

1 A Tool for the Comparison of Transcript Differential Expression
2 Analysis Pipelines

3
4 Stefano Beretta¹, Yuri Pirola¹, Valeria Ranzani², Grazisa Rossetti²
5 Raoul Bonnal², Raffaella Rizzi¹, Gianluca Della-Vedova¹,
6 Massimiliano Pagani², Paola Bonizzoni¹
7

8
9 1 Dipartimento di Informatica, Sistemistica e Comunicazione, Università
10 degli Studi di Milano - Bicocca, Milano, Italy
11 2 Istituto Nazionale di Genetica Molecolare (INGM), Milano, Italy
12

13
14 Corresponding Author:
15 Stefano Beretta¹
16 Email address: stefano.beretta@disco.unimib.it

17 INTRODUCTION

18 Long non-coding RNAs (lncRNAs) have recently gained interest, especially for their involvement in
19 controlling several cell processes, but a full understanding of their role is lacking. Differential
20 Expression (DE) analysis is one of the most important tasks in the analysis of RNA-seq data, since it
21 potentially points out genes involved in the regulation of the condition under study.

22 However, a classical analysis at gene level may disregard the role of Alternative Splicing (AS) in
23 regulating cell conditions. This is the case, for example, when a given gene is expressed in all the
24 different conditions, but the expressed isoform is significantly diverse in the different conditions
25 (that is an isoform switch).

26 A transcript level analysis may better shed light on this case, especially in studies having as goal, for
27 example, a better understanding of the behavior of lncRNAs in lymphocytes T cells, which are
28 fundamental in studies of specific diseases, such as cancer.

29 After Cufflinks/Cuffdiff, several approaches for DE analysis at isoform/transcript level have been
30 proposed. However, their results are often sensitive to the upstream analysis such as read mapping,
31 transcript reconstruction and quantification, and it is often hard to choose "a priori" the most
32 appropriate combination of tools.

33 This work presents a tool for assisting the user in this choice, and poses the bases for a study devoted
34 to the characterization of lncRNAs and the identification of of isoform switch events. Our tool
35 includes a framework for the description and the execution of a set of DE pipelines over the same
36 input dataset, as well a set of tools for reconciling and comparing the results.

37

38 METHOD

39 We designed an automated and easily customizable tool which is able to execute a set of existing
40 pipelines for DE analysis at transcript level starting from RNA-seq data. Our method is built upon
41 Snakemake, a workflow management system, with the specific goal of reducing the complexity of
42 creating workflows. This approach guarantees that the experimentation is fully replicable and easy
43 to customize. Each considered pipeline is structured in three steps: (i) transcript assembly, (ii)
44 quantification, and (iii) DE analysis. By default, our tool builds and compares 9 different pipelines,
45 each taking as input the same set of RNA-seq reads, obtained by combining different state-of-the-
46 art methods to perform the transcript assembly (TA step) with different state-of-the-art methods
47 to perform quantification and differential expression analysis (Q+DE step). More precisely, the 9
48 pipelines are obtained by combining two tools (Cufflinks and StringTie) and a Reference Annotation
49 (Ensembl annotated transcripts) for the TA step, with three tools (Cuffquant+Cuffdiff, StringTie-
50 B+Ballgown and Kallisto+Sleuth) for the Q+DE step. Each pipeline produces for each transcript a p-
51 value, giving an evaluation of the statistical significance of its expression variation among the
52 different conditions (opposed to the null hypothesis of a random variation).

53

54 RESULTS

55 We have tested our tool on 15 datasets of RNA-seq reads consisting of 3 individuals sequenced
56 under 5 different conditions, as a starting point in the characterization of specific lncRNAs. The
57 datasets have been produced by an Illumina HiScanSQ sequencer: each dataset contains on average
58 23.5 million paired-end sequences spanning the entire genome.

59 We have computed the correlation between the two sets of p-values for each pair of pipelines,
60 observing that that the correlation coefficients are larger for some pairs of pipelines using the same
61 approach for the Q+DE step. More precisely, the couples using Cuffquant+Cuffdiff have correlations
62 between 0.86 and 0.89, while those employing StringTie-B+Ballgown have correlations between
63 0.83 and 0.85.

64 The correlation coefficients of all other pairs of pipelines (included those using Kallisto+Sleuth) are
65 smaller than 0.4 (hence much less significant). A likely explanation is that the choice of the Q+DE
66 tools crucially influences the final results, and is more important than the choice of the tool for the
67 TA step. Still, we plan to perform an in-depth analysis of this phenomenon.
68 Moreover, our experiments have confirmed that the datasets contain two specific differentially
69 expressed isoforms of the gene PTPRC, which is known in literature to have a switch event between
70 those isoforms. We have also confirmed other transcripts which are compatible with annotated
71 lncRNAs. A further work is to develop a better method to compute the correlation of the transcripts
72 assembled by the different pipelines, exploiting their predicted intron-exon structure to compute
73 the comparison, and introducing an ad-hoc and robust method to estimate the correlation
74 coefficients. Finally, a future development is to amalgamate the outputs obtained by the different
75 pipelines to produce more reliable predictions.