# The role and robustness of the Gini coefficient as an unbiased tool for the selection of Gini genes for normalising expression profiling data

[1]Marina Wright Muelas*, [1]Farah Mughal, [2,3]Steve O'Hagan, [3,4]Philip J. Day & [1,5]*Douglas B. Kell

[1]Department of Biochemistry, Institute of Integrative Biology, Faculty of Health and Life Sciences, University of Liverpool, Crown Street, Liverpool, L69 7ZB, UK

[2]School of Chemistry, [3]The Manchester Institute of Biotechnology, 131, Princess St, Manchester M1 7DN, UK. [4]Faculty of Biology, Medicine and Health, The University of Manchester M13 9PL, UK.

[5]Novo Nordisk Foundation Centre for Biosustainability, Technical University of Denmark, 10 Building 220, Kemitorvet, 2800 Kgs. Lyngby, Denmark

Emails in order of authorship: m.wright-muelas@liverpool.ac.uk, Farah.Mughal@liverpool.ac.uk, SOhagan@manchester.ac.uk , Philip.J.Day@manchester.ac.uk , dbk@liv.ac.uk

*corresponding authors: m.wright-muelas@liverpool.ac.uk and dbk@liv.ac.uk

17

18

46

47

48

49

## Abstract

51

52    We recently introduced the Gini coefficient (GC) for assessing the expression variation of a particular gene
53    in a dataset, as a means of selecting improved reference genes over the cohort ('housekeeping genes')
54    typically used for normalisation in expression profiling studies. Those genes (transcripts) that we
55    determined to be useable as reference genes differed greatly from previous suggestions based on
56    hypothesis-driven approaches. A limitation of this initial study is that a single (albeit large) dataset was
57    employed for both tissues and cell lines.

58    We here extend this analysis to encompass seven other large datasets. Although their absolute values differ
59    a little, the Gini values and median expression levels of the various genes are well correlated with each
60    other between the various cell line datasets, implying that our original choice of the more ubiquitously
61    expressed low-Gini-coefficient genes was indeed sound. In tissues, the Gini values and median expression
62    levels of genes showed a greater variation, with the GC of genes changing with the number and types of
63    tissues in the data sets. In all data sets, regardless of whether this was derived from tissues or cell lines, we
64    also show that the GC is a robust measure of gene expression stability. Using the GC as a measure of
65    expression stability we illustrate its utility to find tissue- and cell line-optimised housekeeping genes
66    without any prior bias, that again include only a small number of previously reported housekeeping genes.
67    We also independently confirmed this experimentally using RT-qPCR with 40 candidate GC genes in a panel
68    of 10 cell lines. These were termed the Gini Genes.

69    In many cases, the variation in the expression levels of classical reference genes is really quite huge (e.g. 44
70    fold for GAPDH in one data set), suggesting that the cure (of using them as normalising genes) may in some
71    cases be worse than the disease (of not doing so).  We recommend the present data-driven approach for
72    the selection of reference genes by using the easy-to-calculate and robust GC.

73

74    **Keywords:** housekeeping genes – reference genes – Gini index – Gene Expression

75

76

77

## Background

In a recent paper [1], we introduced the Gini index (or Gini coefficient, GC) [2-5] as a very useful, nonparametric statistical measure for identifying those genes whose expression varied least across a large set of samples (when normalised appropriately [6] to the total expression level of transcripts). The GC is a measure that is widely used in economics (e.g. [4, 7-12]) to describe the (in)equality of the distribution of wealth or income between individuals in a population. However, although it could clearly be used to describe the variation in any other property between individual examples [13-16]), it has only occasionally been used in biochemistry [1, 5, 17-22]. Its visualisation and calculation are comparatively straightforward (Fig 1): individual examples are ranked on the abscissa in increasing order of the size of their contribution, and the cumulative contribution is plotted against this on the ordinate. The GC is given by the fractional area mapped out by the resulting 'Lorenz' curve (Fig 1). For a purely 'socialist' system in which all contributions are equal (GC = 0), the curve joins the normalised 0,0 and 1,1 axes, while for a complete 'autocracy', in which the resource or expression is held or manifest by only a single individual (GC=1), the 'curve' follows the two axes (0,0 → 1,0 → 1,1).

Since the early origins of large-scale nucleic acid expression profiling, especially those using microarrays [23-25], it has been clear that expression profiling methods are susceptible to a variety of more or less systematic artefacts within an experiment, whose resolution would require or benefit from some kind of normalisation (e.g. [26-36]). By this ('normalisation of the first kind'), and what is typically done, we mean the smoothing out of genuine artefacts within an arrray or a run, that occur simply due to differences in temperature or melting temperature or dye binding or hybridisation and cross-hybridisation efficiency (and so on) across the surface of the array. This process can in principle use reference genes, but usually exploits smoothing methods that normalise geographically local subsets of the genes to a presumed distribution.

Even after this is done, there is a second level of normalisation, that between chips or experiments, that is usually done separately, not least because it is typically much larger and more systematic, especially because of variations in the total amount of material in the sample analysed or of the overall sensitivity of the detector (much as is true of the within-run versus between-run variations observed in mass spectrometry experiments [37, 38]). This kind of normalising always requires 'reference' genes whose expression varies as little as possible in response to any changes in experimental conditions. The same is true for expression profiling as performed by qPCR [39-44], where the situation is more acute regarding the choice of reference genes since primers must be selected for these *a priori*. Commonly, the geometric mean of the expression levels of that or those that vary the least is selected as the 'reference'. The question then arises as to which are the premium 'reference' genes to choose.

Perhaps surprisingly [45], rather than simply letting the data speak for themselves, choices of candidate reference genes were often made on the basis that reference genes should be 'housekeeping' genes that would simply be assumed ('hypothesised') to vary comparatively little between cells, be involved in nominal routine metabolism and also that they should have a reasonably high expression level (e.g. [46-63]). This is not necessarily the best strategy, and there is in fact (and see below) quite a wide degree of variation of the expression of most standard housekeeping genes between cells or tissues (e.g. [50, 59, 62, 64-76]). Indeed, Lee et al [66] stated explicitly that housekeeping genes may be uniformly expressed in certain cell types but may vary in others, especially in clinical samples associated with disease.

It became obvious that an analysis of the GC of the various genes was actually precisely what was required to assess those 'housekeeping' (or any other) genes that varied least across a set of expression profiles, and

121     we found 35 transcripts for which the GC was 0.15 or below when assessing 56 mammalian cell lines taken
122     from a wide variety of tissues [1]. These we refer to as the 'Gini genes'. Most of these were 'novel' as they
123     had never previously been considered as reference genes, and we noted that their Gini indices were
124     significantly smaller (they were more stably expressed) than were those of the more commonly used
125     reference genes [63]. However, this analysis was done on only one (albeit large) dataset of gene expression
126     profiles. While some of the compilations (e.g. [62, 77]) contain massive amounts of expression profiling
127     data, many of these, especially the older ones, may well be of uncertain quality. Thus, especially since the
128     GC is very prone to being raised by small numbers of large outliers, we decided for present purposes that
129     we should compare our analyses of candidate Gini genes using a smaller but carefully chosen set of
130     expression profiling experiments. The more modern RNA-seq (e.g. [78-82]), in which individual transcripts
131     are simply counted digitally via direct sequencing, is seen as considerably more robust [78, 83, 84] and
132     sensitive [85, 86], and so we selected additional large and recent datasets that used RNA-seq in cell lines
133     and tissues (Table *1*). We note too that the precision of these digital methods (as with other, digital, single-
134     molecule strategies [87-89]), means that the requirement for reasonably high-level expression levels is
135     much less acute.

136     In a similar vein (Table 2), we selected a small number of reasonably detailed studies in which particular
137     housekeeping genes had been proposed as reference genes.

138     To our knowledge, there are no large-scale studies to determine housekeeping genes in large, cell-line
139     cohorts; the present paper serves to provide one. In addition, we include an experimental RT-qPCR analysis
140     of a subset of the Gini genes.

141

# Results

## The Gini Coefficient as a robust measure of gene expression stability in multiple cell-line data sets

We previously identified a number of genes in the Human Protein Atlas (HPA) cell line data set [90] with very low expression variability and thus potential for use as reference genes [91]. However, we did not compare these Gini genes to other genes that have previously been proposed as housekeeping genes. We therefore performed a similar analysis using the potential housekeeping genes we proposed in [91] as well as other reference genes proposed in other studies (Table 2) with additional large RNA-Seq cell line data sets (Table *1*).

Fig 2A shows a plot of the GC of a variety of candidate Gini genes versus their median expression level in the HPA cell lines dataset set [90]. It is clear that genes we identified previously have much lower GC values in the HPA dataset than do any of the others (just two, VPS29 and CHMP2A, were also identified by Eisenberg and Levenson and another, RPL41, by Caracausi). This is not at the expense of an unusually low expression (Fig 2A), a finding broadly confirmed when we look at the median expression levels for the CCLE dataset (Fig 2B) and of the Klijn dataset (Fig 2C).

Fig 3 shows the GC values for the various genes in two other datasets, viz CCLE and Klijn. Our previous Gini genes have a lower GC than that of any of the other housekeeping genes in 25 out of 38 cases in Klijn (all under 0.2) and in 26 out of 40 cases for CCLE (all under 0.22). In confirmation of this, and of the correlation found above between the median expression levels in CCLE and Klijn, the GC values are also well correlated with each other for the two datasets (Fig 3). Thus, although the absolute numbers are slightly larger than are those for the HPA dataset (unsurprisingly, given the much larger number of examples), the trend is still very clear: the GiniGenes remain the best among those variously proposed as reference genes in a variety of large and quite independent datasets. It also suggests that variations in the total amount of mRNA are not an issue either.

Another common statistical measure, more resistant to individual outliers, is the interquartile ratio (the ratio between the 25[th] and 75[th] percentile when expression levels are ranked); by this measure too, the Gini genes that we uncovered previously stand out as being the least varying (Fig 4 A and B). This suggests that, as a measure of gene expression stability, the GC is robust: the GiniGenes have the lowest ratio between their maximum and minimum expression values in the HPA dataset (Fig 4C) and also the lowest interquartile ratio in their levels of expression in all three cell line data sets explored here (Fig. 4B and C) with good correlation between these two datasets.

## Use of the Gini Coefficient to find GiniGenes in an unbiased manner in cell-line data sets

Up to now, our analyses of these data sets have used a set of predefined genes to look at expression stability. We next sought to investigate whether the GC would highlight genes with high expression stability that have been reported by others or by ourselves when performing this analysis in a data-driven manner. To that end, we found 115 genes shared between the three data sets with a GC ≤ 0.2 (Fig. 5, 6). This value for the GC was chosen since reducing this to ≤ 0.15 meant no or very few genes were found in some data sets (e.g. no genes in the CCLE data set had a GC ≤ 0.15) and going above this meant the number of genes were unmanageable (e.g. 1051 genes with a GC ≤ 0.21 in the Klijn data set). Of the 115 genes shared

184    between the datasets with GC <0.2, 13 were GiniGenes and two were housekeeping genes defined by
185    Caracausi and colleagues (Fig. 5 B). When we selected the top 20 expressing genes in each data set, only 13
186    of these were common across these data sets; Table 3 shows some descriptive statistics of 13 of these, with
187    descriptive statistics of all 115 genes found in Supplementary Table S1. Of these genes, two (HNRNPK and
188    PCBP1) are GiniGenes and one (SLC25A3) is a gene previously reported by Caracausi et al. Seven out of the
189    13 genes (HNRNPK, HNRNPC, PCBPB, SF3B1, SRSF3, EDF1 and EIF4H) here share important roles in RNA
190    transcription, translation and stability [92-100], are implicated in a number of diseases, including cancer
191    [92, 95, 101-111], and some, such as SRSF3 are essential for embryo development [112]. Given their pivotal
192    functions, it may be unsurprising that the expression of these genes are tightly regulated across cell lines of
193    different tissue origins, even where these are cancer cell lines. Overall, the distribution, expression stability
194    and important functional roles of these genes suggest that these are excellent housekeeping genes across
195    different cell types.

196    Of particular interest to us was finding one gene encoding a mitochondrial phosphate transporter protein
197    (SLC25A3 [113]) to be within this list of the top expressing stably expressed genes. This might seem logical
198    since mitochondrial ATP synthesis is required by all cell types and tissues.

199    Figure 7 shows the robustness of the GC for the subset of 115 genes common between the three data sets
200    studied here with a low GC  (<0.2). Lower Gini coefficients correlate with lower IQR and Max:median ratios
201    (Fig7: only results for the Klijn data set are shown). The range of IQR values of these genes was smaller in
202    the larger two data sets (CCLE, 1.42-1.67; Klijn, 1.30- 1.64) than in the HPA data set (1.26-1.84) suggesting
203    the measured expression values were more stable in the larger data sets (Supplementary Table S1). This
204    may, however, be due to a larger number of cell lines in these two large datasets (934 and 622 in CCLE and
205    Klijn) compared with the HPA data set (56 cell lines).

206

207    Application of the Gini coefficient to human tissue RNA-Seq data sets
208

209    The results presented thus far are representative of human cell lines. Most reports in the literature
210    regarding housekeeping genes refer to tissue expression data. This may be due to the cell lines being
211    "dedifferentiated" with respect to the tissues from which they are derived [114].

212    In our previous report [1] we also analysed RNA-Seq data from tissues [90] and found 22 genes with a GC <
213    0.15, of which 3 (CHMP2A, VPS29 and PCBP1) were also found in cell line data with a GC <0.15. The median
214    expression level and GC of these and other candidate GiniGenes in this tissue data set are shown in Figure
215    8. As with cell line data, the genes we previously identified (GGs, green dots in Fig 8) have much lower GCs
216    in this tissue data set than do any of the other candidate GiniGenes, with only two of these genes (VPS29
217    and CHMP2A) identified previously by Eisenberg & Levenson [115]. The low GC value of these GiniGenes is
218    not at the expense of low expression: of the 22 GiniGenes, 13 are expressed at a median level of between
219    40 and 200 TPM (see Supplementary Table S2). Moreover, the GC was also representative of the variation
220    in expression of these genes (albeit influenced to a lesser extent by outliers), as shown in Fig. 9 A and B,
221    with all GiniGenes having a GC <  0.15 and the lowest RSD (relative standard deviation), ranging from
222    24.096 % to 28.66 % and IQR (1.26 to 1.44) of this list of housekeeping genes. The expression of other
223    housekeeping genes such as GAPDH, ACTB, RPL13A, SDHA, B2M was quite varied according to these
224    measures. For example, the GC of GAPDH (a commonly used HKG) was 0.33, with a RSD of 72.4 % and IQR
225    of 2.24, and for ACTB (another commonly used HKG) these values were 0.29, 55.24 %, and 2.11.

226　The median expression levels of the proposed reference genes show a similar level of correlation between
227　the data sets as was found with the cell line data (Figure S1 A-C), and GiniGenes displayed a mid-range level
228　of expression. The GC of the tissue GiniGenes we proposed however, tended to be higher and more
229　variable in their GC values than in the HPA dataset (Figure S2 A-C) suggesting that those genes may be
230　representative of the HPA data set only. As an example, in the GTEx dataset only 28 genes had a GC < 0.2,
231　of which the majority (17) were those reported by Caracausi and colleagues, and 7 were GiniGenes. The
232　results here are likely influenced by the number and status (disease or normal) of the tissues analysed in
233　the various data sets compared; for example, the GTEx data come from 53 different, normal human tissues,
234　whereas the HPA tissue data include a mixture of disease and normal tissue samples. In addition, compared
235　to the cell line data where hundreds (in the case of the Cancer Cell Line Encyclopedia) of cell lines were
236　analysed, the number of tissues in these data sets was fewer than 100.

237　In the case of the data set used by Eisenberg and Levanon [115], viz. the Illumina Human Body Map (E-
238　MTAB-513),  10 of the 11 housekeeping genes proposed here (which included 2 Gini Genes, CHMP2A and
239　VPS29) had a GC ≤ 0.2 and were reasonably well expressed (with median expression levels between 50-270
240　TPM, see  Supplementary Table S2 and Supplementary Fig S4). This may be compared to the 5 other GGs
241　with GC < 0.2 in this data set whose expression value was lower, with median expression between 19-35
242　TPM. This suggests that finding suitable HKGs may be dependent on the data set itself, and the type of
243　tissue under investigation.

244　We next sought to perform a more comprehensive and integrative analysis by filtering the tissue data sets
245　to only include genes with a GC ≤ 0.2 to find common genes across these data sets with reasonable
246　expression stability (Supplementary Table S3). As shown in Fig 10 only 15 genes were shared between the
247　four data sets with a GC ≤ 0.2, none of which has been reported previously as a housekeeping genes. Table
248　4 shows some descriptive statistics of these genes. In any case, the names of the proteins encoded by these
249　15 genes suggest these play important and essential roles. The median expression values of these genes
250　varied from around 10-450 TPM, with SNX3 (Sorting nexin-3 (Protein SDP3)) and COX4I1 (Cytochrome c
251　oxidase subunit 4 isoform 1) being consistently the two highest-expressing genes.

252　Sorting nexins are a group of cytoplasmic and membrane-associated proteins involved in the regulation of
253　intracellular trafficking [116]. SNX3 has been reported to play a role in receptor recycling and formation of
254　multivesicular bodies [117], and its dysregulation has been implicated in disorders of iron metabolism and
255　the pathogenesis of some neurodegenerative diseases [118, 119].

256　The COX4I gene encodes the nuclear-encoded cytochrome c oxidase subunit 4 isoform 1, the terminal
257　enzyme in the mitochondrial respiratory chain. Given the key role of the mitochondrial respiratory chain in
258　all human cells (except red blood cells), stable expression of such a gene in all tissues may not be a
259　surprising result. Increased RNA COX4I1 levels have been reported in sperm of an obese male rat model
260　[120] and thus may play a role in obesity-related fertility problems, and reduced expression of this subunit
261　leads to a reduction in mitochondrial respiration as well as sensitising cells to apoptosis [121].

262　The small number of genes shared between these data sets with a GC < 0.2 indicates that the data in these
263　studies are more variable compared to cell lines alone. The cause of this variation may be due to the tissue
264　data having been obtained from different subjects [122]. Moreover, tissues are themselves a mixture of cell
265　types with varying levels of gene expression in each cell type [123], while cell lines are nominally clonal.

266　Our results suggest that in the case of RNA-seq tissue data sets, where gene expression tends to be more
267　variable, an unbiased approach, using the Gini coefficient, may be more fruitful in the search for stably

8

268  expressed genes with which to perform normalisation, than the other commonly used methods used until
269  now [122, 124].

270  RT-qPCR analysis of gene expression stability of some housekeeping genes in 10 cell lines
271

272  In order to illustrate the utility of the GC to find suitable housekeeping genes, we next chose to assess this
273  experimentally by RT-qPCR using a small subset of candidate reference genes (40; top 32 genes from genes
274  ordered by GC and expression value from [91], plus 8 of the most commonly used from the literature,
275  including seven from [63] and one (RPL32) from [125][126], and 10 cell lines from a range of tissues (see
276  Table 5 and 6). We first set a Cq value (which is inversely proportional to expression level) cut-off of 32,
277  above which no expression is observed, and subsequently used the Cq values of genes in cell lines as a
278  relative expression level (Cq cut off/Cq value of gene). Descriptive statistics of the expression of each gene
279  in individual cell lines were then calculated. As a final step, the median expression value of each gene in
280  individual cell lines was used to calculate descriptive statistics, including the GC, of gene expression across
281  these cell lines. Figure 11 illustrates a KNIME workflow [127-129] that we wrote for this purpose. The raw
282  data and descriptive statistics extracted are provided in Supplementary Tables S5 and S6 respectively, and
283  the KMNIME analysis workflow in Supplementary File 1.

284  Fig 12 uses RT-qPCR data to plot the GC of the candidate reference genes analysed here versus their
285  relative median expression level. Three GiniGenes [91] (RBM45, TRNT1 and CNOT2) had very low and
286  variable expression. Most of the other genes analysed showed low GC values with a range of (relative)
287  expression values; the inset in Figure 12 shows genes with a GC < 0.2 including a mix of 35 genes: 26
288  GiniGenes and 6 housekeeping genes referenced by Vandesompele and colleagues [63], one referenced by
289  Caracausi [130] and one by Lee et al [131]. Two of these GiniGenes, HNRNPK and PCBP1, which we also
290  found to be stably expressed in the cell line data suggesting these may be potential stable housekeeping
291  genes.  As shown in Figure 13 and inset, the GC is well correlated with the % RSD.

292  More importantly, the GC of our GiniGenes was particularly low (Fig 12). The low absolute magnitude
293  reflected the fact that Cq value is based on a logarithmic scale. Various commonly used housekeeping
294  genes (HPRT1, GAPDH, ACTB, SDHA, HMBS and B2M) displayed higher % RSDs and GC than other genes
295  studied here in spite of their higher relative expression levels. This was also the case when inspecting the
296  interquartile ratio against the GC of these (Figure S3).

297  The above results suggest that the GC is also applicable to RT-qPCR data, with GiniGenes having good
298  potential (as novel "housekeeping" genes) for the normalisation of such data.

299

# Discussion

Reference genes are commonly used to normalise gene expression data, so as to account for bias resulting from both biological and technical variability, and to enable quantification of gene expression changes or differences in the system under study. It is generally considered that such reference genes should come from pathways that are required for general metabolism, using only one gene per 'pathway' to avoid co-regulation which might make the gene expressions look very stable.

Such reference genes are commonly referred to as 'housekeeping' genes (HKGs) because they are considered to participate in essential cellular functions, are ubiquitously expressed in all cells and tissue types, and their expression is considered to be stable [46-63]). A number of such genes have been proposed over the years, and genes such as GAPDH, ACTB, RPL13A, SDHA, B2M are frequently used in such studies [63]. However, the expression levels of these and other proposed HKGs have in fact been shown to vary widely between cells and tissues (e.g. [50, 59, 62, 64-76]) and their expression has also been reported to be affected by a number of factors relating to the experiment such as cell confluence [132], pathological, experimental and tissue specific conditions [133]. As highlighted by Huggett *et al.* [134], despite the reports of the potential variability of expression of 'classic' references genes such as GAPDH and ACTB, these are still used without mention of any validation processes. Our GiniGenes are selected as reference genes through different, data-driven, criteria.

Various tools have been developed to evaluate and screen reference genes from experimental datasets; these include geNorm [63], NormFinder [135], Best Keeper [136] and the comparative ΔCT finder [49]. RefFinder (http://leonxie.esy.es/RefFinder/#) and RefGenes can integrate these to enable a comparison and ranking of any tested candidate reference genes [137].

These tools assess expression stability of genes in different ways:

- geNorm determines gene stability through a stepwise exclusion or ranking process followed by averaging the geometric mean of the most stable genes from a chosen set. Python implementation: https://eleven.readthedocs.io/en/latest/
- BestKeeper also uses the geometric mean but using raw data rather than copy numbers. BestKeeper [136] can be used as an Excel-based tool. It can accommodate up to 10 housekeeping genes in up to 100 biological samples. Optimal HKGs are determined by pairwise correlation analysis of all pairs of candidate genes, and the geometric mean of the top ranking ones. http://www.gene-quantification.info
- NormFinder measures variation, and ranks potential reference genes between study groups. NormFinder [135] has an add-in for Microsoft Excel and is available as an R programme. It recommends analysis of 5-10 candidate genes and at least 8 samples per group. https://moma.dk/normfinder-software
- The comparative ΔCT finder requires no specialist programmes since this involves comparison of comparisons of ΔCTs between pairs of genes to find a set of genes that show least variability.
- RefGenes allows one to find genes that are stably expressed across tissue types and experimental conditions based on microarray data, and a comparison of results from geNorm, NormFinder and Best Keeper to find a set of reference genes. However, this is not a free service unless one searches for one gene at a time. Furthermore, the site for this tool is no longer available. Moreover, all these tools require the user to make a prior selection of such HKGs (introducing bias and potential errors) and most are cumbersome to understand and calculate.

343   We have here shown how via a simple calculation, the GC, we can find potential reference genes, and
344   illustrated its utility in large-scale cell-line, tissue RNA-Seq data sets and RT-qPCR data. The expression of a
345   number of classical HKGs from a number of carefully selected publications do in fact vary much more
346   substantially between large RNA-Seq data sets, both for tissues and cell lines.

347   Whilst not all studies will involve large data sets such as those we have analysed here, the GC should also
348   be of use for smaller-scale studies to select a subset of genes in a panel of cell lines or tissues relevant to
349   the study in question.

350   Overall we find that (i) two of these genes, HNRNPK and PCBP1, seemed to be particularly robustly and
351   stably expressed at reasonable levels in all cell lines studied, and (ii) a data-driven strategy based on the GC
352   represents a useful and convenient method for normalisation in gene expression profiling and related
353   studies.

354

355

## Methods

The datasets used are described and referenced below. The data, in transcripts per million (TPM) units were downloaded from the EBI expression atlas as a .tsv file. As previously [1], the Gini Index was calculated using the **ineq** package (Achim Zeileis (2014). ineq: Measuring Inequality, Concentration, and Poverty. R package version 0.2-13. https://CRAN.R-project.org/package=ineq) in **R** (https://www.R-project.org/). These calculations were incorporated into KNIME via KNIME's R integration *R Snippet* node. A spreadsheet giving the extracted analyses is provided as supplementary tables (Tables S7 and S8).

Table 1. Studies used for assessing proposed stable reference genes.

| Study short name | Comments | Reference |
|---|---|---|
| GiniGene | Study presenting novel potential housekeeping genes in cells and tissues from the HPA project cell and tissue RNASeq data. | [1] |
| geNorm or Vandesompele | Classic set of reference genes **in tissues** and a means of analysing them | [63] |
| Eisenberg | Very detailed analysis of housekeeping/ reference genes **in tissues using the Illumina Body Map study of** RNA-seq of 16 Human Tissues. E-MTAB-513. | [46] |
| Lee | Two novel reference genes from a detailed analysis of 281 **normal tissue samples from 17 different organs then compares between disease states m and cell lines.** | [131] |
| Caracausi | 646 expression profile data sets from 54 different human **tissues.** | [62] |

Table 2. Studies used for expression profiling data.

| Dataset short name | Comments | Reference |
|---|---|---|
| HPA | RNA-seq-based dataset from the Human Protein Atlas group. Two data sets available: one of 19,628 protein coding genes in 56 **cell lines** (HPA_C) and another of 19,613 protein coding genes in 59 **tissues** (HPA_T). | [90, 91, 138] |
| CCLE | RNA-seq-based dataset (Cancer Cell Line Encyclopedia) of 58,035 genes in 934 human cancer **cell lines** (downloaded from EBI Expression Atlas E-MTAB2770). | [139] |
| Klijn / Genentech | RNA-seq-based analysis of 57,711 genes in 622 human cancer **cell lines** (downloaded from EBI Expression Atlas E-MTAB-2706). | [140] |
| GTEx | RNA-Seq data of 46,711 genes in 53 human **tissue** samples from the | [141] |

12

| | Genotype-Tissue Expression (GTEx) project (downloaded from EBI Expression Atlas E-MTAB-5214). | |
|---|---|---|
| PCAWG | RNA-Seq of 46,816 genes in 76 **tissues**, cancer and normal, from The International Cancer Genome Project: Pan Cancer Analysis of Whole Genomes ((downloaded from EBI Expression Atlas E-MTAB-5200). | Unpublished, may be subject to publication embargo until July 25[th], 2019 https://dcc.icgc.org/pcawg |
| HBM | Illumina Body Map: RNA-seq of 16 Human **Tissues**. E-MTAB-513. Used by Eisenberg and colleagues in their analysis of housekeeping/ reference genes **in tissues**. | [46] |

367

## Cell lines and culture conditions

369 A panel of 10 cell lines were grown in appropriate growth media: K562, PNT2 and T24 in RPMI-1640 (Sigma,
370 Cat No. R7509), Panc1 and HEK293 in DMEM (Sigma, Cat No. D1145), SH-SY5Y in 1:1 mixture of DMEM/F12
371 (Gibco, Cat No. 21041025), J82 and RT-112 in EMEM (Gibco, Cat No. 51200-038), 5637 in Hyclone McCoy's
372 (GE Healthcare, Cat No. SH30270.01) and PC3 in Ham's F12 (Biowest, Cat No. L0135-500). All growth media
373 were supplemented with 10 % fetal bovine serum (Sigma, Cat No. f4135) and 2 mM glutamine (Sigma, Cat
374 No. G7513) without antibiotics. Cell cultures were maintained in T225 culture flasks (Star lab, CytoOne Cat
375 No. CC7682-4225) kept in a 5% $CO_2$ incubator at 37$^{\circ}$C until 70-80 % confluent.

## Harvesting Cells for RNA Extraction

377 Cells from adherent cell lines were harvested by removing growth media and washing twice with 5 mL of
378 pre-warmed phosphate buffered saline (PBS) (Sigma, Cat No. D8537), then incubated in 3 mL of 0.025%
379 trypsin-EDTA solution (Sigma Cat No. T4049) for 2-5 min at 37 $^{\circ}$C. At the end of incubation cells were
380 resuspended in 5-7 mL of respective media when cells appeared detached to dilute trypsin treatment. The
381 cell suspension was transferred to 15 mL centrifuge tubes and immediately centrifuged at 300 x g for 5 min.
382 Suspended cell lines were centrifuged directly from cultures in 50 mL centrifuge tubes and washed with PBS
383 as above. The cell pellets were resuspended in 10-15 mL media and cell count and viability was determined
384 using a Nexcellom Cellometer Auto 1000 Cell Viability Counter (Nexcellom Bioscience) set for Trypan Blue
385 membrane exclusion method. Cells with >95 % viability were used for downstream total RNA extraction.

## RNA Extraction

387 Total RNA was extracted from 2-5 X 10$^6$ cells using the Qiagen RNeasy Mini Kit (Cat No. 74104) and DNAse
388 treated using Turbo DNA-free kit (Invitrogen, Cat No. AM1907) according to the manufacturer's
389 instructions. Briefly, 1 X DNA buffer was added to the extracted RNA prior to adding 2U (1 μL) of DNAse
390 enzyme. The reaction mixture was incubated at 37$^{\circ}$C for 30 min and inactivated for 2 min at room
391 temperature using DNAse inactivating reagent. The mixture was centrifuged at 10,000 x g for 1.5 min and
392 the RNA from the supernatant was transferred to a clean tube. The RNA concentration was determined
393 using a NanoDrop® ND-1000 spectrophotometer and further validated using an Agilent 2100 bio-analyser
394 coupled with 2100 Expert software system. Only RNA samples with an RIN (RNA Integrity Number) between
395 9-10 were selected for cDNA synthesis.

396

397

13

### Reverse Transcription and cDNA Synthesis

1 µg of RNA was reverse transcribed into cDNA. Briefly, a 20 µL reaction was setup by adding 1 µL each of oligodT (50 µM, Invitrogen, cat No. 18418020) and dNTP mix (10 mM, Invitrogen, Cat No. 18427-013) followed by adding an appropriate volume for 1 µg of RNA. Nuclease free water (Ambion, Cat No. AM9937) was then added to make the volume up to 13 µL and incubated at 65°C for 5 min then cooled on ice for 1min. To initiate transcription 4 µL of 5 X first strand buffer (Invitrogen, Cat No. 1889832) and 1 µL each of 0.1 M DTT (Invitrogen, Cat No. 1907572), RNaseOUT™ (Invitrogen, Recombinant RNase Inhibitor, Cat No. 1905432) and SuperScript™ III RT (200 units/µL, Invitrogen, Cat No. 1685475) reverse transcriptase enzyme were added, mixed gently then incubated at 50°C for 60 min followed by inactivation at 70°C for 15 min. The cDNA was diluted 1:100 to be used in RT-qPCR experiment.

### Validation of gene expression by geNorm

A set of candidate reference genes (40; top 32 genes from genes ordered by GC and expression value from [91], plus 8 of the most commonly used from the literature including seven from [63]). RNAseq data were selected for validation of stable gene expression using geNorm [63]. First, a typical qPCR protocol was prepared from a master mix for each gene to be tested per cell line in triplicate. This consisted of 10 µL/well made by adding 0.8 µL of nuclease free water (Ambion), 5 µL of LC480 SYBR Green I Master (2 X conc. Roche, Product No. 04887352001), 0.1 µL each of forward and reverse primers (20 µM) (for primer and amplicon sequences see Supplementary Table S9) and 4 µL of 1:100 diluted cDNA in a 384 well qPCR plate (Starlab Cat. No. E1042-9909-C). The no template controls (NTC) for each gene were produced by replacing cDNA with 4 µL of nuclease free water. Thermal cycling conditions used were: one cycle of 95°C for 10 min followed by 40 cycles of 95°C for 10 sec and 60°C for 30 sec. qPCR was performed using Roche LightCycler LC480 qPCR platform. The fluorescence signals were measured in real time during amplification cycle (Cq) and also during temperature transition for melt curve analysis.

The mean Cq values were converted into relative values for a gene across all cell lines using ΔCq method [142]. Briefly, the lowest Cq value in a panel of cell lines for a gene was subtracted from all the values in that panel using the equation: $R = 2^{(C_{qsample} - C_{qcontrol})}$, where $C_{qsample}$ is the mean Cq value obtained for a gene in each of the cell lines and $C_{qcontrol}$ is the lowest Cq value in that panel. The relative values for each gene in a panel were then obtained by applying $R = 2^{-\Delta C_q}$. These relative values were applied in geNorm Visual Basic applet for Microsoft Excel® [63] that determines the most stable reference genes from a set of genes in a given panel of cell lines.

### Validation of gene expression using the Gini coefficient.

To the raw RT-qPCR data a Cq value (which is inversely proportional to expression level) cut-off of 32 was set, above which no expression is observed. The Cq values of genes in cell lines were subsequently converted to a relative expression level (Cq cut off/Cq value of gene). Descriptive statistics of the expression of each gene in individual cell lines were then calculated. As a final step, the median expression value of each gene in individual cell lines was used to calculate descriptive statistics, including the GC, of gene expression across these cell lines. Figure 11 illustrates a KNIME workflow [127-129] for this purpose. The raw data and descriptive statistics extracted are provided in Supplementary Tables S5 and S6 respectively, and the KMNIME analysis workflow in Supplementary File 1.

440

## Declarations

### Ethics approval and consent to participate

Not applicable.

### Consent for publication

The PCAWG data is under embargo until the WGS pan-cancer consortium publishes its marker paper or until July 25, 2019, whichever is earlier. Methodology papers may be published prior to this embargo, with agreement from the full scientific working group. We have been in email contact with jennifer.jennings@oicr.on.ca who asked that we advise the editor to wait until the July 25 embargo lift.

### Availability of data and materials

All data generated or analysed during this study are included in this published article (and its supplementary information files). The original datasets used are referenced throughout and are summarised in Table 2.

### Competing interests

The authors declare that they have no competing interests.

### Funding

This work is supported by BBSRC Project Grant BB/ P009042/1.

### Authors' contributions

D.B.K. highlighted the utility of the GC as shown in [1]. M.W.M. adapted the Gini method and analyses workflows developed by S.O. from [1] and performed most of the analyses that were done using KNIME. P.J.D. contributed in particular to the analysis of the housekeeping genes. F.M. performed the RT-qPCR analyses. All authors contributed to the writing and approval of the manuscript.

### Acknowledgements

All authors thank the BBSRC (grant BB/P009042/1) and the Novo Nordisk Foundation (grant NNF10CC1016517) for financial support. for financial support.

465

## Legends to figures

467

Fig 1. Graphical indication of the means by which we calculate the Gini coefficient.

Fig 2. Gini coefficient and median expression levels of proposed reference genes in the HPA cell-line dataset. **A**. GC versus median expression level of HPA dataset. **B**. Median expression levels of CCLE vs HPA datasets. Line of best linear fit (in log space) shown is y = 0.991 + 0.827 X ($r^2$=0.606). **C**. Median expression levels of CCLE vs Klijn datasets. Line of best linear fit (in log space) shown is y = 0.998 + 0.804 X ($r^2$=0.593). Colour coding: red, GeneGini reference genes; blue Eisenberg & Levenson; yellow Vandesompele; green Lee; lilac both GeneGini and Eisenberg and Levenson.

Fig 3. Gini coefficient of candidate reference genes in CCLE and Klijn/Genentech cell-line datasets. Left panel shows all proposed housekeeping genes considered in this study, with the right panel showing labels

477 of those genes with a GC < 0.25. The line of best fit is y = -0.171 + 0.829x ($r^2$ = 0.909). Colour code as in Fig
478 2.

479 Fig 4. Robustness of the Gini coefficient. **A.** IQR of different genes in Klijn/Genentech vs HPA cell-line
480 dataset. Left panel shows all genes considered in this study, with right panel showing genes with IQR < 2 in
481 both datasets. Line of best linear fit (in log space) shown is y = 0.01 + 1.11 X ($r^2$=0.937) **B.** IQR of different
482 genes in CCLE vs HPA cell-line dataset. Left panel shows all genes considered in this study, with right panel
483 showing genes with IQR < 2 in both datasets. Line of best linear fit (in log space) shown is y = 0.04 + 0.99 X
484 ($r^2$=0.930). **C.** Max:Mean vs Min expression levels in HPA data set. Colour code as in Fig 2.

485 Fig 5. Shared and unique genes in HPA, CCLE and Klijn/Genentech cell-line data sets. **A.** Genes with a GC <
486 0.2 **B.** Housekeeping genes in Table 2 with GC < 0.2.

487 Fig. 6. GC vs Median for 115 genes in **A.** HPA, **B.** CCLE and **C.** Klijn/Genentech cell-linedata sets. Colour
488 coding: Blue, Caracausi; Green,  GeneGini reference genes; Grey, neither. Shape coding: Circle, other;
489 Triangle, SLC coding gene.

490 Fig. 7. Robustness of GC for finding stably expressed genes using shared genes between HPA, CCLE and
491 Klijn/Genentech cell-line data sets with GC < 0.2. Shown are the results for the Klijn/Genentech dataset **A.**
492 IQR vs GC, **B.** Max:Mean vs Min. Colour coding: Blue, Caracausi; Green,  GeneGini reference genes; Grey,
493 neither. Shape coding: Circle, other; Triangle, SLC coding gene.

494 Fig 8. Gini coefficient and median expression levels of proposed reference genes in the HPA tissue dataset.
495 Colour coding: blue, Caracausi; purple, Eisenberg and Levenson; green, GeneGini reference genes; yellow,
496 both GeneGini and Eisenberg and Levenson; orange, Lee; black, Vandesompele.

497 Fig 9. Robustness of the Gini coefficient in the HPA tissue data set. **A.** RSD versus Gini coefficient of
498 candidate reference genes. Line of best linear fit (in log space) shown is y = 2.45 + 1.24 X ($r^2$=0.938) **B.** IQR
499 versus Gini coefficient of candidate reference genes.  Line of best linear fit (in log space) shown is y = 0.87 +
500 0.96 X ($r^2$=0.566). Colour code as in Fig 8.

501 Fig 10. UpSetR [143] plot showing genes with a GC <0.2 that are variously shared and unique across the
502 PCAWG, HBM, GTEX and HPA tissue data sets. The data underpinning this plot can be found it
503 Supplementary Table S4

504 Fig 11. The KNIME workflow described here to calculate descriptive statistics and the gini coefficient from
505 RT-qPCR data. This workflow can be adapted for use with large RNA-Seq Data sets.

506 Fig 12. Gini coefficient and median expression levels of candidate reference genes assessed by RTqPCR. Left
507 panel shows all genes considered in this study, with right panel showing genes with GC < 0.2. Colour coding:
508 green, GeneGini reference genes; red, both GeneGini and Caracausi reference genes; yellow, GeneGini and
509 Eisenberg and Levenson; orange, Lee, yellow; black, Vandesompele; purple, Zhang and Kriegova.

510 Fig 13. Robustness of the Gini coefficient in assessed experimentally by RT-qPCR using a small subset of
511 proposed reference genes. Left panel shows Gini coefficient vs % RSD for all genes considered in this study,
512 with right panel showing the same with genes with a GC < 0.2 and % RSD < 10. Line of best linear fit shown
513 is y = 0.002 + 0.004x (r2=0.988). Shape coding as in Fig 12.

514 Supplementary Fig S1. Comparison of median expression levels of proposed reference genes between
515 tissue datasets. A. HBM vs HPA tissue datasets. Line of best linear fit (in log space) shown is $\log_{10}$y = 0.35 +

516 (0.74 $\log_{10}$ (x)) (r2=0.472). B. PCAWG vs HPA tissue dataset. Line of best linear fit (in log space) shown is
517 $\log_{10}y$ = 0.46 + (0.73 $\log_{10}$ (x)) (r2=0.500). C. GTEx vs HPA Tissue. Line of best linear fit (in log space) shown is
518 $\log_{10}y$ = 0.45 + (0.68 $\log_{10}$(x)) (r2=0.429). Colour coding: blue, Caracausi reference genes; purple, Eisenberg
519 & Levenson; green, GeneGini; yellow, both GeneGini and Eisenberg and Levenson; orange, Lee; black,
520 Vandesompele.

521 Supplementary Fig S2. Comparison of Gini coefficient of proposed reference genes between tissue
522 datasets. A. HBM vs HPA tissue datasets. Line of best linear fit (in log space) shown is log10y = -0.20 + (0.62
523 $\log_{10}$(x)) (r2=0.392). B. PCAWG vs HPA tissue dataset. Line of best linear fit (in log space) shown is $\log_{10}y$ = -
524 0.15 + (0.59 $\log_{10}$(x)) (r2=0.560). C. GTEx vs HPA Tissue. Line of best linear fit (in log space) shown is $\log_{10}y$ =
525 0.22 + (0.59 $\log_{10}$(x)) (r2=0.388). Colour coding as in Fig S1.

526 Fig S3. Robustness of the Gini coefficient assessed experimentally by RT-qPCR using a small subset of
527 proposed reference genes illustrated with Gini coefficient vs IQR. Left panel shows all 40 genes in Table 6,
528 with right panel showing genes with a GC < 0.2. Colour coding: green, GeneGini reference genes; red, both
529 GeneGini and Caracausi reference genes; yellow, GeneGini and Eisenberg and Levenson; orange, Lee,
530 yellow; black, Vandesompele; purple, Zhang and Kriegova.

531

## List of tables

533

534 Table 1. Studies used for assessing proposed stable reference genes.

535 Table 2. Studies used for expression profiling data.

536 Table 3. Descriptive statistics of 13 genes common across cell-line data sets with GC < 0.2. In addition, the
537 protein name, as well as UniProt ID and function are shown. S/A/O refers to SLC, ABC or Other respectively.

538 Table 4. Descriptive statistics of 15 common genes across tissue data sets with a GC < 0.2. In addition, the
539 protein name, as well as UniProt ID and function are shown.

540 Table 5. Details of human cell lines used for the assessment of expression of candidate reference genes by
541 RT-qPCR.

542 Table 6. Candidate reference genes used to assess expression stability experimentally by RT-qPCR. Included
543 are gene name and UniProt ID, Gini coefficient as calculated using the HPA cell-line data set. S/A/O refers to
544 SLC, ABC or Other respectively.

545 Supplementary Table S1. Descriptive statistics of 115 common genes across cell-line datasets. S/A/O refers
546 to SLC, ABC or Other respectively.

547 Supplementary Table S2. Descriptive statistics and UniProt names and IDs of proposed stable reference
548 genes from Table 1 in tissue datasets. S/A/O refers to SLC, ABC or Other respectively.

549 Supplementary Table S3. Descriptive statistics of common and unique genes across tissue data sets with a
550 GC ≤ 0.2, including gene names and functions. S/A/O refers to SLC, ABC or Other respectively.

551 Supplementary Table S4. Data underpinning UpSetR [143] plot in Figure 10 showing genes with a GC <0.2
552 that are variously shared and unique across the PCAWG, HBM, GTEX and HPA tissue data sets.

17

553 Supplementary Table S5. Raw expression data for candidate reference genes in human cell lines by RT-
554 qPCR.

555 Supplementary Table S6. Descriptive statistics and Gini coefficient data for candidate reference genes in
556 human cell lines by RT-qPCR.

557 Supplementary Table S7. Extracted analyses of cell-line RNA-Seq data sets referenced in Table 2. S/A/O
558 refers to SLC, ABC or Other respectively.

559 Supplementary Table S8. Extracted analyses of tissue RNA-Seq data sets referenced in Table 2. S/A/O
560 refers to SLC, ABC or Other respectively.

561 Supplementary Table S9. Primer and amplicon sequences of candidate reference genes used to assess
562 expression stability experimentally by RT-qPCR. Included are the Gini coefficient and median expression
563 level as found in the HPA cell-line data set. S/A/O refers to SLC, ABC or Other respectively.

## Supplementary Files
564
565

566 Supplementary File 1. KNIME workflow [127-129] that we have written to calculate descriptive statistics,
567 including the GC, of gene expression across cell lines to assess of expression stability of candidate reference
568 genes by RT-qPCR.

569

# References

1. O'Hagan S, Wright Muelas M, Day PJ, Lundberg E, Kell DB: **GeneGini: assessment via the Gini coefficient of reference "housekeeping" genes and diverse human transporter expression profiles** *Cell Syst* 2018, **6**:230-244.

2. Gini C: **Concentration and dependency ratios (in Italian). English translation in: Rivista di Politica Economica, 87 (1997), 769-789.** 1909.

3. Gini C: *Variabilità e Mutabilità. Contributo allo Studio delle Distribuzioni e delle Relazioni Statistiche.* Bologna: C. Cuppini; 1912.

4. Ceriani L, Verme P: **The origins of the Gini index: extracts from Variabilità e Mutabilità (1912) by Corrado Gini.** *J Econ Inequal* 2012, **10**:421-443.

5. Jiang L, Tsoucas D, Yuan GC: **Assessing Inequality in Transcriptomic Data.** *Cell Syst* 2018, **6**:149-150.

6. Wagner GP, Kin K, Lynch VJ: **Measurement of mRNA abundance using RNA-seq data: RPKM measure is inconsistent among samples.** *Theory Biosci* 2012, **131**:281-285.

7. Wilkinson R, Pickett K: *The spirit level: why equality is better for everyone.* London: Penguin Books; 2009.

8. Kondo N, van Dam RM, Sembajwe G, Subramanian SV, Kawachi I, Yamagata Z: **Income inequality and health: the role of population size, inequality threshold, period effects and lag effects.** *J Epidemiol Community Health* 2012, **66**:e11.

9. Pickett KE, Wilkinson RG: **Income inequality and health: a causal review.** *Soc Sci Med* 2015, **128**:316-326.

10. Darkwah KA, Nortey EN, Lotsi A: **Estimation of the Gini coefficient for the lognormal distribution of income using the Lorenz curve.** *Springerplus* 2016, **5**:1196.

11. Kohler TA, Smith ME, Bogaard A, Feinman GM, Peterson CE, Betzenhauser A, Pailes M, Stone EC, Marie Prentiss A, Dennehy TJ, et al: **Greater post-Neolithic wealth disparities in Eurasia than in North America and Mesoamerica.** *Nature* 2017, **551**:619-622.

12. Nishi A, Shirado H, Rand DG, Christakis NA: **Inequality and visibility of wealth in experimental social networks.** *Nature* 2015, **526**:426-429.

13. Damgaard C, Weiner J: **Describing inequality in plant size or fecundity.** *Ecology* 2000, **81**:1139-1142.

14. Sadras V, Bongiovanni R: **Use of Lorenz curves and Gini coefficients to assess yield inequality within paddocks.** *Field Crops Res* 2004, **90**:303-310.

15. Weidlich IE, Filippov IV: **Using the gini coefficient to measure the chemical diversity of small-molecule libraries.** *J Comput Chem* 2016, **37**:2091-2097.

16. Wren JD: **Bioinformatics programs are 31-fold over-represented among the highest impact scientific papers of the past two decades.** *Bioinformatics* 2016, **32**:2686-2691.

17. Ainali C, Valeyev N, Perera G, Williams A, Gudjonsson JE, Ouzounis CA, Nestle FO, Tsoka S: **Transcriptome classification reveals molecular subtypes in psoriasis.** *BMC Genomics* 2012, **13**:472.

18. Tran QN: **Improving the Accuracy of Gene Expression Profile Classification with Lorenz Curves and Gini Ratios.** *Software Tools and Algorithms for Biological Systems* 2011, **696**:83-90.

19. Jiang L, Chen H, Pinello L, Yuan GC: **GiniClust: detecting rare cell types from single-cell gene expression data with Gini index.** *Genome Biol* 2016, **17**:144.

20. Torre E, Dueck H, Shaffer S, Gospocic J, Gupte R, Bonasio R, Kim J, Murray J, Raj A: **A comparison between single cell RNA sequencing and single molecule RNA FISH for rare cell analysis.** *bioRxiv* 2017:138289.

21. Shaffer SM, Dunagin MC, Torborg SR, Torre EA, Emert B, Krepler C, Beqiri M, Sproesser K, Brafford PA, Xiao M, et al: **Rare cell variability and drug-induced reprogramming as a mode of cancer drug resistance.** *Nature* 2017, **546**:431-435.

22. Torre E, Dueck H, Shaffer S, Gospocic J, Gupte R, Bonasio R, Kim J, Murray J, Raj A: **Rare Cell Detection by Single-Cell RNA Sequencing as Guided by Single-Molecule RNA FISH.** *Cell Syst* 2018, **6**:171-179 e175.

| 621 | 23. | Schena M, Shalon D, Heller R, Chai A, Brown PO, Davis RW: **Parallel human genome analysis -** |
| 622 | | **microarray-based expression monitoring of 1000 genes.** *Proc Natl Acad Sci* 1996, **93**:10614-10619. |
| 623 | 24. | Spellman PT, Sherlock G, Zhang MQ, Iyer VR, Anders K, Eisen MB, Brown PO, Botstein D, Futcher B: |
| 624 | | **Comprehensive identification of cell cycle-regulated genes of the yeast *Saccharomyces cerevisiae*** |
| 625 | | **by microarray hybridization.** *Mol Biol Cell* 1998, **9**:3273-3297. |
| 626 | 25. | Schena M, Heller RA, Theriault TP, Konrad K, Lachenmeier E, Davis RW: **Microarrays:** |
| 627 | | **biotechnology's discovery platform for functional genomics.** *Trends Biotechnol* 1998, **16**:301-306. |
| 628 | 26. | Hoyle DC, Rattray M, Jupp R, Brass A: **Making sense of microarray data distributions.** |
| 629 | | *Bioinformatics* 2002, **18**:576-584. |
| 630 | 27. | Quackenbush J: **Microarray data normalization and transformation.** *Nat Genet* 2002, **32** |
| 631 | | **Suppl**:496-501. |
| 632 | 28. | Knight CG, Platt M, Rowe W, Wedge DC, Khan F, Day P, McShea A, Knowles J, Kell DB: **Array-based** |
| 633 | | **evolution of DNA aptamers allows modelling of an explicit sequence-fitness landscape.** *Nucleic* |
| 634 | | *Acids Res* 2009, **37**:e6. |
| 635 | 29. | Walsh CJ, Hu P, Batt J, Santos CC: **Microarray Meta-Analysis and Cross-Platform Normalization:** |
| 636 | | **Integrative Genomics for Robust Biomarker Discovery.** *Microarrays (Basel)* 2015, **4**:389-406. |
| 637 | 30. | Do JH, Choi DK: **Normalization of microarray data: single-labeled and dual-labeled arrays.** *Mol* |
| 638 | | *Cells* 2006, **22**:254-261. |
| 639 | 31. | Steinhoff C, Vingron M: **Normalization and quantification of differential expression in gene** |
| 640 | | **expression microarrays.** *Brief Bioinform* 2006, **7**:166-177. |
| 641 | 32. | Dabney AR, Storey JD: **A new approach to intensity-dependent normalization of two-channel** |
| 642 | | **microarrays.** *Biostatistics* 2007, **8**:128-139. |
| 643 | 33. | Kreil DP, Russell RR: **There is no silver bullet--a guide to low-level data transforms and** |
| 644 | | **normalisation methods for microarray data.** *Brief Bioinform* 2005, **6**:86-97. |
| 645 | 34. | Rahman M, Jackson LK, Johnson WE, Li DY, Bild AH, Piccolo SR: **Alternative preprocessing of RNA-** |
| 646 | | **Sequencing data in The Cancer Genome Atlas leads to improved analysis results.** *Bioinformatics* |
| 647 | | 2015, **31**:3666-3672. |
| 648 | 35. | Lin Y, Golovnina K, Chen ZX, Lee HN, Negron YL, Sultana H, Oliver B, Harbison ST: **Comparison of** |
| 649 | | **normalization and differential expression analyses using RNA-Seq data from 726 individual** |
| 650 | | ***Drosophila melanogaster*.** *BMC Genomics* 2016, **17**:28. |
| 651 | 36. | Li X, Brock GN, Rouchka EC, Cooper NGF, Wu D, O'Toole TE, Gill RS, Eteleeb AM, O'Brien L, Rai SN: **A** |
| 652 | | **comparison of per sample global scaling and per gene normalization methods for differential** |
| 653 | | **expression analysis of RNA-seq data.** *PLoS One* 2017, **12**:e0176185. |
| 654 | 37. | Dunn WB, Broadhurst D, Begley P, Zelena E, Francis-McIntyre S, Anderson N, Brown N, Knowles J, |
| 655 | | Halsall A, Haselden JN, et al: **Procedures for large-scale metabolic profiling of serum and plasma** |
| 656 | | **using gas chromatography and liquid chromatography coupled to mass spectrometry.** *Nat Protoc* |
| 657 | | 2011, **6**:1060-1083. |
| 658 | 38. | Zelena E, Dunn WB, Broadhurst D, Francis-McIntyre S, Carroll KM, Begley P, O'Hagan S, Knowles JD, |
| 659 | | Halsall A, HUSERMET Consortium, et al: **Development of a robust and repeatable UPLC-MS** |
| 660 | | **method for the long-term metabolomic study of human serum.** *Anal Chem* 2009, **81**:1357-1364. |
| 661 | 39. | Heckmann LH, Sørensen PB, Krogh PH, Sørensen JG: **NORMA-Gene: a simple and robust method** |
| 662 | | **for qPCR normalization based on target gene data.** *BMC Bioinformatics* 2011, **12**:250. |
| 663 | 40. | Hruz T, Wyss M, Docquier M, Pfaffl MW, Masanetz S, Borghi L, Verbrugghe P, Kalaydjieva L, Bleuler |
| 664 | | S, Laule O, et al: **RefGenes: identification of reliable and condition specific reference genes for RT-** |
| 665 | | **qPCR data normalization.** *BMC Genomics* 2011, **12**:156. |
| 666 | 41. | Khanna P, Johnson KL, Maron JL: **Optimal reference genes for RT-qPCR normalization in the** |
| 667 | | **newborn.** *Biotech Histochem* 2017:1-8. |
| 668 | 42. | Ling D, Salvaterra PM: **Robust RT-qPCR data normalization: validation and selection of internal** |
| 669 | | **reference genes during post-experimental data analysis.** *PLoS One* 2011, **6**:e17762. |
| 670 | 43. | Sang J, Wang Z, Li M, Cao J, Niu G, Xia L, Zou D, Wang F, Xu X, Han X, et al: **ICG: a wiki-driven** |
| 671 | | **knowledgebase of internal control genes for RT-qPCR normalization.** *Nucleic Acids Res* 2017. |

672 44. Vanhauwaert S, Lefever S, Coucke P, Speleman F, De Paepe A, Vandesompele J, Willaert A: **RT-qPCR**
673         **gene expression analysis in zebrafish: Preanalytical precautions and use of expressed repetitive**
674         **elements for normalization.** *Methods Cell Biol* 2016, **135**:329-342.
675 45. Kell DB, Oliver SG: **Here is the evidence, now what is the hypothesis? The complementary roles of**
676         **inductive and hypothesis-driven science in the post-genomic era.** *Bioessays* 2004, **26**:99-105.
677 46. Eisenberg E, Levanon EY: **Human housekeeping genes, revisited.** *Trends Genet* 2013, **29**:569-574.
678 47. Hoerndli FJ, Toigo M, Schild A, Götz J, Day PJ: **Reference genes identified in SH-SY5Y cells using**
679         **custom-made gene arrays with validation by quantitative polymerase chain reaction.** *Anal*
680         *Biochem* 2004, **335**:30-41.
681 48. Ohl F, Jung M, Xu C, Stephan C, Rabien A, Burkhardt M, Nitsche A, Kristiansen G, Loening SA,
682         Radonic A, Jung K: **Gene expression studies in prostate cancer tissue: which reference gene should**
683         **be selected for normalization?** *J Mol Med (Berl)* 2005, **83**:1014-1024.
684 49. Silver N, Best S, Jiang J, Thein SL: **Selection of housekeeping genes for gene expression studies in**
685         **human reticulocytes using real-time PCR.** *BMC Mol Biol* 2006, **7**:33.
686 50. de Jonge HJM, Fehrmann RSN, de Bont ESJM, Hofstra RMW, Gerbens F, Kamps WA, de Vries EGE,
687         van der Zee AGJ, te Meerman GJ, ter Elst A: **Evidence based selection of housekeeping genes.** *PLoS*
688         *One* 2007, **2**:e898.
689 51. Tatsumi K, Ohashi K, Taminishi S, Okano T, Yoshioka A, Shima M: **Reference gene selection for real-**
690         **time RT-PCR in regenerating mouse livers.** *Biochem Biophys Res Commun* 2008, **374**:106-110.
691 52. Bustin SA, Benes V, Garson JA, Hellemans J, Huggett J, Kubista M, Mueller R, Nolan T, Pfaffl MW,
692         Shipley GL, et al: **The MIQE guidelines: minimum information for publication of quantitative real-**
693         **time PCR experiments.** *Clin Chem* 2009, **55**:611-622.
694 53. Gur-Dedeoglu B, Konu O, Bozkurt B, Ergul G, Seckin S, Yulug IG: **Identification of endogenous**
695         **reference genes for qRT-PCR analysis in normal matched breast tumor tissues.** *Oncol Res* 2009,
696         **17**:353-365.
697 54. Li YL, Ye F, Hu Y, Lu WG, Xie X: **Identification of suitable reference genes for gene expression**
698         **studies of human serous ovarian cancer by real-time polymerase chain reaction.** *Anal Biochem*
699         2009, **394**:110-116.
700 55. Thellin O, ElMoualij B, Heinen E, Zorzi W: **A decade of improvements in quantification of gene**
701         **expression and internal standard selection.** *Biotechnol Adv* 2009, **27**:323-333.
702 56. Chervoneva I, Li Y, Schulz S, Croker S, Wilson C, Waldman SA, Hyslop T: **Selection of optimal**
703         **reference genes for normalization in quantitative RT-PCR.** *BMC Bioinformatics* 2010, **11**:253.
704 57. Wang F, Wang J, Liu D, Su Y: **Normalizing genes for real-time polymerase chain reaction in**
705         **epithelial and nonepithelial cells of mouse small intestine.** *Anal Biochem* 2010, **399**:211-217.
706 58. Zampieri M, Ciccarone F, Guastafierro T, Bacalini MG, Calabrese R, Moreno-Villanueva M, Reale A,
707         Chevanne M, Burkle A, Caiafa P: **Validation of suitable internal control genes for expression**
708         **studies in aging.** *Mech Ageing Dev* 2010, **131**:89-95.
709 59. Casadei R, Pelleri MC, Vitale L, Facchin F, Lenzi L, Canaider S, Strippoli P, Frabetti F: **Identification of**
710         **housekeeping genes suitable for gene expression analysis in the zebrafish.** *Gene Expr Patterns*
711         2011, **11**:271-276.
712 60. Jacob F, Guertler R, Naim S, Nixdorf S, Fedier A, Hacker NF, Heinzelmann-Schwarz V: **Careful**
713         **selection of reference genes is required for reliable performance of RT-qPCR in human normal**
714         **and cancer cell lines.** *PLoS One* 2013, **8**:e59180.
715 61. Oturai DB, Sondergaard HB, Bornsen L, Sellebjerg F, Christensen JR: **Identification of Suitable**
716         **Reference Genes for Peripheral Blood Mononuclear Cell Subset Studies in Multiple Sclerosis.**
717         *Scand J Immunol* 2016, **83**:72-80.
718 62. Caracausi M, Piovesan A, Antonaros F, Strippoli P, Vitale L, Pelleri MC: **Systematic identification of**
719         **human housekeeping genes possibly useful as references in gene expression studies.** *Mol Med*
720         *Rep* 2017, **16**:2397-2410.
721 63. Vandesompele J, De Preter K, Pattyn F, Poppe B, Van Roy N, De Paepe A, Speleman F: **Accurate**
722         **normalization of real-time quantitative RT-PCR data by geometric averaging of multiple internal**
723         **control genes.** *Genome Biol* 2002, **3**:RESEARCH0034.

724  64.  Butte AJ, Dzau VJ, Glueck SB: **Further defining housekeeping, or "maintenance," genes Focus on**
725      **"A compendium of gene expression in normal human tissues".** *Physiol Genomics* 2001, **7**:95-96.
726  65.  Hsiao LL, Dangond F, Yoshida T, Hong R, Jensen RV, Misra J, Dillon W, Lee KF, Clark KE, Haverty P, et
727      al: **A compendium of gene expression in normal human tissues.** *Physiol Genomics* 2001, **7**:97-104.
728  66.  Lee PD, Sladek R, Greenwood CM, Hudson TJ: **Control genes and variability: absence of ubiquitous**
729      **reference transcripts in diverse mammalian expression studies.** *Genome Res* 2002, **12**:292-297.
730  67.  Eisenberg E, Levanon EY: **Human housekeeping genes are compact.** *Trends Genet* 2003, **19**:362-
731      365.
732  68.  Dheda K, Huggett JF, Bustin SA, Johnson MA, Rook G, Zumla A: **Validation of housekeeping genes**
733      **for normalizing RNA expression in real-time PCR.** *Biotechniques* 2004, **37**:112-114, 116, 118-119.
734  69.  Barber RD, Harmer DW, Coleman RA, Clark BJ: **GAPDH as a housekeeping gene: analysis of GAPDH**
735      **mRNA expression in a panel of 72 human tissues.** *Physiol Genomics* 2005, **21**:389-395.
736  70.  Rubie C, Kempf K, Hans J, Su T, Tilton B, Georg T, Brittner B, Ludwig B, Schilling M: **Housekeeping**
737      **gene variability in normal and cancerous colorectal, pancreatic, esophageal, gastric and hepatic**
738      **tissues.** *Mol Cell Probes* 2005, **19**:101-109.
739  71.  Szabo A, Perou CM, Karaca M, Perreard L, Palais R, Quackenbush JF, Bernard PS: **Statistical**
740      **modeling for selecting housekeeper genes.** *Genome Biol* 2004, **5**:R59.
741  72.  Mane VP, Heuer MA, Hillyer P, Navarro MB, Rabin RL: **Systematic method for determining an ideal**
742      **housekeeping gene for real-time PCR analysis.** *J Biomol Tech* 2008, **19**:342-347.
743  73.  Teste MA, Duquenne M, François JM, Parrou JL: **Validation of reference genes for quantitative**
744      **expression analysis by real-time RT-PCR in *Saccharomyces cerevisiae*.** *BMC Mol Biol* 2009, **10**:99.
745  74.  Robinson MD, Oshlack A: **A scaling normalization method for differential expression analysis of**
746      **RNA-seq data.** *Genome Biol* 2010, **11**:R25.
747  75.  Kozera B, Rapacz M: **Reference genes in real-time PCR.** *J Appl Genet* 2013, **54**:391-406.
748  76.  De Spiegelaere W, Dern-Wieloch J, Weigel R, Schumacher V, Schorle H, Nettersheim D, Bergmann
749      M, Brehm R, Kliesch S, Vandekerckhove L, Fink C: **Reference gene validation for RT-qPCR, a note on**
750      **different available software packages.** *PLoS One* 2015, **10**:e0122515.
751  77.  Papatheodorou I, Fonseca NA, Keays M, Tang YA, Barrera E, Bazant W, Burke M, Fullgrabe A,
752      Fuentes AM, George N, et al: **Expression Atlas: gene and protein expression across multiple**
753      **studies and organisms.** *Nucleic Acids Res* 2018, **46**:D246-D251.
754  78.  Mortazavi A, Williams BA, McCue K, Schaeffer L, Wold B: **Mapping and quantifying mammalian**
755      **transcriptomes by RNA-Seq.** *Nat Methods* 2008, **5**:621-628.
756  79.  Wang Z, Gerstein M, Snyder M: **RNA-Seq: a revolutionary tool for transcriptomics.** *Nat Rev Genet*
757      2009, **10**:57-63.
758  80.  Oshlack A, Robinson MD, Young MD: **From RNA-seq reads to differential expression results.**
759      *Genome Biol* 2010, **11**:220.
760  81.  Xu J, Gong B, Wu L, Thakkar S, Hong H, Tong W: **Comprehensive Assessments of RNA-seq by the**
761      **SEQC Consortium: FDA-Led Efforts Advance Precision Medicine.** *Pharmaceutics* 2016, **8**.
762  82.  Bray NL, Pimentel H, Melsted P, Pachter L: **Near-optimal probabilistic RNA-seq quantification.** *Nat*
763      *Biotechnol* 2016, **34**:525-527.
764  83.  Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, Adiconis X, Fan L, Raychowdhury
765      R, Zeng Q, et al: **Full-length transcriptome assembly from RNA-Seq data without a reference**
766      **genome.** *Nat Biotechnol* 2011, **29**:644-652.
767  84.  Schulz MH, Zerbino DR, Vingron M, Birney E: **Oases: robust *de novo* RNA-seq assembly across the**
768      **dynamic range of expression levels.** *Bioinformatics* 2012, **28**:1086-1092.
769  85.  Tang F, Barbacioru C, Wang Y, Nordman E, Lee C, Xu N, Wang X, Bodeau J, Tuch BB, Siddiqui A, et al:
770      **mRNA-Seq whole-transcriptome analysis of a single cell.** *Nat Methods* 2009, **6**:377-382.
771  86.  Macosko EZ, Basu A, Satija R, Nemesh J, Shekhar K, Goldman M, Tirosh I, Bialas AR, Kamitaki N,
772      Martersteck EM, et al: **Highly Parallel Genome-wide Expression Profiling of Individual Cells Using**
773      **Nanoliter Droplets.** *Cell* 2015, **161**:1202-1214.
774  87.  Rissin DM, Walt DR: **Digital concentration readout of single enzyme molecules using femtoliter**
775      **arrays and Poisson statistics.** *Nano Lett* 2006, **6**:520-523.

776  88.   Salehi-Reyhani A, Sharma S, Burgin E, Barclay M, Cass A, Neil MAA, Ces O, Willison KR, Klug DR,
777        Brown A, Novakova M: **Scaling advantages and constraints in miniaturized capture assays for**
778        **single cell protein analysis.** *Lab Chip* 2013, **13**:2066-2074.
779  89.   Hudecova I: **Digital PCR analysis of circulating nucleic acids.** *Clin Biochem* 2015, **48**:948-956.
780  90.   Thul PJ, Åkesson L, Wiking M, Mahdessian D, Geladaki A, Ait Blal H, Alm T, Asplund A, Björk L,
781        Breckels LM, et al: **A subcellular map of the human proteome.** *Science* 2017, **356**.
782  91.   O'Hagan S, Wright Muelas M, Day PJ, Lundberg E, Kell DB: **GeneGini: Assessment via the Gini**
783        **Coefficient of Reference "Housekeeping" Genes and Diverse Human Transporter Expression**
784        **Profiles.** *Cell Syst* 2018, **6**:230-244 e231.
785  92.   Wu Y, Zhao W, Liu Y, Tan X, Li X, Zou Q, Xiao Z, Xu H, Wang Y, Yang X: **Function of HNRNPC in breast**
786        **cancer cells by controlling the dsRNA-induced interferon response.** *The EMBO Journal* 2018,
787        **37**:e99017.
788  93.   Bomsztyk K, Denisenko O, Ostrowski J: **hnRNP K: One protein multiple processes.** *BioEssays* 2004,
789        **26**:629-638.
790  94.   Makeyev AV, Liebhaber SA: **The poly (C)-binding proteins: a multiplicity of functions and a search**
791        **for mechanisms.** *Rna* 2002, **8**:265-278.
792  95.   Huo L-R, Zhong N: **Identification of transcripts and translatants targeted by overexpressed PCBP1.**
793        *Biochimica et Biophysica Acta (BBA)-Proteins and Proteomics* 2008, **1784**:1524-1533.
794  96.   Cho S-J, Jung Y-S, Chen X: **Poly (C)-binding protein 1 regulates p63 expression through mRNA**
795        **stability.** *PloS one* 2013, **8**:e71724-e71724.
796  97.   Lardelli RM, Thompson JX, Yates JR, Stevens SW: **Release of SF3 from the intron branchpoint**
797        **activates the first step of pre-mRNA splicing.** *Rna* 2010.
798  98.   Kfir N, Lev-Maor G, Glaich O, Alajem A, Datta A, Sze Siu K, Meshorer E, Ast G: **SF3B1 Association**
799        **with Chromatin Determines Splicing Outcomes.** *Cell Reports* 2015, **11**:618-629.
800  99.   Effenberger KA, Urabe VK, Prichard BE, Ghosh AK, Jurica MS: **Interchangeable SF3B1 inhibitors**
801        **interfere with pre-mRNA splicing at multiple stages.** *RNA* 2016, **22**:350-359.
802  100.  He X, Zhang P: **Serine/arginine-rich splicing factor 3 (SRSF3) regulates homologous recombination-**
803        **mediated DNA repair.** *Molecular Cancer* 2015, **14**:158.
804  101.  Gallardo M, Lee Hun J, Zhang X, Bueso-Ramos C, Pageon Laura R, McArthur M, Multani A, Nazha A,
805        Manshouri T, Parker-Thornburg J, et al: **hnRNP K Is a Haploinsufficient Tumor Suppressor that**
806        **Regulates Proliferation and Differentiation Programs in Hematologic Malignancies.** *Cancer Cell*
807        2015, **28**:486-499.
808  102.  Barboro P, Repaci E, Rubagotti A, Salvi S, Boccardo S, Spina B, Truini M, Introini C, Puppo P, Ferrari
809        N, et al: **Heterogeneous nuclear ribonucleoprotein K: altered pattern of expression associated**
810        **with diagnosis and prognosis of prostate cancer.** *British Journal Of Cancer* 2009, **100**:1608.
811  103.  Park YM, Hwang SJ, Masuda K, Choi K-M, Jeong M-R, Nam D-H, Gorospe M, Kim HH:
812        **Heterogeneous Nuclear Ribonucleoprotein C1/C2 Controls the Metastatic Potential of**
813        **Glioblastoma by Regulating PDCD4.** *Molecular and Cellular Biology* 2012, **32**:4237.
814  104.  Lee EK, Kim HH, Kuwano Y, Abdelmohsen K, Srikantan S, Subaran SS, Gleichmann M, Mughal MR,
815        Martindale JL, Yang X, et al: **hnRNP C promotes APP translation by competing with FMRP for APP**
816        **mRNA recruitment to P bodies.** *Nature structural & molecular biology* 2010, **17**:732-739.
817  105.  Zarnack K, König J, Tajnik M, Martincorena I, Eustermann S, Stévant I, Reyes A, Anders S, Luscombe
818        Nicholas M, Ule J: **Direct Competition between hnRNP C and U2AF65 Protects the Transcriptome**
819        **from the Exonization of Alu Elements.** *Cell* 2013, **152**:453-466.
820  106.  Wang H, Vardy LA, Tan CP, Loo JM, Guo K, Li J, Lim SG, Zhou J, Chng WJ, Ng SB, et al: **PCBP1**
821        **Suppresses the Translation of Metastasis-Associated PRL-3 Phosphatase.** *Cancer Cell* 2010, **18**:52-
822        62.
823  107.  Zhang T, Huang X-H, Dong L, Hu D, Ge C, Zhan Y-Q, Xu W-X, Yu M, Li W, Wang X, et al: **PCBP-1**
824        **regulates alternative splicing of the CD44 gene and inhibits invasion in human hepatoma cell line**
825        **HepG2 cells.** *Molecular Cancer* 2010, **9**:72.
826  108.  Liu Y, Gai L, Liu J, Cui Y, Zhang Y, Feng J: **Expression of poly(C)-binding protein 1 (PCBP1) in NSCLC**
827        **as a negative regulator of EMT and its clinical value.** *International journal of clinical and*
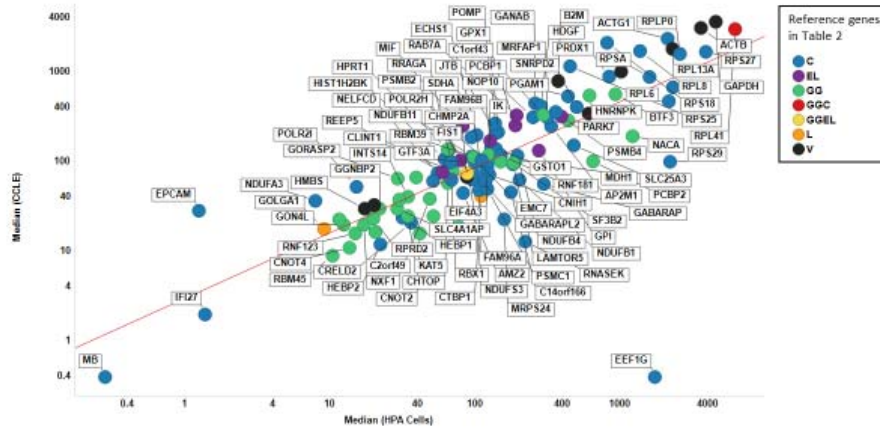828        *experimental pathology* 2015, **8**:7165-7172.

829    109.    Zhang Z-Z, Shen Z-Y, Shen Y-Y, Zhao E-H, Wang M, Wang C-J, Cao H, Xu J: **HOTAIR Long Noncoding**
830            **RNA Promotes Gastric Cancer Metastasis through Suppression of Poly r(C)-Binding Protein (PCBP)**
831            **1.** *Molecular Cancer Therapeutics* 2015, **14**:1162.
832    110.    Wagener R, Aukema SM, Schlesner M, Haake A, Burkhardt B, Claviez A, Drexler HG, Hummel M,
833            Kreuz M, Loeffler M, et al: **The PCBP1 gene encoding poly(rc) binding protein i is recurrently**
834            **mutated in Burkitt lymphoma.** *Genes, Chromosomes and Cancer* 2015, **54**:555-564.
835    111.    Ji F-J, Wu Y-Y, An Z, Liu X-S, Jiang J-N, Chen F-F, Fang X-D: **Expression of both poly r(C) binding**
836            **protein 1 (PCBP1) and miRNA-3978 is suppressed in peritoneal gastric cancer metastasis.**
837            *Scientific reports* 2017, **7**:15488-15488.
838    112.    Jumaa H, Wei G, Nielsen PJ: **Blastocyst formation is blocked in mouse embryos lacking the splicing**
839            **factor SRp20.** *Current Biology* 1999, **9**:899-902.
840    113.    Palmieri F: **The mitochondrial transporter family SLC25: Identification, properties and**
841            **physiopathology.** *Mol Asp Med* 2013, **34**:465-484.
842    114.    Schnabel M, Marlovits S, Eckhoff G, Fichtel I, Gotzen L, Vécsei V, Schlegel J: **Dedifferentiation-**
843            **associated changes in morphology and gene expression in primary human articular chondrocytes**
844            **in cell culture.** *Osteoarthritis and Cartilage* 2002, **10**:62-70.
845    115.    Eisenberg E, Levanon EY: **Human housekeeping genes, revisited.** *Trends in Genetics* 2013, **29**:569-
846            574.
847    116.    Cullen PJ: **Endosomal sorting and signalling: an emerging role for sorting nexins.** *Nature Reviews*
848            *Molecular Cell Biology* 2008, **9**:574.
849    117.    Naslavsky N, Caplan S: **The enigmatic endosome – sorting the ins and outs of endocytic trafficking.**
850            *Journal of Cell Science* 2018, **131**:jcs216499.
851    118.    Chen C, Garcia-Santos D, Ishikawa Y, Seguin A, Li L, Fegan Katherine H, Hildick-Smith Gordon J, Shah
852            Dhvanit I, Cooney Jeffrey D, Chen W, et al: **Snx3 Regulates Recycling of the Transferrin Receptor**
853            **and Iron Assimilation.** *Cell Metabolism* 2013, **17**:343-352.
854    119.    Xu S, Nigam SM, Brodin L: **Overexpression of SNX3 Decreases Amyloid-β Peptide Production by**
855            **Reducing Internalization of Amyloid Precursor Protein.** *Neurodegenerative Diseases* 2018, **18**:26-
856            37.
857    120.    Binder NK, Sheedy JR, Hannan NJ, Gardner DK: **Male obesity is associated with changed**
858            **spermatozoa Cox4i1 mRNA level and altered seminal vesicle fluid composition in a mouse model.**
859            *MHR: Basic science of reproductive medicine* 2015, **21**:424-434.
860    121.    Li Y, Park J-S, Deng J-H, Bai Y: **Cytochrome c oxidase subunit IV is essential for assembly and**
861            **respiratory function of the enzyme complex.** *Journal of Bioenergetics and Biomembranes* 2006,
862            **38**:283-291.
863    122.    Storey JD, Madeoy J, Strout JL, Wurfel M, Ronald J, Akey JM: **Gene-Expression Variation Within and**
864            **Among Human Populations.** *The American Journal of Human Genetics* 2007, **80**:502-509.
865    123.    Lonsdale J, Thomas J, Salvatore M, Phillips R, Lo E, Shad S, Hasz R, Walters G, Garcia F, Young N, et
866            al: **The Genotype-Tissue Expression (GTEx) project.** *Nature Genetics* 2013, **45**:580.
867    124.    Pickrell JK, Marioni JC, Pai AA, Degner JF, Engelhardt BE, Nkadori E, Veyrieras J-B, Stephens M, Gilad
868            Y, Pritchard JK: **Understanding mechanisms underlying human gene expression variation with**
869            **RNA sequencing.** *Nature* 2010, **464**:768.
870    125.    Zhang X, Ding L, Sandford AJ: **Selection of reference genes for gene expression studies in human**
871            **neutrophils by real-time PCR.** *BMC Mol Biol* 2005, **18**:4.
872    126.    Kriegova E, Arakelyan A, Fillerova R, Zatloukal J, Mrazek F, Navratilova Z, Kolek V, du Bois RM,
873            Petrek M.: **PSMB2 and RPL32 are suitable denominators to normalize gene expression profiles in**
874            **bronchoalveolar cells.** *BMC Mol Biol* 2008, **31**:69.
875    127.    Mazanetz MP, Marmon RJ, Reisser CBT, Morao I: **Drug discovery applications for KNIME: an open**
876            **source data mining platform.** *Curr Top Med Chem* 2012, **12**:1965-1979.
877    128.    Fillbrunn A, Dietz C, Pfeuffer J, Rahn R, Landrum GA, Berthold MR: **KNIME for reproducible cross-**
878            **domain analysis of life science data.** *J Biotechnol* 2017.
879    129.    O'Hagan S, Kell DB: **The KNIME workflow environment and its applications in Genetic**
880            **Programming and machine learning.** *Genetic Progr Evol Mach* 2015, **16**:387-391.

24

130. Caracausi M, Piovesan A, Antonaros F, Strippoli P, Vitale L, Pelleri MC: **Systematic identification of human housekeeping genes possibly useful as references in gene expression studies.** *Molecular medicine reports* 2017, **16**:2397-2410.

131. Lee S, Jo M, Lee J, Koh SS, Kim S: **Identification of novel universal housekeeping genes by statistical analysis of microarray data.** *J Biochem Mol Biol* 2007, **40**:226-231.

132. Greer S, Honeywell R, Geletu M, Arulanandam R, Raptis L: **Housekeeping genes; expression levels may change with density of cultured cells.** *Journal of Immunological Methods* 2010, **355**:76-79.

133. Li R, Shen Y: **An old method facing a new challenge: Re-visiting housekeeping proteins as internal reference control for neuroscience research.** *Life Sciences* 2013, **92**:747-751.

134. Huggett J, Dheda K, Bustin S, Zumla A: **Real-time RT-PCR normalisation; strategies and considerations.** *Genes Immun* 2005, **6**:279-284.

135. Andersen CL, Jensen JL, Orntoft TF: **Normalization of real-time quantitative reverse transcription-PCR data: a model-based variance estimation approach to identify genes suited for normalization, applied to bladder and colon cancer data sets.** *Cancer Res* 2004, **64**:5245-5250.

136. Pfaffl MW, Tichopad A, Prgomet C, Neuvians TP: **Determination of stable housekeeping genes, differentially regulated target genes and sample integrity: BestKeeper--Excel-based tool using pair-wise correlations.** *Biotechnol Lett* 2004, **26**:509-515.

137. Xie F, Xiao P, Chen D, Xu L, Zhang B: **miRDeepFinder: a miRNA analysis tool for deep sequencing of plant small RNAs.** *Plant Mol Biol* 2012.

138. Uhlen M, Fagerberg L, Hallstrom BM, Lindskog C, Oksvold P, Mardinoglu A, Sivertsson A, Kampf C, Sjostedt E, Asplund A, et al: **Proteomics. Tissue-based map of the human proteome.** *Science* 2015, **347**:1260419.

139. Barretina J, Caponigro G, Stransky N, Venkatesan K, Margolin AA, Kim S, Wilson CJ, Lehár J, Kryukov GV, Sonkin D, et al: **The Cancer Cell Line Encyclopedia enables predictive modelling of anticancer drug sensitivity.** *Nature* 2012, **483**:603-607.

140. Klijn C, Durinck S, Stawiski EW, Haverty PM, Jiang Z, Liu H, Degenhardt J, Mayba O, Gnad F, Liu J, et al: **A comprehensive transcriptional portrait of human cancer cell lines.** *Nat Biotechnol* 2015, **33**:306-312.

141. Consortium GT: **Human genomics. The Genotype-Tissue Expression (GTEx) pilot analysis: multitissue gene regulation in humans.** *Science* 2015, **348**:648-660.

142. Livak KJ, Schmittgen TD: **Analysis of Relative Gene Expression Data Using Real-Time Quantitative PCR and the 2−ΔΔCT Method.** *Methods* 2001, **25**:402-408.

143. Conway JR, Lex A, Gehlenborg N: **UpSetR: an R package for the visualization of intersecting sets and their properties.** *Bioinformatics* 2017, **33**:2938-2940.
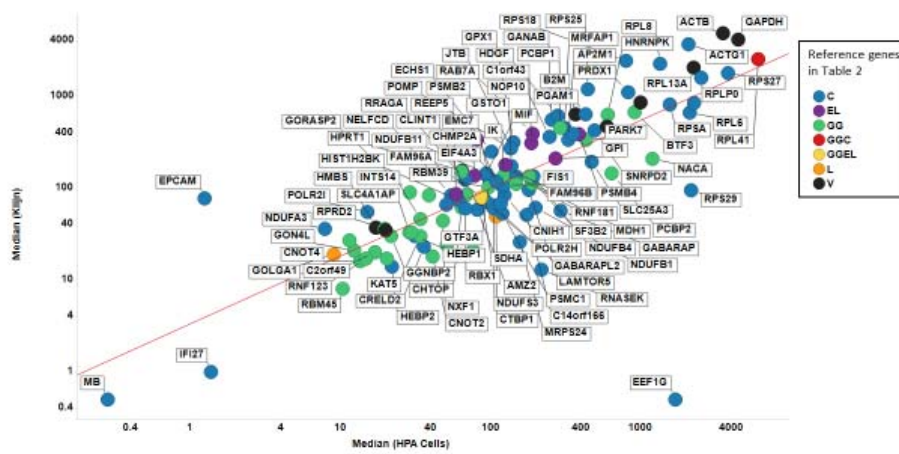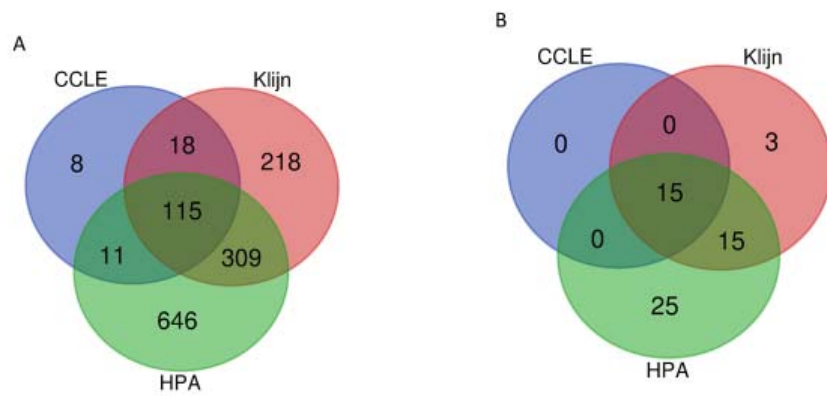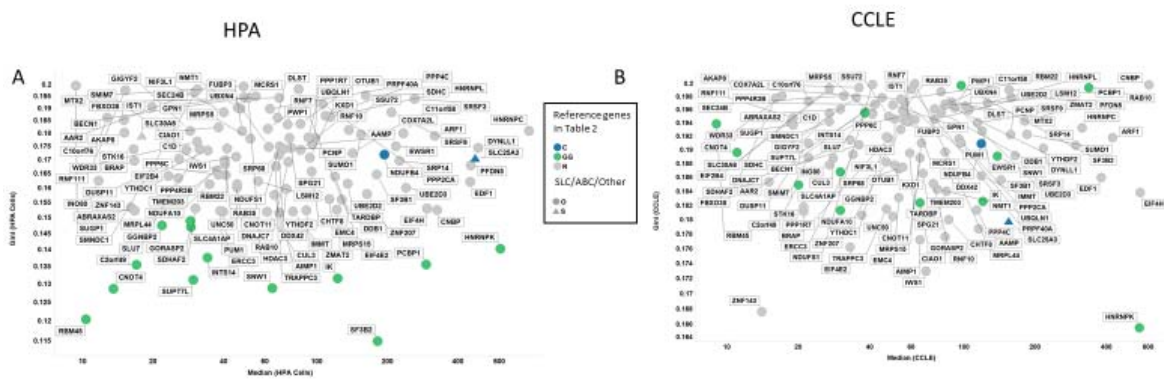
1



916

## 2 A.



917

918



919



920

## 3.



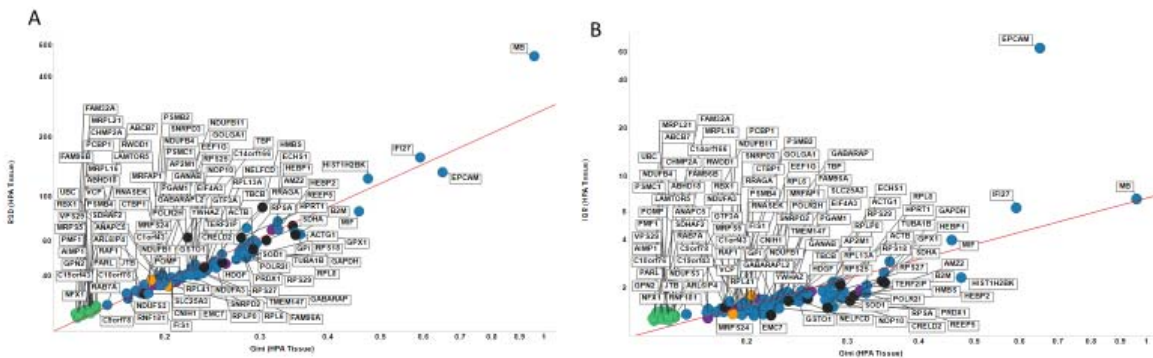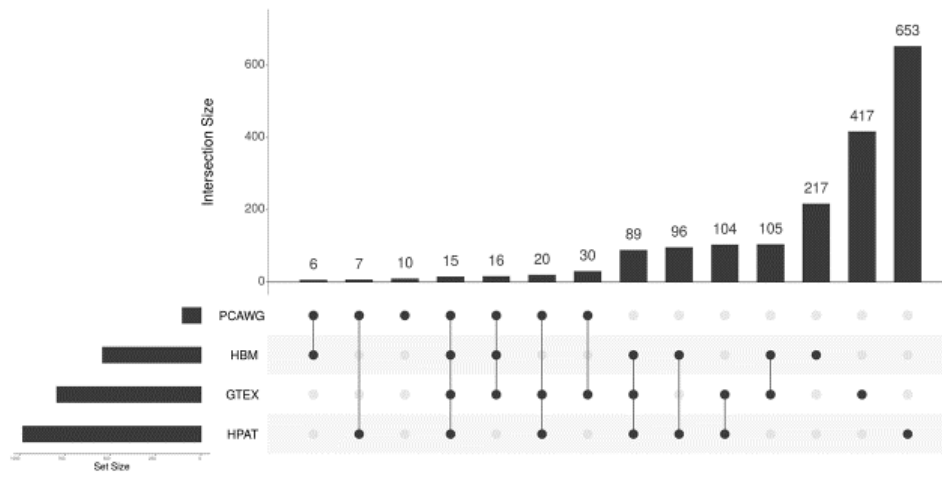921

## 4 A.



922

## 4 B.

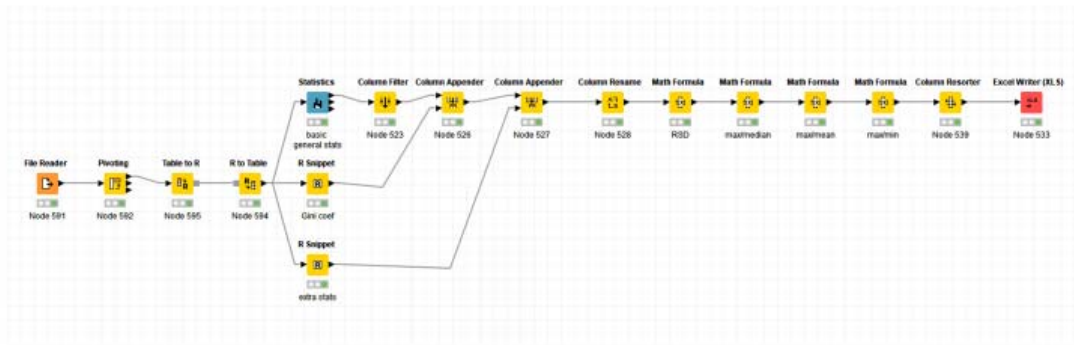

923

## 4 C.



924

5.



925

6



926

6



927

7



928

8



929

9



930

10



931

11



932

12



933

13



934

# 1 A

1 B

1 C

2 A

2 B

2 C

3