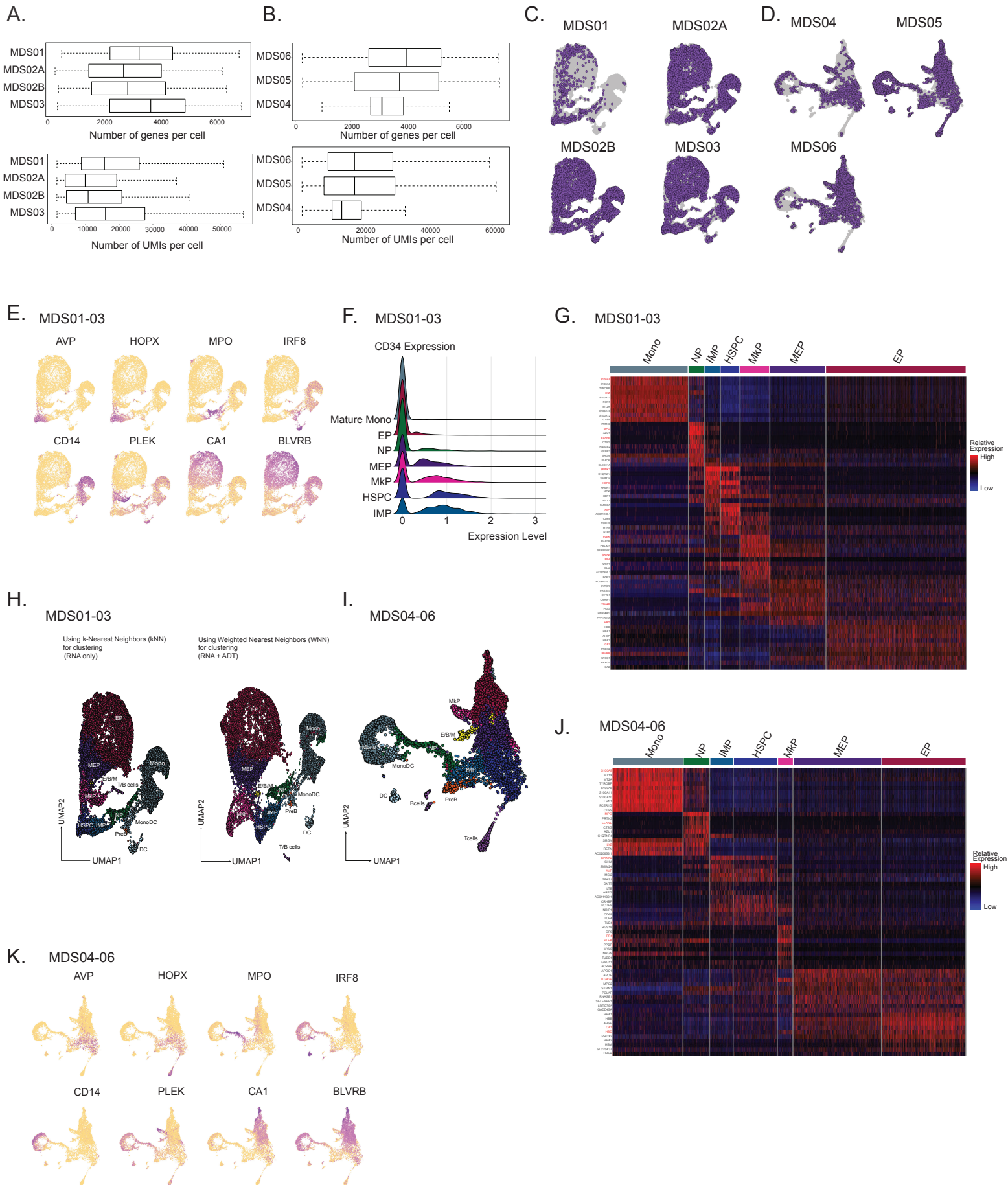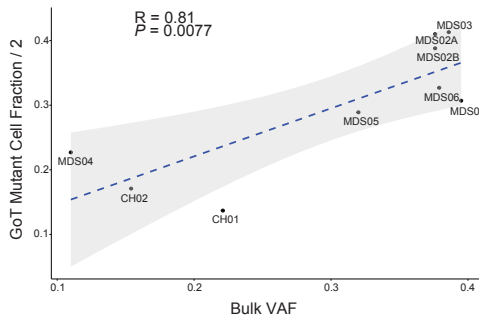# Figure S1

**Figure S1. MDS samples QC, integration, clustering, and cell-type assignment, related to Figure 1.**
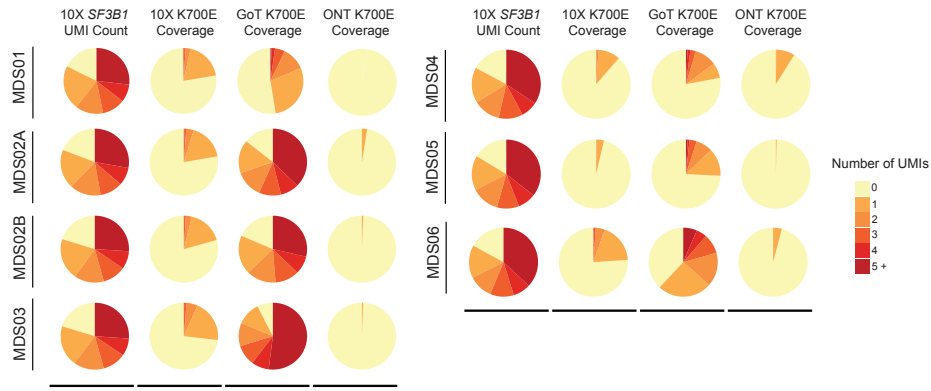
**(A)** Number of genes per cell (top) and number of UMIs per cell (bottom) in CD34+ sorted hematopoietic progenitors from samples MDS01-03 after QC filters, shown by each patient sample. **(B)** Number of genes per cell (top) and number of UMIs per cell (bottom) in CD34+ sorted hematopoietic progenitors from samples MDS04-06 after QC filters, shown by each patient sample. **(C)** UMAP of CD34+ sorted progenitor cells for each individual sample of MDS01-03 after integration using the Seurat package. **(D)** UMAP of CD34+ sorted progenitor cells for each individual sample of MDS04-06 after integration using the Seurat package. **(E)** Expression of lineage-specific genes from Velten *et al.*[S1] scored and projected onto the UMAP representation of cells from MDS01-03. **(F)** CD34 expression per progenitor cell-type of CD34- monocytes among CD34+ sorted hematopoietic progenitors. **(G)** Heatmap of top 10 differentially expressed genes for each progenitor subset for MDS01-03. **(H)** UMAPs comparing the graph-based clustering output when using the k-nearest neighbors (KNN) algorithm to perform clustering of cells with RNA data alone vs. when using the weighted nearest neighbors (WNN) algorithm that allows for the integration of both RNA and ADT data for clustering cells. The cell-type assignments determined from the KNN-RNA based clustering are projected onto the WNN-RNA+ADT clustering for comparison between the two methods. **(I)** UMAP of CD34+ sorted cells ($n$ = 8,879 cells) from samples MDS04-06 with *SF3B1* K700E mutations ($n$ = 3), overlaid with cluster cell-type assignments. HSPC, hematopoietic stem progenitor cells; IMP, immature myeloid progenitors; MkP, megakaryocytic progenitors; MEP, megakaryocytic-erythroid progenitors; EP, erythroid progenitors; NP, neutrophil progenitors; E/B/M, eosinophil/basophil/mast progenitor cells; T/B cell progenitors; Mono, monocyte; DC, dendritic cells; Pre-B, precursors B cells; Mono DC, monocyte/dendritic cell progenitors. **(J)** Heatmap of top 10 differentially expressed genes for each progenitor subset for MDS04-06. **(K)** Expression of lineage-specific genes from Velten *et al.*[S1] scored and projected onto the UMAP representation of cells from MDS04-06.
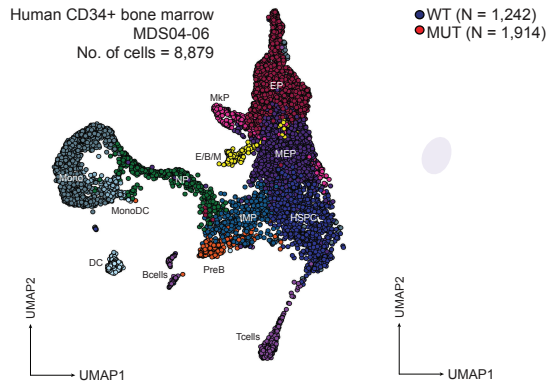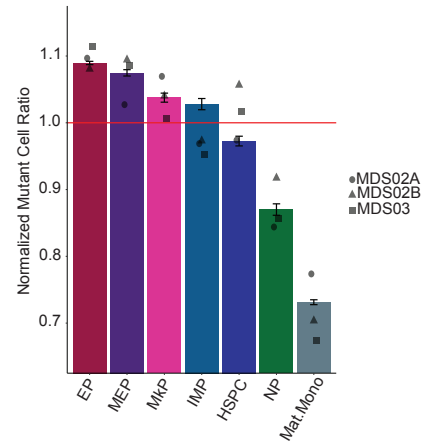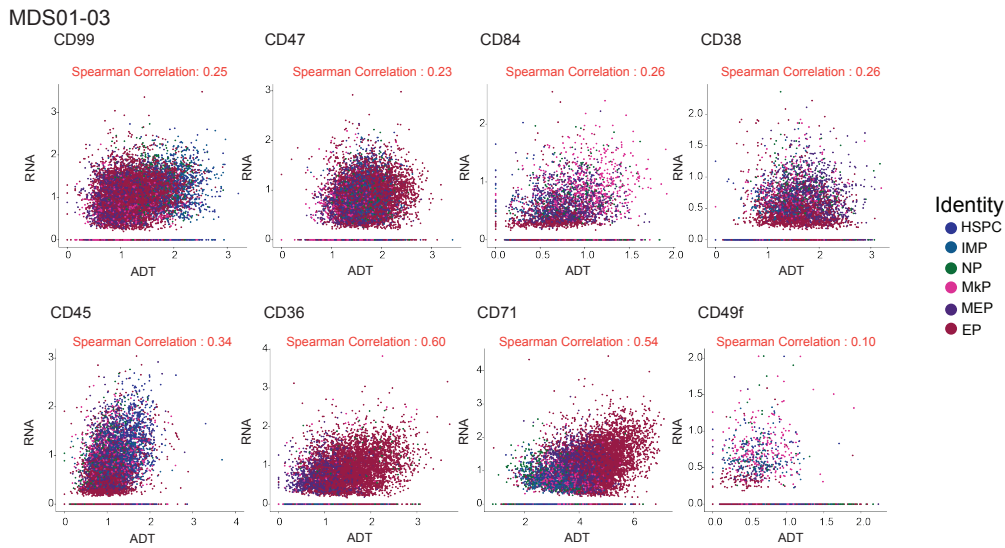
# Figure S2



**A.** R = 0.81, P = 0.0077

**B.** 10X *SF3B1* UMI Count, 10X K700E Coverage, GoT K700E Coverage, ONT K700E Coverage

MDS01, MDS02A, MDS02B, MDS03, MDS04, MDS05, MDS06

Number of UMIs: 0, 1, 2, 3, 4, 5 +

**C.** Human CD34+ bone marrow MDS04-06, No. of cells = 8,879

WT (N = 1,242), MUT (N = 1,914)

**D.** Normalized Mutant Cell Ratio

EP, MEP, MkP, IMP, HSPC, NP, Mat.Mono

MDS02A, MDS02B, MDS03

**E.** MDS01-03

CD99 — Spearman Correlation: 0.25
CD47 — Spearman Correlation : 0.23
CD84 — Spearman Correlation : 0.26
CD38 — Spearman Correlation : 0.26
CD45 — Spearman Correlation : 0.34
CD36 — Spearman Correlation : 0.60
CD71 — Spearman Correlation : 0.54
CD49f — Spearman Correlation : 0.10

Identity: HSPC, IMP, NP, MkP, MEP, EP

**F.** MDS01-03

CD99 — r = 0.66, P = 0.15
CD47 — r = 0.67, P = 0.15
CD84 — r = 0.87, P = 0.025
CD38 — r = 0.87, P = 0.025
CD45 — r = 0.94, P = 0.0047
CD36 — r = 0.98, P = 0.00075
CD71 — r = 0.94, P = 0.0053
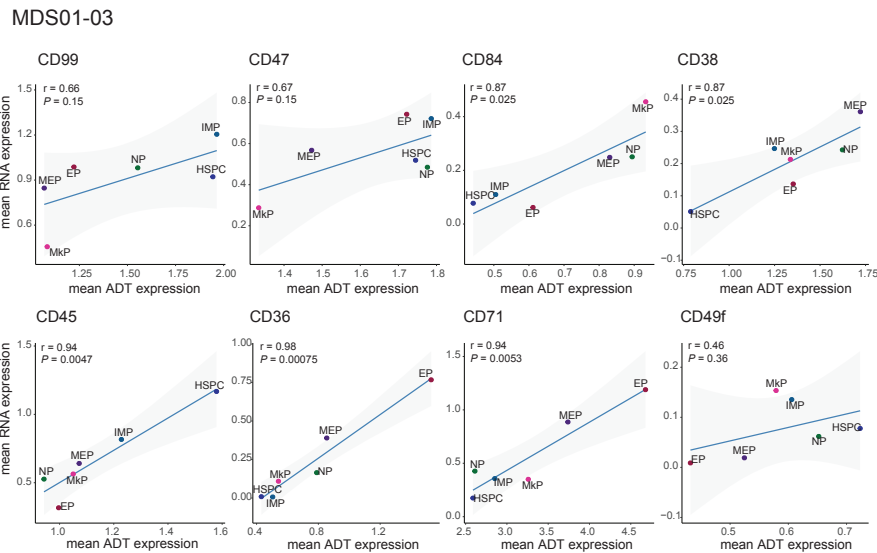CD49f — r = 0.46, P = 0.36

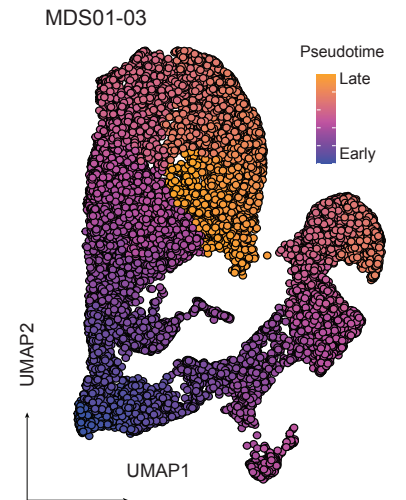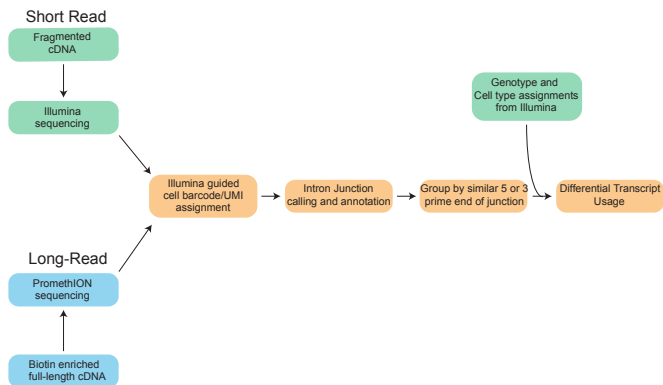**G.** MDS01-03

Pseudotime: Late, Early

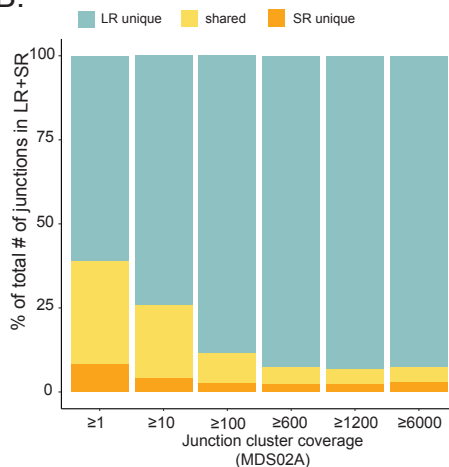**Figure S2. MDS samples GoT, CITE-seq, and pseudotime analyses, related to Figure 2.**

**(A)** *SF3B1* K700E (7) and K666N (1) mutant cell fractions determined by GoT in single cells versus *SF3B1* K700E and K666N mutation variant allele frequencies (VAF) determined in bulk sequencing of matched unsorted bone marrow mononuclear cells (MDS) or matched unsorted stem cell product (CH). **(B)** Fraction of cells in MDS samples by number of *SF3B1* UMIs in standard 10x Genomics data without genotyping information, *SF3B1* UMIs with K700E locus coverage in standard 10x data, *SF3B1* UMIs with K700E locus coverage in GoT amplicon library, and *SF3B1* UMIs with K700E locus coverage in ONT library. **(C)** UMAP of CD34+ sorted cells (*n* = 8,879 cells) from samples MDS04-06 with *SF3B1* K700E mutations (*left*) and density plot of $SF3B1^{mut}$ vs. $SF3B1^{wt}$ cells (*right*). WT, cells with genotype data without *SF3B1* mutation; MUT, cells with genotype data with *SF3B1* mutation. **(D)** Normalized ratio of $SF3B1^{mut}$ cells in progenitor subsets with at least 300 genotyped cells. Bars show aggregate analysis of samples MDS01-03 with mean +/- s.e.m. of 100 downsampling iterations to 1 genotyping UMI per cell. Points represent the mean of *n* = 100 downsampling iteration per sample. **(E)** Comparison of the single-cell expression of markers captured in both CITE-seq (x-axes) and RNA-seq (y-axes) libraries. Correlation coefficient r calculated using Spearman's correlation. Cells are colored by each progenitor subset. **(F)** Comparison of the mean expression per progenitor subset of markers captured in both CITE-seq (x-axes) and RNA-seq (y-axes) libraries. Correlation coefficient r calculated using Spearman's Correlation. *P*-values derived from Student's t-distribution. **(G)** UMAP of progenitor cells from MDS01-03 samples overlaid with pseudotemporal ordering.
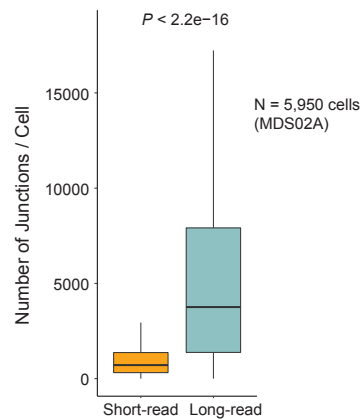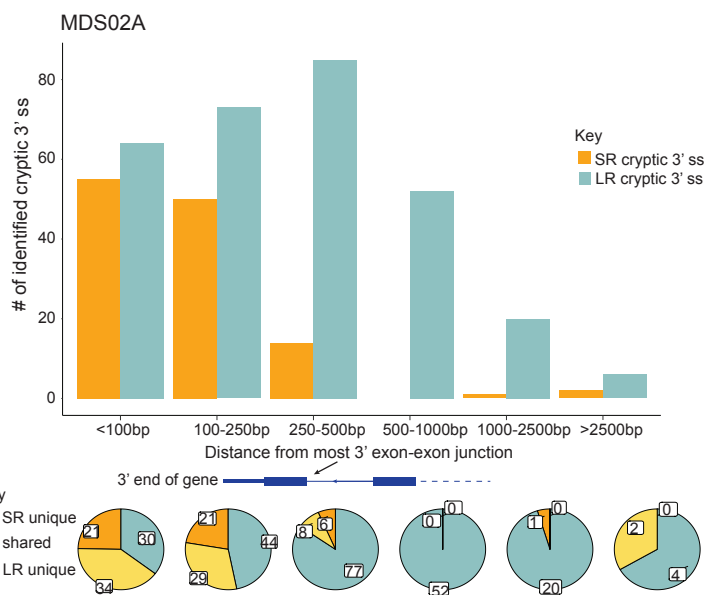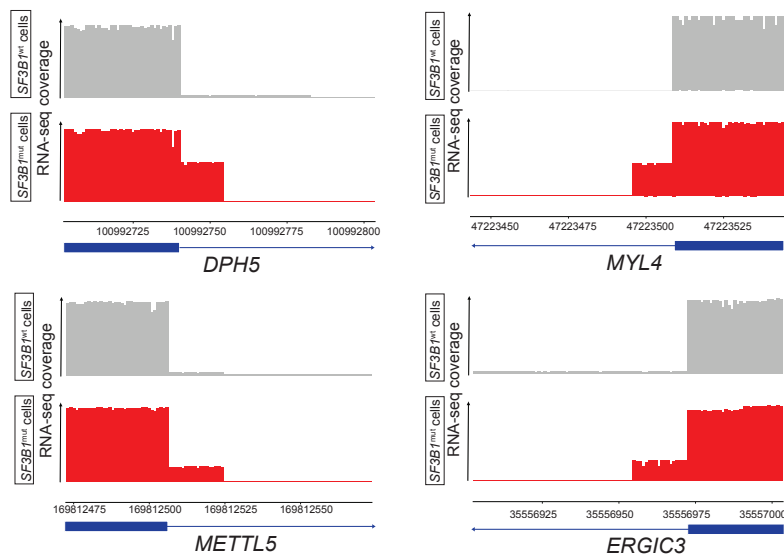
# Figure S3



**A.**

Short Read
Fragmented cDNA → Illumina sequencing

Illumina guided cell barcode/UMI assignment → Intron Junction calling and annotation → Group by similar 5 or 3 prime end of junction → Differential Transcript Usage

Genotype and Cell type assignments from Illumina

Long-Read
PromethION sequencing

Biotin enriched full-length cDNA

**B.**

LR unique / shared / SR unique

% of total # of junctions in LR+SR

Junction cluster coverage (MDS02A): ≥1, ≥10, ≥100, ≥600, ≥1200, ≥6000

**C.**

P < 2.2e−16

Number of Junctions / Cell

N = 5,950 cells (MDS02A)

Short-read / Long-read

**D.**

MDS02A

# of identified cryptic 3' ss

Key
SR cryptic 3' ss
LR cryptic 3' ss

Distance from most 3' exon-exon junction: <100bp, 100-250bp, 250-500bp, 500-1000bp, 1000-2500bp, >2500bp

3' end of gene

Key
SR unique
shared
LR unique

21 / 30 / 34  |  21 / 29 / 44  |  8 / 6 / 77  |  0 / 0 / 52  |  1 / 0 / 20  |  2 / 0 / 4

**E.**

SF3B1^wt cells RNA-seq coverage
SF3B1^mut cells RNA-seq coverage

100992725  100992750  100992775  100992800
*DPH5*

47223450  47223475  47223500  47223525
*MYL4*

169812475  169812500  169812525  169812550
*METTL5*

35556925  35556950  35556975  35557000
*ERGIC3*

**F.**

Density

cryptic 3' ss 1
cryptic 3' ss 2

*PFDN5*

Key
Transcripts with cryptic 3' ss 1 only
Transcripts with cryptic 3' ss 2 only
Transcripts with cryptic 3' ss 1 + 2

Number of cryptic 3' splice site events per gene: 1, 2, 3, 4, 5

**G.**

MDS02

$-\log_{10}(OR)$

A / C / G / T

nucleotide distance from 3' splice site

MDS03

$-\log_{10}(OR)$
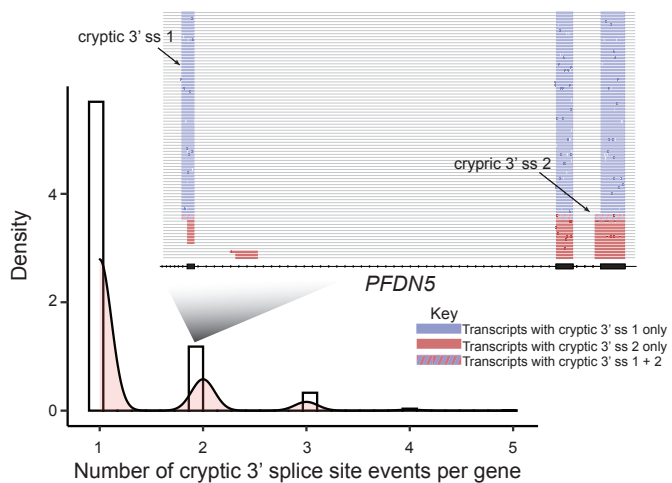
A / C / G / T

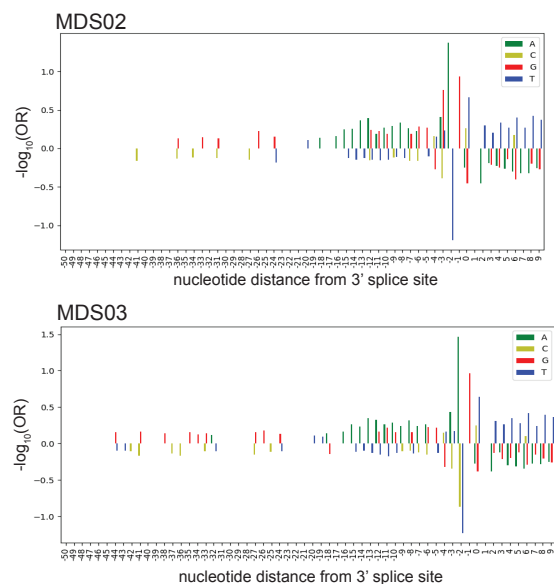nucleotide distance from 3' splice site

**Figure S3. Long-read splicing and motif analysis, related to Figure 3.**

**(A)** Long-read sequencing processing and splicing analysis pipeline. **(B)** Fraction of cryptic junctions detected in short-reads, long-reads or both when defining distinct cluster coverage thresholds. **(C)** Number of junctions per cell detected in 10x short-read compared to ONT full-length long-read (junction cluster coverage >= 1 read). **(D)** Number of detected cryptic 3' splice site junctions (in genes with ~1 read/cell in both short and long-read datasets and in clusters with coverage >= 600 reads) along increasing distances to the most 3' end junction of each transcript (see *inset*) in 10x short versus ONT full-length long-reads. Individual number of events as well as overlaps are itemized in the pie charts (bottom). **(E)** Comparison of the usage of various alternative 3' splice sites found in our MDS *SF3B1^{mut}* cells vs. a CD34+ sample with no *SF3B1* mutation. **(F)** Bar plot of the number of cryptic 3' splice sites identified per gene in MDS. *Inset*: Gene example, *PFDN5*, with 2 unique cryptic 3' splice sites, showing the transcripts that have usage of either site. **(G)** Nucleotide enrichment (measured as log-odds ratio) across the 3' splice site region comparing cryptic vs. canonical sites in MDS02-03 samples.
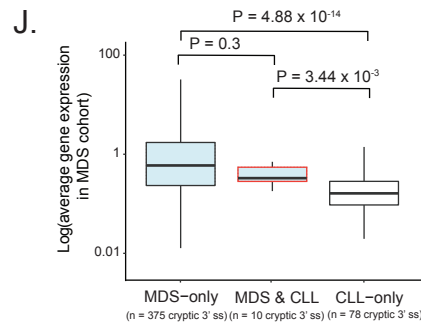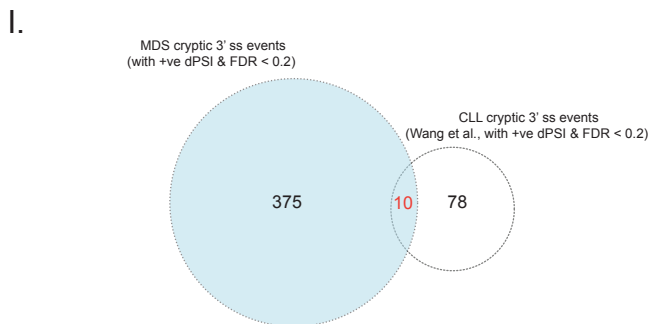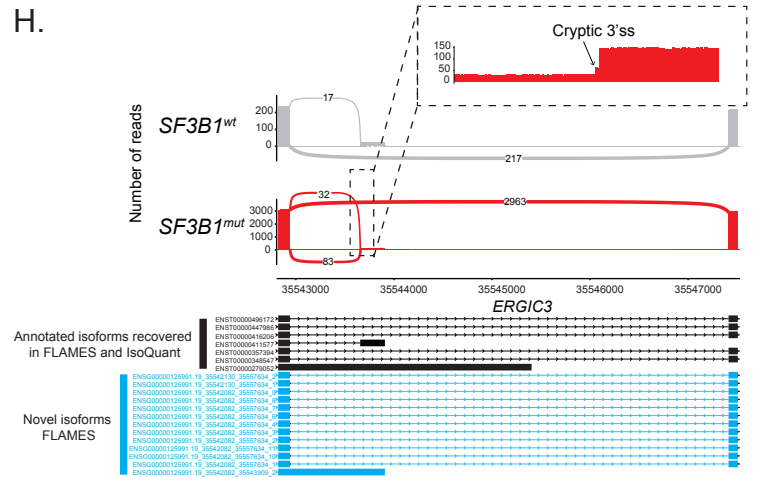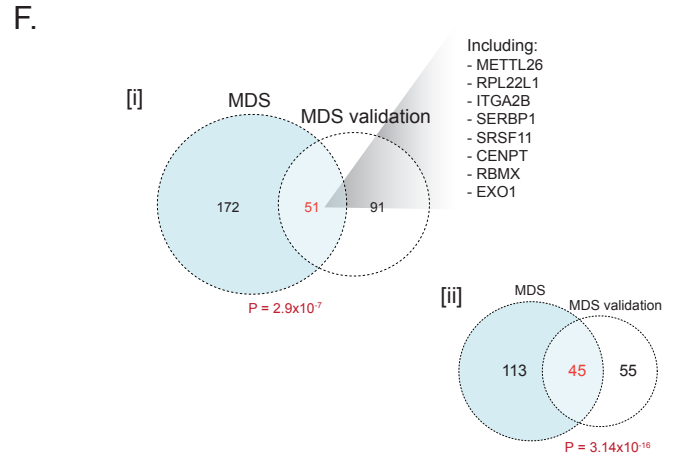
# Figure S4



**A.**

All junctions

GoT−Splice 17000

8900

70000

43000

16000

FLAMES

4700

22000

IsoQuant

Cryptic 3' ss junctions

GoT−Splice 8700

560

57 71
23 96

IsoQuant

1100

FLAMES

**B.**

IsoQuant

ρ = −0.02, P = 0.895, n = 35
y = −0.017 − 0.016 x

FLAMES

ρ = 0.45, P = 0.007, n = 35
y = −0.771 + 0.0808 x

dPSI ($SF3B1^{mut}$ -$SF3B1^{wt}$ cells) estimated from full-length isoform counts

dPSI ($SF3B1^{mut}$ -$SF3B1^{wt}$ bulk Pellagati et al.)

**C.**

MDS cohort - per patient

percentage of significant events (FDR <0.2)

MDS02A
MDS02B
MDS03

+ve dPSI
-ve dPSI

**E.**

MDS validation cohort - per patient

Percentage of significant events (FDR < 0.2)

MDS04
MDS05
MDS06

+ve dPSI
-ve dPSI

**D.**

percentage of significant events (FDR < 0.2)

MDS validation

+ve dPSI
-ve dPSI

Differentially-spliced in $SF3B1^{wt}$ cells ← → Differentially-spliced in $SF3B1^{mut}$ cells

Alternative 3' ss
Cryptic 3' ss (0-100 bp from canonical 3' ss)

−log$_{10}$(FDR)

CYLD
TMEM216
HGF/CTR1
FAM161A
ATPAF2
CKAP2
SERBP1
RHNO1
CDC7
ESYT1
DHX35
TCEA2
EXO1
KDM1A
MYL4
KLHDC1
PPP2R5A
PDGFB

dPSI ($SF3B1^{mut}$ - $SF3B1^{wt}$)

**F.**

[i]

MDS    MDS validation

172    51    91

P = 2.9x10$^{-7}$

Including:
- METTL26
- RPL22L1
- ITGA2B
- SERBP1
- SRSF11
- CENPT
- RBMX
- EXO1

[ii]

MDS    MDS validation

113    45    55

P = 3.14x10$^{-16}$

**G.**

% significant Alternative/Cryptic 3' ss events

IsoQuant
FLAMES

Upregulated
Downregulated

Isoforms with Alternative / Cryptic 3'ss detected in GoT-Splice

Not Shared
Shared

FLAMES

No. of isoforms: 58,753
No. of isoforms containing shared cryptic 3'ss: 7,706

−log$_{10}$(FDR)

FCGRT
RPL22L1
TIMM
PFDN5
SERPINF1
STX10
FCGRT
FCGRT
IMS1

log(FC Isoform proportion [$SF3B1^{mut}/SF3B1^{wt}$])

IsoQuant

No. of isoforms: 2,017
No. of isoforms containing shared cryptic 3'ss: 238

−log$_{10}$(FDR)

FCGRT
FCGRT
NACA
ATP5MC1
RBMX
RPL30

log(FC Isoform proportion [$SF3B1^{mut}/SF3B1^{wt}$])

**H.**

Cryptic 3'ss

$SF3B1^{wt}$

17

217

$SF3B1^{mut}$

32

2963

83

Number of reads

35543000    35544000    35545000    35546000    35547000

*ERGIC3*

Annotated isoforms recovered in FLAMES and IsoQuant

ENST00000496172
ENST00000447966
ENST00000416206
ENST00000411577
ENST00000357394
ENST00000348547
ENST00000279052

Novel isoforms FLAMES

ENSG00000129991_19_35542130_35557634
ENSG00000129991_19_35542130_35557634
ENSG00000129991_19_35557634
ENSG00000129991_19_35542082_35557634
ENSG00000129991_19_35542082_35557634
ENSG00000129991_19_35542082_35557634
ENSG00000129991_19_35542082_35557634
ENSG00000129991_19_35542082_35557634
ENSG00000129991_19_35542082_35557634
ENSG00000129991_19_35542082_35557634
ENSG00000129991_19_35542082_35557699
ENSG00000129991_19_35542082_35553909

**I.**

MDS cryptic 3' ss events (with +ve dPSI & FDR < 0.2)

375    10    78

CLL cryptic 3' ss events (Wang et al., with +ve dPSI & FDR < 0.2)

**J.**

P = 4.88 x 10$^{-14}$

P = 0.3

P = 3.44 x 10$^{-3}$

Log(average gene expression in MDS cohort)

MDS−only (n = 375 cryptic 3' ss)    MDS & CLL (n = 10 cryptic 3' ss)    CLL−only (n = 78 cryptic 3' ss)
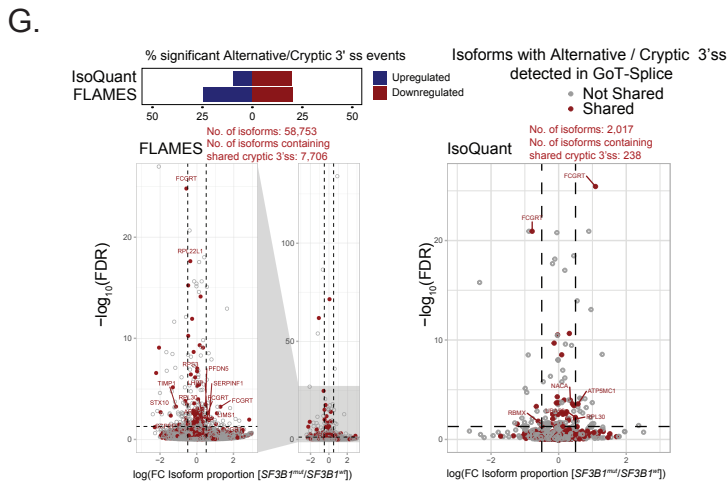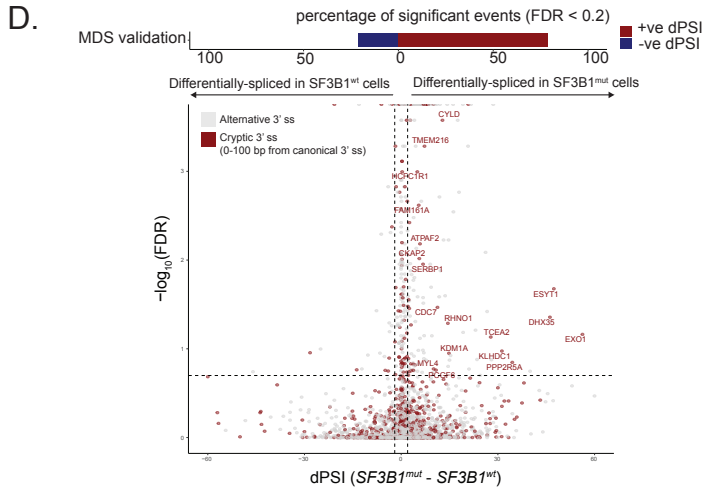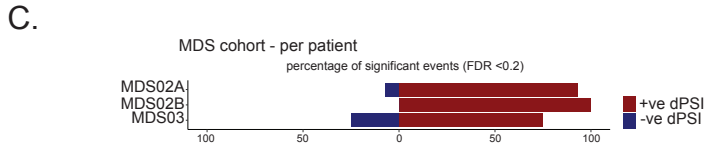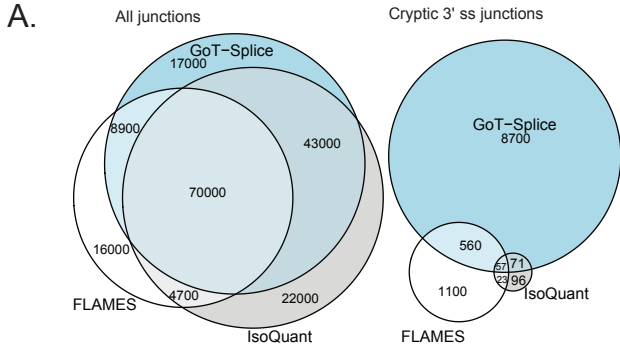
**Figure S4. Comparison of splicing junctions detected with GoT-Splice vs. other tools, MDS samples, and cancer types, related to Figure 4.**

**(A)** *Left*: Overlapping of all splicing junctions detected in GoT-Splice, FLAMES, and IsoQuant. *Right*: Comparison of splicing junctions supporting cryptic 3' ss between the two full-length isoform detection tools and GoT-Splice. **(B)** Estimated dPSI values from long-read tools FLAMES and IsoQuant compared to reported dPSI values measured in bulk data, like **Figure 2g**. **(C)** Bars showing the percentage of genes differentially spliced in $SF3B1^{mut}$ and $SF3B1^{wt}$ cells with BH-FDR adjusted *P*-value < 0.2 in each sample of the MDS cohort (MDS02(A/B)-03). **(D)** Differential splicing analysis between $SF3B1^{mut}$ and $SF3B1^{wt}$ cells across the aggregate of the MDS validation cohort (MDS04-06). Junctions with an absolute delta percent spliced-in (dPSI) > 2 and BH-FDR adjusted *P*-value < 0.2 were defined as differentially spliced. Bars (top) showing the percentage of genes differentially spliced in $SF3B1^{mut}$ and $SF3B1^{wt}$ cells of MDS validation cohort. **(E)** Bars showing the percentage of genes differentially spliced in $SF3B1^{mut}$ and $SF3B1^{wt}$ cells with BH-FDR adjusted *P*-value < 0.2 in each sample of the MDS validation cohort (MDS04-06). **(F)** Venn Diagram for the overlap of differentially spliced genes used more highly in $SF3B1^{mut}$ cells (*P*-values < 0.05, dPSI >0) from the bulk comparison of $SF3B1^{mut}$ vs. $SF3B1^{wt}$ cells in the MDS and MDS validation cohorts [i]. Increasing the read coverage threshold for the differentially spliced genes showed a more significant overlap between cohorts [ii]. *P*-values for the overlap from Fisher's exact test. **(G)** Differential isoform proportion in $SF3B1^{mut}$ vs $SF3B1^{wt}$ cells. Isoforms containing alternative/cryptic 3' ss overlapping with those detected in GoT-Splice are highlighted in red. *Top*: Quantification of significant (FDR adjusted *P*-value < 0.05) full-length isoforms detected containing alternative/cryptic 3'ss, split by direction of the logFC in proportion between MUT and WT cells. **(H)** Example of cryptic 3' ss in the gene *ERGIC3* missed by IsoQuant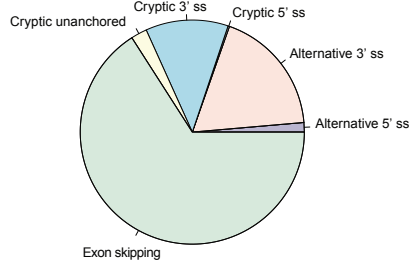 and FLAMES but detected and quantified by GoT-Splice. **(I)** Venn Diagram showing the overlap of genes with cryptic 3' ss splicing events (dPSI > 0 and FDR adjusted *P*-value < 0.2) found in the MDS discovery cohort vs. those found in the bulk RNA sequencing of a set of published $SF3B1^{mut}$ CLL samples. **(J)** Box plot of the average expression levels found across all cells in the MDS discovery cohort for the genes cryptically mis-spliced in the MDS dataset only, in both the MDS and CLL datasets, and the CLL dataset only. *P*-values were computed using Wilcox rank sum test.
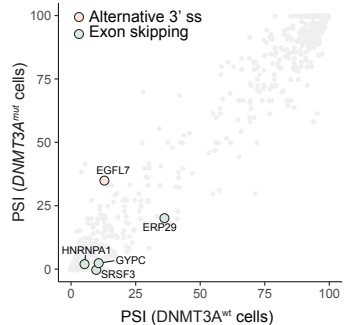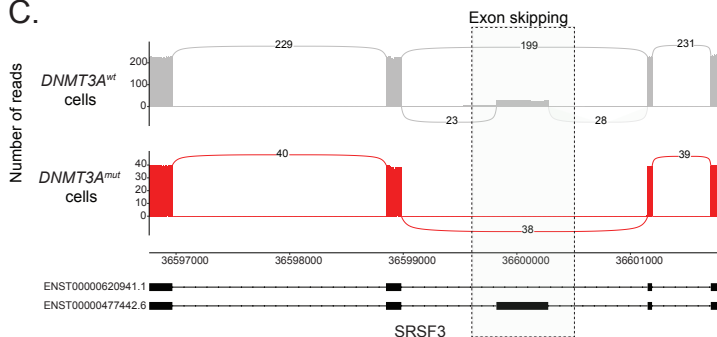
# Figure S5

## A.
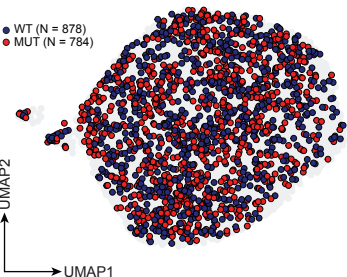Distribution of alternative splicing events in *DNMT3A*^*mut* CH

Cryptic unanchored · Cryptic 3' ss · Cryptic 5' ss · Alternative 3' ss · Alternative 5' ss · Exon skipping

## B.
○ Alternative 3' ss
○ Exon skipping

PSI (*DNMT3A*^*mut* cells) vs PSI (DNMT3A^wt cells)

EGFL7, ERP29, HNRNPA1, GYPC, SRSF3

## C.
Number of reads

*DNMT3A*^*wt* cells: 229, 199, 231, 23, 28 — Exon skipping

*DNMT3A*^*mut* cells: 40, 38, 39

36597000 · 36598000 · 36599000 · 36600000 · 36601000

ENST00000620941.1
ENST00000477442.6

SRSF3

## D.
*U2AF1*^*S34F* AML human sample
CD34+ cells

● WT (N = 878)
● MUT (N = 784)

UMAP2 / UMAP1

## E.
Distribution of alternative splicing events in *U2AF1*^*mut* cells

Cryptic unanchored · Cryptic 3' ss · Cryptic 5' ss · Alternative 3' ss · Alternative 5' ss · Exon skipping

## F.
PSI (*U2AF1*^*mut* cells) vs PSI (*U2AF1*^*wt* cells)

KIN, DAP3

○ Alternative 3' ss
○ Cryptic 3' ss
○ Exon skipping
○ Not significant

## G.
Number of reads

*U2AF1*^*wt* cells: 61, 12, 49, 45 — Exon skipping
*U2AF1*^*mut* cells: 79, 32, 47, 42

7782000 · 7784000 · 7786000

ENST00000379562.9

KIN

*U2AF1*^*wt* cells: 47, 46, 2 — Exon skipping
*U2AF1*^*mut* cells: 27, 28, 8

155728000 · 155728500 · 155729000

ENST00000535183.5
ENST00000471642.6
ENST00000411487.5
ENST00000368336.10
ENST00000343043.7

DAP3

## H.
Z-score dPSI (*SF3B1*^*mut* - *SF3B1*^*wt*)
-2 -1 0 1 2

*red - shared gene with cryptic events in the MDS discovery cohort

EP, MEP, MkP, HSPC, IMP, NP

## I.
Overlap of significant cryptic 3' ss in HSPCs, IMPs, MEPs, EPs between MDS and MDS validation

MDS [MEP+EP] VS MDS validation [HSPC+IMP]
P = 0.46
1.6 %
98.4 %

MDS [MEP+EP] VS MDS validation [MEP+EP]
P = 0.00029
53.2 % · 46.8 %

Key
■ % unvalidated
■ % validated

Including:
- CIA01
- DLST
- FOXRED1
- MED6
- STAU1
- SRSF11

## J.
*UROD*

Number of reads

Short-read: 4000, 3000, 2000, 1000, 0
Long-read: 4000, 2000 — Cryptic 3' ss

45013500 · 45014000 · 45014500 · 45015000

ENST00000652287.1
ENST00000339165.1
ENST00000339598.1
ENST00000371773.1
ENST00000246337.9

*PPOX*

Short-read: 200, 100, 0
Long-read: 400, 200 — Cryptic 3' ss

161168500 · 161169000 · 161169500 · 161170000 · 161170500
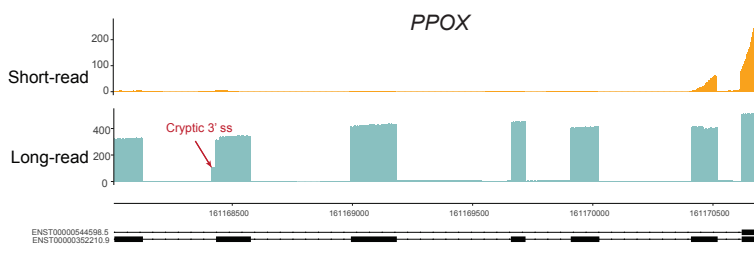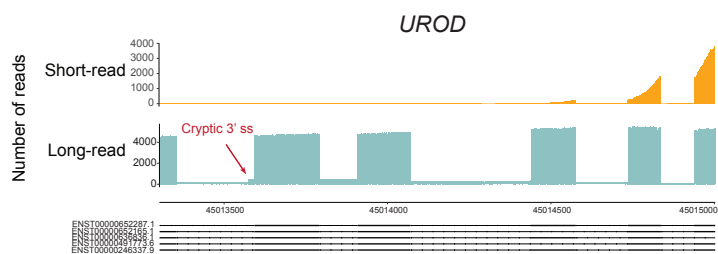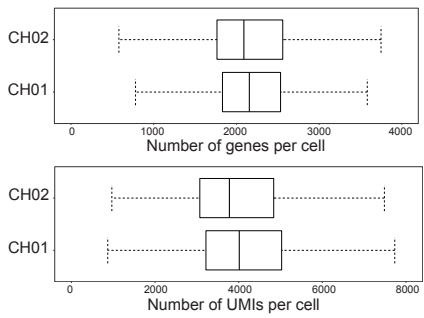
ENST00000544598.5
ENST00000352210.9

**Figure S5. Application of GoT-Splice to other mutations and identification of cell-type specific splicing changes, related to Figure 5.**

**(A)** Pie chart summarizing the distribution of different alternative splicing events detected after junction annotation in the *DNMT3A^mut^* CH sample. **(B)** Comparison of percent spliced-in (dPSI) values of alternative splicing events identified in *DNMT3A^mut^* vs. *DNMT3A^wt^* CH cells. Significant alternative splicing events are highlighted. **(C)** Sashimi plot of *SRSF3* exon skipping event showing the expected increase in the PSI value in *DNMT3A^wt^* cells. **(D)** UMAP of progenitor cells from *U2AF1^S34F^* AML sample overlaid with genotyping data. WT, cells with genotype data without *U2AF1* mutation; MUT, cells with genotype data with *U2AF1* mutation. **(E)** Pie chart summarizing the distribution of different alternative splicing events detected after junction annotation in the *U2AF1^mut^* AML sample. **(F)** Comparison of percent spliced-in (dPSI) values of alternative splicing events identified in *U2AF1^mut^* vs. *U2AF1^wt^* cells. Significant alternative splicing events are highlighted. **(G)** Sashimi plot of *KIN* and *DAP3* exon skipping events showing the expected increase in the exon skipping in *U2AF1^mut^* cells. **(H)** Heatmap of dPSI values between *SF3B1^mut^* and *SF3B1^wt^* cells for cryptic 3' splicing events identified in the main progenitor subsets across MDS validation samples. Rows (z-score normalized) correspond to cryptic 3' junctions found to be differentially spliced in at least one cell-type, with *P*-value <= 0.05 and dPSI >= 2. Columns correspond to cell-type. Genes with an *SF3B1^mut^* associated cryptic 3' splice site found in the MDS cohort are highlighted in red. **(I)** Pie chart showing the percent overlap of cryptically 3' spliced genes unique to MEPs and EPs in the primary MDS cohort that are also cryptically 3' spliced and unique to earlier progenitor cells (HSPCs and IMPs) in the MDS validation cohort (left) as well as the percent overlap with genes cryptically 3' spliced and unique to the MEPs and EPs in the MDS validation cohort (right). *P*-value for the overlap from Fisher's exact test. **(J)** Plots highlighting the differences in sequencing coverage along the transcript, between short- and long-read sequencing, for *PPOX* and *UROD*, demonstrating the ability for long-read sequencing to uncover splicing events than remain undetected in the 3' biased short-read sequencing.
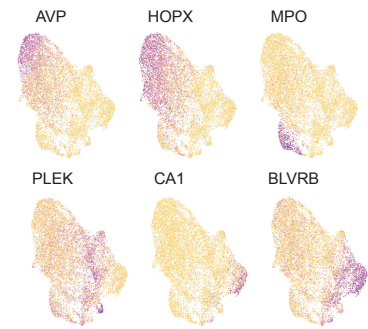
# Figure S6

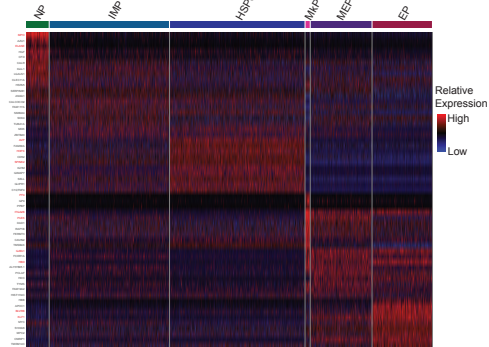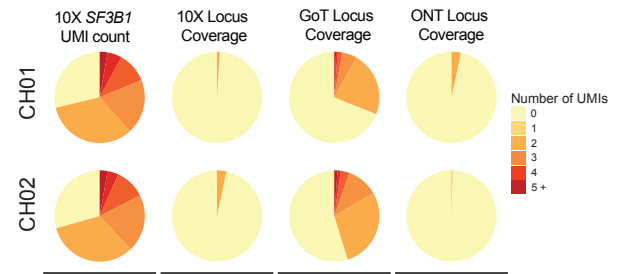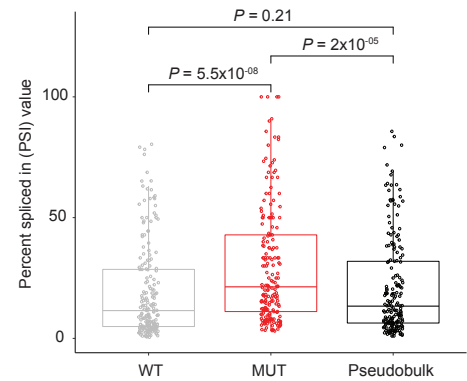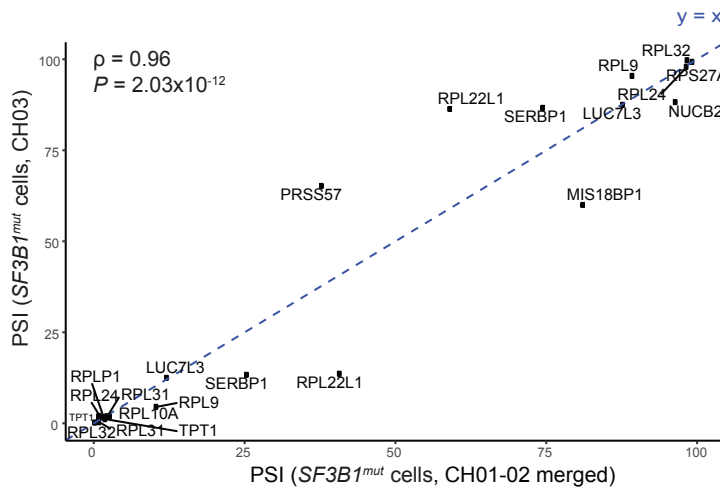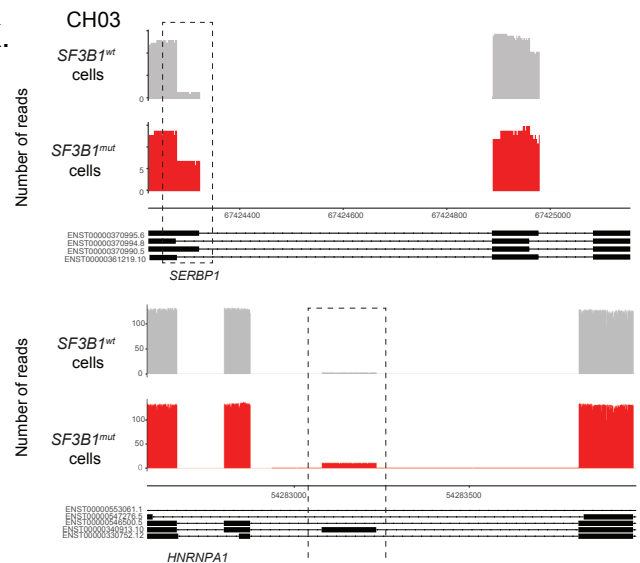**Figure S6. CH cohort QC, integration, and GoT-Splice analyses, related to Figure 6.**

**(A)** Number of genes per cell (top) and number of UMIs per cell (bottom) in CD34+ sorted hematopoietic progenitors from samples CH01-02 after QC filters, shown by each patient sample. **(B)** UMAP of CD34+ sorted progenitor cells for each individual sample of CH01-02 after integration using the Seurat package. **(C)** Expression of lineage-specific genes from Velten *et al.*[S1] scored and projected onto the UMAP representation of cells from CH01-02. **(D)** Heatmap of top 10 differentially expressed genes for each progenitor subset for CH01-02. **(E)** Fraction of cells in CH01-02 by number of *SF3B1* UMIs in standard 10x Genomics data without genotyping information, *SF3B1* UMIs with K666N (CH01) or K700E (CH02) locus coverage in standard 10x data, *SF3B1* UMIs with K666N (CH01) or K700E (CH02) locus coverage in GoT amplicon library, and *SF3B1* UMIs with K666N (CH01) or K700E (CH02) locus coverage in ONT library**. (F)** Per sample heatmap of relative expression of genes ordered by chromosome/chromosomal position following copy number variation analysis using the InferCNV package (see STAR Methods). Cells (y-axis) are stratified by *SF3B1* genotype status. **(G)** Pseudotime in $SF3B1^{mut}$ vs. $SF3B1^{wt}$ cells per CH sample. *P*-value for comparison of means from Wilcoxon rank sum test. **(H)** Normalized ratio of $SF3B1^{mut}$ cells in HSPC and EP cells for CH01 and CH02. Bars show the mean of $n = 100$ downsampling iterations to 1 genotyping UMI per cell. **(I)** Comparison of the PSI values of identified cryptic junctions in WT cells only (gray) vs. MUT cells only (red) vs. all cells in pseudobulk in CH01-02. *P*-values for comparison of means from Wilcoxon rank sum test. **(J)** Comparison of the $SF3B1^{mut}$ cells PSI values for alternative splicing events that were found to be significantly differentially used (*P*-value < 0.05) between $SF3B1^{mut}$ and $SF3B1^{wt}$ in the primary CH cohort and that were also present in the CH validation sample. Spearman's rank correlation coefficient and *P*-value derived from Student's t-distribution are shown. **(K)** Sashimi plot examples of alternative splicing events in *SERBP1* (top, cryptic 3' ss event) and *HNRNPA1* (bottom, exon skipping/inclusion event) found to be significantly differentially spliced between $SF3B1^{mut}$ and $SF3B1^{wt}$ cells in the CH validation sample.

**REFERENCES**

[S1] Velten, L., Haas, S.F., Raffel, S., Blaszkiewicz, S., Islam, S., Hennig, B.P., Hirche, C., Lutz, C., Buss, E.C., Nowak, D., et al. (2017). Human haematopoietic stem cell lineage commitment is a continuous process. Nat Cell Biol *19*, 271–281. 10.1038/ncb3493.