



**QUEEN'S
UNIVERSITY
BELFAST**

Molecular Profiling of RNA Tumors Using High-Throughput RNA Sequencing: From Raw Data to Systems Level Analyses

da Silveira, W. A., Hazard, E. S., Chung, D., & Hardiman, G. (2019). Molecular Profiling of RNA Tumors Using High-Throughput RNA Sequencing: From Raw Data to Systems Level Analyses. *Methods in Molecular Biology*, 1908, 185-204. Advance online publication. https://doi.org/10.1007/978-1-4939-9004-7_13

Published in:
Methods in Molecular Biology

Document Version:
Peer reviewed version

Queen's University Belfast - Research Portal:
[Link to publication record in Queen's University Belfast Research Portal](#)

Publisher rights
© 2019 Springer Nature Switzerland AG.
This work is made available online in accordance with the publisher's policies. Please refer to any applicable terms of use of the publisher.

General rights
Copyright for the publications made accessible via the Queen's University Belfast Research Portal is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy
The Research Portal is Queen's institutional repository that provides access to Queen's research output. Every effort has been made to ensure that content in the Research Portal does not infringe any person's rights, or applicable UK laws. If you discover content in the Research Portal that you believe breaches copyright or violates any law, please contact openaccess@qub.ac.uk.

Open Access
This research has been made openly available by Queen's academics and its Open Research team. We would love to hear how access to this research benefits you. – Share your feedback with us: <http://go.qub.ac.uk/oa-feedback>

1 Molecular Profiling of RNA Tumors Using High-Throughput RNA Sequencing: From Raw Data to 2 Systems Level Analyses

3 Willian A. da Silveira, E. Starr Hazard, Dongjun Chung, and Gary Hardiman

4 Abstract

5 RNAseq is a powerful technique enabling global profiles of transcriptomes in healthy and diseased
6 states. In this chapter we review pipelines to analyze the data generated by sequencing RNA, from
7 raw data to a system level analysis. We first give an overview of workflow to generate mapped reads
8 from FASTQ files, including quality control of FASTQ, filtering and trimming of reads, and alignment
9 of reads to a genome.

10 Then, we compare and contrast three popular options to determine differentially expressed (DE)
11 transcripts (The Tuxedo Pipeline, DESeq2, and Limma/voom). Finally, we examine four tool sets to
12 extrapolate biological meaning from the list of DE genes (Genecards, The Human Protein Atlas,
13 GSEA, and ToppGene). We emphasize the need to ask a concise scientific question and to clearly
14 understand the strengths and limitations of the methods.

15 **Key words:** High-throughput sequencing (HTS), RNAsequencing (RNAseq), FASTQ, Tuxedo pipeline,
16 HTSeq, DESeq 2, Limma/Voom, Gene Set Enrichment Analysis/GSEA, TOPPGENE

17

18 1 Introduction

19

20 Classical techniques in genetics and molecular biology remain the gold standard when one
21 wants to detect the presence and sequence of a gene (i.e., using polymerase chain reaction (PCR)
22 and Sanger sequencing techniques), its mRNA expression level (i.e., using quantitative PCR (qPCR)),
23 and the corresponding protein levels (i.e., using western blots) [1]. The human genome encodes
24 approximately 25,000 genes, with thousands of them expressed in multiple combinations in diverse
25 cellular contexts. Furthermore multiple isoforms exist for the same gene with many possible
26 downstream post-translation modifications [2]. The emergence and application of high throughput
27 approaches in the past decade, the so-called “omics” fields have ignited a revolution in biological
28 research. Classical genomics techniques only allow investigation of a small number of genes and
29 proteins at the same time, where as “omics” approaches (including high throughput RNA
30 sequencing) enable investigation of the entire mRNA content at the same time [1, 2].

31 Transcriptomics is the study of the transcriptome, i.e., the complete set of RNA transcripts
32 that are produced by the genome, under specific circumstances or in a specific cell, using
33 highthroughput methods, such as RNAseq [3]. Comparison of transcriptomes allows the
34 identification of genes that are differentially expressed in distinct cell populations, for example in
35 healthy or tumor tissues, or in response to therapeutic regimes. In this chapter we discuss the
36 methods to analyze and interpret data from human tumor samples generated by RNAseq
37 technology. We cover the key steps and the progression from the “FASTQ” files generated by the
38 sequencing instrument through a list of differentially expressed genes, and system level analyses.

39 Although the methods and workflows described in this chapter are best suited to a Linux
40 environment, there are options to run many of these programs in Windows or Mac OS X as well, e.g.,
41 using Cygwin (<https://www.cygwin.com/>) in Windows. A basic knowledge of the command line
42 environment, file structures and rudimentary coding skills are assumed. More details about these

43 programs and the statistical models underlying them used can be found in the references and/or in
44 the links of the websites in the Materials section. As noted by Mayer-Scho“nberger and Cukier, “the
45 data can reveal secrets to those with the humility, the willingness, and the tools to listen” [4].

46 **2 Materials**

47

48 2.1 Computing Infrastructure

49

50 There are many challenges associated with selecting and implementing the right set of tools.
51 Bioinformatics analyses are complex, multistep processes composed of multiple software
52 applications. Ideally, many of the programs used for RNAseq analysis are designed to show the
53 optimal performance in a high-performance computing environment. However, these programs can
54 still run in a sufficiently powerful and well-configured laptop or desktop machine, e.g., an Intel Core
55 i7 processor with storage RAID 0 configured 2 parallel 1-TB hard disk drives. We recommend at least
56 4GB of RAM, preferably 8GB.

57

58 2.2 Software Tools and Genome Build(s)

59

60 SRA Toolkit 2.9.0 (<https://github.com/ncbi/sra-tools/wiki/Downloads>,
61 <http://www.ncbi.nlm.nih.gov/books/NBK158900/>), Windows/Linux/MAC OS X.

62 FastQC v0.11.7 (<http://www.bioinformatics.babraham.ac.uk/projects/FASTQC/>),
63 Windows/Linux/MAC OS X.

64 FASTX-Toolkit (http://hannonlab.cshl.edu/fastx_toolkit/), Linux/MAC OS X/Web based.

65 Cutadapt (<https://cutadapt.readthedocs.io>), Python Language. Windows/Linux/MAC OS X.

66 Bowtie1 for reads between 35 and 50 bp (<http://bowtie-bio.sourceforge.net/manual.shtml>).
67 Linux/MAC OS X.

68 Bowtie2 for reads greater than 50 bp (<http://bowtie-bio.sourceforge.net/bowtie2/manual.shtml>).
69 Linux/MAC OS X.

70 SAMtools (<http://samtools.sourceforge.net/>) . Windows/Linux/MAC OS X. Samtools also have an “R”
71 compatible version “Rsamtools,” available from the Bioconductor website.

72 Human Reference Genome Sequence, hg38 (<http://hgdownload.cse.ucsc.edu/downloads.html>).

73 TopHat (<https://ccb.jhu.edu/software/tophat/manual.shtml>)). Linux/MAC OS X.

74 Cufflinks package (Cufflinks, CuffMerge, CuffMerge, Cuffdiff)([http://cole-trapnell-
75 lab.github.io/cufflinks/tools/](http://cole-trapnell-lab.github.io/cufflinks/tools/)). Linux/MAC OS X.

76 CummeRBound (<https://bioconductor.org/packages/release/bioc/html/cummeRbund.html>). R
77 Language. Windows/ Linux/MAC OS X.

78 HTSeq (<http://www.huber.embl.de/HTSeq/>). Python Language. Windows/Linux/MAC OS X.

79 Comprehensive gene annotation (GRCh38.p7, “.GTF file”).
80 (<http://www.genecodegenes.org/releases/25.html>).
81 DESeq2 (<https://bioconductor.org/packages/release/bioc/html/DESeq2.html>). R Language.
82 Windows/Linux/MAC OS X.
83 Limma/Voom (<https://bioconductor.org/packages/release/bioc/html/limma.html>). R Language.
84 Windows/Linux/MACOS X.
85 Venny (<https://www.stefanijol.nl/venny>) Web tool.
86 ToppFun (<https://toppgene.cchmc.org/>) Web tool.
87 GSEA tool and Website (<http://software.broadinstitute.org/gsea/index.jsp>). Web tool and a program
88 compatible with Windows/Linux/MAC OS X.

89

90 **3 Methods**

91 A schematic of data analysis pipelines is presented in Fig. 1. To initiate the analyses, it is necessary to
92 have the FASTQ files that contain the information of sequenced reads and the quality score for each
93 nucleotide. These FASTQ files can be obtained directly from the sequencing machine, as in the case
94 depicted in the flowchart, or from the Sequence Read Archive (SRA) repository
95 (<https://www.ncbi.nlm.nih.gov/sra>), where it is required to convert the downloaded “.SRA” file to
96 FASTQ using the SRA Tool kit. Then the quality of the FASTQ file can be checked using the FastQC
97 program and the results from it can be used as parameters in the “FASTX” Tool Kit and/or in
98 Cutadapt to perform trimming and filtering of the reads. Then, the FASTQ file is now ready for the
99 alignment of reads to the genome, which can be attained using the TopHat program with a
100 reference genome, or any other aligner tool. At this point our pipeline bifurcates, the list of DE genes
101 can be obtained using the Tuxedo pipeline (comprised by Cufflinks, Cuffmerge and Cuffdiff) using
102 FPKM (Fragments Per Kilobase of transcript per Million mapped reads). Alternatively, the list of DE
103 genes can also be obtained using the counts of the reads coming from the HTseq program and then
104 processed by DESEQ2 or Limma/Voom. The biological meaning of the DE gene list can be analyzed
105 gene by gene using “The Human Protein Atlas” and “Genecards” websites or at a system level using
106 the GSEA and/or ToppGene Suite. We discuss each step more in detail below.

107

108 **3.1 From FASTQ Files to BAM/SAM Files**

109

110 **3.1.1 Raw Data—FASTQ Files**

111

112 A typical RNAseq data analysis begins with FASTQ files. There is a single FASTQ that
113 corresponds to each RNA sample sequenced. If the study is paired-end sequencing, there are two
114 FASTQ files, each of which corresponds to the left and right reads of each DNA fragment [5]. In the
115 early days of high throughput sequencing, sequenced reads were often stored in the FASTA file
116 format, which is a text file containing a sequence of nucleic acids or amino acids. The FASTQ file is an
117 extension of this file format, i.e., a FASTA that contains both a sequence of nucleic acids and a
118 quality score for that particular sequence. The quality score provides a measure of confidence in the

119 sequencing data [6]. These files are typically large and easily reaching tens to hundreds of gigabytes
120 in size.

121 In the FASTQ file, each raw sequence is described with four lines (Fig. 2a). The first line
122 begins with an “@” and contains an identifier for the sequence. The second line contains the raw
123 sequence itself, as found in the FASTA file. The third line starts with a “+” and optionally repeats the
124 content of the first line. The fourth line contains the phred quality score (Q), a measure of the quality
125 of the identification of the base (Fig. 2b) [6]. Phred was originally developed as a quality score for
126 Sanger sequencing data and was adapted in FASTQ. This quality score is calculated by comparing
127 chemical parameters of the given sequencing process with the parameters of a large dataset of
128 known accuracy [6, 7].

129 In this review, we are using the dataset GSE81167, available from the Gene Expression
130 Omnibus (GEO), for the purpose of illustration. This study evaluated the impact of ZEB1 expression
131 in HCC827 cells, which are human lung cancer cell lines [8]. ZEB1 is one of the principal transcription
132 factors involved in the epithelial-to-mesenchymal transition, a key event in tissue invasion and
133 metastasis [9]. This dataset is available for GEO in the short read archive “.SRA” file format. We used
134 the FASTQ-dump command line program from the SRA Toolkit 2.9.0 to convert files from “.SRA” to
135 “.FASTQ”. In most cases, SRA files downloaded from repositories, as from the Gene Expression
136 Omnibus, can be directly converted to SAM files (Sequence Alignment/Map Format), affording
137 savings to both the user’s time and central processing unit (CPU) cycles. For didactic reasons, we
138 start from the FASTQ files.

139

140 3.1.2 Quality Control of FASTQ Files with FastQC

141

142 Although the FASTQ file provides the quality score for each base in every sequence, checking
143 it manually is time-consuming and impractical. FastQC is a program that allows evaluation of the
144 quality of the FASTQ file as a whole. The program uses the FASTQ file as input, and yields three types
145 of FastQC files: A report in the “.html” format, the same report in the compressed “.zip” file, and the
146 “FastQC” folder with the unzipped version of the files [10]. As presented in Fig. 3, the report analyzes
147 a number of items, basic statistics, per-base sequence quality, and other metrics including sequence
148 content, GC content, sequence duplication and the presence of adapter and overrepresented
149 sequences. In our example, even in what can be considered a good report (Fig. Upper Panel), per
150 base sequence content and K-mer content are flagged as potentially problematic. The per-base
151 sequence content suggests sequences with GC-content are over-represented, which might
152 potentially imply contamination from ribosomal RNA (rRNA). K-mer bias occurs when over-
153 represented sequences result in the K-mers derived from these sequences being highly enriched
154 [10]. In our bad report (Fig. 3 Lower Panel), in addition to problems with the quality of the
155 nucleotides sequenced that occur as the read length grows, we have a warning for per tile sequence
156 quality, sequence duplication level and over-represented sequences and a failure for per base
157 sequence content, per sequence GC content and again the K-mer content. Both the warnings and
158 the failures indicate a problem at the library preparation step [10].

159

160 3.1.3 Preprocessing—Filtering and Trimming

161

162 Once we have the information on sample quality, the next step is to filter out bases and
163 reads with low quality and to extract adaptor, primer, and poly-A tails from the data. To do that, we
164 use the “fastx_trimer” and “fastq_quality_trimmer” programs to filter the adapters, poly-A tails, and
165 PCR primer sequences and “fastq_quality_filter” to filter out low quality bases and/or reads. These
166 programs are a part of the FASTX-Toolkit [11].

167 By taking the FASTQ files as input, a sequential use of the two FASTX-Toolkit programs,
168 “fastx_trimer” and “fastq_quality_trimmer”, will return a “cleaned” FASTQ file. FASTX-Toolkit can be
169 used only in the command line, which requires some basic knowledge of Linux. The information
170 acquired with the FastQC program must be used in the parameters of the “fastq_quality_filter”,
171 when necessary.

172 Alternatively, we can use the program “Cutadapt” [12] for trimming adapter sequences,
173 primers, and poly-A tails from the FASTQ file. Cutadapt accepts the FASTQ as input and operates in
174 any system that runs Python. However the program limits itself to the file cleaning and does not
175 support base quality filtering [12].

176

177 3.1.4 Preprocessing — Read Alignment to the Genome

178

179 We now have good quality reads and a high level of confidence that the original RNA
180 sequence is represented in the FASTQ file. Is it necessary now to align the reads to our reference
181 genome, in this case the Human Genome, version hg38 [13]. We use, TopHat, a part of the Tuxedo
182 pipeline, described by Trapnell and collaborators in 2012 [14]. TopHat is integrated with the Bowtie
183 and SAMtools programs, uses the FASTQ files as input, uses the selected genome as the reference,
184 and generates a SAM file as output. It runs only in the linux shell [14, 15]. To generate the alignment,
185 firstly TopHat uses the Bowtie program to align the reads to the genome. Bowtie is a fast and
186 efficient short read aligner, but unfortunately cannot align reads with large gaps compared with the
187 reference, which makes it unsuitable to align reads that span introns or fusion genes. After the first
188 round of alignment, TopHat breaks the unmapped reads into smaller parts and run a new alignment
189 round using Bowtie. This is one of the key strengths of TopHat, because it permits the identification
190 of splicing variants and fusion genes [14, 15].

191 The output generate a folder with a number of files and the most important of them are: the
192 “align_summary.txt” with a summary of the alignment, “accepted_hit.bam” with a list of read
193 alignments in the SAM format (Fig. 4), (“.bam” is a “binary SAM”), “unmapped.bam” with the
194 information on the reads that could not be aligned, “deletion.bed”, “insertion.bed”, and
195 “junctions.bed”, with the information described in its name in the Browser Extensible Data (BED)
196 format.

197

198 3.2 From BAM/SAM Files to a List of Differentially Expressed Genes

199

200 The need to analyze RNAseq data has given rise to a plethora of methods, with different
201 characteristics and assumptions [16]. Here we describe three widely used tools, the Tuxedo pipeline,
202 DESeq2 and Limma/voom [14, 17–19]. Tuxedo employs FPKM (Fragments Per Kilobase Of Exon Per
203 Million Fragments Mapped) in its analysis pipeline [14], while DESeq2 and Limma/voom accept the

204 counts derived from HTseq analysis as Input [16, 20]. As both HTseq and Cufflinks accept the “.BAM”
205 files, we can use the “accepted_hit.bam” from the previous step [14, 15] as input for these
206 programs.

207

208 3.2.1 The TUXEDO Pipeline

209

210 The TUXEDO package consists of Cufflinks, Cuffmerge, and Cuffdiff programs, more detailed
211 descriptions of these methods are available in the supplementary data of Trapnell et al. [21]. First,
212 the Cufflinks program that uses the “.bam” files from TopHat to assemble the reads into the most
213 probable transcripts and will give a “.gtf” file with the FPKM for the transcripts. Then Cuffmerge
214 takes the “.gtf” files of all the samples and generates a merged “.gtf” file as output. Cuffdiff
215 estimates the differential expression not only of the genes but also of their isoforms and the
216 promoters used (TSS— transcription start site) [14, 21]. The output is tab-delimited text files,
217 including the “.diff” files with the results and other files with information about the analyses, shown
218 in part in Fig. 5. The pipeline also contains the program CummeRbund that can be optionally used to
219 manage, integrate, and visualize the results produced by the Cufflinks package [14].

220 The Tuxedo Pipeline was for many years the de facto RNAseq analysis pipeline, due to its attractive
221 ability to provide gene expression, isoform variation, and TSS use. Nonetheless, it can only be
222 executed in the Linux command line, has multiple command steps, and require expertise in this
223 environment. Furthermore, it has also been reported that the Tuxedo pipeline has a lower precision
224 and sensitivity than DESeq 2 and Limma/Voom [16, 17].

225

226 3.2.2 HTSeq

227

228 HTSeq uses “.bam” files from the samples and a “.gtf” or “.gff” file with the gene models as input. It
229 is scripted in Python and the “HTSeq-count” script is able to count how many aligned reads overlap
230 with the exons of the genes, not considering differential splicing. Reads that align with more than
231 one gene are discarded, excluding fusion genes. The output is a tab-delimited text “.txt” file with two
232 columns: one with the gene name and one with the counts (Fig. 6) [20].

233 The “count” value is exactly as described. It is a measure of how many reads for that gene exist in
234 the “.bam” files. For the nature of the analysis, the “count” is nonnegative integer valued. A read can
235 be counted as belonging to a gene (+1 in the “count” value) or not [20]. This value is not corrected or
236 normalized in any way by HTSeq and downstream analysis will have to address this issue. At this
237 point, it is interesting to note that the same sample sequenced with different depths will have
238 different count values for each gene simply because of the different quantity of reads generated by
239 the sequencer. This issue needs to be considered when working with samples from different origins
240 and/or from multiple datasets [20, 22]. It is possible to introduce a prefiltering step at this point, i.e.,
241 excluding the genes that have zero counts in all samples, to reduce the needed computation time
242 and improve the statistical power as this step can affect the false discovery rate control [23].

243

244 3.2.3 DESeq2

245

246 The DESeq2 package uses unnormalized count data, such as that provided from HTseq, as input. The
247 program internally corrects for library size, so raw counts needs to be used in order to run the
248 analyses [24]. To run the analyses, it is necessary to create a DESeqDataSet object, using the “.txt”
249 counts files and an user-specified design matrix, i.e., assignment of samples to different treatment
250 groups [24]. At this point, the “DESeq” function is invoked to perform the differential analysis. The
251 output is a table containing the gene symbol, the base 2 log-transformed fold change, the pvalue,
252 and the adjusted p-value (or q-value, a measure of false discovery) (Fig. 7). DESeq2 assumes a
253 negative binomial linear model to describe over-dispersed count data from RNAseq data. It uses an
254 empirical Bayes method for more robust and accurate estimation of parameters for dispersion and
255 fold change, by taking into account small numbers of replicates and low read counts in RNAseq data.
256 Finally, DESeq2 uses the Wald test to estimate significance of differential expression, and the
257 Benjamini—Hochberg correction to control the false discovery rate [17, 24].

258

259 3.2.4 Limma/Voom

260

261 Limma is one of the well-known R packages used for differential gene expression analysis. It was
262 initially developed for microarray analyses, prior to the emergence of RNAseq, and has been
263 updated to facilitate analysis of RNASeq data [25]. Limma uses linear models to specify the
264 experimental design, empirical Bayes method to moderate the standard errors between genes, and
265 uses the t-test to calculate the differential expression p-value, while providing multiple choices for
266 adjustment of the p-value for multiple testing, including the Benjamini—Hochberg [19, 25]. The use
267 of Limma is well suited for small numbers of samples per group, as few as two, and also powerful
268 when used in multifactor designed tests [16, 25]. It has often reported that DESeq2 and Limma have
269 similar precision and sensitivity in their analysis results [16, 17].

270 As with DESeq2, Limma accepts the “.txt” count files from HTseq as input but internally transforms
271 the data because Limma assumes the t-distribution for gene expression values as it was originally
272 developed for microarray datasets [25, 26]. In order to address this issue, Limma uses the Voom
273 transformation, i.e., log transformation of counts per million (cpm) with associate precision weights
274 [18]. The Voom transformation is the key step in this analysis, as its log transformation helps gene
275 expression values satisfy the t-distribution assumption of Limma, while using cpm instead of raw
276 count normalizes gene expression values across replicates [18, 25]. Once we have the transformed
277 table, we use the “lmFit” function to fit the data to a linear model, which informs the design of the
278 experiment as in the case of DESeq2. The next step is to call the “eBayes” function that will calculate
279 the moderated t-statistics and log-odds of the differential expression using an empirical Bayes
280 moderation. The output is a table containing the gene symbol, the fold change in log scale with base
281 2, the p-value, and the adjusted p-value as shown in Fig. 8.

282

283 3.3 From a List of Differentially Expressed Genes to Systems Level Analyses

284

285 The final output of the Tuxedo pipeline, DESEQ2 and Limma/voom is a list of differentially expressed
286 genes [14, 24, 25]. Depending on the situation and the established cutoff, the length of this list can

287 vary from tens to thousands genes. In this chapter we consider the analysis comparing two groups of
288 experimental conditions but this can also be easily extended to multiple groups as well. The simplest
289 way to compare the genes from different experimental conditions is using Venn diagrams, which
290 permits to assess what is in common and what is unique in the gene lists of your differential
291 expression analysis [27]. In addition, when only a handful of genes are of interest, deeper
292 interrogation of each gene can be also implemented using Genecards and The Human Protein Atlas
293 [28, 29]. Nonetheless, most of the time, the differentially expressed gene list might be too long for
294 such gene-by-gene analyses and as a result, system level analyses might often be more appropriate,
295 for example, by using GSEA and/or Toppfun [30, 31].

296

297 3.3.1 Venn Diagrams — Venny

298

299 Area-proportional Venn diagrams are a useful graphic approach to compare different analyses.
300 BioVenn is a convenient web application for the comparison and visualization of biological lists [27].
301 For example, we can visualize how the gene lists generated by different DE analysis programs
302 described above are related to each other (Fig. 9). On the other hand, Venn diagrams can also be
303 used to compare the results of two different comparisons between groups and/or experimental
304 conditions. As we can see, the way we analyze our data influence the results. If information on one
305 specific gene being differentially expressed in a given condition is needed, it is often necessary to
306 validate it with qPCR [1].

307

308 3.3.2 Genecards and the Human Protein Atlas

309

310 Once in possession of a list of ranked DE genes, the next task is to make sense of it in a biological
311 context. One possibility is to check the DE gene list on a gene by gene basis, or at least the top
312 ranked genes in the list, for their function and expression in different tissues. For this step there are
313 two web services that can be utilized: Genecards, a repository with information of gene and protein
314 function, expression and known interactions [28] (Fig. 10); and The Human Protein Atlas, with
315 information on RNA and protein levels in different tissue and cancer types, and immunostaining data
316 from images derived from tissue histology [29] (Fig. 11). Although this method is useful to analyze a
317 small list of genes, or the function of the top ranked genes in a list, it is unpractical when we have a
318 list of hundreds or thousands of DE genes. In this case, a system level analysis is required.

319

320 3.3.3 Gene Set Enrichment Analysis (GSEA)

321

322 One of the key assumptions of the gene set enrichment analysis is that a moderate increase in
323 expression of a large number of genes encoding members of a pathway may dramatically alter the
324 flux through the pathway and may be more important than a huge increase in a single gene from
325 that pathway [30, 32]. GSEA is available as online and desktop versions. Both use the same
326 molecular signature database, comprised of eight collections: The Hallmark Gene sets (H), the
327 Positional gene set (C1), the Curated gene sets (C2), the Motif gene sets (C3), the Computational

328 gene sets (C4), the Gene Ontology gene sets (C5), The Oncogenic signatures (C6), and the
329 Immunologic signatures (C7) [30, 32]. The desktop version accepts a “.txt” tab-delimited table as
330 input, although the Limma/voom transformed table is preferable, as the program was originally
331 designed to analyze log₂ values from microarray data. Count values obtained directly from HTseq
332 are also accepted as input, although in this case it is important to check if the data needs to be
333 normalized [33]. In the basic analyses, the samples in the table are separated into two groups, test
334 group against control group. GSEA determines the enrichment of gene sets using a modified
335 Kolmogorov–Smirnov test. Specifically, it calculates the Enrichment Score (ES) taking into
336 consideration if a gene set is over represented at the top or at bottom of the ranked list, and it
337 estimates the p-values using empirical phenotype-based permutation test procedure, along with
338 their FDR [30, 33]. The output contains the table of gene sets that are positively and negatively
339 correlated with the groups, with ES, p values, and FDR values. In addition, GSEA also provides a
340 graphical representation of the analysis, containing a heatmap comparing the two groups in the
341 context of the gene set and an enrichment plot (Fig. 12). We note that in the GSEA analysis, it is
342 important to provide the complete transcriptome as a whole as input, rather than a previously
343 selected list of genes (e.g., from Limma/voom). The desktop version also accepts a preranked list of
344 genes as input for the analyses. In contrast, the online version only accepts a list of previously
345 selected genes as input, and tests the enrichment of gene sets using the hypergeometric test.
346 Specifically, it evaluates if any of the gene sets are over-represented in the provided gene list, and
347 generates a downloadable table with the statistically significant gene sets, the p-value, and the
348 overlap between the gene sets and the provided gene list (Fig. 13).

349

350 3.3.4 ToppGene Suite

351

352 ToppGene suite is a web-based service, with four functionalities: ToppFun, ToppGene, ToppNet, and
353 ToppGenet [31]. ToppFun (1), similar to the GSEA web based tool, accepts a gene list as input and
354 provides a downloadable table with the enriched pathways as output. Additionally it also generates
355 tables for 14 annotation categories including GO terms, pathways, protein–protein interaction,
356 microRNAs, and related diseases [31] (Fig. 14). ToppGene (2) takes two gene lists as input, i.e.,
357 “training gene set” and “test gene set.” In our case, the “test gene set” is the DE gene list while the
358 “training gene set” is a gene list containing gene of our interest. Given these gene lists, ToppGene will
359 run a ToppFun analyses of the “training gene set” to identify the most notable characteristics of this
360 list, and rank the genes in the “test gene set” according to these characteristics. For example, if the
361 “training gene set” is formed by the genes involved in Angiogenesis, the test genes will be ranked
362 based on their relevance to Angiogenesis. Similarly, if the “training gene set” is formed by
363 membrane proteins, the test genes will be ranked based on the degree they are related, or directly
364 interact, with membrane proteins. ToppGene uses all 14 annotation categories to analyze the
365 “training gene set” and rank the “test gene set” based on all of these 14 annotation categories [31].

366 ToppNet (3) and ToppGenet (4) compare “training gene set” with “test gene set” and construct
367 networks of interactions. Specifically, ToppNet is based only on protein–protein interaction (PPI),
368 while ToppGenet uses both PPI and the genes in the neighborhood in order to take in account
369 possible indirect interaction.

370

371

372 **4 Conclusions**

373

374 In this chapter, we described how to analyze RNAseq from raw data to a list of genes and related
375 systems and pathways. RNAseq has gained popularity with the emergence of high-throughput
376 sequencing (HTS) or next-generation sequencing (NGS). These advances in conjunction with
377 improvements in Proteomics and Metabolomics techniques and related analytical methods have
378 introduced a paradigm shift in biomedical research. This big data landscape was unimaginable just a
379 few short years ago [34].

380 This postgenomics era has enabled a new way of thinking, as it becomes clearer and clearer that the
381 events occurring in cells and tissues are much more complex than the sum of their parts [35]. The
382 vast amount of data, from diverse origins, and the need for integration can be overwhelming. Thus it
383 is imperative that the researcher have a clear idea of the question that is being asked, the data and
384 concepts that are used to formulate a hypothesis, and the assumptions that statistical methods are
385 based on, and the inherent limitations of these methods.

386

387 **Acknowledgments**

388

389 This work was conducted with support from start-up funds from the MUSC COM and an award from
390 SC Epscor to GH. The content is solely the responsibility of the authors and does not necessarily
391 represent the official views of the Medical University of South Carolina.

392

393 **References**

394

395 1. Alberts B, Johnson A, Lewis J et al (2015) Analyzing cells, molecules, and systems. In: Molecular
396 biology of the cell, 6th edn. Garland Science, USA, pp 440–528

397 2. Garrett CT (2015) Molecular biology basics in the “Omics” Era: genes to proteins. In: Idowu OM,
398 Dumur IC, Garrett TC (eds) Molecular oncology testing for solid tumors: a pragmatic approach.
399 Springer International Publishing, Cham, pp 3–65

400 3. Wang Z, Gerstein M, Snyder M (2009) RNA-Seq: a revolutionary tool for transcriptomics. *Nat Rev*
401 *Genet* 10(1):57–63

402 4. Mayer-Schoenberger V, Cukier K (2013) Big data: a revolution that will transform how we live,
403 work, and think. Harcourt, Houghton Mifflin

404 5. Sengupta S, Bolin JM, Ruotti V et al (2011) Single read and paired end mRNA-Seq Illumina libraries
405 from 10 nanograms total RNA. *J Vis Exp* (56):e3340. <https://doi.org/10.3791/3340>

406 6. Cock PJ, Fields CJ, Goto N et al (2010) The Sanger FASTQ file format for sequences with quality
407 scores, and the Solexa/Illumina FASTQ variants. *Nucleic Acids Res* 38 (6):1767–1771

- 408 7. Illumina (2011) Quality scores for next generation sequencing: assessing sequencing accuracy
409 using phred quality scoring.
410 https://www.illumina.com/Documents/products/technotes/technote_Q-Scores.pdf
- 411 8. Zhang T, Guo L, Creighton CJ et al (2016) A genetic cell context-dependent role for ZEB1 in lung
412 cancer. *Nat Commun* 7:12231
- 413 9. Lenferink AE (2017) Epithelial-to-Mesenchymal transition (EMT): the good, the bad, and the ugly.
414 In: Wang E (ed) *Cancer systems biology*. CRC Press, Florida
- 415 10. Andrews S (2010) FastQC: a quality control tool for high throughput sequence data.
416 <http://www.bioinformatics.babraham.ac.uk/projects/fastqc>
- 417 11. Gordon A, Hannon G (2010). Fastx-Toolkit. In: FASTQ/A short-reads preprocessing tools.
418 http://hannonlab.cshl.edu/fastx_toolkit/
- 419 12. Martin M (2011) Cutadapt removes adapter sequences from high-throughput sequencing reads.
420 *EMBnetjournal* 17(1):10–12
- 421 13. Speir ML, Zweig AS, Rosenbloom KR et al (2016) The UCSC genome browser database: 2016
422 update. *Nucleic Acids Res* 44(D1): D717–D725
- 423 14. Trapnell C, Roberts A, Goff L et al (2012) Differential gene and transcript expression analysis of
424 RNA-seq experiments with TopHat and cufflinks. *Nat Protoc* 7(3):562–578
- 425 15. Trapnell C, Pachter L, Salzberg SL (2009) TopHat: discovering splice junctions with RNA-Seq.
426 *Bioinformatics* 25(9):1105–1111
- 427 16. Sonesson C, Delorenzi M (2013) A comparison of methods for differential expression analysis of
428 RNA-seq data. *BMC Bioinformatics* 14:91
- 429 17. Love MI, Huber W, Anders S (2014) Moderated estimation of fold change and dispersion for RNA-
430 seq data with DESeq2. *Genome Biol* 15(12):550
- 431 18. Law CW, Chen Y, Shi W et al (2014) Voom: precision weights unlock linear model analysis tools
432 for RNA-seq read counts. *Genome Biol* 15(2):R29
- 433 19. Smyth GK (2005) Limma: linear models for microarray data. In: Gentleman R, Carey VJ, Huber W,
434 Irizarry RA, Dudoit S (eds) *Bioinformatics and computational biology solutions using R and*
435 *bioconductor*. Statistics for biology and health. Springer, New York, NY
- 436 20. Anders S, Pyl PT, Huber W (2015) HTSeq—a python framework to work with highthroughput
437 sequencing data. *Bioinformatics* 31(2):166–169
- 438 21. Trapnell C, Williams BA, Pertea G et al (2010) Transcript assembly and quantification by RNA-Seq
439 reveals unannotated transcripts and isoform switching during cell differentiation. *Nat Biotechnol*
440 28(5):511–515
- 441 22. Conesa A, Madrigal P, Tarazona S (2016) A survey of best practices for RNA-seq data analysis.
442 *Genome Biol* 17(1):181
- 443 23. Bourgon R, Gentleman R, Huber W (2010) Independent filtering increases detection power for
444 high-throughput experiments. *Proc Natl Acad Sci U S A* 107(21):9546–9551

- 445 24. Love M, Anders S, Huber W (2014) Differential analysis of count data—the DESeq2 package.
446 Genome Biol 15(12):550
- 447 25. Ritchie ME, Phipson B, Wu D et al (2015) Limma powers differential expression analyses for RNA-
448 sequencing and microarray studies. Nucleic Acids Res 43(7):e47
- 449 26. Datta S, Nettleton D (2014) Statistical analysis of next generation sequencing data. In: Datta S,
450 Nettleton D (eds) Frontiers in probability and the statistical sciences. Springer International
451 Publishing, Switzerland, pp 1–32
- 452 27. Hulsen T, De Vlieg J, Alkema W (2008) BioVenn—a web application for the comparison and
453 visualization of biological lists using areaproportional Venn diagrams. BMC Genomics 9:488
- 454 28. Stelzer G, Dalah I, Stein TI et al (2011) In-silico human genomics with GeneCards. Hum Genomics
455 5(6):709–717
- 456 29. Uhlen M, Oksvold P, Fagerberg L (2010) Towards a knowledge-based human protein atlas. Nat
457 Biotechnol 28(12):1248–1250
- 458 30. Subramanian A, Tamayo P, Mootha VK et al (2005) Gene set enrichment analysis: a knowledge-
459 based approach for interpreting genome-wide expression profiles. Proc Natl Acad Sci U S A
460 102(43):15545–15550
- 461 31. Chen J, Bardes EE, Aronow BJ et al (2009) ToppGene suite for gene list enrichment analysis and
462 candidate gene prioritization. Nucleic Acids Res 37(Web Server): W305–W311
- 463 32. Mootha VK, Lindgren CM, Eriksson KF et al (2003) PGC-1alpha-responsive genes involved in
464 oxidative phosphorylation are coordinately downregulated in human diabetes. Nat Genet 34(3):267–
465 273
- 466 33. Gene Set Enrichment Analysis (GSEA) User Guide (2010), [http://software.broadinstitute.
467 org/gsea/doc/GSEAUserGuideFrame.html](http://software.broadinstitute.org/gsea/doc/GSEAUserGuideFrame.html)
- 468 34. Wold B, Myers RM (2008) Sequence census methods for functional genomics. Nat Methods
469 5(1):19–21
- 470 35. Alon U (2006) An introduction to systems biology: design principles of biological circuits. CRC
471 Press, Florida

472

473

474

475

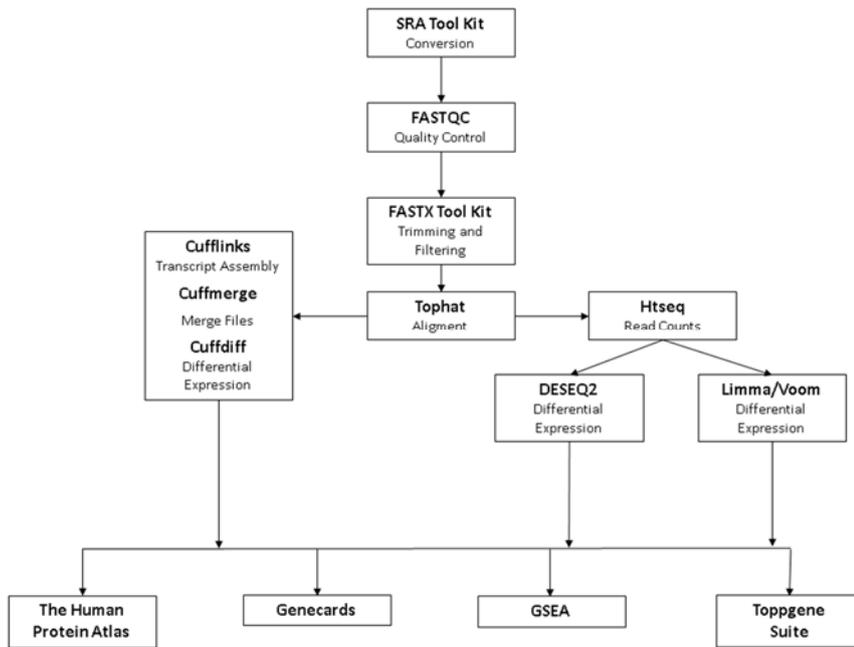
476

477

478

479

480



481

482 Fig. 1 Flowchart summarizing analyses from RNA seq data. The conversion step with SRA tool kit is
 483 generally only used with data from repositories

484

485

486

487

488

489

490

491

492

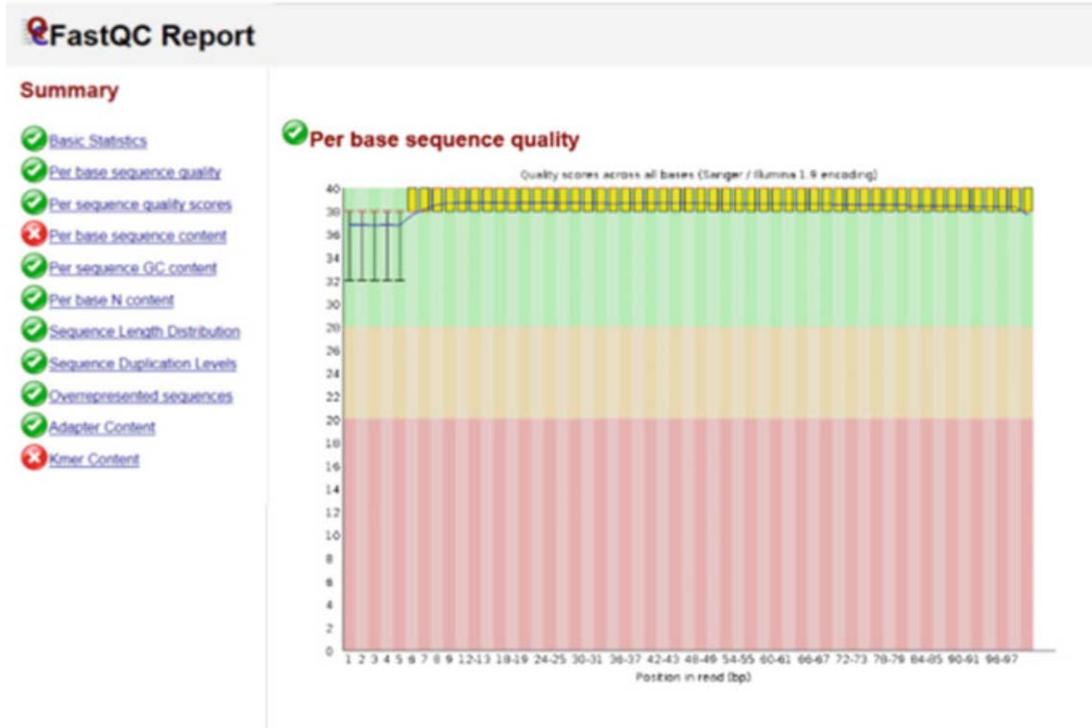
493

494

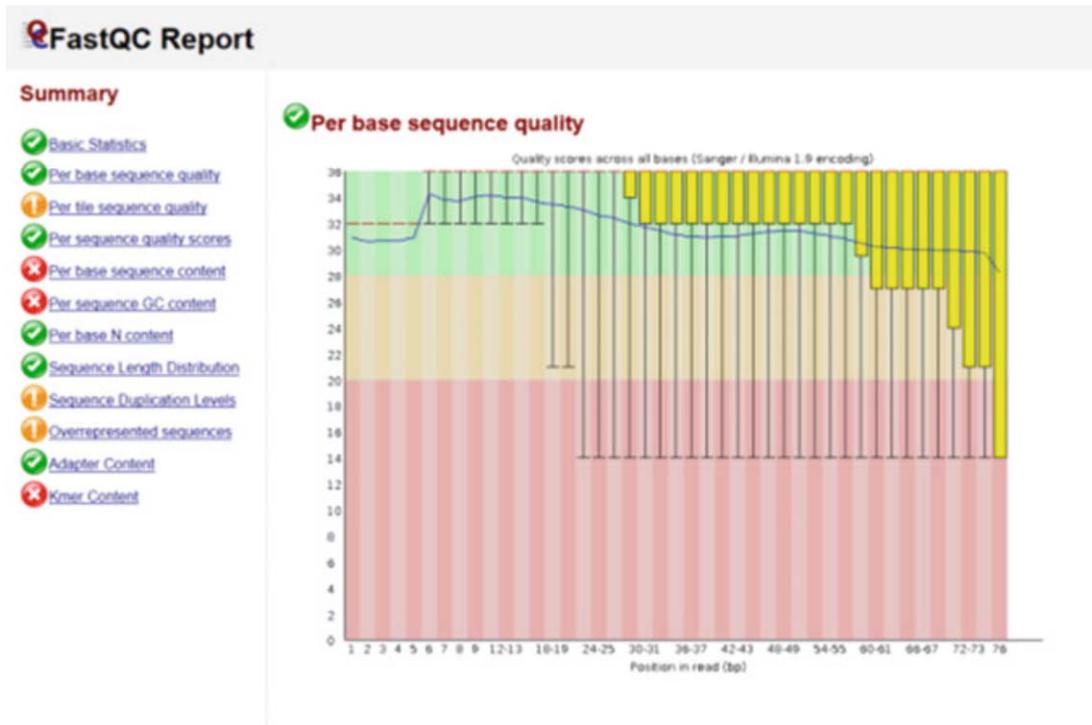
495

496

Good Report



Bad Report



513

514 Fig. 3 FastQC report. Example FastQC report for a good quality sample (Upper Panel) and a poor

515 quality sample (Lower Pane

516

gene_id	gene	locus	Group_1	Group_2	status	value_1	value_2	log2(fold_change)	test_stat	p_value	q_value
XLOC_000126	EFHD2	chr1:15736390-15756839	Control	ZEB1	OK	202.415	95.5723	-1.08265	-2.95758	5.00E-05	0.000518
XLOC_000221	HMG2	chr1:26798901-26803133	Control	ZEB1	OK	161.569	235.038	0.540743	1.63467	0.0131	0.0473428
XLOC_000351	CDC20	chr1:43824625-43828873	Control	ZEB1	OK	225.636	318.984	0.499487	1.62614	0.0133	0.0478431
XLOC_000370	RPS8	chr1:45241245-45244412	Control	ZEB1	OK	2517.32	1969.36	-0.354161	-1.67685	0.01245	0.045781
XLOC_000377	UROD	chr1:45477804-45481341	Control	ZEB1	OK	71.9628	180.767	1.32881	1.89676	0.0012	0.0079017
XLOC_000456	PGM1	chr1:64058946-64125916	Control	ZEB1	OK	142.789	88.7964	-0.685316	-1.38816	0.00725	0.0320071
XLOC_000512	CYR61	chr1:86046443-86049648	Control	ZEB1	OK	701.992	521.97	-0.427489	-1.96439	0.0023	0.0135386
XLOC_000730	PSMB4	chr1:151372040-151374412	Control	ZEB1	OK	254.892	435.399	0.772454	2.2223	0.00425	0.0207689
XLOC_000956	QSOX1	chr1:180123967-180169859	Control	ZEB1	OK	264.486	123.342	-1.10053	-2.63871	5.00E-05	0.000518
XLOC_000970	LAMC2	chr1:183155173-183214262	Control	ZEB1	OK	177.744	13.4283	-3.72645	-4.96459	5.00E-05	0.000518
XLOC_001020	ELF3	chr1:201979689-201986315	Control	ZEB1	OK	123.944	28.7934	-2.10589	-3.44096	5.00E-05	0.000518
XLOC_001074	G0S2	chr1:209848669-209849735	Control	ZEB1	OK	291.39	16.4758	-4.14453	-2.88957	0.00595	0.0278503
XLOC_001153	GALNT2	chr1:230193535-230417876	Control	ZEB1	OK	99.1635	167.39	0.755334	1.55569	0.00275	0.0147362
XLOC_001311	ENO1	chr1:8921058-8939943	Control	ZEB1	OK	2012.97	2502.59	0.314096	1.46414	0.01085	0.0417349
XLOC_001572	SLC2A1	chr1:43391045-43449029	Control	ZEB1	OK	281.788	139.861	-1.01062	-3.46746	5.00E-05	0.000518
XLOC_001573	EBNA1BP2	chr1:43629844-43720029	Control	ZEB1	OK	124.382	212.616	0.773473	1.53116	0.0117	0.0437062
XLOC_001594	PRDX1	chr1:45965855-45988562	Control	ZEB1	OK	528.434	802.64	0.60303	1.88551	0.0008	0.006035
XLOC_001609	PDZK1IP1	chr1:47649260-47655771	Control	ZEB1	OK	358.718	37.0944	-3.27358	-3.42087	0.0076	0.0328067
XLOC_001672	JAK1	chr1:65210777-65432187	Control	ZEB1	OK	116.975	62.8851	-0.895414	-2.55947	5.00E-05	0.000518
XLOC_001756	F3	chr1:94994731-95007413	Control	ZEB1	OK	295.327	72.2413	-2.03142	-4.48758	5.00E-05	0.000518
XLOC_001933	S100A10	chr1:151955385-151966714	Control	ZEB1	OK	337.267	168.961	-0.997202	-2.53059	0.0017	0.0104833

536

537 Fig. 5 Cuffdiff output. Top 20 differentially expressed genes from Cuffdiff analysis

538

539

540

541

542

543

544

Symbol	Counts
A1BG-AS1	1
AAAS	3
AACS	2
AAGAB	5
AAK1	1
AAMP	4
AAR2	1
AARS	5
AASS	1
AATF	3
ABCA1	10
ABL2	11
ACTG1	205
CEP41	1
EEF1A1	419
GAPDH	453
TTC3	12
TTYH3	12
TUBA1B	200
TUBA1C	202

545

546 Fig. 6 HTSeq output. Example of count results from sample GSM2144086

547

548

549

550

551

552

553

554

555

Symbol	GeneID	Count_Control_1	Count_Control_2	Count_Control_3	Count_ZEB1_1	Count_ZEB1_2	Count_ZEB1_3	baseMean	log2FoldChange	lfcSE	stat	pvalue	padj
KRT7	3855	150	164	131	25	23	27	88.32087747	-2.719804965	0.21121846	-12.8767	6.09E-38	2.97E-34
FN1	2335	85	68	75	19	10	17	46.51829	-2.440348331	0.27796973	-8.77919	1.65E-18	4.02E-15
ACTG1	71	222	205	190	85	116	110	155.1586327	-1.136943225	0.13453789	-8.45073	2.89E-17	4.71E-14
MT2A	4502	41	34	20	112	139	128	74.02538642	1.783708582	0.21425014	8.325355	8.41E-17	1.03E-13
LAMC2	3918	50	43	49	4	1	4	25.99212881	-3.852052972	0.46562369	-8.27289	1.31E-16	1.28E-13
KRT18	3875	121	126	101	46	48	43	81.07128371	-1.514812939	0.18871538	-8.02697	9.99E-16	8.13E-13
CKCL1	2919	54	42	28	3	0	0	21.87722077	-4.679926755	0.60298745	-7.76123	8.41E-15	5.87E-12
LCN2	3934	40	51	29	0	1	1	20.97514876	-4.914503581	0.63920532	-7.68845	1.49E-14	9.09E-12
JUP	3728	45	32	35	4	4	2	20.96486972	-3.393701917	0.46634748	-7.2772	3.41E-13	1.66E-10
TNFAIP2	7127	34	42	41	4	3	6	22.22526049	-3.150144598	0.48211452	-7.29007	3.10E-13	1.66E-10
LAMB3	3914	42	28	30	1	2	1	17.94658771	-4.17558093	0.58055798	-7.19236	6.37E-13	2.83E-10
VIM	7431	11	14	10	59	70	64	35.42369692	2.206368812	0.30934126	7.132475	9.86E-13	4.01E-10
CDH1	999	22	38	34	1	1	0	16.59529808	-4.596894063	0.65133486	-7.05765	1.69E-12	6.36E-10
KRT8	3856	115	131	129	69	58	64	94.06276097	-1.152540615	0.1702325	-6.77039	1.28E-11	4.48E-09
LCP1	3936	21	34	24	0	0	0	13.64243193	-5.136608256	0.76190433	-6.7418	1.56E-11	5.09E-09
SAAL	6288	25	29	23	1	0	0	13.4714644	-4.685365375	0.70570744	-6.63925	3.15E-11	9.62E-09
PIP4K2C	79837	94	85	81	37	43	40	63.29305063	-1.287364914	0.20737209	-6.208	5.37E-10	1.54E-07
PLAU	5328	55	55	34	11	15	15	31.04926848	-1.936850937	0.31515463	-6.14572	7.96E-10	2.16E-07
CKCL2	2920	10	20	27	0	0	0	9.925612632	-4.723188071	0.78409378	-6.02375	1.70E-09	4.38E-07
PLEC	5339	70	87	64	36	30	30	52.83260945	-1.365602851	0.23073187	-5.91857	3.25E-09	7.93E-07

556

557 Fig. 7 DESeq2 output. Top 20 differentially expressed genes from DESeq2 analyzes

558

559

560

561

562

563

564

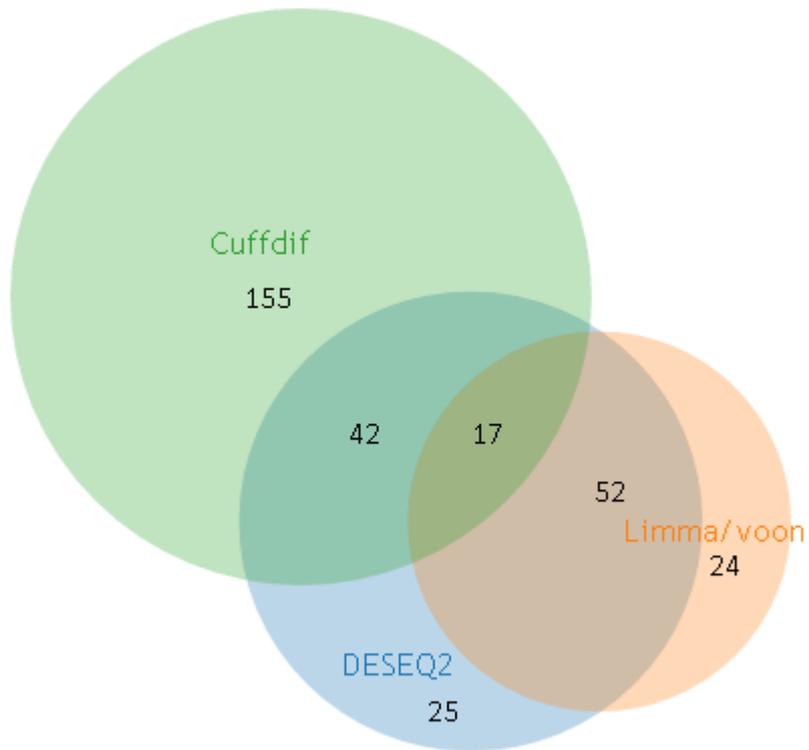
565

Symbol	logFC	AveExpr	t	P.Value	adj.P.Val	B
LCP1	-5.85418	6.573509	-16.9375	2.28E-07	0.001111	7.125779
FAM83A	-4.33212	5.812479	-13.5217	1.22E-06	0.002822	5.874455
ESRP1	-4.31032	5.801578	-12.7464	1.90E-06	0.002822	5.520751
CXCL2	-5.31476	6.303797	-12.0713	2.83E-06	0.002822	5.186808
COL4A2	-4.89651	6.094674	-11.781	3.38E-06	0.002822	5.035188
IGFN1	3.975878	5.77204	11.56598	3.87E-06	0.002822	4.919587
COL4A1	-3.67745	5.485143	-11.419	4.25E-06	0.002822	4.838877
LAMB3	-4.36474	7.659404	-10.9862	5.64E-06	0.002822	4.593108
CYLD	3.482421	5.525311	10.80534	6.36E-06	0.002822	4.486671
LCN2	-5.38404	7.395079	-10.7332	6.68E-06	0.002822	4.443573
LRRC6	-3.58095	5.436892	-10.6768	6.94E-06	0.002822	4.409643
NR5A2	3.939291	5.753746	10.40548	8.37E-06	0.002822	4.243081
CGN	-3.58095	5.436892	-10.3526	8.68E-06	0.002822	4.210012
ECH1	3.402085	5.485143	10.35119	8.69E-06	0.002822	4.209112
SAA1	-5.31258	6.831029	-10.2951	9.04E-06	0.002822	4.173766
GALNT3	-3.40411	5.348473	-10.263	9.24E-06	0.002822	4.153441
SCN9A	3.225246	5.396724	9.984322	1.13E-05	0.003079	3.973631
CDA	-3.70996	5.501396	-9.97544	1.13E-05	0.003079	3.967801
TNFRSF10A	-3.28325	5.288044	-9.67495	1.41E-05	0.003625	3.766785
SPTLC3	3.583743	5.575972	9.609286	1.48E-05	0.003625	3.72187
CDH3	-3.45332	7.44935	-9.45097	1.67E-05	0.00386	3.612084

566

567 Fig. 8 Limma/Voom output. Top 20 differentially expressed genes from Limma/Voom analyses

568



569

570 Fig. 9 Venny output. Proportional Venn diagram comparing the list of differentially expressed genes
 571 uncovered by Cuffdiff, DESeq2 and Limma/Voom, where differential expression was determined at
 572 adjusted p-value of < 0.05 for all the three methods

573

574

GeneCards® HUMAN GENE DATABASE

Free for academic non-profit institutions. Other users need a Commercial license

WEIZMANN INSTITUTE OF SCIENCE

LifeMap

Keywords + Search Term

Advanced

Home User Guide Analysis Tools - News And Views About - My Genes Log In / Sign Up

ZEB1 Gene (Protein Coding)

Zinc Finger E-Box Binding Homeobox 1

★ GCID: GC10P031330 ? GIFIS: 65 ?

Summaries for ZEB1 Gene

Entrez Gene Summary for ZEB1 Gene

This gene encodes a zinc finger transcription factor. The encoded protein likely plays a role in transcriptional repression of interleukin 2. Mutations in this gene have been associated with posterior polymorphous corneal dystrophy-3 and late-onset Fuchs endothelial corneal dystrophy. Alternatively spliced transcript variants encoding different isoforms have been described [provided by RefSeq, Mar 2010]

GeneCards Summary for ZEB1 Gene

ZEB1 (Zinc Finger E-Box Binding Homeobox 1) is a Protein Coding gene. Diseases associated with ZEB1 include *corneal dystrophy, fuchs endothelial, 6* and *corneal dystrophy, posterior polymorphous, 3*. Among its related pathways are *MicroRNAs in cancer* and *ERK Signaling*. GO annotations related to this gene include *nucleic acid binding* and *chromatin binding*. An important paralog of this gene is ZEB2.

UniProtKB/Swiss-Prot for ZEB1 Gene ZEB1_HUMAN.P37275

Acts as a transcriptional repressor. Inhibits interleukin-2 (IL-2) gene expression. Enhances or represses the promoter activity of the ATP1A1 gene depending on the quantity of cDNA and on the cell type. Represses E-cadherin promoter and induces an epithelial-mesenchymal transition (EMT) by recruiting SMARCA4/BRG1. Represses BCL6 transcription in the presence of the corepressor CTBP1. Positively regulates neuronal differentiation. Represses RCOR1 transcription activation during neurogenesis. Represses transcription by binding to the E box (5-CANNTG-3). Promotes tumorigenicity by repressing stemness-inhibiting microRNAs.

Gene Wiki entry for ZEB1 Gene

No data available for Tocris Summary , PharmGKB "VIP" Summary , BRINAb sequence ontologies and piRNA Summary for ZEB1 Gene

Genomics for ZEB1 Gene

Products: Regulatory Element

Genomic Location for ZEB1 Gene

Chromosome: 10
 Start: 31,318,495 bp from pter End: 31,529,814 bp from pter
 Size: 211,320 bases Orientation: Plus strand

Genomic View for ZEB1 Gene

Genes around ZEB1 on UCSC Golden Path with GeneCards custom track

Cytogenetic band: 10p11.22 by Ensembl 10p11.2 by Entrez Gene 10p11.22 by HGNC

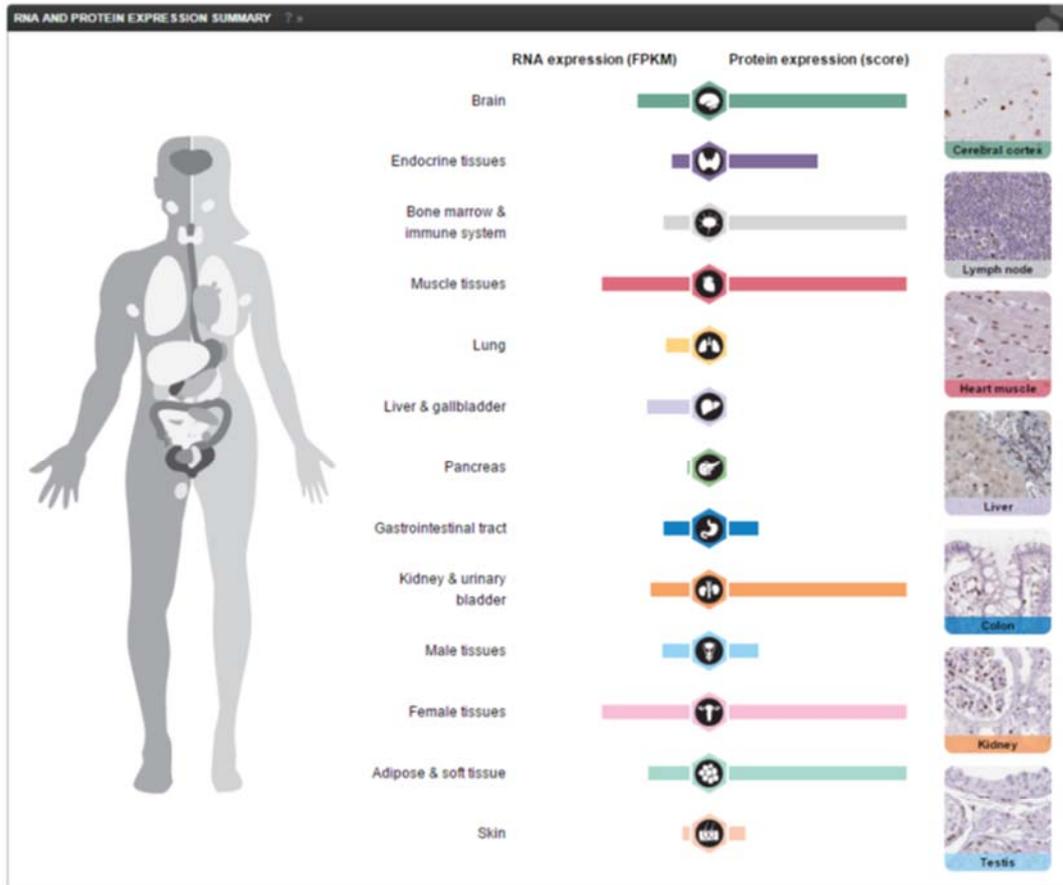
ZEB1 Gene in genomic location: bands according to Ensembl, locations according to GeneLoc (and/or Entrez Gene and/or Ensembl if different)

575

576 Fig. 10 ZEB1 edited description in Genecards

577

ZEB1

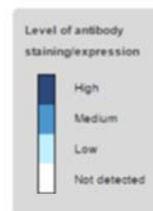


STAINING SUMMARY

HPA027524 CAB058686

Tissue	Cancer staining	Protein expression of normal tissue	Tissue	Cancer staining	Protein expression of normal tissue
Breast cancer	<input type="checkbox"/>	<input type="checkbox"/>	Melanoma	<input type="checkbox"/>	<input type="checkbox"/>
Carcinoid	<input type="checkbox"/>	<input type="checkbox"/>	Ovarian cancer	<input type="checkbox"/>	<input type="checkbox"/>
Cervical cancer	<input type="checkbox"/>	<input type="checkbox"/>	Pancreatic cancer	<input type="checkbox"/>	<input type="checkbox"/>
Colorectal cancer	<input type="checkbox"/>	<input type="checkbox"/>	Prostate cancer	<input type="checkbox"/>	<input type="checkbox"/>
Endometrial cancer	<input type="checkbox"/>	<input type="checkbox"/>	Renal cancer	<input type="checkbox"/>	<input type="checkbox"/>
Glioma	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Skin cancer	<input type="checkbox"/>	<input type="checkbox"/>
Head and neck cancer	<input type="checkbox"/>	<input type="checkbox"/>	Stomach cancer	<input type="checkbox"/>	<input type="checkbox"/>
Liver cancer	<input type="checkbox"/>	<input type="checkbox"/>	Testis cancer	<input type="checkbox"/>	<input type="checkbox"/>
Lung cancer	<input type="checkbox"/>	<input type="checkbox"/>	Thyroid cancer	<input type="checkbox"/>	<input type="checkbox"/>
Lymphoma	<input checked="" type="checkbox"/>	<input type="checkbox"/>	Urothelial cancer	<input type="checkbox"/>	<input type="checkbox"/>

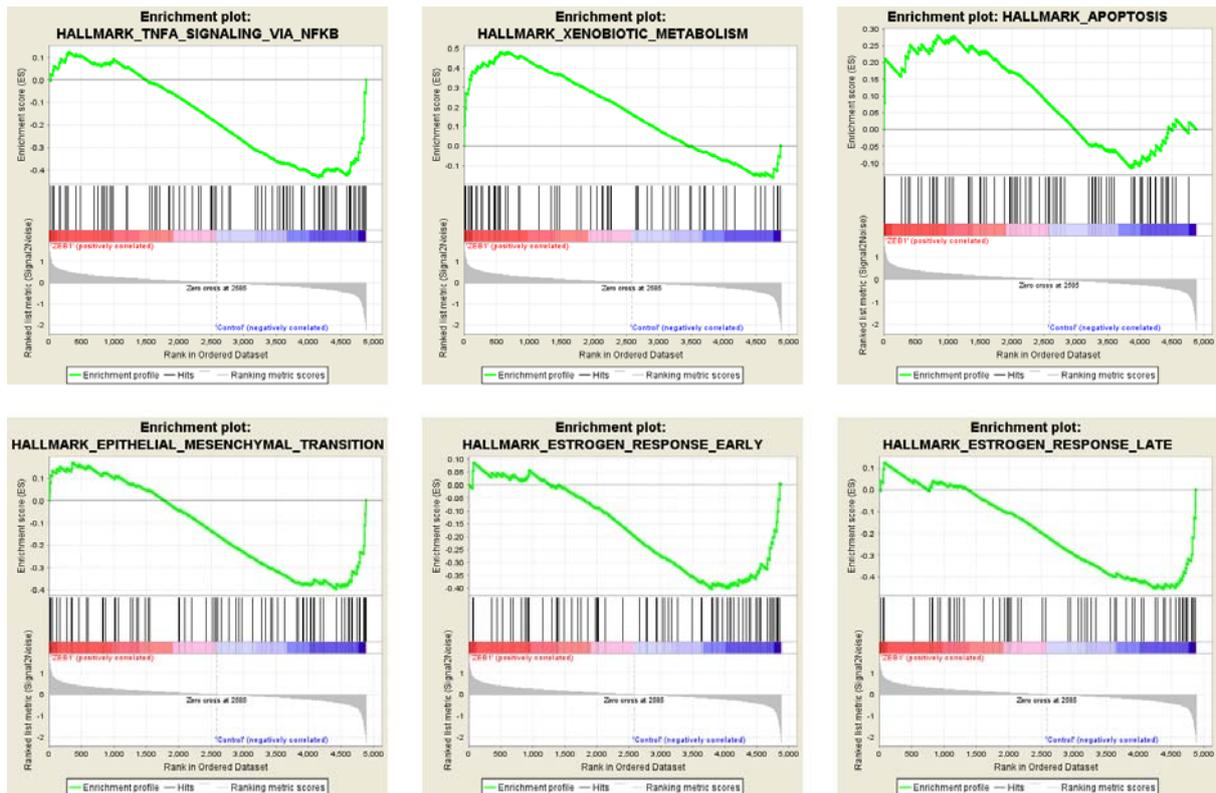
Staining summary: Most cases of gliomas along with several lymphomas and a few melanomas showed moderate to strong nuclear positivity. A single case of high grade urothelial cancer showed strong positivity. Rare liver, renal and testis cancers showed moderate positivity. Remaining cancer tissues were negative.



578

579 Fig. 11 ZEB1 edited description in The Human Protein Atlas website

580



581

582 Fig. 12 Results from GSEA desktop version (using the transformed table from Limma/Voom as input)

583

584

Gene Set Name [# Genes (K)]	Description	# Genes in Overlap (k)	k/K	p-value	FDR q-value
HALLMARK_EPITHELIAL_MESENCHYMAL_TRANSITION [200]	Genes defining epithelial-mesenchymal transition, as in wound healing, fibrosis and metastasis.	14		1.19 e-15	5.96 e-14
HALLMARK_TNFA_SIGNALING_VIA_NFKB [200]	Genes regulated by NF-kB in response to TNF [GeneID=7124].	12		8.77 e-13	2.19 e-11
HALLMARK_ESTROGEN_RESPONSE_EARLY [200]	Genes defining early response to estrogen.	10		4.44 e-10	5.55 e-9
HALLMARK_XENOBIOTIC_METABOLISM [200]	Genes encoding proteins involved in processing of drugs and other xenobiotics.	10		4.44 e-10	5.55 e-9
HALLMARK_ESTROGEN_RESPONSE_LATE [200]	Genes defining late response to estrogen.	9		8.53 e-9	8.53 e-8
HALLMARK_ANDROGEN_RESPONSE [101]	Genes defining response to androgens.	7		2.09 e-8	1.75 e-7
HALLMARK_APOPTOSIS [161]	Genes mediating programmed cell death (apoptosis) by activation of caspases.	8		2.73 e-8	1.95 e-7
HALLMARK_CHOLESTEROL_HOMEOSTASIS [74]	Genes involved in cholesterol homeostasis.	6		8.64 e-8	5.4 e-7
HALLMARK_APICAL_JUNCTION [200]	Genes encoding components of apical junction complex.	8		1.46 e-7	6.63 e-7
HALLMARK_HYPOXIA [200]	Genes up-regulated in response to low oxygen levels (hypoxia).	8		1.46 e-7	6.63 e-7

585

586 Fig. 13 GSEA web tool Analysis. Top 10 Hallmarks gene sets enriched for the DE gene list obtained
587 from DESeq2

588

1: GO: Molecular Function [Display Chart] 557 annotations before applied cutoff / 18661 genes in category

ID	Name	Source	pValue	FDR B&H	FDR B&Y	Bonferroni	Genes from Input	Genes in Annotation
1	GO:0005198	structural molecule activity	8.892E-8	4.953E-5	3.418E-4	4.953E-5	21	762
2	GO:0008201	heparin binding	2.937E-6	8.181E-4	5.645E-3	1.636E-3	9	167
3	GO:0045236	CXCR chemokine receptor binding	5.301E-6	9.841E-4	6.791E-3	2.952E-3	4	17
4	GO:0005200	structural constituent of cytoskeleton	1.312E-5	1.827E-3	1.261E-2	7.310E-3	7	110
5	GO:0005539	glycosaminoglycan binding	2.609E-5	2.906E-3	2.006E-2	1.453E-2	9	219

[Show 22 more annotations](#)**2: GO: Biological Process** [Display Chart] 3331 annotations before applied cutoff / 18623 genes in category

ID	Name	Source	pValue	FDR B&H	FDR B&Y	Bonferroni	Genes from Input	Genes in Annotation
1	GO:0060429	epithelium development	2.597E-10	8.652E-7	7.517E-6	8.652E-7	32	1296
2	GO:0042127	regulation of cell proliferation	5.978E-10	9.956E-7	8.650E-6	1.991E-6	36	1666
3	GO:0040012	regulation of locomotion	1.500E-9	1.500E-6	1.303E-5	4.996E-6	25	866
4	GO:0040011	locomotion	1.801E-9	1.500E-6	1.303E-5	5.999E-6	36	1735
5	GO:0048870	cell motility	2.897E-9	1.609E-6	1.398E-5	9.651E-6	32	1428

[Show 45 more annotations](#)**3: GO: Cellular Component** [Display Chart] 360 annotations before applied cutoff / 19061 genes in category

ID	Name	Source	pValue	FDR B&H	FDR B&Y	Bonferroni	Genes from Input	Genes in Annotation
1	GO:0005925	focal adhesion	3.882E-12	7.361E-10	4.758E-9	1.397E-9	20	393
2	GO:0005924	cell-substrate adherens junction	4.888E-12	7.361E-10	4.758E-9	1.760E-9	20	398
3	GO:0030055	cell-substrate junction	6.134E-12	7.361E-10	4.758E-9	2.208E-9	20	403
4	GO:0005912	adherens junction	2.181E-11	1.963E-9	1.269E-8	7.853E-9	21	484
5	GO:0070161	anchoring junction	4.466E-11	3.216E-9	2.079E-8	1.608E-8	21	503

[Show 45 more annotations](#)

589

590 Fig. 14 Results from TOPPFUN. Gene Ontology terms enriched for the gene list obtained from

591 DESEQ2 are presented