**International Conference on**

# Omics Studies

**September 04-05, 2013   Holiday Inn Orlando International Airport, Orlando, FL, USA**

## The effect of experimental design and analysis tool parameter optimization in RNA-Seq data analysis

**Natalie Abrams, Robert Stephens, Jack Collins, Yvonne Edwards and Eric Stahlberg**
National Institutes of Health, USA

Since its inception, RNA-Seq technology has been transforming cancer research by providing an increasingly detailed, exact and quantitative view of the mammalian transcriptome. In cancer research, it facilitated the identification of multiple cancers subtypes, genetic biomarkers of tumorigenesis and drug response, and somatic driver mutations in cancer-associated genes. Fueled by continuous advances in NGS technologies, the quest continues to pursue new biomarkers, drug targets, genetic susceptibility genes, and causal mechanisms of less common cancers. Despite this remarkable progress, three fundamental obstacles still limit the full potential of RNA-Seq in cancer research: sample quality and quantity, short read length, and poorly defined experimental design principles and quality metrics. To minimize these problems, increasingly sophisticated statistical approaches and computational algorithms have been developed that facilitate analysis, visualization and integration of RNA-Seq data. These state-of-the-art tools can perform increasingly complex queries, which go far beyond simply enumerating transcripts or somatic mutations in tumor samples. Yes the development efforts including quality control software often lag behind technological advances for many NGS applications, so more statistically robust tools are needed to efficiently transform this data into biologically meaningful information. Furthermore, despite the rapid evolution of NGS instruments and data analysis algorithms, there is still a dearth of published studies that focus on experimental design principles and quality control metrics. Experimental design strategies for cancer studies need to be sufficiently robust to match the precision offered by sequencing instruments. In this work, we aim to identify and explore the best experimental design principles in the context of RNA-Seq mammalian studies. To this end, we have examined the effect of experimental design and analysis tool parameter optimization in common software packages including Cufflinks and EdgeR on the quality of RNA-Seq results.

natalie.abrams@nih.gov

J Proteomics Bioinform 2013
ISSN: 0974-276X, JPB an open access journal

Omics Studies-2013
September 04-05, 2013

Volume 6 Issue 8

Page 30