

Methods for Gene Coexpression Analysis

Assessment and Integration for Study of Deregulation in Cancer

O. Griffith¹, E. Pleasance¹, D. Fulton², M. Bilenky¹, G. Robertson¹, S. Montgomery¹, M. Oveis¹, Y. Pan¹, M. Zhang¹, M. Ester², A. Siddiqui¹, and S. Jones¹

**1. Genome Sciences Centre, Vancouver, Canada
2. Simon Fraser University, Burnaby, Canada**

British Columbia Cancer Agency
Vancouver, British Columbia, Canada

Canada's Michael Smith
Genome Sciences Centre
www.bcgsc.ca

1. Abstract

We anticipate that some cases of cancer progression are mediated through changes in genetic regulatory regions that can be detected through gene expression studies and bioinformatics analyses. Co-expressed genes are commonly identified by global analyses of large sets of expression experiments and data from several expression platforms are available. To assess the utility of publicly available expression datasets we have analyzed Homo sapiens data from 1202 cDNA microarray experiments, 242 SAGE libraries and 667 Affymetrix oligonucleotide microarray experiments. The three datasets compared demonstrate significant but low levels of global concordance. Assessment against the Gene Ontology (GO) revealed that all three platforms identified more co-expressed gene pairs with common biological processes than expected by chance, and, as the Pearson correlation for a gene pair increased, it was more likely to be confirmed by GO. The Affymetrix dataset performed best, with gene pairs of correlation 0.9-1.0 confirmed by GO in 74% of cases. However, in all cases, gene pairs confirmed by multiple platforms were more likely to be confirmed by GO, and we have shown that combining results from different expression platforms increases reliability of coexpression. Using this multi-platform/GO approach, we have created an easily extensible database of high-confidence co-expressed genes that currently contains 43,437 gene pairs for 7,103 genes. We are using this data as a high signal-to-noise input for the identification of cis regulatory elements in the cisRED project (www.cisred.org), and we are expanding the database of expression and coexpression data to include new species, platforms, and samples. Currently the database contains 6988 mouse and human samples from five different platforms. In ongoing work, we propose a novel approach to specifically identify mechanisms of gene deregulation in cancer by combining expression data, regulatory element predictions, and chromosomal mutation data.

2. Gene Expression Data

Table 1. Gene expression data in database

Species	Platform	Experiments	Unique genes
<i>H. sapiens</i>	SAGE (short)	243	20283
	Oligo. Array	1640	6613
	cDNA microarray	2852	11962
<i>M. musculus</i>	SAGE (long)	85	5388
	Oligo. Array	1802	6287
	cDNA microarray	366	4721
Total		6988	31185

3. Methods

Figure 1. Gene Coexpression Analysis.

Gene coexpression is determined by calculating a Pearson correlation (r) between each gene pair.

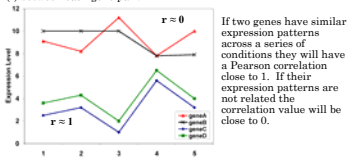


Figure 2. Platform Comparison Analysis.

Platforms are compared by calculating a correlation of correlations (r_c) for all gene pairs.



Figure 3. Gene Ontology (GO) Analysis.

Coexpression measurements can be assessed and calibrated against the Gene Ontology.



4. Platform Comparison Analysis

Figure 4. Affymetrix vs. SAGE

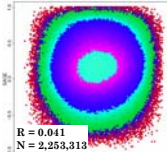
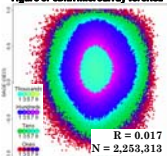
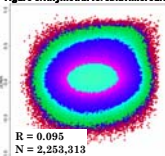


Figure 5. cDNA Microarray vs. SAGE



Figures 4-6: Poor levels of consistency were observed between platforms. Each point on the plots represents a bin of gene pairs, and its coordinates represent the correlation of those pairs for two different datasets. If the different datasets produced the same coexpression results we would expect a correlation of correlations close to 1 and would observe a straight line.

Figure 6. Affymetrix vs. cDNA Microarray



5. Gene Ontology (GO) Analysis

Figure 7. Multi-Platform Assessment

In general, as the Pearson correlation for a gene pair increases it is more likely to share a GO term. Gene pairs confirmed by multiple platforms (higher average Pearson) are more likely to share a GO term than those only coexpressed in a single platform.

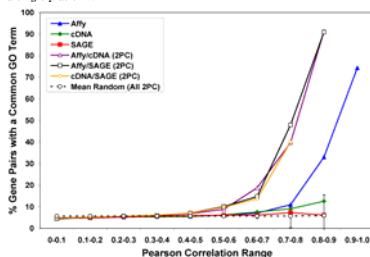
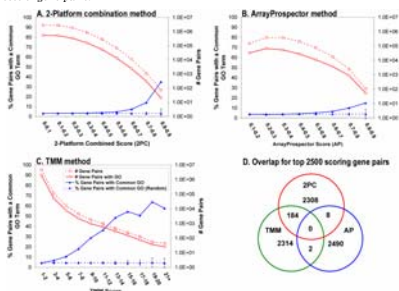


Figure 8. Comparison to other coexpression analysis methods

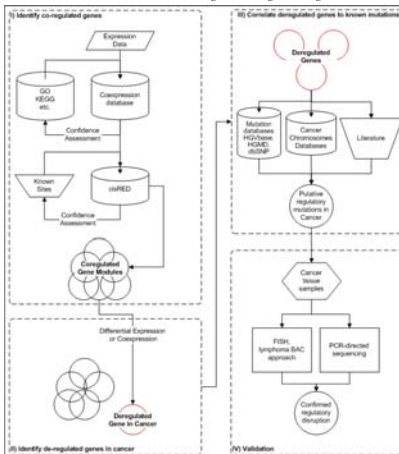
We compared our method of combining global coexpression from different platforms (2PC) to two other recent methods. One identifies experimental subsets separately and employs a 'vote-counting' method to identify gene pairs that appear highly coexpressed in multiple sets (TMM method)³. The second method uses a combination of singular value decomposition and kernel density estimation (ArrayProspector method)⁷. A direct comparison was impossible because the methods utilized different gene sets. Thus, we do not identify the 'best' method but rather show that each method is at least partially effective and we identify reasonable threshold scores for a high-confidence set of coexpressed genes. The Venn diagram indicates that each method identifies almost completely different sets of gene pairs.



6. Gene Deregulation in Cancer

Figure 9. Research plan

Once coexpressed genes are identified they can be used as part of the cisRED pipeline to predict cis regulatory elements (www.cisred.org). These regulatory elements will form the basis of our investigation into gene deregulation in cancer.



7. Conclusions

- Coexpressed genes can be identified based on large-scale gene expression data.
- Direct comparison of correlation values between platforms yields poor correlations ($R < 0.1$).
- Gene pairs identified as coexpressed with a higher Pearson correlation are more likely to share the same GO biological process.
- Gene pairs coexpressed in multiple platforms (higher average Pearson) are more likely to share a GO biological process than pairs coexpressed in only a single platform.
- Using the GO assessment, criteria for a high-confidence set of coexpressed genes can be defined and used for cis-regulatory element prediction.

Acknowledgements

funding | Natural Sciences and Engineering Council of Canada (for OG and EP); Michael Smith Foundation for Health Research (for OG, SJ and EP); CHRM/MSFHR Bioinformatics Training Program (for EP); Killam Trusts (for EP); Genome BC; BC Cancer Foundation

references | 1. Lee et al. 2004. Genome Research. 14:1085-1094; 2. Jensen et al. 2004. Nucleic Acids Research 32:W445-8