1 TITLE:

2 Comprehensive cross-population analysis of high-grade serous ovarian cancer supports
3 no more than three subtypes

4

5 AUTHORS:

6 Gregory P. Way[a,b,c], James Rudd[c,d], Chen Wang[e], Habib Hamidi[f], Brooke L. Fridley[g],
7 Gottfried Konecny[f], Ellen L. Goode[e], Casey S. Greene[b,c,h,1], Jennifer A. Doherty[c,d,2]

8

9 AFFILIATIONS:

10 [a]Genomics and Computational Biology Graduate Program, University of Pennsylvania,
11 Philadelphia, PA 19103, USA
12 [b]Department of Systems Pharmacology and Translational Therapeutics, Perelman School
13 of Medicine, University of Pennsylvania, Philadelphia, PA
14 [c]Quantitative Biomedical Sciences, Norris Cotton Cancer Center, Geisel School of
15 Medicine at Dartmouth, Lebanon, NH
16 [d]Department of Epidemiology, Geisel School of Medicine at Dartmouth, Lebanon, NH
17 [e]Department of Health Sciences Research, Mayo Clinic, Rochester, MN
18 [f]Department of Medicine, David Geffen School of Medicine, University of California,
19 Los Angeles, CA
20 [g]Department of Biostatistics, University of Kansas Medical Center, Kansas City, KS
21 [h]Department of Genetics, Geisel School of Medicine at Dartmouth, Hanover, NH

22

23 CO-CORRESPONDING AUTHORS:

24 [1]10-131 SCTR 34th and Civic Center Blvd, Philadelphia, PA 19104; Phone: 215-573-
25 2991; Fax: 215-573-9135; csgreene@upenn.edu
26 [2]One Medical Center Drive, Lebanon, NH 03766; Phone: 603-653-9065; Fax: 603-653-
27 9093; Jennifer.A.Doherty@Dartmouth.edu

28

29 AUTHOR EMAIL ADDRESSES

30 GW: gregway@upenn.edu
31 JR: james.e.rudd.gr@dartmouth.edu
32 CW: Wang.Chen@mayo.edu
33 HH: HHamidi@mednet.ucla.edu
34 BF: bfridley@kumc.edu
35 GK: GKonecny@mednet.ucla.edu
36 EG: egoode@mayo.edu
37 CG: csgreene@upenn.edu
38 JD: Jennifer.A.Doherty@dartmouth.edu

39

40 CONFLICTS OF INTEREST:

41 The authors do not declare any conflicts of interest.

42

43 OTHER PRESENTATIONS:

44 Aspects of this study were presented at the 2015 AACR Conference and the 2015 Rocky
45 Mountain Bioinformatics Conference.

46

47 RUNNING HEAD:
48 *Cross-population analysis supports no more than three ovarian cancer subtypes*
49
50 KEYWORDS:
51 Ovarian Cancer; Molecular Subtypes; Unsupervised Clustering; Reproducibility
52
53 NOTES:
54 Words: 2,926; Figures: 3; Tables 3; Sup. Figures: 9; Sup. Tables: 6; Sup. Methods

55 AUTHORS' CONTRIBUTIONS
56 Study concept and design: GW, JR, CG, JD. Original data collection and processing:
57 CW, HH, BF, GK, EG. Data analysis: GW, JR, CG, JD. Manuscript drafting and editing:
58 GW, JR, CG, JD. All authors read, commented on, and approved the final manuscript.
59

60

61 ABSTRACT:

62 Four gene expression subtypes of high-grade serous ovarian cancer (HGSC) have been

63 described in several previous studies. In these studies, a fraction of samples that did not fit well

64 into any of the four subtype classifications were excluded. Therefore, we sought to

65 systematically determine the concordance of transcriptomic HGSC subtypes across populations

66 without removing any "hard-to-classify" samples. We created a unified bioinformatics pipeline

67 to independently cluster the five largest mRNA expression datasets using $k$-means and non-

68 negative matrix factorization (NMF). Within each population, we summarized differential

69 expression patterns, which we used to compare clusters across studies. While previous studies

70 reported four HGSC subtypes, our cross-population comparison does not support four subtypes.

71 Because these results contrast with previous reports, we attempted to reproduce the analyses

72 performed in those studies. Our results suggest that early results favoring four subtypes may

73 have been driven, at least in part, by the inclusion of serous borderline tumors. In summary, our

74 analysis suggests that either two or three, but not four, distinct gene expression subtypes are most

75 consistent with the available HGSC data to date.

76

INTRODUCTION:

Invasive ovarian cancer is a heterogeneous disease typically diagnosed at a late stage, with high mortality [1] . The most aggressive and common histologic type is high-grade serous (HGSC) [2], characterized by extensive copy number variation and *TP53* mutation [3]. Given the genomic complexity of these tumors, mRNA expression can be thought of as a summary measurement of these genomic and epigenetic alterations, to the extent that the alterations influence gene expression in either the cancer or stroma.

Four gene expression subtypes with varying components of mesenchymal, proliferative, immunoreactive, and differentiated gene expression signatures have been reported in all studies of HGSC to date [3–7]. Two of these also observed survival differences across subtypes [4,5]. Tothill *et al.* first identified four HGSC subtypes (as well as two other subtypes which largely included low grade serous and serous borderline tumors) in an Australian population using *k*-means clustering. Later, The Cancer Genome Atlas (TCGA) used non-negative matrix factorization (NMF) and also reported four subtypes which were labeled as: 'mesenchymal', 'differentiated', 'proliferative', and 'immunoreactive' [3]. The TCGA group also applied NMF clustering to the Tothill data, and observed concordance with four subtypes [3]. Konecny *et al.* applied NMF to cluster an independent set of HGSC samples and reported four subtypes, which they labeled as C1-C4 [5]. These subtypes were similar to those in the TCGA but a subtype classifier trained on these subtypes better differentiated survival in their own data, and in data from TCGA and Bonome *et al.* [6].

Despite this extensive research in the area, work to date has several limitations. In both TCGA and Tothill *et al.*, ~8-15% of samples were excluded from analyses. A reanalysis of the TCGA data showed that over 80% of the samples could be assigned to more than one subtype

3

100    [8]. In more recent TCGA analyses by the Broad Institute Genome Data Analysis Center

101    (GDAC) Firehose initiative with the largest number of HGSC cases evaluated to date (n = 569),

102    three subtypes fit the data better than four [9,10]. This uncertainty in HGSC subtyping led us to

103    determine if four homogeneous subtypes exist across study populations.

104         To comprehensively characterize subtypes, we analyze data from the five largest

105    independent studies to date, including our own collection of samples, using a standardized

106    bioinformatics pipeline. We apply $k$-means clustering as well as NMF to each population without

107    removing "hard-to-classify" samples. Our goal is to rigorously assess the number of subtypes.

108    These independent and parallel within-dataset analyses followed by cross-dataset comparison

109    sidestep gene expression platform or dataset biases that could affect clustering if under or

110    overcorrected. This contrasts with earlier work that pooled datasets together to identify subtypes

111    [7] and ensures that subtypes identified are not induced by dataset or batch effects. We

112    summarize each subtype's expression patterns and comprehensively characterize correlations

113    between subtype-specific gene expression across populations.

114         Our cross-population comparative analysis does not support that four HGSC subtypes

115    exist; rather the data more strongly support an interpretation that there are either two or three

116    subtypes. We show that the support for four subtypes observed in TCGA's reanalysis of Tothill

117    *et al.* [3] is lost when serous borderline tumors, which have very different genomic profiles and

118    survival than HGSC [11,12], are excluded before clustering. Our work also highlights the impact

119    that a single study can have on the trajectory of subtyping research and suggests the importance

120    of periodic histopathologic review and rigorous reanalysis of existing data for cross-study

121    commonalities.

122

123    METHODS:

124    *Data inclusion*

125    We applied inclusion criteria as described in the supplementary materials using data from

126    the R package, curatedOvarianData [13] and our own novel dataset ("Mayo") [5]

127    (Supplementary Table S1). These criteria selected HGCS samples that were not duplicates from

128    studies including at least 130 HGSC cases assayed on standard microarrays. Data from the new

129    Mayo HGSC samples as well as other samples with mixed histologies and grades, for a total of

130    528 additional ovarian tumor samples, was deposited in NCBI's Gene Expression Omnibus

131    (GEO) [14]; these data can be accessed with the accession number GSE74357

132    (http://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE74357).  All study participants

133    provided written informed consent, and this work was approved by the Mayo Clinic and

134    Dartmouth College Institutional Review Boards.

135    After applying the unified inclusion criteria, our final analytic datasets include: TCGA (n

136    = 499) [3,9]; Mayo (n = 379; GSE74357) [5]; Yoshihara (n = 256; GSE32062.GPL6480) [15];

137    Tothill (n = 242; GSE9891) [4]; and Bonome (n = 185; GSE26712) [6] (Table 1). We restricted

138    analyses to the 10,930 genes measured successfully in all five populations (Supplementary Fig.

139    S1).

140

141    *Clustering*

142    We performed independent clustering within each dataset to avoid potential biases from

143    different platforms or studies.  As detailed in the Supplementary Methods, we identified the

144    1,500 genes with the highest variance from each dataset and used the union of these genes (n =

145    3,698) for clustering. We performed clustering within each dataset using each potential $k$ from 2-

146    8 clusters. We performed *k*-means clustering in each population using the R package "cluster"

147    (version 2.0.1) [16] with 20 initializations. We repeated these analyses using NMF in the R

148    package "NMF" (version 0.20.5) [17] with 100 different random initializations for each *k*. As

149    done in prior studies, we calculated cophenetic correlation coefficients to select appropriate *k* for

150    each dataset after NMF clustering with 10 consensus runs for *k* = 2 through 8.

151

152    *Identification of analogous clusters within and across studies*

153        We performed significance analysis of microarray (SAM) [18,19] analysis on all clusters

154    from each study using all 10,930 genes. This resulted in a cluster-specific moderated *t* statistic

155    for each of the input genes [20]. To summarize the expression patterns of all 10,930 genes for a

156    specific cluster in a specific population, we combined gene-wise moderated *t* statistics into a

157    vector of length 10,930. The TCGA subtype labels have become widely used in the field. To

158    generate comparable labels across *k* and across studies, we mapped our TCGA subtype

159    assignments back to the original TCGA labels to define reference clusters at *k* = 4 (that is,

160    mesenchymal-like, proliferative-like, etc.). Clusters in other populations that were most strongly

161    correlated with the TCGA clusters were assigned the same label.

162

163    *Clustering analysis of randomized data*

164        Any clustering procedure is expected to induce strong correlational structure across

165    clusters within a dataset even if there is no true underlying structure. However, if there is no true

166    underlying structure, clusters across datasets are not expected to be correlated. To assess this, we

167    used the same datasets but shuffled each gene's expression vector to disrupt the correlative

6

168   structure. We performed within and cross-study analyses of cluster identification using this set of

169   data that were parallel to those performed using the non-randomized data.

170

171   *Assessing the reproducibility of single-population studies*

172       We compared our sample assignments at $k = 2 - 4$ to the four subtypes reported in the

173   Tothill, TCGA, and Konecny publications [3–5]. Because the labels that were assigned in

174   TCGA's reanalysis of the Tothill data were not available, we performed NMF consensus

175   clustering of Tothill's data without removing LMP samples in order to generate labels for

176   comparison.

177

178   *Reproducibility of our analyses*

179       We provide software to download the required data and reproduce our analyses. The

180   software is provided under a permissive open source license [21]. Analyses were run in a Docker

181   container, allowing the computing environment to be recreated [22]. Our Docker image can be

182   pulled from here: https://hub.docker.com/r/gregway/hgsc_subtypes/. This allows interested users

183   to freely download the software, reproduce the analyses, and then build on this work.

184

185   RESULTS:

186   *Clustering*

187       To visually inspect the consistency and distinctness of clusters, we compared sample-by-

188   sample correlation heatmaps. For $k = 2$ to 4 within each study, we observed high sample-by-

189   sample correlations within clusters and relatively low sample-by-sample correlations across

190    clusters (Supplementary Fig. S2). Clustering results using NMF were similar to $k$ means results

191    (Supplementary Fig S3.)

192

193    *Correlation of cluster-specific expression patterns*

194          Across datasets, we observed strong positive correlations of moderated $t$ score vectors

195    between analogous clusters in TCGA, Tothill, Mayo, and Yoshihara (Fig. 1; Table 2). However,

196    clustering of the Bonome data did not correlate strongly with clusters identified in the other

197    datasets (Table 2). We believe that we were unable to assign parallel subtypes in Bonome

198    because of either RNA contamination or inappropriate grading assignments. However, more

199    work is required in order to identify exactly why we were unable to classify.

200          To assess our analytical approach, we performed an analysis using randomized data. This

201    showed that within-population correlation structure was induced by clustering, but structure

202    between populations was not (Supplementary Fig. S4). Comparing Figure 1 with S4, we

203    observed much higher correlation across datasets (Fig. 1), which was lost after randomization

204    (Supplementary Fig. S4). For example, for $k = 2$, the TCGA and Mayo cluster correlations for

205    analogous clusters was high (top left panel in Fig. 1). Conversely, the same relationship in

206    randomized data (second row, first column panel in Supplementary Fig. S4) showed correlations

207    near zero. This indicates that the high correlations observed across datasets in Figure 1 are

208    induced by similar underlying structure in the data.

209          Across studies, positive correlations between analogous clusters and negative correlations

210    between non-analogous clusters were stronger for clusters identified when $k = 2$ and $k = 3$ than

211    when $k = 4$ (Fig. 1), with comparable statistical precision (Supplementary Table S2). These

212  cross-population comparisons suggested that two and three subtypes fit HGSC gene expression

213  data more consistently than the four widely accepted subtypes.

214      Within each population, clusters identified by NMF were very similar to those identified

215  using $k$-means clustering (Fig. 2) suggesting that these results were independent of clustering

216  algorithm. With NMF, both positive and negative correlations were stronger for $k = 2$ and $k = 3$

217  than for $k = 4$. Across $k = 3$ and $k = 4$, correlations were strongest for clusters 1 and 2. Sample

218  cluster assignments for both $k$-means and NMF clusters are provided in Supplementary Table S3.

219

220  *Comparison with previously-identified HGSC clusters*

221      Our clustering results for the Tothill, TCGA, and Mayo datasets were highly concordant

222  with the clustering described in the original publications [3–5], as evidenced by the high degree

223  of consistent overlap in sample assignments to the previously-defined clusters (Table 3). Our

224  cross-study cluster 1 was mostly mapped to the "Mesenchymal" label from TCGA, "C1" from

225  Tothill, and "C4" from Mayo. This cluster was the most stable in our analysis within all datasets,

226  across $k = 2$, 3 and 4, and across clustering algorithms. Cross-study cluster 2, which was also

227  observed consistently, was most similar to the "Proliferative" label from TCGA, "C5" from

228  Tothill, and "C3" from Mayo. Cross-study cluster 3 for $k = 3$ was associated with both the

229  "Immunoreactive" and "Differentiated" TCGA labels, "C2" and "C4" in Tothill, and "C1" and

230  "C2" in Mayo. For analyses where $k = 4$, the third cluster was associated with

231  "Immunoreactive", "C2", and "C1" while the fourth cluster was associated with "Differentiated",

232  "C4", and "C2" for TCGA, Tothill, and Mayo respectively.

233

234  *Meta-research into previous HGSC subtyping studies*

9

235 Each of the publications that only considered high-grade samples (TCGA and Konecny *et*

236 *al.*) found clustering coefficients consistent with $k = 2$, $k = 3$, and $k = 4$. Nevertheless, each

237 publication concludes the existence of four subtypes, while our cross-population analysis

238 suggested that two or three clusters fit HGSC data better than four clusters.

239 To compare with previous results, we evaluated the number of subtypes that fit the data

240 best within each study by calculating cophenetic correlation coefficients at $k = 2$ through k=8

241 clusters inclusively. We observed a similar pattern in each population (Supplementary Fig S5 –

242 S7; Fig. 3A) in which the highest cophenetic correlation was reached for two clusters and, based

243 on the heatmaps, appeared to have the highest consensus. In every dataset, four clusters were not

244 observed to represent the data better than two or three. The only results in previous studies that

245 contradicted this work were from TCGA's reanalysis of the Tothill data. According to

246 supplemental figure S6.2 in the TCGA paper, the reanalysis included serous borderline tumors

247 (i.e., tumors with low malignant potential) (n = 18). The inclusion of these tumors in the TCGA

248 HGSC analyses was done even though, in the original Tothill paper, the serous borderline tumors

249 had a unique gene expression patterns and clustered entirely in a group labeled "C3".

250 To assess the extent to which serous borderline tumors inclusion drove the TCGA results,

251 we reproduced TCGA's reanalysis of Tothill *et al.*, including the serous borderline tumors (n =

252 18); we indeed observed that the cophenetic correlation is higher for $k = 4$ than $k = 3$ (Fig. 3A).

253 However, when we appropriately removed these serous borderline tumors we observed an

254 increase in the $k = 3$ cophenetic correlation (Fig. 3B). The results that support four subtypes were

255 generated during clustering of HGSC and serous borderline tumors combined. Subtyping

256 analyses of HGSC alone reveal less than four subtypes. Even after subtyping there remains a

257 complex and nuanced portrait of the disease.

258

259    DISCUSSION:

260    Although prior studies have reported the existence of four molecular subtypes of HGSC

261    ovarian cancer [3–5,9], our analysis suggests the existence of only two or three subtypes. This

262    conclusion is based on our observation that concordance of analogous subtypes across study

263    populations was stronger for two or three clusters as opposed to four. Previous studies used

264    either $k$-means or NMF clustering, and because our results contradicted prior work, we

265    performed analyses using both of these methods. Results for each population were similar for the

266    $k$ means and NMF clustering algorithms suggesting that the clustering algorithm did not drive the

267    observed differences.

268    Because cross-population comparisons suggest that two and three clusters show more

269    consistency than four, we explored within-study heuristics (cophenetic correlation coefficients)

270    that suggested four subtypes in previous research. The cophenetic coefficient measures how

271    precisely a dendrogram retains sample by sample pairwise distances and can be used to compare

272    clustering accuracy [23]. While both Konecny and TCGA reported four subtypes, in both

273    analyses $k = 2$ and $k = 3$ resulted in  higher cophenetic coefficients than $k = 4$ (Konecny Figure

274    2A and TCGA Figure S6.1) [3, 5]. We observed the same patterns in our own reanalysis of

275    TCGA and analysis of the expanded Mayo cohort (Supplementary Figs. S5 and S6). Yoshihara

276    and Tothill did not report cophenetic coefficients, but our analysis of each (Supplementary Fig

277    S7 and Fig 3A) revealed similar patterns to TCGA and Konecny.

278    In the previous literature, the only report to suggest that three subtypes were

279    inappropriate was TCGA's reanalysis of the Tothill *et al.* data (supplemental Figure S6.2 in their

280    publication); the cophenetic coefficient dropped dramatically at $k = 3$ before recovering at $k = 4$

11

281    [3]. Notably, TCGA's figure legend for this supplemental result indicates that they did not

282    remove serous borderline tumors from the Tothill data. Our analysis of Tothill *et al.* differed

283    from TCGA's in that we excluded serous borderline tumors and instead supports the existence of

284    two or three subtypes. To evaluate the influence of these serous borderline tumors in the Tothill

285    data, we repeated our analyses including serous borderline tumors, and observed a drop in the

286    cophenetic coefficient for $k = 3$ relative to $k = 4$ (Fig. 3). This suggests that the four subtypes

287    observed in TCGA's analysis of the Tothill data may be due, in part, to the inclusion of serous

288    borderline tumors.

289        There are several limitations to note in the HGSC data we analyzed. Given the intra-

290    tumor heterogeneity that is likely to exist [24], our approach would be strengthened by having

291    data on multiple areas of the tumors. Additionally, since histology and grade classification have

292    changed over time [25,26], it is unclear whether the populations we studied used comparable

293    guidelines to determine histology and grade. We attempted to exclude all low grade serous and

294    low grade endometrioid samples because they often have very different gene expression patterns

295    and more favorable survival compared to their higher grade counterparts [2]. While the Bonome

296    publication specified that they included only high-grade tumors, grade is not included in the

297    Bonome GSE26712 data set, so we were unable to determine whether the grade distribution

298    differs from the other studies [6]. It is unclear why the Bonome clusters did not correspond to the

299    clusters observed in other populations. Lack of consistency could result from a different

300    distribution of grade or other unreported biological differences.

301        In summary, our study demonstrates that two clusters of HGSC, "mesenchymal-like" and

302    "proliferative-like", are clearly and consistently identified within and between populations. This

303    suggests that there are two reproducible HGSC subtypes that are either etiologically distinct, or

12

304     acquire phenotypically determinant alterations through their development. Our study also

305     suggests that the previously described "immunoreactive-like" and "differentiated-like" subtypes

306     appear more variable across populations, and tend to be collapsed into a single category when

307     three subtypes are specified. These may represent, for example, steps along an immunoreactive

308     continuum or could represent the basis of a third, but more variable subtype.

309          Our analysis also reveals the importance of critically reassessing molecular subtypes

310     across multiple large study populations using parallel analyses and consistent inclusion criteria.

311     New systematic approaches hold promise for the implementation of such analyses [27]. Our

312     results underscore the importance of ovarian cancer histopathology, contradict the four HGSC

313     subtype hypothesis, and suggest that there may be fewer HGSC molecular subtypes with variable

314     immunoreactivity and stromal infiltration.

315

326    American Cancer Society (grant number IRG 8200327 to C.S.G.), and by Norris Cotton Cancer

327    Center Developmental Funds.

328

329    FIGURE LEGENDS:

330

331    **Figure 1.** Significance analysis of microarray (SAM) moderated $t$ score Pearson correlation

332    heatmaps reveal consistency across datasets. (A) Correlations across datasets for $k$ means $k = 2$.

333    (B) Correlations across datasets for $k$ means $k = 3$. (C) Correlations across datasets for $k$ means $k$

334    $= 4$

335

336    **Figure 2.** Significance analysis of microarray (SAM) moderated $t$ score Pearson correlation

337    heatmaps of clusters formed by $k$ means clustering and NMF clustering reveals consistency

338    across clustering methods. Within dataset results are shown for both methods when setting each

339    algorithm to find 2, 3, and 4 clusters.

340

341    **Figure 3.** Comparing NMF consensus clustering in the Tothill dataset**.** Data displays consensus

342    clustering for $k = 2$ to $k = 6$ for 10 NMF initializations alongside the cophenetic correlation

343    results for $k = 2$ to $k = 8$. (A) Tothill dataset (n = 260) with low malignant potential (LMP)

344    samples (n = 18) not removed prior to clustering. (B) Tothill dataset with LMP samples removed

345    (n = 242).

346

347    **Supplementary Figure S1.** Overlapping genes assayed using either the HG-U1133 Affymetrix

348    platform (TCGA, Tothill, Bonome) or the Agilent 4x44K platform (Mayo, Yoshihara).

14

349    Differences across datasets arise from inherent array differences and/or differences in quality

350    control preprocessing.

351

352    **Supplementary Figure S2.** Sample by sample Pearson correlation matrices. Top panel: $k = 2$.

353    Middle panel: $k = 3$. Bottom panel: $k = 4$. The color bars are coded as blue for cluster 1, red for

354    cluster 2, green for cluster 3, and purple for cluster 4. In the matrices, red represents high

355    correlation, blue low correlation, and white intermediate correlation. The scales are slightly

356    different in each population because of different correlational structures. The clusters in the

357    Bonome study are depicted in gray scale because in cross-population analyses to identify

358    analogous clusters, those from Bonome did not correlate with those observed in the four other

359    studies.

360

361    **Supplementary Figure S3.** NMF consensus matrices for datasets when $k = 2$, $k = 3$, and $k = 4$.

362    The first track represents cluster membership for $k$ means clusters and the second track

363    represents silhouette widths. Note that NMF clusters are not ordered in the same way as the $k$

364    means clusters.

365

366    **Supplementary Figure S4.** Significance analysis of microarray (SAM) moderated $t$ score

367    Pearson correlation heatmaps are not consistent across datasets for randomly shuffled gene

368    expression values for $k = 2$, $k = 3$, or $k = 4$. The within dataset correlations are artificially

369    induced because the clustering algorithm will find clusters even without true underlying

370    structure. However, the across dataset clusters are not correlated in the randomized data

371    indicating that the results we observe in Figure 1 are not artifacts of the clustering algorithm.

15

372

373     **Supplementary Figure S5.** Consensus NMF clustering of the TCGA dataset (n = 499) for $k = 2$

374     to $k = 6$ for 10 NMF runs alongside the cophenetic correlation results for $k = 2$ to $k = 8$.

375

376     **Supplementary Figure S6.** Consensus NMF clustering of the Mayo dataset (n = 379 for $k = 2$ to

377     $k = 6$ for 10 NMF runs alongside the cophenetic correlation results for $k = 2$ to $k = 8$.

378

379     **Supplementary Figure S7.** Consensus NMF clustering of the Yoshihara dataset (n = 256) for $k$

380     $= 2$ to $k = 6$ for 10 NMF runs alongside the cophenetic correlation results for $k = 2$ to $k = 8$.

381

382     **Supplementary Figure S8.** Silhouette width plots for $k = 2$, $k = 3$, and $k = 4$ for $k$ means

383     clustering results. Cluster 1 is shown in blue, cluster 2 in red, cluster 3 in green, and cluster 4 in

384     purple.

385

386     **Supplementary Figure S9.** Kaplan-Meier survival curves for $k = 2$, $k = 3$, and $k = 4$ shown for

387     clustering solutions using $k$ means and NMF. Cluster 1 is shown in blue, cluster 2 in red, cluster

388     3 in green, and cluster 4 in purple.

389

390     REFERENCES:

391     1.    Kurman RJ, Shih I-M. The Origin and Pathogenesis of Epithelial Ovarian Cancer: A
392           Proposed Unifying Theory: Am J Surg Pathol. 2010;34: 433–443.
393           doi:10.1097/PAS.0b013e3181cf3d79

394     2.    Vang R, Shih I-M, Kurman RJ. Ovarian Low-grade and High-grade Serous Carcinoma:
395           Pathogenesis, Clinicopathologic and Molecular Biologic Features, and Diagnostic
396           Problems. Adv Anat Pathol. 2009;16: 267–282. doi:10.1097/PAP.0b013e3181b4fffa

397   3.   The Cancer Genome Atlas. Integrated genomic analyses of ovarian carcinoma. Nature.
398        2011;474: 609–615. doi:10.1038/nature10166

399   4.   Tothill RW, Tinker AV, George J, Brown R, Fox SB, Lade S, et al. Novel Molecular
400        Subtypes of Serous and Endometrioid Ovarian Cancer Linked to Clinical Outcome. Clin
401        Cancer Res. 2008;14: 5198–5208. doi:10.1158/1078-0432.CCR-08-0196

402   5.   Konecny GE, Wang C, Hamidi H, Winterhoff B, Kalli KR, Dering J, et al. Prognostic and
403        Therapeutic Relevance of Molecular Subtypes in High-Grade Serous Ovarian Cancer. JNCI
404        J Natl Cancer Inst. 2014;106: dju249–dju249. doi:10.1093/jnci/dju249

405   6.   Bonome T, Levine DA, Shih J, Randonovich M, Pise-Masison CA, Bogomolniy F, et al. A
406        gene signature predicting for survival in suboptimally debulked patients with ovarian
407        cancer. Cancer Res. 2008;68: 5478–5486. doi:10.1158/0008-5472.CAN-07-6595

408   7.   Tan TZ, Miow QH, Huang RY-J, Wong MK, Ye J, Lau JA, et al. Functional genomics
409        identifies five distinct molecular subtypes with clinical relevance and pathways for growth
410        control in epithelial ovarian cancer: A subtyping scheme for epithelial ovarian cancer.
411        EMBO Mol Med. 2013;5: 1051–1066. doi:10.1002/emmm.201201823

412   8.   Verhaak RGW, Tamayo P, Yang J-Y, Hubbard D, Zhang H, Creighton CJ, et al.
413        Prognostically relevant gene signatures of high-grade serous ovarian carcinoma. J Clin
414        Invest. 2012; doi:10.1172/JCI65833

415   9.   Broad Institute TCGA Genome Data Analysis Center. Analysis Overview for Ovarian
416        Serous Cystadenocarcinoma (Primary solid tumor cohort) - 15 July 2014. 2014;
417        doi:10.7908/C1QR4VX4

418   10.  Broad Institute TCGA Genome Data Analysis Center. Clustering of mRNA expression:
419        consensus NMF. 2015; doi:10.7908/C1BR8R71

420   11.  Ouellet V, Provencher DM, Maugard CM, Le Page C, Ren F, Lussier C, et al.
421        Discrimination between serous low malignant potential and invasive epithelial ovarian
422        tumors using molecular profiling. Oncogene. 2005;24: 4672–4687.
423        doi:10.1038/sj.onc.1208214

424   12.  Bonome T, Lee J-Y, Park D-C, Radonovich M, Pise-Masison C, Brady J, et al. Expression
425        profiling of serous low malignant potential, low-grade, and high-grade tumors of the ovary.
426        Cancer Res. 2005;65: 10602–10612. doi:10.1158/0008-5472.CAN-05-2240

427   13.  Ganzfried BF, Riester M, Haibe-Kains B, Risch T, Tyekucheva S, Jazic I, et al.
428        curatedOvarianData: clinically annotated data for the ovarian cancer transcriptome.
429        Database. 2013;2013: bat013–bat013. doi:10.1093/database/bat013

430   14.  Edgar R, Domrachev M, Lash AE. Gene Expression Omnibus: NCBI gene expression and
431        hybridization array data repository. Nucleic Acids Res. 2002;30: 207–210.

432 15. Yoshihara K, Tsunoda T, Shigemizu D, Fujiwara H, Hatae M, Fujiwara H, et al. High-Risk
433    Ovarian Cancer Based on 126-Gene Expression Signature Is Uniquely Characterized by
434    Downregulation of Antigen Presentation Pathway. Clin Cancer Res. 2012;18: 1374–1385.
435    doi:10.1158/1078-0432.CCR-11-2725

436 16. Maechler M, Rousseeuw P, Struyf A, Hubert M, Hornik K. cluster: Cluster Analysis Basics
437    and Extensions. 2014;R package version 1.15.3.

438 17. Brunet J-P, Tamayo P, Golub TR, Mesirov JP. Metagenes and molecular pattern discovery
439    using matrix factorization. Proc Natl Acad Sci. 2004;101: 4164–4169.
440    doi:10.1073/pnas.0308531101

441 18. Tusher VG, Tibshirani R, Chu G. Significance analysis of microarrays applied to the
442    ionizing radiation response. Proc Natl Acad Sci. 2001;98: 5116–5121.
443    doi:10.1073/pnas.091062498

444 19. Schwender H, Krause A, Ickstadt K. Identifying interesting genes with sigenes. RNews.
445    2006;6: 45–50.

446 20. Schwender H. siggenes: Multiple testing using SAM and Efron's empirical Bayes
447    approaches. 2012;R package version 1.40.0.

448 21. Gregory Way, James Rudd, Casey Greene. Analytical Code for "Cross-population analysis
449    of high-grade serous ovarian cancer reveals only two robust subtypes." 2015;
450    doi:10.5281/zenodo.32906

451 22. Boettiger C. An introduction to Docker for reproducible research. ACM SIGOPS Oper Syst
452    Rev. 2015;49: 71–79. doi:10.1145/2723872.2723882

453 23. Sokal RR, Rohlf FJ. The Comparison of Dendrograms by Objective Methods. Taxon.
454    1962;11: 33. doi:10.2307/1217208

455 24. Blagden SP. Harnessing Pandemonium: The Clinical Implications of Tumor Heterogeneity
456    in Ovarian Cancer. Front Oncol. 2015;5. doi:10.3389/fonc.2015.00149

457 25. Silverberg SG. Histopathologic grading of ovarian carcinoma: a review and proposal. Int J
458    Gynecol Pathol Off J Int Soc Gynecol Pathol. 2000;19: 7–15.

459 26. Soslow RA. Histologic Subtypes of Ovarian Carcinoma: An Overview. Int J Gynecol
460    Pathol. 2008;PAP. doi:10.1097/PGP.0b013e31815ea812

461 27. Planey CR, Gevaert O. CoINcIDE: A framework for discovery of patient subtypes across
462    multiple datasets. Genome Med. 2016;8. doi:10.1186/s13073-016-0281-4

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482     **Table 1**: Characteristics of the populations included in the five analytic data sets

| | Bonome *et al.* | Tothill *et al.* |
|---|---|---|
| | GSE26712 | GSE9891 |
| | Affy HGU1133 | Affy HGU1133 |
| | United States | Australia |
| | 195 | 285 |
| | 185 | 242 |
| | 61.5 (11.9) | 60.3 (10.3) |
| | 0 (0%) | 11 (5%) |
| | 0 (0%) | 8 (4%) |
| | 146 (80%) | 178 (83%) |
| | 36 (20%) | 17 (8%) |
| | NA | 80 (37%) |
| | NA | 134 (63%) |
| | 89 (49%) | 132 (62%) |
| | 93 (51%) | 82 (38%) |

19

| | TCGA | Mayo | Yoshihara et al. |
|---|---|---|---|
| GEO | | GSE74357 | GSE32062 |
| Platform | Affy HGU1133 | Agilent 4x44K | Agilent 4x44K |
| Population | United States | United States | Japan |
| Original Sample Size | 578 | 528 | 260 |
| Analytic Sample Size[b] | 499 | 379 | 256 |
| Age [Mean (SD)] | 60.0 (11.6) | 62.9 (11.3) | NA |
| Stage | | | |
| I | 10 (2%) | 7 (3%) | 0 (0%) |
| II | 17 (4%) | 11 (3%) | 0 (0%) |
| III | 351 (80%) | 275 (73%) | 202 (79%) |
| IV | 63 (14%) | 86 (23%) | 54 (21%) |
| Grade | | | |
| 2 | 55 (12%) | 3 (1%) | 130 (51%) |
| 3 | 386 (88%)[a] | 376 (99%) | 126 (49%) |
| Debulking | | | |
| Optimal | 325 (74%) | 287 (76%) | 101 (39%) |
| Suboptimal | 116 (26%) | 87 (23%) | 155 (61%) |

483

484   NA: Data not reported

485   [a]One sample was labeled as 'Grade 4' in TCGA

486   [b]samples without survival data were excluded in survival analyses

487   **Table 2**: SAM moderated *t* score vector Pearson correlations between analogous clusters across

488   populations[a]

| | Cluster 1 | Cluster 2 | Cluster 3 | Cluster 4 |
|---|---|---|---|---|
| $k = 2$[a] | 0.62 – 0.81 | 0.62 – 0.81 | NA | NA |
| $k = 3$[a] | 0.77 - 0.85 | 0.80 - 0.90 | 0.65 - 0.77 | NA |
| $k = 4$[a] | 0.77 - 0.85 | 0.83 - 0.89 | 0.51 - 0.76 | 0.61 - 0.75 |

| | | | | |
|---|---|---|---|---|
| Bonome $k = 2^b$ | -0.08 – 0.24 | -0.08 – 0.24 | NA | NA |
| Bonome $k = 3^b$ | 0.45 – 0.46 | -0.02 - 0.12 | 0.22 - 0.42 | NA |
| Bonome $k = 4^b$ | 0.50 - 0.57 | -0.04 - 0.04 | 0.13 - 0.29 | 0.26 - 0.43 |

489 [a]Correlation ranges for TCGA, Mayo, Yoshihara, and Tothill.

490 [b]Bonome is removed from gene set analyses because of low correlating clusters

491

492

493

494

495

496

497

498

499

500

501 **Table 3:** Distributions of sample membership in the clusters identified in our study by the

502 original cluster assignments in the TCGA, Tothill, and Konecny studies. Clusters identified in

503 our study using $k$-means clustering with $k = 2$, $k = 3$, and $k = 4$

| | TCGA | | | | | Tothill *et al.* | | | | | | | Konecny *et al.* | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $k = 2$ | Mes | Pro | Imm | Dif | NC[a] | C1 | C2 | C3 | C4 | C5 | C6 | NC[a] | C1 | C2 | C3 | C4 | NA[b] |
| Cluster 1 | 98 | 7 | 93 | 68 | 21 | 78 | 39 | 1 | 0 | 0 | 0 | 11 | 36 | 21 | 2 | 26 | 114 |
| Cluster 2 | 1 | 127 | 2 | 60 | 22 | 0 | 5 | 5 | 44 | 35 | 2 | 22 | 6 | 39 | 41 | 0 | 94 |
| $k = 3$ | | | | | | | | | | | | | | | | | |
| Cluster 1 | 98 | 2 | 20 | 11 | 6 | 77 | 22 | 0 | 0 | 0 | 0 | 6 | 16 | 13 | 2 | 26 | 82 |
| Cluster 2 | 1 | 111 | 0 | 11 | 16 | 1 | 0 | 0 | 3 | 35 | 2 | 5 | 0 | 16 | 36 | 0 | 56 |
| Cluster 3 | 0 | 21 | 75 | 106 | 21 | 0 | 22 | 6 | 41 | 0 | 0 | 22 | 26 | 31 | 5 | 0 | 70 |
| $k = 4$ | | | | | | | | | | | | | | | | | |
| Cluster 1 | 97 | 4 | 12 | 12 | 5 | 74 | 0 | 0 | 0 | 0 | 0 | 0 | 7 | 12 | 3 | 25 | 62 |
| Cluster 2 | 1 | 85 | 0 | 0 | 13 | 1 | 0 | 0 | 1 | 34 | 2 | 5 | 0 | 9 | 31 | 0 | 41 |

| | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Cluster 3 | 0 | 5 | 80 | 3 | 12 | 3 | 42 | 0 | 1 | 1 | 0 | 14 | 29 | 6 | 0 | 1 | 57 |
| Cluster 4 | 1 | 40 | 3 | 113 | 13 | 0 | 2 | 6 | 42 | 0 | 0 | 14 | 6 | 33 | 9 | 0 | 48 |

504

505 [a]NC = Samples not clustered in original publication

506 [b]NA = Samples not assessed at the time of the original publication

507 NOTE: The corresponding labels for the generally similar HGSC gene expression subtypes

508 observed in the TCGA, Tothill, and Konecny studies are, respectively: mesenchymal/C1/C4,

509 proliferative/C5/C3, immunoreactive/C2/C1, and differentiated/C4/C2)

22

K-Means (vertical axis label)

NMF (horizontal axis label)

TCGA   Yoshihara   TCGA   Yoshihara   TCGA   Yoshihara

Mayo   Tothill   Mayo   Tothill   Mayo   Tothill

Moderated
t-score Pearson
Correlation
1.0
0.5
0.0
-0.5
-1.0

|  | TCGA | Mayo | Yoshihara | Tothill | Bonome |
|---|---|---|---|---|---|
| k = 2 | | | | | |
| k = 3 | | | | | |
| k = 4 | | | | | |

$k = 2$     $k = 3$     $k = 4$

Moderated t-score Pearson Correlation

1.0
0.5
0.0
-0.5
-1.0

# TCGA

# Mayo

# Yoshihara

## TCGA

1 : 287 | 0.10

2 : 212 | 0.008

1 : 137 | 0.10

2 : 139 | 0.01

3 : 223 | 0.02

1 : 130 | 0.09

2 : 99 | 0.01

3 : 100 | 0.03

4 : 170 | 0.01

## Yoshihara

1 : 97 | 0.13

2 : 159 | 0.03

1 : 92 | 0.12

2 : 44 | 0.003

3 : 120 | 0.05

1 : 88 | 0.11

2 : 30 | 0.03

3 : 74 | 0.04

4 : 64 | 0.04

## Mayo

1 : 199 | 0.14

2 : 180 | 0.01

1 : 139 | 0.10

2 : 108 | 0.008

3 : 132 | 0.04

## Tothill

1 : 129 | 0.16

2 : 113 | 0.02

1 : 109 | 0.07

2 : 81 | 0.03

3 : 93 | 0.05

4 : 96 | 0.01

1 : 105 | 0.14

2 : 46 | 0.04

3 : 91 | 0.04

1 : 73 | 0.14

2 : 43 | 0.03

3 : 67 | 0.02

4 : 59 | 0.01