

#### Improved precision and accuracy for microarrays using updated probe set definitions

Special Topics in Bioinformatics, SS 2010 Andrea Salfinger, 0455957

#### Report about ...

#### **BMC Bioinformatics**

Correspondence



Open Access

# Improved precision and accuracy for microarrays using updated probe set definitions

Rickard Sandberg<sup>1</sup> and Ola Larsson\*<sup>2</sup>

Address: <sup>1</sup>Massachusetts Institute of Technology, Department of biology, 68-211, Cambridge, MA 02139, USA and <sup>2</sup>University of Minnesota, Department of Medicine, MMC 276, Minneapolis, MN 55455, USA

Email: Rickard Sandberg - sandberg@mit.edu; Ola Larsson\* - larss004@tc.umn.edu

\* Corresponding author

Published: 8 February 2007

Received: 26 September 2006 Accepted: 8 February 2007

BMC Bioinformatics 2007, 8:48 doi:10.1186/1471-2105-8-48

This article is available from: http://www.biomedcentral.com/1471-2105/8/48

© 2007 Sandberg and Larsson; licensee BioMed Central Ltd.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<u>http://creativecommons.org/licenses/by/2.0</u>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

## Outline

- Recap: Microarrays
- State of the art: Problems
- Updated Probe Set Definitions
- Motivation
- Purpose of this study
- Study Design
- Methods
- Experimental Data
- Precision
- Accuracy
- Results
- Summary



- When?
  - introduced in 1995
- What?
  - measure extracted mRNA transcripts of a cell
- How?
  - hybridize the amplified and labeled transcripts to template strands on a solid surface (microarray), measure which and how many transcripts have bound
- Why?
  - to measure gene expression levels of a cell











• What is measured?



• How to get gene expression level from this, when only DNA fragments of 25 bp length are measured?

mRNA reference sequence

3'









## State of the art: Problems (1)

- Affymetrix platforms currently in use were designed before genomes were fully sequenced
  - many probes were designed after consensus sequence of clusters of Expressed Sequence Tags (ESTs)
- <u>Problem:</u> design of probes based on (meanwhile) outdated genomic knowledge!

(we know it better now)

# State of the art: Problems (2)

- Problem:
  - accurate probe set definitions essential for integrating probe signals into expression levels!
  - many old Affymetrix chips still in use in research community
  - it takes a lot of time to develop new chips
  - What to do with the studies that have been carried out with older chip-generations?
- Solution:
  - analyze the data using <u>updated probe set definitions</u>

#### **Updated Probe Set Definitions**

- = <u>reannotation of the the existing probes</u> on Affymetrix platforms
  - reannotation better reflects transcript information and gene annotations available today
- Updated probe set definitions map probes to transcript annotations, e.g.
  - EnsEMBL transcripts
  - Refseq
  - Entrez GenelD

#### **Updated Probe Set Definitions**

#### original probe set definitions

#### updated probe set definitions



Probe set 123: Gene X

Probe set 124: Gene X

> Gene expression level of X???



# Motivation (1)

- original Affymetrix probe set definitions:
  - many probe sets map to the same gene
    - integrative microarray studies use ad-hoc heuristics to integrate these values into a single expression estimate
- updated probe set definitions
  - erroneous or non-specific probes are removed
  - probe sets targeting the same gene/transcript are pooled
    - fewer probe sets compared with the original
    - number of probe pairs per probe set is no longer identical
      - better/worse ?!

# Motivation (2)

- Studies showed:
  - Updated probe set definitions affect approximately 20-30% of all probe sets → affect a large portion of gene estimates!
  - Genes identified as differentially expressed using original and updated probe set definitions only show 50% overlap!
  - Improvement of cross-platform reproducibility of microarray experiments when using updated probe set definitions

## Purpose of this study

- Such updated probe set definitions are now available
  - can easily be integrated into bioconductor packages (affy, gcrma)
  - map the platform probe signals to genes, transcripts and even exon expression levels
- Aim of this study
  - evaluate the impact of updated probe set definitons on precision and <u>accuracy</u> in estimated expression levels

# Study Design

- What?
  - Re-analysis of raw data using updated probe set definitions
  - Comparison of results with original Affymetrix probe set definitions
- Why?
  - Does usage of updated probe set definitions improve precision and accuracy of results?
    - <u>Hypothesis</u>: variable number of probe pairs integrated into each probe set → might have negative impact on precision? (fewer probes in some probe sets)

# Methods (1)

- How?
  - Re-analysis of a gene expression data set using
    - the original NetAffx probe set definitions
    - 6 updated probe set definitions (custom CDFs)
  - Dataset:
    - 2 RNA samples which differ in expression of only a few genes
    - both samples were hybridized to 2 HG-U133A Affymetrix arrays each by 5 different labs

#### Methods (2)

Sample A, B







$ \begin{array}{c}                                     $	Methods (2)				
4		Sample A, B	Hybr 2 arr (HG-	ridized to ays -U133A)	
Lab 1	Lab 2	Lab 3	Lab 4	Lab 5	
A1, A2 B1, B2	A1, A2 B1, B2	A1, A2 B1, B2	A1, A2 B1, B2	A1, A2 B1, B2	Raw data files (CEL)







#### **Experimental Data**

#### • 4 arrays per lab

- 2x sample A (A1, A2)
- 2x sample B (B1, B2)
- computed expression values using 7 different probe set definitions (default vs. updated) and 3 different probe set summarization algorithms
- Within each lab:
  - pairs of replicates used to estimate precision and accuracy
- 5 labs performed identical experiment:
  - estimated precision and accuracy obtained in each lab can be summarized to provide a <u>robust assessment</u> of the effects

# Precision (1)

- Measures data reproducibility and variability
- Precision = correlation between the relative log2 expression ratios of the 2 RNA samples using the 2 pairs of replicate pairs:
  - A1/B1 vs. A2/B2
- Clear indication of experiment performance
  - Correlation of  $1 \rightarrow \text{perfect precision}$
  - Correlation of  $0 \rightarrow$  no precision

#### Precision (2)

• Comparison of precision between original probe set definitions and updated probe set definitions:

	ensEMBL exon	ensEMBL gene	ensEMBL transcript_	Entrez	RefSeq	UniGene
MAS5	-0.014	0.013	0.013	0.014	0.019	0.023
	p = 0.094	p = 0.057	p = 0.0092	p = 0.056	p = 0.018	p = 0.052
RMA	-0.035	0.053	0.030	0.059	0.045	0.047
	p = 0.0041	p = 0.00025	p = 0.0071	p = 0.00011	p = 0.00045	p = 2.9E-06
GCRMA	-0.11	0.023	-0.0063	0.040	0.031	0.028
	p = 0.000051	p = 0.045	p = 0.28	p = 0.019	p = 0.062	p = 0.007 l
Precision better for all updated the probe set definitions but exons $\rightarrow$ why?			lean precision different or updated probe set efinition compared w ne original probe set	ence vith Sign (2-ta	ificance of ailed paired	difference t-test) 33

Table 1: Improved precision using update probe set definitions

## Precision (3)

- Possible reason for decrease in precision for probe set definition to ensEMBLE exons?
  - mean number of probes per probe set is lower for ensEMBL exons
  - using fewer probes when estimating an expression level likely increases variance and lowers precision

#### Precision (4)

- Improved precision for other updated probe set definitions could be due to larger number of probes mapping to each probe set
  - mean number of probes higher than for original probe set definition

Probe set definition	Number of probe sets	Mean number of probe pairs per probe set		
NetAffx (original)	22283	11.1		
ensEMBL exon	35191	9.3		
ensEMBL gene	18671	14.0		
ensEMBL transcript	36174	13.9		
Entrez	12132	14.1		
RefSeq	17880	14.9		
UniGene	11694	15.0		

Table 2: Characteristics of probe set definitions

## Precision (5)

 Precision was analyzed as a function of the number of probes used to estimate each probe set

 Positive correlation between the number of probes per probe set and precision found



# Precision (6)

- But:
  - Updated probe set definitions appear to achieve better precision, even when similar number of probes were integrated into signal estimates!
- → Updated probe set definitions have significant improvements in precision!
  - removal of erroneous or non-specific probes that otherwise add noise

# Accuracy (1)

- Accuracy estimates how close microarry estimates are to the <u>"real expression" changes</u>
  - most often "real" expression is measured using RT-PCR (real time PCR)
- To assess accuracy the updated probe set definitions achieved:
  - Differential expression detected with microarrays was compared to those measured with RT-PCR for 16 genes (for the different probe set definitions)

# Accuracy (2)

- Accuracy was defined as the slope after a linear regression between RT-PCR and microarry data
  - An accuracy of 1.0 is optimal
- Difference in accuracy for each lab between the updated probe set definitions and the standard probe set definition was calculated
  - Test: is mean accuracy difference (averaged across the 5 labs) significant?
    - Using paired t-test (2-tailed distribution)

#### Accuracy (3)

 Significant improvements in accuracy were observed (when data was normalized using RMA) for all but UniGene definition

Mean Slope	p-value	Std
0.74		0.02
0.83	0.00040	0.01
0.83	0.00053	0.01
0.78	0.00430	0.01
0.78	0.00381	0.01
0.75	0.07085	0.01
	The slopes estimated from th 5 different labs ar in good	
	Mean Slope 0.74 0.83 0.83 0.78 0.78 0.75	Mean Slope         p-value           0.74         0.00040           0.83         0.00053           0.78         0.00430           0.78         0.00381           0.75         0.07085

agreement.

Table 3: Improved accuracy using updated probe set definitions

#### Accuracy (4)



# Results (1)

 Significant improvement in <u>precision</u> using updated probe set definitions!



## Results (2)



 Significant improvement in accuracy using updated probe set definitions!

#### Summary

- Updated probe set definitions offer expression levels that are more accurately associated to genes and transcripts.
  - It could be shown in this study that they also improve the estimated transcript levels:
    - using updated probe set definitions improves both the <u>precision</u> and <u>accuracy</u> of the relative expression level estimates
- The results of this study encourage a wide spread use of updated probe set definitions.