# Finding subtle mutations with the Shannon human mRNA splicing pipeline

Presentation at the CLC bio Medical Genomics Workshop
American Society of Human Genetics Annual Meeting
November 9, 2012

Peter K Rogan
Western University, London, Ontario   Canada
progan@uwo.ca
info@cytognomix.com

# Motivation

The Shannon pipeline created to address the vexing problem of assessing the many variants of unknown significance (VUS) that are detected in genetic testing, exome, and complete genome sequencing. Which are deleterious and which are benign?

Besides translational effects and modifications, variants may affect promoter regulation, mRNA splicing, long range effects.

Shannon pipeline reanalysis of the Breast Cancer Information Core identified 299 novel mRNA splicing mutations (Mucaki et al. Hum. Mutation 32:735-42, 2011).
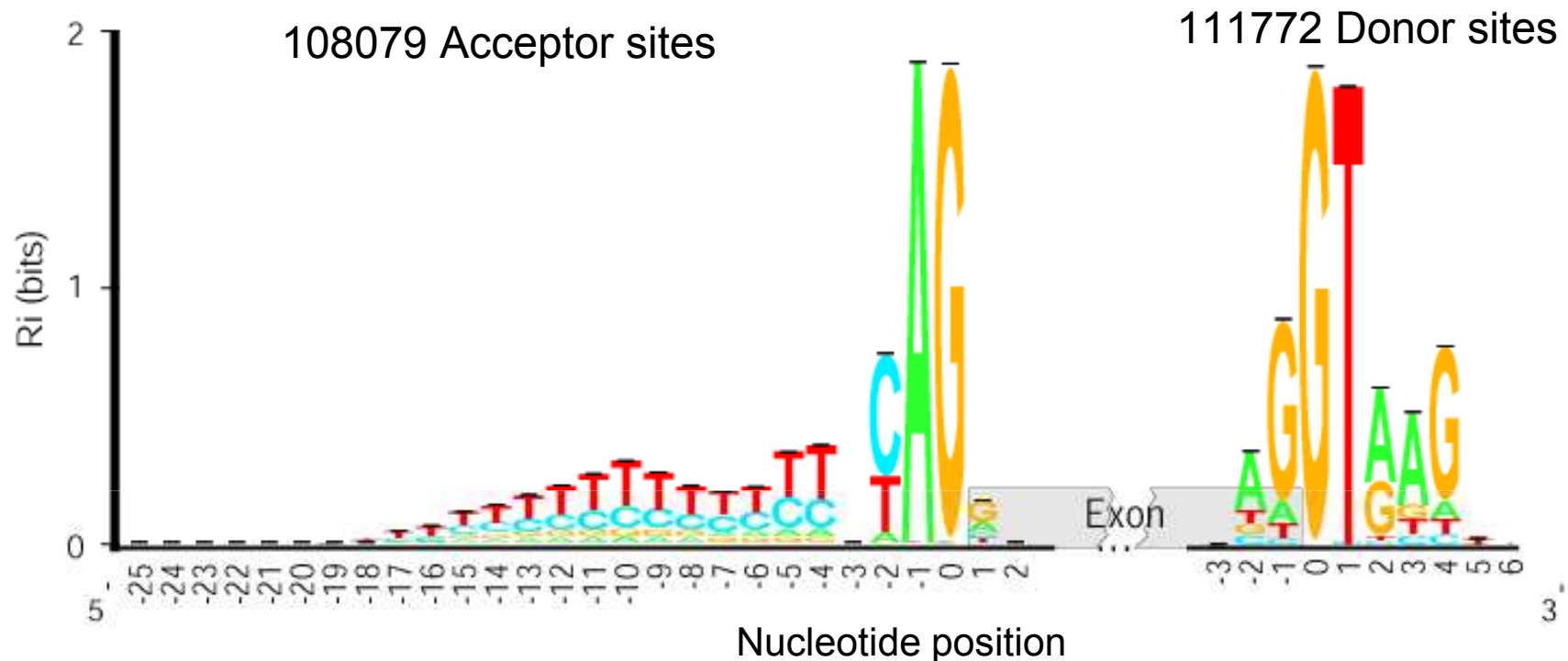
Our goal is to produce software capable of any large-scale, non-coding VUS analysis - regardless of disease, trait or phenotype.

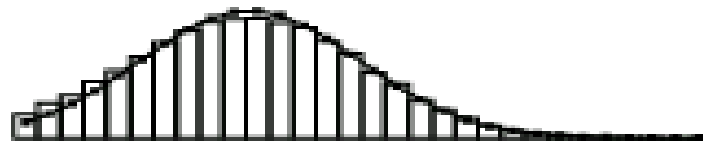# Splicing signals *in vivo* and *in silico*

- Like most nucleic acid binding sites, sequences of splice donor, acceptor, and regulatory sites vary

- The splicing machinery recognizes and processes these combinations of splicing signals.

- **Information theory** provides a *thermodynamic* framework to recognize members of a group of structurally- and functionally-related, variable sequences.

- It models sequence variability inherent in these signals, then predicts which sequence variants are deleterious.

# Human splice junction models
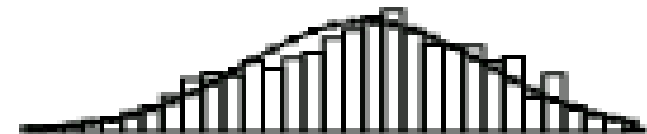## (refined from finished human reference sequence)



Rogan et al. Pharmacogenetics **13**:207-213, 2003

# Effects of mutations on mRNA splicing

# Molecular information theory

The information gained when a nucleic acid is recognized by a molecular machine (ie. the splicesome) is accompanied by a decrease in entropy that occurs upon binding:

$$\Delta S = -k_B \ln(2) R.$$

Second law relates information to the enthalpy of the molecular machine:

constant ➜ $-k_B T \ln 2 < \_q / R_i$     (joules/bit)
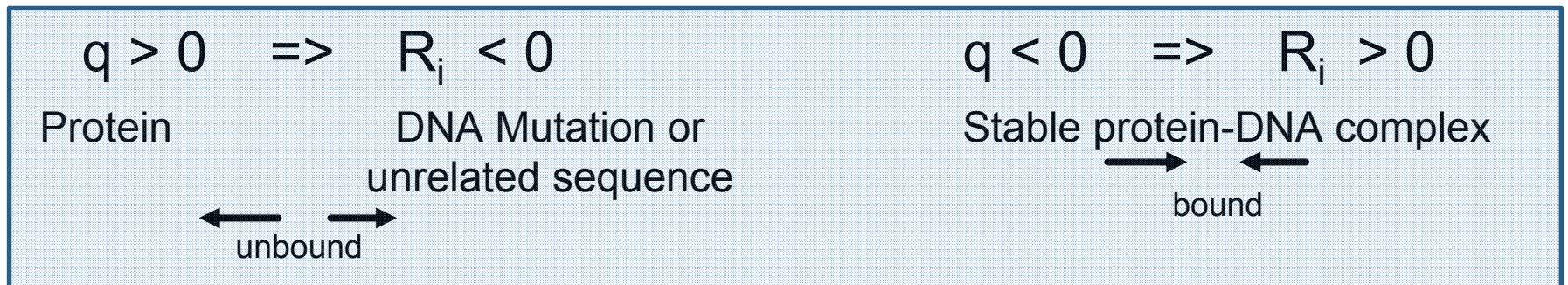
q: work; T: temperature; $R_i$: individual information

$q > 0 \quad => \quad R_i < 0$    $q < 0 \quad => \quad R_i > 0$

Protein                DNA Mutation or           Stable protein-DNA complex
                       unrelated sequence

⟵    ⟶                                    ⟶   ⟵

    unbound                                       bound

Fold change in affinity $\leq 100/2^{\Delta Ri}$

# Mild (or leaky) splicing mutation



| **ΔR$_i$** (bits) | **Fold** | **%** |
|---|---|---|
| 3.1 | 8.6 | 11 |

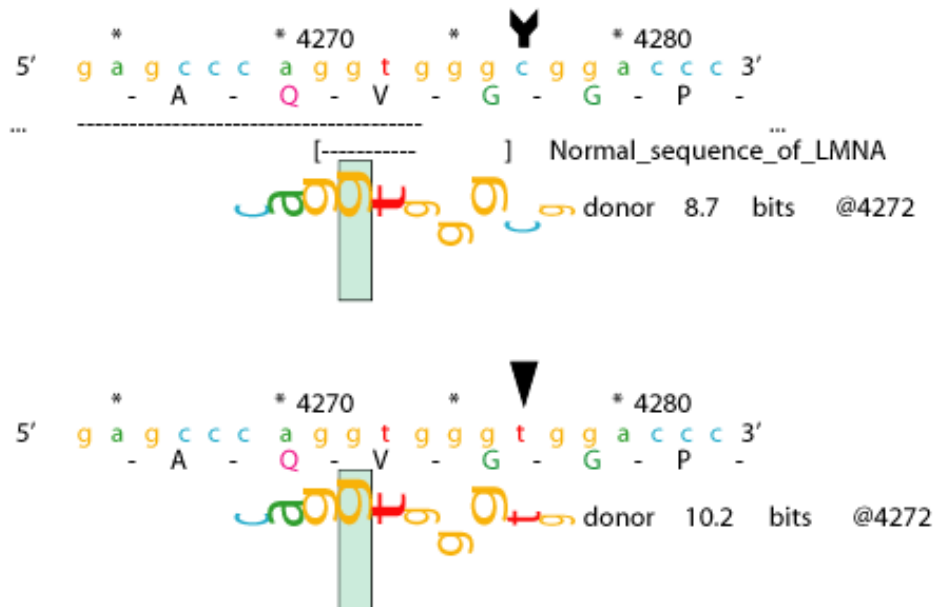A G-> A mutation 1 nucleotide upstream of the exon 8 donor site of the lysosomal lipase gene [LIPA; U04292] results in mild cholesterol ester storage disease with 4-9% enzymatic activity. The reduction in information content is significant even though the Ri value is still much greater than Ri,min.

# Cryptic splicing mutations



A C->T mutation in intron 3 of the iduronidate sulfate synthetase (Mucopolysaccharidosis type II) gene strengthens and activates a cryptic donor site in exon 3 of the gene (Rogan et al. 1998).



A synonymous C>T substitution at codon 608 **strengthens** a cryptic donor splice site in exon 11 the *LMNA* gene in patients with Hutchinson-Gifford progeria (Eriksson et al. Science 2003). The walker, shown below the sequence, indicates a pre-existing 8.7 bit cryptic site that is strengthened by the mutation to 10.2 bits (>=2.8 fold).
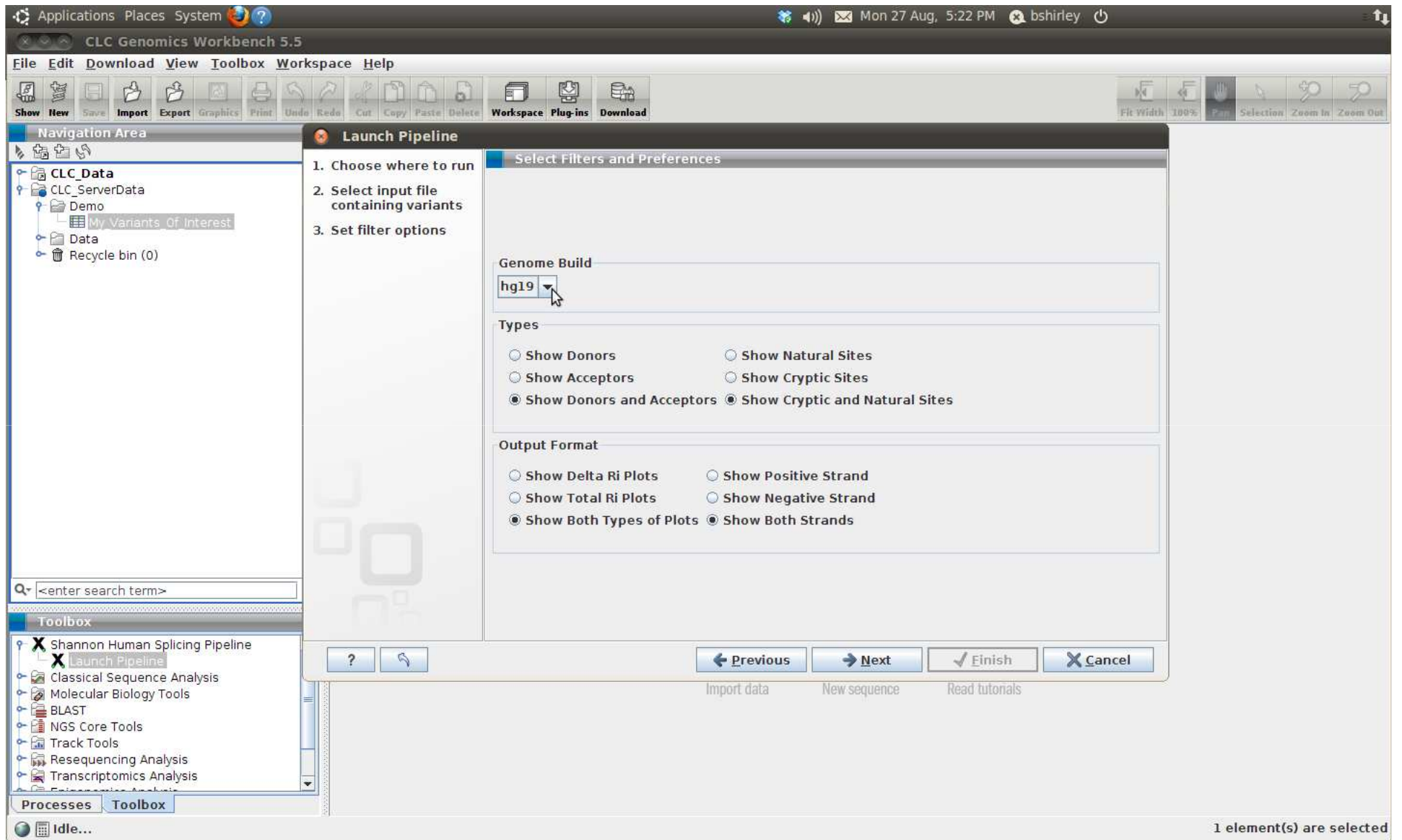
# Shannon Pipeline for mutation analysis

- For prediction of functionally-significant, non-coding variants in genome or exome sequences

- Mutation analysis on a genome-scale

- [Patented](#) and proven information theory-based binding site analysis (US Pat. 5,867,402):

    o Make quantitative predictions - information related to binding affinity.

    o Can distinguish benign, fully and partially inactivating binding site variants

    o Common paradigm for all types of nucleic binding sites (eg. splicing, transcription factors)

- Algorithm has been validated in hundreds of [peer-reviewed research studies of splicing mutations](#)

- Algorithm recommended by the American College of Medical Genetics and Genomics in their published guidelines and standards (*Genet. Med.* 7, 571–583)

- Predecessor [software for single mutation analysis](#) has been designated as a medical device by US FDA (not approved for clinical diagnostics)

# How the plug-in works

- CLC-Genomics Workbench retrieves lists of variants, and either processes the data itself or funnels it to the Shannon pipeline on the Genomics Server.

- Genome-wide information analysis is performed for all variants

- Variants with changes in information content (in bits) are annotated against standard databases (Ensembl Refseq, dbSNP)

- Prospective mutations are categorized and filtered

- Results displayed as exportable chromosome plots, sortable tables, and genome browser tracks.

# Choose genome reference and results to display ...

# Completed run ... Tabular output for cryptic sites below,

**CLC Genomics Workbench 5.5.1**

File   Edit   Download   View   Toolbox   Workspace   Help

Show  New  Save  Import  Export  Graphics  Print  Undo  Redo  Cut  Copy  Paste  Delete  Workspace  Plug-ins  Download   Fit Width  100%  Pan  Selection  Zoom In  Zoom Out

Navigation...

Pos Strand Ac... × | * Track List × | Inactivating ... × | Complete Vari... × | Cryptic Varia... ×

Rows: 22,197      Effect of variants on Ri and other relevant information                Filter:

| Chromoso... | Coordinate | Strand | Ri-initial | Ri-final | ΔRi | Type | Gene Name | Location | Location ... | Loc. Rel. t... | Dist. from... | Loc. of ne... | Ri of near... | Cryptic Ri ... | rsID if ava... | Average |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 2275984... | + | 3.43 | 2.33 | -1.09 | ACCEPTOR | CTD-209... | CRYPTICS... | INTRONIC | - | - | - | - | - | rs3014290 | 0.34260 |
| 1 | 2275987... | + | 3.25 | 2.04 | -1.21 | ACCEPTOR | CTD-209... | CRYPTICS... | INTRONIC | - | - | - | - | - | rs3000786 | 0.29890 |
| 1 | 2275989... | + | 6.27 | 5.05 | -1.22 | ACCEPTOR | CTD-209... | CRYPTICS... | INTRONIC | - | - | - | - | - | rs3000788 | 0.29972 |
| 1 | 2275990... | - | 0.23 | 1.32 | 1.09 | ACCEPTOR | CTD-209... | CRYPTICS... | EXONIC | - | - | - | - | - | rs3000789 | 0.30726 |
| 1 | 2276808... | - | 4.92 | 2.40 | -2.52 | DONOR | RP11-27... | CRYPTICS... | EXONIC | - | - | - | - | - | rs35878... | 0 |
| 1 | 2276943... | + | -1.17 | 0.33 | 1.51 | ACCEPTOR | RP11-27... | CRYPTICS... | EXONIC | - | - | 2276943... | 3.45 | LESS | rs13794... | 0 |
| 1 | 2309065... | - | 0.41 | 3.02 | 2.61 | DONOR | RP11-99J... | CRYPTICS... | INTRONIC | - | - | - | - | - | rs16852... | 0.15434 |
| 1 | 2341629... | - | -2.42 | 0.10 | 2.52 | DONOR | RAC1P7 | CRYPTICS... | EXONIC | - | - | - | - | - | rs12144... | 0.31660 |
| 1 | 2344008... | + | 2.80 | 4.82 | 2.01 | ACCEPTOR | RP4-799... | CRYPTICS... | INTRONIC | - | - | - | - | - | rs559117 | 0.49820 |
| 1 | 2347975... | + | 5.14 | 1.41 | -3.73 | DONOR | RP4-781... | CRYPTICS... | INTRONIC | - | - | - | - | - | rs11803... | 0.32289 |
| 1 | 2348165... | + | 0.91 | 2.19 | 1.28 | ACCEPTOR | RP4-781... | CRYPTICS... | EXONIC | - | - | 2348165... | -3.06 | GREATER | rs482329 | 0.4979 |
| 1 | 2348486... | + | 1.03 | -9.85 | -10.88 | ACCEPTOR | RP4-781... | CRYPTICS... | EXONIC | - | - | 2348484... | 9.12 | LESS | rs486142 | 0.40968 |
| 1 | 2348487... | + | 1.01 | -9.87 | -10.88 | ACCEPTOR | RP4-781... | CRYPTICS... | EXONIC | - | - | 2348484... | 9.12 | LESS | rs508293 | 0.43893 |
| 1 | 2348487... | + | -9.67 | 5.05 | 14.71 | ACCEPTOR | RP4-781... | CRYPTICS... | EXONIC | - | - | 2348484... | 9.12 | LESS | rs508293 | 0.43893 |
| 1 | 2348487... | + | -8.02 | 6.69 | 14.71 | ACCEPTOR | RP4-781... | CRYPTICS... | EXONIC | - | - | 2348484... | 9.12 | LESS | rs10910... | 0.24889 |
| 1 | 2350939... | + | -2.39 | 0.62 | 3.01 | DONOR | RP11-44... | CRYPTICS... | INTRONIC | - | - | 2350931... | 8.68 | LESS | rs2802926 | 0.04996 |
| 1 | 2367063... | - | 4.80 | 5.92 | 1.12 | ACCEPTOR | RP11-38... | CRYPTICS... | INTRONIC | - | - | - | - | - | rs2758175 | 0.46875 |
| 1 | 2367066... | - | -1.34 | 0.32 | 1.66 | ACCEPTOR | RP11-38... | CRYPTICS... | INTRONIC | - | - | - | - | - | rs2243530 | 0.49382 |
| 1 | 2367075... | - | 3.32 | -3.71 | -7.03 | DONOR | RP11-38... | CRYPTICS... | INTRONIC | - | - | - | - | - | rs2758180 | 0.49586 |
| 1 | 2380483... | + | -0.37 | 2.64 | 3.01 | DONOR | RP11-19... | CRYPTICS... | INTRONIC | - | - | - | - | - | rs2298100 | 0.22543 |
| 1 | 2386553... | + | -13.60 | 5.03 | 18.63 | DONOR | RP11-17... | CRYPTICS... | EXONIC | - | - | - | - | - | rs2392861 | 0.46142 |
| 1 | 2422210... | + | 4.70 | -13.92 | -18.62 | DONOR | RP11-32... | CRYPTICS... | INTRONIC | 3'-FLANKI... | 41 | 2422210... | 6.43 | LESS | rs908970 | 0.02666 |
| 1 | 2437090... | + | 4.52 | -14.11 | -18.63 | DONOR | RP11-26... | CRYPTICS... | INTRONIC | 3'-FLANKI... | 147 | 2437089... | 7.18 | LESS | rs1473466 | 0.08869 |
| 1 | 2442458... | - | 4.12 | -7.54 | -11.67 | ACCEPTOR | RP11-27... | CRYPTICS... | INTRONIC | - | - | - | - | - | rs2454230 | 0.44444 |
| 1 | 2442463... | - | -10.29 | 1.38 | 11.67 | ACCEPTOR | RP11-27... | CRYPTICS... | INTRONIC | - | - | - | - | - | rs2454231 | 0.29752 |
| 1 | 2442505... | + | -8.95 | 9.68 | 18.63 | DONOR | RP11-27... | CRYPTICS... | INTRONIC | - | - | - | - | - | rs2500491 | 0.32 |
| 1 | 2445578... | - | -5.81 | 2.97 | 8.78 | ACCEPTOR | RP11-51... | CRYPTICS... | INTRONIC | - | - | - | - | - | rs3127484 | 0.28621 |
| 1 | 2461979... | + | 8.32 | 10.04 | 1.72 | DONOR | RP11-83... | CRYPTICS... | INTRONIC | 3'-FLANKI... | 46 | 2461979... | 8.32 | GREATER | rs10924... | 0.04747 |
| 1 | 2461979... | + | -4.32 | 2.72 | 7.03 | DONOR | RP11-83... | CRYPTICS... | INTRONIC | 3'-FLANKI... | 50 | 2461979... | 8.32 | LESS | rs10924... | 0.04747 |
| 1 | 2466771... | - | 3.74 | 7.33 | 3.59 | DONOR | RP11-69... | CRYPTICS... | EXONIC | - | - | 2466770... | -37.41 | GREATER | rs3120684 | 0.06024 |
| 1 | 2468463... | + | 5.14 | -13.49 | -18.63 | DONOR | RP11-43... | CRYPTICS... | EXONIC | - | - | - | - | - | rs10802... | 0.49561 |
| 1 | 2477994... | - | -2.08 | 0.02 | 2.10 | ACCEPTOR | RP11-97... | CRYPTICS... | INTRONIC | 3'-FLANKI... | 177 | 2477992... | 1.68 | LESS | rs1176039 | 0.5 |
| 1 | 2478353... | + | 1.08 | -0.52 | -1.60 | DONOR | RP11-63... | CRYPTICS... | INTRONIC | - | - | - | - | - | rs1151641 | 0.33369 |
| 1 | 2481382... | + | -17.05 | 1.58 | 18.63 | DONOR | RP11-43... | CRYPTICS... | EXONIC | - | - | - | - | - | rs4451578 | 0.37750 |
| 1 | 2481382... | + | 4.28 | 2.36 | -1.93 | ACCEPTOR | RP11-43... | CRYPTICS... | EXONIC | - | - | - | - | - | rs4451578 | 0.37750 |
| 1 | 2487228... | + | -2.22 | 0.95 | 3.18 | DONOR | RP11-43... | CRYPTICS... | INTRONIC | | | | | | rs78125... | 0.48411 |

Toolbox

Shannon Hu...
Launch P...
Classical Se...
Molecular Bi...
BLAST
NGS Core To...
Track Tools
Resequenci...
Transcriptor...
Epigenomics...
De Novo Se...
Workflows

Processes
Toolbox

Idle...                                                                                        1 element(s) are selected

# Filtering cryptic splice site variants hones in on likely mutations
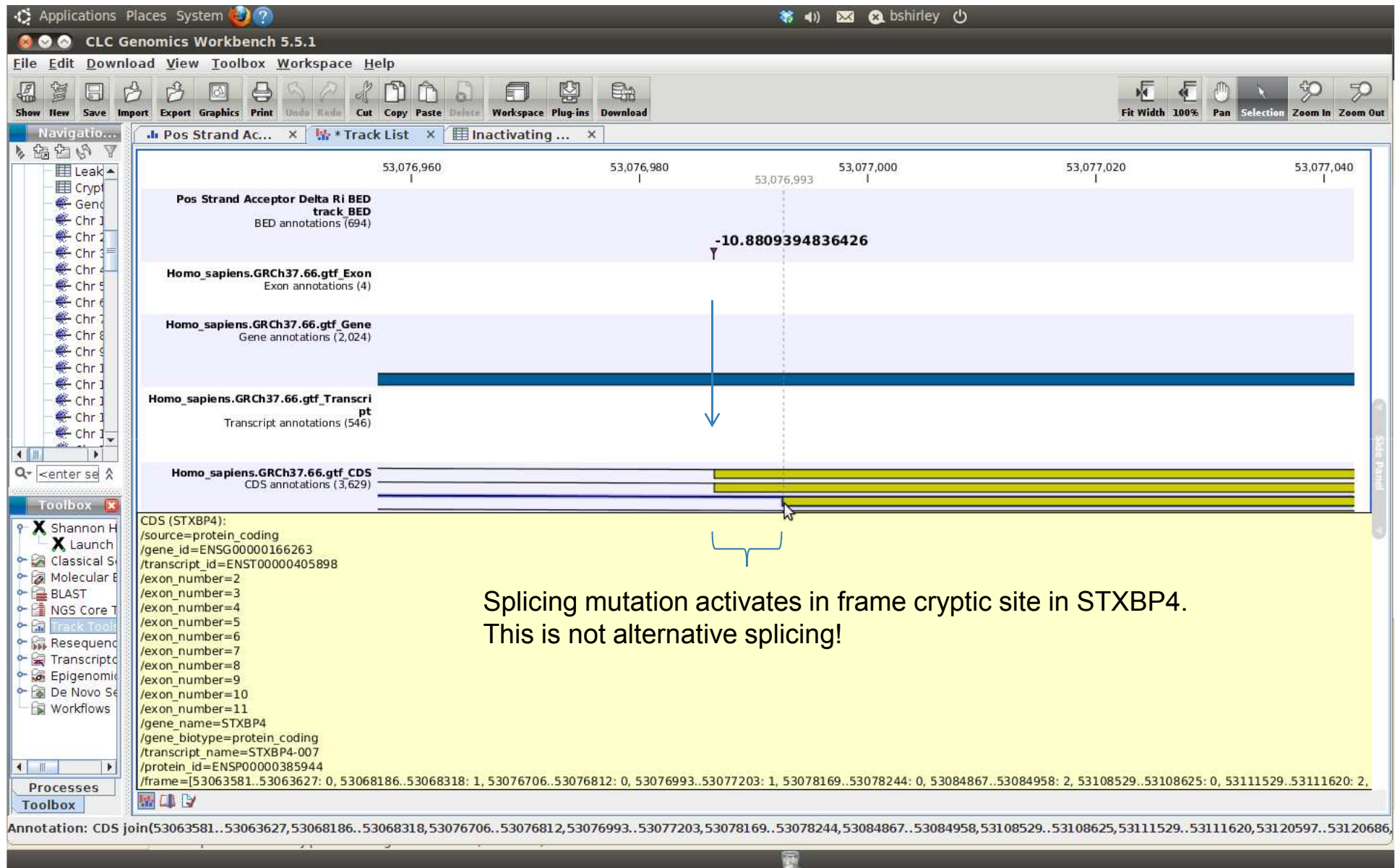
# "Manhattan-like" plots to survey all potential splicing related information changes by chromosome

## Chromosome 1



chr1: 201918763, ΔRi: 11.67, Final Ri: 1.95, rsID: 55868763 : (201918763.000, 11.669)
chr1: 202078797, ΔRi: 11.67, Final Ri: 0.96, rsID: - : (202078797.000, 11.669)
chr1: 203318206, ΔRi: 11.67, Final Ri: 0.64, rsID: - : (203318206.000, 11.669)

# Displaying custom track of $\Delta R_i$ values with Genome Workbench Browser

# Implementation

1. *Shannon Human Splicing Pipeline* has been released for **Linux** and **MacOSX** operating systems supporting Perl and gcc.
2. Installation has been verified with Perl v.5.8.8 and 5.10.1 and gcc v.4.1.2 and v.4.4.3 with the Ubuntu 2.6.32-27 (32 and 64 bit), CentOS 2.6.18-238 (64 bit), and Fedora 16 (32 bit) kernels, and MacOSX (Mountain Lion release version 10.8, Lion release version 10.7.4; gcc v.4.2.1 and Perl 5.12.3 and 5.12.4).
3. Several C libraries determine the information content of a position in the genome before and after a variant is introduced using convolution-style sliding-window computation. Changes in $R_i$ introduced by genomic variation are computed by subtracting the initial $R_i$ value of a position by the sum over a surrounding window, then adding the new value for each position ($\Delta R_i$).
4. Perl scripts wrap these C libraries and annotate data pipeline results. Integration with the CLC-Bio workbench environment was achieved through code written in Java utilizing the CLC-Bio developer API.
5. This software is assembled as a client plugin requiring a connection to the server to execute, a server plugin, and a standalone client plugin. Two additional plugins contain a modified dbSNP135 (Indels and extraneous data removed), Ensembl Exon Data (Build 66), and GRCh37/ NCBI36 respectively.

## Performance

| Number of variants | Running time |
|---|---|
| 100,000 | 37m |
| 211,049 | 1h 12m |
| 290,589 | 1h 17m |
| 314,637 | 1h 20m |

I7 Processor,16 Gb RAM, Ubuntu Linux

# Contributors/ Support