

"RNA-Seq: from computational challenges to biological insights" Valerio Costa, PhD Laboratory of Human Genetics Diseases at the Institute of Genetics and Biophysics "A.Buzzati-Traverso" IGB-CNR, Naples Joint NETTAB 2010 and BBCC 2010 workshops, Naples, Italy





The "next-generation" sequencing era

RNA-Seq allows to:

Characterize organisms' full set of genes

- Detect and quantifiy expression from known genes;

- Find both new coding and non-coding genes;

- Compare genes among organisms (evolution of genomes); Characterize transcript isoforms- Identify and quantify known splice events - Find novel alternative splice isoforms and/or transcript ends (5'-3' UTRs);

Monitor gene expression changes between cells/tissues/organisms or conditions-Identify differential expression between 2 conditions;- Understand the basis of gene expression regulation in a disease;- Identify gene regulatory regions (e.g. coupled with ChIP-Seq);

RNA-Seq and microarrays



Hybridization-basedtechnologies:Background and cross-hybridization issuesOnly transcripts included in the array designSpecific studies requires specific array typesLimited dynamic range

Nowadays much easier to analyze (several software available)- Nowadays still cheaper
 "large" sample production
 Low computational complexity

RNA-Seq:- Low "background signal"-Identification of novel transcribed regions and splice isoforms;- Determination of correct gene boundaries- No upper limit for gene quantification

 Still expensive (sample preparation and sequencing)

 Much more computationally demanding- Still limited amount of software available

RNA source	Characteristics
PolyA RNA	 Polyadenylated fraction of the transcriptome Represents 2 – 5% of all transcripts Coding mRNA only
rRNA-depleted RNA	 Coding & non coding RNA minus the abundant rRNA
Total RNA	 See entire transcriptome Includes abundant rRNAs







Mapping strategy

.csfasta

>852_2042_1999_F3T3201120112302220133211010201103113 2013023321002303

<u>qual</u>

>852_2042_1999_F319 20 14 14 8 5 9 16 11 11 6 14 21 14 11 21 -1 20 11 21 12 22 14 18 14 6 11 16 14 16 5 11 23 13 18 4 6 20 13 15 21 17 18 15 11 4 8 7 5 11

A suitable treatment of the multiple matched reads is fundamental to reduce the bias.



- Quality assessment and filters (quality plot, remove low quality reads, ribosomal RNA reads, sequencing adapters);
- 2. Alignment to a reference genome (genome+junction library)
- "Trim" the rigth-side of the reads and cyclically repeats the step;
- 4. Handle "multiple" reads;



Since the huge number - and the short size - of reads (50 nt in length), using conventional alignment algorithms is not feasible; In addition, not all developed aligners support all .**csfasta** formats from SOLiD platform;

The alignment of reads in RNA-seg is particularly challenging due to the reads spanning across splice-iunctions













- "known" regions (i.e. RefSeq, UCSC annotated genes, transcripts);
- novel transcriptionally active regions (TARs);
- gene boundaries (5' and 3' UTRs analysis);

2) Identification and quantification of <u>splicing isoforms</u>: - "known" transcript isoforms;- detection of new alternative splice isoforms and their quantification

3) Analysis of <u>non-coding RNAs</u> (ncRNAs):

- detection and quantification of "known" ncRNAs;

" " of new ncRNAs and their quantification;

4) Detection of <u>differentially expressed (DE) "features"</u>:
detection of DE "known" genes;
" of DE newly identified genes;- " " of sample/condition-specific isoforms;- " " of DE alternative splicing isoforms
" " of DE nornAs;

Within sample analysis

Between samples analysis

- "known" regions (i.e. Refseq, UCSC annotated genes, transcripts);
- novel transcriptionally active regions (TARs);
- gene boundaries (5' and 3' UTRS analysis);



 10^{3}

106

S.

Quantification based on RefSeq Annotation: Remove ambiguities due to genes overlapping by strand: Use either "exon reads" and "junction reads"; Use unique reads + "uniquely assigned" reads after the "rescue" step.

Expression was measured as the Number of Reads Mapped on the feature i or as **Reads Per Kilobase of transcript per Million of mapped reads (RPKM)**

N = Total number of mapped read

 L_i = Length (bp) of the feature i



Analysis of "extra-genic" transcription



Analysís of "extra-geníc" transcription





- "known" regions (i.e. RefSeq, UCSC annotated genes, transcripts); - novel transcriptionally active regions (TARs);

- gene boundaries (5' and 3' UTRS analysis);

2) Identification and quantification of <u>splicing isoforms</u>: - "known" transcript isoforms;- detection of new alternative splice isoforms and their quantification (in progress)

3) Analysis of <u>non-coding RNAs</u> (ncRNAs):

- detection and quantification of "known" ncRNAs;

" " of new ncRNAs and their quantification;

4) Detection of <u>differentially expressed (DE) "features"</u>:
detection of DE "known" genes;
" of DE newly identified genes;- " " of sample/condition-specific isoforms;- " " of DE alternative splicing isoforms
" " of DE nCRNAS;







- "known" regions (i.e. RefSeq, UCSC annotated genes, transcripts); - novel transcriptionally active regions (TARs);
- gene boundaries (5' and 3' UTRS analysis);

2) Identification and quantification of <u>splicing isoforms</u>: - "known" transcript isoforms;- detection of new alternative splice isoforms and their quantification

3) Analysis of <u>non-coding RNAs</u> (nCRNAs):

- detection and quantification of "known" nCRNAS;

" " of new normas and their quantification (in progress)

4) Detection of <u>differentially expressed (DE) "features"</u>:
- detection of DE "known" genes;
- " of DE newly identified genes;- " " of sample/condition-specific isoforms;- " " of DE alternative splicing isoforms
- " " of DE nCRNAS;





- "known" regions (i.e. RefSeq, UCSC annotated genes, transcripts); - novel transcriptionally active regions (TARs);

- gene boundaries (5' and 3' UTRS analysis);

2) Identification and quantification of <u>splicing isoforms</u>: - "known" transcript isoforms;- detection of new alternative splice isoforms and their quantification

3) Analysis of <u>non-coding RNAs</u> (ncRNAs): - detection and quantification of "known" ncRNAs;

" " of new ncRNAs and their quantification;

4) Detection of <u>differentially expressed (DE) "features"</u>: - detection of DE "known" genes;

 " " of DE newly identified genes;- " " of sample/conditionspecific isoforms (in progress);- " " of DE alternative splicing isoforms (in progress);

" " of DE NCRNAS;

Statistical analysis: Differential Expression

Significant changes in the expression of genes are usually identified by using a **statistical Test** and the results are then corrected for **multiple testing**

Unfortunately one cannot use ordinary tests developed for microarray since RNA-Seq data are count data, and they are heteroscedastic (have no the same finite variance).

Statistical significance has been inferred from total reads count for each RefSeq gene combining 3 tests:
 DEGseq (based on Poisson distribution)
 DESeq and edgeR (based on negative binomial).

- Such tests are based on slightly different assumptions that usually produce a different level of stringency.









Dífferential expression of snoRNAs

46 *SNORD* (3 up- and 43 down-regulated) on a total of 171 expressed (27%); 31 *SNORA* (9 up- and 22 down-regulated) on a total of 95 expressed (32,6%); 9 *SCARNA* (2 up- and 7 down-regulated) on a total of 23 expressed (39%);





Future perspectives (1/2)

- RNA-Seq experiments are a powerful tool for addressing biological questions, although they still require the setup of "sophisticated" computational methods and the development of novel computational/statistical tools;

To develop a probabilistic model which takes into account the uncertainty due to the mapping

To build appropriate gene models to better define & quantify the high level of transcription within yet unannotated extra-genic regions

To reconstruct, and thus further quantifying, multiple isoforms of a transcript (isoform abundance).

TopHat aligns reads to the genomes using <u>**Bowtie</u></u> - an ultrafast short reads aligner - and then analyzes the mapping results to identify splice junctions between exons (both known & newly identified).</u>**

Cufflinks assembles transcripts, estimates their abundances, and tests DE and regulation.

Future perspectives (2/2)

Biological conclusions inferred from the direct comparisons of two samples are however limited;

RNA-Seq experiments can (optimistically) reduce the technical variability, but they do not affect the biological variability.



From statistical significance to biological significance Extend the analysis to a larger number of samples/conditions to increase the detection power for identifying disease-associated genes/features

 Athough our results are very promising, all the capability and information have not been fully extracted from data;

 Further steps of "biological validations" (Real-Time PCR, WesternBlot, RNA interference, etc.) are also required;



Acknowledgements

Experimental design & Sample preparation



Human genetics disease Lab Alfredo Ciccodicola Marianna Aprile Roberta Esposito

SOLiD Sequencing Core Facility Stefania Crispi Luigi Leone Aldo Donizetti Patients' enrolling & Sample preparation Claudio Napoli Raffaele Calabrò Berardo Sarubbi Linda Sommese Amelia Casamassimi



Paola Salvatore

Data analysis



Margherita Mutarelli



Claudia Angelini

Luciana D'Apice

Piergiuseppe De Berardinis



Data validation