

1

2 Full title:

3 Transcriptomic meta-analysis of non-Hodgkin's B-cell lymphomas
4 reveals reliance on pathways associated with the extracellular matrix

5

6 Short title:

7 Transcriptomic meta-analysis of non-Hodgkin's B-cell lymphomas

8

9 **Authors:** Naomi Rapier-Sharman¹, Jeffrey Clancy¹, and Brett E. Pickett^{1,*}

10 ¹ Department of Microbiology and Molecular Biology, Brigham Young University, Provo, UT

11 84602

12 *Corresponding author; email: brett_pickett@byu.edu, phone: 801-422-2506, address: 3141 LSB

13 Provo, UT 84602-1050

14

15

16

17

18

19

20

21

22

23

24

25

26 **Abstract**

27 Approximately 450,000 cases of Non-Hodgkin's lymphoma are diagnosed annually worldwide,
28 resulting in ~240,000 deaths. An augmented understanding of the common mechanisms of
29 pathology among relatively large numbers of B-cell Non-Hodgkin's Lymphoma (BCNHL) patients
30 is sorely needed. We consequently performed a large transcriptomic meta-analysis of available
31 BCNHL RNA-sequencing data from GEO, consisting of 322 relevant samples across ten distinct
32 public studies, to find common underlying mechanisms across BCNHL subtypes. The study was
33 limited to GEO's publicly available human B-cell RNA-sequencing datasets that met our criteria,
34 and limitations may include lack of diversity in ethnicities and age groups. We found ~10,400
35 significant differentially expressed genes (FDR-adjusted p-value < 0.05) and 33 significantly
36 modulated pathways (Bonferroni-adjusted p-value < 0.05) when comparing lymphoma samples to
37 non-diseased samples. Our findings include a significant class of proteoglycans not previously
38 associated with lymphomas as well as significant modulation of extracellular matrix-associated
39 proteins. Our drug prediction results yielded new candidates including ocriplasmin and
40 collagenase. We also used a machine learning approach to identify the BCNHL biomarkers
41 YES1, FERMT2, and FAM98B, novel biomarkers of high predictive fidelity. This meta-analysis
42 validates existing knowledge while providing novel insights into the inner workings and
43 mechanisms of B-cell lymphomas that could give rise to improved diagnostics and/or
44 therapeutics. No external funding was used for this study.

45

46 **Introduction**

47 Lymphomas are cancers of the blood. In 2016, there were 461,000 cases of Non-Hodgkin's
48 lymphoma worldwide, resulting in 240,000 deaths (1). Among non-Hodgkin's lymphomas, only
49 ~10-15% are T-cell lymphomas, while the remaining 85-90% are B-cell malignancies (2). B-cell
50 Non-Hodgkin's Lymphomas (BCNHLs) pose a significant disease burden worldwide. BCNHL
51 subtypes include Burkitt's lymphoma, nodal, extra-nodal, and splenic marginal-zone B-cell

52 lymphomas, follicular lymphoma, diffuse large B-cell lymphoma, mantle cell lymphoma, post-
53 transplantation lymphoproliferative disorders, small lymphocytic lymphoma, lymphoblastic
54 lymphoma, lymphoplasmacytic lymphoma, and lymphomatoid granulomatosis (2). B-cell
55 lymphomas are dependent on their extracellular environment for activation and transformation
56 into malignancies, including antigen activation of the B-cell receptor, canonical B-cell growth
57 signals which are also essential to the maturation of healthy B-cells, and signals delivered by
58 other immune cells in the follicular/germinal center lymphoma microenvironment (3).

59

60 The research community has dedicated extensive effort to identify the attributes that characterize
61 cancers across all subtypes. Specifically, it has been suggested previously that all cancers share
62 the following traits: selective proliferative advantage, altered stress response, vascularization,
63 invasion and metastasis, metabolic rewiring, immune modulation, and an abetting
64 microenvironment (4,5). One example of a molecular mechanism that is common in cancer is
65 malignant development through TP53 mutation, with multiple mutations in the TP53 being
66 associated with hundreds of cancer subtypes (6). Though not every gene-mechanism pairing will
67 be found across malignant cells like TP53, it is logical to assume that identifying shared genes
68 and mechanisms by meta-analyzing previous research in a focused set of related cancer
69 subtypes can be beneficial. We can therefore leverage known mechanisms from well-studied
70 subtypes to enable quicker, less expensive mechanism discovery for understudied subtypes. This
71 approach could potentially enable researchers to develop safe and effective treatments.

72

73 The widespread adoption of RNA-sequencing (RNA-seq) has opened new frontiers in disease
74 research. Rather than identifying and characterizing individual proteins, transcriptomic meta-
75 analyses can provide a mechanistic snapshot of the many upregulated or downregulated genes
76 that are affected in response to a given stimulus, such as lymphoma. Monitoring these
77 transcriptional patterns can aid in the identification of genes that could be worth further
78 experimental investigation due to their selective modulation in diseased samples. Though the

79 RNA-sequencing samples in the current study were previously published, compiling them into a
80 meta-analysis can grant us new insights into disease mechanisms by increasing the signal of
81 significant genes and reducing the statistical “noise” caused by outliers.

82

83 The aim of this study was to identify some of the shared underlying molecular mechanisms and
84 biomarkers of B-cell lymphomas by performing a meta-analysis of transcriptomic data from
85 publicly available B-cell Non-Hodgkin’s Lymphomas (BCNHLs) clinical samples. We expect our
86 analysis to validate past findings of B-cell cancer mechanisms and uncover mechanisms that
87 have not been previously associated with BCNHL.

88

89 **Results**

90 We acquired our BCNHL samples from the NCBI Gene Expression Omnibus (GEO) using the
91 search term, “b-cell lymphoma” with the goal of finding B-cell non-Hodgkin’s lymphoma samples
92 and healthy B-cell controls. We excluded non-human samples, cell lines, formalin-fixed paraffin-
93 embedded tissues, gene expression microarray experiments, single-cell (10X) RNA-sequencing
94 experiments, xenografts, samples known to be infected with EBV and KSHV, and samples which
95 contained more diverse cell types (i.e., whole blood, lymph node, PBMCs, brain, etc.). We
96 intentionally decided to not include multiple myeloma, leukemia, and Hodgkin’s lymphoma
97 samples in favor of focusing on B-cell non-Hodgkin’s lymphomas. We then located more healthy
98 B-cell control samples from BCNHL-unrelated studies to even out case and control numbers, the
99 final three studies cited in Table 1. Our final dataset included a total of 322 samples (134 BCNHL
100 samples and 188 healthy B-cell controls) from ten studies (Fig 1, Table 1, S1 File) (7–18). The
101 samples included in our meta-analysis were all clinical samples. The risk of synthesizing study
102 results and accounting for heterogeneity were reduced due to the lack of treatment metadata and
103 patient outcome data as input to our analysis. Given that the aim of this study was to compare the
104 maximum number of BCNHL samples to healthy B-cells, the only source of heterogeneity that we

105 are concerned with is the distribution of BCNHL samples across the included subtypes. Our study
 106 was limited to GEO's publicly available human B-cell RNA-sequencing datasets that met our
 107 criteria, and limitations may include lack of diversity in ethnicities and age groups.

108 **Figure 1. PRISMA flow diagram for transparent reporting of meta-analysis study selection.**
 109 Contains a study-by-study breakdown of selection criteria. All studies included were retrieved
 110 from the Gene Expression Omnibus (GEO) database hosted by NCBI.

111
 112 **Table 1. Study-based origin of samples included in the meta-analysis.**
 113

Sample Phenotype	Single End or Paired End Reads	GEO Accession #	Relevant Samples
Large B-Cell Lymphoma	Paired End	GSE153437 (7)	25
Diffuse Large B-Cell Lymphoma	Paired End	GSE130751 (8)	63
B-Cell Lymphoma + Healthy	Single End	GSE110219 (9)	2
Diffuse Large B-Cell Lymphoma	Paired End	GSE95013 (10)	28
Follicular Lymphoma + Healthy	Paired End	GSE62241 (11,12)	14
Diffuse Large B-Cell Lymphoma	Paired End	GSE50514 (13)	7
Healthy	Paired End	GSE45982 (14,15)	8
Healthy	Single End	GSE92387 (16)	12
Healthy	Paired End	GSE118254 (17)	147
Healthy	Paired End	GSE110999 (18)	16

114 The samples used in this meta-analysis all originate from publicly-available RNA-sequencing
 115 projects and can be found on NCBI's GEO.

116
 117
 118 We began by trimming, mapping, and quantifying the reads prior to calculating the significant
 119 differential gene expression when comparing the Lymphoma to the non-diseased control
 120 samples. This comparison returned ~13,800 significant differentially expressed genes (DEGs)
 121 (Figs 2 and 3, Table 2, S2 File). We then ranked this list by the FDR-corrected p-value for each
 122 gene. We observed that the top 20 DEGs include accepted biomarkers of various Lymphomas.
 123 Specifically, we confirmed several genes that have previously been explored or characterized in
 124 various subtypes of BCNHL including Apolipoprotein C1 (APOC1; logFC = 6.93, FDR = $8.55 \times$
 125 10^{-117}) and Vascular cell adhesion molecule 1 (VCAM1; logFC = 7.85, FDR = 2.29×10^{-120}) to be
 126 upregulated in BCNHLs. We also found two pathological BCNHL genes, C-C motif chemokine
 127 ligand 18 (CCL18; logFC = 10, FDR = 3.74×10^{-123}) and C-X-C motif chemokine ligand 9

128 (CXCL9; logFC = 11, FDR = 4.31×10^{-141}) to be upregulated in BCNHL as compared to healthy
 129 B-cells.

130 **Figure 2. Visualization of Differentially Expressed Genes and Gene Ontologies.**

131 Differentially expressed gene volcano plot. Green dots represent genes which were not
 132 significantly differentially expressed between healthy B-cells and BCNHL, while the salmon and
 133 blue dots represent underexpressed and overexpressed genes respectively.

134
 135 **Figure 3. Visualization of Gene Ontology Terms.**

136 Gene ontology visualization. Each rectangle represents a gene ontology term found in the KEGG
 137 Brite gene ontology hierarchy. The size of each rectangle corresponds to the number of BCNHL
 138 differentially expressed genes present in that category. The color of each rectangle corresponds
 139 to the average log₂ fold change of the genes included in that gene ontology. No KEGG Brite gene
 140 ontologies were found to be significantly differentially expressed by the bc3net hypergeometric
 141 enrichment.
 142

143 **Table 2. Top 20 significant differentially expressed genes.**
 144

	Gene Symbol	Ensembl ID	Log₂ Fold Change	FDR-corrected p-value
1	LUM	ENSG00000139329	11.1	1.11×10^{-145}
2	CXCL9	ENSG00000138755	11	4.31×10^{-141}
3	C1QC	ENSG00000159189	9.65	2.70×10^{-132}
4	C1QA	ENSG00000173372	9.54	2.03×10^{-123}
5	CCL18	ENSG00000278167	10	3.74×10^{-123}
6	VCAM1	ENSG00000162692	7.58	2.29×10^{-120}
7	C1QB	ENSG00000173369	9.4	8.19×10^{-119}
8	APOC1	ENSG00000130208	6.93	8.55×10^{-117}
9	AL512646.1	ENSG00000203396	-15.6	2.24×10^{-115}
10	CCL19	ENSG00000172724	8.48	1.27×10^{-111}
11	SLAMF8	ENSG00000158714	7.77	4.01×10^{-111}
12	COL3A1	ENSG00000168542	10.1	1.67×10^{-110}
13	TCIM	ENSG00000176907	8.07	7.86×10^{-110}
14	RARRES2	ENSG00000106538	7.25	8.21×10^{-109}
15	CXCL13	ENSG00000156234	8.8	1.72×10^{-107}
16	SPARCL1	ENSG00000152583	7.24	6.42×10^{-107}
17	PTGDS	ENSG00000107317	7.69	1.07×10^{-105}
18	COL1A2	ENSG00000164692	8.33	3.70×10^{-102}
19	CXXC5	ENSG00000171604	-2.73	3.70×10^{-102}
20	C1R	ENSG00000159403	4.7	1.41×10^{-100}

145 The top 20 differentially expressed genes between BCNHL and healthy B-cell samples. This list
 146 includes genes that have been previously researched in conjunction with BCNHL as well as novel
 147 genes.

148 We then examined the highest-ranking differentially expressed genes from our meta-analysis to
 149 identify gene-based mechanisms of disease. The first gene we observed using this approach was
 150 Lumican (LUM), which is a member of the small leucine-rich proteoglycans (SLRPs) (19), and
 151 was substantially upregulated in lymphoma (\log_2 fold change = 11.1, FDR p-value = $1.11 \times$
 152 10^{-145}). In addition, the larger family of SLRPs appears to play a role in BCNHL (Table 2).
 153 Specifically, our data show that 12/18 SLRPs are expressed in cancerous B-cells, and that 11/12
 154 B-cell-expressed SLRPs are significantly differentially expressed in BCNHL samples. We found
 155 that overall, the SLRP fold changes substantially differed (9/12 expressed SLRPs are
 156 upregulated, 2/12 are downregulated, 1/12 had no significant change), with the genes encoding
 157 SLRPs (especially Classes I and V) being well represented in the B-cell lymphoma transcriptome.

158 **Table 3. Differential expression of members of the Small Leucine-Rich Proteoglycan**
 159 **Family (SLRPs).**
 160

SLRP Class	Name	Log ₂ Fold Change	FDR-corrected p-value
Class I	DCN	2.88	1.67×10^{-41}
	BGN	7.88	3.22×10^{-95}
	ASPN	3.09	2.27×10^{-25}
	ECM2	2.1	1.44×10^{-17}
	ECMX	NP	NP
Class II	FMOD	5.71	3.86×10^{-61}
	LUM	11.1	1.11×10^{-145}
	PRELP	0.617	1.54×10^{-4}
	KERA	NP	NP
	OMD	NP	NP
Class III	EPYC	NP	NP
	OPTC	NP	NP
	OGN	NS	NS
Class IV	CHAD	-3.49	1.33×10^{-24}
	NYX	NP	NP
	TSKU	1.37	1.62×10^{-16}
Class V	PODN	1.66	5.70×10^{-10}
	PODNL1	-1.49	3.15×10^{-11}

161 Out of the 18 members of the SLRP family, 12 are expressed in B-cells and 11 are significantly
 162 differentially expressed between BCNHL and healthy B-cells. This is a novel finding.

163 *NS = not significant; NP = not present.
 164

165 Complement proteins are typically regarded as components of the innate immune system, which
 166 bind to antigen-antibody complexes to facilitate the formation of the membrane attack complex.

167 We found that genes encoding the Complement C1q A (C1QA; logFC = 9.54, FDR = 2.03 ×
 168 10⁻¹²³), Complement C1q B (C1QB; logFC = 9.4, FDR = 8.19 × 10⁻¹¹⁹) and Complement C1q C
 169 (C1QC; logFC = 9.65, FDR = 2.7 × 10⁻¹³²) chains were all dramatically and significantly
 170 upregulated in BCNHL.

171

172 We detected AL512646.1 (also known as LOC100128906 and as a WDR45-like pseudogene) as
 173 differentially expressed by B-cell non-Hodgkin's lymphoma samples, a novel observation which
 174 was somewhat unexpected. Though AL512646.1 is annotated as a pseudogene, the RNA-
 175 sequencing data shows that it is downregulated in at least a subset of BCNHLs (log₂FC = -15.1),
 176 and it has not been previously associated with cancer.

177

178 Next, we used the DRIMSeq algorithm to determine which genes had significant differences in
 179 the presence of splice variants between case and control samples. This analysis returned 320
 180 genes for which splice variants were significantly different (Table 4, S3 File). Apolipoprotein E
 181 (APOE) was the most statistically significant splice variant (Lr [likelihood ratio] = 4470, # of
 182 alternate splice variants = 4, adjusted p-value = 0). Specifically, we observed the expression of
 183 APOE transcripts ENST00000252486, ENST00000425718, ENST00000434152,
 184 ENST00000446996, and ENST00000485628 to significantly differ between non-Hodgkin's
 185 lymphoma and non-diseased B-cells.

186 **Table 4. Top 20 most significant splice variants (by gene).**

187

Gene symbol	Ensembl ID	Lr*	# of Alternate Transcripts	Adjusted P-value
APOE	ENSG00000130203	4470	4	0
COL1A1	ENSG00000108821	1520	12	5.56 × 10 ⁻³¹⁵
COL27A1	ENSG00000196739	1060	7	6.71 × 10 ⁻²²⁰
RPL5	ENSG00000122406	1040	10	3.86 × 10 ⁻²¹⁴
KLF6	ENSG00000067082	961	6	7.41 × 10 ⁻²⁰¹
SRSF6	ENSG00000124193	954	5	1.56 × 10 ⁻²⁰⁰
CYBRD1	ENSG00000071967	931	6	2.17 × 10 ⁻¹⁹⁴
PLEKHM1P1	ENSG00000214176	924	5	3.78 × 10 ⁻¹⁹⁴
VCP	ENSG00000165280	912	6	2.37 × 10 ⁻¹⁹⁰

DDX6	ENSG00000110367	872	7	8.01×10^{-181}
THRAP3	ENSG00000054118	846	3	6.63×10^{-180}
FCGR2B	ENSG00000072694	771	4	1.75×10^{-162}
CHI3L1	ENSG00000133048	715	4	2.40×10^{-150}
IFITM3	ENSG00000142089	691	3	2.53×10^{-146}
ADAM28	ENSG00000042980	719	11	4.56×10^{-144}
CIB1	ENSG00000185043	662	2	2.04×10^{-141}
ZNF318	ENSG00000171467	645	3	1.88×10^{-136}
RPS28	ENSG00000233927	621	3	3.10×10^{-131}
CCDC124	ENSG00000007080	549	3	1.14×10^{-115}
ZNF335	ENSG00000198026	545	3	8.02×10^{-115}

188 Significantly differentially expressed splice variants sorted by gene.

189 *Lr = likelihood ratio

190

191 We also observed that Collagen type I alpha 1 chain (COL1A1) had significant splice variants (Lr
192 = 1520, # of alternate splice variants = 12, adjusted p-value = $5.5599999807983 \times 10^{-315}$).

193 Interestingly, our study also found that the COL1A1 gene was significantly upregulated in BCNHL
194 (logFC = 3.73, FDR = 9.78×10^{-48}). We also observed novel significant splice variants in

195 Collagen type XXVII alpha 1 chain (COL27A1), which was found to be significant in BCNHL (Lr =
196 1060, # of alternate splice variants = 7, adjusted p-value = 6.71×10^{-220}).

197

198 We then wanted to determine which functional terms in the Gene Ontology were over-
199 represented by the list of DEGs in BCNHL. The Camera algorithm checked 14,901 terms
200 (including gene ontologies and human phenotypes) for enrichment against the significant
201 differentially expressed genes that we generated with edgeR. Although there were 482 results (p-
202 value < 0.05), none remained significant after multiple hypothesis correction (S4 File). The lack of
203 significant results is somewhat expected given the overall molecular heterogeneity of BCNHL
204 subtypes. To visualize the gene ontology changes, we used a hypergeometric enrichment
205 algorithm that applied a p-value cutoff of 0.05. We then averaged the edgeR fold-change values
206 for the genes of each gene ontology in the KEGG Brite hierarchy and plotted the enrichment
207 results using the R Treemap package to better understand the contribution of various terms to the
208 overall list of DEGs (Fig 3).

209

210 To better understand the results of our analysis at a more mechanistic level, we used the
211 signaling pathway impact analysis (SPIA) algorithm to identify intracellular signaling pathways
212 that play important roles in Lymphoma. This pathway analysis generates a null distribution
213 through bootstrapping to identify pathways that are significantly modulated when comparing sets
214 of samples. Our analysis revealed 33 significantly modulated pathways between lymphoma B-
215 cells and non-diseased B-cells (Table 5, S5 File). Specifically, we observed eight pathways that

216 were involved with the extracellular matrix and connective tissue in general, including Integrin
 217 signaling pathway, Extracellular matrix organization, ECM-receptor interaction, Focal adhesion,
 218 Integrins in angiogenesis, integrin signaling pathway, Collagen formation, and Collagen
 219 degradation. The upregulation of these pathways indicate that BCNHL likely benefits from
 220 modulations to the extracellular matrix.

221 **Table 5. Significant differentially modulated signaling pathways.**
 222

	Name	pSize	NDE	tA	pGFWER	Source Database
1	Integrin signalling pathway	99	86	114.394	2.39×10^{-5}	Panther
2	Extracellular matrix organization	204	180	80.7398395	3.14×10^{-5}	Reactome
3	ECM-receptor interaction	70	64	77.659	4.47×10^{-5}	KEGG
4	Staphylococcus aureus infection	32	29	110.568396	0.000503428	KEGG
5	Complement and coagulation cascades	36	32	37.5461944	0.000674242	KEGG
6	Urokinase-type plasminogen activator (uPA) and uPAR-mediated signaling	28	25	130.821315	0.000740169	NCI
7	Cytokine-cytokine receptor interaction	168	140	98.294	0.000780995	KEGG
8	Focal adhesion	182	150	236.615459	0.001133449	KEGG
9	PI3K-Akt signaling pathway	271	221	260.046197	0.001344307	KEGG
10	Complement cascade	29	27	102.278583	0.001440498	Reactome
11	Systemic lupus erythematosus	17	15	67.4577222	0.001616635	KEGG
12	b cell survival pathway	22	19	26.576	0.00167109	BioCarta
13	Small cell lung cancer	78	64	121.170067	0.001930414	KEGG
14	Integrins in angiogenesis	52	41	146.143424	0.00285984	NCI
15	Olfactory transduction	93	74	-148.8965	0.002966626	KEGG
16	integrin signaling pathway	37	29	77.0156667	0.003336442	BioCarta
17	erk and pi-3 kinase are necessary for collagen binding in corneal epithelia	34	26	166.268917	0.003755165	BioCarta
18	RNA Polymerase I Promoter Clearance	85	72	-40.156	0.004307328	Reactome
19	Initial triggering of complement	15	14	44.508	0.004480759	Reactome
20	RNA Polymerase I Promoter Opening	39	34	-40.907	0.004675938	Reactome
21	RHO GTPases activate PKNs	67	57	39.779	0.004802382	Reactome
22	DNA Damage/Telomere Stress Induced Senescence	61	52	32.7708077	0.004980633	Reactome
23	Creation of C4 and C2 activators	7	7	27.365	0.005633091	Reactome
24	Collagen formation	66	63	26.4272897	0.006045837	Reactome
25	Activated PKN1 stimulates transcription of AR (androgen receptor) regulated genes KLK2 and KLK3	41	35	39.094	0.006675281	Reactome
26	MET activates PTK2 signaling	18	16	63.573	0.007079287	Reactome
27	Collagen degradation	17	15	131.8905	0.008037505	Reactome
28	MET promotes cell motility	28	24	97.5445	0.00808298	Reactome
29	Regulation of IGF Activity by IGFBP	11	10	25.958725	0.008404989	Reactome
30	Classical antibody-mediated complement activation	5	5	27.354	0.008619298	Reactome
31	Serotonin Neurotransmitter Release Cycle	11	9	-13.301889	0.015608196	Reactome
32	Class A/1 (Rhodopsin-like receptors)	81	77	5.908	0.026176599	Reactome
33	Peptide ligand-binding receptors	79	75	5.84	0.040928099	Reactome

223 The significantly differentially modulated pathway results. Included are nine extracellular matrix-
 224 associated pathways.

225 *Abbreviations: psize = number of genes in pathway. NDE = number of genes from pathway
 226 which were differentially expressed. tA = measure of change between healthy and lymphoma
 227 expression; directionality indicates up- or down-regulation. pGFWER = p-value with adjustments
 228 appropriate to a multiplexed interaction network. (20)
 229
 230

231 We next used the Pathways2Targets algorithm to identify potentially novel drug targets for
 232 BCNHL from the signaling pathway results (S6 File). We sorted the results so that drug targets
 233 present in multiple signaling pathways would be ranked higher (Table 6, S7 File). We predicted
 234 the most relevant existing FDA-approved drugs for other indications that could affect the
 235 lymphoma phenotype are Doxycycline, Ocriplasmin, and Collagenase. We also identified ATN-
 236 161 as a candidate drug, but it has only been tested in phase-two trials.

237 **Table 6. Predicted BCNHL drugs based on signaling pathways.**
 238

	Drug Name	Drug ID	Significant Pathways Targeted	Is FDA Approved	Highest Clinical Trial Phase	Has Been Withdrawn
1	OCRIPLASMIN	CHEMBL209522	13	TRUE	4	FALSE
2	ATN-161	CHEMBL429745	10	FALSE	2	FALSE
3	DOXYCYCLINE	CHEMBL120069	10	TRUE	4	FALSE
4	DOXYCYCLINE	CHEMBL1433	10	TRUE	4	FALSE
5	AS-1409	CHEMBL210941	9	FALSE	1	FALSE
6	COLLAGENASE CLOSTRIDIUM HISTOLYTICUM	CHEMBL210870	9	TRUE	4	FALSE
7	FIRATEGRAST	CHEMBL210496	9	FALSE	2	FALSE
8	L19IL2	CHEMBL210960	9	FALSE	3	FALSE
9	L19SIP 131I	CHEMBL210941	9	FALSE	2	FALSE
10	L19TNFA	CHEMBL210958	9	FALSE	2	FALSE
11	VOLOCIXIMAB	CHEMBL210806	9	FALSE	3	FALSE
12	ABITUZUMAB	CHEMBL210962	8	FALSE	2	FALSE
13	AL-78898A	CHEMBL459445	8	FALSE	2	FALSE
14	CILENGITIDE	CHEMBL429876	8	FALSE	3	FALSE

4						
1 5	EPTIFIBATIDE	CHEMBL1174	8	TRUE	4	FALSE
1 6	ETARACIZUMAB	CHEMBL174301 4	8	FALSE	2	FALSE
1 7	HUMAN C1- ESTERASE INHIBITOR	CHEMBL429754 9	8	TRUE	4	FALSE
1 8	INTETUMUMAB	CHEMBL174303 2	8	FALSE	2	FALSE
1 9	PEGCETACOPLA N	CHEMBL429821 1	8	FALSE	3	FALSE

239 The top 19 drug predictions for BCNHL. Included are several drugs currently in use and others
240 that are novel candidates.
241

242 Rather than solely rely on the significant differential expression data to determine biomarkers, we
243 applied a more robust random forest machine learning method to predict biomarkers of BCNHLs.
244 Specifically, the DEG statistics focus on identifying genes that have a large difference in
245 expression between two states, while the random forest approach identifies genes that
246 consistently change across disease vs. healthy samples. Consequently, the random forest
247 approach identifies transcripts that are best capable of differentiating between disease and
248 healthy states. The top three genes identified by our random forest analysis included YES1,
249 FAM98B, and FERMT2 (Table 7, Figs 4A and 4B, S8 File). We then calculated the area under
250 the curve for the receiver-operator characteristic curve, which showed that when the expression
251 values from these three genes are combined they are 99.889% accurate at predicting whether
252 the patient samples had BCNHL (Fig 4C).

253 **Figure 4. Biomarker Prediction Yields Three-Gene Signature with 99% Predictive Ability.**

254 A) Biomarkers ranked by mean decrease in Gini impurity and permutation values using random
255 forest show YES1, FAM98B, and FERMT2 as the highest ranked predictive biomarkers (ranked
256 by mean decrease of Gini impurity score). B) Random forest biomarker prediction for the top
257 three genes in isolation. C) Receiver-operator characteristic curve using only YES1, FAM98B,
258 and FERMT2 shows these three genes are 99.889% accurate at predicting BCNHL status
259 (healthy or diseased).

260

261 **Table 7. BCNHL biomarkers predicted from gene expression using machine learning**

Gene Symbol	Mean Gini Decrease	edgeR Log ₂ Fold Change	edgeR False Discovery Rate	Disease Status	Mean (Read Counts)	Standard Deviation (Read Counts)	Median (Read Counts)
YES1	0.77	2.38	1.98x10 ⁻³⁸	Lymphoma	1151.756	1246.946	629
				Healthy	38.87234	66.01043	11
FAM98B	0.68	1.58	1.48x10 ⁻⁶⁰	Lymphoma	1797.452	1174.797	1456
				Healthy	248.4202	954.1375	32
FERMT2	0.67	2.83	1.46x10 ⁻³⁸	Lymphoma	1246.993	1200.669	841
				Healthy	32.46809	73.10299	4

262 The top three genes identified by our random forest biomarker prediction are high-fidelity
263 biomarkers of BCNHL due to their consistent and extreme upregulation across our 134 clinical
264 BCNHL samples as compared to our 188 healthy B-cell samples. Presented above are statistics
265 that capture the spread of transcriptional levels between BCNHL and healthy groups.
266
267

268 Discussion

269 The goal of this study was to collect publicly available RNA-seq data from GEO and process
270 those data to find differentially expressed genes, pathways, splice variants, and biomarkers. We
271 confirmed several biologically- and clinically relevant biomarkers and pathologic mechanisms that
272 were identified previously, as well as novel entities. We found several key genes that are
273 significantly differentially expressed in BCNHL including LUM and other SLRPs, complement
274 protein components, and the supposed pseudogene AL512646.1. We confirmed that previously
275 characterized biomarkers such as APOC1, VCAM1, CCL18, and CXCL9 are overexpressed in
276 BCNHL, and that 320 genes including APOE, COL1A1, and COL27A1 had differentially
277 expressed splice variants. We additionally found a BCNHL reliance on the upregulation of
278 pathways associated with the extracellular matrix.
279

280 To our knowledge, this is the largest meta-analysis of human samples in the BCNHL field to-date.
281 While some may be concerned that the signals from the individual subtypes of non-Hodgkin's B-
282 cell lymphomas could drown one another out, we believe that including representative samples
283 from multiple BCNHL subtypes augments the signal(s) that are shared among the represented
284 subtypes and could aid in the identification of shared mechanistic insights with reduced bias. One
285 potential limitation of our study was to only evaluate samples from the GEO database for

286 inclusion in our analysis. Our intentional focus on BCNHL excluded multiple myelomas, B-cell
287 leukemias, or Hodgkin's B-cell lymphomas. Promising future directions may include mining
288 additional public databases for similar data and potentially expanding the scope of any future
289 meta-analysis to include all B-cell malignancies. In addition, it is possible that our focus on
290 incorporating multiple publicly available datasets may have introduced potential biases in patient
291 age, gender, or ethnicity.

292

293 Though there is evidence in the literature that directly associate BCNHL to some of our results
294 (e.g. genes, splice variants, and pathways), some of our findings are novel to BCNHL. In cases
295 where no previously published research indicates the relationship between BCNHL and our
296 results, we will appeal to the Hallmarks of Cancer to investigate the relationship between our
297 differentially expressed result and a distinct cancer system (4). We therefore ranked both the
298 accuracy and confidence in our results by their relevance to BCNHL, followed by B-cell cancers in
299 general, all blood cancers, and finally all cancers. We believe that identifying a possible
300 mechanism for a gene that is associated with other cancers, and unresearched in BCNHL is still
301 relevant. We expect that a subset of these findings will justify additional wet lab experimentation.

302

303 **Differentially expressed genes suggest shared underlying** 304 **mechanisms for lymphomas**

305 LUM seems to play a role in the progression or non-progression of several different cancer types.
306 Mahadevan *et al.* previously reported upregulated LUM in both T- and B-cell lymphomas, but
307 offered no insights on potential mechanisms (21). A literature search of parallel systems revealed
308 that in breast cancer, high stromal-cell expression of LUM adjacent to the tumor stalls tumor
309 growth, and lowered stromal expression of LUM correlates with higher breast cancer mortality
310 rates and increased severity (22). In melanoma, LUM in the extracellular matrix halts metastasis
311 through direct interaction with alpha-2-beta-1 integrin (23). Both breast cancer and pancreatic

312 cancer cells have been documented to upregulate LUM, along with many other cancer types (19).
313 Overall, LUM expression by cancer cells seems to correlate with more aggressive cancers and
314 poorer patient outcomes. The massive LUM upregulation illustrated in our samples may be due to
315 the fact that the BCNHL samples available on GEO were mostly from advanced or refractory
316 cases of BCNHL. The prior finding that high LUM expression around tumors is protective against
317 metastasis in several cancer subtypes indicates the potential for LUM as a cancer-stalling
318 therapy.

319

320 Interestingly, a subset of the members in the SLRP protein family have been previously identified
321 in B-cell Non-Hodgkin's lymphomas including DCN (24), BGN (24), ASPN (25), FMOD (26), LUM
322 (21), PRELP (26), and TSKU (27). However, other members within the SLRP family have not
323 been previously considered as lymphoma biomarkers or potential pathology-inducing molecules.
324 Our novel finding is that the SLRPs ECM2, CHAD, PODN, and PODNL1 are differentially
325 expressed in BCNHL. Proteoglycans have been shown to be associated with pro-cancer
326 mechanisms in prostate, breast, colon, lung, ovary, mesothelium, pancreatic, lymphoma, and
327 esophageal cancers (19). Our results show two upregulated pathways in BCNHL that were
328 previously shown to be mechanistically intertwined with proteoglycans in cancer, which are the
329 Focal Adhesion pathway (28) and the PI3K-Akt signaling pathway (29). Taken together, this may
330 suggest a connection between previously established proteoglycan cancer mechanisms and B-
331 cell non-Hodgkin's lymphomas. Additional work is still required to elucidate the role(s) that these
332 entities play in BCNHL.

333

334 In the context of other cancers, increased expression of complement genes C1QA and C1QB at
335 week 16 of mantle cell lymphoma treatment by Venetoclax and Ibrutinib was significantly
336 associated with a worse prognosis (30), illustrating that C1QA and C1QB may be associated with
337 resistance to cancer drugs. Jiang et al. showed via immunohistochemistry that C1QB localizes to
338 the nuclei of gastric cancer cells (31). C1QB's nuclear localization suggests that C1QB may have

339 additional function(s). Upregulation of C1QA, C1QB, and C1QC in peripheral T-cell lymphoma
340 (32) and upregulation of C1QC in Epstein-Barr Virus-positive diffuse large B-cell lymphoma (33)
341 have been reported previously. In other in-vitro and in-vivo cancer models, the whole C1q protein
342 has been shown to mediate metastasis, motility, growth and proliferation, and adhesion (34). Our
343 results add to the growing body of work suggesting a potential alternate function of complement
344 proteins in cancer that warrants further investigation.

345

346 In addition to our novel findings on differentially expressed genes, we were also able to detect
347 statistically significant genes that were previously characterized in at least one subtype of
348 BCNHL. The first of these proteins is Apolipoprotein C1 (APOC1), which we observed to be
349 upregulated in BCNHL. APOC1 is one of three genes whose expression levels are predictive of
350 diffuse large B-cell lymphoma severity (35), and it is also upregulated in late stage lung cancers
351 as compared to early stage lung cancers (36). This suggests that APOC1 may be contributing to
352 cancer pathology across diverse cancers in multiple cell types.

353

354 Our observation that C-C motif chemokine ligand 18 (CCL18), which has a well-recognized role in
355 lymphoma, was upregulated in our BCNHL analysis is relevant since this gene assists large B-
356 cell lymphoma in cell proliferation, the NF-Kappa-B pathway, and the PI3K-AKT pathway (37). Its
357 upregulation in macrophages and dendritic cells from cutaneous T-cell lymphoma lesions was
358 associated with a negative prognosis (38).

359

360 Our finding C-X-C motif chemokine ligand 9 (CXCL9) to be significantly upregulated in our
361 analysis of B-cells is interesting since this gene has been shown to promote the progression of
362 diffuse large B-cell lymphoma by halting degradation of beta-catenin (CTNNB1) and upregulating
363 its initial expression (39). Our findings support this proposed mechanism with CTNNB1 being
364 upregulated in lymphoma ($\log_2FC = 1.1$, $FDR = 1.54 \times 10^{-33}$), while other elements of the
365 CTNNB1 “destruction complex” were mostly downregulated. Specifically, several of the known

366 components of the destruction complex that were detected in our analysis include APC ($\log_2FC =$
367 -0.755 , $FDR = 3.51 \times 10^{-11}$), GSK3B ($\log_2FC = -0.692$, $FDR = 2.62 \times 10^{-3}$), CSNK1A1 (not
368 significant), AXIN1 ($\log_2FC = 0.533$, $FDR = 3.96 \times 10^{-10}$), BTRC (not significant), and FBW11
369 ($\log_2FC = -0.692$, $FDR = 5.60 \times 10^{-20}$).

370

371 We identified several other genes of that may be relevant to pathogenesis. Small but significant
372 upregulation of AXIN1 is of interest for additional investigation due to its ties to CXCL9, and is not
373 known to have multiple heterogenous functions (40). AXIN1 regulates the Wnt and JNK signaling
374 pathways (41), and it regulates the Wnt pathway by degrading CTNNB1 (39). If CTNNB1 isn't
375 degraded by AXIN1, CTNNB1 translocates to the nucleus and interacts with LEF1, which we
376 found to be significantly upregulated, and TCF7 (not significant in this study), causing
377 transcription of Wnt pathway target genes to occur (42,43). Wnt helps to regulate cell cycle and
378 contributes to the increased growth rate of many cancer types (44). AXIN1 activates the JNK
379 signaling pathway by binding to MAP3K1, which we found to be significantly downregulated, or to
380 MAP3K4, which was significantly upregulated (45). Since CTNNB1 has been shown to contribute
381 to apoptosis resistance in multiple myeloma cells (46), it is possible that the inability to stop the
382 destruction of CTNNB1 in lymphoma may share a similar mechanism.

383

384 Finally, VCAM1 upregulation is associated with a poor prognosis for patients with non-Hodgkin's
385 lymphomas, and VCAM1 is under investigation as a potential serum biomarker for assessing
386 disease progression (47). Adhesion molecules such as VCAM1 promote cancer metastasis, or in
387 the case of blood cancers, extravasation, by allowing cancer cells to exit the bloodstream and
388 integrate with healthy tissues throughout the body (48).

389

390 **Splice variants suggest relevance to lymphomas**

391 To better understand the contribution of differentially expressed splice variants to disease, we
392 examined the highest-ranked DRIMseq results. This algorithm calculates statistical significance

393 based on the number of reads mapped to exons that are present in each splice variant. Our
394 observation that Apolipoprotein E (APOE) was the highest-ranking splice variant result validates
395 previous findings that associate this gene with pancreatic cancer pathology (49). In addition,
396 pediatric patients with malignant lymphoma and acute lymphoblastic leukemia who express
397 isoforms E3 and E4 of APOE are at higher risk of developing extreme hypertriglyceridemia (50).
398 Though little research has been done concerning the mechanisms of APOE in BCNHL, we
399 believe that APOE may be contributing to disease by participating in the Regulation of Insulin-like
400 Growth Factor (IGF) activity by Insulin-like Growth Factor Binding Protein (IGFBP) pathway,
401 which is we found to be a significantly modulated pathway that includes APOE. The significance
402 of APOC1 as a DEG in BCNHL, paired with the evidence of significant APOE splice variants
403 suggest that apolipoproteins may be useful targets for future BCNHL treatments.

404

405 Our observation of Collagen type I alpha 1 chain (COL1A1) as a highly ranked splice variant
406 result is novel to the best of our knowledge. However, the literature indicates that the COL1A1-
407 014 transcript regulates the CXCL12-CXCR4 axis in gastric cancer, leading to tumor progression
408 (46). In addition to displaying significant differences in splice variant expression, we also found
409 COL1A1 to be significantly upregulated in BCNHL. COL1A1 has previously been reported to be
410 upregulated in peripheral T-cell lymphoma (21). In Hodgkin's lymphoma, COL1A1 overexpression
411 is associated with epigenetic silencing of the RNA demethylase ALKBH3 and reduced survival
412 (51). COL1A1 is a member of several of our significant upregulated pathways involving the
413 extracellular matrix (ECM-receptor interaction, Focal adhesion, Extracellular matrix organization,
414 and Collagen formation). This involvement in extracellular matrix-related pathways strengthens
415 the case that the mechanism of COL1A1 may involve tumor cell interaction with its outer
416 environment.

417

418 Collagen type XXVII alpha 1 chain (COL27A1) having significant changes among its expressed
419 splice variants in BCNHL is interesting since it was recently reported as being overexpressed in

420 adenoid cystic carcinoma (52). Like COL1A1, COL27A1 is a member of the upregulated
421 Extracellular matrix organization and Collagen formation pathways, suggesting that COL27A1
422 could play a role in BCNHL extravasation.

423

424 **Extracellular matrix-related pathways may contribute to disease**

425 Our signaling pathway enrichment analysis broadened the scope of our analysis and
426 interpretation. Many of our findings supported an interesting reliance of BCNHL on pathways
427 associated with the extracellular matrix. Recent research has suggested the importance of
428 extracellular matrix components in reactivating quiescent cancer cells through the β 1-integrin
429 signaling pathway (53). It would follow that interaction with extracellular matrix components also
430 plays a role in regulating cancer cells. To our knowledge, no studies have reported the integrin
431 signaling pathway to be activated in BCNHL, though it has been reported as activated in the
432 closely-related cancer NK/T-cell lymphoma (21). The activation of these pathways suggests that
433 malignant BCNHL cells may have an advantage by interacting with the extracellular matrix. Such
434 interactions with the extracellular matrix are typically considered to be an important part of
435 metastasis (48). We found this result to be interesting since lymphomas are liquid tumors,
436 unbound by extracellular matrix. This upregulation of pathways allowing interaction with the
437 extracellular matrix may suggest that BCNHL could be invading non-lymphatic and/or non-
438 circulatory tissues.

439

440 The trend of extracellular matrix interaction is also seen in the DEG results, adding support to the
441 idea that interaction with the extracellular matrix is important for BCNHL growth and survival.

442 Additionally, COL1A1 and COL27A1, which are members of extracellular matrix-related
443 pathways, are two of the genes with the most significantly differential expression of splice
444 variants.

445

446 **Drug prediction algorithm returned both tested and novel**
447 **candidates**

448 Of our top drug results, doxycycline is currently in use for ocular B-cell lymphomas (54,55). It is
449 additionally under investigation for diffuse large B-cell lymphoma; recent work found doxycycline
450 suppresses diffuse large B-cell lymphoma growth *in vitro* and *in vivo* via CSN5 inhibition (56).
451 ATN-161 is a novel drug candidate for BCNHL. Though it is only in phase two of clinical trials, it
452 has been a successful drug against refractory solid tumors, making it a promising drug candidate
453 for other susceptible malignancies (57). ATN-161 suppresses cancer via integrin beta1 alpha5
454 antagonism, disabling invasion and metastasis (58). Ocriplasmin reverses vitreomacular
455 adhesion via interaction with fibronectin and laminin (59). Though ocriplasmin has never been
456 used in cancer before, it may be a promising drug candidate due to its ability to modulate
457 adhesion. Collagenase clostridium histolyticum is under investigation for treating collagen-rich
458 uterine fibroids and was successful at reducing the stiffness of the tumors (60).

459

460 **Machine learning predicts novel biomarkers of BCNHL**

461 YES1, FERMT2, and FAM98B are novel biomarkers not previously associated with BCNHL.
462 However, each has well-documented cancer associations. YES1 is a tyrosine kinase which
463 regulates cell cycle and apoptosis *in vitro* and cell growth *in vivo* of tumors with YES1
464 amplification (61). YES1 has been previously identified as a biomarker for non-small cell lung
465 cancer and esophageal adenocarcinoma (62,63) and may be a potential membrane biomarker.
466 YES1 can anchor to the inner membrane with help from peptide SMIM30 (64), but whether it can
467 flip to the outer leaflet has not been investigated. The role of YES1 in BCNHL pathogenesis also
468 needs additional investigation. FERMT2 has previously been pinpointed as a biomarker for other
469 cancers previously including non-small cell lung cancer and prostate cancer (65,66), but not for
470 BCNHL. FERMT2 stabilizes CTNNB1, which is a well-documented activator of oncogene
471 transcription, and is implicated in Wnt pathway regulation (67). Additionally, FERMT2 enhances

472 integrin signaling and mediates migration, invasion, and focal adhesion (68,69). Though FAM98B
473 has been shown to play an important role in the development of multiple cancers, it has not
474 previously been identified as a biomarker for any cancer. FAM98B is an arginine
475 methyltransferase, is utilized in tumorigenesis, and works in tandem with DDX1, a pan-cancer
476 marker, in RNA metabolism/processing (70,71). Like YES1 and FERMT2, FAM98B has not been
477 previously identified as a biomarker for BCNHL. These three genes have substantial diagnostic
478 potential as a liquid biopsy that could be generalizable across B-cell non-Hodgkin's lymphoma
479 subtypes. Further validation is needed to determine whether these are suitable biomarkers for
480 diagnostic or screening application.

481

482

483 In summary, our meta-analysis identified many significant DEGs and pathways that play a role in
484 B-cell non-Hodgkin's lymphomas. Our findings confirm results of previous BCNHL research,
485 indicating that the statistical analyses applied within our computational workflow pipeline are
486 effective at accurately identifying statistically significant genes, splice variants, and pathways with
487 clinical and pathological relevance. Additionally, several of our results are novel, which need
488 additional validation in future experiments. It is likely that at least some of these novel findings
489 were detected due to the ability of our meta-analysis to reduce the statistical "noise" produced by
490 outliers from individual studies and increase the biologically-relevant signal. Specifically, our
491 findings suggest that LUM and 10 other small leucine-rich proteoglycans are significantly
492 differentially expressed in BCNHL, that AL512646.1 is not a pseudo-gene, that APOE, COL1A1,
493 and COL27a1 have significant differentially expressed splice variants in BCNHL, and that BCNHL
494 is strongly reliant on the overexpression of extracellular matrix-associated pathways. The
495 predominant drug prediction results nearly universally targeted extracellular matrix-associated
496 mechanisms, and has yielded several promising new potential drug candidates including
497 ocriplasmin and ATN-161. Our random forest biomarker discovery pinpointed three novel
498 biomarker genes not previously associated with BCNHL, YES1, FERMT2, and FAM98B, which

499 show high fidelity in predicting lymphoma presence based on transcriptional levels. These
500 findings shed additional light on the underlying intracellular mechanisms of BCNHL and could be
501 used in the development of improved diagnostics and therapeutics to further improve human
502 health.

503

504

505 **Methods**

506 **Collecting samples**

507 RNA-sequencing samples were acquired from the National Center for Biotechnology Information
508 (NCBI) Gene Expression Omnibus (GEO) database using the search term, “b-cell lymphoma” to
509 find B-cell non-Hodgkin’s lymphoma samples and healthy B-cell controls. The automatic GEO
510 filters “Homo sapiens” and “high-throughput RNA-sequencing” were applied. Cell lines, formalin-
511 fixed paraffin-embedded tissues, gene expression microarray experiments, single-cell (10X)
512 RNA-sequencing experiments, xenografts, samples known to be infected with EBV and KSHV,
513 and samples which contained more diverse cell types (i.e., whole blood, lymph node, PBMCs,
514 brain, etc.) were manually excluded. All samples that had one or more of these disqualifying
515 attributes were excluded from the dataset prior to analysis, which may have included only a
516 subset of samples from any individual study in the meta-analysis. Multiple myeloma, leukemia,
517 and Hodgkin’s lymphoma samples were intentionally excluded in favor of focusing on B-cell non-
518 Hodgkin’s lymphomas. Records were accepted or rejected based on the standardized exclusion
519 criteria detailed above by our team. To avoid inclusion bias, any sample that could not be
520 excluded by the standardized exclusion criteria was included in the study. While a subset of the
521 healthy control samples was obtained from the same RNA sequencing projects as the BCNHL
522 samples, others were obtained from three unrelated B-cell datasets with healthy controls to
523 create roughly equivalent-sized BCNHL and healthy groups. Final dataset assembly from GEO
524 concluded on October 22, 2020, resulting in a dataset of 322 samples (134 BCNHL samples and

525 188 healthy B-cell controls) from ten studies (7–18). The raw data for these experiments were
526 previously collected by the data providers and conform to the appropriate ethical oversight to
527 protect patient autonomy and patient identity. All 10 primary RNA-sequencing datasets from
528 which samples were gathered for our lymphoma meta-analysis have been published in the peer-
529 reviewed literature, increasing overall confidence that each dataset has acceptable quality (Table
530 1, Fig 1).

531

532 **Preprocessing of RNA-sequencing data**

533 Following the manual curation of the RNA-seq samples, the fastq files were pre-processed as
534 previously described (72). In brief, fastq files containing RNA-sequencing data were downloaded
535 from the Sequence Read Archive (SRA) using the sratools software package. The fastq files, the
536 associated metadata file, and a configuration file for each dataset were then generated and used
537 as input to the Automated Reproducible MOdular Workflow for Preprocessing and Differential
538 Analysis of RNA-seq Data (ARMOR) workflow (73). A configuration file was used by ARMOR to
539 appropriately set up a python-based Snakemake workflow (74). In the ARMOR workflow,
540 adapters and poor-quality regions of reads were trimmed with TrimGalore! (75), quality control
541 metrics were calculated with FastQC (76), reads were mapped to the human GRCh38
542 transcriptome and total gene transcripts quantified with Salmon (77), significant differential gene
543 expression was calculated using a negative binomial distribution implemented in edgeR (78),
544 Gene Ontology enrichment was performed against terms from the MSigDB (79) while adjusting
545 for inter-gene correlation using the Camera algorithm (80), and significant splice variants were
546 predicted with DRIMseq (81). The significant differentially expressed genes from the ARMOR
547 workflow were then used as input to the signaling pathway impact analysis (SPIA) algorithm to
548 enrich differentially expressed genes against intracellular signaling pathways from five databases
549 including KEGG, Panther, BioCarta, Reactome, and NCI (20,82–85). Differentially expressed
550 genes and splice variants calculated by ARMOR and DRIMSeq were evaluated by the effect

551 measures log₂ fold change and likelihood ratio respectively. Confidence in results was
552 accomplished using false discovery-rate adjusted p-values.

553

554 **Additional analysis and visualization of differentially expressed** 555 **genes and gene ontologies**

556 The PRISMA flowchart template was used to generate figure one consistent with transparent
557 reporting of meta-analysis generation and results (86).

558

559 The R package ggplot was used to construct the Fig 2 volcano plot from using FDRs and log₂
560 fold change values for each gene from the edgeR output (87).

561

562 The KEGG ontology was extracted from the Brite Hierarchy using existing code (82). Genes
563 included in the Brite Hierarchy were then computationally matched to their corresponding edgeR
564 log₂ fold change values. A statistical enrichment of the KEGG gene ontologies was performed
565 using the R package bc3net (88) prior to visualizing the bc3net enrichment results with the R
566 package Treemap in Fig 3 (89).

567

568 **Biomarker prediction using differentially expressed gene data**

569 Salmon-derived transcript counts were organized into a tabular format and samples were
570 randomly assigned to either the testing set (30%) or the training set (70%). The R package
571 randomForest was used to run a supervised classification analysis to determine biomarkers (90).

572 The initial results from the whole transcriptome were then reduced to the 3, 5, and 10 best-
573 scoring transcriptional biomarkers, based on the mean Gini impurity decrease values for each of
574 the features. These values were then sorted by size to determine the transcribed genes from the
575 original dataset with the largest association. The area under the curve (AUC) was calculated from

576 the receiver operator characteristic curves that were generated for each set of random forest
577 results to determine the efficacy of the selected biomarkers for disease prediction.

578

579 **Drug prediction using differentially modulated pathways**

580 Drug prediction was conducted by running the significantly modulated pathways file that was
581 generated by SPIA file through Pathways2Targets2.R algorithm (91). To summarize drug
582 findings, the Pathways2Targets output was processed using a custom R script
583 `most_common_treatments_2021_09_19.R` (92).

584

585 **Other information**

586 This meta-analysis was not registered. The review protocol was not prepared separately but is
587 described in detail in the methods section. This meta-analysis received no specific financial
588 support and was supported by general funding from Brigham Young University. Brigham Young
589 University played no role in the ideation, synthesis, or analysis of this transcriptomic meta-
590 analysis. The authors each declare that they have no competing interests. All raw data analyzed
591 in this study can be found in NCBI's GEO. Analytical codes can be found as cited in materials and
592 methods. All processed data are available in the supplementary materials or at DOI:
593 [10.5281/zenodo.4757764](https://doi.org/10.5281/zenodo.4757764).

594

595 **Acknowledgments**

596 We thank the high-performance computing resources hosted by the BYU Research Computing
597 Center. We also gratefully acknowledge those who generated, provided, and submitted the
598 original data.

599

600

601 **References**

602

- 603 1. Global Burden of Disease Cancer Collaboration. Global, Regional, and National Cancer
604 Incidence, Mortality, Years of Life Lost, Years Lived With Disability, and Disability-Adjusted
605 Life-Years for 29 Cancer Groups, 1990 to 2016: A Systematic Analysis for the Global Burden
606 of Disease Study. *JAMA Oncology*. 2018 Nov 1;4(11):1553–68.
- 607 2. Shankland KR, Armitage JO, Hancock BW. Non-Hodgkin lymphoma. *The Lancet*. 2012 Sep
608 1;380(9844):848–57.
- 609 3. Küppers R. Mechanisms of B-cell lymphoma pathogenesis. *Nat Rev Cancer*. 2005
610 Apr;5(4):251–62.
- 611 4. Hanahan D, Weinberg RA. The hallmarks of cancer. *Cell*. 2000 Jan 7;100(1):57–70.
- 612 5. Fouad YA, Aanei C. Revisiting the hallmarks of cancer. *Am J Cancer Res*. 2017;7(5):1016–
613 36.
- 614 6. Ochoa D, Hercules A, Carmona M, Suveges D, Gonzalez-Uriarte A, Malangone C, et al.
615 Open Targets Platform: supporting systematic drug-target identification and prioritisation.
616 *Nucleic Acids Res*. 2021 Jan 8;49(D1):D1302–10.
- 617 7. Faramand R, Jain M, Staedtke V, Kotani H, Bai R, Reid K, et al. Tumor Microenvironment
618 Composition and Severe Cytokine Release Syndrome (CRS) Influence Toxicity in Patients
619 with Large B-Cell Lymphoma Treated with Axicabtagene Ciloleucel. *Clin Cancer Res*. 2020
620 Sep 15;26(18):4823–31.
- 621 8. Li M, Chiang YL, Lyssiotis CA, Teater MR, Hong JY, Shen H, et al. Non-oncogene Addiction to
622 SIRT3 Plays a Critical Role in Lymphomagenesis. *Cancer Cell*. 2019 Jun 10;35(6):916-931.e9.
- 623 9. Porpaczy E, Tripolt S, Hoelbl-Kovacic A, Gisslinger B, Bago-Horvath Z, Casanova-Hevia E, et
624 al. Aggressive B-cell lymphomas in patients with myelofibrosis receiving JAK1/2 inhibitor
625 therapy. *Blood*. 2018 Aug 16;132(7):694–706.
- 626 10. Teater M, Dominguez PM, Redmond D, Chen Z, Ennishi D, Scott DW, et al. AICDA drives
627 epigenetic heterogeneity and accelerates germinal center-derived lymphomagenesis. *Nat*
628 *Commun*. 2018 Jan 15;9(1):222.
- 629 11. Raju S, Kretzmer LZ, Koues OI, Payton JE, Oltz EM, Cashen A, et al. NKG2D-NKG2D Ligand
630 Interaction Inhibits the Outgrowth of Naturally Arising Low-Grade B Cell Lymphoma In
631 Vivo. *J Immunol*. 2016 Jun 1;196(11):4805–13.
- 632 12. Koues OI, Kowalewski RA, Chang LW, Pyfrom SC, Schmidt JA, Luo H, et al. Enhancer
633 sequence variants and transcription-factor deregulation synergize to construct pathogenic
634 regulatory circuits in B-cell lymphoma. *Immunity*. 2015 Jan 20;42(1):186–98.

- 635 13. Rouhigharabaei L, Finalet Ferreiro J, Tousseyn T, van der Krogt JA, Put N, Haralambieva E,
636 et al. Non-IG aberrations of FOXP1 in B-cell malignancies lead to an aberrant expression of
637 N-truncated isoforms of FOXP1. *PLoS One*. 2014;9(1):e85851.
- 638 14. Béguelin W, Popovic R, Teater M, Jiang Y, Bunting KL, Rosen M, et al. EZH2 is required for
639 germinal center formation and somatic EZH2 mutations promote lymphoid
640 transformation. *Cancer Cell*. 2013 May 13;23(5):677–92.
- 641 15. Verma A, Jiang Y, Du W, Fairchild L, Melnick A, Elemento O. Transcriptome sequencing
642 reveals thousands of novel long non-coding RNAs in B cell lymphoma. *Genome Med*. 2015
643 Nov 1;7:110.
- 644 16. Jenks SA, Cashman KS, Zumaquero E, Marigorta UM, Patel AV, Wang X, et al. Distinct
645 Effector B Cells Induced by Unregulated Toll-like Receptor 7 Contribute to Pathogenic
646 Responses in Systemic Lupus Erythematosus. *Immunity*. 2018 Oct 16;49(4):725-739.e6.
- 647 17. Scharer CD, Blalock EL, Mi T, Barwick BG, Jenks SA, Deguchi T, et al. Epigenetic
648 programming underpins B cell dysfunction in human SLE. *Nat Immunol*. 2019
649 Aug;20(8):1071–82.
- 650 18. Wang S, Wang J, Kumar V, Karnell JL, Naiman B, Gross PS, et al. IL-21 drives expansion and
651 plasma cell differentiation of autoreactive CD11c^{hi}T-bet⁺ B cells in SLE. *Nat Commun*. 2018
652 May 1;9(1):1758.
- 653 19. Appunni S, Anand V, Khandelwal M, Gupta N, Rubens M, Sharma A. Small Leucine Rich
654 Proteoglycans (decorin, biglycan and lumican) in cancer. *Clin Chim Acta*. 2019 Apr;491:1–7.
- 655 20. Tarca AL, Draghici S, Khatri P, Hassan SS, Mittal P, Kim JS, et al. A novel signaling pathway
656 impact analysis. *Bioinformatics*. 2009 Jan 1;25(1):75–82.
- 657 21. Mahadevan D, Spier C, Croce KD, Miller S, George B, Riley CJ, et al. Gene Expression
658 Profiling of Peripheral T-Cell Lymphoma (PTCL, NOS) and Comparison to Diffuse Large B-
659 Cell Lymphoma (DLBCL). *Blood*. 2004 Nov 16;104(11):2277–2277.
- 660 22. Troup S, Njue C, Kliewer EV, Parisien M, Roskelley C, Chakravarti S, et al. Reduced
661 Expression of the Small Leucine-rich Proteoglycans, Lumican, and Decorin Is Associated
662 with Poor Outcome in Node-negative Invasive Breast Cancer. *Clin Cancer Res*. 2003 Jan
663 1;9(1):207–14.
- 664 23. Brézillon S, Pietraszek K, Maquart FX, Wegrowski Y. Lumican effects in the control of
665 tumour progression and their links with metalloproteinases and integrins. *FEBS J*. 2013
666 May;280(10):2369–81.
- 667 24. Kotlov N, Bagaev A, Revuelta MV, Phillip JM, Cacciapuoti MT, Antysheva Z, et al. Clinical
668 and Biological Subtypes of B-cell Lymphoma Revealed by Microenvironmental Signatures.
669 *Cancer Discov*. 2021 Jun 1;11(6):1468–89.

- 670 25. Sarkozy C, Chong L, Takata K, Chavez EA, Miyata-Takata T, Duns G, et al. Gene expression
671 profiling of gray zone lymphoma. *Blood Advances*. 2020 Jun 9;4(11):2523–35.
- 672 26. Mikaelsson E, Jeddi-Tehrani M, Å–Sterborg A, Shokri F, Rabbani H, Mellstedt H. Small
673 Leucine Rich Proteoglycans as Novel Tumor Markers In Chronic Lymphocytic Leukemia.
674 *Blood*. 2010 Nov 19;116(21):694.
- 675 27. Huang H, Zhang D, Fu J, Zhao L, Li D, Sun H, et al. Tsukushi is a novel prognostic biomarker
676 and correlates with tumor-infiltrating B cells in non-small cell lung cancer. *Aging (Albany*
677 *NY)*. 2021 Jan 10;13(3):4428–51.
- 678 28. Munesue S, Kusano Y, Oguri K, Itano N, Yoshitomi Y, Nakanishi H, et al. The role of
679 syndecan-2 in regulation of actin-cytoskeletal organization of Lewis lung carcinoma-
680 derived metastatic clones. *Biochem J*. 2002 Apr 15;363(Pt 2):201–9.
- 681 29. Mahtouk K, Tjin EPM, Spaargaren M, Pals ST. The HGF/MET pathway as target for the
682 treatment of multiple myeloma and B-cell lymphomas. *Biochim Biophys Acta*. 2010
683 Dec;1806(2):208–19.
- 684 30. Davis J, Handunnetti SM, Sharpe C, Turner G, Anderson MA, Roberts AW, et al. Long Term
685 Responses to Venetoclax and Ibrutinib in Mantle Cell Lymphoma Are Associated with
686 Immunological Recovery and Prognostic Changes in Inflammatory Biomarkers. *Blood*. 2019
687 Nov 13;134(Supplement_1):2791–2791.
- 688 31. Jiang J, Ding Y, Wu M, Lyu X, Wang H, Chen Y, et al. Identification of TYROBP and C1QB as
689 Two Novel Key Genes With Prognostic Value in Gastric Cancer by Network Analysis. *Front*
690 *Oncol*. 2020;10:1765.
- 691 32. Gao HX, Wang MB, Li SJ, Niu J, Xue J, Li J, et al. Identification of Hub Genes and Key
692 Pathways Associated with Peripheral T-cell Lymphoma. *Curr Med Sci*. 2020 Oct;40(5):885–
693 99.
- 694 33. Yoon H, Park S, Ju H, Ha SY, Sohn I, Jo J, et al. Integrated copy number and gene expression
695 profiling analysis of Epstein-Barr virus-positive diffuse large B-cell lymphoma. *Genes*
696 *Chromosomes Cancer*. 2015 Jun;54(6):383–96.
- 697 34. Bulla R, Tripodo C, Rami D, Ling GS, Agostinis C, Guarnotta C, et al. C1q acts in the tumour
698 microenvironment as a cancer-promoting factor independently of complement activation.
699 *Nat Commun*. 2016 Feb 1;7(1):10346.
- 700 35. Zamani-Ahmadm Mahmudi M, Nassiri SM. Development of a Reproducible Prognostic Gene
701 Signature to Predict the Clinical Outcome in Patients with Diffuse Large B-Cell Lymphoma.
702 *Sci Rep*. 2019 Aug 21;9(1):12198.
- 703 36. Ko HL, Wang YS, Fong WL, Chi MS, Chi KH, Kao SJ. Apolipoprotein C1 (APOC1) as a novel
704 diagnostic and prognostic biomarker for lung cancer: A marker phase I trial. *Thorac Cancer*.
705 2014 Nov;5(6):500–8.

- 706 37. Zhou Q, Huang L, Gu Y, Lu H, Feng Z. The expression of CCL18 in diffuse large B cell
707 lymphoma and its mechanism research. *Cancer Biomark*. 2018;21(4):925–34.
- 708 38. Miyagaki T, Sugaya M, Suga H, Ohmatsu H, Fujita H, Asano Y, et al. Increased CCL18
709 expression in patients with cutaneous T-cell lymphoma: association with disease severity
710 and prognosis. *J Eur Acad Dermatol Venereol*. 2013 Jan;27(1):e60-67.
- 711 39. Ruiduo C, Ying D, Qiwei W. CXCL9 promotes the progression of diffuse large B-cell
712 lymphoma through up-regulating β -catenin. *Biomedicine & Pharmacotherapy*. 2018 Nov
713 1;107:689–95.
- 714 40. Mani M, Chen C, Ambler V, Liu H, Mathur T, Zwicke G, et al. MoonProt: a database for
715 proteins that are known to moonlight. *Nucleic Acids Research*. 2015 Jan 28;43(D1):D277–
716 82.
- 717 41. Kusano S, Raab-Traub N. I-mfa domain proteins interact with Axin and affect its regulation
718 of the Wnt and c-Jun N-terminal kinase signaling pathways. *Mol Cell Biol*. 2002
719 Sep;22(18):6393–405.
- 720 42. Bienz M. TCF: transcriptional activator or repressor? *Current Opinion in Cell Biology*. 1998
721 Jun 1;10(3):366–72.
- 722 43. Escobar G, Mangani D, Anderson AC. T cell factor 1: A master regulator of the T cell
723 response in disease. *Sci Immunol*. 2020 Nov 6;5(53):eabb9726.
- 724 44. Bugter JM, Fenderico N, Maurice MM. Mutations and mechanisms of WNT pathway
725 tumour suppressors in cancer. *Nat Rev Cancer*. 2021 Jan;21(1):5–21.
- 726 45. Luo W, Ng WW, Jin LH, Ye Z, Han J, Lin SC. Axin Utilizes Distinct Regions for Competitive
727 MEK1 and MEK4 Binding and JNK Activation *. *Journal of Biological Chemistry*. 2003 Sep
728 26;278(39):37451–8.
- 729 46. Su N, Wang P, Li Y. Role of Wnt/ β -catenin pathway in inducing autophagy and apoptosis in
730 multiple myeloma cells. *Oncol Lett*. 2016 Dec;12(6):4623–9.
- 731 47. Shah N, Cabanillas F, McIntyre B, Feng L, McLaughlin P, Rodriguez MA, et al. Prognostic
732 value of serum CD44, intercellular adhesion molecule-1 and vascular cell adhesion
733 molecule-1 levels in patients with indolent non-Hodgkin lymphomas. *Leuk Lymphoma*.
734 2012 Jan;53(1):50–6.
- 735 48. Zetter BR. Adhesion molecules in tumor metastasis. *Semin Cancer Biol*. 1993 Aug;4(4):219–
736 29.
- 737 49. Omenn GS, Yocum AK, Menon R. Alternative splice variants, a new class of protein cancer
738 biomarker candidates: findings in pancreatic cancer and breast cancer with systems
739 biology implications. *Dis Markers*. 2010;28(4):241–51.

- 740 50. Tozuka M, Yamauchi K, Hidaka H, Nakabayashi T, Okumura N, Katsuyama T.
741 Characterization of hypertriglyceridemia induced by L-asparaginase therapy for acute
742 lymphoblastic leukemia and malignant lymphoma. *Ann Clin Lab Sci*. 1997 Sep 1;27(5):351–
743 7.
- 744 51. Esteve-Puig R, Climent F, Piñeyro D, Domingo-Domènech E, Davalos V, Encuentra M, et al.
745 Epigenetic loss of m1A RNA demethylase ALKBH3 in Hodgkin lymphoma targets collagen,
746 conferring poor clinical outcome. *Blood*. 2021 Feb 18;137(7):994–9.
- 747 52. Arolt C, Meyer M, Hoffmann F, Wagener-Rydzek S, Schwarz D, Nachtsheim L, et al.
748 Expression Profiling of Extracellular Matrix Genes Reveals Global and Entity-Specific
749 Characteristics in Adenoid Cystic, Mucoepidermoid and Salivary Duct Carcinomas. *Cancers*.
750 2020 Sep;12(9):2466.
- 751 53. Barkan D, El Touny LH, Michalowski AM, Smith JA, Chu I, Davis AS, et al. Metastatic growth
752 from dormant cells induced by a col-I-enriched fibrotic environment. *Cancer Res*. 2010 Jul
753 15;70(14):5706–16.
- 754 54. Kim TM, Kim KH, Lee MJ, Jeon YK, Lee SH, Kim DW, et al. First-line therapy with doxycycline
755 in ocular adnexal mucosa-associated lymphoid tissue lymphoma: A retrospective analysis
756 of clinical predictors. *Cancer Science*. 2010;101(5):1199–203.
- 757 55. Han JJ, Kim TM, Jeon YK, Kim MK, Khwarg SI, Kim CW, et al. Long-term outcomes of first-
758 line treatment with doxycycline in patients with previously untreated ocular adnexal
759 marginal zone B cell lymphoma. *Ann Hematol*. 2015 Apr 1;94(4):575–81.
- 760 56. Pulvino M, Chen L, Oleksyn D, Li J, Compitello G, Rossi R, et al. Inhibition of COP9-
761 signalosome (CSN) deneddylating activity and tumor growth of diffuse large B-cell
762 lymphomas by doxycycline. *Oncotarget*. 2015 Jun 4;6(17):14796–813.
- 763 57. Cianfrocca ME, Kimmel KA, Gallo J, Cardoso T, Brown MM, Hudes G, et al. Phase 1 trial of
764 the antiangiogenic peptide ATN-161 (Ac-PHSCN-NH₂), a beta integrin antagonist, in
765 patients with solid tumours. *Br J Cancer*. 2006 Jun;94(11):1621–6.
- 766 58. Barkan D, Chambers AF. β 1-Integrin: A Potential Therapeutic Target in the Battle against
767 Cancer Recurrence. *Clinical Cancer Research*. 2011 Nov 30;17(23):7219–23.
- 768 59. Baldo BA. Enzymes Approved for Human Therapy: Indications, Mechanisms and Adverse
769 Effects. *BioDrugs*. 2015 Feb 1;29(1):31–55.
- 770 60. Jayes FL, Liu B, Moutos FT, Kuchibhatla M, Guilak F, Leppert PC. Loss of stiffness in
771 collagen-rich uterine fibroids after digestion with purified collagenase *Clostridium*
772 *histolyticum*. *American Journal of Obstetrics and Gynecology*. 2016 Nov 1;215(5):596.e1-
773 596.e8.

- 774 61. Hamanaka N, Nakanishi Y, Mizuno T, Horiguchi-Takei K, Akiyama N, Tanimura H, et al. YES1
775 Is a Targetable Oncogene in Cancers Harboring YES1 Gene Amplification. *Cancer Research*.
776 2019 Nov 15;79(22):5734–45.
- 777 62. Garmendia I, Pajares MJ, Hermida-Prado F, Ajona D, Bértolo C, Sainz C, et al. YES1 Drives
778 Lung Cancer Growth and Progression and Predicts Sensitivity to Dasatinib. *Am J Respir Crit*
779 *Care Med*. 2019 Oct;200(7):888–99.
- 780 63. Liu Z, Jiang F, Tian G, Wang S, Sato F, Meltzer SJ, et al. Sparse Logistic Regression with Lp
781 Penalty for Biomarker Identification. *Statistical Applications in Genetics and Molecular*
782 *Biology* [Internet]. 2007 Feb 10 [cited 2022 Apr 16];6(1). Available from:
783 <http://www.degruyter.com/document/doi/10.2202/1544-6115.1248/html>
- 784 64. Pang Y, Liu Z, Han H, Wang B, Li W, Mao C, et al. Peptide SMIM30 promotes HCC
785 development by inducing SRC/YES1 membrane anchoring and MAPK pathway activation.
786 *Journal of Hepatology*. 2020 Nov 1;73(5):1155–69.
- 787 65. Sun W, Guo J, Cheng Z, Zhang Y, Gao Y. Identification of the Dysregulated Pathways and
788 Key Gene in Prostate Cancer by Transcriptome Analysis and Cell Biology Experiments
789 [Internet]. In Review; 2022 Jan [cited 2022 Apr 16]. Available from:
790 <https://www.researchsquare.com/article/rs-1122595/v1>
- 791 66. Su X, Liu N, Wu W, Zhu Z, Xu Y, He F, et al. Comprehensive analysis of prognostic value and
792 immune infiltration of kindlin family members in non-small cell lung cancer. *BMC Med*
793 *Genomics*. 2021 May 2;14(1):119.
- 794 67. Yu Y, Wu J, Wang Y, Zhao T, Ma B, Liu Y, et al. Kindlin 2 forms a transcriptional complex
795 with β -catenin and TCF4 to enhance Wnt signalling. *EMBO Rep*. 2012 Aug;13(8):750–8.
- 796 68. Kawamura E, Hamilton GB, Miskiewicz EI, MacPhee DJ. Fermitin family homolog-2
797 (FERMT2) is highly expressed in human placental villi and modulates trophoblast invasion.
798 *BMC Developmental Biology*. 2018 Nov 1;18(1):19.
- 799 69. Yasuda-Yamahara M, Rogg M, Frimmel J, Trachte P, Helmstaedter M, Schroder P, et al.
800 FERMT2 links cortical actin structures, plasma membrane tension and focal adhesion
801 function to stabilize podocyte morphology. *Matrix Biology*. 2018 Aug 1;68–69:263–79.
- 802 70. Akter KA, Mansour MA, Hyodo T, Senga T. FAM98A associates with DDX1-C14orf166-
803 FAM98B in a novel complex involved in colorectal cancer progression. *Int J Biochem Cell*
804 *Biol*. 2017 Mar;84:1–13.
- 805 71. Gao B, Li X, Li S, Wang S, Wu J, Li J. Pan-cancer analysis identifies RNA helicase DDX1 as a
806 prognostic marker. *Phenomics*. 2022 Feb 1;2(1):33–49.
- 807 72. Rapier-Sharman N, Krapohl J, Beausoleil EJ, Gifford KTL, Hinatsu BR, Hoffmann CS, et al.
808 Preprocessing of Public RNA-Sequencing Datasets to Facilitate Downstream Analyses of
809 Human Diseases. *Data*. 2021 Jul;6(7):75.

- 810 73. Orjuela S, Huang R, Hembach KM, Robinson MD, Soneson C. ARMOR: An Automated
811 Reproducible MODular Workflow for Preprocessing and Differential Analysis of RNA-seq
812 Data. *G3 (Bethesda)*. 2019 Jul 9;9(7):2089–96.
- 813 74. Köster J, Rahmann S. Snakemake—a scalable bioinformatics workflow engine.
814 *Bioinformatics*. 2012 Oct 1;28(19):2520–2.
- 815 75. Babraham Bioinformatics - Trim Galore! [Internet]. [cited 2021 Jun 7]. Available from:
816 https://www.bioinformatics.babraham.ac.uk/projects/trim_galore/
- 817 76. Babraham Bioinformatics - FastQC A Quality Control tool for High Throughput Sequence
818 Data [Internet]. [cited 2021 Jun 7]. Available from:
819 <https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>
- 820 77. Patro R, Duggal G, Love MI, Irizarry RA, Kingsford C. Salmon provides fast and bias-aware
821 quantification of transcript expression. *Nat Methods*. 2017 Apr;14(4):417–9.
- 822 78. Robinson MD, McCarthy DJ, Smyth GK. edgeR: a Bioconductor package for differential
823 expression analysis of digital gene expression data. *Bioinformatics*. 2010 Jan 1;26(1):139–
824 40.
- 825 79. Liberzon A, Subramanian A, Pinchback R, Thorvaldsdóttir H, Tamayo P, Mesirov JP.
826 Molecular signatures database (MSigDB) 3.0. *Bioinformatics*. 2011 Jun 15;27(12):1739–40.
- 827 80. Wu D, Smyth GK. Camera: a competitive gene set test accounting for inter-gene
828 correlation. *Nucleic Acids Res*. 2012 Sep 1;40(17):e133.
- 829 81. Nowicka M, Robinson MD. DRIMSeq: a Dirichlet-multinomial framework for multivariate
830 count outcomes in genomics. *F1000Res*. 2016;5:1356.
- 831 82. Ogata H, Goto S, Sato K, Fujibuchi W, Bono H, Kanehisa M. KEGG: Kyoto Encyclopedia of
832 Genes and Genomes. *Nucleic Acids Res*. 1999 Jan 1;27(1):29–34.
- 833 83. Mi H, Muruganujan A, Casagrande JT, Thomas PD. Large-scale gene function analysis with
834 the PANTHER classification system. *Nat Protoc*. 2013 Aug;8(8):1551–66.
- 835 84. Nishimura D. BioCarta. *Biotech Software & Internet Report*. 2001 Jun 1;2(3):117–20.
- 836 85. Fabregat A, Jupe S, Matthews L, Sidiropoulos K, Gillespie M, Garapati P, et al. The
837 Reactome Pathway Knowledgebase. *Nucleic Acids Research*. 2018 Jan 4;46(D1):D649–55.
- 838 86. Page MJ, McKenzie JE, Bossuyt PM, Boutron I, Hoffmann TC, Mulrow CD, et al. The PRISMA
839 2020 statement: an updated guideline for reporting systematic reviews. *BMJ*. 2021 Mar
840 29;372:n71.
- 841 87. Wickham H. ggplot2. *WIREs Computational Statistics*. 2011;3(2):180–5.

- 842 88. Simoes R de M, Emmert-Streib F. Bagging Statistical Network Inference from Large-Scale
843 Gene Expression Data. PLOS ONE. 2012 Mar 30;7(3):e33624.
- 844 89. Tennekes M, de Jonge E. TOP-DOWN DATA ANALYSIS WITH TREEMAPS. 2011;6.
- 845 90. Liaw A, Wiener M. Classification and Regression by randomForest. R News. 2002;2(3):18–
846 22.
- 847 91. Scott TM, Jensen S, Pickett BE. A signaling pathway-driven bioinformatics pipeline for
848 predicting therapeutics against emerging infectious diseases. F1000Res. 2021;10:330.
- 849 92. naomi-rapier-sharman. summarize_Pathways2Targets2_output [Internet]. 2022 [cited
850 2022 Apr 15]. Available from: [https://github.com/naomi-rapier-](https://github.com/naomi-rapier-sharman/summarize_Pathways2Targets2_output)
851 [sharman/summarize_Pathways2Targets2_output](https://github.com/naomi-rapier-sharman/summarize_Pathways2Targets2_output)

852
853

854 Supporting information captions

855

856 **S0 File. B-Cell Non-Hodgkin’s Lymphomas Transcriptomic Meta-Analysis: Supplementary**

857 **Materials Description.** This document contains a guide to allow readers to easily navigate the
858 supplementary files.

859 **S1 File. PRISMA 2020 checklist for transparent meta-analysis reporting.** This document
860 contains the PRISMA guidelines for reporting meta-analyses/systematic reviews and the
861 manuscript location of required information.

862 **S2 File. Differentially expressed gene results (edgeR output).** The entire BCNHL differentially
863 expressed gene list produced by edgeR. The file is written in TSV format and can be successfully
864 opened in Excel. Some genes in the original TSV file have more associated data than can fit on
865 one line. The genes in the file are ranked according to FDR, with the smallest FDRs at the top.
866 Please note, the file also contains differential expression results for genes which were not
867 significantly differentially expressed at the bottom. For description of column contents, please
868 see S0 File.

869 **S3 File. Differentially expressed splice variants (by gene; DRIMSeq output).** The BCNHL
870 differentially expressed splice variant results (by gene) produced by DRIMSeq. Genes are ranked
871 according to adjusted p-value, with the lowest adjusted p-value at the top. NOTE: Should you
872 desire to discover which transcripts of a certain gene are present in the BCNHL dataset, see the
873 "tx_ids" column in Supplementary File S1. For description of column contents, please see S0 File.

874 **S4 File. Differentially expressed gene ontology results (Camera output).** The BCNHL
875 differentially expressed gene ontology results produced by Camera. No gene ontologies were
876 significant after performing the FDR correction. For description of column contents, please see
877 S0 File.

878 **S5 File. Differentially regulated pathway results (SPIA output).** The BCNHL differentially
879 expressed pathway results produced by SPIA. For description of column contents, please see S0
880 File.

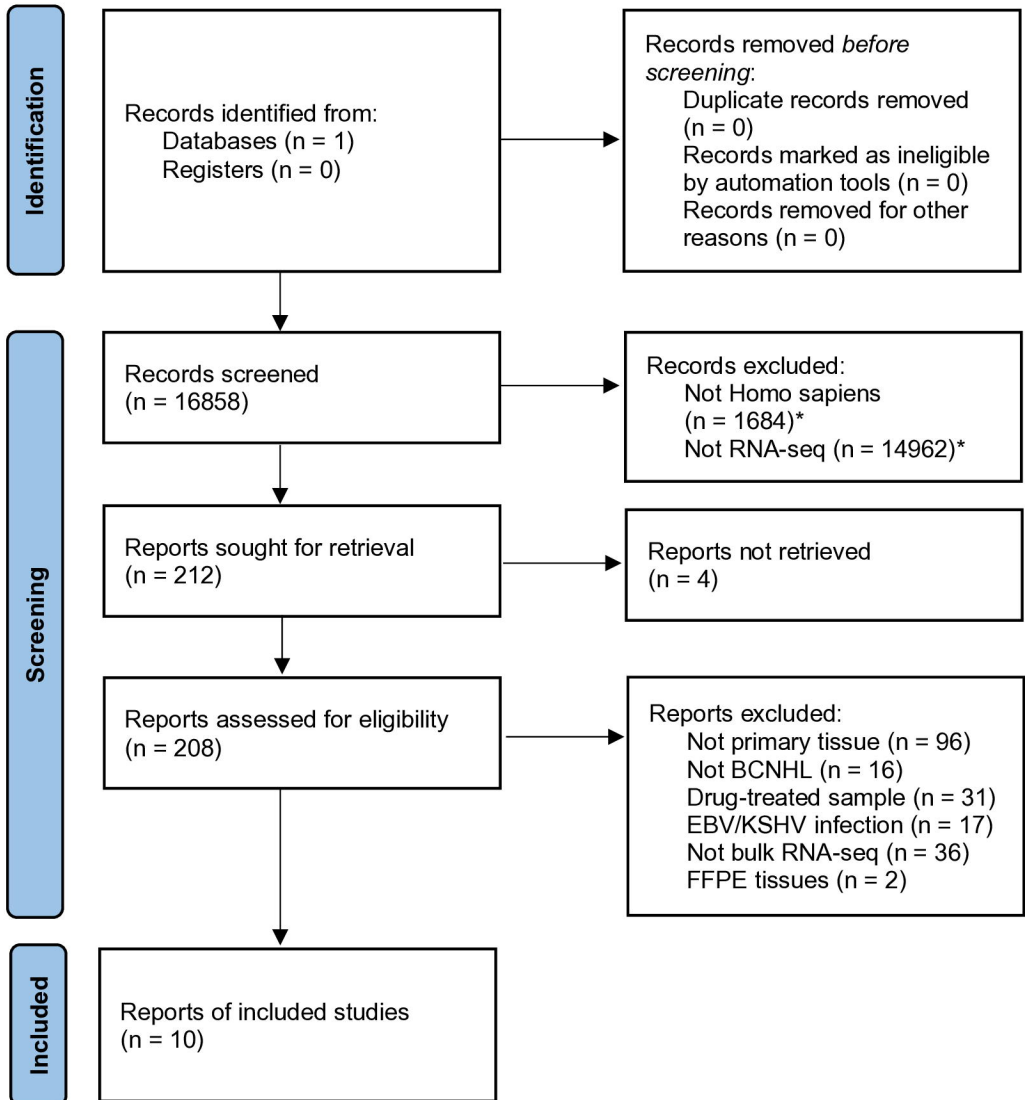
881 **S6 File. Drug prediction results by gene (Pathways2Targets unsorted output).** The raw drug
882 prediction results from Pathways2Targets2.R.

883 **S7 File. Drug prediction results sorted by most significant pathways impacted**
884 **(Pathways2Targets sorted output).** The sorted drug prediction results, ranked according to
885 which drugs impact the highest number of significantly modulated pathways.

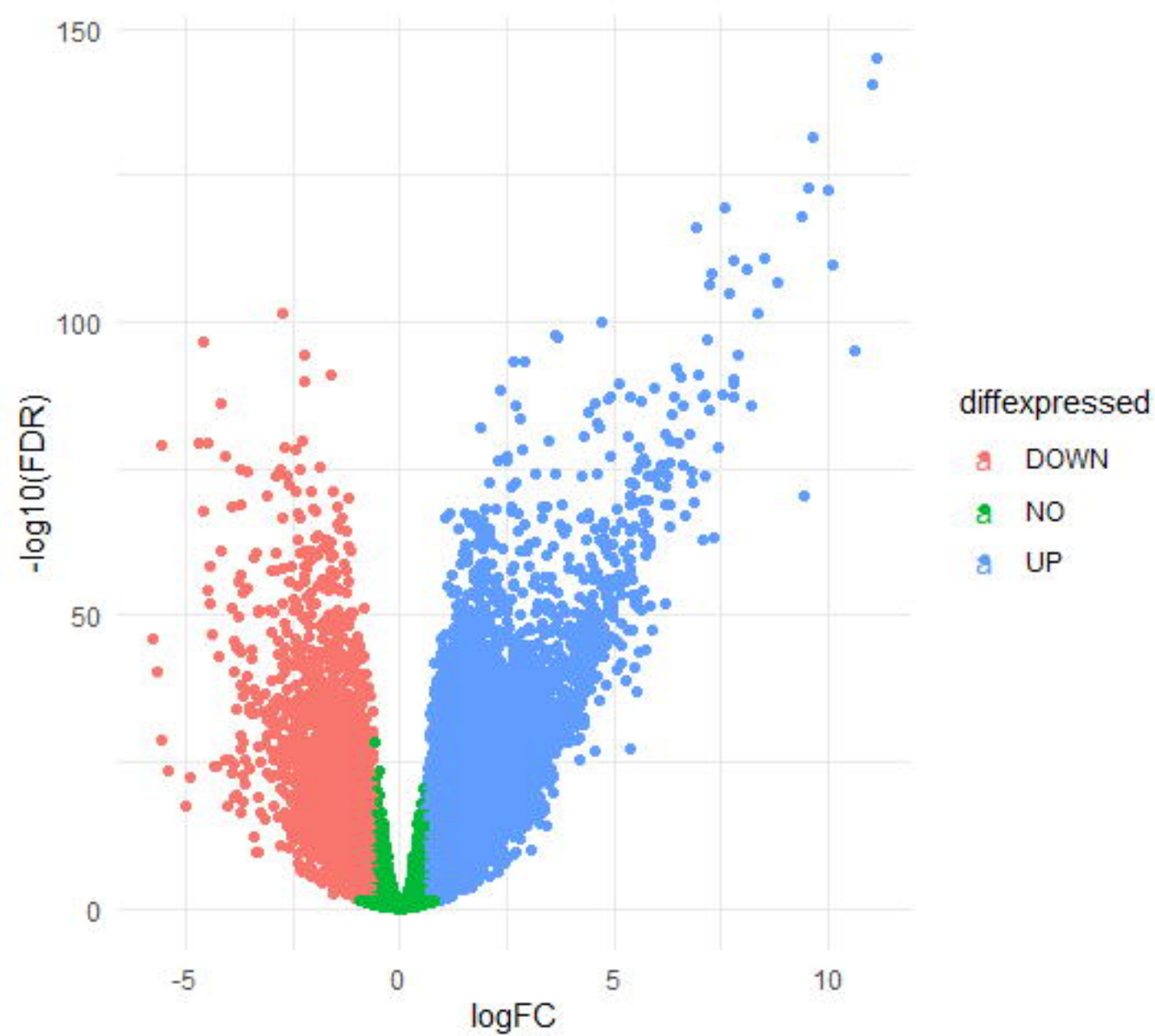
886 **S8 File. Biomarker prediction results (randomForest output).** The ranked random forest
887 biomarker prediction results. Sheet one contains all genes, and sheet two contains the random
888 forest results when the selection was narrowed to the top three genes.

889

Identification of studies via databases and registers

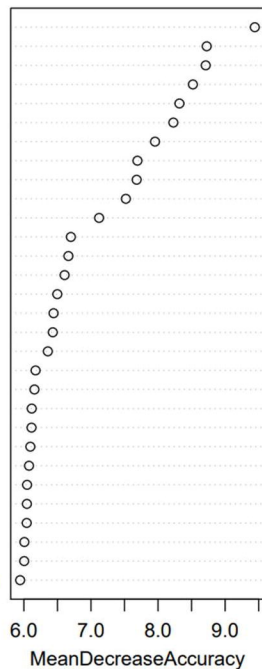


*Excluded by automation. All other excluded records were excluded by a human.

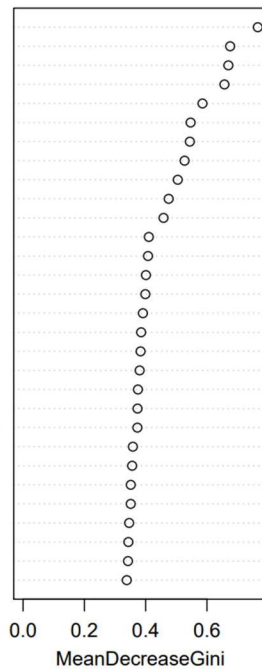


A

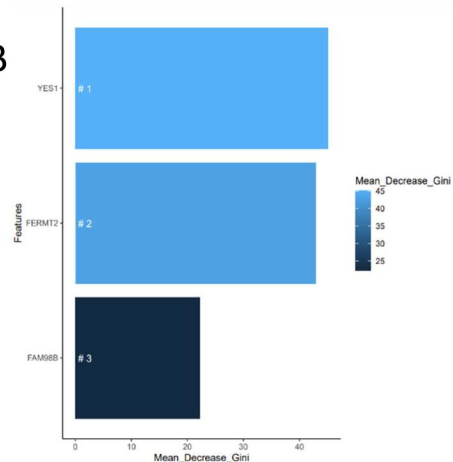
FERMT2
 EDNRB
YES1
FAM98B
 FKBP10
 LAMA4
 GJA1
 RBFOX2
 DCAF12
 CSRP2
 CDK1
 SOCS6
 PEX5
 ANLN
 UBR5
 PSTPIP2
 PBK
 ZNF664
 CENPF
 FBXO45
 KIF14
 TIMM8A
 NAPEPLD
 ZCCHC8
 PSME4
 LPAR1
 BET1
 DMXL2
 COL1A2
 VKORC1L1



YES1
FAM98B
FERMT2
 EDNRB
 DCAF12
 CSRP2
 GJA1
 LAMA4
 RBFOX2
 CDK1
 PEX5
 PBK
 ANLN
 CENPF
 UBR5
 FKBP10
 SOCS6
 FBXO45
 ZNF664
 NAPEPLD
 DMXL2
 VKORC1L1
 TGDS
 ZCCHC8
 FAM83D
 PSME4
 COL1A2
 KIF14
 BET1
 TIMM8A



B



ROC plot

C

