

1 **moGSA: integrative single sample gene-set analysis of** 2 **multiple omics data**

3 Chen Meng¹, Bernhard Kuster^{1,2}, Bjoern Peters³, Aedín C Culhane^{4,5*} and Amin Moghaddas Gholami^{1,6*}

4

5 ¹ Chair of Proteomics and Bioanalytics, Technische Universität München, Freising, Germany

6 ² Center for Integrated Protein Science Munich, Freising, Germany

7 ³ La Jolla Institute for Allergy and Immunology, 9420 Athena Circle, La Jolla, CA 92037, USA

8 ⁴ Department of Biostatistics and Computational Biology, Dana-Farber Cancer Institute, Boston, MA
9 02215, USA.

10 ⁵ Department of Biostatistics, Harvard T.H. Chan School of Public Health, Boston, MA 02215, USA

11 ⁶ Current address: La Jolla Institute for Allergy and Immunology, 9420 Athena Circle, La Jolla, CA 92037,
12 USA

13 * Correspondence: aedin@jimmy.harvard.edu; agholami@lji.org

14 **Abstract**

15 **Background:** The increasing availability of multi-omics datasets has created an opportunity to
16 understand how different biological pathways and molecules interact to cause disease. However, there
17 is a lack of analysis methods that can integrate and interpret multiple experimental and molecular data
18 types measured over the same set of samples.

19 **Result:** To address this challenge, we introduce moGSA, a multivariate single sample gene-set analysis
20 method. It uses multivariate latent variable decomposition to discover correlated global variance
21 structure across datasets and calculates an integrated gene set enrichment score using the most
22 informative features in each data type. Integrating multiple diverse sources of data reduces the impact
23 of missing or unreliable information in any single data type, and may increase the power to discover
24 subtle changes in gene-sets. We show that integrative analysis with moGSA outperforms existing single
25 sample GSA methods on simulated data. We apply moGSA to two studies with real data. First, we
26 discover similarities and differences in mRNA, protein and phosphorylation profiles of induced
27 pluripotent and embryonic stem cell lines. Secondly, we report that three molecular subtypes are
28 robustly discovered when copy number variation and mRNA profiling data of 308 bladder cancers from
29 The Cancer Genome Atlas are integrated using moGSA. Our method provides positive or negative gene-
30 set scores (with p-values) of each gene set in each sample. We demonstrate how to assess the influence
31 of each data type or gene to a moGSA gene set score. With moGSA, there is no requirement to filter
32 data to the intersect of features, therefore, all molecular features on all platforms may be included in
33 the analysis.

34 **Conclusion:** moGSA provides a powerful yet simple tool to perform integrated single sample gene-set
35 analysis. Its latent variable approach is fundamentally different to existing single sample GSA
36 approaches. It is an attractive approach for data integration and is particularly suited to integrated
37 cluster or molecular subtype discovery. It is available in the Bioconductor R package “mogsa”.

38 **Keywords**

39 Gene-set analysis, Multivariate analysis, Data integration, Omics, Bladder cancer, molecular subtype
40 stratification

41 Introduction

42 Technological innovations have enabled the acquisition of unprecedented amounts of multi-scale
43 molecular, genotype and phenotype information. Advances in high-throughput sequencing allow
44 quantification of global DNA variation and RNA expression in tissue or blood samples [1, 2]. Mass
45 spectrometry (MS)-based proteomics has undergone rapid progress in recent years, and systematic MS
46 analyses can now identify and quantify the majority of proteins expressed in a human cell line [3]. More
47 and more studies report comprehensive molecular profiling using multiple different experimental
48 approaches on the same set of biological samples. These data can potentially yield insights into the
49 molecular machinery of biological systems. However, integrating, interpreting and generating biological
50 hypothesis from such complex datasets is a considerable challenge.

51 Our groups and others have described multivariate analysis (MVA) approaches that uncover latent
52 correlated structure within and between omics datasets [4-7]. MVA use extensions of principal
53 component analysis (PCA) to project data onto a lower dimensional space so that trends or relationships
54 between multiple datasets, observations (cases) and features (e.g. genes) can be identified. MVA
55 methods identify global correlated patterns among observations, and therefore do not require pre-
56 filtering of gene identifiers in each dataset to a common intersecting subset of features (genes/proteins).
57 All features whether they have annotation or not can be included in the analysis. This is particularly
58 important when analyzing experimental platforms that include novel genes, or use identifiers that are
59 difficult to be mapped. A further attractive feature of latent variable approaches is that supplementary
60 data such as gene-set information (e.g. Gene Ontology annotations) can be projected onto the MVA to
61 aid interpretation [5, 6, 8].

62 Gene-set analysis (GSA) is widely used in the analysis of genome scale data and is often the first step in
63 the biological interpretation of lists of genes or proteins that are differentially expressed between
64 phenotypically distinct groups [9]. These methods use external biological information to reduce
65 thousands of genes or proteins into short lists of functional related gene-sets (e.g. cellular pathways,
66 subcellular localization, transcription factors or miRNA targets), thus facilitating hypothesis generation.
67 The simplest GSA based methods rely on over-representation analysis and only require a list of genes as
68 input. Hypergeometric tests or Fisher's exact test are often used to identify statistically significant
69 overlap between a shortlist of genes or proteins and a database of gene-sets [10]. Gene-set enrichment
70 analysis (GSEA) and significance analysis of function and expression (SAFE) not only require a list of
71 genes, but also take advantage of quantitative information in omics data [11, 12]. More recently,

72 pathway topology approaches also consider the network structure of biological pathways in over-
73 representation analysis [13]. However, these methods are supervised tests that require predefined
74 groups of samples using known experimental, clinical, phenotypic or conditional data (e.g. tumor vs.
75 normal cases).

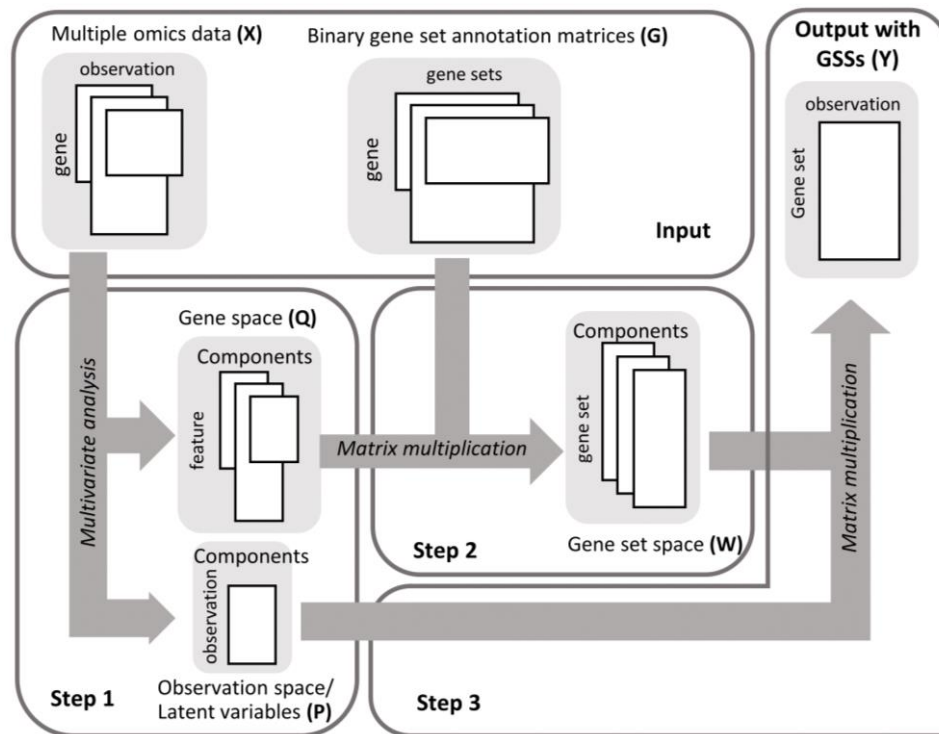
76 Modern omics studies frequently explore a panel of experimental conditions or tissue samples with
77 multiple phenotypes, for example The Cancer Genome Atlas (TCGA), ENCYclopedia of DNA Elements
78 (ENCODE) projects [14] and other studies [15]. Such studies frequently wish to discover new molecular
79 subtypes and thus traditional GSA methods which require known subsets have limited application in
80 such cases. To address this issue, several unsupervised, single sample GSA (ssGSA) methods have been
81 developed [16-19]. These methods do not require prior availability of phenotypic or clinical data. One of
82 the most popular approaches is single-sample GSEA (ssGSEA) that ranks genes according to the empirical
83 cumulative distribution function and calculates a single sample-wise gene-set score by comparing the
84 scores of genes that are inside and outside a gene-set [18]. Another related method described recently,
85 gene-set variation analysis (GSVA), also calculates sample-wise gene set enrichment as a function of the
86 genes that are inside and outside a gene set. GSVA uses a similar Kolmogorov-Smirnov-like rank statistic
87 to assess the enrichment score, but genes are ranked using a kernel estimation of a cumulative density
88 function [16]. Each of these unsupervised single-sample GSA methods are designed for the analysis of a
89 single dataset. To the best of our knowledge no GSA method exists which integrates and calculates a
90 single sample GSA score on multiple datasets simultaneously.

91 Here, we present a novel unsupervised single-sample gene-set analysis that calculates an integrated
92 enrichment score using all of the information in multiple 'omics datasets. We call this approach multiple
93 omics GSA (moGSA). We show that moGSA has higher sensitivity and specificity to detect gene-sets
94 compared to single dataset GSA and demonstrate that moGSA outperforms existing unsupervised GSA
95 methods when applied to simulated data. We apply moGSA to both small and large scale data from
96 multiple omics studies.

97 **Results**

98 moGSA integrates and discovers gene-sets that are enriched in features in two or more omics data
99 matrices obtained on the same set of observations (Figure 1). Omics studies generate multiple data
100 matrices such as RNA sequencing counts of gene expression, measurements of proteins, metabolites,
101 lipids, DNA copy number variations and several other biological molecules that can be mapped to gene-

102 sets. In each, the number of features frequently exceeds the number of observations (rows and columns
103 of the matrix, respectively). In this paper, we refer to genes or other biological molecules as features for
104 simplicity.



105
106 *Figure 1 - Schematic view of the moGSA algorithm. The algorithm requires pairs of matrices as input;*
107 *multiple omics data matrices and corresponding gene-set (GS) annotation matrices. In step 1, the*
108 *multiple matrices are analyzed with a multivariate analysis (MVA) method resulting in an observation*
109 *space and gene space. Next, the gene-set annotation matrices are projected on the same space, and the*
110 *resulting matrix contains the gene-set space. The last step is to reconstruct gene-set-observation through*
111 *multiplying the observation and gene-set spaces.*

112 Figure 1 describes the three steps of the algorithm. Input quantitative or qualitative data matrices must
113 have matched observations but may have different and unmatched features. The number of features
114 may exceed the number of observations. In order to map features to gene-sets, moGSA requires an
115 incidence matrix of gene to gene-set membership associations for each data matrix and in each “gene-
116 set annotation matrix”, a value of 1 indicates that a feature (e.g. gene) is a member of a gene-set. Rows
117 of the gene-set annotation matrix contain the features and each column is an independent annotation
118 vector for a gene-set. A feature may belong to multiple gene-sets simultaneously, that is a row sum may
119 exceed 1.

120 In the first step, several (k) input data matrices are integrated using multiple factor analysis (MFA) [20].
121 MFA is a multiple table extension of principal component analysis (PCA) that is well suited to integrating
122 multiple omics data since it reduces high dimensional omics data to a relatively small number of
123 components that capture the most prominent correlated structure among different datasets [20]. To
124 prevent datasets with more features or different scales to dominate a MFA, each dataset is weighted by
125 dividing it by the first eigenvalue of a decomposition of each individual dataset. MFA generates matrices
126 of latent variables (components) in observation (P) and feature (Q) space. The number of components
127 typically equals the number of observations minus one. We retain and examine the first few
128 components as these represent most of the variance in the data. Approaches for choosing the number
129 of components are discussed later. In the next step (step 2) each gene-set annotation matrix ($G_{1..k}$) is
130 projected as additional information onto the gene-set space ($Q_{1..k}$) generating a score for each gene-set
131 in the same projected space ($W_{1..k}$). In the final step (step 3), moGSA multiplies the latent variables of
132 the observations (P) and latent variables of gene-sets ($W_{1..k}$) to generate a matrix (Y) with a gene-set
133 score (GSS) for each gene-set in each observation (Y).

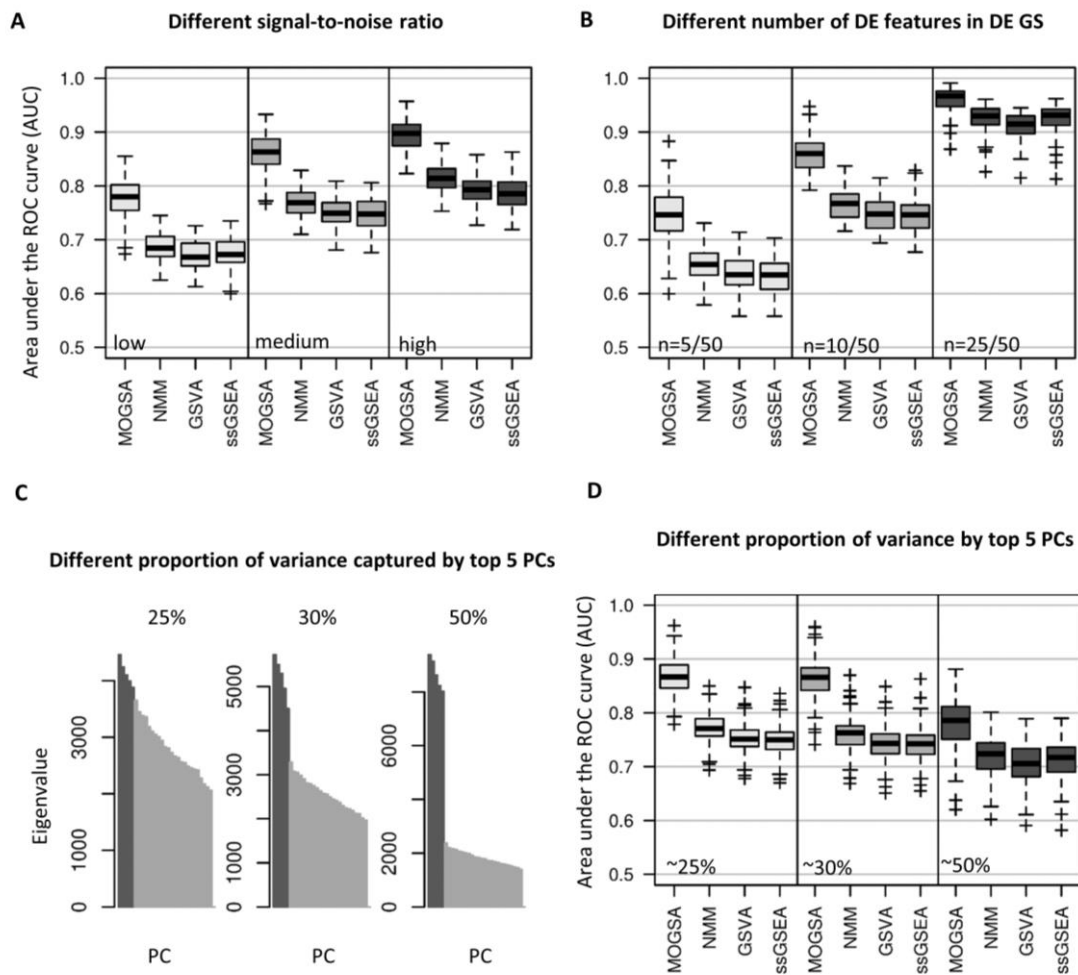
134 A gene-set with a high GSS value has features that explain a large proportion of the global correlated
135 information among data matrices. These features could be from any or all data matrices, and may be
136 non-overlapping, for example a GSS of a gene set with features A-H, could be driven by high levels of
137 gene expression in genes A,B,C, and increased protein levels in proteins C,D,E and amplifications in copy
138 number in gene H. The GSS matrix (Y) may be decomposed with respect to each dataset (X) or latent
139 variable space (P,Q) so that the contribution of each individual dataset or component to the overall
140 score can be evaluated (see Methods).

141 **moGSA outperforms existing single sample GSA methods**

142 Methods to perform integrated ssGSA on multiple 'omics datasets are not yet described. Therefore, we
143 compared the performance of moGSA to ssGSA methods that were developed for analysis of one
144 dataset. One-table ssGSA methods were generally optimized for analysis of gene expression data and
145 include the widely used GSVA and ssGSEA and naïve matrix multiplication (NMM) [16, 18].

146 Figure 2 shows the performance of each method applied to 100 simulated datasets, each run simulated
147 a study of 30 observations with three omics datasets that measured 1,000 features each (Figure S1; see
148 Methods section). Each feature was a member of one of the 20 gene-sets. Each gene-set had 50 genes.
149 The observations were grouped into 6 clusters and each cluster has 5 differentially expressed (DE) gene-
150 sets when compared to the other observations. Within DE gene-sets, 5, 10 and 25 out of 50 genes were

151 randomly simulated to be DE genes (DEG). The triplets were analyzed by moGSA directly, however
 152 matrices were concatenated for NMM, GSVA and ssGSEA as these methods can only accept one matrix
 153 as input.



154

155 *Figure 2 – Comparison of moGSA with NMM, GSVA and ssGSEA. The performance of methods was*
 156 *accessed by their ability to identify differentially expressed gene-sets over 100 simulations in every*
 157 *condition (as indicated by the area under the ROC curve; AUC). (A) Comparison of GSA methods using*
 158 *data with different signal-to-noise ratios. (B) Comparison of data with different number of differentially*
 159 *expressed (DE) genes in each of the DE gene-set. From left to right, 5, 10 and 25 of total 50 genes are*
 160 *differentially expressed in each of the three simulated data matrices if a gene-set is defined as DE gene-*
 161 *sets. (C) Scree plots show representative eigenvalues in each of the conditions in (D). (D) AUCs with*
 162 *different proportion of variance are capture by top 5 components. From left to right, 25%, 30% and 50%*
 163 *of total variance are captured. The darker bars represent the top 5 components.*

164 We anticipated that moGSA might be especially powerful at identifying altered gene-sets in
 165 heterogeneous or noisy data. That is because moGSA, uses only the top few most informative latent

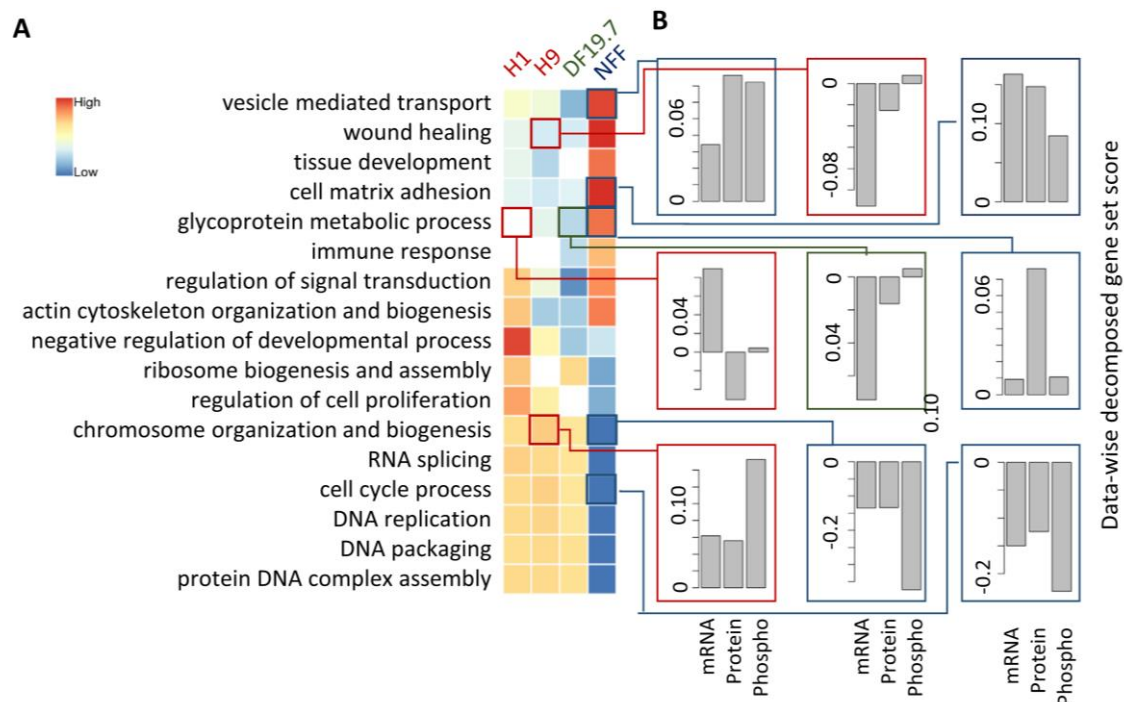
166 variables, thus omitting the signal of many features with little variance, which are potentially noise.
167 Therefore we explored the power of the methods to detect DE gene-sets when there was a strong or
168 weak gene expression signal. First we simulated increasing DEG signal to noise by changing the mean
169 gene expression of DEGs in the cluster and secondly we altered the number of DE genes in a DE gene-set
170 (5, 10 and 25 genes). As expected, the performance of all methods was better when signal-to-noise ratio
171 or the number of DE genes in DE gene-sets increased (Figure 2A and 2B). moGSA consistently
172 outperformed the other methods and the differences were even more apparent when the signal-to-
173 noise ratio was low or when there were few DE genes (5 or 10 of 50 genes) (Figure 2B).

174 Next we compared the performance of each method using data with a simple or complex phenotype. In
175 data with a simple phenotype a few components should easily capture most of the variance in the data.
176 However in data with a complex phenotype for example a heterogeneous tumor dataset, with mixed
177 histology, grade and response to treatment, there are many signals and many latent variables may be
178 required to capture even half of the variance. Specificity and sensitivity of the methods detecting the DE
179 gene-sets (measured as the area under the receiver operating characteristic curve; AUC) were evaluated.
180 In the simulated data, observations were grouped into six clusters, each with highly correlated genes
181 and these six clusters could be captured by the first five components. Therefore we simulated data such
182 that the first 5 components captured 50%, 30% or only 25% of the total variance (Figure 2C). Again,
183 moGSA outperformed the other methods and was relatively robust to changes in the variance retained
184 (Figure 2D). The performance (AUC) of all methods decreased when greater variance was retained,
185 which can be explained by higher intra-cluster correlation that leads to a lower signal-to-noise ratio (see
186 methods).

187 Given the many fundamental differences between moGSA and the other ssGSA methods, we repeated
188 the simulations adjusting for technical aspects of the moGSA approach that might give it an “unfair
189 edge”, but these did little to improve the performance of the others methods. Since, GSVA and ssGSEA
190 were designed for analysis of single datasets, we compared the performance of GSVA and ssGSEA on a
191 single datasets of the triplet compared to the concatenated triplet. Concatenating multiple data
192 matrices neither improved nor decreased the performance compared to analysis of single datasets,
193 most likely because the signal-to-noise ratio increased accordingly with concatenation (Figure S2). In
194 addition, since MFA weights input matrices by their first singular value before moGSA, we examined the
195 effect of data set weighting on the other methods, but found moGSA still outperformed ssGSEA and
196 GSVA when data matrices of the triplet were weighted before concatenation (Figure S3).

197 Application of moGSA to stem cell mRNA and proteomics data

198 We applied moGSA to study a dataset consisting of mRNA, protein and phospho-protein profiling of four
 199 cell lines – two embryonic stem cell lines (ESC; H1 and H9), one induced pluripotent cell line (iPSC;
 200 DF19.7) and a fibroblast cell line (newborn foreskin fibroblast; NFF). Induced pluripotent stem cells (iPSC)
 201 are adult cells that have been reprogrammed to be more like embryonic stem cells (ESC) and have great
 202 potential in the field of regenerative medicine. These cells express ESC markers and can differentiate
 203 into different cell types [21]. Induced pluripotent cells are often derived from NFF cells. The data was
 204 downloaded from [21].



205

206 *Figure 3 – integrative gene-set analysis of iPSC ES 4-plex data. (A) A heatmap shows the gene-set score*
 207 *(GSS) for significantly regulated gene-sets in the cell lines, the white colored blocks/cells indicates the*
 208 *change of gene-sets are non-significant. (B) Data-wise decomposition of the GSS for some of the gene-*
 209 *sets. The contribution of each of the data is represent by a bar. The Y-axis is the data-wise decomposed*
 210 *gene-set score.*

211 After filtering low abundant features, there were 10,961; 5,817; and 7,912 unique mRNAs, proteins and
 212 phosphorylation sites features respectively (see Methods). Principal component analysis (PCA) of each
 213 individual dataset is shown in Figure S4. The strongest signal (first PCs) in all three datasets was the
 214 difference between NFF cells and the stem cell lines, and this difference was particularly apparent in the
 215 proteomics datasets. The second and third components represented subtle differences between iPSC

216 and ESC lines, thus we retained the top 3 components when we applied MFA to transform all of the data
217 onto the same space and scale. The three datasets contributed similarly to the overall variance in the
218 integrated analysis, as indicated by weighting of each dataset in MFA. The first eigenvalues (square of
219 singular values) of each PCA were 0.24, 0.26 and 0.26 for the transcriptome, proteome and phospho-
220 proteome dataset respectively. MFA recapitulated the PCA of the individual datasets. Most of the
221 variance was captured in the first component and it discriminated between NFF and other cell lines. The
222 variance of the molecular differences between the ESC cells (captured on the second component) was
223 greater than the difference between ESC and iPSC cell lines (component 3) (Figure S5).

224 moGSA was used to annotate the features with gene ontology (GO) biological processes. There were
225 228 GO terms (out of 825) that had significant up or down-regulated gene-set scores (GSSs) in at least
226 one cell line (BH corrected p value < 0.01). There was gene overlap among many GO terms and
227 hierarchical clustering analysis (Hamming distance and complete linkage) was used to group the 288 GO
228 terms into 21 broad categories (Table S1). Gene-set scores of representative GO terms from each
229 category are shown in Figure 3A. Biological processes associated with more differentiated cell types
230 were associated with the NFF cells and included up-regulation of vesicle-mediated transport, immune
231 related responses and cell adhesion. In contrast cell proliferation GO terms such DNA replication, and
232 cell cycle processes had significantly higher GGS in the highly proliferative stem cell lines. These results
233 confirm previous findings [21].

234 In integrative analysis of multiple omics data, it is important to evaluate the relative contribution (either
235 concordant or discrepant) of each dataset to the overall GSS. Data-wise decomposition of the GSSs (see
236 Methods) are shown in Figure 3B. The three data sets have concordant contributions to most of the GO
237 terms, including vesicle mediate transport, cell matrix adhesion, cell cycle processes in NFF line;
238 chromosome organization and biogenesis in H9 and NFF cell lines.

239 However, in other GO classes, we also observed differences in the contribution of mRNA, proteins and
240 phosphor-protein data to the GSS. Chromosome organization and biogenesis had significant positive GSS
241 in the stem cells and significant negative GSS in the NFF cells, and was driven by differences in the
242 phosphorylation data. Another case where the mRNA and protein data were incongruent was the GO
243 class “glycoprotein metabolic process”. It had GSS scores of 9.7 ($p < 0.001$), -8.6 ($p < 0.01$), -5.3 ($p < 0.01$)
244 and 0 ($p > 0.05$) in NFF, iPSC, H9 and H1 cells respectively. Up-regulation in NFF mainly reflects up-
245 regulation on the protein level. However, down-regulation in iPSC DF19.7 cells is due to low expression
246 of related mRNAs. The GO term wound healing has previously been shown to be differentially

247 upregulated in fibroblast NFF cells compared to ESC [21]. Consistently, we also found wound healing was
248 upregulated in NFF compared to ESC; the GSS for wound healing were 14.2 ($p < 0.01$), -5.4 ($p < 0.01$), -5.2
249 ($p < 0.01$) and -3.6 ($p < 0.001$) for NFF, iPSC, H9 and H1 cells respectively (Table S1). Down-regulation of
250 wound healing in H9 cell line was dominated by mRNA data, and the two proteomics datasets
251 contributed little to the negative GSS. In contrast to previous studies [21], we did not observe significant
252 differences in wound healing between iPSC and ESC. This difference could be because moGSA is more
253 sensitive (than single data GSA) in detecting gene-sets that have subtle but consistent changes in
254 multiple datasets. More importantly, the contribution of individual gene-set could be evaluated by the
255 decomposition of GSS with respect to datasets

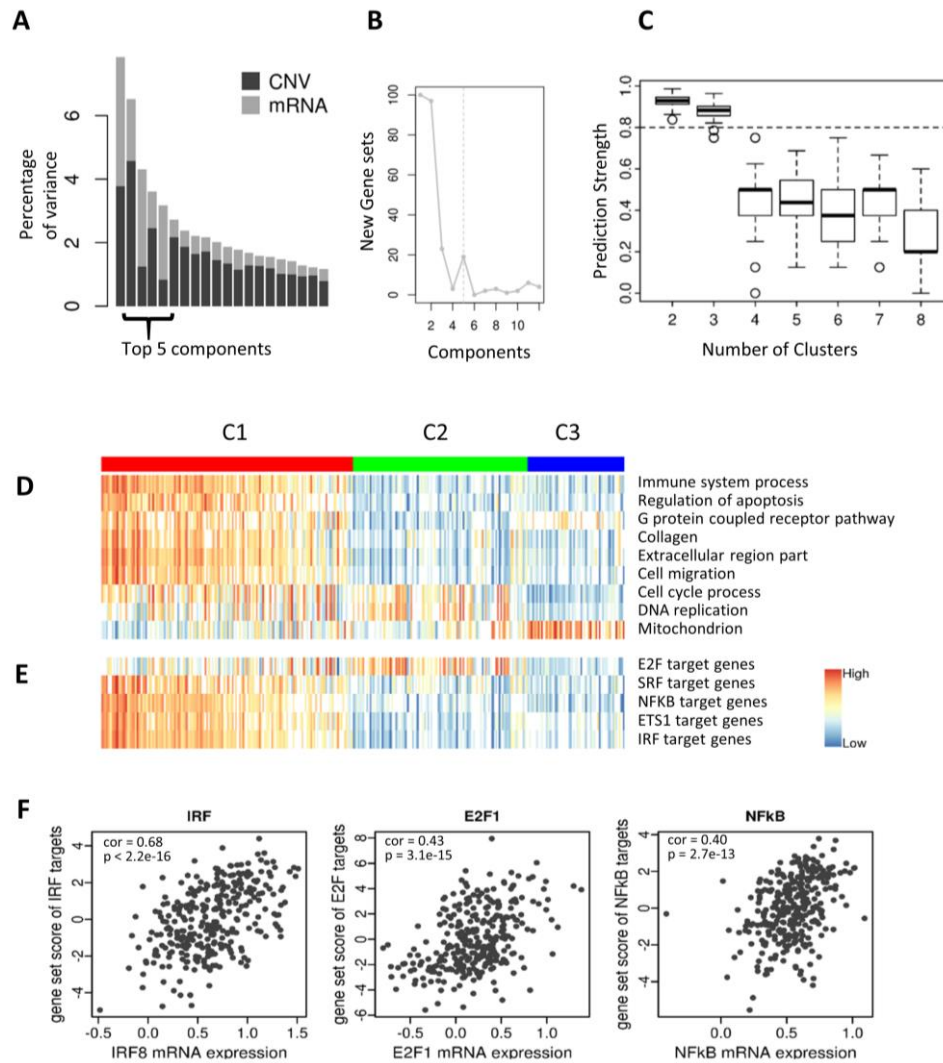
256 **Application of moGSA to TCGA Bladder cancer data analysis**

257 Since moGSA performs unsupervised integrative single sample GSA, it is particularly useful approach for
258 cluster discovery in multi 'omics data. Therefore we applied moGSA to extract an integrative subtype
259 model of BLCA from copy number variation (CNV) and mRNA data of 308 muscle invasive urothelial
260 bladder cancer (BLCA) patients (obtained as part of the TCGA project).

261 BLCA is a molecularly heterogeneous cancer with between 2 and 5 molecular subtypes (reviewed by
262 [38]). Briefly, Sjödaahl et al. first defined five major subtypes termed urobasal A (UroA), UroB,
263 genomically unstable (GU), squamous cell carcinoma-like (SCCL) and 'infiltrated' [22]. The TCGA study
264 defined four expression clusters (I–IV) [23]. The two subtype model consists of basal-like and luminal
265 subtypes [24] which was extended by Choi et al. who defined a 'p53-like' luminal subtype apart from
266 basal-like and luminal subtypes [25].

267 Data were downloaded from the TCGA website and after filtering out features with low variance (see
268 Methods), CNV and RNA-seq mRNA expression data contained 12,447 and 14,710 genes respectively, in
269 which 7,644 genes were common to both datasets (Figure S4). Filtering of features is not required by
270 moGSA but we filter low quality features as they are unlikely to contribute to the analysis. PCA of each
271 individual dataset is shown in Figure S7. From scree plots of the first 10 eigenvalues, an elbow in each
272 plot appears between 4-6 components suggesting this number of components are needed to capture
273 most of the variance (Figure S7), which we anticipated given the known molecular heterogeneity in
274 these data. The first eigenvalue (square of singular value) of the PCA of BLCA mRNA and CNV data are
275 0.0004 and 0.0003 respectively. We applied a preliminary MFA on the data and Figure 4A shows the
276 eigenvalues of the resulting components. The top five components captured a quarter of the total

277 variance and were not dominated by either CNV or mRNA (CNV 50.6%, mRNA 49.4%). Also, these five
 278 components were not correlated with batches (TCGA batch ID), plates, shipping date or tissue source
 279 sites.



280

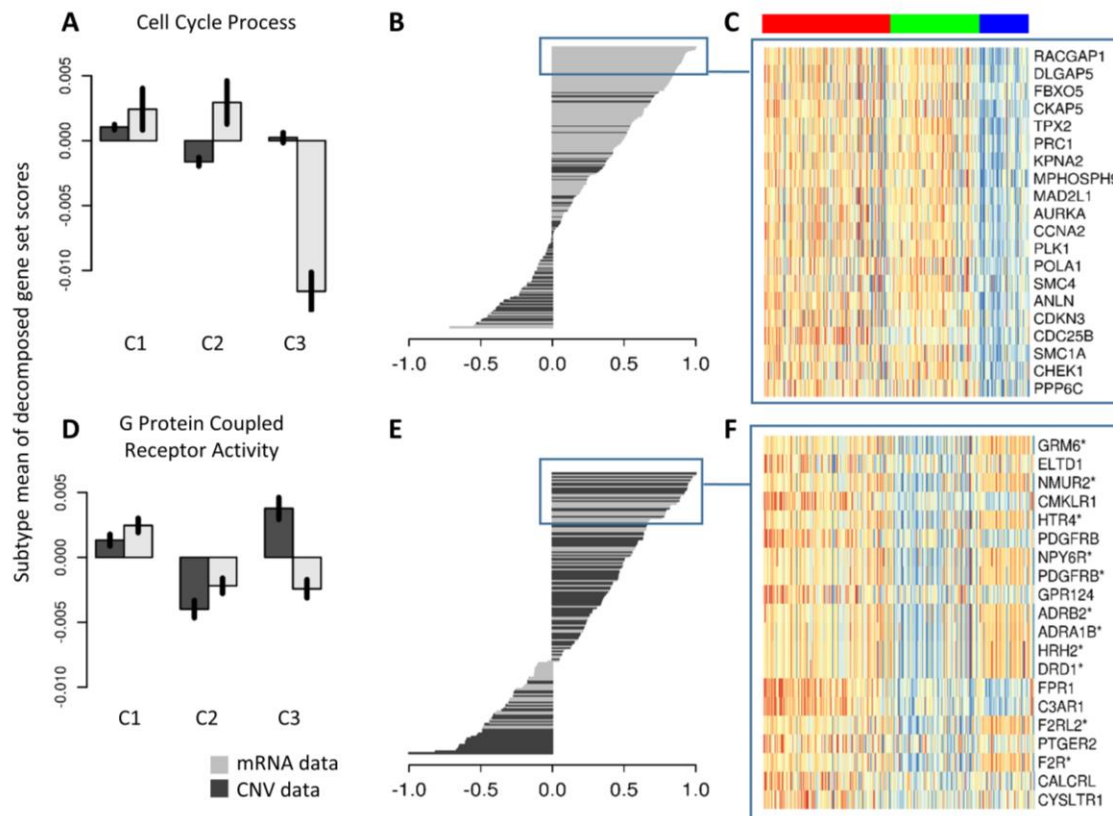
281 *Figure 4 – Data integration with moGSA and integrative subtype defined by latent variables. (A) Bar plot*
 282 *showing the eigenvalues of components defined by MFA. The top 5 components were selected in the*
 283 *analysis. (B) Effect of including additional component (1-12) on the identification of new genesets among*
 284 *the top 100 genesets (C) Prediction strength was used to evaluate the robustness of classification into*
 285 *two to eight subtypes. The boxplot shows the prediction strength of 100 randomizations. Two and Three*
 286 *are relative robust subtype models (prediction strength > 0.8). (D) Gene ontology (GO) and*
 287 *transcriptional target (TFT) gene-sets annotation of tumors. Heatmap showing the GSSs for selected*
 288 *gene-sets. The gene-sets “immune-related, apoptosis, G protein receptor, collagen, extracellular region*
 289 *and cell migration” are strong in the C1 (basal-like) subtype, whereas the mitochondrial related gene-*
 290 *sets are over represented in the C3 (luminal A-like) subtype of tumors. (E) The most significant*

291 *transcriptional factor (TF) target gene-sets. The gene-set scores suggest that 4 out of the 5 TFs are*
292 *hyperactive in the C1 subtype, except E2F family is active in the C2 subtype of cancer. The white spaces in*
293 *(A) and (B) denote non-significant GSSs. (F) The scatter plots display the correlation between gene-set*
294 *scores and the mRNA level of selected TFs. The expression of selected TFs is significantly correlated with*
295 *their gene-set scores (also see Figure S16).*

296 In a typical analysis, we use a scree plot to select the number of components. The scree plot indicated
297 that five components should capture sufficient variance for input to moGSA. We confirmed that this was
298 the optimal number of components as input to moGSA, in the following experiment. We performed
299 moGSA on the BLCA mRNA gene expression and CNV data (n=308) with a number of components ranged
300 from 1 to 12. For each gene-set in the GSS matrix, gene-sets were ranked by the number of tumors in
301 which they were significantly regulated (either positive or negative GSS, $p < 0.05$), such that gene-sets
302 that were significant in most tumors had highest rank. The distribution of the number of tumors in
303 which gene-sets were significant at $p < 0.05$, $p < 0.01$, and $p < 0.001$ is shown Figure S7. No gene-set was
304 significant in all 308 tumors and most gene-sets were insignificant in all tumors (Figure S7). For $p < 0.05$,
305 we examined the 10, 20, 40, 100, 200, 500 and 1000 highest ranked gene-sets and examined the
306 stability of gene-set ranking when additional components were included (Figure S8). Increasing the
307 number of components (from 1 to 5) increased the stability of gene set lists, however there was little
308 additional gain after five components (Figure S9). Among the top 100 ranked gene-sets, few new gene-
309 sets were identified after five components (Figure 4B).

310 Therefore we used moGSA to perform single sample GSA analysis with 1,125 gene-sets on an MFA of the
311 mRNA and CNV BLCA data in which five components were retained. The number of significant gene-sets
312 per patient ($p < 0.05$) ranged from 183 to 595 and these contained both gene-sets with positive and
313 negative GSS. To identify the number of BLCA molecular subtypes, we performed consensus clustering
314 on the five components, which resulted in a three-subtype model (Figure 4B and Figure S10-13). We
315 performed several experiments, to confirm that three subtypes was optimal particularly since between
316 2 and 5 subtypes have been previously reported in BLCA [23]. Whilst consensus clustering analysis
317 indicated high confidence in either two or three subtypes (Figure S10B-D), silhouette analysis (Figure
318 S10E) suggested three subtypes. Stability analysis showed there was no effect when different
319 resampling proportions (50%, 60%, 70%, 80% and 90%) were used in the inner and outer loop of
320 consensus clustering (Figure S11). A recent report highlighted limitations in consensus clustering [26],
321 and therefore in parallel, we also used the “prediction strength” algorithm, to discover the number of
322 stable subtypes that can be predicted from the data [27] (see Methods). Data were divided into training

323 and test, and a KNN classifier was used to iteratively predict the class of each patient. Though no good
 324 choice of K existed (Figure S12), this had minimal influence on the final result, which clearly supported
 325 three subtypes (Figure S13). Therefore using two independent approaches, we determined that the data
 326 (5 components of the integrated analysis) supported three BLCA molecular subtypes.



327

328 *Figure 5 – CNV and mRNA data contribute unequally to defining subtype and gene-set scores. (A) Data-*
 329 *wise decomposition of gene-set scores for “cell cycle process”. The bar plot shows the normalized mean*
 330 *of data-wise decomposed GSSs in each subtype (the black vertical line on the bars show the 95%*
 331 *confidence interval of the mean). (B) The bar plot shows the gene influential scores (GISs) of genes in the*
 332 *“cell cycle process” gene-sets. The expression of the top 30 most influential genes in the gene-set are*
 333 *shown in (C). (D-F) Same as (A-C) for “G protein couple receptor activity”. Gene names in (F) with*
 334 *asterisks indicate genes from CNV data.*

335 The three BLCA subtypes identified in our integrative analysis overlapped with the BLCA subtypes
 336 identified in previous studies (Table S2, Figure S14). Our integrative BLCA subtypes consisted of two
 337 larger subtypes C1, C2 containing 148 and 103 patients respectively, and a smaller group C3 with 57
 338 patients. The smaller subtype, C3, was the most robust (Figure S10E, S11). The integrative subtype C1
 339 harbored a high number of patients in the type III and IV of the TCGA subtypes, the infiltrated and SCCL

340 subtypes of the Sjö Dahl study [22] and the basal-like subtype identified by Damrauer (BH corrected p-
341 value < 0.05, Table S2) [24]. Subtypes C2 and C3 were more similar to the Damrauer luminal subtype.
342 But, the C3 subtype contained more low grade tumors and showed a strong overlap with the UroA
343 subtype of the Sjö Dahl study and type I of the TCGA subtype model. Subtype C2 tumors overlapped with
344 the genomically unstable subtype defined by Sjö Dahl (Table S2). Accordingly, we observed higher
345 mutation rate in the C2 patients (Figure S15). In single sample gene-set analysis with moGSA, C1 patients
346 had more significant GSS ($p < 0.05$) than C2 or C3 (Figure S16).

347 To further characterize BLCA, we focused on gene-sets that were differentially regulated in most
348 patients. There were 73 gene-sets that were significantly regulated (positive or negative GSS, p
349 value < 0.05) in 200 or more of the 308 patients (Table S3 and Figure S17). Alternatively a lower cutoff
350 would include more gene-sets that are regulated fewer tumors, fewer gene-set could be selected using
351 a lower p-value ($p < 0.01$, 0.001) or a supervised analysis could be used to select GSS that most
352 discriminate groups of tumors. Cluster analysis of the GSS matrix (73 selected gene-sets x 308 tumors)
353 revealed 3 clusters of gene-sets. A large cluster of 51 gene-sets had positive GSS scores in C1 but
354 negative scores in C2 or C3. Two smaller clusters of gene-sets of 16 and 6 gene-sets had positive GSS
355 scores in C2 and C3 respectively (Figure S17).

356 The large C1 gene-sets cluster was dominated by 31 gene-sets with terms associated with “immune
357 response” which had significant strongly positive GSS in the C1 basal-like/SCC-like BLCA subtypes.
358 Associations between immune regulation and the basal-like cluster have been previously reported [22].
359 The remaining 20 gene-sets in the C1 cluster of gene-sets included terms associated with “extracellular”,
360 function, cell morphogenesis, migration and muscle cell development, “apoptosis” (2 gene-sets), and “G
361 protein coupled receptor” (6 gene-sets) (Figure S17, S18) and EMT related gene sets (Figure S19), which
362 recent reports that the Basal-like subtype tend to have more muscle-invasive and metastatic disease at
363 presentation [22]. The remaining gene-sets could broadly be defined by biological processes of “cell
364 cycle” (9 gene-sets) and “DNA repair and chromosome related” (7 gene-sets) which had high GSS in C2
365 (and some C1) and “mitochondrion” (4 gene-sets) in C3. A heatmap of the GSSs of representative gene-
366 set of each category is shown in Figure 4C and S17. We found that most of these gene-sets have been
367 associated with subtype of bladder cancer. Increased cell-cycle and DNA repair GSS were associated
368 with the “genomically unstable” luminal C2 cluster [28] (Figure S14, S16). The mitochondrial component
369 has been described in bladder cancer and other cancers previously [28, 29], our study particularly
370 associated this function with C3 low-grade papillary-like subtype in BLCA. However other gene-sets may

371 be associated with C3 that were excluded when GSS were filtered to those that were broadly significant
372 in 200 or more patients.

373 The GSSs clearly distinguished the three BLCA molecular subtypes. The most significant gene-sets,
374 “immune response” and “immune system process” have significant positive or negative GSS in 270 and
375 265 of 308 patients respectively (Table S3). The median GSS for the gene-set “immune system process”
376 was 0.82, -0.75, -0.61 in C1, C2 and C3 respectively (Figure S17, S18) indicating that immune related
377 processes have high gene expression or CNV in the C1 subtype and much lower in C2 and C3. Next, we
378 determined the importance of individual genes in each gene-set by calculating a gene influential score
379 (GIS) using a leave-one-out procedure (see methods). The maximum GIS value for a gene in a gene-set is
380 1, which indicates that gene contributes a high proportion of variance to the overall variance of the GSSs.
381 A GIS close to 1 often suggests a high correlation between the gene expression value and GSS. Gene
382 influential score of the gene-set immune system process in BLCA suggested that the top ranked genes
383 included *ITGB2*, *SPI1*, *DOCK2*, *LILRB2* and *LAT2*. Other highly ranked genes included drug target genes
384 such as *CD4*, *IL6*, the interferon induced proteins *IFITM2* and *IFITM3* and the G protein coupled
385 receptors *GPR183* and *CMKLR1* (Table S4). Top positive influencers in “regulation of apoptosis” were
386 also related to the immune response, such as *STK17A*, *ANXA5* and *BCL2A1*, *STAT1*, Serpin B, *TGFB* and
387 *ANXA1* (Table S4). Moreover, several epithelial to mesenchymal transition (EMT) related gene-sets, such
388 as “collagen” (including *COL6A3*, *COL1A1*, *COL5A1* and *COL3A1*), “extracellular matrix proteins” (e.g.
389 glycoproteins *SRGN* and *FBN1*) and mesenchymal gene-sets were elevated in C1 (Figure S19; Table S4).

390 The C3 subtype tumors had higher GSSs in mitochondrial related gene-set and lower expression of genes
391 related to cell cycle process and DNA replication. GIS analysis suggested that two families of genes,
392 NADH dehydrogenases (NDUFs) and mitochondrial ribosomal proteins (*ABCC1/MRP*) influenced the
393 mitochondrial proteins (Table S4).

394 To identify transcription factors (TF) that may regulate gene expression in the three tumor subtypes, we
395 used transcriptional factor target (TFT) gene-sets to annotate the tumors. Similar to the selection of GO
396 terms, we focused on TFT gene-sets with more than 200 significant GSSs across 308 patients (Table S2).
397 The GSSs of the E2F family target gene-set were significantly different in most of the tumors and are
398 particularly low for the C3 tumors. The rest of the four identified TFs were highly elevated in the C1
399 subtype. Among them, we identified an *MADS* (*MCM1*, Agamous, Deficiens, and *SRF*) box superfamily
400 member, *SRF* and several TFs associated with transactivation of cytokine and chemokine genes,
401 including *NFkB1*, *ETS1* and *IRF1* (Figure 4D). The genes exhibiting the largest GIS in the *IRF1* and *NFkB1*

402 target gene-sets include *ACTN1*, *CXorf21*, *ICAM1*, *MSN*, *TNFSF13B*, *IL12RB1* and *CDK6* (Table S5). Further,
403 we examined the correlations between GSSs and the mRNA expression. All five TFs showed that the TF
404 mRNA and GSSs are significantly correlated (Figure 4E, Figure S20). The boxplot of GSS with respect to
405 subtypes in Figure 4C and D are shown in Figure S18,S21.

406 In order to identify the contribution of each dataset, we decomposed the GSSs with respect to the
407 datasets or components. Figure 5A shows the means of data-wise decomposed GSSs in each subtype for
408 “cell cycle process”, where we found that mRNA expression strongly influenced the GSS, particularly the
409 low GSS of the C3 subtype patients. The gene influential score (GIS) analysis supports this finding as the
410 top 30 most influential genes are all based on mRNA expression (Figure 5B), including *RACGAP1*, *DLGAP5*,
411 *FBXO5*, *AURKA*, *KERA* (*CNA2*) and *CDKN3* (Figure 5C). By contrast, both CNV and mRNA data influenced
412 the gene-set “G protein coupled receptor activity” (Figure 5D) and the GIS analysis shows that the most
413 influential genes include those from both mRNA and CNV data (Figure 5E). However, the CNV and mRNA
414 expression patterns in the C3 subtype shows a clear difference for this gene-set (Figure 5F). Top gene
415 influencers of “G protein couple receptor activity” included CNV of *GRM6*, *NMUR2*, *PDGFRB* and
416 adrenergic receptors, the gene expression of *ADGRL4* (*ELTD1*), *CMKLR1* and *PDGFRB* (Figure 5F). In
417 addition, the data-wise decomposition of GSS identified several GSSs that were only contributed by the
418 mRNA data, including the immune system process, DNA replication and mitochondrion gene-set (Figure
419 S21).

420 Discussion

421 In this paper, we introduced a new multivariate single sample gene-set analysis approach, moGSA that
422 enables discovery of biological pathways with correlated profiles across multiple complex datasets.
423 moGSA uses multivariate latent variable analysis to explore correlated global variance structure across
424 datasets and then extracts the set of gene-sets or pathways with highest variance and most strongly
425 associated with this correlated structure across observations. By combining multiple data types, we can
426 compensate for missing or unreliable information in any single data type so we may find gene-sets that
427 cannot be detected by single omics data analysis alone [4].

428 moGSA uses the maximum variance of the concordant structure across of datasets to calculate the
429 gene-set scores for each observation. This is fundamentally different from other gene-set enrichment
430 analysis methods which use a ‘within observation summarization’ such as the mean or median of gene
431 expression of genes in a gene-set. It has several characteristics that make it attractive for data

432 integration. First moGSA uses MFA, a multi-table extension of PCA to reduce the complexity of the
433 original data by transforming high dimensional data to a small number of components (latent variables).
434 The components with highest eigenvalues (largest variance) capture the most prominent structure
435 among the different datasets. Excluding components with low variance may strengthen the signal-to-
436 noise ratio of data, as it reduces low variant, noise or artifact variance [30, 31]. In moGSA, the entire set
437 of features from each platform is decomposed onto a lower dimension space. The linear combination of
438 feature loadings is used in the calculation of the gene-set scores. Features that contribute low variance
439 contribute little to the score and thus the dimension reduction within moGSA comes with an intrinsic
440 filtering of noise. The advantages of intrinsic variance filtering of features can be clearly seen when we
441 applied moGSA to simulated data. moGSA outperformed ssGSA approaches including ssGSEA and GSVA
442 which do not include a noise-filtering component. Second, data integration of features is achieved at the
443 gene-sets level rather than scoring individual features. This greatly facilitates the biological
444 interpretation among multiple integrated datasets. There is no requirement to pre-filter features in a
445 study or map features from different datasets to a set of common genes. Therefore, moGSA can be used
446 to compare technological platforms that have different or missing features.

447 There is great potential for applying multi-table unsupervised GSA approaches for discovery of new
448 subtypes and pathways in integrated data analysis of complex diseases such as cancer. In this study, we
449 applied moGSA in combination with clustering analysis. Dimension reduction approaches such as moGSA
450 and MFA are well suited to cluster discovery data because these approaches consider the global
451 variance in the data and as such are complementary to hierarchical or k-means clustering approaches
452 which focus on the pair-wise distance between observations [31-33].

453 The number of components is an important input parameter to consider when applying moGSA to gene-
454 set analysis or cluster discovery. Similar to PCA, the optimal number of MFA components may be
455 assessed by examining the variance associated with each component. The first component will capture
456 most variance and the variance associated with subsequent component decreases monotonically. Scree
457 plots (Figure 2C, 4A) may be used to visualize if there is an elbow point in the eigenvalues, allowing one
458 to select the components before the elbow point. Alternatively one may select the number of
459 components that capture a certain proportion of variance (50%, 70%, etc). In addition, one may include
460 components that are of biological interest. For example, in the iPS ES example, there is a clear biological
461 meaning in the third component (ES vs iPS cell line). In analysis of the BLCA data, we examined a range
462 of components (1-12), and show that there is little gain of information once a minimum number of

463 components with high variance are included (Figure 4B). In addition, the variance of retained
464 components should not be dominated by one or a few datasets. To facilitate biological interpretation of
465 components, the GSS could be decomposed with regard to components. In the BLCA example, the
466 second and fourth component are largely contributed by CNV, whereas mRNA is more important in
467 defining the third and fifth components. Including five components ensured that both datasets
468 contributed relatively similar variance to the global variance.

469 An issue might arise with latent variables analysis if components with the large variance capture
470 information unrelated to biological variance [30], such as technical artifacts or batch effects. In practice
471 this is rare in MFA, because it focuses on components that capture global correlation among all datasets.
472 Often batch effects are specific to a platform and thus a component that captures information that is
473 entirely uncorrelated to the global structure will be omitted from the set of highly variant integrated
474 components. However it is still wise to perform careful batch effect control, especially in the large scale
475 omics studies. A more detailed description of batch effect detection is described in [34].

476 Another consideration when applying moGSA, is that it is most efficient in detecting gene-sets that have
477 broad correlation patterns among data types. It may fail to discover gene-sets with few genes,
478 particularly if they had low variances on the selected components.

479

480

481 **Methods**

482 **moGSA algorithm**

483 ***Input data and gene-set annotation matrix***

484 The inputs to moGSA are pairs of multiple matrices ($\mathbf{X}_k, \mathbf{G}_k$). \mathbf{X}_k is a set of matrices, denoted $\mathbf{X}_1, \dots, \mathbf{X}_k, \dots$
485 \mathbf{X}_k , where K is the total number of quantitative matrices. Matrix \mathbf{X}_k is a $p_k \times n$ matrix of quantitative omic
486 data, which contains p_k rows of features (e.g. genes) measured over the same n observations. Each of
487 the matrices $\mathbf{X}_1, \dots, \mathbf{X}_k$ has a corresponding gene-set annotation matrix, $\mathbf{G}_1, \dots, \mathbf{G}_k, \dots, \mathbf{G}_K$. The gene-set
488 annotation matrix \mathbf{G}_k is a $p_k \times m$ binary incidence matrix of gene to gene-set membership associations,
489 where m is the number of gene-sets. The element $g_{k[i,j]}$ in \mathbf{G}_k has the value 1 if the i th feature is a
490 member of the gene-set j and 0 otherwise. \mathbf{G}_k is constructed using predefined gene-set information such
491 as the Gene Ontology [35, 36] GeneSigDb [37] or MSigDB [38].

492

493 ***moGSA step 1 multivariate integration***

494 The first step of the moGSA involves data integration with a multiple table multivariate analysis method.
495 In this study, we use MFA because of its simplicity and computational efficiency. MFA can be viewed as a
496 generalization of principal component analysis (PCA) for a multi-table problem [20]. We briefly describe
497 MFA using the nomenclature of Abdi et al. 2013 [20].

498 When integrating multiple data matrices, one must decide if all datasets should have equal weights, or if
499 some data are “more important”, for example those with higher quality, fewer features, higher variance,
500 etc. Simple tensor decomposition approaches, or PCA on a concatenated matrix, give every dataset
501 equal weights and results are often dominated by the matrix (or matrices) with the large variance or
502 most features. To correct for this, MFA weights datasets by dividing each by their first eigenvalue. The
503 weight of each matrix is expressed as

$$\alpha_k = \frac{1}{\lambda_{k,1}^2} \quad (1)$$

504 Where $\lambda_{k,1}^2$ is the first singular value of data matrix \mathbf{X}_k . For convenience, the weights of matrices are
505 stored in a diagonal matrix \mathbf{A} , whose diagonal elements are

$$\text{diag}\{\mathbf{A}\} = [\text{diag}\{\mathbf{A}_1\}, \dots, \text{diag}\{\mathbf{A}_k\}, \dots, \text{diag}\{\mathbf{A}_K\}] = [\alpha_1 \mathbf{1}_1^T, \dots, \alpha_k \mathbf{1}_k^T, \dots, \alpha_K \mathbf{1}_K^T] \quad (2)$$

506 The transpose of a matrix is denoted by superscript \top . $\mathbf{1}_k^\top$ is a vector of 1 in the length of p_k . As a result, \mathbf{A}
 507 is a $p \times p$ diagonal matrix, the diagonal elements of \mathbf{A} representing the weight of features in $\mathbf{X}_1, \dots, \mathbf{X}_k$.
 508 Similarly, the weight of each observation is an $n \times n$ diagonal matrix, \mathbf{M} . In the present study, we use
 509 $m_{ii}=1/n$, namely, all observations are equally weighted.

510 We then transpose and concatenate all \mathbf{X}_k to a complete $p \times n$ matrix ($p = \sum_k p_k$):

$$\mathbf{X} = [\mathbf{X}_1^\top | \dots | \mathbf{X}_k^\top | \dots | \mathbf{X}_K^\top]^\top \quad (3)$$

511 After deriving the matrix weights, observation weights and the concatenated matrix, MFA is reduced to
 512 an analysis of the triplet $(\mathbf{X}, \mathbf{A}, \mathbf{M})$. The solution of the problem is given by generalized singular value
 513 decomposition (GSVD):

$$\mathbf{X}^\top = \mathbf{P}\mathbf{\Delta}\mathbf{Q}^\top \text{ with the constraint that } \mathbf{P}^\top\mathbf{M}\mathbf{P} = \mathbf{Q}^\top\mathbf{A}\mathbf{Q} = \mathbf{I} \quad (4)$$

514 \mathbf{X} is transpose so that \mathbf{P} is a $n \times r$ matrix, \mathbf{Q} is a $p \times r$ matrix, $\mathbf{\Delta}$ is an $r \times r$ square matrix, the maximum
 515 number of r is the rank of \mathbf{X} . The components of MFA, \mathbf{F} , are given by

$$\mathbf{F} = \mathbf{P}\mathbf{\Delta} \quad (5)$$

516 where \mathbf{F} has the same dimension as \mathbf{P} . In the PCA framework, the matrix \mathbf{P} contains the PCs or latent
 517 variables. We also call it *sample space* in this paper. The column vectors in \mathbf{P} may be plotted on a two
 518 dimensional space to visualize the contribution of each observation to the variance captured by each PC.
 519 The matrix \mathbf{Q} is the loading matrix or *gene space*. Because \mathbf{X} is a concatenation of multiple matrices, the
 520 gene space matrices $\mathbf{Q}_1, \dots, \mathbf{Q}_k$ may also be concatenated or partitioned in the same manner, namely,

$$\mathbf{Q} = [\mathbf{Q}_1^\top | \dots | \mathbf{Q}_k^\top | \dots | \mathbf{Q}_K^\top]^\top \quad (6)$$

521
 522 ***moGSA step 2 project gene-set annotation matrix as supplementary data***
 523 Different gene-sets have different candidate genes, therefore, in order to facilitate the comparison of
 524 gene-set score across gene-sets, we normalized the gene-set annotation matrix so that the sum of each
 525 column in \mathbf{G} equals 1, that is,

$$\hat{g}_{[i,j]} = \frac{\hat{g}_{[i,j]}}{\sum_i \hat{g}_{[i,j]}} \quad (7)$$

526 where $\hat{g}_{[i,j]}$ is the elements on the i th row and j th column in the normalized gene-set annotation matrix
 527 $\hat{\mathbf{G}}$. The gene-set score calculated using un-normalized gene-set annotation matrix for gene-sets in
 528 Figure 3 and 4 are shown in Figures S22 and S23.

529 Next, we project the annotation matrix as supplementary data [35] to generate the gene-set space
 530 matrix $\mathbf{W}_k (m \times r)$, which is calculated as a product of the normalized gene annotation matrix and loading
 531 matrix.

$$\mathbf{W} = \hat{\mathbf{G}}^T \mathbf{A} \mathbf{Q} \text{ where } \hat{\mathbf{G}} = [\hat{\mathbf{G}}_1^T | \dots | \hat{\mathbf{G}}_k^T | \dots | \hat{\mathbf{G}}_K^T]^T \quad (8)$$

532 $\hat{\mathbf{G}}$ is the grand annotation matrix with dimension $p \times m$. The overall gene-set space \mathbf{W} ($m \times r$ matrix) could
 533 also be expressed as the sum of individual $\hat{\mathbf{G}}_k$ and \mathbf{Q}_k , that is,

$$\mathbf{W} = \sum_{k=1}^K \mathbf{W}_k \text{ where } \mathbf{W}_k = \hat{\mathbf{G}}_k^T \mathbf{A}_k \mathbf{Q}_k \quad (9)$$

534

535 **moGSA step 3 reconstruction of gene-set-observation matrix**

536 The main output of MOGSA is a *gene-set score (GSS)* matrix, denoted by \mathbf{Y} , whose rows are m gene-sets
 537 and columns are n observations. It is calculated as

$$\mathbf{Y} = \hat{\mathbf{G}}^T \mathbf{A} \mathbf{Q} \Delta^{[R]} \mathbf{P}^{[R]T} = \mathbf{W}^{[R]} \mathbf{F}^{[R]T} = \hat{\mathbf{G}}^T \mathbf{A} \mathbf{X}^{[R]} \quad (10)$$

538 where $\mathbf{Q}^{[R]}$ and $\mathbf{P}^{[R]}$ are the gene space and observation space within top R components. $\Delta^{[R]}$ is the
 539 diagonal matrix containing top R singular values. As a result, $\mathbf{X}^{[R]}$ is the reconstruction of \mathbf{X} using top R
 540 components. In practice, it is interesting to evaluate the contribution of a dataset or a component to the
 541 overall gene-set score. Therefore, we decompose gene-set scores with respect to data sets and
 542 components. The GSS matrix for dataset \mathbf{X}_k and component r is calculated as

$$\mathbf{Y}_k^r = \mathbf{W}_k^r \mathbf{F}_k^{rT} \quad (11)$$

543 we use superscript r to indicate the r th component and the subscript k to indicate the k th matrix (\mathbf{X}_k).

544 Similarly, \mathbf{W}_k^r denotes the r th dimension of gene-set space of matrix \mathbf{X}_k , \mathbf{F}_k^r is the r th component of the
 545 sample space. The outer product of the two vectors results in a GSS matrix for a specific components
 546 and dataset. Consequently, the overall gene-set score for component r (i.e. component-wise

547 decomposed gene-set scores) is the sum of the gene-set score matrix of the components across all
 548 datasets, that is,

$$\mathbf{Y}^r = \sum_k \mathbf{Y}_k^r = \sum_{k=1}^K \mathbf{W}_k^r \mathbf{F}_k^{rT} \quad (12)$$

549 Similarly, the overall gene-set score matrix by a single dataset (i.e. data-wise decomposed gene-set
 550 scores) is the sum of the matrices by all the components retained.

$$\mathbf{Y}_k = \sum_r \mathbf{Y}_k^r = \sum_{r=1}^R \mathbf{W}_k^r \mathbf{F}_k^{rT} \quad (13)$$

551 Therefore, the contribution of an individual dataset and/or component may be calculated. Finally, the
 552 complete gene-set score matrix is given by

$$\mathbf{Y} = \sum_r \mathbf{Y}^r = \sum_k \mathbf{Y}_k = \sum_{k=1}^K \sum_{r=1}^R \mathbf{W}_k^r \mathbf{F}_k^{rT} \quad (14)$$

553 which is the sum of all contributions by individual components and dataset. In practice, only the
 554 components with greatest variances (highest eigenvalues) should be retained in the analysis. If all
 555 components are retained, the result would be similar or exactly the same as naïve matrix multiplication
 556 (NMM; see later).

557 ***Evaluation of the significance of gene-set scores (calculating p-values)***

558 The expression (7) and (10) say that, for each observation, a gene-set score could be viewed as the mean
 559 of gene expression (in the reconstructed expression values $\mathbf{X}^{[R]}$) of genes in a particular gene-set.

560 If the candidate genes in a gene-set are randomly drawn from all features in $\mathbf{X}^{[R]}$ (null hypothesis), the
 561 distribution of the means of selected genes is given by central limited theorem (CLT),

$$\bar{x} \sim N(\mu, \sigma_{\bar{x}}) \text{ with } \sigma_{\bar{x}} = c \frac{\sigma}{\sqrt{h}} \quad (15)$$

562 Where μ is the mean of a column (observation) in $\mathbf{X}^{[R]}$, $\sigma_{\bar{x}}$ is the sampling standard deviation of means,
 563 σ is the standard deviation of the column in $\mathbf{X}^{[R]}$, h is the number of candidate genes mapped to \mathbf{X} in a
 564 gene-set and $c = \sqrt{(p-h)/(p-1)}$ is the finite population correction factor (p is the number of features in
 565 \mathbf{X}). It is used since each gene was only selected once in one gene-set.

566 **Gene influential score**

567 Gene-sets are composed of genes, and therefore we calculate the contribution of each feature to the
568 GSS, as it is interesting from a biological point of view to identify “driver” genes in a gene-set. In moGSA,
569 feature contribution, denoted by gene influential score (GIS), is calculated via a leave-one-out procedure.

570 The GSS of gene-set i , $\mathbf{Y}_{[i]}$, for all the observations are

$$\mathbf{Y}_{[i]} = \hat{\mathbf{G}}_{[i]}^T \mathbf{A}\mathbf{X}^{[R]} \quad (16)$$

571 where $\hat{\mathbf{G}}_{[i]}$ is the gene-set annotation vector for gene-set i . Correspondingly, the gene-set score for i th
572 gene-set excluding gene g is

$$\mathbf{Y}_{[i]}^{-g} = \hat{\mathbf{G}}_{[i]}^{-gT} \mathbf{A}\mathbf{X}^{[R]} \quad (17)$$

573 Where $\hat{\mathbf{G}}_{[i]}^{-g}$ is the gene-set annotation vector for gene-set i but without gene g . The influence of the
574 gene g is measured by

$$E_{[i]}^g = -\log_2 \frac{sd(\mathbf{Y}_{[i]}^{-g})}{sd(\mathbf{Y}_{[i]})} \quad (18)$$

575 where $sd(\cdot)$ stands for the function of calculating standard deviation. For convenience, the feature
576 influential score then is rescaled, such that the gene with maximum influence always equals 1. Therefore,
577 a positive $E_{[i]}^g$ suggests that gene g tends to have a positive correlation with gene-set score of gene-set i ,
578 whereas a gene with a negative value tends to have a negative correlation.

579 **Data simulation**

580 We simulated 100 multiple ‘omics data projects. Each simulated dataset was a triplet ($K=3$) containing
581 three data matrices (Figure S1), each matrix had the dimension 1000×30 , representing 30 matched
582 observations ($n=30$) and 1,000 features ($p_k=1,000$). Each of dataset of features had an annotation matrix,
583 which assigned each feature to one of 20 non-overlapping "gene-sets". The binary annotation matrix
584 had dimensions 1,000 features \times 20 gene-sets. Each gene-set contained 50 genes.

585 The 30 observations were defined by 6 equal sized clusters with 5 samples per cluster.

586 In each observation, 5 out of 20 gene-sets were simulated as differentially expressed (DE). Within the
587 same cluster, the same set of DE gene-sets were randomly selected as we assume that differentially
588 expressed (DE) gene-sets define the difference between clusters and observations. For a DE gene-set, a

589 number of genes were randomly simulated as DE genes (DEG), denoted as DEG_j . Random selection of
590 DEGs means that the DEGs in different datasets may overlap. In different simulations (Figure 2) we
591 varied the number of DEGs per gene-set (eg 5, 10 and 25 out of 50) or mean signal to noise ratio.

592 We used the following linear additive model adapted from [16], the expression or abundance of gene on
593 i th row and j th column is simulated as

$$x_{ij} = \alpha_i + \beta_l + \gamma_{ij} + \varepsilon_{ij} \quad (19)$$

594 where $i = 1, \dots, p$ is gene specific effect. $\beta_l \sim N(\mu = 0, \sigma = s)$ is the cluster effect. For observations
595 belongs to the same cluster l , the same β_l was applied. The cluster effect factor (categorical variable) is
596 introduced following the hypothesis that observations from the same clusters are driven by some
597 common pathways or “gene-sets” and ensures that observations from the same cluster have a higher
598 within than between cluster correlation. The six correlated clusters in the simulated data are captured
599 by first five components. We adjust the variance of each cluster, so that different variance would be
600 captured by the top five components. The cluster effect $\beta_l \sim N(\mu = 0, \sigma = s)$ is sampled from a
601 distribution with a mean of 0 and standard deviation s . The standard deviation (s) adjusts the correlation
602 between observations in the same cluster, and thus each cluster can have different variances. In this
603 study, we set $s = 0.3, 0.5$ and 1.0 , which lead to 25%, 30% and 50% of total variance are captured by the
604 top 5 components. $\varepsilon \sim N(\mu = 0, \sigma = 1)$ is the noise factor. γ_{ij} is a factor, if a gene is differentially
605 expressed (DE):

$$\gamma_{ij} \begin{cases} \sim N(\mu = m, \sigma = 1) & \text{if } i \in DEG_j \\ = 0 & \text{otherwise} \end{cases} \quad (20)$$

606 Apart from the retained variance, two other parameters are tuned in the simulation study. First is the
607 number of DEGs in a DE gene-set (5, 10 and 25 out of 50 DEGs). The second parameter is different
608 signal-to-noise ratio, which is tuned through modifying m in (20). The candidate m are 0.3, 0.5 and 0.8
609 standing for low, medium and high signal-to-noise ratio. In total, 100 projects of triplet datasets were
610 generated. The three matrix triplets were analyzed by moGSA. NMM, GSVA and ssGSEA, only accept one
611 matrix as input; therefore the three simulated matrices in one triplet set were concatenated. The
612 performance was assessed by the area under the ROC curve (AUC).

613 **Data**

614 **Downloading and Processing of Bladder Cancer TCGA data**

615 Normalized mRNA gene expression, copy number variation (CNV), microRNA (miRNA) expression data
616 and clinical information of BLCA were downloaded from TCGA (Date: 26/09/2014) using TCGA assembler
617 [39]. The processed mRNA gene expression had been obtained on the Illumina HiSeq platform and the
618 MapSplice and RSEM algorithm had been used for the short read alignment and quantification (Referred
619 as RNASeqV2 in TCGA) [40, 41]. The gene level CNV was estimated by the mean of copy number of
620 genomic region of a gene (retrieved by TCGA assembler directly). Patients that were present in both
621 gene expression and the CNV data were included in the analysis (n=308).

622 Before applying moGSA, minimal non-specific filtering of low variance genes was performed on both
623 datasets. RNA sequencing data (normalized count + 1) were logarithm transformed (base 10). Genes
624 were filtered to retain those with a total row sum greater than 300 and median absolute deviation (MAD)
625 greater than 0.1, which retained 14,692 unique genes (out of 20,531 genes). Then, RNA-seq gene
626 expression data were median centered. For the CNV data, genes with standard deviation greater than
627 the median were retained.

628 **Genome instability in TCGA BLCA tumors**

629 GISTIC2.0 [42] data for copy number gains/deletion in 24,776 unique genes were downloaded from
630 TCGA firehouse (<http://gdac.broadinstitute.org/>; download date 2015-03-09). The GISTIC encodes
631 homozygous deletion, heterozygous deletion, low-level gain and high-level amplification as -2, -1, 1 and
632 2 respectively. The four types of events were counted for each of the patients. The total number of
633 events were calculated by sum all four types of events.

634 **Downloading and Processing of the iPS ES 4-plex data**

635 The transcriptomic (RNA-sequencing), proteomic and phosphoproteomics data were downloaded from
636 Stem Cell-Omic Repository (Table S1, S2 and S5 from <http://scor.chem.wisc.edu/data.php>) [21]. In this
637 study, we used the 4-plex data, which consists of 17347 genes, 7952 proteins and 10499 sites of
638 phosphorylation in four cell lines. For the transcriptomics data, the expression levels of genes were
639 represented by RPKM values. Three replicates were available and we used the mean RPKM value of the
640 three replicates. Genes with duplicated symbols and low expression (summed RPKM < 12) were
641 removed. The iTRAQ quantification of protein and phosphorylation sites were performed by TagQuant
642 [43], as describe in [21]. The protein and sites of phosphorylation with low intensity (summed intensity

643 <20) were removed. In the proteomics data, proteins that are not mapped to an official symbol were
644 removed. Finally, all the data were logarithm transformed (base 10). After filtering, 10,961, 5817 and
645 7912 features were retained in the transcriptomic, proteomic and phosphor-proteomic datasets. A few
646 missing values still present and replaced with zero. The enrichment analysis was done on the gene
647 symbol levels, the specific phosphorylation sites were not considered.

648 **Sources of Gene-set annotation**

649 Gene-sets from the Molecular Signature Database MSigDB (version 4.0) [38] were used in this analysis.
650 The following MSigDB categories were included; MSigDB C2 curated pathways, C3 motif pathways
651 which included the transcription factor target (TFT) target gene-set and C5 gene ontology (GO) gene-sets
652 which included biological process (BP), cellular component (CC) and molecular function (MF) GO terms.
653 Among GO gene-sets, there were 825, 233 and 396 gene-sets in the BP, CC and MF categories
654 respectively. There were 617 TFT gene-sets. The pathway databases, Biocarta, KEGG and Reactome had
655 217, 186 and 674 gene-sets respectively. We excluded gene ontology terms that have more than 500
656 genes and less than 5 genes mapped to datasets. For example, in the BLCA analysis, gene-sets (1,454 in
657 total) were filtered to exclude those with less than 5 genes in a list of the concatenated features of CNV
658 and mRNA data resulting in 1,125 retained gene-sets.

659 **Other GSA methods (including NMM)**

660 Single gene-set method, including GSVA and ssGSEA methods were implemented using the
661 R/Bioconductor package GSVA [16]. Default settings were used for these methods. Naïve gene-set score
662 Y_{naive} was calculated through matrix multiplication (NMM).

$$663 \mathbf{Y}_{naive} = \hat{\mathbf{G}}^T \mathbf{X} \quad (21)$$

664 Therefore, the result of NMM is exactly the same as moGSA if all of the axes are retained.

664 **Clustering latent variable**

665 Consensus clustering was used [44, 45] to cluster the top five latent variables with Pearson correlation
666 distance and Ward linkage for the inner loop clustering. Eighty percent of patients were used in the re-
667 sampling step of clustering. In addition, different percentages of patients in the resampling were
668 evaluated. The results suggested the subtype model is robust with regard to different percentages of
669 samples used in resampling (Figure S20). Average agglomeration clustering was used in the final linkage
670 (linkage for consensus matrix) [44].

671 **Prediction strength to determine the optimal number of subtypes**

672 We used the “prediction strength” algorithm to assess the number of subtypes that can be predicted
673 from the data [22]. In prediction strength method, all samples were assigned a “true” subtype label
674 according to the clustering obtained from a given number of clusters. Then, the patients were then
675 divided into “training” and “testing” sets. KNN classifier was used to classify the patients in testing set.
676 Cross-validation suggested that there is no obvious good choice of K (Figure S21), but the number of K
677 does not have a big influence on the result (figure S22). We finally selected to use 9 nearest neighbors
678 (the middle of evaluated numbers). For each test, the agreement in assignment between predicted and
679 true labels were computed. The prediction strength is defined by the lowest proportion among all the
680 subtypes. It indicates the similarity between the true and predicted labels and ranges from 0 to 1, where
681 a value > 0.8 suggests a robust subtype classification [22]. Therefore, the model with the greatest
682 number of subtypes and prediction strength > 0.8 can be considered “optimal”. In this study, we
683 performed 100 random separations of training and testing sets and the prediction strength of each
684 randomization was calculated.

685

686 **List of Abbreviations**

- 687 ANOVA – analysis of variance
- 688 AUC – area under the ROC curve
- 689 BLCA – bladder cancer
- 690 BP – biological process
- 691 CC – cellular component
- 692 CCA – canonical correlation analysis
- 693 CIA – co-inertia analysis
- 694 CLT – central limited theorem
- 695 DE – differentially expressed
- 696 DEGS – differentially expressed gene-set
- 697 EMT – Epithelial to mesenchymal transition
- 698 GIS – gene influential score
- 699 GO – gene ontology
- 700 GS – gene-set
- 701 GSA – gene-set analysis
- 702 GSEA – gene-set enrichment analysis
- 703 GSS – gene-set score
- 704 MAD - median absolute deviation
- 705 MCIA – multiple co-inertia analysis
- 706 MF – molecular function
- 707 MFA – multiple factorial analysis
- 708 MVA – multivariate analysis

- 709 NMM – naïve matrix multiplication
- 710 PCA – principal component analysis
- 711 ROC - Receiver operating characteristic
- 712 SVD – singular value decomposition
- 713 TCGA – the cancer genome atlas
- 714 TF – transcriptional factor
- 715 TFT – transcriptional factor target

716 **Competing interests**

717 The authors declare no conflict of interest

718 **Authors' contribution**

719 AC conceived the study with CM and AMG. AC, CM and AMG developed the concept and experimental
720 design and wrote the manuscript. CM wrote the R code and conducted the experiments. AMG and AC
721 supervised the project. BK and BP had intellectual contribution to both the experimental design and
722 drafting the manuscript.

723 **Description of additional data files**

724 SupplementaryFigures.pdf – 23 supplementary figures

725 Table_S1.xlsx - Table S1 - the gene-set score (GSS) matrix of Gene ontology (GO) for iPS ES 4-plex data.

726 Table_S2.xlsx - Table S2: The Chi square test of association between integrative subtypes and previously
727 published subtypes.

728 Table_S3.xlsx - Table S3 - the gene-set score (GSS) matrix of Gene ontology (GO) and transcriptional
729 factor target (TFT) gene-set with more than 200 significant GSSs for BLCA data.

730 Table_S4.xlsx - Table S4 - the gene influential score (GIS) for selected gene-sets (from Gene Ontology).
731 The document contains GIS analysis for 9 gene-sets.

732 Table_S5.xlsx - Table S5 - the gene influential score (GIS) for selected transcriptional factor gene-sets.
733 The document contains GIS analysis for 2 gene-sets.

734

735 **Acknowledgements**

736 We wish to thank Prof. Joaquim Bellmunt for the insightful discussions about bladder cancer molecular
737 subtypes and treatment. We also thank Dr. Hannes Hanne and Dominic Helm for reading the manuscript
738 and giving the valuable suggestions. Funding for this work was provided by DFCI BCB Research Scientist
739 Developmental Funds, National Cancer Institute at the National Institutes of Health [grant numbers
740 2P50 CA101942-11, 1U19 AI111224-01, 1U19 AI109755-01] and Department of Defense BCRP [award

741 number W81XWH-15-1-0013]. Views and opinions of, and endorsements by the author(s) do not reflect
742 those of the US Army or the Department of Defense

743 Reference

744

- 745 1. Oszolak F, Milos PM: **RNA sequencing: advances, challenges and opportunities.** *Nature reviews*
746 *Genetics* 2011, **12**(2):87-98.
- 747 2. Metzker ML: **Sequencing technologies - the next generation.** *Nature reviews Genetics* 2010,
748 **11**(1):31-46.
- 749 3. Wilhelm M, Schlegl J, Hahne H, Moghaddas Gholami A, Lieberenz M, Savitski MM, Ziegler E,
750 Butzmann L, Gessulat S, Marx H *et al*: **Mass-spectrometry-based draft of the human proteome.**
751 *Nature* 2014, **509**(7502):582-587.
- 752 4. Meng C, Kuster B, Culhane AC, Gholami AM: **A multivariate approach to the integration of**
753 **multi-omics datasets.** *BMC bioinformatics* 2014, **15**:162.
- 754 5. de Tayrac M, Le S, Aubry M, Mosser J, Husson F: **Simultaneous analysis of distinct Omics data**
755 **sets with integration of biological knowledge: Multiple Factor Analysis approach.** *BMC*
756 *genomics* 2009, **10**:32.
- 757 6. Fagan A, Culhane AC, Higgins DG: **A multivariate analysis approach to the integration of**
758 **proteomic and gene expression data.** *Proteomics* 2007, **7**(13):2162-2171.
- 759 7. Le Cao KA, Martin PG, Robert-Granie C, Besse P: **Sparse canonical methods for biological data**
760 **integration: application to a cross-platform study.** *BMC bioinformatics* 2009, **10**:34.
- 761 8. Busold CH, Winter S, Hauser N, Bauer A, Dippon J, Hoheisel JD, Fellenberg K: **Integration of GO**
762 **annotations in Correspondence Analysis: facilitating the interpretation of microarray data.**
763 *Bioinformatics* 2005, **21**(10):2424-2429.
- 764 9. Khatri P, Sirota M, Butte AJ: **Ten years of pathway analysis: current approaches and**
765 **outstanding challenges.** *PLoS computational biology* 2012, **8**(2):e1002375.
- 766 10. Huang da W, Sherman BT, Lempicki RA: **Bioinformatics enrichment tools: paths toward the**
767 **comprehensive functional analysis of large gene lists.** *Nucleic acids research* 2009, **37**(1):1-13.
- 768 11. Barry WT, Nobel AB, Wright FA: **Significance analysis of functional categories in gene**
769 **expression studies: a structured permutation approach.** *Bioinformatics* 2005, **21**(9):1943-1949.
- 770 12. Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, Paulovich A,
771 Pomeroy SL, Golub TR, Lander ES *et al*: **Gene set enrichment analysis: a knowledge-based**
772 **approach for interpreting genome-wide expression profiles.** *Proceedings of the National*
773 *Academy of Sciences of the United States of America* 2005, **102**(43):15545-15550.

- 774 13. Tarca AL, Draghici S, Khatri P, Hassan SS, Mittal P, Kim JS, Kim CJ, Kusanovic JP, Romero R: **A**
775 **novel signaling pathway impact analysis**. *Bioinformatics* 2009, **25**(1):75-82.
- 776 14. Consortium EP: **The ENCODE (ENCyclopedia Of DNA Elements) Project**. *Science* 2004,
777 **306**(5696):636-640.
- 778 15. Gomez-Cabrero D, Abugessaisa I, Maier D, Teschendorff A, Merckenschlager M, Gisel A, Ballestar
779 E, Bongcam-Rudloff E, Conesa A, Tegner J: **Data integration in the era of omics: current and**
780 **future challenges**. *BMC systems biology* 2014, **8 Suppl 2**:11.
- 781 16. Hanzelmann S, Castelo R, Guinney J: **GSVA: gene set variation analysis for microarray and RNA-**
782 **seq data**. *BMC bioinformatics* 2013, **14**:7.
- 783 17. Tomfohr J, Lu J, Kepler TB: **Pathway level analysis of gene expression using singular value**
784 **decomposition**. *BMC bioinformatics* 2005, **6**:225.
- 785 18. Barbie DA, Tamayo P, Boehm JS, Kim SY, Moody SE, Dunn IF, Schinzel AC, Sandy P, Meylan E,
786 Scholl C *et al*: **Systematic RNA interference reveals that oncogenic KRAS-driven cancers require**
787 **TBK1**. *Nature* 2009, **462**(7269):108-112.
- 788 19. Lee E, Chuang HY, Kim JW, Ideker T, Lee D: **Inferring pathway activity toward precise disease**
789 **classification**. *PLoS computational biology* 2008, **4**(11):e1000217.
- 790 20. Abdi H, Williams LJ, Valentin D: **Multiple factor analysis: principal component analysis for**
791 **multitable and multiblock data sets**. *Wiley Interdisciplinary Reviews: Computational Statistics*
792 2013, **5**(2):31.
- 793 21. Phanstiel DH, Brumbaugh J, Wenger CD, Tian S, Probasco MD, Bailey DJ, Swaney DL, Tervo MA,
794 Bolin JM, Ruotti V *et al*: **Proteomic and phosphoproteomic comparison of human ES and iPS**
795 **cells**. *Nature methods* 2011, **8**(10):821-827.
- 796 22. Sjobahl G, Lauss M, Lovgren K, Chebil G, Gudjonsson S, Veerla S, Patschan O, Aine M, Ferno M,
797 Ringner M *et al*: **A molecular taxonomy for urothelial carcinoma**. *Clinical cancer research : an*
798 *official journal of the American Association for Cancer Research* 2012, **18**(12):3377-3386.
- 799 23. Knowles MA, Hurst CD: **Molecular biology of bladder cancer: new insights into pathogenesis**
800 **and clinical diversity**. *Nature reviews Cancer* 2015, **15**(1):25-41.
- 801 24. Damrauer JS, Hoadley KA, Chism DD, Fan C, Tiganelli CJ, Wobker SE, Yeh JJ, Milowsky MI, Iyer G,
802 Parker JS *et al*: **Intrinsic subtypes of high-grade bladder cancer reflect the hallmarks of breast**
803 **cancer biology**. *Proceedings of the National Academy of Sciences of the United States of America*
804 2014, **111**(8):3110-3115.

- 805 25. Choi W, Porten S, Kim S, Willis D, Plimack ER, Hoffman-Censits J, Roth B, Cheng T, Tran M, Lee IL
806 *et al*: **Identification of distinct basal and luminal subtypes of muscle-invasive bladder cancer**
807 **with different sensitivities to frontline chemotherapy.** *Cancer cell* 2014, **25**(2):152-165.
- 808 26. Senbabaoglu Y, Michailidis G, Li JZ: **Critical limitations of consensus clustering in class discovery.**
809 *Sci Rep* 2014, **4**:6207.
- 810 27. Tibshirani R, Walther G: **Cluster Validation by Prediction Strength.** *Journal of Computational*
811 *and Graphical Statistics* 2005, **14**(3):18.
- 812 28. Lindgren D, Frigyesi A, Gudjonsson S, Sjobahl G, Hallden C, Chebil G, Veerla S, Ryden T, Mansson
813 W, Liedberg F *et al*: **Combined gene expression and genomic profiling define two intrinsic**
814 **molecular subtypes of urothelial carcinoma and gene signatures for molecular grading and**
815 **outcome.** *Cancer research* 2010, **70**(9):3463-3472.
- 816 29. Biton A, Bernard-Pierrot I, Lou Y, Krucker C, Chapeaublanc E, Rubio-Perez C, Lopez-Bigas N,
817 Kamoun A, Neuzillet Y, Gestraud P *et al*: **Independent component analysis uncovers the**
818 **landscape of the bladder tumor transcriptome and reveals insights into luminal and basal**
819 **subtypes.** *Cell reports* 2014, **9**(4):1235-1245.
- 820 30. Chang W-C: **On Using Principal Components Before Separating a Mixture of Two Multivariate**
821 **Normal Distributions.** *Journal of the Royal Statistical Society Series C (Applied Statistics)* 1983,
822 **32**(3):9.
- 823 31. Alter O, Brown PO, Botstein D: **Singular value decomposition for genome-wide expression data**
824 **processing and modeling.** *Proceedings of the National Academy of Sciences of the United States*
825 *of America* 2000, **97**(18):10101-10106.
- 826 32. Hastie T, Tibshirani R, Eisen MB, Alizadeh A, Levy R, Staudt L, Chan WC, Botstein D, Brown P:
827 **'Gene shaving' as a method for identifying distinct sets of genes with similar expression**
828 **patterns.** *Genome biology* 2000, **1**(2):RESEARCH0003.
- 829 33. Holter NS, Mitra M, Maritan A, Cieplak M, Banavar JR, Fedoroff NV: **Fundamental patterns**
830 **underlying gene expression profiles: simplicity from complexity.** *Proceedings of the National*
831 *Academy of Sciences of the United States of America* 2000, **97**(15):8409-8414.
- 832 34. Leek JT, Scharpf RB, Bravo HC, Simcha D, Langmead B, Johnson WE, Geman D, Baggerly K,
833 Irizarry RA: **Tackling the widespread and critical impact of batch effects in high-throughput**
834 **data.** *Nature reviews Genetics* 2010, **11**(10):733-739.

- 835 35. Ashburner M, Ball CA, Blake JA, Botstein D, Butler H, Cherry JM, Davis AP, Dolinski K, Dwight SS,
836 Eppig JT *et al*: **Gene ontology: tool for the unification of biology. The Gene Ontology**
837 **Consortium**. *Nat Genet* 2000, **25**(1):25-29.
- 838 36. Gene Ontology C: **Gene Ontology Consortium: going forward**. *Nucleic acids research* 2015,
839 **43**(Database issue):D1049-1056.
- 840 37. Culhane AC, Schroder MS, Sultana R, Picard SC, Martinelli EN, Kelly C, Haibe-Kains B, Kapushesky
841 M, St Pierre AA, Flahive W *et al*: **GeneSigDB: a manually curated database and resource for**
842 **analysis of gene expression signatures**. *Nucleic acids research* 2012, **40**(Database issue):D1060-
843 1066.
- 844 38. Liberzon A, Subramanian A, Pinchback R, Thorvaldsdottir H, Tamayo P, Mesirov JP: **Molecular**
845 **signatures database (MSigDB) 3.0**. *Bioinformatics* 2011, **27**(12):1739-1740.
- 846 39. Zhu Y, Qiu P, Ji Y: **TCGA-assembler: open-source software for retrieving and processing TCGA**
847 **data**. *Nature methods* 2014, **11**(6):599-600.
- 848 40. Li B, Ruotti V, Stewart RM, Thomson JA, Dewey CN: **RNA-Seq gene expression estimation with**
849 **read mapping uncertainty**. *Bioinformatics* 2010, **26**(4):493-500.
- 850 41. Wang K, Singh D, Zeng Z, Coleman SJ, Huang Y, Savich GL, He X, Mieczkowski P, Grimm SA, Perou
851 CM *et al*: **MapSplice: accurate mapping of RNA-seq reads for splice junction discovery**. *Nucleic*
852 *acids research* 2010, **38**(18):e178.
- 853 42. Mermel CH, Schumacher SE, Hill B, Meyerson ML, Beroukhim R, Getz G: **GISTIC2.0 facilitates**
854 **sensitive and confident localization of the targets of focal somatic copy-number alteration in**
855 **human cancers**. *Genome biology* 2011, **12**(4):R41.
- 856 43. Wenger CD, Phanstiel DH, Lee MV, Bailey DJ, Coon JJ: **COMPASS: a suite of pre- and post-search**
857 **proteomics software tools for OMSSA**. *Proteomics* 2011, **11**(6):1064-1074.
- 858 44. Monti S, Tamayo P, Mesirov J, Golub T: **Consensus Clustering: A Resampling-Based Method for**
859 **Class Discovery and Visualization of Gene Expression Microarray Data**. *Machine Learning* 2003,
860 **52**(1-2):28.
- 861 45. Wilkerson MD, Hayes DN: **ConsensusClusterPlus: a class discovery tool with confidence**
862 **assessments and item tracking**. *Bioinformatics* 2010, **26**(12):1572-1573.
- 863