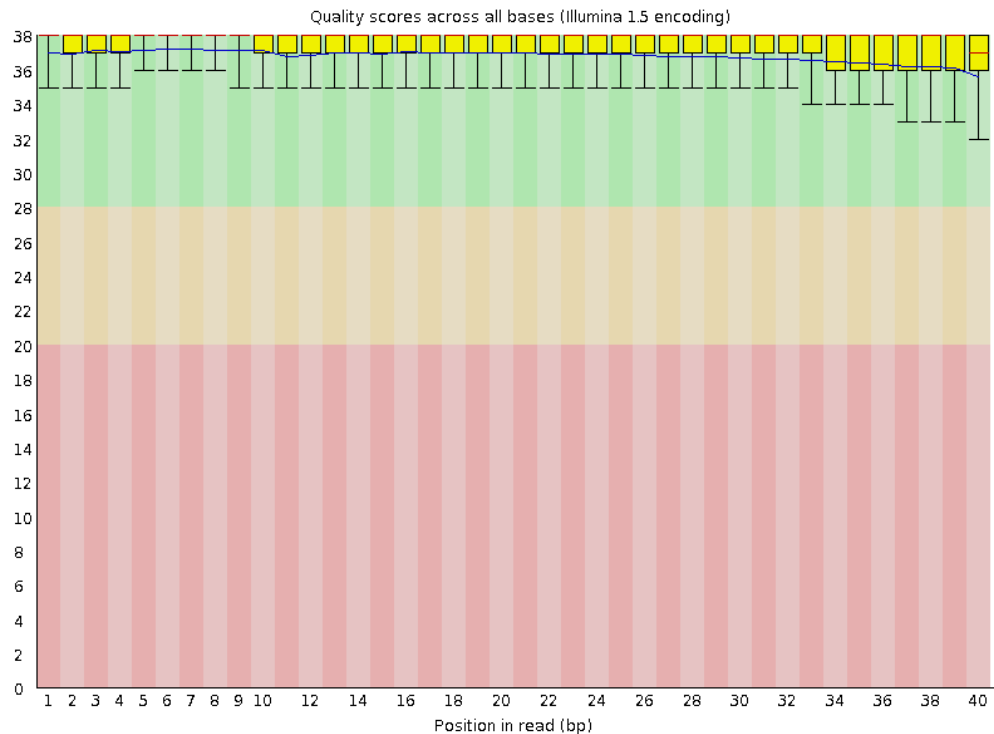
- Avaliar a qualidade das reads
 - Identificar contaminantes;
 - Identificar amostras com baixa performance de sequenciamento.
 - Softwares: **FastQC**, SAMStat, ...
- Download do software FastQC
 - <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>
- Exemplo de uso do software FastQC

```
$ fastqc arquivo.fastq
```
- Resultado do software FastQC
 - Serão gerados arquivos formato .html para visualização em navegador web

Exemplos de resultados - FASTQC

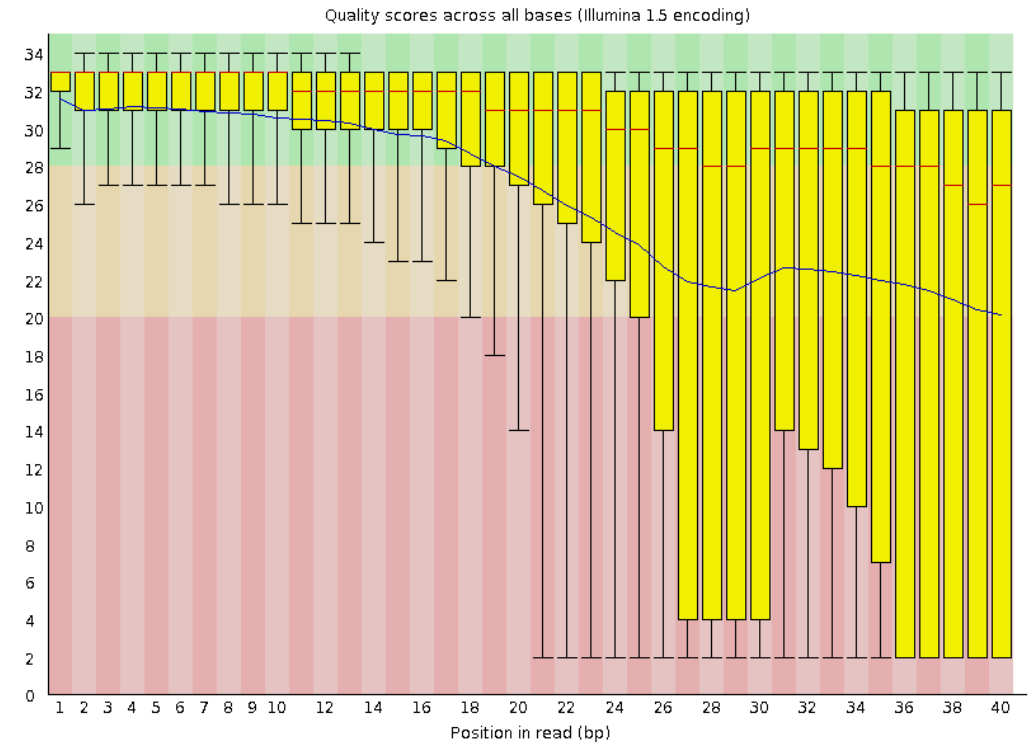
Qualidade boa

✅ Per base sequence quality




Qualidade ruim

❌ Per base sequence quality



Exemplos de resultados - FASTQC

Qualidade boa

 **Overrepresented sequences**
No overrepresented sequences

Qualidade ruim

 **Overrepresented sequences**

Sequence	Count	Percentage	Possible Source
AGAGTTTATCGCTTCCATGACGCAGAAGTTAACACTTTC	2065	0.5224039181558763	No Hit
GATTGGCGTATCCAACCTGCAGAGTTTATCGCTTCCATG	2047	0.5178502762542754	No Hit
ATTGGCGTATCCAACCTGCAGAGTTTATCGCTTCCATGA	2014	0.5095019327680071	No Hit
CGATAAAATGATTGGCGTATCCAACCTGCAGAGTTTAT	1913	0.4839509420979134	No Hit
GTATCCAACCTGCAGAGTTTATCGCTTCCATGACGCAGA	1879	0.47534961850600066	No Hit
AAAAATGATTGGCGTATCCAACCTGCAGAGTTTATCGCT	1846	0.4670012750197325	No Hit
TGATTGGCGTATCCAACCTGCAGAGTTTATCGCTTCCAT	1841	0.46573637449150995	No Hit
AACCTGCAGAGTTTATCGCTTCCATGACGCAGAAGTTAA	1836	0.46447147396328753	No Hit
GATAAAATGATTGGCGTATCCAACCTGCAGAGTTTATC	1831	0.4632065734350651	No Hit
AAATGATTGGCGTATCCAACCTGCAGAGTTTATCGCTTC	1779	0.45005160794155147	No Hit
ATGATTGGCGTATCCAACCTGCAGAGTTTATCGCTTCCA	1779	0.45005160794155147	No Hit
AATGATTGGCGTATCCAACCTGCAGAGTTTATCGCTTCC	1760	0.4452449859343061	No Hit
AAAATGATTGGCGTATCCAACCTGCAGAGTTTATCGCTT	1729	0.4374026026593269	No Hit
CGTATCCAACCTGCAGAGTTTATCGCTTCCATGACGCAG	1713	0.43335492096901496	No Hit
ATCCAACCTGCAGAGTTTATCGCTTCCATGACGCAGAAG	1708	0.43209002044079253	No Hit
CAGAGTTTATCGCTTCCATGACGCAGAAGTTAACACTTT	1684	0.42601849790532476	No Hit
TGCAGAGTTTATCGCTTCCATGACGCAGAAGTTAACACT	1668	0.4219708162150128	No Hit
CAACCTGCAGAGTTTATCGCTTCCATGACGCAGAAGTTA	1668	0.4219708162150128	No Hit
TATCCAACCTGCAGAGTTTATCGCTTCCATGACGCAGAA	1630	0.4123575722005221	No Hit

Aparar “trimar” *reads*

- Processo que permite fazer os cortes e ajustes necessários nas *reads* (leituras)
 - Retirada das sequências de adaptadores;
 - Manutenção de sequências com escore mínimo de qualidade e tamanho mínimo;
 - Softwares: **Trimmomatic**, Prinseq, FASTX-Toolkit, ...



- Download do software Trimmomatic
 - <http://www.usadellab.org/cms/?page=trimmomatic>

- Exemplo de uso do software Trimmomatic para biblioteca *paired-end*

```
$ java -jar trimmomatic-0.36.jar PE -threads 8 -phred33 ARQ_R1.fastq.gz
ARQ_R2.fastq.gz ARQ_R1.paired.fastq.gz ARQ_R1.unpaired.fastq.gz ARQ_R2.paired.fastq.gz
ARQ_R2.unpaired.fastq.gz ILLUMINACLIP:/adapters/TruSeq3-PE-2.fa:2:30:10 LEADING:3
TRAILING:3 SLIDINGWINDOW:4:15 MINLEN:36
```



- Resultado do software Trimmomatic
 - Serão gerados arquivos formato .fastq com *reads* “trimadas”

Exemplos de resultados - Trimmomatic

Arquivo de resultado

```
Input Read Pairs: 28987947
Both Surviving: 27213240 (93.88%)
Forward Only Surviving: 784150 (2.71%)
Reverse Only Surviving: 863068 (2.98%)
Dropped: 127489 (0.44%)
```

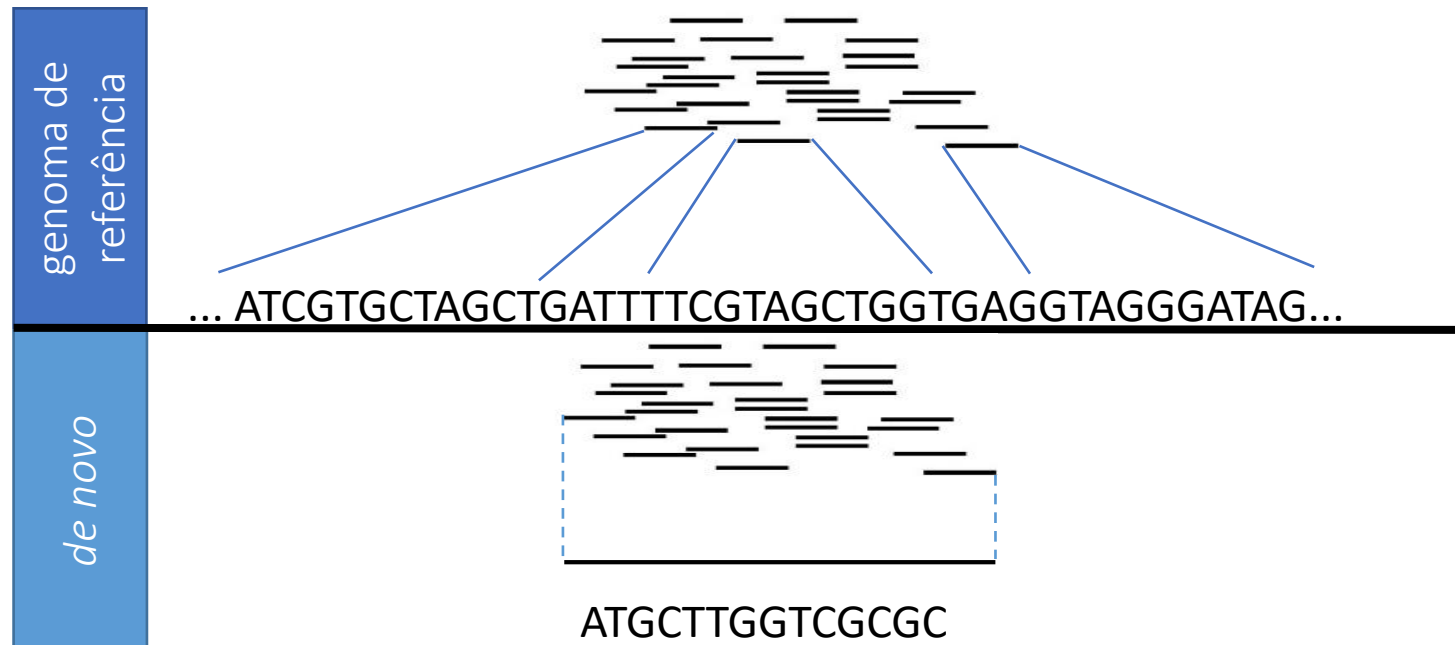
Outros resultados

Sample	Raw reads	High-quality reads
0 hour I (SR)	40,489,531	40,149,007
0 hour II (PE)	28,987,947	27,213,240
0 hour III (PE)	28,893,017	27,136,542
3 hours I (SR)	60,724,079	60,235,207
3 hours II (PE)	67,859,806	63,440,961
3 hours III (PE)	28,182,638	26,438,636
12 hours I (SR)	45,617,478	44,174,487
12 hours II (PE)	30,746,489	28,524,813
12 hours III (PE)	12,463,078	10,635,512

SR=Single-read; PE=Paired-end

Alinhamento das *reads*

- Processo de alinhamento das *reads* sequenciadas (já “trimadas”) no genoma de referência.
 - Quando não houver um genoma de referência, utilizar o método de Alinhamento *de novo*.
 - Softwares: **STAR**, BWA, Bowtie, Bowtie2, Tophat, Tophat2, HISAT, HISAT2, Velvet, ...



Exemplo de uso do software STAR

- Passo 1 - Criando o arquivo de índice do genoma de referência

```
$ STAR --runThreadN 8 --runMode genomeGenerate --genomeDir  
dir_do_genoma --genomeFastaFiles arquivo_supercontigs.fasta --  
sjdbGTFfile arquivo_transcripts.gtf --sjdbOverhang 99
```

GTF

Supercontig21	TR_CBS118892_V2_FINAL_CALLGENES_5	start_codon	14925	14927	.	-	0	gene_id "TERG_00002"; transcript_id "TERG_00002T0";
Supercontig21	TR_CBS118892_V2_FINAL_CALLGENES_5	stop_codon	12841	12843	.	-	0	gene_id "TERG_00002"; transcript_id "TERG_00002T0";
Supercontig21	TR_CBS118892_V2_FINAL_CALLGENES_5	exon	14609	14927	.	-	.	gene_id "TERG_00002"; transcript_id "TERG_00002T0";
Supercontig21	TR_CBS118892_V2_FINAL_CALLGENES_5	CDS	14609	14927	.	-	0	gene_id "TERG_00002"; transcript_id "TERG_00002T0";
Supercontig21	TR_CBS118892_V2_FINAL_CALLGENES_5	exon	13933	14529	.	-	.	gene_id "TERG_00002"; transcript_id "TERG_00002T0";
Supercontig21	TR_CBS118892_V2_FINAL_CALLGENES_5	CDS	13933	14529	.	-	2	gene_id "TERG_00002"; transcript_id "TERG_00002T0";
Supercontig21	TR_CBS118892_V2_FINAL_CALLGENES_5	exon	13702	13869	.	-	.	gene_id "TERG_00002"; transcript_id "TERG_00002T0";
Supercontig21	TR_CBS118892_V2_FINAL_CALLGENES_5	CDS	13702	13869	.	-	2	gene_id "TERG_00002"; transcript_id "TERG_00002T0";
Supercontig21	TR_CBS118892_V2_FINAL_CALLGENES_5	exon	13470	13638	.	-	.	gene_id "TERG_00002"; transcript_id "TERG_00002T0";
Supercontig21	TR_CBS118892_V2_FINAL_CALLGENES_5	CDS	13470	13638	.	-	2	gene_id "TERG_00002"; transcript_id "TERG_00002T0";
Supercontig21	TR_CBS118892_V2_FINAL_CALLGENES_5	exon	12841	13414	.	-	.	gene_id "TERG_00002"; transcript_id "TERG_00002T0";
Supercontig21	TR_CBS118892_V2_FINAL_CALLGENES_5	CDS	12844	13414	.	-	1	gene_id "TERG_00002"; transcript_id "TERG_00002T0";

FASTA

```
>Supercontig21 of Trichophyton rubrum CBS 118892  
ATAGTATTACTTACTACTACTATAAGTCCCTATTATAGAAGGTATGCTCTACTATTAGTA  
TCTAATCTAGATAACTATAAACTATATTTAAATATCTTTAAATACTTAGATATATAGCG  
GTAGTTAAACTACTAGTATACCTCTATTTAGAAGGCGTAGATAGCTTATTTAACTCTCTA  
GATAATACTAAAGTATAGAGATTTAATATTAAGTACTACCTATATTAGCCTTTTTATAAT  
ATTAATAATAATCTTCTATAAAGTAGTAATACCTATAAAAAAGCTATTACTATTCCTATA  
GAGTGTAATAATAATATTATATCTAATTCTAATAGTACTTTAGAGGATATTAATATAAAA  
GATAACTCCACTACTATTTTAAATAGCCCCCTACTATATCTAGACCCCCCTAAGTACCGTAAG  
TATAAGGACGCGTCTCTAATAGAGCTACTATAGCTAAGTAGCTAAAAGTAGATAGGCTA  
ATATCTTATATTATATTAAGGACTAATATAGTACTAAATATAACTAATAAATCTAAAAAA  
GATAGTAAAGGTTTTATATTTAGATTTTAGAATATATTAATATCTAGAGTATTATCTACT
```

Exemplo de uso do software STAR

- Passo 2 – Alinhando as *reads* ao genoma de referência

```
$ STAR --runThreadN 8 --genomeDir dir_do_genoma --readFilesIn  
ARQ_R1_TRIMMED.paired.fastq.gz ARQ_R2_TRIMMED.paired.fastq.gz --  
readFilesCommand zcat --outSAMtype BAM SortedByCoordinate --  
outReadsUnmapped Fastx --outSJfilterReads Unique --twopassMode  
Basic --outFilterMultimapNmax 1 --outFilterType BySJout --  
alignSJoverhangMin 15 --alignSJDBoverhangMin 3 --  
outFilterMismatchNoverLmax 0.06 --quantMode TranscriptomeSAM  
GeneCounts --outFileNamePrefix dir_de_saída
```

- Resultado do software STAR

- Serão gerados arquivos formato .BAM com *reads* alinhadas;
- Caso utilize o parâmetro **GeneCounts**, será gerado o arquivo de contagem de *reads*.

Exemplos de resultados - STAR

[illegible]

SAM/BAM

Number of input reads	46661059
Average input read length	271
UNIQUE READS:	
Uniquely mapped reads number	45099456
Uniquely mapped reads %	96.65%

LOG

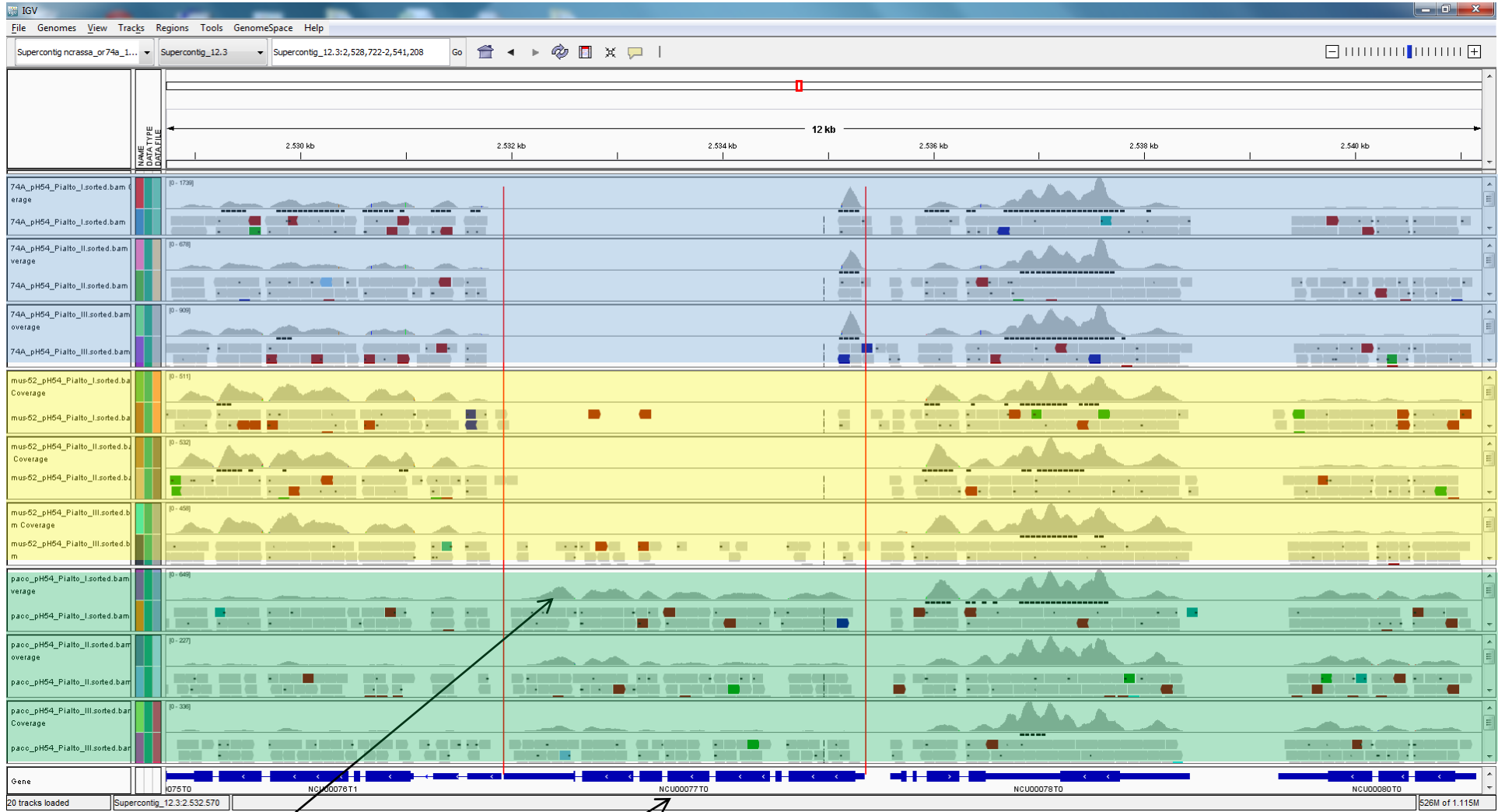
TERG_00002	2441
TERG_00003	7014
TERG_00004	12719
TERG_00008	9056
TERG_00009	1960
TERG_00010	35
TERG_00011	483
TERG_00012	4716
TERG_00013	691
TERG_00014	83
TERG_00015	3444

COUNT

Sample	High-quality reads	Mapped reads STAR	Total mapped reads (%)
Amostra I	46,661,059	45,099,456	96.65
Amostra II	31,721,176	22,577,724	71.18
Amostra III	42,817,030	42,286,629	98.76

Outros resultados

IGV



read counts
+/- 260

Gene X

linhagens
Mutante 1
Mutante 2
Selvagem

Mais sobre seus alinhamentos...

- Conversão BAM → SAM

```
$ samtools view arquivo.bam
```

- Obter uma lista dos nomes das reads sem redundâncias

```
$ samtools view arquivo.bam | cut -f 1 | sort | uniq
```

- Obter a quantidade de reads sem redundância

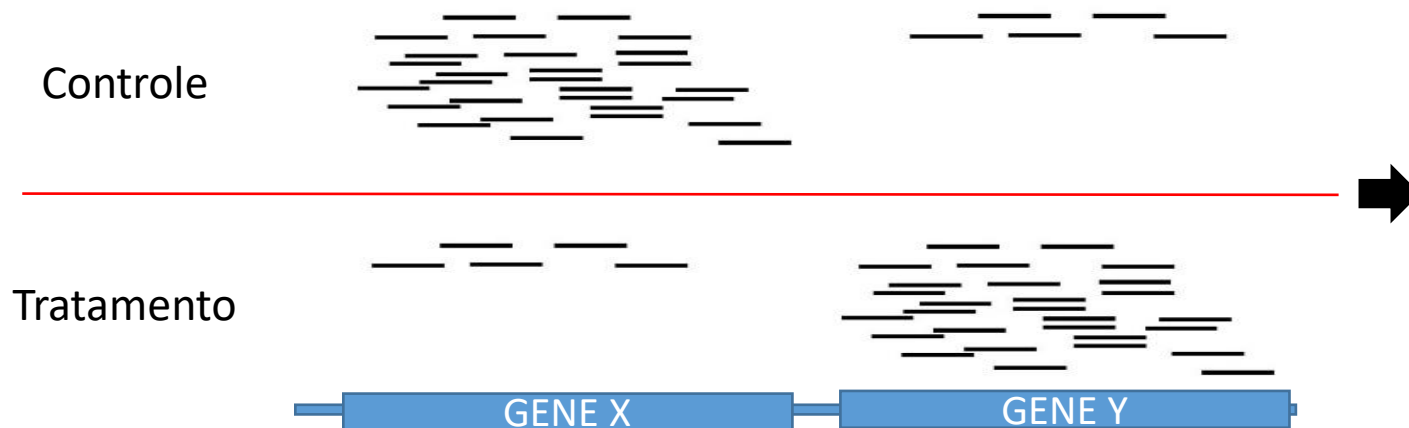
```
$ samtools view arquivo.bam | cut -f 1 | sort | uniq | wc -l
```

- ...

Expressão diferencial

- Ideal que se tenha 3 ou mais sequenciamentos independentes para cada tratamento e para cada grupo controle;
- Podemos usar softwares como: **DESeq2**, EdgeR para normalizar e extrair o diferencial de expressão (estatísticas);
- Ainda outras opções: cuffdiff, limma, baySeq, ...

Representação hipotética



	Controle	Tratamento	FoldChange	\log_2FC
Gene X	28	5	0,18	-2,49
Gene Y	5	28	5,60	2,49

DESeq2

- Baseada em *read counts*
- Pacote do Bioconductor/R
 - <https://www.bioconductor.org/>
- R
 - <https://www.r-project.org/>

The screenshot shows the Bioconductor website for the DESeq2 package. The header includes the Bioconductor logo and navigation links: Home, Install, Help, Developers, and About. A search bar is located in the top right. The main content area for DESeq2 displays various statistics: platforms (all), rank (27 / 1905), posts (283 / 1 / 3 / 42), in Bioc (7.5 years), build (ok), updated (since release), and dependencies (93). It also shows the DOI: 10.18129/B9.bioc.DESeq2 and social media icons for Facebook and Twitter. The description states: "Differential gene expression analysis based on the negative binomial distribution". The Bioconductor version is Release (3.11). The description of the package is: "Estimate variance-mean dependence in count data from high-throughput sequencing assays and test for differential expression based on a model using the negative binomial distribution." The authors listed are Michael Love [aut, cre], Constantin Ahlmann-Eltze [ctb], Simon Anders [aut, ctb], and Wolfgang Huber [aut, ctb]. The maintainer is Michael Love <michaelisaiahlove at gmail.com>. The citation (from within R, enter `citation("DESeq2")`): Love MI, Huber W, Anders S (2014). "Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2." *Genome Biology*, 15, 550. doi: [10.1186/s13059-014-0550-8](https://doi.org/10.1186/s13059-014-0550-8). The installation section provides instructions to install the package in R (version "4.0") and enter the following code:

```
if (!requireNamespace("BiocManager", quietly = TRUE))
  install.packages("BiocManager")

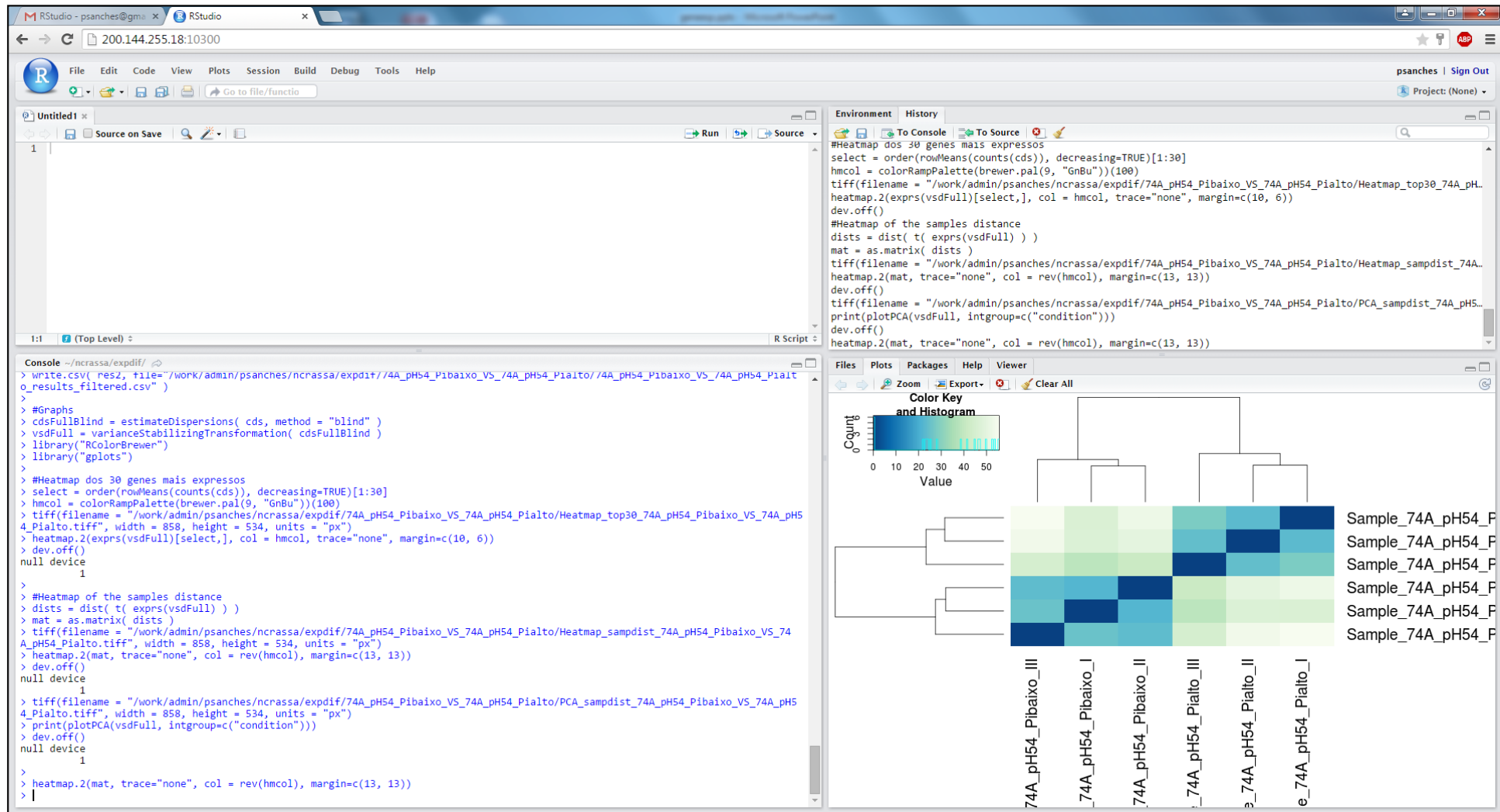
BiocManager::install("DESeq2")
```

For older versions of R, it refers to the appropriate [Bioconductor release](#). The documentation section instructs to view documentation for the version of the package installed in the system, start R and enter:

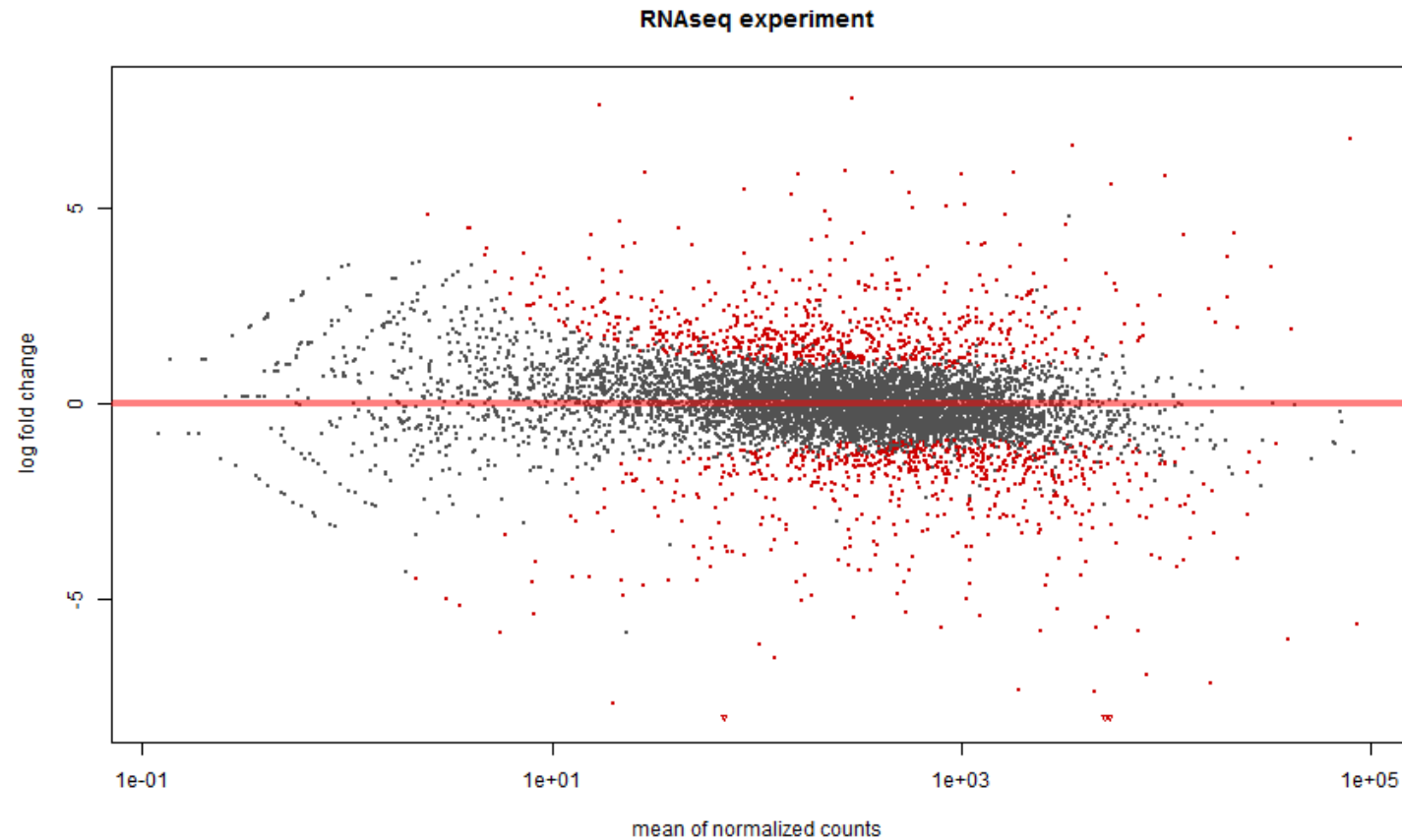
```
browseVignettes("DESeq2")
```

On the right side, there are two boxes: "Documentation" with links to vignettes, workflows, course and conference material, videos, and community resources and tutorials; and "Support" with links to the posting guide, support site, and Bioc-devel mailing list.

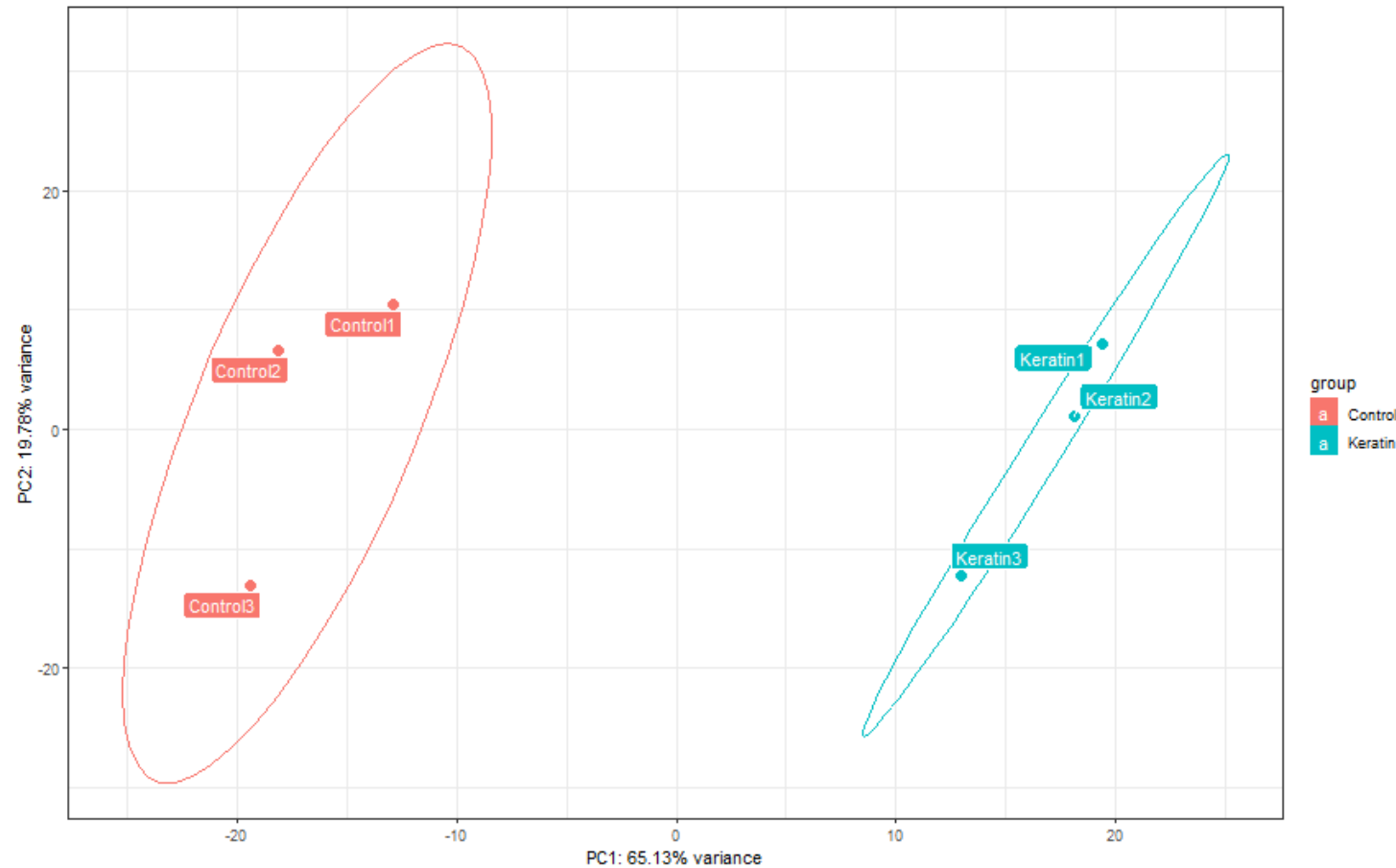
RStudio - <https://rstudio.com/>



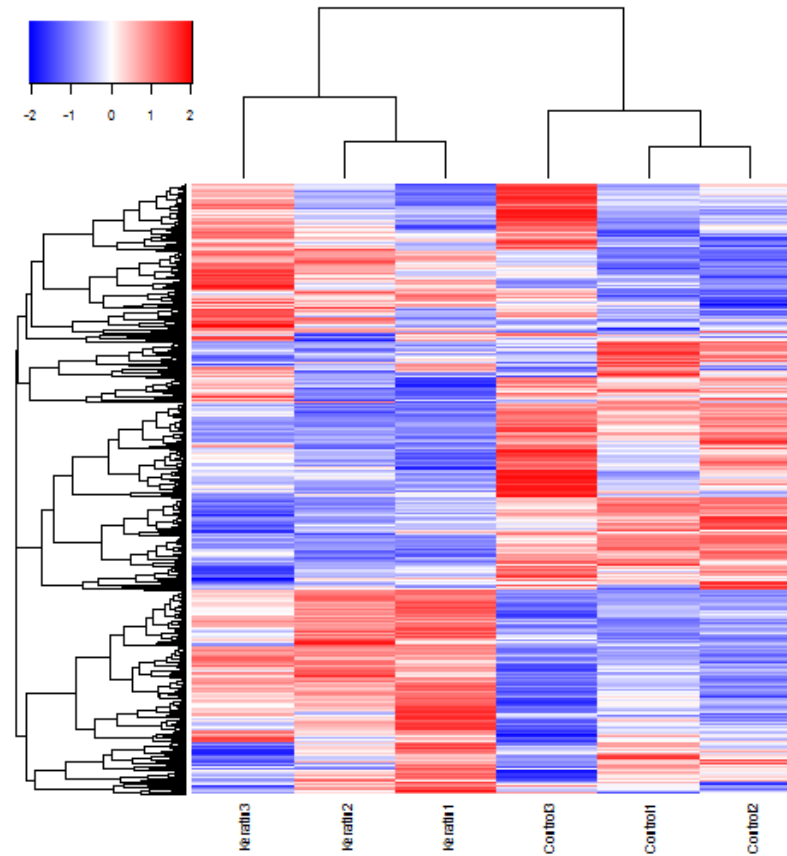
Exemplos de resultados – DESeq2



Exemplos de resultados – DESeq2



Exemplos de resultados – DESeq2



Exemplos de resultados – DESeq2

	gene	baseMean	log2FoldChange	lfcSE	stat	pvalue	padj	Control1	Control2	Control3	Keratin1	Keratin2	Keratin3
1	TERG_08353	5263,696589	-9,002452266	0,457685141	-19,66953142	3,93E-86	3,35E-82	12972,30597	8400,824323	10148,10033	25,66658515	23,4192922	11,86303686
2	TERG_00867	8034,057675	-6,924413959	0,449552077	-15,40291839	1,56E-53	6,65E-50	21224,51933	16861,695	9725,389828	177,182233	117,096461	98,46320595
3	TERG_03274	39068,16881	-6,010855869	0,410694252	-14,63584122	1,66E-48	4,70E-45	75212,20447	101602,9569	54013,45959	911,57775	1348,458193	1320,356003
4	TERG_06242	9819,287482	5,837124878	0,403135796	14,47930184	1,64E-47	3,48E-44	323,0069501	284,3704525	405,4362701	19229,24	24895,9402	13777,73101
5	TERG_01435	1004,302718	5,852167669	0,406323546	14,40272839	4,97E-47	8,46E-44	36,85314196	33,9885003	31,5			
6	TERG_02842	4997,657133	-8,213330898	0,751223985	-10,93326499	7,99E-28	1,13E-24	10492,3063	699,0301561	1869			
7	TERG_05652	1815,001137	5,907089754	0,547900611	10,78131623	4,22E-27	5,12E-24	96,10721335	40,78620035	40,6			
8	TERG_02974	5210,880922	-5,450595199	0,510534625	-10,67624982	1,31E-26	1,40E-23	5386,339611	18580,38016	6599			
9	TERG_04580	85237,93713	-5,637341532	0,533843746	-10,5599093	4,57E-26	4,32E-23	238209,3157	163247,8999	9989			
10	TERG_00856	1096,023656	-4,605000909	0,44760248	-10,28814877	7,97E-25	6,78E-22	2089,067322	3030,641276	1197			
11	TERG_06585	4478,365775	-7,364525546	0,723264717	-10,18233764	2,38E-24	1,84E-21	1282,633862	20739,78288	4685			
12	TERG_04942	3811,654367	-4,396265658	0,437335943	-10,05237672	8,97E-24	6,36E-21	4557,505223	9213,14948	8061			
13	TERG_03223	3494,277166	6,577117485	0,657796482	9,998711852	1,54E-23	1,01E-20	21,6783188	49,84980043	146,			
14	TERG_01841	160,0771112	5,859269453	0,586463982	9,990842806	1,67E-23	1,02E-20	3,613053134	10,19655009	3,04			
15	TERG_03226	5420,83863	5,624314883	0,564531946	9,962792936	2,22E-23	1,26E-20	76,59672643	140,4858012	429,			
16	TERG_02844	16356,73185	-7,115805721	0,721553954	-9,861779122	6,10E-23	3,24E-20	30338,80716	3611,844631	6348			
17	TERG_06548	804,7925492	-5,693311316	0,586368898	-9,709436045	2,75E-22	1,38E-19	1157,622224	2977,392626	602,			
18	TERG_11610	293,7058974	4,113211732	0,427687852	9,617321872	6,76E-22	3,19E-19	26,01398256	35,12145031	35,			
19	TERG_12038	4566,96921	-5,693308736	0,592763145	-9,604694199	7,64E-22	3,42E-19	10717,76082	13017,59561	3147			
20	TERG_05627	465,3771355	5,890375985	0,618818924	9,518739255	1,75E-21	7,45E-19	24,56876131	5,664750049	15,2			
21	TERG_06509	7364,750494	-5,777097642	0,607896269	-9,503426708	2,03E-21	8,23E-19	14262,16594	21570,23524	7566			
22	TERG_04228	12204,28639	-3,970785992	0,42008175	-9,452412526	3,31E-21	1,28E-18	18508,94859	28842,64135	2148			
23	TERG_06807	457,7852227	3,709139395	0,402693409	9,210827178	3,24E-20	1,20E-17	50,58274387	66,84405058	78,2			
24	TERG_02169	561,5525446	5,384599423	0,588892347	9,143605702	6,04E-20	2,14E-17	33,24008883	24,92490022	20,3			

Exemplo de aplicação de filtros

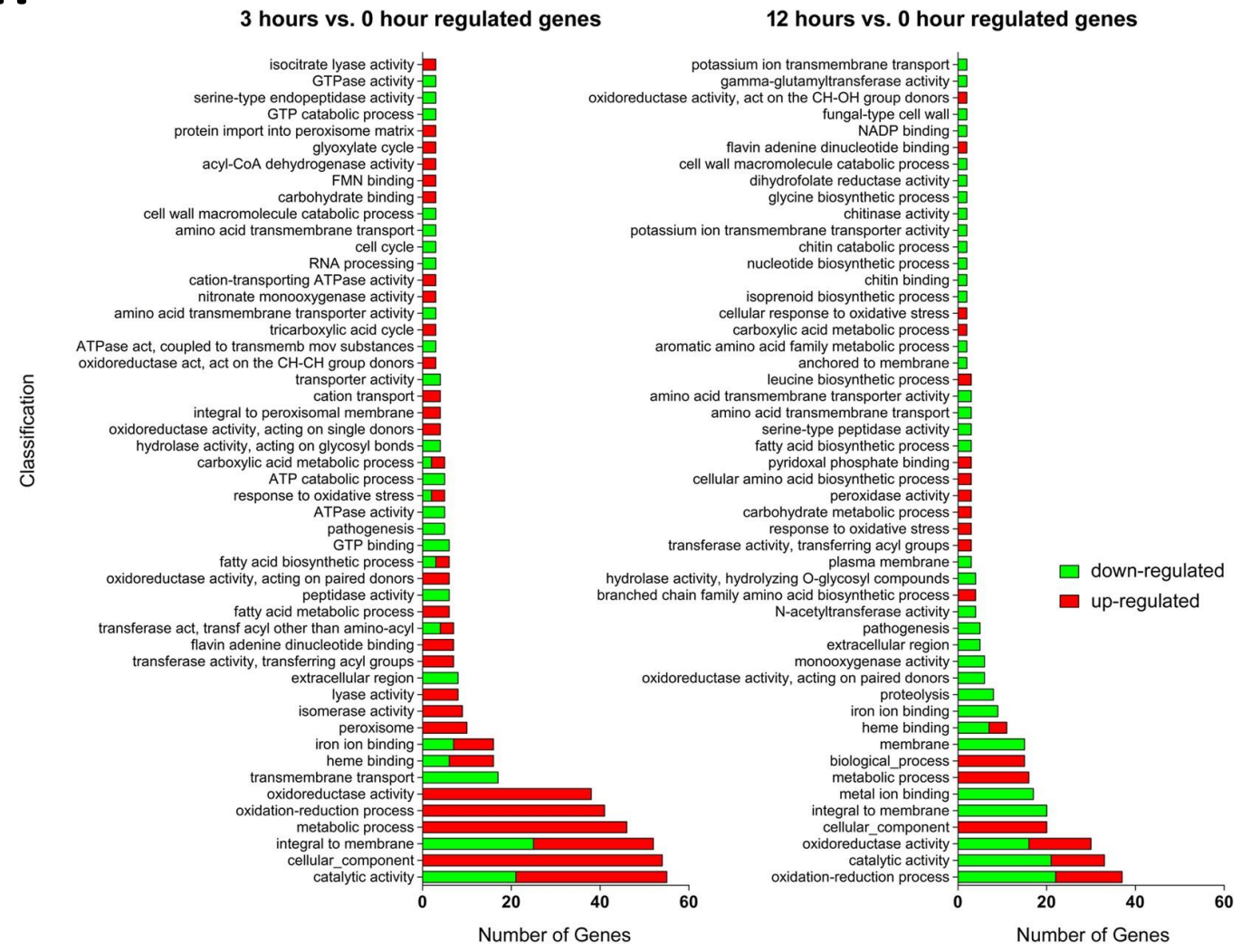
- $pvalue < 0,05$
- $log2FoldChange \geq 1,5$ (*up-regulated*)
- $log2FoldChange \leq -1,5$ (*down-regulated*)

ID	Gene Product Name	log ₂ (Fold change)
24 hours		
Up-regulated		
TERG_01599	hypothetical protein	7.81
TERG_05652	leucine aminopeptidase 1	6.23
TERG_03223	N-acetylglucosamine-6-phosphate deacetylase	6.07
TERG_05627	LysM domain-containing protein (<i>M. canis</i>)	6.04
TERG_06242	glucanase, putative (<i>T. verrucosum</i>)	5.67
TERG_12035	NB-ARC and TPR domain protein (<i>A. benhamiae</i>)	5.64
TERG_07909	isochorismatase family hydrolase, putative (<i>A. benhamiae</i>)	5.11
TERG_01841	hypothetical protein	5.06
TERG_03226	glucosamine-6-phosphate deaminase	4.92
TERG_01435	flavin containing polyamine oxidase, putative (<i>A. benhamiae</i>)	4.78
Down-regulated		
TERG_02842	6-hydroxy-D-nicotine oxidase (<i>T. equinum</i>)	-9.33
TERG_08353	cytochrome P450 55A3 (<i>T. tonsurans</i>)	-8.87
TERG_02746	hypothetical protein	-8.29
TERG_06049	dimethylallyl tryptophan synthase, putative (<i>T. verrucosum</i>)	-8.11
TERG_06054	hypothetical protein	-8.05
TERG_02844	major facilitator superfamily transporter (<i>T. tonsurans</i>)	-7.69
TERG_02412	HHE domain protein (<i>T. verrucosum</i>)	-7.52
TERG_02024	hypothetical protein	-6.95
TERG_11792	hypothetical protein	-6.73
TERG_00867	DUF1212 domain membrane protein (<i>A. benhamiae</i>)	-6.67

Enriquecimento funcional

- Identificar classes de genes que estão super-representadas em um grande conjunto de genes;
- Uso de Bancos de Dados e Ferramentas como:
 - Gene Ontology
 - Blast2GO
 - BayGO
 - FunRich
 - ...

Exemplos de resultados – Enriquecimento funcional

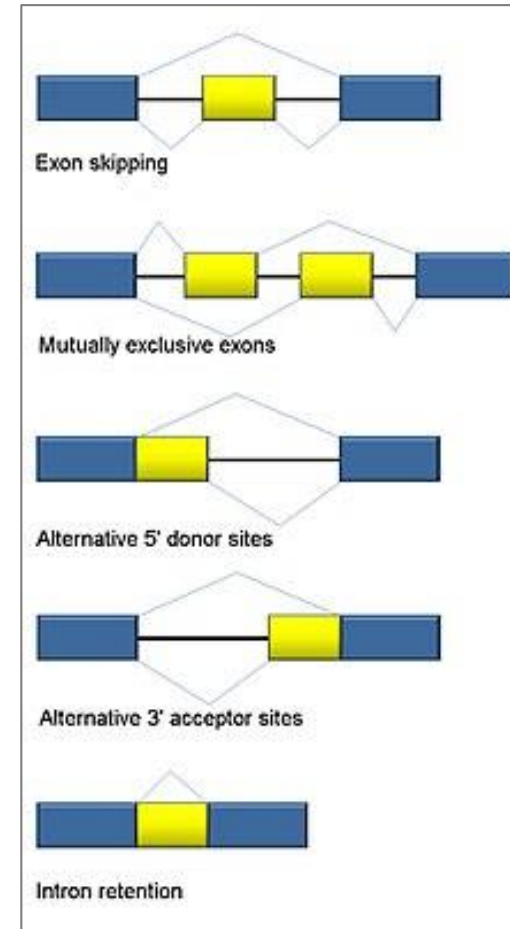


Após enriquecimento funcional → Seleção de genes específicos

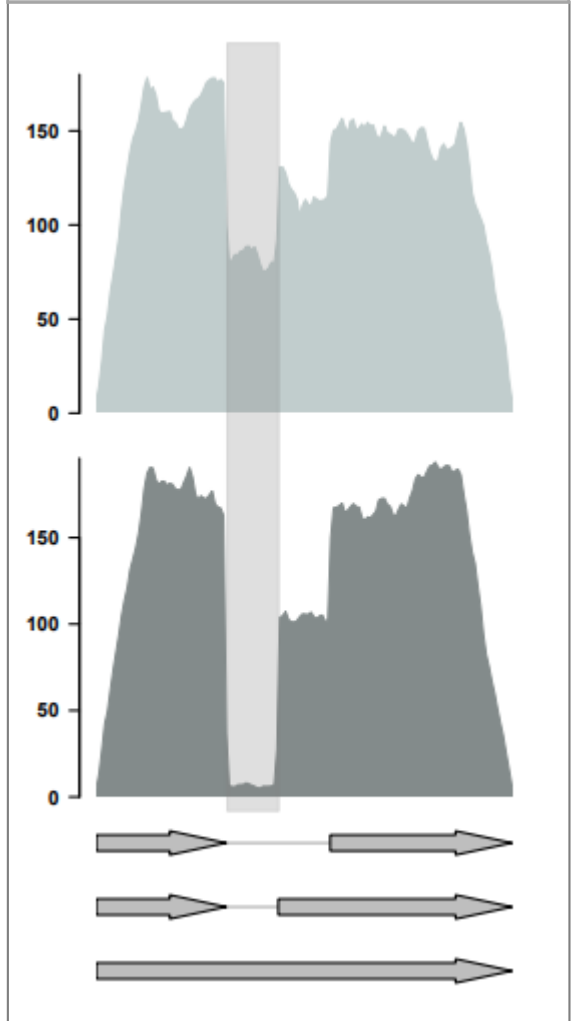
			-3.0	0.0	3.0
	ID	Gene Product Name	3h vs. 0h	12h vs. 0h	
Fatty acid metabolism	TERG_00823	rhannolipids biosynthesis 3-oxoacyl-[acyl-carrier-protein] reductase (<i>T. equinum</i>)	4.22	2.08	
	TERG_00919	Phospholipid:diacylglycerol acyltransferase (<i>T. equinum</i>)	1.61	0.47	
	TERG_01184	fatty acid-binding protein (<i>T. equinum</i>)	1.66	0.12	
	TERG_02250	cyclopropane-fatty-acyl-phospholipid synthase (<i>T. tonsurans</i>)	-1.52	0.84	
	TERG_02909	acyl-CoA oxidase, putative (<i>T. verrucosum</i>)	3.81	1.24	
	TERG_03343	fatty acid desaturase (<i>T. equinum</i>)	-3.38	-1.83	
	TERG_03483	carnitine acetyl transferase (<i>T. tonsurans</i>)	3.54	1.31	
	TERG_03963	mannosyl phosphorylinositol ceramide synthase SUR1 (<i>T. equinum</i>)	-1.36	-1.73	
	TERG_04038	acyl-CoA dehydrogenase (<i>T. tonsurans</i>)	3.67	2.08	
	TERG_04753	3-isopropylmalate dehydratase, large subunit	2.49	3.60	
	TERG_05484	acyl-CoA dehydrogenase (<i>T. tonsurans</i>)	4.02	0.93	
	TERG_07060	acyl-CoA thioesterase II	2.62	0.92	
	TERG_07659	peroxisomal 3-ketoacyl-coA thiolase (Kat1), putative (<i>A. benhamiae</i>)	3.01	1.32	
	TERG_07691	sterol carrier protein (<i>T. tonsurans</i>)	2.73	0.90	
	TERG_08051	3-oxoacyl-(acyl-carrier-protein) reductase (<i>T. tonsurans</i>)	2.46	1.30	
	TERG_08170	fatty acid elongase glg30 (<i>T. equinum</i>)	2.27	1.89	
	TERG_08519	mitochondrial enoyl reductase (<i>T. equinum</i>)	1.98	1.30	
	TERG_08952	long-chain-fatty-acid-CoA ligase (<i>T. equinum</i>)	1.65	0.18	
	TERG_11538	3-oxoacyl-(acyl-carrier-protein) reductase (<i>T. tonsurans</i>)	2.77	0.56	
	TERG_11539	3-oxoacyl-(acyl-carrier-protein) reductase (<i>T. tonsurans</i>)	2.86	0.79	
	TERG_12530	3-ketoacyl-CoA thiolase peroxisomal A (<i>T. tonsurans</i>)	4.19	1.50	
Oxidative Stress	TERG_01252	catalase A	-0.15	2.37	
	TERG_01349	glutathione peroxidase (<i>T. tonsurans</i>)	2.41	1.60	
	TERG_01463	cytochrome c peroxidase (<i>T. tonsurans</i>)	3.17	2.15	
	TERG_02041	glutathione S-transferase (<i>T. equinum</i>)	3.21	-0.07	
	TERG_02747	cytochrome P450 alkane hydroxylase (<i>A. benhamiae</i>)	3.34	0.30	
	TERG_02842	6-hydroxy-D-nicotine oxidase (<i>T. equinum</i>)	2.03	-0.85	
	TERG_03078	cytochrome P450 oxidoreductase OrdA-like, putative (<i>T. verrucosum</i>)	5.63	0.78	
	TERG_03231	cytochrome P450 52A12 (<i>T. tonsurans</i>)	2.34	-0.22	
	TERG_03390	glutathione S-transferase (<i>T. equinum</i>)	3.00	0.95	
	TERG_04960	glutathione S-transferase Ure2-like, putative (<i>A. benhamiae</i>)	6.15	0.81	
	TERG_07282	cytochrome P450 monooxygenase, putative (<i>T. verrucosum</i>)	-2.81	-1.77	
	TERG_08069	glutathione-disulfide reductase	2.05	0.83	
	TERG_08347	cytochrome P450 monooxygenase, putative (<i>A. benhamiae</i>)	1.76	-1.55	
	TERG_08353	cytochrome P450 55A3 (<i>T. tonsurans</i>)	1.67	1.59	
	TERG_08969	cytosolic Cu/Zn superoxide dismutase, putative (<i>A. benhamiae</i>)	-0.87	2.46	
Proteases	TERG_02214	carboxypeptidase 2	-2.37	-2.02	
	TERG_03681	leucine aminopeptidase (<i>T. equinum</i>)	-0.30	-1.90	
	TERG_04324	extracellular metalloproteinase 4	-1.64	-1.39	
	TERG_05735	dipeptidyl peptidase 4	-2.20	-1.53	
	TERG_06361	ATP-dependent protease La	2.13	0.53	
	TERG_06552	aspartic-type endopeptidase (OpsB), putative (<i>T. verrucosum</i>)	-0.99	-1.94	
	TERG_06625	serine protease, putative (<i>A. benhamiae</i>)	-2.42	-2.06	
	TERG_08195	peptidase S41 family protein (<i>M. gypseum</i>)	-1.61	-1.70	
	TERG_08557	carboxypeptidase S1, putative (<i>A. benhamiae</i>)	-2.37	-1.50	
TERG_12606	protease DPPV, putative (<i>A. benhamiae</i>)	-1.32	-2.33		

Análise de *Splicing* Alternativo

- Identificar o uso de isoformas
- Analisar a expressão diferencial de isoformas entre condições distintas
- Uso de softwares como:
 - **ASpli** (pacote Bioconductor/R)
 - DEXSeq (pacote Bioconductor/R)
 - *Scripts* desenvolvidos em linguagem Perl



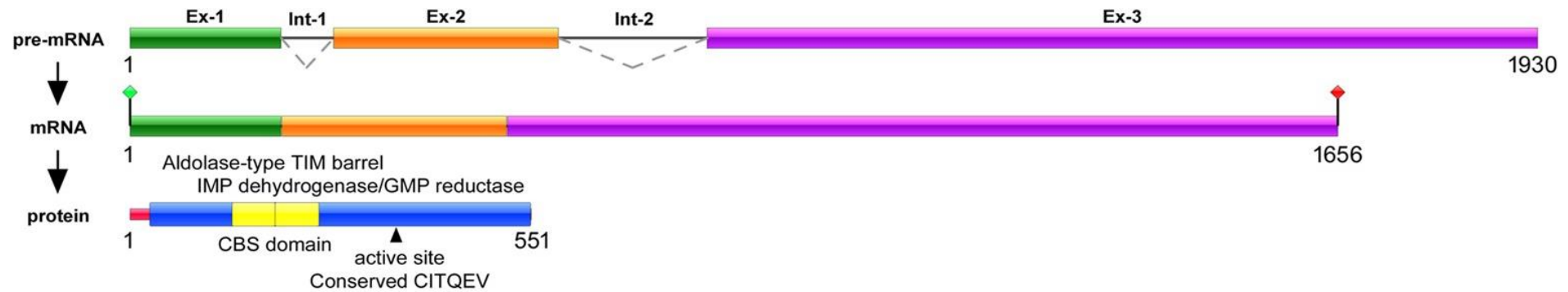
Exemplo de saída do software ASpli para retenção de intron em um gene



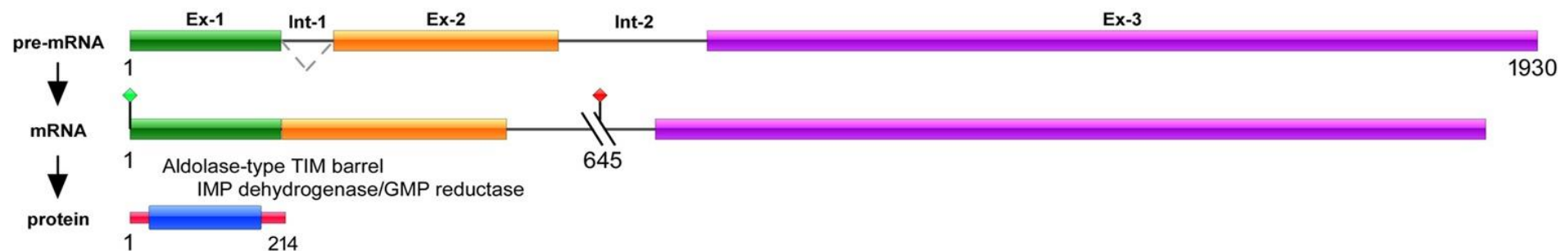
Análise de *Splicing* Alternativos

- Uso de complementos como: pfam, InterPro, UniProt, ExPASy, etc.

Conventional splicing



Retention of intron 2



Habilidades essenciais - Bioinformata

- Conhecimento na área de Biologia Molecular, Computação e Estatística;
- Conhecimento no uso de ferramentas e pacotes de Bioinformática;
- Desejável conhecimento em linguagens de programação;
- Não ter “medo” da interface de linha de comandos (ex. Linux).

Considerações finais

Resultados ruins também são resultados.

O software auxilia o processo de tomada de decisão, mas quem toma a decisão final é você.

“Garbage in, garbage out (GIGO)” -> “lixo entra, lixo sai”.
George Fuechsel (Técnico da IBM)

“se devidamente torturados, os dados contam qualquer coisa”.
Darrel Huff - Como Mentir com Estatísticas.

