Streamlining NGS Data Management & Analysis

Integrated mRNA and microRNA transcriptome analysis in Strand NGS

Veena Hedatale and Rohit Gupta Strand Life Sciences

www.strand-ngs.com

Analyze Visualize Annotate Discover



Application Note



Integrated mRNA and microRNA transcriptome analysis in Strand NGS

Veena Hedatale and Rohit Gupta*

Strand Life Sciences March 2015

* Corresponding Author (rohit@strandls.com)

igh throughput next-generation sequencing has made possible the analysis of mRNA and miRNA expression within the same cellular context. Strand NGS provides an efficient and powerful way to quickly and accurately analyse differential expression and infer the network of miRNA-mRNA interactions in an integrated manner. Using a case study, this paper highlights the application of Strand NGS in the analysis of nasopharyngeal carcinoma transcriptome.

1 Introduction

MicroRNAs are endogenous non-coding RNAs that act by negatively regulating mRNA expression at the post-transcriptional level. The use of high throughput next-generation sequencing has made possible the analysis of mRNA and miRNA expression within the same cellular context. Generally, messenger RNA and microRNA often have an inverse expression with multiple miRNA targets interactions involved [4]. Altered miRNA expression is increasingly identified with major functions in tumor pathogenesis. Since multiple miRNAs may coordinate to regulate targets in a common pathway, analysis of pathways is more important in cancer biology compared to characterization of individual target genes. Furthermore, joint miRNA-mRNA expression profiles, miRNA-target mRNA relationships, and the construction of regulatory networks can provide new insights into complex biological processes.

In this study, we showcase the capabilities for such an integrated transcriptome analysis followed by mRNAmicroRNA regulatory network analysis in Strand NGS, using data published by [5] in nasopharyngeal carcinoma (NPC) model systems.

2 Data Set

The study by [5] characterizes the mRNA and miRNA transcriptome in NPC models to understand transcript regulation in nasopharyngeal carcinomas. The corresponding data set, GSE54174, was downloaded from GEO database. Sequencing was originally done using the Solexa sequencing technology on the Illumina Genome Analyzer IIx to generate 58 bp single end reads. Samples include matched total mRNA and small RNA from HK1 a well differentiated NPC cell line, C666 and X666 two undifferentiated NPC cell lines and NP460 a normal nasal epithelial cell line (table 1).

Cell line	mRNA		smallRNA		
	Sample	Total	Sample	Total	
		Reads		Reads	
NP460	SRR1178332	26,626,486	SRR1121963	26,752,260	
(Control)					
HK1 (Well	SRR1178331	32,699,704	SRR1121962	26,696,805	
differentiated					
EBV-)					
C666 (Undiffer-	SRR1178329	27,076,636	SRR1121960	29,160,516	
entiated EBV+)					
X666 (Undiffer-	SRR1178330	34,508,151	SRR1121961	27,642,905	
entiated EBV+)					

 Table 1: Description of data sets

3 Methods

The overall data analysis workflow for expression and sequence analysis of the RNA-seq and small-RNA samples is outlined in figure 1.

3.1 Alignment of RNA-Seq and small RNA data

The raw reads corresponding to RNA-seq were aligned in Strand NGS against the human genome (hg19) reference and the UCSC transcript model by running the alignment against transcriptome and genome together with novel splices. The raw RNA reads were aligned with minimum



Figure 1: Overall approach for the analysis of RNA-Seq and small RNA data

of 90% identity; maximum of 5% gaps and 25 bp as the minimum aligned read length. The small RNA reads were aligned separately in a small RNA alignment experiment with a mismatch cutoff of 2 and minimum read length of 10 bp. Reads that had an average quality ≤ 10 at the 3' end were trimmed prior to alignment.

3.2 Identifying differentially expressed genes

Samples were grouped based on cell lines with pair-wise interpretations setup for C666 vs. NP460 and HK1 vs. NP460 (cancer vs. normal). Quantification was done to obtain the read densities equivalent to RPKMs as reported in the original study. Differential expression analysis of mRNA and miRNA was performed using the DESeq script [2] via the script editor in Strand NGS with a p-value ≤ 0.05 considered as significant. The original study by [5] is based on the non-adjusted p-values. The adjusted p-value is usually considered significant but we have not applied any correction to be able to compare the results. A fold change ≥ 2.0 was used to define entities as being up- and down- regulated.

3.3 miRNA target prediction

Targets of HK1 and C666 miRNAs were predicted from PITA and TargetScan databases with "Find target genes" using a p-value cut-off 0.05.

3.4 Sequence variant analysis

Strand NGS variant caller is used to call SNPs and short InDels. SNPs were identified in NP460, C666 and HK1 mRNA samples using dbSNP135 annotation and a confidence score >50 and a base quality >10. Please note that we have used dbSNP135 database to compare our results with the original paper, however the latest dbSNP annotation is also available in Strand NGS. Reference locations with coverage less than 10 reads were ignored and only variants seen in at least 2 reads were considered. Using the functionality "SNP Effect analysis", variants were filtered for SNPs with protein coding effects. In addition, SNPs present within 3'UTR and 3' downstream were also identified using SNP effect analysis. Putative somatic and germline variants were investigated using the "Find Somatic Variants" analysis module. The somatic variants were filtered for low coverage (≤ 10 reads), strand bias (≥ 50) and a confidence score (0-800) to focus only on high confidence SNPs.

3.5 Biological Interpretation

A multi-omic pathway analysis was performed on the differentially expressed genes identified from miRNA and mRNA analysis of nasopharyngeal carcinoma cell line C666. This was done using both open source curated pathways like Wiki Reactome, BioCyc and hand created

GPML pathways as well as literature derived Natural Language Processing (NLP)-based pathways. These pathway resources are available in the tool. The biological network was enriched by combining entity interactions from differentially expressed mRNAs, miRNAs and their predicted targets into a miRNA regulatory network.

4 Results

4.1 Data QC

- RNA-Seq Quality of the sequenced reads was assessed by the Pre-Alignment QC module in Strand NGS and found to be satisfactory. The RNA-seq mapping is based on the Phred 64 scale and shows that the mean base qualities range from 63 to 70 with the majority of the reads having a read quality of 49 and above. The post alignment QC shows that most RNA-seq reads have an alignment score >95% and over 65% of the reads are uniquely mapped with a mapping quality of 254.
- small RNA Quality of the sequenced reads was also assessed. The pre alignment quality control on small RNA samples indicates a mean base quality ≥ 34 across the read positions, with most of the reads having an average base quality $\geq Q30$. The base quality by position plot shows a consistently good mean base quality of 38 throughout the read length (58 bp). The post alignment QC shows that 50-70% of the reads are uniquely mapped and hence assigned a mapping quality of 254.

The genic region QC plot (figure 2) shows the distribution of reads across different genic classes that were annotated with miRBase, tRNAscan-SE and Ensembl. The microRNA fractions in these cell lines are 9% in the normal NP460 cell line compared with 6.6% in C666, 15.6% in HK1 and 5% in X666 NPC cell lines. These numbers are very similar for NP460, C666 and X666 and only slightly higher for the HK1 cell line (16% vs 12%), as reported by [5] using miRBase v19.



Figure 2: Small RNA class distribution in genic region QC plot for HK1 cell line

4.2 Alignment Output

The alignment statistics for mRNA and small RNA data is given below:

- RNA-Seq The RNA alignment statistics report (figure 3), shows the number of uniquely mapped and multiply mapped reads for every sample. The percentage of aligned reads ranged from 86% to 96% and are in agreement with those reported by [5] using TopHat [6]. Please refer to table 1 for description of samples.
- small RNA The average read quality of most reads in the data is \geq Q30. The percentages of aligned reads ranged from 88% to 96% (alignment report now shown), and are in agreement with those reported by [5] which used CLC Bio for the small RNA alignment.

4.3 Differentially expressed genes

mRNA differential expression analysis identified 1224 genes in HK1 and 1433 genes in C666 as differentially expressed, when compared to NP460. We compared these differentially expressed entities across the NPC cell lines. A Venn diagram of differentially-expressed mRNAs and miRNAs against the immortalized nasopharyngeal epithelial cell line NP460 are shown in figure 4. We see that 402 mRNAs and 15 miRNAs are commonly expressed in the two NPC cell lines. Of the unique entities, 822 mRNAs and 32 miRNAs show significant expression changes only in HK1. On the other hand, 1031 mRNAs and 29 miRNAs are differentially expressed uniquely in C666.



Figure 4: Comparison of differentially-expressed transcripts in two NPC cell lines

A comparison of the expression analysis in Strand NGS vs. results published by [5] was done (Table 2). In HK1, out of the 1224 RNA genes we identified, 1024 were common, 200 were uniquely identified by Strand NGS and 93 were uniquely identified by [5]. In C666, there

Alignment Statistics								
	SRR1178329	SRR1178330	SRR1178331	SRR1178332				
Total number of reads	27,076,636 (100%)	34,508,151 (100%)	32,699,704 (100%)	26,626,486 (100%)				
Aligned reads	25,701,599 (94.9%)	29,651,780 (85.9%)	31,340,167 (95.8%)	25,553,669 (96%)				
- Uniquely matched reads	22,979,321 (84.9%)	25,877,593 (75%)	27,811,495 (85.1%)	22,845,993 (85.8%)				
- Multiply matched reads	2,722,278 (10.1%)	3,774,187 (10.9%)	3,528,672 (10.8%)	2,707,676 (10.2%)				
Unaligned reads	1,375,037 (5.1%)	4,856,371 (14.1%)	1,359,537 (4.2%)	1,072,817 (4%)				
- No matches found	1,375,037 (5.1%)	4,856,371 (14.1%)	1,359,537 (4.2%)	1,072,817 (4%)				
- Too many matches	NA	NA	NA	NA				
Reads ignored due to failure	0 (0%)	0 (0%)	0 (0%)	0 (0%)				
Reads ignored due to small size	0 (0%)	0 (0%)	0 (0%)	0 (0%)				
Total reads screened	NA	NA	NA	NA				
Maximum read length	58	58	58	58				
Average read length	57	57	57	57				
Alignment Status								
	SRR1178329	SRR1178330	SRR1178331	SRR1178332				
Aligned to transcriptome only	19,000,273	20,578,211	24,919,380	19,839,790				
- Involving known splices	18,973,599	20,539,941	24,892,247	19,815,505				
- Involving novel splices	26,674	38,270	27,133	24,285				
Aligned to genome	6,701,326	9,073,569	6,420,787	5,713,879				
- Involving transcriptome	14,632	19,508	13,680	11,074				
- Without transcriptome exons	6,686,694	9,054,061	6,407,107	5,702,805				
Pairs aligning to same transc	0	0	0	0				

Figure 3: Alignment report for RNA-Seq data

were 1246 common RNA genes with 187 and 85 genes uniquely identified by Strand NGS and [5] respectively. A similar comparison was done for the microRNAs. Differential microRNA expression analysis in Strand NGS identified 47 miRNAs in HK1 and 44 miRNAs in C666 as differentially expressed. In HK1, 41 were common, 6 were unique and 6 entities were not present in Strand NGS. Out of the 44 entities in C666, 37 were in common, while 7 were unique and one was absent in our set. The original study includes both human and EBV microRNAs but we are reporting numbers only from the human dataset comparisons. The comparison shows an overall agreement of differentially expressed genes, with minor differences that are likely due to the non-identical parameters used during alignment and expression analysis.

_				
Γ	Type	Unique Strand	Common	Unique Szeto
Г	HK1 mRNA	93	1024	200
Г	C666 mRNA	187	1246	85
Г	HK1 miRNA	6	41	6
Г	C666 miBNA	7	37	1

 Table 2: A comparison of the expression analysis in Strand

 NGS vs. Szeto et al. [5]

Further, the fold change was plotted against the p-values for a pair of conditions to determine the significant thresholds (Figure 5 and 6). The p-values are uniformly distributed along the y-axis making it difficult to choose an appropriate cut-off. We therefore chose the commonly accepted p-value threshold of 0.05 to select entities for further analysis.

We went ahead with differentially expressed entities



Figure 5: Fold change and p-values for HK1 vs NP460

using non-adjusted p-values to be consistent with the analysis of [5]. Also to illustrate visualisation in Strand NGS, a scatter plot of the signal intensity values with a fold change ≥ 2 shows the entities that are up- and down-regulated in the C666 cell line (Figure 7).

4.4 Target predictions (miRNA-mRNA interactions)

Corresponding targets of miRNAs (Fold change ≥ 2) were predicted in Strand NGS using PITA and TargetScan databases. The original study by [5], predicted 7951 targets for 149 miRNAs identified in any one of the samples from miRanda, PITA and TargetScan databases. Of these 7951 gene targets, 6423 were



Figure 6: Fold change and p-values for C666 vs NP460



Figure 7: Scatter plot of entities that are at least two-fold up-regulated or down-regulated in C666 (SRR1178329) compared to the NP460 (SRR1178332), coloured based on normalized signal intensities

inversely expressed. We have considered only the human miRNAs expressed in two of the cell lines and not included target predictions from miRanda database. Targets of 47 HK1 and 44 C666 miRNAs were predicted from PITA and TargetScan databases. Only the entities supported by target prediction databases with a p-value cut-off 0.05 were selected including 557 HK1 and 5107 C666 miRNA-mRNA interactions. miRNAs negatively regulate mRNA expression and hence, an inverse miRNA and mRNA expression profile is expected for these interactions. In tumor cells, miRNAs can be either downregulated or upregulated and their targets may show consistent or inconsistent deregulation patterns [4].

The output of a "Find Targeted Gene" prediction in smallRNA provides a list of genes predicted against a group of miRNAs. The individual miRNA-mRNA interactions, p-value, regulation and fold change were generated by a script. A pattern score was assigned to prioritize the inversely regulated pairs (up-down, downup). Integration of the transcript expression data and predicted interactions identified 307 HK1 and 197 C666 miRNA-mRNA targets pairs with an inverse expression. Table 3 shows an example of the predicted interactions for two miRNAs of the miR-200 family in C666 that were expressed with at least a 2-fold change. Some of these genes are involved in apoptosis and epithelial to mesenchymal transition. Particularly, the miR-200 family members have been reported to be expressed at very low levels in normal ovarian surface cells and substantially increase in expression in ovarian cancer, whereas expression of ZEB1 and ZEB2 shows the opposite pattern [3]. Furthermore, these miRNAs were validated by qRT-PCR by [5], and found to be consistent with the RNA-seq expression profile.

miRNA	p-val	Gene	FC	FC
Accession			(miRNA)	(Gene)
	0.05	EMP1	10.84	-4.64
	0.01	ITGA6	10.84	-4.51
1	0.02	PRKACB	10.84	-7.15
1	0.03	PTPRG	10.84	-6.74
1	0.00	ZEB1	10.84	-4.23
1	0.03	TBX18	10.84	-16.28
1	0.01	ARMCX2	10.84	-8.15
has miD 141 2m	0.03	LHFP	10.84	-11.01
lisa-mik-141-5p	0.02	LPHN2	10.84	-8.59
1	0.04	RAB30	10.84	-5.11
i	0.01	DNAJC15	10.84	-4.03
1	0.04	PCDH18	10.84	-14.46
1	0.02	FAT4	10.84	-10.36
i	0.01	CHD9	10.84	-2.80
i i	0.02	NETO2	10.84	-18.08
i i	0.03	CSMD3	10.84	-6.75
i i	0.04	ZNF605	10.84	-15.02
	0.05	EMP1	9.73	-4.64
1	0.03	IGF2R	9.73	-3.73
i i	0.01	ITGA6	9.73	-4.51
i i	0.02	PRKACB	9.73	-7.15
i i	0.00	ZEB1	9.73	-4.23
i i	0.01	WASF1	9.73	-3.26
i i	0.03	TBX18	9.73	-16.28
i i	0.03	LHFP	9.73	-11.01
i i	0.03	BASP1	9.73	-10.35
hsa-miR-200c-3p	0.01	FERMT2	9.73	-5.43
i i	0.02	LPHN2	9.73	-8.59
i i	0.01	DNAJC15	9.73	-4.03
i i	0.01	CCNJL	9.73	-6.50
i i	0.02	FAT4	9.73	-10.36
	0.01	CHD9	9.73	-2.80
	0.03	SGPP1	9.73	-4.01
i	0.03	RAB34	9.73	-8.57
i	0.03	CSMD3	9.73	-6.75
i	0.01	CYYR1	9.73	-15.62
1	0.04	ZNF605	9.73	-15.02

Table 3: Predicted targets (PITA and TargetScan) of miR-
NAs that show an inverse expression pattern in
C666

4.5 Sequence variant analysis

SNP detection identified 5,22,147 variants that include substitutions, deletions, insertions and complex variants. Based on a comparison of the normal NP460 against the tumor sample HK1 classified the variants as somatic (1,65,050), germline (2,11,098) and ambiguous (1,46,047). The variants were filtered for low coverage (≤ 10), strand bias (≥ 50) and score (0-800) using the "Region list operations" utility, to remove false positive variant calls that resulted in 29,822 somatic and 16,552 germline SNPs. A subset of 175 of these somatic variants are represented in COSMIC. A total of 37,063 variants have protein coding effects out of which 16,666 were predicted to be damaging variants in 2923 unique genes. From a set of germline SNPs with protein effects (2,71,486) functional risk was predicted to be

Hs_Wnt_Signaling_Pathway_a WikiPathways - Analysis Path	nd_Pluripotency_` ways – Unknown	WP399_68011 - H	omo sapiens		WNT2 WNT26 WNT6	7207 Pm 7208	041	Cytoskeleton	
No normalization_mRNA analysis No normalization_smRNA_Analysis Interpretation C666vs460 (Non-averaged) C666 vs NP460 (Non-averaged) Entity List C666vs460 nonadj pval vla R script Targeted Genes : p-value cut-off =0			RNA_Analysis eraged)	WNT3A WNT3A WNT4 WNT5A WNT56 WNT56 WNT56 WNT56	PZD10		RHCA	- Apoptosis	
Pathway	p-val Matcl	ned Enti Pathway	/ Enti p-value(Match	ed Pathway		LIPPO	AX112	6 Promote	Pluripotency a Differentiation
Hs_Biogenic_Amine_Synthesi.	0.09 2	15	0.240201	15 🔺			1 X	BIO TR: Transcriptional I	Regulation
Hs_SIDS_Susceptibility_Path.	. 2.77 16	166	0.00331 4	166			CSNKE GSK36		
Hs_Wnt_Signaling_Pathway	. 5.09 13	100	6.1659310	100			APC		
Hs_TWEAK_Signaling_Pathw	. 0.00 6	41	1.0 0			MARK		*p Shipu	witewant to Unike own the Complege Fing Enzyme
Hs_Arrhythmogenic_Right_V.	3.87 15	78	1.0 0				Stabilized hote-catenin		PEXW2
Hs_Codeine_and_Morphine	0.27 1	9	0.057652	9			CTINIDI		PAFAH1B1
Hs_miR-targeted_genes_in	0.02 3	17	1.0 0					Usquite Tagged	
Hs_Interferon_type_I_signali.	0.04 5	54	4.536227	54				- month	
Hs_Gastric_cancer_network_	0.31 2	32	0.013823	32		NK NK		Nuclaus 265 Protes	some Degradation
Hs_Synaptic_Vesicle_Pathwa.	0.54 2	51	0.24020 1	51		1953	TCP7 Gesuzho		
Hs_Blood_Clotting_Cascade_	0.18 2	22	1.0 0	22	60000	TR	TOFFL2		
Hs_mRNA_Processing_WP41	0.98 1	127	1.083468	127	TOWER	4 2 1	TOPILI		
Hs_TSLP_Signaling_Pathway_	0.00 6	48	0.013823	48		ESPRE		COND1 CTINID1	
Hs_Notch_Signaling_Pathway	0.00 6	46	0.240201	46	//	,"	NKD2	COND2	
Hs_Hypertrophy_Model_WP	. 5.39 5	20	0.24020 1	20	(POUSF1	aganin NKD1 behacitanin	201803	
Hs_ATM_Signaling_Pathway	0.14 3	41	0.057652	41		90X2		JJN	//
Hs_Electron_Transport_Chai.	0.96 1	104	4.536227	104		<u> </u>		MMP2	
Hs_Drug_Induction_of_Bile	. 0.45 1	17	0.057652	17			PPPRIA PPPRA		
Hs_NAD_Biosynthesis_II_(fro.	0.24 1	8	0.057652	8			PPP2R2A PPP2CA	C044 WNT11	
Hs_Glycerophospholipid_Bio.	0.16 2	34	0.057652	34		Pluringtency	PPP2R2B PPP2CB		
Hs_miRNAs_involved_in_DN.	0.80 1	69	0.057652	69 👻			PPP2R2C PPDRU	Ifferentiation	
Find:	Find Next 🙆 Fin	d Previous 🔲 Match	Case		No permaination mPNI	apakuric (parmalized) Lu			
					r in in in a 2 diori_inicity	anaysis (normalized) [N	normalization_smkivA_Analysis (normalized)	
Concernation of	NIA analusia	C Show	No normalization conDNA An	ale unite	Name	DB	DB ID	SRR1178329	SRR1178332
je provino normalizacion_mi		Je snor	vivo normalización_smktvA_An	aiysis	WNI108	Entrez ID	/480	10.161	3.585
p-value cutoff: 0.5	Min # of matches: 1	p-va	lue cutoff: 0.5 Min # of r	natches: 1	FZD7	Entrez ID	8324	12.48	3.954
			WNI58	Entrez ID	81029	8.441	1		
					LEFI	Entrez ID	511/6	9.374	3
				WN17A	Entrez ID	/4/6	-6.644	7.003	
 Show pathways that pass either filter 				10711	Entrez ID	83439	2.322		
C. Show pathways that page he	C Show pathways that pass both filters				60.44	Entrez ID	/9/6	3.51/	9,408
s briow pacitivays triat pass bt	s o brow pacieways crac pass both ricors				0044	Entrez ID	960	10.301	
					PRKCA	jEntrez ID	81 66	2	10.847

Figure 8: Multi-omics analysis of curated pathway seen in pathway viewer - Wnt signalling pathway.

damaging in 4509 non-synonymous variants involving 767 genes of unknown clinical significance. These germline and somatic variants can be subject to a deeper interpretation to elucidate their effect on protein function.

A novel coding SNP in TP53 (p.118T>P) and another known SNP rs1695 (p.118 I>V) in GSTP1 were reported by [5]. Our analysis did not call out the novel TP53 SNP either as a somatic or germline variant and was instead a heterozygous (T/G) SNP call in the normal NP460 and both NPC cell lines. The other reported SNP, rs1695 (A>G; Ile105Val) is a known heterozygous missense somatic variant in GSTP1 found to be present in both NPC cell lines C666 and HK1. The inheritance of at least one GSTP1 valine-105 allele apparently confers a significantly increased risk of developing therapy-related acute myeloid leukemia post chemotherapy [1].

4.6 Multi-omics pathway analysis

We have focused on the C666 carcinoma cell line for pathway analysis. The differentially expressed genes identified from miRNA and mRNA analysis of nasopharyngeal carcinoma cell line C666 were analyzed in a multi-omic pathway analysis using the open source Wiki Pathways in Strand NGS. In the C666 cell line, the Wnt-signaling (p-val = 0.001), beta-Integrin (p-val = 0.006), Apoptosis (p-val =0.003), EGR-EGFR (p-val =0.113) related pathways entities were significantly represented using a p-value cut-off of 0.05, as shown in the list of pathways in Figure 8. For example, in the Wnt-signaling pathway, WNT10B, ESRRB, FZD7 were upregulated while those downregulated include TCF7L1, PRKCA and PLAU (Figure 8) with established roles in cancer. One of the upregulated genes namely the

estrogen-related receptor beta (ESRRB), involved in pluripotency, is a potential tumor marker with a direct function in cancer progression.

Entities identified by different annotations, for example Entrez Gene ID and RefSeq Transcript ID, are mapped using BridgeDb and only mapped entities are visible on the pathway. Most pathways have limited microRNA information and require extensive curation. We were not able to find any curated pathways that were significantly enriched for microRNAs. Using a miRNA entity list we identified the network targets and regulatory.

A multi-omic analysis was performed comparing mRNA and miRNA entities on the NLP-derived miRNA regulatory network. The NLP-based networks are generated by matching selected entities to entities in the Interaction database in Strand NGS. It then retrieves relations like expression, promoter binding, and regulation between the set of matched entities that satisfy a relation score threshold. Eventually, it displays the results in the form of a graphical network. Figure 9 shows a sub-set of the merged miRNA network run on the C666 entities. The network provides information of the microRNAs belonging to the miR-200 family (miR200a, miR200b and miR200c) and its interacting mRNAs like ZEB1, PROM1 and LARP6 participating in apoptosis and epithelial to mesenchymal transition. Expression patterns of entities belonging to the two experiments for the chosen interpretation along with the relation types based on NLP are highlighted on the network.

An over lay of 3'-UTR SNP information for miRNA



Figure 9: NLP-based merged miRNA regulatory network.



Figure 10: SNP overlay on miR-targeted genes in squamous cell - TarBase (Homo sapiens).

interacting genes is also useful to interpret disease outcomes. SNP Effect analysis identified 38,181 SNPs located in the 3'-UTR/3'-downstream region of 5116 genes. For an imported pathway of miRNA and their target genes in squamous cells, 3' UTR SNP information was overlaid to identify genes with a 3' UTR SNP targeted by miRNAs (Figure 10). There are three 3' UTR SNPs called out in the IL18 gene, one of which is reported to be a disrupted binding pair [5]. However we found SNPs

in low coverage regions or were germline polymorphisms without functional risk predictions. Hence these calls remain inconclusive and require higher coverage for more confident SNP calls.

5 Conclusions

Strand NGS provides an efficient and powerful way to quickly and accurately analyse the integrated transcriptome that includes alignment, quality inspection, expression analysis, sequence analysis and biological interpretation. The purpose of an integrated transcriptomic study is to improve our understanding of miRNA-mRNA interactions in regulatory networks. In addition to identifying microRNA targets computationally, inversely expressed miRNA-mRNAs should aid in establishing the pathways of relevance. Entities identified as differentially expressed in NPC cells were found to be enriched in proliferation, adhesion, survival, and apoptosis pathways. A better understanding of the molecular signalling pathways, such as in NPC, steers the identification of novel diagnostic and prognostic biomarkers and personalized treatment options for cancer.

6 Acknowledgements

We would like to thank Prof. Maria Li Lung and her team from the University of Hong Kong, for the data. We also acknowledge the support of our colleague Anita Sathyanarayanan for helping with figures 5 and 6, which are generated in R programming language.

References

- [1] James M Allan, Christopher P Wild, Sara Rollinson, Eleanor V Willett, Anthony V Moorman, Gareth J Dovey, Philippa L Roddam, Eve Roman, Raymond A Cartwright, and Gareth J Morgan. Polymorphism in glutathione s-transferase p1 is associated with susceptibility to chemotherapy-induced leukemia. *Proceedings of the National Academy of Sciences*, 98(20):11592–11597, 2001.
- [2] Simon Anders and Wolfgang Huber. Differential expression analysis for sequence count data. *Genome biol*, 11(10):R106, 2010.
- [3] Ausra Bendoraite, Emily C Knouf, Kavita S Garg, Rachael K Parkin, Evan M Kroh, Kathy C O'Briant, Aviva P Ventura, Andrew K Godwin, Beth Y Karlan, Charles W Drescher, et al. Regulation of mir-200 family micrornas and zeb transcription factors in ovarian cancer: evidence supporting a mesothelial-to-epithelial transition. *Gynecologic oncology*, 116(1):117–125, 2010.
- [4] Li Guo, Yang Zhao, Sheng Yang, Hui Zhang, and Feng Chen. Integrative analysis of mirna-mrna and mirnamirna interactions. *BioMed research international*, 2014, 2014.
- [5] Carol Ying-Ying Szeto, Chi Ho Lin, Siu Chung Choi, Timothy TC Yip, Roger Kai-Cheong Ngan, George Sai-Wah Tsao, and Maria Li Lung. Integrated mrna and microrna transcriptome sequencing characterizes

sequence variants and mrna-microrna regulatory network in nasopharyngeal carcinoma model systems. *FEBS open bio*, 4:128–140, 2014.

[6] Cole Trapnell, Lior Pachter, and Steven L Salzberg. Tophat: discovering splice junctions with rna-seq. *Bioinformatics*, 25(9):1105–1111, 2009.



New Generation Healthcare

Strand was founded in 2000 by computer science and mathematics professors from India's prestigious Indian Institute of Science who recognized the need to automate and integrate life science data analysis through an algorithmic and computational approach. Strand's segue into the life sciences was through informatics products and services for research biologists, chemists, and toxicologists that combine advanced visualization, predictive systems modeling, data integration and scientific content management - over 2000 research laboratories worldwide (about 30% of global market share) are licensees of Strand's technology products, including leading pharmaceutical and biotechnology companies, research hospitals and academic institutions. With a recent investment by Biomark Capital, Strand has grown its established team to over 200 employees, many with multidisciplinary backgrounds that transcend computation and biology.

Since 2012, Strand has been expanding its focus to include clinical genomics, spanning sequencing, data interpretation, reporting and counseling. Strand operates a 10,000 square foot laboratory space with state-of-the-art clinical genomics capabilities and is also establishing Strand Centers for Genomics and Personalized Medicine in several hospitals around the world to serve as outreach points for genomic counseling. Based on the experience gained from sequencing, analyzing, interpreting and reporting on clinical samples over a wide variety of clinical indications, Strand has developed an end-to-end solution for clinical labs that handles all stagesfrom analysis to reporting. The interpretation and reporting software platform has been designed and developed specifically for the medical professional, ranging from the molecular pathologist to the physician. By enhancing sequence-based diagnostics and clinical genomic data interpretation using a strong foundation of computational, scientific, and medical expertise, Strand is bringing individualized medicine to the world.

For more information about Strand, please visit **www.strandls.com**, or follow us on twitter @StrandLife.

INDIA

5th Floor, Kirloskar Business Park, Bellary Road, Hebbal, Bangalore 560024

548 Market Street, Suite 82804, San Francisco, CA 94104