



ELSEVIER

SCIENCE @ DIRECT®

Genomics xx (2005) xxx – xxx

GENOMICS

www.elsevier.com/locate/ygeno

## Assessment and integration of publicly available SAGE, cDNA microarray, and oligonucleotide microarray expression data for global coexpression analyses

Obi L. Griffith<sup>a</sup>, Erin D. Pleasance<sup>a</sup>, Debra L. Fulton<sup>b</sup>, Mehrdad Oveisi<sup>a</sup>, Martin Ester<sup>c</sup>,  
Asim S. Siddiqui<sup>a</sup>, Steven J.M. Jones<sup>a,\*</sup>

<sup>a</sup>Genome Sciences Centre, British Columbia Cancer Agency, Vancouver, BC, Canada V5Z 4E6

<sup>b</sup>Molecular Biology and Biochemistry, Simon Fraser University, Burnaby, BC, Canada V5A 1S6

<sup>c</sup>School of Computing Science, Simon Fraser University, Burnaby, BC, Canada V5A 1S6

Received 22 March 2005; accepted 16 June 2005

### Abstract

Large amounts of gene expression data from several different technologies are becoming available to the scientific community. A common practice is to use these data to calculate global gene coexpression for validation or integration of other “omic” data. To assess the utility of publicly available datasets for this purpose we have analyzed *Homo sapiens* data from 1202 cDNA microarray experiments, 242 SAGE libraries, and 667 Affymetrix oligonucleotide microarray experiments. The three datasets compared demonstrate significant but low levels of global concordance ( $r_c < 0.11$ ). Assessment against Gene Ontology (GO) revealed that all three platforms identify more coexpressed gene pairs with common biological processes than expected by chance. As the Pearson correlation for a gene pair increased it was more likely to be confirmed by GO. The Affymetrix dataset performed best individually with gene pairs of correlation 0.9–1.0 confirmed by GO in 74% of cases. However, in all cases, gene pairs confirmed by multiple platforms were more likely to be confirmed by GO. We show that combining results from different expression platforms increases reliability of coexpression. A comparison with other recently published coexpression studies found similar results in terms of performance against GO but with each method producing distinctly different gene pair lists.

© 2005 Elsevier Inc. All rights reserved.

**Keywords:** Gene expression; Gene expression profiling; Microarray analysis; cDNA microarray; Oligonucleotide microarray; Coexpression; Serial analysis of gene expression; Gene Ontology

Large-scale expression profiling has become an important tool for the identification of gene functions and regulatory elements. The development of three such techniques, cDNA microarrays [1], oligonucleotide microarrays [2], and serial analysis of gene expression (SAGE) [3] has resulted in a plethora of studies attempting to

elucidate cellular processes by identifying groups of genes that appear to be coexpressed. Our motivation for this study was to explore the fecundity of large extant expression datasets to identify coexpressed genes and their utility as a resource for biological study. Coexpression data are increasingly used for validation and integration with other “omic” data sources such as sequence conservation [4], yeast two-hybrid interactions [5,6], RNA interference [7], and regulatory element predictions [8], to name only a few. If different platforms or datasets produce widely different measures of coexpression it could have significant impacts on the results of such studies. Furthermore, methods to assess these datasets and identify a coherent, consistent picture of coexpression will be needed.

**Abbreviation:** SAGE, serial analysis of gene expression; GEO, Gene Expression Omnibus; GO, Gene Ontology; IEA, inferred electronic annotation; MGC, Mammalian Gene Collection; MCE, minimum number of common experiments;  $r$ , Pearson correlation;  $r_c$ , correlation of Pearson correlations.

\* Corresponding author. Fax: +1 604 876 3561.

E-mail address: sjones@bcgsc.ca (S.J.M. Jones).

48 High degrees of consistency within a platform have been  
 49 reported for cDNA microarrays and Affymetrix oligonucleo-  
 50 tide microarrays [9–11]. The reproducibility of SAGE has  
 51 not been demonstrated as clearly the time and cost required  
 52 to produce individual SAGE libraries are high. However, a  
 53 recent study showed a high degree of reproducibility and  
 54 accuracy for microSAGE (a modification of SAGE) [12] and  
 55 preliminary analysis of SAGE replicates has demonstrated  
 56 high levels of correlation, similar to those seen for  
 57 Affymetrix platforms (A. Delaney, personal communica-  
 58 tion). Cross-platform comparisons of gene expression values  
 59 have found “reasonable” correlations for matched samples,  
 60 especially for more highly expressed transcripts [11,13–19].  
 61 Other comparisons have reported “poor” correlations  
 62 [15,18,20–24]. The correlations reported above were for  
 63 expression levels or expression changes of individual genes,  
 64 not coexpression of gene pairs. To our knowledge, only one  
 65 study has examined the correlation of coexpression results  
 66 from multiple platforms [25]. The authors compared  
 67 matched Affymetrix oligonucleotide chips and spotted  
 68 cDNA microarrays for the NCI-60 cancer cell panel. For  
 69 each platform, the calculation involved determining the  
 70 Pearson correlation ( $r$ ) between expression profiles (across  
 71 60 cell lines) for all pair-wise gene combinations. Then, a  
 72 correlation of correlations ( $r_c$ ) between the two platforms  
 73 was determined. When all gene pairs were considered a  
 74 global concordance of  $r_c = 0.25$  was reported. As the  
 75 correlation cutoff was increased,  $r_c$  improved steadily to 0.92  
 76 at a correlation cutoff of  $r = 0.91$  (but only 28 of 2061 genes  
 77 remained). Thus, for most gene pairs there is poor correlation  
 78 of correlations for global coexpression values.

79 Genome-wide coexpression analyses in *Caenorhabditis*  
 80 *elegans* and *Saccharomyces cerevisiae* have been used with  
 81 some success to identify gene function or genes that are  
 82 coregulated [26–28]. This “guilt-by-association” approach  
 83 has received criticism because of high levels of noise and  
 84 other problems inherent to the methods [29] but still holds  
 85 great interest for biologists. If matched samples display  
 86 questionable levels of consistency between expression  
 87 profiles generated by different platforms the question  
 88 remains as to how effectively unmatched samples from  
 89 many different sources will compare. If two genes are co-  
 90 regulated (i.e., controlled by an identical set of transcription  
 91 factors) they should display similar expression patterns  
 92 across many conditions and be identified as coexpressed.  
 93 This is the basic premise of many gene function and regu-  
 94 lation studies. If true, large datasets from different expression  
 95 platforms should identify the same coexpressed gene pairs  
 96 even if derived from different conditions and tissues.

97 However, it may be that few genes are globally coregulated  
 98 and thus datasets comprising different samples will identify  
 99 different sets of coregulated genes. Similarly, noise and  
 100 biases inherent to the different methods may result in highly  
 101 discordant measures of coexpression, even for genes with  
 102 similar function or under similar regulatory control.

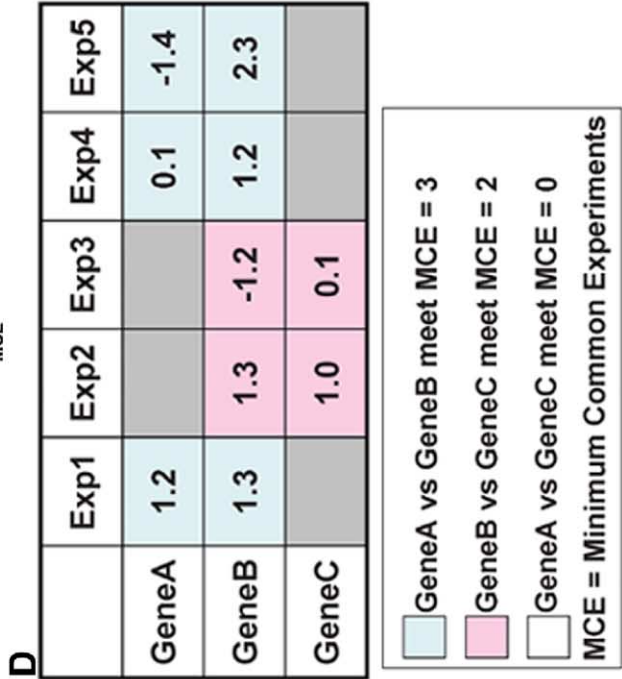
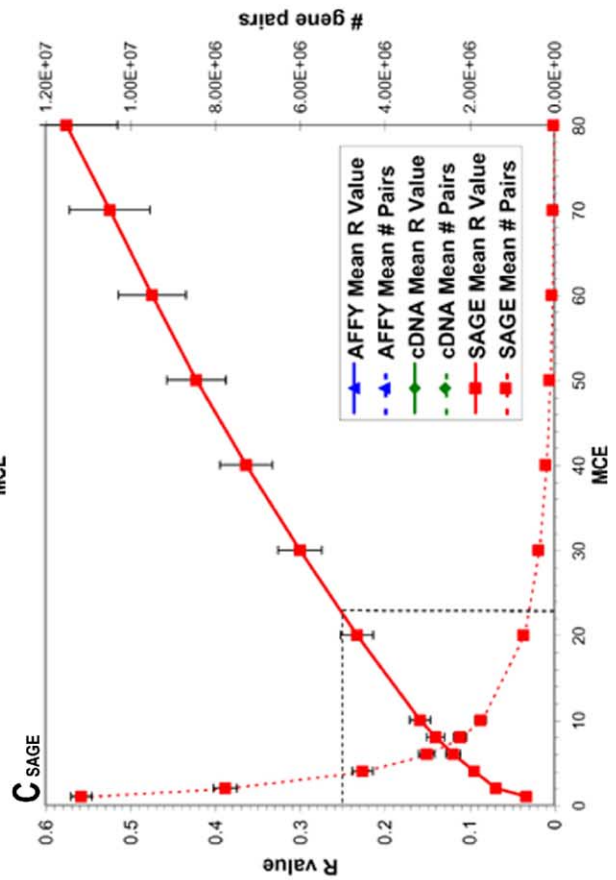
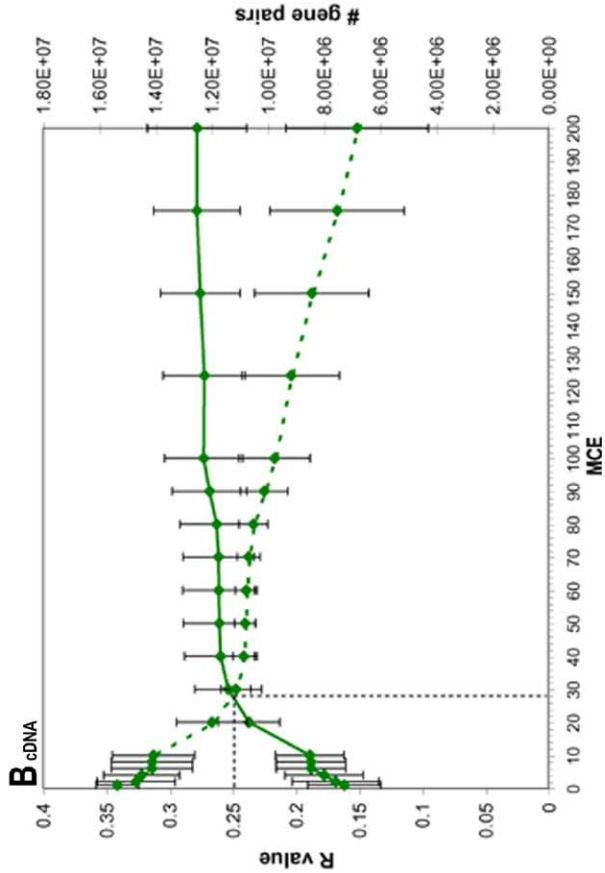
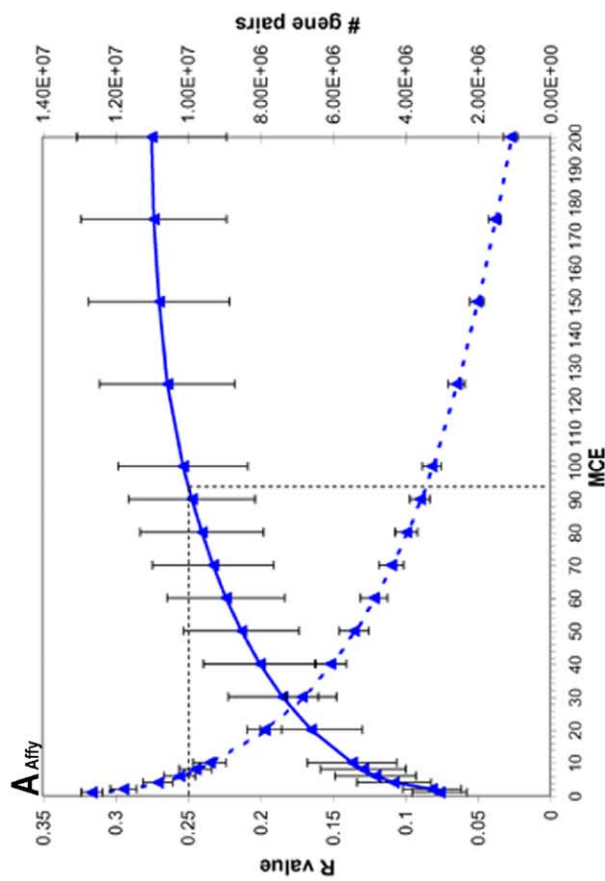
103 The purpose of this study was to assess the differences  
 104 between publicly available expression data for global  
 105 coexpression analyses and investigate the value of combining  
 106 multiple platforms to decrease noise and improve confidence  
 107 in coexpression predictions. We have compared large  
 108 publicly available datasets for SAGE, cDNA microarray  
 109 (cDNA), and Affymetrix oligonucleotide microarray (Affy-  
 110 metrix) platforms (Supplemental Fig. 1). We calculated all  
 111 gene-to-gene Pearson correlation coefficients and assessed  
 112 the platforms for internal consistency, cross-platform con-  
 113 cordance, and agreement with the Gene Ontology. The  
 114 Pearson correlation was chosen as a similarity metric because  
 115 it is one of the most commonly used, with numerous  
 116 published examples for Affymetrix [9,30,31], cDNA  
 117 [5,27,32], and SAGE [33,34]. Because the datasets represent  
 118 unmatched samples, a direct comparison of platforms is  
 119 challenging. Our results indicate that the three platforms  
 120 identify very different measures of coexpression for most  
 121 gene pairs with a very low correlation of correlations between  
 122 platforms. However, coexpression predictions become more  
 123 reproducible with larger datasets and each of the three  
 124 platforms performs better (identifies more gene pairs with  
 125 common GO terms) as the Pearson correlation increases.  
 126 Furthermore, gene pairs confirmed by more than one  
 127 platform (high two-platform average Pearson) were much  
 128 more likely to share a GO term than those identified by only a  
 129 single platform. Other recently published coexpression  
 130 methods (TMM, ArrayProspector) also performed well  
 131 against GO at higher scores but identified very different  
 132 gene pairs. By using the Gene Ontology to choose thresholds  
 133 of high-confidence pairs for each we identify a set of  
 134 coexpressed gene pairs that represents the best of each  
 135 approach.

## 136 Results

### 137 Internal consistency

138 Before performing cross-platform comparisons, it is  
 139 relevant to evaluate each platform individually to determine  
 140 how consistently different experiments from one technology  
 141 identify the same levels of gene coexpression. To this end,

Fig. 1. Internal consistency and minimum common experiments analysis using the pseudo-random division method. For each gene pair, the number of common experiments is determined as the number of experiments for which expression values are available for both genes. On the left axis, MCE is plotted against internal consistency. On the right axis, MCE is plotted against number of gene pairs. In general, as more MCE are required, fewer gene pairs meet the criteria but the internal consistency improves as the correlation is based on more expression data. Notice that the (A) Affymetrix and (B) cDNA datasets appear to level off at approximately  $r_c = 0.3$  with 100 MCE. However, (C) the SAGE correlation continues to improve up to nearly  $r_c = 0.7$  before no genes meet the cutoff and a leveling is not observed. Data represent mean  $r_c$  value and gene pair number of 100 pseudo-random divisions at each MCE. Error bars indicate 1 standard deviation.



142 internal consistency was determined by dividing each of the  
 143 datasets in half and comparing the gene-to-gene Pearson  
 144 correlations for each subset (Figs. 1A–1C). We first divided  
 145 the data in a purely random fashion. To make the internal  
 146 consistency calculation more comparable to the cross-  
 147 platform comparisons, we also devised a pseudo-random  
 148 division, which takes into account the presence of exper-  
 149 imental replicates and very similar experimental conditions  
 150 in the datasets (see Materials and methods).

151 Internal consistency was found to be dependent on the  
 152 minimum number of common experiments (MCE) between  
 153 any two genes on which Pearson correlations are calculated.  
 154 MCE was defined as the minimum required number of  
 155 common or shared experiments for which any two genes  
 156 actually have values available in their respective expression  
 157 profiles (Fig. 1D).

158 Increasing the MCE increased the internal consistency but  
 159 decreased the number of gene pairs considered for both the  
 160 pseudo-random (Fig. 1) and the random (Supplemental Fig.  
 161 2) division methods. With the random division, and an MCE  
 162 of 100, Affymetrix showed the highest average internal  
 163 correlation of 0.925, then cDNA microarray with correlation  
 164 of 0.889, and then SAGE with correlation of 0.776. This  
 165 MCE cutoff was used by the group that provided the cDNA  
 166 microarray data [4] (E. Segal, personal communication). As  
 167 expected, the pseudo-random division, which groups repli-  
 168 cates and experimental datasets, reduced internal consisten-  
 169 cies with values of 0.253 for Affymetrix, 0.273 for the cDNA  
 170 microarray, and 0.660 for SAGE with MCE of 100 (Fig. 1).  
 171 Unfortunately, as the SAGE dataset contains only 242  
 172 samples, division into two groups of approximately 120  
 173 results in relatively few gene pairs that meet the criteria of 100  
 174 MCE (only 1518 pairs on average). Although approximately  
 175 60% of these SAGE libraries are derived from cancer  
 176 samples, we found no evidence of an effect on the  
 177 coexpression results (Supplemental Fig. 3) and therefore  
 178 included them in subsequent analysis.

179 Internal consistency is a measure of the reproducibility or  
 180 robustness of gene coexpression predictions similar to a  
 181 cross-validation test. This is based on the assumption that if a  
 182 gene pair is truly coexpressed based on an expression  
 183 dataset, it should be predicted as coexpressed by random  
 184 subsets of the data. The consistency increases with higher  
 185 MCE but at different rates for the three datasets because of  
 186 their different natures in terms of number of experiments and  
 187 experiment composition. Thus, it would be unfair to compare  
 188 the datasets with MCEs that resulted in different levels of  
 189 reproducibility. Studies generally choose some cutoff for a  
 190 minimum number of common experiments, such as 5, 10, or  
 191 100 [4,30,35]. In an effort to produce an unbiased  
 192 comparison of the three platforms, the pseudo-random  
 193 division was used to determine an appropriate MCE that  
 194 would generate the same internal consistency ( $r_c = 0.25$ ) for  
 195 each (Affymetrix MCE = 95; cDNA MCE = 28; SAGE  
 196 MCE = 23) (Fig. 1). All internal consistency correlations are  
 197 summarized in Table 1.

Table 1

Summary of  $r_c$  values for internal consistency analysis using different  
 sample division methods and MCE cutoffs

Platform	Division	MCE cutoff	Gene pairs	$r_c$ value
Affymetrix	Random	100	4,149,092	0.925
	Pseudo-random	95	3,427,174	0.257
	by GSE series	100	3,260,557	0.253
cDNA microarray	Random	100	10,429,219	0.889
	Pseudo-random	28	11,178,346	0.253
	by author	100	9,747,169	0.273
SAGE	Random	100	2,635	0.776
	Pseudo-random	23	577,820	0.253
	by tissue	100	1,518	0.660

Note that many different divisions are possible for each result (except cancer/normal). Gene pair and  $r_c$  values represent mean values from 100 different random or pseudo-random divisions.

### Cross-platform correlation analysis

199 Considering that the levels of consistency between  
 200 subsets of data from a single platform were relatively low  
 201 (when replicates and similar experiments were kept together)  
 202 it is not surprising that datasets from different platforms  
 203 compared poorly against each other. All comparisons were  
 204 found to have significant but poor positive correlations  
 205 compared to randomly permuted data ( $p < 0.001$ , 1000  
 206 permutations). Affymetrix versus cDNA showed the best  
 207 correlation of 0.102, then Affymetrix versus SAGE with  
 208 0.086, and finally cDNA versus SAGE with 0.041 (Supple-  
 209 mental Fig. 4). A Pearson rank analysis also showed  
 210 significant but poor agreement with only 3–8% better  
 211 performance than randomly permuted data (Supplemental  
 212 Fig. 5).

213 An analysis of correlation at different minimum Pearson  
 214 cutoffs ( $r$  cutoff) for gene pairs was performed as described  
 215 previously [25] (Supplemental Fig. 6). Lee et al. [25] ob-  
 216 served a steady increase in global concordance ( $r_c =$   
 217 correlation of correlations) up to 0.92 at an  $r$  cutoff of  
 218 0.91. Our data did not show such an obvious trend. Global  
 219 concordance stayed close to 0 (or even below) for all three  
 220 pair-wise platform comparisons up to 0.5–0.6 Pearson  
 221 cutoff. The Affymetrix/cDNA correlation did show an  
 222 improvement to  $r_c = 0.163$  ( $p = 0.003$ ,  $n = 289$  gene pairs)  
 223 at an  $r$  cutoff of 0.65. Similarly the Affymetrix/SAGE  
 224 comparison improved to  $r_c = 0.290$  ( $p = 0.028$ ,  $n = 44$  gene  
 225 pairs) at an  $r$  cutoff of 0.7. After these cutoffs, both  
 226 Affymetrix/cDNA and Affymetrix/SAGE comparisons  
 227 returned to  $r_c$  values close to 0 (or below) and were reduced  
 228 to insignificant gene pair numbers. The cDNA/SAGE  
 229 comparison showed no significant increases in  $r_c$  with any  
 230  $r$  cutoff.

### Gene ontology analysis

231 Since the datasets under study demonstrated little  
 232 agreement, we attempted to determine which dataset was  
 233 most “biologically relevant.” GO biological process  
 234

235 domain knowledge [36] was used to evaluate gene  
 236 coexpression predictions for each platform. We hypothe-  
 237 sized that genes that are coexpressed will be more likely to  
 238 be involved in the same biological process. The number of  
 239 gene pairs annotated to the same “most specific” GO  
 240 (Biological Process) term for each platform was deter-  
 241 mined (Supplemental Fig. 7). In general, the datasets from  
 242 all platforms perform better than expected by chance.  
 243 Affymetrix performed best, followed by cDNA microarray  
 244 and SAGE, which performed about equally better than  
 245 randomly permuted data. The analysis was also extended  
 246 up the GO hierarchy to parent and grandparent terms, and  
 247 identical trends and relationships were observed (Supple-  
 248 mental Fig. 8).

249 A second analysis looked at the relationship between the  
 250 Pearson correlation and the performance against GO. For  
 251 each platform, the number of gene pairs annotated to the  
 252 same “most specific term” at different Pearson correlation  
 253 ranges was determined (Fig. 2). Generally, as Pearson  
 254 correlation for a gene pair increases it is more likely to be  
 255 confirmed by GO. With a Pearson value in the range of  
 256 0.3–0.4 or better the platforms always performed signifi-  
 257 cantly better than randomly permuted data ( $p < 0.001$ ,  
 258 1000 permutations). The improvement over randomly  
 259 permuted data was very slight for the cDNA and SAGE  
 260 datasets (2–4%). However, for the Affymetrix data, the  
 261 trend was striking. Gene pairs identified as coexpressed  
 262 with a Pearson correlation of 0.9–1.0 were confirmed by  
 263 GO in 74% of cases. Gene pairs from this list include a  
 264 large set of highly coexpressed protein biosynthesis genes  
 265 as well as a few genes involved in translational elongation  
 266 (a subprocess of protein biosynthesis) and muscle contrac-  
 267 tion. It should be noted that, in the case of the SAGE and  
 268 cDNA datasets, only a few gene pairs had Pearson  
 269 correlations  $>0.9$  (one for cDNA, five for SAGE).

270 A third analysis examined the effect of averaging  
 271 platform results and comparing to individual platforms  
 272 using GO. Requiring coexpression evidence from multiple  
 273 datasets may represent a method of reducing noise and  
 274 increase our confidence that coexpressed genes are actually  
 275 coregulated. The percentage of gene pairs annotated to the  
 276 same most specific term at different average Pearson  
 277 correlation ranges was determined as above. The results  
 278 were again quite striking. With a two-platform combined  
 279 Pearson of 0.4 or greater the combined platforms all  
 280 performed significantly better than randomly permuted data  
 281 ( $p < 0.005$ , 1000 permutations). Furthermore, for any  
 282 platform combination, a gene pair with an average Pearson  
 283 correlation of  $r > 0.6$  was much more likely to share a GO  
 284 term than a gene pair with this level of correlation in only a  
 285 single platform (Fig. 3). For example, a gene pair with a  
 286 two-platform average Pearson of 0.7–0.8 was found to  
 287 share a common GO term 40–50% of the time. Pairs with  
 288 this same Pearson range in individual datasets shared a  
 289 common GO term only 5–10% of the time, only a few  
 290 percent better than expected by chance. Gene pairs

confirmed by multiple datasets ( $r_{\text{avg}} > 0.6$  for any two 291  
 platforms) covered a wide range of GO categories (52 in 292  
 total) (Supplemental Fig. 9). 293

#### 294 *Comparison to other coexpression methods* 294

295 Finally, an analysis was conducted to assess two other 295  
 recent coexpression studies that were published while this 296  
 analysis was in progress. The ArrayProspector method [37], 297  
 the TMM method [35], and our two-platform combination 298  
 method (2PC) were each mapped to UniProt IDs and 299  
 assessed using the same GO analysis as above. In all three 300  
 cases, we observed significantly more gene pairs with 301  
 common GO terms at higher scores (Fig. 4). For our method 302  
 (2PC), the percentage of gene pairs with a common GO 303  
 term rises sharply at a score of approximately 0.6–0.7. For 304  
 ArrayProspector this occurs at a score of approximately 305  
 0.7–0.8 and for TMM at a score of 5–6. At these cutoffs, 306  
 each method represents 2500 to 10,000 gene pairs. Each 307  
 utilizes different genes and identifies different gene pairs as 308  
 highly coexpressed. Thus, a comparison of the highest 309  
 scoring 2500 gene pairs for each found only a minimal 310  
 overlap of less than 10% (Fig. 4D). 311

#### 312 **Discussion** 312

313 We have shown that the genes identified as coexpressed 313  
 are highly dependent on the dataset and expression platform 314  
 used. In general, we find that the more data a correlation is 315  
 based on, the more reproducible it is. When division of 316  
 samples takes similar or replicate experiments into consid- 317  
 eration, Affymetrix and cDNA internal consistencies level 318  
 off at approximately  $r_c = 0.25$  with MCE of about 90 and 319  
 30–40, respectively. The SAGE dataset continued to 320  
 improve to nearly  $r_c = 0.6$  with MCE of 80. This may 321  
 reflect the diverse nature of the SAGE dataset for which 322  
 libraries are rarely constructed from the same or similar 323  
 tissue. In contrast, it is not uncommon for many Affymetrix 324  
 or cDNA experiments to measure expression of a very 325  
 similar series of samples. A recent yeast study found that the 326  
 ability to identify coregulated genes correctly from coex- 327  
 pression analyses is highly dependent on the number of 328  
 experiments, with accuracy leveling off at 50 to 100 329  
 experiments [38]. Our results agree closely with this 330  
 observation for human data and suggest that coexpression 331  
 predictions will be most reproducible if based on 30 to 100 332  
 experiments. Furthermore, global coexpression analysis 333  
 may benefit from a greater representation of tissues and 334  
 conditions rather than greater numbers. 335

336 Given that different experimental subsets of the same 336  
 platform show poor correlation it is perhaps not surprising 337  
 that interplatform comparisons show very poor correlations 338  
 ( $r < 0.11$ ). The fact that none of these datasets agree well 339  
 raises some serious questions about their use for validation 340  
 and integration with other data. There are several possible 341

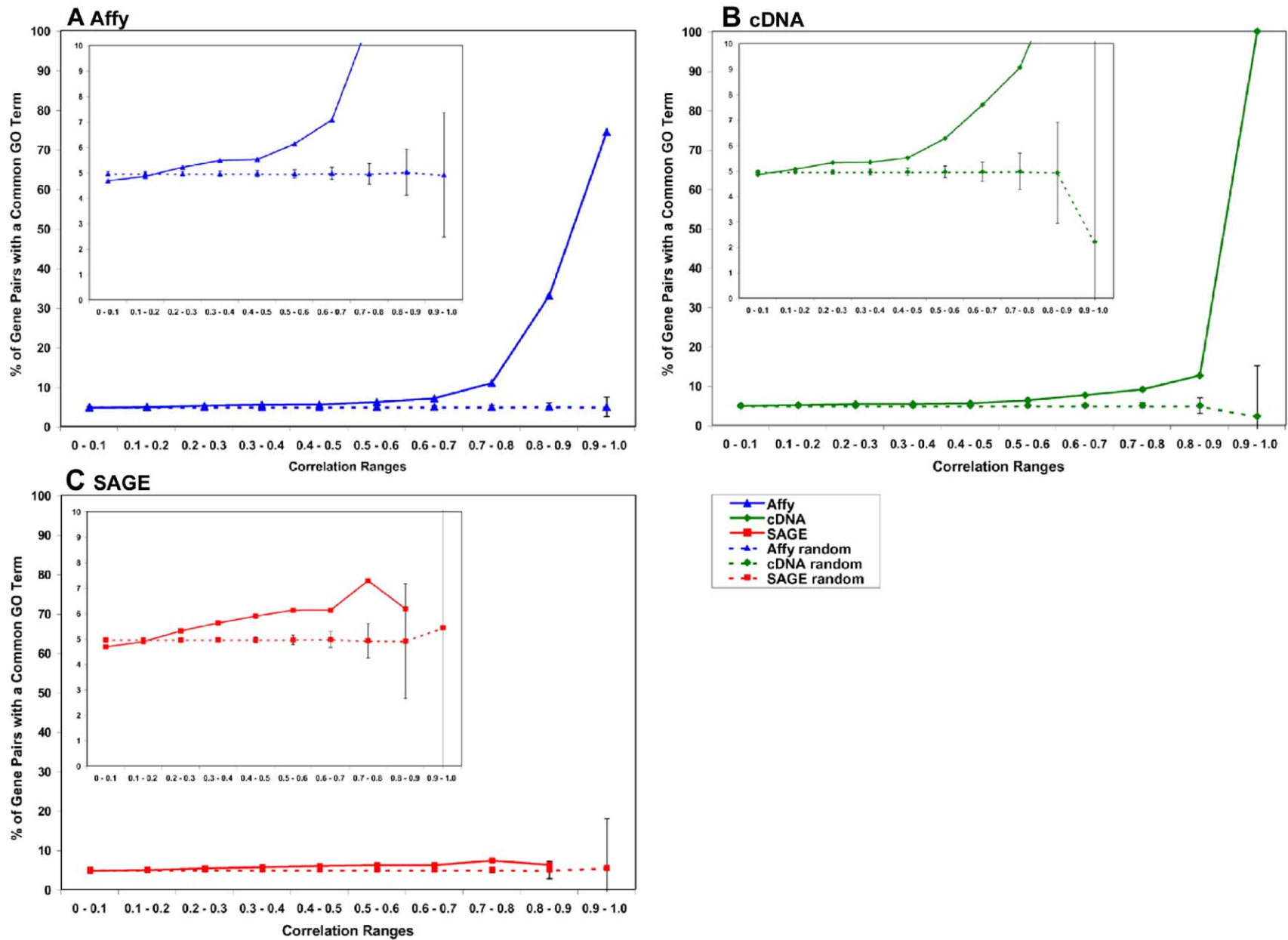


Fig. 2. GO correlation range analysis. At higher Pearson correlations (in particular,  $r > 0.8$ ) gene pairs are more likely to have similar GO biological processes, although very few gene pairs have high correlations in the SAGE and cDNA datasets. 75% of gene pairs with correlation  $>0.9$  calculated from Affymetrix data have the same GO annotation. Interestingly, gene pairs with very low Pearson values are less likely to share a common GO term than randomly permuted data. Random lines represent mean values from 1000 random permutations. Error bars indicate 1 standard deviation.

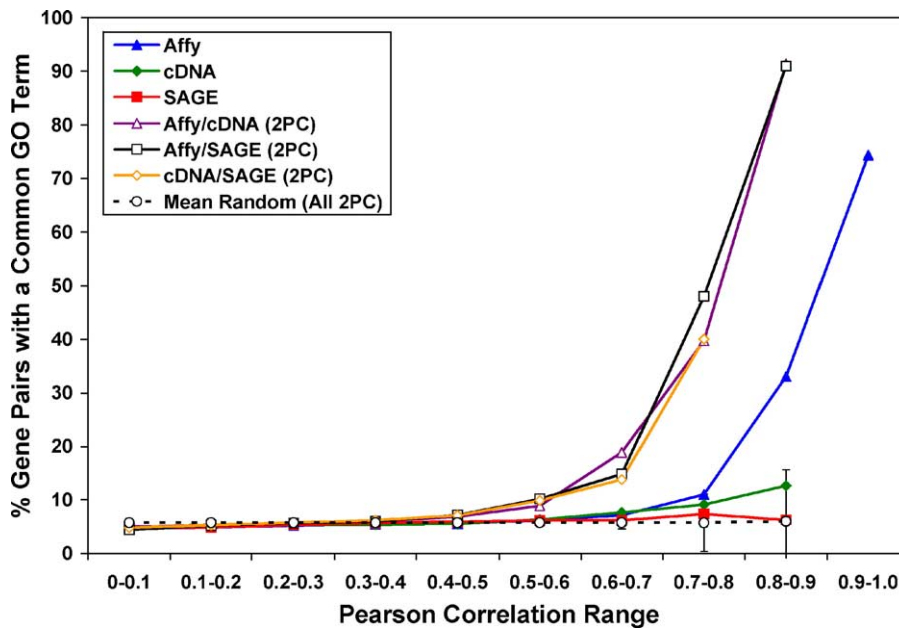


Fig. 3. GO correlation range analysis for multiplatform average. Comparison of two-platform average Pearson to individual platform indicates that gene pairs identified as coexpressed in multiple platforms (higher average Pearson) are much more likely to be confirmed by GO. Random line represents mean values from 1000 random permutations of all two-platform combinations. Error bars indicate 1 standard deviation.

342 explanations for this observation: (1) The data comprising  
 343 these datasets are so noisy as to prevent reliable identi-  
 344 fication of many truly coexpressed genes, (2) the method of  
 345 identifying coexpressed genes is inadequate, (3) the  
 346 unmatched and nonoverlapping nature of the samples that  
 347 make up each dataset results in identification of different  
 348 subsets of truly coexpressed genes, and (4) genes are under  
 349 such complex regulatory control that genes coregulated in  
 350 one cell type or tissue behave in an entirely different manner  
 351 in other cell types or tissues and are therefore not globally  
 352 coexpressed. It is likely that each of the explanations  
 353 outlined above is to some degree responsible for the lack of  
 354 concordance between coexpression analyses produced from  
 355 different datasets and different platforms. It is not the  
 356 purpose of this study to identify which is most important.  
 357 Rather, we wish to make researchers aware that the choice  
 358 of dataset or platform for integration or validation of other  
 359 data could dramatically affect their results, and methods that  
 360 integrate or combine different platforms may be more  
 361 appropriate.

362 The fact that intraplatform comparisons show some  
 363 correlation and improve with number of data points suggests  
 364 that some gene pairs identified are truly coexpressed.  
 365 Furthermore, the GO analysis shows that gene pairs  
 366 identified as highly coexpressed (higher Pearson correla-  
 367 tion) are more likely to share the same biological process  
 368 and thus actually be related. Similarly, gene pairs with lower  
 369 Pearson correlations were as or less likely than random  
 370 chance to share the same biological process. These results  
 371 suggest that the Pearson correlation is a useful metric and  
 372 that both high and low Pearson values have the meaning we  
 373 expect. The GO analysis did not conclusively identify a

single “correct” platform or dataset but it did show that the  
 Affymetrix dataset identified more biologically relevant  
 gene pairs than the cDNA or SAGE dataset. However, gene  
 pairs coexpressed in multiple expression platforms were  
 much more likely to be confirmed by GO. Thus, combining  
 platforms appears to act as a filter, producing high-  
 confidence predictions from noisy datasets. This conclusion  
 is based on the assumption that coexpressed genes are more  
 likely to be biologically relevant if they share common  
 biological processes. Assessments using GO are limited by  
 issues such as the incompleteness of the ontology, the  
 potential for circularity (addressed in Materials and meth-  
 ods), experimental bias toward “well-studied” genes, and  
 inconsistencies in structure and depth. Furthermore, it is  
 likely that some coregulated genes will belong to different  
 biological processes while other genes involved in the same  
 process will not be coregulated. As such, an “absolute”  
 performance against GO is difficult or impossible to define.  
 Despite these issues, we believe GO currently represents  
 one of the best resources for a relative assessment of  
 coexpression platforms or methods.

Recent investigations into the utility of combining  
 expression data from different high-throughput platforms  
 have identified highly variable levels of agreement. Based  
 on an analysis of a small set of matched samples using  
 oligonucleotide arrays, SAGE, and EST data, Haverty and  
 colleagues [39] caution against the combination of plat-  
 forms to confirm expression patterns for specific sets of  
 genes. However, they do suggest that such methods can be  
 used to extract high-confidence subsets of related genes.  
 We agree that for many genes a poor level of agreement  
 between datasets raises questions about their utility.

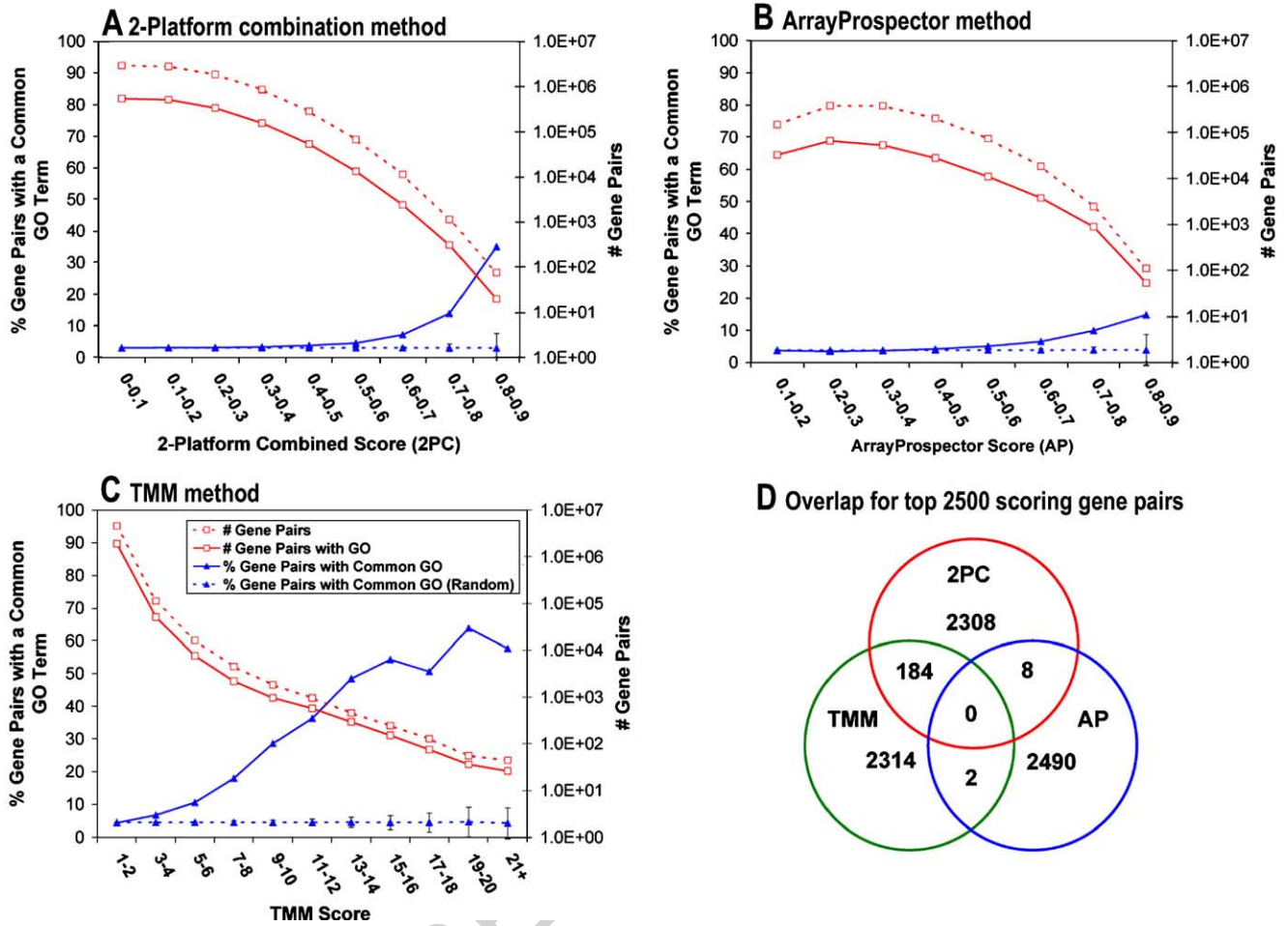


Fig. 4. Comparison of two-platform combination method to other recent coexpression methods. (A–C) For each method, gene pairs with higher scores are more likely to share a common GO term. Lines with open squares represent numbers of gene pairs (right axis). Lines with closed triangles represent percentage of gene pairs with a common GO term. Random lines represent mean values from 1000 random permutations. Error bars indicate 1 standard deviation. (D) Venn diagram indicates overlap between the 2500 top scoring pairs for each method (not required to be in GO). Each method comprises different datasets and has different genes. Therefore a direct comparison of method performance is difficult. Instead, the graphs illustrate that each method is capable of identifying biologically relevant gene-pair relationships and the Venn diagram indicates that they identify very different sets of relationships. Furthermore, the GO analysis provides a means of choosing reasonable score thresholds for each method to generate lists of high-confidence coexpressed genes.

406 However, our results do show that platform combination  
 407 methods can be extended to large sets of unmatched  
 408 publicly available expression data to produce biologically  
 409 meaningful information.

410 As we were nearing completion of our analysis, a similar  
 411 study using multiple microarray datasets (TMM) was  
 412 published [35]. The authors examined 60 microarray data-  
 413 sets (cDNA and Affymetrix oligonucleotide) for gene pairs  
 414 identified as coexpressed in multiple datasets. They report  
 415 that even gene pairs confirmed by only a single dataset have  
 416 better GO similarity scores than random pairs and GO score  
 417 increases steadily with the number of confirmed links. Their  
 418 method differs from ours in that experimental subsets are  
 419 analyzed separately and a “vote-counting” method was used  
 420 to identify gene pairs that appear highly coexpressed (above  
 421 some Pearson cutoff) in multiple sets. Our method combines  
 422 all experimental subsets into a single dataset for each  
 423 expression platform and then averages the global Pearson

424 correlations between platforms. Our method is also the first  
 425 to include SAGE data. A third recently published method  
 426 (ArrayProspector) used a combination of singular value  
 427 decomposition and kernel density estimation [37]. This  
 428 method combines evidence from related arrays and weights  
 429 the contribution of each array according to how well they  
 430 correlate with functional annotation.

431 When attempting to infer function or coregulation from  
 432 coexpression we should consider that it is likely that genes  
 433 are biologically related in a number of different ways and  
 434 therefore different methods will be required to identify each  
 435 type of relationship. For example, one pair of genes might  
 436 be “tightly” coexpressed only under very specific condi-  
 437 tions, whereas another gene pair might be “loosely”  
 438 coexpressed across a broad range of conditions depending  
 439 on the regulatory elements that they share. The three  
 440 methods discussed above (TMM, ArrayProspector (AP),  
 441 and 2PC) represent three different approaches to the



442 problem of identifying high-confidence coexpression for the  
 443 purpose of inferring function or coregulation. Because the  
 444 methods use different datasets and scoring methods and  
 445 comprise different gene sets, a direct comparison of the  
 446 methods is difficult. Therefore, we chose simply to assess  
 447 their respective predictions against GO independently. Thus,  
 448 we do not identify the “best” method but rather show that  
 449 each method is at least partially effective based on  
 450 performance against the Gene Ontology. Furthermore,  
 451 because the highest scoring pairs for each are almost  
 452 completely nonoverlapping we advocate combining the best  
 453 results of each into a single set of high-confidence  
 454 predictions. To this end we have chosen score thresholds  
 455 for each method based on GO performance ( $2PC > 0.65$ ;  
 456  $AP > 0.7$ ;  $TMM > 7$ ) and make available a list of 13,145  
 457 high-confidence coexpressed gene pairs (representing 2979  
 458 unique genes) ([http://www.bcgsc.ca/gc/bomge/coexpression/  
 459 suppl\\_materials](http://www.bcgsc.ca/gc/bomge/coexpression/suppl_materials)) for use in regulatory element prediction or  
 460 other integration studies.

## 461 Materials and methods

### 462 Data sources

463 Human gene expression data for three major expression  
 464 platforms were collected from public sources. We used a  
 465 recently published dataset of 1202 cDNA microarray  
 466 experiments [4] representing 13,595 genes, 242 SAGE  
 467 libraries from the Gene Expression Omnibus (GEO) ([http://  
 468 www.ncbi.nlm.nih.gov/geo/](http://www.ncbi.nlm.nih.gov/geo/)) representing 15,426 genes,  
 469 and 667 Affymetrix HG-U133A oligonucleotide micro-  
 470 array experiments (889 were available but 667 had PMA  
 471 detection calls) representing 8106 genes, also from GEO  
 472 (Supplemental Fig. 1). cDNA microarray genes provided  
 473 by Stuart et al. [4] were identified by LocusLink IDs [40].  
 474 Therefore this identifier was used for the other two  
 475 platforms to allow the gene intersection of the three  
 476 datasets to be determined and used for the subsequent  
 477 analyses.

### 478 Data filtering

479 cDNA microarray data for 13,595 genes were used as  
 480 provided by Stuart et al. [4] except for minor formatting  
 481 changes (see supplementary materials for our data). The  
 482 242 SAGE libraries ranged from 1430 to 308,589 total  
 483 tags in size with an average size of 52,723. SAGE data  
 484 were first filtered to remove tags with less than one count  
 485 in at least 10 libraries, reducing the unique tags from  
 486 609,224 to 87,521 (and total tags from 12,758,981 to  
 487 11,219,373). Next, SAGE tags were mapped to genes by  
 488 the “lowest” sense-strand tag predicted from RefSeq [40]  
 489 or MGC [41] sequences and then mapped to LocusLink  
 490 IDs using the DiscoverySpace software package (Varhol et  
 491 al., unpublished, <http://www.bcgsc.ca/discoveryspace/>),

reducing the tag set further to 47,263 unique tags. 492  
 Generally, the lowest tag corresponded to the canonical 493  
 3′-most *Nla*III anchoring enzyme recognition site (position 494  
 1) expected for the gene sequence. However, if such a 495  
 canonical match was not found, higher position (less 3′) 496  
 mappings were also accepted (see supplementary materials 497  
 for more details). In the event of discrepancy between 498  
 RefSeq and MGC, the former was taken as correct because 499  
 a larger number of tags could be mapped with this 500  
 resource (9568 vs. 6295) and was thus perceived to be 501  
 more complete. Only 297 tags with disagreements between 502  
 RefSeq and MGC are represented in the final gene set 503  
 (~5%) (see supplementary materials for more details). If a 504  
 tag mapped to more than one LocusLink or more than one 505  
 tag mapped to the same LocusLink it was discarded, 506  
 resulting in a final set of 15,426 unique tags (2,762,500 507  
 total tags) confidently mapped to LocusLink IDs. We 508  
 mapped 22,215 Affymetrix probe IDs to 20,577 LocusLink 509  
 IDs using the most current Affymetrix annotation file for 510  
 the HG-U133A chip ([www.affymetrix.com](http://www.affymetrix.com), supplementary 511  
 materials). As with the SAGE tags, probes with ambiguous 512  
 mapping to LocusLink were discarded, resulting in a final 513  
 set of 8106 genes from the Affymetrix dataset. Once 514  
 LocusLink IDs were available for all three platforms, the 515  
 intersection was determined. This subset of 5881 genes, 516  
 present in all three platforms, was used for all subsequent 517  
 analyses. The final 5881 unique SAGE tags represent 518  
 1,173,430 total tags sequenced. 519

### Distance calculations 520

Ratio values for the cDNA microarray data were used as 521  
 is for the Pearson calculation. Affymetrix probe intensities 522  
 were converted to natural log values. All  $\ln(\text{intensity})$  values 523  
 were normalized by subtracting the median and dividing by 524  
 the interquartile range for the experiment [42]. Only 525  
 Affymetrix probe intensities with a “P” call were consi- 526  
 dered ( $p$  value  $< 0.04$ ). Intensities with “A” or “M” calls 527  
 were set to null. To compensate for different library sizes 528  
 SAGE tag counts were normalized to 10,000 tags/library 529  
 and log-transformed as follows [34]: 530

### Tag frequency

$$= \ln((\text{tag count} \times 10000)/\text{total tags in library})$$

SAGE tag counts of 0 were converted to nulls. In all 533  
 platforms, genes are represented by a vector of expression 534  
 values for all the experiments in the dataset. In each case, 535  
 genes have null values if not represented on that array 536  
 (cDNA), no tags were observed (SAGE), or intensity was 537  
 not significantly detected (Affymetrix). Thus, when calcul- 538  
 ating Pearson correlations between gene pairs, the number 539  
 of shared data points varied from 0 to the total number of 540  
 experiments. A minimum number of common experiments 541  
 was required for each gene pair to provide some confidence 542  
 in the value calculated (a Pearson correlation based on 543

544 observations from only two experiments is meaningless). A  
 545 range of MCEs was used for the internal consistency  
 546 analysis (see below) and then one minimum chosen for  
 547 subsequent analyses.

548 A Pearson correlation coefficient was calculated for all  
 549 possible gene pairs for each platform as a measure of  
 550 expression similarity. These calculations were performed by  
 551 a modified version of the C clustering library [43] on 64-bit  
 552 Opteron Linux machines with 8- to 32-GB memory. Please  
 553 see supplementary materials for modified C source code and  
 554 explanation of changes.

#### 555 *Correlation of correlations analysis*

556 Correlation of correlations ( $r_c$ ) for internal consistencies  
 557 and platform comparisons were performed as previously  
 558 described [25] using the Pearson correlation function (cor)  
 559 of the R statistical package (version 1.8.1). This correlation  
 560 involves millions of data points and thus cannot be graphed  
 561 easily. Therefore, data were binned and density plots created  
 562 using the Bioconductor hexbin (version 1.0.3) add-in  
 563 function for R [44].

#### 564 *Internal consistency analysis*

565 To evaluate the consistency of coexpression observed  
 566 within each platform, we divided the experiments available  
 567 and determined coexpression for each subset independently.  
 568 If a platform consistently finds coexpressed genes regardless  
 569 of the exact experiments involved, the  $r_c$  will be close to 1.  
 570 To determine whether the observed  $r_c$  is significant, we  
 571 repeat the procedure with randomly permuted gene expres-  
 572 sion values, expecting an  $r_c$  close to 0.

#### 573 *Pseudo-random division method*

574 Division was performed first randomly and then  
 575 pseudo-randomly. The pseudo-random division was nec-  
 576 essary to prevent artificially high internal consistencies  
 577 resulting from comparing mostly replicates (or very  
 578 similar experiments) in the two subsets. In many cases  
 579 (especially for the Affymetrix data) experimental repli-  
 580 cates or very similar samples exist in the dataset. The  
 581 purpose of coexpression analysis is to identify genes that  
 582 behave similarly across many conditions. The internal  
 583 consistency analysis is meant to measure how consistently  
 584 a series of experiments across different conditions would  
 585 identify the same coexpressed genes. If the two subsets of  
 586 experiments contain replicates, they are more likely to  
 587 identify the same coexpressed genes as the expression  
 588 values of the replicates will be very similar. The cross-  
 589 platform comparisons do not have this advantage because  
 590 they consist of different experiments. Thus, to make the  
 591 internal consistency calculation more comparable to the  
 592 cross-platform comparisons, we used a pseudo-random  
 593 division for subsequent analysis. Experiments were

randomly divided into two subsets but experiments 594  
 belonging to the same experimental series (Affymetrix), 595  
 publication (cDNA), or tissue (SAGE) were required to 596  
 fall into the same subset. 597

#### *Minimum common experiments analysis*

Differences in the number of common experiments 599  
 between any two genes result from missing values in the 600  
 data matrices. In the case of the cDNA microarray data, 601  
 different arrays were used in different experiments, and not 602  
 all genes are present on all the arrays. For SAGE, a tag is 603  
 often observed in one library but will have a 0 tag count in 604  
 other libraries. For Affymetrix oligonucleotide arrays, an 605  
 intensity is always reported for every probe but in some 606  
 cases the Affymetrix statistical software will determine that 607  
 the probe was not reliably detected and assign an absent (A) 608  
 or marginal (M) call instead of a present (P) call for that 609  
 probe. As missing SAGE tags and probes not called P 610  
 represent genes expressed below the detection threshold of 611  
 the SAGE and Affymetrix array experiments, we did not 612  
 include these data in our analysis. Thus, for each dataset, 613  
 there were gene pairs that were rarely represented in the 614  
 same experiment and their Pearson correlations were based 615  
 on very few data points. The effect of number of common 616  
 experiments on internal consistency was determined by 617  
 calculating the internal consistency for a series of datasets 618  
 with different MCE criteria. One hundred different pseudo- 619  
 random divisions were performed to get an average internal 620  
 consistency for each MCE. An MCE was chosen for each 621  
 such that the same internal consistency would result ( $r =$  622  
 0.25) (Fig. 1). Thus, all subsequent analyses were based on 623  
 an MCE of 95 for Affymetrix, 28 for cDNA, and 23 for 624  
 SAGE. Requiring an MCE removes gene pairs from the 625  
 datasets. To maintain an unbiased comparison, only the 626  
 1,173,330 gene pairs common to all three platform datasets 627  
 (after application of MCE criteria) were used in the 628  
 subsequent platform comparisons. 629

#### *Platform comparisons*

As with the internal consistency analysis, a correlation of 631  
 gene correlations was calculated, but was determined for 632  
 each of the three pair-wise platform comparisons instead of 633  
 between subsets of one platform. If the two platforms being 634  
 compared report the same correlation between each gene 635  
 pair, we expect the overall correlation between platforms 636  
 would be near 1. The global concordance ( $r_c$ ) was 637  
 determined for increasing gene correlation cutoffs to 638  
 compare to results obtained in the NCI-60 study [25]. 639

#### *Gene Ontology analysis*

The GO is a controlled vocabulary that describes the 641  
 roles of genes and proteins in all organisms [36]. GO is 642  
 composed of three independent ontologies: biological 643

644 process, molecular function, and cellular component. The  
645 GO descriptive terms are represented as nodes connected by  
646 directed edges that may have more than one parent node  
647 (directed acyclic graph). A gene is annotated to its most  
648 specific GO term description and all ancestor GO terms are  
649 implied.

650 The GO MySQL database dump (release 200402 of  
651 assocdb) was downloaded from [http://www.godatabase.org/  
652 dev/database](http://www.godatabase.org/dev/database). A GO MySQL database was built and a Perl  
653 script was developed to extract three GO information  
654 subspaces from the biological process ontology: (1) the  
655 most specific GO terms for each gene, (2) the most specific  
656 terms along with their associated parent terms, and (3) the  
657 most specific terms along with their associated parent and  
658 grandparent terms. Two categories of annotations were used  
659 for the evaluation of each GO information subspace: (1) gene  
660 annotations that did not include those derived from inferred  
661 electronic annotations (IEAs) (1007 genes found in common  
662 with our dataset) and (2) gene annotations including IEAs  
663 (1426 genes found in common with our dataset). Similar  
664 results were obtained for both non-IEA and IEA analyses.  
665 Only the IEA results are reviewed in the figures and text.

666 One potential issue with our analysis is that of a circular  
667 argument. It is possible that a coexpressed gene pair could  
668 be found to share a common GO term that was annotated for  
669 both genes by a coexpression analysis. Thus, coexpression  
670 data could be confirming coexpression data. To check for  
671 this problem we assessed the degree to which our dataset  
672 depends on annotations inferred from expression profiles  
673 (IEP evidence code). Only 93 of 32,669 biological process  
674 annotations use IEP evidence, corresponding to only 73  
675 genes with 1 or more IEP annotations. Of these, only 1 was  
676 present in our gene set and this gene also had non-IEP  
677 annotations. Therefore the potential for a circular argument  
678 is negligible.

679 Results shown in Supplemental Fig. 7 were extracted  
680 from the gene pair correlation data by enumerating the  
681 number of gene pairs found at common GO terms across a  
682 gene's expression similarity neighborhood for each GO  
683 information subspace. Results shown in Fig. 2 were  
684 extracted by enumerating the number of gene pairs found  
685 at common GO terms for each range of Pearson correlations  
686 from 0 to 1 in increments of 0.1. The results summarized in  
687 Fig. 3 were enumerated in a similar manner but used  
688 average Pearson correlations between two platforms instead  
689 of individual Pearson correlations. One thousand random  
690 permutations of the data were conducted to determine how  
691 often GO confirmation of a gene pair at each neighborhood  
692 or Pearson range would occur by chance. Scripts were  
693 written in Perl and are available at [http://www.bcgsc.ca/gc/  
694 bomge/coexpression/suppl\\_materials](http://www.bcgsc.ca/gc/bomge/coexpression/suppl_materials).

#### 695 *Comparison to other coexpression methods*

696 Results shown in Fig. 4 were generated using the GO  
697 analysis method described above for Figs. 2 and 3. AP

data were obtained by request from the author [37]. Only  
pairs with scores above 0.150 were provided. TMM data  
were downloaded from the authors' supplemental Web  
page (see Web references) [35]. Both negative and positive  
correlations were included and thus a gene pair can appear  
twice. Only pairs with scores of 1 or greater were  
provided. The 2PC method represents all two-platform  
averages (Affymetrix/cDNA, Affymetrix/SAGE, and  
cDNA/SAGE). Thus, a gene pair can appear as many as  
three times if all three pair-wise averages fall within the  
0–1 range graphed. All datasets were converted from their  
respective identifiers to UniProt [45] and the percentage of  
gene pairs found at common GO terms for each range of  
scores was determined. The top 2500 pairs of each were  
examined to determine the overlap in results for high-  
scoring pairs. Thresholds for a high-confidence set of  
coexpressed gene pairs were chosen for each method at the  
approximate respective score at which performance was at  
least three to four times better than random chance (2PC >  
0.65; AP > 0.7; TMM > 7).

#### Acknowledgments

We thank Misha Bilenky, Jaswinder Khattrra, Sheldon  
McKay, Gordon Robertson, Peter Ruzanov, Kim Wong, and  
Scott Zuyderduyn for invaluable help and advice and Allen  
Delaney and Marco Marra for editorial suggestions. Eran  
Segal and Josh Stuart provided useful information regarding  
the published cDNA data. Lars Jensen (Peer Bork lab) and  
Paul Pavlidis generously provided their coexpression data-  
sets through personal communication or supplemental Web  
pages. O.G. and E.P. are Michael Smith Foundation for  
Health Research Trainees and are supported by the Natural  
Sciences and Engineering Research Council of Canada. D.F.  
is supported by the CIHR/MSFHR Bioinformatics Training  
Program. S.J. is a Michael Smith Foundation for Health  
Research Scholar. We also gratefully acknowledge the  
support of Genome BC, the BC Cancer Agency, and the  
BC Cancer Foundation.

#### Appendix A. Supplementary data

Supplementary data associated with this article can be  
found, in the online version, at [doi:10.1016/j.ygeno.2005.  
06.009](https://doi.org/10.1016/j.ygeno.2005.06.009).

#### References

- [1] M. Schena, D. Shalon, R.W. Davis, P.O. Brown, Quantitative monitoring of gene expression patterns with a complementary DNA microarray, *Science* 270 (1995) 467–470.
- [2] D.J. Lockhart, et al., Expression monitoring by hybridization to high-density oligonucleotide arrays, *Nat. Biotechnol.* 14 (1996) 1675–1680.

- 748 [3] V.E. Velculescu, L. Zhang, B. Vogelstein, K.W. Kinzler, Serial  
749 analysis of gene expression, *Science* 270 (1995) 484–487.
- 750 [4] J.M. Stuart, E. Segal, D. Koller, S.K. Kim, A gene-coexpression  
751 network for global discovery of conserved genetic modules, *Science*  
752 302 (2003) 249–255.
- 753 [5] S. Li, et al., A map of the interactome network of the metazoan  
754 *C. elegans*, *Science* 303 (2004) 540–543.
- 755 [6] P. Kemmeren, et al., Protein interaction verification and functional  
756 annotation by integrated analysis of genome-scale data, *Mol. Cell* 9  
757 (2002) 1133–1143.
- 758 [7] A.J. Walhout, et al., Integrating interactome, phenome, and tran-  
759 scriptome mapping data for the *C. elegans* germline, *Curr. Biol.* 12  
760 (2002) 1952–1958.
- 761 [8] E. Segal, et al., Module networks: identifying regulatory modules and  
762 their condition-specific regulators from gene expression data, *Nat.*  
763 *Genet.* 34 (2003) 166–176.
- 764 [9] E.J. Yeoh, et al., Classification, subtype discovery, and prediction of  
765 outcome in pediatric acute lymphoblastic leukemia by gene expression  
766 profiling, *Cancer Cell* 1 (2002) 133–143.
- 767 [10] A. Nimgaonkar, et al., Reproducibility of gene expression across  
768 generations of Affymetrix microarrays, *BMC Bioinform.* 4  
769 (2003) 27.
- 770 [11] P.K. Tan, et al., Evaluation of gene expression measurements from  
771 commercial microarray platforms, *Nucleic Acids Res.* 31 (2003)  
772 5676–5684.
- 773 [12] S. Blackshaw, et al., MicroSAGE is highly representative and  
774 reproducible but reveals major differences in gene expression  
775 among samples obtained from similar tissues, *Genome Biol.* 4  
776 (2003) R17.
- 777 [13] L. Huminiecki, A.T. Lloyd, K.H. Wolfe, Congruence of tissue  
778 expression profiles from Gene Expression Atlas, SAGEmap and  
779 TissueInfo databases, *BMC Genom.* 4 (2003) 31.
- 780 [14] V. Detours, J.E. Dumont, H. Bersini, C. Maenhaut, Integration and  
781 cross-validation of high-throughput gene expression data: comparing  
782 heterogeneous data sets, *FEBS Lett.* 546 (2003) 98–102.
- 783 [15] A.K. Jarvinen, et al., Are data from different gene expression  
784 microarray platforms comparable? *Genomics* 83 (2004) 1164–1168.
- 785 [16] C.A. Iacobuzio-Donahue, et al., Highly expressed genes in pancreatic  
786 ductal adenocarcinomas: a comprehensive characterization and compar-  
787 ison of the transcription profiles obtained from three major  
788 technologies, *Cancer Res.* 63 (2003) 8614–8622.
- 789 [17] H.L. Kim, Comparison of oligonucleotide-microarray and serial  
790 analysis of gene expression (SAGE) in transcript profiling analysis  
791 of megakaryocytes derived from CD34<sup>+</sup> cells, *Exp. Mol. Med.* 35  
792 (2003) 460–466.
- 793 [18] A.T. Rogojina, W.E. Orr, B.K. Song, E.E. Geisert Jr., Comparing  
794 the use of Affymetrix to spotted oligonucleotide microarrays using  
795 two retinal pigment epithelium cell lines, *Mol. Vision* 9 (2003)  
796 482–496.
- 797 [19] M. Ishii, et al., Direct comparison of GeneChip and SAGE on the  
798 quantitative accuracy in transcript profiling analysis, *Genomics* 68  
799 (2000) 136–143.
- 800 [20] S.J. Evans, et al., Evaluation of Affymetrix Gene Chip sensitivity in  
801 rat hippocampal tissue using SAGE analysis: serial analysis of gene  
802 expression, *Eur. J. Neurosci.* 16 (2002) 409–413.
- 803 [21] J. Li, M. Pankratz, J.A. Johnson, Differential gene expression patterns  
804 revealed by oligonucleotide versus long cDNA arrays, *Toxicol. Sci.* 69  
805 (2002) 383–390.
- 806 [22] W.P. Kuo, T.K. Jenssen, A.J. Butte, L. Ohno-Machado, I.S. Kohane,  
807 Analysis of matched mRNA measurements from two different micro-  
808 array technologies, *Bioinformatics* 18 (2002) 405–412.
- 809 [23] N. Mah, et al., A comparison of oligonucleotide and cDNA-based  
810 microarray systems, *Physiol. Genom.* 16 (2004) 361–370.
- 811 [24] J. Lu, A. Lal, B. Merriman, S. Nelson, G. Riggins, A comparison of  
812 gene expression profiles produced by SAGE, long SAGE, and  
813 oligonucleotide chips, *Genomics* 84 (2004) 631–636.
- 814 [25] J.K. Lee, et al., Comparing cDNA and oligonucleotide array data:  
concordance of gene expression across platforms for the NCI-60  
815 cancer cells, *Genome Biol.* 4 (2003) R82. 816
- [26] D.J. Allicco, I.S. Kohane, A.J. Butte, Quantifying the relationship  
817 between co-expression, co-regulation and gene function, *BMC*  
818 *Bioinform.* 5 (2004) 18. 819
- [27] S.K. Kim, et al., A gene expression map for *Caenorhabditis elegans*,  
820 *Science* 293 (2001) 2087–2092. 821
- [28] S.A. Jelinsky, P. Estep, G.M. Church, L.D. Samson, Regulatory  
822 networks revealed by transcriptional profiling of damaged *Saccha-*  
823 *romyces cerevisiae* cells: Rpn4 links base excision repair with  
824 proteasomes, *Mol. Cell. Biol.* 20 (2000) 8157–8167. 825
- [29] J. Quackenbush, Microarrays—Guilt by association, *Science* 302  
826 (2003) 240–241. 827
- [30] E.J. Williams, D.J. Bowles, Coexpression of neighboring genes in  
828 the genome of *Arabidopsis thaliana*, *Genome Res.* 14 (2004)  
829 1060–1067. 830
- [31] B.H. Mechem, et al., Sequence-matched probes produce increased  
831 cross-platform consistency and more reproducible biological results in  
832 microarray-based gene expression measurements, *Nucleic Acids Res.*  
833 32 (2004) e74. 834
- [32] D.T. Ross, et al., Systematic variation in gene expression patterns in  
835 human cancer cell lines, *Nat. Genet.* 24 (2000) 227–235. 836
- [33] M. Nacht, et al., Molecular characteristics of non-small cell lung  
837 cancer, *Proc. Natl. Acad. Sci. USA* 98 (2001) 15203–15208. 838
- [34] D.A. Porter, et al., A SAGE (serial analysis of gene expres-  
839 sion) view of breast tumor progression, *Cancer Res.* 61 (2001)  
840 5697–5702. 841
- [35] H.K. Lee, A.K. Hsu, J. Sajdak, J. Qin, P. Pavlidis, Coexpression  
842 analysis of human genes across many microarray data sets, *Genome*  
843 *Res.* 14 (2004) 1085–1094. 844
- [36] M. Ashburner, et al., Gene Ontology: tool for the unification of  
845 biology. The Gene Ontology Consortium, *Nat. Genet.* 25 (2000)  
846 25–29. 847
- [37] L.J. Jensen, J. Lagarde, C. von Mering, P. Bork, ArrayProspector:  
848 a Web resource of functional associations inferred from microarray  
849 expression data, *Nucleic Acids Res.* 32 (2004) W445–W448. 850
- [38] K. Yeung, M. Medvedovic, R. Bumgarner, From co-expression to co-  
851 regulation: how many microarray experiments do we need? *Genome*  
852 *Biol.* 5 (2004) R48. 853
- [39] P.M. Haverly, L.L. Hsiao, S.R. Gullans, U. Hansen, Z. Weng, Limited  
854 agreement among three global gene expression methods highlights the  
855 requirement for non-global validation, *Bioinformatics* 20 (2004)  
856 3431–3441. 857
- [40] K.D. Pruitt, K.S. Katz, H. Sicotte, D.R. Maglott, Introducing RefSeq  
858 and LocusLink: curated human genome resources at the NCBI, *Trends*  
859 *Genet.* 16 (2000) 44–47. 860
- [41] Mammalian Gene Collection Program Team, et al., Generation  
861 and initial analysis of more than 15,000 full-length human and  
862 mouse cDNA sequences, *Proc. Natl. Acad. Sci. USA* 99 (2002)  
863 16899–16903. 864
- [42] G.S. Davidson, B.N. Wylie, K.W. Boyack, Cluster Stability and the  
865 Use of Noise in Interpretation of Clustering, *IEEE Comput. Soc., Los*  
866 *Alamitos, CA*, 2001, p. 23. 867
- [43] M.J. de Hoon, S. Imoto, J. Nolan, S. Miyano, Open source clustering  
868 software, *Bioinformatics* 20 (2004) 1453–1454. 869
- [44] R. Ihaka, R. Gentleman, A language for data analysis and graphics,  
870 *J. Comput. Graphical Stat.* 5 (1996) 299. 871
- [45] A. Bairoch, et al., The Universal Protein Resource (UniProt), *Nucleic*  
872 *Acids Res.* 33 (2005) D154–D159. 873
- Web site references** 874
- <http://www.ncbi.nlm.nih.gov/geo/>, The Gene Expression Omnibus. 875
- <http://www.r-project.org/>, R Statistical Package Home Page. 876
- <http://www.bioconductor.org/>, Bioconductor Home Page. 877
- 878
- 879
- 880
- 881
- 882
- 883
- 884
- 885
- 886
- 887
- 888
- 889
- 890
- 891
- 892
- 893
- 894
- 895
- 896
- 897
- 898
- 899
- 900
- 901
- 902
- 903
- 904
- 905
- 906
- 907
- 908
- 909
- 910
- 911
- 912
- 913
- 914
- 915
- 916
- 917
- 918
- 919
- 920
- 921
- 922
- 923
- 924
- 925
- 926
- 927
- 928
- 929
- 930
- 931
- 932
- 933
- 934
- 935
- 936
- 937
- 938
- 939
- 940
- 941
- 942
- 943
- 944
- 945
- 946
- 947
- 948
- 949
- 950
- 951
- 952
- 953
- 954
- 955
- 956
- 957
- 958
- 959
- 960
- 961
- 962
- 963
- 964
- 965
- 966
- 967
- 968
- 969
- 970
- 971
- 972
- 973
- 974
- 975
- 976
- 977
- 978
- 979
- 980
- 981
- 982
- 983
- 984
- 985
- 986
- 987
- 988
- 989
- 990
- 991
- 992
- 993
- 994
- 995
- 996
- 997
- 998
- 999
- 1000

- 879 <http://bonsai.ims.u-tokyo.ac.jp/~mdehoon/software/cluster/software.htm>,  
880 The C Clustering Library Home Page. 884
- 881 <http://www.affymetrix.com/>, Affymetrix Home Page. 885
- 882 <http://cmgm.stanford.edu/~kimlab/multiplespecies/Supplement/>, Stuart et  
883 al. Data Home Page. 886
- 884 <http://www.cytoscape.org/>, Cytoscape Home Page. 887
- 885 <http://www.bork.embl.de/ArrayProspector>, ArrayProspector. 888
- 886 <http://microarray.genomecenter.columbia.edu/tmm/>, TMM data. 887
- 887 <http://www.bcgsc.ca/discoveryspace/>, DiscoverySpace Software 888
- 888 References. 888

UNCORRECTED PROOF