

1 **Pan-cancer identification of clinically relevant genomic** 2 **subtypes using outcome-weighted integrative clustering**

3 Arshi Arora¹, Adam B. Olshen^{2,3}, Venkatraman E. Seshan¹, and Ronglai Shen¹

4 ¹**Department of Epidemiology and Biostatistics, Memorial Sloan Kettering Cancer Center, New**
5 **York, NY**

6 ²**Department of Epidemiology and Biostatistics, University of California at San Francisco, CA**

7 ³**Helen Diller Family Comprehensive Cancer Center, University of California at San Francisco, CA**

8

9 **ABSTRACT**

10 Molecular phenotypes of cancer are complex and influenced by a multitude of factors. Conventional
11 unsupervised clustering of heterogeneous cancer patient populations is inevitably driven by the dominant
12 variation from major factors such as cell-of-origin or histology. Drawing from ideas in supervised text
13 classification, we developed survClust, an outcome-weighted clustering algorithm for integrative patient
14 stratification. We show survClust outperforms unsupervised clustering in identifying cancer patient
15 subpopulations characterized by specific genomic phenotypes with more aggressive clinical behavior.
16 The algorithm and tools we developed have direct utility toward clinically relevant patient stratification
17 based on tumor genomics to inform clinical decision-making.

18

19 **KEYWORDS**

20 Integrative Genomics, Supervised Clustering, Cancer Genomics, Statistical Methods, Data Integration

21

22 **INTRODUCTION**

23 Cancer is a complex disease with heterogeneous clinical outcomes. Understanding how patients respond
24 to treatment and what drives disease progression and metastasis is critical for managing and curing the
25 disease. Linking comprehensive molecular profiling data with patient outcome carries great promise in
26 addressing such important clinical questions. This requires innovative statistical and computational

27 methods designed for integrative analysis of multidimensional data sets to model intra-tumor and inter-
28 patient heterogeneity at genomic, epigenetic, and transcriptomic levels. Each of these molecular
29 dimensions is correlated yet characterize the disease in their own unique way. In order to arrive at a
30 comprehensive molecular portrait of the tumor, multiple groups have proposed statistical and
31 computational algorithms to synthesize various channels of information including methods developed by
32 us (iCluster^{1,2}) and others (PARADIGM³, CoCA⁴, SNF⁵, CIMLR⁶) to stratify disease populations. However,
33 the majority of the work has focused on unsupervised clustering, utilizing the molecular data alone.

34
35 Unsupervised learning does not necessarily lead to unique answers as the data are often
36 complex and multi-faceted. Consider the problem of clustering a collection of documents in text mining
37 where multiple structures can be present including authorship, topic, and style. The outcome of the
38 clustering is likely driven by a mixture of these underlying structures. As a result, there is often no single
39 “right” answer in unsupervised clustering problems. In most complex data applications, many local optima
40 exist that poses special challenges in optimization. Xing et al.⁷ proposed a weighted distance metric
41 allowing users to specify what they consider “meaningful” in defining similarity toward a more efficient and
42 local-optima free clustering performance.

43
44 Drawing analogy with the text learning problem described above, the molecular profile of a tumor
45 is influenced by a multitude of factors including tissue-of-origin⁸, histology (e.g., squamous vs.
46 adenocarcinoma), tumor microenvironment (e.g., immune cell infiltration⁹), dedifferentiation states¹⁰, and
47 specific pathway activation¹¹. Conventional unsupervised clustering applied to the most variable features
48 is inevitably driven by the dominant variation from major factors, for example, cell-of-origin⁸ or ancestry¹²
49 (germline variation) in the study cohort. When patient outcome related stratification is of interest, a more
50 directed clustering approach is needed.

51
52 We present *survClust*, an outcome-weighted integrative clustering algorithm for survival
53 stratification based on multi-dimensional omics-profiling data. The algorithm learns a weighted distance
54 matrix that down-weights molecular features with no relevance to the outcome of interest. This method

55 can be used on individual platforms alone, or by integrating various molecular platforms, to mine
56 biological information leading to distinct survival subgroups. We analyzed over 6,000 tumors across 18
57 cancer types. Each disease type was classified by *survClust*, based on six molecular assays – somatic
58 point mutations, DNA copy number, DNA methylation, mRNA expression, miRNA expression, protein
59 expression, and the integration of the six assays. The results have revealed novel survival subtypes not
60 previously identified by unsupervised clustering.

61 RESULTS

62 The *survClust* model: motivation and method overview

63 The molecular profile of a tumor often harbors information on a multitude of factors including cell lineage,
64 tumor microenvironment, cell differentiation and other clinical and histopathological features. Some of
65 these factors are associated with treatment response and/or survival outcome, while others are not. If a
66 particular patient outcome (e.g., patient survival) is of interest, a more supervised approach is needed.
67 We demonstrate this using a simulated data example (**Fig. 1a, Supplementary Fig 1**). In this scenario,
68 we simulated three risk subgroups in a cohort of 300 hypothetical patient samples with distinct survival
69 hazard rates in each subgroup (a median survival of 4, 3, and 2 years respectively). A set of 15 features
70 was then simulated from a mixture Gaussian distribution with different means in the three risk subgroups.
71 Another set of 15 features was simulated in the same way but permuted to disrupt the feature-risk
72 group association. A third group of 270 features were simulated from Gaussian noise. Figure 1b shows
73 that an unsupervised clustering using the K-means algorithm failed to identify the survival subtypes in the
74 context of complex feature variations. To identify outcome-associated clustering solution, *survClust*
75 utilizes a weighted distance metric:

$$d(\mathbf{a}, \mathbf{b}) = \sqrt{(\mathbf{a} - \mathbf{b})^T \mathbf{W}(\mathbf{a} - \mathbf{b})},$$

76 where (\mathbf{a}, \mathbf{b}) denote a pair of sample vectors measured for p features, and \mathbf{W} is a diagonal weight matrix
77 over p features with $\mathbf{W} = \text{diag}\{w_1, \dots, w_p\}$. The weights w_p 's are obtained by fitting a univariate cox
78 proportional hazards model for each feature in the training data with repeated training-test sample splits
79 for cross-validation (see more details in the Methods Section). Figure 1c shows that *survClust* was able to
80 identify the true risk groups with 97.15% accuracy [95% CI = 94% - 100%], whereas the accuracy from an

81 unsupervised clustering was 67.50% without reducing the effect of the survival unrelated and noise
82 features.

83 Our algorithm allows the integration of multiple data modalities. Given m data types measured
84 over respective feature space (**Fig. 1d**), the algorithm learns a weighted distance matrix from each
85 molecular data type incorporating a vector of Cox regression hazard ratio as weights. Each feature is
86 weighed and a pairwise distance matrix is calculated (we refer to this step as **getDist**). This step reduces
87 the computation space considerably by transforming the problem from sample by feature to sample by
88 sample. Note that, different sample sizes across data types are allowed, i.e., a sample can be measured
89 for some but not all platforms. Next, the weighted pairwise distance matrices are integrated by summing
90 over weighted m data types (**combineDist**), which retains all samples with at least one data type
91 available, with complete pairwise information. **survClust** then projects the integrated and weighted
92 distance matrix into a lower dimensional space via multidimensional scaling (MDS) and then clusters
93 sample points into subgroups via the K-means algorithm. More details can be found in the Methods
94 Section.

95

96 ***survClust* is more powerful than unsupervised clustering in identifying clinically relevant** 97 **molecular subtypes**

98 We applied *survClust* to the TCGA data set including 6,209 tumor samples in 18 cancer types to identify
99 survival outcome-associated subtypes defined by somatic mutation, DNA copy number, DNA methylation,
100 mRNA expression, and protein expression, individually and integratively. A summary of the sample sizes
101 and feature space is included in Supplementary Table 1. Supplementary Table 2 compares the survival
102 association (log-rank statistic) for the *survClust* integrated subtypes versus those derived from
103 unsupervised clustering methods commonly used in TCGA studies including COCA and iCluster. The log-
104 rank statistic compares estimates of the hazard functions of each subgroup comparing to the expected
105 values under the null hypothesis (all subgroups have identical hazard functions). Larger log-rank statistic
106 suggests stronger evidence of survival association. By differentially weighting the molecular features by
107 the corresponding survival association in constructing the distance matrix, we show that *survClust* is more
108 powerful in identifying subtypes that are directly relevant to stratify the outcome of interest, leading to

109 substantially more distinct survival subgroups than those existing molecular subclasses obtained by
110 unsupervised clustering. To further demonstrate, we highlight the *survClust* analysis of low-grade glioma
111 and kidney papillary renal cell carcinoma below.

112
113 ***survClust* identifies a poor prognostic *IDH*-mutant low-grade glioma subgroup.** Low Grade Gliomas
114 (LGG) have a unique molecular footprint, characterized by *IDH1/2* mutation status and co-deletion in
115 chromosome 1p and 19q regions of the genome¹³. As shown previously, mutations in *IDH1* and *IDH2*
116 genes are present in a majority of the low-grade gliomas and define a subtype associated with favorable
117 prognosis¹⁴. *IDH*-mutant tumors with chromosome 1p and 19q codeletion (*IDH*mut-codel) exhibit the most
118 prolonged survival times followed by *IDH*-mutant tumors without the codeletion (*IDH*mut-non-codel), with
119 *IDH*-wt tumors demonstrating more aggressive clinical behavior. We performed *survClust* on 6 available
120 molecular platforms (somatic mutation, DNA copy number, DNA methylation, mRNA expression, and
121 protein expression) in 512 LGG samples as profiled by the TCGA. The optimal number of clusters *k* was
122 chosen by assessing *survClust* fits over log-rank test statistics and standardized pooled within-cluster
123 sum-of-squares in cross-validation (see Methods Section). Cross-validation was performed to ensure
124 unbiased estimation of survival association and to avoid over-fitting.

125
126 The integrated *survClust* solution for LGG was optimized at *k*=5, with the *IDH*-mutant-codel (*c*3)
127 and *IDH*-mutant-non-codel (*c*1) subtypes associated with good prognosis as expected (**Fig 2a**). By
128 contrast, the *IDH*-wt subclass (*c*5) showed association with poor survival, enriched for mutations in *EGFR*
129 and *PTEN* gene and concurrent chromosome 7 gain and 10 loss, resembling glioblastomas. Interestingly,
130 *survClust* identified a small *IDH*-mutant subtype characterized by *CDKN2A* deletion (*c*4), which showed
131 markedly worse survival among the *IDH*-mutant tumors, similar to the *IDH*-wt group (*c*5) that tends to
132 behave far more aggressively with prognosis similar to glioblastomas. In addition, a copy number quiet
133 subgroup (*c*2) was identified, showing high expression of mir-1307 and mir-29c (**Supplementary Fig 3**).
134 These results highlight the strength of *survClust* in identifying clinically relevant molecular stratifications
135 and the potential to refine the existing paradigm in glioma subtyping to inform clinical decision-making.

136

137 ***survClust* identifies prognostic subtypes of kidney papillary renal-cell carcinoma (KIRP).** Three
138 survival distinct subtypes were identified using *survClust* integrating DNA copy number, mRNA
139 expression, DNA methylation, miRNA and protein expression assay profiled in 289 tumor samples. The
140 c3 subtype was associated with poor survival (median survival time = 1.63 yrs) (**Fig 2b**), associated with
141 younger age (median age 57 yrs) and more female gender (55%). The defining genomic characteristics
142 include *CDKN2A* loss, arm-level gains in multiple chromosomes including 7, 12, 15 and 17 as described
143 previously¹⁵.

144

145 ***survClust* identifies clinically relevant mutational subgroups across cancer types**

146 *survClust* is a flexible framework and can be applied to individual data types for patient stratification. For
147 example, somatic mutation based stratification is often of interest in a clinical sequencing setting. To
148 illustrate that, we applied *survClust* to mutation data alone using a hazard ratio weighted binary distance-
149 based clustering. A *circomap* plot was created to facilitate annotation and visualization of the results
150 across cancer types (**Fig 3a**). *survClust* identified high TMB subgroups in nearly all cancer types included
151 in this analysis. Correlating mutational signatures¹⁶ with these subtypes in the *circomap* plot further
152 revealed etiology underlying these hypermutated tumors. The smoking signature tracks lung cancer
153 (LUSC and LUAD) and the subset of head and neck cancer (HNSC) with elevated TMB. The DNA
154 mismatch repair (MMR) signature tracks high TMB subgroups in stomach (STAD), endometrial cancer
155 (UCEC), and colon cancer (COAD). The APOBEC signature is prevalent in bladder (BLCA) and cervical
156 cancers (CESC). Finally, the aristolochic acid signature (signature 22) is enriched in a liver cancer
157 subgroup identified by *survClust* (**Supplementary Fig 4e**), which is consistent with aristolochic acid and
158 their derivatives being implicated in liver cancers in Asian populations¹⁷.

159

160 In endometrial cancer, *survClust* confirmed a previously known ultra-high mutated subtype
161 associated with the POLE mutation signature (c2) and a hypermutated microsatellite instability (MSI) (c4)
162 subtype¹⁸ (**Fig 3b**). The *panelmap* in Figure 3b (middle panel) shows that c4 correlated well with clinical
163 MSI status ($P < 0.001$) and predominantly carried mutants in *ARID1A*, *PIK3CA* and *PTEN* genes. The c1
164 subtype, consisting of primarily high-grade serious tumors, was associated with worse outcome with a 5-

165 year survival of 58% compared to 95%, 84%, and 83% for c2 (POLE), c3, and c4 (MMR) respectively,
166 and characterized by higher frequency of mutations in *TP53*, *PPP2R1A* genes, low TMB and older
167 patients with serous endometrial tumors (60%). The c3 subtype was characterized by higher frequency of
168 *CTNNB1* mutants. Immune cell decomposition data derived using the CIBERSORT¹⁹ algorithm was also
169 correlated with the subgroups. Interestingly, high expression of CD8 T-cell immune marker was observed
170 in the POLE (c2) and MSI (c4) subtype ($P < 0.001$) (**Fig 3b**).

171
172 *survClust* stratified the bladder cancer cohort into 3 TMB subgroups – with high (c1), intermediate
173 (c3) and low (c2) mutation burden. The c1 subtype was associated with good outcome, high TMB, high
174 neoantigen load, high APOBEC load, and high expression of the CD8 T-Cell immune marker ($P=0.002$)
175 (**Fig 3c**). The c3 subtype showed intermediate TMB and APOBEC load with a median survival time of
176 3.48 yrs. Patients with a low TMB and low APOBEC load performed the worst in terms of survival with a
177 median survival time of 1.91 yrs.

178
179 A similar pattern emerged when *survClust* was run on colorectal cancer mutation data classifying
180 the disease population into three clusters – two low TMB groups and a MMR-associated high TMB group
181 (c1) (**Supplementary Fig 4b**). c1 was also associated with CD8 T-cell infiltration ($P = 0.004$) and showed
182 concordance with MLH1 silencing status. A similar subdivision of low TMB group by *TP53* mutation status
183 was seen where c3 carried *TP53* mutant samples unlike c2. Correlation with histology revealed significant
184 enrichment of mucinous adenocarcinoma subtype in c1 and c2 (c1, $n=20$, 29%; c2, $n=24$, 20%)
185 compared to c3 ($n=9$, 5%). In addition to the hypermutated subtypes of endometrial, bladder and
186 colorectal cancers, we also observed high TMB subgroups with concurrently high expression of CD8 T-
187 cell markers in cervical cancer c1 subtype (**Supplementary Fig 4a and 5a**), head and neck cancer c4
188 subtype (**Supplementary Fig 4c, 5c**), lung adenocarcinoma c3 subtype (**Supplementary Fig 4f and 5f**),
189 lung squamous cell carcinoma c4 subtype (**Supplementary Fig 4g and 5g**) and stomach cancer c1
190 subtype (**Supplementary Fig 4h and 5h**). There are prior observations that high mutational burden is
191 associated with increased neo-antigen load and activated T-cell infiltration in lung cancer²⁰. Our analysis
192 revealed that such association may be more widely present in multiple cancer types.

193
194 ***survClust* identifies distinct copy number subtypes associated with clinical features across**
195 **cancer types**

196 To identify copy number alterations that define clinically relevant subtypes, segmented data of 18 cancer
197 types was processed via the CBS algorithm²¹ and analyzed with *survClust*. Subtypes characterized by
198 different degrees in the Fraction of Genome Altered (FGA) emerged in various cancer types (**Fig 4**).
199 Interestingly, low FGA was associated with better survival in several cancer types including colon, head
200 and neck, lung adenocarcinoma, soft tissue sarcoma and endometrial cancer (**Supplementary Fig 6 and**
201 **7**).

202
203 The *circomap* plot in Figure 4a also revealed association of subtypes with high-level amplification
204 of major cancer genes including *CCND1* amplification in head and neck cancer (c3), *CCNE1* (c5) and
205 *AKT2*(c6) amplification in ovarian cancer, and *MDM2* amplification (c4) in sarcoma (**Supplementary Fig**
206 **6**). Notably, amplification of 19q13.2 region in ovarian cancer c6 subtype harboring the *AKT2* gene is
207 associated with poor survival (**Supplementary Fig 7f, Supplementary Table 8**) which was consistent
208 with previous findings that *AKT2* amplification is associated with ovarian cancer aggressiveness²².
209 *CCND1* amplified subtype of head and neck cancer (c3) was also associated with poor survival
210 (**Supplementary Fig 7b**). Amplification in the *MYC* gene is broadly present in multiple cancer types (**Fig**
211 **3a circomap**). Among cancer gene deletions, *CDKN2A* loss was observed to define multiple subgroups
212 associated with poor survival including papillary kidney cancer (c1), low-grade glioma (c4), lung
213 adenocarcinoma (c4), and soft tissue sarcoma (c1) (**Supplementary Fig 6 and 7**).

214
215 Colorectal cancer was classified into three varying FGA subtypes with prognostic implications. c1
216 had low FGA and, c2 and c3 carried heavy genome alterations (**Supplementary Fig 6a**). Even though c1
217 and c2 had dissimilar FGA, they performed similar in terms of survival as compared to c3, which had poor
218 outcome with median survival time of 4.5 yrs. (**Supplementary Fig 7a**). Gain in the *MYC* gene was seen
219 throughout the cancer type and c2 was uniquely characterized by loss of the chromosome 20 p-arm,
220 which harbors the hsa-mir-103-2 previously reported to be downregulated in colorectal tumors^{23,24}.

221
222 *survClust* is designed to capture the contribution of survival associated molecular features and
223 reduce the influence from those that are not related to the outcome of interest. Figure 4b provides another
224 example that this approach is better in identifying prognostically relevant subtypes compared to the
225 unsupervised clustering approach applied in the original study²⁵. *survClust* identified 6 unique CN groups
226 in liver cancer, with significant survival differences among subgroups. The c5 subtype was characterized
227 by high FGA and associated with poor outcome with a median survival time of 0.77 yrs. This cluster was
228 distinguished by a loss of chromosome 15. The c2 subtype was associated with the lowest FGA and a
229 median survival time of 6.81 yrs. The c4 subtype was enriched for *CDKN2A* deletion with a median
230 survival time of 2.15 years. By contrast, unsupervised clustering generated subgroups with distinct
231 molecular differences but did not show any separation in terms of survival.

232
233 **Integration of multiple data types enhances the identification of survival distinct subgroups**
234 Figure 5 shows that the integrated *survClust* solution outperformed individual platforms based on the
235 cross validated log rank statistics for multiple cancer types including cervical cancer, head and neck
236 cancer, papillary kidney cancer, lower grade glioma, liver and endometrial cancers. In general, the
237 integrated solutions always emerge at or near the top in performance as compared to the individual
238 platform specific solutions.

239
240 Next, we used the adjusted Rand index (RI) to evaluate the concordance between different
241 solutions. RI is calculated as the proportion of sample pairs that are assigned together in the same cluster
242 in one solution versus another, adjusted for what is expected by random chance. It provides an indirect
243 measure of how much a particular data type contributes to the integrated solution. A non-zero adjusted RI
244 across solutions would suggest shared biology across assay types in some tumors. For example, the
245 mutation subtypes of endometrial cancer (**Fig 5h**) have the highest adjusted RI (0.56) as compared to
246 the integrated solution, which is consistent with the fact that *POLE* and *MSI* are the two major prognostic
247 subtypes that are predominantly defined through mutation burden (**Fig 3b**). Nevertheless, the integrated
248 solution also clearly shows that there is additional information in DNA methylation, DNA copy number,

249 and mRNA expression being effectively incorporated by *survClust* that improved the survival stratification.
250 In bladder cancer, the integrated solution is most concordant with the mRNA cluster solution (adjusted RI
251 = 0.39), indicating influence by mRNA features towards integration (**Fig 5a**). Classification by mutation
252 data type seemed to have little or no overlap between other assays (adjusted RI close to 0), although the
253 integrated solution retained some information. (adjusted RI=0.03).

254
255 The integrated solution classified cervical cancer samples better than rest of the platforms and
256 pointed towards a 5-cluster solution (**Fig 5b**). Interestingly, a high degree of heterogeneity among
257 different platforms was observed as represented by a small adjusted RI across the board. The head and
258 neck cancer integrated solution showed great improvement over individual platforms for $k > 2$ solutions.
259 The $k=4$ integrated solution clearly resulted from effective integration of multiple data types including DNA
260 methylation, DNA copy number, and mRNA expression with an adjusted RI of 0.33, 0.26 and 0.25
261 respectively (**Fig 5c**). In this case, RPPA provided very little information toward the integrated solution.

262
263 The integrated *survClust* analysis stratified papillary kidney cancer type into 3 groups, with CN
264 sharing maximum information with the integrated solution (adjusted RI = 0.32), followed by mRNA (0.31),
265 miRNA (0.24), RPPA (0.23), and Methylation (0.19). Lower grade glioma displayed a wide range of
266 variability among platform type in terms of the logrank statistic (logrank statistic, x-axis from 0-250). The
267 $k=5$ integrated solution performed the best among the 6 platforms with larger contributions from mRNA
268 (RI = 0.63), copy number (RI = 0.62) and mutation (RI = 0.57) (**Fig 5e**). The integrated solution of liver
269 cancer did not show much improvement over individual assay types. Note that we did not use protein
270 data while integrating as more than half of the samples were not assayed with the protein platform
271 (RPPA, $n=182$; integrated $n=371$). miRNA, mRNA and copy number showed high median logrank
272 statistics over rounds of cross-validation demonstrating their role as potential prognostic classifiers.

273
274 **DISCUSSION**
275 We proposed a supervised clustering algorithm, *survClust*, that directly incorporates time to event (e.g.,
276 death, disease progression) information with molecular features to stratify patients into clinically relevant

277 subtypes. We further developed two visualization tools, *circomap* and *panelmap* for displaying and
278 annotating the resulting stratification. As more clinically annotated genomic data is becomes available as
279 a result of clinical sequencing programs^{26,27}, our method will provide a useful tool to facilitate patient
280 stratification for clinical decision making. In this study, we analyzed 18 cancer types in ~ 6200 tumors.
281 Each disease type was classified by *survClust* based on six molecular assays – somatic point mutation,
282 DNA copy number, DNA methylation, mRNA expression, miRNA expression, protein expression and
283 integration of the aforementioned six assays.

284
285 The supervised clustering approach provides a more direct way to identify survival associated
286 molecular subclasses, often leading to substantially more distinct survival subgroups than those existing
287 molecular subclasses obtained by unsupervised clustering. For example. The integrated *survClust*
288 stratification of the hepatocellular carcinomas (LIHC) was associated with a survival log-rank statistic of
289 45.19 (P<0.001) versus 1.69 (P=0.42) under the unsupervised clustering solution (**Supplementary Table**
290 **2, Supplementary Fig 8**), suggesting that *survClust* is a more powerful approach for identifying outcome-
291 associated subtypes. Supplementary Tables 2-7 show comparisons of the log-rank statistics in survival
292 differences across the various integrated and individual platform *survClust* solutions with those from
293 existing molecular clustering solutions reported in the TCGA publications (wherever available). Note that
294 *survClust* solutions have all been cross-validated to avoid overfitting.

295
296 The outcome-weighted learning framework we propose in this study can be extended to model
297 binary outcome types such as treatment response or toxicity (which is an important outcome category in
298 immunotherapy settings). In addition, the integration framework can facilitate the inclusion of other data
299 modalities including histopathological data and radiological images.

300

301 **METHODS**

302 ***survClust* workflow**

303 Let X_m be the m^{th} ($m=1, \dots, M$) data type of dimension N_m (number of samples in m^{th} data type, can vary)
304 by p_m (number of features). Data types may consist of continuous (gene expression, copy number log-

305 ratio, DNA methylation, miRNA, protein expression) or binary (mutation status) data. Overall survival is
306 defined as time from diagnosis to death or last follow-up. The data needs to be pre-processed as
307 described in **Supplementary Information**.

308 For a pair of two samples a and b , **the weighted distance**⁷ is calculated as follows:

$$309 \quad d_w(\mathbf{a}, \mathbf{b}) = \sqrt{(\mathbf{a} - \mathbf{b})^T \mathbf{W} (\mathbf{a} - \mathbf{b})}, \quad (1)$$

310 where, \mathbf{a} and \mathbf{b} are feature vectors of length p for samples a and b respectively, \mathbf{W} is a $p \times p$ diagonal
311 weight matrix with $\mathbf{W} = \text{diag} \{w_1, \dots, w_p\}$. Samples are close to each other when the value of d_w is small
312 and dissimilar when d_w is large.

313
314 The weights w_j ($j = 1, \dots, p$) are obtained by fitting a univariate cox proportional hazards model for each
315 feature:

$$h(t|\mathbf{x}_p) = h_0 \times \exp(\mathbf{x}_j^T * \beta), \quad (2)$$

316 where t represents the survival time, \mathbf{x}_j is the j^{th} column of matrix X of length N , h_0 is the baseline hazard
317 function, β is the regression coefficient and $\exp(\beta)$ is the Hazard Ratio (HR).

318
319 We consider the absolute value of HR on the logarithmic scale as the weight w . An HR=1
320 corresponds to the null that the feature is not associated with survival. This is reflected in a $\log(1) = 0$
321 weighting in the distance matrix. Since \mathbf{W} is a diagonal matrix with diagonal element w_j ($j = 1, \dots, p$), we
322 can simply use euclidean distance for computing distances if we transform the data as follows:

$$\mathbf{X}' = \mathbf{X} * \mathbf{W}^{\frac{1}{2}}, \quad (3)$$

323
324 Euclidean distances are sensitive to scale of the observations. After incorporating weights, we
325 standardize the data by its grand total:

$$326 \quad \frac{\mathbf{X}'}{\sum_i \sum_j x_{ij}'},$$

327

328 where, $\sum_i \sum_j x_{ij}'$ is the grand total of weighted matrix X' , with i rows (N samples) and j columns (p
329 features). Then, one can compute the pairwise distance between samples a ($i = 1$) and b ($i = 2$) as:

330
$$d_w(\mathbf{a}', \mathbf{b}') = d_w(\mathbf{b}', \mathbf{a}') = \sqrt{\sum_{j=1}^p (a_j' - b_j')^2}.$$

331
332 Conversely, a weighted distance matrix \mathbf{D} is calculated for all pairwise samples across M data types. All
333 samples having full survival information are kept, and the union of all samples (N_{union}) across M data
334 types is utilized when analyzing a wide number of samples. Non-overlapping samples in data types are
335 added as *NA* to have a uniform set of N_{union} samples.

336
337 The integrated weighted distance matrix is calculated by averaging over the weighted distance
338 matrices:

$$\mathbf{I}_w = \sum_{m=1}^M \gamma_m \mathbf{D}_m, \quad (4)$$

339 where $\gamma_m = \frac{1}{M} \forall m$. The integrated weighted matrix \mathbf{I}_w , averages the inter- and intra-sample similarity
340 profiles over the M data types. \mathbf{I}_w is then processed by *survClust* via classical multidimensional scaling
341 (MDS)²⁸ and clustered using k-means²⁹. Classical MDS assumes Euclidean distances; however, in cases
342 of non-Euclidean distances, Mardia et al³⁰ provided a method to obtain the resulting positive semidefinite
343 scalar product matrix. Note that \mathbf{I}_w follows the Euclidean norm and hence can be represented in
344 $n - 1$ dimensions. The strong assumption of the Euclidean norm is also important for k-means, as it is
345 essentially trying to assign samples to the closest centroid or calculating the sum of squared deviations
346 from centroids.

347
348 ***Weighted distance metric for mutation data***

349 Somatic mutation data is represented as a binary data matrix where each entry is coded as 1 if the j^{th}
350 gene is mutated in the i^{th} sample, and 0 otherwise. A challenge with the mutation data matrix is the
351 sparsity. It is known that somatic mutation data exhibit a long-tailed distribution in which a relatively small
352 number of variants appear in tumors frequently while the vast majority of variants occur extremely

353 infrequently. We consider genes that are mutated in > 1% of the sample. After incorporating weights, this
354 data is no longer binary, but it still remains sparse. Due to such data sparsity, computing Euclidean
355 distance is not appropriate and may lead to inflated distance measures³¹. To combat this problem, we
356 propose a weighted binary distance metric for such a scenario as described below.

357 Let X'_{mut} be the weighted mutation data matrix (see Equation. 3) of dimension N (samples) by p
358 (genes). Then, the pairwise distance between sample vectors \mathbf{a} and \mathbf{b} is calculated as follows:

359

$$d_w(\mathbf{a}, \mathbf{b}) = d_w(\mathbf{b}, \mathbf{a}) = \frac{w_{01} + w_{10}}{w_{01} + w_{10} + w_{11}},$$

360 where

361 w_{01} = sum of weights of p features that are zero in sample vector \mathbf{a} but non-zero in sample vector \mathbf{b} ;

362 w_{10} = sum of weights of p features that are non-zero in sample vector \mathbf{a} but zero in sample vector \mathbf{b} ;

363 w_{11} = sum of weights of p features that are non-zero in sample vector \mathbf{a} and non-zero in sample vector \mathbf{b} .

364

365 Note that, $d_w(\mathbf{a}, \mathbf{b})$ is a proportion of sum of effect sizes in which only one is non-zero amongst those in
366 which at least one is non-zero.³²

367

368 **Cross-validation**

369 *survClust* classifies sample populations by incorporating outcome information. Resulting clusters are
370 overly optimistic and need to be cross validated to arrive at more generalizable solutions. The
371 *cv.survclust* function provides cross validation for the desired number of folds and outputs cross-validated
372 solution labels. In the results shown above, we performed 5-fold cross validation as follows: (1) Split the
373 data into 5 random partitions, label 4 of them as the training sets and the remaining one as the test set.
374 (2) The weighted distance matrix was calculated from the training data set alone (Eq.1). *survClust*
375 clustering was performed to arrive at outcome weighted labels in the training set. (3) test labels were
376 predicted according to training labels (4) Step 2 was repeated until predictions were made on all 5 test
377 data sets across all 5 folds. (6) clusters were tracked by *centroid relabeling* (**Supplementary Note 1.3**)
378 across folds, and we obtained outcome weighted class labels for our entire dataset. This concluded one

379 round of cross-validation. All results shown here are results from cross-validated labels across 50 rounds
380 of cross-validation. Cluster meaning was preserved across rounds of cross validation via a similar
381 approach to centroid relabeling. The final label for a sample was assigned to a class to which it was
382 predicted in the maximum number of rounds. This is achieved by another function called
383 *consensus.summary*.

384

385 **Choice of the number of clusters k**

386 The logrank test statistic and standardized pooled within-cluster sum of squares were calculated from
387 cross-validated labels to choose an appropriate k .

388

389 **Logrank test statistic**

390 For a particular k cluster solution we have k cross-validated labels. Each class is distinct in survival and
391 we can compare the difference between classes using the logrank test statistic as follows³³:

$$\chi^2 = \frac{\sum_k (O_k - E_k)}{\sqrt{V}},$$

392 where, O_k = observed number of events in the k^{th} group over time, E_k = expected number of events in
393 the k^{th} group over time and $V = \sum Var (O_k - E_k) = \sum V_k$.

394 **Standardized pooled within-cluster sum of squares**

395 Here we calculate the pooled within-cluster sum of squares and standardize it by the total sum of squares
396 similar to methodology used in the gap statistic³⁴ to select the appropriate number of clusters.

397 Suppose that the final labels have clustered the data into k clusters C_1, C_2, \dots, C_k , with C_r denoting the
398 indices of observations in cluster r , and $n_r = |C_r|$. Let

$$w_r = \sum_{\substack{i, j \in C_r \\ i > j}} I_{w_{ij}},$$

399 where w_r is the sum of all pairwise distances in cluster r , $\{ij\}$ represents a pair of samples belonging to a
400 cluster C_r and I_w is calculated from Eq 4. Then the standardized pooled within-cluster sum of squares is
401 calculated as:

402

$$W_s = \sum_{r=1}^k w_r / \sum_{\substack{i \\ i>j}} \sum_j I_{wij} .$$

403
404 Here W_s decreases monotonically as the number of clusters k increases. The optimal number of clusters
405 is where W_s is minimized and creates an ‘elbow’ or a point of inflection, where addition of more clusters
406 does not improve cluster separation. Another property of W_s is that it can be used to compare amongst
407 different datasets as it lies between 0 and 1 after standardization. This is useful in comparing *survClust*
408 runs between individual data types and when we integrate them.

409 410 **Simulation**

411 Continuing from the simulation study presented in **Fig 1**, we go into detail about cross-validation and how
412 to chose k for a *survClust* run. In **Fig 1**, the input matrix was subjected to 50 rounds of 3-fold cross-
413 validation (2/3 training and 1/3 test. The *survClust* fit for a cluster k based on training data from each fold
414 was used to predict cluster membership for the remaining 1/3 test data. Final sample labels were
415 aggregated over all folds and cluster meaning was preserved across folds via centroid relabeling. (**See**
416 **Supplementary Note 1.3**).

417
418 Logrank test statistic and standardized pooled within-cluster sum of squares was calculated for the
419 consolidated test labels over 3-folds for each round. **Supplementary Fig 1(c)** summarizes these metrics
420 for 50 rounds of cross validation for $k=2-7$. We see that logrank is maximized for $k=3$, and the
421 standardized pooled within-cluster sum of squares elbows at $k=3$, pointing to the optimal selection of k at
422 $k=3$. The final class labels are assigned by consolidating solutions across all folds in all rounds of cross
423 validations.

424
425 **Implementation and availability.** *survClust* is freely available as an R package at
426 (<https://github.com/arorarshi/survClust>).

427 For k-means clustering, we used the k-means implementation in the R base package. For
428 multidimensional scaling, we used the *cmdscale* function in base R. The weighted distance metric for

429 binary data was programmed in C++ with R extension using *Rcpp* package, which is computationally fast.
430 Hazard ratios were derived from the cox proportional hazard model came from the R *survival* package.
431 Kaplan Meier curves were plotted using *ggsurvplot* in package *survminer*. Beeswarm plots were made
432 using R package *beeswarm*. Mutation data along with relevant clinical annotations were plotted using
433 *panelmap* (<https://github.com/arorarshi/panelmap>). The *circlize* R package was used to make pan-
434 cancer plots and the code used to plot these is available in a function called *circomap*
435 <https://github.com/arorarshi/panelmap#example---circomap>)

436

437 Below is the workflow of proposed *survClust* method:

- 438 1. **getDist** – Compute a weighted distance matrix across given m data types. Standardization and
439 accounting for non-overlapping samples is also accomplished in this step.
- 440 2. **combineDist** – Integrate m data types by averaging over m weighted distance matrices.
- 441 3. **survClust** and **cv.survclust** – Estimate the *survClust* solution for a given cluster number k based
442 on the weighted and integrated distance matrix. Optimal k is estimated via cross-validation. Use
443 the chosen k and the cross-validated results to arrive at final class labels. Cross-validated results
444 are assessed over the following performance metrics – the logrank statistic, standardized pooled
445 within-cluster sum of squares and cluster solutions with class size less than 5 samples.

446

447 **DECLARATIONS**

448 **Availability of Data and materials**

449 The *survClust* algorithm's software implementation is available at <https://github.com/arorarshi/survClust>.
450 Simulated dataset and simulated survival dataset are also available on GitHub. All genomics and clinical
451 data was downloaded from <https://portal.gdc.cancer.gov/>. *panelmap* and *circomap* are available at
452 <https://github.com/arorarshi/panelmap>

453

454 **Funding**

455 CA008748.

456

457 **Author's contributions**

458 A.A. A.B.O., V.E.S. and R.S. designed the research. A.A. made software implementations and analyzed
459 the data. A.A. A.B.O., V.E.S. and R.S. wrote the paper.

460

461 **Acknowledgements**

462 The research was supported by the National Cancer Institute, Award CA008748.

463

464 **Authors' information**

465 *Department of Epidemiology and Biostatistics, Memorial Sloan Kettering Cancer Center, New York, NY*

466 Arshi Arora, Venkatraman E. Seshan, Ronglai Shen

467

468 *Department of Epidemiology and Biostatistics, University of California at San Francisco, CA*

469 *Hellen Diller Family Comprehensive Cancer Center, University of California at San Francisco, CA*

470 Adam B. Olshen

471

472 **Corresponding authors**

473 Correspondence to Ronglai Shen (shenr@mskcc.org)

474

475 **REFERENCES**

476

477 1 Shen, R. L., Olshen, A. B. & Ladanyi, M. Integrative clustering of multiple genomic data types
478 using a joint latent variable model with application to breast and lung cancer subtype analysis.
479 *Bioinformatics* **25**, 2906-2912, doi:10.1093/bioinformatics/btp543 (2009).

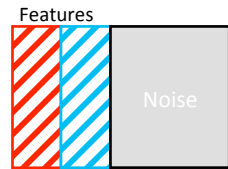
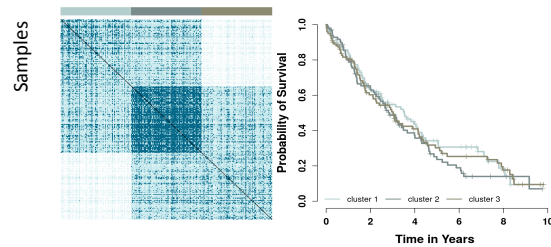
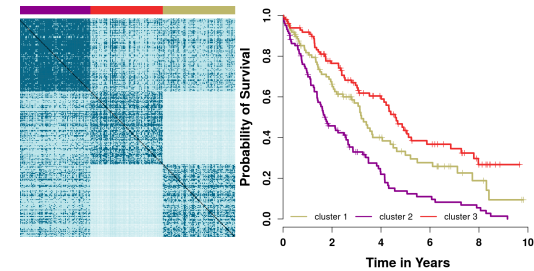
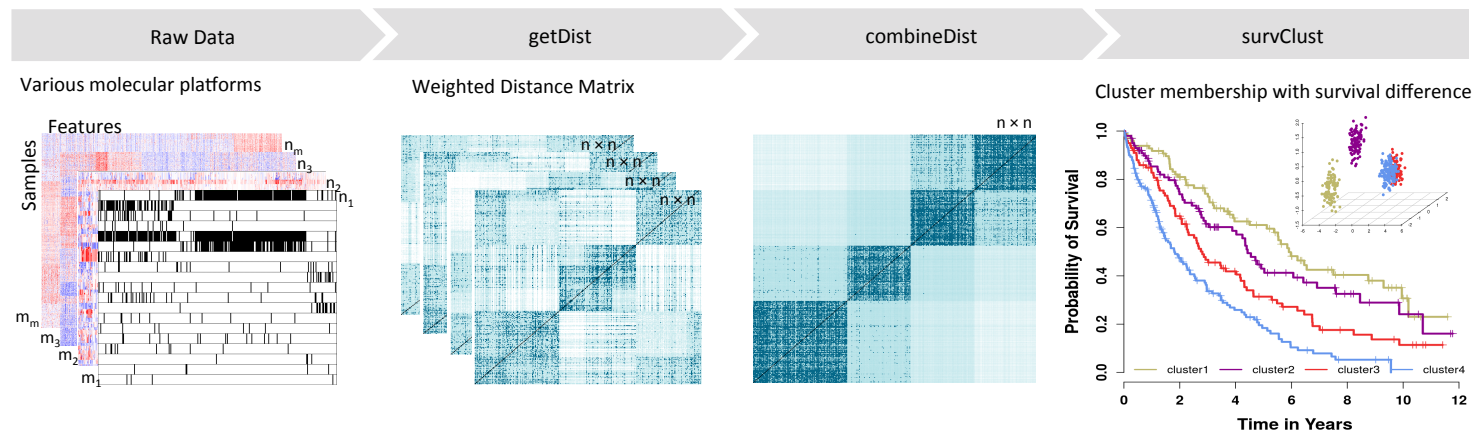
480 2 Mo, Q. *et al.* Pattern discovery and cancer gene identification in integrated cancer genomic data.
481 *Proc Natl Acad Sci U S A* **110**, 4245-4250, doi:10.1073/pnas.1208949110 (2013).

- 482 3 Vaske, C. J. *et al.* Inference of patient-specific pathway activities from multi-dimensional cancer
483 genomics data using PARADIGM. *Bioinformatics* **26**, i237-i245,
484 doi:10.1093/bioinformatics/btq182 (2010).
- 485 4 Hoadley, K. A. *et al.* Multiplatform Analysis of 12 Cancer Types Reveals Molecular Classification
486 within and across Tissues of Origin. *Cell* **158**, 929-944, doi:10.1016/j.cell.2014.06.049 (2014).
- 487 5 Wang, B. *et al.* Similarity network fusion for aggregating data types on a genomic scale. *Nat*
488 *Methods* **11**, 333-337, doi:10.1038/nmeth.2810 (2014).
- 489 6 Ramazzotti, D., Lal, A., Wang, B., Batzoglou, S. & Sidow, A. Multi-omic tumor data reveal
490 diversity of molecular mechanisms that correlate with survival. *Nature Communications* **9**, 4453,
491 doi:ARTN 4453 10.1038/s41467-018-06921-8 (2018).
- 492 7 Xing, E. P., Jordan, M. I., Russell, S. J. & Ng, A. Y. in *Advances in neural information processing*
493 *systems*. 521-528.
- 494 8 Hoadley, K. A. *et al.* Cell-of-Origin Patterns Dominate the Molecular Classification of 10,000
495 Tumors from 33 Types of Cancer. *Cell* **173**, 291-304 e296, doi:10.1016/j.cell.2018.03.022 (2018).
- 496 9 Thorsson, V. *et al.* The Immune Landscape of Cancer. *Immunity* **48**, 812-830 e814,
497 doi:10.1016/j.immuni.2018.03.023 (2018).
- 498 10 Malta, T. M. *et al.* Machine Learning Identifies Stemness Features Associated with Oncogenic
499 Dedifferentiation. *Cell* **173**, 338-354 e315, doi:10.1016/j.cell.2018.03.034 (2018).
- 500 11 Sanchez-Vega, F. *et al.* Oncogenic Signaling Pathways in The Cancer Genome Atlas. *Cell* **173**, 321-
501 337 e310, doi:10.1016/j.cell.2018.03.035 (2018).
- 502 12 Novembre, J. *et al.* Genes mirror geography within Europe. *Nature* **456**, 98-U95,
503 doi:10.1038/nature07331 (2008).
- 504 13 Network, C. G. A. R. Comprehensive, integrative genomic analysis of diffuse lower-grade
505 gliomas. *New England Journal of Medicine* **372**, 2481-2498 (2015).

- 506 14 Yan, H. *et al.* IDH1 and IDH2 mutations in gliomas. *N Engl J Med* **360**, 765-773,
507 doi:10.1056/NEJMoa0808710 (2009).
- 508 15 Network, C. G. A. R. Comprehensive molecular characterization of papillary renal-cell carcinoma.
509 *New England Journal of Medicine* **374**, 135-145 (2016).
- 510 16 Alexandrov, L. B. *et al.* Signatures of mutational processes in human cancer. **500**, 415 (2013).
- 511 17 Zhou, Q. *et al.* Worldwide research trends on aristolochic acids (1957–2017): Suggestions for
512 researchers. **14**, e0216135 (2019).
- 513 18 Getz, G. *et al.* Integrated genomic characterization of endometrial carcinoma. *Nature* **497**, 67-
514 73, doi:10.1038/nature12113 (2013).
- 515 19 Newman, A. M. *et al.* Robust enumeration of cell subsets from tissue expression profiles. *Nat*
516 *Methods* **12**, 453-457, doi:10.1038/nmeth.3337 (2015).
- 517 20 Chae, Y. K. *et al.* Mutations in DNA repair genes are associated with increased neo-antigen load
518 and activated T cell infiltration in lung adenocarcinoma. *Oncotarget* **9**, 7949-7960,
519 doi:10.18632/oncotarget.23742 (2018).
- 520 21 Olshen, A. B., Venkatraman, E. S., Lucito, R. & Wigler, M. Circular binary segmentation for the
521 analysis of array-based DNA copy number data. *Biostatistics* **5**, 557-572,
522 doi:10.1093/biostatistics/kxh008 (2004).
- 523 22 Bellacosa, A. *et al.* Molecular Alterations of the Akt2 Oncogene in Ovarian and Breast
524 Carcinomas. *International Journal of Cancer* **64**, 280-285, doi:DOI 10.1002/ijc.2910640412
525 (1995).
- 526 23 Sheffer, M. *et al.* Association of survival and disease progression with chromosomal instability: a
527 genomic exploration of colorectal cancer. *Proc Natl Acad Sci U S A* **106**, 7131-7136,
528 doi:10.1073/pnas.0902232106 (2009).

- 529 24 Cummins, J. M. *et al.* The colorectal microRNAome. *Proc Natl Acad Sci U S A* **103**, 3687-3692,
530 doi:10.1073/pnas.0511155103 (2006).
- 531 25 Ally, A. *et al.* Comprehensive and integrative genomic characterization of hepatocellular
532 carcinoma. **169**, 1327-1341. e1323 (2017).
- 533 26 Zehir, A. *et al.* Mutational landscape of metastatic cancer revealed from prospective clinical
534 sequencing of 10,000 patients. **23**, 703 (2017).
- 535 27 Micheel, C. M. *et al.* American association for cancer research project genomics evidence
536 neoplasia information exchange: from inception to first data release and beyond—lessons
537 learned and member institutions' perspectives. **2**, 1-14 (2018).
- 538 28 Torgerson, W. S. Theory and methods of scaling. (1958).
- 539 29 Hartigan, J. A. & Wong, M. A. Algorithm AS 136: A k-means clustering algorithm. *J Journal of the*
540 *Royal Statistical Society. Series C* **28**, 100-108 (1979).
- 541 30 Mardia, K. V. & Methods. Some properties of classical multi-dimensional scaling. *J Communications*
542 *in Statistics-Theory* **7**, 1233-1241 (1978).
- 543 31 Legendre, P. & Gallagher, E. D. Ecologically meaningful transformations for ordination of species
544 data. *Oecologia* **129**, 271-280, doi:10.1007/s004420100716 (2001).
- 545 32 Martin, N. & Maes, H. *Multivariate analysis*. (Academic press London, 1979).
- 546 33 Harrington, D. P. & Fleming, T. R. A Class of Rank Test Procedures for Censored Survival-Data.
547 *Biometrika* **69**, 553-566, doi:DOI 10.1093/biomet/69.3.553 (1982).
- 548 34 Tibshirani, R., Walther, G. & Hastie, T. Estimating the number of clusters in a data set via the gap
549 statistic. *Journal of the Royal Statistical Society Series B-Statistical Methodology* **63**, 411-423,
550 doi:Doi 10.1111/1467-9868.00293 (2001).

551

a Simulated data type**b Unsupervised clustering solution****c SurvClust solution****d survClust Workflow****Figure 1: Overview of *survClust*.**

- A simulated data example, consisting of features that define 3 patient subtypes without direct association with survival (shaded in red), features that define 3 patient subtypes with distinct survival outcome (shaded in blue), and random features generated from Gaussian noise (grey).
- Euclidean distance matrix demonstrating patient-level pairwise similarity, with darker blue shade representative of higher similarity. Color panels above the distance matrix show the three class solution obtained by unsupervised algorithm via k-means and the concordance between the simulated 3 survival subtypes (the truth). Kaplan Meier curves for the 3 unsupervised subtypes show no distinction in survival outcome.
- survClust* employs a patient outcome weighted distance matrix to identify the desired subtypes with distinct Kaplan Meier curves.
- survClust* allows integrative analysis of multiple data modalities to identify survival-associated molecular subtypes.

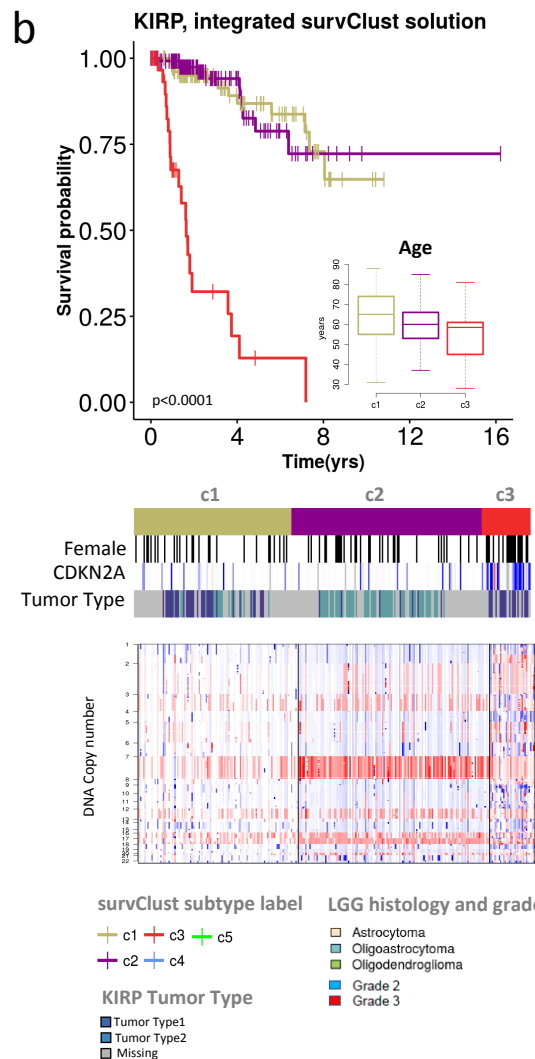
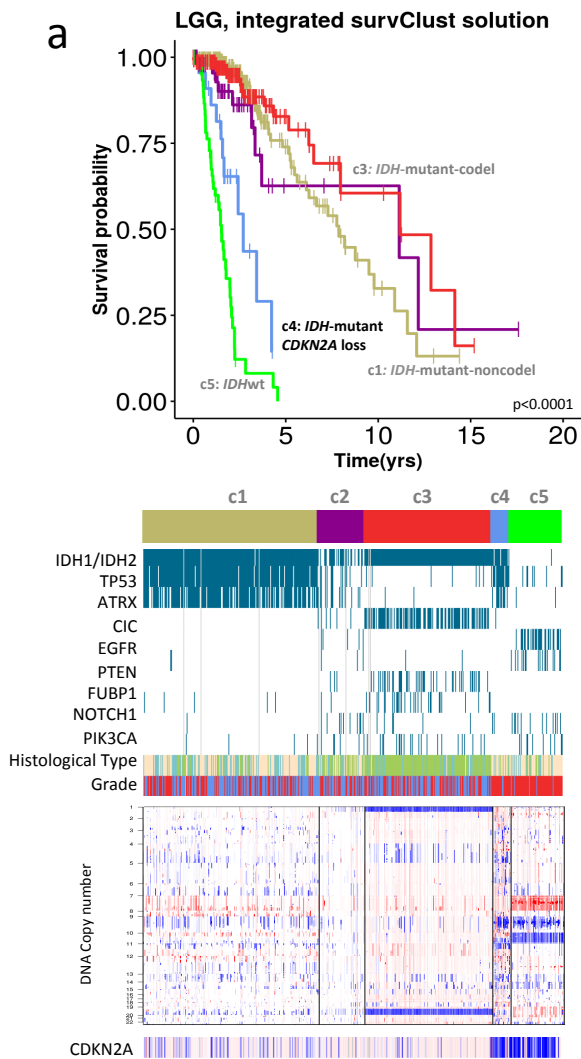


Figure 2: Outcome-weighted integrative clustering of low grade glioma and kidney papillary cell carcinoma using *survClust*.

(a) *survClust* identifies an *IDH*-mutant *CDKN2A*-loss subtype similar to *IDH*-wt tumors in terms of aggressive clinical behavior. Top: Kaplan-Meier curves of the integrated *survClust* subtypes of LGG. Middle: *Panelmap* summarizing major association of mutational and clinical features of the integrated LGG subtypes. Bottom: Copy number profile for each of the integrated subtypes.

(b) *survClust* identifies prognostic kidney papillary renal cell carcinoma (KIRP) subtypes.

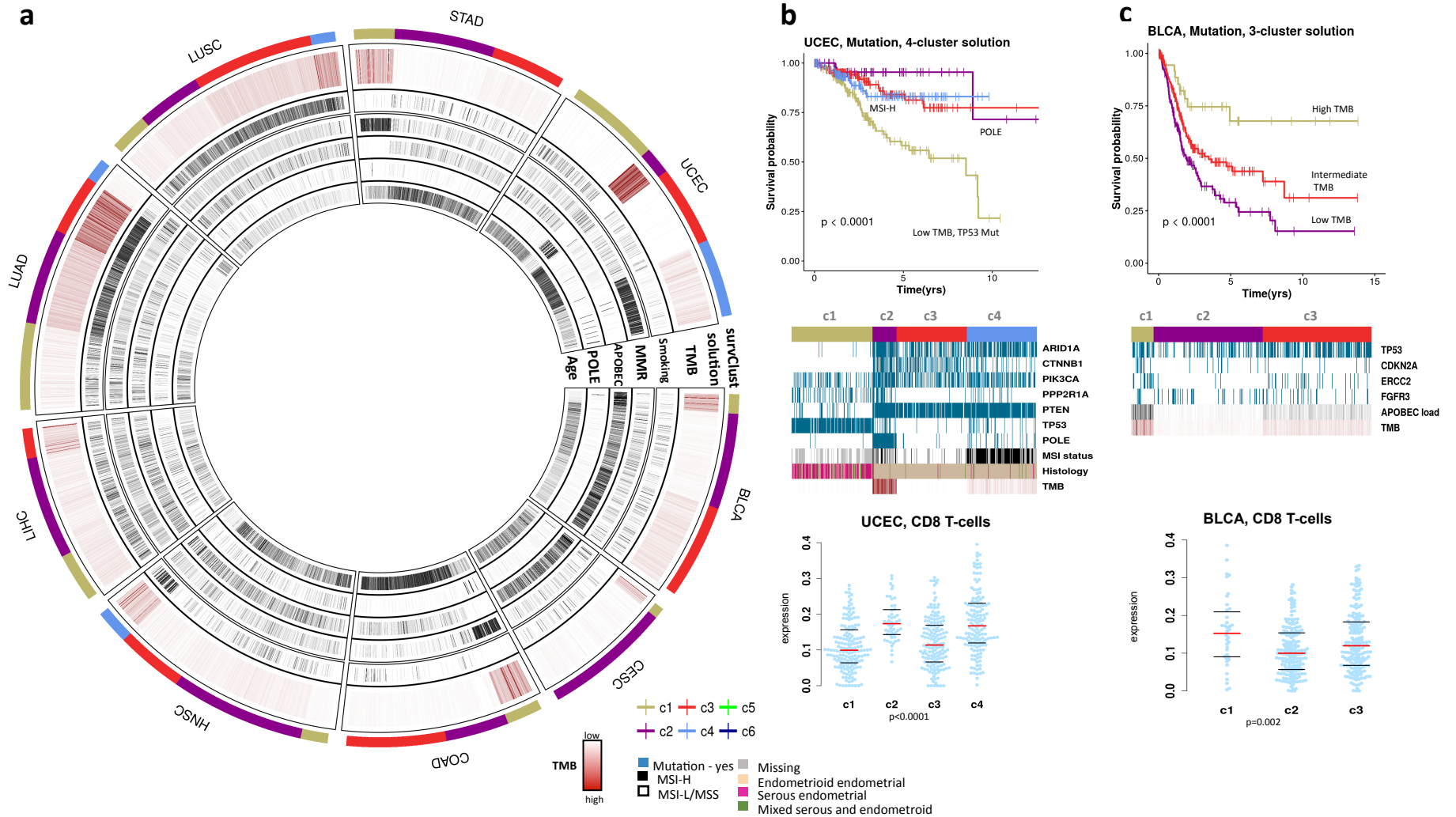


Figure 3: *survClust* identifies mutational subtypes associated with survival across cancer types.

- (a) *Circomap* showing total mutation burden (TMB) in brown color and mutational signatures (smoking, MMR, APOBEC, POLE and aging) in tumors across bladder (BLCA), cervical (CESC), colon (COAD), head and neck (HNSC), liver (LIHC), lung adenocarcinoma (LUAD), lung Squamous Cell (LUSC), stomach (STAD), and endometrial (UCEC) cancers. Outer circle indicates mutation-based *survClust* membership.
- (b) *survClust* mutation subtypes in endometrial cancer. From top to bottom: Kaplan-Meier curves for the 4 mutation subtypes, *panelmap* depicting significantly mutated genes, MSI status, Histology and TMB associated with the subtypes, and beeswarm plot showing CD8 T-cell marker expression (y-axis) across the 4 subtype (x-axis). Red line depicts the median, and top and bottom black bars represent the 25th and 75th percentile respectively.
- (c) *survClust* mutation subtypes in bladder cancer. From top to bottom: Kaplan-Meier curves for the 3 mutation subtypes, *panelmap* depicting significantly mutated genes, Papillary histology (yes – black, no-white), APOBEC load and TMB associated with the 3 subtypes, and beeswarm plot showing CD8 T-cell expression (y-axis) across the 3 subtypes (x-axis).

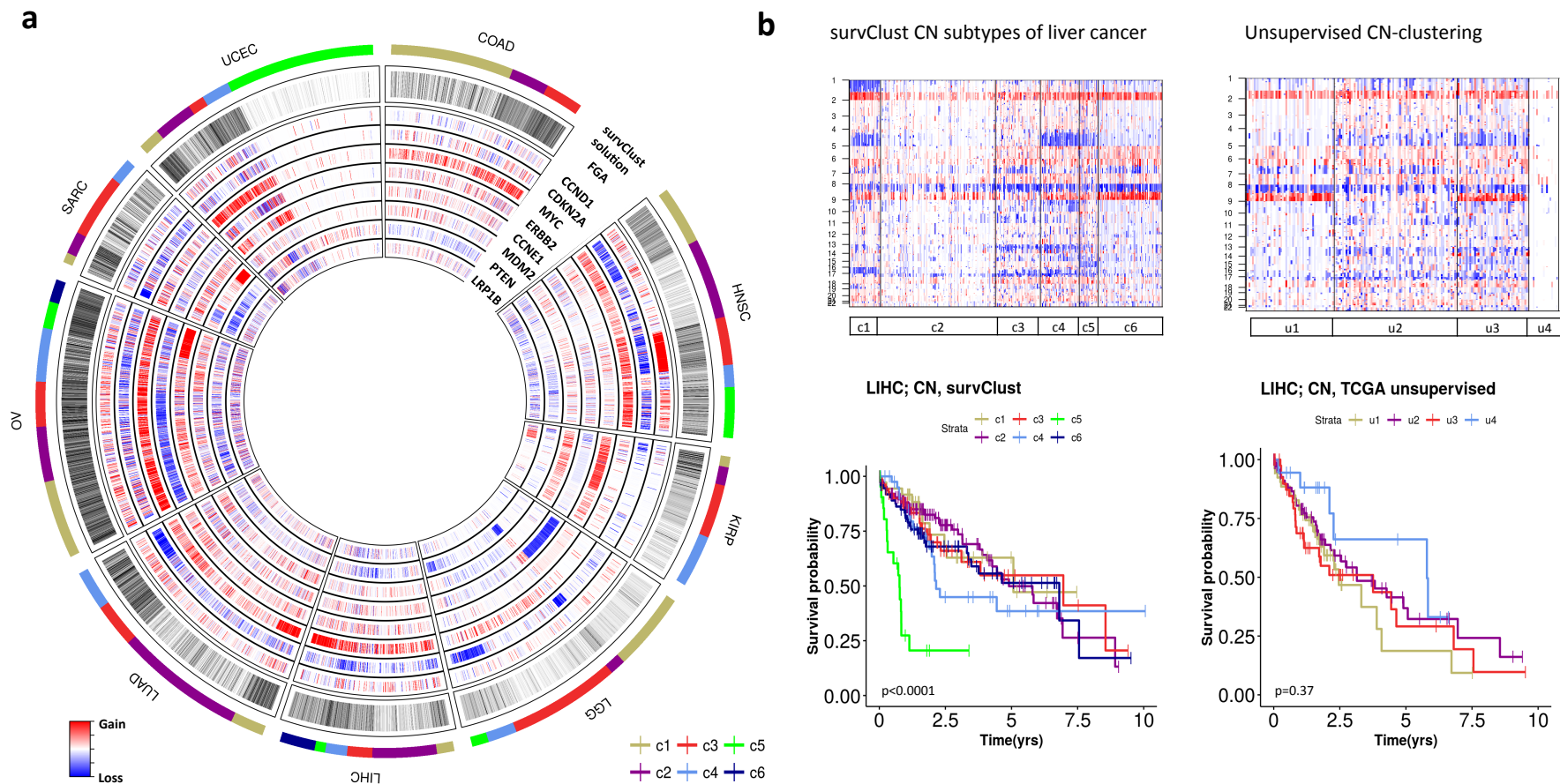


Figure 4: *survClust* identifies copy number patterns associated with patient survival outcome across various cancer types
 (a) *Circomap* showing fraction genome altered (FGA) and gene level copy number alterations in each tumor across colorectal (COAD), head and neck (HNSC,) kidney renal papillary cell carcinoma (KIRP), low grade glioma (LGG), liver (LIHC), lung adenocarcinoma (LUAD), ovarian (OV), soft-tissue sarcoma (SARC) and endometrial (UCEC) cancers. Outer circle indicates the *survClust* membership.
 (b) *survClust* is more powerful than unsupervised clustering in identifying survival-associated copy number subtypes in liver cancer.

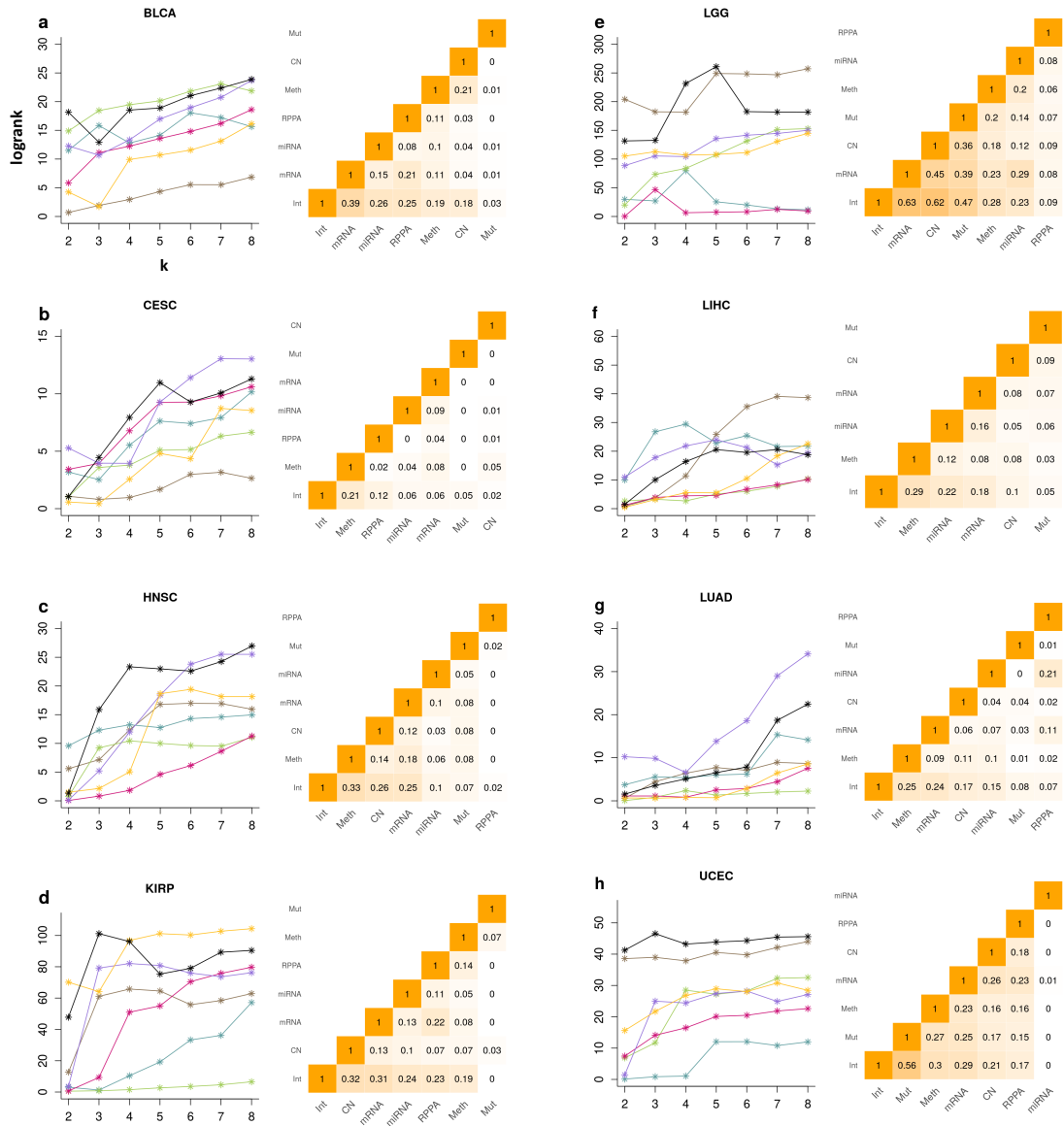


Figure 5: Integration of multiple data types enhances the identification of survival distinct subgroups
a-h: Each panel has two plots: the plot on the left summarizes median cross validated log rank statistic across k=2 to 8 (number of clusters). Each line is a data type (see legend), and the black line represents the *survClust* run on integrating all 6 platforms. Plot on the right summarizes the adjusted rand index between cross validated *survClust* solutions of individual data types and the integration of all. In this comparison, the *survClust* solution was chosen for an appropriate k which maximized logrank statistic and minimized the standardized pooled within sum of squares.

