

Pan-cancer repository of validated natural and cryptic mRNA splicing mutations



Peter K. Rogan^{1,3}, Ben C. Shirley³, Eliseos J. Mucaki¹, Joan H.M. Knoll^{2,3}

University of Western Ontario Departments of Biochemistry¹ and Pathology and Laboratory Medicine², CytoGnomix Inc.³ London, ON, Canada

Abstract

We present a major public resource of mRNA splicing mutations validated according to multiple lines of evidence of abnormal gene expression. Likely mutations present in all tumor types reported in the Cancer Genome Atlas (TCGA) and the International Cancer Genome Consortium (ICGC) were identified based on the comparing counts of splice junction spanning and abundance of transcript reads in RNA-Seq data from matched tissues and tumors lacking these mutations, the majority of which (69.9%) are not present in the Single Nucleotide Polymorphism Database (dbSNP) 150). There are 131,347 unique mutations which strengthen cryptic splice sites, and 222,071 mutations which strengthen cryptic splice sites, and 222,071 mutations which strengthen cryptic splice sites (11,932 affect both simultaneously). dbSNP) were observed in multiple tumor tissue types. Single variants or chromosome ranges can be queried using a Global Alliance for Genomics and Health (GA4GH)-compliant, web-based Beacon "Validated Splicing Mutations" either separately or in aggregate alongside other Beacons through the public Beacon Network (<u>http://www.beacon-network.org/#/search?beacon=cytognomix</u>), as well as through our website (<u>http://validsplicemut.cytognomix.com/</u>).

Introduction

Next generation sequencing continues to reveal large numbers of novel variants whose impact cannot be interpreted from curated variant databases, or through reviews of peer-reviewed biomedical literature¹. This has created a largely unmet need for unequivocal sources of information regarding the molecular phenotypes and potential pathology of variants of unknown significance (VUS) in cancer genomes. Information theory (IT) has been proven to accurately predict impact of mutations on mRNA splicing, and has been used to interpret coding and non-coding mutations that alter mRNA splicing in both common and rare diseases^{2,3}. We have described an ITbased framework for the interpretation and prioritization of non-coding variants of uncertain significance, which has been validated in multiple studies involving variants in patients with history or predisposition to heritable breast and/or ovarian cancer.⁴ The Cancer Genome Atlas (TCGA) Pan-Cancer Atlas is a comprehensive integrated genomic and transcriptomic resource containing data from >10,000 tumors across 33 different tumor types.⁵ Here, we utilized IT-based tools for assessment of high quality sequenced variants in TCGA patients, as well as patients from tumor datasets provided by the International Cancer Genome Consortium (ICGC), for their potential impact on mRNA splicing.⁶ The results of these genome-wide analyses are presented using an online internet resource (<u>http://validsplicemut.cytognomix.com/</u>; Figure 1 [below]) which can also be queried through the Beacon Network.⁷

Results

Examples of IT-predicted and RNAseq-validated mRNA Splicing Mutations in:

Table 3: Validated Splicing Mutations in COSMIC Cancer Gene Census genes

Gene Splice Mutation R. (bits) Tumor Observed Splicing Event

754

| A. GRCh37 1 | 1 : 108214098 G>T | | | Search |
|--|----------------------------|---|-------------------|------------------|
| | Or instead: Query by gene | or range of coordinates (click to expand) | > | |
| VARIANT POSITION | | | | |
| Genomic position (g. notation) | Gene-centric HGVS notatio | n (c. notation) | | |
| chr11:g.108214098G>T LRG_135t1:c.8418G>T; NM_000051.3:c.8418G>T; NM_138292.3:c.4374G>T; XM_005271561.1:c.8418G XM_005271562.1:c.8418G>T; XM_005271563.1:c.8418G>T; XM_005271564.1:c.7374G>T | | | | |
| SPLICE SITE INFORMATION | | | | |
| Splice Site Coordinate | R_i before mutation (1) | R_i after mutation (1) | Splice Type | Site Type |
| 108214099 | 8.6742 | ↓ 5.0805 | DONOR | NATURALSITE |
| VARIANT DATA | | | | |
| Gene | rsID (dbSNP15 |)) | Average Heterozyg | osity (dbSNP150) |
| ATM | rs762744146 | | 0.0000 | |
| INDIVIDUALS | | | | |
| - TCGA-BH-A1ET (BRCA) | | | | |
| View TCGA-BH-A1ET metadata | (NCI Genomic Data Commons) | | | |

| Junction spanning | 0 | 0 | 49 (n=0) | 4 (p=0.1708) | 0 |
|----------------------------|----------------|---------------------|---------------------------|---------------------------|------------------------------------|
| Evidence Type | Cryptic 🕕 | Anti-Cryptic 🕕 | Exon Skipping 🚺 | Intron Inclusion 🕕 | Intron Inclusion with Mutation (1) |
| Veridical validated this m | nutation based | on 112 reads each o | of which either overlap t | he splice boundary or are | wholly contained within an intron. |

| COSIVIL Cancer Gene Census genes |
|----------------------------------|
|----------------------------------|

| Evidence Type | Cryptic | Anti-Cryptic | Exon Skipping | Intron Inclusion | Intron Inclusion with Mutation |
|-------------------|---------|--------------|---|------------------|--------------------------------|
| Junction spanning | 0 | 0 | 0 | 12 (p=0.0002) | 11 (p=0) |
| Read Abundance | 0 | 0 | 0 | 29 (p=0.1813) | 0 |
| chr17 7,577,500 | Ьр | I | 7,577,600 bp | I | 7,577,700 bp |
| [0 - 78] | | | | | |
| | | | | | |
| | | | | | |
| | | ×* | | | |
| | | i i | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | < < < | < < < < < | < < < < + + + + + + + + + + + + + + + + | • • • • • | • • • • • • • • • |

TP53 mutation 17:7577609C>T (c.673-1G>A) abolishes the natural acceptor site of exon 7 (6.0 > -4.9 bits), activating a cryptic site 49nt downstream (R_i = 5.2 bits; blue arrow). Veridical detects 11 reads with mutation which crosses the intron-exon junction (brown), and reports their number significant over controls ($p \approx 0$).



Baylor-Hopkins Centre for Medical Genomics study⁹

| | Cryptic | Anu-Crypuc | Exon Skipping | intron in | clusion | intron inclu | usion with Mutati |
|---|---|--|--|---|--|---|---|
| unction spanning | 0 | 0 | 0 | 12 (p=0.0 | 046) | 9 (p=0) | |
| Read Abundance | 0 | 0 | 0 | 221 (p=0 |) | 0 | |
| chr21 44,272,300 bp | I | 44,272 | 400 bp | 1 | 44,272,500 b | P | 1 |
| [0 - 62] | | | | | | | |
| | | | | | | | |
| | | | - | _ | | | |
| | | | | | | | |
| | | - $($ $($ $($ $))$ | WDR4 | | | | |
| WDR4 21:44272 | 435C>T (c | .976-1G>A) a | abolishes exon | 10 natural | acceptor (| 5.3 > -5.6 | bits). Veridic |
| WDR4 21:442724 finds intron inclu | 435C>T (c ision reads | .976-1G>A) and a mutation | wDR4 abolishes exon a d reads with mu | 10 natural Itation cros | acceptor (| 5.3 > -5.6 | bits). Veridio exon bounda |
| WDR4 21:442724 finds intron inclu (brown) significa | 435C>T (c ision reads int (p≈0). 1 | .976-1G>A) and the mutation | abolishes exon a d reads with mu was found in L | 10 natural Itation cros | acceptor (sing over t t TCGA-55 | 5.3 > -5.6 the intron- -8085 as v | bits). Veridio exon bounda well as BH CN |
| WDR4 21:44272 finds intron inclu (brown) significa patient BH2112 | 435C>T (c ision reads int ($p\approx 0$). 1 2_1 with | .976-1G>A) and s (purple) and The mutation https://gen | abolishes exon d reads with mu was found in L omebiology.bio | 10 natural Itation cros UAD patien medcentral | acceptor (sing over t t TCGA-55 .com/artic | 5.3 > -5.6 the intron -8085 as v cles/10.11 | 5 bits). Veridio -exon bounda well as BH CN 86/s13059-01 |
| WDR4 21:44272 finds intron inclu (brown) significa patient BH2112 0779-xsyndrome | 435C>T (c ision reads int (p \approx 0). 1 2_1 with b. <i>WDR4</i> ha | .976-1G>A) and s (purple) and The mutation https://gen as been linked | wbrad abolishes exon a d reads with mu was found in L homebiology.bio d to disorders with | 10 natural Itation cros UAD patien medcentral | acceptor (sing over t t TCGA-55 .com/artic ephaly ph | 5.3 > -5.6 the intron -8085 as v cles/10.113 enotype. | 5 bits). Veridio -exon bounda well as BH CN 86/s13059-01 |
| WDR4 21:442724 finds intron inclu (brown) significa patient BH2112 0779-xsyndrome | 435C>T (c ision reads int ($p\approx 0$). T 2_1 with b. <i>WDR4</i> ha | .976-1G>A) a s (purple) and The mutation https://gen as been linked | wDR4 abolishes exon a d reads with mu was found in L nomebiology.bio d to disorders wi | 10 natural Itation cros UAD patien medcentral Th a microc | acceptor (sing over t t TCGA-55 .com/artic ephaly ph | 5.3 > -5.6 the intron- -8085 as v cles/10.112 enotype. | 5 bits). Veridio exon bounda well as BH CN 86/s13059-01 |
| WDR4 21:442724 finds intron inclu (brown) significa patient BH2112 0779-xsyndrome Evidence Type Junction spanning | 435C>T (c ision reads int (p≈0). 1 2_1 with b. <i>WDR4</i> ha Cryptic 10 (p=0) | .976-1G>A) and s (purple) and The mutation https://gen as been linked Anti-C | wbR4 abolishes exon a d reads with mu was found in L nomebiology.bio d to disorders wi ryptic Exon Ski | 10 natural Itation cros UAD patien medcentral Th a microc | acceptor (sing over t t TCGA-55 .com/artic ephaly ph n Inclusion 0.7447) | 5.3 > -5.6 the intron- -8085 as v cles/10.112 enotype. Intron Incl 1 (p=0) | 5 bits). Veridio exon bounda well as BH CN 86/s13059-01 |
| WDR4 21:442724 finds intron inclu (brown) significa patient BH2112 0779-xsyndrome Evidence Type Junction spanning Read Abundance | 435C>T (c ision reads int (p≈0). 1 2_1 with b. WDR4 ha Cryptic 10 (p=0) 2.531(p= | .976-1G>A) a s (purple) and The mutation https://gen as been linked <u>Anti-C</u> 0 0.444) 0.041 (| abolishes exon a d reads with mu was found in L nomebiology.bio d to disorders with ryptic Exon Ski 0 p=0.784) 0 | 10 natural Itation cros UAD patien medcentral Th a microc pping Intro 1 (p= 74 (p | acceptor (sing over to t TCGA-55 .com/artic ephaly ph n Inclusion 0.7447) =0.6878) | 5.3 > -5.6 the intron- -8085 as v cles/10.112 enotype. Intron Incl 1 (p=0) 0 | 5 bits). Veridio exon bounda well as BH CN 86/s13059-01 |
| WDR4 21:442724 finds intron inclu (brown) significa patient BH2112 0779-xsyndrome Evidence Type Junction spanning Read Abundance | 435C>T (c ision reads int (p≈0). 1 2_1 with 2. WDR4 ha Cryptic 10 (p=0) 2.531(p=0 42,183,400 bp | .976-1G>A) a s (purple) and the mutation https://gen as been linked Anti-C 0 0.444) 0.041 (| wbrad abolishes exon a d reads with mu was found in L bomebiology.bio d to disorders with ryptic Exon Ski 0 p=0.784) 0 42,18 | 10 natural Itation cros UAD patien medcentral Tth a microc pping Intro 1 (p= 74 (p | acceptor (sing over 1 t TCGA-55 .com/artic ephaly ph n Inclusion 0.7447) =0.6878) | 5.3 > -5.6 the intron -8085 as v cles/10.112 enotype. Intron Incl 1 (p=0) 0 | bits). Veridio exon bounda well as BH CN 86/s13059-01 |
| WDR4 21:442724 finds intron inclu (brown) significa patient BH2112 0779-xsyndrome Evidence Type Junction spanning Read Abundance | 435C>T (c ision reads int (p≈0). 1 2_1 with b. <i>WDR4</i> ha Cryptic 10 (p=0) 2.531(p=0 42,183,400 bp | .976-1G>A) a s (purple) and The mutation https://gen as been linked Anti-C 0 0.444) 0.041 (| wbrad wb | 10 natural Itation cros UAD patien medcentral th a microc pping Intro 1 (p= 74 (p | acceptor (sing over 1 t TCGA-55 .com/artic ephaly ph n Inclusion 0.7447) =0.6878) | 5.3 > -5.6 the intron -8085 as v cles/10.112 enotype. Intron Incl 1 (p=0) 0 | b bits). Veridio exon bounda well as BH CN 86/s13059-01 |
| WDR4 21:442724 finds intron inclu (brown) significa patient BH2112 0779-xsyndrome Evidence Type Junction spanning Read Abundance | 435C>T (c ision reads int (p≈0). 1 2_1 with 0. WDR4 ha Cryptic 10 (p=0) 2.531(p=0 42,183,400 bp | .976-1G>A) a s (purple) and The mutation https://gen as been linked Anti-C 0 0.444) 0.041 (| wbR4 abolishes exon a d reads with mu was found in L bomebiology.bio d to disorders with ryptic Exon Ski 0 p=0.784) 0 42,18 | 10 natural Itation cros UAD patien medcentral Th a microc pping Intro 1 (p= 74 (p | acceptor (sing over 1 t TCGA-55 .com/artic ephaly ph n Inclusion 0.7447) =0.6878) | 5.3 > -5.6 the intron -8085 as v cles/10.113 enotype. Intron Incl 1 (p=0) 0 | bits). Veridio exon bounda well as BH CN 86/s13059-01 |
| WDR4 21:442724 finds intron inclu (brown) significa patient BH2112 0779-xsyndrome Evidence Type Junction spanning Read Abundance | 435C>T (c ision reads int (p≈0). 1 2_1 with 2. WDR4 ha Cryptic 10 (p=0) 2.531(p=1 42,183,400 bp | .976-1G>A) a s (purple) and The mutation https://gen as been linked Anti-C 0 0.444) 0.041 (| wDR4 abolishes exon a d reads with mu was found in L nomebiology.bio d to disorders with ryptic Exon Ski 0 p=0.784) 0 42,18 | 10 natural Itation cros UAD patien medcentral th a microc pping Intro 1 (p= 74 (p | acceptor (sing over 1 t TCGA-55 .com/artic ephaly ph n Inclusion 0.7447) =0.6878) | 5.3 > -5.6 the intron -8085 as v cles/10.112 enotype. Intron Incl 1 (p=0) 0 | 5 bits). Veridio exon bounda well as BH CN 86/s13059-01 |
| WDR4 21:442724 finds intron inclu (brown) significa patient BH2112 0779-xsyndrome Evidence Type Junction spanning Read Abundance | 435C>T (c ision reads int (p≈0). 1 2_1 with 2. WDR4 ha Cryptic 10 (p=0) 2.531(p=0 42,183,400 bp | .976-1G>A) a s (purple) and The mutation https://gen as been linked Anti-C 0 0.444) 0.041 (| wDR4 abolishes exon a d reads with mu was found in L nomebiology.bio d to disorders with ryptic Exon Ski 0 p=0.784) 0 42,18 | 10 natural Itation cros UAD patien medcentral th a microc pping Intro 1 (p= 74 (p | acceptor (sing over 1 t TCGA-55 .com/artic ephaly ph n Inclusion 0.7447) =0.6878) | 5.3 > -5.6 the intron -8085 as v cles/10.112 enotype. Intron Incl 1 (p=0) 0 | bits). Veridio exon bounda well as BH CN 86/s13059-01 |

| IKDKD |
|--|
| IKBKB mutation 8:42183487G>C (c.1981-1G>C) abolishes the natural acceptor of exon 20 (9.9 -> -1.7 |
| bits) while simultaneously strengthening a cryptic acceptor (2.0 > 10.8 bits), resulting in a 2nt deletion |
| (blue reads). This mutation, identified in BLCA patient TCGA-FD-A3NA, was also found in a BH CMG |
| patient with immunodeficiency (BH6169 1), which <i>IKBKB</i> has been linked to. |

| on | CASC5 | <u>15:40942786G>A</u> (c.6212+5G>A) | 4.8 > 1.7 (Natural Site) | AML | The natural donor site of <i>CASC5</i> exon 19 (NM_144508) is weakened, leading to a significant increase in intron inclusion. |
|-------------------|--------|---|---|------|---|
| 1 | DNMT3A | <u>2:25467022A>G</u> (c.1851+2T>C) | 3.6 > -3.5 (Natural Site) | AML | The natural donor site of <i>DNMT3A</i> exon 15 (NM_022552) is abolished, resulting in a increase in total exon skipping and intron inclusion. |
| | STAG2 | <u>X:123176495G>A</u> (c.462G>A) | 6.5 > 3.5 (Natural Site) | BLCA | The natural donor of <i>STAG2</i> exon 6 (NM_006603) is weakened, and a significant amount of exon 6 skipping is observed. |
| | STAG2 | X:123200024G>A (c.2097-1G>A) | 19.5 > 8.6 (Natural Site) | BLCA | The natural acceptor of <i>STAG2</i> exon 21 (NM_006603) is weakened, resulting in a significant increase in exon 21 skipping. |
| | ATM | <u>11:108214098G>T</u> (c.8418G>T) | 8.7 > 5.1 (Natural Site) | BRCA | A natural donor site is weakened, leading to a significant increase in <i>ATM</i> exon 57 (NM_000051) skipping events. Some reads with mutation depict wildtype, leaky splicing. |
| | BARD1 | <u>2:215645882A>T</u> (c.716T>A) | 0.9 > 3.1 (Cryptic Site) | BRCA | The mutation strengthens a cryptic site within <i>BARD1</i> exon 4 (NM_000465). Reads which use this activated cryptic site contain the mutation (one exception). Some reads with mutation depict wildtype, leaky splicing. |
| • | GATA3 | <u>10:8115701G>C</u> (c.1048-1G>C) | 0.9 > -10.7 (Natural Site) | BRCA | Mutation abolishes the natural acceptor of <i>GATA3</i> exon 6 (NM_002051). This increases the use of a pre-existing exonic cryptic splice site (4.2 > 5.6 bits; 8nt deletion) and significantly increases total intron inclusion. |
| / 5 | TP53 | <u>17:7577609C>T</u> (c.673-1G>A) | 6.0 > -4.9 (Natural Site) | BRCA | A natural acceptor site is abolished, activating a cryptic site 49nt upstream (R_i = 5.2 bits) of <i>TP53</i> exon 7 (NM_000546). |
| - | POLD1 | <u>19:50920353A>G</u> (c.3119A>G) | 8.6 > 6.1 (Natural Site) | COAD | The natural donor of <i>POLD1</i> exon 25 (NM_002691) is weakened, leading to a significant increase in overall exon skipping. |
| on | SMAD3 | <u>15:67482748C>G</u> (c.1155-3C>G) | 11.9>3.1 -4.0>7.7 (<u>Natural</u> <u>Cryptic</u>) | COAD | Mutation weakens the natural acceptor of <i>SMAD3</i> exon 9 (NM_005902) and predicts a cryptic site that does not appear to be used. A significant number of intron inclusion reads are observed. A distant pre-existing cryptic acceptor (9.6 bits; 3598nt from natural acceptor) was observed. |
| | PIK3R1 | 5:67591246A>G (c.936-2A>G) | 7.5 > -7.3 (Natural Site) | GBM | The natural acceptor of <i>PIK3R1</i> exon 8 (NM_181504) is abolished, which promotes a significant increase in exon 8 skipping. |
| | FAT1 | <u>4:187521515C>A</u> (c.11641-1G>T) | 5.3 > -2.4 (Natural Site) | HNSC | Natural acceptor of <i>FAT1</i> exon 22 (NM_005245) is abolished, resulting in intron inclusion, and use of 3 cryptic sites; 82nt upstream of acceptor, and 237nt and 234nt downstream from the natural acceptor. |
| | TGFBR2 | <u>3:30729875G>A</u> (c.1397-1G>A) | 8.4 > -2.5 (Natural Site) | HNSC | <i>TGFBR2</i> exon 6 natural acceptor (NM_003242) is abolished, leading to 5 splicing events: intron inclusion, use of three cryptic sites (35nt exonic, 30nt and 972nt intronic), and exon 6 and 7 skipping. |
| | PBRM1 | <u>3:52682355C>G</u> (c.813+5G>C) | 6.8 > 2.9 (Natural Site) | KIRC | The natural donor of <i>PBRM1</i> exon 8 (NM_018313) is weakened, which leads to a significant increase in exon 8 skipping. |
| 1 | PBRM1 | 3:52685756A>G (c.714+2T>C) | 7.7 > 0.7 (Natural Site) | KIRC | The natural donor of <i>PBRM1</i> exon 7 (NM_018313) is abolished, resulting in a significant increase in exon skipping. |
| 7 | SETD2 | <u>3:47079269T>A</u> (c.7239-2A>T) | 9.8 > 2.1 6.4 >9.0 (<u>Natural</u> <u>Cryptic</u>) | KIRC | This mutation both significantly weakens the natural acceptor of <i>SETD2</i> exon 18 (NM_014159) and strengthens a 4nt exonic cryptic site. |
| n G | RB1 | <u>13:49027249T>A</u> (c.1814+2T>A) | 4.9 >-13.7 (Natural Site) | LUAD | The natural donor of <i>RB1</i> exon 18 (NM_000321) is abolished, leading to a significant increase in both exon skipping and intron inclusion. All intron inclusion reads contain the mutation of interest. |
| on | RBM10 | X:47006900G>T (c.17+3G>T) | 7.8 > 4.1 (Natural Site) | LUAD | The natural donor of <i>RBM10</i> exon 2 (NM_005676) is weakened, leading to a significant increase in exon 2 skipping. |
| | RBM10 | X:47028898G>T (c.201+1G>T) | 8.7 > -9.9 (Natural Site) | LUAD | <i>RBM10</i> exon 3 (NM_005676) natural donor is abolished. Reads which overlap the exon-intron junction are observed (with mutation). Use of a cryptic donor (61nt upstream; R_i =1.7 bits) is observed as well. |
| | DDX5 | <u>17:62500098</u> <u>TACAG>T</u> (c.441+2delACAG) | -1.3 > 5.4 (Cryptic Site) | PRAD | The mutation creates a 5.4 bit cryptic donor within <i>DDX5</i> exon 4 (NM_004396), which would lead to a 4nt deletion of exon 4. Note that wildtype splicing is still the dominant isoform observed. |
| | PTEN | <u>10:89690802G>A</u> (c.210-1G>A) | 8.5 > -2.3 (Natural Site) | PRAD | The natural acceptor of <i>PTEN</i> exon 5 (NM_000314) is abolished, leading to an increased amount of exon 5 skipping. |
| | NRAS | <u>1:115258669A>G</u> (c.111+2T>C) | 8.1 > 1.1 (Natural Site) | SKCM | The mutation abolishes the natural donor of <i>NRAS</i> exon 2 (NM_002524), which promotes a significant increase in exon 2 skipping |
| | PPP6C | <u>9:127933364C>T</u> (c.171G>A) | 6.7 > 3.7 (Natural Site) | SKCM | The mutation weakens <i>PPP6C</i> exon 2 (NM_002721) natural donor, leading to increased intron inclusion. All reads crossing the splice junction have the mutation. An intronic cryptic site is also activated (110nt downstream). |
| ich | РРР6С | <u>9:127923119C>G</u> (c.237+1G>C) | 6.8 > -11.8 (Natural Site) | SKCM | This mutation abolishes the natural donor of <i>PPP6C</i> exon 3 (NM_002721), resulting in a significant increase in exon 3 skipping. |
| CM R <i>T1</i> | BAP1 | <u>3:52442512T>C</u> | 1.9 > 5.1 | UVM | A pre-existing cryptic donor within <i>BAP1</i> exon 4 (NM_004656) is strengthened, leading to a significant increase in its use. This mutation |



Figure 1. Screenshot of ATM mutation in ValidSpliceMut. (A) Variant-specific and splice site-specific tabular results are presented under the headings "Splice Site Information" and "Variant Data". Results are organized by TCGA and ICGC sample IDs harboring the mutation within a series of expandable panels. Each panel consists of read counts and p-values by Veridical evidence type. Significant p-values (≤ 0.05) are bolded. (B) An integrative genome viewer (IGV) image showing alignment of expressed sequence reads (C) A dynamically generated histogram presents expression levels of all genes for a selected normal tissue type.

Methods

TCGA and ICGC data acquisition and processing: Controlled-access data was obtained with permission from the Data Access Committee at NIH for TCGA and from ICGC. Patient RNA sequencing BAM files (tumor and normal, when available) and their associated VCF files were obtained from the CancerGenomeHub and Genomic Data Commons. ICGC data was downloaded through the Score client (v1.5.0). Additional VCFs from patients with congenital diseases were obtained from the Baylor Hopkins Center for Mendelian Genomics study (BH CMG; <u>phs000711.v4.p1</u>). Variants which did not pass quality control were not analyzed

Information analysis and RNA-Seq validation of splicing variants: We used the Shannon Pipeline software (SP; applies IT to rapidly perform high-throughput, in silico prediction of variants impact on mRNA splicing) to analyze TCGA and ICGC variants (>168 million TCGA and >41 million ICGC variants) to evaluate their potential impact on splice site binding strength (change in information content, R_i, measured in bits). Veridical software analyzed the genomic variants by comparing the RNA-Seq alignment in the region surrounding the variant with the corresponding interval in control transcriptomes lacking the variant.³ Veridical counts abnormally spliced reads in RNA-Seq data (molecular phenotypes such as cryptic site use, exon skipping, or intron inclusion), and determines the null hypothesis probability that the transformed read count corresponds to normal splicing (p < 0.05 is considered significant). Veridical-flagged mutations found in genes in the Catalogue Of Somatic Mutations In Cancer (COSMIC) Cancer Gene Census (CGC)⁸ are highlighted, as are rare mutations also found in patients of congenital diseases from the BH CMG study (where literature supports the link of congenital disease to the affected gene)⁹. **Development of the ValidSpliceMut database and Beacon:** We created a publicly accessible Application Programming Interface (API) that can be utilized to programmatically query variants passing filter thresholds described above (https://beacon.cytognomix.com). It was built in accordance with the GA4GH Beacon v1.0.0 specification. We also developed the website ValidSpliceMut (Figure 1; below) to serve as a local interface to our Beacon, allowing users to manually search for a variant, by gene name or genome coordinate range.



| Evidence Type | Cryptic | Anti-Cryptic | Exon Skipping | Intron Inclusion | Intron Inclusion with Mutation |
|------------------------------|---------|--------------|---------------------------|------------------|--------------------------------|
| Junction spanning | 0 | 0 | 9 (p=0) | 2 (p=0.0784) | 1(p=0) |
| Read Abundance | 0 | 0 | 0 | 4 (p=0.7828) | 0 |
| chr3 52,682,300 bp | | 1 | 52,682 | 2,400 bp | 52,682,50 |
| [0 - 22] | | | | | |
| | • | | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | _ | |
| | | | | | |
| • • • • • | • • • | К К | F V Q S G P E N Y T K A I | NKALLDIDKAMAHI | SKYSGN |

PBRM1 mutation 3:52682355C>G (c.813+5G>C) weakens the natural donor of exon 7 (6.8 > 2.9 bits). which lead to a significant amount of increased exon skipping (red reads). Veridical detects 9 exon skipping reads over this region, and finds this number significant over controls ($p \approx 0$).

Table 1: Unique Flagged Variants from TCGA and ICGC Datasets

| | Unique TCGA | Unique ICGC Total | | Gene | Splice Mutation | |
|---|--|---|---|---|---|---------------|
| | Variants / Genes | Variants / Genes | lotal | | <u>17:73702087G>A</u> | |
| Total Variants (Passing QC) | > 168 million | > 41 million >209 million | | SAP30BP | (c.661-1G>A) | |
| Met IT-Based Criteria | 1,081,011 / 19,322 | / 19,322 24,101 / 8,151 1,094,749 / | | | (c.664G>A) | |
| Weaken Natural Sites | 373,770 | 11,662 | 380,852 | | | 3 |
| Strengthen Cryptic Sites | 744,848 | 12,547 | 752,472 | JUP | <u>17:39919236T>C</u> | |
| Significant by Veridical (p < 0.05) | 334,202 / 18,249 | 11,961 / 5,761 | 341,486 / 18,386 | | (0.1490A>0) | |
| Flagged for Increased Exon Skipping | 74,426 | 5,700 | 80,126 | | 1:173835706C>A | 8. |
| lagged for Increased Cryptic Site Use | 4,441 | 1,018 | 5,459 | GAS5 | (n.118-1G>T) | |
| % without associated rsID | 70.0% | 42.4% | 69.9% | | | |
| Shared Between Multiple Tissues (Novel variants only) | 9,810 | 246 | 10,125 | CAPN8 | <u>1:223815845C>A</u> (c.427-1G>T) | 4. |
| Shared Between Multiple Tissues Novel or <1% Average Heterozygosity) | 27,581 | 942 | 28,813 | | <u>11:67074511A>T</u> | 6. |
| Shared Molecular Phenotypes for ariants Identified in Multiple Tissues ^a | 20,900 (>50%) 15,793 (>67%) | 779 (>50%) 618 (>67%) | 21,887 (>50%) 16,574 (>67%) | SSH3 | (c.465-2A>T) | |
| Variants (novel or <1% average heterozygosity plicing category (i.e. junction spanning exon sk |) identified in multiple tu ipping) in more than half | mor types, flagged by V (>50%) or over two-thir | eridical for the same rds (>67%) of patients. | TKTL1 | X:153537696G>A (c.253-1G>A) | 7. |
| | Figure 2 Mutation TCGA Pa present in | . Census of Re ns Present in M atients. Predicted multiple tumors | current Splicing ultiple ICGC and splicing mutations that cause splicing | ELN | <u>7:73474706G>C</u> (c.1622G>C) | 6. |
| | abnormalit of validatic of the frac | ties were analyzed to on. Violin plots indic tion of predicted a | SPOP | <u>17:47699430C>A</u> (c.79-1G>T) | 9. | |
| ed shared | mutations present in multiple patients relative to the total number of tumours carrying those mutations in TCGA and ICGA. To achieve statistical | | | | <u>7:54825289T>G</u> (c6-2A>C) | 7. |
| | significance variants sh least 9 ICC | e (95% C.I.), dist ared by both datase GC (left) and 24 T(pared A higher ov | ributions of 1,379 ets and present in at CGA (right) patients | HLA-A | <u>6:29912835G>C</u> (c.1013-1G>C) | 8. |
| ICGC TCGA | mutations (average of | are validated in f 38.6% for ICGC and | the ICGC dataset 27.8% for TCGA). | Example i skipping, | nutations which sigr intron inclusion). Μι | nific utat |

| Evidence Type | Cryptic | Anti-Cryptic | Exon Skipping | Intron Inclusion | Intron Inclusion with Mutation |
|------------------------|---------|---------------|---------------|------------------|--------------------------------|
| lunction spanning | 0 | 0 | 4 (p=0.0032) | 2 (p=0.1033) | 1(p=0) |
| Read Abundance | 0 | 0 | 0 | 99 (p=0.0825) | 0 |
| chr10 69,651,100 bp | I | 69,651,200 bp | I | 69,651,300 bp | 69,651,400 bp |
| 0 - 18] | | | | | |
| | | | | | |
| | _ | | | | |
| | | | | • | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |

| SIRT1 mutation 10:69651312G>A (c.942G>A) weakens exon 4's natural donor (3.2 > 0.1 bits), whic |
|---|
| esulted in a significant amount of exon skipping (red reads; p = 0.032). This mutation, found in SKCN |
| patient TCGA-EE-A29V, was also found in BM CMG patient BH3973_1, with plantar lipomatosis. SIRT |
| nas been linked to soft tissue tumors such as liposarcoma. |

Table 2: Mutations Leading to Multiple Types of Aberrant Splicing

| Gene | Splice Mutation | <i>R_i</i> (bits) | Tumor | Observed Splicing Event |
|--------------------------|--|--|------------------------|---|
| SAP30BP | <u>17:73702087G>A</u> (c.661-1G>A) <u>17:73702091G>A</u> (c.664G>A) | 14.4>3.5 -3.0>2.2 (<u>Natural</u> <u>Cryptic</u>) | BLCA | Two closely situated mutations affect the natural acceptor of <i>SAP30BP</i> exon 10 (NM_013260). The first mutation abolishes the natural acceptor (increasing exon skipping and intron inclusion RNAseq reads), while the second creates a cryptic site 6nt downstream. |
| JUP | <u>17:39919236T>C</u> (c.1496A>G) | 3.1>0.6 0.5>3.7 (<u>Natural</u> <u>Cryptic</u>) | BLCA | Mutation within <i>JUP</i> exon 8 (NM_002230) simultaneously weakens a natural donor while creating a cryptic site. Exon skipping, intron inclusion and cryptic site use are all observed. |
| GAS5 | <u>1:173835706C>A</u> (n.118-1G>T) | 8.1>0.3 1.9>4.3 (<u>Natural</u> <u>Cryptic</u>) | KIRC | The natural acceptor of <i>GAS5</i> exon 4 (NR_002578) is abolished, simultaneously strengthens a cryptic site 9nt downstr. Cryptic site is seen, along with exon skipping and intron inclusion (likely due to the abolished natural site). |
| CAPN8 | <u>1:223815845C>A</u> (c.427-1G>T) | 4.0>-3.7 1.4>3.8 (<u>Natural</u> <u>Cryptic</u>) | LUAD | Simultaneously abolishes the natural acceptor of <i>CAPN8</i> exon 4 (NM_001143962) while strengthening a cryptic site 9nt downstr. Cryptic site use, exon skipping and intron inclusion is significantly increased in this patient. |
| SSH3 | <u>11:67074511A>T</u> (c.465-2A>T) | 6.4>-1.3 3.8>6.6 (<u>Natural</u> <u>Cryptic</u>) | LUAD | <i>SSH3</i> exon 5 natural acceptor (NM_017857) is abolished, while a pre-existing cryptic site 9nt downstr. is strengthened. Exon skipping, total intron inclusion and cryptic use are seen. |
| TKTL1 | <u>X:153537696G>A</u> (c.253-1G>A) | 7.0>-3.9 -2.5>2.7 (<u>Natural</u> <u>Cryptic</u>) | SKCM | The natural acceptor of <i>TKTL1</i> exon 3 (NM_012253) is abolished and a cryptic site 2nt downstream is created, leading to a deletion of 2nt. Cryptic site use, total exon skipping and intron inclusion are detected. |
| ELN | <u>7:73474706G>C</u> (c.1622G>C) | 6.6>4.7 -6.0>2.8 (<u>Natural</u> <u>Cryptic</u>) | SKCM | The variant weakens <i>ELN</i> exon 25 (NM_000501) natural site while creating a cryptic site 3nt away. The cryptic site is observed, as is intron inclusion, however exon skipping is the predominant aberrant splicing isoform detected. |
| SPOP | <u>17:47699430C>A</u> (c.79-1G>T) | 9.9>2.2 -0.4>2.0 (<u>Natural</u> <u>Cryptic</u>) | BLCA | The natural exon 4 acceptor of <i>SPOP</i> is weakened while an equivalently strengthened cryptic site 6nt downstr. is created. Cryptic site use, exon skipping and intron inclusion are seen. |
| SEC61G | <u>7:54825289T>G</u> (c6-2A>C) | 7.7>-10.9 1.4>3.5 (<u>Natural</u> <u>Cryptic</u>) | COAD | The natural acceptor of exon 2 (NM_014302) is abolished and a cryptic acceptor is created. Significant exon skipping, cryptic site use, and intron inclusion is detected by Veridical. |
| HLA-A | <u>6:29912835G>C</u> (c.1013-1G>C) | 8.6>-3.1 1.9>4.3 (<u>Natural</u> <u>Cryptic</u>) | CESC | The natural acceptor of exon 6 of <i>HLA-A</i> (M_002116) is abolished by this variant, and a pre-existing cryptic site 4nt downstream is strengthened. Exon skipping, total intron inclusion, and cryptic site use are detected. |
| Example n skipping, i | nutations which sign ntron inclusion). Mu | ificantly alter splicinations are linked to | ng in all o their p | ways measured by Veridical (i.e. cryptic splice site use, exon bage on https://validsplicemut.cytognomix.com/ . |

Example mutations which alter splicing in tumor-associated genes found in patients with these tumor types. Mutations are hyperlinked to their ValidSpliceMut Beacon page, which provides additional material such as IGV images of the RNAseq evidence for the regions of interest. GRCh37 coordinates are indicated

Conclusions

We present the Validated Splicing Mutation Beacon and web resource (<u>http://validsplicemut.cytognomix.com/</u>) that consists of 341,486 unique mutations from TCGA and ICGC tumor patients predicted to alter splicing with RNA-Seq support. The majority of these mutations (69.9%) are not present in dbSNP 150. In total, 543 Tier 1 COSMIC CGC genes exhibited at least one Veridical-flagged variant present in the ValidSpliceMut database. This resource should significantly contribute to reducing the number of outstanding VUS in tumor (and possibly some germline) genomes, and substantially increases the number of functional variants with previously unappreciated consequences to mRNA splicing, in particular, those which activate cryptic splice sites.

References

1. Foley SB, Rios JJ, Mgbemena VE, et al. Use of Whole Genome Sequencing for Diagnosis and Discovery in the Cancer Genetics Clinic. EBioMedicine. 2015; 2(1):74–81. 26023681. 2. Caminsky N, Mucaki EJ, Rogan PK. Interpretation of mRNA splicing mutations in genetic disease: review of the literature and guidelines for informationtheoretical analysis [version 1; referees: 2 approved]. F1000Res. 2014; 3:282. 25717368. 3. Viner C, Dorman SN, Shirley BC, et al. Validation of predicted mRNA splicing mutations using high-throughput transcriptome data [version] 2; referees: 4 approved]. F1000Res. 2014; 3:8. 4. Mucaki EJ, Caminsky NG, Perri AM, et al. A unified analytic framework for prioritization of non-coding variants of uncertain significance in heritable breast and ovarian cancer. BMC Med Genomics. 2016; 9:19. 5. Hoadley KA, Yau C, Hinoue T, et al. Cell-of-Origin Patterns Dominate the Molecular Classification of 10,000 Tumors from 33 Types of Cancer. Cell. 2018; 173(2):291–304.e6. 6. Shirley BC, Mucaki EJ, Rogan PK. Pan-cancer repository of validated natural and cryptic mRNA splicing mutations [version 2; peer review: 1 approved, 1 approved with reservations]. F1000Research. 2019; 7:1908. 7. Global Alliance for Genomics and Health. GENOMICS: A federated ecosystem for sharing genomic, clinical data. Science. 2016; 352(6291):1278–1280. 8. Tate JG, Bamford S, Jubb HC, et al. COSMIC: the Catalogue Of Somatic Mutations In Cancer. Nucleic Acids Research. 2019; 8;47(D1):D941-D947. 9. Bamshad MJ, Shendure JA, Valle D et al. The Centers for Mendelian Genomics: a new large-scale initiative to identify the genes underlying rare Mendelian conditions. Am J *Med Genet A.* 2012; 158A(7):1523-5.

