

1 DATA AGGREGATION AT THE LEVEL OF MOLECULAR PATHWAYS IMPROVES
2 STABILITY OF EXPERIMENTAL TRANSCRIPTOMIC AND PROTEOMIC DATA

3 Nicolas Borisov^{1,2*}, Maria Suntsova^{3,4}, Andrew Garazha^{4,5}, Ksenia Lezhnina⁵, Olga Kovalchuk^{6,7},
4 Alexander Aliper^{1,5}, Elena Il'nitskaya¹, Maxim Sorokin^{1,2}, Mikhail Korzinkin^{1,3}, Vyacheslav
5 Saenko⁸, Yury Saenko⁸, Dmitry G. Sokov⁹, Nurshat M. Gaifullin^{10,11}, Kirill Kashintsev¹², Valery
6 Shirokorad¹², Irina Shabalina¹³, Alex Zhavoronkov¹⁴, Bhubaneswar Mishra¹⁵, Charles R. Cantor¹⁶,
7 Anton Buzdin^{1,2,3,4}.

8
9 ¹ First Oncology Research and Advisory Center, Moscow, 117997, Russia

10 ² National Research Centre "Kurchatov Institute", Centre for Convergence of Nano-, Bio-,
11 Information and Cognitive Sciences and Technologies, 1, Akademika Kurchatova sq., Moscow,
12 123182, Russia

13 ³ Pathway Pharmaceuticals, Wan Chai, Hong Kong, Hong Kong SAR

14 ⁴ Group for Genomic Regulation of Cell Signaling Systems, Shemyakin-Ovchinnikov Institute of
15 Bioorganic Chemistry, Moscow, 117997, Russia

16 ⁵ Laboratory of Bioinformatics, D. Rogachyov Federal Research Center of Pediatric Hematology,
17 Oncology and Immunology, Moscow, 117198, Russia

18 ⁶ Department of Biological Sciences, University of Lethbridge, 4401 University Drive, Lethbridge,
19 AB, T1K 3M4, Tel: 1-403-394-3916

20 ⁷ Canada Cancer and Aging Research Laboratories, Lethbridge, AB T1K7X8, Canada

21 ⁸ Technological Research Institute S.P. Kapitsa, Ulyanovsk State University, Leo Tolstoy str. 42,
22 Ulyanovsk, 432000, Russia

23 ⁹ Moscow 1st Oncological Hospital, 17/1 Baumanskaya str., Moscow 105005, Russia

24 ¹⁰ Lomonosov Moscow State University, Faculty of Fundamental Medicine, Lomonosovsky Dr., 31-
25 5, Moscow 119192, Russia

26 ¹¹ Russian medical postgraduate academy, Barrikadnaya str., 2/1, Moscow, 123995, Russia

1 ¹² Moscow Oncological Hospital 62, Istra 27, Stepanovskoye, Krasnogorsk region, 143423, Russia

2 ¹³ Petrozavodsk State University, Faculty of Mathematics and Information Technologies, Anohina

3 str., 20, Petrozavodsk 185910, Russia

4 ¹⁴ Insilico Medicine, Inc, ETC, Johns Hopkins University, Baltimore, MD, 21218, USA

5 ¹⁵ Courant Institute, New York University, NY, 10012, USA

6 ¹⁶ Department of Biomedical Engineering, Boston University, Boston, Massachusetts, United States

7 of America

8

9 *The corresponding author

1

2 ABSTRACT

3

4 High throughput technologies opened a new era in biomedicine by enabling massive analysis of
5 gene expression at both RNA and protein levels. Unfortunately, expression data obtained in
6 different experiments are often poorly compatible, even for the same biological samples. Here,
7 using experimental and bioinformatic investigation of major experimental platforms, we show that
8 aggregation of gene expression data at the level of molecular pathways helps to diminish cross- and
9 intra-platform bias otherwise clearly seen at the level of individual genes. We created a
10 mathematical model of cumulative suppression of data variation that predicts the ideal parameters
11 and the optimal size of a molecular pathway. We compared the abilities to aggregate experimental
12 molecular data for the five alternative methods, also evaluated by their capacity to retain
13 meaningful features of biological samples. The bioinformatic method OncoFinder showed optimal
14 performance in both tests and should be very useful for future cross-platform data analyses.

15

16

17

1 INTRODUCTION

2

3 Next generation sequencing (NGS), Microarray hybridization (MH) and high throughput proteomic
4 techniques opened a new era in biomedicine by enabling large scale analysis of gene expression at
5 both the RNA and protein levels [Kumar, 2016]. Multiple experimental platforms based on different
6 principles and utilizing different reagents were developed for these tasks [Kumar, 2016]. According
7 to the International Aging Research Portfolio, over eight billion dollars in government funding have
8 been spent on research projects involving high throughput gene expression analysis since 1993
9 [Zhavoronkov, 2011]. This resulted in tens of thousands of publications. Unfortunately, gene
10 expression data obtained using different experimental platforms are poorly compatible with each
11 other even when obtained using the same biosamples. For example, a generally weak correlation
12 between NGS and microarray gene expression data has been reported [Buzdin, 2014]. Therefore, a
13 new data processing method is badly needed to enable data harmonization among different
14 platforms and experiments [MAQC Consortium, 2006; Zhang, 2013].

15 Recently we showed that aggregation of gene expression data into molecular pathways, each
16 containing dozens or hundreds of gene products, may help to solve the problem of poor data
17 compatibility among different experimental platforms [Buzdin 2014a]. NGS and microarray data
18 obtained for the same transcripts showed generally low correlation (<0.2) when examined at the
19 level of individual genes. However, these correlations improved dramatically, up to 0.9, when
20 activation of 90 molecular pathways was analyzed instead [Buzdin, 2014a]. The output measure
21 was a Pathway Activation Strength (*PAS*), which positively reflects the degree of pathway
22 activation. The *PAS* makes it possible to interrogate, quantitatively, processes such as molecular
23 signaling, metabolism, DNA repair and cytoskeleton reorganization, based on gene expression
24 data. These processes determine cell fate by governing growth, differentiation, proliferation,
25 migration, survival and death [Diderich, 2016; Zhavoronkov, 2014]. Molecular modeling of
26 intracellular pathways has been carried out for more than two decades [Kholodenko, 1999;

1 Hanahan, 2000]. A plethora of molecular pathways have been discovered and catalogued, each
2 containing different numbers of gene products [Haw, 2012; Nakaya et al., 2013]. Pathway
3 activation strength was also found to be a better marker of human tissue types [Borisov, 2014;
4 Lezhnina, 2014] and tumor response to chemotherapy treatment [Zhu 2015; Venkova, 2015;
5 Artemov, 2015]. Several approaches were published by us and others to assess the activation of
6 signaling pathways, basing on large scale molecular data [Khatri, 2012; Buzdin, 2014b,
7 Zhavoronkov, 2014]. These methods take into account different factors like the extent of
8 differential gene expression, architecture of molecular pathways, and the roles of individual gene
9 products in a pathway (e.g., activator/repressor) [Khatri, 2012; Buzdin, 2014b]. For example, a
10 method we used to minimize discrepancies between the NGS and microarray platforms, termed
11 OncoFinder, relies on differential gene expression and the known roles in a pathway, but does not
12 take into account pathway architecture, i.e. the position of a gene product in a pathway [Buzdin,
13 2014b].

14 In spite of this progress it is not known why data aggregation improves expression information
15 stability and what factors influence it. It is also unclear which bioinformatic algorithms provide
16 better *PAS* outputs for cross-platform data stability. Additionally, *PAS* algorithms have not yet been
17 applied to the high throughput proteomic data.

18 In this study, we applied data aggregation methods to transcriptomic information obtained using
19 the Affymetrix HG U133 Plus 2.0, the Illumina HT12 bead array, the Agilent 1M array, the Illumina
20 Genome Analyzer platforms, and to proteomic data from the Orbitrap Velos and XL mass
21 spectrometer platforms. We confirmed that for both transcriptomic and proteomic expression levels,
22 the *PAS* approach provided more stable results than the expression of individual genes. To explain
23 this phenomenon, we created a biomathematical model simulating error acquisition in individual
24 gene expression and in *PAS*-based approaches. In agreement with the experimental data, in the
25 mathematical model *PAS* methods produced significantly more stable results under a majority of
26 conditions. This model also predicts the optimal size of a molecular pathway and ideal parameters

1 of the normalizing (control) set of gene expression data.

2 To test the predictions further of the biomathematical model, we designed a new experimental
3 gene expression array using the CustomArray microchip platform (USA) enabling direct
4 electrochemical synthesis of oligonucleotide probes on a blank array. We compared results for the
5 seven human kidney cancer tissue samples independently profiled by the two laboratories on the
6 this customized array and on the commercial Illumina HT12 bead array platform. In agreement with
7 the theoretical model, gene expression features differed significantly among the platforms for the
8 same biosamples, while *PAS* values remained highly correlated. Therefore, gene expression data
9 aggregated at the *PAS* level appears to be the method of choice for cross-platform data comparisons,
10 including both transcriptomic and proteomic approaches.

11 We next explored the capacity of five most popular *PAS* calculation methods, OncoFinder
12 [Buzdin, 2014b], TAPPA (Topology analysis of pathway phenotype association) [Gao, 2007],
13 Topology-Based Score (TBScore) [Ibrahim, 2012], Pathway-Express [Draghici, 2007], and SPIA
14 (Signal pathway impact analysis) [Tarca, 2009] to generate stable and biologically relevant data.
15 We used the MicroArray Quality Control (MAQC) dataset [MAQC Consortium, 2006] including
16 expression data for four biological samples profiled in fifteen replicates on major commercial
17 microarray platforms. The abilities of the various *PAS* methods to increase correlation between
18 transcriptomic features of the same biosamples examined using different experimental platforms
19 were tested. We also checked whether the *PAS* methods were able to retain biological features after
20 data harmonization using a generally accepted cross-platform harmonization procedure XPN
21 [Shabalina, 2008]. We found that the OncoFinder method showed the optimal performance in both
22 tests.

23

24

25

26

1 RESULTS

2

3 **Cross-platform processing of transcriptomic and proteomic data**

4 We processed transcriptomic and proteomic data to establish pathway activation strength (PAS)
5 profiles corresponding to intracellular molecular pathways. The analysis included 271 molecular
6 pathways (Supplementary dataset S1). For PAS measurements, we applied the OncoFinder method
7 which was previously shown to diminish the cross-platform variation between the MH and NGS
8 data [Buzdin, 2014a]. OncoFinder has previously been applied to many human and non-human
9 systems including cell culture, leukemia and solid cancers, fibrosis, asthma, Hutchinson Gilford and
10 Age-Related Macular Degeneration Disease [Makarev, 2016; Artcibasova, 2016; Alexandorva,
11 2016; Lebedev, 2015]. The PAS for a given pathway (p) is calculated as:

12 $PAS_p = \sum_n ARR_{np} \cdot \log(CNR_n)$ [Buzdin, 2014b], where the functional role of the n^{th} gene product

13 in the pathway is indicated by the *activator/repressor role* (ARR), which equals 1 for an activator, -
14 1 for a repressor, and intermediate values -0,5; 0,5 and 0 for gene products having intermediate
15 repressor, activator, or unknown roles, respectively. The CNR_n value (*case-to-normal ratio*) is the
16 ratio of the expression level of gene n in the sample under investigation to the average expression
17 level in the control samples. A positive PAS value indicates activation of a pathway, and a negative
18 value indicates repression.

19

20 **Building pathway activation profiles and assessment of batch effects**

21 To identify if the OncoFinder technique may improve gene expression analysis by eliminating
22 batch effects, we profiled a set of human clinical bladder cancer tissue samples using the same
23 experimental platform (Illumina human HT 12 v4 bead arrays) in two different laboratories. We
24 investigated gene expression profiles generated from 17 bladder cancer samples and seven normal
25 bladder tissue samples. Eight cancer and four normal samples were analyzed in Dr. Kovalchuk's
26 laboratory in Lethbridge (Canada), and nine cancer and three normal bladder tissue samples were

1 analyzed in Dr. Buzdin's laboratory in Moscow (Russia). The gene expression data were deposited
2 in the GEO database (<http://www.ncbi.nlm.nih.gov/geo/>) with accession numbers GSE52519 and
3 GSE65635.

4 In agreement with previous reports [Lazar, 2013], the normalized gene expression showed
5 significant batch effects with data from different laboratories clearly clustered on a Principal
6 Component Analysis (PCA) plot (Fig.1A). However, the *PAS* data formed a single merged cluster
7 (Fig.1B). The principle component variability was 4-6 times smaller for the *PAS* data (Fig.1A,B).

8 Similarly, using *PAS* values these two sets of samples formed mixed groups on a
9 hierarchical cluster heatmap (Fig.1C). The Canadian samples were labeled 55 - 72; the Russian
10 samples X1 - X8. Some sub-clusters are evidently formed by the samples coming from the different
11 sets, e.g. by samples X5, X8, 69, 68 and X1. (Fig.1C). These data show that data aggregation at the
12 *PAS* level is sufficient to suppress the batch effect in gene expression comparisons.

13

14 **Mathematical modeling of data aggregation effects**

15 We investigated the hypothesis that the apparently higher robustness of OncoFinder *PAS* scoring
16 compared to single gene expression, is due to the cumulative nature of the former. *PAS* is the sum
17 of multiple mathematical terms that correspond to each individual gene product participating in a
18 pathway. Model calculations showed that this cumulative effect is able to reduce stochastic noise.

19 In the model, we included 271 pathways with variable numbers of gene products. We
20 assumed that the expression level of every gene product could be measured using two different
21 methods, say X and Y, corresponding to different experimental platforms (e.g. MH and NGS). Each
22 method introduces errors into the determination of gene expression level, and these errors are
23 independent. A Monte Carlo trial was performed as follows: we simulated both *biased CNR* (with a
24 median value of 1.5) and *unbiased CNR* with a median value of 1. We explored both *noisy* and
25 *exact* expression profiling methods, to allow whether measurement procedures introduce errors in
26 the *true* expression values. The four scenarios of the stochastic simulations (labeled A to D) are

1 shown in Table 1.

2 For each scenario, we calculated the benefit ratio $R = \frac{C_p}{C_g}$, where C_p and C_g are the
3 correlation coefficients between the results obtained using methods X and Y, using *pathway*-based
4 (*PAS*), and *individual* gene product-based log *CNR* values, respectively. For each subset of genes in
5 a pathway, we performed 100 Monte Carlo stochastic simulations and then computed the mean
6 values of C_p and C_g using the R statistical package. The greater $R > 1$, the higher the benefit from
7 using *PAS* instead of individual gene expression for the cross-platform comparisons; $R < 1$ means
8 operating at the individual gene product level is better than the *PAS* level.

9 For *biased* expression profiles, scenarios A and B of Table 1, (Fig. 2), the *PAS* method
10 shows much better agreement between the results obtained using different methods, compared to
11 the individual gene expression levels. The data aggregation advantage of *PAS* is especially strong
12 when both expression methods are *noisy* (scenario A). In scenario B, when one method is *exact*, the
13 benefit of pathway data aggregation is lower. This is caused mainly by higher expression
14 correlation already at the level of individual gene products (Fig. 3). However, the advantages of
15 *PAS* remain considerable for pathways that contain at least 10 gene products (Fig. 2). For shorter
16 pathways, the data aggregation effect is gradually decreased, and the *R* ratio reflecting the benefit of
17 using *PAS* values, trends towards 1.

18 For *unbiased* transcription profiles, with median relative gene expression levels equal to 1,
19 the data aggregation effect is completely lost (scenarios C and D). Here, the mean value for each
20 gene product component of the *PAS* score is zero; consequently, the expected *PAS* is also zero, and
21 the relative data variation is the same at the gene product and the *PAS* level.

22 The simulations clearly elucidate how the cumulative nature of *PAS* suppresses cross-
23 platform data variation and batch effects. They show that there is a significant advantage of using
24 *PAS* to compare platforms, when at least one is *noisy*. This should apply to most if not all existing
25 high throughput experimental platforms, and it should be seen when experimental expression data is
26 compared. The simulations demonstrate that *PAS* calculations are advantageous for *biased*

1 transcriptomes and proteomes and virtually useless for *unbiased* ones. Unbiased data sets are too
2 similar to the control group used as the reference to calculate *CNR* values. This means that the *PAS*
3 approach will be especially useful when the expression signature in the sample under study is very
4 different from that of the control samples. This finding may help to identify appropriate control
5 samples for decreasing expression data noise. Finally, this model shows that the higher is the
6 number of gene products in a pathway, the greater the benefit of shifting from individual
7 gene/protein expression to *PAS* data. For example, the mean number of gene products in the
8 OncoFinder database is 68 per pathway, and the model predicts about a 4.5 –fold decrease in data
9 variation at the *PAS* level in the biased noisy-noisy scenario, which may explain the success of the
10 OncoFinder approach in various applications [Buzdin, 2014a].

11

12 **Experimental model of cross-platform comparisons.**

13 In transcriptomic methods, batch effects arise from errors introduced at the stages of RNA
14 purification, library preparation and amplification, hybridization and reading of arrays [Risso,
15 2011]. We investigated whether the OncoFinder *PAS* algorithm can suppress batch effects
16 introduced by cross-platform comparisons. At the same time we assessed if the algorithm works
17 efficiently for formalin-fixed, paraffin-embedded (FFPE) tissue samples. Seven FFPE tissue blocks
18 isolated from human renal carcinomas were profiled using two independent experimental platforms.
19 The first was the Illumina HT 12 v4 bead array system optimized for FFPE tissues. The second was
20 a customized microchip system developed using the CustomArray (USA) technology of direct on-
21 chip electrochemical oligonucleotide synthesis. The custom arrays had 3775 oligonucleotide probes
22 corresponding to 2214 human gene products involved in 271 intracellular signaling pathways
23 (Supplementary dataset S1). The custom arrays, used the original oligonucleotide probe sequences
24 of the Illumina HT 12 v4 platform, but shortened by 5 nucleotides at the 5' end and by 5 nucleotides
25 at the 3' end. Quantile –normalized gene expression data were deposited into the GEO database
26 with the accession numbers GSE65637 and GSE65639. The differences between the Illumina and

1 the Custom platforms included shorter oligonucleotide probe sequences, different library
2 preparation protocols and different hybridization signal development and reading methods
3 (Supplementary Fig.1). The Custom method for library preparation was quite distinct from Illumina
4 and identical to that used by the Agilent MH platform (Supplementary Fig.1B,C,E) with the sole
5 exception that biotinylated rather than fluorescently labeled DNA is used at the terminal stage
6 (Supplementary Fig.1 B,E). A brief comparison of the protocols used for the Custom and top
7 commercial MH platforms manufactured by the Illumina, Agilent and Affymetrix companies for
8 FFPE tissue profiling is given in Supplementary dataset S2.

9 To compare with the renal carcinoma samples, we used GEO dataset GSE49972 [Karlsson,
10 2014] containing 6 normal kidney samples to normalize the expression data and calculate *PAS*. The
11 normalized *CNR* expression data and *PAS* values are shown in Supplementary dataset S1. At the
12 level of individual gene products, we observed relatively low correlations (0.2-0.3) between the
13 same transcriptomes profiled using the two platforms (Fig. 4; Supplementary dataset S3). In
14 contrast, at the *PAS* level the correlations were strong, varying from 0.84 to 0.91 (Fig. 4;
15 Supplementary dataset S3).

16 These results experimentally confirm the hypothesis that data aggregation at the *PAS* level
17 increases the stability of cross-platform expression data and that the advantage of *PAS* is retained
18 for FFPE samples.

19

20 **Data aggregation effects assessed on different RNA and protein expression profiles**

21 We investigated quantitative aspects of the effect of data aggregation on several datasets where the
22 same samples were profiled using different expression platforms (Tab. 2, Supplementary dataset
23 S4).

24 We observed two trends for the behavior of the benefit ratio $R = \frac{C_p}{C_g}$. In model

25 calculations, we observed a crucial role of expression profile bias between the *case* and *normal*
26 samples for successful data aggregation of genes into pathways (Fig. 2, 3). We introduce a measure

1 of such bias, termed $\beta = \min\left(\frac{|\mu_1|}{\sigma_1}, \frac{|\mu_2|}{\sigma_2}\right)$, where μ_i and σ_i are the mean and standard deviation,
2 respectively, of the set of log *CNR* values obtained for a given sample using the experimental
3 platform *i*. The results of the model calculation (Fig. 2,3, scenarios A and B) suggest that, even for
4 the same values of β , *R* may be different depending on C_g (correlation at the individual gene
5 product level): the higher C_g , the lower *R* at equal β .
6 With a discrimination threshold for C_g chosen as equal to 0.25 between *low-correlated* and the
7 *considerably correlated* samples, we can see the clear clusters of data for data aggregation effect
8 (Fig. 5, blue dots for *low* and red dots for *considerably* correlated samples. Note that the two
9 clusters of data depending on the C_g threshold are seen for both transcriptome-to-transcriptome and
10 transcriptome-to-proteome comparisons.

11 The data obtained suggests that when β is low, the *R* is hardly distinguishable from 1;
12 however, when β exceeds a threshold, the increase of *R* becomes statistically significant. Finally,
13 these results also demonstrate that transcriptomic and proteomic profiles demonstrate more
14 compatible results at the molecular pathway level rather than on the level of individual gene
15 products.

18 **Comparison of *PAS* scoring methods according to their capacities in data aggregation**

19
20 We compared the abilities of five popular *PAS* scoring methods to yield an advantageous
21 the data aggregation effect when the expression of molecular pathways is compared instead of
22 individual gene products. For the seven renal carcinoma samples discussed above, we calculated *R*
23 using alternative *PAS* scoring methods: OncoFinder [Buzdin, 2014b], *topology analysis of pathway*
24 *phenotype association*, TAPPA [Gao, 2007], *topology-based score* (TB) [Ibrahim, 2012], *pathway-*
25 *express* (PE) [Draghici, 2007], and *signaling pathway impact analysis* (SPIA) [Tarca, 2009]

1 methods (Supplementary data set S5). These methods differ in the factors used to evaluate the
2 importance of distinct gene products in pathway activation.

3 Only three of the methods, OncoFinder, PE and SPIA, showed a substantial data aggregation
4 effect (R) ranging from 2-2.3. Other methods showed lack of any positive effect (Fig. 6).

5

6 **Different methods for *PAS* scoring show different properties in retention of biological features**

7 Cross-platform data comparison has the potential to become an extremely useful tool in
8 contemporary biomedicine and bioinformatics. Although the application of *PAS* methods has the
9 ability to restore correlations between different expression data sets, the absolute values of *PAS* may
10 differ between platforms. To overcome this inconsistency, several *cross-platform harmonization*¹
11 methods can be applied ranging from the simplest z-scaling and mean-centering to more
12 sophisticated algorithms utilizing machine-learning/Bayesian harmonization (e.g., [Warnat, 2005;
13 Shabalín, 2008; Hsu, 2014], including the popular harmonization technique XPN [Shabalín, 2008] .
14 In many applications these harmonization methods can diminish the systematic bias introduced by
15 the experimental methods and devices used, but they demonstrate lower efficiencies for routine
16 batch effects like those observed when comparing results obtained using the same platform but on
17 different calendar dates or in different laboratories.

18 This made it of interest to compare the ability of the five *PAS* scoring methods to retain biological
19 features after cross-platform data harmonization with the XPN method.

20 We used the results of the Microarray quality control project (MAQC) [MAQC Consortium,
21 2006] as a model dataset for this study. The MAQC project investigated four types of samples (A-
22 D; each sample profiled in 15 technical replicates) using different microarray devices. Type A
23 samples were taken from the Stratagene Universal Human Reference RNA; type B samples – from
24 the Ambion Human Brain Reference RNA. Type C and D samples were obtained by combining
25 samples A and B in mass ratios 75:25 for C, and 25:75 for D, respectively.

¹ In the current paper, we apply the term *normalization* to any method for *within-platform* batch effect elimination, and *harmonization* when such procedure is performed for the *cross-platform* comparison, although the mathematical methods for both the former and latter procedures may be different.

1 After XPN harmonization of gene expression profiling using the Agilent Whole Human
2 Genome Oligo and Affymetrix Human Genome U133 Plus 2.0 platforms, we applied different
3 methods of *PAS* scoring (Supplementary dataset S6) using the samples of type A as *normal*. The
4 probability densities of the Euclidean distances between the *PAS* vectors calculated for the three
5 samples (B, C, and D) differ greatly depending on the *PAS* scoring method used (Fig.7). In such an
6 assay, an ideal *PAS* scoring method should make distinctions between samples depending primarily
7 on the sample types, rather than on the experimental platform used. A satisfactory *PAS* calculation
8 method, therefore, should yield a unimodal distribution of the *PAS-PAS* distances, without any
9 significant deviations. If the distribution of *PAS-PAS* distances is bimodal or multimodal, this points
10 to the inability to eliminate platform-specific bias even at the pathway level. Only the OncoFinder
11 and TAPPA methods were able to eliminate the cross-platform bias for all three sample types
12 (Fig.7).

13 Hierarchical clustering (dendrograms shown in Supplementary data set S7). demonstrates
14 that only the OncoFinder and TAPPA methods enabled clustering of the *PAS* vectors exclusively
15 according to biological sample type. Thus, among the five *PAS* scoring algorithms tested, only
16 OncoFinder showed effective data aggregation with efficient retention of biological information in
17 three independent tests (Table 3).

19 DISCUSSION

20
21 High throughput gene expression may produce both random and systematic errors, arising from the
22 steps in RNA or protein purification, library preparation and/or amplification, hybridization and
23 sequencing, reading arrays, and mapping and annotation of the reads [Chalaya, 2004; Shugay, 2014;
24 Risso, 2011]. It is generally hard to identify the types of errors and to find out which kind of
25 experimental protocol provides more reliable data. While the measured concentration of each
26 individual gene product may be in error, we show in this report that combining sufficient numbers

1 of these concentrations into a pathway-oriented network apparently generates significantly more
2 stable data. We also tested whether OncoFinder and other *PAS* scoring methods can improve
3 expression data to suppress batch effects, the unwanted variation in gene expression measurements
4 on the same experimental platform made at different times, which frequently originate from the
5 limitation in the number of samples that can be processed at once in a single experiment
6 [Demetrashvili, 2010]. Batch effects also hinder the combination of different experimental datasets.
7 Batch effects are almost inevitable [Lazar, 2012]. By limiting analyses to single data sets, one
8 frequently must use an insufficient number of samples, which leads to high false-negative rates
9 [Lazar, 2012]. Eliminating batch effects enables larger datasets, and provides more statistical
10 power to subsequent analyses [Lazar, 2012].

11 Here, using the Illumina HT12 bead array platform to profile human cancer samples, we
12 demonstrate that the *PAS* scoring technology OncoFinder effectively suppresses batch effects
13 present in the individual gene expression measurements (Fig.1). OncoFinder efficiently increases
14 expression data stability from all major experimental platforms, for both fresh and formalin-fixed,
15 paraffin-embedded tissue samples (Fig.4).

16 Various publicly available repositories of gene expression data embrace the full spectrum of
17 normal and pathological conditions for the majority of known human diseases [Cancer Genome
18 Atlas Research Network, 2008; Jones, 2006]. Unfortunately, batch effects, which bias the
19 expression profiles, hamper the joint analysis of most of this data obtained using different
20 experimental settings.

21 Discrepancies in data obtained on *the same* and *different* experimental platforms, must be
22 addressed by different methods, termed *normalization* and *harmonization*, respectively. For intra-
23 platform normalization, more attention is paid to equilibration of scaling factors, while cross-
24 platform harmonization must address the type of distribution of output intensities for each gene.
25 Existing methods for intra-platform normalization include quantile normalization [Bolstad, 2003]
26 and frozen robust multi-array analysis (FRMA) [McCall, 2010] for microarray data, and the DESeq

1 method [Anders, 2010] for next-generation sequencing.
2 Methods for cross-platform harmonization, such as distance-weighted discrimination (DWD)
3 [Huang, 2012], cross-platform normalization (XPN) [Shabalín, 2008], and platform-independent
4 latent Dirichlet allocation (PLIDA) [Deshwar, 2014], provide deep restructuring or signal intensity
5 redistribution for the entire set of genes profiled. As a rule, the cross-platform harmonization
6 involves data clustering and finding similarity regions among results obtained using different
7 platforms, to strengthen similarity during the harmonization process.

8 Unfortunately, current normalization and harmonization methods hardly distinguish between
9 artifacts introduced by batch effects and the real biological differences. Additional tools are needed
10 to improve normalization and harmonization procedures. We demonstrate here for most major
11 transcriptomic and proteomic commercial platforms that data aggregation at the level of molecular
12 pathways has the potential to reduce greatly the bias in the datasets under comparison. Since each
13 pathway may contain hundreds of different gene products, transition from single gene products to
14 the whole pathway level may restore biologically significant correlations.

15 We propose a term *data aggregation effect* for such restoration of biological correlation at
16 the pathway level. We created a mathematical model that simulates it and identifies the necessary
17 conditions for its applicability. Sample expression profiles must be biased compared to control
18 samples, i.e. the transcriptional signatures of the *case* samples must differ significantly from the
19 normal ones (Fig. 5). The strength of the data aggregation effect grows with the number of gene
20 products in a molecular pathway. The data aggregation effect is especially strong when the initial
21 correlation between the expression data is weak (Fig. 2,3). Finally, the choice of *PAS* scoring
22 method affects the data aggregation effect. On a model data set, the OncoFinder, Pathway-Express
23 and SPIA algorithms result in a considerable data aggregation effect, while TAPPA and TB-Score
24 don't (Fig. 6). Only OncoFinder and TAPPA were able to preserve the biological features on the
25 model dataset MAQC after cross-platform harmonization, while with Pathway-Express, SPIA and
26 TB-Score methods, platform-introduced bias features still dominated the output expression

1 signatures (Fig.7). Thus, among the five *PAS* scoring methods tested here, the OncoFinder
2 algorithm showed the best efficiency and accuracy (Tab.3), which makes OncoFinder a method of
3 choice for many applications using high-throughput analysis of gene expression at the RNA or
4 protein levels.

5 It should be possible in the future to refine *PAS* methods to create universal platform-
6 agnostic analytic tools. These tools have a huge potential to accelerate progress in genetics,
7 physiology, biomedicine, molecular diagnostics and other applications by combining unbiased data
8 from many sources and various experimental platforms.

9 10 MATERIALS AND METHODS

11 12 **Tissue collection and RNA isolation from fresh biosamples**

13 Seven normal bladder and seventeen bladder carcinoma specimens from patients treated at the P.A.
14 Herzen Moscow Oncological Research Institute (HMORI; Moscow, Russia) were analyzed. Of
15 these samples (cancer/normal), nine/three were examined at the Shemyakin-Ovchinnikov Institute
16 of Bioorganic Chemistry (IBC; Moscow, Russia) and eight/four at the University of Lethbridge
17 (UL; Alberta, Canada). All patients provided written informed consent to participate in this study.
18 This study was approved by the local ethical committees at IBC, UL and HMORI. Tumor samples
19 were obtained from patients who had undergone surgery for bladder carcinoma at the HMORI
20 between 2009 and 2013. The median age of the cancer patients at the time of surgical tumor
21 resection was 64 years (range 48–77 years). Tissue samples from non-cancer controls were
22 collected from autopsies at the Department of Pathology at the Faculty of Medicine, Moscow State
23 University. Both the tumors and normal tissues were evaluated by a pathologist to confirm the
24 diagnosis and estimate the tumor cell numbers. All tumor samples used in this study contained at
25 least 80% tumor cells. The median age of the healthy tissue donors was 45 years (range 20–71
26 years). Tissue samples were stabilized in RNAlater (Qiagen, Germany) and then stored at –80°C.

1 Frozen tissue was homogenized in TRIzol Reagent (Life Technologies, USA), and RNA was
2 isolated following the manufacturer's protocol. Purified RNA was dissolved in RNase-free water
3 and stored at -80°C .

4

5 **Microarray profiling of gene expression in fresh biosamples**

6 Total RNA was extracted using TRIzol Reagent and then reverse-transcribed to cDNA and cRNA
7 using the Ambion TotalPrep cRNA Amplification Kit (Invitrogen, USA). The cRNA concentration
8 was quantified and adjusted to 150 ng/ml using an ND-1000 Spectrophotometer (NanoDrop
9 Technologies, USA). 750 ng of each RNA library was hybridized onto the bead arrays.
10 Gene expression experiments were performed by Genoanalytica (Moscow, Russia) and the O.
11 Kovalchuk Laboratory (Lethbridge, Canada) using the Illumina HumanHT-12v4 Expression
12 BeadChip (Illumina, Inc.). This gene expression platform contains more than 25,000 annotated
13 genes and more than 48,000 probes derived from the National Center for Biotechnology
14 Information RefSeq (build 36.2, release 22) and the UniGene (build 199) databases. The expression
15 data were deposited in the GEO database (<http://www.ncbi.nlm.nih.gov/geo/>), accession numbers
16 GSE52519 and GSE65635.

17

18 **Synthesis of microarrays.**

19 A B3 synthesizer (CustomArray, USA) was used for oligonucleotide probe synthesis on the
20 CustomArray ECD 4X2K/12K slides. Synthesis was performed according to the manufacturer's
21 recommendations. At least three replicates of total 3823 unique oligonucleotide probes of 40
22 nucleotides in length for 2278 genes were placed on each chip.

23

24 **Library preparation and hybridization.**

25 RNA was extracted from freshly frozen tissue samples or samples stored in stabilizing buffer
26 solutions using the standard protocol for TRIzol reagent (Life Technologies). RNA extraction from

1 FFPE samples was performed using the RecoverAll™ Total Nucleic Acid Isolation Kit for FFPE.
2 Complete Whole Transcriptome Amplification WTA2 Kit (Sigma) was used for reverse
3 transcription and library amplification. The manufacturer's protocol was modified by adding to
4 amplification reaction a dNTP mix containing biotinylated dUTP, resulting to a final proportion
5 dTTP/biotin-dUTP of 5:1.
6 Hybridization was performed according to the CustomArray ElectraSense™ Hybridization and
7 Detection protocol. The hybridization mix contained 2.5 ug of labeled DNA library, 6X SSPE,
8 0.05% Tween-20, 20mM EDTA, 5x Denhardt solution, 100 ng/ul sonicated calf thymus gDNA, and
9 0,05% SDS. The chip was incubated in the hybridization mix overnight at 50°C. The hybridization
10 efficiency was detected electrochemically using CustomArray ElectraSense™ Detection Kit and
11 ElectraSense™ 4X2K/12K Reader. The chip was designed using the Layout Designer software
12 (CustomArray, USA).

13

14 **Functional annotation of gene expression data**

15 The SABiosciences (<http://www.sabiosciences.com/pathwaycentral.php>) signaling pathways
16 knowledge base was used to determine structures of intracellular pathways, as described previously
17 [Spirin, 2014].

18 *OncoFinder*. We applied the original OncoFinder algorithm [Buzdin, 2014b] for functional
19 annotation of the primary expression data and for calculating *PAS* scores. The microarray gene
20 expression data were quantile normalized according to [Bolstad, 2003]. The formula used to
21 calculate the *PAS* for a given sample and a given pathway *p* is as follows:

$$22 \quad PAS_p = \sum_n ARR_{np} \cdot BTIF_n \cdot \log(CNR_n) \quad (1)$$

23 Here the case-to-normal ratio, CNR_n , is the ratio of the expression level of gene *n* in the sample
24 under investigation to the average expression level of that gene in the control group of samples. The
25 Boolean flag of *BTIF* (beyond tolerance interval flag) equals one or zero when the *CNR* value has

1 simultaneously passed or not passed, respectively, the two criteria that indicate a significantly
2 perturbed expression level from an essentially normal expression level. The first criterion is that the
3 expression level of the sample lies within the tolerance interval, with $p < 0.05$. The second criterion
4 is whether the *CNR* value lies outside the cut-off limits, i.e., either $CNR < 2/3$ or $CNR > 3/2$. ARR_{np} ,
5 the discrete value of the activator/repressor role equals the following fixed values: -1 , when the
6 gene/protein n is a repressor of molecular pathway; 1 , if the gene/protein n is an activator of
7 pathway; 0 , when the gene/protein n is known to be both an activator and a repressor of the
8 pathway; and 0.5 and -0.5 , respectively, tends to be an activator or a repressor of the pathway p ,
9 respectively.

10 Our approach to calculations of *PAS* implies two principal assumptions:

11 1) First, computational modeling of signal transduction processes [Birtwistle, 2007; Borisov, 2009;
12 Kuzmina, 2011] indicates that for most interacting proteins the concentration of their active forms,
13 which are sufficient for downstream signaling, is much lower than the total abundance of the
14 corresponding protein. In other words, signal transduction may be performed even at the very low
15 level for most gene products.

16 2) Second, we stipulate that each pathway graph may be simplified up to the following structure
17 that includes only two chain-like (linear) branches: one for sequential events that promote activation
18 of whole pathway, and another for repressor sequential events. The adequacy of this quite radical
19 approximation was shown before in comparison with the full-scaled kinetic model [Kuzmina,
20 2011], when all protein-protein interactions were described using the mass-action law along each
21 edge of a highly branched pathway graph [Buzdin, 2014].

22 Under these conditions, we presume that all activator/repressor members have equal importance for
23 the whole pathway, and come to the following formula for the overall signal outcome (*SO*) of a

24 given pathway, $SO = \frac{\prod_{i=1}^N [AGEL]_i}{\prod_{j=1}^M [RGEL]_j}$. Here the multiplication is done over all possible activator and

25 repressor proteins in the pathway, $[AGEL]_i$ and $[RGEL]_j$ are relative gene expression levels of

1 activator (i) and repressor (j) members, respectively. To obtain an additive value, it is possible to
 2 take the logarithmic levels of gene expression, and thus come to a function of *PAS*.
 3 The results for 271 pathways were obtained for each sample (see Supplementary Data set S1).
 4 Statistical tests used the R software package.
 5 *TAPPA* (*Topology analysis of pathway phenotype association*). Imagine a pathway graph, $G(V, E)$,
 6 where $V = \{g_1, g_2, \dots, g_n\}$ is the set of graph nodes (vertices), and
 7 $E = \{(g_i, g_j) \mid \text{genes } g_i \text{ and } g_j \text{ interact}\}$ is the set of graph edges [Gao, 2007]. The adjacency
 8 matrix is defined as follows, $a_{ij} = 1$, if $i = j$ or $(g_i, g_j) \in E$, and $a_{ij} = 0$, if $(g_i, g_j) \notin E$. A *centered*
 9 *Z-scoring* procedure was applied to the logarithmic gene expression matrix $x_{is} = (x_{is}^{orig} - \bar{x}_{is}^{orig}) / \sigma_s$.
 10 The adjacency index for a pathway is defined as follows,

$$11 \quad J = \sum_{i=1}^N \sum_{j=1}^N \text{sign}(x_{is} + x_{js}) \sqrt{|x_{is}|} a_{ij} \sqrt{|x_{js}|}, \quad (2)$$

12 where N is the number of genes in the pathway, and the double summation of over the
 13 $\text{sign}(x_{is} + x_{js})$ reveals whether the pathway has more up- or down-regulated genes. The sign (of
 14 what?), indicates whether the whole pathway is up- or down-regulated is calculated as
 15 $\text{TAPPA}_p = J_p - \bar{J}_N$, where \bar{J}_N is the expected value of J over the set of samples that are
 16 considered normal.

17
 18 *TBScore* (*Topology-based score*) [Ibrahim, 2012]. For a pathway p that has N nodes, the value

$$19 \quad \text{TBScore}_p = \sum_{i=1}^N NV_i \cdot NW_i, \text{ where the node value, } NV_i, \text{ equals to zero if all the genes in the node } i \text{ are}$$

20 non-differential genes, or equals to the sum of log-fold-changes of the differential genes in the node
 21 i . The gene is considered differential, if the gene is considered differential in terms of the Boolean
 22 flag *BTIF* (as for the OncoFinder algorithm). The node weight, NW_i , equals the number of
 23 downstream nodes for node i . To determine the value of NW_i , we used the depth-first search

1 method [Even, Sh. Graph Algorithms, Cambridge University Press, 2011] using labeling visited
 2 nodes to avoid the infinite cycling.

3 *Pathway-Express (PE)* [Draghici, 2007]. The PE-score for a pathway K was calculated as follows,

$$4 \quad PE_K = \log(1/p) + \frac{\sum_{g \in K} |PF(g)|}{|\Delta E| N_d(P)}. \quad (3)$$

5 The first term in this sum is the p-value for the probability to obtain the observed or a higher
 6 number N_d of differentially expressed genes (between the pools of case and normal samples) by
 7 random chance, assuming a hypergeometrical distribution for N_d . The second term is a summation
 8 over the perturbation factors (PF) for the all genes g of the pathway K ,

$$9 \quad PF(g) = \Delta E(g) + \sum_{\gamma \in U_g} \beta_{\gamma g} \frac{PF(\gamma)}{n_{down}(\gamma)}. \quad (4)$$

10 Here $\Delta E(g)$ is the signed difference of gene g logarithmic expression in a given sample compared
 11 with the expected value for the pool of normal samples. The latter term expresses the summation
 12 over all the genes γ that belong to the set U_g of the upstream genes for the gene g . The value of
 13 $n_{down}(\gamma)$ denotes the number of downstream genes for gene γ . The weight factor $\beta_{\gamma g}$ indicates the
 14 interaction type between γ and g : $\beta_{\gamma g} = 1$ if γ activates g , and $\beta_{\gamma g} = -1$ when γ inhibits g .
 15 Although the value of PF may be positive or negative, the overall score of PE is obligatory positive.
 16 The search for upstream/downstream genes is performed according to the depth-first search method,
 17 as in the TBScore method.

18 *SPIA (Signal pathway impact analysis)* [Tarca, 2009]. To obtain an estimator for pathway
 19 perturbation that is positive for an up-regulated and negative for a down-regulated pathway, use the
 20 second term in formula (4), resulting in the accuracy value, $Acc(g) = PF(g) - \Delta E(g)$. It can be
 21 shown that [Tarca, 2009] this accuracy vector may be expressed as follows,

$$22 \quad \mathbf{Acc} = \mathbf{B} \cdot (\mathbf{I} - \mathbf{B})^{-1} \cdot \Delta \mathbf{E}, \text{ where} \quad (5)$$

$$\mathbf{B} = \begin{pmatrix} \frac{\beta_{11}}{n_{down}(g_1)} & \frac{\beta_{12}}{n_{down}(g_2)} & \dots & \frac{\beta_{1n}}{n_{down}(g_n)} \\ \frac{\beta_{21}}{n_{down}(g_1)} & \frac{\beta_{22}}{n_{down}(g_2)} & \dots & \frac{\beta_{2n}}{n_{down}(g_n)} \\ \dots & \dots & \dots & \dots \\ \frac{\beta_{n1}}{n_{down}(g_1)} & \frac{\beta_{n2}}{n_{down}(g_2)} & \dots & \frac{\beta_{nn}}{n_{down}(g_n)} \end{pmatrix},$$

2 \mathbf{I} is the identity matrix, and

$$\Delta \mathbf{E} = \begin{pmatrix} \Delta E(g_1) \\ \Delta E(g_2) \\ \dots \\ \Delta E(g_n) \end{pmatrix}.$$

4 The overall score for pathway perturbation calculated as: $SPIA = \sum_g Acc(g)$.

5

6 **Statistical tests**

7 Principal component analyses were performed using the MADE4 package [Culhane, 2005].

8 Hierarchical clustering heat maps with Pearson distances and average linkage were generated using
9 heatmap.2 function from the *gplots* package [Scales, 2015].

10

11 **Mathematical modeling**

12 We performed a Monte Carlo trial to investigate the data aggregation effect. Something may be
13 missing here? We assumed that the number of genes in each pathway is distributed log-normally
14 with the variable median number N . The case-to-normal-ratio (*CNR*) values for each gene were also
15 sampled from the log-normal law, so that the value of $\log CNR$ had a normal distribution. When
16 sampling *CNR*, we distinguished between *biased* and *unbiased* models of gene expression. For the
17 *biased* model, the *CNR* distribution has a median value of 1.5, whereas for the *unbiased* model, the
18 median *CNR* value is 1. The standard deviation of the mean $\log CNR$ value was set to 0.3 for both
19 biased and unbiased models. The independent error produced by an experimental platform was also
20 sampled stochastically. We simulated both the *exact* and *noisy* expression profiling methods. By the

1 definition, *exact* methods did not introduce errors. For *noisy* methods, the error was chosen from
2 the log-normal distribution, with a median value of 1.0. All the calculations were made using the R
3 open source platform (version 3.1.2).

4

5 **Analysis of published transcriptomic and proteomic datasets**

6 Prior to analysis, all the microarray data were quantile normalized [Bolstad, 2003], and the RNA-
7 seq data were normalized using the DESeq package from Bioconductor software [Anders, 2010].
8 All gene products showing zero intensities were skipped to avoid aberrant data variations. Pearson
9 correlation coefficients between the same samples examined using different expression profiling
10 methods (e.g., proteome vs transcriptome or MH vs NGS) were calculated at two levels of data
11 aggregation: first, at the level of distinct genes and gene products – namely for the value of log *CNR*
12 (the so-called C_g correlation value); and, second, at the level of the whole pathways, for the *PAS*
13 value (the C_p correlation coefficient). Then, the ratio $R = \frac{C_p}{C_g}$ was calculated for each sample.

14 **Analysis of biological relevance after cross-platform harmonization.**

15 Transcriptional profiles were obtained using the Agilent Whole Human Genome Oligo and
16 Affymetrix Human Genome U133 Plus 2.0 array platforms. The transcriptomic data were cross-
17 platform harmonized with the XPN method [Shabalin, 2008] using the R package CONOR [Rudy,
18 2011]. Then, the cross-harmonized gene expression patterns between the Agilent and Affymetrix
19 platforms were used as the input data for the *PAS* calculations. For all the calculations, type A
20 samples were used as *normal*, and type B, C and D samples - as *cases*.

21 Euclidean distances between the *PAS* vectors were used to determine whether the resulting
22 *PAS* samples are grouped in agreement with their biological properties (i.e., biological sample types
23 B, C and D compared to A), or according to the experimental platform used to investigate them
24 (i.e., Agilent or Affymetrix microarray platform). The cluster dendrograms and violin plots were
25 drawn using the R packages *dendextend* and *vioplot*, respectively.

26

1 DISCLOSURE DECLARATION

2 The authors declare no conflict of interests.

3

4 ACKNOWLEDGMENTS

5 This work was supported by the Pathway Pharmaceuticals Research Initiative (Hong Kong). The
6 authors thank the First Oncology Research and Advisory Center (Moscow, Russia) for the support
7 in preparation of this manuscript. We would like to thank Alex Kim and ASUS for equipment and
8 support of this research.

9

10

- 1 REFERENCES
- 2 Alexandrova E, Nassa G, Corleone G, Buzdin A, Aliper AM, Terekhanova N, Shepelin D,
3 Zhavoronkov A, Tamm M, Milanesi L. et al. 2016. Large-scale profiling of signaling pathways
4 reveals an asthma specific signature in bronchial smooth muscle cells. *Oncotarget* **7**: 25150-25161.
5 doi: 10.18632/oncotarget.7209.
- 6 Anders S, Huber W. 2010. Differential expression analysis for sequence count data. *Genome Biol*
7 **11**: R106. doi: 10.1186/gb-2010-11-10-r106.
- 8 Artcibasova AV, Korzinkin MB, Sorokin MI, Shegay PV, Zhavoronkov AA, Gaifullin N,
9 Alekseev BY, Vorobyev NV, Kuzmin DV, Kaprin AD et al. 2016. MiRImpact, a new bioinformatic
10 method using complete microRNA expression profiles to assess their overall influence on the
11 activity of intracellular molecular pathways. *Cell Cycle* **15**:689-698. doi:
12 10.1080/15384101.2016.1147633.
- 13 Artemov A, Aliper A, Korzinkin M, Lezhnina K. Jellen L, Zhukov N, Roumiantsev S, Gaifullin N,
14 Zhavoronkov A, Borisov N, Buzdin A. 2015. A method for predicting target drug efficiency in
15 cancer based on the analysis of signaling pathway activation. *Oncotarget* **6**:29347-29356. doi:
16 10.18632/oncotarget.5119.
- 17 Birtwistle MR, Hatakeyama M, Yumoto N, Ogunnaike BA, Hoek JB, Kholodenko BN. 2007.
18 Ligand-dependent responses of the ErbB signaling network: experimental and modeling analyses.
19 *Mol Syst Biol* **3**: 144. doi: 10.1038/msb4100188.
- 20 Bolstad BM, Irizarry RA, Astrand M, Speed TP. 2003. A comparison of normalization methods for
21 high density oligonucleotide array data based on variance and bias. *Bioinformatics* **19**: 185–193.
22 doi: 10.1093/bioinformatics/19.2.185.
- 23 Borisov N, Aksamitiene E, Kiyatkin A, Legewie S, Berkhout J, Maiwald T,
24 Kaimachnikov NP, Timmer J, Hoek JB, Kholodenko BN.. 2009. Systems-level interactions
25 between insulin-EGF networks amplify mitogenic signaling. *Mol Syst Biol*, **5**: 256, 2009. doi:
26 10.1038/msb.2009.19.

- 1 Borisov NM, Terekhanova NV, Aliper SM, Venkova LS, Smirnov PY, Roumiantsev S,
2 Korzinkin MB, Zhavoronkov AA, Buzdin AA. 2014. Signaling pathways activation profiles make
3 better markers of cancer than expression of individual genes. *Oncotarget* **5**: 10198-10205. doi:
4 10.18632/oncotarget.2358.
- 5 Buzdin AA, Zhavoronkov AA, Korzinkin MB, Roumiantsev SA, Aliper AM, Venkova LS,
6 Smirnov PY, Borisov NM. 2014a. The OncoFinder algorithm for minimizing the errors introduced
7 by the high-throughput methods of transcriptome analysis. *Front Mol Biosci* **1**:8. doi:
8 10.3389/fmolb.2014.00008.
- 9 Buzdin AA, Zhavoronkov AA, Korzinkin MB, Venkova LS, Zenin AA, Smirnov PY, Borisov NM.
10 2014b. Oncofinder, a new method for the analysis of intracellular signaling pathway activation
11 using transcriptomic data. *Front Genet* **5**:55. doi: 10.3389/fgene.2014.00055.
- 12 Cabezas-Wallscheid N, Klimmeck D, Hansson J, Lipka DB, Reyes A, Wang Q,
13 Weichenhan D, Lier A, von Paleske L, Renders S et al. 2014. Identification of regulatory networks
14 in HSCs and their immediate progeny via integrated proteome, transcriptome, and DNA methylome
15 analysis. *Cell Stem Cell* **15**:507-522. doi: 10.1016/j.stem.2014.07.005.
- 16 Cancer Genome Atlas Research Network. 2008. Comprehensive genomic characterization defines
17 human glioblastoma genes and core pathways. *Nature* **455**:1061-1068. doi: 10.1038/nature07385.
- 18 Chalaya T, Gogvadze E, Buzdin A, Kovalskaya E, Sverdlov ED. 2004. Improving specificity of
19 DNA hybridization-based methods. *Nucleic Acids Res*, **32**: e130. doi: 10.1093/nar/gnh125.
- 20 Culhane AC, Thioulouse J, Perrière G, Higgins DG. 2005. MADE4: an R package for multivariate
21 analysis of gene expression data. *Bioinformatics* **21**:2789-2790. doi:10.1093/bioinformatics/bti394
- 22 Demetrashvili N, Kron K, Pethe V, Bapat B, Briollais L. 2010. How to deal with batch effect in
23 sequential microarray experiments? *Mol Inform* **29**:387–393. doi: 10.1002/minf.200900019.
- 24 Deshwar AG, Morris Q. 2014. PLIDA: cross-platform gene expression normalization using
25 perturbed topic models. *Bioinformatics* **30**:956-961. doi: 10.1093/bioinformatics/btt574.
- 26 Diederich M, Cerella C. 2016. Non-canonical programmed cell death mechanisms triggered by

- 1 natural compounds. *Semin Cancer Biol* **S1044-579X**:30021-30029. doi:
2 10.1016/j.semcancer.2016.06.001.
- 3 Draghici S, Khatri P, Tarca AL, Amin K, Done A, Voichita C, Georgescu C, Romero R. 2007. A
4 systems biology approach for pathway level analysis. *Genome Res* **17**:1537-1545. doi:
5 10.1101/gr.6202607.
- 6 Gao S, Wang X. 2007. TAPPA: topological analysis of pathway phenotype association.
7 *Bioinformatics* **23**: 3100–3102. doi: 10.1093/bioinformatics/btm460.
- 8 Hanahan D, Weinberg RA. 2000. The hallmarks of cancer, *Cell* **100**: 57-70. doi:10.1016/S0092-
9 8674(00)81683-9.
- 10 Hara Y, Kawasaki N, Hirano K, Hashimoto Y, Adachi J, Watanabe S, Tomonaga T. 2013.
11 Quantitative proteomic analysis of cultured skin fibroblast cells derived from patients with
12 triglyceride deposit cardiomyovasculopathy. *Orphanet J Rare Dis* **8**:197. doi: 10.1186/1750-1172-
13 8-197.
- 14 Croft D, Mundo AF, Haw R, Milacic M, Weiser J, Wu G, Caudy M, Garapati P, Gillespie M,
15 Kamdar MR et al. 2014. The Reactome pathway knowledgebase. *Nucleic Acids Res* **42**: D472-477.
16 doi: 10.1093/nar/gkt1102.
- 17 Huang H1, Lu X, Liu Y, Haaland P, Marron JS. 2012. R/DWD: distance-weighted discrimination
18 for classification, visualization and batch adjustment. *Bioinformatics* **28**: 1182-1183. doi:
19 10.1093/bioinformatics/bts096.
- 20 Hsu MJ1, Chang YC, Hsueh HM. 2014. Biomarker selection for medical diagnosis using the partial
21 area under the ROC curve. *BMC Res Notes* **7**:25. doi: 10.1186/1756-0500-7-25.
- 22 Ibrahim MA, Jassim S, Cawthorne MA, Langlands K. 2012. A topology-based score for pathway
23 enrichment. *J Comput Biol* **19**:563-573. doi: 10.1089/cmb.2011.0182.
- 24 Jones P, Côté RG, Martens L, Quinn AF, Taylor CF, Derache W, Hermjakob H, Apweiler R. 2006.
25 PRIDE: a public repository of protein and peptide identifications for the proteomics community.
26 *Nucleic Acids Res* **34(Database issue)**: D659-663. doi: 10.1093/nar/gkj138.

- 1 Karlsson J, Holmquist Mengelbier L, Ciornei CD, Naranjo A, O'Sullivan MJ, Gisselsson D et al.
2 2014. Clear cell sarcoma of the kidney demonstrates an embryonic signature indicative of a
3 primitive nephrogenic origin. *Genes. Chromosomes Cancer* **53**: 381–391. doi: 10.1002/gcc.22149.
- 4 Khatri P, Sirota M, Butte AJ. 2012. Ten years of pathway analysis: current approaches and
5 outstanding challenges. *PLoS Comput Biol* **8**: e1002375. doi: 10.1371/journal.pcbi.1002375.
- 6 Kholodenko BN, Demin OV, Moehren G, Hoek JB. 1999. Quantification of short term signaling by
7 the epidermal growth factor receptor. *J Biol Chem* **274**: 30169–30181. doi:
8 10.1074/jbc.274.42.30169.
- 9 Kim SC, Jung Y, Park J, Cho Seo C, Kim J, Kim P, Park J, Seo J, Kim J et al. 2013. A high-
10 dimensional, deep-sequencing study of lung adenocarcinoma in female never-smokers. *PLoS One*,
11 **8**: e55596. doi: 10.1371/journal.pone.0055596.
- 12 Kumar D, Bansal G, Narang A, Basak T, Abas T, Dash.D. 2016. Integrating transcriptome and
13 proteome profiling: Strategies and applications. *Proteomics* Jun 25 [Epub ahd of print]. doi:
14 10.1002/pmic.201600140.
- 15 Kuzmina NB and Borisov NM. 2011. Handling complex rule-based models of mitogenic cell
16 signaling (On the example of ERK activation upon EGF stimulation). *Intl Proc Chem Biol Envir*
17 *Engng* **5**: 76-82. doi: 10.7763/ipcbee.2011.v5.17.
- 18 Lazar C, Meganck S, Taminau J, Steenhoff D, Coletta A, Molter C, Weiss-Solís DY, Duque
19 R, Bersini H, Nowé A. 2013. Batch effect removal methods for microarray gene expression data
20 integration: a survey. *Brief Bioinform* **14**: 469-490. doi: 10.1093/bib/bbs037.
- 21 Lebedev TD, Spirin PV, Suntsova MV, Ivanova AV, Buzdin AA, Prokofjeva MM, Rubtsov PM,
22 Prassolov VS. 2015. Receptor tyrosine kinase KIT may regulate expression of genes involved in
23 spontaneous regression of neuroblastoma. *Mol Biol (Mosk)* **49**: 1052-1055. doi:
24 10.7868/S0026898415060154.
- 25 Lezhnina K, Kovalchuk O, Zhavoronkov AA, Korzinkin MB, Zabolotneva AA, Shegay PV,
26 Sokov DG, Gaifullin NM, Rusakov IG, Aliper AM et al. 2014. Novel robust biomarkers for human

- 1 bladder cancer based on activation of intracellular signaling pathways. *Oncotarget* **5**: 9022-9032.
2 doi: 10.18632/oncotarget.2493.
- 3 Makarev E, Izumchenko E, Aihara F, Wysocki PT, Zhu Q, Buzdin A, Sidransky D,
4 Zhavoronkov A, Atala A. 2016. Common pathway signature in lung and liver fibrosis. *Cell Cycle*
5 **15**: 1667-1673. doi: 10.1080/15384101.2016.1152435.
- 6 MAQC Consortium, 2006. The MicroArray Quality Control (MAQC) project shows inter- and
7 intraplatform reproducibility of gene expression measurements. *Nat Biotechnol* **24**: 1151-1161.
- 8 McCall MN, Bolstad BM, Irizarry RA. 2010. Frozen robust multiarray analysis (fRMA).
9 *Biostatistics* **11**: 242-53. doi: 10.1093/biostatistics/kxp059.
- 10 Nakaya A, Katayama T, Itoh M, Hiranuka K, Kawashima S, Moriya Y, Okuda S, Tanaka
11 M, Tokimatsu T, Yamanishi et al. 2013. KEGG OC: a large-scale automatic construction of
12 taxonomy-based ortholog clusters. *Nucleic Acids Res* **41**: D353-D357. doi: 10.1093/nar/gks1239.
- 13 Risso D, Schwartz K, Sherlock G, Dudoit S. 2011. GC-content normalization for RNA-Seq data.
14 *BMC Bioinformatics* **12**: 480. doi: 10.1186/1471-2105-12-480.
- 15 Rudy J, Valafar F. 2011. Empirical comparison of cross-platform normalization methods for gene
16 expression data. *BMC Bioinformatics* **12**: 467. doi: 10.1186/1471-2105-12-467.
- 17 Scales M, Jäger R, Migliorini G, Houlston RS, Henrion MY. 2014. visPIG--a web tool for
18 producing multi-region, multi-track, multi-scale plots of genetic data. *PLoS One* **9**: e107497. doi:
19 10.1371/journal.pone.0107497.
- 20 Shugay M, Britanova OV, Merzlyak EM, Turchaninova MA, Mamedov IZ, Tuganbaev
21 TR, Bolotin DA, Staroverov DB, Putintseva EV, Plevova K. et al. 2014. Towards error-free
22 profiling of immune repertoires. *Nat Methods* **11**: 653-655. doi: 10.1038/nmeth.2960.
- 23 Spirin PV, Lebedev TD, Orlova NN, Gornostaeva AS, Prokofjeva MM, Nikitenko NA, Dmitriev
24 SE, Buzdin AA, Borisov NM, Aliper AM et al. 2014. Silencing AML1-ETO gene expression leads
25 to simultaneous activation of both pro-apoptotic and proliferation signaling. *Leukemia* **28**: 2222-
26 2228. doi: 10.1038/leu.2014.130.

- 1 Tarca AL, Draghici S, Khatri P, Hassan SS, Mittal P, Kim JS, Kim CJ, Kusanovic JP, Romero R.
2 2009. A novel signaling pathway impact analysis, *Bioinformatics* **25**: 75-82.
3 doi:10.1093/bioinformatics/btn577.
- 4 Van Delft J, Gaj S, Lienhard J, Albrecht MW, Kirpiy A, Brauers K, Claessen S, Lizarraga
5 D, Lehrach H, Herwig R, Kleinjans J. 2012. RNA-seq provides new insights in the transcriptome
6 responses induced by the carcinogen benzo[a]pyrene. *Toxicological sciences* **130**: 427–439. doi:
7 10.1093/toxsci/kfs250.
- 8 Venkova L, Aliper A, Suntsova M, Kholodenko R, Shepelin D, [Borisov N](#), Malakhova G, Vasilov
9 R, Roumiantsev S, Zhavoronkov A, Buzdin A. 2015. *Oncotarget* **6**: 27227-32728. doi:
10 10.18632/oncotarget.4507.
- 11 Warnat P, Eils R, Brors B. 2005. Cross-platform analysis of cancer microarray data improves gene
12 expression based classification of phenotypes. *BMC Bioinformatics* **6**: 265. doi: 10.1186/1471-
13 2105-6-265.
- 14 Xu X, Zhang Y, Williams J, Antoniou E, McCombie WR, Wu S, Zhu W, Davidson NO, Denoya
15 P, Li E. 2013. Parallel comparison of Illumina RNA-Seq and Affymetrix microarray platforms on
16 transcriptomic profiles generated from 5-aza-deoxy-cytidine treated HT-29 colon cancer cells and
17 simulated datasets. *BMC Bioinformatics*, **14**: S1. doi: 10.1186/1471-2105-14-S9-S1.
- 18 Yang W, Ramachandran A, You S, Jeong H, Morley S, Mulone MD, Logvinenko T, Kim J, Hwang
19 D, Freeman MR, Adam RM. 2014. Integration of proteomic and transcriptomic profiles identifies a
20 novel PDGF-MYC network in human smooth muscle cells. *Cell Commun Signal* **12**: 44. doi:
21 10.1186/s12964-014-0044-z.
- 22 Zhang L, Zhang J, Yang G, Wu D, Jiang L, Wen Z, Li M. 2013. Investigating the concordance of
23 Gene Ontology terms reveals the intra- and inter-platform reproducibility of enrichment analysis.
24 *BMC Bioinformatics* **14**: 143. doi: 10.1186/1471-2105-14-143.
- 25 Zhavoronkov A, Cantor CR. 2011. Methods for structuring scientific knowledge from many areas
26 related to aging research. *PLoS One* **6**: e22597. doi: 10.1371/journal.pone.0022597.

- 1 Zhavoronkov A, Buzdin AA, Garazha AV, Borisov NM, Moskalev AA. Signaling pathway cloud
2 regulation for in silico screening and ranking of the potential geroprotective drugs. *Front Genet* 5:
3 49. doi: 10.3389/fgene.2014.00049.
- 4 Zhu Q, Izumchenko E, Aliper AM, Makarev E, Paz K, Buzdin AA, Zhavoronkov AA, Sidransky D.
5 2015. Pathway activation strength is a novel independent prognostic biomarker for cetuximab
6 sensitivity in colorectal cancer patients. *Hum Genome Var* 2:15009. doi: 10.1038/hgv.2015.9.
- 7
- 8

1 TABLES

2

3

4 Table 1. Cross-platform comparisons for modeling the data aggregation effect.

	Scenario A	Scenario B	Scenario C	Scenario D
Expression profile	<i>Biased</i>	<i>Biased</i>	<i>Unbiased</i>	<i>Unbiased</i>
Method X	<i>Noisy</i>	<i>Noisy</i>	<i>Noisy</i>	<i>Noisy</i>
Method Y	<i>Noisy</i>	<i>Exact</i>	<i>Noisy</i>	<i>Exact</i>

5

6

1

2 Table 2. Transcriptomic and proteomic datasets used to assess data aggregation effects.

3

Paper reference, Dataset ID	Origin	Case and control samples	Experimental platforms	Number of samples
[van Delft, 2012], GSE36244	HepG2 cells	Cells treated with benzopyrene (<i>cases</i>) vs untreated cells (<i>norms</i>)	Transcriptomes using Affymetrix Human Genome U133 Plus 2.0 arrays and Illumina Genome Analyzer sequencer	4
[Xu, 2013], GSE41588	HT-29 cells	Cells treated with 5-aza-deoxy-cytidine (<i>cases</i>) vs untreated cells (<i>norms</i>)	Transcriptomes using Affymetrix Human Genome U133 Plus 2.0 arrays and Illumina Genome Analyzer sequencer	6
[Kim, 2013] GSE37765	Lung adeno-carcinoma	Tumor samples (<i>cases</i>) vs normal lungs (<i>norms</i>)	Transcriptomes using Agilent 1M CNV arrays and Illumina Genome Analyzer sequencer	6
This study	Renal carcinoma tissue	Tumor samples (<i>cases</i>) vs normal adult kidneys (<i>norms</i>)	Transcriptomes using Illumina Human HT-12 v4 microarrays and Custom microchip platform (see text)	7
[Yang, 2014], GSE52488, PXD000624	Human smooth muscle cells	Cells treated with PDGF served as <i>cases</i> , untreated - as <i>norms</i> .	Transcriptome using Affymetrix Human Gene 1.0 ST arrays and proteome using triplex SILAC at Orbitrap XL mass spectrometer.	2

[Cabezas-Wallscheid, 2014], EMTAB-2262, PXD000572	Murine hemato-poietic stem cells (HSC)	HSC served as <i>norms</i> , multipotent progenitor population 1 (MPP1) – as <i>cases</i> .	Transcriptome using RNA-seq HiSeq2000 (Illumina) and proteome using duplex SILAC at Orbitrap Velos Pro mass spectrometer	4
[Hara, 2013]	Human pathologic skin fibroblasts	Samples from two patients served as <i>cases</i> . Three and two normal samples were used as <i>norms</i> for proteome and transcriptome investigation, respectively	Transcriptome using Affymetrix Human Genome U133 Plus 2.0 arrays and proteome using triplex SILAC at Orbitrap Velos mass spectrometer	2

1

2

1

2

3 Table 3. Comparison of *PAS* scoring methods using functional and statistical tests.

Method	Data aggregation effect	Distance distribution within each sample type	Quality of <i>PAS</i> clustering
OncoFinder	++	+++	+++
TAPPA	--	+++	++
TBScore	-	--	-
Pathway- express	+++	--	--
SPIA	+++	--	--

4

5

6

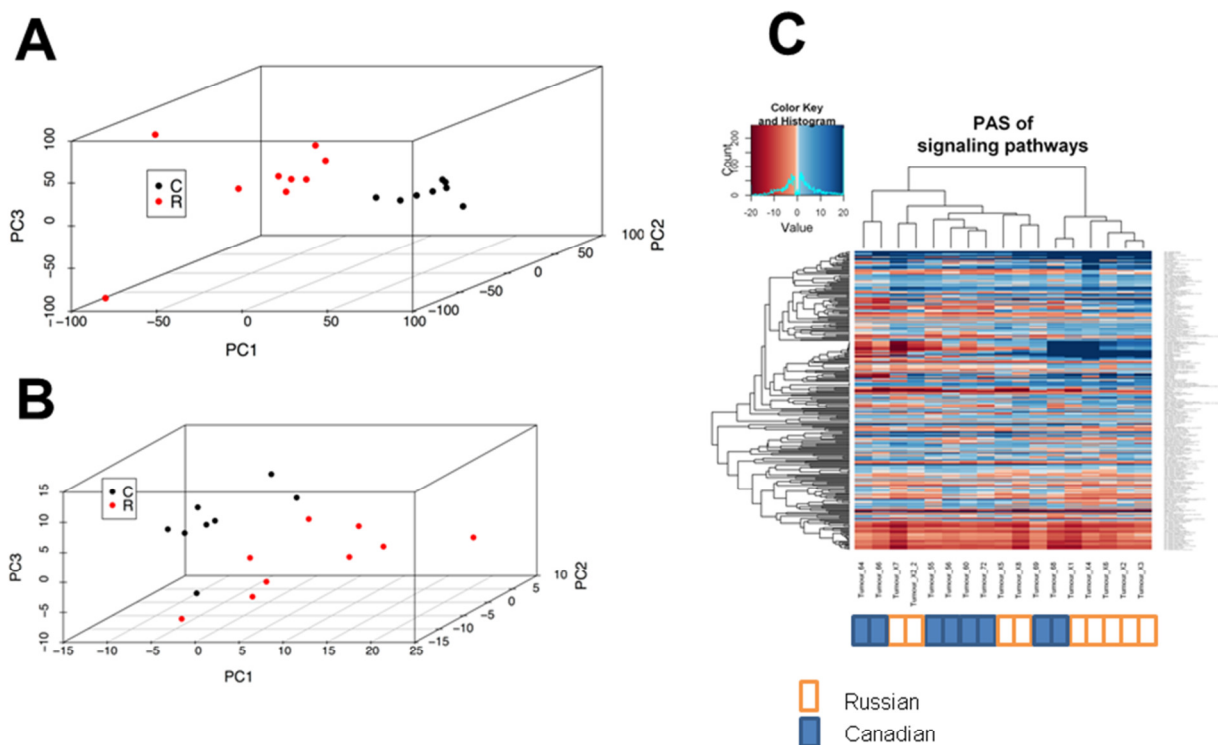
1 FIGURES

2

3

4

5

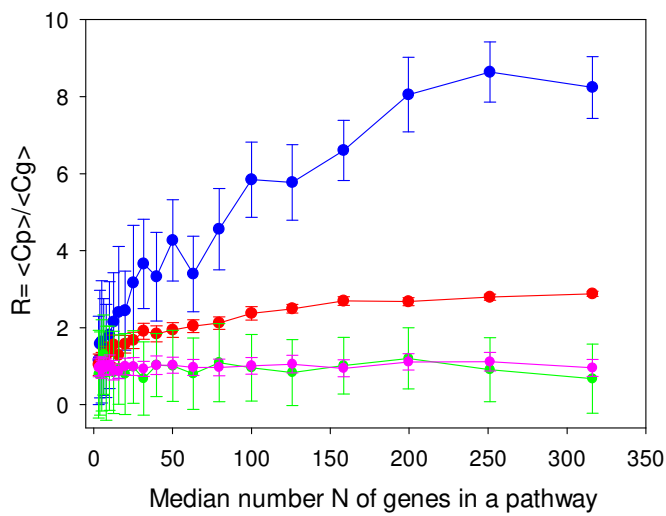


6

7 Figure 1. Renal carcinoma datasets assessed at the level of individual gene expression and pathway
8 activation. *A*, principal component analysis (PCA) plot for transcriptomes from datasets obtained in
9 Russia (red dots) and Canada (black dots), at the level of individual gene expression. *B*, PCA plot at
10 the level of molecular pathway activation. *C*, hierarchical clustering dendrogram of the datasets
11 obtained in Russia (marked white) and Canada (marked blue), at the level of molecular pathway
12 activation.

13

1



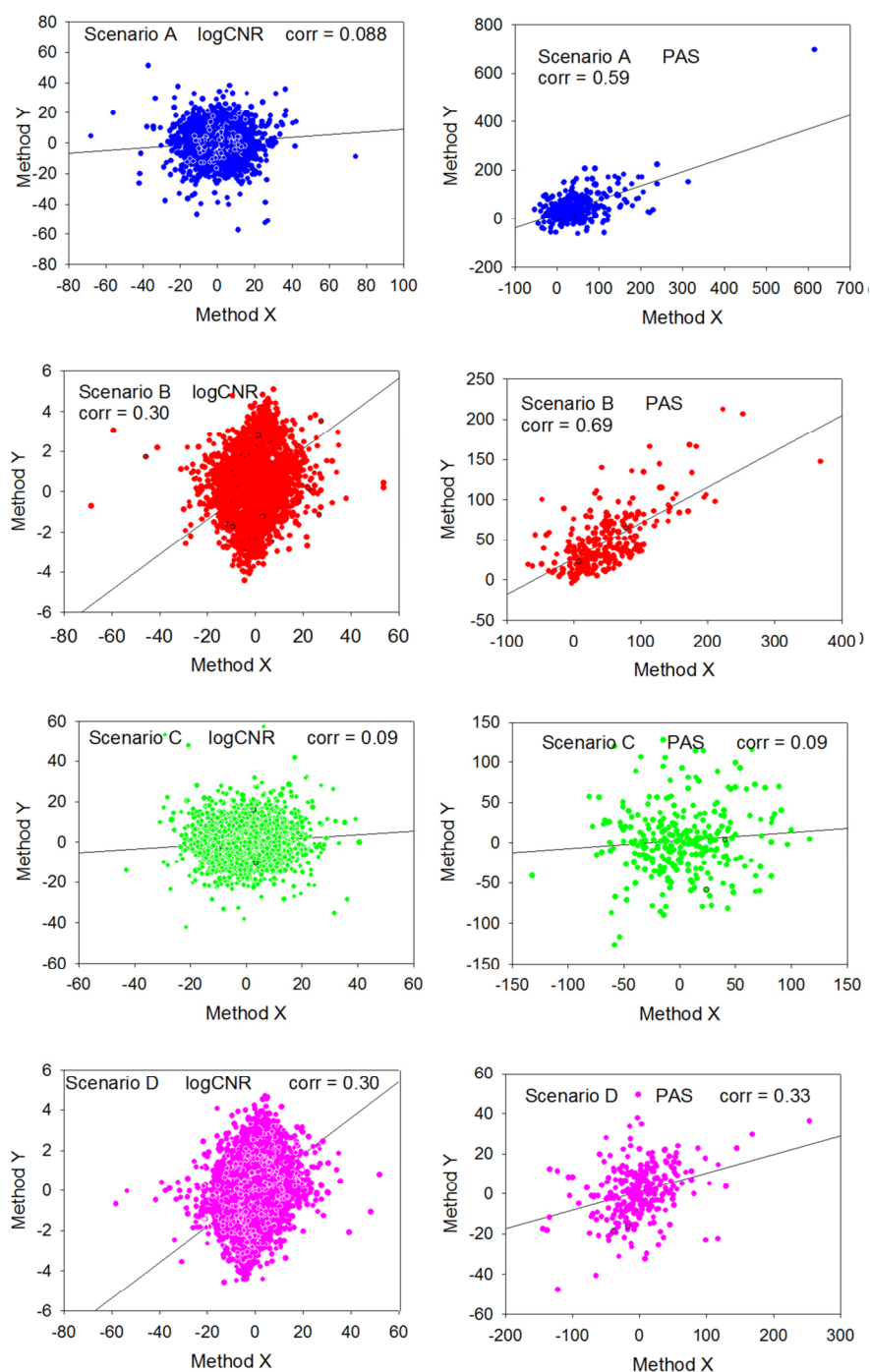
2

3 Figure 2. Ratio of pathway-related and gene-related correlation coefficients between results
4 obtained using hypothetical methods X and Y, as a function of the median gene number, N, in a
5 pathway for four scenarios: A (blue) – *biased* expression profile, *noisy* method Y; B (red) - *biased*
6 expression profile, *exact* method Y; C (green) – *unbiased* expression profile, *noisy* method Y; D
7 (magenta) – *unbiased* expression *exact* method Y. The method X is always considered *noisy*.

8

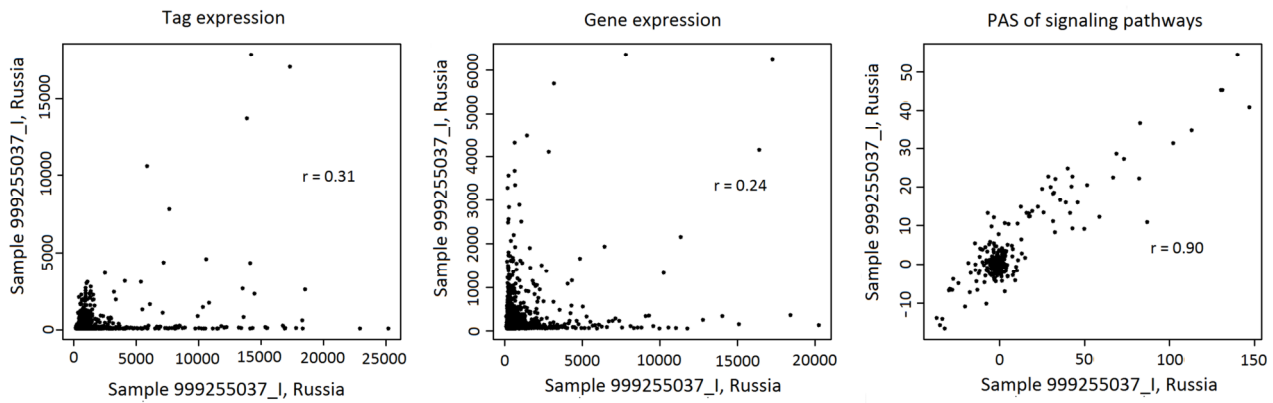
9

10



1

2 Figure 3. Distributions of values obtained during random trials using two different expression
3 profiling methods X (horizontal axis) and Y (vertical axis). Median number of gene products in a
4 pathway is 100. Left column: log *CNR* for individual gene products, method Y vs method X. Right
5 column: *PAS* scoring method Y vs method X. Blue dots: scenario A (*biased* expression profile,
6 *noisy* method Y). Red dots: scenario B (*biased* expression profile, *exact* method Y). Green dots:
7 scenario C (*unbiased* expression profile, *noisy* method Y). Magenta dots: scenario D (*unbiased*
8 expression profile, *exact* method Y). Method X is always considered *noisy*.



1

2

3

4 Figure 4. Correlation between transcriptomic data obtained for the same representative renal

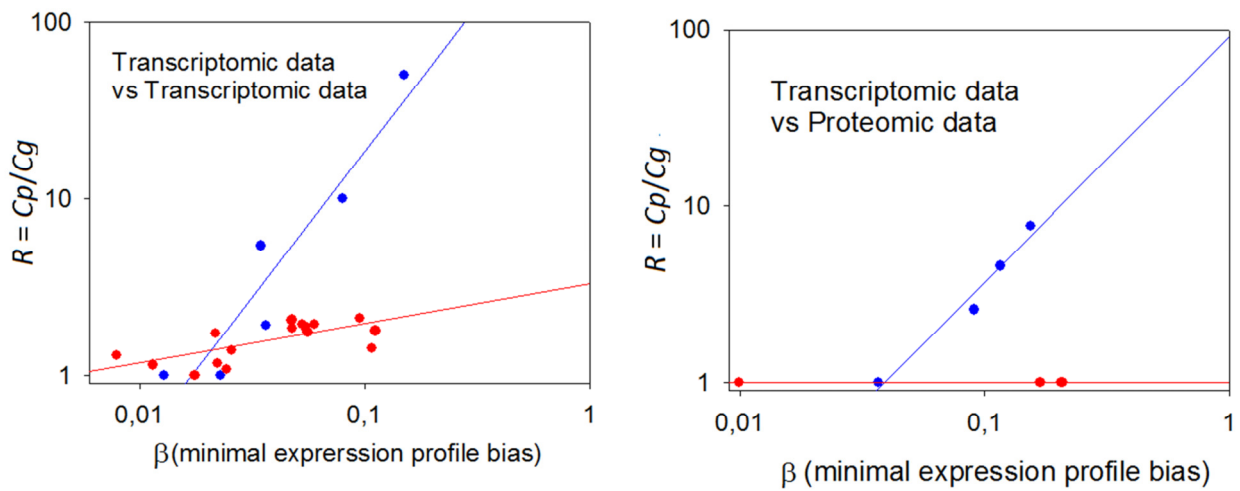
5 carcinoma specimen using the Illumina HT12 (ordinate) and CustomArray (abscissa) microarray

6 platforms. The panels represent (from left to right) correlation between the oligonucleotide

7 expression tags, correlations at the level of individual genes, and correlation at the level of

8 molecular pathways.

9



1

2 Figure 5. Dependence of the data aggregation effect (R) on the minimal expression profile bias β .

3 Left panel: transcriptome-to-transcriptome comparisons for the same samples using different

4 experimental platforms. Right panel: transcriptome-to-proteome comparisons for the same samples.

5 The C_g threshold between the samples *low* and *considerably* correlated at the gp level was chosen as

6 equal to 0.25; blue dots: *low* correlation at gene product level; red dots: *considerable* correlation at

7 gene product level.

8

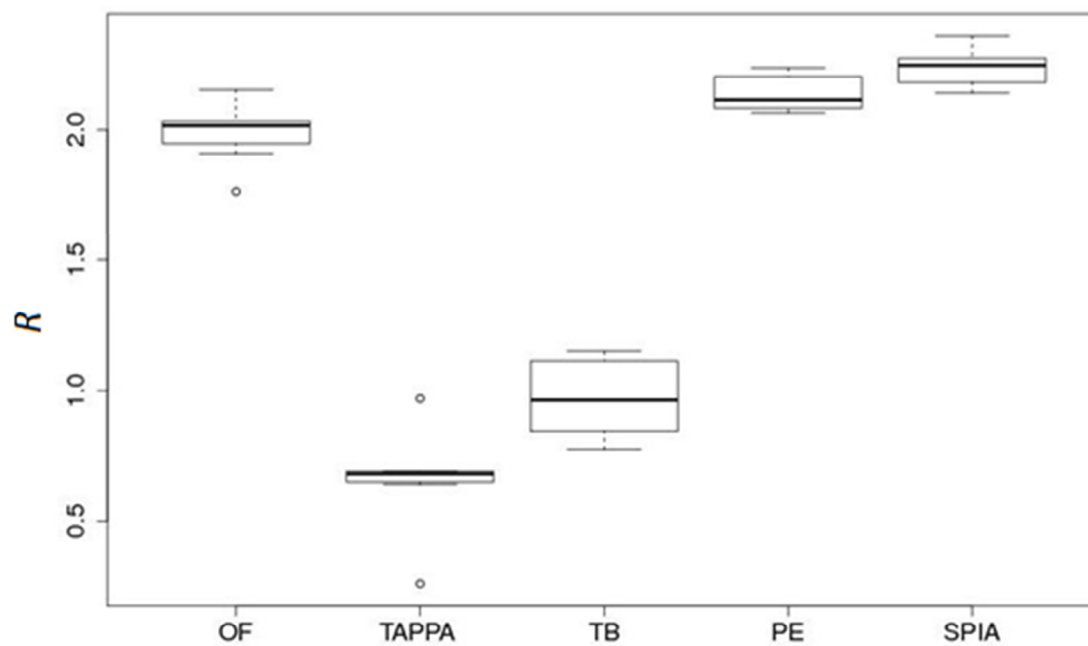
9

10

11

12

1



2

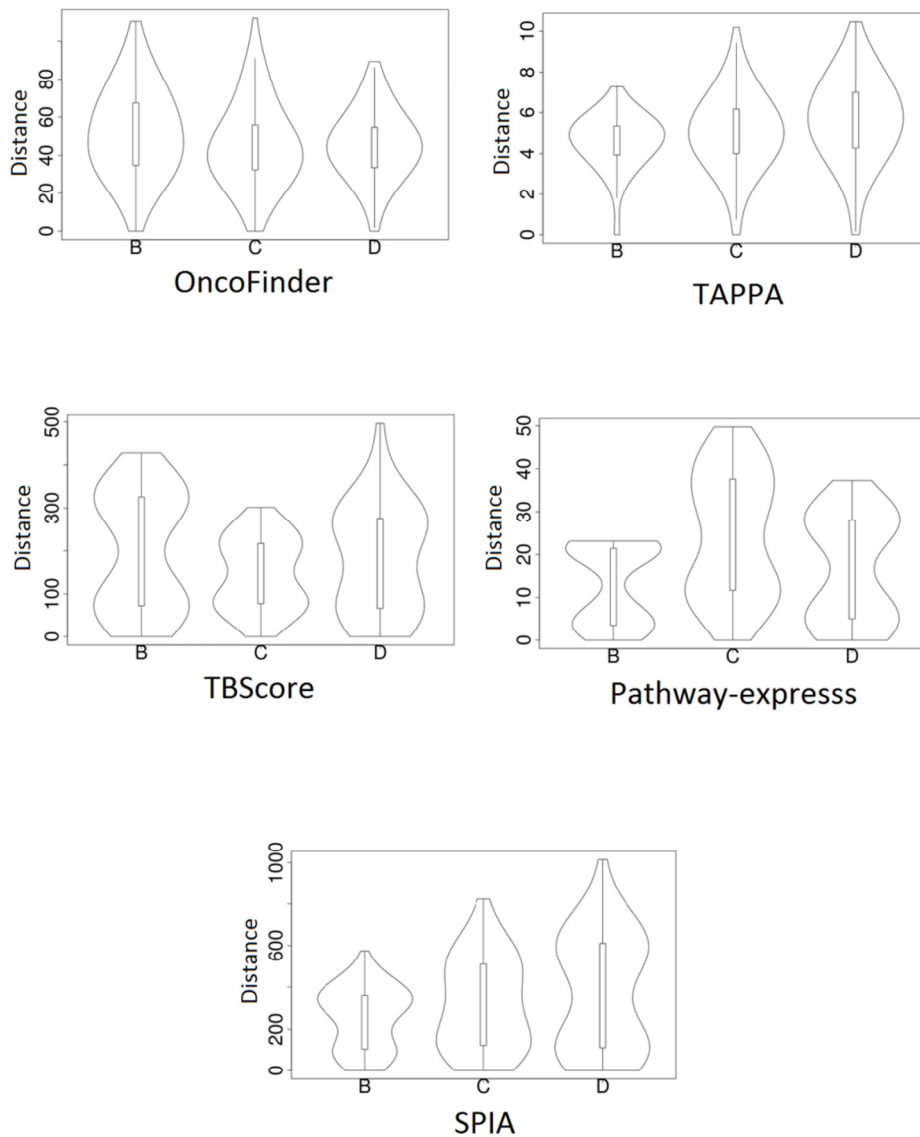
3 Figure 6. Data aggregation effect R for five pathway activation scoring methods (OncoFinder,

4 TAPPA, TBScore (TB), Pathway-Express (PE), and SPIA) on the renal carcinoma dataset.

5

6

7



1

2 Figure 7. Distribution of Euclidean distances between the *PAS* vectors for different sample types
3 taken from the MAQC dataset (marked as B, C and D) using different methods of *PAS* scoring. A
4 *unimodal* distribution indicates lack of significant difference between within-platform and cross-
5 platform distances. A *bimodal* distribution means that the cross-platform *PAS* distance (upper mode
6 in the violin plots) is essentially higher than the within-platform distance. See text for descriptions
7 of the different scoring methods.