

1 **Evaluation of Gene Set Enrichment Analysis (GSEA) tools highlights the value of single**
2 **sample approaches over pairwise for robust biological discovery.**

3

4 Courtney Bull¹, Ryan M Byrne¹, Natalie C Fisher¹, Shania M Corry¹, Raheleh Amirkhah¹,
5 Jessica Edwards¹, Lily Hillson², Mark Lawler¹, Aideen Ryan³, Felicity Lamrock⁴, Philip D
6 Dunne^{1,5,^,*}, Sudhir B Malla^{1,^,*}

7

8 ¹The Patrick G Johnston Centre for Cancer Research, Queen's University Belfast, UK

9 ²School of Cancer Sciences, University of Glasgow, Glasgow, UK

10 ³Discipline of Pharmacology & Therapeutics, School of Medicine, College of Medicine,
11 Nursing and Health Sciences, University of Galway, Ireland.

12 ⁴Mathematical Sciences Research Centre, Queen's University Belfast, UK

13 ⁵Cancer Research UK Scotland Institute, Glasgow, UK

14 [^]Denotes Equal Authorship

15

16 ***Corresponding authors:**

17 Philip D. Dunne; p.dunne@qub.ac.uk, Sudhir B. Malla; s.malla@qub.ac.uk

18

19 **Running Title:**

20 Evaluation and Interpretation of Gene Set Enrichment Analysis Tools

21 **Keywords:** Transcriptional Signatures, GSEA, Molecular Classification, Computational Biology

22 **Acknowledgements:**

23 This work was supported by a CRUK early detection grant (A29834), a CRUK International
24 accelerator programme, ACRCelerate, (A26825), a UK Medical Research Council (MRC)
25 National Mouse Genetics Network programme (MC_PC_21042)

26 **Author Contributions:**

27 **CB:** data analysis, data visualisation, writing-original draft, writing-review and editing, **RB:**
28 writing-review and editing, **NCF:** writing-review and editing, **SMC:** writing-review and
29 editing, **RA:** writing-review and editing, **JE:** writing-review and editing, **LH:** writing-review
30 and editing, **ML:** writing-review and editing, **AR:** writing-review and editing, **FL:** data
31 analysis, writing-review and editing, **PDD:** conceptualisation, resources, supervision, writing-
32 original draft, writing-review and editing, **SBM:** supervision, writing-original draft, writing-
33 review and editing.

34 **Competing Interests:** The authors declare no conflicts of interest.

35 **Abstract**

36 **Background:** Gene set enrichment analysis (GSEA) tools can be used to identify biological
37 insights from transcriptional datasets and have become an integral analysis within gene
38 expression-based cancer studies. Over the years, additional methods of GSEA-based tools
39 have been developed, providing the field with an ever-expanding range of options to choose
40 from. Although several studies have compared the statistical performance of these tools, the
41 downstream biological implications that arise when choosing between the range of pairwise
42 or single sample forms of GSEA methods remain understudied.

43 **Methods:** In this study, we compare the statistical and biological interpretation of results
44 obtained when using a variety of pre-ranking methods and options for pairwise GSEA and
45 fast GSEA (fgGSEA), alongside single sample GSEA (ssGSEA) and gene set variation analysis
46 (GSVA). These analyses are applied to a well-established cohort of n=215 colon tumour
47 samples, using the clinical feature of cancer recurrence status, non-relapse (NR) and relapse
48 (R), as an initial exemplar, in conjunction with the Molecular Signatures Database “Hallmark”
49 gene sets.

50 **Results:** Despite minor fluctuations in statistical performance, pairwise analysis revealed
51 remarkably similar results when deployed using a range of gene pre-ranking methods or
52 across a range of choices of GSEA versus fgGSEA, with the same well-established prognostic
53 signatures being consistently returned as significantly associated with relapse status. In
54 contrast, when the same statistically significant signatures, such as Interferon Gamma
55 Response, were assessed using ssGSEA and GSVA approaches, there was a complete absence
56 of biological distinction between these groups (NR and R).

57 **Conclusions:** Data presented here highlights how pairwise methods can overgeneralise
58 biological enrichment within a group, assigning strong statistical significance to gene sets
59 that may be inadvertently interpreted as equating to distinct biology. Importantly, single
60 sample approaches allow users to clearly visualise and interpret statistical significance
61 alongside biological distinction between samples within groups-of-interest; thus, providing a
62 more robust and reliable basis for discovery research.

63 **Words: 309**

64 **Introduction**

65 Decreasing costs for sequencing, coupled with an increasing adoption of the FAIR principles¹,
66 have provided the cancer research field with a substantial amount of freely available
67 molecular datasets derived from tumour tissue samples. To ensure that these large datasets
68 can reveal important mechanistic insights, increased data availability has been coupled with
69 the development of transcriptional signatures that represent important biological pathways,
70 alongside easy-to-use algorithms that allow users to apply thousands of signatures
71 simultaneously to these data. These are exemplified by the establishment of the Molecular
72 Signatures Database (MSigDB)² and gene set enrichment analysis (GSEA) tools³, providing
73 the field with a stable set of reference templates and methods to compare across cohorts of
74 interest. The success of these approaches has led to a rapid expansion of established
75 signature collections in both human and mouse, most notably the MSigDB biological
76 “Hallmark” collection⁴ and development of programming software-based GSEA tools such as
77 the clusterProfiler⁵ and fast GSEA (fGSEA)⁶ R packages.

78 Given that many tumour cohorts have associated metadata linked to important features,
79 such as clinical outcome, the application of these large collections of signatures to cohorts in
80 conjunction with GSEA can serve as the basis for discovery and validation of biomarkers that
81 represent the biological characteristics of the chosen features, such as prognosis. This
82 approach is referred to as a supervised pairwise analysis, as the groups are known prior to
83 application of the GSEA method, and these tools have been tested extensively in terms of
84 the statistical robustness and performance in this setting^{7,8}. Once identified, these
85 biomarkers can be used as the basis for mechanistic investigations, pre-clinical model
86 development, and/or testing of a therapeutic target.

87 Alongside pairwise GSEA methods, approaches for single sample methods have been
88 developed, which differ from the pairwise approach in that they allow users to apply the
89 same transcriptional signature collections to all samples individually in a cohort, using single
90 sample GSEA (ssGSEA)⁹ and gene set variation analysis (GSVA)¹⁰. While these single sample
91 approaches are based on different statistical models to those in pairwise analyses, the
92 resulting outputs are based on the same gene signatures. Numerous studies have assessed
93 the statistical robustness and performance of this range of pairwise and single sample tools
94 separately^{7,11}. Despite differences being identified between methods when assessed using

95 statistically-driven criteria, few studies have focussed on the consequences in terms of
96 downstream biological approaches. Given that significant pairwise GSEA results can be
97 interpreted as representing the defining biological characteristics of a group of samples, the
98 absence of a comparative study across all approaches means that such an interpretation
99 may be based on incomplete evidence.

100 In this study, we use a fixed set of transcriptional signatures, in conjunction with a fixed
101 clinical feature (relapse status) within a well-characterised colon cancer (CC) transcriptional
102 cohort¹², to perform a series of pairwise and single sample assessments in tandem. Each
103 output is assessed based on the provided statistical values, however the primary focus of
104 this study is to assess how representative and uniform a significant pairwise result is when
105 assessed by single sample methods. Utilising a range of data visualisations and performance
106 measurements, we find that statistical results from a pairwise analysis often do not align
107 with biological distinction when using single sample outputs for the same signature.
108 Moreover, significant signatures identified from pairwise analysis can still be poor predictive
109 biomarkers of the clinical groups they were developed to represent.

110

111 **Words: 560**

112 **Methods**

113 **Datasets**

114 The transcriptional dataset used was previously assembled for the development of the FDA-
115 approved stage II ColDx/GeneFx risk-of-recurrence/relapse assay, consisting of n=215 stage II
116 primary tumours from CC patients profiled on the Almac disease-specific array, and available
117 from ArrayExpress, accession number E-MTAB-863¹². The cohort contained n=73 tumours
118 from patients that went on to develop distant metastasis within 5-year of surgery to remove
119 the primary tumour (relapse) (R) and n=142 tumours from patients that did not experience
120 relapse within five years following surgery (non-relapse) (NR). The E-MTAB-863 CEL files
121 were imported into Partek Genomics Suite (PGS; v6.6) and RMA normalised then log2
122 transformed. The probesets on the array were collapsed by importing the normalised data
123 into R (v3.3.2 or later) and, using the ‘collapseRows’ function from WGCNA (Weighted Gene
124 Coexpression Network Analysis, RRID:SCR_003302) package (v1.61)¹³, selecting the probeset
125 with the highest mean expression per gene.

126 **Differential gene expression analysis**

127 Differential expression analysis (DEA) was performed to measure differentially expressed
128 genes between R and NR CC. DEA was performed using the *limma* R package (v3.54.2).
129 Following DEA, genes were ranked using three different metrics, 1) the *t*-statistic (*t*-stat), 2)
130 the Log Fold Change (LogFC), and 3) the combination of LogFC and p-value (LogFC*-Log10(p-
131 value); hereafter as “combined”). The addition of p-value to LogFC adds statistical
132 significance to the directionality of LogFC. Separately, DEA was also performed for another
133 comparison between tumours classified as PDS1 and PDS3, using the *PDSclassifier* package¹⁴
134 with resulting groups being assessed using the same metrics and thresholds applied to the
135 R/NR analyses.

136 **Pairwise analysis**

137 To perform pairwise analysis two R packages were used, *clusterProfiler* (v4.6.2) and *fgsea*
138 (v1.24.0) and a random seed of 127 was set. Biological pathways were investigated using the
139 Hallmark gene sets from the MSigDB accessed through the *msigdbR* package (v7.5.1). Pre-
140 ranked GSEA was first performed using the GSEA function in *clusterProfiler* with 1000
141 permutations (nPermSimple = 1000, minGSSize = 1, maxGSSize = Inf). Enrichment plots for

142 GSEA were produced using the *gseaplot2* function in the *enrichplot* R package (v1.18.4).
143 GSEA was next conducted using the *fgsea* R package with the same parameters as
144 *clusterProfiler* (nPermSimple = 1000, minSize = 1, maxSize = Inf). Enrichment plots of *fgSEA*
145 were produced using the *plotEnrichment* function from the *fgsea* package. The online tool,
146 GenePattern¹⁵, <https://cloud.genepattern.org>, was also used to perform a pre-ranked
147 pairwise analysis, GSEAPreranked (v7.4.0). The Hallmark gene set collection was selected,
148 'h.all.v2023.2.Hs.symbols.gmt'. Default parameters were set except for 'collapse dataset'
149 which was set to 'FALSE'. Normalised enrichment score (NES) and false discovery rate (FDR)
150 values were recorded for each gene set within the two groups (R vs NR; PDS1 vs PDS3). A
151 gene set with an FDR *q*-value below 0.05 was deemed significant.

152 **Single sample analysis**

153 To perform single sample analysis the R/Bioconductor package *GSVA* (v1.46.0) was used
154 which facilitates ssGSEA⁹ and *GSVA*¹⁰. ssGSEA was performed with Hallmark⁴ gene sets from
155 MSigDB² and method set to "ssgsea". *GSVA* was performed with Hallmark gene sets from
156 MSigDB and the default parameters.

157 **Single sample analysis heatmaps**

158 For both ssGSEA and *GSVA*, matrix was formatted to include only Interferon Alpha Response,
159 Interferon Gamma Response and Epithelial Mesenchymal Transition (EMT), as previously
160 identified to be most significant by GSEA. The single sample scores were converted to Z-
161 scores and were plotted using the *ComplexHeatmap* (v2.14.0) R package and were grouped
162 using their respective groups (R vs NR; PDS1 vs PDS3).

163 **Data visualisation**

164 Additional visualisation R packages used for single sample analysis included: *smpplot2* (v
165 0.1.0), *ggribes* (v 0.5.4), *easyGgplot2* (v 1.0.0.9000), *pROC* (v 1.18.5), *randomForest* (v 4.7 -
166 1.1) and, *waterfalls* (v 1.0.0).

167 **Statistics**

168 The statistical report was generated on RStudio (4.2.2). A Student's *t*-test, from the *stats* (v
169 4.2.2) R package, was used to calculate significance of single sample scores between groups
170 (NR compared to R and PDS1 compared to PDS3). The *cortest* function from the *stats* (v

171 4.2.2) R package, with “pearson” method selected, was used for correlation analysis
172 between single sample enrichment scores for selected significant gene sets. The *cutpointr*
173 function in the *cutpointr* (v 1.1.2) R package was used to find the optimal cutpoint for the
174 single sample scores. Once calculated the single sample scores were centred around the
175 cutpoint resulting in a stratification of high and low scores for each of the gene sets being
176 tested.

177 ***“dualgsea”***

178 The pairwise method, fgSEA¹⁶ and single sample method, ssGSEA⁹ have been combined to
179 create an open source R-based function named *“dualgsea”*,
180 <https://github.com/MolecularPathologyLab/Bull-et-al>. The function enables the user to
181 apply the above statistical analysis and visualisations between two groups-of-interest.

182 **Words: 750**

183

184 Results

185 ***Variations in differential gene expression outputs across a range of methods do not alter*** 186 ***overall GSEA results.***

187 A typical goal when analysing bulk transcriptomic data, is the identification of discriminatory
188 biological signalling cascades that can serve as biomarkers to distinguish between group(s)-
189 of-interest; an output that can rapidly be delivered using transcriptional signatures in
190 conjunction with *in silico* analytical tools, such as pairwise gene set enrichment analysis
191 (GSEA)³ (Figure 1A). The initial step in this GSEA process requires all genes in the expression
192 matrix to be ranked based on their differential expression between the groups-of-interest.
193 For example, when using *limma*¹⁷ for microarray or *DESeq2*¹⁸ for RNA-seq, a ranked list of
194 genes can be produced based on *t*-statistics (*t*-stat) or Log Fold Change (LogFC) values, both
195 of which also provide directionality (up/down) according to the groups used. To assess the
196 outputs from each ranking metric, we compared the ranked order of genes following the
197 application of three approaches based on: 1) *t*-stat, 2) LogFC, and 3) the combination of
198 LogFC and *p*-value (LogFC * $-\text{Log}_{10}(\text{p-value})$; hereafter stated as combined) on expression
199 profiles from *n*=15,723 genes derived from *n*=215 FFPE stage II colon cancer samples (E-
200 MTAB-863)¹², where patients whose cancer relapsed following surgery (*n*=73) compared to
201 those who did not (NR; *n*=142) was used as an exemplar pairwise GSEA comparison (Figure
202 1B). Considering only the top and bottom 100 genes ordered based on *t*-stat (0.6% of genes
203 overall), gene ordering based on LogFC, or the combined rank, remained remarkably stable.
204 The top/bottom ranked genes identified using each method remain highly enriched at the
205 extremes relative to *t*-stat ranking (Figure 1B). When the genes were ranked by logFC the
206 majority (86%) of the top 100 genes fell within the top 500 genes when ranked by *t*-stat and
207 the remaining were represented within the top 2,707 genes. With the combined rank, 100%
208 of the top 100 genes were represented within the top 300 genes when ranked by *t*-stat.

209

210 To test if there were more profound downstream consequences of these small pre-ranking
211 gene order fluctuations, GSEA in clusterProfiler was performed⁵ using each of these ranking
212 metrics on the *n*=50 MSigDB 'Hallmark' gene sets. These analyses revealed that all three
213 ranking methods resulted in remarkably consistent gene sets being returned as significant
214 (FDR adjusted *p*-value < 0.05; *t*-stat =16/50, LogFC = 21/50, combined = 15/50), *n*=14 of the

215 n=22 total significant gene sets identified as common across from all three ranking methods
216 (Figure 1C; Supplementary Figure 1A). When the normalised enrichment score (NES) is
217 assessed to measure directionality, the direction of the n=14 overlapping significant gene
218 sets identified remained entirely consistent (Figure 1D), meaning that regardless of the pre-
219 ranking method used for these GSEA analyses, the biological interpretation will remain the
220 same. Furthermore, when gene sets that were identified as significant by one method but
221 not by the others, these were all enriched with the same directionality yet just below the
222 statistical significance threshold: again, confirming the similarities in outputs for GSEA using
223 all three pre-ranking methods (Supplementary Figure 1A).

224

225 ***Pairwise GSEA methods provide results with consistent downstream interpretation.***

226 As there were minimal differences in the GSEA outcome with the three ranking methods, *t*-
227 stat was used for the remainder of this study. Since the introduction of the original GSEA
228 method, several updated methodologies have been developed and in this study we
229 examined three derivatives of the GSEA method: 1) fast GSEA (fGSEA)⁶, 2) GSEA via
230 clusterProfiler⁵ (as used in Figure 1), which are both R-based tools, and 3) GSEA³ from the
231 Broad Institute GenePattern¹⁵ Server. The GSEA tool from GenePattern performs standard
232 GSEA with default signal-to-noise for ranking genes, however, the server also provides users
233 with a separate module called 'GSEAPreranked', where users can provide their own pre-
234 ranked gene list prior to analysis. To test outputs from each of these GSEA methods, relapse
235 (R) (n=73) and non-relapse (NR) (n=142) groups were compared across the CC cohort
236 previously used (E-MTAB-863), where these methods consistently identify the same
237 common statistically significant gene sets as identified in Figure 1E, additionally the
238 directionality of the NES for gene sets is consistent (Supplementary Figure 1B). Between
239 these three methods, n=3 gene sets were consistently upregulated in the NR group,
240 including Interferon Alpha Response and Interferon Gamma Response, and n=11 gene sets
241 were upregulated in the R group, such as EMT (Figure 1F-H); gene sets that have previously
242 been associated with prognosis in multiple cancer types, including colorectal cancer^{19, 20}.

243

244 ***Single sample GSEA methods provide biological insights that may be masked when using***
245 ***pairwise GSEA alone.***

246 Single sample GSEA (ssGSEA)⁹ has been proposed as an extension of the GSEA method, one
247 which can provide signature enrichment scores for each individual sample, rather than the
248 summarised “average” scores within groups of samples provided by pairwise GSEA, making
249 it suitable for both biological discovery and post-hoc assessments of individual samples
250 within any established groups-of-interest^{21 22}. Therefore, to compare the results obtained
251 from GSEA (Figure 1) with those from the single sample approaches, we explored two such
252 methods: 1) ssGSEA⁹, and 2) gene set variation analysis (GSVA)¹⁰ within our discovery cohort
253 (Figure 2A). Using the top three significant gene sets identified in Figure 1E, namely
254 Interferon Alpha Response, Interferon Gamma Response and EMT, these single sample
255 approaches were run using the GSVA R package by selecting either the “ssGSEA” or “gsva”
256 method. A correlative analysis was performed between the resulting ssGSEA and GSVA
257 scores, which revealed that both single sample methods were highly correlated, with a
258 significantly positive correlation across all three gene sets ($R>0.8$, $p<0.0001$; Figure 2B).
259 These results suggest that while the algorithms are different, the output of either single
260 sample methods provide consistent results.

261

262 Assessment of the ssGSEA and GSVA scores for the three gene sets that were significantly
263 different between the NR and R groups using GSEA, namely Interferon Alpha Response and
264 Interferon Gamma Response and EMT, revealed that there were comparable quantities of
265 high and low expression samples in each group, as indicated by the blue-to-red colours in
266 the heatmap (Figure 2C). To test this, a series of quantitative assessments were performed
267 using scores for the significant signatures using GSEA. Although the two clinical groups may
268 appear statistically significant for these single sample scores (Supplementary Figure 2C-H),
269 both clinical groups fall under the same distribution scale (Figure 2D-I), thus implying in
270 biological terms, they are not distinct for the signatures, which contradicts with GSEA
271 output. The range of ssGSEA scores showed large overlap between R and NR samples,
272 Interferon Alpha Response had 95.3% overlap between R and NR, Interferon Gamma
273 Response had 97.7% overlap between R and NR and EMT had 99.1% overlap between R and
274 NR. With respect to the GSVA results, Interferon Alpha Response scores had 95.8% overlap

275 between R and NR, Interferon Gamma Response had 98.1% overlap and EMT had 98.6%
276 between R and NR. Overall, these data highlight how even the most statistically significant
277 pairwise GSEA results may not be sufficient to identify transcriptional signalling that is
278 discriminatory between samples across two tumour groups.

279

280 **Visualisation of ssGSEA score is essential to ensure that statistical significance between**
281 **sample groups also represents distinct biology.**

282 There are a range of biomarker performance metrics that can be used to objectively test and
283 enumerate how well individual signatures represent the signalling within different groups of
284 samples. Therefore, a series of analyses were conducted to test the predictive value of the
285 most significant signatures identified by pairwise GSEA approaches (n=3) in identifying the
286 specific groups-of-interest that they were enriched in. We performed receiver operating
287 characteristic (ROC) analysis with the ssGSEA/GSVA scores and examined the area under
288 curve (AUC). NR patients displayed statistically significant enrichment in Interferon Alpha
289 and Interferon Gamma Response, implying that these signatures are contributing factors to
290 favourable outcome in NR patients (Supplementary Figure 2C-E), albeit GSVA Interferon
291 Gamma Response did not show any statistically significant enrichment in the NR samples
292 (Supplementary Figure 2F). However, if both interferon response signatures were then to be
293 used to develop a risk stratification tool to predict patient relapse status, the models
294 developed based on these signatures would perform underwhelmingly with the AUC
295 approximately ranging between 0.57 – 0.62 (Figure 3A, 3C). Furthermore, although there are
296 more NR (n=142) than R cases (n=73), when stratified into high and low groups for the
297 Interferon Alpha and Interferon Gamma Response signature scores using both ssGSEA and
298 GSVA, based on the optimal cut-offs defined by the AUROC analyses, ~30-50% of relapse
299 patients have high Interferon Alpha and Interferon Gamma Response scores (Figure 3B, 3D).
300 Likewise, regardless of its statistical significance (Supplementary Figure 2G-H), the EMT
301 ssGSEA and GSVA scores also perform poorly (AUC 0.60), with low sensitivity and specificity
302 as a relapse-specific biological signature for the purpose of risk stratification (Figure 3E-F).

303

304 Taken together, while each of these three signatures have been repeatedly shown to provide
305 statistical significance in terms of association with relapse outcomes, this is primarily due to
306 small (albeit statistically significant) differences in sample distributions, meaning that the
307 biological signalling these signatures are based on cannot be interpreted as reflecting
308 distinct mechanistic phenotypes or biological cascades between the two groups-of-interest.

309

310 ***Pathway-derived subtype serves as an exemplar for performing biological discovery using***
311 ***a single sample approach.***

312 As shown above, pairwise methods comparing relapse and non-relapse tumours can provide
313 users with statistically significant results, however these clinically distinct groups do not
314 represent uniformly biological distinct transcriptional subtypes. Therefore, to test the
315 performance of pairwise and single samples GSEA methodologies in groups of samples that
316 represent biologically distinct entities, we next performed these analyses contrasting
317 tumours based on our recent pathway-derived subtypes (PDS) ¹⁴ which identified three
318 statistically and biologically distinct subtypes; PDS1-3.

319 In this current study we now segregate our transcriptional cohort into these three PDS
320 classes (this dataset was not used in the original study) and perform a series of GSEA/ssGSEA
321 assessments on PDS1 (characterised by high MYC signalling) and PDS3 (characterised by low
322 MYC signalling) in conjunction with the performance metrics and visualisations used so far
323 (Figure 4A). Comparative analysis using the Hallmark gene sets collection and pairwise GSEA,
324 similar to the relapse-based comparisons, highlights a highly significant statistical difference
325 between PDS1 and PDS3 for MYC Targets V1 gene set (hereafter MYC V1; Figure 4B).
326 Importantly, unlike the assessment on R versus NR in the same cohort (Figure 1-3),
327 comparison of PDS1 to PDS3 clearly shows both statistical significance and biological
328 distinction when using single sample approaches (Figure 4C). Most importantly, unlike our
329 earlier analyses based on GSEA results comparing R and NR samples, these new assessments
330 across a known biology, reveal a remarkable difference and minimal overlapping distribution
331 for MYC V1 ssGSEA score, with only 6.7% of ssGSEA scores overlapping between PDS1 and
332 PDS3 (Figure 4D-E), implying that PDS1 and PDS3 can be considered as representing truly
333 distinct biological groups for MYC V1. This is further confirmed using ROC analysis, from both
334 ssGSEA and GSVA MYC V1 scores, which proves a sample will be classified as high MYC V1

335 when the sample is PDS1 with an AUROC = 0.99 (Figure 4F-G). We have created an open
336 source parallel pairwise/single sample R-based function “*dualgsea*”,
337 <https://github.com/MolecularPathologyLab/Bull-et-al>. The function produces multiple
338 visualisations and statistical analysis options that enables users to perform a broad
339 characterisation of their samples and groups-of-interest (Figure 4H).

340

341 **Words: 1875**

342

343

344 **Discussion**

345 In this study, we initially set out to provide a comparison of a number of well-established
346 gene set enrichment analysis (GSEA) methods, with particular emphasis on how choices of
347 standard bioinformatic pipelines can lead to differences in downstream biological
348 interpretation. As an exemplar of this, we assessed how consistent a significant pairwise
349 GSEA result is between pairwise approaches and also when the same signature is assessed
350 using single sample GSEA methods. These analyses highlight concordance *within* pairwise or
351 single sample approaches, however despite similar statistical performance, data presented
352 here provides a clear indication for how vastly different downstream interpretation of results
353 can be derived when using pairwise or single sample methods for the same transcriptional
354 signatures. Pairwise methods provide the user with strong statistical-based evidence of
355 differences in signature expression between two selected groups of samples, however this
356 can result in confusion when interpreting the biological significance of these differences, as
357 illustrated by enrichment scores across individual samples strongly overlapping between and
358 within groups. These results strongly support the use of single sample methods for class
359 discovery and mechanistic biomarker development/testing, given their consistency and
360 robustness in identifying distinct biological signalling between defined groups of samples.
361 Many previous studies have focussed on the statistical advantages and limitations of GSEA
362 methods, providing the field with important information on performance metrics for each
363 algorithm⁷. While these algorithms were developed to identify **statistical** significance
364 between user-selected groups of samples, they can occasionally be interpreted as
365 representing **biologically** distinct groups; a point that becomes even more important if the
366 results from GSEA-based methods are used to guide development of new pre-clinical models
367 that are interpreted as faithfully representing the clinical group-of-interest, or used as the
368 basis of developing prognostic/predictive biomarkers to guide clinical decision-making.

369 Data presented in this paper does not challenge the importance of studies using GSEA
370 methods, as we clearly demonstrate their value in identifying robust statistically distinct
371 groups. Our current study aims to provide an example of the consequence of method
372 selection for biological end-users with a primary interest in using these tools to identify
373 biologically distinct mechanistic signalling between two groups. For such end-users, we
374 propose that emphasis should be placed on more widespread use of visualisation methods

375 at an individual sample resolution, rather than the use of statistical values alone, to ensure
376 there is a clear distinction between the groups being compared²³. This point is particularly
377 important for biomarker discovery, where there is a requirement for the most robust and
378 discriminatory features that can be used to predict tumour groups with high sensitivity and
379 specificity. In addition, the identification of representative biological cascades that are both
380 statistically significant and biologically distinct between the two groups across a cohort of
381 tumours is increasingly important in the era of precision medicine, where interrogation of
382 transcriptional data can be used as the basis for development and testing of subtype-specific
383 therapeutic targets aimed at these patient groups.

384

385 An important feature for performing pairwise GSEA is the ranking of differentially expressed
386 genes. Our analyses highlight that the positions of individual differentially expressed genes
387 in an overall list will vary when using different ranking options. These results provide a clear
388 example of how the use of some of the most widely accepted tools for differential gene
389 expression analyses can lead to different users identifying conflicting biomarkers for the
390 same phenotypes in the exact same datasets. However, we find that the effects on using
391 different pre-ranking methods to rank genes for pairwise approaches have minimal effects
392 on biological interpretation when using downstream pathway analyses with any GSEA
393 method. As such, these data again support the use of pathway-level gene signatures as a
394 more representative way of measuring true biological phenotypes in transcriptional data,
395 over the use of individual gene-level biomarkers that can be undermined by technical biases
396 inherent in method choices for gene ranking. This single sample approach was used as basis
397 for class discovery within our recent pathway-derived subtypes (PDS)¹⁴ study, which used
398 ssGSEA scores to identify three biologically distinct classes of colorectal cancer that was
399 found to have prognostic value.

400 The cancer research field is accustomed to the heavy reliance on statistical thresholds as the
401 primary criteria for significance, as they provide users with a quantitative reference in
402 support of their findings. In data presented here we clearly show that additional
403 visualisation of these same data can lead to questions over the true biological significance of
404 such results. In this setting, if GSEA tools were used for discovery, the biological signalling
405 used as the basis for mechanistic studies could be indistinguishable across samples from

406 these different clinical groups, despite such signalling being based on statistically sound
407 evidence. Moving forward, it is essential to find a balance between statistical significance
408 and biological relevance, utilising visualisation techniques and analysis methods, including
409 distribution plots and ROC curves, to validate and contextualise findings. To ensure users can
410 recapitulate the approaches used here, we have developed an open source parallel
411 pairwise/single sample R-based function, “*dualgsea*”
412 <https://github.com/MolecularPathologyLab/Bull-et-al>, which provides multiple data
413 visualisation outputs and statistical tests, enabling all users to perform a comprehensive
414 assessment of their samples and groups-of-interest as shown in the comparison of PDS1 vs
415 PDS3 (Figure 4H).

416

417 Overall, our study sheds new light on the nuances between established gene set enrichment
418 methods, highlighting the challenges in interpreting results across different methods. The
419 work presented illustrates how a highly significant pairwise result does not always translate
420 to a significant single sample result when the same transcriptional data is analysed using the
421 same gene signatures. By carefully navigating these methods and their implications,
422 researchers can uncover novel meaningful biological insights from transcriptional data.

423

424 **Words: 933**

425

426 **Data Availability Statement**

427 Data is available in a public, open access repository. The “*dualgsea*” scripts used in this
428 current study are publicly available at <https://github.com/MolecularPathologyLab/Bull-et-al>.

429 **Figure Legends:**

430 **Figure 1. Differential gene expression analysis and pairwise analysis of the discovery**
431 **cohort.** (A) Schematic of the differential expression analysis and pairwise analysis. (B)
432 Workflow of differential expression analysis and ranked position of the top 100 differentially
433 expressed genes and bottom 100 genes in NR when ranked by t -stat and the position of
434 these genes when ranked by logFC and combined. (C) Venn diagram of the significant
435 Hallmark signatures ($p_{adj} < 0.05$) from GSEA when genes were ranked by t -stat, logFC, and
436 combined. (D) Significant Hallmark signatures ($p_{adj} < 0.05$) identified from clusterProfiler
437 GSEA when genes were ranked by t -stat, logFC, and logFC combined with the p -value
438 ordered by NES. (E) clusterProfiler GSEA, fgGSEA, and GenePattern pre-ranked GSEA of the
439 significant Hallmark gene sets. (F-H) clusterProfiler GSEA (F), fgGSEA (G), GenePattern (H)
440 comparing NR CC ($n=142$) to R CC ($n=73$) for Interferon Alpha Response, Interferon Gamma
441 Response and EMT.

442 **Figure 2. Comparison of the single sample methods, ssGSEA and GSVA.** (A) Schematic of
443 standard single sample analysis workflow. (B) Scatterplot showing the correlation of ssGSEA
444 scores and GSVA scores for the Hallmark Interferon Alpha Response (Pearson correlation
445 coefficient, $r = 0.835$), Interferon Gamma Response ($r = 0.878$) and EMT ($r = 0.953$). (C)
446 Heatmap of ssGSEA and GSVA scores for Interferon Alpha Response, Interferon Gamma
447 Response and EMT comparing NR and R. (D - I) Distribution of ssGSEA and GSVA scores. (D
448 and E) Distribution of the ssGSEA and GSVA scores for the Interferon Alpha Response
449 signature in the R (orange) and NR (blue) samples depicted using kernel density plots (D) and
450 histograms (E). (F and G) Distribution of the ssGSEA and GSVA scores for the Interferon
451 Gamma Response signature in the R (orange) and NR (blue) samples depicted using kernel
452 density plots (F) and histograms (G). (H and I) Distribution of the ssGSEA and GSVA scores for
453 the EMT signature in the R (orange) and NR (blue) samples depicted using kernel density
454 plots (H) and histograms (I)

455 **Figure 3. Application of single sample analysis as a predictor for relapse.** (A) ROC curve
456 using Interferon alpha response ssGSEA and GSVA scores to predict NR had an AUC ranging
457 between 0.61 – 0.62. True positive rate is when the sample is classified as high Interferon
458 Alpha Response, and the case was a NR. The true negative rate is the proportion of true
459 negatives, when a sample is a NR cases without high Interferon Alpha Response.

460 (B) Interferon Alpha Response waterfall plots show when stratified into high and low groups
461 for the Interferon Alpha Response with both ssGSEA and GSVA scores, we found a greater
462 number of NR patients classed as high (n=98 [75.4%], n=92 [76.0%]) respectively compared
463 to R (n=32 [24.6%], n=29 [24.0%]) respectively. (C) Interferon Gamma Response ROC AUC
464 values ranging between 0.57 – 0.61. True positive rate is when the sample is classified as
465 high Interferon Gamma Response, and the case was a NR. The true negative rate is the
466 proportion of true negatives, when a sample is a NR cases with a low Interferon Gamma
467 Response score. (D) Interferon Gamma Response waterfall plots show when stratified into
468 high and low groups for the Interferon Gamma Response with both ssGSEA and GSVA scores,
469 there were greater number of NR patients classed as high (n=86 [76.1%],; n= 95 [71.4%]
470 respectively compared to R (n=27 [23.9%],; n=38 [28.6%] respectively (E) EMT ROC AUC
471 values of 0.60. True positive rate is when the sample is classified as high EMT, and the case
472 was a relapse. The true negative rate is the proportion of true negatives, when a sample is a
473 NR case with a low EMT score. (F) EMT waterfall plots show when stratified into high and
474 low for EMT for ssGSEA we found that a greater number of R patients classified as high
475 (n=22 [59.5%], compared to NR (n=15 40.5%). GSVA found a higher number of NR patients
476 with a high EMT score (n=50 56.2%) compared to R (n=39 43.8%)

477 **Figure 4. The use of single sample analysis provides distinct biology between groups.** (A)
478 Schematic of application of pathway analysis methods when applied to PDS classification. (B)
479 GSEA revealed MYC targets V1 is enriched in the PDS1 group compared to the PDS3 group.
480 (C) ssGSEA scores show significant difference of MYC targets V1 expression between PDS1
481 and PDS3 groups (**** p-value < 0.0001). (D & E) ssGSEA scores for PDS1 and PDS3 show
482 little overlap of MYC targets V1 expression between groups. (F) ROC curve shows that the
483 MYC V1 scores enable discrimination between PDS1 and PDS3. AUC value of 0.99. True
484 positive rate is when the sample is classified as high MYC V1, and the case was PDS1. The
485 true negative rate is the proportion of true negatives, when a sample is a PDS1 without high
486 MYC V1. (G) Stratification of MYC V1 high and MYC V1 low ssGSEA scores showed that PDS1
487 was classified as high MYC V1 (n=54 [96.4%] and PDS3 contained only samples with a low
488 MYC V1 score (n=63 [100%]).

489 **References:**

- 490 1 Wilkinson, M. D. *et al.* The FAIR Guiding Principles for scientific data management and
491 stewardship. *Sci Data* **3**, 160018 (2016). <https://doi.org/10.1038/sdata.2016.18>
- 492 2 Liberzon, A. *et al.* Molecular signatures database (MSigDB) 3.0. *Bioinformatics* **27**, 1739-1740
493 (2011). <https://doi.org/10.1093/bioinformatics/btr260>
- 494 3 Subramanian, A. *et al.* Gene set enrichment analysis: A knowledge-based approach for
495 interpreting genome-wide expression profiles. *Proceedings of the National Academy of*
496 *Sciences* **102**, 15545-15550 (2005). <https://doi.org/doi:10.1073/pnas.0506580102>
- 497 4 Liberzon, A. *et al.* The Molecular Signatures Database (MSigDB) hallmark gene set collection.
498 *Cell Syst* **1**, 417-425 (2015). <https://doi.org/10.1016/j.cels.2015.12.004>
- 499 5 Wu, T. *et al.* clusterProfiler 4.0: A universal enrichment tool for interpreting omics data. *The*
500 *Innovation* **2**, 100141 (2021). <https://doi.org/https://doi.org/10.1016/j.xinn.2021.100141>
- 501 6 Korotkevich, G. *et al.* Fast gene set enrichment analysis. *bioRxiv*, 060012 (2021).
502 <https://doi.org/10.1101/060012>
- 503 7 Tarca, A. L., Bhatti, G. & Romero, R. A Comparison of Gene Set Analysis Methods in Terms of
504 Sensitivity, Prioritization and Specificity. *PLOS ONE* **8**, e79217 (2013).
505 <https://doi.org/10.1371/journal.pone.0079217>
- 506 8 Maleki, F., Ovens, K., Hogan, D. J. & Kusalik, A. J. Gene Set Analysis: Challenges,
507 Opportunities, and Future Research. *Front Genet* **11**, 654 (2020).
508 <https://doi.org/10.3389/fgene.2020.00654>
- 509 9 Barbie, D. A. *et al.* Systematic RNA interference reveals that oncogenic KRAS-driven cancers
510 require TBK1. *Nature* **462**, 108-112 (2009). <https://doi.org/10.1038/nature08460>
- 511 10 Hänzelmann, S., Castelo, R. & Guinney, J. GSVA: gene set variation analysis for microarray and
512 RNA-Seq data. *BMC Bioinformatics* **14**, 7 (2013). <https://doi.org/10.1186/1471-2105-14-7>
- 513 11 Chang, L.-C., Lin, H.-M., Sibille, E. & Tseng, G. C. Meta-analysis methods for combining
514 multiple expression profiles: comparisons, statistical characterization and an application
515 guideline. *BMC Bioinformatics* **14**, 368 (2013). <https://doi.org/10.1186/1471-2105-14-368>
- 516 12 Kennedy, R. D. *et al.* Development and independent validation of a prognostic assay for stage
517 II colon cancer using formalin-fixed paraffin-embedded tissue. *J Clin Oncol* **29**, 4620-4626
518 (2011). <https://doi.org/10.1200/JCO.2011.35.4498>
- 519 13 Langfelder, P. & Horvath, S. WGCNA: an R package for weighted correlation network analysis.
520 *BMC Bioinformatics* **9**, 559 (2008). <https://doi.org/10.1186/1471-2105-9-559>
- 521 14 Malla, S. B. *et al.* Pathway level subtyping identifies a slow-cycling biological phenotype
522 associated with poor clinical outcomes in colorectal cancer. *Nature Genetics* (2024).
523 <https://doi.org/10.1038/s41588-024-01654-5>
- 524 15 Reich, M. *et al.* GenePattern 2.0. *Nature Genetics* **38**, 500-501 (2006).
525 <https://doi.org/10.1038/ng0506-500>
- 526 16 Korotkevich, G. *et al.* Fast gene set enrichment analysis. (2021).
527 <https://doi.org/10.1101/060012>
- 528 17 Ritchie, M. E. *et al.* limma powers differential expression analyses for RNA-sequencing and
529 microarray studies. *Nucleic Acids Research* **43**, e47-e47 (2015).
530 <https://doi.org/10.1093/nar/gkv007>
- 531 18 Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion for
532 RNA-seq data with DESeq2. *Genome Biology* **15**, 550 (2014). <https://doi.org/10.1186/s13059-014-0550-8>
- 533 19 Guinney, J. *et al.* The consensus molecular subtypes of colorectal cancer. *Nature Medicine* **21**,
534 1350-1356 (2015). <https://doi.org/10.1038/nm.3967>
- 535 20 Corry, S. M. *et al.* Activation of innate-adaptive immune machinery by poly(I:C) exposes a
536 therapeutic vulnerability to prevent relapse in stroma-rich colon cancer. *Gut* **71**, 2502-2517
537 (2022). <https://doi.org/10.1136/gutjnl-2021-326183>
- 538

- 539 21 Wu, S. *et al.* Integrated Machine Learning and Single-Sample Gene Set Enrichment Analysis
540 Identifies a TGF-Beta Signaling Pathway Derived Score in Headneck Squamous Cell
541 Carcinoma. *J Oncol* **2022**, 3140263 (2022). [https://doi.org:10.1155/2022/3140263](https://doi.org/10.1155/2022/3140263)
- 542 22 Yi, M., Nissley, D. V., McCormick, F. & Stephens, R. M. ssGSEA score-based Ras dependency
543 indexes derived from gene expression data reveal potential Ras addiction mechanisms with
544 possible clinical implications. *Scientific Reports* **10**, 10258 (2020).
545 [https://doi.org:10.1038/s41598-020-66986-8](https://doi.org/10.1038/s41598-020-66986-8)
- 546 23 Yanai, I. & Lercher, M. A hypothesis is a liability. *Genome Biology* **21**, 231 (2020).
547 [https://doi.org:10.1186/s13059-020-02133-w](https://doi.org/10.1186/s13059-020-02133-w)
- 548

549 **Supplementary Figure Legends:**

550 **Supplementary figure 1. Ranking metrics of differentially expressed genes for GSEA have**
551 **little impact on GSEA results and GSEA methods have little variation.** (A) 50 Hallmark gene
552 sets from clusterProfiler GSEA when genes were ranked by *t*-stat, logFC, and combined,
553 highlighting the significant ($p_{adj} < 0.05$) hallmarks that are associated with all three ranking
554 methods. (B) clusterProfiler GSEA, fgSEA, and GenePattern pre-ranked GSEA 50 Hallmark
555 gene sets ranked by *t*-stat.

556

557 **Supplementary figure 2. Comparison of single sample analysis methods.** (A) ssGSEA
558 heatmap of 50 Hallmark gene sets. (B) GSVA heatmap of 50 Hallmark gene sets. (C)
559 Significance between NR and R ssGSEA scores for Interferon Alpha Response (** $p < 0.01$) (D)
560 Significance between NR and R GSVA scores for Interferon Alpha Response (** $p < 0.01$). (E)
561 Significance between NR and R ssGSEA scores for Interferon Gamma Response (* $p < 0.05$).
562 (F) No significance between NR and R GSVA scores for Interferon Gamma Response (ns). (G)
563 No significance between NR and R ssGSEA scores for EMT (ns) (H) Significance between NR
564 and R GSVA scores for EMT (* $p < 0.05$).

565

566

567

568

569

570

571

572

573

Figure 1

(which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY 4.0 International license](https://creativecommons.org/licenses/by/4.0/).

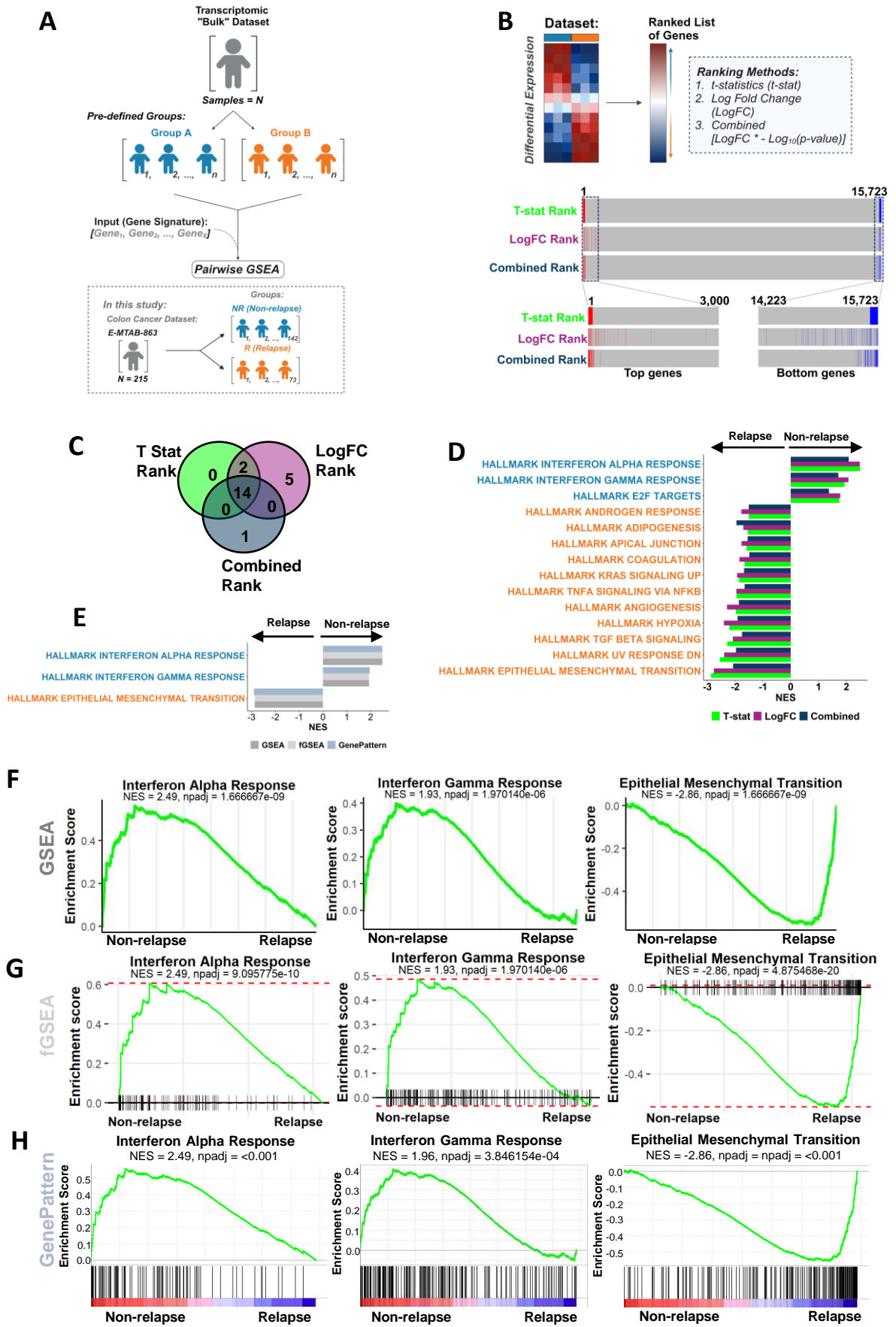
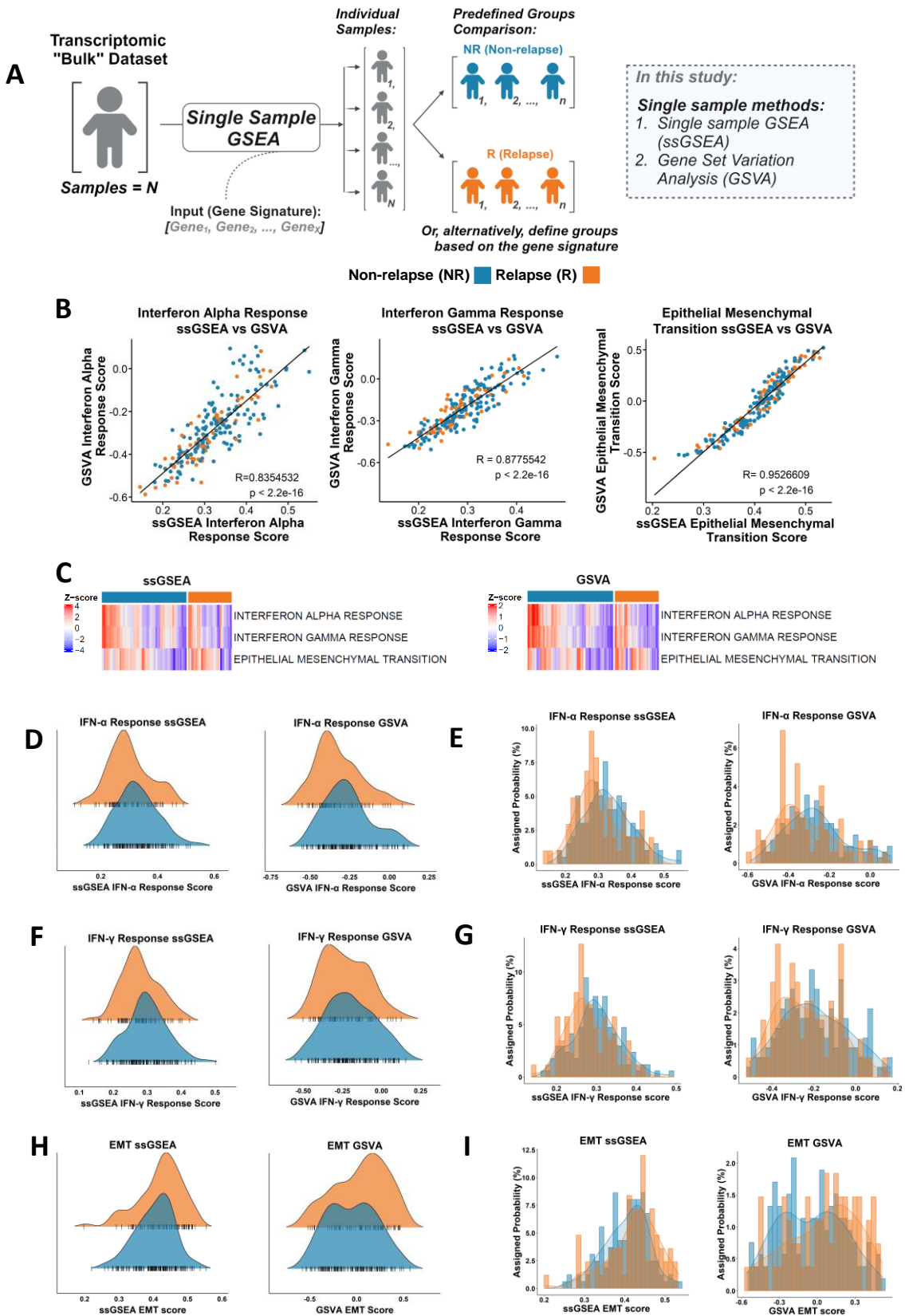


Figure 2

(which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.



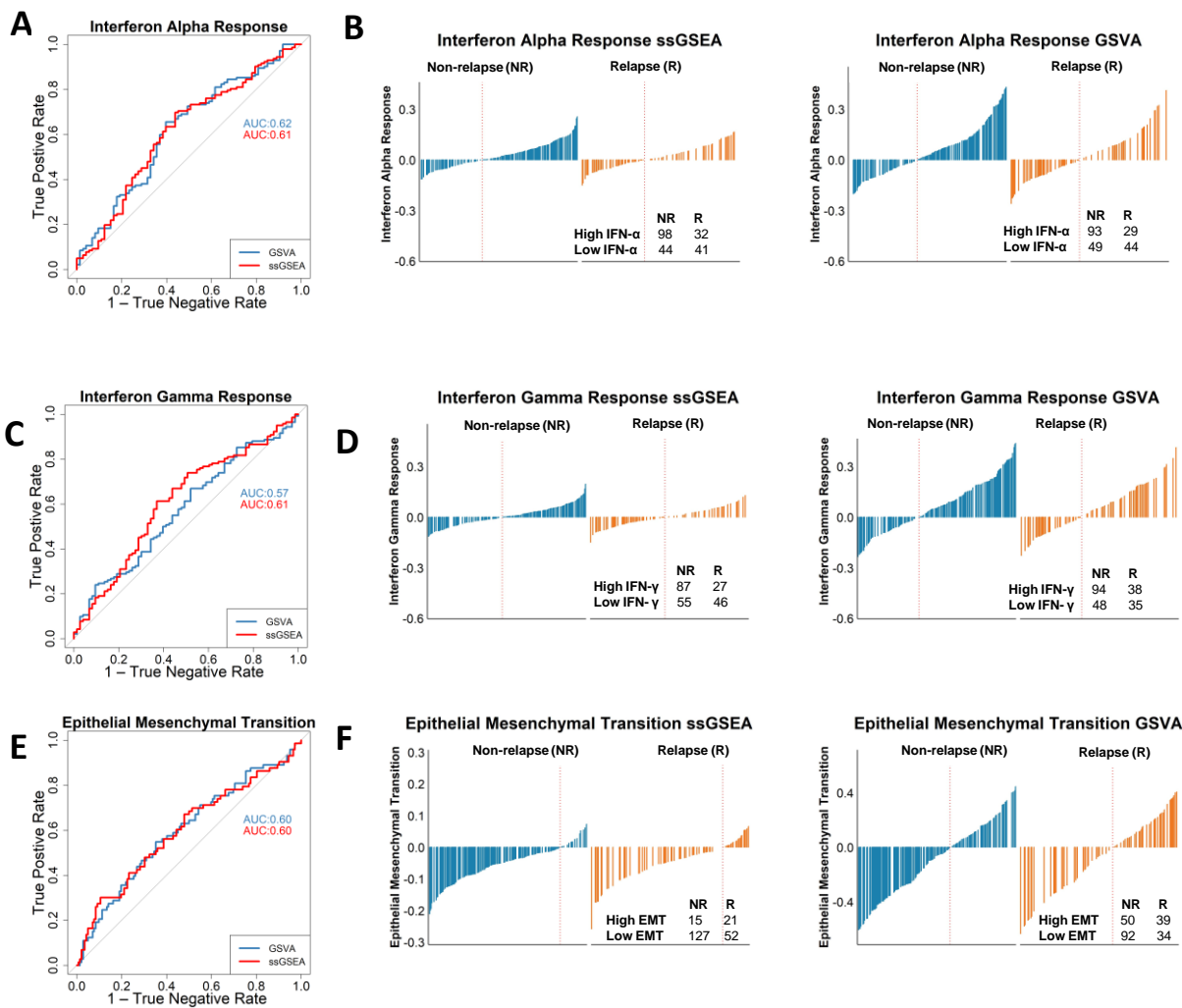


Figure 4

(which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY 4.0 International license](https://creativecommons.org/licenses/by/4.0/).

