

MPAC: a computational framework for inferring cancer pathway activities from multi-omic data

Peng Liu^{1,2}, David Page^{1,2,3†}, Paul Ahlquist^{4,5,6}, Irene M. Ong^{1,2,7,8*}, Anthony Gitter^{1,3,4*}

1. Department of Biostatistics and Medical Informatics, University of Wisconsin–Madison, Madison, Wisconsin, United States of America
2. Carbone Cancer Center, University of Wisconsin–Madison, Madison, Wisconsin, United States of America
3. Department of Computer Sciences, University of Wisconsin–Madison, Madison, Wisconsin, United States of America
4. John and Jeanne Rowe Center for Research in Virology, Morgridge Institute for Research, Madison, Wisconsin, United States of America
5. McArdle Laboratory for Cancer Research, University of Wisconsin–Madison, Madison, Wisconsin, United States of America
6. Institute for Molecular Virology, University of Wisconsin–Madison, Madison, Wisconsin, United States of America
7. Department of Obstetrics and Gynecology, University of Wisconsin–Madison, Madison, Wisconsin, United States of America
8. Center for Human Genomics and Precision Medicine, University of Wisconsin–Madison, Madison, Wisconsin, United States of America

† Present address: Department of Biostatistics & Bioinformatics, Duke University, Durham, North Carolina, United States of America

*Equal contribution

Corresponding authors: irene.ong@wisc.edu, gitter@biostat.wisc.edu

Abstract

Fully capturing cellular state requires examining genomic, epigenomic, transcriptomic, proteomic, and other assays for a biological sample and comprehensive computational modeling to reason with the complex and sometimes conflicting measurements. Modeling these so-called multi-omic data is especially beneficial in disease analysis, where observations across omic data types may reveal unexpected patient groupings and inform clinical outcomes and treatments. We present Multi-omic Pathway Analysis of Cancer (MPAC), a computational framework that interprets multi-omic data through prior knowledge from biological pathways. MPAC uses network relationships encoded in pathways using a factor graph to infer consensus activity levels for proteins and associated pathway entities from multi-omic data, runs permutation testing to eliminate spurious activity predictions, and groups biological samples by pathway activities to prioritize proteins with potential clinical relevance. Using DNA copy number alteration and RNA-seq data from head and neck squamous cell carcinoma patients from The Cancer Genome Atlas as an example, we demonstrate that MPAC predicts a patient subgroup related to immune responses not identified by analysis with either input omic data type alone. Key proteins identified via this subgroup have pathway activities related to clinical outcome as well as immune cell compositions. Our MPAC R package, available at <https://bioconductor.org/packages/MPAC>, enables similar multi-omic analyses on new datasets.

Introduction

Cancer is a complex set of diseases with a great diversity of genomic aberrations and altered signaling pathways (Hanahan, 2022). The Cancer Genome Atlas (TCGA) generated data spanning copy number alteration (CNA), DNA mutation, DNA methylation, mRNA expression, microRNA expression, and protein expression for thousands of tumor samples, leading to many insights into the cancers that were profiled (Hoadley *et al.*, 2018). In addition, this extensive multi-omic data provides clues to tumor regulation, which have led to the development of many computational methods to integrate multi-omic data to obtain comprehensive views on cancer (Picard *et al.*, 2021; Maghsoudi *et al.*, 2022; G. L. Stein-O'Brien *et al.*, 2018).

In particular, biological pathway-based approaches have been demonstrated as a powerful way to integrate multi-omic data (reviewed in Maghsoudi *et al.* 2022). Altered expression or function of different genes in the same pathway can have similar impacts on overall pathway activity. Similarly, diverse alterations of expression or function of the same gene or its protein product—e.g. through DNA mutations, CNAs, or changes in epigenetic modifications, transcript expression, or protein translation, stability, or post-translational modifications—can also suppress, stimulate, or otherwise modulate a particular pathway. These properties allow modeling based on multi-omic inputs to infer pathway activity to more accurately reflect underlying biology than modeling based on a narrow, incomplete view from a single genomic data type. Accordingly, whereas a single data type rarely contains the full explanation for oncogenesis, pathway-based approaches are a particularly advantageous way to understand cancer mechanisms.

Several notable pathway-based methods have demonstrated the benefits of multi-omic data integration for cancer interpretation. For example, Multi-omics Master-Regulator Analysis

identified 112 distinct tumor subtypes and 24 conserved master regulator blocks across 20 TCGA cohorts (Paull *et al.*, 2021). OncoSig delineated tumor-specific molecular interaction signaling maps for the full repertoire of 715 proteins in the COSMIC Cancer Gene Census (Broyde *et al.*, 2021). COSMOS combined signaling, metabolic, and gene regulatory networks to capture crosstalks within and between multi-omics data (Dugourd *et al.*, 2021). Pathway Recognition Algorithm using Data Integration on Genomic Models (PARADIGM) integrates multi-omic data via a factor graph to infer activities of all the proteins in a pathway network (Sedgewick *et al.*, 2013; Vaske *et al.*, 2010). Initially, PARADIGM was successfully applied to breast cancer and glioblastoma patients using CNA and gene expression microarray data to find clinically relevant groups and associated pathways. It was further applied to reveal multiple low-frequency, but high-impact mutations in glioblastoma, ovarian, and lung cancers (Ng *et al.*, 2012) and was incorporated into the standard TCGA analysis pipeline (Hoadley *et al.*, 2018, 2014).

Despite such successes, there are still opportunities to further improve multi-omic modeling. Multi-omics Master-Regulator Analysis and OncoSig focused on direct interactions around master regulators for transcription. The indirect effects of proteins further downstream of the master regulators in biological pathways were not considered. PARADIGM's application across many cancer types focused on grouping patients by their inferred pathway levels or enriched pathways (Hoadley *et al.*, 2018; Berger *et al.*, 2018). But in-depth analysis on the molecular basis of patient grouping, careful interpretation of associated inferred pathway levels (which are abstract quantities that represent neither protein abundance nor any particular post-translational modification and cannot be experimentally measured), and an end-to-end computational process are lacking. Other existing patient stratification methods by multi-omic data either do not utilize biological pathway information (Duan *et al.*, 2021) or rely on

unrealistically small pathways (Zhao *et al.*, 2021). As a result, it is hard to identify key proteins from a broad perspective with meaningful biological interpretation and clinical implication.

Here, we develop a computational framework, named Multi-omic Pathway Analysis of Cancer (MPAC), to integrate multi-omic data for understanding cancer mechanisms. It is built upon the PARADIGM method with notable improvements including providing enhanced insights to the molecular basis and clinical implications of pathway-based patient groups as well as streamlining the whole computational process. In this work, we apply MPAC to Head and Neck Squamous Cell Carcinoma (HNSCC), which accounts for ~500,000 deaths per year worldwide (Mody *et al.*, 2021). First, we describe how MPAC improves upon PARADIGM. Next, we apply MPAC to TCGA HNSCC data and group patients by their significantly altered pathways. Among other results, MPAC predicts a patient group that is enriched with immune response pathways, and this group cannot be predicted from the individual omic data types alone. Investigating this group identifies seven proteins that have activated pathway levels correlated with better overall survival. These findings are validated by a holdout set of TCGA HNSCC samples. Lastly, we present an interactive R Shiny app that lets users explore all the results generated from this work.

Results

An overview of MPAC and its improvements upon PARADIGM

We developed MPAC to integrate multi-omic data to identify key pathways and proteins with biological and clinical implications, and to predict new patient groups associated with distinct pathway alterations. MPAC's workflow contains eight steps (Figure 1; see also Methods): (Step 1) From CNA and RNA-seq data, determine genes' CNA and RNA ternary states (i.e., repressed, normal, or activated). CNA and RNA-seq data are selected as the input

multi-omic data because PARADIGM had shown success with them (Vaske et al. 2010; Sedgewick et al. 2013; Hoadley et al. 2018); (Step 2) Use CNA and RNA state together with a comprehensive biological pathway network from TCGA (Hoadley *et al.*, 2018) to calculate pathway levels with PARADIGM's factor graph model. The TCGA pathway network characterizes interactions at both transcriptional and post-transcriptional levels. PARADIGM's inferred pathway levels are calculated not only for proteins, but also for several other types of pathway entities, such as protein complexes and gene families; (Step 3) Permute input CNA and RNA data for filtering inferred pathway level in the next step. CNA and RNA states are permuted randomly between genes in each patient. Inferred pathway levels for each pathway entity are calculated with PARADIGM from 100 sets of permuted data to build a background distribution representing inferred pathway levels observed by chance; (Step 4) Inferred pathway levels computed from real data are compared with those from permuted data to filter out inferred pathway level changes that could be observed by chance. Because both the real and permuted pathway levels are different between patients, this filtering step creates a patient-specific set of inferred pathway levels representing each patient's unique pathway alteration profiles; (Step 5) From the remaining pathway networks, retain the connected component that contains the largest number of connected pathway entities. This focuses on the main subset of entities that are connected in the pathway network and presumably lead to similar functional alterations; (Step 6) Build patient pathway profiles and predict patient groups. Each patient's pathway profile contains a selected set of 2,085 cancer- and general biology-relevant Gene Ontology (GO) terms (Gene Ontology Consortium *et al.*, 2023; Kuenzi *et al.*, 2020). Each GO term is characterized by an over-representation test p-value on entities selected from the previous step. P-values for all the GO terms are adjusted for multiple hypothesis testing and log-transformed to group patients (Lun *et al.*, 2016); (Step 7) Identify key proteins that all have activated or

repressed inferred pathway levels between patients from the same group; (Step 8) Compare data on key proteins with patients' clinical data to evaluate potential impact of protein inferred pathway levels on patients' clinical outcomes.

MPAC makes multiple improvements over PARADIGM (Vaske *et al.*, 2010), which MPAC runs as a subroutine in Steps 2 and 3. PARADIGM simply divided all the genes into three states with an equal number of entries. In contrast, in Step 1, MPAC defines a gene's RNA state as normal, activated or repressed for each patient by testing the level of the relevant RNA in that tumor sample for significant increase or decrease (two standard deviations from the mean) of that RNA's expression distribution in normal tissue samples. In Steps 3 and 4, as noted above, MPAC filters pathway entities for significant inferred pathway level differences from randomly permuted input. Although some PARADIGM applications also used permutations, permutations were not implemented as part of the software, nor were their results used for downstream analysis (Vaske *et al.*, 2010; Sedgewick *et al.*, 2013). In Step 5, MPAC focuses on the largest patient-specific pathway network subset. This improvement removed from consideration entities in tiny pathways, which were assumed to have less impact on patient pathway alterations and may contribute more noise than signal when predicting patient groups. In Steps 6–8, MPAC provides comprehensive downstream analysis functions to define patient pathway alterations, predict patient groups, and identify key proteins with potential clinical implications. MPAC is available as an R package on Bioconductor (<https://bioconductor.org/packages/MPAC>) to streamline the whole process from preparing the omic input data to identifying key proteins for a patient group.

MPAC predicted an immune response HNSCC group not found by CNA or RNA-seq data alone

We applied MPAC to TCGA HNSCC patients to predict patient groups by their pathway alterations. We selected the 492 patients that had CNA, RNA-seq, and overall survival data available. Of these 492 tumors, 89 carried human papillomavirus DNA (HPV+) and 403 did not (HPV-), a distinction linked to major differences in HNSCC tumor biology and clinical treatment response (Powell *et al.*, 2021). We further randomly divided patients into exploratory sets (71 HPV+ and 322 HPV-) and validation sets (18 HPV+ and 81 HPV-). Our goal was to first tune MPAC and identify pathway-based patient groups in the exploratory set and then test our discoveries in the validation set. We applied MPAC separately to the HPV+ and HPV- exploratory set patients and identified five groups from each based on the patient pathway profiles (Figure 2 and Supplementary Figure 1A). For HPV+ patients (Figure 2), four of the five groups had distinct pathway features. Group I patients had alterations mainly in immune response pathways, groups II and IV in cell cycle pathways, and group V in morphogenesis pathways. Group III had pathway alterations in some patients but did not show an obvious biologically meaningful consensus profile. For HPV- patients (Supplementary Figure 1A), three of the five patient groups had distinct pathway alterations: groups I and IV in cell cycle pathways and group III in immune response pathways. Groups II and V did not show obvious consensus pathway features. The distinct pathway features for many of the patient groups suggested that MPAC is capable of building biomedically relevant patient pathway profiles and predicting patient groups.

MPAC predicted the above patient groups with distinct pathway profiles by integrating multi-omic data. We found that the same groups and pathways cannot be found by examining individual omic data types alone. Starting from either CNA or RNA-seq data, we conducted two

tests: one performing GO enrichment on each single omic data type and then grouping enriched GO terms, like the MPAC workflow, and the other by grouping patients via their single omic data first and then finding commonly enriched GO terms within each group. In the first test, HPV+ patients' CNA data had very few GO terms enriched, and even these were only enriched in a small number of patients (Supplementary Figure 2A). As a result, no group prediction could be made. RNA-seq data was more informative than CNA data, and four patient groups could be predicted (Supplementary Figure 2B). Groups III and IV were related to cell cycle and morphogenesis pathways, respectively, both of which had also been observed in the MPAC results. The immune response patient group predicted by MPAC, however, was not observed from CNA or RNA-seq data, indicating a unique insight from MPAC. For HPV- patients, CNA data did not lead to any patient groups due to insufficient significantly enriched GO terms (Supplementary Figure 1B). RNA-seq data led to six groups. Two of them were related to cell cycle and immune response (Supplementary Figure 1C), which were also observed in MPAC results.

To demonstrate the robustness of this result, we performed another test by grouping patients first and followed by GO enrichment. We applied K-means clustering to the RNA-seq data and divided HPV+ patients into two to six groups. The cluster membership remained stable under different numbers of groups (Supplementary Figure 3A). Therefore, we used six groups (Supplementary Figure 3B) for GO enrichment analysis. Groups I, II, IV, and V did not have any top GO terms related to immune response (Supplementary Figure 3 C–D and F–G). Group III was predominantly enriched with cell cycle-related GO terms with only one GO term (lymphocyte activation) related to immune response (Supplementary Figure 3E). Moreover, this single term was less consistently enriched than the >20 immune response GO terms from MPAC (Figure 2A). Group VI contained only one patient and was skipped. We performed the

same analysis on CNA data. Stable grouping membership was observed (Supplementary Figure 4A), and the six-groups result was used for GO enrichment analysis (Supplementary Figure 4B). No GO term was significantly over-represented in groups I or II, while all the other four groups contained only one patient each and were skipped. In summary, by jointly modeling both CNA and RNA-seq data, MPAC identified a large and unique HPV+ patient group related to immune response that could not be recovered from either individual omic data type alone.

Proteins from the HPV+ immune response group correlated with patient overall survival

Given that MPAC discovered an immune response patient group that could not be found by CNA or RNA-seq data alone, we were interested in pathway submodules and key proteins shared by the eleven patients in this group. We defined a pathway submodule as a pathway subset containing ≥ 5 entities, at least one of which was a protein with input omic data. We required that all submodule entities must have activated or repressed inferred pathway levels in the eleven patients. Four such submodules were identified (Figure 3A). They contained five to twelve entities and collectively eight proteins (Figure 3A, red ovals). Seven of these proteins, CD28, CD86, TYK2, IL12RB1, LCP2, FASLG, and CD247, had activated inferred pathway levels in all eleven group I patients (Figure 3B), suggesting a consensus functional role across patients within this immune response group. Interestingly, prior studies collectively showed that gene expression levels of the seven proteins except for LCP2 are correlated with immune infiltration in HNSCC (Zhu *et al.*, 2022; Wang *et al.*, 2022; He *et al.*, 2022; Chi *et al.*, 2022; Chen *et al.*, 2022; de Vos *et al.*, 2020). Below we show a similar correlation for the patient group analyzed here.

To understand what factors determine pathway levels of these seven proteins, we developed an approach for pathway state visualization in MPAC. We transformed the

continuous-valued inferred pathway levels to discretized pathway states with the values activated, normal, or repressed. The resulting plots presented a protein's direct pathway network interaction partners and all associated pathway state information for that protein and its partners, under the reasoning that a determinant of a functionally-implicated protein's pathway state would have correlated states across all patients (Supplementary Figures 5-7). For example, in all 11 patients, CD86's pathway states agreed with its RNA states as well as with six of its seven downstream interacting complexes (Supplementary Figure 5B), indicating their parallel roles in determining CD86's pathway states. In contrast, CD86's two downstream gene families and one downstream complex had states that disagreed with CD86 in one or three patients, respectively (Supplementary Figure 5B), suggesting a less influential role. CD86's CNA states, to the other extreme, did not agree with CD86's pathway state in any patient (Supplementary Figure 5B). In one patient, TCGA-CR-6487, CD86's CNA state is repressed and its RNA state is activated (Figure 3C). If studying CD86 from individual genomic datasets without any pathway information, it would be hard to determine CD86's functional protein state, illustrating the advantages of our pathway-based approach. Another feature of MPAC's visualization function was showing patient-to-patient variations on pathway determinants. FASLG, for instance, had many upstream and downstream neighbors, of which only two upstream and one downstream complexes had pathway states correlated with FASLG, while all the other neighbors had various states across the eleven patients (Supplementary Figures 7B and 8). Such diverse states of FASLG's neighbors likely reflected subtle cancer mechanism differences within this patient group and MPAC can highlight these differences.

To examine potential clinical implications of these seven proteins, we evaluated their correlation with the patients' overall survival. For proteins from the same submodule, we used their inferred pathway levels to divide the 71 HPV+ exploratory patients into two groups: those

with all relevant proteins in activated pathway levels and those that were not. The overall survival distributions of patients from the two groups were compared and evaluated by a log-rank test. For every submodule, although the improvement was not always statistically significant, the set of patients with proteins with activated pathway levels always had a better survival distribution than the set that did not (Supplementary Figure 9A). In particular, when patients had both CD247 and FASLG with activated pathway levels, their overall survival was significantly better than those that did not (log-rank $P=0.00098$). Moreover, dividing the same set of patients by the activation of all seven proteins (Figure 3D) or individual proteins (Supplementary Figure 9B) also produced overall survival advantages in all cases but with log-rank P ranging from 0.0033 to 0.17. Similar analysis revealed the same trend using the measure of progression-free survival (Supplementary Figure 10). The good correlation with patient overall survival and reduced tumor progression indicated potential clinical implications of these seven proteins individually and collectively.

Since the seven implicated proteins were identified from the immune response patient group, we explored the relationship between these proteins and immune response. We used a bulk RNA-seq deconvolution method, CIBERSORT in absolute mode (Newman *et al.*, 2015), to estimate immune cell composition for the 71 HPV+ exploratory patients and correlate them with inferred pathway levels of the seven proteins. CIBERSORT-inferred cell composition was comparable across cell types within the same patient as well as across patients for the same cell type. As in the survival analysis shown in Figure 3D, patients were stratified by whether or not they had all seven proteins with activated pathway levels. For patients in this 'activated group', the tumor sample always had substantially higher compositions of T follicular helper cells, CD8+ T cells, regulatory T cells, and M1 and M2 macrophages (Figure 3E; Supplementary Figure 11A). Thus, similar to the prior results cited above, this correlation

indicates that patients with the seven proteins with activated pathway levels usually had higher levels of immune cell infiltration and further suggests that inferred pathway levels of the seven proteins can serve as indicators for immune infiltration.

Independent validation set confirmed MPAC's immune response group and key proteins

We used the independent validation set of eighteen HPV+ patients that was held out during MPAC model development and exploratory set analysis to further assess the generality of the seven key proteins identified from the immune response patient group. Thus, we repeated the same MPAC analysis on this validation set, splitting the eighteen validation set patients into two groups. The six patients in the resulting validation group II again had many significantly enriched GO terms related to immune response (Figure 4A). None of the patients with the originally implicated submodule proteins with activated pathway levels died in the interval of record, a notably better overall survival record than those with submodule proteins in normal or repressed pathway levels (Supplementary Figure 12A). Similarly, the overall survival rate of the three patients with all seven proteins in activated pathway levels was always better than the other fifteen patients in the validation set (Figure 4B). The same trend was observed when stratifying patients by individual proteins (Supplementary Figure 12B). The lack of statistically significant differences between two patient groups was due to the small number of patients (Figure 4B; Supplementary Figure 12).

We further examined if the activated pathway levels of the seven proteins also correlated with immune cell infiltration, using the same analysis as for the exploratory set. Validation set patients with the seven proteins in activated pathway levels often had higher fractions of T follicular helper cells, CD8+ T cells, regulatory T cells, M1 and M2 macrophages (Figure 4C), the same as we observed in the exploratory set, although the difference was not significant

most likely due to the small sample size (Supplementary Figure 11B). Altogether, the independent validation set supported MPAC's predictions in the exploratory set and greatly reduced the possibility of bias from using the exploratory set alone.

An interactive MPAC Shiny app supports visualization of results and new analyses

We built an R Shiny app (https://github.com/pliu55/MPAC_Shiny) to display all the results generated from this work and support new analyses of the data. It shows enrichment results from 2,805 pathways, inferred pathway levels of 19,477 pathway entities, CNA and RNA states of 6,251 pathway proteins, and overall survival and immune cell compositions of 492 HNSCC patient samples. Moreover, it illustrates a protein's pathway membership and network neighbors. On the landing page's sidebar, users can choose one of the four TCGA-HNSCC datasets: HPV+ or HPV- combined with an exploratory or validation set. The MPAC app presents results at both the pathway- and protein-level. On the pathway-level page (Figure 5A), Shiny app Panel A displays pathway enrichment results similar to the ones shown in Figure 2, Figure 4A, and Supplementary Figure 1A. Users can enter any pathway(s) of interest to look at their enrichments in MPAC-defined patient groups. To understand which proteins lead to a pathway enrichment, Panel B shows inferred pathway levels of all the proteins from a pathway. For example, in the pathway 'positive regulation of interleukin-2 biosynthetic process', CD28, CD3E, CD4, CD80, CD86, and PTPRC have positive inferred pathway levels in a majority of group I HPV+ exploratory patient samples (Figure 5A), suggesting they are the determinants resulting in this pathway's enrichment in group I patients.

In Panel B, users can enter or select any pathway of interest to examine their proteins' inferred pathway levels. At the protein-level page (Figure 5B), Panel C contains results for a group of user-specified proteins. It has three tabs displaying proteins' inferred pathway levels,

overall survival and immune cell composition of patients stratified by proteins' inferred pathway levels. These figures are similar to Figures 3B, 3D, 3E, 4B, and 4C; Supplementary Figures 9, 11, and 12, with the flexibility of showing the results for any user-specified protein(s) on any of the four TCGA-HNSCC datasets. Panel D shows a heatmap of the CNA, RNA, and pathway states for any user-entered protein as well as pathway states of this protein's pathway network neighbors. It is similar to Supplementary Figures 5-7 with the same flexibility as Panel C. In addition to the interactive data display, this app also contains documentation regarding re-generating figures presented in this manuscript, related papers, future developments, and acknowledgement. In summary, the MPAC Shiny app provides a convenient way to explore all the results generated from this work, especially those not presented as figures in this manuscript.

Discussion

We presented MPAC as a computational framework for calculating inferred pathway levels, predicting patient groups with biologically meaningful pathway profiles, and identifying key proteins with potential clinical associations. One group of HNSCC patients was predicted to have alterations in immune response pathways, and this group could not be identified from CNA or RNA-seq data alone. This finding illustrates the advantages of our pathway-based multi-omic approach. MPAC can use prior knowledge of pathway interactions in the form of a pathway network to integrate CNA and RNA-seq data and infer proteins' pathway behavior. A protein's pathway behavior from MPAC was not solely inferred from its CNA or RNA but also its pathway neighbors (Figure 3C; Supplementary Figures 5-8). Analysis based on CNA and RNA-seq data alone would miss this important biological principle.

Our analysis showed that MPAC can predict patient groups with potentially relevant clinical properties by their pathway profiles. The results presented above (Figure 2 and Supplementary Figure 1) identified an immune response patient group in both HPV+ and HPV- HNSCC, which echoes a new subtype defined by recent studies. One study (Huang *et al.*, 2021) applied a proteogenomic approach on 108 HPV- HNSCC patients. By considering CNA, RNA, miRNA, protein, and phosphopeptide data, the authors defined three subtypes of HNSCC: chromosome instability, basal, and immune. These subtypes mirrored the immune response patient groups MPAC identified in both HPV+ and HPV- patients. Moreover, the authors analyzed immune-hot tumors and revealed the presence of both cytotoxic immune cells (e.g., CD8+ T cells and M1 macrophages) and immunosuppressive cells (e.g., regulatory T cells and M2 macrophages). This is consistent with our analysis of HPV+ tumors stratified by inferred pathway levels of the seven proteins identified from the immune response patient group. Tumor samples with activated pathway levels of any of the seven proteins always had a higher fraction of CD8+ T cells, regulatory T cells, and M1 and M2 macrophages in both the exploratory and validation sets (Figures 3E and 4C). Second, a multi-omic analysis of thirty-three TCGA cancer types (Tiong *et al.*, 2022) identified gene groups enriched by immune response as well as cell cycle, which were also observed in MPAC's results (Figures 2 and 4A; Supplementary Figure 1A). The agreement between these two studies supports MPAC's discovery of the immune response patient group.

Further, not only did the seven proteins identified by MPAC correlate with immune cell composition, but their activated pathway levels also correlated with better overall survival. This was demonstrated in our exploratory patient set (Figure 3D and Supplementary Figure 9) and confirmed by our validation set (Figure 4B and Supplementary Figure 12). The corroboration by the validation set illustrates a major strength of MPAC. To understand these seven proteins'

clinical values and whether they could serve as biomarkers would require a prospective patient cohort, which is not available to us currently (As of May 6, 2024, according to cBioportal (de Bruijn *et al.*, 2023), <https://www.cbioportal.org/>), the only HNSCC dataset that has both CNA and RNA-seq data available is from TCGA, which is used in this work). However, the analyses here demonstrated how the MPAC software could be applied in a prospective setting.

MPAC has several advantages over the PARADIGM algorithm that it calls as a subroutine. MPAC uses a data-driven approach to define each genes' genomic states based on both tumor and normal tissue samples, whereas PARADIGM arbitrarily assigns the top, middle and lower third of omic-ranked genes as activated, normal and repressed (Vaske *et al.*, 2010). MPAC also provides downstream analyses on inferred pathway levels, including built-in permutation testing, defining altered pathways, predicting patient groups, and identifying key proteins with potential clinical implications. All of these functions have been implemented in an R package available through Bioconductor making it easier for others to use in their studies. The MPAC R Shiny app also supports convenient visualizations of the MPAC predictions.

Multi-omic integration methods have been developed for diverse applications (Maghsoudi *et al.*, 2022; Zitnik *et al.*, 2023), such as embedding single-cell data (Ashuach *et al.*, 2023; Argelaguet *et al.*, 2020), clustering cancer samples (Chauvel *et al.*, 2020; Wang *et al.*, 2014), and pathway reconstruction (Tuncbag *et al.*, 2016; Winkler *et al.*, 2022; Paull *et al.*, 2013). Multi-omics analyses have been particularly prominent in cancer, with pathway enrichment (Paczkowska *et al.*, 2020), representation learning (Leng *et al.*, 2022), supervised prediction of cancer subtypes or patient outcomes (Poirion *et al.*, 2021; Choi and Chae, 2023), and biologically interpretable neural networks (Wysocka *et al.*, 2023) as representative areas of study. MPAC's unique role in this methodological landscape is that through PARADIGM it

directly uses pathway interactions to combine information across omic data types, learn protein activities, and conduct downstream analysis with those protein activities.

In this work, we limited the input multi-omic data to CNA and RNA-seq, given PARADIGM's previous success with these two data types. With the availability of many other types of omic data from TCGA and the Clinical Proteomic Tumor Analysis Consortium (Huang *et al.*, 2021) on large cohorts of cancer patients, time course multi-omic data (G. Stein-O'Brien *et al.*, 2018), single-cell RNA-seq (Puram *et al.*, 2017), spatial transcriptomics (Li *et al.*, 2024; Lee *et al.*, 2024), and spatial proteomics (Causer *et al.*, 2023), one of our future goals is to make MPAC compatible with as many omic data types as possible. This requires extending the MPAC software as well as the input biological pathways to include knowledge on the relevant molecules and associated regulatory mechanisms.

Methods

Genomic and clinical datasets

We downloaded the TCGA HNSCC genomic datasets (Cancer Genome Atlas Network, 2015) from NCI GDC Data Portal version 29.0 (<https://portal.gdc.cancer.gov/>), which was released on March 31, 2021. Gene-level copy number scores were used for CNA and $\log_{10}(\text{FPKM}+1)$ values were used for RNA-seq. Patients' HPV status was obtained from their biospecimen manifest files. Patients' clinical data was downloaded from TCGA Pan-Cancer Atlas (Liu *et al.*, 2018) via <https://api.gdc.cancer.gov/data/1b5f413e-a8d1-4d10-92eb-7c4ae739ed81>. 492 HNSCC patients that had CNA, RNA-seq, and clinical data available were stratified by HPV status and then randomly divided into exploratory sets (71 HPV+ and 322 HPV-) and validation sets (18 HPV+ and 81 HPV-). Importantly, only the exploratory set was used for all MPAC algorithm development and refinement.

MPAC's pathway definitions were taken from the TCGA Pan-Cancer Atlas (Hoadley *et al.*, 2018), which compiled interactions from NCI-PID (Schaefer *et al.*, 2009), Reactome (Gillespie *et al.*, 2022), and KEGG (Kanehisa *et al.*, 2023) and superimposed them into a single network. The network included 19,477 entities, including 7,321 proteins, 9,349 complexes, 2,092 families, 591 abstract processes, 15 miRNAs, 82 RNAs, and 27 other types of entities. It also included 45,313 interactions containing 2,133 activations and 401 repressions at the transcript-level, 7,723 activations and 1,083 repressions at the protein-level, 24,870 and 9,103 memberships for complexes and families, respectively. The 2,085 Biological Process GO terms (Gene Ontology Consortium *et al.*, 2023) for characterizing patient or cell line pathway alteration were downloaded from the DrugCell (Kuenzi *et al.*, 2020) GitHub repository (https://github.com/idekerlab/DrugCell/blob/public/data/drugcell_ont.txt). GO terms from DrugCell had more distinct genes between parental and offspring terms because it required a parent to have ≥ 10 genes distinct from all child terms and have ≥ 30 genes more than any child. The root GO term (i.e., the ancestor of all the other GO terms), 'biological process', was not used in this study because it was not a specific functional description.

MPAC workflow

For TCGA data, the signs of CNA focal scores were used to define activated, normal (i.e., focal score is exactly zero), or repressed CNA state as the input for MPAC. To define input RNA state, a gene's RNA-seq expression levels from normal patient samples were fit to a Gaussian distribution. If a gene's expression levels in tumor samples fell within two standard deviations from the mean of this distribution, the gene's RNA state was defined as normal. Otherwise, its RNA state was repressed or activated depending on whether its expression level was below or above the two standard deviations from the mean.

MPAC ran PARADIGM in the default configuration (Vaske *et al.*, 2010) except with a more stringent expectation-maximization convergence criteria of change of likelihood $< 10^{-9}$ under a maximum of 10^4 iterations. To prepare permuted input, paired CNA and RNA states were randomly shuffled between all the genes within the patient. 100 permuted samples were prepared per each real tumor sample resulting in 49,200 permuted samples in total for the 492 patients. This large number of computational jobs were processed through UW-Madison's Center for High Throughput Computing (Center for High Throughput Computing, 2006) with HTCondor (Thain *et al.*, 2005).

A pathway entity's inferred pathway level from a real tumor sample was set to NA if it fell within three median absolute deviations of the inferred pathway levels from the corresponding 100 permuted samples. This filtering helped to remove inferred pathway levels that could be observed by chance. Entities with non-NA inferred pathway levels were mapped to the input pathway network. The largest connected subset of the pathway network with non-NA inferred pathway levels was kept for downstream analysis. Other entities not in this largest subset had their inferred pathway levels set to zero. This allowed us to focus on the entities that act together in pathways.

After the filtering by permuted samples and the largest pathway subset, an entity's pathway state was defined by the sign of its inferred pathway level, where a positive or negative inferred pathway level corresponded to an activated or repressed state, respectively, and a zero inferred pathway level corresponded to a normal state. Based on normal or altered pathway states, GO enrichment was calculated by Fisher's exact test, and the p-values were adjusted by the Benjamini and Hochberg procedure. Similarly, GO enrichment was calculated for the CNA and RNA inputs by their normal or altered states.

Patients were grouped by their adjusted p-values from GO enrichment based on CNA, RNA, or inferred pathway levels. A clustering method originally designed for single-cell RNA-seq analysis was adapted, where a patient tumor sample was treated as a cell and the $|\log_{10}(\text{adjusted } P)|$ was treated as a gene's expression level. Gene variance was modeled by the modelGeneVar function from the scran package (Lun *et al.*, 2016) (version 1.20.1), and the top 100 genes were selected. Patients were grouped by the Louvain method from the igraph R package (Csárdi *et al.*, 2024) (version 1.2.11) with 10 or 20 nearest neighbors for HPV+ or HPV-, respectively.

Changes in newer versions of the igraph R package affected the reproducibility of our results but not the main conclusions from our analyses. Starting in igraph version 1.3, the Louvain method is no longer deterministic and different runs may generate different clustering results (<https://github.com/igraph/rigraph/issues/539>). To evaluate the impact of non-deterministic clustering, we executed three independent batches of 10,000 random Louvain clustering runs. We used the Adjusted Rand Index (ARI) to measure the difference between groups from these random runs with the original groups from igraph version 1.2.11. The median ARIs were 0.92, 1.00, and 0.78 for the HPV+ exploratory set, HPV+ validation set, and HPV- exploratory set, respectively, (HPV- validation set groups were not used in this analysis), indicating small variations in the patient cluster membership. The 11- and 6-patient immune groups in the HPV+ exploratory and validation set were largely unchanged (identical in >9,200 and >8,800 out of 10,000 runs, respectively). For the HPV- exploratory set, the top five grouping results in each batch have either 30- or 32-patient groups and both groups are enriched with immune response GO terms. In summary, despite the randomness introduced in igraph version 1.3, our findings are still maintained.

To summarize pathway features for a group, we plotted heatmaps of \log_{10} (adjusted P) values for GO terms with adjusted P < 0.05 in 100% of the patients from the same group (e.g., Figure 2 and Supplementary Figure 1A). When very few GO terms met this criterion, we lowered the percentage threshold (e.g., $\geq 80\%$ in Figure 4A, $\geq 60\%$ in Supplementary Figure 1C) or by specifying a minimum number (e.g., ≥ 3 in Supplementary Figure 1B) of patients in order to include more GO terms.

For survival analysis, we used the inferred pathway levels of one or multiple proteins to stratify patients into two groups: those with all the protein(s) in activated pathway states (i.e., positive inferred pathway level values) and those not. A log-rank test p-value was calculated to compare the survival distribution of the two groups.

Software and data availability

The MPAC package is available at Bioconductor (<https://bioconductor.org/packages/MPAC>) and archived on Zenodo (<https://doi.org/10.5281/zenodo.10805479>). MPAC's Shiny app is accessible at <https://connect.doit.wisc.edu/content/122/>. The source code for MPAC's Shiny app is available at GitHub (https://github.com/pliu55/MPAC_Shiny) and archived on Zenodo (<https://doi.org/10.5281/zenodo.11623974>). Both the R package and Shiny app are available under the GPL-3.0 license.

Acknowledgements

We acknowledge funding from University of Wisconsin Carbone Cancer Center Support Grant P30 CA014520 and from a joint pilot program of the Morgridge Institute for Research and the NIH/NCI UW Comprehensive Cancer Center Support Grant ; A.G. and I.M.O. acknowledge pilot awards from the Wisconsin Head and Neck Cancer SPORE Career Enhancement Program P50

DE026787; I.M.O. and P.L. acknowledge support from the Wisconsin Head and Neck Cancer SPORE as part of the Biostatistics and Bioinformatics Core P50 CA278595; I.M.O. acknowledges support from the UW Institute for Clinical and Translational Research Pilot Award UL1 TR002373 and KL2 TR002374 from NIH/NCATS; I.M.O. acknowledges support from the NIH National Cancer Institute grant P01 CA250972; P.A. acknowledges support from the NIH National Cancer Institute grant P01 CA022443. This research was partially performed using the compute resources and assistance of the UW–Madison Center for High Throughput Computing in the Department of Computer Sciences. The authors acknowledge the use of the University of Wisconsin Carbone Cancer Center Cancer Informatics Shared Resource, supported by P30 CA014520. A.G. and P.A. are investigators of the Morgridge Institute for Research and an investigator and the director, respectively, of the John and Jeanne Rowe Center for Research in Virology and gratefully acknowledge their support. We thank H. Adam Steinberg for creative contributions to figure design and data visualization, Alicia Williams for organization feedback and manuscript editing, and David Merrell, Deric Wheeler, Mari Iida, Paul Lambert, and Randall Kimple for project feedback.

Figure legends

Figure 1. Overview of the MPAC workflow. MPAC calculates inferred pathway levels (IPLs) from real and permuted CNA and RNA data. It filters real IPLs using the permuted IPLs to remove spurious IPLs. Then, MPAC focuses on the largest pathway subset network with filtered IPLs to compute Gene Ontology (GO) term enrichment, predict patient groups, and identify key group-specific proteins.

Figure 2. MPAC predicted functionally distinct patient groups in the HPV+ exploratory set. Patient groups were derived from GO term enrichment based on inferred pathway levels (IPLs).

Figure 3. Seven proteins identified in the immune response patient group from HPV+ exploratory set correlated with patient overall survival. (A) Consensus pathway submodules in the eleven immune response patient samples from group I. Proteins are colored in red and other pathway entities are black; (B) inferred pathway levels (IPLs) of submodule proteins in the 71 HPV+ exploratory set samples. Except for CD3G, which had both positive and negative IPLs (denoted by *) in group I, the other seven proteins had positive IPLs; (C) CNA, RNA, and pathway states of CD86 as well as pathway states of its pathway network neighbors in a group I patient sample TCGA-CR-6487. B7-1 homodimer/B7-2 (family), B7 family (family), and phospho CD28 B7-1/B7-2 (complex) do not have activated pathway state in all the eleven group I patients (Supplementary Figure 5B) and thus are not included in Figure 3A; (D) Overall survival of the HPV+ exploratory set stratified by IPLs of the seven proteins combined; (E) HPV+ exploratory set immune cell compositions stratified by IPLs of all seven proteins combined.

Figure 4. Independent validation set confirmed MPAC's immune response group and key proteins. (A) Grouping of HPV+ validation set patient samples. The selection threshold was lowered to $\geq 80\%$ in order to include more GO terms; (B) Overall survival of the HPV+ validation stratified by the inferred pathway levels (IPLs) of all the seven submodule proteins combined; (C) HPV+ validation set immune cell compositions stratified by the IPLs of all the seven proteins combined.

Figure 5. Screenshots of an R Shiny app displaying MPAC results from the HPV+

exploratory set. (A) The upper panel shows enrichment of multiple user-selected pathways and the lower panel shows protein inferred pathway levels (IPLs) from a user-selected pathway; **(B)** The upper panel shows IPLs of multiple user-selected proteins and the lower panel shows the pathway states of a user-selected protein and its pathway neighbors as well as its CNA and RNA state.

Supplementary figure legends

Supplementary Figure 1. HPV- exploratory set patient samples grouped by GO term enrichment based on IPL (A), CNA (B), or RNA (C). The selection threshold was lowered to ≥ 3 **(B)** or $\geq 60\%$ **(C)** to include more GO terms to avoid bias.

Supplementary Figure 2. HPV+ exploratory set patient samples grouped by GO term enrichment based on CNA (A) or RNA (B). The selection threshold was lowered to ≥ 2 **(A)** or $\geq 60\%$ **(B)** to include more GO terms to avoid bias.

Supplementary Figure 3. HPV+ exploratory set patient samples grouped by MPAC's input RNA states. (A) Clustering results by K-means under a pre-specified total number of groups ranging from two to six; **(B)** Clustering results by K-means with a pre-specified six groups; **(C–G)** Top significantly enriched GO terms in group I **(C)**, II **(D)**, III **(E)**, IV **(F)**, and V **(G)**. The selection threshold was lowered to $\geq 95\%$ **(F)**, $\geq 70\%$ **(G)**, or $\geq 50\%$ **(C–E)** in order to include more GO terms to avoid bias.

Supplementary Figure 4. HPV+ exploratory set patient samples grouped by MPAC's input

CNA states. (A) Clustering results by K-means under a pre-specified total number of groups

ranging from two to six; **(B)** Clustering results by K-means with a pre-specified six groups.

Group III, IV, V, and VI contained one sample each.

Supplementary Figure 5. CNA, RNA, and pathway states of CD28 (A), CD86 (B), and LCP2

(C), as well as pathway states of their pathway network neighbors in the eleven group I

patients.

Supplementary Figure 6. CNA, RNA, and pathway states of TYK2 (A) and IL12RB1 (B), as

well as pathway states of their pathway network neighbors in the eleven group I patients.

Supplementary Figure 7. CNA, RNA, and pathway states of CD247 (A) and FASLG (B), as

well as pathway states of their pathway network neighbors in the eleven group I patients.

Supplementary Figure 8. CNA, RNA, and pathway states of FASLG as well as pathway

states of its pathway network neighbors in a group I patient TCGA-CR-7385.

Supplementary Figure 9. Overall survival of HPV+ exploratory set patient samples

stratified by the inferred pathway levels (IPLs) of proteins from the same submodule (A)

or individual protein (B). Not all overall survival tests were statistically significant under a $P <$

0.01 cutoff.

Supplementary Figure 10. Progression-free survival of HPV+ exploratory set patient samples stratified by the inferred pathway levels (IPLs) of all the seven proteins combined (A); proteins from the same submodule (B), or individual protein (C). Not all progression-free survival tests were statistically significant under a $P < 0.01$ cutoff.

Supplementary Figure 11. Immune cell compositions stratified by the inferred pathway levels (IPLs) of all the seven proteins combined for HPV+ exploratory (A) and validation (B) set. Adjusted Wilcoxon test P on the cell composition difference between two groups of patient samples were shown for each cell type. Eosinophil was not found in either the exploratory or validation sets. T cell CD4+ naïve was not found in the validation set.

Supplementary Figure 12. Overall survival of HPV+ validation set patient samples stratified by the inferred pathway levels (IPLs) of proteins from the same submodule (A) or individual protein (B). Not all overall survival tests were statistically significant under a $P < 0.01$ cutoff.

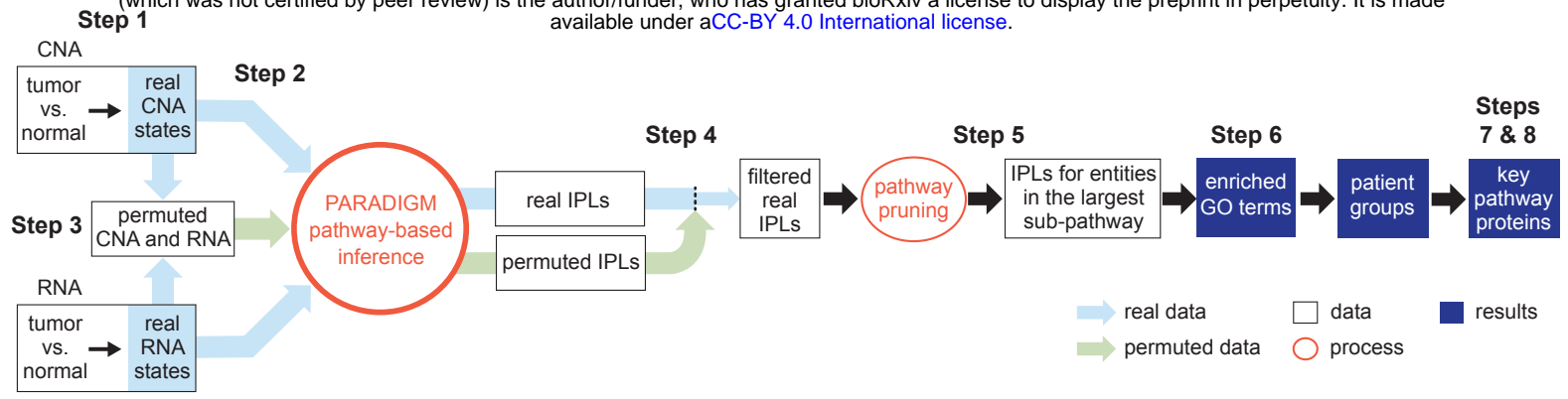
References

- Argelaguet, R. *et al.* (2020) MOFA+: a statistical framework for comprehensive integration of multi-modal single-cell data. *Genome Biol*, **21**, 111.
- Ashuach, T. *et al.* (2023) MultiVI: deep generative model for the integration of multimodal data. *Nat Methods*, **20**, 1222–1231.
- Berger, A.C. *et al.* (2018) A Comprehensive Pan-Cancer Molecular Study of Gynecologic and Breast Cancers. *Cancer Cell*, **33**, 690-705.e9.
- Broyde, J. *et al.* (2021) Oncoprotein-specific molecular interaction maps (SigMaps) for cancer network analyses. *Nat Biotechnol*, **39**, 215–224.
- de Bruijn, I. *et al.* (2023) Analysis and Visualization of Longitudinal Genomic and Clinical Data from the AACR Project GENIE Biopharma Collaborative in cBioPortal. *Cancer Res*, **83**, 3861–3867.
- Cancer Genome Atlas Network (2015) Comprehensive genomic characterization of head and neck squamous cell carcinomas. *Nature*, **517**, 576–582.
- Causser, A. *et al.* (2023) Deep spatial-omics analysis of Head & Neck carcinomas provides

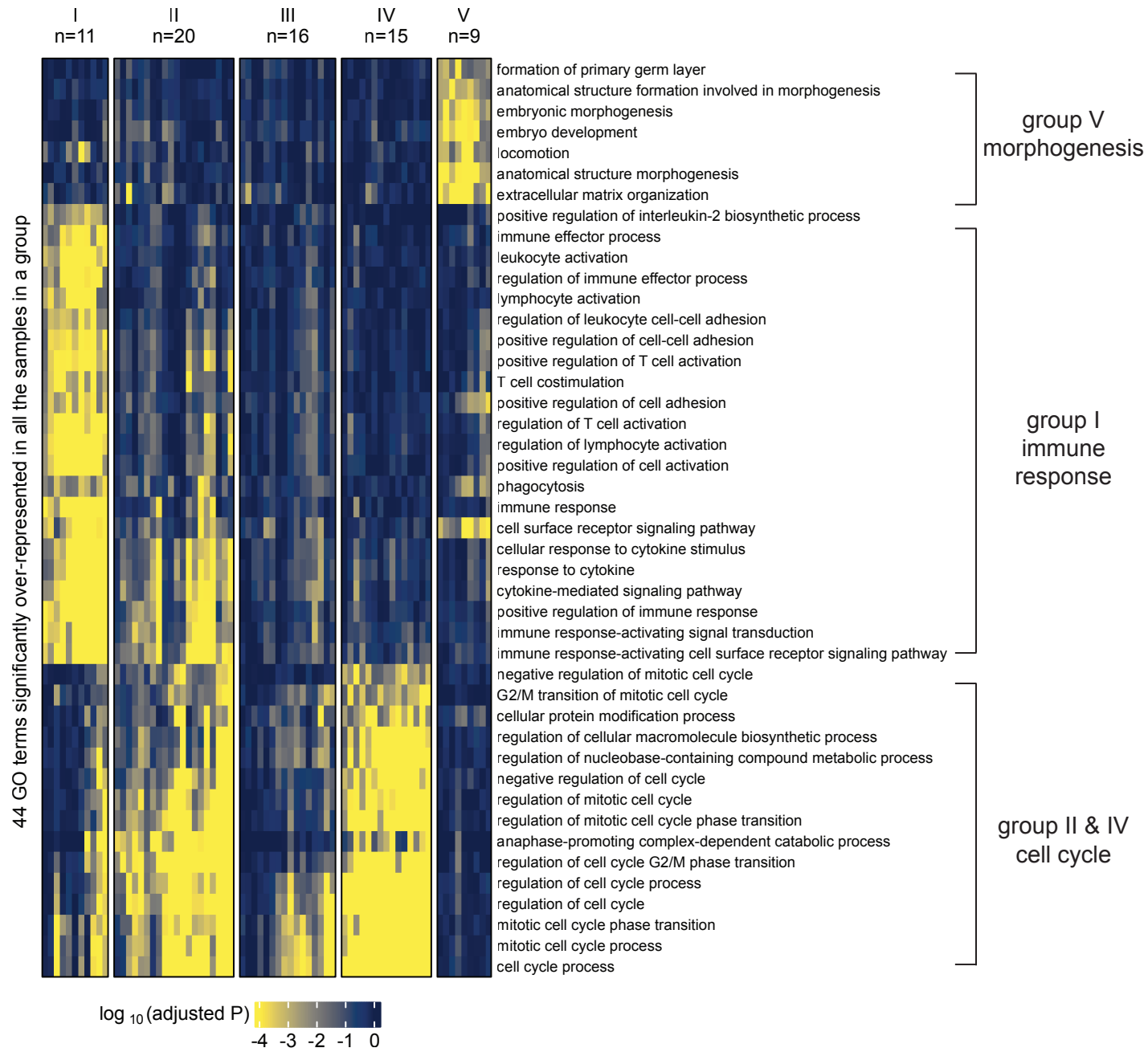
- alternative therapeutic targets and rationale for treatment failure. *NPJ Precis Oncol*, **7**, 89.
- Center for High Throughput Computing (2006) Center for High Throughput Computing. DOI: [10.21231/GNT1-HW21](https://doi.org/10.21231/GNT1-HW21).
- Chauvel,C. *et al.* (2020) Evaluation of integrative clustering methods for the analysis of multi-omics data. *Brief Bioinform*, **21**, 541–552.
- Chen,Z. *et al.* (2022) IL-12RB1: a novel immune prognostic biomarker for oral squamous cell carcinoma and linked to PD-1/PD-L1 expression in the tumor immune microenvironment. *Ann Transl Med*, **10**, 144.
- Chi,H. *et al.* (2022) Natural killer cell-related prognosis signature characterizes immune landscape and predicts prognosis of HNSCC. *Front Immunol*, **13**, 1018685.
- Choi,J.M. and Chae,H. (2023) moBRCA-net: a breast cancer subtype classification framework based on multi-omics attention neural networks. *BMC Bioinformatics*, **24**, 169.
- Csárdi,G. *et al.* (2024) igraph for R: R interface of the igraph library for graph theory and network analysis. <https://zenodo.org/doi/10.5281/zenodo.7682609>.
- Duan,R. *et al.* (2021) Evaluation and comparison of multi-omics data integration methods for cancer subtyping. *PLoS Comput Biol*, **17**, e1009224.
- Dugourd,A. *et al.* (2021) Causal integration of multi-omics data with prior knowledge to generate mechanistic hypotheses. *Mol Syst Biol*, **17**, e9730.
- Gene Ontology Consortium *et al.* (2023) The Gene Ontology knowledgebase in 2023. *Genetics*, **224**, iyad031.
- Gillespie,M. *et al.* (2022) The reactome pathway knowledgebase 2022. *Nucleic Acids Res*, **50**, D687–D692.
- Hanahan,D. (2022) Hallmarks of Cancer: New Dimensions. *Cancer Discov*, **12**, 31–46.
- He,Y. *et al.* (2022) TYK2 correlates with immune infiltration: A prognostic marker for head and neck squamous cell carcinoma. *Front Genet*, **13**, 1081519.
- Hoadley,K.A. *et al.* (2018) Cell-of-Origin Patterns Dominate the Molecular Classification of 10,000 Tumors from 33 Types of Cancer. *Cell*, **173**, 291-304.e6.
- Hoadley,K.A. *et al.* (2014) Multiplatform analysis of 12 cancer types reveals molecular classification within and across tissues of origin. *Cell*, **158**, 929–944.
- Huang,C. *et al.* (2021) Proteogenomic insights into the biology and treatment of HPV-negative head and neck squamous cell carcinoma. *Cancer Cell*, **39**, 361-379.e16.
- Kanehisa,M. *et al.* (2023) KEGG for taxonomy-based analysis of pathways and genomes. *Nucleic Acids Res*, **51**, D587–D592.
- Kuenzi,B.M. *et al.* (2020) Predicting Drug Response and Synergy Using a Deep Learning Model of Human Cancer Cells. *Cancer Cell*, **38**, 672-684.e6.
- Lee,H. *et al.* (2024) Deciphering head and neck cancer microenvironment: Single-cell and spatial transcriptomics reveals human papillomavirus-associated differences. *J Med Virol*, **96**, e29386.
- Leng,D. *et al.* (2022) A benchmark study of deep learning-based multi-omics data fusion methods for cancer. *Genome Biol*, **23**, 171.
- Li,C. *et al.* (2024) Spatial and Single-Cell Transcriptomics Reveal a Cancer-Associated Fibroblast Subset in HNSCC That Restricts Infiltration and Antitumor Activity of CD8+ T Cells. *Cancer Res*, **84**, 258–275.
- Liu,J. *et al.* (2018) An Integrated TCGA Pan-Cancer Clinical Data Resource to Drive High-Quality Survival Outcome Analytics. *Cell*, **173**, 400-416.e11.
- Lun,A.T.L. *et al.* (2016) A step-by-step workflow for low-level analysis of single-cell RNA-seq data with Bioconductor. *F1000Res*, **5**, 2122.
- Maghsoudi,Z. *et al.* (2022) A comprehensive survey of the approaches for pathway analysis using multi-omics data integration. *Brief Bioinform*, **23**, bbac435.
- Mody,M.D. *et al.* (2021) Head and neck cancer. *Lancet*, **398**, 2289–2299.

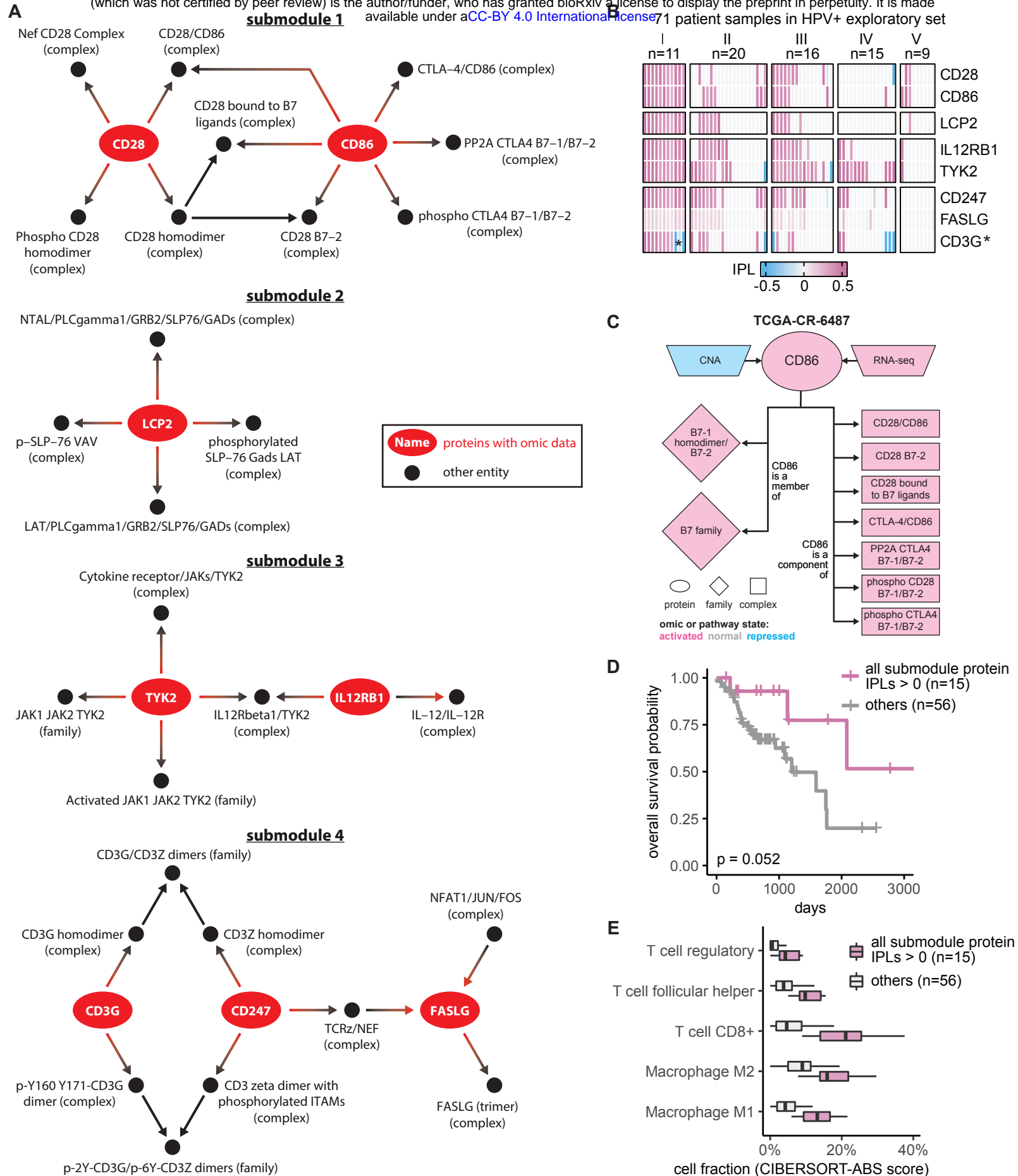
- Newman,A.M. *et al.* (2015) Robust enumeration of cell subsets from tissue expression profiles. *Nat Methods*, **12**, 453–457.
- Ng,S. *et al.* (2012) PARADIGM-SHIFT predicts the function of mutations in multiple cancers using pathway impact analysis. *Bioinformatics*, **28**, i640–i646.
- Paczkowska,M. *et al.* (2020) Integrative pathway enrichment analysis of multivariate omics data. *Nat Commun*, **11**, 735.
- Paull,E.O. *et al.* (2021) A modular master regulator landscape controls cancer transcriptional identity. *Cell*, **184**, 334-351.e20.
- Paull,E.O. *et al.* (2013) Discovering causal pathways linking genomic events to transcriptional states using Tied Diffusion Through Interacting Events (TieDIE). *Bioinformatics*, **29**, 2757–2764.
- Picard,M. *et al.* (2021) Integration strategies of multi-omics data for machine learning analysis. *Comput Struct Biotechnol J*, **19**, 3735–3746.
- Poirion,O.B. *et al.* (2021) DeepProg: an ensemble of deep-learning and machine-learning models for prognosis prediction using multi-omics data. *Genome Med*, **13**, 112.
- Powell,S.F. *et al.* (2021) The Key Differences between Human Papillomavirus-Positive and - Negative Head and Neck Cancers: Biological and Clinical Implications. *Cancers (Basel)*, **13**, 5206.
- Puram,S.V. *et al.* (2017) Single-Cell Transcriptomic Analysis of Primary and Metastatic Tumor Ecosystems in Head and Neck Cancer. *Cell*, **171**, 1611-1624.e24.
- Schaefer,C.F. *et al.* (2009) PID: the Pathway Interaction Database. *Nucleic Acids Res*, **37**, D674-679.
- Sedgewick,A.J. *et al.* (2013) Learning subgroup-specific regulatory interactions and regulator independence with PARADIGM. *Bioinformatics*, **29**, i62-70.
- Stein-O'Brien,G. *et al.* (2018) Integrated time course omics analysis distinguishes immediate therapeutic response from acquired resistance. *Genome Med*, **10**, 37.
- Stein-O'Brien,G.L. *et al.* (2018) Enter the Matrix: Factorization Uncovers Knowledge from Omics. *Trends Genet*, **34**, 790–805.
- Thain,D. *et al.* (2005) Distributed computing in practice: the Condor experience. *Concurrency Computat.: Pract. Exper.*, **17**, 323–356.
- Tiong,K.-L. *et al.* (2022) An integrated analysis of the cancer genome atlas data discovers a hierarchical association structure across thirty three cancer types. *PLoS Digit Health*, **1**, e0000151.
- Tuncbag,N. *et al.* (2016) Network-Based Interpretation of Diverse High-Throughput Datasets through the Omics Integrator Software Package. *PLoS Comput Biol*, **12**, e1004879.
- Vaske,C.J. *et al.* (2010) Inference of patient-specific pathway activities from multi-dimensional cancer genomics data using PARADIGM. *Bioinformatics*, **26**, i237-245.
- de Vos,L. *et al.* (2020) The landscape of CD28, CD80, CD86, CTLA4, and ICOS DNA methylation in head and neck squamous cell carcinomas. *Epigenetics*, **15**, 1195–1212.
- Wang,B. *et al.* (2014) Similarity network fusion for aggregating data types on a genomic scale. *Nat Meth*, **11**, 333–337.
- Wang,Q. *et al.* (2022) A novel immune signature predicts immunotherapy responsiveness and reveals the landscape of the tumor immune microenvironment in head and neck squamous cell carcinoma. *Front Genet*, **13**, 1051051.
- Winkler,S. *et al.* (2022) De novo identification of maximally deregulated subnetworks based on multi-omics data with DeRegNet. *BMC Bioinformatics*, **23**, 139.
- Wysocka,M. *et al.* (2023) A systematic review of biologically-informed deep learning models for cancer: fundamental trends for encoding and interpreting oncology data. *BMC Bioinformatics*, **24**, 198.
- Zhao,Y. *et al.* (2021) Bayesian network-driven clustering analysis with feature selection for high-dimensional multi-modal molecular data. *Sci Rep*, **11**, 5146.

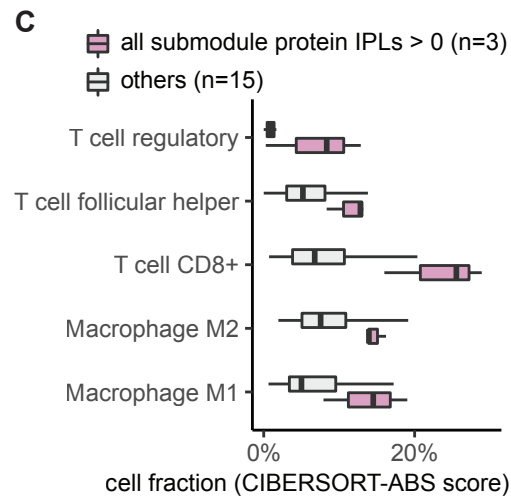
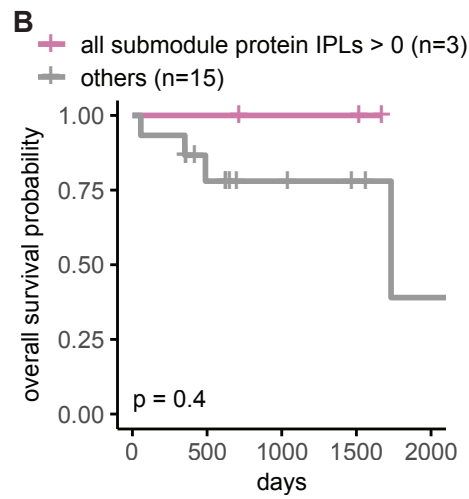
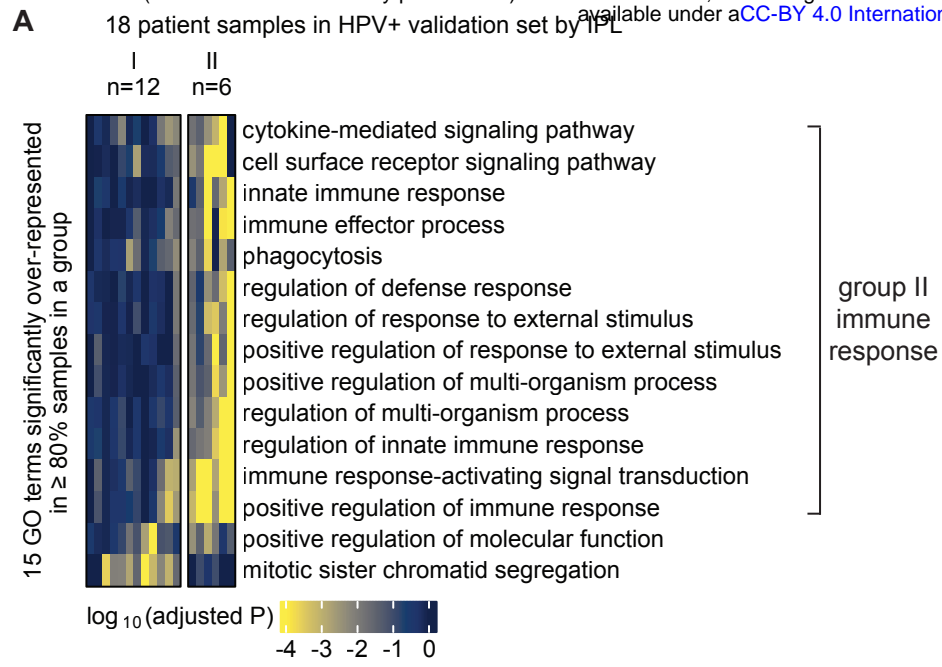
- Zhu,C. *et al.* (2022) Immune Infiltration Characteristics and a Gene Prognostic Signature Associated With the Immune Infiltration in Head and Neck Squamous Cell Carcinoma. *Front Genet*, **13**, 848841.
- Zitnik,M. *et al.* (2023) Current and future directions in network biology. *arXiv:2309.08478*.

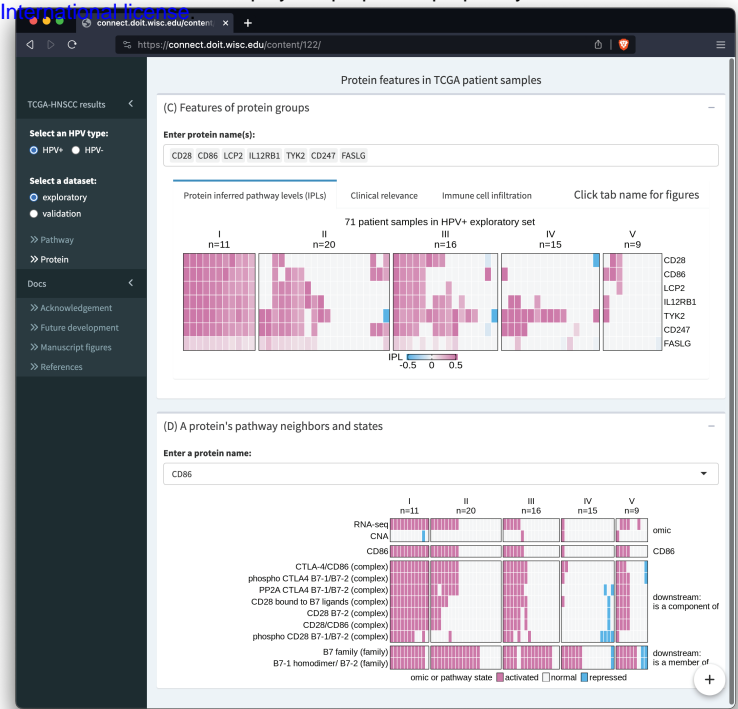
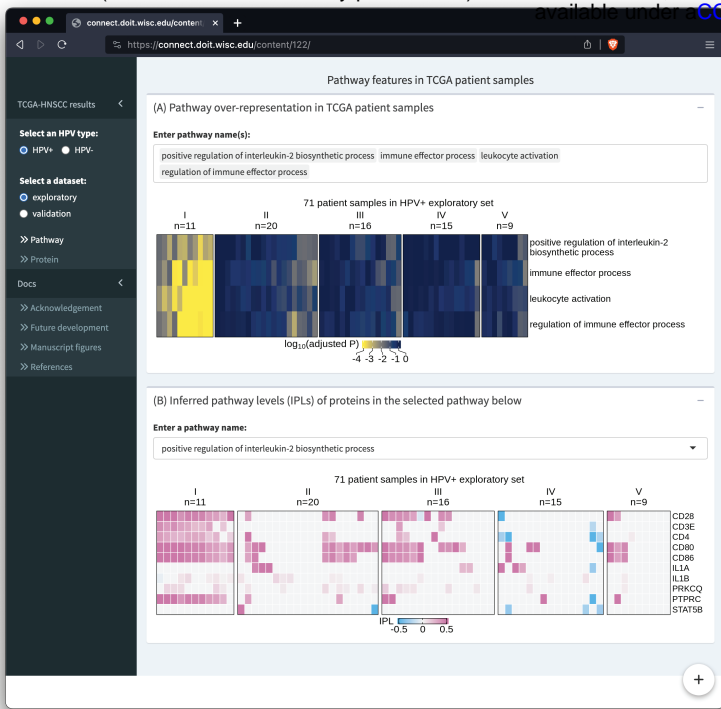


71 patient samples in HIV-1 exploratory set by IFL



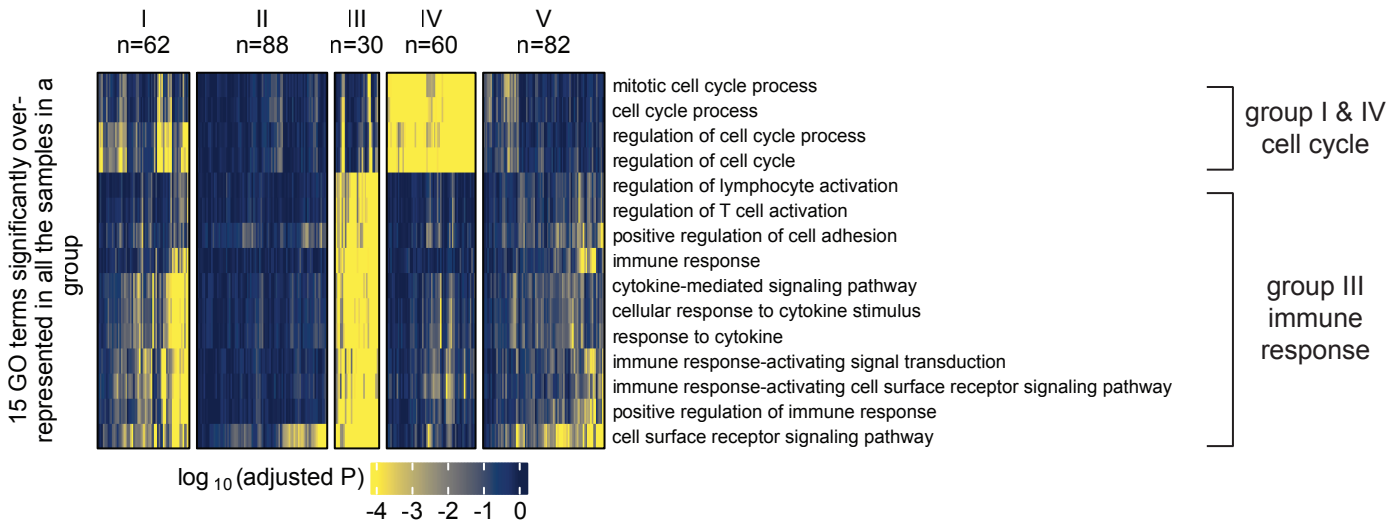






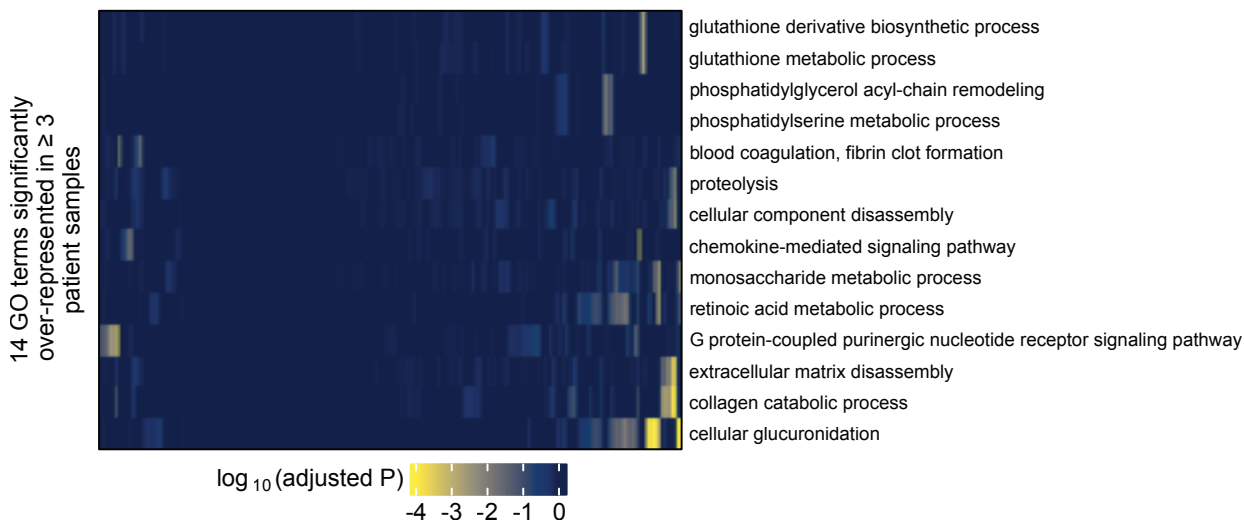
A

322 patient samples in HPV- exploratory set by IPL



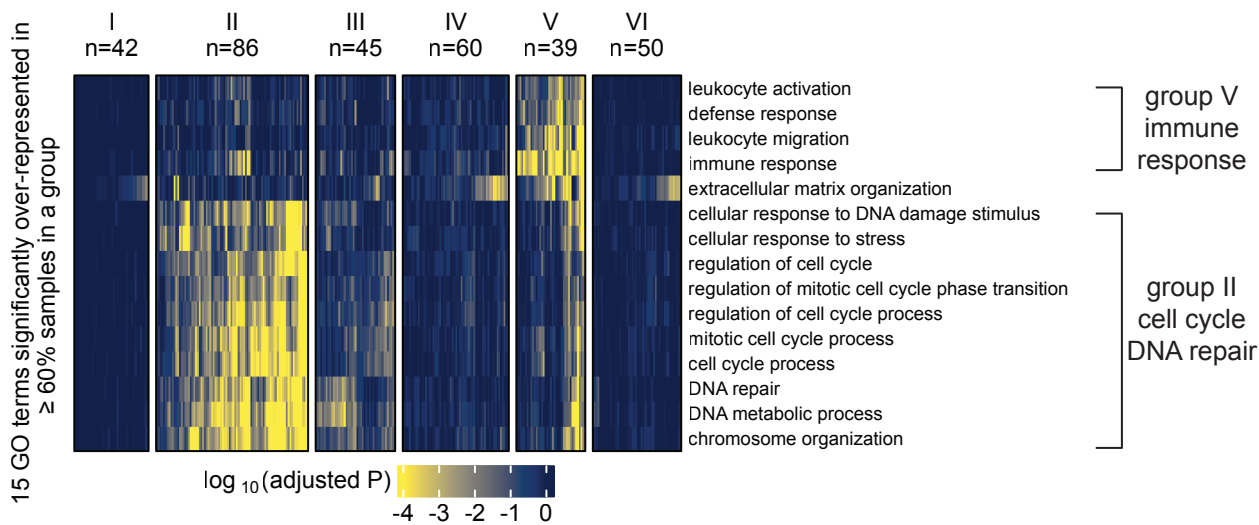
B

322 patient samples in HPV- exploratory set by CNA



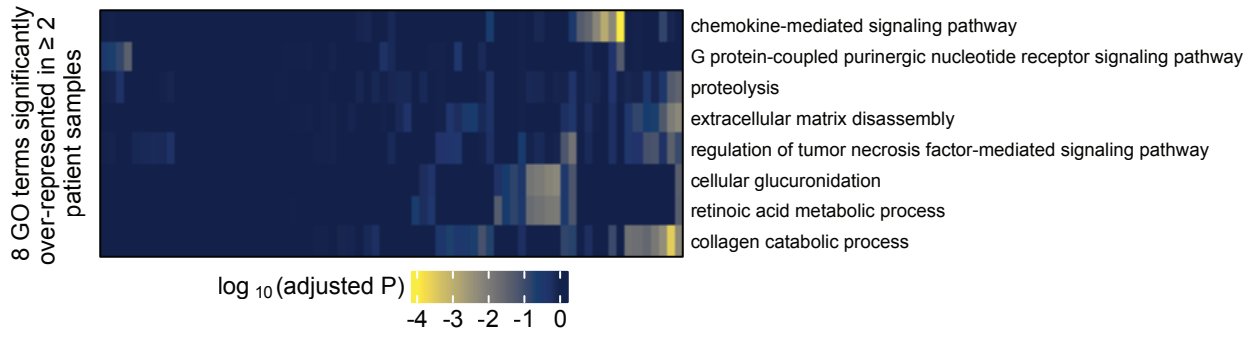
C

322 patient samples in HPV- exploratory set by RNA



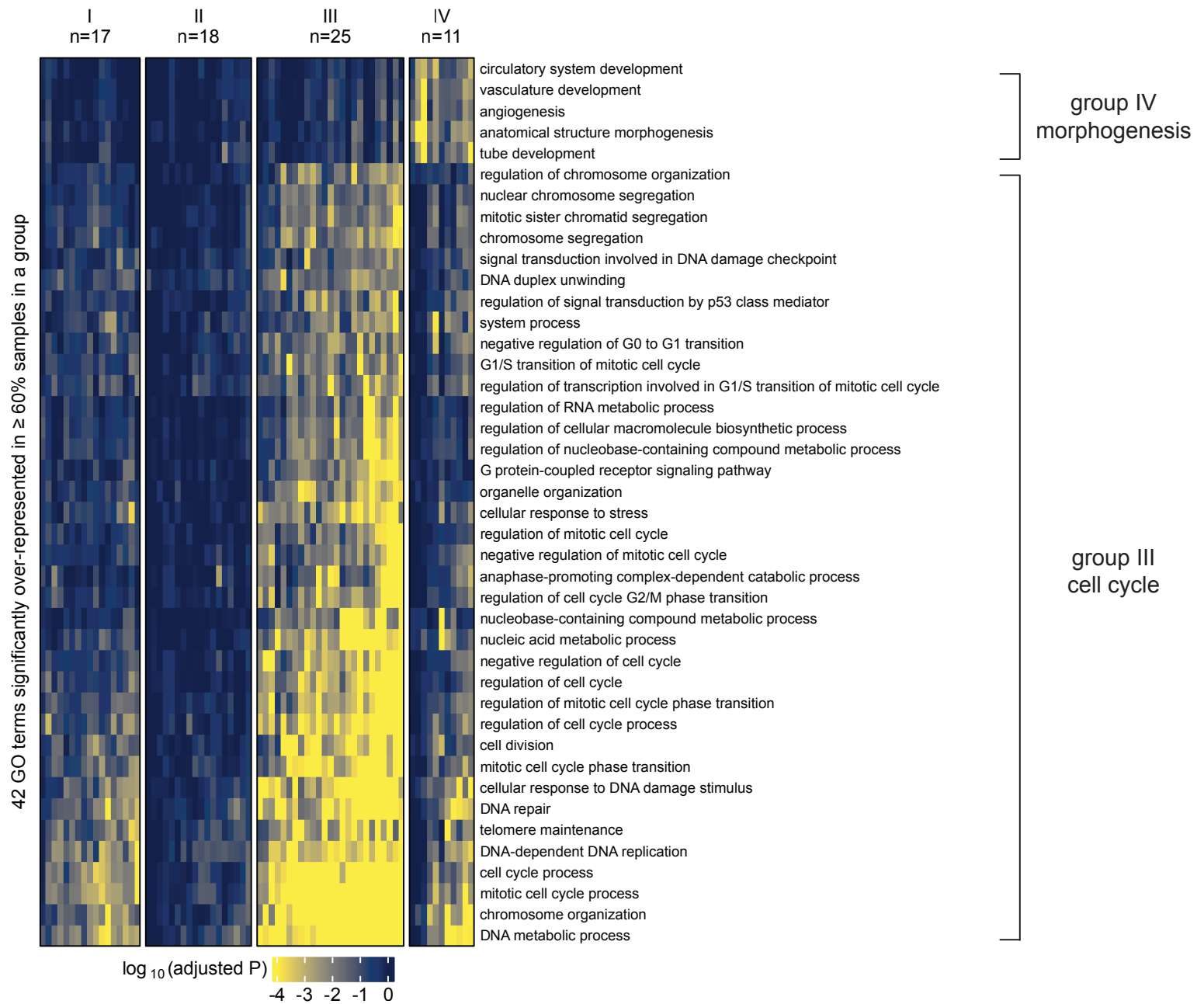
A

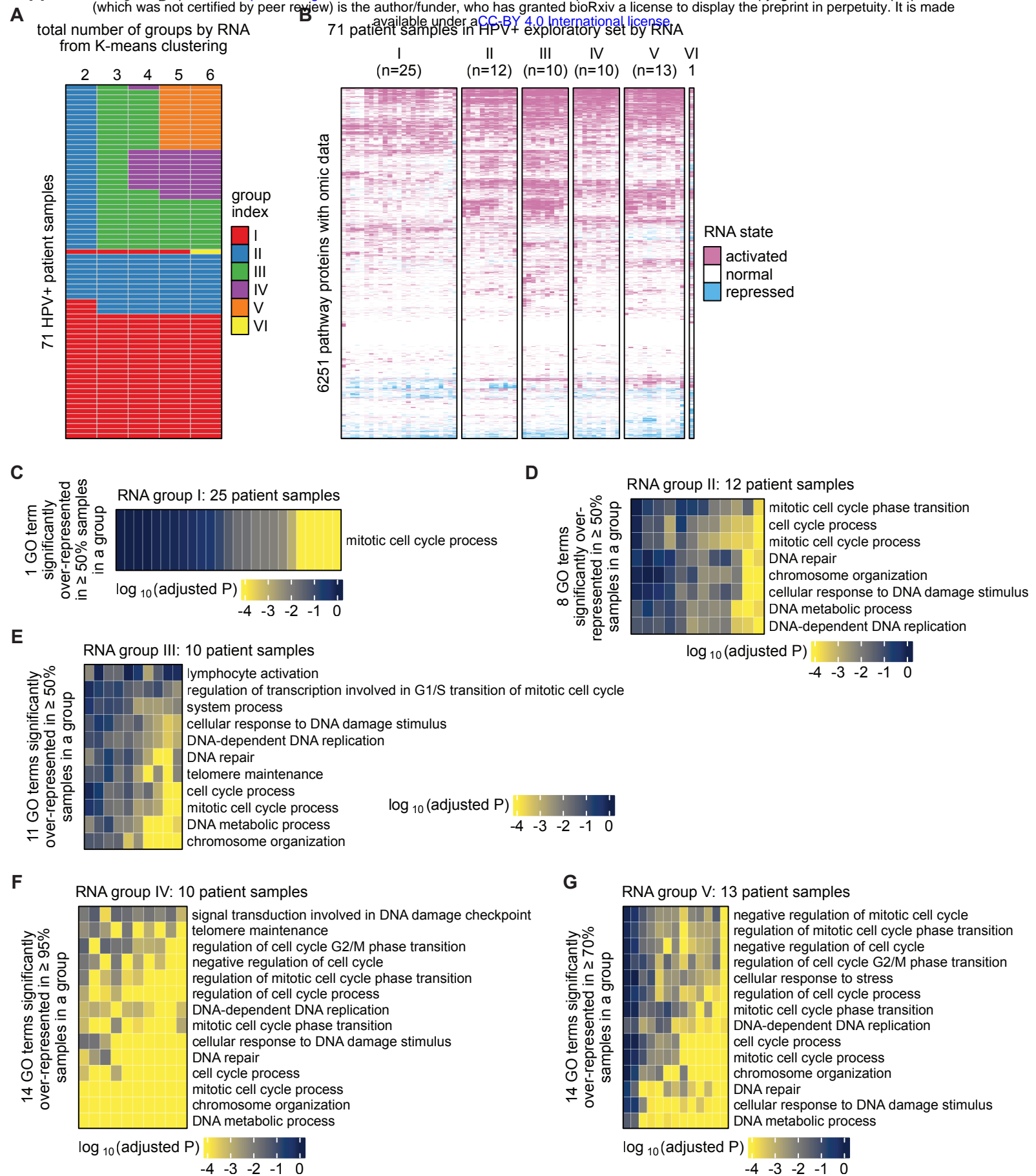
71 patient samples in HPV+ exploratory set by CNA



B

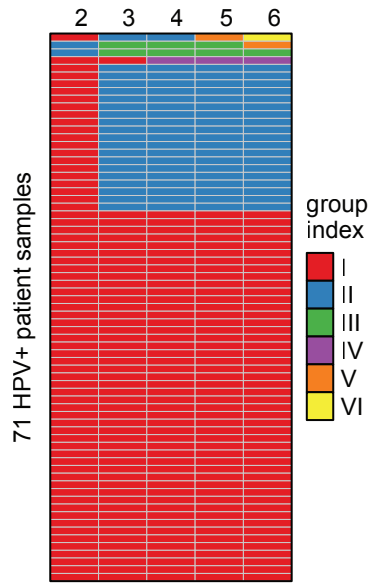
71 patient samples in HPV+ exploratory set by RNA





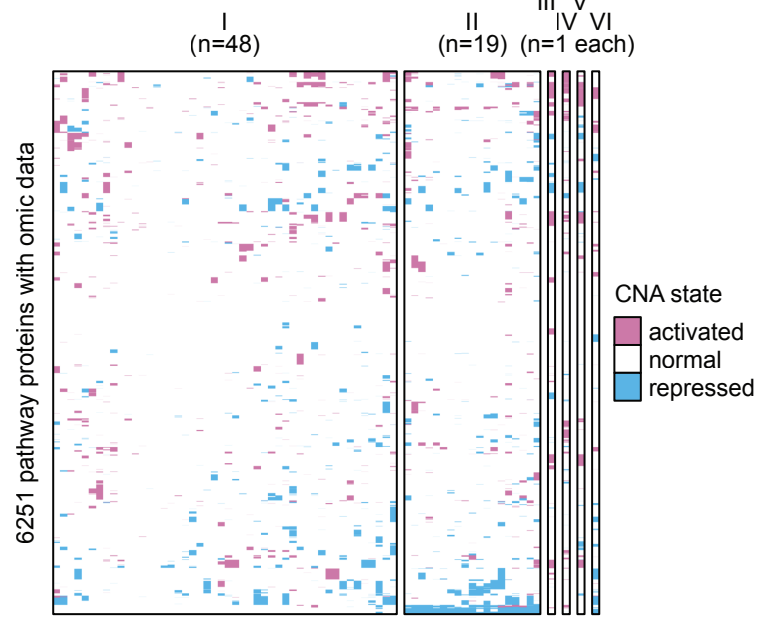
A

total number of groups by CNA
from K-means clustering



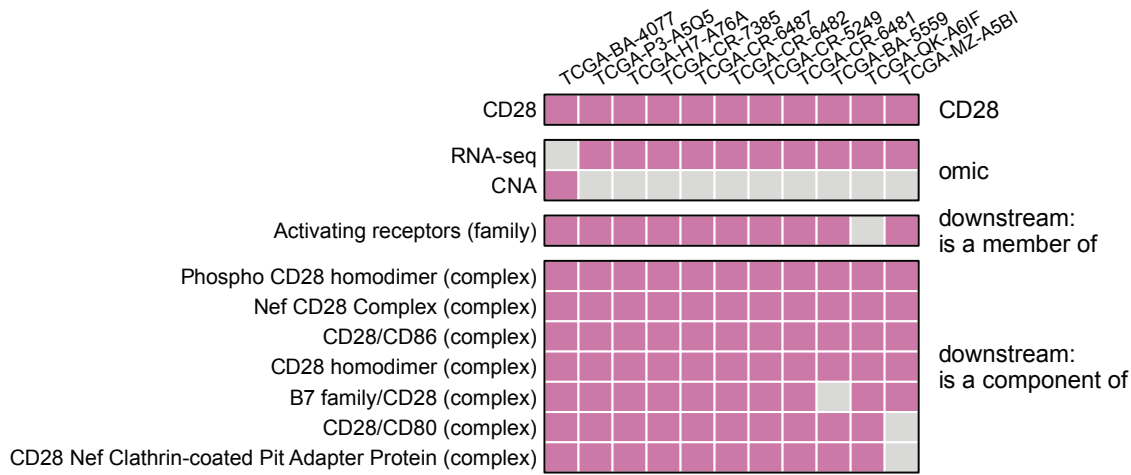
B

71 patient samples in HPV+ exploratory set by CNA



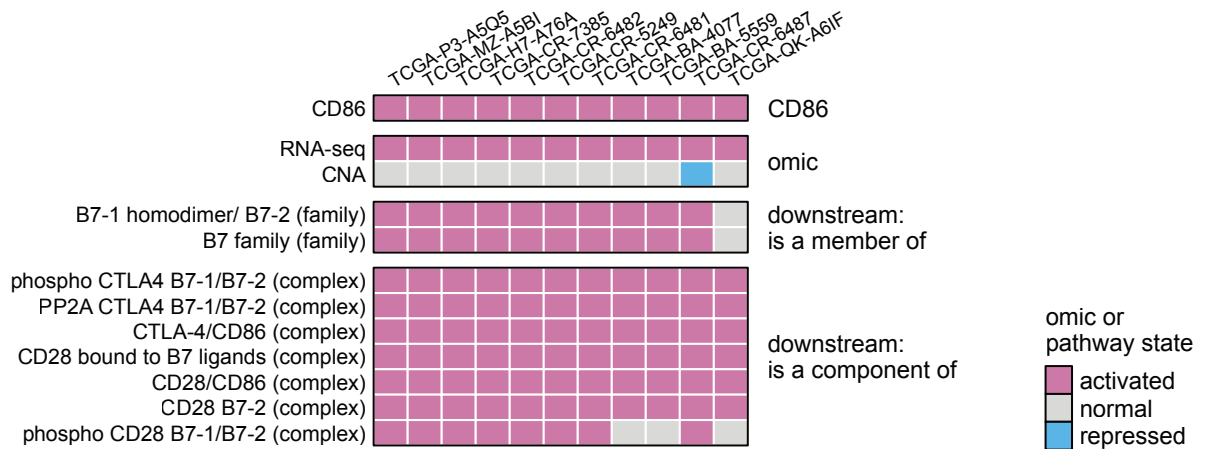
A

11 group I patient samples in HPV+ exploratory set



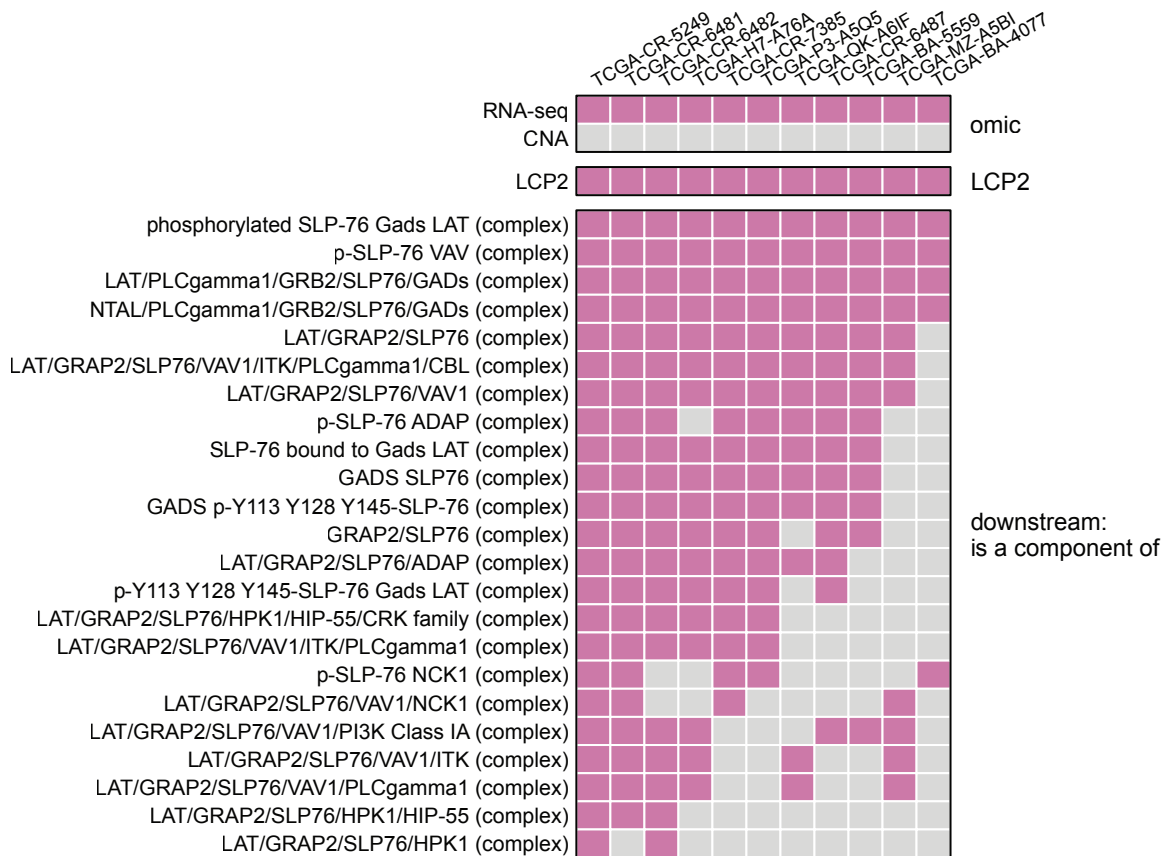
B

11 group I patient samples in HPV+ exploratory set



C

11 group I patient samples in HPV+ exploratory set



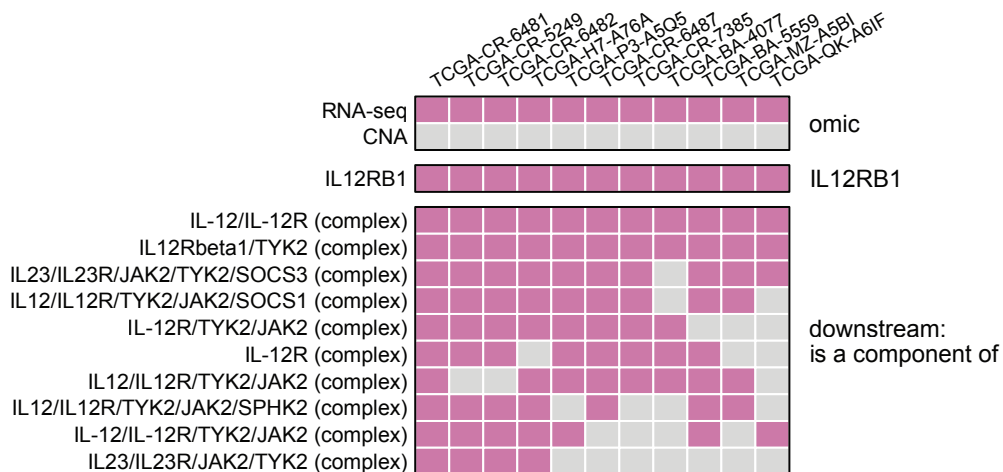
A

11 group I patient samples in HPV+ exploratory set



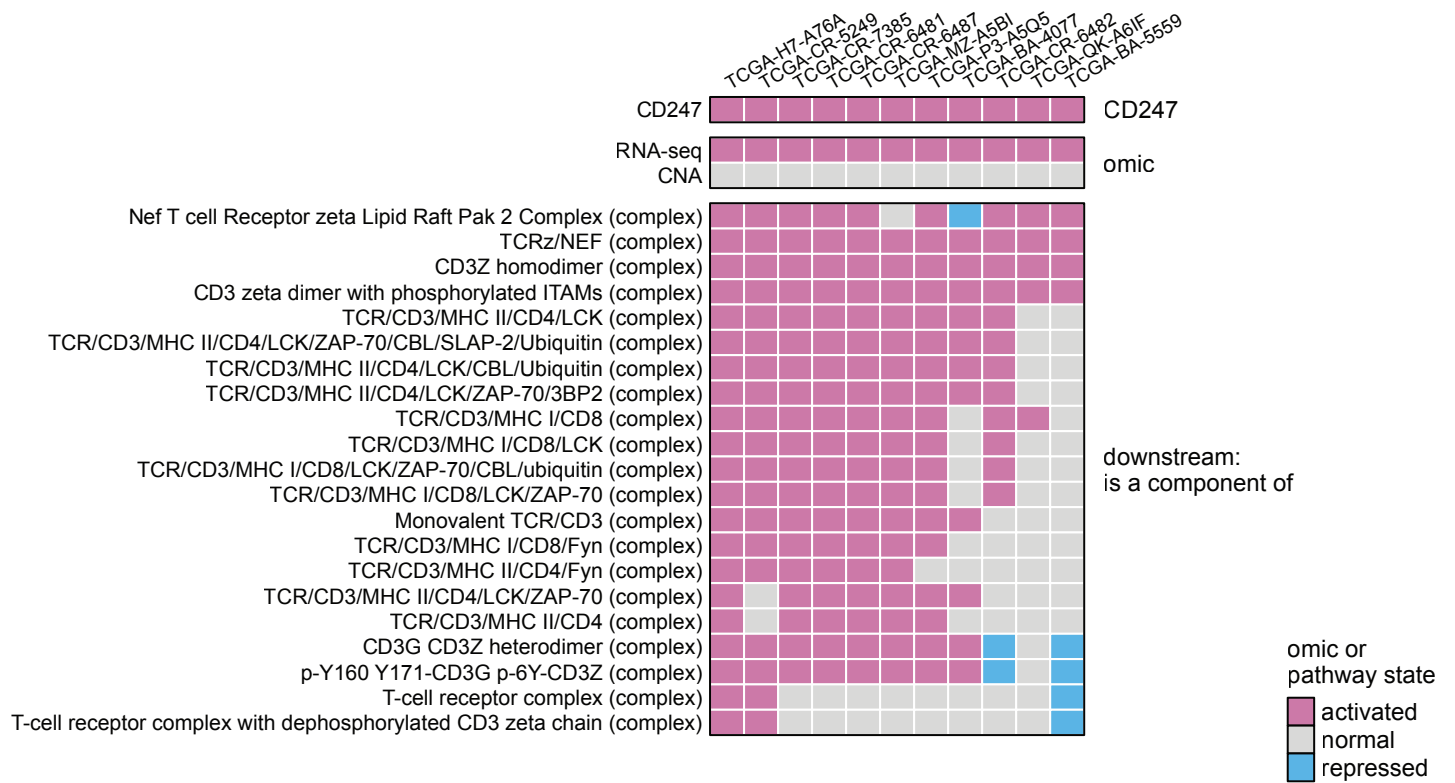
B

11 group I patient samples in HPV+ exploratory set



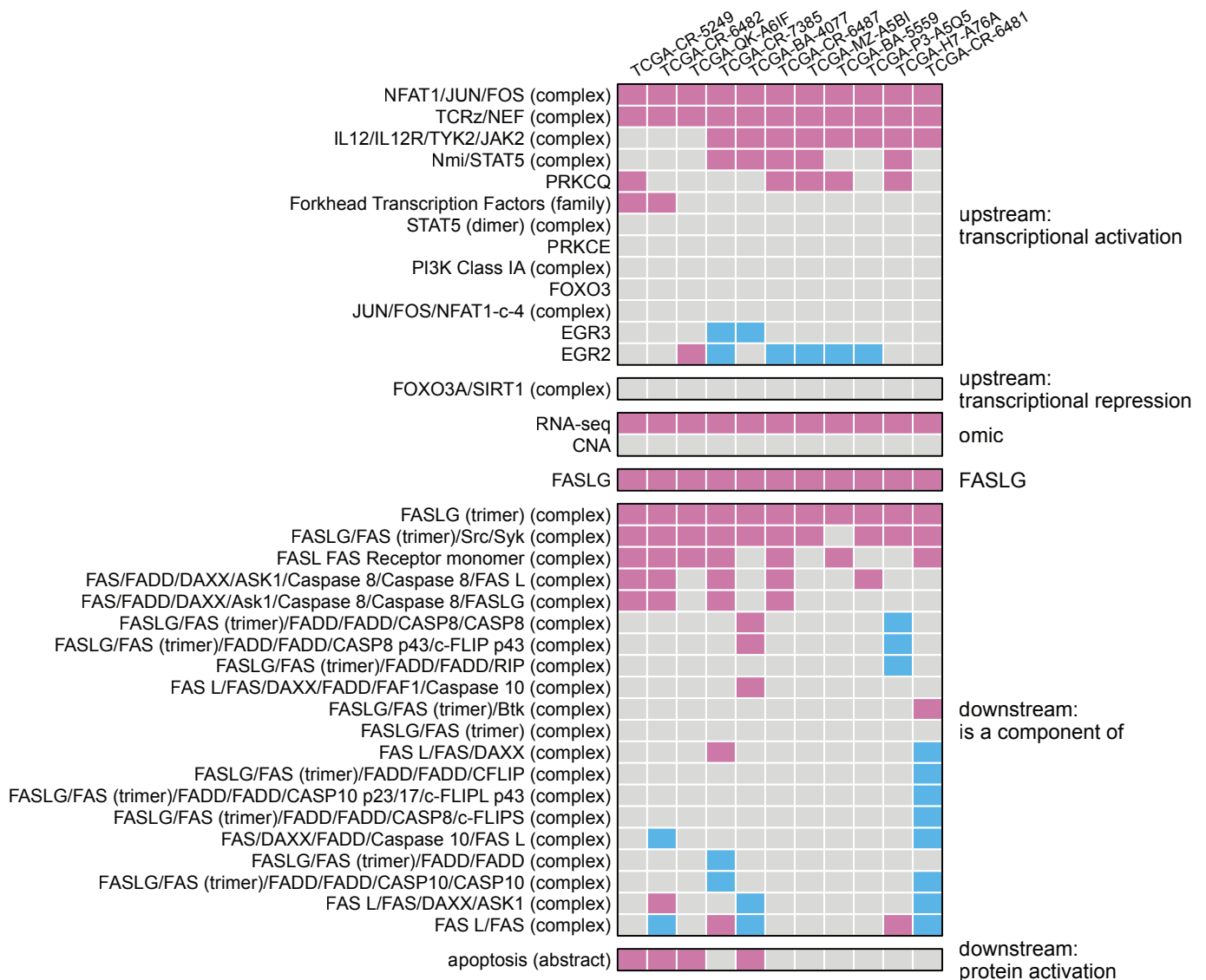
A

11 group I patient samples in HPV+ exploratory set

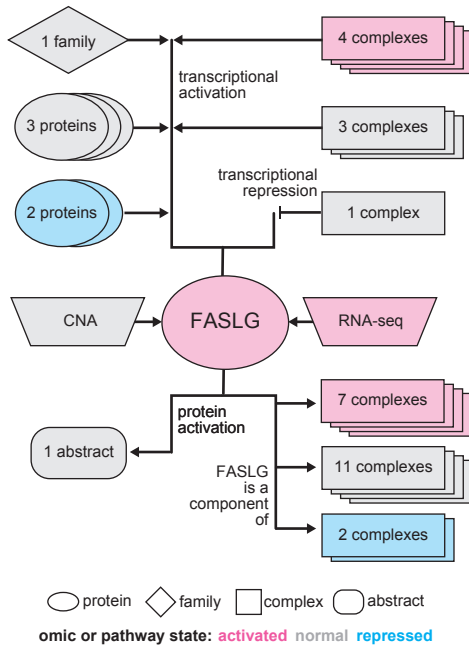


B

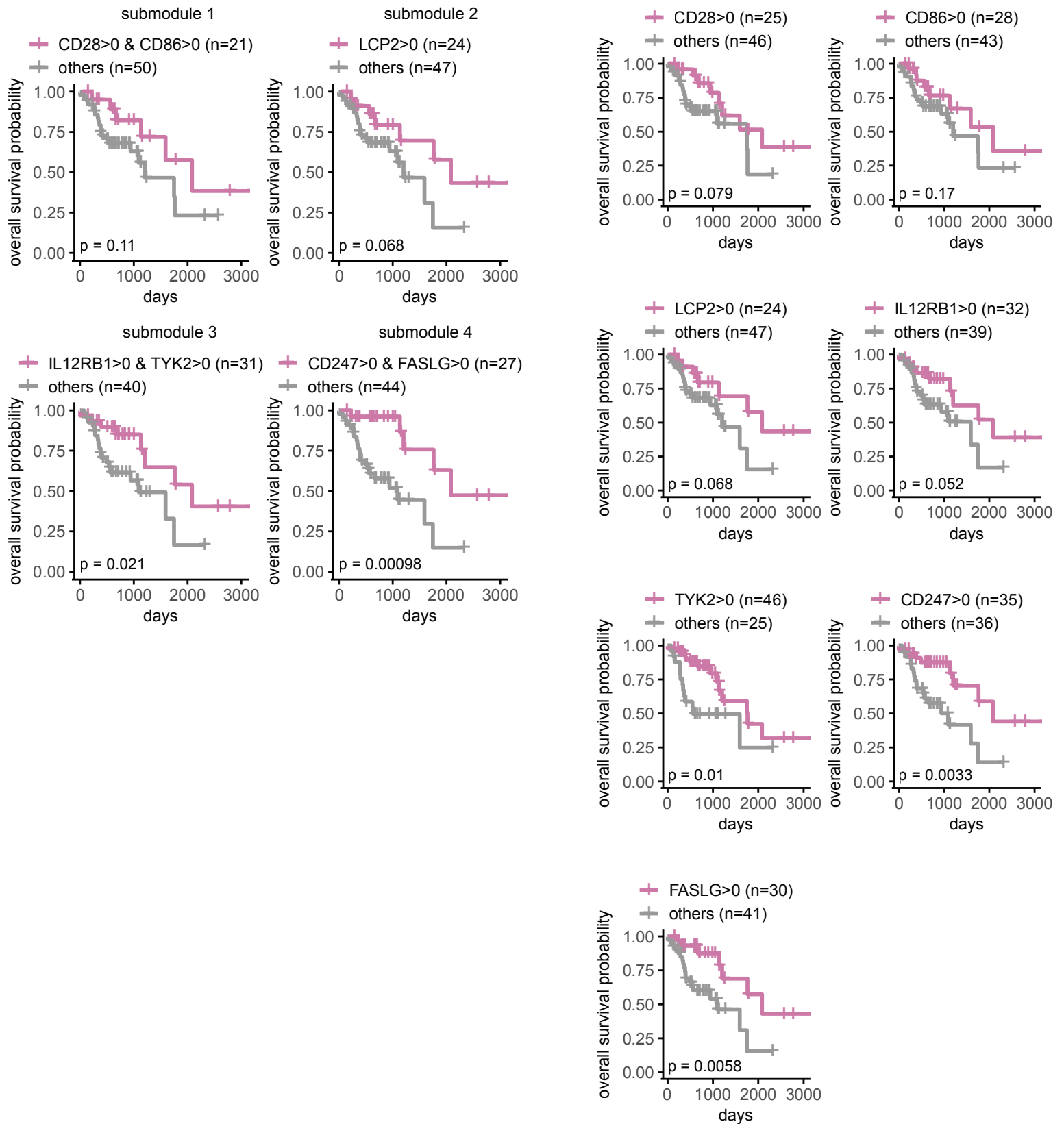
11 group I patient samples in HPV+ exploratory set



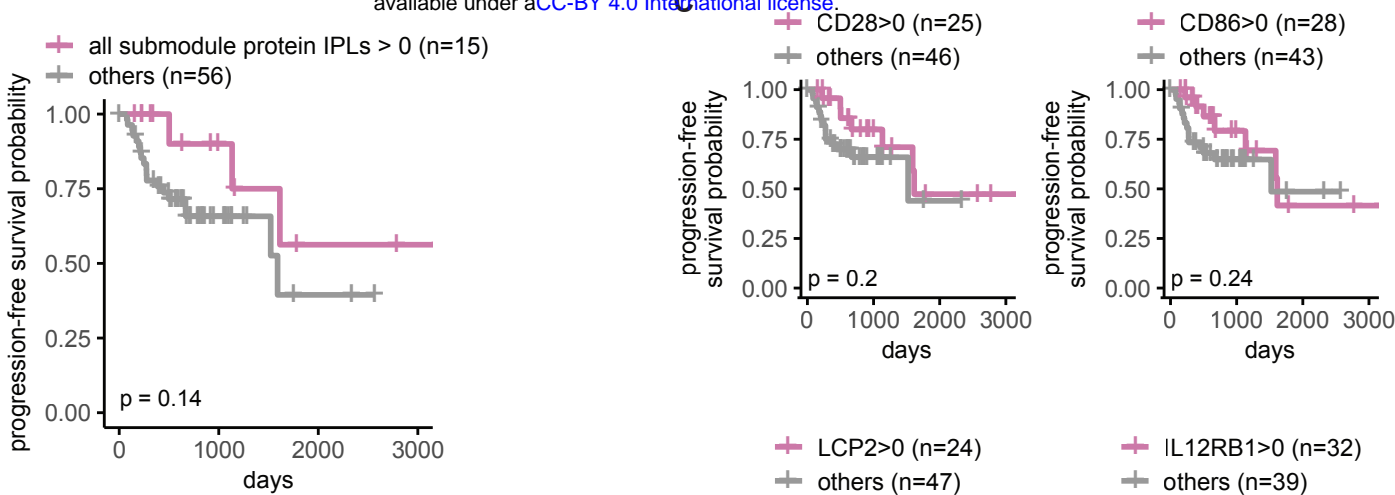
TCGA-CR-7385



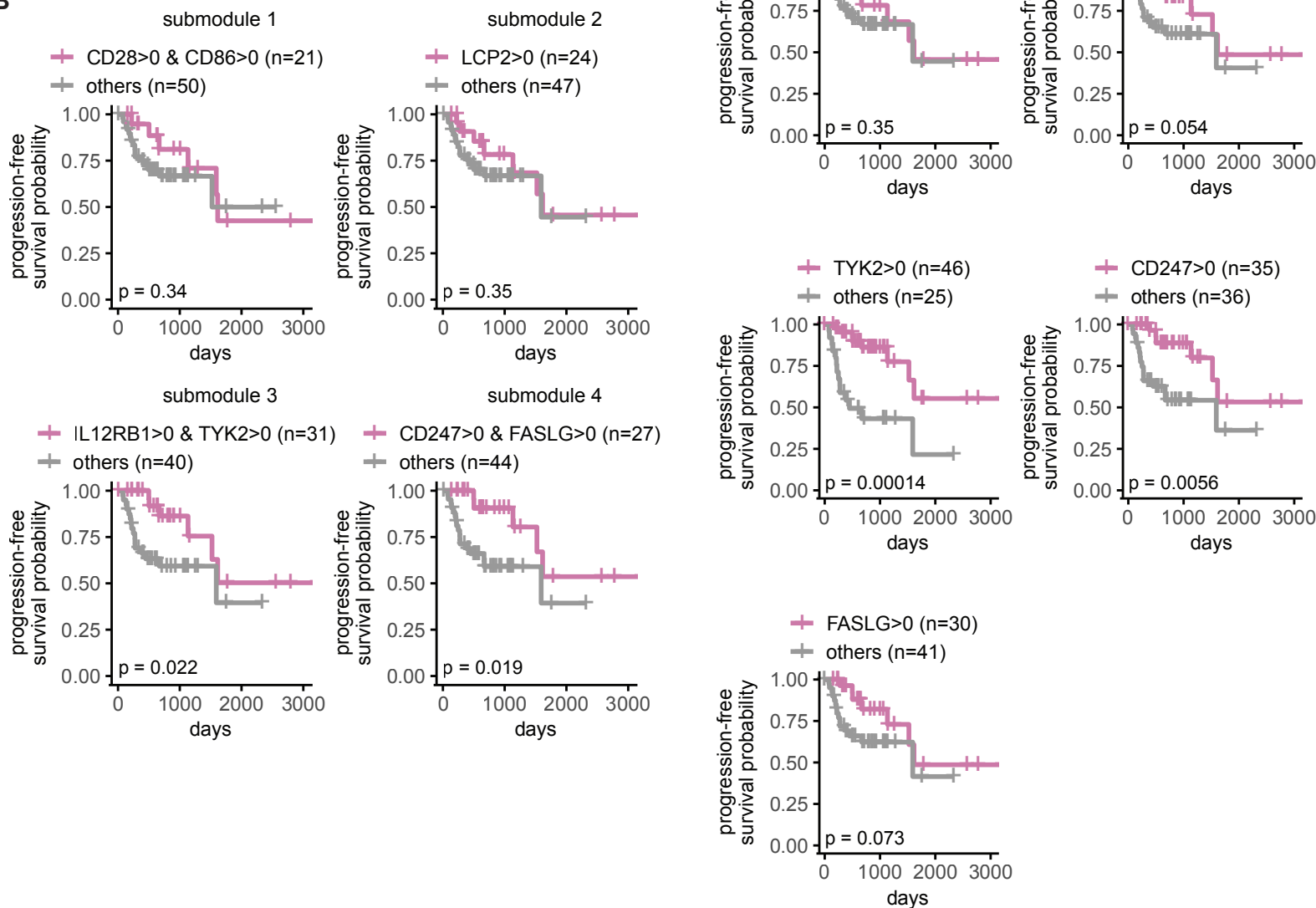
A



A



B

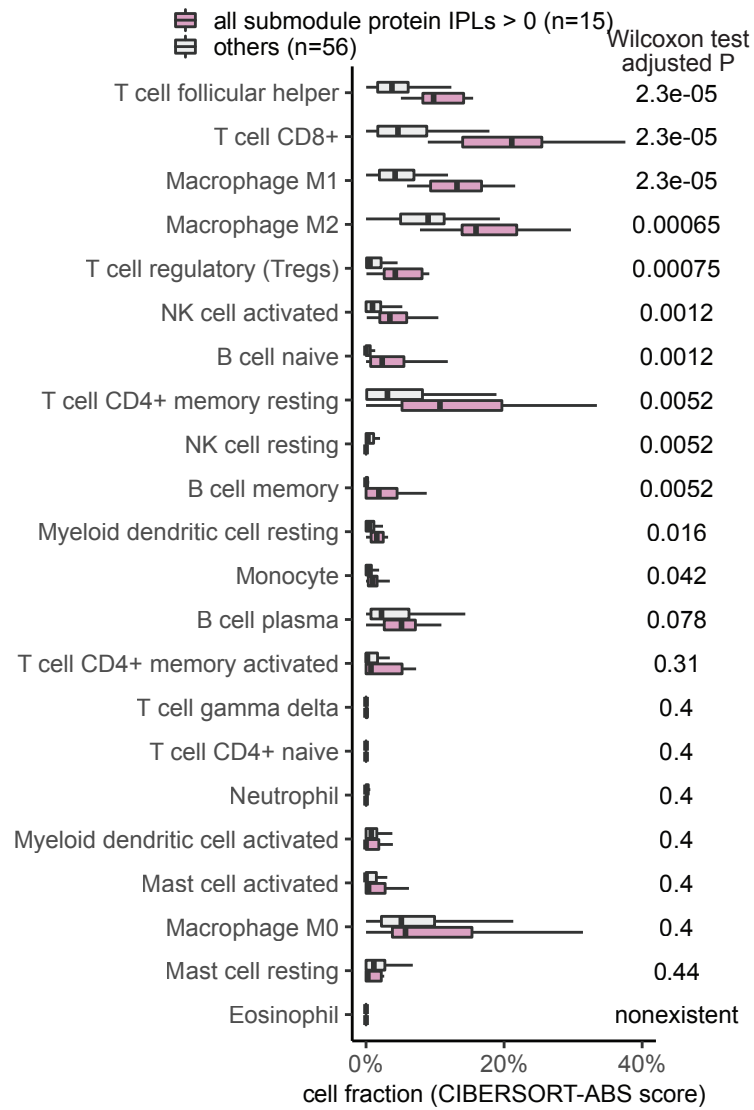


Supplementary Figure 11

bioRxiv preprint doi: <https://doi.org/10.1101/2024.06.15.599113>; this version posted June 17, 2024. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

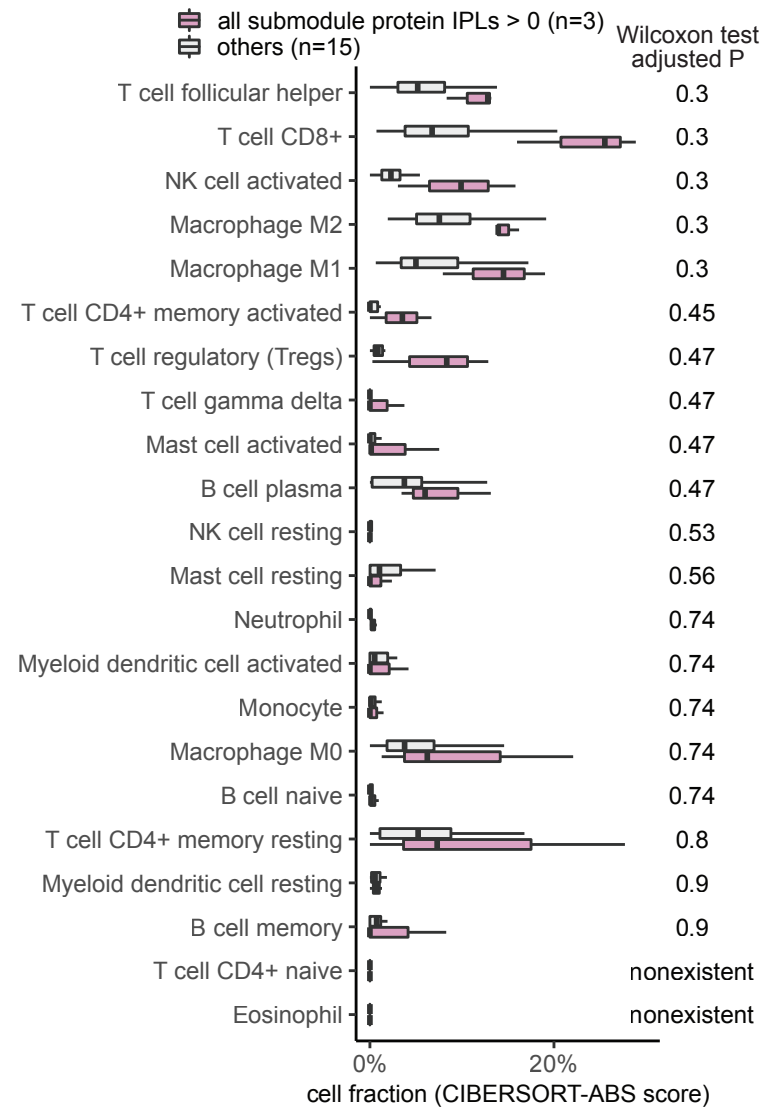
A

71 patient samples in HPV+ exploratory set

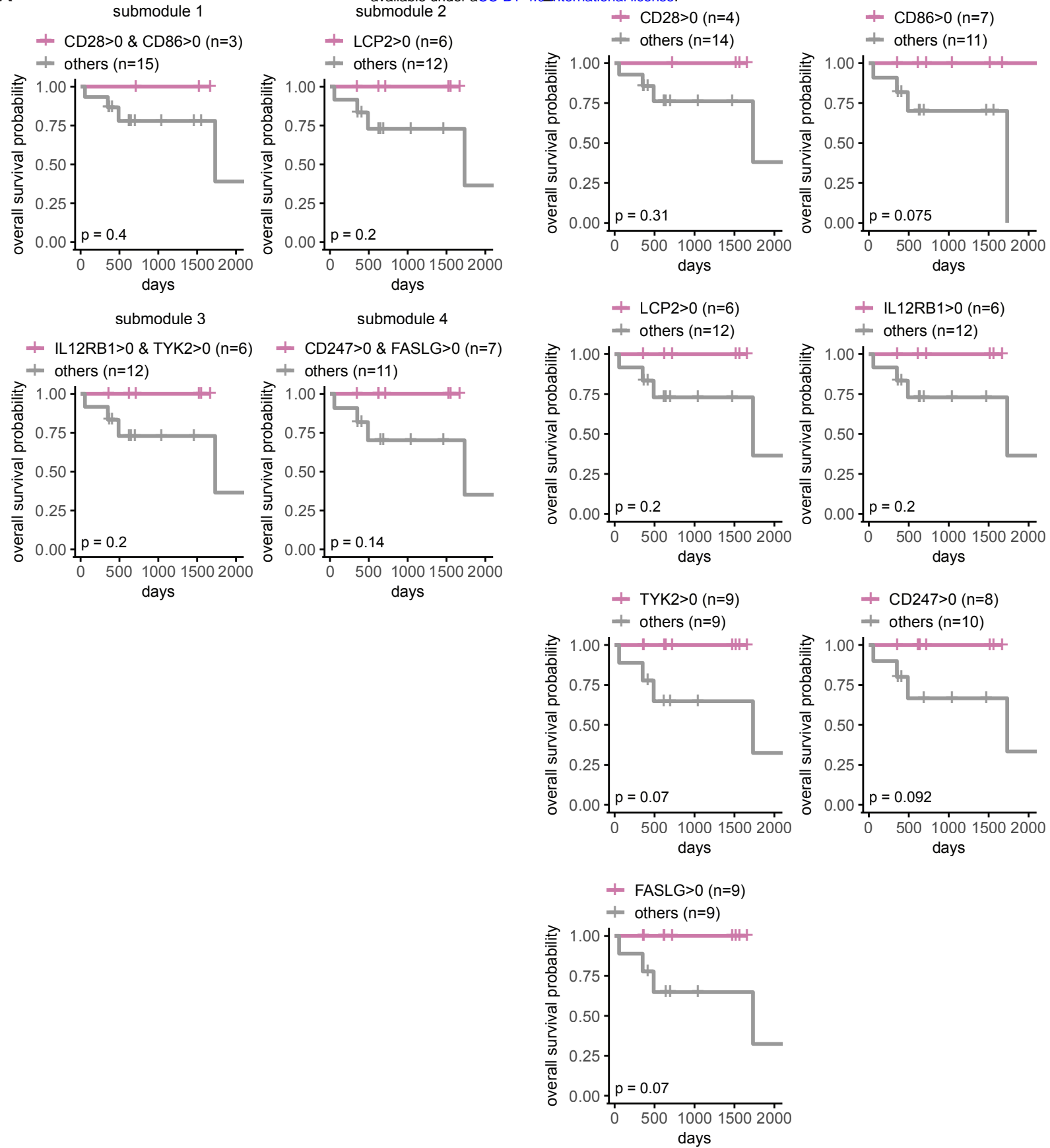


B

18 patient samples in HPV+ validation set



A



B