# Performance Assessment of RNA Sequencing and Expression Arrays for Transcriptome Analysis in Cancer Research

Petr V. Nazarov

petr.nazarov@lih.lu

Luxembourg Institute of Health

**Part I. Comparison of RNA-seq and microarray performance**
- Similar and specific features of the platforms
- Protein coding and long non coding genes
- Gene expression analysis and analysis of alternative splicing

**Part II. Independent Component Analysis (ICA) in transcriptomics**
- The brief introduction to the method
- Deconvolution of biological signals and cell subtypes
- Potential for patient diagnostics in future

# Part I. Comparison of RNA-seq and microarray performance

**Based on** Nazarov et al *BMC Genomics,* **2017**;18(1):443.

## RNA sequencing and transcriptome arrays analyses show opposing results for alternative splicing in patient derived samples

Petr V. Nazarov[1*], Arnaud Muller[1], Tony Kaoma[1], Nathalie Nicot[1], Cristina Maximo[1], Philippe Birembaut[2], Nhan L. Tran[3], Gunnar Dittmar[1] and Laurent Vallar[1]

Majority of comparisons in literature claim that RNA-seq outperforms microarrays. However, comparing RNA-seq with old 3' microarrays... not too fair.
Currently more advanced arrays are available: HTA and its successor Clariom.



**Clariom D Array & HTA 2.0 (WT)**
*Gene modeling probe set*

**Clariom S Array (WT)**
*Constitutive exon probe set*

**U133/MG430/RG230 (3' IVT)**
*Biased probe set*

- How similar are the results obtained by last version arrays and RNA-seq ?
  - protein coding / other biotypes, genes / exons

- What are the differences between platforms?

- Which platform should one use

Image is provided by ThermoFisher Scientific

## Research includes: 1 cancer, 9 patients, 18 samples, 2 platforms

9 patients with lung squamous cell carcinoma (clinical research study)

18 samples:

tumour

adj. normal

**Affymetrix HTA 2.0 Arrays**

100 ng

**Illumina HiSeq 2000**

500 – 1000 ng

**Unité INSERM, University of Reims**
Prof. Ph. Birembaut

- Total RNA extracted using miRNeasy Mini Kit
- **Arrays:** GeneChip® WT Plus Reagent Kit
- **Sequencing:** TruSeq total RNA Sample Preparation Kit v.1.0, polyA selection

**Sequencing depth**

RNA-seq result: 120-280 M paired reads with 77 bp/read.



**Figure.** Number of mapped reads after RNA-seq analysis of the samples, including *TopHat* alignment. In general, normal tissues (green) show more reproducible mapping results than tumours (yellow).

**Data acquisition**

**Microarray analysis with Affymetrix HTA v.2**

CEL

Partek®GS;
Transcriptome analysis console (TAC)

**GC RMA**  normalization, probe summarization

Probe set expression

**R/Bioconductor**  annotation, re-mapping, exon / gene level

**Gene** expression    **Exon** expression    **Junction** expression

… analysis

**RNA-seq analysis with Illumina HiSeq 2000**

FASTQ

alignment  **TopHat**

BAM

counting

**HTSeq**    featureCount **(R/Bioconductor)**

**Gene** expression    **Exon** expression    **Junction** expression

… normalization and analysis

**Important:**
in order to compare the platforms, we re-mapped Affymetrix probesets onto the Ensembl 69 genome using GenomicRanges package of R.

LUXEMBOURG
INSTITUTE
OF **HEALTH**
RESEARCH DEDICATED TO LIFE

## Overlap of features is high



**Main biotypes (Ensembl 69)**

protein_coding
pseudogene
lincRNA
antisense
miRNA
misc_RNA
snRNA
snoRNA
processed_transcript
sense_intronic
rRNA
other

**Protein coding genes**

HTA    RNA-seq

0    20033    13

20033    20046

**lncRNA genes**

HTA    RNA-seq

0    5855    462

5855    6317

Good overlap of the genes and exons

**Protein coding exons**

HTA    RNA-seq

0    511 866    47 358

511 869    559 224

**lncRNA**

HTA    RNA-seq

0    25 640    2 875

25 640    28 515

**Important:**
in order to compare the platforms, we re-mapped
Affymetrix probesets onto the Ensembl 69 genome
using GenomicRanges package of R.

# Clustering

## Coding genes: removable platform effect

**Linearly scaled/centered data**



- Strong effect of tumour/normal condition
- Platform-specific effect can be reduced by simple centring-scaling (standardization)
- lncRNA show similar behavior with, with higher variability

*Clinical research study*

**Coding genes are more correlated than lncRNA**

| Correlation | coding mRNA | lncRNA |
|---|---|---|
| log signal | 0.76 | 0.319 |
| logFC | 0.743 | 0.349 |



Scatter plots showing general tendency in RNA-seq and HTA protein coding gene expressions (orange) and logFC (green). Scatter plots are built by overlap of all available data for SCC patients.

➢ Correlation for protein coding genes is in range of values reported in literature

➢ lncRNA are not so nicely correlated. Reason? ↓

*Clinical research study*

**Gene length matters!**

Protein coding

lncRNA

LUXEMBOURG
INSTITUTE
OF **HEALTH**
RESEARCH DEDICATED TO LIFE

*Clinical research study*

## **Explained variability in the data: better for HTA**

### **protein coding**

### **lncRNA**



unexplained variability

- HTA show less unexplained variability and higher cancer-associated variability

Principal Variance Component Analysis (PVCA) was described in:

Li, J., Bushel, P., Chu, T.-M., and Wolfinger, R.D. (2009) Principal Variance Components Analysis: Estimating Batch Effects in Microarray Gene Expression Data, Batch Effects and Noise in Microarray Experiments: Sources and Solutions, ed. A. Scherer, John Wiley & Sons.

## St.deviation in biological replicates is higher in RNA-seq



**RNA-seq**
**HTA**

Variability between biological replicates is higher for RNA-seq data for both normal and tumour samples, especially for lowly abundant transcripts

# Differential Expression Analysis

## DE gene lists vs TCGA: similar level of confirmation

LUXEMBOURG INSTITUTE OF HEALTH
RESEARCH DEDICATED TO LIFE

protein coding mRNA
**4 777**
1 853 confirmed in TCGA

**1 094**
255 confirmed

**3 683**
1 598 confirmed

**2 490**
658 confirmed

**6 173**
2 256 confirmed in TCGA

RNA-seq

HTA

lncRNA
**892**
28 confirmed in TCGA

516
11 confirmed

**376**
17 confirmed

843
18 confirmed

**1219**
35 confirmed in TCGA

**TCGA LUSC data series: 502 -vs- 31**

PCA for samples by SCC
(23% variability)

PC2, includes 5% variability

Healthy

Cancer

PC1, includes 18% variability

Jaccard index

log₁₀(FDR)

RNA-seq & HTA
RNA-seq & TCGA
HTA & TCGA

➤ More DEG for HTA with FDR<0.01

➤ Comparing with TCGA – similar confirmation rate

➤ Overlapping genes: 1598 of 3683 are found in the top 25% of TCGA

## How to compare "pears" with "apples"?

We proposed considering only significant genes, in order to make the analysis more fair.

| Measure | RNA-seq | HTA |
|---|---|---|
| Lower limit of log expression | -0.80 | 3.83 |
| Higher limit of log expression | 9.20 | 8.89 |
| **Dynamic range of log expression** | **10.00** | **5.06** |
| Lower limit of absolute logFC | 0.67 | 0.17 |
| Lower limit of absolute logFC | 7.55 | 3.58 |
| **Dynamic range of absolute logFC** | **6.87** | **3.41** |

*Values are in log$_2$*

➤ As expected, dynamic range of RNA-seq is higher. But taking into account that HTA allow for detecting genes with smaller fold change - it still can be related to difference in scales.

# Prediction Analysis

## More predictive genes were observed with arrays

Area under ROC curve (AUC) characterizes applicability of a gene to distinguish between 2 groups of samples and, therefore, tells whether a gene can be used as a marker to predict the group.



**AUC – area**
min 0.5, max = 1

➤ AUC constantly shows better values for HTA data

# Gene Set Analysis

## Biological processes (GO:BP) enriched with DE genes

| DE genes (FDR<1e-4) | → | Fisher-based enrichment (FDR<1e-2) | → | ReViGo semantic clustering |

**topGO** package of R/Biocondictor

### biological processes



83    **84**    126

**RNAseq**    **HTA**

-Σlog(FDR) > 100
  -Σlog(FDR) > 10
    -Σlog(FDR) > 2

**RNAseq**
- tissue development
- collagen catabolism
- extracellular matrix organization
- positive regulation of mitotic cell cycle
- cellular component movement
- developmental process
- single-organism cellular process
- single-organism process
- cell proliferation
- multicellular organismal process
- reproduction
- response to alcohol

**common**
- cell cycle process
- cilium organization
- DNA metabolism
- microtubule-based movement
- microtubule-based process
- cell cycle
- cellular component organization or biogenesis
- cell division
- chromosome segregation
- regulation of cell division
- anatomical structure homeostasis
- protein localization to chromosome
- response to ionizing radiation

**HTA**
- protein-DNA complex assembly
- DNA integrity checkpoint
- cellular response to DNA damage stimulus
- RNA transport regulation of ligase activity
- epithelial cilium movement involved in determination of left/right asymmetry
- single-organism metabolism

➢ GO:BP biases are found: extracellular in RNA-seq , DNA-related in HTA
➢ More GO:BP in with HTA analysis

## Cellular components (GO:CC) enriched with DE genes

| DE genes (FDR<1e-4) | → | Fisher-based enrichment (FDR<1e-2) | → | ReViGo semantic clustering |

**topGO** package of R/Biocondictor

### cellular components



17   **63**   40

**RNAseq**   **HTA**

$-\Sigma log(FDR) > 100$
$-\Sigma log(FDR) > 10$
$-\Sigma log(FDR) > 2$

### RNAseq
- proteinaceous extracellular matrix
- extracellular region ciliary tip
- extracellular matrix
- cell-cell junction
- cornified envelope
- intraciliary transport particle
- chaperonin-containing T-complex
- collagen trimer
- intraciliary transport particle B

### common
- microtubule cytoskeleton
- cilium
- extracellular vesicular exosome
- non-membrane-bounded organelle
- organelle part
- membrane-enclosed lumen
- organelle lumen
- organelle
- protein complex
- cytosol
- cytoplasm
- macromolecular complex
- cell projection
- proteasome accessory complex
- vesicle
- midbody
- MCM complex
- desmosome

### HTA
- nucleoplasm
- intracellular part
- intracellular
- intracellular organelle
- membrane-bounded organelle
- DNA packaging complex
- protein-DNA complex
- DNA bending complex
- DNA polymerase complex
- cell
- cell part
- pore complex
- proteasome complex
- envelope

➢ GO:CC biases are found: extracellular in RNA-seq , nucleus in HTA
➢ More GO:CC in with HTA analysis, again

# Gene Set Analysis

## Bias can be linked to RNA abundance



**Figure S6.** Expression of the genes related to cellular component ontologies uniquely identified by RNA-seq (red lines) and HTA (blue lines). The distributions of gene expressions are based on sequencing (A) and microarray (B) data. Both data agree, that genes participating in the functions uniquely found in RNA-seq analysis show higher expression than one of HTA analysis (yellow area).

➢ Abundance of the genes participation in extracellular biofunctions is higher then for nucleus-related genes.

➢ Small bias of the length was seen as well, but it cannot explain the expression differences: checked with *goseq* package (correcting for gene length)

➢ Strong bias is seen only for CC. Only minor for BP

- Linear models are used
- HTA: **DiffSplice** from **limma** package
- RNA-seq: **DEXSeq**

**Challenge:** HTSeq tool does not work for exons – too many overlapped entities (correlation b/w platforms ≈ **0.2**)

**Solution:** Changing counting tool to *featureCount (Rsubread)* improved concordance b/w HTA and RNA-seq: correlation ≈ **0.6-0.7**

# Analysis of Splicing Events

## Low concordance of the results

**a** Based on exon expression

23 934

RNA-seq

20 236

3 698

HTA

23 301

26 999

**b** Based on junction expression

7 063

5 512

1 551

38 833

40 384

protein coding mRNA
4 777
1 853 confirmed in TCGA

RNA-seq

1 094
255 confirmed

HTA

3 383
1 598 confirmed

2 490
658 confirmed

6 173
2 256 confirmed in TCGA

LUXEMBOURG
INSTITUTE
OF HEALTH
RESEARCH DEDICATED TO LIFE

## The 3'-exons and long exons show-up in RNA-seq

The exon parameters distribution among differentially used exons detected by the two platforms



The relative position of the exons within their genes, varying from 5' end (relative position = 0) to 3' end (relative position = 1), shows a 3' bias in RNA-seq (**a**).

Exon length shows that RNA-seq tends to find more significantly splice events among long exons than HTA (**b**).

## 3' bias or length-related bias?

The RNA-seq data show tendency to increase expression at 3'-end...



(ENSG00000049759) (NEDD4L)

chr18:55711599−56068772(+)

HTA − DiffSplice − FDR< 0.05 & abs(log2FC) >= 1.5

HTA

RNASeq (DEXSeq norm. count) − DEXSeq − FDR< 0.05 & abs(log2FC) >= 1.5

RNA-seq

3'

(ENSG00000018408) (WWTR1)

chr3:149235022−149454501(−)

HTA − DiffSplice − FDR< 0.05 & abs(log2FC) >= 1.5

HTA

3'

RNASeq (DEXSeq norm. count) − DEXSeq − FDR< 0.05 & abs(log2FC) >= 1.5

RNA-seq

Length

5'    p5    ce    p3    3'

Abs. lcover

t5    p5    ce    p3    t3

Probably 2 effects play role: the length of 3' exon and poly-A selection. The length bias cannot explain 100% of expression bias

- In our study, HTA showed more reliable results than RNA-seq with 200M reads.

- Length sensitivity makes RNA-seq a difficult technique for non-coding RNA and requires high coverage.

- RNA-seq is very good as a discovery tool!

- Be careful when doing isoform study with any platform!

# Part II. Independent Component Analysis in Transcriptomics

In collaboration with
**Dr. Anke Wienecke** and **Dr. Stephanie Kreis,**
Life Science Research Unit, University of Luxembourg

What did James say?..

## Cell ensemble is as well a "cocktail party"



Endothelial cells

Cancer cells

Normal cells

Invasive cancer cells

Fibroblasts

Immune cells

Hanahan D, Weinberg RA. *Cell* **2011**, 144, 646-74

## The method to solve it…



**I**ndependent
**C**omponent
**A**nalysis

LUXEMBOURG
INSTITUTE
OF **HEALTH**
RESEARCH DEDICATED TO LIFE

## Independent Component Analysis

### Deconvolution of Cell Ensemble

Translational
research study:

**Patient 1**

**Patient 2**

**Patient 3**

**Patient 4**

**Patient 5**

**Original data**
samples

genes

**Metagenes**

genes

components

components

**One
component**

involvement

genes

Can be linked to
biological processes
and cell subpopulations

**Weights of components**

components

samples

**Components
weights in
patients**

Can be linked to patient
groups and survival

Captures & cleans
batch/platform effect

$$X_{gs} \approx S_{gk} \times M_{ks}$$

adapted from Hanahan D,
Weinberg RA. *Cell* **2011**, 144,
646-74

A. Biton et al, Cell Reports 9, 2014
A. Zinovyev et al, Biochem Biophys Res Commun. 2013

## What ICA does and does not

$$X_{gs} \approx S_{gk} \times M_{ks}$$

$g$ – genes
$s$ – samples
$k$ - components

### Pro:
1. Finds **statistically-independent signals** (components) in the expression profiles
2. Identifies the **most important genes** in each component
3. Tells what is the weight of **each component in the samples**
4. Works on data *per se*, **without any additional knowledge**
5. Gives quite **robust answer**… just… reshuffled

### Contra:
1. **No ranking of the components** by importance (not like PCA)
2. Results are **not deterministic** and can to some extent depends on the run
3. **Orientation of the signal is arbitrary** from one run to another
4. If you look for precise estimation of cell fraction – not a good idea (results are qualitative not quantitative)

## Positive and Negative Genes within Components



**Figure S6.** (A) Number of significant positively (red) and negatively (blue) involved genes in metagene of each of the components. (B) Number of enriched GO biological processes found for these genes. For the most cases, only one list of genes is biologically meaningful: either positive (e.g. ic10-ic15) or negative (e.g. ic25, ic28, ic49, ic55).

**LUXEMBOURG INSTITUTE OF HEALTH**
RESEARCH DEDICATED TO LIFE

## ICA for patient classification

*clinical research study*



We use **parallel consensus ICA** that provides quite **robust estimation of the matrices** (based on fastICA package in R)

## **Optimal measure for RNA-seq**



**Raw count**
**DESeq norm**
**FPKM**
**TPM**

## Patient classification in SKCM

SKCM

(skin cutaneous melanoma)

472 samples

5 samples



- SVM & RF work both fine when $n_{comp}$ is small
- For large $n_{comp}$ – RF gives much better predictions (SVM is overtrained)

| Gender | | |
|---|---|---|
| Accuracy | Actual gender | |
| 99.6% | female | male |
| female | 177 | 0 |
| male | 2 | 293 |

| Type | | |
|---|---|---|
| Accuracy | Actual sample type | |
| 78.9% | metastatic | primary |
| metastatic | 177 | 54 |
| primary | 7 | 51 |

| Cluster | | | |
|---|---|---|---|
| Accuracy | Actual cluster | | |
| 90.0% | immune | keratine | MITF-low |
| immune | 160 | 9 | 6 |
| keratine | 9 | 91 | 6 |
| MITF-low | 1 | 2 | 47 |

Here accuracy was estimated using LOOCV

## New samples: mRNA and miRNA

### mRNA level: RNA-seq + RNA-seq

Gender: ● female, ■ male

Sample type:
- ●■ primary tumour
- ●■ metastatic
- ●■ new samples



When ICA is run over new samples and training samples together, it corrects for platform bias.

## New samples: mRNA and miRNA

Gender:  ● female,  ■ male

Sample type:
● ■ primary tumour
● ■ metastatic
● ■ new samples

### miRNA level: RNA-seq + qPCR



miR-146a-3p   miR-205-5p
miR-338-5p   miR-199b-5p
miR-551b-3p   miR-876-5p
miR-598-3p   miR-1266-5p
miR-206   miR-301b-3p
miR-34a-5p   miR-3690
miR-338-3p   miR-365a-3p
miR-146a-5p   miR-125b-1-3p
miR-1269a
miR-573

logtest pv=9.4e−04
LHR=−1.79 (CI = −2.82, −0.75)

MIC1: unknown segregation of samples

When ICA is run over new samples and training samples together, it corrects for platform bias.

## ICA can be used for data integration

**Correlation of weights: mRNA-miRNA-Proteins**



$$score_j = \sum_{i=1}^{k} d_i R_i^2 M_{i,j}^*$$

## ICA helps establishing scores for new samples

$$score_j = \sum_{i=1}^{k} d_i R_i^2 M_{i,j}^*$$

$d_i$ – direction of the component (pos/neg)
$H_i$ – log-hazard of Cox regression
$R_i^2$ – stability of the $i$-th component
$M_{i,j}^*$ – weight of $i$-th component in sample $j$

$$hscore_j = \sum_{i=1}^{k} H_i R_i^2 M_{i,j}^*$$

**logtest pv=1.2e-13**
**LHR=1.08 (CI = 0.79, 1.37)**

- We tested our implementation of consensus ICA
  (before publication, the script is available upon request)

- ICA decomposes large bulk data set into meaningful signals

- New samples are properly mapped in IC-space

- The method allows classifying and scoring new patients
  (clinical research studies)

# Acknowledgements