

1 **Motif-based phosphoproteome clustering improves modeling** 2 **and interpretation**

3
4 **One-sentence Summary:** DDMC is a general and flexible strategy for phosphoproteomic analysis by
5 clustering phosphopeptides using both their phosphorylation abundance and sequence motifs.

6
7 **Authors:** Marc Creixell¹, Aaron S. Meyer^{1,2,3,4*}

8 **Affiliations:**

9 ¹Department of Bioengineering, University of California, Los Angeles, United States of America

10 ²Department of Bioinformatics, University of California, Los Angeles, United States of America

11 ³Jonsson Comprehensive Cancer Center, University of California, Los Angeles, United States of America

12 ⁴Eli and Edythe Broad Center of Regenerative Medicine and Stem Cell Research, University of
13 California, Los Angeles, United States of America

14 *Corresponding author. Email: ameyer@asmlab.org

15

16 **Abstract:** Cell signaling is orchestrated in part through a network of protein kinases and phosphatases.
17 Dysregulation of kinase signaling is widespread in diseases such as cancer and is readily targetable
18 through inhibitors of kinase enzymatic activity. Mass spectrometry-based analysis of kinase signaling can
19 provide a global view of kinase signaling regulation but making sense of these data is complicated by its
20 stochastic coverage of the proteome, measurement of substrates rather than kinase signaling itself, and the
21 scale of the data collected. Here, we implement a dual data and motif clustering strategy (DDMC) that
22 simultaneously clusters substrate peptides into similarly regulated groups based on their variation within
23 an experiment and their sequence profile. We show that this can help to identify putative upstream
24 kinases and supply more robust clustering. We apply this clustering to large-scale clinical proteomic
25 profiling of lung cancer and identify conserved proteomic signatures of tumorigenicity, genetic mutations,
26 and tumor immune infiltration. We propose that DDMC provides a general and flexible clustering
27 strategy for the analysis of phosphoproteomic data.

28 Introduction

29 Cell signaling networks formed by protein kinases dictate cell fate and behavior through protein
30 phosphorylation (1). As such, it is not surprising that kinase dysregulation orchestrates the onset and
31 development of a myriad of diseases, including cancer. Measuring cell signaling by mass spectrometry
32 (MS)-based global phosphoproteomics provides a promising opportunity to direct therapy development
33 (2), particularly given the accessibility of these signaling changes to drug targeting. Nevertheless, despite
34 the rapid accumulation of large-scale phosphoproteomic clinical data, it is still difficult to identify the
35 signaling events leading to observed proteomic alterations and phenotypic outcomes.

36 One approach to make sense of phosphoproteomic measurements has been to infer the activity of
37 upstream kinases. Previously published methods have combined each phosphopeptide with reported
38 kinase-substrate interactions to reconstruct signaling networks. For instance, kinase-substrate enrichment
39 analysis (KSEA) averages the signals of groups of kinase substrates to infer enriched pathways in
40 biological samples (3). Another method, Integrative Inferred Kinase Activity (INKA), infers kinase
41 activity by integrating the scores of two components that compute kinase's overall and activation loop
42 phosphorylation alongside another two components that quantify the phosphorylation abundance of
43 known substrates. Kinase-substrate relationships are either experimentally determined or predicted by
44 NetworKIN, an algorithm that uses sequence motif and protein-protein network information (4–6).
45 Finally, Scansite predicts kinase-substrate interactions using sequence motifs generated from oriented
46 peptide library scanning experiments (7). These methods, sometimes in combination, help to reconstruct
47 signaling pathway activities from phosphoproteomic measurements.

48 Kinase-substrate inference still provides a limited view of signaling network changes, however. Kinase
49 prediction methods are necessarily dependent on having well-characterized kinase-substrate interactions.
50 Unfortunately, the majority of the phosphoproteome remains largely uncharacterized (8). Just 20% of
51 kinases have been shown to phosphorylate 87% of currently annotated substrates and around 80% of
52 kinases have fewer than 20 substrates, with 30% yet to be assigned a single substrate (8). Hence, insights
53 generated by computational methods dependent on this unequal knowledge distribution are less likely to
54 identify understudied protein kinases. An additional major challenge being faced during the analysis of
55 large-scale signaling data is missingness. This is due to two major limitations of discovery-mode
56 multiplexed tandem mass tag (TMT) MS. The technique processes batches of samples with stochastic
57 signaling coverage in each experiment. This means that the portion of the phosphoproteome quantified in
58 the samples of different TMT experiments varies (9). Thus, in the resulting data set, phosphosites are
59 observed in certain groups of samples but not others. Computational tools usually require complete data
60 sets and so a frequent strategy to handle this challenge is either imputing missing values with a
61 representative statistic (e.g. average signal) or throwing out any peptides displaying missing values—at the
62 expense of losing critical information (10, 11). Kinase enrichment and prediction methods are further
63 compromised by this problem. Thus, there is a clear need to develop tailored and unbiased computational
64 methods capable of modeling the entirety of the phosphoproteomic data set despite missingness.

65 Clustering methods such as hierarchical or k-means clustering identify signaling nodes by grouping
66 phosphopeptides based on their co-variation. This clustering criterion results in groups of peptides that
67 display similar activation patterns across conditions, but that may be targeted by sets of different upstream
68 kinases. The residues surrounding phosphorylation sites have had to evolve throughout millions of years
69 to become exquisitely fine-tuned motifs that confer signaling specificity and fidelity (12, 13). Clustering
70 based on motif similarity might, therefore, improve model interpretation by facilitating the identification
71 of upstream kinases modulating particular clusters that display conserved sequence motifs. On the other
72 hand, clustering peptides based on sequence distance may result in groups of proteins that, while sharing
73 the same set of upstream kinases, are differently regulated due to context. Thus, combining

74 phosphorylation status and sequence similarity may enable a balanced characterization of the cell
75 signaling state.

76 Here, we present an algorithm, Dual Data and Motif Clustering (DDMC), that probabilistically and
77 simultaneously models both the peptide phosphorylation variation and peptide sequence motifs of peptide
78 clusters to reconstitute cell signaling networks and identify causal interactions (Fig. 1). To test the utility
79 of our method, we analyze the phosphoproteomes of 110 treatment-naïve lung adenocarcinoma (LUAD)
80 tumors and 101 paired normal adjacent tissues (NATs) from the National Cancer Institute (NCI)'s Clinical
81 Proteomic Tumor Analysis Consortium (CPTAC) LUAD study (11). We characterize the
82 phosphoproteome of patients by identifying those signaling signatures associated with tumorigenesis, the
83 presence of specific mutations, and tumor immune infiltration. In total, we demonstrate DDMC as a
84 general strategy for improving the analysis of phosphoproteomic surveys.

85 **Results**

86 **Constructing an expectation-maximization algorithm tailored for clustering phosphoproteomic** 87 **data**

88 MS-based global phosphoproteomic data provides unparalleled coverage when interrogating kinase
89 signaling networks and their therapeutic implications. However, these data also present challenging issues
90 as a consequence of their incomplete and stochastic coverage, high-content but low-sample throughput,
91 and variation in coverage across experiments. In addressing these issues, we recognized that MS
92 measurements provide two pieces of information: the exact site of phosphorylation on a peptide sequence
93 and some measure of abundance within the measured samples. Both of these pieces of information are
94 critical to the overall interpretation of the data.

95 Based on this observation, we built a mixture model that probabilistically clusters phosphosites based on
96 both their peptide sequence and abundance across samples (Figure S1). In each iteration, DDMC applies
97 an expectation-maximization algorithm to optimize clusters that capture the average features of member
98 sequences and their abundance variation (Figure 1A and S1). Both information sources—peptide
99 abundance and sequence—can be prioritized by a weight parameter. With a weight of 0, DDMC becomes
100 a Gaussian Mixture Model (GMM) that clusters peptides according to their phosphorylation signal. With
101 a very large weight, DDMC exclusively clusters peptides according to their peptide sequences. Clustering
102 both the sequence and abundance measurements ensures that the resulting clusters are a function of both
103 features, which we hypothesized would provide both more meaningful and robust clusters.

104 The resulting clustering provides coordinated outputs that can be used in a few different ways. The cluster
105 centers, by virtue of being a summary for the abundance changes of these peptides, can be regressed
106 against phenotypic responses (e.g., cell phenotypes or clinical outcomes) to establish associations
107 between particular clusters and response (Figure 1B). Regression using the clusters instead of each
108 peptide ensures that the model can be developed despite relatively few samples, with minimal loss of
109 information since each peptide within a cluster varies in a similar manner.

110 In parallel or independently, one can interrogate the resulting Position-Specific Scoring Matrices
111 (PSSMs) to describe the overall sequence features of that cluster. These outputs can be readily compared
112 to other information such as experimentally generated profiles of putative upstream kinases via Position
113 Specific Scanning Libraries (PSPL) (14–18). We extracted a collection of 62 kinase specificity profiles to
114 identify which cluster motifs most resemble the optimal motif of putative upstream kinases (Figure 1C)
115 (17–19). However, as kinase-substrate specificity is also dictated by features outside of the immediate
116 substrate region, we also note that our approach is more general than strictly assembling kinase-substrate
117 predictions as non-enzymatic specificity information may be present in the DDMC sequence motifs.
118 Overall, this overview demonstrates how DDMC can take complex, coordinated signaling measurements

119 and find patterns in the phosphorylation signals to reconstruct signaling networks and associate particular
120 clusters and phenotypes.

121 **Dual data-motif clustering strategy robustly imputes missing values**

122 A major limitation of multiplexed MS-based large-scale phosphoproteomic data is the presence of
123 missing values due to (i) the limited number of samples processed at a time per TMT experiment and (ii)
124 the stochastic signaling coverage in each experiment. Consequently, upon concatenation of the different
125 TMT experiments, many phosphosites are observed in groups of samples. To evaluate the robustness of
126 our combined dual data-motif clustering (DDMC) method in analyzing incomplete data sets, we designed
127 a computational experiment wherein we removed specific observations and predicted them using the
128 cluster centers corresponding to the peptides those missing values belonged to (Figure 2A). The resulting
129 mean squared errors between the actual and predicted values were compared to commonly used
130 imputation strategies such as the peptides' mean or minimum signal, constant zero, or matrix completion
131 by PCA. Furthermore, we evaluated the imputation performance of our method when clustering the data
132 using a different number of clusters. We observed that increasing the number of clusters improved the
133 imputation of missing values (Figure 2B-F). Additionally, we performed the same experiment by
134 clustering the data with different weights. Weight changes barely affected imputation performance,
135 indicating that cluster centers based on sequence only imputed missing values as accurately as when using
136 the phosphorylation signal (Figure 2F-I). These results indicate that DDMC clearly outperforms standard
137 imputation strategies such as using constant zero or the peptides' mean or minimum signal and imputes
138 missing values with similar accuracy to matrix completion by PCA.

139 **DDMC correctly identifies AKT1 and ERK2 as upstream kinases of signaling clusters containing 140 their substrates**

141 DDMC is a tailored method that clusters MS-generated phosphosites using its phosphorylation behavior
142 and sequence information. A major benefit of modeling the sequence information is the construction of
143 cluster motifs which can be useful to infer what putative upstream kinases might preferentially target
144 peptides of a specific cluster. To validate its ability to make upstream kinase predictions, we used DDMC
145 to cluster the phosphoproteomic measurements of MCF7 cells treated with a panel of 61 drug inhibitors
146 reported by Hijazi et al (20). PCA analysis of the resulting cluster centers clearly identified an inverse
147 correlation between the scores of AKT/mTOR targeted inhibitors and the loading of cluster 1, indicating
148 that the cluster's overall signal is attenuated by the presence of these compounds (Figure 3A-B).
149 Additional inhibitors targeting PDK1, FLT3, and S6K were also negatively correlated with cluster 1.
150 While these do not directly inhibit AKT1/mTOR, they are all known regulators of the pathway. A
151 heatmap displaying cluster's 1 phosphorylation signal across treatments corroborates that the abundance
152 of these peptides is substantially decreased when treated with AKT/mTOR/PIK3 inhibitors (Figure 3C).
153 Encouragingly, the specificity profile of AKT—within a collection of 55 different kinase PSPL
154 matrices—most closely matches the PSSM of cluster 1 (Figure 3D). Additionally, NetPhorest identified
155 AKT as the second top scoring upstream kinase of cluster 1, further corroborating DDMC's prediction.

156 Next, we extracted the sequences of ERK2 substrates identified in Carlson *et al* to create an “artificial”
157 ERK2-specific PSSM positive control (ERK2+ motif) (Figure 3F). As expected, ERK2 was predicted to
158 be the upstream kinase with the highest preference for the cluster's motif (Figure 3G). As an additional
159 test, given the consistent enrichment of hydrophobic and polar residues throughout the entire ERK2 target
160 motif (Figure 3F), we asked whether randomly shuffling all cluster PSSM positions surrounding the
161 phosphoacceptor residue would affect the upstream kinase prediction. This experiment led to a 2-fold
162 increase in the distance between ERK2 specificity profile and the ERK2+ motif (Figures 3G and H). We
163 subjected those clusters from the CPTAC data set that were preferentially favored by ERK2 to the same
164 experiment. As expected, we observed a similar decline in specificity between the clusters PSSMs and
165 ERK2 PSPL matrix (Figures 3H). Note that the noticeable difference in prediction between the ERK2+

166 motif and the CTPAC ERK2 motifs is not surprising given that while the former group contains only 26
167 peptides, the CPTAC clusters contain ~500–2000 phosphosites. Overall, this experiment generally shows
168 that despite the homogenous biophysical properties of ERK2 target motif across positions, the relative
169 enrichment of hydrophobic and polar residues in each position determines the extent to which ERK2
170 favors a particular motif (Figures 3G and H). Altogether, these results illustrate two different validation
171 scenarios in which DDMC successfully identifies the upstream kinases regulating clusters.

172 **A dual data-motif strategy improves prediction of different phenotypes and provides more robust** 173 **clustering**

174 As shown later in this study (Figures 5, 6, 7), we utilized DDMC to analyze the phosphoproteomes of 110
175 treatment-naïve LUAD tumors and 101 paired normal adjacent tissues (NATs) from the NCI’s CPTAC
176 LUAD study. We used DDMC with the binomial sequence distance method and 24 clusters (Figure 1,
177 2B). We were able to include 30,561 peptides that were not observed in every tumor through our ability to
178 handle missing data, but still filtered out 11,822 peptides that were only captured in one 10-plex TMT
179 run. We used this fitting result throughout the rest of this study. The resulting 24 cluster motifs can be
180 found in Figure S2.

181 To evaluate the benefit of incorporating the peptide sequence information into the clustering criterion, we
182 asked whether utilizing DDMC with different sequence weights would affect the performance of a
183 regularized logistic regression model that predicts the mutational status of STK11, whether a patient
184 harbors a mutation in EGFR and/or a gene fusion in ALK (EGFRm/ALKf), and the level of tumor
185 infiltration (“Hot” versus “Cold”). We found that for all three phenotypes, when the method only uses the
186 phosphorylation signal (weight=0), the patient samples are classified with lesser accuracy compared with
187 when a combination of both data and sequence is used. In the case of STK11, the use of the largest weight
188 wherein mainly the sequence motifs are used for clustering provided the best prediction performance.
189 Likewise, EGFRm/ALKf samples were best classified with a mix weight of 15 or 50. Finally, the
190 regression model classifying whether a sample is “hot-tumor-enriched” (HTE) or “cold-tumor-enriched”
191 (CTE) showed the best fitness with a weights of 10, 35, and 40. Together, these results indicate that
192 observing the motif information during clustering leads to final clusters that enhance the performance of
193 downstream phenotype prediction models (Figures 4A and S3).

194 Next, we explored how using different weights affects the overall phosphorylation signal and sequence
195 information of the resulting clusters. To do so, we compared the model behavior after clustering the
196 CPTAC data with a weight of 0 (peptide abundance only), 20 (mix), and 50 (mainly sequence). First, we
197 hypothesized that the abundance-only model would generate clusters wherein its members would show
198 less variation in phosphorylation signal and thus a lower mean squared error (MSE). To test this, we
199 computed the average peptide-to-cluster MSE of 2000 randomly selected peptides for each model across
200 all clusters. Although the differences were not significant, we did observe a direct correlation between
201 weight and MSE (Figure 4B). Next, we calculated the cumulative PSSM enrichment by summing the
202 sequence information (bits) of all cluster PSSMs per model. As expected, increasing the weight led to a
203 corresponding increase in the cumulative sequence information (Figure 4C). To further illustrate the
204 clustering behavior, we tracked the phosphosite TBC1D5 S584-p in the three models. Consistent with the
205 general trend, the abundance-only and mixed models generated lower p-signal MSE when compared to its
206 cluster center than the Sequence model whereas weight correlated with the total PSSM enrichment
207 (Figures 4D-E). Next, we quantified whether in addition to an increase in absolute enrichment, the mixed
208 and sequence-only models generated more similar cluster motifs to TBC1D5 S584-p sequence than the
209 abundance-only model. To do so, we computed the mean of all pairwise PAM250 scores between the
210 query sequence and all cluster sequences across models which clearly confirmed that as the sequence
211 prioritization of the model increases, the cluster PSSM is not only more enriched across all positions but
212 also displays a more representative sequence of TBC1D5 phosphosite (Figures 4F-I). These results show
213 that using a mixed weight that similarly prioritizes both information sources—peptide abundance and

214 sequence—leads to more robust clustering of phosphosites through a tradeoff between phosphorylation
215 abundance and sequence motifs.

216 **Widespread, dramatic signaling differences exist between tumor and normal adjacent tissue**

217 We explored whether DDMC could recognize conserved signaling patterns in tumors compared to normal
218 adjacent tissue (NAT). The signaling difference between tumors and NAT samples was substantial,
219 highlighting the significant signaling rewiring that tumor cells must undergo (Figure 5A). Using principal
220 components analysis, we could observe that NAT samples were more similar to one another than to each
221 tumor sample (Figure 5B/C). Nearly every cluster was significantly different in its average abundance
222 between tumor and NAT (Figure 5D). Not surprisingly given these enormous differences, samples could
223 be almost perfectly classified using their phosphopeptide signatures, with or without DDMC (Figures 5E;
224 S4). Using the DDMC clusters, a logistic regression model identified that NAT versus tumor status could
225 be predicted with cluster 11 alone (Figure 5C).

226 With the abundance changes and regression results we observed, we decided to further explore clusters 11
227 and 12. Cluster 11 shows a PSSM motif that might correspond to NEK1, 2, and 4, and an enrichment of
228 peptides involved in gas and oxygen transport, as well as cytoskeleton remodeling or migration-related
229 phenotypes according to a Gene Ontology (GO) analysis (Figure 5G/I). Even though NEKs are a largely
230 understudied family of serine/threonine kinases, NEK1/2 have an established role in the formation and
231 disassembly of cilia and NEK4 has also been implicated in regulating microtubule dynamics and stability
232 (22, 23). The primary cilium serves as a signaling hub via the local expression of cell surface receptors
233 and signaling molecules to sense environmental stimuli and thus promote a handful of phenotypes
234 including adaptation to hypoxia, migration, and escape from apoptosis (24, 25). Cancer cells typically
235 lack cilia which could promote the emergence of these malignant phenotypes. Cluster 11 displays a
236 striking phosphorylation decrease in tumor samples compared with NATs which could be representative
237 of the presence or lack of NEK1/2 signaling, respectively. Within this group of peptides, there is a notable
238 overrepresentation of hemoglobin subunits (HBG1, HBD, HBB, and HBA2) which could illustrate the
239 different oxygenation status of NATs versus malignant tissues. Moreover, several cytoskeletal-
240 remodeling proteins are present in cluster 11 such as PEAK1, FLNA, GAS2L2, MARCKS, PEAK1, and
241 ARHGEF7. The abundance of all these signaling molecules is substantially decreased in tumor compared
242 to NAT samples (Figure 5K).

243 On the other hand, cluster 12 was clearly identified as a CK2-like motif (Figure 5G). This association was
244 also established by NetPhorest which identified multiple experimentally validated CK2 substrates in this
245 cluster (Figure 5J). GO analysis of cluster 12 identified a substantial enrichment of negative regulators of
246 DNA duplex unwinding and pre-replicative complex assembly involved in cell cycle DNA replication
247 (Figure 5G, I-J). DNA duplex unwinding and replication are important processes that play a major role in
248 maintaining genome stability. DNA helicases are the enzymes responsible for unwinding the DNA and
249 thus are essential for DNA replication. As such, they have been widely associated with DNA damage
250 response (DDR) and cancer development (26). CK2 has been widely implicated in modulating DNA
251 repair signaling pathways in response to DNA damage to promote cell survival in cancer (27–29). In fact,
252 a study found that the CK2 inhibitor CX-4945 blocked DDR induced by gemcitabine and cisplatin and
253 synergizes with these compounds in ovarian cancer cell lines (30). Cluster 12 contains several signaling
254 proteins related to DNA replication and genome stability such as MCM3/4, the p53 interactor TP53BP1,
255 BRCA1, ATRX, CENPF, and CDKs whose signal is strikingly decreased in NATs and increased in tumor
256 samples (Figure 5L). These results, therefore, suggest that CK2 might activate signaling molecules within
257 cluster 12 involved in DNA repair pathways to induce the survival of cancer cells. Taken together,
258 DDMC builds phosphoproteomic clusters that present signaling dysregulation common to tumors
259 compared to NATs and identifies putative upstream kinases modulating them. These features can help to
260 interpret phosphoproteomic results and inform the generation of hypotheses for follow up experiments.

261 Genetic driver mutations are associated with more targeted phosphoproteomic rewiring

262 Inactivating somatic mutations in STK11 lead to increased tumorigenesis and metastasis (31). Thus, we
263 aimed to identify the phosphoproteomic aberrations triggered by this genetic event. The majority of
264 clusters were significantly altered, generally toward higher abundances with a mutation (Figure 6A). The
265 cluster centers corresponding to each patient's tumor and NAT samples could successfully predict the
266 STK11 mutational status by regularized logistic regression (Figure 6B). The tumor phosphoproteomic
267 signal of cluster 7 greatly contributed to classify mutant STK11 samples, whereas the tumor signal of 8
268 and 14 helped classify WT STK11 specimens. (Figure 6C). These results motivated further exploration of
269 clusters 7 and 8 which present sequence motifs favored by ERK2, and CK1/BRCA1/PKD, respectively
270 (Figure 6D).

271 Cluster 7 is highly enriched with peptides involved in regulation of the cell cycle by cohesin loading
272 (Figure 6E). Cohesin is a protein complex that mediates sister chromatid cohesion by directly binding
273 with DNA. This interaction holds both chromatids together after DNA replication until anaphase wherein
274 cohesin is removed to facilitate chromosome segregation during cell division. Cluster 7 contains the
275 inhibitor phosphosite of the tumor suppressor RB1 S795-p, the member of the cohesin loading complex
276 NIPBL (S280-p, S280-p;S284-p, and S350-p), and the cohesin release factor WAPL (S221-p and S221-
277 p;S223-p). Studies have shown that RB1 inactivation can lead to defects in chromosome cohesion that in
278 turn compromises chromosome stability (32, 33). Manning et al demonstrated that depletion of WAPL in
279 RB1-deficient cells promoted cohesin association with chromatin (33). Among these phosphosites, we
280 observed strong opposing signals between STK11 WT and mutant patients in NIPBL S280-p; WAPL
281 S221-p, S223-p; and RB1 S795-p (Figure 6E) which reinforces the association between STK11 activity
282 and chromatin instability. Moreover, CDCA5 is key regulator of sister chromatid cohesion by stabilizing
283 cohesin complex association with chromatin and was identified as a prognostic factor of lung cancer
284 through a tumor tissue microarray analysis of 262 non-small cell lung cancer (NSCLC) patients (34).
285 They showed that CDCA5 phosphorylation of S209 by ERK2 enhanced cell proliferation (34). Therefore,
286 these results might suggest that mutations inactivating mutations in STK11 might correlate with signaling
287 defects in sister chromatid cohesion during the cell cycle which in turn lead to chromosome instability
288 and cell cancer growth. In fact, STK11 inactivation has been associated with genomic instability,
289 although the signaling mechanism underlying this phenotypic response remains elusive (35).

290 The signal of phosphosites in cluster 8, specifically in tumor samples, largely contributes to predict the
291 signaling differences between STK11 WT and mutant samples (Figure 6C). This cluster presents a clear
292 enrichment of peptides involved in the regulation of the Golgi apparatus such as GOLGA2-5, GOLGB1,
293 and GOLPH3 (Figure 6F). Cancer cells commonly undergo fragmentation of the Golgi which has been
294 shown to drive several malignant molecular signatures including the hyperactivity of motor proteins and
295 kinase signaling dysregulation (37). Myosin 18A and 1E pertain to cluster 18 and the former has been
296 reported to interact with GOLPH3 to induce Golgi dispersal. Moreover, a series of studies uncovered that
297 GOLPH3 promotes cell proliferation in cancer (38-40). The phosphorylation behavior of GOLPH3,
298 Myosin 18A, and GOLGA2 in STK11 WT compared with STK11 mutant patients shows a dramatic
299 increase of abundance in the latter which supports the association between STK11 activity and an
300 oncogenic role of the Golgi apparatus in these patients (Figure 6E). Together, these results suggest that
301 STK11 mutations in tumor samples could affect the dispersion of the Golgi apparatus compared with
302 STK11 WT samples.

303 Tyrosine kinase inhibitors (TKIs) targeting the receptor tyrosine kinases (RTKs) EGFR and ALK are
304 effective treatments in cancer patients with EGFR mutations and/or ALK translocations (EGFRm/ALKf).
305 However, these treatments are limited by drug resistance which in some cases can be mediated by the
306 concomitant signaling of both RTKs activated by driver mutations (41, 42). Once again, the signaling
307 cluster centers allowed a regularized logistic regression model to more accurately classify samples
308 according to its EGFRm/ALKf status (Figure S5).

309 Finally, we compared the classification performance of four regularized logistic regression models fit to
310 either the DDMC clusters, clusters generated by the standard methods GMM and k-means, or the raw
311 phosphoproteomic data directly. It is worth noting that unlike DDMC, methods such as GMM, k-means,
312 or direct regression cannot handle missing values and thus for these strategies we used the 1,311 peptides
313 that were observed in all samples, whereas DDMC was fit to the entire data set comprising 30,561
314 phosphosites. We found that samples were classified with higher accuracy using DDMC compared to a
315 GMM and with similar performance to k-means, especially with STK11 (Figure S6A). Direct regression
316 to the raw signaling data yielded excellent performance; however, this strategy assigns thousands of
317 coefficients to different peptides that vary every time the model is run, rendering this approach unable to
318 establish a consistent link between mutations and signaling (Figure S6). In contrast, our analysis identifies
319 a consistent association between STK11 activity with two novel phenotypes, namely chromosome
320 cohesion during cell cycle and Golgi fragmentation, and proposes putative signaling mechanisms to
321 support it.

322 **Exploration of immune infiltration-associated signaling patterns in tumors**

323 Immune checkpoint inhibitors (ICIs) have emerged as effective treatment options for NSCLC patients.
324 However, there still is a need to identify or influence which patients will respond to these therapies.
325 Patients that do not respond to ICIs often have tumors with poor immune infiltration either inherently or
326 via an adaptive process after long exposure to the drug. However, the signaling mechanism by which
327 malignant cells prevent tumor infiltration remains elusive. We used our DDMC clusters to explore the
328 shared signaling patterns that differentiate “hot-tumor-enriched” (HTE) from “cold-tumor-enriched” CTE
329 LUAD patients (11, 43). HTE and CTE status per patient was determined using xCell by Gillette et al
330 (11).

331 We observed that four clusters were significantly different in their average abundance between HTE and
332 CTE samples (Figure 7A). Cluster 17, 18, and 20 display significantly higher abundances in HTE
333 compared to CTE samples whereas cluster 21 presents the opposite trend. Samples could be accurately
334 classified using the DDMC clusters (Figure 7B). This predictive performance was mainly explained by a
335 positive association of cluster 2 with HTE status and cluster 6 with CTE. Other clusters contributed to
336 explain the signaling differences between both groups but to a lesser extent (Figure 7C).

337 These results prompted us to further investigate clusters 6, 17, 20, and 21 which our model predicts to be
338 regulated by CK1/PKA, STK11/p38, CK2/STK11, and ERK2, respectively (Figure 7D). When exploring
339 immunologically relevant phenotypes in the GO analysis of each cluster, we observed that clusters 6, 17,
340 and 20 showed a substantial over-representation of immunological processes. Conversely, neither of these
341 were present in the GO analyses of cluster 2 nor cluster 21 wherein the former substantially contributes to
342 predict CTE samples and the latter shows a significant increase of phosphorylation abundance in CTE
343 over HTE samples (Figures 7A and C). A gene ontology analysis indicates that cluster 6 members are
344 particularly involved in mediating B cell homeostasis, but also T cell differentiation, T cell receptor
345 signaling, and regulation of T cell activation. These processes are promoted, at least in part, by ABL1,
346 LCK, PAK1, and DOCK10/11 which show an increased abundance in HTE and are attenuated in CTE
347 samples (Figures 7E & H). Cluster 17 GO analysis unveiled an over-representation of several innate and
348 adaptive immune response pathways possibly involving CD44, SDK1, PKC, PLD1, CAPN1 and GSTP1.
349 For instance, CD44 is expressed in both endothelial and immune cells and its regulation plays a key role
350 in enabling neutrophil and lymphocyte recruitment into tissues (44, 45) (Figures 7F & I). A study found
351 that the osteopontin (OPN)/CD44 interaction is an immune checkpoint that controls CD8+ T cell
352 activation and tumor immune evasion in which elevated expression of OPN correlated with decreased
353 patient survival and conferred host tumor immune tolerance. Cluster 20 is enriched in responses
354 orchestrated by the innate immune system (Figures 7G and J). The transcription factor NFATC crucially
355 promotes T cell activation and proliferation, and several studies show that the predicted upstream kinase
356 of cluster 20 CK2 directly phosphorylates this protein and enhances its gene expression (46, 47). In

357 addition, CK2 has also been shown to phosphorylate Regulators of Calcineurin (RCAN) proteins, which
358 indirectly inhibit NFATC function (48). Several RCAN and NFATC peptides are present in cluster 20,
359 however S210-p and S366-p, respectively show the largest abundance difference between HTE and CTE.
360 Unexpectedly, RCAN1 S210-p shows a higher signal in HTE than in CTE whereas NFATC3 S366-p
361 presents the opposite trend which might indicate that both phosphorylation events are inhibitory.
362 Together, these results reinforce the role of CK2 in promoting immune infiltration in lung cancer patients.
363 Intriguingly, inactivating mutations in STK11 have been reported to promote anti-PD1/PD-L1 resistance
364 in KRAS-mutant LUAD suggesting a key role of STK11 in promoting tumor immune infiltration (49).
365 Overall, these data demonstrate that the presence or lack of tumor immune infiltration can be accurately
366 predicted by the DDMC clusters which in turn help identify putative upstream kinases modulating
367 immune evasion.

368 Discussion

369 Phosphorylation-based cell signaling through the coordinated activity of protein kinases enables cells to
370 swiftly integrate environmental cues and orchestrate a myriad of biological processes. MS-based global
371 phosphoproteomic data provides the unique opportunity to globally interrogate signaling networks to
372 better understand cellular decision-making and its therapeutic implications. However, these data also
373 present challenging issues as a consequence of their incomplete and stochastic coverage, high-content but
374 low-sample throughput, and variation in coverage across experiments. Here, we propose a clustering
375 method, Dual Data and Motif Clustering (DDMC), that untangles highly complex coordinated signaling
376 changes by grouping phosphopeptides based on their phosphorylation behavior and sequence similarity
377 (Figure 1). To test the utility of DDMC, we clustered the phosphoproteomes of LUAD patients and used
378 the resulting groups of peptides to decipher signaling dysregulation common to tumors, genetic
379 backgrounds, and tumor infiltration status (Figures 5, 6, 7).

380 Previous efforts in regressing mass spectrometry-based phosphorylation measurements against
381 phenotypic or clinical data have been based on the ability of certain regression models such as PLSR or
382 LASSO to robustly predict using high-dimensional and correlated data (50). While these models can
383 generally be predictive with such data, they are not easily interpretable (Figure S4B). Hence, we
384 hypothesized that clustering large-scale MS measurements based on biologically meaningful features and
385 utilizing the cluster centers to fit regression methods could enhance the predictive performance of the
386 model while providing highly interpretable results wherein clusters constitute signaling nodes distinctly
387 correlated with cell patient phenotypes. Here, we demonstrate that DDMC enhances model prediction and
388 interpretation (Figures 4A, S6, 3).

389 Model interpretation is enhanced by comparing the resulting cluster PSSMs with kinase specificity data
390 such as PSPL to identify putative upstream kinases modulating signaling clusters. Computational
391 validations showed that DDMC was able to correctly associate AKT1 and ERK2 with clusters of their
392 respective substrates (Figure 3). It is worth noting, however, that kinase specificity is defined by
393 additional features beyond the phosphosite motif such as kinase-substrate co-localization, regulation by
394 phosphosite-binding domains (e.g., SH2, PTB domains), or docking. In addition, a major limitation of
395 PSPL experiments is that since they do not provide docking information, the real affinity between the
396 string of identified peptide residues as key determinants of specificity of a sequence motif and the
397 interacting kinase domain is unknown. This limitation could also compromise kinase-cluster associations
398 established by DDMC. A method combining bacterial surface-display of peptide libraries with next-
399 generation sequencing tackles this limitation by quantifying the specificity of a kinase to virtually all
400 possible motif combinations (51). Thus, as the number of profiled kinases with this technique increases,
401 these measurements could be used to rank cluster peptides by magnitude of specificity to a specific kinase
402 to make better upstream kinase predictions.

403 A key benefit of DDMC is that the identified clusters are not limited to pre-existing motifs and are
404 therefore not dependent on prior experimentally validated kinase-substrate interactions. Thereby, this
405 method could improve our understanding of the signaling effects of understudied kinases. For instance,
406 our model predicts NEK1&2 promote, at least in part, a cluster with strikingly increased signaling in
407 NATs compared to tumors. Further exploration of this cluster led us to hypothesize that the lack of NEK
408 signaling in tumor samples might associated with the absence of ciliogenesis and adaptation to hypoxia in
409 lung tumors (Figure 5G-H). Additionally, we show that cluster 8, which greatly contributes to explain the
410 signaling differences between STK11 WT and mutant samples in tumors (Figure 6C), is enriched with
411 proteins such as GOLPH3 and Myosin 18A that have been shown to promote Golgi fragmentation in
412 cancer (38–40). This prompts us to consider the novel interaction between CK1 and these signaling
413 molecules.

414 An additional major challenge being faced during the analysis of large-scale signaling data is
415 missingness. Given that statistical tools often require complete data sets, researchers use standard
416 methods to impute missing values such as the peptides' mean or minimum signal, constant zero, or PCA
417 imputation only in peptides wherein at least 50% of their samples were required to have non-missing
418 values as excessive missing values can result in poor imputation (10, 11, 52). In this study we show that
419 DDMC can model a data set of 30,561 peptides after filtering out any phosphosites that were not captured
420 in at least 2 TMT (up to ~80% of missingness) by ignoring unobserved values during EM distribution
421 estimation and calculation of GMM probabilities (see methods). Therefore, this method enables clustering
422 of signaling data despite a remarkable number of missing values. Furthermore, DDMC clearly
423 outperforms the imputation performance of using the peptides' mean, minimum signal, or constant zero
424 and provides similar results to PCA imputation. This important feature could offer the possibility of
425 conducting pan-cancer phosphoproteomics studies using readily available large-scale clinical
426 phosphoproteomic data.

427 The benefit of building algorithms combining different information sources is evident in previously
428 published approaches. For instance, INKA predicts active kinases by integrating scores reflecting both
429 phosphorylation status and substrate abundance (53). In another study, Exarchos et al. formulated a
430 decision support system that integrates clinical, imaging, and genomic data to identify the factors that
431 contribute to oral cancer progression and predict relapses. The authors found that combining the more
432 accurate individual predictors yielded better predictions than those generated by other strategies reported
433 in the literature (54). Finally, BOADICEA is a method that allows systematic risk stratification of breast
434 cancer patients by incorporating the effects of lifestyle, hormonal and reproductive risk factors,
435 mammographic density, and of the common breast cancer susceptibility genetic variants into the
436 prediction model (55).

437 In total, in this study we show that combining the information about the sequence features and
438 phosphorylation abundance leads to more robust clustering of global signaling measurements. Use of the
439 DDMC clusters to regress against cell phenotypes led to enhanced model predictions and interpretation.
440 Thus, we propose DDMC as a general and flexible strategy for phosphoproteomic analysis.

441 **Materials and Methods**

442 All analysis was implemented in Python v3.9 and can be found at <https://github.com/meyer-lab/resistance-MS>.

444 **Expectation-maximization (EM) algorithm architecture**

445 We constructed a modified mixture model that clusters peptides based on both their abundance across
446 conditions and sequence. The model is defined by a given number of clusters and weighting factor to
447 prioritize either the data or the sequence information. Fitting was performed using expectation-

448 maximization, initialized at a starting point. The starting point was derived from k-means clustering the
449 abundance data after missing values were imputed by PCA with a component number equal to the number
450 of clusters. During the expectation (E) step, the algorithm calculates the probability of each peptide being
451 assigned to each cluster. In the maximization (M) step, each cluster's distributions are fit using the
452 weighted cluster assignments. The peptide sequence and abundance assignments within the E step are
453 combined by taking the sum of the log-likelihood of both assignments. The peptide log-likelihood is
454 multiplied by the user-defined weighting factor immediately before to influence its importance. Both
455 steps repeat until convergence as defined by the increase in model log-likelihood between iterations
456 falling below a user-defined threshold.

457 **Phosphorylation site abundance clustering in the presence of missing values**

458 We modeled the log-transformed abundance of each phosphopeptide as following a multivariate Gaussian
459 distribution with diagonal covariance. Each dimension of this distribution represents the abundance of
460 that peptide within a given sample. For example, within a data set of 100 patients and 1000 peptides,
461 using 10 clusters, the data is represented by 10 Gaussian distributions of 100 dimensions.
462 Unobserved/missing values were indicated as NaN and ignored during both distribution estimation and
463 when calculating probabilities. Any peptides that were detected in only one TMT experiment were
464 discarded.

465 **Sequence-cluster comparison**

466 *PAM250*

467 During model initialization, the pairwise distance between all peptides in the dataset was calculated using
468 the PAM250 matrix. The mean distance from each peptide to a given cluster could then be calculated by:

$$469 \quad w = \frac{1}{n} (P \cdot v)$$

470 Where P is the $n \times n$ distance matrix, n is the number of peptides in the dataset, v is the probability of
471 each peptide being assigned to the cluster of interest, and w is the log-probabilities of cluster assignment.

472 *Binomial enrichment*

473 We alternatively used a binomial enrichment model for the sequence representation of a cluster based on
474 earlier work (55). Upon model initialization, a background matrix $i \times j \times k$ was created with a position-
475 specific scoring matrix of all the sequences together. Next, an T data tensor i was created where j is the
476 number of peptides, k is the number of amino acid possibilities, and k is the position relative to the
477 phosphorylation site. This tensor contained 1 where an amino acid was present for that position and
478 peptide, and 0 elsewhere.

479 Within each iteration, the cluster motif would be updated using v , the probability of each peptide being
480 assigned to the cluster of interest. First, a weighted count for each amino acid and position would be
481 assembled:

$$482 \quad k = (T^T \cdot v)^T$$

483 Because peptides can be partially assigned to a cluster, the counts of each amino acid and position can
484 take continuous values. We therefore generalized the binomial distribution to allow continuous values
485 using the regularized incomplete Beta function:

$$486 \quad M = B(\| \vec{v} \|_1 - k, k + 1, 1 - G)$$

487 Finally, the log-probability of membership for each peptide was calculated based on the product of each
488 amino acid-position probability.

$$489 \quad w = \log(T \times M)$$

490 We confirmed that this provided identical results to a binomial enrichment model for integer counts of
491 amino acids (55) but allowed for partial assignment of peptides to clusters.

492 **Quantifying the influence of sequence versus data**

493 The magnitude of the weight used to scale the sequence and data scores is arbitrary. We do know that
494 with a weight of 0 the model only uses the phosphorylation measurements. Alternatively, with an
495 enormously large weight the motif information is prioritized. However, we do not know to what extent
496 each information source is prioritized in general. Therefore, to quantify the relative importance of each
497 type of data, we calculated our clustering results at each weighting extreme, and then calculated the
498 Frobenius norm of the resulting peptide assignments between those and the clustering of interest.

499 **Generating Cluster Motifs and Upstream Kinase Predictions**

500 For each cluster we computed a position-specific-scoring matrix (PSSM). To do so, we populated a
501 residue/position matrix with the sum of the corresponding cluster probabilities for every peptide. Once all
502 peptides were accounted for, the resulting matrix was normalized by averaging the mean probability
503 across amino acids and log₂-transformed to generate a PSSM. In parallel, we computed a PSSM
504 including all sequences that served as background to account for the different amino acid occurrences
505 within the data set. Then, we subtracted each cluster PSSM with the background PSSM and limited any
506 large negative numbers to -3. Next, we extracted several kinase specificity profiling results from the
507 literature (16, 18, 18, 19). The distance between PSSM and PSSL motifs was calculated using by the
508 Frobenius norm of the difference. Motif logo plots were generated using logomaker (56).

509 **Evaluate clustering by imputation of values**

510 To evaluate the ability of our model to handle missing values, we removed random, individual TMT
511 experiments for each peptide and used the model to impute these values. The number of missing values
512 per peptide is highly variable. Therefore, in our error quantitation, we stratified peptides by their
513 missingness percentage and computed the average mean squared error between the actual values and
514 predictions—or imputed peptide average—in each group. We calculated the reconstruction error across
515 different combinations of cluster numbers and weights using the same process.

516 **Associating clusters with molecular and clinical features**

517 To find clusters that tracked with specific molecular or clinical features we implemented two different
518 strategies: logistic regression and hypothesis testing. For binary problems such as Tumor vs NAT samples
519 or mutational status we used l1-regularized logistic regression and Mann-Whitney rank tests. In the
520 former, we tried to predict the feature of interest using the phosphorylation signal of the cluster centers,
521 whereas in the latter, for each cluster we split all patients according to their specific feature and tested
522 whether the difference in the median signal between both groups was statistically different. We performed
523 Bonferroni correction on the p-values computed by the Mann-Whitney tests. Gene ontology analysis was
524 performed using the GENEONTOLOGY software (geneontology.org) (57, 58).

525 **References**

- 526 1. T. Hunter, Protein kinases and phosphatases: the yin and yang of protein phosphorylation and
527 signaling. *Cell*. **80**, 225–36 (1995).
- 528 2. M. B. Yaffe, Why geneticists stole cancer research even though cancer is primarily a signaling disease.
529 *Sci Signal*. **12** (2019), doi:10.1126/scisignal.aaw3483.
- 530 3. P. Casado, J.-C. Rodriguez-Prados, S. C. Cosulich, S. Guichard, B. Vanhaesebroeck, S. Joel, P. R.
531 Cutillas, Kinase-substrate enrichment analysis provides insights into the heterogeneity of signaling
532 pathway activation in leukemia cells. *Sci Signal*. **6**, rs6 (2013).
- 533 4. R. Beekhof, C. van Alphen, A. A. Henneman, J. C. Knol, T. V. Pham, F. Rolfs, M. Labots, E.
534 Henneberry, T. Y. Le Large, R. R. de Haas, S. R. Piersma, V. Vurchio, A. Bertotti, L. Trusolino, H. M.
535 Verheul, C. R. Jimenez, INKA, an integrative data analysis pipeline for phosphoproteomic inference of
536 active kinases. *Mol Syst Biol*. **15**, e8981 (2019).
- 537 5. P. V. Hornbeck, J. M. Kornhauser, V. Latham, B. Murray, V. Nandhikonda, A. Nord, E. Skrzypek, T.
538 Wheeler, B. Zhang, F. Gnad, 15 years of PhosphoSitePlus®: integrating post-translationally modified
539 sites, disease variants and isoforms. *Nucleic Acids Res*. **47**, D433–D441 (2019).
- 540 6. R. Linding, L. J. Jensen, G. J. Ostheimer, M. A. T. M. van Vugt, C. Jørgensen, I. M. Miron, F. Diella,
541 K. Colwill, L. Taylor, K. Elder, P. Metalnikov, V. Nguyen, A. Pasculescu, J. Jin, J. G. Park, L. D.
542 Samson, J. R. Woodgett, R. B. Russell, P. Bork, M. B. Yaffe, T. Pawson, Systematic discovery of in vivo
543 phosphorylation networks. *Cell*. **129**, 1415–26 (2007).
- 544 7. J. C. Obenauer, L. C. Cantley, M. B. Yaffe, Scansite 2.0: Proteome-wide prediction of cell signaling
545 interactions using short sequence motifs. *Nucleic Acids Res*. **31**, 3635–41 (2003).
- 546 8. E. J. Needham, B. L. Parker, T. Burykin, D. E. James, S. J. Humphrey, Illuminating the dark
547 phosphoproteome. *Sci Signal*. **12** (2019), doi:10.1126/scisignal.aau8645.
- 548 9. D. L. Tabb, L. Vega-Montoto, P. A. Rudnick, A. M. Variyath, A.-J. L. Ham, D. M. Bunk, L. E.
549 Kilpatrick, D. D. Billheimer, R. K. Blackman, H. L. Cardasis, S. A. Carr, K. R. Clauser, J. D. Jaffe, K. A.
550 Kowalski, T. A. Neubert, F. E. Regnier, B. Schilling, T. J. Tegeler, M. Wang, P. Wang, J. R. Whiteaker,
551 L. J. Zimmerman, S. J. Fisher, B. W. Gibson, C. R. Kinsinger, M. Mesri, H. Rodriguez, S. E. Stein, P.
552 Tempst, A. G. Paulovich, D. C. Liebler, C. Spiegelman, Repeatability and reproducibility in proteomic
553 identifications by liquid chromatography-tandem mass spectrometry. *J Proteome Res*. **9**, 761–76 (2010).
- 554 10. Y.-J. Chen, T. I. Roumeliotis, Y.-H. Chang, C.-T. Chen, C.-L. Han, M.-H. Lin, H.-W. Chen, G.-C.
555 Chang, Y.-L. Chang, C.-T. Wu, M.-W. Lin, M.-S. Hsieh, Y.-T. Wang, Y.-R. Chen, I. Jonassen, F. Z.
556 Ghavidel, Z.-S. Lin, K.-T. Lin, C.-W. Chen, P.-Y. Sheu, C.-T. Hung, K.-C. Huang, H.-C. Yang, P.-Y.
557 Lin, T.-C. Yen, Y.-W. Lin, J.-H. Wang, L. Raghav, C.-Y. Lin, Y.-S. Chen, P.-S. Wu, C.-T. Lai, S.-H.
558 Weng, K.-Y. Su, W.-H. Chang, P.-Y. Tsai, A. I. Robles, H. Rodriguez, Y.-J. Hsiao, W.-H. Chang, T.-Y.
559 Sung, J.-S. Chen, S.-L. Yu, J. S. Choudhary, H.-Y. Chen, P.-C. Yang, Y.-J. Chen, Proteogenomics of
560 Non-smoking Lung Cancer in East Asia Delineates Molecular Signatures of Pathogenesis and
561 Progression. *Cell*. **182**, 226–244.e17 (2020).
- 562 11. M. A. Gillette, S. Satpathy, S. Cao, S. M. Dhanasekaran, S. V. Vasaikar, K. Krug, F. Petralia, Y. Li,
563 W.-W. Liang, B. Reva, A. Krek, J. Ji, X. Song, W. Liu, R. Hong, L. Yao, L. Blumenberg, S. R. Savage,
564 M. C. Wendl, B. Wen, K. Li, L. C. Tang, M. A. MacMullan, S. C. Avanesian, M. H. Kane, C. J. Newton,
565 M. Cornwell, R. B. Kothadia, W. Ma, S. Yoo, R. Mannan, P. Vats, C. Kumar-Sinha, E. A. Kawaler, T.
566 Omelchenko, A. Colaprico, Y. Geffen, Y. E. Maruvka, F. da Veiga Leprevost, M. Wiznerowicz, Z. H.

- 567 Gümüş, R. R. Veluswamy, G. Hostetter, D. I. Heiman, M. A. Wyczalkowski, T. Hiltke, M. Mesri, C. R.
568 Kinsinger, E. S. Boja, G. S. Omenn, A. M. Chinnaiyan, H. Rodriguez, Q. K. Li, S. D. Jewell, M.
569 Thiagarajan, G. Getz, B. Zhang, D. Fenyö, K. V. Ruggles, M. P. Cieslik, A. I. Robles, K. R. Clauser, R.
570 Govindan, P. Wang, A. I. Nesvizhskii, L. Ding, D. R. Mani, S. A. Carr., Proteogenomic Characterization
571 Reveals Therapeutic Vulnerabilities in Lung Adenocarcinoma. *Cell*. **182**, 200–225.e35 (2020).
- 572 12. A. Zarrinpar, S.-H. Park, W. A. Lim, Optimization of specificity in a cellular protein interaction
573 network by negative selection. *Nature*. **426**, 676–80 (2003).
- 574 13. C. S. H. Tan, A. Pasculescu, W. A. Lim, T. Pawson, G. D. Bader, R. Linding, Positive selection of
575 tyrosine loss in metazoan evolution. *Science*. **325**, 1686–8 (2009).
- 576 14. T. Obata, M. B. Yaffe, G. G. Leparc, E. T. Piro, H. Maegawa, A. Kashiwagi, R. Kikkawa, L. C.
577 Cantley, Peptide and protein library screening defines optimal substrate motifs for AKT/PKB. *J Biol*
578 *Chem*. **275**, 36108–15 (2000).
- 579 15. J. Mok, P. M. Kim, H. Y. K. Lam, S. Piccirillo, X. Zhou, G. R. Jeschke, D. L. Sheridan, S. A. Parker,
580 V. Desai, M. Jwa, E. Cameroni, H. Niu, M. Good, A. Remenyi, J.-L. N. Ma, Y.-J. Sheu, H. E. Sassi, R.
581 Sopko, C. S. M. Chan, C. De Virgilio, N. M. Hollingsworth, W. A. Lim, D. F. Stern, B. Stillman, B. J.
582 Andrews, M. B. Gerstein, M. Snyder, B. E. Turk, Deciphering protein kinase specificity through large-
583 scale analysis of yeast phosphorylation site motifs. *Sci Signal*. **3**, ra12 (2010).
- 584 16. B. van de Kooij, P. Creixell, A. van Vlimmeren, B. A. Joughin, C. J. Miller, N. Haider, C. D.
585 Simpson, R. Linding, V. Stambolic, B. E. Turk, M. B. Yaffe, Comprehensive substrate specificity
586 profiling of the human Nek kinome reveals unexpected signaling outputs. *Elife*. **8** (2019),
587 doi:10.7554/elife.44635.
- 588 17. J. E. Hutti, E. T. Jarrell, J. D. Chang, D. W. Abbott, P. Storz, A. Toker, L. C. Cantley, B. E. Turk, A
589 rapid method for determining protein kinase phosphorylation specificity. *Nat Methods*. **1**, 27–9 (2004).
- 590 18. M. J. Begley, C.-h. Yun, C. A. Gewinner, J. M. Asara, J. L. Johnson, A. J. Coyle, M. J. Eck, I.
591 Apostolou, L. C. Cantley, EGF-receptor specificity for phosphotyrosine-primed substrates provides signal
592 integration with Src. *Nat Struct Mol Biol*. **22**, 983–90 (2015).
- 593 19. M. L. Miller, L. J. Jensen, F. Diella, C. Jørgensen, M. Tinti, L. Li, M. Hsiung, S. A. Parker, J.
594 Bordeaux, T. Sicheritz-Ponten, M. Olhovsky, A. Pasculescu, J. Alexander, S. Knapp, N. Blom, P. Bork,
595 S. Li, G. Cesareni, T. Pawson, B. E. Turk, M. B. Yaffe, S. Brunak, R. Linding, Linear motif atlas for
596 phosphorylation-dependent signaling. *Sci Signal*. **1**, ra2 (2008).
- 597 20. M. Hijazi, R. Smith, V. Rajeeve, C. Bessant, P. R. Cutillas, Reconstructing kinase network topologies
598 from phosphoproteomics data reveals cancer-associated rewiring. *Nat Biotechnol*. **38**, 493–502 (2020).
- 599 21. S. M. Carlson, C. R. Chouinard, A. Labadorf, C. J. Lam, K. Schmelzle, E. Fraenkel, F. M. White,
600 Large-scale discovery of ERK2 substrates identifies ERK-mediated transcriptional regulation by ETV3.
601 *Sci Signal*. **4**, rs11 (2011).
- 602 22. L. Moniz, P. Dutt, N. Haider, V. Stambolic, Nek family of kinases in cell cycle, checkpoint control
603 and cancer. *Cell Div*. **6**, 18 (2011).
- 604 23. G. V. Meirelles, A. M. Perez, E. E. de Souza, F. L. Basei, P. F. Papa, T. D. Melo Hanchuk, V. B.
605 Cardoso, J. Kobarg, “Stop Ne(c)king around”: How interactomics contributes to functionally characterize
606 Nek family kinases. *World J Biol Chem*. **5**, 141–60 (2014).

- 607 24. L. Fabbri, F. Bost, N. M. Mazure, Primary Cilium in Cancer Hallmarks. *Int J Mol Sci.* **20** (2019),
608 doi:10.3390/ijms20061336.
- 609 25. O. V. Plotnikova, E. A. Golemis, E. N. Pugacheva, Cell cycle-dependent ciliogenesis and cancer.
610 *Cancer Res.* **68**, 2058–61 (2008).
- 611 26. R. M. Brosh, DNA helicases involved in DNA repair and their roles in cancer. *Nat Rev Cancer.* **13**,
612 542–58 (2013).
- 613 27. J. So, A. Pasculescu, A. Y. Dai, K. Williton, A. James, V. Nguyen, P. Creixell, E. M. Schoof, J.
614 Sinclair, M. Barrios-Rodiles, J. Gu, A. Krizus, R. Williams, M. Olhovsky, J. W. Dennis, J. L. Wrana, R.
615 Linding, C. Jorgensen, T. Pawson, K. Colwill, Integrative analysis of kinase networks in TRAIL-induced
616 apoptosis provides a source of potential targets for combination therapy. *Sci Signal.* **8**, rs3 (2015).
- 617 28. M. V. Bennetzen, D. H. Larsen, J. Bunkenborg, J. Bartek, J. Lukas, J. S. Andersen, Site-specific
618 phosphorylation dynamics of the nuclear proteome during the DNA damage response. *Mol Cell*
619 *Proteomics.* **9**, 1314–23 (2010).
- 620 29. A. J. Rabalski, L. Gyenis, D. W. Litchfield, Molecular Pathways: Emergence of Protein Kinase CK2
621 (CSNK2) as a Potential Target to Inhibit Survival and DNA Damage Response and Repair Pathways in
622 Cancer Cells. *Clin Cancer Res.* **22**, 2840–7 (2016).
- 623 30. A. Siddiqui-Jain, J. Bliesath, D. Macalino, M. Otori, N. Huser, N. Streiner, C. B. Ho, K. Anderes, C.
624 Proffitt, S. E. O'Brien, J. K. C. Lim, D. D. Von Hoff, D. M. Ryckman, W. G. Rice, D. Drygin, CK2
625 inhibitor CX-4945 suppresses DNA repair response triggered by DNA-targeted anticancer drugs and
626 augments efficacy: mechanistic rationale for drug combination therapy. *Mol Cancer Ther.* **11**, 994–1005
627 (2012).
- 628 31. H. Ji, M. R. Ramsey, D. N. Hayes, C. Fan, K. McNamara, P. Kozlowski, C. Torrice, M. C. Wu, T.
629 Shimamura, S. A. Perera, M.-C. Liang, D. Cai, G. N. Naumov, L. Bao, C. M. Contreras, D. Li, L. Chen,
630 J. Krishnamurthy, J. Koivunen, L. R. Chirieac, R. F. Padera, R. T. Bronson, N. I. Lindeman, D. C.
631 Christiani, X. Lin, G. I. Shapiro, P. A. Jänne, B. E. Johnson, M. Meyerson, D. J. Kwiatkowski, D. H.
632 Castrillon, N. Bardeesy, N. E. Sharpless, K.-K. Wong, LKB1 modulates lung cancer differentiation and
633 metastasis. *Nature.* **448**, 807–10 (2007).
- 634 32. A. L. Manning, M. S. Longworth, N. J. Dyson, Loss of pRB causes centromere dysfunction and
635 chromosomal instability. *Genes Dev.* **24**, 1364–76 (2010).
- 636 33. A. L. Manning, S. A. Yazinski, B. Nicolay, A. Bryll, L. Zou, N. J. Dyson, Suppression of genome
637 instability in pRB-deficient cells by enhancement of chromosome cohesion. *Mol Cell.* **53**, 993–1004
638 (2014).
- 639 34. M.-H. Nguyen, J. Koinuma, K. Ueda, T. Ito, E. Tsuchiya, Y. Nakamura, Y. Daigo, Phosphorylation
640 and activation of cell division cycle associated 5 by mitogen-activated protein kinase play a crucial role in
641 human lung carcinogenesis. *Cancer Res.* **70**, 5337–47 (2010).
- 642 35. R. H. Hruban, M. I. Canto, M. Goggins, R. Schulick, A. P. Klein, Update on familial pancreatic
643 cancer. *Adv Surg.* **44**, 293–311 (2010).
- 644 36. R. K. Gill, S.-H. Yang, D. Meerzaman, L. E. Mechanic, E. D. Bowman, H.-S. Jeon, S. Roy
645 Chowdhuri, A. Shakoori, T. Dracheva, K.-M. Hong, J. Fukuoka, J.-H. Zhang, C. C. Harris, J. Jen,
646 Frequent homozygous deletion of the LKB1/STK11 gene in non-small cell lung cancer. *Oncogene.* **30**,
647 3784–91 (2011).

- 648 37. A. Petrosyan, Onco-Golgi: Is Fragmentation a Gate to Cancer Progression? *Biochem Mol Biol J.* **1**
649 (2015), doi:10.21767/2471-8084.100006.
- 650 38. X. Hua, L. Yu, W. Pan, X. Huang, Z. Liao, Q. Xian, L. Fang, H. Shen, Increased expression of Golgi
651 phosphoprotein-3 is associated with tumor aggressiveness and poor prognosis of prostate cancer. *Diagn*
652 *Pathol.* **7**, 127 (2012).
- 653 39. Z. Zeng, H. Lin, X. Zhao, G. Liu, X. Wang, R. Xu, K. Chen, J. Li, L. Song, Overexpression of
654 GOLPH3 promotes proliferation and tumorigenicity in breast cancer via suppression of the FOXO1
655 transcription factor. *Clin Cancer Res.* **18**, 4059–69 (2012).
- 656 40. B.-S. Hu, H. Hu, C.-Y. Zhu, Y.-L. Gu, J.-P. Li, Overexpression of GOLPH3 is associated with poor
657 clinical outcome in gastric cancer. *Tumour Biol.* **34**, 515–20 (2012).
- 658 41. T. Sasaki, J. Koivunen, A. Ogino, M. Yanagita, S. Nikiforow, W. Zheng, C. Lathan, J. P. Marcoux, J.
659 Du, K. Okuda, M. Capelletti, T. Shimamura, D. Ercan, M. Stumpfova, Y. Xiao, S. Weremowicz, M.
660 Butaney, S. Heon, K. Wilner, J. G. Christensen, M. J. Eck, K.-K. Wong, N. Lindeman, N. S. Gray, S. J.
661 Rodig, P. A. Jänne, A novel ALK secondary mutation and EGFR signaling cause resistance to ALK
662 kinase inhibitors. *Cancer Res.* **71**, 6051–60 (2011).
- 663 42. M. Miyawaki, H. Yasuda, T. Tani, J. Hamamoto, D. Arai, K. Ishioka, K. Ohgino, S. Nukaga, T.
664 Hirano, I. Kawada, K. Naoki, Y. Hayashi, T. Betsuyaku, K. Soejima, Overcoming EGFR Bypass Signal-
665 Induced Acquired Resistance to ALK Tyrosine Kinase Inhibitors in ALK-Translocated Lung Cancer. *Mol*
666 *Cancer Res.* **15**, 106–114 (2016).
- 667 43. D. Aran, Z. Hu, A. J. Butte, xCell: digitally portraying the tissue cellular heterogeneity landscape.
668 *Genome Biol.* **18**, 220 (2017).
- 669 44. A. I. Khan, S. M. Kerfoot, B. Heit, L. Liu, G. Andonegui, B. Ruffell, P. Johnson, P. Kubes, Role of
670 CD44 and hyaluronan in neutrophil recruitment. *J Immunol.* **173**, 7594–601 (2004).
- 671 45. J. D. Klement, A. V. Paschall, P. S. Redd, M. L. Ibrahim, C. Lu, D. Yang, E. Celis, S. I. Abrams, K.
672 Ozato, K. Liu, An osteopontin/CD44 immune checkpoint controls CD8+ T cell activation and tumor
673 immune evasion. *J Clin Invest.* **128**, 5549–5560 (2018).
- 674 46. C. M. Porter, M. A. Havens, N. A. Clipstone, Identification of amino acid residues and protein kinases
675 involved in the regulation of NFATc subcellular localization. *J Biol Chem.* **275**, 3543–51 (2000).
- 676 47. W. Yang, S. A. Gibson, Z. Yan, H. Wei, J. Tao, B. Sha, H. Qin, E. N. Benveniste, Protein kinase 2
677 (CK2) controls CD4. *Mucosal Immunol.* **13**, 788–798 (2020).
- 678 48. S. Martínez-Høyer, A. Aranguren-Ibáñez, J. García-García, E. Serrano-Candelas, J. Vilardell, V.
679 Nunes, F. Aguado, B. Oliva, E. Itarte, M. Pérez-Riba, Protein kinase CK2-dependent phosphorylation of
680 the human Regulators of Calcineurin reveals a novel mechanism regulating the calcineurin-NFATc
681 signaling pathway. *Biochim Biophys Acta.* **1833**, 2311–21 (2013).
- 682 49. F. Skoulidis, M. E. Goldberg, D. M. Greenawalt, M. D. Hellmann, M. M. Awad, J. F. Gainor, A. B.
683 Schrock, R. J. Hartmaier, S. E. Trabucco, L. Gay, S. M. Ali, J. A. Elvin, G. Singal, J. S. Ross, D.
684 Fabrizio, P. M. Szabo, H. Chang, A. Sasson, S. Srinivasan, S. Kirov, J. Szustakowski, P. Vitazka, R.
685 Edwards, J. A. Bufill, N. Sharma, S.-H. I. Ou, N. Peled, D. R. Spigel, H. Rizvi, E. J. Aguilar, B. W.
686 Carter, J. Erasmus, D. F. Halpenny, A. J. Plodkowski, N. M. Long, M. Nishino, W. L. Denning, A.
687 Galan-Cobo, H. Hamdi, T. Hirz, P. Tong, J. Wang, J. Rodriguez-Canales, P. A. Villalobos, E. R. Parra, N.
688 Kalhor, L. M. Sholl, J. L. Sauter, A. A. Jungbluth, M. Mino-Kenudson, R. Azimi, Y. Y. Elamin, J. Zhang,

- 689 G. C. Leonardi, F. Jiang, K.-K. Wong, J. J. Lee, V. A. Papadimitrakopoulou, I. I. Wistuba, V. A. Miller,
690 G. M. Frampton, J. D. Wolchok, A. T. Shaw, P. A. Jänne, P. J. Stephens, C. M. Rudin, W. J. Geese, L. A.
691 Albacker, J. V. Heymach, *Cancer Discov.* **8**, 822–835 (2018).
- 692 50. K. Kourou, T. P. Exarchos, K. P. Exarchos, M. V. Karamouzis, D. I. Fotiadis, Machine learning
693 applications in cancer prognosis and prediction. *Comput Struct Biotechnol J.* **13**, 8–17 (2014).
- 694 51. N. H. Shah, M. Löbel, A. Weiss, J. Kuriyan, Fine-tuning of substrate preferences of the Src-family
695 kinase Lck revealed through a high-throughput specificity screen. *Elife.* **7** (2018),
696 doi:10.7554/elife.35190.
- 697 52. B. Deb, P. Sengupta, J. Sambath, P. Kumar, Bioinformatics Analysis of Global Proteomic and
698 Phosphoproteomic Data Sets Revealed Activation of NEK2 and AURKA in Cancers. *Biomolecules.* **10**
699 (2020), doi:10.3390/biom10020237.
- 700 53. R. Beekhof, C. van Alphen, A. A. Henneman, J. C. Knol, T. V. Pham, F. Rolfs, M. Labots, E.
701 Henneberry, T. Y. Le Large, R. R. de Haas, S. R. Piersma, V. Vurchio, A. Bertotti, L. Trusolino, H. M.
702 Verheul, C. R. Jimenez, INKA, an integrative data analysis pipeline for phosphoproteomic inference of
703 active kinases. *Mol Syst Biol.* **15**, e8250 (2019).
- 704 54. K. P. Exarchos, Y. Goletsis, D. I. Fotiadis, Multiparametric decision support system for the prediction
705 of oral cancer reoccurrence. *IEEE Trans Inf Technol Biomed.* **16**, 1127–34 (2011).
- 706 55. A. Antoniou, A. Cunningham, J. Peto, D. G Evans, F Lalloo, S. A. Narod, H. A. Risch, J.E. Eyfjord,
707 J.L. Hopper, M. C. Southey, H. Olsson, O. Johannsson, A. Borg, B. Passini, P. Radice, S Manoukian, D.
708 M Eccles, N. Tang, E. Olah, H. Anton-Culver, E. Warner, J. Lubinski, J. Gronwald, B. Gorski, L.
709 Tryggvadottir, K. Syrjakoski, O.P. Kallioniemi, H. Eerola, H. Nevanlinna, P.D.P Pharoah, D.F. Easton,
710 The BOADICEA model of genetic susceptibility to breast and ovarian cancers: Updates and extensions.
711 *British Journal of Cancer.* **8**, 1457-1466 (2008).
- 712 55. D. Schwartz, S. P. Gygi, An iterative statistical approach to the identification of protein
713 phosphorylation motifs from large-scale data sets. *Nat Biotechnol.* **23**, 1391–8 (2005).
- 714 56. A. Tareen, J. B. Kinney, Logomaker: Beautiful sequence logos in python. *Cold Spring Harbor*
715 *Laboratory* (2019), doi:10.1101/635029.
- 716 57. M. Ashburner, C. A. Ball, J. A. Blake, D. Botstein, H. Butler, J. M. Cherry, A. P. Davis, K. Dolinski,
717 S. S. Dwight, J. T. Eppig, M. A. Harris, D. P. Hill, L. Issel-Tarver, A. Kasarskis, S. Lewis, J. C. Matese,
718 J. E. Richardson, M. Ringwald, G. M. Rubin, G. Sherlock, Gene ontology: tool for the unification of
719 biology. The Gene Ontology Consortium. *Nat Genet.* **25**, 25–9 (2000).
- 720 58. The Gene Ontology resource: enriching a GOLD mine. *Nucleic Acids Res.* **49**, D325–D334 (2021).

721

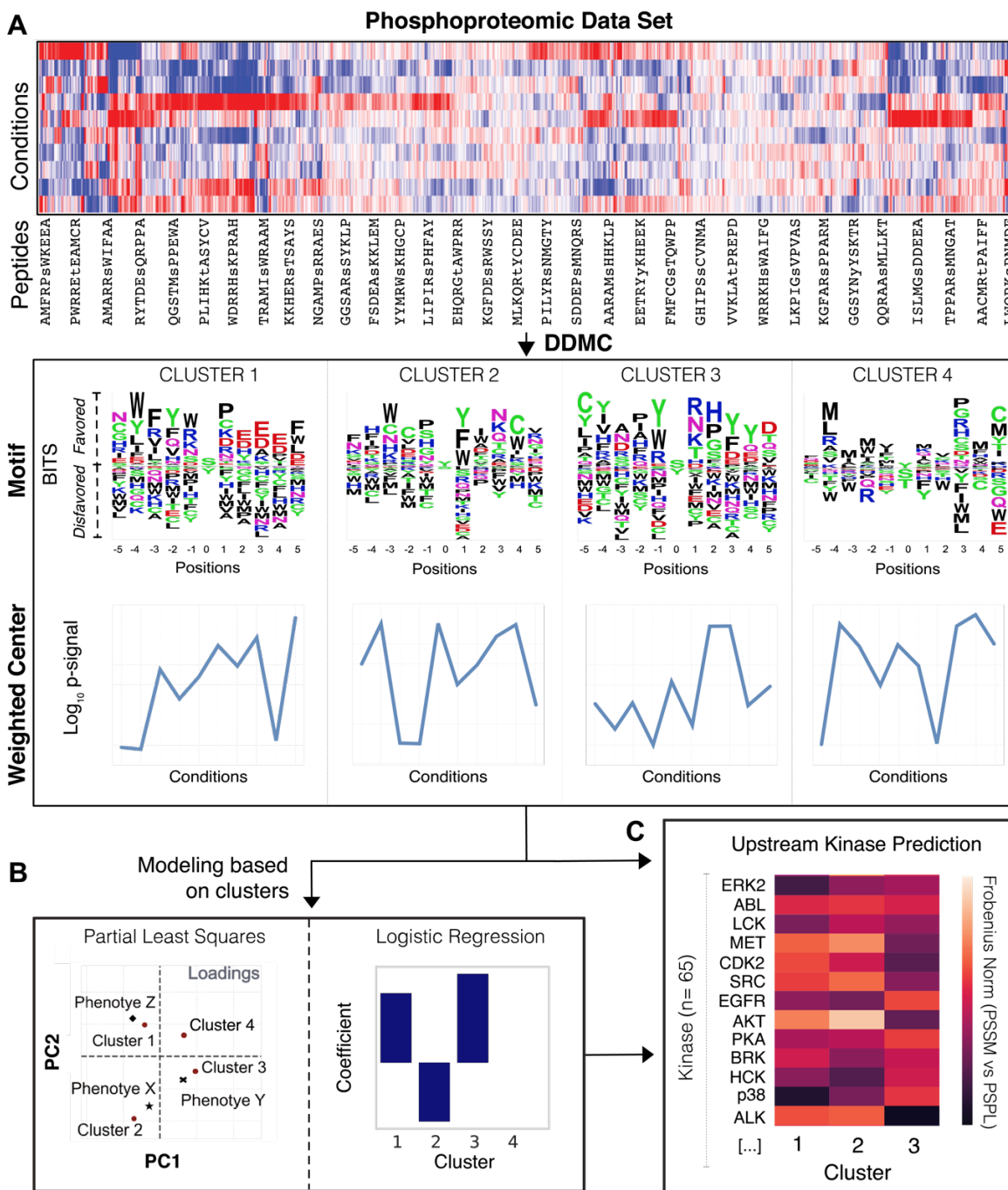
722 **Funding:** This work was supported by NIH U01-CA215709 to A.S.M. and in part by the UCLA Jonsson
723 Comprehensive Cancer Center (JCCC) grant NIH P30-CA016042.

724 **Author contributions:** A.S.M. conceived the project. Both authors performed the analysis. Both authors
725 wrote the manuscript.

726 **Competing interests:** Authors declare that they have no competing interests.

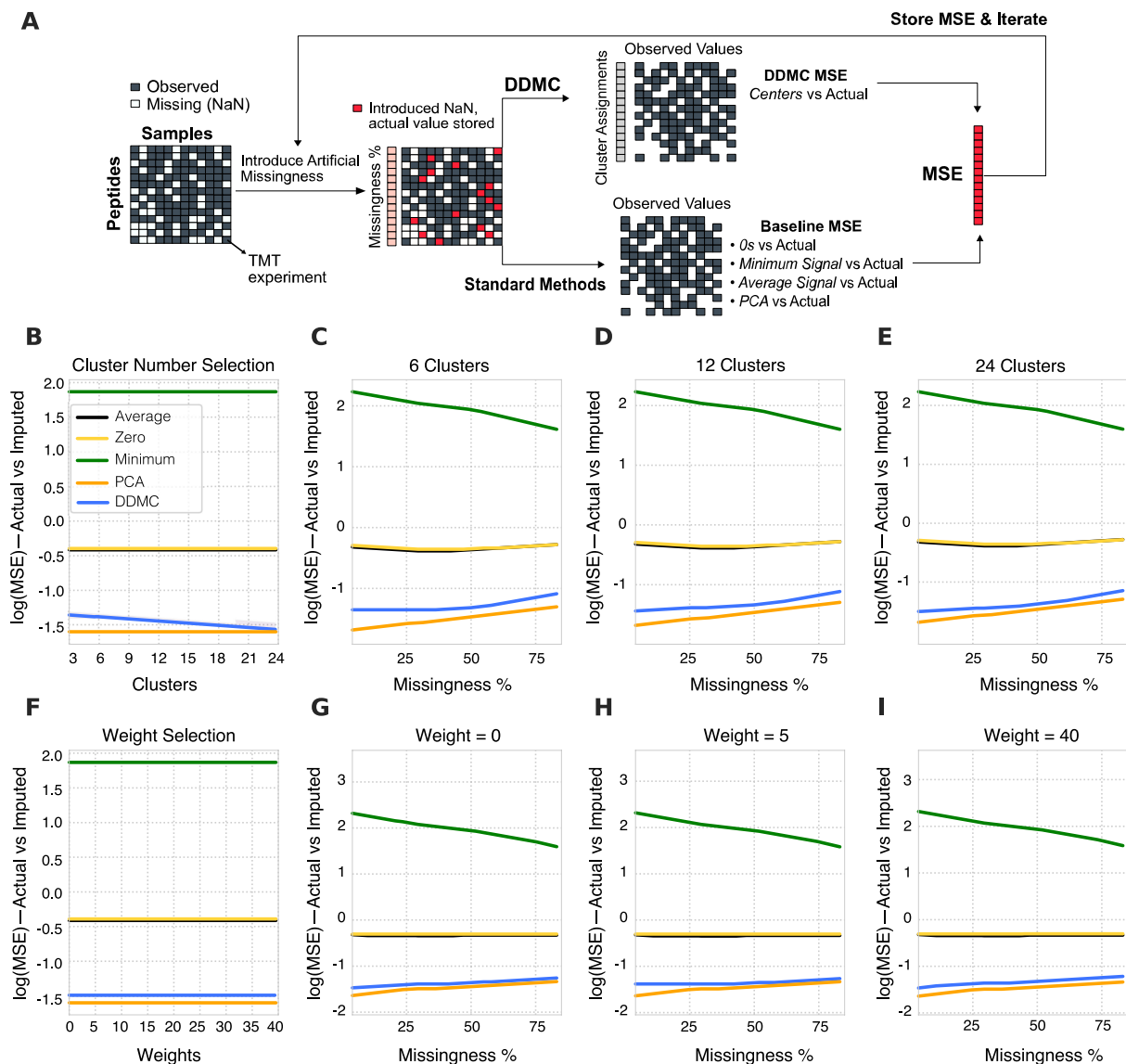
727

728 **Figures**



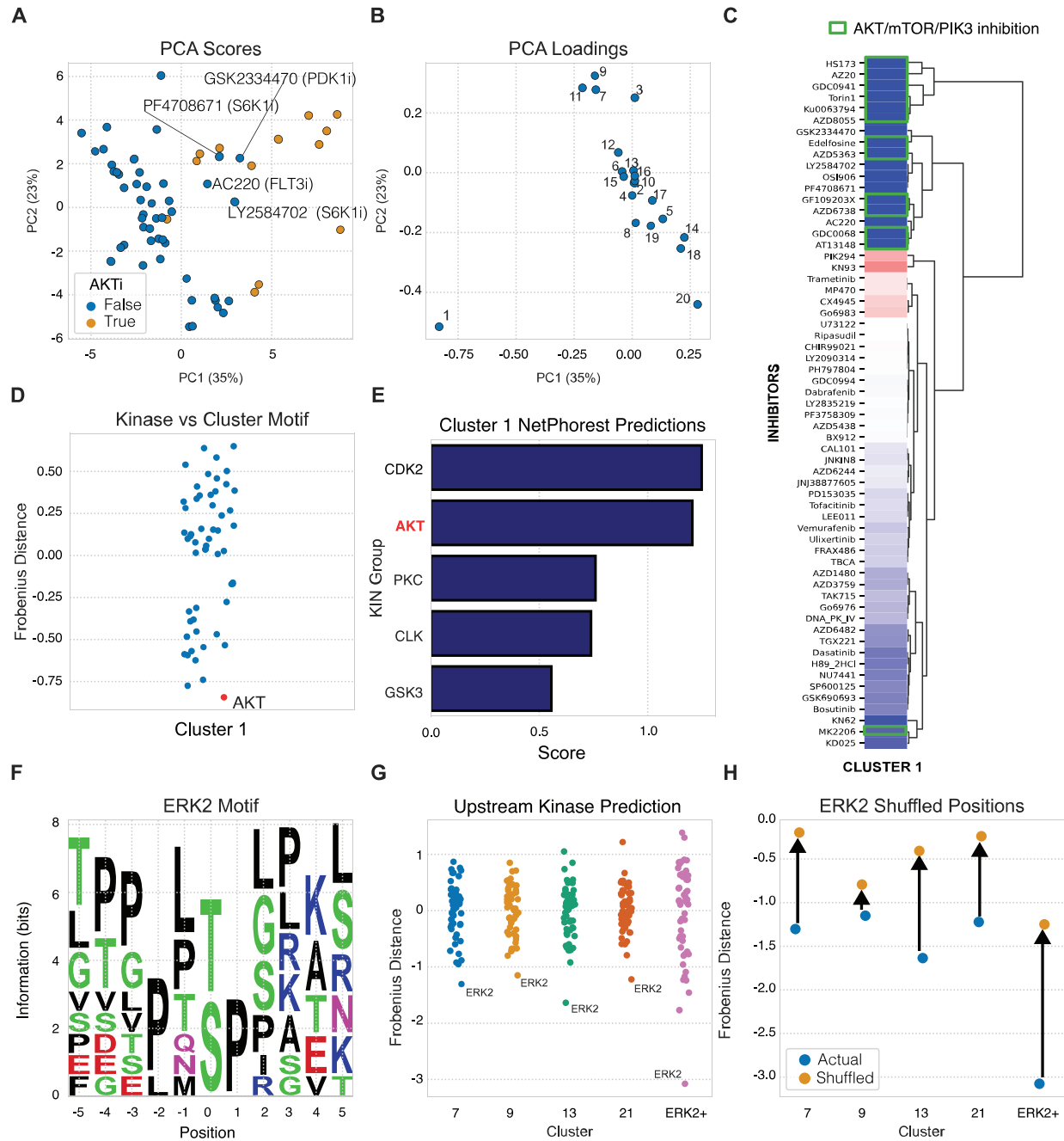
729

730 **Figure 1: Schematic of the DDMC approach to cluster global signaling data and infer upstream**
 731 **kinases driving phenotypes.** A) DDMC is run to cluster an input phosphoproteomic data set to generate
 732 4 clusters of peptides that show similar sequence motifs and phosphorylation behavior. B) Predictive
 733 modeling using clusters allows one to establish associations between specific clusters and features of
 734 interest. C) Putative upstream kinases regulating meaningful clusters can be predicted by computing the
 735 distance between a cluster motif and PSPL PSSM. PSSM; Position-specific scoring matrix, PSPL;
 736 Position scanning peptide library.



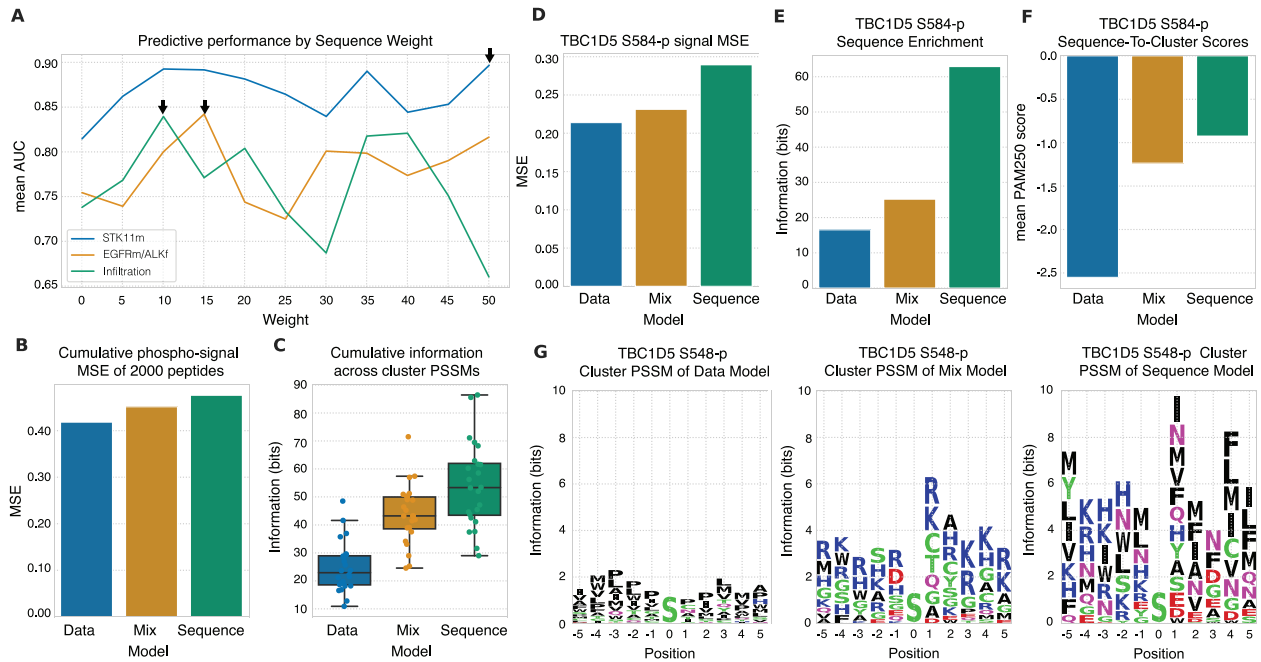
737

738 **Figure 2: Benchmarking the robustness of motif clustering to missing measurements.** A) A
 739 schematic of the process for quantifying robustness to missing values. Any peptides containing less than 7
 740 TMT experiments were discarded. For the remaining 15904 peptides, an entire random TMT experiment
 741 was removed per peptide and these values were stored for later comparison. Next, these artificial missing
 742 values were imputed using either a baseline strategy (peptide mean/minimum signal, constant zero, or
 743 matrix completion by PCA) or the corresponding cluster center. Once a mean squared error was computed
 744 for each peptide, the second iteration repeats this process by removing a second TMT experiment. A total
 745 of 5 random TMT experiments per peptide were imputed by clustering using a different number of
 746 clusters (B-E) or different weights (E-I).



747

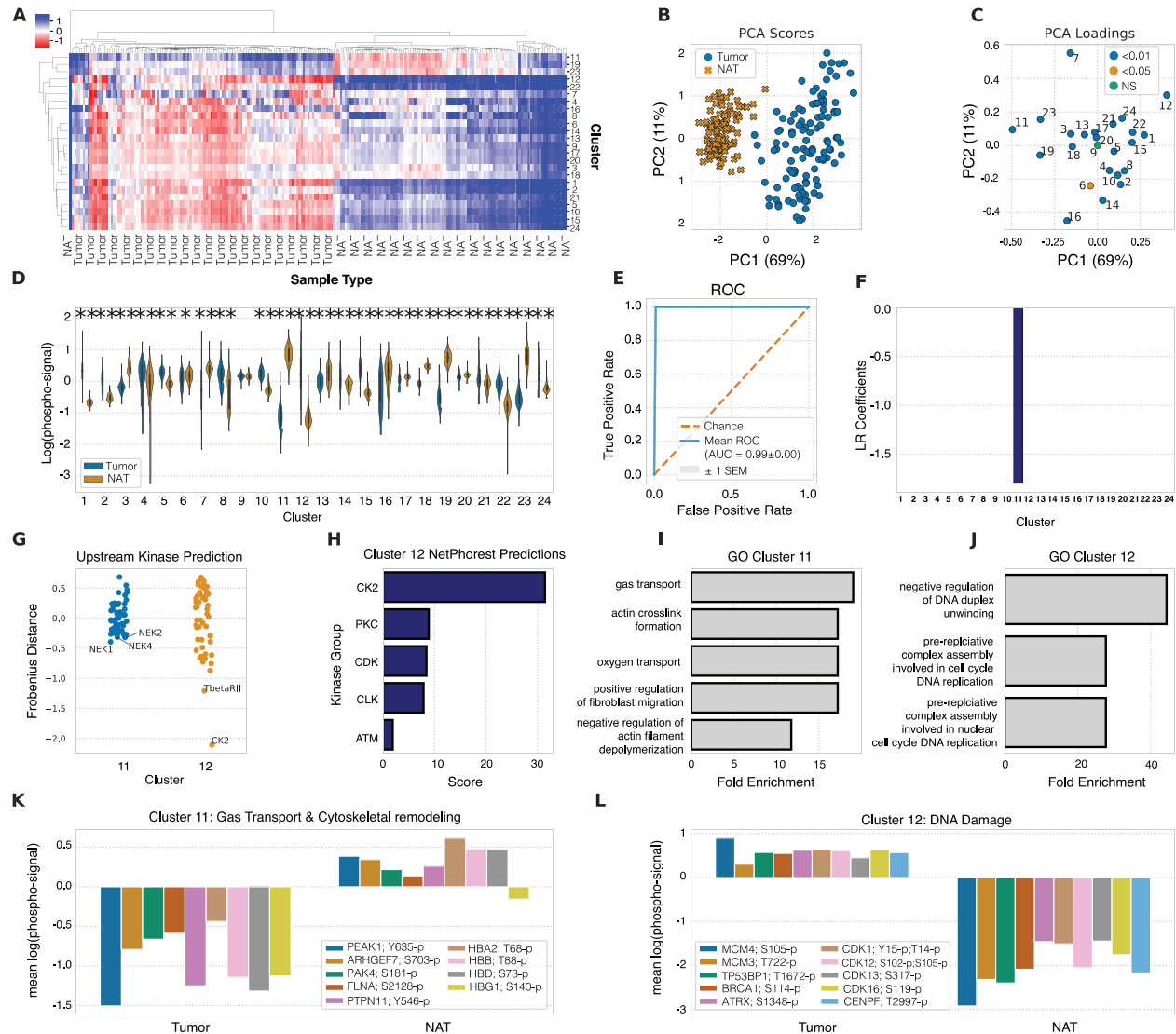
748 **Figure 3: Validation of upstream kinase predictions.** (A-B) PCA analysis of the DDMC
 749 phosphoproteome clusters of MCF7 cells subjected to a drug screen (20). C) Heatmap showing the effect
 750 of inhibitors on the phosphorylation signal of cluster 1. D) DDMC upstream kinase prediction of cluster
 751 1. E) NetPhorest upstream kinase prediction of cluster 1. (F) Resulting PSSM generated using ERK2
 752 substrates reported by Carlson et al (21). (G) Upstream kinase predictions of CPTAC clusters 7, 9, 13,
 753 and 21 in addition to the ERK2 motif shown in (F). H) Upstream kinase predictions of the same PSSMs
 754 after randomly shuffling the motif positions.



755

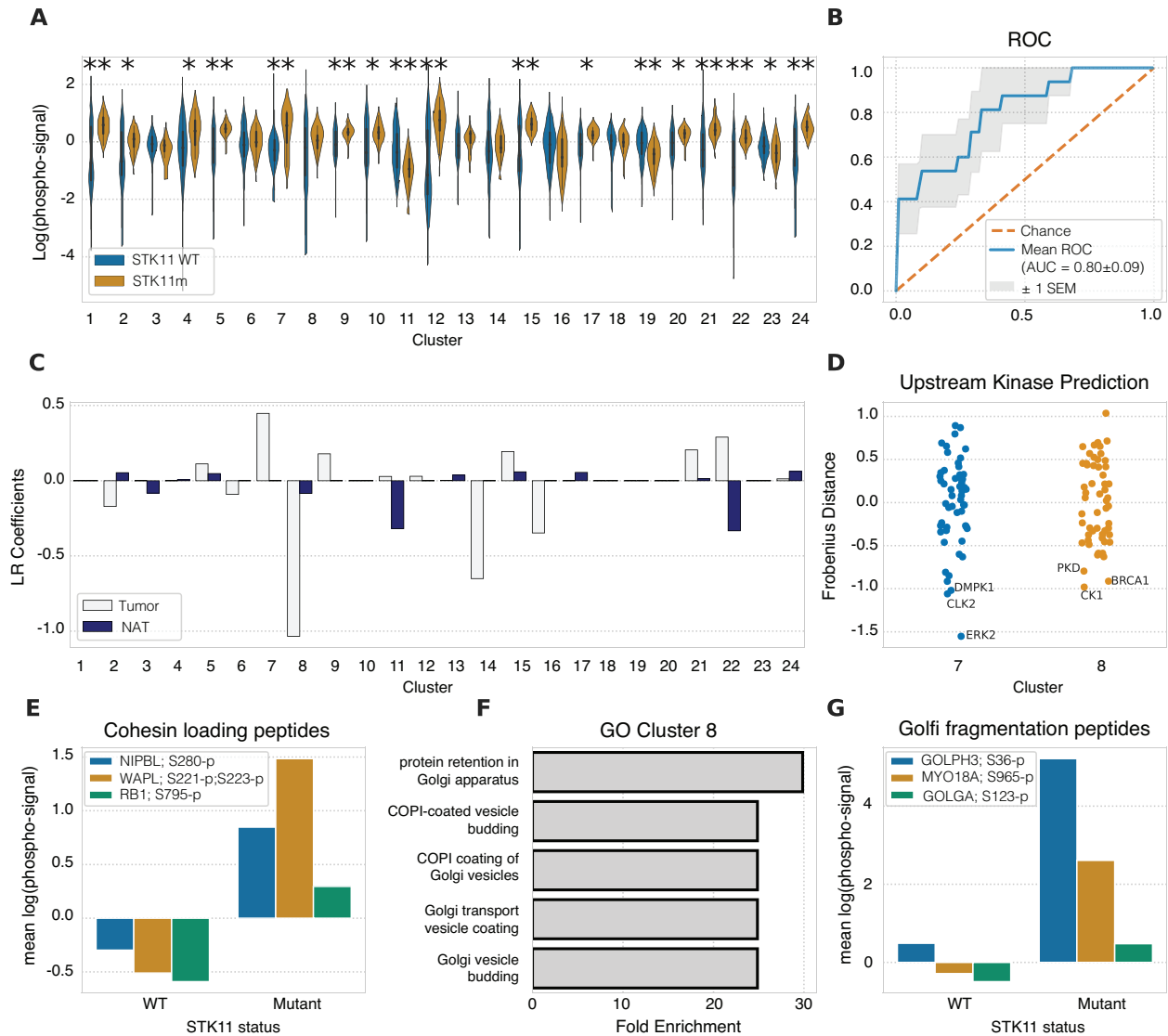
756 **Figure 4: Sequence information enhances model prediction and provides more robust clustering.** A)

757 Performance of a regression model predicting the mutational status of STK11 (blue) EGFR and/or ALK
 758 (yellow) and tumor infiltration (green) in LUAD patients using either only phosphorylation data
 759 (weight=0), mainly sequence information (50), or both ($0 < w < 50$). B) MSE between the
 760 phosphorylation signal of 2000 randomly selected peptides and the center of its assigned clusters using a
 761 weight of 0 (data), 20 (mix), or 50 (sequence). C) Cumulative PSSM enrichment across positions
 762 comparing the data, mix, and sequence clustering strategies. (D-H) TBC1D5 peptide p-signal MSE (D),
 763 cumulative PSSM enrichment (E), and PSSM logo plots (F-H).



764

765 **Figure 5: Conserved tumor differences compared to normal adjacent tissue.** A) Hierarchical
 766 clustering of DDMC cluster centers. B–C) Principal components analysis scores (B) and loadings (C) of
 767 the samples and phosphopeptide clusters, respectively. D) Phosphorylation signal of tumor and NAT
 768 samples per cluster and statistical significance according to a Mann Whitney rank test (* = p-value < 0.05
 769 and ** = p-value < 0.001). E) Receiver operating characteristic curve (ROC) of a regularized logistic
 770 regression model. F) Logistic regression weights per cluster. G) Upstream kinase predictions of clusters
 771 11 and 12. (H) NetPhorest kinase predictions of cluster 12. (I–J) Gene ontology analysis and (K–L)
 772 representative peptides of enriched biological processes of clusters 11 and 12.



773

774 **Figure 6: Phosphoproteomic aberrations associated with STK11 mutational status. A)**

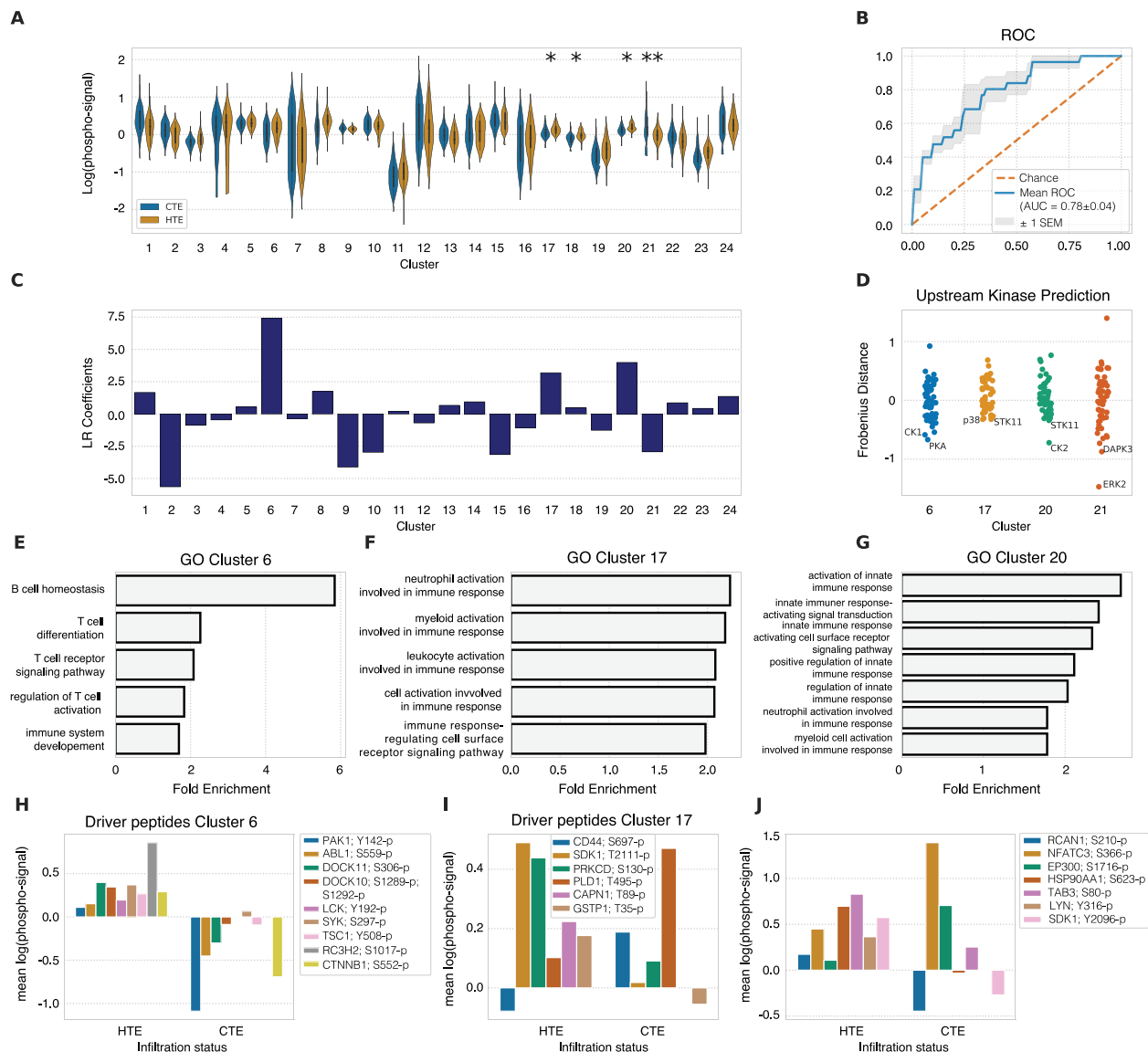
775 Phosphorylation signal of STK11 WT and mutant samples per cluster and statistical significance

776 according to a Mann-Whitney rank test (* = p-value < 0.05 and ** = p-value < 0.001). B)

777 ROC of a logistic regression model predicting the STK11 mutational status and (C)

778 its corresponding weights per sample type. (D) Putative upstream kinases of clusters 7, and 8. (E) Representative cohesin loading

779 peptides in cluster 7. (F-G) GO analysis and representative Golgi fragmentation peptides of cluster 8.



780

781 **Figure 7: Phosphoproteomic signatures driving tumor immune infiltration.** (A) Phosphorylation
 782 abundance of CTE and HTE samples per cluster and statistical significance according to a Mann-Whitney
 783 rank test (* = p-value < 0.05 and ** = p-value < 0.001). (B–C) ROC and coefficients of a logistic
 784 regression model predicting infiltration status—cold-tumor enriched (CTE) versus hot-tumor enriched
 785 (HTE). (D) Putative upstream kinases of clusters 7, 17, 20, and 21. (E–G) GO enrichment analysis of
 786 select clusters. (H–J) Selected peptides driving the GO biological processes in HTE versus CTE samples.