BMC
Bioinformatics

**RESEARCH ARTICLE**
**Open Access**

# Meta-analysis methods for combining multiple expression profiles: comparisons, statistical characterization and an application guideline

Lun-Ching Chang[1], Hui-Min Lin[1], Etienne Sibille[2] and George C Tseng[1,3]*

## Abstract

**Background:** As high-throughput genomic technologies become accurate and affordable, an increasing number of data sets have been accumulated in the public domain and genomic information integration and meta-analysis have become routine in biomedical research. In this paper, we focus on microarray meta-analysis, where multiple microarray studies with relevant biological hypotheses are combined in order to improve candidate marker detection. Many methods have been developed and applied in the literature, but their performance and properties have only been minimally investigated. There is currently no clear conclusion or guideline as to the proper choice of a meta-analysis method given an application; the decision essentially requires both statistical and biological considerations.

**Results:** We performed 12 microarray meta-analysis methods for combining multiple simulated expression profiles, and such methods can be categorized for different hypothesis setting purposes: (1) $HS_A$: DE genes with non-zero effect sizes in all studies, (2) $HS_B$: DE genes with non-zero effect sizes in one or more studies and (3) $HS_r$: DE gene with non-zero effect in "majority" of studies. We then performed a comprehensive comparative analysis through six large-scale real applications using four quantitative statistical evaluation criteria: detection capability, biological association, stability and robustness. We elucidated hypothesis settings behind the methods and further apply multi-dimensional scaling (MDS) and an entropy measure to characterize the meta-analysis methods and data structure, respectively.

**Conclusions:** The aggregated results from the simulation study categorized the 12 methods into three hypothesis settings ($HS_A$, $HS_B$, and $HS_r$). Evaluation in real data and results from MDS and entropy analyses provided an insightful and practical guideline to the choice of the most suitable method in a given application. All source files for simulation and real data are available on the author's publication website.

## Background

Microarray technology has been widely used to identify differential expressed (DE) genes in biomedical research in the past decade. Many transcriptomic microarray studies have been generated and made available in public domains such as the Gene Expression Omnibus (GEO) from NCBI (http://www.ncbi.nlm.nih.gov/geo/) and Array Express from EBI (http://www.ebi.ac.uk/arrayexpress/). From the databases, one can easily obtain multiple studies of a relevant biological or disease hypothesis. Since a single study often has small sample size and limited statistical power, combining information across multiple studies is an intuitive way to increase sensitivity. Ramasamy, et al. proposed a seven-step practical guidelines for conducting microarray meta-analysis [1]: "(i) identify suitable microarray studies; (ii) extract the data from studies; (iii) prepare the individual datasets; (iv) annotate the individual datasets; (v) resolve the many-to-many relationship between probes and genes; (vi) combine the study-specific estimates; (vii) analyze, present, and interpret results". In the first step although theoretically meta-analysis increases the statistical power to detect DE genes, the performance can be deteriorated if

* Correspondence: ctseng@pitt.edu
[1]Department of Biostatistics, Graduate school of Public Health, University of Pittsburgh, Pittsburgh, PA, USA
[3]Department of Human Genetics, Graduate School of Public Health, University of Pittsburgh, Pittsburgh, PA, USA
Full list of author information is available at the end of the article

problematic or heterogeneous studies are combined. In many applications, the data inclusion/exclusion criteria are based on ad-hoc expert opinions, a naïve sample size threshold or selection of platforms without an objective quality control procedure. Kang et al. proposed six quantitative quality control measures (MetaQC) for decision of study inclusion [2]. Step (ii)-(v) are related to data preprocessing. Finally, Step (vi) and (vii) involve the selection of meta-analysis method and interpretation of the result and are the foci of this paper.

Many microarray meta-analysis methods have been developed and applied in the literature. According to a recent review paper by Tseng et al. [3], popular methods mainly combine three different types of statistics: combine *p*-values, combine effect sizes and combine ranks. In this paper, we include 12 popular as well as state-of-the-art methods in the evaluation and comparison. Six methods (Fisher, Stouffer, adaptively weighted Fisher, minimum p-value, maximum p-value and rth ordered p-value) belonged to the p-value combination category, two methods (fixed effects model and random effects model) belonged to the effect size combination category and four methods (RankProd, RankSum, product of ranks and sum of ranks) belonged to the rank combination category. Details of these methods and citations will be provided in the Method section. Despite the availability of many methods, pros and cons of these methods and a comprehensive evaluation remain largely missing in the literature. To our knowledge, Hong and Breitling [4], Campain and Yang [5] are the only two comparative studies that have systematically compared multiple meta-analysis methods. The number of methods compared (three and five methods, respectively) and the number of real examples examined (two and three examples respectively with each example covering 2–5 microarray studies) were, however, limited. The conclusions of the two papers were suggestive with limited insights to guide practitioners. In addition, as we will discuss in the Method section, different meta-analysis methods have different underlying hypothesis setting targets. As a result, the selection of an adequate (or optimal) meta-analysis method depends heavily on the data structure and the hypothesis setting to achieve the underlying biological goal.

In this paper, we compare 12 popular microarray meta-analysis methods using simulation and six real applications to benchmark their performance by four statistical criteria (detection capability, biological association, stability and robustness). Using simulation, we will characterize the strength of each method under three different hypothesis settings (i.e. detect DE genes in "all studies", "majority of studies" or "one or more studies"; see Method section for more details). We will compare the similarity and grouping of the meta-analysis methods based on their DE gene detection results (by using a similarity measure and multi-dimension scaling plot) and use an entropy measure to characterize the data structure to determine which hypothesis setting may be more adequate in a given application. Finally, we give a guideline to help practitioners select the best meta-analysis method under the choice of hypothesis setting in their applications.

## Methods

### Real data sets

Six example data sets for microarray meta-analysis were collected for evaluations in this paper. Each example contained 4–8 microarray studies. Five of the six examples were of the commonly seen two-group comparison and the last breast cancer example contained relapse-free survival outcome. We applied the MetaQC package [2] to assess quality of the studies for meta-analysis and determined the final inclusion/exclusion criteria. The principal component analysis (PCA) bi-plots and the six QC measures are summarized in Additional file 1: Figure S1, Tables S2 and S3. Details of the data sets are available in Additional file 1: Table S1.

### Underlying hypothesis settings

Following the classical convention of Brinbaum [6] and Li and Tseng [7] (see also Tseng et al. [3]), meta-analysis methods can be classified into two complementary hypothesis settings. In the first hypothesis setting (denoted as $HS_A$), the goal is to detect DE genes that have non-zero effect sizes in all studies:

$$H_0: \cap_{k=1}^{K}\{\theta_k = 0\} \text{ versus } H_a: \cap_{k=1}^{K}\{\theta_k \neq 0\} \quad (HS_A)$$

where $\theta_k$ is the effect size of study $k$. The second hypothesis setting (denoted as $HS_B$), however, aims to detect a DE gene if it has non-zero effect size in "one or more" studies:

$$H_0: \cap_{k=1}^{K}\{\theta_k = 0\} \text{ versus } H_a: \cup_{k=1}^{K}\{\theta_k \neq 0\} \quad (HS_B)$$

In most applications, $HS_A$ is more appropriate to detect conserved and consistent candidate markers across all studies. However, different degrees of heterogeneity can exist in the studies and $HS_B$ can be useful to detect study-specific markers (e.g. studies from different tissues are combined and tissue specific markers are expected and of interest). Since $HS_A$ is often too conservative when many studies are combined, Song and Tseng (2012) proposed a more practical and robust hypothesis setting (namely $HS_r$) that targets on DE genes with non-zero effect sizes in "majority" of studies, where majority of studies is defined as, for example, more than 50% of combined

studies (i.e. $r \geq 0.5 \cdot K$). The robust hypothesis setting considered was:

$$H_0: \cap_{k=1}^{K} \{\theta_k = 0\} \text{ versus } H_a: \sum_{k=1}^{K} I\{\theta_k \neq 0\} \geq r \quad (HS_r)$$

A major contribution of this paper is to characterize meta-analysis methods suitable for different hypothesis settings ($HS_A$, $HS_B$ and $HS_r$) using simulation and real applications and to compare their performance with four benchmarks to provide a practical guideline.

### Microarray meta-analysis data pre-processing

Assume that we have $K$ microarray studies to combine. For study $k$ ($1 \leq k \leq K$), denote by $x_{gsk}$ the gene expression intensity of gene $g$ ($1 \leq g \leq G$) and sample $s$ ($1 \leq s \leq S_k$; $S_k$ the number of samples in study $k$), and $y_{sk}$ the disease/outcome variable of sample $s$. The disease/outcome variable can be of binary, multi-class, continuous or censored data, representing the disease state, severity or prognosis outcome (e.g. tumor versus normal or recurrence survival time). The goal of microarray meta-analysis is to combine information of $K$ studies to detect differentially expressed (DE) genes associated with the disease/outcome variable. Such DE genes serve as candidate markers for disease classification, diagnosis or prognosis prediction and help understand the genetic mechanisms underlying a disease. In this paper, before meta-analysis we first applied penalized t-statistic to each individual study to generate $p$-values or DE ranks [8] for a binary outcome. In contrast to traditional t-statistic, penalized t-statistic adds a fudge parameter $s_0$ to stabilize the denominator ($T = (\bar{X} - \bar{Y})/(\hat{s} + s_0)$; $\bar{X}$ and $\bar{Y}$ are means of case and control groups) and to avoid a large t-statistic due to small estimated variance $\hat{s}$. The $p$-values were calculated using the null distributions derived from conventional non-parametric permutation analysis by randomly permuting the case and control labels for 10,000 times [9]. For censored outcome variables, Cox proportion hazard model and log-rank test were used [10]. Meta-analysis methods (described in the next subsection) were then used to combine information across studies and generate meta-analyzed $p$-values. To account for multiple comparison, Benjamini and Hochberg procedure was used to control false discovery rate (FDR) [11]. All methods were implemented using the "MetaDE" package in R [12]. Data sets and all programming codes are available at http://www.biostat.pitt.edu/bioinfo/publication.htm.

### Microarray meta-analysis methods

According to a recent review paper [3], microarray meta-analysis methods can be categorized into three types: combine $p$-values, combine effect sizes and combine ranks. Below, we briefly describe 12 methods that were selected for comparison.

### Combine p-values

**Fisher** The Fisher's method [13] sums up the log-transformed $p$-values obtained from individual studies. The combined Fisher's statistic $\chi^2_{\text{Fisher}} = -2\sum_{i=1}^{k} \log(P_i)$ follows a $\chi^2$ distribution with $2k$ degrees of freedom under the null hypothesis (assuming null $p$-values are un;iformly distributed). Note that we perform permutation analysis instead of such parametric evaluation for Fisher and other methods in this paper. Smaller $p$-values contribute larger scores to the Fisher's statistic.

**Stouffer** Stouffer's method [14] sums the inverse normal transformed $p$-values. Stouffer's statistics $T_{\text{Stouffer}} = \sum_{i=1}^{k} z_i/\sqrt{k}(z_i = \Phi^{-1}(p_i)$, where $\Phi$ is standard normal c.c.f) follows a standard normal distribution under the null hypothesis. Similar to Fisher's method, smaller $p$-values contribute more to the Stouffer's score, but in a smaller magnitude.

**Adaptively weighted (AW) Fisher** The AW Fisher's method [7] assigns different weights to each individual study $T_{\text{AW}} = -\sum_{k=1}^{K} w_k \cdot \log(P_i), w_k = 0 \text{ or } 1$ and it searches through all possible weights to find the best adaptive weight with the smallest derived $p$-value. One significant advantage of this method is its ability to indicate which studies contribute to the evidence aggregation and elucidates heterogeneity in the meta-analysis. Details can be referred to the Additional file 1.

**Minimum $p$-value (minP)** The minP method takes the minimum $p$-value among the $K$ studies as the test statistic [15]. It follows a beta distribution with degrees of freedom $\alpha = 1$ and $\beta = k$ under the null hypothesis. This method detects a DE gene whenever a small $p$-value exists in any one of the $K$ studies.

**Maximum $p$-value (maxP)** The maxP method takes maximum $p$-value as the test statistic [16]. It follows a beta distribution with degrees of freedom $\alpha = K$ and $\beta = 1$ under the null hypothesis. This method targets on DE genes that have small $p$-values in "all" studies.

**r-th ordered $p$-value (rOP)** The rOP method takes the $r$-th order statistic among sorted $p$-values of $K$ combined studies. Under the null hypothesis, the statistic follows a beta distribution with degrees of freedom $\alpha = r$ and $\beta = K - r + 1$. The minP and maxP methods are special cases of rOP. In Song and Tseng [17], rOP is considered a robust form of maxP (where $r$ is set as greater than $0.5 \cdot K$) to identify candidate markers differentially expressed in "majority" of studies.

### Combine effect size

**Fixed effects model (FEM)** FEM combines the effect size across $K$ studies by assuming a simple linear model with an underlying true effect size plus a random error in each study.

**Random effects model (REM)** REM [18] extends FEM by allowing random effects for the inter-study heterogeneity in the model. Detailed formulation and inference of FEM and REM are available in the Additional file 1.

### Combine rank statistics

**RankProd (RP) and RankSum (RS)** RankProd and RankSum are based on the common biological belief that if a gene is repeatedly at the top of the lists ordered by up- or down-regulation fold change in replicate experiments, the gene is more likely a DE gene [19]. Detailed formulation and algorithms are available in the Additional file 1.

**Product of ranks (PR) and Sum of ranks (SR)** These two methods apply a naïve product or sum of the DE evidence ranks across studies [20]. Suppose $R_{gk}$ represents the rank of $p$-value of gene $g$ among all genes in study $k$. The test statistics of PR and SR methods are calculated as $PR_g = \prod_{k=1}^{K} R_{gk}$ and $SR_g = \sum_{k=1}^{K} R_{gk}$, respectively. $P$-values of the test statistics can be calculated analytically or obtained from a permutation analysis. Note that the ranks taken from the smallest to largest (the choice in the method) are more sensitive than ranking from largest to smallest in the PR method, while it makes no difference to SR.

### Characterization of meta-analysis methods

#### MDS plots to characterize the methods

The multi-dimensional scaling (MDS) plot is a useful visualization tool for exploring high-dimensional data in a low-dimensional space [21]. In the evaluation of 12 meta-analysis methods, we calculated the adjusted DE similarity measure for every pair of methods to quantify the similarity of their DE analysis results in a given example. A dissimilarity measure is then defined as one minus the adjusted DE similarity measure and the dissimilarity measure is used to generate an MDS plot of the 12 methods. In the MDS plot, methods that are clustered in a neighborhood indicate that they produce similar DE analysis results.

### Entropy measure to characterize data sets

As indicated in the Section of "Underlying hypothesis settings", selection of the most suitable meta-analysis method(s) largely depends on their underlying hypothesis setting ($HS_A$, $HS_B$ and $HS_r$). The selection of a hypothesis

setting for a given application should be based on the experimental design, biological knowledge and the associated analytical objectives. There are, however, occasions that little prior knowledge or preference is available and an objective characterization of the data structure is desired in a given application. For this purpose, we developed a data-driven entropy measure to characterize whether a given meta-analysis data set contains more $HS_A$-type markers or $HS_B$-type markers [22]. The algorithm is described below:

1. Apply Fisher's meta-analysis method to combine $p$-values across studies to identify the top $H$ candidate markers. Here we used $H = 1,000$, $H$ represents the rough number of DE genes (in our belief) that are contained in the data.
2. For each selected marker, the standardized minus $p$-value score for gene $g$ in the $k$-th study is defined as $l_{gk} = -\log\left(p_{gk}\right)/-\sum_{k=1}^{K}\log\left(p_{gk}\right)$. Note that $0 \le l_{gk} \le 1$, large $l_{gk}$ corresponds to more significant $p$-value $p_{gk}$, and $\sum_{k=1}^{K} l_{gk} = 1$.
3. The entropy of gene $g$ is defined as $e_g = -\sum_{k=1}^{K} l_{gk} \log(l_{gk})$. Box-plots of entropies of the top $H$ genes are generated for each meta-analysis application (Figure 1(b)).
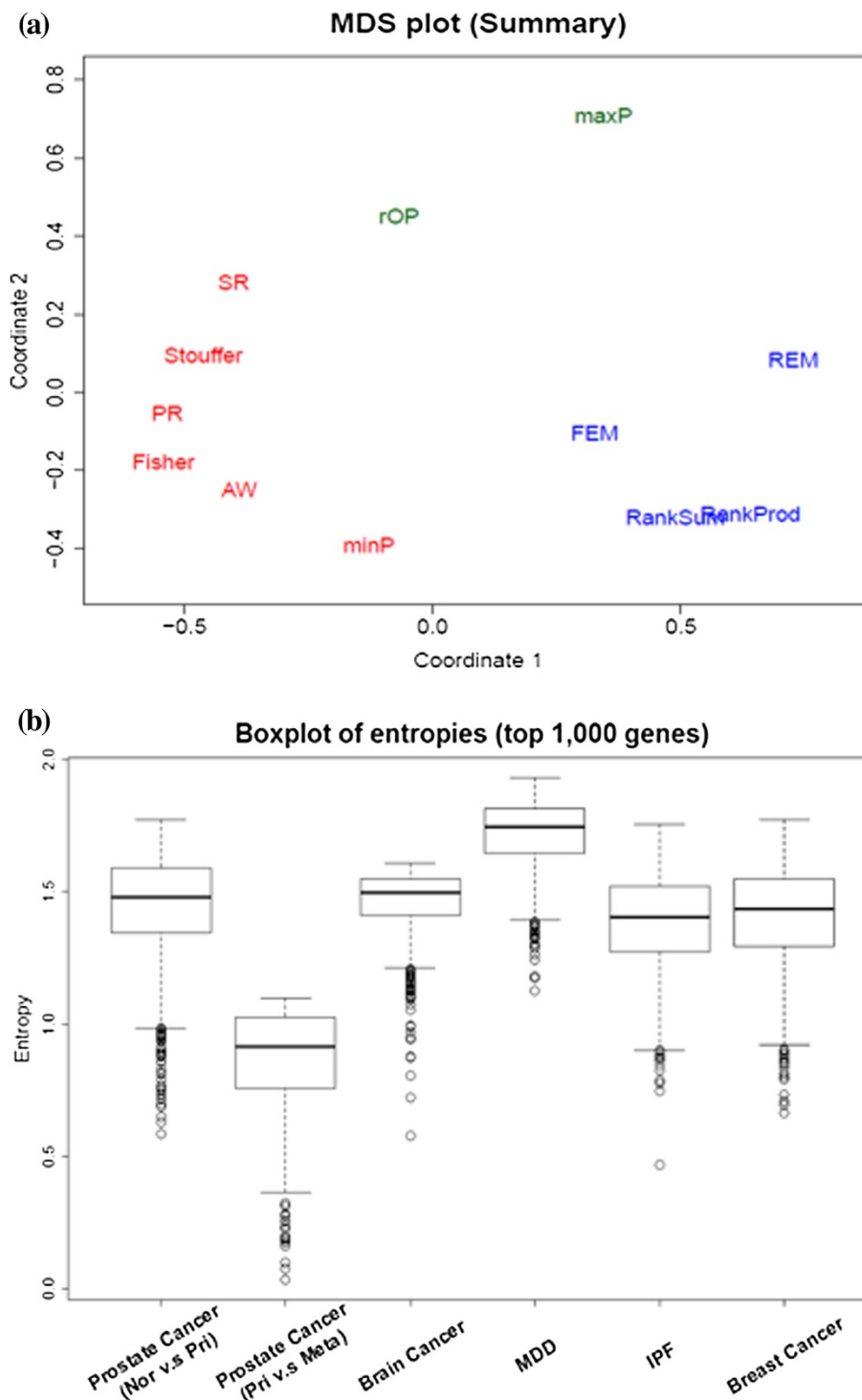
Intuitively, a high entropy value indicates that the gene has small $p$-values in all or most studies and is of $HS_A$ or $HS_r$-type. Conversely, genes with small entropy have small $p$-values in one or only few studies where $HS_B$-type methods are more adequate. When calculating $l_{gk}$ in step 2, we capped $-\log(p_{gk})$ at 10 to avoid contributions of close-to-zero $p$-values that can generate near-infinite scores. The entropy box-plot helps determine an appropriate meta-analysis hypothesis setting if no pre-set biological objective exists.

### Evaluation criteria

For objective quantitative evaluation, we developed the following four statistical criteria to benchmark performance of the methods.

### Detection capability

The first criterion considers the number of DE genes detected by each meta-analysis method under the same pre-set FDR threshold (e.g. FDR = 1%). Although detecting more DE genes does not guarantee better "statistical power", this criterion has served as a surrogate of statistical power in previous comparative studies [23]. Since we do not know the underlying true DE genes, we refer to this evaluation as "detection capability" in this paper. An implicit assumption underlying this

**Figure 1 Characterization of methods and datasets. (a)** Multi-dimensional scaling (MDS) plot of all 12 methods based on the average dissimilarity matrix of six examples and **(b)** The box-plots of entropies in six data sets. Colors (red, green and blue) indicate clusters of methods with similar DE detection ordering. High entropies indicate that high consistency of DE gene detection across studies (e.g. MDD). Low entropies show greater heterogeneity in DE gene detection (e.g. prostate cancer).

criterion is that the statistical procedure to detect DE genes in each study and the FDR control in the meta-analysis are accurate (or roughly accurate). To account for data variability in the evaluation, we bootstrapped (i.e. sampled with replacement to obtain the same number of samples in each bootstrapped dataset) the samples
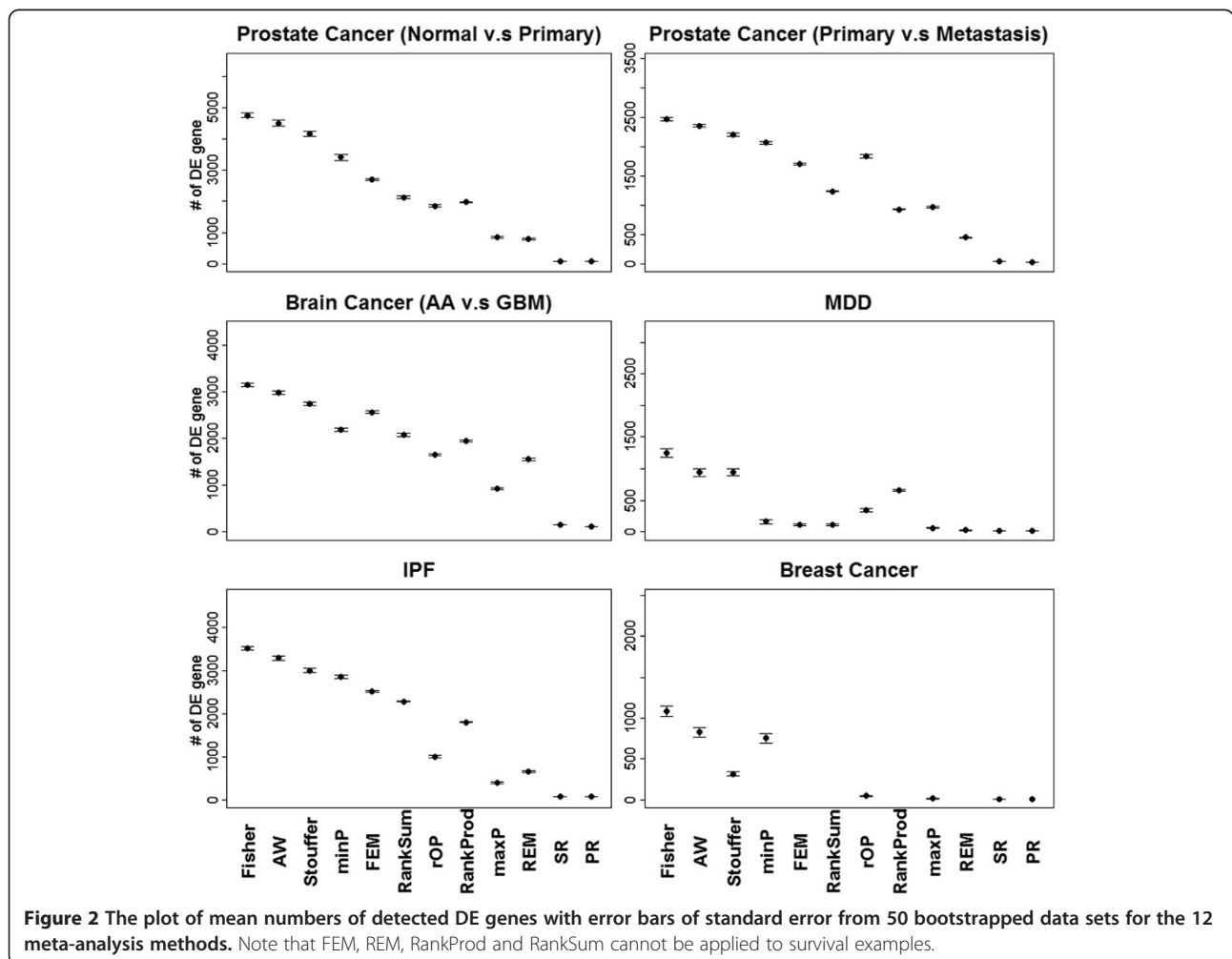
in each study for $B = 50$ times and show the plots of ean with standard error bars. In the bootstrapping, the entire sample is either selected or not so the gene dependence structure is maintained. Denote by $r_{meb}$ the rank of detection capability performance (the smaller the better) of method $m$ ($1 \leq m \leq 12$) in example $e$ ($1 \leq e \leq 6$) and in the $b^{th}$ ($1 \leq b \leq 12$) bootstrap simulation. The mean standardized rank (MSR) for method $m$ and example $e$ is calculated as $MSR_{me} = \sum_{b=1}^{B}$ ($r_{meb}/\#$of methods compared)$/B$ and the aggregated standardized rank (ASR) is calculated as $ASR_m = \sum_{e=1}^{6} MSR_{me}/6$, representing the overall performance of method $m$ across all six examples. Additional file 1: Table S4 shows the MSR and ASR of all 12 methods and Figure 2 (in the Result section) shows plot of mean with standard error bars for each method ordered by ASR. We note that MSR and ASR are both standardized between 0 and 1. The standardization in MSR is necessary because in the breast cancer survival example

we cannot apply FEM, REM, RankSum and RankProd as they are developed only for a two group comparison.

### Biological association

The second criterion requires that a good meta-analysis method should detect a DE gene list that has better association with pre-defined "gold standard" pathways related to the targeted disease. Such a "gold standard" pathway set should be obtained from biological knowledge for a given disease or biological mechanism under investigation. However, since most disease or biological mechanisms are not well-studied, obtaining such "gold standard" pathways is either difficult or questionable. To facilitate this evaluation without bias, we develop a computational and data-driven approach to determine a set of surrogate disease-related pathways out of a large collection of pathways by combining pathway enrichment analysis results from each single study. Specifically, we first collected 2,287 pathways (gene sets) from MSigDB (http://www.broadinstitute.org/gsea/msigdb/): 1,454 pathways from "GO", 186 pathways



**Figure 2 The plot of mean numbers of detected DE genes with error bars of standard error from 50 bootstrapped data sets for the 12 meta-analysis methods.** Note that FEM, REM, RankProd and RankSum cannot be applied to survival examples.

from "KEGG", 217 pathways from "BIOCARTA" and 430 pathways from "REACTOME", respectively. We filtered out pathways with less than 5 genes or more than 200 genes and 2,113 pathways were left for the analysis. DE analysis was performed in each single study separately and pathway enrichment analysis was performed for all the 2,113 pathways by the Kolmogorov-Smirnov (KS) association test. Denote by $p_{uk}$ the resulting pathway enrichment $p$-value from KS test for pathway $u$ ($1 \le u \le 2,113$) and study $k$ ($1 \le k \le K$). For a given study $k$, enrichment ranks over pathways were calculated as $r_{uk} = \text{rank}_u(p_{uk})$. A rank-sum score for a given pathway $u$ was then derived as $S_u = \sum_{k=1}^{K} r_{uk}$. Intuitively, pathways with small rank-sum scores indicate that they are likely associated with the disease outcome by aggregated evidence of the $K$ individual study analyses. We choose the top $|D|$ pathways that had the smallest rank-sum scores as the surrogate disease-related pathways and used these to proceed with the biological association evaluation of meta-analysis methods in the following.

Given the selected surrogate pathways $D$, the following procedure was used to evaluate performance of the 12 meta-analysis methods for a given example $e$ ($1 \le e \le 6$). For each meta-analysis method $m$ ($1 \le m \le M = 12$), the DE analysis result was associated with pathway $u$ and the resulting enrichment $p$-value by KS-test was denoted by $\tilde{P}_{med}$ ($1 \le d \le |D|$). The rank of $\tilde{P}_{med}$ for method $m$ among 12 methods was denoted by $v_{med} = \text{rank}_m (\tilde{P}_{med})$. Similar to the detection capability evaluation, we calculated the mean standardized rank (MSR) for method $m$ and example $e$ as $MSR_{me} = \sum_{d=1}^{D} (v_{med} / \#$ of the methods compared$)/D$ and the aggregated standardized rank (ASR) as $ASR_m = \sum_{e=1}^{6} MSR_{me}/6$, representing the overall performance of method $m$. To select the parameter $|D|$ for surrogate disease-related pathways, Additional file 1: Figure S4 shows the trend of $MSR_{me}$ (on the $y$-axis) versus $|D|$ (on the $x$-axis) as $|D|$ increases. The result indicated that the performance evaluation using different $D$ only minimally impacted the conclusion when $D > 30$. We choose $D = 100$ throughout this paper.

Note that we used KS test, instead of the popular Fisher's exact test because each single study detected variable number of DE genes under a given FDR cutoff and the Fisher's exact test is usually not powerful unless a few hundred DE genes are detected. On the other hand, the KS test does not require an arbitrary $p$-value cutoff to determine the DE gene list for enrichment analysis.

### Stability

The third criterion examines whether a meta-analysis method generates stable DE analysis result. To achieve this goal, we randomly split samples into half in each

study (so that cases and controls are as equally split as possible). The first half of each study was taken to perform the first meta-analysis and generate a DE analysis result. Similarly, the second half of each study was taken to perform a second meta-analysis. The generated DE analysis results from two separate meta-analyses were compared by the adjusted DE similarity measure (to be described in the next section). The procedure is repeated for $B = 50$ times. Denote by $S_{meb}$ the adjusted DE similarity measure of method $m$ of the $b^{th}$ simulation in example $e$. Similar to the first two criteria, MSR and ASR were calculated based on $S_{meb}$ to evaluate the methods.

### Robustness

The final criterion investigates the robustness of a meta-analysis method when an outlying irrelevant study is mistakenly added to the meta-analysis. For each of the six real examples, we randomly picked one irrelevant study from the other five examples, added it to the specific example for meta-analysis and evaluated the change from the original meta-analysis. The adjusted DE similarity measure was calculated between the original meta-analysis and the new meta-analysis with an added outlier. A high adjusted DE similarity measure shows better robustness against inclusion of the outlying study. This procedure was repeated until all irrelevant studies were used. The MSR and ASR are then calculated based on the adjusted DE similarity measures to evaluate the methods.

### Similarity measure between two ordered DE gene lists

To compare results of two DE detection methods (from single study analysis or meta-analysis), a commonly used method in the literature is to take the DE genes under certain $p$-value or FDR threshold, plot the Venn diagram and compute the ratio of overlap. This method, however, greatly depends on the selection of FDR threshold and is unstable. Another approach is to take the generated DE ordered gene lists from two methods and compute the non-parametric Spearman rank correlation [24]. This method avoids the arbitrary FDR cutoff but gives, say, the top 100 important DE genes and the bottom 100 non-DE genes equal contribution. To circumvent this pitfall, Li et al. proposed a parametric reproducibility measure for ChIP-seq data in the ENCODE project [25]. Yang et al. introduced an OrderedList measure to quantify similarity of two ordered DE gene lists [26]. For simplicity, we extended the OrderedList measure into a standardized similarity score for the evaluation purpose in this paper. Specifically, suppose $G_1$ and $G_2$ are two ordered DE gene lists (e.g. ordered by $p$-values) and small ranks represent more significant DE genes. We denote by $O_n(G_1, G_2)$ the number of overlapped genes in the top $n$ genes of $G_1$ and $G_2$. As a result, $0 \le O_n(G_1, G_2) \le n$ and a large $O_n(G_1, G_2)$ value indicates high similarity of the two ordered lists in

the top $n$ genes. A weighted average similarity score is calculated as $S(G_1, G_2) = \sum_{n=1}^{G} e^{-an} \cdot O_n(G_1, G_2)$, where $G$ is the total number of matched genes and the power $\alpha$ controls the magnitude of weights emphasized on the top ranked genes. When $\alpha$ is large, top ranked genes are weighted higher in the similarity measure. The expected value (under the null hypothesis that the two gene rankings are randomly generated) and maximum value of $S$ can be easily calculated: $E_{null}(S(G_1, G_2)) = \sum_{n=1}^{G} e^{-an} \cdot n^2 / G$ and $max(S(G_1, G_2)) = \sum_{n=1}^{G} e^{-an} \cdot n$. We apply an idea similar to adjusted Rand index [27] used to measure similarity of two clustering results and define the adjusted DE similarity measure as

$$S^*(G_1, G_2) = \frac{S(G_1, G_2) - E_{null}(S(G_1, G_2))}{\text{Max}(S(G_1, G_2)) - E_{null}(S(G_1, G_2))}$$

This measure ranges between –1 to 1 and gives an expected value of 0 if two ordered gene lists are obtained by random chance. Yang et al. proposed a resampling-based and ROC methods to estimate the best selection of $\alpha$. Since the number of DE genes in our examples are generally high, we choose a relatively small $\alpha = 0.001$ throughout this paper. We have tested different $\alpha$ and found that the results were similar (Additional file 1: Figure S7).

## Results

### Simulation setting

We conducted simulation studies to evaluate and characterize the 12 meta-analysis methods for detecting biomarkers in the underlying hypothesis settings of $HS_A$, $HS_B$ or $HS_r$. The simulation algorithm is described below:

1. We simulated 800 genes with 40 gene clusters (20 genes in each cluster) and other 1,200 genes do not belong to any cluster. The cluster indexes $C_g$ for gene $g$ ($1 \le g \le 2,000$) were randomly sampled, such that $\sum I\{C_g = 0\} = 1,200$ and $\sum I\{C_g = c\} = 20$, $1 \le c \le 40$.
2. For genes in cluster $c$ ($1 \le c \le 40$) and in study $k$ ($1 \le k \le 5$), we sampled $\sum_{ck}' \sim W^{-1}(\Psi, 60)$, where $\Psi = 0.5I_{20 \times 20} + 0.5J_{20 \times 20}$, $W^{-1}$ denotes the inverse Wishart distribution, $I$ is the identity matrix and $J$ is the matrix with all elements equal 1. We then standardized $\Sigma_{ck}'$ into $\Sigma_{ck}$ where the diagonal elements are all 1's.
3. 20 genes in cluster $c$ was denoted by the index of $g_{c1}, ..., g_{c20}$, i.e. $C_{g_{cj}} = c$, where $1 \le c \le 40$ and $1 \le j \le 20$. We sampled gene expression levels of genes in cluster $c$ for sample $n$ as $\left(X_{g_{c1}nk}', ..., X_{g_{c20}nk}'\right)^T \sim MVN$ $\left(0, \sum_{ck}\right)$ where $1 \le n \le 100$ and $1 \le k \le 5$, and sample

expression level for the gene $g \sim N(0, \sigma_k^2)$ which is not in any cluster for sample $n$, where $1 \le n \le 100$, $1 \le k \le 5$ and $\sigma_k^2$ was uniformly distributed from [0.8, 1.2], which indicates different variance for study $k$.
4. For the first 1,000 genes ($1 \le g \le 1,000$), $k_g$ (the number of studies that are differentially expressed for gene $g$) was generated by sampling $k_g = 1, 2, 3, 4$ and 5, respectively. For the next 1,000 genes ($1,001 \le g \le 2,000$), $k_g = 0$ represents non-DE genes in all five studies.
5. To simulate expression intensities for cases, we randomly sampled $\delta_{gk} \in \{0, 1\}$, such that $\sum_k \delta_{gk} = k_g$. If $\delta_{gk} = 1$, gene $g$ in study $k$ was a DE gene, otherwise it was a non-DE gene. When $\delta_{gk} = 1$, we sampled expression intensities $\mu_{gk}$ from a uniform distribution in the range of [0.5, 3], which means we considered the concordance effect (up-regulated) among all simulated studies. Hence, the expression for control samples are $X_{gnk} = X_{gnk}'$, and case samples are $Y_{gnk} = X_{g(n+50)k}' + \mu_{gk} \cdot \delta_{gk}$, for $1 \le g \le 2,000$, $1 \le n \le 50$ and $1 \le k \le 5$.

**Table 1 The detected number of DE genes (at FDR = 5%), the true FDR, AUC values under $HS_A$ and $HS_B$ and the concluding characterization of targeted hypothesis setting of each method**

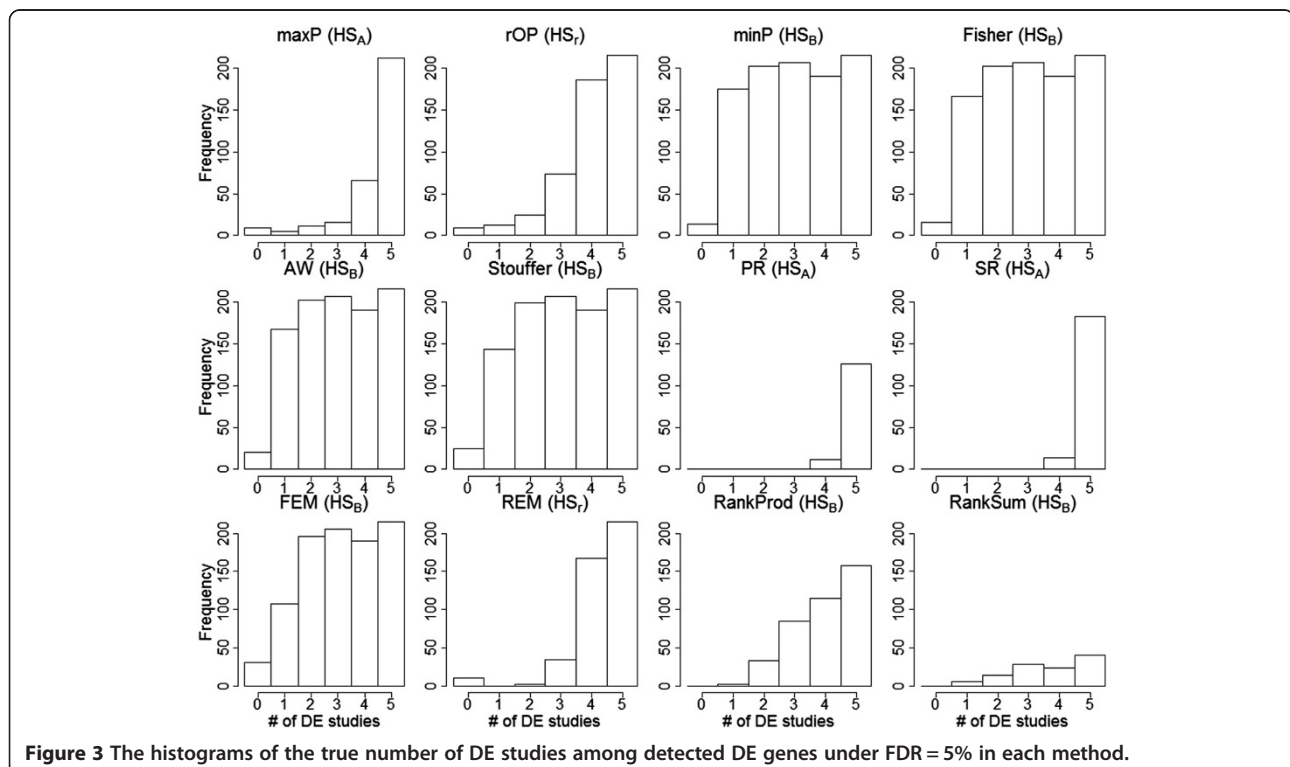| | maxP | rOP | minP | Fisher | AW | Stouffer |
|---|---|---|---|---|---|---|
| Detected # | 321 | 522 | 1005 | 1000 | 1000 | 974 |
| (se) | (2.2) | (2.35) | (0.85) | (1.06) | (1.05) | (1.5) |
| True FDR ($HS_A$) | .068 | .018 | .447 | .444 | .444 | .43 |
| (se) | (.0008) | (.0012) | (.0006) | (.0007) | (.0008) | (.0009) |
| True FDR ($HS_B$) | .007 | .011 | .016 | .017 | .016 | .022 |
| (se) | (.0005) | (.0004) | (.0006) | (.0006) | (.0007) | (.0006) |
| AUC ($HS_A$) | .996 | .964 | .8 | .82 | .79 | .89 |
| (se) | (.0003) | (.0014) | (.0005) | (.0005) | (.0005) | (.0006) |
| AUC ($HS_B$) | .75 | .833 | .99 | .99 | .99 | .99 |
| (se) | (.0013) | (.01) | (.0001) | (.0001) | (.0001) | (.0005) |
| Characterization | $HS_A$ | $HS_r$ | $HS_B$ | $HS_B$ | $HS_B$ | $HS_B$ |
| | **PR** | **SR** | **FEM** | **REM** | **RankProd** | **RankSum** |
| Detected # | 136 | 186 | 948 | 411 | 391 | 105 |
| (se) | (2.51) | (2.3) | (1.75) | (2.86) | (3.31) | (1.514) |
| True FDR ($HS_A$) | .008 | .01 | .415 | .117 | .13 | .389 |
| (se) | (.0003) | (.0004) | (.0009) | (.0015) | (.0014) | (.0008) |
| True FDR ($HS_B$) | 0 | 0 | .022 | .007 | 0 | 0 |
| (se) | (0) | (0) | (.0007) | (.0004) | (0) | (0) |
| AUC ($HS_A$) | .986 | .99 | .917 | .99 | .916 | .504 |
| (se) | (.0003) | (.0002) | (.0009) | (.0002) | (.0011) | (.0046) |
| AUC ($HS_B$) | .981 | .95 | .984 | .92 | .934 | .496 |
| (se) | (.0004) | (.0008) | (.0004) | (.0011) | (.0012) | (.0025) |
| Characterization | $HS_A$ | $HS_A$ | $HS_B$ | $HS_r$ | $HS_B$ | $HS_B$ |

In the simulation study, we had 1,000 non-DE genes in all five studies ($k_g = 0$), and 1,000 genes were differentially expressed in $1 \sim 5$ studies ($k_g = 1, 2, 3, 4, 5$). On average, we had roughly the same number ($\sim 200$) of genes in each group of $k_g = 1, 2, 3, 4, 5$. See Additional file 1: Figure S2 for the heatmap of a simulated example (red colour represents up-regulated genes). We applied the 12 meta-analysis method under FDR control at 5%. With the knowledge of true $k_g$, we were able to derive the sensitivity and specificity for $HS_A$ and $HS_B$, respectively. In $HS_A$, genes with $k_g = 5$ were the underlying true positives and genes with $k_g = 0 \sim 4$ were the underlying true negatives; in $HS_B$, gene with $k_g = 1 \sim 5$ were the underlying true positives and genes with $k_g = 0$ were the true negatives. By adjusting the decision cut-off, the receiver operating characteristic (ROC) curves and the resulting area under the curve (AUC) were used to evaluate the performance. We simulated 50 data sets and reported the means and standard errors of the AUC values. AUC values range between 0 and 1. AUC = 50% represents a random guess and AUC = 1 reaches the perfect prediction. The above simulation scheme only considered the concordance effect sizes (i.e. all with up-regulation when a gene is DE in a study) among five simulated studies. In many applications, some genes may have *p*-value statistical significance in the meta-analysis but the effect sizes are discordant (i.e. a gene is up-regulation in one study but down-regulation in another study). To investigate that

effect, we performed a second simulation that considers random discordant cases. In step 5, the $\mu_{gk}$ became a mixture of two uniform distributions: $\pi_{gk}$ Unif $\cdot [-3, -0.5] + (1 - \pi_{gk}) \cdot$ Unif$[0.5, 3]$, where $\pi_{gk}$ is the probability of gene $g$ ($1 \le g \le 2,000$) in study $k (1 \le k \le 5)$ to have a discordant effect size (down-regulated). We set $\pi_{gk} = 0.2$ for the discordant simulation setting.

## Simulation results to characterize the methods

The simulation study provided the underlying truth to characterize the meta-analysis methods according to their strengths and weaknesses for detecting DE genes of different hypothesis settings. The performances of 12 methods were evaluated by receiver operating characteristic (ROC) curves, which is a visualization tool that illustrates the sensitivity and specificity trade-off, and the resulting area under the ROC curve (AUC) under two different hypothesis settings of $HS_A$ and $HS_B$. Table 1 shows the detected number of DE genes under nominal FDR at 5%, the true FDR and AUC values under $HS_A$ and $HS_B$ for all 12 methods. The values were averaged over 50 simulations and the standard errors are shown in the parentheses.

Figure 3 shows the histogram of the true number of DE studies (i.e. $k_g$) among the detected DE genes under FDR = 5% for each method. It is clearly seen that minP, Fisher, AW, Stouffer and FEM detected $HS_B$-type DE genes and had high AUC values under $HS_B$ criterion (0.98-0.99), compared to lower AUC values under $HS_A$



**Figure 3** The histograms of the true number of DE studies among detected DE genes under FDR = 5% in each method.
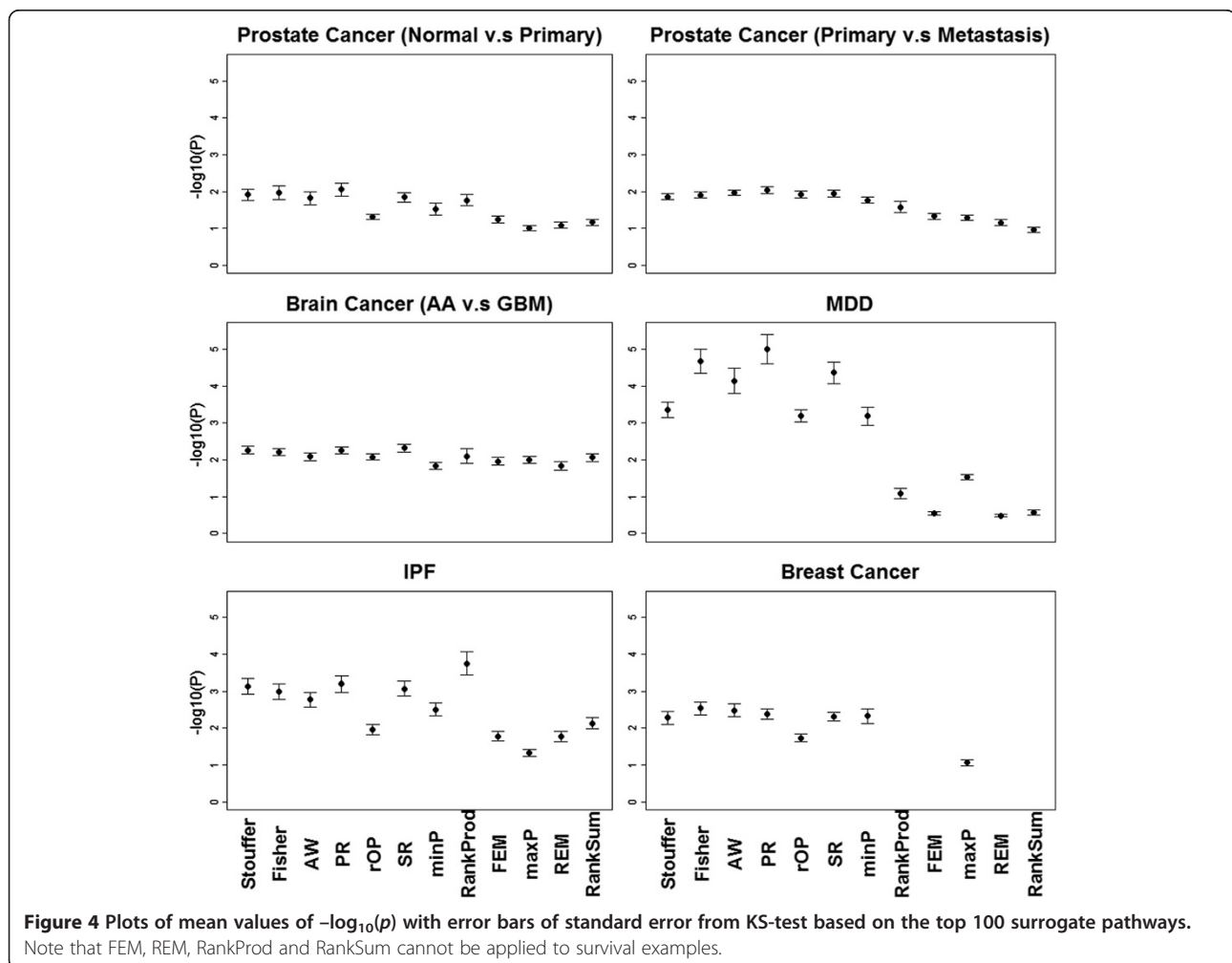
criterion (0.79-0.9). For these methods, the true FDR for $HS_A$ generally lost control (0.41- 0.44). On the other hand, maxP, rOP and REM had high AUC under $HS_A$ criterion (0.96-0.99) (true FDR = 0.068-0.117) compared to $HS_B$ (0.75-0.92). maxP detected mostly $HS_A$-type of markers and rOP and REM detected mostly $HS_r$-type DE genes. PR and SR detected mostly $HS_A$-type DE genes but they surprisingly had very high AUC under both $HS_A$ and $HS_B$ criteria. The RankProd method detected DE genes between $HS_r$ and $HS_B$ types and had a good AUC value under $HS_B$. The RankSum detected $HS_B$-type DE genes but had poor AUC values (0.5) for both $HS_A$ and $HS_B$. Table 1 includes our concluding characterization of the targeted hypothesis settings for each meta-analysis method (see also Additional file 1: Figure S5 of the ROC curve and AUC of $HS_A$-type and $HS_B$-type in 12 meta-analysis methods). Additional file 1: Figure S3 shows the result for the second discordant simulation setting. The numbers of studies with opposite effect size are represented by different colours in histogram plot (green: all studies with concordance effect

size; blue: one study has opposite effect size with the remaining; red: two studies have opposite effect size with the remaining). In summary, almost all meta-analysis methods could not avoid inclusion of genes with opposite effect sizes. Particularly, methods utilizing *p*-values from two-sided tests (e.g. Fisher, AW, minP, maxP and rOP) could not distinguish direction of effect sizes. Stouffer was the only method that accommodated the effect size direction in its z-transformation formulation but its ability to avoid DE genes with discordant effect sizes seemed still limited. Owen (2009) proposed a one-sided correction procedure for Fisher's method to avoid detection of discordant effect sizes in meta-analysis [28]. The null distribution of the new statistic, however, became difficult to derive. The approach can potentially be extended to other methods and more future research will be needed for this issue.

## Results of the four evaluation criteria
### Detection capability
Figure 2 shows the number of DE genes identified by each of the 12 meta-analysis methods (FDR = 10% for MDD



**Figure 4** Plots of mean values of $-\log_{10}(p)$ with error bars of standard error from KS-test based on the top 100 surrogate pathways.
Note that FEM, REM, RankProd and RankSum cannot be applied to survival examples.

and breast cancer due to their weak signals and FDR = 1% for all the others). Each plot shows mean with standard error bars for 50 bootstrapped data sets. Additional file 1: Table S4 shows the MSR and ASR for each method in the six examples. The methods in Figure 2 are ordered according to their ASR values. The top six methods with the strongest detection capability were those that detected $HS_B$-type DE genes from the conclusion of Table 1: Fisher, AW, Stouffer, minP, FEM and RankSum. The order of performance of these six methods was pretty consistent across all six examples. The next four methods were rOP, RankProd, maxP and REM and they targeted on either $HS_r$ or $HS_A$. PR and SR had the weakest detection capability, which was consistent with the simulation result in Table 1.
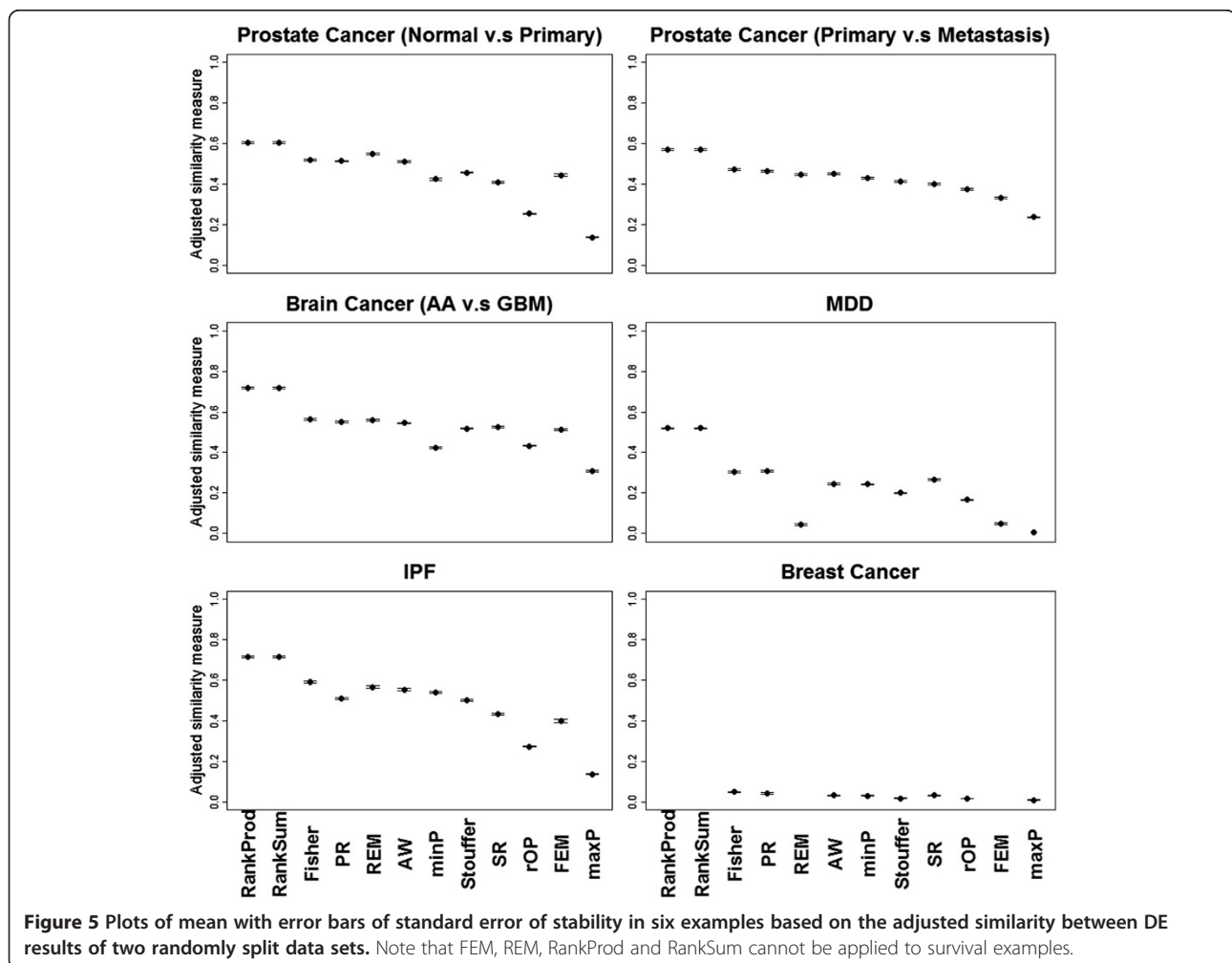
### Biological association

Figure 4 shows plots of mean with standard error bars from the pathway association *p*-values (minus log-transformed) of the top 100 surrogate disease-related pathways for the 12 methods. Additional file 1: Table S5 shows

the corresponding MSR and ASR. We found that Stouffer, Fisher and AW had the best performance among the 12 methods. Surprisingly we found that although PR and SR had low detection capability in simulation and real data, they consistently had relatively high biological association results. This may be due to the better DE gene ordering results these two methods provide, as was also shown by the high AUC values under both hypothesis settings in the simulation.

### Stability

Figure 5 shows the plots of mean with standard error bars of stability calculated by adjusted DE similarity measure. Additional file 1: Table S6 contains the corresponding MSR and ASR. In summary, RankProd and RankSum methods were the most stable meta-analysis methods probably because these two nonparametric approaches take into account all possible fold change calculations between cases and controls. They do not need any distributional assumptions, which provided stability even when sample sizes were small [29]. The maximum



**Figure 5 Plots of mean with error bars of standard error of stability in six examples based on the adjusted similarity between DE results of two randomly split data sets.** Note that FEM, REM, RankProd and RankSum cannot be applied to survival examples.

*p*-value method consistently had the lowest stability in all data sets, which is somewhat expected. For a given candidate marker with a small maximum *p*-value, the chance that at least one study has significantly inflated *p*-values is high when sample size is reduced by half. The stability measures in the breast cancer example were generally lower than other examples. This is mainly due to the weak signals for survival outcome association, which might be improved if larger sample size is available.
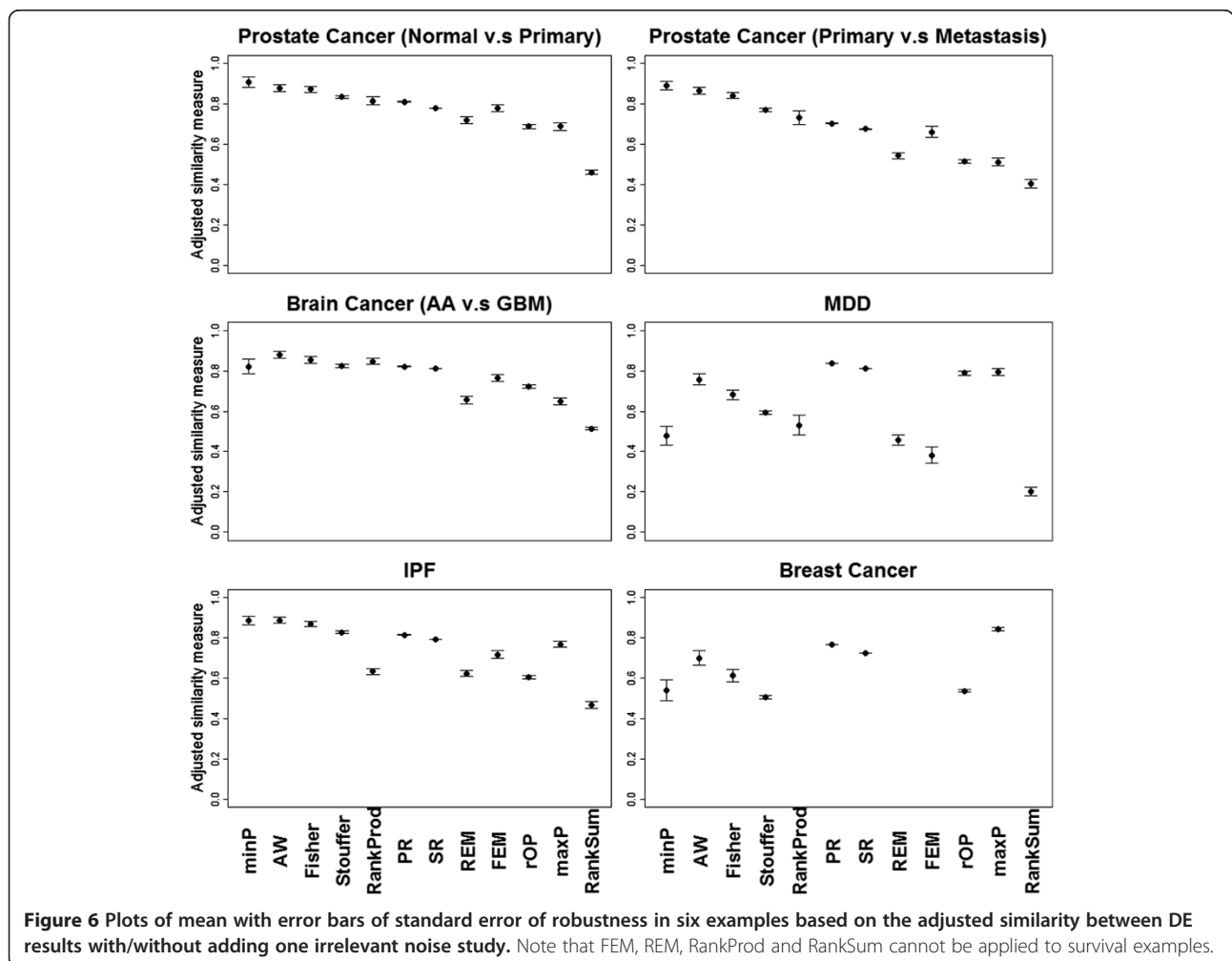
### Robustness

Figure 6 shows the plots of mean with standard error bars of robustness calculated by adjusted DE similarity measure between the original meta-analysis and the new meta-analysis with an added outlier. Additional file 1: Table S7 shows the corresponding MSR and ASR values. In general, methods suitable for $HS_B$ (minP, AW, Fisher and Stouffer) have better robustness than methods for $HS_A$ or $HS_r$ (e.g. maxP and rOP). The trend is consistent in the prostate cancer, brain cancer and IPF examples

but is more variable in the weak-signal MDD and breast cancer examples. RankSum was surprisingly the most sensitive method to outliers, while RankProd performs not bad.

### Characterization of methods by MDS plots

We applied the adjusted DE similarity measure to quantify the similarity of the DE gene orders from any two meta-analysis methods. The resulting dissimilarity measure (i.e. one minus adjusted similarity measure) was used to construct the multidimensional scaling (MDS) plot, showing the similarity/dissimilarity structure between the 12 methods in a two-dimensional space. When two methods were close to each other, they generated similar DE gene ordering. The patterns of MDS plots from six examples generated quite consistent results (Additional file 1: Figure S6). Figure 1(a) shows an aggregated MDS plot where the input dissimilarity matrix is averaged from the six examples. We clearly observed that Fisher, AW, Stouffer, minP, PR and SR were consistently clustered together in all six individual and the aggregated



**Figure 6 Plots of mean with error bars of standard error of robustness in six examples based on the adjusted similarity between DE results with/without adding one irrelevant noise study.** Note that FEM, REM, RankProd and RankSum cannot be applied to survival examples.

MDS plot (labeled in red). This is not surprising given that these methods all sum transformed $p$-value evidence across studies (except for minP). Two methods to combine effect sizes and two methods to combine ranks (FEM, REM, RankProd and RankSum labeled in blue) are consistently clustered together. Finally, the maxP and rOP methods seem to form a third loose cluster (labeled in green).

### Characterization of data sets by entropy measure

From the simulation study, selection of a most suitable meta-analysis method depends on the hypothesis setting behind the methods. The choice of a hypothesis setting mostly depends on the biological purpose of the analysis; that is, whether one aims to detect candidate markers differentially expressed in "all" ($HS_A$), "most" ($HS_r$) or "one or more" ($HS_B$) studies. However, when no biological prior information or preference exists, the entropy measure can be objectively used to determine the choice of hypothesis setting. The analysis identifies the top 1,000 genes from Fisher's meta-analysis method and the gene-specific entropy of each gene is calculated. When the entropy is small, the $p$-values are small in only one or very few studies. Conversely, when the entropy is large, most or all of the studies have small $p$-values. Figure 1(b) shows the box-plots of entropy of the top 1,000 candidate genes identified by Fisher's method in the six data sets. The result shows that prostate cancer comparing primary and metastatic tumor samples had the smallest entropy values, which indicated high heterogeneity across the three studies and that $HS_B$ should be considered in the meta-analysis. On the other hand, MDD had the highest entropy values. Although the signals of each MDD study were very weak, they were rather consistent across studies and application of $HS_A$ or $HS_r$ was adequate. For the other examples, we suggest

using the robust $HS_r$ unless other prior biological purpose is indicated.

## Conclusions and discussions
### An application guideline for practitioners

From the simulation study, the 12 meta-analysis methods were categorized into three hypothesis settings ($HS_A$, $HS_B$ and $HS_r$), showing their strengths for detecting different types of DE genes in the meta-analysis (Figure 3 and the second column of Table 2). For example, maxP is categorized to $HS_A$ since it tends to detect only genes that are differentially expressed in all studies. From the results using four evaluation criteria, we summarized the rank of ASR values (i.e. the order used in Figures 2 and 6) and calculated the rank sum of each method in Table 2. The methods were then sorted first by the hypothesis setting categories and then by the rank sum. The clusters of methods from the MDS plot were also displayed. For methods in the $HS_A$ category, we surprisingly see that the maxP method performed among the worst in all four evaluation criteria and should be avoided. PR was a better choice in this hypothesis setting although it provides a rather weak detection capability. For $HS_B$, Fisher, AW and Stouffer performed very well in general. Among these three methods, we note that AW has an additional advantage to provide an adaptive weight index that indicates the subset of studies contributing to the meta-analysis and characterizes the heterogeneity (e.g. adaptive weight (1,0,…) indicates that the marker is DE in study 1 but not in study 2, etc.). As a result, we recommend AW over Fisher and Stouffer in the $HS_B$ category. For $HS_r$, the result was less conclusive. REM provided better stability and robustness but sacrificed detection capability and biological association. On the other hand, rOP obtained better detection capability and biological association but was neither stable nor robust. In general, since detection

**Table 2 Ranks of method performance in the four evaluation criteria**

| | Targeted HS | Power detection | Biological association | Stability | Robustness | Rank Sum | MDS[*1] |
|---|---|---|---|---|---|---|---|
| **PR** | $HS_A$ | 12 | 4 | 4 | 6 | 26 | 1 |
| SR | $HS_A$ | 11 | 6 | 9 | 7 | 33 | 1 |
| maxP | $HS_A$ | 9 | 10 | 12 | 11 | 42 | 2 |
| **rOP** | $HS_r$ | 7 | 5 | 10 | 10 | 32 | 2 |
| REM | $HS_r$ | 10 | 11 | 5 | 8 | 34 | 3 |
| **Fisher** | $HS_B$ | 1 | 2 | 3 | 3 | 9 | 1 |
| AW | $HS_B$ | 2 | 3 | 6 | 2 | 13 | 1 |
| Stouffer | $HS_B$ | 3 | 1 | 8 | 4 | 16 | 1 |
| minP | $HS_B$ | 4 | 7 | 7 | 1 | 19 | 1 |
| RankProd | $HS_B$ | 8 | 8 | 1 | 5 | 22 | 3 |
| RankSum | $HS_B$ | 6 | 12 | 2 | 12 | 32 | 3 |
| FEM | $HS_B$ | 5 | 9 | 11 | 9 | 34 | 3 |

[*1] Cluster number of methods in the MDS plot (Figure 1(a)).

capability and biological association are of more importance in the meta-analysis and rOP has the advantage to link the choice of $r$ in $HS_r$ with the rOP method (e.g. when $r = 0.7 \cdot K$, we identify genes that are DE in more than 70% of studies), we recommend rOP over REM.

Below, we provide a general guideline for a practitioner when applying microarray meta-analysis. Data sets of a relevant biological or disease hypothesis are firstly identified, preprocessed and annotated according to Step (i) - (v) in Ramasamy et al. Proper quality assessment should be performed to exclude studies with problematic quality (e.g. with the aid of MetaQC as we did in the six examples). Based on the experimental design and biological objectives of collected data, one should determine whether the meta-analysis aims to identify biomarkers differentially expressed in all studies ($HS_A$), in one or more studies ($HS_B$) or in majority of studies ($HS_r$). In general, if higher heterogeneity is expected from, say, heterogeneous experimental protocol, cohort or tissues, $HS_B$ should be considered. For example, if the combined studies come from different tissues (e.g. the first study uses peripheral blood, the second study uses muscle tissue and so on), tissue-specific markers may be expected and $HS_B$ should be applied. On the contrary, if the collected studies are relatively homogeneous (e.g. use the same array platform or from the same lab), $HS_r$ is generally recommended, as it provides robustness and detects consistent signals across the majority of studies. In the situation that no prior knowledge is available to choose a desired hypothesis setting or if the researcher is interested in a data-driven decision, the entropy measure in Figure 1(b) can be applied and the resulting box-plot can be compared to the six examples in this paper to guide the decision. Once the hypothesis setting is determined, the choice of a meta-analysis method can be selected from the discussion above and Table 2.

## Conclusions

In this paper, we performed a comprehensive comparative study to evaluate 12 microarray meta-analysis methods using simulation and six real examples with four evaluation criteria. We clarified three hypothesis settings that were implicitly assumed behind the methods. The evaluation results produced a practical guideline to inform biologists the best choice of method(s) in real applications.

With the reduced cost of high-throughput experiments, data from microarray, new sequencing techniques and mass spectrometry accumulate rapidly in the public domain. Integration of multiple data sets has become a routine approach to increase statistical power, reduce false positives and provide more robust and validated conclusions. The evaluation in this paper focuses on microarray meta-analysis but the principles and messages apply to other types of genomic meta-analysis

(e.g. GWAS, methylation, miRNA and eQTL). When the next-generation sequencing technology becomes more affordable in the future, sequencing data will become more prevalent as well and similar meta-analysis techniques will apply. For these different types of genomic meta-analysis, similar comprehensive evaluation could be performed and application guidelines should be established as well.

## Additional file

**Additional file 1: Supplementary methods of (1) Adaptive weighted (AW) Fisher, (2) Combined statistical estimates (effect size) methods of FEM and REM, (3) Combined rank statistics methods: Rank Product (RankProd) and Rank Sum (RankSum). Figure S1.** MetaQC. **Figure S2.** Heatmap of simulated example (red color represents up-regulated genes). **Figure S3.** The histograms of the true number of DE studies among detected DE genes under FDR = 5% in each method for discordance case. **Figure S4.** Cumulative moving average to determine D = 100. **Figure S5.** The ROC curves and AUC for the hypothesis settings of $HS_A$-type and (red line) $HS_B$-type (black line) in each meta-analysis method. **Figure S6.** Multidimensional scaling (MDS) plots of individual data sets. **Figure S7.** Stability and Robustness plot for α = 0.0001, 0.005 and 0.01. **Table S1.** Detailed data sets description. **Table S2.** MetaQC results. **Table S3.** Data sets and number of matched genes. **Table S4.** Mean standardized rank (MSR) and aggregated standardized rank (ASR) for detection capability. **Table S5.** Mean standardized rank (MSR) and aggregated standardized rank (ASR) for biological association. **Table S6.** Mean standardized rank (MSR) and aggregated standardized rank (ASR) for stability. **Table S7.** Mean standardized rank (MSR) and aggregated standardized rank (ASR) for robustness.

**Authors' contributions**
GCT supervised the whole project. LCC developed all statistical analysis and HML developed partial statistical analysis. GCT and LCC drafted the manuscript. All authors read and approved the final manuscript.

**Author details**
[1]Department of Biostatistics, Graduate school of Public Health, University of Pittsburgh, Pittsburgh, PA, USA. [2]Department of Psychiatry, School of Medicine, University of Pittsburgh, Pittsburgh, PA, USA. [3]Department of Human Genetics, Graduate School of Public Health, University of Pittsburgh, Pittsburgh, PA, USA.

**References**
1. Ramasamy A, Mondry A, Holmes CC, Altman DG: **Key issues in conducting a meta-analysis of gene expression microarray datasets.** *PLoS Med* 2008, **5**(9):e184.
2. Kang DD, Sibille E, Kaminski N, Tseng GC: **MetaQC: objective quality control and inclusion/exclusion criteria for genomic meta-analysis.** *Nucleic Acids Res* 2012, **40**(2):e15.
3. Tseng GC, Ghosh D, Feingold E: **Comprehensive literature review and statistical considerations for microarray meta-analysis.** *Nucleic Acids Res* 2012, **40**(9):3785–3799.
4. Hong F, Breitling R: **A comparison of meta-analysis methods for detecting differentially expressed genes in microarray experiments.** *Bioinformatics* 2008, **24**(3):374–382.

5.    Campain A, Yang YH: **Comparison study of microarray meta-analysis methods.** *BMC Bioinforma* 2010, **11**:408.
6.    Birnbaum A: **Combining independent tests of significance.** *J Am Stat Assoc* 1954, **49**:559–574.
7.    Li J, Tseng GC: **An adaptively weighted statistic for detecting differential gene expression when combining multiple transcriptomic studies.** *Ann Appl Stat* 2011, **5**:994–1019.
8.    Tusher VG, Tibshirani R, Chu G: **Significance analysis of microarrays applied to the ionizing radiation response.** *Proc Natl Acad Sci USA* 2001, **98**(9):5116–5121.
9.    Pesarin F, Salmaso L: *Permutation tests for complex data: theory, applications and software.* Ames, IA 50010: Wiley.com; 2010.
10.   Cox DR: **Regression models and life-tables.** *J R Stat Soc Ser B* 1972, **34**(2):187–220.
11.   Benjamini Y, Hochberg Y: **Controlling the false discovery rate: a practical and powerful approach to multiple testing.** *J R Stat Soc Ser B* 1995, **57**(1):289–300.
12.   Wang X, Kang DD, Shen K, Song C, Lu S, Chang LC, Liao SG, Huo Z, Tang S, Ding Y, *et al*: **An R package suite for microarray meta-analysis in quality control, differentially expressed gene analysis and pathway enrichment detection.** *Bioinformatics* 2012, **28**(19):2534–2536.
13.   Fisher RA: *Statistical methods for research workers.* Edinburgh: Genesis Publishing, Oliver and Boyd; 1925.
14.   Stouffer SA: **A study of attitudes.** *Sci Am* 1949, **180**(5):11–15.
15.   Tippett LHC: **The Methods of Statistics. An introduction mainly for workers in the biological sciences**. In *The Methods of Statistics An Introduction mainly for Workers in the Biological Sciences.* London: Williams &Norgate Ltd; 1931.
16.   Wilkinson B: **A statistical consideration in psychological research.** *Psychol Bull* 1951, **48**(3):156–158.
17.   Song C, Tseng GC: **Order statistics for robust genomic meta-analysis.** *Ann Appl Stat* 2012. Accepted.
18.   Choi JK, Yu U, Kim S, Yoo OJ: **Combining multiple microarray studies and modeling interstudy variation.** *Bioinformatics* 2003, **19**(Suppl 1):i84–i90.
19.   Hong F, Breitling R, McEntee CW, Wittner BS, Nemhauser JL, Chory J: **RankProd: a bioconductor package for detecting differentially expressed genes in meta-analysis.** *Bioinformatics* 2006, **22**(22):2825–2827.
20.   Dreyfuss JM, Johnson MD, Park PJ: **Meta-analysis of glioblastoma multiforme versus anaplastic astrocytoma identifies robust gene markers.** *Mol Cancer* 2009, **8**:71.
21.   Borg I: *Modern multidimensional scaling: theory and applications.* New York, NY 10013: Springer; 2005.
22.   Martin NF, England JW: *Mathematical theory of entropy, vol. 12.* Cambridge CB2 8BS United Kingdom: Cambridge University Press; 2011.
23.   Wu W, Dave N, Tseng GC, Richards T, Xing EP, Kaminski N: **Comparison of normalization methods for CodeLink Bioarray data.** *BMC Bioinforma* 2005, **6**:309.
24.   Spearman C: **The proof and measurement of association between two things.** *Amer J Psychol* 1904, **15**:72–101.
25.   Li Q, Brown JB, Huang H, Bickel PJ: **Measuring reproducibility of high-throughput experiments.** *Ann Appl Stat* 2011, **5**(3):1752–1779.
26.   Yang X, Bentink S, Scheid S, Spang R: **Similarities of ordered gene lists.** *J Bioinform Comput Biol* 2006, **4**(3):693–708.
27.   Hubert L, Arabie P: **Comparing partitions.** *J Classification* 1985, **2**:193–218.
28.   Owen AB: **Karl Pearson's meta-analysis revisited.** *Ann Statistics* 2009, **37**(6B):3867–3892.
29.   Breitling R, Herzyk P: **Rank-based methods as a non-parametric alternative of the T-statistic for the analysis of biological microarray data.** *J Bioinform Comput Biol* 2005, **3**(5):1171–1189.