Constructing Protein-RNA Interactome using RIP-Seq

. Introduction



II. Motivation:

RIP-Seq technology provides an unprecedented opportunity for identifying genome and transcriptome-wide protein-RNA interactions.

III. Challenges:

- Accurately binning the sample mixtures with short barcode sequences.
- Effectively filtering reads mapped to noisy (unknown) background RNAs.
- High sequence duplication and over-representation of k-mers.
- Low mappability to the corresponding genomes and/or transcriptomes.
- Lack of sequence annotations for ncRNA molecules.
- Effectively examining the mapped results.
- Inferring abundance of mapped isoforms with statistical significance.

Yue Li^{1,2}, Dorothy Yanling Zhao^{2,3}, Jack Greenblatt^{2,4} and Zhaolei Zhang^{2,3,4}

¹Department of Computer Science, ²Donnelly Centre for Cellular and Biomolecular Research, ³Department of Molecular Genetics, ⁴Banting and Best Department of Medical Research, University of Toronto

IV. Material



Origin	IP Sample	Adaptor 1 (black) + 4 nucleo
human	CBX2	CTTTCCCTACACGACGC
human	CBX3	CTTTCCCTACACGACGC
•••		
mouse	STK31	CTTTCCCTACACGACGC
mouse	TDRKH	CTTTCCCTACACGACGC

V. Computational Methods



- 1. **Binning** based on 4-nucleotide barcode sequence
- 2. Converting to FASTQ to take into account both read sequence and Phred quality score
- 3. Quality assessment with FASTQC
- 4. Filtering RNA background sequences with *Bowite* aligner: (a) primer/adaptor
- (b) Spike-in virus/bacterial DNA genome (e.g., phiX genome)
- (c) Prevalent ribosomal, mitochondria RNA, actin mRNA, and more to discover
- 5. Mapping against reference genomes with *Bowite* aligner:
- (a) human (hg19) or mouse (mm9) genome
- (b) human or mouse transcriptome (Ensembl cDNA databases)
- (c) human or mouse ncRNA (Ensembl ncRNA databases)
- (d) 7sk and U6 RNA (known positive controls for human CCNT and human CCNT UV) (e) kcnq1ot1 RNA (known positive control for human G9a and human G9a UV)
- 7. Quantification and localization of reads (an effective approach is under development)
- 6. Visualization of mapping results with *Integrative Genome Viewer* (IGV)
- 8. Transcript (de novo) assembly with Cufflinks

VI. Preliminary Results

Per base Phred quality score (left) and sequence duplication levels (right) for GFP UV sample:



PCR

adaptor

ligation

Double strand

Cloonan, et al., Nat Method, 2008

eotide barcode (red) **CTCTTCCGATCTATCC CTCTTCCGATCTTAGC**

CTCTTCCGATCTGACA CTCTTCCGATCTTGAC





- Many unknown background RNA molecules were discovered.

VII. Works in Progress

VIII. Acknowledgment



Donnelly Centre

UNIVERSITY OF TORONTO

Filtering and mapping results:

• As expected, 105 reads from CCNT and < 5 reads from other samples map to 7sk RNA. • A top list of read-enriched cDNA and ncRNA molecules were automatically generated. • Some of them show distinctive enrichments in subsets of the 20 samples. For instance, - PKD2L1 (cDNA on '-' strand of chr10:102,047,903-090,243) only in CCNT and PIWIL4. - AC008541.1 (ncRNA on '+' strand of chr5:122,990,602-1,011) only in the 3 CBX's.

• Filter out highly repeated sequences and newly discovered background RNA for remapping. • Integrate previously published data with our data to distinguish noise and meaningful peaks. • Correlate reads in HEK293 cell lines and mouse testes to gene expression and chromatin state. • Explore frameworks (e.g., *cufflinks*) for (*de novo*) transcript assembly + infer abundances. • Investigate statistical frameworks for significance of read counts *specifically* for RIP-Seq.

