

Q: *How have your customers responded to moving to mRNA-Seq for transcriptome studies?*

SM: I think there were several reasons why our clients were initially slow to evaluate mRNA-Seq. Number one was cost. Until very recently, the cost of mRNA-Seq was still significantly higher than the cost of running expression profiling on microarrays. Illumina solved that with the introduction of HiSeq 2000 and its TruSeq™ v3 reagents. Since most of our customers were more familiar with microarrays, we also felt that many of them were waiting to see more gene expression studies performed with mRNA-Seq.

Q: *Is that why you conducted an mRNA-Seq vs. array comparison study?*

SM: We wanted to help our clients understand the benefits of mRNA-Seq over microarrays for transcriptome analysis, so we decided to conduct real-world experiments and give our clients access to these data sets so they could begin to integrate them into their research. We felt the data would be persuasive in showing how mRNA-Seq could provide additional information about the biology of what they were studying.

WJ: Instead of reference samples, we used samples with biological diversity. We simulated a real biological experiment that people send us every week, performed the analysis on microarrays and the HiSeq 2000 system, and compared the results.

Joel Parker (JM): The design of the study involved 15 breast cancer cell lines, with five unique lines representing each of three cancer subtypes. We created three independent library preparations of each line using the Illumina TruSeq protocol. Forty-five samples were multiplexed into seven pools and run on two HiSeq flow cells. The 15 samples were also assayed on Affymetrix and Illumina microarrays and we compared the performance based on repeatability, sensitivity, specificity of detection, abundance estimation, fold-change, and differential expression. The results demonstrated the uniformity and quality of the data produced with mRNA-Seq. Replicates run at different times demonstrated a negligible decrease in reproducibility as opposed to microarrays, where batch effects due to technical variation, sample handling, and process variances often negatively impact reproducibility.

mRNA-Seq detected more of the content specific to Affymetrix and Illumina microarrays than either of the microarray platforms on the same samples. It also detected genes, isoforms, and differential expression not detected by either one, with approximately one-third of the SNVs detected by mRNA-Seq not known to dbSNP, the public domain archive of all established genetic variation. We found that 25M reads are necessary for repeatability equivalent to or better than the microarray platforms, while 10M may be sufficient for equivalent or improved detection and differential expression.

SM: We generated a pretty compelling data set, initially presenting it in a September webinar¹. It's been a popular presentation, with over 200 downloads and more than 125 people viewing the webinar. We're now sharing this data with customers.

Q: *Were you surprised by these results?*

WJ: We expected the sequencer to perform better than arrays on several different fronts and were happy to see that confirmed. While we weren't surprised that we would obtain more information

from sequencing, we were amazed when the higher resolution of mRNA-Seq showed events occurring that we had not anticipated.

JP: mRNA-Seq identified variations that we couldn't see in the array data. In one cell line, we were able to confirm a known heterozygous mutation in GATA3, a transcription factor involved in growth control and the maintenance of differentiated epithelial cells. Research has shown that GATA3 variants may contribute to tumorigenesis in certain breast tumors².

Unlike microarrays, mRNA-Seq enabled us to see novel expression, specifically isoform switching between two breast cancer stages. There are instances where a gene's isoform expression increases in one cell subtype and decreases in another, while another isoform of the same gene behaves oppositely. For example, CD44 is a cell surface protein that modulates cellular signaling in differentiation, growth, and apoptosis and has been linked with metastasis. The mRNA-Seq data clearly showed increased luminal expression of one variant of CD44 and decreased claudin expression, while another variant of the same gene was just the opposite. When we first saw the switch, we thought it was interesting. Later, we found that another researcher had recently independently published documentation of the same event, and found that it was essential for epithelial-mesenchymal transition and breast cancer progression³.

“While we weren't surprised that we would obtain more information from sequencing, we were amazed when the higher resolution of mRNA-Seq showed events occurring that we had not anticipated.”

We also identified EGFR isoform switching in our samples. While this growth factor has been found to undergo isoform switching in colorectal cancer, this event had not been seen in breast cancer. Further research is needed to understand its potential role in the progression of the disease.

Q: *What other hurdles are there for clients wanting to switch from gene expression microarrays to mRNA-Seq?*

SM: Mainly the fact that computational and analytical tools have not been readily available to assist researchers in analyzing the sequencing data. Most of the gene expression informatics infrastructure has been built around handling array-based data sets, so customer pipelines are not yet in shape to handle RNA-Seq data sets. Up until now, they've also had no way to compare their new mRNA-Seq data with their old microarray data.

Q: *Part of your mRNA-Seq offering is a bioinformatics output that mimics array reports. Can you talk more about your deliverables and the software you developed to do that?*

WJ: We recognized that we'd need to help some of our clients overcome the inertia of switching platforms, especially their concerns about wanting to compare their new data with their old data. So we offer the typical sequencing output you'd expect, with reads



