




# A single-cell strategy for the identification of intronic variants related to mis-splicing in pancreatic cancer

Emre Taylan Duman <sup>1,†</sup>, Maren Sitte <sup>1,†</sup>, Karly Conrads <sup>2,3,4</sup>, Adi Mackay <sup>3,5</sup>, Fabian Ludewig <sup>1</sup>, Philipp Ströbel <sup>3,5</sup>, Volker Ellenrieder <sup>2,3,6</sup>, Elisabeth Hessmann <sup>2,3,6</sup>, Argyris Papantonis <sup>3,5,6,\*</sup> and Gabriela Salinas <sup>1,3,\*</sup>

<sup>1</sup>NGS-Core Unit for Integrative Genomics, Institute of Pathology, University Medical Center, Göttingen, Germany

<sup>2</sup>Clinic of Gastroenterology, Gastrointestinal Oncology and Endocrinology, University Medical Center, Göttingen, Germany

<sup>3</sup>Clinical Research Unit 5002 (CRU5002), University Medical Center, Göttingen, Germany

<sup>4</sup>Institute of Medical Bioinformatics, University Medical Center, Göttingen, Germany

<sup>5</sup>Institute of Pathology, University Medical Center, Göttingen, Germany

<sup>6</sup>Comprehensive Cancer Center Lower Saxony (CCC-N), Göttingen, Germany

\*To whom correspondence should be addressed. Tel: +49 551 39 65734; Email: argyris.papantonis@med.uni-goettingen.de

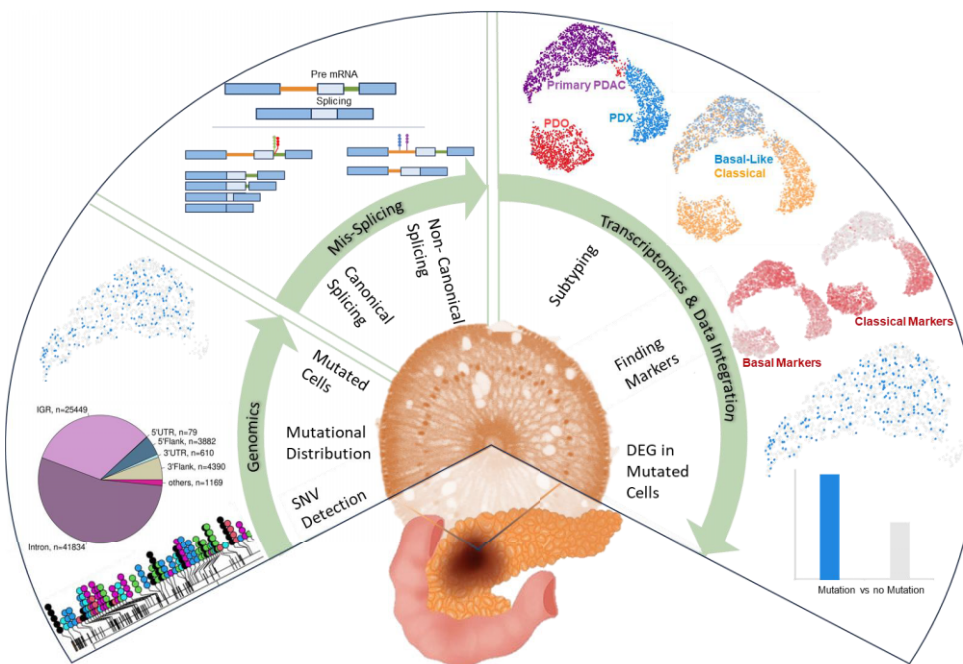
Correspondence may also be addressed to Gabriela Salinas. Tel: +49 551 39 60778; Email: gabriela.salinas@med.uni-goettingen.de

†The first two authors should be regarded as Joint First Authors.

## Abstract

Most clinical diagnostic and genomic research setups focus almost exclusively on coding regions and essential splice sites, thereby overlooking other non-coding variants. As a result, intronic variants that can promote mis-splicing events across a range of diseases, including cancer, are yet to be systematically investigated. Such investigations would require both genomic and transcriptomic data, but there currently exist very few datasets that satisfy these requirements. We address this by developing a single-nucleus full-length RNA-sequencing approach that allows for the detection of potentially pathogenic intronic variants. We exemplify the potency of our approach by applying pancreatic cancer tumor and tumor-derived specimens and linking intronic variants to splicing dysregulation. We specifically find that prominent intron retention and pseudo-exon activation events are shared by the tumors and affect genes encoding key transcriptional regulators. Our work paves the way for the assessment and exploitation of intronic mutations as powerful prognostic markers and potential therapeutic targets in cancer.

## Graphical abstract



Received: February 1, 2024. Revised: April 24, 2024. Editorial Decision: May 6, 2024. Accepted: May 23, 2024

© The Author(s) 2024. Published by Oxford University Press on behalf of NAR Genomics and Bioinformatics.

This is an Open Access article distributed under the terms of the Creative Commons Attribution-NonCommercial License

(<https://creativecommons.org/licenses/by-nc/4.0/>), which permits non-commercial re-use, distribution, and reproduction in any medium, provided the original work is properly cited. For commercial re-use, please contact [journals.permissions@oup.com](mailto:journals.permissions@oup.com)

## Introduction

Introns in human genes can be extraordinarily large (up to 1 Mbp, with ~3400 being longer than 50 kb and ~1200 longer than 100 kb), and account for half of the non-coding human genome. As mutations in introns do not directly affect protein-coding sequences, they are usually overlooked (1–3). As a result, little attention is paid to the importance of intronic-located splicing regulatory elements that control the fidelity of pre-mRNA splicing and transcription timing. This is surprising given that pathogenic variants cause abnormal splicing changes, typically by damaging existing splicing motifs or creating novel splicing motifs and may comprise 15–60% of all human disease variants (1–3).

Recent studies have reported significant numbers of intronic variants and deletions in protein-coding genes that are associated with under- or overexpression of the affected genes or of distant genes interacting in 3D, thus influencing their regulation in normal or pathogenic conditions (4–6). It is also well established that introns contribute to the control of gene expression by including regulatory regions and non-coding (yet functional) genes or even directly by their extensive length (6). A substantial number of pathogenic variants located deep within introns (i.e. >100 bp from an exon-intron boundary) were recently reported, suggesting that sequence analysis of full introns may help to identify causal mutations for many undiagnosed clinical cases (4–9). Moreover, direct associations between intronic mutations and certain diseases have also been reported, albeit sporadically (7–9). These results agree with the findings we present here. For instance, we identified genes that are simultaneously overexpressed in basal-like cells from pancreatic cancer (PDAC) tumors and rank as the most mutated transcripts, particularly when we consider intronic variants. In this regard, PDAC-associated mutations were reported to synergize in tumorigenesis by globally altering the splicing program of cell (10). Moreover, splicing factors were recently shown to either promote the early events in pancreatic tumorigenesis and resistance to chemotherapy or to limit the metastatic potential of PDAC cells (11,12). A very recent publication has identified a splicing signature specific to basal-like cells, distinguishing PDAC subtypes as accurate survival predictors when considered in the overall population of PDAC patients, as well as within homogeneous subtype cohorts, indicating their efficacy as biomarkers (13).

However, there exist many limitations when investigating intronic mis-splicing variants, the main ones being the lack of approaches capable of simultaneously interrogating genomic and transcriptional information, and the lack of guidelines designed for assessing intronic variants and their contribution to abnormal splicing changes in disease. We address these limitations by introducing a novel pipeline that utilizes full-length single-nuclei and bulk RNA sequencing strategy for the ‘deep’ characterization of genetic variability within introns and of their effects on splicing and gene expression in PDAC. Indeed, cancer is one of those diseases, where alternative splicing is the basis for the identification of novel diagnostics, and therapeutic strategies for therapy (e.g. antisense oligonucleotides or small-molecule modulators of spliceosome (9,14,15)). This new approach therefore holds promise for both the elucidation of fundamental biological principles connected to splicing regulation, and the identification of therapeutic targets in human disease.

## Materials and methods

### Patients’ sample information

Utilization and characterization of human PDAC data and samples within the CRU5002 has been approved by the ethical review board of the UMG (11/5/17). Informed consent was obtained from all subjects involved in the study. Sequencing studies and the generation of organoids and PDX models have been performed using tumor tissue from CRU5002 PDAC patients with progressed disease upon histological PDAC confirmation.

### Generation of PDX models

For the generation of PDX models, bulk tumor tissue was subcutaneously transplanted into SHO-*prkdc*<sup>scid</sup>*Hr*<sup>br</sup> mice. Engrafted subcutaneous tumors were passaged in mice for three generations prior to snRNA-sequencing.

### Generation of organoids models

Tumor tissue was minced and digested in Dulbecco’s modified Eagle’s medium containing 5 mg/ml of collagenase XI (Sigma-Aldrich, C9407), DNase final concentration 10 µg/ml (SIGMA D5025-150KU) and Y-27632 final concentration 10.5 µM (Adooq Bioscience, 129830-38-2) and incubated at 37°C for 45 min. The material was further embedded in Matrigel (Corning, New York, USA; Cat#356231) and cultured in human pancreatic cancer complete medium (Wnt3a, R-spondin1, Gastrin, hEGF, A 83-01, hFGF-10, mNoggin, Primocin, N-acetylcystein, Nicotinamide, B27 supplement and Y-27632). For passage, the Matrigel-containing organoid was digested by TrypLE™ Express (Thermo Fisher, 12605-028) with DNase and Y-27632 as described above for 15 min. The sample was centrifuged at 500 × g for 5 min, and the precipitated cells were embedded in GFR Matrigel and cultured in human pancreas organoid complete feeding medium.

### Nuclei extraction from primary PDAC and PDAC models

Nuclei isolation of PDAC was performed according to the ‘Nuclei Isolation from Cell Suspensions & Tissues for Single Cell RNA sequencing’, Document Number CG000124 Rev E, 10× Genomics, (30 June 2021). For organoids, minor modifications were performed. The cells were spun down at 500 × g, lysis took place for 15 minutes and only two washing steps were performed. Following the second wash, the nuclei pellet was resuspended in 1 ml (about 0.03 oz) wash buffer in preparation for the ICELL8 protocol.

### Full-length single-cell RNA-seq using ICELL8

The Takara ICELL8 5184 nano-well chip was used with the full-length SMART-Seq ICELL8 Reagent Kit. Nuclei suspensions were fluorescent-labelled with Hoechst 33342 for 15 min prior to their dispensing into the Takara ICELL8 5184 nano-well chips. Cell Select Software (Takara Bio) was used to visualize and select wells containing single nuclei. Nine 5184 nano-wells chips were used for all samples and 11 084 nuclei were processed for data analysis. Specifically, after quality control, 3416 nuclei were used for primary PDAC, 3037 for PDX and 4631 for organoids respectively. cDNA synthesis and library preparation were done according to description in previous study (16). Libraries were sequenced on the HiSeq

4000 (Illumina) to obtain on average  $\sim 0.3$  Mio reads per nuclei (SE; 50 bp).

### Bulk-RNA-Seq from primary PDAC samples

Total RNA was extracted from FFPE tumor patient samples using the ReliaPrep™ FFPE Total RNA Miniprep System (Promega). RNA Integrity was determined using the Fragment Analyzer. Because of low RNA integrity (sizing from 50 to 140 bp), we performed a modified TruSeq Stranded Total RNA Library Prep Human/Mouse/Rat (Cat. No. 20020596) starting with 200 ng of total RNA. The modifications include (a) ignoring fragmentation step, (b) ligation optimization by adjusting adapters concentration during library preparation, (c) increasing PCR cycles (15 in total) and eliminating primer dimers prior to sequencing (Agencourt AMPure XP magnetic beads, Beckman Coulter). Primary PDAC were sequenced on the NovaSeq6000; S4 flow cell PE 300 cycles generating a data set of 50–400 Mio reads per sample.

### Bulk-RNA-Seq from organoids and PDXs

RNA libraries were prepared starting with 300 ng of total RNA using a non-stranded mRNA Seq (TruSeq RNA Library Preparation Cat. RS-122-2001) from Illumina according to the manufacturer's recommendations. Libraries were sequenced on the Illumina HiSeq 4000 (SE;  $1 \times 50$  bp; 30–35 Mio reads/sample).

### Whole genome sequencing

WGS data from two primary PDAC samples TM56 and TM27 were sequenced at  $\sim 40\times$  coverage on the Illumina NovaSeq 6000 sequencer following the protocol provided by the supplier. Libraries were performed using the PCR Free DNA library preparation from Illumina Cat. No.: 20041794). Alignment, variant calling, and benchmarking were performed using Illumina DRAGEN Germline pipeline 4.2.4.

### Pre-processing of single-nuclei RNA-seq data

Raw sequencing files were processed as described in (16). Briefly, Cogent NGS analysis pipeline (CogentAP) from Takara Bio (v1.0) was applied for de-multiplexing and creating the gene expression matrices from each FASTQ file. Reads were aligned against the human genome GRCh38 v107 ([https://www.ensembl.org/Homo\\_sapiens/Info/Index](https://www.ensembl.org/Homo_sapiens/Info/Index)). Quality control of the data has been done by using CogentDS QC as outlined in (16). QC considered the amount of usable unique reads per nucleus (over 65%), the number of reads generated per nucleus (over 250K), and the median mitochondrial content of the PDAC tumors of 18% (with 26% IQR-3). Intron regions of genes were also included in further evaluation. Generated gene matrices were used as input for the SingleCellExperiment R package (v3.0) (17) to generate SingleCellExperiment objects for the subsequent downstream analysis.

### Identification of disease subtypes and cell types

To identify the tumor subtypes and cell types of which the single cells are composed, two different methods were employed. First, we used the marker-based method AUCell (18) to identify the PDAC subtypes classical and basal-like tumor. This analysis was based on established marker genes (19). AUCell ranks the genes in each cell by decreasing expression value, and marks cells according to their most expressed marker

genes. Secondly, we performed cell type annotation for more refined subpopulations to address the heterogeneity of the tumor and its matched models. We utilized the reference data set as provided by the SeuratData R package (panc8.SeuratData) (20) and the prediction function as implemented in the R package SingleR (21). Downstream analyses performing the UMAP algorithm were done as implemented in CogentDS (v1.0) for dimensionality reduction and data visualization. To determine which genes were differentially expressed between tumor subtypes and cell types in a particular patient, Wilcoxon Rank Sum and Signed Rank Test was used, together with *P* values adjustment with Benjamini–Hochberg method.

### Pseudo-variant calling in snRNA and bulk RNA-seq data

Bam files resulting from CogentDS were used as input for the pseudo bulk variant calling using GATK best practices pipeline for RNAseq variant calling (22). Consistent with the recommendations of GATK, duplicates were removed with Picard MarkDuplicates, and read groups were added with Picard AddOrReplaceReadGroups. Subsequently, Cigar reads were split into exon segment and hard-clip any sequences overhanging into the intronic regions with SplitNCCigarReads. Variant calling was performed by HaplotypeCaller and all variants were then hard filtered by the following criteria:  $FS > 30$  and  $QD < 2$ . Patient specific vcf files were intersected or merged using bcftools-isec to collect shared or unique mutations from the samples. Resulting variants were converted into maf (mutation annotation file) using Ensembl VEP (23) annotation tool for the identification of the intronic, intergenic, and splice junction mutations (Figure 3A). Finally, mutations were grouped into three categories: SS, Proximal, and Deep by using Ensembl GTF file version 107. Total number of SS and Prox mutations were subtracted from total number of mutations in maf files to calculate the number of deep intronic variants. Negative values indicate that the deep intronic mutation affects multiple transcripts of the gene. Pre-Ranked enrichment analysis was performed by using mutation lists from three locational groups. Mutation numbers from SS and Proximal regions were normalized by using total intron number of each gene. Deep mutation numbers were normalized by per unit length (kb). Normalized values are ranked, and overall survival was calculated using Kaplan–Meier analysis based on TCGA-PAAD cohort. The result is shown as Kaplan–Meier plot with *P* value from log rank test generated by the cBioPortal (24).

### Integration of the SpliceAI and Pangolin scoring

To perform scoring on the discovered variant list, hard filtered and intersected VCF files were used. Scores of each variant aggregated on their genes and maximum scores from both algorithms have been taken. Genes that scored with higher than 0.1 one of the tools have been selected as possibly mis-splicing related variants as it stated in their reference manual. SpliceAI v1.3 and Pangolin have been used in our analysis with default parameters with GRCh38 reference genome.

### PubMed annotation of identified genes from different approaches

The role of specific genes as tumor suppressors and their relevance to pancreatic cancer have been analyzed using PubMed data through Entrez Python API (BioPython v-1.83) (25). Au-

tomated PubMed queries have been done, identifying articles that associate each gene with ‘tumor suppressor’ and ‘pancreatic cancer’ in their titles or abstracts related to the gene. The outcomes are then aggregated into a DataFrame and visualized in a heatmap to elucidate patterns of gene involvement in tumor suppression and pancreatic cancer.

### Cell specific variant calling

To perform single cell specific variant calling, bam files were separated into chip- and sample-specific bam files. Separated bam files were used for the variant calling pipeline, which was applied to individual bam file as before, without the base quality recalibration step (BQSR) to be able to get more variants possible per cell. Quality filtering was applied with the same quality filters as mentioned in pseudo-bulk RNA variant calling in the methods section. VCF files for each barcode were used to collect barcode-specific mutations and visualize as a mutation distribution plot using CogentDS UMAP functions.

### Enrichment analysis of non-canonical (Deep) and canonical (SS-Prox) intronic variants

Genes with high Intronic mutation numbers (Table 1) were forwarded to over representation analysis (ORA) to identify enriched pathways and GO terms using the gProfiler online tool (26). Canonical and non-canonical intronic mutations were separately analysed setting the *P*-value threshold to 0.05 under multiple correction method g:SCS.

### Identification of intron retention events from RNA-Seq data

Intron retention analysis was performed using IRFinder 1.3.0 (27). IRFinder calculates IR-ratios to measure IR level reflecting the proportion of intron retaining transcripts. To compare PDAC samples to normal tissue, we downloaded fastq files from nine PDAC and nine healthy human pancreas samples from NCBI GEO database (accession ID: GSE211398). These were pre-processed using the same methods as described in Bulk-RNA-Seq from primary PDAC. Tumor samples from GSE211398 dataset were then analyzed using the IRFinder pipeline to identify IR events. Intronic depth for all samples normalized by CPM due to the library size variations between the public data and CRU samples. Due to the different number of samples between the two dataset, weighted means of IR scores were calculated to balance the variable sample sizes by using 0.1 threshold among 80% of all samples. Normalized intron depths were also used as further filtering to reduce the number of identified possible false negative discoveries.

## Results

### Strategy for the identification of intronic variants affecting splicing in PDAC

One of the most critical post-transcriptional mechanisms reprogramming transcriptional output and proteomic diversity in cancer cells is the loss of splicing precision when removing introns from pre-mRNAs (28). Consequently, many mis-spliced variants are instead targeted for nuclear degradation or for nonsense-mediated mRNA decay (NMD) and thus, only few annotated alternative isoforms correspond to the precursors of the proteins mapped by large-scale proteomics studies (29).

To investigate this mechanism, we applied a new pipeline using a full-length single nuclei RNA-seq (snRNA-seq) approach to three primary PDAC tumors (TM16, TM27 and TM56) and matched tumor-derived preclinical PDAC models (i.e. organoids and subcutaneous patient-derived Xenograft, PDXs). We then sought to determine pathogenic intronic variants causing abnormal splicing in these patient samples (Figure 1A). All samples were part of a unique patient cohort recruited for the Clinical Research Unit 5002 (<https://gcc.umd.edu/en/cru-5002/>).

For our snRNA-Seq experiments, we used the ICELL8 platform previously established in our group connecting genotype to phenotype in individual cells (16). In contrast to the chemistry used by droplet-based platforms (i.e. in 3'-end approaches), ICELL8 is based on the SMART full-length chemistry allowing for the full read coverage of transcripts. Notably, when single nuclei are processed with this platform and chemistry, a strong enrichment in pre-mRNA is observed, including comprehensive coverage of introns and exons along these pre-mRNAs (Figure 1A and Supplemental Figure S1). Furthermore, utilization of full-length chemistry allowed strong detection of intergenic and intronic sequences, as well as of non-coding RNAs, especially long intergenic non-coding RNAs (lincRNAs) (16,30) as demonstrated in Figure 2E.

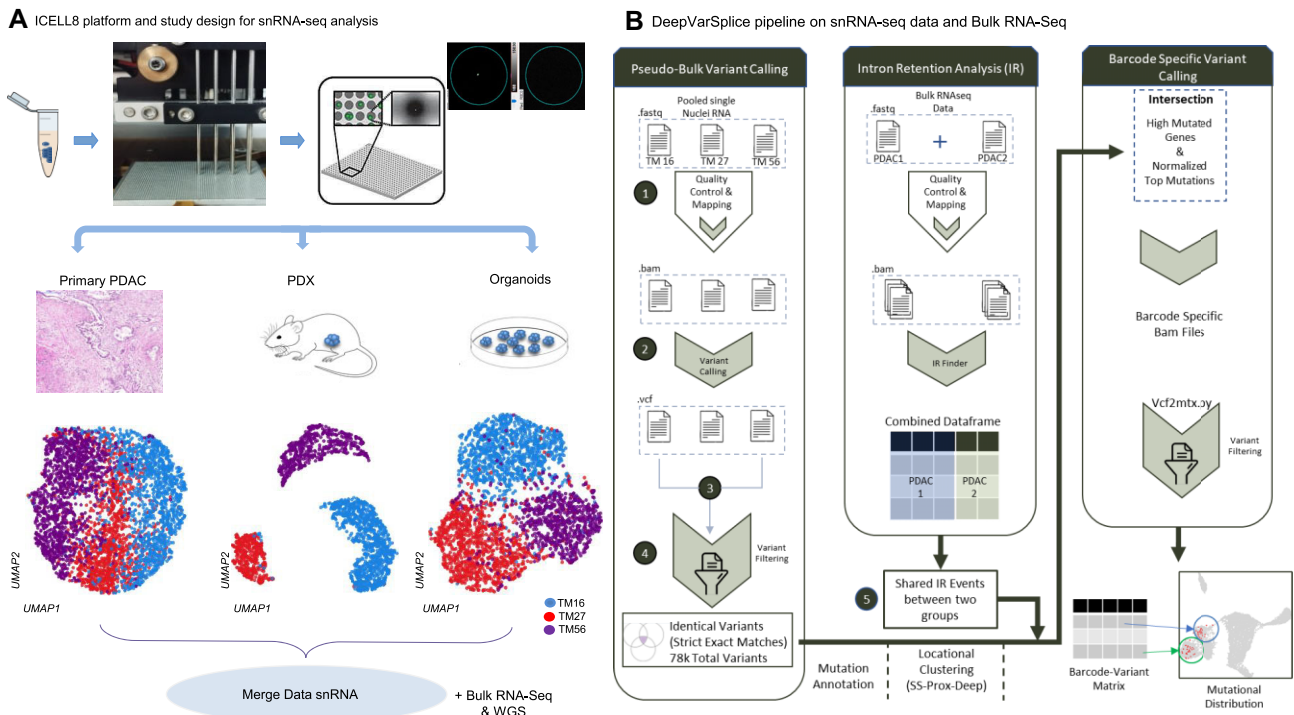
For the identification of intronic mis-splicing variants, we developed a variant-calling pipeline called ‘DeepVarSplice’ that combines (i) snRNA-seq data determining the variant’s position in the genome with (ii) bulk RNA-seq data mainly capturing mRNAs and thus, charting splicing events in our samples (Figure 1B).

The pipeline begins with a pseudo-bulk snRNA-seq variant calling analysis using the per-sample GATK method (Figures 1B (1–4)), followed by an intersection of the individual findings to identify exact mutations present in all samples. The variants are then classified according to their genomic position considering intergenic, intronic and exonic regions using the maf file (Mutation Annotation File). By filtering the intronic variants, the genes were normalized and ranked according to intronic mutation load. The intronic variants are then forwarded to two parallel branches of the pipeline.

The first branch is set up for the investigation of these variants on the single cell level. Therefore, barcode-variant matrices containing the number of mutations for each gene (columns) and cell (rows) have been created to allow further transcriptional analysis related to the variants.

Simultaneously, the second branch of the pipeline connects genes that exhibit intronic mis-splicing to pathogenic splicing events through the investigation of partial or total intron retention (IR) or pseudo-exon activation (PEA) using the IRFinder algorithm (Figure 1B (5)). Finally, the intron mutated genes from the snRNA data were merged with the intron retained genes from bulk RNA-Seq data. As an outcome, we highlight genes showing mis-splicing related intronic variants that contribute to malignancy and thus proposed as potential therapeutic targets in pancreatic cancer.

To ensure the reliability of the variants detected by DeepVarSplice, we conducted comprehensive whole-genome sequencing (WGS) on two primary PDAC samples, namely TM27 and TM56. Our analyses encompassed gold-standard variant discovery, rigorous filtering, and benchmarking, utilizing both WGS and snRNA-Seq data, as depicted in Supplemental Figure S2. Out of the 153 000 shared variants detected by DeepVarSplice (snRNA-Seq data) in the two pri-



**Figure 1.** Strategy for the identification of intronic variants related to splicing in PDAC. **(A)** snRNA-seq using the ICELL8 platform and the SMART full-length chemistry allowing for full transcript coverage was applied to PDAC primary tumors (TM16, TM27 and TM56) and matched organoids and PDX lines. A total of nine 5184-nanowell chips were performed (three for primary PDAC, three for organoids and three for PDX respectively and visualized in the UMAP. Finally, data from all nine nanowell chips were merged. Bulk RNA sequencing was performed additionally for each of the samples. **(B)** DeepVarSplice variant calling pipeline combining snRNA-seq and bulk RNA-seq data determining intronic variant positions in the genome and intron retention events. 1) Gold standard RNA-seq variant calling pipeline pre-processing; 2) sample-level variant calling; 3) combination strategies for variant discovery (Shared-Combined); 4) Further filtering for RNA-editing variant removal; and 5) determination of shared intron retention events with weighted IR-ratio calculation between public and own PDAC samples.

primary PDAC samples, 127 000 (83%) were confirmed using stringent criteria for variant calling in the WGS dataset. Both methods exhibited higher precision, with a Ti/Tv ratio of 2.3 for the snRNA-Seq data and 2.01 for the WGS data, thereby affirming the trustworthiness of the DeepVarSplice method.

### Identification of PDAC intronic variants using DeepVarSplice on snRNA-seq data

UMAP visualization of single nuclei clustering from the primary PDAC tumors, organoids and PDXs (Figure 2A) was strongly influenced by the patient of origin of each sample. To evaluate the consistency in transcriptomics and intronic mis-splicing variants, we merged data from all PDAC models of each patient and regenerated separate UMAP plots for TM16, TM27 and TM56 (Figure 2B). Then, by applying the DeepVarSplice pipeline (Figure 1B) to this snRNA-seq data, we observed that > 95% of the mutations we identified were in non-coding sequences (i.e. in intronic or intergenic regions) and only ~5% were found in protein-coding sequences. The most mutated genes at the level of introns included *RBFOX1*, *CSMD1*, *WWOX*, *CNTNAP2* and *LRP1B* and were notably shared among PDAC tumors and models (Figure 2C).

To identify *bona fide* intronic variants, we determined the Ti/Tv ratio in each primary tumor-, organoid- or PDX-derived dataset. Interestingly, primary tumor data showed higher precision and a more realistic Ti/Tv ratio (2.18) in comparison to organoid- and PDX-derived data (Figure 2D). We therefore decided to focus on the three tumor-derived datasets for all

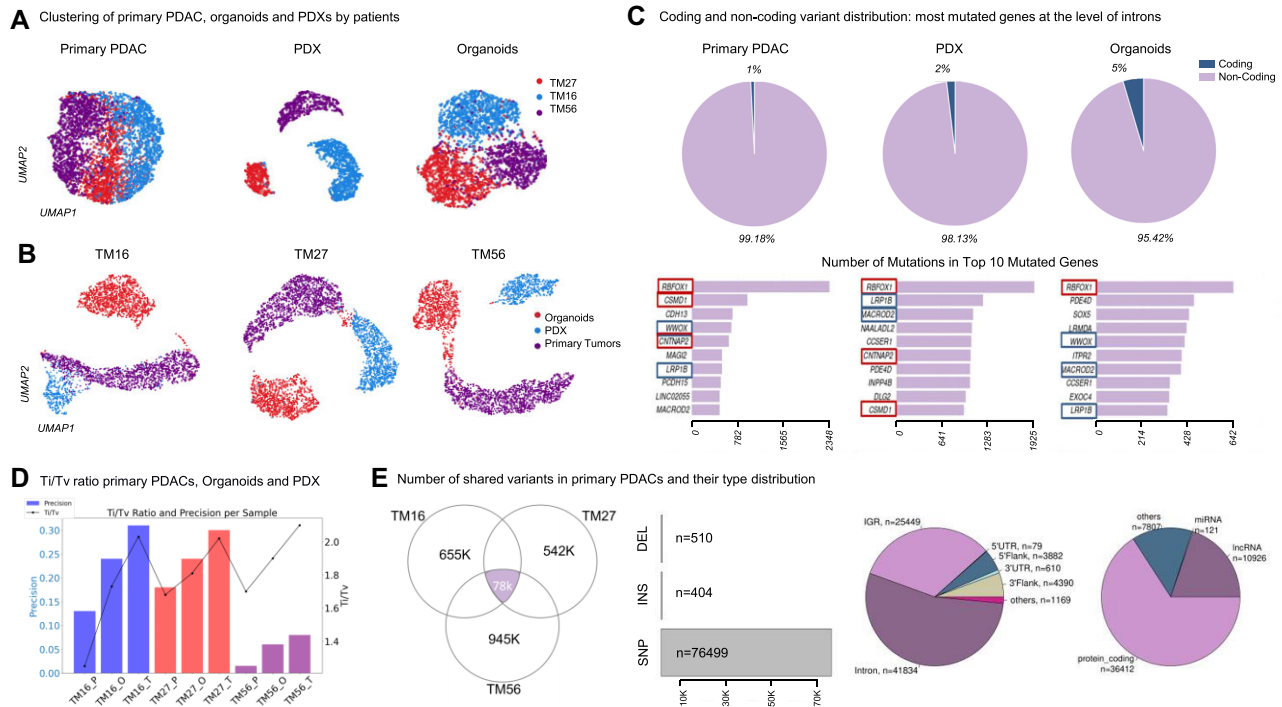
ensuing analyses, which we complemented with bulk RNA-Seq data generated for 24 primary PDAC tumors from the CRU5002 cohort, as well as with public normal and tumor pancreas tissues data (accession ID: GSE211398).

In the end, variants shared between all three-tumor snRNA-seq datasets amounted to ~78 000 non-coding variants with 41 834 in introns, 25 449 in intergenic regions (IGR), 4390 in 3' and 3882 5' gene flanks, 610 in 3' UTRs and 79 in 5' UTRs (Figure 2E). All shared variants detected in the primary PDAC with QUAL > 30 is listed in Supplemental Table S1. Notably, 90% of these variants (69 488) were reported in the dbSNP\_RS database based on WGS data remarking the validity of the variant calling performed with the snRNA Seq data. Moreover, the remaining 7926 intronic variants identified as novel ones underscore the potential of snRNA-Seq datasets to detect both validated intronic variants and novel ones, demonstrating the added value of our approach.

Of these, 98.82% ( $n = 76\,499$ ) qualified as SNPs, 0.66% ( $n = 510$ ) as deletions, and 0.52% ( $n = 404$ ) as insertions. Regarding the genomic annotation, most ( $n = 36\,412$ ) are mapped in protein-coding genes or lncRNAs ( $n = 10\,926$ ) and very few ( $n = 121$ ) in miRNAs (Figure 2E and Supplemental Table S1).

### Classification and validation of intronic mis-splicing variants based on location

Next, to classify intronic variants we first performed a ranking based on the number of variants detected per gene nor-



**Figure 2.** Identification of intronic variants using snRNA-seq data. **(A)** UMAP plot for three primary PDAC tumors and corresponding model's organoids and patient-derived Xenograft, PDXs (TM16, TM27 and TM56). **(B)** UMAP plot showing clustering by patients (primary PDAC and models) for TM16, TM27 and TM56. **(C)** Pie charts of the three PDACs (tumor, organoids and PDX) showing percentage of exonic and intronic mutations (above) and bar charts represent number of top intronic mutated genes (below). **(D)** Ti/Tv ratio of the variant analysis of primary PDAC and their corresponding model's organoids and patient-derived Xenograft, PDXs (TM16, TM27 and TM56). **(E)** Shared intronic variants and their classifications among all primary PDAC: TM16, TM27 and TM 56. Number of totals discovered and quality filtered intronic mutations from each sample (left), distribution of variant types among shared (78K) intronic mutations (middle), pie charts showing the distribution and location of the intronic variants in intronic, intergenic and coding regions, and a pie chart showing the intronic variants detected on coding and non-coding transcripts (right).

malized by the length and number of introns in each gene (top 20 visualized in [Supplemental Figure S3](#)). Then, to link intronic variants to mis-splicing, we stratified mutations located 1–2 and 3–20 nt away from the nearest exon-intron junction, which we classified as donor/acceptor sites (SS) and branchpoint-proximal regions (BPs), polypyrimidine tracts (PPTs). All other variants >20 nt are categorized as 'deep intronic'<sup>5</sup>. (see [Table 1](#) and [Supplemental Table S2](#)). Based on these indexes, many genes showed a high number of deep intronic variants, and some of these were reported as tumor suppressors expressed in specific tissues or only in tumor cells, e.g. *LRP1B*, *CSMD1*, *WWOX*, *FHIT*, *MTUS2*, *MAPK4* and *MAP3K14* (31–38) ([Supplemental Table S2](#)). The most prominent of these was *RBFOX1* and showed 208 deep intronic variants. The *RBFOX* family of RNA-binding proteins is well known to regulate alternative splicing (AS) (39,40). Recently, *RBFOX2* was reported to modulate a metastatic AS signature in pancreatic cancer (41).

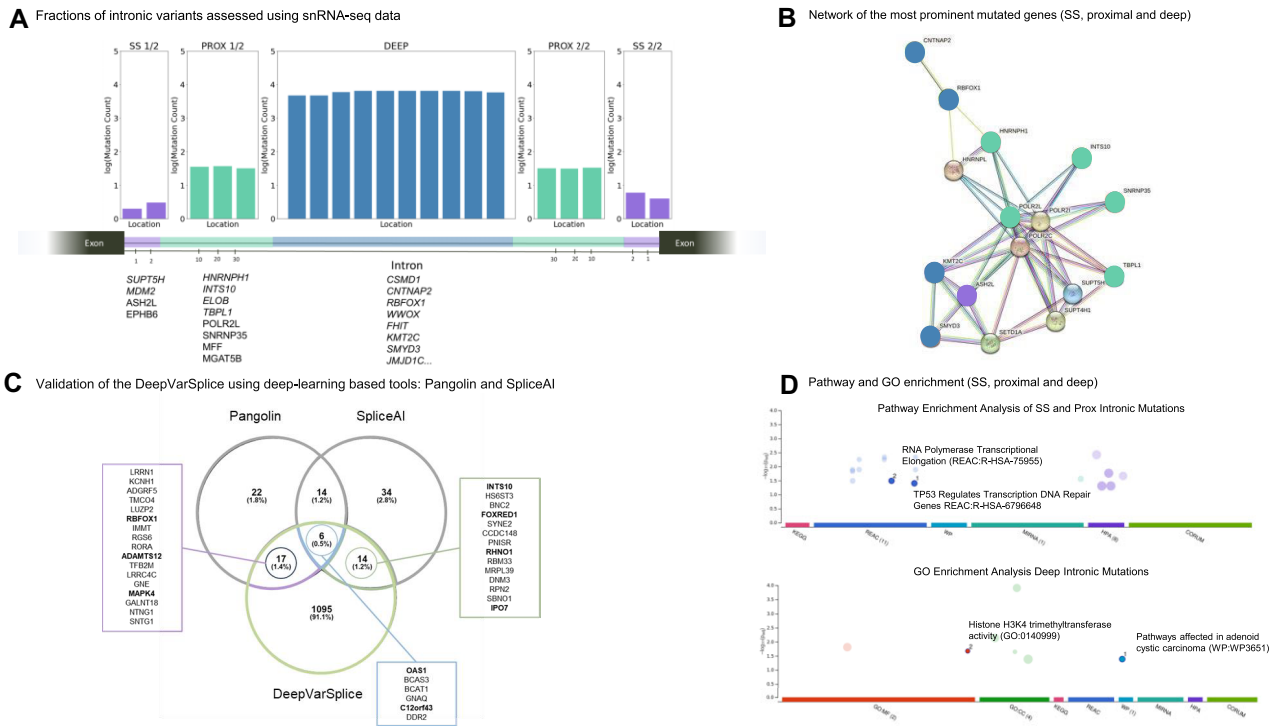
We next examined the fraction of intronic variants over all variants at each position (Figure 3A). 99.6% of them were located deep in introns. Variants located in the proximal region are 0.32% of the total, with just 0.02% near SSs. By conducting functional network analysis on the key genes identified as mutated in the splice site (SS), proximal, and deep intronic regions outlined in [Table 1](#), we discovered connections among genes implicated in mRNA capping (*HNRNPH1*, *INTS10*, *TBPL1*, *POLR2L*, *SNRNP35*) (42), as well as in the regulation of transcription via histone methy-

lation and H3K4-specific histone methyltransferase activity (*SMYD3*, *KMT2C*, *JMJD1C*, *ASH2L*), as illustrated in [Figure 3B](#). Pathway enrichment analysis carried out on genes mutated in the SS and proximal regions revealed enrichment in RNA polymerase II transcription elongation (REAC:R-HSA-75955) and the TP53 regulates transcription DNA repair genes pathway (REAC:R-HSA-6796648) as depicted in [Figure 3D](#). Additionally, gene ontology analysis conducted exclusively on genes mutated deep within introns demonstrated enrichment in histone H3K4 trimethyltransferase activity (GO:0140999) and pathways associated with adenoid cystic carcinoma (WP:WP3651), also shown in [Figure 3D](#).

Finally, visualization of the distribution of intronic variants in top mutated genes *WWOX*, *SMYD3*, *JMJD1C* and *NAV* ([Supplemental Figure S1](#)) exemplifies read coverage from snRNA-seq and bulk RNA-seq in primary PDAC tumors and how these can be superimposed to evaluate splicing events.

Notably, DeepVarSplice identified 1132 genes showing a substantial density of intronic variants that could potentially be linked to abnormal splicing. These genes were selected based on the modified Z-score method. This method is particularly apt for the normalized variant numbers dataset, which exhibits a non-normal distribution. The modified Z-score uses the median and the Median Absolute Deviation (MAD), providing robustness against skewed distributions ([Figure 3A](#)).

To validate our findings, we also employed recently developed deep learning based tools, specifically (i) Pangolin (43)



**Figure 3.** Selection and classification of genes exhibiting highest numbers of intronic variants. **(A)** Mutation percentages of all discovered mutations based on the variant locations on the left (SS, Prox, Deep)  $\log_2$  converted counts of each group represented in y-axis, base-pair distances for SS and PROX regions and locational percentages for intronic region in x-axis. **(B)** Protein network representation with nodes colored according to the location of identified mutations. **(C)** Performance of deep learning-based methods for mis-splicing detection: Validation of DeepVarSplice with SpliceAI and Pangolin. **(D)** Pathway Enrichment Analysis of SS and Prox Intronic Mutations (top) and GO Enrichment Analysis Deep Intronic Mutations (bottom) using the g:Profiler online tool.

and (ii) SpliceAI (44) designed for detecting potential intronic variants causing pathogenic splicing (Figure 3C).

Out of the 78000 variants originally identified (Supplemental Table S1), 295 variants affecting 107 genes were scored as potentially pathogenic by at least one of these tools (Figure 3C). Through a direct comparison involving Pangolin, SpliceAI, and DeepVarSplice, we identified 34 genes depicted as tumor suppressors in pancreatic cancer in (Supplemental Figure S4). One explanation for the heightened sensitivity of variant detection using our approach is its ability to access variants deep within introns from the intronic regions of pre-mRNAs present in sn-RNA data. In fact, a limitation of SpliceAI and Pangolin is their optimization primarily for variants located within 50 bp on the splice site defined as SS and proximal intronic regions, affecting canonical splicing. Furthermore, it's important to acknowledge that both Pangolin and SpliceAI were developed predominantly using bulk-RNA sequencing methods, leading to potential gaps in the dataset due to the absence of intronic sequences. This limitation highlights the advantage of our approach in capturing a more comprehensive spectrum of intronic variants.

### Transcriptional regulation relates to mis-splicing in basal-like tumor cells

We performed transcriptional subtyping of single tumor nuclei into the more aggressive/drug-resistant ‘basal-like’ (BL) or the better prognosis-associated ‘classical’ subtype (CLA) using a ranking markers method described previously on snRNA-seq data (18). The distributions of each subtype in each tumor

showed BL cells highly represented in primary tumors, and partially in matched PDXs. In contrast, CL cells predominated in our organoid models (Figure 4A). It is important to note that the resolution of our analysis of these patient-matched models PDX and organoids undergo quite some selection that can change a substantial part of their transcriptional profiles. Specifically, tissue samples prepared for organoid generation are only small parts of the whole tumor. The higher heterogeneity of PDAC primary tumors questions the reliability of substituting small pieces for whole tumor tissues especially replicating the complexity of the patient-specific environment, e.g. tumor stroma and tumor types possess distinct immune components and different cell quantities affecting the cell composition in the early stage of tumoroid cultures.

Next, we performed differential expression analysis (DEG) between BL and CL tumor cells, where we identified a few hundred markers for both subtypes with absolute  $\log_2FC > 0.5$  and  $P_{adj} < 0.05$  (Supplemental Table S3) and visualized the most prominent ones in UMAP plots (Figure 4B).

Tumor heterogeneity was assessed by evaluating cell types using snRNA-seq from all models at hand. As expected, primary tumors exhibited the highest heterogeneity as regards cell type composition while organoids mostly contained ductal-like cells (Figures 5A-C). Strikingly, *RFXO1*, *CSMD1* and *CNTNAP2*, our topmost intronic mutated genes (Figures 1B and 4C) appear to be simultaneously the genes most upregulated in BL cells (Figure 5A). At the same time, the markers found in CL cells have already been described as biomarkers for pancreatic cancer: *MALAT1*, *NEAT1*, *CEACAM6* or *MUC1* (45–48). For the lncRNA *NEAT1*, two novel

**Table 1.** Classification of high intronic-mutated genes based on variant location

Gene_ID	Number_of_Total	Number_of_SS	Number_of_Prox	Number_of_Deep
RBFOX1	208	0	0	208
CSMD1	125	0	0	125
WWOX	80	0	0	80
BAGE2	75	0	0	75
SMYD3	64	0	0	64
MTUS2	61	0	0	61
CDH13	58	0	0	58
LRP1B	57	0	0	57
MAGI2	57	0	0	57
PTPRN2	52	0	0	52
DLG2	52	0	0	52
EXOC4	50	0	1	49
JMJD1C	48	0	0	48
LINC02055	48	0	0	48
NRG1	47	0	0	47
SLC30A10	47	0	0	47
SNX29	47	0	1	46
FHIT	47	0	0	47
NRXN3	45	0	0	45
HHAT	45	0	0	45
ALK	44	0	0	44
PRKN	43	0	0	43
NOS1AP	43	0	0	43
NAV2	44	0	1	43
CNTNAP2	42	0	0	42
KMT2C	41	0	0	41
THSD4	41	0	0	41
MAPK4	6	0	0	6
GNAQ	15	0	0	15
BCAS3	31	0	1	30
OAS1	4	11	0	-7
SUPT5H	3	1	0	2
ASAH1	3	1	1	1
POU6F2	14	1	0	13
ACTG1	5	2	4	-1
MDM2	4	1	0	3
ASH2L	3	1	0	2
EPHB6	1	5	3	-7
C12orf43	1	1	1	-1
RHNO1	2	1	0	0
HNRNP1	1	0	10	-9
MGAT5B	10	0	7	3
CARD8	1	0	9	-8
MACF1	23	0	5	18
COBLL1	11	0	6	5
MFF	6	0	5	1
MPHOSPH9	22	0	5	17
FOXRED1	1	0	5	-4
TBPL1	1	0	3	-2
POLR2L	2	0	1	1
ELOB	1	0	2	-1
TM4SF5	5	0	4	1
RPL26L1	4	0	5	-1
SNRNP35	1	0	1	0
INTS10	6	0	3	3
RPL27	3	0	2	1
DDR2	18	0	2	16
ADAMTS12	12	0	1	11
BCAT1	8	0	3	5

Classification of highly intronic mutated genes based on the variant location: donor and acceptor splice sites (SS, 1–2 bp away from the nearest exon-intron junction); proximal site (3–20 bp away from the nearest exon-intron junction), and deep intronic (>20 bp away from the nearest exon-intron junction). Total number of intronic mutations based on genes described in [Supplemental Table S2](#).

mutations are reported here for the first time ([Supplemental Table S1](#)). Taken together, CL cells and models are characterized by less abundance and relevance of intronic variants according to the output of our pipeline, thus suggesting that mis-splicing mechanisms are linked to the aggressiveness of BL tumors (Figure 5D).

## Contribution of intronic mis-splicing variants to pathogenic splicing and poor PDAC prognosis

In general, mis-splicing affects the regulation of genes (up- or downregulation) or generates new isoforms in normal and tumor tissues (49–51). However, our focus here is on the identification of potentially pathogenic splicing events in PDAC. It is important to note that the definition of pathogenicity in this study is based on comparisons to human references (via SNP calling) or to non-tumor pancreatic tissues when assessing mis-splicing findings. Nevertheless, we acknowledge that conclusive assertions regarding pathogenicity necessitate the inclusion of functional studies.

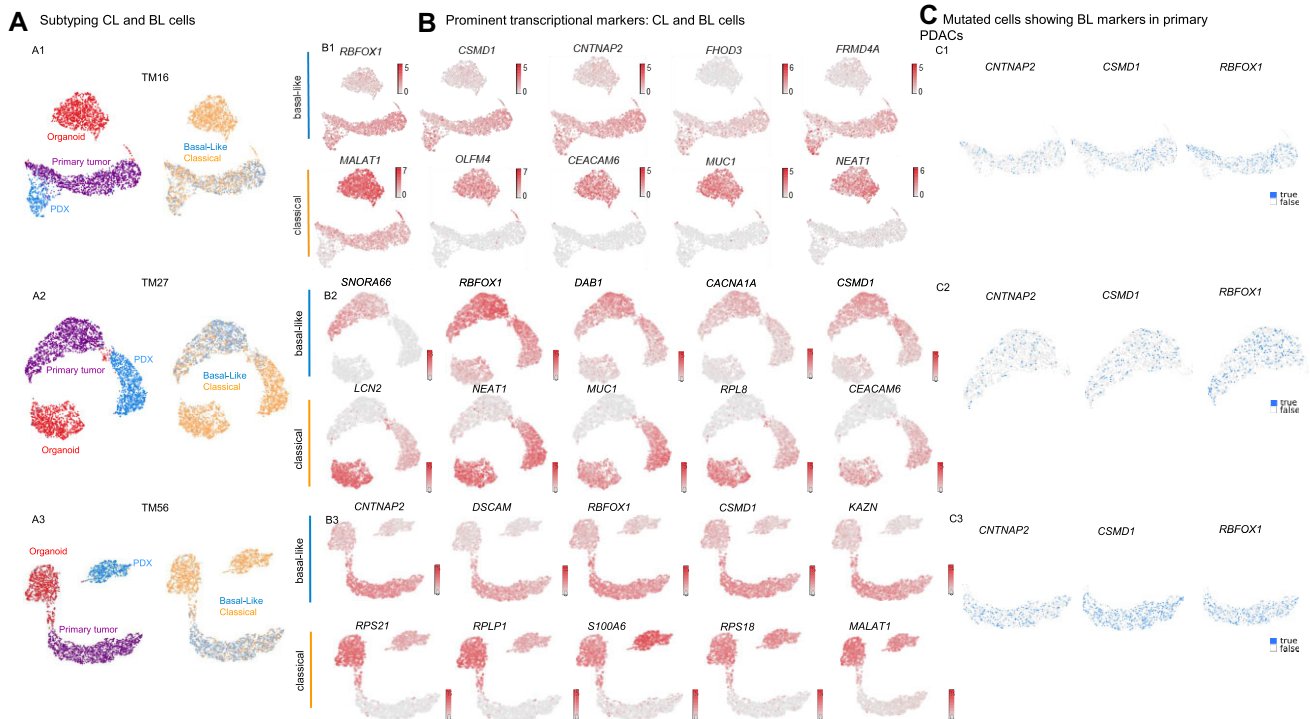
To determine intronic retention (IR), we employed the IRFinder tool by intersecting the findings of bulk RNA-seq from 24 primary PDACs recruited for CRU5002 and 9 public PDACs (GSE211398). To ensure robustness and significance, we included additional healthy and tumor pancreas tissues sourced from credible databases to visualize and verify the IRFinder results (SRA and GEO). A common challenge in utilizing public transcriptomic (bulk RNA-Seq) data for this research field is the limited sequencing depth, which often hinders precise IR event identification. To address this limitation proactively, we sequenced 24 primary PDAC samples from the CRU 5002 with higher depth, guaranteeing a minimum of 100 million reads per sample, specifically for splicing investigations resulting on a total of 2489 shared IR events (listed in [Supplemental Table S4](#)).

Next, we take a closer look into those genes found previously to be (i) most enriched for intronic variants (Table 1) and (ii) simultaneously found overexpressed in BL tumor cells (*RBFOX1*, *CSMD1* and *CNTNAP2* from Figures 1B and 4C). All three are long genes with many (mostly small-, <100 nt, and micro-, <60 nt) exons, and with relatively small introns. Moreover, all intronic variants discovered within these genes were situated deep within the intronic regions, establishing an association with non-canonical splicing regulation. Among these genes we have made a noteworthy discovery of combinatorial abnormal splicing (5) showing multiple splicing events within a single gene. For example, we reported a combinatorial splicing in a small portion of the *CSMD1* gene (chr8: 3187000–3202000) that showed an unusual exon skipping in the middle of the exon (chr8: 3189980–3190067) followed by two PEA events (chr8: 3197495–3197578 and 3200391–3200466; Figure 6A).

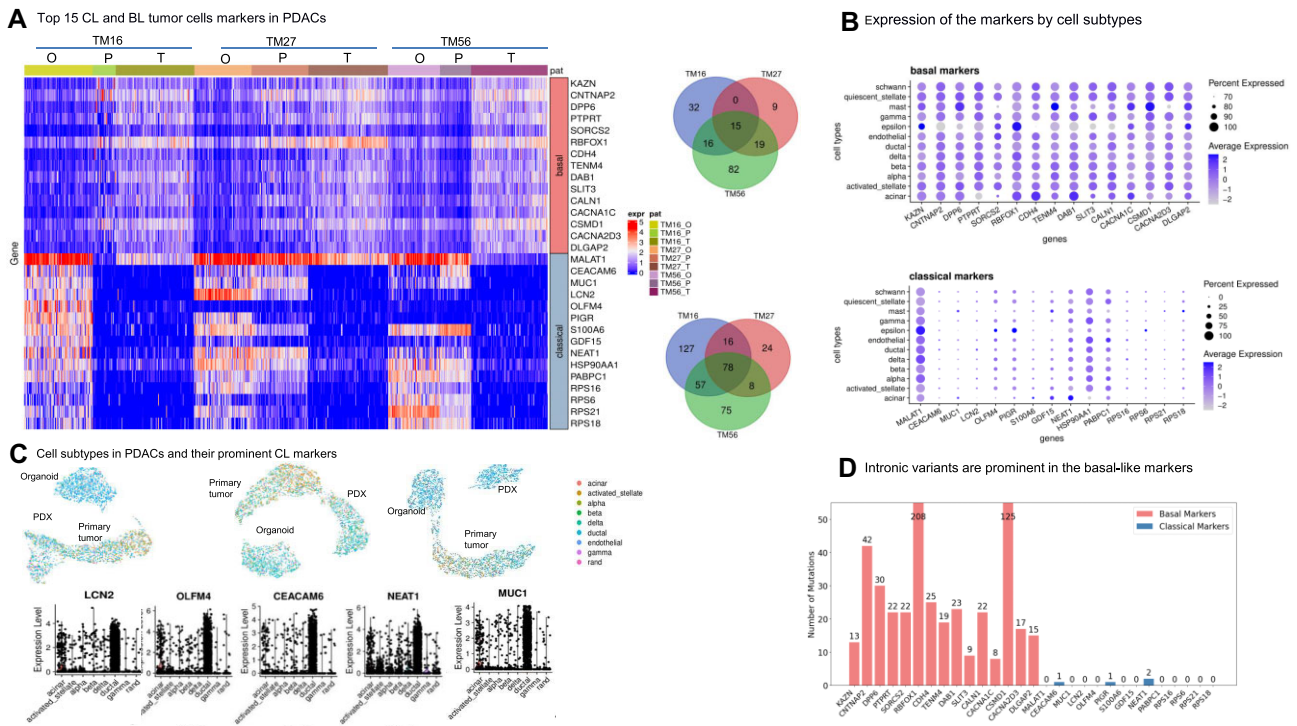
To assess potential non-canonical splicing mechanisms triggered by deep intronic variants in the *CSMD1* gene, we conducted a search for RNA binding motifs near the observed mis-splicing events. The presence of these variants generates motifs, such as SR-protein binding sites, leading to the activation of introns (PEA). Following motif analysis for SR-protein binding sites, we identified two intronic variants in our snRNA-seq data creating two *de novo* SRSF2 motifs (52); i.e. C > G substitution: CACGCT > CAGGCT; p.3194680 and CACGAA > CAGGAA; p.3196860 (Figure 6A). This suggests that intronic variants are located deep in the intron and creating additional SR binding sites that may synergistically contribute to activate PEA events. To this day, few examples of PEA have been reported as caused by deep intronic mutations without directly changing a splice site sequence (53,54).

In addition to genes with significant mutations deep within introns, our findings include other genes exhibiting high mutation rates in the SS and adjacent proximal intronic regions.

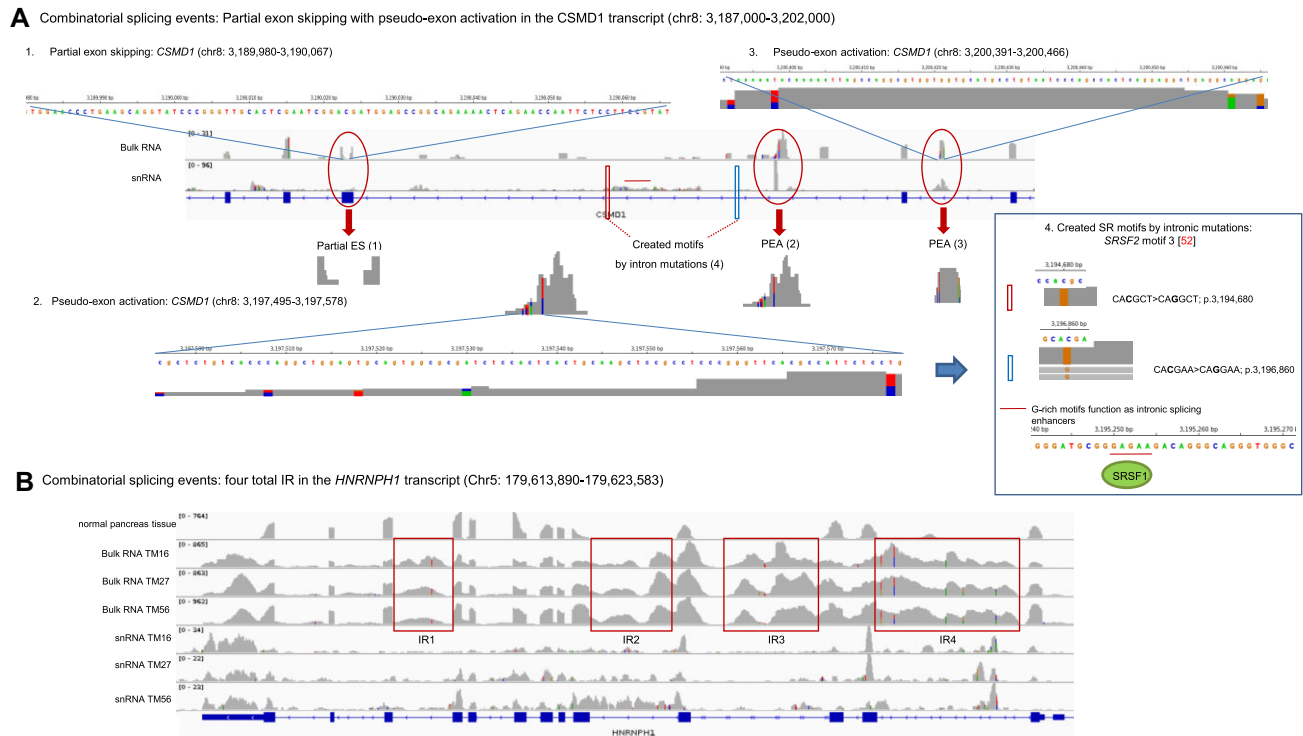




**Figure 4.** Tumor subtype and marker identification at the single-nucleus level. **(A)** UMAP plot of single cells from the three patients TM16, TM27 and TM 56 (colored by tumor model (left) and by inferred tumor subtype on CL and BL tumor pancreatic cells (right)). **(B)** UMAP plot showing the gene expression of the top transcriptional markers comparing CL vs BL cells in primary PDACs, organoids and PDXs. **(C)** UMAP plot showing mutated cells (intronic variants) for *CNTNAP2*, *CSMD1* and *RBFOX1* in primary PDACs.



**Figure 5.** Identification of markers for classical-like (CL) and basal-like (BL) tumor cells. **(A)** Heatmap showing gene expression of the top 15 BL and CL markers overexpressed in all PDACs: 'T' primary PDACs, 'O' organoids and 'P' PDXs. Venn diagrams illustrate the overlap of marker genes between tumor PDAC patient samples. **(B)** Dot plot showing expression of the markers in each of the pancreas cell subtypes. Circle size is proportional to the percentage of cells in each cell type expressing the marker and circle color represents the average marker gene expression in the cell type. **(C)** UMAP plots by patient showing clusters annotated to specific cell subtypes. The violin plots below represent prominent CL tumor pancreas markers across cell subtypes. **(D)** Table showing the number of intronic mutations of the top CL and BL marker genes.



**Figure 6.** Combinatorial splicing events related to intronic variants in PDAC. **(A)** Partial exon skipping with pseudo-exon activation in *CSMD1* (chr8: 3187000–3202000). 1: partial exon skipping in *CSMD1* (chr8: 3189980–3190067); 2: pseudo-exon activation in *CSMD1* (chr8: 3197495–3197578); 3: pseudo-exon activation in the *CSMD1* (chr8: 3200391–3200466). **(B)** Four combinatorial IR events in *HNRNPH1* (Chr5: 179613890–179623583).

One example is *HNRNPH1*, reported by our pipeline for carrying high number of intronic variants in proximal regions (Table 1). For *HNRNPH1* we detected combinatorial splicing involving four IR events (Figure 6B). *HNRNP* nucleoproteins are known to associate with pre-mRNAs in the nucleus and influence their processing and other aspects of mRNA metabolism and transport (55).

Several transcripts encoding members of the Integrator complex (*INTS10*, *INTS3*, *INTS11*) were affected by multiple IR events (see *INTS3* examples in Figure 7C and D). As the Integrator complex interacts with the C-terminal domain of RNA polymerase II to allow processing of U1 and U2 small nuclear RNAs, these splicing alterations could indirectly affect splicing in PDAC. Closer inspection of the two IR events in *INTS3* revealed a G-rich motifs function as intronic splicing enhancers. These give rise to partially overlapping recognition motifs for SRSF1 and SRSF2 (Figure 7E) and could function synergistically as splicing enhancers to compensate for weak PPTs tracts (56). The high intronic mutation load we uncovered using snRNA-seq suggests that the occurrence of such events at the single-cell level might be more frequent than previously presumed (54).

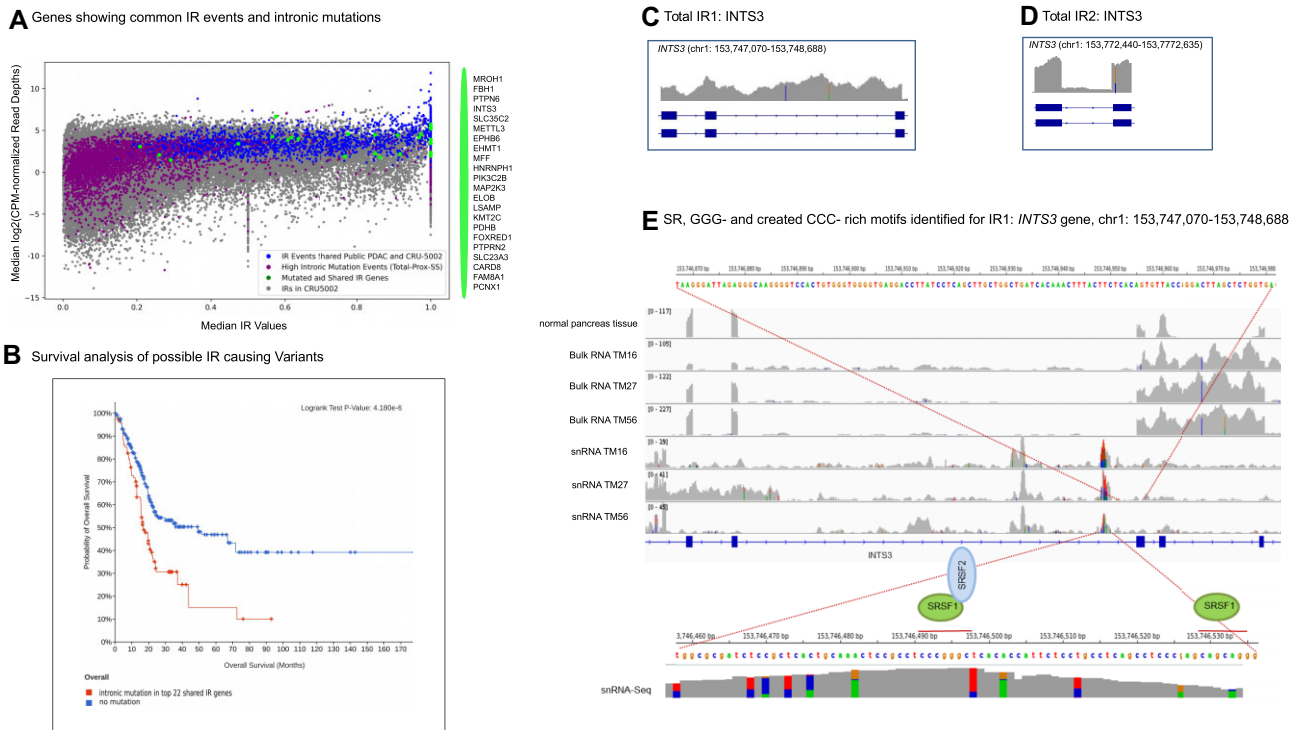
Finally, we assessed the clinical outcomes and significance of intronic variants associated with mis-splicing in pancreatic cancer by combining the sn-RNA-seq from the PDACs with the bulk RNA-Seq data from the PDACs (sourced from the CRU5002 and public data). The scatter plot depicted in Figure 7A illustrates shared intronic retention (IR) events, represented in blue, between the bulk RNA-Seq data from the CRU5002 and public datasets. Additionally, genes with a high mutational rate in introns (SS-Proximal and deep), identified in the sn-RNA-seq data from primary PDACs using our pipeline, are shown in green. Consequently, from this analysis,

we identify and highlight a subset of 22 candidates. Notably, survival analysis performed with the TCGA-PAAD cohort and generated by the cBioportal for the 22 genes revealed their significant association with poor prognosis in PDAC (Figure 7B). Six of these genes (*METTL3*, *HNRNPH1*, *INTS3*, *ELOB*, *EHMT1* and *KMT2C*) are linked to the GO and pathway enrichment of the most mutated intronic genes (SS, proximal and deep variants) related to RNA polymerase II transcriptional elongation, the TP53 regulates transcription DNA repair genes and the histone H3K4 trimethyltransferase activity previously described in Figure 3A and D.

## Discussion

The current study proposes a new strategy for the investigation of intronic mis-splicing variants and their role in promoting pathogenic splicing in PDAC. In contrast to prior investigations utilizing whole-genome sequencing and transcriptomic data, our approach employs a comprehensive full-length snRNA-seq method, offering high-resolution identification of intronic variants. For this new method we developed a comprehensive pipeline, the DeepVarSplice that takes advantage of integrating multi-omics information e.g. variant calling, transcriptomics and splicing within a cell and thus offers a more holistic view of the underlying molecular mechanisms in complex diseases.

Simultaneously, we address the challenges associated with DeepVarSplice in handling low-covered regions, non-uniform read distribution and thus, increased false positive findings. To overcome these limitations, we employed gold standard methodologies, specifically conducting whole-genome sequencing to validate our SNP calling performances using sn-RNA-seq data (57). Additionally, we utilized well-established



**Figure 7.** Analysis of IR related to intronic variants in PDAC. **(A)** Scatter plot of IR ratios from CRU5002 bulk RNA-seq data with their corresponding normalized read depths. Filtered and shared IR events between public PDAC and CRU-5002 data (blue). IR events of high intronic mutated genes (Total-SS-Proximal; green). List of genes exhibiting IR events related to high intronic mutations (right). **(B)** Kaplan–Meier survival analysis generated via the cBioPortal shows overall survival time of PDAC patients with or without mutations in the top 22 IR genes shared by all three tumors. **(C)** Region of high percentage IR event detected by IRFinder for the *INTS3* gene shared in CRU and public PDAC (Mean IRraio: 0.45). **(D)** Second IR event in the *INTS3* gene (mean IRraio: 0.25). **(E)** SR binding motifs identified in *INTS3*. Two CAGCAGG binding sites for *SRSF1* and one CTCCCGG motif for *SRSF2*.

artificial intelligence tools (AI tools), including Pangolin and SpliceAI to affirm our results concerning intronic variants associated with potential pathogenic splicing. Notably, both Pangolin and SpliceAI exhibit limitations when assessing intronic variants deep within introns, as they were primarily optimized for variants in the splice site (SS) and proximal region (within 50 bp of the splice site). Notably, these tools were originally developed using bulk RNA sequencing datasets. For the first time in this study, we employed both tools using our full-length sn-RNA-seq approach. Consequently, our DeepVarSplice pipeline demonstrated enhanced sensitivity in intronic variant detection, capitalizing on the presence of pre-mRNAs in sn-RNA data and achieving full coverage of both intronic and exonic regions within a gene. Lastly, it's crucial to emphasize that, for both splicing analysis and variant detection using sn-RNA-seq approaches, a deep sequencing strategy is vital to enhance sensitivity and reliability of findings. Consequently, we conducted 300 to 400K median reads per nucleus for snRNA-Seq and 100 to 200 million reads per sample, paired-end, with 300 cycles for bulk RNA sequencing.

In this study, we have identified a signature of mis-spliced genes in PDAC primary tumors, linking intronic retention events to potential pathogenic intronic mutations. Our findings reveal a correlation between mis-splicing and BL cells in pancreatic cancer, highlighting BL markers present in primary PDAC tumors and their corresponding models, PDX and organoids. Despite the observed high tumor heterogeneity in PDAC primary tumors to the models, there is a notable consistency in our findings across all PDACs investigated, evident in both transcriptional and mutational analyses, as depicted

in Figures 3C and 5C, respectively. These results align with a recent publication proposing a splicing signature specific to BL cells, thereby distinguishing PDAC subtypes and suggesting their potential utility as biomarkers (13).

Our approach delivers extensive information on the landscape of primary transcripts, especially near potential pathogenic splicing events. Thus, we can begin to decipher the complexity of RNA sequences acting as suppressors or activators of splicing in the context of PDAC. Recent high throughput characterization of exon splicing enhancer and silencer (ESE and ESS) motifs has indicated that pseudo-exons can be discriminated from genuine exons on the basis of their low ESE and high ESS contents (58,59). Moreover, these motifs can be found both in exons and introns, and function via the recruitment of sequence-specific RNA-binding proteins that can dictate splicing choices (60–62). Our analyses suggest that splicing dysregulation in PDAC can be linked to non-canonical sequence signals, i.e. to intronic variants affecting RNA motifs located deep in introns. For instance, we found several motifs for SR RNA-binding proteins (e.g. *SRSF1* and *SRSF2*), enriched for GGG and CCC, distributed in the SS-proximal regions of introns and, in some cases, exhibiting extensive overlap (Figure 7) (52,62). Although RNA-binding proteins like SR proteins bind to RNA with high sequence specificity, it is difficult to obtain well-defined consensus motifs for each of them (52,61,62). However, cooperation and competition between the *SRSF1* and *SRSF2* proteins have been reported in the regulation of alternative splicing events, which are related to synergistic and compensatory interactions with target RNA (49). Still, the exact mechanism by which these variants affect

the splicing machinery, and its downstream choices remains to be dissected, and very few examples of deep-intronic disease-causing mutations have been described to date (53,54).

The complex nature of splicing events and their regulation necessitates the development and implementation of new approaches to uncover nuanced relationships between mutations and splicing choices. We described a pseudo-bulk variant-calling pipeline, the DeepVarSplice exploiting dense and full-length coverage of snRNA-seq datasets to detect putatively pathogenic deep intron variants and link them to mis-splicing events occurring in PDAC tumors, while also gauging cell heterogeneity. As our method generates a multidimensional output that typically serves as the foundation for machine learning models (ranging from Support Vector Machines to Convolutional Neural Networks), we envisage a near-future combination that would vastly improve our understanding of the combinatorial code regulating splicing choices in the context of cancer.

### Data availability

All sequencing data has been deposited in the Gene Expression Omnibus, Bulk RNA-Sequencing GSE228844 and snRNA-Sequencing GSE229007. Any other relevant data is available from the authors upon reasonable request.

### Code availability

Scripts used to analyze the data and generate the figures and tables in this paper are available on GitHub (<https://github.com/UKHG-NIG/DeepVarSplice>) and Zenodo (<https://doi.org/10.5281/zenodo.11147217>).

### Supplementary data

Supplementary Data are available at NARGAB Online.

### Acknowledgements

We are grateful for excellent technical support provided by Jacqueline Fink and Waltraut Kopp, and to Hanibal Bohnenberger for technical advice and for providing patient-derived samples. We also thank Jeniffer Appelhans and Mark Bösherz for the generation of organoids derived from primary PDAC tumors.

*Author contributions:* G.S., A.P. and E.H. conceived the project and experiment design. V.E. and P.S. supported the experimental design and contributed to experimental data. G.S. led the method development. F.L. led the experimental data production. A.M. contributed to experimental data, E.T.D. and M.S. led and performed the data analysis. K.C supported data analysis and A.P. and G.S. wrote the manuscript. All authors have read and agreed to the submitted version of the manuscript.

### Funding

Deutsche Forschungsgemeinschaft [CRU5002 to G.S., P.S., A.P., V.E., E.H.].

### Conflict of interest statement

None declared.

## References

- Shiraishi,Y., Okada,A., Chiba,K., Kawachi,A., Omori,I., Mateos,R.N., Iida,N., Yamauchi,H., Kosaki,K. and Yoshimi,A. (2022) Systematic identification of intron retention associated variants from massive publicly available transcriptome sequencing data. *Nat. Commun.*, **13**, 5357.
- Zhang,P., Philippot,Q., Ren,W., Lei,W.T., Li,J., Stenson,P.D., Palacín,P.S., Colobran,R., Boisson,B., Zhang,S.Y., *et al.* (2022) Genome-wide detection of human variants that disrupt intronic branchpoints. *Proc. Natl. Acad. Sci. U.S.A.*, **119**, e2211194119.
- Park,E., Pan,Z., Zhang,Z., Lin,L. and Xing,Y. (2018) The expanding landscape of alternative splicing variation in human populations. *Am. J. Hum. Genet.*, **102**, 11–26.
- Le Hir,H., Nott,A. and Moore,M.J. (2003) How introns influence and enhance eukaryotic gene expression. *Trends Biochem. Sci.*, **28**, 215–220.
- Jung,H., Lee,K.S. and Choi,J.K. (2021) Comprehensive characterisation of intronic mis-splicing mutations in human cancers. *Oncogene*, **40**, 1347–1361.
- Daguenet,E., Dujardin,G. and Valcárcel,J. (2015) The pathogenicity of splicing defects: mechanistic insights into pre-mRNA processing inform novel therapeutic approaches. *EMBO Rep.*, **16**, 1640–1655.
- Scotti,M.M. and Swanson,M.S. (2016) RNA mis-splicing in disease. *Nat. Rev. Genet.*, **17**, 19–32.
- Qian,X., Wang,J., Wang,M., Igelman,A.D., Jones,K.D., Li,Y., Wang,K., Goetz,K.E., Birch,D.G., Yang,P., *et al.* (2021) Identification of deep-intronic splice mutations in a large cohort of patients with inherited retinal diseases. *Front. Genet.*, **12**, 276.
- Sahin,I., George,A. and Seyhan,A.A. (2021) Therapeutic targeting of alternative RNA splicing in gastrointestinal malignancies and other cancers. *Int. J. Mol. Sci.*, **22**, 11790.
- Escobar-Hoyos,L.F., Penson,A., Kannan,R., Cho,H., Pan,C.H., Singh,R.K., Apken,L.H., Hobbs,G.A., Luo,R., Lecomte,N., *et al.* (2020) Altered RNA splicing by mutant p53 activates oncogenic RAS signaling in pancreatic cancer. *Cancer Cell*, **38**, 198–211.
- Wan,L., Lin,K.T., Rahman,M.A., Ishigami,Y., Wang,Z., Jensen,M.A., Wilkinson,J.E., Park,Y., Tuveson,D.A. and Krainer,A.R. (2023) Splicing factor SRSF1 promotes pancreatitis and KRASG12D-mediated pancreatic cancer. *Cancer Discov.*, **13**, 1678–1695.
- Calabretta,S., Bielli,P., Passacantilli,I., Pilozzi,E., Fendrich,V., Capurso,G., Delle Fave,G. and Sette,C. (2016) Modulation of PKM alternative splicing by PTBP1 promotes gemcitabine resistance in pancreatic cancer cells. *Oncogene*, **35**, 2031–2039.
- Ruta,V., Naro,C., Pieraccioli,M., Leccese,A., Archibugi,L., Cesari,E., Panzeri,V., Allgöwer,C., Arcidiacono,P.G., Falconi,M., *et al.* (2024) An alternative splicing signature defines the basal-like phenotype and predicts worse clinical outcome in pancreatic cancer. *Cell Rep. Med.*, **5**, 101411.
- Raguraman,P., Balachandran,A.A., Chen,S., Diermeier,S.D. and Veedu,R.N. (2021) Antisense oligonucleotide-mediated splice switching: potential therapeutic approach for cancer mitigation. *Cancers (Basel)*, **13**, 5555.
- Peng,Q., Zhou,Y., Oyang,L., Wu,N., Tang,Y., Su,M., Luo,X., Wang,Y., Sheng,X., Ma,J., *et al.* (2022) Impacts and mechanisms of alternative mRNA splicing in cancer metabolism, immune response, and therapeutics. *Mol. Ther.*, **30**, 1018–1035.
- Shomroni,O., Sitte,M., Schmidt,J., Parbin,S., Ludewig,F., Yigit,G., Zelarayan,L.C., Streckfuss-Bömeke,K., Wollnik,B. and Salinas,G. (2022) A novel single-cell RNA-sequencing approach and its applicability connecting genotype to phenotype in ageing disease. *Sci. Rep.*, **12**, 4091–4091.
- Amezquita,R.A., Lun,A.T.L., Becht,E., Carey,V.J., Carpp,L.N., Geistlinger,L., Marini,F., Rue-Albrecht,K., Risso,D., Sonesson,C., *et al.* (2019) Orchestrating single-cell analysis with Bioconductor. *Nat. Methods*, **17**, 137–145.
- Aibar,S., González-Blas,C.B., Moerman,T., Huynh-Thu,V.A., Imrichova,H., Hulselmans,G., Rambow,F., Marine,J.C., Geurts,P.,

- Aerts, J., *et al.* (2017) SCENIC: single-cell regulatory network inference and clustering. *Nat. Methods*, **14**, 1083–1086.
19. Chan-Seng-Yue, M., Kim, J.C., Wilson, G.W., Ng, K., Figueroa, E.F., O’Kane, G.M., Connor, A.A., Denroche, R.E., Grant, R.C., McLeod, J., *et al.* (2020) Transcription phenotypes of pancreatic cancer are driven by genomic events during tumor evolution. *Nat. Genet.*, **52**, 231–240.
  20. Segerstolpe, Å., Palasantza, A., Eliasson, P., Andersson, E.M., Andréasson, A.C., Sun, X., Picelli, S., Sabirsh, A., Clausen, M., Bjursell, M.K., *et al.* (2016) Single-cell transcriptome profiling of human pancreatic islets in health and type 2 diabetes. *Cell Metab.*, **24**, 593–607.
  21. Aran, D., Looney, A.P., Liu, L., Wu, E., Fong, V., Hsu, A., Chak, S., Naikawadi, R.P., Wolters, P.J., Abate, A.R., *et al.* (2019) Reference-based analysis of lung single-cell sequencing reveals a transitional profibrotic macrophage. *Nat. Immunol.*, **20**, 163–172.
  22. Van der Auwera, G.A., Carneiro, M.O., Hartl, C., Poplin, R., del Angel, G., Levy-Moonshine, A., Jordan, T., Shakir, K., Roazen, D., Thibault, J., *et al.* (2013) From FastQ data to high-confidence variant calls: the genome analysis toolkit best practices pipeline. *Curr. Protoc. Bioinformatics*, **43**, 11.10.1–11.10.33.
  23. McLaren, W., Gil, L., Hunt, S.E., Riat, H.S., Ritchie, G.R.S., Thormann, A., Flicek, P. and Cunningham, F. (2016) The Ensembl variant effect predictor. *Genome Biol.*, **17**, 122.
  24. Gao, J., Aksoy, B.A., Dogrusoz, U., Dresdner, G., Gross, B., Sumer, S.O., Sun, Y., Jacobsen, A., Sinha, R., Larsson, E., *et al.* (2013) Integrative analysis of complex cancer genomics and clinical profiles using the cBioPortal. *Sci. Signal*, **6**, p11.
  25. Cock, P.J.A., Antao, T., Chang, J.T., Chapman, B.A., Cox, C.J., Dalke, A., Friedberg, I., Hamelryck, T., Kauff, F., Wilczynski, B., *et al.* (2009) Biopython: freely available Python tools for computational molecular biology and bioinformatics. *Bioinformatics*, **25**, 1422–1423.
  26. Kolberg, L., Raudvere, U., Kuzmin, I., Adler, P., Vilo, J. and Peterson, H. (2023) g:Profiler-interoperable web service for functional enrichment analysis and gene identifier mapping (2023 update). *Nucleic Acids Res.*, **51**, W207–W212.
  27. Middleton, R., Gao, D., Thomas, A., Singh, B., Au, A., Wong, J.J.L., Bomane, A., Cosson, B., Eyraes, E., Rasko, J.E.J., *et al.* (2017) IRFinder: Assessing the impact of intron retention on mammalian gene expression. *Genome Biol.*, **18**, 51.
  28. Bonnal, S.C., López-Oreja, I. and Valcárcel, J. (2020) Roles and mechanisms of alternative splicing in cancer - implications for care. *Nat. Rev. Clin. Oncol.*, **17**, 457–474.
  29. Tress, M.L., Abascal, F. and Valencia, A. (2017) Alternative splicing may not be the key to proteome complexity. *Trends Biochem. Sci.*, **42**, 98–110.
  30. Wang, X., He, Y., Zhang, Q., Ren, X. and Zhang, Z. (2021) Direct comparative analyses of 10X genomics chromium and smart-seq2. *Genomics Proteomics Bioinformatics*, **19**, 253–266.
  31. Kuo, H.C., Kuo, Y.R., Lee, K.F., Hsieh, M.C., Huang, C.Y., Hsieh, Y.Y., Lee, K.C., Kuo, H.L., Lee, L.Y., Chen, W.P., *et al.* (2017) A comparative proteomic analysis of erinacine A’s inhibition of gastric cancer cell viability and invasiveness. *Cell. Physiol. Biochem.*, **43**, 195–208.
  32. Huang, T., Liang, Y., Zhang, H., Chen, X., Wei, H., Sun, W. and Wang, Y. (2021) Csm1 mutations are associated with increased mutational burden, favorable prognosis, and anti-tumor immunity in gastric cancer. *Genes (Basel)*, **12**, 1715.
  33. Chen, H., Chong, W., Wu, Q., Yao, Y., Mao, M. and Wang, X. (2019) Association of LRP1B mutation with tumor mutation burden and outcomes in melanoma and non-small cell lung cancer patients treated with immune check-point blockades. *Front. Immunol.*, **10**, 1113.
  34. Datta, N., Chakraborty, S., Basu, M. and Ghosh, M.K. (2021) Tumor suppressors having oncogenic functions: the double agents. *Cells*, **10**, 46.
  35. Del Mare, S., Salah, Z. and Aqeilan, R.I. (2009) WWOX: its genomics, partners, and functions. *J. Cell. Biochem.*, **108**, 737–745.
  36. Salah, Z., Aqeilan, R. and Huebner, K. (2010) WWOX gene and gene product: tumor suppression through specific protein interactions. *Future Oncol.*, **6**, 249.
  37. Iliopoulos, D., Guler, G., Han, S.Y., Druck, T., Ottey, M., McCorkell, K.A. and Huebner, K. (2006) Roles of FHIT and WWOX fragile genes in cancer. *Cancer Lett.*, **232**, 27–36.
  38. Singh, S.K., Kumar, S., Viswakarma, N., Principe, D.R., Das, S., Sondarva, G., Nair, R.S., Srivastava, P., Sinha, S.C., Grippo, P.J., *et al.* (2021) MAP4K4 promotes pancreatic tumorigenesis via phosphorylation and activation of mixed lineage kinase 3. *Oncogene*, **40**, 6153–6165.
  39. Sun, S., Zhang, Z., Fregoso, O. and Krainer, A.R. (2012) Mechanisms of activation and repression by the alternative splicing factors RBFOX1/2. *RNA*, **18**, 274–283.
  40. Lin, F., Xu, L., Yuan, R., Han, S., Xie, J., Jiang, K., Li, B., Yu, W., Rao, T., Zhou, X., *et al.* (2022) Identification of inflammatory response and alternative splicing in acute kidney injury and experimental verification of the involvement of RNA-binding protein RBFOX1 in this disease. *Int. J. Mol. Med.*, **49**, 32.
  41. Jbara, A., Lin, K.T., Stossel, C., Siegfried, Z., Shqerat, H., Amar-Schwartz, A., Elyada, E., Mogilevsky, M., Raitses-Gurevich, M., Johnson, J.L., *et al.* (2023) RBFOX2 modulates a metastatic signature of alternative splicing in pancreatic cancer. *Nature*, **617**, 147–153.
  42. Cowling, V.H. (2009) Regulation of mRNA cap methylation. *Biochem. J.*, **425**, 295–302.
  43. Zeng, T. and Li, Y.I. (2022) Predicting RNA splicing from DNA sequence using Pangolin. *Genome Biol.*, **23**, 103.
  44. Jaganathan, K., Kyriazopoulou Panagiotopoulou, S., McRae, J.F., Darbandi, S.F., Knowles, D., Li, Y.I., Kosmicki, J.A., Arbelaez, J., Cui, W., Schwartz, G.B., *et al.* (2019) Predicting splicing from primary sequence with deep learning. *Cell*, **176**, 535–548.
  45. Xiong, D.D., Feng, Z.B., Cen, W.L., Zeng, J.J., Liang, L., Tang, R.X., Gan, X.N., Liang, H.W., Li, Z.Y., Chen, G., *et al.* (2017) The clinical value of lncRNA NEAT1 in digestive system malignancies: a comprehensive investigation based on 57 microarray and RNA-seq datasets. *Oncotarget*, **8**, 17665.
  46. Duxbury, M.S., Matros, E., Clancy, T., Bailey, G., Doff, M., Zinner, M.J., Ashley, S.W., Maitra, A., Redston, M. and Whang, E.E. (2005) CEACAM6 is a novel biomarker in pancreatic adenocarcinoma and PanIN lesions. *Ann. Surg.*, **241**, 491–496.
  47. Striefler, J.K., Riess, H., Lohneis, P., Bischoff, S., Kurreck, A., Modest, D.P., Bähr, M., Oettle, H., Sinn, M., Bläker, H., *et al.* (2021) Mucin-1 protein is a prognostic marker for pancreatic ductal adenocarcinoma: results from the CONKO-001 study. *Front. Oncol.*, **11**, 1.
  48. Kumar, S. and Mishra, S. (2022) MALAT1 as master regulator of biomarkers predictive of pan-cancer multi-drug resistance in the context of recalcitrant NRAS signaling pathway identified using systems-oriented approach. *Sci. Rep.*, **12**, 7540.
  49. Wong, J.J.L., Au, A.Y.M., Ritchie, W. and Rasko, J.E.J. (2016) Intron retention in mRNA: No longer nonsense: known and putative roles of intron retention in normal and disease biology. *Bioessays*, **38**, 41–49.
  50. Jung, H., Lee, D., Lee, J., Park, D., Kim, Y.J., Park, W.Y., Hong, D., Park, P.J. and Lee, E. (2015) Intron retention is a widespread mechanism of tumor-suppressor inactivation. *Nat. Genet.*, **47**, 1242–1248.
  51. Ge, Y. and Porse, B.T. (2014) The functional consequences of intron retention: alternative splicing coupled to NMD as a regulator of gene expression. *Bioessays*, **36**, 236–243.
  52. Pandit, S., Zhou, Y., Shiue, L., Coutinho-Mansfield, G., Li, H., Qiu, J., Huang, J., Yeo, G.W., Ares, M. and Fu, X.D. (2013) Genome-wide analysis reveals SR protein cooperation and competition in regulated splicing. *Mol. Cell*, **50**, 223–235.
  53. Davis, R.L., Homer, V.M., George, P.M. and Brennan, S.O. (2009) A deep intronic mutation in FGB creates a consensus exonic splicing enhancer motif that results in afibrinogenemia caused by aberrant

- mRNA splicing, which can be corrected in vitro with antisense oligonucleotide treatment. *Hum. Mutat.*, **30**, 221–227.
54. Homolova, K., Zavadakova, P., Doktor, T.K., Schroeder, L.D., Kozich, V. and Andresen, B.S. (2010) The deep intronic c.903+469T>C mutation in the MTRR gene creates an SF2/ASF binding exonic splicing enhancer, which leads to pseudoexon activation and causes the cblE type of homocystinuria. *Hum. Mutat.*, **31**, 437–444.
55. Vo, T., Brownmiller, T., Hall, K., Jones, T.L., Choudhari, S., Grammatikakis, I., Ludwig, K.R. and Caplen, N.J. (2022) HNRNPH1 destabilizes the G-quadruplex structures formed by G-rich RNA sequences that regulate the alternative splicing of an oncogenic fusion transcript. *Nucleic Acids Res.*, **50**, 6474–6496.
56. Murray, J.I., Voelker, R.B., Henscheid, K.L., Warf, M.B. and Berglund, J.A. (2008) Identification of motifs that function in the splicing of non-canonical introns. *Genome Biol.*, **9**, R97.
57. Brouard, J.S., Schenkel, F., Marete, A. and Bissonnette, N. (2019) The GATK joint genotyping workflow is appropriate for calling variants in RNA-seq experiments. *J. Anim. Sci. Biotechnol.*, **10**, 44.
58. Grellscheid, S.-N. and Smith, C.W.J. (2006) An apparent pseudo-exon acts both as an alternative exon that leads to nonsense-mediated decay and as a zero-length exon. *Mol. Cell Biol.*, **26**, 2237–2246.
59. Han, K., Yeo, G., An, P., Burge, C.B. and Grabowski, P.J. (2005) A combinatorial code for splicing silencing: UAGG and GGGG motifs. *PLoS Biol.*, **3**, e158.
60. Singh, R. (2002) RNA–protein interactions that regulate pre-mRNA splicing. *Gene Expr.*, **10**, 79.
61. Rahman, M.A., Lin, K.T., Bradley, R.K., Abdel-Wahab, O. and Krainer, A.R. (2020) Recurrent SRSF2 mutations in MDS affect both splicing and NMD. *Genes Dev.*, **34**, 413–427.
62. Jankowsky, E. and Harris, M.E. (2015) Specificity and nonspecificity in RNA–protein interactions. *Nat. Rev. Mol. Cell Biol.*, **16**, 533–544.