

# 1 TED toolkit: a comprehensive approach for 2 convenient transcriptomic profiling as a 3 clinically-oriented application

4 Thahmina Ali<sup>1</sup>, Baekdoo Kim<sup>1</sup>, Carlos Lijeron<sup>1</sup>, Olorunseun O.  
5 Ogunwobi<sup>1,2,3</sup>, Raja Mazumder<sup>5,6</sup>, and Konstantinos Krampis<sup>1,2,4</sup>

6 <sup>1</sup>Weill Cornell Medicine - Belfer Research Building, Hunter College of The City  
7 University of New York, New York, NY

8 <sup>2</sup>Department of Biological Sciences, Hunter College of The City University of New York,  
9 NY

10 <sup>3</sup>Joan and Sanford I. Weill Department of Medicine, Weill Cornell Medical College,  
11 Cornell University, New York, NY

12 <sup>4</sup>Department of Physiology and Biophysics, Institute for Computational Biomedicine,  
13 Weill Cornell Medical College, Cornell University, New York, NY

14 <sup>5</sup>The Department of Biochemistry & Molecular Medicine The George Washington  
15 University Medical Center, Washington, DC

16 <sup>6</sup>The McCormick Genomic and Proteomic Center, The George Washington University,  
17 Washington, DC

18 Corresponding author:

19 Konstantinos Krampis<sup>1,2,4</sup>

20 Email address: [kk104@hunter.cuny.edu](mailto:kk104@hunter.cuny.edu)

## 21 ABSTRACT

22 In translational medicine, the technology of RNA sequencing (RNA-seq) continues to prove powerful, and  
23 transforming the RNA-seq data into biological insights has become increasingly imperative. We present  
24 the Transcriptomics profiler for Easy Discovery (TED) toolkit, a comprehensive approach to processing  
25 and analyzing RNA-seq data. TED is divided into three major modules: data quality control, transcriptome  
26 data analysis, and data discovery, with eleven pipelines in total. These pipelines perform the preliminary  
27 steps from assessing and correcting the quality of the RNA-seq data, to the simultaneous analysis of five  
28 transcriptomic features (differentially expressed coding, non-coding, novel isoform genes, gene fusions,  
29 alternative splicing events, genetic variants of somatic and germline mutations) and ultimately translating  
30 the RNA-seq analysis findings into actionable, clinically-relevant reports. TED was evaluated using  
31 previously published prostate cancer transcriptome data where we observed previously studied outcomes,  
32 and also created a knowledge database of highly-integrated, biologically relevant reports demonstrating  
33 that it is well-positioned for clinical applications. TED is implemented on an instance of the Galaxy platform  
34 ( Galaxy page: <http://galaxy.hunter.cuny.edu/u/bioitcore/p/transcriptomics-profiler-for-easy-discovery-ted-toolkit>,  
35 Documentation Manual: <http://ted.readthedocs.io/en/latest/index.html>) as intuitive and reproducible  
36 pipelines providing a manageable strategy for conducting substantial transcriptome analysis in a routine  
37 and sustainable fashion for bioinformatics and clinical researchers alike.

## 38 INTRODUCTION

39 The modern sequencing technology, next generation sequencing (NGS) has expanded the analytical  
40 possibilities of the transcriptome in complete depth, the method known as RNA-sequencing (RNA-seq).  
41 RNA-seq can precisely determine the abundance of transcripts expressed in any RNA sample of study.  
42 Moreover, given the emergence of RNA-seq applications in many biomedical research areas, there are  
43 significant efforts in standardizing the method (1) within clinical settings. In the clinical laboratory,  
44 investigating the transcriptome has uncovered invaluable information of genetic mechanisms within a  
45 RNA sample of a conditioned or diseased individual (2, 3). The thorough view of the transcriptome

46 offered by RNA-seq offers ways for identifying disease causing bio-molecules of an individual that can  
47 serve as potential diagnostic indicators. This is especially applicable to complex diseases like cancer,  
48 where multiple bio-molecules contribute to its abnormal state, and findings through RNA-seq can be used  
49 as a reliable resource for therapeutic targets. In parallel with the considerable RNA-seq applications in  
50 the clinic, analyzing the RNA-seq data is essential, but delivering the biological insights unraveled from  
51 the analysis in the most informative means has become just as crucial. There are various data analysis  
52 programs most notably the Galaxy biomedical research platform (4) that addresses challenges such as  
53 the issues of accessibility and reproducibility. The platform provides an intuitive web based interface  
54 that serves as a workspace for data analysis in which researchers can import their data sets, and apply  
55 bioinformatics tools that are made available from the Galaxy toolshed (5) panel. Galaxy tools can run  
56 as standalone or chained together to create larger analyses transforming entire bioinformatics pipelines  
57 into automated “Galaxy workflows”. By Galaxy offering the ability to create and perform automated  
58 analyses on a user interface fully operational on the web, bioinformatics analyses have become more  
59 approachable in doing all types of data analysis. Yet, there still does not exist a convenient framework  
60 mainstream enough to enable RNA sequencing analysis results in a way that readily lends itself to easy  
61 interpretation. The current approach of performing RNA sequencing analyses is difficult, especially for  
62 non-bioinformatics researchers for the following reasons: (i) analysis methods and protocols are organized  
63 in a non-uniformed manner; (ii) analysis methods dependencies, parameters or supporting data come  
64 across as undocumented (iii) analyses output is in raw file state that consist of incomprehensible results  
65 with no set process to interpret them. These aspects lead to prolonged complexity requiring a learning  
66 curve to understand and tackle them which in turn causes a distraction in performing the actual analysis,  
67 making standardizing RNA sequencing analysis as a diagnostic practice challenging.

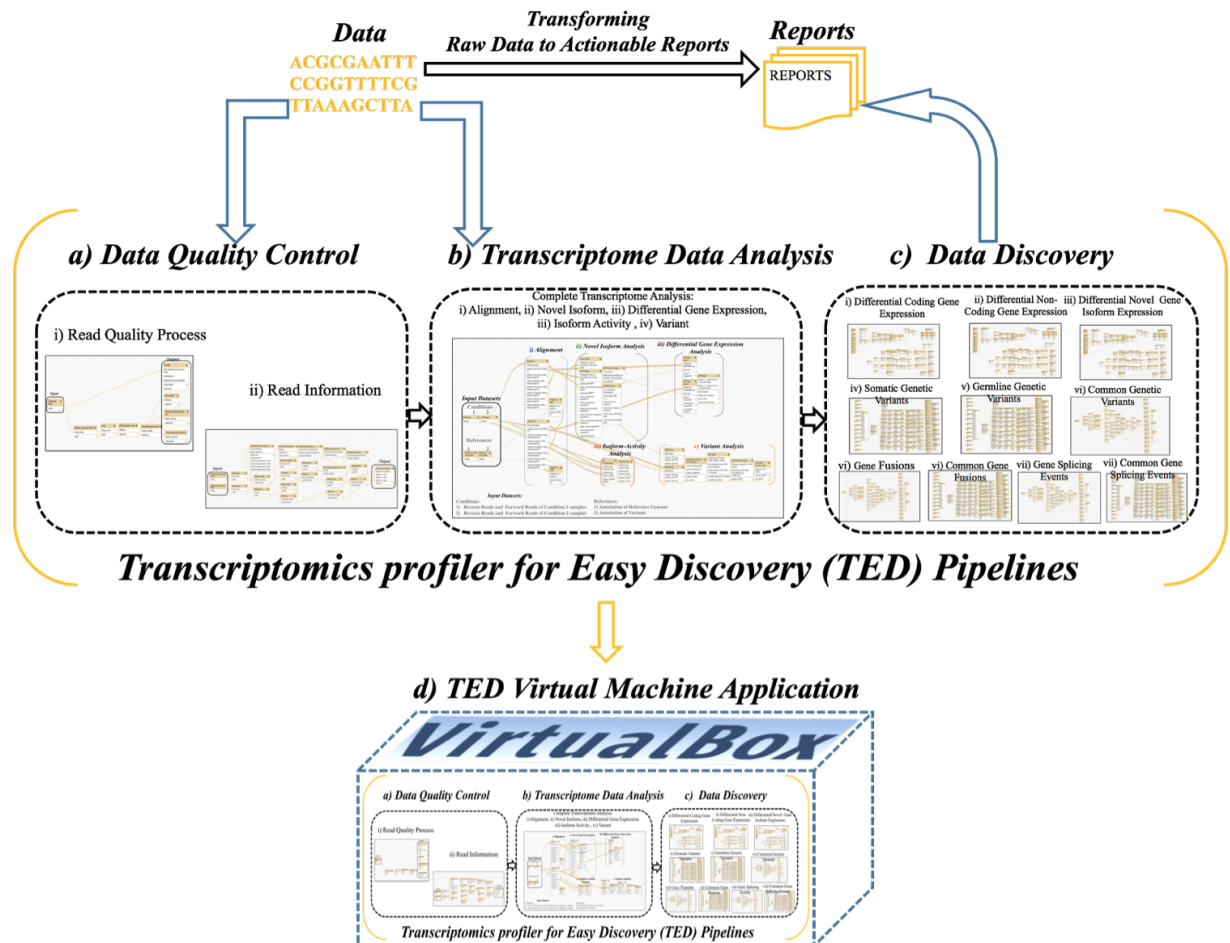
68 The bioinformatics pipelines that have been developed on the Galaxy platform, have had a focus on  
69 automation and standardization, including several pipelines available for transcriptomic data analysis.  
70 For example, the Oqtans (6) workbench performs differential expression and enrichment analysis and  
71 the open pipelines for tumor genome profiling that consist of three separate analyses pipelines: exome,  
72 transcriptome and variant evaluation (7). In addition, the TRAPLINE pipeline (8) performs comparative  
73 transcriptomics analysis, identifying a set of differentially expressed genes and their corresponding protein-  
74 protein interactions. These Galaxy pipelines have accelerated the extensibility in the transcriptome data  
75 analysis, however, in order to visualize the outputs requires importing to external programs. For example,  
76 the TRAPLINE protein-protein interactions output requires the Cytoscape program for visualization, in  
77 which this method does not enable direct interpretation delivered straight from the analysis exclusively.  
78 There are other automated pipelines that are taking initiatives in striving to bring out the most informed  
79 data analysis, by way of a software application approach. RNAseq software methods such as RobiNA (9)  
80 which uses a biostatistical method and Grape (10), both of which provide an environment to analyze and  
81 visualize gene expression data but limited to solely performing differential gene expression analysis. The  
82 Chipster (11) platform houses a comprehensive collection of analysis tools that covers analysis other than  
83 gene expression, such as miRNA, methylation and others, yet has complicated installation procedures,  
84 as well as, technical navigation again requiring a learning curve for non-informatics individuals. There  
85 are methods that function on the web such as MeV (12) which is cloud based that is also limited to  
86 performing differential gene expression analysis and visualization and the functionalities offered stratify  
87 the data analysis with curations that consist of no annotative feature especially with biological content.  
88 Nevertheless, each of these applications still are contributors to the steps towards the potential for  
89 standardizing RNA sequencing within the reach of translational and diagnostic settings.

90 We propose a highly-integrated set of bioinformatics pipelines designed in the form of automated  
91 workflows, which are implemented into the Galaxy platform. The workflows are configured to perform  
92 quality control and analysis on RNAseq data, while also providing beyond the standard analysis in order to  
93 provide data discovery functionality. The entire set of workflows is packaged as a resource toolkit, termed  
94 Transcriptome profiler for Easy Discovery, or TED. TED has three fundamental modules, summarized in  
95 Fig. 1. The first module provides quality control of the RNAseq data which are preprocessing steps, as  
96 well as, acquiring information about the reads such as read length, insert size etc. The second module  
97 carries out analysis of differentially coding, non-coding and novel isoform gene expression, gene fusions,  
98 alternative splicing events, and genetic variants of somatic and germline mutations of the RNAseq data.  
99 And lastly, the third module transforms the analysis results produced from the second module into  
100 detailed, biologically interpreted annotated reports. TED joins these three modules together creating a

101 knowledge database of prioritize biological outcomes, enabling users to obtain a comprehensive insight  
 102 of the transcriptome analyzed from the RNA samples. TED becomes extensible to applications in clinical  
 103 or diagnostic scenarios, allowing the user as a clinician or practitioner to leverage their experience to data  
 104 mine the reports of analyzed results for discovery or indication of biological candidates to examine.

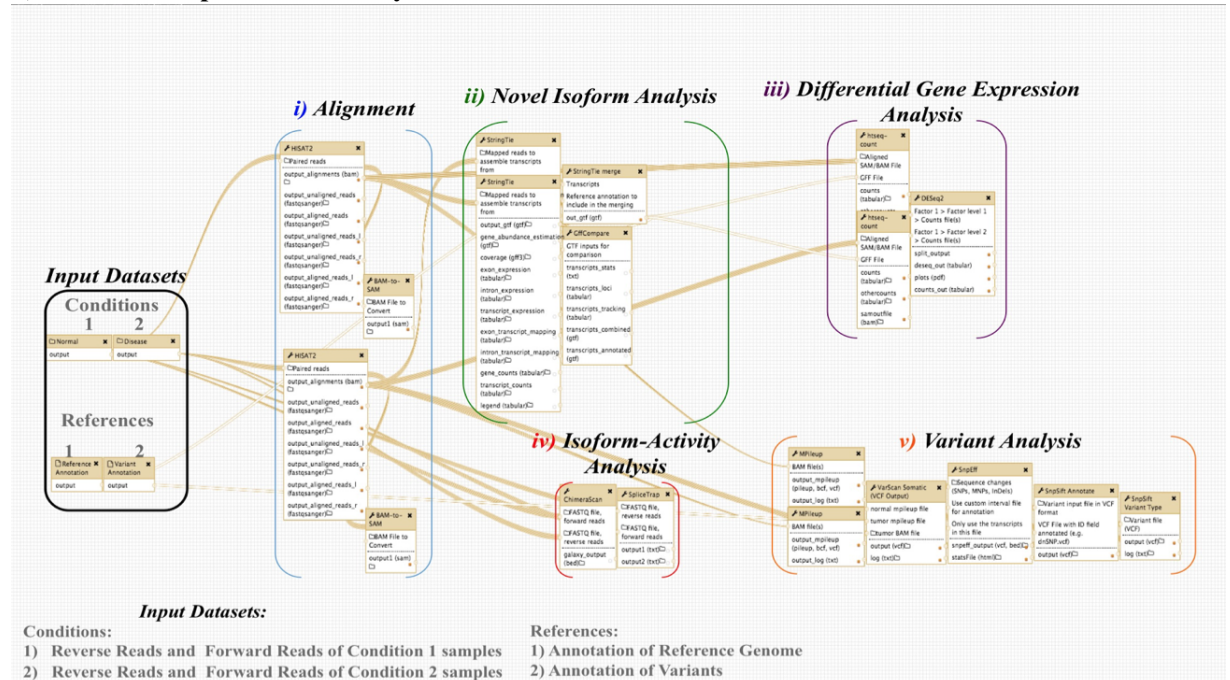
105 We document an example use case of TED with previously published prostate cancer transcriptome  
 106 data (13). We have developed a methodology that can provide the components of data analysis of complex  
 107 RNA-seq datasets through a toolkit interface that is easy to access, handle in addition to a comprehensive  
 108 data processing solution that is reusable and practical for users without extensive bioinformatics expertise.

**Figure 1. Overview of the Transcriptomics Profiler for Easy Discovery (TED) toolkit**



109

## b) TED Transcriptome Data Analysis



Transcriptome Data Analysis (Fig.1b) is the second module of TED comprised of five data analysis pipelines i) Alignment, ii) Novel Isoform, iii) Differential Gene Expression, iv) Isoform-activity and v) Variant analysis. This module consist of 14 bioinformatics tools and 24 steps that will analyze any number of paired-end RNA sequencing data samples from two conditions.

110 **METHODS**111 **Availability**

112 The TED toolkit is freely accessible on our local instance of the Galaxy platform via a url link:  
 113 [http://galaxy.hunter.cuny.edu/workflows/list\\_published](http://galaxy.hunter.cuny.edu/workflows/list_published) or through our custom Galaxy page: <http://galaxy.hunter.cuny.edu/profiler-for-easy-discovery-ted-toolkit>, that contains details of the RNAseq pipeline, datasets, and tutorials  
 114 of the transcriptome analysis as well as described in our documentation manual: <http://ted.readthedocs.io/en/latest/>.  
 115 A user can create an account (14) on our local Galaxy instance in order to have a private workflow  
 116 workspace, then import and run the pipelines directly from the URL links above. Furthermore, for each  
 117 new pipeline run, the results are saved in a separate Galaxy history (15) under the user's account, which  
 118 additionally offers a sharing option of the output through a simple web link. A virtual machine (VM) (16)  
 119 including Galaxy with the TED toolkit is also provided, with the tools and software dependencies prein-  
 120 stalled for download through the Data Libraries on our local Galaxy, under 'TED Virtual Machine (VM)  
 121 Application': (<http://galaxy.hunter.cuny.edu/library/list#folders/Fb56e686e7a485784>) and instructions to  
 122 set up and use the TED VM can be found in our documentation manual mentioned earlier.  
 123

124 **Data Source**

125 A total of 56 RNA-seq datasets were retrieved from the Array Express database of the European Bioin-  
 126 formatics Institute (<https://www.ebi.ac.uk/arrayexpress/experiments/E-MTAB-567/samples/>, EBI). The  
 127 files correspond to 14 sequenced transcriptomes from tumor tissue samples of prostate cancer human  
 128 patients and a technical replicate for each sample (total 28) in addition to 14 sequenced matched sam-  
 129 ples from the healthy tissue adjacent to the tumor tissue with replicates as well (additional 28). The  
 130 samples were collected, prepared and sequenced as described in the study by Ren et al (13). For each  
 131 tumor and healthy sample the dataset sequencing reads are paired-end, with replicates of each forward  
 132 and reverse sequencing read data files also included in the analysis. The EBI RNA-seq datasets are  
 133 also available for download through our local Galaxy Data Libraries, under 'TED toolkit Data Source':  
 134 (<http://galaxy.hunter.cuny.edu/library/list#folders/F862a7cb864998e85>) as well as other supporting data



135 such as the reference genome and reference annotation files.

### 136 **Implementation**

137 The TED toolkit was implemented on our local instance of the Galaxy platform: <http://galaxy.hunter.cuny.edu/>  
138 and freely accessible via a url link as mentioned in the 'Availability' section above. The TED pipelines con-  
139 sist of distinct bioinformatics software components and utilities, in which they were either downloaded and  
140 installed to our Galaxy instance via the public Galaxy toolshed (<https://toolshed.g2.bx.psu.edu/>), or man-  
141 ually integrated (17) in our local Galaxy toolshed in which all of the necessary custom tool scripts and wrap-  
142 pers are published as a repository in the main Galaxy toolshed (<https://toolshed.g2.bx.psu.edu/view/bioitcore/transcriptom>  
143 ) as well as in our public code repository on Github (<https://github.com/BCIL/TED>). All of the pipelines  
144 were assembled on Galaxy's workflow editor, by connecting the tools for the separate stages of the  
145 pipelines. In addition a virtual machine (VM) application was designed and build to include Galaxy  
146 with the TED pipelines, tools and software dependencies pre-installed for download and execution to the  
147 researcher's machine (<http://www.virtualbox.org>).

### 148 **RESULTS**

149 TED is packaged as a toolkit that integrates eleven distinct pipelines, on the Galaxy workflow canvas (18)  
150 and currently supports analysis of paired-end RNA-seq datasets from the Illumina sequencing platform of  
151 the human organism and analyzes the transcriptomes of RNA-seq datasets from two conditions which  
152 outputs are available in the Galaxy history (15) for the user to view and use. Within TED, the set of  
153 pipelines are divided into three fundamental modules based on their functionality that includes Data  
154 Quality Control (Fig. 1a), Transcriptome Data Analysis (Fig. 1b), and Data Discovery (Fig. 1c). All of  
155 TED pipelines are available as published workflows on our Galaxy server mentioned in the 'Availability'  
156 section and can be imported into the workspace of a public or private Galaxy instance by using the  
157 generated workflow links we provide, so that users have the option to run the analysis on their own server  
158 (19). We also provide the TED toolkit, the Galaxy server and all of the required software dependencies  
159 preconfigured as a virtual machine image (Fig. 1d). This is to allow the entire TED toolkit components  
160 and units to operate on any type of physical machine and operating system, by loading and powering  
161 up its appliance image into a virtual machine application. In order to demonstrate the effectiveness and  
162 convenience of our comprehensive analysis toolkit for RNA-seq, TED was used to gain insight into the  
163 molecular pathogenesis of 14 human prostate cancer transcriptomes. Using the TED toolkit we identified  
164 a range of differentially expressed coding, non-coding, novel isoform genes, gene fusions, alternative  
165 splicing events, genetic variants of somatic and germline mutations in these datasets. The following results  
166 below will first describe the Transcriptome Data Analysis module to explain how TED analyzes RNA-seq  
167 datasets and present part of our results for the differentially expressed coding genes and produced from  
168 the Data Discovery pipeline.

Patient	Pvalue	Number of Identified Differentially Expressed Genes		Pearson Correlation of Identified Differentially Expressed Genes			
		Total Genes	Upregulated Genes	Downregulated Genes	Total Regulated Genes	Upregulated Genes	Downregulated Genes
1	p ≤ 0.01	11008	5663	5345	0.93	0.92	0.96
	p ≤ 0.05	11961	6133	5828	0.93	0.92	0.95
2	p ≤ 0.01	10127	4865	5262	0.94	0.97	0.97
	p ≤ 0.05	11287	5421	5866	0.87	0.98	0.97
3	p ≤ 0.01	11513	5973	5540	0.93	0.92	0.96
	p ≤ 0.05	12526	6436	6090	0.94	0.85	0.96
4	p ≤ 0.01	11394	5676	5718	0.72	0.86	0.73
	p ≤ 0.05	12199	6090	6109	0.72	0.86	0.73
5	p ≤ 0.01	11352	5637	5715	0.85	0.95	0.89
	p ≤ 0.05	12184	6074	6110	0.85	0.95	0.89
6	p ≤ 0.01	11129	5533	5596	0.8	0.77	0.91
	p ≤ 0.05	12088	6025	6063	0.8	0.77	0.91
7	p ≤ 0.01	10954	5191	5763	0.97	0.93	0.99
	p ≤ 0.05	11868	5670	6198	0.97	0.93	0.99
8	p ≤ 0.01	11311	5839	5472	0.95	0.95	0.83
	p ≤ 0.05	12126	6220	5906	0.96	0.95	0.83
9	p ≤ 0.01	9910	4815	5095	0.34	0.96	0.35
	p ≤ 0.05	10999	5366	5633	0.35	0.96	0.35
10	p ≤ 0.01	5892	3289	2603	0.85	0.9	0.85
	p ≤ 0.05	6870	3794	3076	0.85	0.9	0.85
11	p ≤ 0.01	11557	5846	5711	0.85	0.93	0.86
	p ≤ 0.05	12510	6371	6139	0.85	0.93	0.86
12	p ≤ 0.01	9881	5121	4760	0.84	0.85	0.94
	p ≤ 0.05	11077	5677	5400	0.83	0.85	0.9
13	p ≤ 0.01	11372	5550	5822	0.85	0.87	0.9
	p ≤ 0.05	12278	6040	6238	0.85	0.87	0.9
14	p ≤ 0.01	9910	5095	4815	0.94	0.96	0.95
	p ≤ 0.05	11883	6271	5612	0.85	0.98	0.94

**Table 1.** Quantitative Summary of Differentially Expressed Genes and Pearson Correlation Statistics

### 169 Transcriptome Data Analysis

170 The Transcriptome Data Analysis (Fig. 1b) is the second module of TED comprised of five data analysis  
 171 pipelines i) Alignment, ii) Novel Isoform, iii) Differential Gene Expression, iv) Isoform-activity and v)  
 172 Variant analysis. This module consist of 14 bioinformatics tools and 24 steps that will analyze any number  
 173 of paired-end RNA sequencing data samples from two conditions.

174 The Alignment (Fig. 1b.i.) pipeline uses the UCSC hg38 reference genome (20), with the HISAT2 (21)  
 175 alignment program. HISAT2 uses an indexing technique to enable faster searches on the genome file,  
 176 which consist of a global index that covers the whole genome and many other small indexes for regions  
 177 that collectively cover the genome, to map whole reads entirely in the exons in which the Bowtie2 aligner  
 178 handles many of the operations required to construct and search the genome indexes. HISAT2 identifies  
 179 reads that span the exonic region as read alignments and the gaps between the spanning exonic regions as  
 180 junction signals. The output from this step is a Binary Alignment file (BAM, (22)), containing the mapped  
 181 exonic reads and their positions in the reference genome. To view the alignment file in text format, we a  
 182 BAM-SAM conversion step is included in the pipeline. All parameters were left as default, except the  
 183 minimum and maximum fragment length which are to be specified by the user that refers to the range  
 184 of the fragment size of the sequencing reads. This information can be found in the Read Information  
 185 pipeline that's part of the TED's first module Data Quality Control (not mentioned in this draft). Once the  
 186 alignment step is complete, the pipeline proceeds in four different analyses paths, the first for the novel  
 187 isoforms of the expressed genes, the second for differentially expressed genes (noncoding and coding),  
 188 the third for alternative splicing events and gene fusions and lastly the fourth for genetic variants of the  
 189 expressed genes.

190 The Novel Isoform analysis pipeline consists of 3 bioinformatics tools performing 4 steps (Fig. 1b.ii.),  
 191 with the Stringtie (23) software at its core, for reconstructing and quantifying the set of transcripts, and  
 192 number of gene isoforms, from the aligned transcriptome read data with the annotations of the reference  
 193 genome. The Stringtie assembler, uses as input the alignment file of mapped exonic reads produced from  
 194 HISAT2. The approach this software takes is, it builds an alternative splice graph from overlapping reads  
 195 in a given locus. This graph will contain nodes that corresponds to exons, and edges that corresponds to

196 reads which connects the exons. Stringtie will identify a path in the generated splice graph that has the  
197 largest number of reads on the edges (highest weight). This selected path will resemble an assembled  
198 transcript and because the edge weight equals to the number of reads, StringTie estimates the coverage  
199 level for this transcript that can be used to estimate the transcript's abundance, thus performs assembly and  
200 quantification simultaneously for every identified transcript. After the procedure of associating the reads  
201 with the assembled transcripts completes, they are then removed and the graph will update to perform the  
202 iteration of the algorithm on the next transcript. Stringtie will generate a separate transcriptome assembly  
203 for each of the HISAT alignment input files in which will then merged together using the Stringtie merge  
204 software. This is to combine redundant transcript structures across the transcriptome assemblies and  
205 identify which transcript structure corresponds to which annotated transcript using a reference annotation  
206 file, from the UCSC hg38 reference annotations (20) in Gene Transfer Format (GTF) (24). The reference  
207 annotation file contains information about known genes and transcripts that will be used to annotate the  
208 origin and nature of each transcript in the transcriptome assemblies. Furthermore in our pipeline, the  
209 gffcompare utility (25) was used to determine the number of assembled transcripts in comparison to  
210 known transcripts in the merged transcriptome assembly. The gffcompare tool will use the same reference  
211 annotation file used in the Stringtie merge step and evaluate the assembled transcripts that matched with  
212 the annotated genes either fully, partially and which ones entirely novel for isoform discovery that are not  
213 annotated.

214 The Differential Gene Expression Analysis pipeline consist of 2 bioinformatics tools and 3 steps (Fig.  
215 1b.iii.) to identify the transcripts that are differentially expressed between the two conditions of the  
216 RNA-seq experiment. This pipeline uses the htseq-count (26) tool from the HTSeq suite for counting the  
217 overlap of reads from the alignment files with annotation features, where each transcript is considered  
218 the union of all its exons. To count how many reads map to each transcript, the alignment files from  
219 HISAT2 are provided and the annotation file generated from the gffcompare tool which represents the  
220 transcripts present within the RNAseq datasets, as well as their location with non-redundant identifiers,  
221 and information regarding the origin. Then we will provide this information to DESeq2 (27) to generate  
222 normalized transcript counts (abundance estimates) and significance testing for differential expression.  
223 The Variant Analysis (Fig. 1b.iii.) includes five tools and performs six steps, with the SAMtools Mpileup  
224 program (28), VarScan Somatic tool (29) and SnpEff suite (30). The pipeline receives as input the  
225 BAM files produced from TopHat2 from the Data Groom and Alignment stage in addition to a reference  
226 human genome, in order for the SAMtools Mpileup to collect summary information the likelihood of  
227 each possible genotype is computed from the data and stored in a file for future reference, in addition  
228 to pileup of read base differences in a binary Variant Call Format (VCF) (31) file of each dataset. The  
229 variant caller tool VarScan Somatic reads the Mpileup output files and produces germline, somatic,  
230 and Loss of heterozygosity (LOH) events at positions where both normal and prostate cancer datasets  
231 have sufficient read coverage. All parameters were left as default for Mpileup and VarScan Somatic,  
232 except for p-value significance threshold set to 0.01 for VarScan Somatic, in order to enable a more  
233 sensitive first-pass algorithm in determining positional variants, that occurred in the supplied normal  
234 and prostate cancer mpileup's. SnpEff is a variant annotation and effect prediction tool. It annotates  
235 and predicts the effects of genetic variants. The Isoform Level Analysis (Fig. 1b.iv.) pipeline consists  
236 of two bioinformatics tools and performs 4 steps that detect chimeric transcripts encoded by a fusion  
237 gene performed by the Chimerascan (32) tool and quantifying alternative splicing events performed by  
238 the SpliceTrap tool (33). The ChimeraScan tool in this pipeline aligns paired-end reads to a combined  
239 genome-transcriptome reference, to identify potential fusion breakpoints from fragments that align to  
240 distinct references, or distance genomic locations of the same reference which are referred to putative  
241 chimeric junction sequences. The junction sequences are then used as reference to realign candidate  
242 junction-spanning reads. Several output files will be produced and the key output file is a tabular text file  
243 named chimeras.bedpe. SpliceTrap detects alternative splicing in paired-end RNA-seq data by using a  
244 Bayesian inference approach, by quantifying for every exon the extent to which it is included, skipped or  
245 subjected to size variations due to alternative 3'/5' splice sites or intron retention.

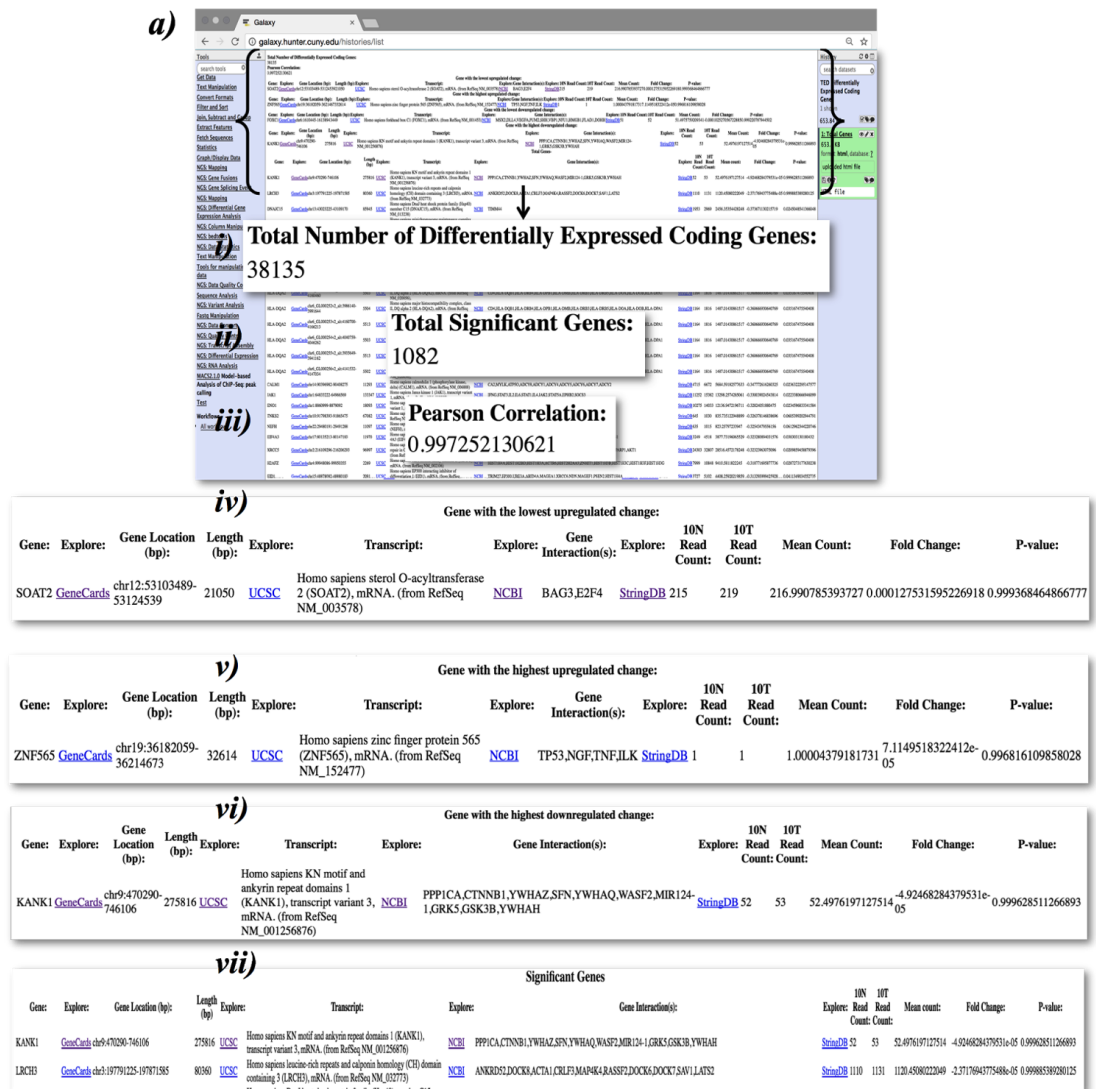
## 246 Data Discovery

247 The Data Discovery module consists of eleven pipelines and utilizes a highly structured approach for  
248 aggregating and summarizing the results produced from the transcriptome data analysis module for easy  
249 assessment, interpretation and downstream discovery. The eleven pipelines in this module generate

250 HyperText Markup Language (HTML) reports which can also be referred to as ‘actionable reports,’ that  
 251 transforms the data results into thorough, concise and intuitive information reports, consist of differential  
 252 coding gene expression (Fig. 1c.i.), differential non-coding gene expression (Fig. 1c.ii.), differential  
 253 novel gene isoform expression (Fig. 1c.iii.), somatic genetic variants (Fig. 1c.iv.), germ line genetic  
 254 variants (Fig. 1c.v.), comparison of genetic variants between samples (Fig. 1vi.), gene fusions (Fig.  
 255 1vii.), comparison of gene fusions between samples (Fig. 1c.viii.), gene splicing events (Fig. 1c.ix.), and  
 256 (Fig. 1c.x.) comparison of gene splicing events between samples. The following below will describe the  
 257 differential expression reports showing an example report of differential coding gene expression in Figure  
 258 2.

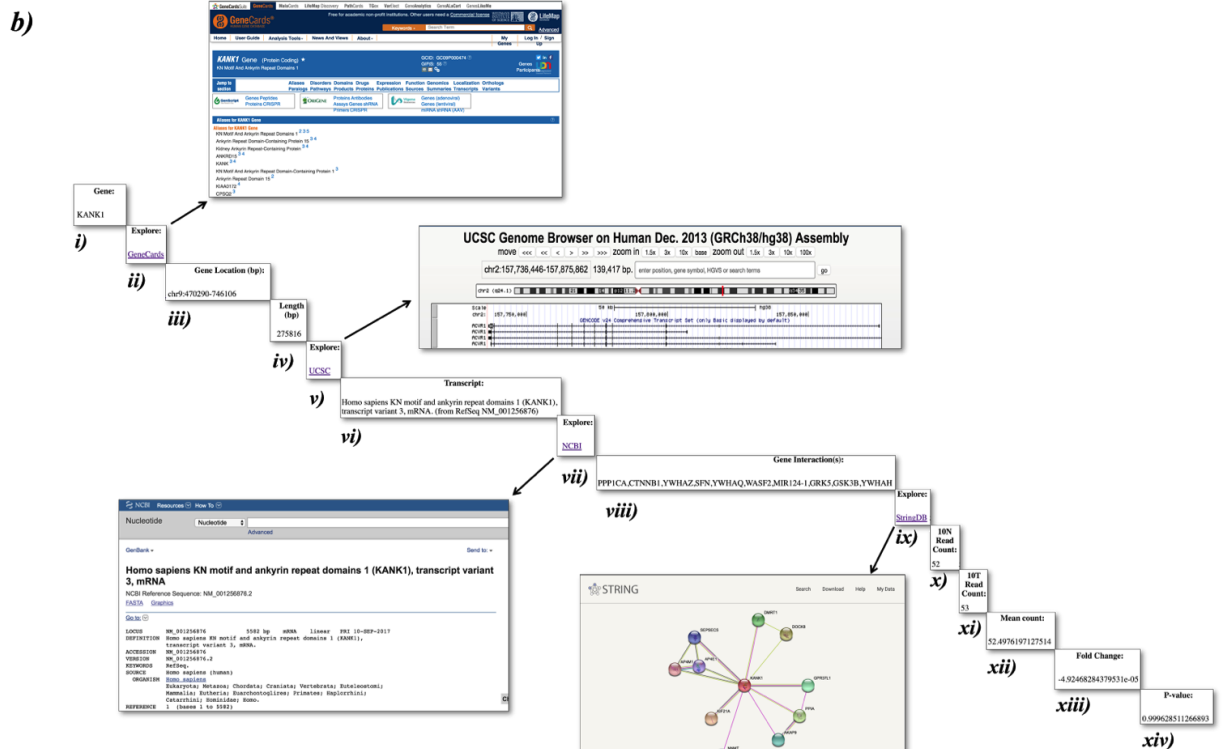
## 259 Differential Expression Report

**Figure 2. TED Data Discovery Analysis Report: Differentially Expressed Coding Genes**



260 For each of the pipelines generating differential coding gene, non-coding gene, and novel gene isoform  
 261 expression reports takes inputs the output data produced from the HTseq tool and Deseq2 tools of the  
 262 transcriptome data analysis module and generates three html reports Upregulated Genes, Downregulated  
 263 and Total regulated genes. Figure 2a illustrates an example of the total regulated html report for differential  
 264 coding gene expression of two sample RNAseq datasets (one sample for each condition), in which the





265 pipeline will populate the report with the following information in 8 parts: i) total number of differentially  
 266 expressed coding genes (Fig. 2a.i.), ii) total number of significant genes (Fig. 2a.ii.), iii) Pearson  
 267 correlation statistic between the genes (Fig. 2a.iii.), iv) the gene with the highest upregulated differential  
 268 fold change in FPKM (Fig. 2a.iv.), v) the gene with the lowest upregulated differential fold change in  
 269 FPKM (Fig. 2a.v.), vi) the gene with the lowest downregulated differential fold change in FPKM (Fig.  
 270 2a.vi.), vii) the gene with the highest differential fold change in FPKM (Fig. 2a.vii.), and the viii) list  
 271 of all the identified differentially expressed genes (Fig. 2a.viii.). The information of the gene with the  
 272 highest differential fold changes, the lowest differential fold changes, and the list of all the differentially  
 273 expressed genes are arranged in 14 columns (Fig. 2b). The 14 columns specify, i) the gene name (Fig.  
 274 2b.i), ii) link to GeneCards (34) a database of predicted human genes that provides concise genomic  
 275 related information, iii) the chromosomal location in which the gene resides, iv) the gene length, v) link to  
 276 UCSC Genome Browser (35) displaying the genomic location and other genomic data, vi) gene transcript  
 277 description name, vii) link to NCBI nucleotide database (36) providing gene and transcript data from  
 278 several sources, viii) list of genes involved in interaction with the differentially expressed gene ix) link to  
 279 StringDB (37) database of known and predicted gene pathway networks displaying direct and indirect  
 280 interactions x) read count of gene for sample 1 in normal condition xi) read count of gene for sample 1  
 281 for experimental condition xii) mean read count of the gene from both samples of both conditions xiii)  
 282 fold change of gene between the conditions and xiv) significant statistic in pvalue.

283 In our study of analyzing the 14 prostate cancer transcriptomes, we identified from the generated  
 284 reports, differentially expressed genes between paired prostate tumor and normal samples based on two  
 285 separate criteria's: pvalue  $\leq 0.05$  and pvalue  $\leq 0.01$  (Table 1). We analyzed at both pvalue cutoffs of  
 286  $\leq 0.01$  and  $\leq 0.05$  patients 1, 3, 8, 10, 11, 12, and 14, exactly half of the patients of the study group,  
 287 portrayed more up regulated genes expressed than down regulated genes at an exhibited correlation  
 288 coefficient of  $0.85 <$  for each of the upregulated genes and downregulated genes. This observation  
 289 emphasizes past findings in differential gene expression prostate cancer studies of the distribution of  
 290 upregulated genes being larger than downregulated genes (38, 39).

## 291 DISCUSSION

292 Compared to other RNA-seq and transcriptome analysis resources (6–9, 11, 12) that has its capabilities  
293 of reaching to a vast number of different scientific settings, the TED toolkit offers potential to reaching  
294 largely to a translational and diagnostic setting searching for a starting point of a preliminary overview of  
295 their RNAseq data that will lead to a discovery process with at most ease and minimal effort. In many  
296 clinical perspectives, RNA-seq delivers specific and sensitive genomic signatures but due to the lack of  
297 easy-to-use pipelines that can process in a transparent and streamlined fashion is limiting the expansion of  
298 RNA-seq from becoming a clinical diagnostic tool. Thus, the TED toolkit was intentionally designed as a  
299 Galaxy webserver since it allows inexperienced users to easily access advanced analysis tools processing  
300 the complex transcriptome analysis that will prepare unified outputs on a versatile workbench. Aside  
301 from TED being hosted on an accessible and intuitive system, it is also framed as a discovery platform  
302 that will structure the analysis results in html reports with analytical statistics and prioritization set with  
303 annotations and resource links to extremely comprehensive databases of disease and non-disease related  
304 information. This methodology offers a basic assessment of a RNA-seq study with initial details that  
305 are coherent as shown earlier with the differential gene expression results and can aid in the direction  
306 of a targeted discovery process to help come further to a conclusive clinical interpretation. Therefore,  
307 the TED toolkit holds strength to be a reliable, convenient and central protocol covering the majority  
308 aspects of transcriptome analytical results that is suitable to cater well within the reach of translational  
309 and diagnostic settings.

## 310 ACKNOWLEDGMENTS

311 The authors would like to thank all members of the Bioinformatics Core Infrastructures and Krampis Lab  
312 for their feedback during manuscript preparation.

## 313 REFERENCES

- 314 1. Kamalakaran,S., Varadan,V., Janevski,A., Banerjee,N., Tuck,D., McCombie,W.R., Dimitrova,N.  
315 and Harris,L.N. (2013) Translating next generation sequencing to practice: Opportunities and necessary  
316 steps. *Mol. Oncol.*, 7, 743–755.
- 317 2. Cummings,B.B., Marshall,J.L., Tukiainen,T., Lek,M., Donkervoort,S., Foley,A.R., Bolduc,V.,  
318 Waddell,L.B., Sandaradura,S.A., O’Grady,G.L., et al. (2017) Improving genetic diagnosis in Mendelian  
319 disease with transcriptome sequencing. *Sci. Transl. Med.*, 10.1126/scitranslmed.aal5209.
- 320 3. Christensen,S.M., Dillon,L.A.L., Carvalho,L.P., Passos,S., Novais,F.O., Hughitt,V.K., Beiting,D.P.,  
321 Carvalho,E.M., Scott,P., El-Sayed,N.M., et al. (2016) Meta-transcriptome Profiling of the Human-  
322 *Leishmania braziliensis* Cutaneous Lesion. *PLoS Negl. Trop. Dis.*, 10, e0004992.
- 323 4. Goecks,J., Nekrutenko,A. and Taylor,J. (2010) Galaxy: a comprehensive approach for supporting  
324 accessible, reproducible, and transparent computational research in the life sciences. *Genome Biol.*, 11,  
325 R86.
- 326 5. Galaxy Tool Shed [Galaxeast Wiki].
- 327 6. Schultheiss,S.J., Jean,G., Behr,J., Drewe,P., Görnitz,N., Kahles,A., Mudrakarta,P., Sreedharan,V.T.,  
328 Zeller,G. and Rättsch,G. (2011) Oqtans: a Galaxy-integrated workflow for quantitative transcriptome  
329 analysis from NGS Data. *BMC Bioinformatics*, 12, A7.
- 330 7. Goecks,J., El-Rayes,B.F., Maithel,S.K., Khoury,H.J., Taylor,J. and Rossi,M.R. (2015) Open  
331 pipelines for integrated tumor genome profiles reveal differences between pancreatic cancer tumors and  
332 cell lines. *Cancer Med.*, 4, 392–403.
- 333 8. Wolfien,M., Rimbach,C., Schmitz,U., Jung,J.J., Krebs,S., Steinhoff,G., David,R. and Wolken-  
334 hauer,O. (2016) TRAPLINE: a standardized and automated pipeline for RNA sequencing data analysis,  
335 evaluation and annotation. *BMC Bioinformatics*, 17, 21.
- 336 9. Lohse,M., Bolger,A.M., Nagel,A., Fernie,A.R., Lunn,J.E., Stitt,M. and Usadel,B. (2012) RobiNA:  
337 a user-friendly, integrated software solution for RNA-Seq-based transcriptomics. *Nucleic Acids Res.*, 40,  
338 W622-7.
- 339 10. Knowles,D.G., Roder,M., Merkel,A. and Guigo,R. (2013) Grape RNA-Seq analysis pipeline  
340 environment. *Bioinformatics*, 29, 614–621.

- 341 11. Kallio,M.A., Tuimala,J.T., Hupponen,T., Klemelä,P., Gentile,M., Scheinin,I., Koski,M., Käki,J.  
342 and Korpelainen,E.I. (2011) Chipster: user-friendly analysis software for microarray and other high-  
343 throughput data. *BMC Genomics*, 12, 507.
- 344 12. Howe,E.A., Sinha,R., Schlauch,D. and Quackenbush,J. (2011) RNA-Seq analysis in MeV. *Bioin-*  
345 *formatics*, 27, 3209–3210.
- 346 13. Ren,S., Peng,Z., Mao,J.-H., Yu,Y., Yin,C., Gao,X., Cui,Z., Zhang,J., Yi,K., Xu,W., et al. (2012)  
347 RNA-seq analysis of prostate cancer in the Chinese population identifies recurrent gene fusions, cancer-  
348 associated long noncoding RNAs and aberrant alternative splicings. *Cell Res.*, 22, 806–821.
- 349 14. Galaxy.
- 350 15. Histories.
- 351 16. Oracle VM VirtualBox.
- 352 17. Adding custom tools to Galaxy.
- 353 18. Data,B.S., Pipelines,A. and Biolinux,C. Galaxy Workflow Composition Canvas. 3.
- 354 19. ToolShed Workflow Sharing.
- 355 20. NCBI - WWW Error 404 Diagnostic.
- 356 21. Pertea,M., Kim,D., Pertea,G.M., Leek,J.T. and Salzberg,S.L. (2016) Transcript-level expression  
357 analysis of RNA-seq experiments with HISAT, StringTie and Ballgown. *Nat. Protoc.*, 11, 1650–1667.
- 358 22. Sam,T., Format,B.A.M. and Working,S. (2015) Sequence Alignment / Map Format Specification.  
359 SAM/BAM Format Specif. Work. Gr., 10.1016/j.ymeth.2012.07.021.
- 360 23. Pertea,M., Pertea,G.M., Antonescu,C.M., Chang,T.-C., Mendell,J.T. and Salzberg,S.L. (2015)  
361 StringTie enables improved reconstruction of a transcriptome from RNA-seq reads. *Nat. Biotechnol.*, 33,  
362 290–295.
- 363 24. GTF2.2: A Gene Annotation Format.
- 364 25. GFF utilities.
- 365 26. Anders,S., Pyl,P.T. and Huber,W. (2015) HTSeq—a Python framework to work with high-  
366 throughput sequencing data. *Bioinformatics*, 31, 166–169.
- 367 27. Love,M.I., Huber,W. and Anders,S. (2014) Moderated estimation of fold change and dispersion  
368 for RNA-seq data with DESeq2. *Genome Biol.*, 15, 550.
- 369 28. Li,H., Handsaker,B., Wysoker,A., Fennell,T., Ruan,J., Homer,N., Marth,G., Abecasis,G. and  
370 Durbin,R. (2009) The Sequence Alignment/Map format and SAMtools. *Bioinformatics*, 25, 2078–2079.
- 371 29. Koboldt,D.C., Zhang,Q., Larson,D.E., Shen,D., McLellan,M.D., Lin,L., Miller,C.A., Mardis,E.R.,  
372 Ding,L. and Wilson,R.K. (2012) VarScan 2: Somatic mutation and copy number alteration discovery in  
373 cancer by exome sequencing. *Genome Res.*, 22, 568–576.
- 374 30. Cingolani,P., Platts,A., Wang,L.L., Coon,M., Nguyen,T., Wang,L., Land,S.J., Lu,X. and Ruden,D.M.  
375 (2012) A program for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff:  
376 SNPs in the genome of *Drosophila melanogaster* strain w1118; iso-2; iso-3. *Fly (Austin)*, 6, 80–92.
- 377 31. Danecek,P., Auton,A., Abecasis,G., Albers,C.A., Banks,E., DePristo,M.A., Handsaker,R.E.,  
378 Lunter,G., Marth,G.T., Sherry,S.T., et al. (2011) The variant call format and VCFtools. *Bioinformatics*,  
379 27, 2156–2158.
- 380 32. Iyer,M.K., Chinnaiyan,A.M. and Maher,C.A. (2011) ChimeraScan: A tool for identifying chimeric  
381 transcription in sequencing data. *Bioinformatics*, 27, 2903–2904.
- 382 33. Wu,J., Akerman,M., Sun,S., McCombie,W.R., Krainer,A.R. and Zhang,M.Q. (2011) SpliceTrap:  
383 a method to quantify alternative splicing under single cellular conditions. *Bioinformatics*, 27, 3010–3016.
- 384 34. Rebhan,M., Chalifa-Caspi,V., Prilusky,J. and Lancet,D. (1997) GeneCards: integrating information  
385 about genes, proteins and diseases. *Trends Genet.*, 13, 163.
- 386 35. Kent,W.J., Sugnet,C.W., Furey,T.S., Roskin,K.M., Pringle,T.H., Zahler,A.M. and Haussler, a. D.  
387 (2002) The Human Genome Browser at UCSC. *Genome Res.*, 12, 996–1006.
- 388 36. NCBI Resource Coordinators (2017) Database Resources of the National Center for Biotechnology  
389 Information. *Nucleic Acids Res.*, 45, D12–D17.
- 390 37. Szklarczyk,D., Franceschini,A., Wyder,S., Forslund,K., Heller,D., Huerta-Cepas,J., Simonovic,M.,  
391 Roth,A., Santos,A., Tsafou,K.P., et al. (2015) STRING v10: protein-protein interaction networks,  
392 integrated over the tree of life. *Nucleic Acids Res.*, 43, D447–52.
- 393 38. Wei,Q., Li,M., Fu,X., Tang,R., Na,Y., Jiang,M. and Li,Y. (2007) Global analysis of differentially  
394 expressed genes in androgen-independent prostate cancer. *Prostate Cancer Prostatic Dis.*, 10, 167–174.

- <sup>395</sup> 39. Savli,H., Szendrői,A., Romics,I. and Nagy,B. (2008) Gene network and canonical pathway  
<sup>396</sup> analysis in prostate cancer: a microarray study. *Exp. Mol. Med.*, 40, 176.