1	Deep transcriptome profiling of multiple myeloma with quantitative measures using the SPECTRA
2	approach
3	
4	Rosalie Griffin Waller, ^{1,2} Heidi A. Hanson, ¹ Brian J. Avery, ¹ Michael J. Madsen, ¹ Douglas W. Sborov, ¹
5	and Nicola J. Camp ¹ *
6	
7	1. Huntsman Cancer Institute and School of Medicine, University of Utah, Salt Lake City, UT, 84112,
8	USA
9	2. Computational Biology, Quantitative Health Sciences, Mayo Clinic, Rochester, MN, 55905, USA
10	
11	Correspondence: * <u>nicola.camp@utah.edu</u>

2

1 ABSTRACT

2 SPECTRA is a new data framework to describe variation in a transcriptome as a set of unsupervised quantitative variables. Spectra variables provide a deep dive into the transcriptome, 3 representing both large and small sources of variance, and are ideal for modeling alongside other 4 5 variables for any outcome of interest. Each spectrum can also be considered a phenotypic trait, providing new avenues for disease characterization or to explore disease risk. We applied the SPECTRA approach 6 to multiple myeloma (MM), the second most common blood cancer. Using RNA sequencing from 7 malignant CD138+ cells, we derived 39 spectra in 767 patients from the MMRF CoMMpass study. We 8 included spectra in prediction models for clinical endpoints, compared to established expression-based 9 risk scores, and used descriptive modeling to identify associations with patient characteristics. Spectra-10 based risk scores added predictive value beyond established clinical risk factors and other expression-11 based risk scores for overall survival, progression-free survival, and time to first-line treatment failure. 12 Significant spectra in models may provide mechanistic insight via gene set enrichment based on their 13 gene weights. Gene set enrichment in CD138+ spectrum S5, which was significant for all prognostic 14 endpoints, indicated enrichment for genes in the unfolded protein response, a mechanism targeted by 15 proteasome inhibitors, common first line agents in MM treatment. We also identified significant 16 17 associations between CD138+ spectra and tumor cytogenetics, race, gender, and age at diagnosis. The SPECTRA approach provides measures of transcriptome variation to deeply profile tumors with greater 18 flexibility to model clinical outcomes and characteristics. 19

3

1 AUTHOR SUMMARY

2 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 3 attractive because it captures genetic, epigenetic, and environmental consequences. Transcriptome studies 4 5 are gaining momentum in genomic epidemiology, and the need to incorporate these data in multivariable models alongside other risk factors brings demands for new approaches. The SPECTRA approach is a 6 new intrinsic quantitative data framework for transcriptomes. A tissue is described by a set of quantitative 7 measures (or 'spectra' variables) to deeply profile gene expression in a tissue. Spectra variables are 8 independent and offer flexibility for use in predictive or descriptive modeling. We applied the SPECTRA 9 approach to multiple myeloma, the second most common blood cancer. A set of 39 spectra variables were 10 derived to represent the myeloma tumors. Outcome modeling provided SPECTRA-based risk scores that 11 12 added predictive value for clinical outcomes beyond established risk factors.

4

1 INTRODUCTION

Numerous factors are involved in risk and prognosis in complex disease. Transcriptomes represent the combined effects of inherited, somatic, and epigenetic variation and can provide insight into genetic and environmental risk factors. As a result, gene expression studies are gaining momentum in genomic epidemiology (Allott et al., 2020; López et al., 2019; Stopsack et al., 2018; Sweeney et al., 2014; Zhang et al., 2018). The need to incorporate transcriptome data in multivariable models alongside other risk factors brings new demands for approaches to describe transcriptomes.

Many current transcriptome approaches are focused on immediate biological interpretation and 8 9 constrained to biological expectations or reduce the data to a single categorical variable (e.g., intrinsic subtypes). While these have advanced knowledge of disease mechanisms (Brunet et al., 2004; Tamayo et 10 al., 1999; Way et al., 2020) and identified important high-level differences in disease (Lapointe et al., 11 12 2004; Perou et al., 2000; Shaughnessy et al., 2007), complementary approaches are needed to go deeper and increase flexibility and application. In common disease, the sources of heterogeneity are many, 13 complex, and often poorly understood. Latent variables, focused on capturing signal and not immediate 14 interpretability, provide the potential for new discoveries. Transcriptome characterization using multiple 15 quantitative variables may provide a meaningful deeper dive into the transcriptome. 16

17 Here we describe the SPECTRA approach, a data workflow and variable derivation method with principal component analysis (PCA) at its core. PCA has many strengths aligned with our intent, such as 18 providing multiple unsupervised variables that optimize coverage of the global variance. The resulting 19 20 latent variables are quantitative, uncorrelated, retain integrity to the original data, and have desirable 21 attributes for subsequent multivariable prediction and descriptive modeling. A limitation of PCA is that if data are not curated adequately, spurious variance will be incorporated in the derived variables, including 22 technical artifacts and missing data structures. A key part of the SPECTRA approach is stringent data 23 culling, quality control, zero-handling, and normalization. 24

For well-curated gene-panel data a PCA-based approach has shown promise. Using a populationbased dataset of breast tumors, we previously used PCA to reduce the 50-gene space of the PAM50 gene panel to five quantitative multi-gene expression variables (Madsen et al., 2018). When implemented as

5

predictor variables in an independent clinical trial dataset, PCA variables were able to predict prognosis 1 and response to paclitaxel; adding value beyond clinical risk factors and outperforming intrinsic subtypes 2 (Camp et al., 2019). When utilized as quantitative tumor phenotypes, PCA variables were superior to the 3 standard PAM50 subtypes for gene mapping (Hanson et al., 2020; Madsen et al., 2018). Here, we 4 5 describe a method to derive a quantitative data framework for whole transcriptomes. 6 Figure 1 uses a color analogy to illustrate the conceptual framework of SPECTRA and contrast 7 our goal of quantitative variables for direct use in outcome modeling with a more conventional categorization approach using hierarchical clustering. Each spectrum in Figure 1 (x_R, x_G, x_B) are 8 independent variables that can be directly used to model any outcome (y_i) , and other covariates/predictors 9 10 can also be easily included (Figure 1d). Conversely, unsupervised hierarchical clustering uses the spectra to categorize patients into groups (Figure 1c), flattening the multiple variables to a single categorical 11 variable which may reduce statistical power. For example, in Figure 1, x_R cannot be represented by any 12 13 group ordering in **Figure 1c**, and associations for that spectra variable would be lost. An alternate convention is to supervise clustering to an outcome. But, while supervised clustering can improve power 14 over unsupervised clustering for the prediction of a single outcome, it also tethers the groups to the 15 16 trained outcome and doesn't facilitate comparison to other outcomes. We illustrate our whole transcriptome SPECTRA approach using bulk RNA sequencing 17 18 (RNAseq) data from CD138+ sorted myeloma cells from the Multiple Myeloma Research Foundation 19 (MMRF) CoMMpass Study (Keats et al., 2013). Multiple myeloma (MM) is a malignancy of the plasma cells with one of the poorest 5-year survival (55.6%) for adult-onset hematological malignancies (SEER, 20 2021). It is most frequently diagnosed at ages 65-74 years (median 69 years) (SEER, 2021). Incidence is 21 22 higher in men (8.7 men vs. 5.6 women per 100,000) and particularly high in patients self-reporting as 23 African American (AA men 16.3, and AA women 11.9 per 100,000) (SEER, 2021). We use the SPECTRA approach to derive quantitative CD138+ transcriptome variables, referred to as spectra, which 24 25 we use in regression models to predict clinical endpoints, derive latent risk groups that may be clinically meaningful, compare to established expression-based molecular risk scores, and describe associations 26 27 with clinical and demographic characteristics.

6

1 **RESULTS**

2 SPECTRA: an approach for deep transcriptome profiling with quantitative measures

The motivation is the derivation of well-behaved, quantitative variables from RNAseq data to 3 capture transcriptome variation that can be used universally as predictors for any outcome, and as novel 4 5 phenotypes. The approach requires an RNAseq dataset to derive the framework of transformations for the SPECTRA variables and multiple spectra are calculated for each individual in the dataset. An overview of 6 the SPECTRA approach is shown in Figure 2 and Figure S1. As an agnostic technique, the goal is to 7 retain only those aspects of the RNAseq data that can represent meaningful variance. Accordingly, genes 8 9 likely to lack precision are removed and only coding genes with sufficient coverage across the dataset are considered. An internal normalization procedure accounts for feature-length, library size, and RNA 10 composition. This normalization avoids the need for reference samples, real or synthetic, and provides the 11 potential for spectra to be ported to follow-up samples and external datasets. Finally, skew and outliers 12 are dealt with before PCA is performed. Specific details are listed in Materials and Methods. 13

PCA is a well-established, unsupervised, data-driven method that, based on the covariance of a dataset, produces a matrix factorization which is a unique solution of linear transformations and transformed values. The linear transformations preserve the variance in the data and provide meaningful comparisons between individuals. The transformed values are orthogonal quantitative variables (linearly uncorrelated), each subsequently going deeper into the global variance. Dimension reduction is achieved by selecting only the first *k* components, for which the proportion of total variance explained can be described. Details of the matrix factorization performed by PCA can be found in Supplemental Methods.

The results of the subsequent PCA are the rotation matrix that describes the multi-gene linear transformations; and the transformed data matrix, the quantitative variables for each individual referred to as <u>SPECTRA variables</u>, or <u>spectra</u>. The set of linear transformations provides a new reduced-dimension <u>framework</u> for the expression space. The SPECTRA variables are linearly independent, each providing additional coverage of the variance, and as unsupervised variables can be used as predictors for any outcome and as novel phenotypes.

7

2 <u>CD138+ Spectra</u>

3	We applied the SPECTRA approach to RNAseq from treatment naïve CD138+ cells collected at
4	diagnosis from 768 patients in the MMRF CoMMpass study (Keats et al., 2013). We used transcript-
5	based expression estimates generated by the CoMMpass study with Salmon (Patro et al., 2017). After
6	quality control, gene expression values for 7,449 genes (56,339 transcripts) in 767 patients were
7	normalized and batch corrected. After PCA, the first 39 components were selected based on a scree test
8	and standardized to create spectra S1-S39. Together these explain 65% of the global expression variation
9	(Figure 3A). As linearly uncorrelated variables each of the 39 CD138+ spectra capture a different source
10	of variance – each spectrum has the potential to describe patient differences. A patient's CD138+
11	transcriptome spectra variables can be visualized with a 39-spectra barcode (Figure 3B).
12	
13	Predictive Modeling – Overall Survival (OS)
14	Predictive Modeling used Cox proportional hazards regression with bootstrap internal validation
15	to adjust for over fitting (Harrell, 2015, pp. 114-116). This strategy maintains the full sample for
16	discovery (referred to as 'apparent' results) and implements bootstrapping to estimate over-fitting
17	(referred to as 'optimism'). The optimism corrected apparent results are referred to as 'adjusted' values.
18	In Cox regression of OS (179 events), nine spectra were significant ($p < 0.05$) and selected for the
19	predictive model (apparent C-index = 0.70). Internal validation (1000 bootstraps) estimated C-optimism =
20	0.04, leading to an adjusted C-index ($C_{adj} = 0.66$ (0.65-0.74). Model coefficients are given in Table 1 .
21	Gaussian mixture modeling on the OS spectra risk scores identified two risk groups ($p = 0.001$) (Figure
22	4) with 88 patients in the high-risk group (58 events, median OS 20.1 months) and 676 patients in the
23	low-risk group (121 events, median OS not reached after 78 months) (Figure 4C). After internal
24	validation, the optimism adjusted hazard ratio (HRadj) for the high- compared to low-risk patients was
25	4.27 (2.31-12.69) (Table 2). This compared to UAMS OS $HR_{adj} = 3.89$ (2.53-5.45) and SBUK OS HR_{adj}
26	= 2.54 (1.91-3.47) in CoMMpass data.

1	The added predictive value (APV) of the quantitative spectra risk score, beyond Revised
2	International Staging System (R-ISS) and age at diagnosis was 0.612 (Table 3). In parallel analyses,
3	UAMS APV was 0.576 and SBUK APV = 0.502, indicating each contains substantial predictive value
4	beyond clinical factors. Compared to a model including R-ISS, age at diagnosis, and UAMS, however,
5	SBUK showed very little added value (APV=0.006, Table 3), indicating its predictive gene expression
6	information is almost all duplicative with UAMS. In contrast, the spectra score $APV = 0.109$ beyond R-
7	ISS, age at diagnosis, and UAMS (Table 3) – also significant by likelihood ratio test (LRT, $p = 7.4 \times 10^{-5}$)
8	- indicating that the spectra risk score contains predictive information beyond clinical factors and the
9	previously established UAMS risk score.
10	Eleven patients had at least three follow-up CD138+ samples, spanning up to six years. To
11	illustrate a potential to track changes over time, the OS spectra risk score was graphed for the initial and
12	longitudinal follow-up samples (Figure 5).
13	
14	Spectra Predictive Modeling – Progression-Free Survival (PFS)
14 15	<u>Spectra Predictive Modeling – Progression-Free Survival (PFS)</u> As the spectra are unsupervised, the same 39 spectra variables can be used to model any outcome.
14 15 16	<u>Spectra Predictive Modeling – Progression-Free Survival (PFS)</u> As the spectra are unsupervised, the same 39 spectra variables can be used to model any outcome. We repeated the same predictive modeling procedure for PFS. Cox regression of PFS (392 events)
14 15 16 17	Spectra Predictive Modeling – Progression-Free Survival (PFS)As the spectra are unsupervised, the same 39 spectra variables can be used to model any outcome.We repeated the same predictive modeling procedure for PFS. Cox regression of PFS (392 events)selected nine spectra in a model with $C_{adj} = 0.60$ (0.60-0.67). Four of the retained spectra were distinct
14 15 16 17 18	Spectra Predictive Modeling – Progression-Free Survival (PFS)As the spectra are unsupervised, the same 39 spectra variables can be used to model any outcome.We repeated the same predictive modeling procedure for PFS. Cox regression of PFS (392 events)selected nine spectra in a model with $C_{adj} = 0.60$ (0.60-0.67). Four of the retained spectra were distinctfrom the OS spectra model (Table 1). Gaussian mixture modeling identified two risk groups (p = 0.001)
14 15 16 17 18 19	Spectra Predictive Modeling – Progression-Free Survival (PFS)As the spectra are unsupervised, the same 39 spectra variables can be used to model any outcome.We repeated the same predictive modeling procedure for PFS. Cox regression of PFS (392 events)selected nine spectra in a model with $C_{adj} = 0.60$ (0.60-0.67). Four of the retained spectra were distinctfrom the OS spectra model (Table 1). Gaussian mixture modeling identified two risk groups (p = 0.001)(Figure S2). Median time to progression was 9.7 months in the high-risk patients (n = 60, 50 events) and
14 15 16 17 18 19 20	Spectra Predictive Modeling – Progression-Free Survival (PFS)As the spectra are unsupervised, the same 39 spectra variables can be used to model any outcome.We repeated the same predictive modeling procedure for PFS. Cox regression of PFS (392 events)selected nine spectra in a model with $C_{adj} = 0.60$ (0.60-0.67). Four of the retained spectra were distinctfrom the OS spectra model (Table 1). Gaussian mixture modeling identified two risk groups (p = 0.001)(Figure S2). Median time to progression was 9.7 months in the high-risk patients (n = 60, 50 events) and35.7 months in the low-risk patients (n = 707, 342 events) (Figure 6A). Between high and low risk
14 15 16 17 18 19 20 21	Spectra Predictive Modeling – Progression-Free Survival (PFS)As the spectra are unsupervised, the same 39 spectra variables can be used to model any outcome.We repeated the same predictive modeling procedure for PFS. Cox regression of PFS (392 events)selected nine spectra in a model with $C_{adj} = 0.60 (0.60-0.67)$. Four of the retained spectra were distinctfrom the OS spectra model (Table 1). Gaussian mixture modeling identified two risk groups ($p = 0.001$)(Figure S2). Median time to progression was 9.7 months in the high-risk patients ($n = 60, 50$ events) and35.7 months in the low-risk patients ($n = 707, 342$ events) (Figure 6A). Between high and low riskgroups, spectra PFS HR _{adj} = 3.08 (1.67-8.44) while UAMS PFS HR _{adj} = 2.40 (1.70-3.31) and SBUK PFS
14 15 16 17 18 19 20 21 22	Spectra Predictive Modeling – Progression-Free Survival (PFS)As the spectra are unsupervised, the same 39 spectra variables can be used to model any outcome.We repeated the same predictive modeling procedure for PFS. Cox regression of PFS (392 events)selected nine spectra in a model with $C_{adj} = 0.60$ (0.60-0.67). Four of the retained spectra were distinctfrom the OS spectra model (Table 1). Gaussian mixture modeling identified two risk groups (p = 0.001)(Figure S2). Median time to progression was 9.7 months in the high-risk patients (n = 60, 50 events) and35.7 months in the low-risk patients (n = 707, 342 events) (Figure 6A). Between high and low riskgroups, spectra PFS HR _{adj} = 3.08 (1.67-8.44) while UAMS PFS HR _{adj} = 2.40 (1.70-3.31) and SBUK PFSHR _{adj} = 1.92 (1.59-2.43) (Table S1). Spectra added information beyond R-ISS, age at diagnosis, and
14 15 16 17 18 19 20 21 22 23	Spectra Predictive Modeling – Progression-Free Survival (PFS) As the spectra are unsupervised, the same 39 spectra variables can be used to model any outcome. We repeated the same predictive modeling procedure for PFS. Cox regression of PFS (392 events) selected nine spectra in a model with $C_{adj} = 0.60$ (0.60-0.67). Four of the retained spectra were distinct from the OS spectra model (Table 1). Gaussian mixture modeling identified two risk groups (p = 0.001) (Figure S2). Median time to progression was 9.7 months in the high-risk patients (n = 60, 50 events) and 35.7 months in the low-risk patients (n = 707, 342 events) (Figure 6A). Between high and low risk groups, spectra PFS HR _{adj} = 3.08 (1.67-8.44) while UAMS PFS HR _{adj} = 2.40 (1.70-3.31) and SBUK PFS HR _{adj} = 1.92 (1.59-2.43) (Table S1). Spectra added information beyond R-ISS, age at diagnosis, and UAMS (APV = 0.24, LRT p = 8.2x10 ⁻⁸) in predicting PFS (Table S2).
14 15 16 17 18 19 20 21 22 23 24	Spectra Predictive Modeling – Progression-Free Survival (PFS) As the spectra are unsupervised, the same 39 spectra variables can be used to model any outcome. We repeated the same predictive modeling procedure for PFS. Cox regression of PFS (392 events) selected nine spectra in a model with $C_{adj} = 0.60$ (0.60-0.67). Four of the retained spectra were distinct from the OS spectra model (Table 1). Gaussian mixture modeling identified two risk groups (p = 0.001) (Figure S2). Median time to progression was 9.7 months in the high-risk patients (n = 60, 50 events) and 35.7 months in the low-risk patients (n = 707, 342 events) (Figure 6A). Between high and low risk groups, spectra PFS HR _{adj} = 3.08 (1.67-8.44) while UAMS PFS HR _{adj} = 2.40 (1.70-3.31) and SBUK PFS HR _{adj} = 1.92 (1.59-2.43) (Table S1). Spectra added information beyond R-ISS, age at diagnosis, and UAMS (APV = 0.24, LRT p = 8.2x10 ⁻⁸) in predicting PFS (Table S2).
14 15 16 17 18 19 20 21 22 23 24 25	Spectra Predictive Modeling – Progression-Free Survival (PFS)As the spectra are unsupervised, the same 39 spectra variables can be used to model any outcome.We repeated the same predictive modeling procedure for PFS. Cox regression of PFS (392 events)selected nine spectra in a model with $C_{adj} = 0.60$ (0.60-0.67). Four of the retained spectra were distinctfrom the OS spectra model (Table 1). Gaussian mixture modeling identified two risk groups (p = 0.001)(Figure S2). Median time to progression was 9.7 months in the high-risk patients (n = 60, 50 events) and35.7 months in the low-risk patients (n = 707, 342 events) (Figure 6A). Between high and low riskgroups, spectra PFS HR _{adj} = 3.08 (1.67-8.44) while UAMS PFS HR _{adj} = 2.40 (1.70-3.31) and SBUK PFSHR _{adj} = 1.92 (1.59-2.43) (Table S1). Spectra added information beyond R-ISS, age at diagnosis, andUAMS (APV = 0.24, LRT p = $8.2x10^{-8}$) in predicting PFS (Table S2).Spectra Predictive Modeling – Time to First-Line Treatment Failure (TTF)

In Cox regression modeling to predict TTF (369 events) ten spectra were selected with the model $C_{adj} = 0.60 (0.59-0.66)$. Four of the retained spectra were distinct from the OS and PFS models (**Table 1**).

9

1	Two risk groups were identified in the TTF spectra score with GMM ($p = 0.008$) (Figure S3). Patients in
2	the spectra high-risk TTF ($n = 31, 25$ events) had median TTF of 9.2 months compared to low-risk
3	patients (n = 736, 344 events) with median TTF of 32.8 months (Figure 6B). Hazards ratios between
4	high- and low-risk groups using TTF spectra risk groups ($HR_{adj} = 3.10 [1.31-5.46]$) outperformed UAMS
5	$(TTF HR_{adj} = 1.98 [1.48-2.68])$ and SBUK $(TTF HR_{adj} = 1.85 [1.52-2.36])$ (Table S3). Spectra provided
6	additional information in predicting TTF beyond R-ISS, age at diagnosis, and UAMS (APV = 0.29, LRT
7	$p = 2.6 \times 10^{-6}$) (Table S4).
8	
9	Spectra Descriptive Modeling
10	Figure 7 illustrates associations between the CD138+ spectra and tumor aberrations or
11	demographic groups with elevated myeloma risk. For descriptive modeling a model containing all 39
12	spectra was fit for each characteristic. In logistic regression models, each tumor aberration showed
13	different significant spectra in the model, with some spectra unique to only one aberration. Linear
14	regression with age at diagnosis as a quantitative outcome highlighted associations with thirteen spectra.
15	Seven spectra were significantly associated with gender and thirteen spectra significant with race (self-

16 reported black or white; other racial categories too small to consider).

10

1 DISCUSSION

2 The promise of personalized prevention, management, and treatment is rooted in an ability to describe an individual's unique experience and model important sources of heterogeneity (Ramón v Cajal 3 et al., 2020). Gene expression in diseased tissue is an established source of heterogeneity (Kwa et al., 4 5 2017). Tools that can take a deeper dive and better characterize expression heterogeneity will be important to advance the promise of personalized medicine. Clinical and epidemiologic studies that wish 6 to model multiple sources of risk in a population, transcriptome variables that can be easily incorporated 7 with other variables and across endpoints of interest are advantageous. The goal of this study was to 8 9 provide a technique to derive an agnostic framework of variables for transcriptome data, to empower multivariable studies, and provide novel quantitative molecular phenotypes. The SPECTRA approach 10 identifies quantitative, orthogonal (non-correlated) variables that capture sources of transcriptome 11 variation for use in subsequent predictive or descriptive modeling or as quantitative phenotypes. 12 The SPECTRA approach provides a measured dive into the transcriptome. Each spectrum 13 iteratively moves quantifiably deeper into the variance of the data (measurable by its corresponding λ). 14 15 Methods that iteratively find independent components (PCA and independent component analysis) have previously been shown to provide superior coverage of transcriptome data (Way et al., 2020). The 16 17 importance of a deeper dive is illustrated in the MM predictive modeling. Spectra representing variance 18 deeper in the data were important in predicting OS, PFS, and TTF (Table 1). For example, S32 (which explains 0.5% variance in the transcriptome) provided the same weight to the PFS risk score as S1 (8.7%) 19 variance). As another illustration, S19 and S26 (both with variance < 1%) are in the OS spectra risk score. 20 Further, twelve patients with poor survival -14% of the high-risk group – would not have been classified 21 in the high-risk group if their S19 and S26 spectra values had been at the population mean. In our MM 22 study, retention of components deep in the data, representing small variances (i.e., deep dives) improved 23 prediction. We suggest that deep transcriptome characterization is a new tool with the potential to identify 24 25 the few patients that respond to a drug, or small groups of individuals with large effects in outcome studies, both relevant to precision medicine. Furthermore, as shown previously for the PAM50 gene 26 panel, superior coverage of data may identify previously overlooked expression differences between 27

1	familial and sporadic tissues, identifying potential for familial components (Madsen et al., 2018), and
2	providing new avenues for gene discovery, exposure and gene by environment studies.
3	In the SPECTRA approach negative gene loadings are embraced. Non-negative matrix
4	factorization (NMF), arguably the leading approach in the computational biology field, restricts all values
5	in the amplitude matrix (equivalent to the eigenvector matrix) and pattern matrix (equivalent to the
6	transformed component values in PCA) to be non-negative. However, non-negative matrix values may
7	not be a natural restriction to systems of genes. Genes in a system may act in opposite directions
8	producing both surplus and deficits of that system. With the non-negativity restriction, NMF is limited to
9	the identification of groups of over-expressed genes (Brunet et al., 2004; Stein-O'Brien et al., 2018), and
10	thus models only neutrality and surplus. Deficits may also be important. While PCA spectra may
11	represent mixtures of different biological mechanisms, these may be important combinations, including
12	genes acting in opposite directions, and may better reflect reality. These differences underscore
13	SPECTRA as a valuable and complementary tool to existing approaches.
14	In our MM study, we illustrated the implementation of 39 CD138+ spectra in predictive modeling
15	for OS, PFS, and TTF. We showed added value beyond established risk scores (UAMS, SBUK) and
16	clinical risk factors (R-ISS, age at diagnosis). The framework of 39 quantitative variables has the
17	potential to provide a bridge to compare many patient or tissue characteristics (Figure 7) and could be
18	used to compare existing categorizations (such as subtypes) of patients, even when no genes overlap in
19	their signatures, or they predict different outcomes (Szalat et al., 2016). In this way, spectra provide an
20	alternate to categorical intrinsic subtyping, a well-established practice for many cancers (Dai et al., 2015).
21	Clinically-relevant stratification may be better represented using risk groups within a
22	transcriptome framework (Camp et al., 2019). To this end, we have described a rigorous strategy to define
23	spectra-based risk-groups using GMM (Figure 2, Figure 4, Table 2). Using our approach to identify risk
24	groups for OS, we classified 27 patients as high risk that UAMS classified as low risk (Figure S4). These
25	27 patients had median survival of 26.7 months. Conversely, the 35 patients classified as low risk by

12

SPECTRA outperform existing classification strategies in MM, while also allowing the flexibility of
 modeling with quantitative variables and multiple outcomes.

Descriptive modeling showed significant associations between spectra and tumor cytogenetics, race, gender, and age a diagnosis (**Figure 7**). As expected for a framework of agnostically derived variables, not all spectra are relevant to every dependent variable. Across the five tumor abnormalities, the number of associated spectra ranged from 2 to 15 with no single spectra associated with all abnormalities (**Figure 7**). Importantly, these examples show the flexibility of the framework as well as how it can support comparisons across different models and outcomes.

9 The potential for increased power using spectra variables is illustrated by the discovery of novel associations between spectra and patient demographic risk groups with known differences in incidence 10 (age, gender, race). A prior CoMMpass study whose goal was to identify differences by race for myeloma 11 tumors used the UAMS score to compare transcriptomes by race and did not identify significant 12 differences (p = 0.662) (Manojlovic et al., 2017). In contrast, our framework of 39 CD138+ spectra did 13 identify 13 spectra that differed significantly by race (Figure 7). Our multivariable results demonstrate 14 that significant differences do exist, but also illustrate that the diseased cells in these demographic groups 15 are not distinct entities; fewer than half the spectra variables differ significantly by these patient 16 17 demographic groups. Focusing on the spectra that do show differences by demographics provides new avenues to explore why incidence varies in these groups, a key to disease prevention, intervention, and 18 control. Because transcriptomes capture both the effects of internal (inherited genetics) and external 19 20 factors (lifestyle, exposures, consequences of access to care), these results support epidemiology and 21 biosociology investigations into such differences. We have provided the variable framework (gene 22 transformations) and the spectra variables for the CoMMpass patients (see Data and Code Availability) to enable further study of spectra in other CoMMpass studies, as well as in other myeloma studies. 23 As for any approach, there are limitations. A key question is representation. For epidemiology 24

studies, for example, spectra should ideally be representative of the entire disease population. This
requires that the derivation dataset is a random sample from that population or based on a known selective
sampling scheme. While there are many publicly available transcriptome datasets (Cancer Genome Atlas

13

Research Network et al., 2013; GTEx Consortium, 2013), most fall short of this ideal. Thus, the spectra variables derived from these will have inherent limitations in representation. An investigator should consider if a derivation dataset is adequate to represent their study goals. We note that the goal of the MMRF CoMMpass study was designed to represent myeloma patients from diagnosis through treatment and is the largest existing cohort of treatment naïve CD138+ transcriptomes, with sampling continuing over time. To limit overfitting in spectra derivation, we used dimension reduction to focus only on the first *k* spectra (largest *k* components of variation), selected using a scree test (Cattell, 1966).

A limitation of agnostic variables is interpretation. To illustrate the ability to gain biological 8 9 insight from spectra, we implemented pre-ranked GSEA (Subramanian et al., 2005). Each spectra variable is a linear transformation based on gene weights from its eigenvector. Gene weights can be positive or 10 negative and order the importance and direction of effect of genes in each spectrum, and thus are an ideal 11 ranking metric for enrichment. Spectrum variable S5 was significant in models for OS, PFS and TTF. We 12 used fast-GSEA and identified enriched Hallmark pathways based on genes ranked by CD138+ S5 13 (Korotkevich et al., 2016; Liberzon et al., 2015). Three pathways were highly significant. Of particular 14 interest was the unfolded protein response (UPR) pathway (normalized enrichment score, NES = 2.25, 15 adjusted- $p=2.5 \times 10^{-8}$). Secreted proteins are processed in the endoplasmic reticulum (ER). Incorrectly 16 17 folded proteins create ER stress and activate the UPR pathway, targeting them for degradation by the proteasome. Bortezomib is a proteasome inhibitor and hence causes build-up of misfolded proteins, 18 increasing ER stress leading to cell-death. Myeloma cells secrete large amounts of incorrectly folded 19 20 immunoglobulins and therefore are near capacity for UPR. Due to this, CD138+ myeloma cells are 21 particularly sensitive to proteasome inhibitors, such as Bortezomib, a common agent in the initial treatment for MM. Thus, spectrum S5 may represent a patient's sensitivity to proteasome inhibition 22 agents, influencing OS, PFS and TTF. 23

Data quality and processing are paramount to derive informative agnostic variables. PCA is a procedure that provides linear transformations of the data to represent variance. If the data have technical artifacts, batch effects, unstable or non-comparable expression measures, the noise can overwhelm authentic variance. Accordingly, the SPECTRA approach intentionally includes strict quality control,

14

1 careful zero-handling, robust normalization, and batch correction (Figure 2). Without these steps, PCA can fail to provide informative variables. An agnostic approach permits stringent data culling because the 2 incentive to retain features based on known functional relevance is removed. The impetus is to only retain 3 features that can contribute to meaningful variance and provide informative variables for modeling. The 4 5 limitation of an agnostic approach is reduced biological interpretation or mechanistic insight of the variables themselves, prior to modeling. However, there are already many approaches that take this 6 alternate goal of intermediate interpretation (Subramanian et al., 2005), whose limitations are instead the 7 reduced flexibility of the variables they produce. Hence, SPECTRA is a complementary approach to the 8 9 current toolset available for all fields. In conclusion, we present a new approach, SPECTRA, to derive an agnostic transcriptome 10 framework of quantitative, orthogonal variables for a dataset. These multi-gene expression variables are 11 designed specifically to capture transcriptome variation, providing variables for flexible modeling, along 12 with other covariates, to better differentiate individuals for any outcome. Spectra may also be used as 13 novel transcriptome phenotypes. Applied to CD138+ transcriptomes for myeloma patients, we defined 14 CD138+ spectra and implemented these in many different outcome models. We illustrated an ability to 15 predict prognostic outcomes and provide new insight into potential differences between tumors and 16 17 patients from demographic groups. Fundamentally, the technique shifts from characterizing a transcriptome using categories to multiple quantitative measures. The SPECTRA approach provides a 18

19 new paradigm and tool for exploring transcriptomes that hold promise for discoveries to advance

20 precision screening, prevention, intervention, and survival studies.

15

1 MATERIALS AND METHODS

2 <u>SPECTRA approach</u>

Our goal was derivation of well-behaved, quantitative variables from RNAseq data to capture the many sources of intrinsic variation in a transcriptome. To represent meaningful variance and enable deep comparison across individuals we model well-behaved genes after normalization and batch correction (Figure 2).

Features in the transcriptome likely to be unduly influenced by poor alignment or lacking precision due to sequencing depth are removed as potentials for introducing spurious and unstable variation. We concentrate on autosomal protein-coding genes. Features with inadequate data for precision, defined as more than 5% of samples with fewer than 100 read counts, are removed. We also remove genes known to be unstable across different RNAseq pipelines (Arora et al., 2020). After the removal of features, individuals are removed from consideration if more than 10% of the remaining genes have fewer than 100 read counts.

Normalization is required for comparisons across genes and individuals and includes adjustment for gene length, sequencing depth (library size), and RNA composition. We use a robust internal (single sample) normalization to obviate the need for a 'reference' sample and to provide the possibility for portability across datasets. The SPECTRA approach is gene-focused, but the workflow handles both gene-based and transcript-based alignment and quantification; the latter has been suggested to be more accurate (Zhao et al., 2015). Normalized gene expression values, e_g , are calculated as follows:

20
$$\mathbf{e}_g = \log_2\left(\frac{\sum_{t=1}^m \frac{\mathbf{c}_t + 1/m}{\mathbf{l}_t}}{\operatorname{median}\left(\sum_{t=1}^m \frac{\mathbf{c}_t + 1/m}{\mathbf{l}_t}\right)}\right)$$

where c_t is the read count for transcript *t*, l_t is the transcript length in kilobases (extracted from the GTF used to align and quantify the RNAseq data), and *m* is the number of transcripts for the gene. Zerohandling is achieved by adding 1/m to the transcript counts ($c_t + 1/m$). Division by l_t corrects for transcript length. Summing the length-corrected transcript counts results in a gene-level count per kilobase (CPK) measure. Gene-level read counts are equivalent to m = 1. Adjustments for sequencing depth and RNA composition (often referred to as the *size factor*) are achieved via the division of each

16

1	gene-based CPK measure by the median of CPK-values for retained features. We note that the more usual
2	upper-quartile adjustment also provides robust internal normalization (Shahriyari, 2019); however, since
3	our implementation is post-QC after numerous features have been removed for low counts, the median is
4	more suitable. Normalized data are log_2 transformed to account for skew. We also truncate outliers to a
5	five standard deviation threshold from the mean of the normalized gene counts in the relevant direction.
6	Batch correction is necessary to correct for potential technical artifacts and spurious variation
7	introduced by differing sequencing protocols. We adjust for sequence batch using ComBat (Johnson et
8	al., 2007) as implemented in the SVA R package (Leek et al., 2012), with covariates included for patient
9	characteristics that are unbalanced by batch.
10	PCA is implemented with the covariance matrix. We use singular value decomposition to perform
11	PCA, and it is necessary to center the expression values first to ensure the MF is performed for the
12	covariance. Expression values (e_g) are centered on the mean across individuals for gene g. The R core
13	function prcomp is used to perform PCA. Eigenvectors contain the gene loadings that define the linear
14	transformation used for each spectrum. A corresponding eigenvalue, λ , indicates the proportion of the
15	global variance represented by the transformed value for an eigenvector. We use a scree test(Cattell,
16	1966) (inflection point of the rank-ordered plot of the λ , or elbow method) to select the k components to
17	retain. The proportion of variance explained by this k-dimensional space $(\sum_{i=1}^{k} \lambda_i / \sum_{\forall i} \lambda_i)$ indicates the
18	depth of the dive into the transcriptome data. Spectra variables are the standardized retained components
19	$(S_1,, S_k).$

- 19
- 20

Myeloma CD138+ Spectra 21

22 Data from the MMRF CoMMpass Study (release IA14) (Keats et al., 2013) were downloaded from the MMRF web portal (https://research.themmrf.org/). Clinical data and CD138+ RNAseq were 23 available for 768 patients prior to treatment at study entry (baseline) and 119 follow-up bone marrow 24 samples. Transcript-based expression estimates processed by Salmon (version 0.7.2) were used. The 25 26 SPECTRA approach was used on baseline samples (n=768) to derive a CD138+ transcriptome framework and SPECTRA variables. While not used for spectra derivation, follow-up samples were batch-corrected 27

1	alongside the baseline samples. Covariates included in batch correction were age, gender, overall
2	survival, progression-free survival, time to first-line treatment failure, and treatment status. The first 39
3	components (spectra variables S1,, S39). were selected based on the scree test (Cattell, 1966).
4	
5	Predictive Modeling of Clinical Outcomes
6	Predictive Modeling used Cox proportional hazards regression, implemented in the survival
7	package in R (Therneau & Grambsch, 2000; Therneau, 2021). An overview of predictive modeling is
8	provided in Figure S1. All 39 spectra were considered as predictor variables for each of the three survival
9	outcomes: overall survival (OS), progression-free survival (PFS), and time to first-line treatment failure
10	(TTF). The spectra were standardized in all analyses. For each outcome, a single step variable selection
11	from the all-spectra model was performed to retain only significant spectra (p<0.05). The linear predictor
12	from the reduced model was used to define a quantitative spectra risk score for the outcome (i.e., the
13	weighted sum of the spectra retained in the model based on their coefficients: $\sum_{j} \beta_{j} S_{j}$).
14	To determine if a quantitative spectra risk score displayed evidence of latent risk groups that
15	could be clinically meaningful, the risk score distribution and normal Quantile-Quantile plots were
16	inspected. The mclust R package was used to statistically assess evidence for risk groups using density
17	estimation by Gaussian finite Mixture Modeling (GMM) (Fraley & Raftery, 2002; Scrucca et al., 2016)
18	assuming equal variances. Bayesian information content (BIC) and bootstrap likelihood ratio tests (LRT)
19	were used to determine the number of risk groups (mclustBIC, mclustBootstrapLRT). The hazard ratio
20	(HR) between risk groups was calculated using relative event rates (Armitage et al., 2002, p. 578) from
21	survdiff in the survival package in R (Therneau, 2021).
22	To address overfitting and lack of a replication sample, bootstrap internal validation was used
23	(Harrell, 2015, pp. 114–116). Bootstrap internal validation involves replicating the entire modeling
24	process on bootstrap resamples of the data to determine and adjust for over-fitting. Here, the process
25	includes both the single step variable selection, GMM density estimation for risk group designation, and
26	adjustment of HRs, C-index and bootstrap confidence limits (Noma et al., 2021). Initial estimates from a

27 model are referred to as 'apparent'. The measure of overfitting is referred to as 'optimism'. The corrected

18

estimates are referred to as 'adjusted' values. Adjusted estimates are more reasonable assessment of effect
 size. If adjusted confidence intervals do not contain 1.0, this indicates validation.

Comparison with existing transcriptome risk scores. We compared spectra-based risk to two 3 previously published risk scores: 1) the first and most widely adopted supervised expression risk score in 4 5 myeloma, from the University of Arkansas for Medical Sciences (UAMS) (Shaughnessy et al., 2007); and 2) a recent supervised risk score, from the Shahid Bahonar University of Kerman (SBUK) (Zamani-6 Ahmadmahmudi et al., 2020). The UAMS risk score was developed in microarray data and tested 54,657 7 probes for association with disease-related survival (Shaughnessy et al., 2007). A total of 70 genes were 8 9 selected (19 under and 51 overexpressed prognostic genes). The UAMS risk score is the ratio of mean expression of the up-regulated to down-regulated genes, with k-means clustering to determine a cutoff for 10 'high-risk' classification (Shaughnessy et al., 2007). The SBUK prognosis score was developed in the 11 CoMMpass RNAseq data. All genes were tested for association with survival, followed by a multi-step 12 process, including univariate Cox analysis, the intersection in six Class Prediction algorithms, and 13 multivariate Cox analysis, was used to select 17 genes consistent across multiple methods (Zamani-14 Ahmadmahmudi et al., 2020). These 17 were entered in a multivariable Cox regression to define the 15 SBUK score, with the 75th percentile used to classify patients into high-risk and low-risk categories 16 17 (Zamani-Ahmadmahmudi et al., 2020). We calculated each patient's UAMS and SBUK risk scores and their risk status (low or high) using the CoMMpass RNAseq data. 18

Added predictive value. Predictive value beyond established clinical risk factors is important to 19 establish. We calculated the added predictive value (APV) (Al-Radi et al., 2007; Califf et al., 1985; 20 21 Harrell, 2015) of each expression risk score (spectra, UAMS, SBUK) beyond risk predicted from established clinical predictors. We used quantitative risk scores, and not risk group status, as these contain 22 more complete risk information and because each gene-expression method employed a different 23 thresholding approach making risk status comparisons potentially misleading. The APV compares the 24 25 LRT statistic for a model including established clinical predictors to a model which also includes each expression risk score. We also determined whether the spectra risk score provided value beyond known 26 clinical predictors and UAMS. APV=0.0 indicates no additional prediction information, i.e., all predictive 27

1	power is duplicative with other factors already in the model. With larger independent predictive
2	contributions, APV increases, to a maximum of 1.0.
3	
4	CD138+ Spectra in Follow-Up Samples
5	To illustrate the potential to track changes over time, OS spectra risk scores were calculated in
6	longitudinal follow-up samples. For these samples, normalized and batch-corrected gene expression
7	values were centered to the derivation set. Spectra values (S_j) were determined using the previously
8	derived linear transformation and scaled to the derivation set. The OS risk score was calculated using the
9	previously established linear combination across spectra ($\sum_{j} \beta_{j} S_{j}$). Spectra variables in all samples
10	(baseline and follow-up) are provided (see Data and Code Availability).
11	
12	Descriptive Modeling of Clinical and Demographic factors
13	We used descriptive modeling to illustrate associations between the CD138+ spectra and clinical
14	and demographic factors with elevated myeloma risk: tumor aberrations (high risk del(17p), t(14;16), gain
15	1q, and t(4;14), and standard risk t(11;14)) (Rajkumar & Kumar, 2020), age, gender, race. Linear or
16	logistic regression was used with all 39 CD138+ spectra modeled as independent variables. Spectra were
17	noted if significant in the descriptive model ($p < 0.05$).
18	
19	Data and Code Availability
20	Processed RNAseq data from the CoMMpass Study can be downloaded from
21	https://research.themmrf.org/. Code used to derive the CD138+ transcriptome spectra and generate the
22	myeloma results is freely available on GitHub: <u>https://github.com/njcamp-lab/MM_spectra</u> . Details of the
23	QC process and the transcriptome framework (linear equations for the gene transformations) necessary to
24	calculate the 39-spectra variables in other studies are also provided. Spectra for individuals in the IA14
25	CoMMpass data can be downloaded from <u>https://github.com/njcamp-lab/MM_spectra/SpectraData</u> .

20

1 ACKNOWLEDGMENTS

- 2 The research reported in this publication was supported by the National Cancer Institute (Award Numbers
- 3 F99CA234943, K00CA234943, K07CA230150, and P30CA042014-29S9), the National Center for
- 4 Advancing Translational Sciences (Award Number UL1TR002538), and the National Library of
- 5 Medicine (Award Number T15LM007124) of the National Institutes of Health. The content is solely the
- 6 responsibility of the authors and does not necessarily represent the official views of the National Institutes
- 7 of Health.

21

1 COMPETING INTERESTS

2 None.

REFERENCES

- Allott, E. H., Shan, Y., Chen, M., Sun, X., Garcia-Recio, S., Kirk, E. L., Olshan, A. F., Geradts, J., Earp, H. S., Carey, L. A., Perou, C. M., Pfeiffer, R. M., Anderson, W. F., & Troester, M. A. (2020).
 Bimodal age distribution at diagnosis in breast cancer persists across molecular and genomic classifications. *Breast Cancer Research and Treatment*, *179*(1), 185–195. https://doi.org/10.1007/s10549-019-05442-2
- Al-Radi, O. O., Harrell, F. E., Caldarone, C. A., McCrindle, B. W., Jacobs, J. P., Williams, M. G., Van Arsdell, G. S., & Williams, W. G. (2007). Case complexity scores in congenital heart surgery: A comparative study of the Aristotle Basic Complexity score and the Risk Adjustment in Congenital Heart Surgery (RACHS-1) system. *The Journal of Thoracic and Cardiovascular Surgery*, *133*(4), 865–875. https://doi.org/10.1016/j.jtcvs.2006.05.071
- Armitage, P., Berry, G., & Matthews, J. N. S. (Eds.). (2002). Statistical Methods in Medical Research. Blackwell Science Ltd. https://doi.org/10.1002/9780470773666
- Arora, S., Pattwell, S. S., Holland, E. C., & Bolouri, H. (2020). Variability in estimated gene expression among commonly used RNA-seq pipelines. *Scientific Reports*, 10(1), 2734. https://doi.org/10.1038/s41598-020-59516-z
- Brunet, J.-P., Tamayo, P., Golub, T. R., & Mesirov, J. P. (2004). Metagenes and molecular pattern discovery using matrix factorization. *Proceedings of the National Academy of Sciences of the United States of America*, 101(12), 4164–4169. https://doi.org/10.1073/pnas.0308531101
- Califf, R. M., Phillips, H. R., Hindman, M. C., Mark, D. B., Lee, K. L., Behar, V. S., Johnson, R. A., Pryor, D. B., Rosati, R. A., Wagner, G. S., & Harrell, F. E. (1985). Prognostic value of a coronary artery jeopardy score. *Journal of the American College of Cardiology*, 5(5), 1055–1063. https://doi.org/10.1016/S0735-1097(85)80005-X
- Camp, N. J., Madsen, M. J., Herranz, J., Rodríguez-Lescure, A., Ruiz, A., Martín, M., & Bernard, P. S. (2019). 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 Research and Treatment*, 175(1), 129–139. https://doi.org/10.1007/s10549-018-05097-5

- Cancer Genome Atlas Research Network, Weinstein, J. N., Collisson, E. A., Mills, G. B., Shaw, K. R. M., Ozenberger, B. A., Ellrott, K., Shmulevich, I., Sander, C., & Stuart, J. M. (2013). The Cancer Genome Atlas Pan-Cancer analysis project. *Nature Genetics*, 45(10), 1113–1120. https://doi.org/10.1038/ng.2764
- Cattell, R. B. (1966). The Scree Test For The Number Of Factors. *Multivariate Behavioral Research*, *1*(2), 245–276. https://doi.org/10.1207/s15327906mbr0102_10
- Dai, X., Li, T., Bai, Z., Yang, Y., Liu, X., Zhan, J., & Shi, B. (2015). Breast cancer intrinsic subtype classification, clinical use and future trends. *American Journal of Cancer Research*, 5(10), 2929– 2943.
- Fraley, C., & Raftery, A. E. (2002). Model-Based Clustering, Discriminant Analysis, and Density Estimation. *Journal of the American Statistical Association*, 97(458), 611–631. https://doi.org/10.1198/016214502760047131
- GTEx Consortium. (2013). The Genotype-Tissue Expression (GTEx) project. *Nature Genetics*, 45(6), 580–585. https://doi.org/10.1038/ng.2653
- Hanson, H. A., Leiser, C. L., Madsen, M. J., Gardner, J., Knight, S., Cessna, M., Sweeney, C., Doherty, J. A., Smith, K. R., Bernard, P. S., & Camp, N. J. (2020). 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*(4), 807–815. https://doi.org/10.1158/1055-9965.EPI-19-0912

- Harrell, J. (2015). Regression Modeling Strategies: With Applications to Linear Models, Logistic and Ordinal Regression, and Survival Analysis (2nd ed. 2015). Springer International Publishing : Imprint: Springer. https://doi.org/10.1007/978-3-319-19425-7
- Johnson, W. E., Li, C., & Rabinovic, A. (2007). Adjusting batch effects in microarray expression data using empirical Bayes methods. *Biostatistics*, 8(1), 118–127. https://doi.org/10.1093/biostatistics/kxj037
- Keats, J. J., Craig, D. W., Liang, W., Venkata, Y., Kurdoglu, A., Aldrich, J., Auclair, D., Allen, K., Harrison, B., Jewell, S., Kidd, P. G., Correll, M., Jagannath, S., Siegel, D. S., Vij, R., Orloff, G., Zimmerman, T. M., Mmrf CoMMpass Network, Capone, W., ... Lonial, S. (2013). Interim Analysis Of The Mmrf Commpass Trial, a Longitudinal Study In Multiple Myeloma Relating Clinical Outcomes To Genomic and Immunophenotypic Profiles. *Blood*, *122*(21), 532–532. https://doi.org/10.1182/blood.V122.21.532.532
- Korotkevich, G., Sukhov, V., Budin, N., Shpak, B., Artyomov, M. N., & Sergushichev, A. (2016). Fast gene set enrichment analysis [Preprint]. Bioinformatics. https://doi.org/10.1101/060012
- Kwa, M., Makris, A., & Esteva, F. J. (2017). Clinical utility of gene-expression signatures in early stage breast cancer. *Nature Reviews. Clinical Oncology*, 14(10), 595–610. https://doi.org/10.1038/nrclinonc.2017.74
- Lapointe, J., Li, C., Higgins, J. P., van de Rijn, M., Bair, E., Montgomery, K., Ferrari, M., Egevad, L., Rayford, W., Bergerheim, U., Ekman, P., DeMarzo, A. M., Tibshirani, R., Botstein, D., Brown, P. O., Brooks, J. D., & Pollack, J. R. (2004). Gene expression profiling identifies clinically relevant subtypes of prostate cancer. *Proceedings of the National Academy of Sciences*, *101*(3), 811–816. https://doi.org/10.1073/pnas.0304146101
- Leek, J. T., Johnson, W. E., Parker, H. S., Jaffe, A. E., & Storey, J. D. (2012). The sva package for removing batch effects and other unwanted variation in high-throughput experiments. *Bioinformatics*, 28(6), 882–883. https://doi.org/10.1093/bioinformatics/bts034

- Liberzon, A., Birger, C., Thorvaldsdóttir, H., Ghandi, M., Mesirov, J. P., & Tamayo, P. (2015). The Molecular Signatures Database (MSigDB) hallmark gene set collection. *Cell Systems*, 1(6), 417– 425. https://doi.org/10.1016/j.cels.2015.12.004
- López, C., Kleinheinz, K., Aukema, S. M., Rohde, M., Bernhart, S. H., Hübschmann, D., Wagener, R., Toprak, U. H., Raimondi, F., Kreuz, M., Waszak, S. M., Huang, Z., Sieverling, L., Paramasivam, N., Seufert, J., Sungalee, S., Russell, R. B., Bausinger, J., Kretzmer, H., ... Siebert, R. (2019). Genomic and transcriptomic changes complement each other in the pathogenesis of sporadic Burkitt lymphoma. *Nature Communications*, *10*(1), 1459. https://doi.org/10.1038/s41467-019-08578-3
- Madsen, M. J., Knight, S., Sweeney, C., Factor, R., Salama, M., Stijleman, I. J., Rajamanickam, V.,
 Welm, B. E., Arunachalam, S., Jones, B., Rachamadugu, R., Rowe, K., Cessna, M. H., Thomas,
 A., Kushi, L. H., Caan, B. J., Bernard, P. S., & Camp, N. J. (2018). Reparameterization of
 PAM50 Expression Identifies Novel Breast Tumor Dimensions and Leads to Discovery of a
 Genome-Wide Significant Breast Cancer Locus at *12q15*. *Cancer Epidemiology Biomarkers & Prevention*, *27*(6), 644–652. https://doi.org/10.1158/1055-9965.EPI-17-0887
- Manojlovic, Z., Christofferson, A., Liang, W. S., Aldrich, J., Washington, M., Wong, S., Rohrer, D., Jewell, S., Kittles, R. A., Derome, M., Auclair, D., Craig, D. W., Keats, J., & Carpten, J. D. (2017). Comprehensive molecular profiling of 718 Multiple Myelomas reveals significant differences in mutation frequencies between African and European descent cases. *PLOS Genetics*, *13*(11), e1007087. https://doi.org/10.1371/journal.pgen.1007087
- Noma, H., Shinozaki, T., Iba, K., Teramukai, S., & Furukawa, T. A. (2021). Confidence intervals of prediction accuracy measures for multivariable prediction models based on the bootstrap-based optimism correction methods. *ArXiv:2005.01457 [Stat]*. http://arxiv.org/abs/2005.01457

Patro, R., Duggal, G., Love, M. I., Irizarry, R. A., & Kingsford, C. (2017). Salmon provides fast and biasaware quantification of transcript expression. *Nature Methods*, 14(4), 417–419. https://doi.org/10.1038/nmeth.4197

Perou, C. M., Sørlie, T., Eisen, M. B., van de Rijn, M., Jeffrey, S. S., Rees, C. A., Pollack, J. R., Ross, D. T., Johnsen, H., Akslen, L. A., Fluge, Ø., Pergamenschikov, A., Williams, C., Zhu, S. X., Lønning, P. E., Børresen-Dale, A.-L., Brown, P. O., & Botstein, D. (2000). Molecular portraits of human breast tumours. *Nature*, 406(6797), 747–752. https://doi.org/10.1038/35021093

Rajkumar, S. V., & Kumar, S. (2020). Multiple myeloma current treatment algorithms. *Blood Cancer Journal*, 10(9), 94. https://doi.org/10.1038/s41408-020-00359-2

Ramón y Cajal, S., Sesé, M., Capdevila, C., Aasen, T., De Mattos-Arruda, L., Diaz-Cano, S. J.,
Hernández-Losa, J., & Castellví, J. (2020). Clinical implications of intratumor heterogeneity:
Challenges and opportunities. *Journal of Molecular Medicine*, *98*(2), 161–177.
https://doi.org/10.1007/s00109-020-01874-2

- Scrucca, L., Fop, M., Murphy, T. B., & Raftery, A. E. (2016). mclust 5: Clustering, Classification and Density Estimation Using Gaussian Finite Mixture Models. *The R Journal*, 8(1), 289–317.
- SEER. (2021). *Myeloma—Cancer Stat Facts https://seer.cancer.gov/statfacts/html/mulmy.html*. SEER. https://seer.cancer.gov/statfacts/html/mulmy.html
- Shahriyari, L. (2019). 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(3), 985–994. https://doi.org/10.1093/bib/bbx153
- Shaughnessy, J. D., Zhan, F., Burington, B. E., Huang, Y., Colla, S., Hanamura, I., Stewart, J. P.,
 Kordsmeier, B., Randolph, C., Williams, D. R., Xiao, Y., Xu, H., Epstein, J., Anaissie, E.,
 Krishna, S. G., Cottler-Fox, M., Hollmig, K., Mohiuddin, A., Pineda-Roman, M., ... Barlogie, B.
 (2007). A validated gene expression model of high-risk multiple myeloma is defined by

27

deregulated expression of genes mapping to chromosome 1. *Blood*, *109*(6), 2276–2284. https://doi.org/10.1182/blood-2006-07-038430

- Stein-O'Brien, G. L., Arora, R., Culhane, A. C., Favorov, A. V., Garmire, L. X., Greene, C. S., Goff, L. A., Li, Y., Ngom, A., Ochs, M. F., Xu, Y., & Fertig, E. J. (2018). Enter the Matrix: Factorization Uncovers Knowledge from Omics. *Trends in Genetics: TIG*, 34(10), 790–805. https://doi.org/10.1016/j.tig.2018.07.003
- Stopsack, K. H., Ebot, E. M., Downer, M. K., Gerke, T. A., Rider, J. R., Kantoff, P. W., & Mucci, L. A. (2018). Regular aspirin use and gene expression profiles in prostate cancer patients. *Cancer Causes & Control: CCC*, 29(8), 775–784. https://doi.org/10.1007/s10552-018-1049-5
- Subramanian, A., Tamayo, P., Mootha, V. K., Mukherjee, S., Ebert, B. L., Gillette, M. A., Paulovich, A., Pomeroy, S. L., Golub, T. R., Lander, E. S., & Mesirov, J. P. (2005). Gene set enrichment analysis: A knowledge-based approach for interpreting genome-wide expression profiles. *Proceedings of the National Academy of Sciences of the United States of America*, 102(43), 15545–15550. https://doi.org/10.1073/pnas.0506580102
- Sweeney, C., Bernard, P. S., Factor, R. E., Kwan, M. L., Habel, L. A., Quesenberry, C. P., Shakespear, K., Weltzien, E. K., Stijleman, I. J., Davis, C. A., Ebbert, M. T. W., Castillo, A., Kushi, L. H., & Caan, B. J. (2014). Intrinsic subtypes from PAM50 gene expression assay in a population-based breast cancer cohort: Differences by age, race, and tumor characteristics. *Cancer Epidemiology, Biomarkers & Prevention: A Publication of the American Association for Cancer Research, Cosponsored by the American Society of Preventive Oncology, 23(5), 714–724.* https://doi.org/10.1158/1055-9965.EPI-13-1023
- Szalat, R., Avet-Loiseau, H., & Munshi, N. C. (2016). Gene Expression Profiles in Myeloma: Ready for the Real World? *Clinical Cancer Research*, 22(22), 5434–5442. https://doi.org/10.1158/1078-0432.CCR-16-0867

- Tamayo, P., Slonim, D., Mesirov, J., Zhu, Q., Kitareewan, S., Dmitrovsky, E., Lander, E. S., & Golub, T.
 R. (1999). Interpreting patterns of gene expression with self-organizing maps: Methods and application to hematopoietic differentiation. *Proceedings of the National Academy of Sciences of the United States of America*, 96(6), 2907–2912. https://doi.org/10.1073/pnas.96.6.2907
- Terry M. Therneau & Patricia M. Grambsch. (2000). *Modeling Survival Data: Extending the Cox Model*. Springer.
- Therneau, T. M. (2021). A Package for Survival Analysis in R. https://CRAN.Rproject.org/package=survival
- Way, G. P., Zietz, M., Rubinetti, V., Himmelstein, D. S., & Greene, C. S. (2020). Compressing gene expression data using multiple latent space dimensionalities learns complementary biological representations. *Genome Biology*, 21(1), 109. https://doi.org/10.1186/s13059-020-02021-3
- Zamani-Ahmadmahmudi, M., Nassiri, S. M., & Soltaninezhad, F. (2020). Development of a RNA sequencing-based prognostic gene signature in multiple myeloma. *British Journal of Haematology*, bjh.16744. https://doi.org/10.1111/bjh.16744
- Zhang, M., Lykke-Andersen, S., Zhu, B., Xiao, W., Hoskins, J. W., Zhang, X., Rost, L. M., Collins, I.,
 Bunt, M. van de, Jia, J., Parikh, H., Zhang, T., Song, L., Jermusyk, A., Chung, C. C., Zhu, B.,
 Zhou, W., Matters, G. L., Kurtz, R. C., ... Amundadottir, L. (2018). Characterising cis-regulatory
 variation in the transcriptome of histologically normal and tumour-derived pancreatic tissues. *Gut*, 67(3), 521–533. https://doi.org/10.1136/gutjnl-2016-313146
- Zhao, S., Xi, L., & Zhang, B. (2015). Union Exon Based Approach for RNA-Seq Gene Quantification: To Be or Not to Be? *PloS One*, 10(11), e0141910. https://doi.org/10.1371/journal.pone.0141910

TABLES

Table 1. Cox regression beta coefficients for overall survival (OS), progression free survival (PFS), and time to treatment failure (TTF). Beta coefficients are by per spectra standard deviation. The quantitative spectra risk score for an outcome is the weighted sum of the spectra retained in its model based on their beta coefficients: $\sum_{j} \beta_{j} S_{j}$.

Spectra	OS	PFS	TTF
1		-0.132	-0.218
2	-0.335	-0.241	-0.212
3	0.243		
4	0.461	0.282	0.182
5	-0.324	-0.271	-0.257
9	-0.335	-0.182	-0.176
10			-0.138
11		0.083	
12	-0.179		
13	0.174		
18			0.131
19	0.111		
24		0.178	0.118
26	0.184	0.101	
31			0.102
32		0.121	0.117

Table 2. Overall survival hazard ratios between high- and low-risk patients using spectra, UAMS, or

SBUK to define patient groups.

Risk Score	isk Score Apparent HR Adjusted HR		Optimism
OS Spectra	6.44 (4.47-14.86)	4.27 (2.31-12.69)	2.16
UAMS	3.98 (2.62-5.54)	3.89 (2.53-5.45)	0.09
SBUK	2.58 (1.95-3.51)	2.54 (1.91-3.47)	0.04

Variable	Base Model	ΔAIC	ΔBIC	APV	LRT <i>P</i> Value
UAMS	R-ISS + Age	-71.7	-68.9	0.576	9.20x10 ⁻¹⁸
SBUK	R-ISS + Age	-52.7	-49.9	0.502	1.42×10^{-13}
OS Spectra	R-ISS + Age	-83.5	-80.7	0.612	2.36x10 ⁻²⁰
UAMS	R-ISS + Age + SBUK	-17.8	-15.0	0.154	8.58x10 ⁻⁶
UAMS	R-ISS + Age + OS Spectra	-1.9	0.8	0.027	0.048
SBUK	R-ISS + Age + UAMS	1.2	4.0	0.006	0.371
SBUK	R-ISS + Age + OS Spectra	1.1	3.8	0.007	0.333
OS Spectra	R-ISS + Age + UAMS	-13.7	-10.9	0.109	7.39x10 ⁻⁵
OS Spectra	R-ISS + Age + SBUK	-29.7	-27.0	0.226	1.77x10 ⁻⁸
OS Spectra	R-ISS + Age + UAMS + SBUK	-12.9	-10.2	0.104	1.13x10 ⁻⁴

Table 3. Overall survival Cox models including covariates and added predictive values.

<u>Key</u>

R-ISS: Revised international staging system

Age at diagnosis

 Δ AIC: Change in Akaike's Information Criterion by adding the variable

 Δ BIC: Change in Bayesian information criterion by adding the variable

APV: Variable added predictive value of the variable

LRT P Value: Likelihood ratio test p-value between the base model with and without the variable

32

FIGURES



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 spectra 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 framework to derive spectra variables and example applications

using spectra variables.



Figure 3. Spectra barcodes in four CoMMpass patients. A) Percent of transcriptome-wide variance between samples captured by each spectrum. Together the spectra capture 65% of the transcriptome-wide variance. B) For each patient, all 39 spectra are illustrated with the value represented by the bar intensity. The color indicates if the patient's spectra value is positive (blue) or negative (green). 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. For spectrum S1, 2497 is more like 1392, for spectrum S15, 1854 is more like 2394.



Figure 4. Overall survival spectra predictive modeling. A) Quantile-quantile plot of the actual spectra overall survival risk score (Cox linear predictors) compared to theoretical values. Elevated tail suggests high-risk group exists. B) Gaussian mixture modeling identified two distributions in the spectra score. Each patient was assigned into an OS risk group based on distribution probability. The blue bars on the x-axis represent patients in the low-risk group and the orange bars are patients in the high-risk group. C) Kaplan-Meier curve of the high/low risk groups. High-risk patients (n = 88, 58 events) had median

survival of 20.1 months. Low-risk patients (n = 679, 121 events) did not reach median survival in the study timeframe.

Figure 5. Aggregate effect of CD138+ spectra for overall survival over time. Spectra scores for overall survival are shown in eleven patients with RNAseq data at multiple time points. Diamonds indicate sequencing events and show the spectra score at that timepoint. The final narrow rectangle indicates timepoint the patient died (filled) or was last known alive (open). Gray horizontal line shows the high/low risk group threshold.

Figure 6. Kaplan-Meier curves of progression free survival and time to first line treatment failure. A) PFS high-risk patients (n = 60, 50 events) had median survival of 9.7 months. PFS low-risk patients (n

months

risk group 🕂 low 🕂 high

ò

= 707, 342 events) had median survival of 35.7 months. B) Patients in the spectra high-risk TTF (n = 31, 25 events) had median TTF of 9.2 months compared to low-risk patients (n = 736, 344 events) with median TTF of 32.8 months.

42

SUPPLEMENTAL METHODS

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

$$X = TQ^T$$
 Equation 1

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

$$\Sigma = Q \Lambda Q^T \qquad Equation 2$$

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

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

$$T_k = XQ_k$$
 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 $(\sum_{s=1}^{k} \lambda_s / \sum_{\forall s} \lambda_s)$ can be used as a measure of coverage.

43

SUPPLEMENTAL TABLES

Risk Score	Apparent HR	Adjusted HR	Optimism
PFS Spectra	4.22 (2.81-9.59)	3.08 (1.67-8.44)	1.14
UAMS	2.44 (1.74-3.35)	2.40 (1.70-3.31)	0.04
SBUK	1.94 (1.61-2.45)	1.92 (1.59-2.43)	0.01

Table S1. Hazard ratios of high vs low risk of progression free survival.

Variable	Base Model	ΔΑΙΟ	ΔBIC	APV	LRT P Value
UAMS	R-ISS + Age	-49.6	-46.1	0.566	6.64x10 ⁻¹³
SBUK	R-ISS + Age	-54.7	-51.1	0.589	5.12×10^{-14}
PFS Spectra	R-ISS + Age	-77.7	-74.2	0.668	4.28x10 ⁻¹⁹
UAMS	R-ISS + Age + SBUK	-1.9	1.7	0.039	0.049
UAMS	R-ISS + Age + PFS Spectra	1.3	4.9	0.005	0.419
SBUK	R-ISS + Age + UAMS	-6.9	-3.4	0.089	0.003
SBUK	R-ISS + Age + PFS Spectra	-1.3	2.3	0.027	0.070
PFS Spectra	R-ISS + Age + UAMS	-26.7	-23.2	0.240	8.26x10 ⁻⁸
PFS Spectra	R-ISS + Age + SBUK	-24.3	-20.8	0.215	2.87x10 ⁻⁷
PFS Spectra	R-ISS + Age + UAMS + SBUK	-20.5	-16.99	0.184	2.07x10 ⁻⁶

Table S2. Progression free survival covariates and added predictive values.

Risk Score	Apparent HR	Adjusted HR	Optimism
TTF Spectra	3.93 (2.13-6.28)	3.10 (1.31-5.46)	0.82
UAMS	2.00 (1.50-2.70)	1.98 (1.48-2.68)	0.02
SBUK	1.86 (1.52-2.36)	1.85 (1.52-2.36)	0.01

Table S3. Hazard ratios of high vs low risk early treatment failure.

Variable	Base Model	ΔΑΙΟ	Δ BIC	APV	LRT <i>P</i> Value
UAMS	R-ISS + Age	-33.5	-30.1	0.478	2.49x10 ⁻⁹
SBUK	R-ISS + Age	-40.6	-37.1	0.523	6.73x10 ⁻¹¹
TTF Spectra	R-ISS + Age	-55.1	-51.6	0.595	4.10×10^{-14}
UAMS	R-ISS + Age + SBUK	0.6	4.1	0.017	0.240
UAMS	R-ISS + Age + TTF Spectra	1.5	5.0	0.005	0.474
SBUK	R-ISS + Age + UAMS	-6.4	-2.9	0.102	0.004
SBUK	R-ISS + Age + TTF Spectra	-0.4	3.1	0.024	0.124
TTF Spectra	R-ISS + Age + UAMS	-20.1	-16.6	0.229	2.61x10 ⁻⁶
TTF Spectra	R-ISS + Age + SBUK	-14.9	-11.4	0.172	3.96x10 ⁻⁵
TTF Spectra	R-ISS + Age + UAMS + SBUK	-13.6	-10.1	0.158	7.98x10 ⁻⁵

Table S4. Time to first-line treatment failure covariates and added predictive values.

SUPPLEMENTAL FIGURES

Figure S1. Diagram of predictive modeling.

50

risk group (spectra–UAMS) 🕂 low–low 🕂 high–high 🕂 low–high 🕂 high–low

Figure S4. Kaplan-Meier curves of overall survival by spectra and UAMS risk-groups. Spectra

identified patients as high-risk with worse survival that UAMS classified as low risk.