



SplicingFactory—splicing diversity analysis for transcriptome data

Benedek Dankó, Péter Szikora, Tamás Pór, Alexa Szeifert, Endre Sebestyén
2021

1

Students: Zahra Ahmad Nezhad, Samie Baniasadi

January 2022

What is splicingFactory?

They created an R package called SplicingFactory to meet different criteria,
They calculated RNA isoforms and examined them under different conditions.

SplicingFactory R :

The SplicingFactory R package uses transcript-level expression values to analyze splicing diversity based on various statistical measures, like Shannon entropy or the Gini index.

Items that have been tested using packages on a variety of isoforms:

- the effect of RNA-seq quantification tools
- quantification uncertainty
- gene expression levels
- isoform numbers

alternative splicing

- is well-known and described in most eukaryotic organisms
- expands the RNA repertoire of most genes

outline

Materials and methods

dataset processing

Comparison

Implementation

Results

Myelodysplastic syndrome dataset processing

- data from SRA (SRP133442, SRP149374) using SRA-tools
- they quantified transcript-level expression using Kallisto and Salmon and the full GENCODE v34 transcriptome annotation.
- Salmon-SAF was run both in alignment-based using the STAR alignments and mapping-based mode using a decoy-aware transcriptome index.
- they run the tool with 100 bootstraps, the unstranded paired-end library option and the additional `-seqBias` and `-gcBias` parameters

Comparison of Kallisto, Salmon and Salmon-SAF results

- they selected the 17 control samples from the SRP133442 dataset.
- calculated the normalized naive-entropy and Gini-index for all genes.
- they selected all genes with non-NA diversity values and calculated their Spearman correlation.

Performance benchmarks

- ▶ using a single core on a server with Intel® Xeon® Gold 4118 2.30 GHz CPU type, a total of 157 GB memory and Ubuntu 18.04.5 LTS operating system.
- ▶ they used the Salmon-based quantification for the benchmarks.
- ▶ Using a fixed number of 60 669 starting genes, they increased the sample number from 10 to 130 using steps of 10, 20, 40, 60, 80, 100 and 130 for the same calculations.

Comparison to other tools

splicehetero

- ran SpliceHetero with parameters `-slb True` and `-prn 17`
- To prepare the input data keeping only those splice junctions that appear in the STAR index `sjdbList.fromGTF.out.tab` file
- converted the filtered STAR junction files
- ran SpliceHetero on these sorted BED files

SEVA

- ran SEVA in RStudio Server to `FALSE`
- To prepare the input data we first performed RPM normalization on the STAR splice junction files

Comparison to other tools

whippet

- ran Whippet with default parameters
- To prepare the input data, we converted the raw fastq files to have only a single '+' character on every third line
- using Whippet's whippet-index.jl script, created an index file based on the GENCODE v34 transcriptome annotation gtf file, and ran Whippet's quantification algorithm with the whippet-quant.jl script with the additional -biascorrect parameter

SF3B1 differential diversity and enrichment analysis

analyzed the SRP149374 dataset

compared the MDS samples without any known mutation to MDS samples with known *SF3B1* somatic mutations

11

calculated the normalized naive-entropy, normalized Laplace-entropy and Gini-index using the Salmon-based quantification

compared the *SF3B1* mutated and *SF3B1* wild-type samples using Wilcoxon-test and performed *P*-value adjustment with the Benjamini-Hochberg method.

performed enrichment analysis using the significant genes ($|\text{mean difference}| > 0.1$ and adjusted *P*-value < 0.05) separately with either a mean diversity increase or decrease between the *SF3B1* mutant and *SF3B1* wild-type groups

used the bone marrow marker gene sets originating

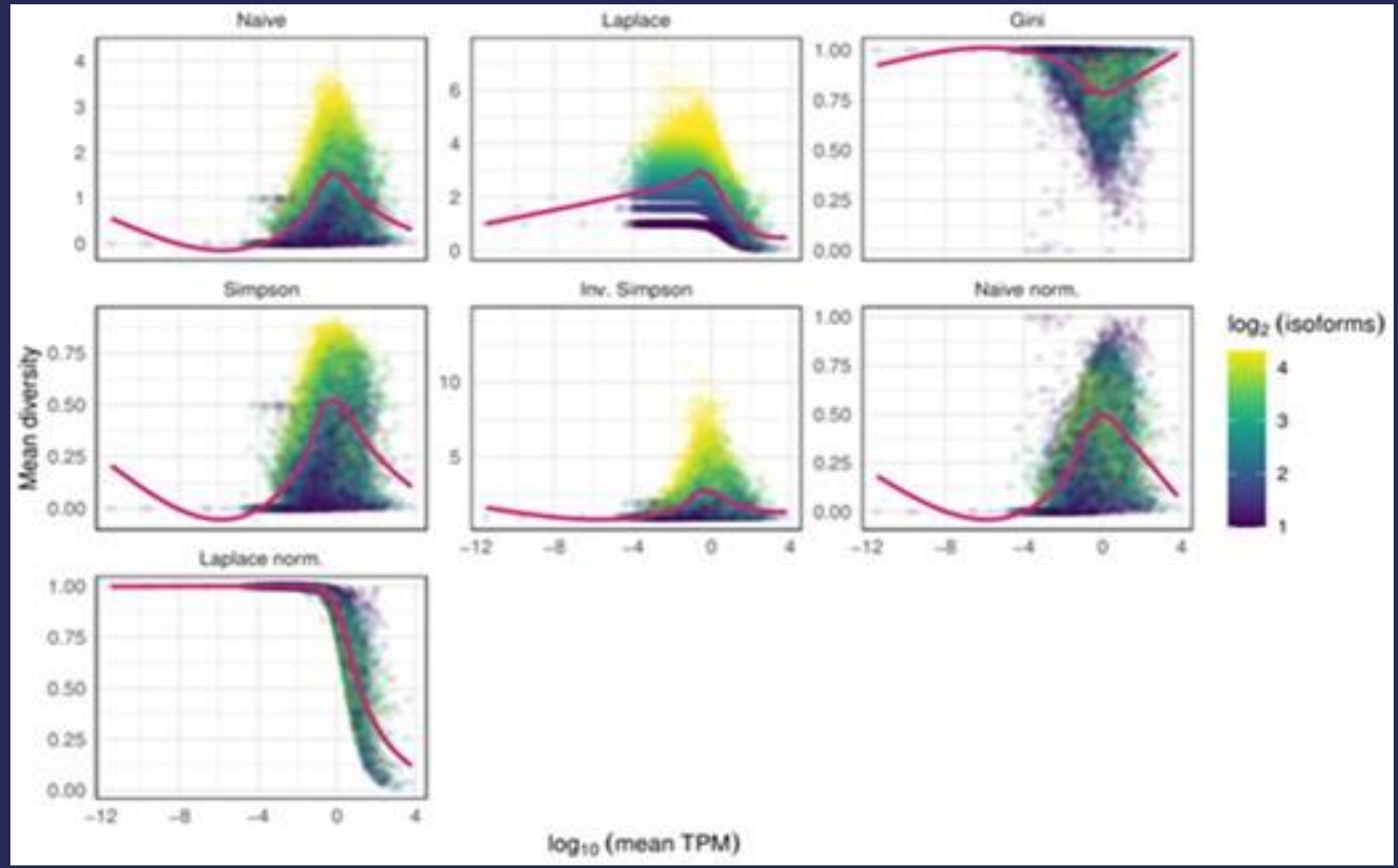
- the package calculates a diversity value for each gene in each sample, using splicing isoform expression values.
- calculates differential diversity results between conditions.

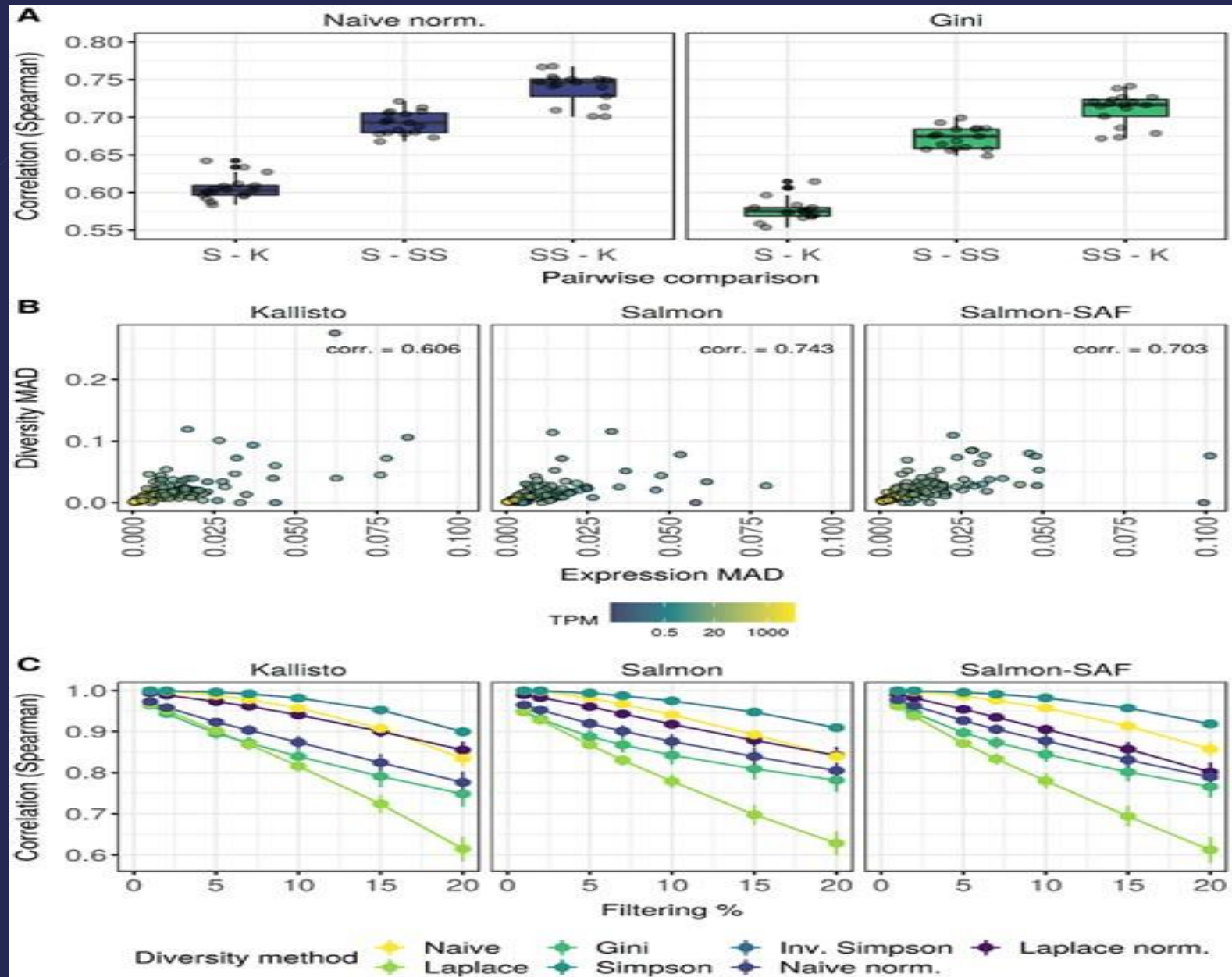
Association of diversity metrics

The Gini-index shows an opposing pattern compared to all others.

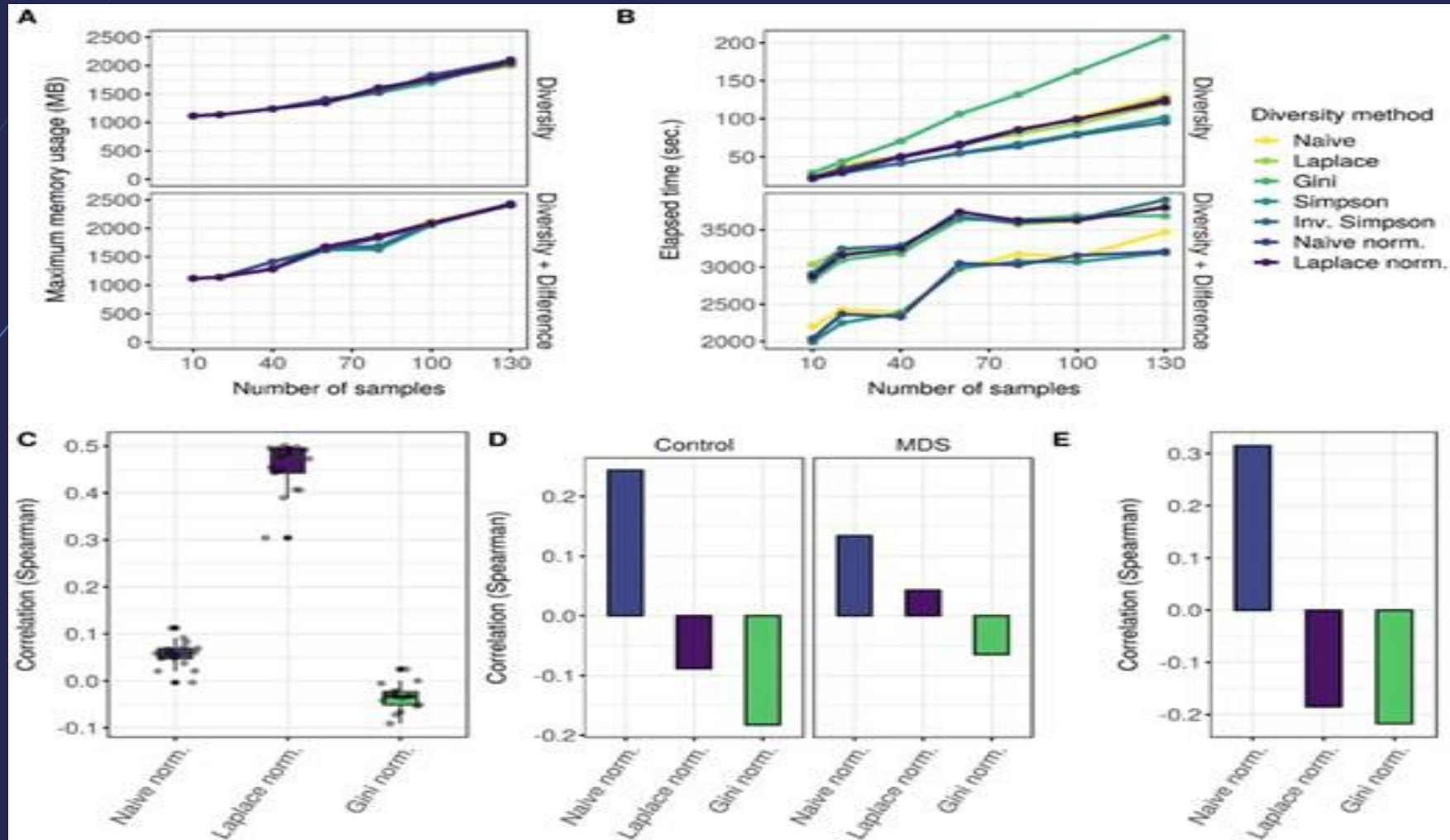
The nonnormalized naive- and Laplace-entropy, besides the Simpson-index and the inverse Simpson-index shows a clear association with isoform number.

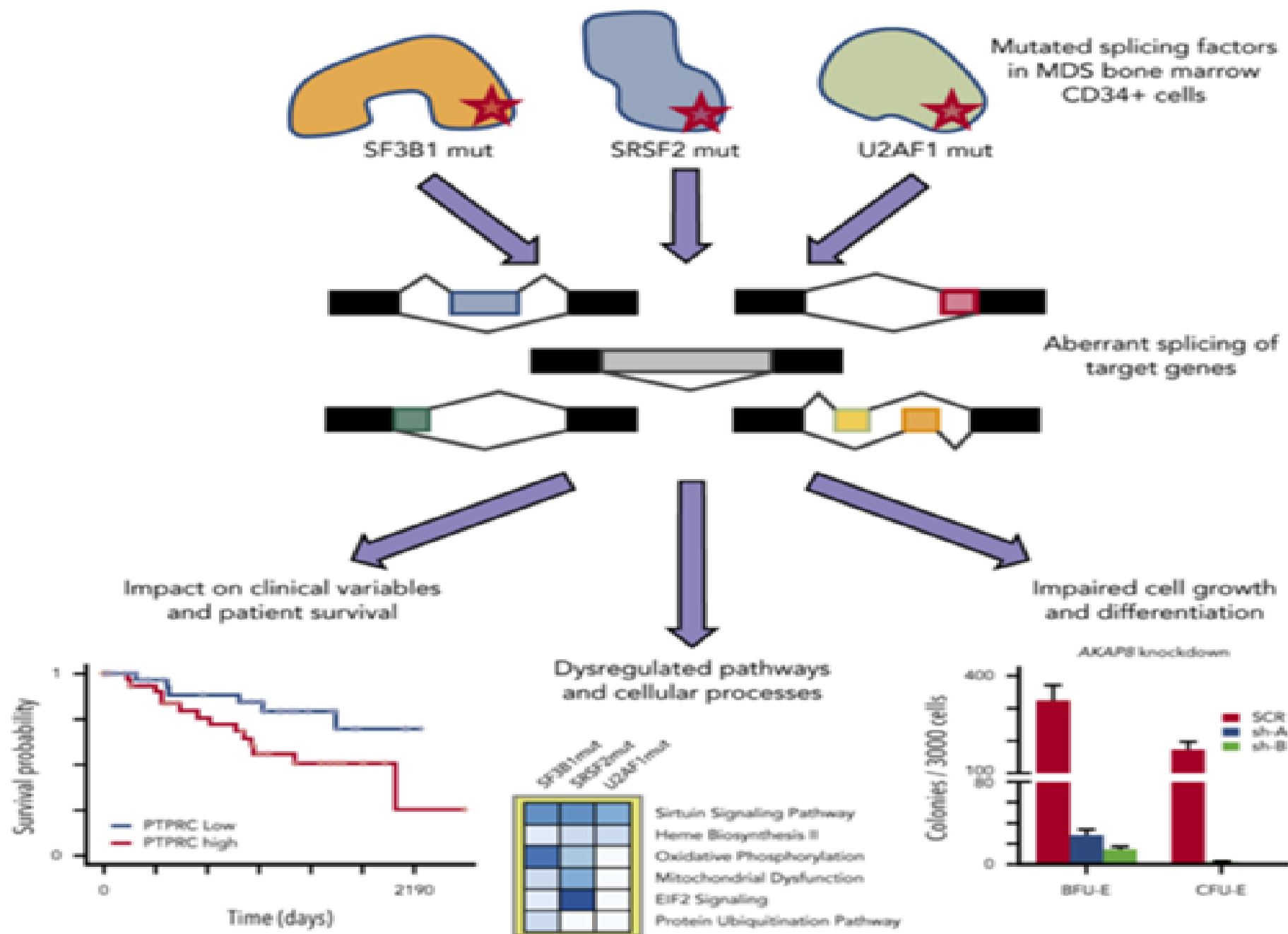
the strongly biased pattern for the normalized Laplace-entropy is again the result of the +1 pseudocount, as all genes with only zero or very low isoform expression levels are given an expression of ~ 1 , leading to the maximum possible entropy of 1.

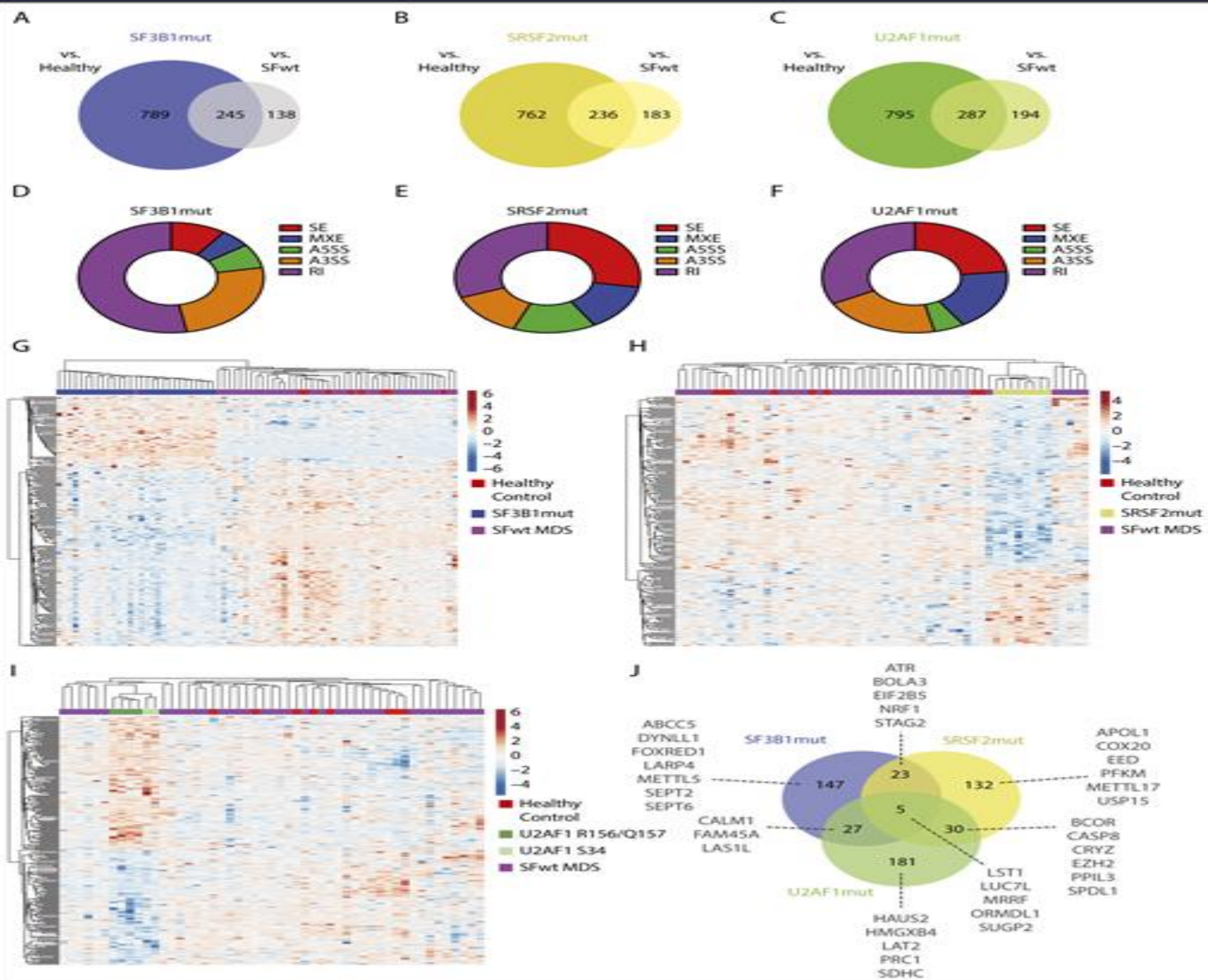


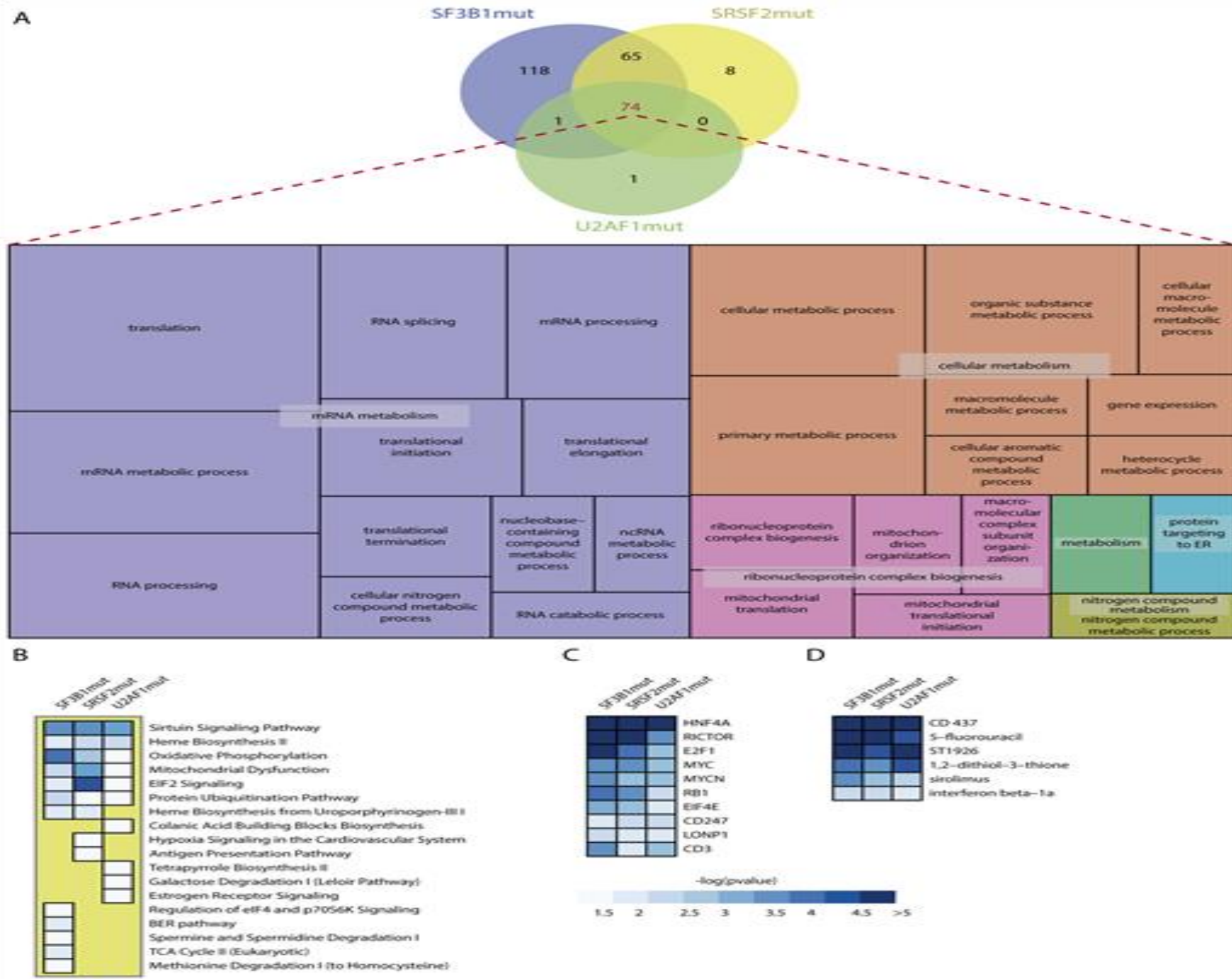


Performance benchmarks









List of aberrant splicing events, identified in *SF3B1*mut MDS, significantly correlated with a clinical variable

20

Event ID	Gene	Event type	Chr	Strand	Start position	End position	Spearman correlation variable	cor_ estimate	P	adj.P
6795	PARVG	A3SS	22	+	44582456	44583758	BM blast %	0.58	4.38 E-07	.003031
5420	RPRD1A	RI	18	-	33605560	33607038	BM blast %	0.58	6.21 E-07	.004292
3146	DOM3Z	RI	6	-	31938382	31938924	ANC	-0.57	1.00 E-06	.006916
12280	CXXC1	SE	18	-	47811694	47811721	ANC	-0.56	1.28 E-06	.00888

3305	AP1 G2	RI	14	–	2403 1170	2403 1624	ANC	–0.5 6	1.68 E-06	.011 579
2362	SNR PN	A3SS	15	+	2521 9434	2521 9603	ANC	–0.5 5	3.14 E-06	.021 689
5460	TCE A2	A3SS	20	+	6270 3210	6270 3294	ANC	0.54	3.24 E-06	.022 393
6627	NIC N1	A3SS	3	–	4946 2381	4946 2579	Pl†	–0.5 4	3.80 E-06	.026 281
6153	ABC C5	A3SS	3	–	1837 0309 1	1837 0324 3	ANC	–0.5 4	4.23 E-06	.029 229
518	ERC C3	RI	2	–	1280 4691 2	1280 4740 0	ANC	–0.5 4	4.71 E-06	.032 551
2359	SNR PN	A3SS	15	+	2521 9457	2521 9603	ANC	–0.5 4	5.02 E-06	.034 68

7563	PPO X	A3SS	1	+	1611 3712 8	1611 3727 6	ANC	0.53	5.99 E-06	.041 411
4975	GPR 108	A3SS	19	–	6730 997	6731 122	ANC	–0.5 3	6.43 E-06	.044 443
5816	PSTPI P1	A3SS	15	+	7732 8142	7732 8276	ANC	–0.5 3	6.49 E-06	.044 832
4728	NIC N1	RI	3	–	4946 2381	4946 2871	Plt	–0.5 3	6.75 E-06	.046 655

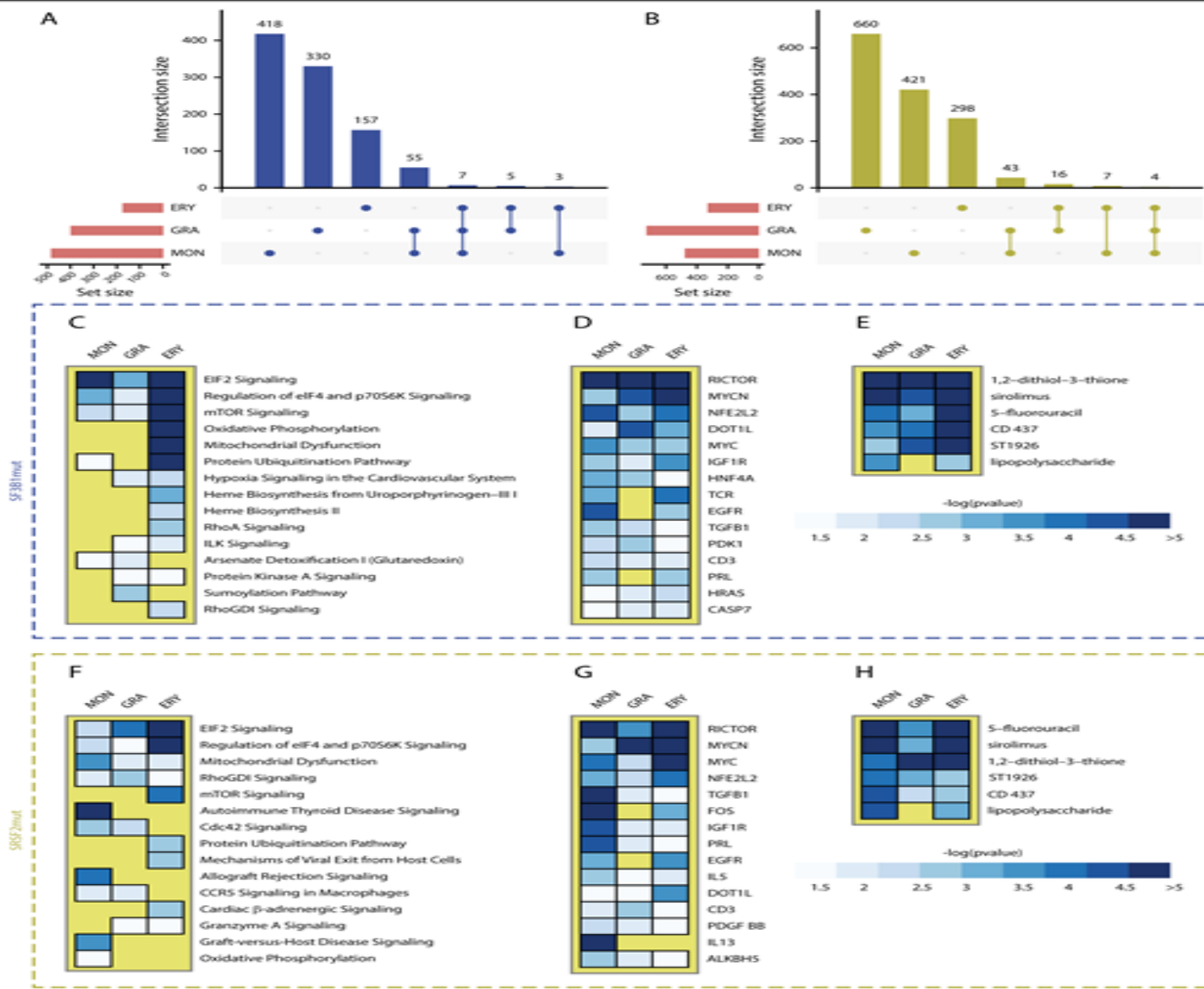
Genes with isoforms that significantly predicted survival in MDS in multivariate models

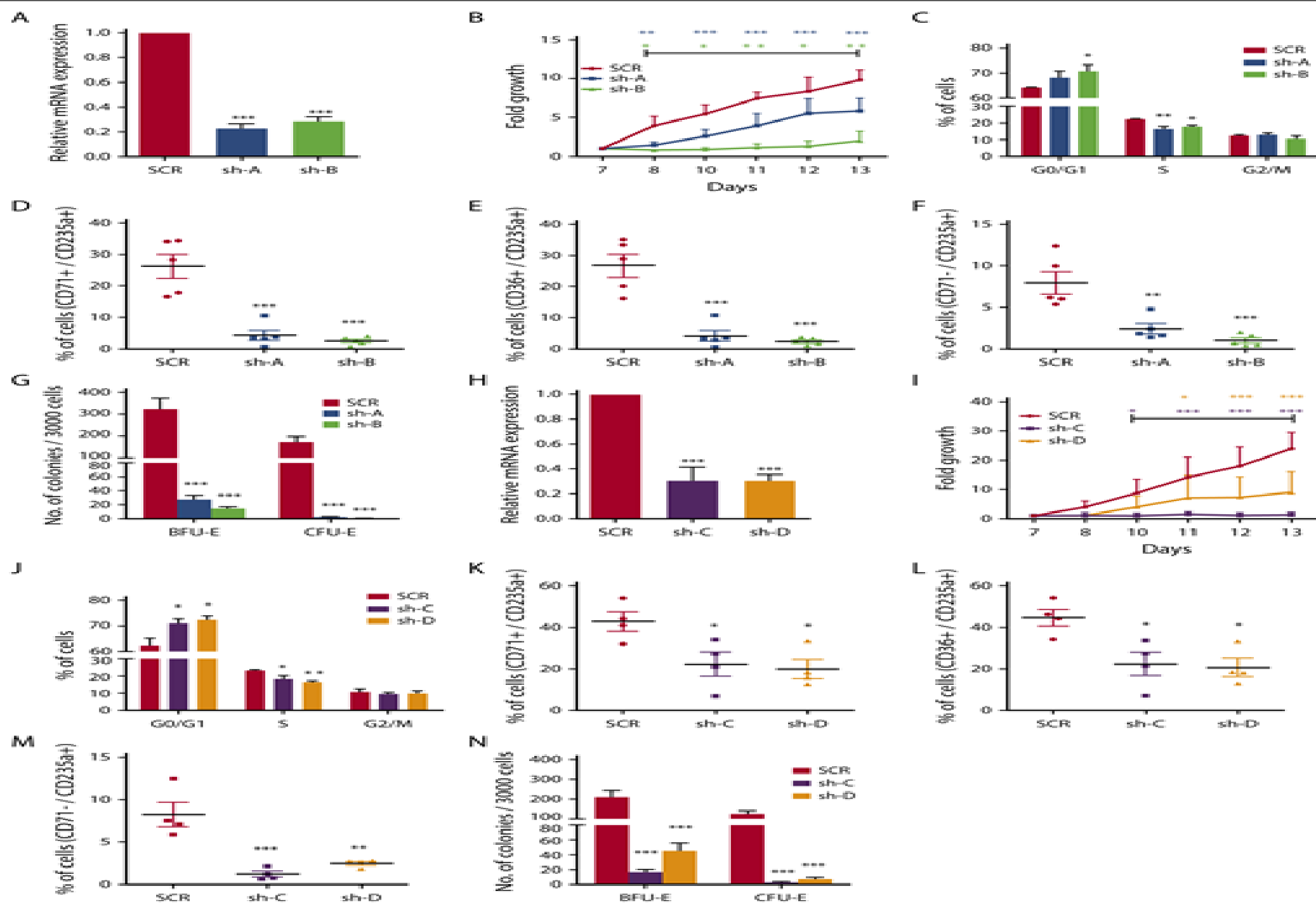
23

Gene	Associated splice factor mutation or mutations	Function/pathway	P (multivariate survival)	
CAP1	SRSF2	Focal adhesion & extracellular exosomes	.0044	
PTPRC	U2AF1 (S34)	Focal adhesion & extracellular exosomes	.0093	
IFI44	SRSF2, U2AF1 (S34)	Interferon, extracellular exosomes	.012	
IFI44L	U2AF1 (S34)	Interferon, extracellular exosomes	.0086	

CD46	U2AF1 (S34)	Focal adhesion & extracellular exosomes	.039
CRTC2	SF3B1	Extracellular exosomes	.035
FCGR2A	U2AF1 (S34)	Extracellular exosomes	.016
PPOX	SF3B1, U2AF1 (R156/Q157)	Heme biosynthesis	.031
AHSA2	SF3B1	HSP90 ATPase	.029
DHP5	SF3B1	Translation elongation factor 2 modification	.026
MECR	U2AF1 (S34 & R156/Q157)	Mitochondrial reductase	.022

NASP	U2AF1 (S34)	HSP90 binding	.014
PFDN5	U2AF1 (R156/Q157)	Prefoldin subunit	.042
PABPC4	U2AF1 (R156/Q157)	NMD mRNA decay	.036





Thank you