

# Splicing mutation analysis reveals previously unrecognized pathways in lymph node-invasive breast cancer.

Stephanie N. Dorman<sup>1</sup>, Coby Viner<sup>2</sup>, Peter K. Rogan<sup>1,2\*</sup>

Departments of Biochemistry<sup>1</sup> and Computer Science<sup>2</sup>, Western University  
\*progan@uwo.ca

## Supplemental Materials

---

<b>I. SomaticSniper Supplementary Materials</b>	<b>Page</b>
– Supplementary Method – Variant Calling Methods	2
– Supplementary Results – SomaticSniper Variant Calling Results	2-3
– Supplementary Data 1 – Variant Summaries by Mutation Type	3
– Supplementary Data 2 – SomaticSniper Variants Compared to TCGA	3
<b>II. Supplemental Figures</b>	
– Supplementary Figure S1 – RNA-Seq Coverage Heatmap by Subtype	4
– Supplementary Figure S2 – Information theory based analysis and corresponding evidence demonstrating abnormal mRNA splicing in predicted mRNA splicing mutations	5-6
– Supplementary Figure S3 – Junction-spanning, cryptic splicing read counts for GATA3 mutation	7
– Supplementary Figure S4 – Intron Inclusion in tumour and normal breast genomes, based on RNA-Seq evidence.	8
– Supplementary Figure S5 – Word Clouds of Overrepresented Pathways by Subtype	9-12
– Supplementary Figure S6 – Flowchart indicating procedure for filtering splicing mutation variants.	13

See page 14 for reference list.

### III. Supplementary Tables (additional xls files)

- Supplementary Table S1 – TCGA Tumour-Normal Pairs Analyzed
- Supplementary Table S2 – Single nucleotide variants annotated by ANNOVAR for codon changes
- Supplementary Table S3 – Insertions/deletions annotated by ANNOVAR for codon changes
- Supplementary Table S4 – Splicing single nucleotide variants predicted by the Shannon Pipeline
- Supplementary Table S5 – Variants Compared to those Previously Published by TCGA
- Supplementary Table S6 – Overrepresentation analysis of TCGA mutations missed by Strelka
- Supplementary Table S7 – MuSiC Results Compared to Significantly Mutated Genes
- Supplementary Table S8 – Protein Coding Mutation Pathway Analysis
- Supplementary Table S9 – Splicing Mutation Pathway Analysis
- Supplementary Table S10 – Pathways Overrepresented by Protein Coding and Splicing Mutations
- Supplementary Table S11 – Pathways Overrepresented by Every Splicing Mutation Type
- Supplementary Table S12 – Pathways Overrepresented by Splicing Mutations in LN- Tumours
- Supplementary Table S13 – Pathways Overrepresented by Splicing Mutations in LN+ Tumours
- Supplementary Table S14 – Comparing Pathways Overrepresented between LN- and LN+ Tumour Mutations
- Supplementary Table S15 – Comparing Grouped Pathways Overrepresented in LN- and LN+ Tumour Mutations
- Supplementary Table S16 – Pathway Analysis of Mutated Genes Unique to LN+ Tumours
- Supplementary Table S17 – Pathway Analysis of Mutated Genes Unique to LN- Tumours
- Supplementary Table S18 – Pathway Analysis of Deleterious Mutations in LN- and LN+ Tumours
- Supplementary Table S19 – Frequency of Mutations in *NCAM1* Pathway Genes
- Supplementary Table S20 – Number of Mutations by Subtype
- Supplementary Table S21 – Pathway Analysis of Mutations by Subtype and Lymph Node Status

## SomaticSniper Supplementary Materials

### Supplementary Methods – Variant Calling Methods

Two independent variant callers, Strelka<sup>1</sup> and SomaticSniper<sup>2</sup>, were evaluated. The main analysis performed using results from Strelka, which has greater sensitivity and ability to detect subclonal mutations, by minimizing reporting of spurious variants and germline polymorphisms<sup>3</sup>. Additionally, the SomaticSniper methods and results are reported below.

Before running SomaticSniper, all DNA sequencing BAM files were realigned using the Genome Analysis Toolkit (GATK) Indel Realigner program<sup>4</sup>. In addition to default parameters, the knownAlleles parameter was used with the well-documented insertions/deletions (indels) files: Mills\_and\_1000G\_gold\_standard.indels.b37.sites.vcf<sup>5</sup> and 1000G\_phase1.indels.b37.vcf<sup>6</sup>, available through the bioinformatic resource Galaxy<sup>7,8</sup>. SomaticSniper data was then post-processed to only include variants with both mapping and somatic qualities of at least 40 (equivalent to running it with `-Q 40 -q 40`).

### Supplementary Results – SomaticSniper Variant Calling Results

SomaticSniper variant predictions are summarized in Table 1. Notably, there were 1,208 variants from SomaticSniper that are predicted to affect both protein coding and splicing 594 genes. In the SomaticSniper data, mutations classified as both protein coding and splicing variants were found in 383 tumours, with 63 of these variants in *PASD1*, 61 in *PRSS3*, 52 in *NF1*. The variants in these genes, as well as others that were highly mutated, are the exact same genomic location and nucleotide change, suggesting that SomaticSniper reported higher numbers of SNPs<sup>3</sup> that were not annotated with dbSNP135 in >1% of the population, which was used to filter out common SNPs. There were 248 variants in 186 tumours from the SomaticSniper set that were classified as silent amino acid changes from ANNOVAR, but were revealed to affect splicing from the Shannon Pipeline predictions.

There was relatively low concordance between the two variant callers, which reported variant lists with less than 50% similarity. There were 21,112 protein coding and 1,811 splicing variants common to both Strelka and SomaticSniper. The predicted variants were compared to the previously reported TCGA Level 2 somatic mutations (Table 2). Strelka showed the highest concordance with TCGA mutations, reporting 82.1% of protein coding mutations, and 86.5% of the splicing variants. Conversely, SomaticSniper predicted 73.4% protein coding and 75.3% splicing variants reported by TCGA.

Both of the somatic variant callers we employed utilize Bayesian methods to elucidate somatic event probabilities. Strelka and SomaticSniper were found to be the two best variant callers in a comparison by Roberts et al 2013. Additionally, these two are a valuable combination, in that SomaticSniper is useful to generate “a variety of candidate SNV sites without any particular drawbacks”, although with a fair amount of false positives, while Strelka is least prone to returning germ-line polymorphisms. The relative stringency of Strelka was our main reason for performing most of our analyses with it, along with the fact that many of its candidates (at probability 0.2) were also returned by other callers. It is worth mentioning that different callers have been found to have poor correlations at the same sites; in particular, Strelka and SomaticSniper were found to have a 0.21 Pearson correlation coefficient in the abovementioned study. Our use of Veridical to validate splicing variants with functional evidence of the mutation significantly resolves the inconsistency between somatic variant callers (for this type of mutation).

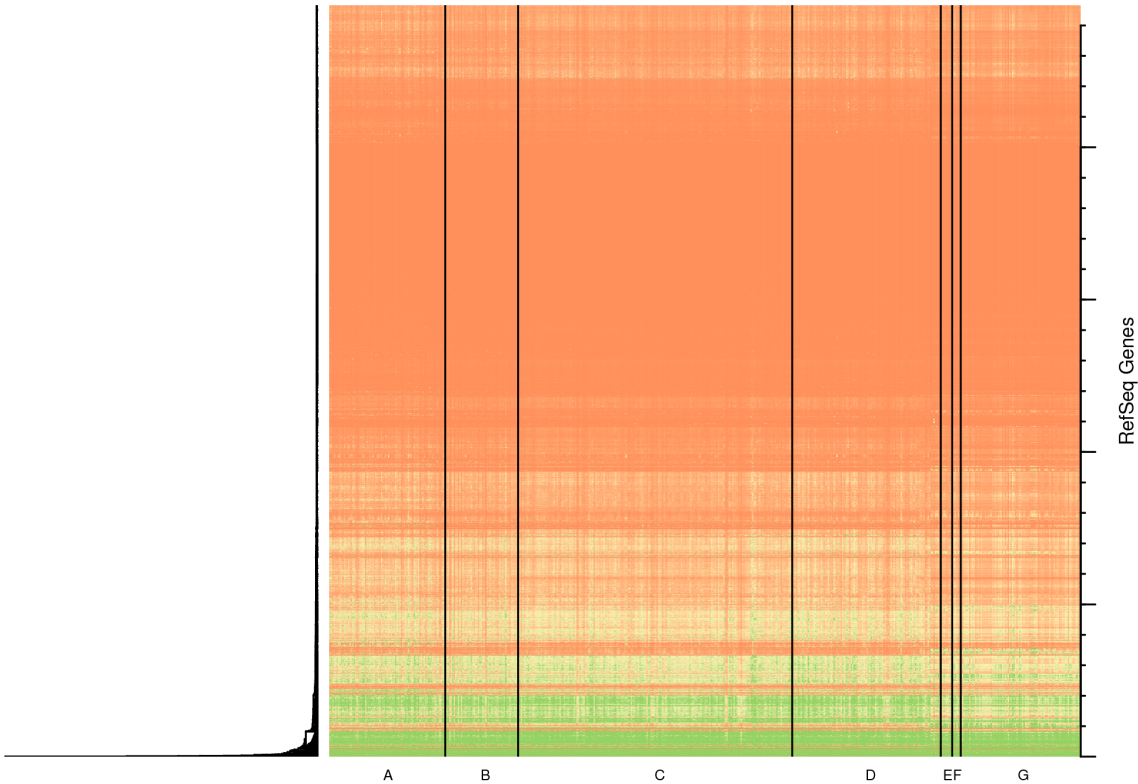
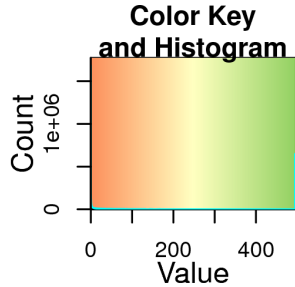
#### Supplementary Data 1 | Variant Summaries by Mutation Type

	<b>SomaticSniper software</b>
<b><i>ANNOVAR protein coding variants</i></b>	
Synonymous	23,458
Nonsynonymous	52,634
Stop gain or loss	2,127
<b>Total protein coding variants</b>	<b>78,219</b>
<b><i>Shannon Pipeline splicing variants</i></b>	
Cryptic	6,441
Inactivating	2,685
Leaky	10,648
<b>Total splicing variants</b>	<b>19,774</b>
Synonymous	248
Nonsynonymous	905
Stop gain or loss	55
<b>Total</b>	<b>1,208</b>
<i>% Synonymous also splicing</i>	1.0572%
<i>% Nonsynonymous also splicing</i>	1.7194%
<i>% Stop gain or loss also splicing</i>	2.5858%

#### Supplementary Data 2 | SomaticSniper Variants compared to TCGA Findings

	<b>Total TCGA</b>	<b>TCGA predicted by SomaticSniper</b>
<b>TCGA Protein Coding Variants</b>		
SNVs Validated	5,557	4,365 (77.3%)
SNVs Not Validated	18,197	13,380 (72.2%)
Indels Validated	125	N/A
Indels Not Validated	1,758	N/A
<b>Total</b>	<b>25,637</b>	<b>17,745 (73.4%)</b>
<b>TCGA Splicing Variants</b>		
SNVs Validated	87	70 (80.5%)
SNVs Not Validated	342	253 (74.0%)
<b>Total</b>	<b>429</b>	<b>323 (75.3%)</b>

**Supplementary Figure S1. RNA-Seq Coverage Heat Map by Subtype.** Heatmap depicting coverage per exonic base of TCGA RNA-Seq tumour and normal data. Expression based on RNA-Seq datasets is shown along the x-axis, with tumours first, ordered by subtype, followed by matched normal breast tissues. These categories are demarcated within the heatmap by black vertical lines, which correspond to the sample types: (A) basal-like; (B) HER2-enriched; (C) luminal A; (D) luminal B; (E) tumour, subtype not available; (F) normal-like tumor; and (G) normal control samples. The y-axis consists of all RefSeq genes (with major and minor tick marks every 5,000 and 1,000 genes, respectively), clustered to form a dendrogram, which is visible on the left side of the graph. Genes with low nominal expression levels were below minimum threshold read counts for analysis by Veridical.



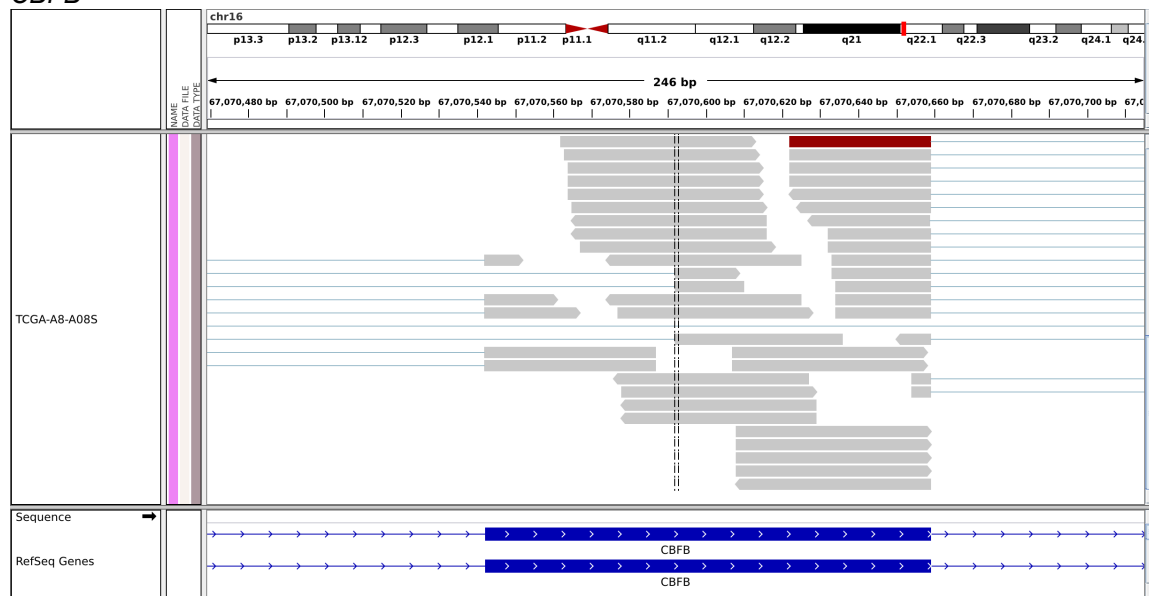
**Supplementary Figure S2. Information theory based analysis and corresponding evidence demonstrating abnormal mRNA splicing in predicted mRNA splicing mutations.** (A) Table indicates the TCGA sample identifier, variant, information analysis and statistical support for the mutation. (B) Screenshots from the Integrative Genomics Viewer (IGV) displaying junction-spanning reads that demonstrate cryptic splicing for mutations predicted by the Shannon Pipeline in the genes *CBFB*, *GATA3*, *PALB2*, and *ABL1*. The normal exonic structure is indicated by blue, with the thick bars representing exons, and the thin lines introns. RNA-Seq reads are shown in grey with the vertical dotted black lines demarcate the location of the cryptic splice site.

(A)

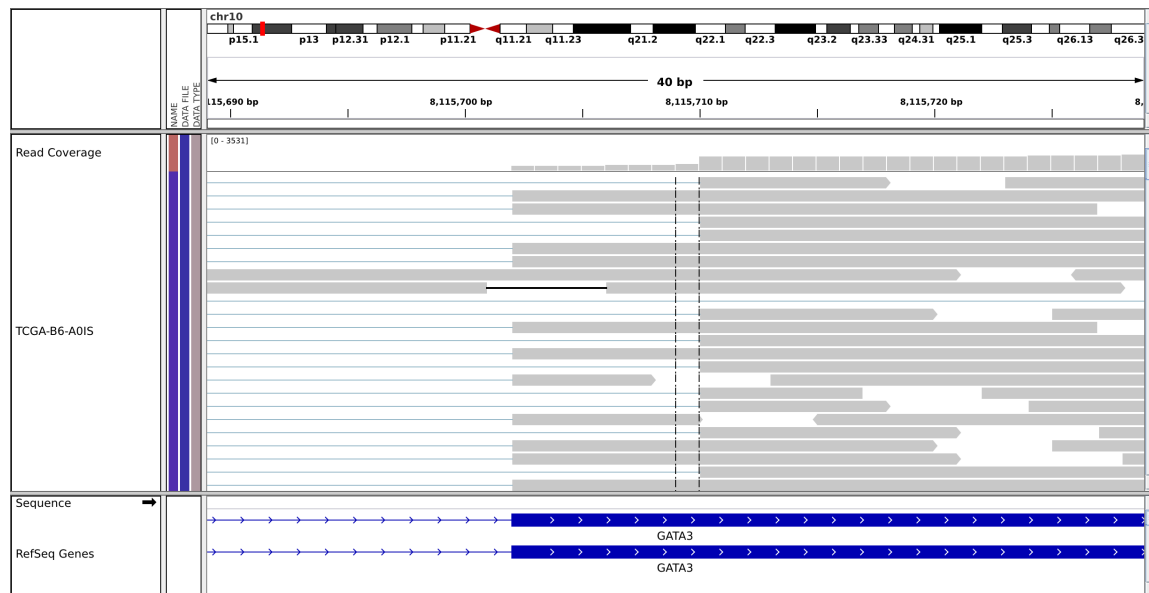
Sample	Gene	Splice Site Coordinate	Variant Coordinate	Ref/Var	$R_i$ -initial	$R_i$ -final	$\Delta R_i$	Cryptic Site Use P-Value	Exon Skipping P-Value
TCGA-A8-A08S	CBFB	chr16:67070591	chr16:67070577	G/T	5.6	7.5	1.9	< 0.005	0.12
TCGA-B6-A015	GATA3	chr10:8115709	chr10:8115702	A/C	4.2	5.9	1.7	< 0.005	NA
TCGA-B6-A0RT	PALB2	chr16:23637694	chr16:23637710	T/A	5.3	7.0	1.7	< 0.005	0.05
TCGA-B6-A0RV	ABL1	chr9:133750256	chr9:133750254	G/C	0.8	9.6	8.8	< 0.005	NA

(B)

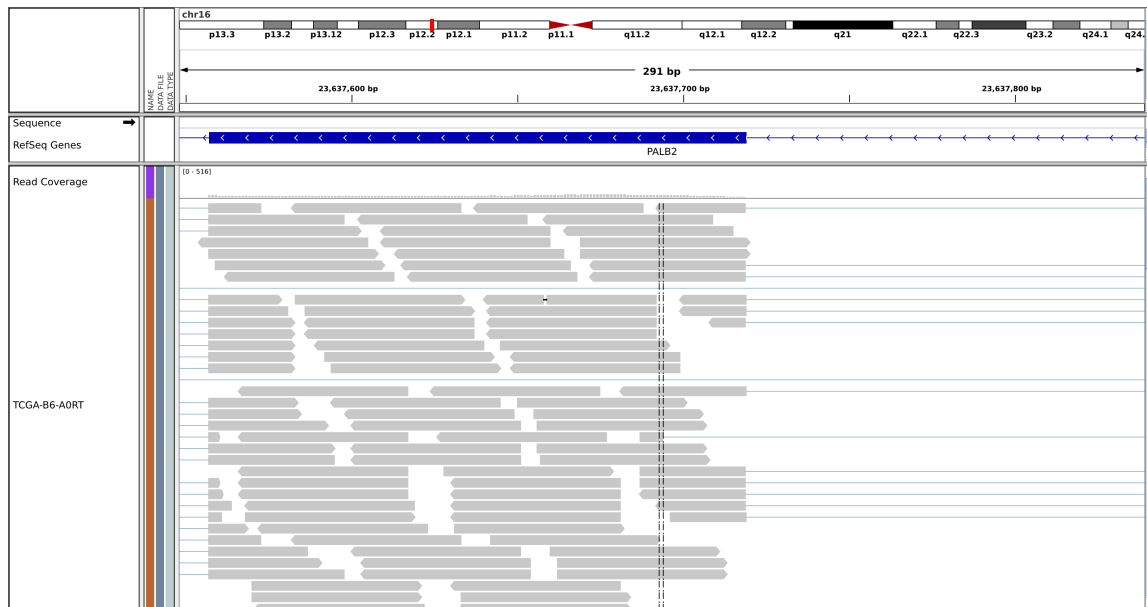
*CBFB*



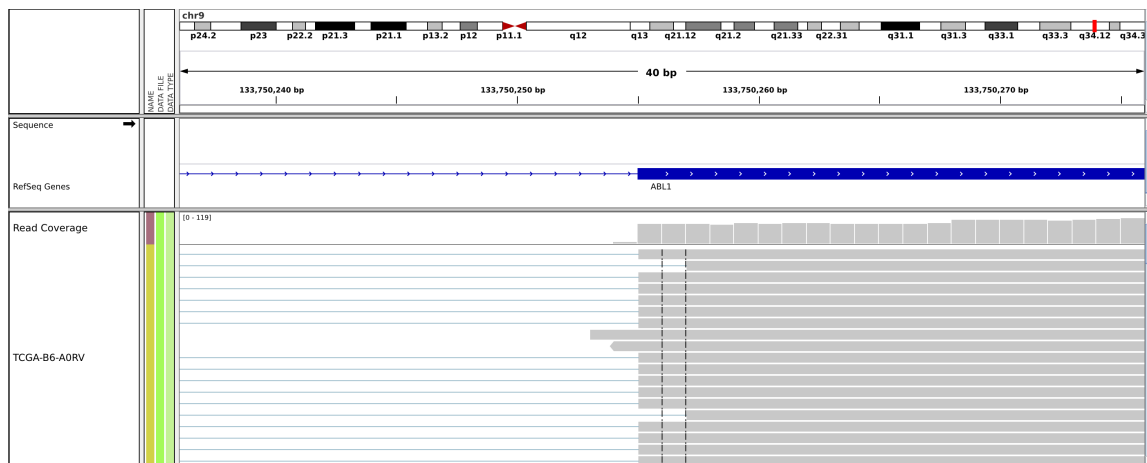
*GATA3*



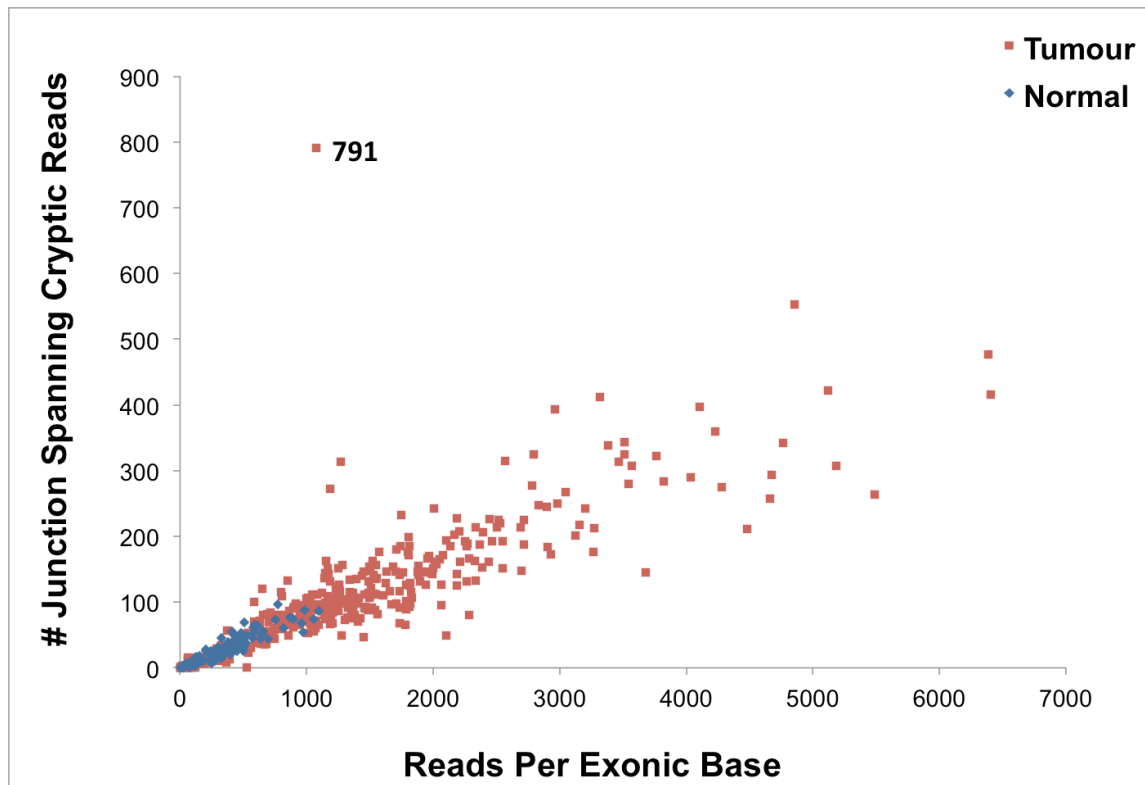
# PALB2



# ABL1

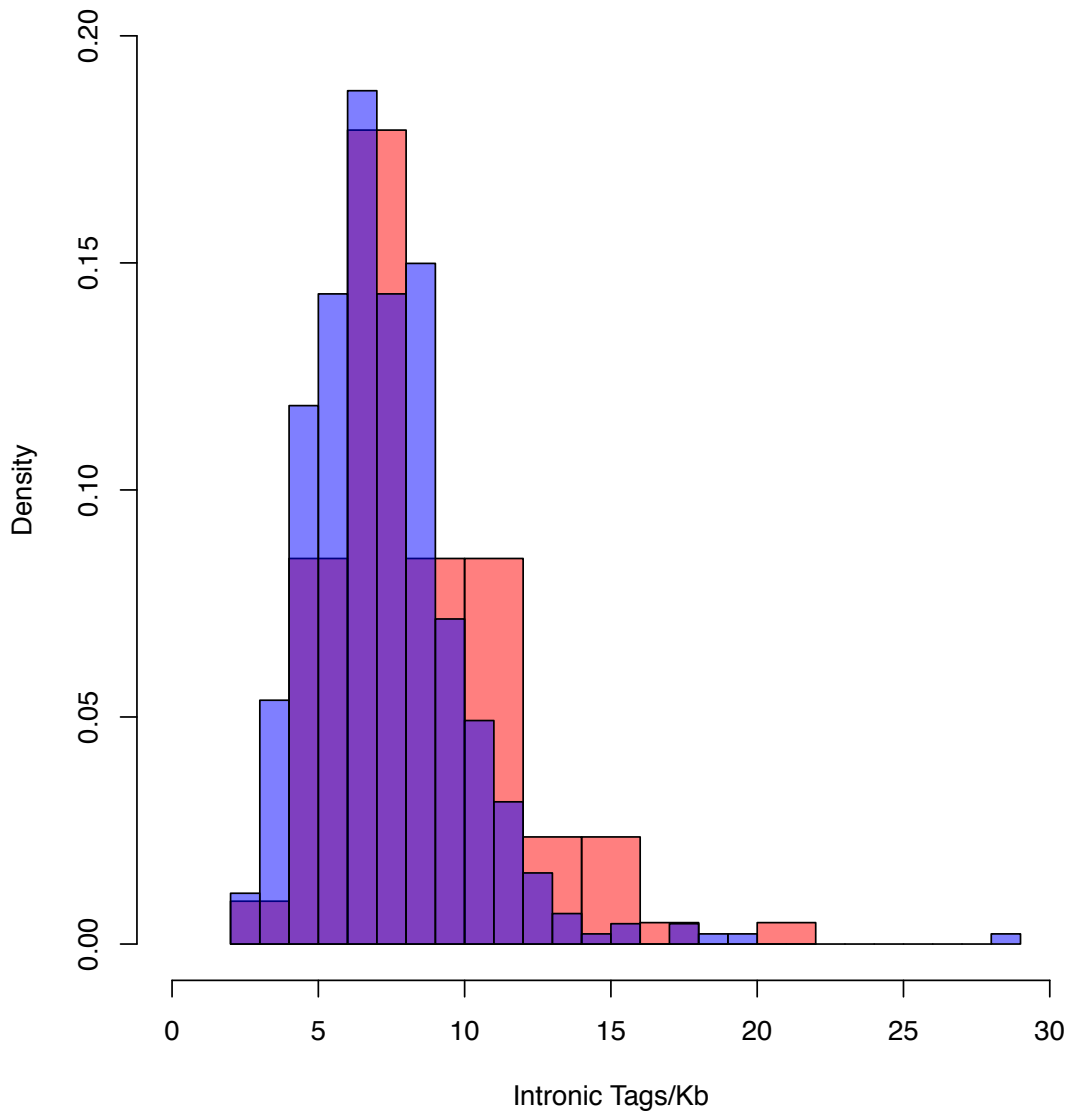


**Supplementary Figure S3. Junction-spanning, cryptic splicing read counts for GATA3 mutation (chr10: g.8115702A>C).** The number of RNA-Seq reads per exonic base were plotted against the number of reads demonstrating GATA3 cryptic splicing in the variant-containing tumours and controls. The variant containing tumour is indicated by the number of cryptic splicing reads (n = 791), tumours that do not contain this variant are in red, and normal controls are in blue. Cryptic splicing in the control samples likely occurs because the cryptic splice site ( $R_i = 4.2$  bits) exceeds the strength of the natural splice site ( $R_i = 0.9$  bits). However, the mutation further weakens the natural splice site (final  $R_i = 0.0$  bits), while simultaneously strengthening the cryptic splice site (final  $R_i = 5.8$  bits), which consistent with the RNA-Seq analysis.



**Supplementary Figure S4. Intron Inclusion in tumour and normal breast genomes, based on RNA-Seq evidence.** Histogram of the density of intronic sequence reads for normal (blue) and tumour (red) RNA-Seq samples. Purple shading represents overlapping components of the two density distributions. Intron inclusion was calculated with RSeQC's ReadDist script and RefSeq's gene annotation.

High levels of unspliced isoforms with intron inclusion were the most frequent outcome of mutations with significant effects on mRNA splicing. Nevertheless, when considering non-specific aberrant splicing across the transcriptome, the numbers of junction-spanning, intron inclusion reads present in normal and tumour samples did not significantly differ ( $p > 0.1$ ). In fact, non-junction-spanning, intronic read-abundance reads of normal controls exceeded those of the tumour samples ( $p < 0.01$ ; Supplementary Figure 4). This suggests that validation events in these tumour samples are not due solely to intron inclusion and aberrant mRNA splicing known to be present in breast tumours<sup>9</sup>. It is notable, however, that the levels of intronic inclusion for validated mutations significantly exceeded the read counts for all controls that did not contain these variants.





**Supplementary Figure S5. Word Clouds of Overrepresented Pathways by Subtype.** Word Clouds of overrepresented Reactome pathways for mutations in breast tumours, stratified by lymph node status (positive or negative) and by breast cancer subtype (basal-like, *HER2*-enriched, Luminal A, or Luminal B). The size of each word is proportional to its frequency in the abstracted list of overrepresented pathways.

**Basal-like lymph node positive**



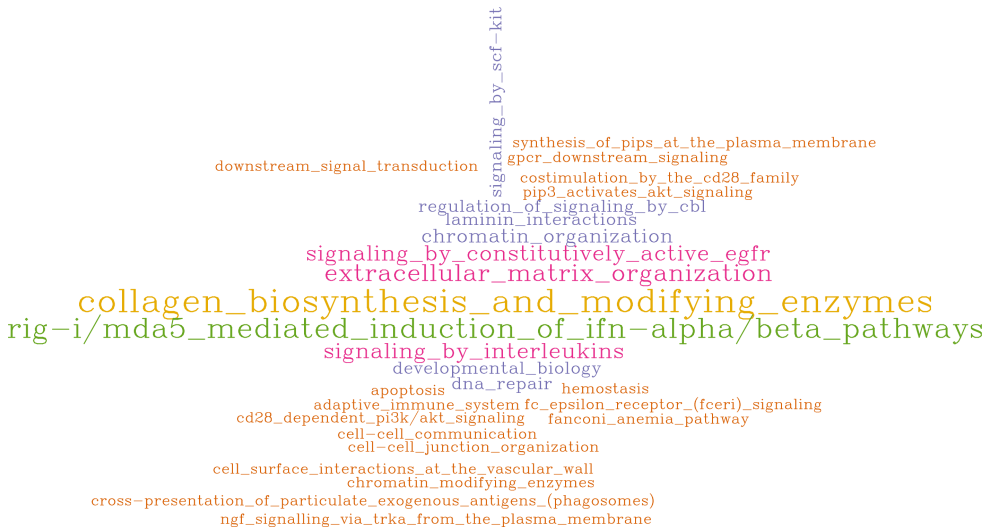
**Basal-like lymph node negative**



## HER2-enriched lymph node positive



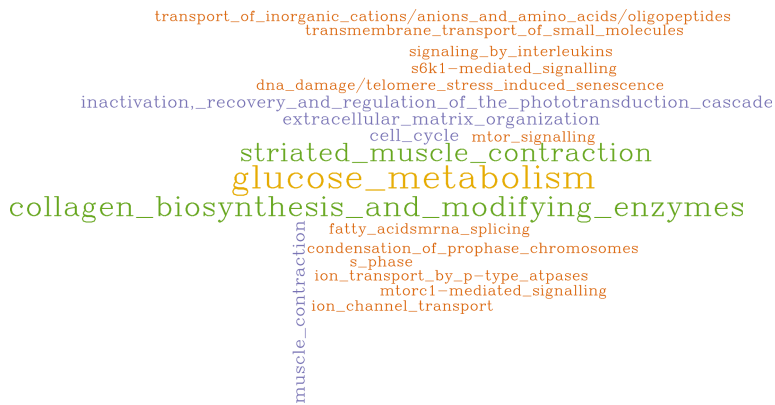
## HER2-enriched lymph node negative



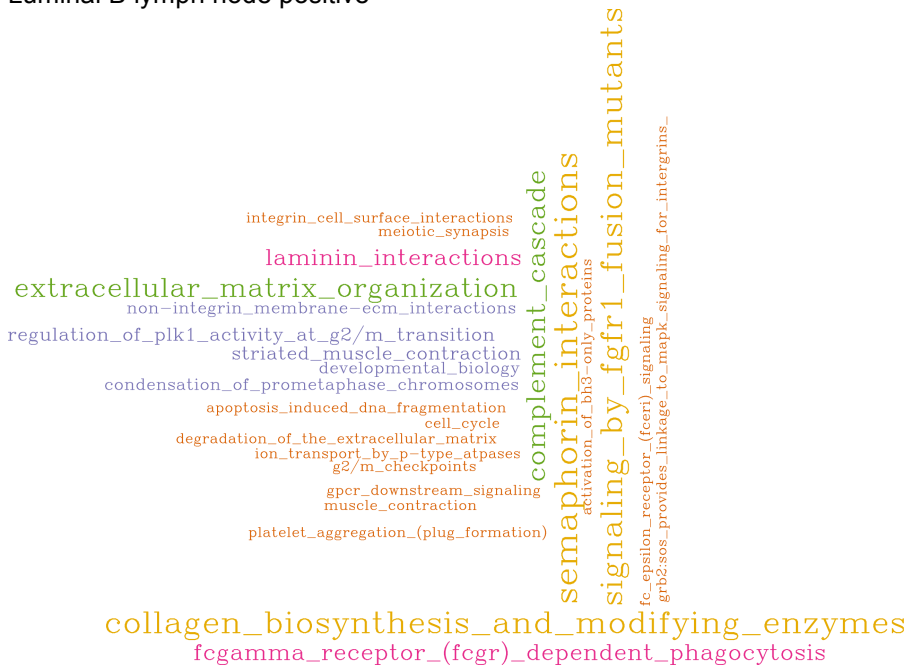
## Luminal A lymph node positive



## Luminal A lymph node negative



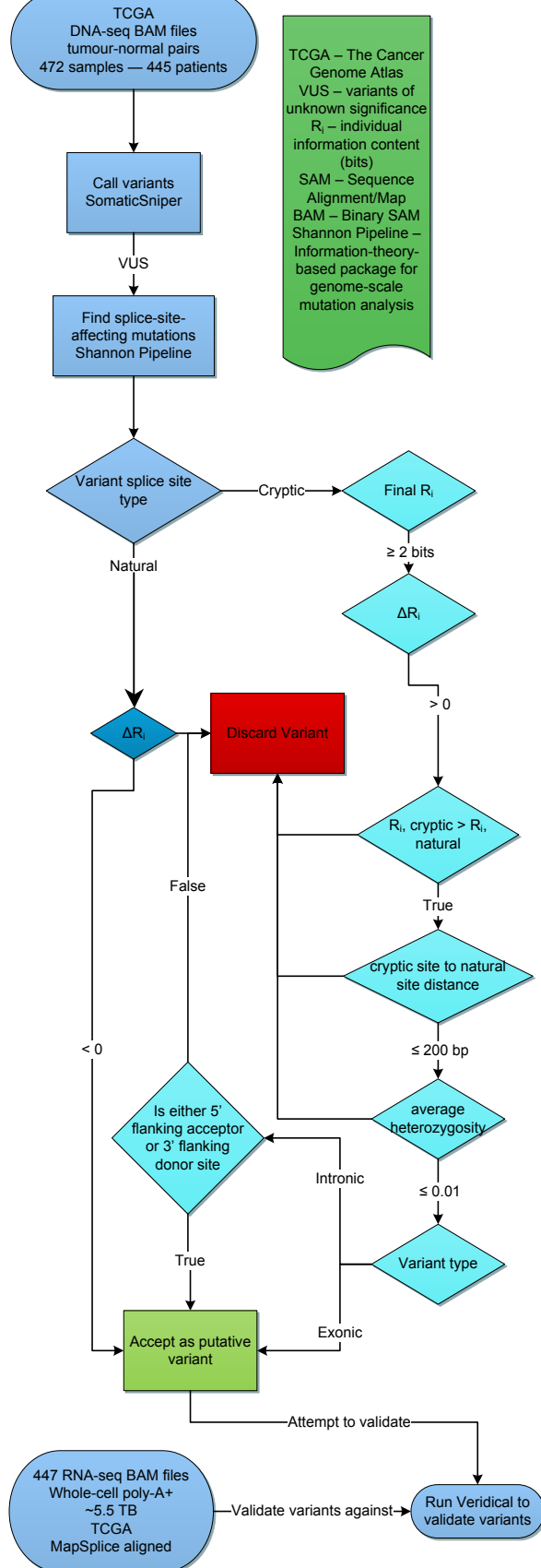
Luminal B lymph node positive



Luminal B lymph node negative



### Supplementary Figure S6



**Supplementary Figure S6. Flowchart indicating procedure for filtering splicing mutation variants.** Shannon pipeline splicing variants output was filtered using the steps shown in this flowchart to identify those variants that are likely to cause aberrant splicing. Upon identifying variants with Strelka (or Somatic Sniper), the VCF files were submitted to the Shannon splicing mutation pipeline, then categorized as either mutations affecting natural splice sites (3' acceptor, or 5' donor) or cryptic splice site strengths. In a small number of cases, both natural and cryptic splice sites were simultaneously altered. Natural sites that were predicted to be abolished were further considered. Predicted leaky splicing mutations were excluded from the present analysis, since the validation methods for such mutations has not yet been assessed. Aside from standard information theory-based mutation criteria, cryptic splicing mutation candidates were also filtered for proximity to the nearest neighboring natural splice site and population frequency. The filtered variant subset ( $n = 5,206$ ) was used for all subsequent analyses.

## References

1. Saunders, C. T. *et al.* Strelka: Accurate somatic small-variant calling from sequenced tumor-normal sample pairs. *Bioinformatics* **28**, 1811-1817 (2012).
2. Larson, D. E. *et al.* Somaticsniper: Identification of somatic point mutations in whole genome sequencing data. *Bioinformatics* **28**, 311-317 (2012).
3. Roberts, N. D. *et al.* A comparative analysis of algorithms for somatic SNV detection in cancer. *Bioinformatics* **29**, 2223-2230 (2013).
4. McKenna, A. *et al.* The genome analysis toolkit: A MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res.* **20**, 1297-1303 (2010).
5.  
[https://usegalaxy.org/library\\_common/ldda\\_info?library\\_id=f9ba60baa2e6ba6d&show\\_deleted=False&controller=library&folder\\_id=709338dd4741c085&use\\_panels=False&id=a78cd737488dc2c1](https://usegalaxy.org/library_common/ldda_info?library_id=f9ba60baa2e6ba6d&show_deleted=False&controller=library&folder_id=709338dd4741c085&use_panels=False&id=a78cd737488dc2c1).
6.  
[https://main.g2.bx.psu.edu/library\\_common/ldda\\_info?library\\_id=f9ba60baa2e6ba6d&show\\_deleted=False&controller=library&folder\\_id=709338dd4741c085&use\\_panels=False&id=93fcddc692e0986f#](https://main.g2.bx.psu.edu/library_common/ldda_info?library_id=f9ba60baa2e6ba6d&show_deleted=False&controller=library&folder_id=709338dd4741c085&use_panels=False&id=93fcddc692e0986f#).
7. Giardine, B. *et al.* Galaxy: a platform for interactive large-scale genome analysis. *Genome Res.* **15**, 1451-1455 (2005).
8. Blankenberg, D. *et al.* Galaxy: a web-based genome analysis tool for experimentalists. *Curr. Protoc. Mol. Biol.* **Chapter 19**, Unit 19.10.1-21 (2010).
9. Eswaran, J. *et al.* RNA sequencing of cancer reveals novel splicing alterations. *Scientific Reports* **3** (2013).