

# Multiple Diseases Discriminated by Quantitation of Blood Transcriptome

S Chao, K Marshall, F El Khettabi, H Tang, C Liew, S Mok

## Citation

S Chao, K Marshall, F El Khettabi, H Tang, C Liew, S Mok. *Multiple Diseases Discriminated by Quantitation of Blood Transcriptome*. The Internet Journal of Genomics and Proteomics. 2009 Volume 6 Number 1.

## Abstract

The steady circulation and physiologically interactive nature of blood ensures that this dynamic system encounters, transmits, and responds to a wide range of biological signals. In this context, we hypothesized that quantitative measurement of the blood transcriptome can enable the identification and validation of RNA transcripts that are specifically associated with the presence of a particular disease or clinical condition. In the present study, we have used 631 blood RNA expression profiling in conjunction with microarray technology to generate highly discriminative panels of 10 pairs of probe sets for each of four separate clinical conditions (gender, colorectal cancer, prostate cancer and osteoarthritis). The robust training set performance for each disease- or condition-specific multi-gene panel was corroborated with an independent test set, with areas under the receiver-operating characteristic curves ranging from 0.87 to 0.93 for each of the four conditions in the test set population. This study demonstrates that quantitative measurement of the blood transcriptome, in conjunction with microarray technology, can be used to generate highly discriminative multi-gene panels for many clinical conditions. This approach has great potential to enable the simultaneous monitoring of multiple disease states or clinical conditions from a single blood sample.

## INTRODUCTION

In 1932, physiologist Walter Cannon penned his classic *The Wisdom of the Body* [1]. This work introduced the concept of homeostasis, the process of autoregulation whereby biological systems self-monitor and self-adjust to preserve steady state equilibrium in a turbulent, ever-changing environment.

The circulating peripheral blood system is a critical integrative force by virtue of the blood's ongoing real-time involvement in the regulation, coordination, metabolism and immune maintenance of essentially all cells, tissues and organs. Functions of blood cells include transporting nutrients, oxygen and biomolecules, and removing cellular wastes. Blood is further intimately involved in immune surveillance throughout the body, and delivery of immune factors and healing mediators to sites of disease, infection and injury. Thus, the steady circulation and physiologically interactive nature of blood ensures that this dynamic system encounters, transmits, and responds to, a wide range of biological signals [2-5].

These dynamic, integrative features of blood, considered in context with the need for maintaining homeostasis, suggest that the presence of a specific disease or clinical condition

will be reflected in specific patterns of gene expression in blood, i.e. transcriptomic signatures. The transcriptome is the complete set of RNA transcripts present in a cell or tissue at any one time. Although a particular cell or tissue's DNA, or genome, is essentially unchanging, its transcriptome will vary according to the current physiological status of the cell or tissue. Thus, we have hypothesized that transcriptomic signatures in blood which are specific to states of health or disease can be identified and used to diagnose such states via transcriptional profiling of blood [2].

Advances of the past decade have made it possible for transcriptomes to be quantitatively profiled and compared on a genome-wide scale using powerful nucleic acid probe microarray technology [reviewed in 6]. Traditional microarray analyses are tissue biopsy-based, which limits the application of array technology to a limited number of clinical situations in which tissue is readily available. By contrast, the use of blood samples enables broadening of the application of transcriptional profiling analysis to a wider range of diseases and clinical conditions. Thus, blood is an ideal sample type which overcomes many of the limitations of traditional microarray studies [7].

In a series of studies, we and others have demonstrated that RNA profiles generated from circulating blood can be used to identify patients with a number of conditions [7,8], including: lung cancer [9] bladder cancer [10], colorectal cancer (CRC) [5], osteoarthritis [4], schizophrenia and bipolar disorder [11,12], kidney diseases [13,14], cardiovascular diseases [15-17], Crohn's disease [18] and diabetes [19]. In the present study, we have extended this approach and used single subject transcriptional signatures from a single blood sample to simultaneously assay for the detection of multiple diseases in a heterogeneous human population.

### MATERIALS AND METHODOLOGY

#### PATIENT SAMPLES

We recruited more than 1500 patients from multiple institutions between January, 2004 and July, 2008. We selected 631 patients for the current study of three distinct diseases. Informed consent was obtained according to the research protocols approved by the research ethics boards of each institution involved.

#### BLOOD COLLECTION AND RNA ISOLATION.

Samples of peripheral whole blood (10 ml) were collected in EDTA Vacutainer™ tubes (Becton Dickinson, Franklin Lakes, N.J.), and stored at 4°C until processing (within 6 hours). RNA was isolated at six different centers according to a standardized protocol. Plasma was removed after centrifugation and a hypotonic buffer (1.6 mM EDTA, 10 mM KHCO<sub>3</sub>, 153 mM NH<sub>4</sub>Cl, pH 7.4) was added at a 3:1 volume ratio to lyse the red blood cells. The mixture was centrifuged to yield a pellet containing predominantly white blood cells, and the pellet was re-suspended into 1.0 mL of TRIzol® Reagent (Invitrogen Corp., Carlsbad, CA) and 0.2 mL of chloroform. RNA quality was assessed on an Agilent 2100 Bioanalyzer RNA 6000 Nano Chip. RNA quantity was determined by absorbance at 260 nm/280 nm in a Beckman-Coulter DU640 Spectrophotometer. The samples were then stored at -80°C at a single center.

#### MICROARRAY HYBRIDIZATION.

Five-microgram samples of purified total RNA were labelled and analyzed using Affymetrix U133Plus 2.0 GeneChip oligonucleotide arrays (Affymetrix; Santa Clara, CA). Hybridization signals were adjusted in the Affymetrix GCOS software (version 1.1.1), using a scaling factor that adjusted the global trimmed mean signal intensity value to 500 for each array. The CEL files were expressed using MAS5 methods. Hybridizations were carried out in batches

across 30 lots of chips (3005291 to 4033799) between 2004 and 2008. All samples passed the recommended quality checks (background, present call, Raw Q, Scale Factor, and 3'/5' ratios for Gapdh and ActB).

#### DATA ANALYSIS

The expression levels from probe sets labelled "present" were log-transformed (base 2). Only the data from probe sets labelled as "present" in all samples across all studies were used (7,226 probe sets). For each study, samples with the condition of interest were labelled as "with condition of interest", while all other samples were labelled as "without condition of interest".

The probe set data for each condition of interest were organized into combinations of 10 pairs of genes. Each combination was evaluated for its discriminative power on the training set by calculating the receiver-operating characteristic (ROC) area under the curve (AUC). The combination that achieved the best ROC AUC was selected as the panel for the condition of interest. The process was repeated for each condition of interest.

Unique discriminative panels were determined for each condition of interest, namely: gender, colorectal cancer, prostate cancer and osteoarthritis. After the training set panels were determined for each condition of interest, a second set of studies was performed using an independent test set to further assess the discriminative power of the panels.

Analysis of the final results and generation of charts was performed using Microsoft Excel and MedCalc ([www.medcalc.be](http://www.medcalc.be)).

#### SAMPLES WITH CONDITION OF INTEREST

##### 1) Gender discrimination

First, we evaluated the reliability of the measurements and data analysis by searching for gene panels that can discriminate between genders. This was conducted in two phases.

##### 1) X-chromosome located genes

In this first phase, we searched for probe sets that exhibit consistent differential expression based on copy number difference.

2) Autosomal genes

In this second phase, as a model of general disease, we searched for discriminatory probe sets for autosomal genes that were effective at discriminating gender.

The training set was composed of 352 subjects (121F:231M). The test set had 279 subjects (96F: 183M).

2) Colorectal cancer

Training set = 80; test set = 68.

3) Prostate cancer

Training set = 80; test set = 63.

4) Osteoarthritis

Training set = 103; test set = 93.

Samples without Condition of Interest

1) No disease

Training set = 30; test set = 18.

2) Ovarian cancer

Training set = 29; test set = 7.

3) Bladder cancer

Training set = 30; test set = 12.

4) Crohn's disease

Training set = 0; test set = 18.

RESULTS

GENDER DISCRIMINATION

Figure 1

Figure 1 Gender discrimination

1A. Gender Training Set			1B. Gender Test Set		
352 Training Samples	Gender Chr-X	Gender Autosome	279 Test Samples	Gender Chr-X	Gender Autosome
121 Female			96 Female		
231 Male			183 Male		
ROCAUC	1.00	0.96	ROCAUC	0.996	0.87
F Accuracy	100.0%	92%	F Accuracy	99%	78%
M Accuracy	99.6%	87%	M Accuracy	98%	82%

Each row represents one sample. The first column represents the panel using x-linked genes, while the second column uses only autosomal genes for discrimination. Dark grey indicates a “female” prediction. Light grey indicates a “male” prediction. Accuracy for each gender is defined as the percentage of correctly predicted subjects from the total number of subjects.

Figure 2

Table 1 Reference Panels

	1A Gender Chr-X	1B Gender Autosome	1C Colorectal Cancer	1D Prostate Cancer	1E Osteoarthritis
	Probe set Gene	Probe set Gene	Probe set Gene	Probe set Gene	Probe set Gene
1	201210_m CCK3X	211902_s_m HNRNPA3	201491_m AHS1	200902_m HIF1A	200902_s_m KDM6A
2	200992_s_m KDM6A	215190_s_m LOC642236	30230_m RGS14	217527_s_m NFATC2IP	224787_s_m RAB18
3	201016_m EIF1AX	212059_m SMCHD1	214367_s_m GPM3	204394_s_m TSC22D2	202120_m KIAA0317
4	212008_s_m VDAC1	229872_s_m LOC100132999	222613_m C12orf4	207777_s_m SP140	200962_s_m KDM6A
5	201018_m EIF1AX	215190_s_m LOC642236	200085_s_m TGFBI	200544_s_m PROXB	200962_s_m KDM6A
6	207040_s_m ST13	216156_s_m TSR1	210151_s_m SEPT7	204900_s_m SAP30	219242_m CEP63
7	201018_m EIF1AX	212095_s_m AIF1	200019_s_m FAU	200403_s_m RPS4X	200850_s_m LDHA
8	210468_s_m SERBP1	208091_s_m RECQL	229876_m BTF3L4	1884747_s_m SEPT2	200962_s_m KDM6A
9	201018_m EIF1AX	226220_m NIACR2	210892_s_m APRT	209907_s_m HNRPOL	200962_s_m KDM6A
10	215440_s_m BEX4	200802_s_m NDRG1	227295_m IKIP	217747_s_m RPS9	227983_m RILPL2
11	210904_s_m IL13RA1	200805_m BCL2	200006_m TNFRSF18	214031_s_m IER3	204009_s_m WDR1
12	201018_m EIF1AX	227305_m	219797_m MGAT4A	217527_s_m NFATC2IP	200962_s_m KDM6A
13	201017_m EIF1AX	207094_m IL6RA	200707_m PGK1	201017_m EIF1AX	200962_s_m KDM6A
14	200862_s_m TTC3	229872_s_m LOC100132999	200190_m NDUFS8	201220_s_m ARIH2	204893_m TRAA2
15	208073_s_m TTC3	228841_m	200805_m TNFRSF18	200803_s_m RPS4X	200962_s_m KDM6A
16	227420_m CXorf15	227308_s_m TNFRSF9	214771_s_m MPRIP	229826_m SPOPL	219474_s_m KCID5
17	201017_m EIF1AX	217819_m SGLG7	217491_s_m COX7C	200962_s_m KDM6A	200962_s_m KDM6A
18	200990_s_m KDM6A	224864_m	200301_s_m INP5D	213257_m GTF2H5	226745_m ERN1
19	208174_s_m ZRSR2	206003_s_m DEFA1	214367_s_m GPM3	210860_m STA75B	200962_s_m KDM6A
20	213318_s_m CXorf40A	212484_s_m HNRPOL	228498_m ZCCHC7	226442_m ABTB1	217984_s_m CPSF3L

The panels of 10 pairs of probe sets with corresponding genes used for gender and disease predictions

### 1) X-CHROMOSOME LOCATED GENES

The final panel of 12 genes represented by 10 pairs of probe sets is detailed in Table 1A. Accuracy was greater than 99% for both male and females in the training set and 98% or greater for both genders in the test set. Several of the probe sets were expressed differentially at 1.4-fold and 1.6-fold which suggests that detection at less than 2-fold is possible with microarray technology.

### 2) AUTOSOMAL GENES

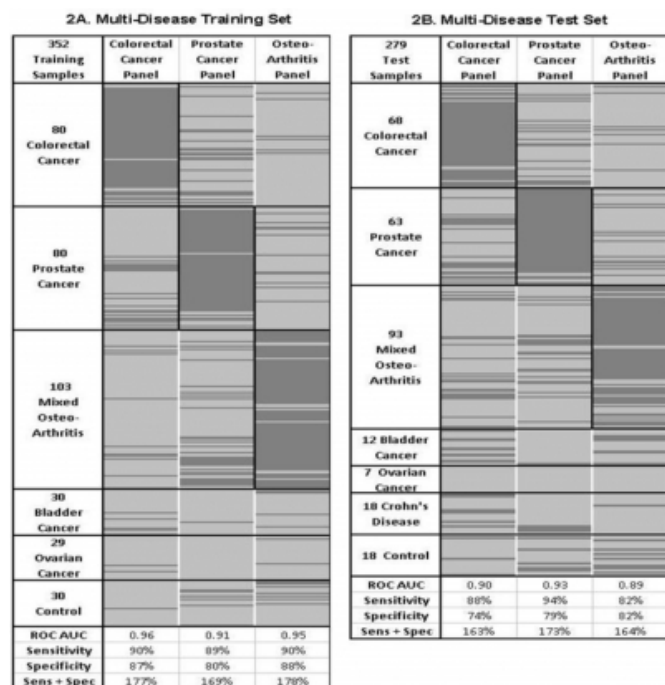
As a general model of disease, we searched for discriminatory probe sets for autosomal genes (Table 1B) and were able to achieve a ROC AUC of 0.96 in the training set (accuracy: 92%F; 87%M) and 0.87 on the test set (accuracy: 78%F; 82%M).

### MULTI-DISEASE DISCRIMINATION

The 20-gene probe set panels for each of the three different diseases are detailed in Table 1C-E. The discriminative power of each of these classification panels is detailed in Figure 2 with samples arranged by rows, grouped by disease type.

Figure 3

Figure 2 Disease discrimination



Each row represents one sample. Each column represents a disease prediction. Dark grey indicates a “positive prediction”; light grey indicates a “negative prediction”. Sensitivity is defined as the percentage of subjects predicted to have the disease of interest from the number of subjects

which actually have the disease of interest. Specificity is defined as the percentage of correctly predicted subjects as not having the disease of interest from all subjects that do not have the disease of interest.

The training set predictions (Figure 2A) achieved ROC AUC values of 0.96, 0.91 and 0.95 for colorectal cancer, prostate cancer and osteoarthritis, respectively. Sensitivity was 90%, 89% and 90% and specificity was 87%, 80% and 88% for colon cancer, prostate cancer and osteoarthritis. Each of the three disease-specific panels was able to reject most of the samples from conditions that were not included in the training phase (no disease, ovarian cancer, bladder cancer).

The independent test set results (Figure 2B) confirmed that each of the three disease-specific panels was effective at discriminating the particular disease it had been trained on, but was not discriminatory for either of the other two diseases nor for any of the three other sets of samples (conditions not of interest, bladder and ovarian cancers and Crohn’s disease). Colorectal cancer had ROC AUC of 0.90, prostate cancer, of 0.93, and osteoarthritis of 0.89; corresponding sensitivities were 88% (colorectal cancer), 94% (prostate cancer), and 82% (osteoarthritis) with specificity at 74% (colorectal cancer), 79% (prostate cancer), and 82% (osteoarthritis).

### DISCUSSION

To our knowledge, all studies to date on blood-based disease biomarkers have focused on identification of biomarkers for single diseases which can, at times, hide poor false positive results when predicting subjects with other diseases. In this study, we have expanded this approach to enable us to detect several diseases at once using one blood sample. The clinical utility of this approach should be immediately apparent. The ability to detect numerous pathologies at one time would simplify population disease screening. As this type of blood based tool becomes refined and applicable to be used in general populations, a patient could be screened at one visit for a range of diseases, for example, colorectal cancer and prostate cancer, rather than using several different and invasive tests.

As we show in this report, using integrated multi-disease analysis our laboratory can identify disease-specific gene expression signatures by quantitative measurement of the blood transcriptome. We have succeeded in markedly reducing crosstalk noise from confounding factors and were able to generate formulae for detecting the presence of certain diseases with specificity > 90% for various types of

organ-specific diseases including cancers. The inclusion of multiple diseases in this study increases the sample variability which can lead to improved performance [3,20].

It should also be pointed out that this approach inherently allows detection of multiple simultaneous conditions as the discrimination is not an either/or decision, but rather a set of independent parallel decisions.

This independent parallel approach is partially supported by preliminary results from a small set of patients with Crohn's disease. It is known that Crohn's disease is associated with a higher incidence of colorectal cancer [21], and it would therefore be beneficial for a colorectal cancer diagnosis to be able to differentiate patients with Crohn's disease that will eventually progress to colorectal cancer from those patients who will not.

At the time of this study, there were not a sufficiently large number of samples with Crohn's disease available to construct a training set. The medical record for these samples was incomplete and the presence of colorectal cancer was unknown. As a result, these samples were only included in the test set to evaluate the performance of the three diagnosis panels.

Six of the 18 samples with Crohn's disease had a positive call on the colorectal cancer panel, of these, one also had a positive call on prostate cancer along with two others and only one had a positive call for osteoarthritis. While these numbers are small, they do show a trend toward more positive calls for colorectal cancer than for prostate cancer or osteoarthritis, which is what would be expected if the panels are actually recognizing a biological signal rather than merely some random noise.

We initially demonstrated that subtle gene expression differences of less than 2-fold can be measured reliably as evidenced by the discrimination of gender differences with 100% accuracy from X-linked genes. We also attempted to discriminate gender using only autosomal genes (ROC AUC: 0.87, accuracy: 80%) in order to demonstrate the level of performance that is likely achievable using this method for disease/clinical condition discrimination.

The generation of unique, 10 pairs of probe sets for each of the three disease conditions of interest (colorectal cancer, prostate cancer, osteoarthritis) resulted in promising discriminatory training set panels for each disease with the ROC AUCs ranging from 0.91 to 0.96. The robust discriminatory capacity for each disease panel was

confirmed by the independent test set (ROC AUC range: 0.89 to 0.93). The close concordance between the training and test set data was reassuring, given that there were a number of potentially confounding variables including: multiple clinical sites used for sample collection; multiple laboratories used for RNA extraction; multiple different chip and reagent lots, the small expression-fold changes seen in blood RNA profiling as compared with the large changes seen in tissue expression profiling, varying durations of RNA storage and multiple microarray hybridizations extending across a four-year span.

We have recently reported on the use of blood RNA profiling with a seven-gene panel utilizing quantitative real-time polymerase chain reaction (qRT-PCR) to discriminate subjects with colorectal cancer from those with no cancer [4]. This approach allows an individual's relative risk of currently having colorectal cancer to be determined, thereby providing clinically actionable information about the need for further investigations such as colonoscopy. Although this method is a powerful tool for improving early detection of disease and provides novel information to enhance clinical decision-making, extending qRT-PCR technology to simultaneously assay for multiple diseases or clinical conditions is unsustainable. This relates to the practical limitations on the number of genes that can be included in a panel using qRT-PCR. The present study suggests that, microarray technology, with its ability to simultaneously measure the activity of a large number of RNA transcripts, can facilitate the application of blood transcriptome profiling to generate and assay multiple disease panels.

## CONCLUSION

In this study we used gender data to show in a straightforward and non-controversial manner the clinical utility of the integrated multi-disease analysis test. The test clearly differentiates male and female (98-99% accuracy) for both male and females even when sex chromosomal factors are excluded. That is, sexes were shown to be different in blood samples using (autosomal) genes, not genes related to the sex chromosomes. The test differentiates male and female regardless of confounding factors such as using samples from different clinics, over several years and use different microarray lots. Similarly, the test methodology is able clearly to indicate the presence or absence of various diseases (colorectal cancer, osteoarthritis, prostate cancer) in the samples. Such a test can be expanded to include other types of cancer and other diseases.

Thus the quantitative transcriptomic approach has significant advantages as a potential tool for personalized medicine.

This approach is not a genetic DNA marker test or polymorphism biomarker test which are the tests currently available and which have been sharply criticised as failing to agree in disease prediction between laboratories and failing to capture genetic contributions to disease risk. [22, 23] Rather, our quantitative measurement of the blood transcriptome reflects in real-time, gene expression alterations occurring over the whole transcriptome as this in turn guides phenotypic disease phenomena.

The present study demonstrates that blood transcriptome profiling in conjunction with microarray technology can be used to generate highly discriminative multi-gene panels for many diseases. This approach has great potential to enable the simultaneous monitoring of multiple disease states or clinical conditions from a single blood sample, which could, as the process is refined and developed, hold great promise as a population multiple disease screening tool.

### ACKNOWLEDGEMENTS

We would like to thank Dimitri Stamatiou, and Jay Ying for their technical assistance and Ma Jun for helpful comments and criticism of this manuscript. CC Liew, Samuel Chao, Faysal El Khettabi, Hongchang Tang and K Wayne Marshall are all employed by GeneNews Ltd, who sponsored this research.

### References

1. Cannon WB. The wisdom of the body. WW Norton: New York 1932.
2. Liew CC. 2002. Method for the detection of gene transcripts in blood and uses thereof. US Patent No. 7,598,031, filed: October 9, 2002, and issued October 6, 2009.
3. Dumeaux V, Olsen KS, Nuel G, et al: Deciphering Normal Blood Gene Expression Variation – the NOWAC Postgenome Study. *PloS Genetics*; 2010; 6:e1000873.
4. Marshall KW, Zhang, H, Yager T, et al: Blood-based biomarkers for detecting mild osteoarthritis in the human knee. *Osteoarthritis Cartilage*; 2005; 13: 861-871.
5. Marshall KW, Mohr S, El Khettabi F, et al: A Blood-based biomarker panel for stratifying current risk for colorectal cancer *Int J Cancer*; 2010;126:1177-1186.
6. Hocquette JF: Where are we in genomics? *J Physiol Pharmacol*; 2005; 56 S3: 37-70.
7. Mohr S, Liew CC: The peripheral blood transcriptome: new insights into disease and risk assessment. *Trends Mol Med*; 2007; 13: 422-432.
8. Edelman LB, Toia G, Geman D, Zhang W, Price ND: Two-transcript gene expression classifiers in the diagnosis and prognosis of human diseases. *BMC Genomics*; 2009; Dec 5; 10: 583. [epub ahead of print]
9. Showe MK, Vachani A, Kossenkov AV, et al: Gene expression profiles in peripheral blood mononuclear cells can distinguish patients with non-small cell lung cancer from patients with nonmalignant lung disease. *Cancer Res*; 2009; 69: 9202-10.
10. Osman I, Bajorin D, Sun TT, et al: Novel blood biomarkers of human urinary bladder cancer. *Clin Cancer Res*; 2006; 12(11 Pt 1): 3374-80.
11. Tsuang MT, Nossova N, Yager T, et al : Assessing the validity of blood-based gene expression profiles for the classification of schizophrenia and bipolar disorder: A preliminary report. *Am J Med Genet B Neuropsychiatr Genet*; 2005; 133B: 1-5.
12. Glatt SJ, Everall IP, Kremen WS, et al: Comparative gene expression analysis of blood and brain provides concurrent validation of SELENBP1 up-regulation in schizophrenia. *Proc Natl Acad Sci USA*; 2005; 102: 15533-8.
13. Alcorta D, Preston G, Munger W, et al: Microarray studies of gene expression in circulating leukocytes in kidney diseases. *Exp Nephrol*; 2002; 10: 139-49.
14. Twine NC, Stover JA, Marshall B, et al: Disease-associated expression profiles in peripheral blood mononuclear cells from patients with advanced renal cell carcinoma. *Cancer Res*; 2003; 63: 6069-75.
15. Bull TM, Coldren CD, Moore M, et al: Gene microarray analysis of peripheral blood cells in pulmonary arterial hypertension. *Am J Respir Crit Care Med*; 2004; 170: 911-9.
16. Ma J, Dempsey AA, Stamatiou D, Marshall KW, Liew CC: Identifying leukocyte gene expression patterns associated with plasma lipid levels in human subjects. *Atherosclerosis*; 2007; 191: 63-72.
17. Deng MC, Eisen HJ, Mehra MR, et al: Noninvasive discrimination of rejection in cardiac allograft recipients using gene expression profiling. *Am J Transplant*; 2006; 6: 150-160.
18. Burakoff R, Hande S, Ma J, et al: Differential Regulation of Peripheral Leukocyte Genes in Patients with Active Crohn's Disease and Crohn's Disease in Remission. *J Clin Gastroenterol*; 2010; 44:120-126.
19. Takamura T, Honda M, Sakai Y, et al; Gene expression profiles in peripheral blood mononuclear cells reflect the pathophysiology of type 2 diabetes. *Biochem Biophys Res Comm*; 2007; 361: 379-384.
20. Shen-Orr SS, Tibshirani R, Khatri P, et al: Cell type-specific gene expression differences in complex tissues. *Nature Methods*; 2010;7:287-9.
21. Freeman HJ. Colorectal cancer risk in Crohn's disease. *World J. Gastroenterol* 2008; 14(12): 1810-1811
22. Ng PC, Murray SS, Levy S, Venter JC: An agenda for personalized medicine. *Nature*; 2009; 461: 724-6.
23. Samuel P, Dickson, Kai Wang, Ian Krantz, Hakon Hakonarson, David B. Goldstein: Rare variants create synthetic genome-wide associations. *PloS Biology*; 2010 Jan. 8:1.

**Author Information**

**Samuel Chao**

GeneNews Ltd

**K Wayne Marshall**

GeneNews Ltd

**Faysal El Khettabi**

GeneNews Ltd

**Hongchang Tang**

GeneNews Ltd

**C.C. Liew**

GeneNews Ltd

**Samuel Mok**

M.D. Anderson Cancer Center