

SPECTRA: AGNOSTIC EXPRESSION VARIABLES FOR FLEXIBLE TRANSCRIPTOME MODELING IN COMPLEX DISEASE

Rosalie G. Waller*, Heidi A. Hanson*, Michael J. Madsen, Brian Avery, Douglas Sborov, Nicola J. Camp
University of Utah School of Medicine

Abstract. We describe SPECTRA, a novel approach to measure variation in the transcriptome, providing unsupervised quantitative variables to model with any clinical, demographic, or biological endpoint. Complex diseases, including cancer, are highly heterogeneous, and large molecular datasets are increasingly part of describing an individual's unique experience. Gene expression is particularly attractive because it captures both genetic and environmental consequences. SPECTRA provides a framework of agnostic multi-gene linear transformations to calculate variables tuned to the needs of complex disease studies. SPECTRA variables are not supervised to an outcome and are quantitative, linearly uncorrelated variables that retain integrity to the original data and cumulatively explain the majority of the global population variance. Together these variables represent a deep dive into the transcriptome, including both large and small sources of variance. The latter is often overlooked but holds the potential for the identification of smaller groups of individuals with large effects, important for developing precision strategies. Each spectrum is a quantitative variable that can also be considered a phenotypic outcome, providing new avenues to explore disease risk. As a set, SPECTRA variables are ideal for modeling alongside other predictors for any clinical outcome of interest. We demonstrate the flexibility of SPECTRA variables for multiple endpoints, and the potential to out-perform existing methods, using 767 myeloma patients in the CoMMpass study. SPECTRA provides an approach to incorporate deep, transcriptome variability in studies to advance research in precision screening, prevention, intervention, and survival.

1 **Introduction**

2 To identify risk and prognostic factors and understand complex diseases in a population,
3 numerous data types are often collected on study participants. Transcriptomes represent the
4 combined effects of inherited and somatic insults as well as epigenetic and non-genetic factors
5 and thus appeal to researchers interested in both genetic and environmental risk factors. For
6 this reason, gene expression studies are gaining momentum in new fields, such as genomic
7 epidemiology.¹⁻¹³ The need to incorporate transcriptomes in multivariable models alongside
8 other risk factors brings new demands for techniques designed with this in mind.

9 Here we develop the SPECTRA approach to determine multi-gene transformations and
10 calculate transcriptome variables with desirable attributes for multivariable modeling.
11 Specifically, variables that optimize coverage of the global variance, are quantitative,
12 uncorrelated and that retain integrity to the original data. A 'variable' has more power for
13 modeling if it represents the variance in the population. Multiple variables may be required for
14 broad coverage (a 'deep dive'). Knowledge of how much variance transcriptome variables
15 represent is also important to understand the limitations of a study. Furthermore, truly
16 quantitative variables can achieve greater power than if discretized.¹⁴ Uncorrelated variables
17 provide parsimony in penalized modeling and are often simpler to interpret in multivariable
18 analyses. The integrity of variables to the original data (preservation of 'distance' between
19 samples) is important for interpretation. Our goal for such attributes contrasts with the more
20 common strategy to use transcriptome data to categorize samples or patients, reducing the
21 transcriptome data to a single variable consisting of mutually exclusive categories, often called
22 'subtypes'.¹⁵

23 We focus on agnostic derivation. Current techniques for characterizing transcriptomes
24 largely have a computational biology emphasis, interwoven with and constrained by biological
25 knowledge. While these have had great success advancing our understanding of mechanism
26 and pathway,¹⁶⁻²³ there remains room for complementary approaches. Sources of heterogeneity
27 are complex and we require methods that match that complexity. Common diseases, and

1 cancers, in particular, are multifactorial, where a wealth of other covariates may be equally
2 important to an endpoint. New approaches that can embrace this complexity will enhance the
3 toolset available for interrogating transcriptomes. Conceptually, the advantage of an agnostic
4 data-driven approach is the liberty to discover signals that may challenge conventional wisdom
5 or defy “known” rules. Our agnostic approach is complementary and adds to current approaches
6 for the analysis of tumor etiology, risk, treatment, and mortality.

7 Finally, our motivation is for universal measures, and hence our approach is
8 unsupervised. SPECTRA produces a framework of multi-gene transformations to describe the
9 expression space. The calculated SPECTRA variables can be used for different outcomes,
10 providing the flexibility to explore the same variables across several different models (e.g.
11 overall survival, progression-free survival, and time-to-treatment-failure). This can support
12 interpretation and comparisons, improving the ability to decipher the true nature of associations
13 and explore differences. Furthermore, the same framework of transformations can be
14 implemented in external studies, increasing our ability to compare findings across multiple
15 studies.

16 To satisfy these ideals, the core of our approach utilizes principal components analysis
17 (PCA). PCA is an agnostic and unsupervised procedure that provides an isometry to provide a
18 new set of orthogonal (linearly uncorrelated) variables that optimize the representation of the
19 variance. In simple terms, PCA reveals the internal structure of the data in a way that best
20 explains the variance in the data. Paramount in this approach is careful attention to quality
21 control, normalization, and batch correction to ensure the variables capture meaningful
22 variance. The results of the subsequent PCA are the rotation matrix that describes the multi-
23 gene linear transformations; and the transformed data matrix, the quantitative variables for each
24 individual that we refer to as a SPECTRA variables, or simply, spectra. Each measure is a
25 spectrum that combined are spectra. The set of linear transformations provides a new reduced-
26 dimension framework for the expression space. The SPECTRA variables are linearly
27 independent, each providing additional coverage of the variance.

1 We previously used PCA to define a framework for the PAM50, a targeted and
2 standardized gene-panel for breast cancer.^{1,24} Using a population cohort of breast tumors, we
3 used PCA to reduce the 50-gene space to five multi-gene expression variables. When the
4 framework was implemented in an external dataset of tumors from high-risk breast cancer
5 pedigrees, these quantitative PCA variables as phenotypes proved superior to the standard
6 PAM50 subtypes for gene mapping.^{1,25} Further, when implemented in a second external clinical
7 trial dataset, PCA variables were able to predict response to paclitaxel.²⁴ Here, we extend the
8 approach to the whole transcriptome.

9 Improved representation of an individual's tissue (normal or diseased) will be vital to
10 improving our ability to identify expression characteristics that are important phenotypic traits
11 (predict disease risk) and/or important expression variables to predict patient outcomes. **Figure**
12 **1** uses a color analogy to illustrate the conceptual shift of SPECTRA, contrasting our goal of
13 quantitative variables for direct use in outcome modeling with a more conventional
14 categorization approach using hierarchical clustering. In our approach, each spectrum in Figure
15 1 (x_R, x_G, x_B) are independent variables that can be directly used to model any outcome (y_i), and
16 other covariates/predictors can also be easily included (Figure 1d). Conversely, unsupervised
17 hierarchical clustering uses the spectra to categorize patients into groups (Figure 1c), flattening
18 the multiple spectra to a single categorical variable and reducing statistical power. For example,
19 in Figure 1 analogy, x_R cannot be captured by any group ordering and associations for that
20 spectra variable would be lost. An alternate convention is to supervise clustering to an outcome.
21 But, while supervised clustering can improve power over unsupervised clustering for prediction
22 of a single outcome, it also tethers the groups to the particular trained outcome and doesn't
23 facilitate comparison to other outcomes.

24 We illustrate SPECTRA using the Multiple Myeloma Research Foundation (MMRF)
25 CoMMpass Study data.²⁶ We derive the gene transformations (framework) for bulk whole
26 transcriptome RNA sequencing (RNAseq) data from CD138+ myeloma cells. As a proof-of-
27 concept, we utilize the spectra variables in various regression models to identify associations

1 with several outcomes, including established risk scores, patient characteristics, and clinical
2 endpoints.

3

4 **Results**

5 **SPECTRA, a quantitative transcriptome approach**

6 The motivation is the derivation of well-behaved, quantitative variables from RNAseq
7 data to capture transcriptome variation that can be used universally as predictors for any
8 outcome, and as novel phenotypes. The approach requires a dataset to derive the framework of
9 transformations for the SPECTRA variables. Then multiple spectra are calculated for each
10 individual in the dataset. An overview of the SPECTRA approach is shown in **Figure 2**. As an
11 agnostic technique, the goal is to retain only those aspects of the RNAseq data that can
12 represent meaningful variance. Accordingly, rigorous quality control (QC), normalization, and
13 batch correction are performed before the derivation of the variables. Genes likely to lack
14 precision are removed. Only coding genes with sufficient coverage across the dataset are
15 considered. An internal normalization procedure accounts for feature-length, library size, and
16 RNA composition. This normalization avoids the need for reference samples, real or synthetic,
17 and provides the potential for spectra to be ported to follow-up samples and external datasets.
18 Finally, skew and outliers are dealt with before PCA is performed. Specific details are listed in
19 the Methods.

20 PCA is a well-established, data-driven method that, based on the covariance of a
21 dataset, produces a matrix factorization which is a unique solution of linear transformations
22 (*framework*, rotation matrix) and transformed values (*spectra*, transformed data). The linear
23 transformations preserve the variance in the data, i.e., the transformed values preserve the
24 distance between the sample data for individuals. Integrity to the original data provides
25 meaningful comparisons between individuals. The resulting transformed values are quantitative
26 variables that are orthogonal (linearly uncorrelated). For dimension reduction, components are
27 ordered according to the amount of global variance they explain and the first k (S_1, \dots, S_k)

1 selected, for which the proportion of total variance explained can be described. This reduces
2 attention from 60,000+ expression features in a transcriptome to a handful of spectra
3 specifically derived to represent independent components of the natural global variation across
4 the dataset studied; precisely the type of variables with power to identify differences and
5 important for prediction. The procedure is unsupervised, describing only intrinsic variance in the
6 data, hence the spectra can be incorporated into modeling any outcome in an unbiased way,
7 and the framework of transformations can be implemented in external datasets.

8 When modeling, the significance and effect size for each individual spectrum as a
9 predictor can be determined. In addition, the aggregate effect of all spectra in the model can be
10 determined to describe the impact of the transcriptome as a whole. We define the aggregate
11 effect of all spectra in a model as the poly-spectra liability (**PSL**) score for the outcome. This is
12 the weighted sum of the spectra based on the model.

13

14 **Illustrative case study: CD138+ spectra in multiple myeloma**

15 The ultimate value of SPECTRA will be its use in the discovery of novel tissue
16 phenotypes and predictors or outcomes in etiological studies. Here, as a proof-of-concept that a
17 SPECTRA framework can capture meaningful information, we present associations between a
18 set of spectra to several well-established outcomes or risk groups for multiple myeloma, across
19 several different model types. These are not presented to suggest spectra could replace current
20 clinical tests, but to illustrate the flexibility of SPECTRA to provide a universal transcriptome
21 framework and set of variables for use in various models with disparate outcomes. We applied
22 SPECTRA to transcriptome data for CD138+ cells from the MMRF CoMMpass study.²⁶ We
23 investigated associations of CD138+ spectra with 1) existing expression-based risk scores;^{27,28}
24 2) clinically-relevant DNA aberrations; 3) clinical prognostic outcomes, and 4) patient
25 demographic groups with elevated myeloma risk. Also, we illustrate the potential to track
26 transcriptome changes over time.

1 The CoMMpass dataset is the most extensive sequencing effort in multiple myeloma
2 patients to date. Multiple myeloma is a malignancy of plasma cells (CD138+ cells). The publicly
3 available transcriptome data (IA14) comprised RNAseq data for 887 CD138+ samples on 794
4 unique patients. Here, data for 768 patients with treatment naïve samples collected at diagnosis
5 were the focus. We used transcript-based expression estimates from Salmon,²⁹ generated by
6 the CoMMpass study (<https://research.themmf.org>). From the total 54,324 features, 7,436
7 genes and 767 patients' treatment-naïve CD138+ RNAseq data met quality control. The
8 transcriptome framework and spectra were derived in a quality controlled, normalized and batch
9 corrected data from these treatment naïve samples. The first $k = 39$ spectra (S1—S39) were
10 selected, based on a scree test, which captured $v = 65\%$ of the global variation. No spectra
11 showed association with batch (F-statistic, all $p > 0.8$). The details from each step of the QC
12 process, the linear transformations necessary to calculate the 39 spectra, and the individual-
13 level spectra variables for the patients in the IA14 CoMMpass data are provided in
14 Supplemental Data. R markdown notebooks containing the code used to generate CD138+
15 spectra in the IA14 dataset, full model analyses, and results are provided in the Supplemental
16 Materials.

17 As linearly uncorrelated variables, each of the 39 CD138+ spectra captures a different
18 source of variance, and hence any spectrum has the potential to explain patient differences and
19 provide insight. **Figure 3** shows spectra charts for 4 patients and illustrates that while patients
20 may be similar at a high-level (overall patterning), that individual spectra may not follow that
21 apparent high-level similarity.

22 Common approaches to prediction modeling include penalized or stepwise techniques to
23 address concerns about multicollinearity and improve fit and parsimony. Here, we included all
24 39 spectra into each model for simplicity and to ease comparison across results. Association
25 results for the full 39-spectra models for several different outcomes are described below. Overall
26 model significance and the significance for individual spectra in those models are summarized
27 in **Figure 4**.

1
2 *CD138+ spectra and established expression-based risk scores.* We illustrate that spectra
3 capture expression variability fundamental to previously established supervised risk scores. We
4 compare to the state-of-the-art risk scores: 1) the most widely adopted and first supervised
5 expression risk score in myeloma, from the University of Arkansas for Medical Sciences
6 (UAMS-70);²⁷ and 2) one of the most recent supervised risk scores, from the Shahid Bahonar
7 University of Kerman (SBUK-17).²⁸ Both risk scores were derived from the entire transcriptome,
8 reducing dimensionality by restricting to selected genes, and providing a multi-gene risk metric
9 based on those genes.

10 The UAMS-70 risk score was developed in microarray data and tested 54,657 probes for
11 association with disease-related survival.²⁷ A total of 70 genes were selected (19 under and 51
12 over expressed prognostic genes). The UAMS-70 risk score is the ratio of mean expression of
13 the up-regulated to down-regulated genes, and k-means clustering was used to determine a
14 cutoff for 'high-risk' classification.²⁷ Using the established classifier, we calculated each
15 CoMMpass patient's UAMS-70 risk score and their risk status (low or high).

16 The SBUK-17 prognosis score was developed in RNA sequencing data.²⁸ All genes in
17 the entire transcriptome were tested for association with survival.²⁸ A multi-step process,
18 including univariate Cox analysis, intersection in six Class Prediction algorithms, and
19 multivariate Cox analysis, was used to select 17 genes consistent across multiple methods.²⁸
20 These were entered in a multivariable Cox regression to define the SBUK-17 score.²⁸ The 75th
21 percentile was used to classify patients into a high-risk and low-risk categories.²⁸ Using the
22 established classifier, we calculated each CoMMpass patient's SBUK-17 prognostic score and
23 risk status (low or high).

24 To illustrate the ability for unsupervised spectra to maintain variation fundamental to the
25 UAMS-70 and SBUK-17 scores, we considered each score as a quantitative outcome and
26 performed linear regression to predict them using spectra. The UAMS-70 score could be
27 predicted with high accuracy and extreme significance ($R^2 = 0.86$, $F_{39,727} = 118.1$, $p < 10^{-50}$).

1 Thirty spectra were individually significant ($p < 0.05$, **Figure 4**) in the model. The SBUK-17 score
2 could also be predicted with excellent accuracy and significance ($R^2 = 0.93$, $F_{39,272} = 252.9$, $p <$
3 10^{-50}). Twenty-five spectra were individually significant (**Figure 4**). **Figure 5** illustrates the high
4 correlation between the spectra predictions (model PSL scores) and the previously established
5 risk scores. The CD138+ spectra can recapitulate previously established supervised expression
6 risk scores, indicating the spectra framework can capture previously identified prognostic
7 signals.

8

9 CD138+ spectra and disease course. Without constraints to previous risk scores, we used Cox
10 proportional hazards analysis to predict overall survival (**OS**) using spectra. Spectra significantly
11 predicted OS (179 events, likelihood ratio test, $p = 3.1 \times 10^{-17}$, C-statistic = 0.74), with 12 spectra
12 individually significant. To view using Kaplan Meier curves, patients were split into three equal
13 tertiles based on PSL scores for OS (**Figure 7**). For OS, patients in OS-PLS tertile 3 had hazard
14 ratios (HR) and 95% confidence intervals (CI) of 6.7 (2.9-15.3) and 8.8 (5.1-15.3) at 1 year and
15 3 years, respectively, compared to patients in tertile 1.

16 To assess ability of spectra-based modeling to differentiate patients with contrasting
17 survival trajectories we compared to UAMS-70 and SBUK-17 scores based on their
18 classification to high and low risk groups. For this comparison, we combined PSL tertiles 1 and
19 2 vs tertile 3. The CD138+ PSL provided larger and more significant HRs than UAMS-70 and
20 SBUK-17 for OS (**Table 1**). Similar trends in HRs were observed between the PSL and UAMS-
21 70 risk score, but PSL was more significant at 1, 3, and 5-years. SBUK-17 did not perform as
22 well as either PSL or UAMS-70 and did not appear well suited to predict survival beyond 3
23 years. SPECTRA was not supervised on OS, but outperforms existing state-of-the-art
24 expression scores in predicting OS.

25 Using the same 39-spectra framework, we used multivariable Cox proportional hazards
26 analysis to predict time to first-line treatment failure (**TTF**). Spectra were also significant for
27 predicting TTF (369 events, likelihood ratio test, $p = 7.9 \times 10^{-10}$, C-statistic = 0.66), with 8 spectra

1 individually significant. We note that spectra significant in both OS and TTF models had effects
2 in the same direction (**Figure 4**). As above, patients were categorized into three equal tertiles
3 based on PSL scores for TTF. Kaplan Meier curves for these three equal groups for TTF are
4 shown in **Figure 7**. Comparing TTF-PLS tertile 3 to tertile 1, HR and 95% CIs were 4.8 (3-7.7)
5 and 7.2 (4.3-12) at 1 year and 3 years, respectively. These results indicate spectra can capture
6 signals and differentiate patients for disease course.

7
8 *CD138+ spectra and clinical risk.* Large somatic chromosomal DNA aberrations detected by
9 cytogenetics are used clinically to define prognostic risk groups in myeloma.³⁰ Clinical risk
10 categories defined by mSMART³¹ include: high risk (del(17p) and t(14;16)); intermediate risk
11 (amp(1q) and t(4;14)); and standard risk (t(11;14)). Models for each of these five chromosomal
12 aberrations (**Figure 4**) showed different spectra individually significant, with some spectra
13 unique to only one aberration. Interestingly, while the models for all three translocations and
14 amp(1q) were highly significant (all $p < 2 \times 10^{-10}$), the full 39-spectra model for del(17p) was not.
15 To investigate the possibility that the model was over-parameterized, we repeated the del(17p)
16 analysis using a stepwise procedure. This produced a significant model containing only 3
17 spectra ($p = 0.014$, Supplemental Material). These results indicate transcriptome spectra
18 capture signals from DNA chromosomal changes in CD138+ cells (**Figure 6a-b**).

19 The international staging system (**ISS**) for myeloma is also used to classify and stratify
20 patients at diagnosis, based on somatic cytogenetics, levels of beta-2 microglobulin, albumin,
21 and lactate dehydrogenase in the blood.³² In an ordinal logistic regression with the ISS stage at
22 diagnosis as the outcome, 13 spectra were significant, providing a model that significantly
23 differentiated the three clinical stages (**Figure 6c**). These results indicate spectra can capture
24 signals for the disease stage.

25
26 *CD138+ spectra and demographic risk groups.* Myeloma is an adult-onset malignancy, most
27 frequently diagnosed at ages 65-74 years (median 69 years).³³ Incidence is higher in men (8.7

1 men vs. 5.6 women per 100,000) and patients self-reporting as African American (AA men 16.3,
2 and AA women 11.9 per 100,000). Linear regression with age at diagnosis as a quantitative
3 outcome was significant ($p = 2 \times 10^{-14}$), with 15 individually-significant spectra (**Figure 6d**).
4 Logistic regression models for gender, race (self-reported black or white; other racial categories
5 too small to consider) and Hispanic status were all significant ($p = 4 \times 10^{-9}$, $p = 9 \times 10^{-10}$ and $p =$
6 1×10^{-3} , respectively) (**Figure 4**). The state-of-the-art study whose goal was to identify
7 differences by race for myeloma tumors did not have a variable framework for the transcriptome
8 with which to perform comparisons. Instead, they used the UAMS-70 score to compare
9 transcriptomes by race and did not identify significant differences ($p = 0.662$).³⁹ In contrast, our
10 CD138+ spectra model identified 13 spectra that differed significantly by race, and a multi-
11 spectra model (PSL) that highlights significant differences (**Figure 6e**). We note that
12 associations found for demographic risk factors may be complex, as such factors involve social
13 constructs, e.g. race and ethnicity. Transcriptomes can harbor the effects of genetic, epigenetic,
14 lifestyle, and environmental factors. These results indicate spectra can capture signals
15 originating from demographic risk.

16

17 CD138+ spectra for tracking changes over time. The SPECTRA framework provides the
18 transformations for the spectra variables such that they can be calculated in follow-up samples
19 and tracked over time. We illustrate this potential in the eleven MM patients for whom at least
20 three longitudinal CD138+ samples were available in the CoMMpass study. **Figure 8** shows a
21 line graph of the PSL score for OS for these eleven patients over 80 months. In this example,
22 the potential for tracking a patient's hazard over time is illustrated using the OS PSL score.

23

24 **Discussion**

25 The promise of personalized prevention, management, and treatment is rooted in an
26 ability to describe an individual's unique experience and model important sources of
27 heterogeneity.³⁴ In complex diseases, and cancer specifically, gene expression in diseased

1 tissue may be an established source of heterogeneity.³⁵ Tools that can take a deeper dive and
2 characterize multiple sources of expression heterogeneity will be important to advance the
3 promise of personalized medicine. In particular, for human studies and domains such as
4 epidemiology wishing to model multiple sources of risk in a population, transcriptome variables
5 that can be easily incorporated with other variables are needed. The goal of this study was to
6 provide a technique to derive an agnostic framework of variables for transcriptome data, to
7 empower multivariable studies, and provide novel molecular phenotypes. SPECTRA identifies
8 quantitative, orthogonal variables (non-correlated) that capture sources of transcriptome
9 variation for use in subsequent modeling or as quantitative phenotypes. Many applications can
10 benefit from the qualities of spectra variables, and this new framework has the potential to
11 provide utility to numerous study designs and many outcome types.

12 Data quality and processing are paramount in the quest to derive informative variables.
13 PCA itself is a simple procedure that provides linear transformations of the data to best
14 represent variance. If the data have technical artifacts, batch effects, unstable or non-
15 comparable expression measures, the noise will overwhelm authentic variance. Accordingly, our
16 technique intentionally includes strict quality control, zero-handling and normalization
17 procedures, and batch correction (Figure 2). Without these steps, PCA can fail to provide
18 variables with the desired qualities. An agnostic approach permits stringent data culling because
19 the incentive to retain features based on known functional relevance is removed. The impetus is
20 to only retain features that can contribute to meaningful variance and provide informative
21 variables for modeling (quantitative, orthogonal, variance-representing). Of course, the limitation
22 of an agnostic approach is reduced biological interpretation or insight into the mechanism of the
23 variables before modeling. However, there are already many approaches that take this alternate
24 goal of intermediate interpretation,¹⁸ whose limitations are instead the flexibility of the variables
25 they produce. Hence, SPECTRA offers a complementary approach to the current toolset
26 available for all fields.

1 Beyond the agnosticism taken by our proposed technique, other potential advantages of
2 SPECTRA include its unique solution within a dataset, such that the rank of the dimension
3 reduction can be post-hoc and does not influence the definition of retained dimensions. As a
4 statistically rigorous technique, it also provides a measured dive into the transcriptome. Each
5 dimension (eigenvector q_s) iteratively moves quantifiably deeper into the variance of the data
6 (measurable by λ_s). Methods that iteratively find independent components (PCA and
7 independent component analysis) have previously been shown to provide superior coverage of
8 transcriptome data.²³ Retention of components deep in the data, representing small variances
9 (i.e., deep dives) provide potential and power to identify small groups of individuals with large
10 effects in outcome studies, such as a molecular phenotype that hones-in on a rare Mendelian
11 form of cancer, or the few patients that respond to a drug. These findings could be the 'low
12 hanging fruit' scenarios where the precision translation is more straight-forward. SPECTRA also
13 embraces negative weights. The allowance of negative values in its matrix factorization (**MF**) is
14 often given as a criticism of PCA,^{16,36} argued as a conceptual source of its lack of biological
15 interpretability; a premise that components may mix biological processes due to a focus on
16 variance. Non-negative matrix factorization (**NMF**), arguably the leading approach in the
17 computational biology field, restricts all values in the amplitude (equivalent to Q^T in PCA) and
18 pattern (equivalent to T^T in PCA) matrices to be non-negative. Reasoned as beneficial and a
19 natural restriction because expression values themselves cannot be negative. Also, because
20 NMF transformed values represent the proportions of each factor and thus provide a simple
21 interpretation. However, non-negative values may not be a natural restriction to *systems of*
22 *genes*, and over-simplicity may not adequately represent true complexity. With a non-negativity
23 restriction, NMF limits itself to the identification of groups of over-expressed genes,^{16,36} modeling
24 only neutrality and surplus. Deficits may also be important. So, while PCA spectra may
25 represent mixtures of different biological mechanisms, these may be important combinations,
26 including genes acting in opposite directions, and may better reflect reality. By embracing
27 negative values, PCA can also capture gene systems in deficit, which may be more difficult to

1 interpret, but may equally be just as important to recognize. These differences underscore
2 SPECTRA's value as a complementary tool to existing approaches.

3 Our myeloma case study illustrated derivation of a transcriptome framework and spectra
4 variables for CD138+ cells, and the application of these in various models (linear, logistic and
5 Cox regression, and ANOVA) with many different outcomes. We showed that the set of 39
6 unsupervised, agnostic spectra could significantly capture signals corresponding to published
7 expression-based risk scores from traditional supervised approaches, known clinical DNA-
8 based risk factors, disease stage, disease progression, survival, and demographic risk groups.
9 We also illustrated the potential to track tissue changes using PSL scores over time. As
10 expected for a framework of agnostically derived variables, not all spectra are relevant to every
11 outcome. Across the 14 models presented, the number of individually significant spectra in a
12 model ranged from 3 to 30, and only one spectra-variable (S30) showed no association with any
13 model. Importantly, these examples show the flexibility of the framework as well as how it can
14 support comparisons across different models and outcomes. For example, our results illustrated
15 how two previously established gene-expression prognostic scores could be captured using a
16 single framework and illustrate that they are similar (**Figure 4**). Hence, the framework has the
17 potential to provide a bridge to compare various existing categorizations (subtypes) of patients,
18 even when no genes overlap in their signatures, or they predict different outcomes.³⁷ In this
19 way, spectra provide an alternate to categorical intrinsic subtyping, a well-established practice
20 for many cancers.³⁸ The ability to predict OS and TTF suggests spectra hold utility in clinical
21 studies predicting disease course. Clinically-relevant stratification may be better represented
22 using thresholds within a transcriptome framework.²⁴

23 The potential for increased power using spectra variables is illustrated by the discovery
24 of novel associations between spectra and patient demographic risk groups with known
25 differences in incidence (age, gender, race). Prior studies, using the UAMS 70-gene panel and
26 a Ki67 proliferation index, were not able to identify gene expression differences in CD138+ cells
27 from self-reported AA and white patients.³⁹ Our multivariable results demonstrate that significant

1 differences do exist, but also illustrate that the diseased cells in these demographic groups are
2 not distinct entities; fewer than half the spectra variables differ significantly by these patient
3 demographic groups. Focusing on the spectra that do show differences by demographics
4 provides new avenues to explore why incidence varies in these groups; a key to disease
5 prevention, intervention, and control. In particular, because transcriptomes capture both the
6 effects of internal (inherited genetics) and external factors (lifestyle, exposures, consequences
7 of access to care), these results could also support epidemiology and biosociology
8 investigations into such differences. We provide the variable framework (gene transformations)
9 and the spectra variables for the CoMMpass patients in Supplementary material to enable
10 further study of spectra in other CoMMpass studies, as well as in other myeloma studies.

11 There are numerous potential applications beyond those undertaken here that could
12 benefit from a statistically rigorous transcriptome framework of expression variables. As shown
13 previously for the PAM50 panel in breast tumors, differences can be observed between familial
14 and sporadic tissues, suggesting familial components,¹ and defining powerful new phenotypes
15 for genetic, exposure, and gene-environment studies. Future avenues for spectra as
16 quantitative phenotypes could include expression-quantitative trait locus analyses, Mendelian
17 randomization, seeding machine-learning applications,⁴⁰ tissue measures for pre-clinical
18 models, and corollary studies in clinical trials.

19 As for any approach, there are limitations. A key question is one of representation. For
20 epidemiology studies, for example, spectra should ideally be representative of the entire
21 disease population. This requires that the derivation dataset is a random sample from that
22 population, or based on a known selective sampling scheme. While there are many publicly
23 available transcriptome datasets,^{41,42} most fall short of this ideal. Thus, the spectra variables
24 derived from these will have inherent limitations in representation. An investigator should
25 consider if a derivation dataset is adequate to represent their study goals. We note that the goal
26 of the MMRF CoMMpass study was intentionally designed to represent myeloma patients from
27 diagnosis through treatment, and is the largest existing cohort of treatment-naïve CD138+

1 transcriptomes, with sampling continuing over time. However, the demographic representation
2 of patients was not achieved, and this remains a limitation of that study. Another limitation is
3 that, as a simple variance-based procedure, PCA models all sources of variance in the dataset.
4 If artifacts remain in the data, the resulting spectra will also represent these. To minimize this
5 issue, we employed a strict data quality and batch correction process in our workflow,
6 concentrating only on a subset of genes for which PCA is likely to be meaningful: well-mapped,
7 stable, with sufficient depth, and with batch correction. We also removed genes known to be
8 unstable across different RNAseq pipelines.⁴³ A third limitation is the ability to use the
9 framework of spectra in external studies. As a data-driven technique, the complete PCA
10 decomposition is overfitted to the derivation dataset. To limit this, we use dimension reduction
11 and focus on the first k spectra (largest k components of variation), selected using a scree test⁴⁴
12 to be those before decreasing marginal returns. Last, SPECTRA is intentionally agnostic,
13 designed for modeling, and dimensions are not pre-interpreted for functional relevance. Hence,
14 post-hoc analyses will be required to uncover the mechanism/s that underlie the associations
15 identified.

16 In conclusion, we present a new technique, SPECTRA, to derive an agnostic
17 transcriptome framework of quantitative, orthogonal variables for a dataset. These multi-gene
18 expression variables are designed specifically to capture transcriptome variation, providing new
19 transcriptome phenotypes and variables for flexible modeling, along with other covariates, to
20 better differentiate individuals for any outcome. Applied to CD138+ transcriptomes for myeloma
21 patients, we defined CD138+ spectra and implemented these in many different outcome
22 models. We illustrated an ability to predict prognosis, survival, clinical risk, and provide new
23 insight into potential differences between patients from demographic groups. Fundamentally,
24 the technique shifts from categorization to multiple quantitative measures. SPECTRA variables
25 provide a new paradigm and toolset for exploring transcriptomes that hold promise for
26 discoveries to advance precision screening, prevention, intervention, and survival studies.

27

1 **Methods**

2 **Conceptual construction**

3 Here we establish the matrix factorization (**MF**) natural for individual-based outcome
4 modeling. Data matrices, X and T , are oriented with individuals as subjects (n rows) and genes
5 as variables (g columns). Given a $n \times g$ design matrix, X (mean-centered expression values for
6 n individuals on g genes), PCA is the MF

7

$$8 \quad X = TQ^T \quad \text{Equation 1}$$

9

10 where T contains the transformed values (the dimension variables), and Q is the PCA 'rotation'
11 matrix. Each row in $Q^T = (q_1, q_2, \dots, q_g)^T$ is an orthogonal eigenvector (or component) which
12 holds the coefficients for the linear model to transform the observed gene values into the
13 spectra variables. The set of linear transformations are the transcriptome framework. The
14 rotation matrix can be derived from the eigen decomposition of the covariance matrix, Σ

15

$$16 \quad \Sigma = Q\Lambda Q^T \quad \text{Equation 2}$$

17

18 where Σ is proportional to $X^T X$, and Λ is the diagonal matrix of eigenvalues. Each eigenvalue,
19 λ_s , is a scalar indicating the proportion of the global variance represented by the transformed
20 value defined by the s^{th} eigenvector, q_s , in Q . Eigenvalues are ranked, such that the first PC,
21 defined by q_1 captures the most variance, q_2 the next highest, and so on. We note that there
22 can only be $\min(n, g)$ non-zero eigenvalues, because by definition, beyond this no variance
23 remains. In most, if not all, existing RNAseq studies, there are more genes than individuals and
24 hence n is the limiting rank.

25 Dimensionality can be reduced to k dimensions by utilizing Q_k ; only the first k columns
26 (PCs) of Q . After selection of k PCs, transformed values are represented as:

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27

$$T_k = XQ_k \quad \text{Equation 3}$$

We note that PCA is deterministic and therefore the selection of k is a post-procedure decision that does not influence the MF. The proportion of variance explained by the retained dimensions ($\sum_{s=1}^k \lambda_s / \sum_{v \in S} \lambda_s$) can be used as a measure of coverage.

SPECTRA workflow

Careful attention to quality control, normalization, and batch correction are used to ensure the spectra capture meaningful variation. Gene expression counts from bulk RNAseq are the input data. The four steps in the workflow are (1) quality control; (2) internal normalization; (3) correction for batch effects; (4) PCA and dimension reduction (**Figure 2**).

Quality control. QC is essential to ensure the transcriptome dimensions capture meaningful variation across the individuals. Features in the transcriptome likely to be unduly influenced by poor alignment or lacking precision due to sequencing depth were removed as potentials for introducing spurious and unstable variation. Accordingly, we removed all non-autosomal and non-protein-coding genes as well as features with low counts. A feature was considered to have inadequate data for precision if more than 5% of samples had fewer than 100 read counts. After the removal of features, individuals were removed from consideration if more than 10% of the remaining features had fewer than 100 read counts.

Normalization. This is required for comparisons across genes and individuals and includes adjustment for gene length, sequencing depth (library size), and RNA composition. Zero-handling is also necessary to appropriately incorporate counts of zero for a feature or transcript). We chose to use a robust internal (single sample) normalization to obviate the need for a 'reference' sample and to provide the possibility for portability across datasets. While our

1 technique is gene-focused, our processing is designed to handle transcript-based alignment and
2 quantification because these have been suggested to be more accurate.⁴⁵ Normalized gene
3 expression estimates, e_g , were calculated according to the following procedure:

4

$$5 \quad e_g = \log_2 \left(\frac{\sum_{t=1}^m \frac{c_t + 1/m}{l_t}}{\text{median} \left(\sum_{t=1}^m \frac{c_t + 1/m}{l_t} \right)} \right) \quad \text{Equation 4}$$

6

7 where c_t is the read count for transcript t , l_t is the transcript length in kilobases (extracted from
8 the GTF used to align and quantify the RNAseq data), and m is the number of transcripts for the
9 gene. Zero-handling is achieved by adding $1/m$ to the transcript counts: $c_t + 1/m$. Division by l_t
10 corrects for transcript length. Summing the length-corrected transcript counts results in a gene-
11 level count per kilobase (CPK) measure. **Equation 4** may also be used for gene-level read
12 counts (equivalent to $m = 1$). Adjustments for sequencing depth and RNA composition (often
13 referred to as the *size factor*) is achieved via division of each gene-based CPK measure by the
14 median of CPK-values for retained features. We note that the more usual upper-quartile
15 adjustment also provides robust internal normalization,⁴⁶ however, since our implementation is
16 post-QC after numerous features have been removed for low counts, the median is more
17 suitable. Normalized data are \log_2 transformed to account for skew. We also truncate outliers
18 beyond the five standard deviation thresholds from the mean of the normalized gene counts to
19 the relevant threshold value.

20

21 Batch correction. Sequencing is often generated in batches, and it is necessary to correct for
22 the potential of technical artifacts and any spurious variation introduced. We adjust for sequence
23 batch using ComBat⁴⁷ as implemented in the sva R package,⁴⁸ with patient characteristics that
24 are unbalanced by batch included as covariates.

25

1 PCA. We implement PCA with the covariance matrix. For functions that use singular value
2 decomposition to perform PCA, it is necessary to center the expression values first to ensure
3 the MF is performed for the covariance. Expression values (e_g) are centered on the mean
4 across individuals for gene g . These centered data represent the design matrix, X ($n \times g$)
5 (**Equation 1**) for the PCA. The R core function $prcomp(x = X, center = TRUE, scale =$
6 $FALSE, retx = TRUE)$ was used to perform PCA. We use a scree test⁴⁴ (the inflection point of
7 the rank-ordered plot of λ_s , or elbow method) to select the k spectra to retain. The proportion of
8 variance explained by this k -dimensional space ($\sum_{s=1}^k \lambda_s / \sum_{v_s} \lambda_s$) indicates the depth of the dive
9 into the transcriptome data.

10

11 **CD138+ spectra in myeloma**

12 Data were generated as part of the MMRF CoMMpass Study (release IA14)²⁶ and
13 downloaded from the MMRF web portal (<https://research.themmr.org/>). Clinical data and
14 CD138+ RNAseq were available for 781 patients at baseline (newly diagnosed bone marrow
15 samples) and 123 follow-up bone marrow samples. Transcript-based expression estimates
16 processed by Salmon (version 0.7.2) were used. The 768 baseline samples were used in the
17 PCA to derive the CD138+ transcriptome framework and SPECTRA variables. Covariates
18 included in batch correction were age, gender, overall survival, progression-free survival, and
19 time to first-line treatment failure. The first 39 components were selected based on the scree
20 test.⁴⁴ All 39 spectra were forced variables in all regression models.

21 To illustrate the flexibility of the transcriptome framework, linear, logistic, and Cox
22 regression were performed for several different clinical outcomes and demographic risk groups.
23 In each analysis, all 39 CD138+ spectra were entered into the model as independent, predictor
24 variables. No model fitting was performed. An individual spectrum was considered significant in
25 a model if its model coefficient was significantly different from 1.0 ($p < 0.05$). A likelihood ratio
26 test comparing the full 39-spectra model to the null model was used to determine the
27 significance of the overall model fit. To illustrate the aggregate effect of all spectra in the model

1 we used a poly-spectra liability (**PSL**) score. This score is the weighted sum of the spectra
2 values based on the spectra coefficients in the model. In our illustrations here, the PSL scores
3 contain all 39 spectra. In other applications, such as penalized modeling, a PSL score may
4 include only those spectra retained in the model.

5 To illustrate the potential to track longitudinal changes, spectra and PSL scores were
6 calculated for follow-up longitudinal samples. To enable this, batch corrected gene-level
7 measures for the follow-up samples were centered on the mean of the baseline data (\bar{e}), and
8 then multiplied with the rotation matrix (Q_k) which holds the linear transformations for the
9 spectra framework.

10

11 **Data availability.** Processed RNAseq data from the CoMMpass Study can be downloaded from
12 <https://research.themmr.org/>. Dimension variables for the IA14 CoMMpass data are provided in
13 the Supplement. We also provide the details of the QC process and the transcriptome
14 framework (linear equations for the gene transformations) necessary to calculate the 39-spectra
15 variables in other studies in the Supplement.

16

17 **Code availability.** R markdown notebooks used to derive the CD138+ transcriptome spectra
18 and generate the myeloma results are included in the Supplement.

19

20 **Funding.** The research reported in this publication was supported by the National Cancer
21 Institute (Award Numbers F99CA234943, K00CA234943, K07CA230150, and P30CA042014-
22 29S9), the National Center for Advancing Translational Sciences (Award Number
23 UL1TR002538), the National Library of Medicine (Award Number T15LM007124) of the National
24 Institutes of Health. The content is solely the responsibility of the authors and does not
25 necessarily represent the official views of the National Institutes of Health.

References

1. Madsen, M. J. *et al.* Reparameterization of PAM50 Expression Identifies Novel Breast Tumor Dimensions and Leads to Discovery of a Genome-Wide Significant Breast Cancer Locus at *12q15*. *Cancer Epidemiol Biomarkers Prev* **27**, 644–652 (2018).
2. Sweeney, C. *et al.* Intrinsic subtypes from PAM50 gene expression assay in a population-based breast cancer cohort: differences by age, race, and tumor characteristics. *Cancer Epidemiol. Biomarkers Prev.* **23**, 714–724 (2014).
3. Stopsack, K. H. *et al.* Regular aspirin use and gene expression profiles in prostate cancer patients. *Cancer Causes Control* **29**, 775–784 (2018).
4. Wang, S. *et al.* Gene expression in triple-negative breast cancer in relation to survival. *Breast Cancer Res. Treat.* **171**, 199–207 (2018).
5. Bhattacharya, A. *et al.* A framework for transcriptome-wide association studies in breast cancer in diverse study populations. *Genome Biol.* **21**, 42 (2020).
6. Caan, B. J. *et al.* Intrinsic subtypes from the PAM50 gene expression assay in a population-based breast cancer survivor cohort: prognostication of short- and long-term outcomes. *Cancer Epidemiol. Biomarkers Prev.* **23**, 725–734 (2014).
7. Allott, E. H. *et al.* Bimodal age distribution at diagnosis in breast cancer persists across molecular and genomic classifications. *Breast Cancer Res. Treat.* **179**, 185–195 (2020).
8. Troester, M. A. *et al.* Racial Differences in PAM50 Subtypes in the Carolina Breast Cancer Study. *J. Natl. Cancer Inst.* **110**, (2018).
9. Huo, D. *et al.* Comparison of Breast Cancer Molecular Features and Survival by African and European Ancestry in The Cancer Genome Atlas. *JAMA Oncol* **3**, 1654–1662 (2017).

10. Millstein, J. *et al.* Prognostic gene expression signature for high-grade serous ovarian cancer. *Ann. Oncol.* (2020) doi:10.1016/j.annonc.2020.05.019.
11. López, C. *et al.* Genomic and transcriptomic changes complement each other in the pathogenesis of sporadic Burkitt lymphoma. *Nat Commun* **10**, 1459 (2019).
12. Zhu, B. *et al.* Immune gene expression profiling reveals heterogeneity in luminal breast tumors. *Breast Cancer Res.* **21**, 147 (2019).
13. Zhang, M. *et al.* Characterising cis-regulatory variation in the transcriptome of histologically normal and tumour-derived pancreatic tissues. *Gut* **67**, 521–533 (2018).
14. Altman, D. G. & Royston, P. The cost of dichotomising continuous variables. *BMJ* **332**, 1080 (2006).
15. Perou, C. M. *et al.* Molecular portraits of human breast tumours. *Nature* **406**, 747–752 (2000).
16. Brunet, J.-P., Tamayo, P., Golub, T. R. & Mesirov, J. P. Metagenes and molecular pattern discovery using matrix factorization. *Proc. Natl. Acad. Sci. U.S.A.* **101**, 4164–4169 (2004).
17. Reich, M. *et al.* The GenePattern Notebook Environment. *Cell Syst* **5**, 149-151.e1 (2017).
18. Subramanian, A. *et al.* Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proc. Natl. Acad. Sci. U.S.A.* **102**, 15545–15550 (2005).
19. Tamayo, P. *et al.* Interpreting patterns of gene expression with self-organizing maps: methods and application to hematopoietic differentiation. *Proc. Natl. Acad. Sci. U.S.A.* **96**, 2907–2912 (1999).
20. Zhang, S., Li, X., Lin, Q., Lin, J. & Wong, K.-C. Uncovering the key dimensions of high-throughput biomolecular data using deep learning. *Nucleic Acids Res.* **48**, e56 (2020).

21. Chen, K. M. *et al.* PathCORE-T: identifying and visualizing globally co-occurring pathways in large transcriptomic compendia. *BioData Min* **11**, 14 (2018).
22. Sompairac, N. *et al.* Independent Component Analysis for Unraveling the Complexity of Cancer Omics Datasets. *Int J Mol Sci* **20**, (2019).
23. Way, G. P., Zietz, M., Rubinetti, V., Himmelstein, D. S. & Greene, C. S. Compressing gene expression data using multiple latent space dimensionalities learns complementary biological representations. *Genome Biol.* **21**, 109 (2020).
24. Camp, N. J. *et al.* Re-interpretation of PAM50 gene expression as quantitative tumor dimensions shows utility for clinical trials: application to prognosis and response to paclitaxel in breast cancer. *Breast Cancer Res Treat* **175**, 129–139 (2019).
25. Hanson, H. A. *et al.* Family Study Designs Informed by Tumor Heterogeneity and Multi-Cancer Pleiotropies: The Power of the Utah Population Database. *Cancer Epidemiology, Biomarkers & Prevention: A Publication of the American Association for Cancer Research, Cosponsored by the American Society of Preventive Oncology* **29**, 807–815 (2020).
26. Keats, J. J. *et al.* Interim Analysis Of The Mmrf Commpass Trial, a Longitudinal Study In Multiple Myeloma Relating Clinical Outcomes To Genomic and Immunophenotypic Profiles. *Blood* **122**, 532–532 (2013).
27. Shaughnessy, J. D. *et al.* A validated gene expression model of high-risk multiple myeloma is defined by deregulated expression of genes mapping to chromosome 1. *Blood* **109**, 2276–2284 (2007).
28. Zamani-Ahmadmahmudi, M., Nassiri, S. M. & Soltaninezhad, F. Development of a RNA sequencing-based prognostic gene signature in multiple myeloma. *Br J Haematol* bjh.16744 (2020) doi:10.1111/bjh.16744.

29. Patro, R., Duggal, G., Love, M. I., Irizarry, R. A. & Kingsford, C. Salmon provides fast and bias-aware quantification of transcript expression. *Nature Methods* **14**, 417–419 (2017).
30. Paner, A., Patel, P. & Dhakal, B. The evolving role of translocation t(11;14) in the biology, prognosis, and management of multiple myeloma. *Blood Rev.* 100643 (2019)
doi:10.1016/j.blre.2019.100643.
31. Mikhael, J. R. *et al.* Management of Newly Diagnosed Symptomatic Multiple Myeloma: Updated Mayo Stratification of Myeloma and Risk-Adapted Therapy (mSMART) Consensus Guidelines 2013. *Mayo Clinic Proceedings* **88**, 360–376 (2013).
32. Palumbo, A. *et al.* Revised International Staging System for Multiple Myeloma: A Report From International Myeloma Working Group. *JCO* **33**, 2863–2869 (2015).
33. Myeloma - Cancer Stat Facts. <https://seer.cancer.gov/statfacts/html/mulmy.html>.
34. Ramón y Cajal, S. *et al.* Clinical implications of intratumor heterogeneity: challenges and opportunities. *J Mol Med* **98**, 161–177 (2020).
35. Kwa, M., Makris, A. & Esteva, F. J. Clinical utility of gene-expression signatures in early stage breast cancer. *Nat Rev Clin Oncol* **14**, 595–610 (2017).
36. Stein-O'Brien, G. L. *et al.* Enter the Matrix: Factorization Uncovers Knowledge from Omics. *Trends Genet.* **34**, 790–805 (2018).
37. Szalat, R., Avet-Loiseau, H. & Munshi, N. C. Gene Expression Profiles in Myeloma: Ready for the Real World? *Clin Cancer Res* **22**, 5434–5442 (2016).
38. Dai, X. *et al.* Breast cancer intrinsic subtype classification, clinical use and future trends. *Am J Cancer Res* **5**, 2929–2943 (2015).

39. Manojlovic, Z. *et al.* Comprehensive molecular profiling of 718 Multiple Myelomas reveals significant differences in mutation frequencies between African and European descent cases. *PLoS Genet* **13**, e1007087 (2017).
40. Ma, S. & Dai, Y. Principal component analysis based methods in bioinformatics studies. *Briefings in Bioinformatics* **12**, 714–722 (2011).
41. GTEx Consortium. The Genotype-Tissue Expression (GTEx) project. *Nat. Genet.* **45**, 580–585 (2013).
42. Cancer Genome Atlas Research Network *et al.* The Cancer Genome Atlas Pan-Cancer analysis project. *Nat. Genet.* **45**, 1113–1120 (2013).
43. Arora, S., Pattwell, S. S., Holland, E. C. & Bolouri, H. Variability in estimated gene expression among commonly used RNA-seq pipelines. *Sci Rep* **10**, 2734 (2020).
44. Cattell, R. B. The Scree Test For The Number Of Factors. *Multivariate Behav Res* **1**, 245–276 (1966).
45. Zhao, S., Xi, L. & Zhang, B. Union Exon Based Approach for RNA-Seq Gene Quantification: To Be or Not to Be? *PLoS ONE* **10**, e0141910 (2015).
46. Shahriyari, L. Effect of normalization methods on the performance of supervised learning algorithms applied to HTSeq-FPKM-UQ data sets: 7SK RNA expression as a predictor of survival in patients with colon adenocarcinoma. *Briefings in Bioinformatics* **20**, 985–994 (2019).
47. Johnson, W. E., Li, C. & Rabinovic, A. Adjusting batch effects in microarray expression data using empirical Bayes methods. *Biostatistics* **8**, 118–127 (2007).

48. Leek, J. T., Johnson, W. E., Parker, H. S., Jaffe, A. E. & Storey, J. D. The sva package for removing batch effects and other unwanted variation in high-throughput experiments. *Bioinformatics* **28**, 882–883 (2012).

Table 1. Hazard Ratios for overall survival (95% CI)

	1 year	3 year	5 year
CD138+ spectra OS-PSL score*	4.8 (2.7-8.5)	6.1 (4.0-9.3)	8.5 (3.4-21.7)
UAMS-70 risk score²⁷	3.1 (1.8-5.5)	3.6 (2.3-5.6)	4.3 (1.5-12.7) [†]
SBUK-17 prognosis score²⁸	2.7 (1.6-4.7)	3.2 (2.2-4.9) [‡]	1.3 (0.6-2.8)

Notes: *Overall survival poly-spectra liability score tertile 3 vs tertiles 1 and 2. The cutoff has not been optimized for OS. We used the UAMS-70 and SBUK-17 high vs low score cutoff as published.^{27,28} All HRs are based on the CoMMpass data. The OS-PSL is shown in training data. For comparison, UAMS-70 published [†]HR = 5.2 at 5 years²⁷ and SBUK-17 published [‡]HR = 3.7 (2.5-5.7) at 3 years²⁸ in training data.

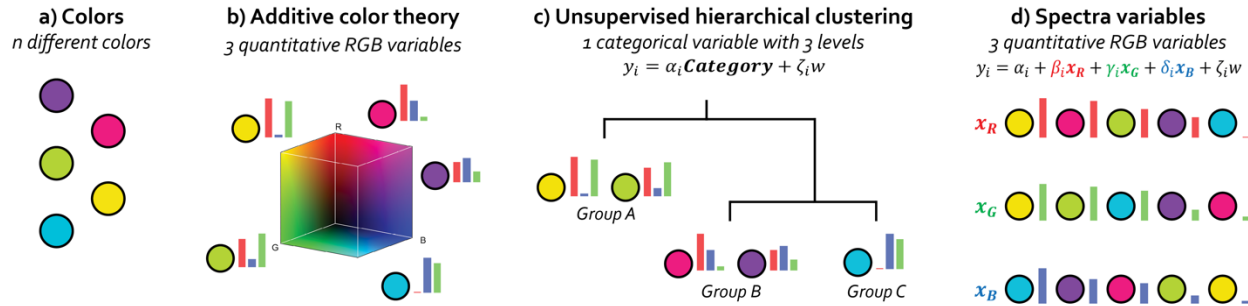


Figure 1. A color analogy to illustrate the advantages of spectra variables for modeling. a) Individual observations of color. b) Dimension Reduction (additive color theory), all colors can be represented using 3 quantitative RGB variables. c) Standard-use, modeling on the 3 RGB variables used to identify structure across samples using hierarchical clustering. This derives groups based on the complete 3-variable RGB profile to derive one polychotomous meta-variable (different groups are non-ordinal levels). d) Multivariable modeling implementation of spectra variables, multiple separate spectrum integrated directly into a multivariable analysis. Each uncorrelated variable can be assessed separately for its predictive value for an outcome. This implementation retains the full resolution of the initial data because the variables are quantitative and retain integrity to the initial data. Note, lower-resolution versions of x_B and x_G can be achieved using hierarchical groups but the loss of quantification will likely also lose power. x_R cannot be captured by any group ordering and associations for this spectrum would be lost using hierarchical groups.

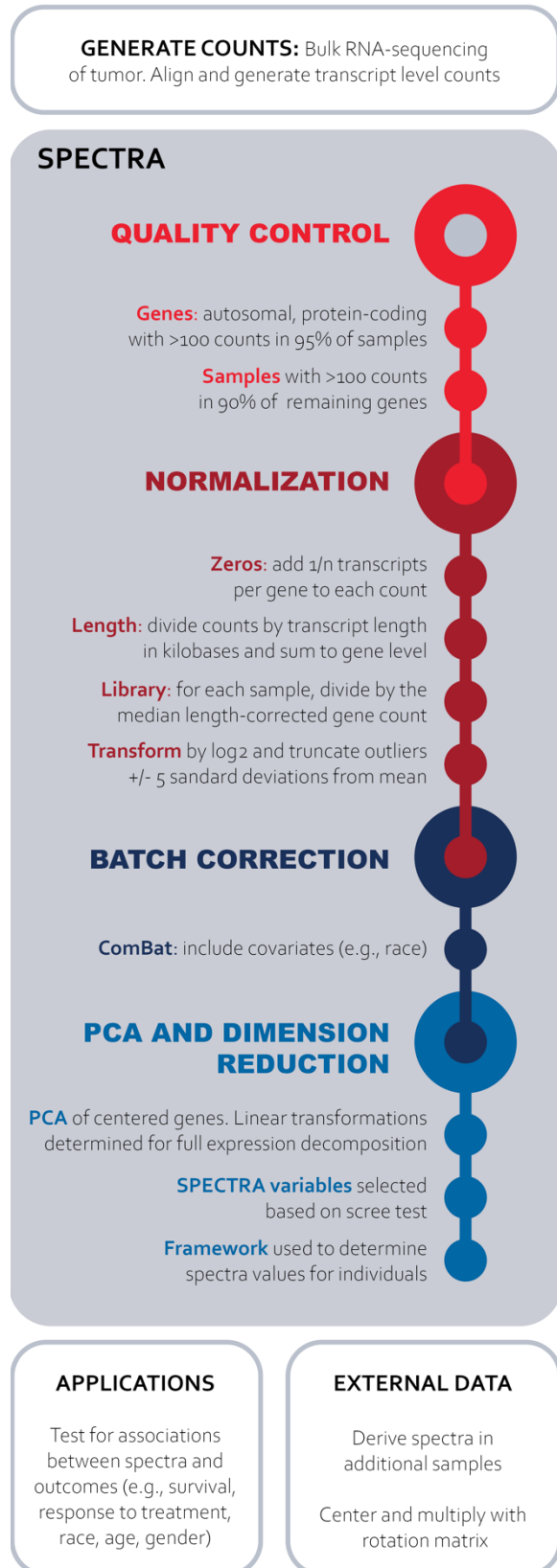


Figure 2. Overview of SPECTRA workflow.

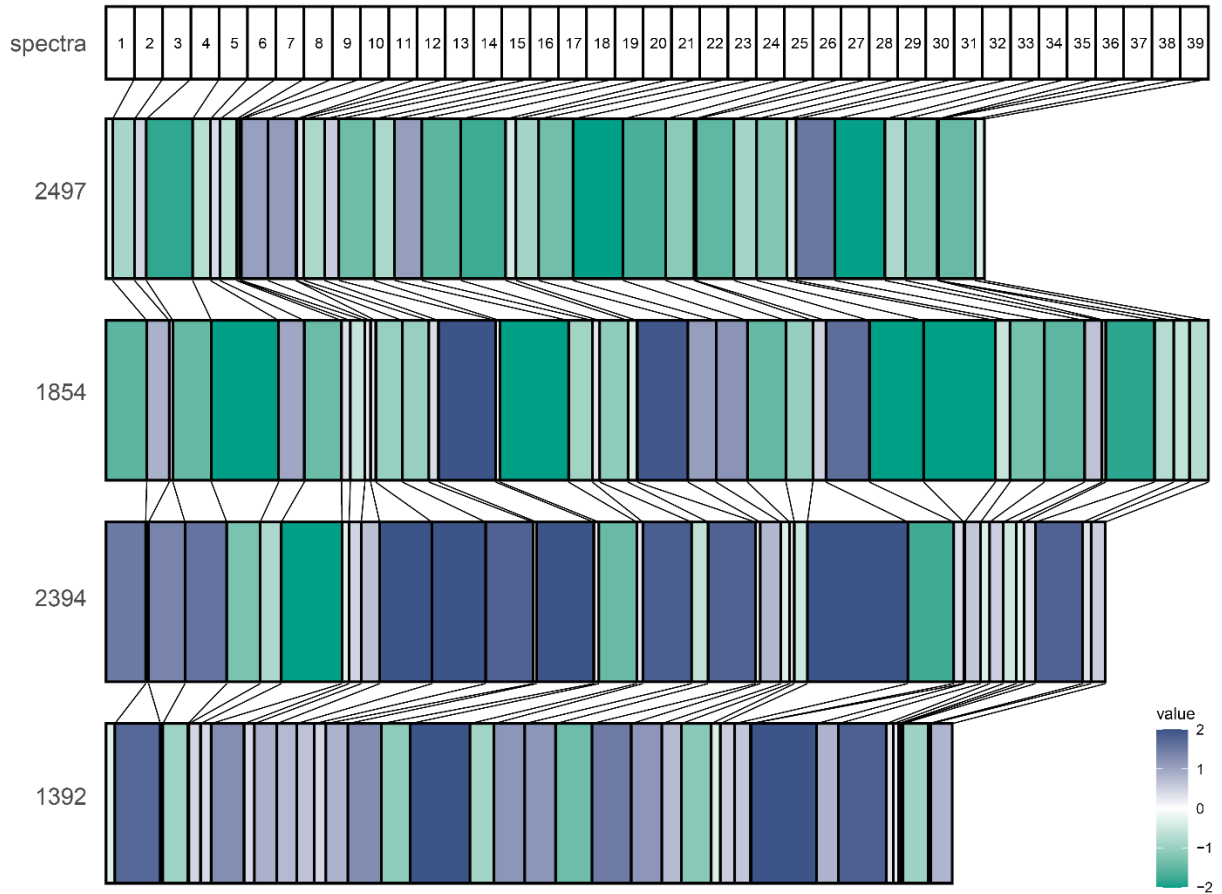


Figure 3. Spectra charts in four CoMMpass patients. For each patient, all 39 spectra are illustrated with the value represented by the bar width and intensity. The color indicates if the patient's spectra value is positive or negative. Each patient has a unique profile across the 39 spectra. At a high-level patients 2497 & 1854 may appear most similar (mostly green) and 2394 & 1392 (mostly blue). However, at a finer resolution, similarities vary. For example, for spectrum S1 and S2, patients 2497 and 1854 are quite different. In fact, for spectrum S1, 2497 is more similar to 1392, for spectrum S15, 1854 is more similar to 2394.

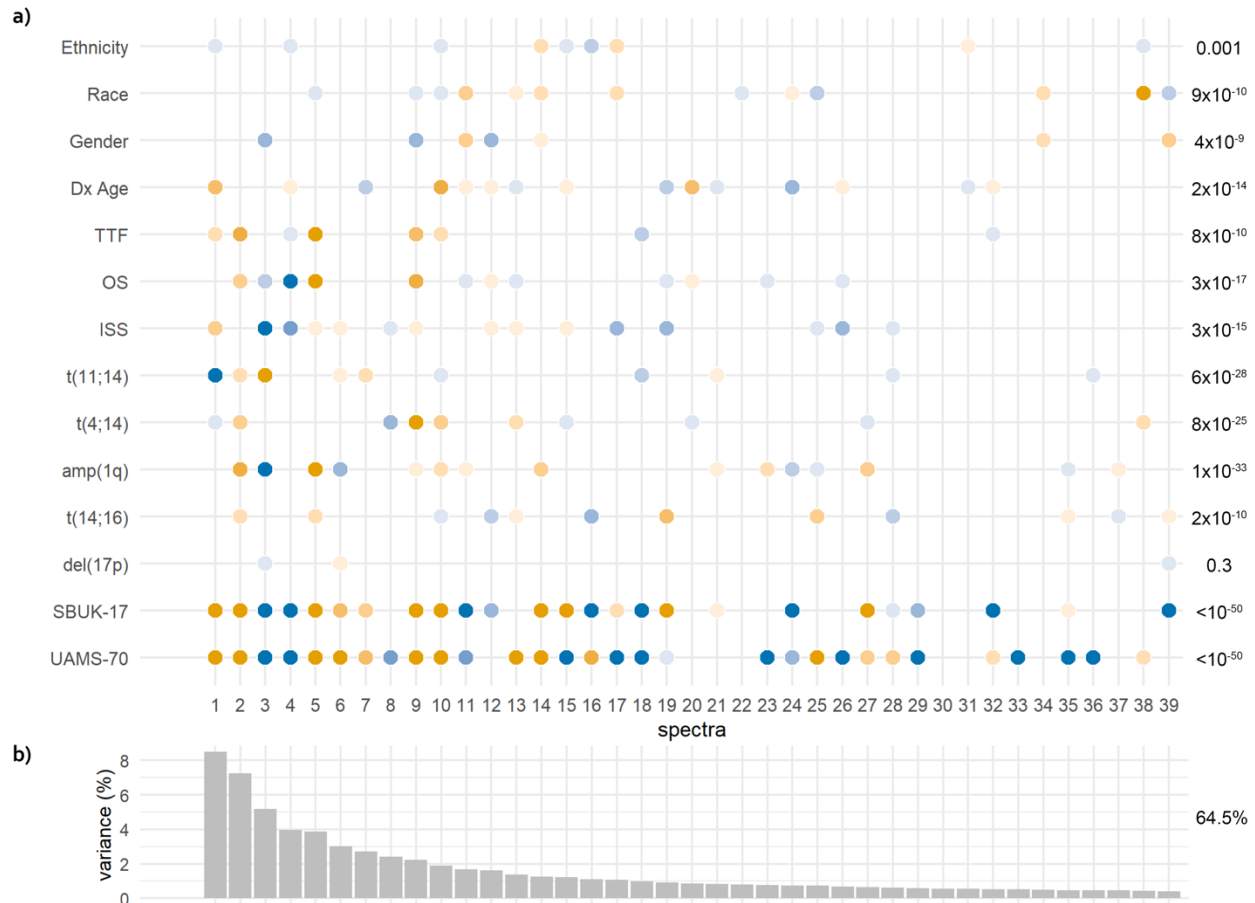


Figure 4. Overview of CD138+ spectra outcome modeling. a) Multivariable modeling results. The outcome of interest (left y-axis) and the significance of the full 39-spectra model (right y-axis). Spectra variables are illustrated on the x-axis. The significance and direction of each spectrum are indicated: blue negative beta coefficient, orange positive beta coefficient. No dot is shown if a spectrum was not significantly associated at $p < 0.05$ level. **b) Percent of the global variance captured by each spectra variable.** The total variance captured by all 39 spectra is shown at the right.

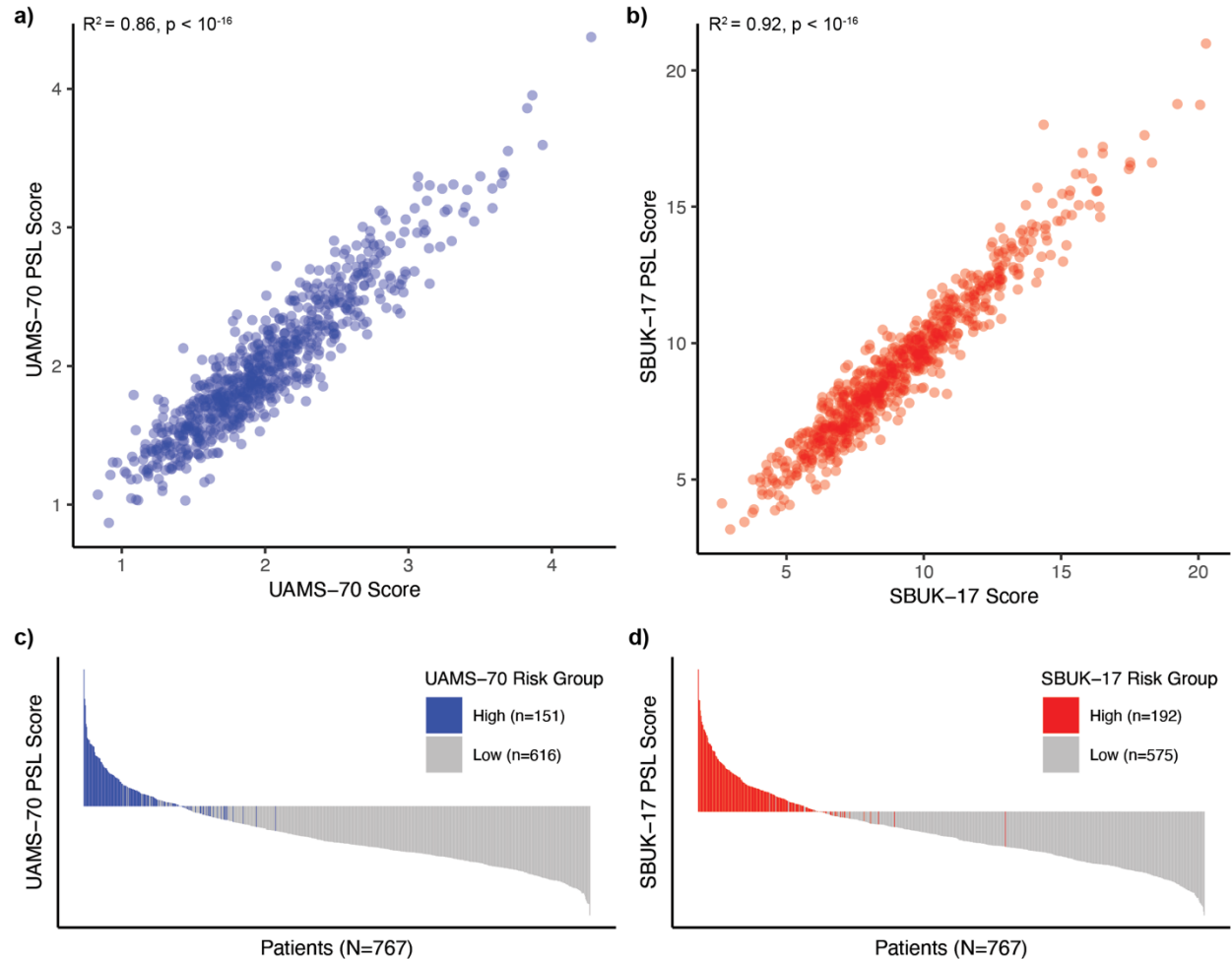


Figure 5. CD138+ spectra and two established expression scores. Correlation of the 39-spectra PSL score and established gene expression profiles from **a)** University of Alabama School of Medicine 70 gene risk score (UAMS-70) and **b)** Shahid Bahonar University of Kerman 17-gene prognostic score (SBUK-17). Waterfall plots for **c)** UAMS-70 and **d)** SBUK-17 low- and high-risk scores. Patients were ordered by their PSL score and colored by high/low risk as predicted by the UAMS-70 score. (UAMS-70 high-risk cutoff determined by clustering).

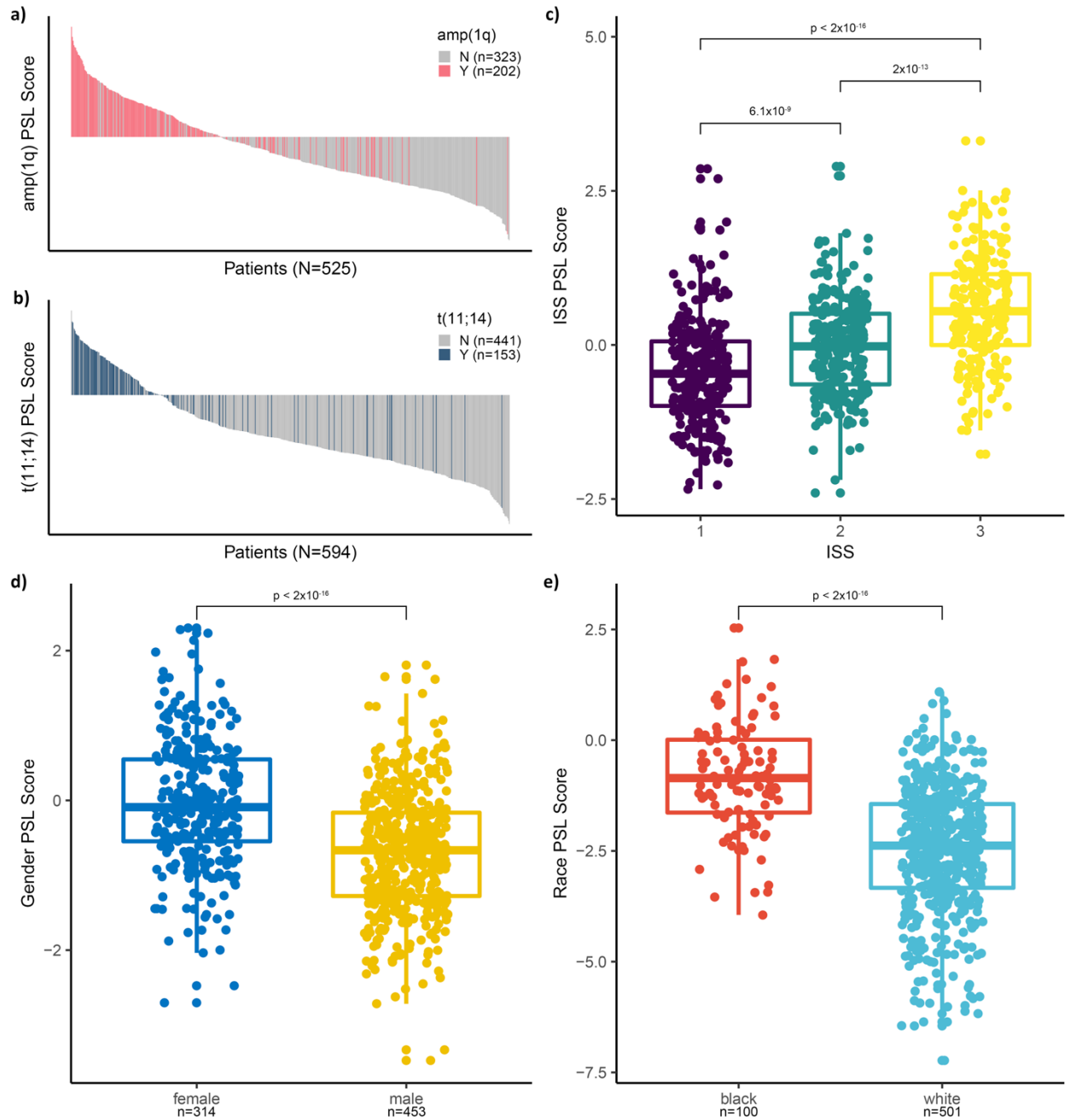
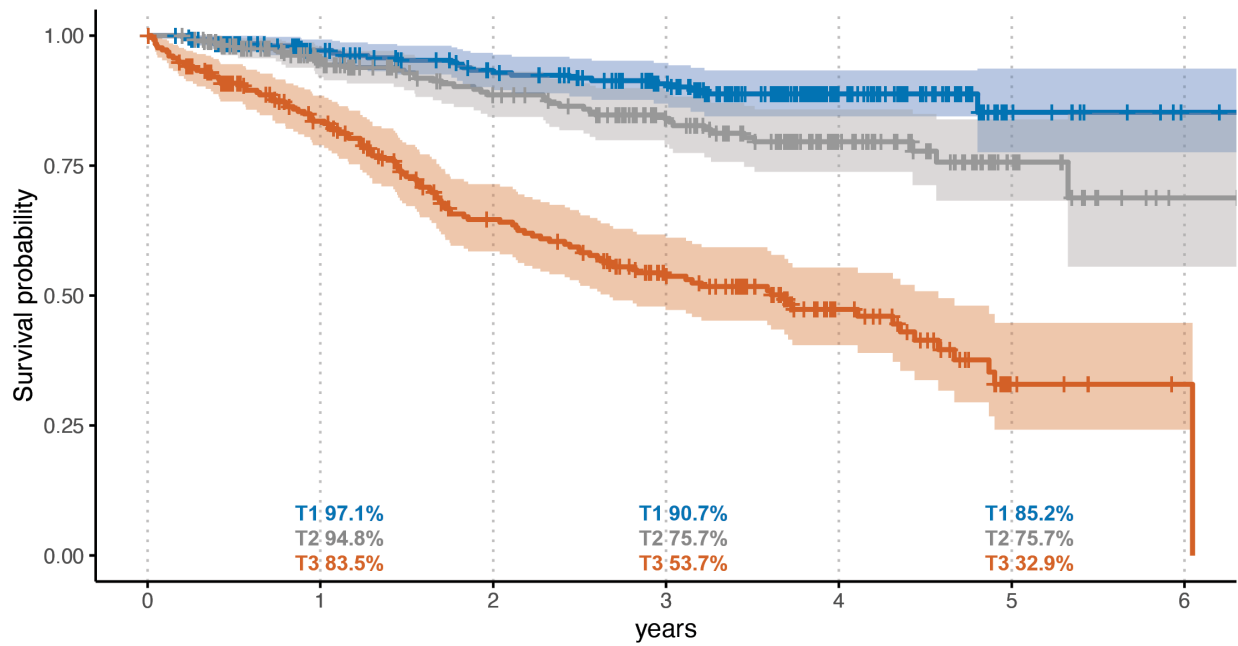


Figure 6. CD138+ spectra and clinical or demographic risk. Waterfall plots for PSL scores for models of **a)** tumor amplification chr1q, **b)** tumor translocation chr11;16. Box and Whisker plots for PSL scores for models of **c)** international tumor stage at diagnosis, **d)** gender, and **e)** self-reported black or white race.

a) Overall Survival, 179 events, $\chi^2_{\text{logrank}} = 120$, $p = 8.6 \times 10^{-27}$



b) Time to first-line treatment failure, 369 events, $\chi^2_{\text{logrank}} = 95.7$, $p = 1.7 \times 10^{-21}$

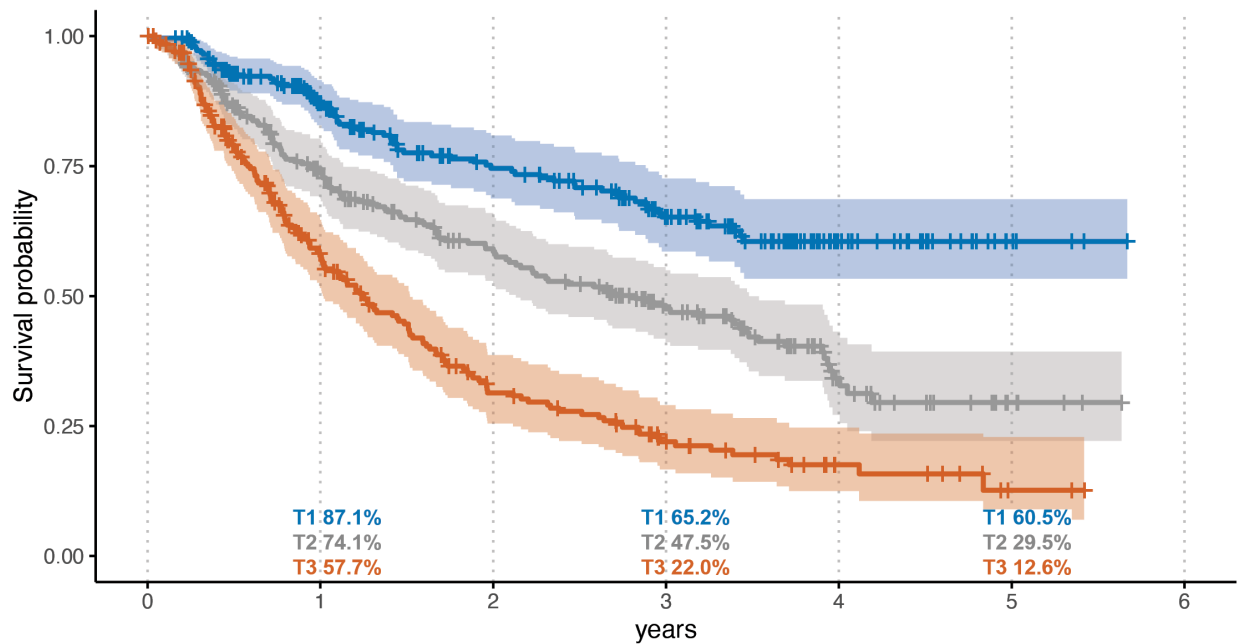


Figure 7. CD138+ spectra and disease course. Cox proportional hazards models were generated for overall survival (OS) and time to first-line treatment failure (TTF). From the models, PSL scores were generated and split into three equal tertiles. Kaplan-Meier curves of the PSL scores by tertile are shown for **a)** overall survival and **b)** time to first-line treatment failure.

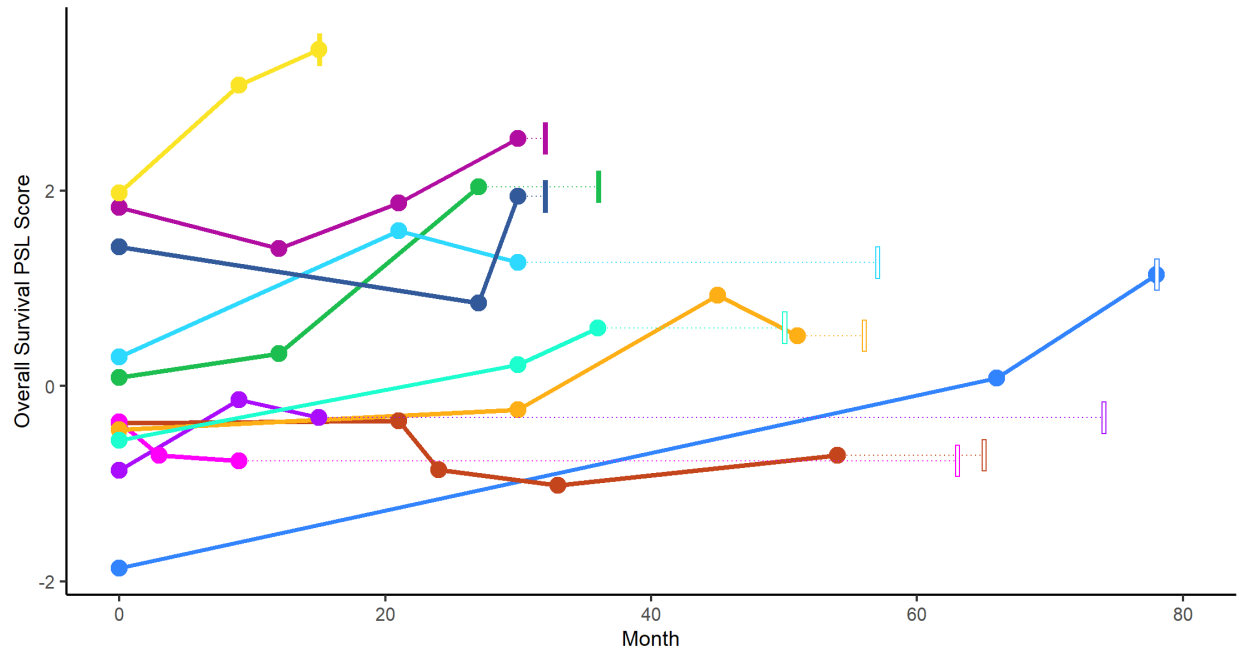


Figure 8. Aggregate effect of CD138+ spectra for overall survival over time. Poly-spectra liability scores for overall survival are shown for eleven patients with RNAseq data at multiple time points. Dots indicate sequencing events and show the PSL score at that timepoint. The final narrow rectangle indicates the month after diagnosis the patient died (filled) or was last known alive (open).