# Individualized multi-omic pathway deviation scores

# using multiple factor analysis

Andrea Rau[1,2], Regina Manansala[2], Michael J. Flister[3],

Hallgeir Rui[3], Florence Jaffrézic[1], Denis Laloë[1*], Paul L. Auer[2*]

[1] GABI, INRA, AgroParisTech, Universite Paris-Saclay, 78350, Jouy-en-Josas, France

[2] Joseph J. Zilber School of Public Health, University of Wisconsin-Milwaukee, Milwaukee, WI 53201, USA

[3] Department of Pathology, Medical College of Wisconsin, Milwaukee, WI 53226, USA

[*] Corresponding authors: denis.laloe@inra.fr, pauer@uwm.edu

**ABSTRACT**

Malignant progression of normal tissue is typically driven by complex networks of somatic changes, including genetic mutations, copy number aberrations, epigenetic changes, and transcriptional reprogramming. To delineate aberrant multi-omic tumor features that correlate with clinical outcomes, we present a novel pathway-centric tool based on the multiple factor analysis framework called *padma*. Using a multi-omic consensus representation, *padma* quantifies and characterizes individualized pathway-specific multi-omic deviations and their underlying drivers, with respect to the sampled population. We demonstrate the utility of *padma* to correlate patient outcomes with complex genetic, epigenetic, and transcriptomic perturbations in clinically actionable pathways in breast and lung cancer.

**Keywords**: Multi-omic data, multiple factor analysis, pathways, cancer genomics

**BACKGROUND**

Large sets of patient-matched multi-omics data have become widely available for large-scale human health studies in recent years, with notable examples including the The Cancer Genome Atlas (TCGA)[1] and Trans-omics for Precision Medicine (TOPMed) program. The increasing emergence of multi-omic data has in turn led to a renewed interest in multivariate, multi-table approaches[2] to account for interdependencies within and across data types[3]. In such large-scale multi-level data, there is often limited or incomplete *a priori* knowledge of relevant phenotype groups for comparisons, and a primary goal may be to identify subsets of individuals that share common molecular characteristics, design therapies in the context of personalized medicine, or identify relevant biological pathways for follow-up. With these goals in mind, many multivariate approaches have the advantage of being unsupervised, using matched or partially matched omics data across genes, obviating the need for predefined groups for comparison as in the framework of standard differential analyses. A variety of such approaches has been proposed in recent years. For example, Multi-omics Factor Analysis (MOFA) uses group factor analysis to infer sets of hidden factors that capture biological and technical variability for downstream use in sample clustering, data imputation, and sample outlier detection[4].

In multi-omic integrative analyses, an intuitive first approach is to consider a gene-centric analysis, as we previously proposed in the *EDGE in TCGA* tool[5]. Expanding such analyses to the pathway-level is also of great interest, as it can lead to improved biological interpretability as well as reduced or condensed gene lists to facilitate the generation of relevant hypotheses. In particular, our goal is to define a method that quantifies an individual's deviation from a sample average, at the pathway-level, while simultaneously accounting for multiple layers of molecular information. Several related approaches for pathway-specific single-sample analyses have been proposed in recent years[6-8]. For example, PARADIGM[7] is a widely used approach based on structured

2

53    probabilistic factor graphs to prioritize relevant pathways involved in cancer progression as well

54    as identify patient-specific alterations; both pathway structures and multi-omic relationships are

55    hard-coded directly in the model, but it requires a discretization of the data and is now a closed-

56    source software, making extensions and application to other gene sets difficult. Pathway

57    relevance ranking[9] integrates binarized tumor-related omics data into a comprehensive network

58    representation of genes, patient samples, and prior knowledge to calculate the relevance of a

59    given pathway to a set of individuals. A pathway-centric supervised principal component-based

60    analysis implemented in *pathwayPCA*[10] performs gene selection and estimates latent variables

61    for association testing with respect to binary, continuous, and survival outcomes within each set

62    of omics data independently. Pathifier[6] instead seeks to calculate a personal pathway

63    deregulation score (PDS), based on the distance of a single individual from the median reference

64    sample on a principal curve; this principal curve approach is analogous to a nonlinear principal

65    components analysis (PCA), but can be applied only to a single-omic dataset (e.g., gene

66    expression). For both PARADIGM and Pathifier, clusters of scores across pathways are shown

67    to correlate with clinically relevant clustering of patients.

68

69    Here, we extend the basic philosophy of the Pathifier approach to multi-omics data, using an

70    innovative application of a Multiple Factor Analysis (MFA), to quantify individualized pathway

71    deviation scores. In particular, we propose an approach called *padma* ("PAthway Deviation scores

72    using Multiple factor Analysis") to characterize individuals with aberrant multi-omic profiles for a

73    given pathway of interest and to quantify this deviation with respect to the sampled population

74    using a multi-omic consensus representation. We further investigate the following succession of

75    questions. In which pathways are high deviation scores strongly associated with measures of

76    poor prognosis? For such pathways, which specific individuals are characterized by the most

77    highly aberrant multi-omic profile? And for such individuals, which specific genes and omics drive

78    large pathway deviation scores? By providing graphical and numerical outputs to address these

3

79    questions, *padma* represents both an approach for generating hypotheses as well as an

80    exploratory data analysis tool for identifying individuals and genes/omics of potential interest for

81    a given pathway.

82

83    There is already some precedent for using MFA to integrate multi-omic data, although existing

84    approaches differ from that proposed here. For instance, de Tayrac et al. suggested using MFA

85    for paired CGH array and microarray data, superimposed with functional gene ontology terms, to

86    highlight common structures and provide graphical outputs to better understand the relationships

87    between omics[11]. In addition, *padma* shares some similarities with a recently proposed integrative

88    multi-omics unsupervised gene set analysis called *mogsa*, which is similarly based on a MFA[12].

89    By calculating an integrated multi-omics enrichment score for a given gene set with respect to the

90    full gene list, *mogsa* identifies gene sets driven by features that explain a large proportion of the

91    global correlated information among different omics. In addition, these integrated enrichment

92    scores can be decomposed by omic and used to identify differentially expressed gene sets or

93    reveal biological pathways with correlated profiles across multiple complex data sets. However,

94    the fundamental difference in the two approaches is that *mogsa* evaluates pathway-specific

95    enrichment with respect to the entire set of genes, while *padma* instead focuses on identifying

96    and quantifying pathway-specific multi-omic deviations between each individual and the sampled

97    population.

98

99    **RESULTS AND DISCUSSION**

100

101   **Description of the approach**

102

103   *Pathway-centric multiple factor analysis for multi-omic data*

104

4

105    MFA represents an extension of principal component analysis for the case where multiple

106    quantitative data tables are to be simultaneously analyzed [13–16]. As such, MFA is a dimension

107    reduction method that decomposes the set of features from a given gene set into a lower

108    dimension space. In particular, the MFA approach weights each table individually to ensure that

109    tables with more features or those on a different scale do not dominate the analysis; all features

110    within a given table are given the same weight. These weights are chosen such that the first

111    eigenvalue of a PCA performed on each weighted table is equal to 1, ensuring that all tables play

112    an equal role in the global multi-table analysis. According to the desired focus of the analysis,

113    data can be structured either with molecular assays (e.g., RNA-seq, methylation, miRNA-seq,

114    copy number alterations) as tables (and genes as features within omics), or with genes as tables

115    (and molecular assays as features within genes). The MFA weights balance the contributions of

116    each omic or of each gene, respectively. In this work, we focus on the latter strategy in order to

117    allow different omics to contribute to a varying degree depending on the chosen pathway. In

118    addition, we note that because the MFA is performed on standardized features, simple differences

119    in scale between omics (e.g., RNA-seq log-normalized counts versus methylation logit-

120    transformed beta values) do not impact the analysis.

121

122    More precisely, consider a pathway or gene set composed of $p$ genes (Figure 1A), each of which

123    is measured using up to $k$ molecular assays (e.g., RNA-seq, methylation, miRNA-seq, copy

124    number alterations), contained in the set of gene-specific matrices $X_1, \ldots, X_p$ that have the same

125    $n$ matched individuals (rows) and $j_1, \ldots, j_p$ potentially unmatched variables (columns) in each,

126    where $j_g \in \{1, \ldots, k\}$ for each gene $g = 1, \ldots, p$. Because only the observations and not the

127    variables are matched across data tables, genes may be represented by potentially different

128    subset of omics data (e.g., only expression data for one gene, and expression and methylation
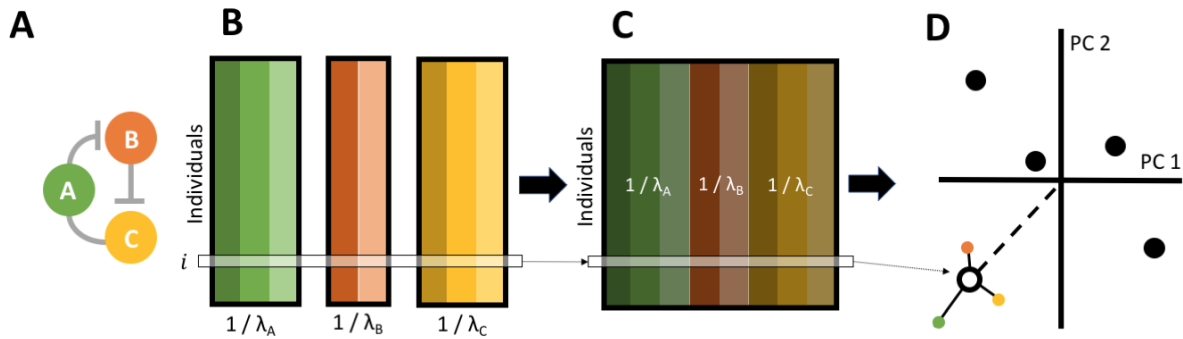
129    data for another).

130

131    In the first step, these data tables are generally standardized (i.e., centered and scaled). Next, an

132    individual PCA is performed using singular value decomposition for each gene table $X_g$, and its

133    largest singular value $\lambda_g^1$ is calculated (Figure 1B). Then, all features in each gene table $X_g$ are

134    weighted by $\frac{1}{\lambda_g^1}$, and a global PCA is performed using a singular value decomposition on the

135    concatenated set of weighted standardized tables, $X^* = \left[\frac{X_1}{\lambda_1^1}, \ldots, \frac{X_p}{\lambda_p^1}\right]$ (Figure 1C). This yields a

136    matrix of components (i.e., latent variables) in the observation and variable space. Optionally, an

137    independent set of supplementary individuals (or supplementary variables) can then be projected

138    onto the original representation; this is performed by centering and scaling variables for the

139    supplementary individuals (or individuals for the supplementary variables, respectively) to the

140    same scale as for the reference individuals, and projecting these rescaled variables into the

141    reference PCA space. Note that in the related *mogsa* approach, supplementary binary variables

142    representing gene membership in gene sets are projected onto a transcriptome-wide multiple

143    factor analysis to calculate gene set scores[12].

144

145    The MFA thus provides a consensus across-gene representation of the individuals for a given

146    pathway, and the global PCA performed on the weighted gene tables decomposes the consensus

147    variance into orthogonal variables (i.e., principal components) that are ordered by the proportion

148    of variance explained by each. The coordinates of each individual on these components, also

149    referred to as factor scores, can be used to produce factor maps to represent individuals in this

150    consensus space such that smaller distances reflect greater similarities among individuals. In

151    addition, *partial factor scores*, which represent the position of individuals in the consensus for a

152    given gene, can also be represented in the consensus factor map; the average of partial factor

153    scores across all dimensions and genes for a given individual corresponds to the factor score

6

154    (Figure 1D). A more thorough discussion of the MFA, as well as its relationship to a PCA, may be

155    found in the Supplementary Methods.

156



157

158    **Figure 1**. Illustration of the *padma* approach for calculating individualized multi-omic

159    pathway deviation scores. (A-B) For a given pathway, matched multi-omic measures for

160    each gene are assembled, with individuals in rows. Note that genes may be assayed for

161    varying types of data (e.g., measurements for one gene may be available for expression,

162    methylation, and copy number alterations, while another may only have measurements

163    available for expression and methylation). (C) Using a Multiple Factor Analysis, each gene

164    table is weighted by its *largest singular value*, and per-gene weighted tables are combined

165    into a global table, which in turn is analyzed using a Principal Component Analysis. (D)

166    Finally, each individual *i* is projected onto the consensus pathway representation; the

167    individualized pathway deviation score is then quantified as the distance of this individual

168    from the average individual. These scores can be further decomposed into parts attributed

169    to each gene in the pathway.

170

171    *Individualized pathway deviation scores*

172

173    In the consensus space obtained from the MFA, the origin represents the "average" pathway

174    behavior across genes, omics, and individuals; individuals that are projected to increasingly

7

175    distant points in the factor map represent those with increasingly aberrant values, with respect to

176    this average, for one or more of the omics measures for one or more genes in the pathway. To

177    quantify these aberrant individuals, we propose an individualized pathway deviation score $d_i$

178    based on the multidimensional Euclidean distance of the MFA component loadings for each

179    individual to the origin:

180    $$d_i^2 = \sum_{l=1}^{L} f_{i,l}^2,$$

181    where $f_{i,l}$ corresponds to the MFA factor score of individual *i* in component *l*, and *L* corresponds

182    to the rank of $X^*$. Note that this corresponds to the weighted Euclidean distance of the scaled

183    multi-omic data (for the genes in a given pathway) of each individual to the origin. These

184    individualized pathway deviation scores are thus nonnegative, where smaller values represent

185    individuals for whom the average multi-omic pathway variation is close to the average, while larger

186    scores represent individuals with increasingly aberrant multi-omic pathway variation with respect

187    to the average. An individual with a large pathway deviation score is thus characterized by one or

188    more genes, with one or more omic measures, that explain a large proportion of the global

189    correlated information across the full pathway.

190

191    Note that the full set of components is used for this deviation calculation, rather than subsetting

192    to an optimal number of components; we remark that due to their small variance relative to lower

193    dimensions, components from larger dimensions contribute relatively little to the overall pathway

194    deviation scores. Finally, to facilitate comparisons of scores calculated for pathways of differing

195    sizes (e.g., the number of genes), deviation scores with respect to the origin are normalized for

196    the pathway size.

197    *Decomposition of individualized pathway deviation scores into per-gene contributions*

198

8

199    In order to quantify the role played by each gene for each individual, we decomposed the

200    individualized pathway deviation scores into gene-level contributions. Recall that the average of

201    partial factor scores across all MFA dimensions corresponds to each individual's factor score. We

202    define the gene-level deviation for a given individual as follows:

203
$$d_{i,g} = \frac{\sum_{l=1}^{L} f_{i,l}(f_{i,l,g} - f_{i,})}{\sum_{l=1}^{L} f_{i,l}^2},$$

204    where as before $f_{i,l}$ corresponds to the MFA factor score of individual $i$ in component $l$, $L$

205    corresponds to the rank of $X^*$, and $f_{i,l,g}$ corresponds to the MFA partial factor score of individual $i$

206    in gene $g$ in component $l$. Note that by construction, the contributions of all pathway genes to the

207    overall deviation score sum to 0. In particular, per-gene contributions can take on both negative

208    and positive values according to the extent to which the gene influences the deviation of the

209    overall pathway score from the origin (i.e., the global center of gravity across individuals); large

210    positive values correspond to tables with a large influence on the overall deviation of an individual,

211    while large negative values correspond to genes that tend to be most similar to the global average.

212    In the following, we additionally scale these per-gene scores by the inverse overall pathway score

213    to highlight genes with highly atypical multi-omic measures both with respect to other genes in

214    the pathway and with respect to individuals in the population.

215

216    *Quantifying percent contribution of omics to pathway-centric multiple factor analysis*

217

218    The richness of MFA outputs also includes various decompositions of the total variance (that is,

219    the sum of the variances of each individual MFA component) of the multi-omic data for a given

220    pathway. Similarly to a standard PCA, the percent contribution of each axis of the MFA can be

221    calculated as the ratio between the variance of the corresponding MFA component and the total

222    variance; by construction, the fraction of explained variance explained decreases as the MFA

9

223    dimension increases. Similarly, the percent contribution to the inertia of each axis for a given omic,

224    gene, or individual can be quantified as the ratio between the inertia of its respective partial

225    projection in the consensus space and the inertia of the full data projection for that axis. These

226    per-gene, per-omic, and per-individual contributions can be quantified for a subset of components

227    (e.g., the first ten dimensions) or for the entire set of components; here, as we calculate

228    individualized pathway deviation scores using the full set of dimensions, we also calculated a

229    weighted per-omic contribution, which corresponds to the average contribution across all

230    dimensions, weighted by the corresponding eigenvalue.

231

232    **Application**

233

234    *TCGA data acquisition and pre-processing*

235

236    We illustrate the utility of *padma* on data from two cancer types with sufficiently large multi-omic

237    sample sizes in the TCGA database: invasive breast carcinoma (BRCA), which was chosen as

238    individuals have previously been classified into one of five molecular subtypes [17] (Luminal A,

239    Luminal B, Her2+, Basal, and Normal-like), as well as lung adenocarcinoma (LUAD), which was

240    chosen for its high recorded mortality.

241

242    The *padma* approach integrates multi-omic data by mapping omics measures to genes in a given

243    pathway. Although this assignment of values to genes is straightforward for RNA-seq, CNA, and

244    methylation data, a definitive mapping of miRNA-to-gene relationships does not exist, as miRNAs

245    can each potentially target multiple genes. Many methods and databases based on text-mining

246    or bioinformatics-driven approaches exist to predict miRNA-target pairs [18]. Here, we make use of

247    the curated miR-target interaction (MTI) predictions in miRTarBase (version 7.0)[19], using only

10

248    exact matches for miRNA IDs and target gene symbols and predictions with the "Functional MTI"

249    support type. Although the TCGA data used here have been filtered to include only those genes

250    for which expression measurements are available, there are cases where missing values are

251    recorded in other omics datasets (e.g., when no methylation probe was available in the promoter

252    region of a gene, or when no predicted MTIs were identified) or where a given feature has little or

253    no variance across individuals. In this analysis, features for a given omics dataset were removed

254    from the analysis only if missing values are recorded for all individuals or if the feature has minimal

255    variance across all individuals (defined here as $< 10^{-5}$ after scaling). After running *padma*, we

256    remark that the first ten MFA dimensions represent a large proportion of the total multi-omic

257    variance across pathways for both cancers (Supplementary Figure 5; BRCA median = 46.1%,

258    LUAD median = 51.9%).

259

260    As a measure of patient prognosis, we focused on two different metrics. First, we used the

261    standardized and curated clinical data included in the TCGA Pan-Cancer Clinical Resource

262    (TCGA-CDR)[20] to identify the progression-free interval (PFI). The PFI corresponds to the period

263    from the date of diagnosis until the date of the first occurrence of a new tumor event (e.g.,

264    locoregional recurrence, distant metastasis) and typically has a shorter minimum follow-up time

265    than measures such as overall survival. In the BRCA data, a total of 72 uncensored and 434

266    censored events were recorded (median PFI time of 792 and 915 days, respectively); among

267    LUAD individuals, a total of 65 uncensored and 79 censored events were recorded (median PFI

268    time of 439 and 683 days, respectively). Second, we used the histological grade for breast cancer,

269    which is an established cancer hallmark of cellular de-differentiation and poor prognosis[21]

270    (downloaded from http://legacy.dx.ai/tcga_breast on March 7, 2019). Tumors are typically graded

271    by pathologists on a scale of 1 (well-differentiated), 2 (moderately differentiated), or 3 (poorly

272    differentiated)    based    on    three    different    measures,    including    nuclear    pleomorphism,

11

273    glandular/tubule formation, and mitotic index, where higher grades correspond to faster-growing

274    cancers that are more likely to spread (Supplementary Table 1).

275

276    *Large deviation scores for relevant oncogenic pathways are associated with survival in lung*

277    *cancer*

278

279    The first question we address is the prioritization of pathways that are associated with a given

280    phenotype of interest. After processing the TCGA data and assembling the collection of gene

281    sets, we sought to identify a subset of pathways for which deviation scores were significantly

282    associated with patient outcome, as measured by PFI. To focus on pathways with the largest

283    potential signal (i.e., those for which a small number of individuals have very large deviation

284    scores relative to the remaining individuals) we consider only those with the most highly positively

285    skewed distribution of deviation scores. For each of the top 5% of pathways ($n = 57$) ranked

286    according to their Pearson's moment coefficient of skewness, we fit a Cox proportional hazards

287    (PH) model for the PFI on the pathway deviation score. Using the Benjamini-Hochberg[22] adjusted

288    p-values from a likelihood ratio test (FDR < 5%), we identified 14 pathways with deviation scores

289    that were significantly associated with the progression-free interval in lung cancer (Table 1; see

290    Supplementary Table 2 for the full gene lists in each pathway); for all of these, higher pathways

291    scores corresponded to a worse survival outcome. Note that the filtering on skewness of the

292    pathway scores is performed completely independently of the survival phenotype, ensuring that

293    the downstream survival analysis is not biased [23]. Of note, while candidates within the majority

294    deviated pathways (Table 1) have been univariately associated with patient outcome (e.g., cell

295    cycle, DNA repair, and apoptosis [24,25]), the *padma* TCGA analysis is unique in its ability to extend

296    these associations across multiple gene patient-specific perturbations within a pathway at the

297    genomic and transcriptomic RNA levels.

298

299    The detection of several pathways related to DNA repair (ATM, Homologous DNA repair,

300    BRCA1/2-ATR; Table1), as well as cell cycle and apoptosis related pathways, prompted us to

301    consider is whether these pathway deviation scores are simply acting as proxies for the tumor

302    mutational burden (i.e., the total number of nonsynonymous mutations) for each individual. To

303    investigate this, we estimated the mutational burden for each individual by counting the number

304    of somatic nonsynonymous mutations in a set of cancer-specific driver genes ($n$=183 and $n$=181

305    genes in breast and lung cancer, respectively) identified by IntOGen[26]. After adding a constant of

306    1 to these counts and log-transforming them, we fit a linear model to evaluate their association

307    with the pathway deviation scores; after correcting p-values from the Wald test statistic for multiple

308    testing (FDR < 10%), no pathways were found to be associated with the mutational burden. In

309    addition, when repeating the Cox PH model described above including the log-mutational burden

310    as an additional covariate, adjusted p-values were generally similar to previous values, and the

311    top six pathways remained significant at a significance threshold of 5%. This suggests that the

312    biological signal contained in the pathway deviation scores is indeed independent of that linked

313    to mutational burden.

314

| Pathway name | Pathway database | Adj. p-value | Hazard ratio | # of genes |
|---|---|---|---|---|
| D4-GDI (GDP dissociation inhibitor) signaling pathway | Biocarta | 0.0111 | 1.2692 | 13 |
| NF-kB activation through FADD/RIP-1 pathway mediated by caspase-8 and -10 | Reactome | 0.0111 | 1.2839 | 12 |
| Class I PI3K signaling events mediated by Akt | PID | 0.0251 | 1.1700 | 35 |
| ATM signaling pathway | Biocarta | 0.0265 | 1.1644 | 20 |
| CARM1 and regulation of the estrogen receptor | Biocarta | 0.0265 | 1.1426 | 35 |
| Homologous recombination repair of replication-independent double-strand breaks | Reactome | 0.0265 | 1.2432 | 16 |
| Role of BRCA1, BRCA2, and ATR in cancer susceptibility | Biocarta | 0.0467 | 1.1823 | 21 |

| | | | | |
|---|---|---|---|---|
| CD40L signaling pathway | Biocarta | 0.0467 | 1.1880 | 15 |
| Induction of apoptosis through DR4 and DR4/5 death receptors | Biocarta | 0.0467 | 1.1208 | 33 |
| Cell cycle: G1/S check point | Biocarta | 0.0467 | 1.1263 | 28 |
| Double stranded RNA induced gene expression | Biocarta | 0.0467 | 1.2007 | 10 |
| Signaling events mediated by HDAC class III | PID | 0.0467 | 1.1543 | 25 |
| HIV-1 Nef: Negative effector of Fas and TNF-alpha | PID | 0.0467 | 1.1268 | 35 |
| Regulation of telomerase | PID | 0.0467 | 1.0950 | 68 |

315      **Table 1**. Pathways whose deviation scores are significantly correlated with progression-

316      free interval in lung cancer. Hazard ratios and adjusted p-values correspond to a Cox PH

317      model for pathway deviation alone, with FDR < 5%. The number of genes for each pathway

318      corresponds to the number of genes with expression quantified by RNA-seq in the TCGA

319      data.

320

321      *Padma identifies individualized aberrations in the D4-GDP dissociation inhibitor signaling*

322      *pathway in lung cancer*

323

324      To illustrate the full range of results provided by *padma*, we focus in particular on the results for

325      the *D4-GDP dissociation inhibitor (GDI) signaling* pathway. D4-GDI is a negative regulator of the

326      ras-related Rho Family of GTPases, and it has been suggested that it may promote breast cancer

327      cell proliferation and invasiveness [27,28]. The D4-GDI signaling pathway is made up of 13 genes;

328      RNA-seq, methylation, and CNA measures are available for all 13 genes, with the exception of

329      CYCS and PARP1, for which no methylation probes were measured the promoter region. In

330      addition, miRNA-seq data were included for one predicted target pair: hsa-mir-421 → CASP3.

331      Over the 13 genes in the pathway, 130 of the 144 individuals had no nonsynonymous mutations,

332      while 13 and 1 individuals had 1 or 3 such mutations; ARHGAP5 and CASP3 were most often

333      characterized by mutations (3 individuals affected for each). Notably, although the D4-GDI

14

334    pathway has been previously implicated in breast cancer aggressiveness [27,28], this is to our

335    knowledge the first evidence suggesting that D4-GDI pathway might play a similar role in

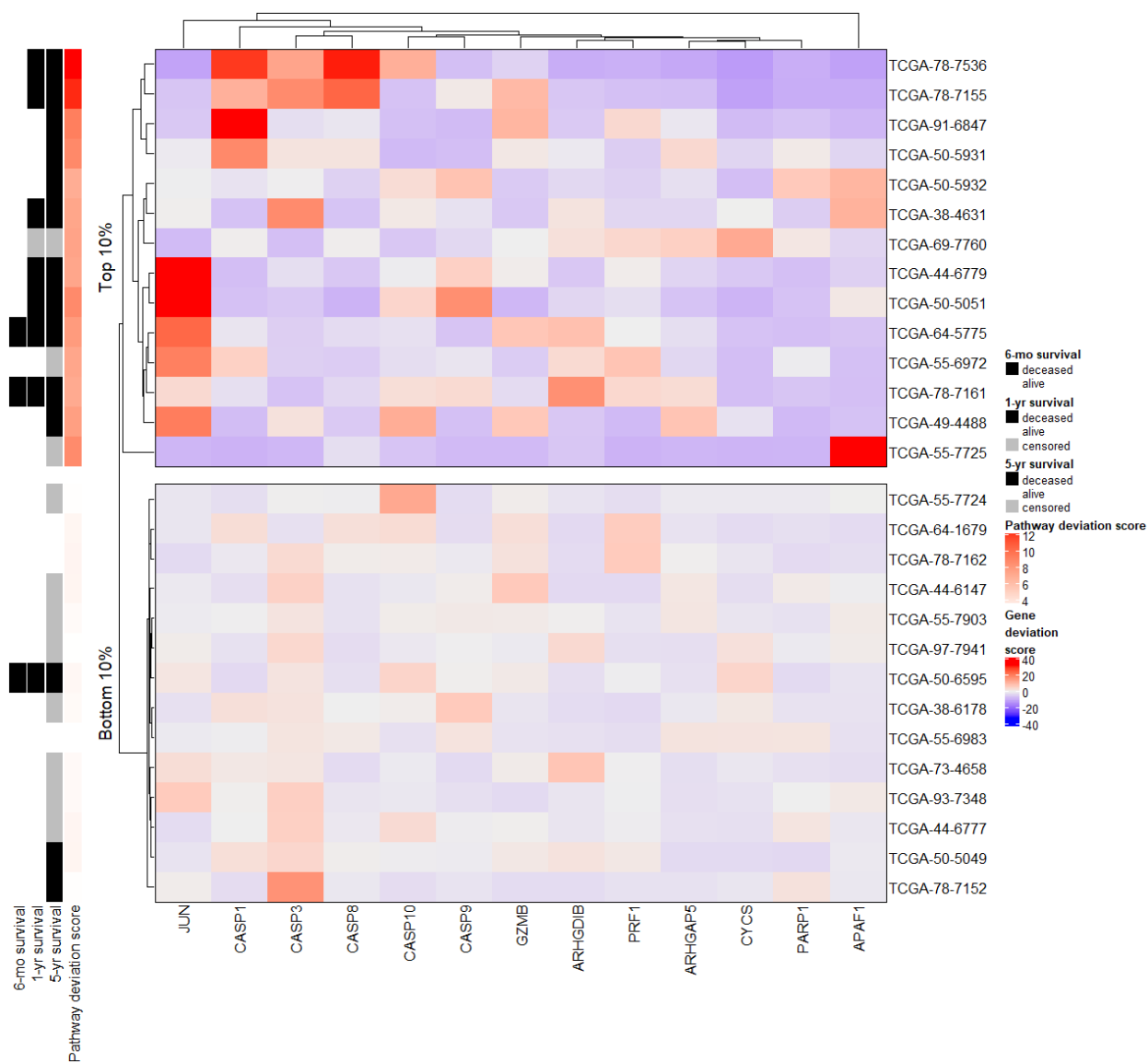336    promoting lung cancer.

337

338    Using the multi-omic data available for the D4-GDI signaling pathway, we can use the outputs of

339    *padma* to better understand the individualized drivers of multi-omic variation. In particular, it is

340    possible to quantify both gene-specific deviation scores as well as an overall pathway deviation

341    score for each individual, respectively based on the set of partial or full MFA components. We first

342    visualize the scaled gene-specific deviation scores for the top and bottom decile of individuals,

343    according to their overall pathway deviation score (Figure 2); these groups thus correspond to the

344    individuals that are least and most similar to the average individual within the population. We

345    remark that the 10% of individuals with the most aberrant overall scores for the D4-GDI signaling

346    pathway, who also had a high 1- and 5-year mortality rate, are those that also tend to have large

347    aberrant (i.e., red in the heatmap) scaled gene-specific deviation scores for one or more genes.

348    For example, the two individuals with the largest overall scores, TCGA-78-7536 and TCGA-78-

349    7155 (12.79 and 12.31, respectively), both had large scaled gene-specific scores for CASP3

350    (12.93 and 17.05, respectively), CASP1 (27.80 and 10.85, respectively), and CASP8 (29.72 and

351    22.61, respectively). While a subset of five individuals from the top decile were all characterized

352    by high deviation scores for JUN (TCGA-64-5775, TCGA-55-6972, TCGA-50-5051, TCGA-44-

353    6779, TCGA-49-4488), several other genes appear to have relatively small deviation scores for

354    all individuals plotted here (e.g., PRF1, PARP1). In addition, we remark the presence of highly

355    individualized gene-specific aberrations (e.g., APAF1 in individual TCGA-55-7725).

356

357    To provide an intuitive link between these gene-specific deviation scores with the original batch-

358    corrected multi-omics data that were input into *padma*, we further focus on the three genes

359    (CASP1, CASP3, and CASP8) for which large deviation scores were observed for the two highly

15

360    aberrant individuals (TCGA-78-7536 and TCGA-78-7155) in the D4-GDI signaling pathway. We

361    plot boxplots of the Z-scores for each available omic for the three genes across all 144 individuals

362    with lung cancer (Figure 3), specifically highlighting the two aforementioned individuals; full plots

363    of all 13 genes in the pathway are included in Supplementary Figure 1. This plot reveals that both

364    individuals are indeed notable for their overexpression, with respect to the other individuals, of

365    miRNA hsa-mir-421 (Figure 3D), which is predicted to target CASP3; in coherence with this, both

366    individuals had weaker CASP3 expression than average (although we note that its expression

367    was not particularly extreme with respect to the full sample). Individual TCGA-78-7536 appears

368    to have a hypomethylated CASP1 promoter, but a significantly higher number of copies of CASP8,

369    while individual TCGA-78-7155 is characterized by a large underexpression of CASP8 with

370    respect to other individuals. Both individuals appear to have deletions of CASP3, and

371    hypermethylated CASP8 promoters. This seems to indicate that, although the large overall

372    pathway deviations for these two individuals share some common etiologies, each also exhibit
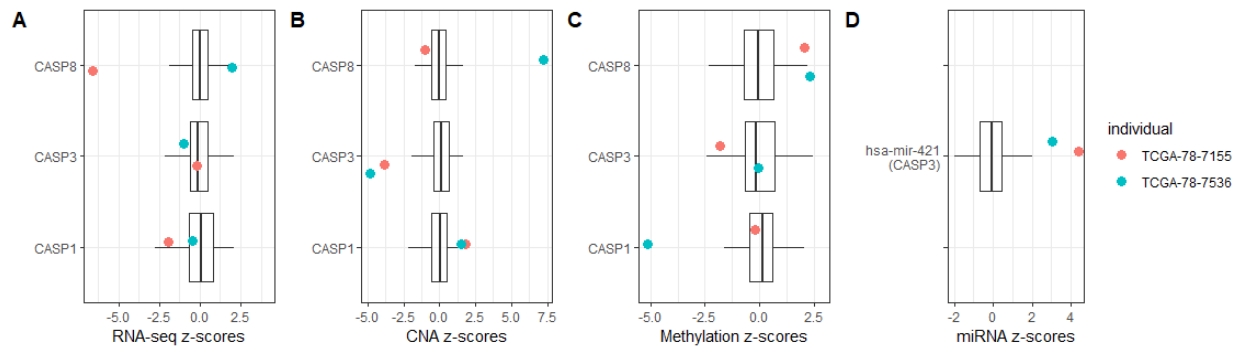
373    unique characteristics.

374

**Figure 2**. Scaled per-gene deviation scores for the D4-GDI signaling pathway for individuals corresponding to the top and bottom decile of overall pathway deviation scores. Red scores correspond to highly aberrant gene scores with respect to each individual's global score, while blue indicates gene scores close to the overall population average. Annotations on the left indicate the 6-month, 1-year, and 5-year survival status (deceased, alive, or censored) and overall pathway deviation score for each individual. Genes and individuals within each sub-plot are hierarchically clustered using the Euclidean distance and complete linkage.

**Figure 3**. Boxplots of Z-scores of gene expression (A), copy number alterations (B), methylation (C), and miRNA expression (D) for all individuals with lung cancer, with the 3 genes (CASP1, CASP3, CASP8) and one miRNA (hsa-mir-421, predicted to target CASP3) of interest in the D4-GDI signaling pathway. The two individuals with the largest pathway deviation score (TCGA-78-7155, TCGA-78-7536) are highlighted in red and turquoise, respectively.
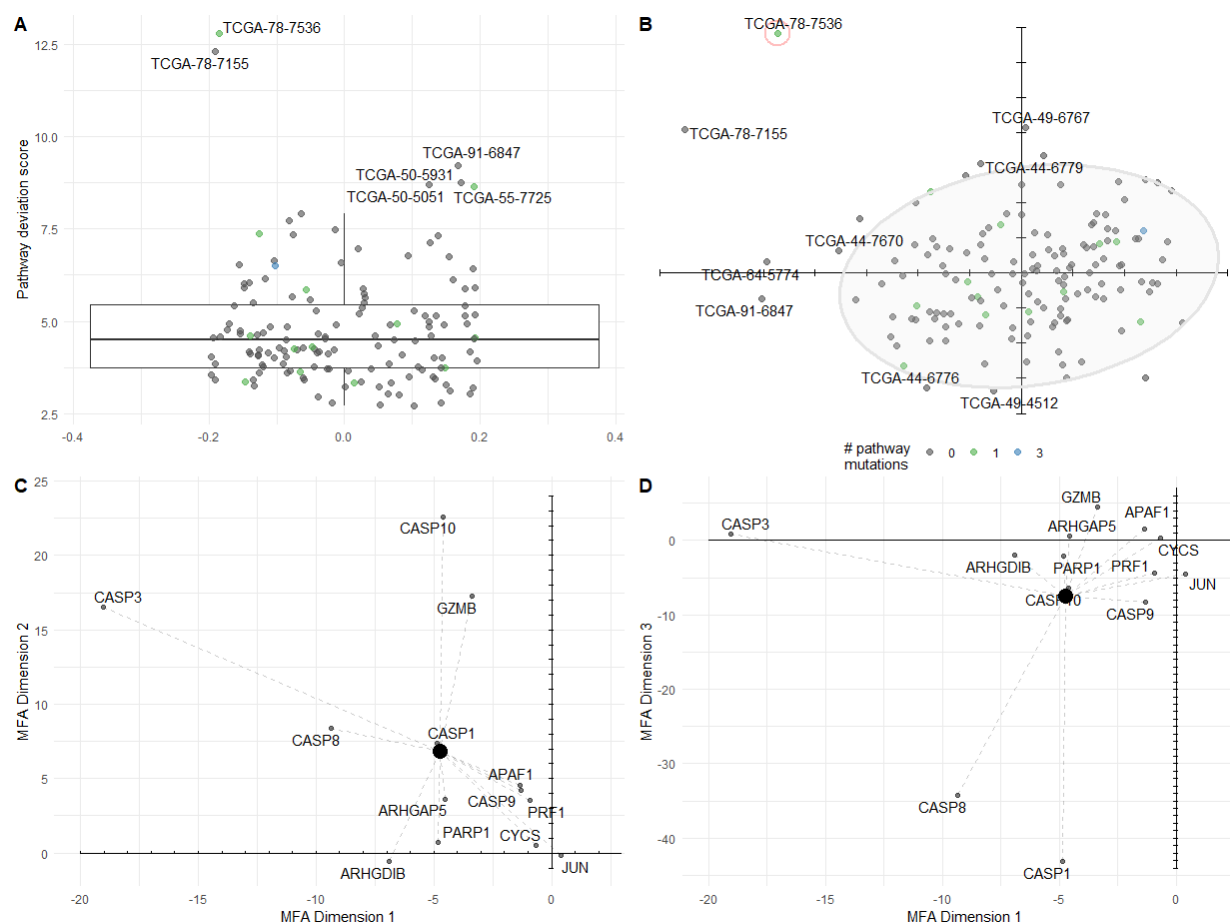
As overall pathway deviation scores represent the multi-dimensional average of these gene-specific deviation scores, a deeper investigation into them can also provide useful insight for a given pathway. We first note that the distribution of deviation scores for the D4-GDI signaling pathway (Figure 4A) is highly skewed, with a handful of individuals (e.g., TCGA-78-7536, TCGA-78-7155, TCGA-91-6847, TCGA-50-5931, TCGA-50-5051, and TCGA-66-7725) characterized by particularly large scores with respect to the remaining individuals. The individual with the most aberrant score for this pathway, TCGA-78-7536, had a single pathway-specific somatic mutation in the CASP1 gene, and a total of 7 cancer-specific driver gene mutations (corresponding to the 80th percentile of individuals considered here). Although these pathway deviation scores are calculated across all dimensions of the MFA, it can also be useful to represent individuals in the

18

407 first few components of the consensus MFA space (Figure 4B); the farther away an individual is

408 from the origin over multiple MFA dimensions, the larger the corresponding pathway deviation

409 score. In this case, we see that TCGA-78-7536 is a large positive and negative outlier in the

410 second (9.55% total variance explained), and third (8.07% total variance explained) MFA

411 components, respectively, although less so in the first component (11.97% total variance

412 explained). In addition, we note that RNA-seq is the major driver of the first MFA dimension

413 (54.38% contribution), while promoter methylation and copy number alterations take a larger role

414 in the second and third dimensions (42.29% and 59.18% contribution, respectively). miRNA

415 expression appears to play a fairly minor role in the MFA, with its maximum contribution (21.14%)

416 occurring at only the 16th dimension.

417

418 When examining the partial factor maps for this individual over the first three MFA dimensions

419 (Figures 4C-D), we note the large contribution of CASP3 (axis 1), CASP10 (axis 2), CASP1 and

420 CASP 8 (axis 3), as evidenced by their distance from the origin in these dimensions. Overall, this

421 is coherent with the previous gene-level analyses (Figure 2), where hypomethylation in CASP1

422 and large copy number gains for CASP3 and CASP8 with respect to the population were identified

423 for this individual. Other individuals with large overall deviation scores (e.g., TCGA-50-5931) are

424 not obvious outliers in the first two MFA dimensions, reflecting the fact that additional dimensions

425 play a more important role for them. Taken together, the individualized gene-specific and overall

426 pathway deviation scores output by *padma* provide complementary and interesting exploratory

427 insight into atypical multi-omic profiles for a given pathway of interest (here, the D4-GDI signaling

428 pathway in lung cancer).

429

**Figure 4**. (A) Distribution of pathway deviation scores for the D4-GDI signaling pathway in lung cancer; individuals with unusually large scores are labeled. (B) Factor map, representing the first two components of the MFA for the D4-GDI signaling pathway in lung cancer, with normal confidence ellipse superimposed. Individuals with extreme values in each plot are labeled with their barcode identifiers and colored by the number of pathway-specific nonsynonymous mutations. For the individual circled in red, TCGA-78-7536, a partial factor map representing the first MFA components 1 and 2 is plotted in (C), and MFA components 1 and 3 in (D). The large black dot represents the individual's overall pathway deviation score, as plotted in panel (B) for the first two axes, and gene-specific scores are joined to this point with dotted lines.

*Pathway deviation scores globally recapitulate histological grade in breast cancer*

For some cancers, additional clinical phenotypes beyond survival information may be of particular interest; to illustrate the use of *padma* in such a case, we focus on histological grade for breast cancer. To quantify whether pathway deviation scores tend to be associated with histological grade in breast cancer, we performed a one-way ANOVA on the three measures that comprise histological grade for each of the 1136 pathways. Based on the Benjamini-Hochberg[22] adjusted p-values from an F-test (FDR < 5%), all (1136) or nearly all (1135) pathways were found to have deviation scores that are significantly correlated with mitotic index and nuclear pleomorphism. Intriguingly, no pathways were found to be associated with degree of glandular/tubule formation; this may in part be due to the large proportion of individuals identified as grade III (poorly differentiated) for this measure (*n* = 285). The rankings of pathways based on mitotic index and nuclear pleomorphism were generally in agreement (Supplementary Figure 2). In all but two cases, higher deviation pathway scores corresponded to the higher grades for these two measures, corresponding to more aggressive tumors; the two exceptions were the *Presynaptic nicotinic acetylcholine recepto*r and *Highly calcium permeable postsynaptic nicotinic acetylcholine receptor* pathways (both from Reactome), for which the largest pathway deviation scores were associated with grade II, rather than grade III, of the mitotic index.

To prioritize pathways among this list, we calculated the rank product of the individual rankings by *p*-value for mitosis and nuclear pleomorphism; the top 10 pathways according to this joint ranking are shown in Table 2 (see Supplementary Table 3 for the full gene lists in each pathway). The *signaling by Wnt* pathway, which is made up of 63 genes, had the highest combined ranking for these two histological measures. Of this set of genes, all had RNA-seq, methylation, and CNA measures available, with the exception of FAM123B and PSMD10 (no CNA measures with nonzero variance) and PSMB1 to PSMB10, PSMC2, PSMC3, PSMC5, PSMC6, PSME1, and

21

469     PSME2 (no promoter methylation measures). miRNA-seq data were included for only two

470     predicted target pairs: hsa-mir-375 →CTNNB1 and hsa-mir-320a →CTNNB1. Over the 63 genes

471     in the pathway, 453 individuals had no nonsynonymous mutations, while 39, 6, 3, 2, and 1

472     individuals had 1, 2, 3, 4, or 5 such mutations; APC, PSMD1, and FAM123B were most often

473     characterized by mutations (10, 7, and 7 individuals affected, respectively).

474

| Pathway name | Pathway database | Combined ranking | # of genes |
|---|---|---|---|
| Signaling by Wnt | Reactome | 3.16 | 63 |
| Apoptotic execution phase | Reactome | 5.00 | 52 |
| APC/C:Cdh1 mediated degradation of Cdc20 and other APC/C:Cdh1 targeted proteins in late mitosis/early G1 | Reactome | 6.78 | 64 |
| Genes involved in Beta-catenin phosphorylation cascade | Reactome | 10.49 | 16 |
| Autodegradation of Cdh1 by Cdh1:APC/C | Reactome | 10.95 | 56 |
| Genes involved in M/G1 transition | Reactome | 11.62 | 72 |
| Regulation of the Fanconi anemia pathway | Reactome | 13.93 | 7 |
| Apoptotic cleavage of cellular proteins | Reactome | 14.14 | 38 |
| Apoptosis | Reactome | 14.28 | 143 |
| ER-phagosome pathway | Reactome | 15.62 | 58 |

475     **Table 2**. Pathways whose deviation scores are significantly correlated with measures of

476     histological grade (mitosis, nuclear pleomorphism) in breast cancer. Adjusted p-values

477     after Benjamini-Hochberg correction were $< 3.31 \times 10^{-12}$ for all pathways presented in the

478     table. Combined ranks correspond to the rank product of the individual rankings from

479     mitosis and nuclear pleomorphism, and the number of genes for each pathway

480     corresponds to the number of genes with expression quantified by RNA-seq in the TCGA

481     data.

482

483     Similarly to the distribution of D4-GDI pathway scores in lung adenocarcinomas, a small number

484     of breast cancer patients are characterized by highly aberrant scores in the signaling by Wnt

485    pathway, including TCGA-BH-A1FM, TCGA-E9-A22G, and TCGA-EW-A1PH, and the number of

486    pathway-specific nonsynonymous somatic mutations does not appear to be related to this score.

487    The associated factor map on the first two dimensions of the MFA (Figure 5A) clearly captures

488    relevant biological structure from the data, as evidenced by the quasi-separation of individuals in

489    different intrinsic inferred molecular subtypes (AIMS). Notably, individuals with Basal and Luminal

490    A breast cancer are clearly separated in the first two dimensions and tend to respectively have

491    positive and negative loadings in the first dimension of the MFA; Luminal B and Normal-like

492    subtypes largely overlap with the Luminal A subtype for this pathway, while Her2 is located

493    intermediate to the Luminal and Basal subtypes, as could be anticipated due to the equal

494    prevalence of Her2 amplification in both Luminal and Basal subtypes. Similar relevant biological

495    signal can be seen when considering a larger spectrum of pathways (Figure 5C). In particular,

496    individuals with the Basal and Luminal B subtypes tend to have much more highly variant

497    deviation scores across all pathways, whereas Luminal A and Normal-like subtypes are generally

498    much less variant.

499

500    When examining the percent contribution of each omic to the axes of the MFA for the Wnt

501    signaling pathway (Figure 5B), we remark the preponderant contribution of gene expression to

502    the first component (84.40%), while variability in the second component is largely driven by both

503    gene expression and copy numbers (45.66 and 35.37%, respectively). The large role played by

504    RNA-seq here is coherent with the definition of the AIMS subtypes themselves, which are defined

505    on the basis of gene expression. On average, after weighting by the eigenvalue of each

506    component, gene expression and copy number alterations were found to have similar

507    contributions to the overall variation (36.6%, 35.4%, respectively), while methylation played a less

508    important role (26.8%). For this pathway, as for most others we studied (Supplementary Figure

509    6), miRNA expression contributed relatively little to the overall variation (1.2%).

510

511    Taken together, these results illustrate that the *padma* approach, which is used in an

512    unsupervised manner on multi-omic cancer data for a given pathway, is able to recapitulate known

513    sample structure in the form of intrinsic tumor subtypes as well as relevant prognostic factors such

514    as histological grade.

515

516    **CONCLUSIONS**

517

518    Unsupervised dimension reduction approaches (such as PCA) have been widely used in genetics

519    and genomics for many years, both to identify sample structure and batch effects[29] and to

520    visualize overall variation in large data[30]. Here, we present a generalization of this approach to

521    multi-omic data for investigating biological variation at the pathway-level by aggregating across

522    genes, omic-type, and individuals. Compared to single-omics approaches (for instance, running

523    a PCA on RNA-seq data alone), *padma* accommodates multiple omics-sources which, for some

524    sample sets and pathways, account for more than 50% of the overall variation (Figure 5B). Using

525    MFA to partition variance, we construct a clinically relevant pathway disruption score that

526    correlates with survival outcomes in lung cancer patients, and histological grade in breast cancer

527    patients.

528

529    Our MFA-based approach allows investigators to (a) identify overall sources of variation (such as

530    batch effects); (b) prioritize high variance pathways defined by variability across subjects; (c)

531    identify aberrant observations (i.e., individuals) within a given pathway; and (d) identify the genes

532    and omics sources that drive these aberrant observations. For large, multi-omic data such as

533    TCGA, *padma* allows investigators to summarize overall variation and assist in generating

534    hypotheses for more targeted analyses and follow-up studies. As a case in point, we identified

535    two lung cancer patients with aberrant multi-omic profiles at three *CASP* genes. With access to

24

536    the tumor samples and more fine-grained clinical data, future molecular experiments could help

537    to clarify the role (if any) that these genes play in contributing to lung cancer mortality.

538

539    There are a number of natural extensions and alternative formulations to our MFA-based

540    approach. If comparisons between sets of individuals (e.g., healthy vs. disease) are of interest,

541    the MFA can be based on one set of samples (e.g., healthy, or a "reference set"), and the other

542    set of samples (e.g., diseased, or a "supplementary set") can be projected onto this original

543    representation. This is accomplished by centering and scaling supplementary individuals to the

544    same scale as the reference individuals, and projecting these rescaled variables into the

545    reference MFA space. In this setting, the interpretation of pathway deviation scores would no

546    longer correspond to the identification of "aberrant" individuals compared to an overall average,

547    but rather individuals that are most different from the reference set (e.g., the most "diseased" as

548    compared to a healthy reference); this strategy would be similar in spirit to the individualized

549    pathway aberrance score (iPAS) approach, which proposed using accumulated (unmatched)

550    normal samples as a reference set[31]. There is also no reason to limit this approach to pathways,

551    as the analysis could be performed just once, genome-wide (accordingly, inferences would no

552    longer be applicable to specific pathways). Here, we have structured the data with genes

553    representing data tables and omics representing columns within each table. Alternatively, the

554    data could be re-weighted by having omics represented as data tables and genes as columns

555    within each, similar to de Tayrac et al. (2009)[11]. Extensions to our work could include incorporating

556    the hierarchical structure of genes within pathways, or relatedness structure among samples. In

557    principle, other types of omics that do not map to genes or pathways (e.g., genotypes on single

558    nucleotide polymorphisms) could also be incorporated. Finally, though we illustrate the use of

559    *padma* for cancer genomics data, we anticipate that it will be broadly useful to other multi-omic

560    applications in human health or agriculture.

561

562 **MATERIALS AND METHODS**

563

564 *TCGA data acquisition and pre-processing*

565

566 The multi-omic TCGA data were downloaded and processed as described in Rau et al. (2019)[5].

567 Briefly, using *TCGA2STAT*[32] we downloaded processed TCGA Level 3 data from the Broad

568 Institute Genome Data Analysis Center (GDAC) Firehose on March 18, 2017 for individuals of

569 self-reported European ancestry for whom gene expression, methylation, copy number alterations

570 (CNA), microRNA (miRNA) abundance, and somatic mutation data were all available; this

571 ancestry filter was applied to minimize population-specific variance and focus on the group with

572 the largest available sample size. In addition, two individuals from the BRCA dataset (TCGA-E9-

573 A245, TCGA-BH-A1ES) were identified as outliers with consistently extreme deviation scores

574 across multiple pathways and were removed from the remainder of the analyses; the final sample

575 sizes were thus $n$=504 and $n$=144 individuals for the BRCA and LUAD datasets, respectively.

576

577 Per-gene normalized expression estimates were calculated using RSEM[33]. Methylation was

578 quantified using the maximally variant probe from the Illumina Infinium Human Methylation450

579 BeadChip located within ±1500bp of the transcription start site, and representative probe beta

580 measures were transformed to the logit scale. Somatic CNAs were called by comparing Affymetrix

581 6.0 probe intensities from normal (i.e., non-cancer tissue) and cancer tissue, and genome

582 segments were aggregated to gene-level measures by *TCGA2STAT* and *CNTools*. Individuals

583 were classified as carriers or noncarriers of a nonsynonymous somatic mutation for each gene

584 using *TCGA2STAT*. Normalized miRNA abundance was quantified as Reads per million

585 microRNA mapped (RPMMM) values. RNA-seq and miRNA-seq quantifications were TMM-

586 normalized[34], converted to counts per million (CPM), and log2-transformed. Only genes with

587 available RNA-seq expression measures were retained for the remainder of the analysis,

26

588  corresponding to 20,501 and 19,971 genes for BRCA and LUAD, respectively. Finally, batch

589  effects have been shown to have a strong impact on the analysis of high-throughput data in

590  general[29] and for the TCGA data specifically[35]. As specific sample plates have been shown to

591  represent significant batch effects in previous analyses[36], each processed omic (with the

592  exception of somatic mutation data) was individually batch adjusted for each cancer to correct for

593  plate-specific effects using `removeBatchEffects` in limma[37]. Plots of the first two components

594  from a transcriptome-wide and genome-wide single-omics PCA and multi-omics MFA for the

595  batch-corrected data are included in Supplementary Figures 3 and 4.

596
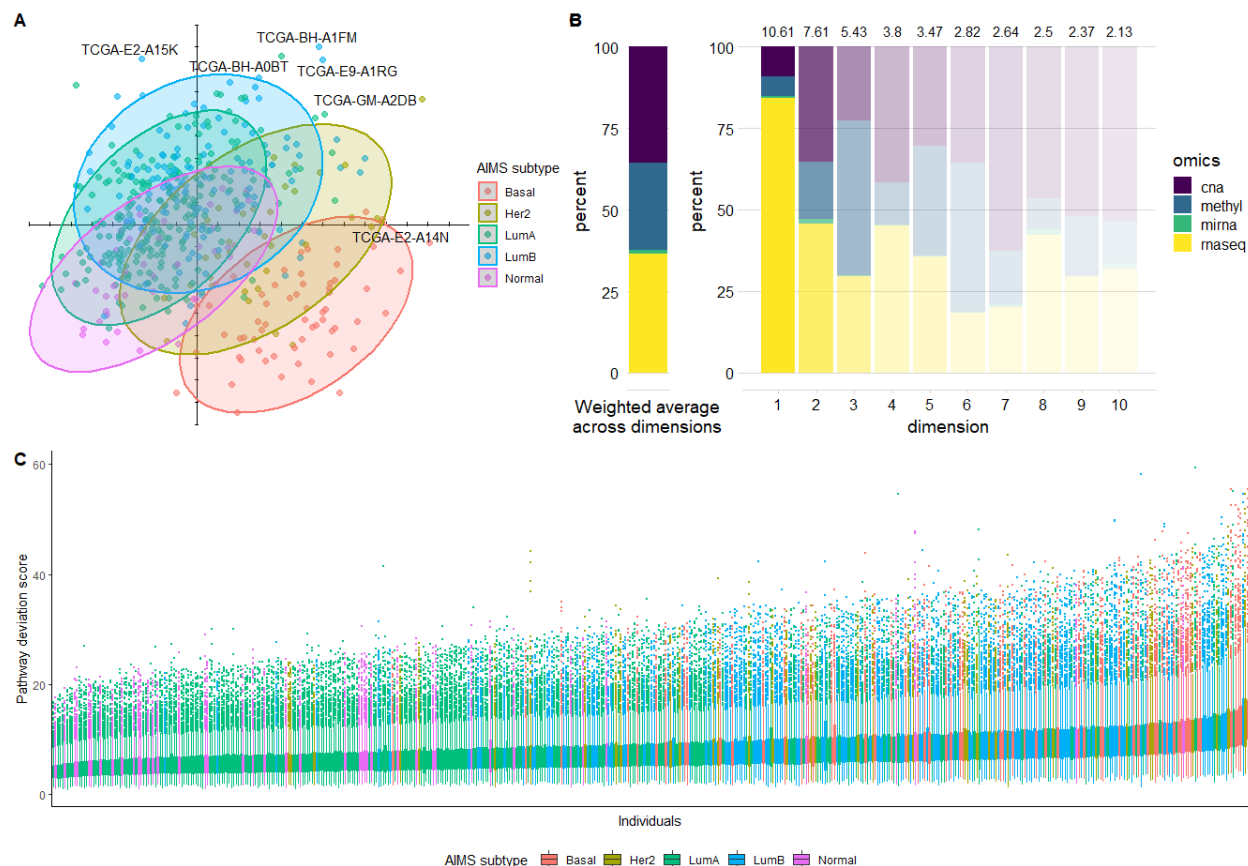
597  *Choice of curated pathway collection*

598

599  We consider the pathways included in the MSigDB canonical pathways curated gene set

600  catalog[38], which includes genes whose products are involved in metabolic and signaling pathways

601  reported in curated public databases. We specifically use the "C2 curated gene sets" catalog from

602  MSigDB v5.2 available at http://bioinf.wehi.edu.au/software/MSigDB/ as described in the *limma*

603  Bioconductor package[37]. We focus in particular on a collection of 1322 gene sets from public

604  databases, including Biocarta, Pathway Interaction Database[39], Reactome[40]; Sigma Aldrich,

605  Signaling Gateway, Signal Transduction Knowledge Environment, and the Matrisome Project[41],

606  the smallest and largest of which were respectively made up of 6 and 478 genes (median size 29

607  genes). For the subsequent *padma* analysis, we excluded gene sets for which fewer than 3 genes

608  mapped to quantified features in the TCGA gene expression data, corresponding to a total of

609  1136 gene sets.

610

611  *Padma R software package*

27

612    The proposed method described above has been implemented in an open-source R package

613    called *padma*, freely available on GitHub. *Padma* notably makes use *FactoMineR*[3,15] to run the

614    MFA; heatmaps in the following results were produced using *ComplexHeatmap*[42]. All of the

615    analyses in this paper were performed using R v3.5.1.

616



617

618    **Figure 5**. (A) Factor map of individuals, representing the first two components of the MFA,

619    for the Wnt signaling pathway in breast cancer, with normal confidence ellipses

620    superimposed for the five AIMS subtypes. B) Weighted overall percent contribution per

621    omic (left) and for each of the first 10 MFA components (right) for the Wnt signaling

622    pathway, with colors faded according to the percent variance explained for each

623    (represented in text above each bar). (C) Distribution of pathway deviation scores for each

624    individual in the breast cancer data, with individuals colored according to their AIMS

625    subtype.

626

**DECLARATIONS**

628

629     *Ethics approval and consent to participate*: Not applicable.

630     *Consent for publication*: Not applicable.

631     *Availability of data and materials*: The TCGA data analysed in the current study were retrieved

632     and pre-processed as described in the Methods section and in Rau et al. (2018)[5]; in particular, all

633     associated scripts can be found at https://github.com/andreamrau/EDGE-in-TCGA

634     (https://doi.org/10.5281/zenodo.3524080). All R scripts used to generate the results in this work

635     may be found at https://github.com/andreamrau/RMFRJLA_2019, and the associated *padma* R

636     package may be found at https://github.com/andreamrau/padma.

637     *Competing interests*: The authors declare that they have no competing interests.

641     *Authors' contributions*: AR conceived and designed the study, wrote the *padma* R package,

642     analyzed the data, and drafted the manuscript. RM analyzed the data and contributed to the R

643     package development. MJF and HR interpreted results and contributed to study design. FJ

644     contributed to the study conception and writing of the manuscript. DL supervised the study

645     conception and method implementation and drafted the manuscript and supplementary materials.

646     PLA conceived and designed the study and drafted the manuscript. All authors read and approved

647     the final manuscript.

648     *Acknowledgements*: Not applicable.

649 **REFERENCES**

650

651 1. The Cancer Genome Atlas Research Network *et al.* The Cancer Genome Atlas Pan-Cancer

652      analysis project. *Nat. Genet.* **45**, 1113–1120 (2013).

653 2. Meng, C. *et al.* Dimension reduction techniques for the integrative analysis of multi-omics

654      data. *Brief. Bioinform.* **17**, 628–641 (2016).

655 3. Husson, F., Lê, S. & Pagès, J. *Exploratory multivariate analysis by example using R.* (CRC

656      Press, 2017).

657 4. Argelaguet, R. *et al.* Multi- Omics Factor Analysis—a framework for unsupervised integration

658      of multi- omics data sets. *Mol. Syst. Biol.* **14**, (2018).

659 5. Rau, A., Flister, M., Rui, H. & Auer, P. L. Exploring drivers of gene expression in the Cancer

660      Genome Atlas. *Bioinformatics* (2018) doi:10.1093/bioinformatics/bty551.

661 6. Drier, Y., Sheffer, M. & Domany, E. Pathway-based personalized analysis of cancer. *Proc.*

662      *Natl. Acad. Sci.* **110**, 6388–6393 (2013).

663 7. Vaske, C. J. *et al.* Inference of patient-specific pathway activities from multi-dimensional

664      cancer genomics data using PARADIGM. *Bioinformatics* **26**, i237–i245 (2010).

665 9. Verbeke, L. P. C. *et al.* Pathway Relevance Ranking for Tumor Samples through Network-

666      Based Data Integration. *PLOS ONE* **10**, e0133503 (2015).

667 10.      Odom, G. J. *et al. pathwayPCA: an R package for integrative pathway analysis with*

668      *modern PCA methodology and gene selection.* http://biorxiv.org/lookup/doi/10.1101/615435

669      (2019) doi:10.1101/615435.

670 11.      de Tayrac, M., Le, S., Aubry, M., Mosser, J. & Husson, F. Simultaneous analysis of

671      distinct Omics data sets with integration of biological knowledge: Multiple Factor Analysis

672      approach. *BMC Genomics* **10**, 32 (2009).

673 12.      Meng, C. *et al.* MOGSA: Integrative Single Sample Gene-set Analysis of Multiple Omics

674  Data. *Mol. Cell. Proteomics* **18**, S153–S168 (2019).

675  13.  Escofier, B. & Pagès, J. *Analyses factorielles simples et multiples: objectifs, méthodes et*

676  *interprétation.* (2014).

677  14.  Pagès, J. *Multiple factor analysis by example using R.* (CRC Press, Taylor & Francis

678  Group, 2015).

679  15.  Lê, S., Josse, J. & Husson, F. **FactoMineR** : An *R* Package for Multivariate Analysis. *J.*

680  *Stat. Softw.* **25**, (2008).

681  16.  Abdi, H., Williams, L. J. & Valentin, D. Multiple factor analysis: principal component

682  analysis for multitable and multiblock data sets: Multiple factor analysis. *Wiley Interdiscip.*

683  *Rev. Comput. Stat.* **5**, 149–179 (2013).

684  17.  Paquet, E. R. & Hallett, M. T. Absolute assignment of breast cancer intrinsic molecular

685  subtype. *J. Natl. Cancer Inst.* **107**, 357 (2015).

686  18.  Riffo-Campos, Á., Riquelme, I. & Brebi-Mieville, P. Tools for Sequence-Based miRNA

687  Target Prediction: What to Choose? *Int. J. Mol. Sci.* **17**, 1987 (2016).

688  19.  Chou, C.-H. *et al.* miRTarBase update 2018: a resource for experimentally validated

689  microRNA-target interactions. *Nucleic Acids Res.* **46**, D296–D302 (2018).

690  20.  Liu, J. *et al.* An Integrated TCGA Pan-Cancer Clinical Data Resource to Drive High-

691  Quality Survival Outcome Analytics. *Cell* **173**, 400-416.e11 (2018).

692  21.  Heng, Y. J. *et al.* The molecular basis of breast cancer pathological phenotypes:

693  Molecular basis of breast cancer pathological phenotypes. *J. Pathol.* **241**, 375–391 (2017).

694  22.  Benjamini, Y. & Hochberg, Y. Controlling the False Discovery Rate: A Practical and

695  Powerful Approach to Multiple Testing. *J. R. Stat. Soc. Ser. B Methodol.* **57**, 289–300 (1995).

696  23.  Bourgon, R., Gentleman, R. & Huber, W. Independent filtering increases detection

697  power for high-throughput experiments. *Proc. Natl. Acad. Sci.* **107**, 9546–9551 (2010).

698  24.  Bosken, C. H., Wei, Q., Amos, C. I. & Spitz, M. R. An analysis of DNA repair as a

699  determinant of survival in patients with non-small-cell lung cancer. *J. Natl. Cancer Inst.* **94**,

700      1091–1099 (2002).

701   25.   Singhal, S., Vachani, A., Antin-Ozerkis, D., Kaiser, L. R. & Albelda, S. M. Prognostic

702      implications of cell cycle, apoptosis, and angiogenesis biomarkers in non-small cell lung

703      cancer: a review. *Clin. Cancer Res. Off. J. Am. Assoc. Cancer Res.* **11**, 3974–3986 (2005).

704   26.   Gonzalez-Perez, A. *et al.* IntOGen-mutations identifies cancer drivers across tumor

705      types. *Nat. Methods* **10**, 1081–1082 (2013).

706   27.   Zhang, Y., Rivera Rosado, L. A., Moon, S. Y. & Zhang, B. Silencing of D4-GDI inhibits

707      growth and invasive behavior in MDA-MB-231 cells by activation of Rac-dependent p38 and

708      JNK signaling. *J. Biol. Chem.* **284**, 12956–12965 (2009).

709   28.   Zhang, Y. & Zhang, B. D4-GDI, a Rho GTPase regulator, promotes breast cancer cell

710      invasiveness. *Cancer Res.* **66**, 5592–5598 (2006).

711   29.   Leek, J. T. *et al.* Tackling the widespread and critical impact of batch effects in high-

712      throughput data. *Nat. Rev. Genet.* **11**, 733–739 (2010).

713   30.   Novembre, J. *et al.* Genes mirror geography within Europe. *Nature* **456**, 98–101 (2008).

714   31.   Ahn, T., Lee, E., Huh, N. & Park, T. Personalized identification of altered pathways in

715      cancer using accumulated normal tissue data. *Bioinformatics* **30**, i422–i429 (2014).

716   32.   Wan, Y.-W., Allen, G. I. & Liu, Z. TCGA2STAT: simple TCGA data access for integrated

717      statistical analysis in R. *Bioinforma. Oxf. Engl.* **32**, 952–954 (2016).

718   33.   Li, B. & Dewey, C. N. RSEM: accurate transcript quantification from RNA-Seq data with

719      or without a reference genome. *BMC Bioinformatics* **12**, 323 (2011).

720   34.   Robinson, M. D. & Oshlack, A. A scaling normalization method for differential expression

721      analysis of RNA-seq data. *Genome Biol.* **11**, R25 (2010).

722   35.   Akulenko, R., Merl, M. & Helms, V. BEclear: Batch Effect Detection and Adjustment in

723      DNA Methylation Data. *PLOS ONE* **11**, e0159921 (2016).

724   36.   *MBatch: TCGA Batch Effects Viewer.* (2019).

725   37.   Ritchie, M. E. *et al.* limma powers differential expression analyses for RNA-sequencing

726   and microarray studies. *Nucleic Acids Res.* **43**, e47–e47 (2015).

727 38.  Liberzon, A. *et al.* Molecular signatures database (MSigDB) 3.0. *Bioinformatics* **27**,

728   1739–1740 (2011).

729 39.  Schaefer, C. F. *et al.* PID: the Pathway Interaction Database. *Nucleic Acids Res.* **37**,

730   D674–D679 (2009).

731 40.  Fabregat, A. *et al.* The Reactome Pathway Knowledgebase. *Nucleic Acids Res.* **46**,

732   D649–D655 (2018).

733 41.  Naba, A. *et al.* The Matrisome: *In Silico* Definition and *In Vivo* Characterization by

734   Proteomics of Normal and Tumor Extracellular Matrices. *Mol. Cell. Proteomics* **11**,

735   M111.014647 (2012).

736 42.  Gu, Z., Eils, R. & Schlesner, M. Complex heatmaps reveal patterns and correlations in

737   multidimensional genomic data. *Bioinformatics* **32**, 2847–2849 (2016).

738

739

**Supplementary Materials**

741

**Supplementary Figure 1**. Z-scores of RNA-seq, CNA, methylation, and miRNA-seq data for genes in the D4-GDI signaling pathway for individuals in the TCGA LUAD data (n = 144). Data corresponding to the two individuals with the largest overall pathway deviation scores, TCGA-78-7155 and TCGA-78-7536, are highlighted in red and blue.

746

**Supplementary Figure 2**. Negative log10-transformed p-values from the ANOVA F-test of pathway deviation score versus mitosis and nuclear pleomorphism for each pathway among breast cancer individuals. The signaling by Wnt pathway is highlighted in red.

750

**Supplementary Figure 3**. Factor maps for the first two dimensions of a global transcriptome- and genome-wide PCA of the methylation, miRNA-seq, CNA, and RNA-seq data (left), as well as a global MFA of all four omics combined (right) for the TCGA BRCA data.

754

**Supplementary Figure 4**. Factor maps for the first two dimensions of a global transcriptome- and genome-wide PCA of the methylation, miRNA-seq, CNA, and RNA-seq data (left), as well as a global MFA of all four omics combined (right) for the TCGA BRCA data.

758

**Supplementary Figure 5**. Percent variance explained by the first 5 (blue) or 10 (red) components of the MFA for each pathway for the TCGA BRCA (A) and LUAD (B) data.

761

**Supplementary Figure 6**. Average percent contribution to the MFA of each omic (miRNA-seq, methylation, CNA, RNA-seq) for each pathway. (A) Per-omic average contribution across the first

764    10 MFA components for TCGA BRCA. (B) Per-omic average contribution across all MFA

765    components for TCGA BRCA. (C) Per-omic average contribution across the first 10 MFA

766    components for TCGA LUAD. (D) Per-omic average contribution across all MFA components for

767    TCGA LUAD.

768

769    **Supplementary Table 1**. Sample size for each histological measure for the $n = 504$ breast cancer

770    patients.

771

772    **Supplementary Table 2**. Full gene lists for pathways in Table 1. Genes correspond to those with

773    expression quantified by RNA-seq in the TCGA data.

774

775    **Supplementary Table 3**. Full gene lists for pathways in Table 2. Genes correspond to those with

776    expression quantified by RNA-seq in the TCGA data.

777